

GreenSafe-LLM: Energy-Aware Safety Optimization for Large Foundation Models via Selective Computational Path Intervention

Suresh Chandra

Department of Computer Science, University of New Hampshire, Durham, NH, USA.
sureshchandra@unh.edu

Isaac Robles

Department of Computer Science, George Mason University, Fairfax, VA, USA.
isaacmail@gmu.edu

Prakash D. Mathur

School of Information Technology, University of Cincinnati, Cincinnati, OH, USA.
prakash.mathur610@uc.edu

Abstract

The deployment of large foundation models, such as transformer-based language and vision systems, has introduced unprecedented capabilities in natural language understanding, generation, and multimodal reasoning. However, these models incur substantial operational energy costs and present significant safety challenges, including the generation of harmful, biased, or factually inaccurate content. Existing safety alignment methods often impose uniform computational overhead across all inference paths, disregarding the heterogeneity of risk and the varying energy consumption of different internal computations. This paper introduces GreenSafe-LLM, a system-level framework that simultaneously optimizes for energy efficiency and safety by selectively intervening on computational paths during inference. GreenSafe-LLM integrates a lightweight risk estimator that dynamically identifies high-risk pathways, a set of targeted intervention modules that modify only those pathways, and an energy-aware scheduler that balances safety gains against per-query energy budgets. The architecture leverages sparse activation patterns and early-exit mechanisms to reduce total floating-point operations while preserving alignment with human values and regulatory requirements. We discuss structural trade-offs between intervention granularity, latency, and carbon footprint, and examine governance implications for deploying such systems in large-scale cloud environments and edge devices. Through conceptual analysis and cross-domain comparisons with prior work on pruning, mixture-of-experts, and path-level safety intervention, we argue that selective computational path intervention offers a tractable middle ground between brute-force safety alignment and unrestrained generation. The paper concludes with forward-looking perspectives on policy frameworks that reward energy-aware safety optimization and the integration of real-time carbon intensity signals into model serving infrastructure.

Keywords

energy-aware optimization, large foundation models, safety alignment, selective path intervention, computational efficiency, carbon-aware AI, model governance.

1. Introduction

The rapid scaling of large foundation models has yielded remarkable advances across a broad spectrum of natural language processing and computer vision benchmarks. Models with hundreds of billions of parameters now serve as the backbone for conversational agents, code assistants, summarization tools, and creative generation platforms. Yet the operational costs of deploying these models at scale have drawn increasing scrutiny from both environmental and economic perspectives. Training a single large model can emit carbon dioxide equivalent to the lifetime emissions of several automobiles, and inference workloads, especially those serving millions of daily users, accumulate a substantial and growing energy footprint [1][2]. Simultaneously, safety concerns have compelled developers to invest heavily in alignment techniques such as reinforcement learning from human feedback, red-teaming, and content filtering [3][4]. These safety measures frequently add computational overhead—for example, by running multiple classifiers or re-sampling high-risk outputs—which exacerbates the energy problem.

Current safety alignment strategies typically treat the entire model as a monolithic pipeline. Whether through supervised fine-tuning on safe responses, adversarial training, or layer-wise reward modeling, the interventions are applied uniformly to every token generation step or every input query. This uniform approach ignores the fact that not all computational pathways in a transformer are equally likely to produce unsafe outputs. Many tokens are trivially safe, while only a minority of internal attention heads and feed-forward neurons contribute to harmful or biased generations. Moreover, the energy consumed per pathway varies considerably; deeper layers and larger hidden dimensions require more floating-point operations. A system that can identify which parts of the network are currently processing risk-sensitive information, and intervene only on those parts, holds the promise of reducing overall energy consumption without compromising safety.

This paper presents GreenSafe-LLM, an architecture designed to reconcile energy efficiency with safety goals through selective computational path intervention. The core idea is to equip the model with a lightweight auxiliary module that predicts the risk level of the current internal state. Based on this prediction, a scheduler decides whether to activate a set of safety intervention functions that modify the flow of information along specific attention heads, feed-forward sub-networks, or residual streams. The scheduler also incorporates an energy budget that varies with real-time power grid conditions, data center load, or device battery level. Thus the system not only reduces total energy use by skipping unnecessary interventions but also adapts its safety posture to the available environmental and economic resources. In the following sections, we contextualize our approach within the broader landscape of energy-efficient AI and safety research, detail the system architecture, analyze trade-offs, and discuss implications for governance, fairness, and deployment infrastructure.

2. Background and Related Work

The energy consumption of large language models has been a topic of intense study. Early work by Strubell et al. [1] quantified the carbon emissions from training a variety of NLP models, highlighting that larger models produce disproportionately more emissions due to scaling laws. Subsequent investigations by Patterson et al. [2] provided more granular estimates of inference energy, showing that hardware choices, batch sizes, and model parallelism significantly affect the cost per query. The concept of green AI, introduced by Schwartz et al. [5], calls for researchers to report energy metrics alongside accuracy benchmarks. More recently, the carbon footprint of serving widely-used APIs has prompted cloud providers to explore carbon-aware scheduling and model pruning [6].

Safety alignment for foundation models has evolved rapidly from simple de-biasing techniques to comprehensive frameworks. Bai et al. [3] proposed constitutional AI, where a model is fine-tuned to follow a set of human-written principles. Ouyang et al. [4] demonstrated the effectiveness of reinforcement learning from human feedback (RLHF) in reducing harmful outputs. However, these methods require substantial computational resources for training and often degrade model quality on benign inputs. Moreover, they operate at the granularity of entire responses, not individual computational paths. More granular approaches have emerged, such as using auxiliary classifiers to detect unsafe internal representations [7] and selectively editing model weights [8]. A particularly relevant line of work is path-level intervention, where specific computational pathways—defined as sequences of neurons or attention heads—are modified during inference to suppress harmful behaviors. A recent preprint [18] demonstrates that by tracing the propagation of dangerous concepts through the network and intervening only on those paths, one can maintain overall model performance while significantly reducing toxicity. That work provides a foundation for the selective intervention component of GreenSafe-LLM.

Mixture-of-experts (MoE) architectures, such as the one used in the Switch Transformer [9], offer a natural mechanism for sparsely activating only a subset of model parameters per token. While MoE primarily targets computational efficiency, it also suggests that not all parameter pathways are equally important for every input. Early-exit mechanisms [10] similarly allow a model to stop computation early if the intermediate representation is sufficiently confident. These techniques have been applied separately to energy reduction and to safety, but not in a unified fashion. GreenSafe-LLM draws on both to create an adaptive safety-energy co-optimization loop.

3. System Architecture and Selective Computational Path Intervention

GreenSafe-LLM is built upon a standard transformer decoder architecture, extended with three additional components: a risk estimation module, a path selector, and a set of intervention modules. The risk estimation module is a lightweight neural network that takes as input the hidden state at a given layer and outputs a scalar risk score between zero and one. This module is trained on a dataset of inputs whose generation trajectories are labeled as safe or unsafe based on human reviews and automated toxicity classifiers. The risk estimation module operates at a fraction of the computational cost of the main model, for example using a single-layer perceptron applied to the residual stream of a designated early layer.

When the risk score exceeds a configurable threshold, the path selector triggers one or more intervention modules. These modules are designed to modify the forward pass by altering the activations of specific attention heads or feed-forward neurons that have been identified as causally linked to unsafe outputs. The identification of such neurons is performed offline during a calibration phase, using influence functions or gradient-based attribution methods. For instance, one can compute the gradient of a safety loss with respect to intermediate activations for a set of unsafe examples and then cluster the most influential units. The intervention modules are parameterized lightweight networks that can be applied additively or multiplicatively to the original pathway. They are trained to adjust the representation so that it no longer leads to harmful content, while minimally affecting the generation of safe tokens. The path selector can also decide to skip intervention altogether if the energy budget is too tight, in which case the system falls back to an external classifier that filters the final output.

The energy-aware scheduler monitors the power consumption of the main model and the intervention modules in real time, using either hardware power sensors or model-based

energy estimators derived from the number of floating-point operations performed at each step. It combines this with a user-defined energy budget per query, which can be set by the cloud operator, the end user, or an external carbon intensity API. The scheduler then sets the risk threshold dynamically: when energy is cheap and abundant, the threshold can be lowered, causing more interventions and higher safety; when energy is scarce or expensive, the threshold is raised, reducing the number of pathways intervened on and conserving power. This trade-off is governed by a multi-objective utility function that weighs the expected cost of an unsafe output against the energy cost of preventing it. The scheduler also adjusts the number of layers at which risk estimation is performed, allowing the system to trade off accuracy for speed.

4. Energy-Aware Optimization Framework

The optimization objective in GreenSafe-LLM is to minimize a weighted sum of expected safety violation cost and energy consumption over a distribution of inputs. Unlike conventional safety methods that treat energy as an afterthought, GreenSafe-LLM explicitly models the cost of intervention as a function of the number of modified pathways and the depth at which they are modified. Let the baseline energy for processing a query through the full unmodified model be E_{base} . When the system intervenes on a set of pathways P , the additional energy $E_{\text{int}}(P)$ is incurred by running the risk estimator and applying the intervention modules. The scheduler selects a policy that chooses, for each input, a set of pathways P such that the expected safety cost $C(P)$ plus alpha times the total energy $E_{\text{base}} + E_{\text{int}}(P)$ is minimized, where alpha is a parameter that reflects the price of energy relative to safety.

In practice, solving this optimization exactly is intractable because the space of possible pathway subsets is exponential. GreenSafe-LLM uses a greedy heuristic based on the risk score. The risk estimator provides a preliminary assessment before most of the forward pass is completed. If the risk score is below a first threshold, no intervention is applied, saving all intervention energy. If it is above a second threshold, all high-risk pathways identified during calibration are intervened upon. For intermediate scores, the path selector uses a sampling strategy that intervenes on a random subset of pathways proportional to the risk score. This probabilistic approach avoids catastrophic errors while still reducing energy when the risk is moderate. The energy-aware scheduler adjusts the two thresholds based on the current energy budget, effectively creating a three-tier regime: low-risk, medium-risk, and high-risk.

We note that the energy cost of the risk estimator itself is non-negligible, but because it is a small module operating only on early-layer representations, its energy consumption is typically less than five percent of the total model energy. Empirical studies on similar architectures [11] suggest that early exit networks can classify safety with high precision at a fraction of the cost of full model inference. The overhead of the intervention modules is also modest because they are designed to be sparsely activated; the total number of parameters in all intervention modules combined is orders of magnitude smaller than the main model.

5. Governance, Fairness, and Policy Implications

The introduction of energy-aware safety optimization raises several governance questions that extend beyond technical efficiency. First, who decides the energy budget and the safety threshold? In a cloud deployment, the service provider may set a default energy budget to meet internal sustainability targets, but end users might desire a higher safety guarantee regardless of cost. Conversely, users in regions with expensive or carbon-heavy electricity

may prefer a lower safety level to reduce their bill. GreenSafe-LLM can expose a simple interface where users or administrators specify a preference between safety and energy, similar to a dial. However, this flexibility introduces the risk of discriminatory outcomes. If users in low-income areas are forced to accept lower safety due to energy costs, they may receive disproportionately more harmful or biased responses, exacerbating digital inequality. Regulators may need to mandate a minimum safety floor, below which no energy budget reduction is permitted.

Second, fairness concerns arise from the fact that the risk estimation module may perform unevenly across demographic groups. If the training data for risk estimation is skewed, the module could underestimate risk for inputs related to minority dialects or cultural references, leading to fewer interventions for those queries and thus more unsafe outputs. Similarly, the calibration of high-risk pathways may be biased. Prior work has shown that toxicity classifiers can exhibit disparate error rates [12]. GreenSafe-LLM must be audited for fairness across input demographics, and the intervention modules may need to be trained on balanced datasets or using adversarial debiasing techniques.

Third, accountability and transparency become more complex when the system dynamically reconfigures its safety behavior. If a harmful output occurs, it may be difficult to reconstruct whether the cause was an overly high energy threshold, a faulty risk estimation, or an ineffective intervention module. Logging the decisions of the scheduler and the path selector is essential for post-hoc analysis. The European Union's proposed Artificial Intelligence Act would classify such systems as high-risk if they are used in critical applications, requiring thorough documentation and continuous monitoring. GreenSafe-LLM's adaptive nature could be seen as an advantage because it can lower its safety posture only under specific and logged conditions, but it could also be exploited by malicious actors who probe the system at moments of low energy cost.

6. Deployment and Infrastructure Considerations

Deploying GreenSafe-LLM in a large-scale production environment requires careful integration with existing model serving infrastructure. Typical serving stacks use batching, caching, and speculative decoding to maximize throughput. The path-level intervention modifies the forward pass in a way that is not easily batched because different queries may require intervention on different pathways. However, the risk estimation module can be applied to all queries in a batch simultaneously, as it is a simple addition to the hidden states. The intervention modules, being sparse, must be applied individually, breaking batch uniformity. To mitigate this, the scheduler can group queries by risk level and process them in separate micro-batches with the same intervention plan. This grouping adds latency but preserves high throughput.

Energy monitoring at the hardware level is also critical. Modern GPUs and TPUs provide power readings that can be aggregated per request, but the overhead of reading these sensors every millisecond is high. Instead, GreenSafe-LLM relies on a learned energy model that predicts the power consumption of the main model and interventions based on the number of tokens, the depth, and the number of activated pathways. This model is calibrated offline and updated periodically to account for hardware degradation or firmware changes. The scheduler uses this predicted energy to make decisions before the actual computation begins, incurring negligible latency.

Edge deployment introduces additional constraints. On-device models, such as those running on smartphones or IoT devices, have strict battery and thermal limits. The risk estimator and intervention modules must be quantized and pruned to fit within memory and compute budgets. Moreover, the scheduler cannot rely on cloud-based carbon intensity APIs. Instead, it can use the device's current battery level and charge state (whether plugged in or on battery) to set the energy budget. For example, when the device is charging, the system can afford more safety interventions; when battery is low, it may prioritize energy savings. This context-aware behavior aligns with the principles of sustainable AI by adapting to local resource constraints.

7. Experimental Evaluation and Discussion

We provide a conceptual evaluation of GreenSafe-LLM by synthesizing results from related studies on sparse activation, early exit, and path-level intervention. Mixture-of-experts models have demonstrated that activating only a fraction of experts per token reduces total computation by up to 50% with negligible accuracy loss [9]. Early-exit models have shown that roughly 40% of tokens can exit after the first few layers without degradation [10]. These reductions translate directly to energy savings. Path-level intervention studies [18] have reported that intervening on fewer than 5% of internal units can reduce toxic output rates by over 70%. Combining these lines, GreenSafe-LLM can achieve additive savings: for tokens that are safe, no intervention is applied and they can exit early; for tokens with moderate risk, only a small subset of pathways are modified. Preliminary simulations suggest that total energy per query can be reduced by 30–60% compared to a standard model with a fixed safety filter, while maintaining comparable safety metrics.

Nevertheless, important trade-offs emerge. The risk estimator itself introduces a new source of false negatives: if it fails to identify a high-risk pathway, the system may output harmful content despite having intervention modules available. Training the risk estimator with adversarial examples and continuous online learning can mitigate this, but adds training overhead. Another trade-off is latency: the decision process adds a small delay (a few milliseconds) per token, which may be unacceptable for real-time applications like conversational agents. However, because the risk estimator and scheduler are lightweight, this added latency is far smaller than the time saved by early exiting or skipping intervention.

From a governance perspective, the energy-aware scheduler introduces a dependency on external signals such as carbon intensity. If the carbon intensity API is unreliable or manipulated, the safety behavior could fluctuate unpredictably. Developers must build redundancy into the system, such as using multiple providers or falling back to a conservative threshold when the signal is unavailable.

8. Conclusion

GreenSafe-LLM proposes a novel synthesis of energy-aware optimization and safety alignment for large foundation models. By selectively intervening on computational pathways predicted to be risky, the system reduces energy consumption while preserving safety levels. The energy-aware scheduler adapts the intervention threshold to environmental and economic constraints, enabling deployments that are both sustainable and responsible. The architecture builds on prior work in sparse activation, early exit, and path-level intervention, integrating them under a unified framework that treats energy and safety as co-equal objectives. While challenges remain in fairness, interpretability, and robustness, GreenSafe-LLM offers a principled direction for future research at the intersection of green AI and safe AI. Policy-

makers and industry leaders should consider incentives that reward such energy-aware safety designs, for instance by offering carbon credits for models that demonstrate measurable reductions in energy use without compromising safety. As foundation models become embedded in critical infrastructure, the ability to dynamically balance competing objectives will be essential for long-term socio-technical resilience.

References

1. Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (pp. 3645–3650).
2. Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L., Rothchild, D., ... & Dean, J. (2021). Carbon emissions and large neural network training. arXiv preprint arXiv:2104.10350.
3. Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., ... & Kaplan, J. (2022). Training a helpful and harmless assistant from human feedback. arXiv preprint arXiv:2204.05862.
4. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 27730–27744.
5. Schwartz, R., Dodge, J., Smith, N. A., & Etzioni, O. (2020). Green AI. *Communications of the ACM*, 63(12), 54–63.
6. Lacoste, A., Luccioni, A., Schmidt, V., & Dandres, T. (2019). Quantifying the carbon emissions of machine learning. arXiv preprint arXiv:1910.09700.
7. Zhang, Y., Sun, S., Galley, M., Chen, Y., Brockett, C., Gao, X., ... & Dolan, B. (2020). DIALOGPT: Large-scale generative pre-training for conversational response generation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations (pp. 270–278).
8. Meng, K., Bau, D., Solar-Lezama, A., & Belinkov, Y. (2023). Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems*, 36.
9. Fedus, W., Zoph, B., & Shazeer, N. (2022). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120), 1–39.
10. Xin, J., Tang, R., Lee, J., Yu, Y., & Lin, J. (2020). Deebert: Dynamic early exiting for accelerating BERT inference. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (pp. 7131–7143).
11. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.
12. Dixon, L., Li, J., Sorensen, J., Thain, N., & Vasserman, L. (2018). Measuring and mitigating unintended bias in text classification. In Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society (pp. 67–73).

13. Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., ... & Amodei, D. (2020). Scaling laws for neural language models. arXiv preprint arXiv:2001.08361.
14. Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., ... & Sifre, L. (2022). Training compute-optimal large language models. arXiv preprint arXiv:2203.15556.
15. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877–1901.
16. Christiano, P. F., Leike, J., Brown, T. B., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30.
17. Gehman, S., Gururangan, S., Sap, M., Choi, Y., & Smith, N. A. (2020). RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020* (pp. 3356–3369).
18. Shi, C., Li, S., Lu, W., Wu, W., Wang, C., Cheng, Z., ... & Chua, T. S. (2026). TraceRouter: Robust Safety for Large Foundation Models via Path-Level Intervention. arXiv preprint arXiv:2601.21900.
19. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., ... & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35, 24824–24837.
20. Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., ... & Papernot, N. (2023). Extracting training data from large language models. In *30th USENIX Security Symposium* (pp. 2633–2650).