We use **mechanistic understanding** to compress **proofs of model performance** on toy transformers

Compact Proofs of Model Performance via Mechanistic Interpretability

We **formalize** post-hoc mechanistic interpretability as proving worst-case generalization bounds on the performance of models.

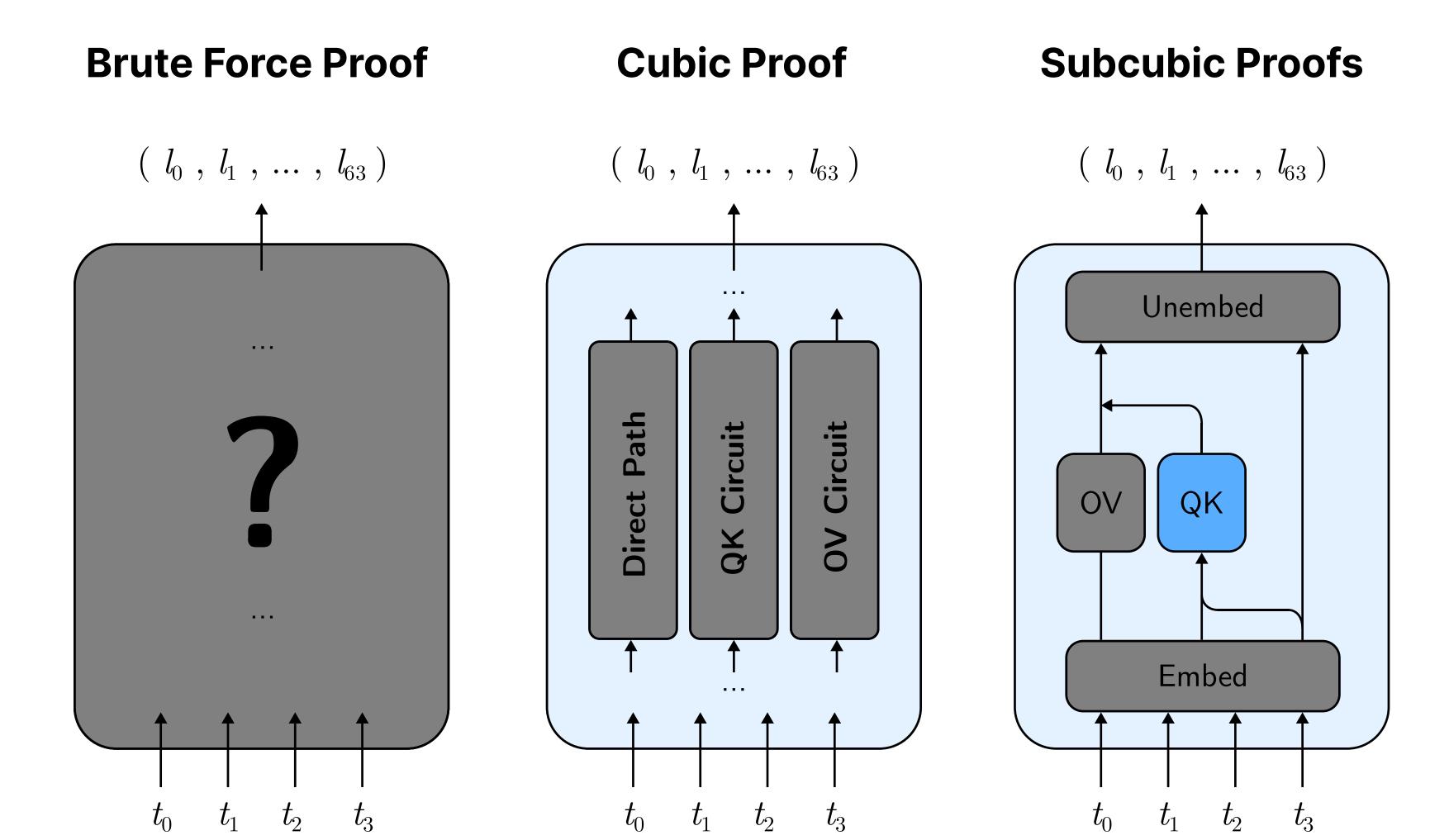
 $\mathbb{E}_{\mathbf{x}\sim\mathcal{D}}[f(\mathcal{M}(\mathbf{x}))]\geq b$

We define a bound approximation algorithm; a proof is the trace of running the algorithm along with an explanation that it provides valid lower bounds on every input.

Algorithm 1 Counting Correct Sequences By Brute Force

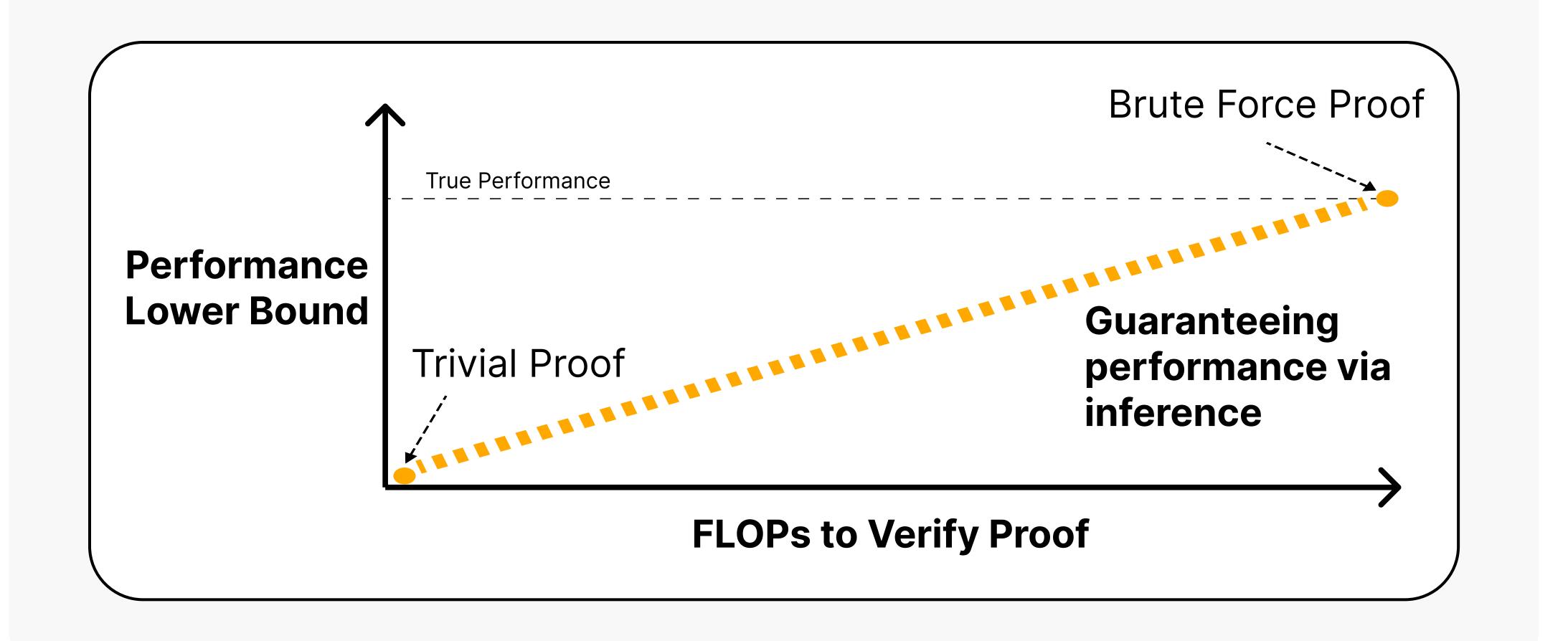
- 1: function CORRECTNESS(\mathcal{M} , input-sequence)
- 2: return MODEL-BEHAVIOR(\mathcal{M} , input-sequence) == MAX(input-sequence) 3: end function
- 4: function BRUTE-FORCE $(d_{\text{vocab}}, n_{\text{ctx}}, \mathcal{M})$
- 5: return $\frac{1}{d_{\text{vocab}}n_{\text{ctx}}}$ SUM(CORRECTNESS(\mathcal{M} , tokens) for tokens $\in (\text{RANGE}(d_{\text{vocab}}))^{n_{\text{ctx}}})$ 6: end function

We implement many proof strategies, which incorporate varying degrees of **mechanistic detail**...



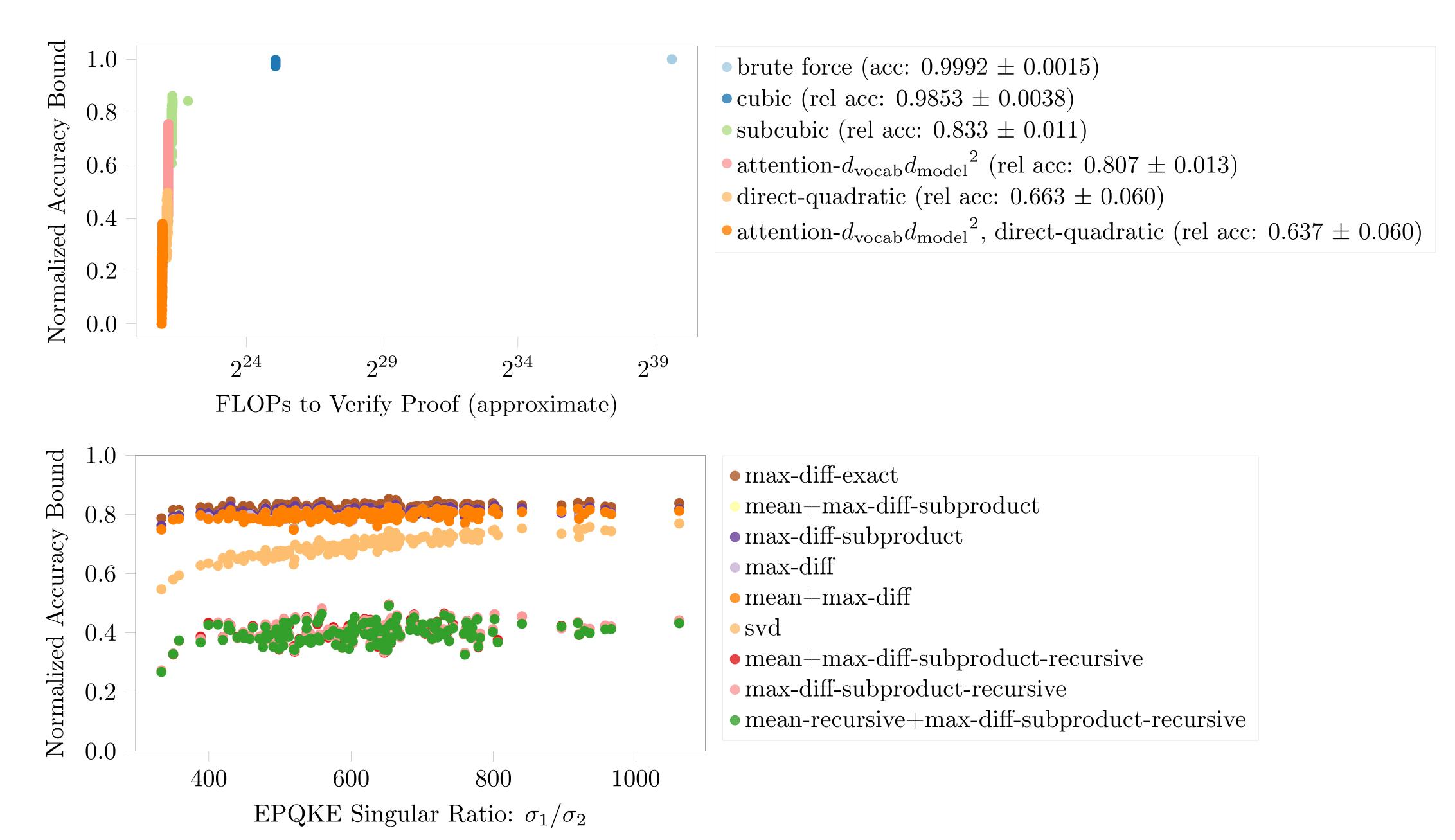
... and find that more mechanistic detail results in **tighter performance bounds**.

We can run inference on various inputs to guarantee model performance on those inputs. Can mechanistic understanding provide a **compression** of model behavior that beats this inefficient baseline?



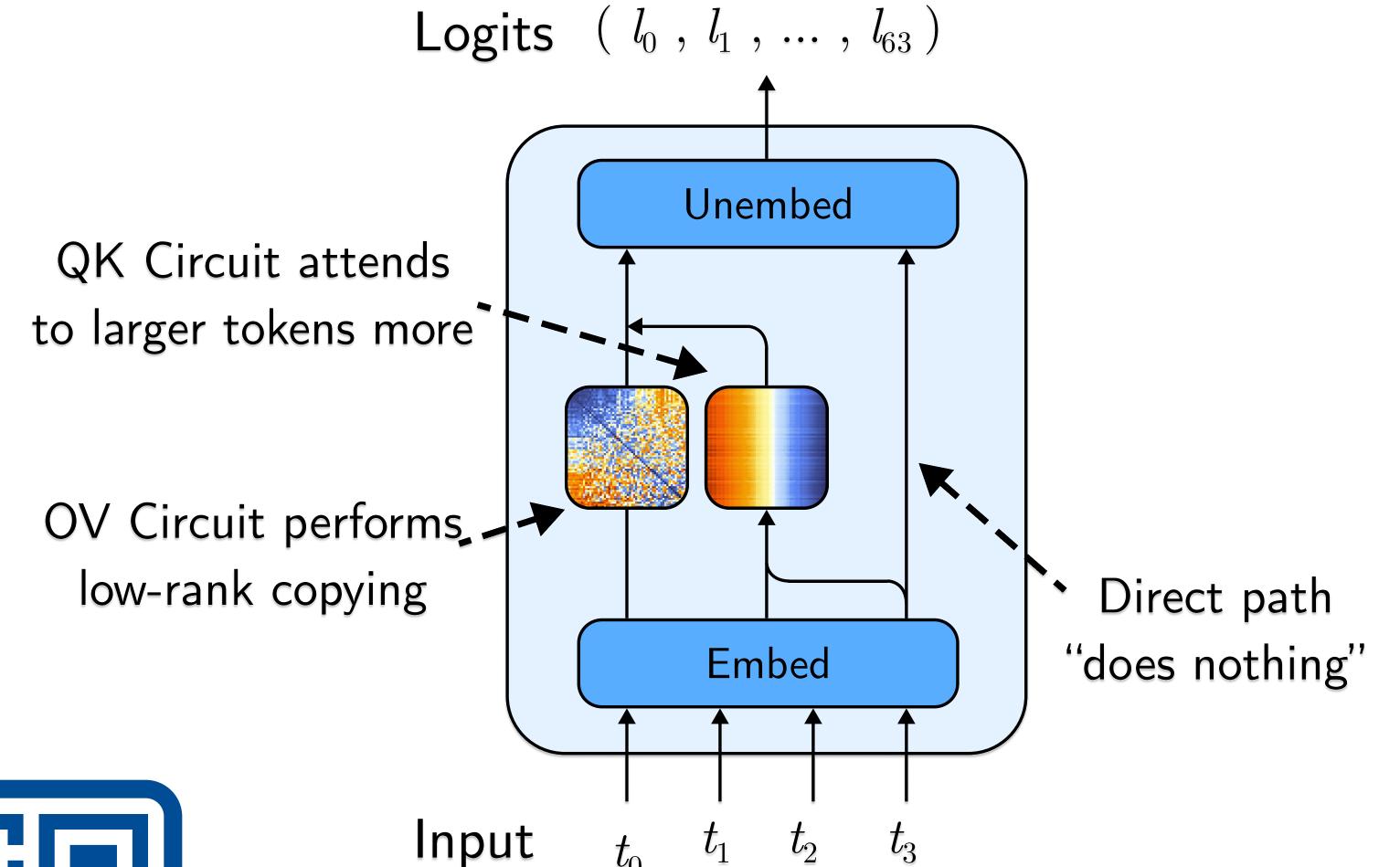
We prototype the proofs-based approach to

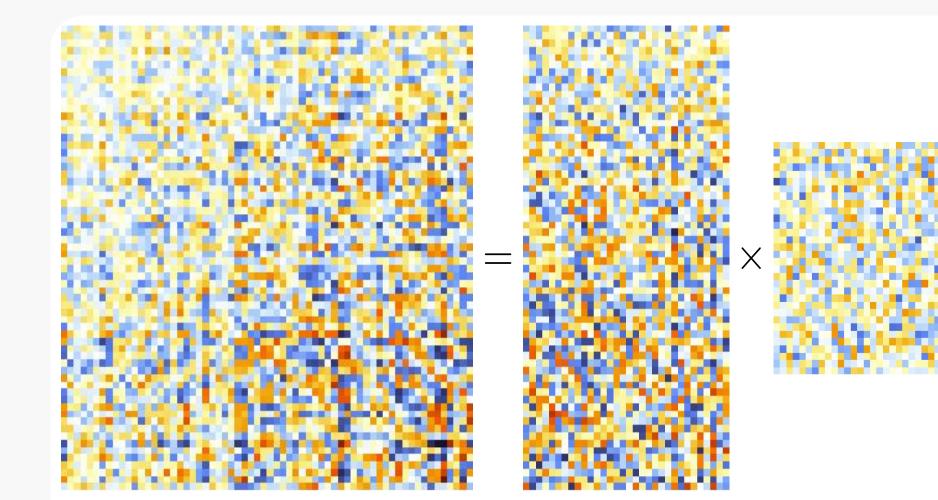
Moreover, for any given proof length, proving a tighter performance bound requires more mechanistic detail.

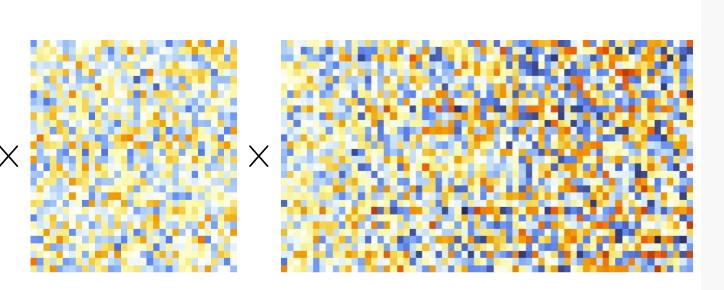


We identify **compounding approximation error** in post-hoc analysis as a key challenge to proving worst-case bounds.

formal interpretability on transformer models trained on the the max_of_k task.







Approximation Strategy Result
(exact) max row diff ≈ 1.8
$2 \cdot (\text{max abs value}) \approx 2.0$
max row diff on subproduct ≈ 5.7
recursive max row diff ≈ 97

Complexity
$(\mathcal{O}({d_{\mathrm{vocab}}}^2 d_{\mathrm{model}}))$
$(\mathcal{O}({d_{\mathrm{vocab}}}^2 d_{\mathrm{model}}))$
$(\mathcal{O}(d_{\mathrm{vocab}} d_{\mathrm{model}}^2))$
$(\mathcal{O}(d_{\mathrm{vocab}}d_{\mathrm{model}}))$



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