# On the Origin of Algorithmic Progress in AI

**Hans Gundlach**[†,*]
hansgund@mit.edu

**Alex Fogelson**[†]
fogelson@mit.edu

**Jayson Lynch**[†]
jaysonl@mit.edu
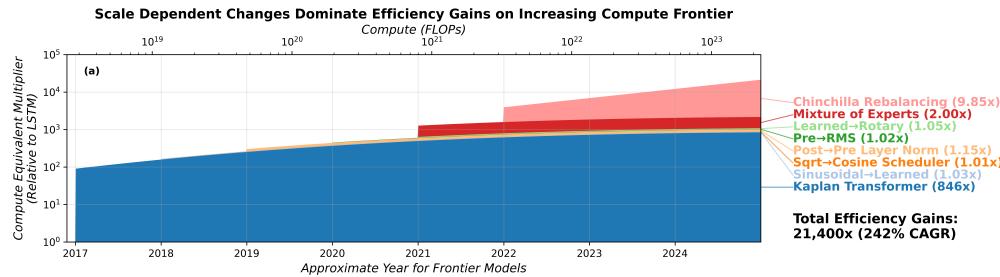
**Ana Trišović**[†]
ana_tris@mit.edu

**Jonathan Rosenfeld**[†]
jonsr@csail.mit.edu

**Anmol Sandhu**[‡]
asandhu@alumni.olin.edu

**Neil Thompson**[†,*]
neil_t@mit.edu

[†]MIT FutureTech, CSAIL     [‡]Olin College
[*]Corresponding authors

## Abstract

Algorithms have been estimated to increase AI training FLOP efficiency by a factor of $22,000$ between 2012 and 2023 [Ho et al., 2024]. Running small-scale ablation experiments on key innovations from this time period, we are able to account for less than $10\times$ of these gains. Surveying the broader literature, we estimate that additional innovations not included in our ablations account for less than $10\times$, yielding a total under $100\times$. This leads us to conduct scaling experiments, which reveal that much of this efficiency gap can be explained by algorithms with scale-dependent efficiency improvements. In particular, we conduct scaling experiments between LSTMs and Transformers, finding exponent differences in their compute-optimal scaling law while finding little scaling difference for many other innovations. These experiments demonstrate that – contrary to standard assumptions – an algorithm's efficiency gains are tied to compute scale. Using experimental extrapolation and literature estimates, we account for $6,930\times$ efficiency gains over the same time period, with the scale-dependent LSTM-to-Transformer transition accounting for the majority of gains. Our results indicate that algorithmic progress for small models has been far slower than previously assumed, and that measures of algorithmic efficiency are strongly reference-dependent.

A preview of our main plot Figure 6. This shows our decomposition of algorithmic progress for models scaling at the frontier as defined by Epoch AI [2025]. Importantly, the scale-dependent algorithmic change from LSTM to Transformer comprises most of this progress.

Preprint.

# 1 Introduction

Progress in AI over the past decade has been driven by two tightly coupled forces: rapidly increasing compute budgets and a succession of algorithmic efficiency innovations in architectures, optimization, and other training practices. While compute growth is relatively straightforward to measure, we still lack a clear accounting of algorithmic progress—what concrete changes drive efficiency gains, how large those gains are, and whether they persist across compute scales.

Recent work by Ho et al. [2024] has started to address this question. By analyzing hundreds of historical language models, they estimate that in the last 10 years, algorithmic progress has contributed more than 4 orders of magnitude in so-called "effective compute" while compute scaling has increased 7 orders of magnitude based on analyzing historic AI literature. Specifically, models have gotten roughly $22,000\times$ more efficient across the sum total of innovations in algorithms, supposedly allowing the same performance level with drastically fewer FLOPs.

However, a precise decomposition of this progress remains largely unknown, and many questions about the origins of algorithmic advancement remain underexplored. For instance, how do algorithmic advances interact? Are algorithmic advances driven by a series of small improvements, or a few large ones? Do algorithmic improvements follow smooth, Moore's law-like trends, or are they driven by punctuated equilibria? Answering these questions has important implications for the future of AI, as it could enhance the predictability of future algorithmic gains, guide innovation efforts towards specific areas of improvement, and inform our expected returns to compute scaling.

In order to answer these questions, we take three complementary approaches: we conduct ablation experiments on important algorithmic advances in language models; we conduct scaling experiments to measure differences in optimal scaling across architectures; and we perform theoretical analysis of transitions in data and parameter scalings. Our experiments focus primarily on architectural algorithmic modifications, but we try to address other methods to improve efficiency, like data improvements, using literature estimates. Our analysis points to three core conclusions:

1. **We find most algorithmic innovations we experimentally evaluate have small, scale-invariant efficiency improvements** with less than $10\times$ compute efficiency gain overall, and representing less than $10\%$ of total improvements extrapolated to the 2025 compute frontier ($2 \times 10^{23}$ FLOPs). This suggests that scale-invariant algorithmic progress contributes only a minor share of overall efficiency improvements (Section 3).

2. **We find two strongly scale-dependent algorithmic innovations: LSTMs to Transformers, and Kaplan to Chinchilla re-balancing.** Together, these account for $91\%$ of total efficiency gains when extrapolating to the 2025 compute frontier. This implies that algorithmic progress for small-scale models is several orders of magnitude smaller than previously thought (Section 4).

3. **We show that in the presence of scale-dependent innovations, not only do efficiency gains require continued compute investment, but the rate of algorithmic progress strongly depends on your choice of reference algorithm.** In other words, the rate of progress in successive models can appear exponential relative to one baseline algorithm, yet be zero relative to another (Section 5.2).

Collectively, these findings suggest that gains in algorithmic progress may be inherently scale-dependent, requiring ever-increasing compute to realize its benefits, and implying that algorithmic progress has benefited larger model builders more than smaller players in AI development.

## 1.1 Previous Work

There is a growing body of work studying the impact of algorithmic advances in machine learning. Hernandez and Brown [2020] studied computer vision models between 2012 and 2019 and found the number of training FLOPs required to reach AlexNet performance decreased $44\times$. Ho et al. [2024] conducted a literature-based analysis of FLOP efficiency and found algorithmic gains in language models from 2012-2023 on the order of $22,000\times$. However, these gains have started to come under criticism. Whitfill [2025] found that correcting for observational bias, in particular the correlation between compute scale and algorithmic efficiency, algorithmic progress may be an order of magnitude smaller. In addition, there has been relatively little experimental work on the extent of

overall historical algorithmic progress. Sanderson et al. [2025] did a small selection of experiments identifying impact of innovations like rotary embeddings and layer normalization, finding that both of these innovations contributed efficiency gains of $1.7\times$.

To understand the effects of algorithmic choices, it is important to understand how innovations scale with greater computational resources. If algorithmic advances have scale-dependent efficiency gains, this would lead to changes in scaling exponents. Algorithmic advances are generally thought to have little effect on scaling exponents: Hestness et al. [2017] found no data scaling exponent difference between LSTM and RHN (Recurrent Highways Networks) architectures, nor between Adam and SGD optimizers. Bansal et al. [2022] measured the data scaling exponent between Encoder-Decoder, Decoder only, and mixed LSTM-Transformer architectures and found little difference. Similarly, the theory literature has focused on the intrinsic properties of training data as the most important factor in determining the exponent in neural scaling laws [Rosenfeld, 2021, Michaud et al., 2023, Bahri et al., 2024].

However, there is a growing body of literature pointing to the effects of algorithmic advances on scaling. Droppo and Elibol [2021] conducted scaling studies on acoustic models and found that Transformers have a steeper scaling exponent than LSTMs. Sanderson et al. [2025] examines the efficiency gains from selected papers and speculates that some algorithmic advances, like the LSTM-Transformer transition, had scale-dependent effects. There are also new AI architectures like KANs (Kolmogorov-Arnold Networks) that are purported to have a larger/steeper scaling exponent [Liu et al., 2024].

It is in this light that we conduct experiments to understand what algorithmic changes have driven the growth in language models and the potentially scale-dependent nature of these changes.

## 2    Expanding the Compute Equivalent Gains Framework

In order to more precisely discuss algorithmic progress, we introduce a generalized notion of the standard compute equivalent gains (CEG) framework [Davidson et al., 2023, Ho et al., 2024, Sanderson et al., 2025]. Our approach varies from much of the existing literature by allowing algorithmic innovations to have different efficiency gains across compute scales. More precisely, we define the **CEG function** and **CEG multiplier** distinctly:

> *Definition:* For two algorithms $A$, $A'$ and a fixed performance metric (e.g. loss), the **compute equivalent gain function** (CEG function) from $A$ to $A'$ is a function $f$ such that a model trained with $A$ using compute $C$ and $A'$ using compute $C/f(C)$ reach equivalent performance, for all $C$ in the domain of $f$.

Fixing the compute budget for one of the algorithms, we can recover the traditional definition of a CEG multiplier:

> *Definition:* Let $M$ be a model trained with algorithm $A$ using compute $C$. Let $f$ be the CEG function between $A$ and some $A'$. Then we call $f(C)$ the **CEG multiplier** of $A'$ *relative to $M$*.

Intuitively, these definitions refer to two distinct categories: CEG functions compare training algorithms to training algorithms, defining the efficiency gains across varying compute scales; CEG multipliers compare models to models by fixing a performance level.[1] This distinction will become central to our analysis, as we show that merely measuring CEG multipliers may exhibit misleading trends when algorithms exhibit scale dependence.[2]

Note that in specific cases where an algorithm plateaus in performance, the CEG function may not be well defined, requiring "infinite" compute to reach a given performance level [Erdil and Besiroglu, 2022]. Finally, if the CEG function is *known* to be constant between two algorithms, one may simply choose a preferred performance threshold compatible with both algorithms to calculate that constant value, and thereby recover the entire function. For our ablation experiments, we indeed argue that

---

[1]"Training algorithms" includes everything except compute i.e., architecture, optimizers, data-parameter balancing, etc.

[2]When ablating algorithmic innovations experimentally, we report CEG multipliers between algorithms using the compute optimal performance within our hyperparameter search.

the ablated algorithmic innovations admit constant CEG functions (i.e., are scale invariant with respect to our baselines), and we choose the minimum negative log-likelihood threshold compatible between our ablations when computing this constant (namely $NLL = 5.3$, see Section 2).

## 3 Scale Invariant Algorithms

We first analyze the effects of individual algorithms by running a large series of ablation experiments to paint a fine-grained picture of algorithmic improvements. For instance, by running a transformer with older activation functions like GeLU and newer activation functions like SwiGLU, we can quantify efficiency gains by measuring differences in compute required to reach a fixed performance level (see Section 2). In addition, we try to estimate the joint effects of algorithms by ablating multiple algorithmic innovations and comparing this to the multiplicative effects of combining individual ablations.

We examine activation functions, positional encodings, normalization techniques, learning rate schedules, and optimizers. We leave out tokenizers and data quality enhancements as they are difficult to evaluate using perplexity, though we examine the literature that has benchmarked these contributions. In general, we find that estimates of algorithmic improvements from the primary literature (i.e., from the papers that proposed the algorithmic improvement itself) are significantly higher than those from the secondary literature (i.e., estimates from papers that include multiple previous algorithms as benchmarks), as well as from our own experiments.

In Section 4, we demonstrate near-identical scaling exponents between our baseline transformer and a transformer with all innovations ablated (Figure 4). This constitutes strong evidence that each of these innovations is indeed scale-invariant within the compute regime of our experiments.

### 3.1 Ablation Experiments

For our experimental approach, we examine a small 3.6M parameter transformer with rotary-based encodings, a constant width-to-depth ratio, and a GeLU activation function. Our baseline is identical to our "modern" transformer, described in our scaling experiments (Section 4.1.1). For each ablation, we run learning rate tunes by testing the default learning rate ($10^{-2.25}$) and $[10^{-2.5}, 10^{-2}]$.[3] As mentioned previously, we choose the lowest loss threshold that all ablated models are capable of (in our case, a loss of 5.3) to determine the CEG multiplier. All of our ablation experiments use a token-to-parameter ratio of 20 as in Hoffmann et al. [2022].
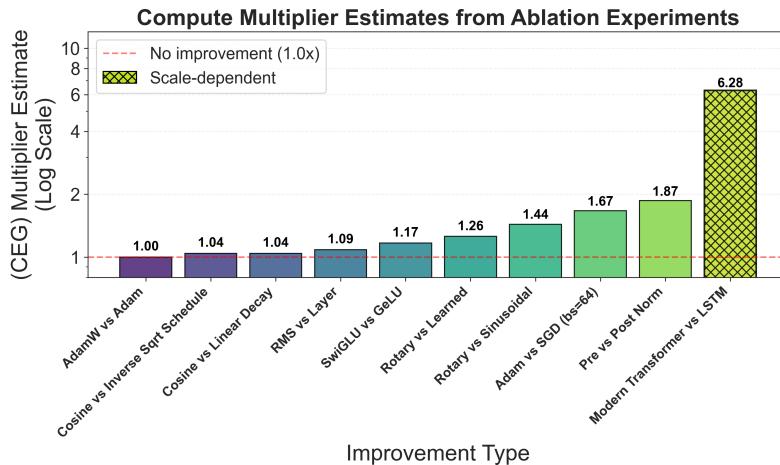


Figure 1: Compute Equivalent Gain multiplier for algorithms measured on a 3.6M parameter transformer model using ablation experiments. Many recent training advancements have a small impact on training efficiency. Hatched bar represent improvements we believe are scale-dependent.

---

[3]If optimal learning rate is different then default, we test learning rates outside this range by factors of $\sqrt{10}$

**Activation Functions**   We test the impact of new activation functions like SwiGLU over the traditional GeLU activation function. In line with Shazeer [2020], we find that SwiGLU, has a measurable $1.17\times$ efficiency gain over GeLU even taking into account its increased parameter requirements.

**Positional Encoding**   We find that positional encodings might have some impact on training compute efficiency. We test three different types of positional encoding: rotary, sinusoidal, and learned positional encodings. We estimate that rotary encoding constitutes an improvement of $44\%$ in training efficiency over sinusoidal encoding originally employed in Vaswani et al. [2017]. This is in contrast to Sanderson et al. [2025], which estimates a $70\%$ gain, and consistent with Biderman et al. [2021], which sees a $10\% - 30\%$ training efficiency gain. However, we think such efficiency gains might be highly dependent on sequence length, as rotary encodings were originally developed to incorporate larger and more flexible context windows rather than as a general training efficiency improvement.

**Learning Rate Schedules**   Our test consists of three prominent learning rate schedules used in the last 10 years. These include a linear warmup/inverse square-root learning rate decay schedule used in T5 [Raffel et al., 2020] and the Transformer [Vaswani et al., 2017]. We also implement a linear warmup/linear decay schedule as used in BERT [Devlin et al., 2019] and linear-warmup with cosine decay, which is the more recent default used in GPT-3 [Brown et al., 2020]. We observe very small training efficiency gains between any learning rate schedules ($< 5\%$). These are roughly in line with Kaplan et al. [2020], which finds minimal difference between state-of-the-art learning rate schedulers.

**Normalization Techniques**   We test three different normalization variations: a pre-RMSNorm model, a pre-layernorm model, and a post-RMSNorm model. Pre-layernorm(RMSNorm) applies layernorm before multi-head attention in the residual stream. This is generally seen as more stable than its predecessor post-layernorm [Xiong et al., 2020]. RMSNorm is an update to layernorm that removes the mean centering step and can have runtime improvements of $7\% - 64\%$ [Zhang and Sennrich, 2019]. We observe substantial improvements of $87\%$ transitioning from post-RMSNorm to pre-RMSNorm, with more modest improvements of $9\%$ switching from a pre-layernorm model to a pre-RMSNorm model.

**Optimizers**   Prior work has shown that SGD performs notably worse on transformers in comparison to Adam [Ahn et al., 2023, Zhao et al., 2024]. This is theorized to be due to Hessian heterogeneity [Zhang et al., 2024]. More recent results suggest that this gap becomes small when SGD uses high momentum and very small batch sizes [Srećković et al., 2025]. Further, Srećković et al. [2025] argues that the optimal batch size for SGD is 1 and that at this scale, there would be a negligible gap between the two optimizers. In addition, Zhang et al. [2019] finds that Adam is compatible with much higher critical batch sizes than SGD.

At batch sizes above 128, SGD has very poor performance. However, with our default batch size of 64 along with the 0.98 momentum and no weight decay, recommended by Srećković et al. [2025], we observe smoother scaling behavior and a sizable $1.87\times$ gap in training efficiency compared to AdamW. In addition, we further examine the scaling effects of optimizer choice in Appendix C and find little scale-dependent change in compute efficiency gain. This is consistent with scaling studies of new optimizers like Muon, which similarly report minimal changes in scaling exponent [Liu et al., 2025]. We compare Adam to AdamW but find a negligible difference for our model.

## 3.2   Interaction Between Algorithms

Conventional theory of algorithmic gains presupposes that most algorithmic gains are multiplicative [Ho et al., 2024, Davidson et al., 2023]. For instance, if algorithm $A$ increases training efficiency by $2\times$ and algorithm $B$ increases training efficiency by $3\times$, respectively, then the combination of algorithm $A$ and $B$ should have a training efficiency gain of $6\times$. However, across architectures this is not the case. For example, Adam provides a computational efficiency gain for transformers, yet most state-of-the-art LSTMs relied on RMSProp or plain SGD over Adam [Merity et al., 2017, Melis et al., 2017, Keskar and Socher, 2017].

**Modern Transformer Improvements**

Compute Equivalent Multiplier

Total: 3.55x

SwiGLU etc 1.32x

Rotary Encoding 1.44x

Post to Pre 1.87x

1.33x

Theoretical: Post to Pre Norm × Rotary Encoding × SwiGLU and Other Innovations
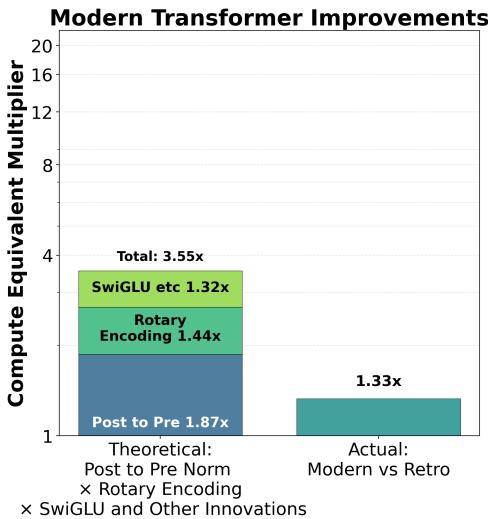
Actual: Modern vs Retro

Figure 2: Comparison of the total effect of all measured post transformer algorithmic changes (left bar) vs multiplying ablation estimated effects together.



**Transformer Evolution from LSTM**

Compute Equivalent Multiplier

Total: 16.66x

SwiGLU etc 1.32x

Rotary Encoding 1.44x

Post to Pre 1.87x

LSTM to Retro Transformer 4.69x

6.28x

Theoretical: Retro × Post to Pre Norm × Rotary Encoding × SwiGLU and Other Innovations
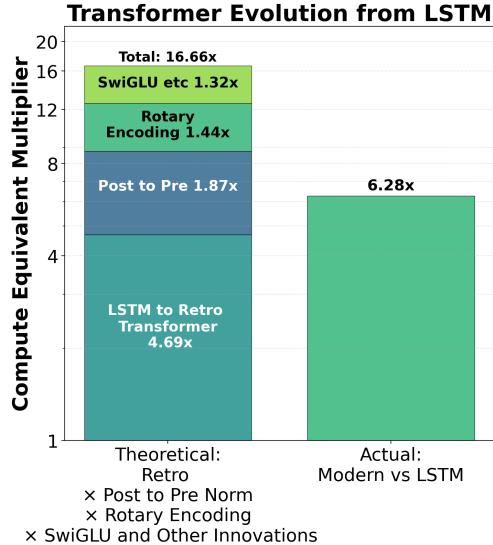
Actual: Modern vs LSTM

Figure 3: Comparison of the total effect of all measured changes, including transformer to (left bar) vs multiplying ablation estimated effects together.

To help identify these interactions, we train a transformer model similar to the original architecture developed by Vaswani et al. [2017]. In this "Retro" Transformer, four algorithmic improvements are reverted. We switch from SwiGLU to GeLU, pre-RMSNorm to post-layernorm, rotary encoding to sinusoidal encodings, and cosine decay learning rate schedule to square root decay schedule. Using the estimates in Figure 1 and assuming multiplicative improvement we would expect CEG gains of $3.43\times$. However, the efficiency gain is only $1.33\times$. In addition, we compare our default or "Modern" transformer (defined in Section 4.1.1) to our LSTM baseline. At a loss threshold of $5.3$, we find a combined efficiency gain of $6.28\times$, whereas multiplying the individual improvements together, along with the LSTM to Retro Transformer improvements, yields a substantially higher predicted gain of $16.66\times$.[4]

This suggests that multiplying individual improvements, estimated from ablation experiments, overestimates the impact of algorithmic improvements. We hypothesize the following explanation by analogy: consider replacing the engine of a modern car with a *Model T* engine. A *Model T* engine in an otherwise modern car might perform far worse than it did in the original *Model T* due to system incompatibilities (similarly with any other component). As a result, one would overestimate the proportion of progress for each of these inventions independently. Similarly, it may be that replacing algorithmic components individually appears to cause greater efficiency loss than replacing all components at once. In this sense, algorithmic improvements are sub-multiplicative and slightly non-orthogonal.

### 3.3 Literature Estimates

**Why Include Literature Estimates?** In addition to performing ablation experiments, we include estimates from the literature for Mixture-of-Experts (MoE) models and tokenization methods. MoE models are challenging to implement and measure at small scales due to complex routing mechanisms and loss-balancing requirements. Moreover, there is already substantial literature comparing dense and MoE transformers.

Tokenizers also present significant challenges. Perplexity is inherently tokenizer-dependent. Therefore, to measure tokenizer influence using perplexity, one would need to evaluate validation loss

---

[4]In the case where we have scaling fits, we use the CEG between scaling fits to be consistent with Figure 4. If we use individual training curves, we would get a Modern-to-Retro Transformer difference of $2.7\times$ and a Modern-to-LSTM difference of $5.9\times$

with a fixed tokenizer while training with a different tokenizer—a clearly unprincipled approach. Although tokenizer effects could be studied experimentally using other approaches like perplexity per byte [Choe et al., 2019] or non-perplexity evaluation metrics (e.g., MMLU), these are beyond the current scope of this paper.

**Mixture of Experts**   MoE is generally considered a substantial inference improvement, but MoE architectures do significantly improve training performance as well. However, we estimate that this is at or below a $2\times$ factor in FLOP efficiency holding the number of experts constant. Riquelme et al. [2021] scales both mixture of experts and dense vision models and shows a computational efficiency gain of slightly less than $2\times$. He [2024] looks at isoflop curves for dense and MoE models, finding their MoE model has an optimal perplexity of 22, while the dense model with double the compute has a perplexity around 18. This suggests that their MoE model has a computational efficiency gain significantly less than $2\times$. Li et al. [2025] does not explicitly train models Chinchilla-optimally, but mentions that their MoE model approaches the performance of dense models with twice the compute. Muennighoff et al. [2024] finds that MoE improves training compute by $3\times$ but does not do an optimal isoflop comparison, which makes controlled efficiency comparisons challenging. Tay et al. [2023] demonstrate similar scaling behavior between dense and MoE transformers, though their results exhibit noisy performance. Given the uncertainty surrounding the scaling effects of MoE training, we leave detailed investigation to future work. Where applicable, we approximate MoE efficiency gains as $2\times$.

**Tokenizers**   Early language models used word-level or character-level tokenization, which led to too large vocabularies or too large sequence lengths, respectively. In the mid-2010s subword tokenization schemes like BPE, WordPiece, and SentencePiece became popular, which mitigated many of the disadvantages of word-level and character-level language models. Despite these advances, Ali et al. [2024] finds that inefficient tokenization leads to an additional training cost of $68\%$. We avoid reporting these improvements in our aggregations due to the uncertainty in these measurements and the lack of strong historical comparison studies with standardized evaluations.

### 3.4   Scale Invariant Algorithms Show Small, Highly Unequal Distribution

Experimentally, we find that the total measured efficiency gains switching from LSTM to a Modern Transformer of $6.28\times$, while switching from an LSTM to our Retro Transformer has efficiency gains of $4.69\times$. This is much smaller than estimates from Ho et al. [2024], which estimate a $60\times$ from the LSTM to Transformer transition. Although we see moderate $\sim 2\times$ gains for some innovations like Adam and or post-layernorm to pre-RMSNorm, the majority of the innovations we measure have small gains.

Interestingly, we see evidence of a highly skewed distribution of efficiency improvements. While all ablated innovations in our study are under $4\times$, the distribution of multipliers is uneven and concentrated in a few notable improvements, such as Adam or pre-layernorm. In this light, algorithmic progress looks much more disjoint than previously thought, whereas many years of incremental changes are followed by larger algorithmic transitions.

The small level of efficiency gain for many of these improvements motivates the second part of our experimental study, where we compare the effects of algorithmic changes across scales, revealing the effect of these changes on neural scaling laws.

## 4   Scale-Dependent Algorithms

We saw in Section 3.4 that algorithmic gains are relatively small at small scales. This raises the natural question, are algorithmic gains larger at larger scales? To answer this question, we conduct scaling experiments across architectures, optimizers, and the innovations from Section 3.1 to get a better sense of how their efficiency gains scale with compute. In addition, we exhibit a mathematical analysis of the Kaplan to Chinchilla scaling shift, which gives a particularly interesting example of a scale-dependent algorithmic change that is not an exponent shift.

### 4.1 Scaling Experiments: LSTM to Transformer

#### 4.1.1 Experiments Setup for Scaling

We implement two main model architectures to serve as baselines in our experiments: LSTMs and transformers. In addition, we implement two transformer variants. The first is our Modern Transformer, which incorporates all the recent improvements we study and serves as our default transformer model. The second is our Retro Transformer introduced in Section 3.2, which reverts these modifications. We also define **"Kaplan" Transformers** to be our Retro Transformer adjusted analytically to approximate Kaplan scaling, representing our ablations as well as suboptimal data scaling.

**Modern Transformer** We first standardize our notion of a "modern" transformer from which we perform ablations. We use a vanilla transformer with rotary-based encodings and keep a constant width-to-depth aspect ratio. Kaplan et al. [2020] used a $N : L \approx 100$ width-depth ratio [Dey et al., 2025], but at the scale of our experiments, we choose a ratio of $N : L \approx 16$. This gives us depth to width-to-depth ratio similar to that of other small transformers like Lan et al. [2019]. Attention heads number is commonly set such that the dimension of each head stays at a constant size, such as 64 or 128 [Hoffmann et al., 2022, Vaswani et al., 2017]. Given the scale of our experiments, we choose a constant attention head dimension of 16. In our experience, setting a larger head dimension and therefore reducing the number of attention heads degraded performance at small scales. We measure FLOP numbers by running PyTorch profiler over one minibatch and scaling according to gradient accumulation and number of steps. Our initialization, unless stated otherwise, has 0 bias and all weights are sampled from $\mathcal{N}(0, 0.02)$ as is done in BERT [Lan et al., 2019]. We use pre-layer normalization [Xiong et al., 2020] with RMSNorm.[5]

**Retro Transformer** Our Retro Transformer is the same as our modern Transformer model, but with four key algorithmic components reverted. We switch from SwiGLU to GeLU, pre-RMSNorm to post-layernorm, rotary encoding to sinusoidal encodings, and cosine decay learning rate schedule to square root decay schedule.[6]

**LSTM Model** For our LSTM model we choose a vanilla LSTM and do not analyze more recent variations developed after 2018 (e.g., Melis et al. [2019] and Beck et al. [2024]). We also exclude more-recent variations of LSTMs such as GRUs. Greff et al. [2016] finds that these variants do not improve on the standard LSTM. We develop a setup based on state-of-the-art LSTMs before the invention of the Transformer using hyperparameters from Melis et al. [2017]. We focus our efforts on a 1-layer LSTM, which seems to perform the best in practice for an LSTM with 10M parameters [Melis et al., 2017]. In addition, it is noticeably better than a 2-layer LSTM at the largest and smallest sizes we test, and represents the recommendation based on the aspect ratio rule in [Droppo and Elibol, 2021]. We use initialization consisting of Xavier-uniform embedding weights, orthogonal recurrent weights, and 0 bias except for a forget gate bias of $+1$. We include a variational dropout with separate dropout rates for embedding matrices vs hidden matrices [Melis et al., 2017, Merity et al., 2017]. This most closely follows [Melis et al., 2017] except for the choice of initialization and our decision not to tie LSTM gates.

Training an LSTM on internet-scale data under Chinchilla-optimal conditions is somewhat unnatural, as most historical LLMs were trained on much smaller datasets where overfitting is a large issue. In our implementation, we set all LSTM dropout parameters equal to 0, as is the practice for transformer models in this setting Hoffmann et al. [2022]. This improves our LSTM loss from 6.5 to 5.9 for our compute-optimal 3.3M parameter LSTM model.[7]

We scale two different transformer model variants as well as our default LSTM. For each, we examine training runs of models with hidden dimension $[32, 64, 96, 128, 160, 192, 224, 256]$. This corresponds to transformer models with between 1.6M parameters and 29M parameters and LSTM models with between 1.6M and 11.2M parameters. To determine the optimal learning rate, we use a procedure similar to Bi et al. [2024], which is outlined in Appendix A.1.

---

[5]For more details on our default setup, please see Appendix 2.

[6]For more details on our Retro Transformer setup, please see Appendix 2.

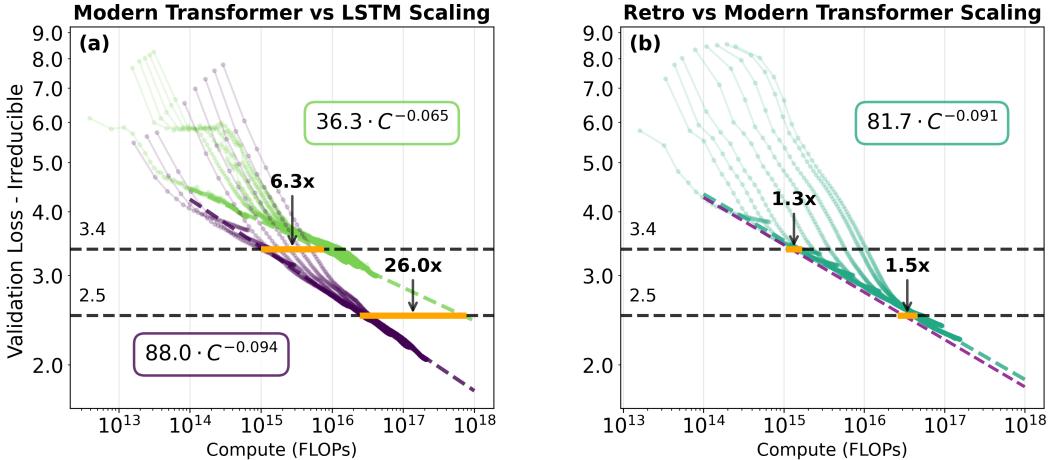[7]See Appendix 2 for more details on our hyperparameter choices.

Figure 4: Figure (a) depicts the scaling difference between a Modern Transformer in purple and a standard LSTM in green. Figure (b) depicts the scaling difference between a Modern Transformer in purple and a Retro Transformer in blue, where all post-2017 innovations are ablated. LSTM seem to have significantly different scaling exponents, while post-2017 transformers have minimal effect on scaling. All graphs depict the training curve for models with hidden dimensions between 32 and 256 with all other hyperparameter scaled proportionately.

Our LSTM model is based on Melis et al. [2017], where we use the parameters from their word-level 10M parameter model. Optimal layer number seems to scale very moderately with LSTM size [Melis et al., 2017], so we keep the layer number constant while scaling the hidden dimension for the range of scale we examine (see also the LSTM scaling study done by Droppo and Elibol [2021]).

We scale our default Modern Transformer and Retro Transformer with algorithmic changes ablated. For each setup, we follow conventional scaling practices, keeping depth-to-width ratio constant.[8]

### 4.1.2 Transformers are a Scale-Dependent Improvement to LSTMs

We plot the difference in scaling between LSTMs and a Modern Transformers in Figure 4A and between modern and Retro Transformers in Figure 4B. For each, we subtract the irreducible loss component (1.9) (see Appendix D) and fit a power law to all Pareto-optimal (frontier) points from $10^{16}$ FLOP onward.[9]

Our scaling graphs suggest that improvements in neural network architecture are not scale invariant and have increasing returns to scale effects. Interestingly, our LSTM model at the smallest scales we measure is only $6.28\times$ less compute-efficient than our transformer model (using a loss threshold of 5.3 or above, see Section 3.1). While at a validation threshold of 4.4, the efficiency gap is $26\times$

Meanwhile, the choice between Adam and SGD, as well as the difference between a modern and a Retro Transformer, seems to be relatively scale-invariant. This provides direct evidence that optimizers may exhibit constant CEG functions, and strong evidence that our ablated algorithms would not individually show scale-dependencies.

### 4.2 Transition from Kaplan to Chinchilla Scaling Laws

In our discussion so far, we have focused on model training choices and have paid relatively little attention to how models are scaled, assuming they are scaled optimally for a set of algorithms. Here, we mathematically examine the shift from Kaplan to Chinchilla scaling practices, which are a particularly interesting example of scale-dependent algorithmic progress that is not an exponent shift. Kaplan et al. [2020] advocates for a smaller level of data scaling as compute is increased in comparison to Hoffmann et al. [2022]. Subsequent literature found that Kaplan misestimated

---

[8]See Appendix A and the following discussion for more details on our experimental setup.
[9]The scaling experiments for Adam-to-SGD are described in Appendix C.

scaling law exponents due to small sizes, not accounting for embedding parameters, and using a constant warmup size [Porian et al., 2024, Pearce and Song, 2024]. In contrast, Hoffmann et al. [2022] accounted for all parameters and found that data and parameters should be scaled equally while scaling compute (so-called "Chinchilla-scaling"). Our analysis provides a closed form for the CEG function between Kaplan and Chinchilla scaling, and displays larger gains than previous estimates [Ho et al., 2024].

To formalize this switch from Kaplan to Chinchilla laws, we assume that compute scaling laws follow a Chinchilla-optimal form, while model trainers allocate parameters and data according to misspecified Kaplan scaling laws (see Appendix B). This form of algorithmic improvement has clear scale-dependent effects, though not a simple exponent change. Training efficiency initially converges between the two scaling laws (within the regime tested by Kaplan), but diverges at larger scales. Furthermore, the Chinchilla training efficiency gains are large but remain below 10x for much of the scale used in contemporary machine learning (i.e., $10^{16} - 10^{24}$ FLOPs).
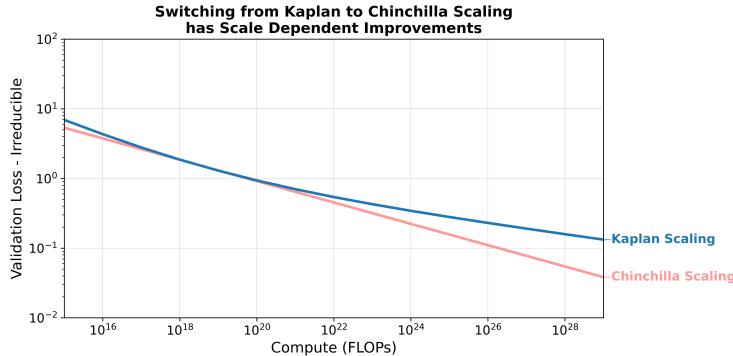


Figure 5: Analytical difference in performance between transformer models scaled with Kaplan versus Chinchilla recommendations. Interestingly, the efficiency gap first converges, then diverges.

## 5  Algorithmic Progress Depends Strongly on Compute, Reference Points

Our theoretical and experimental results allow us to decompose algorithmic progress into individual innovations introduced over time. This leads to two main observations. First, many measured CEG gains may be a consequence of compute scaling, rather than a multitude of algorithmic innovations. Second, when innovations are scale dependent, the growth rates of CEG multipliers depends entirely on the chosen reference model, meaning that progress may appear unbounded with respect to one algorithm while non-existent with respect to another.

### 5.1  Compute (Not Time) May Explain Most Algorithmic Progress

Existing estimates of algorithmic progress measure the rate of progress in algorithms as a function of time [Ho et al., 2024]. However, our analysis opens the question as to whether this time-dependence is instead driven by steady increases in compute investment. Estimates show the compute budgets of frontier models have exponentially increased at a rate of $4.2\times$ per year [Epoch AI, 2023]. Thus as compute budgets increase exponentially, the rate of algorithmic progress may in fact be driven more by the regularity of compute scaling rather than the discovery of new innovations.

We stack our empirical results, as well as literature estimates for mixture of experts and theoretical results, to extrapolate the expected CEG gains over varying compute scales (with respect to LSTMs). To compare our results to literature estimates (which track gains over time), we use the exponential relationship between time and compute among notable models [Epoch AI, 2023] and refer to this as the "compute frontier". This allows us to calculate the cumulative CEG function with respect to compute, and do so over historically relevant time intervals.

When stacking our ablated, scale-invariant innovations, we scale down their total CEG multiplier according to our interaction experiments. When reporting the decomposition, the results remain log-proportional to the independent multipliers. Moreover, when reporting gains from Kaplan Trans-
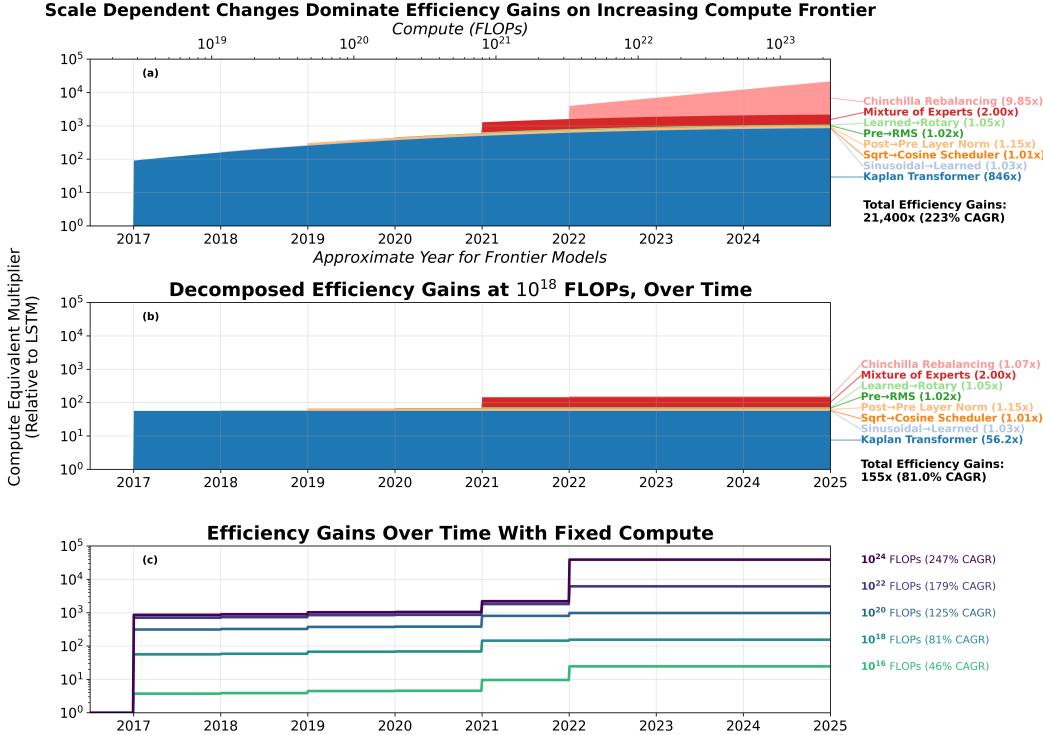
Figure 6: LSTM to Transformer transition has been responsible for most of the algorithmic progress out of the innovations we review (repeated from Introduction). We adjust scale invariant transitions to be log-proportional to our overall ablation results, given the interaction effects we demonstrate in Section 3.2.

formers, we adjust our ablated transformer according to the analytically derived gains from the Kaplan-to-Chinchilla CEG function.

Importantly, between 2017 and 2025, we find the overwhelming majority of algorithmic progress can be accounted for by two scale-dependent innovations: the switch from LSTM to Kaplan Transformers and the rebalancing to Chinchilla scaling. Of the total measured progress of $21,400\times$ (relative to an LSTM), we find that $846\times$ is achieved through LSTMs to Kaplan Transformers, and nearly $10\times$ is attributable to Chinchilla rebalancing. Together, these comprise $91\%$ of total relative efficiency gains, the vast majority of our measured progress. These results are stylized in Figure 6.

Though our experiments do not claim to be exhaustive, we compare our findings with estimates from the literature. Namely, between 2012 to 2023, Ho et al. [2024] found a doubling time of 8 months, or $2.83\times$ per year, for a total efficiency gain of $22,000\times$. In contrast, the growth rate of our CEG multiplier is approximately $2.23\times$ annually, for a total of $6,930\times$, of which $2,700\times$ ($89\%$) is due to scale-dependent changes. This leaves a gap of $3.18\times$ from our estimates, which could be from data selection, tokenizer advancements, or a long tail of innovations not captured in our analysis.

## 5.2 CEG Multipliers Depend Strongly on Reference Algorithms

Our findings decomposing algorithmic progress suggest that scale-dependent changes may have disproportionate effects on the overall measured progress, contributing orders of magnitude more than their scale-invariant counterparts. *However, a deeper problem exists once scale-dependent algorithms are introduced.*

Consider a sequence of Modern Transformer models $M_1, M_2, ..., M_t$, each enhanced with a mixture-of-experts architecture, and trained with exponentially increasing compute. If we want to measure the growth rate of CEG multipliers (i.e., the rate of algorithmic progress), we need to choose a reference algorithm, such as baseline LSTMs or dense transformers. Measuring the CEG multiplier
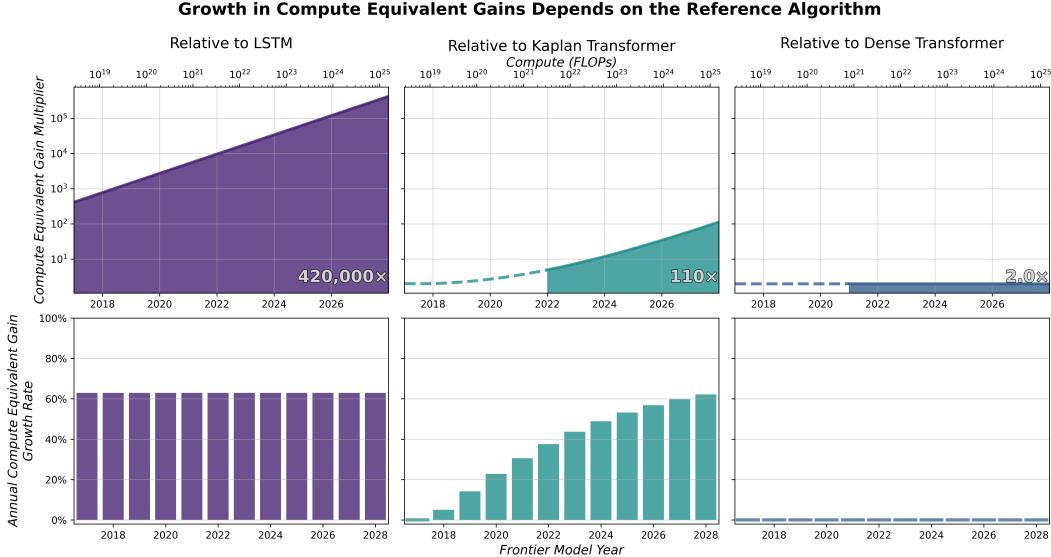
Figure 7: The right figure depicts compute equivalent gain multiplier growth rate of a Modern Transformer with respect to different reference algorithms extrapolated to 2028. Depending on the choice, growth rates vary dramatically. The left figure depicts the overall CEG gain with respect to different reference algorithms. Years are aligned with approximate frontier model sizes.

with respect to LSTMs for each model $M_t$ yields an exponential in $t$, with a growth rate of about 63% annually (Figure 7). However, measuring with respect to the dense Transformers yields a constant multiplier of $2\times$ (Section 3.3), and therefore a 0% growth rate. Thus choosing one reference point yields an exponential growth rate in algorithmic efficiency, while a different reference point yields *zero growth*. Figure 7 also shows the apparent growth in CEG multiplier from the Kaplan-Chinchilla transition, showing that these innovations can even appear to have variable growth rates over compute scales.

Only once scale-dependent algorithms are considered, selecting different reference algorithms for computing CEG multipliers can determine whether growth appears exponential, non-existent, or somewhere in between. In this sense, algorithmic progress with scale-dependent innovations is strongly contextual: even if you are completely certain you will measure exponential growth rates in the future with respect to some *fixed algorithm*, this offers no guarantee of any progress at all *relative to current models*. Importantly, this *is not* merely the claim that "we cannot be assured algorithmic progress will continue." Rather it is a statement about the measurement framework: **one can continue to *measure* arbitrarily large rates of progress from some fixed point while realizing none of these gains in practice.**

## 6 Limitations and Further Work

### 6.1 Impacts of Data Selection

We do not do any dataset ablation experiments. In general, it is hard to compare performance between datasets as perplexity is a dataset-relative benchmark. However, it is generally thought that datasets play an important role in determining neural scaling law exponents [Michaud et al., 2023]. Sorscher et al. [2022] found that neural scaling laws can even be exponential rather than power laws under the right data pruning metric. Nevertheless, in practice, while a better dataset may improve power laws, the difference between most datasets on evaluation metrics is moderate. Li et al. [2024] did a large-scale study comparing pretraining efficiency on new and historical datasets. They claim a $10\times$ efficiency gain for their DCLM baseline dataset but most common web datasets are within a factor of $3\times$ FLOP efficiency of one another in their study at scales of $10^{20}$ to $10^{22}$ FLOPs.

## 6.2 Our Experiments are Relatively Small and Non-Exhaustive

Our experiments are conducted at small scales compared to more recent scaling studies (i.e., Hoffmann et al. [2022]). However, our studies are done at slightly lower but similar scales to papers that originally explored many of these innovations, like Melis et al. [2017].

Finally, there are many things we did not test experimentally. These include mixture of experts, tokenizers, and improvements in data. These are important contributors to algorithmic progress, but we are only able to give literature-based estimates. We also do not test innovations involved in the creation of LSTMs. However, many LSTM variations have little benefit [Greff et al., 2016], and we use a generic architecture that would have been available in 2012. Further, we did not run scaling interventions on all the algorithms in our study individually, though when we ablate all algorithmic improvements, we find very small exponent differences in scaling exponents.

## 6.3 The Burden of Optimality

Another key limitation is the fickle nature of hyperparameters. For instance, we found it nearly impossible to scale a transformer trained with SGD and a batch size over 128. We are more confident in our results because our performance and exponents aligns well with other studies like Porian et al. [2024] and Melis et al. [2017]. Notably, Porian et al. [2024] conducts chinchilla optimal scaling and finds scaling exponents from .092 to .106. Using the same tokenizer but slightly different data than Porian et al. [2024], we achieve a scaling exponent of .094. Hoffmann et al. [2022] achieves a scaling exponent of .154 but uses both a different tokenizer and different data.

However, our study leaves open the possibility of a "golden" hyperparameters that could significantly lower the gap between the algorithmic improvements we measure or even change the experimental scaling relationship. For instance, SGD was long thought to be inherently inferior to Adam for transformer training, but studies like Srećković et al. [2025] demonstrate a very small performance gap with the right hyperparameters. We found very little literature on scaling LSTMs. Therefore, we could find improved scaling for LSTMs if extensive hyperparameter sweeps were done at all scales. This becomes computationally infeasible with much larger models, as determining optimal scaling for transformers took many years of progress in the machine learning community. Finding such an optimal scaling for other architectures and settings might require similar effort.

## 6.4 Beyond Data and Architecture: Reasoning Models

Reasoning-optimized language models have recently emerged as a new pathway for improving capabilities. Over the past year, model developers have been increasing AI capability by optimizing models to "think longer". Ho and Berg [2025] estimates that the shift from conventional to reasoning models was equivalent to a $10\times$ gain in pretraining compute. However, utilizing reasoning models is more costly due to the dramatically larger inference costs, which adds an additional dimension to measures of efficiency gains. Such improvements fall outside of our study's original 2012-2023 time period. Yet, it illustrates that methods beyond conventional architectural or data improvements can radically improve language model capabilities.

## 6.5 Beyond FLOP Efficiency: The Multidimensional Nature of Progress

Our focus on training FLOP efficiency necessarily omits many critical innovations. Numerous improvements we tested (RMSNorm, learning rate schedules, etc.) showed minimal FLOP efficiency gains but substantially improved training stability, convergence reliability, or extrapolation capabilities. Other innovations like Flash Attention reduce wall-clock time and energy consumption without affecting total FLOP operations. Capability-enabling innovations like instruction fine-tuning [Wei et al., 2021] and constitutional AI fundamentally change model utility without improving pretraining efficiency. There are also technical improvements that make large-scale training feasible in the first place, including pipeline and tensor parallelism techniques [Barnett, 2025].

These observations suggest that algorithmic progress in AI may be only partially characterized by scalar FLOP efficiency metrics alone. In analogy to transportation, AI algorithmic improvements are less well described as increases in fuel efficiency but instead as new types of vehicles and new roads to travel on. A comprehensive understanding of algorithmic progress requires multidimensional evaluation across efficiency, stability, capability, and deployment axes.

# 7 Discussion

Our experimental and theoretical analysis reveals that algorithmic progress in language models exhibits fundamentally different behavior across compute scales, with implications for understanding both historical progress and future trajectories.

## 7.1 Reconciling Small-Scale Measurements with Large-Scale Progress

At performance thresholds comparable to Transformers trained with $10^{15}$ FLOPs (i.e NLL= 5.3), we measure total algorithmic efficiency gains of approximately $6.28\times$ relative to baseline LSTMs when accounting for all tested innovations. Combined with literature estimates – MoE estimated at ($\sim 2\times$), improved datasets at ($\sim 3\times$), and tokenizers at ($\sim 1.6\times$) – this totals $60.3\times$, well under $100\times$. This contrasts sharply with Ho et al. [2024]'s estimate of $22,000$ efficiency gains the last decade, leaving an apparent gap of $350\times$ or more.

However, this discrepancy resolves when accounting for scale-dependent innovations. Extrapolating our measured scaling exponents to the 2023 frontier compute budgets ($1.3\times10^{22}$ FLOPs), the LSTM to Kaplan Transformer transition alone accounts for $725\times$ of the efficiency gains. Combined with the Chinchilla rebalancing gains of $3.7\times$ and $2.6\times$ for scale-invariant innovations, we estimate total efficiency gains of approximately $6,930\times$, substantially closing the gap with literature estimates over the same time period.

Further, our extrapolated estimates of LSTM to transformer efficiency gain in 2017 are around $91\times$. This is slightly larger than previous estimates from the literature at these scales. Ho et al. [2024] attributes a $60\times$ gain to this transition. Sanderson et al. [2025] estimates that this transition had a $20 - 50\times$ at scales of around $10^{19}$ FLOPs.

Critically, our framework exposes that "algorithmic progress" as a single number is ill-defined without specifying both the reference algorithm and target compute scale. The same sequence of innovations yields $100\times$ gains measured at small scale, but $10,000\times$ gains at frontier scale relative to the same LSTM baseline. This scale-dependence fundamentally challenges conventional metrics of algorithmic efficiency.

### 7.1.1 Algorithmic Versus Hardware Improvements

Algorithmic progress is traditionally seen as much faster than hardware progress. For instance, Ho et al. [2024] estimates algorithmic progress of $22,000$ over a ten year period while Moore's law over a ten year period would only amount to efficiency gains of $32\times$. However, at small scales, for instance using $10^{18}$ FLOPs algorithmic, progress is only $20\times$ – less than hardware progress (See Figure 6).

## 7.2 The Concentration of Progress in Architectural Transitions

Our decomposition reveals that algorithmic progress exhibits extreme concentration: the LSTM-to-Transformer architectural shift comprises $68\%$ of measured efficiency gains at frontier scales, with the Kaplan-to-Chinchilla rebalancing accounting for most of the remainder. The distribution of scale-invariant innovations is highly skewed, with a handful of improvements approaching or exceeding $2\times$, while most innovations contribute less than $1.5\times$.

This concentration has two important implications. First, algorithmic progress appears more punctuated than smooth exponential trends would suggest, with long periods of incremental refinement interrupted by rare architectural transitions. Second, the small magnitude of most scale-invariant improvements ($< 10\times$ cumulatively) suggests that future progress may similarly depend on discovering fundamentally new architectures rather than incremental refinements of existing ones.

## 7.3 Implications for the Future of AI

Many forecasts of AI capabilities assume smooth exponential gains in training efficiency over time [Kokotajlo et al., 2025, Erdil et al., 2025] regardless of compute level. Our results suggest three critical corrections to this view.

### 7.3.1 Larger Players Gain More From Algorithms

If efficiency improvements are scale-dependent, then algorithmic progress have not been uniform across small and large algorithmic providers. Model trainers with a large computational budget see much greater improvements with the introduction of new algorithms than smaller builders. Model builders scaling at the frontier see gains in the tens of thousands, while smaller model builders see algorithmic progress at levels much more similar to hardware progress. This mimics progress in traditional algorithms as seen in Sherry and Thompson [2021]. This means that new architectures or scaling paradigms could accelerate rather than reduce inequality between model producers.

### 7.3.2 Reference-dependence of Progress Measurements.

In the presence of scale-dependent algorithms, the measured growth rate of algorithmic efficiency depends entirely on the choice of reference algorithm. For example, computing CEG multipliers for a sequence of modern, MoE-enhanced transformers at the compute frontier yields $63\%$ annual growth relative to LSTMs. Computing CEG multipliers for this same sequence of models yields $0\%$ growth with respect to dense transformers. Once scale dependence is introduced, gains in algorithmic efficiency over time require first fixing an algorithm for reference. Throughout our analysis, we predominantly compare models to compute-optimal LSTMs to compare with literature estimates. However, existing estimates for growth in algorithmic progress may not accurately represent the pace of progress with respect to more recent algorithms, even at the frontier.

### 7.3.3 Difficulty of Early Identification

Scale-dependent innovations may appear modest at small scales, making it difficult to identify transformative improvements before committing substantial compute resources. The LSTM-to-Transformer transition showed only $6\times$ gains at $10^{15}$ FLOPs but more than $100\times$ gains at $10^{23}$ FLOPs. This suggests that limits to compute scaling pose obstacles not only to realizing efficiency gains but also to discovering them. This opens the possibility of algorithms that decrease performance at small scales, but which would improve performance if we scaled. Therefore, researchers must examine the scale-dependent nature of their improvements to identify the real impact of their innovations.

### 7.3.4 Will Algorithmic Progress Slow Down?

These findings have important implications for society and AI governance. If efficiency gains primarily arise from scale-dependent innovations, then limits to compute scaling—whether from energy constraints, semiconductor supply chains, or regulatory restrictions—may substantially slow AI algorithmic progress. In addition, if algorithmic progress cannot be meaningfully separated from compute investment, then the current exponential growth in AI capabilities may be more brittle than commonly assumed, requiring sustained increases in both computational resources and algorithmic breakthroughs.

## Acknowledgements

## References

Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Linear attention is (maybe) all you need (to understand transformer optimization). *arXiv preprint arXiv:2310.01082*, 2023.

Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, et al. Tokenizer choice for llm training: Negligible or crucial? In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3907–3924, 2024.

Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27), June 2024. ISSN 1091-6490. doi: 10.1073/pnas.2311878121. URL http://dx.doi.org/10.1073/pnas.2311878121.

Yamini Bansal, Behrooz Ghorbani, Ankush Garg, Biao Zhang, Colin Cherry, Behnam Neyshabur, and Orhan Firat. Data scaling laws in nmt: The effect of noise and architecture. In *International Conference on Machine Learning*, pages 1466–1482. PMLR, 2022.

Peter Barnett. Compute requirements for algorithmic innovation in frontier ai models. *arXiv preprint arXiv:2507.10618*, 2025.

Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xlstm: Extended long short-term memory. *Advances in Neural Information Processing Systems*, 37:107547–107603, 2024.

Tamay Besiroglu, Ege Erdil, Matthew Barnett, and Josh You. Chinchilla scaling: A replication attempt. *arXiv preprint arXiv:2404.10102*, 2024.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024.

Stella Biderman, Sid Black, Charles Foster, Leo Gao, Eric Hallahan, Horace He, Ben Wang, and Phil Wang. Rotary embeddings: A relative revolution. blog.eleuther.ai/, 2021. [Online; accessed ].

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Dokook Choe, Rami Al-Rfou, Mandy Guo, Heeyoung Lee, and Noah Constant. Bridging the gap for tokenizer-free language models. *arXiv preprint arXiv:1908.10322*, 2019.

Tom Davidson, Jean-Stanislas Denain, Pablo Villalobos, and Guillem Bas. Ai capabilities can be significantly improved without expensive retraining, 2023. URL https://arxiv.org/abs/2312.07413.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.

Nolan Dey, Bin Claire Zhang, Lorenzo Noci, Mufan Li, Blake Bordelon, Shane Bergsma, Cengiz Pehlevan, Boris Hanin, and Joel Hestness. Don't be lazy: Completep enables compute-efficient deep transformers. *arXiv preprint arXiv:2505.01618*, 2025.

Jasha Droppo and Oguz Elibol. Scaling laws for acoustic models. *arXiv preprint arXiv:2106.09488*, 2021.

Epoch AI. Key trends and figures in machine learning, 2023. URL https://epoch.ai/trends. Accessed: 2025-11-03.

Epoch AI. Data on ai models, 07 2025. URL https://epoch.ai/data/ai-models. Accessed: 2025-10-14.

Ege Erdil and Tamay Besiroglu. Algorithmic progress in computer vision. *arXiv preprint arXiv:2212.05153*, 2022.

Ege Erdil, Andrei Potlogea, Tamay Besiroglu, Edu Roldan, Anson Ho, Jaime Sevilla, Matthew Barnett, Matej Vrzla, and Robert Sandler. Gate: An integrated assessment model for ai automation. *arXiv preprint arXiv:2503.04941*, 2025.

Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. Lstm: A search space odyssey. *IEEE transactions on neural networks and learning systems*, 28(10): 2222–2232, 2016.

Xu Owen He. Mixture of a million experts. *arXiv preprint arXiv:2407.04153*, 2024.

Danny Hernandez and Tom B Brown. Measuring the algorithmic efficiency of neural networks. *arXiv preprint arXiv:2005.04305*, 2020.

Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*, 2017.

Anson Ho and Arden Berg. Quantifying the algorithmic improvement from reasoning models, 2025. URL https://epoch.ai/gradient-updates/quantifying-the-algorithmic-improvement-from-reasoning-models. Accessed: 2025-11-17.

Anson Ho, Tamay Besiroglu, Ege Erdil, David Owen, Robi Rahman, Zifan C Guo, David Atkinson, Neil Thompson, and Jaime Sevilla. Algorithmic progress in language models. *Advances in Neural Information Processing Systems*, 37:58245–58283, 2024.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022. URL https://arxiv.org/abs/2203.15556.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. URL https://arxiv.org/abs/2001.08361.

Nitish Shirish Keskar and Richard Socher. Improving generalization performance by switching from adam to sgd. *arXiv preprint arXiv:1712.07628*, 2017.

Daniel Kokotajlo, Scott Alexander, Thomas Larsen, Eli Lifland, and Romeo Dean. AI 2027: A scenario forecasting the impact of superhuman ai over the next decade. Website, 2025. URL https://ai-2027.com/. Published April 3, 2025.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.

Houyi Li, Ka Man Lo, Ziqi Wang, Zili Wang, Wenzhen Zheng, Shuigeng Zhou, Xiangyu Zhang, and Daxin Jiang. Can mixture-of-experts surpass dense llms under strictly equal resources? *arXiv preprint arXiv:2506.12119*, 2025.

Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training sets for language models. *Advances in Neural Information Processing Systems*, 37:14200–14282, 2024.

Jingyuan Liu, Jianlin Su, Xingcheng Yao, Zhejun Jiang, Guokun Lai, Yulun Du, Yidao Qin, Weixin Xu, Enzhe Lu, Junjie Yan, et al. Muon is scalable for llm training. *arXiv preprint arXiv:2502.16982*, 2025.

Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle, James Halverson, Marin Soljačić, Thomas Y Hou, and Max Tegmark. Kan: Kolmogorov-arnold networks. *arXiv preprint arXiv:2404.19756*, 2024.

Gábor Melis, Chris Dyer, and Phil Blunsom. On the state of the art of evaluation in neural language models. *arXiv preprint arXiv:1707.05589*, 2017.

Gábor Melis, Tomáš Kočiskỳ, and Phil Blunsom. Mogrifier lstm. *arXiv preprint arXiv:1909.01792*, 2019.

Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing lstm language models. *arXiv preprint arXiv:1708.02182*, 2017.

Eric Michaud, Ziming Liu, Uzay Girit, and Max Tegmark. The quantization model of neural scaling. *Advances in Neural Information Processing Systems*, 36:28699–28722, 2023.

Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36:50358–50376, 2023.

Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, et al. Olmoe: Open mixture-of-experts language models. *arXiv preprint arXiv:2409.02060*, 2024.

Tim Pearce and Jinyeop Song. Reconciling kaplan and chinchilla scaling laws. *arXiv preprint arXiv:2406.12907*, 2024.

Tomer Porian, Mitchell Wortsman, Jenia Jitsev, Ludwig Schmidt, and Yair Carmon. Resolving discrepancies in compute-optimal scaling of language models. *Advances in Neural Information Processing Systems*, 37:100535–100570, 2024.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

Albert Reuther, Jeremy Kepner, Chansup Byun, Siddharth Samsi, William Arcand, David Bestor, Bill Bergeron, Vijay Gadepally, Michael Houle, Matthew Hubbell, Michael Jones, Anna Klein, Lauren Milechin, Julia Mullen, Andrew Prout, Antonio Rosa, Charles Yee, and Peter Michaleas. Interactive supercomputing on 40,000 cores for machine learning and data analysis. In *2018 IEEE High Performance extreme Computing Conference (HPEC)*, page 1–6. IEEE, 2018.

Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.

Jonathan S Rosenfeld. Scaling laws for deep learning. *arXiv preprint arXiv:2108.07686*, 2021.

Jack Sanderson, Teddy Foley, Spencer Guo, Anqi Qu, and Henry Josephson. Rethinking llm advancement: Compute-dependent and independent paths to progress. *arXiv preprint arXiv:2505.04075*, 2025.

Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.

Yash Sherry and Neil C Thompson. How fast do algorithms improve?[point of view]. *Proceedings of the IEEE*, 109(11):1768–1777, 2021.

Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems*, 35:19523–19536, 2022.

Teodora Srećković, Jonas Geiping, and Antonio Orvieto. Is your batch size the problem? revisiting the adam-sgd gap in language modeling. *arXiv preprint arXiv:2506.12543*, 2025.

Department of Computer Science Stanford University. Lecture 3: Architecture. https://github.com/stanford-cs336/spring2025-lectures/blob/e9cb2488fdb53ea37f0e38924ec3a170192cef3/nonexecutable/2025%20Lecture%203%20-%20architecture.pdf, March 2025. URL https://stanford-cs336.github.io/spring2025/index.html#coursework. Lecture slides for course CS 336 (Spring 2025).

Yi Tay, Mostafa Dehghani, Samira Abnar, Hyung Chung, William Fedus, Jinfeng Rao, Sharan Narang, Vinh Tran, Dani Yogatama, and Donald Metzler. Scaling laws vs model architectures: How does inductive bias influence scaling? In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12342–12364, 2023.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.

Parker Whitfill. Note on selection bias in observational estimates of algorithmic progress. *arXiv preprint arXiv:2508.11033*, 2025.

Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *International conference on machine learning*, pages 10524–10533. PMLR, 2020.

Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in neural information processing systems*, 32, 2019.

Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George Dahl, Chris Shallue, and Roger B Grosse. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model. *Advances in neural information processing systems*, 32, 2019.

Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhiquan Luo. Why transformers need adam: A hessian perspective. *Advances in neural information processing systems*, 37:131786–131823, 2024.

Rosie Zhao, Depen Morwani, David Brandfonbrener, Nikhil Vyas, and Sham Kakade. Deconstructing what makes a good optimizer for language models. *arXiv preprint arXiv:2407.07972*, 2024.

# A    Model Architecture

## A.1    Learning Rate Tunes and Extrapolation

To scale both our transformers and LSTM, we run a learning rate tune at least 4 model sizes over the range: $[10^{-3}, 10^{-2.5}, 10^{-2}, 10^{-1.5}, 10^{-1}]$. We then do a more fine-grained learning rate tune in intervals of $10^{0.25}$, ensuring that equal hyperparameter tuning resources are used for all our models (Modern Transformer, Retro Transformer, LSTM). Afterward, we fit a power-law trend between the learning rate and model hidden dimension and use that trend to determine the optimal learning rate at all sizes.

## A.2    Further Discussion of Hyperparameter Choices

We choose to use no dropout for both our LSTM and Transformer as both perform considerably better without dropout. To check this for LSTM we compared an LSTM with hidden dimension 64 with the dropout settings recommended by Melis et al. [2017] to an LSTM with no dropout and tune learning rates for both. The best performing LSTM with dropout has a validation loss of 6.5, while the best performing LSTM without dropout has a validation loss of 5.9.

## A.3    Data Selection, Loading, and Size

We choose to use a subsample of the C4 dataset for both the transformer and LSTM. The size of dataset we use for training is based on the model size and our token-to-parameter ratio. In order to better outline the optimal compute frontier, we choose to use a token-to-parameter ratio of 40 in our scaling experiments. This is to ensure that we do not stop early before compute-optimal level of data. Stopping early would change the frontier envelope, while stopping late would not change the frontier, as overtrained models would pass the frontier. While for our ablation experiments, we use a token to parameter ratio of 20 as recommended by Hoffmann et al. [2022].

Table 1: Base hyperparameters used to train the Modern Transformer Model. Taken from Porian et al. [2024] and optimal recommendations from Stanford University [2025].

| Hyperparameter | Value |
|---|---|
| Width/Depth Aspect Ratio | 16 |
| Norm type | RMSNorm prenorm |
| Feedforward / model dimension | 4 |
| Activation Function | SwiGLU |
| Positional Encoding | Rotary |
| Sequence length | 128 |
| Stride | 128 |
| Batch size | 64 |
| Dropout | 0.0 |
| Weight decay | 0.01 |
| Gradient Clipping | 1.0 |
| Learning Rate Schedule | cosine annealing |
| Min learning rate | 0.1x max lr |
| Warmup fraction | 10% |
| Default Optimizer | AdamW |
| Initialization | GPT-2/BERT init |
| Embedding Weight Tying | True |
| Tokenizer | GPT-2 Tokenizer |
| Vocabulary | 50257 |
| Dataset | C4 |
| Validation Set Size | 500k tokens |
| Token to parameter ratio | 40 |

Table 2: Architecture and hyperparameters for our Retro Transformer. This is the same in Table 2 but with many algorithmic choices reverted.

| Hyperparameter | Value |
|---|---|
| Width/Depth Aspect Ratio | 16 |
| Norm type | layer postnorm |
| Feedforward / model dimension | 4 |
| Activation Function | GeLU |
| Positional Encoding | Sinusoidal |
| Sequence length | 128 |
| Stride | 128 |
| Batch size | 64 |
| Dropout | 0.0 |
| Weight decay | 0.01 |
| Gradient Clipping | 1.0 |
| Learning Rate Schedule | inverse square root decay |
| Min learning rate | 0.1x max lr |
| Warmup fraction | 10% |
| Default Optimizer | Adam |
| Initialization | GPT-2/BERT init |
| Embedding Weight Tying | True |
| Tokenizer | GPT-2 Tokenizer |
| Vocabulary | 50257 |
| Dataset | C4 |
| Validation Set Size | 500k tokens |
| Token to parameter ratio | 40 |

Table 3: Base hyperparameters used to train LSTM Model, mostly taken from Melis et al. [2017].

| Hyperparameter | Value |
| --- | --- |
| Number of layers | 1 |
| TBPTT length | 32 |
| TBPTT stride | 32 |
| Sequence length | 128 |
| Batch size | 64 |
| Input dropout | 0.0 |
| Hidden dropout | 0.0 |
| Output dropout | 0.0 |
| Weight decay | $10^{-4}$ |
| Gradient Clipping | 1.0 |
| Learning Rate Schedule | cosine annealing |
| Min learning rate | 0.1x max lr |
| Warmup fraction | 10% |
| Default Optimizer | AdamW |
| Initialization | LSTM (orthogonal hidden matrices) |
| Embedding Weight Tying | True |
| Tokenizer | GPT-2 Tokenizer |
| Vocabulary | 50257 |
| Dataset | C4 |
| Validation Set Size | 500k tokens |
| Token to parameters ratio | 40 |

## B Kaplan-Style Scaling When the True Frontier is Chinchilla-Optimal

In Figure 5, we show loss curves from using Kaplan and Chinchilla's recommendations for compute optimal scaling (specifically optimal scaling by analytically solving the full parameter model). We take the Chinchilla functional form for model loss $L(N, D)$ as ground-truth, only modifying the allocation of parameters and dataset size between each curve. Unsurprisingly, Chinchilla scaling is a perfect power law.

We use the constants $E = 1.69$, $A = 406.4$, $B = 410.7$, $\alpha = .34$, and $\beta = .28$ and include the functional form below for the reader's convenience.

$$L_{\text{Chinchilla}}(N, D) = \frac{A}{N^\alpha} + \frac{B}{D^\beta} + E \tag{1}$$

$$= \frac{406.4}{N^{0.34}} + \frac{410.7}{D^{0.28}} + 1.69 \tag{2}$$

For the compute allocation, Chinchilla scaling explicitly defines optimal scaling of $N$ and $D$ to be

$$N_{\text{Chinchilla}}(C) = G\left(\frac{C}{6}\right)^{\beta/(\alpha+\beta)} \tag{3}$$

$$D_{\text{Chinchilla}}(C) = \left(\frac{1}{G}\right)\left(\frac{C}{6}\right)^{\alpha/(\alpha+\beta)} \tag{4}$$

where $G = \left(\frac{\alpha A}{\beta B}\right)^{1/(\alpha+\beta)}$. By construction, this gives a precise power law with a well defined inverse:

$$\text{ChinchillaCompute}(L) = \frac{5.4 \times 10^{19}}{(L - 1.69)^{6.5}} \tag{5}$$

Finally, we infer the full Kaplan scaling recommendation: Figure 14 in Kaplan et al. [2020] gives a functional form for compute optimal dataset scaling, and after applying unit conversions and using the relationship $C = 6ND$, we get the following:

$$N_{\text{Kaplan}}(C) = (3.6 \times 10^{-6}) C^{0.73} \tag{6}$$

$$D_{\text{Kaplan}}(C) = (4.6 \times 10^{4}) C^{0.27} \tag{7}$$

Together, these give us a closed form for the compute multiplier when rebalancing from Kaplan to Chinchilla scaling for a fixed $C$. We call this function $M$, and compute it as follows:

$$M(C) = \frac{C}{\text{ChinchillaCompute}(L_{\text{Chinchilla}}(N_{\text{Kaplan}}(C), D_{\text{Kaplan}}(C)))} \tag{8}$$

$$= (1.85 \times 10^{-20})(C) \left( \frac{28737.9}{C^{0.2482}} + \frac{20.337}{C^{0.0756}} \right)^{6.512} \tag{9}$$

As $C$ gets large (roughly $10^{23}$ FLOPs), we can use a power law to approximate this expression:

$$M(C) = (6.13 \times 10^{-12})(C^{0.508})$$

## C  Transformer SGD Scaling Study

In order to train a Transformer with SGD we use a momentum of .98 and 0 weight decay as recommended by Srećković et al. [2025]. We use the same scaling procedure as our default Transformer, and all other hyperparameters are held constant. We find little scaling difference between Transformers trained with SGD and Transformers trained with AdamW. However, in our experiments, transformers trained with SGD are notably more unstable. For instance, our SGD training curve have a slight concavity early on.
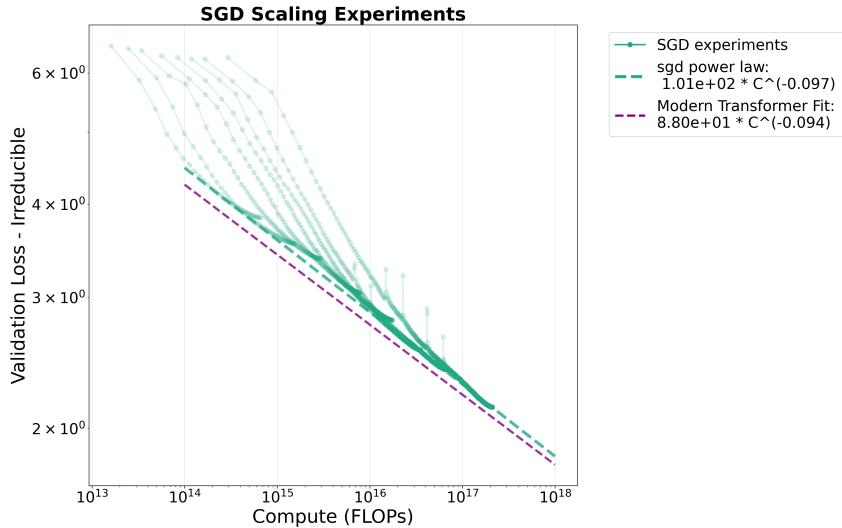


Figure 8: Optimizers have little effect on scaling exponents. The blue dashed line represents the compute-optimal scaling fit for our Modern Transformer trained with SGD. The purple dashed line represents the experimental fit for our Transformer trained with AdamW.

## D  Irreducible Loss

The irreducible loss represents the self-entropy of the data distribution. This means it varies from dataset to dataset but should be constant across models. In our paper, we set the irreducible loss to 1.9. In contrast, Hoffmann et al. [2022] fits an irreducible loss of 1.69 using MassiveText (similar to C4) Besiroglu et al. [2024] finds an irreducible loss of 1.8 using Hoffmann et al. [2022] data. Porian et al. [2024] did a replication of Kaplan and Chinchilla scaling on RefinedWeb and OpenWebText2.

They find irreducible loss fits from 1.68 to 2.01. Muennighoff et al. [2023] finds an irreducible loss of 1.87 for Transformers trained on C4. In light of this variation, we use scikit learn to fit the optimal frontier to the form $L = E + AC^{-\alpha}$, restricting E to lie between 1.3 and 2.2. For our Modern Transformer, 2017 transformer, and LSTM, this approach yields an irreducible loss close to 1.9 for all these models.

# E    Reproducibility Statement

All code used in model training and model analysis, along with setup instructions, is hosted here: https://github.com/hansgundlach/Experimental_Progress. The main body of our text contains details on our training and analysis procedure. A detailed list of hyperparameters used is in Appendix A. All experiments were done using 8 V100 GPUs on MIT Supercloud [Reuther et al., 2018].