



miniCOIL

A Sparse Neural Retriever

Haystack
2025





Vector Search

An essential part of the AI Transformation

INDUSTRIES

HR-Tech Ad-Tech Online Dating Gig Economy

E-Commerce Law-Tech Fashion Med-Tech

Ed-Tech Media & News Biometrics

Agriculture Manufacturing Streaming Services

Marketplaces Anti-fraud



NN Encoders + Vector Database



SEARCH SYSTEMS



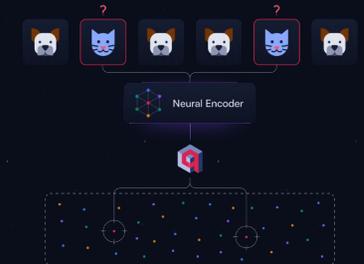
RAG / INFORMATION ASSISTANTS



RECOMMENDATIONS



ANOMALY DETECTION







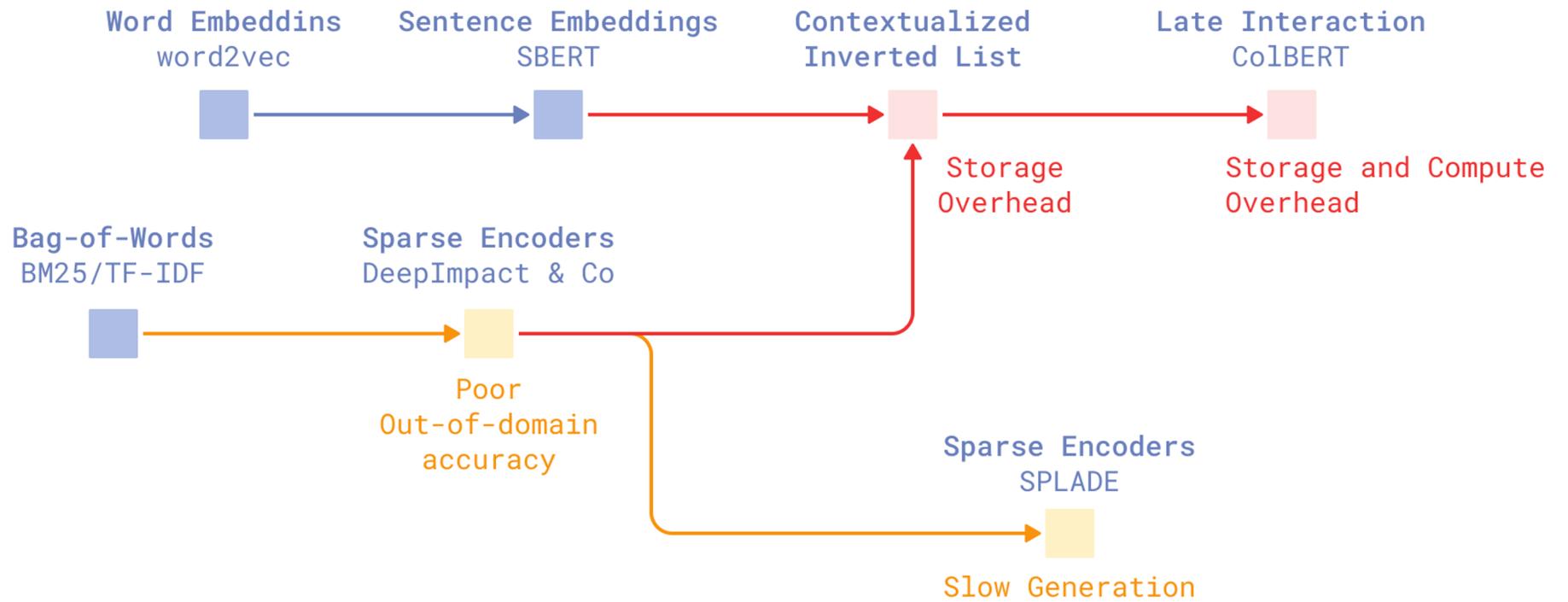
You shall know a word

by the company it keeps



Agenda

1. Motivation & BM25 recap
2. miniCOIL concept & 4-D vectors
3. Architecture & training pipeline
4. Advantages & benchmarks
5. Demo + how to extend
6. Q&A



$$\text{score}(D, Q) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \left[\frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdl}}\right)} + \delta \right]$$





For example:

$$BM25 = \sum_{t \in q} \log \left[\frac{N}{df(t)} \right] \cdot \frac{(k_1 + 1) \cdot tf(t, d)}{k_1 \cdot [(1 - b) + b \cdot \frac{dl(d)}{dl_{avg}}] + tf(t, d)}$$

- k_1, b - parameters
- $dl(d)$ - length of document d
- dl_{avg} - average document length

Relevance Score = <some formula based on word statistics>

Query: What is Qdrant

Qdrant is a vector database

1 1 1 1 1

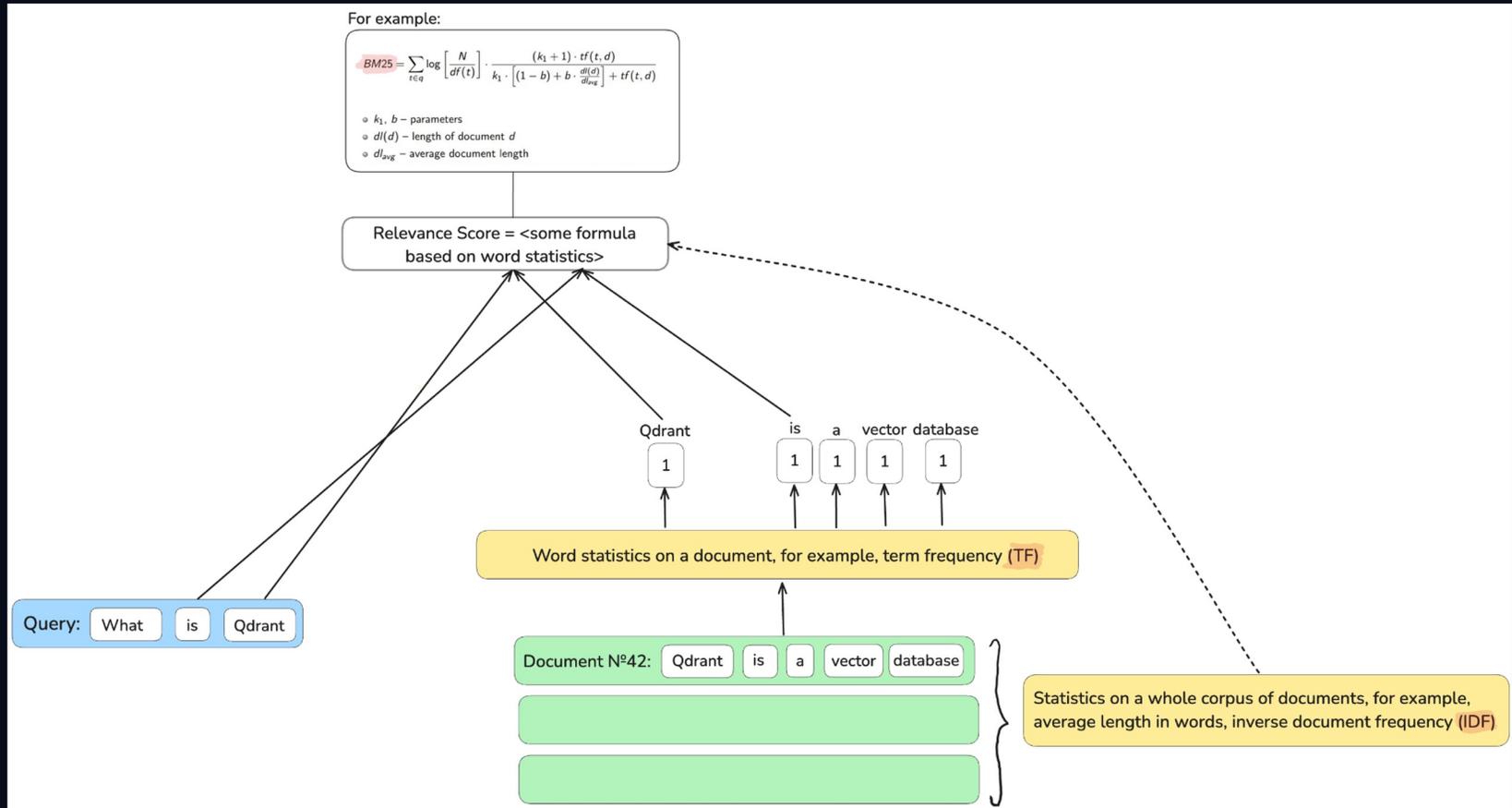
Word statistics on a document, for example, term frequency (TF)

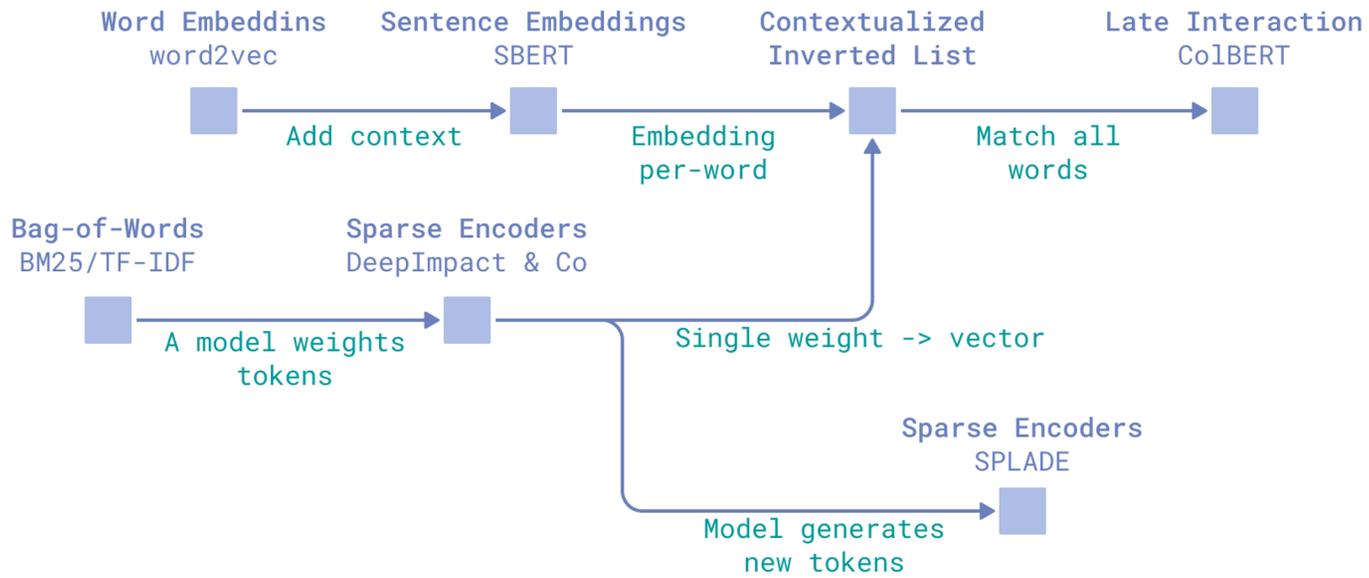
Document N°42: Qdrant is a vector database

[Redacted]

[Redacted]

Statistics on a whole corpus of documents, for example, average length in words, inverse document frequency (IDF)







miniCOIL

BM25 formula

×

4-D context vectors



Homographs

“fruit bat” ≈ “baseball bat”

“calculator” ≈ “calculating”



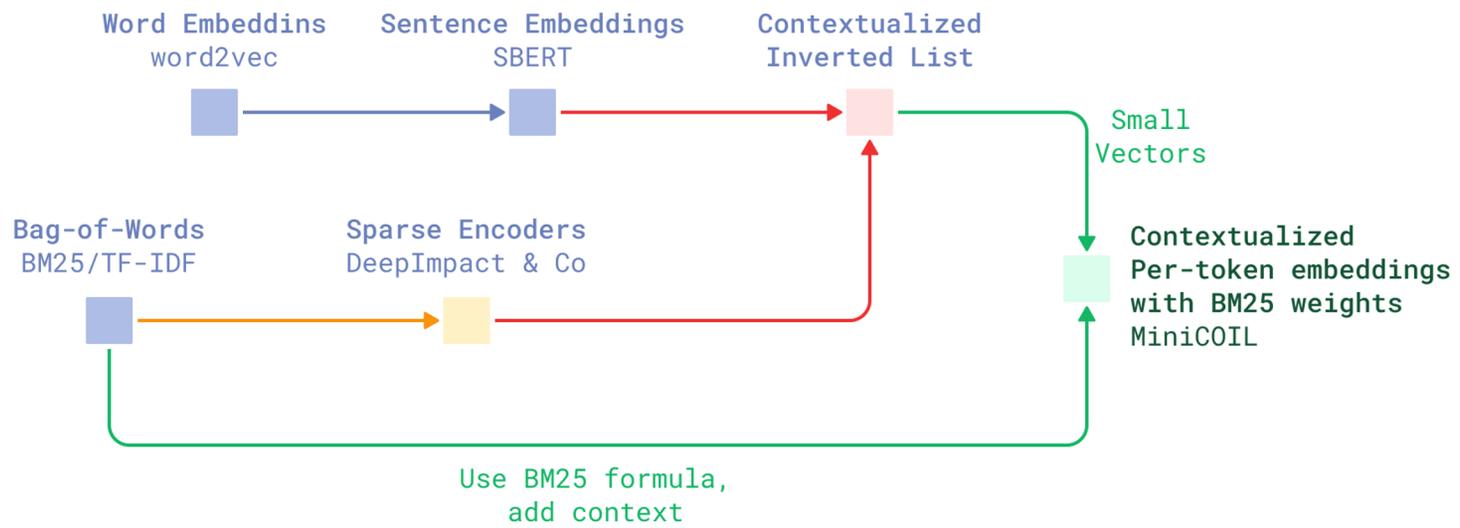
WELCOME TO THE FAMILY

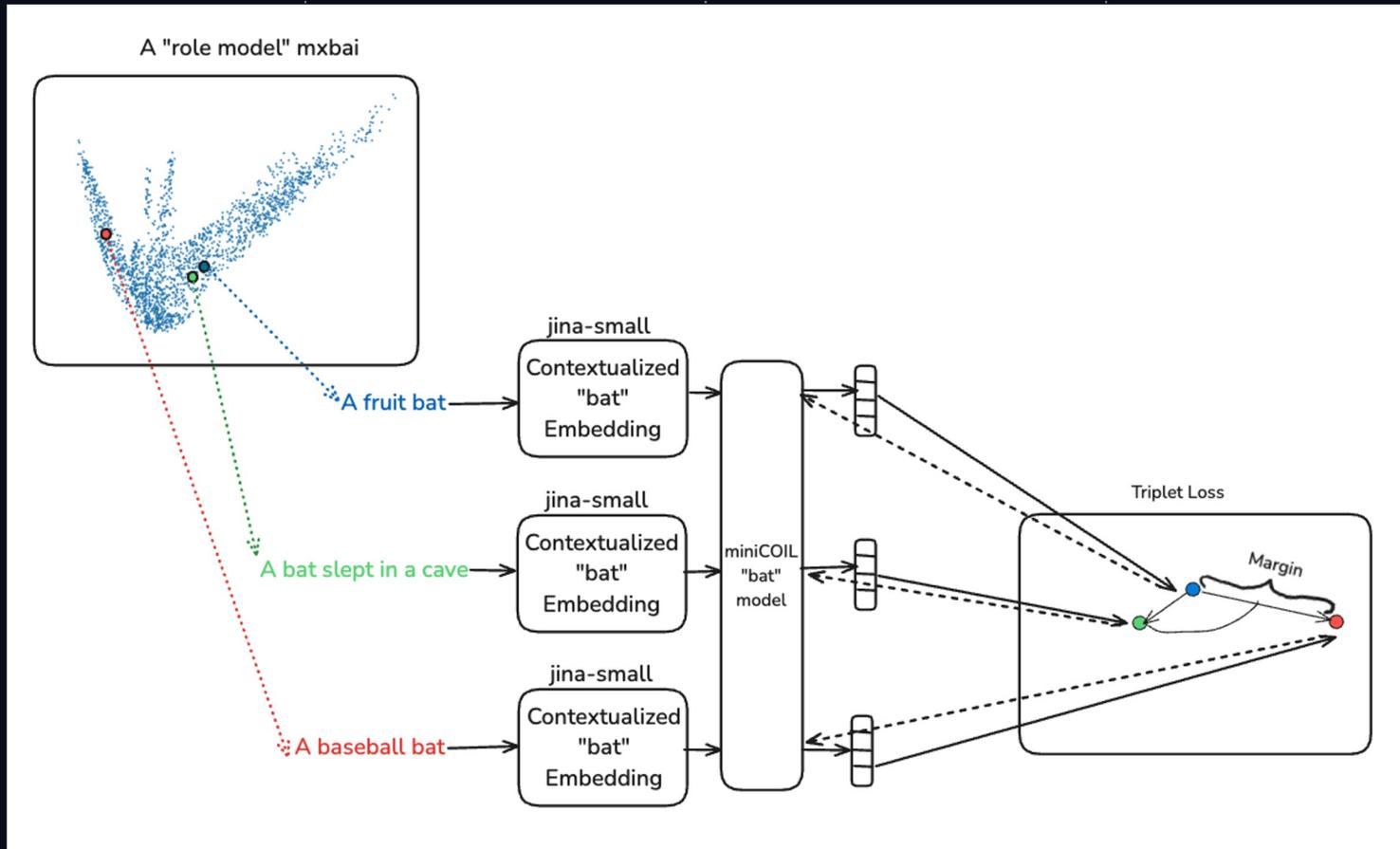
MARVEL STUDIOS
the Fantastic 4
FIRST STEPS

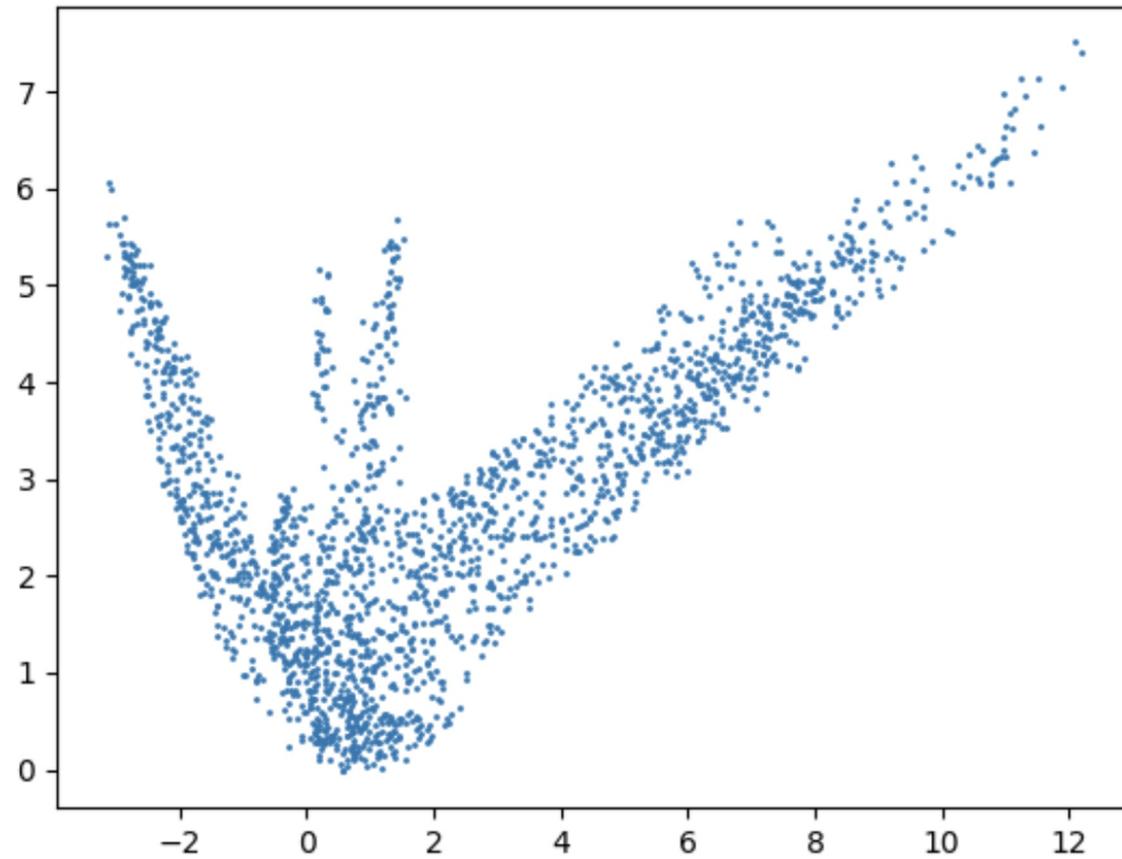
#theFantastic4

ONLY IN CINEMAS JULY

© 2015 The Marvel Family. TM & © 2015 Marvel.









Architecture

Context Encoder:

We use a transformer model (we used Jina, but it is agnostic) to generate contextualized token embeddings.



Architecture

Word Mapping:

miniCOIL works with real words or stems, not subword tokens – like “retriever” instead of “re ##trie ##ver”.



Tokens

“retriever”

- [re, ##trie, ##ver] (tokens)
- “retriever” (miniCOIL)



Architecture

Downprojection Model:

Each word in the vocabulary has its own lightweight model that projects the contextualized embedding to a 4-dimensional vector.



Architecture

Sparse Scoring:

Each 4D vector is multiplied by the BM25 score for that word, creating a new sparse representation that's both semantic and interpretable.



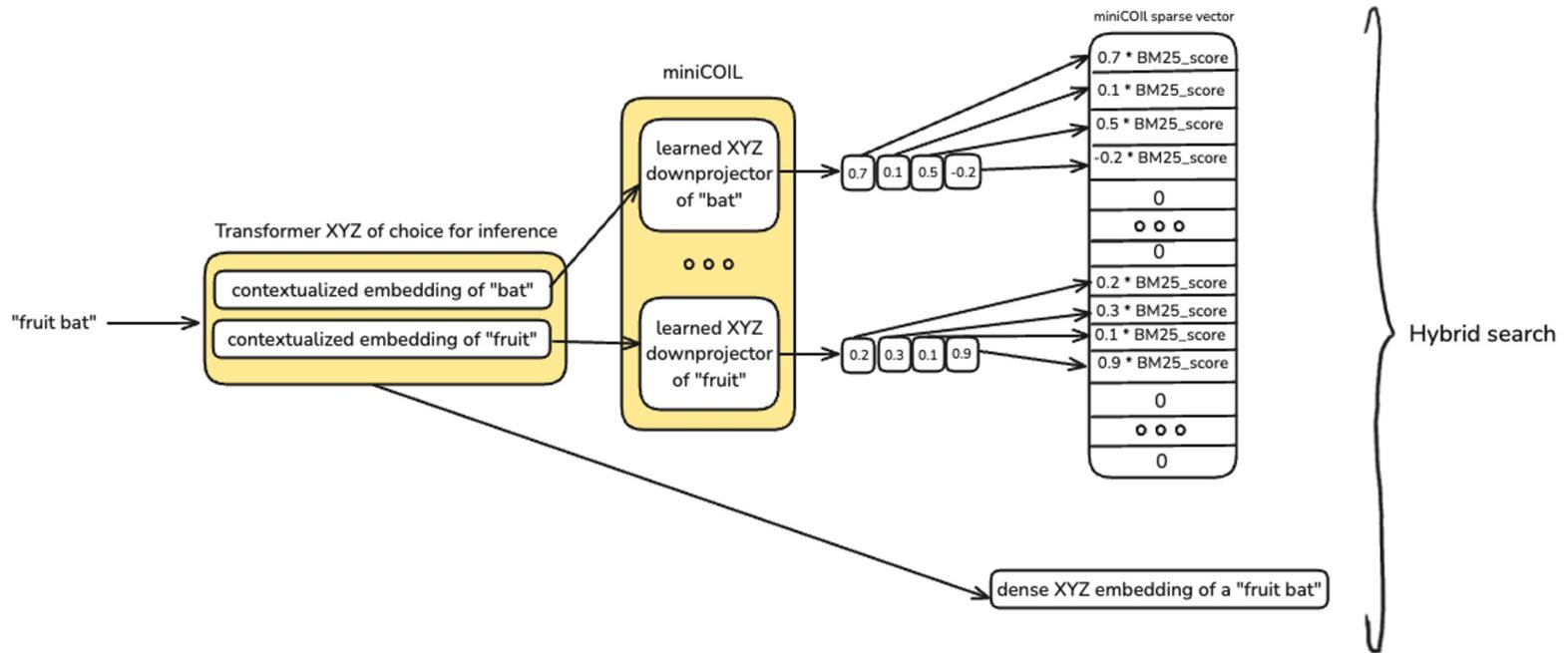
Architecture

Hybrid Use:

The resulting miniCOIL vector works seamlessly in hybrid search alongside dense vectors.



- Validation loss closely matches the theoretical lower bound
- Model generalizes across examples





MS MARCO

BM25 vs miniCOIL

NDCG@10

BM25: 0.237

miniCOIL: 0.244



**Not about
benchmarks...**



DEMO

miniCOIL Strengths

Flexible Transformer Support

Compatible with various transformers used for inference in hybrid search

Simple Architecture

1 word = 1 small trainable model, enabling extremely fast inference

No Labeled Training Data

Avoids relevance objectives, reducing data collection and overfitting risks

Vocabulary Extendable

New words can be added to the vocabulary without retraining existing models



Thierry Damiba

DATA SCIENTIST AT QDRANT



@ptdamiba



thierrypdamiba



thierrydamiba



td@qdrant.com

