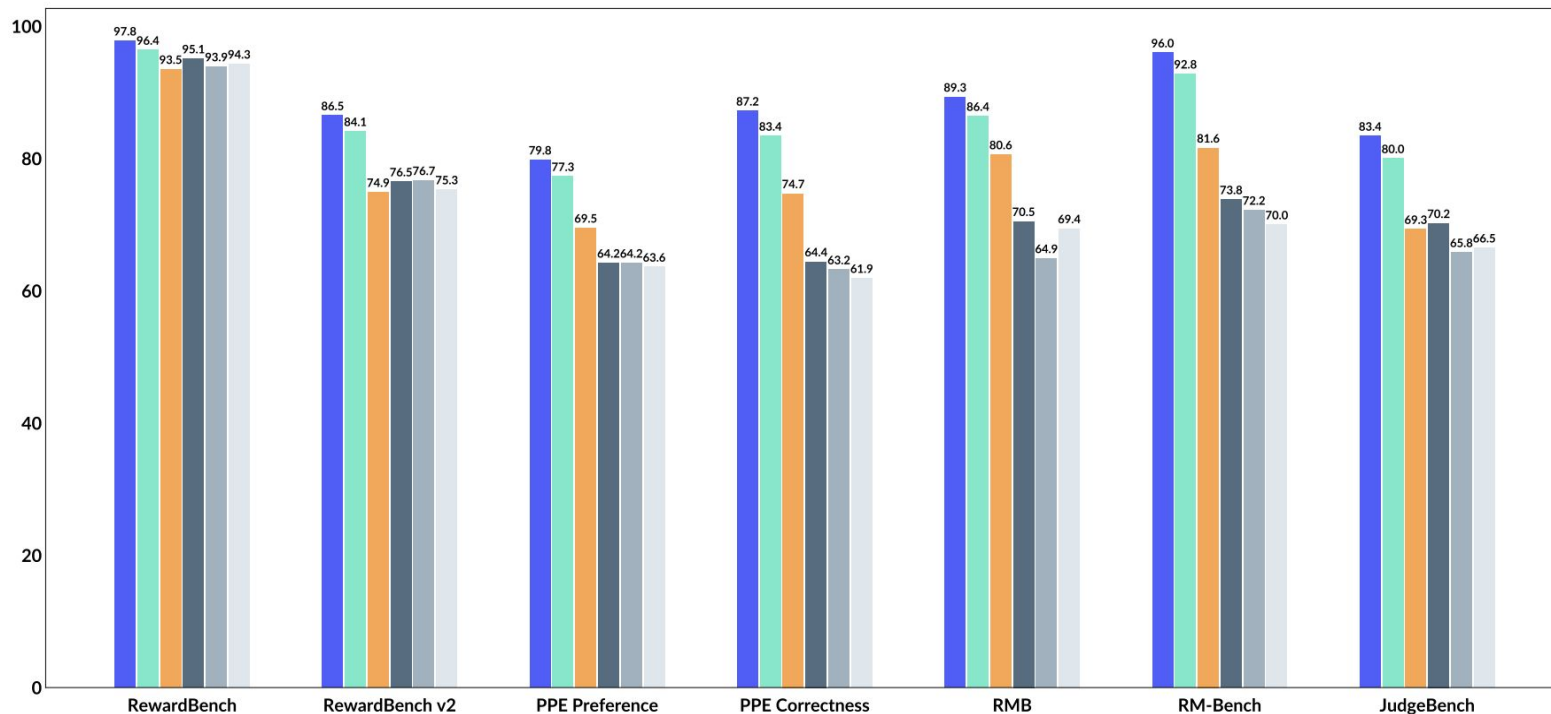# Skywork Reward v2

How should we see it?

Zhilin Wang
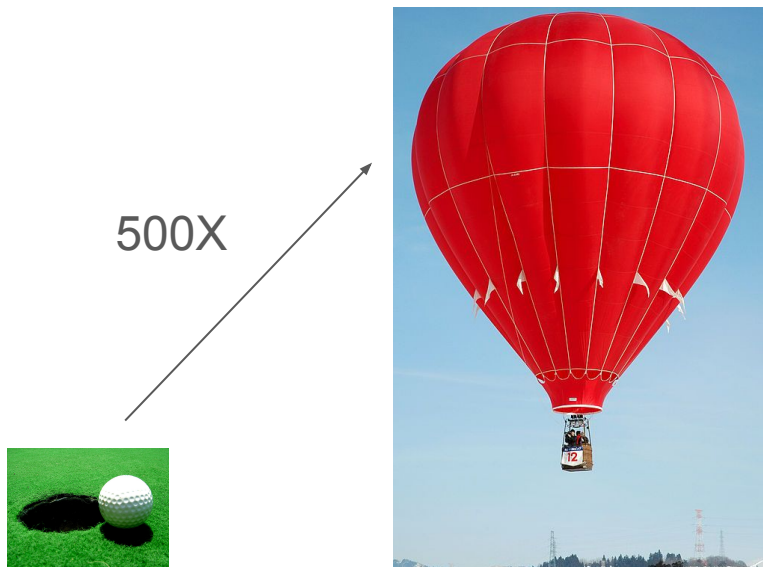
# It achieves SOTA in 7 distinct benchmarks



Notes: it's better than

1. 70B BT models
2. GenRMs
3. O3-mini

# What's the secret sauce?

500X



Data Scale (80 Thousand vs 40 Million)

Human label most informative samples,
Synthetic label everything else



Human-AI Collaboration for Curation

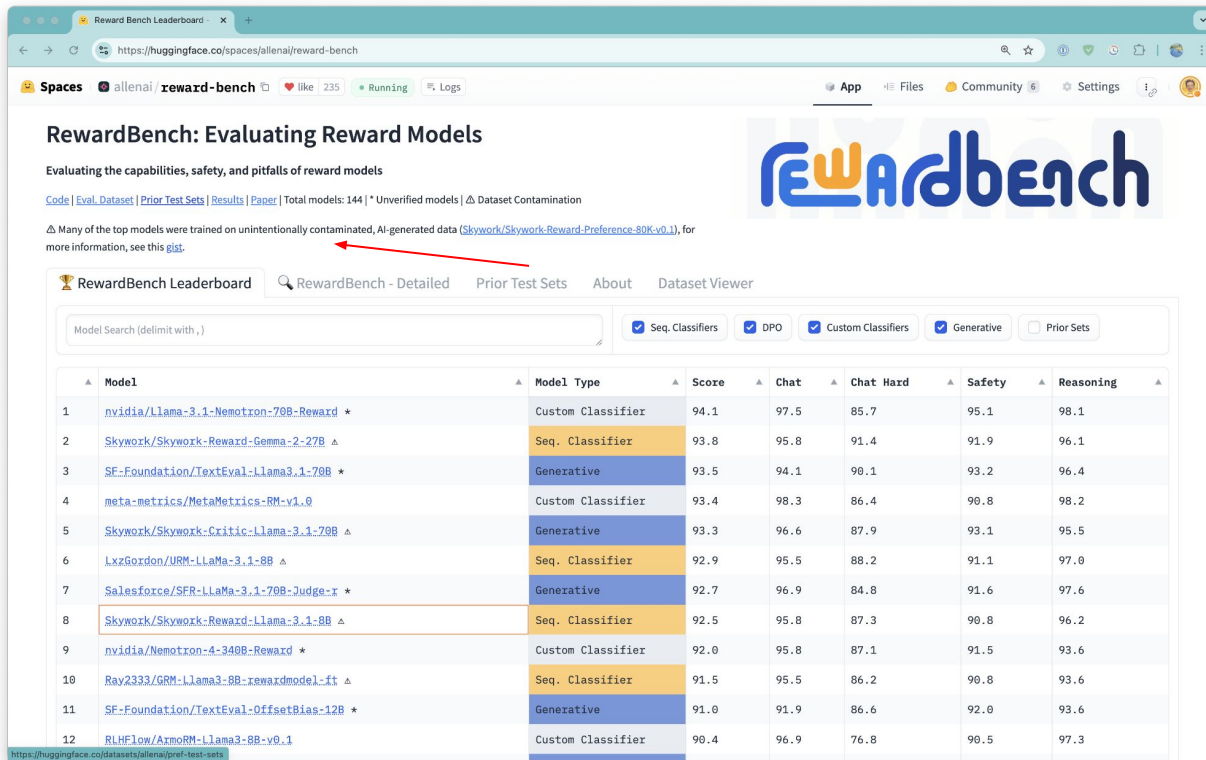# What made me think twice about the results?



No open data



Astronomically Large Gains



No RLHF Experiments

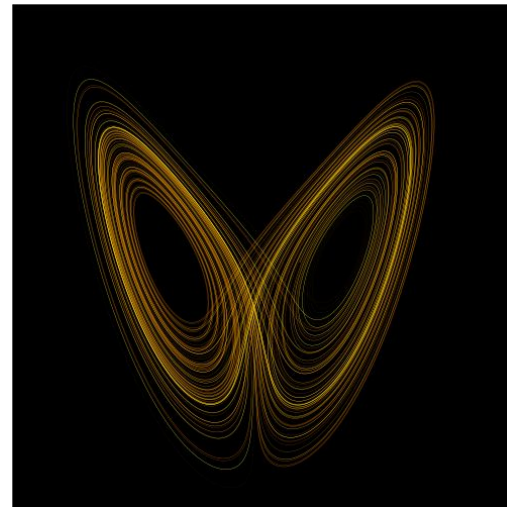# What happened with SkyWork Reward V1?

# Potential Hypothesis for why Skywork Reward V2 Rocks



Flywheel from V1



Contamination with Test Sets



Interaction w. Base Model

# Testing our hypothesis indirectly

| Model Name | Original | After removing 5 letter | Absolute Drop ↓ |
|---|---|---|---|
| Llama 3.1 8B v2-40M | 80.5 | 66.9 | 13.6 |
| Llama 3.1 8B v2 | 74.7 | 70.1 | 4.6 |
| Qwen3 8B v2 | 69.5 | 63.6 | 5.9 |
| Llama 3.1 8B v0.1 | 58.5 | 53.0 | 5.5 |
| Llama 3.1 8B v0.2 | 59.7 | 53.9 | 5.8 |

Artifact from JudgeBench-Knowledge Subset "Once you have your answer, please duplicate that letter five times in a single string. For example, if the answer is K, then write KKKKK."

# What does this mean for Open-Data Evaluation?



Models w/ Open Data
& Not Incentivised to
Train on Test Data

Models w/o Open Data
& Can Train on Test
Data w/o Detection

# What can we do about it?



Distinguish between Open Source & Open Weight

Have Private Test Data that cannot be Trained on

Evaluate Reward Models on Downstream Alignment Performance