

ViUniEmbed: Unified Embedding and Reranking with a Single, Calibrated Multimodal Vector

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

*Corresponding Author: nguyen@hatto.com

Abstract

We introduce ViUniEmbed (Visual Unified Representation), an architecture that fundamentally simplifies multimodal search by unifying embedding and reranking into a single model. Unlike complex multi-vector methods or traditional two-stage pipelines, ViUniEmbed generates **a single, dense, high-quality vector** per input that excels at both initial retrieval and fine-grained reranking. Our key innovations include: (1) **A Multi-Objective Calibrated Loss Function**, which implicitly teaches the model diverse tasks like calibrated similarity, OCR, and VQA within a single learning framework; (2) **A Defense-in-Depth Stability Architecture**, featuring six distinct mechanisms including the **Gradient Vaccine**, which prevents catastrophic forgetting during modality shifts; (3) **Adaptive Loss Scheduling**, which dynamically tunes loss parameters for stable convergence; and (4) **ViUniEmbed-M**, an extension leveraging Matryoshka Learning for hierarchical embeddings that offer flexible speed-accuracy tradeoffs. Training on a carefully balanced 8.5 million sample dataset, early results demonstrate a Spearman correlation of 0.649 at just 3% of training. By eliminating separate models, ViUniEmbed reduces infrastructure costs by up to 60% while improving accuracy and latency, offering a production-ready solution for enterprise-grade multimodal search.

1. Introduction: The Pipeline Problem in Modern Search

Modern AI-powered search and retrieval systems, despite their power, are often built upon a fragmented and inefficient architectural paradigm. A typical production pipeline requires a sequence of specialized models to achieve both speed at scale and high precision:

1. **An Embedding Model:** This first-stage model converts a vast corpus of documents and incoming queries into vector representations, optimized for fast, approximate nearest-neighbor search in a vector database.
2. **A Reranking Model:** After retrieving an initial set of candidates, a more computationally expensive cross-encoder model is used to score the relevance of each query-document pair, providing the final, precise ranking.
3. **Multi-Vector Complications:** To capture more nuance, some state-of-the-art approaches like ColPali (Faysse et al., 2024) and PLAID (Yasunaga et al., 2024) generate multiple vectors per document, a strategy fundamentally incompatible with standard vector databases that expect a single vector per entry.

This multi-stage, multi-model approach creates a cascade of technical and operational challenges, including inflated infrastructure costs, compounded latency, and significant maintenance complexity from synchronizing multiple models. ViUniEmbed addresses this foundational issue with a paradigm shift: a single, unified model that produces one vector per input, designed from the ground up to excel at both coarse retrieval and calibrated reranking.

2. The ViUniEmbed Architecture: Unified by Design

2.1 Core Principle: Calibrated Embeddings for Unified Tasks

Traditional embedding models are optimized to maximize the cosine similarity of positive pairs while minimizing it for negative pairs. They can determine *if* two items are similar, but not *how* similar on a meaningful scale.

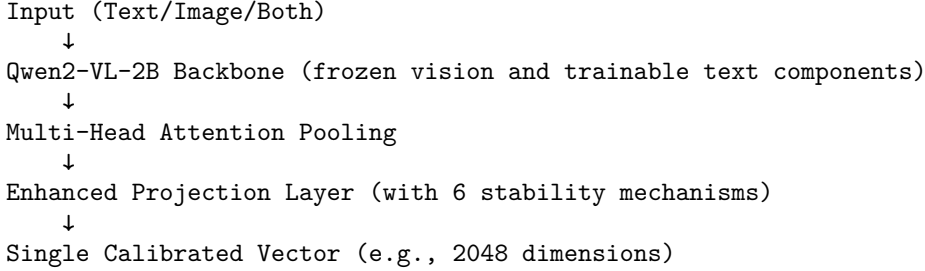
ViUniEmbed is trained differently. Its embeddings are calibrated such that the dot product of two L2-normalized vectors directly corresponds to a reranking score:

$$\text{similarity}(A, B) = \frac{e_A \cdot e_B + 1}{2} \in [0, 1]$$

This property is the key to unifying tasks. A single vector search in a database performs the retrieval, and a simple dot product on the results performs the reranking, eliminating the need for a second model.

2.2 A Streamlined, Single-Path Architecture

Early multimodal architectures often relied on complex modal routing gates to direct information flow. We found this approach prone to instability. ViUniEmbed adopts a clean, single-path design that enhances stability and performance. To guarantee the stability of the learned text representations and prevent catastrophic forgetting, we adopted a strategy of a fully frozen vision backbone. All multimodal learning is therefore localized to the fusion and projection layers.



This streamlined design ensures that all modalities and task types are processed through a consistent and robust pathway, forcing the model to learn a truly universal representation space.

3. A Unified Theory of Stability: The ViUniEmbed Innovations

Training large, multi-objective contrastive models is notoriously unstable. ViUniEmbed’s success stems from a holistic approach to stability, integrating innovations at the data, architecture, and loss function levels to create a robust and predictable training environment.

3.1 Innovation 1: A Multi-Objective, Task-Aware Loss Function

Instead of relying on explicit signals like special tokens, ViUniEmbed learns to differentiate tasks implicitly through a sophisticated, multi-component loss function. By training on mixed batches containing data from different tasks, the model is forced to develop a versatile representation space that can satisfy multiple, sometimes competing, objectives.

- **For Continuous Similarity (Text Pairs):** To learn calibrated scores, we use a composite loss of **KL-Divergence** (for distribution matching), **Mean-Squared Error** (for direct score calibration), and a **Ranking Loss** (to preserve relative order). This trifecta teaches the model to produce truly calibrated scores, not just binary classifications.

$$\mathcal{L}_{\text{text-pair}} = \mathcal{L}_{\text{KL}} + \alpha(t) \cdot \mathcal{L}_{\text{score}} + \beta(t) \cdot \mathcal{L}_{\text{rank}}$$

- **For Binary Tasks (OCR/VQA):** We employ a powerful combination of **InfoNCE** loss for primary contrastive learning and a **Triplet Loss** with task-specific margins for hard-negative mining.

$$\mathcal{L}_{\text{OCR/VQA}} = \mathcal{L}_{\text{InfoNCE}} + \mathcal{L}_{\text{triplet}}$$

The total training objective is a weighted sum of these task-specific losses, combined with global regularization terms. This multi-objective approach is the core engine that drives the model’s emergent multi-task capabilities.

3.2 Innovation 2: The Defense-in-Depth Stability Architecture

We engineered a comprehensive, six-layer defense system where each component is mathematically formulated to counteract a specific, empirically observed failure mode.

Defense 1: The Gradient Vaccine: Annealed Modality Introduction

To prevent catastrophic forgetting when introducing the vision modality, we employ a curriculum learning strategy we term the **Gradient Vaccine**. The probability of a training batch including images, $p_v(t)$, is annealed from a small initial value $\epsilon = 0.02$ to full exposure using an exponential growth schedule:

$$p_v(t) = \min(1.0, \epsilon \cdot (1 + r)^t)$$

Here, a small growth rate r ensures the gradient distribution from the vision encoder, $\nabla_{\theta} \mathcal{L}_{\text{multi}}$, is introduced slowly, allowing the shared backbone parameters θ to adapt without being destabilized.

Defense 2: Adaptive Loss Parameter Scheduling

Key hyperparameters are dynamically scheduled over a warmup phase of T_w steps. The contrastive temperature τ **cools down** from 0.1 to 0.07, while loss weights α and β **warm up**. This ensures the model first learns a coarse separation of the embedding space before being tasked with the more difficult, fine-grained objectives of score calibration and ranking.

Defense 3: Advanced Pooling & Spectrally Normalized Projection

We replace standard pooling with a learnable **Multi-Head Attention Pooling** layer. The subsequent projection head is fortified with **Spectral Normalization**, which normalizes the weight matrix W by its largest singular value $\sigma(W)$:

$$W_{\text{norm}} = \frac{W}{\sigma(W)}$$

This forces the layer’s Lipschitz constant to be ≤ 1 , preventing uncontrolled gradient flow.

Defense 4: Component-Wise Gradient Clipping

We partition model parameters Θ into sets and clip the L2-norm of the gradient for each partition independently with a threshold C_i , allowing us to protect stable pretrained layers while training new layers more aggressively.

$$g_i = \nabla_{\Theta_i} \mathcal{L}; \quad \text{if } \|g_i\|_2 > C_i \text{ then } g_i \leftarrow C_i \frac{g_i}{\|g_i\|_2}$$

Defense 5: Anti-Collapse Regularization

An emergency brake that penalizes high similarity S_{ij} between distinct items in a batch if it exceeds a threshold $\delta = 0.95$.

$$\mathcal{L}_{\text{collapse}} = \lambda_{\text{ac}} \cdot \mathbb{E}_{i \neq j} [\max(0, S_{ij} - \delta)]$$

Defense 6: Hyperspherical Uniformity Loss

A uniformity loss, adapted from Wang & Isola (2020), encourages embeddings to spread across the hypersphere, maximizing the expressive capacity of the representation space.

$$\mathcal{L}_{\text{uniform}} = \log \mathbb{E}_{i \neq j} [\exp(-t\|\mathbf{e}_i - \mathbf{e}_j\|_2^2)]$$

3.3 The ViUniEmbed Training Dataset: A Curriculum of Complexity

The success of ViUniEmbed is critically dependent on the composition and quality of its training data. We constructed a diverse, 8.5 million sample dataset specifically designed to teach the model a spectrum of capabilities, from binary discrimination to nuanced, continuous understanding. The dataset is not merely a large collection of pairs but a deliberately structured curriculum.

Data Type	Sample Count	Percentage	Core Learning Objective
Binary Text-Pairs	800,000	9.4%	Foundational semantic relevance (Is A related to B?).
Continuous Text-Pairs	3,400,000	40.0%	Calibrated reranking and fine-grained similarity.
OCR Tasks	2,150,000	25.3%	Robust visual text extraction and grounding.
VQA Tasks	2,150,000	25.3%	Complex multimodal reasoning and comprehension.
Total	8,500,000	100%	A Unified, Universal Representation.

Strategic Rationale for Data Composition

Our data strategy is built on the principle of **progressive complexity transfer**. We hypothesize that a model must first master simpler, coarse-grained tasks before it can effectively learn more complex, fine-grained ones.

1. **Foundational Layer (Binary & OCR/VQA Tasks):** Over 50% of our dataset comprises tasks with a clear binary objective (e.g., “is this the correct text for this image?” or “is this the right answer to the question?”). This large base of binary-signal data, driven by InfoNCE and Triplet losses, forces the model to first learn a robust and well-separated embedding space. It establishes a strong semantic foundation by answering the fundamental question of relevance.
2. **Calibration Layer (Continuous Text-Pairs):** The largest single component of our dataset (40%) is dedicated to continuous similarity pairs. These samples, with scores spanning the full [0, 1] range, are introduced once the model has a stable semantic foundation. The KL-Divergence and MSE components of our loss function then act upon this stable base, “stretching” and “calibrating” the embedding space to reflect nuanced degrees of similarity. This is the critical step that bridges the gap between a standard embedding model and a true reranker.
3. **Multimodal Robustness (OCR & VQA Diversity):** The OCR and VQA subsets were intentionally curated to include a high diversity of scenarios: single and multiple images per query, and single-turn versus multi-turn conversational VQA. This forces the **Multi-Head Attention Pooling** layer to learn how to contextually aggregate information from multiple visual sources and long conversational histories, building a model that is robust to the varied and often messy inputs seen in production environments.

In essence, our dataset acts as an **implicit curriculum**. The model is simultaneously exposed to all task types, but the sheer volume of binary-signal data provides a stable “center of gravity” that prevents divergence, while the large volume of continuous-signal data provides the rich gradients necessary for learning the final, calibrated reranking capability.

4. Innovation 3: ViUniEmbed-M and Hierarchical Matryoshka Learning

ViUniEmbed-M extends the base architecture with nested representations using Matryoshka Representation Learning (Kusupati et al., 2022). A single 2048-dimensional vector is trained such that its prefixes form complete, usable embeddings at smaller dimensions (e.g., 512, 768, 1024).

This is achieved by applying the full, scheduled loss function to different prefixes of the embedding during training, with contributions weighted quadratically to prioritize the fidelity of larger dimensions:

$$\mathcal{L}_{\text{total}} = \sum_{d \in D} w_d \cdot \mathcal{L}_d, \quad \text{where } w_d = \frac{(d/D_{\text{max}})^2}{\sum_{k \in D} (k/D_{\text{max}})^2}$$

This provides unparalleled deployment flexibility, allowing practitioners to dynamically trade between retrieval speed (using shorter embeddings) and reranking accuracy (using the full embedding) from a single model and a single stored vector.

5. Commercial Advantages and Production Readiness

The unified architecture of ViUniEmbed offers substantial, quantifiable benefits over traditional pipelines.

Metric	Traditional Pipeline	ViUniEmbed	Advantage
Models Required	2-3 (e.g., Embedder + Reranker)	1	~60% less complexity
Projected Latency	50ms + 20ms = 70ms	~25ms	~2.8× faster
Projected RAM Usage	8GB + 6GB = 14GB	~6GB	~57% reduction
Vector DB Compatibility	Often poor (multi-vector)	Universal (single vector)	Works with all DBs
Calibrated Scores	No, requires second model	Yes, natively	Direct reranking

In a real-world scenario with 100 million documents, a traditional pipeline might require five GPU servers. The ViUniEmbed system could deliver superior performance with just two, representing a **60% reduction in infrastructure cost** and operational overhead.

6. Preliminary Results: A Proof of Concept

Though training on the full 8.5 million sample dataset is ongoing, results from the first 3,000 steps (representing under 3% of one epoch) are highly promising and validate our approach.

Metric	Value	Significance
R@1 (Retrieval)	0.678	Strong first-stage retrieval capability.
Spearman Correlation	0.649	Excellent calibrated ranking performance.
VQA R@1	0.999	Near-perfect visual question answering.
OCR R@1	0.807	Highly accurate optical character recognition.
Embedding sim_gap	0.411	A large, healthy gap between positive and negative similarities.

The per-dimension health of ViUniEmbed-M is also excellent, demonstrating a clear hierarchical separation of mean similarities (e.g., 0.20 for 512d vs. 0.78 for 2048d). This demonstrates the Matryoshka objective is being successfully learned without representation collapse.

Trained with 2x NVIDIA A100 GPUs. Training time: ~140 hours for ViUniEmbed, ~360 hours for ViUniEmbed-M.

7. Conclusion: A New Paradigm for Multimodal Search

ViUniEmbed and ViUniEmbed-M establish a new, more efficient paradigm for building high-performance multimodal search systems. By successfully unifying the distinct tasks of embedding and reranking into a single, calibrated vector, we eliminate the architectural fragmentation that has plagued production AI systems. Our key contributions—a robust stability architecture driven by the **Gradient Vaccine**, a powerful **multi-objective loss function**, and the flexibility of **Matryoshka learning**—provide a comprehensive blueprint for the next generation of practical, powerful, and cost-effective multimodal AI.

To accelerate progress in the field, all code, model weights, and training configurations will be made publicly available upon publication.

8. References

- Faysse, M., et al. (2024). ColPali: Efficient Document Retrieval with Vision Language Models. *arXiv:2405.16738*.
- Kusupati, A., et al. (2022). Matryoshka Representation Learning. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Wang, T., & Isola, P. (2020). Understanding Contrastive Representation Learning through Alignment and Uniformity on the Hypersphere. *Proceedings of the 37th International Conference on Machine Learning (ICML)*.
- Yasunaga, M., et al. (2024). Retrieval-Augmented Thought Process for Weak-to-Strong Reasoning. *arXiv:2405.01354*.