

Paper Review: Training Compute-Optimal Large Language Models

Summary: This paper addresses a fundamental question in large-scale language model (LLM) training: given a fixed compute budget (C), how should one balance model size (N) and number of training tokens (D) to achieve optimal performance? The authors observe that most prior large models, such as GPT-3 and Gopher, scaled up parameters dramatically while keeping training data roughly constant (around 300 B tokens). This imbalance left those models substantially under-trained relative to their compute expenditure, leading to inefficient use of resources and sub-optimal downstream performance. To investigate this, they conduct a large-scale empirical study, training and analyzing over 400 transformer language models ranging from 70 M to 16 B parameters. They model the final pre-training loss, $L(N, D)$, and minimize it under a fixed FLOP constraint, $\text{FLOPs}(N, D) = C$. Three complementary approaches are used to estimate the optimal model and dataset sizes: fixing N while varying D , constructing IsoFLOP profiles by fixing FLOPs and varying N , and fitting a parametric form of $L(N, D)$. All three methods converge to a consistent scaling relationship: for compute-optimal training, model size and number of training tokens should increase proportionally, that is, doubling model parameters should be accompanied by doubling the training data. This finding contrasts sharply with Kaplan et al. (2020), who suggested model size should scale much faster than data. To validate their new law, the authors trained Chinchilla (70 B parameters, 1.4 T tokens) using the same compute as Gopher (280 B, 300 B tokens). Both were trained on the MassiveText dataset, and all the analysis was performed on TPUv3/v4. Despite being four times smaller, Chinchilla outperformed Gopher, GPT-3, and MT-NLG on 130+ benchmarks (MMLU 67.5 % vs 60 %). Its compactness also cuts inference and fine-tuning costs. Overall, the work establishes a revised scaling law emphasizing balanced data-model growth and greater training efficiency.

Evaluation: The paper succeeds in achieving its stated goal of identifying compute-optimal strategies for training large language models. Its significance is high, as it challenges the long-standing assumption that increasing parameter count alone drives performance, and instead demonstrates that most existing LLMs are under-trained relative to their compute. This shift in perspective has substantial implications for efficiency, accessibility, and environmental impact. The solution is both valid and novel, the authors ground their conclusions in an exceptionally broad empirical base of over 400 training runs and reinforce them with three independent analytical methods that converge on the same result. The Chinchilla-Gopher comparison provides strong, real-world validation under equal compute budgets. Experimental evaluation is detailed, covering diverse benchmarks with transparent methodology and well-controlled variables. The paper is clear, well-structured, and easy to follow, effectively balancing mathematical concepts with practical interpretation.

Main Takeaways

1. Compute-optimal training requires equal scaling of model size and dataset size.
2. Many existing large models (for example: GPT-3, Gopher) are under-trained relative to their compute.
3. Smaller, well-trained models can outperform much larger ones at lower cost.

Strengths

1. The study draws on over 400 transformer models trained across a broad range of sizes and compute budgets, giving the results strong empirical grounding and reducing the chance that the observed patterns are coincidental.
2. Three independent approaches, varying data, IsoFLOP profiling, and parametric modeling, all lead to the same scaling law, reinforcing the reliability and robustness of the findings.

Weaknesses

1. Variations such as AdamW optimization, mixed-precision weights, and tokenizer updates in Chinchilla may have influenced performance, making it hard to fully isolate the impact of data-parameter scaling.
2. Limited transparency on dataset composition and filtering raises reproducibility, bias, and privacy concerns given the trillion-token scale.

Discussion: Would Gopher's performance converge with Chinchilla's if trained using AdamW and identical preprocessing?