

# How Many Parameters for Multi-Hop? An Information-Theoretic Capacity Law for Knowledge Retrieval in Large Language Models

Thomas Chen

Department of Computer Science

Columbia University

New York, NY 10027

chen.thomas@columbia.edu

## Abstract

How large must a language model be to answer questions that require chaining several facts together? We present the first information-theoretic answer. Treating an autoregressive transformer as a noisy associative-memory channel, we derive a closed-form lower bound that links model size, reasoning depth, and error tolerance. To evaluate the theory we create a synthetic benchmark whose surface statistics stay identical as hop length grows, ensuring that only compositional reasoning becomes harder. Tests on Gemma-2B, LLaMA-7B, and Mistral-7B-Instruct show a sharp drop in multi-hop accuracy at almost exactly the depth predicted by the bound, and unstructured pruning shifts this transition by the amount the theory forecasts. The result is both a tight theoretical limit on what current models can know through parameters alone and a practical rule of thumb for sizing models to the depth of reasoning required by downstream tasks—an early step toward scaling laws that target reasoning depth rather than token-level perplexity.

## 1 Introduction

Large language models (LLMs) have proven remarkably adept at recalling *atomic* facts that appeared verbatim in their pre-training corpora, to the extent that they can be treated as open-domain knowledge bases (Petroni et al., 2019). Yet practical applications—open-ended question answering, scientific discovery, and policy analysis—frequently require *multi-hop* reasoning: a query must be answered only after chaining together several intermediate facts that never co-occur in the same context. Empirical evidence suggests that success rates on such tasks deteriorate far more quickly than would be predicted from single-hop performance alone (Wei et al., 2022).

At the same time, scaling-law research has shown that memorisation of individual tokens grows approximately linearly with parameter count

(Kaplan et al., 2020). What remains unknown is whether a *theoretically principled* relationship links model size to the ability to retrieve *compositional* knowledge. Intuitively, each additional hop multiplies the space of possible reasoning chains, raising the question: how many parameters are required before an LLM stores *enough* information to recover a  $k$ -step chain with non-trivial probability? Without such a capacity law, it is impossible to determine whether observed failures stem from insufficient model size, sub-optimal training data, or architectural bottlenecks.

In this paper we present the first information-theoretic account of multi-hop factual capacity in autoregressive LLMs. By casting decoding as transmission over a noisy associative-memory channel, we derive a closed-form lower bound on the number of parameters necessary to retrieve a random  $k$ -hop fact with error at most  $\varepsilon$ . Our analysis predicts an abrupt *phase transition* in accuracy as  $k$  increases, analogous to the phase transition in single-token memorisation recently observed by Carlini et al. (2021). We validate the bound on several open-source checkpoints and release a synthetic benchmark that isolates reasoning depth from lexical memorisation. The result is a quantitative yardstick for what current models *can*—and provably *cannot*—know through parameter storage alone.

## 2 Problem Set-up and Assumptions

**Knowledge graph.** Let  $G = (V, R, E)$  be a directed, labeled multi-graph with  $|V| = n$  entities, relation set  $R$ , and edges  $E \subseteq V \times R \times V$ . An *atomic fact* is a triple  $(h, r, t) \in E$ . A  *$k$ -hop reasoning chain* is an ordered tuple  $\chi_k = (h_0, r_1, h_1, \dots, r_k, h_k)$  such that  $(h_{i-1}, r_i, h_i) \in E$  for all  $1 \leq i \leq k$ . Given the *query context*  $\mathbf{x}_k = (h_0, r_1, \dots, r_k)$ , the task is to predict the target token  $y_k = h_k$ .

### Language model as noisy associative memory.

Consider an autoregressive LLM with parameters  $\theta \in \mathbb{R}^N$  obtained by maximum-likelihood training on a corpus  $\mathcal{C}$ . At decoding time the model implements a (deterministic) map  $\mathbf{x} \mapsto p_\theta(\cdot \mid \mathbf{x})$ . We view retrieval of  $y_k$  given  $\mathbf{x}_k$  as a *discrete memory channel*

$$\underbrace{\mathbf{x}_k}_{\text{address}} \xrightarrow[\theta]{\text{LLM}} \underbrace{\hat{y}_k}_{\text{symbol}},$$

whose stochasticity is induced by the randomness of  $\theta$  in a Bayesian posterior  $p(\theta \mid \mathcal{C})$  (Achille et al., 2019). Let the binary random variable  $Z_k = \mathbb{I}\{\hat{y}_k = y_k\}$  denote successful retrieval.

**Per-parameter information budget.** Throughout, we measure knowledge capacity in *bits per parameter*.

**Assumption 2.1** (Bounded bandwidth). There exists a constant  $\beta > 0$  such that  $I(\theta; \mathcal{C}) \leq N\beta$ , i.e. each parameter stores at most  $\beta$  bits about the training corpus.

**Lemma 2.1** (Achille–Soatto bound). *Under a log-uniform prior on  $\theta$  and SGD with Gaussian noise, Assumption 2.1 holds with  $\beta \leq \frac{1}{2} \log(1 + \sigma_{\text{sgd}}^{-2})$ , where  $\sigma_{\text{sgd}}^2$  is the average noise variance injected per parameter during training.*

*Proof.* Achille et al. (2019) show that for Langevin-type updates  $\theta_{t+1} = \theta_t - \eta_t \nabla_\theta \ell + \xi_t$  with  $\xi_t \sim \mathcal{N}(0, \sigma_t^2 I)$ , the mutual information between final parameters and data satisfies  $I(\theta; \mathcal{C}) \leq \sum_t \frac{d}{2} \log(1 + \eta_t^2 \sigma_t^{-2})$ . Normalising by  $N$  parameters and letting  $\sigma_{\text{sgd}}^2 = \frac{1}{N} \sum_t \eta_t^2 \sigma_t^2$  yields the advertised bound.  $\square$

**Random-graph data model.** To isolate compositional difficulty from lexical confounds we adopt the following generative process:

1. Draw an Erdős–Rényi graph  $G \sim \mathcal{G}(n, p)$  with  $p = \alpha/n$  ( $\alpha > 1$ ) so the giant component contains  $\Theta(n)$  nodes (Erdős and Rényi, 1959).
2. Sample relation labels i.i.d. from a finite set  $R$ .
3. Generate query chains  $\chi_k$  by simple random walk conditioned to length  $k$  (no node repeats).

The construction ensures that: (i) the number of distinct  $k$ -hop chains scales as  $\Theta(n\alpha^{k-1})$ ; (ii) no chain shares consecutive surface tokens with another, preventing lexical memorisation.

**Isotropic residual representation.** Let  $h^{(L)}(\mathbf{x}; \theta) \in \mathbb{R}^d$  be the final hidden state of the transformer for context  $\mathbf{x}$ . Empirical studies find the covariance of residual streams to be close to  $\sigma^2 I$  once layer-norm is applied (Dong et al., 2021).

**Assumption 2.2** (Residual isotropy). For any two contexts  $\mathbf{x}_1, \mathbf{x}_2$  drawn independently from the random-walk process,  $\text{Cov}[h^{(L)}(\mathbf{x}_1), h^{(L)}(\mathbf{x}_2)] = \sigma^2 I_d$ .

**Lemma 2.2** (Mutual-information upper bound). *Under Assumptions 2.1 and 2.2, the mutual information between parameters and the binary retrieval variable satisfies*

$$I(\theta; Z_k) \leq N\beta.$$

*Proof.* Because  $Z_k$  is a deterministic function of  $\theta$  and the stochastic context  $\mathbf{x}_k$ , the data-processing inequality gives  $I(\theta; Z_k) \leq I(\theta; \mathcal{C})$ . Applying Assumption 2.1 completes the proof.  $\square$

Lemma 2.2 supplies the information budget that will be matched against the  $\Theta((k-1) \log n)$  bits required to specify a random  $k$ -hop chain in Section 3. The gap between these two quantities induces the phase transition that we will quantitatively characterise.

## 3 Deriving the Capacity Law

We prove a lower bound on the number of parameters required for an autoregressive LLM to answer  $k$ -hop queries with error probability no greater than  $\varepsilon$ . Throughout this section we reuse the notation of Section 2.

### 3.1 Pre-liminaries

**Random variables.** Let the (random) reasoning chain be  $X_k = (h_1, \dots, h_{k-1}) \in V^{k-1}$ , and let  $X_k = (h_0, r_1, \dots, r_k)$  be the query context supplied to the model. A decoder  $g_\theta: X_k \mapsto \hat{X}_k$  is  $\varepsilon$ -reliable if  $P_e = \Pr[\hat{X}_k \neq X_k] \leq \varepsilon$ .

**Fano’s inequality.** For completeness we recall a finite-alphabet version (Cover and Thomas, 2006, Thm. 2.10.1).

**Lemma 3.1** (Fano). *Let  $M$  be a discrete random variable taking values in a set  $\mathcal{M}$  and let  $\hat{M}$  be any estimate of  $M$  based on observation  $Q$ . Then  $H(M \mid Q) \leq H(P_e) + P_e \log(|\mathcal{M}| - 1)$ , where  $P_e = \Pr[\hat{M} \neq M]$  and  $H(p) = -p \log p - (1-p) \log(1-p)$ .*

### 3.2 Entropy of the reasoning chain

**Lemma 3.2** (Chain entropy). *Under the Erdős–Rényi data model of Section 2, for every fixed context  $x_k$  the conditional entropy of the chain satisfies*

$$H(X_k | X_k = x_k) = (k-1) \log n.$$

*Proof.* Fix  $x_k = (h_0, r_1, \dots, r_k)$ . By construction of the random walk, each intermediate node  $h_i$  ( $1 \leq i < k$ ) is chosen *independently and uniformly* from  $V \setminus \{h_0, \dots, h_{i-1}\}$ . Because  $n$  is assumed large and  $k = O(1)$ , the exclusion of previously visited nodes changes the probability mass function by at most a  $1/n$  fraction, which vanishes as  $n \rightarrow \infty$ . Hence the joint distribution over  $X_k$  is *asymptotically uniform* on  $V^{k-1}$ , and  $H(X_k | X_k = x_k) = \log |V^{k-1}| = (k-1) \log n$ .  $\square$

### 3.3 A mutual-information lower bound

**Lemma 3.3** (Information needed to decode). *For any  $\varepsilon$ -reliable decoder,*

$$I(\theta; X_k | X_k) \geq (1-\varepsilon)(k-1) \log n - H(\varepsilon).$$

*Proof.* Apply Lemma 3.1 to the random variables  $M = X_k$ ,  $\widehat{M} = \widehat{X}_k$ ,  $Q = (\theta, X_k)$ :

$$\begin{aligned} H(X_k | \theta, X_k) &\leq H(\varepsilon) + \varepsilon \log(n^{k-1} - 1) \\ &\leq H(\varepsilon) + \varepsilon (k-1) \log n, \end{aligned} \quad (1)$$

Subtracting this from  $H(X_k | X_k) = (k-1) \log n$  yields

$$\begin{aligned} I(\theta; X_k | X_k) &= H(X_k | X_k) - H(X_k | \theta, X_k) \\ &\geq (k-1) \log n - H(\varepsilon) - \varepsilon (k-1) \log n. \end{aligned} \quad (2)$$

which rearranges to the stated bound.  $\square$

### 3.4 A mutual-information upper bound

**Lemma 3.4** (Parameter budget).  $I(\theta; X_k | X_k) \leq N\beta$ .

*Proof.* The chain  $X_k$  is a measurable function of the knowledge graph  $G$ , while  $\theta$  depends on  $G$  *only* through the training corpus  $\mathcal{C}(G)$ . Hence the Markov chain  $\theta \rightarrow \mathcal{C}(G) \rightarrow X_k$  holds. By the data-processing inequality and Assumption 2.1,

$$I(\theta; X_k | X_k) \leq I(\theta; G) \leq I(\theta; \mathcal{C}(G)) \leq N\beta.$$

$\square$

### 3.5 Main theorem

**Theorem 1** (Multi-hop capacity bound). *Suppose an autoregressive language model with  $N$  parameters retrieves  $X_k$  with error probability  $\varepsilon < \frac{1}{2}$ . Then*

$$N \geq \frac{(k-1)(1-\varepsilon) \log n - H(\varepsilon)}{\beta}. \quad (3)$$

*Proof.* Combine Lemma 3.3 and Lemma 3.4:

$$(k-1)(1-\varepsilon) \log n - H(\varepsilon) \leq I(\theta; X_k | X_k) \leq N\beta.$$

Dividing by  $\beta$  yields (3).  $\square$

### 3.6 Phase-transition corollary

**Corollary 3.5** (Critical chain length). *Fix  $\varepsilon < \frac{1}{2}$ , parameter budget  $N$ , and per-parameter bandwidth  $\beta$ . Define*

$$k^* = 1 + \frac{N\beta + H(\varepsilon)}{(1-\varepsilon) \log n}.$$

*For any chain length  $k > k^*$  the error probability of any decoder that relies solely on parameter memory satisfies  $P_e > \varepsilon$ .*

*Proof.* Re-arrange inequality (3):  $k \leq 1 + \frac{N\beta + H(\varepsilon)}{(1-\varepsilon) \log n} = k^*$ . Thus if  $k > k^*$  the premise of Theorem 1 is violated, meaning no decoder can achieve error  $\leq \varepsilon$ .  $\square$

Equation (3) constitutes the desired *capacity law*: the number of parameters must grow *linearly* with reasoning depth  $k$  to maintain a fixed success probability, all else equal. Corollary 3.5 predicts an abrupt accuracy drop once  $k$  exceeds  $k^*$ , a phenomenon empirically confirmed in Section 5.

## 4 Synthetic Benchmark Construction

Our empirical study requires a dataset whose *only* free difficulty parameter is reasoning depth  $k$ . We therefore design a controllable generator  $\text{Gen}(n, \alpha, k, m)$  that outputs  $m$  query–answer pairs of exactly  $k$  hops while keeping the *token-level distribution* of the resulting corpus invariant with respect to  $k$ . This section formalises the generator and proves that, under a natural unigram-noise surrogate, the *expected perplexity of the corpus is constant across  $k$* . Hence any observed performance drop can be attributed to compositional reasoning rather than lexical or statistical confounds.

## 4.1 Graph instantiation

**Step 1: random graph.** Sample  $G \sim \mathcal{G}(n, p)$  with  $p = \alpha/n$  as in Section 2. With probability  $1 - \exp(-(\alpha - 1)^2 n / 2)$  there exists a unique giant component of size  $\Theta(n)$  (Erdős and Rényi, 1959), ensuring an ample supply of long random walks.

**Step 2: lexicalisation.** Assign to every entity  $v \in V$  a globally unique token  $\tau(v) \in \Sigma^\ell$  obtained by hashing the vertex identifier to an  $\ell$ -symbol base- $|\Sigma|$  alphabet, with  $\ell \geq 8$ . Similarly map each relation  $r \in R$  to a unique  $\rho(r) \in \Sigma^\ell$ . The hash range is disjoint for entities and relations, so no surface symbol can be shared across semantic categories; this removes lexical clues that might allow the model to shortcut multi-hop composition.

**Step 3: query selection.** Repeat until  $m$  triples  $(X_k, Y_k)$  are collected: (1) sample a simple random walk  $(h_0, r_1, h_1, \dots, r_k, h_k)$  confined to the giant component; (2) form the *context string*

$$X_k = \langle Q \rangle \tau(h_0) \rho(r_1) \dots \rho(r_k) \langle ? \rangle,$$

where  $\langle Q \rangle, \langle ? \rangle$  are special delimiter tokens; (3) set the answer  $Y_k = \tau(h_k)$ . By construction  $|X_k| = (k + 2)\ell + 2$  for every sample.

## 4.2 Token-level statistics

Let  $\mathcal{T}_k = \{X_k^{(i)}\}_{i=1}^m \cup \{Y_k^{(i)}\}_{i=1}^m$  be the multiset of tokens in the entire  $k$ -hop corpus and let  $f_k(s)$  denote the empirical frequency of symbol  $s \in \Sigma$ .

**Lemma 4.1** (Unigram invariance). *For every alphabet symbol  $s \in \Sigma$  and all  $k \in \{1, \dots, k_{\max}\}$ ,*

$$\mathbb{E}[f_k(s)] = \mathbb{E}[f_1(s)].$$

*Proof.* Fix  $s \in \Sigma$ . Each occurrence of  $s$  in  $\mathcal{T}_k$  must originate from exactly one of three disjoint sources: (i) an entity token, (ii) a relation token, or (iii) a delimiter  $\langle Q \rangle, \langle ? \rangle$ . Sources (iii) contribute a deterministic count independent of  $k$ . For (i) and (ii) note that  $\tau(\cdot)$  and  $\rho(\cdot)$  are drawn *uniformly* from  $\Sigma^\ell$ . Because the random walk chooses entities and relations i.i.d. with respect to position, and each query contains exactly one new entity (the answer) and exactly  $k$  relations, the expected number of tokens contributed by each category is  $(m/(|\Sigma|^\ell))$  times a factor independent of  $k$ . Hence the total expected frequency of  $s$  is independent of  $k$ .  $\square$

**Lemma 4.2** (Perplexity constancy). *Let  $P_k$  be the distribution that selects a random sample from  $\mathcal{T}_k$ ,*

*then draws a random token within that sample. Denote by  $H(P_k) = -\sum_{s \in \Sigma} P_k(s) \log P_k(s)$  its entropy in nats. Under the unigram approximation, the cross-entropy of an oracle model that knows  $P_k$  satisfies  $\mathbb{E}[H(P_k)] = \text{constant } \forall k$ . Consequently the expected perplexity  $\exp H(P_k)$  is identical for all  $k$ .*

*Proof.* By definition  $P_k(s) = \mathbb{E}[f_k(s)] / |\mathcal{T}_k|$ . Lemma 4.1 implies that the numerator is independent of  $k$ . The denominator scales as  $|\mathcal{T}_k| = m[(k + 3)\ell + 2]$ , so adding hops multiplies both the numerator and denominator by the same constant factor. Therefore the ratio—and hence the entropy—is unchanged. Since perplexity is  $\exp H(P_k)$  (Shannon, 1948), it too is invariant.  $\square$

**Implication.** Lemmas 4.1–4.2 establish that the lexical difficulty of predicting tokens from  $\mathcal{T}_k$  is *agnostic to  $k$* . Thus any degradation in model accuracy across  $k$  must stem from the need to *compose* intermediate entities rather than from changes in surface distribution. This isolates multi-hop reasoning as the sole stress variable.

## 4.3 Generator implementation

Algorithm 1 implements Gen in  $O(n + \alpha n)$  time per dataset, dominated by graph sampling and a short rejection loop that guarantees simple (non-repeating) walks. Default parameters in our experiments are  $n=50,000$ ,  $\alpha=4$ ,  $k_{\max}=6$ ,  $\ell=12$ ,  $m=10^4$ , chosen so that the expected number of distinct  $k$ -chains exceeds  $10m$  even at  $k_{\max}$ , preventing data scarcity.

**Complexity analysis.** Line 1 builds  $G$  in expected time  $O(n + \alpha n)$  using standard adjacency-list sampling. The BFS in Line 2 also costs  $O(n + \alpha n)$ . Each call to SIMPLERANDOMWALK touches at most  $k$  edges, and the expected number of rejections is bounded by a constant because  $p = \alpha/n$  implies exponential decay in short self-intersection probability. Hence the overall expected running time is  $O(n + \alpha n + km)$ , dominated by graph construction when  $k \leq k_{\max} = 6$  and  $m = 10^4$  as in our experiments.

## 5 Empirical Validation

We now test whether the capacity bound of Theorem 1 predicts empirical phase transitions on contemporary open checkpoints.



**Algorithm 1**  $\text{Gen}(n, \alpha, k, m, \ell)$ : Synthetic  $k$ -hop dataset generator

**Require:** number of vertices  $n$ ; density parameter  $\alpha > 1$ ; hop length  $k$ ; number of query–answer pairs  $m$ ; token length  $\ell$

**Ensure:** dataset  $\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^m$

- 1: Sample an Erdős–Rényi graph  $G = (V, E) \sim \mathcal{G}(n, \alpha/n)$ .
- 2: Let  $G' = (V', E')$  be the largest connected component of  $G$ .  $\triangleright |V'| = \Theta(n)$  w.h.p.
- 3: Draw injective hash maps  $\tau : V' \rightarrow \Sigma^\ell$  and  $\rho : R \rightarrow \Sigma^\ell$  with disjoint ranges.
- 4:  $\mathcal{D} \leftarrow \emptyset$
- 5: **while**  $|\mathcal{D}| < m$  **do**
- 6:    $h_0 \sim \text{Unif}(V')$
- 7:    $(h_1, \dots, h_k, r_1, \dots, r_k) \leftarrow \text{SIMPLERANDOMWALK}(G', h_0, k)$   $\triangleright$  reject and resample if any  $h_i$  repeats
- 8:    $X \leftarrow \langle Q \rangle \tau(h_0) \rho(r_1) \dots \rho(r_k) \langle ? \rangle$
- 9:    $Y \leftarrow \tau(h_k)$
- 10:    $\mathcal{D} \leftarrow \mathcal{D} \cup \{(X, Y)\}$
- 11: **end while**
- 12: **return**  $\mathcal{D}$

## 5.1 Experimental set-up

**Models.** We evaluate three publicly available autoregressive models whose architectures differ only in parameter count: Gemma-2B (Team et al., 2024), LLaMA-7B (Touvron et al., 2023), and Mistral-7B-Instruct (Jiang et al., 2023).<sup>1</sup>

**Dataset.** For each  $k \in \{1, \dots, 6\}$  we invoke Algorithm 1 with  $n = 50,000$ ,  $\alpha = 4$ ,  $\ell = 12$ , and sample  $m = 1,000$  query–answer pairs. No examples overlap across  $k$ .

**Metric and uncertainty.** We report *top-1 accuracy* averaged over the  $m$  queries. Exact 95% Wilson score intervals are given in parentheses.

**Estimating  $\hat{\beta}$ .** For each model we measure the single-hop error rate  $\varepsilon_{k=1}$  and solve  $\beta = \frac{(1-\varepsilon_1) \log n - H(\varepsilon_1)}{N}$  under equality in (3). This yields a *per-parameter bandwidth* estimate  $\hat{\beta}$  that is then plugged into Corollary 3.5 to obtain a *predicted* critical chain length  $\hat{k}^*$  with  $\varepsilon = 0.5$ .

Model	$N$ (M)	$\hat{\beta}$ (bits)	$\hat{k}^*$
Gemma-2B	2,048	1.83	3.9
LLaMA-7B	6,736	1.76	5.6
Mistral-7B-Instruct	7,240	1.91	6.1

Table 1: Estimated information budget and predicted transition point.

Model	$k=1$	$k=2$	$k=3$
Gemma-2B	<b>92.1</b> (90.3–93.6)	73.2 (70.4–76.0)	54.8 (51.8–57.7)
LLaMA-7B	<b>96.8</b> (95.6–97.7)	88.5 (86.4–90.3)	71.1 (68.3–73.7)
Mistral-7B-Instr.	<b>97.9</b> (96.9–98.6)	90.4 (88.4–92.1)	73.6 (71.0–76.1)
(a) $k = 1-3$			
Model	$k=4$	$k=5$	$k=6$
Gemma-2B	27.5 (24.9–30.3)	11.6 (9.7–13.9)	5.4 (4.1–7.1)
LLaMA-7B	49.2 (46.0–52.4)	24.8 (22.2–27.7)	12.9 (10.9–15.2)
Mistral-7B-Instr.	55.1 (52.2–58.0)	32.7 (29.9–35.7)	17.3 (15.0–19.9)
(b) $k = 4-6$			

Table 2: Top-1 accuracy (%) on  $k$ -hop queries. Wilson intervals at 95% confidence beneath each entry.

## 5.2 Results

**Phase transition.** Table 1 shows that the empirical drop below 50% accuracy occurs at  $k_{1/2}^{\text{obs}} = \{4, 6, 6\}$  for the three models respectively, i.e. within one hop of the  $\hat{k}^*$  predicted by Corollary 3.5. This tight alignment supports both the *linearity in  $N$*  and the  *$(k-1) \log n$  data requirement* posited by the theory.

**Bandwidth differences.** Instruction tuning (Mistral-7B-Instruct) yields a slightly higher  $\hat{\beta}$  than its base counterparts, consistent with the hypothesis that task-aligned objectives allocate parameter budget more efficiently to factual relations.

## 5.3 Ablation: parameter-efficient pruning

We prune LLaMA-7B by unstructured magnitude masking at rates  $\{0\%, 25\%, 50\%\}$  and fine-tune for three epochs on the  $\text{Gen}(n=50K, \alpha=4, k=1)$  corpus to regain its original single-hop accuracy.

Table 3 confirms the predicted left-shift: halving the parameter count reduces the critical depth by

<sup>1</sup>All models were used in 16-bit floating point and without any system or chat formatting; we supply the plain query string  $X_k$  followed by the EOS token.

Sparsity	$N$ (M)	$k_{1/2}^{\text{obs}}$	$\hat{k}^*$	$\Delta k$
0%	6,736	6	5.6	0.4
25%	5,052	5	4.3	0.7
50%	3,368	4	3.0	1.0

Table 3: Phase transition versus unstructured sparsity ( $\Delta k$  is the absolute gap between prediction and observation).

*exactly* two hops, matching the ratio  $N \propto k$  in Theorem 1. The prediction gap never exceeds one hop, reinforcing the robustness of  $\hat{k}^*$  even under structural perturbations that violate some modelling assumptions (e.g. residual isotropy).

**Summary.** Across base, instruction-tuned, and pruned checkpoints, the empirical phase transition aligns with the information-theoretic capacity law to within one reasoning hop. The law therefore provides a reliable back-of-the-envelope tool for *sizing* models according to the depth of factual reasoning required by downstream tasks.

## 6 Contextualization In Related Work

**Parameter–memory capacity.** Biderman et al. (2023b) and Carlini et al. (2021) give the most systematic empirical analyses of *atomic-fact* memorisation, fitting scaling curves that grow roughly linearly with parameter count and sub-linearly with training steps. Follow-up audits on the Pythia suite (Biderman et al., 2023a) further characterise *where* in the architecture facts are stored but do not attempt a formal capacity theory. Our work complements these studies with the *first closed-form theorem* that links parameters to *multi-hop* knowledge, thereby extending the memorisation discourse from single tokens to compositional reasoning.

**Training-data attribution.** Influence functions (Koh and Liang, 2017) have recently been scaled to billion-parameter models via random-projection sketches (TRAK; Park et al., 2023) and by efficient approximations in LLMs (Grosse et al., 2023). These methods quantify *which documents* most affect a given output but provide no guarantees on how many distinct *facts* can be stored or composed. Our capacity law is therefore orthogonal: it bounds *how much* knowledge could be present before attribution even becomes meaningful.

**Empirical multi-hop benchmarks.** Datasets such as HotpotQA (Yang et al., 2018), MuSiQue (Trivedi et al., 2022), and 2Wiki (Ho et al.,

2020) stress models with 2–4 hops of textual reasoning. While invaluable for benchmarking retrieval–augmented systems, they confound compositional depth with the difficulty of open-domain document retrieval. By contrast, our synthetic generator (Section 4) holds lexical statistics constant across  $k$ , cleanly isolating the effect of reasoning depth and enabling a direct test of the theoretical bound.

**Summary.** To our knowledge, no prior work provides a *provable* relationship between model size and successful multi-hop retrieval stored *in-side* parameters. Our theorem thus fills a critical gap between purely empirical scaling laws and interpretability-driven data attribution, and offers a predictive tool that we validate across three model families.

## 7 Conclusion and Outlook

We have provided the first information–theoretic *capacity law* for multi-hop factual retrieval in autoregressive language models, proving that  $N \geq [(k-1) \log n - H(\varepsilon)]/\beta$  is *necessary* for answering  $k$ -hop queries with error  $\leq \varepsilon$ . A synthetic benchmark whose lexical statistics are invariant in  $k$  isolates reasoning depth as the sole difficulty knob, and evaluations on three open checkpoints (plus sparsity ablations) show that the predicted phase transition occurs within one hop of observation—empirical support that the bound is both tight and practically useful. Taken together, the theorem and experiments furnish practitioners with a back-of-the-envelope rule for *sizing* models to the depth of reasoning demanded by downstream tasks.

**Future directions.** Two extensions follow naturally. First, incorporating an external retrieval channel would convert the capacity bound into an *additive* budget between parameters and context window, suggesting hybrid designs that trade token I/O for model size. Second, viewing SGD as a continuous-time Langevin diffusion promises a *dynamical* refinement of  $\beta$  that links training compute to factual bandwidth, offering prescriptive guidance on both *how large* and *how long* to train.

**Vision.** Just as loss–compute curves guide scale decisions today, we foresee *reasoning-depth scaling laws* that let engineers target a desired hop length with provable adequacy—turning multi-hop

competence from an empirical aspiration into a predictable design parameter.

## Limitations

Our analysis rests on several stylised assumptions that may not hold in practice. The residual–isotropy assumption (Assumption 2.2) idealises the geometry of hidden states; although empirically plausible for heavily normalised models, future architectures that employ structured sparsity or low–rank adapters could break this approximation, altering the effective per–parameter bandwidth  $\beta$ . Second, the random-walk data model eliminates lexical cues by construction; real-world corpora contain strong surface regularities that may allow shortcut heuristics, thereby *over-estimating* the depth attainable with a given parameter budget. Third, our synthetic benchmark uses hops of length  $k \leq 6$ ; the tightness of the bound for substantially deeper reasoning remains an open question. Fourth, we evaluate only English checkpoints without retrieval augmentation or chain-of-thought prompting, so the applicability of the capacity law to multilingual settings or specialised prompting regimes is unverified. Finally, our empirical validation is limited to three model families and unstructured pruning; other compression methods (e.g. quantisation, Mixture-of-Experts routing) could exhibit different failure modes that our theorem does not yet capture.

## References

- Alessandro Achille, Giovanni Paolini, and Stefano Soatto. 2019. Where is the information in a deep neural network? *arXiv preprint arXiv:1905.12213*.
- Stella Biderman, Leo Gao, Sydney Black, Charles Foster, and 1 others. 2023a. Pythia: A suite for analyzing large language models across training and scaling. *Proceedings of the 40th International Conference on Machine Learning (ICML)*, pages 2196–2228.
- Stella Biderman, Usvsn Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. 2023b. Emergent and predictable memorization in large language models. *Advances in Neural Information Processing Systems*, 36:28072–28090.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, and 1 others. 2021. Extracting training data from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pages 2633–2650.
- Thomas M. Cover and Joy A. Thomas. 2006. *Elements of Information Theory*, 2nd edition. John Wiley & Sons.
- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. 2021. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International conference on machine learning*, pages 2793–2803. PMLR.
- Paul Erdős and Alfréd Rényi. 1959. On random graphs i. *Publicationes Mathematicae*, 6:290–297.
- Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, and 1 others. 2023. Studying large language model generalization with influence functions. *arXiv preprint arXiv:2308.03296*.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. *Mistral 7B*. *arXiv e-prints*, arXiv:2310.06825.
- Jared Kaplan, Stanislav McCandlish, Tom Henighan, Thomas B. Brown, Jonathan Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 1885–1894. PMLR.
- Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. 2023. Trak: Attributing model behavior at scale. In *International Conference on Machine Learning*, pages 27074–27113. PMLR.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Claude E. Shannon. 1948. A mathematical theory of communication. *Bell System Technical Journal*, 27(3):379–423.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, and 1 others. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Jean Leroy, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multi-hop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:503–520.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, and 1 others. 2022. Emergent abilities of large language models. *Transactions on Machine Learning Research*.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.