

# LeanVec: Searching vectors faster by making them fit

<sup>1</sup> DataStax  
<sup>2</sup> Intel Labs  
<sup>3</sup> Microsoft

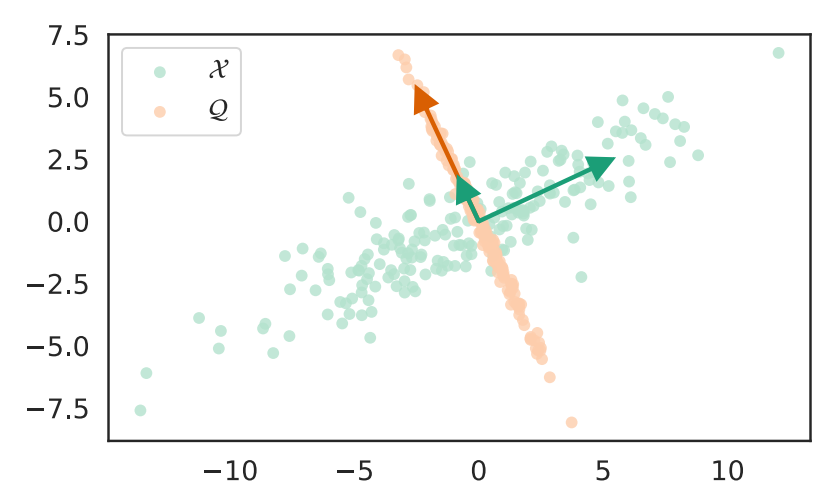
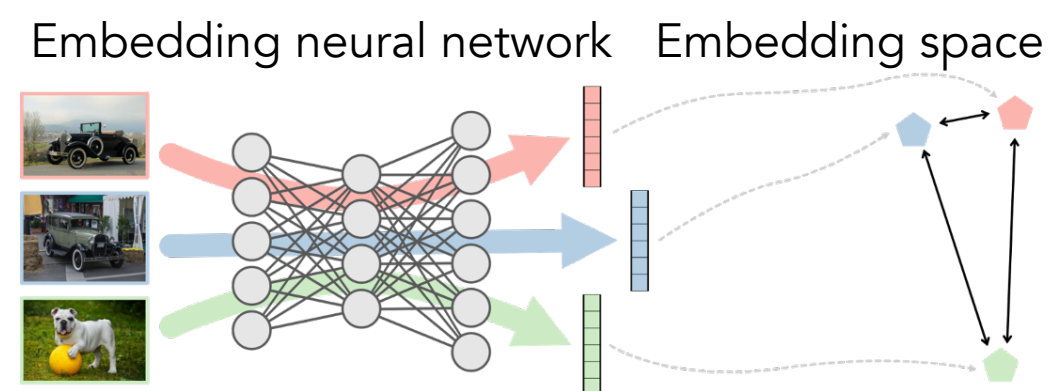
Mariano Tepper<sup>1\*</sup>, Ishwar Singh Bhati<sup>2</sup>, Cecilia Aguerrebere<sup>2</sup>, Mark Hildebrand<sup>3\*</sup>, Ted Willke<sup>1\*</sup>

*\* Work done while at Intel Labs*

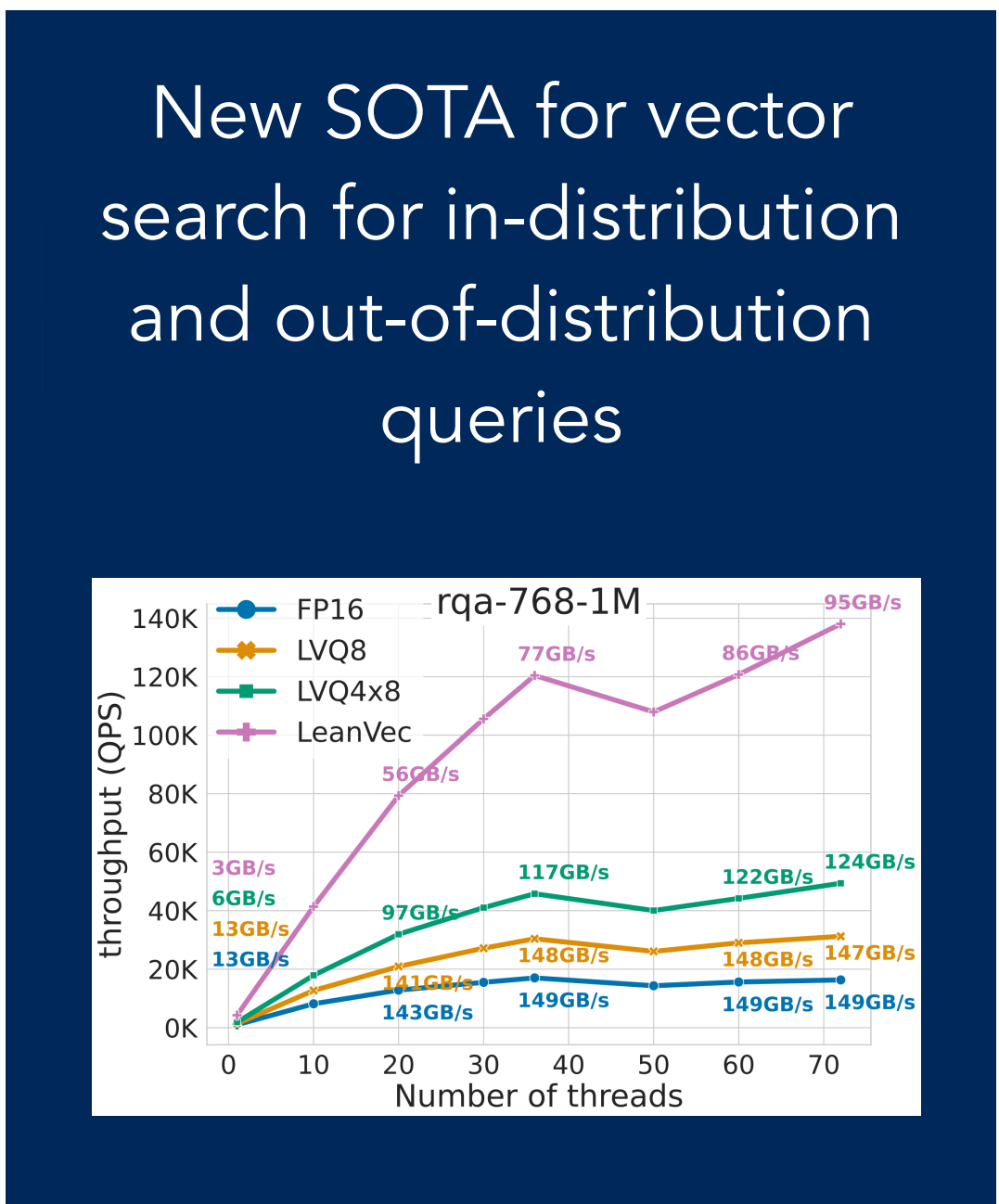
## Problem Statement

Finding semantically related vectors among millions or billions is crucial for applications like recommender systems, RAG, NLP, among many others.

Graph-based indices are SOTA, but random access patterns and increasing vector dimensionality limit their potential.



Considering the distribution of queries is key for dimensionality reduction



## LeanVec

$$\min_{\mathbf{A}, \mathbf{B} \in \text{St}(D, d)} \left\| \mathbf{Q}^\top \mathbf{A}^\top \mathbf{B} \mathbf{X} - \mathbf{Q}^\top \mathbf{X} \right\|_F^2$$

$\mathbf{A}, \mathbf{B}$  row-orthonormal matrices  
 $\mathbf{Q}$  matrix with query vectors  
 $\mathbf{X}$  matrix with database vectors

We propose two solvers:

Frank-Wolfe block-coordinate descent	Eigen-search
Convex relaxation of the orthonormality constraints on $\mathbf{A}, \mathbf{B}$	Finds the best convex combination of $\mathbf{Q}^\top \mathbf{Q}$ and $\mathbf{X}^\top \mathbf{X}$
✓ Theoretical guarantees	✓ Fast algorithm

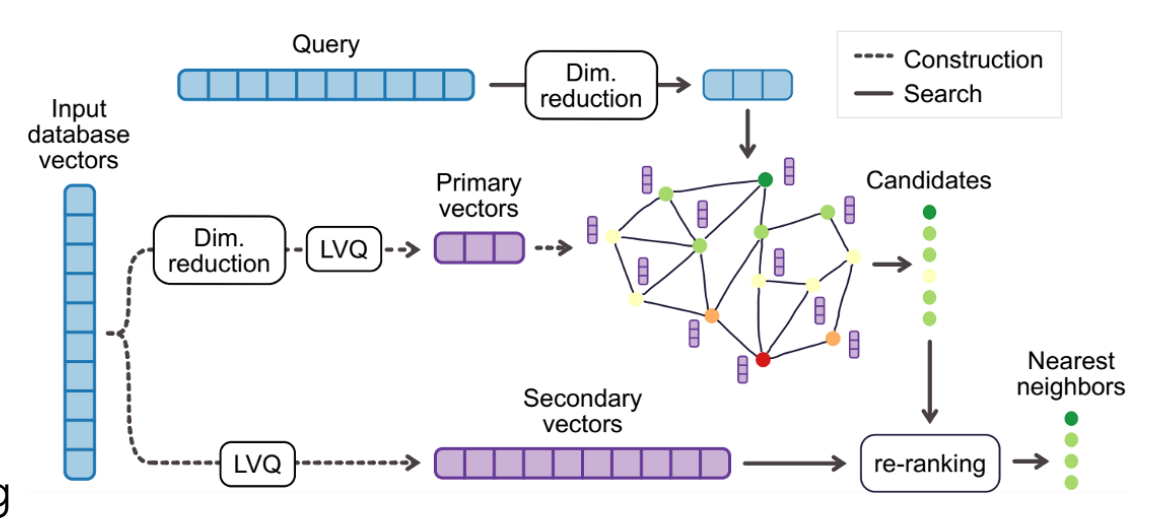
### Using LeanVec for graph search

- Two level for speed and accuracy

$$\langle \mathbf{q}, \mathbf{x} \rangle \approx \langle \mathbf{A} \mathbf{q}, \text{quant}(\mathbf{B} \mathbf{x}) \rangle$$

Computed at search.  
Linear approx for speed

Precomputed during indexing

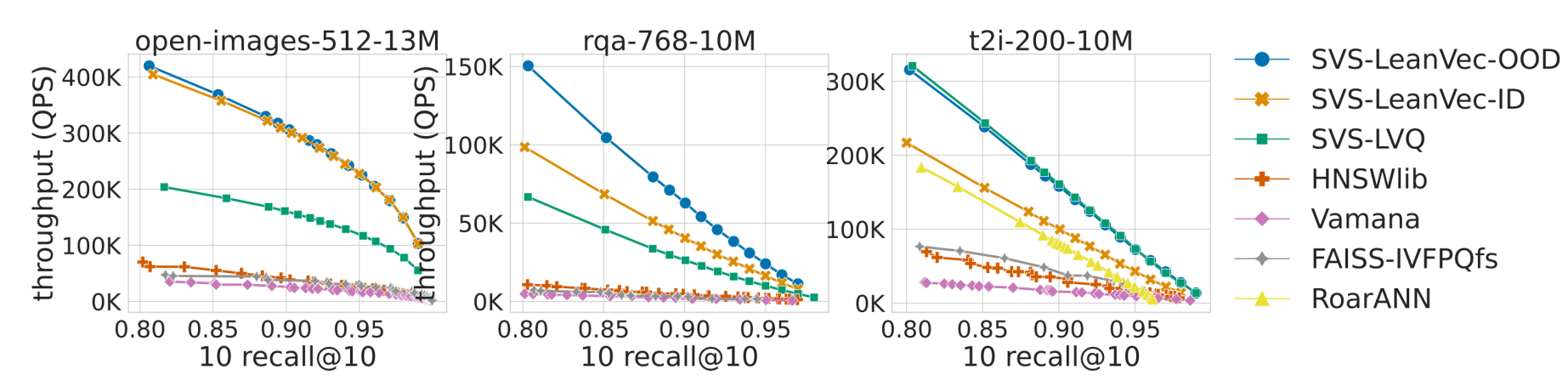


## Contributions

- High-quality graphs up to 4.9x faster than SOTA
- LeanVec-ID (query distr. agnostic) up to 3.6x faster search
- LeanVec-OOD, new optimal projection with two fast optimization algorithms, up to 5.4x faster search
- Integrated into Scalable Vector Search for reproducibility

## Results

### Faster search



### Faster construction

