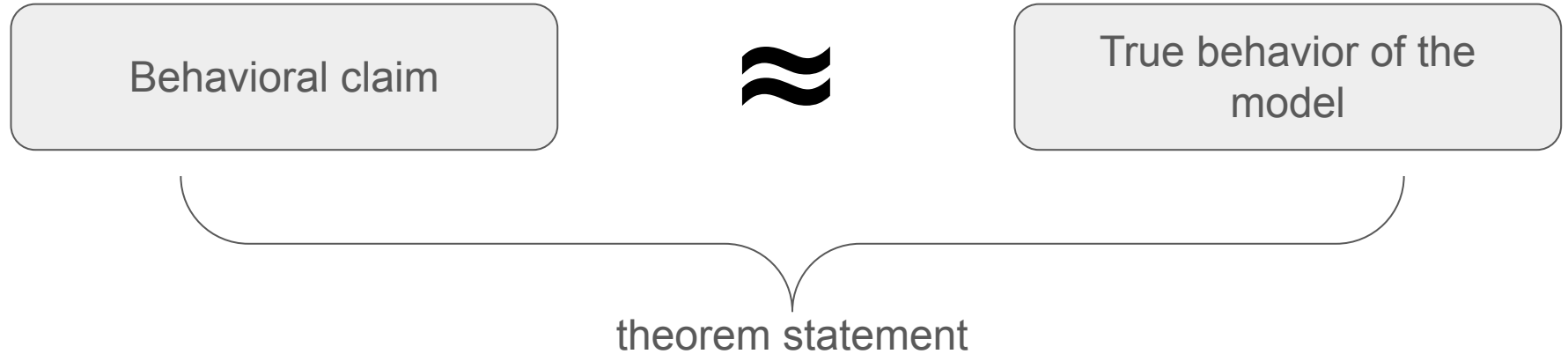# Compact Proofs of Model Performance via Mechanistic Interpretability

## TLDR

Proof length can be a metric on mech interp

# Formalizing proof length to quantify compression

Behavioral claim ≈ True behavior of the model

theorem statement

Proof = sound computation of worst-case error (divergence in behavior)

Length of proof = cost of running computation

# Formalizing proof length to quantify compression

Behavioral claim ≈ True behavior of the model

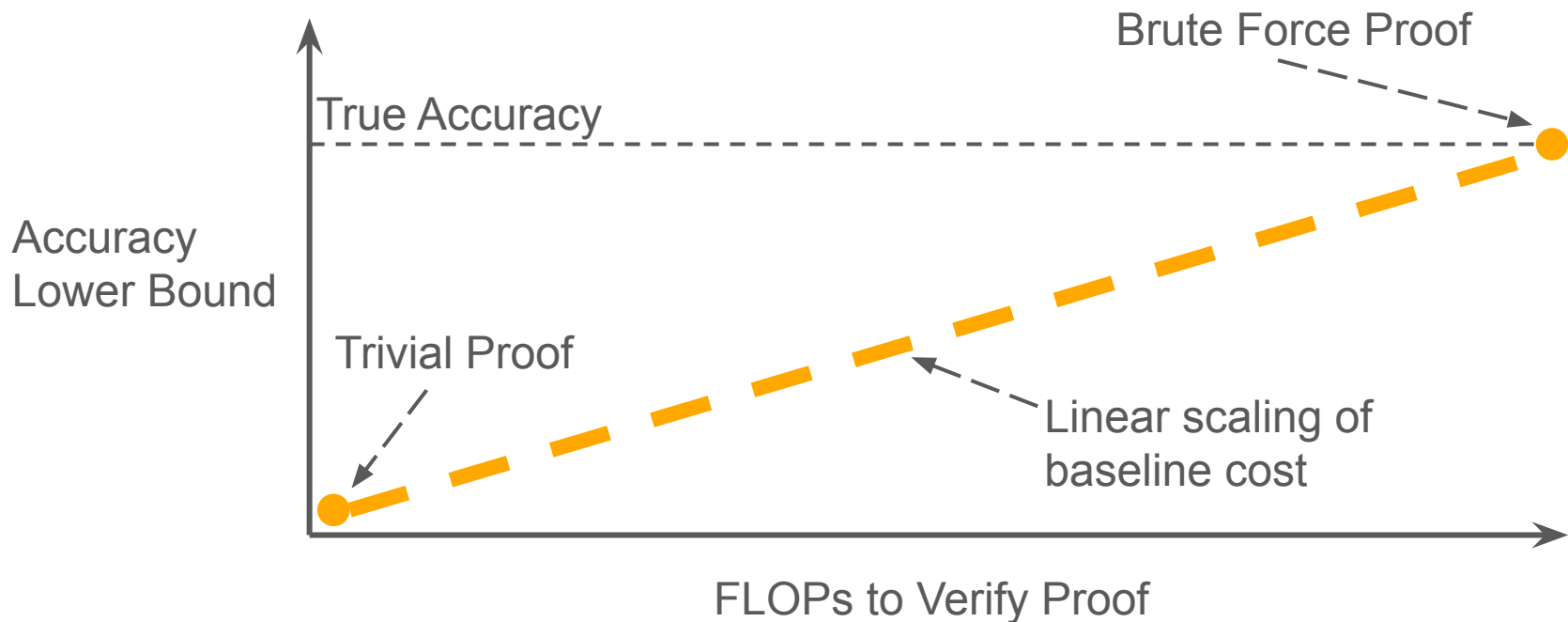$$\mathbb{E}_{(x,y)\sim\mathcal{D}}\left[f(y, M(x))\right] \geq b$$

theorem statement

Proof = sound computation of worst-case error (divergence in behavior)
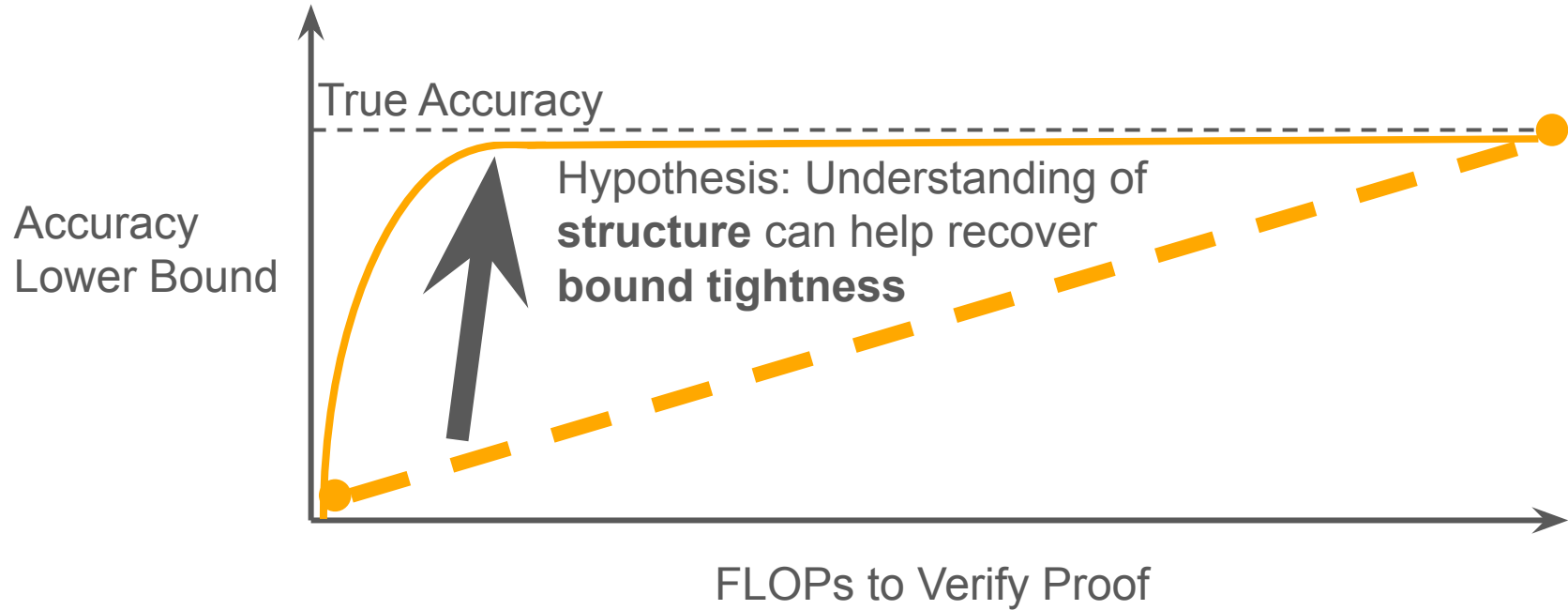
Length of proof = cost of running computation

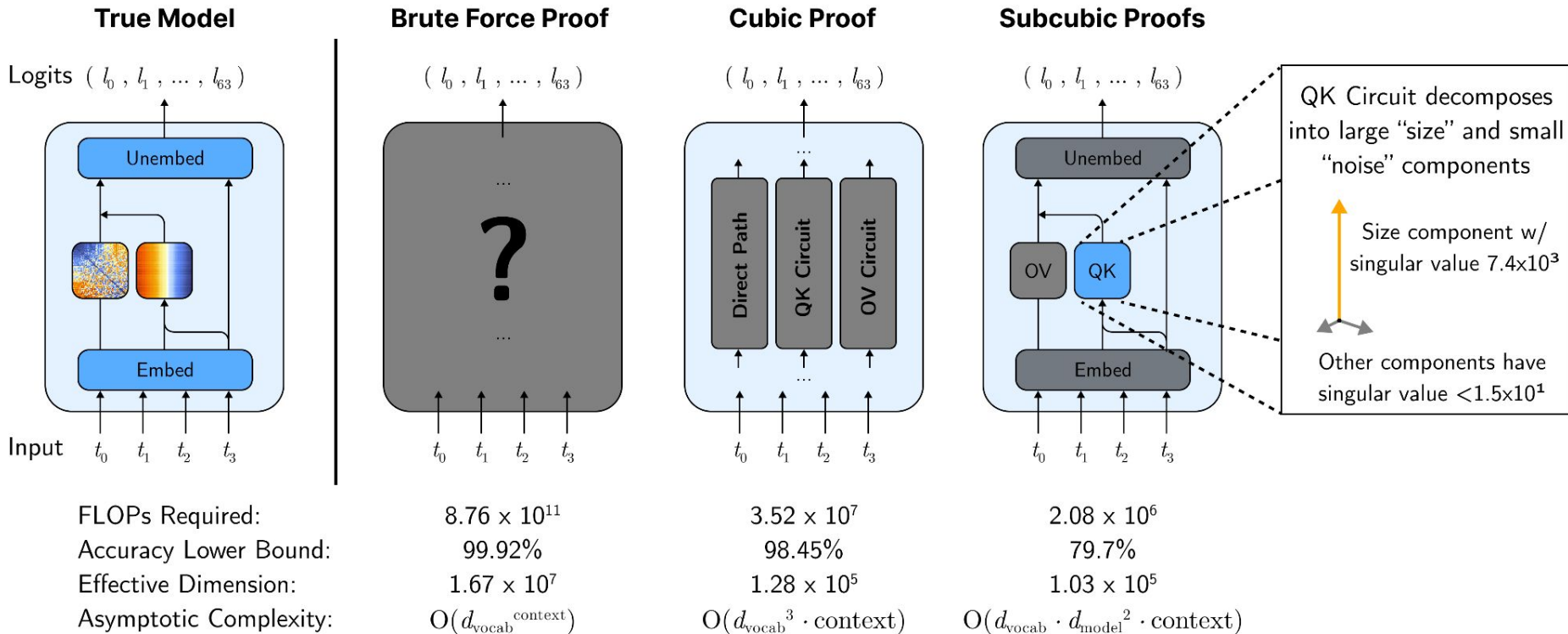# Quantifying the compute-cost of explanations

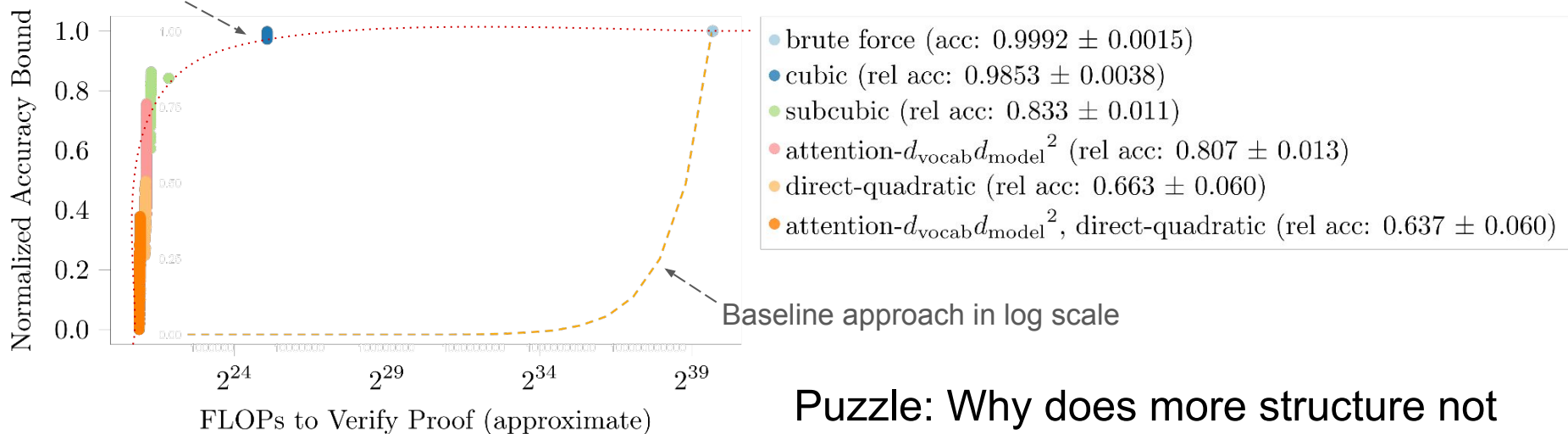# Does understanding improve upon the linear baseline?

True Accuracy

Accuracy
Lower Bound

Hypothesis: Understanding of **structure** can help recover **bound tightness**

FLOPs to Verify Proof

# Proofs with varying mechanistic understanding



**True Model**

Logits $(l_0, l_1, \ldots, l_{63})$

Unembed

Embed

Input $t_0$ $t_1$ $t_2$ $t_3$

**Brute Force Proof**

$(l_0, l_1, \ldots, l_{63})$

...

**?**

...

$t_0$ $t_1$ $t_2$ $t_3$

**Cubic Proof**

$(l_0, l_1, \ldots, l_{63})$

...

Direct Path | QK Circuit | OV Circuit

...

$t_0$ $t_1$ $t_2$ $t_3$

**Subcubic Proofs**

$(l_0, l_1, \ldots, l_{63})$

Unembed

OV | QK

Embed

$t_0$ $t_1$ $t_2$ $t_3$

QK Circuit decomposes into large "size" and small "noise" components

Size component w/ singular value $7.4 \times 10^3$

Other components have singular value $< 1.5 \times 10^1$

| | Brute Force Proof | Cubic Proof | Subcubic Proofs |
|---|---|---|---|
| FLOPs Required: | $8.76 \times 10^{11}$ | $3.52 \times 10^7$ | $2.08 \times 10^6$ |
| Accuracy Lower Bound: | 99.92% | 98.45% | 79.7% |
| Effective Dimension: | $1.67 \times 10^7$ | $1.28 \times 10^5$ | $1.03 \times 10^5$ |
| Asymptotic Complexity: | $\mathrm{O}(d_{\text{vocab}}{}^{\text{context}})$ | $\mathrm{O}(d_{\text{vocab}}{}^3 \cdot \text{context})$ | $\mathrm{O}(d_{\text{vocab}} \cdot d_{\text{model}}{}^2 \cdot \text{context})$ |

# We found an empirical "pareto frontier"

Pareto frontier from incorporating mechanistic understanding



Baseline approach in log scale

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

# Compounding errors from lack of structure



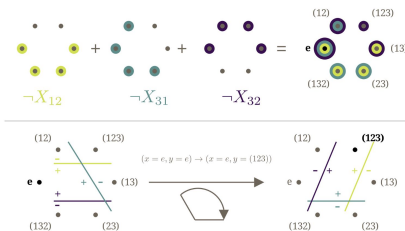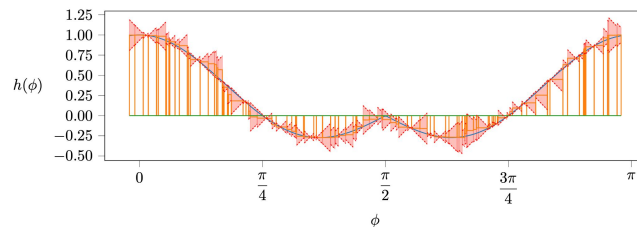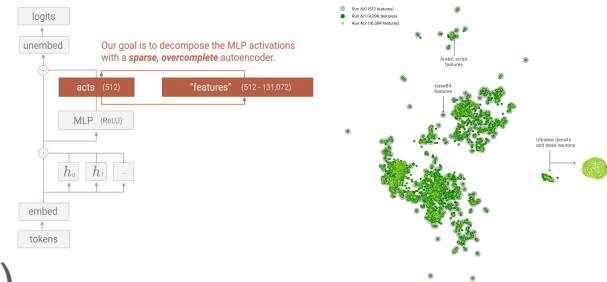| Approximation Strategy | Result | Complexity |
|---|---|---|
| (exact) max row diff | ≈ 1.8 | $(\mathcal{O}(d_{\text{vocab}}{}^2 d_{\text{model}}))$ |
| 2 · (max abs value) | ≈ 2.0 | $(\mathcal{O}(d_{\text{vocab}}{}^2 d_{\text{model}}))$ |
| max row diff on subproduct | ≈ 5.7 | $(\mathcal{O}(d_{\text{vocab}}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

# Applying Compact Proofs



- Optimization targets for representation search (SAEs)



- Compressing MLPs (integration)

- Ground truth for comparing mech interp approaches (groups)

Some images from https://transformer-circuits.pub/2023/monosemantic-features

# Open Problems for Scaling Compact Proofs



- Fix compounding errors
    - Fine-tuning; or heuristic arguments; or sampling
- Suppress exponential in # layers
    - Toy model: induction heads
- Autoformalize proofs
    - AlphaProof
- Autointerp
    - Step 2: ???
- Step 3: Profit

Images by GPT-4o

https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/