

SCALE-TIME EQUIVALENCE IN NEURAL NETWORK TRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural networks have demonstrated remarkable performance improvements as model size, training time, and data volume have increased, but the relationships among these factors remain poorly understood. We develop a theoretical framework to investigate the interplay between model size and training time, introducing the concept of *scale-time equivalence*: the idea that in certain training regimes, a small model trained for a long time can achieve similar performance to a larger model trained for a shorter time. To analyze this, we model neural network training as a gradient flow in random subspaces of larger non-linear models. In this setting, we show theoretically that neural network scale and training time can be traded off with each other. Empirically, we validate scale-time equivalence on MLPs and CNNs trained with gradient descent across standard vision benchmarks, CIFAR-10, SVHN, and MNIST.

We then investigate the consequences of scale-time equivalence on double descent, the phenomenon where model performance changes non-monotonically with respect to training data volume, model scale, and training time. In regimes where scale-time equivalence holds, we show that double descent with respect to training time and model scale may share a common cause: namely, overfitting to noise early during training. Through this, we provide potential explanations for several previously unexplained empirical phenomena: reduced data requirements for generalization in larger models, heightened sensitivity to label noise in overparameterized models, and instances where increasing model scale does not necessarily enhance performance.

1 INTRODUCTION

The remarkable progress in artificial intelligence (AI) has been driven by the rapid growth of model and dataset sizes. Scaling laws have provided useful insights into how test error decreases with increases in these factors, often following a power-law relationship. However, a key dimension remains underexplored: the role of training time in these scaling dynamics (Kaplan et al., 2020; Bahri et al., 2021; Rosenfeld et al., 2020; Sharma & Kaplan, 2022). Understanding this relationship is essential for optimizing performance under fixed compute budgets—for instance, whether smaller models trained for more epochs can match the performance of larger models trained briefly.

To better understand the relationship between model performance and training time, we introduce the concept of *scale-time equivalence*, which posits that increasing model size and extending training time are interchangeable to a degree. Using a theoretical model of gradient flow in random subspaces of larger nonlinear networks, we demonstrate that scaling up a network is functionally equivalent to training a smaller network for longer. Empirically, we validate this equivalence across standard vision benchmarks using MLPs and CNNs.

This framework enables us to shed light on another key phenomenon in learning: double descent. This behavior — where test error changes non-monotonically with respect to model size, training data volume, and training time — is associated with poor generalization near the interpolation threshold, where the model fits the training set exactly (Belkin et al., 2019; Nakkiran et al., 2021; Advani & Saxe, 2017). Despite its theoretical appeal, double descent remains inconsistently observed in practice, and existing theories leave several empirical phenomena unexplained. By leveraging scale-time equivalence, we provide new insights into double descent, suggesting it arises from the

interplay of model size, training time, and noise. Our framework explains key phenomena such as reduced data requirements for larger models, heightened sensitivity of overparameterized models to label noise, and why larger models can underperform their smaller counterparts in certain scenarios. The contributions of this paper are three-fold:

- Under a nonlinear, random subspace model of neural networks trained with gradient flow, we theoretically show *scale-time equivalence*: scaling model size is equivalent to increasing training time in a proportional manner. This equivalence holds exactly in the limit of very large networks trained for very small amounts of time.
- We empirically validate that scale-time equivalence can arise in MLPs and CNNs trained with gradient descent on standard vision benchmarks, demonstrating that network performance can be predicted by training much smaller models over extended training periods.
- Using scale-time equivalence, we propose a new explanation for parameter-wise double descent, linking it to noisy feature acquisition in smaller models that effectively “train slower” than larger counterparts.

2 RELATED WORK

2.1 SCALING LAWS

Neural network scaling laws describe how generalization error scales with data and model size. A number of works have observed power-law scaling with respect to data and model size (Kaplan et al., 2020; Rosenfeld et al., 2020; Clark et al., 2022), which has been explained theoretically (Bahri et al., 2021; Sharma & Kaplan, 2022; Hutter, 2021; Paquette et al., 2024).

However, other work demonstrates that model scale may not be sufficient to predict model performance (Tay et al., 2022), and casts doubt on power laws as the best model of error rate scaling (Alabdulmohsin et al., 2022; Bansal et al., 2022; Mahmood et al., 2022). Moreover, scaling with respect to training time, holding data volume fixed, remains poorly understood. These observations highlight the need for a more general framework that can predict model performance under many settings.

2.2 DOUBLE DESCENT

Double descent is an empirically observed phenomenon in which generalization error of machine learning models with respect to training data volume, model size, and training time exhibits an initial decrease, followed by a brief, sharp increase followed by a final decrease (Belkin et al., 2019; Nakkiran et al., 2021). Double descent with respect to model and data size has been theoretically understood as occurring due to a high degree of overfitting at the *interpolation threshold*, the point at which the model size is just sufficient to interpolate the training data (Adlam & Pennington, 2020; D’Ascoli et al., 2020; Belkin et al., 2019; Advani & Saxe, 2017). Typically, this work uses tools from random matrix theory to explain double descent for random feature models, in which linear regression maps a random, fixed feature pool to the desired output (Simon et al., 2024; Atanasov et al., 2024; Adlam et al., 2022; Bordelon et al., 2024; Maloney et al., 2022; Mei & Montanari, 2019; Ali et al., 2019; Lin et al., 2024).

However, these models generally do not explain double descent behavior in terms of training time. Originally, Nakkiran et al. (2019; 2021) hypothesized that training models for longer results in an effective increase in model size, thus allowing time-wise double descent to be explained in the same way as scale-wise double descent. Recent works have provided an alternate explanation: time-wise double descent occurs due to data features being learned at varying scales (Pezeshki et al., 2022; Heckel & Yilmaz, 2021; Stephenson & Lee, 2021). Error rises as models overfit to quickly-learned noisy features, but then falls as models more slowly learn signal features. In this work, we unify this explanation of time-wise double descent with the traditional account of double descent.

3 SCALE-TIME EQUIVALENCE IN NEURAL NETWORKS

In this section, we demonstrate that under certain training regimes, model size and training time may be traded off with each other. This result is consistent with and generalizes prior results demonstrat-

ing that models learn functions of increasing complexity over time (Nakkiran et al., 2019). We establish the result theoretically in a simplified model and validate it empirically in neural networks across several datasets and architectures.

3.1 RANDOM SUBSPACE MODEL VS. RANDOM FEATURES MODEL

We use a simplified model of neural networks, which we call the *random subspace model*, to theoretically analyze the relationship between model size and training time. Traditional analyses of scaling laws often treat neural networks as random feature models in which the network learns a linear combination of a fixed set of random features. Unfortunately, these random feature models produce network outputs that are *linear functions of the model parameters*, thus not fully capturing the nonlinear nature of neural networks.

Instead, we use a different set of assumptions that treats neural networks as random linear subspaces of larger nonlinear models. Crucially, the output of the model is a *nonlinear function of model parameters*. This framework avoids constraints on the architecture of a network.

Our framework operates by considering three nested linear subspaces: a large P dimensional space, a smaller p dimensional subspace, and an even smaller r dimensional subspace. All points in the P dimensional space correspond to parameters β of a large *nonlinear* neural network. We suppose that points in the smaller p dimensional subspace correspond to a smaller nonlinear neural network architecture (such that each point corresponds to a particular choice of parameters); we assume we have control over the parameters of this smaller network. This smaller network corresponds to the actual neural network we would like to model; we assume the larger P -dimensional network is unobserved. Note that smaller neural networks can naturally be viewed as linear subspaces of larger neural networks that contain them architecturally (i.e. have a larger width and depth) (Hall & Li, 1993). This relationship is common in model distillation, for example. We denote the parameters of our controllable p -dimensional neural network as $\theta \in \mathbb{R}^p$, where:

$$\beta = R\theta + \beta_0 \quad (1)$$

for a random matrix $R \in \mathbb{R}^{P \times p}$ with elements drawn iid from a unit Gaussian and with $\beta_0 \in \mathbb{R}^P$ fixed. Observe that any p dimensional affine subspace of \mathbb{R}^P may be represented in the form above for some choice of R and β_0 . Finally, suppose that the function represented by the smaller, p -dimensional neural network depends only on a low-dimensional linear projection of β of r dimensions. This is reasonable for neural networks: it is well known that neural networks trained by stochastic gradient descent tend towards flat minima of their loss landscapes, thus revealing many redundant dimensions in the network (i.e. only a low-dimensional subspace affects the network output). Specifically, we denote the low-dimensional projection as $\alpha \in \mathbb{R}^r$ where $r < p$ which is constructed as:

$$\alpha = K\beta \quad (2)$$

where $K \in \mathbb{R}^{r \times P}$ is a fixed projection matrix. We assume that the output of the p -dimensional network on any point is a *nonlinear* function of α (and therefore of the parameters θ as well). To keep our analysis tractable, we make the critical assumption that K is fixed over training, although we note that it may be more realistic to treat K as time-varying. We will represent time t with subscript t (i.e. β_0 denotes β at time 0).

To summarize our framework, we treat the observed p -dimensional neural network as lying in between two subspaces: the larger subspace of P -dimensional neural networks (which is unobserved), and the smaller subspace of functions represented by the network. Under these assumptions, we can show that under gradient descent, increasing the scale p of the model is approximately equivalent to increasing training time:

Theorem 3.1. *We denote the loss as a function of α : $L \in \mathbb{R}^r \rightarrow \mathbb{R}$. Suppose L has Lipschitz constant l and its second derivative has Lipschitz constant h . Suppose that continuous time gradient flow is applied to θ with learning rate η from initialization $\theta = 0$. We denote $\alpha_t = K(R\theta_t + \beta_0)$ where θ_t are the parameters at time t . Denote $A_t \in \mathbb{R}^r$ as the solution to:*

$$\dot{A}_t = -\eta K K^T \nabla L(A_t) \quad (3)$$

with initial condition $A_0 = K\beta_0$. Note that A_t does not depend on p . Then, with probability $1 - \epsilon$:

$$\|\alpha_t - A_{pt}\| \leq \frac{l\sqrt{\|K\|_F^4 + \|K K^T\|_F^2}}{h\sqrt{p\epsilon}\|K K^T\|} (e^{\eta p t h \|K K^T\|} - 1) \quad (4)$$

See Appendix A for a proof. The theorem implies closeness between the function implemented by the network, represented by α_t , and another quantity A_{pt} which *only depends on the product* pt . In other words, the learned model only depends on the product of the number of parameters p and time t (up to some error). Thus, we may interpret the product pt as representing the distance of a model along the training trajectory; larger pt implies more training progress, or equivalently, larger effective model capacity. Increasing the number of parameters by a scale factor is equivalent to increasing the training time by the same scale factor and vice versa: scale is *equivalent* to time. Intuitively, this is because each parameter allows the function to learn at a fixed rate; thus, adding more parameters linearly increases the effective learning rate.

The result also reveals when such a scale-time equivalence *cannot* be made. The error bound implies that when training progress pt is fixed, as p grows, scale and time become increasingly equivalent (the bound approaches zero). The equivalence breaks down for small p since here, the randomly chosen subspace of the model may or may not align well with K ; as p grows, the amount of alignment becomes less stochastic. Moreover, as training progress pt grows, the bound grows exponentially because small perturbations to the model early in the training trajectory lead to exponentially larger changes later in the trajectory. Thus, scale-time equivalence holds exactly in the limit as 1) model size p is large and 2) training progress pt is small.

Note that our result holds under standard neural network parameterizations where the learning rate is fixed and the gradient of the model output with respect to each parameter does not scale with p . In other parameterizations such as Neural Tangent Kernel (Jacot et al., 2018) or with variable learning rates, we expect a different form of scale-time equivalence.

We first validate our prediction in a simple linear model in which a varying fraction of the model parameters are controllable. See Appendix C for details. In Appendix D Figure 6, we find lines of 1:1 proportionality between model scale and the number of iterations required to reach a fixed loss level, validating our theory.

3.2 EMPIRICAL VALIDATION IN NEURAL NETWORKS

We next turn to examine whether scale-time equivalence can be found empirically in neural networks. We conduct experiments on MNIST (Deng, 2012), CIFAR-10 (Krizhevsky, 2009), and SVHN (Goodfellow et al., 2013) training a 7-layer convolutional neural network (CNN) and a 6-layer multilayer perceptron (MLP) with stochastic gradient descent (SGD). In these settings, the test loss is initially nearly constant, then suddenly starts decreasing after a set amount of training. Following this observation, we assess scale-time equivalence by measuring the *minimum* amount of training time required to achieve non-zero generalization under various network widths and dataset sizes (by subsampling the full dataset). Scale-time equivalence predicts that wider networks will require less time to generalize in a systematically predictable way. See Appendix C for further experimental details.

As observed in Figure 1, in all settings, we see a clear tradeoff curve between scale and training time: increasing scale by a fixed factor is nearly equivalent to reducing training time by another fixed factor. Importantly, the scale here is set as the *effective* number of network parameters, defined as the maximum number of training points that can be fit by the network, not the absolute number of parameters. We set the effective parameter count as the *cube root* of the number of parameters. Appendix B provides a heuristic argument for this scaling rate, although we highlight that effective parameter count can vary for different network architectures and learning settings.

Under this choice of scale, we find a systematic and predictable relationship between scale and training time, demonstrating that scale-time equivalence can be empirically observed. We also note that this phenomenon has been observed in prior literature (although not quantitatively characterized); for instance, Nakkiran et al. (2019) find that patterns of double descent are similar with respect to training epochs and network size. We emphasize that these results are limited to gradient descent with a fixed learning rate (the setting of our theory). With Adam optimization (Kingma & Ba, 2015), the number of epochs for generalization first decreases with scale, then increases (see Appendix D Figure 7). We hypothesize that since Adam has an adaptive learning rate, it is using a smaller effective learning rate for very large networks. As learning rate shrinks, more epochs are needed to generalize, leading to the observed results.

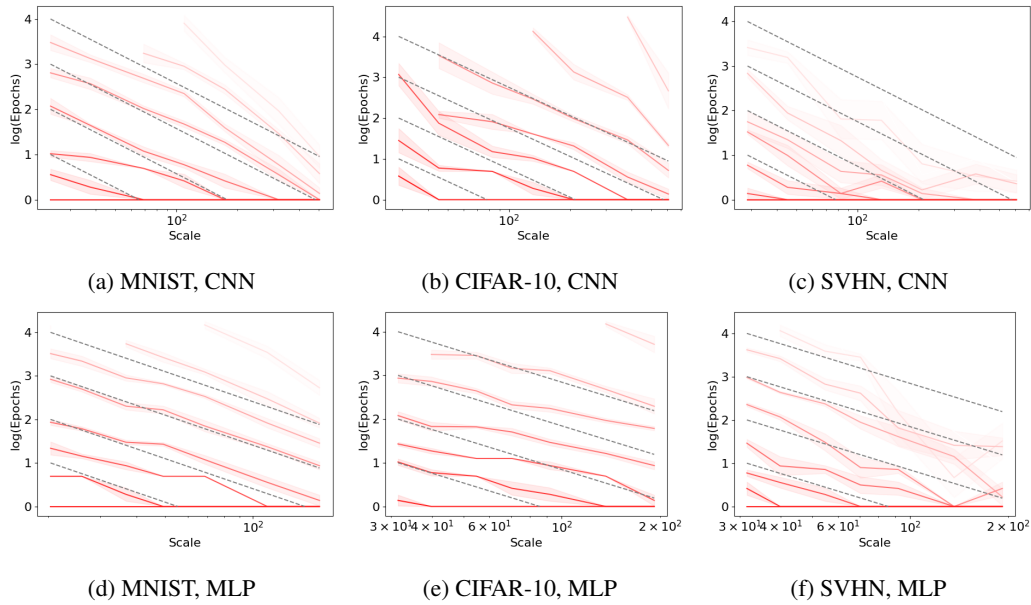


Figure 1: **Proportional trade-off between model scale and training time: testing the prediction on neural networks.** Red lines indicate tradeoff curves between number of training epochs and network scale for different datasets and architectures trained with SGD. Different curves indicate different amounts of training data; darker lines indicate more data. Curves are computed by, for each network scale, measuring the minimum amount of training time necessary to achieve non-zero generalization. Margins indicate standard errors over 5 trials. Grey curves are lines of 1:1 proportionality between scale and training epochs.

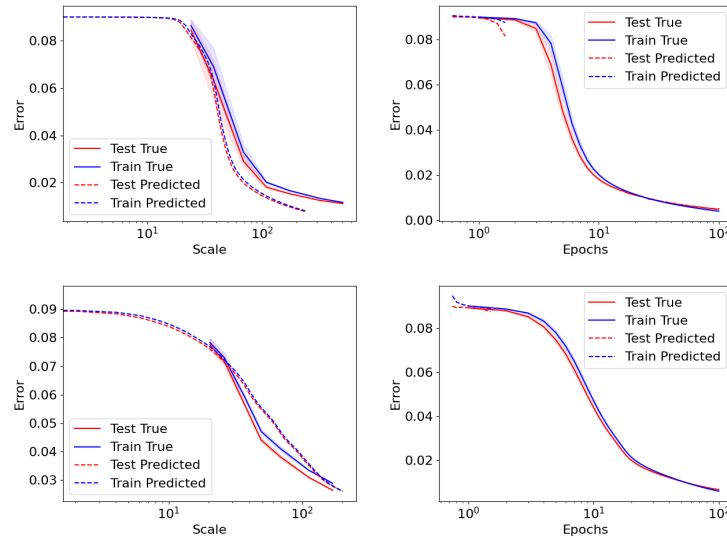


Figure 2: **By scale-time equivalence, small models trained for long times predict performance of large models trained for small times and vice versa: test of prediction.** Predicted and true test and train error of a CNN (top row) and MLP (bottom row) trained on MNIST. Column 1: predicting the performance of larger models over a few epochs by training smaller models for up to 100 epochs. Column 2: predicting performance of smaller models over many epochs by training larger models for 1 epoch. We use scale-time equivalence to predict the equivalent scale or number of epochs for each prediction. Margins indicate standard errors over 5 trials.

4 PREDICTING OPTIMAL NETWORK SCALE AND TRAINING TIME

In settings where scale-time equivalence holds, it should be possible to predict the performance of large models by training small models for many epochs and vice versa. This allows us to predict the optimal network scale and training time for a given dataset and base architecture. On benchmark datasets and architectures, we conduct two experiments: 1) predict performance under varying model scales from a small network trained for long training times, 2) predict performance over long training times by using larger networks trained for just 1 epoch. See Appendix C for experimental details and Figures 10 and 11 for full results. Figure 2 illustrates that scale-time equivalence can indeed be used to extrapolate performance on large scales and training times. Predictions of test and train performance under large model scales are particularly close to the true performance; notably, we can closely predict the scale at which generalization starts to occur. However, there is a small discrepancy between predicted and actual performance; we believe this can be corrected with dataset and model-specific tuning of the scale-time trade-off curve. Nevertheless, our findings reveal that scale-time equivalence can be used to predict optimal network scale and training time.

5 IMPLICATIONS OF SCALE-TIME EQUIVALENCE ON DOUBLE DESCENT

We next leverage scale-time equivalence to obtain a more unified understanding of the phenomenon of double descent, with respect to training time, parameter count and training set size.

5.1 ERROR SCALING OVER TIME IN A LINEAR MODEL

To provide context for the following section, we first summarize the existing results of Pezeshki et al. (2022); Heckel & Yilmaz (2021); Stephenson & Lee (2021); Schaeffer et al. (2023) by presenting a simple linear model that explains double descent with respect to time. Notably, these existing results have established that their analysis empirically predicts behavior in nonlinear neural networks as well, despite modeling them as linear. Consider a linear student-teacher setting in which training set outputs $Y \in \mathbb{R}^n$ are constructed as:

$$Y = Xw + \varepsilon \quad (5)$$

where training data $X \in \mathbb{R}^{n \times m}$, noise $\varepsilon \in \mathbb{R}^n$, $w \in \mathbb{R}^m$ is the unknown true model, n is the number of training points and m is the model dimensionality. We assume w is drawn independently from ε . Then, the parameters θ_t of a linear model learned after t time of gradient flow on mean squared error (with learning rate η) can be expressed as:

$$\theta_t = X^\dagger (I - e^{-\eta X X^T t}) Y \quad (6)$$

The resulting prediction error $x^T \theta_t - x^T w$ on a test point x can be expressed as:

$$x^T [X^\dagger (I - e^{-\eta X X^T t}) X - I] w + x^T X^\dagger (I - e^{-\eta X X^T t}) \varepsilon \quad (7)$$

We then use the singular value decomposition $U \Sigma V^T$ of X to express the prediction error:

$$\begin{aligned} x^T V [\Sigma^\dagger (I - e^{-\eta \Sigma \Sigma^T t}) \Sigma - I] V^T w + x^T V \Sigma^\dagger (I - e^{-\eta \Sigma \Sigma^T t}) V^T \varepsilon \\ = \sum_{i=1}^m -(x^T V)_i (V^T w)_i e^{-\eta \sigma_i^2 t} + (x^T V)_i (V^T \varepsilon)_i \frac{1 - e^{-\eta \sigma_i^2 t}}{\sigma_i} \end{aligned} \quad (8)$$

where σ_i are the singular values of X , and we denote $\sigma_i = 0$ for $i > n$ when $n < m$ (in this case, $\frac{1 - e^{-\eta \sigma_i^2 t}}{\sigma_i}$ denotes 0). Using the independence of w and ε , we finally may simply express the expected *squared* prediction error as:

$$\mathbb{E}[(x^T w - x^T \theta_t)^2] = \mathbb{E}[(\sum_{i=1}^m S_i e^{-\eta \sigma_i^2 t})^2] + \mathbb{E}[(\sum_{i=1}^m N_i \frac{1 - e^{-\eta \sigma_i^2 t}}{\sigma_i})^2] \quad (9)$$

where $S_i = -(x^T V)_i (V^T w)_i$, $N_i = (x^T V)_i (V^T \varepsilon)_i$. The first, signal term captures how well w can be learned in the absence of noise. In the underparameterized regime ($n > m$), this term approaches 0 as $t \rightarrow \infty$: without noise, the model can be learned perfectly. Observe that in general,

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

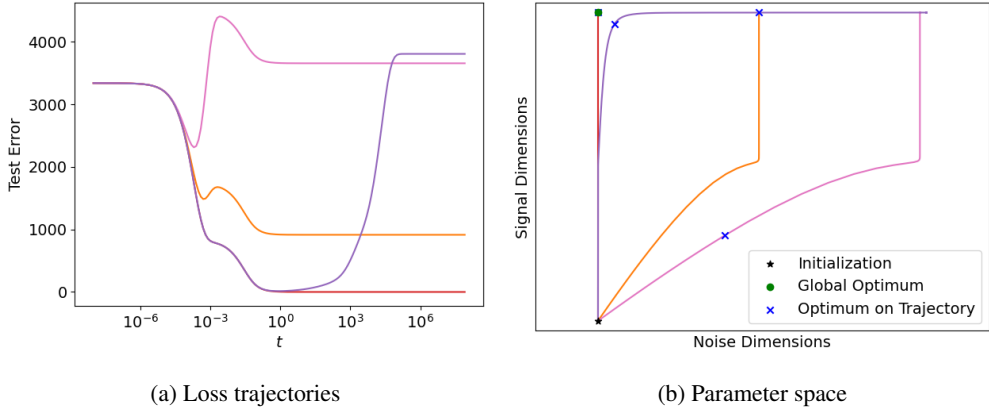


Figure 3: **Predicted variations in double-descent behavior, depending on training noise profile.** Schematic loss trajectories (a) and corresponding parameter space trajectories (b) of linear regression under various noise settings (different color curves). Depending on the noise profile, parameters may experience a temporary increase in error resembling an interpolation threshold.

the prediction error is *not* predicted to decay directly as a power law with t : instead, it decays or grows following a combination of exponential curves (though a combination of exponential decays at different rates can mimic a power law (Reed & Hughes, 2002)).

The second, noise term initially starts at 0 and grows over time. Notably, the size of the noise term is largest near the interpolation threshold (when $m \approx n$) since the smallest singular values σ_i will take small, non-zero values. As n or m grows larger than the other (distance from the interpolation threshold increases), the size of the noise term decreases. However, the noise term’s magnitude monotonically increases with t . For any given singular component of X , if the noise components are large relative to signal components (i.e. $|N_i| > |S_i|$), we expect that error will increase at the corresponding timescale, around $t = \frac{1}{\eta\sigma_i^2}$ and vice versa. Thus, if the small singular value components are noisy, error will increase later during training, while if large singular value components are noisy, error will increase early during training. Double-descent occurs when noisy components are acquired first (increase in error) followed by signal components (decrease in error). Figure 3(a) illustrates how different acquisition rates of noise vs. signal components may yield different loss trajectories, with some corresponding to double descent. We may also view these trends geometrically in the parameter space; Figure 3(b) shows that if signal and noise correspond to orthogonal parameter dimensions, then double descent (orange curve) corresponds to a setting in which noise dimensions are learned rapidly before signal dimensions. The optimal training time depends on which point in the parameter trajectory is closest to the true optimum and may occur either at the end of training or at an intermediate point.

5.2 EARLY NOISE ACQUISITION CAN EXPLAIN PARAMETER-WISE DOUBLE DESCENT

Next, we hypothesize in settings where scale-time equivalence holds, parameter-wise double descent can occur *due to the same mechanism* as time-wise double descent. Namely, in settings where noise is acquired before signal, smaller scale models (which require effectively larger training time by scale-time equivalence) fail to acquire signal. More precisely, by combining Equation 9 and scale-time equivalence, we propose the following novel scaling law in terms of the number of parameters p and time t by substituting t with pt :

$$error^2 = \mathbb{E}[(\sum_{i=1}^{\infty} S_i e^{-\eta\sigma_i^2 pt})^2] + \mathbb{E}[(\sum_{i=1}^{\infty} N_i \frac{1 - e^{-\eta\sigma_i^2 pt}}{\sigma_i})^2] \tag{10}$$

where σ_i depends on the number of training points n and S_i and N_i are random variables determining the strength of signal and noise respectively. This scaling law simultaneously explains double descent in terms of both p and t . Moreover, it can explain double descent in terms of data volume n as well: if n is set such that there are several small non-zero values of σ_i , the noise term becomes amplified. This explanation for double descent in n follows prior literature (Advani & Saxe, 2017).

378 The explanation for parameter-wise double descent differs from conventional wisdom in which dou-
 379 ble descent occurs due to the same reason as double descent in n : namely, when the number of
 380 model parameters is close to n , the model is highly sensitive to noisy directions in the training data
 381 (corresponding to small σ_i in a linear model), thus severely overfitting. How can we distinguish this
 382 explanation from ours? We propose three tests to separate the two hypotheses.
 383

384 5.3 LESS DATA REQUIRED FOR GENERALIZATION WITH MODEL SCALE

385 The interpolation threshold can be defined as
 386 the point when the number of data points equals
 387 the effective complexity of a model (Nakkiran
 388 et al., 2019). Thus, for any model, the loca-
 389 tion of the interpolation threshold with respect
 390 to data volume is the model’s effective capac-
 391 ity: higher capacity models have a rightward-
 392 shifted interpolation threshold with respect to
 393 data volume. Conventional double descent theo-
 394 ry argues that this point corresponds to a sharp
 395 decrease in the test set error as the amount of
 396 data grows. Starting from zero data, we would
 397 therefore expect that larger models (which have
 398 a more effective capacity), would experience a
 399 sharp decrease in test set error at *higher data*
 400 *volumes*: larger models require *more* data to
 401 generalize. By contrast, under our explanation,
 402 as model size grows, the amount of data needed
 403 to generalize *decreases*: since model size cor-
 404 responds to training time, generalization occurs
 405 more easily with larger models.

406 To test this, we conduct experiments on bench-
 407 mark datasets and architectures. Figure 4 re-
 408 veals that larger models indeed require less data
 409 to generalize, thus supporting our hypothesis. (See Appendix C for experimental details and Figure 8
 410 for full results.) Indeed, the *training set errors decrease* with data volume which is at odds with con-
 411 ventional double descent theory in which more data is always harder to interpolate. This is because
 412 conventional double descent theory typically assumes complete training convergence, which does
 413 not explain phenomena under the practically relevant setting of a fixed training budget. In summary,
 414 the increased ease of generalization with model scale does not neatly fit standard double descent
 415 theory, but readily fits in our explanation.

416 5.4 PERSISTENT VS LOCAL EFFECTS OF NOISE ON ERROR CURVES

417 Another key distinction is how the two explanations behave under varying levels of noise. Under
 418 the conventional explanation, as the noise ($|N_i|$) increases, the error increases mostly *locally around*
 419 $p \approx n$, near the interpolation threshold. This is because in Equation 9, the noise coefficients N_i
 420 multiply terms that are largest near the interpolation threshold and decrease as either data volume or
 421 model scale grow larger than the other (Schaeffer et al., 2023). In contrast, under our hypothesis, as
 422 the noise increases, the error monotonically grows with both model scale and time (see Equation 10).
 423 In other words, when the noise level changes, our hypothesis predicts a global performance change
 424 with respect to model scale while the conventional explanation predicts a primarily local change.
 425

426 We conduct experiments on benchmark datasets and architectures to validate this. We add label
 427 noise to the training set and evaluate performance as a function of epochs, model scale and data size.
 428 Appendix C includes experimental details. Noise leads to an increase in error in all cases, Figure 5
 429 (See Appendix D for additional plots), with a persistent or growing error with training time (epochs)
 430 and model scale. However, it leads to only localized increases in error in terms of the dataset size:
 431 in the case of MNIST and SVHN, there is nearly no change in error at the highest data volume.

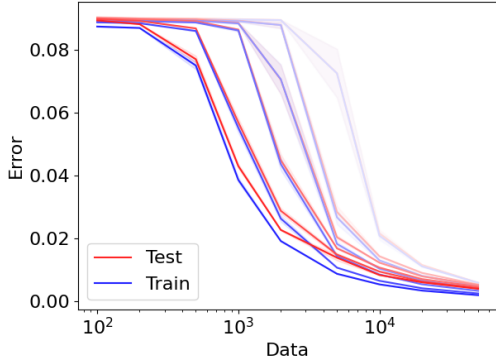


Figure 4: **Larger models require less data to interpolate: test of prediction.** Test and train mean squared error of CNN models trained on MNIST under varying levels of data. Different curves indicate different model scales; darker colors indicate larger models. Margins indicate standard errors over 5 trials.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

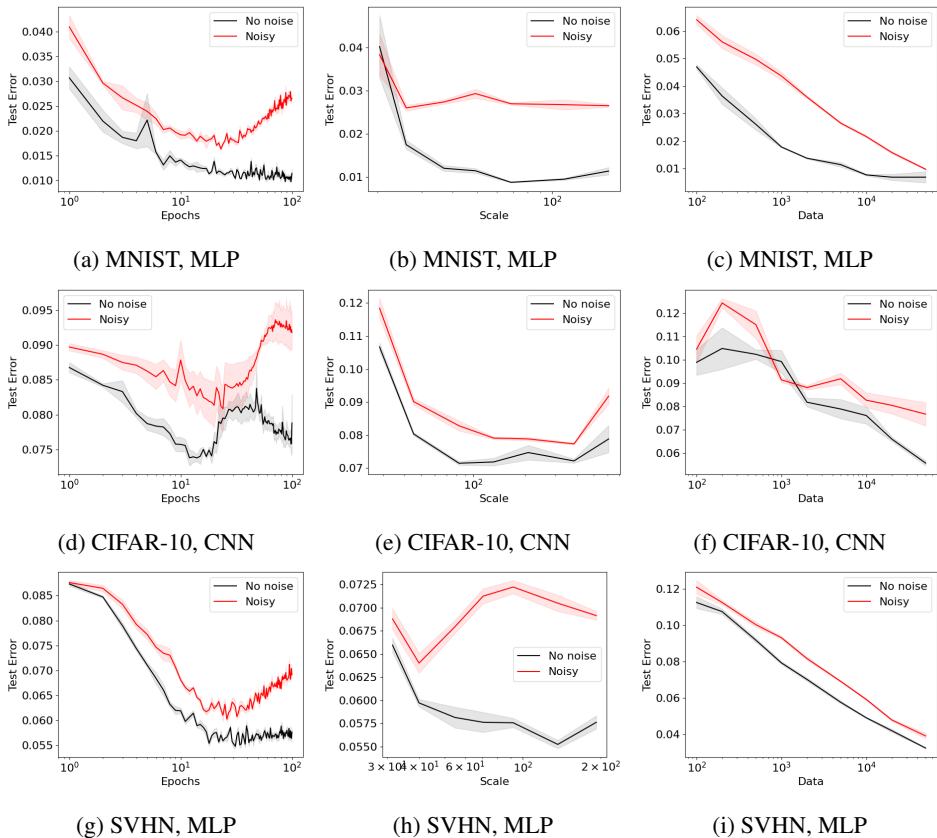


Figure 5: **Noise induces a persistent or growing error with training time and model scale, but not with dataset size.** Test mean squared error vs. number of epochs, model scale and training data under noisy and noise-free labels. Each row indicates a different combination of dataset and architecture. Margins indicate standard errors over 5 trials.

5.5 SHAPE OF LOSS CURVE: UPTURN WITH MODEL SIZE

A final key distinction between the two explanations is their prediction of the shape of the loss curve. Under the conventional explanation of double descent, past the interpolation threshold, increasing model size always improves performance. However, under our hypothesis, increasing model size may in certain cases *worsen performance*, if signal features are learned before noisy features. In Figure 5, we find that CIFAR-10 and SVHN indeed exhibit a U-shaped error curve with respect to model scale. This departure from a conventional double descent curve is consistent with our model.

6 DISCUSSION

We introduce scale-time equivalence during neural network training as phenomenon that emerges in regimes of large models trained for small amounts of time. In settings where this equivalence holds, we find that parameter-wise double descent arises from the same mechanism as epoch-wise double descent—namely, the early acquisition of noisy features during training. This perspective explains key empirical phenomena: larger models require less data for generalization, label noise significantly impacts overparameterized models, and increasing scale does not always improve performance. Despite its promise, our work has limitations. Like a number of scaling studies, our experiments focus on vision benchmarks, and extending the framework to language models is an important next step. Moreover, we use a heuristic to estimate the effective parameter count of neural networks; this heuristic warrants further investigation. Nevertheless, we believe our findings provide a step towards a more comprehensive understanding of the tradeoff between training time and model size, particularly in the regime of large models.

REFERENCES

- 486
487
488 Ben Adlam and Jeffrey Pennington. Understanding double descent requires a fine-grained bias-
489 variance decomposition. In *NeurIPS*, 2020. <https://arxiv.org/abs/2011.03321>.
- 490 Ben Adlam, Jake Levinson, and Jeffrey Pennington. A random matrix perspective on mixtures of
491 nonlinearities for deep learning. *AISTATS*, 2022.
- 492 Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural
493 networks. *arXiv preprint*, 2017.
- 495 Ibrahim Alabdulmohsin, Behnam Neyshabur, and Xiaohua Zhai. Revisiting neural scaling laws in
496 language and vision. In *NeurIPS*, 2022.
- 497 Alnur Ali, J Zico Kolter, and Ryan J Tibshirani. A continuous-time view of early stopping for least
498 squares regression. In *AISTATS*, 2019.
- 500 Alexander Atanasov, Blake Bordelon, Sabarish Sainathan, and Cengiz Pehlevan. The onset of
501 variance-limited behavior for networks in the lazy and rich regimes. *ICLR*, 2024.
- 502 Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural
503 scaling laws. *arXiv preprint*, 2021.
- 505 Yamini Bansal, Behrooz Ghorbani, Ankush Garg, Biao Zhang, Colin Cherry, Behnam Neyshabur,
506 and Orhan Firat. Data scaling laws in nmt: The effect of noise and architecture. In *ICML*, pp.
507 1466–1482. PMLR, 2022.
- 508 Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-
509 learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy*
510 *of Sciences*, 116(32):15849–15854, 2019.
- 511 Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural scaling
512 laws. *ICML*, 2024.
- 514 Aidan Clark, Diego de Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoff-
515 mann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, et al. Unified scaling
516 laws for routed language models. In *International Conference on Machine Learning*, pp. 4057–
517 4086. PMLR, 2022.
- 518 Stéphane D’Ascoli, Maria Refinetti, Giulio Biroli, and Florent Krzakala. Double trouble in double
519 descent: Bias and variance(s) in the lazy regime. In *ICML*, 2020.
- 521 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE*
522 *Signal Processing Magazine*, 29(6):141–142, 2012.
- 523 Ian J Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. Multi-digit number
524 recognition from street view imagery using deep convolutional neural networks. *arXiv preprint*,
525 2013.
- 526 Peter Hall and Ker-Chau Li. On almost Linearity of Low Dimensional Projections from High Di-
527 mensional Data. *The Annals of Statistics*, 21(2):867 – 889, 1993. doi: 10.1214/aos/1176349155.
528 URL <https://doi.org/10.1214/aos/1176349155>.
- 529 Reinhard Heckel and Fatih Furkan Yilmaz. Early stopping in deep networks: Double descent
530 and how to eliminate it. In *ICLR*, 2021. URL <https://openreview.net/forum?id=t1V90jvZbw>.
- 531
532
533 Marcus Hutter. Learning curve theory. *arXiv preprint*, 2021.
- 534 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gener-
535 alization in neural networks. *NeurIPS*, 31, 2018.
- 536 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
537 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
538 models. *arXiv preprint*, 2020.
- 539

- 540 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
541
- 542 Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
543
- 544 Licong Lin, Jingfeng Wu, Sham M Kakade, Peter L Bartlett, and Jason D Lee. Scaling laws in linear
545 regression: Compute, parameters, and data. *NeurIPS*, 2024.
546
- 547 Rafid Mahmood, James Lucas, David Acuna, Daiqing Li, Jonah Philion, Jose M Alvarez, Zhiding
548 Yu, Sanja Fidler, and Marc T Law. How much more data do i need? estimating requirements for
549 downstream tasks. In *CVPR*, pp. 275–284, 2022.
- 550 Alexander Maloney, Daniel A Roberts, and James Sully. A solvable model of neural scaling laws.
551 *arXiv preprint*, 2022.
552
- 553 Song Mei and Andrea Montanari. The generalization error of random features regression: Precise
554 asymptotics and double descent curve. *arXiv preprint*, 2019.
555
- 556 Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L Edelman, Fred
557 Zhang, and Boaz Barak. Sgd on neural networks learns functions of increasing complexity.
558 *NeurIPS*, 2019.
- 559 Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep
560 double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics:
561 Theory and Experiment*, 2021(12):124003, 2021.
562
- 563 Elliot Paquette, Courtney Paquette, Lechao Xiao, and Jeffrey Pennington. 4+ 3 phases of compute-
564 optimal neural scaling laws. *NeurIPS*, 2024.
565
- 566 Mohammad Pezeshki, Amartya Mitra, Yoshua Bengio, and Guillaume Lajoie. Multi-scale fea-
567 ture learning dynamics: Insights for double descent. In Kamalika Chaudhuri, Stefanie Jegelka,
568 Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th Inter-
569 national Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning
570 Research*, pp. 17669–17690. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.
571 press/v162/pezeshki22a.html](https://proceedings.mlr.press/v162/pezeshki22a.html).
- 572 William J. Reed and Barry D. Hughes. Power-law distributions from exponential processes : An
573 explanation for the occurrence of long-tailed distributions in biology and elsewhere. *Physical
574 Review E*, 66:4, 2002.
575
- 576 Jonathan S Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction
577 of the generalization error across scales. In *ICLR*, 2020.
578
- 579 Rylan Schaeffer, Mikail Khona, Zachary Robertson, Akhilan Boopathy, Kateryna Pistunova, Ja-
580 son W Rocks, Ila Rani Fiete, and Oluwasanmi Koyejo. Double descent demystified: Identifying,
581 interpreting & ablating the sources of a deep learning puzzle. *arXiv preprint*, 2023.
- 582 Utkarsh Sharma and Jared Kaplan. Scaling laws from the data manifold dimension. *JMLR*, 23:9–1,
583 2022.
584
- 585 James B Simon, Dhruva Karkada, Nikhil Ghosh, and Mikhail Belkin. More is better in modern
586 machine learning: when infinite overparameterization is optimal and overfitting is obligatory.
587 *ICLR*, 2024.
588
- 589 Cory Stephenson and Tyler Lee. When and how epochwise double descent happens. *arXiv preprint*,
590 2021.
591
- 592 Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan
593 Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale efficiently: Insights from
pre-training and fine-tuning transformers. In *ICLR*, 2022.

A PROOF OF THEOREM 1

First, observe that gradient descent applied to θ corresponds to:

$$\dot{\theta} = -\eta R^T K^T \nabla L(\alpha) \quad (11)$$

This yields a change in α of:

$$\dot{\alpha} = -\eta K R R^T K^T \nabla L(\alpha) \quad (12)$$

Properties of $K R R^T K^T$ First, we investigate the properties of the $r \times r$ random matrix $K R R^T K^T$. Since the elements of R are drawn iid from a unit Gaussian, $\mathbb{E}[R R^T] = pI$. Thus,

$$\mathbb{E}[K R R^T K^T] = p K K^T \quad (13)$$

Next, consider the element at the i th row and j th column of $K R R^T K^T$. This may be expressed as:

$$K_{i,:} R R^T K_{:,j}^T = \sum_l K_{i,:} R_{:,l} R_{l,:}^T K_{:,j}^T = \sum_l (K_{i,:} R_{:,l})(K_{j,:} R_{:,l}) \quad (14)$$

where we express different columns of R as $R_{:,l}$. Note that each term in the summand is independent from one another since R has iid elements. Observe that $K_{i,:} R_{:,l} = \sum_k K_{i,k} R_{k,l}$ and each term $K_{i,k} R_{k,l}$ is an independent mean zero Gaussian with variance $K_{i,k}^2$. Thus, $K_{i,:} R_{:,l}$ is a mean zero Gaussian with variance $K_{i,:} K_{:,i}^T$. Moreover, since the expectation of $K_{i,:} R R^T K_{:,j}^T$ is $p K_{i,:} K_{:,j}^T$, $K_{i,:} R_{:,l}$ and $K_{j,:} R_{:,l}$ are jointly Gaussian with covariance $K_{i,:} K_{:,j}^T$.

Now, consider the expectation of $(K_{i,:} R R^T K_{:,j}^T)^2$:

$$\mathbb{E}[(K_{i,:} R R^T K_{:,j}^T)^2] = \mathbb{E}\left[\sum_l (K_{i,:} R_{:,l})(K_{j,:} R_{:,l}) \sum_{l'} (K_{i,:} R_{:,l'})(K_{j,:} R_{:,l'})\right] \quad (15)$$

since $(K_{i,:} R_{:,l})(K_{j,:} R_{:,l})$ and $(K_{i,:} R_{:,l'})(K_{j,:} R_{:,l'})$ are independent for $l \neq l'$:

$$\mathbb{E}[(K_{i,:} R R^T K_{:,j}^T)^2] = \mathbb{E}\left[\sum_l (K_{i,:} R_{:,l})^2 (K_{j,:} R_{:,l})^2\right] + p(p-1) \mathbb{E}[(K_{i,:} R_{:,l})(K_{j,:} R_{:,l})]^2 \quad (16)$$

$\mathbb{E}[(K_{i,:} R_{:,l})(K_{j,:} R_{:,l})]$ is simply $K_{i,:} K_{:,j}^T$. Since $K_{i,:} R_{:,l}$ and $K_{j,:} R_{:,l}$ are jointly Gaussian, we may reparameterize them as:

$$K_{i,:} R_{:,l} = a z_i \quad (17)$$

$$K_{j,:} R_{:,l} = b z_i + c z_j \quad (18)$$

where z_i and z_j are independent unit Gaussians, $a = \sqrt{K_{i,:} K_{:,i}^T}$, $b = \frac{K_{i,:} K_{:,j}^T}{a}$, $c = \sqrt{K_{j,:} K_{:,j}^T - b^2}$.

Then, $(K_{i,:} R_{:,l})^2 (K_{j,:} R_{:,l})^2$ may be expressed as:

$$(K_{i,:} R_{:,l})^2 (K_{j,:} R_{:,l})^2 = (a z_i)^2 (b z_i + c z_j)^2 = a^2 z_i^2 (b z_i + c z_j)^2 = a^2 b^2 z_i^4 + 2 a^2 b c z_i^3 z_j + a^2 c^2 z_i^2 z_j^2 \quad (19)$$

Taking the expectation:

$$\mathbb{E}[a^2 b^2 z_i^4 + 2 a^2 b c z_i^3 z_j + a^2 c^2 z_i^2 z_j^2] = a^2 b^2 \mathbb{E}[z_i^4] + 2 a^2 b c \mathbb{E}[z_i^3 z_j] + a^2 c^2 \mathbb{E}[z_i^2 z_j^2] \quad (20)$$

Using the moments of unit Gaussians, we have $\mathbb{E}[z_i^4] = 3$, $\mathbb{E}[z_i^3 z_j] = 0$, $\mathbb{E}[z_i^2 z_j^2] = 1$. Thus, we have

$$\begin{aligned} \mathbb{E}[(K_{i,:} R_{:,l})^2 (K_{j,:} R_{:,l})^2] &= 3 a^2 b^2 + a^2 c^2 = a^2 (3 b^2 + c^2) \\ &= a^2 (K_{j,:} K_{:,j}^T + 2 \frac{(K_{i,:} K_{:,j}^T)^2}{a^2}) = K_{i,:} K_{:,i}^T K_{j,:} K_{:,j}^T + 2 (K_{i,:} K_{:,j}^T)^2 \end{aligned} \quad (21)$$

Finally, we may express the expectation of $(K_{i,:} R R^T K_{:,j}^T)^2$ as:

$$\mathbb{E}[(K_{i,:} R R^T K_{:,j}^T)^2] = p (K_{i,:} K_{:,i}^T K_{j,:} K_{:,j}^T + 2 (K_{i,:} K_{:,j}^T)^2) + p(p-1) (K_{i,:} K_{:,j}^T)^2 \quad (22)$$

Decomposition into expectation and noise We now decompose KRR^TK^T into an expectation term and a mean-zero noise term:

$$KRR^TK^T = pKK^T + N \quad (23)$$

where $N \in \mathbb{R}^{r \times r}$ is a mean-zero matrix. Note that each element of N has variance:

$$\mathbb{E}[N_{i,j}^2] = \mathbb{E}[(K_{i,:}RR^TK_{:,j}^T)^2] - p^2(K_{i,:}K_{:,j}^T)^2 = p(K_{i,:}K_{:,i}^TK_{j,:}K_{:,j}^T + (K_{i,:}K_{:,j}^T)^2) \quad (24)$$

Thus, the squared Frobenius norm of N has expectation:

$$\mathbb{E}[\|N\|_F^2] = \sum_{i,j} \mathbb{E}[N_{i,j}^2] = p \sum_{i,j} K_{i,:}K_{:,i}^TK_{j,:}K_{:,j}^T + (K_{i,:}K_{:,j}^T)^2 = p(\|K\|_F^4 + \|KK^T\|_F^2) \quad (25)$$

We express the dynamics of α as:

$$\dot{\alpha} = -\eta pKK^T\nabla L(\alpha) - \eta N\nabla L(\alpha) \quad (26)$$

Now suppose we have two copies of α : one copy ($\alpha^{(1)}$) with noise-free dynamics, and a second copy with noise:

$$\dot{\alpha}^{(1)} = -\eta pKK^T\nabla L(\alpha^{(1)}) \quad (27)$$

$$\dot{\alpha}^{(2)} = -\eta pKK^T\nabla L(\alpha^{(2)}) - \eta N\nabla L(\alpha^{(2)}) \quad (28)$$

We define the discrepancy between them as $\delta = \alpha^{(2)} - \alpha^{(1)}$, which has dynamics:

$$\dot{\delta} = -\eta pKK^T[\nabla L(\alpha^{(2)}) - \nabla L(\alpha^{(1)})] - \eta N\nabla L(\alpha^{(2)}) \quad (29)$$

Consider the rate of change of the ℓ_2 norm of δ :

$$\frac{d}{dt}\|\delta\| \leq \|\dot{\delta}\| \leq \eta p\|KK^T[\nabla L(\alpha^{(2)}) - \nabla L(\alpha^{(1)})]\| + \eta\|N\nabla L(\alpha^{(2)})\| \quad (30)$$

Using the Lipschitz bounds on ∇L and L :

$$\frac{d}{dt}\|\delta\| \leq \eta p h\|KK^T\|\|\delta\| + \eta\|N\|l \quad (31)$$

where the matrix norms in the expression denote ℓ_2 operator norm.

Now, returning to $\|N\|_F^2$, note that by Markov's inequality, with probability $1 - \epsilon$:

$$\|N\|_F^2 \leq \frac{p}{\epsilon}(\|K\|_F^4 + \|KK^T\|_F^2) \quad (32)$$

Using the fact that $\|N\| \leq \|N\|_F$, with probability $1 - \epsilon$:

$$\frac{d}{dt}\|\delta\| \leq \eta p h\|KK^T\|\|\delta\| + \eta \frac{\sqrt{p}}{\sqrt{\epsilon}} \sqrt{\|K\|_F^4 + \|KK^T\|_F^2} l \quad (33)$$

This is a differential inequality in $\|\delta\|$. Observe that $\|\delta\|$ takes its maximum possible trajectory at equality. Assuming $\|\delta\| = 0$ at time $t = 0$, this differential inequality may be solved by simply solving the equality case and setting the solution as the upper bound on δ_t (with subscript denoting time t):

$$\|\delta_t\| \leq \frac{\eta \frac{\sqrt{p}}{\sqrt{\epsilon}} \sqrt{\|K\|_F^4 + \|KK^T\|_F^2} l}{\eta p h\|KK^T\|} (e^{\eta p h\|KK^T\|t} - 1) = \frac{\sqrt{\|K\|_F^4 + \|KK^T\|_F^2} l}{\sqrt{p\epsilon} h\|KK^T\|} (e^{\eta p h\|KK^T\|t} - 1) \quad (34)$$

Solving the noise-free dynamics Now, we consider the noise-free dynamics:

$$\dot{\alpha}^{(1)} = -\eta pKK^T\nabla L(\alpha^{(1)}) \quad (35)$$

Note that A_{pt} solves the dynamical equation for $\alpha^{(1)}$ by the chain rule:

$$\frac{d}{dt}A_{pt} = p(-\eta KK^T\nabla L(A_{pt})) = -\eta pKK^T\nabla L(A_{pt}) \quad (36)$$

Thus, we may simply express $\alpha_t^{(1)}$ as A_{pt} .

Final result Combining and reexpressing our previous results, we may write that with probability $1 - \epsilon$:

$$\|\alpha_t^{(2)} - A_{pt}\| = \|\alpha_t^{(2)} - \alpha_t^{(1)}\| = \|\delta_t\| \leq \frac{l\sqrt{\|K\|_F^4 + \|KK^T\|_F^2}}{h\sqrt{p\epsilon}\|KK^T\|} (e^{\eta p t h \|KK^T\|} - 1) \quad (37)$$

Thus, the true dynamics deviate from A_{pt} by a bounded amount.

B MEASURING EFFECTIVE PARAMETER COUNT

In this section, we aim to provide a justification for why the *cube root* of the absolute number of parameters of a network is a good proxy for the *effective* number of parameters. Finding the effective number of parameters requires defining the parameter count at which the interpolation threshold occurs; this number is the effective parameter count. Unfortunately, in most models, training data is not *perfectly* interpolated by the model; thus defining the interpolation threshold location exactly is nontrivial. Instead, we identify a property of the interpolation threshold in the training error of linear models, then identify at which parameter count this property also holds in nonlinear models.

Training error in linear models Suppose we are provided a training set $X \in \mathbb{R}^{n \times p}$ and corresponding training labels $Y \in \mathbb{R}^n$ where n is the number of data points. Assume that X has elements drawn uniformly from a unit Gaussian. We aim to find a parameter $\theta \in \mathbb{R}^p$ such that:

$$Y \approx X\theta \quad (38)$$

The minimum norm solution minimizing the mean squared error is:

$$\theta = X^\dagger Y \quad (39)$$

where \dagger denotes pseudoinverse. Denoting the predicted training labels as \hat{Y} , the mean squared error is then:

$$\frac{1}{n} \|Y - \hat{Y}\|_2^2 = \frac{1}{n} \|Y - XX^\dagger Y\|_2^2 = \frac{1}{n} Y^T (I - XX^\dagger) Y \quad (40)$$

Finally, assuming that $Y_i^2 = 1$ for all i , we may write the mean training error as:

$$\frac{1}{n} \|Y - \hat{Y}\|_2^2 = \max(1 - \frac{p}{n}, 0) \quad (41)$$

Thus, for a linear model, the training error decreases linearly at rate $-\frac{1}{n}$ with respect to p before the interpolation threshold ($p = n$), and then is zero after the interpolation threshold. We extract a key property around the interpolation threshold from the linear model: the training error is $\frac{1}{n}$ at $p = n - 1 = O(n)$. Note that for any $\alpha < 1$, the training error is $\frac{1}{n^\alpha}$ before $O(n)$.

Power law training rate decay Next, we consider power-law decays of training rate error and characterize which power laws are consistent with the interpolation threshold properties outlined above. Consider a power law decay of the mean training error of $\frac{n^\alpha}{p^\beta}$. Note that in order to satisfy condition (2), we must have:

$$\frac{1}{n} = O\left(\frac{n^\alpha}{n^\beta}\right) \quad (42)$$

where we set $p = O(n)$. Thus, $\beta - \alpha = 1$. Power-law decays of the form $\frac{n^\alpha}{p^{\alpha+1}}$ are consistent with the interpolation threshold property.

Error scaling in neural networks Next, we turn to model training error scaling in neural networks. We make the following heuristic argument: to fit n training points, a network of width m needs to encode $O(mn)$ numbers corresponding to $m + 1$ numbers for each training point (to represent the features and labels of each training point). On the other hand, the network has $O(m^2)$ parameters since its intermediate layer weights have m^2 parameters. We expect that scaling both the network's capacity of $O(m^2)$ and the required capacity of $O(mn)$ at the same rate will not change the training error. Thus, we expect the training error to be a function of $\frac{mn}{m^2} = \frac{n}{m}$.

756 To find *which* function of $\frac{n}{m}$ models the training error, we introduce another argument. We assume
 757 that the network can be approximated as an ensemble of $O(m^2)$ submodels, each with $O(1)$ param-
 758 eters. Suppose the model output is the average of the outputs of the submodels. Then, assuming that
 759 the distribution of the submodel outputs on any given training point has variance $O(1)$ over different
 760 model initializations, the variance of the ensemble output is $O(\frac{1}{m^2})$. Thus, over model initialization,
 761 the model output at any given training point will center around a mean value with deviations on the
 762 order of $O(\frac{1}{m})$. Assuming that as the network capacity goes to ∞ , the training error goes to 0, the
 763 mean value of the predicted output on a training point must be the true value; thus, the predicted
 764 output differs from the true output by $O(\frac{1}{m})$. This corresponds to a mean squared error on the train-
 765 ing points scaling as $O(\frac{1}{m^2})$. Finally, to get the dependence on n , we use the observation that the
 766 error must depend on only the fraction $\frac{n}{m}$. Thus, the mean squared training error scales as $O(\frac{n^2}{m^2})$.
 767

768 Finally, we know by the argument above that power law error rate scaling of the form $\frac{n^\alpha}{p^{\alpha+1}}$ are
 769 consistent with the interpolation threshold property, where p is the effective number of parameters.
 770 Equating this with the neural network scaling result of $O(\frac{n^2}{m^3})$, we must have that $\alpha = 2$, yielding:
 771

$$772 \quad O\left(\frac{n^2}{p^3}\right) = O\left(\frac{n^2}{m^3}\right) \quad (43)$$

773
 774 Equating the denominators, we have $p^3 = m^3$, or $p = m$. In other words, the effective number
 775 of parameters is m . Since the absolute number of parameters in the network is m^2 , the effective
 776 number of parameters is the square root of the absolute parameter count.
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809

810 C EXPERIMENTAL DETAILS

811 C.1 MODEL TRAINING

812 We used three benchmark datasets: CIFAR-10, MNIST (Creative Commons license), and SVHN
813 (non-commercial use only). Each dataset was subject to preprocessing involving standard normal-
814 ization. For CIFAR-10 and SVHN, the normalization was performed using means and standard
815 deviations of (0.5, 0.5, 0.5). For MNIST, the normalization used a mean and standard deviation of
816 0.5.

817 Two types of neural network architectures were evaluated:

- 818 • Convolutional Neural Network (CNN): Our CNN architecture consisted of six ReLU acti-
819 vated convolutional layers, with kernel sizes (3, 3, 3, 3, 3, 2 for CIFAR-10/SVHN or 1 for
820 MNIST), strides (2, 1, 2, 1, 1, 1), and number of filters (5s, 10s, 20s, 40s, 80s) where s is a
821 width parameter. This was followed by a fully connected layer.
- 822 • Multilayer Perceptron (MLP): The MLP architecture comprised six fully connected layers
823 with hidden layer width 10s where s is a scale parameter.

824 The width parameter was set to values of 1, 2, 5, 10, 20, 50, and 100. When we refer to "model
825 scale", we quantify this as the cube root of the number of network parameters rather than the value
826 of the width parameter. We tested two learning rates (0.001 and 0.01) in combination with two
827 optimizers (Adam and SGD respectively). Label noise was introduced at levels of 0.0 (no noise)
828 and 0.2 (20% noise) to evaluate the robustness of the models. The number of training samples
829 was varied among 100, 200, 500, 1000, 2000, 5000, 10000, 20000, and 50000. The samples were
830 randomly selected from their respective base training sets. A constant batch size of 32 and 100
831 training epochs were used across all experiments. Mean Squared Error (MSE) was used as the loss
832 function. Model performance was assessed using Mean Squared Error (MSE) on both training and
833 test sets. For reproducibility, we set a manual seed for the PyTorch random number generator. Five
834 different seeds (101, 102, 103, 104, 105) were used to assess the variance in results due to random
835 initialization.

836 C.2 SCALE-TIME TRADEOFF VALIDATION

837 **Linear model** We construct a loss function as:

$$838 L(\theta) = \|K(R\theta + \beta_0) - \alpha^*\|^2 \quad (44)$$

839 where α^* is a target value of α . K , R , β_0 and α^* are all independently sampled from unit Gaussians,
840 and θ is initialized from a unit Gaussian. We set $P = 1000$ and $r = 3$. We train θ via gradient
841 descent with learning rate 10^{-6} and evaluate the number of iterations required to reach various
842 values of the loss for varying values of p . This setup corresponds to a linear classifier trained with
843 gradient descent on mean squared error using a training batch of $k = 3$ points where only p of the
844 $P = 1000$ parameters are controllable.

845 **Neural networks** To evaluate the scale time tradeoff, for each model, we compute the minimum
846 number of epochs necessary to achieve a test set MSE below 0.09.

847 C.3 SHIFTING INTERPOLATION THRESHOLD

848 Experimental results are shown for networks trained with SGD.

849 C.4 IMPACT OF NOISE AND SHAPE OF LOSS CURVE

850 Experimental results are shown for networks trained with Adam; since Adam trains faster, this
851 allows us to examine training trends at later points in training.

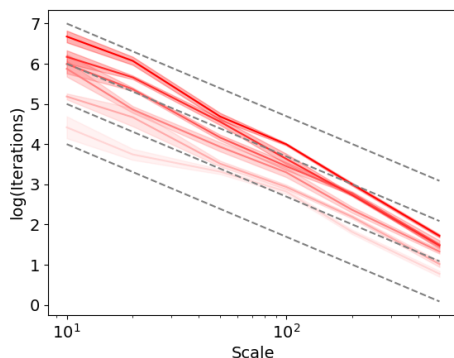
864 C.5 OPTIMAL SCALE PREDICTIONS

865
866 To predict error over scale using error over time or vice versa, we treat the quantity pt as a predictor
867 of performance, where p is the effective number of network parameters. Thus, if we wish to know
868 the error of a network at (p_1, t_1) and we know the error of the network for all p_0 , then we may
869 predict the error at (p_1, t_1) as the error at $(p_0, (\frac{p_1}{p_0})t_1)$.

871 C.6 COMPUTING INFRASTRUCTURE

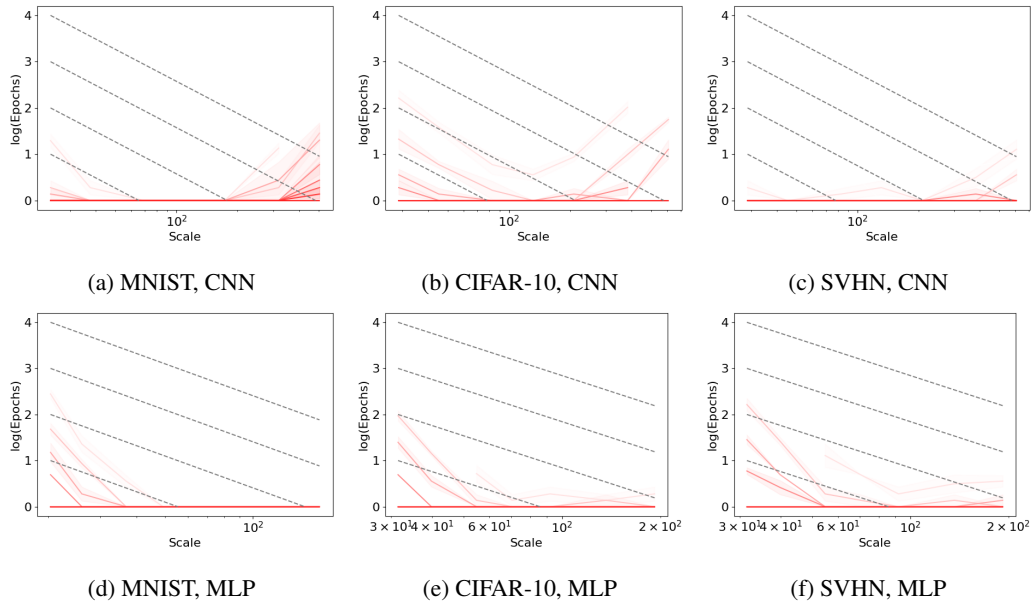
872 Experiments were run on a computing cluster with GPUs ranging in memory size from 11 GB to 80
873 GB.
874

876 D ADDITIONAL EXPERIMENTS



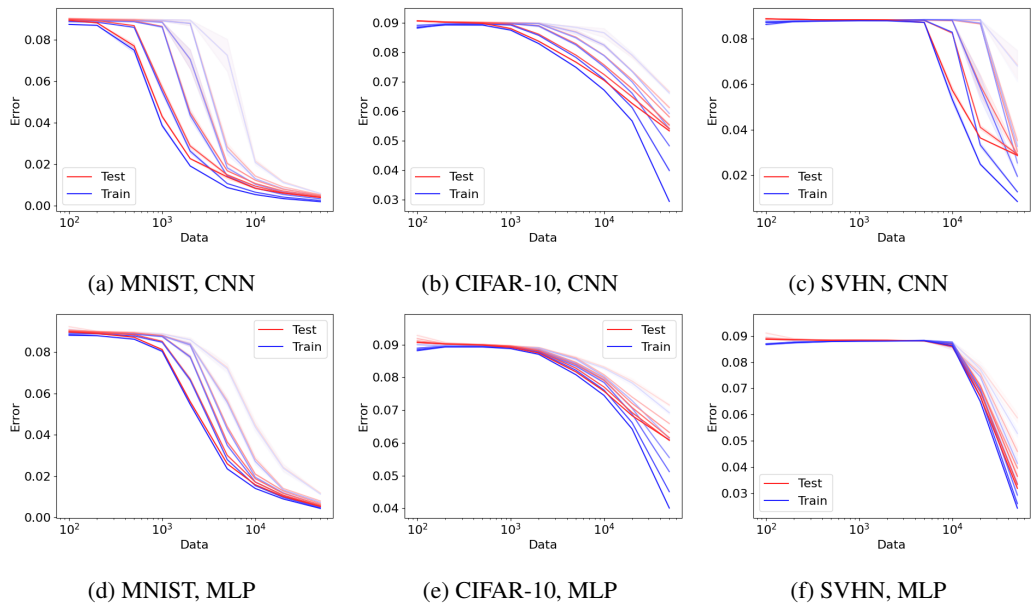
880
881
882
883
884
885
886
887
888
889
890
891
892 **Figure 6: Proportional trade-off between model scale and training time: testing the prediction**
893 **on a linear model.** Red lines indicate tradeoff curves between number of training iterations and
894 model size. Curves are computed by, for each model size, measuring the minimum amount of train-
895 ing time necessary to achieve different loss levels. Different curves indicate different performance
896 thresholds; darker lines indicate a smaller error threshold. Margins indicate standard errors over 5
897 trials. Grey dashed lines represent 1:1 proportionality between scale and training iterations.

900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917



939
940
941
942
943
944
945
946
947

Figure 7: Red lines indicate tradeoff curves between number of training epochs and network scale for different datasets and architectures trained with Adam. Different curves indicate different amounts of training data darker lines indicate more data. Curves are computed by, for each network scale, measuring the minimum amount of training time necessary to achieve non-zero generalization. Margins indicate standard errors over 5 trials. Grey curves are lines of 1:1 proportionality between scale and training epochs.



968
969
970
971

Figure 8: Test and train mean squared error of MLP and CNN models trained on benchmark datasets under varying levels of data. Different curves indicate different model scales; darker colors indicate larger models. Margins indicate standard errors over 5 trials.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

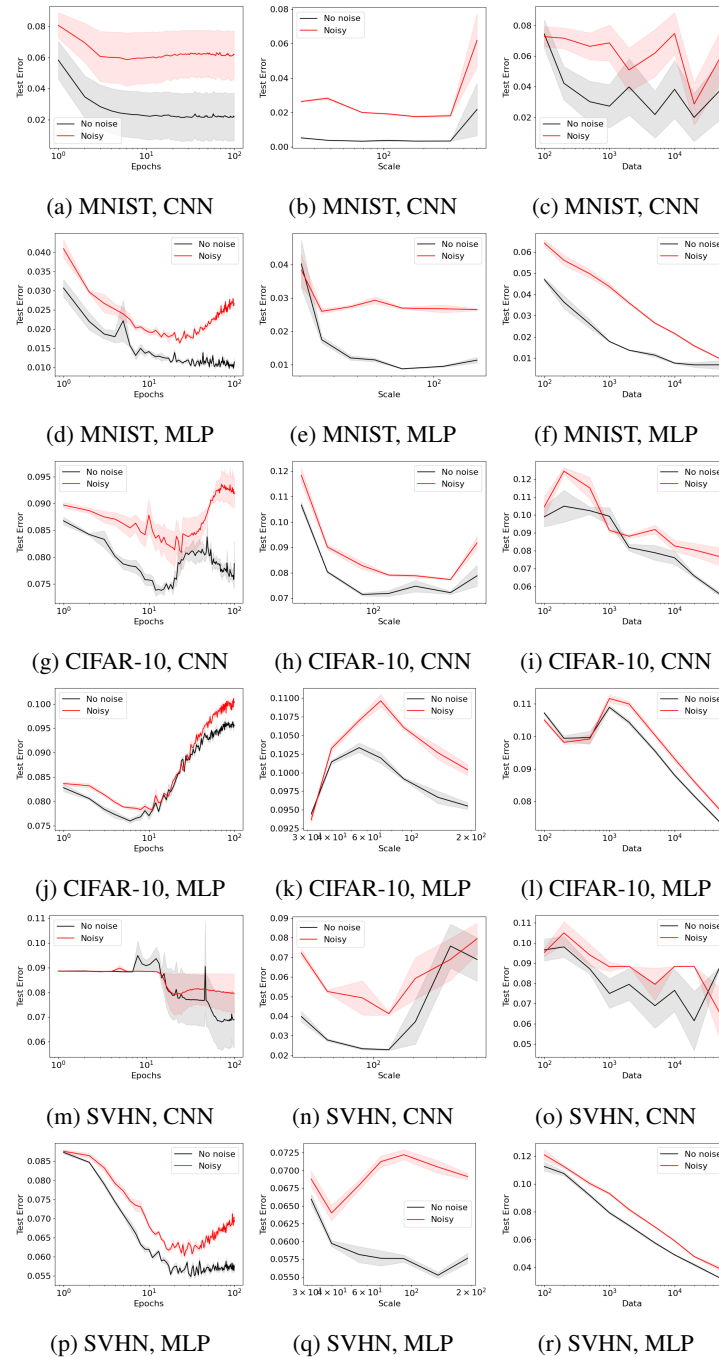


Figure 9: Test mean squared error vs. number of epochs, model scale and training data under noisy and noise-free labels. Each row indicates a different combination of dataset and architecture. Margins indicate standard errors over 5 trials.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

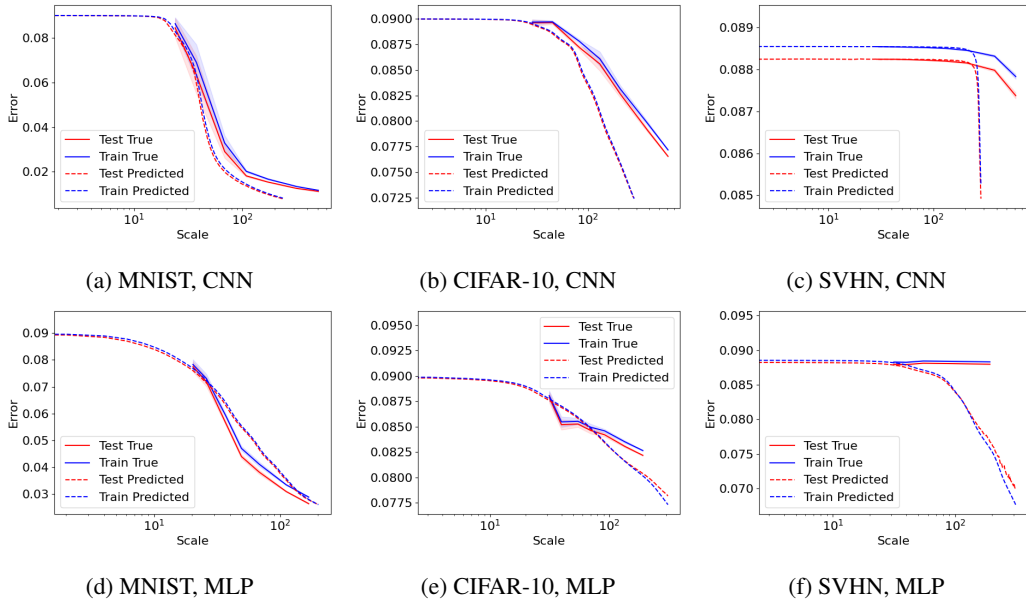


Figure 10: Predicted and true test and train mean squared error of MLP and CNN models trained on benchmark datasets under varying model widths for 10 epochs. Margins indicate standard errors over 5 trials. Predictions are generated by training a small model for 100 epochs and using scale-time equivalence to predict the equivalent scale.

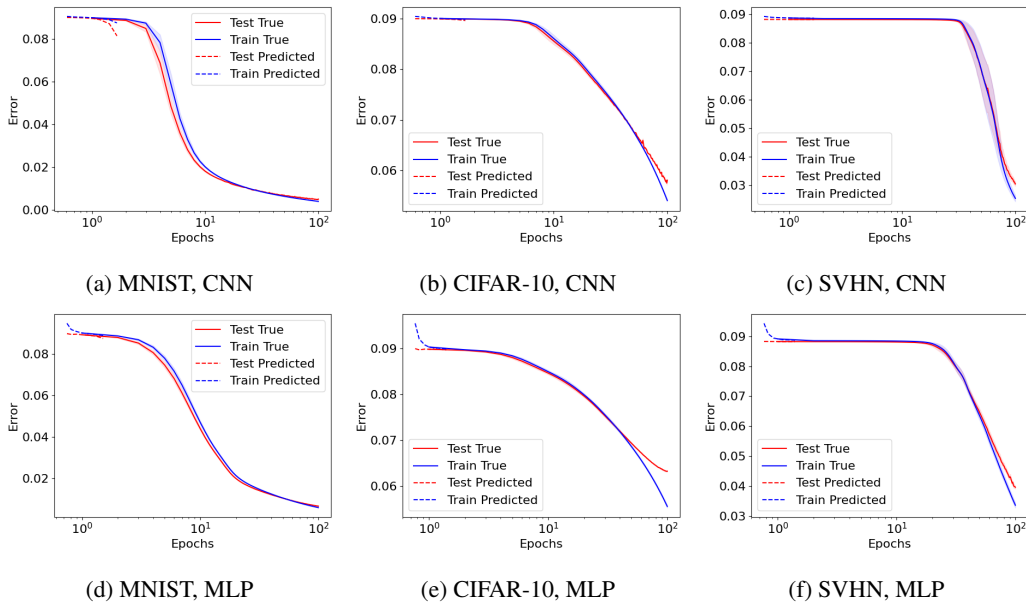


Figure 11: Predicted and true test and train mean squared error of MLP and CNN models trained on benchmark datasets under over training time for a medium-sized model. Margins indicate standard errors over 5 trials. Predictions are generated by training models of varying sizes for 1 epoch and using scale-time equivalence to predict the equivalent number of training epochs.