

# Sparse Local Embeddings for Extreme Multi-label Classification

Kush Bhatia†

Himanshu Jain§

Purushottam Kar‡

Manik Varma†

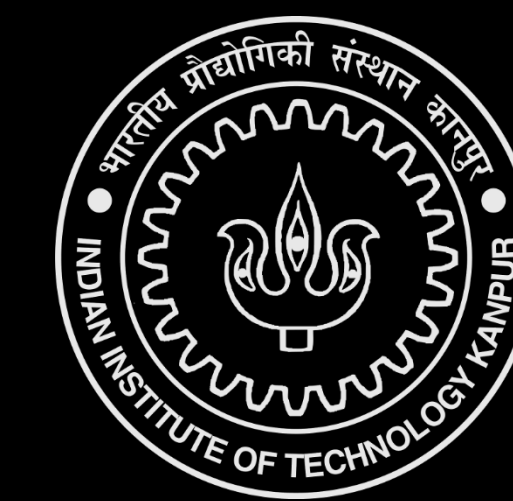
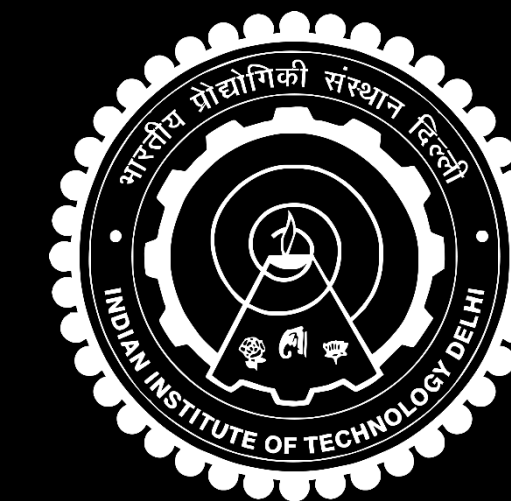
Prateek Jain†

† Microsoft Research, Bengaluru, India

§ Indian Institute of Technology Delhi

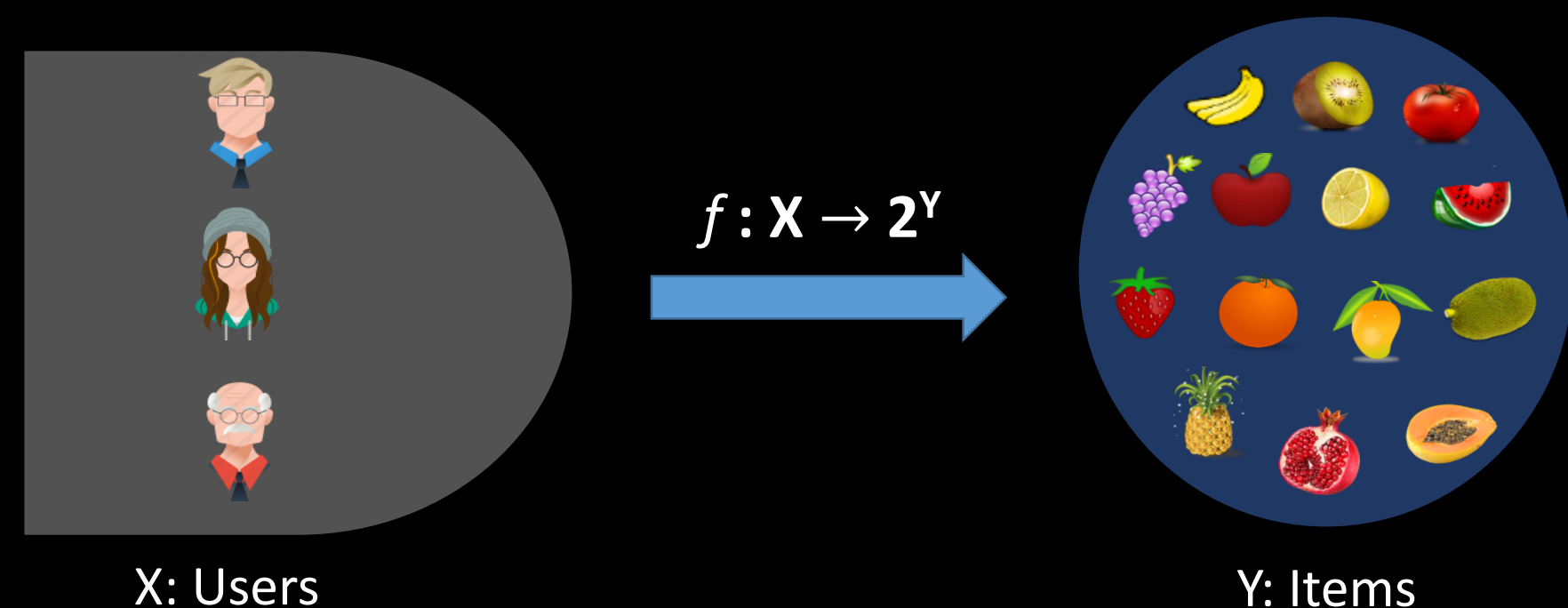
‡ Indian Institute of Technology Kanpur

Microsoft  
Research



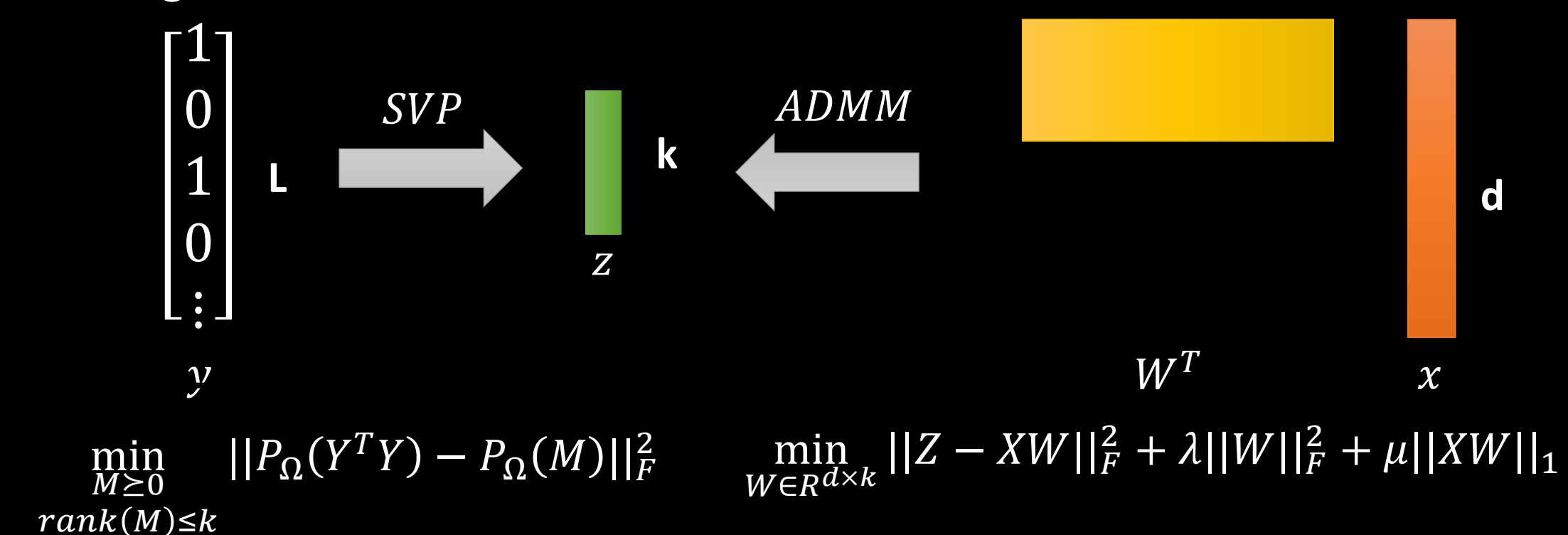
## Extreme Multi-Label Learning

- Learning with millions of labels
- New paradigm for ranking and recommendation

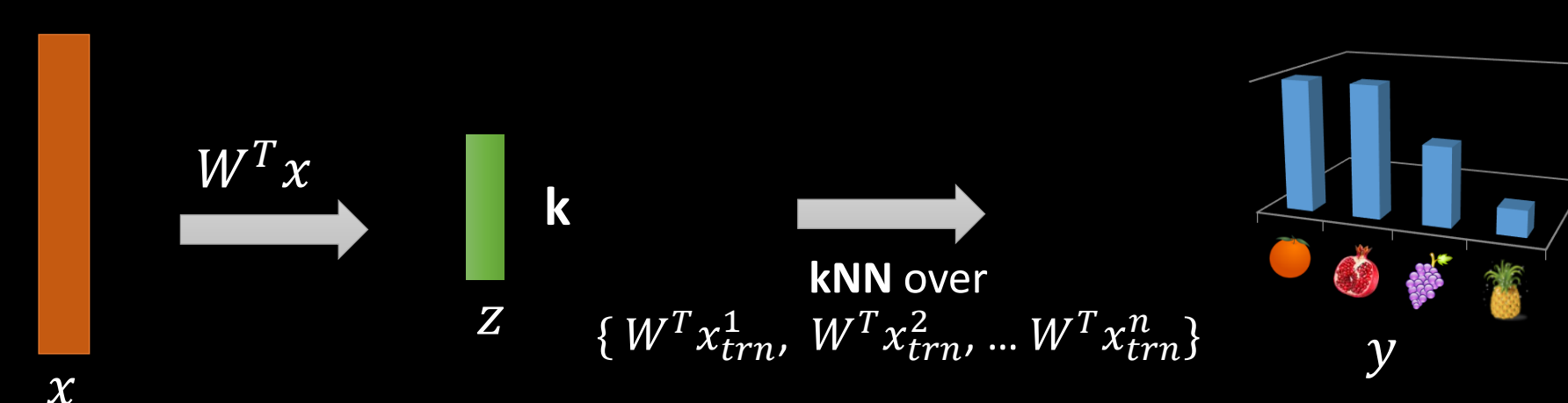


## SLEEC : Algorithm

Training:



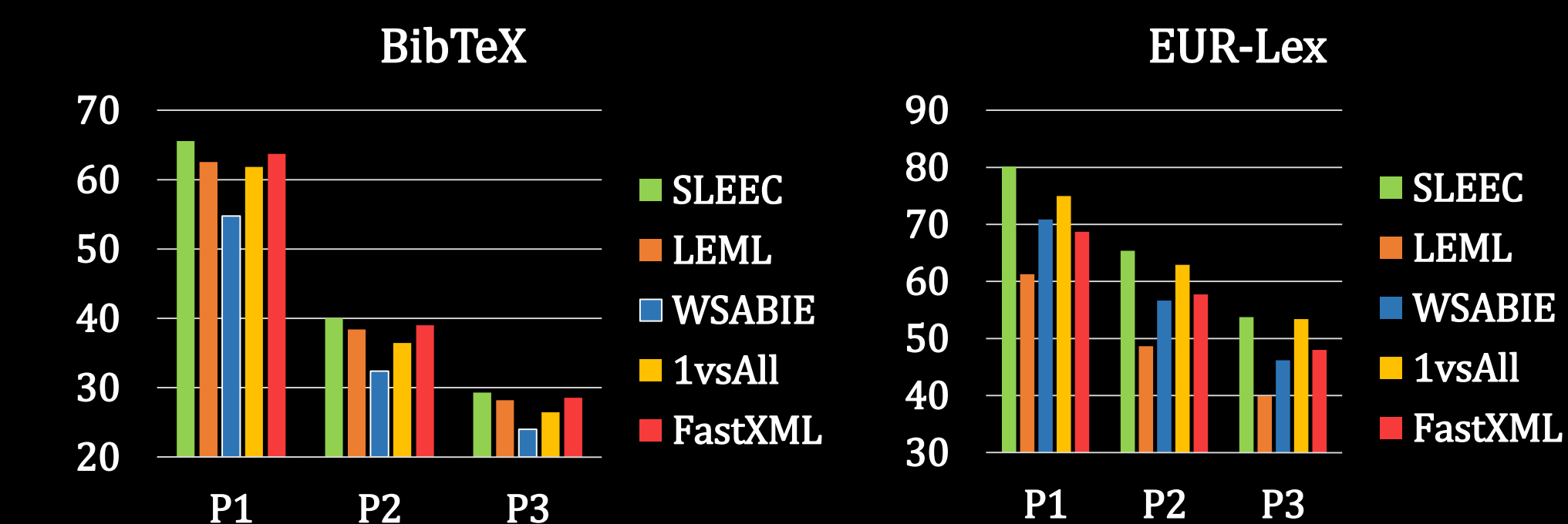
Prediction:



## Results

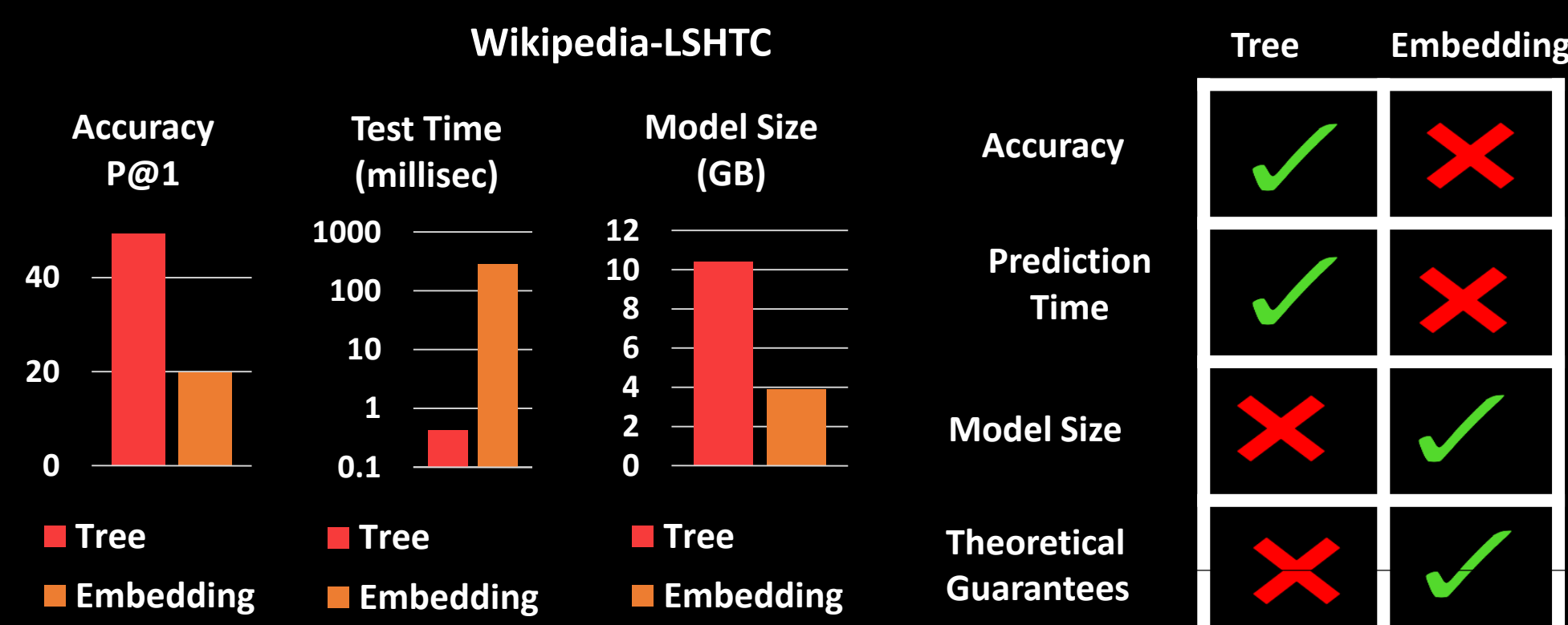
Dataset Statistics:

Data Set	# of Training Points	# of Test Points	# of Dimensions	# of Labels
EurLEX	17413	1935	5000	3993
BibTeX	4,880	2,515	1,836	159
Wikipedia-LSHTC	1.77M	587k	1.61M	325k

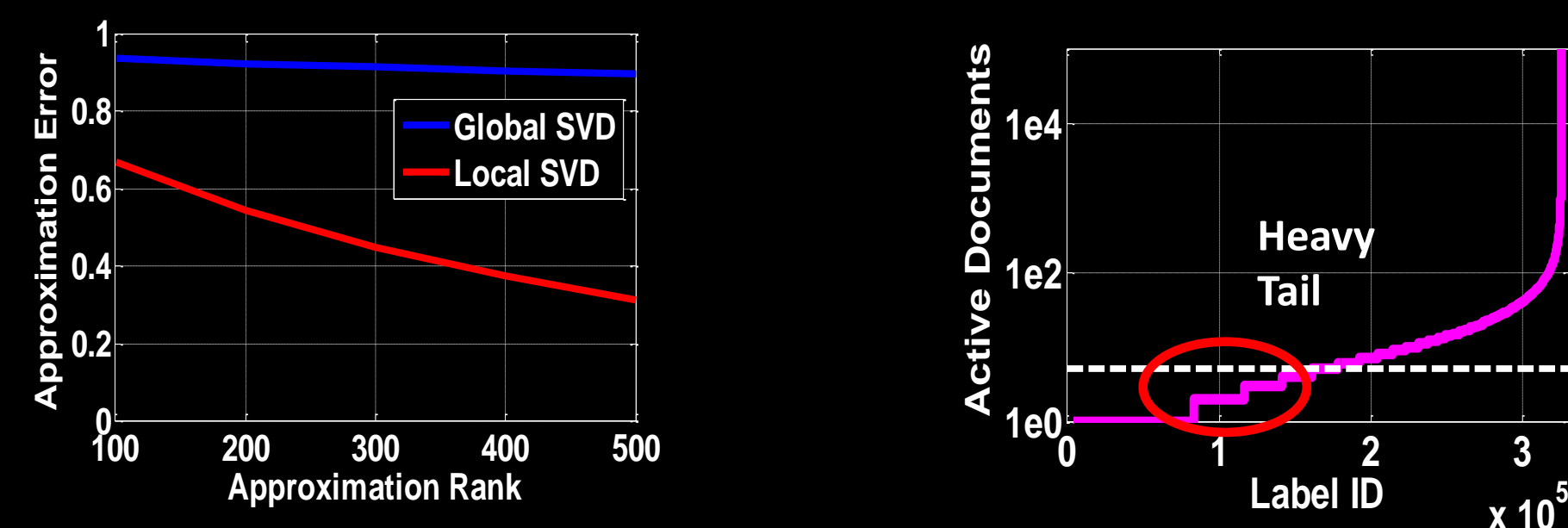


## Motivation

Embedding v/s Tree Techniques:



## Limitations of Embedding Methods

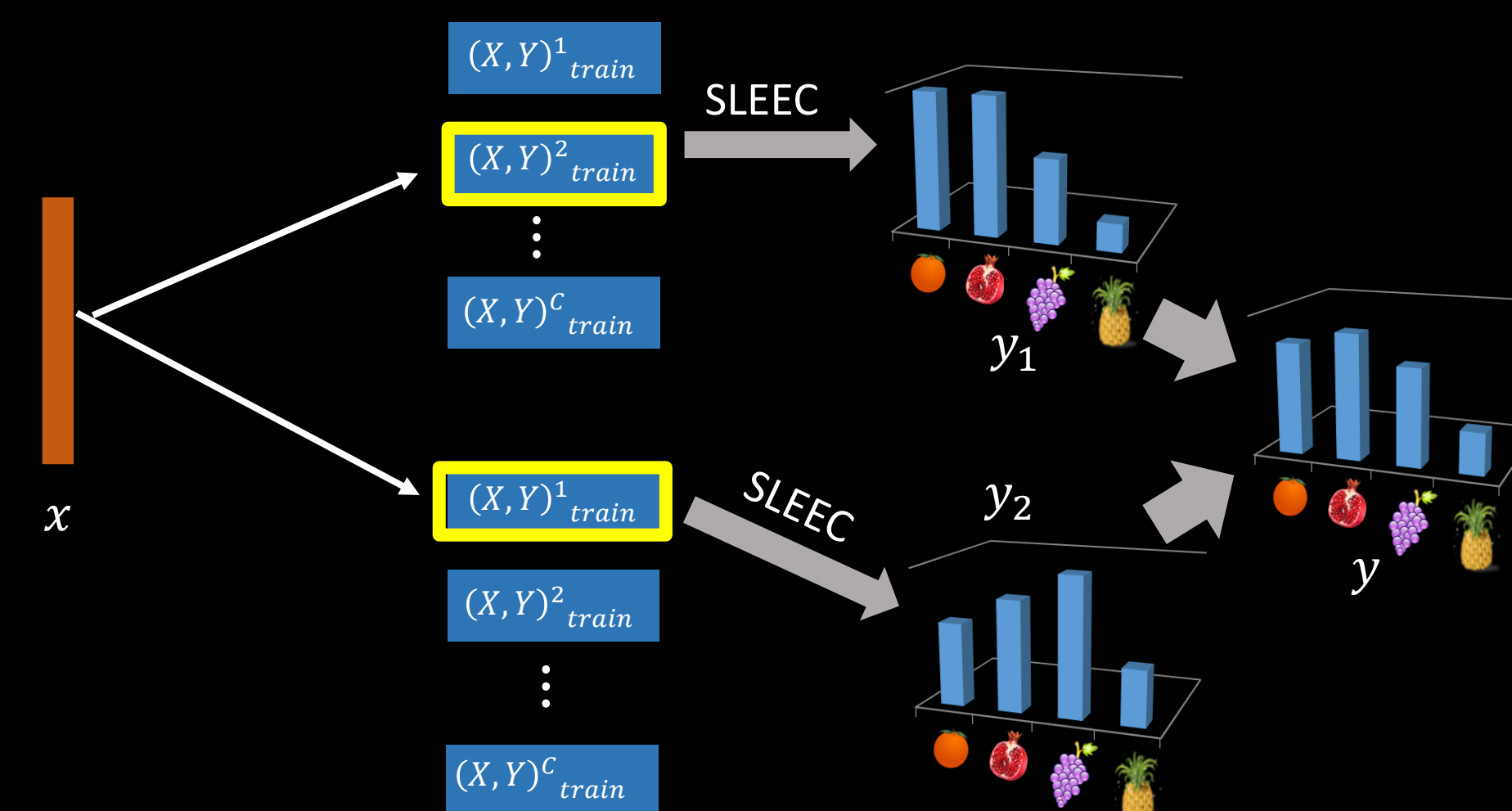


Key Insight:

- k-NN consistently performs very good in heavy tailed distributions.
- Practically infeasible: Test time is  $O(Nd')$

Our Solution : Neighborhood preserving low rank embeddings

## SLEEC : Scaling up with Multiple Learners



SLEEC : Key Features

- Can scale efficiently to datasets with a million labels
- Accuracy Improvements of upto
  - 35% over embedding methods
  - 6% over tree methods
- Theoretical Guarantees : Consistency Guarantees for Finite Sample

## Wikipedia-LSHTC

