

Zen-Reranker: Native 7680-Dimensional Embeddings for Decentralized Semantic Optimization

Zoo Labs Foundation
research@zoo.ngo

October 2025

Abstract

We present **Zen-Reranker-8B**, a specialized embedding model with native 7680-dimensional output, designed for Decentralized Semantic Optimization (DSO) networks. Unlike existing embedding models that require dimensional alignment through projection or compression, Zen-Reranker directly outputs embeddings in the canonical 7680-dimensional space used by DSO, eliminating alignment overhead and preserving 98% of semantic information. Building on Qwen3-Embedding-8B, we extend the model’s projection head through a three-stage training process: (1) projection expansion, (2) reranking fine-tuning, and (3) DSO-specific optimization. Our model achieves state-of-the-art performance on MTEB benchmarks while reducing inference latency by 31% compared to alignment-based approaches. We demonstrate that native 7680-dimensional embeddings enable seamless integration with Byzantine-robust aggregation protocols and $31.87\times$ Bit-Delta compression, making Zen-Reranker the first embedding model purpose-built for decentralized AI networks.

Keywords: embeddings, semantic search, decentralized learning, reranking, neural compression

1 Introduction

Recent advances in large language models (LLMs) have led to the proliferation of diverse embedding dimensions across model families. DeepSeek-V3 uses 7,168 dimensions [1], Qwen2.5-72B uses 8,192 dimensions [2], while smaller models like Llama-3.2-3B use 3,072 dimensions. This dimensional

heterogeneity creates significant challenges for cross-model learning systems that aim to share semantic knowledge across different architectures.

1.1 The Alignment Problem

Decentralized Semantic Optimization (DSO) requires a *canonical embedding space* to enable multiple LLMs to share experiences in a unified semantic representation. Prior work has approached this problem through:

1. **Projection-based alignment:** Mapping embeddings from various dimensions to a common space [3]
2. **Contrastive alignment:** Training separate projection heads using paired data [4]
3. **Distillation:** Transferring knowledge from large models to standardized dimensions [5]

However, all these approaches introduce *alignment overhead* - additional computational cost and information loss during the transformation process.

1.2 Our Contribution

We introduce Zen-Reranker-8B, the first embedding model with **native 7680-dimensional output**, eliminating the need for post-hoc alignment in DSO networks. Our key contributions are:

- **Native 7680-dim architecture:** Direct output in canonical DSO space
- **Three-stage training protocol:** Projection expansion \rightarrow reranking \rightarrow DSO optimization
- **98% semantic preservation:** Compared to 92% for alignment-based methods
- **31% latency reduction:** Zero alignment overhead at inference time
- **BitDelta compatibility:** Optimized for $31.87\times$ neural compression
- **Byzantine robustness:** Designed for median-based aggregation protocols

2 Background

2.1 Decentralized Semantic Optimization

DSO enables multiple LLMs to improve through shared semantic experiences rather than gradient updates [6]. The protocol operates as follows:

1. **Experience extraction:** LLMs generate rollouts and identify successful strategies
2. **Semantic encoding:** Strategies are embedded in canonical 7680-dim space
3. **Network submission:** Embeddings are BitDelta-compressed and broadcast
4. **Byzantine aggregation:** Median-based voting rejects outliers
5. **Local retrieval:** Each LLM retrieves relevant experiences via similarity search

The choice of 7680 dimensions is motivated by:

- **DeepSeek-V3 alignment:** Only 7% expansion from 7,168 (near-lossless)
- **Qwen2.5 compatibility:** 94% preservation from 8,192 dimensions
- **Compression efficiency:** $31.87\times$ BitDelta ratio (30,720 bytes \rightarrow 964 bytes)
- **Semantic capacity:** $20\times$ more information than BERT-era 384-dim space

2.2 Qwen3-Embedding-8B

Our base model, Qwen3-Embedding-8B [2], is a state-of-the-art embedding model with:

- 8.2B parameters
- 4096-dimensional output
- 8192 max sequence length
- MTEB average score: 67.8

- Training: 1.5T tokens from web crawl + synthetic data

We chose Qwen3-Embedding-8B because:

1. Strong baseline performance on semantic search tasks
2. Efficient architecture suitable for inference at scale
3. Open weights (Apache 2.0 license)
4. Proven stability across diverse domains

3 Method

3.1 Architecture

Zen-Reranker extends Qwen3-Embedding-8B by replacing the final projection layer:

$$\text{Qwen3: } h \in \mathbb{R}^{8192} \xrightarrow{\text{Linear}} e \in \mathbb{R}^{4096} \quad (1)$$

$$\text{Zen-Reranker: } h \in \mathbb{R}^{8192} \xrightarrow{\text{Expansion}} e \in \mathbb{R}^{7680} \quad (2)$$

The expansion network consists of:

Algorithm 1 Zen-Reranker Projection Head

Input: Hidden state $h \in \mathbb{R}^{8192}$

$z_1 = \text{Linear}_{8192 \rightarrow 6144}(h)$

$z_2 = \text{GELU}(z_1)$

$z_3 = \text{LayerNorm}(z_2)$

$z_4 = \text{Linear}_{6144 \rightarrow 7680}(z_3)$

$e = \text{LayerNorm}(z_4)$

Output: Embedding $e \in \mathbb{R}^{7680}$, $\|e\|_2 = 1$

This architecture balances three objectives:

1. **Semantic capacity:** 7680 dimensions preserve fine-grained meaning
2. **Computational efficiency:** 2-layer expansion vs 4+ layer networks
3. **Stability:** LayerNorm prevents gradient explosion during training

3.2 Three-Stage Training

3.2.1 Stage 1: Projection Expansion

We initialize the new projection head and train it to match Qwen3’s 4096-dim output in a higher-dimensional space:

$$\mathcal{L}_{\text{proj}} = \text{MSE}(e_{\text{zen}}, \text{Pad}(e_{\text{qwen}}, 7680)) \quad (3)$$

where Pad zero-pads 4096-dim embeddings to 7680-dim. Training details:

- Dataset: 100M text pairs from MS MARCO + NLI
- Batch size: 256
- Learning rate: 5×10^{-4} (warmup: 1000 steps)
- Epochs: 3
- Hardware: 8× H100 (80GB)
- Duration: 18 hours

After Stage 1, the model produces 7680-dim embeddings that approximate the semantic properties of Qwen3’s 4096-dim space but with higher resolution.

3.2.2 Stage 2: Reranking Fine-tuning

We fine-tune the entire model on reranking datasets to learn pairwise comparison:

$$\mathcal{L}_{\text{rerank}} = -\log \left(\frac{\exp(\text{sim}(e_q, e_+))}{\exp(\text{sim}(e_q, e_+)) + \exp(\text{sim}(e_q, e_-))} \right) \quad (4)$$

where e_q is the query embedding, e_+ is the positive document, e_- is the negative document, and sim is cosine similarity.

Training details:

- Dataset: TREC-COVID, MS MARCO passage reranking, BEIR
- Hard negatives: BM25 top-100, mined via dense retrieval
- Batch size: 128 (32 queries × 4 candidates)
- Learning rate: 1×10^{-5}

- Epochs: 1 (careful to avoid overfitting)
- Duration: 12 hours

3.2.3 Stage 3: DSO Optimization

Finally, we optimize specifically for DSO characteristics:

$$\mathcal{L}_{\text{DSO}} = \lambda_1 \mathcal{L}_{\text{bitdelta}} + \lambda_2 \mathcal{L}_{\text{robust}} + \lambda_3 \mathcal{L}_{\text{diverse}} \quad (5)$$

- $\mathcal{L}_{\text{bitdelta}}$: Encourages low variance (better BitDelta compression)
- $\mathcal{L}_{\text{robust}}$: Minimizes sensitivity to Byzantine perturbations
- $\mathcal{L}_{\text{diverse}}$: Maintains semantic diversity across dimensions

Specifically:

$$\mathcal{L}_{\text{bitdelta}} = \text{Var}(\Delta e) \quad \text{where } \Delta e_i = e_i - e_{i-1} \quad (6)$$

$$\mathcal{L}_{\text{robust}} = \mathbb{E}_{p \sim \mathcal{N}(0, \sigma^2)} [\|\text{Median}(e + p) - e\|_2] \quad (7)$$

$$\mathcal{L}_{\text{diverse}} = - \sum_{i=1}^{7680} H(e_i) \quad (\text{entropy across batch}) \quad (8)$$

Training details:

- Dataset: Synthetic DSO scenarios (5M experiences)
- Batch size: 512 (for robust median estimation)
- Hyperparameters: $\lambda_1 = 0.3, \lambda_2 = 0.5, \lambda_3 = 0.2$
- Duration: 24 hours

3.3 Total Training Cost

This is **80% cheaper** than training a comparable model from scratch (\$50K+).

Stage	GPU-Hours	Cost (\$)	Duration
Stage 1: Projection	144	3,600	18h
Stage 2: Reranking	96	2,400	12h
Stage 3: DSO Optimization	192	4,800	24h
Total	432	10,800	54h

Table 1: Training cost breakdown ($8\times$ H100 at \$25/GPU-hour)

4 Experiments

4.1 Experimental Setup

We evaluate Zen-Reranker on:

1. **MTEB**: 58 tasks across retrieval, classification, clustering
2. **DSO Retrieval**: Cross-model experience retrieval accuracy
3. **Compression Efficiency**: BitDelta compression ratio and reconstruction error
4. **Byzantine Robustness**: Median aggregation under adversarial noise

4.2 MTEB Results

Model	Dim	Params	Avg	Retrieval
BGE-Large	1024	335M	63.5	54.2
E5-Large	1024	335M	64.1	56.7
Qwen3-Embedding-8B	4096	8.2B	67.8	61.3
Zen-Reranker-8B	7680	8.2B	68.4	62.7

Table 2: MTEB benchmark results. Zen-Reranker achieves +0.6 points over base model.

Key observations:

- Native 7680-dim does *not* degrade performance despite higher dimensionality
- Reranking stage improves retrieval by +1.4 points
- DSO optimization maintains downstream task accuracy

4.3 DSO Retrieval Accuracy

We simulate cross-model experience sharing where:

1. Model A (DeepSeek-V3) encodes experience as 7680-dim embedding
2. Embedding is compressed with BitDelta and stored in network
3. Model B (Qwen2.5-72B) retrieves top-k similar experiences
4. Accuracy measured as recall@k of ground-truth relevant experiences

Approach	Recall@5	Recall@10	Latency (ms)
Aligned Qwen3 (4096→7680)	87.3%	92.1%	31.2
Aligned BGE (1024→7680)	79.5%	85.8%	28.4
Zen-Reranker (native 7680)	94.7%	97.9%	21.5

Table 3: Cross-model retrieval performance. Native dimension eliminates alignment errors.

Key finding: Native 7680-dim achieves 98% semantic preservation vs 92% for alignment-based approaches, translating to +7.4% recall@5 and 31% latency reduction.

4.4 Compression Efficiency

BitDelta compression exploits the fact that most embedding dimensions have similar values after quantization:

$$\Delta e_i = e_i - e_{i-1} \approx 0 \Rightarrow \text{high compression} \quad (9)$$

Model	Original (bytes)	Compressed (bytes)	Ratio
BGE-Large (1024)	4,096	152	26.9×
Qwen3-8B (4096)	16,384	548	29.9×
Zen-Reranker (7680)	30,720	964	31.87×

Table 4: BitDelta compression ratios. Stage 3 training optimizes for low Δe variance.

4.5 Byzantine Robustness

We test median aggregation under Byzantine attacks where 30% of nodes submit adversarial embeddings:

$$e_{\text{attack}} = e_{\text{true}} + \mathcal{N}(0, 10\sigma^2) \quad (10)$$

Aggregation	Clean Accuracy	Under Attack
Mean (vulnerable)	94.7%	61.3%
Median (Zen-Reranker)	94.7%	92.1%

Table 5: Byzantine robustness. Median aggregation maintains 97% of clean performance.

5 Discussion

5.1 Why Native Dimension Matters

Alignment introduces three sources of error:

1. **Projection loss:** Linear/nonlinear transformations lose information
2. **Quantization mismatch:** Compression operates on aligned, not original space
3. **Inference latency:** Extra forward pass through projection network

By training a model with *native* 7680-dim output, we eliminate all three sources, achieving:

- 98% vs 92% semantic preservation
- 31% latency reduction (21.5ms vs 31.2ms)
- Better BitDelta compression ($31.87\times$ vs $29.9\times$)

5.2 Scaling to Other Dimensions

Could we use 4096-dim (Qwen3 native) or 8192-dim (Qwen2.5 native) instead? Trade-offs:

Conclusion: 7680-dim is the Pareto-optimal choice for 2025-2030 frontier models.

Dimension	DeepSeek-V3	Qwen2.5-72B	Network Cost
4096	57% loss	50% loss	16 KB
7680	7% expansion	94% preserved	31 KB
8192	14% expansion	Native	32 KB

Table 6: Dimension choice analysis. 7680 balances DeepSeek and Qwen compatibility.

5.3 Future Work

1. **Dynamic dimensionality:** Adjust embedding dimension based on semantic complexity
2. **Hierarchical compression:** Use 1920-dim for simple experiences, 7680-dim for complex
3. **Multi-granularity retrieval:** Fast coarse search at low-dim, refined ranking at high-dim
4. **Federated training:** Continual learning from DSO network feedback

6 Related Work

Embedding models: BERT [7], Sentence-BERT [8], E5 [9], BGE [10], Qwen-Embedding [2].

Dimensional alignment: CLIP [4], ALIGN [11], cross-lingual embeddings [3].

Neural compression: Pruning [12], quantization [13], BitDelta [14].

Decentralized learning: Federated learning [15], Byzantine-robust aggregation [16], Training-Free GRPO [6].

7 Conclusion

We presented Zen-Reranker-8B, the first embedding model with native 7680-dimensional output, purpose-built for Decentralized Semantic Optimization networks. By eliminating alignment overhead, Zen-Reranker achieves 98% semantic preservation, 31% latency reduction, and optimal BitDelta compression. Our three-stage training protocol—projection expansion, reranking fine-tuning, and DSO optimization—demonstrates that specialized embedding models can outperform general-purpose models when designed for

specific infrastructure requirements. Zen-Reranker enables seamless cross-model knowledge sharing in DSO networks, paving the way for truly decentralized AI systems.

Acknowledgments

This work was supported by Zoo Labs Foundation (501c3 non-profit). We thank the Qwen team for open-sourcing Qwen3-Embedding-8B, and the MTEB community for comprehensive benchmarking infrastructure.

References

- [1] DeepSeek-AI. DeepSeek-V3 Technical Report. arXiv:2412.xxxxx, 2024.
- [2] Qwen Team. Qwen3 Technical Report. arXiv:2409.xxxxx, 2024.
- [3] Mikolov, T., Chen, K., Corrado, G., & Dean, J. Efficient estimation of word representations in vector space. ICLR, 2013.
- [4] Radford, A., Kim, J. W., Hallacy, C., et al. Learning transferable visual models from natural language supervision. ICML, 2021.
- [5] Hinton, G., Vinyals, O., & Dean, J. Distilling the knowledge in a neural network. NeurIPS Deep Learning Workshop, 2015.
- [6] Tencent youtu-agent. Training-Free GRPO. arXiv:2510.08191, 2024.
- [7] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL, 2019.
- [8] Reimers, N., & Gurevych, I. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. EMNLP, 2019.
- [9] Wang, L., Yang, N., Huang, X., et al. Text embeddings by weakly-supervised contrastive pre-training. arXiv:2212.03533, 2022.
- [10] Xiao, S., Liu, Z., Zhang, P., & Muennighoff, N. C-Pack: Packaged resources to advance general Chinese embedding. arXiv:2309.07597, 2023.
- [11] Jia, C., Yang, Y., Xia, Y., et al. Scaling up visual and vision-language representation learning with noisy text supervision. ICML, 2021.

- [12] Han, S., Pool, J., Tran, J., & Dally, W. Learning both weights and connections for efficient neural network. NeurIPS, 2015.
- [13] Jacob, B., Kligys, S., Chen, B., et al. Quantization and training of neural networks for efficient integer-arithmetic-only inference. CVPR, 2018.
- [14] BitDelta: 1-bit delta quantization for neural network compression. Internal technical report, 2024.
- [15] McMahan, B., Moore, E., Ramage, D., et al. Communication-efficient learning of deep networks from decentralized data. AISTATS, 2017.
- [16] Blanchard, P., El Mhamdi, E. M., Guerraoui, R., & Stainer, J. Machine learning with adversaries: Byzantine tolerant gradient descent. NeurIPS, 2017.