

# Large Language Models: Vector embeddings

Nathanaël Fijałkow  
CNRS, LaBRI, Bordeaux



**LaBRI**

université  
de BORDEAUX

# VECTOR EMBEDDINGS

A ***vector embedding*** is a model which turns inputs into dense vectors, in a way that preserves the information the input contains

This can be used for many things, here are two examples:

# EXAMPLE 1: RECOMMENDATION SYSTEMS

We have a dataset of stuff:

- Texts
- Images
- Videos
- ...

and we would like to be able to retrieve *similar* stuff

## EXAMPLE 2: DEDUPLICATION

We have a dataset of stuff (typically for a training dataset) and we would like to remove (approximate) duplicates

# WHAT IS A TEXT EMBEDDER?

A Transformer model trained on

- **Masked language modeling:** taking a sentence, we randomly mask 15% of the words in the input then run the entire masked sentence through the model and predict the masked words
- **Next sentence prediction:** the model concatenates two masked sentences as inputs during pretraining. Sometimes they correspond to sentences that were next to each other in the original text, sometimes not. The model then has to predict if the two sentences were following each other or not

# EXAMPLES OF TEXT EMBEDDERS

- BERT: most popular language model on Hugging Face.
- SBERT: Also known as sentence BERT and sentence transformers
- DistilBERT: A lightweight BERT variant
- RoBERTa: short for “robustly optimized BERT pre training approach”, RoBERTa refined the BERT training procedure to optimize its performance

# THE NOTIONS OF SIMILARITY

There are many, but two main options:

- **Cosine similarity (IP = Inner Product)**, which becomes dot product if the vectors are normalized
- **Euclidean distance (L2)**

# SIMILARITY IS DIFFERENT FROM HASHING

Both **hashing** and **similarity search** are about mapping objects to dense vectors, but:

- Hashing is about avoiding collisions for different objects
- Similarity search is about enforcing collisions for similar objects



# IN PRACTICE

A very efficient, well maintained, and SOTA library is

FAISS: Facebook AI Similarity Search

# BASILINE: EXHAUSTIVE SEARCH

We have a dataset of  $N$  vectors and a query vector

**Exhaustive search** computes similarity scores between the query and each vector in the dataset:

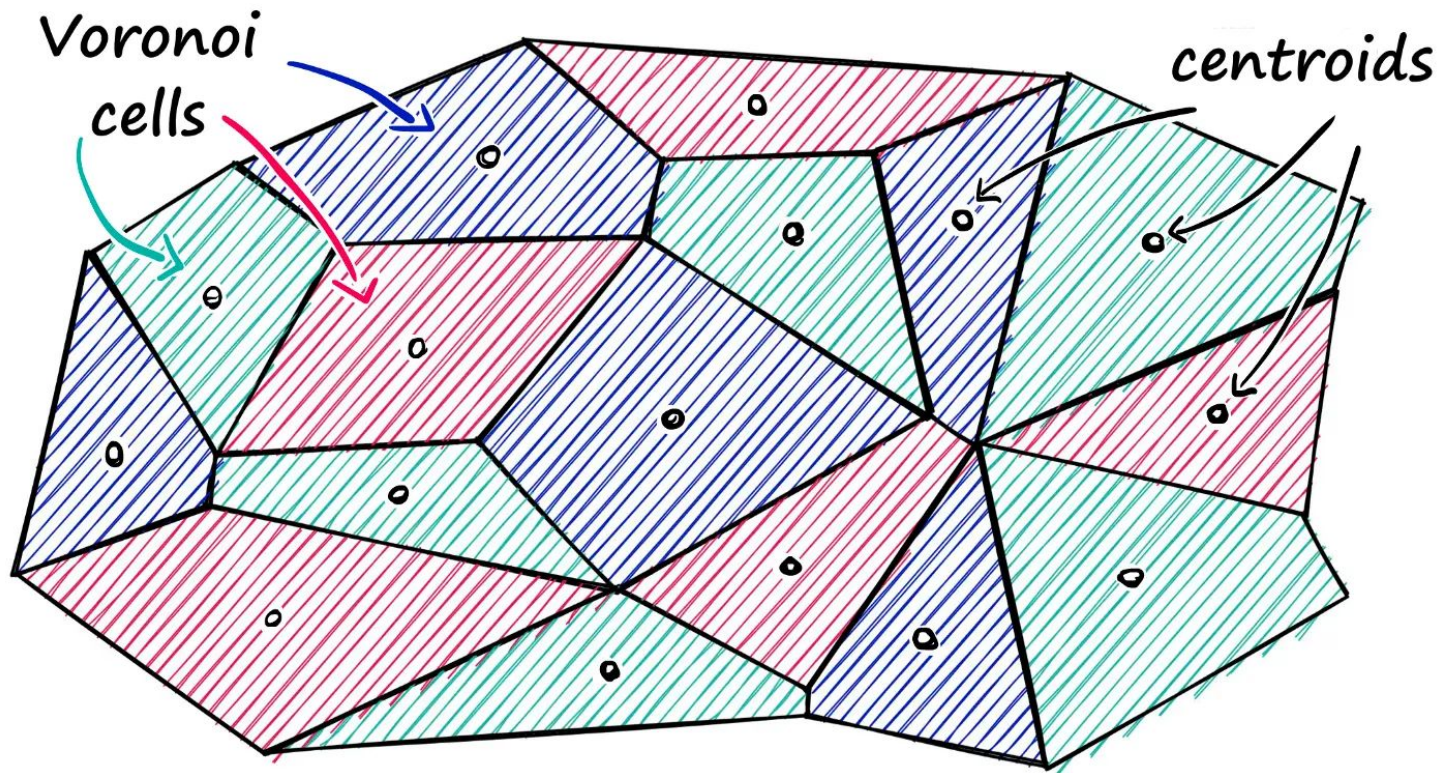
- Very accurate
- But slow and heavy on memory

*IndexFlatL2* and *IndexFlatIP* implement exhaustive search

# FIRST OPTIMIZATION: INVERTED FILE INDEX

*Inverted File Index* (IVF) reduces search space through clustering:

# FIRST OPTIMIZATION: INVERTED FILE INDEX



# FIRST OPTIMIZATION: INVERTED FILE INDEX

**Key idea:** we compare the query vector first only to centroids (*nlist*), pick a number of cells (*nprobes*), and then to all vectors in the selected cells

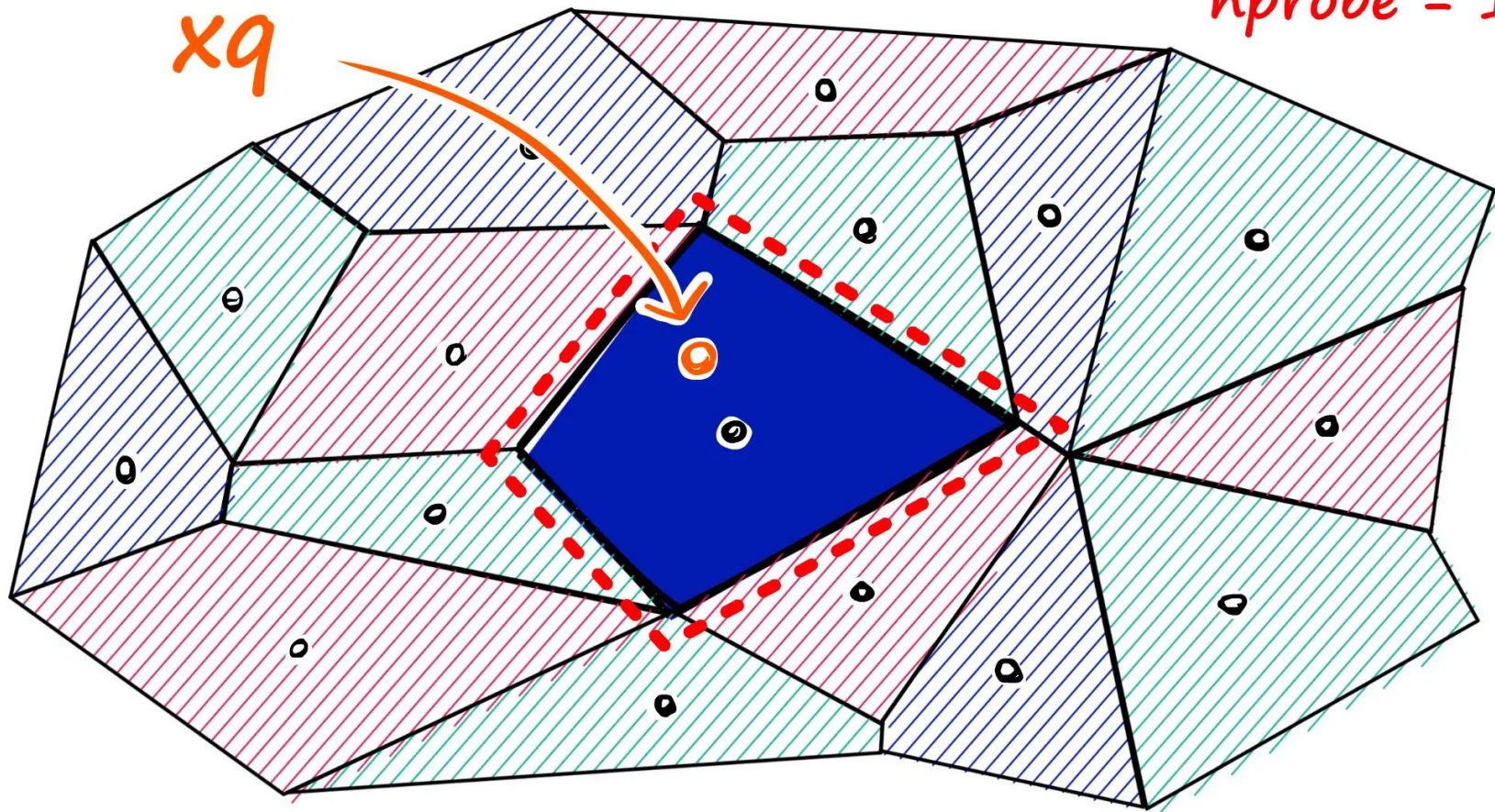
**Important:** gives approximate results, but a lot faster and less memory. Two parameters:

- Number of centroids (*nlist*)
- Number of cells considered (*nprobes*)

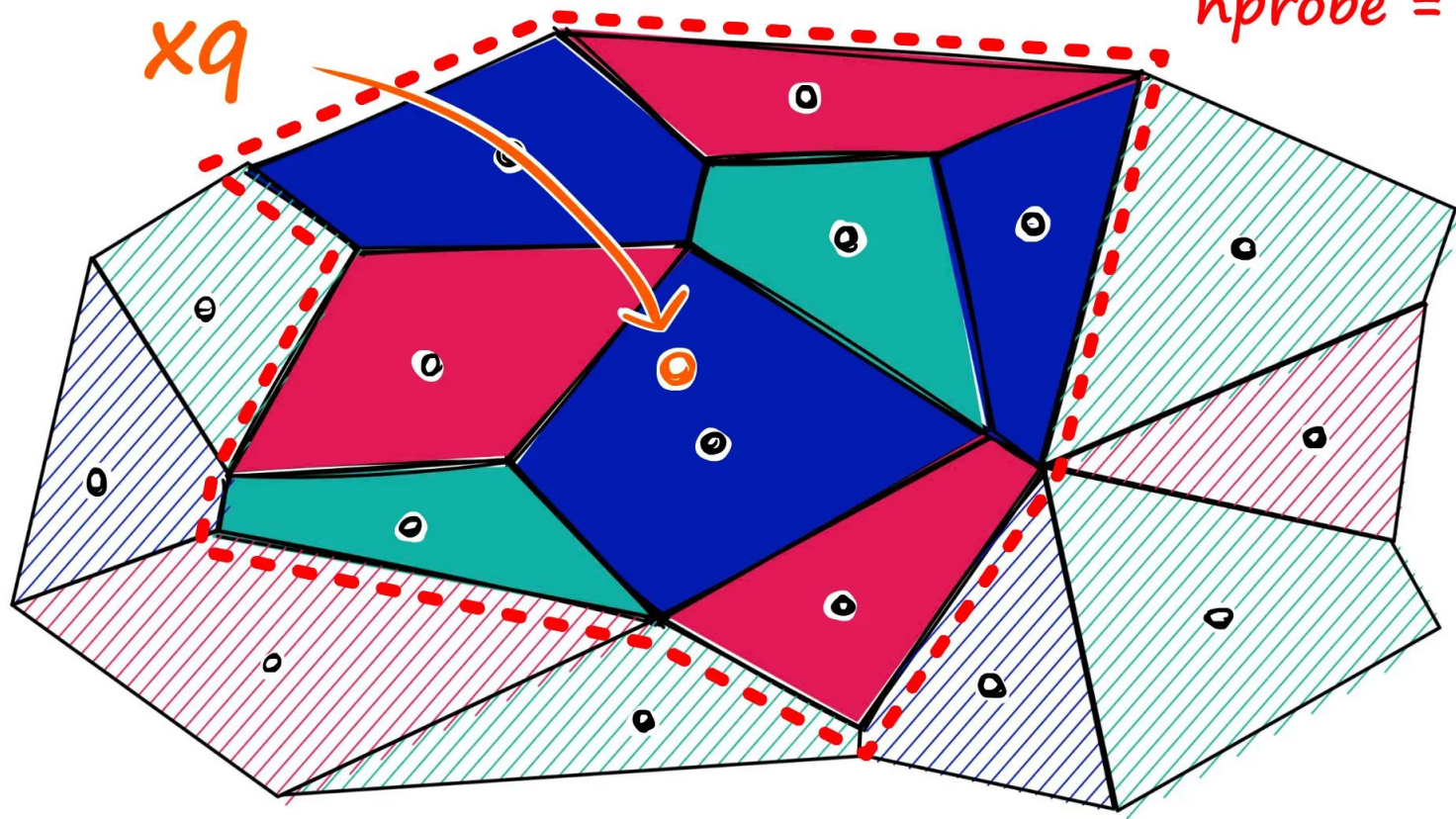
*IndexIVFFlat* implements the index partition



search scope  
 $n_{\text{probe}} = 1$

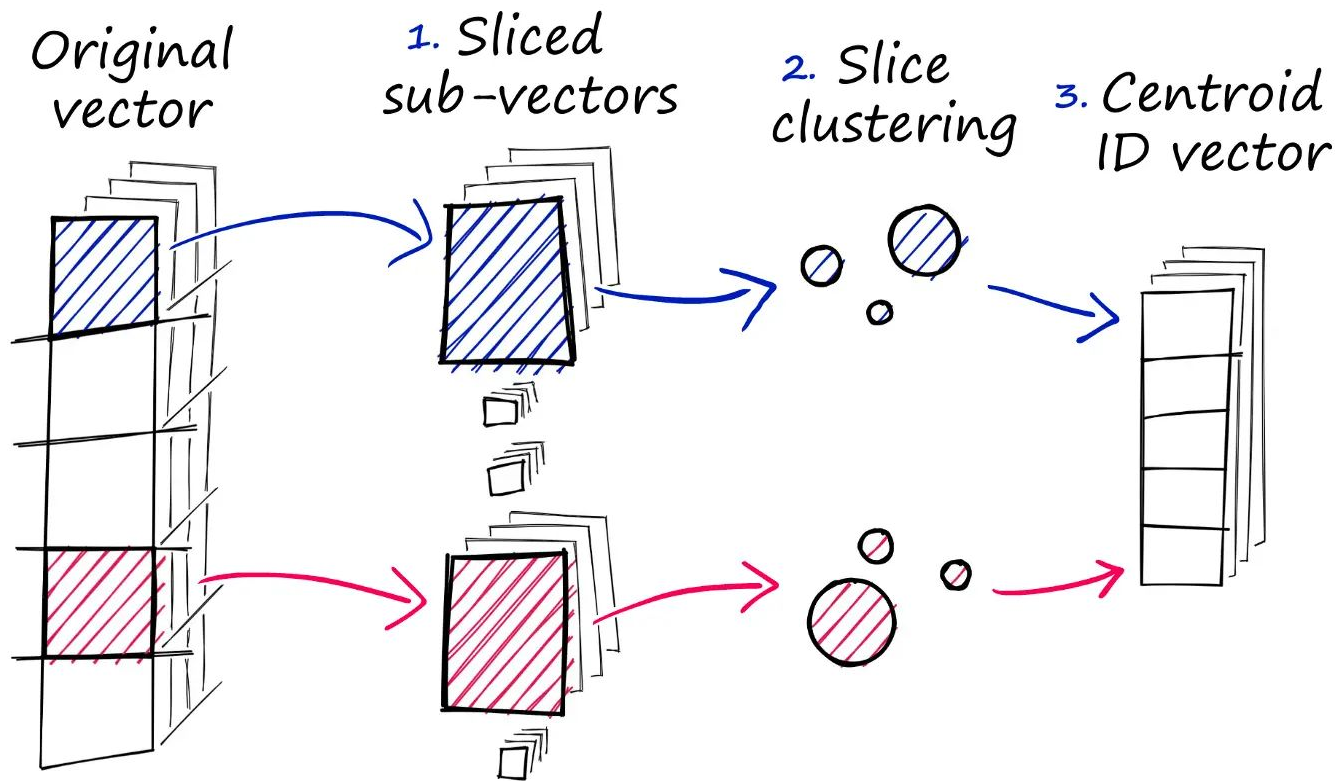


search scope  
 $n_{\text{probe}} = 8$





# SECOND OPTIMIZATION: PRODUCT QUANTIZATION





# SECOND OPTIMIZATION: PRODUCT QUANTIZATION

**Objective:** approximating similarity score computation

**Key idea:**

- We split a vector in dimension  $d$  into  $m$  “subvectors” of dimension  $d/m$
- In each subspace of dimension  $m$  we run a clustering algorithm
- We represent a subvector by the (id of the) closest centroid

*IndexPQ* implements both optimizations

# TWO OPTIMIZATIONS COMBINED

## **Key idea:**

- First, an IVF index to reduce to a small number of cells
- Then, a PQ index in a much smaller set of vectors

*IndexIVFPQ* implements both optimizations

# LOCALITY SENSITIVE HASHING

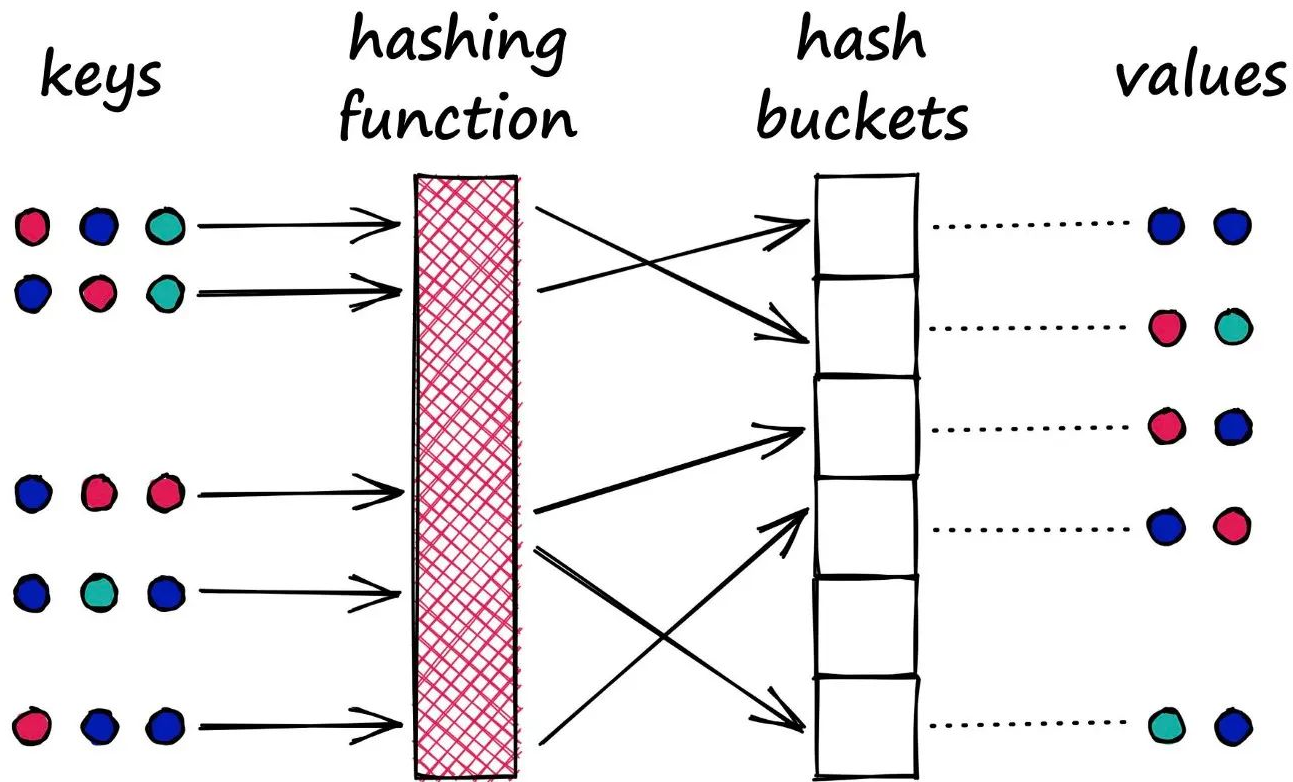
---

# TWO STEP PROCESS

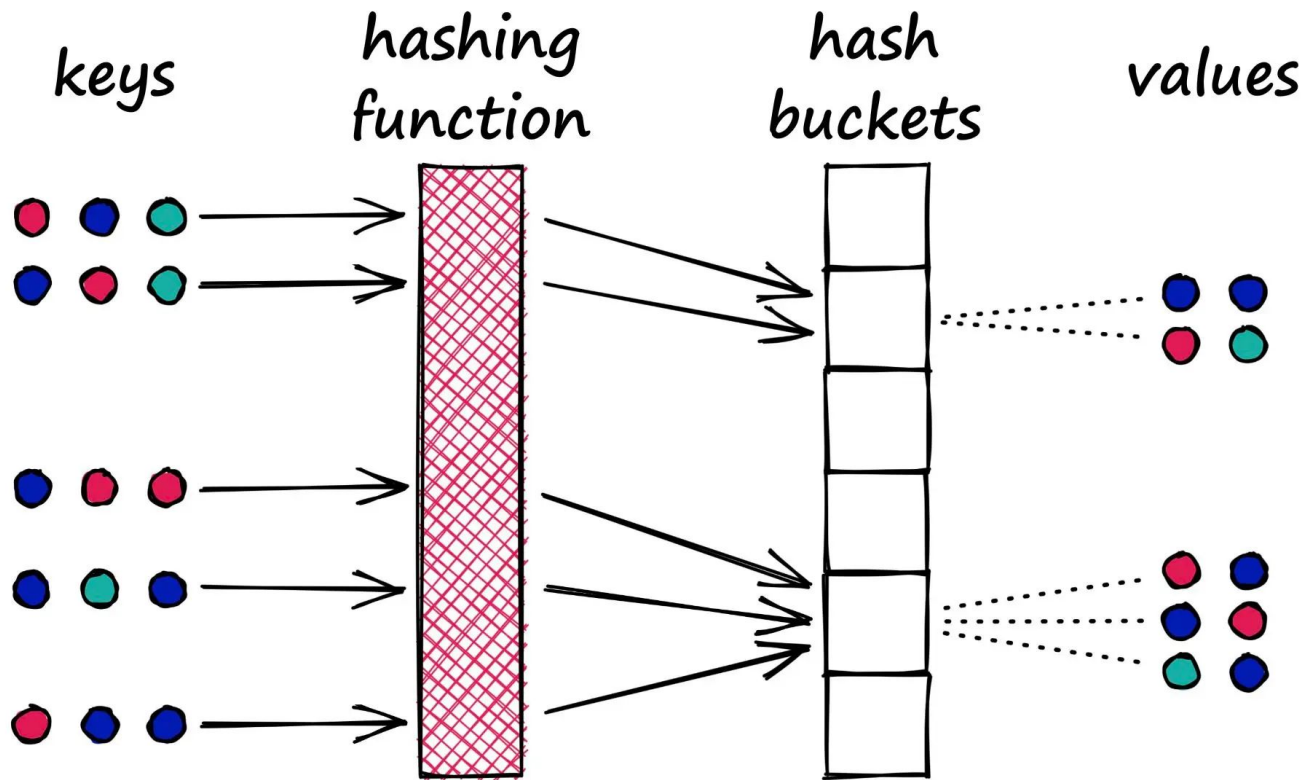
Step 1: identify candidate pairs

Step 2: exhaustive search on candidates

# HASHING: MINIMISING COLLISIONS



# LSH: MAXIMISING COLLISIONS FOR SIMILAR INPUTS

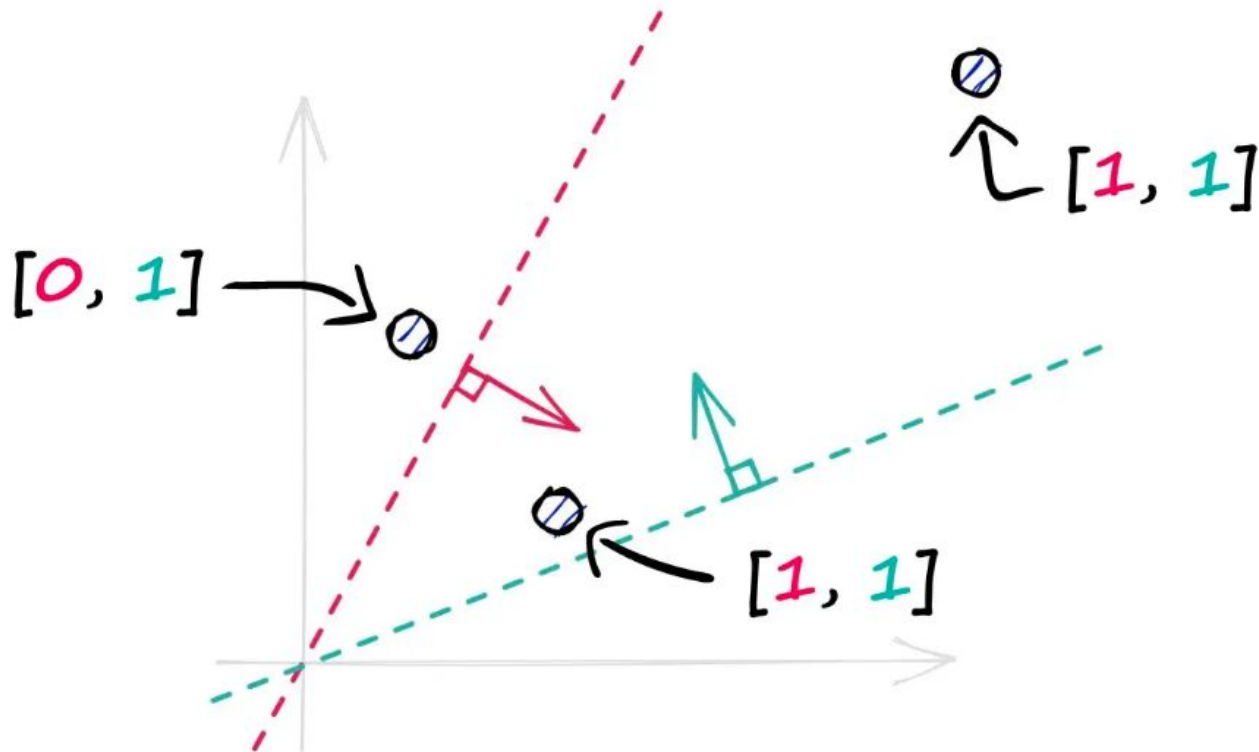


# LSH IS DIFFERENT FROM HASHING

There are two main approaches:

- Random projections
- Shingling, MinHashing, Banding

# LSH WITH RANDOM PROJECTIONS



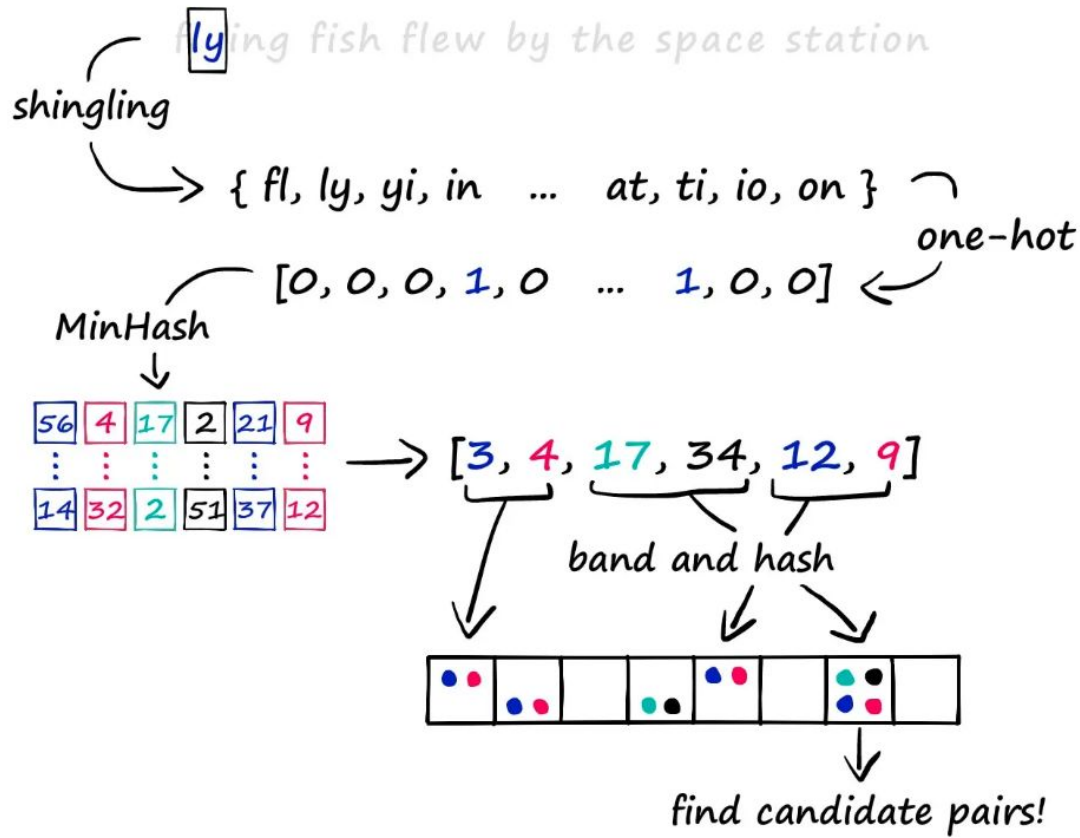


# LSH WITH RANDOM PROJECTIONS

## Key idea:

- We define *nbits* random hyperplanes in embedding space, each represented by the normal vector
- We represent a vector by the boolean vector expressing on which side of each hyperplane the vector falls on, determined by the sign of dot-product
- We compute Hamming distance between boolean vectors to determine similarity

# LSH WITH SHINGLING, MINHASHING, BANDING



# A GREAT REFERENCE

[Faiss: The Missing Manual](#)