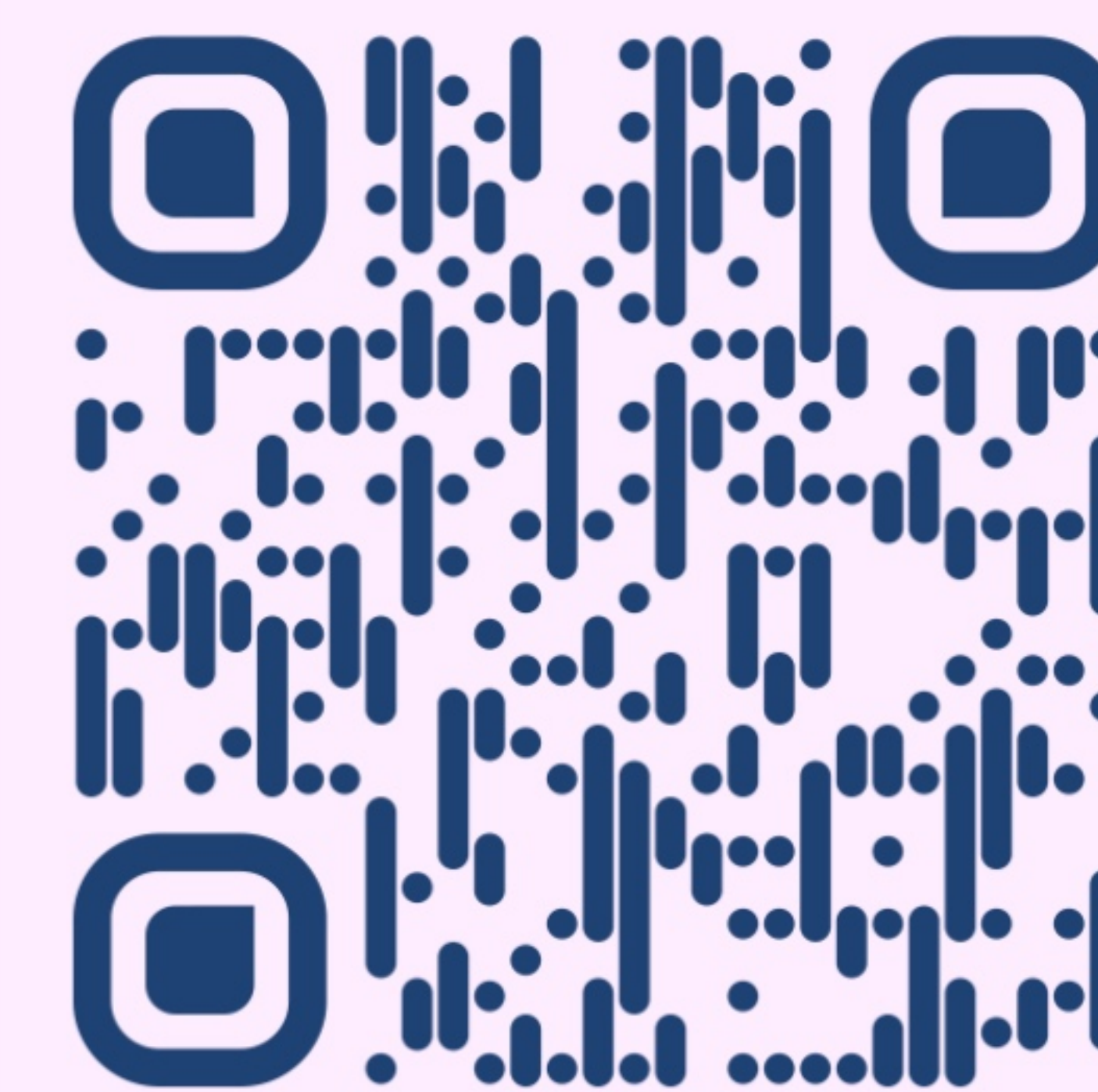


# EAGer: Entropy-Aware Generation for Adaptive Inference-Time Scaling

Daniel Scalena<sup>1,2</sup>, Leonidas Zotos<sup>2</sup>, Elisabetta Fersini<sup>1</sup>, Malvina Nissim<sup>2</sup>, Ahmet Üstün<sup>3</sup>



## Departure: Redundancy

Test-time scaling samples  $M$  parallel reasoning chains per prompt regardless of difficulty, wasting compute on easy problems while starving hard ones of exploration.

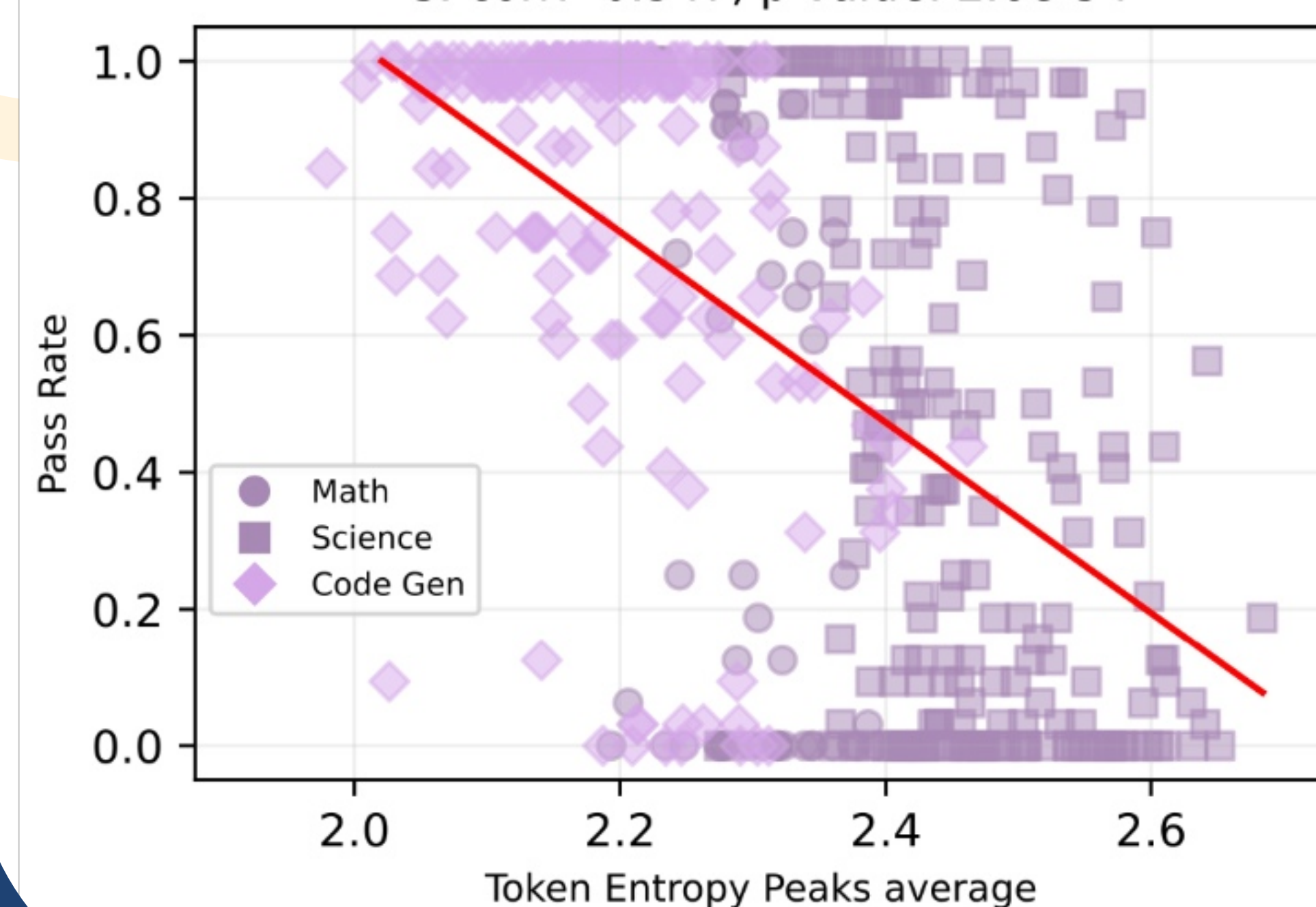
- CoT reasoning is dominated by predictable, low-information tokens;
- Parallel sampling re-generates near-identical prefixes across chains;
- Fixed budget  $M$  ignores that prompt difficulty varies widely.

## Entropy Peak Lookout

Token-level entropy spikes during generation correlate with lower Pass Rate, making peak entropy a cheap, real-time proxy for problem difficulty.

- Average peak entropy (99.9th percentile) negatively correlates with Pass Rate ( $\rho \approx -0.55$ ,  $p < 1e-30$ );
- High-entropy moments  $\rightarrow$  genuine decision points worth exploring;
- Low entropy  $\rightarrow$  redundant continuation.

Pass Rate VS Token Entropy Peaks average  
S. corr: -0.547, p-value: 2.6e-34



## Token Savings Platform

EAGer-init alone uses less than half the tokens of Full Parallel sampling across all models and benchmarks, before any performance gain is even counted.

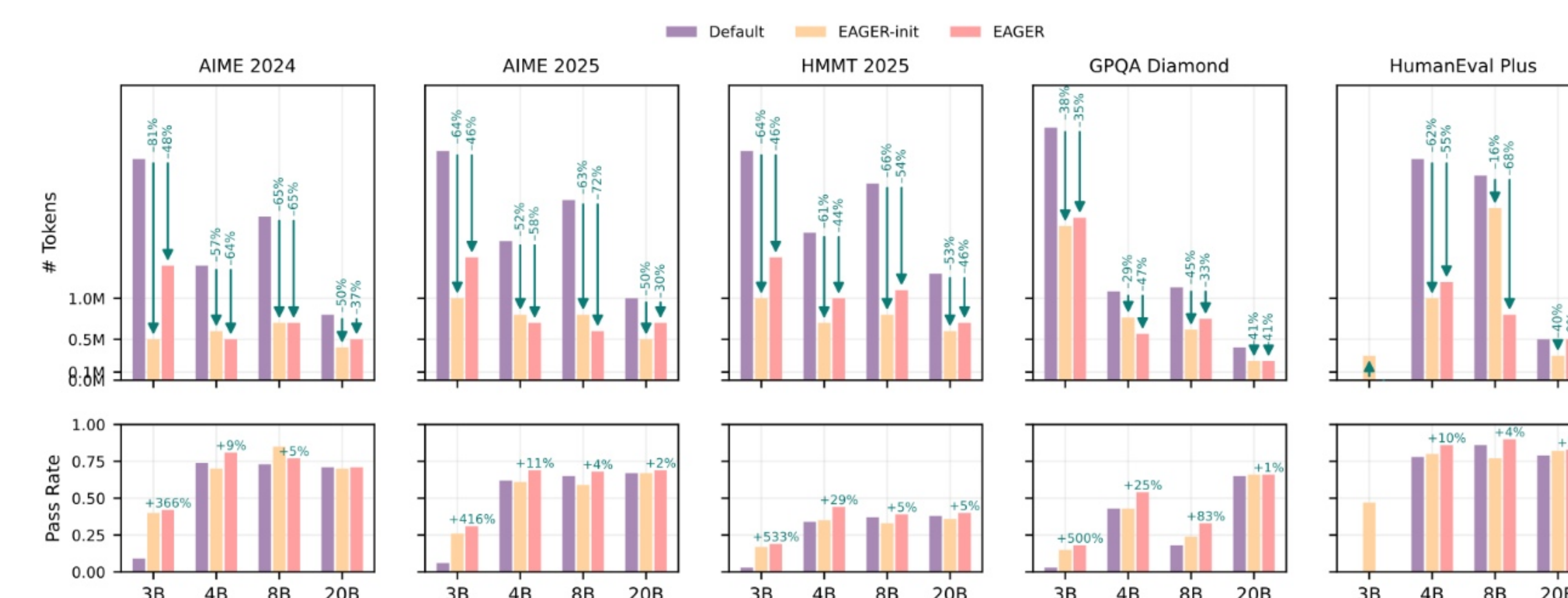


Figure 3. Compute and performance trade-offs of EAGer-init and EAGer. Across all benchmarks and model size, the efficiency of EAGer-init and EAGer consistently outperforms FULL PARALLEL sampling, requiring only half as many tokens in most cases (top). In addition, they achieve higher pass rate accuracy (bottom). For issues specific to the smallest 3B model, see Appendix D.

## Accuracy Terminal

Pass Rate, Pass@k, and Cons@k:  
• up to +37% Pass@k with labels,  
• +12% without,  
while cutting tokens by up to 64%.

## Benchmark Depot

Evaluated across 4 open models (3B–20B) and 5 reasoning benchmarks spanning math, science, and code.

## Reallocation Interchange

Unused budget from easy prompts (capped at 2M) is redirected to harder prompts — either saturating prompts (label-free: EAGer-adapt) or prompts with Pass@k=0 (label-available: full EAGer).

- Total token usage still stays below fixed-budget FULL PARALLEL baseline.

## Pareto Express

As budget  $M$  grows, EAGer shifts the entire Pareto frontier outward.

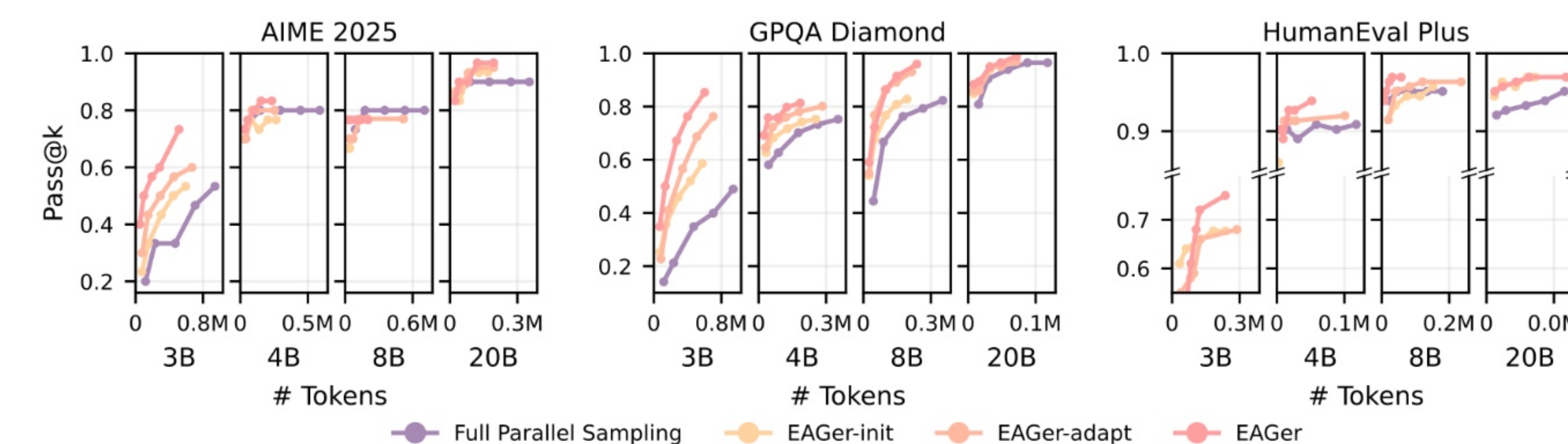


Figure 4. Performance comparison with scaling the total allowed sequences for generating ( $M \in \{1, 4, 8, 16, 24, 32\}$ ). As  $M$  increases (line's markers), EAGer consistently improves Pass@k (y-axis) while reducing the number of tokens needed to find the correct solution (x-axis), further shifting the Pareto frontier of the performance–efficiency trade-off.

## Terminus

EAGer consistently beats Full Parallel sampling requires no fine-tuning, generalizes across model families/sizes and task domains, and needs only a single threshold  $\theta$  calibrated from 10 examples.

## The Branch Switch

Branch a sequence into two candidates only when token entropy exceeds threshold  $\theta$ .

- Branch = pick top-2 most likely tokens, spawn new candidate continuation;
- Shared prefix tokens are never regenerated  $\rightarrow$  direct token savings.

