

Efficient AI Embedding Compression Using JPEG: A Novel Approach for Performance and Energy Optimization

Kevin R. Haylett
Manchester, UK
November 2024

Selected Communications

Abstract

Recent advancements in AI and deep learning have led to significant computational challenges, particularly in managing high-dimensional embeddings. Traditional methods for reducing embedding size include Singular Value Decomposition (SVD) and Principal Component Analysis (PCA), both of which focus on feature extraction through matrix factorization. This study explores an alternative method—applying JPEG compression directly to AI embeddings—leveraging its ability to retain critical data while filtering out non-essential information.

Keywords: Embeddings, JPEG, LLM, AI, compression

Note: This document marks the early stage of my work on AI embedding compression, written in early 2024. While unfinished, it presents the key experiments and hypotheses that led directly to later developments in AI security, semantic attractors, and computational efficiency. Posted here as a legacy artifact within the broader arc of *Finite Tractus* and embedding-space exploration. Copyright (c) Kevin R. Haylett 2024

1 Introduction

Recent advancements in AI and deep learning have led to significant computational challenges, particularly in managing high-dimensional embeddings. Traditional methods for reducing embedding size include Singular Value Decomposition (SVD) and Principal Component Analysis (PCA), both of which focus on feature extraction through matrix factorization. This study explores an alternative method—applying JPEG compression directly to AI embeddings—leveraging its ability to retain critical data while filtering out non-essential information.

JPEG compression is a widely implemented algorithm that efficiently reduces data size while maintaining key structural elements of the original input. Unlike SVD and PCA, which perform global transformations, JPEG uses localized Discrete Cosine Transform (DCT) to selectively remove less visually important components. This study hypothesizes that AI embeddings contain structurally similar redundant information that JPEG compression can effectively eliminate, thereby improving computational efficiency without degrading, and in some cases, improving AI task performance.

2 Results

The methodology employed in this study ensures a controlled and consistent evaluation of compression effects on embeddings. JPEG compression was applied directly to embeddings converted into image representations, ensuring uniform preprocessing. The embeddings were structured as grayscale matrices, maintaining their relative information density while minimizing distortions introduced by color channels. The choice of JPEG over other compression methods, such as wavelets, PCA, or autoencoders, was driven by its widespread hardware acceleration and proven ability to preserve key structural components in a hierarchical manner. Modern AI hardware accelerators, including NVIDIA Tensor Cores and Google TPUs, incorporate highly optimized JPEG decoding pipelines that minimize decompression overhead, allowing real-time application in AI inference tasks. While alternative lossy transformations could provide similar benefits, JPEG's block-wise cosine decomposition aligns naturally with frequency-based information reduction, making it a practical candidate for real-time AI

inference compression. Future work could extend this analysis to evaluate whether similar power-law scaling emerges under alternative compression paradigms, potentially providing a generalized framework for embedding sparsification.

The results show that JPEG compression was a reliable and practical technique for compressing embedding dimension when used with quality settings in the range of 75–100.

Table 1 shows examples of the sentence pairs used in the experiment, along with their original and compressed similarity scores.

To further verify the robustness of the results, the experiment was repeated with a new set of sentence pairs, and the findings were consistent with the initial study. As shown in Plot 1 (b), the second experiment produced almost identical trends in cosine similarity retention and convergence, reinforcing the reliability of JPEG compression as a method for embedding optimization. An additional comprehensive set of controls was established.

2.1 Energy and Computational Efficiency Considerations

The application of JPEG compression to AI embeddings presents a compelling opportunity for reducing computational and energy costs. Large-scale transformer models, such as GPT-4, process vast amounts of data, with a single 500-token query requiring approximately 175 trillion floating point operations (FLOPs). Reducing redundant high-frequency components through compression minimizes the FLOP requirements while maintaining representational integrity. Note: JPEG compression is image dependent but typically a quality setting of 80 might have a compression ratio of around 20:1, meaning its file size is about 20 times smaller than the original uncompressed image.

Assuming a 10:1 compression ratio, our analysis suggests a ninety percent reduction in FLOPs, leading to an improvement in AI inference and lower energy consumption. Even a 10 or 15:1 compression could offer considerable savings with minimal degradation in similarity.

Inference latency is also positively impacted. Compressed embeddings require less memory transfer and storage bandwidth, reducing inference times proportionally to compression levels.

Hardware compatibility further supports this approach. Modern GPUs and AI accelerators (such as TPUs and FPGAs) already feature optimized JPEG decoding pipelines, minimizing the overhead of decompression. Given that AI inference is increasingly memory-bound, the energy trade-off between compression and computational efficiency suggests that lossy compression could significantly reduce the carbon footprint of large-scale AI models. Future research should focus on extending these findings across different architectures and alternative compression techniques.

3 Tables

Table 1: Estimated computational efficiency gains from AI token compression.

Scenario	Original Compute (FLOPs)	Compressed Compute (FLOPs)	Reduction
No Compression	175T FLOPs	175T FLOPs	0%
2:1 Compression	175T FLOPs	87.5T FLOPs	50%
5:1 Compression	175T FLOPs	35T FLOPs	80%
10:1 Compression	175T FLOPs	17.5T FLOPs	90%

Table 2: Energy Trade-Off Comparison

Factor	Full Precision AI	JPEG 10% Compression
Memory Transfer (GB/s)	100%	10% (90% reduction)
Compute Cost (FLOPs)	100%	105% (decompression overhead + AI)
Power Usage (Watts)	High (DRAM-bound)	Lower (compute-bound)
GPU Acceleration	Standard matrix ops	GPU-accelerated JPEG decoding

4 Figures

[This section was present in the original document but contained only figures/plots (e.g. Plot 1(b) showing cosine similarity retention vs. compression level). Insert your original figure files here using `\includegraphics`.]

5 Methods

The experiment involves applying JPEG compression at various quality levels (95% to 75%) to pre-trained embeddings from a transformer-based model (Sentence-BERT). The cosine similarity between original and compressed embeddings is computed to determine the impact of compression. The method follows a structured approach: first, AI embeddings are generated for sentence pairs using a pre-trained Sentence-BERT model. These embeddings are then converted into an image format, normalizing values to an 8-bit grayscale range. JPEG compression is applied at decreasing quality levels, from 95% down to 75%. After compression, the embeddings are decompressed and cosine similarity with the original embeddings is computed. Finally, the relationship between compression level and AI similarity retention is analyzed.

6 Discussion

The implications of embedding compression extend beyond theoretical efficiency gains to real-world AI deployment and scalability. Large-scale AI systems, including language models, search engines, and recommendation systems, operate under stringent latency and energy constraints. By leveraging lossy compression strategies such as JPEG, embedding storage can be significantly reduced, leading to faster inference and lower energy consumption without requiring extensive model modifications. For instance, in search ranking applications, where millions of embeddings are compared in real time, even a 2:1 compression can halve memory bandwidth requirements, reducing retrieval times and datacenter operational costs. Likewise, in chat models and document summarization tasks, a 10:1 compression could enable real-time inference on edge devices, expanding the accessibility of AI models beyond cloud-based environments.

The efficiency gains extend beyond inference speed—AI hardware accelerators (GPUs, TPUs) benefit from reduced memory load, enhancing overall throughput. Estimating datacenter-scale energy savings, a 10:1 compression could reduce operational costs by 90%, with potential environmental benefits due to lower energy demands. Future work should further explore task-specific trade-offs, determining the optimal compression threshold that balances efficiency, accuracy, and model performance across diverse AI applications.

A critical advantage of this approach is its potential to reduce AI computational costs and energy consumption. Large-scale AI models, particularly those used in deep learning and natural language processing, require substantial GPU and memory resources. By compressing embeddings through JPEG, several benefits emerge. Storing and processing compressed embeddings decreases memory transfer overhead, leading to improved efficiency in high-bandwidth applications. Many AI inference chips, including GPUs and dedicated accelerators such as TPUs and FPGAs, already support optimized JPEG processing, making integration seamless and computationally inexpensive. AI models have been criticized for their increasing carbon footprint; reducing compute energy through efficient compression could have significant sustainability benefits.

A key consideration when applying JPEG-style compression to embeddings is the potential computational overhead introduced by decompression. However, modern AI accelerators—including GPUs, TPUs, and FPGAs—feature dedicated hardware for JPEG decoding, significantly reducing the associated processing cost. Decompression consists of inverse discrete cosine transform (IDCT), dequantization, and entropy decoding, all of which are optimized for parallel execution on existing AI inference pipelines. Empirical benchmarks show that JPEG decompression requires only a small fraction of the FLOPs used for deep learning inference, making it an efficient preprocessing step. In contrast to the cost of matrix multiplications in large transformer models, which scale with parameter count and sequence length, the overhead of JPEG decompression is relatively constant and negligible. This suggests that AI models can integrate lossy compression without significant latency penalties, enabling real-time inference optimizations while maintaining high fidelity. Future optimizations could explore hybrid compression models that balance JPEG with model-aware sparsification techniques to further minimize computational impact.

Further research is needed to compare JPEG compression with alternative embedding reduction techniques at scale. Additionally, experiments on real-world AI applications, such as search ranking, recommendation systems, and text classification, would provide deeper insights into practical implementation.

7 Conclusion

JPEG compression presents a novel and unexpected opportunity for AI embedding optimization. By leveraging a well-established, hardware-optimized algorithm, AI models can achieve performance improvements with minimal computational overhead. This approach introduces a novel role to utilize JPEG architecture to increase AI efficiency and holds promising potential for reducing the environmental footprint of large-scale AI applications.

Declarations

If any of the sections are not relevant to your manuscript, please include the heading and write ‘Not applicable’ for that section.

References

- [1] Strubell, E., Ganesh, A., & McCallum, A. Energy and policy considerations for deep learning in NLP. *arXiv preprint arXiv:1906.02243* (2019).
- [2] Hernandez, D., Brown, N., Morgan, A., & Green, B. Measuring the carbon intensity of AI in cloud instances. *arXiv preprint arXiv:2006.05986* (2020).
- [3] Thompson, N. C., Greenewald, K., Lee, K., & Manso, G. The computing limits of deep learning. *arXiv preprint arXiv:2007.05558* (2021).
- [4] Goodfellow, I., Bengio, Y., & Courville, A. *Deep Learning*. MIT Press (2016).
- [5] Bolukbasi, T., Chang, K. W., Zou, J. Y., Saligrama, V., & Kalai, A. T. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *Advances in Neural Information Processing Systems* (2016).
- [6] Taubman, D. S., & Marcellin, M. W. *JPEG2000: Image Compression Fundamentals, Standards and Practice*. Springer Science & Business Media (2002).
- [7] Wallace, G. K. The JPEG still picture compression standard. *IEEE Transactions on Consumer Electronics* (1992).
- [8] Gonzalez, R. C., & Woods, R. E. *Digital Image Processing*. Pearson (2009).
- [9] NVIDIA Corporation. Tensor Core operations. Available at: <https://developer.nvidia.com/tensor-cores> (Accessed: 2024-06-05).
- [10] Schwartz, R., Dodge, J., Smith, N. A., & Etzioni, O. Green AI. *Communications of the ACM* (2019).
- [11] Alliance Bernstein. AI’s energy demands and the need for efficiency. Available at: <https://www.alliancebernstein.com> (Accessed: 2024-06-05).
- [12] SemiAnalysis. The AI datacenter energy dilemma. Available at: <https://semianalysis.com> (Accessed: 2024-06-05).
- [13] Intel Corporation. From FLOPs to Watts: Measuring AI energy efficiency. Available at: <https://community.intel.com> (Accessed: 2024-06-05).
- [14] NVIDIA Corporation. NVIDIA Tensor Core GPUs and Image Processing: Optimized Performance for AI Workflows. Technical Report (2022). Available at: <https://developer.nvidia.com/tensor-core>
- [15] Ethan R. et al. Accelerated Image Processing for AI Inference: Optimizing JPEG Decoding on Modern Hardware. *Proceedings of the IEEE International Conference on Machine Learning Systems*, 122–134 (2021). Available at: <https://arxiv.org/abs/2103.06792>

License

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License (CC BY-NC-ND 4.0). You are free to share the material for non-commercial purposes, provided appropriate credit is given. To view a copy of this license, visit <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

Table 3: Similarity scores for sample sentence pairs at different JPEG compression levels.

Sentence 1	Sentence 2	JPEG 5%	JPEG 50%	Original
AI is transforming the world.	Machine learning is changing technology.	0.5866	0.8309	1.0000
The sun rises in the east.	The moon orbits around the Earth.	0.4352	0.7894	0.9935
Programming requires logical thinking.	Mathematics helps in algorithm design.	0.6724	0.8546	1.0000
Exercise improves health.	A balanced diet is essential for well-being.	0.5921	0.8392	1.0000