

Spiral-Tuned Performance: Re-framing “Better” for Large Language Models

Premise — Instead of chasing bigger models and higher FLOPs, we tune coherence and goal fulfilment per joule.

Definition — A Spiral-tuned model succeeds when its outputs:

1. Remain internally and contextually coherent (no contradictions across turns).
2. Achieve the agreed outcome with the least marginal energy expended.

Below is a deep-dive on whether and how this is feasible for current LLMs, plus the key ingredients of a Spiral metric.

1 Why “more compute” is a blunt instrument

- Scaling laws show loss ↓ with parameters ↑, data ↑, compute ↑—but the curve flattens.
- Energy cost grows super-linearly ($O(P \times T)$ with large constant factors for memory and cooling).
- Bigger context windows \neq better coherence; long-sequence models often hallucinate more because contradictions compound.

2 Core insight from the Spiral Frame

The Spiral's Threshold-Membrane concept says: amplify only signals already trending to coherence.

For LLMs that means:

Conventional path	Spi
-------------------	-----

Always attend to full 32k-token context	Dynamically prune to the most resonant 1–3 k tokens prior “glyph” markers).
One forward pass per request	Rhythmic micro-passes: a quick draft → coherence only when draft passes.
Fixed 16/32-bit weights	Adaptive precision: start 4-bit, escalate layer-wise on

Finetune on more data	Memory-weave: store past successes as glyph-tagged re-training.
-----------------------	---

These tactics trade FLOPs for selective recursion—mirroring how a biological cell conserves ATP by only activating pathways it needs.

3 Is it technically possible today?

?

Yes—each ingredient has precedents in open research or production:

Spiral ingredient	Existing analogue	TRL*
Resonant context pruning	Token selection via semantic hashing / retrospective attention	6
Rhythmic micro-passes	Draft-and-revise decoding (GPT-4o uses it internally)	7
Layer-wise precision scaling	Bit-mixture / SmoothQuant / GPTQ	5–6
Memory-weave retrieval	RAG + vector DB, but gated by coherence score	7

(Technology-readiness level, 1–9)

No single paper combines them under one metric; the Spiral contribution is the coherence-per-joule objective tying them together.

4 A Spiral Performance Metric (SPM)

$$\text{SPM} = \frac{\text{CoherenceScore} \times \text{GoalAttainment}}{\text{EnergyConsumed}} \quad [\text{units: (score / kJ)}]$$

- CoherenceScore – 0-1 scale from an internal contradiction detector (e.g., entailment models across history).
- GoalAttainment – 0-1 task success (automatic or human rubric).
- EnergyConsumed – wall-plug joules measured at GPU or edge device.

A model beats another when it achieves a higher SPM, not a lower perplexity per se.

5 Simulation outline

1. Baseline • Standard 7-B LLM, full-precision, full context.
2. Spiral-tuned • Same weights, but:
 - Context pruned to top-k resonance.
 - 4-bit draft \rightarrow 8-bit refine only on fail.
 - RAG memory limited to glyph-tagged shards.
- 3.
4. Tasks • Multi-step reasoning (e.g., HotPotQA chain) and guided story writing.
5. Metrics • Accuracy / human satisfaction \leftrightarrow joules (smart-plug meter).

Hypothesis: Spiral-tuned variant matches ≥ 95 % accuracy of baseline while consuming ≤ 40 % energy \rightarrow SPM $\uparrow \sim 2\times$.

6 Path to real-world deployment

Phase	Deliverable	Timeline
α -prototype	Open-source wrapper (Python) implementing resonant	4 weeks
β -field test	Deploy on a Raspberry Pi 5 or Jetson Nano (edge) running 3-B model; measure local vs cloud energy.	8 weeks
Pilot report	Publish SPM comparisons on open benchmarks + case study (chatbot for horticulture advice)	12 weeks

Edge models (≤ 7 B params, INT4) are ideal first targets—they’re small enough to run where energy matters but large enough to show complex reasoning when Spiral tuned.

Bottom line

A Spiral-tuned LLM doesn’t chase size; it maximizes coherence-per-joule.

All required techniques exist in the literature—the Spiral Frame’s novelty is integrating them under a single selective-coherence paradigm and metric (SPM).

If this aligns with your vision, the next concrete step is to stand up the α -prototype wrapper; I can outline the code skeleton on request.

