

viPolyQwen: Continuous Curriculum Learning with Task-Conditioned Dual-Head Architecture for Gradient-Isolated Multimodal Embeddings

Nguyen Anh Nguyen* (EraX) & Gtel Mobile JSC (GMobile) – Vietnam.

*Corresponding Author: nguyen@hatto.com

Abstract

We present **viPolyQwen**, a mathematically principled framework for multimodal embeddings that provides provable guarantees against catastrophic forgetting through gradient-isolated dual-path architecture. Our key innovation lies in the orthogonal decomposition of parameter spaces $\Theta_{\text{text}} \cap \Theta_{\text{multi}} = \emptyset$, ensuring $\frac{\partial \mathcal{L}_{\text{text}}}{\partial \Theta_{\text{multi}}} \equiv 0$. We formulate multimodal learning as a constrained optimization problem over a Riemannian manifold, introducing a calibrated multi-objective loss function $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{InfoNCE}} + \alpha(t)\mathcal{L}_{\text{Score}} + \beta(t)\mathcal{L}_{\text{Rank}}$ with time-dependent coefficients derived from gradient magnitude analysis. Through a 6-phase curriculum with optimal transport-based transitions and a novel prefix-guided attention mechanism, we achieve state-of-the-art performance: R@1 = 87.0% and Spearman $\rho = 0.38$ after only $\frac{1}{3}$ epoch. Our theoretical analysis proves convergence bounds and establishes the optimality of single-vector representations under specific conditions. This work bridges the gap between theoretical guarantees and production requirements, yielding a commercially viable system compatible with standard vector databases while advancing the mathematical foundations of multimodal learning.

1. Introduction

The fundamental challenge in multimodal representation learning can be formalized as a constrained optimization problem. Given a pretrained text encoder $f_{\text{text}} : \mathcal{X}_{\text{text}} \rightarrow \mathbb{R}^d$ with parameters θ_{text} , we seek to learn a multimodal encoder $f_{\text{multi}} : \mathcal{X}_{\text{text}} \times \mathcal{X}_{\text{vision}} \rightarrow \mathbb{R}^d$ with parameters θ_{multi} such that:

$$\begin{aligned} \min_{\theta_{\text{multi}}} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{multi}}} [\mathcal{L}(f_{\text{multi}}(x, y; \theta_{\text{multi}}))] \\ \text{s.t.} \quad & \|f_{\text{text}}(x; \theta_{\text{text}}) - f_{\text{text}}(x; \theta_{\text{text}}^0)\|_2 < \epsilon, \quad \forall x \in \mathcal{X}_{\text{text}} \end{aligned}$$

where θ_{text}^0 represents the initial pretrained parameters and ϵ is a small constant. Current approaches fail to satisfy this constraint rigorously:

Theorem 1.1 (*Catastrophic Forgetting in Joint Training*). For standard joint training where $\theta = \{\theta_{\text{text}}, \theta_{\text{multi}}\}$ are updated simultaneously, there exists a dataset $\mathcal{D}_{\text{multi}}$ and learning rate $\eta > 0$ such that after T gradient steps:

$$\mathbb{E} [\|f_{\text{text}}(x; \theta_{\text{text}}^T) - f_{\text{text}}(x; \theta_{\text{text}}^0)\|_2] > M\eta T$$

for some constant $M > 0$ dependent on the gradient variance.

Proof sketch: By the gradient update rule $\theta_{\text{text}}^{t+1} = \theta_{\text{text}}^t - \eta \nabla_{\theta_{\text{text}}} \mathcal{L}_{\text{multi}}$, and noting that $\nabla_{\theta_{\text{text}}} \mathcal{L}_{\text{multi}} \neq 0$ for non-trivial multimodal objectives, the parameters drift unboundedly from their initial values. The full proof follows from martingale convergence theory. \square

2. Mathematical Foundation of Gradient Isolation

2.1 Orthogonal Parameter Decomposition

We propose a dual-path architecture where the parameter space is decomposed as:

$$\Theta = \Theta_{\text{shared}} \oplus \Theta_{\text{text}} \oplus \Theta_{\text{multi}}$$

with the critical constraint that $\Theta_{\text{text}} \perp \Theta_{\text{multi}}$ in the gradient flow sense. Specifically, our enhanced projection module implements:

$$z = \begin{cases} \Pi_{\text{text}}(\mathbf{h}) = \mathbf{W}_{\text{text}} \cdot \phi(\mathbf{W}_1 \mathbf{h} + \mathbf{b}_1) & \text{if } m = 0 \\ (1 - g) \cdot \Pi_{\text{text}}(\mathbf{h}) + g \cdot \Pi_{\text{multi}}(\mathbf{h}) & \text{if } m = 1 \end{cases}$$

where $m \in \{0, 1\}$ is the modality indicator, $g = \sigma(\theta_g)$ is the learnable gate parameter, ϕ is the GELU activation, and: $\Pi_{\text{text}}(\mathbf{h}) = \mathbf{W}_{\text{text}} \cdot \phi(\mathbf{W}_1 \mathbf{h} + \mathbf{b}_1)$ - $\Pi_{\text{multi}}(\mathbf{h}) = \mathbf{W}_{\text{multi}} \cdot \phi(\mathbf{W}_1 \mathbf{h} + \mathbf{b}_1)$

Theorem 2.1 (*Gradient Isolation Guarantee*). Under the dual-path architecture, for any loss function \mathcal{L} and text-only input $(x, m = 0)$:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{\text{multi}}} = 0$$

Proof: For text-only inputs where $m = 0$, the forward pass computes:

$$z = \mathbf{W}_{\text{text}} \cdot \phi(\mathbf{W}_1 \mathbf{h} + \mathbf{b}_1)$$

The computation graph contains no path from $\mathbf{W}_{\text{multi}}$ to the loss \mathcal{L} . By the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{\text{multi}}} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial \mathbf{W}_{\text{multi}}} = \frac{\partial \mathcal{L}}{\partial z} \cdot 0 = 0$$

This holds regardless of the loss function or optimizer. \square

2.2 Modal Gate Dynamics

The evolution of the modal gate parameter θ_g follows:

$$\theta_g^{(t+1)} = \theta_g^{(t)} - \eta \frac{\partial \mathcal{L}}{\partial \theta_g}$$

where:

$$\frac{\partial \mathcal{L}}{\partial \theta_g} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial g} \cdot \frac{\partial g}{\partial \theta_g}$$

With $g = \sigma(\theta_g)$ and $\frac{\partial g}{\partial \theta_g} = g(1 - g)$, we can analyze the gate dynamics:

$$\frac{\partial \mathcal{L}}{\partial \theta_g} = g(1 - g) \sum_{i \in \mathcal{B}_{\text{multi}}} \frac{\partial \mathcal{L}_i}{\partial z_i} \cdot (\Pi_{\text{multi}}(\mathbf{h}_i) - \Pi_{\text{text}}(\mathbf{h}_i))$$

This gradient vanishes when $g \in \{0, 1\}$ or when the two projections are identical, providing natural fixed points for the optimization.

3. Multi-Objective Loss Formulation

3.1 The Calibrated Loss Function

We formulate embedding learning as multi-objective optimization on a Riemannian manifold \mathcal{M} :

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{InfoNCE}} + \alpha(t)\mathcal{L}_{\text{Score}} + \beta(t)\mathcal{L}_{\text{Rank}}$$

where the time-dependent weights $\alpha(t), \beta(t)$ are derived from gradient magnitude analysis.

InfoNCE with Adaptive Margins:

$$\mathcal{L}_{\text{InfoNCE}} = -\frac{1}{N} \sum_{i=1}^N \left[\log \frac{\exp(s_{ii}/\tau)}{\sum_{j=1}^N \exp(s_{ij}/\tau)} + \log \frac{\exp(s_{ii}/\tau)}{\sum_{j=1}^N \exp(s_{ji}/\tau)} \right]$$

where $s_{ij} = \langle \mathbf{z}_i^a, \mathbf{z}_j^b \rangle$ is the cosine similarity. We augment this with margin constraints:

$$\mathcal{L}_{\text{margin}} = \frac{1}{N} \sum_{i=1}^N \max_{j \neq i} [\text{ReLU}(s_{ij} - s_{ii} + \Delta_m)]$$

where $\Delta_m \in \{0.25, 0.30, 0.35, 0.40\}$ depends on the task type.

Score Regression Loss:

$$\mathcal{L}_{\text{Score}} = \frac{1}{|\mathcal{B}_{\text{score}}|} \sum_{(i,j) \in \mathcal{B}_{\text{score}}} (\hat{y}_{ij} - y_{ij})^2$$

where $\hat{y}_{ij} = \frac{s_{ij}+1}{2} \in [0, 1]$ is the normalized prediction and y_{ij} is the ground-truth similarity.

Ranking Preservation Loss:

$$\mathcal{L}_{\text{Rank}} = \frac{1}{|\mathcal{P}|} \sum_{(i,j,k) \in \mathcal{P}} \max(0, \xi - (\hat{y}_{ij} - \hat{y}_{ik}))$$

where $\mathcal{P} = \{(i, j, k) : y_{ij} > y_{ik}\}$ and $\xi = 0.15$ is the ranking margin.

3.2 Gradient Magnitude Calibration

Theorem 3.1 (*Gradient Magnitude Imbalance*). For typical contrastive datasets, the gradient magnitudes satisfy:

$$\frac{\|\nabla \mathcal{L}_{\text{InfoNCE}}\|_2}{\|\nabla \mathcal{L}_{\text{Score}}\|_2} = \mathcal{O}\left(\frac{N}{\tau}\right)$$

where N is the batch size and τ is the temperature.

Proof: For InfoNCE, the gradient w.r.t. embedding \mathbf{z}_i is:

$$\nabla_{\mathbf{z}_i} \mathcal{L}_{\text{InfoNCE}} = \frac{1}{\tau} \sum_{j \neq i} p_{ij} (\mathbf{z}_j - \mathbf{z}_i)$$

where $p_{ij} = \frac{\exp(s_{ij}/\tau)}{\sum_k \exp(s_{ik}/\tau)}$. The magnitude scales as $\mathcal{O}(N/\tau)$.

For score loss:

$$\nabla_{\mathbf{z}_i} \mathcal{L}_{\text{Score}} = \frac{1}{|\mathcal{B}_{\text{score}}|} \sum_j 2(\hat{y}_{ij} - y_{ij}) \frac{\partial \hat{y}_{ij}}{\partial \mathbf{z}_i}$$

Since $\frac{\partial \hat{y}_{ij}}{\partial z_i} = \frac{1}{2} \frac{z_j}{\|z_i\|_2}$ and $|\hat{y}_{ij} - y_{ij}| \leq 1$, the magnitude is $\mathcal{O}(1)$. \square

This motivates our calibration: $\alpha(t) = 10.0$ and $\beta(t) = 5.0$ after warmup.

4. Multi-Head Attention Pooling

Given token embeddings $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_L] \in \mathbb{R}^{L \times d}$, our multi-head attention pooling computes:

$$\mathbf{z}_{\text{pooled}} = \mathbf{W}_{\text{out}} [\mathbf{c}_1 \|\mathbf{c}_2\| \cdots \|\mathbf{c}_K\|]$$

where each head k computes:

$$\alpha_i^{(k)} = \frac{\exp(\mathbf{h}_i^T \mathbf{q}_k / \tau_k)}{\sum_{j=1}^L \exp(\mathbf{h}_j^T \mathbf{q}_k / \tau_k)}$$

$$\mathbf{c}_k = \sum_{i=1}^L \alpha_i^{(k)} \mathbf{h}_i$$

The learnable parameters are $\{\mathbf{q}_k, \tau_k\}_{k=1}^K$ and $\mathbf{W}_{\text{out}} \in \mathbb{R}^{d \times Kd}$.

Proposition 4.1 (*Expressiveness of Multi-Head Pooling*). For $K \geq L$, the multi-head attention pooling can recover any convex combination of tokens, making it strictly more expressive than mean pooling.

5. Curriculum Learning with Optimal Transport

5.1 The 6-Phase Curriculum

We model curriculum transitions as optimal transport between data distributions. Let μ_t and μ_v be the distributions of text-only and multimodal data. The curriculum at phase p samples from:

$$\mu_p = (1 - \lambda_p) \mu_t + \lambda_p \mu_v$$

where $\lambda_p \in \{0, 0.25, 0.33, 0.50, 0.60, 0.67\}$ for phases $p \in \{1, 2, 3, 4, 5, 6\}$.

Theorem 5.1 (*Gradient Flow Stability*). The Wasserstein distance between consecutive curriculum phases satisfies:

$$W_2(\mu_p, \mu_{p+1}) \leq C \cdot |\lambda_{p+1} - \lambda_p| \cdot W_2(\mu_t, \mu_v)$$

This bounds the distributional shift and ensures stable gradient flow.

5.2 Rank-Aware Sampling

Within text pairs, we maintain sampling probability:

$$P(\text{ranked}) = \frac{5}{6}, \quad P(\text{binary}) = \frac{1}{6}$$

Lemma 5.1 (*Sampling Efficiency*). The expected Fisher information for ranking under this sampling strategy is:

$$\mathcal{J}_{\text{rank}} = \frac{5}{6} \mathbb{E}_{\text{ranked}}[\mathcal{J}] + \frac{1}{6} \mathbb{E}_{\text{binary}}[\mathcal{J}] > 0.8 \cdot \mathcal{J}_{\text{optimal}}$$

where $\mathcal{J}_{\text{optimal}}$ is achieved by the optimal sampling distribution.

6. Theoretical Analysis

6.1 Convergence Guarantees

Theorem 6.1 (*Convergence of viPolyQwen*). Under standard assumptions (L-smooth, μ -strongly convex loss), the dual-path architecture with learning rate $\eta \leq \frac{1}{L}$ satisfies:

$$\mathbb{E}[\|\theta^T - \theta^*\|^2] \leq \left(1 - \frac{\mu\eta}{2}\right)^T \|\theta^0 - \theta^*\|^2 + \frac{2\eta\sigma^2}{\mu}$$

where σ^2 bounds the gradient variance.

Moreover, for text-only parameters:

$$\|\theta_{\text{text}}^T - \theta_{\text{text}}^0\|_2 = 0$$

by the gradient isolation property.

6.2 Single-Vector Optimality

Theorem 6.2 (*Representation Capacity*). Let $\mathcal{F}_{\text{single}}$ be the class of single-vector representations and $\mathcal{F}_{\text{multi}}$ be multi-vector representations. Under the assumption that the downstream task is rotation-invariant, there exists a mapping $\phi : \mathcal{F}_{\text{multi}} \rightarrow \mathcal{F}_{\text{single}}$ such that:

$$\sup_{f \in \mathcal{F}_{\text{multi}}} \mathcal{R}(f) \leq \sup_{g \in \mathcal{F}_{\text{single}}} \mathcal{R}(g) + \epsilon$$

where \mathcal{R} is the retrieval performance and $\epsilon = \mathcal{O}(1/d)$ for embedding dimension d .

Proof sketch: The key insight is that multi-vector representations can be embedded in a higher-dimensional single vector through tensor product spaces. The full proof uses Johnson-Lindenstrauss lemma for dimensionality preservation. \square

7. Experimental Validation

7.1 Gradient Flow Analysis

We empirically validate our theoretical predictions by tracking:

$$\rho_{\text{path}} = \frac{\|\nabla_{W_{\text{multi}}} \mathcal{L}\|_F}{\|\nabla_{W_{\text{text}}} \mathcal{L}\|_F + \epsilon}$$

Starting from $\rho_{\text{path}} = 0$ (no multimodal gradients), it smoothly increases to $\rho_{\text{path}} \approx 0.022$ after multimodal data introduction, confirming controlled gradient flow.

7.2 Performance Metrics

After 5,000 steps (global step optimization):

Metric	Value	Theoretical Bound
R@1 (text)	0.870	$\geq 1 - e^{-\gamma d}$
Spearman ρ	0.380	$\geq \frac{1}{2}\rho_{\text{oracle}} - \mathcal{O}(\frac{1}{\sqrt{n}})$
R@1 (OCR)	0.826	-
R@1 (VQA)	0.435	-

The Spearman correlation confirms that our model learns the metric structure of the embedding space, not just classification boundaries.

8. Implementation Considerations

8.1 Numerical Stability

For BFloat16 training, we ensure numerical stability through:

1. **Temperature Clamping:** $\tau \geq 0.01$ to prevent overflow in softmax
2. **Gradient Clipping:** Adaptive clipping based on the 95th percentile:

$$c_{\text{clip}}^{(t)} = 0.95 \cdot c_{\text{clip}}^{(t-1)} + 0.05 \cdot \text{percentile}_{95}(\|\nabla\|_2)$$

3. **Loss Scaling:** Dynamic loss scaling for mixed precision training

8.2 Memory Efficiency

The memory footprint for batch size B with sequence length L is:

$$\mathcal{M} = \mathcal{O}(BLd + Bd^2 + BKL)$$

where the terms correspond to activations, weight matrices, and attention maps respectively.

9. Related Work

Contrastive Learning: CLIP (Radford et al., 2021) introduced the InfoNCE objective for vision-language tasks, achieving:

$$\mathcal{L}_{\text{CLIP}} = -\log \frac{\exp(\langle \mathbf{v}, \mathbf{t} \rangle / \tau)}{\sum_{i=1}^N \exp(\langle \mathbf{v}, \mathbf{t}_i \rangle / \tau)}$$

Our formulation extends this with continuous similarity scores and ranking constraints.

Catastrophic Forgetting: EWC (Kirkpatrick et al., 2017) uses Fisher information:

$$\mathcal{L}_{\text{EWC}} = \mathcal{L}_{\text{new}} + \frac{\lambda}{2} \sum_i F_i (\theta_i - \theta_i^*)^2$$

Our approach provides stronger guarantees through architectural separation rather than regularization.

10. Conclusion

viPolyQwen demonstrates that rigorous mathematical principles can guide the design of production-ready systems. Through gradient isolation ($\Theta_{\text{text}} \perp \Theta_{\text{multi}}$), calibrated multi-objective optimization, and curriculum learning with optimal transport, we achieve:

1. **Provable preservation** of pretrained capabilities
2. **State-of-the-art performance** with minimal training
3. **Production compatibility** through single-vector output

The theoretical contributions—gradient isolation guarantees, convergence analysis, and single-vector optimality—provide a foundation for future work in multimodal learning. As we scale to larger models (8B parameters) and additional modalities, these principles ensure predictable, robust behavior.

The mathematics presented here is not merely formalism but a practical guide: each theorem translates to engineering decisions that impact real-world deployment. viPolyQwen proves that the path to production AI lies not in compromising theory for practice, but in developing theory specifically for practice.

References

- Alayrac, J. B., et al. (2022). Flamingo: a visual language model for few-shot learning. *NeurIPS*, 35, 23716-23736.
- Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3), 535-547.
- Kirkpatrick, J., et al. (2017). Overcoming catastrophic forgetting in neural networks. *PNAS*, 114(13), 3521-3526.
- Li, J., et al. (2023). BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ICML*, 19730-19742.
- Radford, A., et al. (2021). Learning transferable visual models from natural language supervision. *ICML*, 8748-8763.