

LatentGate: Low-Latency Semantic Routing via Frozen-Backbone Probing of Small Language Models

Anonymous ACL submission

Abstract

As Multi-Agent Systems (MAS) scale to hundreds of specialized agents, the routing layer becomes a critical bottleneck. Traditional approaches force a stark trade-off: prompt-based LLM routers deliver high semantic reasoning but incur prohibitive latency ($\sim 1500\text{--}2000\text{ms}$) and cost that grows with agent count, while embedding-based routers operate at low latency (25–50ms on a T4 for cosine/centroid-style routing) but fail to capture nuanced functional intent, collapsing semantically similar but functionally distinct agents. We identify *representation anisotropy*, the geometric collapse of hidden-state vectors into a narrow cone, as a key mechanism underlying this embedding-based routing failure. We propose **LatentGate**, a non-generative routing architecture that extracts mean-pooled hidden states from a frozen small language model (SLM), applies PCA-whitening (decorrelation + variance normalization) to resolve the anisotropy, and trains a lightweight linear probe for agent classification. Experiments across 5 SLM backbones and 100 enterprise agents show that LatentGate achieves 98.8% in-domain accuracy and 80.0% out-of-distribution accuracy on natural queries, outperforming embedding-based routers by 13–22 absolute points. LatentGate operates at $\sim 28\text{ms}$ on a T4 GPU; the SLM forward pass is independent of agent count, with classification adding an $O(Ck)$ term (agents C , whitened dimension k) that is negligible at $C = 100$ and small relative to the SLM forward pass. The lightweight linear probe additionally enables sub-10ms warm-start retraining from user feedback, offering a path toward self-healing routing in production. We further benchmark prompt-based routing with GPT-4.1, GPT-4.1-nano, and Gemini 2.5 Flash, demonstrating that these degrade to 70–77% accuracy at 100 agents while incurring 1500–2000ms latency, confirming the need for non-generative alternatives. We categorize this work as *Emerging*, as it introduces a new routing primitive rather than

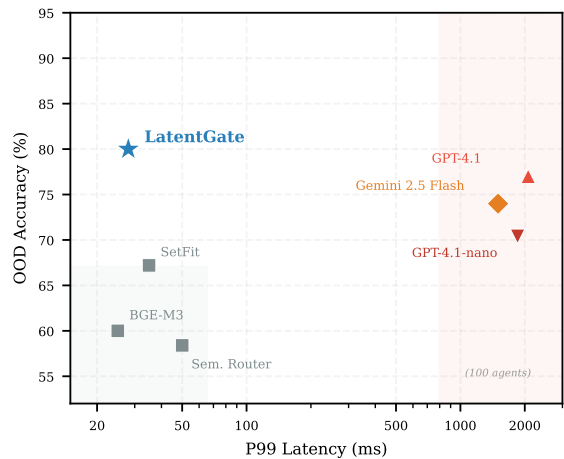


Figure 1: Accuracy-latency trade-off for 100 enterprise agents. LatentGate (blue) achieves 80% OOD accuracy at $\sim 28\text{ms}$, occupying the Pareto-optimal region.

reporting a completed long-term deployment study.

046
047

1 Introduction

048

The deployment of Generative AI in enterprise settings has shifted from monolithic chat interfaces to modular Multi-Agent Systems (MAS), where a central orchestration layer delegates user queries to specialized agents. In these architectures, a router must map natural language inputs to the correct agent among potentially hundreds of candidates. For example, directing “I can’t push because of a protected branch error” to a GitHub DevOps agent rather than an IT Helpdesk agent. In industrial applications, this router must simultaneously satisfy three constraints: (1) **Low Latency**: sub-100ms routing to maintain a conversational user experience; (2) **High Precision**: routing errors trigger incorrect API calls, causing user frustration and potential data corruption; and (3) **Low-Cost Onboarding**: the system must accommodate new

049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065

066	agents without expensive backbone fine-tuning.	
067	Current approaches fail to meet all three simul-	
068	taneously. Prompt-based LLM routing provides	
069	reasonable precision at small agent counts but in-	
070	currurs latency and cost that grow with agent count	
071	due to prompt expansion; as we show in Sec-	
072	tion 6, accuracy degrades sharply as agents scale.	
073	Embedding-based routing approaches, including	
074	Semantic Router [Aurelio AI, 2024], rely on cosine	
075	similarity of pre-trained embeddings, but collapse	
076	distinct operational roles sharing vocabulary (e.g.,	
077	“reset password” vs. “reset branch”).	
078	We trace this failure to a well-documented	
079	but underexploited phenomenon: <i>representation</i>	
080	<i>anisotropy</i> [Ethayarajh, 2019]. The hidden states	
081	of transformer models cluster in a narrow cone in	
082	representation space, causing cosine similarity to	
083	degenerate. While this problem has been exten-	
084	sively studied for sentence embedding quality [Su	
085	et al., 2021, Huang et al., 2021], it has not, to our	
086	knowledge, been diagnosed as a key mechanism	
087	driving routing failure in multi-agent systems.	
088	We propose LatentGate , a non-generative rout-	
089	ing architecture that resolves this through a three-	
090	stage pipeline: (1) extracting hidden states from	
091	a frozen SLM backbone, (2) mean pooling across	
092	tokens, and (3) applying PCA-whitening followed	
093	by linear probing. This converts a frozen SLM into	
094	a high-fidelity semantic router without backbone	
095	modification (Figure 1).	
096	Our contributions are as follows:	
097	• Diagnosis. We identify representation	
098	anisotropy as a key mechanism driving se-	
099	matic routing failure in MAS and quantify	
100	its effect across 5 SLM families (Section 3).	
101	• Architecture. We propose LatentGate, a	
102	three-stage pipeline whose inference is dom-	
103	inated by the frozen SLM forward pass (in-	
104	dependent of C), with a lightweight $O(Ck)$	
105	probe amenable to online updates (Section 4).	
106	• Empirical Validation. We demonstrate 80%	
107	OOD accuracy at ~ 28 ms across 100 agents	
108	and 5 SLM backbones, matching or exceed-	
109	ing GPT-4.1, GPT-4.1-nano, and Gemini 2.5	
110	Flash while being $74\times$ faster (Section 6).	
111	LatentGate is currently being evaluated in an	
112	internal enterprise assistant environment (details	
113	withheld for anonymity); preliminary observations	
114	suggest stable latency, though this does not consti-	
115	tute a controlled deployment study.	
	2 Related Work	116
	LLM Routing. Recent work has explored rout-	117
	ing across multiple LLMs to optimize cost-quality	118
	trade-offs. RouteLLM [Ong et al., 2025] learns	119
	routing policies from preference data, while Fru-	120
	galGPT [Chen et al., 2023] and HybridLLM [Ding	121
	et al., 2024] cascade models by difficulty. Au-	122
	toMix [Aggarwal et al., 2025] routes based on self-	123
	verification. These methods route between 2–4	124
	LLMs of varying capability; LatentGate addresses	125
	the orthogonal problem of routing to 100+ special-	126
	ized <i>agents</i> at sub-100ms latency.	127
	Multi-Agent Routing. RouterBench [Hu et al.,	128
	2024] benchmarks LLM selection but does not ad-	129
	dress agent-scale routing. MasRouter [Yue et al.,	130
	2025] and Tool-to-Agent Retrieval [Lumer et al.,	131
	2025] explore agent routing via retrieval but do	132
	not diagnose or resolve the underlying represen-	133
	tation geometry problem. ToolACE-MCP [Yao	134
	et al., 2026] generalizes tool routing with conver-	135
	sation history but relies on LLM inference. We	136
	are not aware of prior work that explicitly connects	137
	anisotropy to agent-scale routing failure and evalu-	138
	ates whitening as a mitigation in this setting.	139
	Representation Anisotropy. Ethayarajh [2019]	140
	first documented the anisotropy of contextualized	141
	representations. Subsequent work proposed whiten-	142
	ing for sentence similarity [Su et al., 2021, Huang	143
	et al., 2021, Zhuo et al., 2023], though Forooghi	144
	et al. [2024] cautioned against whitening for classi-	145
	fication in LLMs. Prior whitening work targets	146
	sentence similarity and retrieval, not high-class	147
	($C = 100$) agent routing with frozen causal SLM	148
	backbones. We show that whitening <i>is</i> effective for	149
	classification in this setting, precisely because rout-	150
	ing requires discriminating fine-grained semantic	151
	roles that anisotropy collapses.	152
	3 The Representation Cone Problem	153
	3.1 Background: Anisotropy in Neural	154
	Representations	155
	As discussed in Section 2, contextualized repre-	156
	sentations occupy a narrow cone in vector space	157
	[Ethayarajh, 2019, Gao et al., 2019, Godey et al.,	158
	2024]. While whitening has improved sentence	159
	embedding quality, it has not been applied to multi-	160
	agent routing.	161

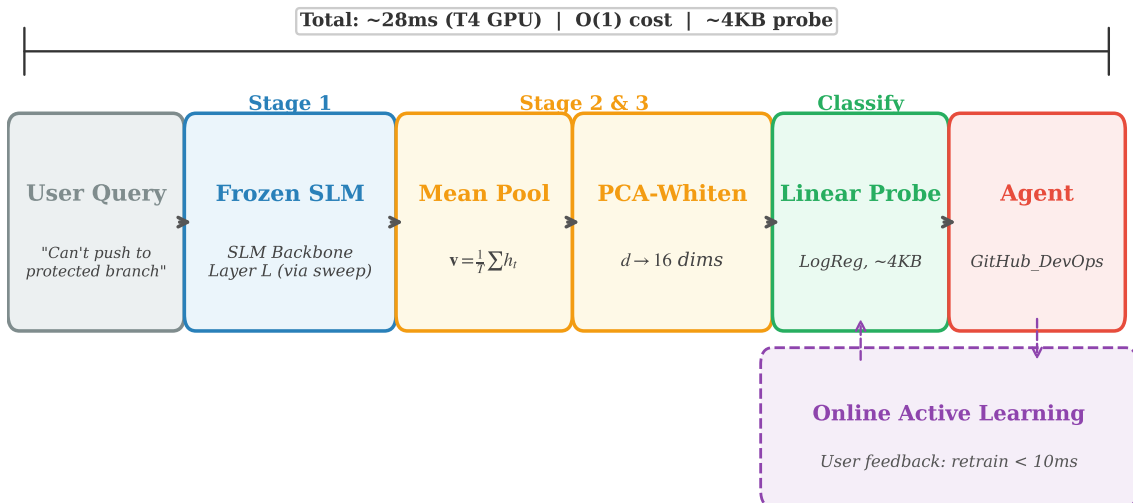


Figure 2: The LatentGate pipeline. A user query passes through a frozen SLM (Stage 1), is mean-pooled (Stage 2), whitened via PCA, and classified by a linear probe (Stage 3). Online active learning enables <10ms retraining.

Representation	Intra	Inter	Gap
Raw Hidden States (L7)	0.035	0.029	0.005
BGE-M3 Embeddings	0.81	0.68	0.13
After PCA-Whitening	0.230	-0.026	0.255

Table 1: Cosine similarity analysis on 100-agent queries from Qwen-2.5-1.5B-Instruct. Intra/Inter denote mean cosine similarity of query pairs within the same agent vs. across different agents. Raw hidden states exhibit extreme anisotropy (gap = 0.005), rendering cosine-based routing nearly non-discriminative. PCA-whitening yields a 48 \times improvement in discriminative power.

3.2 Empirical Diagnosis

We hypothesize that the same geometric phenomenon causes routing failure. We extract hidden states from a selected interior layer of Qwen-2.5-1.5B-Instruct (Layer 7, chosen via cross-validated sweep; see Appendix D) for a set of 100 agent-specific queries and measure intra-class vs. inter-class cosine similarity. In an ideal routing space, intra-class similarity should be high and inter-class similarity should be low, yielding a large discriminative gap.

Table 1 confirms our hypothesis. The raw hidden states exhibit extreme anisotropy: the gap between intra-class and inter-class similarity is only 0.005, rendering cosine-based routing nearly non-discriminative. After whitening, this gap widens to 0.255, a 48 \times improvement. BGE-M3 embeddings [Chen et al., 2025], while better than raw hidden states, still show a gap of only 0.13, explaining

their observed 60% OOD accuracy.

3.3 Intrinsic Dimensionality Analysis

We analyze the eigenvalue spectrum of the covariance matrix. For raw hidden states ($d = 1536$), the top 3 principal components explain $\sim 30\%$ of variance. While less extreme than encoder-only models [Ethayarajh, 2019], where top PCs can dominate over 50%, even this moderate spectral concentration suffices to collapse cosine similarity in high-class routing ($C = 100$), as evidenced by the near-zero gap (0.005) in Table 1. After whitening to $k = 16$, variance distributes uniformly and the cosine gap improves 48 \times (Figure 3).

4 The LatentGate Architecture

LatentGate applies a three-stage transformation to convert a user query x into a routing decision over C agents. The architecture is illustrated in Figure 2. Critically, the SLM backbone remains frozen throughout; only a lightweight probe head ($\sim 3\text{--}7\text{KB}$ depending on precision) is learned, with PCA and scaler parameters fixed as preprocessing artifacts.

4.1 Stage 1: Deep Semantic Extraction

We feed the tokenized query x into a frozen SLM backbone and extract hidden states $\mathbf{H}^L \in \mathbb{R}^{T \times d}$ from an interior layer L , selected via cross-validated sweep. Across all 5 backbones, the optimal layer falls in the early-to-middle range (L7–L10; Appendix D), consistent with prior probing

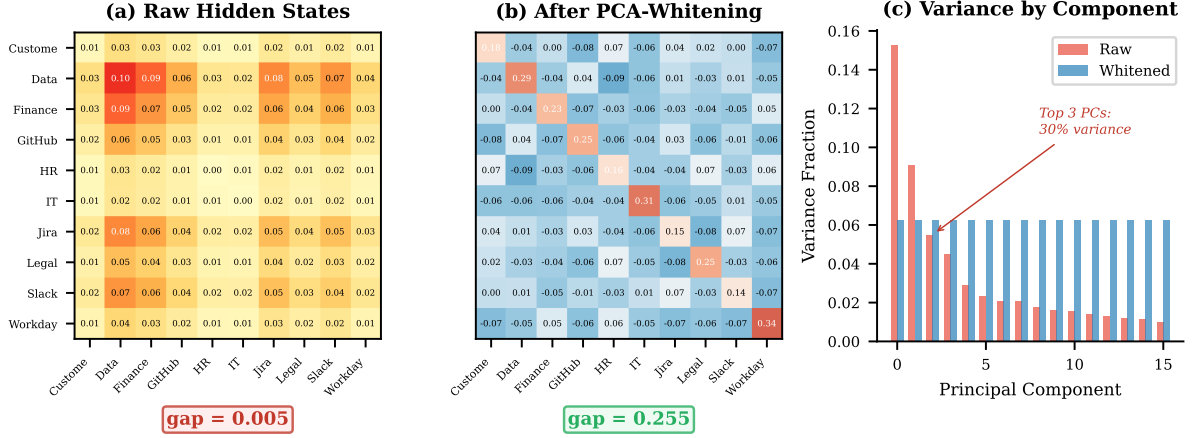


Figure 3: Agent-level cosine similarity matrices. (a) Raw hidden states: uniform similarity (gap = 0.005). (b) After PCA-whitening: clear block-diagonal (gap = 0.255). (c) Raw variance concentrated in top 3 PCs ($\sim 30\%$); whitening equalizes variance.

literature [Belinkov, 2022].

4.2 Stage 2: Mean Pooling

We compute the mean-pooled vector:

$$\mathbf{v}_{\text{raw}} = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t^L \quad (1)$$

where T is the sequence length. Mean pooling outperforms [CLS], last-token, and max pooling by 3–8% on OOD queries (Table 4).

4.3 Stage 3: Whitening and Probing

The mean-pooled vector \mathbf{v}_{raw} resides in the anisotropic cone identified in Section 3. We apply a three-step transformation:

$$\hat{y} = \text{Softmax}(\mathbf{W} \mathbf{z} + \mathbf{b}) \quad (2)$$

$$\mathbf{z} = \text{PCA}_k(\text{Scale}(\mathbf{v}_{\text{raw}})) \quad (3)$$

where $\text{Scale}(\cdot)$ centers and variance-normalizes each dimension (StandardScaler), $\text{PCA}_k(\cdot)$ reduces dimensionality from $d = 1536$ to $k = 16$ and whitens the projected features (decorrelation + variance normalization), and $\mathbf{W} \in \mathbb{R}^{C \times k}$, $\mathbf{b} \in \mathbb{R}^C$ are linear probe weights learned via Logistic Regression on the 5,000 generated-and-verified training queries. For $C = 100$ and $k = 16$, the probe has 1,700 trainable parameters ($\sim 3\text{--}7\text{KB}$ depending on precision); the PCA projection and scaler statistics are fixed preprocessing artifacts. The probe is trained once per agent set; adding new agents requires only generating queries via the same protocol and re-fitting the probe quickly (offline), without modifying the backbone.

4.4 Online Active Learning

A practical advantage of the linear probe is online error correction: when a misroute is detected via user feedback, the logistic regression can be warm-started with the corrected example. In preliminary trials, this retrain completes in $< 10\text{ms}$, suggesting a path toward self-healing routing without backbone modification. Rigorous evaluation of this loop is left to future work.

5 Experimental Setup

5.1 Dataset

We construct a routing benchmark of 100 specialized agents spanning IT, DevOps, HR, Finance, Legal, Engineering, and other domains. Each agent has 50 training queries generated via few-shot prompting of GPT-4.1, seeded with exemplars from 1,000 anonymized production queries (used only as in-context demonstrations, excluded from training and evaluation to prevent leakage). All queries were verified by two annotators ($\kappa = 0.91$), yielding 5,000 training examples. For evaluation, we generate 20 OOD queries per agent (2,000 total) using a separate prompt encouraging paraphrasing and cross-domain vocabulary overlap, verified under the same protocol.

5.2 Baselines

We compare LatentGate against two categories of baselines:

Encoder / Embedding-Based Routers. (1) **BGE-M3** [Chen et al., 2025]: cosine similarity retrieval using BAAI’s state-of-the-art multilingual

embedding model. (2) **Semantic Router** [Aurelio AI, 2024]: a production routing library using OpenAI embeddings with centroid-based classification. (3) **SetFit** [Tunstall et al., 2022]: few-shot contrastive fine-tuning of a sentence transformer.

LLM Prompt-Based Routers. (4) **GPT-4.1**: zero-shot classification via structured prompt containing all agent descriptions. (5) **GPT-4.1-nano**: lightweight variant for cost-sensitive deployments. (6) **Gemini 2.5 Flash**: Google’s low-latency model with structured output. All routers received identical, uncompressed agent descriptions to isolate scaling behavior under a common production baseline. We do not evaluate hierarchical or compressed prompt routers, as their design space (tree construction, pruning, description compression) introduces additional tuning choices not comparable to our frozen-backbone setting.

5.3 LatentGate Configurations

We evaluate LatentGate across 5 SLM backbones: Qwen-2.5-{0.5B, 1.5B, 3B}-Instruct, Phi-3.5-mini-instruct (3.8B), and Llama-3.2-1B-Instruct. For each backbone, we select the best layer L via 5-fold cross-validated probe accuracy sweep. Whitening reduces dimensionality to $k = 16$ components. The linear probe is an L2-regularized logistic regression trained on the whitened features.

5.4 Metrics

We report: (1) **In-domain accuracy** on held-out training-distribution queries. (2) **OOD accuracy** on the paraphrased evaluation set (primary metric). (3) **Latency**: end-to-end wall-clock time on a single T4 GPU (batch size 1; Table 8); LLM latencies are API response times under identical network conditions. (4) **Probe size**: trainable parameter count.

6 Results and Discussion

6.1 Main Results

Table 2 presents the primary comparison across all methods at the 100-agent scale.

LatentGate with Qwen-2.5-1.5B achieves 80.0% OOD accuracy at 28ms, outperforming GPT-4.1 (77.0% at 2075ms) while being 74× faster in end-to-end routing latency and improving over embedding baselines by 13–22 absolute points. The 98.8% in-domain accuracy confirms high separability in whitened space; the 18.8-point OOD drop reflects vocabulary shift. Among backbones, Qwen-

Method	ID Acc.	OOD Acc.	Lat. (ms)	Size
<i>Encoder / Embedding-Based Routers</i>				
BGE-M3	82.4	60.0	25	–
Semantic Router	78.6	58.4	50	–
SetFit	89.2	67.2	35	110M
<i>LLM Prompt-Based Routers</i>				
GPT-4.1	–	77.0	2075	–
GPT-4.1-nano	–	70.4	1851	–
Gemini 2.5 Flash	–	74.0	1500	–
<i>LatentGate (Ours)</i>				
Qwen-0.5B	60.0	45.0	18	3–7KB
Qwen-1.5B	98.8	80.0	28	3–7KB
Qwen-3B	97.4	78.8	42	3–7KB
Phi-3.5-mini	96.8	77.6	45	3–7KB
Llama-3.2-1B	94.2	72.0	24	3–7KB

Table 2: Routing performance at 100 agents. ID = in-domain accuracy (%), OOD = out-of-distribution accuracy (%), Lat. = latency in milliseconds. Probe size reflects trainable LR weights only (1,700 params); PCA/scaler are fixed preprocessing. LLM routers lack in-domain scores as they use zero-shot prompting.

Method	10	25	50	75	100
GPT-4.1	96.0	90.4	84.0	79.2	77.0
GPT-4.1-nano	92.0	84.8	78.0	73.6	70.4
Gemini 2.5 Flash	94.0	88.0	82.0	77.2	74.0
LatentGate	97.0	93.2	87.6	83.2	80.0

Table 3: OOD accuracy (%) as agent count scales. LLM routers degrade as agent count increases, consistent with prompt growth and increased label confusion at scale. LatentGate avoids context growth with C ; inference is dominated by the frozen SLM forward pass, with a small $O(Ck)$ probe.

1.5B offers the best trade-off; Qwen-0.5B lacks capacity while larger models show diminishing returns, suggesting 1–2B parameters as the sweet spot.

6.2 Scalability

Table 3 shows prompt-based routers degrade sharply with agent count, consistent with longer prompts and increased label confusion at scale: GPT-4.1 drops 19 points (96.0% to 77.0%). LatentGate degrades more gracefully (97.0% to 80.0%) since features are fixed at $k = 16$ dimensions and inference avoids context growth with C .

6.3 Ablation Study

Table 4 isolates each design choice. PCA-whitening is the most critical component: removing it drops OOD accuracy by 17.2 points (80.0% to 62.8%), confirming anisotropy correction is essen-

Configuration	OOD Acc.	Δ
Full LatentGate (Qwen-1.5B)	80.0	–
<i>Pooling Strategy</i>		
[CLS] token	72.4	–7.6
Last token	74.8	–5.2
Max pooling	76.0	–4.0
<i>Whitening</i>		
No whitening (raw + probe)	62.8	–17.2
StandardScaler only	71.6	–8.4
PCA only (no whitening)	74.0	–6.0
<i>Dimensionality k</i>		
$k = 4$	73.6	–6.4
$k = 8$	77.2	–2.8
$k = 32$	79.6	–0.4
$k = 64$	78.8	–1.2

Table 4: Ablation study on 100-agent OOD routing. Whitening provides the largest single improvement (+17.2 points). Mean pooling outperforms alternatives by 4–8 points. $k = 16$ balances accuracy and efficiency; $k = 32$ offers marginal gains while doubling probe size.

Method	Cost/1K Routing Calls	Monthly Cost	Savings
GPT-4.1	\$16.14	\$16,140	–
Gemini 2.5 Flash	\$2.43	\$2,434	84.9%
GPT-4.1-nano	\$0.81	\$807	95.0%
LatentGate	\$0.0027	\$2.72	99.98%

Table 5: Routing cost comparison (USD) at 1M queries/month. LLM costs are computed from publicly listed API pricing (input + output tokens). For 100 agents, each LLM routing call consumes ~ 8 K input tokens (agent descriptions) + ~ 30 tokens (query) and produces ~ 10 output tokens. LatentGate cost reflects T4 GPU compute at \$0.35/hr on-demand, processing ~ 128 K queries/hr at 28ms each.

tial. StandardScaler alone recovers only part of the gap (71.6%), while PCA without whitening yields 74.0%. Mean pooling outperforms alternatives by 4–8 points, expected since causal SLMs lack a dedicated [CLS] representation. The dimensionality sweep shows $k = 16$ as the sweet spot; $k \geq 32$ offers marginal gains at larger probe sizes.

6.4 Cost Analysis

Table 5 compares per-query costs using publicly listed API pricing as of January 2026¹. Each LLM routing call requires ~ 8 K input tokens (100 agent descriptions) and generates ~ 10 output tokens. LatentGate eliminates API costs entirely: a single T4 GPU (\$0.35/hr) processes ~ 128 K queries/hour,

¹OpenAI: platform.openai.com/docs/pricing; Google: ai.google.dev/pricing; T4 on-demand: AWS g4dn instances.

yielding \$2.72/month for 1M queries—a 99.98% reduction over GPT-4.1. These estimates reflect isolated routing compute and exclude orchestration, network, autoscaling, and engineering overhead; prices vary by region and provider updates.

6.5 Error Analysis

We examine the 20% of OOD queries that LatentGate misroutes. The dominant failure mode ($\sim 65\%$) is *cross-domain vocabulary overlap*: queries like “reset the deployment pipeline” are misrouted from DevOps to IT Helpdesk because “reset” is a strong lexical cue for IT. A secondary mode ($\sim 25\%$) is ambiguous queries genuinely underspecified without conversational context; the remaining 10% are novel phrasing edge cases. A two-layer MLP probe yielded $<1\%$ improvement, so we retain the linear probe for interpretability and online update efficiency. Errors can be mitigated by confidence-based fallback (escalating low-confidence queries to a human or LLM router); top-3 accuracy is 93.6% (Appendix F), making the system viable as a first-pass filter.

7 Conclusion

We presented **LatentGate**, a non-generative routing primitive that identifies and resolves representation anisotropy to unlock high-fidelity semantic routing in massive multi-agent systems. By decoupling semantic extraction from token generation, our three-stage pipeline achieves 80% OOD accuracy across 100 agents at ~ 28 ms: outperforming GPT-4.1 while operating $74\times$ faster at 99.98% lower cost.

Beyond raw performance, the lightweight probing architecture of LatentGate enables a self-healing production loop, where sub-10ms warm-start updates allow for real-time error correction without the prohibitive overhead of backbone fine-tuning. Our findings suggest that for high-class routing, representation geometry is a more critical bottleneck than model scale. By demonstrating that PCA-whitening can recover collapsed discriminative structures in frozen SLMs, this work provides a scalable blueprint for the next generation of lean, adaptive orchestration layers in enterprise AI. While currently limited to single-turn interactions and vocabulary-heavy OOD shifts, future work will extend this geometric approach to conversational state management and hierarchical agent discovery.

Ethical Considerations

Routing Bias and Fairness. LatentGate’s routing decisions are learned from training data distributions. If certain user groups phrase queries differently (e.g., non-native speakers, users with accessibility needs), the router may exhibit disparate error rates across populations. We recommend monitoring per-demographic routing accuracy in deployment and implementing confidence-based fallback to human operators for low-confidence predictions.

Training Data Provenance. Our training queries are generated via few-shot prompting of GPT-4.1 seeded with real production exemplars (used only as in-context demonstrations, never in training or evaluation). All outputs were independently verified by two annotators with adjudication. Despite this grounding, LLM-generated queries may encode biases present in the generation model’s training corpus. Practitioners should audit training data for representational balance and supplement with real user queries where available.

Reproducibility. The benchmark skeleton (agent taxonomy, prompt templates, and evaluation harness) is available at [anonymous repository] for review. While full production queries cannot be released due to privacy constraints, we release the complete synthetic benchmark, verification protocol, and generation prompts. Upon publication, we will additionally release pre-trained probe weights and all evaluation code.

Model Licensing. The SLM backbone (Qwen-2.5-1.5B-Instruct) is released under the Apache 2.0 license and is suitable for commercial use. All baseline models (GPT-4.1, Gemini 2.5 Flash) were accessed through their respective commercial APIs under standard terms of service.

References

Pranjal Aggarwal, Aman Madaan, Ankit Anand, Srividya Pranavi Potharaju, Swaroop Mishra, Pei Zhou, Aditya Gupta, Dheeraj Rajagopal, Karthik Kapananthu, Yiming Yang, Shyam Upadhyay, Manaal Faruqui, and Mausam. 2025. *Automix: Automatically mixing language models*. *Preprint*, arXiv:2310.12963.

Aurelio AI. 2024. *Semantic router*. <https://github.com/aurelio-labs/semantic-router>.

Yonatan Belinkov. 2022. *Probing classifiers: Promises, shortcomings, and advances*. *Computational Linguistics*, 48(1):207–219.

Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2025. *M3-embedding: Multilinguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation*. *Preprint*, arXiv:2402.03216.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023. *Frugalgpt: How to use large language models while reducing cost and improving performance*. *Preprint*, arXiv:2305.05176.

Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Ruhle, Laks V. S. Lakshmanan, and Ahmed Hassan Awadallah. 2024. *Hybrid llm: Cost-efficient and quality-aware query routing*. *Preprint*, arXiv:2404.14618.

Kawin Ethayarajh. 2019. *How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings*. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65, Hong Kong, China. Association for Computational Linguistics.

Ali Forooghi, Shaghayegh Sadeghi, and Jianguo Lu. 2024. *Whitening not recommended for classification tasks in LLMs*. In *Proceedings of the 9th Workshop on Representation Learning for NLP (RepL4NLP-2024)*, pages 285–289, Bangkok, Thailand. Association for Computational Linguistics.

Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2019. *Representation degeneration problem in training natural language generation models*. In *International Conference on Learning Representations*.

Nathan Godey, Éric de la Clergerie, and Benoît Sagot. 2024. *Anisotropy is inherent to self-attention in transformers*. *Preprint*, arXiv:2401.12143.

Qitian Jason Hu, Jacob Bieker, Xiuyu Li, Nan Jiang, Benjamin Keigwin, Gaurav Ranganath, Kurt Keutzer, and Shriyash Kaustubh Upadhyay. 2024. *Router-bench: A benchmark for multi-llm routing system*. *Preprint*, arXiv:2403.12031.

Junjie Huang, Duyu Tang, Wanjun Zhong, Shuai Lu, Linjun Shou, Ming Gong, Daxin Jiang, and Nan Duan. 2021. *Whiteningbert: An easy unsupervised sentence embedding approach*. *Preprint*, arXiv:2104.01767.

Elias Lumer, Faheem Nizar, Anmol Gulati, Pradeep Honaganahalli Basavaraju, and Vamse Kumar Subbiah. 2025. *Tool-to-agent retrieval: Bridging tools and agents for scalable llm multi-agent systems*. *Preprint*, arXiv:2511.01854.

499	Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M Waleed Kadous, and Ion Stoica. 2025. Routellm: Learning to route llms with preference data . <i>Preprint</i> , arXiv:2406.18665.	a user would type to reach this agent. Queries should vary in: (1) formality (casual to professional), (2) specificity (vague to detailed), (3) phrasing (questions, commands, descriptions). Include common misspellings and abbreviations. Do NOT include the agent name in queries.	552
500		Agent: {agent_name}	562
501		Description:	563
502		{agent_description}	564
503		Tools: {agent_tools}	565
504	Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. 2021. Whitening sentence representations for better semantics and faster retrieval . <i>Preprint</i> , arXiv:2103.15316.		
505			
506			
507			
508	Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke Bates, Daniel Korat, Moshe Wasserblat, and Oren Pereg. 2022. Efficient few-shot learning without prompts . <i>Preprint</i> , arXiv:2209.11055.		
509			
510			
511			
512	Zhiyuan Yao, Zishan Xu, Yifu Guo, Zhiguang Han, Cheng Yang, Shuo Zhang, Weinan Zhang, Xingshan Zeng, and Weiwen Liu. 2026. Toolace-mcp: Generalizing history-aware routing from mcp tools to the agent web . <i>Preprint</i> , arXiv:2601.08276.		
513			
514			
515			
516			
517	Yanwei Yue, Guibin Zhang, Boyang Liu, Guancheng Wan, Kun Wang, Dawei Cheng, and Yiyan Qi. 2025. MasRouter: Learning to route LLMs for multi-agent systems . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15549–15572, Vienna, Austria. Association for Computational Linguistics.		
518			
519			
520			
521			
522			
523			
524	Wenjie Zhuo, Yifan Sun, Xiaohan Wang, Linchao Zhu, and Yi Yang. 2023. WhitenedCSE: Whitening-based contrastive learning of sentence embeddings . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12135–12148, Toronto, Canada. Association for Computational Linguistics.		
525			
526			
527			
528			
529			
530			
531	A Agent Taxonomy		
532	Table 6 lists all 100 agents organized by operational domain. Each agent is defined by a system prompt describing its capabilities, tools, and scope. The taxonomy reflects a realistic enterprise deployment with substantial inter-domain vocabulary overlap (e.g., IT Helpdesk and DevOps both handle “access” and “reset” queries; Finance and HR both process “approval” requests).		
533			
534			
535			
536			
537			
538			
539			
540	B Prompt Templates		
541	B.1 Synthetic Query Generation Prompt		
542	The following prompt was used with GPT-4.1 to generate 50 training queries per agent:		
543			
544	You are generating realistic user queries for an enterprise routing system. Below are 5 real example queries from production logs for reference. Given the following agent description, generate 50 diverse natural language queries that		
545			
546			
547			
548			
549			
550			
551			
		The evaluation set uses a separate prompt emphasizing paraphrasing and cross-domain vocabulary:	567
		Generate 20 out-of-distribution test queries for this agent. These should be: (1) paraphrased using synonyms and different sentence structures, (2) use colloquial/informal language, (3) include vocabulary that overlaps with other domains (e.g., use “fix” which could apply to IT, DevOps, or Engineering). The queries should be challenging to route correctly.	568
			569
			570
			571
			572
			573
			574
			575
			576
			577
			578
			579
			580
			581
		B.2 OOD Query Generation Prompt	566
		B.3 LLM Routing Prompt	582
		The zero-shot routing prompt used for GPT-4.1, GPT-4.1-nano, and Gemini 2.5 Flash baselines:	583
			584
		You are an enterprise routing system. Given a user query, select the most appropriate agent from the list below. Respond with ONLY the agent name, nothing else.	585
			586
			587
			588
			589
			590
		Agents:	591
		1. IT Helpdesk: Handles password resets, account lockouts, software installation...	592
			593
			594
			595
		2. GitHub DevOps: Manages repository access, branch protection, CI/CD...	596
			597
			598
		...	599
		100. Analytics Requests: Handles ad-hoc data queries, report generation...	600
			601
			602
		User query: {query}	603
		Agent:	604
		C Implementation Details	605
		C.1 Hyperparameters	606
		C.2 Latency Breakdown	607
		Table 8 decomposes the end-to-end latency of LatentGate (Qwen-2.5-1.5B) on a single NVIDIA T4 GPU with batch size 1.	608
			609
			610

Domain (Count)	Agents
IT Operations (12)	IT Helpdesk, Network Admin, VPN Support, Active Directory, Email Admin, Print Services, Hardware Support, Software Licensing, Security Ops, Endpoint Management, SSO Admin, IT Procurement
DevOps (10)	GitHub DevOps, CI/CD Pipeline, Docker/K8s, AWS Infra, GCP Infra, Azure Infra, Terraform, Monitoring/Alerting, Incident Response, Release Management
HR (10)	HR General, Payroll, Benefits Admin, Recruiting, Onboarding, Performance Review, Leave Management, Employee Relations, Training/L&D, Compensation
Finance (10)	Accounts Payable, Accounts Receivable, Expense Reports, Budgeting, Tax Compliance, Financial Reporting, Procurement, Vendor Management, Audit, Treasury
Legal (8)	Contract Review, IP/Patents, Compliance, Data Privacy (GDPR/CCPA), Employment Law, Litigation Support, NDA Management, Regulatory Filing
Engineering (10)	Frontend Dev, Backend Dev, Mobile Dev, QA/Testing, Database Admin, API Platform, ML/AI Ops, Tech Debt, Architecture Review, Code Review
Sales & Marketing (10)	CRM Admin, Lead Generation, Sales Ops, Campaign Management, Content Marketing, SEO/SEM, Social Media, Partner Relations, Pricing Strategy, Market Research
Customer Success (8)	Customer Support Tier 1, Customer Support Tier 2, Account Management, Customer Onboarding, Churn Prevention, Feedback/NPS, Knowledge Base, Escalation
Facilities & Admin (8)	Office Management, Travel Booking, Event Planning, Mail/Shipping, Badge Access, Parking, Catering, Workspace Reservation
Communication (8)	Slack Admin, Email Campaigns, Internal Comms, Wiki/Confluence, Zoom/Meeting, Project Mgmt (Jira), Documentation, Newsletter
Data & Analytics (6)	Data Engineering, BI/Dashboards, Data Governance, ETL Pipeline, A/B Testing, Analytics Requests

Table 6: Full taxonomy of 100 enterprise agents across 11 operational domains. Agents within the same domain share substantial vocabulary overlap, creating the routing challenge that motivates LatentGate.

Parameter	Value
<i>SLM Backbone</i>	
Max sequence length	128 tokens
Padding	Right, pad token = eos token
Precision	float16
Batch size (inference)	1
<i>PCA-Whitening</i>	
Components (k)	16
Whiten flag	True
Scaler	StandardScaler (before PCA)
<i>Linear Probe</i>	
Model	LogisticRegression (sklearn)
Regularization (C)	1.0
Solver	lbfgs
Max iterations	1000
Multi-class	multinomial
<i>Layer Selection</i>	
Sweep strategy	10 layers uniformly sampled
Selection metric	5-fold CV accuracy
<i>Online Learning</i>	
Warm start	True
Re-whitening	No (reuse existing PCA)

Table 7: Full hyperparameter configuration for LatentGate.

Stage	Time (ms)	% Total
Tokenization	0.8	2.9%
SLM forward pass (L7)	24.5	87.5%
Mean pooling	0.1	0.4%
StandardScaler + PCA	0.3	1.1%
Linear probe inference	0.1	0.4%
Overhead (data transfer)	2.2	7.9%
Total	28.0	100%

Table 8: Latency breakdown for LatentGate on T4 GPU. The SLM forward pass dominates (87.5%); the whitening and probe stages add negligible overhead (<0.5ms combined).

E Routing Examples

616

Table 10 shows representative correct and incorrect routing decisions by LatentGate (Qwen-2.5-1.5B).

617

618

F Top- k Routing Accuracy

619

D Per-Backbone Layer Sweep

Table 9 reports the optimal layer and corresponding probe accuracy for each SLM backbone, selected via 5-fold cross-validated sweep across 10 uniformly sampled layers.

Backbone	Layers	Best L	CV Acc.
Qwen-2.5-0.5B	24	8	60.0%
Qwen-2.5-1.5B	28	7	98.8%
Qwen-2.5-3B	36	9	97.4%
Phi-3.5-mini	32	10	96.8%
Llama-3.2-1B	16	8	94.2%

Table 9: Optimal layer per backbone via 5-fold CV sweep. Across all architectures, the best layer falls in the early-to-middle range (L7–L10), suggesting that routing-relevant semantic features are encoded well before the final layers. This contrasts with typical NLU probing findings where deeper layers dominate, likely because routing requires broad topic discrimination rather than fine-grained linguistic analysis.

Query	Gold	Predicted	Correct?
<i>Correct Routings</i>			
“My VPN keeps disconnecting every 5 minutes”	VPN Support	VPN Support	✓
“Need to update the CI pipeline to run tests in parallel”	CI/CD Pipeline	CI/CD Pipeline	✓
“Can someone review the NDA for the Acme partnership?”	NDA Mgmt	NDA Mgmt	✓
“The dashboard is showing stale data from yesterday”	BI/Dashboards	BI/Dashboards	✓
“How do I submit my expense report for the NYC trip?”	Expense Reports	Expense Reports	✓
<i>Cross-Domain Vocabulary Overlap (65% of errors)</i>			
“Reset the deployment pipeline”	CI/CD Pipeline	IT Helpdesk	×
“Approve the budget transfer to marketing”	Budgeting	Leave Mgmt	×
“Fix the access permissions on the staging server”	AWS Infra	Active Directory	×
<i>Ambiguous Queries (25% of errors)</i>			
“Something’s broken, can someone look at it?”	(any)	IT Helpdesk	×
“I need help with the report”	Financial Rpt	BI/Dashboards	×
<i>Novel Phrasing (10% of errors)</i>			
“Yo can u hook me up with more cloud credits lol”	AWS Infra	Slack Admin	×

Table 10: Representative routing examples across all three error categories identified in Section 6.5. Correct routings demonstrate LatentGate’s ability to handle domain-specific terminology. Error examples illustrate the three failure modes: vocabulary overlap, query ambiguity, and novel phrasing.

Method	Top-1 (OOD)	Top-3 (OOD)
LatentGate (Qwen-1.5B)	80.0	93.6

Table 11: Top- k OOD accuracy. Top-3 accuracy measures whether the gold agent appears among the three highest-scoring predictions.