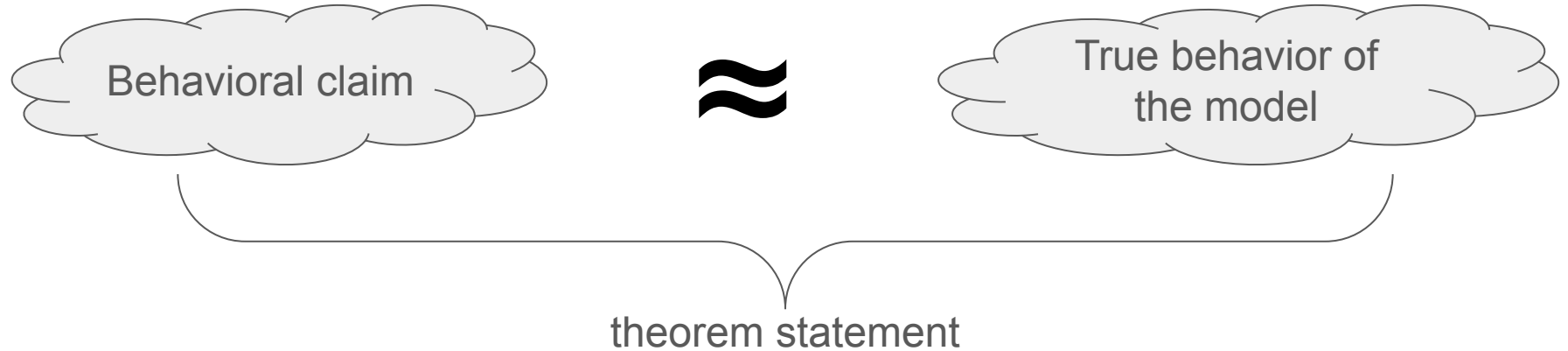


Compressing explanations

Compact Proofs of Model Performance via
Mechanistic Interpretability

*Jason Gross, Rajashree Agrawal, Thomas Kwa, Euan Ong, Chun Hei Yip,
Alex Gibson, Soufiane Noubir, Lawrence Chan*

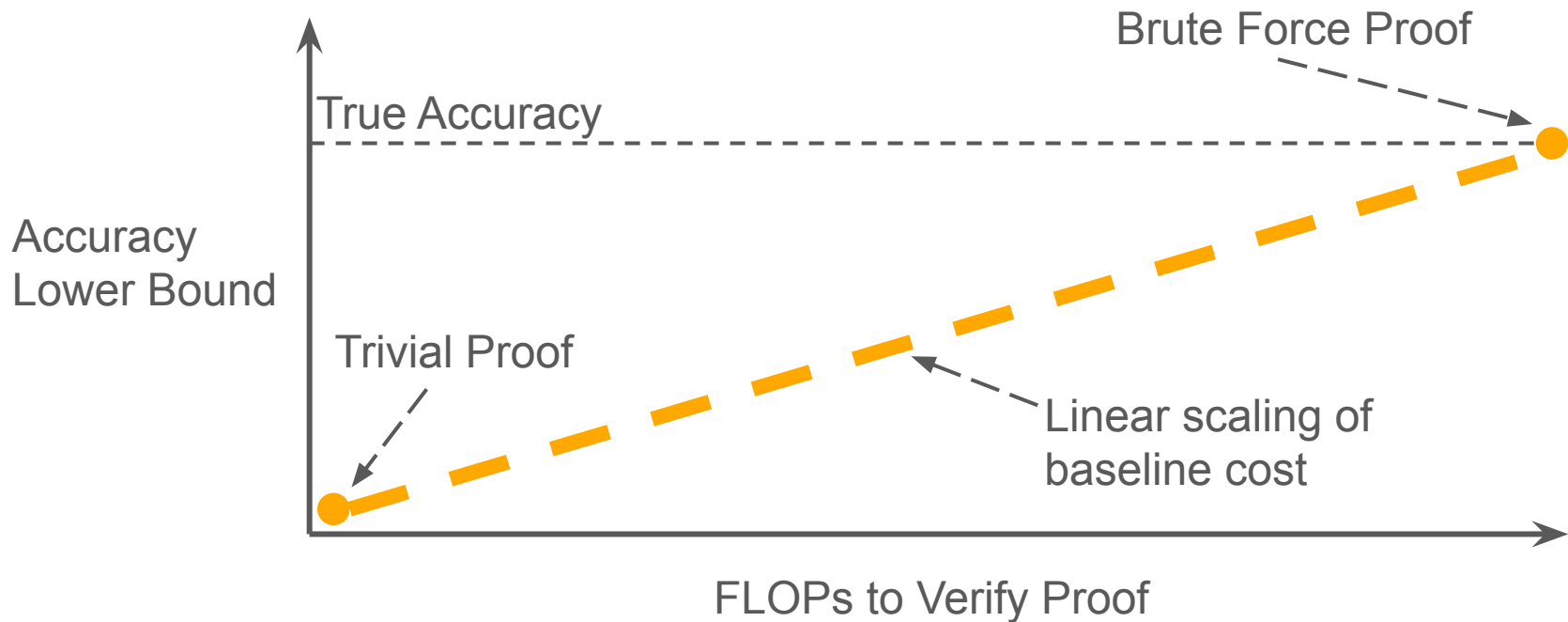
Formalizing proof length to quantify compression



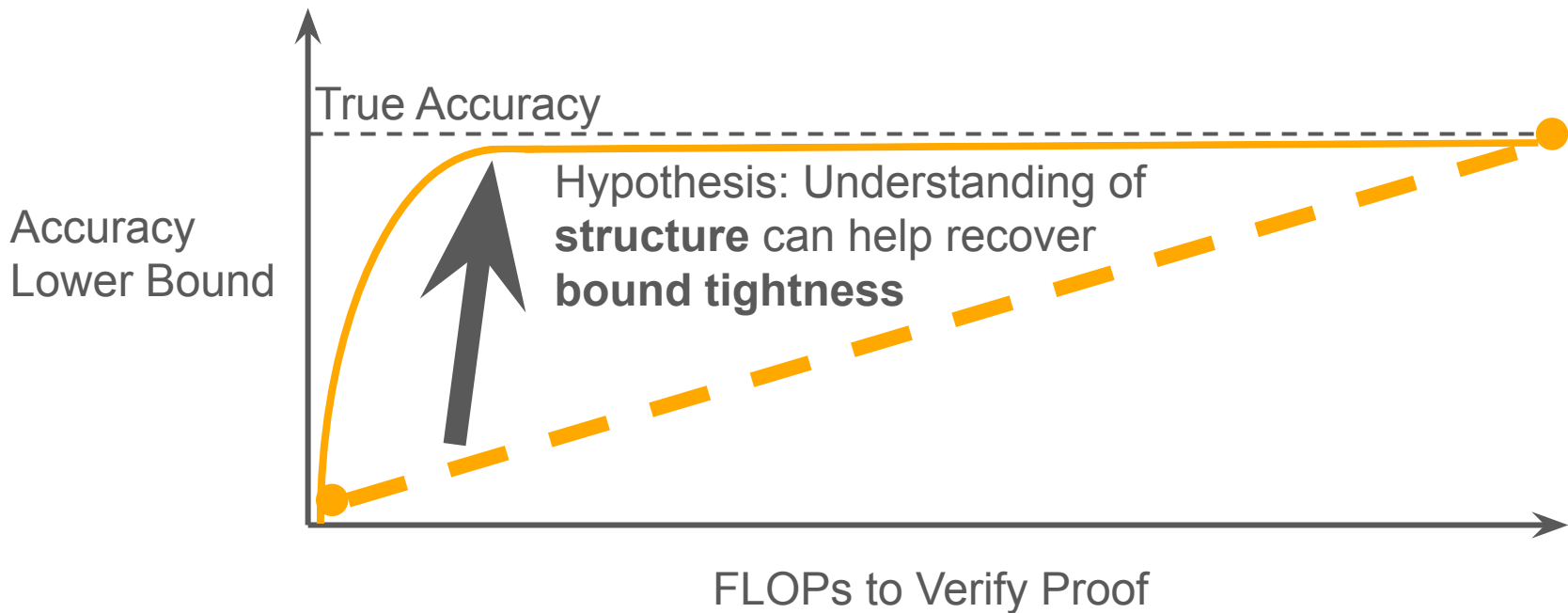
Proof = sound computation of worst-case error

Length of proof = cost of running computation

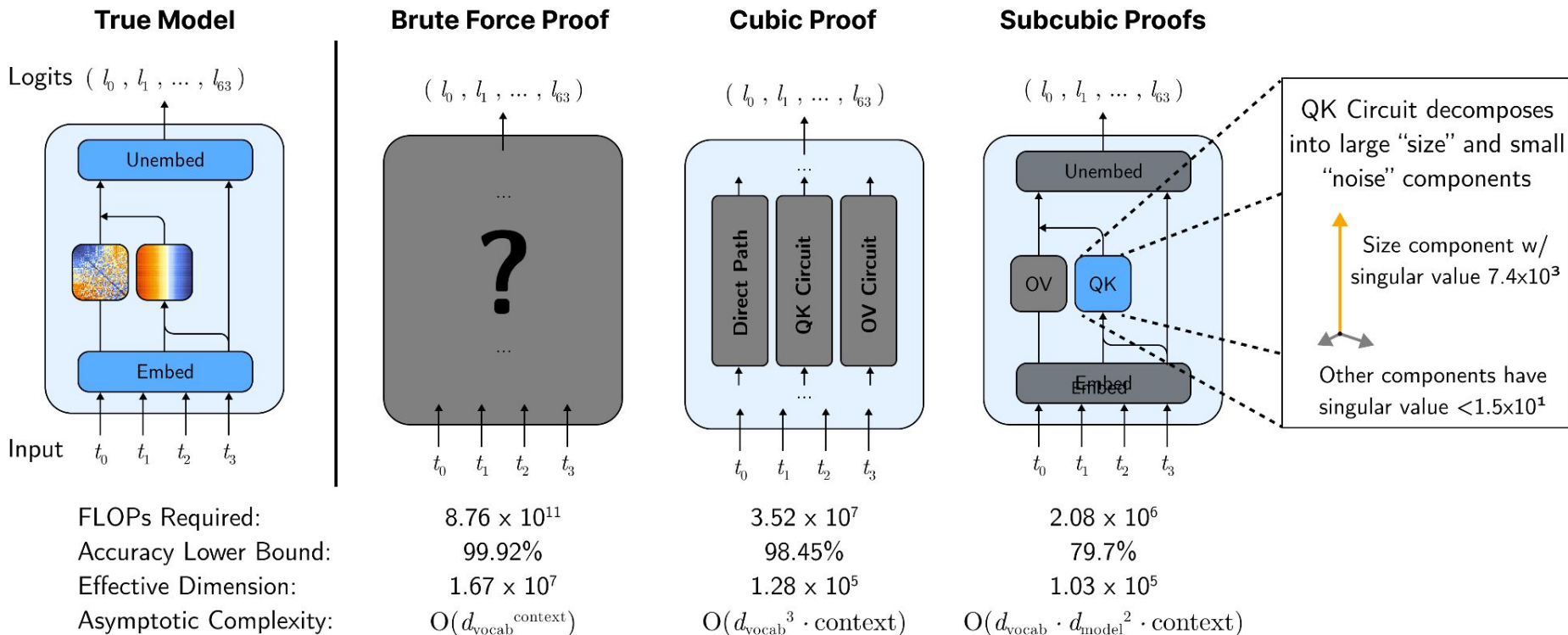
Quantifying the compute-cost of explanations



Does understanding improve upon the linear baseline?

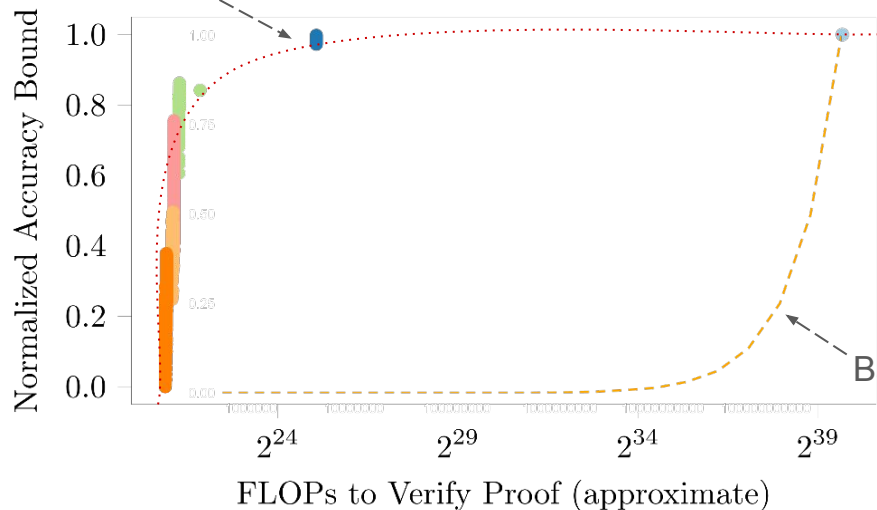


Proofs with varying mechanistic understanding



We found an empirical “pareto frontier”

Pareto frontier from
incorporating mechanistic
understanding

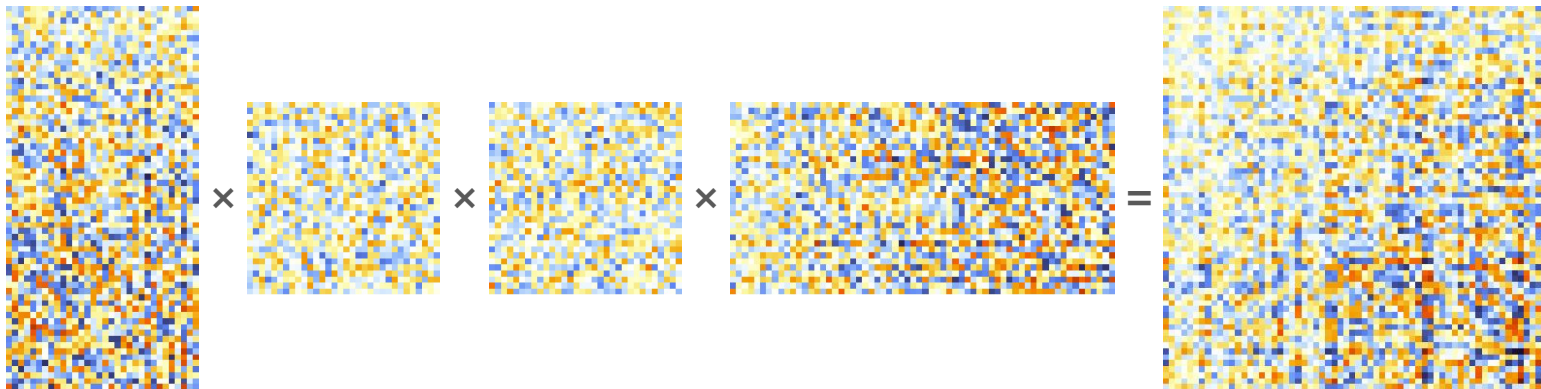


- brute force (acc: 0.9992 ± 0.0015)
- cubic (rel acc: 0.9853 ± 0.0038)
- subcubic (rel acc: 0.833 ± 0.011)
- attention- $d_{\text{vocab}}d_{\text{model}}^2$ (rel acc: 0.807 ± 0.013)
- direct-quadratic (rel acc: 0.663 ± 0.060)
- attention- $d_{\text{vocab}}d_{\text{model}}^2$, direct-quadratic (rel acc: 0.637 ± 0.060)

Baseline approach in log scale

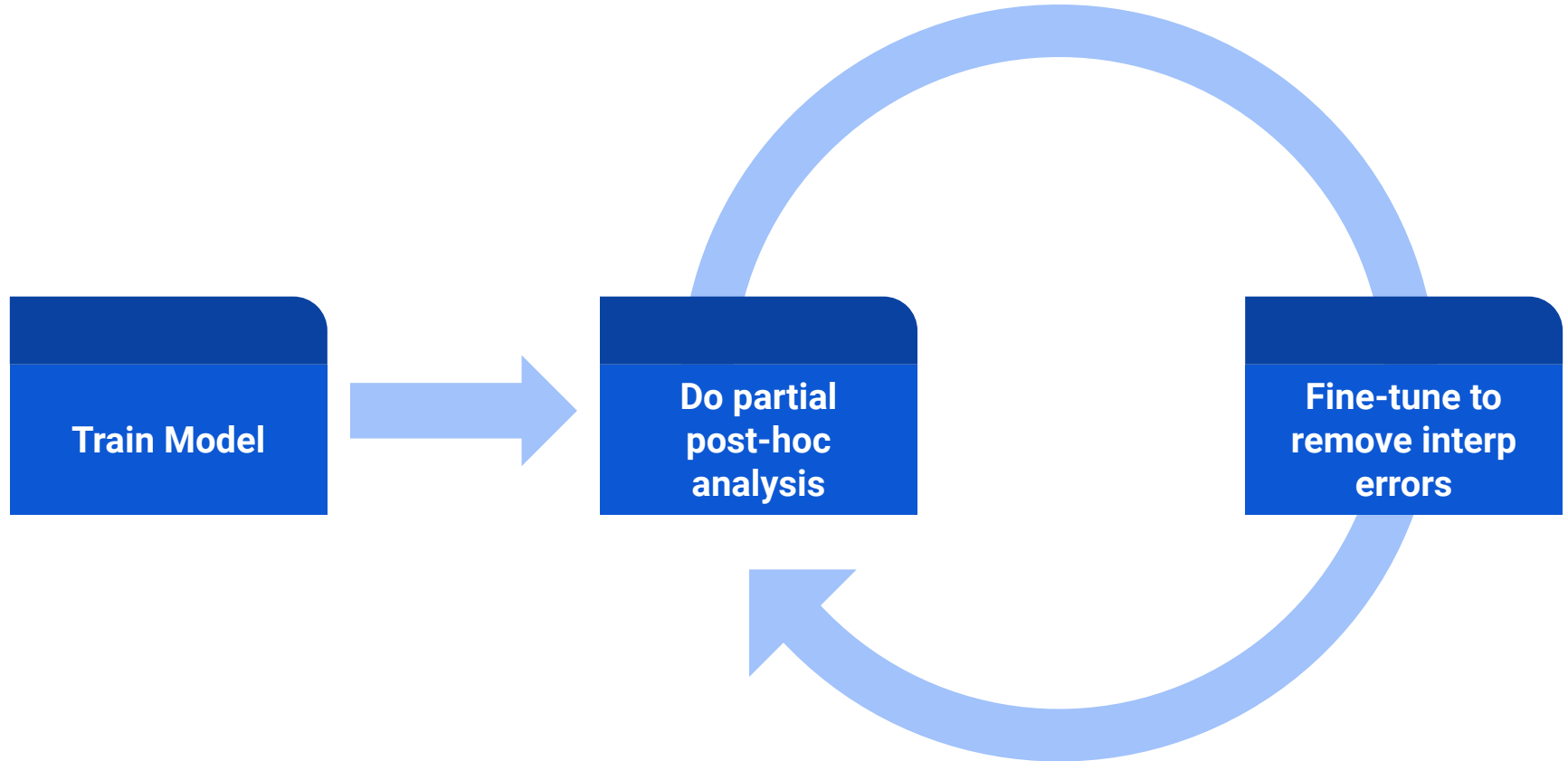
Puzzle: Why does more structure not always mean better bound?

Compounding errors from lack of structure



| Approximation Strategy | Result | Complexity |
|----------------------------------|---------------|--|
| (exact) max row diff | ≈ 1.8 | $(\mathcal{O}(d_{\text{vocab}}^2 d_{\text{model}}))$ |
| $2 \cdot (\text{max abs value})$ | ≈ 2.0 | $(\mathcal{O}(d_{\text{vocab}}^2 d_{\text{model}}))$ |
| max row diff on subproduct | ≈ 5.7 | $(\mathcal{O}(d_{\text{vocab}} d_{\text{model}}^2))$ |
| recursive max row diff | ≈ 97 | $(\mathcal{O}(d_{\text{vocab}} d_{\text{model}}))$ |

Wanted: Compression of highly expressive systems



Check out our poster!

Scan for paper



Scan for Poster

