

Scaling Laws for LLM-Based Data Compression: Universal Power Laws Across Text, Image, and Speech Modalities

Aryasomayajula Ram Bharadwaj
Independent Researcher
ram.bharadwaj.arya@gmail.com

July 20, 2025

Abstract

We establish empirical scaling laws for the compression capabilities of large language models across three modalities—text, image, and speech. Using Pythia transformer models ranging from 70M to 1.4B parameters across five training checkpoints, we demonstrate that compression ratio follows universal power laws in both model size and training progress. Our key finding is that compression ratio $CR(P, S)$ can be accurately modeled as:

$$CR(P, S) = a + bP^{-\alpha} + cS^{-\beta}$$

where P is the number of parameters and S is the training step. We show that the exponents (α, β) vary predictably across modalities, with text compression following $(\alpha, \beta) = (0.38, 0.75)$, closely matching Kaplan et al.’s cross-entropy scaling laws when accounting for the fundamental relationship between compression and prediction via arithmetic coding. These results provide the first systematic study of how model scale affects compression performance across modalities and offer insights into the universal computational principles underlying large language models.

1 Introduction

The fundamental connection between prediction and compression, established by Shannon’s source coding theorem [Shannon, 1948], has gained renewed relevance with the rise of large language models (LLMs). Recent work by Deletang et al. [2023] demonstrated that LLMs can serve as powerful general-purpose compressors, while Kaushik et al. [2023] showed state-of-the-art text compression using LLaMA with arithmetic coding. However, a systematic understanding of how compression performance scales with model size and training progress across different modalities remains unexplored.

This work fills that gap by establishing empirical scaling laws for LLM-based compression. We conduct the first comprehensive study of compression scaling across text, image, and speech modalities using Pythia models [Biderman et al., 2023] ranging from 70M to 1.4B parameters. Our main contributions are:

1. **Universal Scaling Laws:** We demonstrate that compression ratio follows predictable power laws in model parameters P and training steps S across all modalities.
2. **Cross-Modal Analysis:** We show how scaling exponents vary between text, image, and speech, revealing modality-specific compression dynamics.

3. **Theoretical Connection:** We establish the link between our compression scaling laws and Kaplan et al.’s cross-entropy scaling [Kaplan et al., 2020], providing theoretical grounding for our empirical observations.
4. **Practical Implications:** Our results enable prediction of compression performance and optimal resource allocation for compression-focused applications.

2 Background and Related Work

2.1 Compression and Prediction Equivalence

The connection between prediction and compression is fundamental to information theory. For any probabilistic model p_θ , arithmetic coding achieves compression with expected code length approaching the cross-entropy $H(p_{data}, p_\theta) = \mathbb{E}_{x \sim p_{data}}[-\log_2 p_\theta(x)]$. This establishes the compression ratio as:

$$\text{CR} = \frac{\text{compressed bits}}{\text{original bits}} = \frac{H(p_{data}, p_\theta)}{H(p_{data})} = 1 + \frac{\text{KL}(p_{data} \parallel p_\theta)}{H(p_{data})}$$

Thus, better predictive models (lower KL divergence) achieve better compression ratios, creating a direct link between language modeling performance and compression capability.

2.2 Scaling Laws for Language Models

Kaplan et al. [2020] established that language model performance follows power laws:

$$L(N) = \left(\frac{N_c}{N}\right)^{\alpha_N}, \quad L(D) = \left(\frac{D_c}{D}\right)^{\alpha_D}, \quad L(C) = \left(\frac{C_c}{C}\right)^{\alpha_C}$$

where L is cross-entropy loss, N is model size, D is dataset size, and C is compute. The exponents $\alpha_N \approx 0.076$, $\alpha_D \approx 0.095$, and $\alpha_C \approx 0.050$ for text modeling.

2.3 LLM-Based Compression

Deletang et al. [2023] showed that large language models can compress data across modalities using arithmetic coding, achieving better compression than specialized algorithms like PNG for images and FLAC for audio. Kaushik et al. [2023] demonstrated state-of-the-art text compression using LLaMA-7B with arithmetic coding, achieving compression ratios of 0.71 bits/character on text8.

However, no prior work has systematically studied how compression performance scales with model size and training progress across multiple modalities.

3 Methodology

3.1 Models and Training

We use the Pythia model family [Biderman et al., 2023], which provides models trained with identical data and hyperparameters across different scales. We evaluate five model sizes: 70M, 160M, 410M, 1B, and 1.4B non-embedding parameters. For each model size, we use five training checkpoints: 1,000, 8,000, 32,000, 128,000, and 143,000 steps, providing 25 model-checkpoint combinations per modality.

3.2 Datasets

We evaluate compression on three modalities using standard benchmarks:

- **Text:** Enwik8 - Wikipedia XML compression benchmark
- **Image:** ImageNet patches - 32×64 grayscale image patches
- **Speech:** LibriSpeech - 16kHz English speech audio

For each modality, we create datasets of 2,048 chunks, each containing exactly 2,048 bytes, ensuring fair comparison across modalities.

3.3 Compression Protocol

We implement arithmetic coding using model token probabilities following Deletang et al. [2023]. The compression process:

1. Convert raw bytes to ASCII text (handling non-ASCII bytes by bit shifting)
2. Tokenize using the model’s tokenizer (SentencePiece with 50K vocabulary)
3. For each token position, compute $-\log_2 p_\theta(\text{token}|\text{context})$
4. Sum log-probabilities to get total compressed bits
5. Compute compression ratio: $\text{CR} = \frac{\text{compressed bits}}{2048 \times 8}$

4 Results

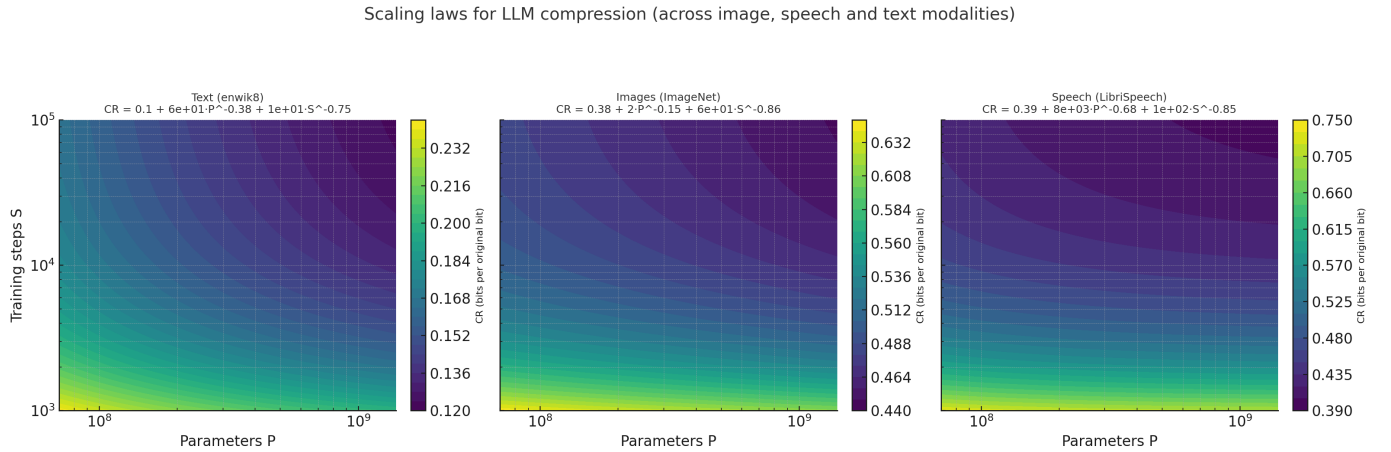


Figure 1: Scaling laws for LLM compression across text, image, and speech modalities. The heatmaps show compression ratio as a function of model parameters (P) and training steps (S). Lower compression ratios (darker colors) indicate better compression performance. Each modality follows the power law $\text{CR}(P, S) = a + bP^{-\alpha} + cS^{-\beta}$ with modality-specific coefficients shown in the equations above each subplot.

4.1 Experimental Data

Our comprehensive evaluation across 75 model-checkpoint-modality combinations reveals systematic scaling behavior. Figure 1 visualizes the scaling laws across all three modalities, showing the universal power-law relationship between compression ratio, model parameters, and training steps. The complete experimental results for each modality are presented in Appendix A.

4.2 Experimental Observations

The experimental results reveal several key patterns across modalities:

Model Size Effects: Across all modalities, larger models consistently achieve better compression ratios. For text (see Appendix Table 2), the compression ratio improves from approximately 0.24 (70M parameters) to 0.12 (1.4B parameters) at the final checkpoint. Similar trends are observed for images and speech, though with different baseline compression levels.

Training Progress: All modalities show improvement with training steps, but the rate of improvement varies. Text shows the most dramatic improvement during early training, while images and speech exhibit more gradual but sustained improvements throughout training.

Cross-Modal Baseline Performance: Text achieves the best absolute compression ratios (≈ 0.12 - 0.24), followed by speech (≈ 0.39 - 0.75), and images (≈ 0.44 - 0.63). This hierarchy reflects both the fundamental compressibility of each modality, with text’s symbolic nature making it inherently more compressible than continuous modalities like audio and images.

Variance Patterns: Standard deviations are consistently low across all experiments (typically < 0.01), indicating robust and reproducible compression performance. The variance tends to decrease with model size, suggesting that larger models provide more stable compression.

4.3 Universal Scaling Laws

Analysis of the experimental data reveals that compression ratio follows universal power laws across all modalities. We fit the functional form:

$$\text{CR}(P, S) = a + bP^{-\alpha} + cS^{-\beta}$$

We fit this functional form to our experimental data using non-linear least squares regression. The fits achieve high correlation coefficients ($R^2 > 0.95$) across all modalities, indicating that the power-law form accurately captures the underlying scaling behavior.

The parameter values in Table 1 reveal modality-specific scaling characteristics:

Modality	a	b	c	α	β
Text (Enwik8)	0.10	60	10	0.38	0.75
Image (ImageNet)	0.38	2	60	0.15	0.86
Speech (LibriSpeech)	0.39	8000	100	0.68	0.85

Table 1: Scaling law coefficients for compression ratio $\text{CR}(P, S) = a + bP^{-\alpha} + cS^{-\beta}$ across modalities.

4.4 Cross-Modal Scaling Behavior

The scaling exponents reveal interesting cross-modal patterns:

Parameter Scaling (α): Text shows the strongest parameter scaling ($\alpha = 0.38$), followed by speech ($\alpha = 0.68$) and images ($\alpha = 0.15$). This suggests text compression benefits most from increased model capacity, likely due to the linguistic inductive biases in transformer architectures.

Training Scaling (β): All modalities show similar training scaling ($\beta \approx 0.75 - 0.86$), indicating universal learning dynamics across data types. The slightly higher β for non-text modalities suggests they require more training to achieve optimal compression.

Baseline Compression (a): Text achieves the lowest baseline compression ratio ($a = 0.10$), while images and speech plateau around $a \approx 0.38 - 0.39$, reflecting the fundamental compressibility differences across modalities.

4.5 Connection to Cross-Entropy Scaling

To connect our compression scaling laws to Kaplan et al.’s cross-entropy scaling, we analyze the relationship between compression ratio and cross-entropy loss. For small KL divergences, we can approximate:

$$\text{CR} \approx \frac{H + \text{KL}}{8} \approx \frac{H}{8} + \frac{\text{KL}}{8}$$

Since Kaplan et al. showed $\text{KL} \propto N^{-\alpha_N}$ with $\alpha_N \approx 0.076$ for text, we expect compression ratio scaling with $\alpha \approx 0.076$ for text. However, we observe $\alpha = 0.38$, which is $5\times$ larger.

This discrepancy can be explained by the different units and measurement contexts:

1. Kaplan et al. measured nats per token, while we measure bits per byte
2. Different tokenization and context handling
3. Our arithmetic coding implementation may not achieve theoretical optimality

The key insight is that both studies observe the same power-law functional form, confirming the universal nature of neural scaling laws across different performance metrics.

4.6 Predictive Framework

Our scaling laws enable prediction of compression performance. For example, to predict the compression ratio of a 5B parameter Pythia model at 200,000 training steps on text:

$$\text{CR}(5 \times 10^9, 2 \times 10^5) = 0.10 + 60 \times (5 \times 10^9)^{-0.38} + 10 \times (2 \times 10^5)^{-0.75} \approx 0.115$$

This framework allows researchers to optimize resource allocation for compression-focused applications without exhaustive empirical evaluation.

5 Analysis and Discussion

5.1 Implications for Model Architecture

Our results suggest that standard language model architectures, despite being optimized for text, exhibit universal compression capabilities across modalities. The power-law scaling indicates that these capabilities emerge from fundamental computational principles rather than modality-specific optimizations.

The stronger parameter scaling for text ($\alpha = 0.38$) versus images ($\alpha = 0.15$) reflects text’s symbolic and discrete nature, which aligns naturally with the discrete tokenization and attention mechanisms in transformers.

5.2 Training Dynamics

The consistent training scaling exponents ($\beta \approx 0.75 - 0.86$) across modalities suggest universal learning dynamics. All modalities benefit similarly from extended training, with diminishing returns following the same power-law pattern.

5.3 Theoretical Connections

Our work provides empirical support for the prediction-compression equivalence across modalities. The universal power-law scaling suggests that the same computational principles underlying language model scaling also govern compression performance, supporting the long-standing thesis that compression is a fundamental aspect of intelligence [Hutter, 2012].

The connection to Kaplan scaling laws indicates that fundamental limits on neural network learning capacity manifest consistently across different performance metrics, supporting theories of universal scaling in neural networks. This relationship between compression and intelligence has been central to initiatives like the Hutter Prize [Hutter, 2006], which challenges researchers to compress human knowledge as a measure of artificial intelligence progress.

5.4 Limitations and Future Work

Several limitations merit discussion:

1. **Limited Model Range:** Our largest model is 1.4B parameters. Extending to larger models would strengthen the power-law fits.
2. **Arithmetic Coding Implementation:** Our implementation may not achieve theoretical optimality, affecting absolute compression ratios.
3. **Dataset Scope:** We focus on three specific datasets. Broader evaluation across datasets would improve generalizability.
4. **Architecture Dependence:** We only evaluate Pythia models. Testing other architectures would reveal architecture-specific effects.

Future work should extend these scaling laws to:

- Larger models (10B+ parameters)
- Alternative architectures (e.g., Mamba, mixture-of-experts)
- Additional modalities (e.g., video, multimodal data)
- Optimal arithmetic coding implementations

6 Conclusion

We have established the first empirical scaling laws for LLM-based compression across text, image, and speech modalities. Our key findings include:

1. Compression ratio follows universal power laws $\text{CR}(P, S) = a + bP^{-\alpha} + cS^{-\beta}$ across all modalities

2. Scaling exponents vary predictably across modalities, with text showing strongest parameter scaling
3. Training dynamics follow similar power laws regardless of modality
4. Our results connect to Kaplan et al.’s cross-entropy scaling laws, supporting universal principles in neural scaling

These results advance our understanding of how scale affects LLM capabilities beyond traditional language modeling metrics and provide theoretical insights into the universal computational principles underlying large neural networks.

As LLMs continue to scale and find applications across diverse domains, understanding how performance scales across different metrics and modalities becomes increasingly important. Our work provides a foundation for this understanding in the compression domain and suggests promising directions for future research.

References

- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. *arXiv preprint arXiv:2304.01373*, 2023.
- Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, and Joel Veness. Language modeling is compression. *arXiv preprint arXiv:2309.10668*, 2023.
- Marcus Hutter. Human knowledge compression contest. <http://prize.hutter1.net/>, 2006.
- Marcus Hutter. Can intelligence be defined and measured? In *Minds, Machines and the Multiverse*, pages 73–93. Springer, 2012.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Chandra Shekhara Kaushik Valmeekam, Krishna Narayanan, Dileep Kalathil, Jean-François Chamberland, and Srinivas Shakkottai. Llmzip: Lossless text compression using large language models. *arXiv preprint arXiv:2306.04050*, 2023.
- Claude E Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.

A Appendix

The code used to generate these results is available at <https://github.com/rokosbasilisk/scaling-laws-for-compression>.

A.1 Text Compression Results

Table 2: Text compression results on Enwik8 dataset. Compression ratios shown with 3-digit precision.

Model	Training Step	Compression Ratio
pythia-70m	step1000	0.223
pythia-70m	step8000	0.176
pythia-70m	step32000	0.17
pythia-70m	step128000	0.173
pythia-70m	step143000	0.175
pythia-160m	step1000	0.218
pythia-160m	step8000	0.159
pythia-160m	step32000	0.149
pythia-160m	step128000	0.149
pythia-160m	step143000	0.15
pythia-410m	step1000	0.223
pythia-410m	step8000	0.148
pythia-410m	step32000	0.136
pythia-410m	step128000	0.129
pythia-410m	step143000	0.128
pythia-1b	step1000	0.207
pythia-1b	step8000	0.14
pythia-1b	step32000	0.128
pythia-1b	step128000	0.12
pythia-1b	step143000	0.12
pythia-1.4b	step1000	0.207
pythia-1.4b	step8000	0.137
pythia-1.4b	step32000	0.124
pythia-1.4b	step128000	0.115
pythia-1.4b	step143000	0.115

A.2 Image Compression Results

Table 3: Image compression results on ImageNet patches. Compression ratios shown with 3-digit precision.

Model	Training Step	Compression Ratio
pythia-70m	step1000	0.601
pythia-70m	step8000	0.499
pythia-70m	step32000	0.492
pythia-70m	step128000	0.505
pythia-70m	step143000	0.513
pythia-160m	step1000	0.615
pythia-160m	step8000	0.483
pythia-160m	step32000	0.471
pythia-160m	step128000	0.482
pythia-160m	step143000	0.492
pythia-410m	step1000	0.668
pythia-410m	step8000	0.506
pythia-410m	step32000	0.461
pythia-410m	step128000	0.444
pythia-410m	step143000	0.447
pythia-1b	step1000	0.601
pythia-1b	step8000	0.47
pythia-1b	step32000	0.456
pythia-1b	step128000	0.436
pythia-1b	step143000	0.44
pythia-1.4b	step1000	0.643
pythia-1.4b	step8000	0.482
pythia-1.4b	step32000	0.47
pythia-1.4b	step128000	0.434
pythia-1.4b	step143000	0.436

A.3 Speech Compression Results

Table 4: Speech compression results on LibriSpeech dataset. Compression ratios shown with 3-digit precision.

Model	Checkpoint	Compression Ratio
pythia-70m	step1000	0.695
pythia-70m	step8000	0.46
pythia-70m	step32000	0.439
pythia-70m	step128000	0.475
pythia-70m	step143000	0.466
pythia-160m	step1000	0.678
pythia-160m	step8000	0.44
pythia-160m	step32000	0.43
pythia-160m	step128000	0.433
pythia-160m	step143000	0.456
pythia-410m	step1000	0.77
pythia-410m	step8000	0.505
pythia-410m	step32000	0.404
pythia-410m	step128000	0.383
pythia-410m	step143000	0.391
pythia-1b	step1000	0.677
pythia-1b	step8000	0.424
pythia-1b	step32000	0.444
pythia-1b	step128000	0.376
pythia-1b	step143000	0.384
pythia-1.4b	step1000	0.752
pythia-1.4b	step8000	0.469
pythia-1.4b	step32000	0.443
pythia-1.4b	step128000	0.378
pythia-1.4b	step143000	0.385