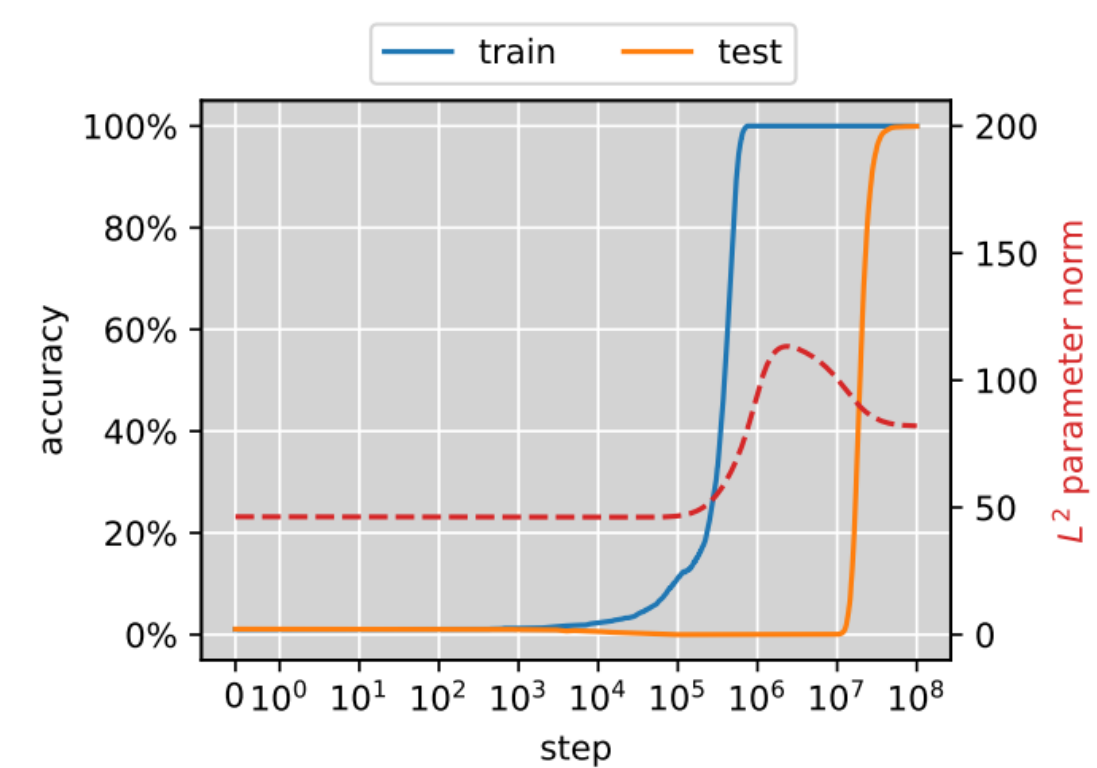


Dichotomy of Early and Late Phase Implicit Biases Can Provably Induce Grokking

Kaifeng Lyu*, Jikai Jin*, Zhiyuan Li, Simon S. Du, Jason D. Lee, Wei Hu

¹Princeton University ²Stanford University ³Toyota Technological Institute at Chicago ⁴University of Washington ⁵University of Michigan

What is Grokking?



For some algorithmic datasets...

- **Example (Modular Addition):** Given one-hot representations of a, b , output $(a + b) \bmod p$
- **Other Examples:** learning group operations, sparse parity, GCD...

training neural nets leads to grokking (Power et al., 2022):

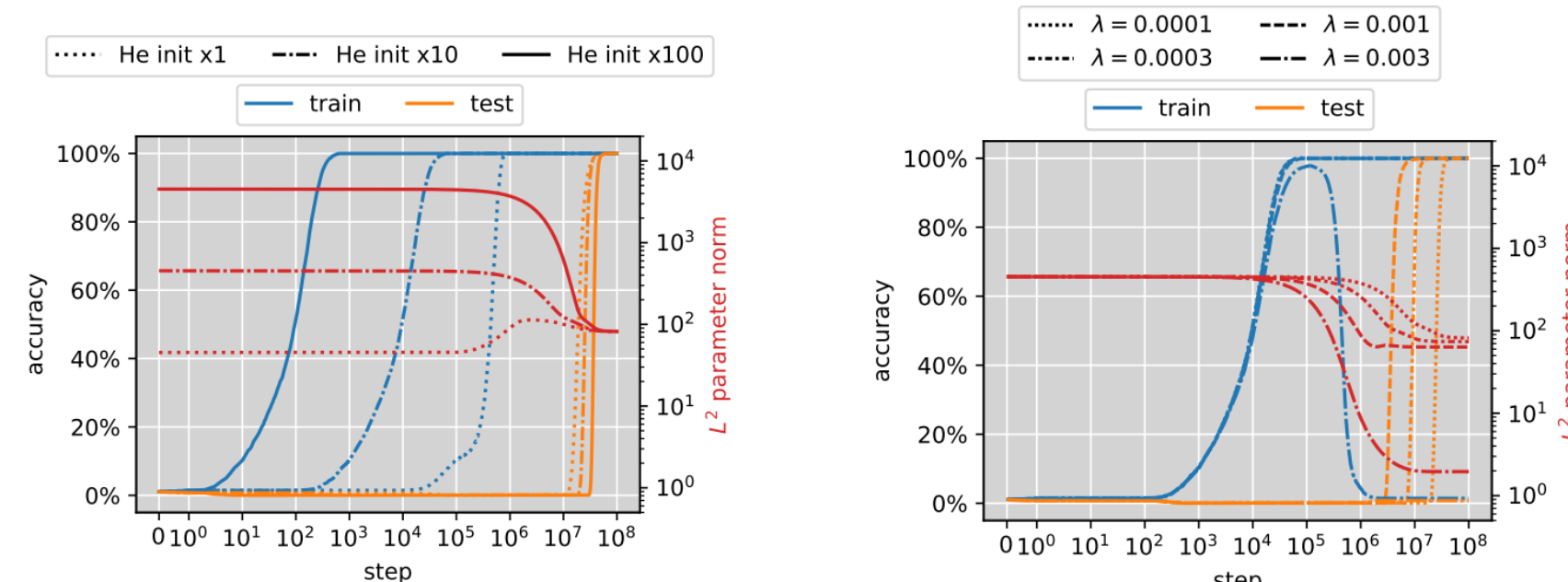
- **Phase 1:** train acc = 100% but test acc is very low
- But after training for sufficiently longer ...
- **Phase 2:** test acc “suddenly” jumps to 100%

Many recent works aim to explain grokking: Liu et al. (2023); Varma et al. (2023); Thilak et al. (2022); Notsawo Jr et al. (2023); Nanda et al. (2023); Chughtai et al. (2023); Gromov (2023); Žunkovič & Ilievski (2022); Levi et al. (2024).

But ...

1. No theoretical analysis for neural nets
2. No quantitative explanation for why the transition is sharp in grokking

Ablation Studies: Init Scale and WD



Grokking time is delayed with larger init

Grokking time is delayed with smaller WD

“large init & small WD makes grokking more significant”

Prior work (Liu et al., 2023): With large init and small WD, grokking may happen even on natural tasks (e.g., image/sentiment classification)

Main Results

Our Goal: Give some examples of grokking with precise theoretical characterization.

Classification with L2 Regularization (Weight Decay):

$$\mathcal{L}_\lambda(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\theta; x_i)) + \frac{\lambda}{2} \|\theta\|_2^2$$

We use the exponential loss $\ell(q) = e^{-q}$ for simplicity

Gradient Flow:

$$\frac{d\theta}{dt} = -\nabla \mathcal{L}_\lambda(\theta)$$

(= gradient descent with infinitesimal learning rate)

Large Initialization, Small Weight Decay

“large init & small WD makes grokking more significant”
 \Rightarrow Take $\theta(0) = \alpha \bar{\theta}$, $\alpha \rightarrow +\infty$, $\lambda \rightarrow 0$.

Asymptotics: $\lambda = \Theta(\alpha^{-p})$ for some positive $p = \Theta(1)$.

Assumption 1: Homogeneous Nets (Lyu & Li, 2020; Ji & Telgarsky, 2020, ...)

The model is L -homogeneous wrt θ :

$$f(c\theta; x) = c^L f(\theta; x) \text{ for all } c > 0$$

Assumption 2: Smoothness

$f(\theta; x)$ is \mathcal{C}^2 -smooth wrt θ .

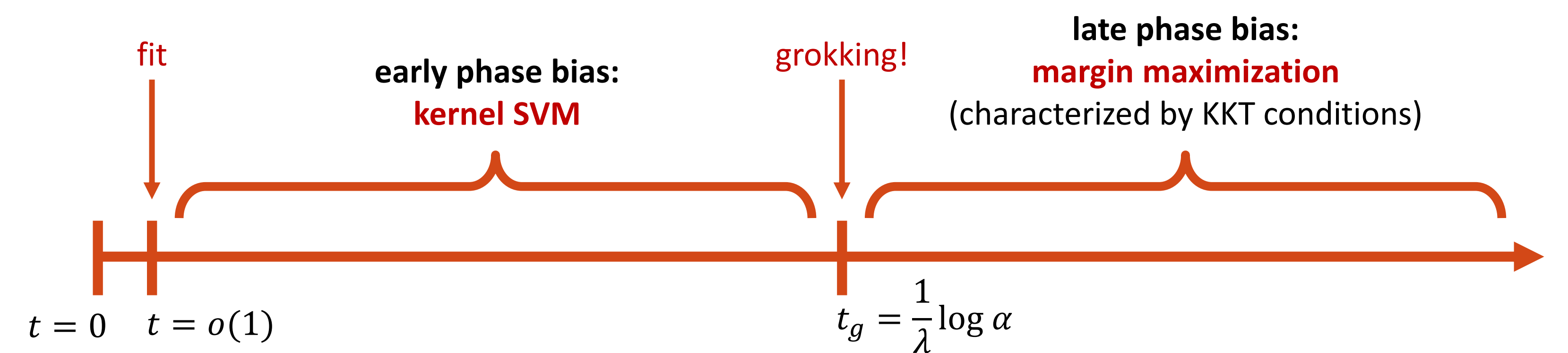
(just for simplicity; the proof should be extendable to non-smooth cases)

Assumption 3: Zero initial output

$\bar{\theta}$ satisfies $f(\bar{\theta}; x) = 0$ for all x .

Justification: A common assumption for studying NTK. Make the proof much simpler (Chizat et al., 2019; Hu et al., 2020).

- can be done by symmetrized init or “difference trick”
- NTK init + width $\rightarrow \infty \Rightarrow$ approximately true. Proof is usually extendable to this case.



Theorem 1. For $c \in (0,1)$, at $t = (1-c)t_g$, the model outputs the same as the kernel SVM with $K(x, x') = \langle \nabla f(\bar{\theta}; x), \nabla f(\bar{\theta}; x') \rangle$ when $\alpha \rightarrow +\infty$.

$$\min \|h\|_2^2 \quad \text{s.t.} \quad y_i \langle \nabla f(\bar{\theta}; x_i), h \rangle \geq 1$$

Theorem 2. For $c > 0$, at $t = (1+c)t_g$, the model attains first-order optimal conditions (KKT conditions) for the following margin maximization problem:

$$\min \|\theta\|_2^2 \quad \text{s.t.} \quad y_i f(\theta; x_i) \geq 1$$

Key Insights

Kernel regime:

- Unlike the usual kernel regime, where $\theta(t) \approx \theta(0)$, here the parameter itself changes a lot but direction does not, i.e., $\frac{\theta(t)}{\|\theta(t)\|_2} \approx \frac{\theta(0)}{\|\theta(0)\|_2}$.
- Analyze the change in direction and norm very carefully.
- Some calculation \Rightarrow As long as the norm is not $o\left(\left(\log \frac{1}{\lambda}\right)^{1/L}\right)$, θ is in the kernel regime.

Rich regime:

- Once the norm decays to this level, only $o\left(\frac{1}{\lambda} \log \frac{1}{\lambda}\right)$ time is needed to reach the KKT.
- Turns out to be very short compared to the time to decay the norm!

Example: Two-layer Diagonal Nets

Two-layer Diagonal Net: A reparameterization of linear model (Woodworth et al., 2020)

$$f(\theta; x) = \langle u \odot u - v \odot v, x \rangle, \quad \theta = (u, v).$$

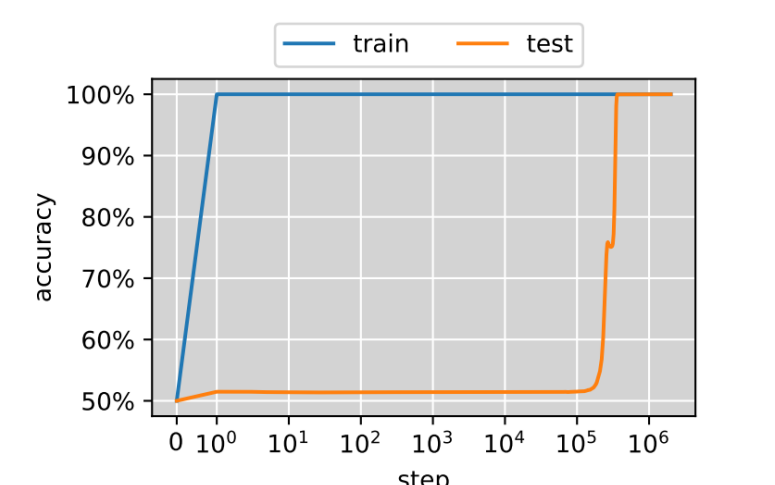
Alternatively: $w = u \odot u - v \odot v$, $f(\theta; x) = \langle w, x \rangle$.

Initialization: $u = v = (1, 1, \dots, 1)$.

Kernel SVM = L2 max-margin linear classifier

Max-margin solution = L1 max-margin linear classifier (encouraging sparsity)

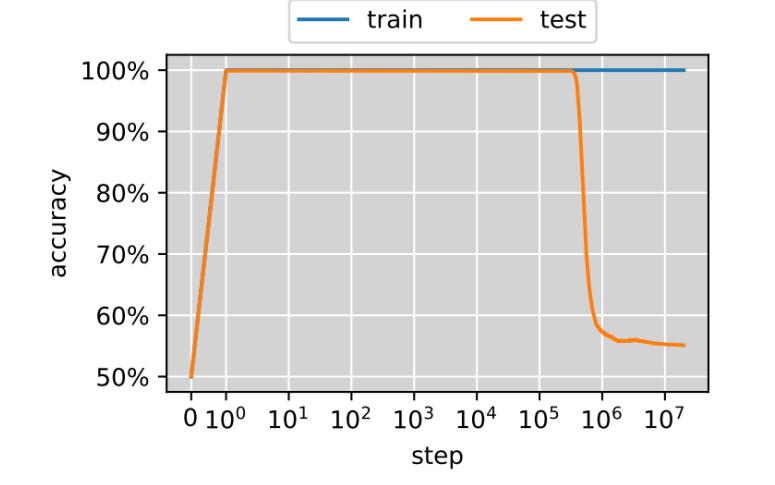
Sparse Linear Regression
 \Rightarrow Grokking



$$x \sim U(\{\pm 1\}^d)$$

$$y = \text{sgn}(x_1 + x_2 + x_3)$$

Labels are generated by a linear classifier with large L2 margin
 \Rightarrow “Misgrokking”



$$z \sim \mathcal{N}(0, I)$$

$$y = \text{sgn}(\langle z, w^* \rangle), x = z + \frac{Y}{2} y w^*$$

Example: Completing Multiplication Tables

(Our results can be extended to regression settings; see our paper)

Overparameterized Matrix Completion:

Parameterize $W = UU^T - VV^T$, $U, V \in \mathbb{R}^{d \times d}$.

Use MSE loss on observed entries.

Related to learning two-layer nets with quadratic activation.

