

# PENCIL: Long Thoughts with Short Memory

Chenxiao Yang

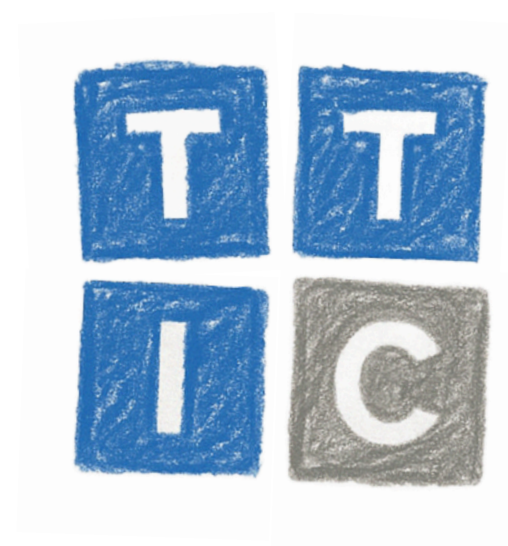
@ASAP

May 28th, 2025

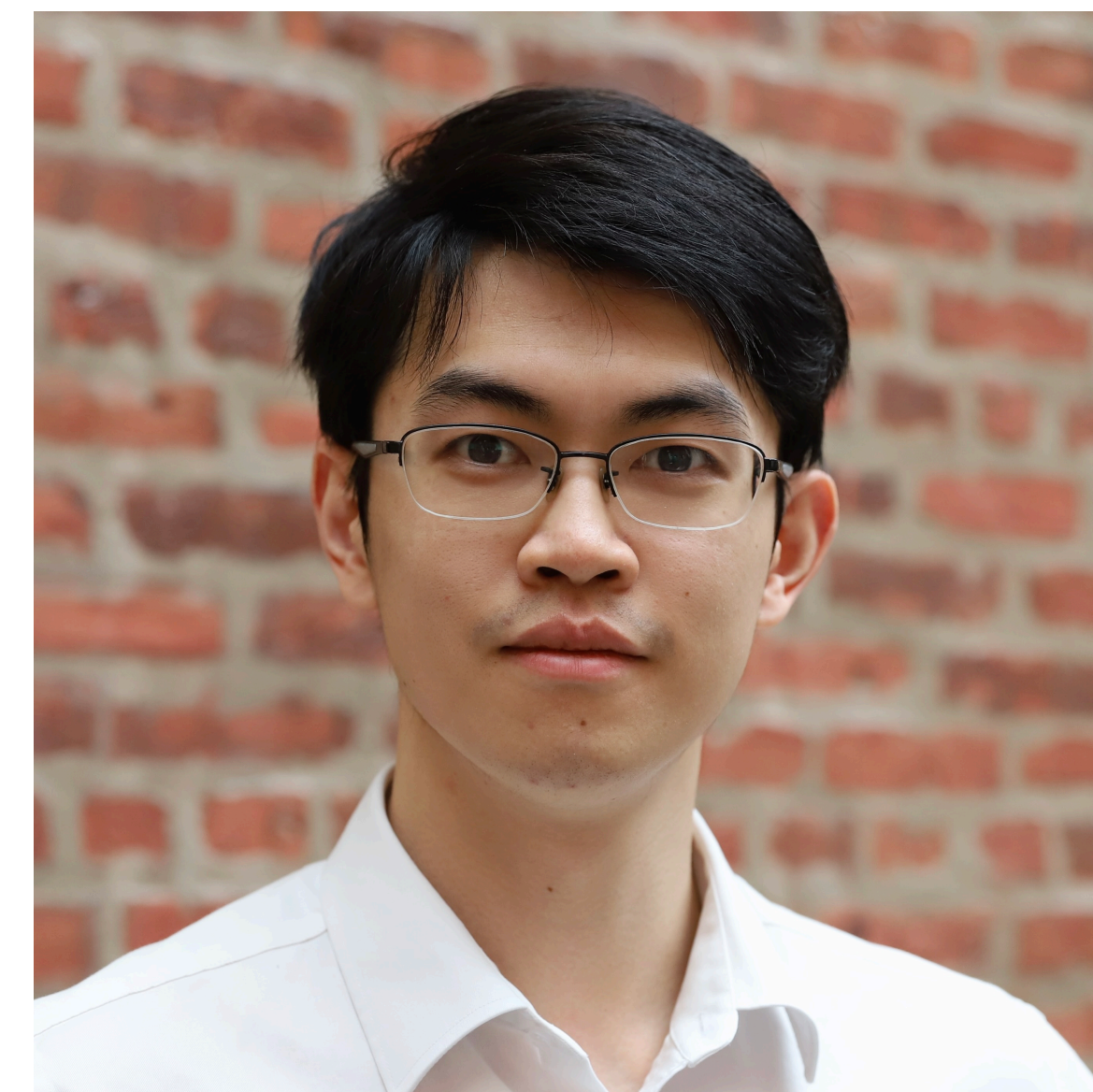
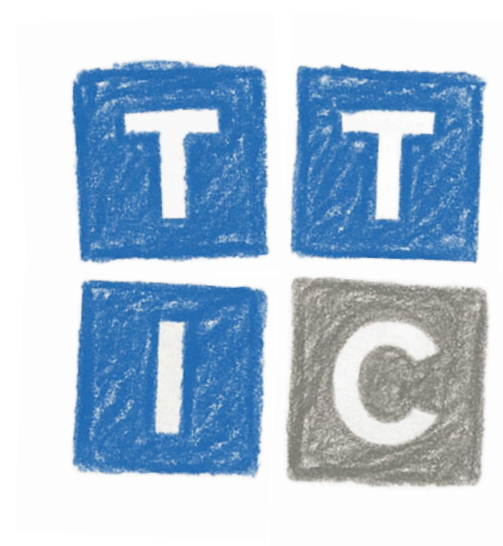




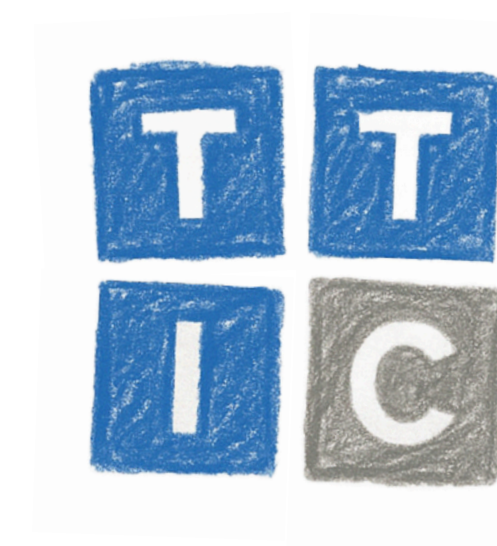
Nathan Srebro



David McAllester



Zhiyuan Li





# Chain-of-Thought (CoT)

## CommonsenseQA

Q: Sammy wanted to go to where the people were. Where might he go?

Options: (a) race track (b) populated areas (c) desert (d) apartment (e) roadblock

A: The answer must be a place with a lot of people. Race tracks, desert, apartments, and roadblocks don't have a lot of people, but populated areas do. So the answer is (b).

## StrategyQA

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about  $0.6 \text{ g/cm}^3$ , which is less than water. Thus, a pear would float. So the answer is no.

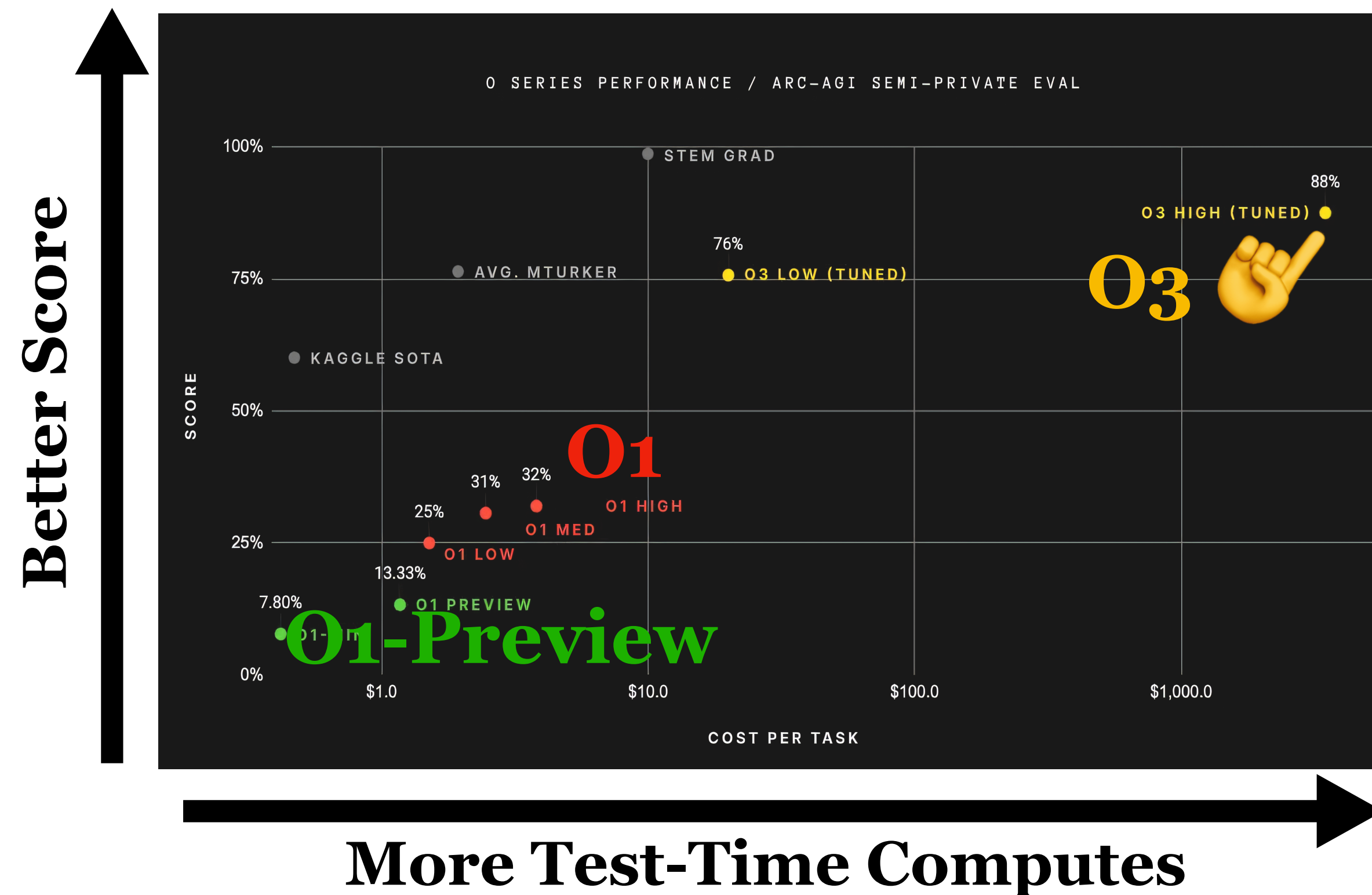
## Date Understanding

Q: The concert was scheduled to be on 06/01/1943, but was delayed by one day to today. What is the date 10 days ago in MM/DD/YYYY?

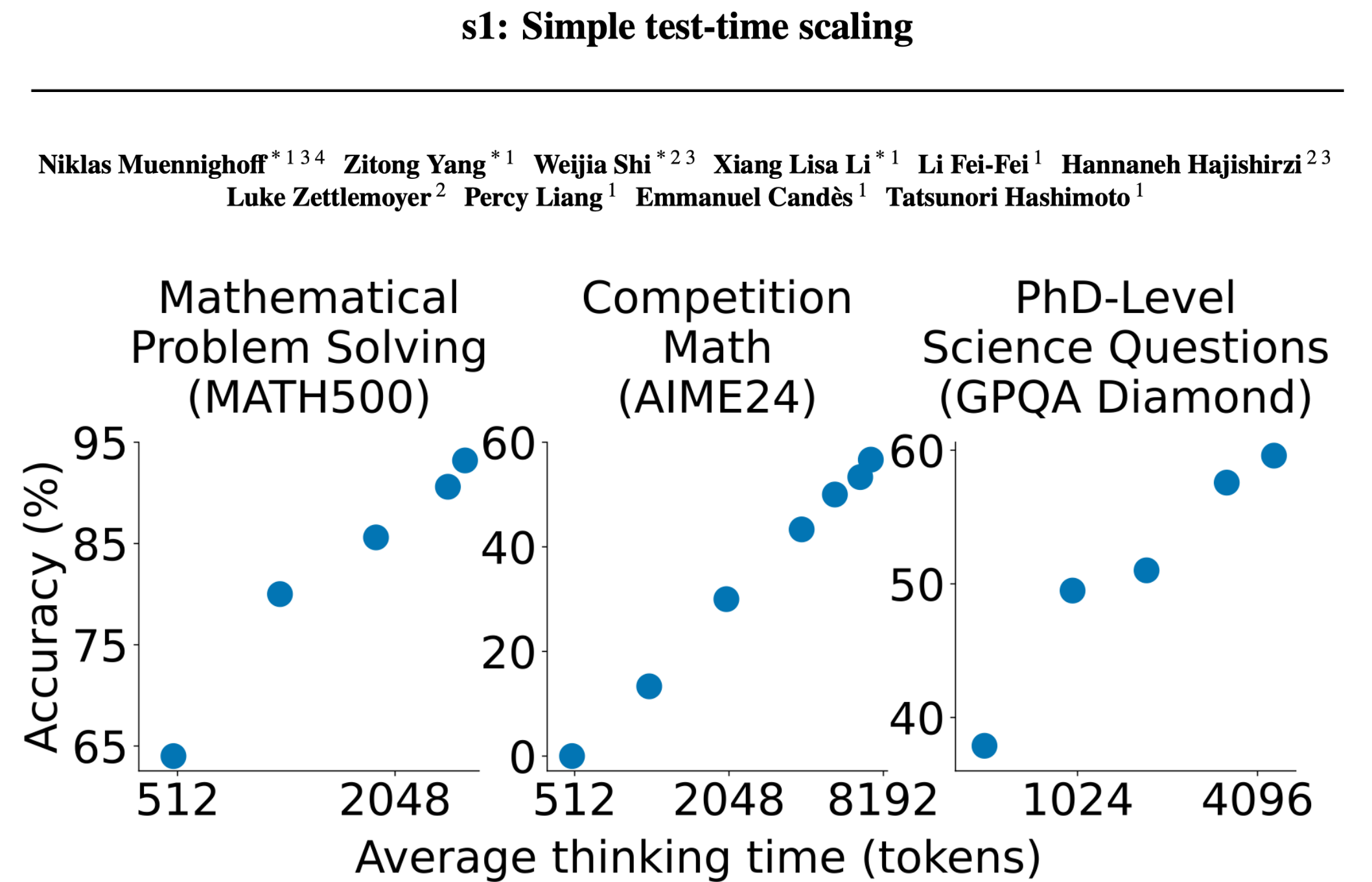
A: One day after 06/01/1943 is 06/02/1943, so today is 06/02/1943. 10 days before today is 05/23/1943. So the answer is 05/23/1943.

Chain-of-Thought generates a series of **thoughts** before providing the final answer.

# The Power of Long CoT



*Performance of O3 on ARC-AGI-Pub*

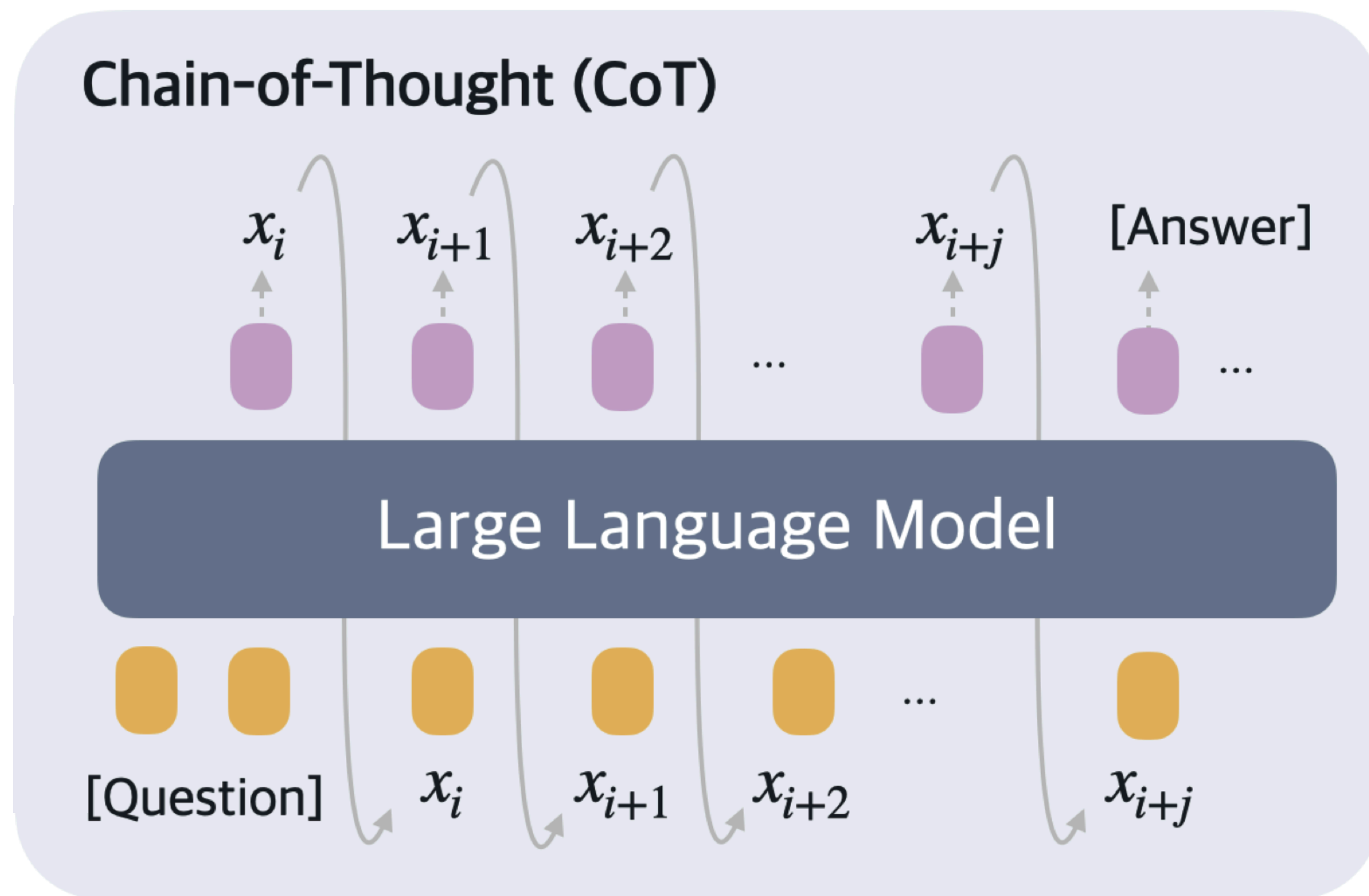


*S1, Muennighoff et al., 2025*

Long CoT empowers transformers to solve computationally more intensive tasks (Feng et al. 2024, Merrill et al. 2024, Li et al. 2024, etc.)



# Long CoT is great! BUT ...



*Image adapted from Hao et al., 2024*

- **Out-of-Memory:** Accumulated thoughts will exceed model's effective context windows.
- **Large Cost:** Self-attention makes generating each new token prohibitively expensive.
- **Low Accuracy:** Relevant information becomes buried and hard to retrieve (Liu et al. 2024, Ye et al. 2025, etc.)

# Most Intermediate Steps are “Erased”

## • Math

### **Theorem.**

There are infinitely many prime numbers.

### **Lemma.**

Every integer  $n > 1$  has a prime  $p$  divisor.

**Proof of Lemma.** Take the smallest divisor  $p > 1$  of  $n$ . If  $p$  were composite, it would factor as  $ab$  with  $1 < a, b < p$ , contradicting minimality. Hence  $p$  is prime and divides  $n$ . ■

No longer useful

**Proof of Theorem.** Assume, for contradiction, that the complete list of primes is  $p_1, \dots, p_k$ . Set  $N = p_1 p_2 \dots p_k + 1$ . **By the lemma**,  $N$  has some prime divisor  $q$ . But every  $p_i$  leaves remainder 1 when dividing  $N$ , so none of them divides  $N$ ; hence  $q$  is a new prime, contradicting completeness. Therefore primes are infinite. ■

## • LLM Reasoning (from DeepSeek-R1 run)

**Question:** “Your mother has 4 children: North, South, East—what’s the 4th child’s name?”

Alright, since the three children already named—North, South, and East—correspond precisely to three of the four cardinal directions, it seems perfectly logical to complete the compass pattern by concluding that the fourth child must therefore be named **West**.

No longer useful

**Wait** — maybe it’s a trick: “your mother” means \*you\* are the fourth child.

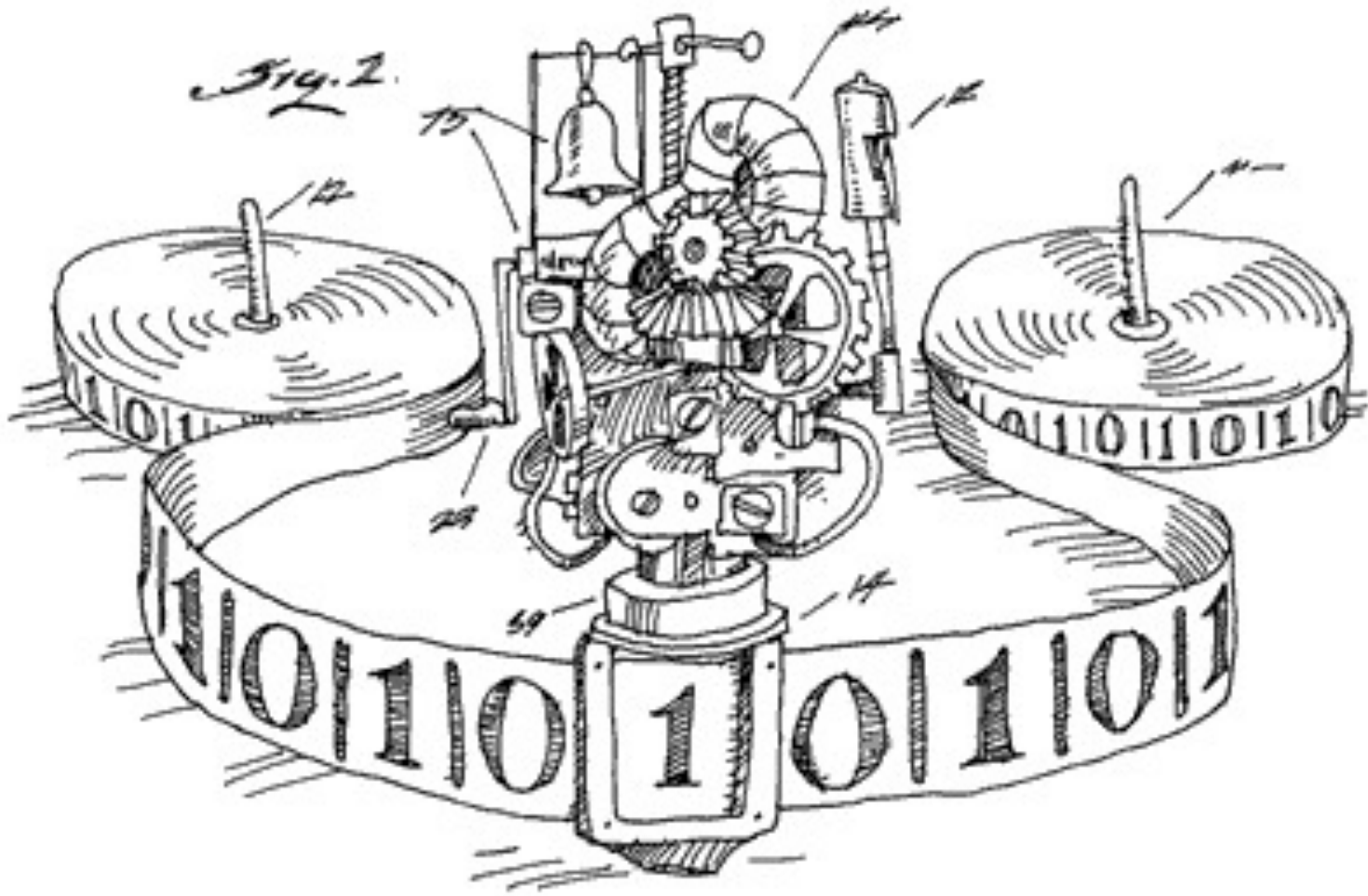
**Wait** — the question wants a \*name\*, so the answer is actually \*your own name\*.

→ **Final answer:** your name.



# “Erasure” is Fundamental to Computation

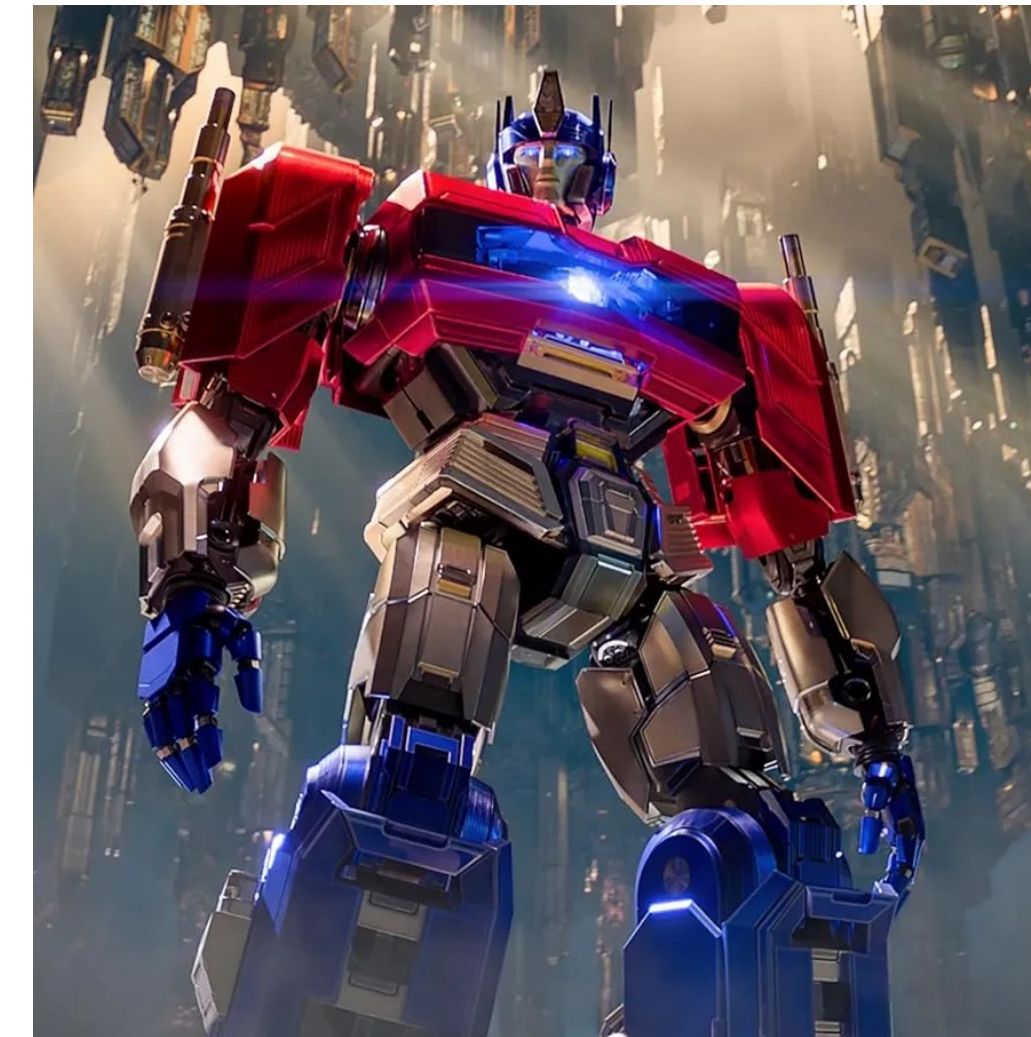
One Step of Thinking  $\Leftrightarrow$  One Step of Computation



Turing Machine (1936)



Modern Computer



Transformer + **PENCIL**



# PENCIL

## PENCIL ENables Context-efficient Inference and Learning

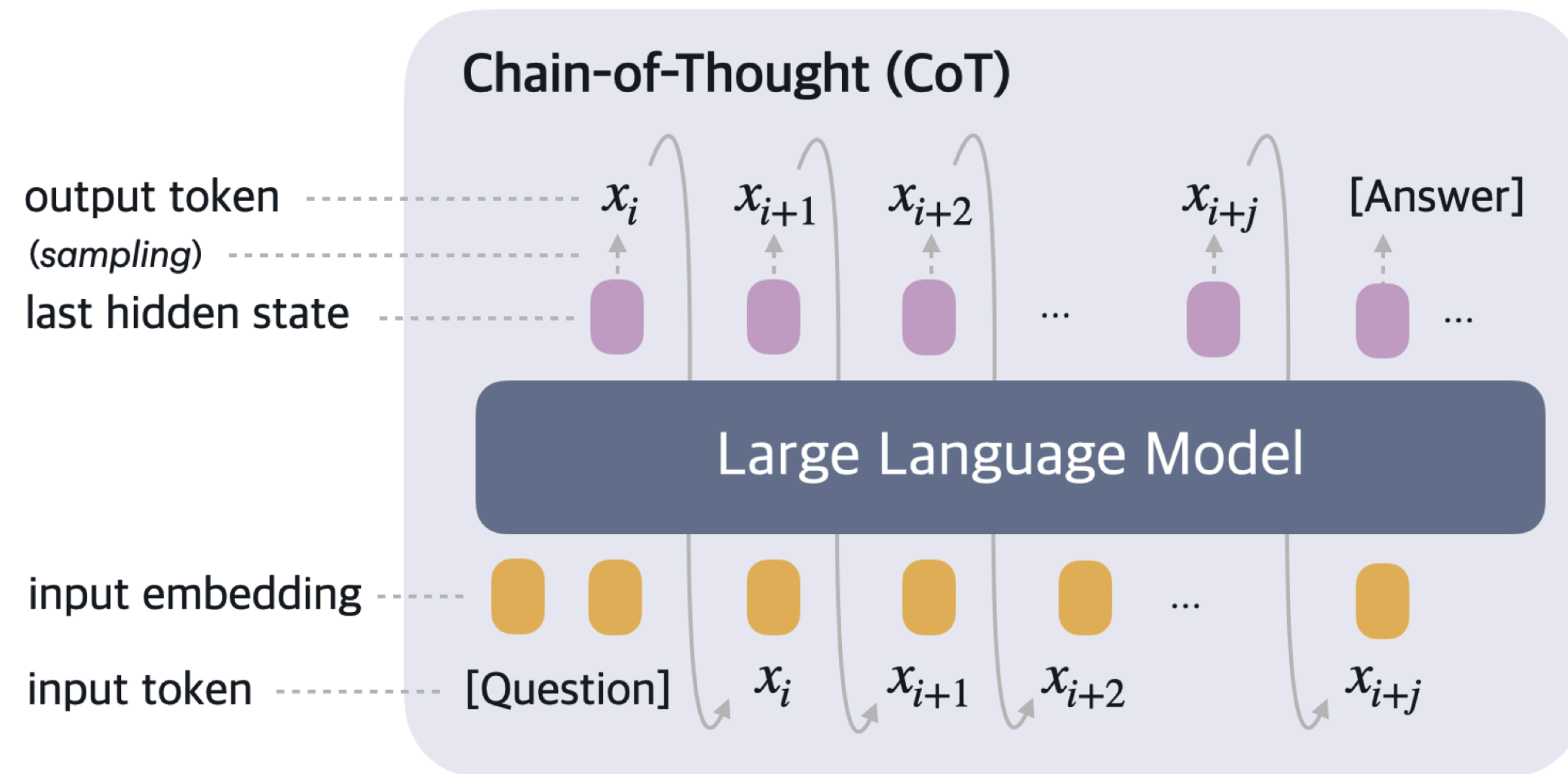




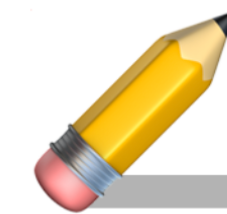
# What is PENCIL?



## Model Generation (Write)



Autoregressive Next-token Generation



## Reduction Rule (Erase)

```
examples.lisp
;; Common Lisp examples.

(defun hello-world ()
  "Print 'hello, world' message."
  (format t "hello, world~%"))

(defun factorial (n)
  "Compute factorial of n."
  (if (= n 0)
      1
      (* n (factorial (- n 1)))))

(defun fibonacci (n)
  "Compute nth Fibonacci number."
  (if (< n 2)
      n
      (+ (fibonacci (- n 1)) (fibonacci (- n 2)))))

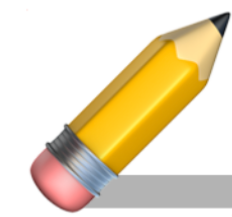
; SLIME 2.26
CL-USER> (hello-world)
hello, world
NIL
CL-USER> (factorial 6)
720
CL-USER> (trace factorial)
(FACTORIAL)
CL-USER> (factorial 6)
0: (FACTORIAL 6)
1: (FACTORIAL 5)
2: (FACTORIAL 4)
3: (FACTORIAL 3)
4: (FACTORIAL 2)
5: (FACTORIAL 1)
6: (FACTORIAL 0)
6: FACTORIAL returned 1
5: FACTORIAL returned 1
4: FACTORIAL returned 2
3: FACTORIAL returned 6
2: FACTORIAL returned 24
1: FACTORIAL returned 120
0: FACTORIAL returned 720
720
CL-USER> (fibonacci 6)
8
CL-USER>
```

Functional Programming

# What is PENCIL?

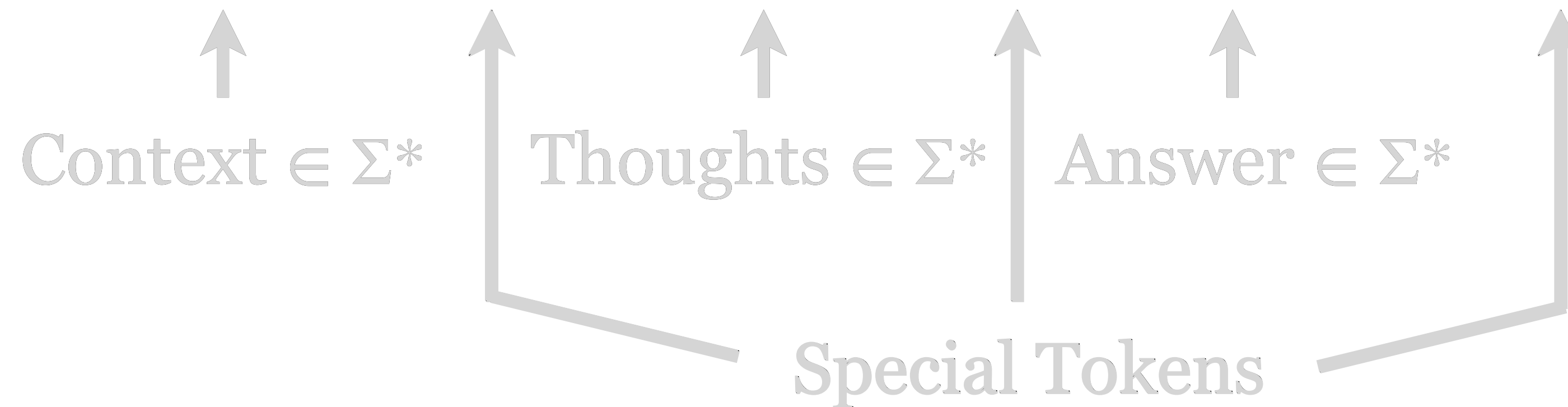




Model Generation (Write)



Reduction Rule (Erase)

**C** [CALL] **T** [SEP] **A** [RETURN]  $\Rightarrow$  **C A**



- Reduction is triggered when the sequence matches the pattern
- PENCIL iteratively  generate thoughts and  triggers reduction.



# Example 1: Arithmetic Expression Evaluation

A toy store put together party bags for a birthday. They made 3 blue bags with 5 toys each and 2 red bags with 4 toys each. How many toys were used in total? **There were 23 toys used in total.** own into parts! **There were 15 toys in all blue bags.** ~~There were 8 toys in all red bags. They made 3 bags and they were toys in each bag so they had 15 and 8 toys.~~ **[SEP]** ~~There were 23 toys in all blue bags.~~ **[SEP]** There were 23 toys used in total. **[RETURN]** 8 toys in all red bags. **[RETURN]**

**C** [CALL] **T** [SEP] **A** [RETURN]  $\Rightarrow$  **C A**

# Example 1: Arithmetic Expression Evaluation

**Prompt:** A toy store put together party bags for a birthday. They made 3 blue bags with 5 toys each and 2 red bags with 4 toys each. How many toys were used in total?

**PENCIL**



A toy store put together party bags for a birthday. They made 3 blue bags with 5 toys each and 2 red bags with 4 toys each. How many toys were used in total?



# Example 1: Arithmetic Expression Evaluation

**Prompt :** A toy store put together party bags for a birthday. They made 3 blue bags with 5 toys each and 2 red bags with 4 toys each.  
How many toys were used in total?

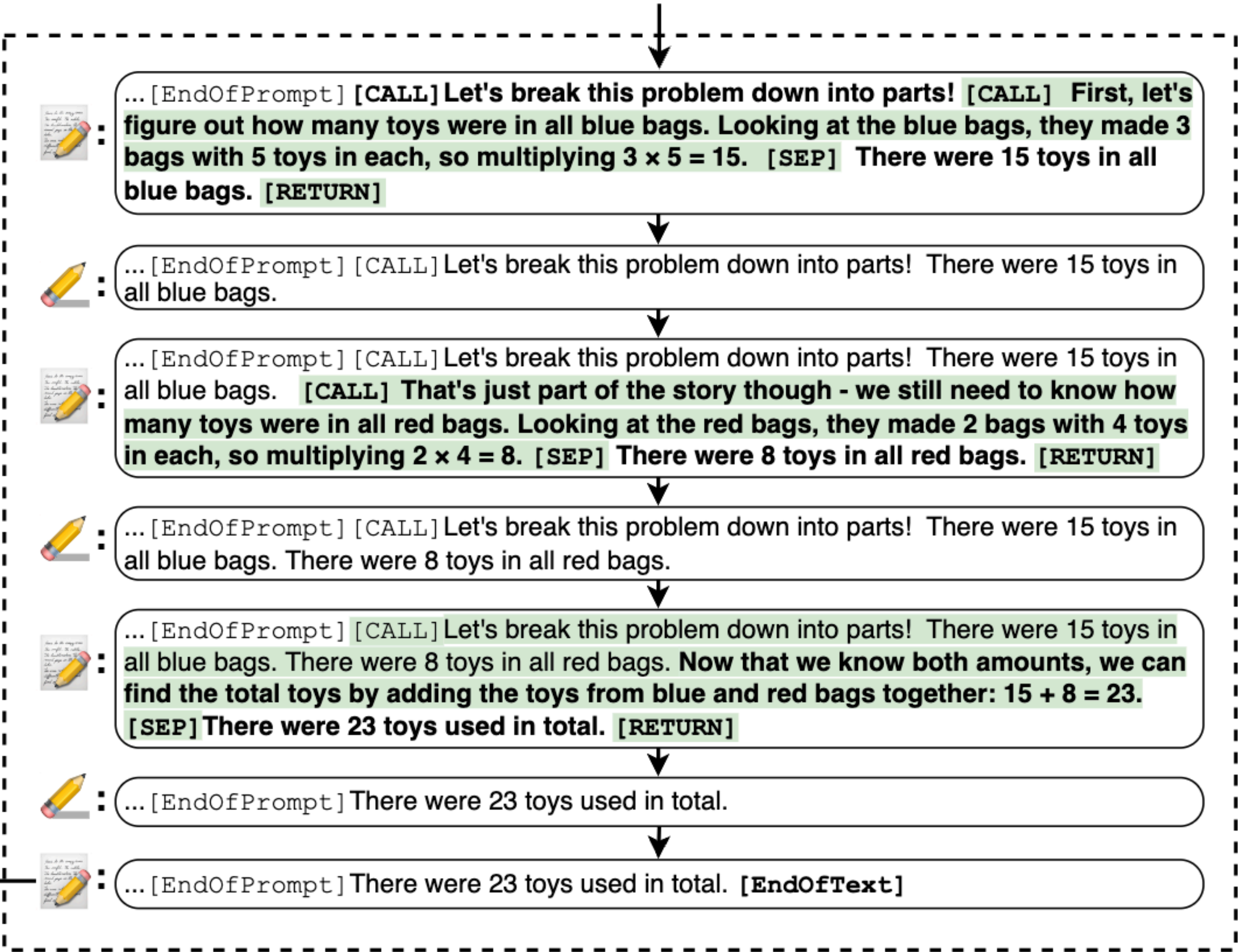
## Chain-of-Thought

**Response :** Let's break this problem down into parts! First, let's figure out how many toys were in all blue bags. Looking at the blue bags, they made 3 bags with 5 toys in each, so multiplying  $3 \times 5 = 15$ . There were 15 toys in all blue bags. That's just part of the story though - we still need to know how many toys were in all red bags. Looking at the red bags, they made 2 bags with 4 toys in each, so multiplying  $2 \times 4 = 8$ . There were 8 toys in all red bags. Now that we know both amounts, we can find the total toys by adding the toys from blue and red bags together:  $15 + 8 = 23$ . There were 23 toys used in total.

## PENCIL

**Response :** There were 23 toys used in total.

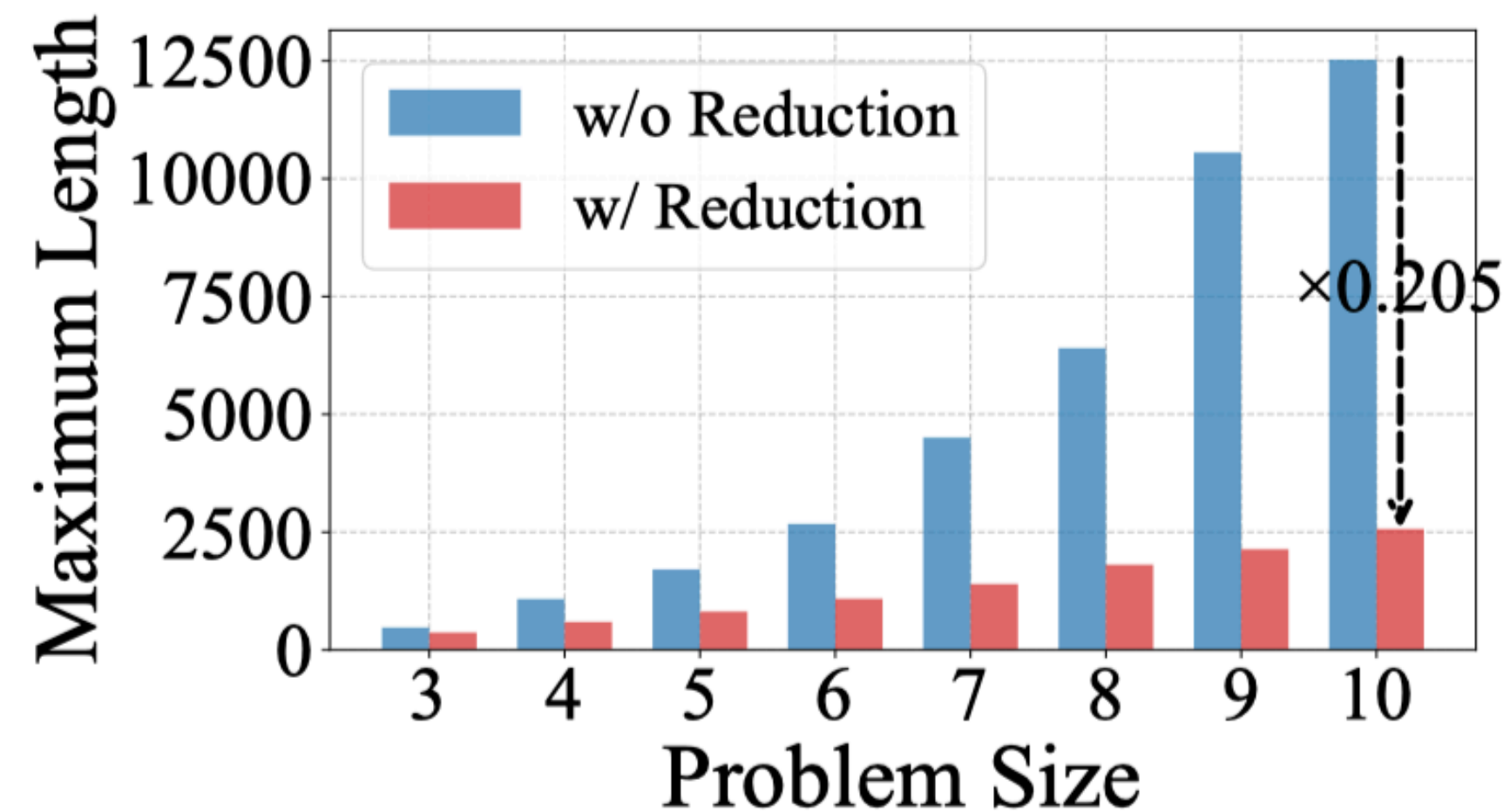
Hidden CoT 



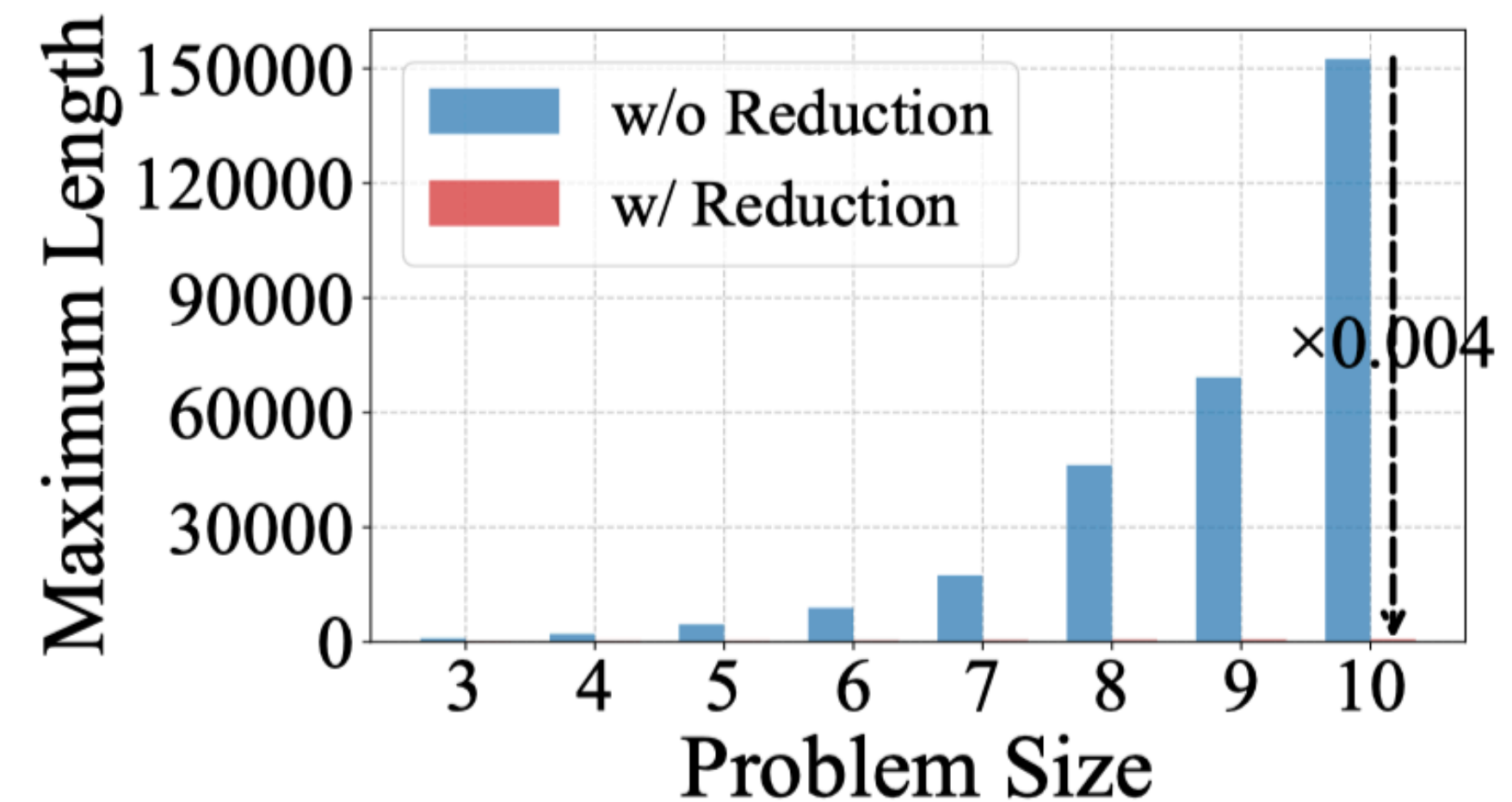




# Max Sequence Length Comparison (CoT v.s. PENCIL)



**3-SAT**  
(a special case of QBF)



**QBF**

**PENCIL** significantly reduces the maximal context length during inference  
( $\times 0.004$  when  $n=10$  on QBF)



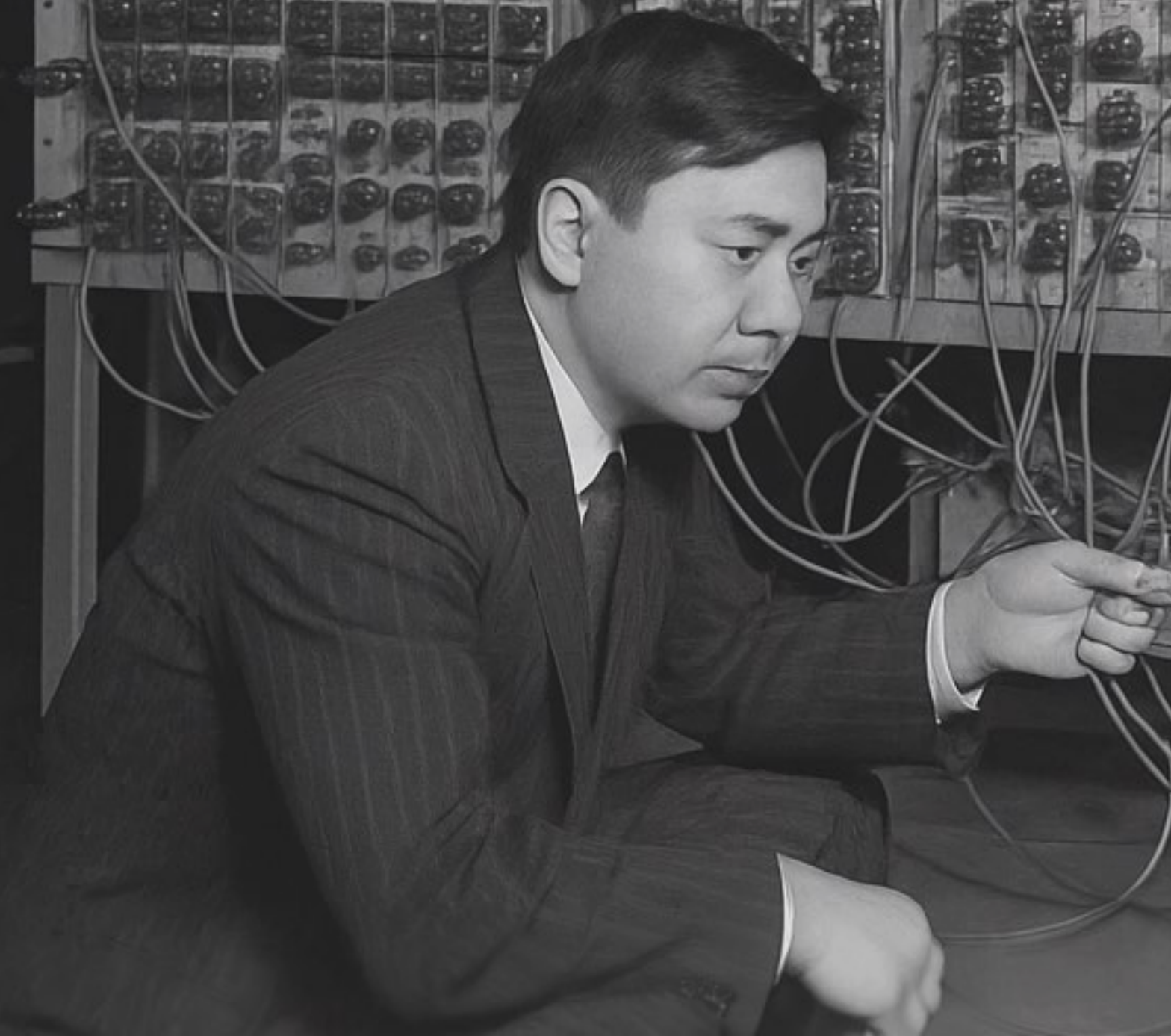
[illegible]



# PENCIL 2025

Long Thoughts with Short Memory

How does **PENCIL** perform?





# Experimental Setting: Training and Inference

**Training:** Key difference between CoT and PENCIL  $\Rightarrow$  **Data Generation\***

$$\mathcal{L}_{\text{CoT}} = - \sum \log p(\text{next token} \mid \text{complete sequence})$$

$$\mathcal{L}_{\text{PENCIL}} = - \sum \log p(\text{next token} \mid \text{reduced sequence})$$

We train a small transformer (25M parameter, 2048 context length) from scratch.

**Inference:**                    **C** [CALL] T [SEP] **A** [RETURN]  $\Rightarrow$  **C A**

Preserve the KV cache of **Context (C)** and recompute that for **Answer (A)**

*\* Datasets are generated by running specialized algorithms*



# Performance Comparison on 3-SAT and QBF

n =	3	4	5	6	7	8	9	10
Baseline	66	57	46	51	46	51	49	51
CoT	100	100	100	99	84	63	54	50
PENCIL	100	100	100	99	99	100	100	100

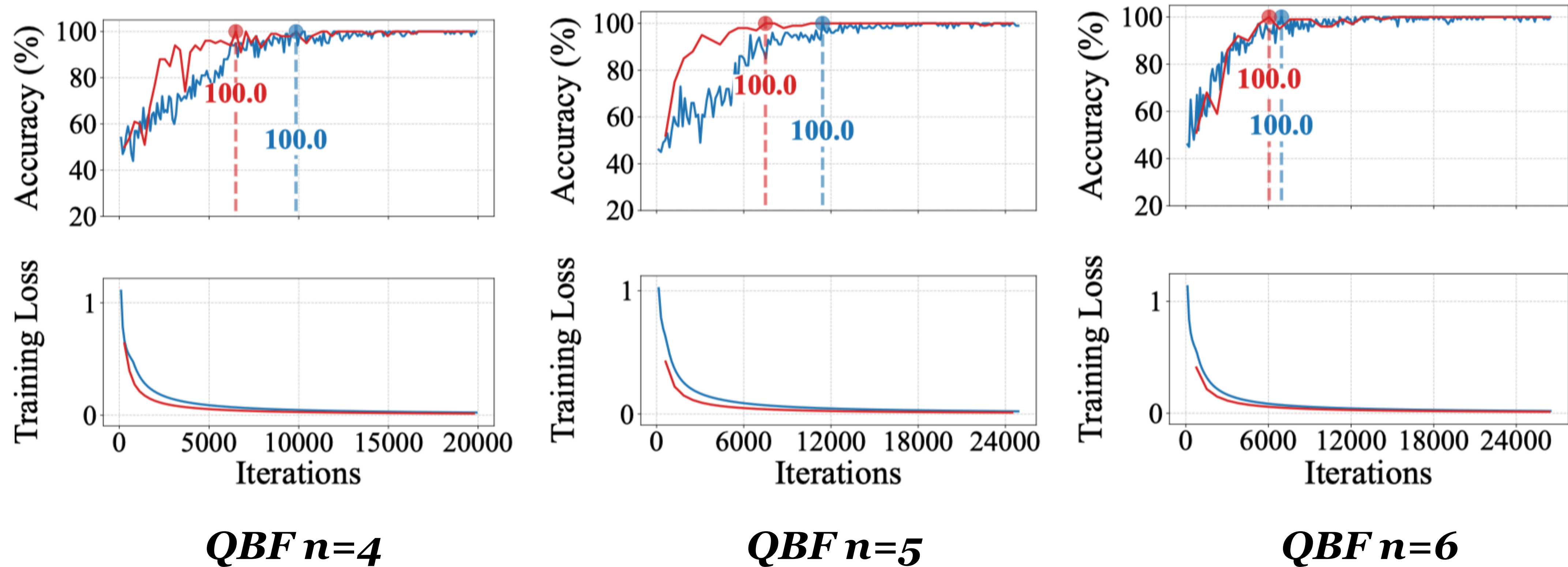
***3-SAT***

n =	3	4	5	6	7	8	9	10
Baseline	90	82	85	68	60	69	71	66
CoT	100	100	97	94	74	72	69	73
PENCIL	100	100	100	100	100	100	100	100

***QBF***

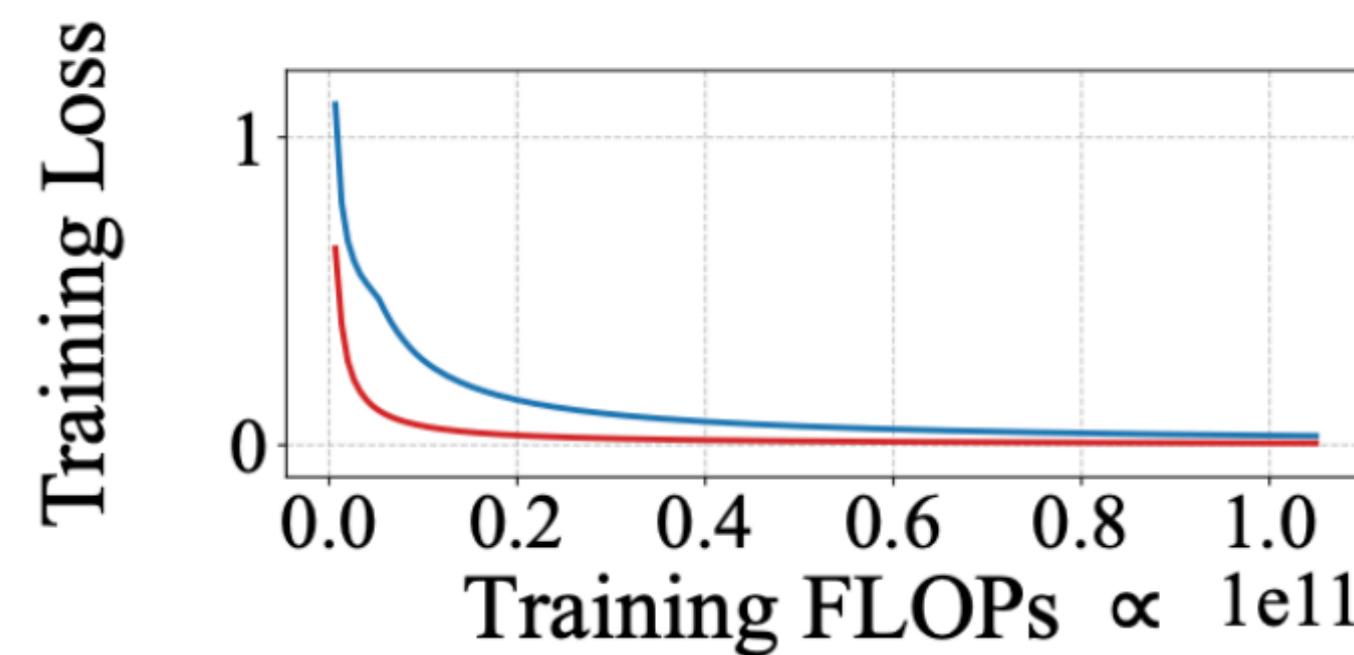
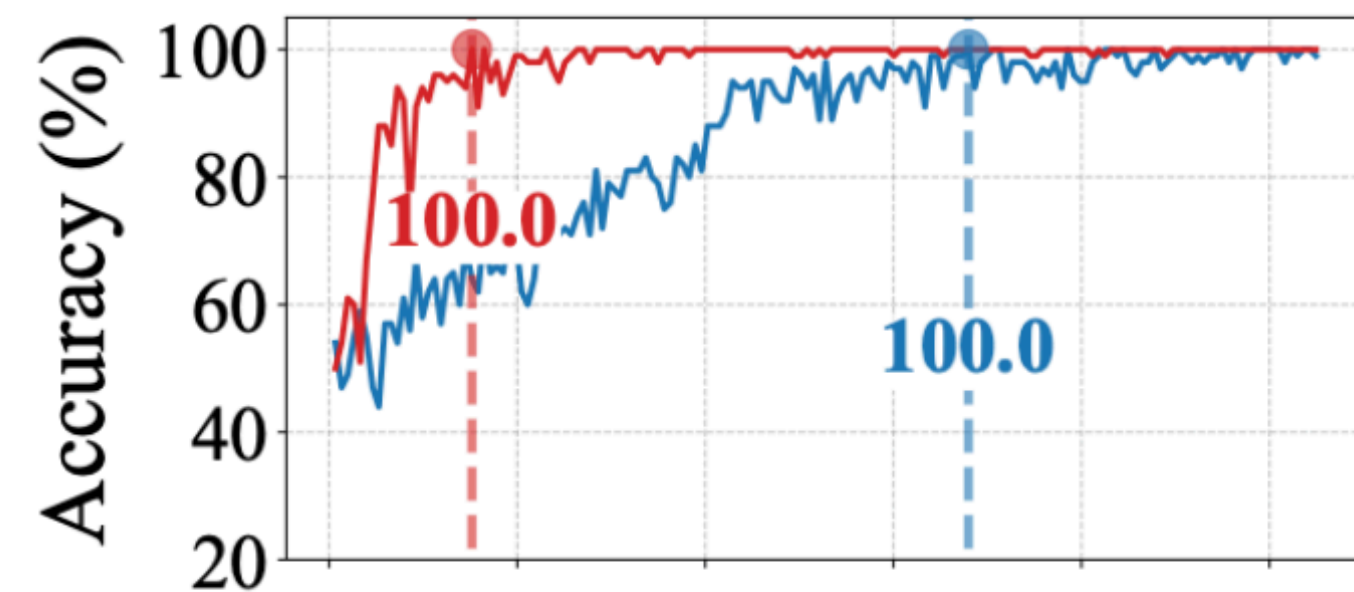
 **PENCIL** significantly outperforms  **CoT** on NP-hard tasks SAT and QBF (i.e. almost perfect v.s. random guessing).

# Convergence Speed (**CoT** v.s. **PENCIL**)

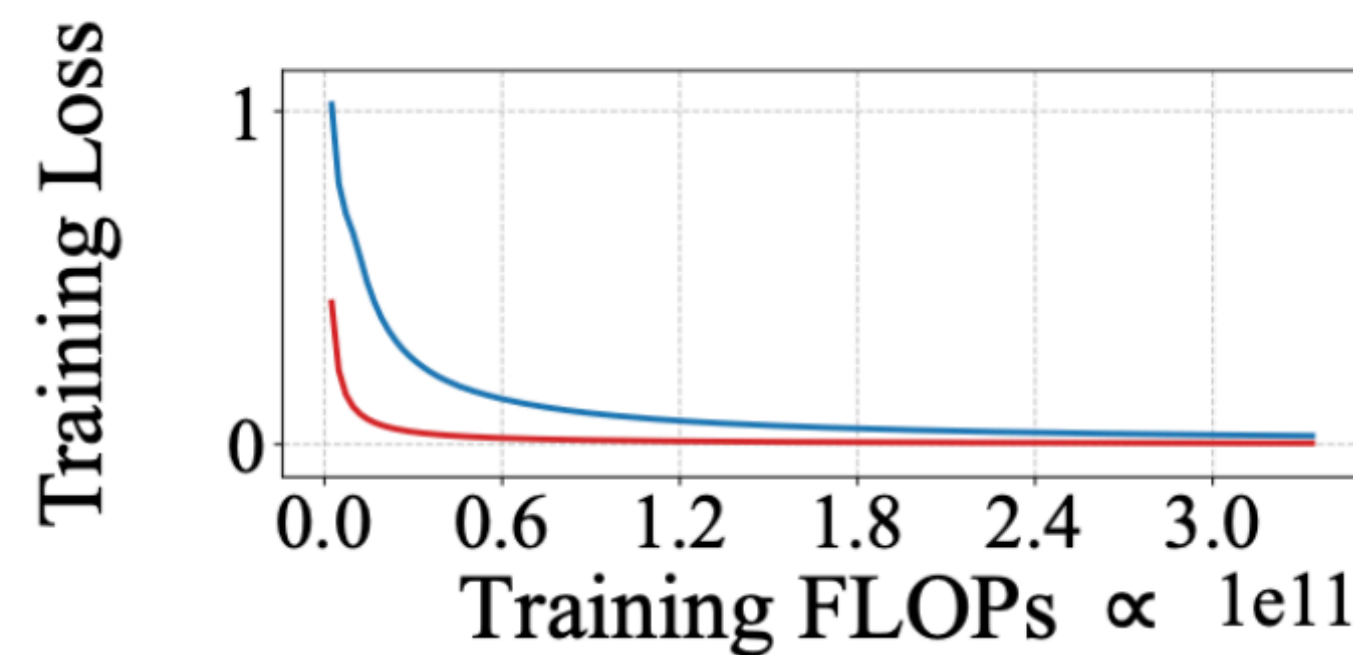
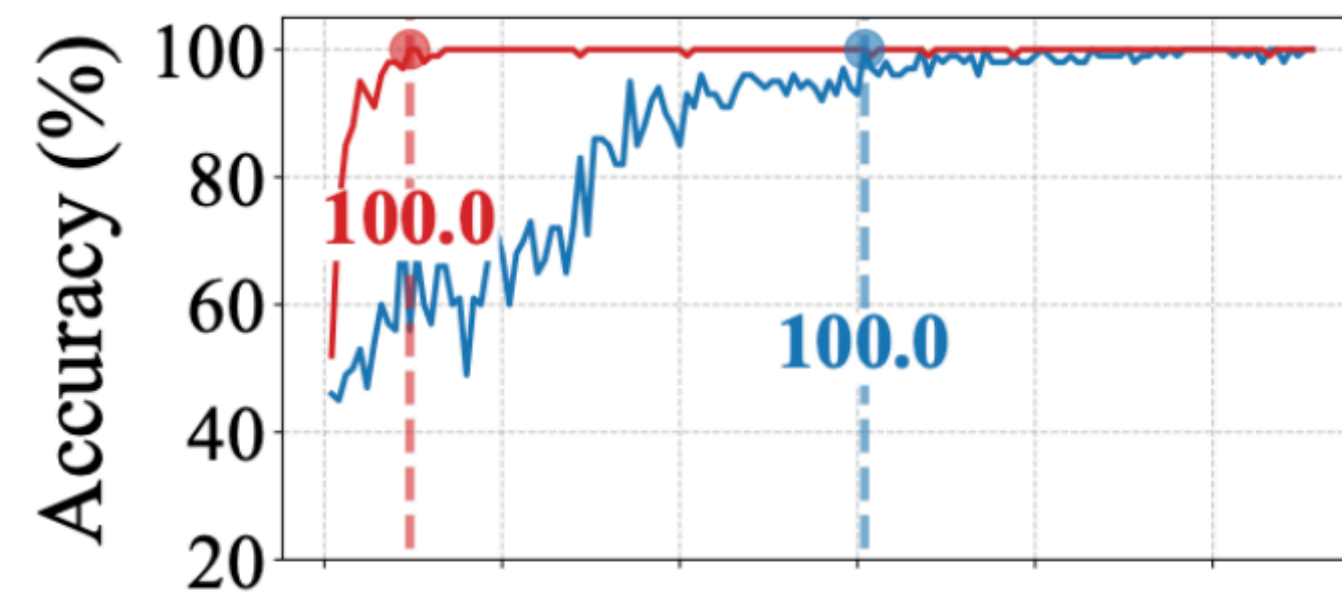


 **PENCIL** converges faster and is more sample efficient.

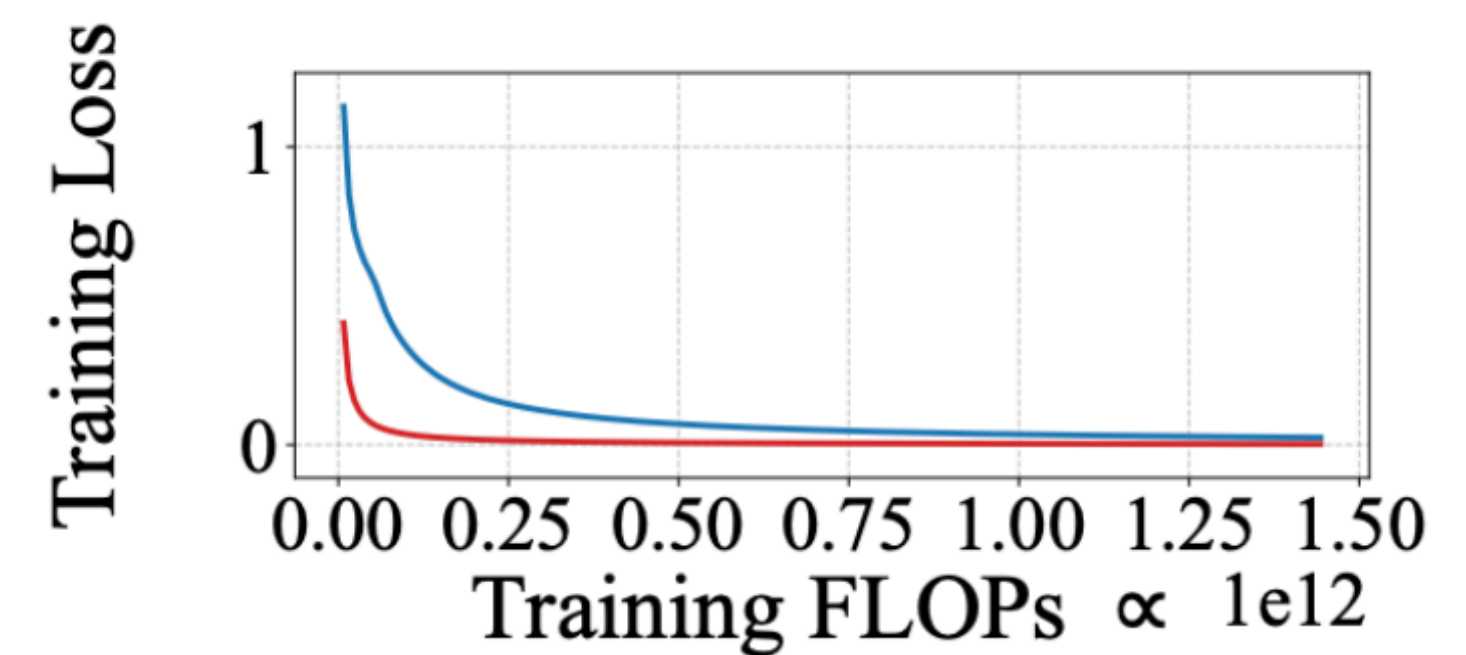
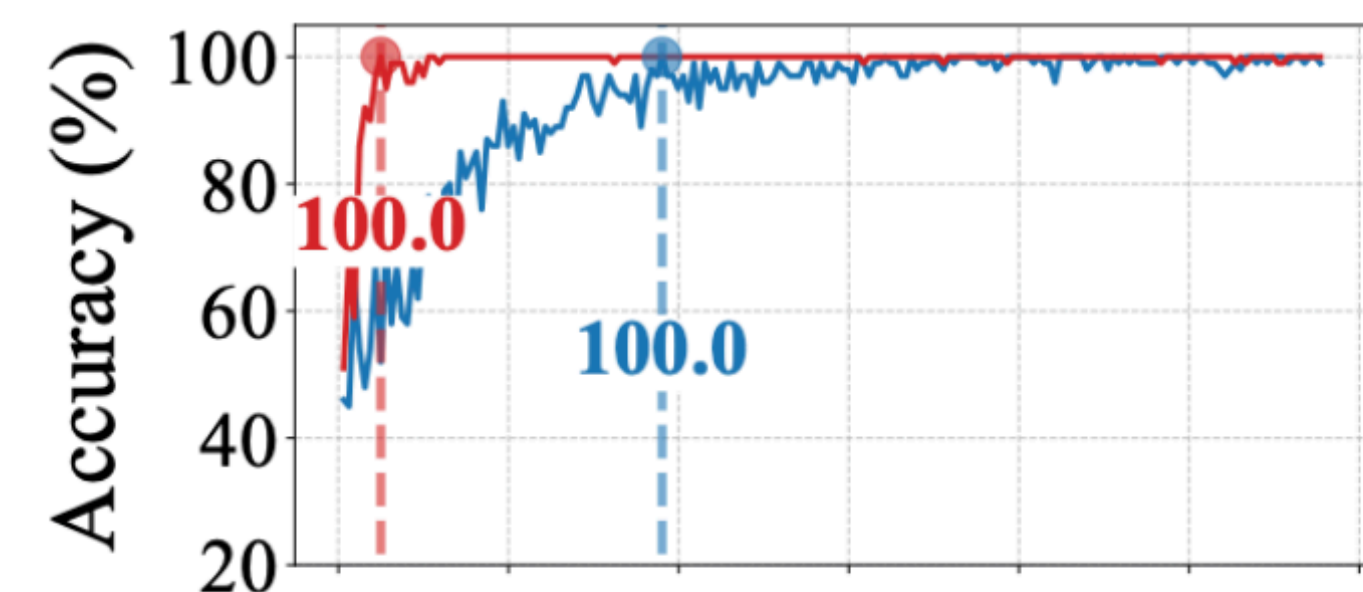
For Each Token,  $\text{len}(\text{prefix for } \text{PENCIL}) \ll \text{len}(\text{prefix for } \text{CoT})$



***QBF  $n=4$***



***QBF  $n=5$***



***QBF  $n=6$***

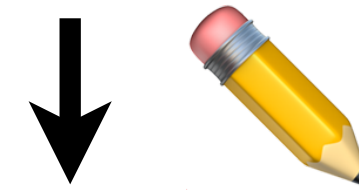
**PENCIL** is computationally more efficient than **CoT**.



# Example 3: Einstein's Puzzle

- **Constraint 1** : The green house is immediately to the right of the one who keeps birds
- **Constraint 2** : The Brit is immediately to the right of the German
- **Constraint 3** : The one who keeps dogs is the same house as the red house
- **Constraint 4** : The one who keeps birds is immediately to the right of the Swede

**Question:** who owns the fish?



**Solution :**

House #	1	2	3
Color	Red	Blue	Green
Nationality	Swede	German	Brit
Pet	Dogs	Birds	Fish

**Answer:** the Brit owns the fish



# Example 3: Einstein's Puzzle

Prompt :

- Constraint 1 : The one who keeps birds is immediately to the right of the German
- Constraint 2 : The Brit is immediately to the right of the German
- Constraint 3 : The one who keeps dogs is the same house as the red house
- Constraint 4 : The one who keeps birds is immediately to the right of the Swede

Who owns the fish?

PENCIL



Solution :

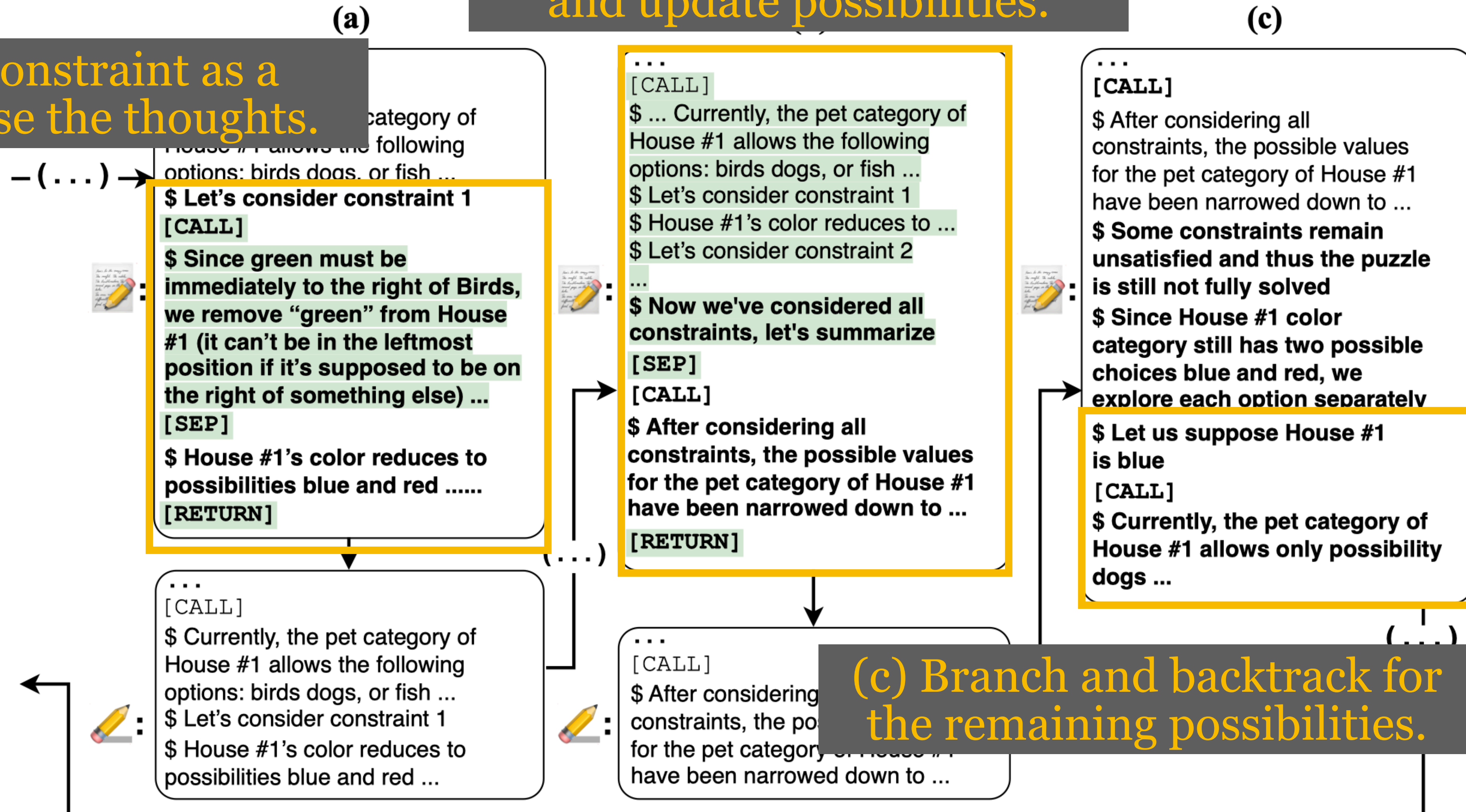
House #	1	2	3
Color	Red	Blue	Green
Nationality	Swede	German	Brit
Pet	Dogs	Birds	Fish

The Brit owns the fish

(a) Solve each constraint as a subtask and erase the thoughts.

(b) Summarize state changes and update possibilities.

(c) Branch and backtrack for the remaining possibilities.





# Special Usage: Summarization

**C** [CALL] **T** [SEP] [CALL] **T'** [RETURN]  $\Rightarrow$  **C** [CALL] **T'**

                  ↑                                  ↑  
Long Thoughts      Summarized Thoughts

**Tail recursion** in functional programming:

python

 Copy

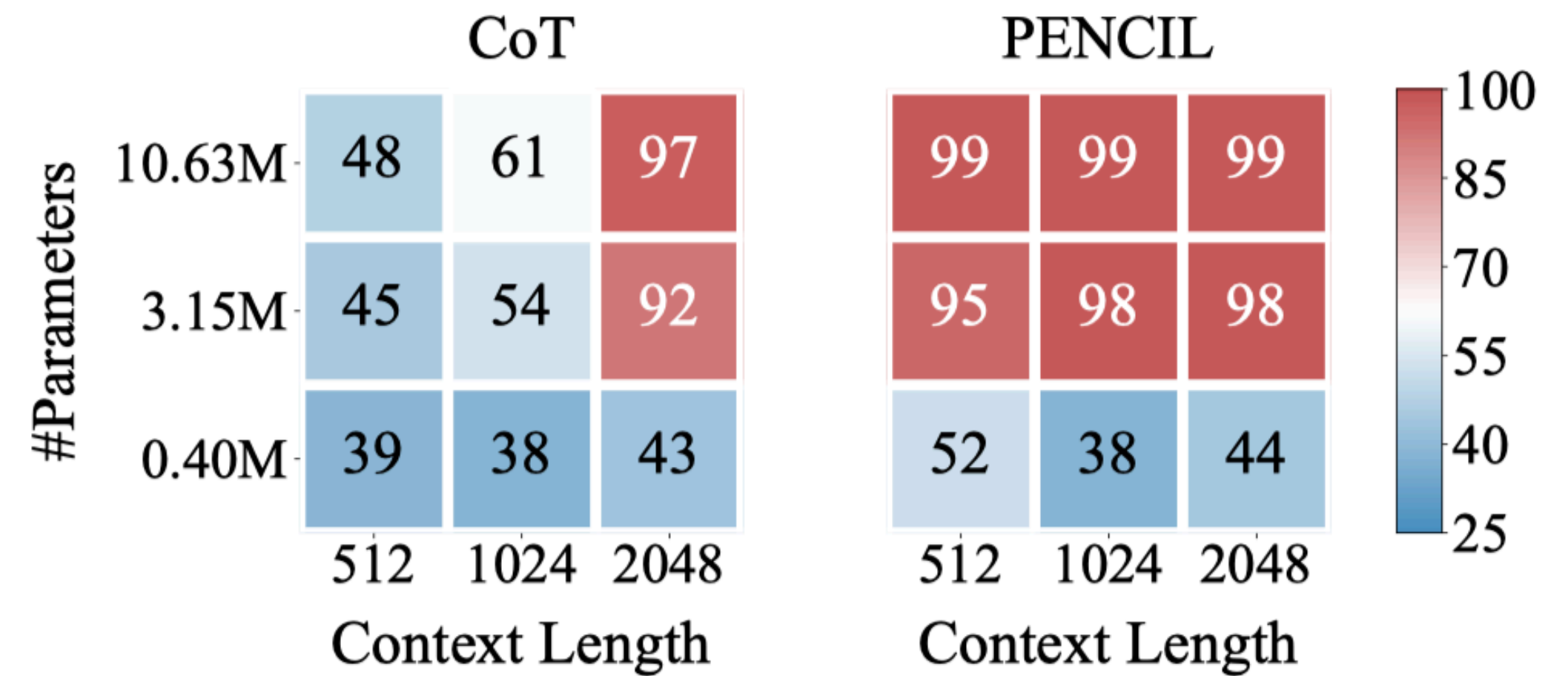
```
def factorial_tail(n, acc=1):  
    """Tail-recursive factorial: the recursive call is the final action."""  
    if n == 0:  
        return acc # base case  
    return factorial_tail(n-1, acc*n) # tail call  
  
print(factorial_tail(5)) # → 120
```

The “returned value” is another “function call”  $\Rightarrow$  **A (Answer) = [CALL] T'**



# Performance on Einstein's Puzzle

Puzzle Size	CoT	PENCIL
$5 \times 5$	25	97
$4 \times 4$	34	100
$3 \times 3$	99	99

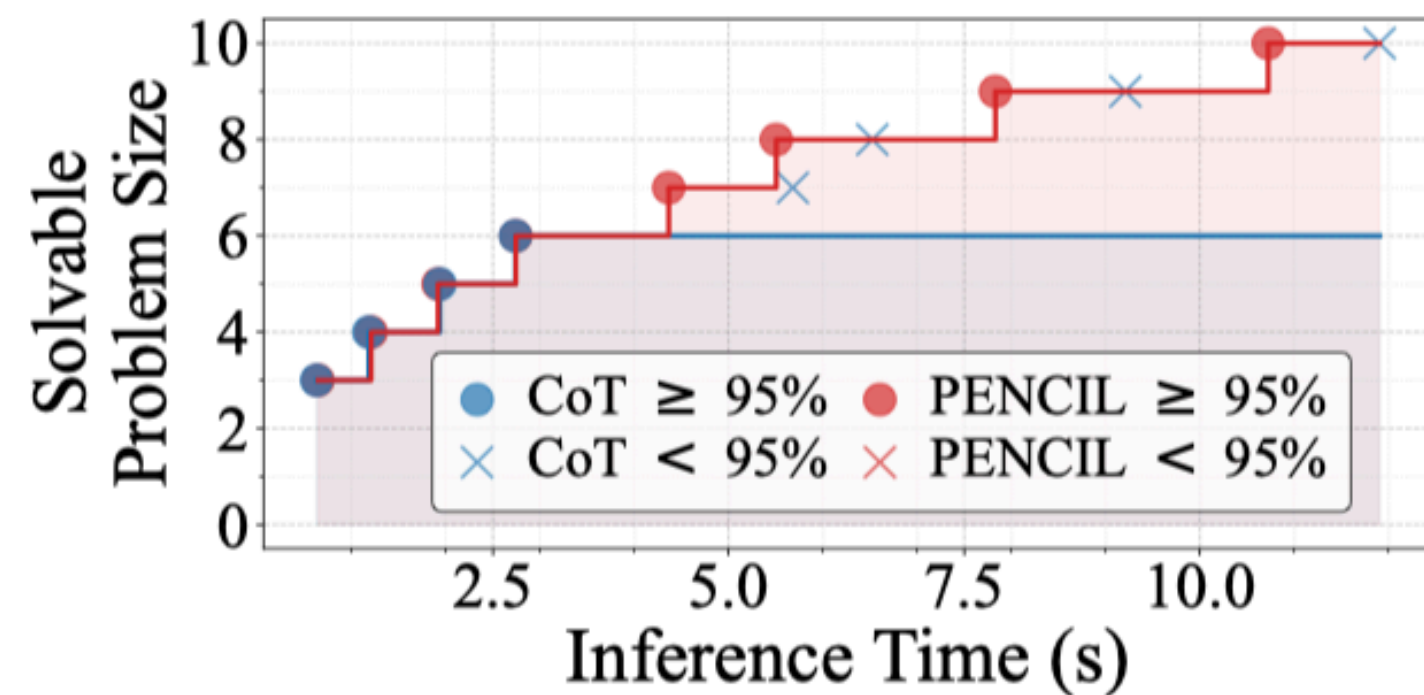


 **PENCIL** solves Einstein's puzzle almost perfectly – a logic puzzle that even GPT-4 struggles with.

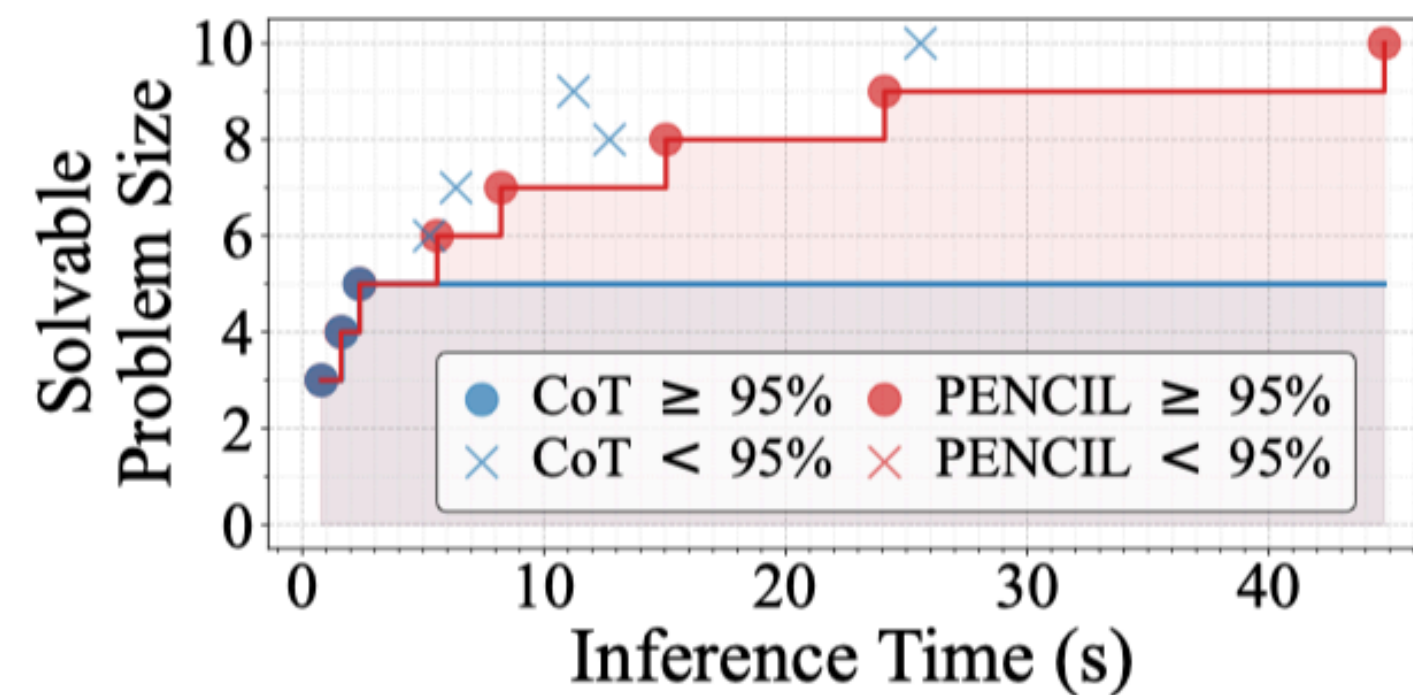
(Max context length, **CoT** = 151,192 **PENCIL** = 3,335)



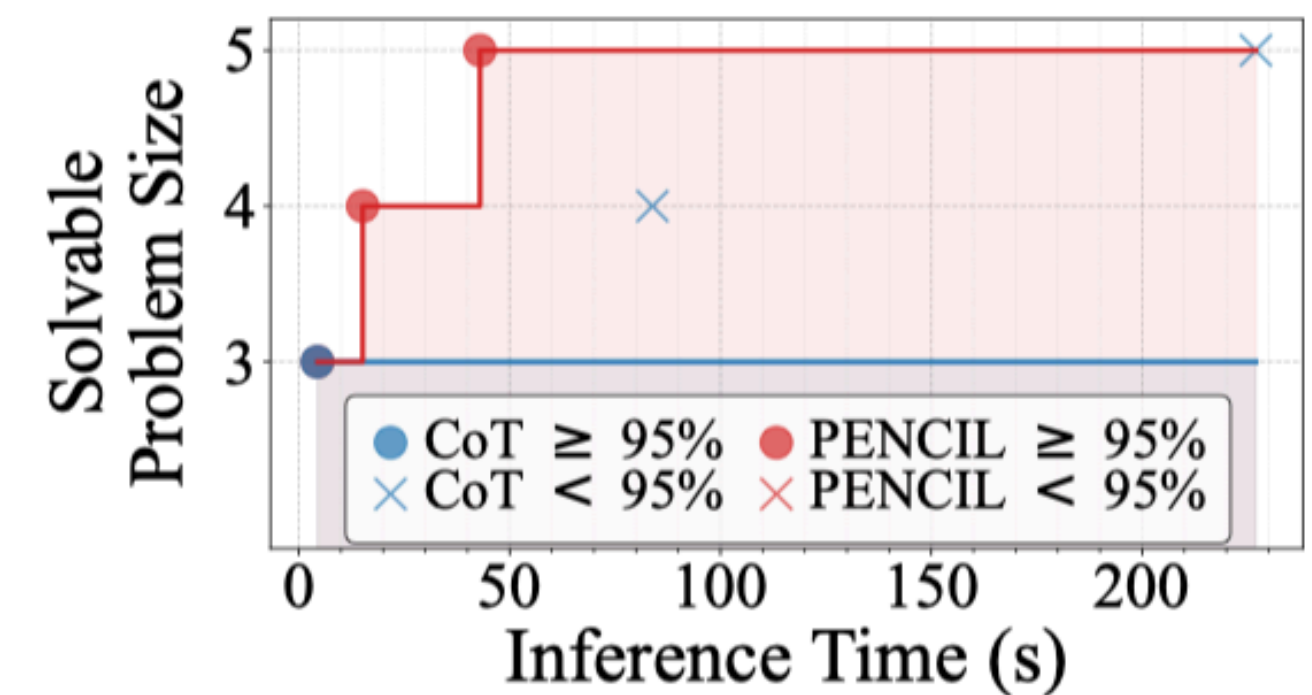
# Test-Time Scalability



***3-SAT***



***QBF***



***Einstein's Puzzle***

Given more inference time,  **PENCIL** can solve larger-sized problems.

**How about other tasks beyond 3-SAT, QBF, Einstein's Puzzle?**



A large yellow pencil is being lowered by a crane in a factory. The pencil has a red eraser at the top and a sharpened lead tip at the bottom. In the background, a large rocket engine is visible, and several workers are standing on a yellow platform. The factory has a high ceiling with a complex steel structure.

# How Powerful is **PENCIL**?

## **PENCIL** Can Perform Universal Space / Time- Efficient Computation!

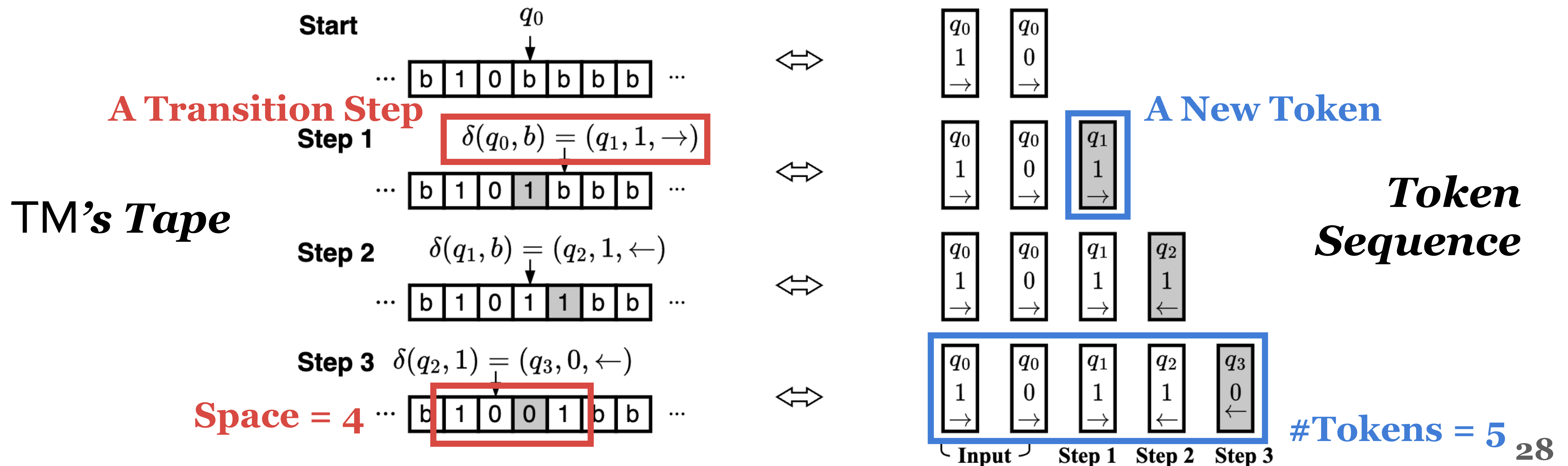


# ✍ CoT is Turing-Complete, but Inefficiently

**Theorem (Merrill et al. 24, Joshi et al. 25, etc.)**

For any Turing machine TM, there exists a finite-size decoder-only transformer such that ✍ CoT with this transformer simulates the **Turing machine** with

- **Total number of generated tokens = Maximal context length =  $\mathcal{O}(T)$**





# Universal Efficient Computation Power of PENCIL

## Theorem (Main, Informal)\*

For any Turing machine, there exists a finite-size decoder-only transformer such that for any input, on which **Turing machine** uses  $T$  steps and  $S$  space to compute,  **PENCIL** with this transformer computes the same output with

1. **Total number of generated tokens** =  $\mathcal{O}(T)$
2. **Maximal context length** =  $\mathcal{O}(S)$

- For complex problems, typically  $S \ll T$
- PENCIL is Turing-complete with optimal time and space complexity
- PENCIL can solve ANY computable tasks efficiently

*\*finite size and finite parameter precision, but infinite precision in forward pass. Also assumes average-hard attention.*

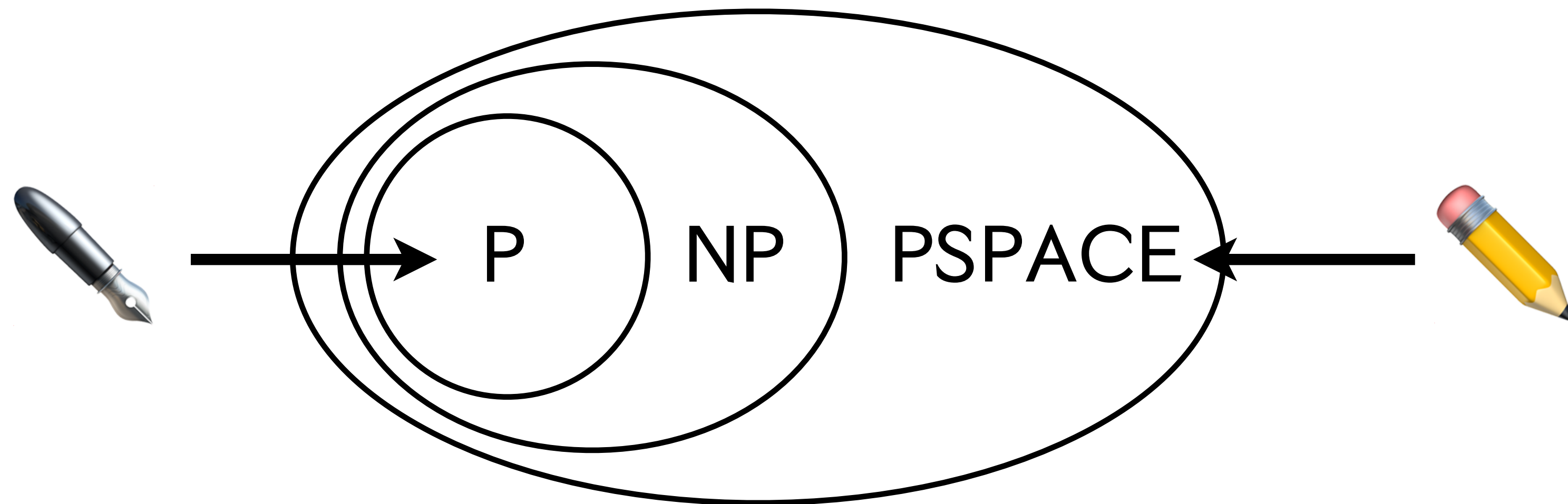


# Universal Efficient Computation Power of PENCIL

## Corollary (Informal)

With  $\text{poly}(n)$  context length,  PENCIL can solve all problems in PSPACE, while standard  CoT can only solve problems in P.

- P : Problems solvable in **polynomial time**.
- PSPACE : Problems solvable using polynomial space, **regardless of time**.





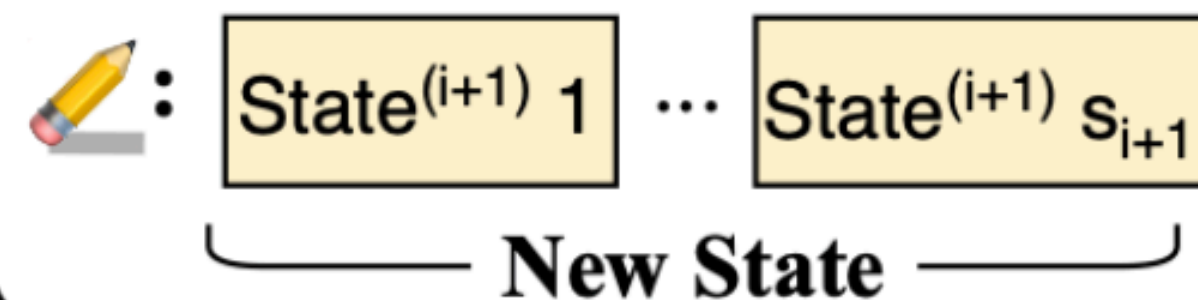
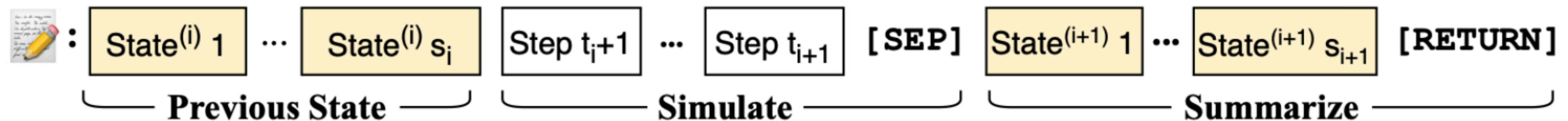
# Strategy: Iterative “Think” and “Summarize”

## Chain-of-Thought



Total Steps :  $\mathcal{O}(\text{Time})$       Max Length :  $\mathcal{O}(\text{Time})$

## PENCIL



Total Steps :  $\mathcal{O}(\text{Time})$       Max Length :  $\mathcal{O}(\text{Space})$

- **Step** : simulating a computation step of TM
- **State** : the current configuration of TM (written symbols)



# When to Summarize?

*Length exceeds twice the actual needed space.*

When?	Never (✍️)	Every Step (🖋️)	$2 T'  <  T $ (🖋️)
# Tokens Generated	$O(\text{Time})$ ✅	$O(\text{Time} \times \text{Space})$ ❌	$O(\text{Time})$ ✅
Max Context Length	$O(\text{Time})$ ❌	$O(\text{Space})$ ✅	$O(\text{Space})$ ✅

Recall we use [CALL] T [SEP] [CALL] T' [RETURN]  $\Rightarrow$  [CALL] T' to summarize.

**BUT, can *transformers* automatically detect when to summarize / erase ?**



# Proof Technique: FASP (Full-Access Sequence Processing)

## Lemma (Informal)

Programs in FASP = All finite-size transformer functions

A FASP program describes a process of constructing transformers

- **Each Variable** = a transformer
- **Each Line of Code** = an operator from simpler transformers to a more complex transformer
- **Returned Variable** = the target transformer one aims to construct

**This is the Proof !**

```
# Detect separator token
is_sep = (get_token = onehot([SEP]))
exist_sep = seq_or(is_sep)
# Phase masks to distinguish between simulation and summarization phases
sim_phase_mask = not exist_sep
sum_phase_mask = exist_sep and (not is_sep)
current_sum_pos = seq_sum(get_move and sum_phase_mask)
# Position tracking for Simulation, frozen in SUMMARIZATION (after [SEP] is generated)
next_sim_pos = seq_sum(get_move and sim_phase_mask)
current_sim_pos = next_sim_pos - (get_move and sim_phase_mask)
max_pos = seq_max(current_sim_pos)
min_pos = seq_min(current_sim_pos)
expected_sum_len = max_pos - min_pos + ReLU(max_pos - next_sim_pos - 1) + 1
summary_symbol = rightmost_best_match(current_sum_pos + min_pos, current_sim_pos, get_symbol)
# SIMULATION Phase
summary_step = get_state ⊗ summary_symbol ⊗ onehot(next_move)
# Get current symbol at head position
current_symbol = rightmost_exact_match(next_sim_pos, current_sim_pos, get_symbol, onehot(b))
# Compute next step based on transition function
simulation_step = transition(get_state, current_symbol)

# MAIN - Select appropriate action based on current phase
result = if_then_else(exist_sep, summary, simulation)
```

**returned variable**




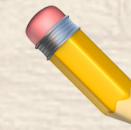

# Future Directions & Open Questions

Long Thoughts with Short Memory

- How to incorporate PENCIL into real-world LLM systems? How to more effectively teach them to reason in a structured way?
- Does there exist other “erasing” mechanisms (e.g. other reduction rules) that are even more efficient than PENCIL?
- Are there any other perspectives from which theories in TCS can help guide practice?



# Takeaways

1. We propose  **PENCIL**, a new LLM reasoning paradigm that iteratively generates and erases thoughts using the **reduction rule**:  
$$\mathbf{C} \text{ [CALL] } \mathbf{T} \text{ [SEP] } \mathbf{A} \text{ [RETURN] } \Rightarrow \mathbf{C} \mathbf{A}$$
3. Empirically,  **PENCIL** enables longer and deeper thinking using shorter context, and thus can scale up to handle more complicated tasks.
4. Theoretically,  **PENCIL** is Turing-complete with optimal space and time complexity, and thus can solve arbitrary computable problems efficiently.

## Thanks !