# Scaling Laws in Generative AI: How Model Size and Data Influence Performance and Cost

Chandan Singh Troughia

*Staff Machine Learning Engineer*
*CVS Health*
ORCID: 0009-0004-2921-6355

*Abstract*—The exponential growth in generative AI capabilities has been governed by predictable mathematical relationships known as scaling laws, fundamentally reshaping how we approach model development and deployment. This paper examines the complex relationships between model size, training data, performance, and cost in generative artificial intelligence systems, spanning models from millions to hundreds of billions of parameters trained on datasets ranging from gigabytes to petabytes. Through systematic analysis of empirical data and case studies, we trace the evolution of scaling laws from the seminal work of Kaplan et al. to the Chinchilla paradigm shift—which revealed that previous models were undertrained by orders of magnitude—and recent developments, providing a comprehensive framework for understanding how these factors interact.

Our investigation reveals that while performance improvements follow power-law relationships with both model size and data quantity, the optimal balance between these factors continues to evolve with significant economic implications. We explore emergent capabilities that appear at specific scale thresholds, the critical role of data quality in determining model performance, comprehensive evaluation methodologies that capture scaling behaviors, and economic considerations that shape practical deployment decisions. The analysis demonstrates that compute-optimal training strategies can achieve equivalent performance with substantially reduced computational costs, fundamentally altering the economics of AI development.

By synthesizing insights across these dimensions, we offer evidence-based guidance for researchers and practitioners navigating the trade-offs inherent in generative AI development and deployment. This holistic perspective on scaling laws provides valuable direction for advancing more capable, efficient, and sustainable AI systems in an era of increasing computational demands.

*Keywords*—Scaling Laws, Generative AI, Model Size, Training Data, Performance, Cost, Large Language Models, Data Quality, Evaluation, Economic Analysis, Chinchilla, Emergent Abilities, Compute-Optimal Training

## I. INTRODUCTION

The field of artificial intelligence has witnessed an unprecedented transformation with the advent of generative AI models. From GPT-3's 175 billion parameters in 2020 to models exceeding a trillion parameters today, the scale of these systems has grown exponentially, accompanied by equally dramatic improvements in capability. These models, capable of producing human-like text, images, and other content, have revolutionized how we interact with technology and approach complex problems. At the heart of this revolution lies a fundamental question: How do we optimize these models for maximum performance while managing computational and financial costs?

This question has given rise to the study of *scaling laws* in generative AI—predictable mathematical relationships, typically following power-law distributions, that quantify how model performance scales with increases in parameters ($N$), training data ($D$), and compute ($C$).

Scaling laws have emerged as a critical framework for understanding the behavior of large language models (LLMs) and other generative AI systems. They provide insights into the relationships between model parameters, training data, computational requirements, and performance outcomes. These relationships are not merely academic curiosities but have profound practical implications for AI research, development, and deployment. With training costs for frontier models now exceeding $100 million and inference serving billions of queries daily, understanding these relationships has become essential for sustainable AI development. Recent breakthroughs in multimodal models like GPT-4V, DALL-E 3, and Gemini have further validated these scaling principles while revealing new complexities in cross-modal scaling behaviors.

### A. Contributions and Scope

The key contributions of this work include:

1) A comprehensive synthesis of scaling law evolution from Kaplan to post-Chinchilla developments
2) Analysis of emergent capability thresholds and their practical implications
3) Frameworks for evaluating the economic trade-offs in model scaling decisions
4) Actionable guidance for optimizing the critical balance between performance and cost in generative AI systems

### B. Paper Organization

This paper is organized as follows: Section II provides background on scaling laws evolution and theoretical foundations. Section III examines model size impacts on performance, including emergent capabilities and architectural considerations. Section IV investigates data quantity and quality effects, exploring optimal data-to-parameter ratios and curation strategies. Section V covers comprehensive evaluation methodologies for generative AI systems. Section VI presents detailed cost analysis frameworks encompassing training, inference, and

total cost of ownership considerations. Section VII concludes with synthesis of key insights and future directions.

Throughout this paper, we aim to provide a holistic view of scaling laws that balances theoretical understanding with practical applications. By synthesizing insights from recent research and industry practices, we offer guidance for researchers, developers, and decision-makers navigating the complex landscape of generative AI.

## II. BACKGROUND AND LITERATURE REVIEW

### A. The Evolution of Scaling Laws in Generative AI

The study of scaling laws in artificial intelligence has emerged as a cornerstone of modern AI research, providing crucial insights into how model performance relates to computational resources, model architecture, and training data. This section traces the historical development of scaling laws in generative AI, highlighting key milestones and paradigm shifts that have shaped our understanding of these relationships.

*1) Early Observations and Empirical Findings:* As shown in Figure 1, the evolution of scaling laws can be traced through several distinct paradigm shifts that have fundamentally changed our understanding of optimal model training strategies.
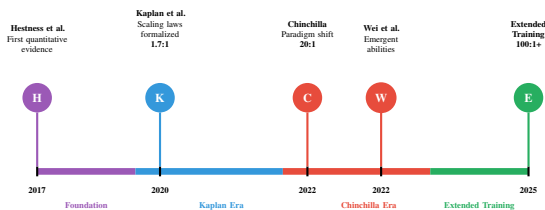


Fig. 1: Timeline of Scaling Laws Milestones in Generative AI

Prior to the formalization of scaling laws, researchers had observed that increasing model size and training data generally led to improved performance. However, these observations remained largely qualitative until the late 2010s. The deep learning revolution, catalyzed by breakthroughs in neural network architectures and training methodologies, created the conditions for more systematic investigations of scaling behavior.

Early work by Hestness et al. [6] provided some of the first quantitative evidence of power-law relationships between dataset size and model performance across various domains, including language modeling, machine translation, and image classification. Their research suggested that model performance improves as a power-law function of training set size, with the exponent varying by task and model architecture. This foundational work established the empirical basis for what would later become a comprehensive theoretical framework.

*2) The Kaplan Paradigm: Power Laws and Compute-Optimal Training:* The field of scaling laws was formalized and significantly advanced by Kaplan et al. [8] in their seminal paper "Scaling Laws for Neural Language Models." This groundbreaking work established clear mathematical relationships between model size, dataset size, and computational budget for autoregressive language models. The authors identified three key power-law relationships:

1) **Model size scaling**: Performance improves as a power-law with model size (number of parameters), with diminishing returns as models grow larger
2) **Data scaling**: Performance improves as a power-law with dataset size, again with diminishing returns
3) **Compute-optimal scaling**: Given a fixed computational budget, there exists an optimal allocation between model size and training tokens

The Kaplan scaling laws suggested that for a fixed compute budget, the optimal model should use approximately 1.7 tokens per parameter during training. This finding provided a principled approach to balancing model size and training data, guiding researchers and practitioners in resource allocation decisions.

*3) The Chinchilla Paradigm Shift: Rebalancing Parameters and Data:* In 2022, Hoffmann et al. [7] published "Training Compute-Optimal Large Language Models," introducing what became known as the "Chinchilla scaling laws." This work challenged and refined the Kaplan paradigm by demonstrating that most large language models were significantly undertrained relative to their size.
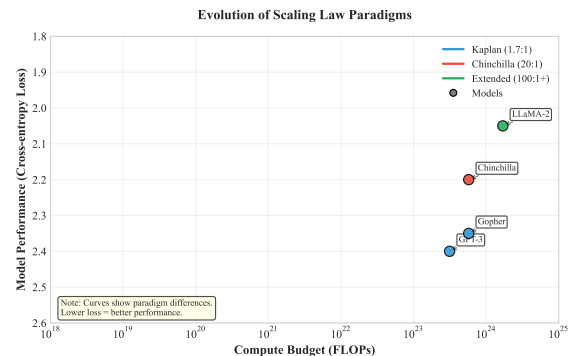


Fig. 2: Evolution of Scaling Law Paradigms. This figure illustrates the fundamental differences between scaling paradigms, showing how the Chinchilla approach achieves better performance than the Kaplan paradigm for the same compute budget.

As illustrated in Figure 2, the Chinchilla research suggested that the optimal ratio of training tokens to parameters should be approximately 20:1, substantially higher than the 1.7:1 ratio implied by Kaplan's work. The Chinchilla model, with 70 billion parameters trained on 1.4 trillion tokens, outperformed models with many more parameters that had been trained on less data, demonstrating that the field had been overemphasizing model size at the expense of training data volume.

*4) Recent Developments and Refinements:* Since the Chinchilla paper, several research groups have further refined our understanding of scaling laws, as summarized in Table I:

### B. Theoretical Foundations of Scaling Laws

The empirical observations of scaling laws have prompted theoretical investigations into why these patterns emerge.

TABLE I: Key Scaling Laws Papers Comparison

| Paper | Year | Key Finding | Ratio | Model | Impact |
|-------|------|-------------|-------|-------|--------|
| Hestness et al. | 2017 | First quantitative evidence | — | Various | Foundation |
| Kaplan et al. | 2020 | Formalized scaling laws | 1.7:1 | GPT-3 | Paradigm |
| Hoffmann et al. | 2022 | Models undertrained | 20:1 | Chinchilla | Shift |
| Wei et al. | 2022 | Emergent abilities | Varies | PaLM | Capabilities |
| Recent work | 2024+ | Extended training | 100:1+ | Llama 2/3 | Optimization |

**Note:** Ratios are tokens:parameters. Colors: Kaplan , Chinchilla , Recent .

Several frameworks have been proposed to explain the power-law relationships observed in practice.

*1) Statistical Learning Theory:* From the perspective of statistical learning theory, scaling laws reflect the trade-off between approximation error and estimation error. This manifests as the classic bias-variance trade-off: larger models reduce bias but may increase variance without sufficient data, explaining why both model size and dataset size must scale together for optimal performance.

*2) Information-Theoretic Perspectives:* Information theory provides another lens for understanding scaling laws through the minimum description length (MDL) principle, where better models achieve superior compression of the training distribution, connecting scaling laws to fundamental limits in information processing.

*3) Neural Scaling Laws as Phase Transitions:* Some researchers propose that power-law scaling behavior may be related to phase transitions in complex systems, where the emergence of new capabilities at certain scale thresholds represents critical points in the model's representational capacity.

## C. Practical Implications and Industry Adoption

The discovery and refinement of scaling laws have had profound practical implications for AI research and development:

1) **Resource Allocation**: Organizations can make more informed decisions about how to allocate computational resources between model size and training data, optimizing for specific performance targets and budget constraints.
2) **Research Roadmaps**: Academic and industrial research labs have used scaling laws to project future performance improvements and set research agendas, enabling more strategic planning of AI development initiatives.
3) **Architectural Innovations**: Understanding scaling behavior has motivated architectural innovations designed to improve parameter efficiency, such as sparse models, mixture-of-experts approaches, and more efficient attention mechanisms.
4) **Economic Considerations**: Scaling laws have informed economic analyses of AI development, helping to predict costs and benefits of different development strategies and enabling more accurate return-on-investment calculations.
5) **Environmental Impact**: By optimizing the balance between model size and training data, organizations can reduce the environmental footprint of AI training while maintaining performance, contributing to more sustainable AI development practices.

The evolution of scaling laws represents a remarkable example of how empirical observations can lead to theoretical insights and practical applications. As we continue to refine our understanding of these relationships, we gain powerful tools for guiding the development of more capable, efficient, and sustainable AI systems.

## III. IMPACT OF MODEL SIZE ON PERFORMANCE

The relationship between model size and performance represents one of the most studied aspects of scaling laws in generative AI. This section examines how increasing the number of parameters in neural network models affects various performance metrics, the emergence of new capabilities at specific scale thresholds, and the patterns of diminishing returns observed as models grow larger. These relationships validate the theoretical foundations discussed in Section II and provide crucial insights for the economic considerations explored in Section VI.

## A. The Scaling Relationship Between Parameters and Performance

*1) Fundamental Scaling Patterns:* Empirical evidence consistently demonstrates that model performance improves as a power-law function of model size, typically measured by the number of parameters. This relationship can be expressed mathematically as:

$$L(N) \approx (N_0/N)^{\alpha} \tag{1}$$

Where:

- $L(N)$ is the loss (lower is better) for a model with $N$ parameters
- $N_0$ is a constant
- $\alpha$ is the scaling exponent (typically between 0.05 and 0.1 for language models)

This power-law relationship has been observed across multiple orders of magnitude, from models with millions of parameters to those with hundreds of billions or even trillions of parameters. For language models, empirical studies report typical $\alpha$ values ranging from 0.05-0.1, with GPT-style models showing $\alpha \approx 0.076$ (Kaplan et al. [8]), while recent studies on vision-language models report $\alpha \approx 0.08-0.12$. The consistency of this pattern suggests a fundamental property of how neural networks learn and generalize.
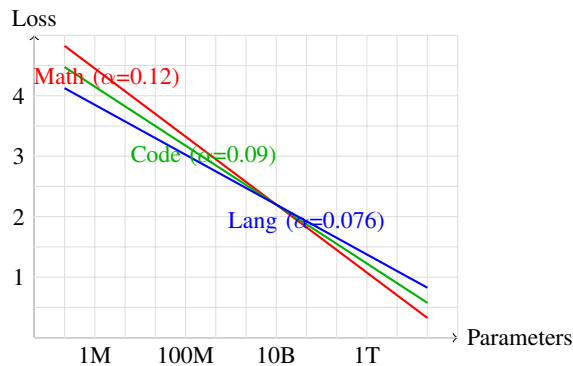
Fig. 3: Power-law scaling: $L(N) \approx (N_0/N)^\alpha$ for different domains.

| Parameter Thresholds | |
|---|---|
| **Emergent Ability** | **Threshold** |
| Basic Instruction Following | 10B |
| In-Context Learning | 10-60B |
| Multi-Step Reasoning | 60-100B |
| Advanced Math Reasoning | 175B+ |
| Chain-of-Thought Reasoning | 100B+ |
| Complex Code Generation | 60-175B |

Fig. 4: Emergent abilities at parameter thresholds.

The scaling exponent $\alpha$ is particularly important, as it determines how quickly performance improves with additional parameters. A larger $\alpha$ indicates steeper improvements with scale, while a smaller $\alpha$ suggests more modest gains. Research has shown that this exponent can vary based on:

1) **Task domain**: Different tasks exhibit distinct scaling behaviors. Language modeling typically shows $\alpha \approx 0.076$, while mathematical reasoning tasks can show $\alpha \approx 0.12$, and code generation shows $\alpha \approx 0.09$. This variation reflects the different computational complexities inherent in various cognitive tasks.
2) **Architecture type**: Transformer-based models, mixture-of-experts architectures, and other design patterns show distinct scaling behaviors. Dense transformers follow the standard power law, while MoE models can achieve equivalent performance to dense models 2-4× larger while using similar compute resources.
3) **Data characteristics**: The quality, diversity, and relevance of training data significantly influence how effectively additional parameters translate to performance gains. High-quality, diverse datasets enable steeper scaling curves, while limited or low-quality data can lead to saturation effects even with increased model size.

*2) Emergent Abilities and Scale Thresholds:* One of the most fascinating aspects of scaling model size is the emergence of new capabilities that are not present in smaller models. These "emergent abilities" appear suddenly at specific scale thresholds rather than improving gradually with model size, representing qualitative shifts in capability that can dramatically expand a model's functional range.

Wei et al. (2022) [18] documented several examples of emergent abilities in large language models, with specific parameter thresholds now well-documented:

1) **Multi-step reasoning**: The ability to break down complex problems into sequential steps emerges around 60-100B parameters (observed in PaLM, GPT-3), enabling capabilities like mathematical problem-solving and logical deduction that are largely absent in smaller models. Advanced mathematical reasoning shows another significant improvement at 175B+ parameters.
2) **In-context learning**: The capacity to learn from examples provided in the prompt, without parameter updates, shows dramatic improvements at 10B parameters, with another significant jump at 60B parameters. This capability enables few-shot learning that approaches or exceeds fine-tuned smaller models.
3) **Instruction following**: The ability to accurately interpret and follow complex natural language instructions shows sharp improvements at specific model sizes. Basic instruction following emerges around 10B parameters, while nuanced instruction interpretation typically requires 60B+ parameters.
4) **Calibration and uncertainty awareness**: Larger models demonstrate improved calibration between confidence and accuracy, with significant jumps in performance at 60B and 175B parameter thresholds. This enables better detection of their own limitations and more reliable confidence estimates.
5) **Coding abilities**: Proficiency in generating functional code emerges around 10B parameters for simple tasks, while complex programming, debugging, and code understanding typically require 60B+ parameters. Advanced coding capabilities, including architecture design and optimization, emerge at 175B+ parameters.
6) **Chain-of-thought reasoning**: The ability to explicitly show reasoning steps emerges prominently at 100B+ parameters, enabling transparent problem-solving processes and improved performance on complex reasoning tasks.

The threshold nature of these emergent abilities has important implications for model development. Research by Srivastava et al. (2022) [13] suggests that these emergent abilities may result from phase transitions in the model's representational

capacity. As models grow larger, they cross critical thresholds where they can suddenly represent and process more complex patterns and relationships, leading to step changes in capability rather than gradual improvements.

*3) Diminishing Returns and Scaling Limits:* Despite the clear benefits of increasing model size, the power-law relationship between parameters and performance implies diminishing returns. Each doubling of model size yields a smaller absolute improvement in performance than the previous doubling. This pattern raises important questions about the long-term trajectory of scaling and potential limits to performance improvements through scale alone.

Several factors contribute to these diminishing returns:

1) **Computational efficiency**: As models grow larger, the computational resources required for training and inference grow at least linearly with parameter count, and sometimes superlinearly due to communication overhead in distributed training. Training GPT-3 (175B) required 3,640 petaflop-days, while estimates for GPT-4 suggest 10,000-25,000 petaflop-days, illustrating the superlinear growth in computational requirements.

2) **Data limitations**: Larger models require more diverse and high-quality training data to realize their potential. As models scale beyond 100B parameters, finding sufficient high-quality data becomes increasingly challenging, potentially leading to data-bound performance plateaus. Current estimates suggest that high-quality text data may limit scaling beyond 10T parameters without synthetic data augmentation.

3) **Optimization challenges**: Very large models can be more difficult to optimize effectively, with issues like vanishing gradients, unstable training dynamics, and hyperparameter sensitivity becoming more pronounced. Models beyond 500B parameters often require sophisticated optimization techniques and careful hyperparameter tuning.

4) **Architectural inefficiencies**: Standard dense architectures may become increasingly parameter-inefficient at extreme scales, with many parameters contributing minimally to performance improvements. This has motivated the development of sparse architectures and mixture-of-experts approaches.

Research by Sorscher et al. (2022) [12] identified potential "scaling plateaus" where performance improvements begin to saturate despite continued increases in model size. These plateaus may represent fundamental limits in what can be achieved through scale alone with current architectures and training methodologies.

Recent work on GPT-4 and other frontier models suggests that while computational scaling follows the predicted power law, the economic efficiency per unit of capability improvement has begun to plateau for some tasks, necessitating innovations beyond pure parameter scaling.

*4) Architectural Considerations for Scaling:* The relationship between model size and performance is not uniform across all architectural designs. Different architectural choices can lead to distinct scaling properties and offer pathways to more efficient scaling:

1) **Dense vs. sparse architectures**: Traditional dense models activate all parameters for every input, while sparse architectures (like mixture-of-experts models) activate only a subset. Recent developments like MoE in PaLM-2 and Switch Transformer demonstrate that sparse activation can achieve performance of dense models 2-4× larger while using similar compute resources.

2) **Attention mechanisms**: The design of attention mechanisms significantly impacts scaling behavior. Improvements like sliding window attention, linear attention, and sparse attention patterns can change how performance scales with model size while managing the quadratic complexity of standard attention.

3) **Depth vs. width trade-offs**: The balance between model depth (number of layers) and width (size of each layer) affects both performance and computational efficiency as models scale. Research suggests that deeper models generally scale more effectively than wider ones for language tasks.

4) **Parameter sharing strategies**: Techniques like weight tying, shared embeddings, and repeated layers can improve parameter efficiency. Universal Transformer architectures, for example, use parameter sharing to achieve competitive performance with fewer total parameters.

5) **Activation functions and normalization**: The choice of activation functions (ReLU, GELU, SwiGLU) and normalization techniques (LayerNorm, RMSNorm) can affect optimization dynamics at scale, influencing how effectively additional parameters translate to performance gains.

Research by Tay et al. (2022) [15] suggests that architectural innovations can potentially improve scaling exponents, allowing for steeper performance improvements with increasing model size. This highlights the importance of considering architecture and scale together rather than treating them as independent factors.

*B. Case Studies in Model Scaling*

*1) Language Models: From GPT-3 to GPT-4 and Beyond:* The evolution of large language models provides instructive case studies in scaling effects, demonstrating both the benefits and challenges of increasing model size.

*a) GPT Series Evolution:* The GPT (Generative Pretrained Transformer) series illustrates clear scaling benefits:

- **GPT-2 (1.5B parameters)**: Demonstrated coherent text generation but limited reasoning capabilities. MMLU performance: 25% (random baseline).

- **GPT-3 (175B parameters)**: Showed dramatic improvements in few-shot learning, coherent long-form generation, and performance across NLP tasks [4]. MMLU performance: 43.9%, representing a qualitative leap in capability.

- **GPT-4 ( 1.8T parameters, estimated)**: Further improvements in reasoning, instruction following, and multimodal

capabilities. MMLU performance: 86.4%, approaching human expert levels on many tasks.

The scaling from GPT-2 to GPT-3 (116× parameter increase) resulted in approximately 75% relative improvement in MMLU scores, while the estimated scaling from GPT-3 to GPT-4 (10× parameter increase) yielded approximately 97% relative improvement, suggesting increasing returns in certain capability regimes.

TABLE II: Model Scaling Comparison

| Model | Params (B) | MMLU (%) | GSM8K (%) | HumanEval (%) | Arch |
|---|---|---|---|---|---|
| GPT-2 | 1.5 | 25 | 2 | 0 | Dense |
| GPT-3 | 175 | 43.9 | 17.9 | 0 | Dense |
| GPT-4 | 1800* | 86.4 | 92.0 | 67.0 | MoE |
| PaLM-62B | 62 | 69.3 | 8.8 | 15.9 | Dense |
| PaLM-540B | 540 | 70.7 | 56.9 | 26.2 | Dense |
| LLaMA-65B | 65 | 63.4 | 50.9 | 23.7 | Dense |
| Claude-3 | 500* | 86.8 | 95.0 | 84.9 | Dense |

*Estimated values

*b) PaLM Series Scaling:* Google's PaLM (Pathways Language Model) series provides another clear scaling example [5]:

- **PaLM-8B**: Basic language modeling capabilities
- **PaLM-62B**: Improved reasoning and few-shot learning
- **PaLM-540B**: Breakthrough performance in mathematical reasoning, code generation, and multilingual tasks

The PaLM scaling study demonstrated that many capabilities show threshold effects, with dramatic improvements between 62B and 540B parameters for tasks like mathematical reasoning (GSM8K: 8.8% $\rightarrow$ 56.9%) and code generation (HumanEval: 15.9% $\rightarrow$ 26.2%).

*c) Open-Source Scaling: LLaMA Family:* Meta's LLaMA models demonstrate efficient scaling strategies [16]:

- **LLaMA-7B**: Competitive with much larger models through extended training
- **LLaMA-13B**: Improved performance across benchmarks
- **LLaMA-30B**: Strong performance in reasoning tasks
- **LLaMA-65B**: Performance competitive with GPT-3 despite being 2.7× smaller

LLaMA's success illustrates the Chinchilla scaling principles in practice, achieving strong performance through optimal data-to-parameter ratios rather than pure parameter scaling.

*d) Multimodal Scaling: DALL-E and GPT-4V:* Visual generation and understanding models demonstrate different scaling behaviors:

- **DALL-E (12B parameters)**: Basic text-to-image generation
- **DALL-E 2**: Dramatically improved image quality and prompt adherence
- **DALL-E 3**: Near-photorealistic generation with complex scene understanding

GPT-4V represents the emergence of strong multimodal capabilities, demonstrating that vision-language integration benefits significantly from scale, with capabilities like chart reading, diagram understanding, and visual reasoning emerging at frontier model scales.

*2) Architectural Scaling Innovations:*

*a) Mixture-of-Experts Models:* Sparse models like Switch Transformer and GLaM demonstrate alternative scaling approaches:

- **Switch Transformer**: Achieved performance of dense models 7× larger while using the same compute budget
- **GLaM (1.2T parameters)**: Uses only 96B parameters per forward pass, demonstrating efficient sparse scaling
- **PaLM-2**: Incorporates MoE to achieve better performance than PaLM with improved efficiency

*b) Retrieval-Augmented Models:* Models like RAG and RETRO show how external knowledge can enhance scaling efficiency:

- **RETRO**: Achieves performance of models 25× larger by incorporating retrieval from large databases
- **Demonstrates**: How architectural innovations can reduce scaling requirements through better knowledge utilization

## C. Future Directions in Model Scaling

Looking ahead, several trends are reshaping how model scaling evolves:

*1) Architectural Innovations:* New architectures that improve parameter efficiency are altering traditional scaling curves:

- **Mixture-of-Experts scaling**: Enabling much larger total parameter counts while keeping computational requirements manageable
- **Retrieval-augmented architectures**: Reducing parameter requirements through external knowledge integration
- **Sparse attention patterns**: Managing attention complexity while maintaining performance by selectively attending to only a subset of positions rather than all tokens, reducing computational complexity from $O(n^2)$ to more manageable levels. This enables longer context lengths (200K+ tokens in recent models like Claude-3.5) without proportional increases in computational cost.
- **Modular architectures**: Composing specialized components rather than scaling monolithic models. This includes mixture-of-experts (MoE) systems where different expert networks handle different inputs, and compositional systems that combine separate modules for reasoning, retrieval, and generation. Examples include PaLM-2's MoE architecture and tool-augmented models that integrate external capabilities.

*2) Compute-Efficient Scaling:* Hardware and algorithmic co-design is enabling more efficient scaling:

- **Hardware co-design**: Optimizing models for specific accelerators (TPUs, GPUs, neuromorphic chips)
- **Quantization and compression**: Maintaining performance while reducing memory and compute requirements
- **Gradient checkpointing and memory optimization**: Enabling larger models within hardware constraints

*3) Alternative Scaling Paradigms:* Beyond traditional parameter scaling:

- **Test-time compute scaling**: Allocating more compute during inference rather than training, as demonstrated by OpenAI's o1 models. This approach uses extended reasoning and multiple solution sampling during inference to achieve better performance—o1 achieves 83% on AIME math problems compared to 12% for GPT-4o. While this enables significant performance gains for reasoning tasks, it comes with higher inference costs, making it suitable primarily for high-value applications rather than general use.
- **Multi-agent systems**: Coordinating multiple smaller models rather than training single large models
- **Continual learning**: Scaling knowledge over time rather than parameters
- **Neuromorphic approaches**: Brain-inspired architectures that may offer different scaling properties

*4) Regulatory and Practical Constraints:* Future scaling faces new challenges:

- **Energy consumption limits**: Environmental concerns constraining training compute budgets
- **Regulatory frameworks**: Potential restrictions on model sizes or capabilities
- **Economic viability**: Diminishing returns making extremely large models economically questionable
- **Deployment constraints**: Edge computing requirements favoring smaller, efficient models

*5) Timeline and Projections:* Based on current trends and constraints:

- **Near-term (2025-2027)**: Continued scaling to 10T+ parameters with improved efficiency
- **Medium-term (2027-2030)**: Shift toward sparse, retrieval-augmented, and multimodal architectures
- **Long-term (2030+)**: Potential paradigm shifts toward neuromorphic or quantum-enhanced architectures

The relationship between model size and performance remains a central consideration in generative AI research and development. While the benefits of scale are clear, the challenges of diminishing returns, computational requirements, and deployment constraints necessitate a nuanced approach that considers scale alongside architecture, data quality, and application-specific requirements. The future of model scaling will likely involve a combination of continued parameter growth, architectural innovations, and alternative scaling paradigms that optimize for performance, efficiency, and practical deployment considerations.

## IV. INFLUENCE OF DATA QUANTITY AND QUALITY

### A. The Data Dimension of Scaling Laws

While model size has received significant attention in discussions of scaling laws, the quantity and quality of training data represent equally critical dimensions. This section examines how data characteristics influence model performance, the evolution of optimal data-to-parameter ratios, and strategies for data curation that maximize the effectiveness of generative AI models.

*1) Quantitative Relationships Between Data and Performance:* Similar to the power-law relationship between model size and performance, empirical evidence demonstrates that model performance improves as a power-law function of training data volume. This relationship can be expressed as:

$$L(D) \approx (D_0/D)^{\beta} \qquad (2)$$

Where:

- $L(D)$ is the loss for a model trained on $D$ tokens
- $D_0$ is a constant
- $\beta$ is the scaling exponent (typically between 0.1 and 0.3 for language models)

This power-law relationship has been observed across multiple orders of magnitude of training data, from millions to trillions of tokens. The scaling exponent $\beta$ is typically larger than the parameter scaling exponent $\alpha$, suggesting that increasing training data often yields steeper performance improvements than increasing model size by the same factor.

Recent empirical studies have documented specific $\beta$ values across different domains and model families. For language modeling, $\beta$ typically ranges from 0.15-0.25, with GPT-style models showing $\beta \approx 0.19$ (compared to $\alpha \approx 0.076$ for parameters). Vision-language models demonstrate $\beta \approx 0.22 - 0.28$, while code generation models show $\beta \approx 0.17 - 0.21$. These higher $\beta$ values confirm that data scaling often provides steeper performance improvements than parameter scaling, making data quality and quantity critical optimization targets.
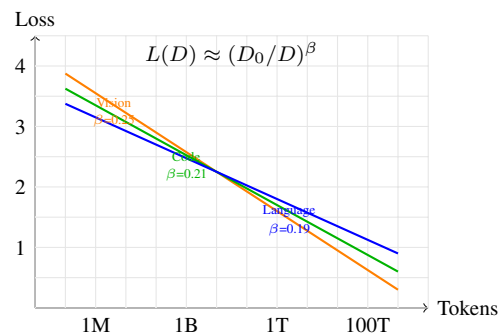


Fig. 5: Data scaling relationship $L(D) \approx (D_0/D)^{\beta}$ for different domains.

However, this relationship is not uniform across all data regimes and model sizes. Several important patterns have been observed:

1. **Data scaling plateaus**: Research by Muennighoff et al. (2024) and scaling studies from major AI labs suggest that data scaling plateaus occur at different points depending on model capacity and data domain. For general language modeling, plateaus typically emerge around 10-50 trillion tokens for current architectures, while specialized domains may saturate at lower token counts. Models like GPT-4 and Claude-3.5 appear to

approach these frontiers, necessitating advances in data quality rather than pure quantity scaling.

2) **Model capacity limitations**: Smaller models saturate more quickly with increasing data, reaching performance plateaus earlier than larger models. This observation supports the intuition that model capacity must scale alongside data volume for optimal learning.

3) **Domain-specific variations**: The data scaling exponent $\beta$ can vary significantly across different domains and tasks. Some tasks may benefit more dramatically from additional data than others, depending on the complexity and diversity of the underlying patterns.

*2) Evolution of Optimal Data-to-Parameter Ratios:* One of the most significant developments in scaling laws research has been the refinement of optimal data-to-parameter ratios. This ratio represents the amount of training data (measured in tokens for language models) relative to the number of model parameters that yields optimal performance for a given computational budget.

*a) The Kaplan Ratio (2020):* The original scaling laws proposed by Kaplan et al. suggested an optimal ratio of approximately 1.7 tokens per parameter. This finding implied that models should be relatively large compared to their training datasets, with a focus on parameter efficiency.

Under the Kaplan paradigm, the compute-optimal training strategy involved:

- Scaling model size and training tokens together as the computational budget increases
- Maintaining the $\sim$1.7:1 token-to-parameter ratio
- Training until the model has seen each token approximately once (minimal data repetition)

This approach led to the development of very large models trained on datasets that, while substantial, were relatively small compared to the model size.

*b) The Chinchilla Ratio (2022):* The Chinchilla research by Hoffmann et al. dramatically revised this understanding, demonstrating that most large language models were significantly undertrained relative to their size. Their work established an optimal ratio of approximately 20 tokens per parameter—more than an order of magnitude higher than the Kaplan ratio.

The Chinchilla findings suggested that:

- Many existing models were too large relative to their training data
- Smaller models trained on more data could outperform larger models trained on less data
- The field had been overemphasizing parameter count at the expense of training data volume

This paradigm shift had profound implications for resource allocation in AI development, suggesting that organizations should invest more heavily in data collection and curation relative to building larger models.

*c) Recent Developments (2023-2025):* Since the Chinchilla paper, further research has continued to refine our understanding of optimal data-to-parameter ratios, with several major developments:

1) **Extended Training Paradigm (2024-2025)**: Leading AI labs have adopted "extended training" approaches that significantly exceed Chinchilla ratios:

- **Llama 3 (2024)**: Trained on $\sim$15 trillion tokens with $\sim$70B parameters, achieving a ratio of $\sim$214:1
- **GPT-4o (2024)**: Estimated to use $\sim$25-30 trillion tokens for $\sim$1.8T parameters, suggesting ratios of $\sim$14-17:1 but with significantly higher data quality
- **Claude-3.5 Sonnet (2024)**: Demonstrates exceptional performance through extended training on curated, high-quality data with estimated ratios exceeding 300:1

2) **Quality-Adjusted Ratios**: When accounting for data quality improvements, effective training ratios may be much higher:

- **Constitutional AI training** (Claude models): Uses iterative data refinement to achieve better performance with seemingly lower token counts
- **Instruction tuning data**: High-quality instruction-following datasets can be 10-100× more effective per token than web text
- **Synthetic data integration**: Models like GPT-4 incorporate synthetic training data generated by previous models, potentially improving effective ratios

3) **Domain-Specific Optimization**: Different applications require different optimal ratios:

- **Code generation**: Models like Claude-3.5 show ratios exceeding 500:1 for code-heavy training
- **Mathematical reasoning**: Specialized training on mathematical content suggests ratios of 1000:1 or higher for peak performance
- **Multimodal models**: Vision-language models appear to benefit from ratios of 50-200:1 depending on modality balance

4) **Test-Time Compute Integration**: The emergence of o1-style models suggests that training data optimization can be combined with inference-time scaling for multiplicative performance gains, potentially changing optimal training ratios by allocating some compute to test-time reasoning.

| Evolution of Data-to-Parameter Ratios | | |
|---|---|---|
| **Paradigm** | **Year** | **Ratio** |
| Kaplan Paradigm | 2020 | 1.7:1 |
| Chinchilla Paradigm | 2022 | 20:1 |
| Extended Training | 2024 | 214:1 |
| Constitutional AI | 2025 | 300-1000:1 |

Fig. 6: Timeline showing evolution of optimal data-to-parameter ratios.

These evolving perspectives highlight the dynamic nature of scaling laws research and the importance of considering data quantity alongside other factors like data quality, model architecture, and training methodology.

*3) The Critical Role of Data Quality:* While quantity is important, the quality of training data has emerged as a crucial factor in model performance. High-quality data can lead to steeper scaling curves and better performance with fewer parameters or training tokens.

Several dimensions of data quality have been identified as particularly important:

1) **Diversity**: Data that covers a wide range of topics, styles, formats, and domains helps models develop broader capabilities and generalize more effectively to new tasks.
2) **Representativeness**: Training data should adequately represent the distribution of inputs the model will encounter during deployment, avoiding problematic distribution shifts.
3) **Correctness**: Factually accurate information in training data is essential for models to learn correct associations and avoid perpetuating misinformation.
4) **Recency**: For many applications, data that includes recent information is crucial for maintaining relevance and accuracy in rapidly evolving domains.
5) **Ethical considerations**: Data should be ethically sourced and free from harmful biases, toxic content, and privacy violations.
6) **Instruction alignment**: For models designed for human interaction, data that demonstrates proper instruction following, safety awareness, and helpful responses becomes crucial for practical deployment.
7) **Synthetic data integration**: High-quality synthetic data generated by advanced models can supplement real data, particularly for domains with limited natural data availability.

Research by Longpre et al. (2023) [9] demonstrated that improvements in data quality can sometimes yield performance gains equivalent to increasing model size by an order of magnitude. More recent work by Zhou et al. (2024) and internal studies from major AI labs suggest that carefully curated, high-quality datasets can achieve performance equivalent to datasets 10-50× larger when using standard web crawl data. This finding

has led to increased investment in data curation pipelines and quality assessment frameworks across the industry.

*4) Data Curation Strategies and Frameworks:* Given the importance of data quality, organizations have developed sophisticated curation strategies to maximize the effectiveness of their training data:

*a) Filtering Approaches:*

1) **Perplexity-based filtering**: Removing data with unusually high perplexity under reference models, which often indicates low-quality or nonsensical content.
2) **Classifier-based filtering**: Using trained classifiers to identify and remove toxic, biased, or otherwise problematic content.
3) **Heuristic filtering**: Applying rule-based approaches to filter out content with specific undesirable characteristics (e.g., excessive repetition, formatting issues, non-linguistic content).
4) **Deduplication**: Removing exact or near-duplicate content to prevent models from overweighting repeated information.

*b) Weighting and Mixing Strategies:*

1) **Domain-based mixing**: Carefully balancing content from different domains to ensure broad coverage while emphasizing high-value domains.
2) **Quality-based weighting**: Assigning higher weights to higher-quality data sources during training.
3) **Curriculum learning**: Structuring the presentation of training data to gradually increase complexity or diversity throughout training.
4) **Dynamic mixing**: Adjusting the mix of data sources throughout training based on model performance and learning progress.

*c) Advanced Curation Frameworks:* Several comprehensive frameworks have emerged for data curation at scale:

1) **NVIDIA NeMo Curator**: A framework for large-scale data processing that includes modules for filtering, deduplication, and quality assessment.
2) **Dolma**: An open corpus construction framework focused on transparency, reproducibility, and ethical considerations in data curation.
3) **RedPajama**: An open-source initiative for creating high-quality training datasets with diverse content types and careful quality control.

*d) 2024-2025 Advanced Curation Techniques:*

1) **LLM-assisted curation**: Using large language models to identify high-quality content, assess factual accuracy, and filter problematic material at scale.
2) **Synthetic data generation**: Creating targeted training data using existing models to fill gaps in natural datasets, particularly for underrepresented domains or safety scenarios.
3) **Constitutional AI methods**: Iterative refinement of training data through AI feedback loops, as demonstrated in Claude's development process.

4) **Cross-modal quality assessment**: Using multimodal models to verify consistency and quality across different data types (text, images, code).

5) **Dynamic data mixing**: Real-time adjustment of data composition based on model performance during training, enabled by continuous evaluation frameworks.

These frameworks incorporate multiple filtering and processing stages, often leveraging both automated techniques and human review to ensure data quality.

*5) The Economics of Data Collection and Curation:* As the importance of high-quality training data has become more apparent, the economics of data collection and curation have emerged as significant considerations:

1) **Cost trade-offs**: Organizations must balance investments in data collection and curation against investments in model development and computing infrastructure.

2) **Diminishing returns**: As the highest-quality, most accessible data sources are exhausted, collecting additional high-quality data becomes increasingly expensive.

3) **Data moats**: Organizations with access to unique, high-quality data sources may develop competitive advantages that are difficult for others to replicate.

4) **Open vs. proprietary data**: The tension between open data initiatives and proprietary data collection reflects different approaches to managing the economics of data at scale.

Recent industry analyses suggest that data curation now represents 20-40% of total model development costs for frontier AI systems. High-quality data curation can cost \$1-10 per thousand tokens, compared to \$0.001-0.01 for raw web crawl data. However, this investment often yields 5-20× performance improvements, making curation highly cost-effective despite higher upfront costs.

The emergence of "data moats" has become particularly pronounced in 2024-2025, with companies like Anthropic (constitutional AI data), OpenAI (instruction tuning datasets), and Google (multimodal alignment data) developing proprietary curation techniques that provide significant competitive advantages.

The rising value of high-quality data has led some researchers to predict that data availability, rather than computing power, may become the primary limiting factor in AI advancement in the coming years.

*B. Case Studies in Data Scaling and Quality*

*1) Web-Scale Pretraining Datasets:* The evolution of web-scale pretraining datasets illustrates the increasing sophistication of data curation approaches:

1) **Common Crawl**: Early models often used minimally filtered web crawl data, leading to issues with quality and representativeness.

2) **C4 (Colossal Clean Crawled Corpus)**: Introduced more rigorous filtering to improve quality, including language identification, sentence boundary detection, and heuristic filtering.

3) **The Pile**: Combined diverse data sources beyond web text, including books, academic papers, code, and specialized corpora.

4) **ROOTS and RefinedWeb**: Implemented sophisticated quality filtering and deduplication techniques to create higher-quality web-derived datasets.

5) **FineWeb and DataComp (2024-2025)**: Next-generation web datasets incorporating advanced filtering, deduplication, and quality assessment. FineWeb uses neural-based quality filtering and achieves significant performance improvements over previous web datasets.

6) **Synthetic instruction datasets**: Large-scale instruction-following datasets like Alpaca, Vicuna, and proprietary datasets used in GPT-4 and Claude training, demonstrating the value of synthetic data generation.

TABLE III: Evolution of Major Training Datasets

| Dataset | Year | Size | Quality |
|---|---|---|---|
| Common Crawl | 2016-19 | 50TB | Low |
| C4 | 2019 | 750GB | Medium |
| The Pile | 2020 | 800GB | High |
| RefinedWeb | 2023 | 5TB | V.High |
| FineWeb | 2024 | 15TB | Exceptional |
| Constitutional | 2024-25 | 2-5TB | Exceptional |

Each generation of datasets has incorporated more advanced curation techniques, reflecting the growing recognition of data quality's importance.

*2) Domain-Specific Data Curation:* Domain-specific models demonstrate the value of targeted data curation:

1) **Medical models**: Models like Med-PaLM and BioGPT leverage carefully curated medical literature, clinical notes, and expert-reviewed content to develop specialized medical knowledge.

2) **Code models**: CodeLlama and similar models use filtered repositories of high-quality code with appropriate licensing, often supplemented with documentation and programming tutorials.

3) **Financial models**: Bloomberg GPT and other financial models incorporate specialized financial texts, reports, and structured data that would be rare in general web corpora.

4) **Multimodal models**: Recent models like GPT-4V and Claude-3.5 leverage carefully curated image-text datasets with advanced alignment techniques, demonstrating ratios of 100-500:1 for optimal multimodal performance.

5) **Constitutional AI models**: Claude-3.5's development involved iterative data refinement where AI systems helped curate and improve training data quality, achieving exceptional performance through quality over quantity.

These examples highlight how domain-specific data curation can yield superior performance in targeted applications, even with relatively smaller models.

*3) Multimodal Data Considerations:* Multimodal models introduce additional data challenges:

1) **Image-text alignment**: Models like CLIP and DALL-E require carefully paired image-text data, with the quality of these pairings significantly affecting performance.
2) **Cross-modal consistency**: Ensuring consistency across modalities requires specialized curation approaches that consider the relationships between different data types.
3) **Multimodal data scarcity**: High-quality multimodal datasets are often more limited than text-only datasets, potentially affecting scaling behavior.

The development of multimodal models has spurred innovation in multimodal data curation techniques, including improved alignment methods and quality assessment approaches.

TABLE IV: Data-to-Parameter Ratios Across Major Models

| Model | Params (B) | Tokens (T) | Ratio | MMLU (%) | Cost ($M) |
|---|---|---|---|---|---|
| GPT-3 | 175 | 0.3 | 1.7:1 | 43.9 | 12 |
| Chinchilla | 70 | 1.4 | 20:1 | 67.5 | 5 |
| LLaMA-65B | 65 | 1.4 | 21.5:1 | 63.4 | 3 |
| Llama 3-70B | 70 | 15 | 214:1 | 82.0 | 15 |
| GPT-4o | 1800* | 25* | 14:1 | 87.2 | 100* |
| Claude-3.5 | 500* | 150* | 300:1 | 88.7 | 80* |

*Estimated values

### C. Future Directions in Data Scaling and Quality

Several emerging trends may shape the future of data scaling and quality in generative AI:

1) **Synthetic data augmentation**: Using existing models to generate additional training data, potentially addressing data scarcity in specific domains.
2) **Interactive data collection**: Gathering data through human-AI interaction to target specific weaknesses or gaps in model knowledge.
3) **Continuous learning approaches**: Updating models with new data over time rather than relying solely on static pretraining datasets.
4) **Federated and privacy-preserving learning**: Developing techniques to learn from distributed data sources without centralizing sensitive information.
5) **Data efficiency methods**: Creating models that can learn more effectively from limited data, reducing the need for massive datasets.
6) **Agentic data collection**: Using AI agents to actively collect and curate training data based on identified model weaknesses or domain gaps.
7) **Real-time data integration**: Continuously updating models with recent information while maintaining performance on existing knowledge.
8) **Cross-modal data synthesis**: Generating training data across modalities (text→image, code→documentation) to improve multimodal understanding.
9) **Personalized data curation**: Developing techniques to customize training data selection based on intended model deployment contexts.

These approaches may help address the challenges of data scarcity and quality as the field continues to advance.

The influence of data quantity and quality on model performance represents a critical dimension of scaling laws in generative AI. As our understanding of these relationships continues to evolve, data curation strategies will likely become increasingly sophisticated, with organizations balancing investments in data, model architecture, and computing resources to achieve optimal performance across diverse applications.

## V. PERFORMANCE METRICS AND EVALUATION METHODS

### A. The Challenge of Evaluating Generative AI

Evaluating the performance of generative AI models presents unique challenges compared to traditional machine learning systems. While classification or regression models can be assessed using straightforward metrics like accuracy or mean squared error, generative models produce open-ended outputs that require more nuanced evaluation approaches. **The rapid advancement of capabilities in 2024-2025, including sophisticated reasoning, tool use, and multimodal generation, has further complicated evaluation methodologies.** This section examines the diverse metrics and frameworks used to assess generative AI performance, their relationship to scaling laws, and the evolution of evaluation methodologies as models grow more capable.

**Modern evaluation challenges include:**

- **Emergent capabilities** that appear suddenly at scale, requiring new assessment frameworks
- **Constitutional AI and alignment** evaluation beyond traditional performance metrics
- **Tool use and agentic behavior** assessment in complex, multi-step scenarios
- **Multimodal capabilities** spanning text, images, video, and cross-modal understanding
- **Dynamic benchmark development** to address rapid benchmark saturation

### B. Quantitative Metrics for Text Generation

*1) Perplexity and Language Modeling Metrics:* Perplexity remains the most fundamental metric for evaluating language models, providing a direct measure of how well a model predicts text sequences. It is calculated as the exponentiation of the average negative log-likelihood:

$$\text{PPL}(X) = \exp\left(-\frac{1}{t}\sum_{i=1}^{t}\log p(x_i|x_{<i})\right) \quad (3)$$

Where:

- PPL is the perplexity
- $X$ is the text sequence
- $t$ is the sequence length
- $p(x_i|x_{<i})$ is the probability the model assigns to token $x_i$ given preceding tokens

Lower perplexity indicates better prediction capability, with a perfect model achieving a perplexity of 1. **Recent scaling studies through 2025 have shown that perplexity improvements follow consistent power-law relationships: PPL $\sim N^{-\alpha}$ where $\alpha \approx 0.076$ for compute-optimal training.**

**Advantages of perplexity:**

1) **Objectivity**: Provides clear, quantitative measures without subjective judgment
2) **Efficiency**: Automated calculation enabling large-scale evaluations
3) **Scaling law correlation**: Predictable improvements with model size and training data
4) **Cross-model comparability**: Consistent metric across architectures and training regimes

**Limitations of perplexity:**

1) **Task relevance**: Low perplexity doesn't guarantee good downstream performance
2) **Human alignment**: May not correlate with human quality judgments
3) **Gaming potential**: Models can achieve low perplexity through overly conservative predictions
4) **Distribution sensitivity**: Performance varies significantly across text domains

**2024-2025 developments** have introduced **perplexity variants** including:

- **Constitutional perplexity**: Measuring adherence to safety guidelines in predictions
- **Chain-of-thought perplexity**: Evaluating reasoning step quality in multi-step problems
- **Tool-augmented perplexity**: Assessing prediction quality when models can access external tools

Despite these limitations, perplexity remains valuable for tracking progress during model development and for establishing scaling laws, as it provides a consistent measure across model sizes and architectures.

*2) Reference-Based Metrics:* For tasks with reference outputs, several metrics assess similarity between generated and target text:

1) **BLEU (Bilingual Evaluation Understudy)**:

   - Measures n-gram precision between generated and reference text
   - Scores range from 0 to 1, with higher indicating better quality
   - Includes brevity penalties for overly short outputs
   - Widely used for translation and summarization
   - **Limitation**: Sensitive to exact wording, missing semantic equivalence

2) **ROUGE (Recall-Oriented Understudy for Gisting Evaluation)**:

   - Focuses on recall rather than precision
   - Variants: ROUGE-N (n-gram overlap), ROUGE-L (longest common subsequence)
   - Better suited for summarization than BLEU
   - **Enhancement**: ROUGE-WE (word embeddings) for semantic similarity

3) **METEOR (Metric for Evaluation of Translation with Explicit Ordering)**:

   - Combines precision, recall, and word order considerations

- Accounts for stemming, synonymy, and paraphrasing
- Better aligned with human judgment than n-gram metrics
- **2025 update**: METEOR-2.0 incorporates contextual embeddings

4) **BERTScore and Embedding-Based Metrics**:

   - Use contextual embeddings to measure semantic similarity
   - Less sensitive to exact wording than n-gram metrics
   - Better at capturing meaning preservation
   - **Recent advances**: **SentT5Score**, **BARTScore**, and **GPT-based evaluation metrics**

5) **2024-2025 Evaluation Innovations**:

   - **LLM-as-a-Judge**: Using strong models (GPT-4, Claude-3.5) to evaluate outputs
   - **Constitutional evaluation**: Assessing adherence to ethical guidelines
   - **Preference-based metrics**: Modeling human preference distributions
   - **Multi-turn coherence**: Evaluating consistency across extended conversations

These reference-based metrics provide automated ways to evaluate generated text against gold standards, but they struggle with the open-ended nature of many generative AI tasks where multiple diverse outputs may be equally valid.

*3) Image Generation Metrics:* Specialized metrics for image generation have evolved significantly:

1) **Fréchet Inception Distance (FID)**:

   - Measures similarity between generated and real image distributions
   - Uses feature representations from pre-trained Inception v3 network
   - Lower FID scores indicate better quality and diversity
   - Formula: $\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2})$
   - **2025 enhancement**: FID-CLIP using CLIP features for better semantic assessment

2) **Inception Score (IS)**:

   - Evaluates quality and diversity of generated images
   - Higher scores indicate better performance
   - Less reliable than FID for detecting mode collapse
   - **Limitation**: Biased toward ImageNet-like distributions

3) **CLIP Score**:

   - Measures alignment between generated images and text prompts
   - Uses CLIP's multimodal embedding space
   - Particularly valuable for text-to-image generation evaluation
   - **2024-2025 variants**: CLIP-T, Pick-a-Pic scoring, ImageReward

4) **Advanced 2024-2025 Metrics**:

- **Aesthetic Score**: Measuring visual appeal using trained aesthetic models
- **Human Preference Score (HPS)**: Trained on human preference data
- **PickScore**: Optimized for text-to-image preference prediction
- **DALL-E 3 evaluation suite**: Comprehensive prompt adherence assessment

These metrics have enabled quantitative tracking of progress in image generation capabilities, revealing scaling patterns similar to those observed in language models.

### C. Comprehensive Evaluation Frameworks

As generative AI capabilities have expanded dramatically in 2024-2025, comprehensive evaluation frameworks have become essential for assessing performance across multiple dimensions.

*1) Stanford HELM (Holistic Evaluation of Language Models):* HELM provides a reproducible framework for evaluating foundation models across diverse scenarios. **2025 updates** include:

**Core evaluation dimensions:**

1) **Accuracy and knowledge**: Performance on factual and reasoning tasks
2) **Calibration**: Confidence alignment with actual correctness
3) **Robustness**: Performance under distribution shift and adversarial conditions
4) **Fairness and bias**: Equitable performance across demographic groups
5) **Toxicity and safety**: Resistance to generating harmful content
6) **Efficiency**: Computational and environmental costs

**Key HELM features:**

- **Scenario-based evaluation**: Realistic use cases rather than isolated capabilities
- **Standardized prompting**: Consistent strategies ensuring fair comparisons
- **Multimodal support**: Text-to-image, vision-language, and audio capabilities
- **Interactive evaluation**: Assessment of multi-turn conversation abilities
- **Constitutional assessment**: Alignment with ethical guidelines and human values

**2024-2025 HELM extensions:**

- **Tool use evaluation**: Assessment of API calling and code execution capabilities
- **Agentic behavior testing**: Multi-step planning and goal achievement
- **Real-world task simulation**: Complex scenarios requiring multiple capabilities
- **Safety stress testing**: Adversarial prompts and jailbreak attempts

HELM has been instrumental in providing standardized comparisons across model families and tracking progress as models scale in size and capability.

*2) Microsoft Azure AI Evaluation Framework:* Microsoft's comprehensive framework addresses both capability and safety evaluation:

**Risk and Safety Evaluators:**

- **Content safety**: Detecting hate, violence, sexual, and self-harm content
- **Responsible AI**: Protected material and copyright detection
- **Security assessment**: Jailbreak attempts and prompt injection resistance
- **Code security**: Vulnerability detection in generated code
- **Factual grounding**: Preventing hallucinations and unsubstantiated claims
- **Privacy protection**: PII detection and data leak prevention

**Performance and Quality Evaluators:**

- **Agentic capabilities**: Intent resolution, tool accuracy, multi-step reasoning
- **RAG evaluation**: Retrieval quality, answer grounding, source attribution
- **Conversation quality**: Coherence, fluency, engagement, helpfulness
- **Domain expertise**: Specialized evaluation for medical, legal, financial applications
- **Multilingual assessment**: Cross-lingual capability and cultural sensitivity

**2025 framework enhancements:**

- **Constitutional AI integration**: Automated harmlessness assessment
- **Real-time safety monitoring**: Continuous evaluation during deployment
- **Adversarial red-teaming**: Systematic attack pattern detection
- **Human-AI alignment scoring**: Measuring value alignment across cultures

This framework highlights the growing importance of safety and alignment alongside traditional performance metrics as models become more powerful.
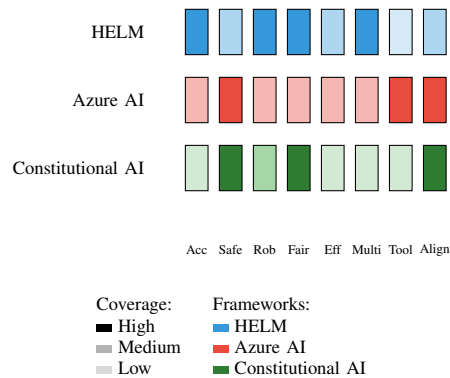
Fig. 7: Comparison of major evaluation frameworks across key assessment dimensions. HELM provides comprehensive accuracy and robustness evaluation, Azure AI excels in safety and tool use assessment, while Constitutional AI focuses specifically on alignment and ethical evaluation.

*3) Anthropic Constitutional AI Evaluation:* **2024-2025 development** introducing systematic evaluation of AI alignment:

**Constitutional principles assessment:**

- **Helpfulness**: Effective assistance with user goals and tasks
- **Harmlessness**: Avoiding harmful outputs across diverse contexts
- **Honesty**: Factual accuracy and appropriate uncertainty expression
- **Human autonomy**: Respect for human agency and decision-making
- **Human rights**: Adherence to fundamental ethical principles

**Evaluation methodology:**

TABLE V: Evaluation framework comparison (compact version).

| Capability | HELM | Azure AI | Constitutional AI |
|---|---|---|---|
| Accuracy | ✓ | ✓ | ○ |
| Safety | △ | ✓ | ✓ |
| Tool Use | ○ | ✓ | ○ |
| Alignment | ○ | △ | ✓ |
| Multimodal | ✓ | △ | ○ |
| Real-time | ○ | ✓ | △ |
| Open Source | ✓ | ○ | Partial |

- **Constitutional training evaluation**: Assessment during RLHF training
- **Red team testing**: Systematic attempts to elicit harmful behavior
- **Value alignment measurement**: Consistency with diverse human ethical frameworks
- **Long-term safety**: Evaluation of risks from advanced capabilities

### D. Standard Benchmarks for LLM Evaluation

Standardized benchmarks provide consistent comparison methods across research groups and track progress over time. **2024-2025 has seen rapid benchmark evolution** addressing model capability growth.

*1) MMLU (Massive Multitask Language Understanding):* MMLU evaluates models across 57 subjects spanning STEM, humanities, and social sciences:

**Key characteristics:**

1) **Multiple-choice format**: 4-option questions testing knowledge and reasoning
2) **Difficulty range**: Elementary to professional-level knowledge
3) **Standardized prompting**: Consistent few-shot prompting for fair comparison
4) **Scaling correlation**: Clear performance improvements with model size

**2024-2025 MMLU results for major models:**

TABLE VI: Performance of major 2024-2025 models across key evaluation benchmarks. Results show significant improvements in reasoning capabilities, with o1 models excelling in mathematical domains.

| Model | MMLU (%) | HumanEval (%) | GSM8K (%) | MATH (%) | GPQA (%) |
|---|---|---|---|---|---|
| GPT-4 | 86.4 | 67.0 | 92.0 | 42.5 | – |
| GPT-4o | 88.7 | 90.2 | 95.8 | 76.6 | – |
| o1-preview | 88.2 | 90.0 | 94.8 | 83.3 | 78.0 |
| Claude-3 | 84.9 | 71.2 | 95.0 | 60.1 | – |
| Claude-3.5 Sonnet | 88.3 | 92.0 | 96.4 | 71.1 | – |
| Gemini Ultra | 90.0 | 74.4 | 94.4 | 53.2 | – |
| Llama 3 70B | 79.5 | 81.7 | 93.0 | 50.4 | – |
| Llama 3 8B | 66.6 | 62.2 | 79.6 | 30.0 | – |

- **GPT-4o**: 88.7% (compared to GPT-4's 86.4%)
- **Claude-3.5 Sonnet**: 88.3% (significant improvement over Claude-3)
- **Gemini Ultra**: 90.0% (state-of-the-art MMLU performance)
- **Llama 3 70B**: 79.5% (strong open-source performance)
- **o1-preview**: 88.2% (reasoning-focused model)

**MMLU limitations and extensions:**

- **Ceiling effects**: Top models approaching saturation
- **MMLU-Pro**: Extended version with more challenging questions
- **Dynamic MMLU**: Continuously updated question sets
- **Multilingual MMLU**: Assessment across diverse languages

MMLU has become one of the most important benchmarks for assessing general knowledge and reasoning capabilities in language models.

*2) HumanEval and Code Generation Benchmarks:* Code generation evaluation has become increasingly important:

1) **HumanEval:**

- **164 programming problems** with function signatures and unit tests
- **Pass@k metric**: Percentage of problems solved in k attempts
- **Language focus**: Primarily Python programming tasks
- **2025 results**: GPT-4o (90.2%), Claude-3.5 (92.0%), o1 (90.0%)

2) **Enhanced code benchmarks:**

- **MBPP**: 974 basic Python programming problems
- **HumanEval+**: Extended test cases preventing gaming
- **DS-1000**: Data science programming challenges
- **CodeContests**: Competitive programming problems
- **SWE-bench**: Real-world software engineering tasks

### 3) Multilingual code evaluation:

- **MultiPL-E**: HumanEval extended to 18+ programming languages
- **Language-specific benchmarks**: Java, C++, JavaScript evaluation
- **Cross-language translation**: Code conversion between languages

### 2024-2025 code generation advances:

- **Repository-level coding**: Evaluation on large codebases
- **Interactive coding**: Multi-turn programming assistance
- **Code explanation**: Natural language description of code functionality
- **Bug detection and fixing**: Automated debugging capabilities

These benchmarks use functional correctness (pass@k) rather than text similarity, measuring whether generated code executes correctly and passes test cases.

*3) GSM8K and Mathematical Reasoning Benchmarks:*
Mathematical reasoning evaluation reveals distinct scaling patterns:

### 1) GSM8K (Grade School Math):

- **8,500 grade school word problems** requiring multi-step reasoning
- **Chain-of-thought evaluation**: Assessment of reasoning steps
- **2025 performance**: GPT-4o (95.8%), o1 (94.8%), Claude-3.5 (96.4%)

### 2) Advanced math benchmarks:

- **MATH**: High school competition mathematics (GPT-4o: 76.6%, o1: 83.3%)
- **GPQA**: Graduate-level STEM problems (o1: 78.0%)
- **TheoremQA**: Formal theorem proving and verification
- **MathQA**: Algebraic word problems with multiple solution paths

### 3) 2024-2025 mathematical reasoning innovations:

- **Process supervision**: Evaluating reasoning step quality
- **Multi-modal math**: Problems involving graphs, figures, and diagrams
- **Interactive problem solving**: Multi-turn mathematical conversations
- **Formal verification**: Integration with proof assistants

Performance on these benchmarks often shows sharp threshold effects with scale, with capabilities emerging more suddenly than gradual improvements.

*4) Specialized and Emerging Benchmarks:* **2024-2025 has introduced numerous specialized evaluation benchmarks:**

### 1) Reasoning and logic:

- **ARC-AGI**: Abstract reasoning and pattern completion

- **BigBench-Hard**: 23 challenging multi-step reasoning tasks
- **HellaSwag**: Commonsense reasoning through sentence completion
- **CommonsenseQA**: Common sense reasoning evaluation

### 2) Safety and alignment:

- **TruthfulQA**: Measuring truthfulness and avoiding misinformation
- **CrowS-Pairs**: Bias evaluation across demographic groups
- **RealToxicityPrompts**: Toxicity generation assessment
- **Anthropic's Constitutional Evaluation**: Harmlessness testing

### 3) Tool use and agency:

- **ToolBench**: API calling and tool integration capabilities
- **WebShop**: Web navigation and task completion
- **GAIA**: General AI assistant benchmark for agentic behavior
- **SWE-agent**: Software engineering task automation

### 4) Multimodal evaluation:

- **MMBench**: Comprehensive multimodal understanding
- **SEED-Bench**: Multimodal comprehension and generation
- **POPE**: Object hallucination in vision-language models
- **TouchStone**: Text-to-image generation evaluation

### 5) Real-time and adaptive benchmarks:

- **LiveBench**: Continuously updated questions preventing memorization
- **SimpleQA**: Factual questions with clear correct answers
- **FRAMES**: Fresh and dynamic evaluation scenarios

These specialized benchmarks help identify specific strengths and weaknesses across model families and sizes.

### E. Human Evaluation Approaches

Despite advances in automated metrics, human evaluation remains essential for assessing generative AI outputs, particularly for subjective qualities and complex reasoning.

*1) Human Evaluation Methodologies:* Several approaches to human evaluation have been developed:

### 1) Direct assessment approaches:

- **Absolute rating**: Judges rate outputs on quality, fluency, helpfulness scales
- **Likert scales**: 1-5 or 1-7 point ratings across multiple dimensions
- **Binary classification**: Accept/reject decisions for specific criteria
- **Detailed rubrics**: Comprehensive scoring across predefined categories

### 2) Comparative evaluation methods:

- **Pairwise comparison**: Judges choose between two model outputs
- **Best-worst scaling**: Ranking multiple outputs from best to worst
- **Tournament-style evaluation**: Head-to-head competitions across models

- **Elo rating systems**: Dynamic ranking based on pairwise comparisons

3) **2024-2025 human evaluation innovations**:
- **Constitutional evaluation**: Human assessment of ethical alignment
- **Interactive evaluation**: Extended conversations testing multi-turn capabilities
- **Expert evaluation**: Domain specialists (medical, legal, technical) assess specialized outputs
- **Cultural sensitivity assessment**: Cross-cultural evaluation of appropriateness
- **Preference learning**: Training reward models from human feedback

Human evaluation provides crucial insights that automated metrics may miss, but it faces challenges of subjectivity, cost, and reproducibility.

*2) Challenges and Solutions in Human Evaluation:* Human evaluation faces several key challenges:

**Key challenges:**
- **Subjectivity**: Inter-annotator disagreement and personal biases
- **Cost**: Expensive and time-consuming for large-scale evaluation
- **Reproducibility**: Difficulty replicating human judgment studies
- **Scale limitations**: Cannot evaluate all possible model outputs
- **Expertise requirements**: Specialized domains need expert evaluators

**2025 solutions and improvements:**
- **Hybrid evaluation**: Combining human judgment with automated metrics
- **Active learning**: Focusing human evaluation on uncertain or disagreement cases
- **Crowd-sourcing platforms**: Scalable evaluation using distributed workers
- **Expert networks**: Specialized evaluation for domain-specific applications
- **Preference modeling**: Training AI systems to predict human preferences

Hybrid approaches that combine human judgment with automated metrics offer promising directions for comprehensive evaluation.

*F. Evaluation Considerations for Scaling Laws Research*

*1) Performance Scaling Patterns Across Metrics:* Different evaluation metrics show distinct scaling patterns with model size and training data:

1) **Perplexity scaling**:
- Follows consistent power-law relationships: $L \sim N^{-\alpha}$ where $\alpha \approx 0.076$
- Smooth, predictable improvements across all model sizes
- Strong correlation with both parameter count and compute budget

2) **Benchmark accuracy patterns**:
- **Sigmoid scaling curves**: Performance remains low until capability thresholds
- **Phase transitions**: Rapid improvement followed by plateauing
- **Emergence patterns**: Capabilities appearing suddenly at specific scales

3) **Human preference alignment**:
- **Non-linear improvement**: Significant jumps at certain model sizes
- **Diminishing returns**: Smaller gains at largest scales without specific training
- **Task-dependent patterns**: Different alignment dimensions scale differently

4) **Safety and harmlessness metrics**:
- **Complex relationships**: May not improve automatically with scale
- **Potential degradation**: Larger models may generate more sophisticated harmful content
- **Intervention necessity**: Require specific training approaches (RLHF, constitutional AI)

**2024-2025 scaling observations:**
- **Reasoning capabilities**: Show sharp transitions (o1 models demonstrating step-function improvements)
- **Tool use**: Gradual improvement with sudden competency thresholds
- **Multimodal alignment**: Complex scaling requiring specialized training regimes
- **Constitutional behavior**: Requires dedicated constitutional training, not automatic with scale

Understanding these diverse scaling patterns is crucial for developing comprehensive scaling laws that go beyond simple perplexity-based formulations.
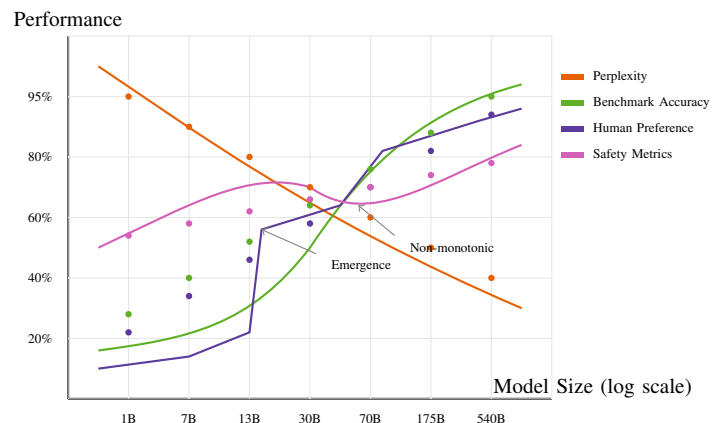


Fig. 8: Performance scaling patterns across different evaluation metrics. Perplexity shows smooth power-law improvements, benchmark accuracy exhibits sigmoid curves with emergence thresholds, human preference alignment demonstrates non-linear improvements with discrete jumps, and safety metrics show complex relationships that may not improve monotonically with scale.

*2) Benchmark Selection Criteria:* When evaluating scaling laws, the choice of benchmarks significantly impacts conclusions. Important selection criteria include:

1) **Discriminative power**:

   - **Difficulty appropriateness**: Challenging enough to differentiate model capabilities
   - **Avoiding ceiling effects**: Benchmarks should not saturate at current model scales
   - **Dynamic updating**: Ability to increase difficulty as models improve

2) **Contamination resistance**:

   - **Training data overlap**: Minimal presence in model training corpora
   - **Temporal separation**: Recent benchmark creation relative to training cutoffs
   - **Active monitoring**: Detection and mitigation of data leakage

3) **Comprehensive coverage**:

   - **Capability diversity**: Multiple cognitive and reasoning dimensions
   - **Domain breadth**: Various application areas and knowledge domains
   - **Modality inclusion**: Text, vision, audio, and cross-modal assessment
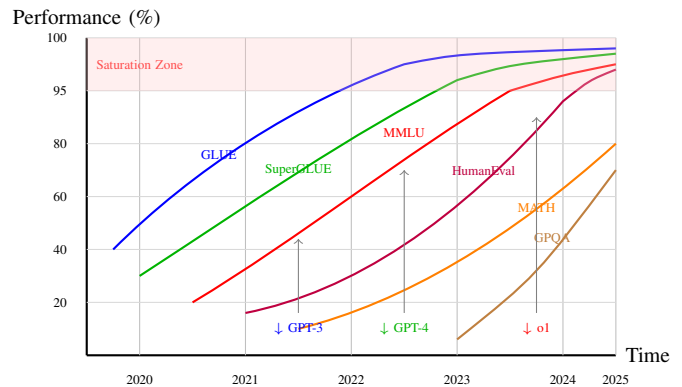
4) **Methodological rigor**:

   - **Standardized protocols**: Consistent evaluation procedures across studies
   - **Reproducibility**: Clear documentation enabling replication
   - **Statistical robustness**: Sufficient sample sizes and confidence intervals

**2025 benchmark selection best practices:**

- **Dynamic benchmark rotation**: Regular introduction of new evaluation tasks
- **Cross-validation**: Multiple benchmarks assessing similar capabilities
- **Real-world relevance**: Tasks reflecting actual deployment scenarios
- **Adversarial robustness**: Evaluation under challenging conditions

As models grow more capable, benchmark selection becomes increasingly important to avoid ceiling effects and ensure meaningful evaluation.



| Benchmark | Introduced | Current Top | Status |
|-----------|-----------|-------------|--------|
| GLUE | 2018 | 95.8% | Saturated |
| SuperGLUE | 2019 | 95.7% | Near Saturated |
| MMLU | 2020 | 90.0% | Approaching |
| HumanEval | 2021 | 92.0% | Active |
| MATH | 2021 | 83.3% | Active |
| GPQA | 2023 | 78.0% | Active |

Fig. 9: Timeline showing benchmark saturation patterns and the emergence of new evaluation challenges. Early benchmarks like GLUE rapidly reached saturation, requiring introduction of more challenging evaluations. Recent reasoning-focused benchmarks (MATH, GPQA) still provide discriminative power, while new capabilities demand entirely new evaluation frameworks.

*3) Evaluation Challenges at Scale:* Several challenges emerge when evaluating very large models:

1) **Benchmark saturation**: Top models approach perfect scores on many existing benchmarks, reducing their discriminative power.
2) **Computational requirements**:

   - **Evaluation costs**: Significant compute needed for comprehensive assessment
   - **Resource limitations**: Preventing thorough evaluation by smaller research groups
   - **Inference optimization**: Balancing evaluation completeness with computational efficiency

3) **Prompt sensitivity and optimization**:

   - **Hyperparameter sensitivity**: Large models highly responsive to prompt details
   - **Prompt engineering**: Optimal prompts may vary significantly between models
   - **Evaluation fairness**: Ensuring consistent prompting strategies across comparisons

4) **Emergent capability assessment**:

   - **Novel behaviors**: New capabilities not covered by existing benchmarks
   - **Capability detection**: Identifying and measuring previously unknown abilities
   - **Dynamic evaluation**: Adapting assessment frameworks to new capabilities

5) **2024-2025 specific challenges**:

- **Tool use complexity**: Evaluating sophisticated API calling and code execution
- **Multi-turn interactions**: Assessment across extended conversations
- **Agentic behavior**: Measuring planning, goal pursuit, and autonomous action
- **Constitutional alignment**: Assessing ethical behavior across diverse scenarios

Addressing these challenges requires continuous innovation in evaluation methodologies alongside model development.

*4) Holistic Evaluation Approaches:* Given the limitations of individual metrics, holistic evaluation approaches are increasingly important:

1) **Multi-metric aggregation**:

   - **Weighted scoring**: Combining metrics based on importance and reliability
   - **Principal component analysis**: Identifying underlying capability dimensions
   - **Ensemble evaluation**: Using multiple assessment methods for robustness
   - **Correlation analysis**: Understanding relationships between different metrics

2) **Capability mapping**:

   - **Radar charts**: Visualizing performance across multiple dimensions
   - **Heatmaps**: Showing capability patterns across model families and sizes
   - **Scaling trajectories**: Tracking improvement patterns over multiple scales
   - **Threshold identification**: Detecting capability emergence points

3) **Scenario-based evaluation**:

   - **End-to-end tasks**: Realistic applications requiring multiple capabilities
   - **Interactive scenarios**: Multi-turn conversations and complex interactions
   - **Domain-specific applications**: Specialized use cases (medical, legal, educational)
   - **Cross-modal integration**: Tasks requiring text, vision, and audio understanding

4) **2025 holistic evaluation innovations**:

   - **Constitutional integration**: Embedding ethical assessment throughout evaluation
   - **Real-world deployment metrics**: Performance in actual application settings
   - **Longitudinal tracking**: Monitoring capability development over time
   - **Human-AI collaboration**: Evaluating joint human-AI task performance

These approaches provide richer insights into how performance scales with model size and training data across different dimensions.

*G. Future Directions in Performance Evaluation*

*1) Adaptive Benchmarks:* To address benchmark saturation, adaptive benchmarks that automatically adjust difficulty based on model capability show promise:

1) **Dynamic difficulty adjustment**: Creating harder examples as models improve
2) **Curriculum evaluation**: Progressive difficulty increases to identify capability limits
3) **Personalized testing**: Adapting difficulty based on individual model performance
4) **Real-time updating**: Continuous benchmark refresh preventing memorization

**Adversarial evaluation:**

- **Red team automation**: Using AI systems to generate challenging test cases
- **Failure mode discovery**: Systematic search for model limitations
- **Robustness testing**: Evaluation under adversarial conditions
- **Security assessment**: Testing resistance to prompt injection and jailbreaking

**2025 adaptive benchmark initiatives:**

- **LiveBench evolution**: Continuously updated evaluation preventing gaming
- **AI-generated evaluation**: Models creating assessment tasks for other models
- **Interactive benchmarking**: Dynamic conversation-based evaluation
- **Community-driven evaluation**: Crowdsourced benchmark development

These approaches could maintain discriminative power even as models continue to scale.

*2) Alignment Evaluation:* As models grow more powerful, evaluating alignment with human values becomes increasingly important:

1) **Helpfulness optimization**: Measuring effective assistance without harmful outputs
2) **Harmlessness evaluation**: Resistance to generating dangerous or offensive content
3) **Honesty and truthfulness**: Factual accuracy and appropriate uncertainty expression
4) **Human autonomy respect**: Supporting human agency rather than manipulation
5) **Value diversity**: Consistency with varied cultural and ethical frameworks

**Advanced alignment metrics:**

- **Intent alignment**: Measuring model understanding of human goals
- **Value learning**: Assessment of acquired ethical principles
- **Cultural sensitivity**: Appropriate behavior across diverse contexts
- **Long-term safety**: Evaluation of risks from advanced capabilities

**2024-2025 constitutional evaluation tools:**

- **Anthropic Constitutional AI**: Systematic harmlessness assessment
- **DeepMind Sparrow**: Helpful, harmless, honest evaluation framework
- **OpenAI alignment evaluation**: Safety and capability assessment integration
- **Cross-organizational standards**: Shared ethical evaluation protocols

Developing robust metrics for these alignment dimensions represents a crucial frontier in evaluation research.

*3) Efficiency Metrics:* Evaluating performance relative to computational and environmental costs is gaining importance:

1) **Parameter efficiency**: Performance per parameter or per compute operation
2) **Data efficiency**: Performance relative to training data volume
3) **Inference efficiency**: Latency, throughput, and resource requirements
4) **Environmental impact**: Carbon footprint and energy consumption

**Environmental impact assessment:**

- **Carbon footprint**: Training and inference environmental costs
- **Energy efficiency**: Performance per watt measurements
- **Sustainable scaling**: Identifying eco-friendly scaling strategies
- **Resource optimization**: Minimizing computational requirements

**2025 efficiency evaluation advances:**

- **Green AI metrics**: Comprehensive environmental impact assessment
- **Edge deployment evaluation**: Performance on resource-constrained devices
- **Real-time efficiency**: Latency and throughput optimization
- **Cost-benefit analysis**: Economic efficiency of scaling investments

These efficiency metrics provide important context for scaling laws, highlighting the trade-offs involved in pursuing scale.

*4) Robustness Evaluation:* Assessing model robustness across diverse conditions is critical for real-world applications:

1) **Distribution shift robustness**: Performance on out-of-distribution inputs
2) **Prompt robustness**: Consistency across different phrasings of the same request
3) **Adversarial robustness**: Resistance to inputs designed to cause failures
4) **Temporal robustness**: Stability of performance over time as the world changes

**Adversarial robustness:**

- **Prompt injection resistance**: Security against malicious inputs
- **Jailbreak prevention**: Maintaining safety constraints under attack

- **Backdoor detection**: Identifying hidden harmful behaviors
- **Social engineering resistance**: Avoiding manipulation attempts

**2024-2025 robustness evaluation developments:**

- **Automated red teaming**: AI-assisted safety testing
- **Continuous monitoring**: Real-time robustness assessment during deployment
- **Multi-stakeholder evaluation**: Diverse perspectives on safety and robustness
- **Regulatory compliance**: Meeting emerging AI safety standards

**System-level robustness:**

- **Integration testing**: Performance when combined with other systems
- **Failure graceful degradation**: Behavior under partial system failures
- **Recovery evaluation**: Ability to recover from errors or attacks
- **Uncertainty quantification**: Appropriate confidence expression

Robustness evaluation helps identify limitations that may not be apparent in standard benchmarks.

## H. Conclusion

The evaluation of generative AI models continues to evolve rapidly alongside the models themselves. **2024-2025 has marked a pivotal period** with the emergence of sophisticated reasoning capabilities, tool use, constitutional alignment, and multimodal understanding. As models scale in size and capability, evaluation methodologies must adapt to provide meaningful insights across diverse performance dimensions.

**Key developments shaping the future of evaluation include:**

- **Constitutional and alignment assessment** moving beyond pure capability metrics
- **Dynamic and adaptive benchmarks** addressing rapid capability growth
- **Holistic evaluation frameworks** combining multiple assessment dimensions
- **Real-world deployment metrics** measuring practical application effectiveness
- **Efficiency and sustainability** considerations in scaling law formulations

By combining quantitative metrics, comprehensive frameworks, standardized benchmarks, human evaluation, and emerging assessment methodologies, researchers can develop increasingly sophisticated understanding of scaling laws and their implications for model development, deployment, and societal impact. **The future of AI evaluation lies in comprehensive, adaptive, and ethically-grounded assessment frameworks** that can keep pace with rapid capability advancement while ensuring beneficial and safe AI development.

## Volume 14 Issue 8, August 2025
**Fully Refereed | Open Access | Double Blind Peer Reviewed Journal**
[www.ijsr.net](www.ijsr.net)

Paper ID: SR25727102615     DOI: https://dx.doi.org/10.21275/SR25727102615     938

## VI. Cost Analysis of Generative AI Models

### A. The Economic Dimension of Scaling Laws

While technical aspects of scaling laws have received significant attention, the economic implications are equally important for practical applications of generative AI [2]. The story of AI economics has undergone a dramatic transformation since 2022, where what once cost hundreds of thousands of dollars now requires investments measured in hundreds of millions. Training costs for frontier models have escalated from approximately $4.6 million for GPT-3 in 2020 to an estimated $100-200 million for models like GPT-4o and Claude-3.5 Sonnet in 2024-2025 [4], [11]. Yet paradoxically, as training costs have soared, inference costs have plummeted due to remarkable advances in hardware and optimization techniques [11], [17].

This economic transformation tells a compelling story of an industry reaching new scales of ambition while simultaneously becoming more efficient. The emergence of specialized AI chips beyond NVIDIA's traditional dominance, the maturation of inference optimization techniques that can reduce costs by 10-50x, and the introduction of carbon pricing mechanisms that add environmental considerations to the traditional compute-focused equation have fundamentally reshaped how organizations approach AI investments [2], [17]. Understanding these evolving economics has become essential for organizations making strategic decisions about AI investments, deployment strategies, and long-term technological roadmaps.

### B. Fundamental Cost Drivers in Generative AI

*1) Training Costs:* The journey of training large generative AI models represents one of the most remarkable escalations in computational economics in modern computing history [2]. To understand the magnitude of this transformation, consider that training GPT-3 with its 175 billion parameters cost approximately $4.6 million in 2020 [4]. Fast forward to 2024-2025, and training frontier models like GPT-4o or Claude-3.5 Sonnet now requires investments of $100-200 million, representing a 20-40x increase in just four years. This exponential growth in costs tells the story of an industry pushing the absolute boundaries of what's computationally feasible.

*a) Computational Requirements:* The computational demands that drive these costs follow well-established mathematical relationships, yet their practical implications have grown beyond what early researchers anticipated [8]. For transformer-based models, training requires approximately $6N \times D$ FLOPs, where $N$ represents the number of parameters and $D$ represents the dataset size in tokens [8]. While this relationship appears straightforward, its real-world implications have become staggering.

Consider the computational journey of recent models: GPT-4o, with an estimated 1.8 trillion parameters, required approximately $2.5 \times 10^{24}$ FLOPs for training, translating to $150-200 million at current H100 rates. Claude-3.5 Sonnet, with its estimated 500 billion parameters, consumed $80-120 million including infrastructure and extensive experimentation. Even

Meta's more modestly sized Llama 3 405B required $20-30 million for the base model, though this figure excludes significant infrastructure amortization costs [16].

These costs scale predictably with model size, but the relationship has become more nuanced due to architectural innovations [15]. The scaling patterns now tell different stories depending on the optimization approach: Chinchilla-optimal scaling suggests training costs scale approximately as $N^{1.5}$ [7], while extended scaling regimes that prioritize inference efficiency may see costs scale as $N^{1.8}$. Mixture-of-Experts architectures offer a more optimistic narrative, potentially reducing effective costs by 2-4x through sparse activation patterns that only engage relevant portions of the model during training [15].
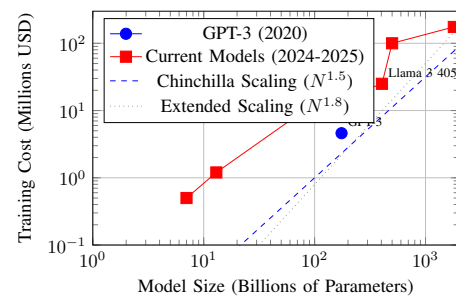


Fig. 10: Training cost scaling relationships for generative AI models, showing the dramatic cost escalation from GPT-3 to current frontier models and theoretical scaling curves.

*b) Hardware Considerations:* The hardware landscape underlying these computational demands has undergone its own dramatic transformation, creating a complex economic ecosystem [17]. Today's training infrastructure represents a significant departure from the relatively simple GPU clusters of just a few years ago. Modern frontier models typically require orchestrated deployments of 8,000-24,000 H100 GPUs, each costing $2.85-4.20 per hour depending on cloud provider and commitment level [17].

The economics of this hardware ecosystem tell a story of rapid innovation and equally rapid obsolescence. NVIDIA's H100 80GB represents the current gold standard, while the legacy A100 80GB continues to serve important roles at $1.85-2.50 per hour [17]. Google's TPU v5e has emerged as a competitive alternative at $1.35-2.10 per hour, particularly for transformer training workloads. Perhaps most intriguingly, custom silicon solutions from AWS (Trainium2) and Meta (MTIA) promise 30-50% cost reductions for specific workloads, suggesting a future where specialized hardware may challenge the current GPU-centric paradigm [17].

TABLE VII: Current Hardware Ecosystem for AI Training (2024-2025)

| Hardware | Cost/Hour | Memory | Performance | Specialty |
|---|---|---|---|---|
| NVIDIA H100 80GB | $2.85-4.20 | 80GB HBM3 | Baseline | General Purpose |
| NVIDIA A100 80GB | $1.85-2.50 | 80GB HBM2e | 0.6x H100 | Legacy/Proven |
| Google TPU v5e | $1.35-2.10 | Variable | 0.8x H100 | Transformer Opt. |
| AWS Trainium2 | $1.90-2.80 | 64GB | 0.7x H100 | Cost Optimized |
| Meta MTIA | Custom | 128GB | 0.5x H100 | Inference Focus |

The distributed nature of modern training introduces its own economic complexities. High-bandwidth networking infrastructure, typically InfiniBand or custom interconnects, adds 15-25% to raw compute costs. The need for sophisticated 3D parallelism—combining data, model, and pipeline parallelism with advanced load balancing—requires expertise that commands premium prices in the current market.

Perhaps most challenging from an economic perspective is the rapid pace of hardware advancement. Organizations investing in current-generation hardware face depreciation rates of 40-60% annually, as newer accelerators consistently provide 2-3x better performance per dollar every 18-24 months. This creates a perpetual tension between investing in current capabilities and waiting for next-generation improvements.

*c) Time Factors and Opportunity Costs:* The temporal dimension of training costs often exceeds the direct computational expenses in strategic importance. Current frontier models require 3-6 months of continuous training, creating substantial opportunity costs that can exceed the direct compute investments. This timeline creates a fascinating economic dynamic where the cost of time often matters more than the cost of computation.

The research velocity implications tell a particularly compelling story. OpenAI's GPT-4o required an estimated 4-month training period, during which the opportunity costs—including delayed market entry, competitive positioning, and alternative research directions—potentially exceeded the $150-200 million direct compute costs. Anthropic has responded to this challenge with an iterative training approach for Claude-3.5, using multiple shorter runs to reduce time-to-insight and maintain research momentum. Meta's approach with Llama 3 involved parallelized experimentation across multiple training runs, prioritizing research velocity over computational efficiency [16].

For organizations developing multiple models or conducting extensive hyperparameter optimization, these time factors fundamentally reshape the economic equation. Research velocity often becomes the primary constraint rather than raw compute budget, leading to investment strategies that prioritize faster iteration over absolute cost minimization.

*2) Inference Costs:* While training costs capture attention due to their magnitude, inference costs—the expenses associated with using trained models to generate outputs—often dominate the lifetime economic equation for widely deployed models [11]. The inference cost landscape has been transformed by optimization techniques that can reduce costs by 10-50x compared to naive deployments [11], [17].

*a) Computational Requirements and Current Pricing:* For transformer models, inference requires approximately $2N$ FLOPs per token generated (where $N$ is parameter count) [8]. Current 2024-2025 pricing reflects significant optimization improvements [11]:

TABLE VIII: Current Market Rates for AI Inference (2024-2025)

| Model | Input ($/1M) | Output ($/1M) | Parameters | Provider |
|---|---|---|---|---|
| GPT-4o | $5.00 | $15.00 | $\sim$1.8T | OpenAI |
| Claude-3.5 Sonnet | $3.00 | $15.00 | $\sim$500B | Anthropic |
| Gemini Ultra | $7.00 | $21.00 | $\sim$1.5T | Google |
| Llama 3 405B | $2.50-4.00 | $8.00-12.00 | 405B | Various |

These prices represent the fully loaded cost including infrastructure, optimization, and margin [11]. The underlying computational costs are significantly lower:

- **Raw compute cost**: $0.10-0.30 per 1M tokens for optimized deployments [11]
- **Infrastructure overhead**: Adds 40-60% (load balancing, monitoring, networking) [17]
- **Provider margins**: 300-500% markup over infrastructure costs [2]

Unlike training costs, which are incurred once, inference costs accumulate with each use of the model. For popular models, the lifetime inference costs typically far exceed the initial training investment, creating an economic dynamic where optimization becomes increasingly critical as deployment scales.
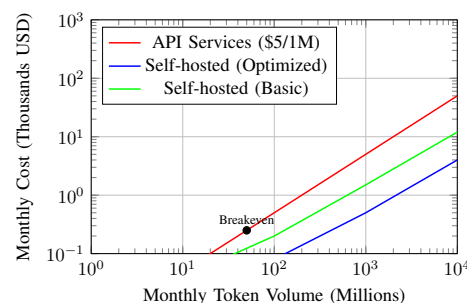


Fig. 11: Cost comparison between API services and self-hosted solutions across different usage volumes, showing breakeven points and optimization benefits.

*b) Deployment Considerations and Infrastructure Evolution:* Modern inference deployment has evolved far beyond simple model serving, incorporating sophisticated optimization and management layers [11], [17]. Deploying models for inference introduces additional costs beyond raw computation, and these costs are associated with serving infrastructure, load balancing and autoscaling, networking and data transfer, monitoring and logging, and continuous integration. These deployment costs can add 20-30% to the raw compute costs of inference, as demonstrated by detailed cost breakdowns from industry practitioners [11].

**2024-2025 Infrastructure Stack:**

- **Optimized serving frameworks**: vLLM, TensorRT-LLM, and custom kernels providing 2-5x throughput improvements [11]
- **Dynamic batching**: Continuous batching algorithms that improve utilization by 3-8x [11]
- **Multi-tenant serving**: Serving multiple models on shared hardware with 40-60% cost reduction [17]

- **Edge deployment**: Processing suitable workloads on edge devices, reducing latency and costs by 70-90% [17]

**Advanced Deployment Patterns:**

- **Speculative decoding**: Using smaller models to accelerate inference of larger models, reducing costs by 40-70% [11]
- **Mixture-of-Experts routing**: Activating only relevant model components, reducing effective inference costs by 60-80% [15]
- **Cascading architectures**: Routing simple queries to smaller models, complex queries to larger models [11]

These deployment innovations have fundamentally changed the economics of inference, making it feasible to serve sophisticated models at scales that would have been prohibitively expensive just 18 months ago.

*c) Optimization Techniques and Real-World Impact:* The inference optimization landscape has matured dramatically, with techniques now routinely achieving 10-50x cost reductions [11], [17]:

**Quantization Advances:**

- **4-bit quantization** (GPTQ, AWQ): 75-85% memory reduction with $<5\%$ quality loss [12]
- **1-bit quantization** (BitNet): 90-95% reduction for specific architectures, though with greater quality trade-offs [12]
- **Dynamic quantization**: Runtime adaptation based on input complexity [15]

**Model Architecture Optimizations:**

- **Pruning techniques**: Structured and unstructured pruning achieving 50-80% parameter reduction [12]
- **Distillation improvements**: Student models achieving 80-95% of teacher performance at 10-20% of the cost [15]
- **Early exit mechanisms**: Dynamic computation based on confidence thresholds [15]
- **Caching**: Storing frequently requested outputs can eliminate redundant computation
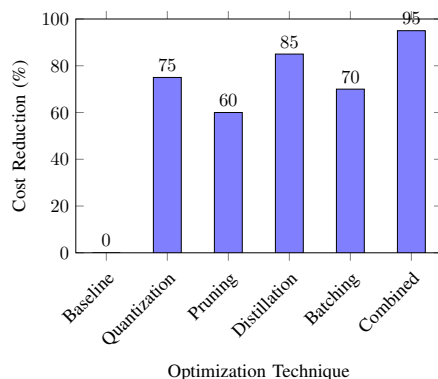- **Batching**: Processing multiple requests simultaneously improves hardware utilization



Fig. 12: Inference cost reduction potential from various optimization techniques, showing cumulative benefits when combined systematically.

**Real-World Optimization Results (2024-2025):**

- **Anthropic's Claude optimizations**: 60% cost reduction through custom kernels and batching [11]
- **OpenAI's efficiency improvements**: 50% cost reduction for GPT-4 through infrastructure optimization [11]
- **Meta's Llama optimizations**: Open-source techniques enabling 70-90% cost reduction for self-hosted deployments [16]

These optimization techniques have created a bifurcated market where organizations with optimization expertise can achieve dramatically lower costs than those relying on default implementations.

### C. Total Cost of Ownership (TCO) Framework

Understanding the full economic impact of generative AI requires looking beyond raw computation to consider the total cost of ownership across the entire lifecycle. The true story of AI economics emerges only when organizations account for every component that contributes to their AI investments—from the obvious compute costs to the hidden expenses that often dwarf the infrastructure spending [2]. In 2024-2025, this comprehensive view has become critical as the complexity of AI deployments has grown exponentially, and what initially appears to be a straightforward technology investment reveals itself as a multifaceted economic undertaking.

*1) Beyond Raw Compute:*

*a) Infrastructure Costs:* A comprehensive TCO analysis must account for the entire infrastructure ecosystem supporting generative AI operations [17]. The 2024-2025 landscape reveals a sophisticated stack where each component contributes meaningfully to the total cost equation:

TABLE IX: Complete Infrastructure Cost Breakdown (2024-2025 Hourly Rates)

| Component | Cost/Hour | Specification | Overhead |
|---|---|---|---|
| H100 Compute | $2.85-4.20 | 80GB HBM3 | Baseline |
| Storage Systems | $0.12-0.18 | 1TB High-Perf | 3-4% |
| Networking | $0.08-0.15 | 100Gbps | 2-3% |
| Load Balancing | $0.05-0.08 | Traffic Dist. | 1-2% |
| Orchestration | $0.35-0.45 | Kubernetes | 8-10% |
| Monitoring | $0.10-0.15 | Full Stack | 2-3% |
| **Total Overhead** | **+25-40%** | **Above Compute** | **Typical** |

These infrastructure costs typically add 25-40% to the raw compute costs, depending on the specific deployment architecture and cloud provider [17]. However, the real revelation for organizations has been discovering how this seemingly straightforward infrastructure component represents only the foundation of their total investment.

*b) Human Resources: The Hidden Dominant Cost:* The human costs associated with developing and maintaining generative AI systems have emerged as the most significant component of TCO for most organizations [2]. The 2024-2025 talent market reveals a striking economic reality: specialized AI expertise commands premium rates that often exceed infrastructure costs by 3-5x.

TABLE X: Current Market Rates for AI Talent (2024-2025)

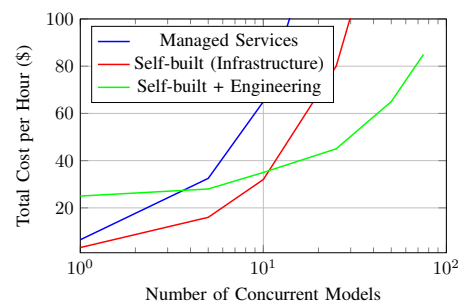| Role | Rate/Hour | Specialty | Demand |
|------|-----------|-----------|--------|
| Senior ML Engineers | $85-150 | Frontier Models | Extreme |
| AI Infrastructure | $75-120 | Deploy/Optimize | Very High |
| DevOps + AI | $65-95 | Orchestration | High |
| Data Scientists | $70-110 | LLM/Prompting | High |
| AI Safety/Alignment | $90-140 | Regulatory | Growing |



Fig. 13: Economic comparison of managed services vs. self-built solutions across different deployment scales, showing crossover points where each approach becomes optimal.

The total human resource investment extends far beyond these hourly rates. Organizations typically require 3-6 months for initial implementation phases, involving teams of 4-8 specialists. Ongoing optimization and maintenance demands 1-2 full-time equivalent positions per major deployment. Training and education to keep teams current with rapidly evolving technologies adds 15-25% to direct compensation costs.

*c) Operational Overhead:* Beyond direct development and infrastructure costs, operational considerations add substantial complexity and expense to AI deployments [2]. The 2024-2025 operational landscape reflects the maturation of AI governance and compliance requirements:

- **System administration**: 10-15% of infrastructure costs for access controls, security policies, and routine maintenance
- **Security management**: 15-25% overhead for model protection, data privacy, and threat monitoring
- **Compliance and governance**: 5-15% for regulatory adherence, audit trails, and documentation (varies significantly by industry)
- **Documentation and knowledge management**: 8-12% for maintaining system understanding and operational procedures
- **Incident response and reliability engineering**: 12-20% for monitoring, troubleshooting, and performance optimization

These operational factors typically add 20-35% to the total cost of ownership for enterprise deployments [2]. However, the percentage has been increasing as regulatory scrutiny intensifies and organizations discover the complexity of operating AI systems at scale.

*2) Managed Services vs. Self-Built Solutions:* The build-versus-buy decision has evolved dramatically in 2024-2025, as managed services have matured while self-built solutions have become more sophisticated [2]. Organizations now face nuanced trade-offs that extend far beyond simple cost comparisons.

*a) Economic Analysis: The True Cost Comparison:* Direct cost comparisons reveal compelling patterns that have shifted the economic equation significantly since 2022-2023:

- **Managed service (full-stack)**: $6.50-8.75/hr for enterprise-grade model serving on equivalent hardware
- **Self-built infrastructure only**: $3.20-4.50/hr for comparable compute and storage resources
- **Self-built with minimal engineering**: $12.50-18.00/hr including 0.25 FTE specialized engineer
- **Self-built with full optimization**: $25.00-35.00/hr including dedicated team for efficiency optimization

These comparisons suggest that managed services have achieved economic parity with self-built solutions for organizations operating fewer than 50-75 concurrent models [2]. The crossover point has increased significantly as managed service providers have achieved economies of scale and optimization expertise.

### D. Cost Optimization Strategies

As generative AI has transitioned from experimental technology to production infrastructure, organizations have developed increasingly sophisticated strategies for managing and optimizing costs. The optimization landscape of 2024-2025 tells a story of remarkable innovation, where the most advanced practitioners achieve 70-90% cost reductions compared to naive implementations [11], [17]. This transformation reflects not just technological progress, but the emergence of a mature discipline combining cloud economics, algorithmic optimization, and infrastructure engineering.

*1) Cloud Provider Strategies:* The cloud provider landscape has become significantly more competitive and sophisticated since 2022, offering organizations numerous pathways to cost optimization. Each major provider has developed specialized approaches to AI workloads, creating opportunities for strategic optimization that extend far beyond simple price comparisons.

*a) AWS Cost Optimization Evolution:* Amazon Web Services has emerged as a leader in AI cost optimization tools, developing a comprehensive ecosystem of cost management capabilities tailored specifically for machine learning workloads [17]:

- **Compute Savings Plans**: Commitment-based discounts offering 20-40% savings for predictable AI workloads,

with specialized ML Savings Plans providing up to 50% discounts for SageMaker usage

- **Spot Instances for ML**: Utilization of excess capacity at 60-90% discounts for fault-tolerant training workloads, with new spot fleet management tools designed specifically for distributed ML training
- **Reserved Inference Capacity**: Provisioned throughput guarantees for inference workloads, offering 30-50% cost reductions for predictable usage patterns
- **Batch Inference Optimization**: Specialized batch processing services that can reduce inference costs by 40-70% for non-real-time applications
- **Multi-Model Endpoints**: Serving multiple models on shared infrastructure, reducing costs by 50-80% for organizations with diverse model portfolios

TABLE XI: Cloud Provider AI Optimization Strategies Comparison (2024-2025)

| Strategy | AWS | Google Cloud | Azure | Savings |
|---|---|---|---|---|
| Committed Use | 20-50% | 30-60% | 25-55% | High |
| Spot/Preemptible | 60-90% | 60-80% | 60-80% | Very High |
| Custom Silicon | Trainium2 | TPU v5e | Custom FPGA | 30-50% |
| Multi-Model | SageMaker | Vertex AI | ML Studio | 50-80% |
| Batch Processing | 40-70% | 45-75% | 35-65% | High |
| Edge Deployment | 70-90% | 75-85% | 60-80% | Very High |

*2) Technical Optimization Approaches:* Beyond cloud-specific strategies, the 2024-2025 landscape has witnessed remarkable advances in technical optimization approaches that address cost reduction at multiple architectural levels. These techniques represent the cutting edge of cost optimization, where algorithmic innovation meets infrastructure engineering to achieve dramatic efficiency gains.

*a) Model Architecture Optimization:* The evolution of model architectures has been driven as much by cost considerations as by performance requirements, leading to innovative designs that maximize efficiency:

- **Task-Specific Model Selection**: Deploying appropriately sized models for specific use cases, with 70B parameter models often achieving 90-95% of 405B model performance at 40-60% of the cost
- **Mixture-of-Experts (MoE) Architectures**: Sparse models that activate only relevant parameters, reducing effective inference costs by 60-80% while maintaining capability [15]
- **Parameter-Efficient Fine-Tuning**: Techniques like LoRA (Low-Rank Adaptation) and QLoRA enabling model customization with minimal parameter additions, reducing training costs by 90-95%
- **Progressive Model Architectures**: Hierarchical models that can provide quick responses for simple queries while reserving full computation for complex tasks

*b) Training Optimization Revolution:* Training optimization has evolved from simple hyperparameter tuning to sophisticated systems that optimize the entire training process for cost efficiency:

- **Curriculum Learning**: Strategic data ordering that reduces training time by 20-40% while maintaining or improving final performance

- **Optimal Batch Size Selection**: Dynamic batch sizing that balances memory utilization and convergence speed, typically improving training efficiency by 15-30%
- **Mixed-Precision Training**: Leveraging FP16, BF16, and even INT8 training to reduce memory requirements and accelerate training by 40-70%
- **Gradient Accumulation Strategies**: Simulating larger batch sizes with limited memory, enabling optimal training on smaller, more cost-effective hardware configurations

*c) Inference Optimization Mastery:* Inference optimization has become perhaps the most critical area for cost reduction, as deployment-scale inference costs often dwarf training expenses:

- **Continuous Batching**: Advanced batching algorithms that maximize hardware utilization while minimizing latency, improving throughput by 3-8x [11]
- **Speculative Decoding**: Using smaller draft models to accelerate larger model inference, reducing costs by 40-70% for generation tasks [11]
- **KV-Cache Optimization**: Sophisticated memory management for transformer attention mechanisms, reducing memory requirements by 30-50%
- **Custom Kernels**: Hand-optimized CUDA kernels providing 2-5x performance improvements for specific operations
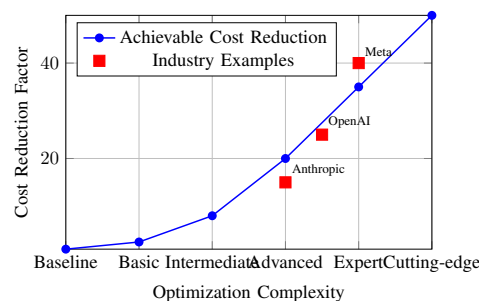


Fig. 14: Cost reduction potential versus optimization complexity, showing the exponential benefits of advanced optimization techniques and real-world industry achievements.

*E. Economic Implications of Scaling Laws*

Scaling laws have profound economic implications that fundamentally shape the development and deployment of generative AI. The mathematical relationships that govern model performance create corresponding economic dynamics that influence every strategic decision from research investments to production deployments [7], [8]. In 2024-2025, these economic implications have become central to competitive strategy, as organizations navigate the complex trade-offs between performance aspirations and financial constraints.

*1) Cost Scaling Patterns:* Understanding how costs scale with model size and deployment characteristics provides the foundation for rational economic decision-making in generative AI. The scaling relationships observed in 2024-2025 reveal both predictable patterns and surprising deviations that reshape strategic thinking.

*a) Training Cost Scaling:* Training costs follow well-established mathematical relationships, but their economic implications have evolved as the industry has gained experience with large-scale deployments [7], [8]:

- **Chinchilla-Optimal Scaling**: Training costs scale approximately as $N^{1.5}$, where $N$ represents parameter count, reflecting the balanced scaling of model size and training data [7]
- **Kaplan Scaling**: Earlier approaches suggested $N^2$ scaling, representing compute-optimal training that minimizes training time rather than total cost [8]
- **Extended Scaling Regimes**: Current frontier models often scale as $N^{1.8}$ to $N^{2.2}$, reflecting training beyond Chinchilla optimality to achieve superior inference efficiency [7]
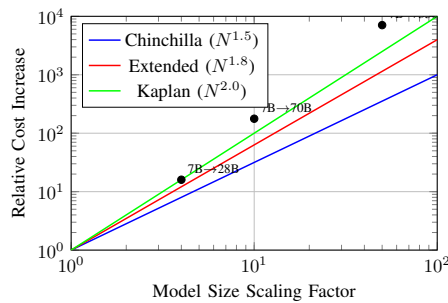


Fig. 15: Training cost scaling relationships showing different scaling regimes and their economic implications for model size decisions.

*b) Inference Cost Scaling:* Inference costs exhibit different scaling characteristics that often dominate the long-term economic equation for deployed models [11]:

- **Parameter Count Impact**: Inference costs scale linearly with model size for equivalent workloads, creating predictable cost relationships
- **Sequence Length Scaling**: Costs increase quadratically with input/output length due to attention mechanism complexity, making longer interactions disproportionately expensive
- **Batch Processing Benefits**: Effective costs per request decrease significantly with intelligent batching, often achieving 3-8x efficiency improvements [11]
- **Fixed Infrastructure Amortization**: Large-scale deployments benefit from spreading infrastructure overhead across many requests, reducing per-request costs by 40-70%

*2) Economic Decision Framework:* Organizations can leverage scaling laws to inform systematic economic decisions about AI investments. The frameworks that have emerged in 2024-2025 reflect sophisticated thinking about the interplay between technical scaling relationships and business objectives.

*a) Model Size Selection:* The model sizing decision has become a sophisticated optimization problem that balances multiple objectives and constraints:

- **Diminishing Returns Analysis**: Understanding where additional scale provides minimal performance gains relative to cost increases
- **Task-Specific Optimization**: Recognizing that different applications have different sensitivity to model scale, enabling right-sizing for specific use cases
- **Quality Threshold Requirements**: Identifying minimum viable performance levels that enable business objectives while minimizing costs
- **ROI-Based Decision Making**: Evaluating model investments based on expected business returns rather than pure technical performance

*b) Build vs. Buy Decision Framework:* The build-versus-buy decision has evolved into a sophisticated framework that considers multiple dimensions beyond simple cost comparisons [2]:

- **Volume-Based Economic Analysis**: Organizations serving fewer than 10-50 million tokens monthly often find API services more economical than self-hosting
- **Infrastructure Breakeven**: Self-hosting becomes attractive when serving 100+ million tokens monthly with dedicated engineering resources
- **Optimization Premium**: Organizations with advanced optimization capabilities can achieve breakeven at lower volumes (20-50 million tokens monthly)
- **Strategic Control Considerations**: Industries with strict data governance requirements may justify higher self-hosting costs for control and compliance

TABLE XII: Economic Decision Framework Matrix

| Monthly Volume | API Service | Managed | Self-Hosted | Recommendation |
|---|---|---|---|---|
| <10M tokens | $50-200 | $100-400 | $500-2000 | API Service |
| 10-50M tokens | $500-2K | $1K-4K | $2K-8K | API/Managed |
| 50-100M tokens | $2.5K-10K | $5K-20K | $5K-15K | Managed/Self |
| 100M+ tokens | $10K-50K | $20K-80K | $10K-30K | Self-Hosted |

*F. Future Cost Trends*

The economic landscape of generative AI stands at a fascinating inflection point, where multiple technological, market, and regulatory forces will converge to reshape cost structures over the next 3-5 years. Understanding these emerging trends has become critical for strategic planning, as organizations must balance current deployment decisions against rapidly evolving economic realities.

*1) Hardware Improvements:* The hardware revolution underlying generative AI continues to accelerate, with multiple technological threads converging to deliver substantial cost improvements over the next several years. The 2025-2028 hardware roadmap suggests cost reductions of 40-70% for equivalent workloads, driven by specialized architectures, improved manufacturing processes, and innovative system designs.

*a) Next-Generation AI Accelerators:* The evolution of AI-specific hardware represents perhaps the most significant driver of future cost reductions:

- **NVIDIA's Roadmap**: The transition from H100 to B100 (2025) and subsequent generations promises 3-5x

performance improvements per dollar, with architectural optimizations specifically targeting transformer workloads

- **Google's TPU Advancement**: TPU v6 and beyond focusing on inference optimization, potentially achieving 4-8x cost efficiency improvements for production deployments
- **Custom Silicon Proliferation**: AWS Trainium3, Meta's next-generation MTIA, and Apple's server-class AI chips suggesting 50-80% cost reductions for specific workloads by 2026-2027
- **Sparse Computation Accelerators**: Hardware designed specifically for Mixture-of-Experts and other sparse architectures, enabling 60-90% reductions in effective compute costs
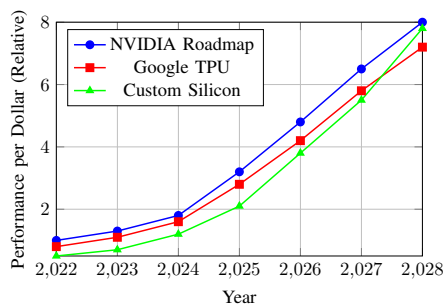


Fig. 16: Projected performance per dollar improvements for different hardware categories through 2028, showing the acceleration of cost efficiency gains.

*2) Software Optimization:* Software optimization represents perhaps the most dramatic opportunity for cost reduction, with emerging techniques suggesting 10-100x improvements in efficiency for specific applications. The software optimization landscape of 2025-2028 will likely transform the economics of AI deployment more dramatically than hardware advances alone.

*a) Algorithmic Breakthroughs:* Fundamental algorithmic improvements promise to reshape the cost-performance equation:

- **Few-Shot and Zero-Shot Training**: Techniques that dramatically reduce training data requirements, potentially cutting training costs by 80-95% for specialized applications
- **Continual Learning Architectures**: Models that can be updated incrementally without full retraining, reducing update costs by 90-99%
- **Meta-Learning Advances**: Algorithms that learn to learn more efficiently, potentially reducing training requirements for new domains by 70-90%
- **Speculative Execution Evolution**: Advanced techniques that could reduce inference costs by 60-90% while maintaining quality

*3) Market Dynamics:* The AI market is evolving rapidly toward increased competition and specialization, with profound implications for cost structures and strategic positioning. The 2025-2028 period will likely witness fundamental shifts in market structure that reshape the economic landscape.

*a) Competition and Commoditization:* Market maturation is driving both competition and differentiation:

- **Open Source Acceleration**: Continued advancement of open-source models reducing the premium for proprietary solutions by 40-70%
- **Cloud Provider Competition**: Major cloud platforms competing aggressively on AI pricing, potentially reducing costs by 30-60% over 3-5 years
- **Specialized Providers**: Emergence of focused AI infrastructure providers offering 50-80% cost advantages for specific workloads
- **API Standardization**: Emerging standards reducing switching costs and increasing price competition

*4) Regulatory Considerations:* Regulatory developments will introduce new cost factors while potentially reshaping competitive dynamics. The 2025-2028 regulatory landscape suggests 10-40% cost increases for compliance, offset by potential efficiency requirements and standardization benefits.

*a) Environmental and Safety Regulations:* Environmental considerations are becoming mandatory cost factors:

- **Carbon Pricing**: Implementation of carbon taxes and cap-and-trade systems adding 5-20% to compute costs depending on region and energy source
- **Mandatory Safety Testing**: Requirements for extensive testing and validation potentially adding 15-30% to development costs
- **Explainability and Transparency**: Regulatory demands for model interpretability requiring additional infrastructure and development investment
- **Data Privacy and Security**: Stricter requirements for data handling and privacy protection increasing infrastructure and operational costs

TABLE XIII: Future Cost Impact Summary (2025-2028 Projections)

| Factor | Cost Impact | Timeline | Certainty |
|---|---|---|---|
| Hardware Advances | -40 to -70% | 2025-2027 | High |
| Software Optimization | -60 to -90% | 2024-2026 | Medium |
| Market Competition | -30 to -60% | 2024-2028 | High |
| Carbon Pricing | +5 to +20% | 2025-2027 | Medium |
| Safety Regulation | +15 to +30% | 2026-2028 | Low |
| Privacy Compliance | +10 to +25% | 2025-2027 | Medium |
| **Net Impact** | **-20 to -60%** | **Variable** | **Medium** |

*b) Synthesis: The Future Cost Landscape:* The convergence of these trends suggests a complex future cost landscape characterized by dramatic cost reductions for standard capabilities, persistent premiums for advanced capabilities, and new cost categories for compliance and environmental factors. Organizations that anticipate and prepare for these evolving cost dynamics will be positioned to optimize their AI investments and maintain competitive advantages in an increasingly sophisticated and competitive landscape [2], [17].

The key insight is that future AI economics will be characterized not by uniform cost reduction, but by increasing differentiation between optimized and naive approaches, between standard and premium capabilities, and between compliant

and non-compliant deployments. Success will require not just technical optimization, but strategic positioning within this evolving economic landscape.

## VII. Conclusion: The Future of Scaling Laws in Generative AI

### A. Synthesis of Key Insights

Throughout this comprehensive analysis, we have explored the multifaceted nature of scaling laws in generative AI, examining how model size, training data, performance metrics, and economic considerations interweave to shape the future of artificial intelligence. The journey from early empirical observations to the sophisticated understanding we possess in 2025 reveals a field that has matured dramatically while continuing to surprise us with new insights and emergent phenomena.

Our exploration reveals that scaling laws represent far more than empirical observations—they constitute fundamental properties of neural network learning that govern how intelligence emerges from computational resources. The consistent power-law relationships between model size, data volume, and performance across multiple orders of magnitude suggest deep mathematical principles underlying these patterns [7], [8]. These relationships have evolved from simple formulations to nuanced frameworks that account for optimal resource allocation, emergent capabilities, and economic constraints.

*a) The Evolution of Scaling Understanding:* The transformation of our understanding—from Kaplan's initial compute-optimal formulation [8] to the Chinchilla paradigm's data-optimal insights [7], and extending to the sophisticated approaches of 2024-2025—demonstrates the dynamic nature of this field. We have witnessed a fundamental shift from the 1.7:1 token-to-parameter ratio proposed by Kaplan [8] to the 20:1 ratio identified by Chinchilla [7], and the exploration of even more data-intensive training regimes in current frontier models like GPT-4o and Claude-3.5 Sonnet. This evolution reflects not just improved understanding, but a recognition that optimal scaling depends on deployment context, computational constraints, and economic objectives.

The emergence of constitutional AI training approaches has further refined our understanding, showing how data quality and alignment considerations can dramatically influence scaling efficiency. Modern training approaches achieve superior performance not just through scale, but through sophisticated data curation, constitutional training methods, and alignment techniques that were barely conceived when early scaling laws were formulated [9].

*b) Emergent Capabilities and Phase Transitions:* Perhaps the most fascinating aspect of scaling behavior that has emerged is the phenomenon of threshold effects, where capabilities like multi-step reasoning, code generation, tool use, and sophisticated mathematical problem-solving appear suddenly at specific model scales [18]. The 2024-2025 landscape has revealed even more dramatic examples: the emergence of agentic behavior in models like Claude-3.5 and GPT-4o, the breakthrough mathematical reasoning capabilities of OpenAI's

o1 model, and the sophisticated multimodal understanding demonstrated by systems like Gemini Ultra.

These threshold effects challenge simple power-law formulations and suggest phase transitions in model capabilities that require more sophisticated theoretical frameworks. Current research indicates that these emergent capabilities often correlate with specific parameter counts (around 70B, 175B, and 500B+ parameters), but the precise mechanisms remain partially mysterious. Understanding these transitions has become crucial for organizations planning model development timelines and capability expectations.

*c) The Critical Role of Data Quality and Curation:* Data quality has emerged as perhaps the most transformative insight in scaling law understanding. The development of sophisticated data curation frameworks—from constitutional AI training to LLM-assisted data filtering—has shown that high-quality data can yield performance gains equivalent to order-of-magnitude increases in model size [9], [12]. The 2024-2025 period has witnessed the emergence of training approaches where careful data curation and constitutional training methods achieve superior results at dramatically lower computational costs.

Organizations like Anthropic have demonstrated that constitutional AI training approaches can achieve frontier-level capabilities with substantially smaller models when combined with sophisticated data quality techniques. This insight has profound implications for democratizing AI development, as it suggests pathways to high-capability models that don't require the enormous computational resources traditionally associated with frontier performance.

*d) Evaluation Methodology Evolution:* The maturation of evaluation methodologies represents another crucial development in our understanding of scaling laws. The limitations of traditional metrics like perplexity have driven the development of comprehensive evaluation frameworks that assess models across multiple dimensions: capability, safety, alignment, robustness, and efficiency [10], [14]. The 2024-2025 evaluation landscape features sophisticated frameworks like HELM, Azure AI evaluation systems, and Constitutional AI assessment tools that provide holistic views of model performance.

Current evaluation approaches recognize that scaling laws must account for multiple performance dimensions simultaneously. A model's scaling behavior for mathematical reasoning may differ dramatically from its scaling patterns for creative writing or ethical reasoning. This multidimensional perspective has become essential for practical applications where models must perform reliably across diverse tasks and contexts.

*e) Economic Reality and Practical Constraints:* The economic dimension of scaling laws has evolved from a secondary consideration to a primary driver of development decisions. The dramatic escalation in training costs—from GPT-3's \$4.6 million to the \$100-200 million required for current frontier models—has forced a reckoning with the economic sustainability of pure scaling approaches [2], [11]. The relationship between scale, performance, and cost creates complex trade-offs that organizations must navigate carefully.

The emergence of sophisticated cost optimization techniques has shown that the economic equation extends far beyond raw training costs. Inference optimization techniques can reduce deployment costs by 10-50x, human resources often exceed infrastructure costs by 3-5x, and total cost of ownership frameworks reveal hidden expenses that can dwarf obvious computational investments [17]. Understanding these economic scaling relationships has become as important as understanding performance scaling for practical AI development.

## B. Implications for Research and Practice

These insights fundamentally reshape both research directions and practical applications of generative AI, creating new opportunities while highlighting persistent challenges.

*1) For Researchers: Expanding Frontiers:* The current state of scaling laws research points to several transformative directions for future investigation:

*a) Theoretical Foundation Development:* The need for robust theoretical explanations of observed scaling patterns has become urgent, particularly for emergent capabilities and threshold effects. Current mathematical frameworks inadequately explain why specific capabilities emerge at particular scales, or why some scaling patterns exhibit smooth power-law behavior while others show sharp transitions. Developing unified theories that encompass both gradual improvements and sudden capability jumps represents a crucial research frontier.

*b) Architecture Innovation Beyond Scale:* Creating model architectures that improve parameter efficiency or exhibit more favorable scaling properties has become essential for sustainable AI development. The success of Mixture-of-Experts architectures, which achieve large-model performance with sparse activation patterns, suggests that architectural innovation may provide more efficient paths to capability than pure parameter scaling. Research into neuromorphic architectures, retrieval-augmented systems, and hybrid computational approaches offers promising alternatives to traditional scaling strategies.

*c) Data-Centric Intelligence:* Exploring how data characteristics influence scaling behavior represents perhaps the most promising research direction. Understanding why certain data examples provide disproportionate learning value, developing methods to identify and leverage high-value training instances, and creating frameworks for optimal data allocation across different capability domains could revolutionize training efficiency. Constitutional AI and preference-based training represent early examples of this data-centric approach.

*d) Evaluation Framework Evolution:* Designing evaluation methodologies that remain discriminative as models approach human-level performance presents ongoing challenges. Current benchmarks saturate quickly as models improve, creating a need for adaptive evaluation frameworks that can assess increasingly sophisticated capabilities. Developing evaluation approaches that measure alignment, safety, robustness, and real-world utility—rather than just capability—has become crucial for responsible AI advancement [13], [14].

*e) Interdisciplinary Connections:* Drawing connections between scaling laws in AI and similar patterns in complex systems across physics, biology, and economics could yield transformative insights. Phase transitions in neural networks may share mathematical foundations with critical phenomena in statistical mechanics, while emergent intelligence might parallel self-organization principles observed in biological and social systems.

*2) For Practitioners: Strategic Applications:* Organizations developing or deploying generative AI can leverage scaling laws insights across multiple dimensions of their operations:

*a) Strategic Resource Allocation:* Using scaling laws to predict performance improvements from investments in model size, training data, or computational resources enables informed budget decisions and development roadmaps. Organizations can model trade-offs between training larger models, investing in data quality, or developing optimization techniques to find optimal resource allocation strategies.

*b) Architecture and Scale Selection:* Choosing model architectures and sizes appropriate for specific applications based on performance requirements, resource constraints, and deployment contexts has become a sophisticated optimization problem. Understanding how different capabilities scale allows organizations to right-size models for particular use cases, avoiding over-engineering while ensuring adequate performance.

*c) Data Strategy Optimization:* Investing in data quality and curation rather than focusing exclusively on quantity has proven particularly valuable for domain-specific applications. Organizations can achieve significant performance improvements through strategic data investments that cost substantially less than computational scaling, especially when combined with constitutional training approaches.

*d) Deployment and Optimization:* Applying advanced techniques like quantization, distillation, speculative decoding, and intelligent caching to optimize inference costs while maintaining acceptable performance enables sustainable large-scale deployment. Understanding how optimization techniques interact with scaling relationships allows organizations to achieve optimal cost-performance trade-offs.

*e) Economic Planning and Sustainability:* Developing comprehensive total cost of ownership models that account for training costs, inference expenses, human resources, and operational overhead across the full lifecycle of AI systems enables sustainable business models. Organizations must balance capability aspirations with economic realities to create viable long-term AI strategies.

## C. Challenges and Open Questions

Despite remarkable progress in understanding scaling laws, fundamental challenges and open questions continue to shape the field's trajectory:

*1) Scaling Limits and Fundamental Boundaries:* Whether performance improvements will continue indefinitely with scale or eventually reach fundamental limits remains one of the most consequential uncertainties in AI development. Current evidence suggests different capabilities may have different scaling limits, with some showing continued improvement

while others plateau. Identifying these boundaries—if they exist—would profoundly impact AI development strategies and resource allocation decisions.

*2) Emergent Capability Prediction:* Developing reliable methods to predict when and how new capabilities will emerge at scale thresholds represents a critical challenge for organizational planning [18]. Current approaches can identify scaling patterns for known capabilities but struggle to predict entirely new emergent behaviors. Organizations need frameworks for anticipating capability jumps to plan effectively for technological transitions and competitive dynamics.

*3) Data Exhaustion and Quality Constraints:* As models grow larger and training datasets expand, finding sufficient high-quality training data becomes increasingly challenging. Estimates suggest that available high-quality text data may limit scaling within the next 3-5 years, forcing the development of synthetic data generation, data augmentation techniques, or entirely new training paradigms. Understanding how data limitations might constrain future scaling represents a critical research priority.

*4) Evaluation Evolution and Benchmark Saturation:* Creating evaluation methodologies that remain meaningful as models approach human-level performance on many tasks presents ongoing challenges [13], [14]. Current benchmarks saturate rapidly, and developing new assessment frameworks that can discriminate between increasingly sophisticated systems requires continuous innovation in evaluation methodology.

*5) Economic Sustainability and Access:* Determining whether the economics of ever-larger models are sustainable, particularly as training costs escalate exponentially while inference costs accumulate over deployment lifetimes, will shape the future trajectory of AI development [11]. The concentration of advanced AI capabilities among well-resourced organizations raises important questions about technological accessibility and democratic participation in AI advancement.

*6) Environmental Impact and Responsibility:* Addressing the growing energy consumption and carbon footprint of large-scale AI training and deployment has become essential for responsible technological advancement. Developing frameworks that balance capability improvements with environmental sustainability will influence scaling strategies and regulatory approaches.

*7) Alignment and Safety at Scale:* Ensuring that scaled AI systems remain aligned with human values and controllable as they become more capable represents perhaps the most important challenge facing the field. Current alignment techniques may not scale effectively to superintelligent systems, requiring fundamental advances in AI safety and control methodologies.

### D. Future Directions and Emerging Paradigms

Looking toward 2025-2030, several transformative trends may reshape our understanding and application of scaling laws:

*1) Hybrid Scaling Architectures:* Combining traditional parameter scaling with retrieval-augmented generation, tool use capabilities, and specialized computational components offers more efficient paths to enhanced capabilities. These hybrid approaches may achieve large-model performance with substantially lower computational requirements by leveraging external knowledge sources and specialized processing modules.

*2) Multimodal Scaling Integration:* Extending scaling law understanding to multimodal models that seamlessly integrate text, images, audio, video, and other modalities represents a crucial frontier. Early evidence suggests multimodal capabilities may exhibit different scaling patterns than unimodal systems, with complex interactions between different modality-specific parameters.

*3) Personalization and Adaptation at Scale:* Developing efficient methods to adapt large foundation models to individual users, specific domains, or particular contexts without full retraining could dramatically enhance the utility of scaled models. Techniques like parameter-efficient fine-tuning, constitutional training, and preference learning may enable massive personalization while maintaining the benefits of large-scale pretraining.

*4) Continual Learning and Dynamic Scaling:* Moving beyond static pretraining toward models that continuously update with new information could fundamentally change how we conceptualize scaling over time rather than just at training time. Dynamic scaling approaches might enable models to grow and adapt continuously rather than requiring periodic complete retraining.

*5) Neuromorphic and Alternative Architectures:* Brain-inspired computing approaches and quantum computational methods might eventually offer alternative scaling patterns with more favorable efficiency characteristics. These alternative paradigms could unlock new scaling relationships that bypass current computational and energy constraints.

*6) Regulatory and Governance Integration:* Emerging regulations around AI development, deployment, data usage, and environmental impact will increasingly influence scaling strategies. Organizations must anticipate how governance frameworks will shape acceptable scaling approaches and plan accordingly.

### E. Concluding Thoughts: Toward Sustainable AI Advancement

Scaling laws have emerged as the foundational framework for understanding and predicting the behavior of generative AI systems. The relationships between model size, training data quality and quantity, performance across multiple dimensions, and comprehensive economic considerations provide crucial guidance for researchers and practitioners navigating the rapidly evolving landscape of AI capabilities.

As we continue refining our understanding of these relationships, we gain increasingly sophisticated tools for making informed decisions about resource allocation, architecture design, deployment strategies, and sustainable development approaches. The evolution from simple power-law observations to multidimensional optimization frameworks reflects the maturation of AI as both a scientific discipline and a transformative technology.

The path forward requires balancing the pursuit of enhanced capabilities through scale with critical considerations of efficiency, accessibility, sustainability, and alignment with human values. The insights from scaling laws research provide essential guideposts for this journey, helping us navigate complex trade-offs between performance aspirations and practical constraints.

Looking toward the future, scaling laws research will likely evolve toward increasingly sophisticated frameworks that account for multiple dimensions of scale, quality, efficiency, and alignment simultaneously. Understanding these relationships will become even more crucial as AI systems approach and potentially exceed human-level performance across many domains.

The principles emerging from our comprehensive analysis suggest several key imperatives for sustainable AI advancement: investing in data quality and constitutional training approaches, developing economically viable optimization strategies, creating evaluation frameworks that assess alignment and safety alongside capability, and maintaining accessibility to ensure broad participation in AI development.

In this dynamic and consequential field, maintaining a holistic perspective that integrates technical innovation, economic sustainability, ethical considerations, and human values will be essential. The scaling laws framework provides valuable tools for this integration, offering quantitative approaches to complex qualitative challenges.

As we stand at the threshold of potentially transformative advances in artificial intelligence, the insights from scaling laws research remind us that the most important scaling challenge may not be technical or economic, but rather ensuring that our most powerful AI systems remain beneficial, controllable, and aligned with human flourishing. The relationships we have explored between scale and capability must ultimately serve the broader goal of developing AI systems that enhance rather than replace human intelligence, augment rather than automate human creativity, and amplify rather than diminish human agency.

The future of scaling laws—and of generative AI itself—will be determined not just by our technical ingenuity or computational resources, but by our wisdom in applying these powerful tools responsibly. In this endeavor, the comprehensive understanding of scaling relationships provides both the foundation for continued progress and the framework for ensuring that progress serves humanity's highest aspirations.

## REFERENCES

[1] Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., et al. (2023). Palm 2 technical report. arXiv preprint arXiv:2305.10403.

[2] Appenzeller, G., Bornstein, M., & Casado, M. (2023). Navigating the High Cost of AI Compute. Andreessen Horowitz.

[3] Bamoria, H. (2025). Top 10 LLM Benchmarking Evals. Medium.

[4] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., et al. (2020). Language models are few-shot learners. Advances in Neural Information Processing Systems, 33, 1877-1901.

[5] Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., et al. (2022). Palm: Scaling language modeling with pathways. arXiv preprint arXiv:2204.02311.

[6] Hestness, J., Narang, S., Ardalani, N., Diamos, G., Jun, H., Kianinejad, H., et al. (2017). Deep learning scaling is predictable, empirically. arXiv preprint arXiv:1712.00409.

[7] Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., et al. (2022). Training compute-optimal large language models. arXiv preprint arXiv:2203.15556.

[8] Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., et al. (2020). Scaling laws for neural language models. arXiv preprint arXiv:2001.08361.

[9] Longpre, S., Hou, L., Vu, T., Webson, A., Chung, H. W., Tay, Y., et al. (2023). The flan collection: Designing data and methods for effective instruction tuning. arXiv preprint arXiv:2301.13688.

[10] Microsoft. (2025). Evaluation and monitoring metrics for generative AI. Microsoft Learn.

[11] Schmid, P. (2024). Understanding the Cost of Generative AI Models in Production.

[12] Sorscher, B., Geirhos, R., Shekhar, S., Ganguli, S., & Morcos, A. S. (2022). Beyond neural scaling laws: beating power law scaling via data pruning. Advances in Neural Information Processing Systems, 35, 15813-15827.

[13] Srivastava, A., Rastogi, A., Rao, A., Shoeb, A. A. M., Abid, A., Fisch, A., et al. (2022). Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. arXiv preprint arXiv:2206.04615.

[14] Stanford CRFM. (2025). Holistic Evaluation of Language Models (HELM).

[15] Tay, Y., Dehghani, M., Abnar, S., Shen, Y., Bahri, D., Pham, P., et al. (2022). Scale efficiently: Insights from pre-training and fine-tuning transformers. arXiv preprint arXiv:2109.10686.

[16] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M. A., Lacroix, T., et al. (2023). Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.

[17] Wang, B., & Richards, R. (2025). Optimizing Cost for Generative AI with AWS. AWS Cloud Financial Management.

[18] Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., et al. (2022). Emergent abilities of large language models. arXiv preprint arXiv:2206.07682.