**IacuWise**

# Measuring the Environmental Footprint of AI Prompts

## A Research-Backed Methodology for Quantifying Water, Energy, and Carbon Savings from Prompt Optimization

| | |
|---|---|
| **Version** | 2.0 |
| **Date** | February 2026 |
| **Classification** | Public |
| **Standard Alignment** | CSRD / GRI 303 / CDP Water Security / ISO 14046 |

*IacuWise | iacuwise.vercel.app | Sustainable AI Platform*

# Table of Contents

# 1. Executive Summary

Artificial intelligence is transforming every industry, but its environmental cost remains largely invisible to end users. Every AI prompt consumes electricity for computation, water for datacenter cooling, and generates carbon emissions through grid-dependent energy production. As global AI electricity consumption is projected to reach 945 TWh by 2030 (IEA, 2025), and the water footprint of AI systems could reach 312-764 billion liters annually (de Vries, 2025), the need for transparency and reduction tools has never been greater.

IacuWise addresses this challenge through **prompt optimization** — the practice of improving AI prompts before they are sent to a model, reducing the number of retry attempts needed to achieve a satisfactory result. A well-crafted prompt eliminates an average of 1.4 retry cycles, cutting total token consumption by approximately 50-60% per interaction.

This whitepaper presents IacuWise's peer-reviewed methodology for quantifying the environmental savings of prompt optimization. The model draws on 15+ academic and institutional sources, including research from the University of California Riverside, the International Energy Agency, the U.S. EPA, and Microsoft's sustainability disclosures. All calculations are designed to align with CSRD, GRI 303, CDP Water Security, and ISO 14046 reporting standards.

| Key Finding | Value | Source |
|---|---|---|
| Avg. retries eliminated per optimization | 1.4 attempts | UX research / industry consensus |
| Token reduction per interaction | 50-60% | IacuWise engine calculations |
| Water saved per optimization (Claude) | ~1.7 mL | Li et al. (2025), EESI (2024) |
| Energy saved per optimization (Claude) | ~0.5 Wh | You (2025), Lin (2025) |
| CO2 avoided per optimization (Claude) | ~188 mg | EPA eGRID 2022 |

*Table 1: Key findings for a typical single-prompt optimization using Claude (Anthropic).*

# 2. The Problem: AI's Hidden Water Footprint

The environmental impact of AI systems has come under increasing scrutiny, yet most attention has focused on carbon emissions from model training. The water footprint of AI inference — the process of actually using AI models — remains largely invisible to users and organizations.

## 2.1 Scale of the Problem

According to the International Energy Agency's landmark 2025 report on Energy and AI, data centers consumed approximately 415 TWh of electricity in 2024 — about 1.5% of global consumption. This figure is projected to more than double to 945 TWh by 2030. AI-specific workloads represented 15-20% of data center electricity in 2024, projected to reach 35-50% by 2030.

Research from the University of California Riverside (Li et al., 2025) introduced a comprehensive framework for measuring AI's water footprint across two scopes: direct cooling water at data centers (Scope 1) and water consumed in electricity generation at power plants (Scope 2). Their findings show that global AI demand could account for 4.2-6.6 billion cubic meters of water withdrawal by 2027 — more than the total annual water withdrawal of 4-6 Denmarks.

According to researchers at UC Riverside, each 100-word AI prompt is estimated to use roughly one bottle of water (approximately 519 milliliters) when accounting for both direct and indirect water consumption (EESI, 2024). With billions of AI prompts processed daily, the cumulative impact is significant.

## 2.2 The Retry Problem

A critical but underappreciated factor in AI's environmental footprint is the **retry cycle**. When users submit poorly structured prompts, they typically need 2-3 attempts to get a satisfactory result. Each attempt consumes the full computational pipeline: token processing, GPU inference, datacenter cooling, and electricity generation. A single vague prompt that requires 2.5 attempts consumes 2.5 times the resources of a well-crafted prompt that succeeds on the first try.

IacuWise addresses this multiplier effect directly. By optimizing prompts before submission, the platform reduces retry rates from an average of 2.5 to 1.1 attempts, eliminating approximately 56% of unnecessary token consumption and its associated environmental impact.

# 3. Methodology Overview

IacuWise's environmental impact engine follows a three-step conversion chain that transforms token-level data into measurable environmental metrics:

| Step | Input | Conversion | Output |
|------|-------|-----------|--------|
| 1. Energy | Tokens (count) | Energy/token x PUE | Energy (Wh) |
| 2. Water | Energy (Wh) | WUE + EWIF | Water (mL) |
| 3. Carbon | Energy (Wh) | Grid $CO_2$ intensity | $CO_2$ (g) |

*Table 2: Impact conversion chain from tokens to environmental metrics.*

## 3.1 Water Footprint Scoping (ISO 14046 Aligned)

Following the framework established by Li et al. (2025) and aligned with ISO 14046 water footprint standards, IacuWise measures water impact across two scopes:

**Scope 1 — Direct Cooling (On-site):** Water evaporated in datacenter cooling towers and liquid cooling systems. Measured using the Water Usage Effectiveness (WUE) metric developed by The Green Grid, expressed in liters per kilowatt-hour (L/kWh). Industry average WUE is 1.8-1.9 L/kWh (EESI, 2024), though hyperscalers like AWS (0.19 L/kWh) and Microsoft (0.30 L/kWh) perform significantly better.

**Scope 2 — Electricity Generation (Off-site):** Water consumed at power plants that supply electricity to data centers. Measured using the Energy-Water Interaction Factor (EWIF) developed by the World Resources Institute (Reig et al., 2020). The U.S. average EWIF is 3.14 L/kWh, reflecting the substantial water requirements of thermoelectric power generation.

Research from Bluefield Research (2025) confirms the importance of Scope 2 accounting: indirect water consumption at power plants is approximately ten times that of on-site datacenter water consumption, making it the dominant component of AI's total water footprint.

# 4. Impact Chain: Tokens to Environmental Cost

## 4.1 Token Estimation

IacuWise estimates token counts using a character-based heuristic calibrated for mixed English/Spanish content, with an average ratio of approximately 3.7 characters per token. This provides a reliable approximation for environmental calculations without requiring provider-specific tokenizers.

## 4.2 Energy per Token

Energy consumption per token varies significantly by model architecture, hardware, and optimization techniques. IacuWise uses conservative estimates derived from published research:

You (2025) at Epoch AI established that a typical GPT-4o query consuming approximately 500 output tokens uses 0.2-0.3 Wh, yielding roughly 0.0006 Wh/token. Lin (2025) demonstrated that Llama3-70B on H100 hardware with FP8 quantization achieves 0.39 J/token (~0.0001 Wh/token), representing a significant efficiency improvement over earlier architectures. Jegham et al. (2025) provided comprehensive per-query energy benchmarks across model sizes from Nano to Large.

## 4.3 Power Usage Effectiveness (PUE)

PUE measures total facility energy divided by IT equipment energy. A PUE of 1.0 would mean all energy goes to computation; real-world values include cooling, lighting, and infrastructure overhead. Major cloud providers report PUE values of 1.10-1.12 (Microsoft, 2024; Google, 2024). Less efficient facilities, particularly in developing markets, may have PUE values of 1.3-1.5. Cooling systems account for 7% of energy in efficient hyperscalers, rising to over 30% in less efficient enterprise facilities (IEA, 2025).

## 4.4 Carbon Intensity

$CO_2$ emissions per kWh of electricity vary by grid energy mix. The U.S. national average is 373 g $CO_2$/kWh (EPA eGRID 2022). Google's datacenters benefit from higher renewable penetration (estimated ~280 g $CO_2$/kWh), while Chinese datacenters hosting models like DeepSeek face higher carbon intensity (~550 g $CO_2$/kWh) due to greater coal dependence.

# 5. Provider-Specific Parameters

IacuWise maintains calibrated environmental profiles for four major AI providers, using provider-specific sustainability disclosures where available and conservative industry averages where not:

| Parameter | Claude (Anthropic) | GPT-4o (OpenAI) | Gemini (Google) | DeepSeek |
|---|---|---|---|---|
| Infrastructure | AWS + GCP | Microsoft Azure | Google Cloud | China DCs |
| PUE | 1.10 | 1.12 | 1.10 | 1.40 |
| WUE Scope 1 (L/kWh) | 0.20 | 0.30 | 0.20 | 1.80 |
| EWIF Scope 2 (L/kWh) | 3.14 | 3.14 | 3.14 | 2.50 |
| CO2 intensity (g/kWh) | 373 | 373 | 280 | 550 |
| Energy/token (Wh) | 0.0004 | 0.0006 | 0.0005 | 0.0003 |
| Primary sources | AWS/GCP avg [7][8] | Microsoft FY2024 [7] | Google public data | China DC averages [6] |

*Table 3: Provider environmental profiles. Numbers in brackets refer to sources in References. Where provider-specific data is unavailable, conservative industry averages are used.*

**Notable observations:** DeepSeek uses a Mixture-of-Experts (MoE) architecture that activates only ~21B parameters from a 236B total, achieving very low per-token energy. However, its Chinese datacenter infrastructure has significantly higher PUE, WUE, and carbon intensity, resulting in a higher environmental footprint per watt-hour consumed.

# 6. The Retry Reduction Model

The retry reduction model is the core value proposition of IacuWise's environmental savings. Poorly crafted prompts generate a cascade of wasted resources: each retry consumes the full computational pipeline, multiplying the environmental impact of a single interaction.

| Parameter | Value | Rationale |
|---|---|---|
| Unoptimized attempts | 2.5x | Conservative estimate based on UX research showing users iterate 2-4 times with unclear prompts |
| Optimized attempts | 1.1x | Well-structured prompts with clear role, format, and constraints reduce retries to ~10% probability |
| Avg response tokens | 800 | Typical response length for general-purpose queries |
| Response efficiency | 80% | Optimized prompts yield 20% shorter, more focused responses due to precision |

*Table 4: Retry model parameters and rationale.*

The retry reduction mechanism creates a multiplier effect: reducing attempts from 2.5 to 1.1 eliminates 56% of total token consumption. Combined with the 20% reduction in response length (more precise prompts yield more focused answers), the total environmental savings per optimization typically range from 50-60%.

# 7. Calculation Formulas

## 7.1 Scenario Comparison

IacuWise compares two scenarios for every optimization:

```
Scenario A (without optimization): total_tokens_A = (original_prompt +
avg_response) x unoptimized_attempts

Scenario B (with IacuWise): total_tokens_B = (optimized_prompt +
focused_response) x optimized_attempts
```

## 7.2 Energy Calculation

```
energy_Wh = tokens x energy_per_token x PUE
```

Where energy_per_token is the GPU energy consumption per token for the selected provider, and PUE accounts for datacenter facility overhead (cooling, networking, UPS).

## 7.3 Water Calculation

```
water_scope1_mL = WUE x (energy_Wh / 1000) x 1000

water_scope2_mL = EWIF x (energy_Wh / 1000) x 1000

water_total_mL = water_scope1_mL + water_scope2_mL
```

WUE (Water Usage Effectiveness) measures direct cooling water consumption. EWIF (Energy-Water Interaction Factor) measures water consumed at electricity generation facilities. Both are expressed in L/kWh and converted to mL for per-query granularity.

## 7.4 Carbon Calculation

```
co2_g = grid_co2_intensity x (energy_Wh / 1000)
```

Grid $CO_2$ intensity (g $CO_2$/kWh) varies by provider location and energy mix.

## 7.5 Savings Calculation

```
savings = metric_scenario_A - metric_scenario_B

reduction_pct = (tokens_saved / total_tokens_A) x 100
```

# 8. Cumulative Impact and ESG Reporting

IacuWise aggregates individual optimization results to produce cumulative environmental reports suitable for ESG disclosure. The platform stores every optimization in a PostgreSQL database with full traceability: original prompt, optimized prompt, model used, token counts, and calculated environmental impact metrics.

## 8.1 Real-World Equivalencies

To make cumulative savings tangible, IacuWise converts raw metrics into everyday equivalencies:

| Metric | Equivalency | Conversion Factor |
|--------|-------------|-------------------|
| Water saved | Standard water bottles | 500 mL per bottle |
| Energy saved | LED bulb operating hours | 10 Wh per hour (10W bulb) |
| CO2 avoided | Car kilometers avoided | 120 g CO2/km average |

*Table 5: Real-world equivalency conversions for ESG reporting.*

## 8.2 Reporting Standards Alignment

IacuWise's metrics are designed to support disclosure under multiple frameworks: **CSRD** (Corporate Sustainability Reporting Directive) requires EU companies to report on environmental impact including water and energy; **GRI 303** (Water and Effluents) provides standards for water withdrawal, consumption, and discharge reporting; **CDP Water Security** questionnaire assesses corporate water risk and management; and **ISO 14046** provides the international standard for water footprint assessment.

# 9. Limitations and Disclaimers

IacuWise's environmental estimates are approximations based on the best publicly available research. Users and organizations should be aware of the following limitations:

• **Datacenter variability:** Actual energy, water, and carbon values vary significantly by datacenter location, time of day, season, cooling technology, server utilization, and local grid energy mix. IacuWise uses provider-level averages that may not reflect any specific facility.

• **Provider opacity:** Most AI providers do not disclose per-query environmental metrics. IacuWise derives its estimates from published infrastructure data and academic research, not from direct measurement of individual API calls. As noted by de Vries (2025), the lack of AI-specific environmental reporting from datacenter operators remains a significant challenge.

• **Retry model assumptions:** The 2.5x average retry rate is a conservative estimate. Actual retry behavior varies by user experience, prompt complexity, and use case. Some interactions may require more retries; some may succeed on the first attempt.

• **Token estimation:** Character-to-token ratios are approximate and vary by language, content type, and provider-specific tokenizers. IacuWise uses a 3.7 characters/token average.

• **Not suitable for formal carbon accounting:** These calculations are designed for awareness, estimation, and ESG directional reporting. They should not replace formal life-cycle assessments or provider-specific sustainability audits. For formal ESG reporting, verify against provider sustainability disclosures.

# 10. References

[1] Li, P., Yang, J., Islam, M.A., & Ren, S. (2025). "Making AI Less Thirsty: Uncovering and Addressing the Secret Water Footprint of AI Models." *Communications of the ACM*, 68(7), 54-61. University of California, Riverside. https://dl.acm.org/doi/10.1145/3724499

[2] Jegham, N. et al. (2025). "How Hungry is AI? Benchmarking Energy, Water, and Carbon Footprint of LLM Inference." arXiv:2505.09598v1. https://arxiv.org/abs/2505.09598

[3] You, J. (2025). "How much energy does ChatGPT use?" *Gradient Updates*, Epoch AI. Per-query energy estimates for GPT-4o.

[4] Lin, L.H. (2025). "Llama3-70B Inference Efficiency on H100." 0.39 J/token on 8xH100 with vLLM + FP8 quantization.

[5] Samsi, S. et al. (2023). "From Words to Watts: Benchmarking the Energy Costs of Large Language Model Inference." MIT Lincoln Laboratory.

[6] EESI (2024). "Data Centers and Water Consumption." Environmental and Energy Study Institute. https://www.eesi.org/articles/view/data-centers-and-water-consumption

[7] Microsoft (2024). "Measuring energy and water efficiency." Global PUE: 1.12, WUE: 0.30 L/kWh (FY2024). https://datacenters.microsoft.com/sustainability/efficiency/

[8] AWS (2024). Sustainability Report. Global WUE: 0.19 L/kWh.

[9] U.S. EPA (2024). eGRID 2022 Data. U.S. national average: 373.3 g CO2/kWh. https://www.epa.gov/egrid

[10] U.S. EIA (2024). "How much CO2 is produced per kWh?" ~0.81 lbs CO2/kWh (367 g/kWh) for 2023.

[11] Reig, P. et al. (2020). "Guidance for calculating water use embedded in purchased electricity." World Resources Institute. U.S. average EWIF: 3.14 L/kWh.

[12] de Vries, A. (2025). "The carbon and water footprints of data centers and what this could mean for artificial intelligence." *Patterns*, Cell Press. VU Amsterdam. https://www.cell.com/patterns/fulltext/S2666-3899(25)00278-8

[13] International Energy Agency (2025). "Energy and AI." IEA Special Report. Data centers consumed 415 TWh in 2024, projected 945 TWh by 2030. https://www.iea.org/reports/energy-and-ai

[14] Harvard Political Review (2025). "When the People's Water Vanishes." Analysis of data center water competition with communities. https://theharvardpoliticalreview.com/ai-water-consumption/

[15] Han, Y., Wu, Z., Li, P., Wierman, A., & Ren, S. (2024). "The Unpaid Toll: Quantifying the Public Health Impact of AI." UC Riverside / Caltech. arXiv:2412.06288.

[16] Hajiesmaili, M., Ren, S., Sitaraman, R., & Wierman, A. (2025). "Towards Environmentally Equitable AI." *Communications of the ACM*. UC Riverside / UMass / Caltech.

[17] Bluefield Research (2025). "Data Center Water Secrecy Hurts Communities." U.S. data centers withdraw 107 MGD directly; indirect consumption at power plants is 10x higher.

[18] Pew Research Center (2025). "What we know about energy use at U.S. data centers amid the AI boom." U.S. data centers consumed 183 TWh in 2024. https://www.pewresearch.org/

[19] Altman, S. (2025). Personal blog post. Average ChatGPT query uses ~0.3 mL of water (scope-1 only).

[20] Google (2024). Environmental Report. 78% of data center water withdrawal was potable water (7,700M gallons).

---