

# How Do Transformers “Do” Math?

Subhash Kantamneni, Matt McManus, Alex Quach, Daniel (Dong Young) Kim  
MIT 6.S986, LLM and Beyond, Spring 2024



## Background & Motivation

How do transformers compute mathematical quantities? We study this question through “intermediates”

$$Y = wX$$



$w$  is an intermediate ( $I$ ) because it is not directly inputted/outputted by the model

but what if the model was using  $\exp(\log(w) + \log(x))$  or  $\sqrt{w^2 x^2}$  ?

Key Questions:

- How can we find if a quantity is represented in a transformer?
- How can we prove that a model is using method  $g$  with an intermediate  $I$  (e.g.,  $g = wx$ ,  $I = w$ )
- How can we apply this to non-trivial problems?

## Experimental Setup

**Model Problem:**

$$Y = wX$$

- Sample 5000 values of  $w \in [-0.75, 0.75]$ 
  - Sample 65 (x, y) points for given  $w$
- Train transformer with  $L = \{1, 2, 3, 4, 5\}$  and  $H = \{2, 4, 8, 16, 32\}$

## Interpretability Techniques:

**Linear Probe**

Linear( $f(I)$ ,  $HS$ ) finds  $W$  s.t.  
 $f(I) = W \times HS$

**Nonlinear Probe (Taylor Probe)**

$$f(I) = a_1 I + a_2 I^2 + \dots + a_n I^n$$

**Reverse Probe**

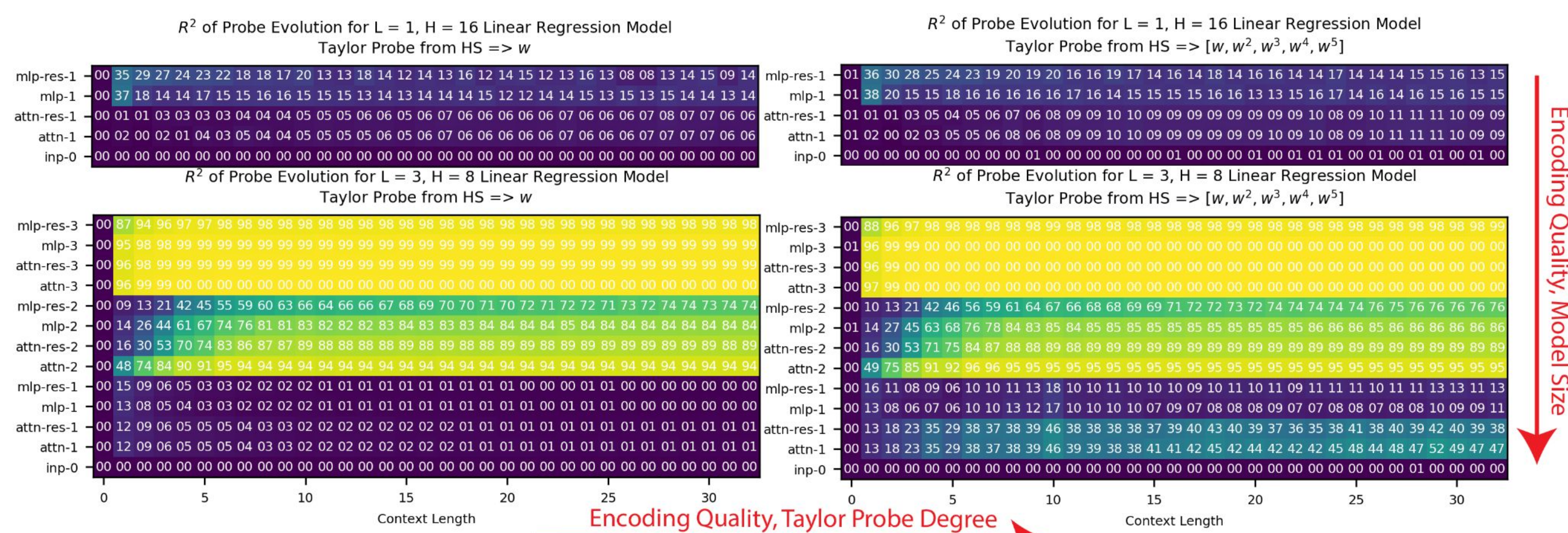
Determine proportion of hidden state represented by  $I$

**Intervening**

Change model output from using  $w \Rightarrow w'$

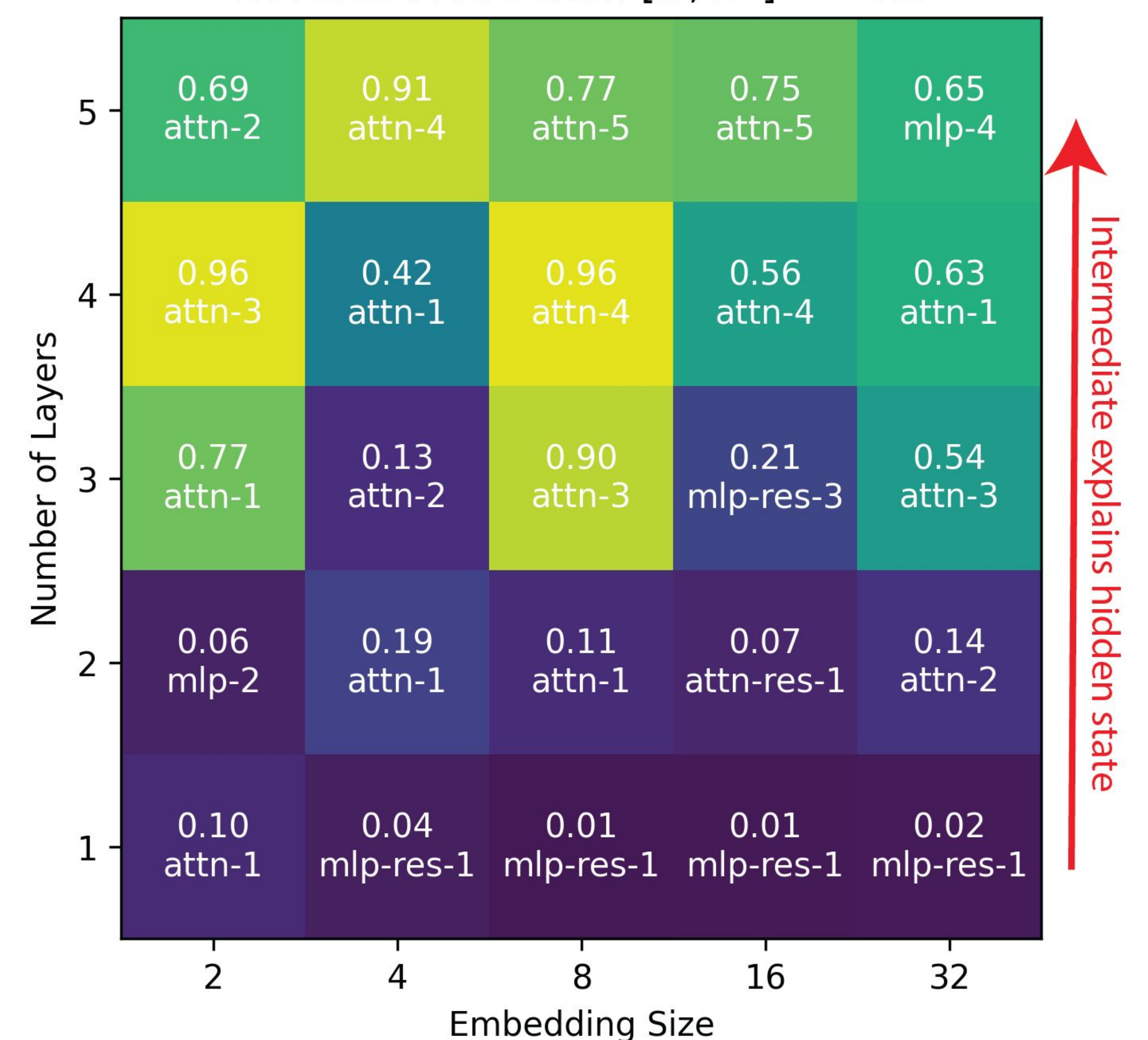
## Results & Discussion

I. If a model uses a method  $g$ , its hidden state should encode  $I$



III. Iff the model uses  $g$ , we expect some hidden state's variance to be almost fully explained by  $I$

$\max(R^2)$  for Linear Regression Models  
Reverse Probe from  $[w, w^2] = HS$



IV. Iff the model uses  $g$ , we can intervene on hidden states to change  $I \Rightarrow I'$  and predictably change the model output from  $g(X, I) \Rightarrow g(X, I')$

II. If a model uses a method  $g$ , model performance should improve if  $I$  is better represented

