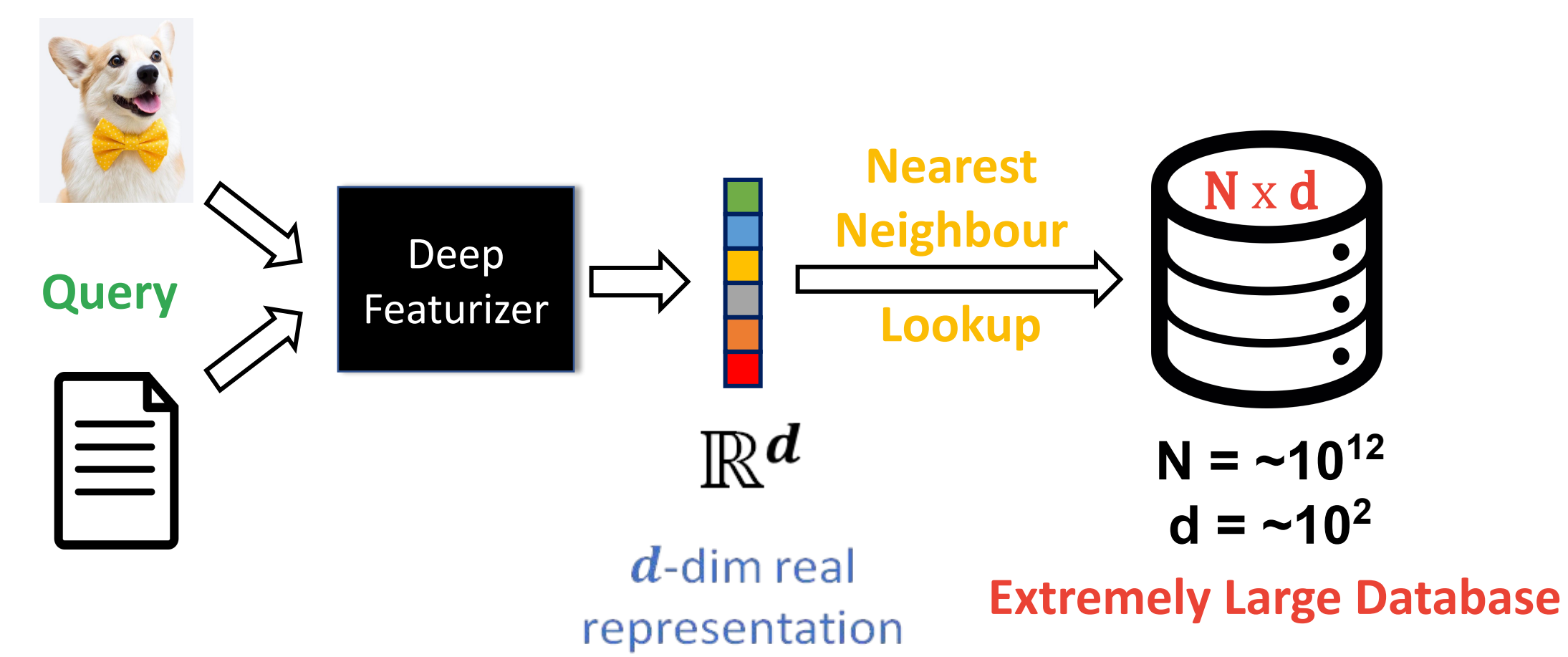# Matryoshka Representation Learning

Aditya Kusupati*, Gantavya Bhatt*, Aniket Rege*,
Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen,
Sham Kakade, Prateek Jain, Ali Farhadi

/RAIVNLab/MRL
bit.ly/mrl-paper

## Motivation: Query-based Retrieval

Query

Deep Featurizer

Nearest Neighbour Lookup

$\mathbb{R}^d$

$d$-dim real representation

N x d

N = ~$10^{12}$
d = ~$10^2$

Extremely Large Database

Applicable for large-scale classification with millions of labels
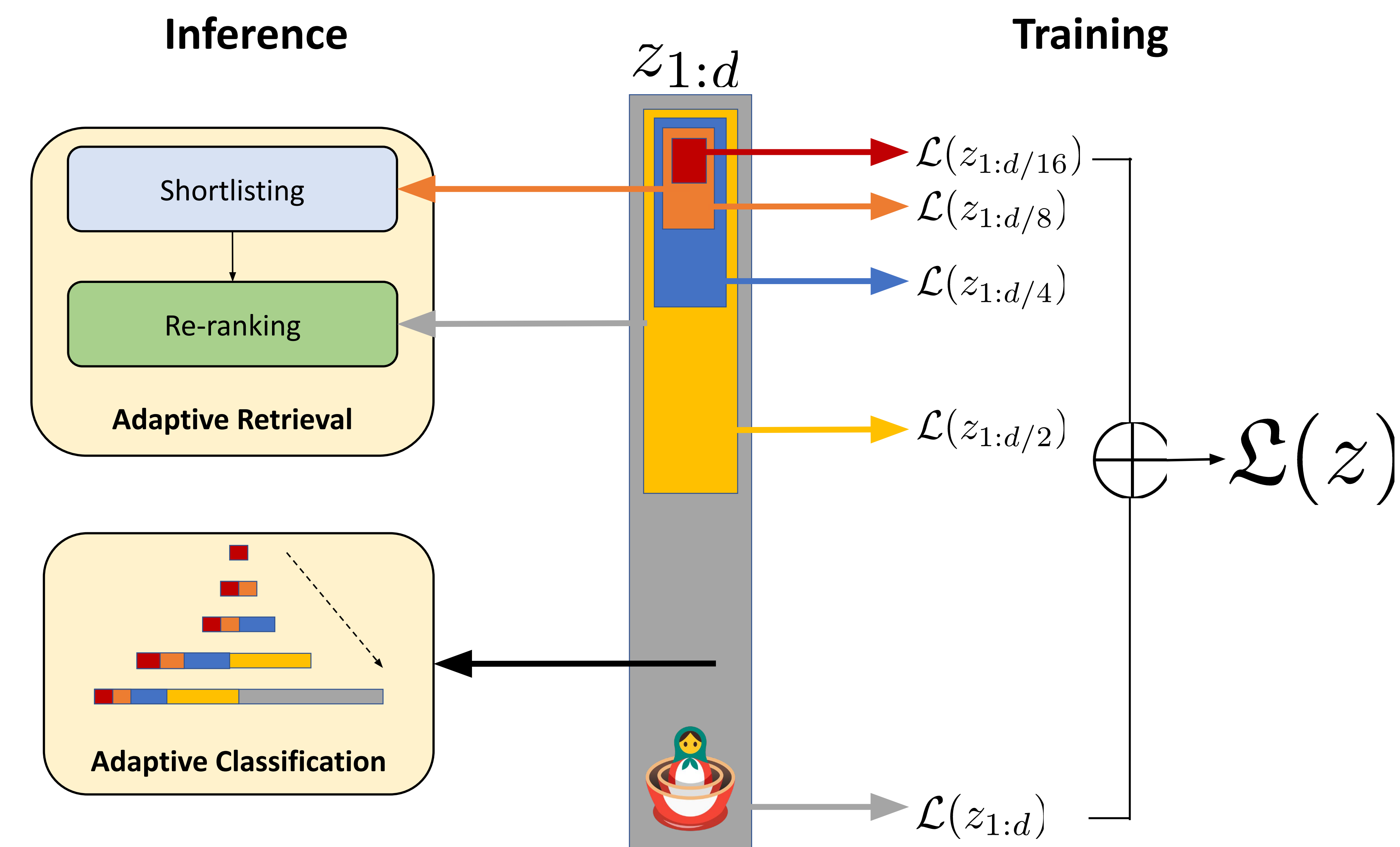
## Web-scale Challenges

- Extremely large databases – **100s of TB**
  - Linear dependence on representation size ($d$)
  - Embedding look-up much more expensive than featurization

- Require Approximate Nearest Neighbour Search (ANNS)
  - Post-hoc compressed index

- **Incapable** of Multi-Granularity
  - Use same *high-d* embedding for all tasks
  - Retrain a model for *low-d* based on deployment constraints
  - Eg: *2048-d* ResNet50 image representation for all tasks

## Adaptive Deployment – Goals

- **One representation** vector for all downstream tasks
  - No post-hoc compression or expensive feature selection
  - No retraining for specific resource constraints

- Accurate and efficient *low-d* embeddings
  - Baked within the *high-d* embedding – **Free**
  - **Reduced costs** for expensive & high-recall shortlisting
  - As **accurate** as independently trained counterparts

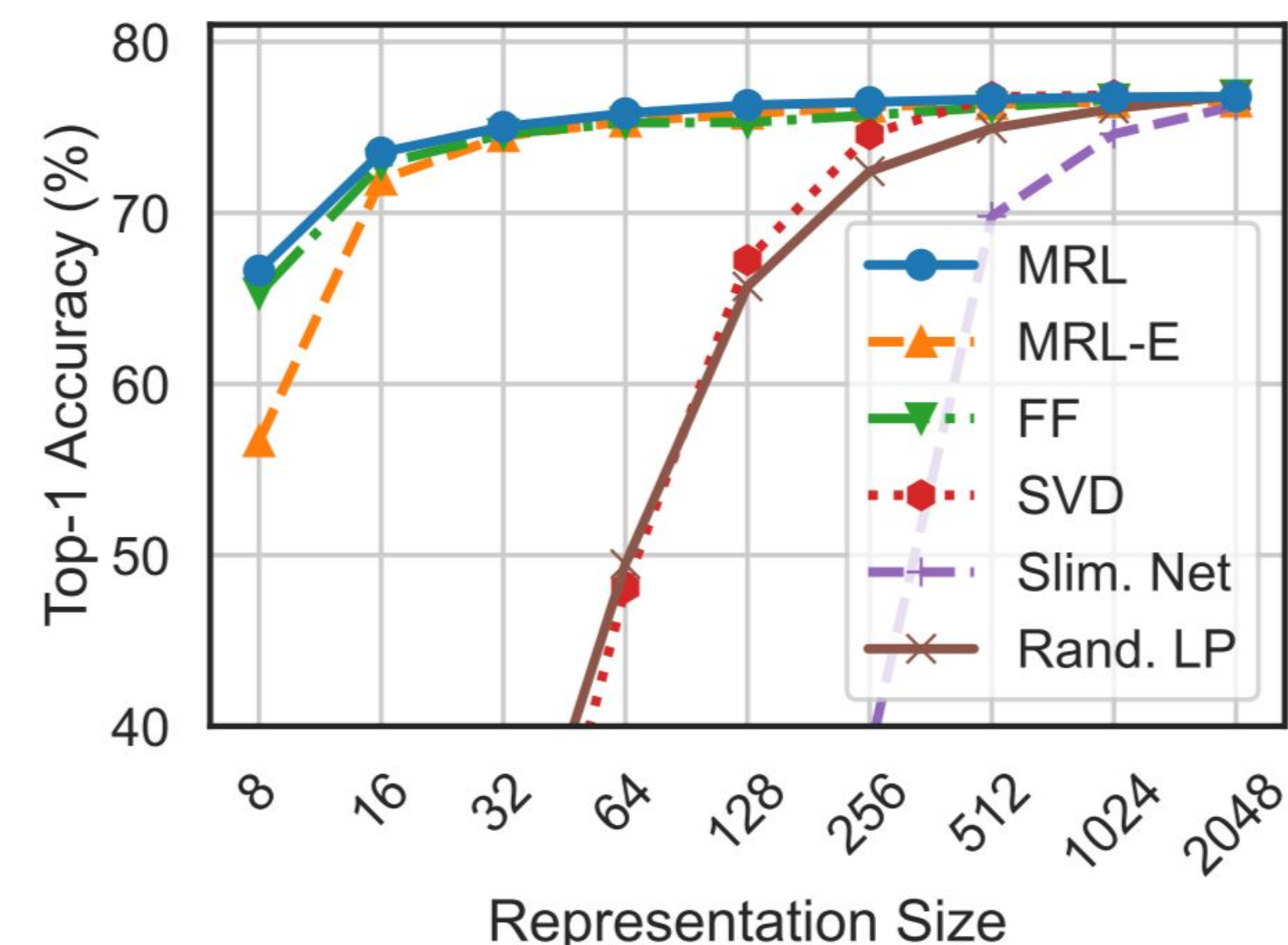- *High-d* embedding for **cheap** & precise re-ranking

## Matryoshka Representation Learning - MRL

- Solve the same learning task at **multiple granularities** ($log(d)$)

- Easily adaptable to any representation learning setup
  - Scale, modality and task agnostic – **1B images with ease**

- **First $k$ dims** form the required *low-d* embeddings
  - As accurate as retrained *low-d* counterparts

- Enable **adaptive** deployment
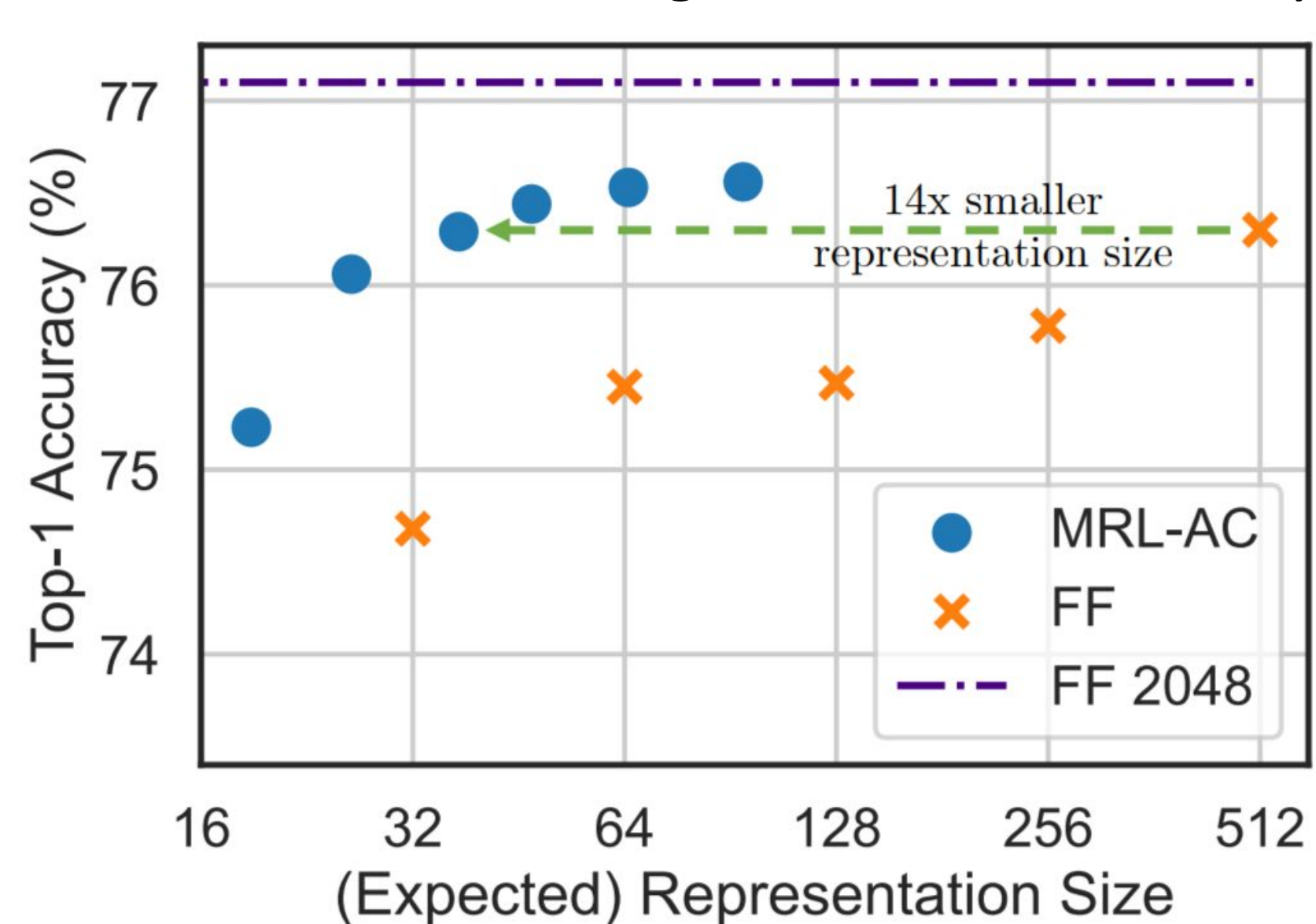  - Accurate large-scale classification & retrieval based on constraints

**Inference**

Shortlisting

Re-ranking

**Adaptive Retrieval**

**Adaptive Classification**

**Training**

$z_{1:d}$

$\mathcal{L}(z_{1:d/16})$
$\mathcal{L}(z_{1:d/8})$
$\mathcal{L}(z_{1:d/4})$
$\mathcal{L}(z_{1:d/2})$

$\bigoplus \mathfrak{L}(z)$

$\mathcal{L}(z_{1:d})$

## Classification Accuracy
### *ImageNet* OVA

- ResNet50: Same accuracy as independently trained *low-d* models (FF)



Top-1 Accuracy (%) vs Representation Size — MRL, MRL-E, FF, SVD, Slim. Net, Rand. LP

## Adaptive Classification
### *ImageNet-1K*

- ResNet50-MRL model with cascades
- 14x smaller embedding size for same accuracy



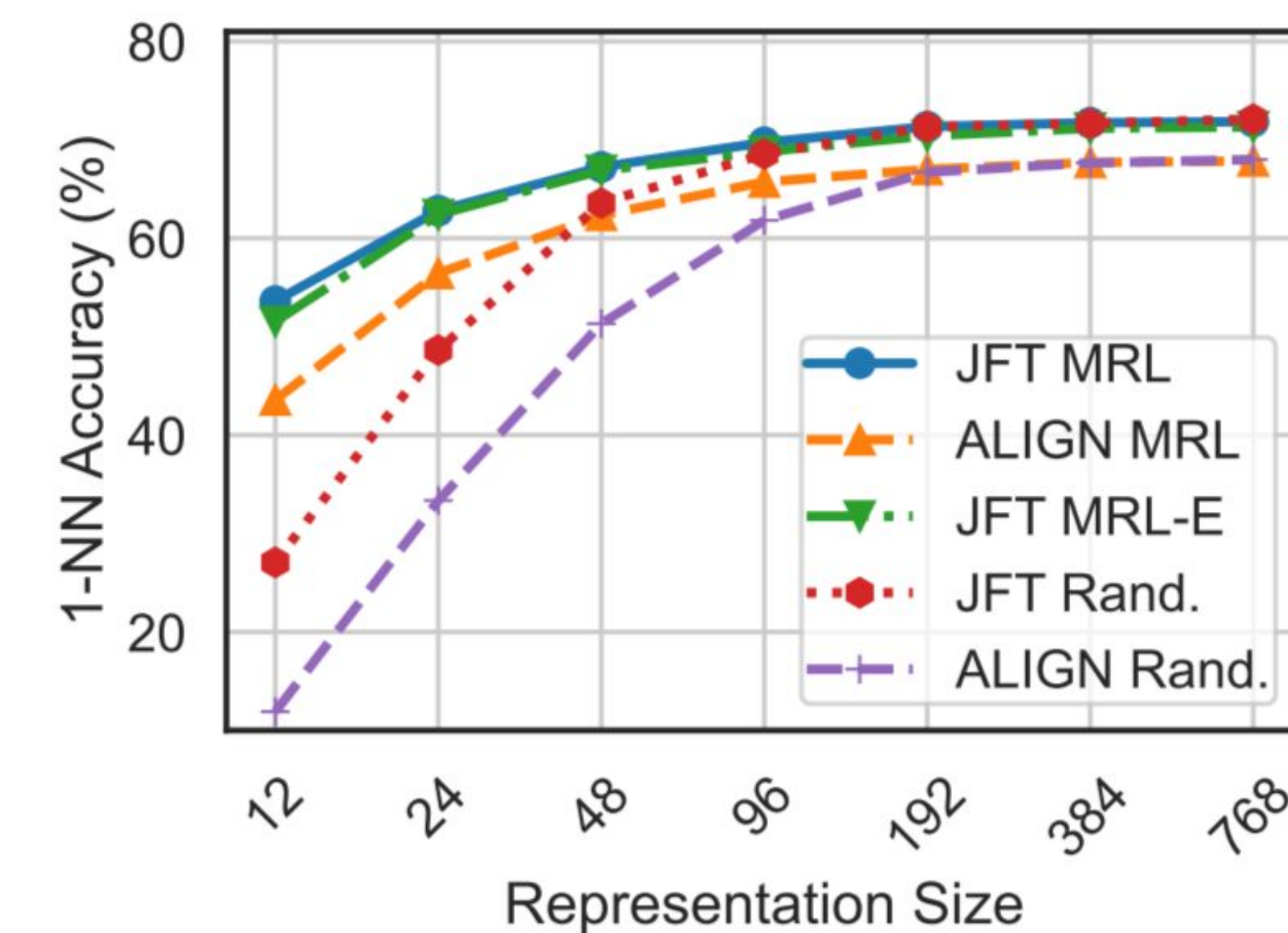Top-1 Accuracy (%) vs (Expected) Representation Size — MRL-AC, FF, FF 2048; *14x smaller representation size*

## Representation Quality

### *ImageNet* k-NN

- ResNet50 models trained on ImageNet-1K
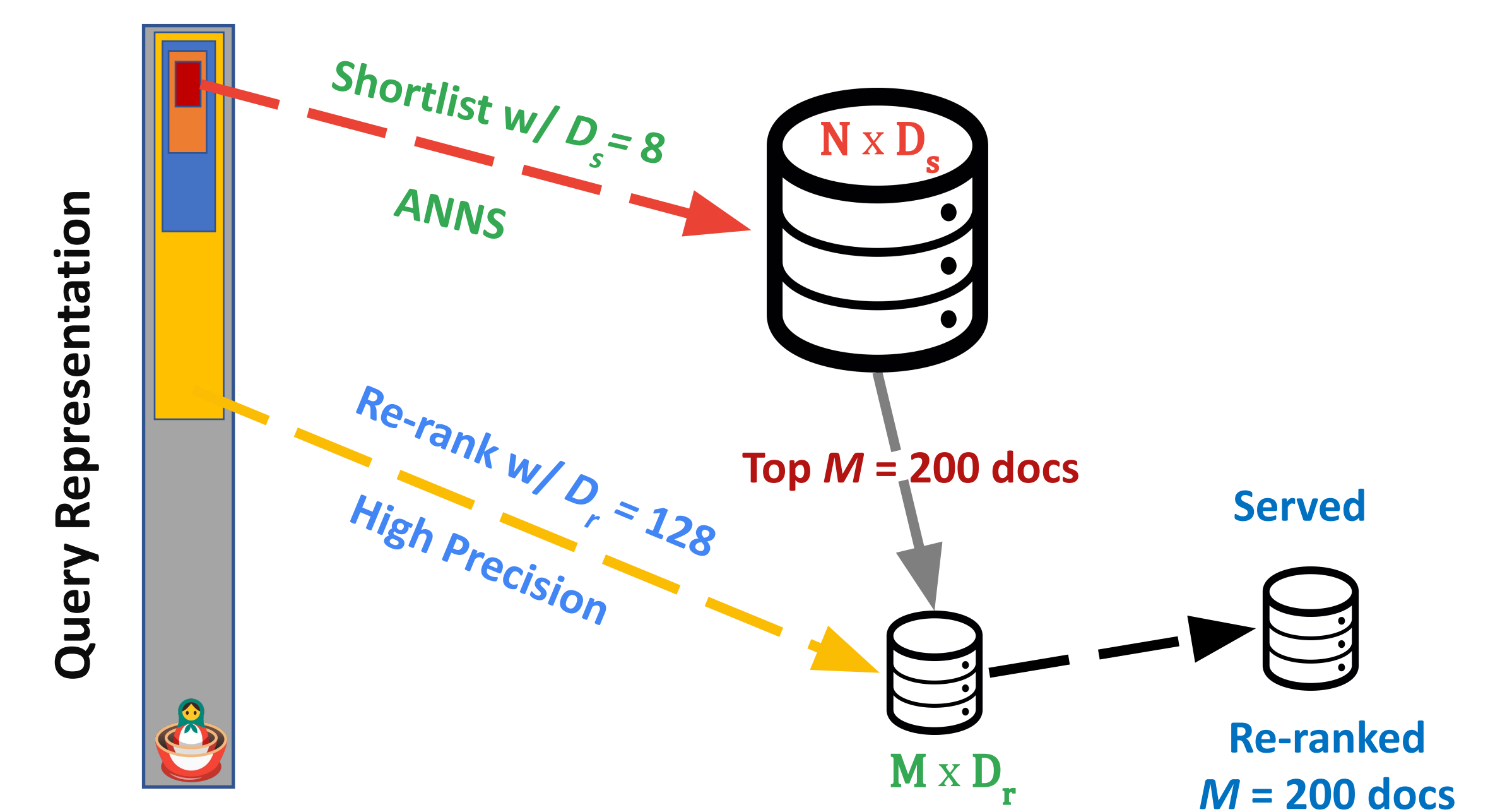- Other baselines fall off drastically at *low-dimensions*



1-NN Accuracy (%) vs Representation Size — MRL, MRL-E, FF, SVD, Slim. Net, Rand. FS

### *ViT+JFT & ALIGN* k-NN

- ViT-B/16 models trained on JFT-300M and ALIGN (V+L)
- Scales to **1B images** w/o accuracy drop



1-NN Accuracy (%) vs Representation Size — JFT MRL, ALIGN MRL, JFT MRL-E, JFT Rand., ALIGN Rand.

## Adaptive Retrieval

Query Representation

Shortlist w/ $D_s$ = 8
ANNS

N x $D_s$

Re-rank w/ $D_r$ = 128
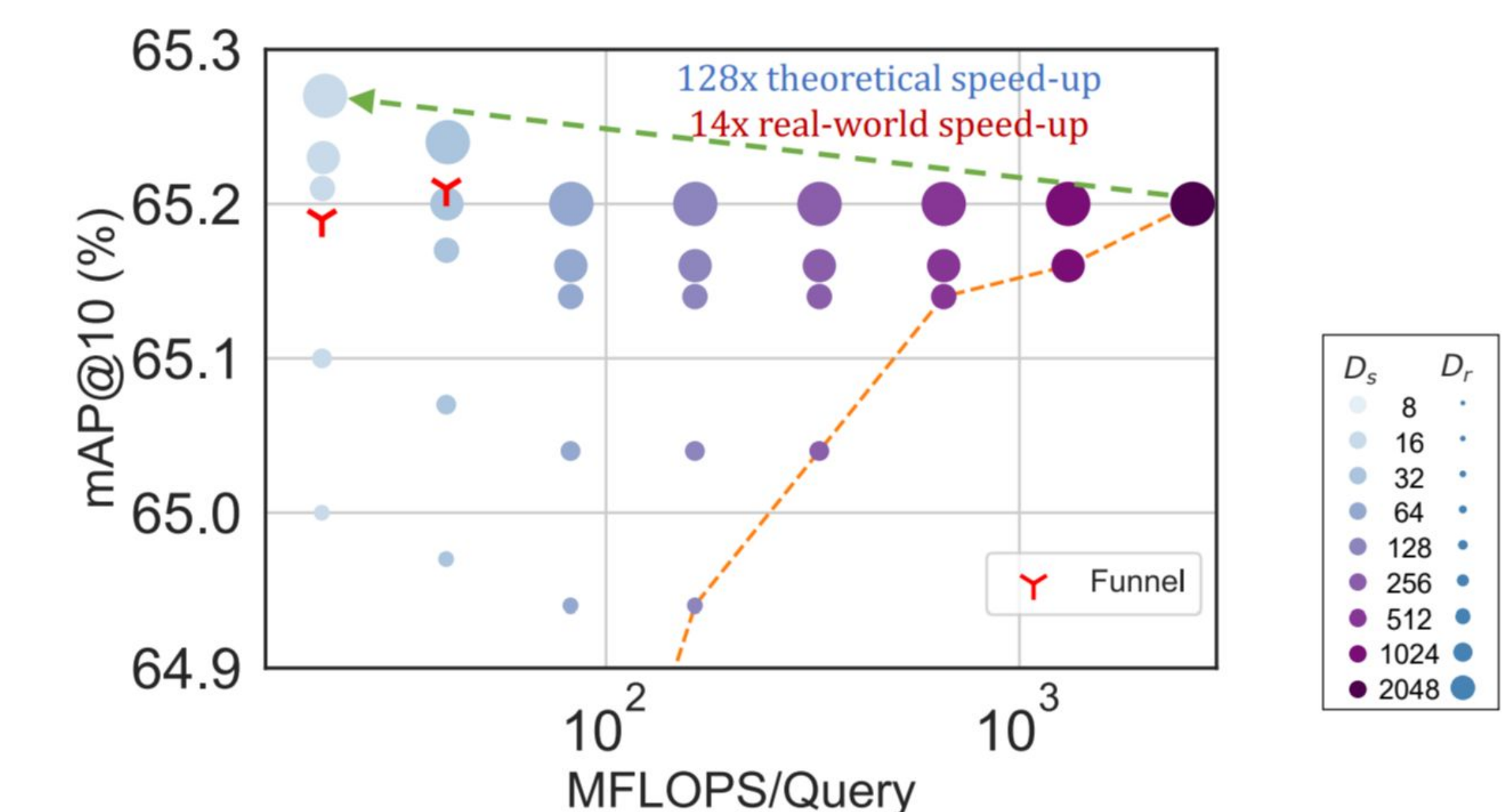High Precision

Top $M$ = 200 docs

M x $D_r$

Served

Re-ranked
$M$ = 200 docs

### *ImageNet-1K*

- 14x real-world speed-up for the best mAP@10
- All real-world implementations use HNSW for shortlisting



mAP@10 (%) vs MFLOPS/Query — 128x theoretical speed-up, 14x real-world speed-up; Funnel

### *ImageNet-4K* (Try it!)

- 6x real-world speed-up for the best mAP@10
- Funnel retrieval alleviates the need for optimal $D_s$ & $D_r$



mAP@10 (%) vs MFLOPS/Query — 6x real-world speed-up, 32x theoretical speed-up; Funnel