



# Plan and Budget: Effective and Efficient Test-Time Scaling on Large Language Model Reasoning

Junhong Lin<sup>1</sup>, Xinyue Zeng<sup>2</sup>, Jie Zhu<sup>2</sup>, Song Wang<sup>3</sup>  
Julian Shun<sup>1</sup>, Jun Wu<sup>4</sup>, Dawei Zhou<sup>2</sup>

<sup>1</sup>Massachusetts Institute of Technology, <sup>2</sup>Virginia Tech

<sup>3</sup>University of Virginia, <sup>4</sup>Michigan State University

April 7, 2026

# Reasoning Miscalibration



# Reasoning Improves Accuracy - But at What Cost?

- **Recent Advances in Test-Time Scaling**

- Chain-of-Thought (CoT)
- Self-Consistency
- Tool-augmented reasoning

→ Higher asymptotic accuracy

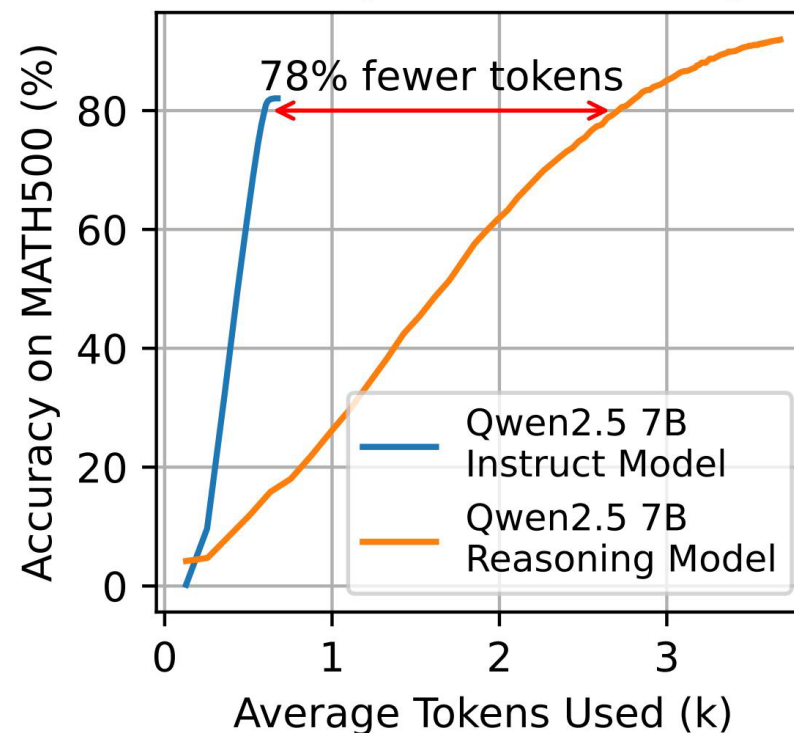
- **But...**

- Accuracy improves slowly with more tokens
- Token usage grows rapidly
- Diminishing marginal returns

- **Key Question**

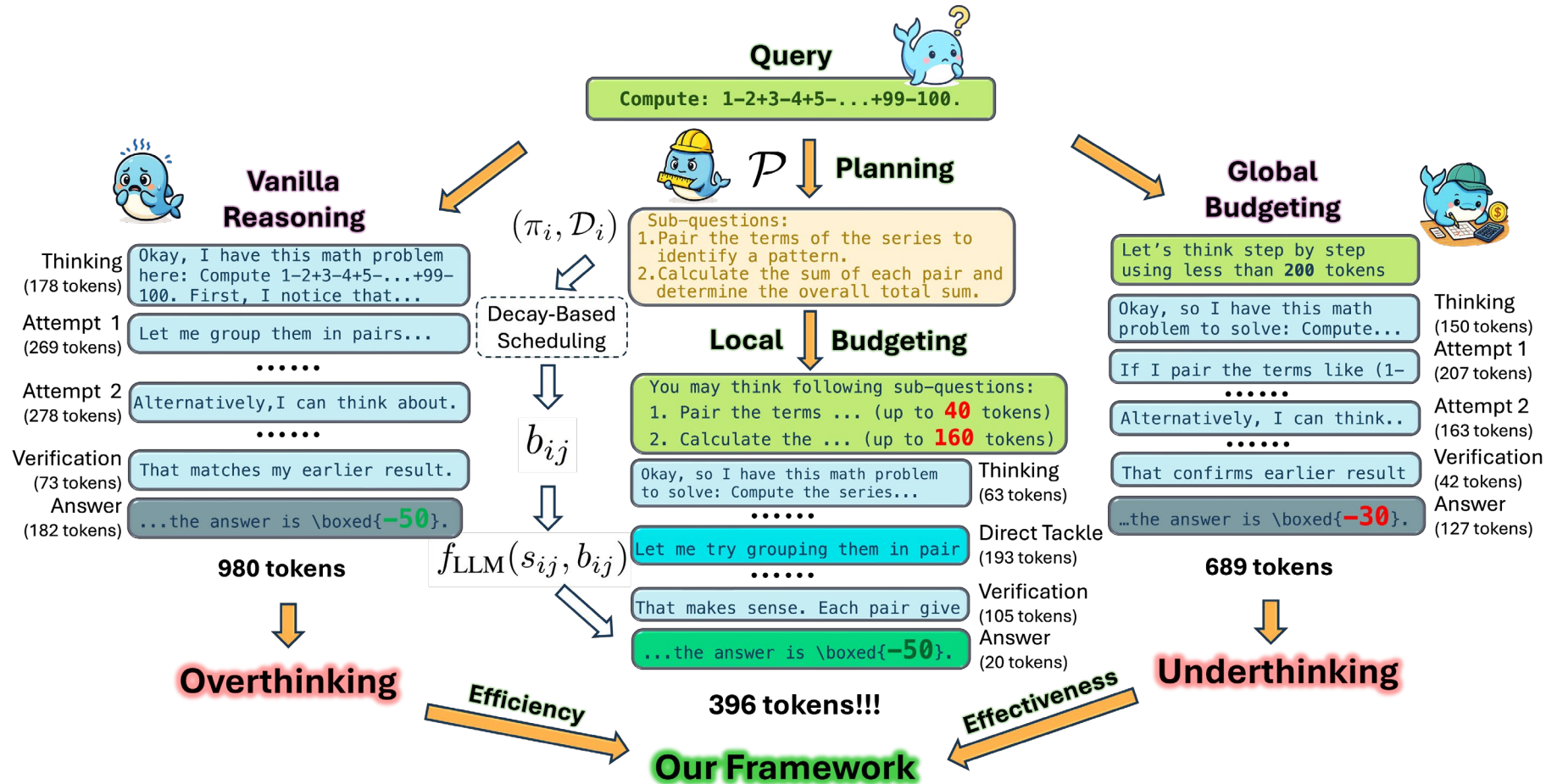
Can we achieve deep reasoning without paying the full token cost?

Accuracy vs Token Usage



**Reasoning improves accuracy, but in a highly inefficient way .**

# Reasoning Miscalibration



● **Overthinking**

→ verbose, tangential reasoning

● **Underthinking**

→ premature termination

# When Does Reasoning Go Wrong?

## ■ Case 1 — Trivial but Ambiguous → *Overthinking*

Example: “Is 2 + 3 even?”

Model may busy:

- Define parity
- Consider modular arithmetic
- Provide multiple formulations
- Re-verify (let me count my fingers)
- ....

**Model keeps reducing something that cannot be reduced further.**

## ■ Case 2 — Hard and Rare → *Underthinking*

Example: “How many 4-digit numbers can be formed using digits {1,2,3,4,5}, if no digit is repeated and the number is even? Answer in **5 seconds**”

Model may say:

“There are 5 digits, choose 4 →  $5P_4 = 120$ , I think half of them are even → 60”

- No verification
- No intermediate check
- Stops early

**The model fails to invest where uncertainty is reducible.**

It's not about too many tokens or too few tokens.  
It's about allocating compute incorrectly.

# When Does Reasoning Go Wrong?

## ■ Case 1 — Trivial but Ambiguous → *Overthinking*

Example: “Is 2 + 3 = 5?”

Model may busy:

- Define parity
- Consider mod
- Provide multip
- Re-verify (let m
- ....

Model keeps reducing something that cannot be reduced further.

## ■ Case 2 — Hard and Rare → *Underthinking*

ers can be formed  
epeated and the  
ls”

120, I think half

- Stops early

The model fails to invest where uncertainty is reducible.

**How can we guide LLMs to allocate computation adaptively based on task complexity?**

We reinterpret reasoning as a **resource allocation problem under uncertainty.**

It's not about too many tokens or too few tokens.  
It's about allocating compute incorrectly.

# Two Types of Uncertainty in Reasoning

- When a model is unsure, why is it unsure?

$$\mathcal{U}(x) = \mathcal{U}_{\text{epistemic}}(x) + \mathcal{U}_{\text{aleatoric}}(x)$$

## ■ Epistemic Uncertainty (Reducible)

The model doesn't know enough —  
but more reasoning could help.

**Query: “Compute  $37 \times 48$ .”**

If the model hesitates, it's because:

- It hasn't completed the multiplication
- It hasn't derived the intermediate steps

**More reasoning → uncertainty decreases.**

## ■ Aleatoric Uncertainty (Irreducible)

The ambiguity comes from the input itself.  
More thinking doesn't help.

**Query: “Is the bank open?”**

- Which bank?
- Which city?
- Which time?

**Even infinite reasoning can't resolve missing info.**

# Reasoning as a Sequence of Sub-Problems

- LLM reasoning is **multi-step** — implicitly solving a chain of sub-questions.
- Different sub-questions have **different uncertainty profiles**:
  - Some steps are structural / strategy (high epistemic → benefit from more compute)
  - Some steps are routine execution (low epistemic → extra tokens are wasteful)
- **Miscalibration happens when we spend too much compute on easy steps, and too little on hard steps.**

**Compute:**  $1 - 2 + 3 - 4 + \dots + 99 - 100$

**Sub-questions:**

- Pair terms to identify a pattern
- Compute the sum of each pair
- Aggregate to get the final total

# From Theory to Practice: **Plan-And-Budget**



# Budget Allocation Model (BAM)

- We start with a simple assumption inspired by neural scaling laws
  - Allocating more tokens reduces epistemic uncertainty, but with diminishing returns.
  - Uncertainty decreases roughly as a power law with compute:

$$\mathcal{U}_{\text{epistemic}}(s_{ij} | b_{ij}) = \frac{c_{ij}}{b_{ij}^{\beta_{ij}}}$$

$s$ : sub-questions

$b$ : tokens allocated

$c$ : initial uncertainty

$\beta$ : difficulty parameter

- Given a total token budget  $B$ , the optimal allocation of budget across subquestion is:

$$b_{ij} = B_i \cdot \frac{(c_{ij} \beta_{ij})^{\frac{1}{\beta_{ij}+1}}}{\sum_k (c_{ik} \beta_{ik})^{\frac{1}{\beta_{ik}+1}}}$$

Allocate more tokens to sub-questions that:

- Start with high uncertainty (large  $c$ )
- And where additional compute is still effective (moderate  $\beta$ )
- Extremely easy steps get few tokens.
- Extremely hard steps also get few tokens.
- Moderately difficult steps get the most.

# Efficiency-aware Effectiveness Evaluation Score ( $\mathcal{E}^3$ )

- Many existing works measure efficiency as accuracy per token:

$$A/T$$

- We propose to weights the A/T metric by accuracy to favor methods that maintain high correctness:

$$\mathcal{E}^3 = A^2 / T$$

## Example

For **Model A** (Acc = **90%**, Tokens = **2000**), and **Model B** (Acc = **45%**, Tokens = **1000**), they achieve the **same A/T score**:

$$A/T = 0.00045$$

## Example

For Model A :

$$\mathcal{E}^3 = \frac{0.9^2}{2000} = 0.000405$$

For Model B:

$$\mathcal{E}^3 = \frac{0.45^2}{1000} = 0.000203$$

Now the scores are different: the high-accuracy model is strongly preferred.

# Plan-and-Budget:

## Structured & Adaptive Reasoning

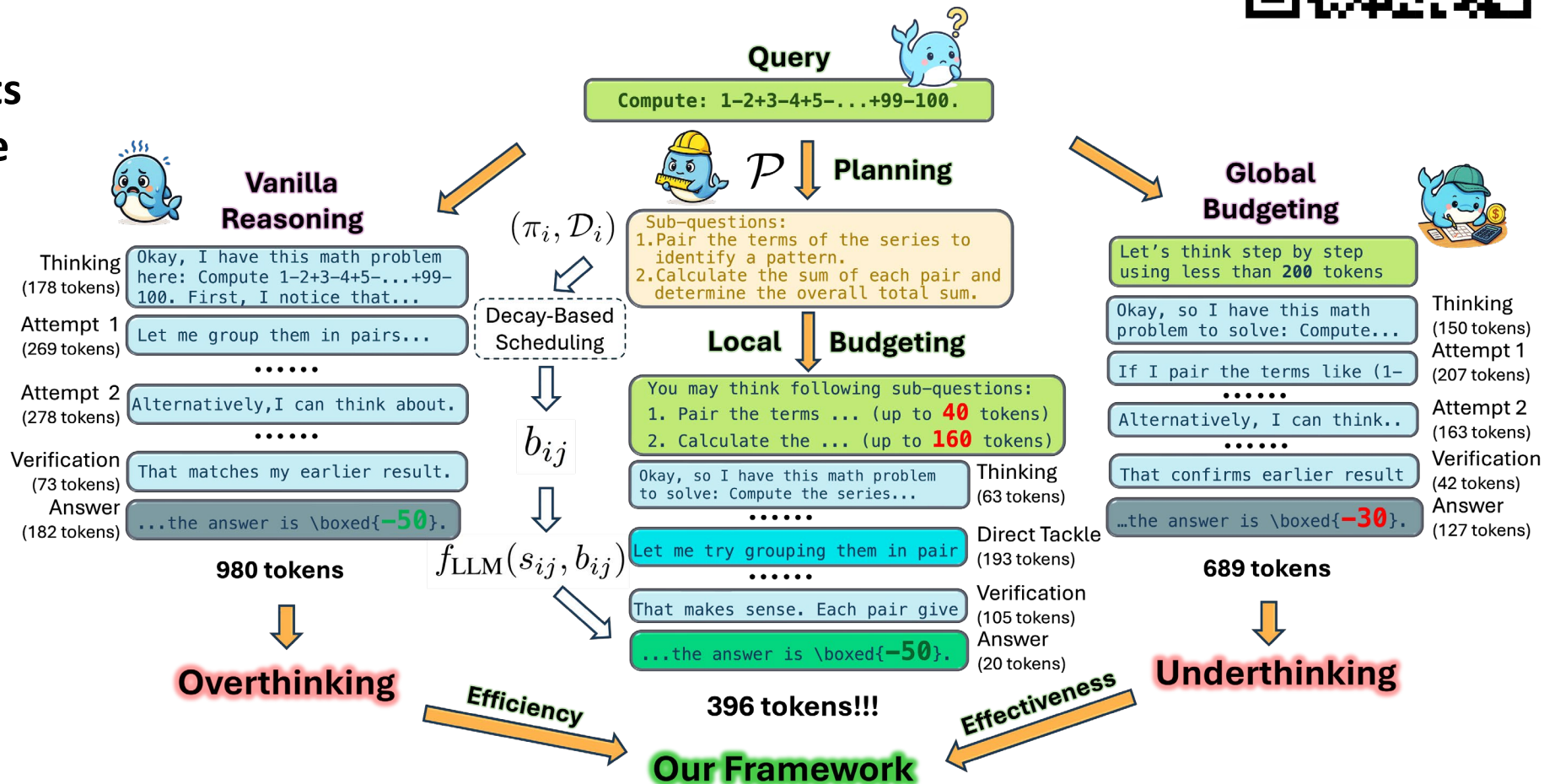
arXiv 2505.16122

License MIT



Plan-and-Budget implements uncertainty-guided compute allocation:

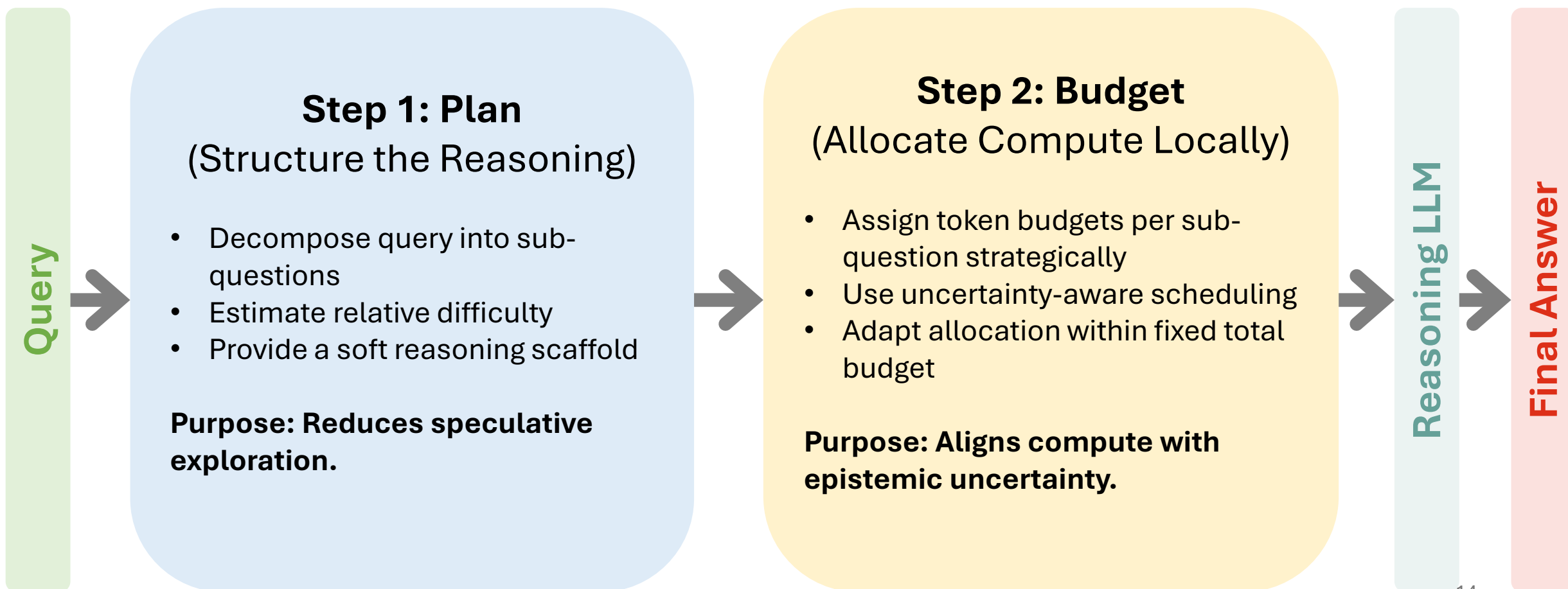
- **Planning:** Decompose queries into sub-questions
- **Local Budgeting:** Allocate tokens adaptively across steps
- **Uncertainty-Aware:** Reduces epistemic uncertainty efficiently
- **Training-Free:** Works with any LLM at test time



# Plan-and-Budget: General Workflow

arXiv 2505.16122

License MIT



# Plan-and-Budget:

## Practical Budget Allocation via Decay Scheduling

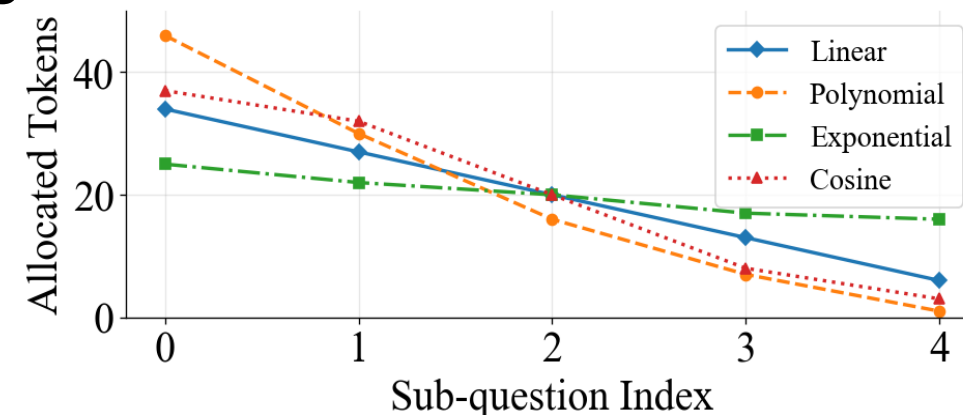
arXiv 2505.16122

License MIT



We approximate optimal allocation using simple decay schedules:

- Allocate more tokens to early (high-uncertainty) steps.
- Gradually reduce budget across sub-questions.
- Different decay shapes reflect different assumptions.
- All schedules are lightweight and training-free.



Strategy	Formula of $d_{ij}$	Description
Non-decay	1	Equal priority for all sub-questions; budget follows $w_{ij}$ .
Linear decay	$m - j$	Decreases priority linearly with $j$ ; emphasizes early steps.
Polynomial decay	$(m - j)^p$	Stronger emphasis on early steps; steeper with higher $p > 1$ .
Exponential decay	$\gamma^j$	Exponentially favors earlier sub-questions; controlled by $\gamma \in (0, 1)$ .
Cosine annealing	$0.5 \left( 1 + \cos \left( \frac{\pi j}{m-1} \right) \right) + \epsilon$	Smooth decay with mid-sequence flexibility; $\epsilon$ adds stability.

# Results



# Experiment Setup

## ■ Tasks

3 key tasks covering math, general reasoning, and long-horizon planning:

- **MATH-500**
  - Multi-step symbolic reasoning
  - Exact match accuracy
- **NaturalInstructions**
  - Instruction following
  - ROUGE score
- **TravelPlanner**
  - Constraint-heavy planning
  - Hard constraint pass rate

## ■ Models

4 main-stream reasoning LLM without fine-tuning:

- **DS-Qwen-32B**
- **QwQ-32B**
- **DS-LLaMA-70B**
- **o4-mini**
- **Vanilla reasoning**
- **Global budget constraints**
- **Planned global budgeting**
- **Plan-and-Budget (ours)**

## ■ Metrics

- Accuracy / ROUGE / Pass rate
- Average tokens (completion tokens)
- $E^3 = \frac{A^2}{T}$

# Extensive Results

- We conduct extensive experiment covering 3 kinds of tasks X 4 LLMs
- And you are probably too busy to read through them...

Table 3: Experiment results across different reasoning models on MATH-500. Acc denotes accuracy.

Models →	DeepSeek-R1-Distill-Qwen-32B			QwQ-32B			DeepSeek-R1-Distill-Llama-70B			o4-mini			
	Methods↓	Acc (%)↑	Avg. Tok.↓	$\mathcal{E}^3$ ↑	Acc (%)↑	Avg. Tok.↓	$\mathcal{E}^3$ ↑	Acc (%)↑	Avg. Tok.↓	$\mathcal{E}^3$ ↑	Acc (%)↑	Avg. Tok.↓	$\mathcal{E}^3$ ↑
Direct	Vanilla	89.76 $\pm$ 0.26	2105.12 $\pm$ 31.94	3.83	84.88 $\pm$ 1.18	3523.72 $\pm$ 97.42	2.04	90.44 $\pm$ 0.61	2286.63 $\pm$ 26.42	3.58	<b>93.16</b> $\pm$ 0.89	711.20 $\pm$ 8.31	12.20
	Global Budget	89.60 $\pm$ 0.88	1526.15 $\pm$ 10.09	5.26	<b>90.56</b> $\pm$ 0.33	2565.18 $\pm$ 37.10	3.20	90.80 $\pm$ 0.62	1810.83 $\pm$ 51.64	4.55	91.84 $\pm$ 0.48	636.41 $\pm$ 8.14	13.25
Planned	Vanilla	<b>91.04</b> $\pm$ 0.62	1883.73 $\pm$ 63.82	4.40	85.30 $\pm$ 1.56	3309.69 $\pm$ 18.06	2.20	92.12 $\pm$ 1.16	2022.38 $\pm$ 28.74	4.20	91.88 $\pm$ 1.36	539.36 $\pm$ 18.94	15.65
	Global Budget	91.24 $\pm$ 1.34	1552.62 $\pm$ 29.93	5.36	88.20 $\pm$ 1.17	2671.60 $\pm$ 15.02	2.91	92.56 $\pm$ 0.71	1661.24 $\pm$ 34.43	5.16	91.84 $\pm$ 0.75	586.18 $\pm$ 6.50	14.39
PLAN-AND-BUDGET	+ Uniform	90.16 $\pm$ 0.74	1440.70 $\pm$ 47.55	5.64	88.68 $\pm$ 0.58	2397.16 $\pm$ 23.01	3.28	92.28 $\pm$ 0.41	1575.04 $\pm$ 29.68	5.41	91.36 $\pm$ 0.85	525.53 $\pm$ 18.88	15.88
	+ Weighted	90.48 $\pm$ 0.46	1485.99 $\pm$ 45.63	5.51	87.45 $\pm$ 0.66	2479.46 $\pm$ 39.21	3.08	92.64 $\pm$ 0.68	1557.64 $\pm$ 47.71	5.51	91.64 $\pm$ 1.21	538.22 $\pm$ 5.30	15.60
	+ Linear	90.04 $\pm$ 0.46	<b>1336.27</b> $\pm$ 31.18	<b>6.07</b>	88.13 $\pm$ 0.90	2346.35 $\pm$ 25.33	3.31	92.32 $\pm$ 0.88	1529.98 $\pm$ 45.35	5.57	90.56 $\pm$ 0.73	534.45 $\pm$ 7.64	15.34
	+ Exponential	90.80 $\pm$ 0.68	1389.75 $\pm$ 61.06	5.93	87.90 $\pm$ 1.27	2320.04 $\pm$ 72.33	3.33	93.04 $\pm$ 0.22	<b>1469.29</b> $\pm$ 73.77	<b>5.89</b>	90.88 $\pm$ 0.36	525.51 $\pm$ 11.70	15.72
	+ Polynomial	90.04 $\pm$ 0.26	1371.59 $\pm$ 21.75	5.91	88.27 $\pm$ 0.99	2346.94 $\pm$ 17.73	3.32	91.92 $\pm$ 1.15	1514.43 $\pm$ 47.94	5.58	90.36 $\pm$ 0.83	525.00 $\pm$ 9.15	15.55
	+ Cosine	89.88 $\pm$ 1.72	1365.51 $\pm$ 44.92	5.92	88.60 $\pm$ 0.28	<b>2306.83</b> $\pm$ 24.11	<b>3.40</b>	<b>92.88</b> $\pm$ 0.46	1487.83 $\pm$ 61.78	5.80	91.32 $\pm$ 0.94	<b>522.89</b> $\pm$ 10.01	<b>15.95</b>

Table 4: Experiment results across different reasoning models on NaturalInstructions.

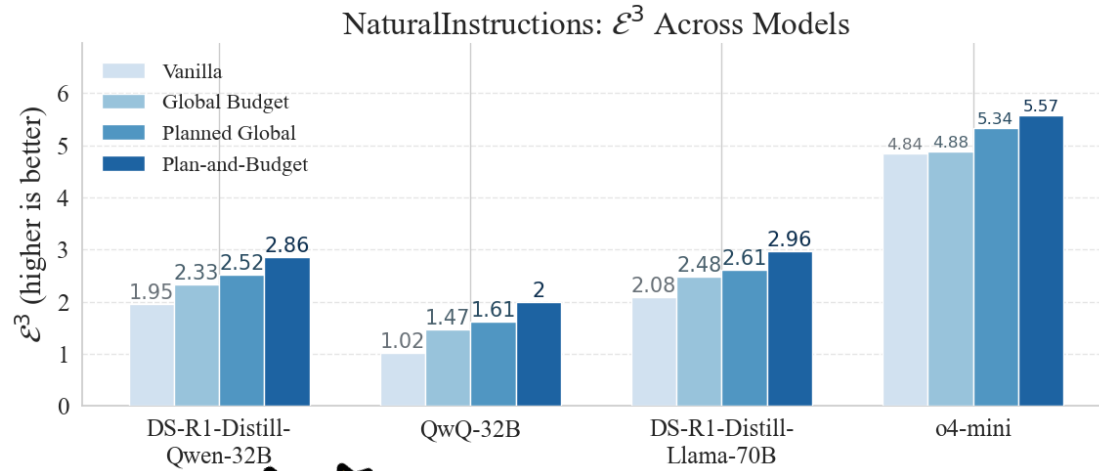
Models →	DeepSeek-R1-Distill-Qwen-32B			QwQ-32B			DeepSeek-R1-Distill-Llama-70B			o4-mini			
	Methods↓	ROUGE (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑	ROUGE (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑	ROUGE (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑	ROUGE (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑
Direct	Vanilla	<b>43.47</b> $\pm$ 0.52	968.17 $\pm$ 44.78	1.95	43.16 $\pm$ 1.12	1818.34 $\pm$ 24.99	1.02	43.13 $\pm$ 0.76	894.46 $\pm$ 50.69	2.08	<b>47.24</b> $\pm$ 0.31	460.99 $\pm$ 11.31	4.84
	Global Budget	42.81 $\pm$ 0.39	787.25 $\pm$ 58.17	2.33	<b>44.77</b> $\pm$ 0.73	1360.49 $\pm$ 101.64	1.47	<b>43.80</b> $\pm$ 1.28	772.98 $\pm$ 47.44	2.48	45.39 $\pm$ 1.27	422.20 $\pm$ 56.78	4.88
Planned	Vanilla	42.48 $\pm$ 0.67	860.85 $\pm$ 49.58	2.10	44.24 $\pm$ 0.67	1426.74 $\pm$ 52.92	1.37	43.40 $\pm$ 0.18	821.27 $\pm$ 21.85	2.29	43.78 $\pm$ 1.47	<b>344.99</b> $\pm$ 14.44	5.56
	Global Budget	42.50 $\pm$ 0.36	717.98 $\pm$ 36.28	2.52	45.13 $\pm$ 0.56	1265.78 $\pm$ 33.23	1.61	42.48 $\pm$ 0.33	691.79 $\pm$ 12.18	2.61	43.78 $\pm$ 0.96	358.84 $\pm$ 14.44	5.34
PLAN-AND-BUDGET	+ Uniform	41.03 $\pm$ 0.55	644.87 $\pm$ 46.34	2.61	44.47 $\pm$ 0.35	996.91 $\pm$ 31.31	1.98	43.06 $\pm$ 0.33	665.94 $\pm$ 47.22	2.78	44.08 $\pm$ 0.81	348.74 $\pm$ 8.13	<b>5.57</b>
	+ Weighted	41.29 $\pm$ 0.50	663.9427.29	2.57	44.40 $\pm$ 0.61	1025.02 $\pm$ 24.91	1.92	43.05 $\pm$ 0.39	626.37 $\pm$ 19.46	<b>2.96</b>	43.72 $\pm$ 1.00	371.85 $\pm$ 9.53	5.14
	+ Linear	41.56 $\pm$ 0.50	633.79 $\pm$ 34.17	2.73	44.22 $\pm$ 0.66	1003.24 $\pm$ 26.23	1.95	42.05 $\pm$ 0.99	<b>613.05</b> $\pm$ 33.68	2.88	44.21 $\pm$ 0.44	363.65 $\pm$ 13.70	5.37
	+ Exponential	41.44 $\pm$ 0.50	650.19 $\pm$ 31.35	2.64	43.99 $\pm$ 0.22	1026.89 $\pm$ 8.51	1.88	42.73 $\pm$ 0.24	622.72 $\pm$ 33.58	2.93	43.68 $\pm$ 1.06	364.86 $\pm$ 10.81	5.23
	+ Polynomial	41.44 $\pm$ 0.78	<b>600.04</b> $\pm$ 40.52	<b>2.86</b>	44.66 $\pm$ 0.68	<b>995.95</b> $\pm$ 14.43	<b>2.00</b>	43.19 $\pm$ 0.44	641.62 $\pm$ 32.22	2.91	44.63 $\pm$ 1.04	363.16 $\pm$ 11.71	5.48
	+ Cosine	41.43 $\pm$ 1.01	628.20 $\pm$ 36.63	2.73	44.53 $\pm$ 0.54	1000.64 $\pm$ 17.85	1.98	42.83 $\pm$ 0.63	657.93 $\pm$ 59.06	2.79	44.36 $\pm$ 1.06	363.05 $\pm$ 16.72	5.42

Table 5: Experiment results on TravelPlanner. Rate denotes the hard constraint pass rate.

Models →	DeepSeek-R1-Distill-Qwen-32B			QwQ-32B			DeepSeek-R1-Distill-Llama-70B			o4-mini			
	Methods↓	Rate (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑	Rate (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑	Rate (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑	Rate (%)↑	Avg. Tokens ↓	$\mathcal{E}^3$ ↑
Direct	Vanilla	14.33 $\pm$ 2.17	1430.14 $\pm$ 43.73	0.14	34.89 $\pm$ 3.20	3432.33 $\pm$ 78.66	0.35	26.22 $\pm$ 1.82	1361.37 $\pm$ 47.93	0.50	11.58 $\pm$ 2.15	1559.65 $\pm$ 8.84	0.086
	Global Budget	13.78 $\pm$ 1.20	1158.81 $\pm$ 20.23	0.16	30.78 $\pm$ 2.06	2530.04 $\pm$ 40.87	0.37	24.33 $\pm$ 2.30	1215.29 $\pm$ 35.05	0.49	8.33 $\pm$ 1.71	<b>1248.53</b> $\pm$ 26.97	0.056
Planned	Vanilla	20.22 $\pm$ 1.01	1343.67 $\pm$ 62.44	0.30	<b>37.22</b> $\pm$ 1.80	3669.88 $\pm$ 42.09	0.38	30.67 $\pm$ 2.17	1464.50 $\pm$ 65.40	0.64	<b>12.20</b> $\pm$ 2.47	1640.46 $\pm$ 95.33	0.091
	Global Budget	22.56 $\pm$ 2.41	1241.19 $\pm$ 54.66	0.41	35.22 $\pm$ 4.85	3199.58 $\pm$ 63.14	0.39	30.67 $\pm$ 1.73	1220.41 $\pm$ 32.22	0.77	7.19 $\pm$ 2.43	1392.11 $\pm$ 31.05	0.037
PLAN-AND-BUDGET	+ Uniform	20.67 $\pm$ 1.20	1227.99 $\pm$ 68.55	0.35	36.00 $\pm$ 2.79	2854.24 $\pm$ 44.87	0.45	31.56 $\pm$ 2.20	1232.98 $\pm$ 34.16	0.81	11.00 $\pm$ 1.62	1345.32 $\pm$ 58.88	0.090
	+ Weighted	<b>23.33</b> $\pm$ 1.11	1222.09 $\pm$ 40.69	0.45	33.89 $\pm$ 2.22	2842.74 $\pm$ 77.68	0.40	29.67 $\pm$ 3.01	1197.32 $\pm$ 10.78	0.74	10.91 $\pm$ 3.01	1353.67 $\pm$ 37.64	0.088
	+ Linear	19.56 $\pm$ 2.47	<b>1136.18</b> $\pm$ 54.92	0.34	34.55 $\pm$ 2.65	2671.70 $\pm$ 67.97	0.45	31.67 $\pm$ 2.32	1162.24 $\pm$ 43.31	0.86	11.66 $\pm$ 1.96	1306.54 $\pm$ 55.05	0.103
	+ Exponential	21.44 $\pm$ 2.98	1156.64 $\pm$ 30.52	0.40	35.44 $\pm$ 2.06	2724.23 $\pm$ 41.87	0.46	32.00 $\pm$ 2.14	1187.85 $\pm$ 36.57	0.86	9.91 $\pm$ 1.96	1307.87 $\pm$ 40.83	0.075
	+ Polynomial	23.11 $\pm$ 2.14	1148.53 $\pm$ 37.33	<b>0.47</b>	35.00 $\pm$ 3.35	2511.35 $\pm$ 84.18	0.49	<b>32.67</b> $\pm$ 2.06	<b>1148.14</b> $\pm$ 59.00	<b>0.93</b>	11.49 $\pm$ 1.31	1266.11 $\pm$ 28.48	<b>0.104</b>
	+ Cosine	20.22 $\pm$ 2.34	1140.79 $\pm$ 6.68	0.36	36.18 $\pm$ 3.00	<b>2496.46</b> $\pm$ 40.10	<b>0.52</b>	31.67 $\pm$ 2.22	1173.96 $\pm$ 44.22	0.85	9.79 $\pm$ 1.57	1252.06 $\pm$ 80.85	0.077

# Adaptive Allocation Improves Efficiency

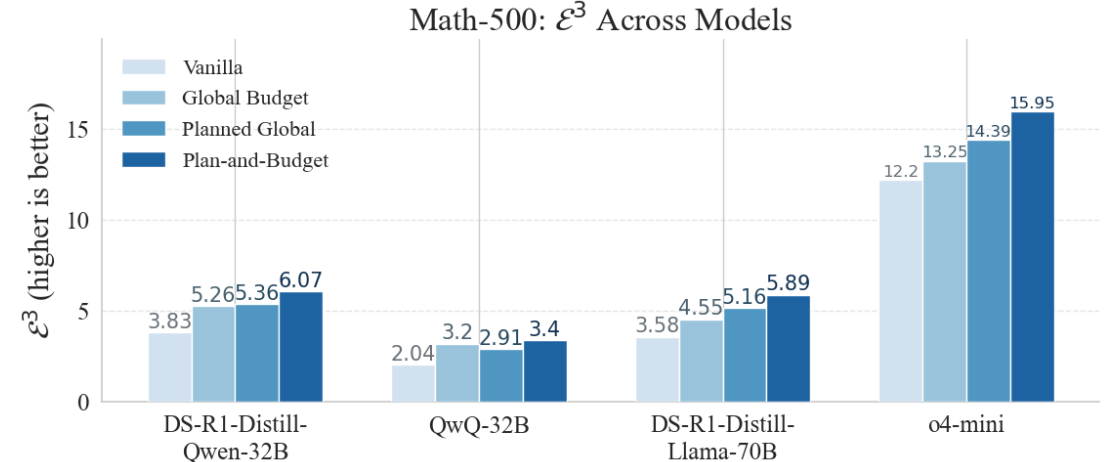
- Consistent gains across datasets
- Largest improvement: **+193.8%**
- Improvement holds across model scales



**Up to +193.8%  
 $\mathcal{E}^3$  Gain!**

**Up to +70%  
Accuracy Gain!**

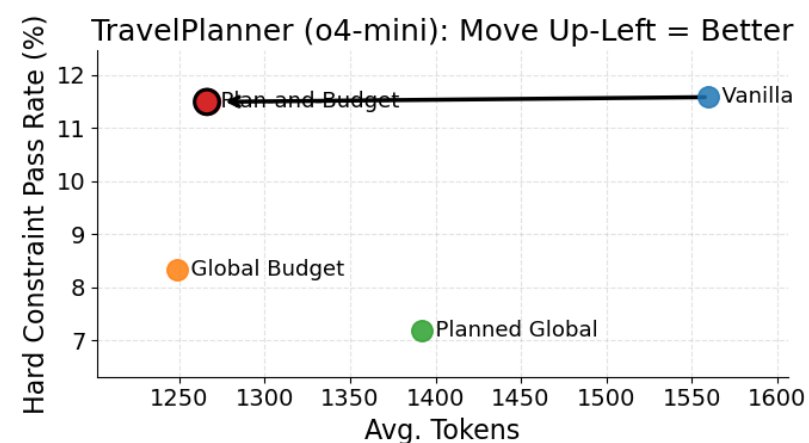
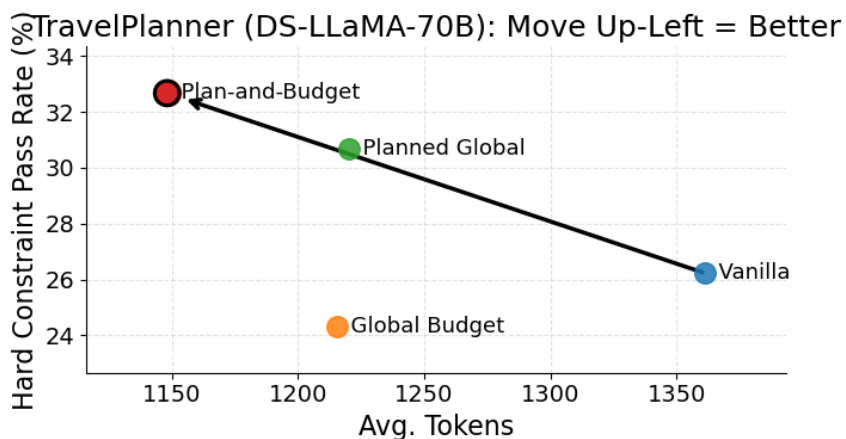
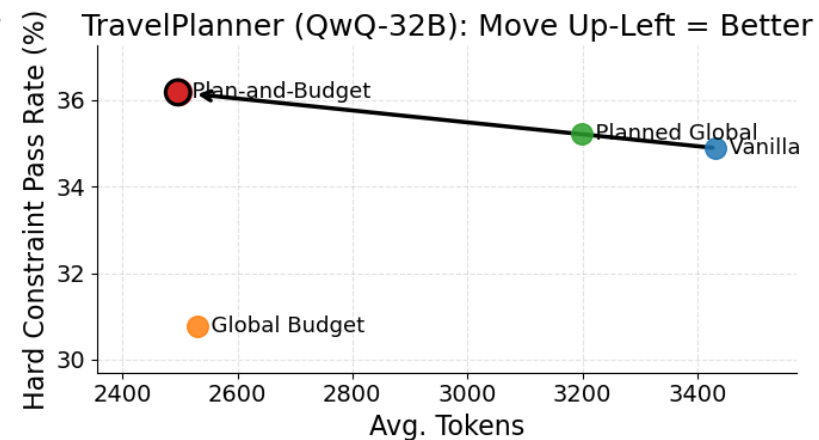
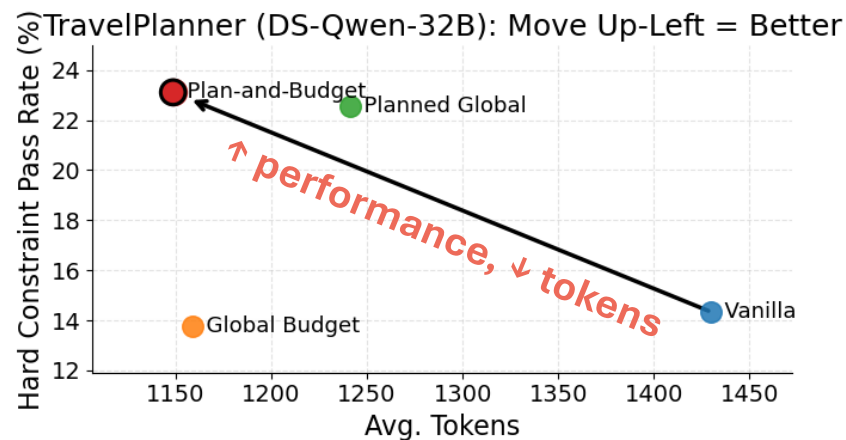
**Up to -39%  
Token Reduction!**



# Adaptive Allocation Improves Performance at Lower Cost

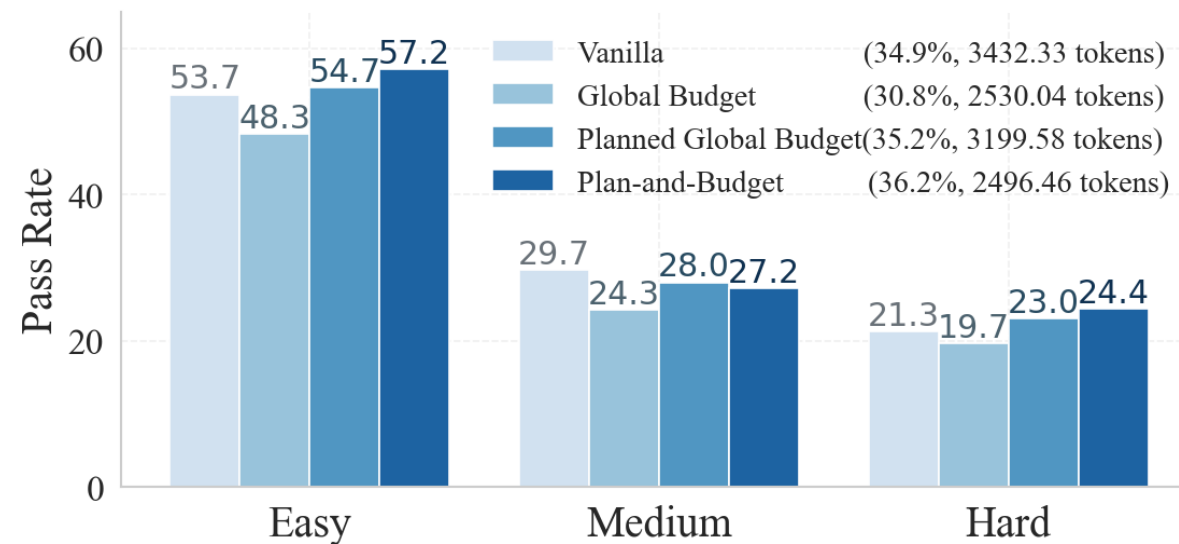
- Left = **fewer** ↓ tokens
- Up = **higher** ↑ correctness
- Global budget  
→ left but down (underthink)
- Plan-and-Budget  
→ **left and up** (calibrated reasoning)

Efficient reasoning  
requires **calibrated**  
**compute allocation**, not  
simply reducing tokens.



# Global Budgets Misallocate Compute Across Difficulty

- Global budget uniformly reduces tokens  
→ performance **drops across difficulty** levels
- Plan-and-Budget reallocates compute  
→ improves Easy and Hard queries
- Effects vary across models and datasets  
→ consistent trend is adaptive redistribution



**The problem is not difficulty alone, it is uniform compute allocation across heterogeneous queries.**

# Key Takeaways

arXiv 2505.16122

License MIT



- 1. Reasoning Miscalibration is a Core Bottleneck:** LLMs often misallocate compute at test time:
  - overthinking on simple queries,
  - underthinking on complex ones.
- 2. Adaptive Allocation Improves Efficiency and Accuracy:** Plan-and-Budget:
  - decomposes reasoning into structured sub-problems,
  - allocates tokens based on uncertainty,
  - achieves better performance with fewer tokens.
- 3. Structured Control Enables Efficient Reasoning**

**Compute should be allocated based on uncertainty,  
not applied uniformly across reasoning steps.**