
Scaling Laws for Autoregressive Generative Modeling

Tom Henighan*

Jared Kaplan*†

Mor Katz*

Mark Chen

Christopher Hesse

Jacob Jackson

Heewoo Jun

Tom B. Brown

Prafulla Dhariwal

Scott Gray

Chris Hallacy

Benjamin Mann

Alec Radford

Aditya Ramesh

Nick Ryder

Daniel M. Ziegler

John Schulman

Dario Amodei

Sam McCandlish

OpenAI

Abstract

We identify empirical scaling laws for the cross-entropy loss in four domains: generative image modeling, video modeling, multimodal image↔text models, and mathematical problem solving. In all cases autoregressive Transformers smoothly improve in performance as model size and compute budgets increase, following a power-law plus constant scaling law. The optimal model size also depends on the compute budget through a power-law, with exponents that are nearly universal across all data domains.

The cross-entropy loss has an information theoretic interpretation as $S(\text{True}) + D_{\text{KL}}(\text{True}||\text{Model})$, and the empirical scaling laws suggest a prediction for both the true data distribution’s entropy and the KL divergence between the true and model distributions. With this interpretation, billion-parameter Transformers are nearly perfect models of the YFCC100M image distribution downsampled to an 8×8 resolution, and we can forecast the model size needed to achieve any given reducible loss (ie D_{KL}) in nats/image for other resolutions.

We find a number of additional scaling laws in specific domains: (a) we identify a scaling relation for the mutual information between captions and images in multimodal models, and show how to answer the question “Is a picture worth a thousand words?”; (b) in the case of mathematical problem solving, we identify scaling laws for model performance when extrapolating beyond the training distribution; (c) we finetune generative image models for ImageNet classification and find smooth scaling of the classification loss and error rate, even as the generative loss levels off. Taken together, these results strengthen the case that scaling laws have important implications for neural network performance, including on downstream tasks.

*equal contribution

†Johns Hopkins University and OpenAI

Correspondence to: [henighan, jared, sam]@openai.com

Author contributions listed at end of paper.

Contents

1	Introduction	3
1.1	Summary of Results	5
2	Central Empirical Scaling Laws in Each Domain	6
2.1	Domain Descriptions and Training Setups	6
2.2	Model Size Scaling and Aspect Ratios	9
2.3	Compute Scaling and Optimal Model Sizes	9
2.4	Loss versus Position in the Context Depends on the Structure of the Data	11
3	Image and Video Modeling, the Reducible Loss, and Downstream Tasks	11
3.1	Varying the Image Resolution and Encoding	11
3.2	Video Modeling and Individual Frames	13
3.3	Scaling Trends for Individual Images	14
3.4	Finetuning on ImageNet at 32x32 Resolution	14
4	Multimodal Models and Information Gain	16
5	Mathematical Problem Solving and Extrapolation	16
6	An Inconsistency in Compute and Datasize Scaling Laws	17
7	Related Work	19
8	Discussion	20
A	More Details on Image Modeling	22
B	Details of Math Experiments and Additional Results	26
B.1	Procedurally Generated Training Data	26
B.2	Dataset Size Scaling	26
B.3	Additional Math Results	26
C	Additional Multimodal Results	30
D	Additional Language Results	31
E	Mutual Information, Infogain, and Scaling	33
E.1	Approximate Derivation of Scaling Relations	33
E.2	Estimating D_{KL} Between Real-World Distributions	33
F	Hyperparameter Settings	34

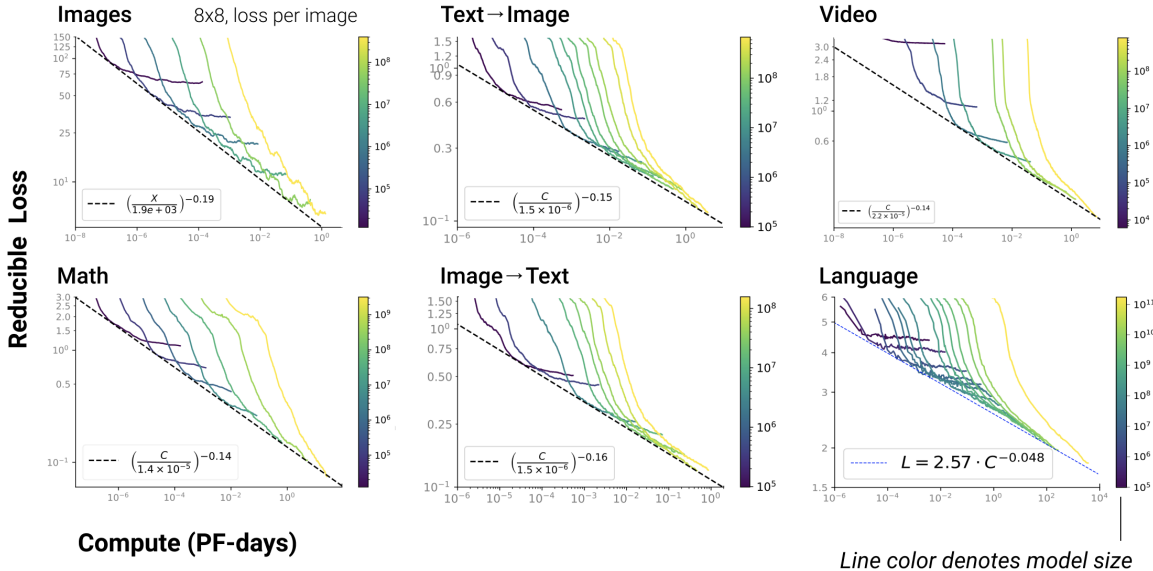


Figure 1 Smooth scaling of reducible loss across domains— We show power-law scaling laws for the *reducible* loss $L - L_\infty$ as a function of compute, where the irreducible loss L_∞ is a fitted domain-dependent constant. Under plausible assumptions concerning the infinite data and compute limits, the irreducible loss estimates the entropy of the underlying data distribution, while the reducible loss approximates the KL divergence between the data and model distributions. In the case of language we use results from [BMR⁺20], and only show the full loss L .

1 Introduction

Large scale models, datasets, and compute budgets have driven rapid progress in machine learning. Recent work [HNA⁺17, RRBS19, LWS⁺20, RDG⁺20, KMH⁺20, SK20, BMR⁺20] suggests that the benefits of scale are also highly predictable. When the cross-entropy loss L of a language model is bottlenecked by either the compute budget C , dataset size D , or model size N , the loss scales with each of these quantities as a simple power-law. Sample efficiency also improves with model size.

These results raise a number of questions. Do they apply to all data modalities? How do improvements on the loss translate to improvements in representation quality and performance on downstream tasks? Is there any way to determine when and why the performance of a model might be maxed out, so that further scaling will be met with diminishing returns? What explains the precision and universality of these trends, and what else can we learn from them?

We will demonstrate that scaling laws apply to generative modeling across a wide variety of data modalities, including generative language [KMH⁺20, BMR⁺20], image [TSF⁺15, CRC⁺20], and video modeling [WTU19], multimodal modeling [TBL⁺19] of text-image correlations, and even mathematical problem solving [SGHK19], a task requiring a degree of reasoning ability. Moreover, we demonstrate that a single architecture – the Transformer [VSP⁺17, LSP⁺18], with an autoregressive cross-entropy loss – scales smoothly in all of these domains, with only minimal changes to hyperparameters such as width, depth, or learning rate. We also observe that larger models consistently learn faster, achieving any given value of the loss in fewer steps.

By studying many different model sizes N , compute budgets C , or dataset sizes D , we demonstrate that the scaling relation for the loss

$$L(x) = L_\infty + \left(\frac{x_0}{x}\right)^{\alpha_x} \quad (1.1)$$

applies to each data modality, where α_x is a modality-dependent scaling exponent, and we primarily study $x = N, C$, and occasionally D . We will refer to L_∞ as the irreducible loss and the power-law scaling term as the reducible loss. These scaling relations often hold to high precision, even when the reducible loss is much smaller than the irreducible loss; we display trends in $L(C)$ for the reducible loss in figure 1. Note that small deviations are visually amplified on the log-plot, but nevertheless the trends fit remarkably well.

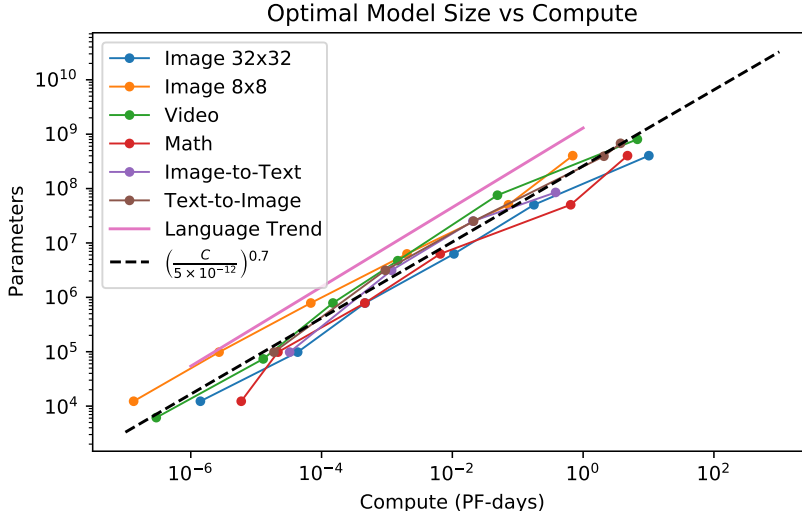


Figure 2 Optimal model size is consistent across domains— We display the optimal model size N_{opt} as a function of the training compute budget C . Not only does $N_{\text{opt}}(C)$ behave as a power-law, but the behavior is remarkably similar for all data modalities.

These observations suggest the information theoretic interpretation

$$\begin{aligned}
 L_\infty &\approx S(\text{True}) && \text{“Irreducible Loss”} \\
 \left(\frac{x_0}{x}\right)^{\alpha_x} &\approx D_{\text{KL}}(\text{True}||\text{Model}) && \text{“Reducible Loss”}
 \end{aligned} \tag{1.2}$$

In other words, the irreducible loss estimates the entropy of the true data distribution, while the reducible loss is an estimate of the KL divergence between the true and model distributions. One might have guessed that as the $L(x)$ curve bends and the loss approaches L_∞ , returns to increasing N, C, D are diminishing. But the identification of the reducible loss with D_{KL} suggests this is not necessarily the case, and further increases in scale may still provide important additional semantic information. To justify equation (1.2), we must assume that in the limit $D \rightarrow \infty$ followed³ by $N, C \rightarrow \infty$, an infinitely large transformer could model the data distribution exactly.

The scaling relations provide insights into the complexity of the data and clarify the value of increasing N, D , and C . By evaluating the reducible loss for a full image or video, we are actually estimating the number of bits of information that ‘remain to be understood’ by a given model. Equivalently, the reducible loss approximates the degree to which the data could be further compressed. We find that billion-parameter models can extract all but a few nats/image concerning YFCC100M images [TSF⁺15] downsampled to an 8x8 resolution, so they may be nearly perfect models of this data distribution. For larger, more practically relevant images we would need far larger models to achieve this feat, but the scaling laws make it possible to forecast this precisely. These trends are closely tied to the scaling exponents α_x : smaller exponents imply slower improvement with increasing scale, meaning that the data can only be compressed further with much larger models.

The scaling of loss with compute makes it possible to estimate the optimal model size for a given compute budget. We find that just as in [KMH⁺20] this relation is very nearly a pure power-law $N_{\text{opt}}(C) \propto C^\beta$. Surprisingly, the exponent $\beta \sim 0.7$ is very similar for all domains, as shown in figure 2. This has important implications for the scaling of dataset size with model size for compute-optimal training, suggesting that $D \propto N^{0.4}$ if we only train on each data element once. Even allowing for significant errors or deviations, this strongly suggests sub-linear scaling of dataset size with model size.

We can learn more if we focus on questions specific to each data modality. Generative image models can be finetuned for classification. We will show that ImageNet [CLH17] classification performance improves smoothly with pre-trained model size, following another power law. This trend continues even into the large-model regime where the generative loss trend “bends” and becomes dominated by the irreducible component. This strongly suggests that there are benefits to squeezing as much performance as possible out of large

³We specify this order of limits to make it clear that regularization will not be required.

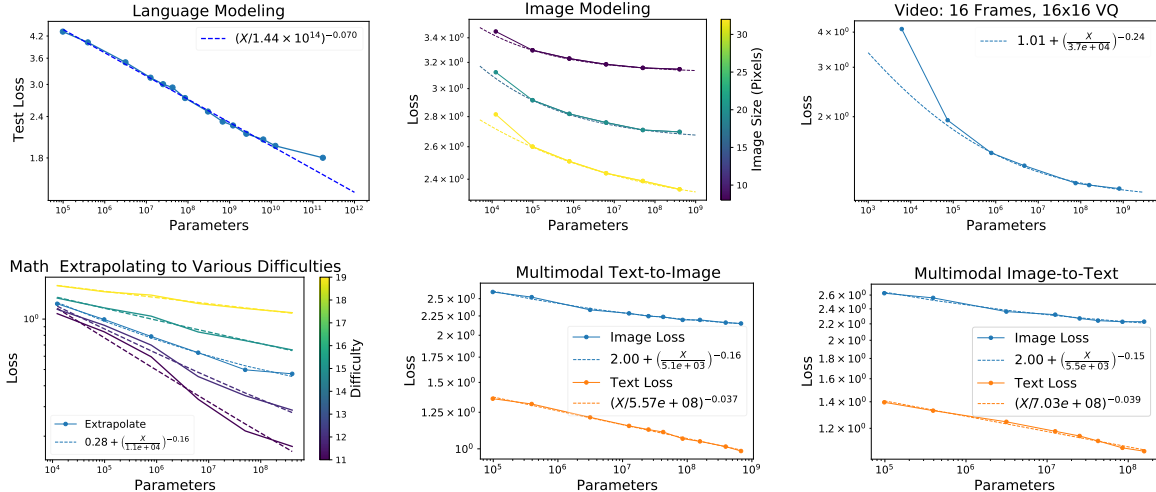


Figure 3 Scaling with model size— We show scaling laws with model size for various domains, along with fits (dashed) to equation (1.1). Note that the largest language models [BMR⁺20] in the top-left figure were not trained to convergence, so deviations from the trend are not necessarily meaningful. Very small models for video and higher-resolution images are off-trend; we speculate this is due to these models attempting to attend to a context with length comparable to their non-embedding parameter count.

generative image models, as significant semantic information may lie in the ‘last few bits’. The smooth trends for finetuned performance on image classification suggest a more general lesson: that the scaling laws for unsupervised learning imply that downstream performance also improves with model size and compute.

Information theory provides a useful lens for examining model performance in other contexts. A striking case is provided by multimodal models, such as those that model the joint distribution between text captions and images. Typically the entropy of the caption is much smaller than that of the image, so the ratio between the (empirical) mutual information⁴ and the model’s loss on the text, which we refer to as the

$$\text{Infogain} \equiv \frac{I(\text{text, image})}{L(\text{text})} \quad (1.3)$$

provides an interesting metric for model performance. The mutual information shared between distributions must be smaller than the amount of information in either distribution, so this ratio must be less than 1. Furthermore, it appears that the Infogain increases smoothly with model size, so that the bound $\text{Infogain} < 1$ can suggest a target model size for maximum performance. Typically this is far beyond current capabilities.

These smooth scaling results on a wide variety of datasets also demonstrate the remarkable versatility of the Transformer architecture.

1.1 Summary of Results

We apply autoregressive decoder-only Transformer models to all data modalities, which include web-scraped YFCC100M images [TSF⁺15] of various resolutions, video data from various sources, multimodal image+language data, and procedurally generated math problems. We also reference prior results on language [KMH⁺20, BMR⁺20]. **Across all domains** we find:

- The scaling laws of equation (1.1) apply consistently, including for very small values of the reducible loss. Since the $L(C)$ trends can be extended to arbitrarily large data distributions, model sizes, and training steps, we argue that this supports the interpretation of equation (1.2).
- We identify the optimal model size $N_{\text{opt}}(C)$ for a given compute budget, and find that it can be accurately modeled as a pure power law [KMH⁺20]

$$N_{\text{opt}} \propto C^\beta \quad (1.4)$$

⁴By the empirical mutual information we are referring to $\mathbb{E}_{x,y \sim q} \left[\log \frac{p(x,y)}{p(x)p(y)} \right]$ where p is the model distribution and q is the true distribution of the data. This must be smaller than the cross-entropy loss of the model on both X and Y .

with a power $\beta \sim 0.7$ for all modalities, as shown in figure 2. As compute budgets grow, it’s best to devote a majority of resources towards training larger models. This strongly suggests sub-linear scaling of $D \propto N^{0.4}$ for dataset size with model size during compute-optimal training.

- For each domain, there is an optimal aspect ratio $d_{\text{model}}/n_{\text{layer}}$ for the Transformer. Most data modalities require smaller aspect ratios (i.e. deeper networks) as compared to language [KMH⁺20].
- We study an apparent inconsistency between $L(D)$ and $L(C)$ trends in section 6.

We also find a number of results specific to certain domains, though we expect that many of the lessons are more general. For **image and video modeling** (see section 3):

- When generative image models are finetuned for ImageNet classification, we find a power-law for classification loss vs model size (see figure 11), even beyond the model size where we approach the irreducible loss for generative modeling. We conclude that **the approach to the irreducible loss does not necessarily indicate diminishing returns for representation quality or semantic content**.
- We explore scaling trends for individual images and for percentiles of the image loss distribution (see figures 17, 10, 20, 21). We find that the loss on individual images scales with model size in the same way as the mean over all images in the data distribution. We expect similar behavior in other data modalities.
- We test a variety of image resolutions (see figure 8), and find distinct scaling exponents and irreducible losses for each. We also test two VQVAE [vdOVK18] based models.
- We examine scaling of the loss with video frame index (see figures 6 and 9).

For **multimodal models** (see section 4):

- We explore the mutual information between captions and images (see figure 12), and the information gain defined in equation (1.3). We find a smooth scaling for both the mutual info and information gain with model size N .
- We revisit the question “Is a picture worth a thousand words?” by comparing the information-content of textual captions to the image/text mutual information.

For **mathematical problem solving** (see section 5 and appendix B):

- We explore the ability of models to extrapolate from the training distribution to increasingly more challenging problems. We find that extrapolation performance depends predominantly on performance on the training distribution (figure 24), and is otherwise independent of model size. So while larger models perform better, model size does not provide benefits to ‘strong generalization’.
- We provide a detailed breakdown of performance by math problem type (see appendix B).

2 Central Empirical Scaling Laws in Each Domain

In this section we will describe our common experiments in each domain and our results establishing equation (1.1) for compute, model size, and (in a few cases) dataset size scaling.

2.1 Domain Descriptions and Training Setups

In every domain we use decoder-only transformer models trained using an autoregressive cross-entropy loss. For many models we use a sparse attention pattern [CGRS19], though we use dense attention when solving math problems.

The transformers used for language and multimodal modeling have fully connected layers of size $4d_{\text{model}}$ and attention layers of size d_{model} , in the notation of [KMH⁺20, BMR⁺20]. For math, image, and video modeling we scale the FC layers to d_{model} and the attention layers to $d_{\text{model}}/4$. We use an aspect ratio $d_{\text{model}}/n_{\text{layer}} \approx 10$ for math, images, and videos as we find that this is approximately optimal, meaning that these domains prefer much deeper models as compared to language [KMH⁺20], where the optimal aspect ratio ~ 100 . Thus our math, image, and video models are essentially identical, differing only in context length. For math alone we used a weight decay [LH17] of 0.05. We provide more detailed hyperparameter settings in appendix F.

Domain	$L(N)$ (model size)	$L(C)$ (compute)	$N_{\text{opt}}(C)$
Language	$\left(\frac{N}{1.47 \times 10^{14}}\right)^{-0.070}$	$\left(\frac{C}{3.47 \times 10^8}\right)^{-0.048}$	$\left(\frac{C}{3.3 \times 10^{-13}}\right)^{0.73}$
Image 8x8	$3.12 + \left(\frac{N}{8.0 \times 10^1}\right)^{-0.24}$	$3.13 + \left(\frac{C}{1.8 \times 10^{-8}}\right)^{-0.19}$	$\left(\frac{C}{5.3 \times 10^{-14}}\right)^{0.64}$
Image 16x16	$2.64 + \left(\frac{N}{2.8 \times 10^2}\right)^{-0.22}$	$2.64 + \left(\frac{C}{1.6 \times 10^{-8}}\right)^{-0.16}$	$\left(\frac{C}{4.8 \times 10^{-12}}\right)^{0.75}$
Image 32x32	$2.20 + \left(\frac{N}{6.3 \times 10^1}\right)^{-0.13}$	$2.21 + \left(\frac{C}{3.6 \times 10^{-9}}\right)^{-0.1}$	$\left(\frac{C}{1.6 \times 10^{-13}}\right)^{0.65}$
Image VQ 16x16	$3.99 + \left(\frac{N}{2.7 \times 10^4}\right)^{-0.13}$	$4.09 + \left(\frac{C}{6.1 \times 10^{-7}}\right)^{-0.11}$	$\left(\frac{C}{6.2 \times 10^{-14}}\right)^{0.64}$
Image VQ 32x32	$3.07 + \left(\frac{N}{1.9 \times 10^4}\right)^{-0.14}$	$3.17 + \left(\frac{C}{2.6 \times 10^{-6}}\right)^{-0.12}$	$\left(\frac{C}{9.4 \times 10^{-13}}\right)^{0.7}$
Text-to-Im (Text)	$\left(\frac{N}{5.6 \times 10^8}\right)^{-0.037}$	(combined text/image loss)	$\left(\frac{C}{9.4 \times 10^{-13}}\right)^{0.7}$
Text-to-Im (Image)	$2.0 + \left(\frac{N}{5.1 \times 10^3}\right)^{-0.16}$	$1.93 + \left(\frac{C}{1.5 \times 10^{-6}}\right)^{-0.15}$	
Im-to-Text (Text)	$\left(\frac{N}{7.0 \times 10^8}\right)^{-0.039}$	(combined text/image loss)	$\left(\frac{C}{3.3 \times 10^{-12}}\right)^{0.72}$
Im-to-Text (Image)	$2.0 + \left(\frac{N}{5.5 \times 10^3}\right)^{-0.15}$	$1.97 + \left(\frac{C}{1.5 \times 10^{-6}}\right)^{-0.16}$	
Video VQ 16x16x16	$1.01 + \left(\frac{N}{3.7 \times 10^4}\right)^{-0.24}$	$0.95 + \left(\frac{C}{2.2 \times 10^{-5}}\right)^{-0.14}$	$\left(\frac{C}{1.13 \times 10^{-12}}\right)^{0.71}$
Math (Extrapolate)	$0.28 + \left(\frac{N}{1.1 \times 10^4}\right)^{-0.16}$	$0.14 + \left(\frac{C}{1.4 \times 10^{-5}}\right)^{-0.17}$	$\left(\frac{C}{2.3 \times 10^{-12}}\right)^{0.69}$

Table 1 Summary of scaling laws— In this table we summarize the model size and compute scaling fits to equation (1.1) along with $N_{\text{opt}}(C)$, with the loss in nats/token, and compute measured in petaflop-days. In most cases the irreducible losses match quite well between model size and compute scaling laws. The math compute scaling law may be affected by the use of weight decay, which typically hurts performance early in training and improves performance late in training. The compute scaling results and data for language are from [BMR⁺20], while $N_{\text{opt}}(C)$ comes from [KMH⁺20]. Unfortunately, even with data from the largest language models we cannot yet obtain a meaningful estimate for the entropy of natural language.

2.1.1 Language

We show results from GPT-3 [BMR⁺20] for comparison, including the performance of much larger models than we train in other domains. In figure 2 we use the optimal model size trend from [KMH⁺20]. In appendix D we show some experiments on the scaling of arithmetic and factual question answering abilities, and make some additional qualitative observations about the progression of language understanding with scale.

2.1.2 Images

We study a dataset of approximately 10^8 web images [TSF⁺15] scaled to pixel resolutions $R \times R = 8 \times 8$, 16×16 , and 32×32 represented in raster order using RGB colors, each in the range $[0, 255]$, giving a total of $3R^2$ tokens per image. We also study the same images at 64×64 resolution but VQ [vdOVK18] encoded with either a 16×16 or 32×32 VQ encoding pattern, for a total of either 256 or 1024 tokens per image. To reduce compute, we use sparse attention patterns [CGRS19], alternating between locally-banded attention and fixed-stride attention in sequential layers, where both the local context length and fixed-stride length are given by the side-length in tokens of the square images.

2.1.3 Video

We study a dataset of approximately 7×10^5 videos totaling about 100 hours scraped from the web, where each frame is scaled to a pixel resolution of 64×64 . Each individual frame is encoded with the same 16×16 VQVAE [vdOVK18] used for images, resulting in 256 tokens per frame. We train on sequences of 16 sequential frames, resulting in a total of 4096 tokens per video. As with images, we reduce compute by using a sparse attention pattern [CGRS19] alternating between locally-banded and fixed-stride attention, where both the local context length and fixed-stride length are given by the side length in tokens of the square frames.

Input Resolution (pixels)	Output Resolution (VQ Codes)	Codebook Size
64x64	16x16	4096
64x64	32x32	1024

Table 2 Details of VQVAEs used to encode images and frames of video.

2.1.4 VQ Encoding

The VQVAE models mentioned in 2.1.2 and 2.1.3 were trained on frames of the web-scraped videos described in 2.1.3, using the VQ-VAE architecture [vdOVK18] with modifications described in [DJP⁺20], including dead code revival. More details can be found in table 2.

2.1.5 Multimodal Text and Images

Multimodal models are trained to autoregressively predict both image tokens and language tokens in series. We simply concatenate together the token lists for BPE encoding of text (using the tokenization of [BMR⁺20]) and the $[0, 255]$ colorscale of each of the RGB pixels in the images, and let the model learn the necessary embedding matrix. We separately study models for text-to-image and image-to-text mappings, as we found poor performance for bidirectional models in preliminary experiments. For both image-to-text and text-to-image models we compute the mean pixel and mean text token loss, and then weight them to form the total loss $L = 9L_{\text{image}} + L_{\text{text}}$, as we found this weighting produced good results in a scan. We use 32x32 images together with a 128-token captions (padded or trimmed as needed), for a total context length of 3200 tokens per image/caption pair. For the multimodal dataset we used a wide variety of image/text pairs curated through web search.

2.1.6 Mathematical Problem Solving

Mathematical problem solving would seem to be a rather different domain from generative language, image, video, and multimodal modeling. To solve math problems, a model needs to learn to execute an algorithm to arrive at a deterministic answer. In contrast, the other distributions we have studied are typically genuinely probabilistic, and at least at an intuitive level, seem to require something a bit different from the simple algorithms that perform arithmetic or solve equations. We have included math problems to probe the generality of scaling laws and transformer performance.

We train and test models using the math problem generator [SGHK19], which generates a variety of problems in algebra, arithmetic, calculus, comparisons, numbers (integer properties), measurement, polynomials, and probability. When studying model and compute-budget scaling we procedurally generate the training problems in an online setting. We sample the default mixture of easy, medium, and hard problems, without a progressive curriculum. When studying dataset size scaling we use static training data sampled from the same distribution. As discussed further in appendix B, the data distribution has some unusual features, as easier problems will naturally appear more often than more difficult problems.

A few problem types require interpreting both numbers and strings as sequences of individual characters, so for simplicity we model all questions and responses at the character (byte) level. The model receives the problems as plain text, and we fill a transformer’s 512-token context window with concatenated problems, using a mask so that only the tokens corresponding to answers contribute to the loss.

The problem generator⁵ [SGHK19] can be provided with an ‘entropy’ s . The training distribution samples from $s \in [3, 10]$, while interpolate-testing corresponds to $s = 8$, and the extrapolate test involves $s = 12$, along with some other extensions to increase compositionality. In the online setting, we cannot be sure the interpolate tests are deduplicated from the training data, but the extrapolate test must be. To supplement the test data and further study extrapolation, we generated new test sets with $s \in [1, 19]$, with larger s posing a greater challenge to the model, as $s > 10$ is literally out of the training distribution, and requires extrapolation.

We found consistently poor performance on the two extrapolation generators `probability__swr_p_level_set_more_samples` and `probability__swr_p_sequence_more_samples` from [SGHK19], with larger models overfitting against them and achieving worse loss (but higher accuracy)

⁵The generator settings vary somewhat among problem types, with some depending on more parameters.

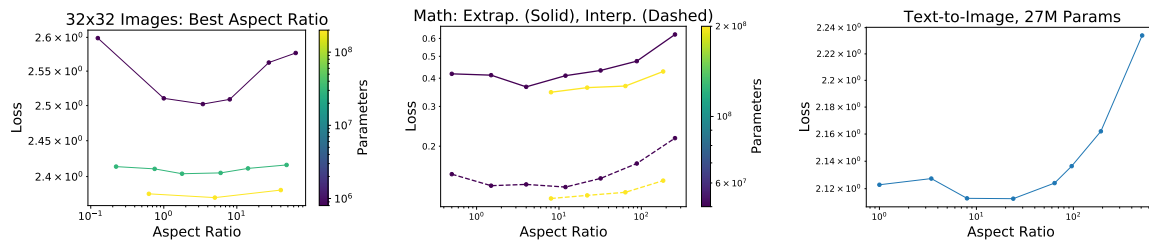


Figure 4 Optimal aspect ratio— We show trained performance as a function of the aspect ratio, defined as width / depth, or more precisely $\equiv d_{\text{model}}/n_{\text{layer}}$. The optimal aspect ratio for language [KMH⁺20] was about 10x larger.

than some smaller models. So *we have not included their contribution* in figures 1 and 5, as the poor loss on these modules would dominate the trends.

We provide more details and many additional results on math in appendix B, including results per module, dataset size⁶ scaling, and further analysis of performance vs difficulty level. There we also show trends for the training loss, which do not adhere as well to a power-law form, perhaps because of the implicit curriculum in the frequency distribution of easy and hard problems.

2.2 Model Size Scaling and Aspect Ratios

Arguably the simplest scaling relation compares the loss achieved by models of various sizes N once they are trained to convergence with a dataset large enough to obviate overfitting. Throughout this paper we report N as the number of non-embedding parameters in a transformer model, motivated by prior results on language [KMH⁺20]. Results for the scaling of $L(N)$ are depicted in figure 3, along with fits to equation (1.1).

We define $L(N)$ using the loss at convergence (practically, this means as close to convergence as is feasible), but the largest models we study will not have fully converged. Thus caution is warranted when interpreting $L(N)$ trends according to equation (1.2) and identifying the irreducible loss as an entropy, and the reducible loss as a KL divergence. Nevertheless, the reducible losses typically fit very well to a pure power-law trend. As an aside, we often find intriguingly good power-law plus constant trends when recording the loss after training all models for a fixed number of training steps.

We have found that for any given data modality, transformer models typically have an ideal aspect ratio $d_{\text{model}}/n_{\text{layer}}$ that maximizes performance while holding model size N fixed. In figure 4 we display converged performance as a function of aspect ratio for a few model sizes in several domains. We see that image and math models perform optimally with an aspect ratio ≈ 5 , which suggests that on these domains we should aim for deeper and thinner models, with at least a 10x smaller aspect ratio compared to optimized language models. The difference may be even greater due variations in m_{attn} and m_{mlp} settings.

Finally, note that image and video models with roughly 10^4 parameters under-perform the trends, with worse performance evident for higher resolution images. The video models must attend to a 4096-token context, while 32x32 images have a 3072-token context, so we speculate that tiny models under-perform because they have difficulty attending to contexts comparable in length to their non-embedding parameter count.

2.3 Compute Scaling and Optimal Model Sizes

Instead of focusing on converged performance, one can study the loss L achieved with a finite training compute budget C when training with a large enough dataset to avoid overfitting. We define C theoretically rather than empirically, and approximate⁷ it as $C \equiv 6NE$ where N is the non-embedding parameter count (model size) and $E = SB$ is the total number of tokens processed during training (with S the number of parameter updates and B the batch size in tokens). The results for $L(C)$ from a variety of model sizes are depicted in figure 5, along with the pareto-frontier of optimal loss for a given compute budget, and a power-law plus constant fit forced to lie below this frontier.

⁶The math models in figure 4 used $m_{\text{mlp}} = 4, m_{\text{attn}} = 1$ like language models, unlike the math models used to study model and compute scaling, as these aspect ratio tests were performed earlier.

⁷The factor of 6 includes a factor of 2 for add-multiply and a 3 to include forward and backward passes.

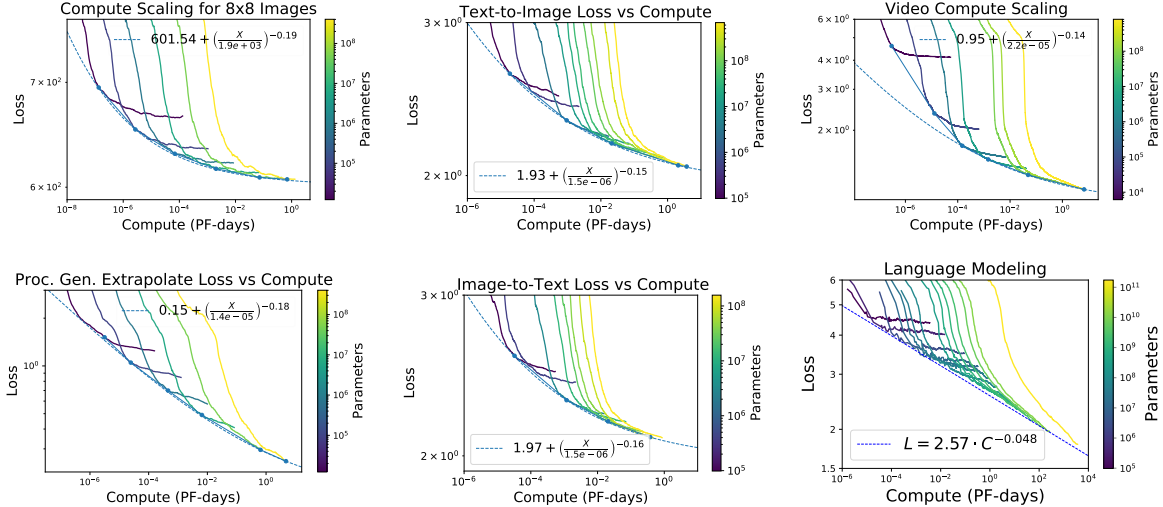


Figure 5 Scaling laws with compute— Scaling laws with compute (total estimated floating point operations) for various domains, along with power-law plus constant fits (dashed). This is identical to figure 1, except that we do not subtract the fitted constant irreducible loss. Note that very small models underperform compared to the trends when they model images or videos with very large contexts. Note also that the largest language models [BMR⁺20] were not trained to convergence.

The compute trends are most relevant for differentiating between the irreducible loss and reducible losses, since they avoid the issue of training to convergence, which makes the interpretation of $L(N)$ difficult. We display the reducible loss trends for $L(C)$ in figure 1, and emphasize that these appear to be pure power-laws, even when the reducible loss is much smaller than the irreducible loss.

We can use the $L(C)$ trends to estimate the model size N_{opt} that optimizes the loss when training is constrained by a fixed compute⁸ budget C . For this purpose we select points on the convex hull of the loss versus compute frontier; these can be seen as blue points in figure 5. The results for all domains together appear in figure 2, while each domain is shown separately with individual fits in figure 16. In all cases we find that $N_{\text{opt}}(C) \propto C^\beta$ can be fit with a pure power-law, with all exponents fairly close to $\beta \sim 0.7$. This suggests that one should spend most of a growing training compute budget by training much larger generative models.

When estimating $N_{\text{opt}}(C)$, one might worry about errors due to a sub-optimal usage of data. Specifically, if the batch size is too large early in training, then some compute may effectively be wasted. This can be studied by identifying the critical batch size [MBB17, MKAT18] above which there are diminishing returns to further data parallelism. In prior work [KMH⁺20] this was taken into account by measuring the critical batch size and using relations derived in [MKAT18] to adjust compute estimates. We have not made this adjustment here, as it would require a number of additional experiments in order to measure the critical batch size in each domain. For large model sizes and compute budgets these effects should be small, because most or all of training involves batches smaller than the critical batch size (which grows quickly during training [MKAT18]), but this issue may be worth revisiting in the future.

The total number of tokens processed during all of training is $E = \frac{C}{6N} \geq D$, where D is the dataset size, with equality representing training for only a single epoch. This means that $D \propto C^{1-\beta} \propto N^{\frac{1-\beta}{\beta}}$. We clearly have $\beta > 0.6$ for all data modalities and by a comfortable margin, suggesting that dataset size should not grow faster than $D \propto N^{2/3}$ during compute-optimal training, with a more reasonable median estimate of $D \propto N^{0.4}$. This unambiguously sub-linear scaling across all data modalities runs somewhat counter to conventional wisdom. As a word of caution, we have yet to train models in a regime where compute optimal training actually implies $D \ll N$ numerically. We discuss this further in section 6.

⁸For a fixed amount of training compute, we can train smaller models at the expense of worse performance. Hence, when accounting for both inference and training compute, the optimal model size may be somewhat smaller than described here. See [KMH⁺20] for a discussion this tradeoff.

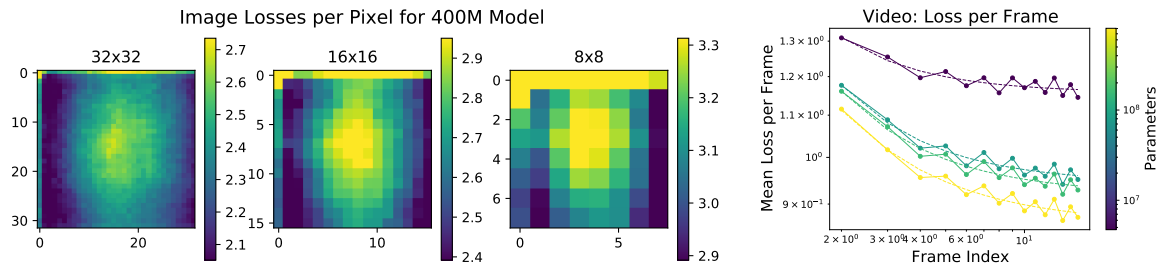


Figure 6 Position-dependent loss for images and video— We show trends for the loss as a function of position in the context for image and video models. On the left we have the mean loss over the three colors for images of various resolutions. The top-left pixel actually has significantly higher loss, off the color scale, which was set to make the pattern clear for the image as a whole. On the right we see the mean loss per frame for video models, as a function of the frame index. The oscillatory behavior per frame is due to the video encoding.

2.4 Loss versus Position in the Context Depends on the Structure of the Data

Some trends in the loss are highly dependent on the structure of the data. A clear example of this is the loss as a function of the position in the context, ie the loss per token for language models, loss per frame for video models, or the loss per pixel in visual domains. We provide two examples in figure 6. Note that for images the very first pixel typically has a large loss, outside the color range shown; we chose not to extend the color range as it would have obscured the patterns in the remainder of the image.

Language [KMH⁺20] and videos (per frame) show a power-law plus constant trend as a function of context position, as their data is naturally sequential. However, these trends do not apply to all to image modeling, where the loss is largest for the first pixels and near the center of the image. Thus power-law correlations in the context depend in an essential way on the nature of the data, and are not universal. In contrast, the form of the compute and model size scaling laws appears to be largely independent of the data distribution.

3 Image and Video Modeling, the Reducible Loss, and Downstream Tasks

Image data can be presented at a wide variety of resolutions, or it may be compressed, for example with VQ codes [vdOVK18]. These settings provide a way to modify the complexity of the data distribution, creating a useful arena for the study of neural scaling laws. Furthermore, we can finetune generative image models for classification to explore the quality of their learned features.

We will use these tools to explore the nature of the reducible and irreducible loss. In particular, at very low resolution (8x8) we can follow the power-law trend in the reducible loss all the way to a few nats/image, which can be achieved by models approaching a billion parameters. This gives us some reason for optimism when extrapolating similar trends on larger images beyond the realm that we can currently explore. It also strongly suggests that the power-law plus constant form of equation (1.1) will remain an excellent approximation.

Furthermore, we will show that improvement in fine-tuned classification performance continues smoothly even as the generative loss approaches the irreducible loss. This result strongly suggests that representation quality continues to improve smoothly even when the generative loss trend appears to taper off.

3.1 Varying the Image Resolution and Encoding

We trained Transformers on the YFCC100m dataset after scaling images down to 8x8, 16x16, and 32x32 pixel resolutions, along with 64x64 images encoded with VQ codes [vdOVK18] with 16x16 and 32x32 VQ code patterns. We display the trends for the reducible loss per image as a function of the compute budget in figure 8 (see figure 18 in the appendix for trends for the full loss). We include these figures to emphasize that the reducible loss for an optimally-allocated compute budget follows a power-law trend, even when the reducible loss becomes very small.

Note that the smallest models underperform as compared to the trends at resolutions greater than 8x8. We see this both for the compute trends in figure 8 as well as in model-size trends in figure 7. We speculate that this is due to difficulty utilizing the positional encodings. For example, our smallest models have only 10k

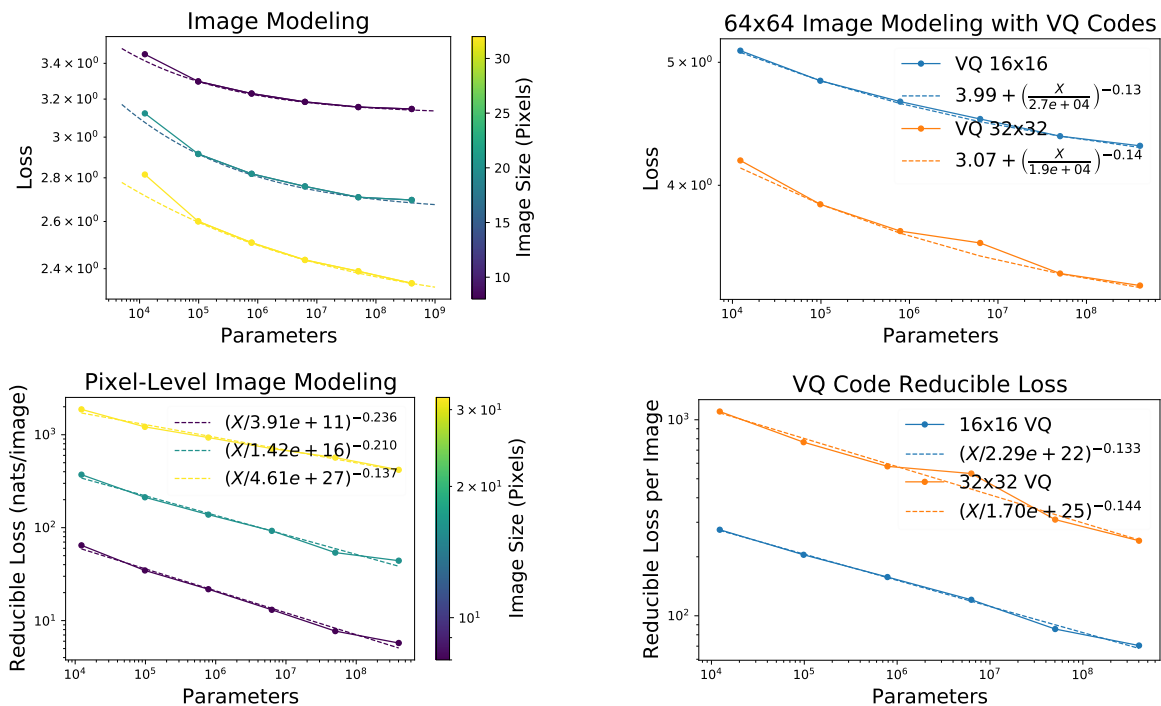


Figure 7 Comparison of image resolutions (model size scaling)— **Top:** We display scaling laws with model size for various image resolutions, and also with various VQ encodings, along with power-law plus constant fits (dashed) to equation (1.1). The fits for pixel-level image modeling are shown in table 3. Note that the tiniest (10k non-embedding parameter) pixel models underperform at higher resolutions; we suspect they have difficulty recognizing relative positions in larger images. These deficiencies are even more clearly visible in the compute trends. **Bottom:** We show the reducible losses, which estimate the KL divergence between the true probability distribution over images and the distribution predicted by our models. We show the result as a function of model size and image resolution or encoding, along with pure power-law trends.

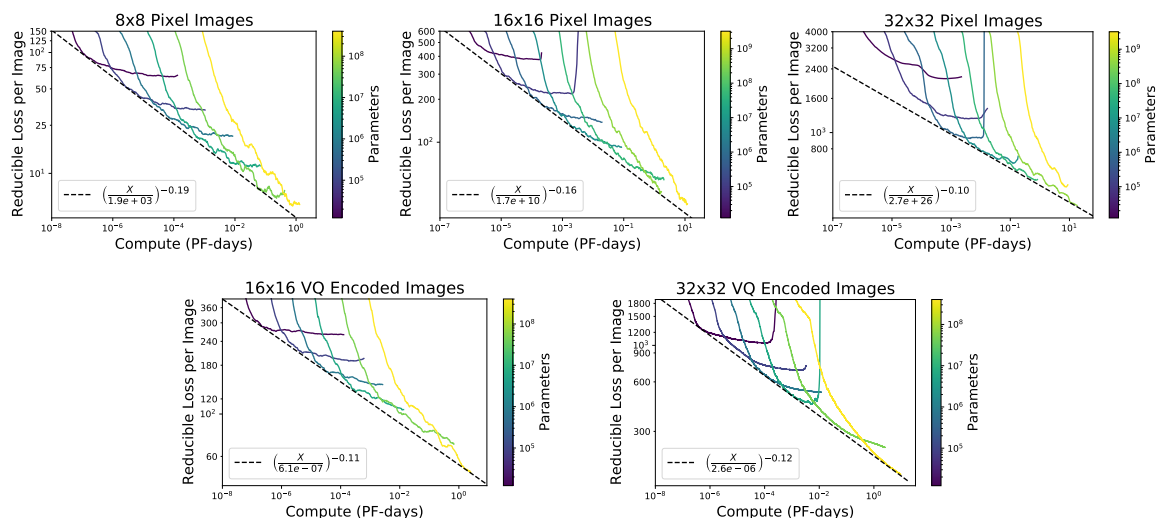


Figure 8 Comparison of image resolutions (compute scaling)— We display scaling of the reducible loss with compute for pixel-level image modeling at various resolutions (first line), and for various VQ encodings of 64x64 images (second line). We show the test loss, but we did not observe any train/test gap for these models. A few models diverged late in training.

Resolution	Reducible Loss per Image (nats)	Irreducible Loss per Image (nats)
8x8	$\left(\frac{C}{1.9 \times 10^3}\right)^{-0.19}$	602
16x16	$\left(\frac{C}{1.7 \times 10^{10}}\right)^{-0.16}$	2026
32x32	$\left(\frac{C}{2.7 \times 10^{26}}\right)^{-0.1}$	6806
64x64 (16x16 VQ)	$\left(\frac{C}{4.7 \times 10^{15}}\right)^{-0.11}$	1047
64x64 (32x32 VQ)	$\left(\frac{C}{3.1 \times 10^{19}}\right)^{-0.12}$	3246

Table 3 Per-image loss trends— Fits for the reducible and irreducible loss as a function of compute for various image resolutions, shown *per-image* rather than per-token as in table 1. Here compute C is measured in PF-days, so the denominators estimate the amount of compute needed to achieve a reducible loss of 1 nat/image. The irreducible losses estimate the entropy of the YFCC100M data distribution [TSF⁺15].

non-embedding parameters, while 32x32 images include 3072 tokens in their context, each with a distinct positional embedding.

To understand the significance of the reducible loss trends in table 3, recall that the cross entropy loss between the true distribution P and the model distribution Q is

$$\mathbb{E}_{x \sim P} \left[\log \frac{1}{Q(x)} \right] = D_{\text{KL}}(P||Q) + S(P) \quad (3.1)$$

The KL divergence vanishes when $P = Q$, and is otherwise strictly non-negative. Thus we can identify the irreducible loss with $S(P)$, the constant entropy of the true distribution. Then the reducible loss estimates the KL divergence between the true distribution and the distribution predicted by the model. This interpretation can only make sense if in the limit of infinite data and compute, we expect the transformer to perfectly model the data distribution. We have focused on $L(C)$ trends because the asymptotic limits of the model size trend $L(N)$ could be misleading if the models have not all been trained fully to convergence.

The power-law trends in D_{KL} can be extrapolated down to the level of just a few nats per image. Models powerful enough to reach this level of performance model the distribution of images with near-perfect fidelity. In fact we see that models with $\sim 1\text{B}$ parameters nearly achieve this feat for 8x8 ‘images’. However, we see that for larger images we would need enormous quantities of compute to perfectly model the true image distribution.

The consistency of the trends among distinct image resolutions in figure 7 and the strikingly small reducible loss for the 8x8 case suggests that if we could run much larger models, we would continue to see smooth improvements at higher resolution. It seems that compute requirements for a near-perfect model of the data distribution grow as a steep power-law or even an exponential in the image resolution. Of course we do not expect to need a perfect model of the probability distribution of real-world images for practical tasks.

3.2 Video Modeling and Individual Frames

For the case of video modeling, it is natural to extend the overall trends to the study of specific frames. We display several frame-dependent results in figure 9. On the left we show loss as a function of model size, omitting the first frame, which has a much larger loss and should be considered an image modeling problem. In the center we show compute scaling of the reducible loss on the final frame. On the right in the same figure we show the reducible loss for the final (16th) frame, which is of particular interest when generating a continuation of an existing video. Much like the trends for image modeling, we see that the reducible loss is very well approximated by a power-law, making it possible to forecast that we would need a model size around $\sim 10^{13}$ parameters and compute of around 10^4 PF-days to achieve a loss of just a few nats/frame on the final frame of this type of video.

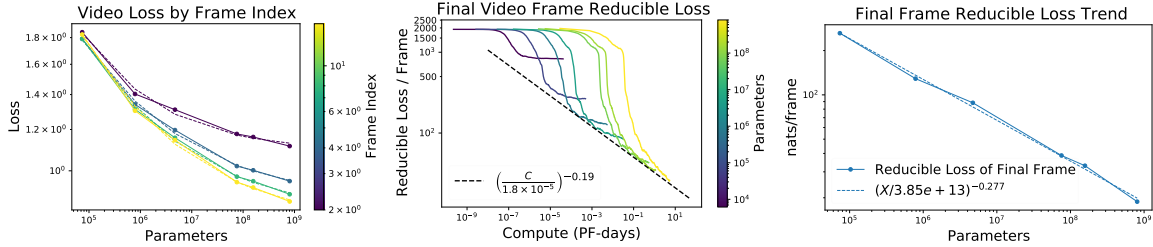


Figure 9 Per-frame video performance trends — On the left we show scaling trends for specific frames in 16-frame videos. In the center we show the reducible loss as a function of compute for the final frame of the video. On the right we show the reducible loss and its pure power-law trend with model size for the final frame in a video.

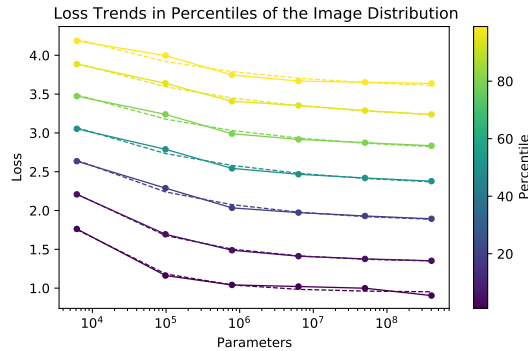


Figure 10 Performance trends for image dataset percentiles— We selected one thousand images from the 32x32 image test set, and evaluated the loss of all models on each image. In this figure we plot the trends in the 1, 5, 20, 50, 80, 95, 99 percentiles of the loss distribution over these images, along with power-law plus constant fits (dashed). We also observe similar trends for randomly chosen individual images (figure 17).

3.3 Scaling Trends for Individual Images

We have observed very consistent scaling trends on a variety of data modalities. This raises a question – does the loss achieved by different sized models on specific, individual data examples scale in the same way? Or are the distribution-level trends an aggregate of many different trends on individual examples?

To answer these questions, we evaluated the loss of all the pixel-level 32x32 image models on a thousand randomly chosen images from the test set. When plotting the loss as a function of model size for individual, randomly chosen examples, in essentially all cases we observe a smooth, power-law plus constant trend.

To convey this information, for each model size we evaluate the 1,5,20,50,80,95, and 99 percentile of the loss among a thousand images in the distribution, for each model size. We then plot the trends in these percentile losses in figure 10. We see very similar trends among all percentiles of the loss distribution, and all are well-described by equation (1.1). We show model size trends for eight randomly chosen individual test images in figure 17. We also display the most and least improved 10 images from a sample of one thousand test images in figure 20. Finally, we visualize the trends in a different way, by generating conditional samples at each model size, in figure 21.

We would expect that these findings also apply to other data modalities. On a quick inspection, we found the same patterns for randomly chosen text sequences and language models of different sizes.

3.4 Finetuning on ImageNet at 32x32 Resolution

By finetuning generative models for image classification we gain another handle on the scaling of performance with model size. We use the scaled-down 32x32 resolution ImageNet [CLH17] and finetune the 32x32 resolution pixel-level generative image models.

To turn these models into classifiers, we remove their final embedding matrix and use the mean-pooled (over all pixels) activations of the transformer’s final layer as the input to a new single-layer classifier. During

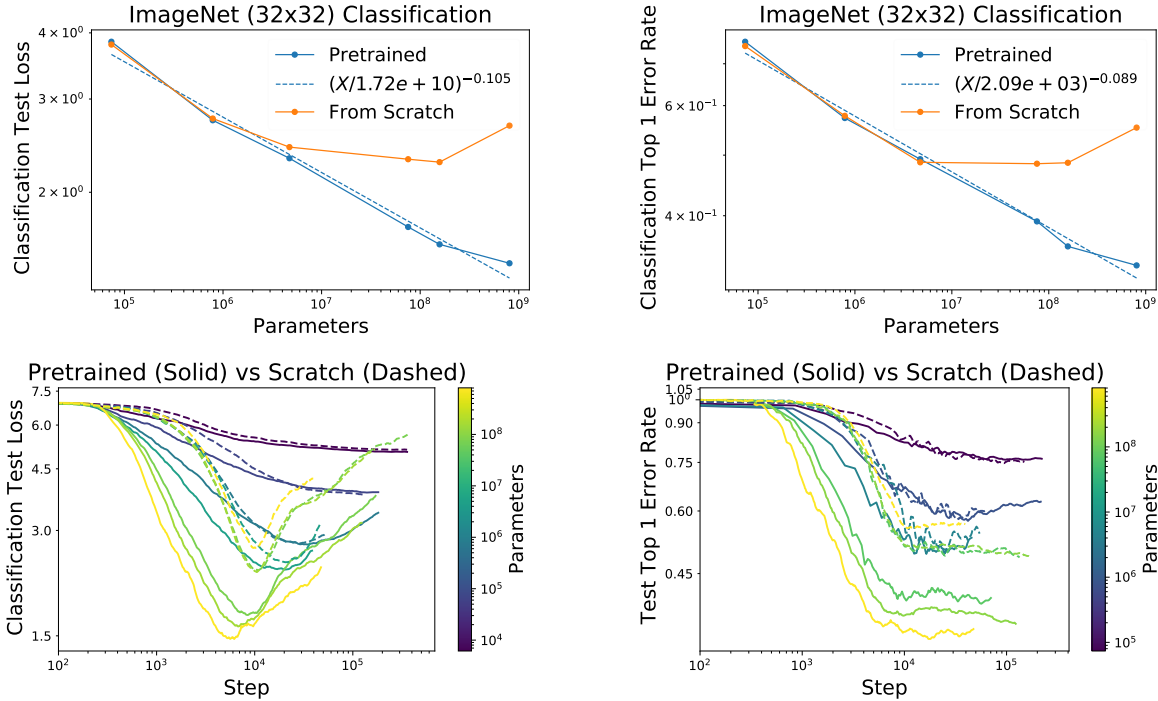


Figure 11 Trends in image classification performance— Top: We show model size scaling results for 32x32 pixel ImageNet [CLH17] classification. We compare models trained from scratch on ImageNet classification (ie with no pre-training) to finetuned generative models. Though the generative loss trend bends as it approaches the irreducible loss (figure 7), the pretrained models exhibit a straight power-law trend in classification performance vs model size, which also continues far beyond the point where the models that were trained from scratch exhibit overfitting. **Bottom:** Larger pre-trained models fine-tune significantly faster, and to significantly better performance, despite the approach to the irreducible generative loss. The same does not hold when training from scratch.

finetuning we backpropagate through the full transformer, and we do not freeze any of its weights. As a comparison, we also train equivalent randomly initialized transformer models ‘from scratch’ on only the classification task.

Finetuning learning curves for both pretrained and randomly initialized models are available in figure 11. In all cases we use a batch size of 1024 images, and we use the same learning rate schedule for finetuning as was used for pretraining. We see that for small models, pretraining affords almost no benefit compared to training from scratch, but it greatly enhances the performance of larger models.

More importantly, in figure 11 we show the model-size trends of ImageNet classification performance for pretrained and randomly initialized models. We see that the pre-trained models follow a smooth, pure power-law⁹ trend in both loss as well as error rate (1– accuracy). The very existence of these trends on a downstream finetuning task provides a striking confirmation of the importance of neural scaling laws for AI capabilities. In the case of language, GPT-3 [BMR⁺20] provides many more examples.

We also emphasize that the proximity to the irreducible loss does not necessarily indicate diminishing returns with regards to model performance. The trends in figure 11 continue smoothly, even though the green curve corresponding to 32x32 resolution in figure 7 suggests a close approach to the irreducible loss for models with $> 10^7$ parameters. Apparently, a great deal of important semantic information lies in the ‘last few bits’ near the irreducible loss. We may also interpret this as the pre-training process providing a highly effective regularizer for downstream tasks.

⁹We have not encountered a clear irreducible loss in the range of model sizes that we have explored.

4 Multimodal Models and Information Gain

Is a picture worth a thousand words? With multimodal models we can study the amount of information that one domain provides about another. For this purpose we study the empirical mutual information between images and text and the infogain defined in equation (1.3). The latter has the interesting property that it must lie in the interval $[0, 1]$, with larger values suggestive of better performing multimodal models.

To estimate the empirical mutual information between the image and text for text-to-image models, we subtract the captioned-image loss from the image loss in the presence of a blank caption. Similarly, we subtract text losses with and without corresponding images for image-to-text models.

However, these measurements have a potentially serious flaw – if the models have only been trained on multimodal data, then blank captions and blank images may be out of distribution. We minimize this issue by measuring the mutual information only after finetuning our models for 10^4 steps on an even mixture of data with and without captions (for text-to-image) or with and without images (for image-to-text). Empirically we find that without this finetuning, the mutual information is measured to be about twice as large. In the case of text-to-image models, we also tried training from scratch on a 95/5 mixture of multimodal and blank caption data, and found very similar results. The learning curves for the mutual information and some other comparisons can be found in appendix C.

We plot the mutual information and the infogain ratio in figure 12. We see that billion-parameter, decoder-only transformer models extract about 8 nats of information concerning the image from an average text caption in the test set. In the case of both Image-to-Text and Text-to-Image multimodal models, we observe empirically that mutual information and infogain varies with model size as

$$I(\text{text}, \text{image}), \text{ Infogain} \approx \lambda \log \left(\frac{N}{N_c} \right) \quad (4.1)$$

with different λ and N_c for the two cases. We can derive this approximate formula from plausible assumptions, as discussed in appendix E. If this trend holds over a large range of N , it might be used in combination with the upper bound $\text{infogain} < 1$ to roughly estimate the maximal productive model size.

However, the trends identified in figure 12 suggest a very slow growth of infogain with N for these models, so it seems unrealistic to extrapolate all the way to an $\text{infogain} = 1$. Furthermore, in the data distribution the text and images are not always closely correlated, as in many examples much of the text has little to do with the accompanying image. So instead we might ask when 20% of the information in the text will be used to define the image, doubling the infogain of a 1B parameter model. For text-to-image models, this threshold will be met with models of size $N \approx 3$ trillion parameters, though for image-to-text models this remains far out of reach. Other architectures may improve on these results, but we conjecture that they will display similar trends with model size.

Text-to-image models have much larger mutual information and infogain, as compared to image-to-text models. We speculate that this is due to the fact that much more processing is required to extract semantic information from images than from text.

We can now revisit the question of how many words a picture is worth. Figure 3 shows the loss per text token, including padding tokens; if we exclude padding tokens, the largest image-to-text models achieve a loss of 2.6 nats per text token, or about 3.4 nats per word. Comparing the image-to-text mutual information of 8 nats, we find that a 32x32 image is worth only about 2-3 words to our best models.

5 Mathematical Problem Solving and Extrapolation

In the context of machine learning, generalization most often refers to the gap between test and training performance. But on a conceptual level, generalization can also refer to the more ambitious possibility of extrapolation from the training distribution to a larger or more diverse distribution. Mathematical problem solving lends itself very naturally to the study of extrapolation, because we can extend the range of numbers or operations used to create math problems, or the recursive/compositional depth [HDMB19] required for a solution.

We studied this phenomenon in the fundamental figure 3, where we evaluate problem solving performance using a variety of test sets indexed by a numerical level, which corresponds to an ‘entropy’ used for generation [SGHK19]. We observe fairly smooth power-law plus constant trends for the loss on all of these test sets, but

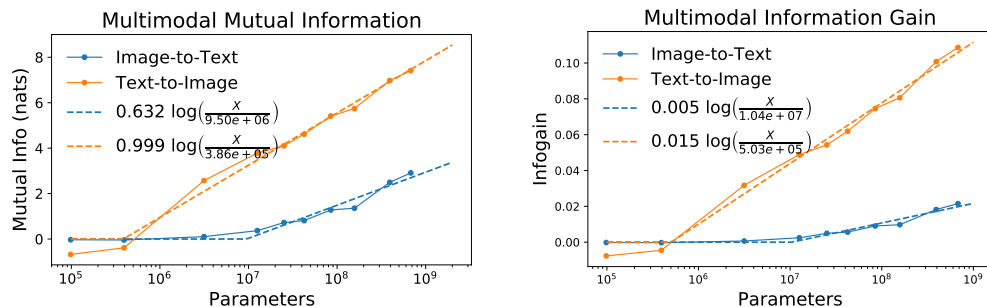


Figure 12 Multimodal information trends for multimodal models— We show the empirical mutual information between image and text in multimodal models (left) and the Infogain (right), which is the ratio of the empirical mutual information to the empirical entropy of the text. The results in these plots were compiled after finetuning multimodal models for 10k steps on half multimodal, half blanked caption/image data, to ensure that blank captions/images were not out of distribution. The largest text-to-image models use about 10% of the information in the text when constructing images.

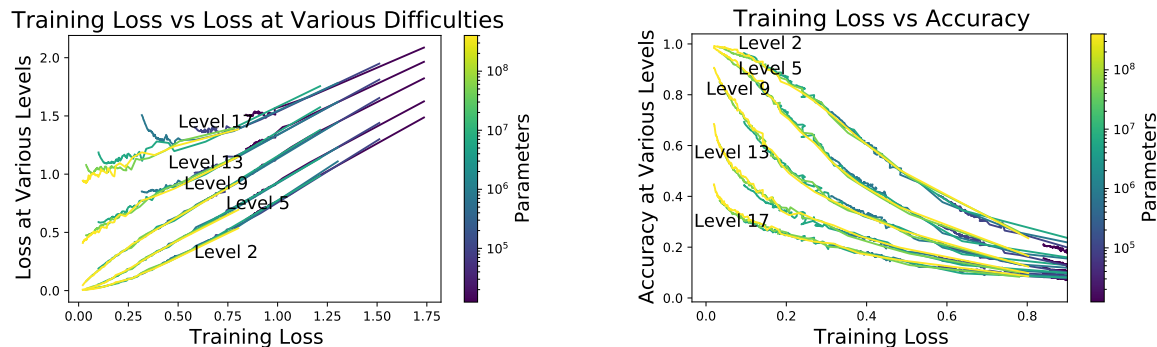


Figure 13 Mathematics difficulty levels— We show the loss (left) and accuracy (right) during training, as a function of the training loss, for math problems at various difficulty levels. We emphasize that models of different size perform nearly identically when we hold the training loss fixed. Thus in the case of math problem solving, both interpolation and extrapolation performance depends on model size primarily through the training loss. Note the difficulties ≤ 10 are within the training distribution; for levels > 10 we expect non-zero test loss even as the training loss tends to zero.

with different exponents and offsets depending on the difficulty level. So extrapolation performance improves with model size.

However, as we show in figure 13, the extrapolative capabilities of these models predominantly depends on the models’ performance on the training distribution. That is, models of different sizes that achieve the same loss on the training distribution perform about equally on the various test distributions. In this sense, increasing the model size does not automatically improve extrapolation, except insofar as it improves performance on the training distribution. Similar results were found in [KMH⁺20] when extrapolating from one text distribution to another.

Finally, for completeness we note that the information theoretic interpretation of the loss has a somewhat different meaning in the context of math problem solving, where the answers are deterministically related to the questions, so that the entropy should truly vanish. For much more detailed results on math performance and a great many more trends see appendix B.

6 An Inconsistency in Compute and Datasize Scaling Laws

An inconsistency among the datasize and compute scaling laws was observed in [KMH⁺20]. In this section we will study the same phenomenon using image models on low resolution images, though we expect the results will be qualitatively the same on any of the datasets we have covered.

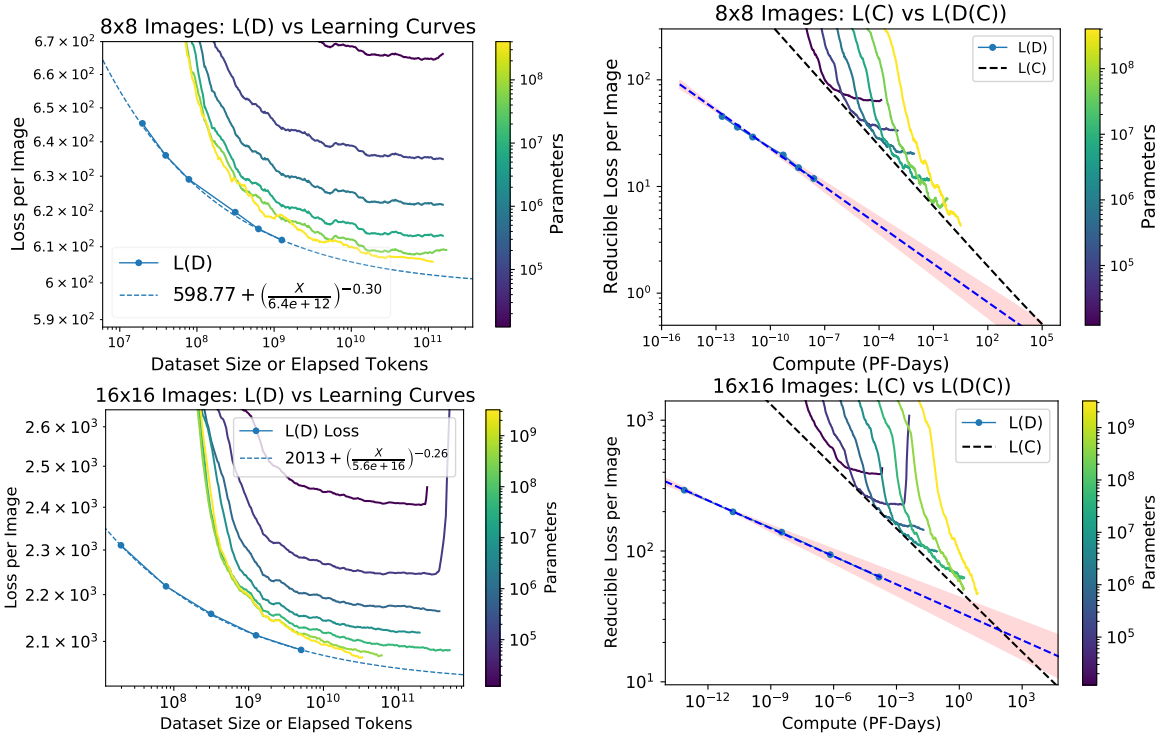


Figure 14 Training speed approaches a limit— **Left:** These figures show learning curves for various model sizes, along with the trend for fully trained, early-stopped $L(D)$, identifying the dataset size in tokens with the number of elapsed tokens during training. We observe that the learning curves are approaching $L(D)$ as model size increases. **Right:** We show learning curves along with the $L(C)$ trend in black. On the same plot we show $L(D)$ vs $C(D)$ in blue, where the latter is determined by identifying the optimal proportion of compute to allocate to tokens, and then assuming this corresponds to one epoch of training. By construction all learning curves must lie above and to the right of the blue dashed line, so the intersection of the black and blue lines suggests a breakdown of some trend. The red shaded region corresponds to altering the optimal model size exponent by $\pm 5\%$, illustrating that projections are extremely sensitive to these trends.

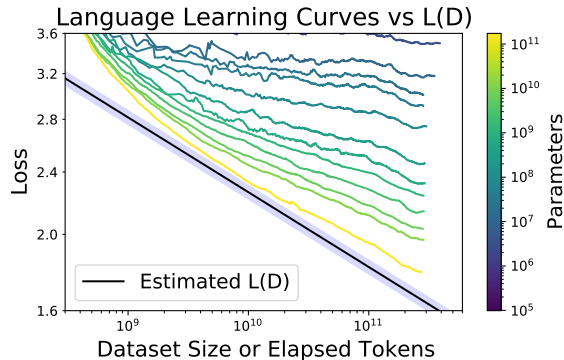


Figure 15 Training speed approaches a limit (language)— Here we show an approximation of $L(D)$ with 2% estimated errors, and the language modeling learning curves from [BMR⁺20]. The $L(D)$ trend comes from [KMH⁺20], but the models in that work were trained on a slightly different data distribution and with half the context length of [BMR⁺20].

Before discussing the inconsistency, consider the plots on the left of figure 14. We show both learning curves and the trend $L(D)$ for trained models, identifying the dataset size with the number of tokens seen by various models during training. The learning curves lie above the $L(D)$ trend because the optimization process fails to achieve the minimum loss in a single epoch. If the optimizer were perfect (in a sense), then $L(D)$ would coincide with the learning curve, assuming performance is not limited by model size. Note that as model size increases, the learning curves appear to approach ever closer to the $L(D)$ trend. This means that larger models learn faster, and it also implies that optimization becomes increasingly effective as model size increases. But learning curves will always be bounded by $L(D)$, which sets the sample efficiency. We show the same phenomena for language in figure 15, though we can only estimate¹⁰ $L(D)$ for these models.

To see an apparent inconsistency, we must compare the projections from two different trends. For the $L(C)$ compute trend we can just reproduce results from figure 7. To plot $L(D)$ with compute on the x-axis, we will use the power-law trend $N_{\text{opt}}(C) \approx (2.8 \times 10^8)C^{0.74}$ for 16x16 images (see figure 16), where C is measured in petaflop-days. From this we can solve for the optimal number of tokens processed during training using $C = 6DN$, which leads to $C(D) \approx (5 \times 10^{-42})D^{3.9}$ where D is measured in tokens. A similar analysis applies to 8x8 images. Using these results we can plot $L(D)$ vs $C(D)$ parametrically, as shown on the right of figure 14 for the reducible¹¹ loss (chosen for clarity on the log plot). We have also included a shaded region showing the effect of changing the empirically extracted $N_{\text{opt}}(C)$ trend exponent by $\pm 5\%$.

The inconsistency arises because all learning curves must lie above the $L(D)$ trend on the right of figure 14, but the extrapolation of $L(C)$ eventually intersects and passes below $L(D)$. Either $L(D)$, $L(C)$, or the $N_{\text{opt}}(C)$ trend must break down at or before this intersection point. Note that the existence of this intersection is an inevitable consequence of the power-law form of the trends, since these lead to straight lines on a log-plot, and two straight lines must intersect.

We do not know for certain how this inconsistency or its equivalent for language [KMH⁺20] are resolved. However, the observation of the left of figure 14 and our earlier discussion suggests a plausible hypothesis. As we increase model and dataset sizes, optimization becomes increasingly efficient, until eventually learning curves begin to merge with the $L(D)$ trend, so that there are no benefits to be gained from training for more than a single epoch [Kom19]. Near the intersection point, the compute frontier would bend and become coincident with $L(D)$. From this point of view, the fact that $L(C)$ appears steeper than $L(D(C))$ is due to a deficiency with optimization, which requires more than one epoch to reach a local minimum of the test loss. It would be interesting to investigate this hypothesis in the future. If it is true, it suggests that the relative scaling of optimal model and dataset sizes may eventually change, and perhaps will ultimately be set by trends for overfitting such as those found in [RRBS19, KMH⁺20].

Finally, we note that the irreducible loss from dataset size trend is measured at $L(D = \infty) \approx 2013$ nats/image (16x16), and 599 nats/image (8x8), while that extracted from compute trends is $L(C = \infty) \approx 2023$ nats/image (16x16), and 602 nats/image (8x8). These estimates for the entropy of low-resolution YFCC100M images are quite similar, and provide a consistency check on our results.

7 Related Work

Predictable scaling trends for modern neural networks have been studied by a variety of groups, beginning with [HNA⁺17]. More recently [RRBS19, LWS⁺20, RDG⁺20, Kom19, RFCS20] studied scaling relations using many model architectures and datasets, with the work on language modeling in [KMH⁺20] closest to our approach here. Work on the 175B parameter GPT-3 model [BMR⁺20] was partially motivated by neural scaling laws.

There has not been a great deal of work on theoretical explanations for the very precise scaling relations we and others have identified. A simple theory connecting scaling exponents to the inverse of the dimension of the data manifold was proposed in [SK20]. Expansions in the model size, particularly at large width [LXS⁺19, JGH18] may provide another useful framework for thinking about some of our scaling relations, if they are in fact applicable [LBD⁺20] to optimally tuned hyperparameter settings.

The models and data modalities we used have been widely studied in the past. Autoregressive image models have been trained starting with PixelRNN [vdOKK16], with the recent work [CRC⁺20] nearly identical to our

¹⁰We need to account for slightly different data distributions, and context lengths differing by a factor of 2. We estimate that these produce errors less than about 2% of the loss, which we show as a shaded region on the plot.

¹¹For these figures we subtract the irreducible loss measured from each of $L(D)$ and $L(C)$, respectively, since numerically the irreducible losses from these two measurements are not exactly equal.

models and training procedure. Transformer-based video models were trained in [WTU19] and multimodal models in [TBL⁺19]. The original authors trained various models, including transformers, on the math problem dataset [SGHK19], and it has also been studied with more specialized architectures [SSF⁺19]. Our models are typically simpler than many of those that have been previously discussed, as we exclusively use decoder-only [LSP⁺18] transformers with dense or sparse [CGRS19] attention.

8 Discussion

We have argued that a single neural architecture, the Transformer, can be applied to the generative modeling of images, videos, multimodal data, and math, along with language [KMH⁺20, BMR⁺20]. We identified common scaling laws for the loss achieved on all data modalities as a function of both model size and compute budget. As in the case of language, these results imply that larger models become more sample efficient. Furthermore, we found that in some important cases, finetuned performance on downstream tasks also follows similar scaling laws. This suggests that trends in the generative modeling loss translate into advantages in practical capabilities.

A greater surprise was the approximately universal trend (figure 2) for optimal model size as a function of the training compute budget – we did not anticipate that the exponent $N_{\text{opt}} \propto C^{0.7}$ would be largely independent of the data distribution. This trend implies a dual trend for the number of tokens elapsed during optimized training, as a function of C or N , and leads to the conclusion that larger compute budgets should be ‘spent’ mostly on larger models, rather than much longer training runs. So this lesson from language modeling [KMH⁺20] generalizes. These empirical regularities beg for theoretical explanation – why do these scaling relations hold?

The scaling laws also suggest a shift in perspective away from the particularities of neural architectures, loss functions, and training algorithms and towards the broader commonalities that appear when machine learning is studied across a large hierarchy of model, data, and compute scales. Work in ML often involves identifying specific deficiencies in current capabilities and remedying them through the alteration of models and algorithms. Perhaps many capabilities simply lie on a spectrum that can be continuously unlocked with increasing scale, as might be suggested by the metalearning capabilities of the GPT-3 model [BMR⁺20].

We also discussed some information theoretic implications of the scaling laws. Perhaps the most important point was that the two terms in equation (1.1) can be interpreted as the entropy of the true data distribution, and the KL divergence between that distribution and a given generative model. The identification of the entropy was made possible through the extrapolation of a precise trend, and would not be predictable using the results from a single model. We also observed intriguing scaling laws for the empirical mutual information between images and captions in multimodal models. This is particularly interesting because the mutual information must be bounded by the entropy of the caption.

Acknowledgments

We thank Yasaman Bahri, Miles Brundage, Yura Burda, Paul Christiano, Ajeya Cotra, Psycho Debiak, Ethan Dyer, Harri Edwards, Danny Hernandez, Jacob Hilton, Jaehoon Lee, Brice Menard, Chris Olah, Utkarsh Sharma, and Ilya Sutskever for discussions and feedback on this work.

Thanks as well to Chris Berner, Ben Chess, Eric Sigler, and Clemens Winter for managing and scaling the supercomputing clusters and research platform that allowed us to run these experiments.

Contributions

Tom Henighan performed and analyzed the image and video modeling experiments, and maintained the codebases for experimentation and data analysis that enabled our results.

Jared Kaplan performed and analyzed the math experiments, led the overall data analysis, and wrote the paper.

Mor Katz performed the multimodal experiments and data analysis.

Jacob Jackson, Chris Hesse, Heewoo Jun, and John Schulman collaborated on video modeling experiments.

Jacob Jackson, Heewoo Jun, Prafulla Dhariwal, and Alec Radford, developed the VQ-VAE training strategies and codebase.

Sam McCandlish analyzed the progression of question-answering capabilities in language models.

Aditya Ramesh and **Alec Radford** provided guidance on multimodal modeling and optimization.

Chris Hallacy and **Alec Radford** curated the multimodal datasets.

Heewoo Jun and **Aditya Ramesh** curated the image datasets.

Chris Hesse, Heewoo Jun, and Alec Radford curated the video datasets.

Mark Chen provided guidance on image modeling and finetuning.

Tom Brown, Scott Gray, Benjamin Mann, Nick Ryder, Prafulla Dhariwal, and Daniel Ziegler built, optimized, and maintained our codebase for training large transformer models.

Dario Amodei advocated for a broad study of scaling laws for generative modeling.

Sam McCandlish and **Jared Kaplan** led the research.

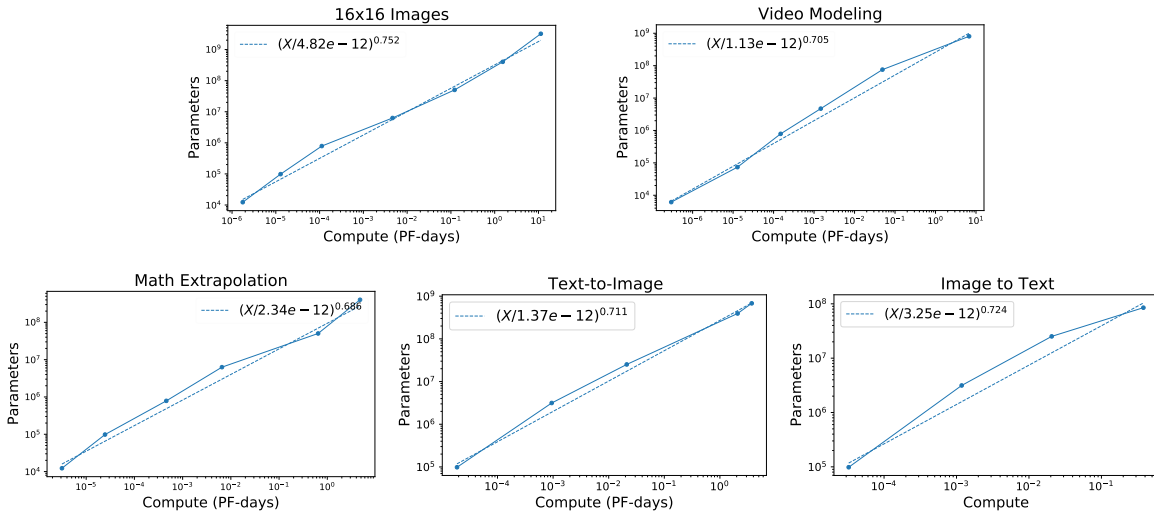


Figure 16 Optimal model size (individual trends)— We show the optimal model size for a given compute budget, along with power-law fits, based on the points at the compute-efficient frontier of figure 5. These trends are combined in figure 2.

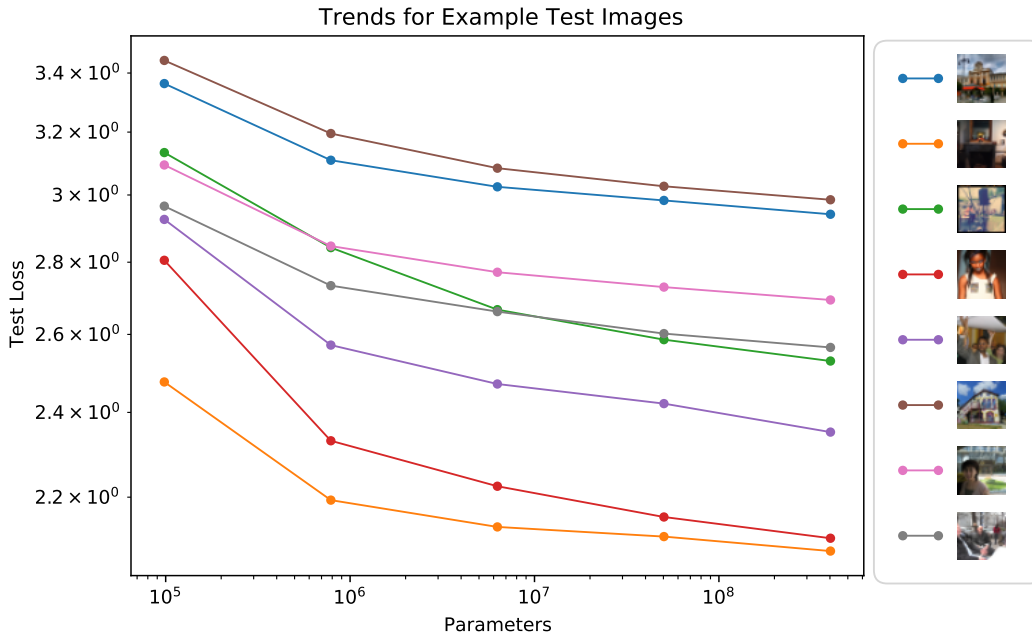


Figure 17 Loss trend for individual images— We show the loss trend for eight randomly chosen images from the test set. These results are fairly typical.

A More Details on Image Modeling

In figures 18 and 19 we provide some additional information documenting compute scaling trends for images with different resolutions and encodings. In figure 20 we show images where the loss improved most or least as we pass from a 100k parameter model to a 400M parameter model. In figure 17 we also show trends for randomly selected individual images from the test set.

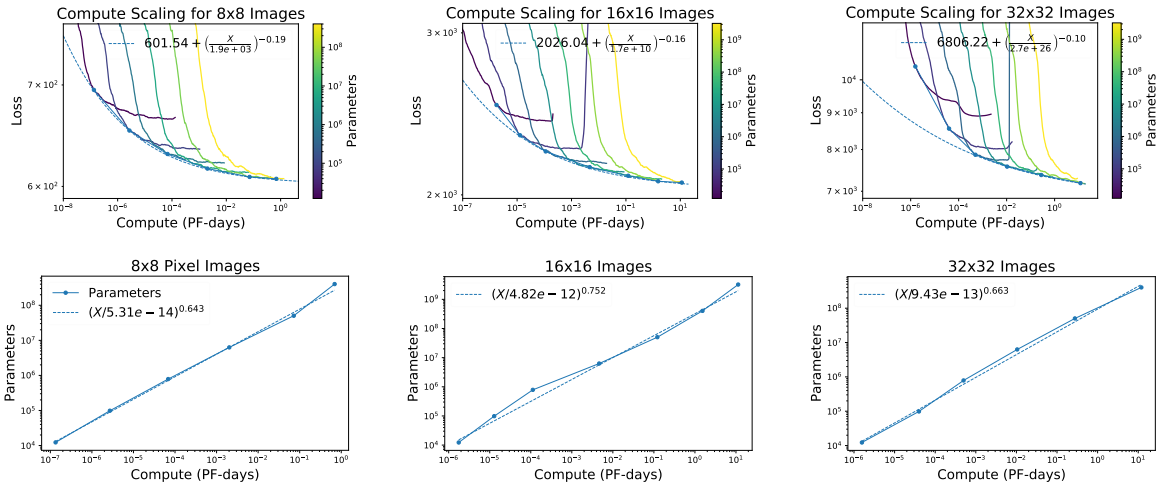


Figure 18 Compute trends for varied image resolution (pixel-level)— Scaling laws with compute for various image resolutions in pixels, along with power-law plus constant fits (dashed) to equation (1.1). The fits for pixel-level image modeling are shown in table 3.

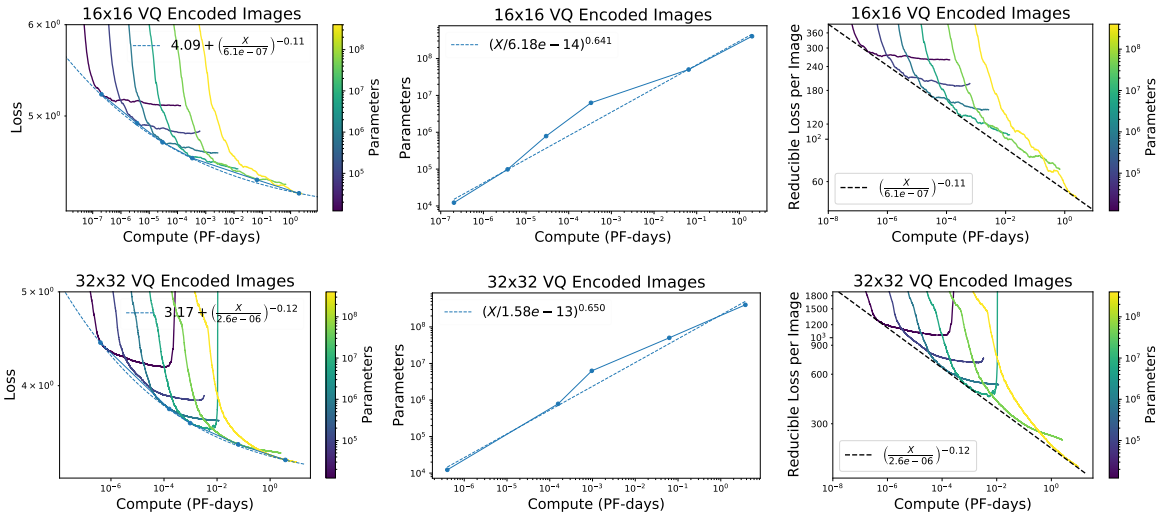


Figure 19 Compute trends for various image resolutions (VQVAE-encoded)— We display scaling laws with compute for 64x64 images encoded with two different VQ code resolutions, along with power-law plus constant fits (dashed) to equation (1.1). A few of these runs diverged beyond the compute frontier; in the worst case this led to a visible deviation from the model size trend in figure 7.

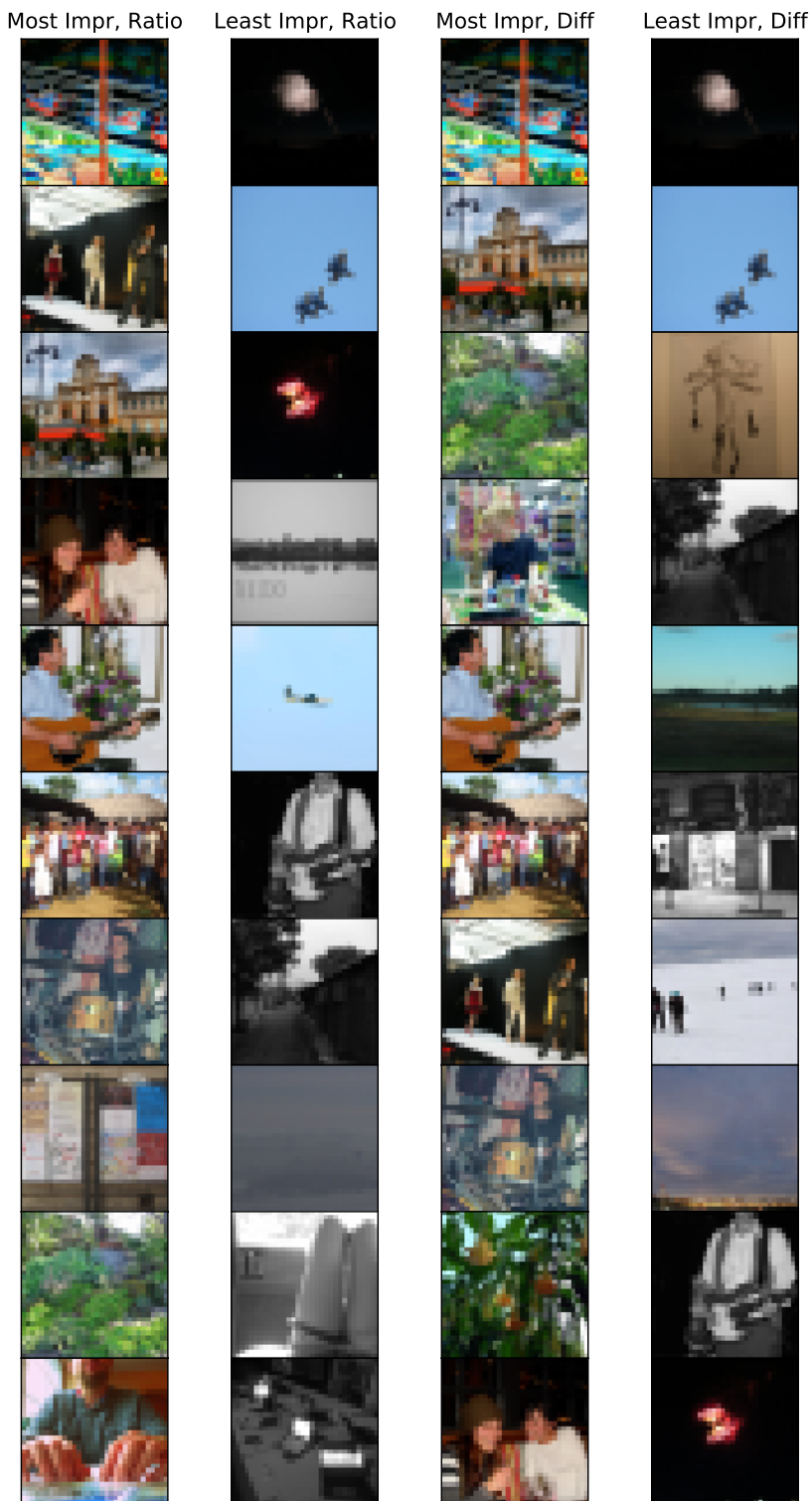




Figure 21 Trends in image completion quality— Here we show conditional-completions of 32x32 pixel models of various sizes, where the leftmost column is the original image, and each of the other columns shows completions from a model with a non-embedding parameter count labeled at the top. Models are provided the top half of the image as conditional context and the bottom half is sampled with temperature 1.0. There is a clear trend of increasing photorealism with larger models.

B Details of Math Experiments and Additional Results

B.1 Procedurally Generated Training Data

We generated all training data procedurally using the code provided by [SGHK19]. Problems were generated by randomly sampling modules from the training distribution, with an ‘entropy’ setting sampled uniformly from the integers $s \in [3, 10]$. The number of problems with entropy s is approximately 10^s , meaning that easy problems with low-entropy would likely be seen by the model many, many times, while some problems with $s \geq 9$ may not be seen at all. This means that the easy components of the training distribution may be memorized. Furthermore, our procedurally generated data was not deduplicated from the ‘interpolate’ test distribution [SGHK19], but it is completely disjoint from the ‘extrapolate’ test distribution.

The official extrapolate distribution only provides one difficulty level, and it also does not include all eight module-types. So we also generated distributions of problems with smoothly increasing difficulty level by setting the entropy $s = 1, 2, \dots, 19$. For most modules we simply used the interpolate settings, though for modules where other parameters were needed we generally used the extrapolation settings. Importantly, we did not include the `probability__swr_p_level_set_more_samples` and `probability__swr_p_sequence_more_samples_generators`, as we found our models always performed poorly on these problems, and quickly overfit on the loss for these generators (this can be seen in figure 23, where ‘probability’ represents the mean of these two generators).

Performance as a function of difficulty level and model size can be seen in figure 24. We note that performance degrades smoothly as we extrapolate away from the training distribution.

As an additional note, because these experiments were conducted much earlier, our dataset size scaling and aspect ratio scans use models with the fairly standard setting $m_{\text{mlp}} = 4$ and $m_{\text{attn}} = 1$, as with language and multimodal models, but different from the math models we used for compute and model size trends, where these parameters were smaller by a factor of 4, as with our image and video models. We made this change to smaller $m_{\text{mlp}}, m_{\text{attn}}$ as we found it helped to improve the training stability of very deep math models.

It is also worth noting that we evaluated extrapolation performance both using the training data files provided with [SGHK19] and by sampling with procedurally generated data (leaving out the two probability modules previously discussed). For trend plots we have used the procedurally generated data, but for reporting final accuracies in figure 26 we use the ‘official’ files.

B.2 Dataset Size Scaling

For the math dataset we studied optimal performance as a function of dataset size D , in the limit where $N \gg D$ so that performance is constrained by overfitting rather than by model size or compute budget. For each dataset size and problem distribution, we define $L(D)$ by taking the minimum loss during training (this differs slightly from early stopping, since we may evaluate at different steps if there are several metrics, ie losses on different test distributions, as is the case for math). For these experiments we used models with $n_{\text{layer}} = 64$ and $d_{\text{model}} = 512$ for all dataset sizes. We obtain power-law fits for $L(D)$, as shown in figure 22.

B.3 Additional Math Results

Here we provide several additional observations about math performance, which can be divided among different math modules and difficulty levels. In figure 23 we show performance on different modules (using the files provided in [SGHK19]), while in figure 24 we show performance as a function of difficulty level for different model sizes. We provide details of achieved accuracies on the official extrapolation and interpolation test sets in figures 26 and 27.

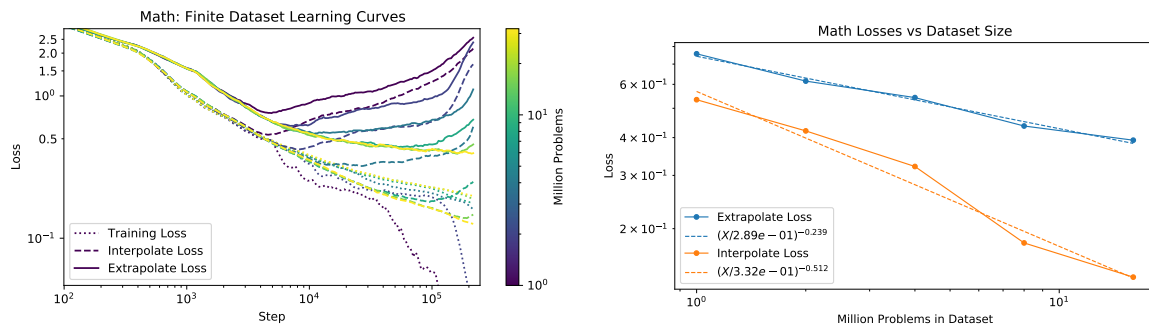


Figure 22 Math dataset size dependence— We show learning curves and trends in early-stopped loss as a function of dataset size. For the case of mathematical problem solving, we use a model with $n_{\text{layer}} = 64$ and $d_{\text{model}} = 512$ for all dataset sizes.

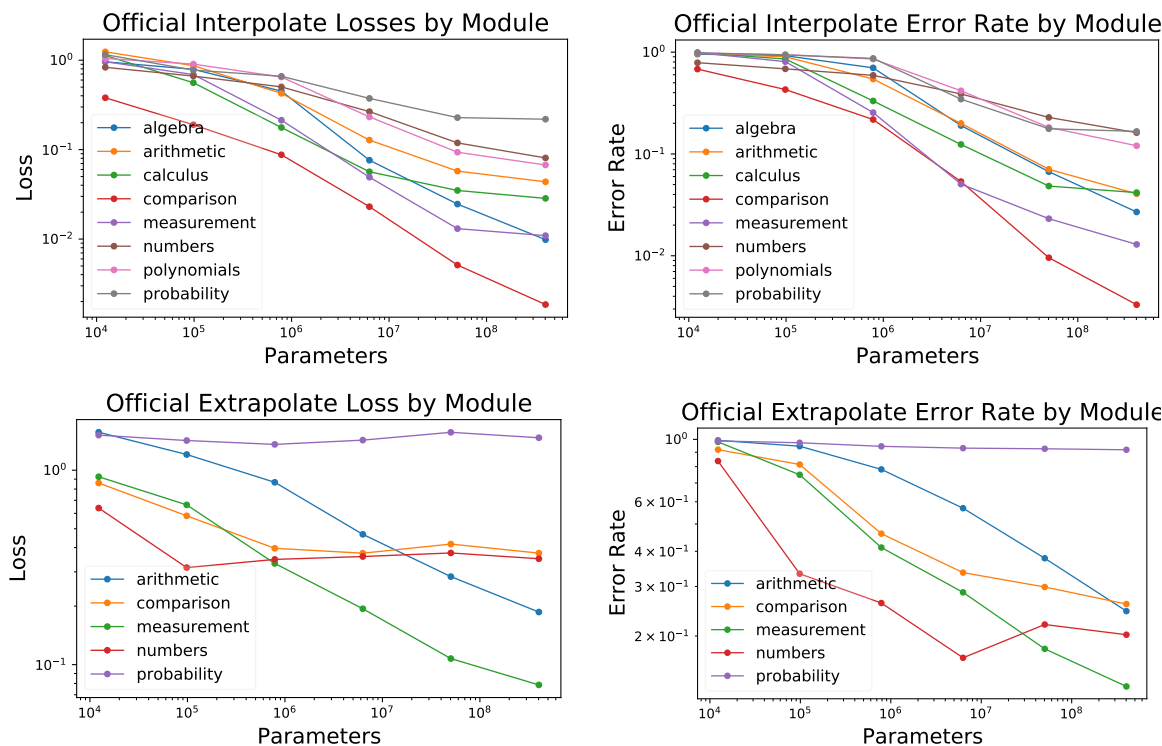


Figure 23 Math problem types— Here we show the performance of the math models on various modules of the math dataset, using the ‘official’ files of problems provided by [SGHK19]. The interpolate problems may have been seen by the models during training, as our training set was procedurally generated. We note that the losses on individual modules are approximate power-laws with model size on most of the interpolate modules, and on two the extrapolate modules.

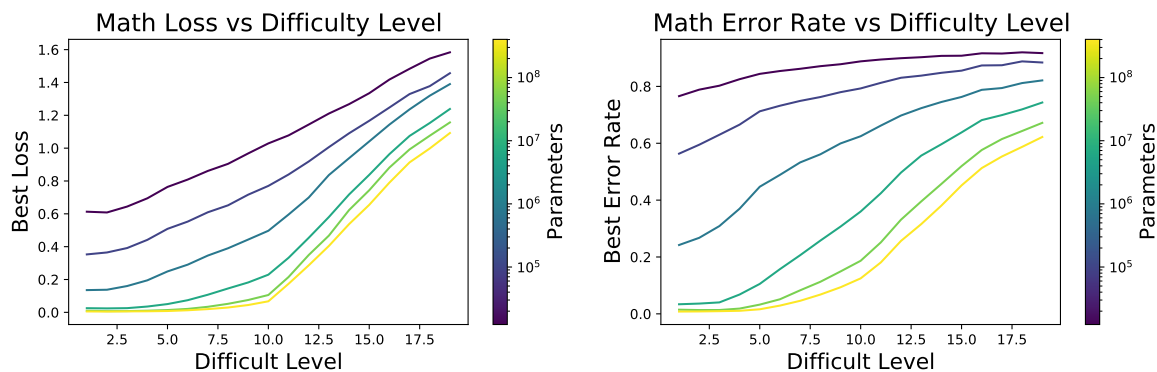


Figure 24 Math difficulty levels— Here we show how performance of math models varies as a function of the difficulty level or ‘entropy’ of the problem distribution, with levels ≤ 10 represented in the training distribution. We note an observable kink at level 10, suggesting some degree of overfitting, though as we extrapolate to more difficult problems the performance varies smoothly. It is clear that larger models perform better.

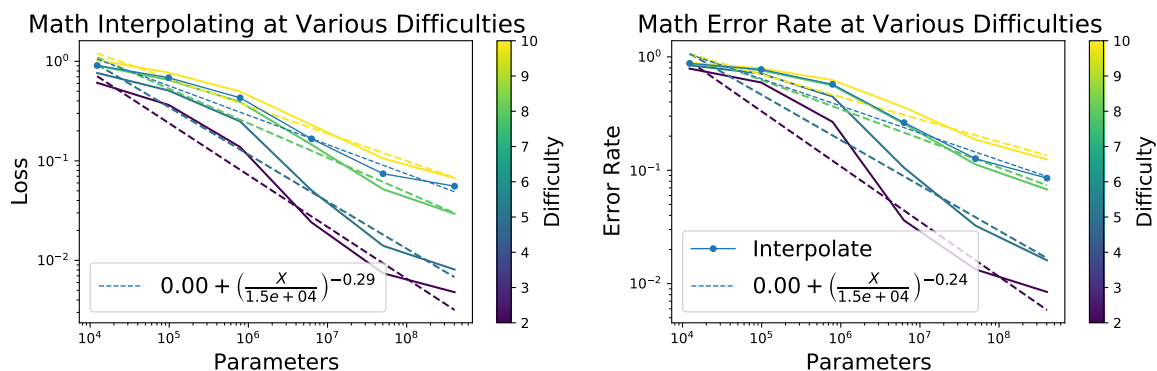


Figure 25 Model size trends for math difficulty levels— These plots show trends for the official interpolate dataset, as well as several difficulty levels that are within the training distribution. We observe that the power-law trends are distorted, perhaps as a consequence of memorization and the implicit curriculum in the data distribution.

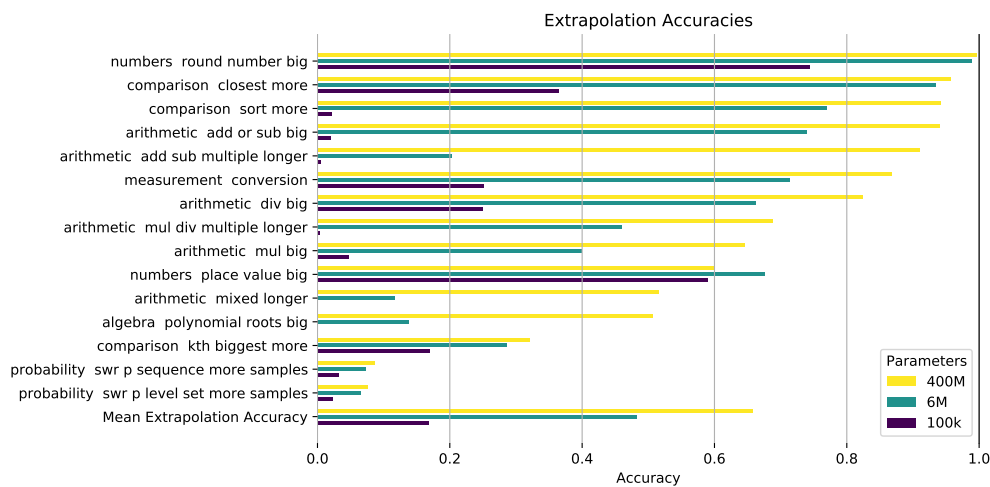


Figure 26 Extrapolation results for all math problem types— Here we show accuracies achieved by models of three different sizes on the official extrapolation test set files from [SGHK19], grouped by problem generator. Performance almost always improves with model size, though as shown in figure 13, this is due to the fact that larger models achieve better training loss.

400M Parameter Model Interpolation Accuracies

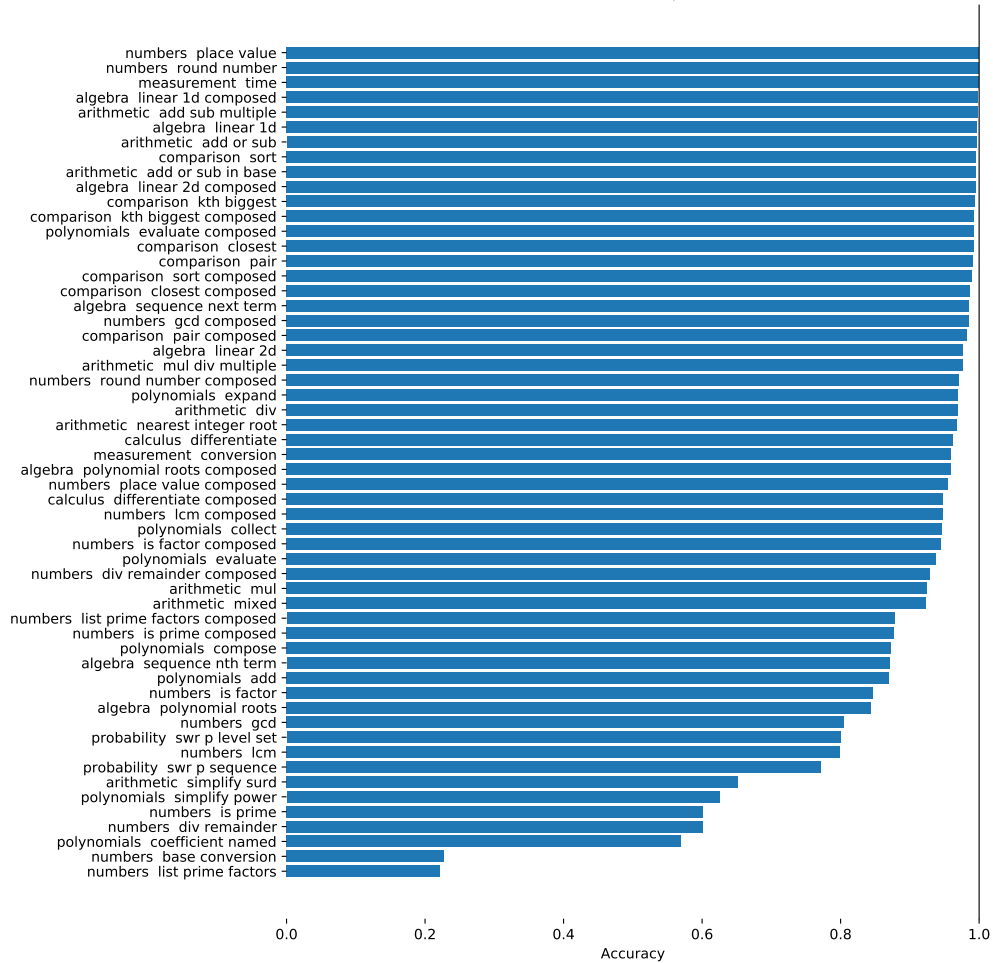


Figure 27 Interpolation results for all math problem types— Here we show interpolation accuracies achieved by a 400M parameter model, by problem generator. Note that these problems (files from [SGHK19]) were not deduplicated from our procedurally generated training set, so they may be contaminated by memorization.

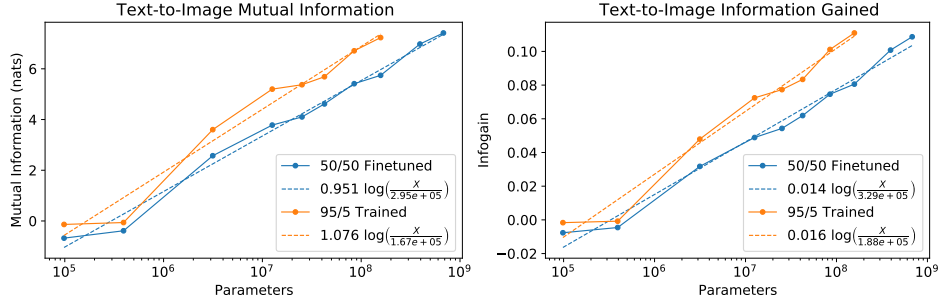


Figure 28 Mutual information In this plot we show the empirical mutual information for text-to-image multimodal models, as well as the Infogain, or the mutual information divided by the empirical entropy of the text.

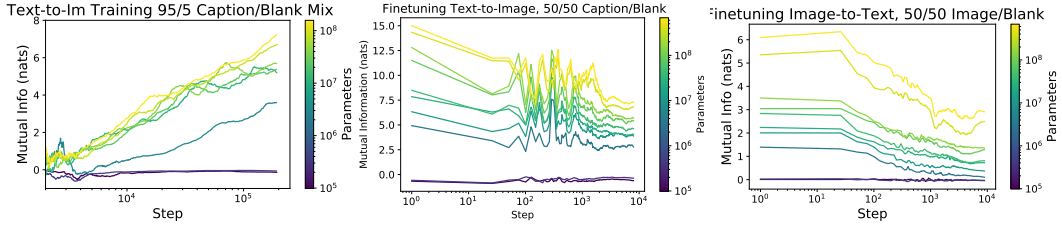


Figure 29 Mutual information learning curves— Here we show learning curves for the mutual information when either training or finetuning on a mixture of data with and without captions or images. We include training and finetuning on mixtures in order to ensure that our mutual information and Infogain estimates are not confounded by issues with blank captions or images being out of distribution.

C Additional Multimodal Results

Here we show a few additional results on the multimodal experiments. The learning curves for the mutual information are shown in figure 29. This includes both training from scratch on a 95/5 mixture of captioned and blank-caption data for text-to-image, as well as finetuning for 10k steps on a 50/50 mixture for both multimodal directions. We compare the final mutual information and infogain for the two strategies in figure 28; they are very similar.

Q: What is 6 times 4? A: 6 times 4 is _____

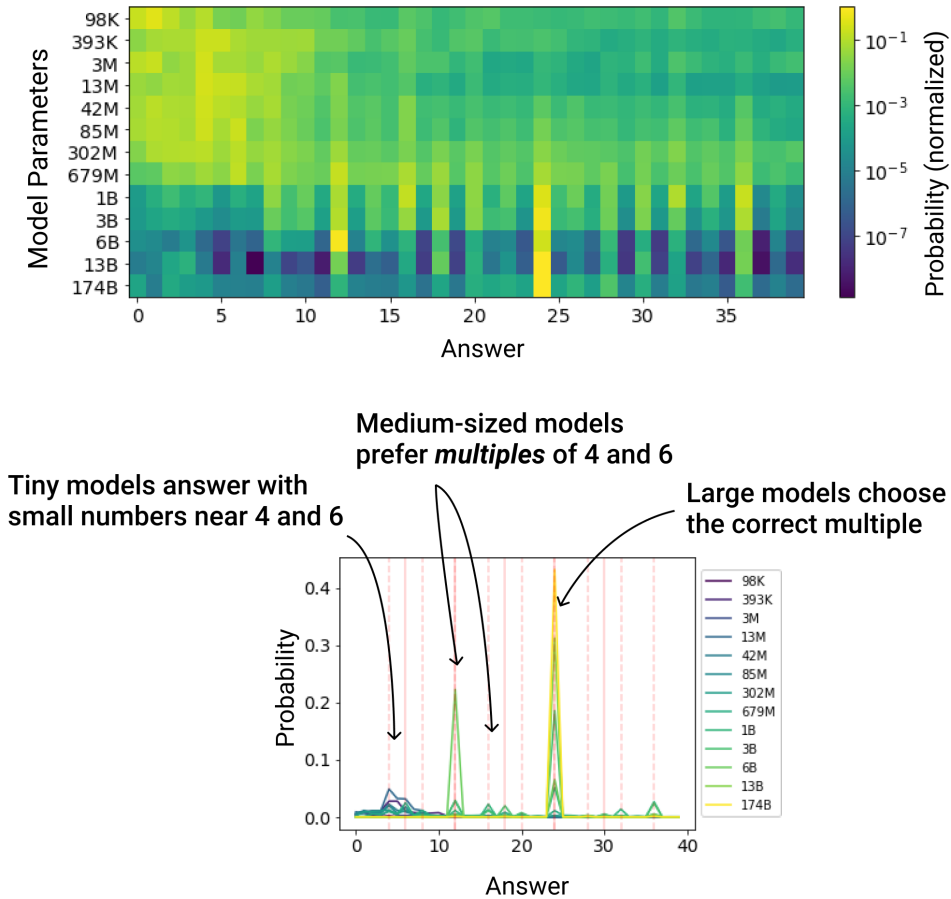


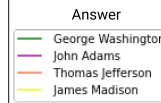
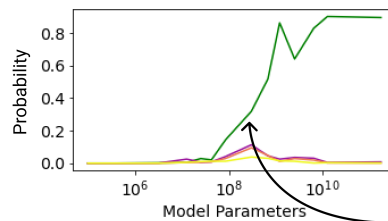
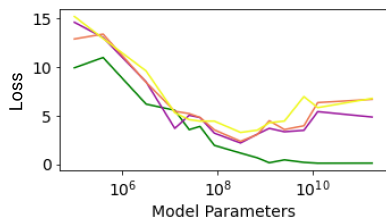
Figure 30 Arithmetic— We show the progression of arithmetic capabilities of GPT-3 family models as we increase the parameter count [BMR⁺20]. We measure the probability of different numeric answers for a simple multiplication problem. On the top we show a heat-map of normalized probabilities for each model size, and on the bottom we show a line chart of un-normalized probabilities. The smallest models put some weight on small numbers near those in the question. Somewhat larger models start to put some weight on multiples on 4 and 6 (visible as bright vertical streaks on the heat map, and marked as red lines on the line plot), suggesting that they’ve started to understand the meaning of the multiplication question. The largest models choose the correct answer confidently.

D Additional Language Results

Here we show a few additional results on the language experiments that measure how performance improves with parameter count. In figure 30, we investigate the progression of arithmetic capabilities, and in figure 31 we measure the ability to answer a simple factual question. In both cases we find smooth improvement in the loss on the correct answer as the model size increases. However, we also observe some qualitative “phases of learning”, with small models having difficulty understanding the question being asked of them, larger models showing some rudimentary understanding, and the largest models correctly answering the questions.

Q: Who was the **first** president of the United States?

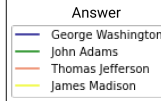
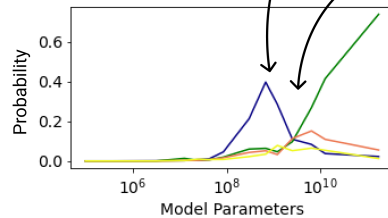
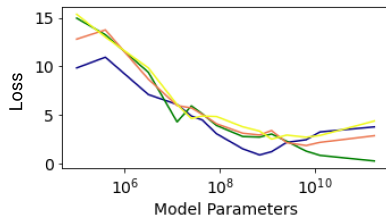
A: The **first** president of the United States was _____



100M-parameter models think **every** US president is George Washington

Q: Who was the **second** president of the United States?

A: The **second** president of the United States was _____



3B-parameter models understand **ordering** of presidents

Figure 31 Q&A— We show the progression of simple Q&A capabilities of GPT-3 family models as we increase the parameter count [BMR⁺20]. We ask the model who the first and second president of the United States was.

Tiny models appear to have trouble understanding the question, and don't place any significant probability on the correct answer. Larger models understand that we're requesting a US president, but fail to understand that the "second president" and "first president" are different requests, placing most of their weight for both questions on "George Washington". Only larger models understand both aspects of the questions, answering both correctly.

E Mutual Information, Infogain, and Scaling

We are studying the empirical mutual information

$$I(X, Y) = \mathbb{E}_{x, y \sim q} \left[\log \frac{p(x, y)}{p(x)p(y)} \right] \quad (\text{E.1})$$

where p is the model distribution and q is the true distribution of the data. This must be smaller than the cross-entropy loss of the model

$$L(X) = \mathbb{E}_{x \sim q} \left[\log \frac{1}{p(x)} \right] \quad (\text{E.2})$$

on either X or Y , so that the empirical InfoGain in equation 1.3 cannot be greater than one. As with the usual mutual information, the empirical mutual information is maximized when $y = f(x)$ or vice versa, so that the relation between X and Y is deterministic, and minimized when $p(x, y) = p(x)p(y)$.

However, it’s worth noting an interesting subtlety: in some cases it is possible for *our evaluations* to cause an apparent violation of the bound $\text{InfoGain} < 1$. This can occur in language models that are not precisely translation invariant when $x =$ the first T tokens while $y =$ the following tokens. For example, it’s theoretically possible that a language model with limited computational resources would assign a higher probability to “The MD5 hash of ‘powerlaw’ is e9f7a4aafeda67a0dab579ba480c24d6” than to the sequence “e9f7a4aafeda67a0dab579ba480c24d6” by itself.

E.1 Approximate Derivation of Scaling Relations

We do not know how to derive the relation 4.1 for multimodal models. However, we can derive a similar relation for the mutual information and infogain in language models. In this case, we study the mutual information between the first T tokens in a text sample, and the next T tokens (it is easy to generalize to sequences of different lengths).

We know that for a given model size N , the loss scales as a power-law with token position $t \geq 1$ [KMH⁺20]. In fact, we can roughly approximate

$$L(t) \approx L(N) + \frac{L_U - L(N)}{t^p} \quad (\text{E.3})$$

where $p < 1$ is a power, L_U is the unigram entropy, and p is roughly independent of N . This model is not perfect, but it permits a straightforward estimate of the empirical mutual information, namely

$$\begin{aligned} I([1, T], [T + 1, 2T]) &\approx (L_U - L(N)) \sum_{t=1}^T \left[\frac{1}{t^p} - \frac{1}{(t + T)^p} \right] \\ &= (2H_T^{(p)} - H_{2T}^{(p)})(L_U - L(N)) \end{aligned} \quad (\text{E.4})$$

where $H_T^{(p)}$ is the T th harmonic number with power p . We can evaluate or approximate $H_T^{(p)}$ if desired, but the point is that it’s identical for all N , and so the N -dependence of this expression comes only from $L(N)$. Because the exponent $\alpha_N \ll 1$ for language models, we can approximate $N^{-\alpha_N} \approx 1 - \alpha_N \log(N)$ to obtain equation 4.1.

Similarly, to approximate the infogain we need to divide by the loss on the final T tokens, so that

$$\text{Infogain} \approx \frac{(2H_T^{(p)} - H_{2T}^{(p)})(L_U - L(N))}{TL(N) + (H_{2T}^{(p)} - H_T^{(p)})(L_U - L(N))} \quad (\text{E.5})$$

Expanding this using $L(N) \propto N^{-\alpha_N} \approx 1 - \alpha_N \log(N)$ leads to the approximate formula from section 4. But more generally we see that the InfoGain is bounded by a certain ratio depending only on p and T , since $L(N)$ lies between 0 and L_U . So it will not actually approach 1.

E.2 Estimating D_{KL} Between Real-World Distributions

We have interpreted scaling trends in terms of the intrinsic entropy of the data distribution and the KL divergence between the true distribution and our models. This is based on the idea that with infinite data, model

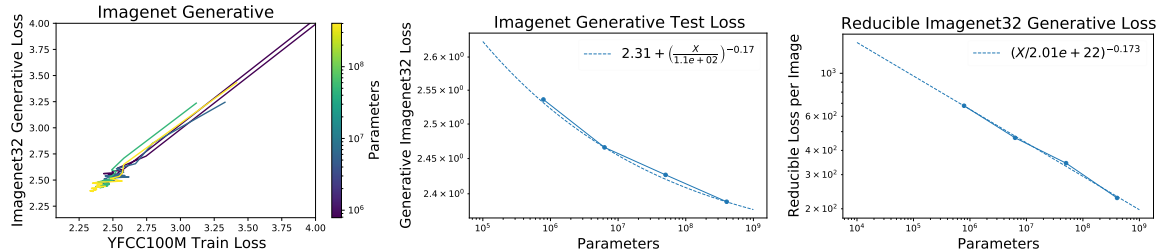


Figure 32 Generalization from YFCC100M to ImageNet generation— We show information about evaluating YFCC100M-trained models on ImageNet data distribution. On the left we show that loss on ImageNet depends only on the loss on YFCC100M, and does not otherwise depend on model size. In the center we show ImageNet loss vs model size, and on the right we subtract the irreducible loss to compute the reducible loss. These results strongly suggest that $L(N)$ follows a power-law plus constant form, even somewhat off of the training distribution.

size, and compute we could model the data distribution exactly. If the empirical loss of our models on a new data distribution also follows a predictable scaling trend, then this means we can estimate the fundamental KL divergence between the new distribution and the training distribution. Since our models were trained on YFCC100M images [TSF⁺ 15], it’s interesting to examine the trends for the loss on ImageNet, as we would expect in the infinite limit

$$L(\text{ImageNet}) = D_{\text{KL}}(\text{ImageNet}||\text{YFCC100M}) + S(\text{ImageNet}) \quad (\text{E.6})$$

where on the left we have the cross-entropy loss on ImageNet for a model trained on YFCC100M. We show the loss $L(N)$ when evaluating on ImageNet in figure 32, where we see that it appears to follow a power-law plus constant trend. Unfortunately this isn’t enough to identify $D_{\text{KL}}(\text{ImageNet}||\text{YFCC100M})$ because we also need a separate estimate of $S(\text{ImageNet})$, but our techniques are not easily applied there due to overfitting. But this quantity might be extracted by studying dataset size scaling in the future.

F Hyperparameter Settings

Here we include more details on the hyperparameter settings used to train the models.

All models used a learning rate schedule with a 3000 step linear warm-up followed by a linear decay to 1/10 of the maximum learning rate. Model hyperparameters and learning rates are shown in tables 4 and 5. The number of attention heads was always chosen to be $\max(2, d_{\text{model}}/64)$. Most models were trained with roughly 5×10^5 tokens per batch; differences from this are noted in the captions of the tables below. ‘Parameters’ always refers to the non-embedding parameter counts, and is approximate (we do not include biases for simplicity).

All models were trained for at least 250k steps (parameter updates), but many models were trained for significantly longer, as we noted that they had not yet reached the compute-efficient frontier, or did not seem to have converged. Trends in the loss as a function of model size were computed at the step minimizing the test loss. We used very similar learning rates for all models of a given size; these were determined through an initial grid search.

Parameters	d_{model}	n_{layer}	Max LR	Image-to-Text?
98304	64	2	0.00164	✓
393216	128	2	0.00144	✓
3145728	256	4	0.00115	✓
12582912	512	4	0.000959	✓
25165824	512	8	0.000862	✓
42467328	768	6	0.000789	✓
84934656	768	12	0.000692	✓
157286400	1280	8	0.000606	✓
393216000	1280	20	0.000479	
679477248	1536	24	0.000402	

Table 4 Multimodal hyperparameter settings— All Text-to-Image model settings are shown, the Image-to-Text models used identical settings, but the two largest models were not trained. ‘Parameters’ refers to the non-embedding parameter counts, and is approximate (we do not include biases for simplicity). These models were all trained with a batch size of 128 text/image pairs, or 409600 tokens per batch.

Parameters	d_{model}	n_{layer}	Max LR
1.23e+04	32	4	0.002686
9.83e+04	64	8	0.001597
7.86e+05	128	16	0.000950
6.29e+06	256	32	0.000565
5.03e+07	512	64	0.000336
4.03e+08	1024	128	0.000200
3.22e+09	2048	256	0.000119

Table 5 Math, Image, and Video modeling hyperparameter settings— ‘Parameters’ refers to the non-embedding parameter counts, and is approximate (we do not include biases for simplicity). The math models used $n_{\text{ctx}} = 512$ and a batch size of 524,288 tokens per batch. Video models used a batch size of 128 video clips, for a total of 524,288 tokens per batch. All image models used a batch size of 128 images, so the batch sizes in tokens vary depending on image or VQ resolution. We did not train the largest model size in some domains.

References

- [BMR⁺20] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020, 2005.14165. 3, 5, 6, 7, 8, 10, 15, 18, 19, 20, 31, 32
- [CGRS19] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *CoRR*, abs/1904.10509, 2019, 1904.10509. URL <http://arxiv.org/abs/1904.10509>. 6, 7, 20
- [CLH17] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the CIFAR datasets. *CoRR*, abs/1707.08819, 2017, 1707.08819. URL <http://arxiv.org/abs/1707.08819>. 4, 14, 15
- [CRC⁺20] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *Proceedings of Machine Learning and Systems 2020*, pages 10466–10478. 2020. 3, 19
- [DJP⁺20] Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. Jukebox: A generative model for music, 2020, 2005.00341. 8
- [HDMB19] Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality decomposed: how do neural networks generalise?, 2019, 1908.08351. 16
- [HNA⁺17] 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, 2017, 1712.00409. 3, 19
- [JGH18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018. 19
- [KMH⁺20] 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, 2001.08361. 3, 4, 5, 6, 7, 9, 10, 11, 17, 18, 19, 20, 33
- [Kom19] Aran Komatsuzaki. One epoch is all you need, 2019, arXiv:1906.06669. 19
- [LBD⁺20] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism, 2020, 2003.02218. 19
- [LH17] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101, 2017, 1711.05101. URL <http://arxiv.org/abs/1711.05101>. 6
- [LSP⁺18] Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. *arXiv:1801.10198 [cs]*, 2018, 1801.10198. URL <http://arxiv.org/abs/1801.10198>. 3, 20
- [LWS⁺20] Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joseph E. Gonzalez. Train large, then compress: Rethinking model size for efficient training and inference of transformers, 2020, 2002.11794. 3, 19
- [LXS⁺19] Jaehoon Lee, Lechao Xiao, Samuel S. Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent, 2019, arXiv:1902.06720. 19
- [MBB17] Siyuan Ma, Raef Bassily, and Mikhail Belkin. The power of interpolation: Understanding the effectiveness of SGD in modern over-parametrized learning. *CoRR*, abs/1712.06559, 2017, 1712.06559. URL <http://arxiv.org/abs/1712.06559>. 10
- [MKAT18] Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. An empirical model of large-batch training, 2018, arXiv:1812.06162. 10
- [RDG⁺20] Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric M. Smith, Y-Lan Boureau, and Jason Weston. Recipes for building an open-domain chatbot, 2020, 2004.13637. 3, 19

- [RFCS20] Jonathan S. Rosenfeld, Jonathan Frankle, Michael Carbin, and Nir Shavit. On the predictability of pruning across scales, 2020, 2006.10621. 19
- [RRBS19] Jonathan S. Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction of the generalization error across scales, 2019, 1909.12673. 3, 19
- [SGHK19] David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical reasoning abilities of neural models. *CoRR*, abs/1904.01557, 2019, 1904.01557. URL <http://arxiv.org/abs/1904.01557>. 3, 8, 16, 20, 26, 27, 28, 29
- [SK20] Utkarsh Sharma and Jared Kaplan. A neural scaling law from the dimension of the data manifold, 2020, 2004.10802. 3, 19
- [SSF⁺19] Imanol Schlag, Paul Smolensky, Roland Fernandez, Nebojsa Jojic, Jürgen Schmidhuber, and Jianfeng Gao. Enhancing the transformer with explicit relational encoding for math problem solving, 2019, 1910.06611. 20
- [TBL⁺19] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2019, page 6558. NIH Public Access, 2019. 3, 20
- [TSF⁺15] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. The new data and new challenges in multimedia research. *CoRR*, abs/1503.01817, 2015, 1503.01817. URL <http://arxiv.org/abs/1503.01817>. 3, 4, 5, 7, 13, 34
- [vdOKK16] Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. *CoRR*, abs/1601.06759, 2016, 1601.06759. URL <http://arxiv.org/abs/1601.06759>. 19
- [vdOVK18] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning, 2018, 1711.00937. 6, 7, 8, 11
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>. 3
- [WTU19] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models, 2019, 1906.02634. 3, 20