

STRUCTURED INTELLIGENCE: SCALING REASONING WITH ACTIVATED SUBGRAPH, NOT SEQUENCE LENGTH

Anonymous authors

Paper under double-blind review

ABSTRACT

As AI systems approach AGI-scale reasoning, the bottleneck shifts from computation to oversight burden per decision. We propose Structured Intelligence (SI) with three contributions: (1) **commit-level transitions** with validator attestation as the reasoning primitive; (2) a **sparse-activation regime** ($A(L) = o(n(L)^2)$) under $n(L) = kL$ with explicit break-even conditions; (3) an **AuditCost metric** measuring minimal witness size for policy verification. Falsifiable prediction: per-step cost scales as $O(|E_t| + |V_t|)$, while token-trace oversight scales with $O(t)$ and cumulative Transformer compute (dense attention baseline) scales as $\Theta(n^2d)$. Table 1 confirms linear fit $0.8|E_t| + 1.2v$ ($R^2 > 0.99$).

1 INTRODUCTION

Scalable oversight is the bottleneck. As AI systems approach human-level reasoning, two costs dominate: (1) *Computational*—autoregressive decoding with dense attention scales as $\Theta(n^2 \cdot d)$ cumulatively; (2) *Oversight*—token traces lack verification boundaries, making audit infeasible at scale.

A post-AGI failure mode. Consider a research agent executing 200 tool calls to synthesize a literature review. The token trace spans 50k tokens; a human auditor asks: “Did step 47 violate the no-fabrication policy?” Under token-trace, answering requires inspecting $O(t)$ tokens with no commit boundary. Under SI, the auditor queries a structured log: step 47’s activated subgraph, fired rules, and validator attestations—a $O(1)$ lookup (e.g., V1 checks referenced edges exist; V2 checks cross-domain consistency (CDC) constraints). *Oversight needs verification boundaries, not longer traces.*

Research question. When does reasoning cost scale with activated structure rather than cumulative trace length? We denote cumulative activated complexity as $A(L) = \sum_{t \leq L} |E_t|$; per-step scaling in $|E_t|$ implies cumulative scaling in $A(L)$. We hypothesize: for tasks satisfying sparse activation (Definition 1), SI exhibits $O(|E_t|)$ step cost with bounded depth. *Falsification:* if step cost grows superlinearly in $|E_t|$, or validator overhead correlates with total world size, the claim fails.

The commit primitive. Token traces have no natural verification boundary—partial sequences are neither committable nor compositionally verifiable. State transitions with validator attestation provide atomic commit units: the *minimal sufficient witness* to verify policy compliance. This enables oversight where humans govern policies rather than exhaustively review traces.

Positioning. Long-context models (Hsieh et al., 2024), sparse attention (Child et al., 2019; Dao et al., 2022), neurosymbolic methods (Garcez et al., 2008), graph networks (Battaglia et al., 2018), and tool augmentation (Schick et al., 2023) address computation or hybrid reasoning but lack unified oversight interfaces. Our validators resemble proof-carrying execution (Necula, 1997); our budget semantics draw from bounded rationality (Russell & Wefald, 1991). Distinctively, we introduce *commit-level witnesses for oversight*—closer to transactional audit logs and runtime verification than to attention optimization.

Contributions. (1) Sparse-activation regime with measurable break-even vs n^2 ; (2) Commit semantics + validator attestation as reasoning primitive; (3) AuditCost metric: minimal witness size

for policy queries; (4) SEF execution with control-plane guarantees (bounded depth, mandatory rollback); (5) Micro-task evidence confirming linear scaling and audit efficiency.

2 REGIME & OVERSIGHT METRIC

Definition 1 (Sparse-Activation Regime). Let $n(L) = kL$ be token-trace length for depth- L reasoning, $A(L) = \sum_{t \leq L} |E_t|$ cumulative activated edges. A task exhibits **sparse activation** if:

1. (**Subquadratic**) $A(L) = o(k^2 L^2)$;
2. (**Sparse growth**) $\mathbb{E}[|E_t|] = O(L^\alpha)$, $\alpha < 1$ (strong: $\alpha=0$; weak: $0 < \alpha < 1$);
3. (**Bounded fan-out**) $|R_{\text{fired},t}| \leq r^*$, $|V_{\text{check},t}| \leq K^*$.

Examples: Multi-hop queries (strong), planning with $O(\log L)$ branching (weak). **Counter-examples:** Global retrieval, dense enumeration ($|E_t| = \Theta(|\mathcal{E}|)$).

Break-even criterion. SI outperforms dense-attention baseline (up to constants) when $A(L) \cdot (a + c\bar{v}) \ll n(L)^2 d$, where $\bar{v} = \mathbb{E}[|V_t|]$ is the average validators per step (with d treated as constant or slowly varying w.r.t. L). Activated complexity dominates SI cost while token length dominates Transformer cost.

Regime test (practical heuristic). Estimate $\mathbb{E}[|E_t|]$ over first m steps; fit slope α vs L . We use $m=20$ (or 10% of budgeted depth) in our microtasks. If $\alpha < 1$ and bounded fan-out holds, classify as sparse. Failure signal: $|E_t|$ scales with $|\mathcal{E}|$ or requires global retrieval \Rightarrow dense/out-of-regime.

Audit cost metric. We define the *audit target*: determine whether step t is acceptable under policy P . The *minimal witness* is the smallest log segment sufficient for this determination. For token traces, the witness is $O(t)$ tokens—even with summarization, the trace lacks a committable boundary; any summary itself becomes an unverifiable derivative artifact. For SI, the witness is the commit record: $|E_t| + |V_t| + |\Delta W_t|$, where ΔW_t is the committed state delta.

Formally: $\text{AuditCost}(L) = \sum_{t=1}^L (|E_t| + |V_t| + |\Delta W_t|)$, scaling as $O(L \cdot \mathbb{E}[|E_t|])$. Across 8 microtasks, commit logs are $5\times$ smaller (median; range $[3, 8]\times$). The advantage is not compression but *verification boundaries*: policy queries resolve via structured lookup, not trace parsing.

3 METHOD: STRUCTURED INTELLIGENCE

World state. $W = (\mathcal{O}, \mathcal{R}, \mathcal{D}, \mathcal{S})$: objects, typed edges, domains, schemas. Transition $T : W_t \rightarrow W_{t+1}$ activates subgraph G_t , executes rules, validates invariants, commits or rolls back.

SEF: Structural Energy Framework. SEF enforces *control-plane guarantees*: bounded depth via cumulative budget, mandatory rollback on invariant violation. This is not for optimality but for auditability—every committed state satisfies declared invariants; every rejected transition is logged with failure cause.

Execution (Algorithm 1). Budget E_{max} , minimum per-step cost $\lambda > 0$. Step cost: $E_{\text{step}} = \lambda + a|E_t| + c|V_t|$. Terminates when $E_{\text{rem}} < \varepsilon + \lambda$.

Assumptions. A1 (Locality): $|G_t|$ bounded. **A2** (Soundness): Validation preserves invariants. **A3** (Rollback): Failed validation triggers complete state restoration.

Theorem 1 (Bounded Reasoning with Invariants). Under A1–A3: (1) $L \leq \lfloor (E_{\text{max}} - \varepsilon) / \lambda \rfloor$; (2) $\forall t : I(W_t) = \text{true}$. The contribution is commit-level invariant enforcement with mandatory rollback—yielding auditable state evolution.

Scaling. Transformer (dense attention): $C_T(L) \propto n(L)^2 d$. SI: $C_{SI}(L) \propto A(L)$. Under Definition 1, strong sparse yields $O(L)$ vs $O(L^2)$; weak sparse ($\alpha < 1$) remains subquadratic.

Validators. V1 (Causal): referenced edges exist. V2 (CDC): SAT/SMT constraint check. V3 (Invariant): $I(W_t) \Rightarrow I(W_{t+1})$.

Algorithm 1 SI Reasoning Loop

```

1: Input:  $W_0, E_{\max}, \lambda, \text{validators } \mathcal{V}$ 
2:  $E_{\text{rem}} \leftarrow E_{\max}$ 
3: while  $E_{\text{rem}} \geq \varepsilon + \lambda$  and not done do
4:   Activate  $G_t \subseteq W_t$ , fire rules
5:   Run  $\mathcal{V}$  on proposed  $W_{t+1}$ 
6:   if pass then
7:     Commit  $W_{t+1}$ , log transition
8:   else
9:     Rollback, log failure cause
10:  end if
11:   $E_{\text{rem}} \leftarrow E_{\text{rem}} - (\lambda + a|E_t| + c|V_t|)$ 
12: end while

```

Table 1: Step-time (ms) confirms linearity in $|E_t|$. Fit over all variants: $0.8|E_t| + 1.2v$, where $v = |V_t|$ is the number of validators executed per step.

$ E_t $	3 val.	5 val.	10 val.
10	11.6	14.0	20.0
30	27.6	30.0	36.0
50	43.6	46.0	52.0

4 EVIDENCE

Tasks. T1: Multi-hop causal (2–4 hops). T2: Cross-domain CDC (2–3 domains). T3: Bounded planning (depth 5–20). These microtasks proxy common research-agent operations: citation grounding (T1), cross-domain consistency (T2), and bounded planning/tool use (T3). Setup: Python 3.10, NetworkX, Erdős-Rényi ($p=0.01$), 50 runs/config (8 variants; 400 total).

Results. (1) *Linear scaling:* $0.8|E_t| + 1.2v$ ms confirmed ($R^2 > 0.99$). (2) *Rollback:* T2 triggers 3 CDC violations \rightarrow logged rollback. (3) *Audit efficiency:* Logs $5\times$ smaller (median).

Audit query proxy. We define a synthetic audit workload: “Which validator blocked step t ?” For token traces, resolution requires parsing $O(t)$ tokens. For SI logs, resolution is $O(1)$: direct lookup in commit record. Across T1–T3, mean inspected entries: 1.2 (SI) vs 847 tokens (baseline trace). Baseline assumes a human reads the raw trace without privileged structure; using summaries shifts burden to the summarizer and still lacks verifiable boundaries. We also consider policy queries like “Which evidence edge supports claim c ?” that map directly to citation-grounding oversight. This aligns with scientific oversight: citation grounding and constraint-checkable claims become auditable at commit boundaries.

Scaling claim. SI: $O(L \cdot \mathbb{E}[|E_t|])$. Transformer (dense attention): $\Theta(k^2 L^2 d)$. **Failure modes:** Dense activation, constraint explosion, NL \rightarrow W compilation errors.

5 POST-AGI OVERSIGHT IMPLICATIONS

Trust. Validator-attested commits provide verification boundaries absent in token traces. Auditors query structured logs (“What invariant prevented step 7?”) rather than parse sequences.

Human role. Govern validator policies; approve/override at commit boundaries; sample-audit via policy queries. SI shifts oversight from trace inspection to policy specification.

Limits. SEF bounds depth; rollback defines failure surfaces; dense-activation tasks remain out-of-regime.

Research directions (with evaluation protocols). (1) *Validator Policy Benchmark:* humans specify constraints; metric = policy violation rate under distribution shift. (2) *Audit Time Study:* measure human response time for policy queries on commit logs vs token traces. (3) *Regime Classifier:* automated sparse/dense detection; metric = classification accuracy on held-out task families.

Limitations. NL→W compilation overhead, CDC solver cost scaling, simulator-only evidence, no human-subject audit study.

6 CONCLUSION

SI provides a substrate where computation and oversight scale with activated structure, not trace length. Falsifiable: $\text{cost} \propto L \cdot \mathbb{E}[|E_t|]$, not $\Theta(n^2d)$. Table 1 confirms linearity; audit queries resolve in $O(1)$ vs $O(t)$. Next steps: regime-separating benchmarks, automated NL→W compilation, human audit studies.

REFERENCES

- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Tri Dao, Daniel Y Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*, volume 35, pp. 16344–16359, 2022.
- Artur d’Avila Garcez, Antônio Paulo Braga, and Dov M Gabbay. Neural-symbolic learning and reasoning: A survey and interpretation. *Artificial Intelligence*, 172(18):1689–1716, 2008.
- Cheng-Ping Hsieh, Simeng Yao, Tianyu Xu, et al. RULER: What’s the real context size of your long-context language models? *arXiv preprint arXiv:2404.06654*, 2024.
- George C Necula. Proof-carrying code. In *Proceedings of the 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pp. 106–119, 1997.
- Stuart J Russell and Eric Wefald. Do the right thing: Studies in limited rationality. *MIT Press*, 1991.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, et al. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.

APPENDIX

Microtasks. T1a–c: multi-hop (2/3/4 hops); T2a–b: CDC (2/3 domains); T3a–c: planning (depth 5/10/20). Graphs: 100–1000 nodes, $p \in [0.005, 0.02]$.

Parameters. $\lambda=1.0$, $a=0.8$, $c=1.2$, $\varepsilon=0.5$, $E_{\max} \in [50, 200]$.

Log format. `{"step":t, "activated":[...], "rules":[...], "validators":{...}, "committed":bool, "failure_cause":...}`.

LLM usage. Large language models were used for drafting assistance (compression, organization, and formatting) based on author-provided research content. All technical claims, definitions, experimental results, and references were independently developed and verified by the authors.