

# Implicit Geometry of Next-token Prediction: From Language Sparsity to Model Representations

Christos Thrampoulidis (UBC)

November 5, 2024

CIMI Workshop, Toulouse

# Disclaimer

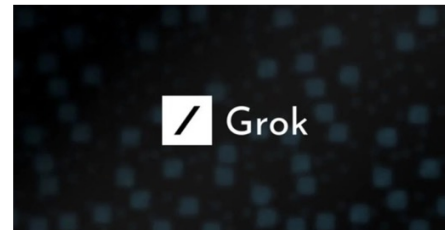
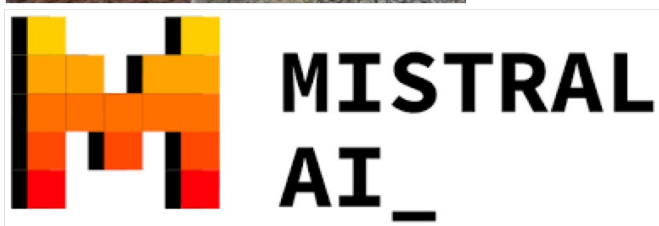
---

- Today's talk is not explicitly about "*statistical*" 😞
- But, it might be implicitly, and is "*beyond classical regimes*" 😊

# New Sheriff in town

---

- ❑ DL success “started” with **image classification** task
- ❑ Today’s “hot” topic: **Language modeling**
- ❑ LLMs: revolution in natural-language processing and generation



# New Sheriff in town

❑ DL success “started” with image classification task

❑ Today’s “hot” topic: Language modeling

- LLMs: revolution in natural-language processing and generation

## Key ingredients:

### 1. Architecture: **Transformer**

- Parallelizable + trainable to huge scale ( $\sim \Theta(B)$  parameters)
- Self-Attn: leverage long-range context info

Mamba: Linear-Time Sequence Modeling with Selective State Spaces

Albert Gu<sup>1</sup> and Tri Dao<sup>2</sup>

<sup>1</sup>Machine Learning Department, Carnegie Mellon University

<sup>2</sup>Department of Computer Science, Princeton University

agu@cs.cmu.edu, tri@tridao.me

Abstract

Foundation models, now powering most of the exciting applications in deep learning, are almost universally based on the Transformer architecture and its core attention module. Many subquadratic-time architectures such as linear attention, gated convolution and recurrent models, and structured state space models (SSMs) have been developed to address Transformers’ computational inefficiency on long sequences, but they have not performed as well as attention on important modalities such as language. We identify that a key weakness of such models is their inability to perform content-based reasoning, and make several improvements. First, simply letting the SSM parameters be functions of the input addresses their weakness with discrete modalities, allowing the model to selectively propagate or forget information along the sequence length dimension depending on the current token. Second, even though this change prevents the use of efficient convolutions, we design a hardware-aware parallel algorithm in recurrent mode. We integrate these selective SSMs into a simplified end-to-end neural network architecture without attention or even MLP blocks (Mamba). Mamba enjoys fast inference (5x higher throughput than Transformers) and linear scaling in sequence length, and its performance improves on real data up to million-length sequences. As a general sequence model backbone, Mamba achieves state-of-the-art performance across several modalities such as language, audio, and genomics. On language modeling, our Mamba-3B model outperforms Transformers of the same size and matches Transformers twice its size, both in pretraining and downstream evaluation.

xLSTM: Extended Long Short-Term Memory

Maximilian Beck<sup>1,2</sup> Korbinian Pöppel<sup>1,2</sup> Markus Spanring<sup>1</sup>  
Andreas Auer<sup>1,2</sup> Oksandra Prudnikova<sup>1</sup> Michael Kopp  
Günter Klambauer<sup>1,2</sup> Johannes Brandstetter<sup>1,2,3</sup> Sepp Hochreiter<sup>1,2,3</sup>

<sup>1</sup>Equal contribution

<sup>1</sup>ELLIS Unit, LIT AI Lab, Institute for Machine Learning, JKU Linz, Austria  
<sup>2</sup>NXAI Lab, Linz, Austria. <sup>3</sup>NXAI GmbH, Linz, Austria

Abstract

In the 1990s, the constant error carousel and gating were introduced as the central ideas of the Long Short-Term Memory (LSTM). Since then, LSTMs have stood the test of time and contributed to numerous deep learning success stories, in particular they constituted the first Large Language Models (LLMs). However, the advent of the Transformer technology with parallelizable self-attention at its core marked the dawn of a new era, outpacing LSTMs at scale. We now raise a simple question: How far do we get in language modeling when scaling LSTMs to billions of parameters, leveraging the latest techniques from modern LLMs, but mitigating known limitations of LSTMs? Firstly, we introduce exponential gating with appropriate normalization and stabilization techniques. Secondly, we modify the LSTM memory structure, obtaining: (i) xLSTM with a scalar memory, a scalar update, and new memory mixing, (ii) mLSTM that is fully parallelizable with a matrix memory and a covariance update rule. Integrating these LSTM extensions into residual block backbones yields xLSTM blocks that are then residually stacked into xLSTM architectures. Exponential gating and modified memory structures boost xLSTM capabilities to perform favorably when compared to state-of-the-art Transformers and State Space Models, both in performance and scaling.

1.04517v1 [cs.LG] 7 May 2024

# New Sheriff in town

---

❑ DL success “started” with **image classification** task

❑ Today’s “hot” topic: **Language modeling**

- LLMs: revolution in natural-language processing and generation

## Key ingredients:

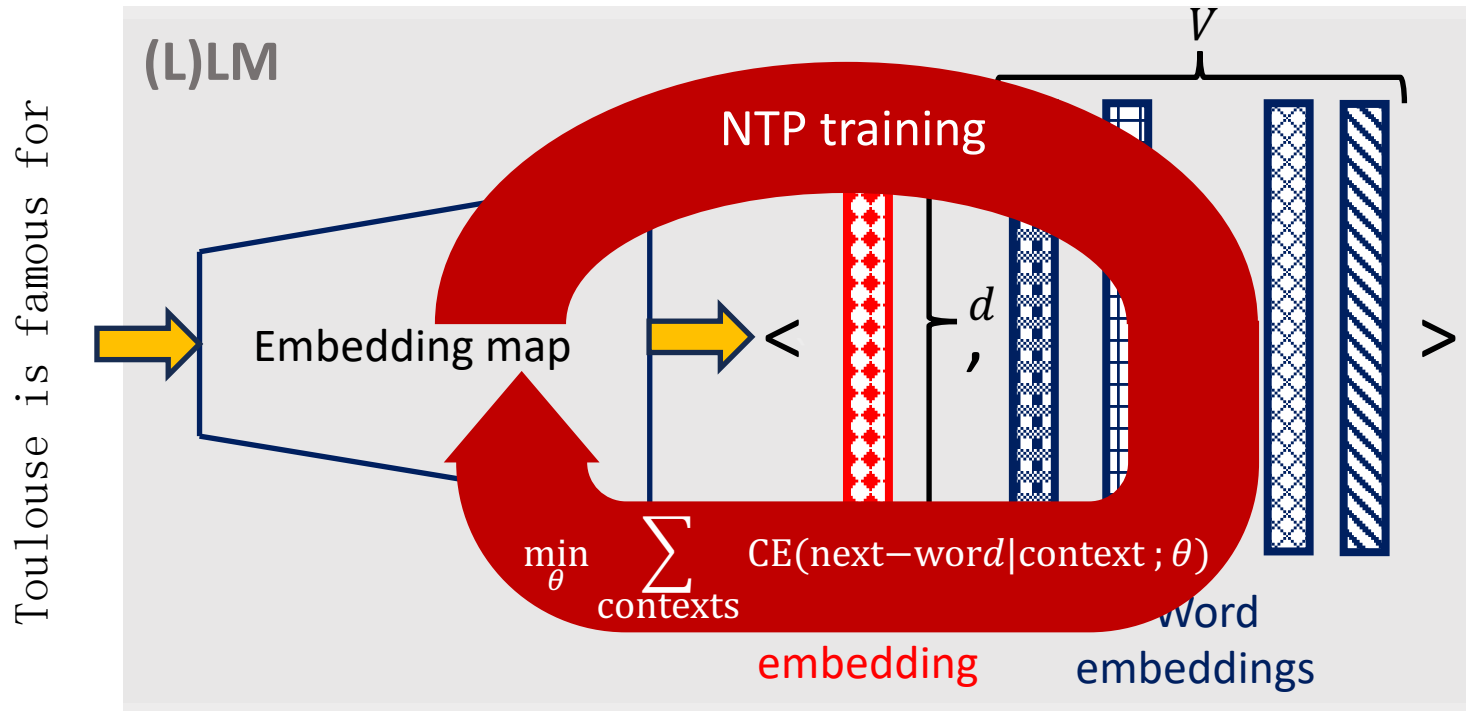
1. Architecture: **Transformer**

- Parallelizable + trainable to huge scale ( $\sim\Theta(B)$  parameters)
- Self-Attn: leverage long-range context info

2. Training: Autoregressive **next-token prediction (NTP)**

- (Pre)Train to sequentially predict **next-token** in a sequence
- Unsupervised method with supervised flavor

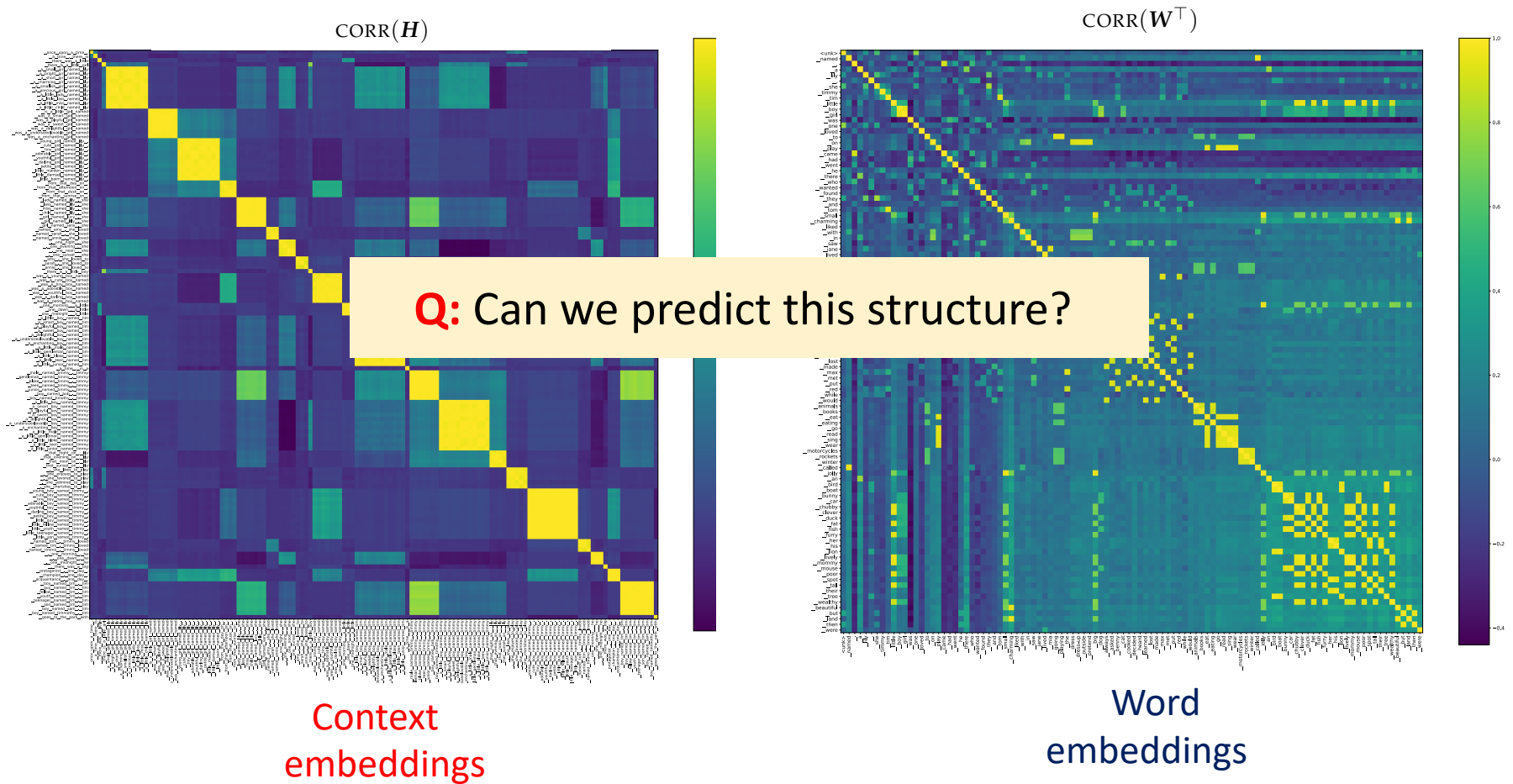
# Focus: NTP



**Q:** How do the learnt context/word representations encode the statistics of the data they are trained on?

# Implicit Geometry

**Q:** How distances and angles of the model representations correlate with linguistic patterns at the end of training?



# Why?

---

- ❑ Interpretability: **transparency** on the inner workings of LLMs
- ❑ Identify/mitigate **sources of errors/biases**
- ❑ **Algorithm improvements** upon vanilla NTP paradigm and optimizers
- ❑ Enhance our **grasp of language itself**



# Challenges & key message

- ❑ Representations are outputs of training **complicated models** (architecture, size) over **complex datasets** (source, size, tokenization) with **varying choices of optimization hyperparameters** (learning rate, weight decay, number of iterations)

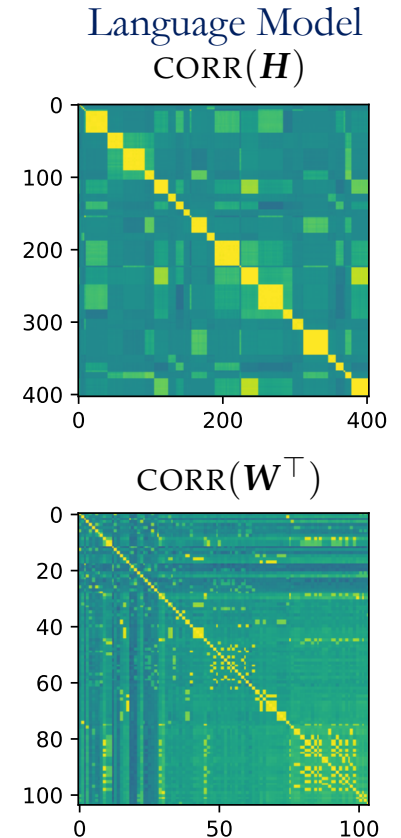
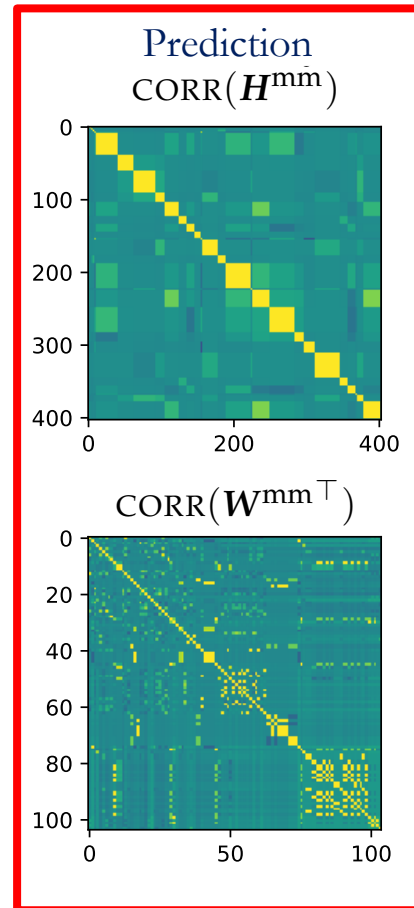
- ❑ Message:

In a doubly-asymptotic regime of  
**Large model + Long training**

context/word embeddings are  
*matrix factorization of a logit matrix*

$$L_{\text{sparse}} + \infty \cdot L_{\text{low-rank}}$$

with components determined solely  
by **patterns of text training data**



# Key Ingredients

---

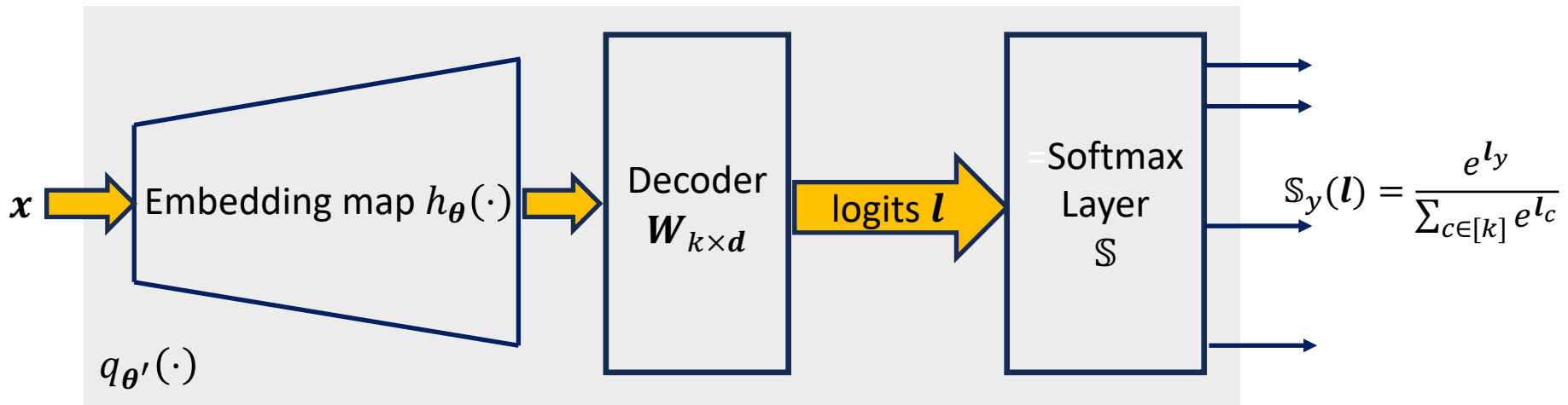
1. Correctly **framing** the **next token prediction** training task
2. Leveraging the technical framework of **implicit optimization bias**
3. Assuming large model with **unconstrained features**

# Traditional Setting

## □ One-hot multiclass classification

- **Training Data:**  $\mathcal{T}_n \triangleq (\mathbf{x}_i, y_i)_{i \in [n]}$ ,  $y_i \in [k] \triangleq \{1, 2, \dots, k\}$
- **Training Loss:**  $\min_{\theta'} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(y_i, q_{\theta'}(\mathbf{x}_i))$

$$\min_{\theta'=(\mathbf{W}, \theta)} \left\{ \text{CE}(\mathbf{W}, \theta) \triangleq \frac{1}{n} \sum_{i=1}^n -\log \left( S_{y_i}(\mathbf{W} h_{\theta}(\mathbf{x}_i)) \right) \right\}$$



# Traditional setting

---

## □ One-hot multiclass classification

- **Training Data:**  $\mathcal{T}_n \triangleq (\mathbf{x}_i, y_i)_{i \in [n]}$ ,  $y_i \in [k] \triangleq \{1, 2, \dots, k\}$
- **Training Loss:**  $\min_{\theta'} \frac{1}{n} \sum_{i \in [n]} \mathcal{L}(y_i, q_{\theta'}(\mathbf{x}_i))$

$$\min_{\theta' = (\mathbf{W}, \theta)} \left\{ \text{CE}(\mathbf{W}, \theta) \triangleq \frac{1}{n} \sum_{i=1}^n -\log \left( \mathbb{S}_{y_i}(\mathbf{W} h_{\theta}(\mathbf{x}_i)) \right) \right\}$$

□ Trained with first-order **gradient-based optimizers**, e.g. (S)GD

□ **Overparameterization**  $\implies$  Interpolation (Separability)  
 $\implies \inf_{\theta'} \text{CE}(\theta') = 0$

## Questions:

1. Does GD lead to zero loss?
2. Among the many possible solutions, which one it “prefers”?

# “Textbook” result

## □ Linear model

- **Fixed** embeddings  $\mathbf{h}_i \triangleq h_\theta(\mathbf{x}_i)$ .
- **Trainable** decoder  $\mathbf{W} \in \mathbb{R}^{k \times d}$

## □ Linearly-separable data (e.g. $d > n$ )

$$\exists \mathbf{W} : \mathbf{w}_{y_i}^T \mathbf{h}_i - \mathbf{w}_c^T \mathbf{h}_i > 0, \forall i \in [n], c \neq y_i \in [k]$$

$$\iff (\mathbf{e}_{y_i} - \mathbf{e}_c)^T \mathbf{W} \mathbf{h}_i = \langle \mathbf{W}, (\mathbf{e}_{y_i} - \mathbf{e}_c) \mathbf{h}_i^T \rangle > 0$$

$$\mathbf{e}_\ell = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 1 \\ 0 \\ \dots \\ 0 \end{bmatrix} \leftarrow \ell^{\text{th}} \text{ entry}$$

# “Textbook” result [Soudry et al.’18]

**Thm.** Assume separability. Run GD with  $\eta \leq 2/L$ .

Then,  $\lim_{k \rightarrow \infty} \text{CE}(\mathbf{W}_k) = 0$ .

Moreover,  $\lim_{k \rightarrow \infty} \|\mathbf{W}_k\| = \infty$  and

$$\lim_{k \rightarrow \infty} \left\langle \frac{\mathbf{W}_k}{\|\mathbf{W}_k\|}, \frac{\mathbf{W}^{\text{mm}}}{\|\mathbf{W}^{\text{mm}}\|} \right\rangle = 1$$

**Defn. (max-margin)** Let  $\mathbf{W}^{\text{mm}}$  be the max-margin classifier

$$\begin{aligned} \mathbf{W}^{\text{mm}} &= \operatorname{argmin}_{\mathbf{W}} \|\mathbf{W}\| \\ &\text{subj. to } \langle \mathbf{W}, (\mathbf{e}_{y_i} - \mathbf{e}_c) \mathbf{h}_i^T \rangle \geq 1, \quad \forall c \neq y_i \in [k], j \in [m] \end{aligned}$$

# Why nice?

---

## □ Insights on:

✓ what GD learns (impact of architectures/initializations)

[RZH03,SHN+18,JT18,GLSS18,JDST20,LL20,JT20] ++++

✓ role of optimizers (e.g. adaptive / mirror-descent)

[NLG+19,ALH21,PPVF21,SATA22,AF22] ++++

✓ **stepping stone to generalization (benign overfitting)**

[BLLT19,MRSY19,DKT19,MVS19,DL20,DL21,KZSS21,WT21,TPT21,CCBG22] ++++

✓ **loss design and hyperparameter tuning (imbalanced data)**

[KPOT21,CLB21,BKVT22]

.....

# Stepping stone to generalization

---

? How well does the solution found by GD generalize?



Implicit bias of GD

? How well does SVM solution generalize?

$$\mathbb{P}_{(\mathbf{h}, y) \sim D} \left( \min_{c \neq y} (\mathbf{e}_y - \mathbf{e}_c)^T \mathbf{W}^{\text{mm}} \mathbf{h} > 0 \right)$$



# Stepping stone to generalization

---

? How well does SVM solution generalize?

- ❑ **Catch:** Overparameterization ( $d > n$ ) makes “classical” margin-based bounds vacuous
  
- ❑ **Rescue:** modern\* tools from HD-stats/RMT and universality
  - Approximate Message Passing (AMP) [DMM09,MM12++]
  - Gordon’s comparison inequalities [Gor88,RV08,Sto09,CRPW’12++]  
[Sto13+,TOH15 ++]

[...]

\* Developed for **compressed-sensing**

<?> Can we push this storyline  
and (eventually) its implications  
to “new” setting of NTP in LMs

# Next-token Prediction (NTP)

---

## □ Training data

- Vocabulary  $\mathcal{V} \triangleq [V]$  of tokens/words
- (many many)  $n$  sequences  $(z_{i1}, z_{i2}, \dots, z_{iT})_{i \in [n]}$ ,  $z_{it} \in \mathcal{V}$

## □ Training loss

- $\min_{(\mathbf{W}, \theta)} \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \mathcal{L} \left( \underbrace{z_{it}}_{\text{next-word}}, \mathcal{S} \left( \mathbf{W} h_{\theta} \left( \underbrace{\mathbf{z}_{i, < t}}_{\text{Context } (z_{i,1}, \dots, z_{i,t-1})} \right) \right) \right)$
- For simplicity: focus on **last-token**

Denote  $(\mathbf{z}_{i, < T}, z_{iT}) \triangleq (\mathbf{x}_i, z_i)$

$$\min_{(\mathbf{W}, \theta)} \left\{ \text{CE}(\mathbf{W}, \theta) \triangleq \frac{1}{n} \sum_{i=1}^n -\log \left( \mathcal{S}_{z_i}(\mathbf{W} h_{\theta}(\mathbf{x}_i)) \right) \right\}$$

# NTP vs one-hot classification

---

## ❑ Ansatz #1:

- a. Contexts repeat
- b. **Multiple** possible next-tokens with **varying frequencies** after each distinct context.

[Shannon48]

	restaurants	0.05
	mountains	0.1
	rain	0.4
	UBC	0.01
Example: Vancouver is famous for its	_____	...
	culture	0
	sun	0
	affordability	0

## ❑ Ansatz #2: [Sparsity]

Not all vocabulary tokens are possible next-tokens per distinct context

# NTP training is sparse soft-label classification

---

## □ Data

- $m < n$  **distinct** contexts  $\mathbf{x}_j$  each with frequency  $\hat{\pi}_j$
- Each associated with **sparse probabilistic label**  $\hat{\mathbf{p}}_j \in \Delta^V$
- support set  $\mathcal{S}_j$  of  $\hat{\mathbf{p}}_j$ :  $|\mathcal{S}_j| < V$

## □ Loss wrt distinct contexts

$$\min_{(\mathbf{W}, \boldsymbol{\theta})} \left\{ \text{CE}(\mathbf{W}, \boldsymbol{\theta}) \triangleq - \sum_{j \in [m]} \hat{\pi}_j \sum_{z \in \mathcal{S}_j} \hat{p}_{j,z} \log \left( \mathcal{S}_z(\mathbf{W} h_{\boldsymbol{\theta}}(\mathbf{x}_j)) \right) \right\}$$

## Questions:

1. Does GD lead to the loss lower bound?
2. Among the many possible solutions, which one it “prefers”?

# Entropy Lower-bound

---

- Empirical (T-gram) entropy:

$$\mathcal{H} \triangleq \mathbb{E}_{(x,z) \sim \mathcal{T}_n} [-\log(\hat{p}(z|\mathbf{x}))] = - \sum_{j \in [m]} \sum_{z \in \mathcal{S}_j} \hat{\pi}_j \hat{p}_{j,z} \log(\hat{p}_{j,z})$$

- $\text{CE}(\boldsymbol{\theta}') = \mathcal{H} + \text{KL}(\hat{p}||q_{\boldsymbol{\theta}'}) \Rightarrow \text{CE}(\boldsymbol{\theta}') \geq \mathcal{H}$

# When is the lower bound reached?

---

➤ Consider linear model (fixed embeddings)

**Lemma.** The NTP loss reaches its lower bound if and only if the following two conditions hold.

**Defn. (NTP <sub>$\mathcal{H}$</sub> —compatibility)** There exists matrix  $\mathbf{W}^p$  such that for all  $j \in [m]$ :

$$\langle \mathbf{W}^p, (\mathbf{e}_z - \mathbf{e}_{z'}) \mathbf{h}_j^T \rangle = \log \left( \frac{\hat{p}_{j,z}}{\hat{p}_{j,z'}} \right) \quad \forall z \neq z' \in \mathcal{S}_j$$

**Defn. (NTP—separability)** There exists matrix  $\mathbf{W}^d$  such that for all  $j \in [m]$ :

$$\begin{aligned} \langle \mathbf{W}^d, (\mathbf{e}_z - \mathbf{e}_{z'}) \mathbf{h}_j^T \rangle &= 0 \quad \forall z \neq z' \in \mathcal{S}_j \\ \langle \mathbf{W}^d, (\mathbf{e}_z - \mathbf{e}_v) \mathbf{h}_j^T \rangle &\geq 1 \quad \forall z \in \mathcal{S}_j, v \notin \mathcal{S}_j \end{aligned}$$

# NTP compatibility & separability

□ Need  $\text{KL}(\hat{p} || q_{\theta'}) = 0$



1.  $S_z(\mathbf{W}\mathbf{h}_j) = \hat{p}_{j,z}, \forall z \in \mathcal{S}_j$



**NTP<sub>H</sub>—compatibility**

$$\langle \mathbf{W}^p, (\mathbf{e}_z - \mathbf{e}_{z'})\mathbf{h}_j^T \rangle = \log \left( \frac{\hat{p}_{j,z}}{\hat{p}_{j,z'}} \right)$$

2.  $S_v(\mathbf{W}\mathbf{h}_j) = 0, \forall v \notin \mathcal{S}_j$



**NTP—separability**

$$\begin{aligned} \langle \mathbf{W}^d, (\mathbf{e}_z - \mathbf{e}_{z'})\mathbf{h}_j^T \rangle &= 0 \\ \langle \mathbf{W}^d, (\mathbf{e}_z - \mathbf{e}_v)\mathbf{h}_j^T \rangle &\geq 1 \end{aligned}$$

**Lemma (Overparameterization).** If  $d > m$  and generic embeddings, then the two conditions hold.

$$\begin{aligned} S_z(\mathbf{a}) &= \frac{\exp(\mathbf{e}_z^T \mathbf{a})}{\sum_{v \in \mathcal{V}} \exp(\mathbf{e}_v^T \mathbf{a})} \\ &= \frac{1}{\sum_{v \in \mathcal{V}} \exp(-(\mathbf{e}_z - \mathbf{e}_v)^T \mathbf{a})} \end{aligned}$$



# Implicit bias

$$\mathcal{F} = \text{span}\{(\mathbf{e}_z - \mathbf{e}_{z'})\mathbf{h}_j^T : z \neq z' \in \mathcal{S}_j, j \in [m]\} \subseteq \mathbb{R}^{V \times d}$$

**Thm.** Assume NTP compatibility and separability. Run GD with  $\eta \leq 2/L$ .

Then, (i)  $\lim_{k \rightarrow \infty} \text{CE}(\mathbf{W}_k) = \mathcal{H}$ .

(ii)  $\lim_{k \rightarrow \infty} \mathbb{P}_{\mathcal{F}}(\mathbf{W}_k) = \mathbf{W}^*$

(iii)  $\lim_{k \rightarrow \infty} \|\mathbb{P}_{\perp}(\mathbf{W}_k)\| = \infty$  with  $\lim_{k \rightarrow \infty} \left\langle \frac{\mathbf{W}_k}{\|\mathbf{W}_k\|}, \frac{\mathbf{W}^{\text{mm}}}{\|\mathbf{W}^{\text{mm}}\|} \right\rangle = 1$

**Defn. (subspace component)**  $\mathbf{W}^* \in \mathcal{F}$  is the unique solution of:

$$\langle \mathbf{W}^*, (\mathbf{e}_z - \mathbf{e}_{z'})\mathbf{h}_j^T \rangle = \log \left( \frac{\hat{p}_{j,z}}{\hat{p}_{j,z'}} \right) \quad \forall z \neq z' \in \mathcal{S}_j, j \in [m]$$

**Defn. (orthogonal component)**  $\mathbf{W}^{\text{mm}} \in \mathcal{F}^{\perp}$  is the unique solution of:

$$\mathbf{W}^{\text{mm}} = \text{argmin}_{\mathbf{W}} \|\mathbf{W}\|$$

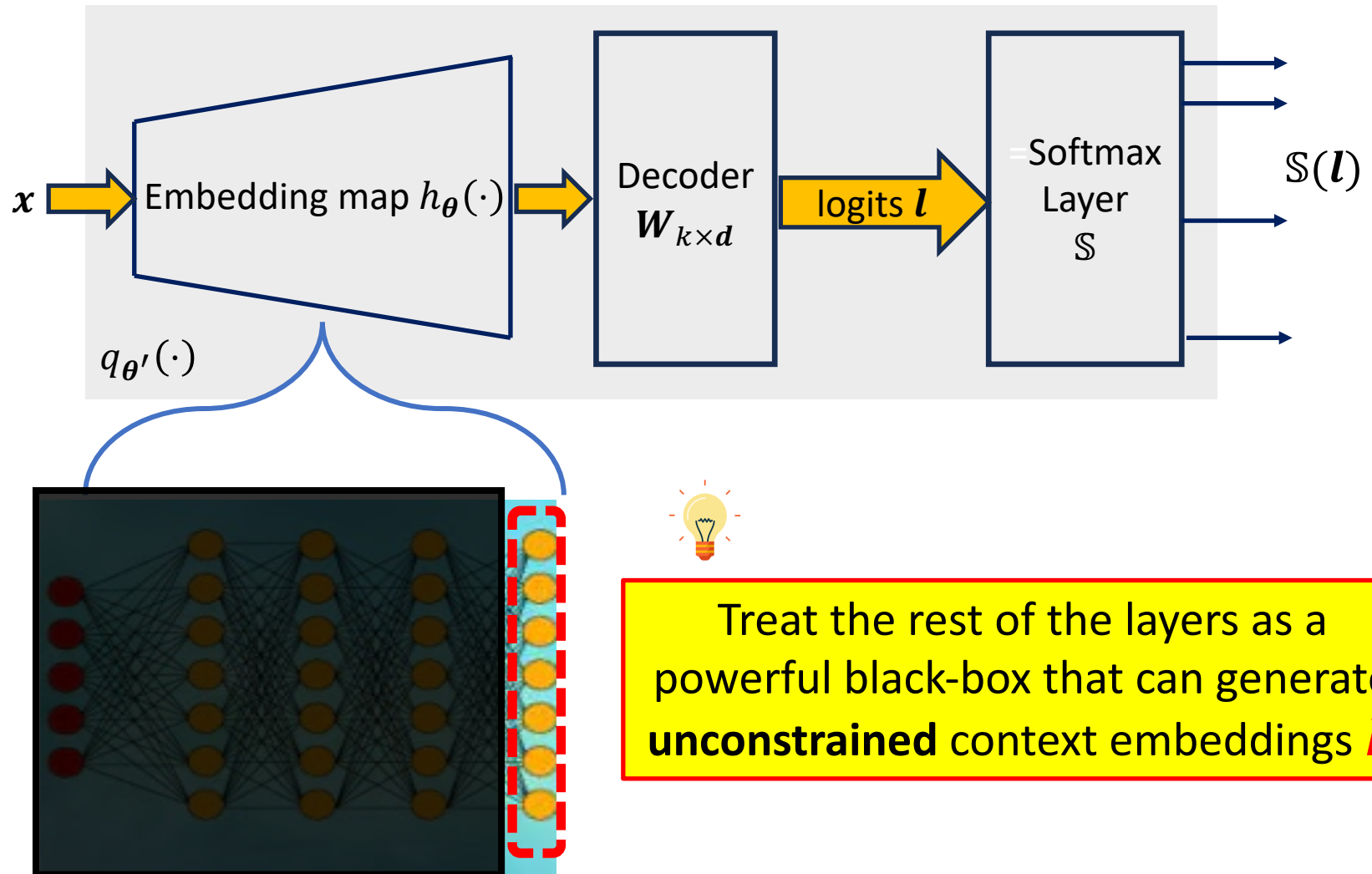
(NTP-SVM)

$$\text{subj. to } \langle \mathbf{W}, (\mathbf{e}_z - \mathbf{e}_{z'})\mathbf{h}_j^T \rangle = 0, \quad \forall z \neq z' \in \mathcal{S}_j$$

$$\langle \mathbf{W}, (\mathbf{e}_z - \mathbf{e}_v)\mathbf{h}_j^T \rangle \geq 1, \quad \forall z \in \mathcal{S}_j, v \notin \mathcal{S}_j, j \in [m]$$

<?> How to go beyond  
linear models?

# Unconstrained features



# NTP-UFM

---

□ **Unconstrained-features model (UFM):**

$$\min_{(\mathbf{W}, \mathbf{H})} \left\{ \text{CE}(\mathbf{W}, \mathbf{H}) \triangleq - \sum_{j \in [m]} \hat{\pi}_j \sum_{z \in \mathcal{S}_j} \hat{p}_{j,z} \log \left( S_z(\mathbf{W} \mathbf{h}_j) \right) \right\}$$

□  $\mathbf{W} \in \mathbb{R}^{V \times d}$ : **word** embeddings

□  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_m] \in \mathbb{R}^{d \times m}$ : **context** embeddings

□  $\mathbf{P} = [\{\hat{p}_{j,z}\}] \in \mathbb{R}^{V \times m}$ : **sparse** next-token probability matrix

## □ Unconstrained-features model (UFM):

What is the geometry of context/word embeddings in terms of the language statistics as encoded in the sparse conditional probability matrix  $P$ ?

□  $W \in \mathbb{R}^{V \times d}$ : word embeddings

□  $H \in \mathbb{R}^{d \times m}$ : context embeddings

If I were to optimize the log-bilinear NTP-UFM model, where does GD converge?

Input: sparse conditional probabilities matrix  $P \in \mathbb{R}^{V \times m}$

and corresponding support-set matrix  $S \in \{0,1\}^{V \times m}$

Output: The **implicit geometry** of word/context embeddings

i.e., angles between ctx-ctx / word-word / ctx-word vectors in  $\mathbb{R}^d$

# Proxy: Regularization path

---

$$\min_{(\mathbf{W}, \mathbf{H})} \left\{ - \sum_{j \in [m]} \hat{\pi}_j \sum_{z \in \mathcal{S}_j} \hat{p}_{j,z} \log \left( \mathcal{S}_z(\mathbf{W} \mathbf{h}_j) \right) + \lambda (\|\mathbf{W}\|^2 + \|\mathbf{H}\|^2) \right\}$$

□ **Goal:** Compute the solution as  $\lambda \rightarrow 0$

□ A proxy for GD-path (" $\lambda \rightarrow 0$ "  $\equiv$  " $k \rightarrow \infty$ ")

- Formal equivalence in linear settings [Ji et al. 20] [Rosset et al. '03]

# Logit-space relaxation

$$\min_{(\mathbf{W}, \mathbf{H})} \left\{ - \sum_{j \in [m]} \hat{\pi}_j \sum_{z \in \mathcal{S}_j} \hat{p}_{j,z} \log \left( \mathcal{S}_z(\mathbf{W} \mathbf{h}_j) \right) + \lambda (\|\mathbf{W}\|^2 + \|\mathbf{H}\|^2) \right\}$$

□  $\mathbf{L} = \mathbf{W}\mathbf{H} \in \mathbb{R}^{V \times d}$ : **logit** matrix

**Lemma.** The following relaxation to the  $\mathbb{R}^{V \times m}$  **logit-space** is tight:

$$\min_{\mathbf{L} : \text{rank}(\mathbf{L}) \leq d} \left\{ - \sum_{j \in [m]} \hat{\pi}_j \sum_{z \in \mathcal{S}_j} \hat{p}_{j,z} \log \left( \mathcal{S}_z(\mathbf{l}_j) \right) + \lambda \|\mathbf{L}\|_* \right\}$$

If  $\mathbf{L}_\lambda$  has SVD  $\mathbf{L}_\lambda = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , then for partially orthogonal matrix  $\mathbf{R} \in \mathbb{R}^{r \times d}$   
 $\mathbf{W}_\lambda = \mathbf{U}\sqrt{\mathbf{\Sigma}}\mathbf{R}$  and  $\mathbf{H}_\lambda = \mathbf{R}^T\sqrt{\mathbf{\Sigma}}\mathbf{V}$

# Large embedding space

---

$$\min_{L \in \mathbb{R}^{V \times d}, \text{rank}(L) \leq d} \left\{ - \sum_{j \in [m]} \hat{\pi}_j \sum_{z \in \mathcal{S}_j} \hat{p}_{j,z} \log \left( \mathbb{S}_z(\mathbf{l}_j) \right) + \lambda \|L\|_* \right\}$$

## □ Assumption: $d \geq V$

- Under this we can characterize regularization path
- Limiting but nontrivial:
  1. “# of contexts  $m$ ”  $\gg$  “dimension  $d$ ”
  2. how geometry depends on language statistics?



# Regularization-path of NTP-UFM

$$\mathcal{F} = \text{span}\{(\mathbf{e}_z - \mathbf{e}_{z'}) \tilde{\mathbf{e}}_j^T : z \neq z' \in \mathcal{S}_j, j \in [m]\} \subseteq \mathbb{R}^{V \times m}$$

**Thm.** Assume  $d \geq V - 1$ .

Then, (i)  $\lim_{\lambda \rightarrow 0} \text{CE}(\mathbf{L}_\lambda) = \mathcal{H}$ .

(ii)  $\lim_{\lambda \rightarrow 0} \mathbb{P}_{\mathcal{F}}(\mathbf{L}_\lambda) = \mathbf{L}^*$

(iii)  $\lim_{\lambda \rightarrow 0} \|\mathbb{P}_{\perp}(\mathbf{L}_\lambda)\| = \infty$  with  $\lim_{\lambda \rightarrow 0} \left\langle \frac{\mathbf{L}_\lambda}{\|\mathbf{L}_\lambda\|}, \frac{\mathbf{L}^{\text{mm}}}{\|\mathbf{L}^{\text{mm}}\|} \right\rangle = 1$

**Defn.**  $\mathbf{L}^* \in \mathcal{F}$  is the unique solution of:

$$\mathbf{L}_{z,j} - \mathbf{L}_{z',j} = \log \left( \frac{\hat{p}_{j,z}}{\hat{p}_{j,z'}} \right) \quad \forall z \neq z' \in \mathcal{S}_j, j \in [m]$$

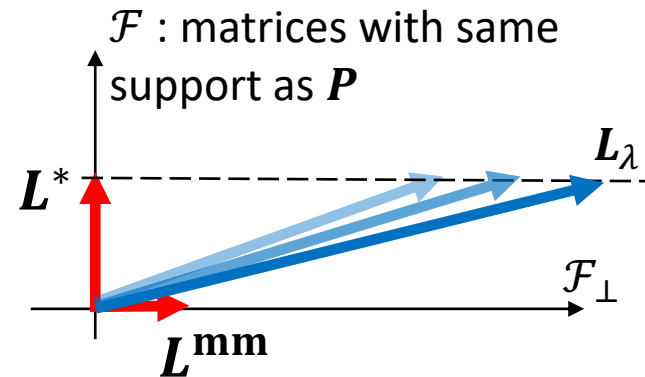
**Defn.**  $\mathbf{L}^{\text{mm}} \in \mathcal{F}^{\perp}$  is a solution of:

$$\begin{aligned} & \min_{\mathbf{L}} \quad \|\mathbf{L}\|_* && \text{(NTP-SVM)} \\ & \text{subj. to} \quad \mathbf{L}_{z,j} - \mathbf{L}_{z',j} = 0, \quad \forall z \neq z' \in \mathcal{S}_j \\ & \quad \quad \quad \mathbf{L}_{z,j} - \mathbf{L}_{v,j} \geq 1, \quad \forall z \in \mathcal{S}_j, v \notin \mathcal{S}_j, j \in [m] \end{aligned}$$

# Regularization-path of NTP-UFM

As  $\lambda \rightarrow 0$ , for some  $\rho(\lambda) \rightarrow \infty$ :

$$L_\lambda \approx L_{\text{sparse}} + \rho(\lambda) \cdot L_{\text{low-rank}}$$



□  $L^* \stackrel{\text{def}}{=} L_{\text{sparse}}$  inherits **sparsity** of  $P$  and depends on frequencies of in-support tokens

□  $L^{\text{mm}} \stackrel{\text{def}}{=} L_{\text{low-rank}}$  minimizes nuclear-norm promoting **low-rankness** and only depends on sparsity pattern  $S$  (not on frequencies)

**Dominant as  $\lambda \rightarrow 0$**

# NTP max-margin logits

□ In some special cases, can compute  $L^{\text{mm}}$  in closed form

**Prop.** Suppose  $\mathcal{S}$  contains all  $m = \binom{V}{k}$  support sets of size  $k$ .

Then, (i)  $L^{\text{mm}} = (\mathbf{I}_V - \mathbf{1}\mathbf{1}^T)\mathcal{S} \stackrel{\text{def}}{=} \bar{\mathcal{S}}$ .

(ii) Word embeddings form equiangular tight frame

(iii) Context embeddings are equinorm and  $\mathbf{h}_j$  is colinear to  $\sum_{z \in \mathcal{S}_j} w_z$

$$\mathcal{S} = \begin{bmatrix} 1 & 1 & & & & 0 \\ 1 & 0 & \dots & 1 & & 0 \\ 0 & 1 & & 1 & & 0 \\ \cdot & \cdot & & \cdot & & \cdot \\ \cdot & \cdot & & \cdot & \dots & \cdot \\ \cdot & \cdot & & \cdot & & \cdot \\ 0 & 0 & \dots & 0 & & 1 \\ 0 & 0 & & 0 & & 1 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} 1 & 1 & & & & 0 \\ 1 & 0 & \dots & 1 & & 0 \\ 0 & 1 & & 1 & & 0 \\ \cdot & \cdot & & \cdot & & \cdot \\ \cdot & \cdot & & \cdot & \dots & \cdot \\ \cdot & \cdot & & \cdot & & \cdot \\ 0 & 0 & \dots & 0 & & 1 \\ 0 & 0 & & 0 & & 1 \end{bmatrix}} \right\} V$$

$$m = \binom{V}{2}$$

# NTP max-margin logits

□ In some special cases, can compute  $L^{\text{mm}}$  in closed form

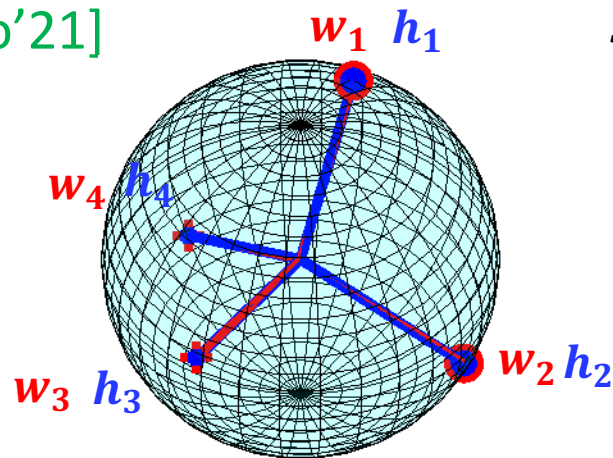
**Prop.** Suppose  $\mathcal{S}$  contains all  $m = \binom{V}{k}$  support sets of size  $k$ .

Then, (i)  $L^{\text{mm}} = (\mathbf{I}_V - \mathbf{1}\mathbf{1}^T)\mathcal{S} \stackrel{\text{def}}{=} \bar{\mathcal{S}}$ .

(ii) Word embeddings form equiangular tight frame

(iii) Context embeddings are equinorm and  $\mathbf{h}_j$  is colinear to  $\sum_{z \in \mathcal{S}_j} w_z$

□ Special case  $k = 1$ , recovers the **Neural Collapse** geometry by [Papayan, Han, Donoho'21]



$$\mathcal{S} = \begin{bmatrix} 1 & 1 & & & & 0 \\ 1 & 0 & \dots & 1 & & 0 \\ 0 & 1 & & 1 & & 0 \\ \cdot & \cdot & & \cdot & & \cdot \\ \cdot & \cdot & & \cdot & \dots & \cdot \\ \cdot & \cdot & & \cdot & & \cdot \\ 0 & 0 & \dots & 0 & & 1 \\ 0 & 0 & & 0 & & 1 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} 1 & 1 & & & & 0 \\ 1 & 0 & \dots & 1 & & 0 \\ 0 & 1 & & 1 & & 0 \\ \cdot & \cdot & & \cdot & & \cdot \\ \cdot & \cdot & & \cdot & \dots & \cdot \\ \cdot & \cdot & & \cdot & & \cdot \\ 0 & 0 & \dots & 0 & & 1 \\ 0 & 0 & & 0 & & 1 \end{bmatrix}} \right\} V$$

$$m = \binom{V}{2}$$

# NTP max-margin logits

---

□ In some special cases, can compute  $\mathbf{L}^{\text{mm}}$  in closed form

**Prop.** Suppose  $\mathcal{S}$  contains all  $m = \binom{V}{k}$  support sets of size  $k$ .

Then, (i)  $\mathbf{L}^{\text{mm}} = (\mathbf{I}_V - \mathbf{1}\mathbf{1}^T)\mathcal{S} \stackrel{\text{def}}{=} \bar{\mathcal{S}}$ .

(ii) Word embeddings form equiangular tight frame

(iii) Context embeddings are equinorm and  $\mathbf{h}_j$  is colinear to  $\sum_{z \in \mathcal{S}_j} w_z$

□ In general, need to solve SDP.

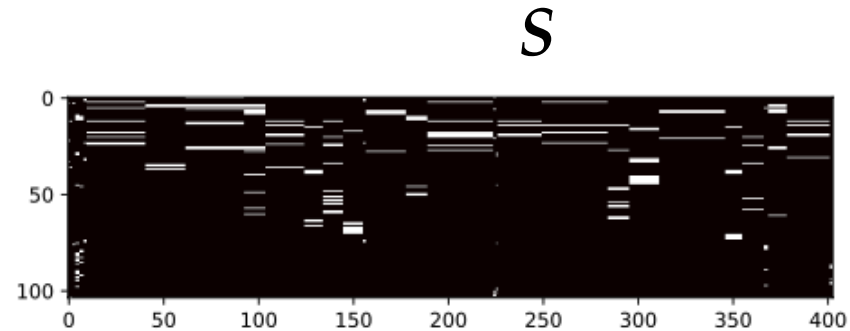
□ But, experimentally  $\bar{\mathcal{S}}$  is a good “proxy”

# Experiment

---

## Data:

- Synthetic extracted from TinyStories\*
- $n = 3050$  contexts of length  $T = 5$
- $m = 400$  *distinct* contexts
- $V = 104$



\*

“a little girl named lily” .... {“and”, “was”, “found”, “had”, “went”, “.”}

“there was a little boy” .... {“named”, “called”, “and”, “.”, “who”, “had”, “with”}

# Experiment

## Data:

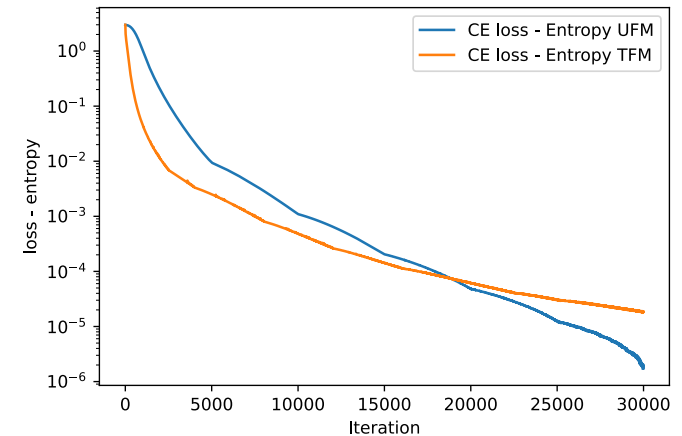
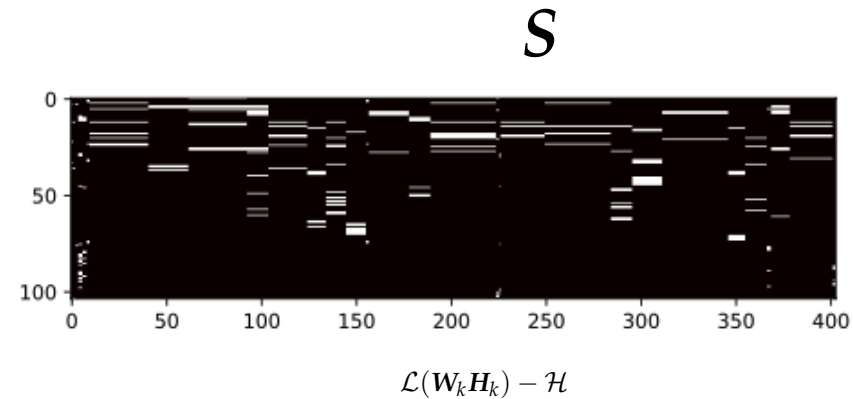
- Synthetic extracted from TinyStories
- $n = 3050$  contexts of length  $T = 5$
- $m = 400$  *distinct* contexts
- $V = 104$

## Language Model:

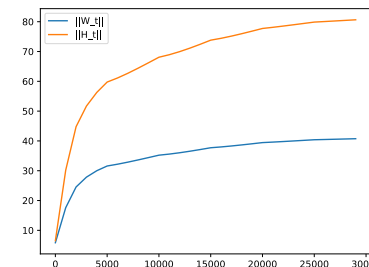
- 4-layer TF
- $d = 128$
- Trained with Adam-W for 30k epochs

## Analysis Model:

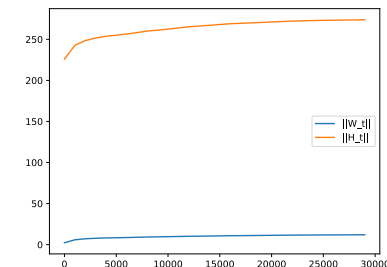
- Unconstrained features model (**UFM**)
- $d = 128$
- Trained with Adam-W for 30k epochs



UFM on Tinstory Context



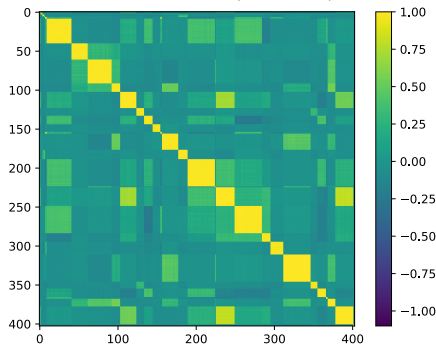
Transformer on Tinstory Context



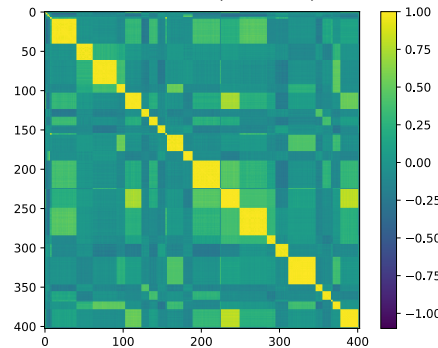
# Numerical Example

- ✓ UFM is a good proxy
- ✓ Eigenfactors of NTP-max-margin  $L^{\text{mm}}$  predict geometry
- ✓ Support overlaps already give a good proxy for geometry

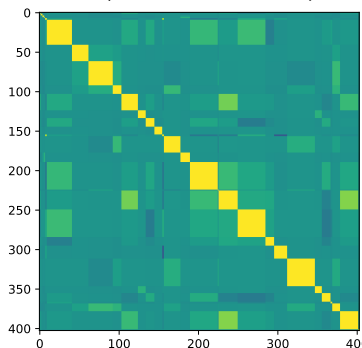
TF:  $\cos(H, H)$



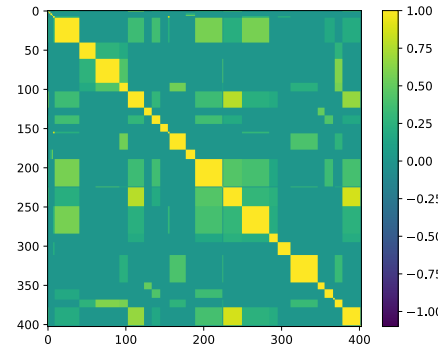
UFM:  $\cos(H, H)$



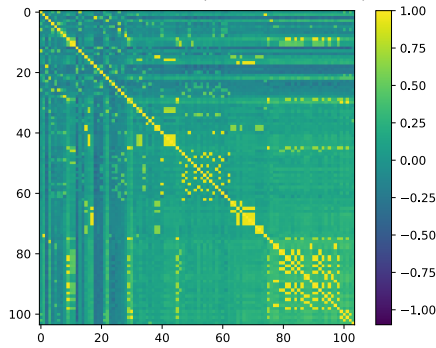
$\cos(H^{\text{mm}}, H^{\text{mm}})$



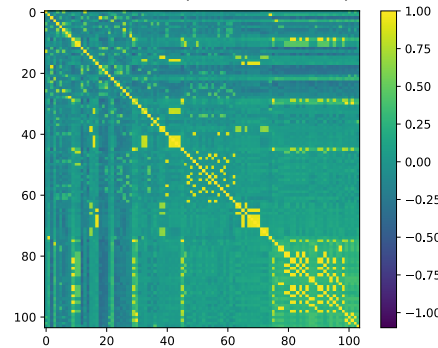
$\cos(S, S)$



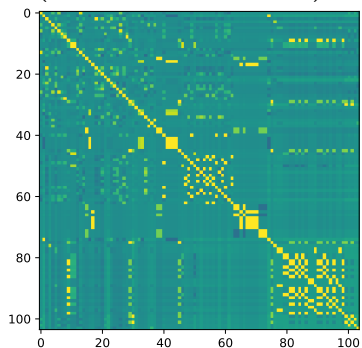
TF:  $\cos(W^\top, W^\top)$



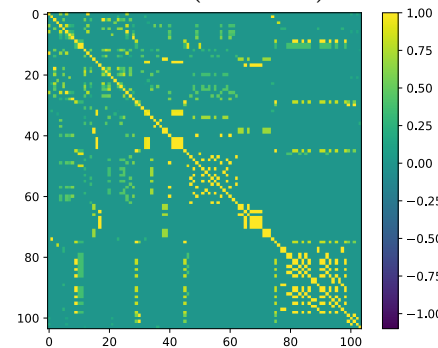
UFM:  $\cos(W^\top, W^\top)$



$\cos(W^{\text{mm}\top}, W^{\text{mm}\top})$



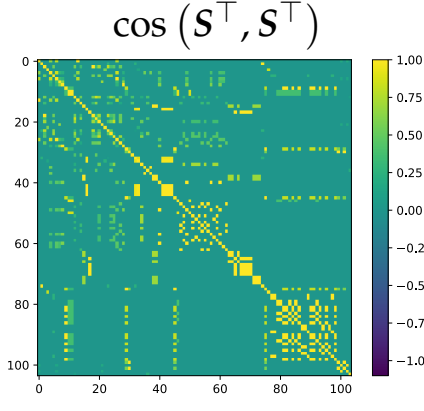
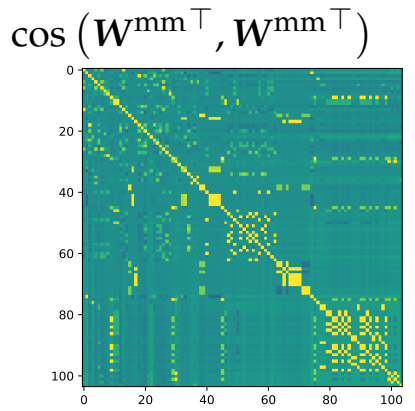
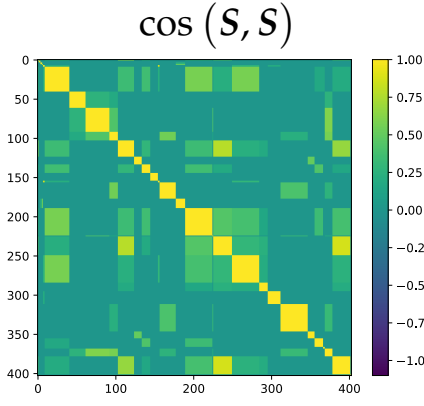
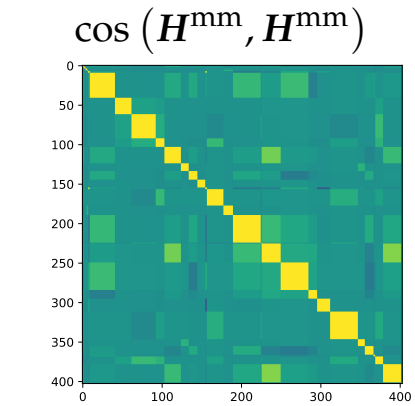
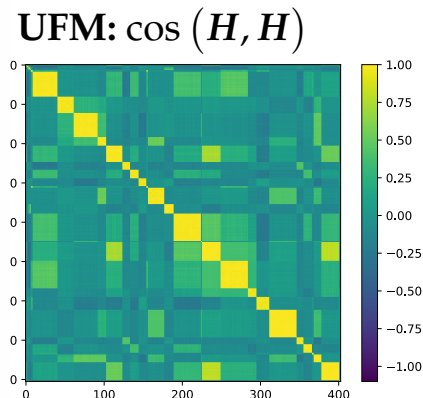
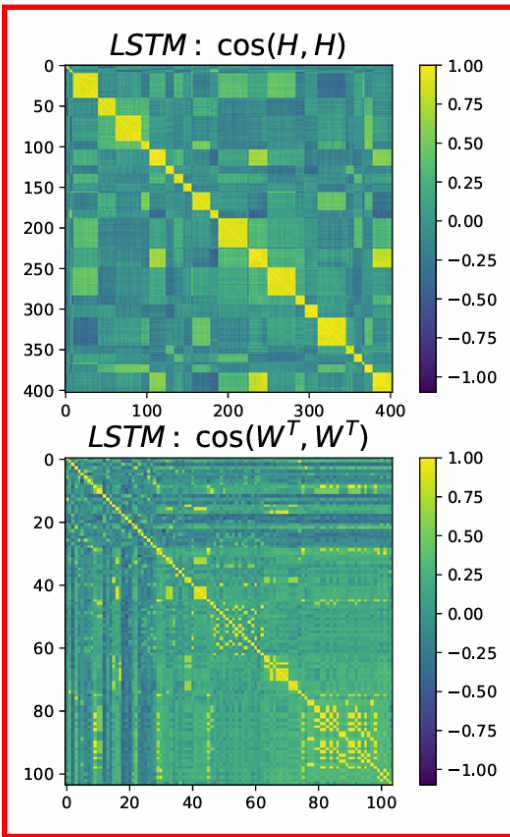
$\cos(S^\top, S^\top)$





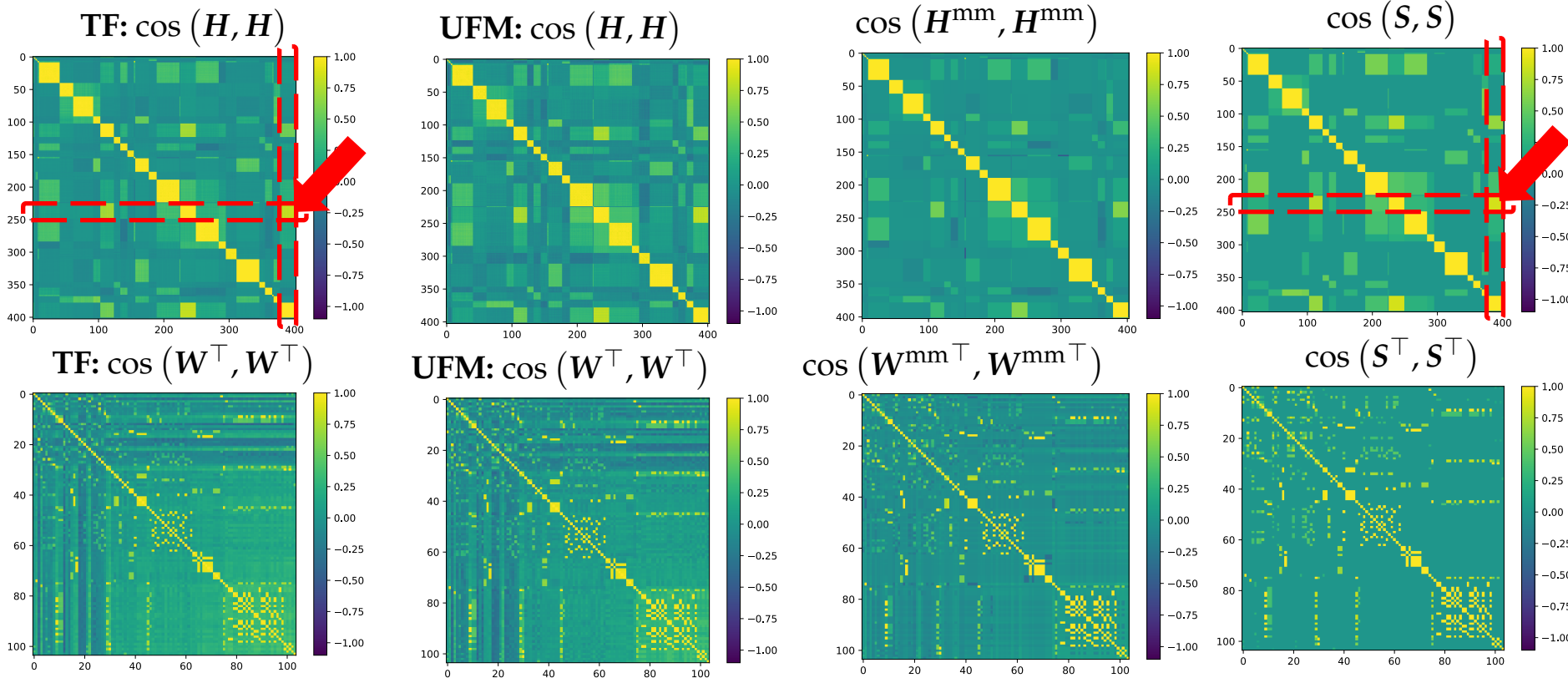
# Numerical Example

- ✓ UFM is a good proxy
- ✓ Eigenfactors of NTP-max-margin  $L^{\text{mm}}$  predict geometry
- ✓ Support overlaps already give a good proxy for geometry



# Numerical Example

- ✓ UFM is a good proxy
- ✓ Eigenfactors of NTP-max-margin  $L^{\text{mm}}$  predict geometry
- ✓ Support overlaps already give a good proxy for geometry



Context 1: "boy named timmy . timmy"  
Context 2: "kid called lilly. she"

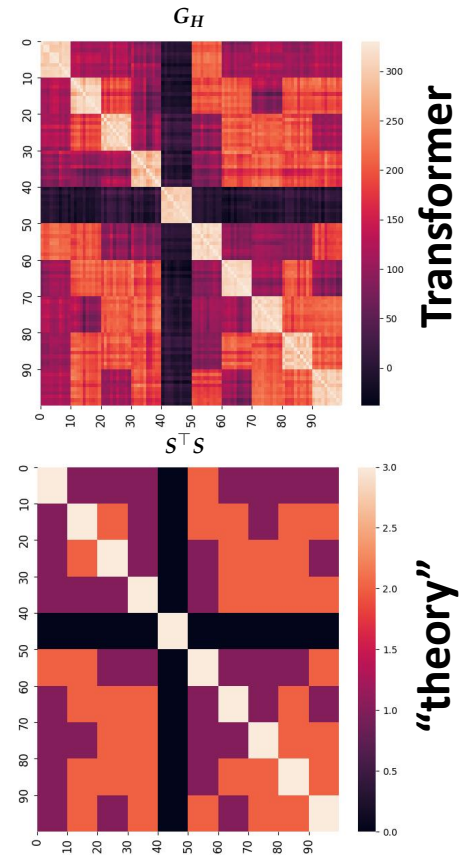
**Subspace collapse:** Embeddings of ctxs whose supports overlap collapse

# Summary

A framework for **mapping language patterns to embeddings geometry** via:

1. Framing NTP as **sparse soft labels** classification
2. Applying **unconstrained features**
3. Leveraging **implicit bias** viewpoint

→ Word/context embeddings as mtX factorization of  $L_\lambda \approx L_{\text{sparse}} + \rho(\lambda) \cdot L_{\text{low-rank}}$



# NTP: Open questions

---

- ❑ Directly-related questions:
  - ? **Gradient-descent** convergence
  - ? What is the impact of Zipf-law imbalances on convergence?
  - ?  **$d < V$ : Do linguistics sparsity patterns lead to low-rank solutions?**
  - ? Geometry at higher layers of linguistic understanding, e.g. concepts
  
- ❑ The setting is clearly “statistical”:
  - ? How do these optimization results inform **generalization**?
  - ? What is the statistical role of margin btwn in/out-of-support tokens?
  - ? What are good data models to study these
  - ? When is it good to train long or is better to stop early?

# Summary

A framework for **mapping language patterns to embeddings geometry** via:

1. Framing NTP as **sparse soft labels** classification
2. Applying **unconstrained features**
3. Leveraging **implicit bias** viewpoint

→ Word/context embeddings as mtv factorization of  $L_\lambda \approx L_{\text{sparse}} + \rho(\lambda) \cdot L_{\text{low-rank}}$

1. CT, Implicit Bias of Next-token Prediction, NeurIPS 2024
2. Zhao, Behnia, Vakilan, CT, Implicit Geometry of Next-token Prediction: From Language Sparsity Patterns to Model Representations, COLM 2024.

**Thank you!**

