

GreenMind: Smarter AI, Lower Carbon

The Problem: AI's Growing Energy Footprint

The energy consumption of **Large Language Models (LLMs)** is **skyrocketing**, with inference alone accounting for **65% of total energy use**—far exceeding training costs.

By **2030**, global data centers' electricity consumption is expected to **double** due to the rising demands of Generative AI.

User Behavior Insights:

- **70%** of users are aware of AI's **environmental impact**, and **80%** are willing to adopt **sustainable solutions**.
- **35%** don't manually select models and **20%** always default to the latest (and most energy-consuming) versions, often unnecessarily.

The Consequence

Inefficient model selection **wastes power, increases costs, and drives up CO₂ emissions**—all while many tasks could be handled by smaller, more efficient models.

The GreenMind Solution: Adaptive AI Model Selection

GreenMind is an **intelligent routing framework** that **classifies prompt complexity** and dynamically selects the most **energy-efficient LLM** without sacrificing performance.

How it Works:

✓ **Custom Fine-Tuned BERT Classifier** → **Analyzes prompt complexity** based on reasoning, context depth, and creativity.
✓ **Adaptive Chat Framework** → **Reroutes each prompt** to the optimal model based on complexity:

- **Low Complexity** → **LLaMA 3.1 - 8B** (Fastest, Lowest Energy Use)
- **Medium Complexity** → **LLaMA 3.1 - 70B** (Balanced Power & Accuracy)
- **High Complexity** → **LLaMA 3.1 - 405B** (Max Performance, Highest Cost)
 - ✓ **AWS Bedrock Integration** → Ensures **scalable & seamless deployment**.

By dynamically selecting the **right-sized model for the job**, GreenMind **cuts AI energy waste and saves costs** while maintaining high-quality responses.

Impact & Key Benefits

- ♦ **Up to 5× Lower CO₂ Emissions** – Reduces energy waste by running lighter models for simpler tasks.
- ♦ **Lower AI Compute Costs** – Avoids unnecessary use of high-power models, saving enterprises thousands.
- ♦ **No User Friction** – Works **automatically** with an override option for advanced users.
- ♦ **Minimal Latency** – Faster response times by skipping heavier models when unnecessary.

Real-World Example

A **medium-complexity prompt**, when classified correctly, can run on **LLaMA 70B instead of 405B**, saving:

- ✓ **2 kWh per 1,000 tokens**
 - ✓ **\$0.04 per 1,000 tokens**
 - ✓ **900g CO₂e per 1,000 tokens** (equivalent to 3 miles driven in a gas car)
-

UI & User Experience

- 1 **Enter Prompt** – Users input queries normally.
 - 2 **Automated Model Selection** – GreenMind **classifies & routes** the request.
 - 3 **Live Carbon Savings Indicator** – Users see **real-time impact** when a smaller model is selected
-

Next Steps & Roadmap

- ♦ **Expand BERT Classifier Dataset** – Improve **classification accuracy** to optimize routing.
 - ♦ **Optimize Model Routing** – Enhance **latency and efficiency** of model selection.
 - ♦ **Deploy Interactive Dashboard** – Provide **personalized CO₂ savings insights** for users.
 - ♦ **Launch GreenMind as a Browser Extension & Mobile App** – Make it easy to use across platforms.
-

Join GreenMind: Smarter AI, Lower Carbon.

-  **Reduce AI energy waste today!**
 -  **Let's build a more sustainable AI ecosystem together.**
-