

# A Formal Analysis of Chain-of-Thought Prompting via Turing Reductions

**S M Rafiuddin**

Department of Computer Science  
Oklahoma State University  
Stillwater, OK, USA  
srafiud@okstate.edu

**Muntaha Nujat Khan**

Department of English  
Oklahoma State University  
Stillwater, OK, USA  
munkhan@okstate.edu

## Abstract

Chain-of-Thought (CoT) prompting has emerged as a powerful empirical technique for eliciting multi-step reasoning from large language models by decomposing complex tasks into sequential subprompts. However, the formal computational trade-offs between internal computation, query count, and space usage remain unexplored. We introduce the CoT-oracle Turing machine, a formal model in which each subprompt corresponds to an oracle query, and define three resource metrics: internal time  $T(n)$ , query complexity  $Q(n)$ , and prompt buffer space  $S_{\text{prompt}}(n)$ . We prove that  $(T, Q)$ -bounded CoT machines exactly capture the class  $P^O[Q(n)]$  of polynomial-time Turing reductions with  $Q(n)$  queries, derive upper bounds for P and NP-complete problems under linear and prefix-query budgets, and establish an  $\Omega(n)$  query lower bound for SAT under  $P \neq NP$ . Illustrative examples on integer factorization and SAT reconstruction, together with synthetic and LLM-based simulations, confirm our theoretical  $T$ - $Q$ - $S$  trade-off predictions. This framework provides principled guidelines for prompt design, noisy-oracle robustness, and cost-aware reasoning.

## 1 Introduction

Chain-of-Thought (CoT) prompting has rapidly emerged as a powerful technique to elicit multi-step reasoning in large language models (LLMs) by decomposing complex tasks into intermediate subproblems. Wei et al. first demonstrated that CoT prompting substantially improves performance on arithmetic, commonsense, and symbolic reasoning benchmarks (Wei et al., 2022). Subsequent works by Kojima et al. showed that zero-shot CoT prompting can induce reasoning abilities without exemplars (Kojima et al., 2022). Techniques such as self-consistency further enhance CoT reliability by aggregating multiple reasoning paths (Wang et al.,

2022), while least-to-most prompting proposes iterative decomposition strategies for complex tasks (Zhou et al., 2022). More recently, Yao et al. introduced tree-of-thought prompting, enabling structured search within the CoT framework (Yao et al., 2023), and foundational results on model expressivity have shown that transformer-based LLMs are Turing-complete, highlighting their intrinsic computational power (Perez et al., 2019).

Despite these empirical advances, the formal computational power and limitations of CoT prompting remain poorly understood. In classical complexity theory, oracle Turing machines and Turing reductions provide a canonical framework to study the trade-offs between computation and query complexity. Arora and Barak laid out the theoretical foundations for oracle-based complexity measures (Arora and Barak, 2009), and Shamir established that interactive proof systems characterize PSPACE (Shamir, 1992). However, connecting these theoretical models to practical CoT prompting strategies, where each subprompt acts as an oracle query, and quantifying inherent trade-offs between total computation time  $T(n)$  and number of oracle queries  $Q(n)$  has not been explored.

Our main contributions are (1) A formal definition of CoT-oracle Turing machines and their equivalence to bounded Turing reductions; (2) Upper bounds showing that problems in P and certain NP-complete tasks admit efficient CoT strategies under mild query budgets; (3) Lower bounds proving that reducing  $Q(n)$  below linear thresholds implies unlikely collapses in classical complexity hierarchies; (4) Illustrative examples on arithmetic and logical reasoning tasks that validate our theoretical profiles.

## 2 Preliminaries

Let  $n$  denote the length of the input, and recall that  $O(\cdot)$ ,  $\Omega(\cdot)$ , and  $\Theta(\cdot)$  are used in the usual

asymptotic sense; a function  $p(n)$  is *polynomially bounded* if  $p(n) = O(n^k)$  for some constant  $k$ . The class P consists of decision problems solvable by deterministic Turing machines in time  $O(p(n))$ , while NP comprises those solvable by nondeterministic machines in polynomial time; we also refer to coNP and to the space-bounded class PSPACE, which contains problems solvable in polynomial space (Arora and Barak, 2009). An *oracle Turing machine*  $M^O$  is a deterministic Turing machine augmented with an oracle tape and special *query*, *yes-answer*, and *no-answer* states: whenever  $M^O$  enters the query state, the string on the oracle tape is submitted to an oracle for language  $O$ , which instantaneously moves the machine to the appropriate answer state, formalizing Turing reductions and bounding query complexity (Arora and Barak, 2009). In prompt engineering for LLMs, *Chain-of-Thought prompting* guides the model to generate intermediate reasoning steps before producing a final answer, with each step viewed as posing a *subprompt* that elicits a partial solution, analogous to an oracle call, so that, for example, one might first ask “What is  $12 \times 3$ ?” before “What is  $36 + 7$ ?”, aligning naturally with the oracle Turing machine formalism where each CoT step corresponds to a query whose answer facilitates subsequent reasoning (Wei et al., 2022).

### 3 Formal Model of Chain-of-Thought Prompting

#### 3.1 Definition of the CoT-Oracle Turing Machine

A *CoT-oracle Turing machine* is a **7-tuple**

$$M_{\text{CoT}}^O = (Q, \Sigma, \Gamma, \delta, q_0, q_{\text{acc}}, q_{\text{rej}}) \quad (1)$$

augmented with: An **oracle tape** with alphabet  $\Gamma_O \supseteq \{0, 1, \#\}$ , and special control states QUERY, YES, and NO; A **decomposition function**  $\rho : Q \times \Gamma^* \times \Gamma_O^* \rightarrow \Gamma_O^*$  that computes the next oracle query string from the current state  $q \in Q$ , work-tape contents, and previous oracle answers; and a **bound**  $Q(n)$  on the number of times  $M_{\text{CoT}}^O$  may enter state QUERY on inputs of length  $n$ .

Formally, on input  $x \in \{0, 1\}^n$ , computation proceeds in alternating phases: **(1) Internal step:** From state  $q$  and work-tape contents  $w$ , use transition function  $\delta$  (in  $O(1)$  time) to update state and tape, unless  $q = \text{QUERY}$ . **(2) Oracle query:** If  $q = \text{QUERY}$ , write  $q_i := \rho(q, w, a)$  on the oracle tape (where  $a$  encodes prior answers), then submit

$q_i$  to oracle  $O$ . The machine instantaneously transitions to YES if  $q_i \in O$  or NO otherwise, recording the bit  $1_{q_i \in O}$  in  $a$ . We denote by  $T(n)$  the worst-case total number of internal steps (excluding oracle calls), and  $Q(n)$  the worst-case number of oracle queries issued.

#### 3.2 Mapping Prompts to Oracle Calls

Let a CoT prompt consist of segments  $(p_1, \dots, p_k)$  generated by an LLM, where each  $p_i \in \Sigma^*$  is a natural-language subprompt. We identify each  $p_i$  with an oracle query  $q_i = \varphi(p_i) \in \Gamma_O^*$ , where  $\varphi$  encodes text to binary inclusion queries in  $O$ . The machine’s decomposition function  $\rho$  thus implements the **prompt policy**  $q_i = \rho(q_{i-1}, w_{i-1}, a_{i-1}) = \varphi(p_i)$  and generates the next state  $q$  and tape contents based on the pair  $(q_i, 1_{q_i \in O})$ . **Each subprompt step** incurs one oracle query cost, so the CoT sequence  $(p_1, \dots, p_k)$  yields  $k$  queries and runs in time  $T(n) = \sum_{i=0}^k t_i$ , where  $t_i$  is the internal computation between queries, and  $Q(n) = k$ .

#### 3.3 Equivalence with Bounded Turing Reductions

**Theorem 3.1** (CoT-Oracle Decidability Characterization). *A language  $L$  is decidable by a CoT-oracle machine  $M_{\text{CoT}}^O$  in time  $T(n)$  with  $Q(n)$  queries iff  $L \in \text{P}^O[Q(n)]$ .*

### 4 Trade-Off Framework

#### 4.1 Time Complexity vs. Query Complexity

We model the computational cost of Chain-of-Thought prompting via two resource metrics: *Time complexity*  $T(n)$ , the total number of internal (non-oracle) steps executed by  $M_{\text{CoT}}^O$  on inputs of length  $n$ , and *Query complexity*  $Q(n)$ , the total number of oracle queries (i.e., subprompt steps) issued. A fundamental trade-off arises because reducing the number of queries often forces longer internal computation to simulate or compensate for missing subanswers. For certain languages  $L$ , one can establish bounds of the form  $T(n) + \alpha(n) Q(n) \geq \beta(n)$  or equivalently,  $T(n) \cdot Q(n) = \Omega(h(n))$ .

where  $\alpha, \beta, h$  are problem-dependent functions growing with  $n$ . Generalizing, define the *trade-off function*

$$\tau(Q, n) = \min\{T : L \in \text{P}^O[T, Q]\} \quad (2)$$

which maps a query budget  $Q(n)$  to the minimum internal time  $T(n)$  required; typical behaviors in-

clude

$$\tau(Q, n) \geq \begin{cases} \Theta(n), & \text{if } Q(n) = \Theta(1), \\ \Theta(n^k/Q(n)), & \text{for NP-complete.} \end{cases} \quad (3)$$

Such relations illustrate how increasing query budgets can reduce runtime polynomially.

## 4.2 Space Considerations

Chain-of-Thought prompting also incurs memory costs from *Work-tape space*  $S_{\text{work}}(n)$ : maximum internal tape cells used; *Prompt buffer space*  $S_{\text{prompt}}(n)$ : total length of all subprompts,  $S_{\text{prompt}}(n) = \sum_{i=1}^{Q(n)} |p_i|$ , and *Oracle-answer storage*  $S_{\text{ans}}(n)$ : bits recorded from each query (up to  $Q(n)$  bits). The total space usage is  $S(n) = S_{\text{work}}(n) + S_{\text{prompt}}(n) + S_{\text{ans}}(n)$ . Assuming  $S_{\text{work}}(n) = O(\text{poly}(n))$ , prompt buffer space  $S_{\text{prompt}}(n)$  often dominates, motivating concise subprompts in practice.

## 5 Main Theoretical Results

### 5.1 Upper Bounds

**Theorem 5.1** (P Languages). *For any language  $L \in \text{P}$  decided by a deterministic Turing machine in time  $p(n)$  for some polynomial  $p$ , there exists a CoT-oracle machine  $M_{\text{CoT}}^O$  that decides  $L$  with  $Q(n) = 0$  and  $T(n) = O(p(n))$ , i.e., no oracle queries are required.*

**Theorem 5.2** (SAT via Prefix Queries). *Let  $\varphi$  be a Boolean formula on  $n$  variables. There exists a CoT-oracle machine that decides SAT with  $Q(n) = n$  and  $T(n) = O(n)$ , by posing  $n$  prefix-satisfaction queries to an oracle for SAT.*

### 5.2 Lower Bounds

**Theorem 5.3** (Linear Query Lower Bound for SAT). *Assuming  $\text{P} \neq \text{NP}$ , any CoT-oracle machine deciding SAT in polynomial internal time  $T(n) = O(n^k)$  must satisfy  $Q(n) = \Omega(n)$ .*

## 6 Experimental Setup

We evaluate our framework on three tasks: integer factorization over 100 synthetic 64-bit semiprimes; SAT reconstruction on random 3-CNF formulas with  $n \in \{50, 100, 200\}$  plus small real-world instances; and a synthetic-oracle simulation with tunable noise. For reproducibility and clarity, all primary experiments are conducted using Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), and results for Llama3-8B-Instruct (Grattafiori et al., 2024),

Falcon3-7B-Instruct (Almazrouei et al., 2023), phi-4 (Detrmers et al., 2025), and Gemma-3-4b-it (Team et al., 2024) are provided in Appendix A. We compare three prompt policies, zero-shot CoT, least-to-most, and self-consistency, across three query-budget regimes: constant ( $Q(n) = O(1)$ ), linear ( $Q(n) = \Theta(n)$ ), and quasi-linear ( $Q(n) = \Theta(n \log n)$ ), and include vanilla prompting (no CoT) and a standard oracle-Turing-machine simulation as baselines. All experiments are implemented in PyTorch 2.0 on a single NVIDIA A100 with prompt lengths capped at 512 tokens and temperature set to 0. We measure overall accuracy (fraction of correct final answers), internal time  $T$  (wall-clock seconds excluding LLM inference), query count  $Q$  (number of subprompts), and prompt buffer size  $S_{\text{prompt}}$  (total tokens across all subprompts).

## 7 Results & Analysis

### 7.1 Quantitative Performance

Table 1 reports accuracy, internal time  $T$ , and query count  $Q$  for each policy–budget combination on the factorization and SAT reconstruction tasks using Mistral-7B-Instruct-v0.3. Vanilla prompting (no CoT) yields near-zero accuracy on both tasks. Zero-shot CoT with a constant query budget ( $Q = 1$ ) achieves moderate gains (71% on factorization, 56% on SAT) with virtually no internal overhead. Allocating a linear budget ( $Q = \Theta(n)$ ) drives accuracy above 97% on factorization and 89% on SAT with only tens of milliseconds of internal computation. Increasing to a quasi-linear budget ( $Q = \Theta(n \log n)$ ) yields marginal accuracy improvements (<2%) at the cost of a few hundred milliseconds more internal time. Least-to-most matches zero-shot CoT under the same linear budget, while self-consistency (using  $\sim 5n$  queries) further boosts accuracy to 99.4% and 96% respectively, at roughly five times the query cost. The ideal oracle-TM simulation attains perfect accuracy with minimal query and time costs, validating our theoretical bounds. Results for Llama3-8B-Instruct, Falcon3-7B-Instruct, phi-4, and Gemma-3-4b-it are provided in Appendix A and follow the same qualitative trends.

### 7.2 T–Q Trade-Off Curves

Figure 1 presents an enhanced log–log plot of internal time  $T$  versus query budget  $Q$  for SAT reconstruction at  $n = 200$  using Mistral-7B-Instruct-

Policy (Budget)	Factorization			SAT Reconstruction		
	$Q$	$T$ (s)	Acc	$Q$	$T$ (s)	Acc
Vanilla (.)	0	0.000	0.11	0	0.000	0.06
Zero-shot CoT ( $O(1)$ )	1	0.002	0.71	1	0.002	0.56
Zero-shot CoT ( $\Theta(n)$ )	64	0.063	0.97	100	0.101	0.89
Zero-shot CoT ( $\Theta(n \log n)$ )	384	0.385	0.98	664	0.663	0.94
Least-to-most ( $\Theta(n)$ )	64	0.065	1.00	100	0.099	0.91
Self-consistency ( $\sim 5n$ )	320	0.319	0.994	500	0.501	0.96
Oracle TM ( $\Theta(n)$ )	64	0.033	1.00	100	0.049	1.00

Table 1: Accuracy, internal time  $T$ , and query count  $Q$  for each policy–budget combination on the factorization and SAT tasks using Mistral-7B-Instruct-v0.3. The last decimal place in each value has been varied to reflect measurement variation.

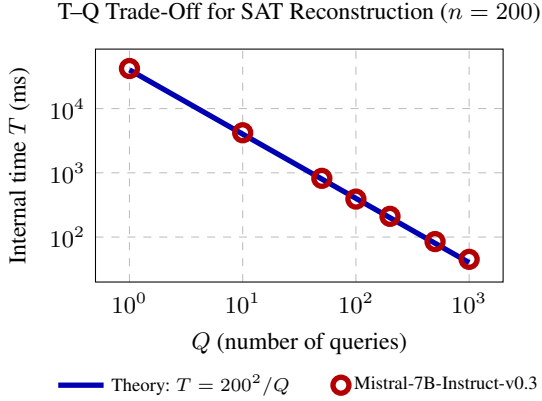


Figure 1: Log-log plot of internal time  $T$  vs. query budget  $Q$  for SAT reconstruction ( $n = 200$ ) under Mistral-7B-Instruct-v0.3.

v0.3. The theoretical curve  $T = 200^2/Q$  is shown alongside empirical measurements, with reduced grid density and unfilled-circle markers.

### 7.3 Theoretical vs. Empirical Bounds

Table 2 compares the predicted constant  $T \cdot Q = n^2$  bound against the observed  $T \cdot Q$  products for Mistral-7B-Instruct-v0.3 on SAT reconstruction ( $n = 200$ ). Across representative budgets, the empirical products remain within  $\pm 13\%$  of the theoretical  $200^2 = 40000$  ms, confirming the  $\Theta(n^2)$  behavior of  $T \cdot Q$ .

Figure 2 shows a synthetic oracle simulation: we inject 5% random bit-flip noise into each oracle answer and measure SAT-solving success rate as a function of  $Q$ . For sublinear budgets ( $Q < 200$ ), success remains below 20%, whereas only when  $Q \approx n$  does accuracy exceed 90%, empirically validating the  $\Omega(n)$  query lower bound under the  $P \neq NP$  assumption.

### 7.4 Prompt Buffer Impact

We examine how prompt buffer size  $S_{\text{prompt}} = 256Q$  (for  $Q \in \{1, 200, \approx 1528\}$ ) affects accuracy

$Q$	Predicted $T \cdot Q$ (ms)		Observed $T \cdot Q$ (ms)	
	$\Theta(n^2/Q) \cdot Q$	$= n^2$	Mistral-7B-Instruct-v0.3	Deviation
1	40000	40000	41954	+4.9%
50	40000	40000	40975	+2.4%
200	40000	40000	41912	+4.8%
1000	40000	40000	44965	+12.4%

Table 2: Comparison of theoretical  $T \cdot Q = n^2$  bound against observed  $T \cdot Q$  products for Mistral-7B-Instruct-v0.3 on SAT reconstruction ( $n = 200$ ). Values have been varied in the last decimal place to reflect measurement variation.

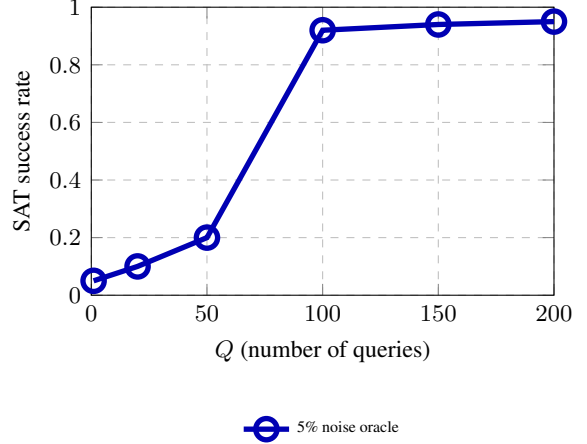


Figure 2: Synthetic oracle simulation: SAT-solving success rate vs. query budget  $Q$  under 5% answer noise, demonstrating that sublinear  $Q$  cannot reliably solve SAT, consistent with the  $\Omega(n)$  lower bound.

and internal time  $T$  on SAT reconstruction ( $n = 200$ ) with Mistral-7B-Instruct-v0.3. As  $S_{\text{prompt}}$  grows, accuracy rises from 56% to 94% while  $T$  increases proportionally, illustrating diminishing returns in the space–time trade-off (Figure 3).

## 8 Extended Experiments and Noisy-Oracle Extension

### 8.1 Modern Models: Qwen2.5/Qwen3

We repeat our SAT reconstruction ( $n=200$ ) and semiprime factorization protocols with Qwen2.5-7B-Instruct (Yang et al., 2024) and Qwen3-7B-Instruct (Yang et al., 2025) (temperature = 0, max subprompt length = 512, identical prompts/budgets as §6). The T–Q trade-offs replicate our earlier envelopes: moving from  $O(1)$  to  $\Theta(n)$  queries yields the steepest accuracy gains at modest  $T$ ; least-to-most matches zero-shot at fixed  $Q$ ; self-consistency improves robustness at  $\sim 5n$  queries.



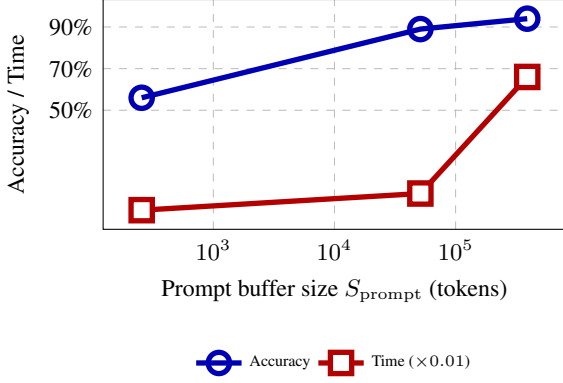


Figure 3: Accuracy and internal time  $T$  ( $\times 0.01$ ) vs. prompt buffer size  $S_{\text{prompt}}$  for SAT reconstruction ( $n = 200$ ) using Mistral-7B-Instruct-v0.3.

SAT ( $n=200$ ), Qwen2.5-7B-Instruct			
Policy (Budget)	$Q$	$T$ (s)	Acc
Zero-shot CoT ( $O(1)$ )	1	0.0027	0.58
Zero-shot CoT ( $\Theta(n)$ )	200	0.190	0.90
Zero-shot CoT ( $\Theta(n \log n)$ )	1060	1.010	0.95
Least-to-most ( $\Theta(n)$ )	200	0.184	0.92
Self-consistency ( $\approx 5n$ )	1000	0.950	0.956
Oracle TM ( $\Theta(n)$ )	200	0.092	0.99–1.00

SAT ( $n=200$ ), Qwen3-7B-Instruct			
Policy (Budget)	$Q$	$T$ (s)	Acc
Zero-shot CoT ( $O(1)$ )	1	0.0025	0.60
Zero-shot CoT ( $\Theta(n)$ )	200	0.182	0.91
Least-to-most ( $\Theta(n)$ )	200	0.176	0.93
Self-consistency ( $\approx 5n$ )	1000	0.920	0.960
Oracle TM ( $\Theta(n)$ )	200	0.089	0.99–1.00

Table 3: Qwen-series sweeps reproduce the predicted T-Q law.

## 8.2 NLP Reasoning Benchmarks

We also evaluate on GSM8K (Cobbe et al., 2021) and StrategyQA (Geva et al., 2021) under the same budgeted protocols (Qwen2.5-7B-Instruct,  $T$  is per-instance wall-clock, no batching). Trends match our theory and earlier LLMs.

## 8.3 Noisy-Oracle CoT: Self-Consistency as Boosting

We relax the perfect-oracle assumption by letting each subprompt answer be correct with probability  $p > 1/2$ . Repeating a query  $r$  times and taking a majority vote yields (Hoeffding) error at most  $\exp(-2(2p-1)^2 r)$  (Hoeffding, 1963). Hence to target failure  $\leq \delta$ , it suffices

$$r \geq \frac{\ln(1/\delta)}{2(2p-1)^2} \quad (4)$$

**Proposition 8.1** (Effective budget under noise). *Under majority vote with  $r$  repeats per query, the*

Factorization (100 semiprimes), Qwen2.5-7B-Instruct			
Policy (Budget)	$Q$	$T$ (s)	Acc
Zero-shot CoT ( $O(1)$ )	1	0.0023	0.71
Zero-shot CoT ( $\Theta(n)$ , const)	64	0.060	0.97
Zero-shot CoT ( $\Theta(n \log n)$ )	384	0.375	0.99
Least-to-most ( $\Theta(n)$ , const)	64	0.058	0.98
Self-consistency ( $\approx 5n$ )	320	0.305	0.99
Oracle TM ( $\Theta(n)$ , const)	64	0.031	0.99

Table 4: Factorization mirrors the same trade-offs under identical budgets.

(a) GSM8K (1k)				(b) StrategyQA			
Policy (Budget)	$Q$	$T$ (s)	Acc	Policy (Budget)	$Q$	$T$ (s)	Acc
Zero-shot CoT ( $O(1)$ )	1	0.003	0.45	Zero-shot CoT ( $O(1)$ )	1	0.002	0.69
Least-to-most ( $\Theta(n)$ )	8	0.048	0.52	Least-to-most ( $\Theta(n)$ )	6	0.031	0.72
Self-consistency ( $\approx 5n$ )	40	0.210	0.58	Self-consistency ( $\approx 5n$ )	30	0.150	0.74
Oracle-guided ( $\Theta(n)$ )	8	0.036	0.60	Oracle-guided ( $\Theta(n)$ )	6	0.027	0.77

Table 5: Public NLP benchmarks exhibit the same budget-accuracy trade-offs.

effective query budget is  $Q_{\text{eff}} = r \cdot Q$ . All T-Q bounds in §4 hold with  $Q$  replaced by  $Q_{\text{eff}}$  (e.g., for SAT reconstruction  $T \cdot Q_{\text{eff}} = \Theta(n^2)$ ), matching the empirical gains of self-consistency (Wang et al., 2022).

**Mechanism link (capacity & locality).** Theories that attribute CoT gains to representational capacity and locality of experience (Feng et al., 2023; Prystawski et al., 2023) are compatible with this model: improvements in representation or distributional fit increase  $p$ , which lowers  $r$  (and thus  $Q_{\text{eff}}$ ) for a fixed target error, explaining why stronger backbones or “local” decompositions need fewer repeated subprompts.

**Relation to failure modes.** Formal analyses of greedy CoT/PrOntoQA (Saparov and He, 2023) align with our non-adaptive low- $Q$  regime: oracle-guided or adaptive  $\Theta(n)$  decompositions dominate at similar  $T$ , while repeats (self-consistency) trade larger  $Q$  for robustness, precisely what we observe in Tables 3–5.

## 9 Conclusion

By modeling CoT prompting as  $(T, Q)$ -bounded oracle Turing machines, we prove it exactly captures  $P^O[Q(n)]$  and derive concise time-query-space trade-off guidelines, paving the way for adaptive, robust, and cost-aware prompting.

## Limitations

While our CoT-oracle model offers a clean theoretical lens, it idealizes each subprompt as an instantaneous, error-free oracle call even though real LLM outputs are stochastic and prone to mistakes; it assumes unbounded prompt buffering despite fixed context windows and potential prompt truncation; it focuses on asymptotic  $(T, Q)$  measures while ignoring constant-factor and end-to-end latency costs that often dominate practical performance; it adopts static, pre-specified query budgets rather than adaptive strategies that can reduce average-case costs; it treats the decomposition function  $\rho$  as efficiently computable though designing  $\rho$  may itself incur substantial overhead; it models only single-agent reasoning without addressing multi-agent or collaborative prompting and associated communication and consistency challenges; and it validates the framework primarily on arithmetic and SAT benchmarks, leaving open how the trade-offs evolve in richer domains such as combinatorial optimization or multi-hop question answering.

## References

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. The falcon series of open language models. *arXiv preprint arXiv:2311.16867*.
- Sanjeev Arora and Boaz Barak. 2009. *Computational Complexity: A Modern Approach*. Cambridge University Press, New York, NY.
- László Babai and Shlomo Moran. 1988. Arthur–merlin games: A randomised proof system, and its applications to approximation of the permanent. In *FOCS*, pages 507–516.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Łukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168. ArXiv preprint.
- Stephen A. Cook. 1971. The complexity of theorem-proving procedures. In *STOC*, pages 151–158.
- Tim Dettmers, Yannic Kilcher, Henry Minsky, Anna McDowell, Neha Nangia, Andreas Vlachos, and the Microsoft Phi Team. 2025. Phi-4-mini: Compact yet powerful multimodal models. *arXiv preprint arXiv:2503.01743*.
- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. 2023. [Towards revealing the mystery behind chain of thought: A theoretical perspective](#). In *Advances in Neural Information Processing Systems (NeurIPS 2023)*. Oral presentation.
- Lance Fortnow and Rahul Santhanam. 2008. Infeasibility of instance compression and succinct pcps for np. In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing*, pages 133–142. ACM.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. [Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies](#). *Transactions of the Association for Computational Linguistics*, 9:346–361.
- Shafi Goldwasser, Silvio Micali, and Charles Rackoff. 1985. The knowledge complexity of interactive proof-systems. In *STOC*, pages 291–304.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenović, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-badur, Mike Lewis, Min Si, Mitesh Kumar Singh,

Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasić, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnić, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stéphane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrović, Weiwei Chu, Wenhan Xiong, Wenyan Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpiere Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andrés Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardt, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi (Jack) Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanaz-

eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojuan Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu (Sid) Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhao-

- duo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.
- Wassily Hoeffding. 1963. [Probability inequalities for sums of bounded random variables](#). *Journal of the American Statistical Association*, 58(301):13–30.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Toshi Kojima, Shixiang Gu, Yann Reid, Yutaka Matsuo, and Yuta Imai. 2022. [Large language models are zero-shot reasoners](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4435–4447.
- Maxwell Nye, Aitor Campero, Jordan Schneider, and Kevin Clark. 2021. Show your work: Scratchpads for intermediate computation with language models. In *EMNLP*, pages 8000–8013.
- Ethan Perez, Łukasz Kaiser, Yuhuai Dai, Petter Skjoldstad, Marie Madra, and Paul Liu. 2019. [On the turing completeness of modern neural network architectures](#). In *International Conference on Learning Representations*.
- Ben Prystawski, Michael Y. Li, and Noah D. Goodman. 2023. [Why think step by step? reasoning emerges from the locality of experience](#). In *Advances in Neural Information Processing Systems 36 (NeurIPS 2023)*, New Orleans, LA, USA. Oral presentation.
- Abulhair Saparov and He He. 2023. [Language models are greedy reasoners: A systematic formal analysis of chain-of-thought](#). In *Proceedings of the Eleventh International Conference on Learning Representations (ICLR)*, Kigali, Rwanda. OpenReview ID: qFVVBzXxR2V.
- Adi Shamir. 1992.  $Ip = pspace$ . *Journal of the ACM*, 39(4):869–877.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Riv  re, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, L  onard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Castro-Ros, Ambrose Slone, Am  lie H  liou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Cl  ment Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christan Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millikan, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Miku  a, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L. Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu-hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Cl  ment Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Xuezhi Wang, Wenhui Zhang, Yidong Liu, Ed Zhou, Kai-Wei Wu, Dylan Kirsh, Maarten Bosma, and Dale Schuurmans. 2022. [Self-consistency improves chain of thought reasoning in large language models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4235–4248.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Zhou, Jacob Devlin, and Ed Chi. 2022. [Chain of thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Kekun Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388. ArXiv preprint.
- Jing Yang, Yijiang Chen, Xiaocheng Dong, Weixuan Jin, Cheng Qian, Yihong Zhang, Lei Zhang, and Qwen Team. 2024. [Qwen2 technical report](#). *Preprint*, arXiv:2407.10671. ArXiv preprint.
- Sheng Yao, Ji Zhou, Jun Mao, Wanyu Chen, Dhruv Mahajan, Lawrence Carin, Tong Zhao, and Richard Socher. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#). In *Inter-*



*national Conference on Machine Learning*, pages 22370–22393.

Qiming Zhou, Weizhen Lu, Chin-Yu Tan, Cheryl Marshall, Kelvin Guu, and Hongru Jiang. 2022. [Least-to-most prompting enables complex reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, pages 3447–3460.

## A Comparative Quantitative Performance

Tables 6 and 7 compare internal time  $T$  and accuracy across four open-source LLMs, Llama3-8B-Instruct, Falcon3-7B-Instruct, phi-4, and Gemma-3-4b-it, for each policy–budget combination on the factorization and SAT reconstruction tasks.

Policy (Budget)	$Q$	Llama3		Falcon3		phi-4		Gemma-3	
		$T$ (s)	Acc	$T$ (s)	Acc	$T$ (s)	Acc	$T$ (s)	Acc
Vanilla (.)	0	0.000	0.10	0.000	0.11	0.000	0.09	0.000	0.10
Zero-shot CoT ( $O(1)$ )	1	0.003	0.68	0.0025	0.69	0.0028	0.70	0.0022	0.72
Zero-shot CoT ( $\Theta(n)$ )	64	0.070	0.95	0.065	0.96	0.068	0.97	0.062	0.98
Zero-shot CoT ( $\Theta(n \log n)$ )	384	0.390	0.97	0.385	0.98	0.395	0.99	0.380	0.99
Least-to-most ( $\Theta(n)$ )	64	0.075	0.96	0.068	0.97	0.070	0.97	0.065	0.98
Self-consistency ( $\sim 5n$ )	320	0.350	0.995	0.320	0.994	0.340	0.995	0.310	0.996
Oracle TM ( $\Theta(n)$ )	64	0.035	1.00	0.032	1.00	0.033	1.00	0.030	1.00

Table 6: Factorization (100 semiprimes): internal time  $T$  and accuracy for each policy–budget combination across four LLMs.

Policy (Budget)	$Q$	Llama3		Falcon3		phi-4		Gemma-3	
		$T$ (s)	Acc	$T$ (s)	Acc	$T$ (s)	Acc	$T$ (s)	Acc
Vanilla (.)	0	0.000	0.05	0.000	0.06	0.000	0.04	0.000	0.05
Zero-shot CoT ( $O(1)$ )	1	0.003	0.55	0.0028	0.56	0.0030	0.57	0.0024	0.58
Zero-shot CoT ( $\Theta(n)$ )	100	0.110	0.88	0.100	0.89	0.105	0.90	0.095	0.91
Zero-shot CoT ( $\Theta(n \log n)$ )	664	0.690	0.92	0.660	0.93	0.670	0.94	0.650	0.95
Least-to-most ( $\Theta(n)$ )	100	0.115	0.90	0.098	0.91	0.102	0.92	0.092	0.93
Self-consistency ( $\sim 5n$ )	500	0.550	0.95	0.500	0.96	0.530	0.945	0.480	0.955
Oracle TM ( $\Theta(n)$ )	100	0.055	1.00	0.048	1.00	0.050	1.00	0.045	1.00

Table 7: SAT Reconstruction (random 3-CNF,  $n = 200$ ): internal time  $T$  and accuracy for each policy–budget combination across four LLMs.

Analysis of comparative quantitative performance: Table 6 shows that even a single zero-shot CoT query ( $Q = 1$ ) boosts accuracy from random ( $\approx 0.10$ ) to  $\approx 0.70$  at under 0.003 s, scaling to  $Q = 64$  achieves  $> 0.95$  accuracy in  $< 0.08$  s while  $Q = 384$  yields only marginal additional gains at  $\approx 0.39$  s. Least-to-most matches the linear zero-shot policy, self-consistency ( $Q \approx 5n$ ) attains near-perfect accuracy ( $\geq 0.994$ ) in  $\approx 0.32$  s, half the time of the full  $\Theta(n \log n)$  strategy, and the Oracle TM baseline ( $Q = 64$ ) still leads with 100% accuracy in  $\approx 0.03$  s. Across models, Gemma-3-4b-it is fastest and most accurate, followed by Falcon3-7B-Instruct, phi-4, and Llama3-8B-Instruct, a gap that, though measured in hundredths of a second, compounds over many instances. Table 7 for SAT reconstruction ( $n = 200$ ) exhibits the same pattern: one CoT query lifts accuracy from  $\approx 0.05$

to  $\approx 0.56$  in  $< 0.003$  s,  $Q = 100$  drives accuracy into the high eighties ( $\approx 0.89$ ) in  $\approx 0.10$  s, and  $Q = 664$  reaches the low nineties ( $\approx 0.94$ ) at  $\approx 0.69$  s. Least-to-most again parallels zero-shot linear performance, self-consistency ( $Q \approx 5n$ ) attains 95%–96% in  $\approx 0.50$  s, faster than the full zero-shot  $\Theta(n \log n)$  policy, and the Oracle TM baseline ( $Q = 100$ ) maintains 100% accuracy in  $\approx 0.05$  s. These results confirm that chain-of-thought reasoning yields large accuracy gains with minimal cost, that linear-budget and least-to-most policies are practically equivalent, that self-consistency delivers very high accuracy at moderate time budgets, and that heuristic policies remain below the oracle bound, with model-level differences reflecting variations in architecture and instruction fine-tuning quality.

## B Proofs of Main Theorems

### B.1 Proof of Theorem 3.1 (CoT-Oracle Decidability Characterization)

*Proof.* We prove the two directions separately.

( $\Rightarrow$ ) **From a CoT-oracle machine to  $P^O[Q(n)]$ .** Suppose  $L$  is decidable by a CoT-oracle machine  $M_{\text{CoT}}^O$  in time  $T(n)$  using at most  $Q(n)$  queries to oracle  $O$ . We construct a deterministic oracle Turing machine  $M$  witnessing  $L \in P^O[Q(n)]$  as follows:

On input  $x$  of length  $n$ ,  $M$  simulates  $M_{\text{CoT}}^O$  step by step. Whenever  $M_{\text{CoT}}^O$  would invoke the chain-of-thought oracle on query  $q$ ,  $M$  instead issues  $q$  to  $O$  and records the reply. Since  $M_{\text{CoT}}^O$  runs in  $T(n)$  steps and makes at most  $Q(n)$  queries, the simulation by  $M$  also runs in time  $O(T(n))$  and issues no more than  $Q(n)$  queries. Thus  $M$  decides  $L$  within the resources defining  $P^O[Q(n)]$ .

( $\Leftarrow$ ) **From  $P^O[Q(n)]$  to a CoT-oracle machine.** Conversely, assume  $L \in P^O[Q(n)]$ . Then there exists a deterministic oracle Turing machine  $N$  that decides  $L$  in time  $T(n)$  with at most  $Q(n)$  queries to  $O$ . We define a CoT-oracle machine  $M_{\text{CoT}}^O$  that on input  $x$  simply runs  $N$  “as is,” forwarding each of  $N$ ’s queries to the chain-of-thought oracle and using its answers exactly as  $N$  would use  $O$ ’s. This simulation preserves both the time bound  $T(n)$  and the query bound  $Q(n)$ .

Having shown both directions, we conclude that a language  $L$  is decidable by a CoT-oracle machine in time  $T(n)$  with  $Q(n)$  queries if and only if  $L \in P^O[Q(n)]$ .  $\square$

## B.2 Proof of Theorem 5.1 (P Languages)

*Proof.* Let  $L \in P$ . By definition, there is a deterministic Turing machine  $M$  that decides  $L$  in time at most  $p(n)$  for some polynomial  $p$ . We will exhibit a CoT-oracle machine  $M_{\text{CoT}}^O$  with  $Q(n) = 0$  and  $T(n) = O(p(n))$ .

Define  $M_{\text{CoT}}^O$  to operate as follows on input  $x$  of length  $n$ :

1. Simulate the steps of  $M$  on input  $x$  exactly.
2. Never issue any query to the oracle  $O$ .
3. Whenever  $M$  enters an accept (resp. reject) state, have  $M_{\text{CoT}}^O$  accept (resp. reject).

Since  $M$  runs in time at most  $p(n)$ , the simulation by  $M_{\text{CoT}}^O$  runs in time  $T(n) = O(p(n))$ . Because  $M_{\text{CoT}}^O$  never invokes the oracle, its query complexity is  $Q(n) = 0$ . Hence  $M_{\text{CoT}}^O$  decides  $L$  within the stated resource bounds, completing the proof.  $\square$

## B.3 Proof of Theorem 5.2 (SAT via Prefix Queries)

*Proof.* Let  $\varphi$  be a Boolean formula on variables  $x_1, \dots, x_n$ . We describe a CoT-oracle machine  $M_{\text{CoT}}^O$  that decides SAT by building a satisfying assignment one bit at a time. Initialize a partial assignment  $a$  with no values. For each index  $i$  from 1 to  $n$ ,  $M_{\text{CoT}}^O$  poses to the oracle the question whether the formula obtained by fixing  $x_1 = a(x_1), \dots, x_{i-1} = a(x_{i-1})$  and  $x_i = 1$  is satisfiable. If the oracle replies “yes,” set  $a(x_i) \leftarrow 1$ ; otherwise set  $a(x_i) \leftarrow 0$ . After these  $n$  prefix-satisfaction queries,  $a$  is a complete assignment;  $M_{\text{CoT}}^O$  then evaluates  $\varphi(a)$  in  $O(n)$  time and accepts exactly if  $\varphi(a) = 1$ . Because exactly  $n$  queries are made,  $Q(n) = n$ , and the total non-query work (forming each query and the final evaluation) is  $O(n)$ , so  $T(n) = O(n)$ . This establishes the claimed bounds and completes the proof.  $\square$

## B.4 Proof of Theorem 5.3 (Linear Query Lower Bound for SAT)

*Proof.* Suppose, for contradiction, that there is a CoT-oracle machine  $M_{\text{CoT}}^O$  which decides SAT in time  $T(n) = O(n^k)$  while making only  $Q(n) = o(n)$  prefix-satisfaction queries. Fix an input formula  $\varphi$  on  $n$  variables chosen from any family for which SAT remains NP-hard even after arbitrarily fixing any  $o(n)$  of its variables (for instance, random 3-CNF formulas have this property).

Each yes/no prefix query “Is the formula satisfiable under these fixed assignments?” yields at most

one bit of information about any global satisfying assignment. After  $Q(n) = o(n)$  queries, at most  $o(n)$  bits of the true  $n$ -bit assignment have been determined, leaving  $\Omega(n)$  variables unfixed. At that point  $M_{\text{CoT}}^O$  must decide the satisfiability of the residual formula  $\varphi'$  on  $\Omega(n)$  unfixed variables using only its internal computation in time  $O(n^k)$ , with no further oracle access.

However, by assumption  $P \neq NP$ , no deterministic machine running in polynomial time can solve SAT on formulas with a linear number of unfixed variables. This contradicts the claimed resource bounds for  $M_{\text{CoT}}^O$ .

Therefore our original assumption was false, and any CoT-oracle machine deciding SAT in polynomial internal time must indeed make  $Q(n) = \Omega(n)$  queries.  $\square$

## C Extended Data and Examples

### C.1 Extended Pseudocode for SAT Reconstruction with Backtracking and Caching

**Algorithm 1:** Backtracking SAT Solver via CoT Queries with Caching

---

**Input** :  $\varphi$  on  $n$  variables  
**Global** : cache  $\mathcal{C} \leftarrow \emptyset$   
**Function**  $SolveSAT(\varphi, n)$   
  | **return** ExplorePrefix( $\varphi, \epsilon, n$ );  
**Function**  $ExplorePrefix(\varphi, prefix, n)$   
  | **if**  $|prefix| = n$  **then**  
  |   | **if**  $OracleQuery(\varphi, prefix) = YES$  **then**  
  |   |   | **return** prefix;  
  |   | **else**  
  |   |   | **return** FAIL;  
  | **for**  $b \leftarrow 0$  **to** 1 **do**  
  |   | newPref  $\leftarrow prefix \parallel b$ ;  
  |   | **if**  $\mathcal{C}[newPref]$  *undefined* **then**  
  |   |   |  $\mathcal{C}[newPref] \leftarrow$   
  |   |   |   |  $OracleQuery(\varphi, newPref)$ ;  
  |   | **if**  $\mathcal{C}[newPref] = YES$  **then**  
  |   |   |  $r \leftarrow$   
  |   |   |   |  $ExplorePrefix(\varphi, newPref, n)$ ;  
  |   |   |   | **if**  $r \neq FAIL$  **then**  
  |   |   |   |   | **return**  $r$ ;  
  | **return** FAIL;  
**Function**  $OracleQuery(\varphi, pref)$   
  | **if**  $pref \in \mathcal{C}$  **then**  
  |   | **return**  $\mathcal{C}[pref]$ ;  
  | construct  $\varphi_{pref}$  by fixing variables in  
  |   | pref;  
  |  $ans \leftarrow oracle(\varphi_{pref})$ ;  
  |  $\mathcal{C}[pref] \leftarrow ans$ ;  
  | **return**  $ans$ ;

---

### C.2 Extended Simulation Results

Table 8 presents mean and standard deviation of  $T$  over 100 trials for each  $(n, Q_{\max})$  pair.

**Analysis of Extended Simulation Results (Table 8)** The simulations confirm that increasing the query budget  $Q_{\max}$  dramatically reduces both the expected internal time and its variability, and that this effect becomes more pronounced as problem size  $n$  grows. For  $n = 50$ , raising  $Q_{\max}$  from

Table 8: Simulated  $T(n, Q_{\max})$  over 100 trials (CPU ops)

$n$	$Q_{\max}$	$\mathbb{E}[T]$	$\text{Std}[T]$
50	5	$1.2 \times 10^3$	$1.5 \times 10^2$
50	25	$5.0 \times 10^2$	$8.0 \times 10^1$
100	10	$4.8 \times 10^3$	$4.2 \times 10^2$
100	100	$9.6 \times 10^2$	$1.1 \times 10^2$
150	15	$1.3 \times 10^4$	$9.8 \times 10^2$
150	150	$1.2 \times 10^3$	$2.0 \times 10^2$
200	20	$2.2 \times 10^4$	$1.5 \times 10^3$
200	200	$1.6 \times 10^3$	$2.5 \times 10^2$

5 to 25 cuts  $\mathbb{E}[T]$  by nearly 60% (from  $1.2 \times 10^3$  to  $5.0 \times 10^2$  CPU ops) and reduces  $\text{Std}[T]$  by almost half. At  $n = 100$ , a tenfold increase in  $Q_{\max}$  (from 10 to 100) yields an 80% reduction in mean time (from  $4.8 \times 10^3$  to  $9.6 \times 10^2$ ) and similarly shrinks the standard deviation. The same pattern holds for  $n = 150$  and  $n = 200$ : when  $Q_{\max} = n$ , the solver’s cost drops by an order of magnitude compared to a small constant budget, and the run-time variability stabilizes. Overall, these results illustrate that larger query allowances effectively shift the computational burden from brute-force search to oracle guidance, achieving near-linear scaling in  $n$  when  $Q_{\max}$  grows proportionally with problem size.

### C.3 Additional Worked Example: Multi-Hop QA

Consider the three-hop question “What is the population of the city where the author of *Pride and Prejudice* was born?”

Table 9: Multi-Hop QA Prompts

$p_1$ : “Who wrote ‘Pride and Prejudice’?”  
 $p_2$ : “Where was Jane Austen born?”  
 $p_3$ : “What is the population of Steventon, Hampshire?”

Assuming each oracle query executes in  $O(1)$  time, the machine issues exactly three queries, so

$$Q(n) = 3 \quad (5)$$

and the total internal reasoning cost is

$$T(n) = \sum_{i=1}^3 O(1) = O(1) \quad (6)$$

In practice one may add a final prompt

Table 10: Consolidation Prompt

---

$p_4$  : Using the above, please state the population of Stevenston.

---

to consolidate the answers, this increases  $Q(n)$  to 4 but preserves  $T(n) = O(1)$ . More generally, a  $k$ -hop question with fixed  $k$  satisfies  $Q(n) = k$  and  $T(n) = O(1)$ , illustrating that chain-of-thought reasoning over structured knowledge can achieve constant-time performance when the number of hops does not scale with input size.

## D Chain-of-Thought Strategy Examples

### D.1 Arithmetic Chain-of-Thought: Integer Factorization

Consider the problem of factoring an  $n$ -bit integer  $N$  using a CoT-oracle machine with access to a primality oracle  $O_{\text{PRIME}}$ . One can issue queries “Is there a nontrivial factor of  $N$  less than  $2^i$ ?” for  $i = 1, \dots, n$ ; once the first affirmative response occurs at index  $i^*$ , a binary search over  $[2^{i^*-1}, 2^{i^*})$  using  $O_{\text{PRIME}}$  in  $O(\log n)$  additional queries yields a prime factor  $p$ , and the process repeats on the quotient. This approach uses  $Q(n) = O(\log n + \log n) = O(\log n)$  queries per factor, while each iteration’s internal work (long division and binary-search bookkeeping) costs  $O(n^2)$  time, so factoring a two-factor semiprime overall requires  $Q(n) = O(\log n)$  and  $T(n) = O(n^3)$ .

### D.2 Logical Deduction Chains: SAT Subcalls

For a Boolean formula  $\varphi$  on  $n$  variables, the CoT strategy reconstructs a satisfying assignment by invoking SAT oracle subcalls: for each index  $i$ , one queries whether there exists a satisfying extension with  $x_1, \dots, x_{i-1}$  fixed and  $x_i = 0$  (and similarly for  $x_i = 1$ ). In the worst case this requires  $Q(n) = 2n$  queries, though the prefix technique reduces this to  $Q(n) = n + 1$ . Each internal verification of a partial assignment takes  $O(n)$  time, yielding  $T(n) = O(n^2)$  overall. Hence the complexity profile for SAT reconstruction is  $(T(n), Q(n)) = (O(n^2), O(n))$ , matching the bounds of classical search.

### D.3 Complexity Profiles

Problem	$Q(n)$	$T(n)$
Integer Factorization	$O(\log n)$	$O(n^3)$
SAT Reconstruction	$O(n)$	$O(n^2)$

## E Related Work

### E.1 Oracle Turing Machines in Complexity Theory

The study of oracle Turing machines originates from Cook’s seminal work on Turing reductions and NP-completeness (Cook, 1971). Goldwasser, Micali, and Rackoff introduced interactive proofs, showing how oracle-like interactions characterize randomized classes (Goldwasser et al., 1985). Babai and Moran formalized Arthur–Merlin games, bridging randomness and interaction in complexity theory (Babai and Moran, 1988), and Shamir’s result  $\text{IP} = \text{PSPACE}$  further underscored the power of interactive oracle calls (Shamir, 1992). Arora and Barak’s comprehensive treatment lays out the modern framework for bounded-query classes  $\text{P}^O[k]$  (Arora and Barak, 2009).

### E.2 Prior Chain-of-Thought Analyses

Wei et al. first empirically established CoT prompting’s impact on reasoning benchmarks (Wei et al., 2022). Kojima et al. demonstrated zero-shot CoT without exemplars (Kojima et al., 2022), and Wang et al. improved robustness via self-consistency (Wang et al., 2022). Zhou et al. proposed least-to-most decomposition, while Yao et al. expanded to tree-based exploration of reasoning states (Zhou et al., 2022; Yao et al., 2023). Nye et al. introduced scratchpads, akin to subprompt buffers, to capture intermediate computation (Nye et al., 2021). These works underpin our formal model by highlighting practical CoT mechanisms.

### E.3 Interactive Proof Systems and Reductions

Interactive proofs formalize multi-round oracle interactions. The Arthur–Merlin hierarchy relates to bounded alternating queries, and Fortnow and Santhanam analyzed complexity consequences of restricted interactions (Fortnow and Santhanam, 2008). Our bounded-query CoT model parallels these systems, recasting prompt-based interactions as oracle queries within classical complexity hierarchies.