# Pretrained transformer efficiently learns low-dimensional target functions in-context

Kazusato Oko, Yujin Song, Taiji Suzuki, Denny Wu

UC Berkeley   UTokyo   AIP   NEW YORK UNIVERSITY   FLATIRON INSTITUTE

# In-Context Learning (ICL)

- Pretrained transformers can recognize patterns from prompts *without updating model parameters*

- A very short context can be sufficient

Please guess the number that fits in the '?'.

context
```
1,1 -> 2
2,3 -> 5
8,13 -> 21
6,0 -> 6
10,1 -> 11
```
query  5,27 -> ?

The pattern in the given pairs of numbers appears to be the sum of the two numbers.

So, the number that fits in the '?' is 32.

Question

ChatGPT (GPT-4)

# Prior Works

- Known fact: linear transformers can emulate *linear regression on the context* in its forward pass [ACDS23, MHM23, ZFB23...]

  - requires the same context length ($N$) as the amount of data needed for linear regression

  - higher vector dimension (of $x$)
    → higher required context length

**Q: Can TF outperform learning algorithms applied directly to the prompt?**

(in terms of context length)

prompt $t$

$$x_1 \rightarrow y_1 = x_1^\top \beta^t$$
$$x_2 \rightarrow y_2 = x_2^\top \beta^t$$
$$\vdots$$
$$x_N \rightarrow y_N = x_N^\top \beta^t$$

context

$$x^q \rightarrow y^q = ?$$

query

# Our Main Message

Q: Can TF outperform learning algorithms working directly on prompt?

A: Yes, by *adapting to the problem structure* during pretraining

# Problem Setting

- Learning single-index functions in-context

- On $t$-th prompt, $\boldsymbol{x} \sim N(0, \boldsymbol{I}_d) \in \mathbb{R}^d$ and

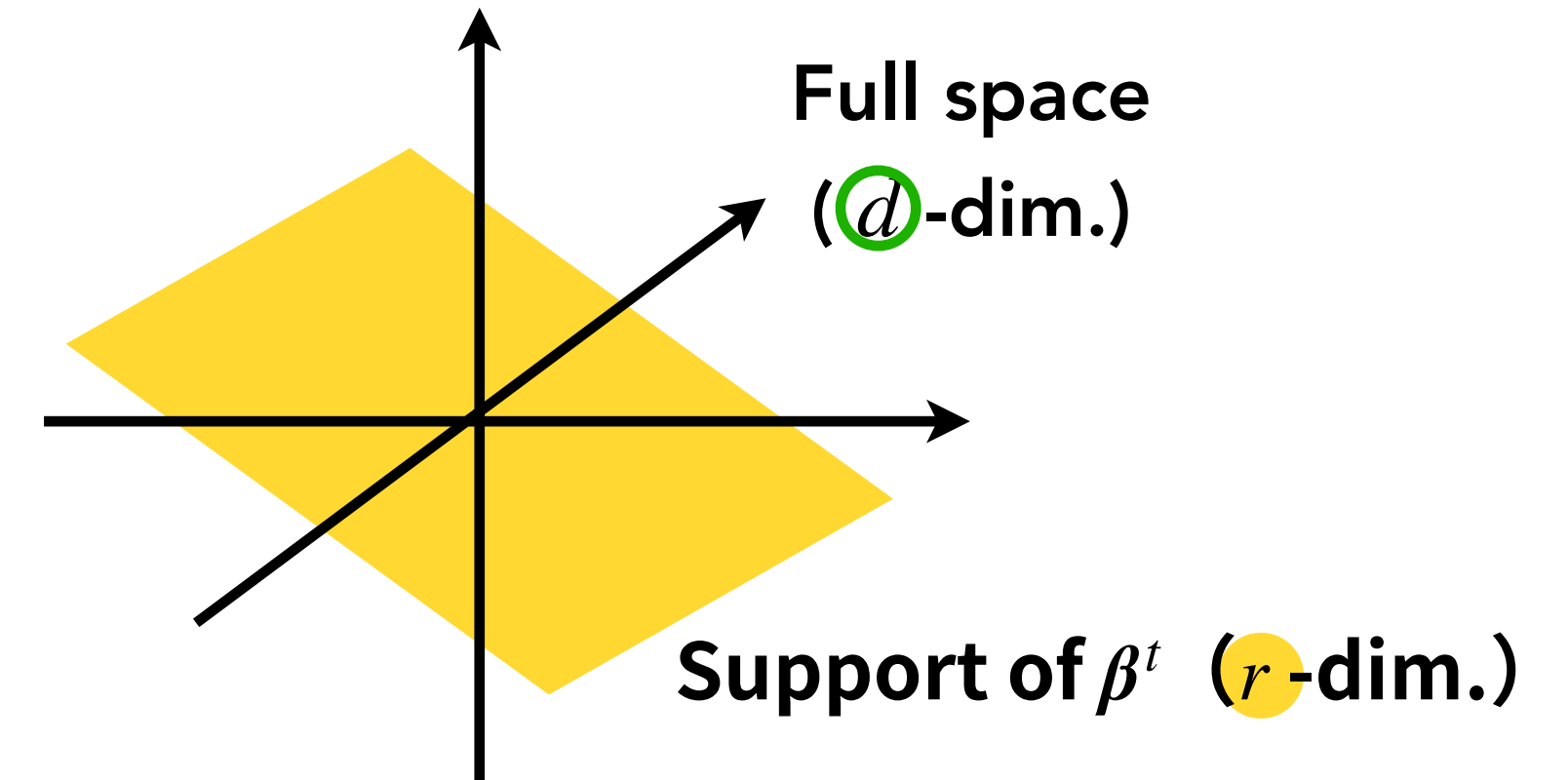$$y = \sigma_*^t(\boldsymbol{x}^\top \boldsymbol{\beta}^t)$$

> $y$ depends only on the direction of $\boldsymbol{\beta}^t$

- $\sigma_*^t$ : random polynomial of degree $P$ (nonlinear)

- $\boldsymbol{\beta}^t \in \mathbb{R}^d$ : random vector drawn from
  $r \ll d$-dimensional subspace of $\mathbb{R}^d$

**Problem distribution is low-dimensional**

- Learning algorithms (kernel, NN…) on the test prompt need $\mathrm{poly}(d)$ samples

…Can pretrained TF outperform them?

[GMMM21, BAGJ21]

**prompt $t$**

context

$$\boldsymbol{x}_1 \to y_1 = \sigma_*^t(\boldsymbol{x}_1^\top \boldsymbol{\beta}^t)$$

$$\boldsymbol{x}_2 \to y_2 = \sigma_*^t(\boldsymbol{x}_2^\top \boldsymbol{\beta}^t)$$

$$\vdots$$

$$\boldsymbol{x}_N \to y_N = \sigma_*^t(\boldsymbol{x}_N^\top \boldsymbol{\beta}^t)$$

query

$$\boldsymbol{x}^q \to y^q = ?$$

**Full space**
($d$-dim.)

**Support of $\beta^t$** ($r$-dim.)

# Our Main Result

- Consider pretraining a single-layer transformer (nonlinear MLP+attention) on $d^{\Theta(Q)}$ tasks with a prompt length of $d^{\Theta(Q)}$ ($Q$: *lowest* degree of $\boldsymbol{\sigma}_*$ in $y = \sigma_*(\boldsymbol{x}^\top \boldsymbol{\beta})$).

Polynomial

Context $\boldsymbol{X}, \boldsymbol{y}$
Query $\boldsymbol{x}^q$

$\boldsymbol{X} = (\boldsymbol{x}_1 \cdots \boldsymbol{x}_N)$

$y = (y_1 \cdots y_N)^\top$

MLP
(width $\gtrsim r^{\Theta(P)}$)

Linear attention

$\to \hat{y}^q(\boldsymbol{X}, \boldsymbol{y}, \boldsymbol{x}^q)$

Estimation of $y^q$

Pretraining

I. One-step gradient descent on MLP weight

II. Ridge regression on attention matrix

- **Theorem** TF pretrained above achieves low test error ($\mathbb{E}[|\hat{y}^q - y^q|] = o_d(1)$) if context length $N^*$ at test prompt satisfies $N^* \gtrsim r^{4P}$ ($P$: *highest* degree of $\boldsymbol{\sigma}_*$)
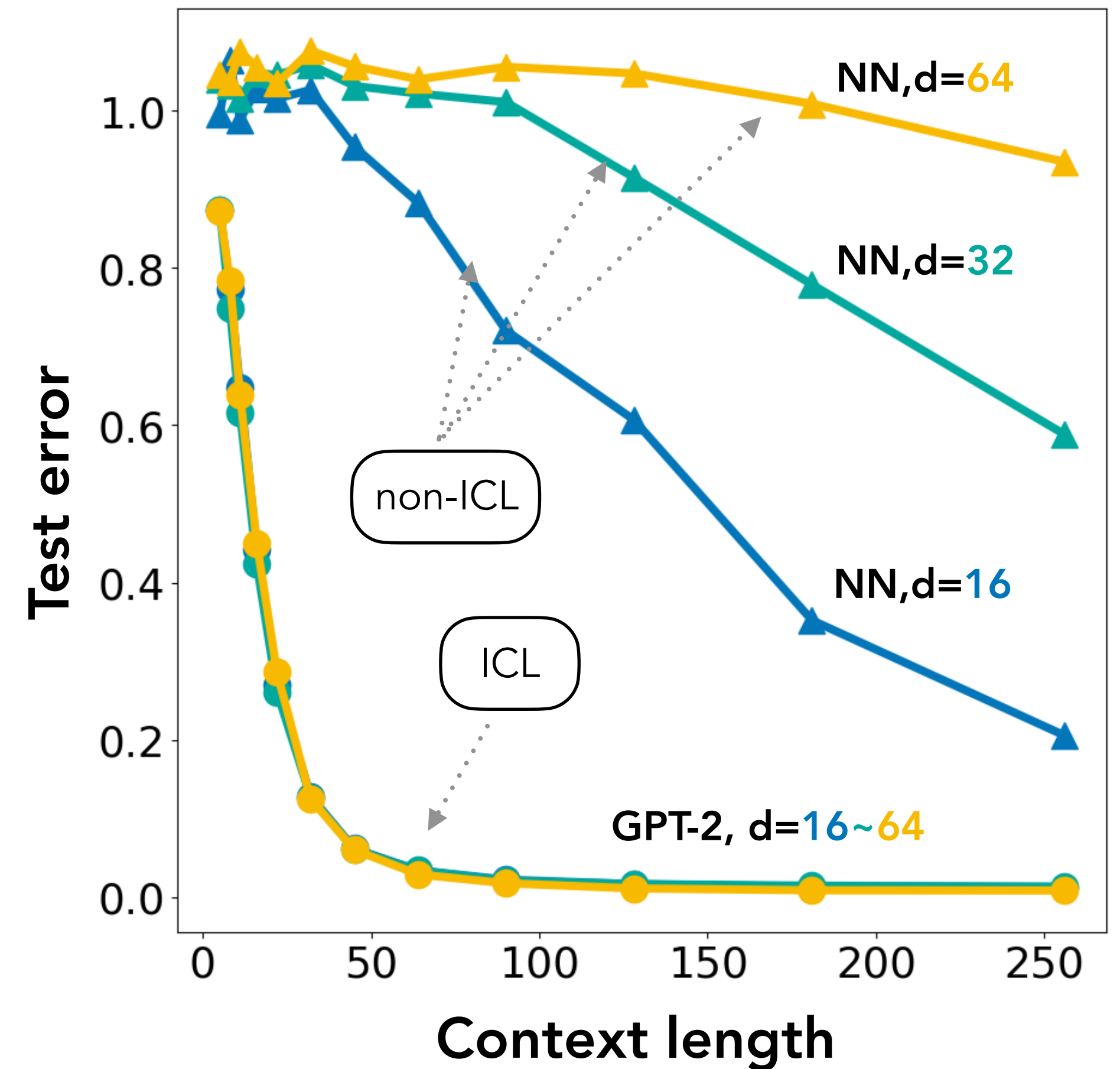
**Required prompt length only depends on the inner dimension $r$**

Baseline algorithms (kernel, NN) require $d$-dependent amount of data $\to$ superiority under $r \ll d$

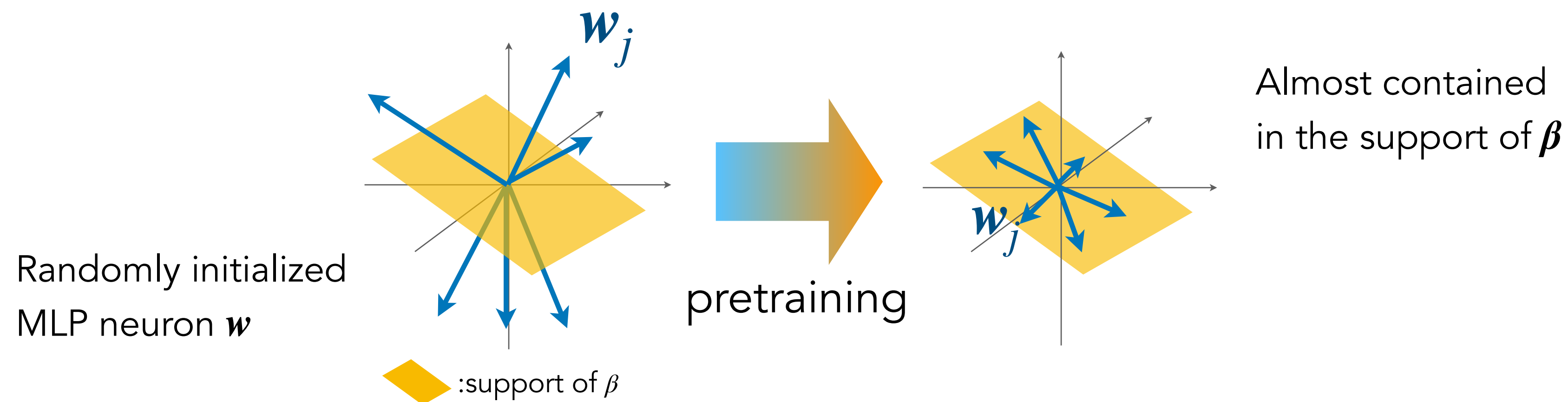*Pretraining is nonconvex optimization… end-to-end optimization & generalization analysis

# Experiment

- We fix the inner dimension $r = 8$, while altering the ambient dimensiom $d$ from 16 to 64, for the problem $y = \sigma_*^t(\boldsymbol{x}^\top \boldsymbol{\beta}^t)$.

- NN performance deteriorates with increasing $d$

- GPT-2 achieves low test error even when $d$ is high

# Takeaway & Mechanism

- Takeaway: TF can adapt to the prior distribution of problems via pretraining

- Mechanism: pretrained MLP neurons align with the $r$-dimensional subspace



Randomly initialized
MLP neuron $\boldsymbol{w}$

$\boldsymbol{w}_j$

pretraining

$\boldsymbol{w}_j$

Almost contained
in the support of $\boldsymbol{\beta}$

◆ :support of $\beta$

- This "memorization" of the prior distribution of problems results in $d$-free context length complexity

See you in Vancouver!

preprint: https://arxiv.org/abs/2411.02544