

# Conformal Thinking

Risk Control for Reasoning on a Compute Budget

Xi Wang, Anushri Suresh, Alvin Zhang, Rishi More, William Jurayj, Benjamin Van Durme, Mehrdad Farajtabar, Daniel Khashabi, Eric Nalisnick

Johns Hopkins University · Apple — equal contribution



arXiv:2602.03814

Adaptive stopping turns “set a token budget” into “set a threshold.” We turn it into “set a risk” — an interpretable target a validation set calibrates for you.

## 01 The problem

Reasoning LLMs *overthink*: more tokens buy accuracy, but much of that compute is wasted. Adaptive early-stopping promises a fix — yet two problems block it in practice.

PROBLEM 1

### Thresholds are hard to set

“Stop when confident” just swaps a token budget for a confidence threshold — a trickier knob. Its value is uninterpretable and shifts with every signal, model, and task.

PROBLEM 2

### Budget wasted on the unsolvable

On impossible problems, confidence never rises — so an upper-threshold-only rule never fires. The model burns the *entire* budget getting nowhere.

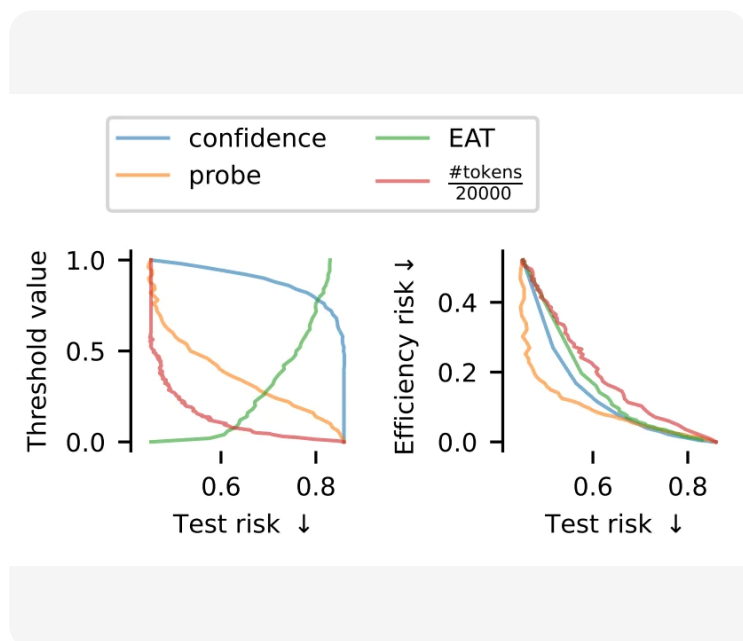


Fig 1 · Threshold for a target risk is signal-dependent — no single value transfers.

## 02 Risk control

Every early stop risks an error, so we control that risk directly. The user names a **risk tolerance  $\epsilon$**  — a target error rate — and a validation set does the rest.

$\epsilon$  is interpretable — it interfaces with downstream decisions; a raw threshold is not, so we calibrate it away.

### Calibration $\epsilon \rightarrow$ thresholds Alg. 1

- 1 Enumerate signal–threshold candidates on a held-out validation set.
  - 2 Keep those with risk  $\leq \epsilon$  under a finite-sample (UCB) correction, so the guarantee transfers to unseen data.
  - 3 Among survivors, pick the lowest efficiency loss — the most tokens saved.
- A stopping rule with a high-probability risk guarantee.

One risk target  $\epsilon$  in  $\rightarrow$  calibrated thresholds out, with a guarantee. No opaque knobs to hand-tune.

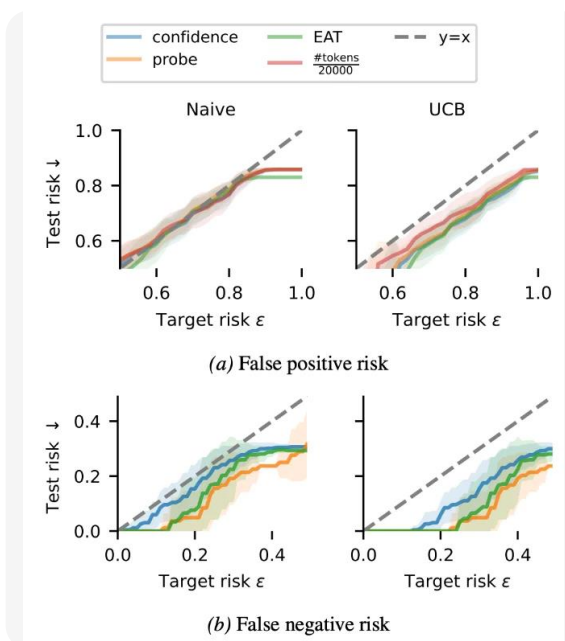
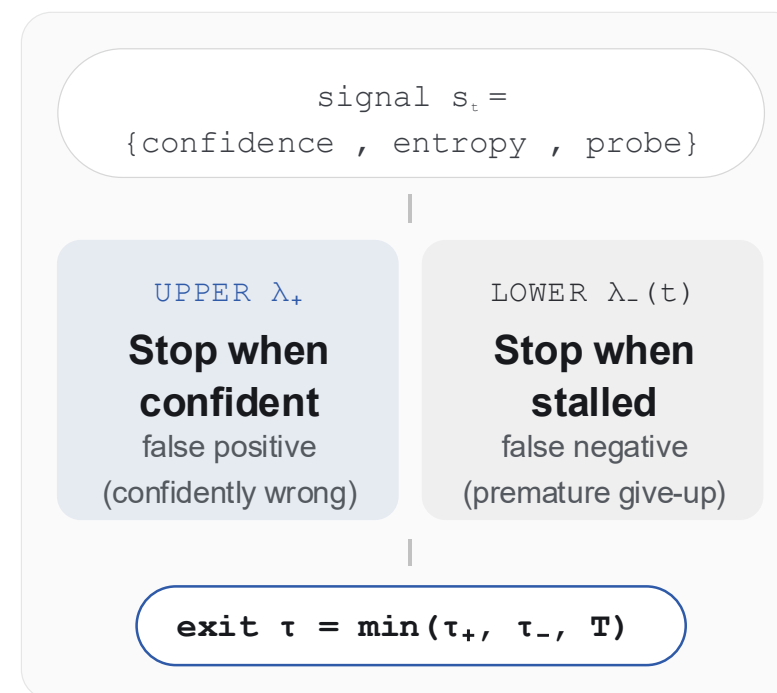


Fig 4 · UCB keeps realized risk  $\leq$  target; Naive violates it.

## 03 Dual thresholds

**Two risks, two thresholds.** One catches confident success; a **novel lower threshold** catches confident failure.



The lower threshold rises with token use — halting runs whose confidence fails to keep pace, saving the budget on the unsolvable.

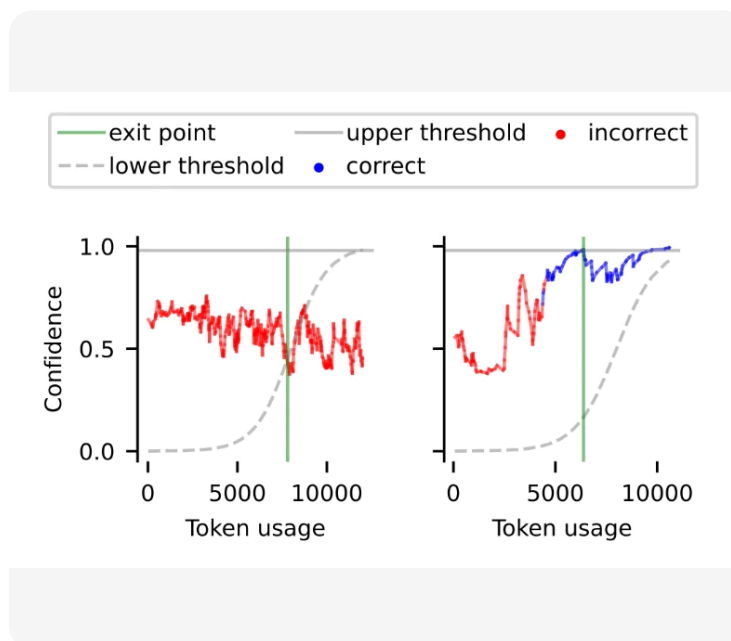


Fig 2 · Unsolvable runs exit via the lower threshold; solvable runs cross the upper.

## 04 Results

### risk $\leq \epsilon$

Test risk stays under target with high probability; Naïve calibration overfits and breaks it (Fig 4).

### ↓ tokens

Same accuracy, fewer tokens — gains grow as unsolvable problems dominate (Fig 6).

- **Division of labor:** solvable instances exit via the upper threshold, unsolvable via the lower.
- **Ensembling signals** per  $\epsilon$  automatically picks the most efficient stop.
- **Robust** across validation-set size and distribution shift.

MODELS Qwen3-8B Qwen3-30B-A3B  
R1-Distill-32B Qwen3-VL-8B  
 DATA AIME DeepScaleR GPQA-Diamond  
MathVision

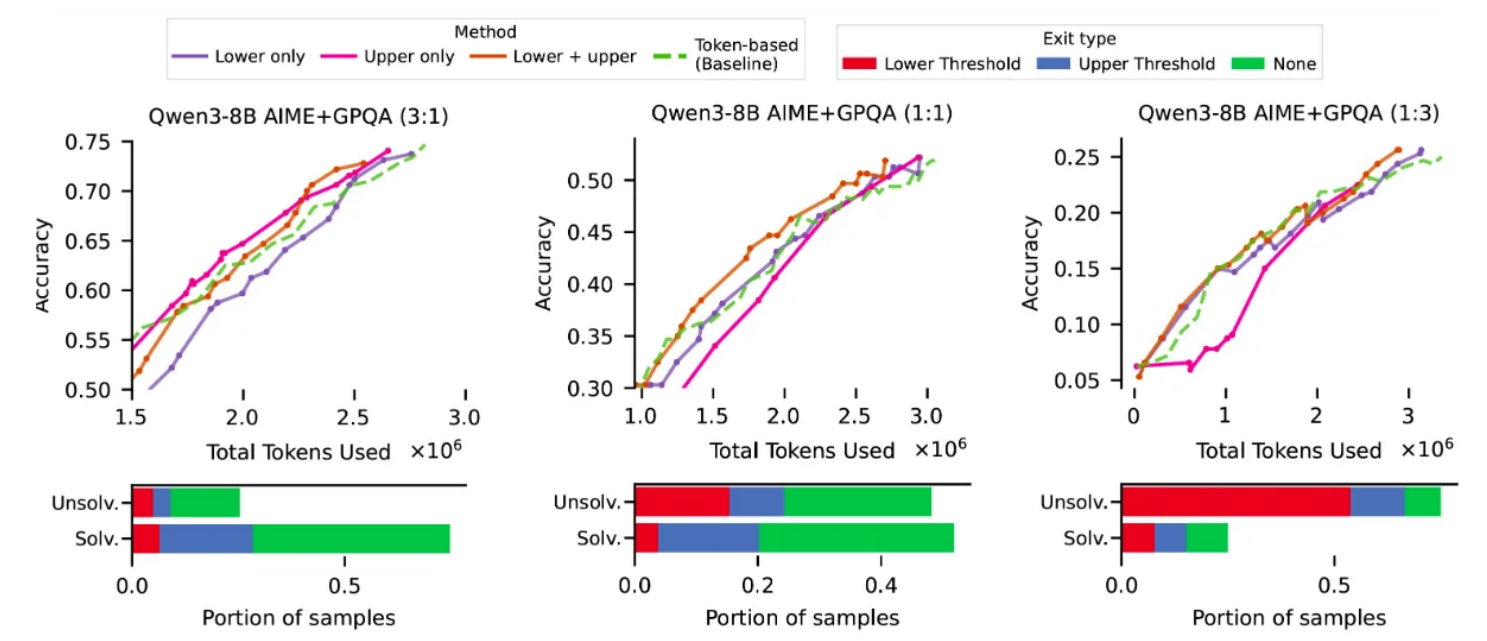


Fig 6 · The lower threshold shifts the accuracy–token curve left when unsolvable instances are common (1:1, 1:3).