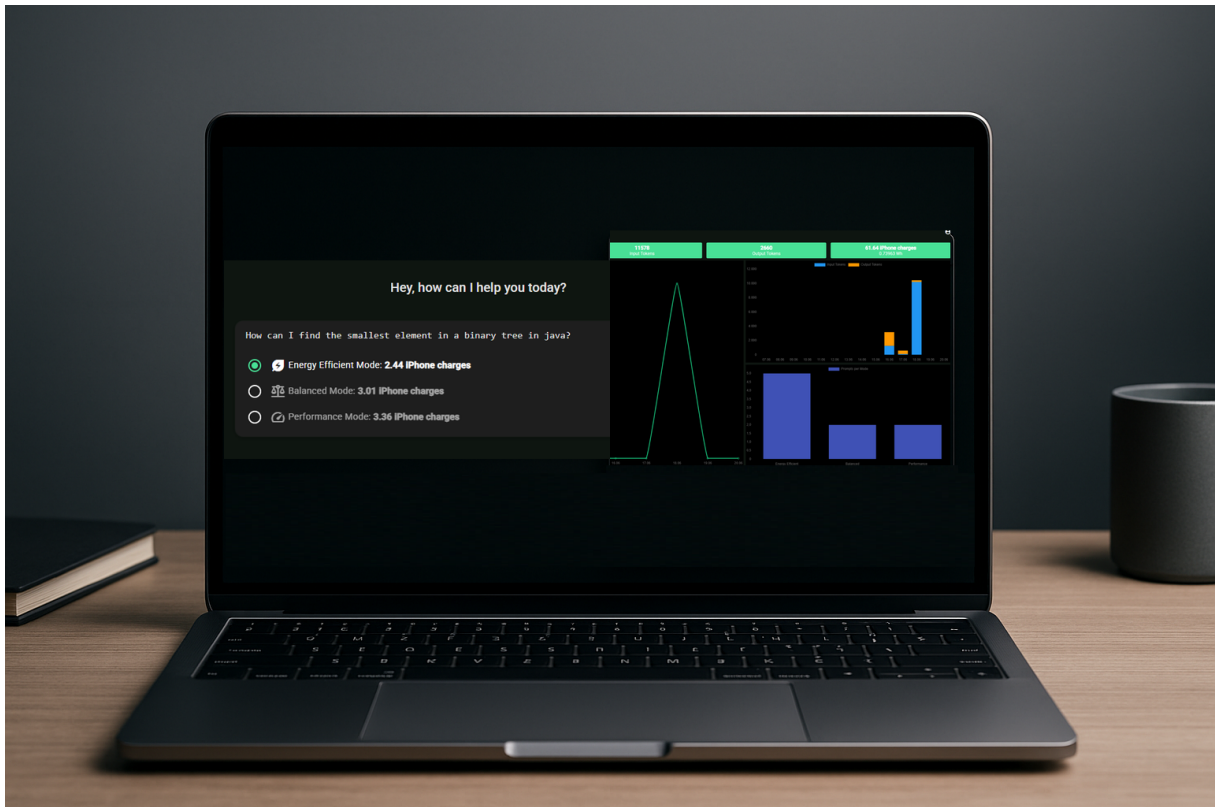


# UI Strategies for Reducing Conversational AI Energy Consumption

Jack Gläser and Simon Lüscher

Windisch, August 2025



Students	Jack Gläser and Simon Lüscher
Expert	Romano Roth
Supervisors	Prof. Martin Kropp and Dr. Nitish Patkar
Stakeholder	Fachhochschule Nordwestschweiz - Hochschule für Technik
Project number	IMVS-24
Fachhochschule Nordwestschweiz, Hochschule für Technik	

## Abstract

The rapid adoption of large language models (LLMs) for everyday information-seeking is shifting energy demand from highly optimized search engines to far more power-hungry conversational AI systems [1], [2]. Yet this cost remains largely invisible to end users.

This thesis investigates whether user interface (UI)-only interventions can increase user awareness and lead to measurable reductions in the energy consumption associated with chatbot usage.

We first conducted a baseline survey ( $n = 50$ ), which revealed both awareness deficits and strong support for transparency features such as energy indicators and low-power modes. Guided by these findings and prior work in human computer interaction (HCI), we developed a full-stack ChatGPT-style prototype (“The Botter”) incorporating five sustainability-focused UI features: (1) a three-mode toggle, (2) prompt-level energy prediction, (3) per-response Energy-Notes, (4) a usage metrics dashboard and (5) personalized energy analogies.

In a five-day field experiment involving eleven frequent LLM users, we observed a substantial increase in energy awareness (peaking at  $M = 4.44$  on a 5 point Likert scale), alongside high usability ratings (all features above 4/5). Importantly, higher awareness correlated with more energy-efficient usage behavior: Participants in the top awareness quartile selected the energy-saving mode for 72% of their prompts, compared to 34% in the bottom quartile. Overall, more than half of all prompts were routed through the energy-efficient mode, yielding an estimated 35% reduction in energy consumption relative to a performance-mode baseline, without degrading user experience.

These findings demonstrate that lightweight, frontend design interventions can effectively nudge sustainable behavior in conversational AI interfaces. This contributes to addressing a key challenge in the design of human-AI interaction and highlights the role of interface design in mitigating the environmental impact of large-scale language model usage.

## Acknowledgement

We would like to express our sincere gratitude to Prof. Martin Kropp and Dr. Nitish Patkar for their steady supervision, constructive feedback, and constant encouragement throughout this project.

We also thank Dr. Pooja Rani (University of Zurich) for her insightful advice on the experimental design and consent procedures, which markedly improved the study's quality.

Finally, a heartfelt thank-you to all who took part in our research: The 50+ survey respondents who provided the baseline insights and the eleven experiment participants who devoted a full week to testing the prototype. Your time and input made this thesis possible.

Thank you all.

# Contents

<b>List of Figures</b>	<b>vi</b>
<b>List of Tables</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Initial Situation . . . . .	1
1.2 Problem Statement . . . . .	1
1.3 Overall Goal and Requirements . . . . .	2
1.4 Structure of the Thesis . . . . .	2
<b>2 Methodology</b>	<b>4</b>
2.1 Phase I – Literature Review . . . . .	4
2.2 Phase II – Survey . . . . .	4
2.3 Phase III – Prototype Development . . . . .	5
2.4 Phase IV – Controlled Experiment . . . . .	5
2.5 Phase V – Data Preparation and Analysis . . . . .	6
<b>3 State of the Art</b>	<b>8</b>
3.1 Perceptual Gap: User Awareness and Behavior . . . . .	8
3.2 Survey results . . . . .	8
3.3 Technical Gap: Measurement and Mitigation Strategies . . . . .	10
3.4 Design Gap: UI Feedback, Tools, and Market Landscape . . . . .	11
3.5 Conclusion and Implications . . . . .	11
<b>4 Conceptual Solution</b>	<b>13</b>
4.1 Solution Approach . . . . .	13
4.2 Baseline Chatbot . . . . .	13
4.3 Features and Functionalities . . . . .	13
4.4 Definitive Features . . . . .	16
4.5 Estimation Model of Energy Consumption . . . . .	19
<b>5 Implementation</b>	<b>23</b>
5.1 System Architecture . . . . .	23
5.2 Frontend . . . . .	23
5.3 Backend . . . . .	26

5.4	Solution Structure . . . . .	27
5.5	Feature Implementation . . . . .	29
5.6	Data Storage: Cosmos DB . . . . .	35
5.7	Non functional requirements . . . . .	37
5.8	Limitations and constraints . . . . .	40
<b>6</b>	<b>Validation and Results</b>	<b>41</b>
6.1	Results from Daily Check-in and Final Questionnaire (a & b) . . . . .	41
6.2	Final Questionnaire . . . . .	42
6.3	Behavioral Results from Application Data . . . . .	43
6.4	Summary . . . . .	48
6.5	Hypothesis Validation Summary . . . . .	48
<b>7</b>	<b>Discussion</b>	<b>50</b>
7.1	Interpretation and Evaluation . . . . .	50
7.2	Evaluation of Energy consumption . . . . .	50
7.3	Answering the Research Questions . . . . .	50
7.4	Limitations . . . . .	51
7.5	Future Work and Development . . . . .	51
7.6	Concluding Remarks . . . . .	52
<b>8</b>	<b>Conclusion</b>	<b>53</b>
	<b>Declaration of honesty</b>	<b>61</b>
<b>A</b>	<b>Appendix</b>	<b>62</b>
A.1	Survey Form and Results . . . . .	62
A.2	Experiment Forms and Results . . . . .	74
A.3	JSON Schema . . . . .	87

## List of Figures

2.1	Overview of the five-phase thesis methodology . . . . .	4
4.1	Baseline ChatGPT-inspired chatbot without sustainability features . . . . .	13
5.1	Application architecture overview . . . . .	23
5.2	Frontend project structure . . . . .	24
5.3	app.routes.ts . . . . .	25
5.4	Proxy configuration file . . . . .	26
5.5	Interface abstraction for OpenAI integration . . . . .	28
5.6	Users opening the application will land on /chat and be able to start a new dialog instantly	29
5.7	Previous conversations can be revisited to continue prompting or collect information . .	30
5.8	Users can delete or rename conversations via the menu icon . . . . .	30
5.9	Mode can be switched for every prompt . . . . .	30
5.10	Real-time energy prediction for current prompt . . . . .	31
5.11	Workflow: Predicting output tokens and applying energy calculation . . . . .	31
5.12	Prediction of output token count based on input tokens . . . . .	31
5.13	Energy usage calculation based on token counts . . . . .	32
5.14	Post-response energy note with real-world analogy . . . . .	32
5.15	User metrics dashboard with energy, token, and mode usage breakdown . . . . .	33
5.16	Guessing energy analogies . . . . .	33
5.17	Users can choose their preferred energy analogy . . . . .	34
5.18	User settings . . . . .	34
5.19	Logical database schema (conceptual relationships) . . . . .	35
6.1	Total prompts sent per mode per day . . . . .	45
6.2	Average normalized prompt growth curve . . . . .	47
6.3	Scatter plot of input vs. output tokens with LOWESS trend and prediction function . . .	47

## List of Tables

3.1	Awareness-related Likert Items (1 = Strongly Disagree, 5 = Strongly Agree) . . . . .	9
4.1	Technical configuration features . . . . .	14
4.2	Awareness and behavioral features . . . . .	14
4.3	Gamification and feedback mechanisms . . . . .	15
4.4	Decision matrix of all evaluated UI features . . . . .	15
4.5	Overview of metrics shown on the user dashboard . . . . .	17
4.6	Model pricing per million tokens and cost-based scaling . . . . .	20
4.7	Reported energy consumption per query and per token . . . . .	20
4.8	Comparison of estimated coefficients for GPT-4o . . . . .	21
4.9	Final energy model coefficients per mode . . . . .	21
5.1	All important frontend packages . . . . .	24
5.2	Overview of implemented backend Azure Functions . . . . .	27
6.1	Mean awareness per day “ <i>At this moment I’m aware of the energy cost of the prompts I sent today</i> ”, (1 = Strongly Disagree, 5 = Strongly Agree) . . . . .	41
6.2	Behavioral items, (1 = Strongly Disagree, 5 = Strongly Agree) . . . . .	42
6.3	Usability ratings: “ <i>Easy to understand and effect clear</i> ”, (1 = Strongly Disagree, 5 = Strongly Agree) . . . . .	42
6.4	Energy unit preferences . . . . .	43
6.5	Metrics visits per user and day . . . . .	44
6.6	Share of prompts, tokens, and energy consumption per mode . . . . .	44
6.7	Energy consumption per 1,000 input tokens by mode . . . . .	45
6.8	Aggregated prompts per user and per mode . . . . .	46
8.1	External tools used during thesis work . . . . .	61

## Glossary

**LLM** Large Language Model. AI model trained on vast text data to generate human-like language, e.g., GPT-4.

**UI** User Interface. The means by which a user interacts with a computer system or software.

**SPA** Single-Page Application. A web application that loads a single HTML page and dynamically updates content.

**CO<sub>2</sub>** Carbon Dioxide. A greenhouse gas, often referenced in the context of emissions from computational processes.

**Wh** Watt-hour. A unit of energy used to quantify electricity consumption.

**RAG** Retrieval-Augmented Generation. An AI technique that combines information retrieval with generative models.

**Azure Functions** A serverless compute service provided by Microsoft Azure for running event-driven code.

**Cosmos DB** A globally distributed, multi-model database service by Microsoft Azure.

**ChatGPT** A conversational AI model developed by OpenAI, based on the GPT architecture.

**Energy-Efficient Mode** A system mode designed to minimize energy consumption, often by using smaller models or reduced computation.

**Performance Mode** A system mode prioritizing response quality and speed, typically using larger models and more computation.

**Balanced Mode** A compromise between energy efficiency and performance in system operation.

**Token** In NLP, a unit of text (word, subword, or character) processed by language models.

**Inference** The process of running a trained AI model to generate predictions or responses.

**Prompt** The input text provided to an AI model to elicit a response.

**Dashboard** A user interface element that displays key metrics and statistics, such as energy usage.

**Awareness** User knowledge of the energy and environmental impact of their AI usage.

**HCI** Human-Computer Interaction. The study and design of how people interact with computers and technology.

**FLOPs** Floating Point Operations per Second. A measure of computational performance, often used to describe the processing power required by AI models.

**Heuristic** A practical method or approach to problem-solving that is not guaranteed to be optimal but is sufficient for reaching an immediate goal.

**Latency** The delay between a user's action and the system's response, often measured in milliseconds.

**Modalities** Different forms or types of input/output data (e.g., text, image, audio) that an AI system can process.

**Prompt Engineering** The practice of designing and refining input prompts to optimize the output of language models.

**Sustainability** Meeting present needs without compromising the ability of future generations to meet their own needs, often referenced in the context of environmental impact.

**Tokenization** The process of breaking text into smaller units (tokens) for processing by language models.

**Usability** The ease with which users can effectively and efficiently interact with a system or product.

**LOWESS** Locally Weighted Scatterplot Smoothing. A non-parametric regression method that fits simple models to localized subsets of data to create a smooth curve through points in a scatterplot.

**Empiric/Empirical** Based on observation or experience rather than theory or pure logic; often refers to data or results derived from experiments or real-world evidence.

**Regression** A statistical method for modeling the relationship between a dependent variable and one or more independent variables.

**Non-parametric** Refers to statistical methods that do not assume a specific distribution for the data.



# 1 Introduction

## 1.1 Initial Situation

Conversational AI consumes a lot of energy during its entire life cycle [3]. A significant part of this energy is required for the development and training of the model [4], [5]. But the actual conversational AI service also consumes a great amount of energy during inference. Although a single prompt does not have a significant impact on the required energy, the cumulative effect grows rapidly with a large user base. In fact, inference energy usage now outweighs training in many large-scale deployments: If chatbot-based interfaces were to completely replace traditional web search, global electricity demand could increase by several tens of terawatt hours annually [3]. During one month, the inference energy can already exceed that of the model training [6]. Bond Capital reports that ChatGPT reached 800 million weekly active users just 17 months after launch, implying trillions of inference calls per month [7]. These figures underscore why inference energy has become a central environmental concern.

## 1.2 Problem Statement

Despite this growing footprint, most end users remain unaware of the hidden environmental costs of their interactions with large language models (LLMs). Traditional searches have well-optimized, transparent energy profiles, but LLM-based chat incurs substantial energy consumption. In what follows, we outline the key dimensions of this problem.

### 1.2.1 Perceptual Gap: Lack of User Awareness

Most users underestimate the energy implications of AI services. Our survey (Section 3.2) shows participants vastly misjudge per-query consumption, echoing Sustainable HCI findings that even tech-savvy users lack calibrated mental models for digital energy use [8]. Without awareness, users cannot factor sustainability into everyday choices.

### 1.2.2 Technical Gap: Barriers to Live Transparency

Accurate, per-query energy feedback is technically complex. Current estimation tools are not very accurate, offline, or not integrated into user interfaces [9], [10]. Live token-level attribution remains an open challenge behind proprietary, cloud-based APIs. Without precise, real-time estimates, transparency features risk being symbolic rather than informative.

### 1.2.3 Design Gap: Platform Incentives and UI Constraints

Commercial providers often deprioritize energy transparency when performance, engagement, and monetization are at stake. Defaults prioritize speed and fluency over efficiency, and although model-selection or prompt-optimization can cut inference energy by up to 60 % [11], [12], these options are rarely exposed to users. The design ecosystem lacks strong nudges toward sustainable behavior.

### 1.2.4 Problem Synthesis

Taken together, these facets reveal three interrelated gaps that this thesis addresses:

1. **Perceptual Gap** — Users are unaware of the true environmental cost of LLM usage.
2. **Technical Gap** — Live, accurate energy estimation is difficult to implement and deploy.
3. **Design Gap** — Current UIs and provider incentives offer little support for sustainability-aware decisions.

### 1.2.5 Opportunities for Intervention

Research in adjacent domains shows that real-time, contextual feedback can shift behavior, for example smart-meter dashboards and mobile energy apps have driven measurable reductions in consumption [13], [14]. UI-only strategies such as dashboards, mode toggles, and energy analogies can make backend energy innovations visible and actionable at the user level. This thesis builds on these insights to explore how purely interface-based interventions can raise awareness and nudge more sustainable use of conversational AI.

## 1.3 Overall Goal and Requirements

We strive to reduce the overall energy consumption associated with conversational AI by improving user awareness through targeted UI-based interventions. By researching and contributing to the state of the art, this thesis aims to understand how interface design alone can effectively influence user behavior regarding energy consumption. We will identify and evaluate UI features that increase user awareness, providing transparency and actionable insights into the energy impact of their interactions with conversational AI. Ultimately, our goal is to empower users to make informed decisions when and how to use conversational AI, fostering more sustainable and energy-efficient usage patterns without compromising the user experience.

### 1.3.1 Research Questions

Therefore, our work is guided by the following three research questions:

- RQ1:** *To what extent are users currently aware of the energy implications associated with their chatbot interactions?*
- RQ2:** *How can UI-based features most effectively increase user awareness regarding the energy consumption of conversational AI?*
- RQ3:** *How strongly does increased user awareness correlate with reductions in conversational AI energy consumption?*

## 1.4 Structure of the Thesis

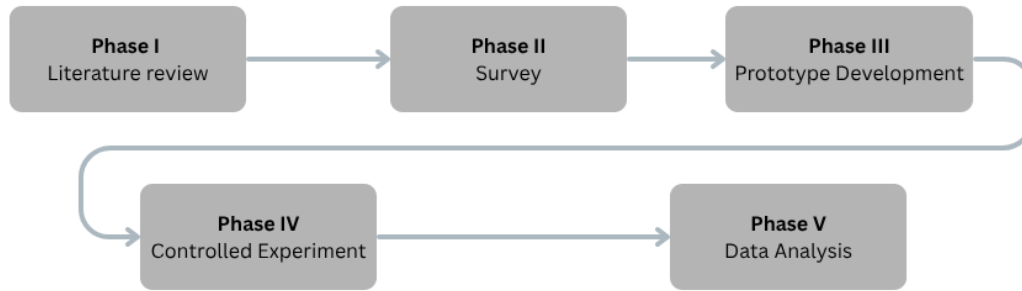
This thesis is organized into eight core chapters, followed by a declaration of honesty and a comprehensive appendix.

- **Chapter 1 - Introduction** situates the study, delineates the problem space, formulates the research questions, and specifies the goals and requirements that guide the work.
- **Chapter 2 - Methodology** details the five-phase research design: Literature review, the baseline survey, prototype development, controlled experiment, and data-analysis plan.
- **Chapter 3 - State of the Art** reviews prior work on user awareness, energy footprints of conversational AI, UI-based sustainability interventions, and identifies the resulting knowledge gaps addressed in this thesis including our own survey results.
- **Chapter 4 - Conceptual Solution** outlines the high-level approach, describes the baseline chatbot and the five UI features (three-mode switch, metrics dashboard, prompt prediction, energy note, and energy analogies) and presents the underlying energy-estimation model.
- **Chapter 5 - Implementation** explains the concrete realisation of the system architecture on Azure, covering frontend, backend, data storage, non-functional requirements, and CI/CD pipeline.
- **Chapter 6 - Validation and Results** reports quantitative and qualitative findings from the user study(experiment), including awareness trajectories, behavioral metrics derived from server logs, and usability assessments.

- **Chapter 7 - Discussion** interprets the empirical results in light of the research questions and hypotheses, reflects on methodological and technical limitations, and outlines avenues for further research and development.
- **Chapter 8 - Conclusion** synthesises the key contributions and practical implications of the thesis.

## 2 Methodology

The following section outlines the five phases of the methodology that guided the development of this thesis. Each phase is briefly described to provide an overview of the overall approach.



**Figure 2.1:** Overview of the five-phase thesis methodology

### 2.1 Phase I – Literature Review

We conducted a structured literature review to identify existing research and solutions related to sustainable AI and energy awareness in conversational systems. The process included:

- **Databases Searched:** Google Scholar, Google search and AI (ChatGPT deepsearch functionality) for literature searches related to the subject.
- **Keywords Used:** Our search included combinations of terms such as “chatbot”, “LLM”, “large language models”, “sustainable AI”, “green AI”, “energy efficiency”, “carbon footprint AI”, “eco-feedback”, “energy”, “consumption”, “behavior change”, “energy literacy”, “gamification”, “peer comparison”, “eco-nudges”, “token efficiency”, “budget-constrained inference”, “carbon tracking widgets”, “HCI for sustainability”, and “sustainability metrics”. We adapted and expanded the keyword set iteratively as themes emerged.
- **Selection Criteria:** We prioritized peer-reviewed papers, highly cited publications, and industry reports published between 2016 and 2025. Additionally, we considered emerging tools (such as browser extensions) and preprints to gain insight into current trends and early-stage research on energy feedback and user-facing interventions in conversational AI.
- **Screening Process:** Abstracts and tool descriptions were screened for relevance to user awareness, energy consumption measurement, and UI-driven strategies for promoting sustainable behavior or raising awareness in chatbot interfaces or similar fields.

We totally collected 67 documents, articles or other tools and sources related to this keywords and used 45 directly in this thesis.

This methodology ensured that our review was both systematic and focused, providing a solid foundation for the analysis and proposed solutions presented in subsequent sections. The results of this review presented throughout the whole thesis directly shaped the structure of multiple sections such as the problem analysis and state-of-the-art, as well as the selection of features for the prototype development.

### 2.2 Phase II – Survey

To complement the literature review, we conducted a survey to assess the current user awareness and attitudes toward energy consumption in conversational AI. The survey methodology included:

- **Instrument:** The survey comprised fourteen closed items on a 5-point Likert scale [15] and three open questions.
- **Recruitment:** Participants were recruited via E-Mail, WhatsApp, LinkedIn, and in-person. Fifty valid responses were collected.
- **Data Collection:** Responses were collected anonymously through Google Forms, with an option for participants to add their email to volunteer for subsequent experiment.
- **Analysis Plan:** Descriptive statistics and correlation analysis was conducted to investigate the users awareness and examine their acceptance for energy consumption reducing UI-features.

The survey results are reported in Section 3.2 and informed the design of subsequent phases.

## 2.3 Phase III – Prototype Development

The development of the proof-of-concept chatbot (further called: *The Botter*) reproduces the ChatGPT interface and adds five sustainability modules which we wanted to test and research the effectiveness of:

- a) **Energy note** below each response (Wh and user-chosen equivalent, e.g. “0.53 Wh which equals to 43 s laptop use”).
- b) **Prompt prediction** shows the estimated energy cost before sending a query.
- c) **3-mode toggle** (Energy efficient, Balanced, Performance) with token-based cost prediction shown before sending to indicate that more powerful modes equal more energy consumption.
- d) **Usage-metrics dashboard** aggregating per-day Wh, modes share and more.
- e) **Energy-analogies** describe energy consumption in relatable units.

## 2.4 Phase IV – Controlled Experiment

### 2.4.1 Experimental Design

To evaluate the impact of our sustainability-focused UI features, we conducted a five-day controlled user study using our prototype (*The Botter*). The experiment was designed as a single-group field study, maximizing statistical power within the constraints of a small sample size [16]. All sustainability features were enabled by default, but participants could disable them at any time. System interactions, self-reports (daily check-ins and a final questionnaire at the end of the experiment) and background telemetry were logged for evaluation.

### 2.4.2 Hypotheses

This experiment was guided by the research questions defined in the section 1.3.1. Based on these questions, we formulated two hypotheses:

- H1.** Showing per-prompt consumption (*Energy-Note*) and predicted consumption (*Prompt Prediction + Three-Mode Switch*) increases the average awareness score from pre- to post-study.
- H2.** Individual awareness scores (end-of-study) are positively correlated with positive sustainability behavior logged during the study.

### 2.4.3 Participants and Duration

Eleven frequent LLM users (self-reported usage of  $\geq 15$  minutes per day in the baseline survey) were recruited from the survey pool. All participants live and work in Switzerland and demonstrated English proficiency, which is required because the prototype chatbot operates in English only. The sample spanned an age range of 25–35 years. Most participants were employed in IT-related roles or business services (e.g. software engineering, support, mortgage and financial consulting); one participant worked

in child care, providing some occupational diversity beyond the tech and finance sectors. All participants signed a consent document before starting the experiment, see A.2.1. The study ran across five consecutive workdays (Monday 23 June 2025 to Friday 27 June 2025). Each participant interacted with the prototype in their usual professional or academic context.

#### 2.4.4 Instrumentation and Logged Metrics

Data collection combined three sources:

- Daily check-ins consisting of four Likert scale items on awareness and feature salience, see Appendix A.2.2
- A final questionnaire covering awareness, behavior, usability, and open-ended feedback, see Appendix A.2.3
- Backend telemetry including prompts, input/output token counts, selected mode, energy usage estimates, feature toggles, and navigation events

Logs were stored in Azure Cosmos DB and retained only until 31 December 2025. All infrastructure was hosted in the EU.

#### 2.4.5 Procedure

The procedure of the whole experiment was the following:

**Day 0 – Onboarding:** Participants completed an informed consent form detailing data categories, retention period, and their right to withdraw. They also filled out a baseline questionnaire and received a remote walkthrough of the prototype.

**Days 1–5 – Usage period:** Participants used the chatbot during their daily work, received a morning reminder, and completed a four-item check-in survey in the evening.

**Day 5 – Wrap-up:** Logging was disabled and participants completed a comprehensive final questionnaire.

The study was approved through internal review and followed ethical guidelines for digital research involving personal usage data and behavioral tracking.

### 2.5 Phase V – Data Preparation and Analysis

After the experiment period, raw logs and database contents were exported as JSON and CSV. Then they were manually cleaned and combined in an Excel spreadsheet. Data sources included token counts, selected modes, timestamps, prompt content, page visits, and user preferences. To analyze the data Microsoft Excel and MATLAB were used.

The following methods were central to our data analysis:

- **Grouped Analysis:** Aggregated prompt counts, token usage, and energy consumption in different groups, per user, day and chat mode.
- **Daily Trends:** Prompt and page visit activity tracked over five days to observe usage patterns.
- **Prompt Growth Curves:** Normalized input-editing traces to assess user behavior prior to prompt submission.
- **Energy Efficiency:** Calculated Wh per 1,000 input tokens per mode (excluding fixed overhead) to compare efficiency.
- **LOWESS Regression:** Used to smooth input/output token scatter plots and analyze prediction trends.

The methodology combines survey research, prototype engineering and a five-day experiment to examine the impact of different UI-driven features regarding increasing energy awareness in conversational AI. The following chapter reviews existing work in related areas to provide context for our approach and highlight the research gaps this thesis aims to fill.

### 3 State of the Art

Building on the challenges identified in the problem analysis, in this section we review current research/literature and practical developments related to the three core gaps: User awareness (Perceptual Gap), technical constraints (Technical Gap), and design incentives and implementation (Design Gap). For each dimension, we synthesize insights from academic literature, existing tools, and our own survey to establish the research context and identify areas for innovation.

#### 3.1 Perceptual Gap: User Awareness and Behavior

A growing body of Sustainable-HCI and digital energy literacy research shows that users consistently underestimate the energy consumption and environmental impact of digital services. For instance, Preist et al. [8] found that only 17% of UK participants could estimate the energy cost of uploading a photo to the cloud within an order of magnitude. Walters et al. [17] demonstrated that awareness can be improved through visible, real-time carbon-footprint feedback via the *Purple* interface, although this did not reliably lead to sustainable behavioral changes. Similarly, Chen et al. [18] observed that anthropomorphic cues in chatbots increase user engagement and affect visual attention patterns. While their study did not examine sustainability directly, it suggests that enhancing engagement through design does not automatically translate into sustainable behavior.

These findings are echoed also in broader HCI research. Geelen et al. [13] showed that real-time feedback can improve energy awareness but only has an effect when the information is easy to understand and directly actionable. HCI literature consistently finds that *just-in-time*, personalized, and contextual feedback is more effective than generic sustainability messaging [14], [19].

Our own survey confirms the presence of a significant “intention–action” gap. While most participants express concern for sustainability, only a minority are willing to set personal usage limits or pay for energy offsets. However, strong support exists for transparency and optional eco-modes, suggesting that users are open to behavioral nudges when paired with credible and interpretable feedback mechanisms.

Research outside the AI domain highlights additional intervention mechanisms:

- *Real-time feedback*: Smart-meter apps displaying live usage data increased energy awareness, though savings depended on usability [13].
- *Social comparison*: Opower’s neighbour reports reduced household energy use by 2–4% [20], [21].
- *Gamification and eco-badges*: Reward systems increased pro-environmental behavior in mobile contexts [22].

These insights provide a foundation for our UI design, which aims to make energy data salient, actionable, and easy to interpret.

#### 3.2 Survey results

To further assess current user awareness of energy consumption in conversational AI, we conducted our own survey. A total of **50 participants** completed the questionnaire. The sample consisted predominantly of digitally literate respondents: 66% identified as technically advanced (e.g., developers, researchers), and 74% reported regular or heavy use of LLM chatbots such as ChatGPT, Gemini or Claude. Most worked in IT-related roles and used LLMs primarily for research, learning, or coding purposes. The majority accessed chatbots via desktop or laptop, and 38% reported using a paid license. While not representative of the general population, this technically savvy sample reflects the early-adopter audience most likely to engage with LLM technologies and therefore most exposed to their sustainability implications.



### 3.2.1 General Awareness and Attitudes

Participants expressed only moderate concern about environmental impacts: The statement “*I am concerned about the environmental impact of AI usage*” averaged 3.26 on a 5-point Likert scale as seen in the table Table 3.1. However, agreement was much higher for statements related to transparency and accountability. Items such as “*LLM chatbots should disclose energy usage*” and “*AI companies don’t disclose enough energy data*” each received a mean rating of 4.32, indicating strong support for disclosure even in the absence of precise knowledge.

Item	Mean	Std. Dev.
LLM chatbots should disclose energy usage	4.32	0.84
AI companies don’t disclose enough energy information	4.32	1.02
I am concerned about the environmental impact of AI usage	3.26	1.19
It is important to see energy consumption information	3.20	1.18
Sustainability is important in chatbot design	2.82	1.12

**Table 3.1:** Awareness-related Likert Items (1 = Strongly Disagree, 5 = Strongly Agree)

These results were further supported by correlation analysis. For instance, support for an “Eco Mode” was strongly correlated with a preference for low-carbon chatbot variants ( $r = 0.76$ ), and the perceived importance of showing energy usage data correlated with general environmental concern ( $r = 0.66$ ). This suggests that once awareness is present, users tend to support adaptive or energy-saving features.

### 3.2.2 Comparative Estimation and Perception

When asked to compare the energy consumption of ChatGPT to other familiar technologies, responses varied substantially. Most participants estimated that a single ChatGPT query uses around 10 times more energy than a Google search, and that 20–100 queries would equal the energy cost of fully charging an iPhone 14. However, answers spanned several orders of magnitude from 10 to 10,000 times more energy, or from 20 to 1000 queries per charge indicating highly wrong beliefs.

This variation points not only inaccurate estimations but also fundamentally different mental models of digital energy use. Such inconsistencies are well-documented in Sustainable HCI literature: Most users lack the mental calibration to estimate the energy cost of digital services within even an order of magnitude [8]. Crucially, this uncertainty appears to be self-recognised. The strong support for transparency (Table 3.1) suggests that participants are aware of their own informational gaps even when they self-report being “aware.”

This epistemic uncertainty is not entirely misplaced. The per-query energy cost of LLMs is itself an estimate, typically cited at around 0.30–0.34 Wh for GPT-4o [23], [24], while a Google search is estimated at 0.04 Wh [25]. Charging an iPhone 14 requires approximately 12.68 Wh [26], placing the true equivalence at roughly 4–40 prompts. In this light, the range of user responses reflects not just cognitive bias, but also a rational response to the lack of accessible and verifiable information. These findings highlight the value of visual and contextual feedback tools that bridge the perception gap regardless of precise technical measurements.

### 3.2.3 Qualitative Insights

Open-ended responses further illustrated the general uncertainty and lack of strategies. While some participants proposed fallback mechanisms (e.g., “route simple queries to Google”) or limitations on

excessive use, most responses were left blank or non-informative (e.g., “–”, “n/a”). This supports the interpretation that even among technically proficient users, baseline conceptual understanding of energy implications is limited.

### 3.2.4 Implications

The survey results confirm a significant gap between perceived and actual awareness. While users believe they understand the environmental costs of chatbot usage, their comparative estimations are widely scattered. At the same time, they strongly support transparency and feature based nudges suggesting a readiness to act if actionable information were available. This supports our design rationale: By providing contextualized, UI-driven feedback (e.g., dashboards, equivalency visualizations, or mode switching), we aim to raise user awareness without requiring expert knowledge or technical detail.

## 3.3 Technical Gap: Measurement and Mitigation Strategies

The energy consumption of LLMs arises from training, fine-tuning, and especially inference. While early studies focused on training costs, inference now dominates total energy use due to its frequent and large-scale deployment [3]. With nearly a billion users projected by 2025 [27], the aggregate energy burden is expected to grow rapidly.

Current benchmarks estimate that a single GPT-4o query consumes 0.30–0.34 Wh [23], [24], making it roughly eight times more energy-intensive than a Google search [25]. If chatbot usage replaced search engines at scale their energy consumption could increase by 60–70 times [1].

Mitigation strategies exist at multiple levels:

- (a) **Model-level:** Sparsification, quantisation, distillation, and retrieval-augmented generation (RAG) reduce FLOPs per token [28].
- (b) **System-level:** Dynamic model selection, GPU frequency scaling, and batch scheduling cut inference energy by 40–60% [11], [12], [29].
- (c) **Infrastructure-level:** Renewable energy sourcing, hardware reuse, and waste-heat recycling reduce lifecycle emissions [30].

### 3.3.1 Recent Advances

Additional tools such as MELODI [31], LLMCO2 [32], and EnergyMeter [33] enable more accurate energy attribution, though real-time, token-level feedback for end-users remains an open challenge [9], [10]. Research also explores prompt engineering and decoding strategies that can reduce energy consumption by up to 99% in some cases without compromising output quality [34], [35]. The past two years have seen a surge of research on the energy and carbon footprint of LLMs, with a focus on both measurement and mitigation. Key findings from recent studies include:

- **Life-cycle and System-level Solutions:** [3] identifies eight major life-cycle phases for LLM-powered chatbots, with hardware manufacturing and training as the most energy-intensive. System-level solutions include dynamic reporting, extended producer responsibility, and management standards.
- **Inference Optimization:** Frameworks like SPROUT [29] and DynamoLLM [12] demonstrate 40–60% carbon savings by controlling output verbosity, dynamic instance scaling, and GPU frequency adaptation. Workload-based models [36] and energy-aware routing [37] further optimize inference energy.

- **Measurement and Modeling:** Tools such as MELODI [31] and LLMCO2 [32] provide accurate, real-time energy monitoring and prediction, revealing that larger models can consume 100x more energy per token than smaller ones. EnergyMeter [33] and other benchmarks [28] highlight the dominant role of GPU and the impact of batch size, quantization, and model architecture.
- **Prompt Engineering and Decoding:** Prompt engineering [34], [38] and decoding strategies [35] can reduce inference energy by up to 99% in some configurations, with minimal or even positive effects on accuracy. Assisted Decoding and stochastic methods offer practical trade-offs between quality and energy use.
- **Code Generation and Use Cases:** Studies on LLM-generated code [39], [40] show that model size, architecture, and usage patterns (e.g., concurrency, streaming) significantly affect energy efficiency. Only a small fraction of generated code is actually used, highlighting waste.
- **Hardware and Deployment:** Using older GPUs in low-carbon regions [30] and extending hardware lifetimes can reduce embodied emissions. On-device generation is far less efficient than remote inference [41].
- **Gamification and Education:** Prompt-based games [42] and eco-badges can increase awareness and promote sustainable behavior.

These advances collectively demonstrate that both technical and behavioral interventions are needed for sustainable LLM deployment.

Despite these advances, technical constraints persist: Most accurate energy measurement tools are offline, coarse-grained, or limited to academic settings. As a result, platforms rarely offer meaningful energy feedback to users. Open-source libraries such as CarbonTracker [43] and CodeCarbon [44] estimate emissions post-hoc. However, live token-level feedback for chatbots remains an open challenge [9], [10].

### 3.4 Design Gap: UI Feedback, Tools, and Market Landscape

Although technical and behavioral research on sustainable LLM usage has advanced, few practical applications have reached end-users. Existing tools illustrate emerging approaches:

**ScaleDown** [45]: Displays per-prompt CO<sub>2</sub> feedback and incentivizes short, efficient prompts with a visual badge.

**AI Wattch** [46]: Adds a live energy gauge and dashboard to ChatGPT without backend access.

Both operate purely on the frontend and rely on static assumptions for energy estimation [47]. Their adoption and effectiveness have not yet been studied in peer-reviewed contexts.

UI interventions from other domains offer transferrable patterns: Gamification, comparative feedback, equivalency visualizations, and nudges have proven effective in domains like smart metering and mobile energy management. Prompt-based games and eco-badges [22], [42] exemplify how sustainability cues can be embedded without requiring infrastructure changes.

### 3.5 Conclusion and Implications

While UI-based sustainability features are well established in domains such as smart metering and mobile apps, their integration into conversational AI remains nascent. Existing research and third-party tools provide valuable design patterns and initial evidence, but several gaps persist:

- There is a lack of precise, real-time energy feedback for LLM users.
- The effective energy consumption for the major chatbots is not disclosed and there is no unanimous formula or method to calculate or predict the energy consumption of inference with such a chatbot.

- The impact of UI-only interventions on actual energy consumption and user behavior is not well understood.
- Corporate incentives and technical barriers limit the adoption of transparency features.

Our work addresses these gaps by investigating UI-only strategies to improve energy transparency and promote sustainable user behavior in conversational AI. Informed by prior research and existing tools, we introduce a user-centered approach that includes a lightweight method for estimating energy consumption suitable for UI integration. The next chapter presents the conceptual framework that underpins this solution

## 4 Conceptual Solution

### 4.1 Solution Approach

To encourage users to reduce their energy consumption, a new chatbot prototype was developed that incorporates various user interface techniques. These techniques are designed either to directly reduce the energy required for inference or to indirectly impact it by increasing users awareness of their personal energy usage.

Based on insights from our user survey, academic literature, and a market analysis, several features were defined to pursue this goal.

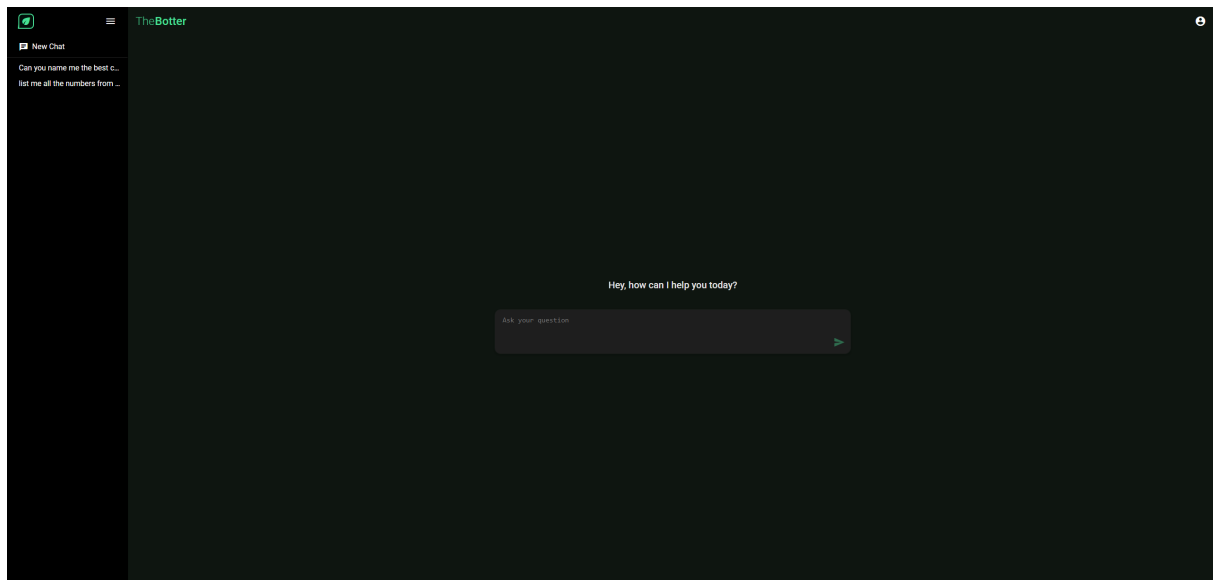
### 4.2 Baseline Chatbot

To assess the impact of energy-awareness features, it is crucial to establish a robust and widely recognized baseline. Among existing LLM-based chatbots, **ChatGPT**, developed by OpenAI, stands out as one of the most well-known and widely adopted systems worldwide at the moment of this thesis.

As detailed by Dam et al. [48], ChatGPT plays a central role in the current LLM ecosystem and is frequently cited as a benchmark for usability, response quality, and system integration. Its global prominence and maturity make it an ideal reference point for designing and evaluating alternative chatbot features.

National-level data further reinforces this decision. According to a 2024 Swiss survey by Comparis, over two-thirds of the population have already used ChatGPT or comparable tools such as Google Gemini [49]. This widespread familiarity among users ensures the relevance and realism of employing ChatGPT as the baseline system.

The baseline version of our prototype replicates the core user experience of ChatGPT, including standard features such as deleting conversations and editing conversation titles and the design.



**Figure 4.1:** Baseline ChatGPT-inspired chatbot without sustainability features

### 4.3 Features and Functionalities

To promote more sustainable usage of large language models, this work explores a variety of user interface features aimed at increasing awareness of energy consumption and influencing user behavior.

The proposed features were derived through a combination of methods: A user survey conducted during the research (see Section 3.2), analysis of related academic literature, and brainstorming informed by behavioral design principles and digital sustainability research.

This section first categorizes all brainstormed features into three conceptual groups: (1) technical configuration features, (2) awareness and behavioral nudges, and (3) feedback, limits, and gamification mechanisms. Afterwards, a final subset of selected features for implementation is presented.

#### 4.3.1 Technical Configuration Features

These features shown in the table Table 4.1 focus on reducing energy consumption through technical adjustments that are exposed to the user via the UI. They allow the system to use fewer computational resources by choosing more efficient processing strategies or limiting input complexity.

Feature	Description
Eco-Mode Toggle	Switch to a smaller model or activate multiple eco-features at once
Model Selector	Choose between performance and efficiency focused models
Ignore Context Option	Skip conversation history to reduce token processing
Alternative Tool Suggestion	Redirect queries to efficient tools like Google or DeepL
Document Splitting Warning	Ask users whether only parts of large uploads should be processed instead of defaulting to the whole document always

**Table 4.1:** Technical configuration features

#### 4.3.2 Awareness and Behavioral Nudges

This category includes features as shown in Table 4.2 that aim to influence users through information, education, or subtle nudges during interaction. The goal is to increase awareness of energy usage and foster more deliberate, sustainable usage patterns in general.

Feature	Description
Prompt Energy Prediction	Estimate energy consumption before sending a message [50]
Hints	Display short energy-saving tips and nudges while chatting [51]
Prompt Suggestions	Recommend more efficient or specific prompts [51]
Footprint Awareness Tips	Relate usage to real-world analogies (e.g. phone charges) [50]
Onboarding Tutorial	Show introductory guide to raise awareness before first use

**Table 4.2:** Awareness and behavioral features

#### 4.3.3 Feedback, Limits, and Gamification

These features in Table 4.3 provide real-time or retrospective feedback, integrate social comparison mechanisms, and apply gamified incentives. They are designed to reinforce sustainable behavior through motivation, reflection, and competition.

Feature	Description
Usage Meter	Track energy usage per message, session, or time period [50]
Saving Meter	Show energy saved by using efficiency features [50]
Personal Limits	Allow users to set usage caps or warning thresholds [51]
Eco-Score	Assign a sustainability score based on user behavior
User Ranking	Compare individual usage with anonymized peer data [21]
Challenges	Set daily, weekly, or monthly goals for eco-friendly use
Fee System	Introduce symbolic cost to reflect energy consumption

**Table 4.3:** Gamification and feedback mechanisms

#### 4.3.4 Decision Matrix - Overview of Features

The following table Table 4.4 shows a simplified scoring overview of all features we had in mind and their scores/priority. Each feature was scored based on its expected impact (awareness and energy), research, survey support and implementation complexity (note: The scores are speculative):

Feature	Awareness	Energy	Survey	Research	Complexity	Note	Total
Eco-Mode Toggle	3	3	1	1	-2	[1]	6
Alternative Tool Suggestion	2	2	1	1	-1	[2]	5
Usage Meter	3	1	1	1	-1		5
Footprint Awareness Tips	3	1	1	1	-2		4
Prompt Suggestions	2	2	0	1	-1		4
Prompt Prediction	3	1	0	1	-2	[3]	3
Document Splitting Warning	2	2	1	0	-2	[4]	3
User Ranking	2	1	1	1	-2		3
Hints	2	1	1	1	-2		3
Model Selector	2	2	0	0	-2	[5]	2
Saving Meter	2	1	0	1	-2		2
Personal Limits	2	1	0	1	-2		2
Challenges	2	1	0	1	-2		2
Onboarding Tutorial	2	1	1	0	-2		2
Eco-Score	2	1	0	0	-2		1
Ignore Context Option	1	2	0	0	-3		0
Fee System	1	1	0	0	-3		-1

**Table 4.4:** Decision matrix of all evaluated UI features

**Note:**

1. Combines model selector; intuitive, low-effort toggle.
2. Idea seems good, although extensive analysis of the prompt might add an overhead.
3. Prediction/estimation of the actual consumption or some sort of measure is a critical feature for many others including statistical features, ranking or fee.
4. Comes with overhead effort since document prompt input must be supported by the bot.
5. Can be combined with the mode toggle; may be redundant as a separate feature.

Based on the decision matrix and a feasibility analysis, five features were selected and refined to form a coherent and complementary set. These features were chosen due to their expected high impact, their contribution to the effectiveness of other selected features and their alignment with user expectations as identified in the survey. Feasibility within the scope of this thesis was also a key consideration.

## 4.4 Definitive Features

In the following sections, we will describe each of the final chosen features in detail.

### 4.4.1 Three-Mode Switch (Eco-Mode toggle)

The Three-Mode Switch is a central UI element that allows users to consciously select between three energy profiles before sending a prompt. It builds upon the idea of an eco-mode toggle and integrates functionality from a previously proposed model selector. It also includes restrictions about the included history length. This feature enables users to balance environmental impact against output quality in an intuitive way.

Each of the three selectable modes corresponds to a different underlying model and restrictions concerning history length.

- **Energy-Efficient Mode:** Utilizes *GPT-4.1-nano* (1.1 billion parameters), limited to the last 2 messages of context. Ideal for simple or repetitive queries.
- **Balanced Mode:** Uses *GPT-4o-mini* (8 billion parameters), with the last 5 messages of context. Offers a compromise between resource usage and conversational depth.
- **Performance Mode:** Employs *GPT-4o* (175 billion parameters), including the last 10 messages of context. Suited for complex reasoning or in-depth queries.

Unlike systems that auto-switch based on prompt type, this implementation leaves the decision entirely to the user, ensuring transparency and control. The switch appears inline during prompting and can be adjusted for each individual message.

### 4.4.2 Metrics Dashboard

The Metrics Dashboard is a central feature designed to enhance user awareness of the environmental footprint associated with their chatbot usage. It consolidates key energy and behavior related statistics into a unified visual interface. This feature responds directly to survey feedback indicating a strong desire for transparency and supports reflection on usage patterns.

#### Purpose and Motivation

The dashboard serves a dual purpose: (1) to increase awareness of energy consumption through quantified feedback and (2) to support behavior change through actionable insights and motivational design. Research has shown that eco-feedback systems can lead to energy savings of 4–5% on average [19], and that real-time feedback can enhance these effects by an additional 5% [14]. To increase user engagement, the system also integrates elements of gamification such as emojis, color-coded badges and progress indicators, an approach shown to positively influence sustainable behavior [22].

#### Displayed Metrics

The dashboard provides live and historical insights into user behavior, with metrics structured around energy tracking, token-level statistics and usage distribution across different modes. These metrics are visualized via bar charts, line graphs, and text summaries, and are computed from both client-side and backend-collected data. Table 4.5 describes each metric and its intended user-facing function.



No.	Metric	Purpose
1	Total energy usage	Provides a high-level perspective on cumulative energy consumption over time.
2	Total number of conversations	Indicates overall usage intensity, including deleted entries.
3	Total number of prompts	Helps users understand the frequency of input, beyond just the number of conversations.
4	Total input tokens	Directly correlated with energy usage; reflects the user's prompt length.
5	Total output tokens	Complements input tokens; both are used in energy calculations.
6	Delta compared to yesterday	Displays day-over-day change in energy consumption, reinforced by gamified emoji to promote behavior change [22].
7	Daily energy usage	Visualizes trends in energy consumption over time.
8	Usage per chat mode	Helps users identify which modes contribute most to their energy footprint.
9	Prompts per chat mode	Highlights the distribution of usage across different LLM configurations.
10	Prompts per day	Tracks how frequently users engage with the system.
11	Average usage per prompt	Designed to support prompt optimization by showing the typical energy cost per message.
12	Tokens per day (input/output)	Provides token-level insight into daily usage behavior.

**Table 4.5:** Overview of metrics shown on the user dashboard

### User Interface of the Metrics Dashboard

The dashboard is accessible via the navigation drawer and is available at any time during a session. The layout prioritizes simplicity and readability: Key metrics are displayed numerically, supported by trend indicators and visual charts. Comparative badges and emoji feedback are used sparingly to avoid overload, while still encouraging reflection and sustainable engagement.

#### 4.4.3 Prompt Prediction

Prompt prediction is a key anticipatory feedback mechanism that estimates the energy consumption of a prompt before it is sent. This feature directly supports the broader goals of the Metrics Dashboard by enabling more reflective and energy-aware behavior at the point of decision-making. Rather than relying solely on post hoc analytics, users are provided with information that helps them evaluate the impact of their actions *before* the energy is consumed.

This aligns with findings from behavioral science which suggest that predictive or real-time eco-feedback is more effective than delayed feedback in promoting sustainable behavior [14]. By surfacing environmental costs early in the interaction, the system empowers users to rephrase or reconsider their input whether by simplifying, shortening, or selecting a more energy-efficient mode.

### Functionality

While text is entered in the prompt input (but before the submission), the system constantly calculates and displays an estimated energy cost based on:

- The selected mode (Energy-Efficient, Balanced, Performance)
- The input token count (estimated via pre-tokenization)

- Expected output token length

This estimate is visualized for all chat modes in real time as the user types, which enables the users to compare the modes for their specific prompt. The calculations and the energy model coefficients used in the estimation process will be described in subsection 4.5.

### **Benefits and Integration**

Prompt prediction not only enhances transparency but also enables the user to take proactive steps toward lower-impact usage. The feature is tightly integrated with both the Three-Mode Switch (for dynamic updates based on model selection) and the Metrics Dashboard (which records both predicted and actual energy usage for comparison). Over time, this may help users build an intuition for the cost-effectiveness of different prompt styles.

#### **4.4.4 Energy-Note (Per-Response Footprint)**

The Energy-Note feature provides immediate, post-interaction feedback about the energy consumed to generate a specific response. In contrast to the Prompt Prediction (which estimates energy usage before sending), the Energy-Note confirms and reinforces the environmental cost after the system has generated an output.

### **Purpose**

The goal of the Energy-Note is to foster reflection through concrete, real-time feedback. After each assistant message, users are shown the exact energy cost for the current exchange, expressed both in watt-hours (Wh) and in a relatable real-world comparison further explained in 4.4.5. This reinforces the connection between digital interactions and physical resource usage.

Such fine-grained, moment-of-use feedback is designed to make energy consumption more tangible and personal, encouraging behavior change over time. The feature complements cumulative metrics in the dashboard and enhances user understanding of the impact of individual prompts and shows the impact of the history length to the energy consumption of a prompt.

### **User Interface of the Energy-Note**

The Energy-Note appears below each assistant message in the form of a compact text element, styled subtly to maintain visual hierarchy. A tooltip or expandable section optionally provides additional context or references. This lightweight integration ensures the information is visible but not intrusive, maintaining conversational flow while offering transparency.

#### **4.4.5 Energy Analogies**

This feature originates from the concept of “Footprint Awareness Tips” and aims to make energy consumption more relatable by expressing watt-hour values in terms of everyday devices and scenarios. Pure numeric feedback in Wh may lack intuitive meaning for most users; thus, this feature provides personalized analogies to improve comprehension and emotional resonance.

It integrates with other features such as Prompt Prediction and Energy-Note to display real-time feedback using the user’s selected analogy unit. It is also embedded into the onboarding flow to encourage early reflection and engagement.

### **Design and Workflow**

To find an analogy the user can relate to a short on-boarding process which is done by every user after the registration. During this process users are presented with several real-world energy anchors and are asked to estimate the energy cost of each. These initial guesses are used to determine which units users understand best. After the estimation users are presented with the actual consumption in Wh and their deviation in percentage, allowing them to pick the one that feels most intuitive.

### Available Analogy Units

The following energy anchors are provided during onboarding and remain available for later customization:

- **iPhone 14 full charge** (12.68 Wh) [52]
- **One minute powering a refrigerator** (0.21 Wh) [53]
- **One minute working on a laptop** (0.75 Wh) [54]
- **One hour of a 500-lumen LED lamp** (6 Wh) [55]
- **One minute of PlayStation 5 gaming** (3.5 Wh) [56]

Once selected, this energy unit is used consistently throughout the whole app wherever energy feedback is shown, including next to each response (Energy-Note), in the Prompt Prediction label, and in the Metrics Dashboard summaries. This consistency ensures that users develop an intuitive mental model of their cumulative impact.

Users can change their analogy unit at any time via the settings menu, allowing them to adjust the framing as their preferences evolve.

## 4.5 Estimation Model of Energy Consumption

In order to calculate and predict the energy consumption of prompts we built a formula based on the current research and existing benchmarks. To quantify the energy consumption  $E$  (in watt-hours, Wh) for a single prompt–response interaction, we adopt the following linear model:

$$E = \alpha \cdot T_{\text{in}} + \beta \cdot T_{\text{out}} + \zeta$$

where  $T_{\text{in}}$  and  $T_{\text{out}}$  are the number of input and output tokens respectively. The coefficients  $\alpha$  and  $\beta$  (in Wh/token) represent the energy cost per token, while  $\zeta$  is a fixed overhead cost per request accounting for memory allocation, logging, and general infrastructure load.

### 4.5.1 Empirical Justification

Recent research confirms the viability of token-based energy estimation. Poddar et al. highlight a strong correlation between token count and inference energy [9], while Fernandez et al. criticize FLOP-based approaches and advocate for token-level modeling instead [10]. These results support the use of static coefficients  $\alpha$ ,  $\beta$ , and  $\zeta$ .

Several empirical studies have shown that the energy cost of generating output tokens significantly exceeds that of processing input tokens. This is primarily due to the autoregressive nature of large language models, where each output token requires sequential computation and additional memory operations. Hardware-level profiling and benchmark studies consistently report that output tokens consume between 4 and 5 times more energy than input tokens [57].

OpenAI’s own API pricing reflects this asymmetry: Output tokens are generally priced at 4–5× the rate of input tokens. For instance, GPT-4o charges \$5.00/M (M = millions) input tokens vs. \$20.00/M output tokens which is a 4× ratio. [58] Based on both empirical findings and pricing structure, we adopt a fixed scaling factor of:

$$\beta = 4\alpha \tag{4.1}$$

This conservative ratio balances observed energy benchmarks and commercial token economics, while simplifying model calibration.

### 4.5.2 Scaling Across Modes

To generalize the model across chat modes, we use OpenAI’s official pricing as a proxy for computational cost. Table 4.6 shows the cost per million tokens and derived scaling ratios.

Model	Input Price (USD)	Output Price (USD)	Input Ratio	Output Ratio
GPT-4o	\$5.00 /M	\$20.00 /M	1×	1×
GPT-4o-mini	\$0.60 /M	\$2.40 /M	0.12×	0.12×
GPT-4.1-nano	\$0.10 /M	\$0.40 /M	0.02×	0.02×

**Table 4.6:** Model pricing per million tokens and cost-based scaling

### 4.5.3 Reported Energy Benchmarks

Table 4.7 lists published values for energy use per query, showing consistent differences in efficiency across models.

Source	Model	Energy (Wh)	Token Length Assumption
EpochAI (2024) [59]	GPT-4o	0.30	“Average ChatGPT query” (unspecified token count)
Jegham et al. (2025) [57]	GPT-4o	0.43	“Short prompt” (no token details)
Jegham et al. (2025) [57]	GPT-4.1-nano	$\approx 70\times$ less than GPT-4o	Relative comparison for long prompt
Altman (2025) [24]	GPT-4o	0.34	No token count given

**Table 4.7:** Reported energy consumption per query and per token

### 4.5.4 Assumptions

To estimate  $\alpha$  and  $\beta$  from a single interaction, we assume a typical short chat prompt with  $T_{\text{in}} = 150$  input tokens and  $T_{\text{out}} = 300$  output tokens. This is a typical assumption also made in prior benchmarks [57], [59]. While OpenAI’s tokenizer documentation [60] provides illustrative examples, it does not specify statistically representative token lengths.

While exact distributions vary, multiple benchmark studies report typical ChatGPT interactions with input lengths ranging from 100–200 tokens and outputs often between 200–400 tokens. The chosen values reflect a balanced, realistic estimate for a brief user prompt and a moderately long model-generated response [57], [59].

The fixed overhead  $\zeta = 0.02$  Wh accounts for infrastructure-related energy use that occurs independently of token processing, such as authentication, API gateway routing, logging, and request serialization. These operations introduce a baseline energy cost per query that does not scale with token count or model size.

We select this value based on empirical calibration: For example, EpochAI reports 0.30 Wh per GPT-4o query [59], making 0.02 Wh a plausible 7% share. Additionally, Jegham et al. emphasize that total query energy includes non-model components like orchestration and monitoring [57], further supporting a fixed infrastructure cost in end-to-end modeling.

This fixed overhead avoids underestimating energy use in short prompts and lightweight models while preserving a simple, transparent estimation framework. As more information about the underlying system and infrastructure is exposed this coefficients could be fine-tuned.

### 4.5.5 Coefficient Estimation

We derive the coefficients  $\alpha$ ,  $\beta$ , and  $\zeta$  using two independent benchmarks for GPT-4o: The public estimate from EpochAI [59] and the infrastructure-aware analysis by Jegham et al. [57]. We assume  $T_{\text{in}} = 150$ ,  $T_{\text{out}} = 300$ , and a fixed overhead  $\zeta = 0.02$  Wh. The total energy is modeled as:

$$E = \alpha \cdot T_{\text{in}} + \beta \cdot T_{\text{out}} + \zeta, \quad \text{with } \beta = 4\alpha$$

**EpochAI** reports  $E = 0.30$  Wh. Solving:

$$\begin{aligned} 0.30 &= \alpha(150 + 4 \cdot 300) + 0.02 \Rightarrow \alpha = \frac{0.28}{1350} \approx 0.000207 \text{ Wh/token} \\ &\Rightarrow \alpha \approx 0.21 \text{ mWh/token}, \quad \beta = 4\alpha \approx 0.83 \text{ mWh/token} \end{aligned}$$

**Jegham et al.** report  $E = 0.43$  Wh. Solving:

$$\begin{aligned} 0.43 &= \alpha(150 + 4 \cdot 300) + 0.02 \Rightarrow \alpha = \frac{0.41}{1350} \approx 0.000304 \text{ Wh/token} \\ &\Rightarrow \alpha \approx 0.30 \text{ mWh/token}, \quad \beta = 4\alpha \approx 1.22 \text{ mWh/token} \end{aligned}$$

Method	$\alpha$ (mWh/token)	$\beta$ (mWh/token)	$\zeta$ (Wh)
EpochAI	0.21	0.83	0.02
Jegham	0.30	1.22	0.02

**Table 4.8:** Comparison of estimated coefficients for GPT-4o

We present both estimates to illustrate the range of plausible energy use under different assumptions. While Jegham et al.’s value reflects a more infrastructure-aware upper bound, we adopt the coefficients derived from EpochAI for several reasons: They align better with short, interactive usage typical in chat applications, rely on transparent and reproducible assumptions, and avoid potential overestimation due to infrastructure inefficiencies. Moreover, the resulting model remains simple and interpretable for end-user communication.

$$\alpha = 0.21 \text{ mWh/token}, \quad \beta = 0.83 \text{ mWh/token}, \quad \zeta = 0.02 \text{ Wh}$$

### 4.5.6 Final Coefficients per Mode

We scale the token-based coefficients  $\alpha$  and  $\beta$  across chat modes using the pricing-derived ratios shown in Table 4.6. However, the fixed overhead  $\zeta$  is assumed to remain constant across all modes, reflecting infrastructure-level costs such as authentication, logging, and request routing that do not scale proportionally with model size. This yields the final coefficients shown in Table 4.9.

Mode	Model Name	Scaling	$\alpha$ (Wh/token)	$\beta = 4\alpha$ (Wh/token)	$\zeta$ (Wh)
Performance	gpt-4o	1.00×	0.00021	0.00083	0.020
Balanced	gpt-4o-mini	0.12×	0.0000252	0.0001008	0.020
Energy Efficient	gpt-4.1-nano	0.02×	0.0000042	0.0000168	0.020

**Table 4.9:** Final energy model coefficients per mode

#### 4.5.7 Predicting Output Tokens

To enable energy estimation prior to execution, the number of output tokens must be predicted from the input. While actual output length depends on prompt complexity, model behavior and context, studies show that outputs of large language models (LLMs) like ChatGPT are typically at least as long as their inputs and often even longer especially in open-ended tasks such as summarization or explanation [58], [61]. To interpret the benchmarks we assumed a input-output-ratio of 150/300 respectively 150/500, however those values were chosen under the assumption that inference does not include previous context. To account for a conversational scenario, we define the output token count as:

$$T_{\text{out}} = T_{\text{in}}$$

This identity function avoids overestimation and makes it simple. In multi-turn interactions, users frequently re-prompt or refine previous queries, causing the cumulative input (i.e., context/history) to grow. Since the frequency and extent of such re-prompting vary greatly and are hard to predict, assuming  $T_{\text{out}} = T_{\text{in}}$  provides a balanced and stable estimate. Prior heuristics using  $2\text{--}3\times$  the input may overfit to specific scenarios and lack the ability to be generalized.

## 5 Implementation

In this chapter we present the detailed composition and implementation of the prototype.<sup>1 2</sup>

### 5.1 System Architecture

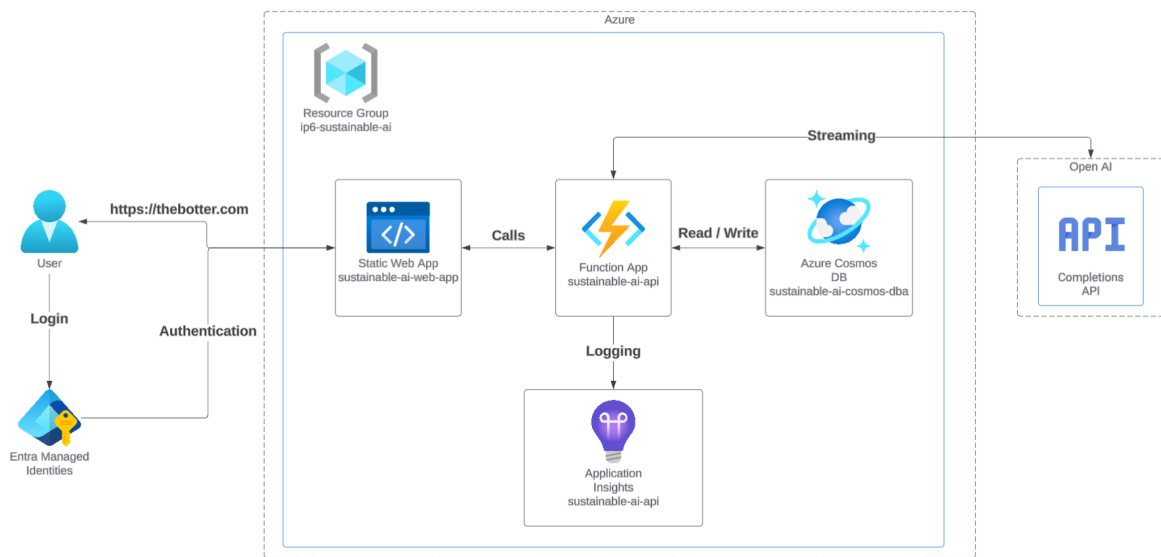
The prototype is designed as a web application, accessible at <https://thebotter.com>, and implemented as a cloud-native system, fully hosted and scalable within the Microsoft Azure Cloud [62]. Continuous integration and deployment (CI/CD) are realized via GitHub Actions [63], ensuring seamless and automated deployments.

The frontend is a single-page application (SPA) [64] deployed via Azure Static Web Apps [65]. The backend is implemented using serverless Azure Functions [66]. Backend logic accesses Azure Cosmos DB Containers for storage and interacts with the OpenAI API [67]. User authentication is handled via Microsoft Entra ID using Easy Auth [68].

All user data is fully partitioned to ensure isolation between tenants. Administrative monitoring and logging are implemented using Azure Application Insights [69]. Manual logging is additionally implemented for experimental evaluation.

Two environments are available: A production environment for the experiment and a development environment for local testing.

The following figure gives a full overview of all the components and their interaction.



**Figure 5.1:** Application architecture overview

### 5.2 Frontend

The frontend is implemented as a single-page application (SPA) using Angular 19 [70], following the official Angular Style Guide and deployed using Azure Static Web Apps. The application is optimized for desktop use and adopts a green and black color scheme to visually reinforce the sustainability theme, with black also reducing power draw on OLED displays. [71]

<sup>1</sup>Frontend Repository: [github.com/simonluescherfhnw/ip6-sustainable-ai-frontend](https://github.com/simonluescherfhnw/ip6-sustainable-ai-frontend)

<sup>2</sup>Backend Repository: [github.com/simonluescherfhnw/ip6-sustainable-ai-backend](https://github.com/simonluescherfhnw/ip6-sustainable-ai-backend)

All user routes are protected by Angular route guards. Authentication is handled via Azure Static Web Apps using Microsoft Entra ID (formerly Azure Active Directory). Authenticated user information is available via the `/.auth/me` endpoint, with login and logout triggered through `/.auth/login/aad` and `/.auth/logout` respectively [72].

### 5.2.1 Libraries and Packages

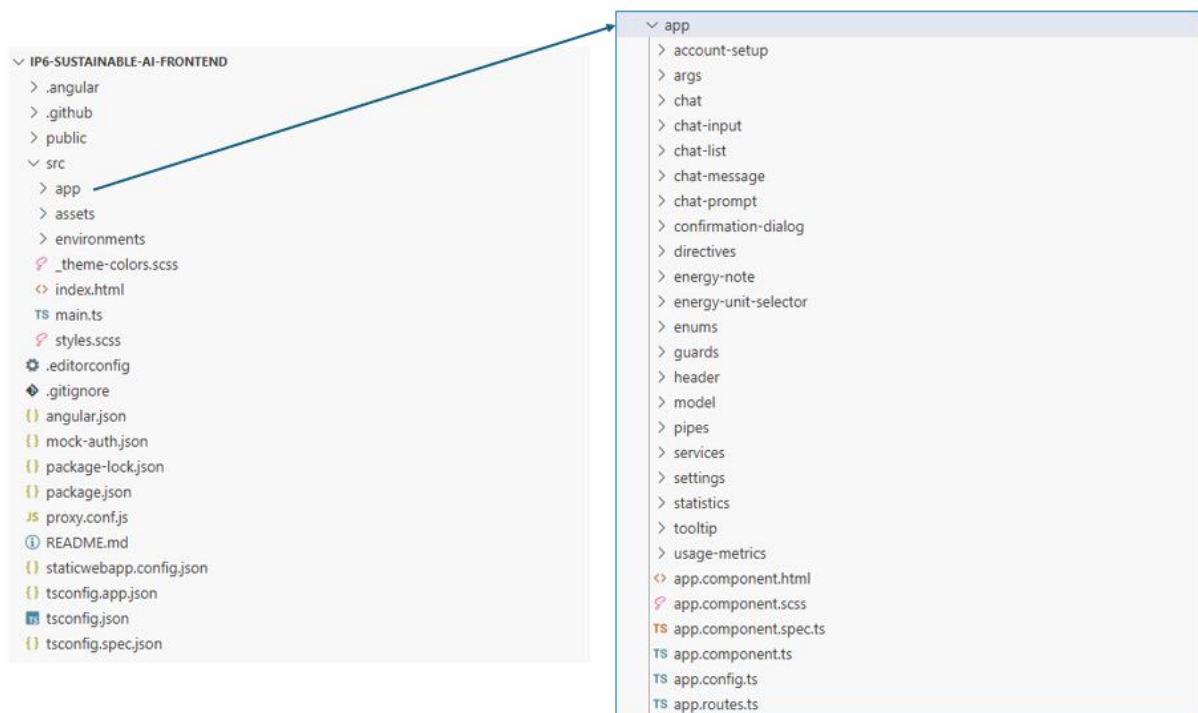
Key libraries as seen in Table 5.1 used in the frontend include Angular Material for UI components and Chart.js for data visualization:

Package	Version	Description
@angular/*	19.1.0	Core Angular packages [73]
@angular/material	19.1.0	Used as component library and for theming [74]
chartjs	4.4.9	Rendering graphs for the metrics pages [75]
marked	15.0.12	Converting markdown chat responses to HTML [76]
rxjs	7.8.0	Implements observer pattern in JavaScript, used for state management and interaction with Angular reactive forms [77]
jasmine-core	5.5.0	Jasmine test framework core library
karma	6.4.0	Test runner for Angular and Jasmine

**Table 5.1:** All important frontend packages

### 5.2.2 Code Structure

Figure 5.2 shows the folder structure of the frontend. The UI components are all directly in the app folder and split for separation of concerns.



**Figure 5.2:** Frontend project structure



### 5.2.3 Routing

There are four main areas of the web application accessible via navigation. `/` or `/chat` and `/chat/id` are redirecting the user to the main area, the chat section. `/metrics` shows the user's dashboard and `/settings` the settings page. After registrations users get redirected to `/account-setup` where the account initialization process starts. Figure 5.3 shows the detailed configuration of the routes in the Angular project.

```
export const routes: Routes = [
  { path: '', redirectTo: '/chat', pathMatch: 'full' },
  { path: 'chat', redirectTo: '/chat/', pathMatch: 'full' },
  { path: 'chat/:id', component: ChatComponent, canActivate: [authGuard] },
  { path: 'metrics', component: UsageMetricsComponent, canActivate: [authGuard] },
  { path: 'settings', component: SettingsComponent, canActivate: [authGuard] },
  { path: 'account-setup', component: AccountSetupComponent, canActivate: [authGuard] }
];
```

Figure 5.3: app.routes.ts

### 5.2.4 API Connection

On the productive environment the static web app is connected to the azure functions backend directly via configuration. It is setup so that all calls to `/api` are redirected to the functions url. This also makes it possible to automatically authenticate the user for the api and attach the user claims to the request headers automatically [72]. To replicate this behavior when developing in a local environment some additional steps are necessary.

A proxy (Figure 5.4) is redirecting all requests to `/api` to `http://localhost:7028/api` where the functions are running locally. Additionally, the `x-ms-client-principal` header must be set manually for every request since it won't be added automatically.

```

module.exports = {
  "/.auth/me": {
    "target": "http://localhost:4200",
    "secure": false,
    "bypass": function (req, res, proxyOptions) {
      if (req.url === "/.auth/me") {
        res.setHeader("Content-Type", "application/json");
        res.end(require("fs").readFileSync("mock-auth.json"));
        return true;
      }
    }
  },
  '/api': {
    target: 'http://localhost:7028/api', // Azure Function runtime
    changeOrigin: true,
    logLevel: 'debug',
    pathRewrite: {'^/api': ''},
    onProxyReq: (proxyReq, req, res) => {
      // Forward the auth header if it exists
      const principal = req.headers['x-ms-client-principal'];
      if (principal) {
        proxyReq.setHeader('x-ms-client-principal', principal);
      }
    }
  }
}

```

Figure 5.4: Proxy configuration file

## 5.3 Backend

The backend for the prototype is implemented using Azure Functions in C#. This decision was made based on Azure Functions scalability, ease of integration with other Azure services and the serverless nature of the framework which eliminates the need for infrastructure management. The backend communicates with the Angular frontend via a clean HTTP interface, is testable both locally and through the Azure Portal and integrates seamlessly with GitHub, Cosmos DB and Application Insights for logging and monitoring.

### 5.3.1 Functions Overview

Azure Functions offer a lightweight and scalable way to build backend logic without managing dedicated servers. In the context of a time-boxed student project, this offered several advantages:

- **No server management:** Reduces setup and DevOps effort.
- **Auto-scaling:** Adjusts to load during the experiment and reduces costs when idle.
- **Built-in GitHub integration:** Allows fast deployment with CI/CD pipelines.
- **Easy to test:** Each function can be tested in isolation via HTTP-Client or the Azure Portal.

The following HTTP-triggered functions in Table 5.2 define the main capabilities of the backend:

FunctionName	Params	Returns	Description
GetAppData	None	DtoAppData	Returns base app data including energy units and mode configurations.
GetConversations	GetConversationsArgs	ICollection <DtoConversation>	Returns user's conversations.
GetPrompts	GetPromptsArgs	ICollection <DtoPrompt>	Returns prompts in a conversation.
GetUsageStatistics	GetUsageStatisticsArgs	DtoStatistics	Returns user's statistics for the last week.
UpdateConversation	UpdateConversationArgs	DtoConversation	Updates conversation name.
UpdateUser	UpdateUserArgs	DtoUser	Updates user preferences.
DeleteConversation	DeleteConversationArgs	None	Soft-deletes a conversation.
LogPageVisit	LogPageVisitArgs	None	Logs visited page.
PredictPromptUsage	PredictPromptUsageArgs	DtoPrediction	Predicts energy usage for a prompt.
SendPrompt	SendPromptArgs	IAsyncEnumerable <DtoMessageResponsePart>	Sends prompt and returns model response stream.

**Table 5.2:** Overview of implemented backend Azure Functions

### 5.3.2 Code Structure and Modularity

## 5.4 Solution Structure

The solution follows a modular folder structure that aligns with best practices for organizing .NET applications as recommended by Microsoft [78], [79]. This organization supports maintainability, testability, and separation of concerns.

```
IP6-Sustainable-AI/
|-- .gitignore
|-- README.md
|-- SustainableAI.sln
|-- .github/
|   |-- workflows/
|       |-- main_sustainable-ai-api.yml
|-- .vscode/
|-- SustainableAI.Api/
|   |-- chatmode.configuration.json
|   |-- EnvironmentHelpers.cs
|   |-- host.json
|   |-- local.settings.json
|   |-- Program.cs
|   |-- SustainableAI.Api.csproj
|   |-- Common/
|   |-- Configuration/
|   |-- Data/
|   |-- Functions/
|   |-- Properties/
|   |-- Repositories/
```

```
|   |-- Service/
|-- SustainableAI.Api.Tests/
|   |-- SustainableAI.ApiTests.csproj
|   |-- TestData.cs
|   |-- Mock/
|   |-- UnitTests/
```

### Key Components:

- `SustainableAI.sln`: The main solution file that aggregates all project references.
- `.github/workflows/`: Contains GitHub Actions CI/CD pipelines for automated testing and deployment.

**SustainableAI.Api/** is the core Azure Functions project and contains:

- Entry and configuration files like `Program.cs`, `host.json`, and `local.settings.json`.
- `chatmode.configuration.json`: Custom configuration to define available chat modes.
- `Common/`, `Configuration/`, `Data/`, `Repositories/`, `Service/`: Represent logical layers of the application, adhering to the Clean Architecture pattern.
- `Functions/`: Contains Azure Function entry points representing external APIs.

**SustainableAI.Api.Tests/** is a dedicated project for automated tests:

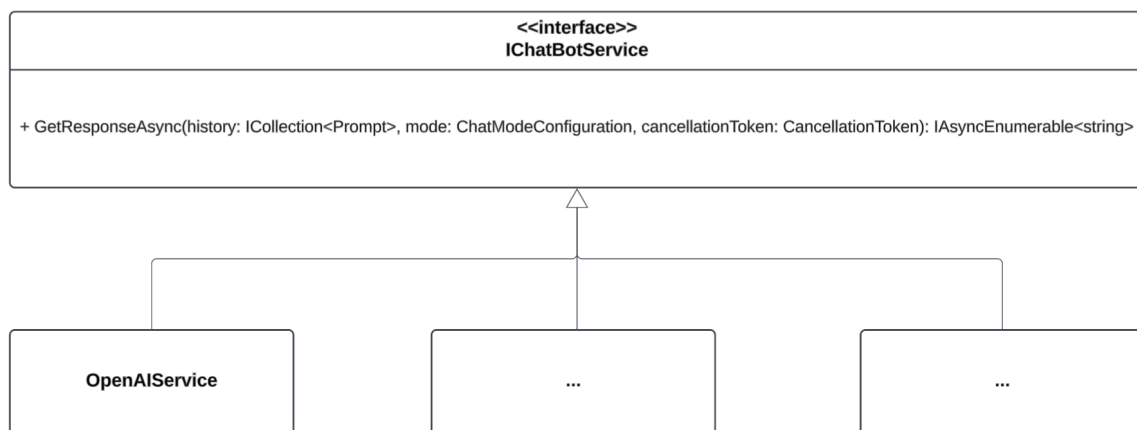
- `Mock/`: Provides mocked dependencies for test isolation.
- `UnitTests/`: Contains unit tests for individual services, repositories, or functions.
- `TestData.cs`: Supplies reusable test data objects.

This structure enables clear boundaries between infrastructure, business logic, and API endpoints, improving readability and testability of the solution.

#### 5.4.1 OpenAI Integration

To replicate ChatGPT-like behavior, the backend interacts with OpenAI's chat completion API [67]. An abstraction layer is implemented via the `IChatBotService` interface to decouple model-specific logic from the core application.

- Decoupling from the vendor: The backend could be switched to e.g. Azure OpenAI or open models.
- Streaming [80] responses using server-sent events (SSE) [81], improving UX with incremental output.



**Figure 5.5:** Interface abstraction for OpenAI integration

## 5.5 Feature Implementation

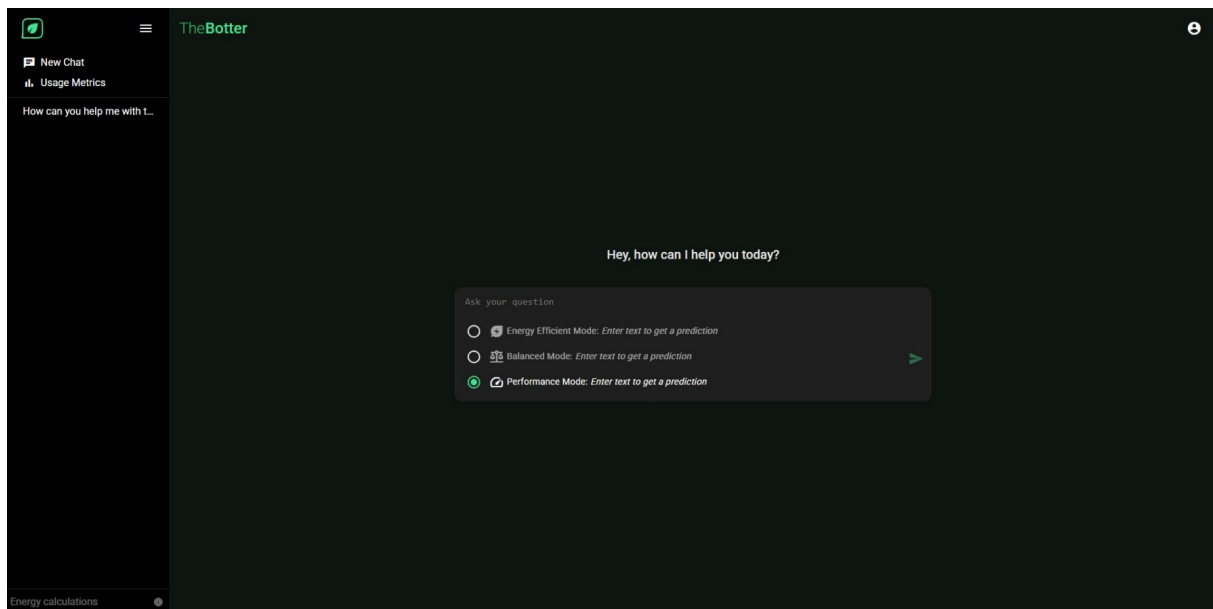
In the following sections we will explain technical implementation points of our core features.

### 5.5.1 Baseline Chatbot Features

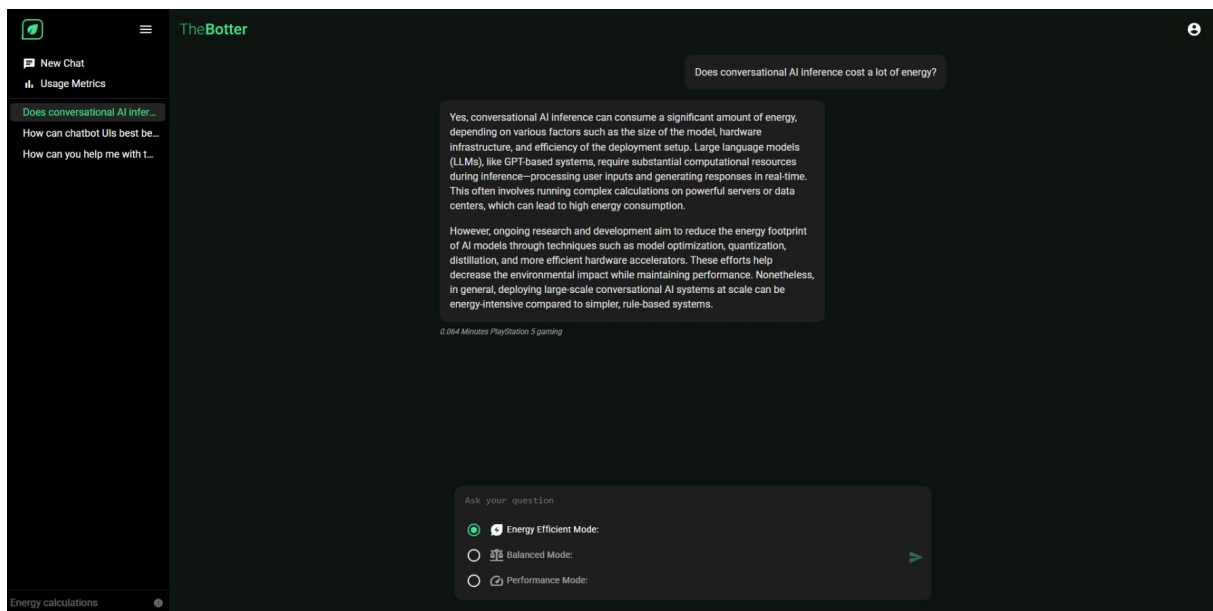
The prototype replicates core ChatGPT-like functionality to serve as a reference system for measuring the effect of sustainability features. These baseline capabilities include starting new chats, renaming, and deleting existing conversations as well as prompting in existing chats or just start a new chat with a new prompt.

Users begin each session on the main chat screen (Figure 5.6) where they can initiate a new conversation. Previously saved conversations can be selected and revisited (Figure 5.7).

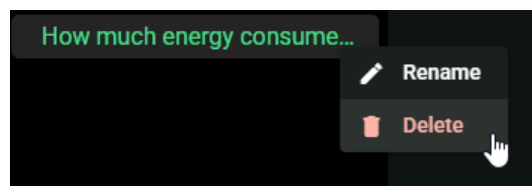
Each conversation includes a contextual menu that appears on hover, enabling users to rename or delete it (Figure 5.8). These changes are reflected in the Cosmos DB backend via the `UpdateConversation` and `DeleteConversation` functions.



**Figure 5.6:** Users opening the application will land on /chat and be able to start a new dialog instantly



**Figure 5.7:** Previous conversations can be revisited to continue prompting or collect information

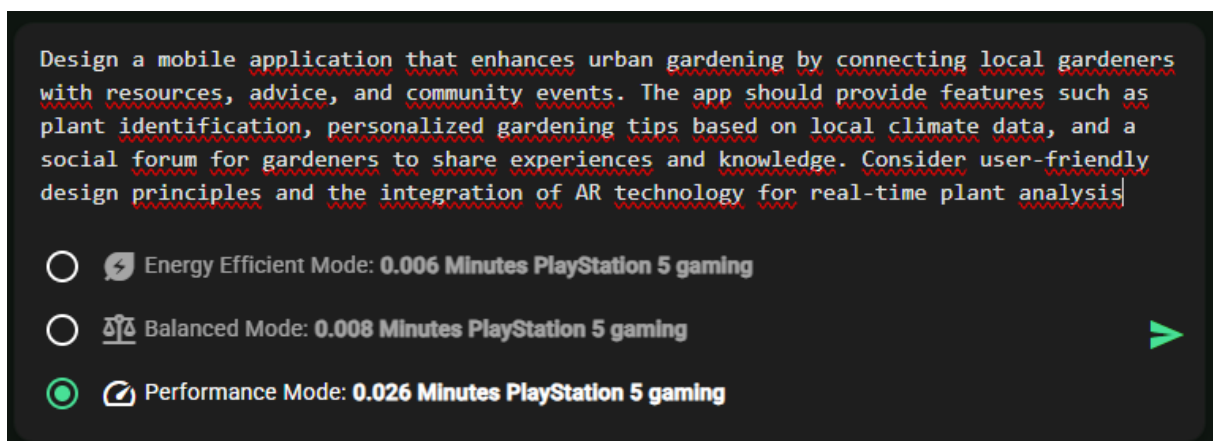


**Figure 5.8:** Users can delete or rename conversations via the menu icon

### 5.5.2 Three-Mode Switch

Implemented in the Angular chat component, the mode switch allows users to select between "Energy-Efficient", "Balanced", and "Performance" modes. Each mode corresponds to a model preset and context window configuration. Mode changes are reflected in the frontend UI and stored per-prompt (see Figure 5.9).

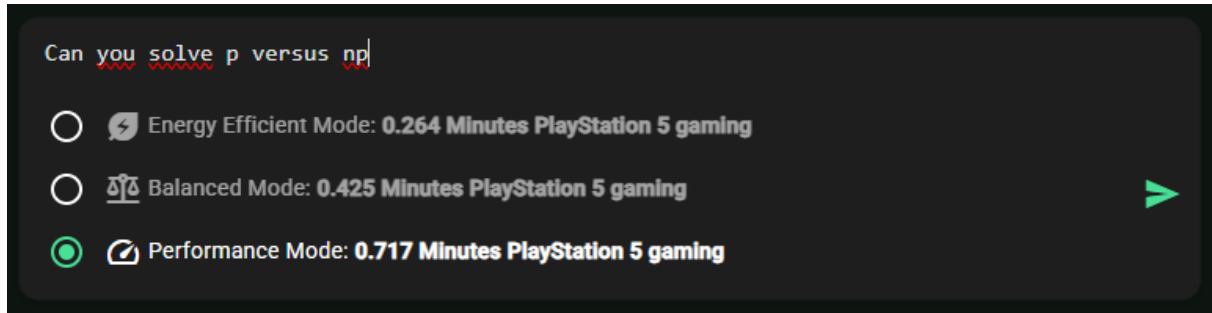
The configuration is loaded from the backend via `GetAppData` and used for both prompt sending and prediction logic.



**Figure 5.9:** Mode can be switched for every prompt

### 5.5.3 Prompt Prediction

Prompt prediction is implemented via reactive data binding in the chat input component. As the user types, the frontend triggers a call to the backend function `PredictPromptUsage`, which estimates the energy consumption (in Wh) based on the input token length across all available chat modes. The resulting prediction is visualized in real time within the UI, as shown in Figure 5.10.



**Figure 5.10:** Real-time energy prediction for current prompt

The underlying estimation logic resides in the `CalculationService` class, which applies the energy estimation model introduced in Section 4.5. This service computes predicted usage by applying a linear function to the input and predicted output token counts, using coefficients provided by a corresponding `ChatModeConfiguration`. The configuration files are read from an environment variable and can be modified without redeploying the backend [82].

When a user enters a prompt, the process begins with a call to the `PredictOutputTokens` method (Figure 5.12), which estimates the number of output tokens. These values are then passed to the `CalculateUsageInWh` function (Figure 5.13), as depicted in the overall flow shown in Figure 5.11.

```
public Usage PredictUsageInWh(int inputTokens, ChatModeConfiguration cfg)
{
    if (inputTokens < 0)
    {
        throw new ArgumentOutOfRangeException("Input tokens cannot be below zero");
    }
    var outputTokens = PredictOutputTokens(inputTokens, cfg);
    return CalculateUsageInWh(inputTokens, outputTokens, cfg);
}
```

**Figure 5.11:** Workflow: Predicting output tokens and applying energy calculation

```
private static int PredictOutputTokens(int inputTokens, ChatModeConfiguration cfg)
{
    // Simple heuristic: 1:1
    return Math.Min(inputTokens, cfg.MaxOutputTokens ?? int.MaxValue);
}
```

**Figure 5.12:** Prediction of output token count based on input tokens



```

public Usage CalculateUsageInWh(int inputTokens, int outputTokens, ChatModeConfiguration cfg)
{
    if (inputTokens < 0)
    {
        throw new ArgumentOutOfRangeException("Input tokens cannot be below zero");
    }

    if (outputTokens < 0)
    {
        throw new ArgumentOutOfRangeException("Output tokens cannot be below zero");
    }

    var energyWh = cfg.AlphaWhPerInputToken * inputTokens +
                  cfg.BetaWhPerOutputToken * outputTokens +
                  cfg.ZetaConstWh;

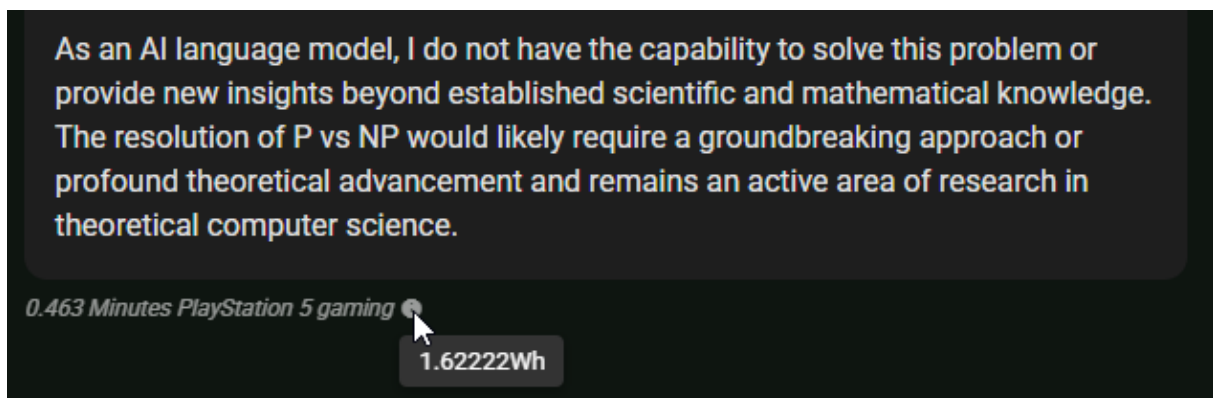
    return new Usage
    {
        NumberOfInputTokens = inputTokens,
        NumberOfOutputTokens = outputTokens,
        UsageInWh = energyWh
    };
}

```

**Figure 5.13:** Energy usage calculation based on token counts

#### 5.5.4 Energy Note

After the chatbot's response is received, the actual energy usage is calculated using the function `CalculateUsageInWh` and displayed beneath the response message. In addition to the numeric value, the UI also presents a relatable real-world analogy based on the user's selected energy unit (e.g., minutes of LED light usage), as illustrated in Figure 5.14.



**Figure 5.14:** Post-response energy note with real-world analogy

#### 5.5.5 Metrics Dashboard

The dashboard visualizes aggregated usage metrics using Chart.js. The statistics are retrieved via `GetUsageStatistics` and rendered into graphs showing energy usage, token counts, prompt counts, and usage distribution.



The frontend uses a `StatisticsService` to manage API interaction and local state. Charts include bar, line, and pie diagrams with dynamic annotations. More details on the specific metrics can be found in 4.4.2.

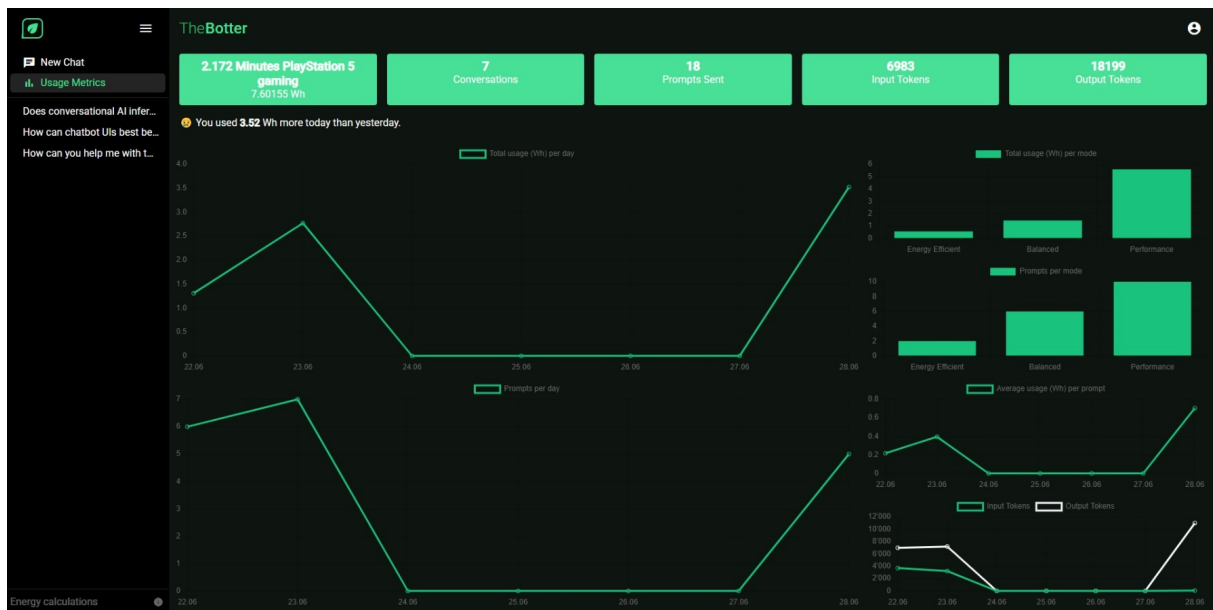


Figure 5.15: User metrics dashboard with energy, token, and mode usage breakdown

### 5.5.6 Energy Analogies and Onboarding

During account setup, users are shown a guessing game where they estimate the energy of real-world devices. Their responses are scored based on the deviation and used to recommend an energy anchor (e.g., LED lamp, iPhone charge). Screenshots shown in Figure 5.16 and 5.17.

The onboarding screen for 'TheBotter' includes a 'Welcome!' message and a guessing game to estimate energy consumption. The instructions are: 'Let's personalize your experience. Please estimate the energy consumption of the following items (in Wh):'.

1. iPhone 14 charging\*
2. Minute powering a fridge\*
3. Minute working on a laptop\*
4. LED spot approx. 500lm (1h)\*
5. Minute PlayStation 5 gaming\*

A 'Next' button is located at the bottom right of the form.

Figure 5.16: Guessing energy analogies

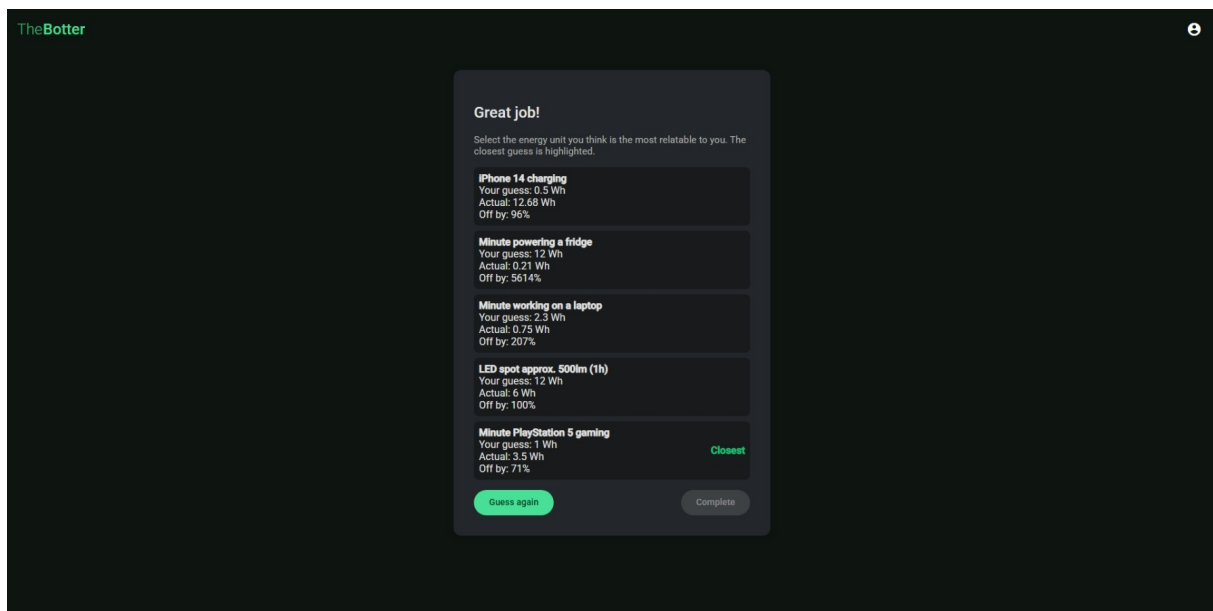


Figure 5.17: Users can choose their preferred energy analogy

This selected analogy unit is stored via `UpdateUser` and used throughout the app in predictions, notes, and the dashboard.

### 5.5.7 Settings Page

The settings page allows users to change their energy unit or disable features. Toggles trigger backend updates via `UpdateUser` and are also logged via log type `SustainabilityModeChange`.

Figure 5.18 shows the final implementation.

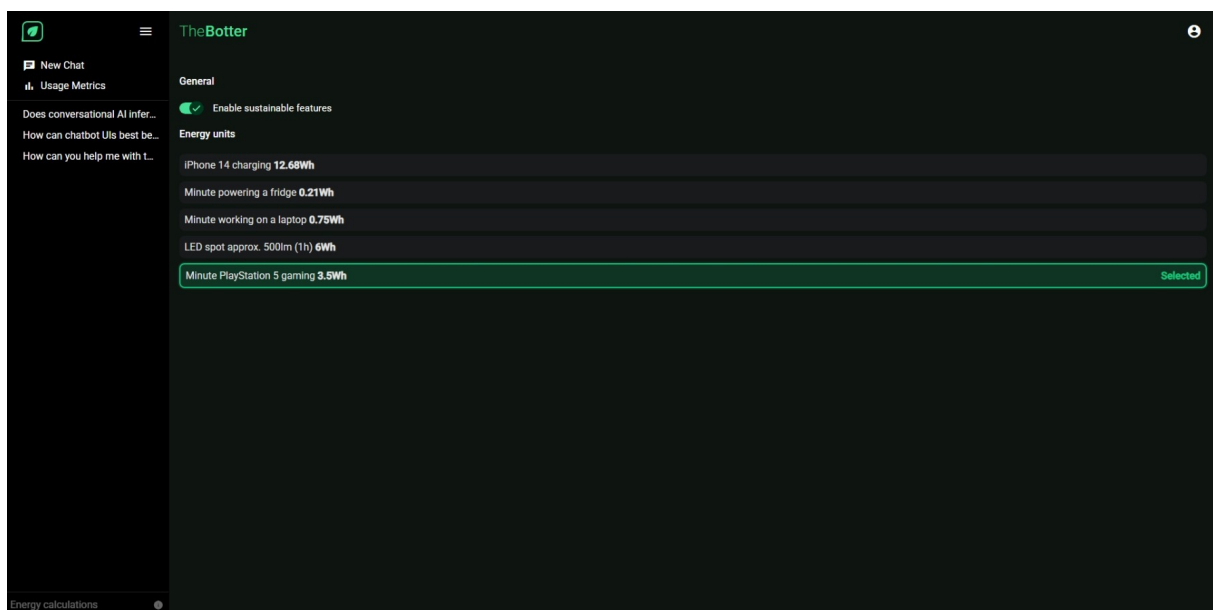


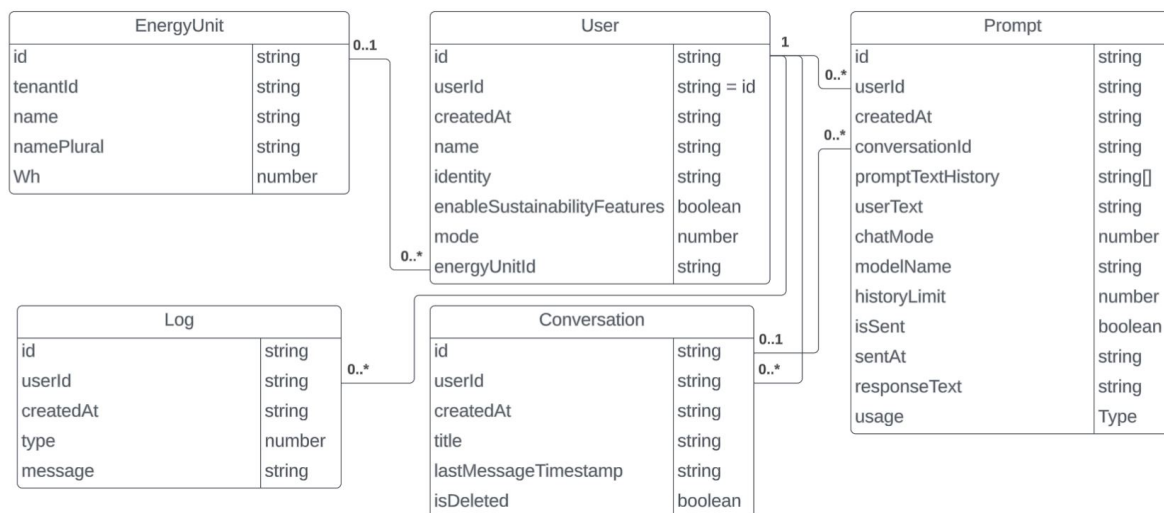
Figure 5.18: User settings

## 5.6 Data Storage: Cosmos DB

Azure Cosmos DB was chosen as the primary data storage solution due to its flexible, schema-less design, low-latency access, and scalability. This makes it particularly well-suited for storing chat-related data such as user prompts, responses, and metadata, especially in scenarios where the structure may evolve over time, such as future integration of Retrieval Augmented Generation (RAG) features.

As a NoSQL document database, Cosmos DB allows for heterogeneous document structures within the same container. This flexibility enables efficient storage of evolving prompt data, user metadata, and interaction histories. For example, arrays such as previous messages or model-specific metadata can be stored directly as `string[]` without requiring complex schema changes.

The application organizes its data into multiple containers. Each container is logically associated with an entity type, and partitioning is primarily based on the user ID to isolate user data and support horizontal scalability. The database schema is visualized in the following diagram. Note that foreign key relationships shown in the diagram are for conceptual clarity only, Cosmos DB does not support referential integrity across containers [83].



**Figure 5.19:** Logical database schema (conceptual relationships)

All containers, except the one for energy units, are partitioned by user ID. This design allows for scalable performance, simplified query filtering by user and logical data isolation. Although this does not enforce access control at the database level, it enables efficient per-user access filtering at the application level. The energy units container, by contrast, is globally shared across users and uses a constant partition key to avoid redundant definitions.

### 5.6.1 User

Each document in this container represents a user entity. The `id` field corresponds to the authenticated user ID from Microsoft login. This `id` is also used as the partition key, meaning each user has their own partition. Additionally, this user ID is referenced as the partition key in other containers.

### 5.6.2 Conversation

Represents a user conversation. Each conversation belongs to one user. Deletion is handled using a soft-delete strategy by marking records as deleted, rather than removing them from the database.

### 5.6.3 Prompt

Stores individual prompt–response pairs. A draft prompt document is created immediately when the user starts typing. If a prediction is generated, the record is persisted in the database, even if the user later cancels or changes their input. Once a prompt is sent, it is associated with a conversation, marked as "sent", and enriched with additional metadata such as the model response, estimated energy usage, chat mode, model name, and the included history size. This allows accurate reconstruction of the context for analysis or replay, even if configuration settings change over time.

### 5.6.4 Log

Used primarily for tracking user behavior during experiments. Each log entry is written manually and assigned a type. Current log types include: `Unknown`, `PageVisit`, and `SustainabilityModeChange`.

### 5.6.5 EnergyUnit

Represents an energy unit or analogy such as an "iPhone 14 charge". Since Cosmos DB requires a partition key on all containers, this container uses a fixed partition key called `tenantId` with the constant value `"0"`. This pseudo-partitioning avoids data duplication while remaining compliant with platform requirements.

## 5.7 Non functional requirements

Beyond implementing the core features of the web application, various non-functional aspects have been considered to ensure the prototype is maintainable, scalable and secure. This chapter discusses the system's architecture and operational qualities, focusing on deployment processes, authentication, performance capabilities and technical limitations.

While the application is not yet in a production environment, many architectural decisions were made with future extensibility in mind. The use of Azure-native services such as Static Web Apps, Functions, and Cosmos DB allows for scalable and secure operation with minimal configuration. Furthermore, modern development workflows such as continuous integration and automated deployment contribute to a high level of maintainability.

The following sections provide a detailed overview of the system's non-functional characteristics.

### 5.7.1 Logging and monitoring

The prototype incorporates basic logging and monitoring functionality using Azure-native tools. Azure Application Insights is connected to the Azure Functions application and provides real-time observability into backend operations. It collects telemetry such as request rates, failure rates, response times, and traces, which can be explored through dashboards, Kusto queries, or alerts.

On the database level, monitoring is available via the Azure Portal within the Cosmos DB account. Built-in tools allow developers to inspect metrics such as throughput consumption (RU/s), latency, and storage usage. In addition, logs and query performance data can be viewed to identify performance bottlenecks or data inconsistencies.

These monitoring capabilities are crucial for debugging, performance tuning, and maintaining operational health, especially in scalable, distributed environments like those offered by Azure.

### 5.7.2 Security

Security within the developed prototype is centered around authentication and access control, leveraging Azure Static Web Apps' built-in identity management. By default, Azure Static Web Apps supports authentication providers such as Microsoft, GitHub, Google, and Twitter. This project uses Microsoft Entra ID (formerly Azure Active Directory) as the sole identity provider [84].

When users access the application, they are redirected to a secure Microsoft login page. Upon successful authentication, Azure injects identity information into the request via the `x-ms-client-principal` header. This header is base64-encoded and contains user attributes such as ID, email, and assigned roles.

On the frontend, route-level protection is implemented using an Angular `AuthGuard`. For any navigation to sensitive routes such as chat, settings, metrics, or account setup, the guard first checks if the current user is authenticated. If the user is not authenticated, the guard redirects them to the login endpoint (`/.auth/login/aad`). This ensures that only authenticated users can access protected areas of the application.

In development mode, where authentication is not enforced by Azure, a proxy configuration is used to simulate authentication. The frontend serves mock identity data via the `/.auth/me` endpoint using a local file (`mock-auth.json`). This approach enables developers to test authenticated behavior locally without logging in via Microsoft Entra ID.

The backend, implemented using Azure Functions, independently verifies the presence and validity of the `x-ms-client-principal` header. A custom `AuthorizationService` extracts and deserializes this header, ensuring that only requests from authenticated users are processed. Unauthorized

requests, either due to a missing header, invalid deserialization, or lack of the “authenticated” role are rejected early with detailed error logging. This additional server-side validation ensures that even if a route were misconfigured on the frontend or Azure Static Web Apps, access to backend functions would remain protected.

Access control is further enforced declaratively in the `staticwebapp.config.json` configuration file, which restricts the `/api/*` route to authenticated users only. This creates a defense-in-depth model: Configuration-level protection via Azure, route guarding on the frontend, and validation on the backend.

Currently, the system adopts a flat authorization model: All authenticated users are granted identical access privileges. No role-based or fine-grained access control mechanisms have been implemented at this stage. If the application evolves toward a multi-role system (e.g., users vs. admins), Azure Static Web Apps supports assigning roles [85], and additional logic could be introduced on both frontend and backend to differentiate access.

Overall, the prototype achieves a strong level of security appropriate for its scope by combining platform-native authentication with layered access control and explicit validation. Improvements for future production deployment may include token expiration handling, detailed audit logging, support for additional identity providers, and fine-grained role enforcement.

### 5.7.3 Testing

Automated tests have been implemented for both the backend and frontend to ensure functional correctness and maintainability of the application over time.

**Backend:** The backend tests are located in a dedicated test project within the backend’s repository and cover both unit and integration levels. The primary focus is on the business logic of the application. Unit tests isolate individual components and validate their behavior under different conditions. Integration tests verify that multiple components such as services and controllers work together as expected, without relying on real database connections or external services. In both cases, data access and external dependencies are mocked to ensure test speed, reliability, and independence from the deployment environment [86], [87].

**Frontend:** The frontend project includes a comprehensive suite of automated tests for all components containing extensive TypeScript logic and for all the services [88]. These tests are implemented using Jasmine [89] and executed with Karma [90]. The test suite is integrated into the CI/CD pipeline, ensuring that all frontend logic is verified before deployment. This approach helps maintain high code quality as the application evolves.

This test structure allows both backend and frontend logic to be verified in isolation and supports a robust, maintainable development process. Automated tests are also integrated into the continuous integration pipeline, which ensures that code changes are verified before deployment.

### 5.7.4 Scalability and Availability

Although scalability and availability were not primary goals of this prototype, the chosen architecture inherently provides a foundation that could support these qualities in a production environment.

The frontend is hosted using Azure Static Web Apps, which is globally distributed by default. While the deployment is configured with the location *West Europe* for backend integrations, the static content itself is served via Azure’s global edge network, ensuring low-latency access from most geographic regions [91]. Currently, no advanced configuration such as custom CDN integration or regional routing has been applied. However, the platform supports enhancements such as Azure Front Door [92] and

enterprise-grade edge caching that could significantly improve performance and availability in high-traffic scenarios.

The backend logic and data storage are deployed in the Azure region *Switzerland North*. This includes the Azure Functions used to handle API requests and the Azure Cosmos DB instance managing per-user data. Azure Functions is currently running on the default consumption plan, allowing automatic scaling based on traffic demand [93]. This suits the low and intermittent workload of the prototype, although cold-start latency may become an issue in production. Premium plans are available to mitigate this by enabling pre-warmed instances and improved scaling limits [94].

Azure Cosmos DB is provisioned in single-region mode in Switzerland North. While this simplifies the setup, it also introduces a regional single point of failure. Cosmos DB supports multi-region replication with automatic failover and global distribution, which could be enabled to increase resilience and reduce read latency across geographies [95]. Additionally, Cosmos DB's autoscale and partitioning features remain unused but offer a clear path to horizontal scalability as the dataset or user base grows [96].

The application relies on the OpenAI Completions API [97], which is hosted outside of Azure and geographically decoupled from the application's infrastructure. This introduces some dependency on public internet connectivity and may be affected by regional latency or rate limits. The prototype operates well within the default OpenAI API quotas. However, for production scenarios, implementing retry logic, monitoring API usage, and applying for higher quotas would be essential [98].

In summary, while the current deployment operates under default configurations suitable for prototyping, the system is cloud-native and inherently scalable. All critical services such as Azure Functions, Cosmos DB and Static Web Apps support scaling mechanisms, multi-region deployment, and high-availability configurations that could be activated with minimal architectural changes. These capabilities would become crucial if the prototype evolves into a production-grade system [99]. Enhancing scalability and availability would primarily involve enabling geo-replication, autoscaling configurations, and CDN integration, all of which are natively supported within the Azure ecosystem [99].

### 5.7.5 Continuous Integration and Deployment (CI/CD)

The project incorporates two separate CI/CD pipelines. One for the frontend and one for the backend using GitHub Actions in combination with Azure services. These pipelines enable continuous integration and deployment with minimal manual intervention, ensuring that the latest versions of both the user interface and backend logic are reliably deployed to Azure upon each commit to the main branch.

The **frontend pipeline** is defined in the file `azure-static-web-apps-green-mud-04afae203.yml`. It is automatically triggered when changes are pushed to the frontend directory. The pipeline installs dependencies, builds the Angular application, runs the tests and deploys it to Azure Static Web Apps. This deployment process is fully integrated into GitHub via the Azure Static Web Apps GitHub Action [100], which abstracts away much of the manual configuration typically required.

The **backend pipeline**, defined in `main_sustainable-ai-api.yml`, is responsible for building the .NET project, running all automated tests, and deploying the application to Azure Functions. This pipeline ensures that only tested code reaches the production environment. It uses standard GitHub Actions for .NET development, such as building, testing, and publishing the application, along with the Azure Functions GitHub Action for deployment [101].

While the CI/CD setup ensures automated deployments, the project currently does not include any form of **Infrastructure as Code (IaC)**. All Azure resources have been provisioned manually via the Azure Portal.

This decision reflects the prototypical nature of the project. Since the goal was rapid iteration and experimentation, setting up full IaC automation would have added complexity and overhead without immedi-

ate benefit. Manual provisioning allowed for flexible configuration and faster development during early stages which was beneficial in a rapidly evolving prototype setting. However, for production scenarios, IaC is considered a best practice to enable reproducible, version-controlled, and automated infrastructure management [102], [103].

Overall, the CI/CD approach implemented in this project provides a pragmatic balance between automation and flexibility. It ensures reliable deployments while keeping the setup lightweight and maintainable for a prototype.

## 5.8 Limitations and constraints

The current prototype of the web application is subject to several limitations and constraints, which result both from deliberate design decisions and from the technical scope of the implementation.

At this stage, the web application is optimized for desktop use and has not been adapted for mobile or tablet devices. Furthermore, the user interface and chatbot interactions are available exclusively in English, which may limit accessibility for non-English-speaking users.

Interaction with the chatbot is currently limited to text input and output. Features such as voice input, file upload, or multimodal interfaces are not supported. This constraint aligns with the prototype's goal of keeping the interaction model simple while focusing on core functionality.

User authentication is restricted to login via Entry ID, with no support for alternative identity providers or anonymous access. While sufficient for testing and evaluation, this would need to be extended in a production environment to support broader identity management options.

The web application does not implement pagination or lazy loading for chat history. As a result, only the most recent 100 conversations and the last 100 messages (prompts and responses) are displayed to the user. This simplification was chosen to reduce implementation complexity but may limit usability for users with extensive chat histories.

The chatbot's responses are streamed in real-time to the frontend using server-sent events. This approach improves user experience by reducing perceived latency but limits compatibility to a subset of models that support streaming. While it would be possible to support non-streaming models by implementing a custom version of the `IChatBotService` interface that yields responses as an `IAsyncEnumerable`, such an adaptation is not yet implemented.

On the data persistence layer, the use of Azure Cosmos DB introduces another important constraint: As a NoSQL database with a flexible schema, Cosmos DB does not enforce structural consistency across documents. To mitigate this, the application ensures schema validity at the serialization level by marking key properties as required. This ensures that documents with missing or malformed data trigger exceptions during deserialization, preventing faulty runtime behavior. However, this approach only identifies the problems but it does not resolve them. A more robust solution would be to implement lazy migrations adjusting documents on-the-fly as they are loaded to conform to the latest schema version. Such an approach would improve resilience and forward compatibility and could be considered in the future if the prototype evolves into a production-grade application.



## 6 Validation and Results

This chapter evaluates to what extent the UI-only interventions proposed in Section 4.3 achieve their intended goals of (1) raising user awareness regarding the energy footprint of LLM queries and (2) nudging users towards lower-impact behavior. We triangulate three independent data sources:

- (a) Five *daily awareness check-ins* ( $n = 11$ , 52 responses, 95 % completion),
- (b) A *post-study questionnaire* ( $n = 11$ , 100 % completion), and
- (c) Detailed *user data and behavioral logs* captured by the prototype backend.

### 6.1 Results from Daily Check-in and Final Questionnaire (a & b)

In this section we report the survey-based outcomes from the five-day experiment.<sup>3</sup> We focus on two self-report surveys:

- The **daily awareness check-in** (five small surveys,  $n = 11$ ; 52 responses; 95 % completion rate)
- The **final comprehensive questionnaire** ( $n = 11$ ; 100 % completion rate).

All questions used a 5-point Likert scale [15] with 1 meaning *strongly disagree* and 5 *strongly agree*.

#### 6.1.1 Daily Awareness Trajectory

Table 6.1 shows the mean self-reported awareness<sup>4</sup> for each study day. The small mid-week dip (Days 3–4) matches qualitative feedback that participants “stopped paying attention once the novelty wore off,” suggesting repetition alone is insufficient, more salient reminders were needed until the dashboard visit on Day 5 (see below). This upward trend from Day 1 ( $M = 3.27$ ) to Day 5 ( $M = 4.44$ ) supports **Hypothesis H1**, which states that visibility of energy information increases user awareness over time.

Day	Mean	Std. Dev.	$\Delta$ vs. Day 1
1	3.27	0.90	—
2	3.64	0.81	+0.37
3	3.09	1.14	−0.18
4	3.10	0.99	−0.17
5	<b>4.44</b>	0.73	<b>+1.17</b>

**Table 6.1:** Mean awareness per day “At this moment I’m aware of the energy cost of the prompts I sent today”, (1 = Strongly Disagree, 5 = Strongly Agree)

#### 6.1.2 Feature Salience

Across all 52 check-ins, participants rated the *Energy-Note* and the *Mode-Toggle* as the most visible interventions ( $M = 3.19$  for both), whereas active dashboard use occurred on 59.6 % of participant-days. The spike in awareness on Day 5 coincides with the highest dashboard traffic (Table 6.5), indicating that summary feedback still plays a crucial consolidation role even after inline cues are present. This further corroborates **Hypothesis H1**, since increased visibility of features correlates with the highest awareness score.

<sup>3</sup>Behavioral log data are analyzed separately in Section 6.3.

<sup>4</sup>“At this moment I’m aware of the energy cost of the prompts I sent today”; 5-point Likert.

## 6.2 Final Questionnaire

### 6.2.1 Self-Reported Behavioral Change

Table 6.2 summarizes the behavioral section. Nine of eleven participants (81.8%) reported *regularly* switching to Energy-Efficient-Mode when answer accuracy was “not mission-critical” ( $M = 4.18$ ). Prompt shortening received a neutral rating ( $M = 3.00$ ), and the detailed examination of the log data in Section 6.3 confirms there was no substantial reduction in input length on average (see Figure 6.2). Distraction by sustainability cues was low ( $M = 1.45$ ), suggesting the interventions did not impair chat-bot usability.

Item	Mean	Std. Dev.
Energy-Efficient-mode when accuracy uncritical	<b>4.18</b>	0.75
Shortened / reduced prompts	3.00	1.34
Dashboard visited $\geq 2\times$	0.77 (Yes)	n.a.
Energy info distracted me	<b>1.45</b>	0.82
Want similar features elsewhere	<b>3.91</b>	1.14

**Table 6.2:** Behavioral items, (1 = Strongly Disagree, 5 = Strongly Agree)

Preliminary day-level correlations show that higher awareness is associated with a greater Energy-Efficient-mode share and shorter prompts, which supports **Hypothesis H2**, stating that awareness is positively correlated with pro-sustainability behavior.

### 6.2.2 Per-Feature Usability

All four implemented features scored above 4 (“easy to understand and effect was clear”), with the *inline three-mode toggle* leading ( $M = 4.64$ ). Qualitative answers (Section 6.2.3) emphasized participants’ desire for more ephemeral reminders and aggregate figures in familiar units.

Feature	Mean
3-mode toggle	<b>4.46</b>
Placement of toggle	<b>4.64</b>
Energy note (per response)	4.09
Metrics dashboard	4.00

**Table 6.3:** Usability ratings: “Easy to understand and effect clear”, (1 = Strongly Disagree, 5 = Strongly Agree)

### 6.2.3 Qualitative Feedback – Text Answers

Four of the eleven participants wrote substantive comments, each referring to different ideas or improvements. Three themes emerged:

- T1 Proactive reminders:** Two participants requested “occasional push notifications” to surface energy consumption without requiring a dashboard visit.
- T2 Real-time aggregate figures in familiar units:** Three comments called for a constantly visible “current usage” badge expressed in the chosen energy unit (e.g., *Laptop-minutes*) instead of raw Wh values, which is currently only visible on the dashboard.

**T3 Contextual framing of impact:** Respondents stated that numbers alone “don’t really help me that much” and proposed a concise “wake-up call” every 100 prompts, optionally linking to a deeper explanation or video.

One participant (original German comment) also suggested: “*Bei sehr kurzen Prompts das Verwenden einer klassischen Google Suche anstelle von LLM vorschlagen*” (“For very short prompts, recommend a classic Google search rather than an LLM”), which was also a potential feature considered but not implemented.

#### 6.2.4 Summary of Awareness Trajectory

- Awareness increased over the week, reaching  $M = 4.44$  on Day 5
- Self-reported Energy-Efficient-mode adoption was high and distraction minimal.
- All usability means exceeded 4/5, with the mode toggle scoring highest.
- Preliminary day-level correlations show that higher awareness is associated with a greater Energy-Efficient-mode share and shorter prompts, supporting **Hypothesis H2**.

These findings motivate a deeper behavioral analysis based on the collected server log data presented next.

### 6.3 Behavioral Results from Application Data

#### 6.3.1 Data Set Results Total

A total of 278 prompt–response pairs were recorded over the study’s duration, totaling 683,580 processed tokens (231,779 input, 451,801 output) and 168.0 Wh estimated energy. In the established energy analogies this would correspond to:

- 13.2 full charges of an iPhone 14
- 13.3 hours of refrigerator usage
- 3.7 hours of laptop use
- 28 hours of a 500-lumen LED lamp
- 48 minutes of PlayStation 5 gaming

#### 6.3.2 Energy-Unit Selection

Eight of eleven participants selected the *iPhone 14 charge* analogy, two preferred *minute powering a fridge*, and only one chose *minute PlayStation 5 gaming*, no one selected the laptop or LED-lamp options (Table 6.4). Participants thus gravitate toward anchors that (i) represent a *complete, familiar action* (charging a phone) and (ii) occur in their everyday routine. Time-based slices of low-power devices were perceived as too abstract. Future designs should therefore favor whole-device cycles e.g. “electric-toothbrush charge” over “x-minutes of use” to maximize intuitiveness.

Mode	Used by
iPhone 14 charging	8 Users
Minute working on a laptop	0 Users
Minute powering a fridge	2 Users
LED spot approx. 500lm (1h)	0 Users
Minute PlayStation 5 gaming	1 User

**Table 6.4:** Energy unit preferences

### 6.3.3 Dashboard Engagement

Across the five-day study, the dashboard was opened **41** times ( $\bar{x} = 3.7$  per participant, median = 3), see Table 6.5 for the overview. Two participants never visited it, while the top user (ID 8) alone generated **29%** of all visits; the three most active users (IDs 8, 2, 10) accounted for **59%** (24 / 41). Day-level activity peaked on Day 1 (14 visits, 34 %), dropped by half on Day 2, briefly rebounded on Day 3, and tapered off to four visits on Day 5:

$$\text{Day totals} = \langle 14, 7, 10, 6, 4 \rangle$$

This pattern suggests initial curiosity followed by habituation; sustained reflection therefore requires additional prompts (e.g., weekly summaries or push reminders) rather than a purely on-demand dashboard. The usage pattern supports **H1** by illustrating that dashboard-based summary feedback temporarily raised awareness levels.

	Day 1	Day 2	Day 3	Day 4	Day 5	Total
User 1	1	0	2	0	1	4
User 2	1	2	0	3	0	6
User 3	1	0	1	0	0	2
User 4	0	1	2	0	0	3
User 5	0	1	2	0	0	3
User 6	1	0	1	0	2	4
User 7	0	0	0	0	0	0
User 8	7	2	1	2	0	12
User 9	0	0	0	1	0	1
User 10	3	1	1	0	1	6
User 11	0	0	0	0	0	0

**Table 6.5:** Metrics visits per user and day

### 6.3.4 Mode Comparison

The three-mode toggle distributed user traffic across models with varying energy demands. Table 6.6 compares the share of prompts, token counts, and energy consumption across modes.

Mode	Prompts	Tokens in	Tokens out	Energy
Energy-Efficient	155 (55.8%)	60473 (26.1%)	213350 (47.2%)	6.98Wh (4.1%)
Balanced	54 (19.4%)	83213 (35.9%)	82272 (18.2%)	11.51Wh (6.8%)
Performance	69 (24.8%)	88093 (38.0%)	156179 (34.6%)	149.55Wh (89.0%)
<b>Total</b>	278 (100%)	231779 (100%)	451801 (100%)	168.04Wh (100%)

**Table 6.6:** Share of prompts, tokens, and energy consumption per mode

Performance mode accounted for 24.8% of prompts but 89.0% of total energy usage. Energy-Efficient mode processed the majority of prompts (55.8%) but consumed only 4.1% of the total.

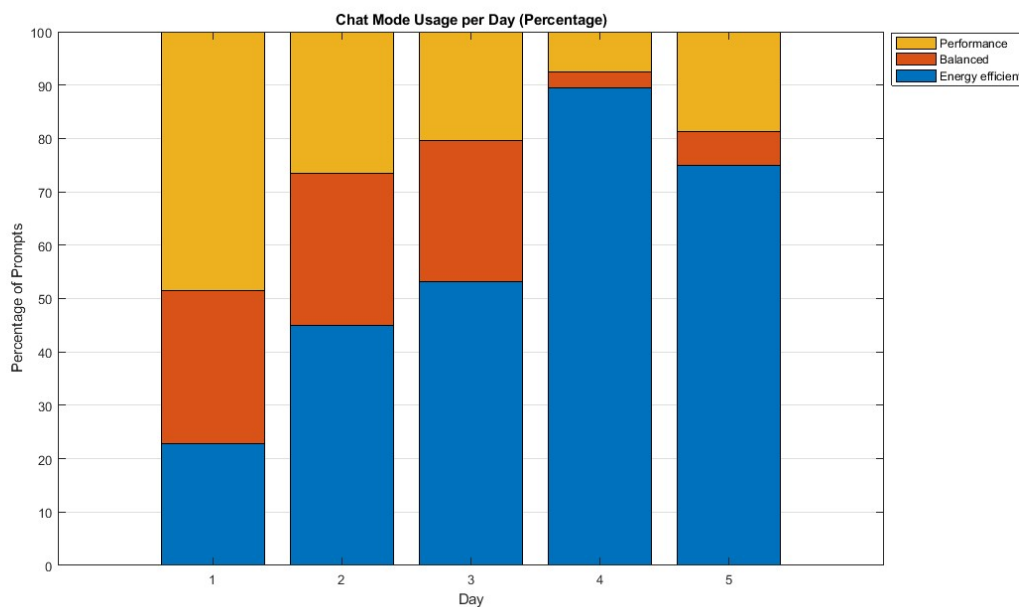
To control for prompt length, Table 6.7 compares each mode’s energy consumption normalized by input token count only. This isolates energy per 1,000 input tokens and excludes any weighting of output. Additionally, the constant overhead of (0.020 Wh) has been excluded from the usage.

Performance mode consumed approximately 13.4 times more energy per input token than balanced mode. Balanced mode consumed about 1.95 times more than Energy-Efficient mode.

Mode	Tokens in	Total	Overhead	Token based	per 1k Tokens in
Energy-Efficient	60,473	6.98Wh	3.1Wh	3.88Wh	0.064Wh
Balanced	83,213	11.51Wh	1.08Wh	10.43Wh	0.125Wh
Performance	88,093	149.55Wh	1.38Wh	148.17Wh	1.682Wh

**Table 6.7:** Energy consumption per 1,000 input tokens by mode

Figure 6.1 illustrates the daily distribution of prompts per mode. Performance mode dominated on Day 1, then decreased during the week with a small spike on Day 5, while Energy-Efficient mode increased steadily, peaking on Day 4, with a small decrease on Day 5. This indicates that users adapted their behavior over the week. Balanced mode remained relatively stable until Day 3, then dropped to a lower level on Days 4 and 5.



**Figure 6.1:** Total prompts sent per mode per day

### 6.3.5 User-Level Patterns

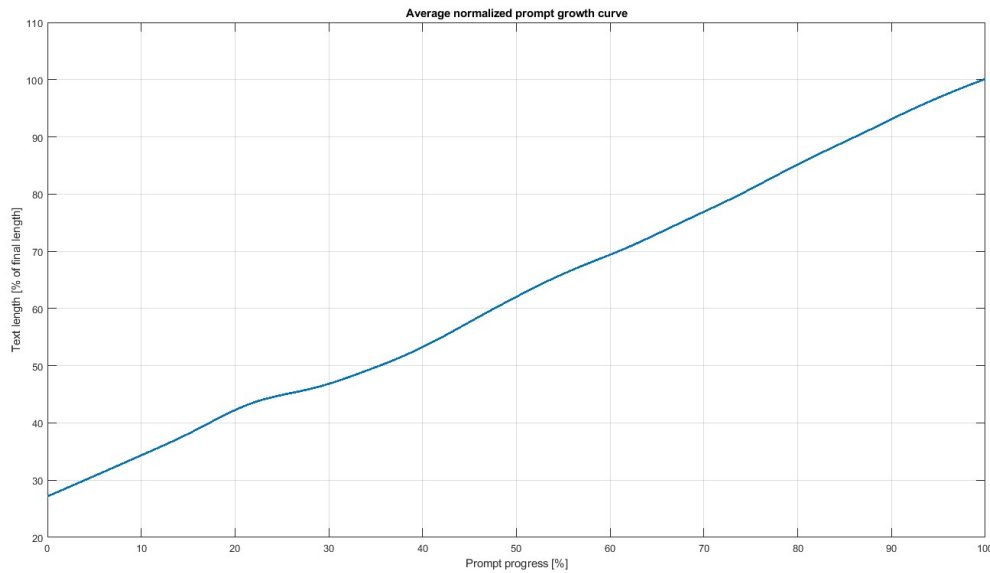
Three heavy users (IDs 8, 5, 3) generated **77%** of the study's total energy footprint as shown in Table 6.8; ID 8 alone accounted for nearly half of it, owing to a 64% reliance on Performance mode. Conversely, eight participants sent a majority of prompts in Energy-Efficient-mode, thus reaching a much lower total energy consumption. This supports **H2**, indicating that higher individual awareness is associated with energy-saving behavior, as seen in both mode choices and per-user footprints.

User	Mode	Prompts	%	In	Out	Usage (Wh)
1	Energy efficient	9	47.368	1513	7588	0.31383
	Balanced	8	42.105	1313	3694	0.56544
	Performance	2	10.526	272	2748	2.378
	Total	19	100	3098	14030	3.2572
2	Energy efficient	26	74.286	7297	37641	1.183
	Balanced	7	20	4407	8760	1.1341
	Performance	2	5.7143	2239	4150	3.9547
	Total	35	100	13943	50551	6.2718
3	Energy efficient	15	53.571	5298	27434	0.78314
	Balanced	0	0	0	0	0
	Performance	13	46.429	27227	27878	29.116
	Total	28	100	32525	55312	29.9
4	Energy efficient	4	57.143	166	5905	0.1799
	Balanced	3	42.857	264	1299	0.19759
	Performance	0	0	0	0	0
	Total	7	100	430	7204	0.37749
5	Energy efficient	9	31.034	12321	10920	0.4152
	Balanced	13	44.828	68652	36767	5.6961
	Performance	7	24.138	4850	15288	13.848
	Total	29	100	85823	62975	19.959
6	Energy efficient	17	60.714	3137	23811	0.7532
	Balanced	6	21.429	3941	6732	0.8979
	Performance	5	17.857	3800	11202	10.196
	Total	28	100	10878	41745	11.847
7	Energy efficient	0	0	0	0	0
	Balanced	4	100	671	4900	0.59083
	Performance	0	0	0	0	0
	Total	4	100	671	4900	0.59083
8	Energy efficient	10	18.868	7853	15612	0.49526
	Balanced	9	16.981	2753	14733	1.7345
	Performance	34	64.151	48302	79649	76.932
	Total	53	100	58908	109990	79.162
9	Energy efficient	65	98.485	22888	84439	2.8147
	Balanced	1	1.5152	5	48	0.024964
	Performance	0	0	0	0	0
	Total	66	100	22893	84487	2.8397
10	Energy efficient	0	0	0	0	0
	Balanced	0	0	0	0	0
	Performance	6	100	1403	15264	13.084
	Total	6	100	1403	15264	13.084
11	Energy efficient	0	0	0	0	0
	Balanced	3	100	1207	5339	0.62859
	Performance	0	0	0	0	0
	Total	3	100	1207	5339	0.62859

**Table 6.8:** Aggregated prompts per user and per mode

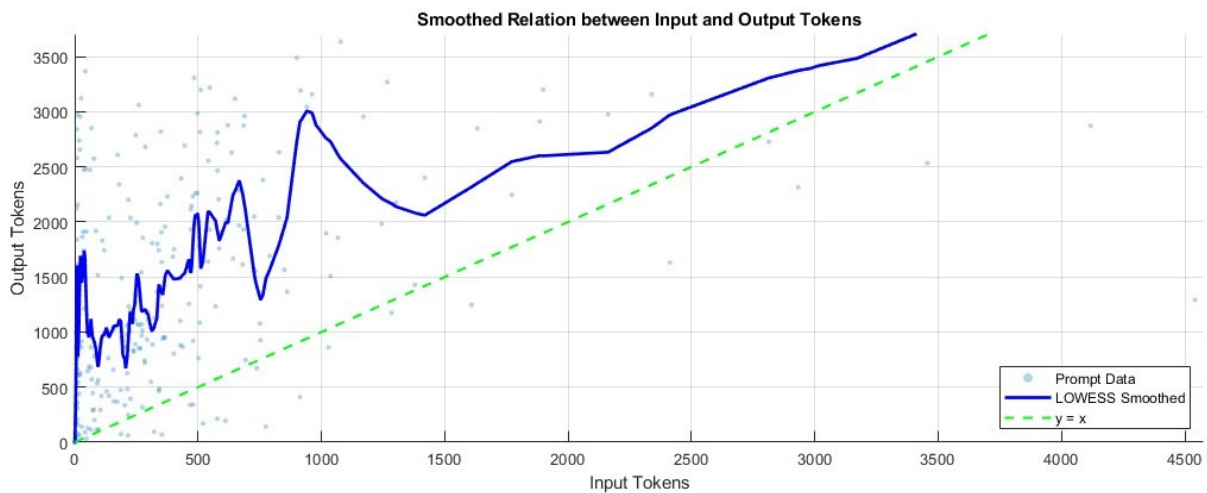
### 6.3.6 Token Length Observations

The average input length remained stable throughout the week, suggesting that participants did not adjust their prompts to reduce length over time. Participant comments indicated that “rewriting takes too much time,” implying that selecting a different model was perceived as a more efficient strategy. Figure 6.2 shows that users rarely edited their input during typing, particularly with regard to shortening prompts in response to model predictions.



**Figure 6.2:** Average normalized prompt growth curve

Furthermore, Figure 6.3 offers an additional perspective by presenting a scatter plot of input versus output tokens (95th percentile), including both the output prediction function  $f(x) = x$  defined in section 4.4.3 and a LOWESS [104] trend of the actual prediction.



**Figure 6.3:** Scatter plot of input vs. output tokens with LOWESS trend and prediction function

The plot shows that the prediction function did not accurately estimate output tokens in most cases. On average, the number of output tokens was higher than the number of input tokens. The distribution of

prompts per conversation shows a mean of 4.03, a median of 3, and a mode of 1, indicating that while some users sent multiple prompts, many conversations consisted of only one.

## 6.4 Summary

### 6.4.1 Self-Reported Awareness and Behavior

Self-report ratings reveal the following patterns in participants awareness and energy-saving behavior:

- Awareness increased over the five days, peaking on Day 5 ( $M = 4.44$ ), aligned with highest dashboard activity. [Supports H1]
- 81.8% of participants regularly used Energy-Efficient-mode when accuracy was not critical ( $M = 4.18$ ). [Supports H2]
- Prompt shortening was rated neutrally ( $M = 3.00$ ); log data confirmed no substantial input length reduction.
- All features scored above 4/5 in usability, with the inline *mode toggle* rated highest ( $M = 4.64$ ).
- Qualitative feedback highlighted needs for:
  - **T1:** Proactive push reminders,
  - **T2:** Real-time usage feedback in relatable units (e.g., *Laptop-minutes*),
  - **T3:** Periodic contextual framing or “wake-up calls”.

### 6.4.2 Behavioral Log Findings

System logs show how each mode’s use affected energy consumption and how participants choices shifted over time:

- Despite handling only 24.8% of prompts, **Performance mode** caused 89% of energy usage.
- **Energy-Efficient-mode** handled 55.8% of prompts with just 4.1% of the energy footprint.
- Energy consumption per 1,000 input tokens:
  - Energy-Efficient-mode: **0.064 Wh**
  - Balanced: **0.125 Wh** ( $1.95\times$  more than Energy-Efficient)
  - Performance: **1.682 Wh** ( $13\times$  more than Balanced)
- Energy-Efficient-mode usage increased over the week, indicating behavioral adaptation. [Supports H2]
- Dashboard usage was front-loaded (Day 1: 14 visits), then tapered off, suggesting that sustained reflection requires push-based reminders. [Supports H1]
- Three heavy users (IDs 8, 5, 3) caused 77% of total energy usage, with ID 8 alone responsible for 47%.
- Input lengths remained stable; participants favored mode selection over prompt rewriting due to time efficiency.

## 6.5 Hypothesis Validation Summary

**H1. Showing per-prompt consumption and predicted consumption increases awareness.**

**Supported.**

- Mean awareness increased from 3.27 (Day 1) to 4.44 (Day 5).
- Awareness peaks coincided with dashboard usage and high feature visibility.
- Inline and summary cues (Energy-Note, Mode-Toggle, Dashboard) visibly contributed to raising awareness.



**H2. Awareness scores are positively correlated with pro-sustainability behavior.**

**Supported.**

- Higher awareness levels (Days 2 and 5) aligned with increased Energy-Efficient-mode usage.
- Participants with higher self-reported awareness also demonstrated lower per-user energy footprints.
- Prompt length remained stable; mode-switching was the primary behavioral adjustment.

## 7 Discussion

This chapter reflects on the implications of the findings, evaluates the hypotheses in light of user behavior, addresses limitations of the study, and outlines directions for future work.

### 7.1 Interpretation and Evaluation

The experiment showed that even lightweight, non-intrusive UI changes can nudge users toward more energy-conscious decisions in LLM-based interactions. Awareness increased measurably, and most participants adopted energy-saving behaviors without being forced to. However, the extent of behavioral change varied widely between users. The results suggest that the success of such interventions depends not only on their design but also on user intent, attention, and usage style.

While mode switching was a low-effort action with high uptake, reducing prompt length appeared cognitively costly and was largely avoided. This asymmetry reveals an important insight: Users are more likely to adopt energy-efficient behavior when the intervention aligns with their workflow and requires minimal effort. In contrast, behaviors that demand extra cognitive investment like rewriting queries which may need stronger incentives or automation.

Moreover, usage concentration among a few participants highlights the “heavy-user bias”: A small group can dominate total energy use regardless of system nudges. For these users, awareness alone may not be enough, additional strategies such as adaptive defaults or tailored recommendations may be needed.

### 7.2 Evaluation of Energy consumption

The total number of prompts and the overall energy consumption observed during the experiment were lower than initially expected. Although the underlying formula used for estimating consumption is largely based on empirical observations rather than disclosed usage data, it raises the valid question of whether the absolute energy consumption per user is truly significant.

When compared to familiar energy analogies such as gaming on a PlayStation 5 (168 Wh  $\approx$  48 minutes of playtime) or using a laptop (3.7 hours) the consumption appears relatively low and arguably negligible. However, this perspective shifts when considering the energy saved through users actively selecting the energy-efficient mode. This behavioral change highlights the potential of UI-based interventions to reduce energy usage.

If all interactions had occurred exclusively in Performance Mode, total energy consumption would have increased from approximately 168 Wh to an estimated 429 Wh, calculated as follows:

$$231,779 \cdot 0.00021 + 451,801 \cdot 0.00083 + 278 \cdot 0.020 \approx 429 \text{ Wh}$$

This reinforces the relevance of optimizing user interfaces for energy awareness. Moreover, as this experiment was limited to a single prototype, the actual energy footprint across other conversational AI tools exceeds the measured values, further underscoring the importance of sustainable design in AI applications.

### 7.3 Answering the Research Questions

RQ1: To what extent are users currently aware of the energy implications associated with their chatbot interactions?

Initial awareness was moderate and grew over time, confirming that awareness is improvable via simple interventions. Participants responded well to inline cues and dashboard summaries. Still, awareness fluctuated depending on salience and novelty, suggesting that reminders need to be periodically refreshed.

RQ2: How can UI-based features most effectively increase user awareness regarding the energy consumption of conversational AI?

Features were most effective when integrated seamlessly into the workflow. The inline toggle and energy note stood out due to constant visibility and clarity. In contrast, less frequently accessed features like the dashboard had high value but required user initiative. Real-time, ambient cues appear more sustainable than static overviews.

RQ3: How strongly does increased user awareness correlate with reductions in conversational AI energy consumption?

Awareness was linked to energy-saving behavior, particularly through mode selection. However, this correlation was dampened by individual differences: Some users consistently optimized their usage, while others reverted to high-consumption modes for complex tasks. Prompt length remained largely unaffected. The relationship is therefore present but not uniform.

## 7.4 Limitations

Although the research questions could be answered and the hypothesis were at least partially supported this results should be seen in light of the constraints and limitations of this project:

- **Sample size and bias:** With only 11 participants, mostly with technical backgrounds, the sample is too narrow for generalization.
- **Energy approximation:** Due to lack of public data from OpenAI, energy usage was estimated via heuristic models. This introduces uncertainty in all consumption metrics.
- **Output prediction:** The output token length estimation heuristic proved to be an imprecise predictor, especially for longer prompts. This affects the accuracy of energy feedback and weakens user trust in numerical feedback. More sophisticated models like regression or transformer-specific predictors are needed for future iterations. This might also be a factor why participants did not change their prompt length.
- **Self-selection bias:** Participants were self-motivated and likely more aware of sustainability issues, which may not reflect the general population.
- **Limited feature set:** The prototype focused on a few UI changes; other potential interventions (e.g., features discussed in the conceptual solution in the decision matrix Table 4.4) were not tested.
- **Short duration:** Five days are sufficient for observing short-term change but not long-term habits or retention.
- **Limited granularity:** Some behaviors (e.g., mental effort, editing hesitations) are not captured in the logs, leaving gaps in interpretation.

## 7.5 Future Work and Development

### 7.5.1 Research and Validation

This thesis has demonstrated both the potential and the necessity for more extensive research in this emerging field. Several aspects remain insufficiently explored, and notable blind spots persist. To address these gaps, future research should:

- Incorporate real energy data from API providers to validate energy estimations.
- Expand the participant pool with non-technical and more diverse user groups.
- Run longitudinal studies to investigate lasting behavioral change and intervention fatigue.
- Introduce richer behavior tracking to detect prompt refinement, hesitation, or adaptation patterns.

### 7.5.2 Prototype Improvements

While the current prototype successfully demonstrated the feasibility of UI-based energy awareness features, several enhancements could further increase its effectiveness, usability, and adaptability. Future iterations should aim to:

- Add user-configurable modes or AI-guided suggestions based on usage patterns.
- Improve energy and output predictions with data-driven models.
- Implement features proposed in Section 4, such as real-time “current usage” displays and periodic “wake-up calls.”
- Add detection of trivial prompts and recommend simpler alternatives like a Translation Service or, calculator etc.

### 7.5.3 User Guidance and Framing

The experiment revealed that participants responded positively to clear and relatable feedback, especially when energy consumption was contextualized through comparisons (e.g., charging a smartphone). Building on this insight, future designs could enhance user engagement and awareness through:

- Emotional and contextual framing (“this equals powering a fridge for X hours”).
- Push reminders with usage summaries or tailored advice.
- More interactive dashboards showing impact over time.

## 7.6 Concluding Remarks

UI-only strategies show real potential for raising awareness and enabling more sustainable behavior in conversational AI interfaces. However, their impact depends on user context, feature salience and the perceived effort–benefit ratio. Future work should aim to make these interventions adaptive, personalized and grounded in real data.

## 8 Conclusion

This thesis explored the potential of UI-only strategies to increase awareness and reduce energy consumption in conversational AI usage. Motivated by the growing inference-related energy demand of large language models, we investigated whether lightweight frontend interventions could measurably influence user behavior without compromising usability.

The findings of our five-day field experiment with frequent LLM users indicate that lightweight frontend interventions can have a measurable impact. Participants showed a significant increase in awareness ( $M = 4.44/5$ ), and those with higher awareness scores consistently opted for more energy-efficient settings. More than half of all prompts were routed through the energy-saving mode and the overall usage resulted in an estimated 35% reduction in energy consumption compared to a performance-mode baseline. In particular, these reductions were achieved without negatively affecting user experience or perceived output quality, suggesting that sustainability and usability can coexist in the same interface.

Beyond the empirical results, this work contributes to a growing discourse within Sustainable Human-Computer Interaction. While prior research has highlighted the perceptual, technical, and design gaps that hinder energy-conscious digital behavior, our work demonstrates a concrete, scalable way to address them through intentional interface design. The study also reinforces the idea that users are not indifferent to energy concerns; instead, they often lack accessible, interpretable feedback to guide their decisions. By bridging this awareness-action gap, UI interventions can act as effective nudges in everyday systems.

However, several limitations must be acknowledged. Our study involved a relatively small technically literate user group, and energy estimates were derived using an approximation model based on token counts and pricing-based coefficients. While this approach is sufficient for behavioral nudging, it may not reflect actual backend power consumption with high accuracy. Furthermore, long-term behavioral retention and broader user demographics remain unexplored. These factors limit generalisability and should be addressed in future work.

Looking ahead, this thesis opens the door for further research into sustainability-aware design strategies for AI-powered applications. Opportunities include integrating live energy telemetry from commercial providers, exploring adaptive interfaces that respond dynamically to carbon intensity, and embedding these ideas in educational or workplace contexts.

Ultimately, as AI systems become embedded in the daily workflows of millions, energy use will no longer be just a backend concern, it will be a user-facing issue. By giving users the tools and information to make responsible choices, interface designers have a critical role to play in shaping a more sustainable digital future. This thesis has shown that even modest UI changes can lead to meaningful environmental benefits, setting the stage for a new kind of interaction design: One that is not only human-centered, but energy-aware.

## References

- [1] W. Vanderbauwhede, *Estimating the increase in emissions caused by ai-augmented search*, 2025. arXiv: 2407.16894 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2407.16894>.
- [2] A. de Vries, „The growing energy footprint of artificial intelligence“, *Joule*, vol. 7, no. 10, pp. 2191–2194, Oct. 18, 2023, Publisher: Elsevier, ISSN: 2542-4785. DOI: 10.1016/j.joule.2023.09.004. [Online]. Available: <https://doi.org/10.1016/j.joule.2023.09.004> (visited on Jul. 14, 2025).
- [3] P. Jiang, C. Sonne, W. Li, F. You, and S. You, „Preventing the immense increase in the life-cycle energy and carbon footprints of llm-powered intelligent chatbots“, *Engineering*, vol. 40, pp. 202–210, Sep. 2024. DOI: 10.1016/j.eng.2024.04.002. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2095809924002315>.
- [4] E. Strubell, A. Ganesh, and A. McCallum, „Energy and policy considerations for deep learning in nlp“, in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3645–3650. [Online]. Available: <https://arxiv.org/abs/1906.02243>.
- [5] P. Henderson, J. Hu, J. Romoff, E. Brunskill, D. Jurafsky, and J. Pineau, *Towards the systematic reporting of the energy and carbon footprints of machine learning*, 2022. arXiv: 2002.05651 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2002.05651>.
- [6] X. Wang, C. Na, E. Strubell, S. Friedler, and S. Luccioni, „Energy and carbon considerations of fine-tuning bert“, in *Findings of the Association for Computational Linguistics: EMNLP 2023*, Association for Computational Linguistics, 2023, pp. 9058–9069. DOI: 10.18653/v1/2023.findings-emnlp.607. [Online]. Available: <http://dx.doi.org/10.18653/v1/2023.findings-emnlp.607>.
- [7] B. Capital, „Trends in artificial intelligence“, Bond Capital, Tech. Rep., 2025, ChatGPT reached 800 million weekly active users in 17months. [Online]. Available: [https://www.bondcap.com/report/pdf/Trends\\_Artificial\\_Intelligence.pdf](https://www.bondcap.com/report/pdf/Trends_Artificial_Intelligence.pdf).
- [8] C. Preist, D. Schien, and E. Blevis, „Understanding and mitigating the effects of device and cloud service design decisions on the environmental footprint of digital infrastructure“, in *Proceedings of the 34th Annual ACM Conference on Human Factors in Computing Systems (CHI’16)*, CHI 2016: Conference on Human Factors in Computing Systems, San Jose, CA, USA, Association for Computing Machinery, May 2016, pp. 1324–1337, ISBN: 9781450333627. DOI: 10.1145/2858036.2858378. [Online]. Available: <https://doi.org/10.1145/2858036.2858378>.
- [9] S. Poddar, P. Koley, J. Misra, S. Podder, N. Ganguly, and S. Ghosh, *Towards sustainable nlp: Insights from benchmarking inference energy in large language models*, 2025. arXiv: 2502.05610 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2502.05610>.
- [10] J. Fernandez, C. Na, V. Tiwari, Y. Bisk, S. Luccioni, and E. Strubell, *Energy considerations of large language model inference and efficiency optimizations*, 2025. arXiv: 2504.17674 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2504.17674>.
- [11] B. Li, Y. Jiang, V. Gadepally, and D. Tiwari, *Toward sustainable genai using generation directives for carbon-friendly large language model inference*, 2024. arXiv: 2403.12900 [cs.DC]. [Online]. Available: <https://arxiv.org/abs/2403.12900>.
- [12] J. Stojkovic, C. Zhang, Í. Goiri, J. Torrellas, and E. Choukse, *Dynamollm: Designing llm inference clusters for performance and energy efficiency*, 2024. arXiv: 2408.00741 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2408.00741>.
- [13] D. Geelen, R. Mugge, S. Silvester, and A. Bulters, „The use of apps to promote energy saving: A study of smart-meter-related feedback in the netherlands“, *Energy Efficiency*, vol. 12, no. 6, pp. 1635–1660, 2019. DOI: 10.1007/s12053-019-09777-z. [Online]. Available: <https://doi.org/10.1007/s12053-019-09777-z>.

- [14] S. Darby, „The effectiveness of feedback on energy consumption“, *A Review for DEFRA of the Literature on Metering, Billing and direct Displays*, vol. 486, no. 3, pp. 93–109, Jan. 2006, immediate feedback = +5% savings. [Online]. Available: <https://smartgridawareness.org/wp-content/uploads/2016/05/effectiveness-of-feedback-on-energy-consumption-darby-2006.pdf>.
- [15] R. Likert, „A technique for the measurement of attitudes“, *Archives of Psychology*, vol. 22, no. 140, pp. 1–55, 1932.
- [16] B. Kitchenham, S. L. Pfleeger, and L. Madeyski, „Guidelines for performing experiments in software engineering“, Keele University & Durham University, Joint Technical Report, Staffordshire, UK, Tech. Rep. EBSE-TR-2013-003, 2013, Version 2.0, updated January 2013. [Online]. Available: <https://www.ebse.org.uk/>.
- [17] J. Walters, A. Nair, P. Mastronardi, A. Li, and X. Bai, „Purple: Combining individual and collective action to increase online sustainability“, in *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, ser. CHI EA 25, New York, NY, USA: Association for Computing Machinery, 2025, ISBN: 9798400713958. DOI: 10.1145/3706599.3720304. [Online]. Available: <https://doi.org/10.1145/3706599.3720304>.
- [18] J. Chen, G. Fu, Z. Ren, M. Li, and J. Ham, „Effects of anthropomorphic design cues of chatbots on users’ perception and visual behaviors“, *International Journal of Human–Computer Interaction*, vol. 40, no. 14, pp. 3636–3654, Jul. 2024. DOI: 10.1080/10447318.2023.2193514. [Online]. Available: <https://doi.org/10.1080/10447318.2023.2193514>.
- [19] J. Kelly and W. Knottenbelt, *Does disaggregated electricity feedback reduce domestic electricity consumption? a systematic review of the literature*, 2016. arXiv: 1605.00962 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/1605.00962>.
- [20] H. Allcott and T. Rogers, „The short-run and long-run effects of behavioral interventions: Experimental evidence from energy conservation“, *American Economic Review*, vol. 104, no. 10, pp. 3003–37, Oct. 2014. DOI: 10.1257/aer.104.10.3003. [Online]. Available: <https://www.aeaweb.org/articles?id=10.1257/aer.104.10.3003>.
- [21] S. LeVine, *How peer pressure can help save the planet*, <https://www.axios.com/2019/01/30/harvard-business-review-peer-pressure-energy-conservation>, Accessed: 2025-07-06, 2019.
- [22] B. Huseynli, „Gamification in energy consumption: A model for consumers’ energy saving“, *International Journal of Energy Economics and Policy*, vol. 14, no. 1, pp. 312–320, 2024, proposes gamified energy-saving model. [Online]. Available: <http://dx.doi.org/10.32479/ijeep.14395>.
- [23] Epoch AI, „How much energy does chatgpt use?“ (Feb. 10, 2025), [Online]. Available: <https://epoch.ai/gradient-updates/how-much-energy-does-chatgpt-use> (visited on Jun. 11, 2025).
- [24] S. Altman, *The gentle singularity*, <https://blog.samaltman.com/the-gentle-singularity>, 2025.
- [25] E. Masanet, A. Shehabi, N. Lei, S. Smith, and J. Koomey, „Recalibrating global data center energy-use estimates“, *Environmental Research Letters*, vol. 13, no. 1, p. 014003, 2018, Relevant for recalculating per-search energy based on efficiency gains. DOI: 10.1088/1748-9326/aa9671. [Online]. Available: [https://datacenters.lbl.gov/sites/default/files/Masanet\\_et\\_al\\_Science\\_2020.full\\_.pdf](https://datacenters.lbl.gov/sites/default/files/Masanet_et_al_Science_2020.full_.pdf).
- [26] AppleInsider, *Apple’s iphone 14 battery capacities revealed in filing*, 2022.
- [27] N. Kumar, „65 chatbot statistics for 2025 — new data released“. DemandSage. (Jan. 28, 2025), [Online]. Available: <https://www.demandsage.com/chatbot-statistics/> (visited on Mar. 17, 2025).

- [28] S. Samsi, D. Zhao, J. McDonald, *et al.*, *From words to watts: Benchmarking the energy costs of large language model inference*, arXiv preprint 2310.03003, 2023. DOI: 10.48550/arXiv.2310.03003. [Online]. Available: <https://arxiv.org/abs/2310.03003>.
- [29] B. Li, Y. Jiang, V. Gadepally, and D. Tiwari, „Sprout: Green generative AI with carbon-efficient LLM inference“, in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds., Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 21 799–21 813. DOI: 10.18653/v1/2024.emnlp-main.1215. [Online]. Available: <https://aclanthology.org/2024.emnlp-main.1215/>.
- [30] T. Shi, Y. Wu, S. Liu, and Y. Ding, *Greenllm: Disaggregating large language model serving on heterogeneous gpus for lower carbon emissions*, 2024. arXiv: 2412.20322 [cs.AR]. [Online]. Available: <https://arxiv.org/abs/2412.20322>.
- [31] E. J. Husom, A. Goknil, L. K. Shar, and S. Sen, *The price of prompting: Profiling energy use in large language models inference*, 2024. arXiv: 2407.16893 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2407.16893>.
- [32] Z. Fu, F. Chen, S. Zhou, H. Li, and L. Jiang, *Llmco2: Advancing accurate carbon footprint prediction for llm inferences*, 2024. arXiv: 2410.02950 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2410.02950>.
- [33] M. Argerich and M. Patiño-Martínez, „Measuring and improving the energy efficiency of large language models inference“, *IEEE Access*, vol. PP, pp. 1–1, Jan. 2024. DOI: 10.1109/ACCESS.2024.3409745.
- [34] R. Rubei, A. Moussaid, C. di Sipio, and D. di Ruscio, *Prompt engineering and its implications on the energy consumption of large language models*, 2025. arXiv: 2501.05899 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2501.05899>.
- [35] A. Nik, M. A. Riegler, and P. Halvorsen, „Energy-conscious llm decoding: Impact of text generation strategies on gpu energy consumption“, 2025. DOI: 10.48550/arXiv.2502.11723. arXiv: 2502.11723 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2502.11723>.
- [36] G. Wilkins, S. Keshav, and R. Mortier, *Offline energy-optimal llm serving: Workload-based energy models for llm inference on heterogeneous systems*, 2024. arXiv: 2407.04014 [cs.DC]. [Online]. Available: <https://arxiv.org/abs/2407.04014>.
- [37] G. Wilkins, S. Keshav, and R. Mortier, *Hybrid heterogeneous clusters can lower the energy consumption of llm inference workloads*, 2024. arXiv: 2407.00010 [cs.DC]. [Online]. Available: <https://arxiv.org/abs/2407.00010>.
- [38] M. Adamska, D. Smirnova, H. Nasiri, Z. Yu, and P. Garraghan, *Green prompting*, 2025. arXiv: 2503.10666 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2503.10666>.
- [39] L. Solovyeva, S. Weidmann, and F. Castor, *Ai-powered, but power-hungry? energy efficiency of llm-generated code*, 2025. DOI: 10.48550/arXiv.2502.02412. arXiv: 2502.02412 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2502.02412>.
- [40] T. Coignion, C. Quinton, and R. Rouvoy, *Green my llm: Studying the key factors affecting the energy consumption of code assistants*, 2024. DOI: 10.48550/arXiv.2411.11892. arXiv: 2411.11892 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2411.11892>.
- [41] V. Nguyen, H. Huynh, V. Dhopate, *et al.*, „On-device or remote? on the energy efficiency of fetching llm-generated content“, in *CAIN 2025: Architecting and Testing AI Systems*, Empirical comparison showing remote server inference uses 4–9× less client energy than on-device LLMs, 2025.
- [42] A. Isaza-Giraldo, P. Bala, P. Campos, and L. Pereira, „Prompt-gaming: A pilot study on llm-evaluating agent in a meaningful energy game“, May 2024, pp. 1–12. DOI: 10.1145/3613905.3650774.



- [43] L. Anthony, B. Kanding, and R. Selvan, „Carbontracker: Tracking and predicting the carbon footprint of training deep learning models“, in *Proceedings of the ICML 2020 Workshop on Energy Efficient ML*, 2020. [Online]. Available: <https://arxiv.org/abs/2007.03051>.
- [44] B. Courty, V. Schmidt, S. Luccioni, *et al.*, *Mlco2/codecarbon: V2.4.1*, version v2.4.1, May 2024. DOI: 10.5281/zenodo.11171501. [Online]. Available: <https://doi.org/10.5281/zenodo.11171501>.
- [45] ScaleDown Team, *Scaledown, Ai productivity & sustainability suite — chrome extension*, Browser extension, version 0.2.4, Available on the Chrome Web Store, Apr. 29, 2025. [Online]. Available: <https://chromewebstore.google.com/detail/scaledown/jofapkamgblhjaajlppnaiomcjhnllhnd> (visited on Jun. 11, 2025).
- [46] Pascal, *Ai wattch, Track chatgpt's carbon emissions — chrome extension*, Browser extension, version 1.5, Available on the Chrome Web Store, Jun. 2, 2025. [Online]. Available: <https://chromewebstore.google.com/detail/ai-wattch-%E2%80%93-track-chatgpt/meacendfnhnjbkmfbfobgmekkhnamffn> (visited on Jun. 11, 2025).
- [47] ScaleDown Team. „Methodology & accuracy notes“. Accessed 11 Jun 2025. (), [Online]. Available: <https://scaledown.ai/methodology> (visited on Jun. 11, 2025).
- [48] S. K. Dam, C. S. Hong, Y. Qiao, and C. Zhang, *A complete survey on llm-based ai chatbots*, 2024. arXiv: 2406.16937 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2406.16937>.
- [49] comparis.ch AG, *Medienmitteilung: Zwei drittel der schweizerinnen und schweizer haben schon chatgpt oder gemini genutzt*, Presseportal, Zürich, 18. März 2025, Mar. 2025. [Online]. Available: <https://www.presseportal.ch/de/pm/100003671/100929653>.
- [50] L. F. W. Anthony, B. Kanding, and R. Selvan, *Carbontracker: Tracking and predicting the carbon footprint of training deep learning models*, 2020. arXiv: 2007.03051 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2007.03051>.
- [51] T. Rist and M. Masoodian, „Promoting sustainable energy consumption behavior through interactive data visualizations“, *Multimodal Technologies and Interaction*, vol. 3, no. 3, 2019, ISSN: 2414-4088. DOI: 10.3390/mti3030056. [Online]. Available: <https://www.mdpi.com/2414-4088/3/3/56>.
- [52] CHEMTREC®, „Apis & bpis report (updated 09 sep 2022)“. Battery capacity reference used for iPhone 14. (2022), [Online]. Available: [https://app.clickdimensions.com/blob/chemtreccom-ajrbl/files/apis\\_bpisreportupdated20220909.pdf](https://app.clickdimensions.com/blob/chemtreccom-ajrbl/files/apis_bpisreportupdated20220909.pdf) (visited on Jul. 1, 2025).
- [53] Coolblue. „Stromverbrauch kühlsschrank: Wie viel watt braucht ein kühlsschrank?“ (2023), [Online]. Available: <https://www.coolblue.de/beratung/verbrauch-kuehlsschrank.html> (visited on Jul. 1, 2025).
- [54] U. Blog. „Stromverbrauch von laptops: Wie viel watt verbraucht ein laptop?“ (2024), [Online]. Available: <https://de.ugreen.com/blogs/powerstation/stromverbrauch-von-laptops-wie-viel-watt-verbraucht-ein-laptop> (visited on Jul. 1, 2025).
- [55] Lampen&Leuchten. „Nachhaltigkeit und umwelt: Alles über led-beleuchtung“. (2024), [Online]. Available: <https://www.lampenundleuchten.ch/c/beratung/alles-uber-led-beleuchtung/nachhaltigkeit-und-umwelt> (visited on Jul. 1, 2025).
- [56] Verivox. „Ps5: So hoch ist der stromverbrauch“. (2024), [Online]. Available: <https://www.verivox.de/strom/ratgeber/ps5-so-hoch-ist-der-stromverbrauch-1120593/> (visited on Jul. 1, 2025).
- [57] N. Jegham, M. Abdelatti, L. Elmoubarki, and A. Hendawi, „How hungry is ai? benchmarking energy, water, and carbon footprint of llm inference“, 2025. arXiv: 2505.09598 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2505.09598>.
- [58] OpenAI, *Openai api pricing and tokenization*, <https://openai.com/pricing>, 2023.

- [59] EpochAI, *Energy efficiency benchmarking of llms: Gpt-4, claude, gemini, and more*, <https://epochai.org/blog/energy-benchmarking-llms>, Accessed: 2025-07-01, 2024.
- [60] OpenAI, *Tokenizer documentation*, <https://platform.openai.com/tokenizer>, 2023.
- [61] M. Levy, A. Jacoby, and Y. Goldberg, *Same task, more tokens: The impact of input length on the reasoning performance of large language models*, 2024. arXiv: 2402.14848 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2402.14848>.
- [62] Microsoft, *Cloud native apps on azure*, <https://learn.microsoft.com/en-us/dotnet/architecture/cloud-native/>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [63] GitHub, *Github actions for ci/cd*, <https://docs.github.com/en/actions/automating-builds-and-tests/about-continuous-integration>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [64] Microsoft Corporation, *Auswählen zwischen herkömmlichen webanwendungen und single-page-webanwendungen (spas)*, <https://learn.microsoft.com/de-de/dotnet/architecture/modern-web-apps-azure/choose-between-traditional-web-and-single-page-apps>, Accessed: 2025-07-06, Jun. 2023.
- [65] Microsoft, *Host a static website in azure storage*, <https://learn.microsoft.com/en-us/azure/storage/blobs/storage-blob-static-website>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [66] Microsoft, *Azure functions overview*, <https://learn.microsoft.com/en-us/azure/azure-functions/functions-overview>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [67] OpenAI, *Openai api documentation*, <https://platform.openai.com/docs/introduction>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [68] Microsoft, *Microsoft entra id documentation*, <https://learn.microsoft.com/en-us/azure/active-directory/fundamentals/active-directory-what-is>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [69] Microsoft, *Monitor azure functions with application insights*, <https://learn.microsoft.com/en-us/azure/azure-functions/functions-monitoring>, Zuletzt abgerufen am 28.06.2025. (visited on Jun. 28, 2025).
- [70] Angular Team, *Angular v19 overview*, <https://v19.angular.dev/overview>, Zugriff am 5. Juli 2025, 2025.
- [71] M. Maheshwari and S. Agarwal, „Energy consumption analysis of display technologies on mobile devices“, *International Journal of Computer Applications*, vol. 176, no. 35, pp. 6–11, 2020. DOI: 10.5120/ijca2020919984.
- [72] Microsoft Docs. „Authentication and authorization in azure static web apps“. Reference for `/auth/me`, `/auth/login/aad`, and `/auth/logout`. (2024), [Online]. Available: <https://learn.microsoft.com/en-us/azure/static-web-apps/authentication-authorization> (visited on Jul. 1, 2025).
- [73] A. Team, *@angular/core*, <https://www.npmjs.com/package/@angular/core>, Version 17 (aktuell bei Abfrage), 2025. (visited on Jun. 28, 2025).
- [74] A. Team, *@angular/material*, <https://www.npmjs.com/package/@angular/material>, Material Design components for Angular, 2025. (visited on Jun. 28, 2025).
- [75] C. Contributors, *Chart.js*, <https://www.npmjs.com/package/chart.js>, Simple yet flexible JavaScript charting library, 2025. (visited on Jun. 28, 2025).
- [76] M. P. Contributors, *Marked*, <https://www.npmjs.com/package/marked>, A markdown parser and compiler, 2025. (visited on Jun. 28, 2025).
- [77] R. Contributors, *Rxjs*, <https://www.npmjs.com/package/rxjs>, Reactive Extensions Library for JavaScript, 2025. (visited on Jun. 28, 2025).

- [78] Microsoft Learn, *Common web application architectures*, Accessed: 2025-07-09, 2023. [Online]. Available: <https://learn.microsoft.com/en-us/dotnet/architecture/modern-web-apps-azure/common-web-application-architectures>.
- [79] Microsoft Learn, *Unit testing best practices*, Accessed: 2025-07-09, 2023. [Online]. Available: <https://learn.microsoft.com/en-us/dotnet/core/testing/unit-testing-best-practices>.
- [80] OpenAI, „Chat - streaming“. Zuletzt abgerufen am 28.06.2025. (2024), [Online]. Available: <https://platform.openai.com/docs/api-reference/chat-streaming> (visited on Jun. 28, 2025).
- [81] Mozilla Developer Network. „Server-sent events - web apis — mdn“. Zuletzt abgerufen am 28.06.2025. (2024), [Online]. Available: [https://developer.mozilla.org/en-US/docs/Web/API/Server-sent\\_events](https://developer.mozilla.org/en-US/docs/Web/API/Server-sent_events) (visited on Jun. 28, 2025).
- [82] Microsoft. „Configure azure functions app settings“. Zuletzt abgerufen am 28.06.2025. (2024), [Online]. Available: <https://learn.microsoft.com/en-us/azure/azure-functions/functions-how-to-use-azure-function-app-settings> (visited on Jun. 28, 2025).
- [83] Microsoft, *Databases, containers, and items – azure cosmos db*, <https://learn.microsoft.com/azure/cosmos-db/resource-model>, 2024.
- [84] Microsoft, *Authentication and authorization in azure static web apps*, <https://learn.microsoft.com/en-us/azure/static-web-apps/authentication-authorization>, 2024.
- [85] Microsoft, *Manage roles and access control in azure static web apps*, <https://learn.microsoft.com/en-us/azure/static-web-apps/authentication-authorization?tabs=github-actions#roles>, 2024.
- [86] M. Fowler, *Unit test*, <https://martinfowler.com/bliki/UnitTest.html>, 2024.
- [87] Microsoft, *Azure functions - unit testing guidance*, <https://learn.microsoft.com/en-us/azure/azure-functions/functions-test-a-function>, 2024.
- [88] Angular, *Testing angular applications*, <https://angular.dev/guide/testing>, 2024.
- [89] Jasmine Team, *Jasmine – behavior-driven javascript*, Accessed: 2025-07-31, 2025. [Online]. Available: <https://jasmine.github.io/>.
- [90] Karma Team, *Karma: Test runner for javascript*, <https://www.npmjs.com/package/karma>, version 6.4.4, Last published a year ago; accessed on 2025-07-31, 2025.
- [91] Microsoft, *Azure static web apps documentation*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/static-web-apps/overview>.
- [92] Microsoft, *Azure front door documentation*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/frontdoor/front-door-overview>.
- [93] Microsoft, *Azure functions scale and hosting options*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/azure-functions/functions-scale>.
- [94] Microsoft, *Azure functions limits*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/azure-functions/functions-scale#limits>.
- [95] Microsoft, *Globally distributed applications with azure cosmos db*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/cosmos-db/global-distribution>.
- [96] Microsoft, *Autoscale provisioned throughput in azure cosmos db*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/cosmos-db/provisioned-throughput-autoscale>.

- [97] OpenAI, *Openai api documentation*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://platform.openai.com/docs/guides/completions>.
- [98] OpenAI, *How openai scales chatgpt*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://openai.com/blog/scaling-chatgpt>.
- [99] Microsoft, *Resiliency patterns and best practices in azure*, Accessed: 2025-07-05, 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/architecture/resiliency/>.
- [100] Microsoft, *Build and deploy a static web app using github actions*, <https://learn.microsoft.com/en-us/azure/static-web-apps/deploy-github-actions>, 2024.
- [101] Microsoft, *Azure functions github action*, <https://github.com/Azure/functions-action>, 2024.
- [102] Microsoft, *Infrastructure as code on azure*, <https://learn.microsoft.com/en-us/azure/devops/learn/devops-at-microsoft/infrastructure-as-code>, 2024.
- [103] Microsoft, *Iac strategy, best practices, and tools*, <https://learn.microsoft.com/en-us/azure/cloud-adoption-framework/ready/azure-best-practices/infra-as-code/>, 2024.
- [104] MathWorks, *Lowess Smoothing – MATLAB & Simulink*, <https://www.mathworks.com/help/curvefit/lowess-smoothing.html>, Accessed: 2025-07-12, MathWorks, 2025.

## Declaration of honesty

I hereby declare that any individual/pair/team- work submitted for assessment is entirely the product of my own, my partner's and my team's effort that we have correctly cited all text passages that do not originate from us, in accordance with standard academic citation rules (e.g. IEEE), and that we have clearly mentioned all sources used; that we have declared in the table Table 8.1 all aids used; that we have acquired all intangible rights to any materials we may have used, such as images or graphics, or that these materials were created by us; that the topic, the thesis or parts of it have not been used in an assessment of another module, unless this has been expressly agreed with the lecturer in advance and is stated as such; that we are aware that our work may be checked for plagiarism and for third-party authorship of human or technical origin; that we are aware that the FHNW School of Engineering and Environment (FHNW School of Computer Science) will pursue a violation of this declaration of authenticity and that disciplinary consequences may result from this.

Tool	Usage
AI-Chatbots(ChatGPT, Gemini)	Creativity & Brainstorming assistance and paraphrasing
Copilot & Codex	Coding assistance
MatLab	Generating graphics and statistical analysis
Microsoft Excel	Statistical analysis and cleaning data
Overleaf	Writing thesis and support with LaTeX
Visual Studio Code	Developing web application

**Table 8.1:** External tools used during thesis work

Windisch, 31. Juli 2025

**Name:** Jack Gläser

**Signature:**



**Name:** Simon Lüscher

**Signature:**



## **A Appendix**

### **A.1 Survey Form and Results**

# Conversational AI & Energy: How Much Do You Know?

Conversational AI and LLM chatbots are everywhere, making life easier and conversations smarter. But have you ever wondered about the energy they consume? Every chat, every request, and every response requires computing power—and that means real-world energy use.

**Did you know** that with just 5 prompts you could power a 1 W device like a small speaker for 1 hour 💡

**That makes us curious.**

How aware are you of conversational AI's resource consumption? This short survey will take **less than five minutes**, and your insights will help us understand public awareness about conversational AI's environmental impact.

This survey is a part of the research project "**Designing towards higher user awareness: UI strategies for reducing conversational AI energy consumption**" at the School of Computer Science (HSI) of the University of Applied Sciences and Arts Northwestern Switzerland (FHNW).

Our aim is to better understand how people use conversational AI or LLM chatbots, shed light on their environmental impact, and discover what features might help to use them more sustainably. Therefore, your voice matters!

*The survey will be open for a month, **from 18.03.2025 to 18.04.2025**.*

*You can complete it anonymously; if you'd like to see the results, please leave your email address at the end.*

Note that we will use the terms "**conversational AI**" and "**LLM chatbot**" interchangeably in this survey.

👉 Start the survey now!

---

\* Gibt eine erforderliche Frage an

## Usage Patterns - Part 1

***First, tell us about your experience and habits when using conversational AI tools.***

1. How would you describe your overall technical proficiency? \*

*Markieren Sie nur ein Oval.*

- ☐ Beginner (little to no technical knowledge)
- ☐ Intermediate (comfortable using the technology for basic tasks)
- ☐ Advanced (can code or work with AI tools)

2. Which of the following best describes your familiarity with conversational AI (e.g., ChatGPT, Bard, Claude, etc.)? \*

*Markieren Sie nur ein Oval.*

- ☐ I've never heard of them (no familiarity) *Fahren Sie mit Frage 11 fort*
- ☐ I've heard of them but never used them (aware, but not a user)  
*Fahren Sie mit Frage 11 fort*
- ☐ I've used them occasionally (infrequent user) *Fahren Sie mit Frage 3 fort*
- ☐ I'm a regular user (frequent user) *Fahren Sie mit Frage 3 fort*
- ☐ I'm an advanced user (heavy/integrated usage) *Fahren Sie mit Frage 3 fort*

### Usage Patterns - Part 2

*Now, we would like to get to know your conversational AI usage in more detail.*

3. How frequently do you use conversational AI tools? \*

*Markieren Sie nur ein Oval.*

- ☐ Several times a day
- ☐ A few times a week
- ☐ A few times a month
- ☐ Once a few months



4. What types of conversational AI services do you use? (please select all that apply) \*

*Wählen Sie alle zutreffenden Antworten aus.*

- ☐ Text-based AI (e.g., ChatGPT, Google Bard)
- ☐ Image / Video generation AI (e.g., DALL·E, MidJourney, Runway)
- ☐ Code assistance (e.g., GitHub Copilot, Tabnine)
- ☐ AI-driven search engines (e.g., Perplexity, Google Search AI)
- ☐ Speech recognition/assistants (e.g., Siri, Alexa, Google Assistant)
- ☐ Sonstiges: \_\_\_\_\_

5. How do you **mainly** work with LLMs? \*

*Markieren Sie nur ein Oval.*

- ☐ I use tools, like ChatGPT, Claude etc directly
- ☐ I directly interact with the LLMs via API or custom integrations
- ☐ Sonstiges: \_\_\_\_\_

6. What are your primary reasons for using LLM chatbots? (please select all that apply) \*

*Wählen Sie alle zutreffenden Antworten aus.*

- ☐ Searching for information
- ☐ Research & learning
- ☐ Creative writing/content generation
- ☐ Coding assistance
- ☐ Entertainment & casual conversations
- ☐ Personal organization (e.g., summarizing, scheduling)
- ☐ Sonstiges: \_\_\_\_\_

7. Roughly how much time do you spend **actively** using such tools or interacting with LLMs **on a typical day**? \*

*Markieren Sie nur ein Oval.*

- ☐ Less than half an hour
- ☐ Half - One hour
- ☐ 1 - 2 hours
- ☐ More than 2 hours
- ☐ I am not sure how to measure my usage time
- ☐ I don't see the necessity to track my time spent

8. Which device do you **primarily** use to access LLM chatbots? \*

*Markieren Sie nur ein Oval.*

- ☐ Desktop/Laptop
- ☐ Smartphone
- ☐ Tablet

9. Do you use a paid license for access to an LLM chatbot (e.g., ChatGPT Plus, Copilot Pro)? \*

*Markieren Sie nur ein Oval.*

- ☐ Yes
- ☐ No, but I thought about it
- ☐ No, I don't see the reason for it currently
- ☐ No, I never will

10. How important are LLM chatbots in your daily workflow? \*

*Markieren Sie nur ein Oval.*

- ☐ Essential
- ☐ Useful but not critical
- ☐ Nice to have but rarely needed
- ☐ Not important

### Environmental Awareness

*Now, we would like to explore your knowledge and perceptions of the energy use and environmental impact of conversational AI*

11. On a scale of 1–5, how concerned are you about the environmental impact of technology in general? \*

*Markieren Sie nur ein Oval.*

- 1   2   3   4   5
- Not ☐ ☐ ☐ ☐ ☐ Extremely concerned

12. Are you aware that LLM chatbots consume significant energy and have a massive environmental impact? \*

*Markieren Sie nur ein Oval.*

- ☐ Yes, I am well aware
- ☐ I am somewhat aware
- ☐ No, I did not know

13. Compared to a typical Google search, how much energy do you believe a single ChatGPT prompt consumes roughly? \*

*Markieren Sie nur ein Oval.*

- ☐ Less energy than a typical Google search
- ☐ About the same amount of energy as a Google search
- ☐ Around 2× more energy
- ☐ Around 5× more energy
- ☐ Around 10× more energy
- ☐ I'm not sure / I can't estimate

14. How many ChatGPT queries do you think use roughly the same amount of energy as to fully charge an iPhone 14 (from 0% to 100%)? \*

*Markieren Sie nur ein Oval.*

- ☐ About 20 queries
- ☐ About 50 queries
- ☐ About 100 queries
- ☐ About 1000 queries
- ☐ About 10'000 queries

15. Do you agree that LLM chatbots should generally be optimised to reduce energy consumption? \*

*Markieren Sie nur ein Oval.*

	1	2	3	4	5	
Stro	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Strongly agree

16. Currently, AI companies don't disclose a lot of information about the energy consumption of their models. Do you agree that AI companies should be more transparent about the environmental impact of their models and products? \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

Stro ☐ ☐ ☐ ☐ ☐ Strongly agree

### Reimagine Conversational AI/chatbots

***Finally, we'd love your input on potential sustainability features and how they might influence your usage.***

17. On a scale of 1–5, how important is the environmental impact of conversational AI in your decision to use any such services? \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

Not ☐ ☐ ☐ ☐ ☐ Extremely important

18. Would you like LLM chatbots to provide an "Eco Mode" that reduces computational power for less demanding queries? \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

Not ☐ ☐ ☐ ☐ ☐ Extremely interested

19. Would you prefer to use an LLM chatbot that demonstrates a smaller carbon footprint, even if it is slower or less feature-rich? \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

Do not ☐ ☐ ☐ ☐ ☐ Strongly prefer

20. How important is it for you to see energy consumption information related to your conversational AI usage? \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

Not ☐ ☐ ☐ ☐ ☐ Extremely important

21. If such usage information was provided, would it influence how you use LLM chatbots? \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

It would not ☐ ☐ ☐ ☐ ☐ It would strongly influence

22. Would you like to set limits on your LLM chatbot usage based on environmental impact? \*

*Markieren Sie nur ein Oval.*

☐ Yes

☐ No

☐ Unsure

23. Would you agree to pay a small fee or donation to offset the carbon emissions of your conversational AI usage? \*

*Markieren Sie nur ein Oval.*

	1	2	3	4	5	
Stro	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Strongly agree

24. Do you have any other thoughts on how conversational AI could be optimized for sustainability?

---

---

---

---

---

### About you

***Before wrapping everything up, please tell us a bit about yourself.***

25. Which age group do you belong to? \*

*Markieren Sie nur ein Oval.*

- ☐ 18 - 29
- ☐ 30 - 44
- ☐ 45 - 60
- ☐ above 60

26. What is your gender? \*

*Markieren Sie nur ein Oval.*

- ☐ Male
- ☐ Female
- ☐ Diverse
- ☐ Sonstiges: \_\_\_\_\_

27. Which of the following best describes your current role in your organization? \*

*Markieren Sie nur ein Oval.*

- ☐ Student / Intern
- ☐ Academic / Researcher / Educator
- ☐ Developer / Engineer
- ☐ Data Scientist / Analyst
- ☐ Manager / Team Lead
- ☐ Consultant / Advisor
- ☐ Operations / Administrative
- ☐ Non-technical / Business-focused
- ☐ Sonstiges: \_\_\_\_\_



28. In which business domain are you primarily working currently? \*

*Markieren Sie nur ein Oval.*

- ☐ Education / Academia
- ☐ Software / IT Services
- ☐ AI / Machine Learning
- ☐ Finance / Banking
- ☐ Healthcare / Life Sciences
- ☐ Manufacturing / Industry
- ☐ Government / Public Sector
- ☐ Marketing / E-commerce
- ☐ Media / Entertainment
- ☐ Nonprofit / Social Sector
- ☐ Sonstiges: \_\_\_\_\_

29. ***That's it, thanks a ton!***

If you're interested in receiving the results of this survey please leave your email below.

\_\_\_\_\_

---

Dieser Inhalt wurde nicht von Google erstellt und wird von Google auch nicht unterstützt.

Google

Formulare

## **A.2 Experiment Forms and Results**

### **A.2.1 Participant's Consent Document**

*Participant Consent Form for Participation of Experiment*  
**Bachelor thesis**  
**UI strategies for reducing conversational AI energy consumption**

You are invited to participate in a study for our bachelor thesis: UI strategies for reducing conversational AI energy consumption that is being conducted by Jack Gläser and Simon Lüscher, supervised by Prof. Martin Kropp and Dr. Nitish Patkar from University of Applied Sciences Northwestern Switzerland.

**Contact email:**

[jack.glaeser@students.fhnw.ch](mailto:jack.glaeser@students.fhnw.ch)

[simon.luescher@students.fhnw.ch](mailto:simon.luescher@students.fhnw.ch)

**1 Purpose**

We are investigating whether specific user-interface features (energy note, 3-mode toggle, usage-metrics dashboard) raise users' awareness of the energy cost of large-language-model (LLM) prompts and nudge more sustainable usage patterns.

**2 What participation involves**

- **Duration:** 5 working days (Mon–Fri).
- **Use of prototype:** Chat with our web-based LLM assistant; all sustainability features are enabled by default.
- **Daily check-in:** One very short question micro-survey (< 1 min).
- **Final survey:** Single questionnaire (~5 min) at the end of Day 5.  
There are no interviews, screen recordings, audio or video captures

**3 Data collected**

**Usage logs:** prompt & response text, token counts, selected mode, feature toggles, timestamps. Encrypted on FHNW Azure servers; deleted 31 Dec 2025; anonymised stats kept 5 y

**Surveys:** daily check-ins, final questionnaire (awareness, usability, etc)

**4 Voluntary participation**

Taking part is entirely voluntary. You may skip questions or withdraw at any time without penalty; all identifiable data will then be erased.

**5 Confidentiality**

Data are stored under a coded participant-ID. Only the research team has access. Publications will report only aggregated, non-identifiable results.

**6 Risks & benefits**

The study poses minimal risk. You may gain insight into the energy footprint of everyday AI usage and help design more sustainable chatbots.

**7 Results dissemination**

Findings will appear in the bachelor thesis and may be submitted to academic venues. An anonymised dataset may be shared openly for reproducibility.

*Experiment by Jack Gläser & Simon Lüscher*

I have read and understood this information. I had the opportunity to ask questions and received satisfactory answers. I voluntarily agree to:

1. Participate in this study.
2. Allow my anonymised data to be analysed and published as described.

**Name and Signature**

**Date**

**A.2.2 Daily check ins**

# IP6 - Experiment: Daily check-in

Daily check-in survey for participants of the Experiment.

Thanks a ton for filling it out! :-)

\* Gibt eine erforderliche Frage an

1. E-Mail-Adresse \*

---

2. At this moment I'm aware of the energy cost of the prompts I sent today. \*

Markieren Sie nur ein Oval.

1 2 3 4 5

stro ☐ ☐ ☐ ☐ ☐ strongly agree

3. The energy note shown under each response made me think about energy use. \*

Markieren Sie nur ein Oval.

1 2 3 4 5

stro ☐ ☐ ☐ ☐ ☐ strongly agree

4. The energy estimate and 3-mode toggle influenced my choice of mode for today's prompts. \*

Markieren Sie nur ein Oval.

1 2 3 4 5

stro ☐ ☐ ☐ ☐ ☐ strongly agree

5. I looked at the usage-metrics dashboard at least once today. \*

*Markieren Sie nur ein Oval.*

☐ yes

☐ no

---

Dieser Inhalt wurde nicht von Google erstellt und wird von Google auch nicht unterstützt.

Google

**A.2.3 Final Questionnaire**



# IP6 - Experiment: Final Questionnaire

Hi there,

Thank you for completing the five-day study with our sustainable-AI chatbot. We have one final favour to ask: please fill out the **end-of-study questionnaire** linked below. It is a little longer than the daily check-ins (about 5–10 minutes) and will give us the insights we need to understand how the features worked for you.

We really appreciate the time and effort you've invested in this project.

Many thanks!

Jack & Simon

*\* Gibt eine erforderliche Frage an*

1. E-Mail-Adresse \*

---

2. I am aware that LLM AI-chatbots consume significant amount of energy \*

*Markieren Sie nur ein Oval.*

1   2   3   4   5

stro ☐ ☐ ☐ ☐ ☐ strongly agree

3. Compared with a Google search, a single ChatGPT-style prompt consumes more energy. \*

*Markieren Sie nur ein Oval.*

1   2   3   4   5

stro ☐ ☐ ☐ ☐ ☐ strongly agree

4. It is important for me to know the energy cost of my AI-chatbot usage. \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

stro ☐ ☐ ☐ ☐ ☐ strongly agree

#### Behaviour & feature experience

5. Knowing the rough energy costs of my prompts influences how I formulate prompts. \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

stro ☐ ☐ ☐ ☐ ☐ strongly agree

6. I consciously shortened or reduced prompts after seeing energy-related information. \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

stro ☐ ☐ ☐ ☐ ☐ strongly agree

7. I actively choose Eco-mode when high accuracy was not required. \*

*Markieren Sie nur ein Oval.*

1 2 3 4 5

---

stro ☐ ☐ ☐ ☐ ☐ strongly agree

8. I checked the sustainability dashboard(Usage metrics) at least 2 times within the period. \*

Markieren Sie nur ein Oval.

- ☐ yes  
☐ no  
☐ not sure

9. The energy-related information distracted me from completing my main task. \*

Markieren Sie nur ein Oval.

	1	2	3	4	5	
stro	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	strongly agree

10. I would like to see similar sustainability features in our mainstream AI-chatbots. \*

Markieren Sie nur ein Oval.

	1	2	3	4	5	
stro	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	strongly agree

11. Which of the experiences features within the experiment would you like to see adopted by mainstream AI-Chatbots?

Wählen Sie alle zutreffenden Antworten aus.

- ☐ Energy info-note on the responses  
☐ Different models / being able to switch models fast based on input  
☐ Usage metrics / some sort of dashboard with usage /energy consumption data  
☐ Estimation on energy costs while typing out a prompt  
☐ None of them

Usability & satisfaction

12. The energy note (referring to how many X-units it took) was easy to understand and its effect was clear. \*

*Markieren Sie nur ein Oval.*

	1	2	3	4	5	
strongly disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	strongly agree

13. The 3-mode toggle was easy to understand and its effect was clear. \*

*Markieren Sie nur ein Oval.*

	1	2	3	4	5	
strongly disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	strongly agree

14. The usage metrics dashboard was easy to understand and its effect was clear \*

*Markieren Sie nur ein Oval.*

	1	2	3	4	5	
strongly disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	strongly agree

15. The usage metrics dashboard was easy to understand \*

*Markieren Sie nur ein Oval.*

	1	2	3	4	5	
strongly disagree	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	strongly agree

16. The usage metrics dashboard helped to increase my overall knowledge on energy consumption of AI-Chatbot usage \*

Markieren Sie nur ein Oval.

1 2 3 4 5

strongly ☐ ☐ ☐ ☐ ☐ strongly agree

17. The placement of the “choose your AI mode” was well placed and easy to understand. \*

Markieren Sie nur ein Oval.

1 2 3 4 5

strongly ☐ ☐ ☐ ☐ ☐ strongly agree

18. The features overall fit well with how I usually work with AI-chatbots and did not further distract me a lot from using the AI-chatbot. \*

Markieren Sie nur ein Oval.

1 2 3 4 5

strongly ☐ ☐ ☐ ☐ ☐ strongly agree

### Open questions

19. Do you have suggestion on raising energy usage or making energy usage more actionable in future AI-chatbot tools? / Do you have a feature idea? 😊 \*

---

---

---

---

---

20. Roughly how many AI-chatbot prompts per day do you make *outside* this prototype? (*numeric*)

\*

---

---

Dieser Inhalt wurde nicht von Google erstellt und wird von Google auch nicht unterstützt.

Google

## A.3 JSON Schema

```

1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "User",
4   "type": "object",
5   "description": "Represents an application user with sustainability
6     settings and preferences.",
7   "properties": {
8     "id": {
9       "type": "string",
10      "description": "Unique identifier of the user entity."
11    },
12    "userId": {
13      "type": "string",
14      "description": "Partition key: ID linking this entity to the user
15        account."
16    },
17    "createdAt": {
18      "type": "string",
19      "format": "date-time",
20      "description": "Timestamp when the user record was created."
21    },
22    "name": {
23      "type": "string",
24      "description": "Display name of the user."
25    },
26    "identity": {
27      "type": "string",
28      "minLength": 1,
29      "description": "External or internal identity identifier."
30    },
31    "enableSustainabilityFeatures": {
32      "type": "boolean",
33      "description": "Indicates whether sustainability features are enabled
34        ."
35    },
36    "mode": {
37      "type": "number",
38      "enum": [ 0, 1, 2 ],
39      "description": "Chat mode preference for the user. 0=EnergyEfficient,
40        1=Balanced, 2=Performance"
41    },
42    "energyUnitId": {
43      "type": [ "string", "null" ],
44      "description": "Reference ID to the user's selected energy unit."
45    }
46  },
47  "required": [
48    "id",
49    "userId",
50    "createdAt",
51    "name",
52    "identity",
53    "enableSustainabilityFeatures",
54    "mode"
55  ],

```

```

52 "additionalProperties": false
53 }

```

```

1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "Conversation",
4   "type": "object",
5   "description": "Represents a user's conversation with the chat bot.",
6   "properties": {
7     "id": {
8       "type": "string",
9       "description": "Unique identifier of the conversation."
10    },
11    "userId": {
12      "type": "string",
13      "description": "Partition key: User ID that owns the conversation."
14    },
15    "createdAt": {
16      "type": "string",
17      "format": "date-time",
18      "description": "Timestamp when the conversation was created."
19    },
20    "title": {
21      "type": "string",
22      "description": "Title for the conversation."
23    },
24    "lastMessageTimestamp": {
25      "type": [ "string", "null" ],
26      "format": "date-time",
27      "description": "Timestamp of the most recent message."
28    },
29    "isDeleted": {
30      "type": [ "boolean", "null" ],
31      "description": "Whether the conversation is marked as deleted."
32    }
33  },
34  "required": [ "id", "userId", "createdAt", "title" ],
35  "additionalProperties": false
36 }

```

```

1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "Prompt",
4   "type": "object",
5   "description": "Represents a user prompt and its associated metadata.",
6   "properties": {
7     "id": {
8       "type": "string",
9       "description": "Unique identifier of the prompt entity."
10    },
11    "userId": {
12      "type": "string",
13      "description": "Partition key: ID of the user who submitted the
14        prompt."
15    },
16    "createdAt": {
17      "type": "string",

```



```

17     "format": "date-time",
18     "description": "Timestamp when the prompt was created."
19 },
20 "conversationId": {
21     "type": [ "string", "null" ],
22     "description": "ID of the conversation this prompt belongs to."
23 },
24 "promptTextHistory": {
25     "type": [ "array", "null" ],
26     "items": {
27         "type": "string"
28     },
29     "description": "List of previous (not sent) texts of this prompt."
30 },
31 "userText": {
32     "type": [ "string", "null" ],
33     "description": "Prompt text provided by the user."
34 },
35 "chatMode": {
36     "type": "number",
37     "enum": [ 0, 1, 2 ],
38     "description": "Operational mode used for this prompt. 0=
        EnergyEfficient, 1=Balanced, 2=Performance"
39 },
40 "modelName": {
41     "type": "string",
42     "description": "Name of the model used to generate the response."
43 },
44 "historyLimit": {
45     "type": "integer",
46     "description": "Maximum number of previous prompts included in the
        request."
47 },
48 "isSent": {
49     "type": [ "boolean", "null" ],
50     "description": "Indicates whether the prompt was successfully sent."
51 },
52 "sentAt": {
53     "type": [ "string", "null" ],
54     "format": "date-time",
55     "description": "Timestamp when the prompt was sent."
56 },
57 "responseText": {
58     "type": [ "string", "null" ],
59     "description": "Text response generated for the prompt."
60 },
61 "usage": {
62     "type": [ "object", "null" ],
63     "description": "Represents token usage and associated energy
        consumption.",
64     "properties": {
65         "numberOfInputTokens": {
66             "type": "integer",
67             "description": "Number of input tokens provided by the user."
68         },
69         "numberOfOutputTokens": {
70             "type": "integer",

```

```

71     "description": "Number of output tokens generated in response."
72   },
73   "usageInWh": {
74     "type": "number",
75     "description": "Estimated energy usage in watt-hours (Wh)."
76   }
77 },
78 "required": [ "numberOfInputTokens", "numberOfOutputTokens", "
    usageInWh" ],
79 "additionalProperties": false
80 }
81 },
82 "required": [ "id", "userId", "createdAt", "chatMode", "modelName", "
    historyLimit" ],
83 "additionalProperties": false
84 }

```

```

1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "Log",
4   "type": "object",
5   "description": "Represents a log entry.",
6   "properties": {
7     "id": {
8       "type": "string",
9       "description": "Unique identifier of the log entry."
10    },
11    "userId": {
12      "type": "string",
13      "description": "Partition key: User ID associated with the log."
14    },
15    "createdAt": {
16      "type": "string",
17      "format": "date-time",
18      "description": "Timestamp when the log was created."
19    },
20    "type": {
21      "type": "number",
22      "enum": [ 0, 1, 2 ],
23      "description": "Type/category of the log message. 0 = Unknown, 1 =
        PageVisit, 2 = SustainabilityModeChange."
24    },
25    "message": {
26      "type": "string",
27      "description": "Descriptive log message content."
28    }
29  },
30  "required": [ "id", "userId", "createdAt", "type", "message" ],
31  "additionalProperties": false
32 }

```

```

1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "EnergyUnit",
4   "type": "object",
5   "description": "Defines a unit for measuring energy consumption.",
6   "properties": {

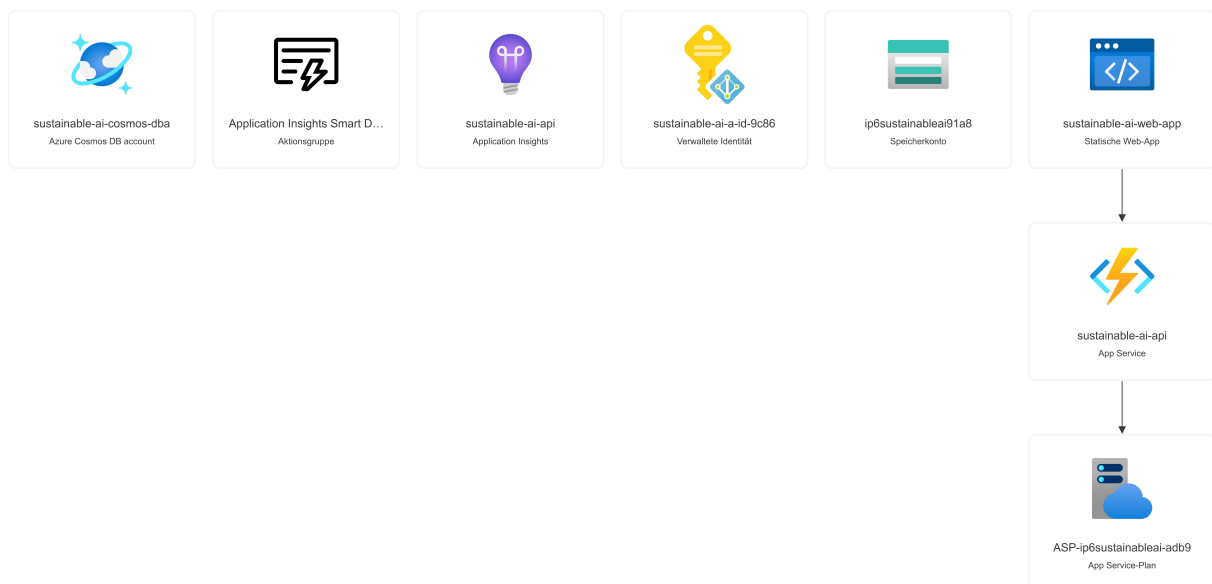
```

```

7  "id": {
8    "type": "string",
9    "description": "Unique identifier for the energy unit."
10 },
11 "tenantId": {
12   "type": "string",
13   "const": "0",
14   "description": "Partition key: Fixed value \"0\""
15 },
16 "name": {
17   "type": "string",
18   "description": "Singular name of the energy unit (e.g., 'phone charge'
19   ')."
20 },
21 "namePlural": {
22   "type": "string",
23   "description": "Plural form of the energy unit name (e.g., 'phone
24   charges')."
25 },
26 "Wh": {
27   "type": "number",
28   "minimum": 0,
29   "description": "Conversion factor to watt-hours (Wh)."
30 },
31 "required": [ "id", "tenantId", "name", "namePlural", "Wh" ],
32 "additionalProperties": false
33 }

```

### A.3.1 Azure Infrastructure



```

1  {
2    "$schema": "https://schema.management.azure.com/schemas/2019-04-01/
3    deploymentTemplate.json#",
4    "contentVersion": "1.0.0.0",
5    "parameters": {

```

```

5     "sites_sustainable_ai_api_name": {
6         "defaultValue": "sustainable-ai-api",
7         "type": "String"
8     },
9     "components_sustainable_ai_api_name": {
10        "defaultValue": "sustainable-ai-api",
11        "type": "String"
12    },
13    "staticSites_sustainable_ai_web_app_name": {
14        "defaultValue": "sustainable-ai-web-app",
15        "type": "String"
16    },
17    "serverfarms_ASP_ip6sustainableai_adb9_name": {
18        "defaultValue": "ASP-ip6sustainableai-adb9",
19        "type": "String"
20    },
21    "storageAccounts_ip6sustainableai91a8_name": {
22        "defaultValue": "ip6sustainableai91a8",
23        "type": "String"
24    },
25    "databaseAccounts_sustainable_ai_cosmos_dba_name": {
26        "defaultValue": "sustainable-ai-cosmos-dba",
27        "type": "String"
28    },
29    "actionGroups_Application_Insights_Smart_Detection_name": {
30        "defaultValue": "Application Insights Smart Detection",
31        "type": "String"
32    },
33    "userAssignedIdentities_sustainable_ai_a_id_9c86_name": {
34        "defaultValue": "sustainable-ai-a-id-9c86",
35        "type": "String"
36    },
37    "workspaces_DefaultWorkspace_37baa613_2be4_4804_8a7e_c2c4a19538a5
38        _CHN_externalid": {
39        "defaultValue": "/subscriptions/37baa613-2be4-4804-8a7e-c2c4a19
40            538a5/resourceGroups/DefaultResourceGroup-CHN/providers/
41            Microsoft.OperationalInsights/workspaces/DefaultWorkspace-37
42            baa613-2be4-4804-8a7e-c2c4a19538a5-CHN",
43        "type": "String"
44    }
45 },
46 "variables": {},
47 "resources": [
48     {
49         "type": "Microsoft.DocumentDB/databaseAccounts",
50         "apiVersion": "2024-12-01-preview",
51         "name": "[parameters('
52             databaseAccounts_sustainable_ai_cosmos_dba_name')]",
53         "location": "Switzerland North",
54         "tags": {
55             "defaultExperience": "Core (SQL)",
56             "hidden-workload-type": "Development/Testing",
57             "hidden-cosmos-mmsspecial": ""
58         },
59         "kind": "GlobalDocumentDB",
60         "identity": {
61             "type": "None"

```

```

57     },
58     "properties": {
59         "publicNetworkAccess": "Enabled",
60         "enableAutomaticFailover": true,
61         "enableMultipleWriteLocations": false,
62         "isVirtualNetworkFilterEnabled": false,
63         "virtualNetworkRules": [],
64         "disableKeyBasedMetadataWriteAccess": false,
65         "enableFreeTier": false,
66         "enableAnalyticalStorage": false,
67         "analyticalStorageConfiguration": {
68             "schemaType": "WellDefined"
69         },
70         "databaseAccountOfferType": "Standard",
71         "enableMaterializedViews": false,
72         "capacityMode": "Serverless",
73         "defaultIdentity": "FirstPartyIdentity",
74         "networkAclBypass": "None",
75         "disableLocalAuth": false,
76         "enablePartitionMerge": false,
77         "enablePerRegionPerPartitionAutoscale": false,
78         "enableBurstCapacity": false,
79         "enablePriorityBasedExecution": false,
80         "defaultPriorityLevel": "High",
81         "minimalTlsVersion": "Tls12",
82         "consistencyPolicy": {
83             "defaultConsistencyLevel": "Session",
84             "maxIntervalInSeconds": 5,
85             "maxStalenessPrefix": 100
86         },
87         "locations": [
88             {
89                 "locationName": "Switzerland North",
90                 "failoverPriority": 0,
91                 "isZoneRedundant": false
92             }
93         ],
94         "cors": [],
95         "capabilities": [],
96         "ipRules": [],
97         "backupPolicy": {
98             "type": "Periodic",
99             "periodicModeProperties": {
100                 "backupIntervalInMinutes": 240,
101                 "backupRetentionIntervalInHours": 8,
102                 "backupStorageRedundancy": "Geo"
103             }
104         },
105         "networkAclBypassResourceIds": [],
106         "diagnosticLogSettings": {
107             "enableFullTextQuery": "None"
108         },
109         "capacity": {
110             "totalThroughputLimit": 4000
111         }
112     }
113 },

```

```

114     {
115         "type": "microsoft.insights/actionGroups",
116         "apiVersion": "2024-10-01-preview",
117         "name": "[parameters('
118             actionGroups_Application_Insights_Smart_Detection_name')]",
119         "location": "Global",
120         "properties": {
121             "groupShortName": "SmartDetect",
122             "enabled": true,
123             "emailReceivers": [],
124             "smsReceivers": [],
125             "webhookReceivers": [],
126             "eventHubReceivers": [],
127             "itsmReceivers": [],
128             "azureAppPushReceivers": [],
129             "automationRunbookReceivers": [],
130             "voiceReceivers": [],
131             "logicAppReceivers": [],
132             "azureFunctionReceivers": [],
133             "armRoleReceivers": [
134                 {
135                     "name": "Monitoring Contributor",
136                     "roleId": "749f88d5-cbae-40b8-bcfc-e573ddc772fa",
137                     "useCommonAlertSchema": true
138                 },
139                 {
140                     "name": "Monitoring Reader",
141                     "roleId": "43d0d8ad-25c7-4714-9337-8ba259a9fe05",
142                     "useCommonAlertSchema": true
143                 }
144             ]
145         }
146     },
147     {
148         "type": "microsoft.insights/components",
149         "apiVersion": "2020-02-02",
150         "name": "[parameters('components_sustainable_ai_api_name')]",
151         "location": "switzerlandnorth",
152         "kind": "web",
153         "properties": {
154             "Application_Type": "web",
155             "Flow_Type": "Redfield",
156             "Request_Source": "IbizaAIExtensionEnablementBlade",
157             "RetentionInDays": 90,
158             "WorkspaceResourceId": "[parameters('
159                 workspaces_DefaultWorkspace_37baa613_2be4_4804_8a7e_c2c4
160                 a19538a5_CHN_externalid')]",
161             "IngestionMode": "LogAnalytics",
162             "publicNetworkAccessForIngestion": "Enabled",
163             "publicNetworkAccessForQuery": "Enabled"
164         }
165     },
166     {
167         "type": "Microsoft.ManagedIdentity/userAssignedIdentities",
168         "apiVersion": "2025-01-31-preview",
169         "name": "[parameters('
170             userAssignedIdentities_sustainable_ai_a_id_9c86_name')]",

```

```

167     "location": "switzerlandnorth"
168   },
169   {
170     "type": "Microsoft.Storage/storageAccounts",
171     "apiVersion": "2024-01-01",
172     "name": "[parameters('storageAccounts_ip6sustainableai91a8_name')]",
173     "location": "switzerlandnorth",
174     "sku": {
175       "name": "Standard_LRS",
176       "tier": "Standard"
177     },
178     "kind": "Storage",
179     "properties": {
180       "defaultToOAuthAuthentication": true,
181       "publicNetworkAccess": "Enabled",
182       "allowCrossTenantReplication": false,
183       "minimumTlsVersion": "TLS1_2",
184       "allowBlobPublicAccess": false,
185       "networkAcls": {
186         "bypass": "AzureServices",
187         "virtualNetworkRules": [],
188         "ipRules": [],
189         "defaultAction": "Allow"
190       },
191       "supportsHttpsTrafficOnly": true,
192       "encryption": {
193         "services": {
194           "file": {
195             "keyType": "Account",
196             "enabled": true
197           },
198           "blob": {
199             "keyType": "Account",
200             "enabled": true
201           }
202         },
203         "keySource": "Microsoft.Storage"
204       }
205     }
206   },
207   {
208     "type": "Microsoft.Web/serverfarms",
209     "apiVersion": "2024-04-01",
210     "name": "[parameters('serverfarms_ASP_ip6sustainableai_adb9_name')]",
211     "location": "Switzerland North",
212     "sku": {
213       "name": "Y1",
214       "tier": "Dynamic",
215       "size": "Y1",
216       "family": "Y",
217       "capacity": 0
218     },
219     "kind": "functionapp",
220     "properties": {
221       "perSiteScaling": false,

```

```

222         "elasticScaleEnabled": false,
223         "maximumElasticWorkerCount": 1,
224         "isSpot": false,
225         "reserved": false,
226         "isXenon": false,
227         "hyperV": false,
228         "targetWorkerCount": 0,
229         "targetWorkerSizeId": 0,
230         "zoneRedundant": false
231     }
232 },
233 {
234     "type": "Microsoft.Web/staticSites",
235     "apiVersion": "2024-04-01",
236     "name": "[parameters('staticSites_sustainable_ai_web_app_name')
237         ]",
238     "location": "West Europe",
239     "sku": {
240         "name": "Standard",
241         "tier": "Standard"
242     },
243     "properties": {
244         "repositoryUrl": "https://github.com/simonluescherfhnw/ip6-
245             sustainable-ai-frontend",
246         "branch": "main",
247         "stagingEnvironmentPolicy": "Enabled",
248         "allowConfigFileUpdates": true,
249         "provider": "GitHub",
250         "enterpriseGradeCdnStatus": "Disabled"
251     }
252 },
253 {
254     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases",
255     "apiVersion": "2024-12-01-preview",
256     "name": "[concat(parameters('
257         databaseAccounts_sustainable_ai_cosmos_dba_name'), '/
258         Development')]",
259     "dependsOn": [
260         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
261             parameters('
262                 databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
263     ],
264     "properties": {
265         "resource": {
266             "id": "Development"
267         }
268     }
269 },
270 {
271     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases",
272     "apiVersion": "2024-12-01-preview",
273     "name": "[concat(parameters('
274         databaseAccounts_sustainable_ai_cosmos_dba_name'), '/
275         Production')]",
276     "dependsOn": [

```



```

269         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
                parameters('
                    databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
270     ],
271     "properties": {
272         "resource": {
273             "id": "Production"
274         }
275     }
276 },
277 {
278     "type": "Microsoft.DocumentDB/databaseAccounts/
                sqlRoleDefinitions",
279     "apiVersion": "2024-12-01-preview",
280     "name": "[concat(parameters('
                databaseAccounts_sustainable_ai_cosmos_dba_name'), '/0000000
                0-0000-0000-0000-000000000001'))]",
281     "dependsOn": [
282         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
                parameters('
                    databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
283     ],
284     "properties": {
285         "roleName": "Cosmos DB Built-in Data Reader",
286         "type": "BuiltInRole",
287         "assignableScopes": [
288             "[resourceId('Microsoft.DocumentDB/databaseAccounts',
                parameters('
                    databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
289         ],
290         "permissions": [
291             {
292                 "dataActions": [
293                     "Microsoft.DocumentDB/databaseAccounts/
                        readMetadata",
294                     "Microsoft.DocumentDB/databaseAccounts/
                        sqlDatabases/containers/executeQuery",
295                     "Microsoft.DocumentDB/databaseAccounts/
                        sqlDatabases/containers/readChangeFeed",
296                     "Microsoft.DocumentDB/databaseAccounts/
                        sqlDatabases/containers/items/read"
297                 ],
298                 "notDataActions": []
299             }
300         ]
301     }
302 },
303 {
304     "type": "Microsoft.DocumentDB/databaseAccounts/
                sqlRoleDefinitions",
305     "apiVersion": "2024-12-01-preview",
306     "name": "[concat(parameters('
                databaseAccounts_sustainable_ai_cosmos_dba_name'), '/0000000
                0-0000-0000-0000-000000000002'))]",
307     "dependsOn": [

```



```

347         "Microsoft.DocumentDB/databaseAccounts/tables/
348             containers/entities/read"
349     ],
350     "notDataActions": []
351 }
352 ]
353 },
354 {
355     "type": "Microsoft.DocumentDB/databaseAccounts/
356         tableRoleDefinitions",
357     "apiVersion": "2024-12-01-preview",
358     "name": "[concat(parameters('
359         databaseAccounts_sustainable_ai_cosmos_dba_name'), '/0000000
360         0-0000-0000-0000-000000000002')]",
361     "dependsOn": [
362         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
363             parameters('
364                 databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
365     ],
366     "properties": {
367         "roleName": "Cosmos DB Built-in Data Contributor",
368         "type": "BuiltInRole",
369         "assignableScopes": [
370             "[resourceId('Microsoft.DocumentDB/databaseAccounts',
371                 parameters('
372                     databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
373         ],
374         "permissions": [
375             {
376                 "dataActions": [
377                     "Microsoft.DocumentDB/databaseAccounts/
378                         readMetadata",
379                     "Microsoft.DocumentDB/databaseAccounts/tables/*",
380                     "Microsoft.DocumentDB/databaseAccounts/tables/
381                         containers/*",
382                     "Microsoft.DocumentDB/databaseAccounts/tables/
383                         containers/entities/*"
384                 ],
385                 "notDataActions": []
386             }
387         ]
388     }
389 },
390 {
391     "type": "microsoft.insights/components/
392         ProactiveDetectionConfigs",
393     "apiVersion": "2018-05-01-preview",
394     "name": "[concat(parameters('components_sustainable_ai_api_name
395         '), '/degradationindependencyduration')]",
396     "location": "switzerlandnorth",
397     "dependsOn": [
398         "[resourceId('microsoft.insights/components', parameters('
399             components_sustainable_ai_api_name'))]"
400     ],
401     "properties": {

```

```

389         "ruleDefinitions": {
390             "Name": "degradationindependencyduration",
391             "DisplayName": "Degradation in dependency duration",
392             "Description": "Smart Detection rules notify you of
                                performance anomaly issues.",
393             "HelpUrl": "https://docs.microsoft.com/en-us/azure/
                                application-insights/app-insights-proactive-
                                performance-diagnostics",
394             "IsHidden": false,
395             "IsEnabledByDefault": true,
396             "IsInPreview": false,
397             "SupportsEmailNotifications": true
398         },
399         "enabled": true,
400         "sendEmailsToSubscriptionOwners": true,
401         "customEmails": []
402     },
403 },
404 {
405     "type": "microsoft.insights/components/
                                ProactiveDetectionConfigs",
406     "apiVersion": "2018-05-01-preview",
407     "name": "[concat(parameters('components_sustainable_ai_api_name
                                '), '/degradationinserverresponsetime')]",
408     "location": "switzerlandnorth",
409     "dependsOn": [
410         "[resourceId('microsoft.insights/components', parameters('
                                components_sustainable_ai_api_name'))]"
411     ],
412     "properties": {
413         "ruleDefinitions": {
414             "Name": "degradationinserverresponsetime",
415             "DisplayName": "Degradation in server response time",
416             "Description": "Smart Detection rules notify you of
                                performance anomaly issues.",
417             "HelpUrl": "https://docs.microsoft.com/en-us/azure/
                                application-insights/app-insights-proactive-
                                performance-diagnostics",
418             "IsHidden": false,
419             "IsEnabledByDefault": true,
420             "IsInPreview": false,
421             "SupportsEmailNotifications": true
422         },
423         "enabled": true,
424         "sendEmailsToSubscriptionOwners": true,
425         "customEmails": []
426     }
427 },
428 {
429     "type": "microsoft.insights/components/
                                ProactiveDetectionConfigs",
430     "apiVersion": "2018-05-01-preview",
431     "name": "[concat(parameters('components_sustainable_ai_api_name
                                '), '/digestMailConfiguration')]",
432     "location": "switzerlandnorth",
433     "dependsOn": [

```

```

434         "[resourceId('microsoft.insights/components', parameters('
435             components_sustainable_ai_api_name'))]"
436     },
437     "properties": {
438         "ruleDefinitions": {
439             "Name": "digestMailConfiguration",
440             "DisplayName": "Digest Mail Configuration",
441             "Description": "This rule describes the digest mail
442                 preferences",
443             "HelpUrl": "www.homail.com",
444             "IsHidden": true,
445             "IsEnabledByDefault": true,
446             "IsInPreview": false,
447             "SupportsEmailNotifications": true
448         },
449         "enabled": true,
450         "sendEmailsToSubscriptionOwners": true,
451         "customEmails": []
452     },
453     {
454         "type": "microsoft.insights/components/
455             ProactiveDetectionConfigs",
456         "apiVersion": "2018-05-01-preview",
457         "name": "[concat(parameters('components_sustainable_ai_api_name
458             '), '/extension_billingdatavolumedailyspikeextension')]",
459         "location": "switzerlandnorth",
460         "dependsOn": [
461             "[resourceId('microsoft.insights/components', parameters('
462                 components_sustainable_ai_api_name'))]"
463         ],
464         "properties": {
465             "ruleDefinitions": {
466                 "Name": "extension_billingdatavolumedailyspikeextension",
467                 "DisplayName": "Abnormal rise in daily data volume (
468                     preview)",
469                 "Description": "This detection rule automatically
470                     analyzes the billing data generated by your
471                     application, and can warn you about an unusual
472                     increase in your application's billing costs",
473                 "HelpUrl": "https://github.com/Microsoft/
474                     ApplicationInsights-Home/tree/master/SmartDetection/
475                     billing-data-volume-daily-spike.md",
476                 "IsHidden": false,
477                 "IsEnabledByDefault": true,
478                 "IsInPreview": true,
479                 "SupportsEmailNotifications": false
480             },
481             "enabled": true,
482             "sendEmailsToSubscriptionOwners": true,
483             "customEmails": []
484         }
485     },
486     {
487         "type": "microsoft.insights/components/
488             ProactiveDetectionConfigs",

```

```

478     "apiVersion": "2018-05-01-preview",
479     "name": "[concat(parameters('components_sustainable_ai_api_name'
480     ' '), '/extension_canaryextension')]",
481     "location": "switzerlandnorth",
482     "dependsOn": [
483         "[resourceId('microsoft.insights/components', parameters('
484         components_sustainable_ai_api_name'))]"
485     ],
486     "properties": {
487         "ruleDefinitions": {
488             "Name": "extension_canaryextension",
489             "DisplayName": "Canary extension",
490             "Description": "Canary extension",
491             "HelpUrl": "https://github.com/Microsoft/
492             ApplicationInsights-Home/blob/master/SmartDetection/
493             ",
494             "IsHidden": true,
495             "IsEnabledByDefault": true,
496             "IsInPreview": true,
497             "SupportsEmailNotifications": false
498         },
499         "enabled": true,
500         "sendEmailsToSubscriptionOwners": true,
501         "customEmails": []
502     },
503     {
504         "type": "microsoft.insights/components/
505         ProactiveDetectionConfigs",
506         "apiVersion": "2018-05-01-preview",
507         "name": "[concat(parameters('components_sustainable_ai_api_name'
508         ' '), '/extension_exceptionchangeextension')]",
509         "location": "switzerlandnorth",
510         "dependsOn": [
511             "[resourceId('microsoft.insights/components', parameters('
512             components_sustainable_ai_api_name'))]"
513         ],
514         "properties": {
515             "ruleDefinitions": {
516                 "Name": "extension_exceptionchangeextension",
517                 "DisplayName": "Abnormal rise in exception volume (
518                 preview)",
519                 "Description": "This detection rule automatically
520                 analyzes the exceptions thrown in your application,
521                 and can warn you about unusual patterns in your
522                 exception telemetry.",
523                 "HelpUrl": "https://github.com/Microsoft/
524                 ApplicationInsights-Home/blob/master/SmartDetection/
525                 abnormal-rise-in-exception-volume.md",
526                 "IsHidden": false,
527                 "IsEnabledByDefault": true,
528                 "IsInPreview": true,
529                 "SupportsEmailNotifications": false
530             },
531             "enabled": true,
532             "sendEmailsToSubscriptionOwners": true,
533             "customEmails": []
534         }
535     }
536 ]

```

```

522     }
523   },
524   {
525     "type": "microsoft.insights/components/
      ProactiveDetectionConfigs",
526     "apiVersion": "2018-05-01-preview",
527     "name": "[concat(parameters('components_sustainable_ai_api_name
      '), '/extension_memoryleakextension')]",
528     "location": "switzerlandnorth",
529     "dependsOn": [
530       "[resourceId('microsoft.insights/components', parameters('
      components_sustainable_ai_api_name'))]"
531     ],
532     "properties": {
533       "ruleDefinitions": {
534         "Name": "extension_memoryleakextension",
535         "DisplayName": "Potential memory leak detected (preview
          )",
536         "Description": "This detection rule automatically
          analyzes the memory consumption of each process in
          your application, and can warn you about potential
          memory leaks or increased memory consumption.",
537         "HelpUrl": "https://github.com/Microsoft/
          ApplicationInsights-Home/tree/master/SmartDetection/
          memory-leak.md",
538         "IsHidden": false,
539         "IsEnabledByDefault": true,
540         "IsInPreview": true,
541         "SupportsEmailNotifications": false
542       },
543       "enabled": true,
544       "sendEmailsToSubscriptionOwners": true,
545       "customEmails": []
546     }
547   },
548   {
549     "type": "microsoft.insights/components/
      ProactiveDetectionConfigs",
550     "apiVersion": "2018-05-01-preview",
551     "name": "[concat(parameters('components_sustainable_ai_api_name
      '), '/extension_securityextensionspackage')]",
552     "location": "switzerlandnorth",
553     "dependsOn": [
554       "[resourceId('microsoft.insights/components', parameters('
      components_sustainable_ai_api_name'))]"
555     ],
556     "properties": {
557       "ruleDefinitions": {
558         "Name": "extension_securityextensionspackage",
559         "DisplayName": "Potential security issue detected (
          preview)",
560         "Description": "This detection rule automatically
          analyzes the telemetry generated by your application
          and detects potential security issues.",
561         "HelpUrl": "https://github.com/Microsoft/
          ApplicationInsights-Home/blob/master/SmartDetection/
          application-security-detection-pack.md",

```

```

562         "IsHidden": false,
563         "IsEnabledByDefault": true,
564         "IsInPreview": true,
565         "SupportsEmailNotifications": false
566     },
567     "enabled": true,
568     "sendEmailsToSubscriptionOwners": true,
569     "customEmails": []
570 }
571 },
572 {
573     "type": "microsoft.insights/components/
        ProactiveDetectionConfigs",
574     "apiVersion": "2018-05-01-preview",
575     "name": "[concat(parameters('components_sustainable_ai_api_name
        '), '/extension_traceseveritydetector')]",
576     "location": "switzerlandnorth",
577     "dependsOn": [
578         "[resourceId('microsoft.insights/components', parameters('
        components_sustainable_ai_api_name'))]"
579     ],
580     "properties": {
581         "ruleDefinitions": {
582             "Name": "extension_traceseveritydetector",
583             "DisplayName": "Degradation in trace severity ratio (
                preview)",
584             "Description": "This detection rule automatically
                analyzes the trace logs emitted from your
                application, and can warn you about unusual patterns
                in the severity of your trace telemetry.",
585             "HelpUrl": "https://github.com/Microsoft/
                ApplicationInsights-Home/blob/master/SmartDetection/
                degradation-in-trace-severity-ratio.md",
586             "IsHidden": false,
587             "IsEnabledByDefault": true,
588             "IsInPreview": true,
589             "SupportsEmailNotifications": false
590         },
591         "enabled": true,
592         "sendEmailsToSubscriptionOwners": true,
593         "customEmails": []
594     }
595 },
596 {
597     "type": "microsoft.insights/components/
        ProactiveDetectionConfigs",
598     "apiVersion": "2018-05-01-preview",
599     "name": "[concat(parameters('components_sustainable_ai_api_name
        '), '/longdependencyduration')]",
600     "location": "switzerlandnorth",
601     "dependsOn": [
602         "[resourceId('microsoft.insights/components', parameters('
        components_sustainable_ai_api_name'))]"
603     ],
604     "properties": {
605         "ruleDefinitions": {
606             "Name": "longdependencyduration",

```



```

607         "DisplayName": "Long dependency duration",
608         "Description": "Smart Detection rules notify you of
609         performance anomaly issues.",
610         "HelpUrl": "https://docs.microsoft.com/en-us/azure/
611         application-insights/app-insights-proactive-
612         performance-diagnostics",
613         "IsHidden": false,
614         "IsEnabledByDefault": true,
615         "IsInPreview": false,
616         "SupportsEmailNotifications": true
617     },
618     "enabled": true,
619     "sendEmailsToSubscriptionOwners": true,
620     "customEmails": []
621 },
622 {
623     "type": "microsoft.insights/components/
624     ProactiveDetectionConfigs",
625     "apiVersion": "2018-05-01-preview",
626     "name": "[concat(parameters('components_sustainable_ai_api_name
627     '), '/migrationToAlertRulesCompleted')]",
628     "location": "switzerlandnorth",
629     "dependsOn": [
630         "[resourceId('microsoft.insights/components', parameters('
631         components_sustainable_ai_api_name'))]"
632     ],
633     "properties": {
634         "ruleDefinitions": {
635             "Name": "migrationToAlertRulesCompleted",
636             "DisplayName": "Migration To Alert Rules Completed",
637             "Description": "A configuration that controls the
638             migration state of Smart Detection to Smart Alerts",
639             "HelpUrl": "https://docs.microsoft.com/en-us/azure/
640             application-insights/app-insights-proactive-
641             performance-diagnostics",
642             "IsHidden": true,
643             "IsEnabledByDefault": false,
644             "IsInPreview": true,
645             "SupportsEmailNotifications": false
646         },
647         "enabled": false,
648         "sendEmailsToSubscriptionOwners": true,
649         "customEmails": []
650     }
651 },
652 {
653     "type": "microsoft.insights/components/
654     ProactiveDetectionConfigs",
655     "apiVersion": "2018-05-01-preview",
656     "name": "[concat(parameters('components_sustainable_ai_api_name
657     '), '/slowpageloadtime')]",
658     "location": "switzerlandnorth",
659     "dependsOn": [
660         "[resourceId('microsoft.insights/components', parameters('
661         components_sustainable_ai_api_name'))]"
662     ],

```

```

652     "properties": {
653       "ruleDefinitions": {
654         "Name": "slowpageloadtime",
655         "DisplayName": "Slow page load time",
656         "Description": "Smart Detection rules notify you of
657           performance anomaly issues.",
658         "HelpUrl": "https://docs.microsoft.com/en-us/azure/
659           application-insights/app-insights-proactive-
660           performance-diagnostics",
661         "IsHidden": false,
662         "IsEnabledByDefault": true,
663         "IsInPreview": false,
664         "SupportsEmailNotifications": true
665       },
666       "enabled": true,
667       "sendEmailsToSubscriptionOwners": true,
668       "customEmails": []
669     },
670     {
671       "type": "microsoft.insights/components/
672         ProactiveDetectionConfigs",
673       "apiVersion": "2018-05-01-preview",
674       "name": "[concat(parameters('components_sustainable_ai_api_name
675         '), '/slowserverresponsetime')]",
676       "location": "switzerlandnorth",
677       "dependsOn": [
678         "[resourceId('microsoft.insights/components', parameters('
679           components_sustainable_ai_api_name'))]"
680       ],
681       "properties": {
682         "ruleDefinitions": {
683           "Name": "slowserverresponsetime",
684           "DisplayName": "Slow server response time",
685           "Description": "Smart Detection rules notify you of
686             performance anomaly issues.",
687           "HelpUrl": "https://docs.microsoft.com/en-us/azure/
688             application-insights/app-insights-proactive-
689             performance-diagnostics",
690           "IsHidden": false,
691           "IsEnabledByDefault": true,
692           "IsInPreview": false,
693           "SupportsEmailNotifications": true
694         },
695         "enabled": true,
696         "sendEmailsToSubscriptionOwners": true,
697         "customEmails": []
698       }
699     },
700     {
701       "type": "Microsoft.ManagedIdentity/userAssignedIdentities/
702         federatedIdentityCredentials",
703       "apiVersion": "2025-01-31-preview",
704       "name": "[concat(parameters('
705         userAssignedIdentities_sustainable_ai_a_id_9c86_name'), '/
706         simonluescherfhnw-ip6-sustainable-ai-backend-a416')]",
707       "dependsOn": [

```

```

697         "[resourceId('Microsoft.ManagedIdentity/
           userAssignedIdentities', parameters('
           userAssignedIdentities_sustainable_ai_a_id_9c86_name'))]"
698     ],
699     "properties": {
700         "issuer": "https://token.actions.githubusercontent.com",
701         "subject": "repo:simonluescherfhnw/ip6-sustainable-ai-
           backend:ref:refs/heads/main",
702         "audiences": [
703             "api://AzureADTokenExchange"
704         ]
705     }
706 },
707 {
708     "type": "Microsoft.Storage/storageAccounts/blobServices",
709     "apiVersion": "2024-01-01",
710     "name": "[concat(parameters('storageAccounts_ip6sustainableai91
           a8_name'), '/default')]",
711     "dependsOn": [
712         "[resourceId('Microsoft.Storage/storageAccounts',
           parameters('storageAccounts_ip6sustainableai91a8_name'))]"
713     ],
714     "sku": {
715         "name": "Standard_LRS",
716         "tier": "Standard"
717     },
718     "properties": {
719         "cors": {
720             "corsRules": []
721         },
722         "deleteRetentionPolicy": {
723             "allowPermanentDelete": false,
724             "enabled": false
725         }
726     }
727 },
728 {
729     "type": "Microsoft.Storage/storageAccounts/fileServices",
730     "apiVersion": "2024-01-01",
731     "name": "[concat(parameters('storageAccounts_ip6sustainableai91
           a8_name'), '/default')]",
732     "dependsOn": [
733         "[resourceId('Microsoft.Storage/storageAccounts',
           parameters('storageAccounts_ip6sustainableai91a8_name'))]"
734     ],
735     "sku": {
736         "name": "Standard_LRS",
737         "tier": "Standard"
738     },
739     "properties": {
740         "protocolSettings": {
741             "smb": {}
742         },
743         "cors": {

```

```

744         "corsRules": []
745     },
746     "shareDeleteRetentionPolicy": {
747         "enabled": true,
748         "days": 7
749     }
750 }
751 },
752 {
753     "type": "Microsoft.Storage/storageAccounts/queueServices",
754     "apiVersion": "2024-01-01",
755     "name": "[concat(parameters('storageAccounts_ip6sustainableai91a8_name'), '/default')]",
756     "dependsOn": [
757         "[resourceId('Microsoft.Storage/storageAccounts',
758             parameters('storageAccounts_ip6sustainableai91a8_name'))]"
759     ],
760     "properties": {
761         "cors": {
762             "corsRules": []
763         }
764     },
765     {
766         "type": "Microsoft.Storage/storageAccounts/tableServices",
767         "apiVersion": "2024-01-01",
768         "name": "[concat(parameters('storageAccounts_ip6sustainableai91a8_name'), '/default')]",
769         "dependsOn": [
770             "[resourceId('Microsoft.Storage/storageAccounts',
771                 parameters('storageAccounts_ip6sustainableai91a8_name'))]"
772         ],
773         "properties": {
774             "cors": {
775                 "corsRules": []
776             }
777         },
778         {
779             "type": "Microsoft.Web/sites",
780             "apiVersion": "2024-04-01",
781             "name": "[parameters('sites_sustainable_ai_api_name')]",
782             "location": "Switzerland North",
783             "dependsOn": [
784                 "[resourceId('Microsoft.Web/serverfarms', parameters('serverfarms_ASP_ip6sustainableai_adb9_name'))]"
785             ],
786             "kind": "functionapp",
787             "properties": {
788                 "enabled": true,
789                 "hostNameSslStates": [
790                     {
791                         "name": "[concat(parameters('sites_sustainable_ai_api_name'), '.azurewebsites.net')]",

```

```

792         "sslState": "Disabled",
793         "hostType": "Standard"
794     },
795     {
796         "name": "[concat(parameters('
797             sites_sustainable_ai_api_name'), '.scm.
798             azurewebsites.net')]",
799         "sslState": "Disabled",
800         "hostType": "Repository"
801     }
802 ],
803 "serverFarmId": "[resourceId('Microsoft.Web/serverfarms',
804     parameters('serverfarms_ASP_ip6sustainableai_adb9_name')
805 )]",
806 "reserved": false,
807 "isXenon": false,
808 "hyperV": false,
809 "dnsConfiguration": {},
810 "vnetRouteAllEnabled": false,
811 "vnetImagePullEnabled": false,
812 "vnetContentShareEnabled": false,
813 "siteConfig": {
814     "numberOfWorkers": 1,
815     "acrUseManagedIdentityCreds": false,
816     "alwaysOn": false,
817     "http20Enabled": false,
818     "functionAppScaleLimit": 200,
819     "minimumElasticInstanceCount": 0
820 },
821 "scmSiteAlsoStopped": false,
822 "clientAffinityEnabled": false,
823 "clientCertEnabled": false,
824 "clientCertMode": "Required",
825 "hostNamesDisabled": false,
826 "ipMode": "IPv4",
827 "vnetBackupRestoreEnabled": false,
828 "customDomainVerificationId": "12AB493E38ABD001261CF21640AB
829     413738832F9779A1C425D06A0380761D8BA9",
830 "containerSize": 1536,
831 "dailyMemoryTimeQuota": 0,
832 "httpsOnly": true,
833 "endToEndEncryptionEnabled": false,
834 "redundancyMode": "None",
835 "publicNetworkAccess": "Enabled",
836 "storageAccountRequired": false,
837 "keyVaultReferenceIdentity": "SystemAssigned"
838 }
839 },
840 {
841     "type": "Microsoft.Web/sites/basicPublishingCredentialsPolicies",
842     "apiVersion": "2024-04-01",
843     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
844         '/ftp')]",
845     "location": "Switzerland North",
846     "dependsOn": [

```

```

841         "[resourceId('Microsoft.Web/sites', parameters('
842             sites_sustainable_ai_api_name'))]"
843     ],
844     "properties": {
845         "allow": false
846     }
847 },
848 {
849     "type": "Microsoft.Web/sites/basicPublishingCredentialsPolicies",
850     "apiVersion": "2024-04-01",
851     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
852         '/scm')]",
853     "location": "Switzerland North",
854     "dependsOn": [
855         "[resourceId('Microsoft.Web/sites', parameters('
856             sites_sustainable_ai_api_name'))]"
857     ],
858     "properties": {
859         "allow": false
860     }
861 },
862 {
863     "type": "Microsoft.Web/sites/config",
864     "apiVersion": "2024-04-01",
865     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
866         '/web')]",
867     "location": "Switzerland North",
868     "dependsOn": [
869         "[resourceId('Microsoft.Web/sites', parameters('
870             sites_sustainable_ai_api_name'))]"
871     ],
872     "properties": {
873         "numberOfWorkers": 1,
874         "defaultDocuments": [
875             "Default.htm",
876             "Default.html",
877             "Default.asp",
878             "index.htm",
879             "index.html",
880             "iisstart.htm",
881             "default.aspx",
882             "index.php"
883         ],
884         "netFrameworkVersion": "v8.0",
885         "requestTracingEnabled": false,
886         "remoteDebuggingEnabled": false,
887         "httpLoggingEnabled": false,
888         "acrUseManagedIdentityCreds": false,
889         "logsDirectorySizeLimit": 35,
890         "detailedErrorLoggingEnabled": false,
891         "publishingUsername": "REDACTED",
892         "scmType": "GitHubAction",
893         "use32BitWorkerProcess": false,
894         "webSocketsEnabled": false,
895         "alwaysOn": false,
896         "managedPipelineMode": "Integrated",

```

```

892     "virtualApplications": [
893         {
894             "virtualPath": "/",
895             "physicalPath": "site\\wwwroot",
896             "preloadEnabled": false
897         }
898     ],
899     "loadBalancing": "LeastRequests",
900     "experiments": {
901         "rampUpRules": []
902     },
903     "autoHealEnabled": false,
904     "vnetRouteAllEnabled": false,
905     "vnetPrivatePortsCount": 0,
906     "publicNetworkAccess": "Enabled",
907     "cors": {
908         "allowedOrigins": [
909             "https://portal.azure.com",
910             "https://green-mud-04afae203.6.azurestaticapps.net"
911         ],
912         "supportCredentials": false
913     },
914     "localMySqlEnabled": false,
915     "ipSecurityRestrictions": [
916         {
917             "ipAddress": "Any",
918             "action": "Allow",
919             "priority": 2147483647,
920             "name": "Allow all",
921             "description": "Allow all access"
922         }
923     ],
924     "scmIpSecurityRestrictions": [
925         {
926             "ipAddress": "Any",
927             "action": "Allow",
928             "priority": 2147483647,
929             "name": "Allow all",
930             "description": "Allow all access"
931         }
932     ],
933     "scmIpSecurityRestrictionsUseMain": false,
934     "http20Enabled": false,
935     "minTlsVersion": "1.2",
936     "scmMinTlsVersion": "1.2",
937     "ftpsState": "FtpsOnly",
938     "preWarmedInstanceCount": 0,
939     "functionAppScaleLimit": 200,
940     "functionsRuntimeScaleMonitoringEnabled": false,
941     "minimumElasticInstanceCount": 0,
942     "azureStorageAccounts": {}
943 },
944 {
945     "type": "Microsoft.Web/sites/deployments",
946     "apiVersion": "2024-04-01",

```

```

948     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
949       '/076c043f429d4b8ab61a0a615df58b07')]",
950     "location": "Switzerland North",
951     "dependsOn": [
952       "[resourceId('Microsoft.Web/sites', parameters('
953         sites_sustainable_ai_api_name'))]"
954     ],
955     "properties": {
956       "status": 4,
957       "author_email": "N/A",
958       "author": "N/A",
959       "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
960       "message": "{ \"type\": \"deployment\", \"sha\": \"419f7ebeckb79
961         48f46bf9078a0528fc7514af3f81\", \"repoName\": \"
962         simonluescherfhnw/ip6-sustainable-ai-backend\", \"actor\":
963         \"simonluescherfhnw\", \"slotName\": \"production\" }",
964       "start_time": "2025-06-23T15:01:11.4681421Z",
965       "end_time": "2025-06-23T15:01:12.8275307Z",
966       "active": false
967     }
968   },
969   {
970     "type": "Microsoft.Web/sites/deployments",
971     "apiVersion": "2024-04-01",
972     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
973       '/2338922e33cf4e73bec6c487c2719e5b')]",
974     "location": "Switzerland North",
975     "dependsOn": [
976       "[resourceId('Microsoft.Web/sites', parameters('
977         sites_sustainable_ai_api_name'))]"
978     ],
979     "properties": {
980       "status": 4,
981       "author_email": "N/A",
982       "author": "N/A",
983       "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
984       "message": "{ \"type\": \"deployment\", \"sha\": \"4ca2c2e9bc10
985         1412c924390186decd0ee39e4a08\", \"repoName\": \"
986         simonluescherfhnw/ip6-sustainable-ai-backend\", \"actor\":
987         \"simonluescherfhnw\", \"slotName\": \"production\" }",
988       "start_time": "2025-06-22T20:06:26.235006Z",
989       "end_time": "2025-06-22T20:06:27.7987243Z",
990       "active": false
991     }
992   }
993 },
994 {
995   "type": "Microsoft.Web/sites/deployments",
996   "apiVersion": "2024-04-01",
997   "name": "[concat(parameters('sites_sustainable_ai_api_name'),
998     '/23bb4f51bfa244c9b15bb32eee71b5f1')]",
999   "location": "Switzerland North",
1000   "dependsOn": [
1001     "[resourceId('Microsoft.Web/sites', parameters('
1002       sites_sustainable_ai_api_name'))]"
1003   ],
1004   "properties": {
1005     "status": 4,

```



```

993     "author_email": "N/A",
994     "author": "N/A",
995     "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
996     "message": "{\\"type\\":\\"deployment\\",\\"sha\\":\\"1b590ce92de3
11c631be7d64d636f067b64a405f\\",\\"repoName\\":\\"
simonluescherfhnw/ip6-sustainable-ai-backend\\",\\"actor\\"
:\\simonluescherfhnw\\",\\"slotName\\":\\"production\\"}",
997     "start_time": "2025-06-22T21:22:34.294617Z",
998     "end_time": "2025-06-22T21:22:35.5289937Z",
999     "active": false
1000   }
1001 },
1002 {
1003   "type": "Microsoft.Web/sites/deployments",
1004   "apiVersion": "2024-04-01",
1005   "name": "[concat(parameters('sites_sustainable_ai_api_name'),
'/435e372ef56847438fe34fe62aafeb84')]",
1006   "location": "Switzerland North",
1007   "dependsOn": [
1008     "[resourceId('Microsoft.Web/sites', parameters('
sites_sustainable_ai_api_name'))]"
1009 ],
1010   "properties": {
1011     "status": 4,
1012     "author_email": "N/A",
1013     "author": "N/A",
1014     "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
1015     "message": "{\\"type\\":\\"deployment\\",\\"sha\\":\\"e6c8c842e11a
4e1ffdf5f04518dc099f8489d676\\",\\"repoName\\":\\"
simonluescherfhnw/ip6-sustainable-ai-backend\\",\\"actor\\"
:\\simonluescherfhnw\\",\\"slotName\\":\\"production\\"}",
1016     "start_time": "2025-06-22T18:13:00.8034236Z",
1017     "end_time": "2025-06-22T18:13:02.335107Z",
1018     "active": false
1019   }
1020 },
1021 {
1022   "type": "Microsoft.Web/sites/deployments",
1023   "apiVersion": "2024-04-01",
1024   "name": "[concat(parameters('sites_sustainable_ai_api_name'),
'/60c0227b3c3b45b7bcdffb5c6b0c2e3e')]",
1025   "location": "Switzerland North",
1026   "dependsOn": [
1027     "[resourceId('Microsoft.Web/sites', parameters('
sites_sustainable_ai_api_name'))]"
1028 ],
1029   "properties": {
1030     "status": 4,
1031     "author_email": "N/A",
1032     "author": "N/A",
1033     "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
1034     "message": "{\\"type\\":\\"deployment\\",\\"sha\\":\\"2532b70ad018
ca5fbf8c8da5e015a965cd3ed196\\",\\"repoName\\":\\"
simonluescherfhnw/ip6-sustainable-ai-backend\\",\\"actor\\"
:\\simonluescherfhnw\\",\\"slotName\\":\\"production\\"}",
1035     "start_time": "2025-06-22T15:06:30.2846682Z",
1036     "end_time": "2025-06-22T15:06:31.7455107Z",

```

```

1037         "active": false
1038     }
1039 },
1040 {
1041     "type": "Microsoft.Web/sites/deployments",
1042     "apiVersion": "2024-04-01",
1043     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1044         '/731c60269c7d4c478be50128397d5f3a')]",
1045     "location": "Switzerland North",
1046     "dependsOn": [
1047         "[resourceId('Microsoft.Web/sites', parameters('
1048             sites_sustainable_ai_api_name'))]"
1049     ],
1050     "properties": {
1051         "status": 4,
1052         "author_email": "N/A",
1053         "author": "N/A",
1054         "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
1055         "message": "{ \"type\": \"deployment\", \"sha\": \"c15f9163235f
1056             02835ad786a184571713476ee7ad\", \"repoName\": \"
1057             simonluescherfhnw/ip6-sustainable-ai-backend\", \"actor\":
1058             \"simonluescherfhnw\", \"slotName\": \"production\" }",
1059         "start_time": "2025-06-22T13:05:18.8184065Z",
1060         "end_time": "2025-06-22T13:05:20.4590363Z",
1061         "active": false
1062     }
1063 },
1064 {
1065     "type": "Microsoft.Web/sites/deployments",
1066     "apiVersion": "2024-04-01",
1067     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1068         '/77373ccd6d6b41468229de0efcc5f418')]",
1069     "location": "Switzerland North",
1070     "dependsOn": [
1071         "[resourceId('Microsoft.Web/sites', parameters('
1072             sites_sustainable_ai_api_name'))]"
1073     ],
1074     "properties": {
1075         "status": 4,
1076         "author_email": "N/A",
1077         "author": "N/A",
1078         "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
1079         "message": "{ \"type\": \"deployment\", \"sha\": \"9fee899ff89c
1080             096f832a0fe66c993dd6df09b391\", \"repoName\": \"
1081             simonluescherfhnw/ip6-sustainable-ai-backend\", \"actor\":
1082             \"simonluescherfhnw\", \"slotName\": \"production\" }",
1083         "start_time": "2025-06-22T15:36:23.1319342Z",
1084         "end_time": "2025-06-22T15:36:24.7725459Z",
1085         "active": false
1086     }
1087 },
1088 {
1089     "type": "Microsoft.Web/sites/deployments",
1090     "apiVersion": "2024-04-01",
1091     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1092         '/c46f976169fd44ec9fcd4498c602cffc')]",
1093     "location": "Switzerland North",

```

```

1083     "dependsOn": [
1084         "[resourceId('Microsoft.Web/sites', parameters('
            sites_sustainable_ai_api_name'))]"
1085     ],
1086     "properties": {
1087         "status": 4,
1088         "author_email": "N/A",
1089         "author": "N/A",
1090         "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
1091         "message": "{\\"type\\":\\"deployment\\",\\"sha\\":\\"3e471ef8240
            eb2b06363ac53bd75717beb281f59\\",\\"repoName\\":\\"
            simonluescherfhnw/ip6-sustainable-ai-backend\\",\\"actor\\"
            :\\"simonluescherfhnw\\",\\"slotName\\":\\"production\\"}",
1092         "start_time": "2025-06-22T20:53:56.0463655Z",
1093         "end_time": "2025-06-22T20:53:57.5153299Z",
1094         "active": false
1095     }
1096 },
1097 {
1098     "type": "Microsoft.Web/sites/deployments",
1099     "apiVersion": "2024-04-01",
1100     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
            '/d7386b0db21443178515e424abbffda3')]",
1101     "location": "Switzerland North",
1102     "dependsOn": [
1103         "[resourceId('Microsoft.Web/sites', parameters('
            sites_sustainable_ai_api_name'))]"
1104     ],
1105     "properties": {
1106         "status": 4,
1107         "author_email": "N/A",
1108         "author": "N/A",
1109         "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",
1110         "message": "{\\"type\\":\\"deployment\\",\\"sha\\":\\"3766cbf208f6
            95f7ade3216e16fc2363650bc154\\",\\"repoName\\":\\"
            simonluescherfhnw/ip6-sustainable-ai-backend\\",\\"actor\\"
            :\\"simonluescherfhnw\\",\\"slotName\\":\\"production\\"}",
1111         "start_time": "2025-06-22T18:22:19.429819Z",
1112         "end_time": "2025-06-22T18:22:21.682451Z",
1113         "active": false
1114     }
1115 },
1116 {
1117     "type": "Microsoft.Web/sites/deployments",
1118     "apiVersion": "2024-04-01",
1119     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
            '/f978abb7e53b40b9bfba2cb195bb857e')]",
1120     "location": "Switzerland North",
1121     "dependsOn": [
1122         "[resourceId('Microsoft.Web/sites', parameters('
            sites_sustainable_ai_api_name'))]"
1123     ],
1124     "properties": {
1125         "status": 4,
1126         "author_email": "N/A",
1127         "author": "N/A",
1128         "deployer": "GITHUB_ZIP_DEPLOY_FUNCTIONS_V1",

```

```

1129         "message": "{\\"type\\":\\"deployment\\",\\"sha\\":\\"9a8019effc5d
1130             17f37d528e0f2d3f75c6fd7d6942\\",\\"repoName\\":\\"
1131             simonluescherfhnw/ip6-sustainable-ai-backend\\",\\"actor\\"
1132             :\\"simonluescherfhnw\\",\\"slotName\\":\\"production\\"}",
1133     },
1134 },
1135 {
1136     "type": "Microsoft.Web/sites/functions",
1137     "apiVersion": "2024-04-01",
1138     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1139         '/DeleteConversation')]",
1140     "location": "Switzerland North",
1141     "dependsOn": [
1142         "[resourceId('Microsoft.Web/sites', parameters('
1143             sites_sustainable_ai_api_name'))]"
1144     ],
1145     "properties": {
1146         "script_href": "https://sustainable-ai-api.azurewebsites.
1147             net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1148         "test_data_href": "https://sustainable-ai-api.azurewebsites
1149             .net/admin/vfs/data/Functions/sampleddata/
1150             DeleteConversation.dat",
1151         "href": "https://sustainable-ai-api.azurewebsites.net/admin
1152             /functions/DeleteConversation",
1153         "config": {
1154             "name": "DeleteConversation",
1155             "entryPoint": "DeleteConversation.Run",
1156             "scriptFile": "SustainableAI.Api.dll",
1157             "language": "dotnet-isolated",
1158             "functionDirectory": "",
1159             "bindings": [
1160                 {
1161                     "name": "req",
1162                     "type": "httpTrigger",
1163                     "direction": "In",
1164                     "authLevel": "Anonymous",
1165                     "methods": [
1166                         "post"
1167                     ]
1168                 },
1169                 {
1170                     "name": "$return",
1171                     "type": "http",
1172                     "direction": "Out"
1173                 }
1174             ]
1175         },
1176         "invoke_url_template": "https://sustainable-ai-api.
1177             azurewebsites.net/api/deleteconversation",
1178         "language": "dotnet-isolated",
1179         "isDisabled": false
1180     }
1181 }

```

```

1176     "type": "Microsoft.Web/sites/functions",
1177     "apiVersion": "2024-04-01",
1178     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1179         '/GetAppData')]",
1179     "location": "Switzerland North",
1180     "dependsOn": [
1181         "[resourceId('Microsoft.Web/sites', parameters('
1182             sites_sustainable_ai_api_name'))]"
1182     ],
1183     "properties": {
1184         "script_href": "https://sustainable-ai-api.azurewebsites.
1185             net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1186         "test_data_href": "https://sustainable-ai-api.azurewebsites
1187             .net/admin/vfs/data/Functions/sampledData/GetAppData.dat",
1188         "href": "https://sustainable-ai-api.azurewebsites.net/admin
1189             /functions/GetAppData",
1190         "config": {
1191             "name": "GetAppData",
1192             "entryPoint": "GetAppData.Run",
1193             "scriptFile": "SustainableAI.Api.dll",
1194             "language": "dotnet-isolated",
1195             "functionDirectory": "",
1196             "bindings": [
1197                 {
1198                     "name": "req",
1199                     "type": "httpTrigger",
1200                     "direction": "In",
1201                     "authLevel": "Anonymous",
1202                     "methods": [
1203                         "get"
1204                     ]
1205                 },
1206                 {
1207                     "name": "$return",
1208                     "type": "http",
1209                     "direction": "Out"
1210                 }
1211             ]
1212         },
1213         "invoke_url_template": "https://sustainable-ai-api.
1214             azurewebsites.net/api/getappdata",
1215         "language": "dotnet-isolated",
1216         "isDisabled": false
1217     }
1218 },
1219 {
1220     "type": "Microsoft.Web/sites/functions",
1221     "apiVersion": "2024-04-01",
1222     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1223         '/GetConversations')]",
1224     "location": "Switzerland North",
1225     "dependsOn": [
1226         "[resourceId('Microsoft.Web/sites', parameters('
1227             sites_sustainable_ai_api_name'))]"
1228     ],
1229     "properties": {

```

```

1224     "script_href": "https://sustainable-ai-api.azurewebsites.
1225         net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1226     "test_data_href": "https://sustainable-ai-api.azurewebsites
1227         .net/admin/vfs/data/Functions/sampledData/
1228         GetConversations.dat",
1229     "href": "https://sustainable-ai-api.azurewebsites.net/admin
1230         /functions/GetConversations",
1231     "config": {
1232         "name": "GetConversations",
1233         "entryPoint": "GetConversations.Run",
1234         "scriptFile": "SustainableAI.Api.dll",
1235         "language": "dotnet-isolated",
1236         "functionDirectory": "",
1237         "bindings": [
1238             {
1239                 "name": "req",
1240                 "type": "httpTrigger",
1241                 "direction": "In",
1242                 "authLevel": "Anonymous",
1243                 "methods": [
1244                     "post"
1245                 ]
1246             },
1247             {
1248                 "name": "$return",
1249                 "type": "http",
1250                 "direction": "Out"
1251             }
1252         ]
1253     },
1254     "invoke_url_template": "https://sustainable-ai-api.
1255         azurewebsites.net/api/getconversations",
1256     "language": "dotnet-isolated",
1257     "isDisabled": false
1258 },
1259 {
1260     "type": "Microsoft.Web/sites/functions",
1261     "apiVersion": "2024-04-01",
1262     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1263         '/GetPrompts')]",
1264     "location": "Switzerland North",
1265     "dependsOn": [
1266         "[resourceId('Microsoft.Web/sites', parameters('
1267             sites_sustainable_ai_api_name'))]"
1268     ],
1269     "properties": {
1270         "script_href": "https://sustainable-ai-api.azurewebsites.
1271             net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1272         "test_data_href": "https://sustainable-ai-api.azurewebsites
1273             .net/admin/vfs/data/Functions/sampledData/GetPrompts.dat"
1274         ,
1275         "href": "https://sustainable-ai-api.azurewebsites.net/admin
1276             /functions/GetPrompts",
1277         "config": {
1278             "name": "GetPrompts",
1279             "entryPoint": "GetPrompts.Run",

```

```

1270         "scriptFile": "SustainableAI.Api.dll",
1271         "language": "dotnet-isolated",
1272         "functionDirectory": "",
1273         "bindings": [
1274             {
1275                 "name": "req",
1276                 "type": "httpTrigger",
1277                 "direction": "In",
1278                 "authLevel": "Anonymous",
1279                 "methods": [
1280                     "post"
1281                 ]
1282             },
1283             {
1284                 "name": "$return",
1285                 "type": "http",
1286                 "direction": "Out"
1287             }
1288         ]
1289     },
1290     "invoke_url_template": "https://sustainable-ai-api.
        azurewebsites.net/api/getprompts",
1291     "language": "dotnet-isolated",
1292     "isDisabled": false
1293 }
1294 },
1295 {
1296     "type": "Microsoft.Web/sites/functions",
1297     "apiVersion": "2024-04-01",
1298     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
        '/GetUsageStatistics')]",
1299     "location": "Switzerland North",
1300     "dependsOn": [
1301         "[resourceId('Microsoft.Web/sites', parameters('
        sites_sustainable_ai_api_name'))]"
1302     ],
1303     "properties": {
1304         "script_href": "https://sustainable-ai-api.azurewebsites.
        net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1305         "test_data_href": "https://sustainable-ai-api.azurewebsites
        .net/admin/vfs/data/Functions/sampledData/
        GetUsageStatistics.dat",
1306         "href": "https://sustainable-ai-api.azurewebsites.net/admin
        /functions/GetUsageStatistics",
1307         "config": {
1308             "name": "GetUsageStatistics",
1309             "entryPoint": "GetUsageStatistics.Run",
1310             "scriptFile": "SustainableAI.Api.dll",
1311             "language": "dotnet-isolated",
1312             "functionDirectory": "",
1313             "bindings": [
1314                 {
1315                     "name": "req",
1316                     "type": "httpTrigger",
1317                     "direction": "In",
1318                     "authLevel": "Anonymous",
1319                     "methods": [

```





```

1370         "invoke_url_template": "https://sustainable-ai-api.
1371             azurewebsites.net/api/logpagevisit",
1372         "language": "dotnet-isolated",
1373         "isDisabled": false
1374     },
1375     {
1376         "type": "Microsoft.Web/sites/functions",
1377         "apiVersion": "2024-04-01",
1378         "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1379             '/PredictPromptUsage')]",
1380         "location": "Switzerland North",
1381         "dependsOn": [
1382             "[resourceId('Microsoft.Web/sites', parameters('
1383                 sites_sustainable_ai_api_name'))]"
1384         ],
1385         "properties": {
1386             "script_href": "https://sustainable-ai-api.azurewebsites.
1387                 net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1388             "test_data_href": "https://sustainable-ai-api.azurewebsites
1389                 .net/admin/vfs/data/Functions/sampleddata/
1390                 PredictPromptUsage.dat",
1391             "href": "https://sustainable-ai-api.azurewebsites.net/admin
1392                 /functions/PredictPromptUsage",
1393             "config": {
1394                 "name": "PredictPromptUsage",
1395                 "entryPoint": "PredictPromptUsage.Run",
1396                 "scriptFile": "SustainableAI.Api.dll",
1397                 "language": "dotnet-isolated",
1398                 "functionDirectory": "",
1399                 "bindings": [
1400                     {
1401                         "name": "req",
1402                         "type": "httpTrigger",
1403                         "direction": "In",
1404                         "authLevel": "Anonymous",
1405                         "methods": [
1406                             "post"
1407                         ]
1408                     },
1409                     {
1410                         "name": "$return",
1411                         "type": "http",
1412                         "direction": "Out"
1413                     }
1414                 ]
1415             },
1416             "invoke_url_template": "https://sustainable-ai-api.
1417                 azurewebsites.net/api/predictpromptusage",
1418             "language": "dotnet-isolated",
1419             "isDisabled": false
1420         }
1421     }
1422 ]

```

```

1418     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1419         '/SendPrompt')]",
1420     "location": "Switzerland North",
1421     "dependsOn": [
1422         "[resourceId('Microsoft.Web/sites', parameters('
1423             sites_sustainable_ai_api_name'))]"
1424     ],
1425     "properties": {
1426         "script_href": "https://sustainable-ai-api.azurewebsites.
1427             net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1428         "test_data_href": "https://sustainable-ai-api.azurewebsites
1429             .net/admin/vfs/data/Functions/sampledData/SendPrompt.dat"
1430         ,
1431         "href": "https://sustainable-ai-api.azurewebsites.net/admin
1432             /functions/SendPrompt",
1433         "config": {
1434             "name": "SendPrompt",
1435             "entryPoint": "SendPrompt.Run",
1436             "scriptFile": "SustainableAI.Api.dll",
1437             "language": "dotnet-isolated",
1438             "functionDirectory": "",
1439             "bindings": [
1440                 {
1441                     "name": "req",
1442                     "type": "httpTrigger",
1443                     "direction": "In",
1444                     "authLevel": "Anonymous",
1445                     "methods": [
1446                         "post"
1447                     ]
1448                 },
1449                 {
1450                     "name": "$return",
1451                     "type": "http",
1452                     "direction": "Out"
1453                 }
1454             ]
1455         },
1456         "invoke_url_template": "https://sustainable-ai-api.
1457             azurewebsites.net/api/sendprompt",
1458         "language": "dotnet-isolated",
1459         "isDisabled": false
1460     }
1461 },
1462 {
1463     "type": "Microsoft.Web/sites/functions",
1464     "apiVersion": "2024-04-01",
1465     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1466         '/UpdateConversation')]",
1467     "location": "Switzerland North",
1468     "dependsOn": [
1469         "[resourceId('Microsoft.Web/sites', parameters('
1470             sites_sustainable_ai_api_name'))]"
1471     ],
1472     "properties": {
1473         "script_href": "https://sustainable-ai-api.azurewebsites.
1474             net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",

```

```

1465     "test_data_href": "https://sustainable-ai-api.azurewebsites
      .net/admin/vfs/data/Functions/sampledData/
      UpdateConversation.dat",
1466     "href": "https://sustainable-ai-api.azurewebsites.net/admin
      /functions/UpdateConversation",
1467     "config": {
1468         "name": "UpdateConversation",
1469         "entryPoint": "UpdateConversation.Run",
1470         "scriptFile": "SustainableAI.Api.dll",
1471         "language": "dotnet-isolated",
1472         "functionDirectory": "",
1473         "bindings": [
1474             {
1475                 "name": "req",
1476                 "type": "httpTrigger",
1477                 "direction": "In",
1478                 "authLevel": "Anonymous",
1479                 "methods": [
1480                     "post"
1481                 ]
1482             },
1483             {
1484                 "name": "$return",
1485                 "type": "http",
1486                 "direction": "Out"
1487             }
1488         ]
1489     },
1490     "invoke_url_template": "https://sustainable-ai-api.
      azurewebsites.net/api/updateconversation",
1491     "language": "dotnet-isolated",
1492     "isDisabled": false
1493 }
1494 },
1495 {
1496     "type": "Microsoft.Web/sites/functions",
1497     "apiVersion": "2024-04-01",
1498     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
      '/UpdateUser')]",
1499     "location": "Switzerland North",
1500     "dependsOn": [
1501         "[resourceId('Microsoft.Web/sites', parameters('
      sites_sustainable_ai_api_name'))]"
1502     ],
1503     "properties": {
1504         "script_href": "https://sustainable-ai-api.azurewebsites.
      net/admin/vfs/site/wwwroot/SustainableAI.Api.dll",
1505         "test_data_href": "https://sustainable-ai-api.azurewebsites
      .net/admin/vfs/data/Functions/sampledData/UpdateUser.dat"
1506     },
1507     "href": "https://sustainable-ai-api.azurewebsites.net/admin
      /functions/UpdateUser",
1508     "config": {
1509         "name": "UpdateUser",
1510         "entryPoint": "UpdateUser.Run",
1511         "scriptFile": "SustainableAI.Api.dll",
1512         "language": "dotnet-isolated",

```

```

1512         "functionDirectory": "",
1513         "bindings": [
1514             {
1515                 "name": "req",
1516                 "type": "httpTrigger",
1517                 "direction": "In",
1518                 "authLevel": "Anonymous",
1519                 "methods": [
1520                     "post"
1521                 ]
1522             },
1523             {
1524                 "name": "$return",
1525                 "type": "http",
1526                 "direction": "Out"
1527             }
1528         ],
1529         "invoke_url_template": "https://sustainable-ai-api.
1530             azurewebsites.net/api/updateuser",
1531         "language": "dotnet-isolated",
1532         "isDisabled": false
1533     },
1534 },
1535 {
1536     "type": "Microsoft.Web/sites/hostnameBindings",
1537     "apiVersion": "2024-04-01",
1538     "name": "[concat(parameters('sites_sustainable_ai_api_name'),
1539         '/', parameters('sites_sustainable_ai_api_name'), '.
1540         azurewebsites.net')]",
1541     "location": "Switzerland North",
1542     "dependsOn": [
1543         "[resourceId('Microsoft.Web/sites', parameters('
1544             sites_sustainable_ai_api_name'))]"
1545     ],
1546     "properties": {
1547         "siteName": "sustainable-ai-api",
1548         "hostNameType": "Verified"
1549     }
1550 },
1551 {
1552     "type": "Microsoft.Web/staticSites/basicAuth",
1553     "apiVersion": "2024-04-01",
1554     "name": "[concat(parameters('
1555         staticSites_sustainable_ai_web_app_name'), '/default')]",
1556     "location": "West Europe",
1557     "dependsOn": [
1558         "[resourceId('Microsoft.Web/staticSites', parameters('
1559             staticSites_sustainable_ai_web_app_name'))]"
1560     ],
1561     "properties": {
1562         "applicableEnvironmentsMode": "SpecifiedEnvironments"
1563     }
1564 },
1565 {
1566     "type": "Microsoft.Web/staticSites/customDomains",
1567     "apiVersion": "2024-04-01",

```

```

1563         "name": "[concat(parameters('
1564             staticSites_sustainable_ai_web_app_name'), '/thebotter.com')
1565         ]",
1566         "location": "West Europe",
1567         "dependsOn": [
1568             "[resourceId('Microsoft.Web/staticSites', parameters('
1569                 staticSites_sustainable_ai_web_app_name'))]"
1570         ],
1571         "properties": {}
1572     },
1573     {
1574         "type": "Microsoft.Web/staticSites/customDomains",
1575         "apiVersion": "2024-04-01",
1576         "name": "[concat(parameters('
1577             staticSites_sustainable_ai_web_app_name'), '/www.thebotter.
1578             com'))]",
1579         "location": "West Europe",
1580         "dependsOn": [
1581             "[resourceId('Microsoft.Web/staticSites', parameters('
1582                 staticSites_sustainable_ai_web_app_name'))]"
1583         ],
1584         "properties": {}
1585     },
1586     {
1587         "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/
1588             containers",
1589         "apiVersion": "2024-12-01-preview",
1590         "name": "[concat(parameters('
1591             databaseAccounts_sustainable_ai_cosmos_dba_name'), '/
1592             Development/Conversation')]",
1593         "dependsOn": [
1594             "[resourceId('Microsoft.DocumentDB/databaseAccounts/
1595                 sqlDatabases', parameters('
1596                     databaseAccounts_sustainable_ai_cosmos_dba_name'), '
1597                     Development'))]",
1598             "[resourceId('Microsoft.DocumentDB/databaseAccounts',
1599                 parameters('
1600                     databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1601         ],
1602         "properties": {
1603             "resource": {
1604                 "id": "Conversation",
1605                 "indexingPolicy": {
1606                     "indexingMode": "consistent",
1607                     "automatic": true,
1608                     "includedPaths": [
1609                         {
1610                             "path": "/*"
1611                         }
1612                     ],
1613                     "excludedPaths": [
1614                         {
1615                             "path": "/\"_etag\"/?"
1616                         }
1617                     ]
1618                 }
1619             },
1620             "partitionKey": {

```

```

1606         "paths": [
1607             "/userId"
1608         ],
1609         "kind": "Hash",
1610         "version": 2
1611     },
1612     "uniqueKeyPolicy": {
1613         "uniqueKeys": []
1614     },
1615     "conflictResolutionPolicy": {
1616         "mode": "LastWriterWins",
1617         "conflictResolutionPath": "/_ts"
1618     },
1619     "computedProperties": []
1620 }
1621 }
1622 },
1623 {
1624     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/containers",
1625     "apiVersion": "2024-12-01-preview",
1626     "name": "[concat(parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Production/Conversation')]",
1627     "dependsOn": [
1628         "[resourceId('Microsoft.DocumentDB/databaseAccounts/sqlDatabases', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Production')]",
1629         "[resourceId('Microsoft.DocumentDB/databaseAccounts', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1630     ],
1631     "properties": {
1632         "resource": {
1633             "id": "Conversation",
1634             "indexingPolicy": {
1635                 "indexingMode": "consistent",
1636                 "automatic": true,
1637                 "includedPaths": [
1638                     {
1639                         "path": "/*"
1640                     }
1641                 ],
1642                 "excludedPaths": [
1643                     {
1644                         "path": "/\"_etag\"/?"
1645                     }
1646                 ]
1647             },
1648             "partitionKey": {
1649                 "paths": [
1650                     "/userId"
1651                 ],
1652                 "kind": "Hash",
1653                 "version": 2
1654             },

```

```

1655         "uniqueKeyPolicy": {
1656             "uniqueKeys": []
1657         },
1658         "conflictResolutionPolicy": {
1659             "mode": "LastWriterWins",
1660             "conflictResolutionPath": "/_ts"
1661         },
1662         "computedProperties": []
1663     }
1664 }
1665 },
1666 {
1667     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/
        containers",
1668     "apiVersion": "2024-12-01-preview",
1669     "name": "[concat(parameters('
        databaseAccounts_sustainable_ai_cosmos_dba_name'), ' /
        Development/EnergyUnit')]",
1670     "dependsOn": [
1671         "[resourceId('Microsoft.DocumentDB/databaseAccounts/
        sqlDatabases', parameters('
        databaseAccounts_sustainable_ai_cosmos_dba_name'), '
        Development')]",
1672         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
        parameters('
        databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1673     ],
1674     "properties": {
1675         "resource": {
1676             "id": "EnergyUnit",
1677             "indexingPolicy": {
1678                 "indexingMode": "consistent",
1679                 "automatic": true,
1680                 "includedPaths": [
1681                     {
1682                         "path": "/*"
1683                     }
1684                 ],
1685                 "excludedPaths": [
1686                     {
1687                         "path": "/\"_etag\"/?"
1688                     }
1689                 ]
1690             },
1691             "partitionKey": {
1692                 "paths": [
1693                     "/tenantId"
1694                 ],
1695                 "kind": "Hash",
1696                 "version": 2
1697             },
1698             "uniqueKeyPolicy": {
1699                 "uniqueKeys": []
1700             },
1701             "conflictResolutionPolicy": {
1702                 "mode": "LastWriterWins",
1703                 "conflictResolutionPath": "/_ts"

```

```

1704         },
1705         "computedProperties": []
1706     }
1707 }
1708 },
1709 {
1710     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/containers",
1711     "apiVersion": "2024-12-01-preview",
1712     "name": "[concat(parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), '/Production/EnergyUnit')]",
1713     "dependsOn": [
1714         "[resourceId('Microsoft.DocumentDB/databaseAccounts/sqlDatabases', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Production')]",
1715         "[resourceId('Microsoft.DocumentDB/databaseAccounts', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1716     ],
1717     "properties": {
1718         "resource": {
1719             "id": "EnergyUnit",
1720             "indexingPolicy": {
1721                 "indexingMode": "consistent",
1722                 "automatic": true,
1723                 "includedPaths": [
1724                     {
1725                         "path": "/*"
1726                     }
1727                 ],
1728                 "excludedPaths": [
1729                     {
1730                         "path": "/\"_etag\"/?"
1731                     }
1732                 ]
1733             },
1734             "partitionKey": {
1735                 "paths": [
1736                     "/tenantId"
1737                 ],
1738                 "kind": "Hash",
1739                 "version": 2
1740             },
1741             "uniqueKeyPolicy": {
1742                 "uniqueKeys": []
1743             },
1744             "conflictResolutionPolicy": {
1745                 "mode": "LastWriterWins",
1746                 "conflictResolutionPath": "/_ts"
1747             },
1748             "computedProperties": []
1749         }
1750     }
1751 },
1752 {

```



```

1753     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/
1754         containers",
1755     "apiVersion": "2024-12-01-preview",
1756     "name": "[concat(parameters('
1757         databaseAccounts_sustainable_ai_cosmos_dba_name'), '/
1758         Development/Log')]",
1759     "dependsOn": [
1760         "[resourceId('Microsoft.DocumentDB/databaseAccounts/
1761             sqlDatabases', parameters('
1762                 databaseAccounts_sustainable_ai_cosmos_dba_name'), '
1763                 Development')]",
1764         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
1765             parameters('
1766                 databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1767     ],
1768     "properties": {
1769         "resource": {
1770             "id": "Log",
1771             "indexingPolicy": {
1772                 "indexingMode": "consistent",
1773                 "automatic": true,
1774                 "includedPaths": [
1775                     {
1776                         "path": "/*"
1777                     }
1778                 ],
1779                 "excludedPaths": [
1780                     {
1781                         "path": "/\"_etag\"/?"
1782                     }
1783                 ]
1784             },
1785             "partitionKey": {
1786                 "paths": [
1787                     "/userId"
1788                 ],
1789                 "kind": "Hash",
1790                 "version": 2
1791             },
1792             "uniqueKeyPolicy": {
1793                 "uniqueKeys": []
1794             },
1795             "conflictResolutionPolicy": {
1796                 "mode": "LastWriterWins",
1797                 "conflictResolutionPath": "/_ts"
1798             },
1799             "computedProperties": []
1800         }
1801     }
1802 },
1803 {
1804     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/
1805         containers",
1806     "apiVersion": "2024-12-01-preview",
1807     "name": "[concat(parameters('
1808         databaseAccounts_sustainable_ai_cosmos_dba_name'), '/
1809         Production/Log')]",

```

```

1799     "dependsOn": [
1800         "[resourceId('Microsoft.DocumentDB/databaseAccounts/
            sqlDatabases', parameters('
            databaseAccounts_sustainable_ai_cosmos_dba_name'), '
            Production')]",
1801         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
            parameters('
            databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1802     ],
1803     "properties": {
1804         "resource": {
1805             "id": "Log",
1806             "indexingPolicy": {
1807                 "indexingMode": "consistent",
1808                 "automatic": true,
1809                 "includedPaths": [
1810                     {
1811                         "path": "/*"
1812                     }
1813                 ],
1814                 "excludedPaths": [
1815                     {
1816                         "path": "/\"_etag\"/?"
1817                     }
1818                 ]
1819             },
1820             "partitionKey": {
1821                 "paths": [
1822                     "/userId"
1823                 ],
1824                 "kind": "Hash",
1825                 "version": 2
1826             },
1827             "uniqueKeyPolicy": {
1828                 "uniqueKeys": []
1829             },
1830             "conflictResolutionPolicy": {
1831                 "mode": "LastWriterWins",
1832                 "conflictResolutionPath": "/_ts"
1833             },
1834             "computedProperties": []
1835         }
1836     },
1837 },
1838 {
1839     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/
            containers",
1840     "apiVersion": "2024-12-01-preview",
1841     "name": "[concat(parameters('
            databaseAccounts_sustainable_ai_cosmos_dba_name'), '/
            Development/Prompt')]",
1842     "dependsOn": [
1843         "[resourceId('Microsoft.DocumentDB/databaseAccounts/
            sqlDatabases', parameters('
            databaseAccounts_sustainable_ai_cosmos_dba_name'), '
            Development')]",

```

```

1844         "[resourceId('Microsoft.DocumentDB/databaseAccounts',
1845             parameters('
1846                 databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1847     ],
1848     "properties": {
1849         "resource": {
1850             "id": "Prompt",
1851             "indexingPolicy": {
1852                 "indexingMode": "consistent",
1853                 "automatic": true,
1854                 "includedPaths": [
1855                     {
1856                         "path": "/*"
1857                     }
1858                 ],
1859                 "excludedPaths": [
1860                     {
1861                         "path": "/" + "_etag" + "/" + "?"
1862                     }
1863                 ]
1864             },
1865             "partitionKey": {
1866                 "paths": [
1867                     "/userId"
1868                 ],
1869                 "kind": "Hash",
1870                 "version": 2
1871             },
1872             "uniqueKeyPolicy": {
1873                 "uniqueKeys": []
1874             },
1875             "conflictResolutionPolicy": {
1876                 "mode": "LastWriterWins",
1877                 "conflictResolutionPath": "/" + "ts"
1878             },
1879             "computedProperties": []
1880         }
1881     },
1882     {
1883         "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/containers",
1884         "apiVersion": "2024-12-01-preview",
1885         "name": "[concat(parameters('
1886             databaseAccounts_sustainable_ai_cosmos_dba_name'), ' /
1887             Production/Prompt')]",
1888         "dependsOn": [
1889             "[resourceId('Microsoft.DocumentDB/databaseAccounts/
1890                 sqlDatabases', parameters('
1891                     databaseAccounts_sustainable_ai_cosmos_dba_name'), '
1892                     Production')]",
1893             "[resourceId('Microsoft.DocumentDB/databaseAccounts',
1894                 parameters('
1895                     databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1896         ],
1897         "properties": {
1898             "resource": {

```

```

1891         "id": "Prompt",
1892         "indexingPolicy": {
1893             "indexingMode": "consistent",
1894             "automatic": true,
1895             "includedPaths": [
1896                 {
1897                     "path": "/*"
1898                 }
1899             ],
1900             "excludedPaths": [
1901                 {
1902                     "path": "/\"_etag\"/?"
1903                 }
1904             ]
1905         },
1906         "partitionKey": {
1907             "paths": [
1908                 "/userId"
1909             ],
1910             "kind": "Hash",
1911             "version": 2
1912         },
1913         "uniqueKeyPolicy": {
1914             "uniqueKeys": []
1915         },
1916         "conflictResolutionPolicy": {
1917             "mode": "LastWriterWins",
1918             "conflictResolutionPath": "/_ts"
1919         },
1920         "computedProperties": []
1921     }
1922 },
1923 {
1924     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/containers",
1925     "apiVersion": "2024-12-01-preview",
1926     "name": "[concat(parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Development/User')]",
1927     "dependsOn": [
1928         "[resourceId('Microsoft.DocumentDB/databaseAccounts/sqlDatabases', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Development')]",
1929         "[resourceId('Microsoft.DocumentDB/databaseAccounts', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1930     ],
1931     "properties": {
1932         "resource": {
1933             "id": "User",
1934             "indexingPolicy": {
1935                 "indexingMode": "consistent",
1936                 "automatic": true,
1937                 "includedPaths": [
1938                     {
1939

```

```

1940         "path": "/*"
1941     },
1942 ],
1943 "excludedPaths": [
1944     {
1945         "path": "/\"_etag\"/?"
1946     }
1947 ],
1948 },
1949 "partitionKey": {
1950     "paths": [
1951         "/userId"
1952     ],
1953     "kind": "Hash",
1954     "version": 2
1955 },
1956 "uniqueKeyPolicy": {
1957     "uniqueKeys": []
1958 },
1959 "conflictResolutionPolicy": {
1960     "mode": "LastWriterWins",
1961     "conflictResolutionPath": "/_ts"
1962 },
1963 "computedProperties": []
1964 }
1965 }
1966 },
1967 {
1968     "type": "Microsoft.DocumentDB/databaseAccounts/sqlDatabases/containers",
1969     "apiVersion": "2024-12-01-preview",
1970     "name": "[concat(parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Production/User')]",
1971     "dependsOn": [
1972         "[resourceId('Microsoft.DocumentDB/databaseAccounts/sqlDatabases', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'), 'Production')]",
1973         "[resourceId('Microsoft.DocumentDB/databaseAccounts', parameters('databaseAccounts_sustainable_ai_cosmos_dba_name'))]"
1974     ],
1975     "properties": {
1976         "resource": {
1977             "id": "User",
1978             "indexingPolicy": {
1979                 "indexingMode": "consistent",
1980                 "automatic": true,
1981                 "includedPaths": [
1982                     {
1983                         "path": "/*"
1984                     }
1985                 ],
1986                 "excludedPaths": [
1987                     {
1988                         "path": "/\"_etag\"/?"

```

```

1989         }
1990     ],
1991     },
1992     "partitionKey": {
1993         "paths": [
1994             "/userId"
1995         ],
1996         "kind": "Hash",
1997         "version": 2
1998     },
1999     "uniqueKeyPolicy": {
2000         "uniqueKeys": []
2001     },
2002     "conflictResolutionPolicy": {
2003         "mode": "LastWriterWins",
2004         "conflictResolutionPath": "/_ts"
2005     },
2006     "computedProperties": []
2007 }
2008 }
2009 },
2010 {
2011     "type": "Microsoft.Storage/storageAccounts/blobServices/
containers",
2012     "apiVersion": "2024-01-01",
2013     "name": "[concat(parameters('storageAccounts_ip6sustainableai91
a8_name'), '/default/azure-webjobs-hosts')]",
2014     "dependsOn": [
2015         "[resourceId('Microsoft.Storage/storageAccounts/
blobServices', parameters('storageAccounts_ip6
sustainableai91a8_name'), 'default')]",
2016         "[resourceId('Microsoft.Storage/storageAccounts',
parameters('storageAccounts_ip6sustainableai91a8_name'))
]"
2017     ],
2018     "properties": {
2019         "immutableStorageWithVersioning": {
2020             "enabled": false
2021         },
2022         "defaultEncryptionScope": "$account-encryption-key",
2023         "denyEncryptionScopeOverride": false,
2024         "publicAccess": "None"
2025     }
2026 },
2027 {
2028     "type": "Microsoft.Storage/storageAccounts/blobServices/
containers",
2029     "apiVersion": "2024-01-01",
2030     "name": "[concat(parameters('storageAccounts_ip6sustainableai91
a8_name'), '/default/azure-webjobs-secrets')]",
2031     "dependsOn": [
2032         "[resourceId('Microsoft.Storage/storageAccounts/
blobServices', parameters('storageAccounts_ip6
sustainableai91a8_name'), 'default')]",
2033         "[resourceId('Microsoft.Storage/storageAccounts',
parameters('storageAccounts_ip6sustainableai91a8_name'))
]"

```

```

2034     ],
2035     "properties": {
2036         "immutableStorageWithVersioning": {
2037             "enabled": false
2038         },
2039         "defaultEncryptionScope": "$account-encryption-key",
2040         "denyEncryptionScopeOverride": false,
2041         "publicAccess": "None"
2042     }
2043 },
2044 {
2045     "type": "Microsoft.Storage/storageAccounts/fileServices/shares"
2046     ,
2047     "apiVersion": "2024-01-01",
2048     "name": "[concat(parameters('storageAccounts_ip6sustainableai91a8_name'), '/default/sustainable-ai-api-35026538')]",
2049     "dependsOn": [
2050         "[resourceId('Microsoft.Storage/storageAccounts/fileServices', parameters('storageAccounts_ip6sustainableai91a8_name'), 'default')]",
2051         "[resourceId('Microsoft.Storage/storageAccounts', parameters('storageAccounts_ip6sustainableai91a8_name'))]"
2052     ],
2053     "properties": {
2054         "accessTier": "TransactionOptimized",
2055         "shareQuota": 102400,
2056         "enabledProtocols": "SMB"
2057     }
2058 },
2059 {
2060     "type": "Microsoft.Storage/storageAccounts/fileServices/shares"
2061     ,
2062     "apiVersion": "2024-01-01",
2063     "name": "[concat(parameters('storageAccounts_ip6sustainableai91a8_name'), '/default/sustainable-ai-api-350265382f84')]",
2064     "dependsOn": [
2065         "[resourceId('Microsoft.Storage/storageAccounts/fileServices', parameters('storageAccounts_ip6sustainableai91a8_name'), 'default')]",
2066         "[resourceId('Microsoft.Storage/storageAccounts', parameters('storageAccounts_ip6sustainableai91a8_name'))]"
2067     ],
2068     "properties": {
2069         "accessTier": "TransactionOptimized",
2070         "shareQuota": 102400,
2071         "enabledProtocols": "SMB"
2072     }
2073 },
2074 {
2075     "type": "Microsoft.Storage/storageAccounts/fileServices/shares"
2076     ,
2077     "apiVersion": "2024-01-01",
2078     "name": "[concat(parameters('storageAccounts_ip6sustainableai91a8_name'), '/default/sustainable-ai-apib12e')]",
2079     "dependsOn": [

```

```

2077         "[resourceId('Microsoft.Storage/storageAccounts/
           fileServices', parameters('storageAccounts_ip6
           sustainableai91a8_name'), 'default')]",
2078         "[resourceId('Microsoft.Storage/storageAccounts',
           parameters('storageAccounts_ip6sustainableai91a8_name'))
           ]"
2079     ],
2080     "properties": {
2081         "accessTier": "TransactionOptimized",
2082         "shareQuota": 102400,
2083         "enabledProtocols": "SMB"
2084     }
2085 },
2086 {
2087     "type": "Microsoft.Storage/storageAccounts/tableServices/tables",
2088     "apiVersion": "2024-01-01",
2089     "name": "[concat(parameters('storageAccounts_ip6sustainableai91
           a8_name'), '/default/AzureFunctionsDiagnosticEvents202507')]",
2090     "dependsOn": [
2091         "[resourceId('Microsoft.Storage/storageAccounts/
           tableServices', parameters('storageAccounts_ip6
           sustainableai91a8_name'), 'default')]",
2092         "[resourceId('Microsoft.Storage/storageAccounts',
           parameters('storageAccounts_ip6sustainableai91a8_name'))
           ]"
2093     ],
2094     "properties": {}
2095 },
2096 {
2097     "type": "Microsoft.Web/staticSites/linkedBackends",
2098     "apiVersion": "2024-04-01",
2099     "name": "[concat(parameters('
           staticSites_sustainable_ai_web_app_name'), '/backend1')]",
2100     "location": "West Europe",
2101     "dependsOn": [
2102         "[resourceId('Microsoft.Web/staticSites', parameters('
           staticSites_sustainable_ai_web_app_name'))]",
2103         "[resourceId('Microsoft.Web/sites', parameters('
           sites_sustainable_ai_api_name'))]"
2104     ],
2105     "properties": {
2106         "backendResourceId": "[resourceId('Microsoft.Web/sites',
           parameters('sites_sustainable_ai_api_name'))]",
2107         "region": "switzerlandnorth"
2108     }
2109 },
2110 {
2111     "type": "Microsoft.Web/staticSites/userProvidedFunctionApps",
2112     "apiVersion": "2024-04-01",
2113     "name": "[concat(parameters('
           staticSites_sustainable_ai_web_app_name'), '/backend1')]",
2114     "location": "West Europe",
2115     "dependsOn": [
2116         "[resourceId('Microsoft.Web/staticSites', parameters('
           staticSites_sustainable_ai_web_app_name'))]",

```



```

2117         "[resourceId('Microsoft.Web/sites', parameters('
2118             sites_sustainable_ai_api_name'))]"
2119     ],
2120     "properties": {
2121         "functionAppResourceId": "[resourceId('Microsoft.Web/sites'
2122             , parameters('sites_sustainable_ai_api_name'))]",
2123         "functionAppRegion": "switzerlandnorth"
2124     }
2125 ]
2126 }

```

### A.3.2 Matlab Scripts

```

1  %% === 0. Load files ===
2  file_prompts = 'Data_Pseudonym.xlsx';
3
4  % Read sheets into tables
5  prompts = readtable(file_prompts, 'Sheet', 'Prompts');
6  logs = readtable(file_prompts, 'Sheet', 'Logs');
7  conversations = readtable(file_prompts, 'Sheet', 'Conversations');
8  usersTable = readtable(file_prompts, 'Sheet', 'Users');
9
10 %% === 1. Preprocessing ===
11 % Interpret enums as categorical variables
12 prompts = prompts(prompts.isSent == true, :);
13 % Define custom chat mode order (including 'Total')
14 prompts.chatMode = categorical(prompts.chatMode, [0 1 2], {'Energy
    efficient', 'Balanced', 'Performance'});
15
16 % Convert timestamps to datetime
17 prompts.sentAt = datetime(prompts.sentAt, 'InputFormat', 'yyyy-MM-dd''T''HH
    :mm:ss');
18 prompts.createdAt = datetime(prompts.createdAt, 'InputFormat', 'yyyy-MM-dd
    ''T''HH:mm:ss');
19 %% === 2. Create new columns ===
20 prompts.responseLength = strlength(string(prompts.responseText));
21
22 %% === 3. Grouped evaluation (with all user-mode combos and totals) ===
23
24 % Get all unique users and all modes
25 allUsers = unique(prompts.userId);
26 allModes = categories(prompts.chatMode);
27
28 % Create full combination of users and modes
29 [U, M] = ndgrid(allUsers, allModes);
30 comboTable = table;
31 comboTable.userId = reshape(U, [], 1);
32 comboTable.chatMode = categorical(reshape(M, [], 1), allModes);
33
34 % Group actual data
35 G = findgroups(prompts.userId, prompts.chatMode);
36 T_actual = table;
37
38 T_actual.userId = splitapply(@(x) x(1), prompts.userId, G);
39 T_actual.chatMode = splitapply(@(x) x(1), prompts.chatMode, G);

```

```

40 T_actual.NumberOfPrompts = splitapply(@sum, prompts.id, G);
41 T_actual.InputTokens = splitapply(@sum, prompts.usage_numberOfInputTokens,
    G);
42 T_actual.OutputTokens = splitapply(@sum, prompts.usage_numberOfOutputTokens
    , G);
43 T_actual.TotalUsageWh = splitapply(@sum, prompts.usage_usageInWh, G);
44 T_actual.TotalUsageWhCorrected = splitapply(@sum, prompts.
    usageInWhCorrected, G);
45
46 % Join full combination with actual data to ensure 0s are included
47 T = outerjoin(comboTable, T_actual, ...
48     'Keys', {'userId', 'chatMode'}, ...
49     'MergeKeys', true);
50
51 % Replace NaNs with 0 for numeric columns
52 T.NumberOfPrompts(isnan(T.NumberOfPrompts)) = 0;
53 T.InputTokens(isnan(T.InputTokens)) = 0;
54 T.OutputTokens(isnan(T.OutputTokens)) = 0;
55 T.TotalUsageWh(isnan(T.TotalUsageWh)) = 0;
56 T.TotalUsageWhCorrected(isnan(T.TotalUsageWhCorrected)) = 0;
57
58 % Add total row per user
59 G_user = findgroups(T.userId);
60 T_userTotal = table;
61 T_userTotal.userId = splitapply(@(x) x(1), T.userId, G_user);
62 T_userTotal.chatMode = categorical(repmat("Total", height(T_userTotal), 1),
    ...
63     [allModes; "Total"]);
64 T_userTotal.NumberOfPrompts = splitapply(@sum, T.NumberOfPrompts, G_user);
65 T_userTotal.InputTokens = splitapply(@sum, T.InputTokens, G_user);
66 T_userTotal.OutputTokens = splitapply(@sum, T.OutputTokens, G_user);
67 T_userTotal.TotalUsageWh = splitapply(@sum, T.TotalUsageWh, G_user);
68 T_userTotal.TotalUsageWhCorrected = splitapply(@sum, T.
    TotalUsageWhCorrected, G_user);
69
70 % Combine mode rows and total rows
71 T = [T; T_userTotal];
72
73 % Sort nicely by user, then mode
74 T = sortrows(T, {'userId', 'chatMode'});
75
76 % Add % column (only for modes, not total)
77 % Compute total prompts per user (one row per user)
78 userTotalPrompts = groupsummary(T, "userId", "max", "NumberOfPrompts");
79 userTotalPrompts.Properties.VariableNames{'max_NumberOfPrompts'} = '
    TotalPromptsPerUser';
80
81 % Join this summary back into T
82 T = outerjoin(T, userTotalPrompts(:, {'userId', 'TotalPromptsPerUser'}),
    ...
83     'Keys', 'userId', 'MergeKeys', true);
84 T.PctPromptsPerUserMode = 100 * (T.NumberOfPrompts ./ T.TotalPromptsPerUser
    );
85 T.TotalPromptsPerUser = []; % Remove helper column from final table
86
87 T.PctPromptsPerUserMode(T.chatMode == "Total") = 100; % blank for totals
88

```

```

89 %% === 4. Display grouped results ===
90 disp(T);
91
92 %% === 5. Plot: Number of prompts per mode ===
93 figure('Name','Prompts per Mode');
94 modes = categories(T.chatMode);
95 counts = zeros(numel(modes),1);
96 for i = 1:numel(modes)
97     counts(i) = sum(T.NumberOfPrompts(T.chatMode == modes{i}));
98 end
99 bar(categorical(modes), counts);
100 ylabel('Number of Prompts');
101 title('Total Prompts per Mode');
102 grid on;
103
104 %% === 6. Plot: Energy usage per user and mode ===
105 figure('Name','Energy Usage per User and Mode');
106 barData = unstack(T(:, {'userId', 'chatMode', 'TotalUsageWh'}), '
    TotalUsageWh', 'chatMode');
107 bar(categorical(barData.userId), barData{:,2:end});
108 xlabel('User ID');
109 ylabel('Energy Usage (Wh)');
110 title('Total Energy Usage per User & Mode');
111 legend(barData.Properties.VariableNames(2:end), 'Location', 'northwest');
112 grid on;
113
114 %% === 11. Display final table ===
115 disp('Evaluation per User and Mode:');
116 disp(T);
117
118 %% === 12. Aggregated evaluation per mode ===
119 [Gm, modes] = findgroups(T.chatMode);
120 Agg = table;
121
122 Agg.chatMode = modes;
123 Agg.InputTokens = splitapply(@sum, T.InputTokens, Gm);
124 Agg.OutputTokens = splitapply(@sum, T.OutputTokens, Gm);
125 Agg.TotalUsageWh = splitapply(@sum, T.TotalUsageWh, Gm);
126 Agg.TotalUsageWhCorrected = splitapply(@sum, T.TotalUsageWhCorrected, Gm);
127 Agg.TotalPrompts = splitapply(@sum, T.NumberOfPrompts, Gm);
128
129 % Calculate total prompts and total energy usage (use MAX because every
130 % mode also has a total)
131 totalPromptsAll = max(Agg.TotalPrompts);
132 totalWhAll = max(Agg.TotalUsageWh);
133 totalWhCorrectedAll = max(Agg.TotalUsageWhCorrected);
134
135 % Calculate percentage columns
136 Agg.PctPromptsPerMode = (Agg.TotalPrompts ./ totalPromptsAll) * 100;
137 Agg.PctUsagePerMode = (Agg.TotalUsageWh ./ totalWhAll) * 100;
138 Agg.PctUsageCorrectedPerMode = (Agg.TotalUsageWhCorrected ./
    totalWhCorrectedAll) * 100;
139
140 disp('Aggregated values per Mode:');
141 disp(Agg);
142
143 %% === Prompts per User and Day (Split by Mode and Total) ===

```

```

144 % Group by userId, day, and chatMode
145 G = findgroups(prompts.userId, prompts.day, prompts.chatMode);
146 summaryTable = table;
147 summaryTable.userId = splitapply(@(x) x(1), prompts.userId, G);
148 summaryTable.day = splitapply(@(x) x(1), prompts.day, G);
149 summaryTable.chatMode = splitapply(@(x) x(1), prompts.chatMode, G);
150 summaryTable.numPrompts = splitapply(@numel, prompts.id, G);
151
152 % Get all unique combinations
153 users = unique(prompts.userId);
154 days = unique(prompts.day);
155 modes = categories(prompts.chatMode);
156
157 % Generate variable names for each day-mode and day-total
158 varNames = {};
159 for d = days'
160     for m = modes'
161         varNames{end+1} = sprintf('Day%d_%s', d, matlab.lang.makeValidName(
162             char(m)));
163     end
164     varNames{end+1} = sprintf('Day%d_Total', d);
165 end
166
167 % Initialize result table
168 result = array2table(zeros(numel(users), numel(varNames)), ...
169     'VariableNames', varNames);
170 result.userId = users;
171
172 % Fill in prompt counts per user/day/mode
173 for i = 1:height(summaryTable)
174     uid = summaryTable.userId(i);
175     day = summaryTable.day(i);
176     mode = summaryTable.chatMode(i);
177     n = summaryTable.numPrompts(i);
178
179     rowIdx = find(result.userId == uid);
180     colName = sprintf('Day%d_%s', day, matlab.lang.makeValidName(char(mode)
181         ));
182     result{rowIdx, colName} = result{rowIdx, colName} + n;
183
184     % Update total
185     totalCol = sprintf('Day%d_Total', day);
186     result{rowIdx, totalCol} = result{rowIdx, totalCol} + n;
187 end
188
189 % Move userId to the front
190 result = movevars(result, 'userId', 'Before', 1);
191
192 % Display result
193 disp('Prompt Matrix per User, Day, and Mode with Totals:');
194 disp(result);
195
196 %% === Prompts per User and Day ===
197 % Group prompts by userId and day, count number of prompts
198 G_day = findgroups(prompts.userId, prompts.day);
199 userDayTable = table;

```

```

199 userDayTable.userId = splitapply(@(x) x(1), prompts.userId, G_day);
200 userDayTable.day = splitapply(@(x) x(1), prompts.day, G_day);
201 userDayTable.NumberOfPrompts = splitapply(@numel, prompts.id, G_day);
202
203 % Display the table
204 disp('Number of prompts per User and Day:');
205 disp(userDayTable);
206
207 % Unique users and days
208 users = unique(userDayTable.userId);
209 days = unique(userDayTable.day);
210
211 % Build prompt matrix: rows = users, columns = days
212 promptMatrix = zeros(numel(users), numel(days));
213 for i = 1:numel(users)
214     for j = 1:numel(days)
215         idx = userDayTable.userId == users(i) & userDayTable.day == days(j)
216             ;
217         if any(idx)
218             promptMatrix(i,j) = userDayTable.NumberOfPrompts(idx);
219         end
220     end
221 end
222
223 % === Grouped Bar Chart: Per User (X) and Day (grouped bars) ===
224 figure('Name','Prompts per User and Day (Bar Chart)');
225 bar(users, promptMatrix, 'grouped');
226 xlabel('User ID');
227 ylabel('Number of Prompts');
228 title('Number of Prompts per User and Day');
229 legend(arrayfun(@(d) sprintf('Day %d', d), days, 'UniformOutput', false),
230         ...
231         'Location', 'northeastoutside');
232 grid on;
233
234 % === Line Chart: Per User (X) and Day (one line per day) ===
235 figure('Name', 'Prompts per User and Day (Line Chart)');
236 hold on;
237
238 for j = 1:numel(days)
239     plot(users, promptMatrix(:,j), '-o', 'DisplayName', sprintf('Day %d',
240         days(j)));
241 end
242
243 % Plot average across days
244 avgPrompts = mean(promptMatrix, 2);
245 plot(users, avgPrompts, '-k', 'LineWidth', 2, 'DisplayName', 'Average');
246
247 xlabel('User ID');
248 ylabel('Number of Prompts');
249 title('Number of Prompts per User and Day with Average');
250 legend('Location', 'northeastoutside');
251 grid on;
252 hold off;

```

```

253 %% === Plot: Chat Mode Usage per Day ===
254
255 % Group by day and chat mode
256 G_dayMode = findgroups(prompts.day, prompts.chatMode);
257 modeDayTable = table;
258 modeDayTable.day = splitapply(@(x) x(1), prompts.day, G_dayMode);
259 modeDayTable.chatMode = splitapply(@(x) x(1), prompts.chatMode, G_dayMode);
260 modeDayTable.NumberOfPrompts = splitapply(@numel, prompts.id, G_dayMode);
261
262 % Prepare matrix: rows = days, columns = modes
263 days = unique(modeDayTable.day);
264 modes = categories(prompts.chatMode);
265 modeMatrix = zeros(numel(days), numel(modes));
266
267 for i = 1:numel(days)
268     for j = 1:numel(modes)
269         idx = (modeDayTable.day == days(i)) & (modeDayTable.chatMode ==
270             modes{j});
271         if any(idx)
272             modeMatrix(i, j) = modeDayTable.NumberOfPrompts(idx);
273         else
274             modeMatrix(i, j) = 0;
275         end
276     end
277 end
278
279 % Plot stacked bar chart
280 figure('Name','Chat Mode Usage per Day');
281 bar(days, modeMatrix, 'stacked');
282 xlabel('Day');
283 ylabel('Number of Prompts');
284 title('Chat Mode Usage per Day');
285 legend(modes, 'Location', 'northeastoutside');
286 grid on;
287
288 %% === Plot: Chat Mode Usage per Day (Percentage) ===
289
290 % Normalize modeMatrix to percentages
291 modeMatrixPct = modeMatrix ./ sum(modeMatrix, 2) * 100;
292
293 % Handle potential division by zero (in case a day has no prompts)
294 modeMatrixPct(isnan(modeMatrixPct)) = 0;
295
296 % Plot stacked percentage bar chart
297 figure('Name','Chat Mode Usage per Day (Percentage)');
298 bar(days, modeMatrixPct, 'stacked');
299 xlabel('Day');
300 ylabel('Percentage of Prompts');
301 title('Chat Mode Usage per Day (Percentage)');
302 legend(modes, 'Location', 'northeastoutside');
303 grid on;
304
305 %% === Load metrics ===
306 metrics = logs(strcmp(logs.message, '/metrics'), :);
307
308 %% === Convert datetime day to numeric weekday 1=Monday ... 5=Friday ===
309 wday = weekday(metrics.day);

```

```

309 wday_adj = wday - 1;          % Monday=1 ... Sunday=0
310 wday_adj(wday_adj == 0) = 7; % Sunday=7
311
312 % Keep only Monday to Friday
313 validDaysIdx = wday_adj >= 1 & wday_adj <= 5;
314 metrics = metrics(validDaysIdx, :);
315 metrics.dayNum = wday_adj(validDaysIdx);
316
317 %% === Prepare full user-day grid (all users    days 1 to 5) ===
318 days = (1:5)';
319 [U, D] = ndgrid(allUsers, days);
320 combo = table;
321 combo.userId = reshape(U, [], 1);
322 combo.dayNum = reshape(D, [], 1);
323
324 %% === Group metrics by userId and dayNum ===
325 G = findgroups(metrics.userId, metrics.dayNum);
326 T_visits = table;
327 T_visits.userId = splitapply(@(x) x(1), metrics.userId, G);
328 T_visits.dayNum = splitapply(@(x) x(1), metrics.dayNum, G);
329 T_visits.PageVisits = splitapply(@numel, metrics.message, G);
330
331 %% === Outer join full grid with actual counts ===
332 T_full = outerjoin(combo, T_visits, ...
333     'Keys', {'userId', 'dayNum'}, ...
334     'MergeKeys', true);
335
336 % Replace missing visits with zero
337 T_full.PageVisits(isnan(T_full.PageVisits)) = 0;
338
339 %% === Pivot to wide format: one row per user, columns Day1...Day5 ===
340 T_wide = unstack(T_full, 'PageVisits', 'dayNum', 'VariableNamingRule', '
    preserve');
341
342 % Rename columns for clarity
343 dayCols = strcat("Day", string(days));
344 T_wide.Properties.VariableNames(2:end) = dayCols;
345
346 %% === Display result ===
347 disp('Page Visits per Day and User (Monday=1 to Friday=5):');
348 disp(T_wide);

```

```

1 % Load Excel data
2 file = 'Data_Pseudonym.xlsx';
3 sheet = 'Prompts';
4 data = readtable(file, 'Sheet', sheet);
5
6 % Extract input and output tokens
7 x = data.usage_numberOfInputTokens;
8 y = data.usage_numberOfOutputTokens;
9
10 % Filter out invalid entries
11 valid = ~isnan(x) & ~isnan(y) & x > 0 & y > 0;
12 x = x(valid);
13 y = y(valid);
14
15 % Sort by input tokens for smooth plotting

```

```

16 [x_sorted, idx] = sort(x);
17 y_sorted = y(idx);
18
19 % LOWESS smoothing to reveal real relationship
20 y_smooth = smooth(x_sorted, y_sorted, 0.05, 'lowess'); % 0.05 = smoothing
    span (adjust as needed)
21
22 % Reference line y = x
23 y_ref = x_sorted;
24
25 % Optional: Crop axes to the 95th percentile to avoid outliers
26 x_max = prctile(x, 95);
27 y_max = prctile(y, 94.83); % Skip one point to clean the graph image
28
29 % Plot
30 figure;
31 scatter(x, y, 10, 'filled', 'MarkerFaceAlpha', 0.3);
32 hold on;
33 plot(x_sorted, y_smooth, 'b-', 'LineWidth', 2); % LOWESS trend
34 plot(x_sorted, y_ref, 'g--', 'LineWidth', 1.5); % y = x reference
35
36 % Formatting
37 xlim([0, x_max]);
38 ylim([0, y_max]);
39 xlabel('Input Tokens');
40 ylabel('Output Tokens');
41 title('Smoothed Relation between Input and Output Tokens');
42 legend('Prompt Data', 'LOWESS Smoothed', 'y = x', 'Location', 'southeast');
43 grid on;

```

```

1 % Load Excel file and sheet
2 file = 'Data.xlsx'; % Adjust path if needed
3 sheet = 'Prompts';
4
5 % Read the table
6 data = readtable(file, 'Sheet', sheet, 'ReadVariableNames', true);
7 lengthsRaw = data.promptTextHistoryLengths;
8 numPrompts = height(data);
9
10 % Number of interpolation points (e.g., representing 0 100 %)
11 nInterp = 100;
12
13 % Matrix to store all normalized prompt curves
14 normalizedCurves = NaN(numPrompts, nInterp);
15
16 row = 1; % Row counter for valid prompts
17
18 for i = 1:numPrompts
19     % Convert string to numeric array
20     str = lengthsRaw{i};
21     nums = sscanf(str, '%d,', Inf); % read comma-separated values
22     if isempty(nums)
23         parts = split(str, ',');
24         nums = str2double(parts);
25     end
26     nums = nums(:)'; % ensure row vector
27

```



```

28     % Skip if not enough points for interpolation
29     if length(nums) < 2
30         continue;
31     end
32
33     % Normalize x-values to range [0, 1] (prompt progress)
34     steps = linspace(0, 1, length(nums));
35
36     % Normalize y-values to final length (range [0, 1])
37     finalLen = nums(end);
38     if finalLen == 0 || any(isnan(nums))
39         continue;
40     end
41     yNorm = nums / finalLen;
42
43     % Interpolate to a fixed number of x-values (0 1 scale)
44     xInterp = linspace(0, 1, nInterp);
45     yInterp = interp1(steps, yNorm, xInterp, 'linear', 'extrap');
46
47     % Store in matrix
48     normalizedCurves(row, :) = yInterp;
49     row = row + 1;
50 end
51
52 % Compute average curve across prompts
53 meanCurve = nanmean(normalizedCurves, 1);
54
55 % Smooth the average curve (optional)
56 smoothCurve = smooth(meanCurve, 0.2, 'loess'); % 0.2 = smoothing factor
57
58 % Plot
59 figure;
60 plot(linspace(0, 100, nInterp), smoothCurve * 100, 'LineWidth', 2);
61 xlabel('Prompt progress [%]');
62 ylabel('Text length [% of final length]');
63 title('Average normalized prompt growth curve');
64 grid on;
65
66 % Optional: show how many prompts were used
67 fprintf('Used prompts: %d of %d\n', row-1, numPrompts);

```

### A.3.3 Project Documents from FHNW

Windisch, 24.03.25

## **Project proposal for IP6, 25FS\_IMVS24**

### **Designing towards higher user awareness: UI strategies for reducing conversational AI energy consumption**

**Students:** Jack Gläser  
Simon Lüscher

**Supervisors:** Martin Kropp  
Nitish Patkar

**Stakeholder:** Fachhochschule Nordwestschweiz FHNW

**Project duration:** 17.02.2025 - 14.08.2025

Version	Note	Autor
1	Initial version	Jack Gläser, Simon Lüscher
2	Draft	Jack Gläser, Simon Lüscher
3	Final	Jack Gläser, Simon Lüscher

## 1. Initial situation

Conversational AI consumes a lot of energy during its whole lifecycle. A significant part of this energy is required for the development and training of the model. But the actual conversational AI service consumes a great amount of energy as well. While a single prompt does not have a significant impact on the required energy it grows rapidly with a higher number of users. Over the period of a month the overall energy usage is already more than that of the training of the LLM. <sup>1</sup> The number of users is constantly increasing and already reached 987 million<sup>2</sup> in 2025. This adds even more significance to the inference part of the lifecycle of conversational AI since with a higher number of users the usage will increase even more.

Furthermore, the phase where users infer with the model also has an impact on previous lifecycle phases. If for example a user can be directed to prompt in a more optimized way the models could be trained in an even more specific direction and with a less generic dataset. Ideally this would further reduce the energy consumption. <sup>3</sup> Although the effect size is not clear, and this should be further investigated it adds to the weight of the inference phase.

Another issue that we experienced ourselves and have seen with other people is that the trend is more to ask a conversational AI for an answer than using a search engine. Even for trivial questions. We hypothesize that users are not aware of the high energy consumption of a single inference with a conversational AI. Which is estimated to be at around 0.005 kWh per inference for GPT-4 for example. In comparison a search on google would only consume about 0.0003 kWh which is more than 16 times less.

---

<sup>1</sup> <https://www.sciencedirect.com/science/article/pii/S2095809924002315#s0010>

<sup>2</sup> <https://www.demandsage.com/chatbot-statistics/>

<sup>3</sup> <https://www.sciencedirect.com/science/article/pii/S2095809924002315#s0010>

## 2. Problem statement

The increasing adoption of conversational AI over traditional search engines is driving a significant rise in energy consumption. A 2023 study estimates that if users were to replace Google searches with interactions using large language models (LLMs), the energy demand for Google's AI alone would reach 29.3 TWh per year—equivalent to the electricity consumption of an entire country like Ireland.<sup>4</sup> This shift raises concerns about the sustainability of AI-powered search and the need for energy-efficient solutions.

Despite these implications, most users are unaware of the significant energy consumption associated with AI-driven interactions. Unlike traditional searches, which have well-optimized infrastructure with relatively low energy costs per query, LLM-based responses require substantial computational power. Since this energy usage remains largely invisible to end users, there is little public awareness or discussion on how AI adoption impacts global energy demand. This lack of awareness may slow down efforts to develop and implement more energy-efficient alternatives.

Furthermore, there is a lack of research specifically focusing on the energy consumption of AI inference in chatbots. While training LLMs is known to be highly energy-intensive, the long-term impact of frequent real-time inference remains underexplored. Companies developing these AI models do not transparently disclose their energy usage, making it difficult to assess the true environmental cost of conversational AI. This lack of data hinders informed decision-making and the development of more sustainable AI architectures.

## 3. Project vision

We strive to reduce the overall energy consumption associated with conversational AI by enhancing user awareness through targeted UI-based interventions. By researching and contributing to the state of the art, our project aims to understand how interface design alone can effectively influence user behavior regarding energy consumption. We will identify and evaluate UI features that increase user awareness, providing transparency and actionable insights into the energy impact of their interactions with conversational AI. Ultimately, our goal is to empower users to make informed decisions when and how to use conversational AI, fostering more sustainable and energy-efficient usage patterns without compromising user experience.

### Research contributions:

We aim to present the state-of-the-art research on user awareness of the energy consumption associated with conversational AI. Additionally, we seek to conduct our own study to address existing knowledge gaps in this area.

### Impact of UI functionalities on user awareness:

We want to research the potential impact that purely UI-based functionalities can/can't have in increasing user awareness and influencing user behavior to lower energy usage.

### Increasing awareness in conversational AI:

Through implementing those functionalities and strategies, we intend to elevate overall user consciousness about their energy implications of their interactions with conversational AI.

---

<sup>4</sup> [https://www.cell.com/joule/fulltext/S2542-4351\(23\)00365-3](https://www.cell.com/joule/fulltext/S2542-4351(23)00365-3)

#### **4. Goals and research questions**

##### **4.1 Goals**

1. Identify and evaluate methods to effectively measure user awareness of conversational AI energy consumption
2. Design and evaluate UI-only features that successfully enhance user awareness of energy consumption
3. Influence user behavior through increased awareness, resulting in decreased chatbot energy consumption.
4. Ensure all interventions maintain or enhance the user experience, preserving performance, responsiveness, and ease of use
5. (optional): Develop predictive capabilities to estimate a prompt's energy consumption based on user inputs and interactions.
6. (optional): Quantify potential savings (energy, tokens, API calls, CPU cycles) achievable through increased user awareness

##### **4.2 Research questions:**

- A. To what extent are users currently aware of the energy implications associated with their chatbot interactions?
- B. (Exploratory question): How can UI-based features most effectively increase user awareness regarding the energy consumption of conversational AI?
- C. How strongly does increased user awareness correlate with reductions in conversational AI energy consumption?

## 5. Methodology

### 1. Literature Review

- Conduct a systematic review of existing **research** on:
  - Energy consumption in AI and conversational AI models.
  - Current strategies to enhance user awareness about energy consumption.
  - Existing UI-based interventions and their effectiveness in influencing user behavior.

### 2. Initial User Awareness Analysis

- **Survey:**
  - Develop and distribute a questionnaire to gauge users' current awareness of the energy impact of conversational AI.
  - Identify baseline awareness levels, knowledge gaps, and perceptions.

### 3. Design and Implementation of UI-based Interventions

- **Prototyping:**
  - Generate multiple UI design concepts aimed at increasing user awareness (e.g., visual energy indicators, impact scores, prompt suggestions...).
  - Prototype selected UI features, emphasizing transparency, ease of understanding, and actionability.
  - Implement selected UI interventions into a conversational AI environment, enabling real or simulated interactions.

### 4. Experimental Evaluation

- **User Testings:**
  - Designing and conducting a controlled experiment comparing user behavior with and without UI-based interventions.
  - Utilize metrics like the number of prompts, prompt length, frequency, and complexity of interactions to evaluate behavioral changes.
  - Assess UX impact through usability testing and user satisfaction surveys to ensure interventions do not degrade the conversational AI experience (performance, ease of use, responsiveness).
- **Energy Consumption Measurement:**
  - Deploy a monitoring system to measure and record energy usage, tokens consumed, API calls, and CPU cycles during interaction sessions.
  - Compare energy metrics before and after deploying UI-based interventions.
  - Use collected data to build and validate predictive models estimating the energy consumption of user prompts based on historical and behavioral data.

### 5. Data Analysis and Interpretation

- **Analysis:**
  - Perform statistical analysis to identify significant differences and relationships between user awareness, behavior, and energy consumption.
  - Analyze feedback from surveys and interviews using thematic analysis to identify qualitative insights into user experiences and perceptions of UI features.

## 6. Planning

We will follow an iterative approach throughout the project. Tasks will be created in GitLab and managed on a Board.<sup>5</sup> Our workflow is structured into two-week sprints, with a biweekly planning meeting to review the previous sprint and plan the next one. Coaches will receive updates either in person, if a meeting is scheduled, or via Teams in our dedicated channel. Apart from the planning meetings, all discussions and conversations will take place on Teams, with no additional meetings required.

The most important milestones and deadlines for the project are visible on GitLab<sup>6</sup> or in the following listing.

### Deadlines:

Start project:	17.02
Finish project proposal:	18.03
Research, Literatur & Survey:	17.02– 30.03
Weg zur Thesis Workshop 1	15.03
Infrastructure setup:	15.03 to 31.03 (fully closed on 15.05 for Nitish)
Implementation phase:	01.04– 25.05
Weg zur Thesis Workshop 2	10.04 or 12.04
Experiment / User testings:	26.05– 08.06
Analysis of results from Experiment & Finish report:	09.06– 22.06
Weg zur Thesis Workshop 3	5.06 or 7.06
Sent first draft of report to Coaches:	22.06
Feedback and last adjustments:	23.06– 13.07
Hand-In Thesis:	~18.07
Theoretical Deadline:	14.08
Work on presentation & defence:	14.08 – presentation date (TBD)
Defence:	~ 01.09 bis 12.09

<sup>5</sup> [https://gitlab.fhnw.ch/groups/25fs\\_imvs\\_ip6/-/boards](https://gitlab.fhnw.ch/groups/25fs_imvs_ip6/-/boards)

<sup>6</sup> [https://gitlab.fhnw.ch/groups/25fs\\_imvs\\_ip6/-/milestones](https://gitlab.fhnw.ch/groups/25fs_imvs_ip6/-/milestones)

## 7. Risiko Assessment

Risk	Level	Impact	Prevention / Strategy
Too high workload	Medium	Overloading of team members, loss of quality in work results	Iterative Approach
Absence of a team member	Low	Delays in project implementation, increased work pressure for remaining team member	Knowledge transfer through documentation of work processes and task distribution as well as reducing the scope.
Miscommunication in team	Low - Medium	Misunderstandings, inefficient collaboration	Regular meetings to clarify questions and ensure a smooth exchange of information
Insufficient communication with coaches	Low	Misunderstandings, unclear requirements, lack of support	Biweekly meetings and regular communication on Teams
Technical challenges	Medium	Delays in development	Read literature and support each other in learning new things, seek expert advice.
Time management problems	Medium	Delays in completion, increased time pressure	Realistic planning and stick to planned meetings and sprints
Topic is relatively new territory, finding research is difficult	Medium	Motivation drops, work stagnates	Take a break in case of blockages and ask your coaches for advice



## 25FS\_IMVS24: Better ways to manage prompts in LLM Chatbots

**Advisor:** [Martin Kropp](#)  
[Nitish Patkar](#)

**Work scope:** Priority 1  
P5 oder P6 ---  
**Team size:** Priority 2  
2er Team ---

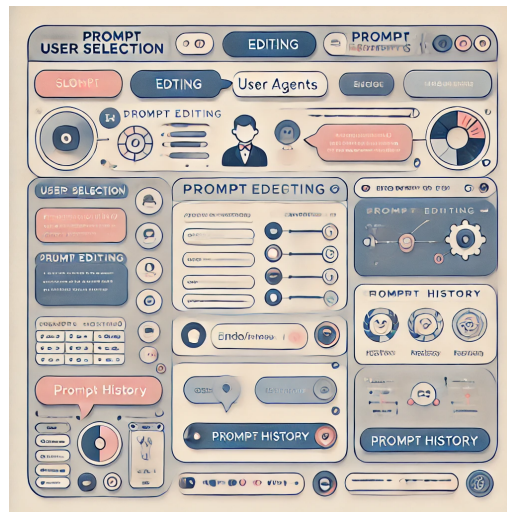
**Languages:** German or English  
**Study course:** Computer Science

### Initial position

AI chatbots are widely used tools today for various routine purposes. Research shows that the UI and interaction design of most chatbots lack essential features, such as search functionality within history or prompt autocompletion. Additionally, each prompt contributes to an environmental toll, emphasizing the need for efficient prompt management. Currently, most chatbots handle past prompts in a similar, limited way, offering few options to revisit and branch from specific past prompts for conversations in new directions.

### Objective

The overarching objective of this project is to advance the concept and implementation of our proof-of-concept AI chatbot. We aim to rethink how to make AI chatbots more user-friendly and improve the overall user experience. Specifically, within IP5, we want to explore interaction mechanisms that enable users to work effectively with past prompts, reducing the need for repeated re-prompting. This involves experimenting with visualization alternatives, addressing design and implementation challenges, and validating the effectiveness of the developed solutions. For IP6, students can also explore prompt speculation- speculating next prompt based on previous interaction.



Better ways for prompt management

### Problem statement

In today's prompt management for LLMs, UI/UX faces several challenges. Usability suffers due to complex prompt tuning and the lack of intuitive interfaces for users to refine prompts without technical expertise. Traceability is limited, as it's often unclear how specific prompts influence responses, complicating debugging and optimization efforts. Additionally, there is a lack of historicization tools to track prompt versions, changes over time, and their impacts on outputs, which hinders iterative improvements and makes it difficult to maintain consistent performance across updates. Together, these issues make prompt management cumbersome and less accessible for users.

The following questions should be answered:

- What features do current popular AI chatbots support for interaction and prompt management?
- What visualization and interaction mechanisms can enhance interaction with past prompts?

### Technologies/Technical emphasis/References

Web development, preferably Python, Angular, SQLite, to be discussed.

### Note

There has been earlier work on a this [https://nitishspatkar.github.io/pdfs/IP5\\_FS24\\_Simon\\_Jack.pdf](https://nitishspatkar.github.io/pdfs/IP5_FS24_Simon_Jack.pdf) which can be used as inspiration by the students.