

Scaling Large Language Models: Performance, Infrastructure, and the Path to Artificial General Intelligence

1. Introduction: The Significance of LLM Scaling

Large Language Models (LLMs) represent a significant advancement in artificial intelligence, demonstrating remarkable capabilities in understanding and generating human-like text ¹. These sophisticated deep learning models, primarily built on the transformer architecture, have revolutionized various natural language processing tasks, including text generation, translation, question answering, and code generation ¹. Their increasing integration into diverse real-world applications, such as customer service chatbots, content creation platforms, and software development tools, underscores their growing importance and impact across industries ¹.

The concept of scaling in the context of AI refers to the practice of increasing the fundamental resources involved in developing and training these models, namely, the size of the model (quantified by the number of parameters), the volume and diversity of the training data, and the computational resources dedicated to the training process ³. The central premise driving this approach is that augmenting these resources leads to improved model performance and the potential emergence of novel and more advanced capabilities ³. These improvements are often empirically observed and described by scaling laws, which provide a framework for predicting how increases in scale correlate with enhancements in model performance, typically following a power-law relationship ³.

This report aims to provide a comprehensive analysis of the scaling phenomenon in LLMs. It will delve into the foundational principles governing these scaling laws, explore the practical implications of scaling for model performance and resource demands, examine the role of cutting-edge hardware in supporting these advancements, discuss strategies for efficient inference and token output management, and ultimately consider the potential trajectory towards achieving Artificial General Intelligence (AGI) through continued scaling efforts. The subsequent sections will systematically address these aspects, providing a detailed and insightful perspective on the current state and future directions of LLM scaling.

The remarkable progress observed in LLMs is fundamentally linked to the ability to scale them effectively. These models, built upon deep learning architectures, inherently benefit from access to larger datasets and the capacity offered by larger models to learn intricate patterns within that data. The groundbreaking performance of models like GPT-3 and its successors has served as a powerful validation of the scaling hypothesis. Furthermore, scaling is not simply about increasing the magnitude of individual resources; it necessitates a careful and intricate balance between the number of model parameters, the quality and quantity of the training data, and the computational resources employed. Research, such as the Chinchilla scaling law, has highlighted that optimizing the interplay between these factors, rather than solely maximizing one in isolation, is crucial for achieving optimal performance within given computational constraints ³. Understanding these scaling laws and their nuances is therefore paramount for guiding the future development and strategic allocation of resources in the pursuit of

increasingly capable LLMs ².

2. The Core Principles: Scaling Laws and Performance in LLMs

The concept of scaling laws in artificial intelligence has its roots in earlier machine learning theories and has gained significant prominence with the advent of deep learning ³. Initial investigations explored the relationship between model complexity and the ability of a model to generalize to unseen data. OpenAI's seminal work in 2020 provided substantial empirical evidence for the notion that for large deep learning models, a consistent increase in model size, dataset size, and the amount of computational resources used for training leads to a predictable reduction in model loss, following a power-law relationship ³. This research laid the groundwork for a quantitative understanding of the performance gains achievable through scaling.

The fundamental power-law relationship that often describes scaling in AI can be expressed mathematically as $y \propto x^{^b}$, where y represents a measure of model performance (or conversely, loss), x represents a scaling factor such as the number of parameters or dataset size, and b is a constant exponent ³. This relationship signifies that improvements in performance are proportional to a power of the increase in the scaling factor. A key implication of this is that to achieve increasingly significant performance gains, exponentially larger increases in the scaling factor might be required.

DeepMind's Chinchilla scaling law, introduced in 2022, further refined the understanding of optimal scaling strategies by emphasizing the crucial need to balance the number of model parameters with the size of the training dataset for a given computational budget ³. Their research demonstrated that, under comparable computational constraints, a smaller model (with 70 billion parameters) trained on a significantly larger and more carefully curated dataset could outperform a much larger model (with 175 billion parameters) that was trained on a comparatively smaller dataset ³. This finding underscored the importance of data efficiency and suggested that, within a given compute budget, prioritizing the size and quality of the training data can often be more beneficial than simply maximizing the number of model parameters.

While the initial focus of scaling law research was primarily on predicting and minimizing training loss, it has become increasingly apparent that training loss alone is not always a perfect indicator of a model's performance on specific downstream tasks or its overall utility in real-world applications ². Various factors, including the specific architecture of the model, the distribution of the training data, the tokenization process used, and the precision of the computations, can influence the actual performance achieved on tasks of interest ⁴. To address this, researchers have developed more direct methods for predicting downstream performance. For example, the concept of a "Performance Law" has been proposed as an empirical equation that can directly predict the MMLU score (a widely used metric for evaluating the general conversational and application capabilities of LLMs) based on a few key hyperparameters of the LLM architecture and the size of the training data ⁴. Similarly, the "FLP" (Flops Loss Performance) framework offers a two-stage approach. It first uses scaling laws derived from smaller "sampling" language models to estimate the relationship between the amount of computational resources (measured in FLOPs) used during pre-training and the resulting pre-training loss. Subsequently, it maps this pre-training loss to the expected performance on downstream tasks, taking into account the "emergent abilities" that often manifest in LLMs once they surpass certain computational thresholds ². An extension of this, known as "FLP-M," further

refines the prediction process by specifically addressing the practical need to integrate datasets from multiple sources during pre-training, such as blending general corpora with code data ². A significant challenge in accurately predicting downstream performance arises from the phenomenon of emergent abilities in LLMs ². These are capabilities that models only begin to exhibit once they reach a certain scale or have been trained with a sufficient amount of compute, often appearing rather abruptly and being difficult to extrapolate from the performance of smaller models.

Despite their utility, scaling laws are not without limitations and practical considerations ³. One notable aspect is the phenomenon of diminishing returns. As both the size of the model and the dataset increase, the incremental improvements in performance tend to become smaller ³. This suggests that simply doubling the model size or the dataset might not always yield a doubling of the performance gain. Furthermore, the quality of the training data is of paramount importance for realizing the full benefits of scaling ¹. Datasets that are noisy, biased, or of low quality can lead to suboptimal model performance, regardless of the scale. Therefore, ensuring that the training data is diverse, relevant, and of high quality is crucial for maximizing the positive impact of scaling. Finally, the concept of a compute-optimal balance among model size, dataset size, and the amount of compute used is essential for efficient scaling ³. Rather than simply maximizing one of these factors in isolation, the most effective approach often involves finding a proportional adjustment of all three to achieve the best performance within a given computational budget.

The evolution of scaling law research reflects a growing understanding that while increasing the scale of LLMs is a crucial driver of progress, it is not a monolithic endeavor. The field is moving towards a more nuanced perspective that considers the interplay of various factors and aims to optimize not just for theoretical loss reduction but for tangible improvements in real-world task performance. The development of tools and frameworks that can predict this downstream performance is vital for guiding the efficient allocation of resources and accelerating the development of increasingly capable language models.

3. The Impact of Model Size: Parameters and Capabilities

The size of a Large Language Model (LLM) is primarily determined by the number of parameters it contains ⁸. These parameters are essentially the adjustable internal settings within the model, optimized during the training process to enable the model to accurately predict the next token in a sequence of text ⁸. These parameters, which can be thought of as the weights and biases in the underlying neural network, are where the model stores the knowledge it learns from the vast amounts of textual data it is trained on ⁷.

To illustrate the concept of parameters, one can consider a simplified model for estimating house prices ⁸. In such a model, inputs like the square footage of the house and the number of bedrooms might be multiplied by certain parameters, and these results are then combined to produce an estimated price. The parameters in this case are the numerical weights assigned to each input factor. LLMs operate on a similar principle, but with billions of such adjustable parameters and a far more intricate underlying architecture ⁸. During the training phase, an LLM is exposed to massive quantities of text data. As it processes this data, it makes predictions about the subsequent words in the sequence. It then compares these predictions to the actual text and, through an iterative process, adjusts its internal parameters to minimize the

discrepancy between its predictions and the ground truth, thereby gradually enhancing its predictive accuracy ⁹.

Generally, an LLM with a greater number of parameters possesses a larger capacity to learn and represent the complex relationships that exist within human language ⁷. This increased capacity often translates to superior performance across a wide spectrum of natural language processing tasks, encompassing improved language understanding, more fluent and coherent text generation, and enhanced reasoning capabilities ⁷. For example, larger models tend to be better at comprehending the subtle nuances of language, including idioms, metaphors, and intricate semantic connections between words ¹¹. The advanced reasoning capabilities exhibited by models like GPT-4 are largely attributed to their significantly larger parameter counts compared to their predecessors ¹¹. Furthermore, larger models are typically capable of handling larger context windows, meaning they can effectively "remember" and consider more of the preceding text in a conversation or a document when generating a response ¹¹. This ability to maintain context over longer sequences leads to more coherent and contextually relevant outputs, whereas smaller models might struggle to retain information from earlier parts of a lengthy input ¹¹.

However, it is a common misconception that simply increasing the number of parameters will invariably lead to a better-performing LLM ⁷. While the parameter count is undoubtedly a critical factor influencing a model's potential, it is not the sole determinant of its success. The quality, richness, and relevance of the data used to train the model play an equally, if not more, vital role in shaping its capabilities ⁷. Indeed, a smaller language model that has been trained on a high-quality and well-curated dataset can often outperform a larger model that has been trained on a dataset of lower quality or relevance ⁹. The success of models like Mistral 7B and LLaMA 3, which have demonstrated performance levels comparable to much larger models despite having significantly fewer parameters, underscores the importance of training data and the techniques used to train the models ⁷. These examples suggest that architectural innovations and the strategic application of training methodologies are just as crucial as sheer parameter size ³.

There are also significant trade-offs associated with developing and deploying LLMs with very large numbers of parameters ⁷. Training these models demands substantial computational resources, including access to powerful GPUs, large amounts of memory, and considerable training time ⁷. Similarly, running these large models for inference (generating responses) also requires significant computational power and can lead to longer processing times, making them expensive to operate and potentially limiting their accessibility to individuals and organizations with constrained resources ⁹. Moreover, the sheer size of these models can pose challenges for deployment in real-world applications, often necessitating specialized infrastructure that may not be universally available ⁷. Finally, the considerable energy consumption associated with both training and running very large LLMs raises important concerns about their environmental impact ⁷.

To provide a practical way to gauge the potential performance of LLMs, the AI community often categorizes them into approximate size classes based on their number of parameters, such as 2B (2 billion), 7B, 70B, 175B, and so on ¹². These size classes offer a general indication of a model's capacity and expected capabilities, with larger models typically exhibiting more advanced performance. However, it is important to note that these classifications are not rigid

industry standards, and the actual performance of models within the same size class can vary significantly depending on other factors, such as the specific architecture of the model and the data it was trained on ¹². For instance, two different LLMs, both having 7 billion parameters, might demonstrate different levels of accuracy, fluency, and reasoning ability based on these other factors. Resources like the HuggingFace Open LLM Leaderboard provide a more granular and comparative evaluation of LLM performance across a range of benchmarks, taking into account various factors beyond just the number of parameters ¹².

While a larger number of parameters in an LLM generally provides a greater capacity for learning and can lead to improved performance across various language tasks, it is crucial to recognize that parameter count is just one piece of a larger puzzle. The quality and characteristics of the training data, the efficiency of the model's architecture, and the sophistication of the training techniques employed all play equally vital roles in determining the overall effectiveness and capabilities of the model. Furthermore, the benefits of increased parameter size must be weighed against the associated costs in terms of computational resources, accessibility, and environmental impact. The trend towards developing smaller, more efficient models that can achieve comparable performance to their larger counterparts highlights the ongoing efforts to optimize this balance and democratize access to advanced language technologies.

4. Boosting Performance at Inference: Test-Time Scaling Explained

Inference scaling, also referred to as test-time scaling, represents an innovative approach in the field of language modeling that aims to enhance the performance of pre-trained Large Language Models (LLMs) by strategically employing additional computational resources during the inference phase, which is when the model is generating responses to prompts ¹³. Unlike traditional methods that focus on improving performance through more extensive training or by increasing the size of the model, inference scaling achieves gains without requiring any further training or modifications to the model's underlying parameters ¹⁴.

The fundamental principle behind inference scaling is that by allowing the LLM to perform more computations or to explore a wider range of possibilities at the moment it is generating a response, its reasoning abilities and the overall quality of its output can be significantly improved ¹⁴. This concept is analogous to providing a human with more time to deliberate and think through a complex problem before arriving at a solution. This paradigm marks a notable shift from the conventional emphasis on scaling up the resources used during training to achieve better performance, suggesting that substantial improvements can also be realized by optimizing the process through which the model generates its responses ¹⁴.

Several mechanisms are employed to implement test-time scaling. One common technique involves **generating multiple potential answers (often referred to as "Best of N" sampling)** for a given prompt and then selecting the answer that is deemed most plausible or of the highest quality based on certain criteria ¹³. This allows the model to explore different avenues of response and choose the one that appears most coherent, accurate, or aligned with the desired characteristics. A variation of this is the **"Best of N Weighted"** approach, where identical responses that are generated multiple times are given more weight in the selection process ¹⁷.

Chain-of-thought (CoT) reasoning, a technique where the model is prompted to explicitly

articulate its reasoning steps leading to an answer, can be significantly amplified through inference scaling ¹³. By increasing the computational budget available during the inference stage, the model can explore a more extensive set of potential reasoning pathways, potentially uncovering more accurate and robust solutions to complex problems.

Process Reward Models (PRMs) represent another approach used in conjunction with inference scaling ¹³. A separate, smaller model trained to evaluate the quality of reasoning steps can be used to guide the LLM's generation process. The PRM assesses the intermediate steps of reasoning produced by the LLM and provides feedback, enabling the model to identify and rectify potential errors or unproductive lines of thought as it progresses towards a final answer.

Iterative self-refinement is a method where the LLM generates an initial response to a prompt and then iteratively revises and improves upon its own output in a step-by-step manner ¹⁷. The model uses its previously generated text as additional context to refine subsequent parts of the response, leading to a potentially more polished and coherent final output.

For more intricate problems, **parallel sampling** or **tree-search methods** can be utilized ¹⁷. These techniques allow the model to explore multiple potential answers or reasoning strategies concurrently. This is particularly beneficial when the initial response generated by the model might not be on the right track, as it enables the exploration of a broader range of possibilities. **Beam search** is one such method that generates a set of the most likely next tokens at each step, effectively exploring multiple potential output sequences simultaneously ¹⁷.

Finally, **"budget forcing"** is a simpler yet effective technique for controlling the amount of computational resources used during inference ¹⁹. This method involves either forcefully terminating the model's thought process if it exceeds a predefined limit of generated tokens or encouraging it to continue thinking for longer by appending special tokens like "Wait" to its current reasoning trace. This can prompt the model to double-check its work or explore alternative approaches.

Research has demonstrated that the application of inference scaling techniques can lead to substantial improvements in performance on standard benchmark metrics, often exceeding 20 percentage points in accuracy ¹³. Notably, smaller LLMs that employ inference scaling have been shown to outperform much larger models that do not utilize these methods ¹³. For instance, IBM's 8 billion parameter Granite 3.2 model, when using inference scaling techniques, has achieved state-of-the-art results on math reasoning tasks, even surpassing the performance of much larger proprietary models like GPT-4o and Claude 3.5 Sonnet ¹³. Similarly, studies have found that a relatively small PaLM 2-S model, when augmented with test-time compute, can outperform a model that is 14 times larger ¹⁷, and Llama-3.2 with just 3 billion parameters, when subjected to 256 iterations of test-time compute, has been shown to achieve better results than Llama-3.1 with 70 billion parameters ¹⁷. In some scenarios, the performance gains realized through inference scaling can be comparable to or even greater than those achieved by further increasing the size of the model or the amount of compute used during pre-training, suggesting that it can be a more efficient pathway to enhancing model capabilities for certain types of tasks ¹⁵. This is particularly relevant for organizations that may have limitations on the resources they can dedicate to large-scale pre-training runs ¹⁶.

Researchers are also actively exploring **compute-optimal scaling strategies** for inference,

which involve adaptively allocating the amount of computational resources used at test time based on the perceived difficulty of the prompt or the task at hand ¹⁵. The underlying idea is to expend more computational effort on challenging prompts that are likely to require more complex reasoning, while using fewer resources for simpler prompts where the model can arrive at a satisfactory answer more readily. These adaptive strategies aim to improve the overall efficiency of inference scaling compared to approaches that use a fixed amount of compute for every prompt, and studies have indicated that they can lead to significant improvements in efficiency, such as a fourfold increase in performance for math reasoning problems ¹⁵.

Inference scaling is increasingly being recognized as a crucial and complementary approach to the more traditional methods of scaling AI models through pre-training (increasing data, model size, and compute for the initial training phase) and post-training (using techniques like reinforcement learning and human feedback to further refine the model's capabilities) ¹⁶. As the performance gains from simply scaling up pre-training resources begin to show signs of plateauing or become prohibitively expensive, test-time scaling offers a promising new direction for enhancing the capabilities of LLMs ¹⁶. It potentially allows developers to focus on training smaller, more efficient base models that possess strong reasoning cores, and then selectively scale up the computational resources used during inference for specific tasks that are economically valuable or particularly demanding ¹⁶. This could lead to a more sustainable and cost-effective approach to building and deploying high-performing AI systems.

The emergence of test-time scaling signifies a fundamental shift in how computational resources are leveraged in the development of advanced AI. It suggests that intelligence and the ability to solve complex problems effectively can arise not only from the vast amount of knowledge encoded within a model's parameters during training but also from the capacity to reason and explore potential solutions more thoroughly at the moment of inference. This approach holds the potential to democratize access to high-performing AI, as it might enable smaller organizations with limited training budgets to achieve state-of-the-art results by strategically investing in inference-time computation. Furthermore, the development of more sophisticated inference scaling techniques could represent a significant step towards achieving artificial general intelligence by providing a mechanism for models to dynamically improve their problem-solving abilities and potentially overcome some of the inherent limitations of static pre-trained knowledge.

5. The Hardware Foundation: Latest GPUs for LLM Workloads

Graphics Processing Units (GPUs), initially conceived for rendering intricate three-dimensional graphics, have become indispensable for modern artificial intelligence, particularly for the computationally intensive tasks involved in both training and deploying Large Language Models (LLMs) ²¹. The architecture of LLMs, which involves massive neural networks with billions of parameters, relies heavily on matrix operations. GPUs are exceptionally well-suited for these types of parallel computations, allowing them to perform the necessary calculations at a significantly faster rate than traditional Central Processing Units (CPUs) that process data sequentially ²¹. This inherent parallel processing capability of GPUs is a critical factor in accelerating both the extensive training processes and the real-time inference speeds required for LLMs.

Several key specifications of a GPU are particularly important for handling LLM workloads ²¹.

GPU memory (VRAM) is paramount, as the entire LLM model, along with the intermediate data generated during computations, needs to reside within this memory during both training and inference. Larger LLMs with more parameters necessitate greater VRAM capacity. Insufficient VRAM can lead to performance bottlenecks, such as "Out of Memory" errors or a substantial slowdown in processing as data must be constantly transferred between the GPU and the system's main memory²⁵. **Memory bandwidth** is another crucial factor, referring to the speed at which data can be moved between the GPU's memory and its processing cores. For LLM inference, where models are often memory-bound, high memory bandwidth is essential for quickly loading model parameters and processing tokens, thus reducing latency²¹. Finally, **compute power**, often measured in TeraFLOPS (TFLOPS), indicates the raw processing capability of the GPU. For LLMs, the performance in FP16 (half-precision floating-point) operations is particularly relevant as many models are trained and run using this format to achieve a balance between speed and accuracy²¹. NVIDIA GPUs also incorporate **Tensor Cores**, specialized units designed to accelerate matrix multiplications, which are fundamental to the computations in deep learning²¹. The number and the generation of these Tensor Cores significantly impact the overall performance of LLMs on NVIDIA hardware.

Among the latest NVIDIA GPUs, several stand out for their capabilities in handling LLM workloads²¹. The **NVIDIA H100** is a top-tier data center GPU that offers exceptional performance for training and inferencing very large LLMs. It boasts a substantial memory capacity (up to 80GB of high-bandwidth HBM3 memory), extremely high memory bandwidth (up to 3 terabytes per second), and powerful fourth-generation Tensor Cores²¹. However, its high cost and power consumption necessitate specialized infrastructure²¹. The **NVIDIA A100**, a previous generation high-performance GPU, remains a popular choice in data centers and cloud platforms. It offers large memory options (40GB or 80GB of HBM2e memory) and high memory bandwidth (up to 2 terabytes per second), along with third-generation Tensor Cores²¹. While it is not as powerful as the H100, it provides a strong balance of performance and cost-effectiveness. Interestingly, the consumer-grade **NVIDIA RTX 4090** has emerged as a surprisingly capable option for LLM inference, particularly for smaller to medium-sized models. It features a significant number of CUDA and fourth-generation Tensor Cores, along with a respectable amount of memory (24GB of GDDR6X) and high memory bandwidth (1 terabyte per second)²¹. Its excellent performance relative to its price makes it a popular choice for development and smaller-scale deployments, even outperforming some enterprise-grade GPUs in cost-adjusted benchmarks³³. The **NVIDIA L40S** is a professional-grade GPU that offers a good balance of compute power and memory (48GB of GDDR6) with a memory bandwidth of 864 gigabytes per second. It includes fourth-generation Tensor Cores and is more energy-efficient than the H100 or A100, making it well-suited for various LLM inference tasks²¹. Benchmarks using tools like llama.cpp and NVIDIA's TensorRT-LLM demonstrate the varying performance levels (in tokens per second) achieved by these different GPUs when running various LLM models, often highlighting the strong performance of the RTX 4090 for models that fit within its memory constraints²⁵.

AMD has also become a significant player in the AI hardware landscape, with GPUs like the **AMD MI300X** presenting a strong challenge to NVIDIA's dominance, especially in the realm of LLM inference²⁴. The MI300X has been shown to outperform NVIDIA's H100 in certain LLM inference benchmarks due to its larger memory capacity (192 GB compared to 80/94 GB) and higher memory bandwidth (5.3 terabytes per second versus 3.3–3.9 terabytes per second)²⁴. This makes the MI300X particularly well-suited for handling very large LLMs on a single GPU,

potentially eliminating the need for multi-GPU setups in some cases. Benchmarks have indicated that the MI300X can nearly double the request throughput and significantly reduce latency compared to the H100 when running models like Llama 3 70B and Mixtral ²⁴.

Estimating the minimum GPU memory required for LLM inference can be done using a simplified formula: $M = P * Z * 1.2$, where M is the GPU memory in gigabytes, P is the model size in billions of parameters, and Z is the quantization factor in bytes per parameter (e.g., 2 for FP16) ³⁵. The 1.2 factor accounts for additional overhead. For instance, a 70 billion parameter model using FP16 precision would require approximately $70 * 2 * 1.2 = 168$ GB of GPU memory. When the memory capacity of a single GPU is insufficient to accommodate a large LLM, multi-GPU setups are typically employed. Techniques such as tensor parallelism distribute the model's weights and the computational workload across multiple GPUs, effectively increasing the total available memory ²⁴. For example, running a 70 billion parameter Llama 3 model might necessitate the use of two or more high-end GPUs like the A100 or H100 ²⁴.

To further enhance inference speed and reduce the memory footprint of LLMs, various optimization techniques are employed ¹². **Quantization** involves reducing the precision of the model's weights (e.g., from 16-bit floating-point numbers to 8-bit or even 4-bit integers), which can significantly decrease memory usage and potentially increase inference speed with minimal impact on accuracy ¹². **Pruning** is a technique that aims to reduce the model's size by removing less important connections (weights) from the neural network ¹². **Knowledge distillation** involves training a smaller, more efficient model to mimic the behavior and performance of a larger, pre-trained model ¹². **Low-Rank Adaptation (LoRA)** is a fine-tuning method that freezes the weights of a pre-trained model and introduces small, low-rank matrices to adapt the model for specific downstream tasks, thereby reducing the memory overhead associated with fine-tuning ¹². **Batching** is a strategy where multiple independent input sequences are grouped together and processed simultaneously by the GPU, which can improve the overall throughput of the system and increase the utilization of the hardware, although it might slightly increase the latency for individual requests ²⁴. LLM inference typically occurs in two main stages: **prefill**, where the input prompt is processed in parallel, and **decoding**, where the output text is generated one token at a time in an autoregressive manner ³¹. Understanding the performance characteristics of these two stages is crucial for optimizing the overall inference process.

The performance of LLM inference is significantly influenced by the complex interplay between a GPU's memory capacity, its memory bandwidth, and its computational power. The specific bottleneck can vary depending on the size of the LLM being used and the batch size for inference. For smaller batch sizes, the process is often limited by memory bandwidth, meaning the speed at which the model's parameters can be loaded from memory to the processing units is the determining factor. For larger models or when using larger batch sizes, the raw computational power of the GPU might become the primary limitation. The emergence of competitive GPUs from AMD, such as the MI300X, is introducing more options and potentially driving down costs in the AI hardware market. The superior memory and bandwidth of the MI300X, for example, provide a compelling alternative to NVIDIA's high-end GPUs for certain LLM inference workloads. Ultimately, the continuous advancements in GPU technology are a critical enabler for the ongoing scaling of LLMs and the development of increasingly sophisticated AI applications. The growing availability of high-performance GPUs, combined with various optimization techniques, is making it feasible to deploy and run more powerful

models in a wider range of environments.

Table 1: Comparison of Key NVIDIA GPUs for LLM Inference

GPU	CUDA Cores	Tensor Cores (Generation)	Memory (GB)	Memory Bandwidth (TB/s)	TDP (W)
H100 (SXM5)	18,432	576 (4th)	80	3.0	700
A100 (SXM)	6,912	432 (3rd)	80	2.0	400
RTX 4090	16,384	512 (4th)	24	1.0	450
L40S	18,176	568 (4th)	48	0.864	300
A10 (PCIe)	9,216	288 (3rd)	24	0.6	150

6. Planning for Output: AI Inference Tokens and Efficiency

Tokens serve as the fundamental units of text that Large Language Models (LLMs) process ⁹. A single token can represent a complete word, a part of a word (known as a subword), or even an individual character, depending on the specific method of tokenization employed by the model ⁹. For instance, the sentence "The evil that men do lives after them" might be broken down into a sequence of numerical tokens, each corresponding to a specific word or subword ¹⁰. The precise manner in which an LLM segments text into tokens is determined by its tokenizer, which can vary from model to model ³⁷. As a general guideline for English text, one token is often considered to be approximately equivalent to 0.75 of a word ³⁵.

Inference latency is a critical metric that measures the time delay between when an LLM receives a user's prompt and when it generates a response ³¹. This is a key factor in determining the user experience, particularly for interactive applications such as chatbots, where low latency is highly desirable ³⁷. Several metrics are used to quantify latency, including the **Time To First Token (TTFT)**, which measures the delay until the very first token of the response is produced; the **Time Per Output Token (TPOT)**, which indicates the average time taken to generate each subsequent token; and the **Total Generation Time**, which is the entire duration from the initiation of the request to the completion of the response ³¹. **Throughput**, on the other hand, measures the efficiency of the inference system, typically expressed as the number of **Tokens Per Second (TPS)** that the model can generate or the number of **Requests Per Second (QPS)** that the system can handle ²⁶. A higher throughput signifies a more efficient

inference setup.

Various factors can influence the performance of LLM inference ³¹. The **size of the model** itself plays a significant role, with larger models generally exhibiting higher latency due to the increased computational demands of generating each token ³¹. The **length of both the input prompt and the desired output sequence** also affects performance; longer sequences necessitate the processing of more tokens, which can lead to increased latency ³¹. As discussed in Section 5, the **GPU hardware** used for inference has a substantial impact, with more powerful GPUs that have greater memory capacity and higher memory bandwidth typically resulting in lower latency and higher throughput ²¹. The **batch size**, which refers to the number of independent requests processed simultaneously, can improve overall throughput by more effectively utilizing the GPU's parallel processing capabilities, although it might slightly increase the latency for individual requests ²⁴. **Quantization**, a technique that reduces the precision of the model's weights, can lead to lower memory usage and faster computations, thereby improving both latency and throughput ¹². Various **inference optimization techniques**, such as caching previously computed information, can also enhance efficiency ²⁶. Finally, the **size of the context window**, which determines how much preceding text the model can consider, can also impact memory requirements and processing time, potentially affecting latency ³⁵.

The cost associated with using LLMs for inference is often determined by the number of tokens processed, including both the input prompt and the generated output ⁴⁵. Pricing is typically structured around a cost per million tokens and can vary considerably depending on the specific model being used and the provider offering the service ⁴⁵. Notably, the cost of LLM inference has been on a rapid decline due to advancements in hardware technology, model optimization strategies, and the use of quantization techniques ⁴⁶. This downward trend in cost is making LLMs a more economically viable option for an expanding range of applications. When planning to use LLMs, it is important to consider the desired length of the output, as reducing the number of generated tokens can directly lower the cost per request and also potentially improve the latency of the response ³⁹. Ultimately, when designing applications that utilize LLMs, it is crucial to strike a balance between the desired performance characteristics (in terms of latency and throughput) and the associated costs. Selecting an appropriately sized model, carefully tuning inference parameters such as batch size, and exploring optimization techniques like quantization can all contribute to achieving the necessary performance within a given budgetary constraint ³⁶.

Achieving optimal LLM inference performance necessitates a comprehensive approach that takes into account not only the characteristics of the model itself but also the capabilities of the underlying hardware, the application of various software optimizations, and the specific performance requirements of the intended use case. The decreasing cost of LLM inference is a significant factor that is driving the broader adoption of this technology across diverse industries. As the price per token continues to decline, the economic feasibility of using LLMs for a wide array of tasks, from content creation to customer support, increases substantially. This trend is likely to continue to fuel innovation and integration of LLMs into more and more products and services, ultimately transforming how we interact with technology.

7. The AGI Horizon: Projecting Achievement Through Scaling

Artificial General Intelligence (AGI) is a theoretical form of artificial intelligence that would

possess the ability to understand, learn, and apply knowledge across a broad range of tasks, at a level comparable to or exceeding that of human cognitive abilities. Unlike the specialized AI systems prevalent today, which are designed for specific functions, an AGI would exhibit general-purpose intelligence.

A significant hypothesis in the field of AI research posits that the continued scaling of Large Language Models (LLMs) – in terms of the number of parameters, the volume and diversity of training data, and the computational resources allocated to training – represents a potential pathway toward the realization of AGI¹⁶. The remarkable emergent abilities that have been observed in very large LLMs, such as their capacity for complex reasoning, their understanding of nuanced language, and their ability to perform tasks they were not explicitly trained on, lend support to this idea³. These unexpected capabilities suggest that scale might be a fundamental factor in unlocking more general forms of intelligence.

Current trends indicate a rapid and sustained increase in the computational resources being used to train leading-edge machine learning models, with the amount of compute growing exponentially over the past decade⁴⁷. This trend reflects a continued belief within the research community in the power of scaling. However, there is also a growing awareness that relying solely on scaling up pre-training resources might eventually encounter diminishing returns or become economically unfeasible. This has led to an increased focus on post-training techniques, such as reinforcement learning, and, more recently, on test-time scaling as potential methods to further enhance the capabilities of AI systems¹⁶. The emergence of test-time scaling, where AI models dynamically allocate computational resources during the inference phase to engage in more extensive "thinking" and explore a wider range of potential solutions, is viewed by some as a promising step towards achieving truly general intelligence¹⁴. This approach suggests that intelligence might arise not only from the vast amount of knowledge encoded within a model's parameters but also from its ability to reason and solve problems effectively by leveraging available computational resources when they are needed.

Estimating the number of parameters that might be required to achieve AGI is a highly speculative endeavor and remains a topic of ongoing debate within the AI community. Some researchers hypothesize that extremely large models, potentially containing trillions or even quadrillions of parameters, could be necessary. However, the development of test-time scaling techniques offers an alternative perspective. If AI models can significantly enhance their reasoning and problem-solving abilities through more sophisticated inference processes, it is conceivable that AGI could be achieved with models that have a more moderate number of parameters but are capable of more intensive computation and exploration at test time¹⁵. Furthermore, advancements in the efficiency of algorithms, such as the development of more effective model architectures and training methodologies, could also play a crucial role in reaching AGI with potentially lower parameter counts or reduced computational demands³. The observed trend of a decreasing amount of physical compute needed to achieve a certain level of performance in language models, due to algorithmic improvements, is an encouraging sign in this regard⁴⁷.

The realization of AGI will likely necessitate unprecedented levels of computational power and energy efficiency²². While current GPU technology is incredibly powerful, it might face limitations in meeting the comprehensive demands of true general intelligence, particularly concerning scalability and energy consumption for such massive computations²³. Specialized

AI hardware, such as Tensor Processing Units (TPUs) developed by Google and custom-designed chips from companies like Cerebras Systems and Groq, are being actively explored as potential solutions for the intensive computational requirements of AGI²². These specialized processors are often optimized for the specific types of operations that are prevalent in AI workloads, potentially offering greater efficiency and performance compared to general-purpose GPUs. The development of hardware that can efficiently support emerging AI architectures, such as Mixture of Experts (MoE) and neuro-symbolic systems, will also be crucial for making progress towards AGI²³. Additionally, achieving real-time interaction with the world, a key characteristic of general intelligence, will require hardware capable of ultra-low latency and high throughput.

While the scaling of LLMs has led to remarkable progress in artificial intelligence, the achievement of AGI remains a highly complex and long-term objective. There are numerous challenges that extend beyond simply increasing the scale of models. These include endowing AI systems with genuine understanding, consciousness, robust common-sense reasoning, and the ability to seamlessly generalize knowledge across vastly different domains and tasks. Current "reasoner models" that leverage test-time compute to achieve impressive performance on specific benchmarks, while representing a significant step forward, are not yet considered to be AGI¹⁷. These models often excel in tasks with clearly defined and objectively verifiable steps, such as mathematics and coding, but may struggle with more open-ended or subjective tasks that require a broader understanding of the world. There is an ongoing debate within the AI community regarding whether continued scaling alone will ultimately lead to AGI, or if fundamental breakthroughs in our understanding of intelligence and the development of entirely new AI paradigms will be necessary. While some researchers maintain that sufficient scaling will eventually unlock AGI, others argue that more profound conceptual advancements are required to bridge the gap.

The pursuit of AGI through scaling and other means carries profound societal and ethical implications that demand careful consideration. As AI systems become increasingly capable, it is crucial to proactively address issues related to safety, potential biases embedded in the data or models, and the broader impact on the job market and human society as a whole.

8. Conclusion: Navigating the Future of LLM Scaling

This report has provided an in-depth exploration of the critical role that scaling plays in the advancement of Large Language Models (LLMs). The analysis has shown how scaling laws offer a valuable framework for understanding the relationship between the resources invested in model development (size, data, compute) and the resulting performance. While increasing the scale of LLMs generally leads to enhanced capabilities, the quality of the training data and the efficiency of the training techniques employed are equally important considerations. The emergence of test-time scaling represents a significant advancement, offering a method to further boost the performance of LLMs during inference without the need for additional training. The latest advancements in GPU technology, particularly from industry leaders like NVIDIA and AMD, are fundamental to enabling the training and deployment of increasingly sophisticated LLMs. Careful planning for efficient token output and strategic management of inference costs are essential for the practical and sustainable application of these powerful models. Finally, this report has considered the potential of continued scaling, in conjunction with other innovations, to

contribute to the long-term and ambitious goal of achieving Artificial General Intelligence.

The optimal performance of LLMs, and ultimately the potential for achieving AGI, is not determined by any single factor related to scale. Instead, it is the result of a complex and interconnected interplay between the size of the model (number of parameters), the quality and diversity of the data used for training, the amount of computational resources that are available, and the strategies that are implemented during both the training and inference phases, including innovative techniques like test-time scaling. The key to maximizing the benefits of scaling lies in finding the most effective balance between these various elements.

The field of LLM scaling is characterized by rapid evolution, with ongoing research continually pushing the boundaries of what is possible. Future directions of research include a deeper investigation into the critical role of data quality, a more thorough exploration of the characteristics and limits of emergent abilities that arise in very large models, and the application of scaling principles to multimodal models that can process and generate information across different types of data. Furthermore, there is a persistent drive to enhance the efficiency of algorithms and to develop novel training strategies that can lead to improved performance with more judicious use of resources. The continued innovation in hardware technology, resulting in more powerful and efficient computing infrastructure, will also be essential for supporting the development and deployment of the next generation of AI models.

In conclusion, the sustained scaling of LLMs has undeniably propelled us closer to realizing more intelligent and capable artificial intelligence systems. The emergent abilities that have been observed in these large-scale models offer promising indications of the potential for even more general forms of AI. While increasing the scale alone might not be the complete solution for achieving AGI, it is highly likely to be a critical component of that complex and challenging journey. The development of innovative techniques such as test-time scaling, coupled with continued advancements in hardware and a deeper understanding of the fundamental principles of intelligence, holds significant promise for unlocking even greater AI capabilities in the years ahead. However, it is also crucial to maintain a perspective grounded in realistic expectations and to carefully consider the profound ethical and societal implications that arise from the development and deployment of such powerful technologies.

Works cited

1. Are Large Language Models Reliable? How To Improve Accuracy - Secoda, accessed March 14, 2025, <https://www.secoda.co/blog/are-large-language-models-reliable-how-to-improve-accuracy>
2. SCALING LAWS FOR PREDICTING DOWNSTREAM PERFORMANCE IN LLMS - Amazon Science, accessed March 14, 2025, <https://assets.amazon.science/9b/28/e05036cc446e90544a9b6e33a1b9/scaling-laws-for-predicting-downstream-performance-in-llms.pdf>
3. Scaling Laws in Large Language Models | HackerNoon, accessed March 14, 2025, <https://hackernoon.com/scaling-laws-in-large-language-models>
4. [2408.09895] Performance Law of Large Language Models - arXiv, accessed March 14, 2025, <https://arxiv.org/abs/2408.09895>
5. [2410.08527] Scaling Laws for Predicting Downstream Performance in LLMs - arXiv, accessed March 14, 2025, <https://arxiv.org/abs/2410.08527>

6. Scaling Laws for Predicting Downstream Performance in LLMs - OpenReview, accessed March 14, 2025, <https://openreview.net/forum?id=BDixnHzRL>
7. Optimizing LLM Performance: The Impact of Data Quality and Model ..., accessed March 14, 2025, <https://medium.com/@souhailguennouni/optimizing-llm-performance-the-impact-of-data-quality-and-model-size-95b988cdd4ae>
8. What's a parameter in an LLM?. How to think about a billion parameters - Catherine Breslin, accessed March 14, 2025, <https://catherinebreslin.medium.com/what-is-a-parameter-3d4b7736c81d>
9. The Role of Parameters in LLMs - Alexander Thamm, accessed March 14, 2025, <https://www.alexanderthamm.com/en/blog/the-role-of-parameters-in-llms/>
10. LLM Model Parameter & Memory Required for Training and Inference - Medium, accessed March 14, 2025, <https://medium.com/@plthiyagu/llm-model-parameter-memory-required-for-training-and-inference-634963b36b59>
11. LLM Model Size: Parameters, Training, and Compute Needs in 2025 ..., accessed March 14, 2025, <https://labeleyourdata.com/articles/llm-model-size>
12. Understand LLM sizes | web.dev, accessed March 14, 2025, <https://web.dev/articles/llm-sizes>
13. Reasoning in Granite 3.2 using inference scaling - IBM Research, accessed March 14, 2025, <https://research.ibm.com/blog/inference-scaling-reasoning-ai-model>
14. The State of LLM Reasoning Models - Sebastian Raschka, accessed March 14, 2025, <https://sebastianraschka.com/blog/2025/state-of-llm-reasoning-and-inference-scaling.html>
15. Scaling LLM Test-Time Compute Optimally Can be More Effective than Scaling Parameters for Reasoning | OpenReview, accessed March 14, 2025, <https://openreview.net/forum?id=4FWAwZtd2n>
16. Test-Time Scaling: The New Frontier for AI | CDOTrends, accessed March 14, 2025, <https://www.cdotrends.com/story/4376/test-time-scaling-new-frontier-ai>
17. The Rise Of Reasoner Models: Scaling Test-Time Compute - DEV Community, accessed March 14, 2025, <https://dev.to/rogjia/the-rise-of-reasoner-models-scaling-test-time-compute-33e3>
18. Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters | by Eleventh Hour Enthusiast | Medium, accessed March 14, 2025, <https://medium.com/@EleventhHourEnthusiast/scaling-llm-test-time-compute-optimally-can-be-more-effective-than-scaling-model-parameters-19a0c9fb7c44>
19. s1: Simple test-time scaling - arXiv, accessed March 14, 2025, <https://arxiv.org/html/2501.19393v1>
20. s1: Simple test-time scaling - arXiv, accessed March 14, 2025, <https://arxiv.org/pdf/2501.19393>
21. Top NVIDIA GPUs for LLM Inference | by Bijit Ghosh | Medium, accessed March 14, 2025, <https://medium.com/@bijit211987/top-nvidia-gpus-for-llm-inference-8a5316184a10>
22. What Hardware Is Needed for AI? - Multimodal, accessed March 14, 2025, <https://www.multimodal.dev/post/what-hardware-is-needed-for-ai>
23. AGI and the Explosion of AI Hardware: Tailoring Compute to Intelligence - FRANKI T, accessed March 14, 2025, <https://www.francescatabor.com/articles/2024/12/31/agi-and-the-explosion-of-ai-hardware-tailoring-compute-to-intelligence>
24. AMD GPU Performance for LLM Inference: A Deep Dive - Valohai, accessed March 14,

- 2025, <https://valohai.com/blog/amd-gpu-performance-for-llm-inference/>
25. LLM Inference - Consumer GPU performance - Puget Systems, accessed March 14, 2025, <https://www.pugetsystems.com/labs/articles/llm-inference-consumer-gpu-performance/>
26. Benchmarking NVIDIA TensorRT-LLM - Jan.ai, accessed March 14, 2025, <https://jan.ai/post/benchmarking-nvidia-tensorrt-llm>
27. XiongjieDai/GPU-Benchmarks-on-LLM-Inference: Multiple NVIDIA GPUs or Apple Silicon for Large Language Model Inference? - GitHub, accessed March 14, 2025, <https://github.com/XiongjieDai/GPU-Benchmarks-on-LLM-Inference>
28. LLM Inference Sizing and Performance Guidance - VMware Cloud Foundation (VCF) Blog, accessed March 14, 2025, <https://blogs.vmware.com/cloud-foundation/2024/09/25/llm-inference-sizing-and-performance-guidance/>
29. Hardware Recommendations for Generative AI - Puget Systems, accessed March 14, 2025, <https://www.pugetsystems.com/solutions/ai-and-hpc-workstations/generative-ai/hardware-recommendations/>
30. Hardware Recommendations for Machine Learning / AI - Puget Systems, accessed March 14, 2025, <https://www.pugetsystems.com/solutions/ai-and-hpc-workstations/machine-learning-ai/hardware-recommendations/>
31. LLM Inference Performance Engineering: Best Practices | Databricks Blog, accessed March 14, 2025, <https://www.databricks.com/blog/llm-inference-performance-engineering-best-practices>
32. A guide to LLM inference and performance | Baseten Blog, accessed March 14, 2025, <https://www.baseten.co/blog/llm-transformer-inference-guide/>
33. LLM Benchmarking: Cost Efficient Performance - Nosana, accessed March 14, 2025, https://nosana.com/blog/llm_benchmarking_cost_efficient_performance
34. ninehills/llm-inference-benchmark - GitHub, accessed March 14, 2025, <https://github.com/ninehills/llm-inference-benchmark>
35. Lenovo LLM Sizing Guide, accessed March 14, 2025, <https://lenovopress.lenovo.com/lp2130-lenovo-llm-sizing-guide>
36. Understanding performance benchmarks for LLM inference | Baseten Blog, accessed March 14, 2025, <https://www.baseten.co/blog/understanding-performance-benchmarks-for-llm-inference/>
37. A Guide to LLM Inference Performance Monitoring | Symbl.ai, accessed March 14, 2025, <https://symbl.ai/developers/blog/a-guide-to-llm-inference-performance-monitoring/>
38. What is a context window? - IBM, accessed March 14, 2025, <https://www.ibm.com/think/topics/context-window>
39. Latency optimization - OpenAI API, accessed March 14, 2025, <https://platform.openai.com/docs/guides/latency-optimization>
40. Benchmarking LLMs: TPS, TTFT, GPU Usage | Medium - Ruman, accessed March 14, 2025, <https://rumn.medium.com/benchmarking-llm-performance-token-per-second-tps-time-to-first-token-ttft-and-gpu-usage-8c50ee8387fa>
41. Optimizing inference - Hugging Face, accessed March 14, 2025, https://huggingface.co/docs/transformers/main/llm_optims
42. Why larger LLM context windows are all the rage - IBM Research, accessed March 14, 2025, <https://research.ibm.com/blog/larger-context-window>
43. Guide to Context in LLMs | Symbl.ai, accessed March 14, 2025,

<https://syml.ai/developers/blog/guide-to-context-in-llms/>

44. How Does The Context Window Size Affect LLM Performance? - Deepchecks, accessed March 14, 2025,

<https://www.deepchecks.com/question/how-does-context-window-size-affect-llm-performance/>

45. Observations About LLM Inference Pricing - LessWrong, accessed March 14, 2025,

<https://www.lesswrong.com/posts/mRKd4ArA5fYhd2BPb/observations-about-llm-inference-pricing>

46. Welcome to LLMflation - LLM inference cost is going down fast | Andreessen Horowitz, accessed March 14, 2025, <https://a16z.com/llmflation-llm-inference-cost/>

47. Machine Learning Trends - Epoch AI, accessed March 14, 2025, <https://epoch.ai/trends>