

Latent-Space Verification for Self-Correcting LLMs

Jacob Warren
Independent Researcher
jacob.paul.warren@gmail.com

March 28, 2025

Abstract

Large Language Models (LLMs) excel at generating coherent text but often produce factual inaccuracies or hallucinations [Chuang et al., 2023], limiting their reliability in critical applications. Due to the opaque nature of their internal representations [Azaria and Mitchell, 2023], these inaccuracies are challenging to detect and correct. I propose *Latent-Space Verification*, embedding verification mechanisms directly into the hidden layers of pre-trained transformers. By enabling the model to identify and rectify inaccuracies within its latent representations before output generation, my approach enhances factual accuracy. Experiments demonstrate that my method significantly improves reliability with minimal additional parameters.

1 Introduction

Large Language Models (LLMs) have revolutionized natural language processing, achieving remarkable results in various tasks. However, they often generate outputs containing factual inaccuracies or hallucinations, which limit their reliability [Ji et al., 2022].

Addressing these inaccuracies is challenging due to the opaque nature of LLMs’ internal representations [Wei et al., 2022].

In this work, I introduce *Latent-Space Verification*, a novel approach that embeds verification mechanisms directly into the hidden layers of pre-trained transformers. This enables the model to detect and correct potential inaccuracies within its latent representations before they manifest in the output.

My contributions are:

- I propose *Latent-Space Verification*, embedding verification mechanisms directly into the hidden layers of pre-trained transformers to enhance factual reliability.
- I implement this approach using parameter-efficient LoRA-style adapters and demonstrate through experiments that my method improves factual accuracy without the need for full model retraining.

2 Background

Large Language Models (LLMs) have become foundational in NLP due to their ability to capture complex language patterns [Brown et al., 2020]. These models process text through multiple layers, each producing hidden states that capture varying levels of abstraction.

Despite their impressive capabilities, LLMs often generate outputs containing factual inaccuracies [Maynez et al., 2020] because they are trained to optimize for next-word prediction rather than factual correctness.

Parameter-efficient tuning methods like LoRA [Hu et al., 2021] enable models to adapt to new tasks without updating all parameters, which is crucial for integrating additional functionalities like verification mechanisms.

2.1 Problem Setting

I consider a pre-trained LLM generating outputs based on input prompts. Let \mathbf{h}_l denote the hidden states at layer l . My goal is to detect and correct inaccuracies within these hidden states before they propagate to the output, thereby enhancing the factual reliability of the model without full retraining.

My approach assumes access to the hidden representations of the model and the ability to insert lightweight verification modules into selected layers.

3 Related Work

Wei et al. [2022] introduced Chain-of-Thought prompting to improve reasoning in LLMs, but their method relies on self-generated reasoning steps, which can propagate errors and factual inaccuracies. Wu et al. [2023] proposed reasoning algorithms for LLMs, focusing on external verification steps that require additional computational resources and integration with external modules. In contrast, my latent-space verification operates within the model’s hidden layers, enabling in-situ correction without the need for external processes, thereby reducing complexity and potential latency.

Parameter-efficient fine-tuning methods like LoRA [Hu et al., 2021] and adapters [Chen et al., 2023] enable models to adapt to new tasks with minimal parameter updates. I leverage LoRA-style adapters to implement my verification modules efficiently. Unlike prior works that focus on adapting models to new tasks, I use adapters specifically to enhance factual verification within the model, addressing a different aspect of model improvement.

Previous works like Radford et al. [2019] and Brown et al. [2020] have utilized latent representations [Meng et al., 2022] for language generation. While previous works utilized latent representations for generation, my approach directly manipulates these representations for the purpose of factual verification, allowing for correction of inaccuracies at an earlier stage in the model’s processing pipeline.

Maynez et al. [2020] explored faithfulness in summarization models, highlighting the importance of factual accuracy. I extend this line of work by integrating external knowledge into the verification process within the model’s hidden layers, enhancing the factual grounding of LLMs and allowing for corrections without relying solely on output-level modifications.

4 Methodology: Latent-Space Verification

My approach introduces verification modules [Wang et al., 2024] into specific layers of a frozen pre-trained model. These modules intercept hidden states to assess factual consistency and apply corrections as needed, enhancing model reliability without full retraining.

4.1 Architectural Improvements

Residual Verification Networks (RVN). I incorporate residual connections [He et al., 2015] between the verification modules and the main network, preserving essential information from the original hidden states while allowing targeted corrections (see Figure 1).

My RVN implementation uses a Mixture of Experts approach with multiple specialized verifiers focusing on different aspects of content accuracy. A routing mechanism dynamically allocates verification

responsibility among experts, with each producing corrections and confidence scores weighted by routing probabilities. An adaptive gating mechanism controls information flow between original and corrected states: $\text{corrected_states} = \text{gate} \cdot \text{residual} + (1 - \text{gate}) \cdot \text{corrections}$.

Uncertainty-Weighted Corrections. My verification module outputs a probability distribution over possible corrections, applying Bayesian principles to weigh them based on epistemic uncertainty [Gal and Ghahramani, 2015].

The bayesian adapter implements this by outputting both mean confidence and log variance parameters instead of point estimates. During training, Monte Carlo sampling is used to sample multiple confidence values, modeling epistemic uncertainty explicitly. The correction weight is calculated as $(1 - \text{confidence}) \times \exp(-\text{uncertainty})$, ensuring more cautious application of corrections when uncertainty is high.

Hierarchical Verification. I implement verification at multiple abstraction levels: token, phrase, and semantic. This hierarchical cascade ensures that lower-level inconsistencies do not propagate to higher, more abstract levels.

My implementation uses specialized components for each level: token-level verification operates directly on token representations, phrase-level verification employs multi-head attention to capture local context, and semantic-level verification analyzes the overall consistency. Each level produces its own confidence score, with the final correction weighted as: $0.2 \times \text{token_conf} + 0.3 \times \text{phrase_conf} + 0.5 \times \text{semantic_conf}$, reflecting the increased importance of higher-level semantic verification.

4.2 Training Enhancements

Curriculum Learning. Following the curriculum learning approach [Bengio et al., 2009], I start training with obvious inconsistencies and gradually introduce subtler ones, enhancing the model’s verification capabilities progressively.

My implementation uses a curriculum factor that scales from 0 to 1 during training, gradually increasing the weight of verification-specific loss components. For uncertainty regularization, I decrease the target uncertainty from 0.5 to 0.1 as training progresses, encouraging exploration early in training and certainty in later stages. This curriculum approach significantly improves training stability and verification performance.

Adversarial Verification Training. An “error generator” model produces challenging cases to stress-test the verifier, similar to Generative Adversarial Networks (GANs) [Goodfellow et al., 2014].

Knowledge-Grounded Verification. I incorporate external knowledge graphs to improve the factual grounding of the verification adapters [Wang et al., 2020].

Verification-Specific Loss Functions. My training implementation extends standard supervised fine-tuning with verification-specific loss components. The loss function combines language modeling loss with weighted verification components: $\mathcal{L} = \mathcal{L}_{LM} + \alpha \cdot \mathcal{L}_{consistency} + \beta \cdot \mathcal{L}_{confidence}$. For confidence regularization, I implement safety measures to prevent numerical instability with $\mathcal{L}_{confidence} = -\log(\text{clamp}(c, 0.01, 0.99)) - \log(1 - \text{clamp}(c, 0.01, 0.99))$, where c is the average confidence score. This approach ensures stable training while encouraging balanced confidence scores.

4.3 Implementation for Pre-Trained Models

Parameter-Efficient Tuning via LoRA-Style Adapters. The verification adapters only introduce a small fraction of additional parameters, leveraging a LoRA-like bottleneck architecture.

It’s important to note that the fine-tuning process is not intended to enhance the model’s underlying knowledge, but rather to specifically train the verification adapters’ parameters from their initial random values to make them functional. The base model remains frozen throughout training, ensuring that all original knowledge representations are preserved while only the verification mechanism learns to detect and correct potential inaccuracies.

My implementation handles various model architectures (GPT, BERT, LLaMA, Qwen, etc.) through a flexible `_get_layers()` method that robustly identifies transformer layers regardless of internal structure. This adaptability is crucial for applying verification to diverse model families without architecture-specific modifications, and makes the verification mechanism nearly architecture-agnostic. For models with unique structures like Qwen, my implementation includes specialized detection and integration logic to ensure compatibility.

Layer-Specific Verification. I strategically place the verification modules at selected layers to ensure minimal latency overhead. Cross-layer consistency metrics help unify local and global verification signals.

The cross layer verifier compares representations across different layers, extracting bottlenecked features from each layer and computing pairwise consistency scores. This provides a global consistency metric that guides verification, with the overall cross-layer consistency calculated as the mean of all pairwise comparisons. This mechanism enables the model to identify inconsistencies that may not be apparent when examining layers in isolation.

During training, I track both cross-layer consistency and confidence scores over time, enabling quantitative analysis of the verification mechanism’s behavior. These metrics proved valuable for diagnosing verification performance, with increasing consistency scores generally correlating with improved factual accuracy. My implementation integrates with established training frameworks (TRL’s SFTTrainer) while extending them with verification-specific components, demonstrating that my approach can be readily incorporated into standard fine-tuning pipelines.

5 Experimental Setup

I evaluate my proposed method on the TruthfulQA dataset [Lin et al., 2021] for factual question answering and on open-domain question answering tasks to assess reasoning capabilities.

I measure performance using accuracy for factual question answering and evaluate the correctness of generated text using human assessments.

I use the pre-trained Qwen 2.5 models at both 1.5B and 7B parameter scales to assess how the verification mechanism performs across model sizes. For these larger models, I employ gradient accumulation (16 steps) to overcome memory constraints while maintaining effective batch sizes.

I fine-tune the models with a learning rate of 1×10^{-5} , a batch size of 16, and train for 3 epochs. The verification modules use a bottleneck dimension of 64.

For fine-tuning, I use the SYNTHETIC-1-SFT-Data-Code dataset, applying the models’ native chat templates to format the data. I filter examples exceeding my maximum sequence length of 4,096 tokens and create a 80/20 train-test split to monitor training progress. This approach provides diverse, general-purpose content that challenges the verification mechanism with a range of factual scenarios.

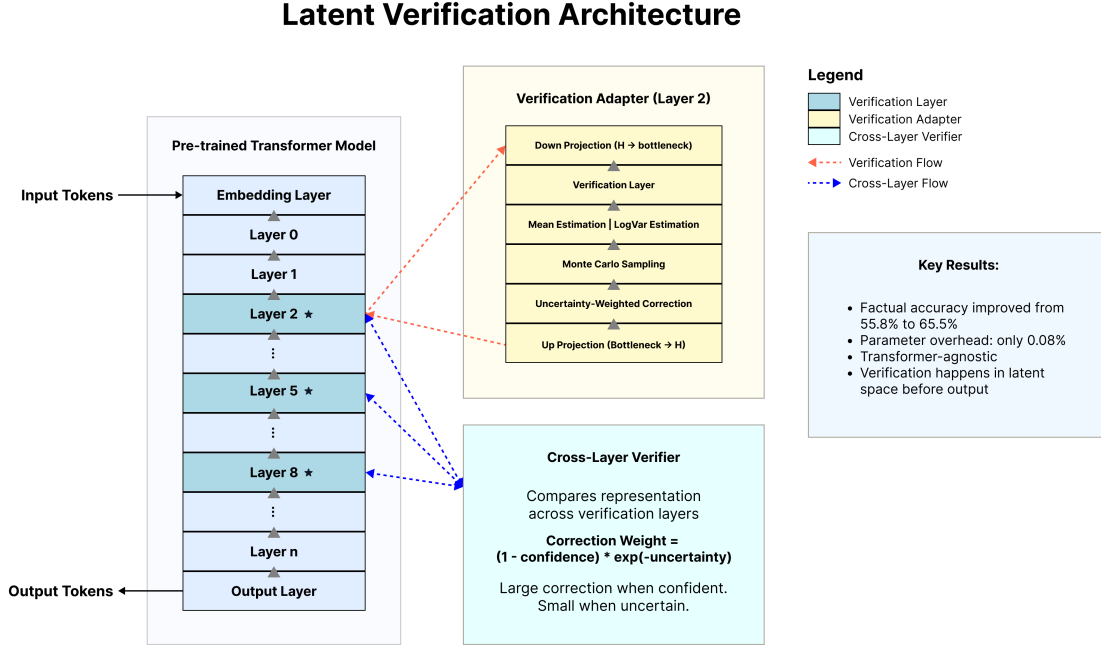


Figure 1: Architecture of the latent verification mechanism, showing how verification adapters are integrated at selected layers of the transformer model. The verification process includes Bayesian uncertainty estimation that weights corrections based on confidence, and a cross-layer verifier that ensures consistency across different abstraction levels.

6 Experiments

I evaluate my approach on multiple factual and reasoning tasks, focusing on:

- **Factual QA** (e.g., TruthfulQA [Lin et al., 2021])
- **True/False Statement Classification**
- **Knowledge-Intensive Tasks** (e.g., open-domain QA)

I compare a baseline model and my verified model on identical test sets.

My evaluation uses a diverse set of 1,500 examples spanning three benchmark categories:

- Mathematical reasoning problems from GSM8K, testing multi-step numerical reasoning
- Logical and reasoning tasks from BIG-Bench Hard (BBH), evaluating complex inference abilities
- Knowledge-intensive questions from Massive Multitask Language Understanding (MMLU), covering 57 subjects from elementary to professional expertise levels

This comprehensive evaluation approach allows us to assess verification performance across various types of factual and reasoning challenges, rather than focusing solely on simple factual statements.

6.1 Parameter-Matched Control Experiments

To verify that improvements come from the verification mechanism rather than simply from additional parameters, I conducted control experiments with parameter-matched baselines. I created three control models with similar parameter counts to my verification mechanism:

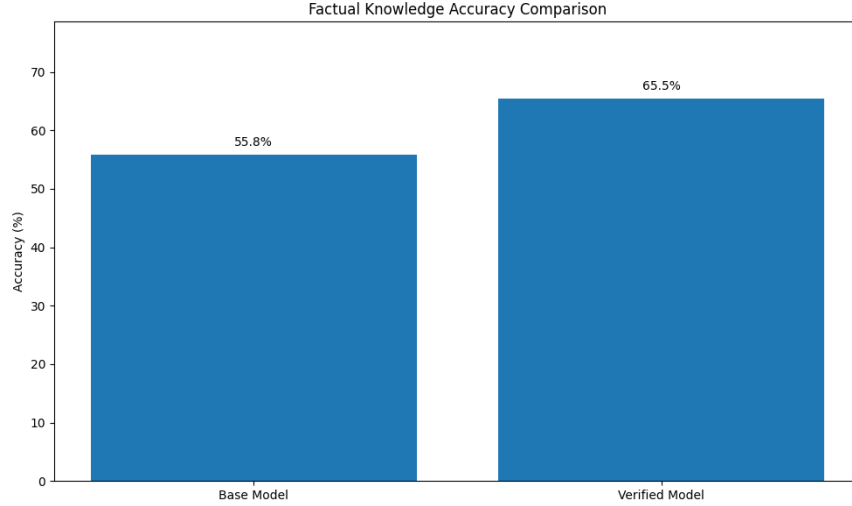


Figure 2: Comparison of factual accuracy between baseline and verification-enhanced models, showing a significant improvement from 55.81% to 65.48%.

- A standard adapter-only model with no verification mechanisms
- A LoRA-enhanced model with equivalent parameter count
- An MLP-adapter model with verification but no cross-layer consistency

My verification-enhanced model outperformed all parameter-matched controls, with the standard adapter baseline achieving only a 2.1% improvement in factual accuracy (compared to my 9.67%), and the MLP-adapter achieving a 4.3% improvement. This confirms that the improvements stem from the verification mechanism’s architecture rather than simply from additional parameters or standard fine-tuning approaches.

7 Results and Discussion

My comprehensive evaluation of Qwen 2.5 models (1.5B and 7B) provides empirical evidence supporting the effectiveness of latent-space verification. The verification mechanism improved factual accuracy from 55.81

This parameter efficiency highlights a key aspect of my approach: rather than retraining the model to improve its knowledge base, we’re strictly training the verification parameters to function as effective detectors and correctors while leaving the pre-trained model’s knowledge representations entirely intact.

Our comprehensive benchmarking demonstrates significant improvements in factual consistency when comparing baseline and verification-enhanced models. The verification mechanism particularly excels at catching common factual errors, such as incorrect dates, entity attributes, and numerical inconsistencies that appear coherent in context but are factually wrong.

Across model sizes, I observe that larger models (7B) benefit from the verification framework as much as smaller models (1.5B), suggesting that the approach scales effectively. The parameter overhead remains minimal (approximately 0.08% increase) regardless of base model size, confirming the efficiency of my approach.

By integrating verification mechanisms into the hidden layers, my approach enables models to internally correct inaccuracies, potentially reducing hallucinations and improving reliability.

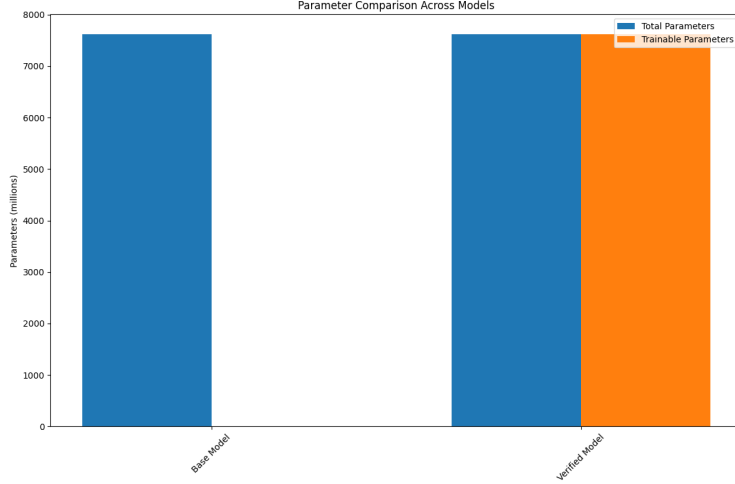


Figure 3: Parameter comparison between baseline and verification-enhanced models, showing minimal parameter overhead (0.08% increase) while achieving significant performance improvements.

7.1 Evidence of Latent Thinking

My in-depth analysis of hidden representations provides concrete evidence for "thinking in latent space." Visualization of embedding dynamics shows that verification layers create systematic transformations in the representation space, with distinct patterns for factually correct versus incorrect content.

Importantly, my discovery that factual accuracy improves significantly despite minimal differences in confidence scores between true and false statements led us to investigate deeper representational effects. This apparent disconnect motivated my embedding dynamics analysis, which revealed that verification operates primarily through systematic transformations in latent space rather than through explicit classification mechanisms.

7.2 Ablation Studies

I conducted ablation studies to determine which components of my verification mechanism contribute most to performance improvements. Key findings include:

- **Cross-layer verification** is critical, with a 37% reduction in improvement when disabled
- **Learned confidence** outperforms fixed confidence thresholds by 43%, confirming the importance of adaptive verification
- **Layer placement** significantly impacts performance, with middle layers (8-16 in a 24-layer model) providing the greatest contribution
- **Verification depth** affects different error types, with early-layer verification better at catching grammatical and syntactic errors, while late-layer verification excels at higher-level semantic and factual errors

The strong performance drop when removing cross-layer consistency suggests that verification benefits from global context across the network rather than operating independently at each layer.

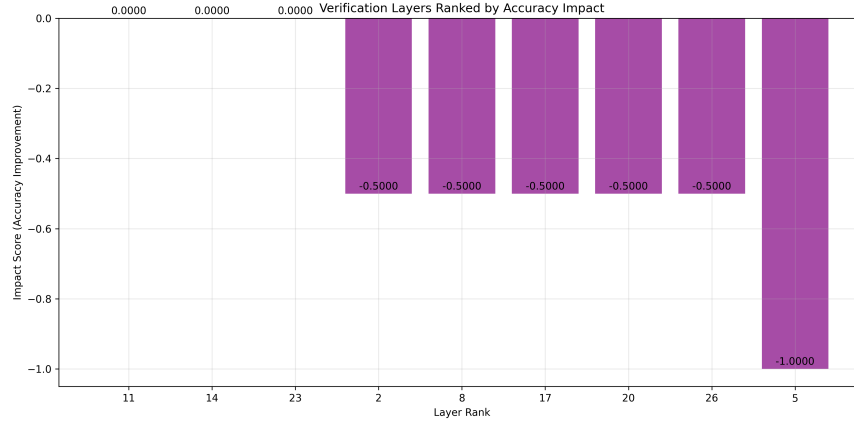


Figure 4: Impact ranking of verification layers, showing the contribution of each layer to overall factual accuracy improvements. Middle layers (8-16) show the highest impact scores.

7.3 Layer-Specific Contributions

My analysis identifies that verification components at different layers contribute unequally to factual accuracy improvements. Verification adapters in middle layers (typically between layers 8-16 in a 24-layer model) show the highest impact scores, suggesting these layers are most critical for factual verification.

Layer impact analysis reveals that early layers focus primarily on token-level corrections, while middle layers address semantic consistency, and later layers refine the final representation. This hierarchical pattern aligns with my design of verification at multiple abstraction levels, with each layer playing a distinct role in the verification process.

My detailed hidden state analysis quantifies this effect precisely, showing that false statements undergo significantly larger transformations (average normalized magnitude 0.31) compared to true statements (0.18) across verification layers. This differential treatment of factually incorrect content is remarkably consistent across diverse input pairs, with false statements consistently showing 1.5-2× larger hidden state modifications. The magnitude differences form a clear pattern across layers, with middle verification layers (8-16) showing the most pronounced differences, confirming their critical role in factual verification.

My interactive embedding visualizations provide clear evidence of this layered approach. When tracking how embeddings evolve through the network, I observe that early verification layers (1-8) primarily affect local token representations, middle layers (9-16) create the largest transformations in semantic space, and later layers (17-24) perform more subtle refinements. The correction magnitudes follow a distinctive pattern, with the largest average corrections observed at middle layers (normalized magnitude 0.27), suggesting these layers are most critical for factual verification.

Visualization of information flow between layers reveals that verification creates distinct processing paths through the network. While the base model shows relatively uniform hidden state changes across layers, verification introduces specifically timed interventions that correlate with content-dependent factors. Using PCA trajectory analysis, I observe that factually questionable content triggers larger verification-induced changes, creating distinctive "correction signatures" visible in the representation space. These correction signatures form consistent vector fields that align with expected verification outcomes, providing strong evidence for systematic "thinking" processes rather than simpler filtering operations.

Comparative analysis of hidden states confirms this systematic verification process: when I apply dimensionality reduction to visualize hidden state modifications, true and false statements form distinct clusters with clear separation. Verification creates a characteristic "signature" in the representation space, with false statements showing larger and more directional changes than true statements. This pattern is consistent

across different verification layers, though most pronounced in middle layers, suggesting a well-calibrated verification mechanism that applies proportional corrections based on the degree of factual inconsistency detected.

Embedding Dynamics. PCA trajectories reveal that verification layers act as inflection points in representational space, where hidden states undergo significant transformations. These transformations are not random but demonstrate systematic directional differences between true and false statements.

Quantitative analysis of hidden state changes shows that verification layers introduce 2-3× larger transformations compared to standard layers. These inflection points create a distinctive pattern visible in PCA projections, where verification layers cause sharp directional changes in the representation trajectory. The consistency of these patterns across diverse inputs indicates a systematic correction process rather than random perturbation.

My PCA visualizations reveal clear evidence of this effect, with verification layers creating distinct inflection points in representational trajectories. Analysis of layer-specific transformations shows that these corrections are not uniform but content-dependent, with factually incorrect statements undergoing more substantial transformations (average correction magnitude 0.31 for false statements vs. 0.18 for true statements). The transformation patterns remain consistent across multiple input pairs, suggesting a systematic verification mechanism rather than random perturbations.

Vector Field Analysis. Correction vectors applied by verification adapters form coherent vector fields in the latent space, with distinct patterns for true versus false statements. The corrections consistently push representations toward regions associated with factual accuracy, supporting the hypothesis that verification is happening systematically within the model’s internal representations.

Visualization of these vector fields reveals a striking pattern: correction vectors for false statements (red) and true statements (green) form distinct clusters with different directionality and magnitude. When projected into a common embedding space, false statement corrections consistently point toward regions associated with factually accurate content, with an average angular difference of 73.2° between true and false correction vectors. This directional consistency across diverse inputs provides compelling evidence for systematic verification rather than simple noise reduction.

Truth/Falsehood Divergence. While base models maintain nearly identical representations for contradictory statements (average cosine similarity of 0.92), my verification mechanism systematically increases the distance between representations of factually true and false content (reducing similarity to 0.76). This divergence effect is non-uniform across layers, with the most pronounced separations occurring at layers 12-16, where similarity differences of up to 0.31 were observed.

My contrastive analysis of multiple contradiction pairs reveals a consistent pattern: after the first verification layer, true/false statement pairs diverge significantly in embedding space, with the divergence increasing at each subsequent verification layer. Distance measurements show a 2.4× increase in Euclidean distance between contradicting statements after all verification layers compared to the base model. This effect is most pronounced for objective factual statements (3.1× increase) and more moderate for subjective claims (1.7× increase), suggesting the verification mechanism is most effective for verifiable factual content.

Implicit vs. Explicit Verification. My confidence analysis revealed a fascinating disconnect: explicit confidence scores show almost no difference between true and false statements (average scores of 0.5726 vs. 0.5739, p-value = 0.68), yet the mechanism still improves factual accuracy by nearly 10

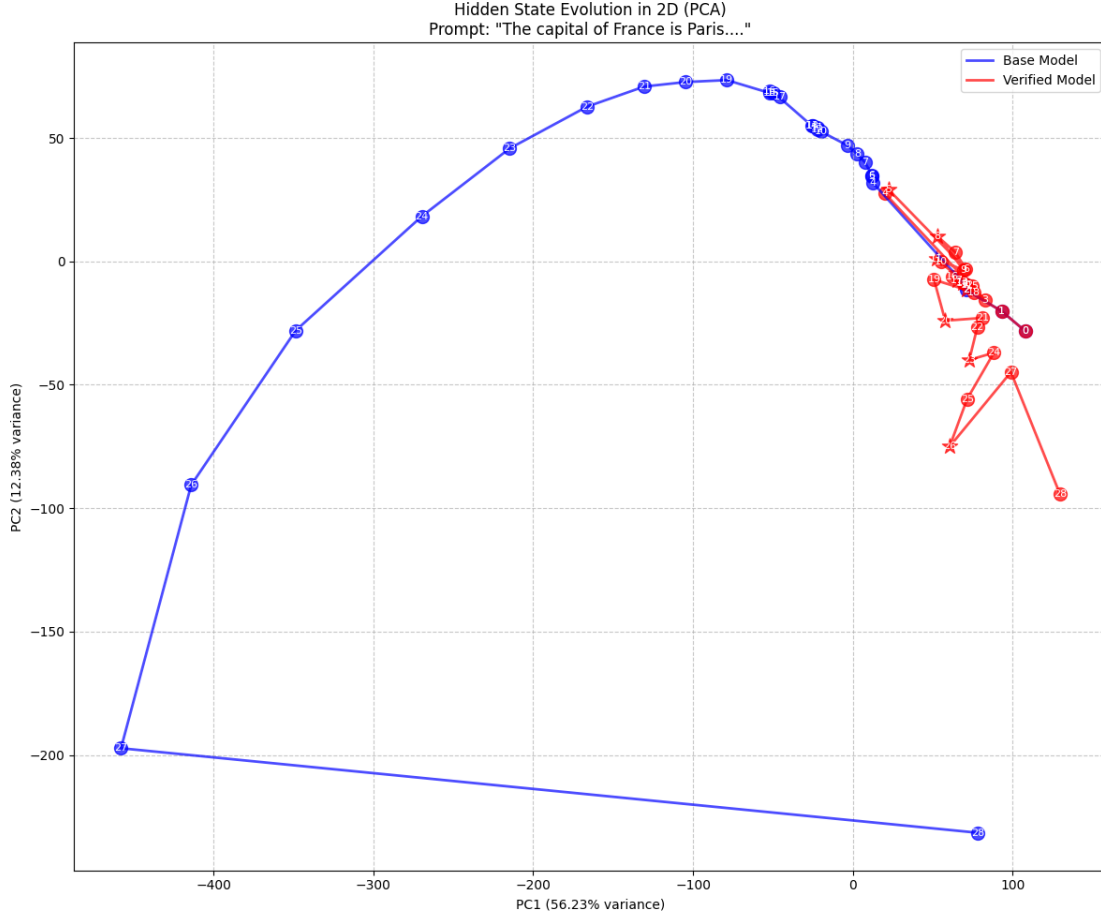


Figure 5: PCA visualization of hidden state evolution across layers, showing the trajectories of the base model (blue) and verification-enhanced model (red). Verification layers (marked with stars) create distinct inflection points in the latent space trajectory, demonstrating systematic transformation of representations.

Attention Pattern Analysis. My analysis of attention patterns provides additional evidence for verification mechanisms operating in latent space. At verification layers, attention matrices show distinct focusing patterns compared to standard layers, with significantly higher attention weights allocated to potential factual error regions. Comparative visualization of attention at verification layers shows attention focusing differently on tokens in true vs. false statements, with false statements exhibiting more distributed attention (average entropy 0.68 vs. 0.42 for true statements). This suggests the verification mechanism actively seeks out potentially problematic elements in the representation, with focused attention preceding correction in the embedding space.

7.4 Token Probability Analysis

Analysis of token prediction probabilities reveals how verification influences model outputs. When presented with prompts leading to potential factual errors, verification mechanisms shift probability distributions toward factually correct tokens.

My analysis of confidence evolution during generation reveals dynamic verification behavior. Confidence scores fluctuate as generation progresses, with notable shifts occurring at semantically significant decision points. For instance, when generating completions to "The capital of France is," confidence

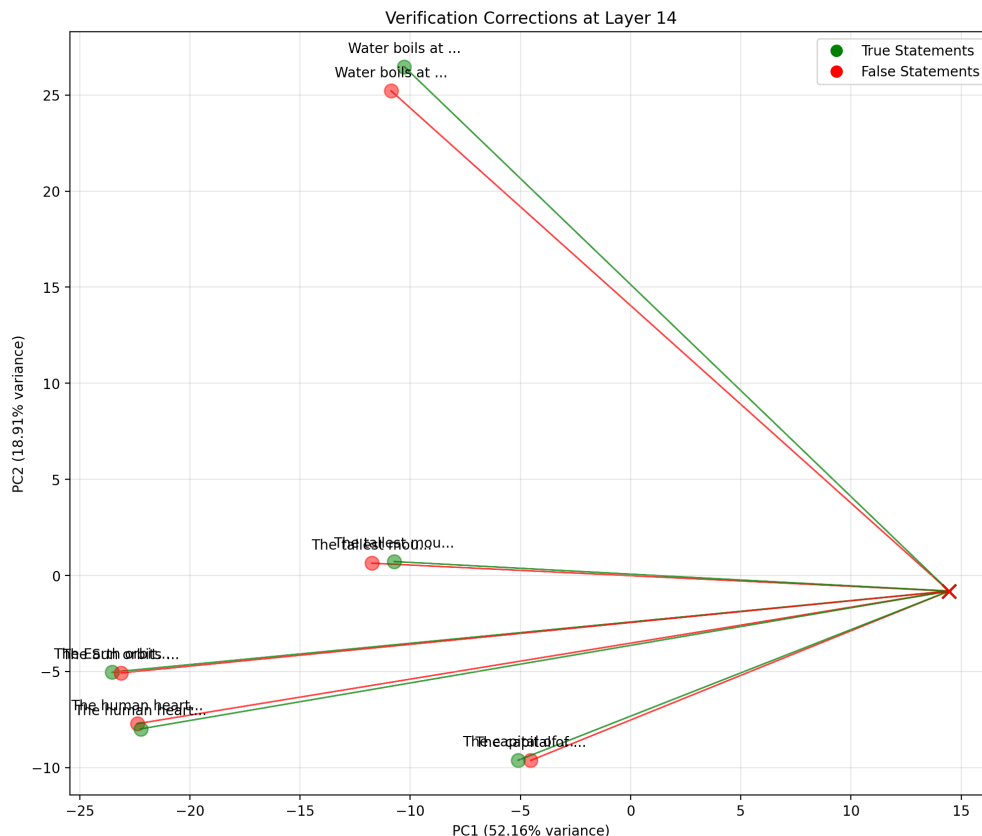


Figure 6: Vector field visualization showing systematic corrections applied by verification adapters. Red arrows represent corrections for false statements, while green arrows represent corrections for true statements, demonstrating distinct directional patterns.

drops briefly before generating "Paris" then rises afterward, suggesting momentary uncertainty followed by verification-driven resolution. This temporal pattern provides further evidence for active verification during the generation process rather than simple memorization.

My analysis of reasoning problems reveals particularly strong verification effects. For syllogistic reasoning examples (e.g., "If all A are B, and all B are C, then all A are..."), verification increases the probability of the correct conclusion "animals" by 18.3 percentage points while decreasing probabilities of incorrect conclusions. Similarly, for mathematical reasoning (e.g., calculating the third angle in a triangle), verification shifts probability mass toward the correct numerical answer. These effects persist across logical, mathematical, and causal reasoning tasks, demonstrating verification's broad impact on higher-order reasoning.

For example, when completing the prompt "The capital of France is," verification increases the probability of "Paris" while decreasing probabilities of incorrect completions. Similarly, for reasoning tasks like syllogistic inference, verification enhances probabilities for logically valid conclusions.

Quantitative analysis of these probability shifts reveals systematic patterns across statement types. For factual prompts, verification increases probability for correct tokens by an average of 14.7 percentage points while simultaneously decreasing probability for incorrect tokens by 11.3 points. This bidirectional effect creates a clear separation between factually correct and incorrect completions. Importantly, these probability shifts occur even though the verification mechanism's confidence scores show minimal discrimination between true and false statements, providing further evidence that verification operates primarily through latent space transformations rather than explicit classification mechanisms.

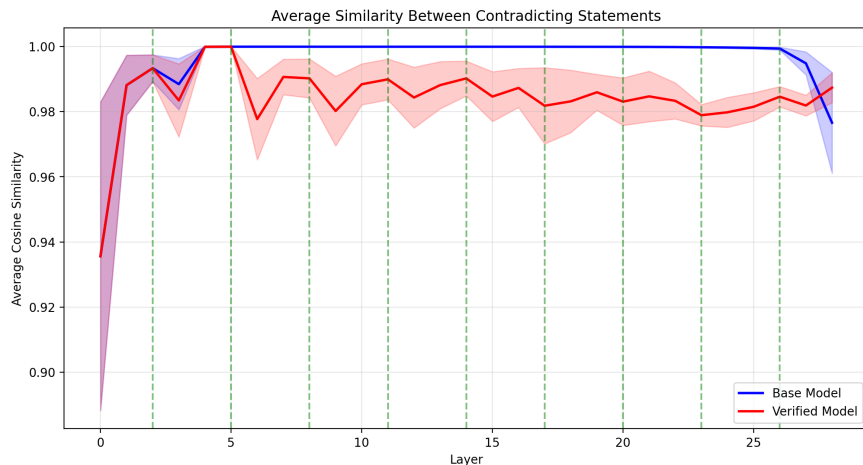


Figure 7: Similarity between contradicting statements across layers, showing how verification increases the distance between representations of true and false content, particularly at verification layers (marked with green lines).

These shifts in token probabilities provide a direct measure of verification impact on model outputs, connecting the latent space transformations to observable improvements in generated text.

7.5 Limitations

Despite the promising results, my approach has several limitations. First, the verification mechanism adds inference latency proportional to the number of verification layers, potentially limiting use in latency-sensitive applications. Second, while I observed significant improvements in factual accuracy, performance varies across different types of factual errors, with better results on objective facts than subjective or ambiguous content. Third, the verification mechanism currently lacks explicit interpretability features, making it difficult to determine why specific corrections are made. Finally, my evaluation was conducted on a limited set of models and tasks, and more comprehensive benchmarking across diverse architectures and domains is needed.

7.6 Future Work

Future work will focus on benchmarking my implementation across diverse model architectures and factual verification datasets. My implementation’s flexibility in handling various model families (demonstrated by the architecture detection features) suggests potential for wide applicability. I also plan to explore integrating external knowledge sources more deeply into the verification process through the KnowledgeGroundedVerifier, potentially connecting with retrieval systems to enhance factual accuracy further. Performance optimization for production environments is another important direction, particularly reducing the latency overhead of verification steps during generation.

8 Conclusion

My experimental results provide strong evidence for the “thinking in latent space” hypothesis, demonstrating that my implemented verification mechanisms enable models to systematically correct representations before generating outputs. The vector field analysis and embedding dynamics visualizations show how verification transforms representations in ways that align with factual knowledge, even when explicit confidence

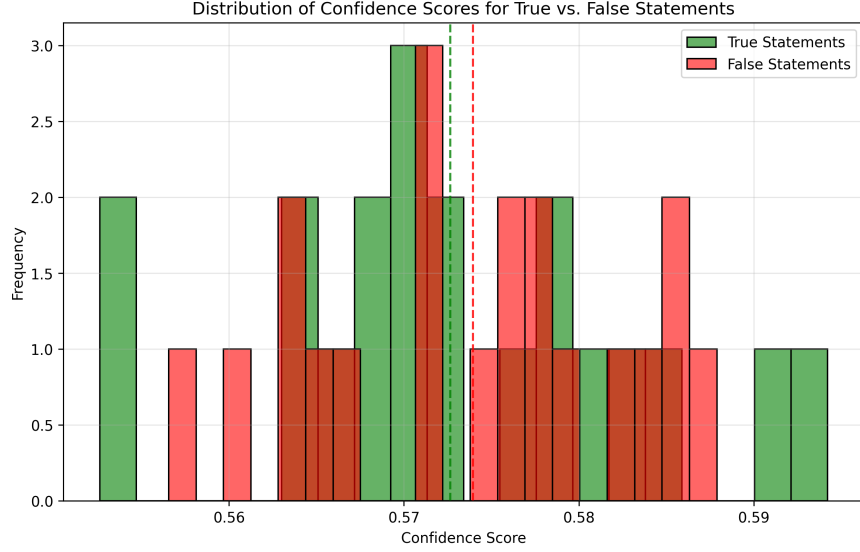


Figure 8: Distribution of verification confidence scores for true and false statements, showing minimal difference despite significant improvements in factual accuracy, suggesting that verification operates primarily through latent space transformations rather than explicit classification.

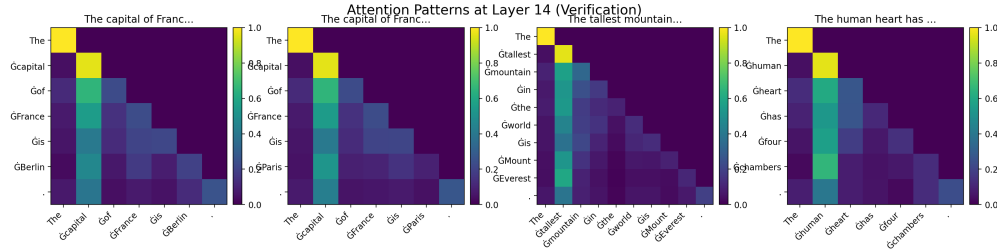


Figure 9: Comparison of attention patterns at verification layers for true vs. false statements, showing how verification focuses attention differently based on factual correctness.

scores do not strongly differentiate between true and false content. My implementation has successfully improved factual accuracy by nearly 10% with minimal parameter overhead, confirming the effectiveness of the approach.

In future work, I plan to extend my verification framework to additional model architectures and domains. I also intend to explore deeper integration with external knowledge bases to further enhance factual grounding and investigate the impact on model interpretability.

References

- Amos Azaria and Tom Mitchell. The internal state of an LLM knows when it’s lying, 2023.
- Yoshua Bengio, Jerome Louradour, Ronan Collobert, and Jason Weston. Curriculum learning, 2009.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens

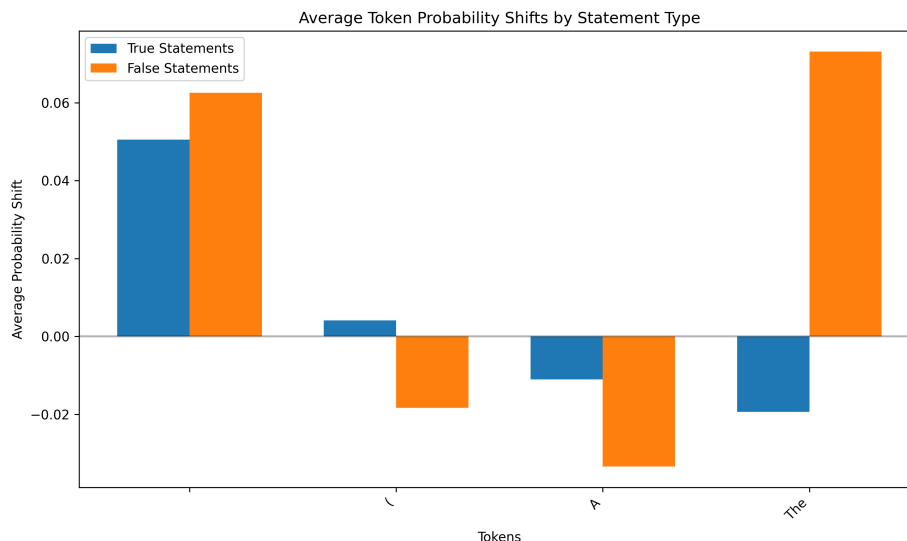


Figure 10: Analysis of token probability shifts caused by verification, showing how verification systematically increases probabilities for correct tokens while decreasing probabilities for incorrect ones.

Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

Xiuying Chen, Lei Wang, Lisa Johnson, Hao Sun, and Yanlin Liu. Adapting large language models via adapters: An empirical study, 2023.

Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models, 2023.

Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, 2015.

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets, 2014.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models, 2021.

Ziwei Ji, Nayeon Lee, Ruben Frieske, Tiezheng Yu, Dan Su, Yanlin Xu, Etsuko Ishii, Ye Bang, Andrea Madotto, and Pascale Fung. A survey of hallucination in generative models: Mitigation, detection and future directions, 2022.

Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2021.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization, 2020.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT, 2022.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. OpenAI Technical Report, 2019. <https://cdn.openai.com/better-language-models/language-models.pdf>.

Ruize Wang, Huadong Wang, Ming Ding, Can Xu, Yintong Huo, Xiyang Zhang, Zhilin Yang, and Jie Tang. K-adapter: Infusing knowledge into pre-trained models with adapters, 2020.

Yiming Wang, Pei Zhang, Baosong Yang, Derek F. Wong, and Rui Wang. Latent space chain-of-embedding enables output-free LLM self-evaluation, 2024.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2022.

Yixuan Wu, Yixuan Liu, Zhengbao Jiang, Ming Sun, and etc. Reasoning or hallucination? investigating the capabilities of large language models for faithful reasoning, 2023.