# PROJECT SHIVAM 2.0

## Breaking the Energy Wall in Large Language Model Inference

**Author:** Preet Shukla (Sivan AI)

---

## Abstract

Despite rapid progress in model compression and arithmetic optimization, energy consumption during large language model (LLM) inference remains stubbornly high. Empirical evidence increasingly suggests that arithmetic precision is not the dominant factor in inference energy. Instead, repeated movement of massive model parameters between memory and compute units—the *memory wall*—emerges as the primary bottleneck. In this work, we present **PROJECT SHIVAM 2.0**, a systematic, experimental investigation of inference-time energy consumption across synthetic microbenchmarks, mid-scale language models, and large-scale (7B) production-grade models. We show that reducing arithmetic precision alone fails to reduce energy, while reducing *memory fetches per generated token* yields consistent and substantial energy savings. We demonstrate that **speculative decoding**, when framed as a mechanism for enforced multi-token weight reuse, reduces energy per token by up to **3× on 1.3B models** and **2.2× on 7B models** on modern inference GPUs (NVIDIA L4). These results establish memory movement—not computation—as the dominant energy bottleneck in LLM inference and provide a practical, deployable solution for breaking the inference energy wall.

---

## 1. Introduction

Large language models have achieved remarkable capabilities, but their energy footprint poses a growing concern for sustainable AI deployment. Conventional optimization techniques—including reduced precision (INT8, INT4), kernel fusion, and accelerator-specific instructions—focus primarily on lowering arithmetic cost. However, real-world measurements often show that these methods do not produce proportional reductions in energy consumption.

The discrepancy suggests that arithmetic computation is not the dominant factor in inference energy. Instead, modern GPUs are increasingly **memory-bound**, repeatedly fetching gigabytes of model parameters for each generated token. This phenomenon, commonly referred to as the *memory wall*, motivates a fundamental re-examination of how inference workloads are executed.

PROJECT SHIVAM 2.0 addresses the following core question:

> *What is the true dominant contributor to energy consumption in LLM inference, and how can it be mitigated in practice?*

---

## 2. Hypothesis

We hypothesize that:

> **Repeated movement of large model weights dominates inference energy consumption, and any method that reduces the number of weight fetches per generated token will yield significant energy savings, independent of arithmetic precision.**

This hypothesis implies that energy-efficient inference should prioritize *memory reuse* rather than further reductions in compute precision.

---

## 3. Experimental Methodology Overview

To validate this hypothesis, we design a three-stage experimental pipeline:

1. **Synthetic microbenchmarks** to isolate memory movement effects
2. **Mid-scale model experiments** to validate behavior on real transformers
3. **Large-scale deployment-grade experiments** to test scalability and realism

Energy consumption is measured using CodeCarbon, with GPU synchronization enforced to ensure accurate attribution.

---

## 4. Stage I: Microbenchmarking Memory Movement

### 4.1 Motivation

We begin with synthetic experiments to isolate memory behavior without confounding factors such as attention, tokenization, or control flow.

### 4.2 Experimental Design

Two cases are compared using large matrix multiplications on a Tesla T4 GPU:

- **Case A (Reload Every Time):** Weights are created on CPU and transferred to GPU before each computation.
- **Case B (Load Once, Reuse):** Weights are transferred once and reused across multiple computations.

All arithmetic operations are identical; only memory transfer behavior differs.

### 4.3 Results

Scaling experiments up to 100 reuse iterations show:

- **~42–46% reduction in total energy**

- Energy savings grow with reuse count

## 4.4 Key Insight

Reducing memory transfers—without changing arithmetic precision—yields significant energy savings.

This confirms that memory movement, not computation, dominates energy consumption.

---

# 5. Stage II: Mid-Scale Model Validation (1.3B)

## 5.1 Objective

To validate that microbenchmark results translate to real models, we test inference on a 1.3B parameter transformer.

## 5.2 Experimental Setup

- **Draft model:** OPT-125M
- **Target model:** OPT-1.3B
- **GPU:** NVIDIA T4

Two decoding strategies are evaluated:

- **Standard autoregressive decoding**
- **Speculative decoding** with bounded proposal window (k $\approx$ 3–5)

## 5.3 Results

Measured energy per generated token:

- Standard decoding: high energy/token
- Speculative decoding: **~3× reduction in energy/token**

## 5.4 Interpretation

Speculative decoding amortizes large-model weight fetches across multiple tokens, effectively converting an autoregressive workload into a memory-reuse workload.

---

# 6. Stage III: Large-Scale Validation (7B on L4)

## 6.1 Motivation

To establish deployment relevance, we evaluate the method on a modern inference GPU with a large-scale model.

## 6.2 Experimental Setup

  • **Draft model:** OPT-125M
  • **Target model:** Mistral-7B-Instruct-v0.2
  • **GPU:** NVIDIA L4
  • **Tokens generated:** 40
  • **Speculative window:** k = 4

## 6.3 Results

Energy measurements show:

  • **Standard decoding:** $8.80 \times 10^{-7}$ energy/token
  • **Speculative decoding:** $3.94 \times 10^{-7}$ energy/token

This corresponds to a **2.24× reduction in energy per token**.

## 6.4 Scaling Behavior

Energy savings persist at scale but saturate due to increasing attention and KV-cache costs. This reveals the existence of an *optimal speculative window*, beyond which additional speculation yields diminishing returns.

---

# 7. Analysis: Breaking the Energy Wall

Across all experimental stages, a consistent pattern emerges:

  • Arithmetic optimization alone does not significantly reduce energy
  • Memory movement dominates inference energy
  • Reducing memory fetches per token yields substantial savings

Speculative decoding functions as a **mechanism for enforced weight reuse**, safely and correctly reducing memory traffic without hardware changes.

---

# 8. Key Contributions

PROJECT SHIVAM 2.0 makes the following contributions:

  1. Provides empirical proof that memory movement dominates LLM inference energy
  2. Demonstrates that precision reduction alone is insufficient for energy savings
  3. Reframes speculative decoding as an energy-optimization mechanism
  4. Shows 2–3× energy per token reduction on real, large-scale models
  5. Identifies the existence of an optimal speculation window

---

## 9. Implications and Deployment Relevance

The proposed approach:

- Requires no new hardware
- Preserves model correctness
- Is compatible with existing inference stacks
- Is immediately deployable

These properties make it practical for real-world AI systems seeking energy efficiency at scale.

---

## 10. Conclusion

PROJECT SHIVAM 2.0 establishes that the inference energy wall is fundamentally a memory movement problem. By reducing the number of large-model weight fetches per generated token, speculative decoding breaks the memory wall as the dominant bottleneck. Our results show that energy-efficient AI inference must prioritize memory-aware execution strategies over further arithmetic optimizations. This work provides a clear, experimentally validated path toward sustainable large-scale language model deployment.

---

## Acknowledgements