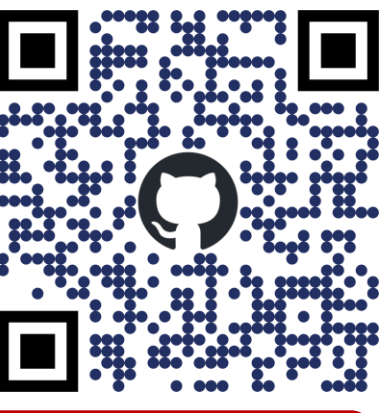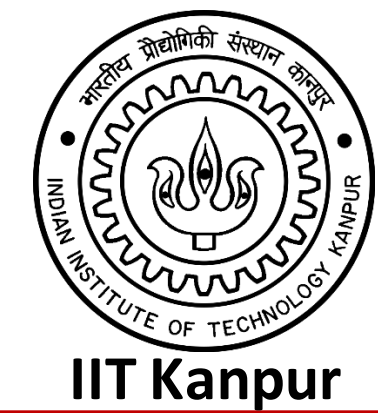# Deep Encoders with Auxiliary Parameters for Extreme Classification

Kunal Dahiya, Sachin Yadav, Sushant Sondhi, Deepak Saini, Sonu Mehta, Jian Jiao, Sumeet Agarwal, Purushottam Kar, Manik Varma
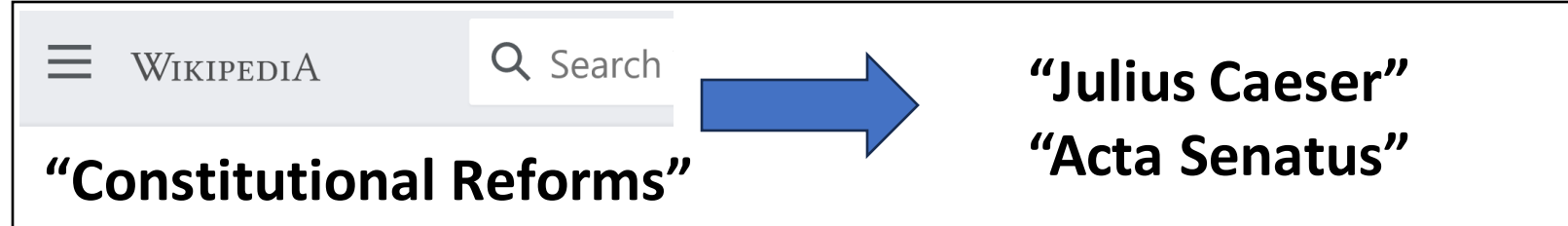
IIT Delhi    Microsoft    IIT Kanpur

**Goal:** Annotate a data point with the *most relevant subset* of labels from an extremely large set.

**DEXA:** Deep Encoders with Auxiliary Parameters for Extreme Classification

**Theoretical Results:** Provable accurate training and crisp generalization bounds

**Results:** Significant gains in offline evaluation; Readily incorporate with existing architectures

## Extreme Classification (XC)

WIKIPEDIA  Q Search

"Constitutional Reforms" → "Julius Caeser" "Acta Senatus"

### Applications
- Related web-page recommendation
- Matching user queries to advertiser keywords
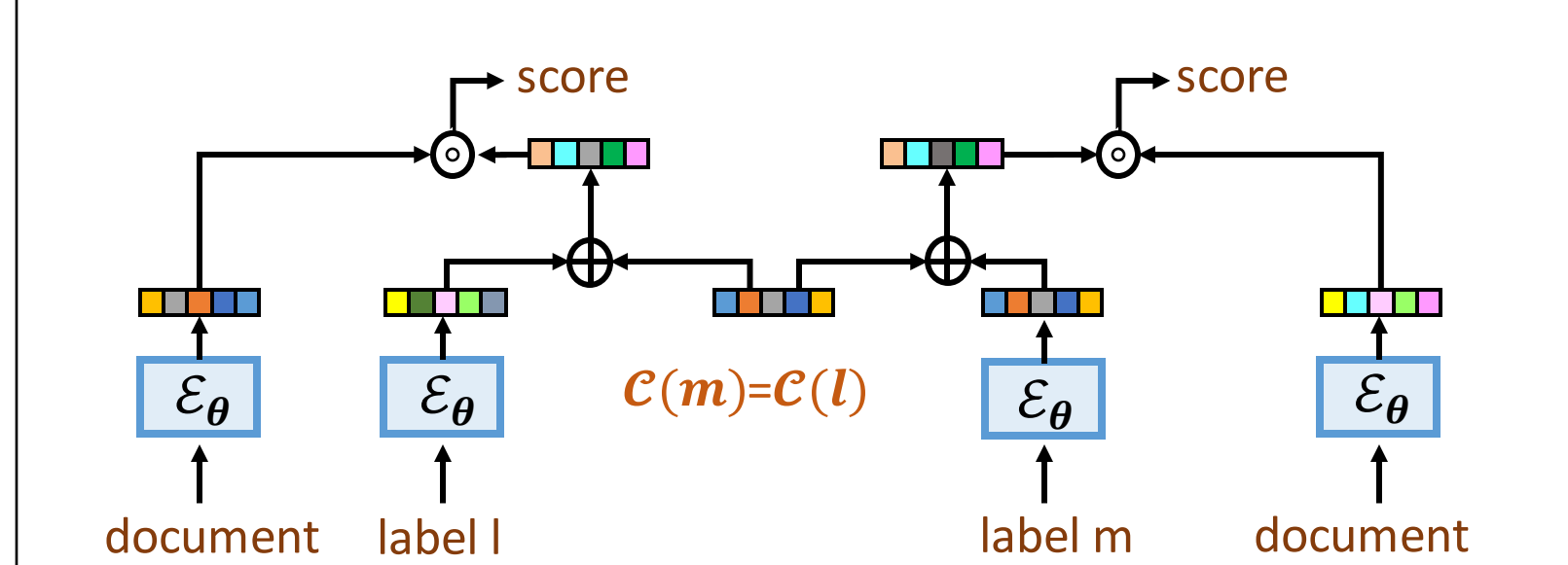- Product-to-product recommendation

### Semantic Gap in Siamese Networks
- Label text is unable to capture full meaning in case of short-text applications
- This leads to distortions in encoder training and sub-optimal accuracies
- Naïve solution to add $L$ free vectors improve accuracies but does not scale
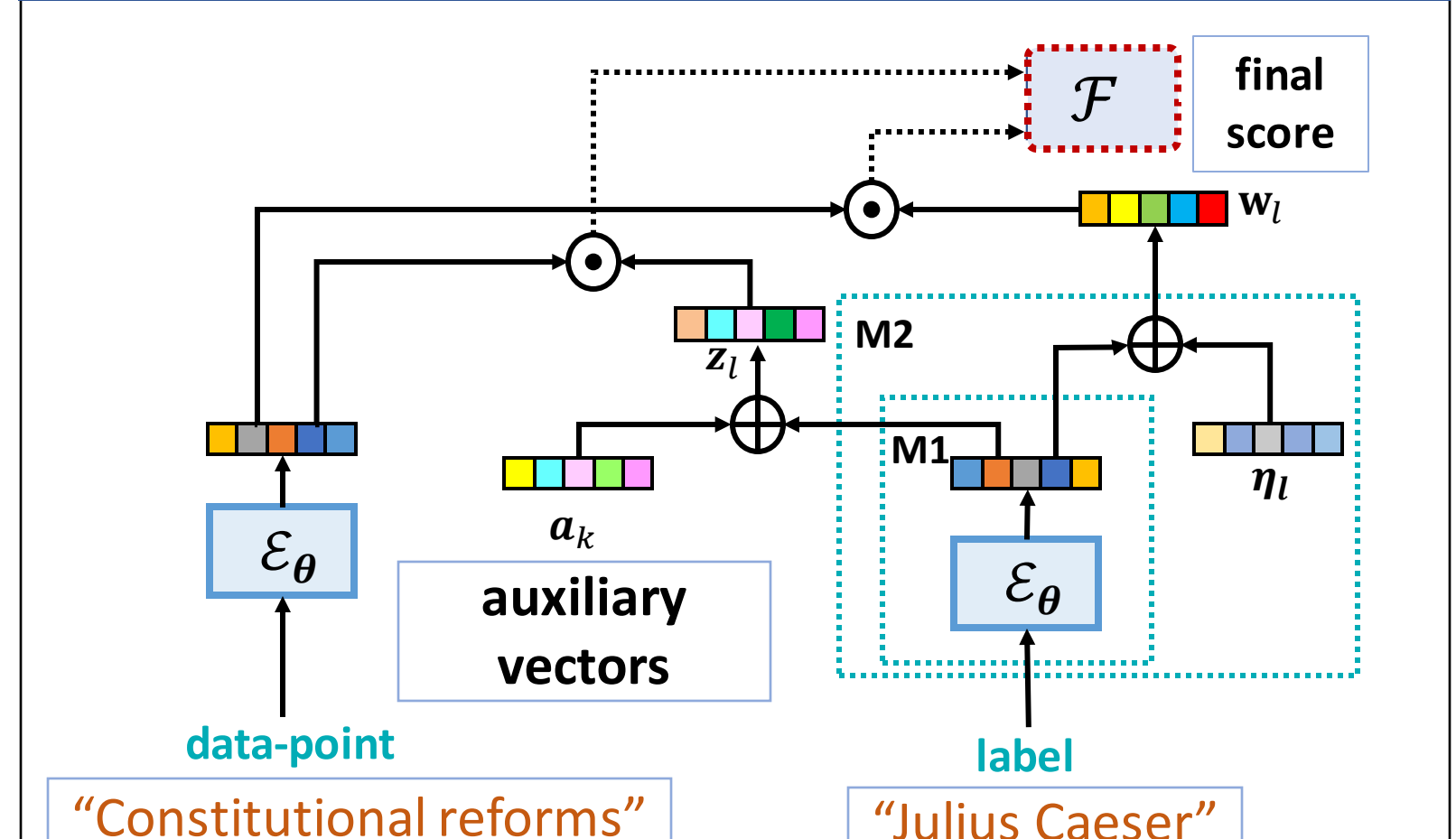
## DEXA: Foundations

**Dataset:** $\{x_i, y_i\}_{i=1}^N, \{z_l\}_{l=1}^L y_i \in \{-1,1\}^L$ where $x_i, z_l \in \mathcal{X}$ are doc/label text respectively

**Goal:** Learn params $\theta$ for embedding architecture $\mathcal{E}_\theta : \mathcal{X} \to \mathbb{R}^D$, One-vs.-All (OvA) classifiers $w_l \in \mathbb{R}^D$, one for each label, to minimize triplet loss $\mathcal{L}$.



$\mathcal{C}(m) = \mathcal{C}(l)$

document    label l    label m    document

Related labels may have similar correction terms.

## Architecture



$\mathcal{F}$  final score

$w_l$    M2    $z_l$    M1    $\eta_l$

$\mathcal{E}_\theta$    $a_k$ **auxiliary vectors**    $\mathcal{E}_\theta$

data-point    label

"Constitutional reforms"    "Julius Caeser"

## Modular Training & Efficient Prediction

**Module I (*Encoder with Auxiliary Parameters*):**
- Cluster $L$ labels into $K \ll L$ clusters
- Introduce $K$ auxiliary parameters ($a_k$)
- Use $\mathcal{E}_\theta(z_l) + a_k$ as surrogate for $w_l$ and train embed. arch. $\hat{\theta}$ using a Siamese loss

**Module II (Extreme Classifiers):** Train $\eta_l$ in $\mathcal{O}(ND \log L)$ time using +ve & hard –ve labels

**Prediction:** Efficient procedure, taking $\mathcal{O}(D^2 + D \log L)$ time per point

## Illustrative Example

**Data Point:** Constitutional reforms of Julius Caesar..

| Method | Top-5 predictions |
|---|---|
| DEXA | Acta Senatus ✔<br>Centuria ✔<br>Roman Law ✔<br>Interrex ✔<br>Byzantine Senate ✔ |
| NGAME | Julius Caeser ✘<br>Assassination of Julius Caesar ✘<br>Caesarism ✘<br>Constitution of the Roman Republic ✘<br>Caesar's civil war ✘ |

## Provable Accurate Training

**Lemma**: Consider a linear encoder parametrized by $E$ & DEXA with auxiliary params $A$, the gradient norms at optimal value $E^*$:

$$\left\| \Delta_E \mathcal{L}(E^*) \right\|_2 \leq 2\|E_*\|_2^2 \sqrt{\frac{1}{L}\sum \|\Delta_l\|_2^2}$$

$$\left\| \Delta_E \mathcal{L}(E^*, A) \right\|_2 \leq 4\|E_*\|_2^2 \sqrt{\frac{1}{L}\sum \sigma_k^2}$$

where, $\Delta_l$ is semantic gap for label $l$, and $\sigma_k^2$ is intra-cluster variance in semantic gap

- DEXA offers smaller encoder gradient, indicating a more faithful recovery of true encoder parameters, if $\sigma_k^2 \ll \sum_{l \in C_k} \|\Delta_l\|_2^2$

- Even if the individual $\Delta_l$ are large in a cluster, DEXA offers faithful encoder recovery so long that semantic gaps are similar to each other

## Generalization Bounds

**Theorem**: Suppose DEXA is used with an encoder parameters $\theta$ and auxiliary parameters $A$, then with probability $1-\delta$, we have
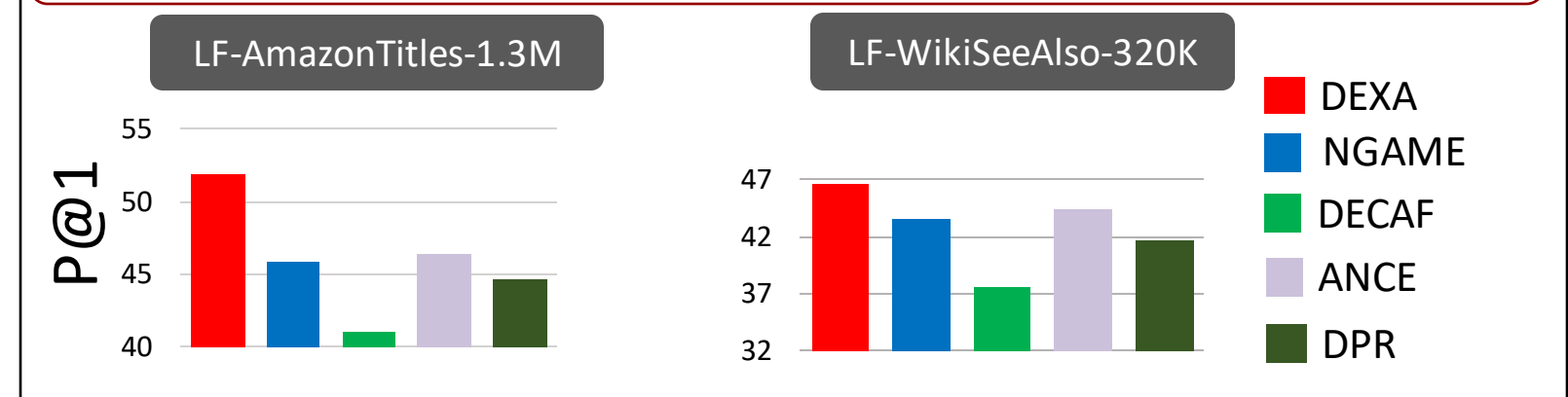
$$\ell(\theta, A) \leq \hat{\ell}_N(\theta, A) + \epsilon(N) + \sqrt{\frac{\ln\frac{1}{\delta}}{N}} + \frac{\Delta \ln N}{\sqrt{N}}$$

where, $\Delta = \mathcal{O}(\ln(DK))$ apart from numerical constants independent of $L$
- $\epsilon(N)$ captures the dependence of the excess risk on the encoder parameter characteristics
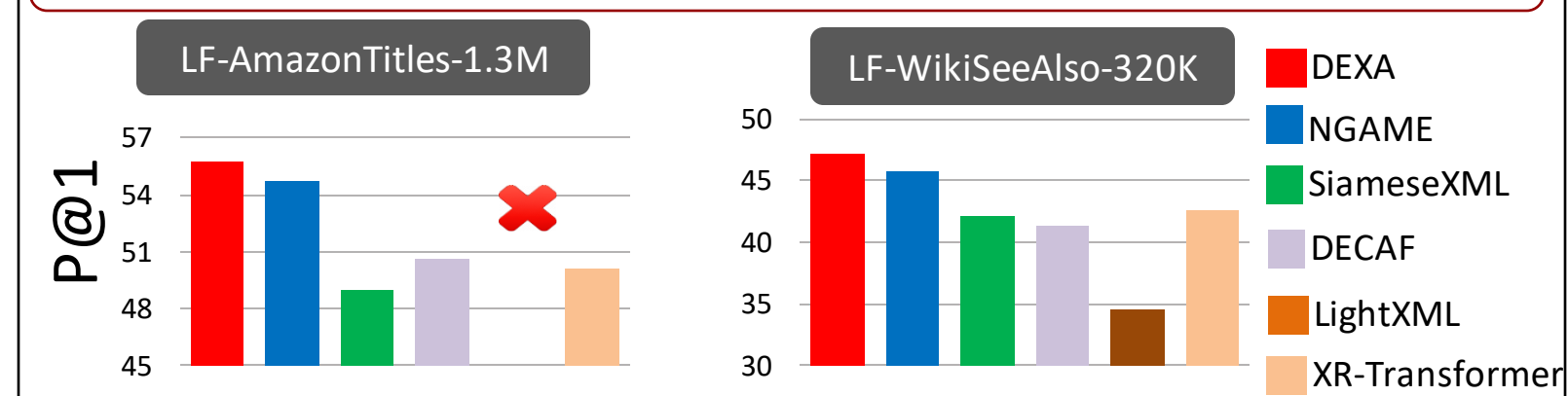- Independent of $L$ in favour of $\mathcal{O}(\ln(DK))$

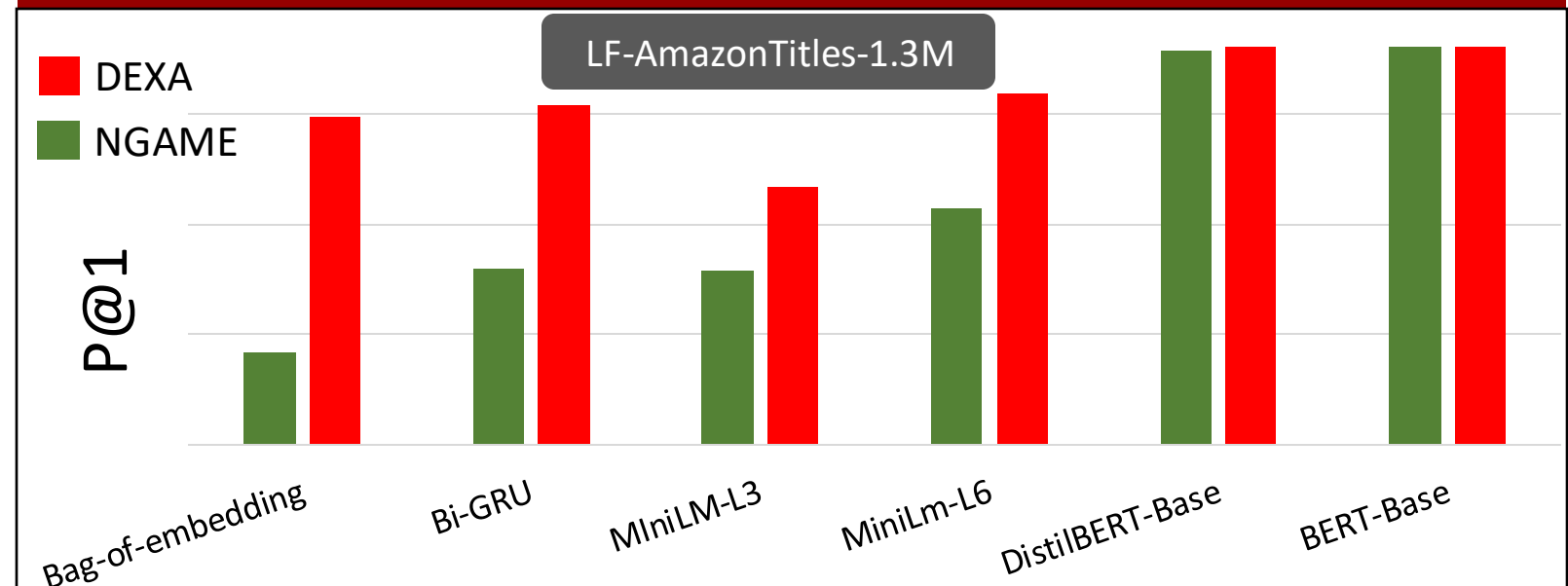## Embeddings on Benchmark Datasets

Up to 11% more accurate embeddings



LF-AmazonTitles-1.3M    LF-WikiSeeAlso-320K

DEXA, NGAME, DECAF, ANCE, DPR

## End-to-end Results on Benchmark Datasets

Up to 12% gains over leading existing XC methods



LF-AmazonTitles-1.3M    LF-WikiSeeAlso-320K

DEXA, NGAME, SiameseXML, DECAF, LightXML, XR-Transformer

## Ablation over architectures



DEXA, NGAME

LF-AmazonTitles-1.3M

Bag-of-embedding, Bi-GRU, MIniLM-L3, MiniLm-L6, DistilBERT-Base, BERT-Base

## Sponsored Search-40M

7-15% gains in matching queries to keywords



MiniLM-L3-v2, DistilBERT-base

#Auxiliary vectors (K) (millions)