



## Key Contributions

- We formalize pattern matching as functional-equivalence-based generalization with a precise, testable boundary.
- We prove **tight data scaling laws**: 2-hop training data for complete ID generalization grows as  $N_{\text{req}} = \tilde{\Theta}(n^c)$  with exponent  $2.0 < c \leq 2.5$ , which robustly holds across 20x model scaling and architectures.
- We identify **path ambiguity as a structural barrier**: when a variable influences output via multiple paths, models fail to form unified representations, even with CoT.

## Our Formalization of Pattern Matching

### Key Insight: Functional Equivalence

If two input fragments consistently produce the same output in all shared contexts observed during training, a learner can substitute one for the other to predict unseen combinations.

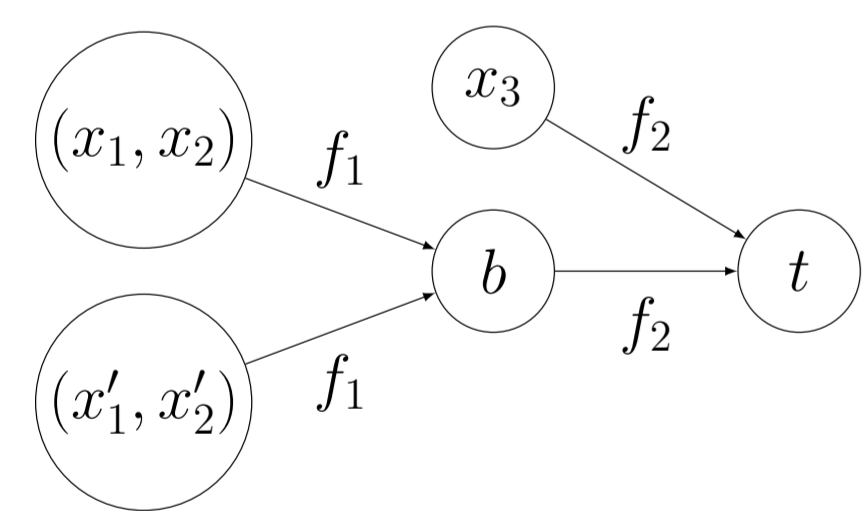
#### Definition: Functional $k$ -Equivalence

Two input subsequences  $a, a'$  at positions  $I$  are **functionally  $k$ -equivalent in  $D$**  if there exist  $\geq k$  contexts  $b$  (at the remaining positions  $I^c$ ) such that both  $(a, b)$  and  $(a', b)$  appear in  $D$  and **always yield the same output**:  $f(a, b) = f(a', b)$ .

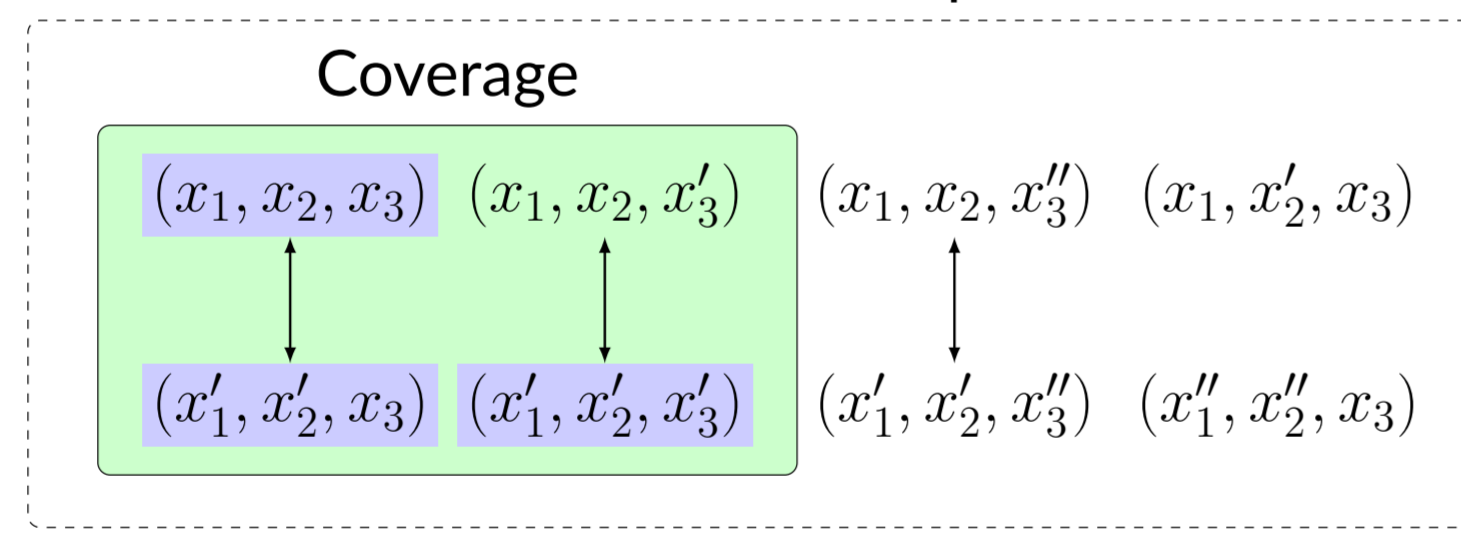
#### Definition: $k$ -Coverage ( $\leftarrow$ Boundary of Pattern Matching)

The  **$k$ -coverage** of  $D$  is the set of all inputs reachable from observed data through chains of functionally  $k$ -equivalent substitutions in the **substitution graph**.

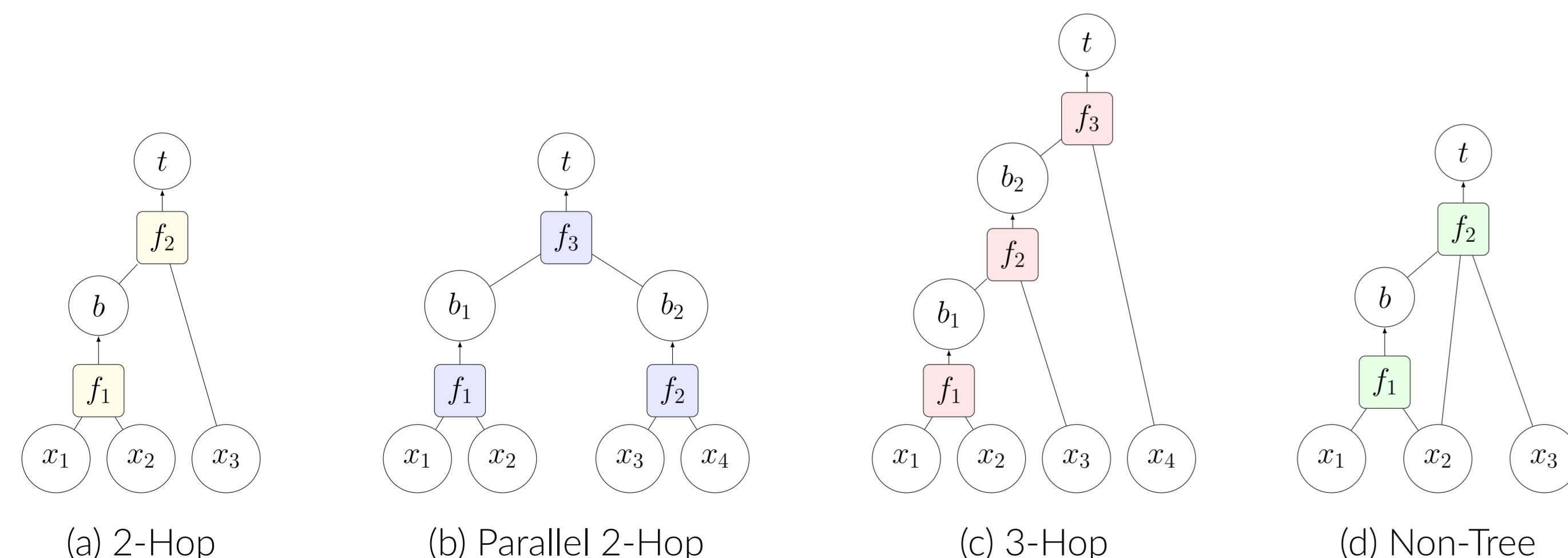
#### Functional Equivalence Observation



#### Substitution Graph



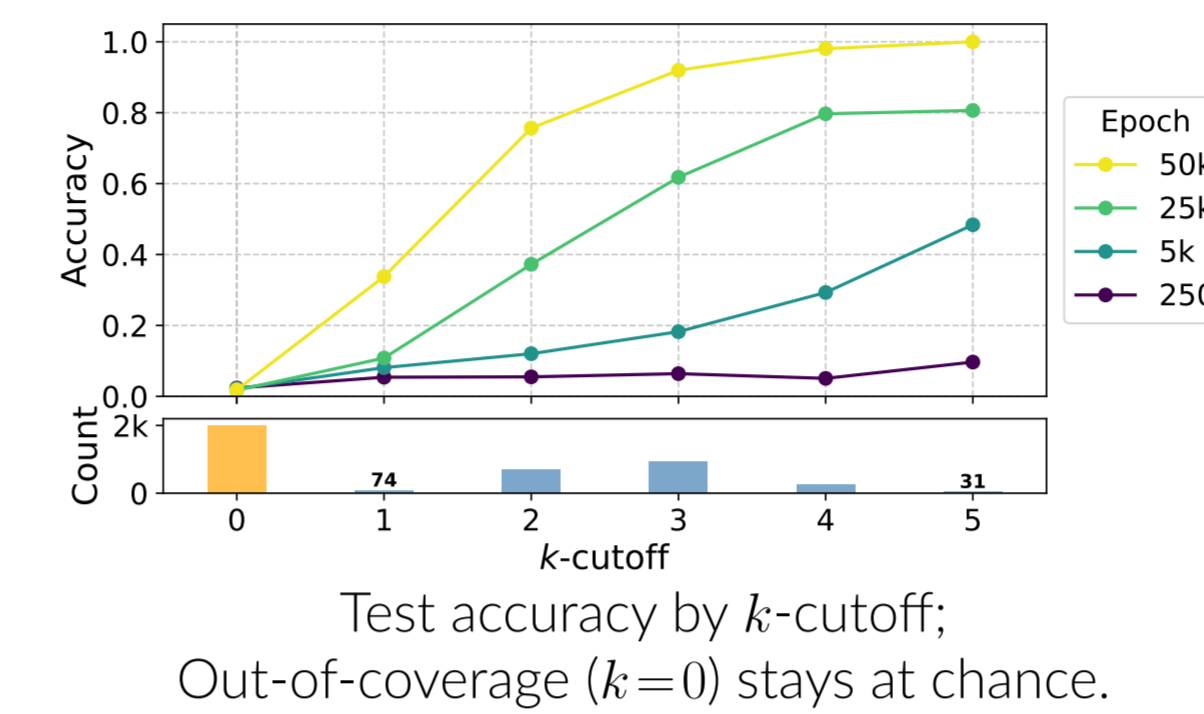
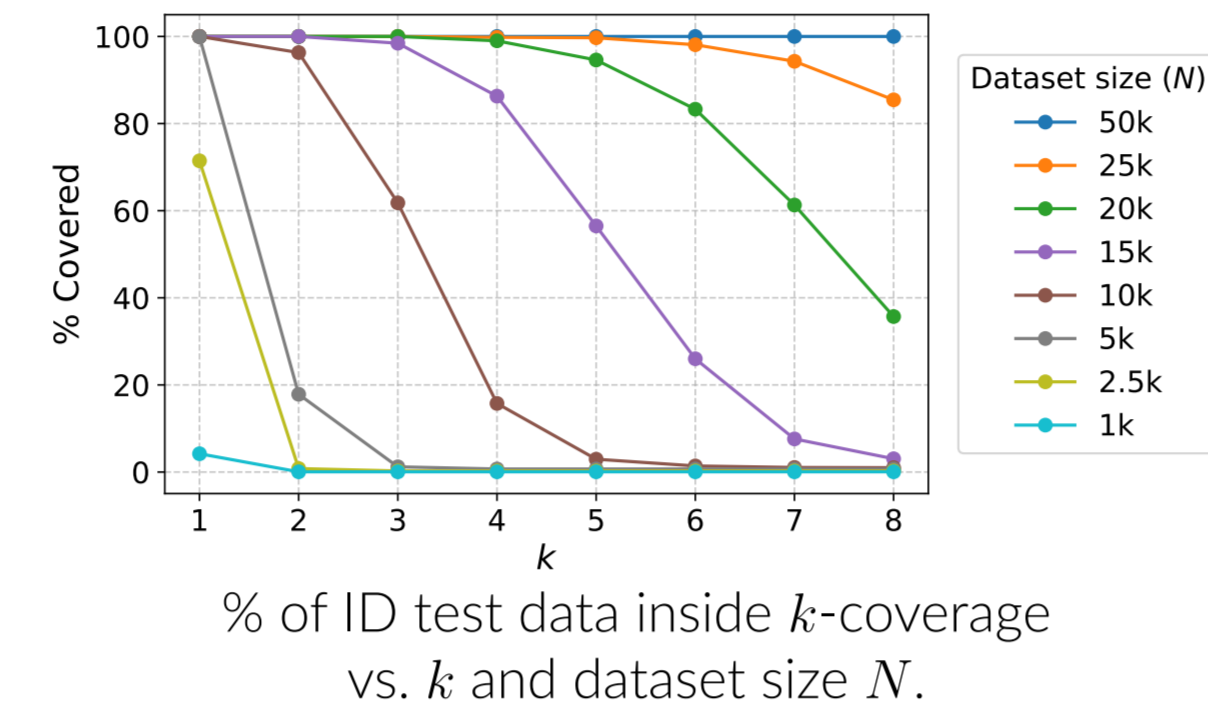
## Experimental Setup: Four Task Structures



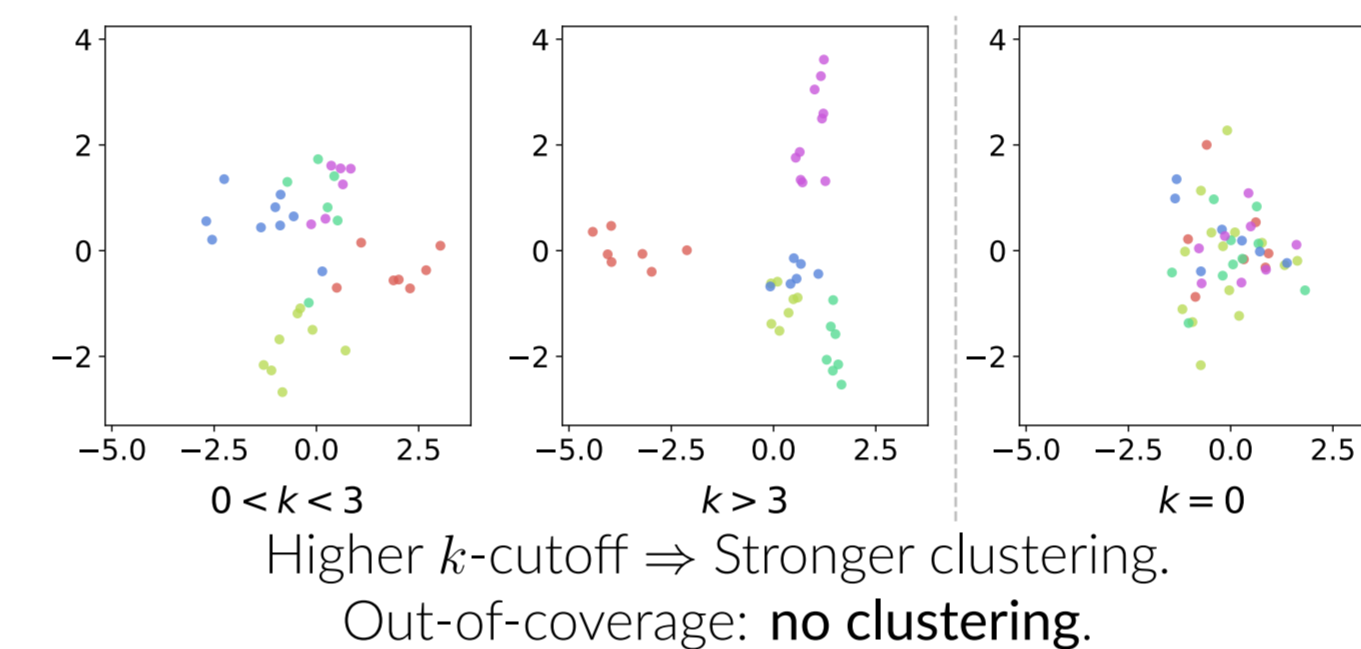
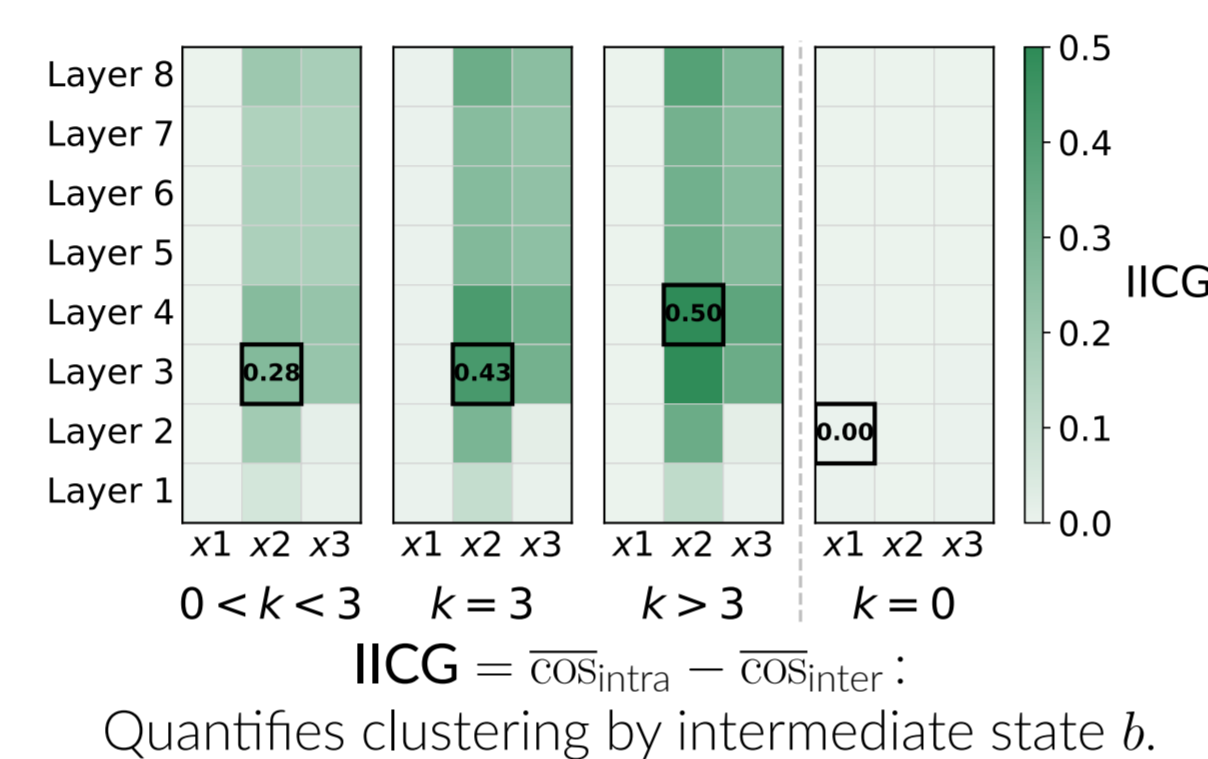
- Primitive functions  $f_1, f_2, \dots$  are **randomly generated** to isolate pattern matching from algebraic shortcuts
- Train GPT-2 (68M–1.5B) and Mamba from scratch;  $|\mathcal{X}| \in \{50, \dots, 200\}$

## Evidence Strength Predicts Pattern-Matching Success

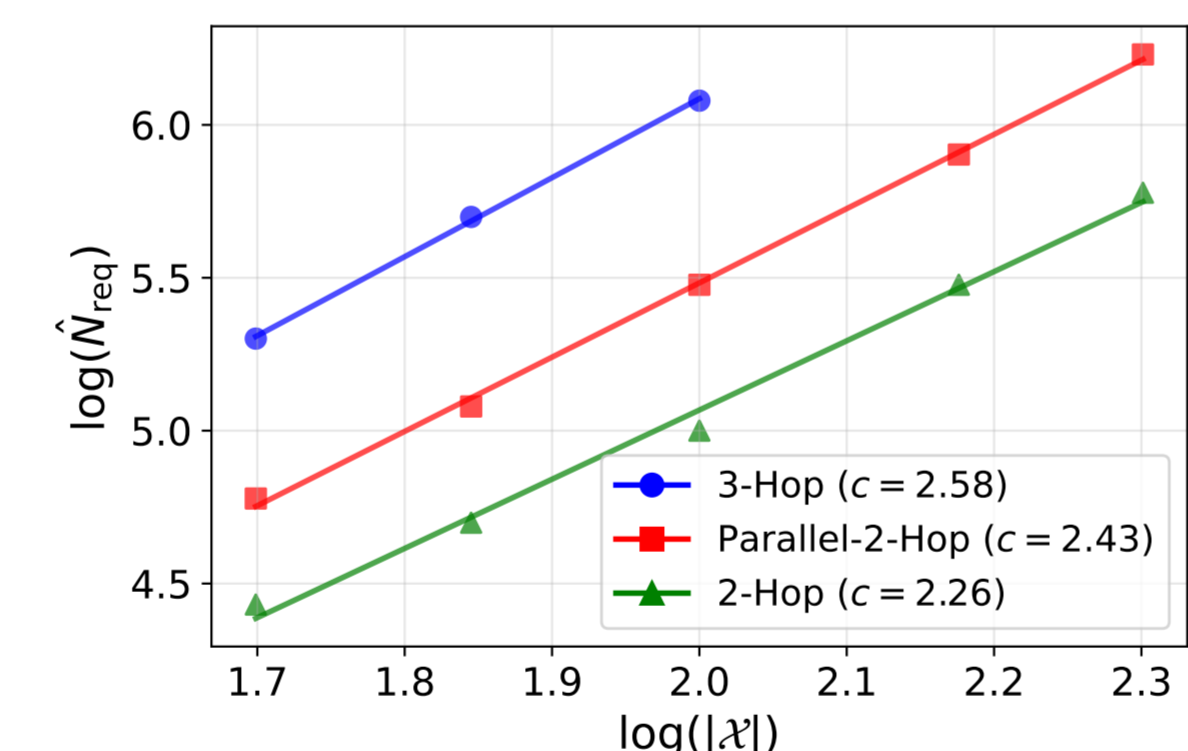
Stronger evidence  $\Rightarrow$  faster & more reliable generalization!



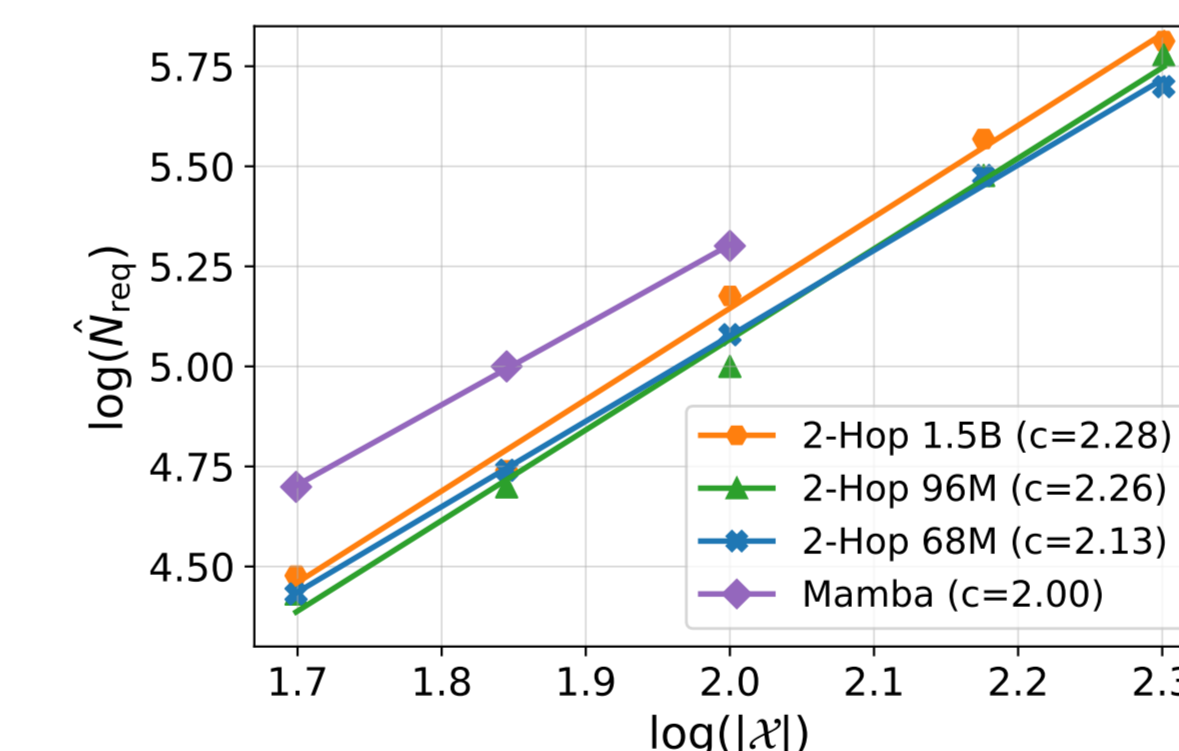
## Representation Clusters Drive Pattern Matching



## Data Scaling Law: Empirical Confirmation

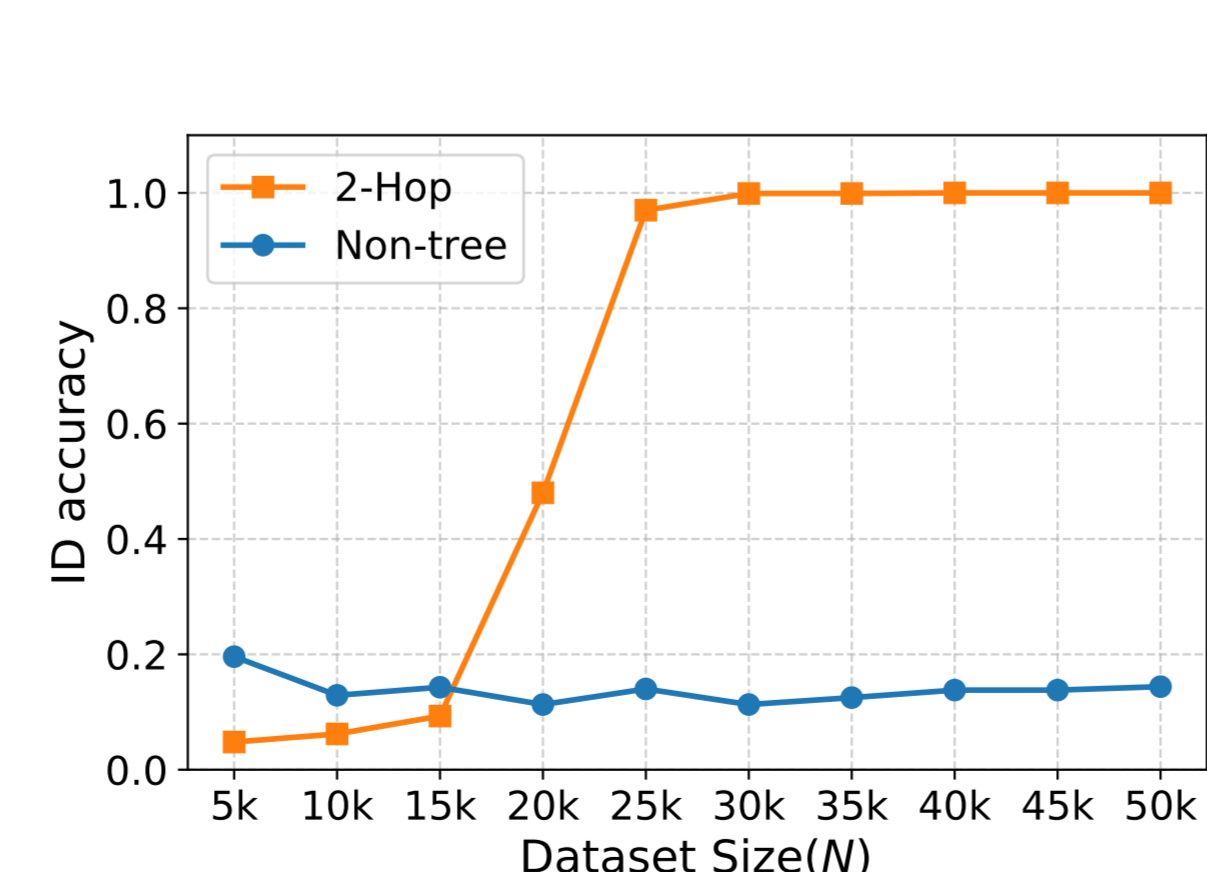


Our theory explains the empirical 2-Hop exponent  $c=2.26$ ; Deeper structures scale steeper.

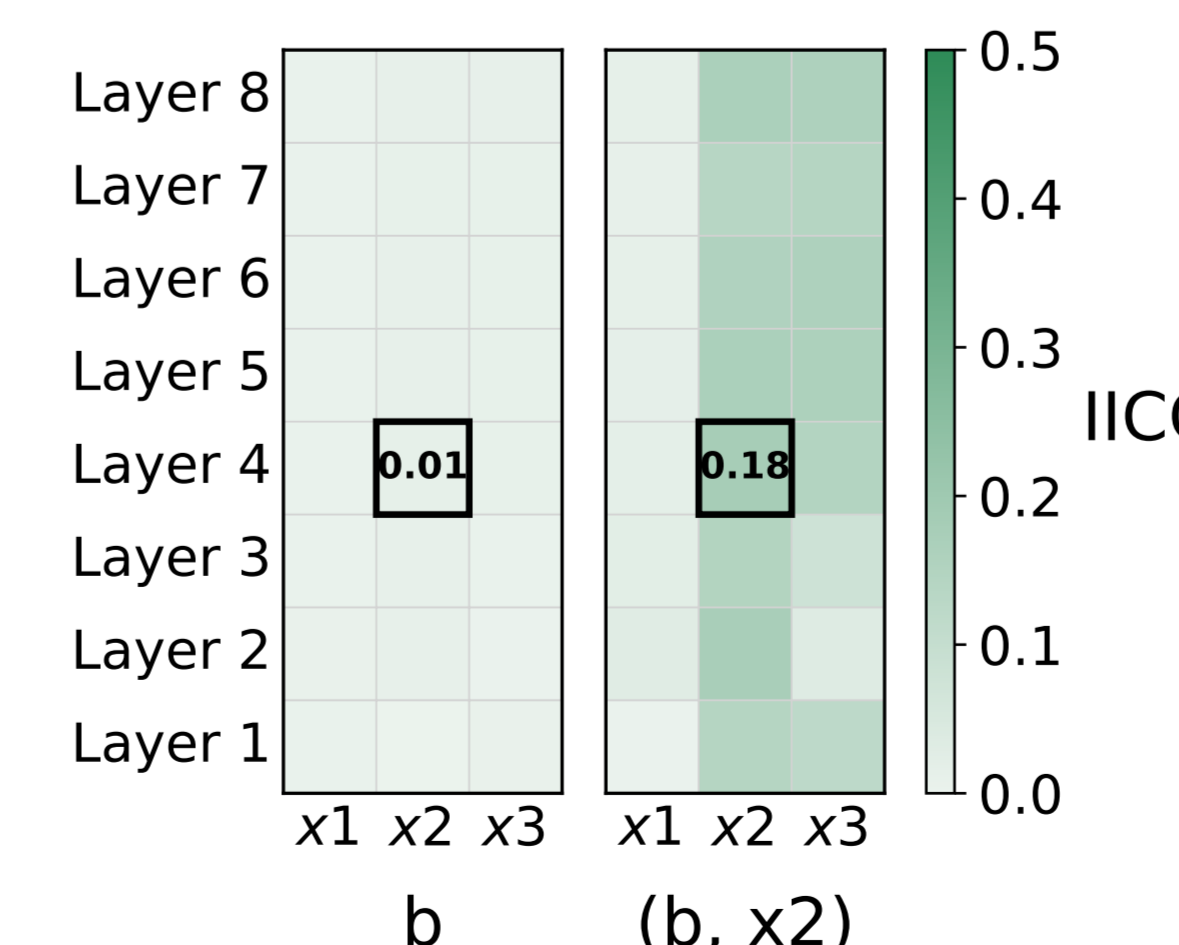


Exponents **invariant across 20x model scaling** and architectures  $\Rightarrow$  a **data property**, not capacity!

## Path Ambiguity in 'Non-Tree' Task & Chain-of-Thought



2-Hop achieves full ID generalization; Non-Tree fails even with near-exhaustive data.



Model forms **context-dependent** representations (i.e., high IICG for  $(b, x_2)$ , near-zero for  $b$ ).

CoT "flattens" multi-hop (improving scaling laws) but **does not fully resolve** path ambiguity.

## Theorem 6.1 (Tight Sample Complexity for 2-Hop Task)

For a 2-Hop task with token set size  $n$  and an i.i.d. train data of size  $|D| = N_{\text{tr}}$ , a learner generalizing within  $k$ -coverage achieves **perfect ID generalization** w.h.p. if

$$N_{\text{tr}} \geq \tilde{O}(n^c) \quad \text{with } c = 2.5 - \frac{0.5}{k}.$$

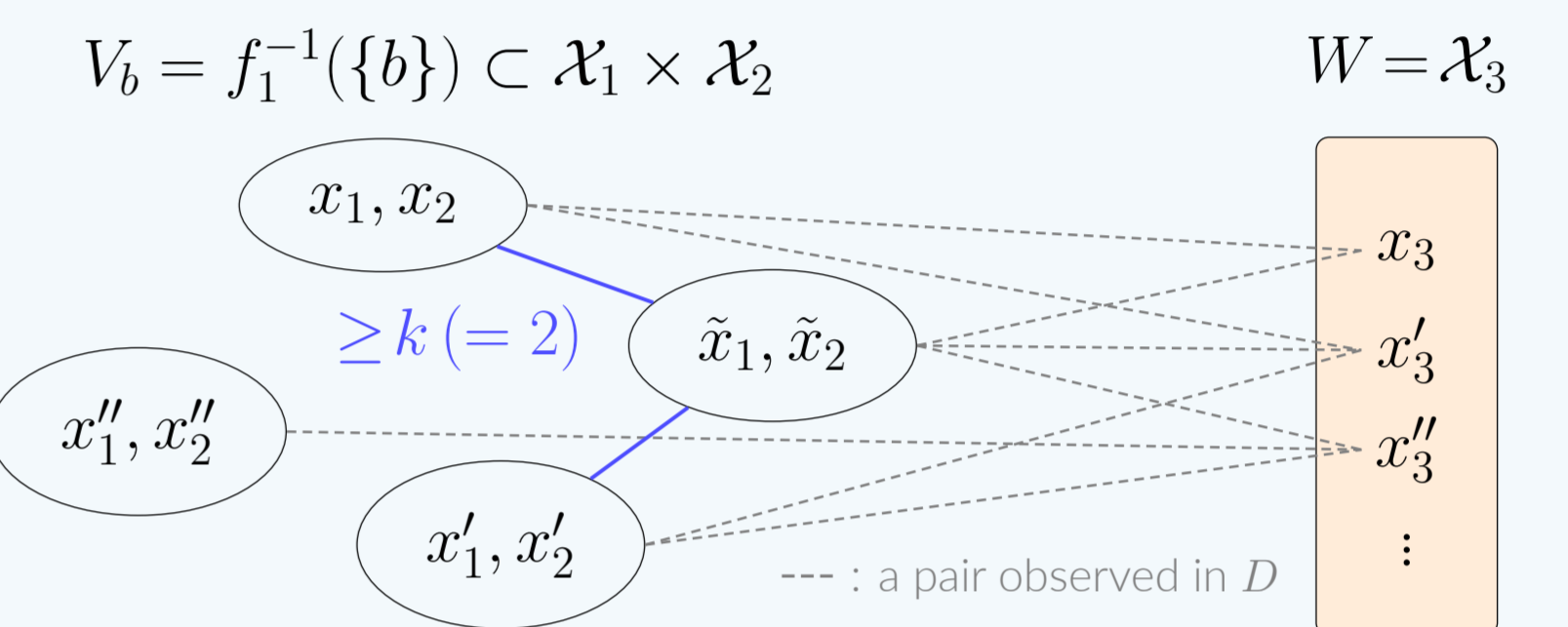
Conversely ( $k \geq 2$ ), it **cannot** for some 2-Hop task if  $n^2 \lesssim N_{\text{tr}} \leq \tilde{\Omega}(n^c)$ . (We omit  $\text{polylog}(n)$ .)

## Proof Sketch of Upper Bound (Lower bound proof is similar)

$$(x_1, x_2) \xrightarrow{f_1} b \in (\text{intermediate states}); \quad (b, x_3) \xrightarrow{f_2} y \in (\text{outputs}). \quad (2\text{-Hop task})$$

### Step 1. Reduce to graph connectivity

For each intermediate state  $b$ , define the  $b$ -evidence graph  $G_b = (V_b, E_b)$ :



**Key Lemma:** If  $G_b$  is connected ( $\forall b$ ), then  $\text{Cover}_k(D)$  covers all in-domain data. (proof omitted.)

$\Rightarrow$  We want  $\Pr[\exists b: \text{disconnected } G_b]$  small. (randomness: from i.i.d. sampling of  $D$ .)

### Step 2. Apply Poissonization [Kac49, Ald89, Joh98] for graph independence

Under i.i.d. sampling,  $G_b$ 's are **dependent** across  $b$ .  $\because \#\{(x_1, x_2, x_3) \in D\} \sim \text{Multinomial}(N_{\text{tr}}; p = \frac{1}{n^3})$ .

- We break this dependency by replacing  $\#\{(x_1, x_2, x_3) \in D\} \stackrel{\text{fixed}}{\sim} \text{Poisson}((N_{\text{tr}} - \epsilon)/n^3)$ .
- Using monotonicity of 'graph disconnectedness', we can prove:

$$\Pr_{\text{Multinomial}}[\exists b: \text{disconnected } G_b] \leq \sum_b \Pr_{\text{Poisson}}[G_b \text{ is disconnected}] + e^{-\epsilon^2/(2N_{\text{tr}})}.$$

### Step 3. Identify random $k$ -intersection graphs ( $k$ -RIG) [Sin95, ZYG14]

Under Poissonization,  $\mathbf{1}\{(x_1, x_2, x_3) \in D\} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$  with  $p = 1 - e^{-(N_{\text{tr}} - \epsilon)/n^3} \approx N_{\text{tr}}/n^3$ .

- Each  $G_b$  becomes a **binomial  $k$ -RIG  $\mathcal{G}^{(k)}(n_b, m, p)$** , with  $n_b = |V_b| \approx n$  and  $m = |\mathcal{X}_3| = n$ .
- Apply the connectivity threshold [Sin95, ZYG14]: for  $m = \Omega(n_b)$  and  $k \geq 2$ ,

$$p \geq \left(\frac{k! (2 \ln n_b)}{n_b}\right)^{\frac{0.5}{k}} \cdot \frac{1}{\sqrt{m}} \Rightarrow \Pr_{\text{Poisson}}[\mathcal{G}^{(k)}(n_b, m, p) \text{ is disconnected}] \rightarrow 0 \text{ as } n_b \rightarrow \infty.$$

We conclude the proof by plugging  $p \approx N_{\text{tr}}/n^3$  and  $n_b \approx n$ , and solving for  $N_{\text{tr}}$ .

## Discussion & Practical Implications

- **Long-tail:** Low-frequency combinations receive low  $k \Rightarrow$  effectively outside coverage
  - **Reversal curse:** "A is B" gives no functional equivalence evidence for "B is A"
  - **Planning:** Multi-path state tracking encounters path ambiguity
- $\Rightarrow$  Future work: Data augmentation to "maximize  $k$ -coverage" by diversifying shared contexts

## References

- [Ald89] David Aldous. *Probability Approximations via the Poisson Clumping Heuristic*. New York, NY: Springer New York: Imprint: Springer, 1989.
- [Joh98] Kurt Johansson. The nearly increasing subsequence in a random permutation and a unitary random matrix model. *Mathematical Research Letters*, 5(1):68–82, 1998.
- [Kac49] Marc Kac. On deviations between theoretical and empirical distributions. *Proceedings of the National Academy of Sciences*, 35(5):252–257, 1949.
- [Sin95] Karen B. Singer. *Random intersection graphs*. PhD thesis, The Johns Hopkins University, 1995.
- [ZYG14] Jun Zhao, Osman Yağan, and Virgil Gligor. On  $k$ -connectivity and minimum vertex degree in random  $s$ -intersection graphs. In *2015 Proceedings of the Twelfth Workshop on Analytic Algorithmics and Combinatorics (ANALCO)*, pages 1–15. SIAM, 2014.