

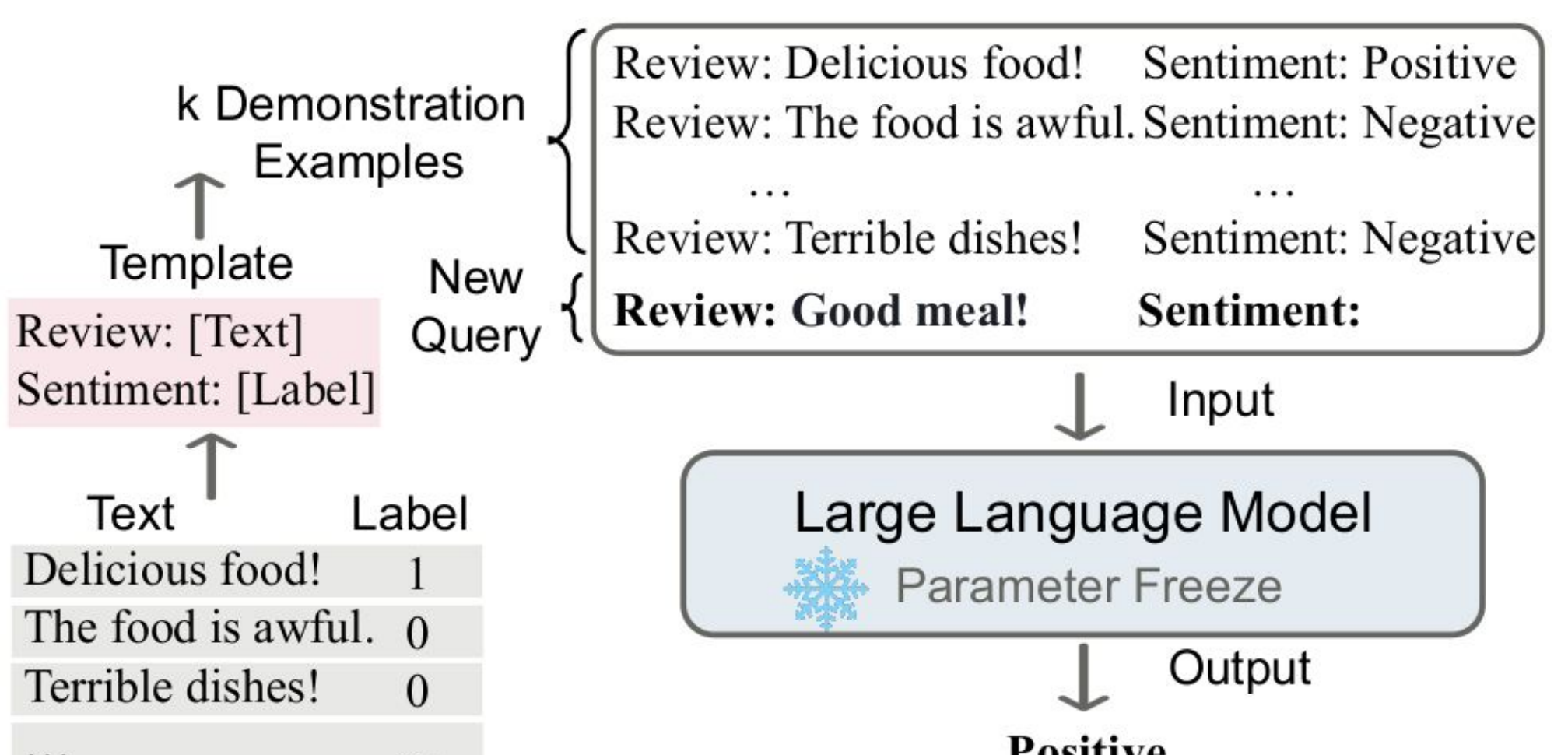
Why Larger Language Models Do In-context Learning Differently?

Zhenmei Shi, Junyi Wei, Zhuoyan Xu, Yingyu Liang

Background

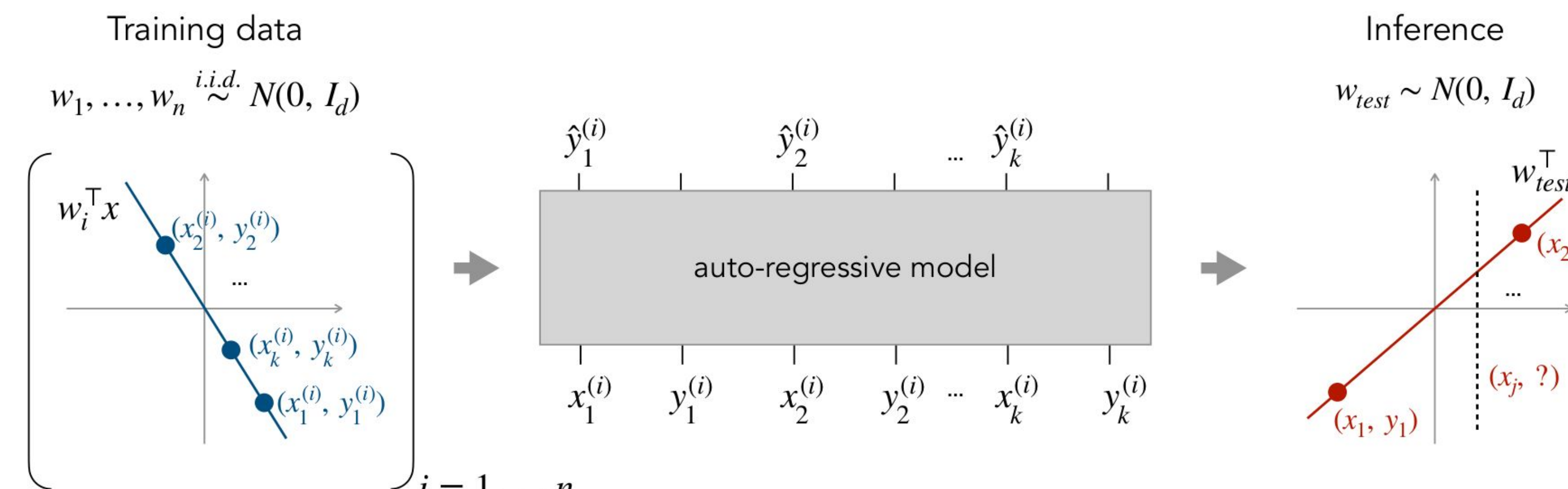
In-context Learning

Source "A Survey on In-context Learning." (arXiv 2023)



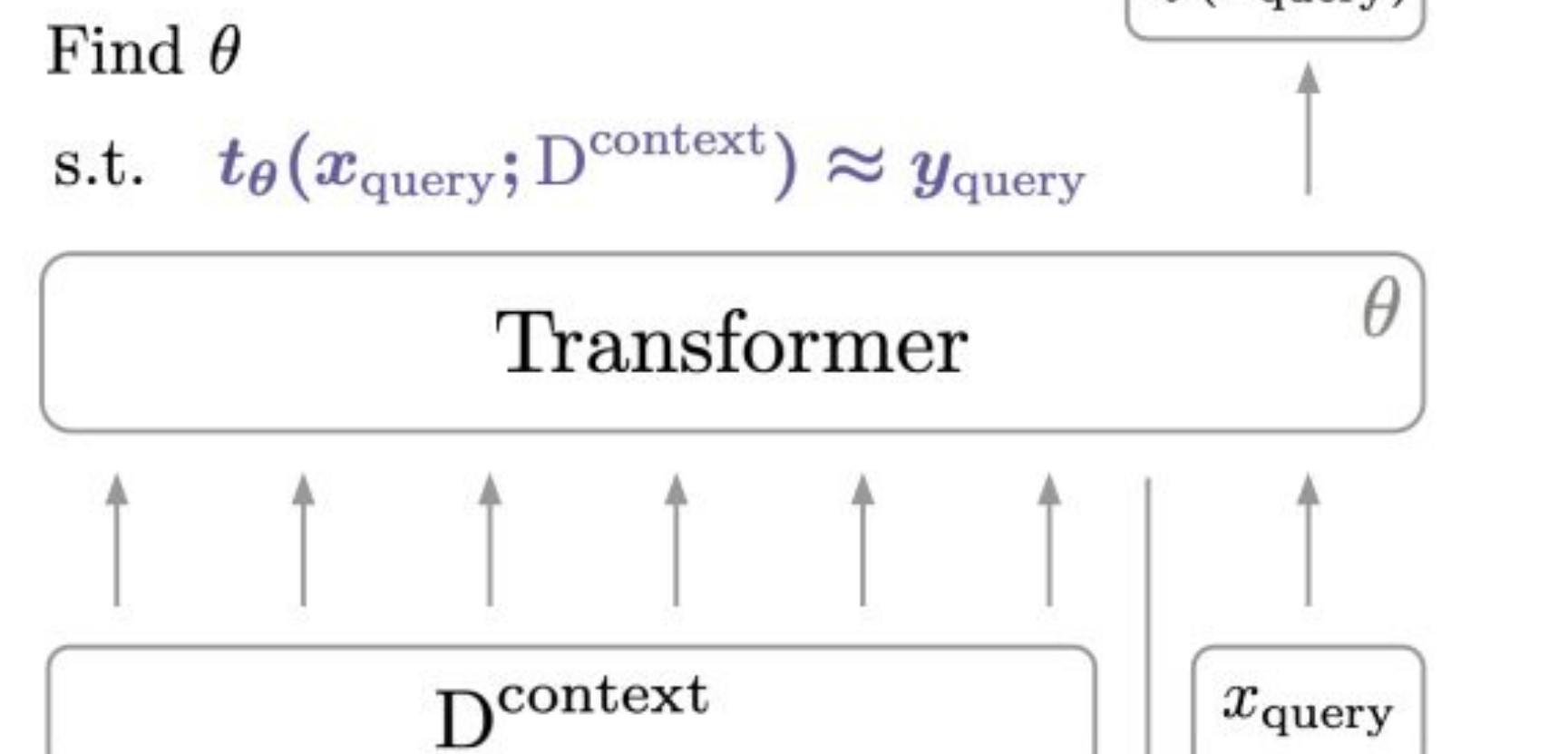
Linear Regression In-context Learning Data Model

Source "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes." (NeurIPS 2022)

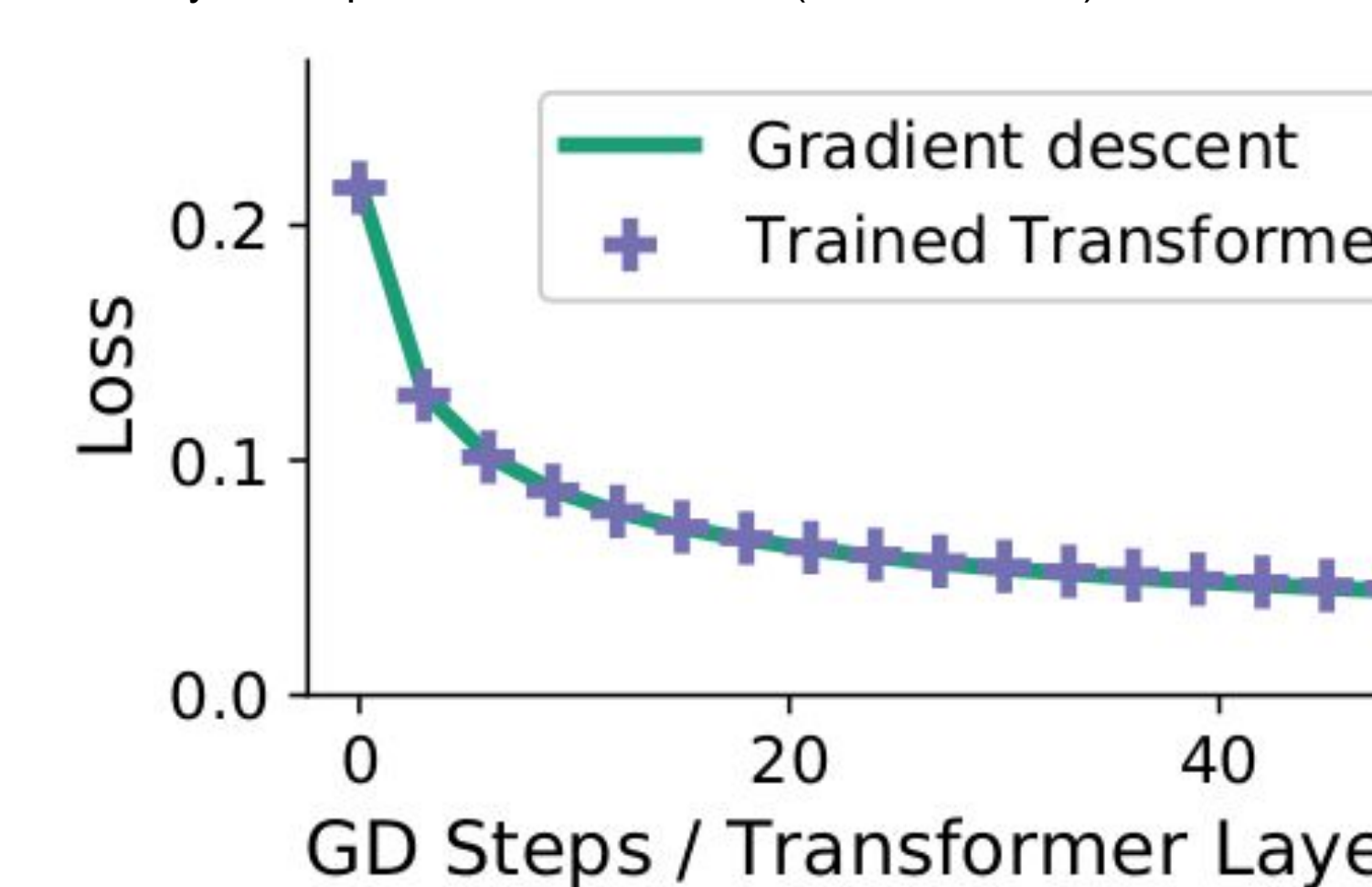


Find θ

s.t. $t_\theta(x_{\text{query}}; D^{\text{context}}) \approx y_{\text{query}}$



GD Steps / Transformer Layers



Gradient-based Optimization and Attention-based In-context Learning

Source "Transformers Learn In-Context by Gradient Descent." (ICML 2023)

Motivation

Regular ICL

Natural language targets: {Positive/Negative} sentiment

Contains no wit [...] \n Negative
Very good viewing [...] \n Positive
A smile on your face \n _____

Language Model → Positive

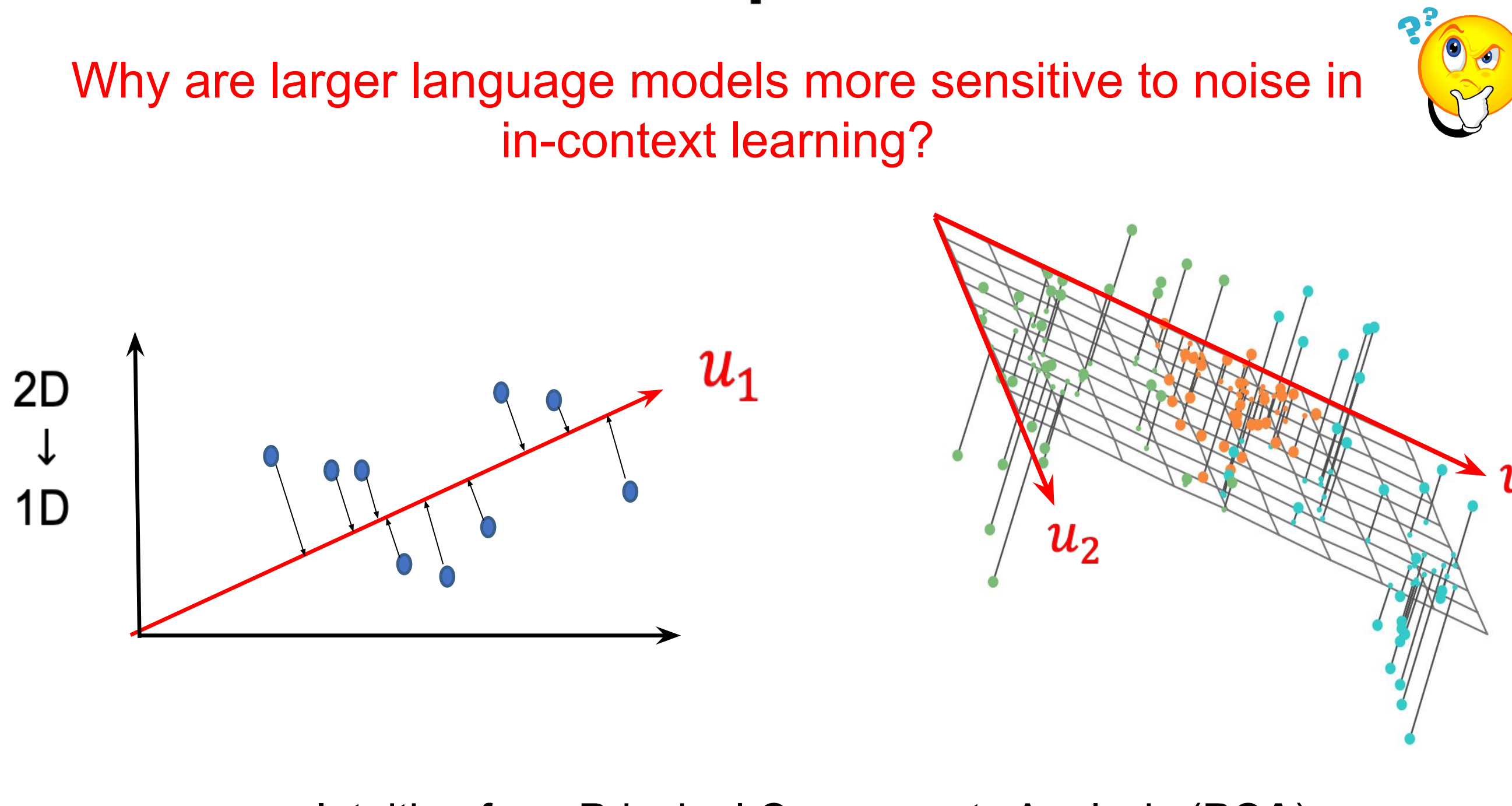
Flipped-Label ICL

Flipped natural language targets: {Negative/Positive} sentiment

Contains no wit [...] \n Positive
Very good viewing [...] \n Negative
A smile on your face \n _____

Language Model → Negative

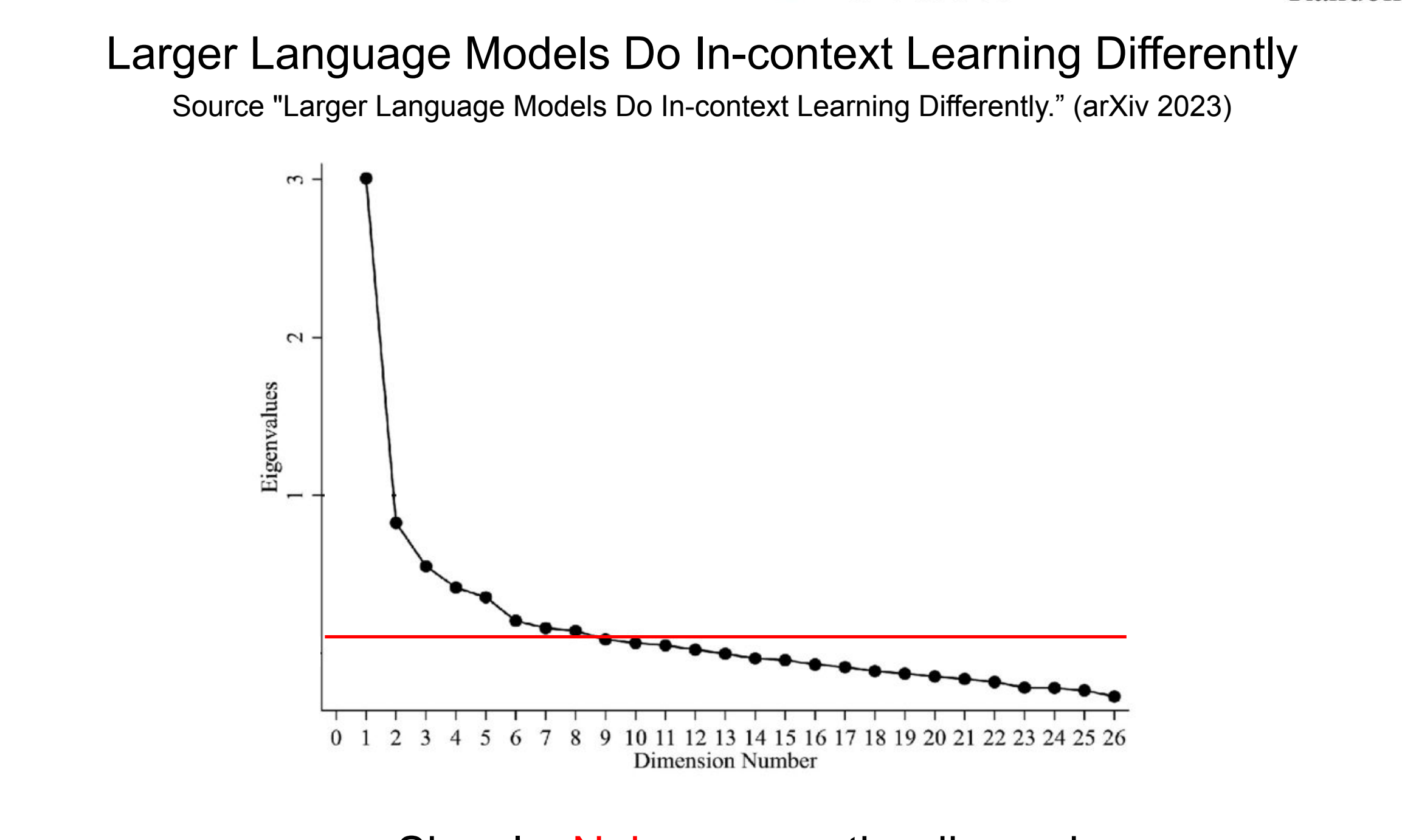
Why are larger language models more sensitive to noise in in-context learning?



Intuition from Principal Components Analysis (PCA)

Larger Language Models Do In-context Learning Differently

Source "Larger Language Models Do In-context Learning Differently." (arXiv 2023)



Signals+Noise across the dimensions

Setup: In-context Learning and Linear Self-Attention

• Test prompt:

N tokens as the context → From a new task

target label to be predicted

• Training prompts:

From task 1, From task 2, ..., From task T

Tasks from some meta distribution

• Formally, a prompt from a task τ is formed by N examples $(x_{\tau,1}, y_{\tau,1}), \dots, (x_{\tau,N}, y_{\tau,N})$ and a query input $x_{\tau,q}$ for prediction, $y_{\tau,q}$ is the target label

$$E_\tau := \begin{pmatrix} x_{\tau,1} & x_{\tau,2} & \dots & x_{\tau,N} & x_{\tau,q} \\ y_{\tau,1} & y_{\tau,2} & \dots & y_{\tau,N} & 0 \end{pmatrix} \quad y_{\tau,q}$$

Attention Mechanism in LLMs

- Compute projections (key, value, query)
- Compute attention scores using softmax on inner products
- Compute the weighted sum of values weighted by attention

$$Q = XW^Q, K = XW^K, V = XW^V$$

$$Y = \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

• We consider the Linear Self-attention Networks (only one-layer attention, and remove softmax):

$$f_{\text{LSA},\theta}(E) = \left[E + W^{PV} E \cdot \frac{E^T W^{KQ} E}{\rho} \right]$$

Model size: rank r of the weight matrix

Main Results: ICL Linear Regression

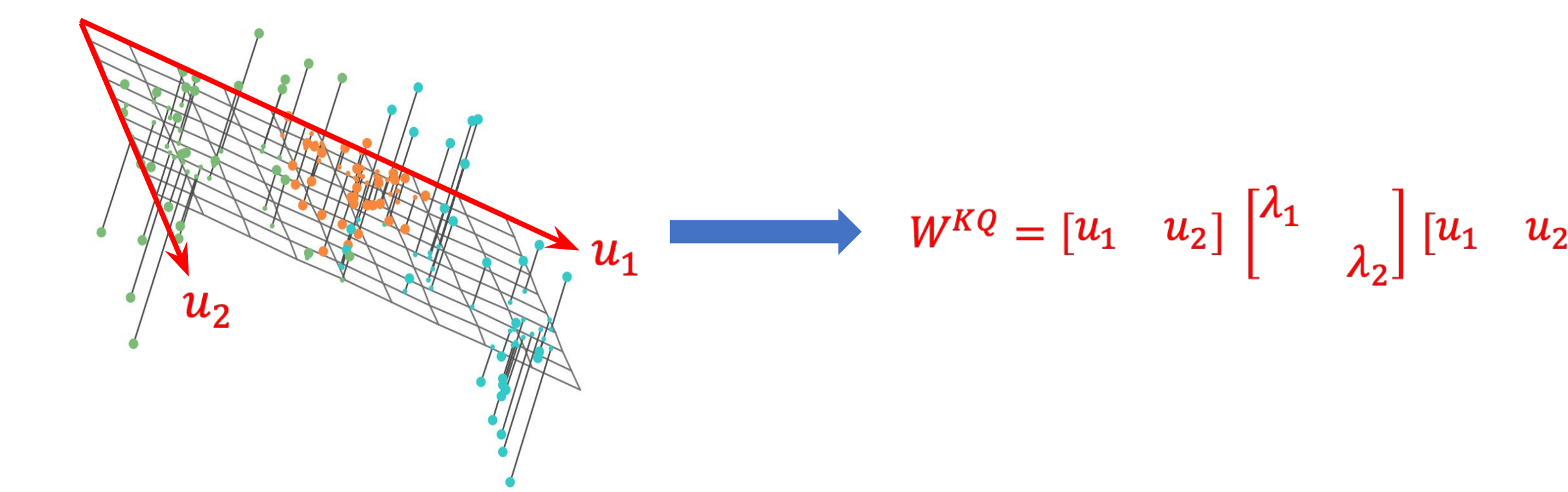
- Task distribution: $w_\tau \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_{d \times d})$
- $x_{\tau,i}, x_{\tau,q} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Lambda)$
- $y_{\tau,i} = \langle w_\tau, x_{\tau,i} \rangle + \epsilon_i$ where $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$

Theorem (Optimal rank- r solution, informal)

The optimal rank- r weight matrix W^{KQ} lies in S_r , the span of the top- r eigenvectors of the data covariance. It is the truncated version (low rank approximation) of the optimal full-rank solution.

Theorem (Behavior difference, special case and informal)

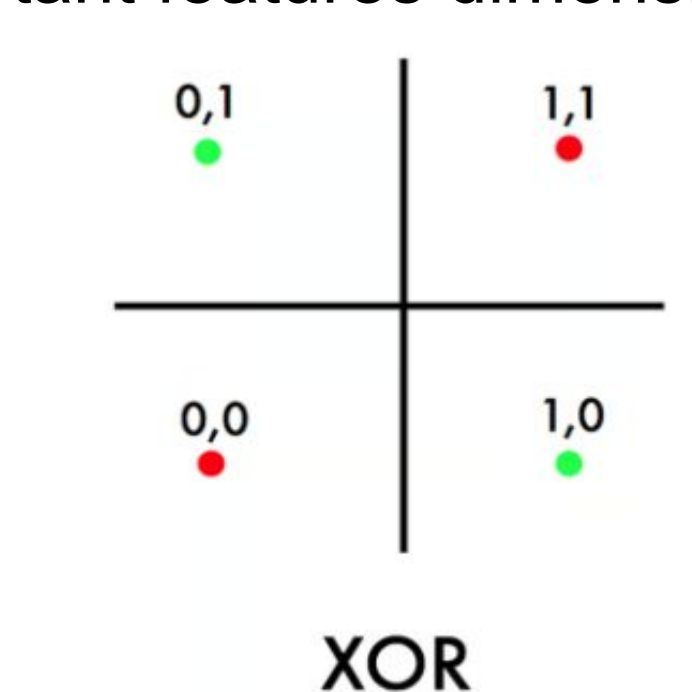
For $r_1 < r_2$, let f_1 be the optimal rank- r_1 model, and f_2 be the optimal rank- r_2 model. Assume the true regression weight lies in S_{r_1} . Let s be the projection of the true weight in S_{r_1} .

$$\mathcal{L}(f_2) - \mathcal{L}(f_1) \approx \frac{r_2 - r_1}{M} \|s\|_D^2 + \frac{r_2 - r_1}{M} \sigma^2$$


Main Results: Sparse Parity Classification

Setup

- Consider random pick 2 dimension as XOR input
- Important features dimensions and less-important features dimensions
- Linear self-attention with ReLU MLP
- Consider Hinge Loss
- Number of heads to measure mode size

$$g(X, y, x_q) = \sum_{i \in [m]} a_i \sigma \left[\frac{y^T X}{N} W^{(i)} x_q \right]$$


Theorem (Optimal solution for parity)

The optimal solution of smaller number of heads model will mainly encode the **important features**, while larger number of heads model will encode **all features**.

Theorem (Behavior difference for parity)

During evaluation, we can decompose the input into two parts: signal and noise. Both the larger model and smaller model can capture the signal part well. However, the smaller model has a much smaller influence from noise than the larger model.

Take Home Messages

- Larger models are more sensitive to noise injected in-context
- We provide theoretical insights by analyzing two stylized settings
- Main insight: larger models cover more feature directions, but also covers more injected noise, thus are more sensitive to noise

Further thoughts:

- How to "regularize" the model to be more robust in-context?
- The same insight can explain the scaling law of LLMs?