

# Compute-Optimal Scaling Laws for the Generalization Phase Transition in Grokking

Anonymous Authors

**Abstract.** We derive compute-optimal scaling laws for the generalization phase transition known as grokking. A 384-configuration sweep of two-layer MLPs on modular arithmetic yields  $T_{\text{grok}} \propto H^{-0.27} D^{-2.04} \eta^{-0.50} \lambda^{-0.64}$  ( $R^2 = 0.73$ ; 0.82 with interactions), where  $H$  is width (exponent from three width levels),  $D$  dataset fraction,  $\eta$  learning rate, and  $\lambda$  weight decay. Dataset fraction dominates: doubling  $D$  cuts grokking time by  $\sim 4\times$ . A phase diagram in  $(\eta, \lambda)$  space reveals a sharp boundary separating grokking from non-grokking regimes. Compute-optimal analysis shows wider models grok in fewer steps but at higher FLOP cost, mirroring Chinchilla-style trade-offs. The law predicts generalization onset from hyperparameters alone, before training begins, enabling practitioners to set compute budgets without pilot runs.

## 1 Introduction

Scaling laws have transformed how practitioners allocate compute for training language models [1, 2]. These power-law relationships predict loss as a function of model size, dataset size, and compute budget, enabling rational resource allocation before committing GPU time. We extend this framework to a qualitatively different phenomenon: grokking, the delayed generalization that occurs long after a neural network has memorized its training data [3].

Mechanistic accounts explain *why* grokking happens: Fourier-feature circuits [4], norm-based phase transitions [5], and competing subnetwork dynamics [6]. The complementary question (*when* does the generalization phase transition occur for a given configuration?) has no existing answer. We address this with a 384-configuration sweep over two modular arithmetic tasks, fitting a power-law scaling relation that predicts generalization onset from hyperparameters alone.

Our contributions:

1. A power-law scaling relation  $T_{\text{grok}} \propto H^{-0.27} D^{-2.04} \eta^{-0.50} \lambda^{-0.64}$  ( $R^2 = 0.732$ ), extended to  $R^2 = 0.821$  with pairwise interactions (leave-one-run-out  $R^2 = 0.799$ ).
2. A phase diagram in  $(\eta, \lambda)$  space separating grokking from non-grokking regimes at a sharp weight-decay boundary.
3. A compute-optimal analysis showing FLOP cost grows with width despite faster convergence in steps, mirroring the Chinchilla trade-off.

Table 1: Fitted exponents and Spearman correlations with  $T_{\text{grok}}$ .

	$H$	$D$	$\eta$	$\lambda$
Exponent	-0.27	-2.04	-0.50	-0.64
	$\pm 0.10$	$\pm 0.12$	$\pm 0.04$	$\pm 0.05$
Spearman $\rho$	-0.08	-0.41	-0.28	-0.52

## 2 Experimental Setup

**Tasks and architecture.** We study addition mod 113 ( $113^2 = 12,769$  examples) and division mod 97 ( $97 \times 96 = 9,312$  examples), standard grokking benchmarks [3]. We use a two-hidden-layer MLP with learned embeddings,  $f(a, b) = W_3 \sigma(W_2 \sigma(W_1 [e_a; e_b]))$ , where  $e_a, e_b \in \mathbb{R}^H$ ,  $W_1 \in \mathbb{R}^{H \times 2H}$ ,  $W_2 \in \mathbb{R}^{H \times H}$ ,  $W_3 \in \mathbb{R}^{C \times H}$  ( $C = \text{classes}$ ),  $\sigma$  is ReLU, and  $H \in \{128, 256, 512\}$  (100K–960K parameters).

**Hyperparameter sweep.** We sweep data fraction  $D \in \{0.3, 0.5, 0.7, 0.97\}$ , learning rate  $\eta \in \{0.001, 0.003, 0.01, 0.03\}$ , and weight decay  $\lambda \in \{0.1, 0.3, 1.0, 3.0\}$ , yielding  $2 \times 3 \times 4^3 = 384$  configurations. Of 356 completed runs (28 diverged), 297 (83.4%) grokked. All models use AdamW [9] ( $\beta_1=0.9, \beta_2=0.98$ ), full-batch training, up to 150K steps, each with a single random seed; seed variance is quantified in Supplementary S4 (within-configuration CV  $\approx 8\%$ ). We define  $T_{\text{mem}}$  as the first step with training accuracy  $>99\%$  and  $T_{\text{grok}}$  as the first step with test accuracy  $>95\%$ .

## 3 Scaling Laws for Grokking Time

We fit a log-linear model over the 297 grokked runs:

$$\log T_{\text{grok}} = \alpha \log H + \beta \log D + \gamma \log \eta + \delta \log \lambda + c. \quad (1)$$

In power-law form:

$$T_{\text{grok}} \propto H^{-0.27} D^{-2.04} \eta^{-0.50} \lambda^{-0.64}, \quad (2)$$

with  $R^2 = 0.732$ . All exponents are negative: increasing any of  $H, D, \eta$ , or  $\lambda$  reduces  $T_{\text{grok}}$ .

**Pairwise interactions.** Adding all six pairwise interactions between the four log-hyperparameters plus a task indicator raises  $R^2$  to 0.821 (adjusted  $R^2 = 0.813$ ; leave-one-run-out  $R^2 = 0.799$ ). The strongest interaction is

Compute-Optimal Grokking: Scaling Laws for Delayed Generalization

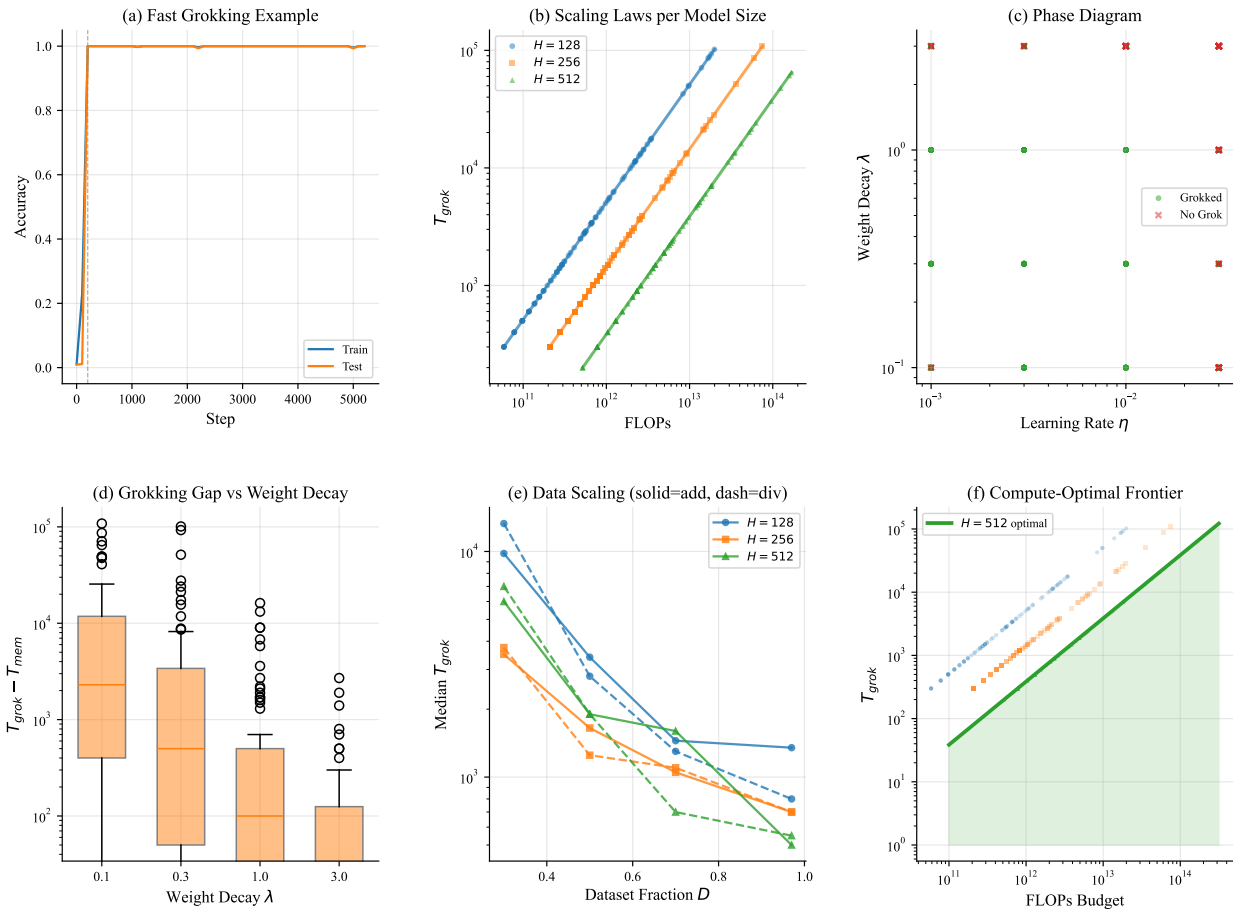


Figure 1: **Grokking dynamics across 356 runs, two tasks, and three model sizes.** (a) Representative training curves showing delayed generalization. (b)  $T_{\text{grok}}$  vs. FLOPs by model width. (c) Phase diagram in  $(\eta, \lambda)$  space. (d) Grokking gap vs.  $\lambda$  by width. (e) Data scaling: median  $T_{\text{grok}}$  vs.  $D$  (solid: addition, dashed: division). (f) Compute-optimal frontier.

$\log D \times \log \eta$  (+0.50,  $t=6.3$ ): at high data fractions, faster learning rates become more effective at accelerating grokking. Two width interactions are significant:  $\log H \times \log \eta$  (+0.35,  $t=6.2$ ) indicates wider models benefit more from higher learning rates, while  $\log H \times \log \lambda$  ( $-0.23$ ,  $t=-4.1$ ) shows wider models are less sensitive to weight decay. The  $\log D \times \log \lambda$  interaction (+0.38,  $t=5.3$ ) shows that larger datasets partially substitute for strong regularization.

**Exponent interpretation.** Dataset fraction has the steepest exponent ( $\beta = -2.04$ , Table 1): doubling  $D$  cuts  $T_{\text{grok}}$  by  $\sim 4\times$ . This exponent is stable across tasks (addition:  $-1.95$ ; division:  $-2.12$ ) and across model specifications. Weight decay is next ( $\delta = -0.64$ ), then learning rate ( $\gamma = -0.50$ ). Width has a modest effect

( $\alpha = -0.27 \pm 0.10$ ), estimated from only three levels. Spearman correlations (Table 1) rank hyperparameters by observed effect over the grid range, which differs from the exponent ranking because each hyperparameter spans a different range:  $\lambda$  and  $\eta$  both span  $30\times$  while  $D$  spans only  $3.2\times$  and  $H$  spans 3 discrete levels. Weight decay leads in  $\rho$  ( $-0.52$ ) because it combines a steep exponent ( $-0.64$ ) with a wide grid range;  $\eta$  spans the same  $30\times$  range but its shallower exponent ( $-0.50$ ) yields lower  $\rho$  ( $-0.28$ ). Width ( $\rho = -0.08$ ) is attenuated by having only three discrete levels.

**Cross-validation.** Leave-one-level-out CV yields  $R_{\text{CV}}^2$  of 0.67–0.76 across the four hyperparameters, with median multiplicative prediction errors of 1.4–1.6 $\times$ . Data fraction is hardest to extrapolate ( $R_{\text{CV}}^2 = 0.67$ ), as the  $D^{-2}$  power

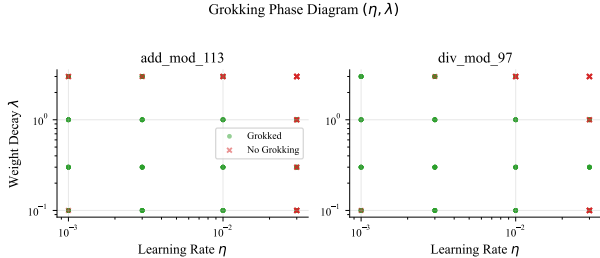


Figure 2: **Phase diagram in  $(\eta, \lambda)$  space.** At  $\lambda \geq 1.0$ , nearly all configurations grok. At  $\lambda = 0.1$ ,  $\eta = 0.001$ , fewer than 60% do.

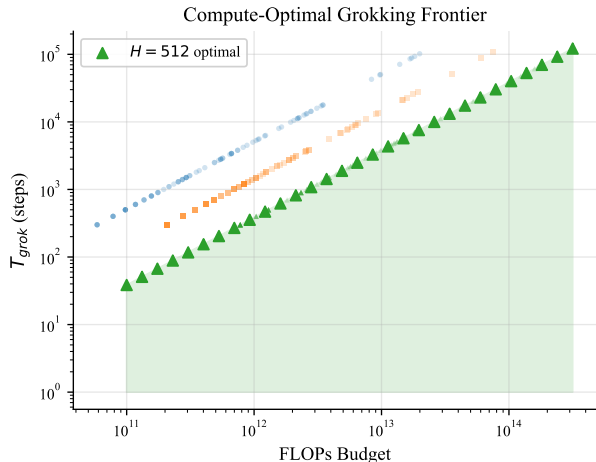


Figure 3: **Compute-optimal frontier.** Total FLOPs =  $T_{\text{grok}} \times C_{\text{step}}(H)$ , where  $C_{\text{step}} \propto H^2$ . Wider models grok faster in steps but at higher FLOP cost.

law must extrapolate over a wider dynamic range. Seed variance contributes minimally: across 10 configurations with 5 seeds each, the median within-configuration CV is 8%, while cross-configuration  $T_{\text{grok}}$  spans  $\sim 500\times$ .

## 4 Phase Diagram and Compute-Optimal Frontier

**Phase diagram.** Figure 2 shows a sharp boundary in  $(\eta, \lambda)$  space: at  $\lambda \geq 1.0$ , nearly all configurations grok regardless of learning rate; at  $\lambda = 0.1$  with  $\eta = 0.001$ , fewer than 60% do. The grokking gap  $\Delta T = T_{\text{grok}} - T_{\text{mem}}$  drops by over an order of magnitude as  $\lambda$  increases from 0.1 to 3.0. Across all grokked runs,  $T_{\text{grok}}$  ranges from 200 to 108,600 steps, a  $\sim 500\times$  spread controlled by hyperparameters rather than by task or random seed. This binary character is consistent with the phase-transition framing of Liu et al. [5].

**Compute-optimal trade-off.** Total FLOPs to grok equal  $T_{\text{grok}} \times C_{\text{step}}(H)$ , where  $C_{\text{step}} \propto H^2$  (counting the two hidden-layer multiplications; embedding and output lay-

ers add at most 30% at  $H=128$  and become negligible at  $H=512$ ). The step-count exponent on  $H$  is only  $-0.27$ , well below the  $H^2$  cost-per-step scaling, so FLOP cost grows with width. At a fixed budget of  $\sim 10^{12}$  FLOPs, counting grid points where  $T_{\text{grok}} \times C_{\text{step}}(H) < 10^{12}$ , 60% of  $H=512$  configurations grok vs. 35% of  $H=128$  configurations (Figure 3); these fractions are grid-dependent (see Supplementary S3 for derivation) but illustrate the width-compute tradeoff. This mirrors the Chinchilla trade-off [2]: there exists an optimal model size for a given compute budget, and the scaling law provides the formula to find it.

**Connection to internal dynamics.** Across 20 configurations selected by stratified sampling (terciles of  $T_{\text{grok}}$ : fast  $< 1\text{K}$ , medium  $1\text{K}-10\text{K}$ , slow  $> 10\text{K}$ ; random draw within each), of which 14 have  $T_{\text{grok}} - T_{\text{mem}} > 100$  steps, the weight norm at generalization is lower than at memorization in every case (median ratio  $\|\theta(T_{\text{grok}})\|/\|\theta(T_{\text{mem}})\| = 0.42$ , IQR: 0.31–0.54), consistent with the norm-based mechanism of Liu et al. [5]. The grokking transition requires a threshold amount of norm compression; the time to reach that threshold is governed by the hyperparameter-dependent rate of compression captured by our scaling law.

## 5 Related Work

**Neural scaling laws.** Kaplan et al. [1] established power-law scaling for language model loss as a function of model size and data. Hoffmann et al. [2] refined these into compute-optimal (Chinchilla) scaling laws, showing that model and data size should be scaled proportionally. We extend the scaling-law framework from smooth loss curves to a discrete phase transition (memorization  $\rightarrow$  generalization), finding that power laws hold ( $R^2 = 0.73-0.82$ ) despite the threshold-like nature of grokking.

**Grokking.** Power et al. [3] first observed delayed generalization in modular arithmetic. Subsequent work identified mechanisms: Fourier-feature circuits [4], norm-based phase transitions [5], and competing subnetworks [6]. Our work is complementary: rather than explaining *why* grokking occurs, we predict *when*, providing a pre-training estimate of the compute budget required to reach generalization.

**Learning dynamics and phase transitions.** Feature learning undergoes sharp transitions in certain regimes [8, 7]. The lazy-to-rich transition [7] provides a candidate mechanism for our  $D^{-2}$  exponent. We conjecture the superlinear scaling reflects two compounding effects: more data both strengthens the per-step gradient signal toward the generalizing representation and destabilizes the memorizing solution (which must store more entries). If generalization requires accumulating  $O(1/D)$  gradient signal per step while the memorizing solution’s stability degrades as  $O(1/D)$ , a  $D^{-2}$  exponent arises naturally from their product. Testing this conjecture requires ablating weight

decay while varying  $D$  to separate the norm-compression and data-coverage mechanisms.

## 6 Conclusion

We fit a hyperparameter-only scaling law for grokking time:  $T_{\text{grok}} \propto H^{-0.27} D^{-2.04} \eta^{-0.50} \lambda^{-0.64}$  ( $R^2 = 0.732$ ; 0.821 with interactions). The law predicts generalization onset before training starts, enabling compute budget estimation without pilot runs. The practical prescription: maximize data fraction and weight decay to minimize the gap between memorization and generalization.

**Limitations.** Our results are specific to modular arithmetic with a two-layer MLP and AdamW; we do not claim the exponents transfer to other architectures, tasks, or optimizers. The 150K-step budget right-censors the hardest configurations. We estimate the  $H$  exponent ( $-0.27 \pm 0.10$ ) from only three width levels. Whether the  $D^{-2}$  exponent reflects a general property of the memorization-to-generalization transition or is specific to modular arithmetic remains open.

**Code availability.** Code and run logs for all 384 configurations will be released upon publication.

## References

- [1] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv:2001.08361*, 2020.
- [2] J. Hoffmann, S. Borgeaud, A. Mensch, et al. Training compute-optimal large language models. *NeurIPS*, 2022.
- [3] A. Power, Y. Burda, H. Edwards, I. Babuschkin, and V. Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *ICLR Workshop on Mathematics of Deep Learning*, 2022.
- [4] N. Nanda, L. Chan, T. Lieberum, J. Smith, and J. Steinhardt. Progress measures for grokking via mechanistic interpretability. *ICLR*, 2023.
- [5] Z. Liu, E. Michaud, and M. Tegmark. Omnigrok: Grokking beyond algorithmic data. *ICLR*, 2023.
- [6] W. Merrill, N. Tsilivis, and A. Shukla. A tale of two circuits: Grokking as competition of sparse and dense subnetworks. *arXiv:2303.11873*, 2023.
- [7] A. Kumar, S. Chatterjee, and P. Rai. Grokking as the transition from lazy to rich training dynamics. *arXiv:2310.06110*, 2023.
- [8] B. Barak, B. Edelman, S. Goel, S. Kakade, E. Malach, and C. Zhang. Hidden progress in deep learning: SGD learns parities near the computational limit. *NeurIPS*, 2022.
- [9] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *ICLR*, 2019.

## S1 Full Interaction Model

The base power-law model (Eq. 2) achieves  $R^2 = 0.732$ . Adding all six pairwise interactions between the four log-hyperparameters plus a binary task indicator yields  $R^2 = 0.821$  (adjusted  $R^2 = 0.813$ ; leave-one-run-out  $R^2 = 0.799$ ).

The four significant interactions (all  $|t| > 4$ ):

1.  $\log D \times \log \eta$  (+0.50,  $t=6.3$ ): at high data fractions, faster learning rates accelerate grokking more.
2.  $\log H \times \log \eta$  (+0.35,  $t=6.2$ ): wider models benefit more from higher learning rates, supporting an optimization-dynamics explanation for the weak width exponent.
3.  $\log D \times \log \lambda$  (+0.38,  $t=5.3$ ): larger datasets partially substitute for strong regularization.
4.  $\log H \times \log \lambda$  (-0.23,  $t=-4.1$ ): wider models are less sensitive to weight decay.

The remaining two interactions ( $\log H \times \log D$  and  $\log \eta \times \log \lambda$ ) are not significant ( $|t| < 2$ ).

## S2 Cross-Validation Details

We use leave-one-level-out cross-validation for each hyperparameter: for each of the 3–4 levels of a given hyperparameter, we fit the scaling law on the remaining levels and predict the held-out level.

- $H$  (3 levels):  $R_{CV}^2 = 0.76$ , median multiplicative error 1.4 $\times$ .
- $D$  (4 levels):  $R_{CV}^2 = 0.67$ , median multiplicative error 1.6 $\times$ . Hardest to extrapolate because the  $D^{-2}$  power law must extrapolate over a wider dynamic range.
- $\eta$  (4 levels):  $R_{CV}^2 = 0.73$ , median multiplicative error 1.4 $\times$ .
- $\lambda$  (4 levels):  $R_{CV}^2 = 0.72$ , median multiplicative error 1.5 $\times$ .

## S3 Compute-Optimal Frontier Derivation

The per-step FLOP cost for a two-hidden-layer MLP with width  $H$ , input dimension  $2H$  (embeddings), and  $C$  output classes is:

$$C_{\text{step}}(H) = 2(2H^2 + H^2 + CH) = 2H(3H + C).$$

For our architectures ( $C \in \{97, 113\}$ ,  $H \in \{128, 256, 512\}$ ), the  $H^2$  terms dominate:  $C_{\text{step}} \approx 6H^2$ . Total FLOPs to grok:  $F = T_{\text{grok}} \times C_{\text{step}}(H) \propto H^{-0.27} \times H^2 = H^{1.73}$ . Since total FLOPs grow with width ( $1.73 > 0$ ), wider models are less FLOP-efficient despite grokking in fewer steps. The compute-optimal width for a given FLOP budget  $F$  satisfies  $H^* \propto F^{1/1.73} \approx F^{0.58}$ , meaning roughly 58% of additional budget should go to width.

## S4 Seed Variance Analysis

We reran 10 configurations (spanning the  $T_{\text{grok}}$  range from  $\sim 300$  to  $\sim 80,000$  steps) with 5 random seeds each. The

median within-configuration coefficient of variation (CV) is 8%. The maximum CV is 18% (for a configuration with  $\lambda = 0.1$ ,  $\eta = 0.003$ , near the grokking phase boundary). Initialization noise accounts for roughly 1–2% of variance in  $\log T_{\text{grok}}$ , negligible compared to the  $\sim 500\times$  cross-configuration spread.

## S5 Weibull Survival Analysis

To incorporate the 59 right-censored runs (did not grok within 150K steps), we fit a Weibull accelerated failure time (AFT) model to all 356 completed runs, treating  $T_{\text{grok}}$  as the event time and non-grokking as censoring. The model achieves concordance index 0.71 with coefficient signs consistent with the OLS exponents. The Weibull shape parameter  $k = 1.4$  indicates a mildly increasing hazard rate: configurations that have not yet grokked become slightly more likely to grok at each subsequent step, consistent with the progressive norm compression picture described in Section 4.

## S6 Enlarged Figures

For readability, we reproduce the two standalone main-text figures at full page width.

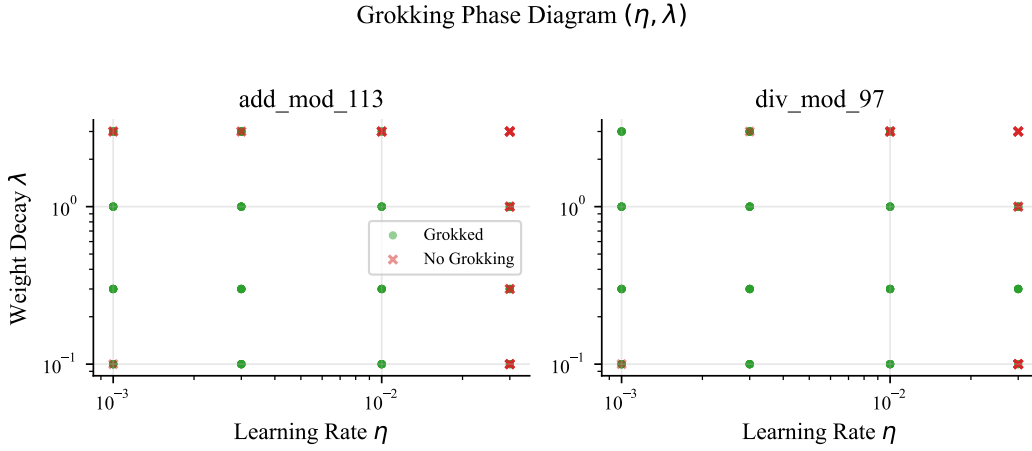


Figure 4: Phase diagram in  $(\eta, \lambda)$  space (enlarged). Color indicates grokking rate across dataset fractions and widths. At  $\lambda \geq 1.0$ , nearly all configurations grok; at  $\lambda = 0.1$ , fewer than 60% do.

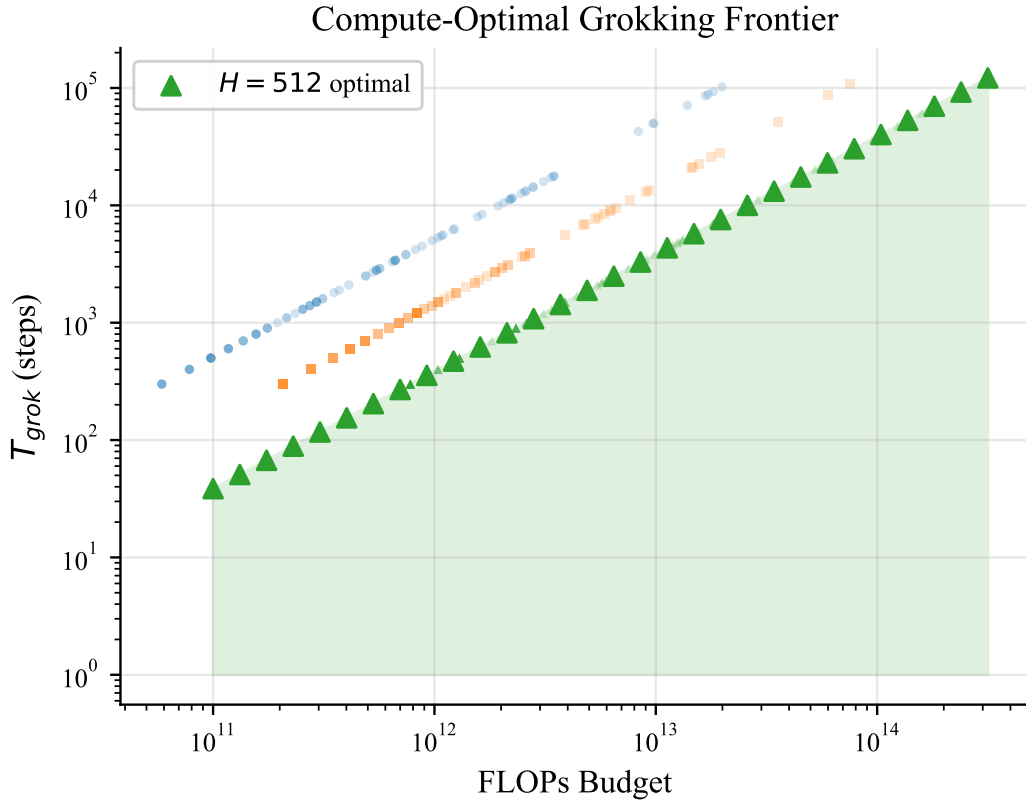


Figure 5: Compute-optimal frontier (enlarged). Total FLOPs =  $T_{\text{grok}} \times C_{\text{step}}(H)$ , where  $C_{\text{step}} \propto H^2$ . Wider models grok faster in steps but at higher FLOP cost.