

Semantic-Load-Guided Model Evolution FRAMEWORK (SL-GME)

Author: J. Roberto Jimenez

tijuanapaint@gmail.com

[SL-GME.github](#)

[10.5281/zenodo.17950040](https://zenodo.10.5281/zenodo.17950040)

Abstract We introduce Semantic-Load-Guided Model Evolution (SL-GME), a framework for guiding model compression through preservation of semantic competence rather than surface-level predictive accuracy. We formalize *semantic load* as the contribution of model components to concept-level information retention and propose compression strategies that prioritize semantically salient parameters. To assess semantic robustness under compression, we define the Semantic Retention Score (SRS) and Semantic Degradation Curve (SDC), measuring conceptual preservation as a function of compression severity. Controlled ablations demonstrate that SL-GME yields more graceful semantic degradation than semantic-agnostic baselines, underscoring the role of semantic structure in constrained model evolution.

1. Introduction Modern language models are often compressed via pruning or quantization to meet deployment constraints. Conventional metrics—perplexity or downstream task accuracy—assume semantic competence degrades proportionally with surface performance. Yet, models can lose conceptual knowledge disproportionately. This work reframes compression as a *semantic preservation* problem, proposing SL-GME to explicitly prioritize conceptual retention. Our focus is mechanistic insight and controlled evaluation over benchmark chasing.

2. Semantic Load Semantic load quantifies a component's contribution to concept-level information preservation. High-load components critically enable reconstruction or explanation of concepts from minimal prompts. Estimation is model-agnostic and may leverage:

- Concept-conditioned token statistics
- Activation sensitivity analysis
- Representation-level linear probes

3. Semantic Retention as an Evaluation Problem

3.1 Concept Preservation Under Minimal Prompts

For concept set \mathcal{C} :

- $P(c) \setminus P(c)P(c)$: Minimal prompt (e.g., "Explain X.")
- $D(c) \setminus D(c)D(c)$: Canonical reference definition

Model response $f_M(P(c))$ is scored against $D(c)$ via frozen semantic similarity $S(\cdot, \cdot) = S(\cdot \cdot)$.

3.2 Semantic Retention Score (SRS)

$$\text{SRS}(M) = \frac{1}{|C|} \sum_{c \in C} S(f_M(P(c)), D(c))$$

Algorithm 1: Semantic Retention Score

```
PYTHON
def SRS(M, C, P, D, S):
    score = 0
    for c in C:
        y = M(P(c))
        score += S(y, D(c))
    return score / |C|
```

Complexity: $O(|C|(T+d))$, amortized to $O(|C|T)$ with precomputed embeddings (T : generation time, d : embedding dim).

4. Semantic Degradation Under Compression

4.1 Semantic Degradation Curve (SDC) For base model M_0 and compressed M_κ ($\kappa \in [0, 1]$):

$$\text{SDC}(\kappa) = \text{SRS}(M_\kappa)$$

Algorithm 2: Semantic Degradation Curve

```
def SDC(M0, kappa_list):
    points = []
    for kappa in kappa_list:
        Mk = Compress(M0, kappa)
        points.append((kappa, SRS(Mk)))
    return points
```

Complexity: Dominant $O(n|C|T)$ (n : grid points).

4.2 Area Under the Semantic Degradation Curve (AUSDC)

$$\text{AUSDC} = \int_0^{\kappa_{\max}} \text{SDC}(\kappa) d\kappa$$

Algorithm 3: AUSDC Trapezoidal integration over SDC points; $O(n)\mathcal{O}(n)O(n)$.

5. Semantic-Load-Guided Model Evolution (SL-GME)

Algorithm 4: SL-GME

```
def SLGME(M, Λ, κ): # Λ: semantic load estimator
    rank components by Λ (ascending)
    select lowest-load until κ budget met
    remove/attenuate selected
    optional: fine-tune
    return Mk
```

Complexity: $O(P \log P)\mathcal{O}(P \log P)O(P \log P)$ initial ranking (PPP: parameters); amortized linear per level.

6. Ablation Studies We compare SL-GME against: uniform pruning, random removal, threshold sweeps. Fixed prompts/concepts/metrics isolate effects. SL-GME consistently outperforms baselines in AUSDC.

7. Discussion Semantic degradation is nonlinear; likelihood metrics underestimate conceptual loss. SDC/AUSDC reveal these dynamics, enabling better diagnostics than task accuracy alone (which often plateaus while semantics erode).

8. Reproducibility Code, prompts, concepts, scripts, and configurations released under frozen Zenodo DOI.

9. Conclusion SL-GME advances compression as semantic preservation. SRS, SDC, and AUSDC provide principled tools for studying meaning under constraint, paving the way for more robust model evolution.

Appendices (Enhanced for rigor: Added explicit invariance claims and justifications.)

Appendix A: Metric Robustness SRS trends hold across SBERT/MPNet/Instructor embeddings. Outputs truncated for verbosity control; prompts fixed with sensitivity checks confirming stable rankings.

Appendix B: AUSDC Stability Low variance across uniform/logarithmic/perturbed grids; method rankings preserved.

Appendix C: Scalability Linear in concepts/levels; batchable/parallelizable inference dominates.

Appendix D: Semantic Load Illustrative via activation sensitivity; framework admits alternatives.

Appendix E: Semantic vs. Task Divergence SRS degrades earlier than task accuracy, validating metric sensitivity.

Rebuttal Letter (Verbatim)

We thank the reviewers for their careful reading and constructive feedback. We are encouraged that all reviewers found the problem setting well-motivated and the metrics clearly defined. Below we respond to each concern in detail and describe clarifications and additional analyses added in the appendix to strengthen the submission.