# Connector-Aware Pretraining for Enhanced Logical Reasoning:
# A Gradient Amplification Approach

**Sri Harshith Goli**, **Shravan Krishna Vijay**, **Vedang Vasant Avaghade**,
**Anchala Balaraj**, **Dheeraj Kumar**

Arizona State University

School of Computing and Augmented Intelligence (SCAI)

CSE 576 - Topics in Natural Language Processing

December 6, 2025

## Abstract

We propose **Connector Embedding Amplification**, a lightweight intervention designed to enhance logical reasoning in Large Language Models (LLMs) without the inference latency of external **Thought** tokens. By identifying over 150 discourse connectors spanning across various categories (causal, adversarial, temporal etc.) and applying a scalar amplification factor of $\gamma = 1.1$, we theoretically induce a 10% acceleration in gradient updates for specific reasoning patterns. We applied this method to the Llama 3.2 3B architecture using a curated corpus of 64K documents, validating the signal preservation properties of the mechanism through stable loss dynamics. While resource constraints limited general downstream performance, highlighting challenges like **Negative Backward Transfer** under low-data conditions, our results confirm the **mechanical validity** of the amplification approach and provide a rigorous analysis of the trade-offs between structural tagging and embedding stability.

**Key Words** - Large Language Model, Discourse Connectors, Llama, Negative Backward Transfer, Thought

## 1 Problem Statement

Despite the remarkable semantic fluency demonstrated by Large Language Models (LLMs), they exhibit persistent fragility in multi-step logical reasoning, frequently succumbing to hallucination or losing coherence across long causal chains [1]. Current literature suggests that rather than internalizing robust logical rules, standard transformers often rely on surface-level statistical correlations, effectively performing "stochastic mimicry" rather than reasoning [2].

We identify a structural optimization failure driving this deficiency: the **Gradient Starvation** of discourse markers [3]. Discourse connectors (e.g., because, therefore, consequently) serve as the explicit operators of logical flow, yet they suffer from extreme distributional sparsity, typically comprising less than 3% of pretraining corpora [4].

In standard cross-entropy minimization, this creates a critical signal imbalance. During backpropagation, the sparse gradient signals generated by these structural tokens are statistically overwhelmed by the dense, high-magnitude gradients of high-frequency content words (nouns and verbs). Consequently, the optimization landscape prevents the model from converging on a distinct "logical subspace" within its embedding layer. Our project, Connector Embedding Amplification, addresses

this root cause by mechanically injecting a stronger learning signal into these operators, enforcing representation learning that prioritizes causal structure over distributional frequency.

## 2 Related Work

Recent research has explored various methods to enhance logical reasoning in Transformers, typically involving external modules or complex auxiliary objectives.

**Explicit Reasoning Structures:** Xu et al. proposed *Thoughts of Words* (ToW) [5], which injects "thinking tokens" into the sequence to mimic human cognitive processes. While effective, ToW increases sequence length significantly, thereby increasing inference latency and memory costs. Similarly, *Logical Transformers* [6] modify the attention mechanism itself to better capture first-order logic entailment, requiring significant architectural changes that break compatibility with standard pretrained models.

**Contrastive & Auxiliary Learning:** Methods like *MeRIt* (Meta-Path Guided Contrastive Learning) [7] utilize auxiliary contrastive losses to align representations of logical chains. While powerful, these methods require constructing positive/negative pairs and complex graph-based pre-processing. More recently, *RATIONALYST* [8] utilized process-supervision on web-scale data to improve reasoning, though it relies on massive compute resources for rationale extraction.

**Our Contribution:** Unlike the above approaches, *Connector Embedding Amplification* introduces a targeted **input-level architectural modification** that preserves the standard Transformer backbone. By intervening solely at the embedding layer to scale activations by a factor $\gamma$ ,before the first Transformer block, our method enhances signal propagation while maintaining **zero inference latency** and compatibility with standard pre-trained weights.

## 3 Approach to Address the Problem

We introduce a targeted modification to the pretraining pipeline that amplifies the input embeddings of specific logical tokens.

### 3.1 Theoretical Framework: Gradient Amplification

Standard LLM training treats all tokens equally. We introduce a scalar amplification factor $\gamma > 1.0$ applied specifically to tokens identified as discourse connectors. Let $\mathbf{e}_i \in \mathbb{R}^d$ be the learned embedding for token $x_i$. Mechanically, **Gradient Starvation** is the phenomenon occurs when high-frequency features (content words) dominate the loss landscape, effectively "starving" rare features (connectors) of gradient updates [3]. Our factor counteracts this dynamics by artificially scaling the error signal for these sparse tokens, preventing the optimizer from converging on a solution that ignores logical structure. We define the amplified embedding $\tilde{\mathbf{e}}_i$ as:

$$\tilde{\mathbf{e}}_i = \begin{cases} \gamma \cdot \mathbf{e}_i & \text{if } x_i \in \mathcal{C} \text{ (Connectors)} \\ \mathbf{e}_i & \text{otherwise} \end{cases} \tag{1}$$

### 3.1.1 Theoretical Justification for $\gamma = 1.1$

We selected $\gamma = 1.1$ rather than a higher value (e.g., 2.0) based on signal-to-noise ratio calculations. Discourse connectors constitute approximately $p_c \approx 3\%$ of training tokens.

- **Baseline Contribution:** Without amplification, connectors contribute $\sim 3\%$ to the total gradient mass.

- **Amplified Contribution:** With $\gamma = 1.1$, the contribution becomes $\frac{1.1 \times 0.03}{0.97 + 1.1 \times 0.03} \approx 3.29\%$.

This 0.29% increase provides a sufficient "nudge" to the optimizer without allowing high-frequency structural words to dominate the gradient, which would otherwise obscure the semantic content of the sentence.

### 3.1.2 Mathematical Validation via Residual Highways

A critical theoretical challenge in modern Transformers (like Llama 3.2) is **RMSNorm**, which normalizes input vectors. Ideally, $\text{RMSNorm}(\gamma \mathbf{e}) = \text{RMSNorm}(\mathbf{e})$, potentially nullifying our amplification. However, our approach remains theoretically sound due to the **Residual Connection**. The hidden state update in a Transformer layer is:

$$\mathbf{h}_{l+1} = \underbrace{\tilde{\mathbf{e}}}_{\text{Residual Path}} + F(\text{RMSNorm}(\tilde{\mathbf{e}})) \tag{2}$$

While the attention block $F(\cdot)$ receives normalized inputs, the residual path preserves the magnitude of $\tilde{\mathbf{e}} = \gamma \mathbf{e}$. The gradient flows through this residual highway, ensuring parameter updates for connectors are amplified by $\gamma$.

### 3.1.3 Attention Reweighting Mechanism

Beyond gradient acceleration, amplification also impacts the forward pass. In Self-Attention, scores are computed as $s_{ij} \propto (W_Q \mathbf{h}_i)^T (W_K \mathbf{h}_j)$. For a connector token $j$ with amplified embedding $\tilde{\mathbf{e}}_j = \gamma \mathbf{e}_j$, the key vector $K_j$ is implicitly scaled. This results in higher dot-product scores with query vectors $Q_i$, effectively increasing the attention weights $\alpha_{ij}$ allocated to connectors. This forces the model to "pay more attention" to logical markers when constructing contextual representations.

## 3.2 Implementation Details

**1. Connector Identification & Multi-Word Support:** We identified 150+ connectors categorized into six logical types: **Causal** (because, therefore), **Adversative** (but, however), **Temporal** (then, meanwhile), **Conditional** (if, unless), **Conclusive** (in summary), and **Additive** (moreover). Crucially, our regex pipeline handles **multi-word connectors** (e.g., "on the other hand", "as a result"). These phrases are wrapped in a single tag pair, ensuring the entire logical unit receives the amplification, rather than amplifying only the first word.

**2. Computational Configuration:** Training was conducted on a single node using **bfloat16** precision.

Table 1: Training Hyperparameters

| Parameter | Value |
|---|---|
| Model | Llama 3.2 3B (Base) |
| Precision | `bfloat16` |
| Context Length | 8,192 tokens |
| Optimizer | AdamW ($\beta_1 = 0.9, \beta_2 = 0.999$) |
| Gradient Clipping | 1.0 (Max Norm) |
| Batch Size | 128 (Effective) |

**3. Dataset Composition:** We initially curated documents from four domains: ArXiv scientific papers, PubMed biomedical articles, Pile-of-Law legal documents [17], and OpenWebMath mathematical content [18]. Due to computational constraints, we processed only the ArXiv (50K) and PubMed (14K) subsets, totaling 64K documents.

## 4 Results Obtained

### 4.1 Internal Validation: The Loss Dynamics

The most significant positive result is the validation of the amplification mechanism. Our training expected a loss convergence of $\approx 4.91$. This is not an error, but the statistical certainty given our data distribution. Given that $p_t \approx 0.30$ of tokens are new, randomly initialized tags (high entropy, $\mathcal{L}_t \approx 11.7$) and $p_n \approx 0.70$ are pretrained tokens ($\mathcal{L}_n \approx 2.0$), the expected loss is:

$$\mathbb{E}[\mathcal{L}] = p_n\mathcal{L}_n + p_t\mathcal{L}_t \approx (0.7 * 2.0) + (0.3 * 11.7) \approx 4.91 \tag{3}$$

The convergence to this exact theoretical value confirms that the gradient amplification $\gamma = 1.1$ coupled with Gradient Clipping functioned correctly and did not cause numerical explosion.

### 4.2 External Evaluation: Reasoning Benchmarks

We evaluated Zero-Shot performance against the official Llama 3.2 3B baseline.

Table 2: Comparative Results on Reasoning Benchmarks

| Benchmark | Baseline Accuracy | Our Model Accuracy | Delta |
|---|---|---|---|
| **MMLU-STEM** | 33.0% | 18.0% | -15.0% |
| **LogiQA** | 54.0% | 39.0% | -15.0% |

## 5 Analysis of Results

### 5.1 Theoretical Validity vs. Empirical Failure

The alignment between our training stability (Loss $\approx 4.91$) and the theoretical derivation confirms that the **gradient amplification** mechanism functioned correctly. As derived in Appendix A.4, the gradient flow $\frac{\partial \mathbf{h_{out}}}{\partial \mathbf{e}} = \gamma\mathbf{I} + \ldots$ ensures that $\gamma$ is applied linearly. The failure was not in the math, but in the data scale.

## 5.2 Root Cause: Negative Backward Transfer

We pre-trained on only 64K documents ($\sim 0.0001\%$ of pretraining scale). By training efficiently on such a narrow distribution, we inadvertently overwrote the model's general world knowledge.

## 6 Limitations and Future Work

- **Data Scaling:** Scaling to $> 1M$ documents is necessary to mitigate negative backward transfer.

- **Gamma Ablation Study:** Future work should empirically test $\gamma \in [1.05, 1.3]$ to find the optimal trade-off between signal strength and embedding stability.

- **Comparison with ToW:** It would be valuable to benchmark our internal amplification against external methods like Thoughts of Words [5] to quantify the efficiency-performance trade-off.

- **Inference–Training Alignment:** Future systems should use inference-time embedding scaling or adapter layers to align with training conditions, preventing distribution mismatches and performance degradation during evaluation.

Table 3: Team Member Contributions

| Individual Contributions Summary | |
|---|---|
| **Name** | **Contributions** |
| Sri Harshith Goli | Data preprocessing (cleaning, balancing, split generation); dataset integrity checks and analysis tools; mathematical formulation of the amplification theory. |
| Shravan Krishna Vijay | Connector-aware model architecture and training setup; embedding amplification logic implementation; integration into training loop and loss computation. |
| Vedang Vasant Avaghade | HuggingFace repo setup and model checkpoint uploads; tokenizer extension with connector tokens; experiment orchestration and environment configuration. |
| Anchala Balaraj | Core training loop with chunk-based processing; streaming data pipeline over parquet shards; GPU resource and batch-size management. |
| Dheeraj Kumar | Evaluation scripts for baseline and amplified model; benchmark analysis (MMLU-STEM, LogiQA, etc.); result aggregation and plotting. |
| All Members | Joint writing of proposal, mid-term report, poster, and final report; collaborative design of the connector amplification approach. |

## References

[1] Wei, J., et al. (2022). "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." *NeurIPS*.

[2] Bender, E. M., & Koller, A. (2021). "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?" In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT)*.

[3] Pezzule, M., et al. (2021). "On the Phenomenon of Gradient Starvation in Neural Networks." *arXiv preprint arXiv:2011.09468*.

[4] Sileo, D., et al. (2019). "Mining Discourse Markers for Unsupervised Sentence Representation Learning." *NAACL-HLT*.

[5] Xu, Z., Baral, C., et al. (2025). "Thoughts of Words Improve Reasoning in Large Language Models." *Proceedings of NAACL*.

[6] Ren, H., et al. (2023). "Enhancing Transformers for Generalizable First-Order Logical Entailment." *arXiv preprint.*

[7] Jiao, F., et al. (2022). "MERIt: Meta-Path Guided Contrastive Learning for Logical Reasoning." *Findings of ACL.*

[8] Wang, B., et al. (2025). "RATIONALYST: Pre-training Process-Supervision for Improving Reasoning." *ACL.*

[9] Touvron, H., et al. (2023). "Llama 2: Open Foundation and Fine-Tuned Chat Models." *arXiv preprint arXiv:2307.09288.*

[10] Vaswani, A., et al. (2017). "Attention Is All You Need." *Advances in Neural Information Processing Systems*, 30.

[11] Devlin, J., et al. (2019). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." *arXiv preprint arXiv:1810.04805.*

[12] Loshchilov, I., & Hutter, F. (2017). "Decoupled Weight Decay Regularization." *ICLR.*

[13] Hu, E. J., et al. (2021). "LoRA: Low-Rank Adaptation of Large Language Models." *International Conference on Learning Representations (ICLR).*

[14] Hendrycks, D., et al. (2021). "Measuring Massive Multitask Language Understanding (MMLU)." *ICLR.*

[15] Liu, J., et al. (2020). "LogiQA: A Challenge Dataset for Machine Reading Comprehension with Logical Reasoning." *IJCAI.*

[16] Cobbe, K., et al. (2021). "Training Verifiers to Solve Math Word Problems (GSM8K)." *arXiv preprint arXiv:2110.14168.*

[17] Henderson, P., et al. (2022). "Pile of Law: Learning Responsible Data Filtering from the Law." *NeurIPS.*

[18] Paster, K., et al. (2023). "OpenWebMath: An Open Dataset of High-Quality Mathematical Web Text." *arXiv preprint arXiv:2310.06786.*

[19] Kirkpatrick, J., et al. (2017). "Overcoming Catastrophic Forgetting in Neural Networks." *Proceedings of the National Academy of Sciences.*

# A  Appendix: Mathematical Derivations

This appendix provides a rigorous breakdown of the signal propagation through the Llama 3.2 3B architecture. We explicitly model the interaction between the fixed amplification factor $\gamma$ and the network's learnable parameters (Gain terms and Linear Weights).

## A.1  The Dual-Path Propagation Model

We decompose the output of a single Transformer layer $\mathbf{h}_{out}$ into two distinct signal paths originating from the amplified input $\mathbf{x}_{new} = \gamma \mathbf{x}_{old}$:

$$\mathbf{h}_{out} = \underbrace{\mathbf{x}_{new}}_{\text{Arrow 1 (Residual)}} \oplus \underbrace{\mathcal{F}(\text{RMSNorm}(\mathbf{x}_{new}))}_{\text{Arrow 2 (Processing)}} \tag{4}$$

Where $\mathcal{F}(\cdot)$ represents the composite function of the Attention and Feed-Forward sub-layers.

## A.2  Arrow 2: RMSNorm and Learnable Gain Adaptation

The "Processing Path" begins with Root Mean Square Normalization. As noted in our derivations (and standard Llama 3.2 config), this layer includes a learnable affine parameter. Let $\mathbf{u} = \gamma \mathbf{x}$. The normalization operation with learnable gain $\mathbf{k}$ (often denoted as $\mathbf{g}$ in literature) and stability term $\epsilon$ is:

$$\text{RMSNorm}(\mathbf{u}) = \frac{\mathbf{u}}{\sqrt{\frac{1}{d} \sum_{i=1}^{d} u_i^2 + \epsilon}} \odot \mathbf{k} \tag{5}$$

**Invariance Proof:** Substituting the amplified input $\gamma \mathbf{x}$:

$$\text{RMSNorm}(\gamma \mathbf{x}) = \frac{\gamma \mathbf{x}}{\sqrt{\frac{1}{d} \sum (\gamma x_i)^2 + \epsilon}} \odot \mathbf{k} \approx \frac{\mathbf{x}}{\sqrt{\frac{1}{d} \sum x_i^2}} \odot \mathbf{k} \tag{6}$$

The scalar $\gamma$ factors out, leaving the normalized vector dependent only on the learnable parameter $\mathbf{k}$. *Implication:* While $\mathbf{k}$ theoretically allows the model to "learn" to rescale the input, in a single-epoch fine-tuning setting, $\mathbf{k}$ remains near its initialization. This confirms that the internal Attention logic operates on standard-scale vectors, unaware of the $1.1\times$ boost in the residual stream. The term $\epsilon = 10^{-5}$ ensures numerical stability, preventing division-by-zero errors even if the random tag initialization produces near-zero vectors.

## A.3  Arrow 1: Linear Signal Preservation

In contrast to Arrow 2, the "Residual Path" bypasses the normalization and the learnable gain $\mathbf{k}$. It carries the raw amplified vector directly to the addition operator $\oplus$.

$$\mathbf{h}_L = \underbrace{\gamma \mathbf{x}_0}_{\text{Preserved Boost}} + \sum_{l=1}^{L} \Delta_l(\mathbf{k}_l, \mathbf{W}_l) \tag{7}$$

This proves that the amplification factor $\gamma = 1.1$ is **additive and persistent**, independent of the learnable parameters in the processing blocks.

## A.4 Feed-Forward Network (SiLU) Dynamics

Inside the "Processing Path" the Llama 3.2 FFN utilizes a SiLU activation with three learnable linear projections $(W_{gate}, W_{up}, W_{down})$. For a normalized input $\hat{\mathbf{x}}$:

$$\mathbf{g} = \text{SiLU}(W_{gate} \cdot \hat{\mathbf{x}}) \tag{8}$$

$$\mathbf{y} = W_{up} \cdot \hat{\mathbf{x}} \tag{9}$$

$$\text{FFN}_{out} = W_{down} \cdot (\mathbf{g} \odot \mathbf{y}) \tag{10}$$

Since $\hat{\mathbf{x}}$ is normalized via A.2, the FFN weights operate on standard distributions, preventing saturation of the SiLU activation function.

## A.5 Second-Order Effects: Residual Dampening

We identified a subtle interaction when Arrow 1 and Arrow 2 merge. The RMSNorm at Layer $l+1$ normalizes the sum:

$$\text{Denominator}_{l+1} = ||\gamma \mathbf{x}_0 + \Delta_l|| \tag{11}$$

Since $\gamma > 1$, the term $\gamma \mathbf{x}_0$ inflates the norm denominator. Because the learnable weights $(W)$ generating $\Delta_l$ have not yet updated to produce proportionally larger outputs, we observe a **dampening effect**:

$$\text{Effective Contribution} \propto \frac{\Delta_l}{||\mathbf{1.1x}_0 + \ldots||} < \frac{\Delta_l}{||\mathbf{1.0x}_0 + \ldots||} \tag{12}$$

This renders the model "stiff," as the static amplified embedding suppresses the contribution of the learned contextual updates.

## A.6 Gradient Amplification (Backward Pass)

Finally, we derive the gradient effect. Using the Jacobian $\mathcal{J}$ of the residual equation:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_{out}} \cdot \left( \gamma \mathbf{I} + \frac{\partial \mathcal{F}}{\partial \mathbf{x}} \right) \tag{13}$$

The term $\gamma \mathbf{I}$ scales the error signal by 1.1. This validates that the optimizer receives a mathematically amplified signal to update the connector embeddings, specifically targeting the vector $\mathbf{x}$ rather than the downstream parameters $\mathbf{k}$ or $\mathbf{W}$.