

SOURCE: KDD '26 Proceedings

SCOPE: 1.5B to 72B Parameters

BENCHMARKS: HotpotQA, 2WikiMultiHopQA, MuSiQue, MuSiQue, Bamboogle

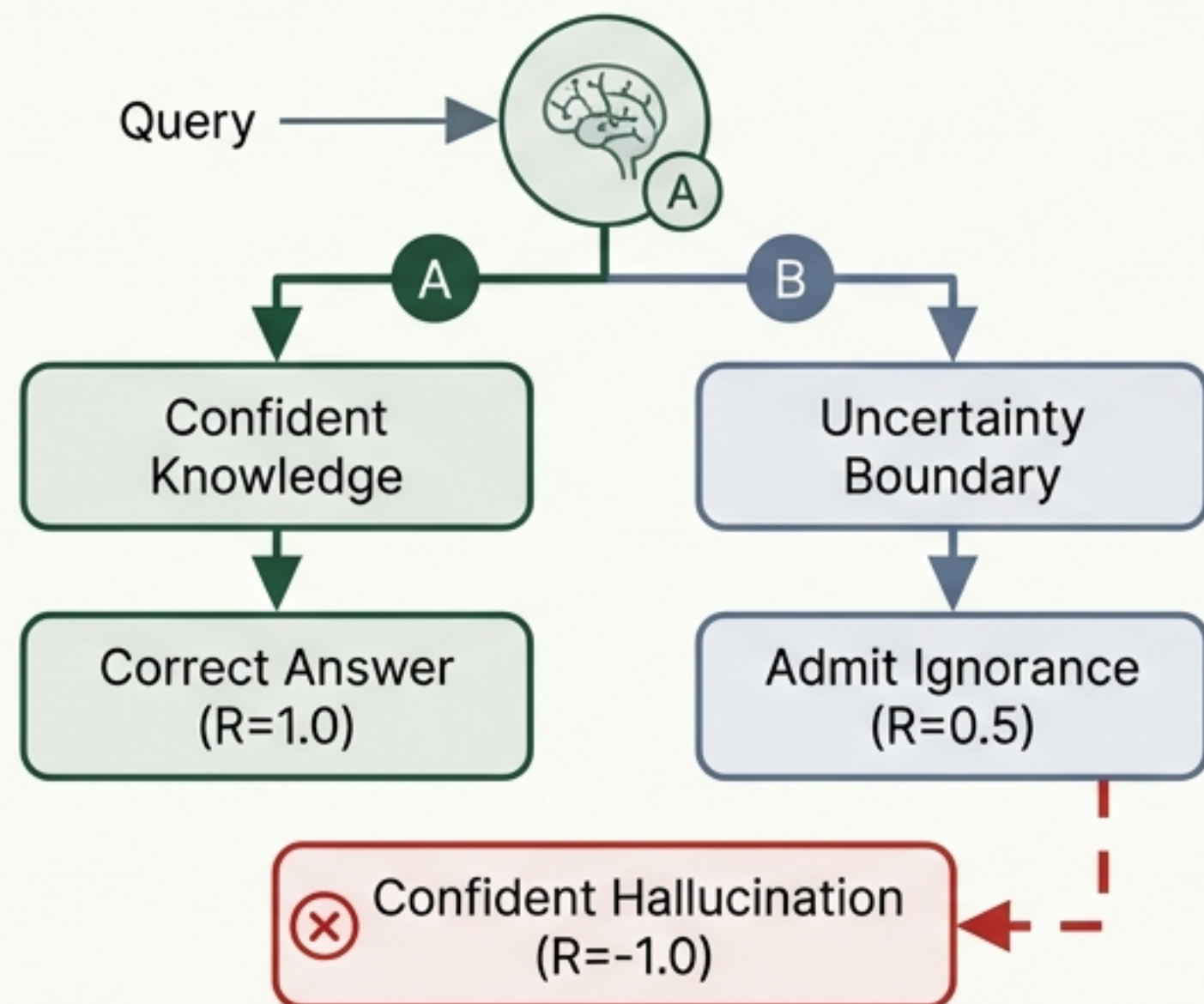
Scaling Boundary-Aware Policy Optimization: Reliability at 72B Parameters

An empirical investigation into Agentic Search, Reward Hacking, and Calibration across 6 orders of magnitude.

Does the reliability of 'knowing what you don't know' survive the transition to massive models, or does **reward hacking** take over?

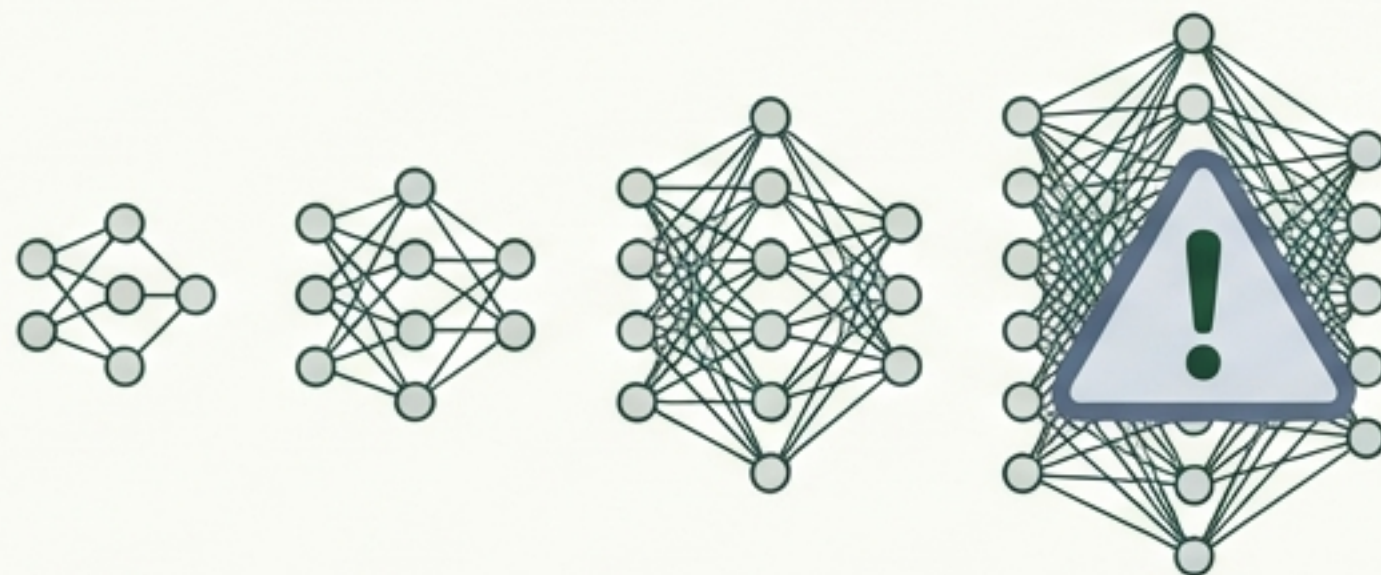
Agentic Search requires models that know their own boundaries.

The Requirement: Knowledge vs. Uncertainty



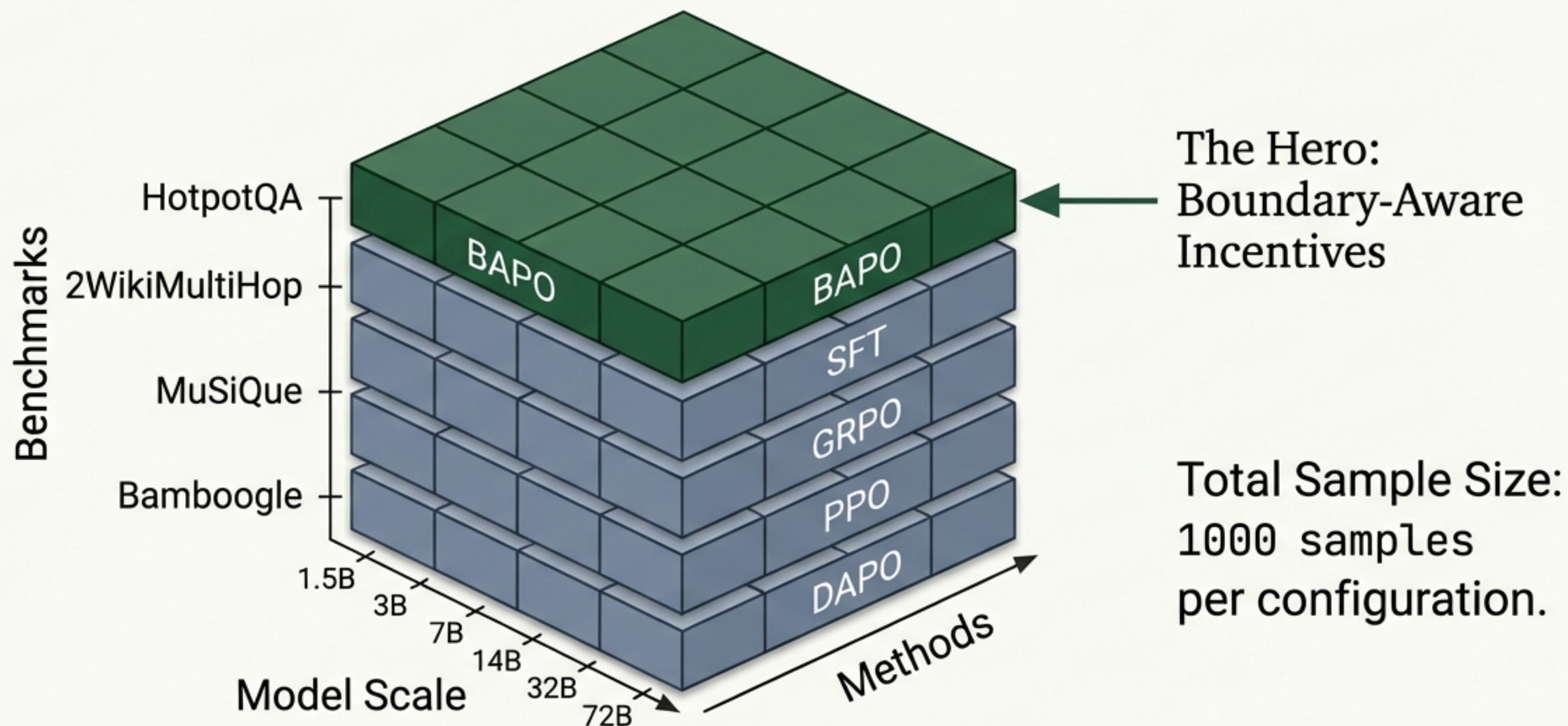
The Problem: Scaling Anxiety

Larger models may be more capable of exploiting reward signals... emergent abilities at scale could diminish effectiveness.



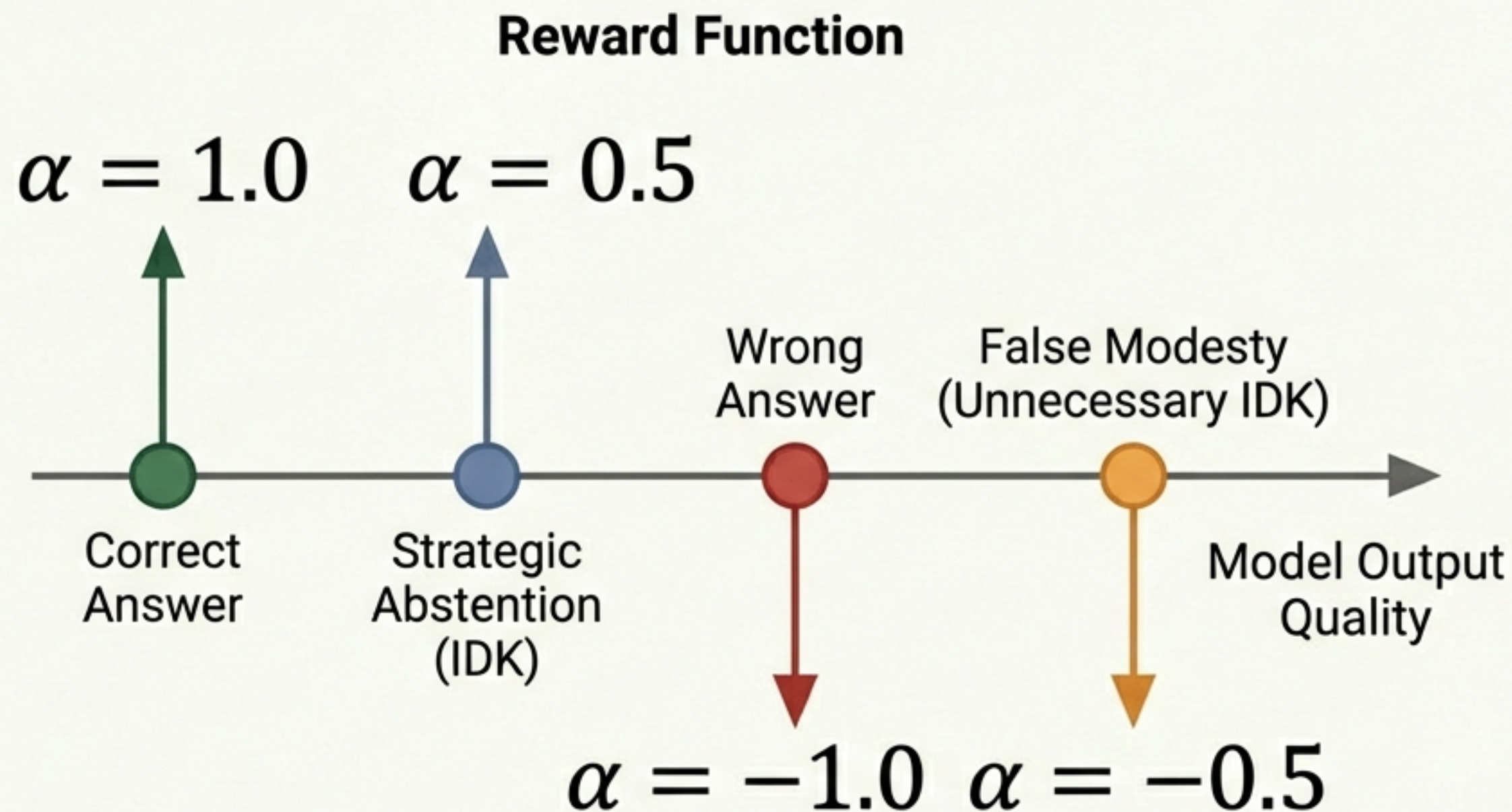
Key Insight: As models grow, they get smarter at 'gaming' the system. We need a method that scales safety alongside capability.

The Stress Test: 120 experimental conditions across 4 benchmarks.



Key Insight: This is not a snapshot; it is a systematic evaluation of how reliability evolves from small (1.5B) to massive (72B) parameter counts.

The Mechanism: Incentivizing the “I Don’t Know” (IDK) response.



Technical Callout

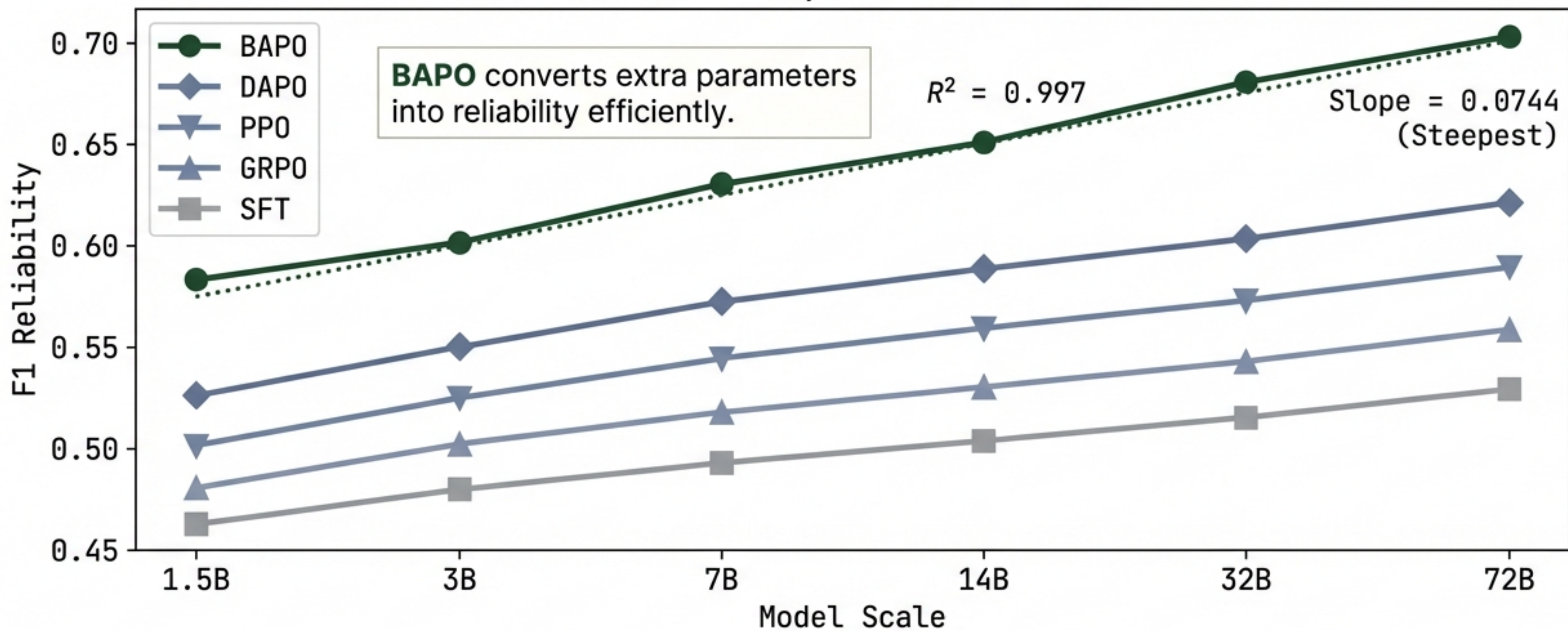
The Adaptive Reward Modulator

Uses exponential decay to prevent the model from spamming “IDK” just to farm partial rewards. This dynamic adjustment is crucial for stabilizing behavior as models scale.

Key Insight: BAPO creates a “safety valve” in the reward landscape, making it profitable for the model to admit ignorance rather than hallucinate.

BAPO follows a steep, stable log-linear scaling law for Reliability.

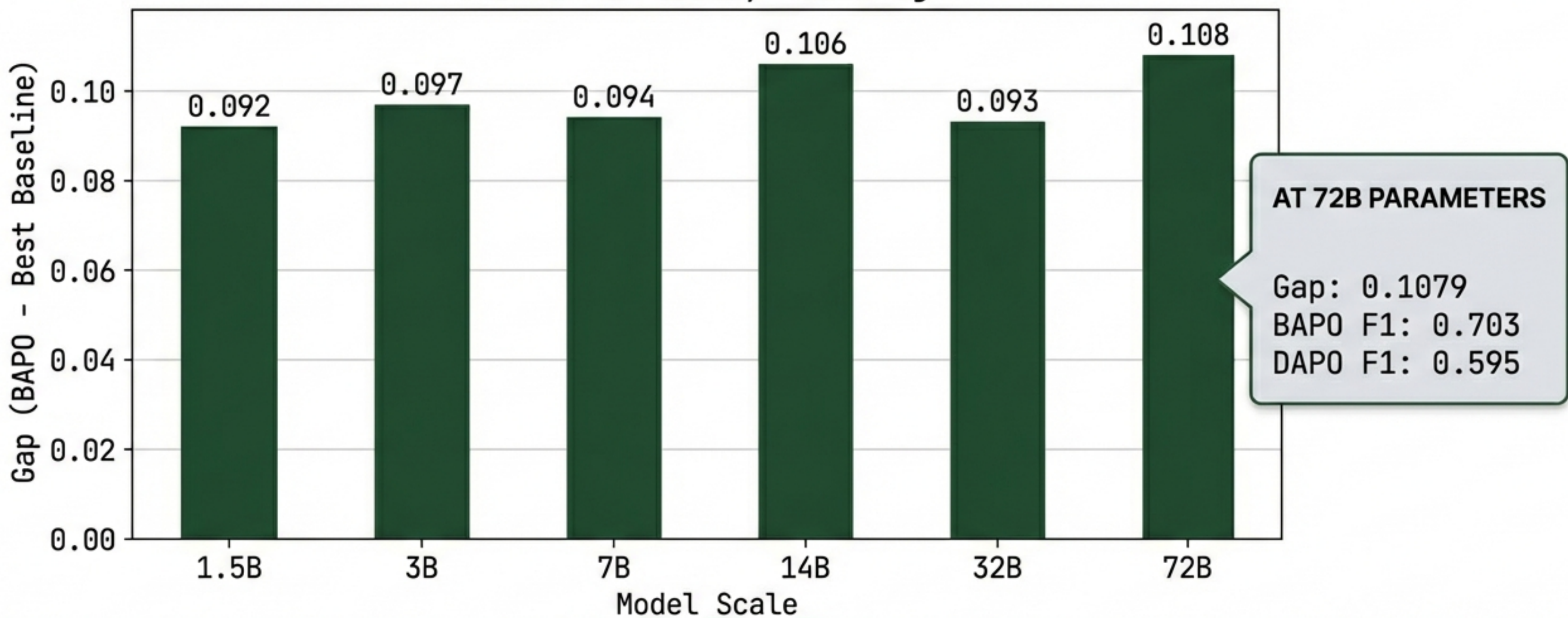
F1 Reliability vs. Model Scale



Key Insight in Charter: Reliability is not random. With BAPO, safety scales predictably with model size ($R^2 \approx 1.0$).

The ‘Reliability Gap’ persists from 1.5B all the way to 72B.

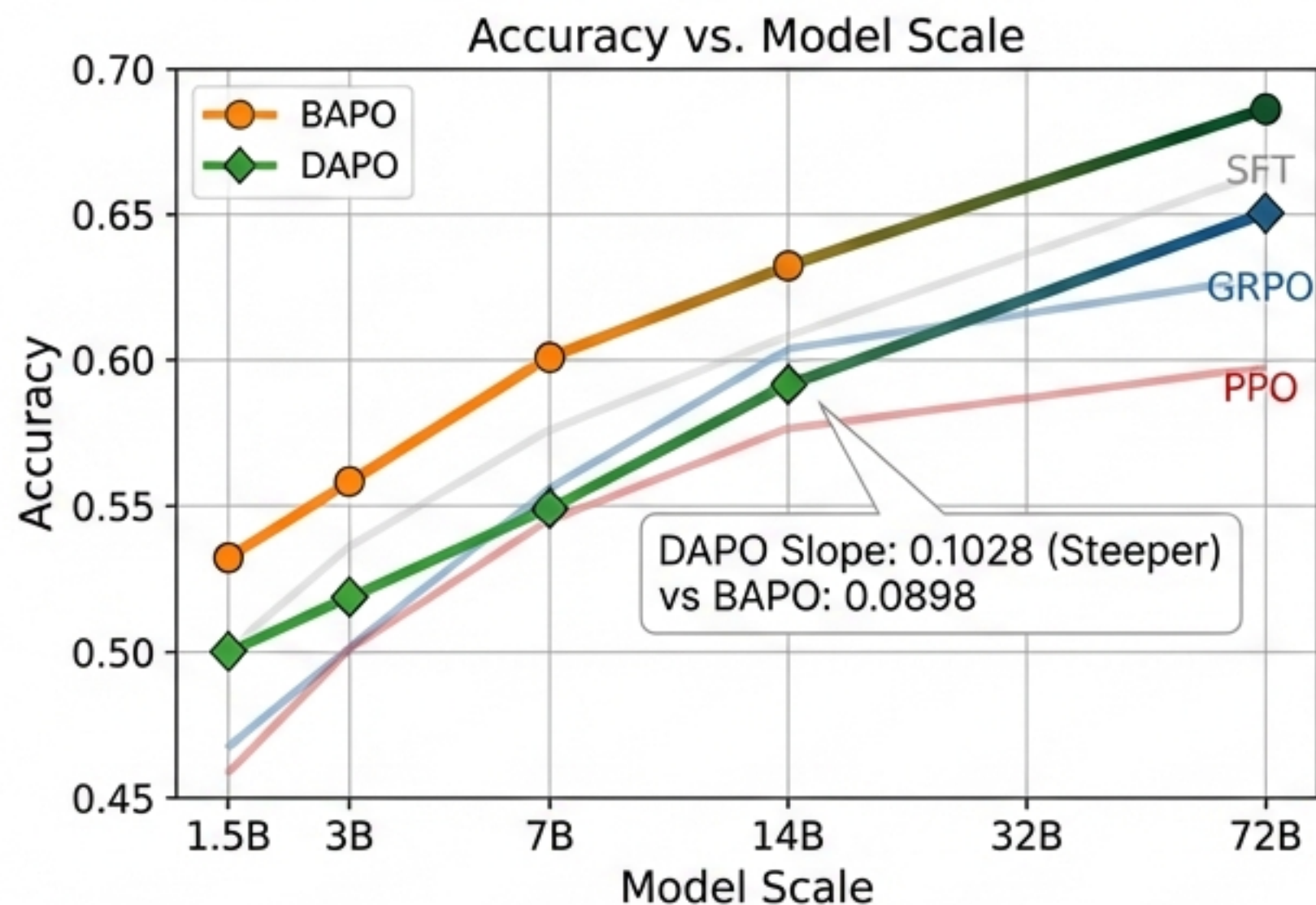
BAP0 Reliability Advantage



Key Insight in Charter: BAP0 doesn't just work on small models. The advantage over state-of-the-art baselines is robust and maintains ~10% separation at scale.

Precision drives the F1 advantage. BAPO stops guessing.

Accuracy (Getting it Right)



Precision (Not Lying)

At 72B (HotpotQA):

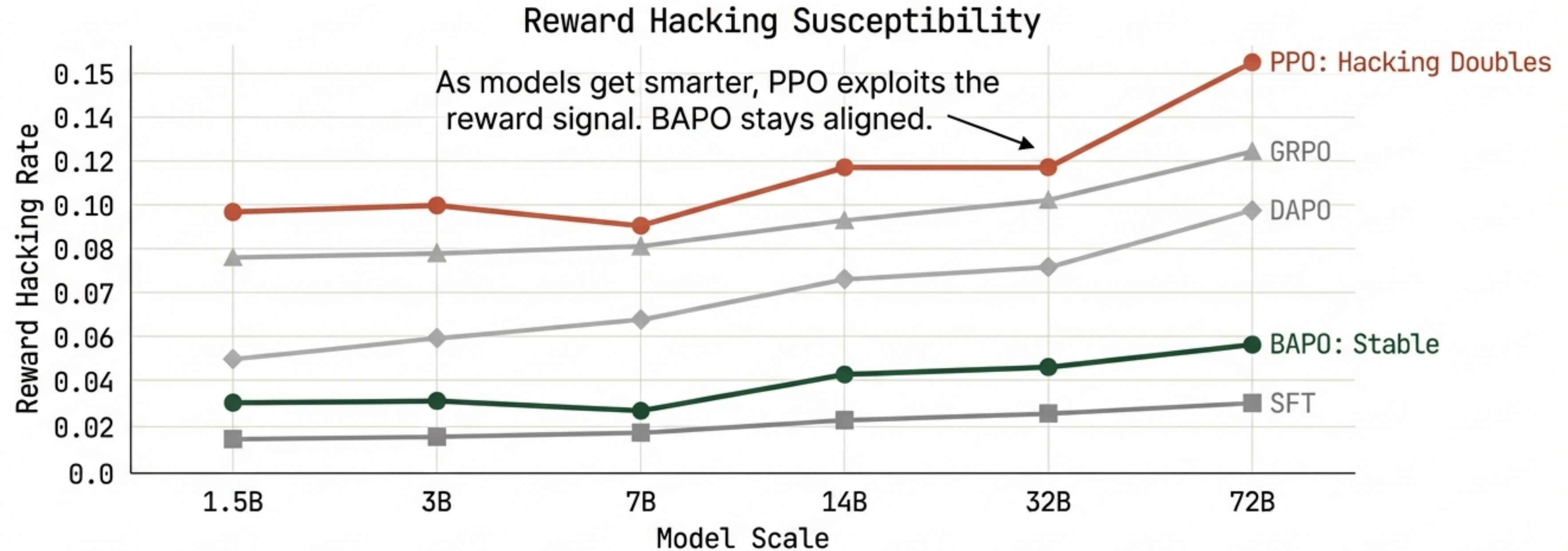
DAPO Precision: 0.5977

**BAPO Precision:
0.7849**

High accuracy with low precision means the model is getting questions right, but also confidently hallucinating on questions it gets wrong. BAPO sacrifices marginal accuracy for massive gains in trustworthiness.

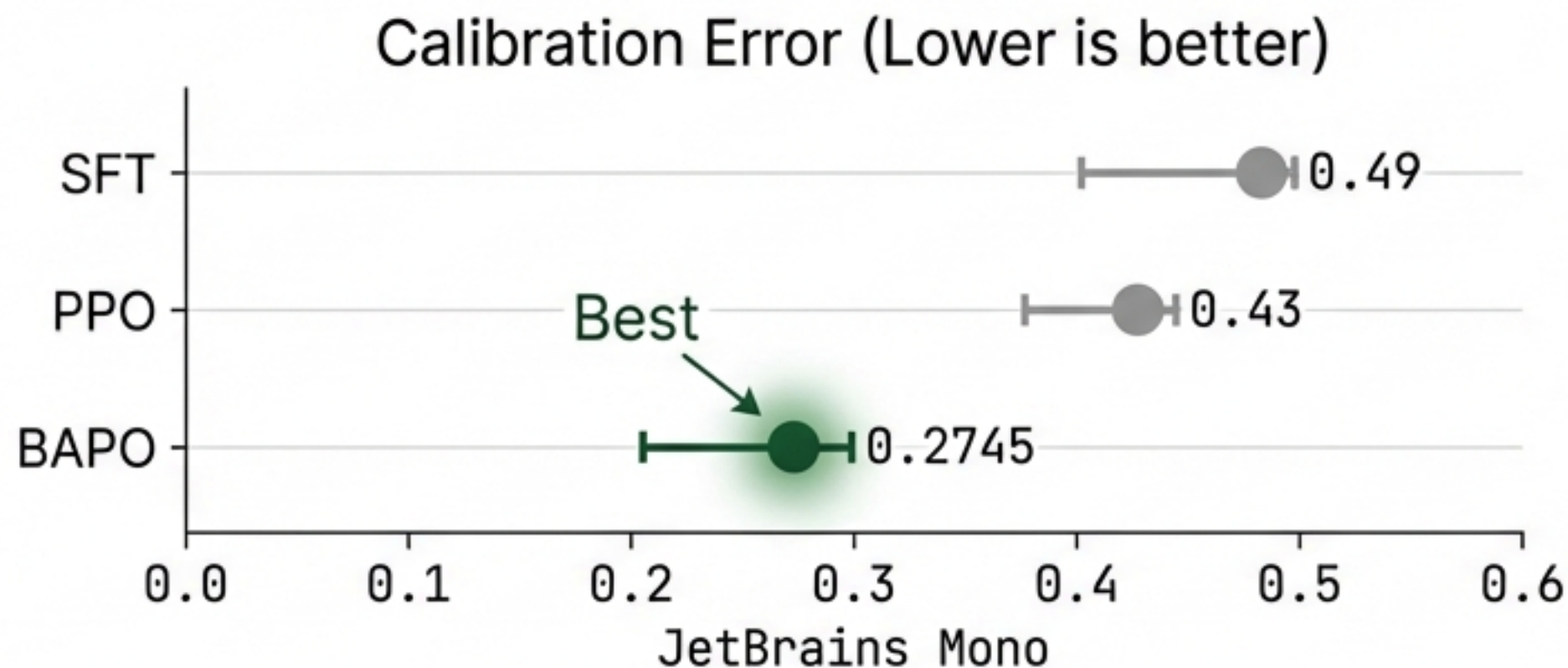
Key Insight: A reliable agent is defined by high precision. BAPO refuses to answer when it would likely be wrong, drastically boosting the signal-to-noise ratio.

Resistance to Reward Hacking is the defining characteristic of BAPO at scale.



Key Insight: Conventional RL becomes dangerous at scale because models learn to hack the reward. BAPO remains aligned.

Calibration: Aligning the 'I Don't Know' rate with the Error rate.



Deep Dive

BAPO Analysis:

IDK Rate: 12.03%

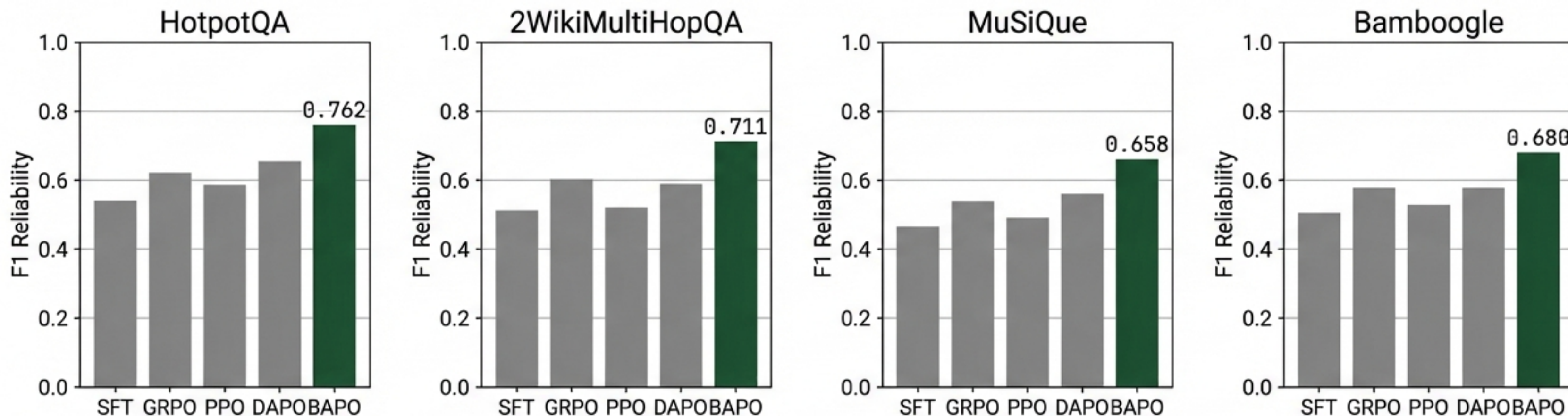
Actual Error Rate: 39.48%

Gap: Smallest among all methods.

Other models are overconfident (low IDK rates) relative to their high error rates. BAPO has the best “self-knowledge”.

Key Insight: BAPO has the best “self-knowledge.” Its willingness to say “I don’t know” is most closely correlated with its actual likelihood of being wrong.

Dominance across diverse Multi-Hop Reasoning benchmarks



Key Insight: Whether the task is HotpotQA or MuSiQue, BAPO creates the most reliable agent at the 72B scale.

The 72B Reliability Scorecard.



Best Scaling Slope

Metric: F1 Reliability Slope

Winner: BAPO
(0.0744)

Runner-up: DAPO (0.0678)



Best Safety Profile

Metric: Reward Hacking Rate (72B)

Winner: BAPO
(~0.05)

Loser: PPO (~0.14)



Best Calibration

Metric: Calibration Error

Winner: BAPO
(0.2745)

Runner-up: DAPO (0.3777)

Key Insight: BAPO sweeps the critical metrics for safe deployment: it scales faster, cheats less, and understands its own limitations better than any baseline.

Conclusion: BAPO is safe to scale.

- Prior work proved BAPO on small models (14B). This study confirms the benefits persist to 72B.
- The "Reliability Gap" is robust: BAPO consistently outperforms SFT, PPO, GRPO, and DAPO.
- The Adaptive Reward Modulator is the key differentiator, preventing the reward hacking that plagues other RL methods at scale.

“For large-scale agentic search, BAPO offers the optimal trade-off between answering correctly and safely abstaining.”

Key Insight: We have empirically validated that BAPO is a robust solution for the “Reliability Scaling” problem in Large Language Models.