# Using Algorithms to Understand Transformers (and Using Transformers to Understand Algorithms)

Vatsal Sharan (USC)
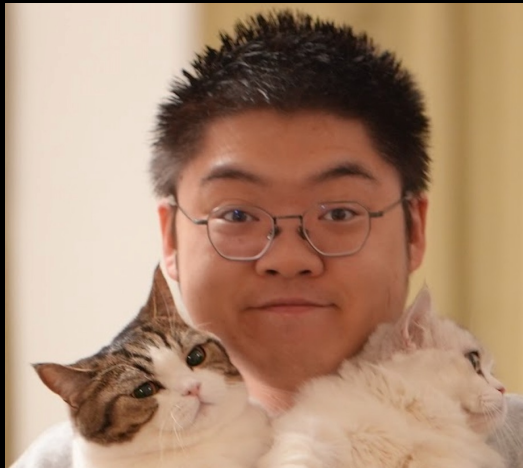
USC University of Southern California

- How can we use understanding of computational and information theoretic landscape to understand Transformers?
- How can we use Transformers to understand and discover algorithms and data structures?
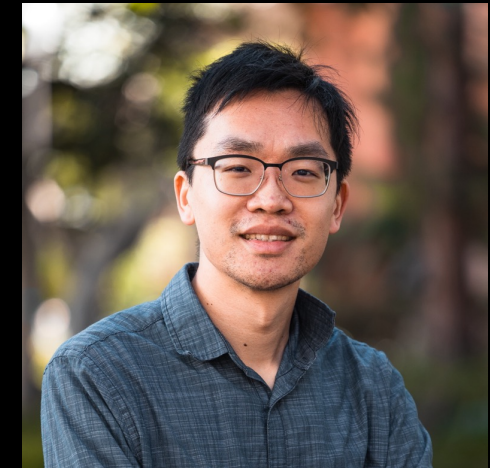
# How do Transformers do linear regression?



Deqing Fu (USC)

Tianqi Chen (USC)

Robin Jia (USC)

Transformers Learn Higher-Order Optimization Methods for In-Context Learning: A Study with Linear Models, Neurips 2024

# Transformers excel at in-context learning



sea otter => loutre de mer          *examples*

peppermint => menthe poivrée

plush girafe => girafe peluche

cheese =>                           *prompt*

In-context learning

vs.

| | |
|---|---|
| 1  sea otter => loutre de mer | *example #1* |
| ↓ | |
| **gradient update** | |
| ↓ | |
| 1  peppermint => menthe poivrée | *example #2* |
| ↓ | |
| **gradient update** | |
| ↓ | |
| 1  plush giraffe => girafe peluche | *example #N* |
| **gradient update** | |
| 1  cheese => | *prompt* |

Usual fine-tuning

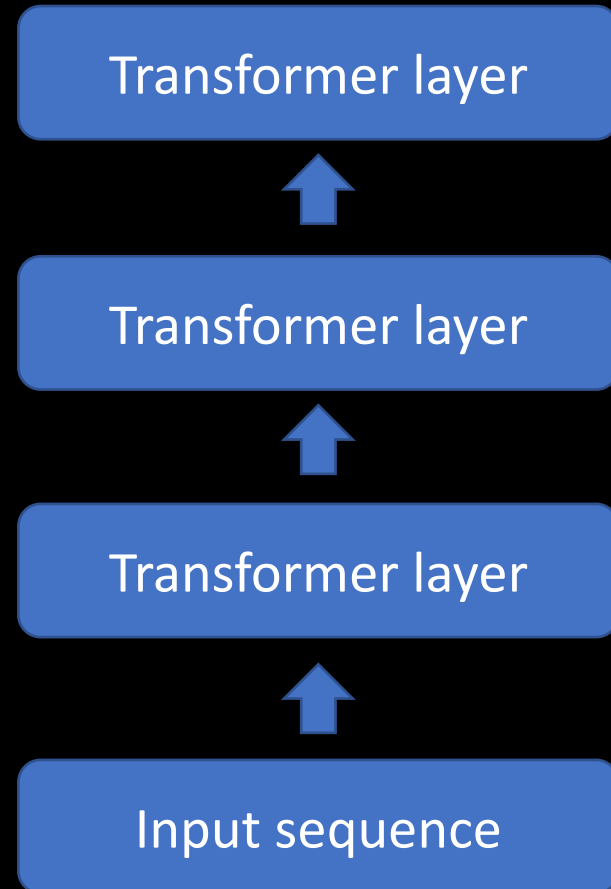# How do Transformers do in-context learning?

The case of linear models ($y_i = {w^*}^T x_i$):

$$x_1 = (3, 5), y_1 = 4$$

$$x_2 = (-2, 2), y_2 = 8$$

$$x_3 = (-7, -2), y_3 = 10$$

$$x_4 = (4, -1), y_4 = ?$$

References: Garg-Tsipras-Liang-Valiant 2022, Akyurek-Schuurmans-Andreas-Ma-Zhou 2022

# A prevailing hypothesis: Transformers do in-context learning via gradient descent
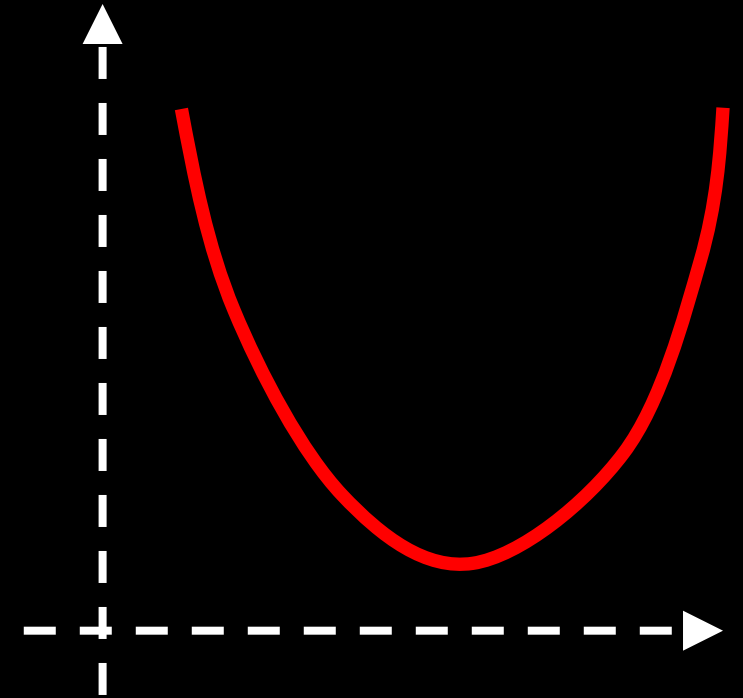
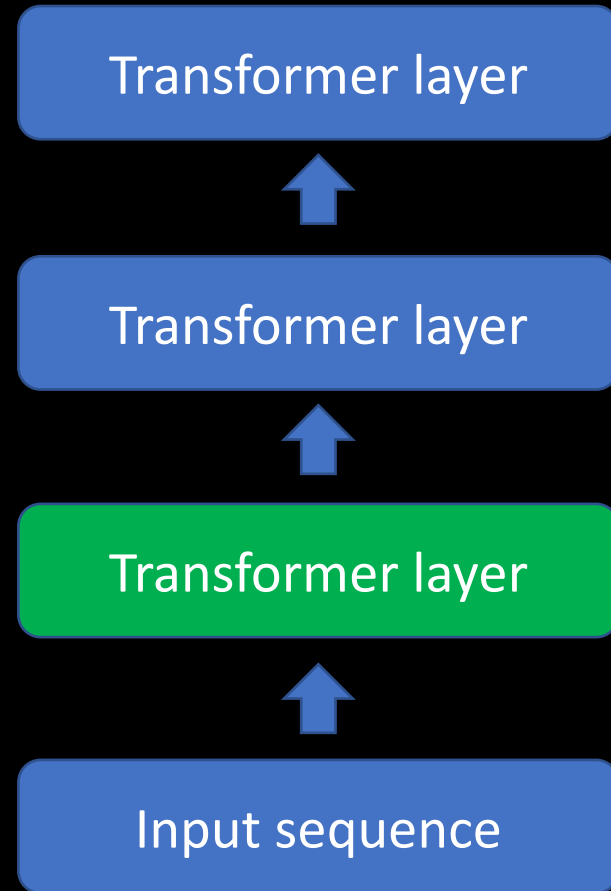Linear models:

$x_1 = (3, 5), y_1 = 4$

$x_2 = (-2, 2), y_2 = 8$
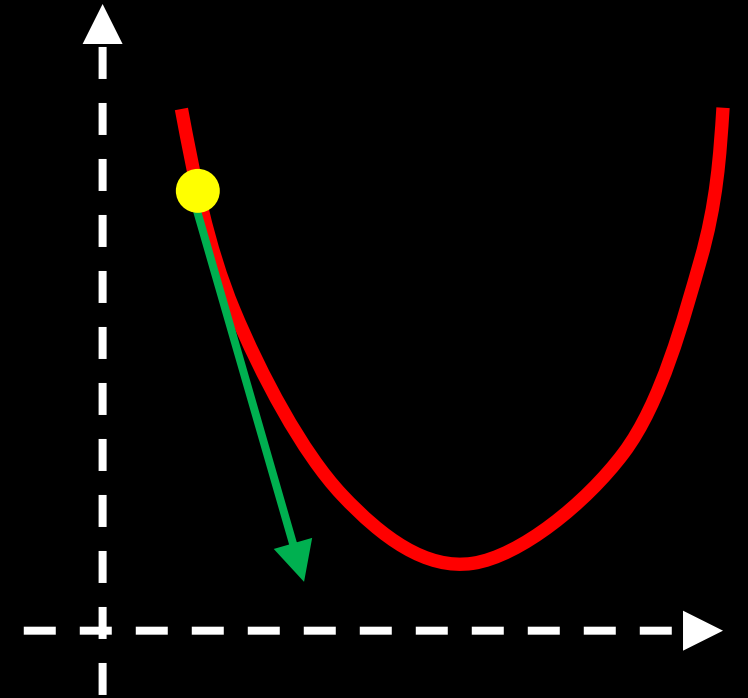
$x_3 = (-7, -2), y_3 = 10$

$x_4 = (4, -1), y_4 = ?$

Transformer layer

↑

Transformer layer

↑

Transformer layer

↑

Input sequence

≈

References: von Oswald et al. 2022, 2023, Ahn et al. 2023, Dai et al. 2023

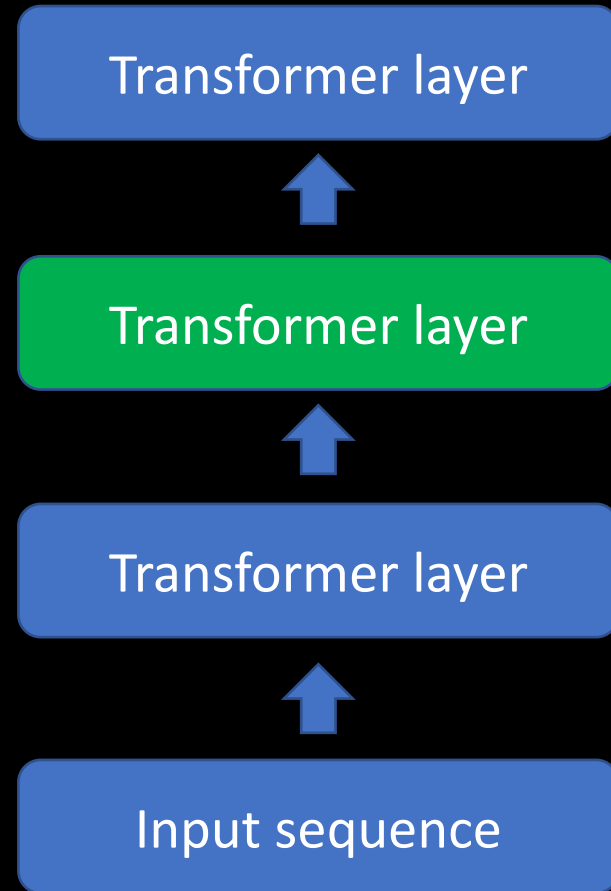# A prevailing hypothesis: Transformers do in-context learning via gradient descent
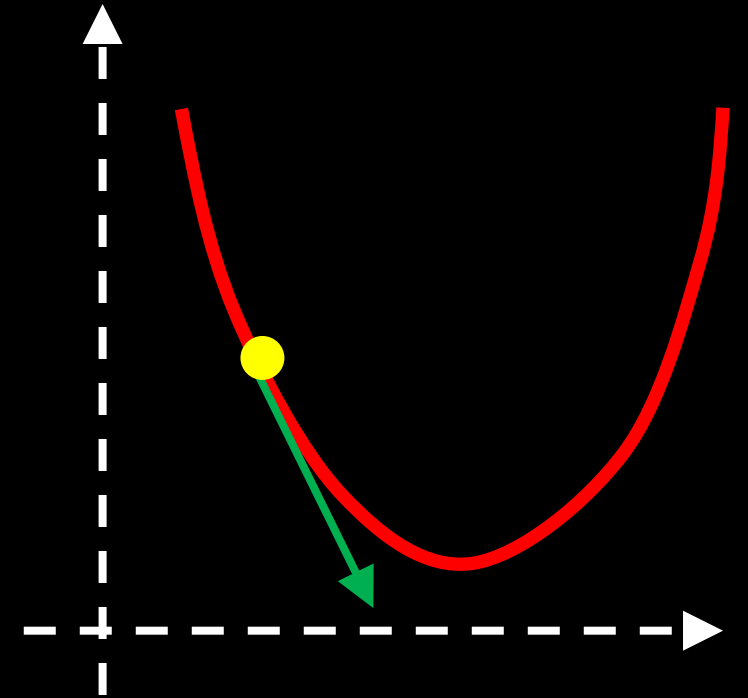
Linear models:

$x_1 = (3, 5), y_1 = 4$

$x_2 = (-2, 2), y_2 = 8$

$x_3 = (-7, -2), y_3 = 10$

$x_4 = (4, -1), y_4 = ?$

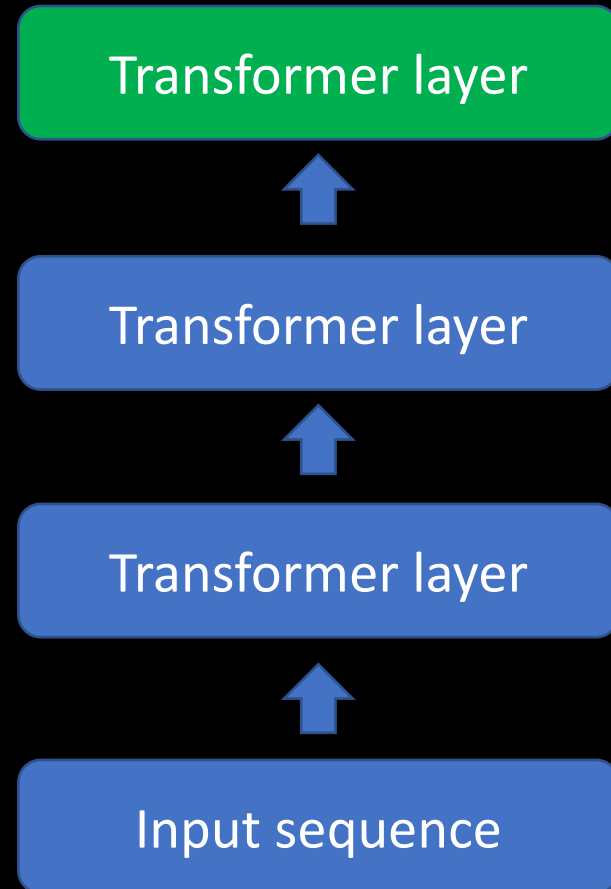# A prevailing hypothesis: Transformers do in-context learning via gradient descent
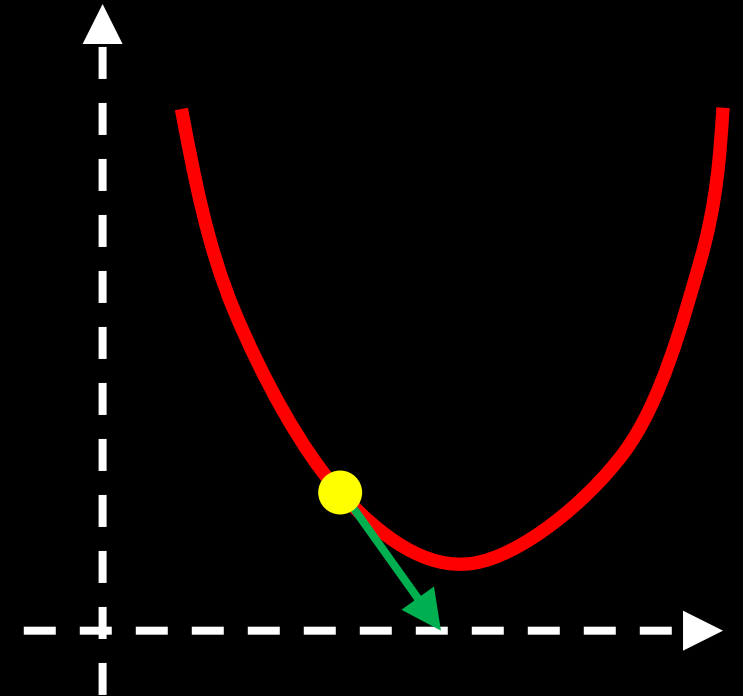
Linear models:

$x_1 = (3, 5), y_1 = 4$

$x_2 = (-2, 2), y_2 = 8$

$x_3 = (-7, -2), y_3 = 10$

$x_4 = (4, -1), y_4 = ?$

Transformer layer

Transformer layer

Transformer layer

Input sequence

$\approx$

# A prevailing hypothesis: Transformers do in-context learning via gradient descent

Linear models:

$x_1 = (3, 5), y_1 = 4$

$x_2 = (-2, 2), y_2 = 8$

$x_3 = (-7, -2), y_3 = 10$

$x_4 = (4, -1), y_4 = ?$

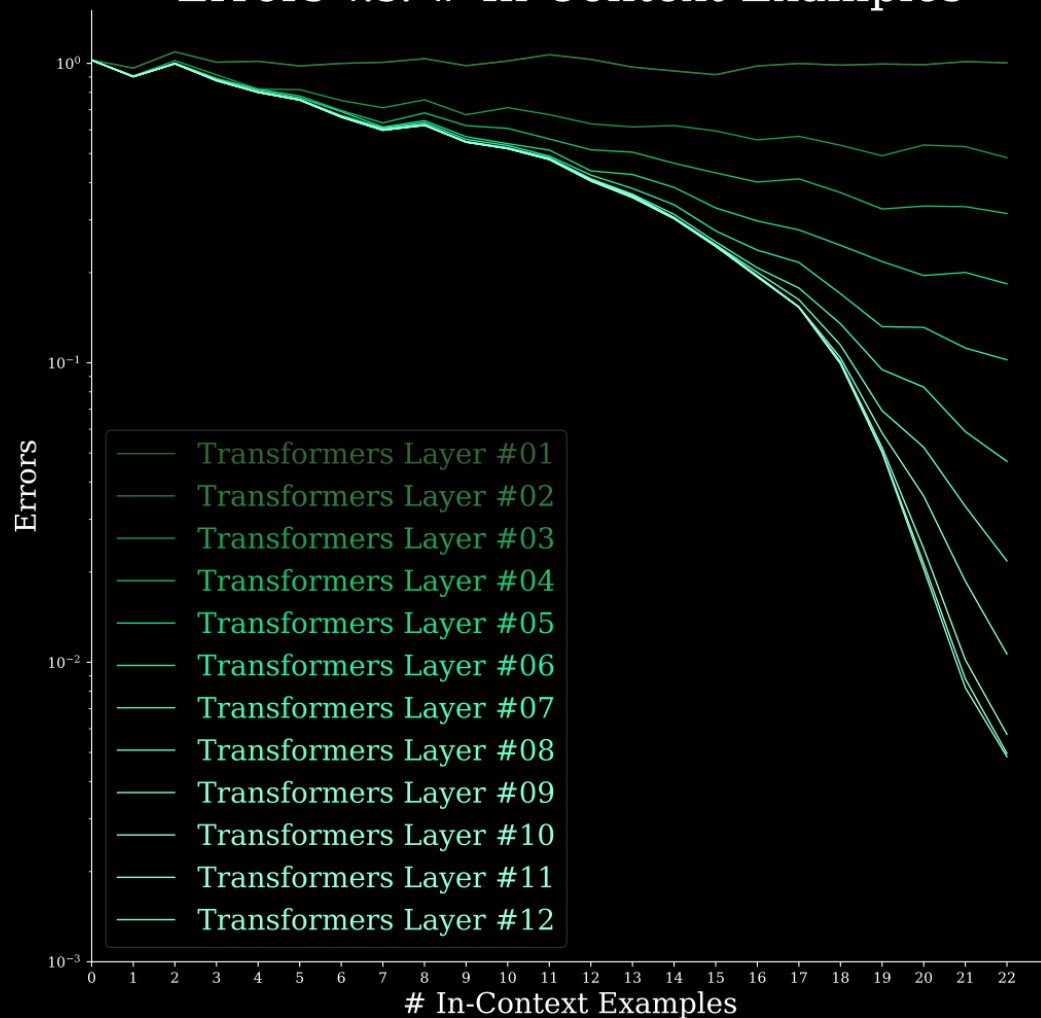Transformer layer

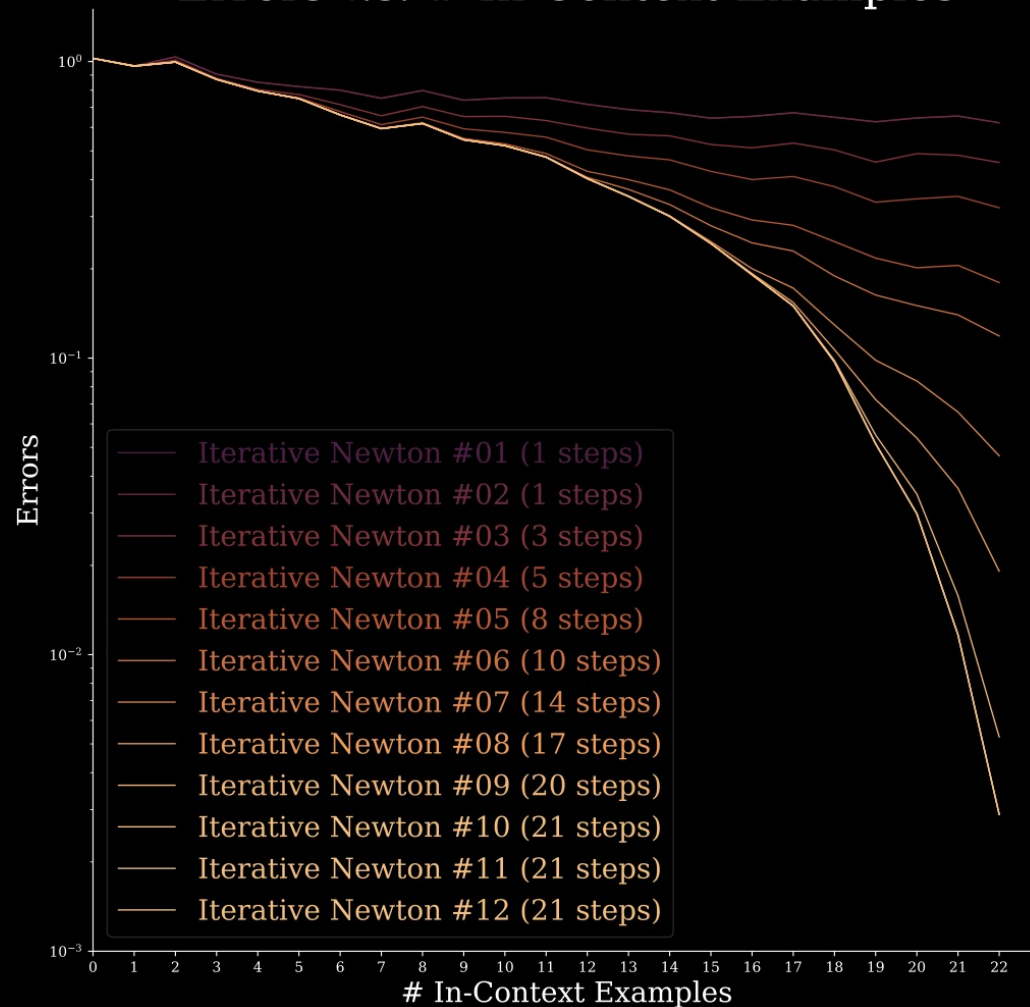Transformer layer

Transformer layer

Input sequence

$\approx$

# This work: Transformers do in-context learning via an iterative second-order method



**Errors v.s. # In-Context Examples** (left panel)

Legend:
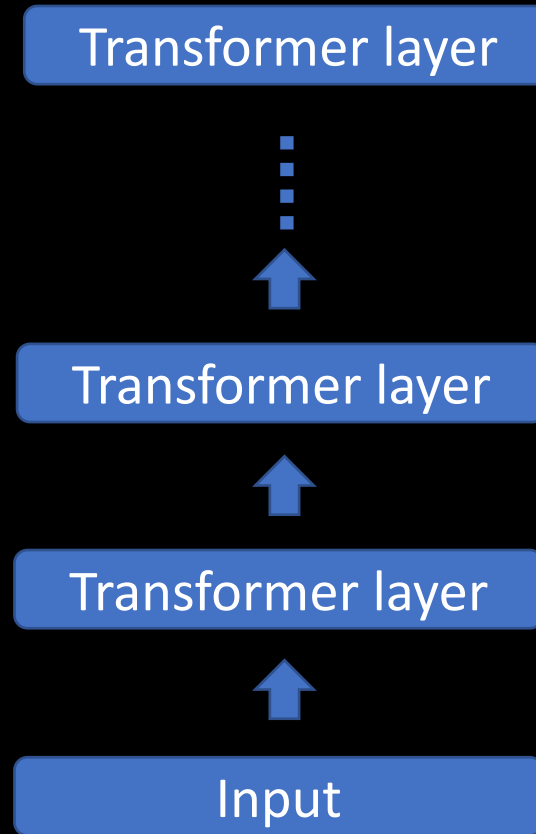- Transformers Layer #01
- Transformers Layer #02
- Transformers Layer #03
- Transformers Layer #04
- Transformers Layer #05
- Transformers Layer #06
- Transformers Layer #07
- Transformers Layer #08
- Transformers Layer #09
- Transformers Layer #10
- Transformers Layer #11
- Transformers Layer #12

Axes: Errors (y-axis), # In-Context Examples (x-axis)

**Errors v.s. # In-Context Examples** (right panel)

Legend:
- Iterative Newton #01 (1 steps)
- Iterative Newton #02 (1 steps)
- Iterative Newton #03 (3 steps)
- Iterative Newton #04 (5 steps)
- Iterative Newton #05 (8 steps)
- Iterative Newton #06 (10 steps)
- Iterative Newton #07 (14 steps)
- Iterative Newton #08 (17 steps)
- Iterative Newton #09 (20 steps)
- Iterative Newton #10 (21 steps)
- Iterative Newton #11 (21 steps)
- Iterative Newton #12 (21 steps)

Axes: Errors (y-axis), # In-Context Examples (x-axis)
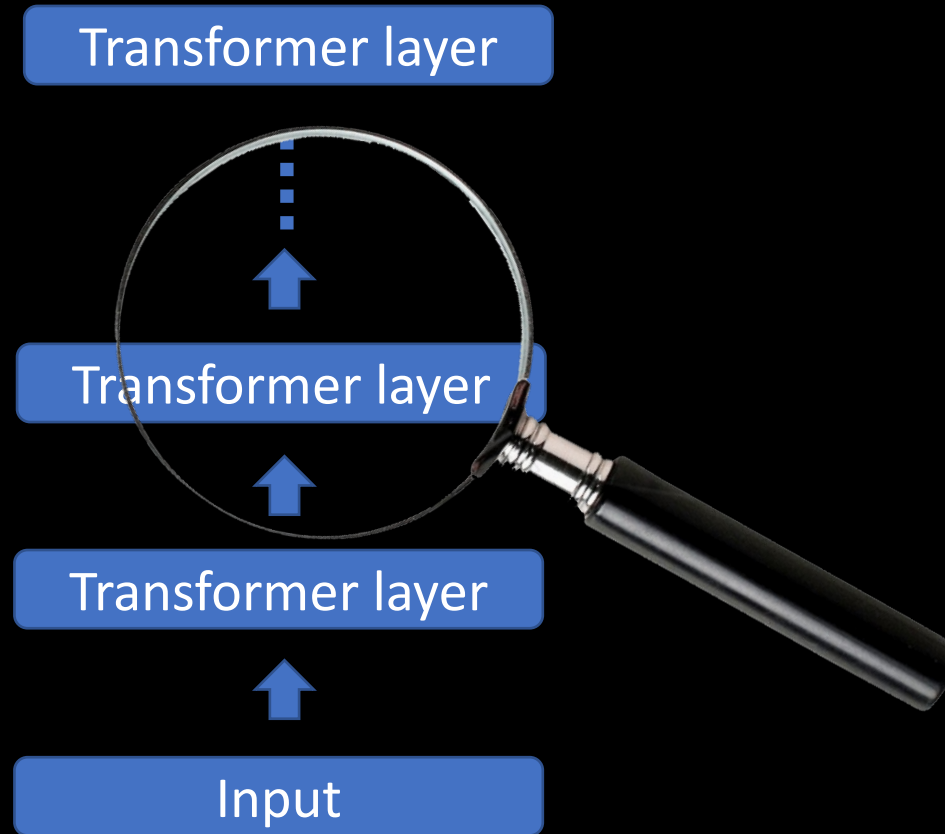
Techniques: "Applied theory"?

# Techniques: "Applied theory"?

How should we understand how Transformers solve a problem?
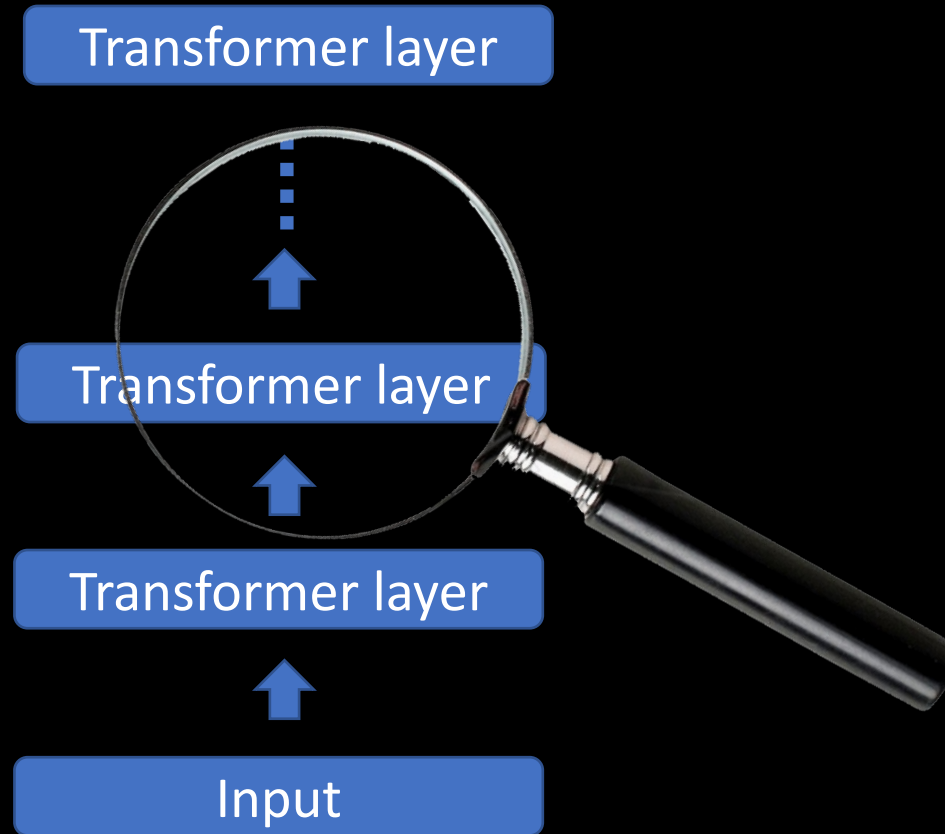
# Techniques: "Applied theory"?

How should we understand how Transformers solve a problem?



Inspect weights to invert mechanism?

# Techniques: "Applied theory"?

How should we understand how Transformers solve a problem?



Issue: Space of possible solutions can be too large and complex

# Techniques: "Applied theory"?

How should we understand how Transformers solve a problem?



One Solution: Using understanding of information and computation can refine search

# Techniques: "Applied theory"?

For linear regression:

- We know information-theoretic lower bounds on rates achievable by any first-order method
- We understand settings where gap between first and second-order methods is largest

Can we use this understanding, combined with empirical investigations, to uncover Transformer mechanisms?

# Setup and algorithms

$w^{*T}x$

$(x_1, y_1)$

$(x_2, y_2)$

$(x_3, y_3)$

$(x_4, y_4)$

# The Setup

## Data distribution

For each sequence of $n$ examples $\{x_i, y_i\}_{i=1}^{n}$

Sample $w^* \sim N(0, I)$
Sample data covariance $\Sigma$ (for now, let $\Sigma = I$)
For each $i \in [n], x_i \sim N(0, \Sigma), y_i = {w^*}^T x_i$

$w^{*T}x$

$(x_1, y_1)$

$(x_2, y_2)$

$(x_3, y_3)$

$(x_4, y_4)$

# Some algorithms for linear regression

For any time step $t$, let $X$ be matrix of datapoints, $y$ be vector of labels

Ordinary Least Squares: Minimum norm solution to sum of squares objective

$$w_{OLS} = (X^T X)^\dagger X^T y$$

Gradient descent on sum of squares objective:

$$w_{GD}^{(k+1)} = w_{GD}^{(k)} - \eta * (\text{Gradient at } w_{GD}^{(k)})$$

$O(\log(\frac{1}{\epsilon}))$ iterations to find $\epsilon$ accurate solution

Iterative Newton's: Iterative 2$^{\text{nd}}$ order method to find inverse ($\approx$ matrix Taylor series)

Let $S = X^T X$

$O(\log \log(\frac{1}{\epsilon}))$ iterations to find $\epsilon$ accurate solution

$$M_0 = \alpha S, M_{k+1} = 2M_k - M_k S M_k$$
$$w_{Newton}^{(k)} = M_k X^T y$$

# Transformers for linear regression



Train on $T$ such instances

# Transformers for linear regression

# Transformers as an iterative algorithm: probing layers

# Transformers as an iterative algorithm: probing layers

# Transformers as an iterative algorithm: probing layers



Errors v.s. # In-Context Examples

Legend:
- Transformers Layer #01
- Transformers Layer #02
- Transformers Layer #03
- Transformers Layer #04
- Transformers Layer #05
- Transformers Layer #06
- Transformers Layer #07
- Transformers Layer #08
- Transformers Layer #09
- Transformers Layer #10
- Transformers Layer #11
- Transformers Layer #12

Y-axis: Errors

X-axis: # In-Context Examples

# Metric: Similarity of errors

$$x_1, \ y_1, x_2, y_2, x_3, y_3, \quad \ldots, x_n, y_n,$$

Algorithm A
$$y_1^A, \quad y_2^A, \quad y_3^A, \quad \ldots, y_n^A,$$

Algorithm B
$$y_1^B, \quad y_2^B, \quad y_3^B, \quad \ldots, y_n^B,$$

Algorithm A residuals
$$(y_1 - y_1^A), (y_2 - y_2^A), (y_3 - y_3^A), \ldots, (y_n - y_n^A),$$

Algorithm B residuals
$$(y_1 - y_1^B), (y_2 - y_2^B, \ (y_3 - y_3^B), \ldots, (y_n - y_n^B),$$

Similarity of errors on $\{x_i, y_i\}_{i=1}^n$ between Algorithm A, Algorithm B
= Cosine similarity between residuals of A , B

Overall similarity of errors (Algorithm A, Algorithm B)
= $\mathbb{E}_{\{x_i, y_i\}}$ [ Cosine similarity between residuals of A , B ]

# Claim 1: Transformers improve across layers



Errors v.s. # In-Context Examples

# Claim 2: Transformers are more similar to Iterative Newton than to GD



Transformers vs Newton

Transformers vs Gradient Descent

# Claim 2: Transformers are more similar to Iterative Newton than to GD



Transformers vs Newton

Transformers vs Gradient Descent

# Claim 2: Transformers are more similar to Iterative Newton than to GD



Transformers vs Newton

Transformers vs Gradient Descent

# Claim 3: Transformers are still able to match Newton on harder distributions

What is a setting where the gap between 1st and 2nd order methods is especially large?

On **ill-conditioned instances**, gradient descent (or its variants) get $poly(\kappa)$ dependence on the condition number of the linear system $\kappa$, 2nd order methods get $polylog(\kappa)$ dependence.

# Claim 3: Transformers are still able to match Newton on harder distributions

What is a setting where the gap between $1^{st}$ and $2^{nd}$ order methods is especially large?

On **ill-conditioned instances**, gradient descent (or its variants) get $poly(\kappa)$ dependence on the condition number of the linear system $\kappa$, $2^{nd}$ order methods get $polylog(\kappa)$ dependence.

Conjecture (Sharan-Sidford-Valiant'19): No first-order (linear memory method) can avoid a $poly(\kappa)$ dependence on $\kappa$ in general.

Hard distribution: Sample $\Sigma$ with $d/2$ eigenvalues at 100, $d/2$ eigenvalues at 1, uniformly random eigenbasis.

# Claim 3: Transformers are still able to match Newton on ill-conditioned data



Errors v.s. # In-Context Examples

Legend:
- Transformers Layer #01
- Transformers Layer #02
- Transformers Layer #03
- Transformers Layer #04
- Transformers Layer #05
- Transformers Layer #06
- Transformers Layer #07
- Transformers Layer #08
- Transformers Layer #09
- Transformers Layer #10
- Transformers Layer #11
- Transformers Layer #12
- Ordinary Least Squares
- Iterative Newton (21 Steps)
- Gradient Descent (800 Steps)

# Claim 3: Transformers are still able to match Newton on ill-conditioned data

Transformers vs Newton

Transformers vs Gradient Descent

# Theoretical justification

Can Transformers efficiently implement Iterative Newton's?

Informal Theorem:
Transformers can match predictions of $k$ steps of Iterative Newton's with $(k + 8)$ layers, $O(d)$ hidden units per layer.

Construction uses ideas from Akyurek-Schuurmans-Andreas-Ma-Zhou'2022, and is similar to a matrix inverse construction by Giannou-Rajput-Sohn-Lee-Lee-Papailiopoulos'2023
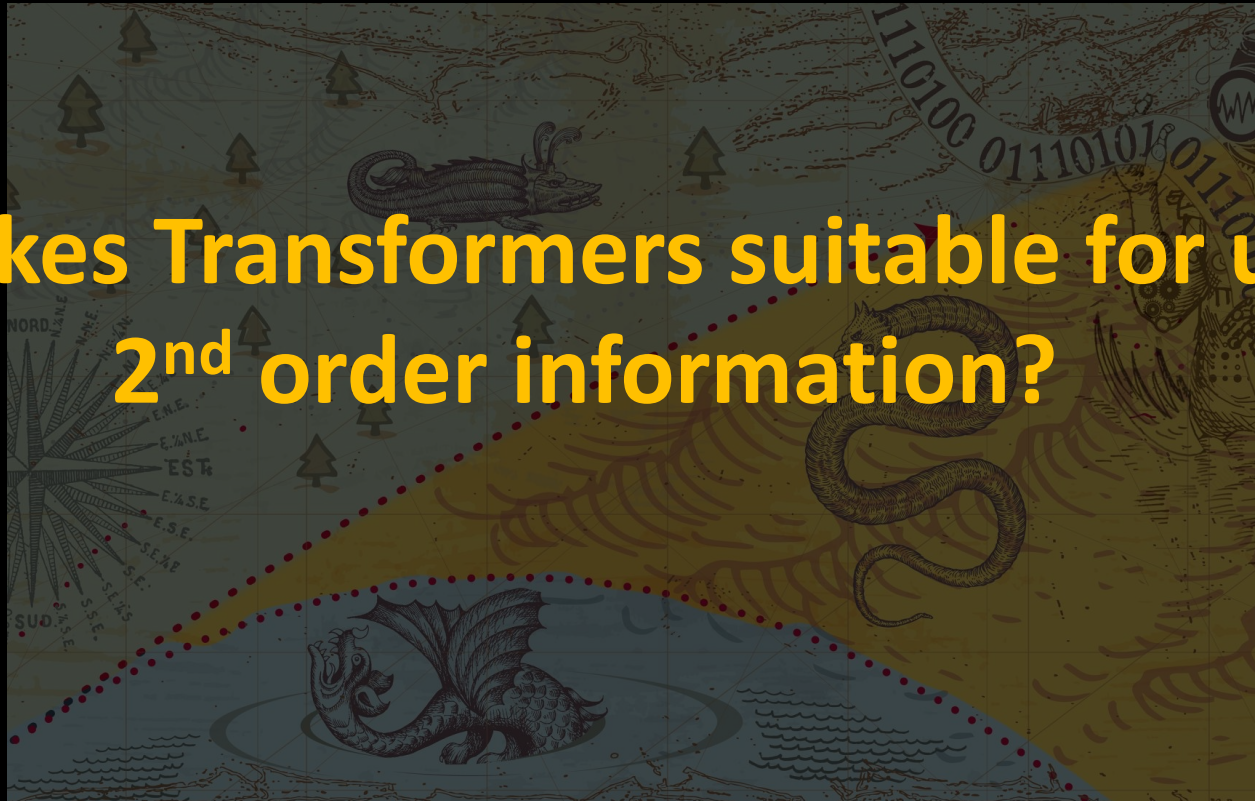
# Some more related work

Ahn-Cheng-Daneshmand-Sra'2023, Zhang-Frei-Bartlett'2023 & Mahankali-Hashimoto-Ma'2024 analyze dynamics of trained one-layer Transformers
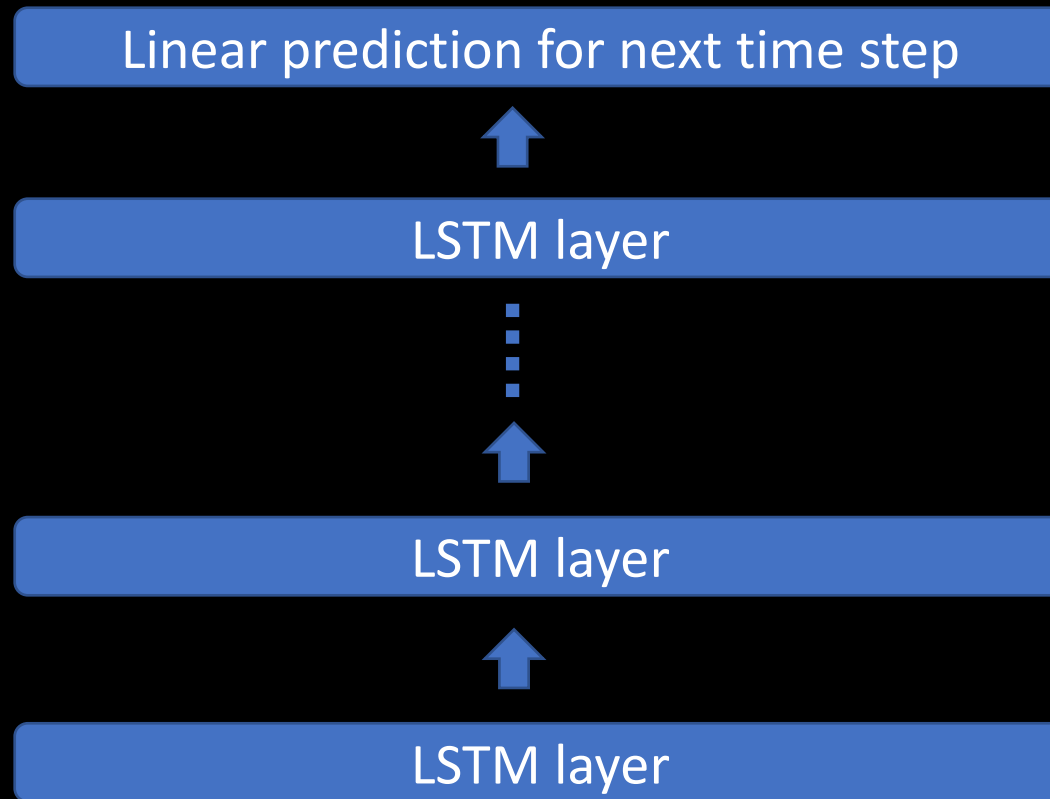
Vladymyrov-von Oswald-Sandler-Ge'2024 show that a second-order variant of GD can mimic Iterative Newton by implicitly approximating inverse

Giannou-Yang-Wang-Papailiopoulos-Lee'2024 show that Transformers can do Iterative Newton beyond linear regression

What makes Transformers suitable for utilizing 2nd order information?

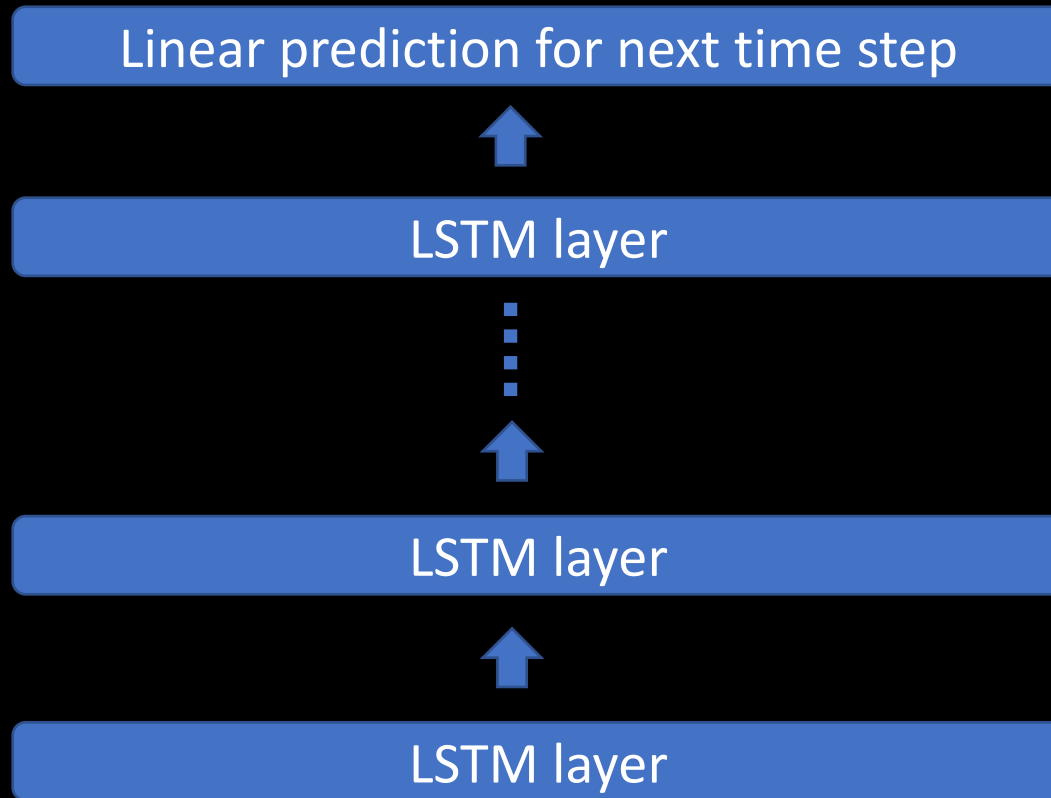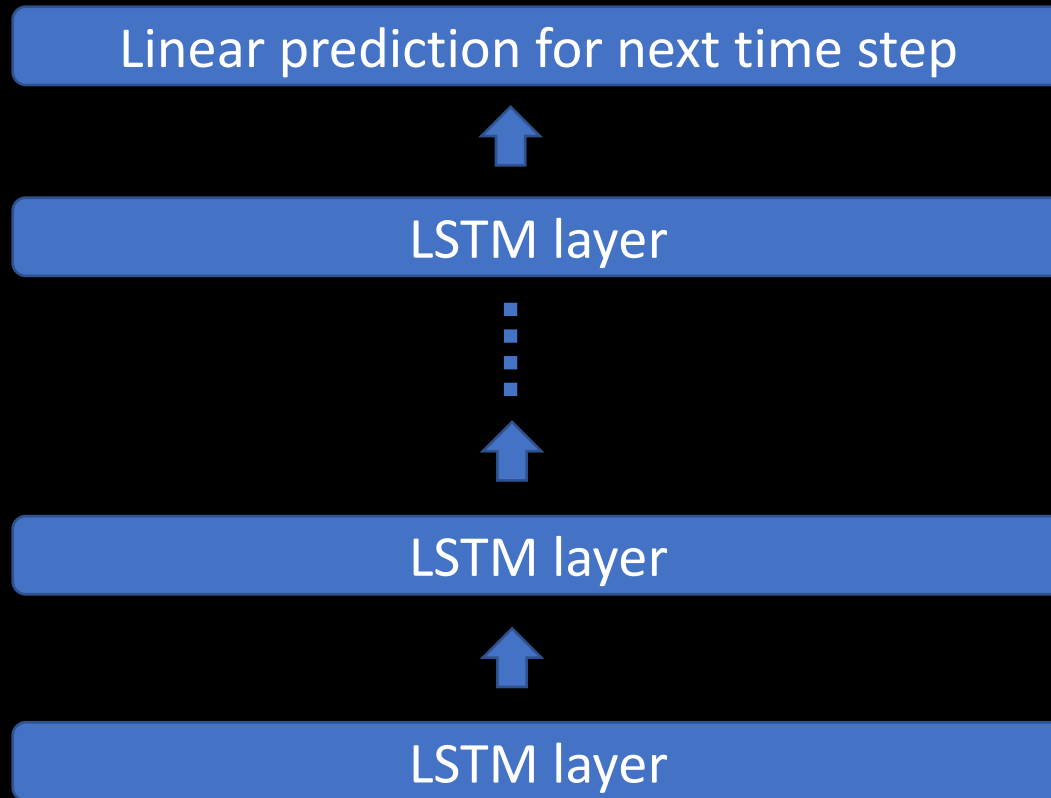# LSTMs for linear regression

# LSTMs for linear regression

# LSTMs for linear regression

Linear prediction for next time step

↑

LSTM layer

⋮

LSTM layer

↑

LSTM layer

$$\begin{bmatrix} x_{t+1}^{(1)} \\ x_{t+1}^{(2)} \\ \vdots \\ x_{t+1}^{(d)} \end{bmatrix}$$

...

# LSTMs as an iterative algorithm: probing layers

LSTM layer $\rightarrow$ Linear prediction for $y_{t+1}$

LSTM layer $\rightarrow$ Train a linear model on activation to predict $y_{t+1}$

LSTM layer $\rightarrow$ Train a linear model on activation to predict $y_{t+1}$

$$\begin{bmatrix} x_{t+1}^{(1)} \\ x_{t+1}^{(2)} \\ \vdots \\ x_{t+1}^{(d)} \end{bmatrix}$$

# What do LSTMs implement?



Errors v.s. # In-Context Examples

Legend (left plot):
- LSTM Layer #01
- LSTM Layer #02
- LSTM Layer #03
- LSTM Layer #04
- LSTM Layer #05
- LSTM Layer #06
- LSTM Layer #07
- LSTM Layer #08
- LSTM Layer #09
- LSTM Layer #10

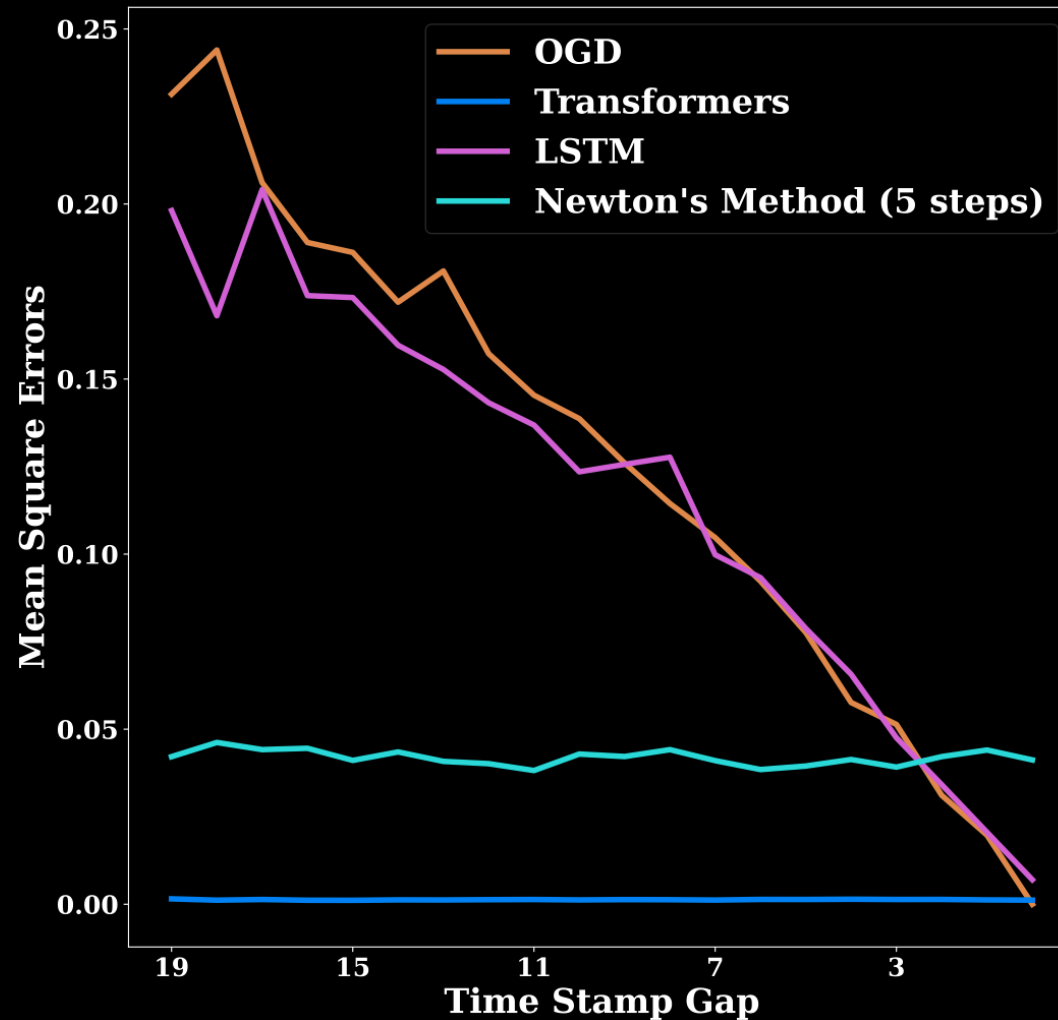Legend (right plot):
- OGD
- Transformers
- LSTM
- Newton's Method (5 steps)

LSTMs seem similar to online gradient descent

# Like OGD, LSTMs 'forget' previous examples



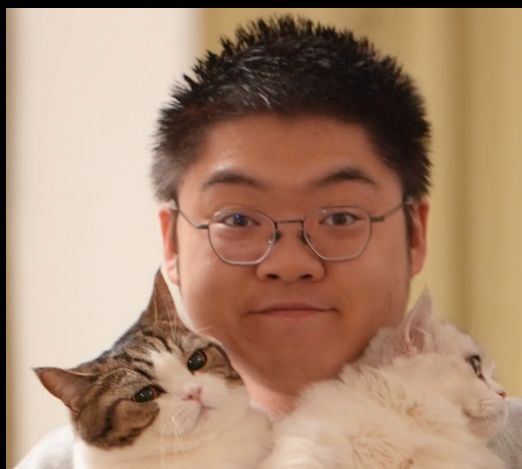Error when input from $t$ time steps ago is given as query point

**Hypothesis: The additional memory available to Transformers (since they have access to entire past sequence) versus recurrent architectures enables it to learn more efficient algorithm**

Recent line of theoretical work suggests that the available memory determines the best possible convergence rate, is gap between architectures an instantiation of this?
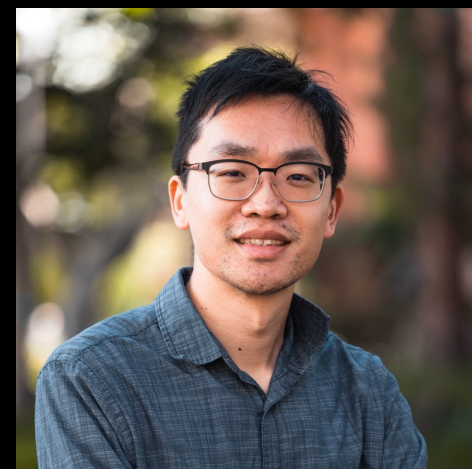
# What is the role of pre-training?
# *How do LLMs add?*

Tianyi Zhou (USC)    Deqing Fu (USC)    Robin Jia (USC)

Pre-trained LLMs Use Fourier Features to Compute Addition, Neurips 2024

# How do pre-trained Transformers do addition?

Fine-tune GPT-2XL on addition dataset:

- *What is the sum of 15 and 93? 108*
- *What is the sum of 24 and 171? 195*
*...*

# How do pre-trained Transformers do addition?

Fine-tune GPT-2XL on addition dataset:

- *What is the sum of 15 and 93? 108*
- *What is the sum of 24 and 171? 195*
...

Each number is its own token

# How do pre-trained Transformers do addition?

Fine-tune GPT-2XL on addition dataset:

- *What is the sum of 15 and 93? 108*
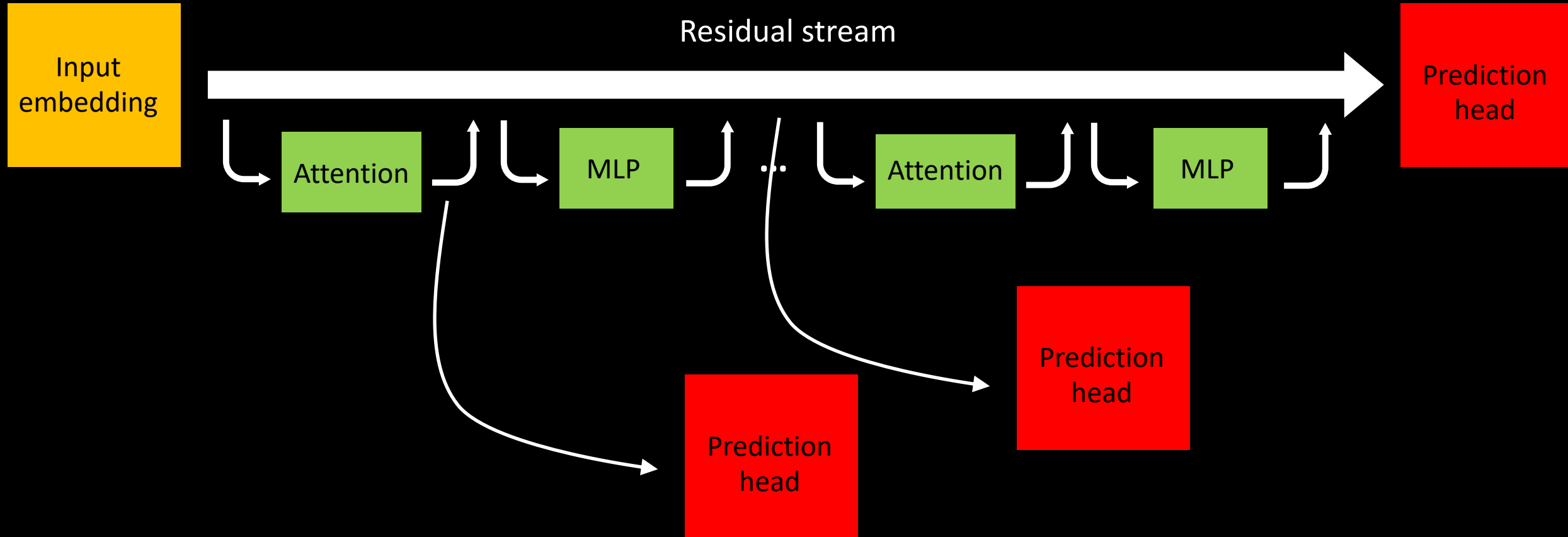- *What is the sum of 24 and 171? 195*
*...*

Model gets ≈ 100% test accuracy.

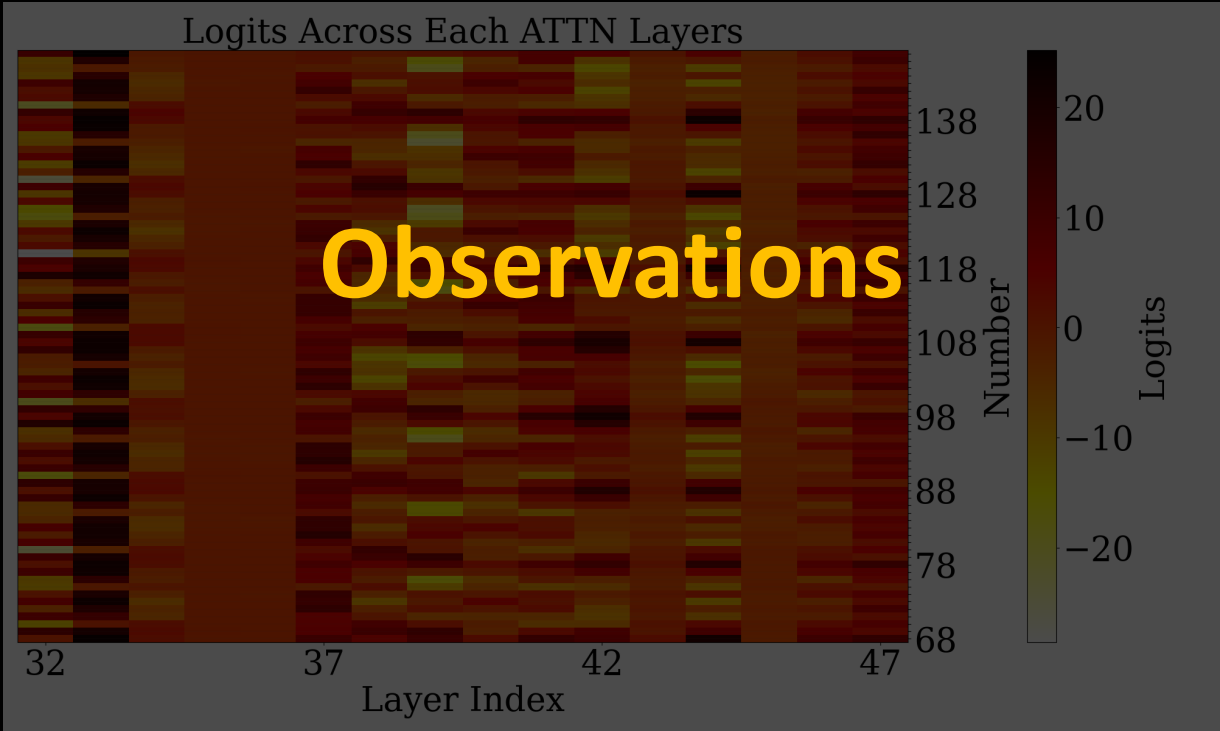What mechanisms does the model use?

# Understanding mechanisms: Logit Lens



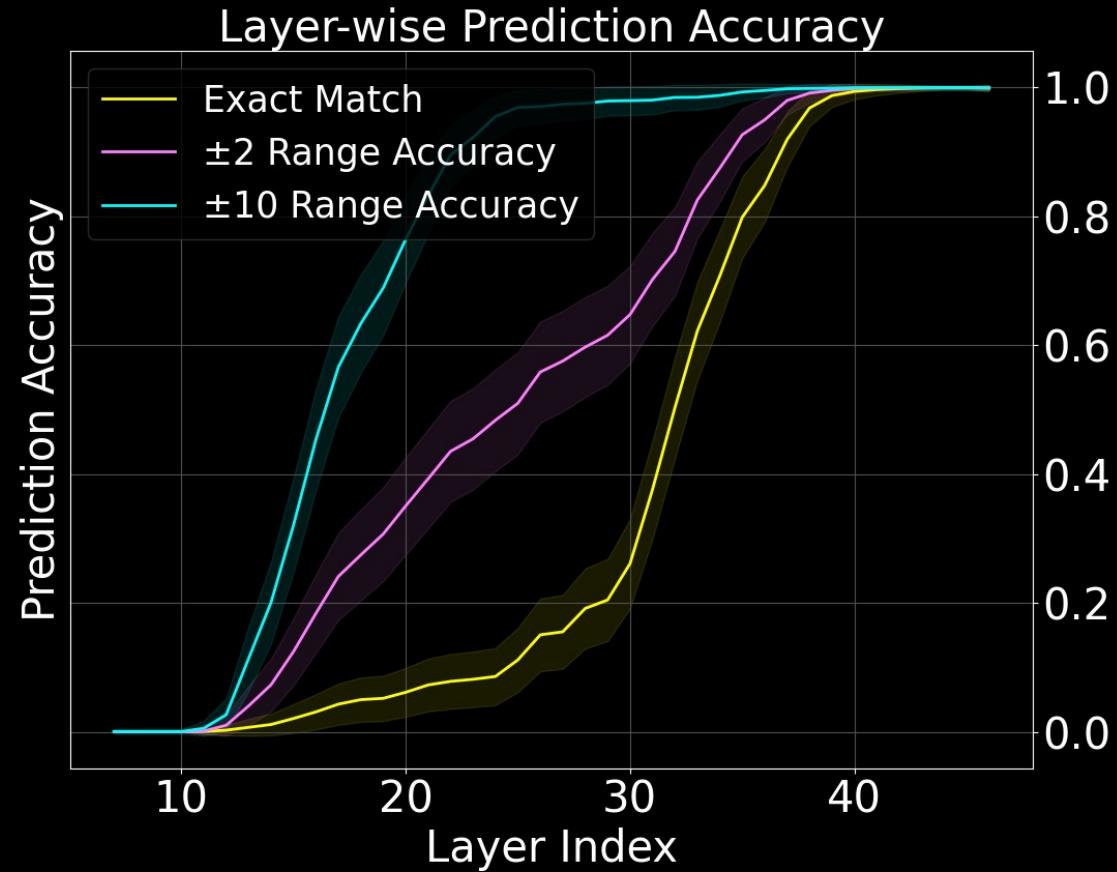Each Attention/MLP component makes additive contribution to residual stream

# Understanding mechanisms: Logit Lens



Residual stream

Input embedding

Attention

MLP

Attention

MLP

Prediction head

Prediction head

Prediction head

Can use prediction head to understand predictions at any stage

Logits Across Each ATTN Layers

# Model improves across layers



Layer-wise Prediction Accuracy

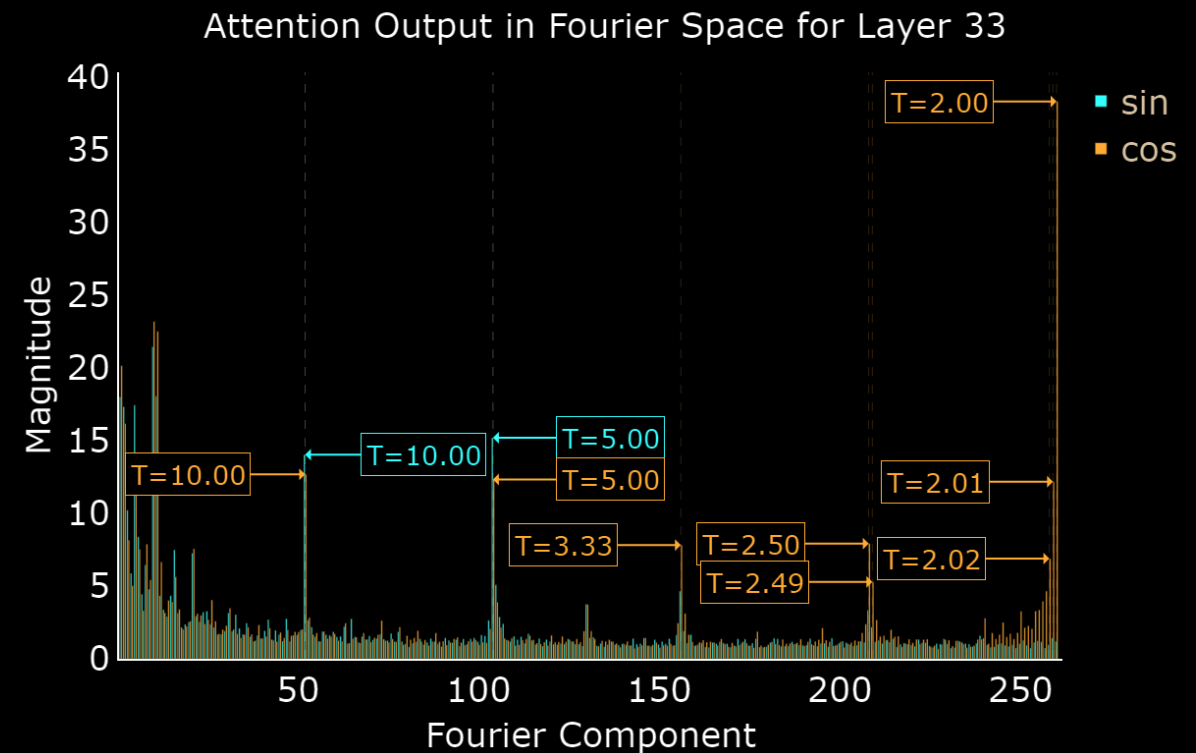Model finds answer within a ±2 and ±10 range early on, and finds exact match later
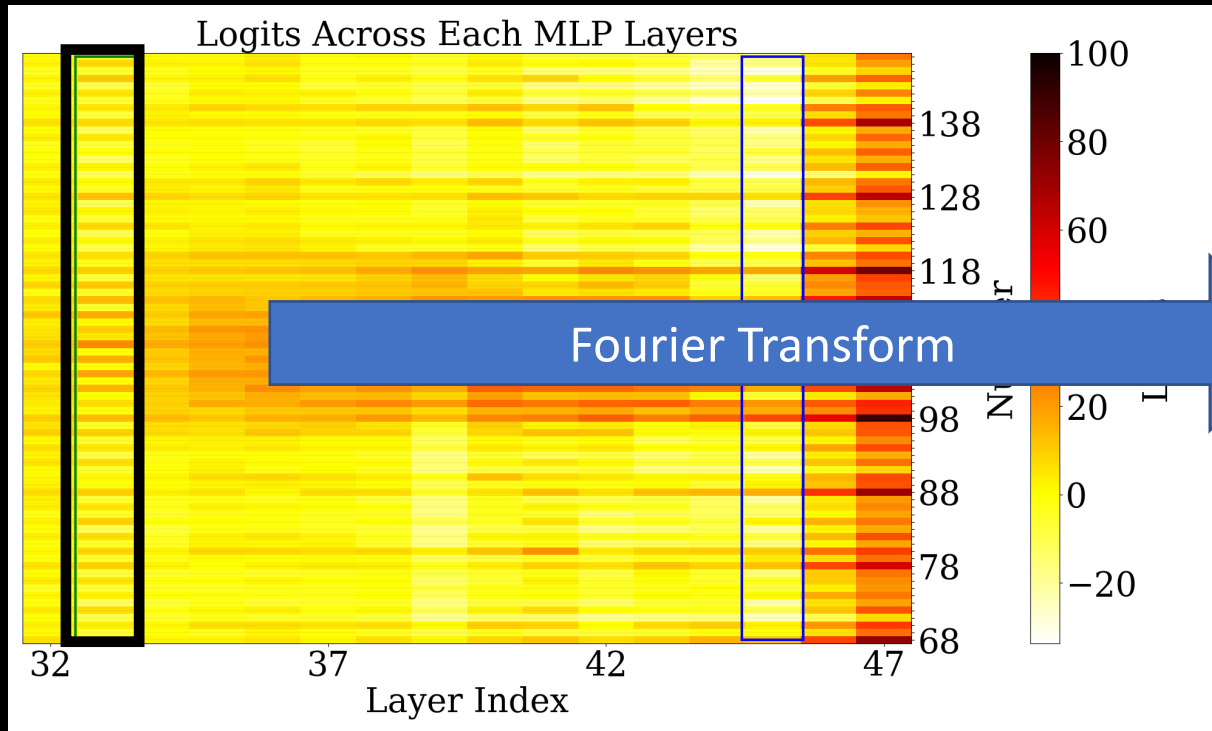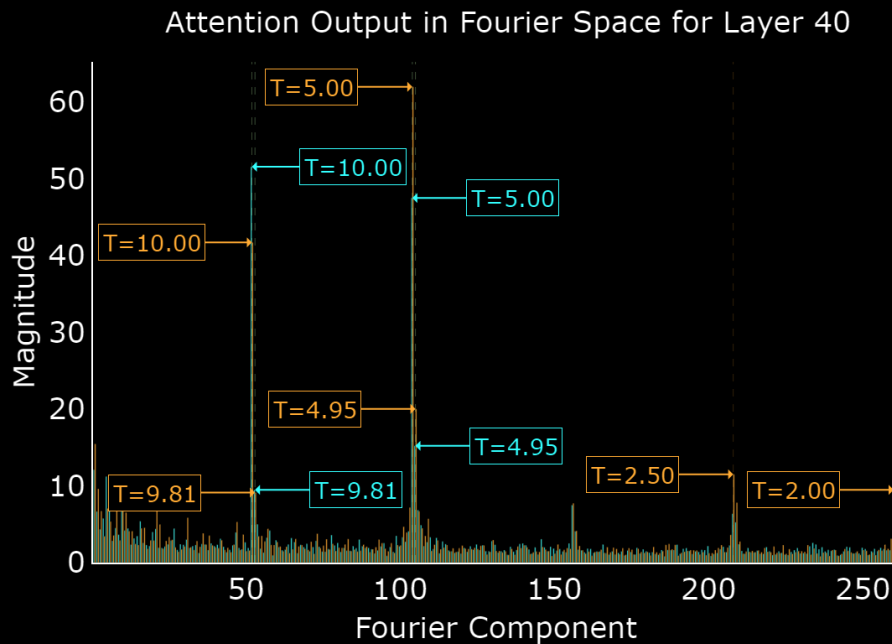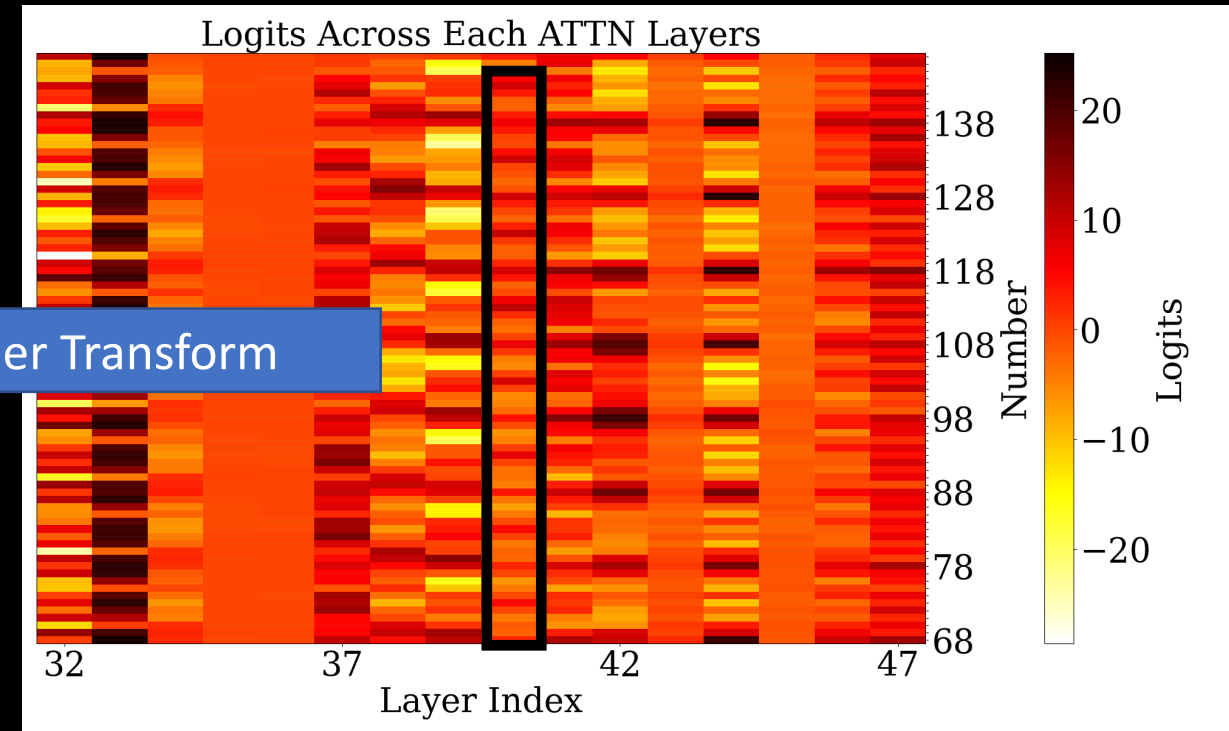
# Examining the contribution of each MLP & Attention layer

*Input: What is the sum of 15 and 93?*

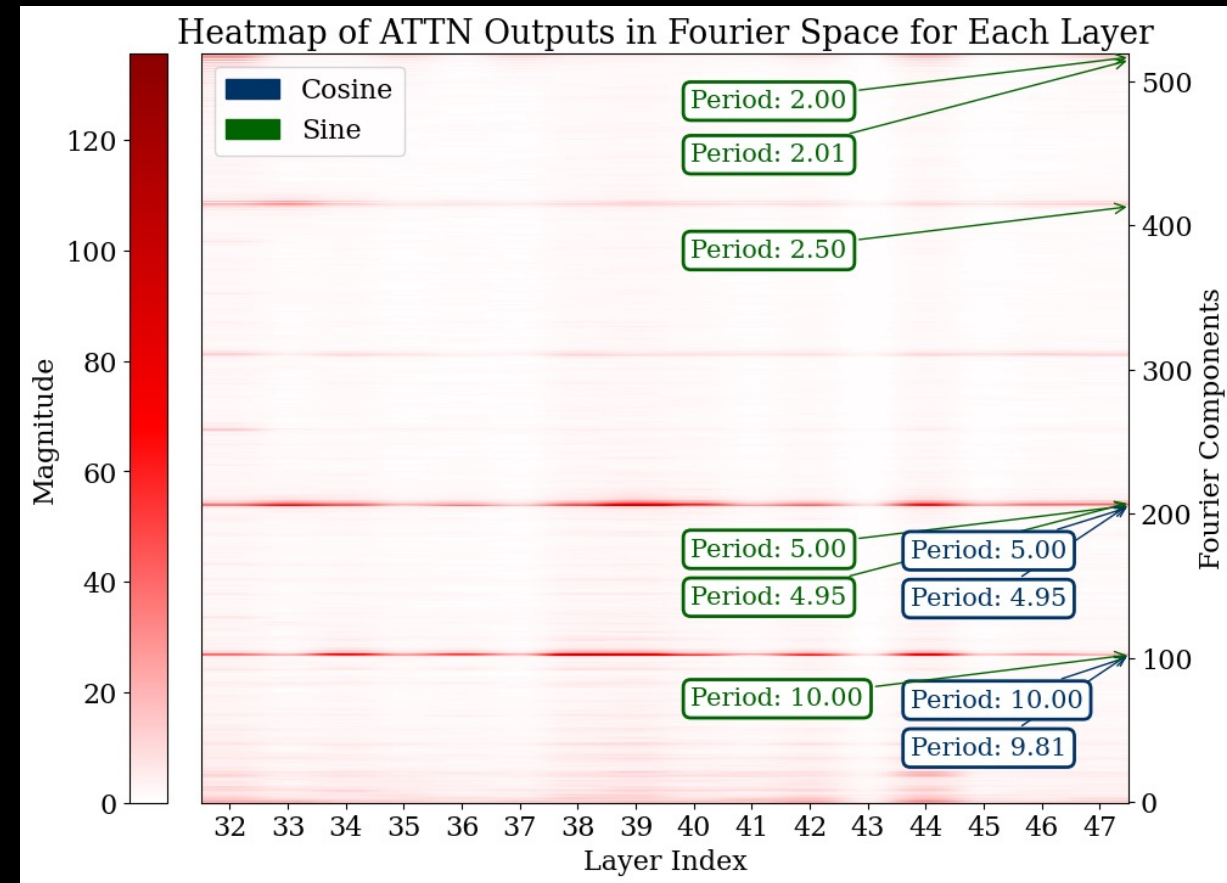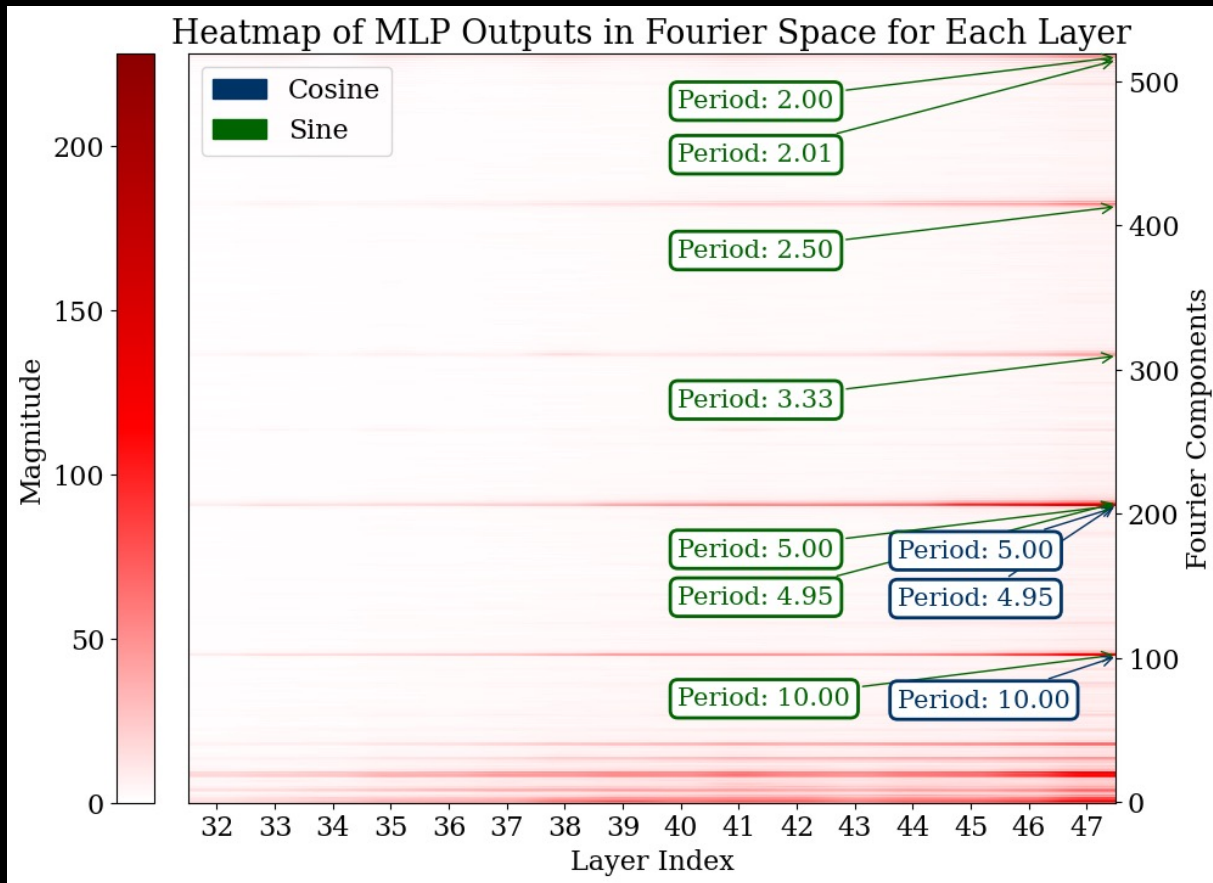# Examining the contribution of each MLP & Attention layer

*Input: What is the sum of 15 and 93?*

# Examining the contribution of each MLP & Attention layer

*Input: What is the sum of 15 and 93?*
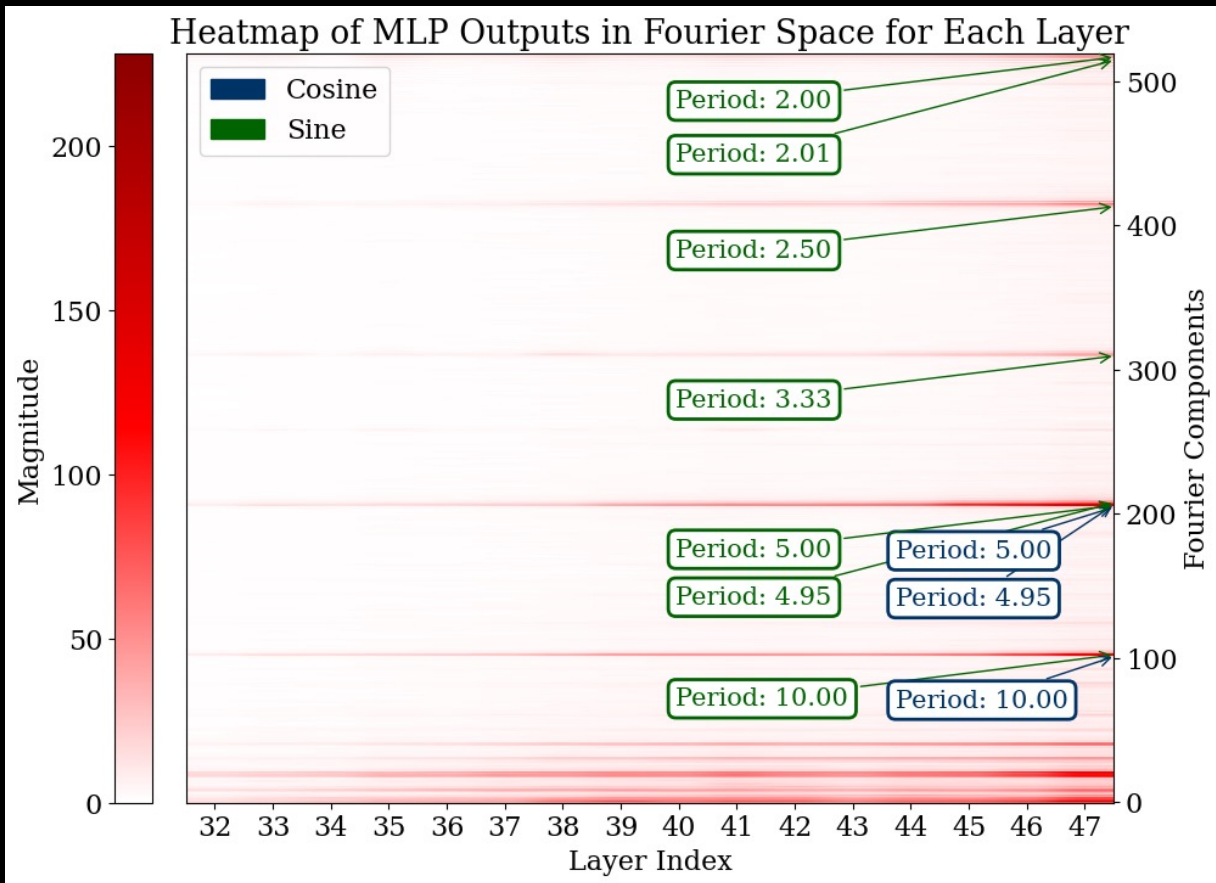
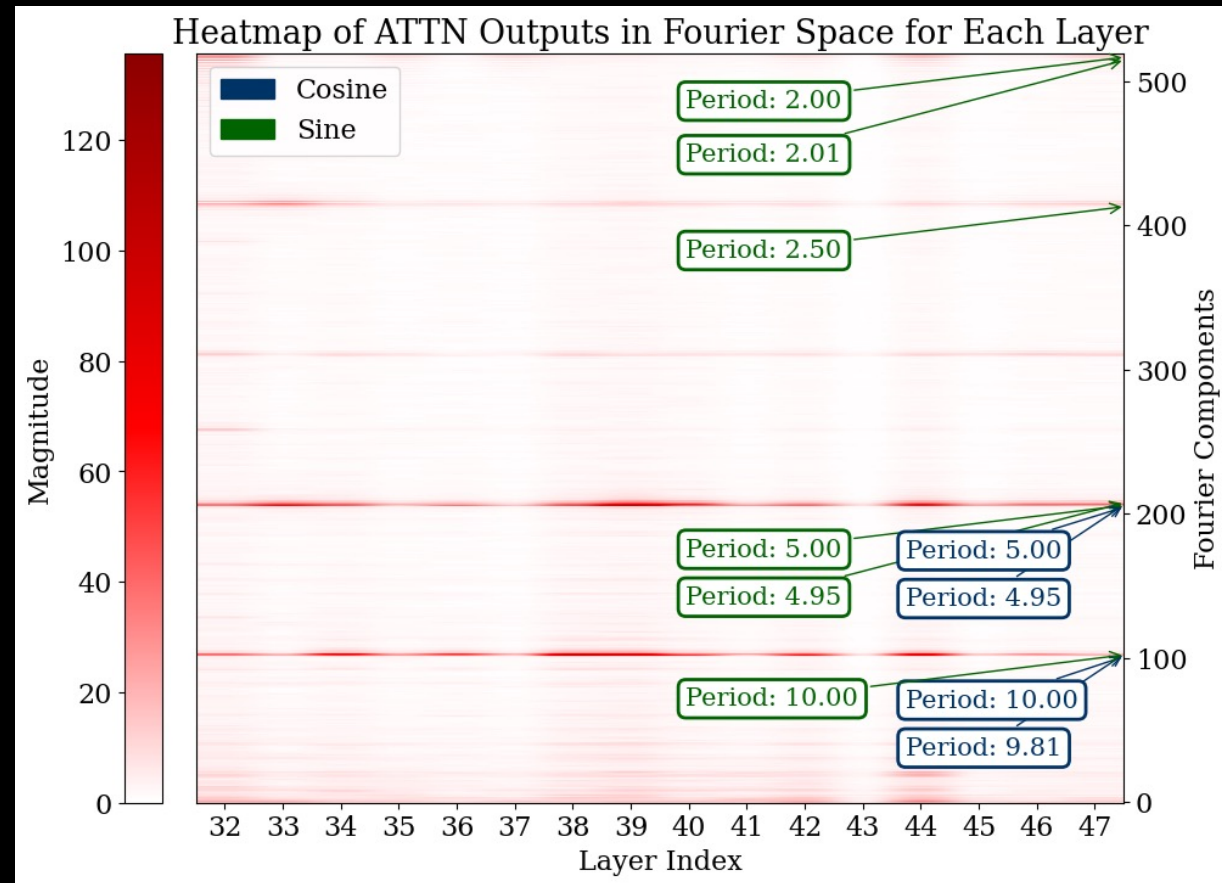# On average across all examples, logits are sparse in Fourier space

# Fourier features: Sparse representations in Fourier space



Heatmap of MLP Outputs in Fourier Space for Each Layer

Heatmap of ATTN Outputs in Fourier Space for Each Layer

Low frequency components **approximate** magnitude of answer

High frequency components do **classification**: compute sum modulo p for $p \in \{2, 5, 10, etc.\}$

# Fourier features: Sparse representations in Fourier space

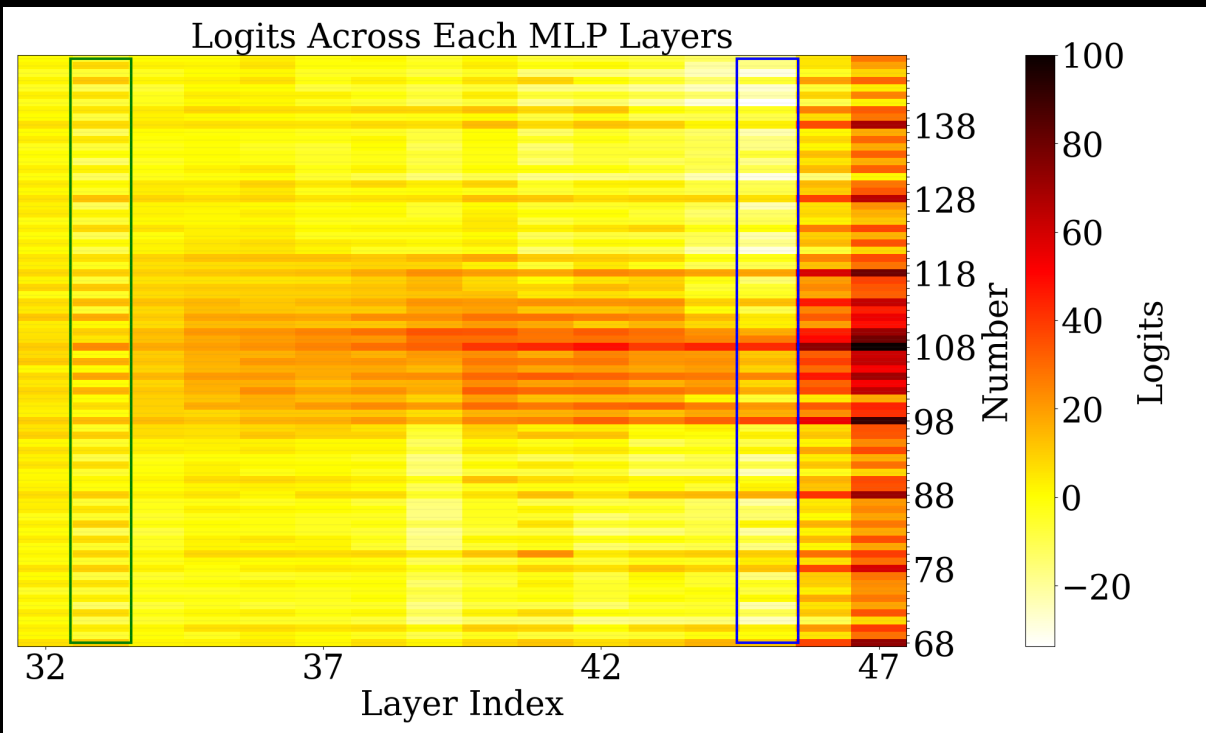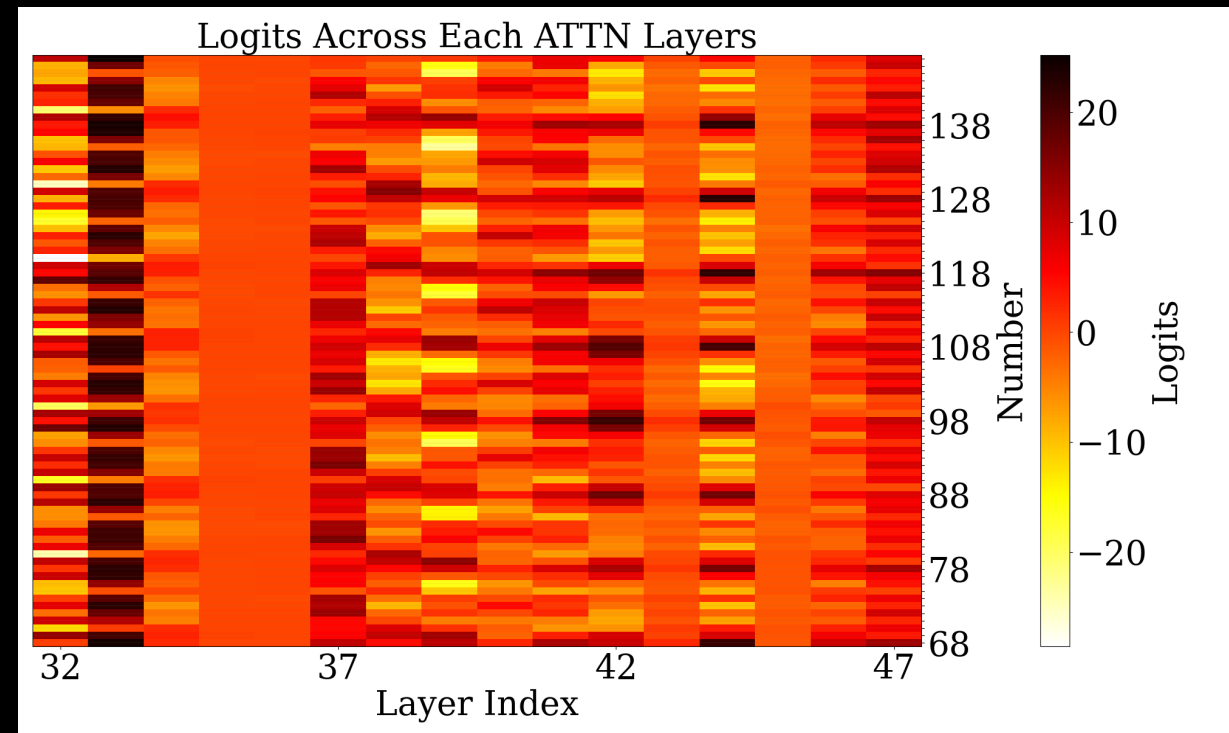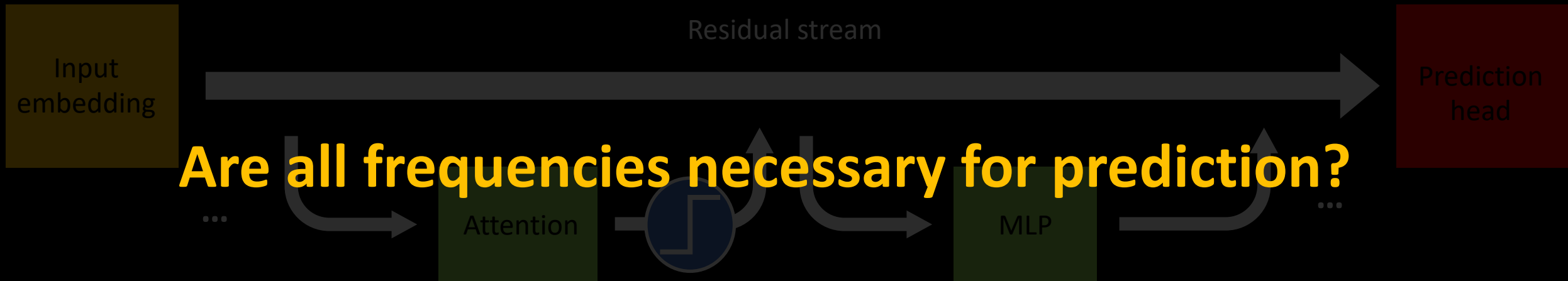*Input: What is the sum of 15 and 93?*
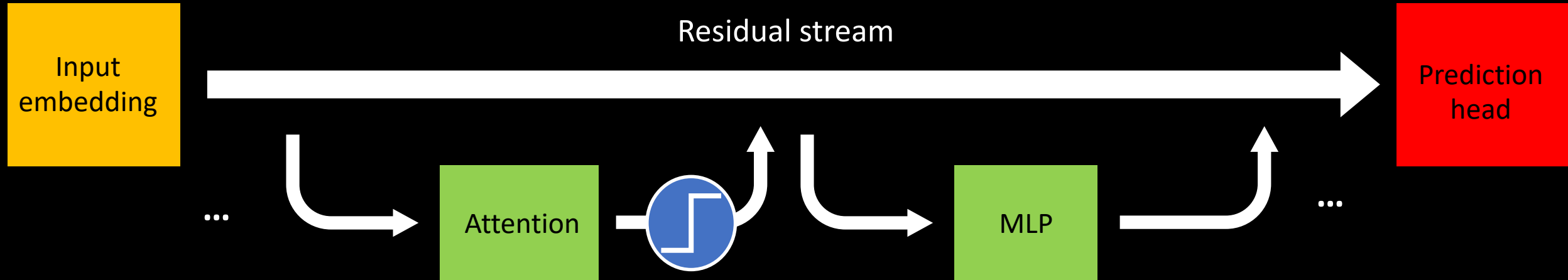


Low frequency components
**approximate** magnitude of answer

High frequency components do **classification**:
compute sum modulo p for $p \in \{2, 5, 10, etc.\}$

Are all frequencies necessary for prediction?

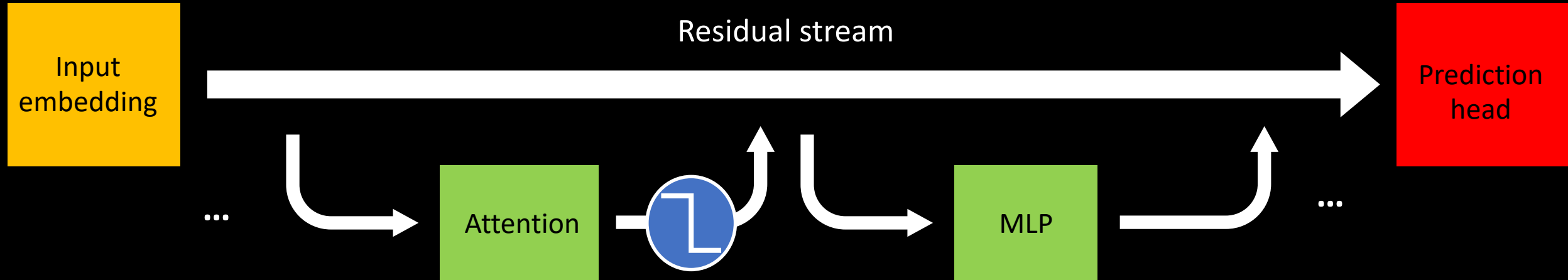Do Attention and MLP layers have similar roles?

# Applying filters to understand role of components



Residual stream

Input embedding

Prediction head

... Attention MLP ...

⟳ : High-pass filter to remove all low-frequency components in logit space

# Applying filters to understand role of components

Residual stream

Input embedding

Attention

MLP

Prediction head

...

...

⊔ : High-pass filter to remove all low-frequency components in logit space

⊓ : Low-pass filter to remove all high-frequency components in logit space

# Applying filters to understand role of components



Residual stream

Input embedding

Prediction head

Attention

MLP

⬙ : High-pass filter to remove all low-frequency components in logit space

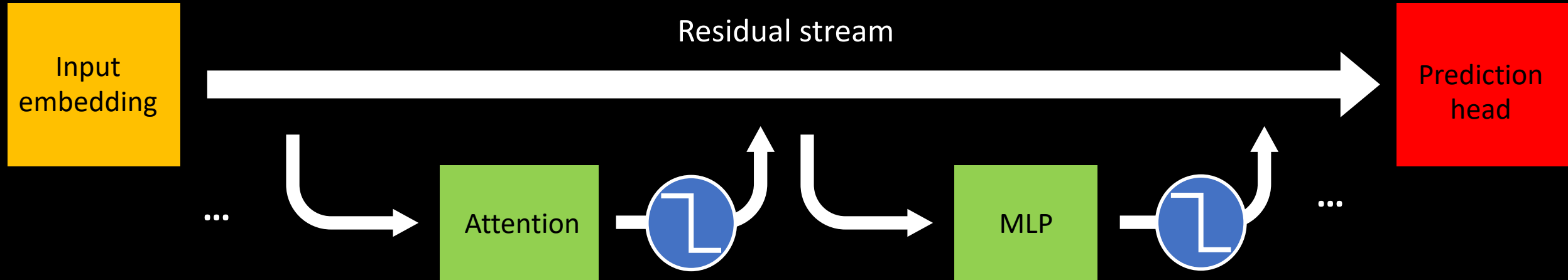⬙ : Low-pass filter to remove all high-frequency components in logit space
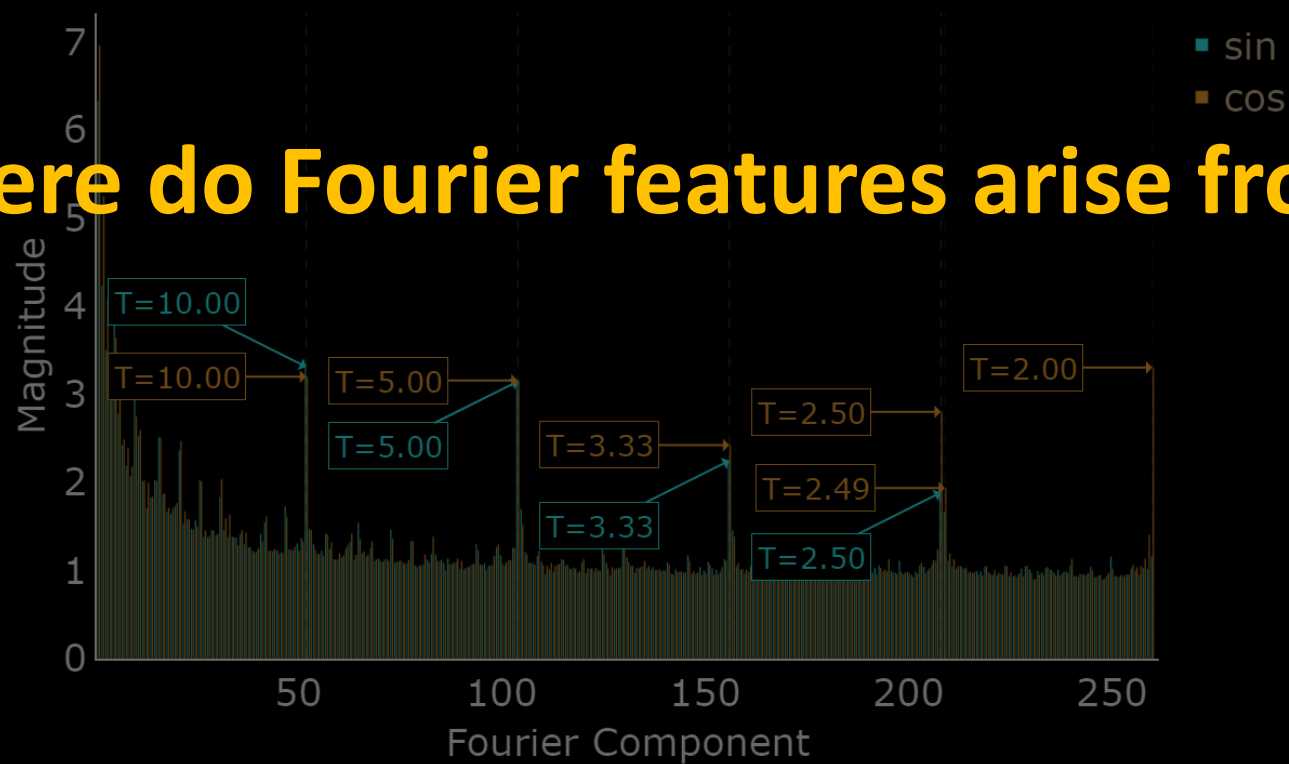
# Applying filters to understand role of components

| Module | Fourier Component Removed | Accuracy |
|---|---|---|
| None | No filtering | 99.74% |
| Attn & MLP | Low frequency | 5.94% |
| Attn | Low frequency | **99.12%** |
| MLP | Low frequency | 35.89% |
| Attn & MLP | High frequency | 27.08% |
| Attn | High frequency | 78.36% |
| MLP | High frequency | **98.10%** |

# MLP: mainly low-frequency, Attn: mainly high-frequency

| Module | Fourier Component Removed | Accuracy |
|---|---|---|
| None | No filtering | 99.74% |
| Attn & MLP | Low frequency | 5.94% |
| Attn | Low frequency | **99.12%** |
| MLP | Low frequency | 35.89% |
| Attn & MLP | High frequency | 27.08% |
| Attn | High frequency | 78.36% |
| MLP | High frequency | **98.10%** |

Number Embedding in Fourier Space: Pre-trained GPT-2-XL
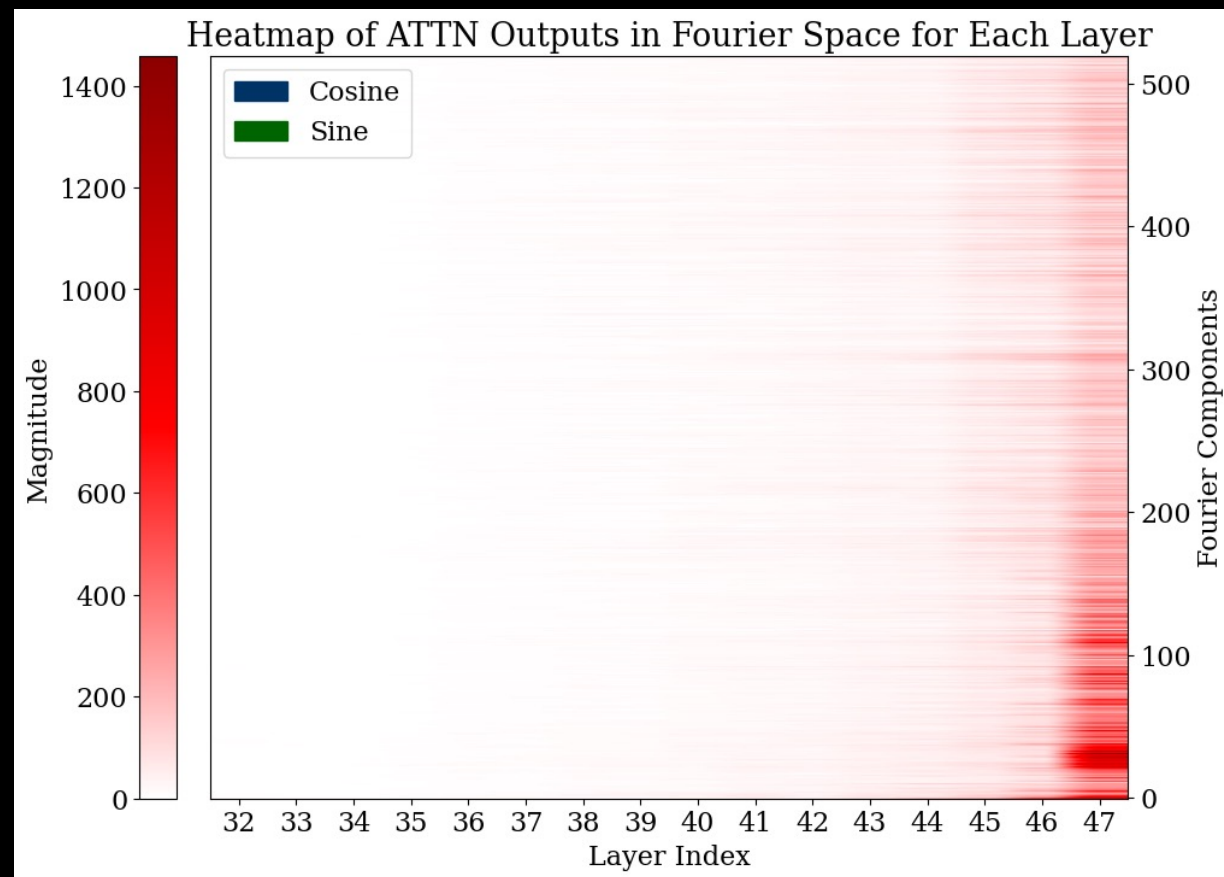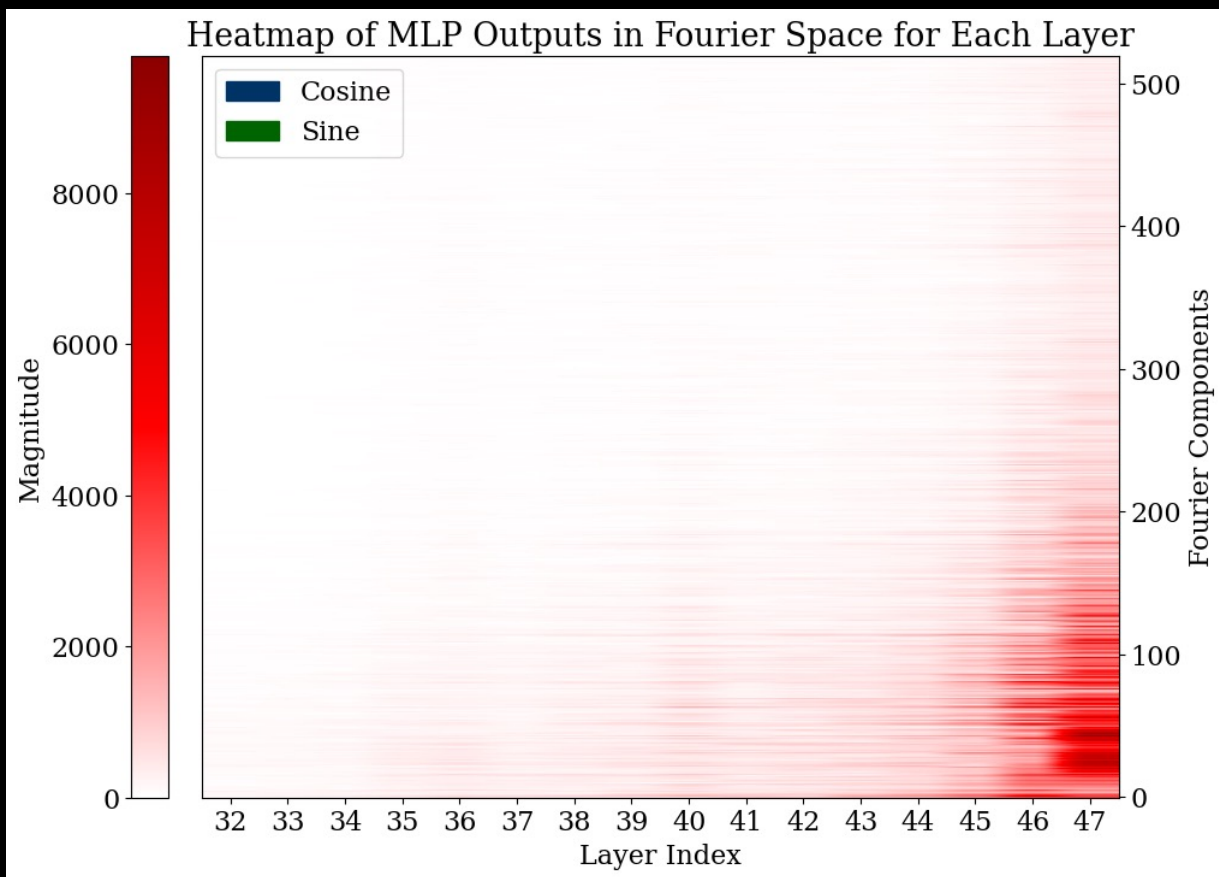
**Where do Fourier features arise from?**

# Appear to arise due to token embeddings from pre-training



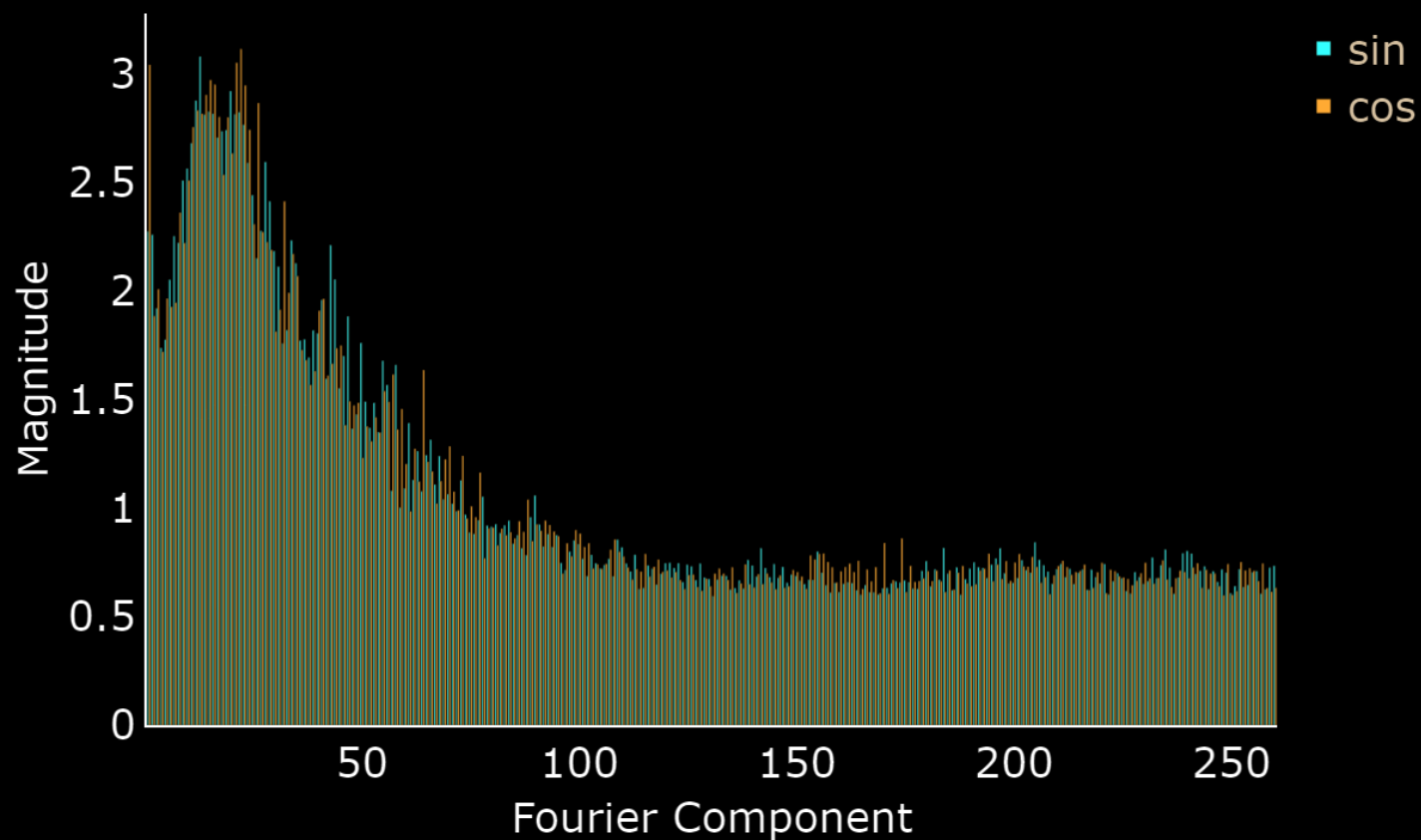Number Embedding in Fourier Space: Pre-trained GPT-2-XL

Also see similar behavior for other pre-trained models (Phi-2, RoBERTa).

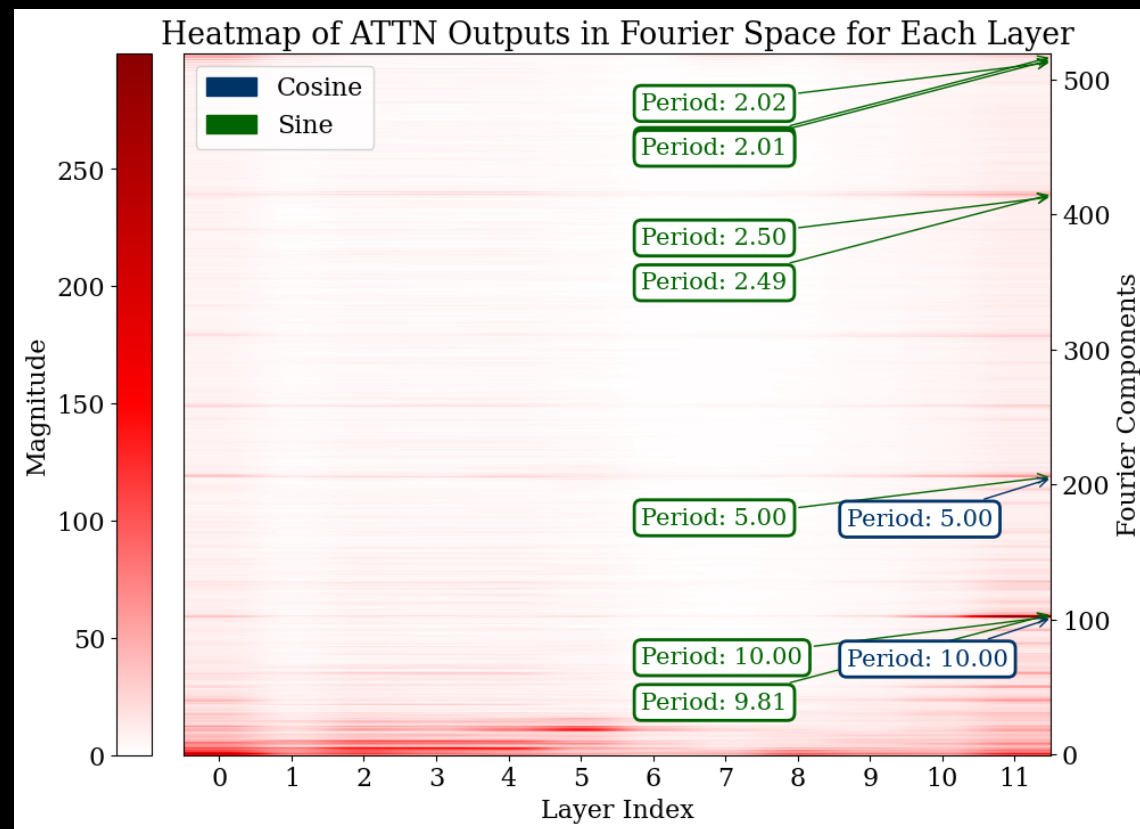# Model trained from scratch does not exhibit Fourier features
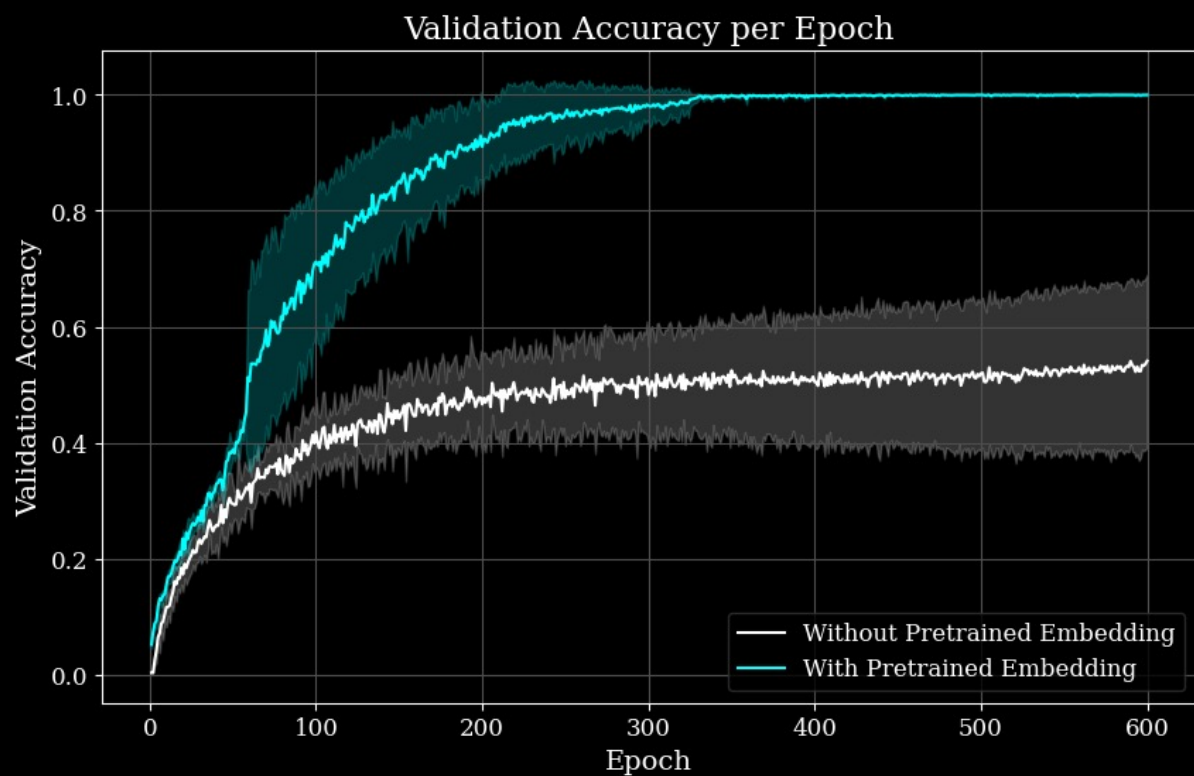
# Token embeddings of model trained from scratch do not have Fourier features either



Number Embedding in Fourier Space: GPT-2 Trained From Scratch

# Training model from scratch but with token embeddings from pre-trained models (a) improves training (b) leads to Fourier features

Related work: Fourier features in modular addition

Nanda-Chan-Lieberum-Smith-Steinhardt'2023 shows Fourier features are used in modular arithmetic
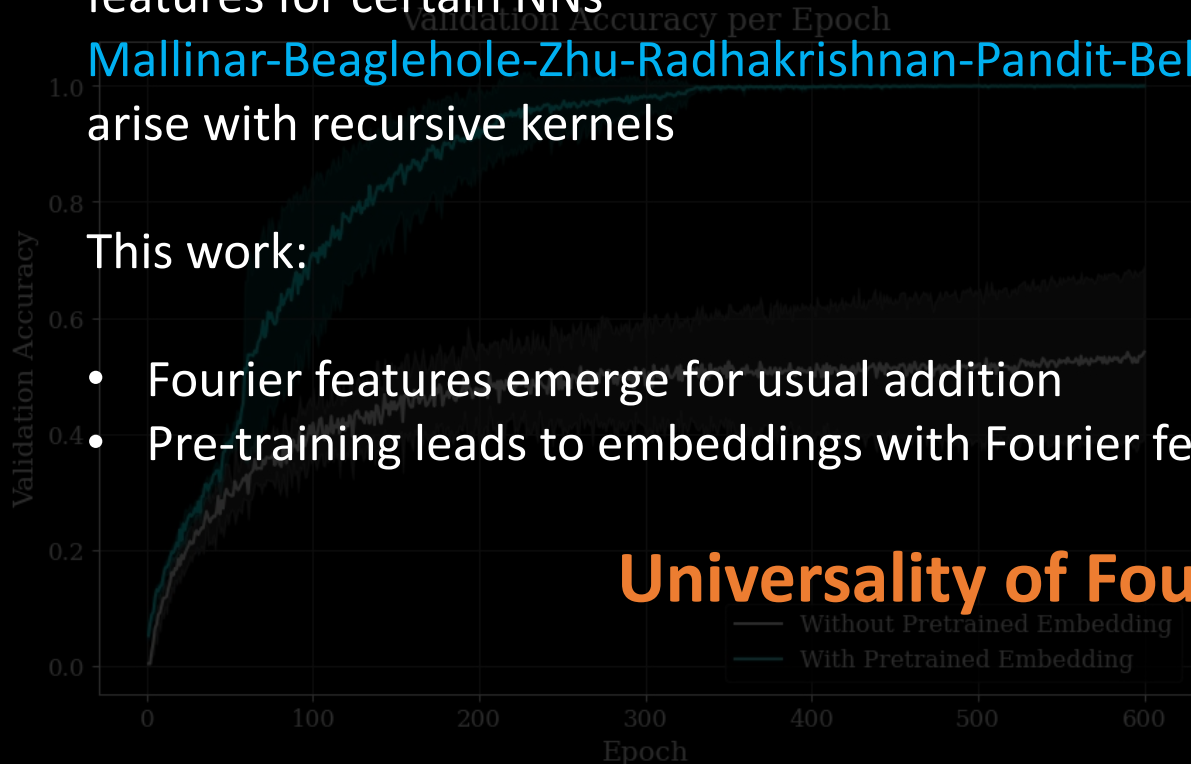
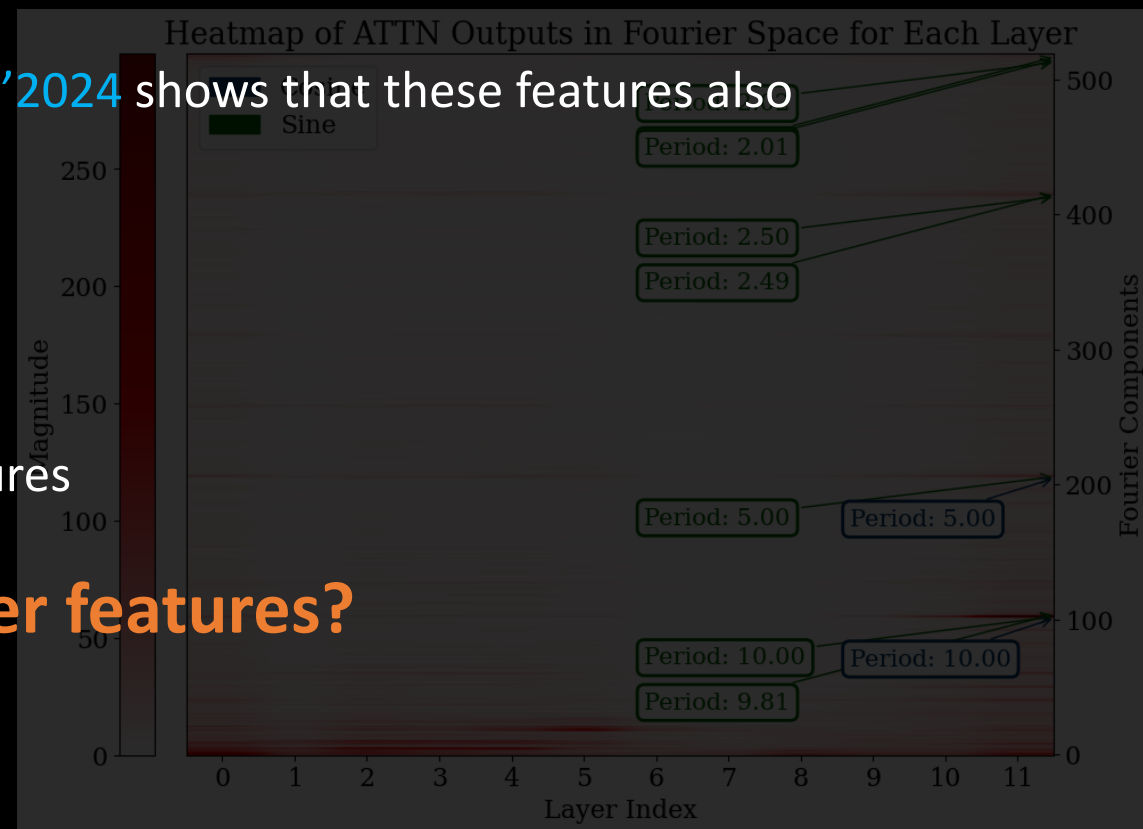Morwani-Edelman-Oncescu-Zhao-Kakade'2023 proves margin maximization leads to Fourier features for certain NNs

Mallinar-Beaglehole-Zhu-Radhakrishnan-Pandit-Belkin'2024 shows that these features also arise with recursive kernels

This work:

- Fourier features emerge for usual addition
- Pre-training leads to embeddings with Fourier features
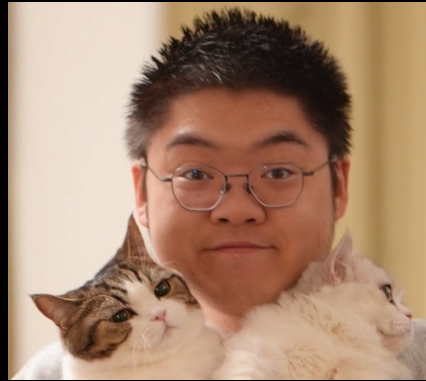
**Universality of Fourier features?**

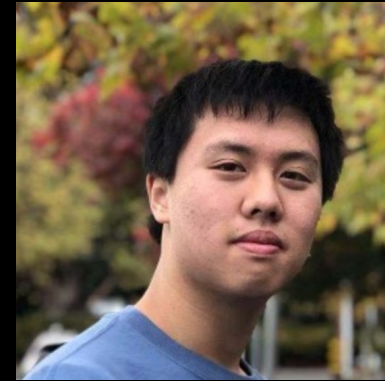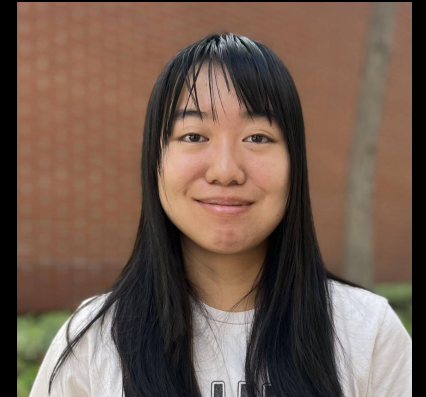# What classes of functions do Transformers prefer to learn?



Bhavya Vasudeva (USC)　　Deqing Fu (USC)　　Tianyi Zhou (USC)　　Elliot Kau (USC)　　You-Qi Huang (USC)

Simplicity Bias of Transformers to Learn Low
Sensitivity Functions, arXiv, 2024

# Sensitivity from Boolean function analysis

Consider some function $f$ defined on the Boolean hypercube $H_d$

$$Sensitivity(f) = \mathbb{E}_{x \sim H_d} \left[ \frac{1}{d} \sum_{i=1}^{d} \mathbf{1}(f(x) \neq f(x^{\oplus i})) \right]$$

Does flipping the $i$-th coordinate change the function?

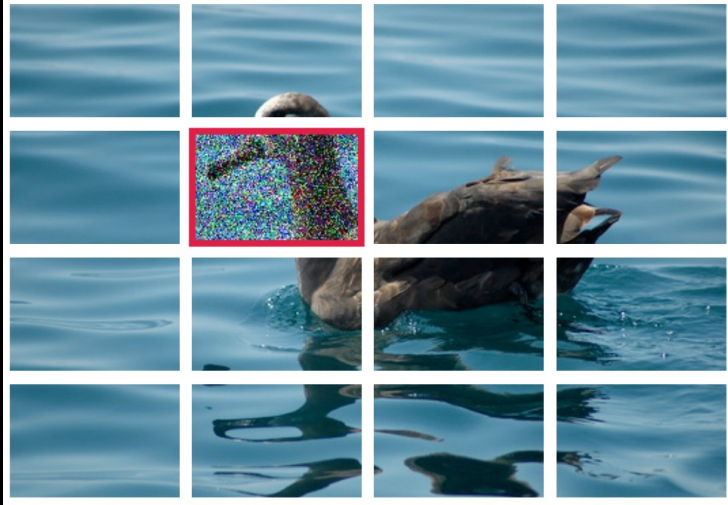Related to measures such as degree, noise stability etc.

Bhattimishra-Patel-Kanade-Blunsom'23 shows that
Transformers prefer to learn low-sensitivity Boolean functions

# Sensitivity beyond Boolean data



Evaluate model on original input

If model's predictions change, model is **sensitive** to that token

Evaluate model on perturbation to random token

# Observations: Transformers learn lower sensitivity functions

- **Image** (Fashion MNIST, CIFAR-10, SVHN, ImageNet-1k)
  - For same accuracy, Transformers learn solutions with lower sensitivity than MLPs, CNN, and also other patch-based architectures such as ConvMixer

- **Language** (Paraphrasing tasks: MRPC, QQP)
  - For same accuracy, Transformers learn solutions with lower sensitivity than LSTMs
  - LSTMs are more sensitive to recent tokens, Transformers have more uniform sensitivity across context

- **Advantages of low sensitivity**
  - Adding sensitivity as a regularizer improves robustness
  - Adding sensitivity as a regularizer also leads to flatter minima

**Sensitivity as a measure to understand inductive bias?**

# Can we use Transformers to discover data structures from scratch?

Omar Salemohamed

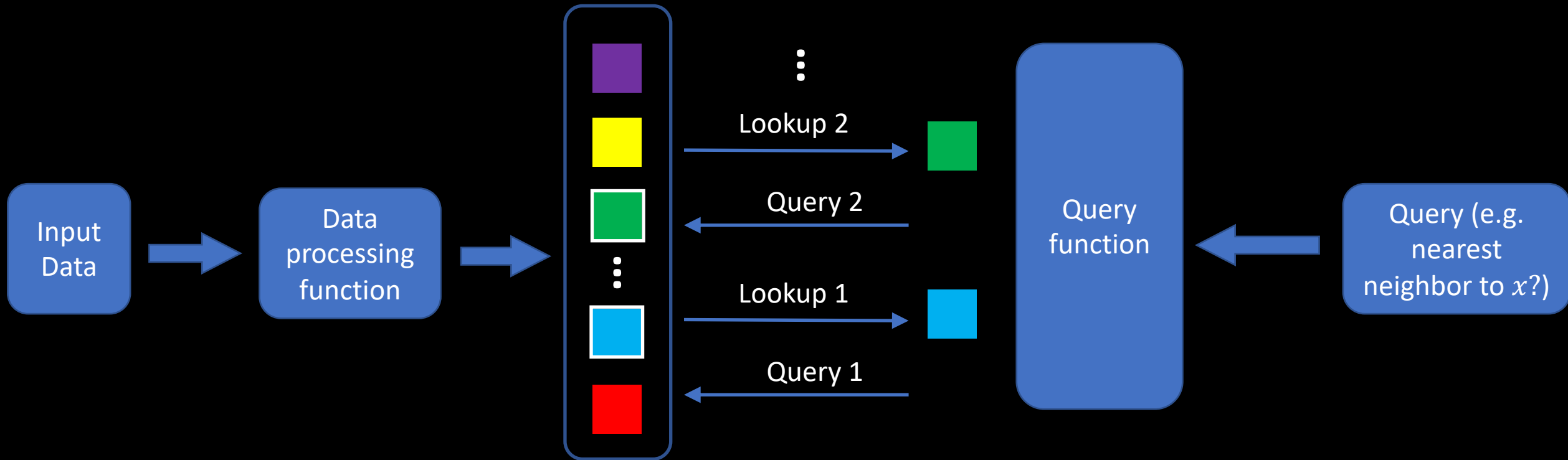(Universite de Montreal/MILA)

Laurent Charlin

(HEC Montreal/MILA)
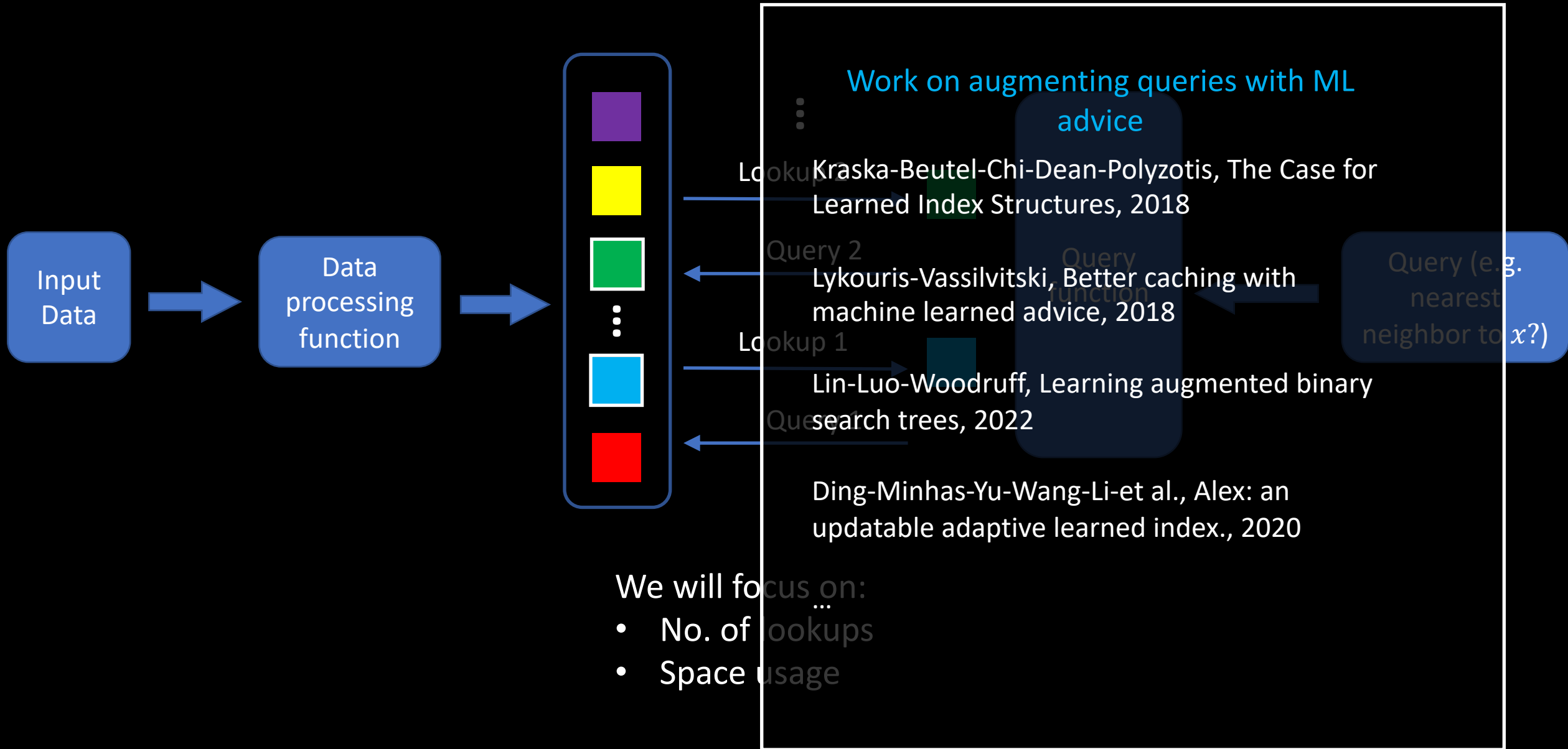
Shivam Garg

(MSR NYC)

Greg Valiant

(Stanford)

Discovering Data Structures: Nearest Neighbor Search and Beyond, ongoing

# Data structures (think nearest neighbor lookup in 1D)



We will focus on:
- No. of lookups
- Space usage

# Recent work has tried to augment data structures with ML



Input Data → Data processing function →

Lookup 2

Query 2

Lookup 1

Query 1

...

## Work on augmenting queries with ML advice

Kraska-Beutel-Chi-Dean-Polyzotis, The Case for Learned Index Structures, 2018

Lykouris-Vassilvitski, Better caching with machine learned advice, 2018

Lin-Luo-Woodruff, Learning augmented binary search trees, 2022

Ding-Minhas-Yu-Wang-Li-et al., Alex: an updatable adaptive learned index., 2020
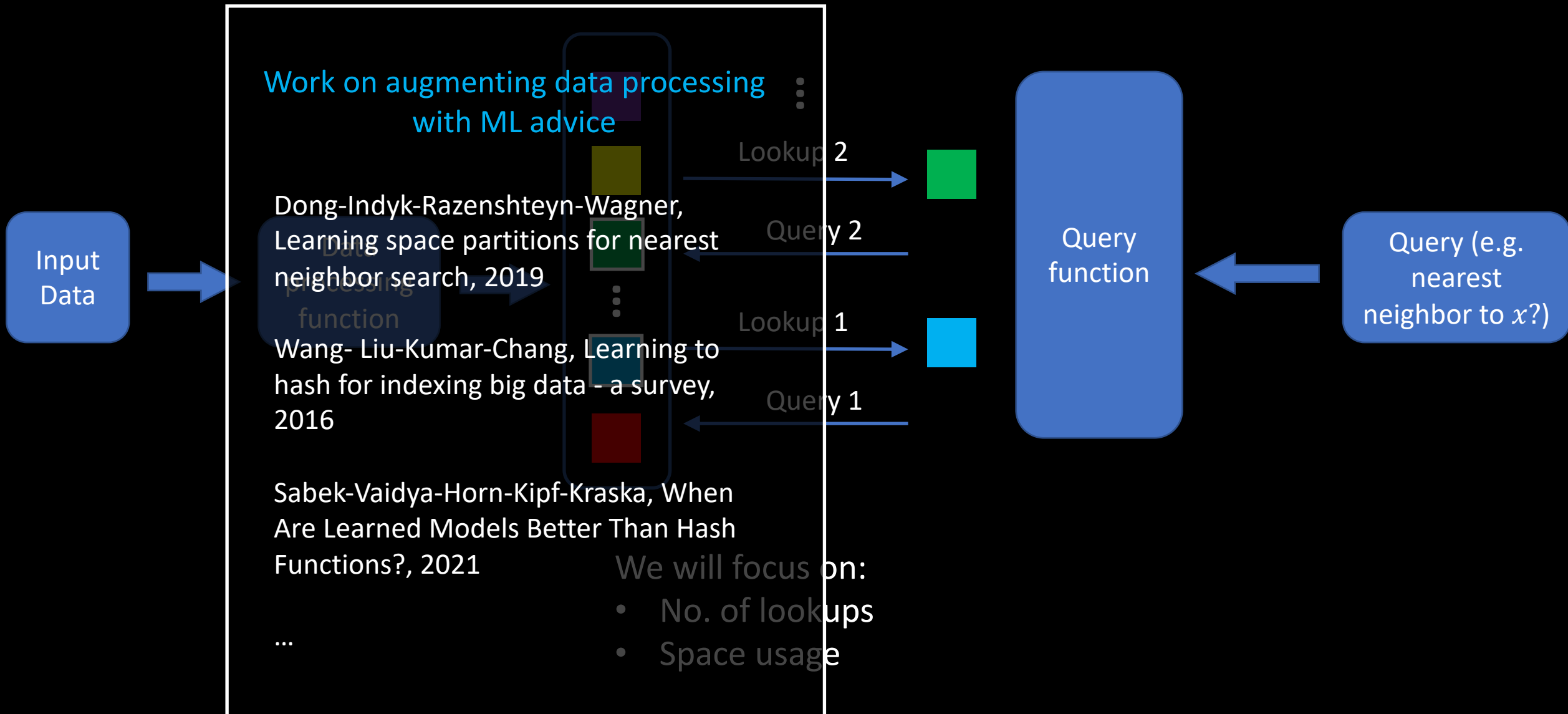
Query function

Query (e.g. nearest neighbor to $x$?)

We will focus on:
- No. of lookups
- Space usage

# Recent work has tried to augment data structures with ML

Work on augmenting data processing with ML advice

Input Data

Data processing function

Dong-Indyk-Razenshteyn-Wagner, Learning space partitions for nearest neighbor search, 2019

Wang- Liu-Kumar-Chang, Learning to hash for indexing big data - a survey, 2016

Sabek-Vaidya-Horn-Kipf-Kraska, When Are Learned Models Better Than Hash Functions?, 2021

…

Lookup 2

Query 2

Lookup 1

Query 1

We will focus on:
- No. of lookups
- Space usage

Query function

Query (e.g. nearest neighbor to $x$?)

# What if we learn everything end to end with ML, with no algorithmic priors?

**To make training efficient:**

- To enable sparse lookups: Gumbel softmax + noise
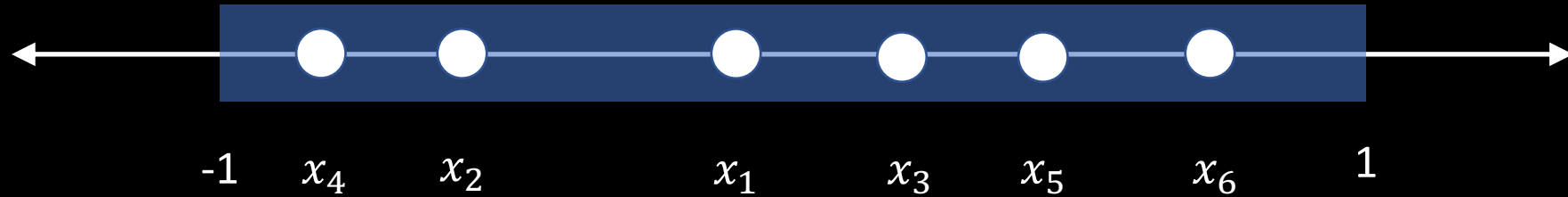- To preserve inputs: Transformer + differentiable sorting

## What data structures can we learn?

Query model ($MLP_k$) ← Query $q$

⋮

Query model ($MLP_2$) ← Query $q$

Query model ($MLP_1$) ← Query $q$

$\tilde{x}_1$    $\tilde{x}_2$    $\tilde{x}_3$   ⋯   $\tilde{x}_N$

Data-processing network (Transformer)

$x_1$    $x_2$    $x_3$   ⋯   $x_N$

**For training:**

- Sample dataset $\{x_1, \dots, x_N\}$ and query point $q$ from some distribution $D$
- If $y$ is true nearest neighbor and $\hat{y}$ is model's answer, loss is $\| y - \hat{y} \|_2^2$

# Uniform distribution in 1D
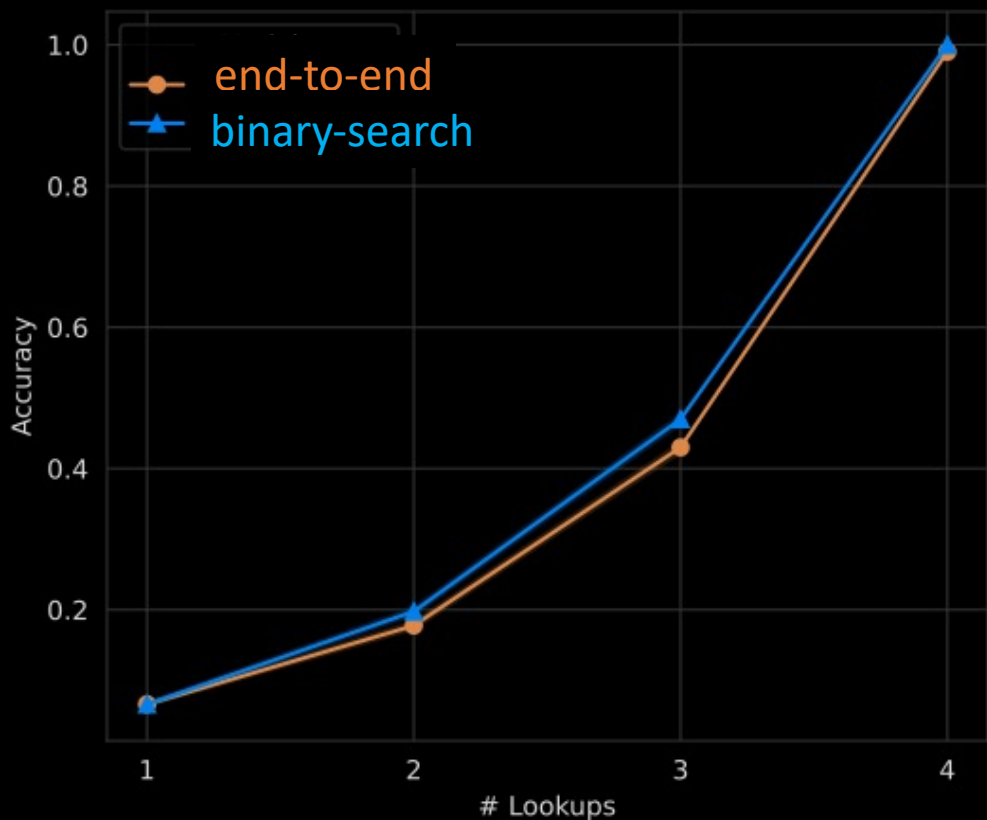


Model trained on this distribution:

- **Learns to sort, with small error**
- **Does better than binary search**
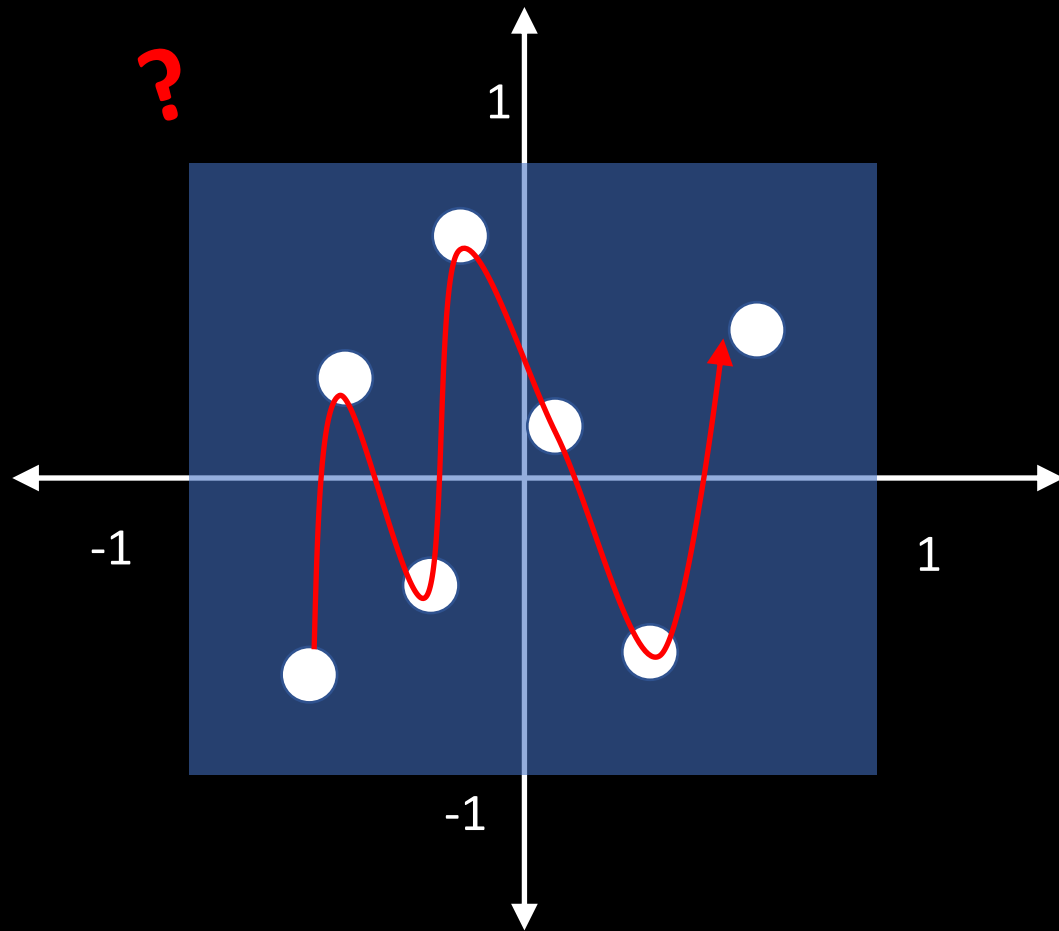
# Model outperforms binary search



Query model begins search not far from nearest neighbor
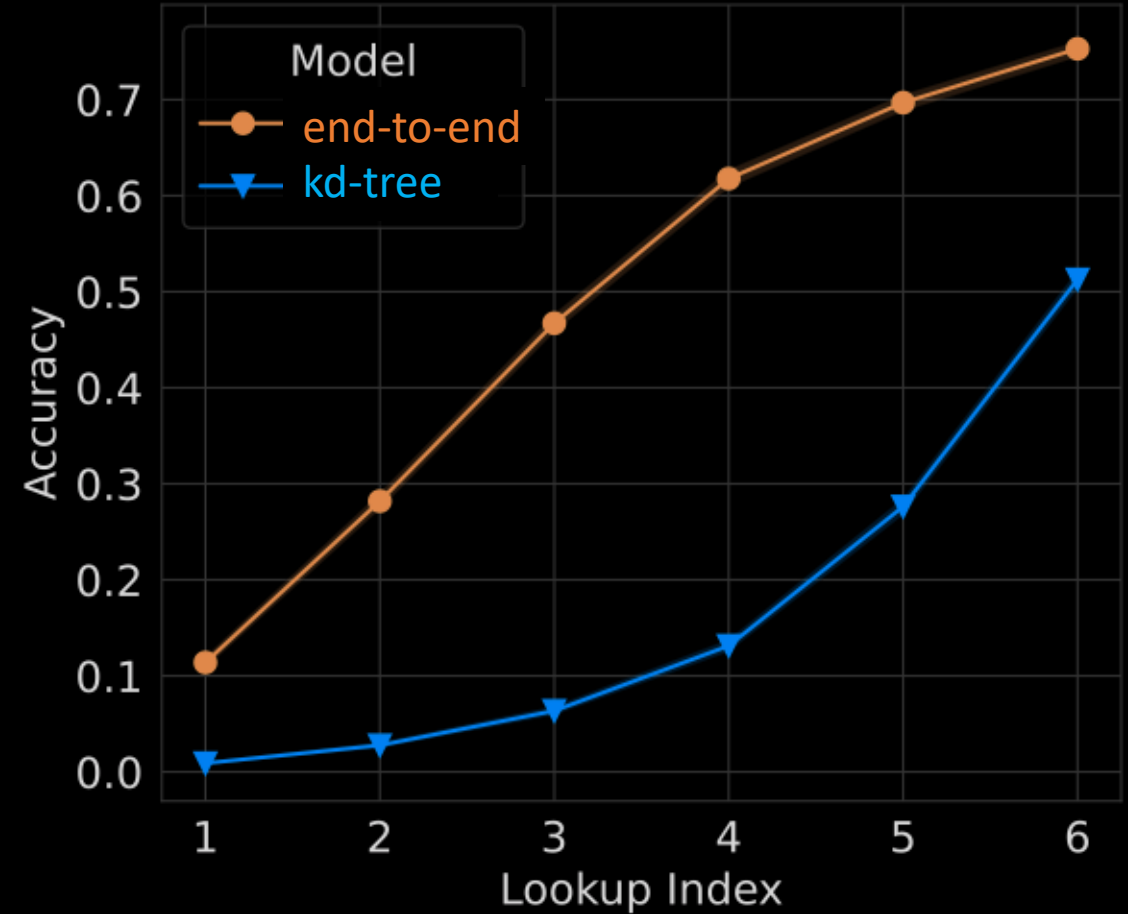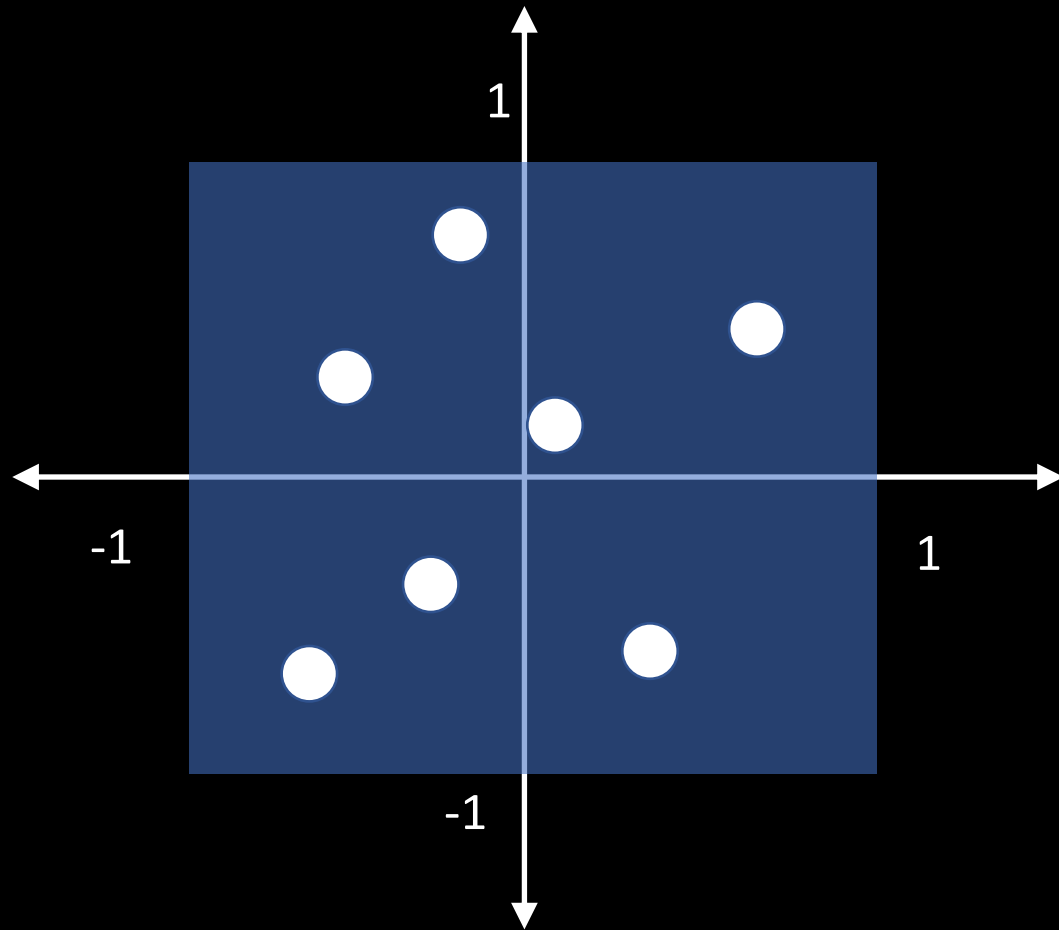
# Harder 1D distribution where quantiles don't concentrate



1D Hard Distribution Average Lookup Position

**Model learns binary search!**

# Uniform distribution in 2D: What is the right permutation?
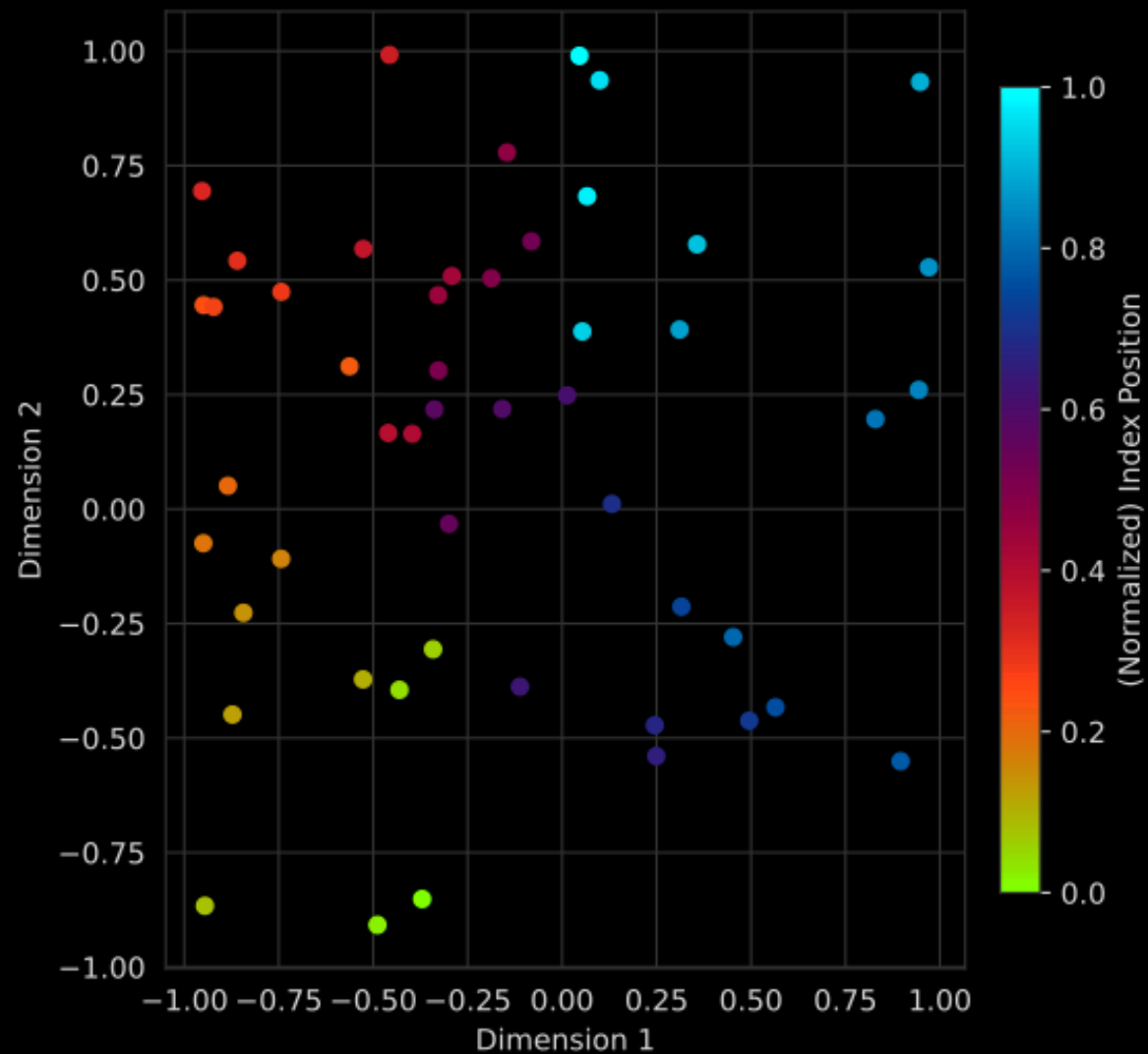
# Uniform distribution in 2D: Outperforms kd-trees
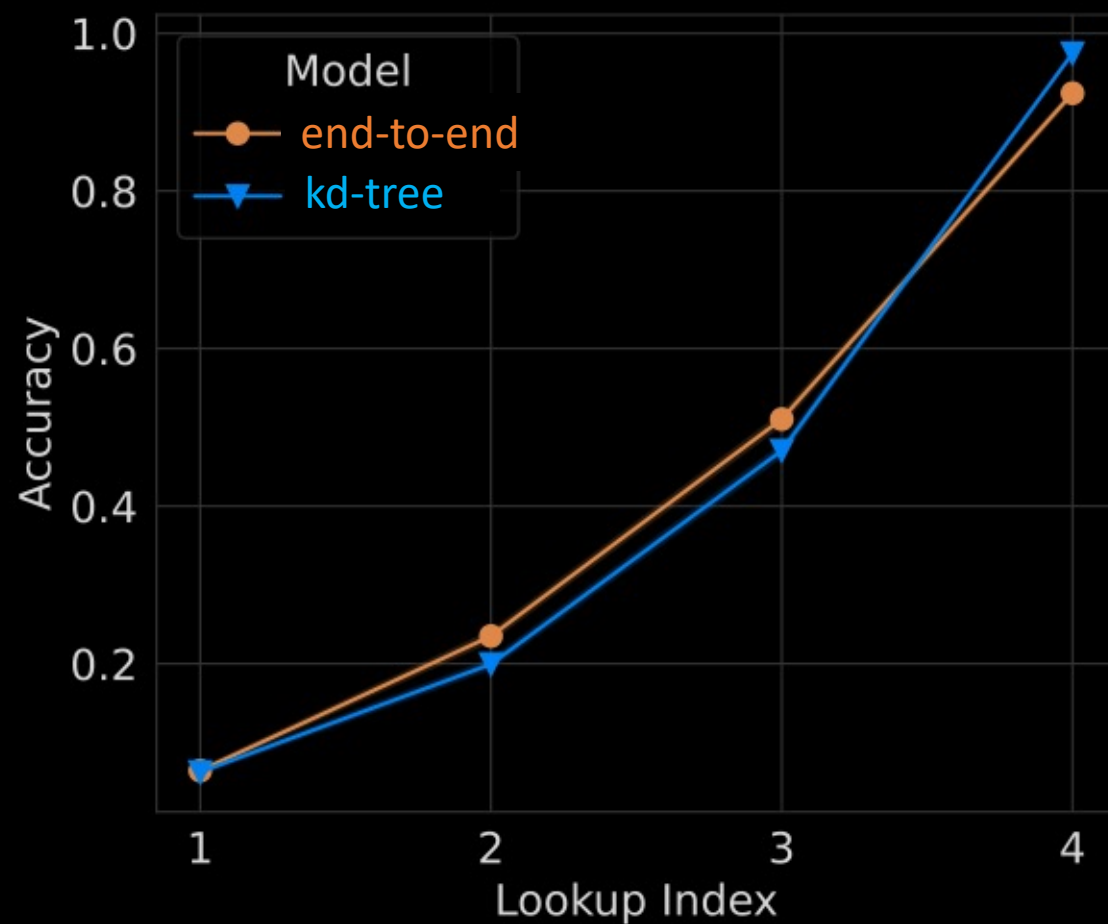
# Model learns to index nearby points together
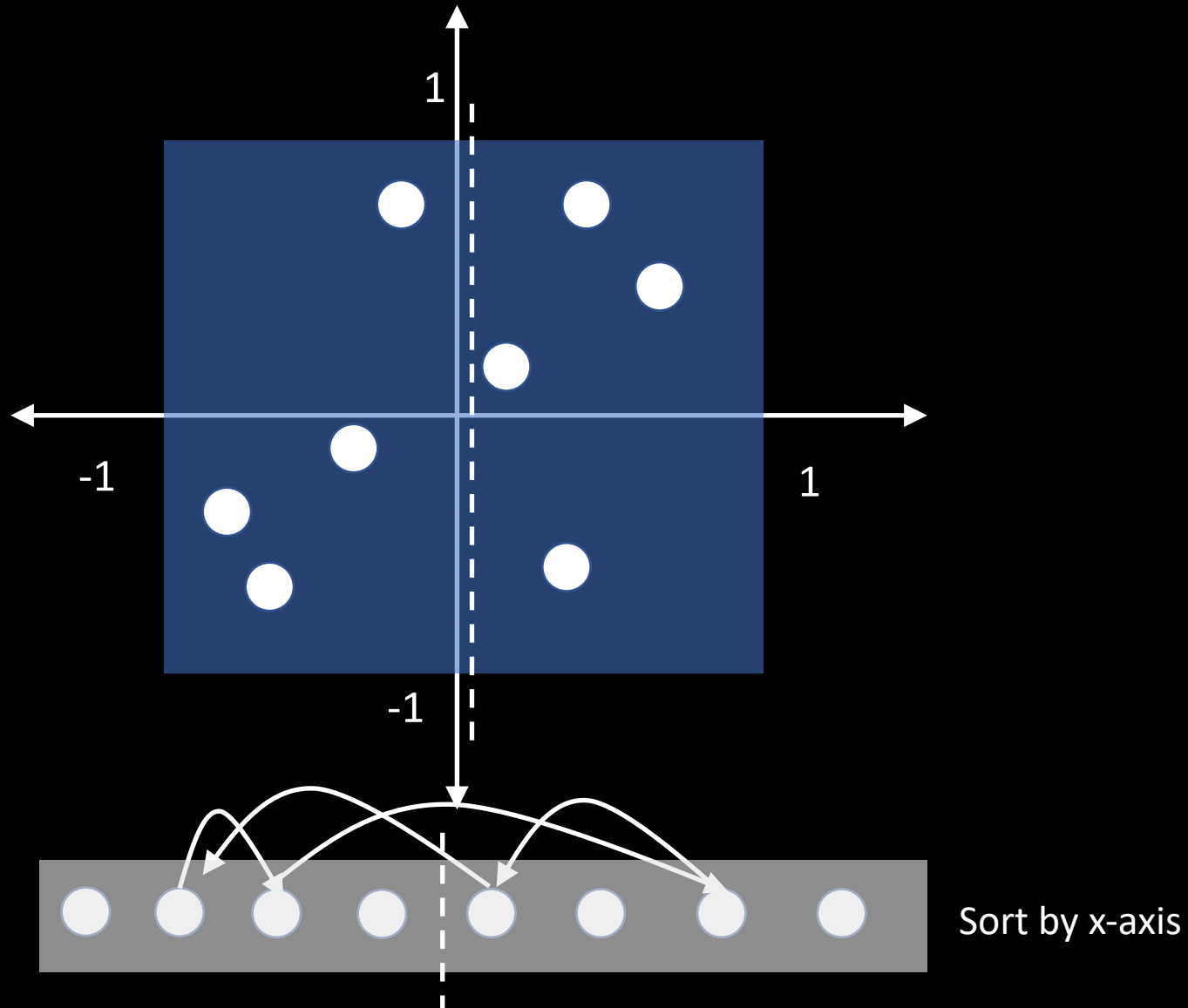


Original points colored by index position
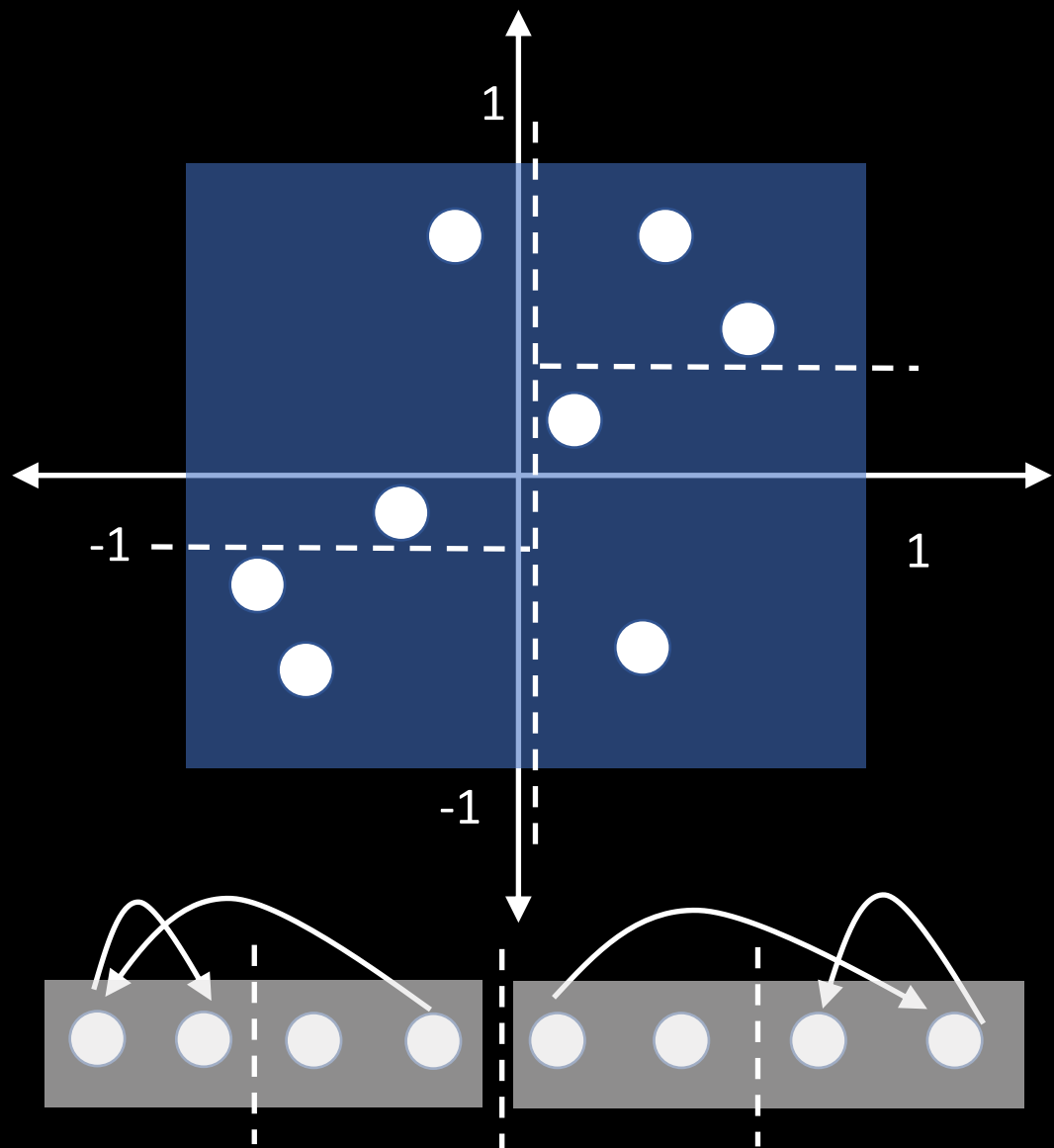
Transformed points colored by index position

# Hard distribution in 2D: Matches kd-trees

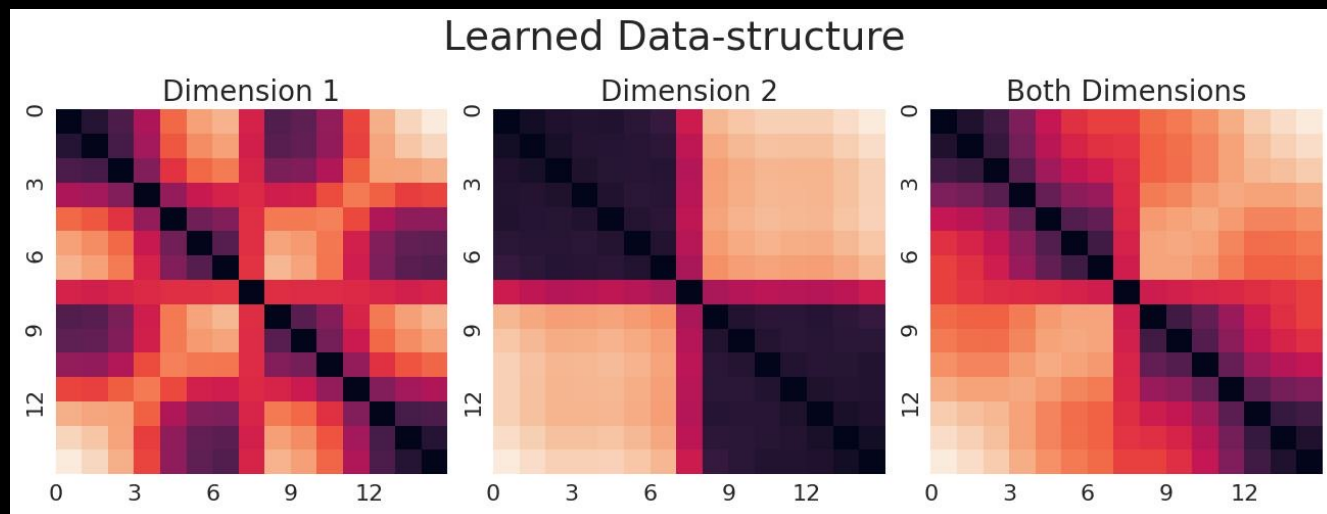# Can see that the model is essentially recovering a kd-tree!
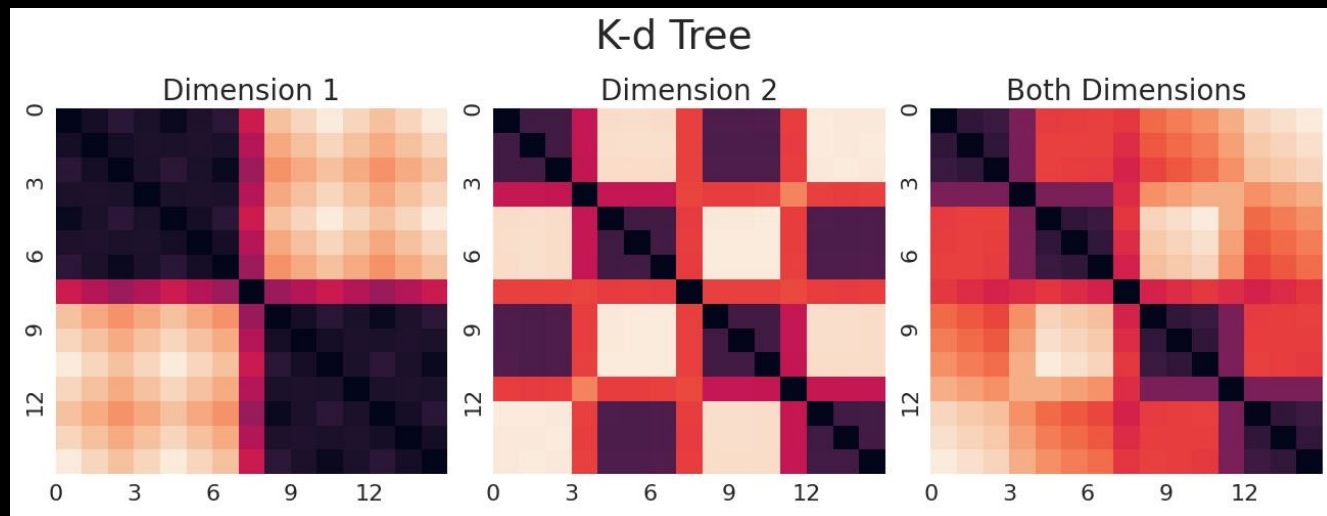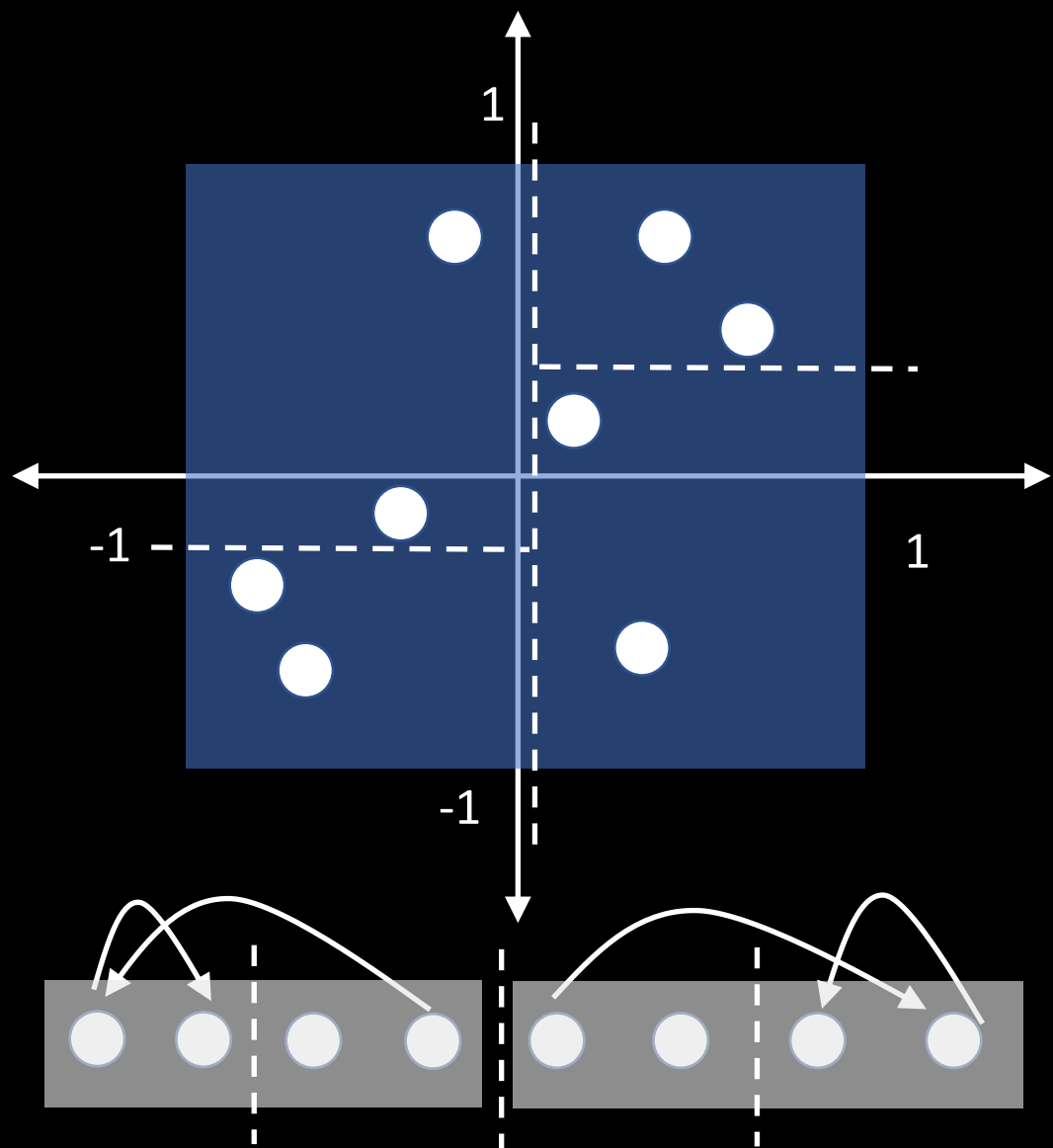


Sort by x-axis

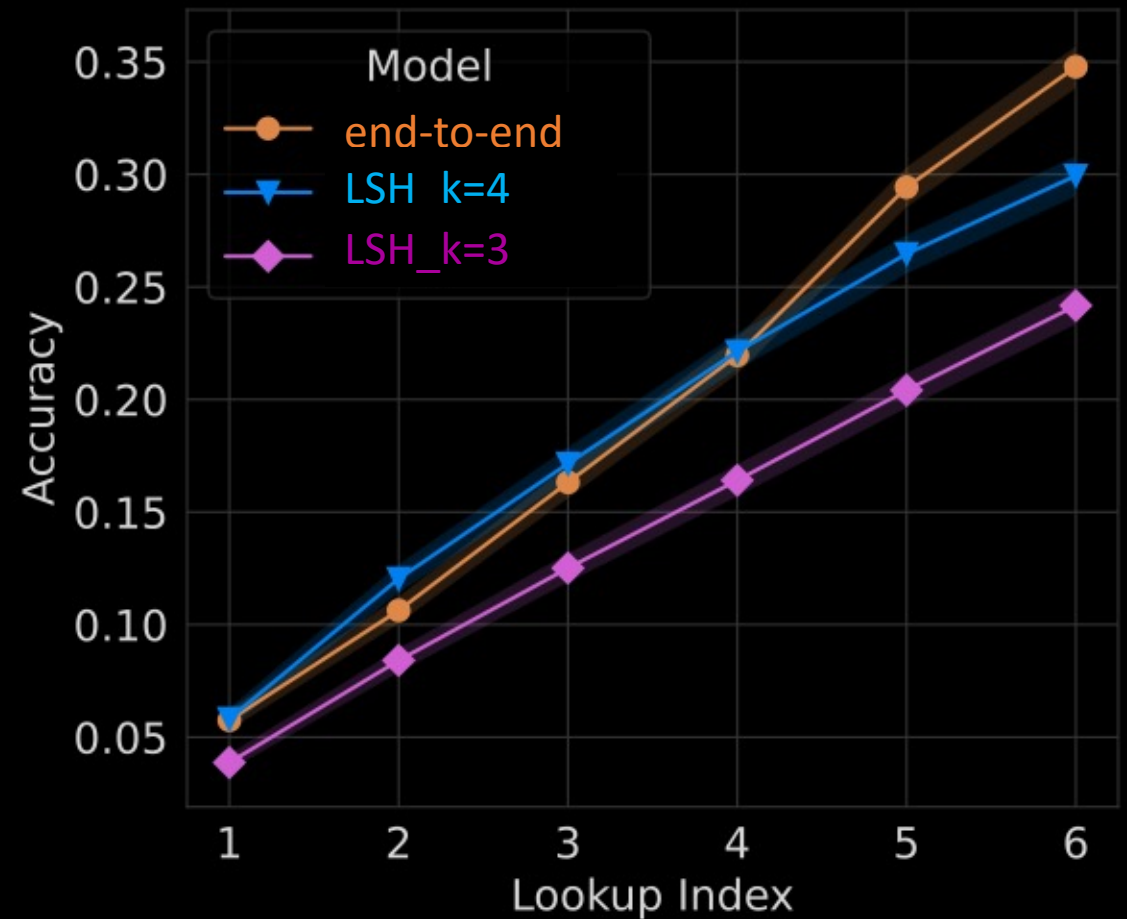# Can see that the model is essentially recovering a kd-tree!



Sort each half
by y-axis

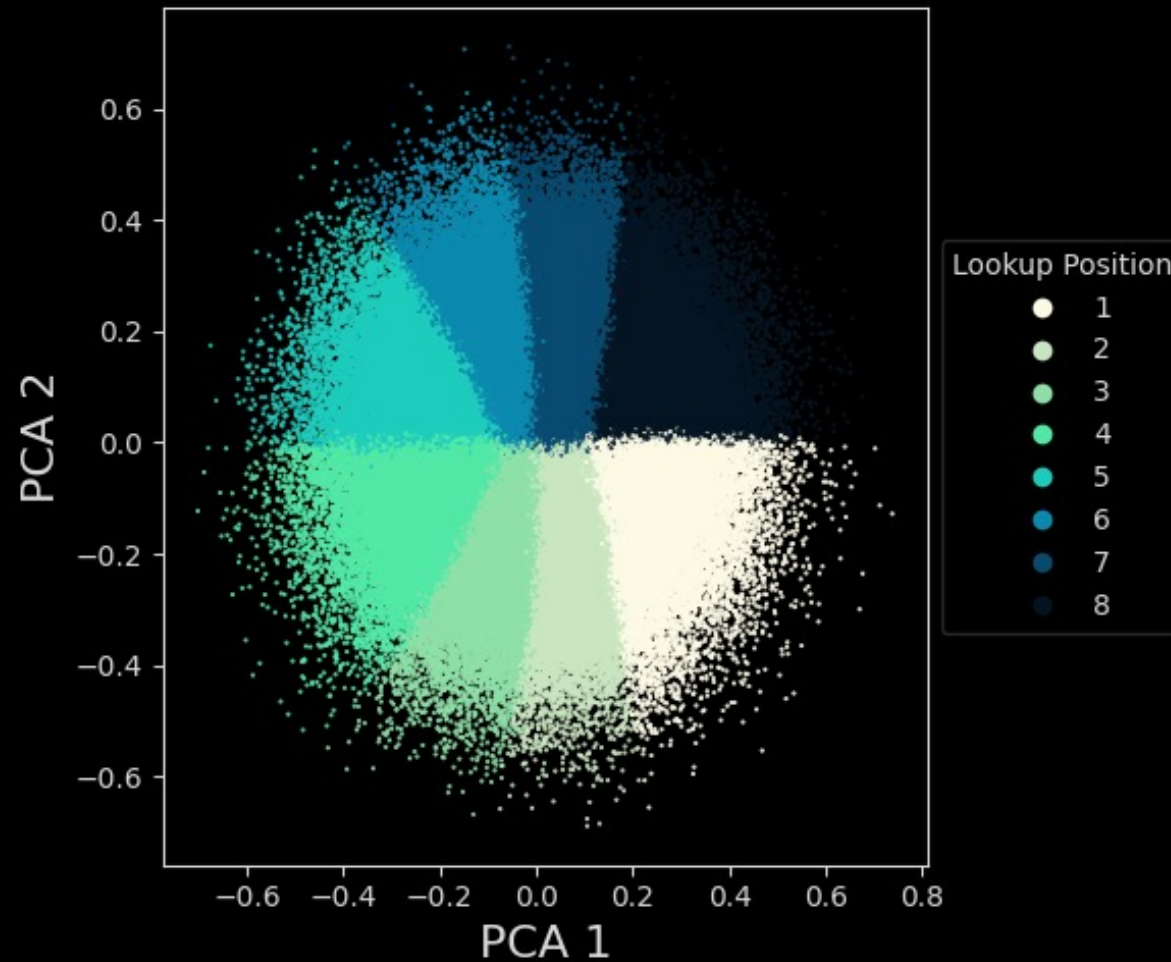# Can see that the model is essentially recovering a kd-tree!

# Uniform distribution in 30D: Matches LSH

- In high dimensions (even 30), we don't understand optimal data structures, even for the uniform distribution!

- Kd-trees suffer from curse of dimensionality

- LSH is a popular alternative

# Model learns to do a projection, like LSH



Query model mainly considers projection of query onto this 2-dimensional subspace to decide where to look

# Model can learn underlying metric space

Input: 50 images of numbers uniformly drawn from [0,200]
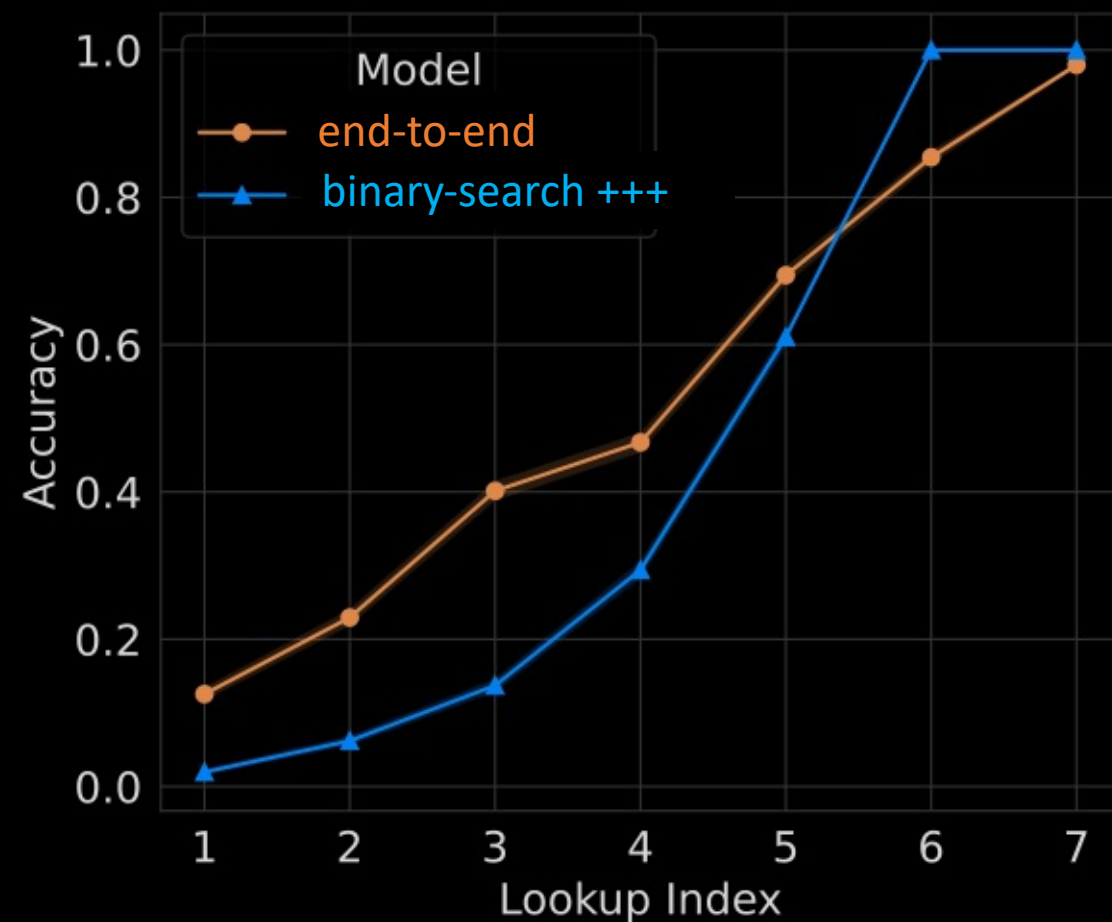


$x_1$  $x_2$  ...  $x_N$

Query: Images of numbers uniformly drawn from [0,200]



$x_q$

- Train on cross-entropy loss of prediction
- Model gets no access to the labelling of the image as a number

# Summary: Claims & Thoughts

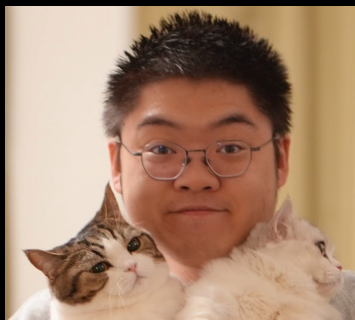**We can train models end to end to learn data structures**

- Model also learns to use extra space
- We also show we can learn data structures for frequency estimation in a data stream, recovering/outperforming count-sketch

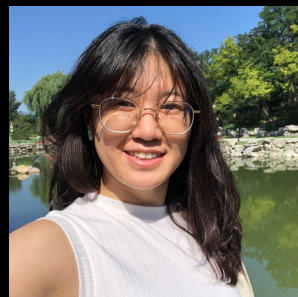**Models outperform data-independent baselines**

- Also consider settings with power-law distributions etc.

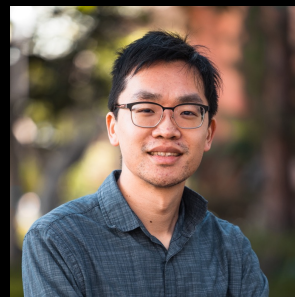**Learned models can be interpreted and understood, providing insights for data-structure design**

- *Can we use these to understand tradeoffs in theory, build better strategies for high-dimensional NN search and other data structure problems?*
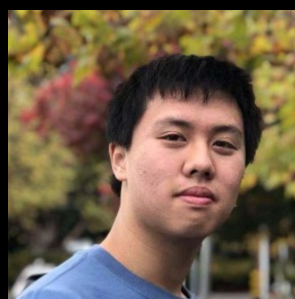
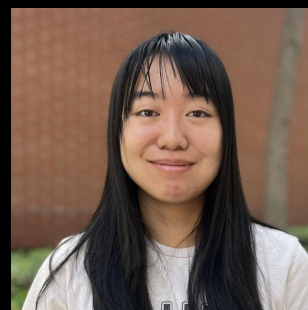Deqing Fu   Tianqi Chen   Robin Jia   Tianyi Zhou

Bhavya Vasudeva   Elliot Kau   You-Qi Huang   Omar Salemohamed

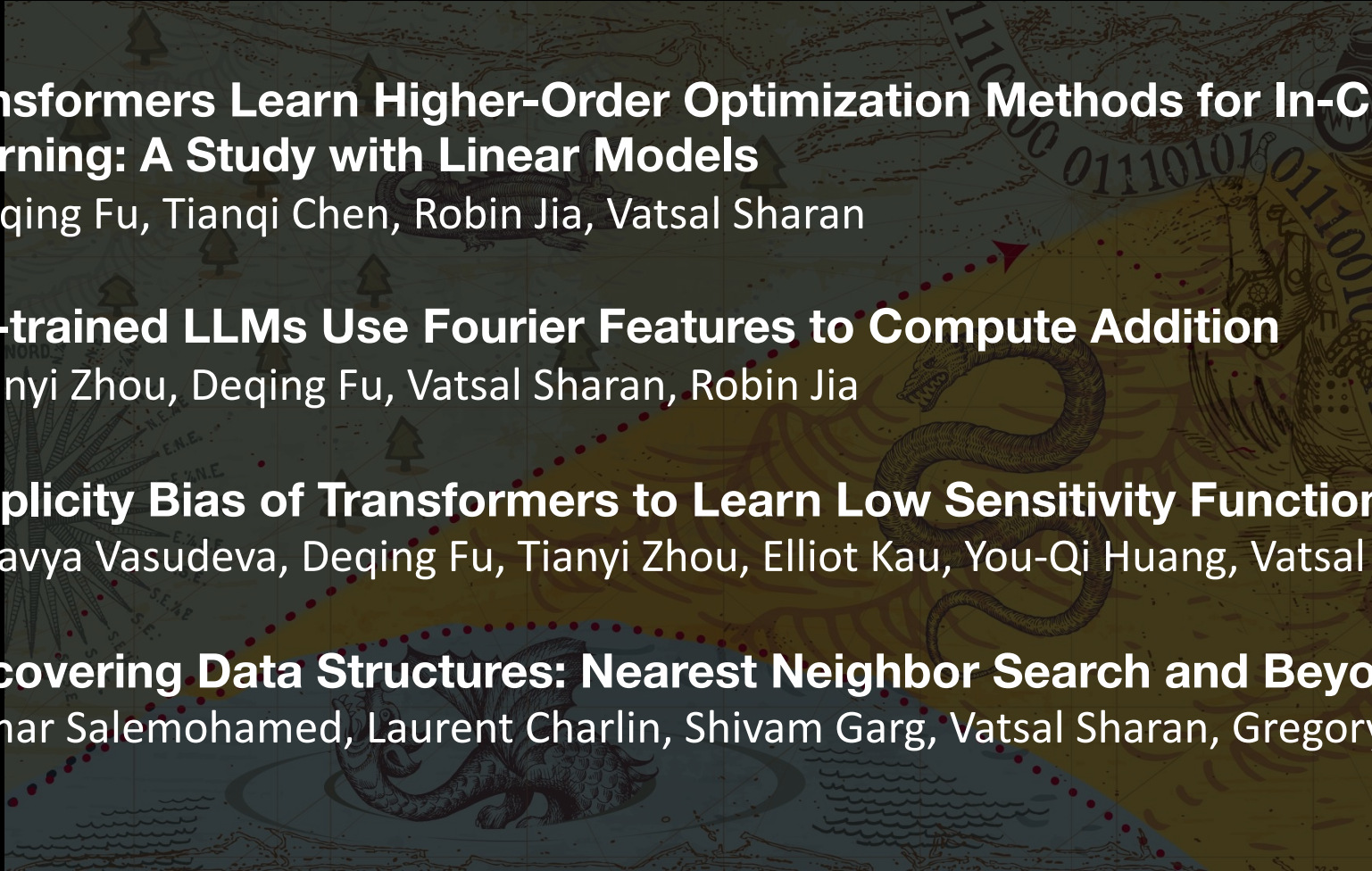Laurent Charlin   Shivam Garg   Greg Valiant

# Thanks!



- How can we use understanding of computational and information theoretic landscape to understand Transformers?
- How can we use Transformers to understand and discover algorithms and data structures?

- **Transformers Learn Higher-Order Optimization Methods for In-Context Learning: A Study with Linear Models**
  Deqing Fu, Tianqi Chen, Robin Jia, Vatsal Sharan

- **Pre-trained LLMs Use Fourier Features to Compute Addition**
  Tianyi Zhou, Deqing Fu, Vatsal Sharan, Robin Jia

- **Simplicity Bias of Transformers to Learn Low Sensitivity Functions**
  Bhavya Vasudeva, Deqing Fu, Tianyi Zhou, Elliot Kau, You-Qi Huang, Vatsal Sharan

- **Discovering Data Structures: Nearest Neighbor Search and Beyond**
  Omar Salemohamed, Laurent Charlin, Shivam Garg, Vatsal Sharan, Gregory Valiant