# Information Retrieval

Sept 30, 2025 @ Introduction to Human Language Technology

Eugene Yang eugene.yang@jhu.edu

# What is Information Retrieval?

**(relevant)**

**Retrieve information from a storage based on user's information need**

# Don't we have Google?

**Yes, but Google is not all.**

**Still need to understand the intent (it's about MLB)**

**User still look for references**

**Still need retrieval for source information**

Google

world series champions prediction for 2025

AI Mode · All · News · Images · Short videos · Videos · Forums · More · Tools

✦ AI Overview

While predictions vary, recent odds and reports from late September 2025 suggest the **Philadelphia Phillies** are the favorites to win the World Series, having overtaken the Dodgers in the betting market, according to ESPN and MLB.com. Other strong contenders mentioned with favorable odds include the Los Angeles Dodgers, New York Yankees, and Atlanta Braves.

**Betting Favorites:**

**Philadelphia Phillies:**

The Phillies are now the favorites to win the 2025 World Series, holding the top spot in betting odds from major sportsbooks like BetMGM.

**Los Angeles Dodgers:**

The Dodgers were initially favored but have now fallen behind the Phillies in the betting odds, though they remain a top contender with strong odds.

**New York Yankees:**

The Yankees are consistently listed among the top few favorites, holding strong odds throughout the season.

**Atlanta Braves:**

The Braves also remain a significant threat, with their name appearing on lists of top contenders and favorable odds.

**Other Potential Contenders:**

**Milwaukee Brewers:**

The Brewers have secured the top seed and home-field advantage throughout the 2025 playoffs, making them a strong team to watch.

**Seattle Mariners:**

One insider has boldly predicted the Mariners will win the World Series.

World Series odds: Phillies stake claim as favorites - ESPN
Sep 22, 2025 — According to ESPN BET, the Dodgers remained the favorite to win the Worl...
Ⓔ ESPN ⋮

MLB World Series predictions: Expert picks for 2025 champion, ...
Mar 27, 2025 — 2025 World Series odds * Dodgers: +230. * Yankees: +900. * Braves: +9...
CBS Sports ⋮

Sports betting odds favor Phillies for 2025 title - MLB.com
Sep 24, 2025 — This article was contributed by BetMGM. For more sports betting insights, che...
MLB.com ⋮

Show all

# Google Search is just one implementation

Google trained us well!

- Even faster?

- Smarter?

- Cross language?

# Hard Matching Problem

- Text to text
  - Search in notes
  - Cross language search
  - Cross domain search
- Text to other modalities
  - Image search
  - Video search

# Different Search Process

- Iterative search
  - e.g., electronic discovery and systematic review
- Conversational search
  - Alexa search
- Recommendation systems
  - Implicit queries

# Core Problem

- Rank relevant document at top
- Do it fast

**Ranked List**

# Design Space

**Effectiveness**

- Definition of relevancy

- How to model relevancy

**Efficiency**

- How fast

- Fast at what stage

**Downstream Customer**

- Human or LLM?

- What kind of presentation?

# Agenda

- What is information retrieval?
- Retrieval modeling and pipeline
  - statistical and neural
- Evaluation
- Retrieval-Augmented Generation
- Research problems

# Retrieval Modeling and Pipeline

Modeling relevancy and do it fast

# Three main modeling strategies

- **Pointwise**
- Pairwise
- Listwise

- And combinations of them



https://medium.com/vptech/learning-to-rank-at-veepee-ed420fd828e5

# Statistical Models

$$\text{score}(D, Q) = \sum \boxed{\text{How important the term is}} \ \text{x} \ \boxed{\text{How often the term appear in the D}}$$

For each query term

$$\text{score}(D, Q) = \sum \boxed{\text{Inverted document frequency}} \ \text{x} \ \boxed{\text{Term frequency}}$$

For each query term

**TF-IDF**

$$\text{score}(D, Q) = \sum_{i=1}^{n} log \frac{N}{n_t} \times \log(f(q_i, D) + 1)$$

**BM25**

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

# How to make it fast?

- "Fast" in responding to queries

- Better data structure

- Preprocess the data

# Inverted Index

Term Index

| | Doc 1 | Doc 2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | Doc 7 | Doc 8 | | Postings |
|---|---|---|---|---|---|---|---|---|---|---|
| aid | | | | 1 | | | | 1 | | 4, 8 |
| all | | 1 | | 1 | | 1 | | | | 2, 4, 6 |
| back | 1 | | 1 | | | | 1 | | | 1, 3, 7 |
| brown | 1 | | 1 | | 1 | | 1 | | | 1, 3, 5, 7 |
| come | | 1 | | 1 | | 1 | | 1 | | 2, 4, 6, 8 |
| dog | | | 1 | | 1 | | | | | 3, 5 |
| fox | | | 1 | | 1 | | 1 | | | 3, 5, 7 |
| good | | 1 | | 1 | | 1 | | 1 | | 2, 4, 6, 8 |
| jump | | | 1 | | | | | | | 3 |
| lazy | 1 | | 1 | | 1 | | 1 | | | 1, 3, 5, 7 |
| men | | 1 | | 1 | | | | 1 | | 2, 4, 8 |
| now | | 1 | | | | 1 | | 1 | | 2, 6, 8 |
| over | 1 | | 1 | | 1 | | 1 | 1 | | 1, 3, 5, 7, 8 |
| party | | | | | | 1 | | 1 | | 6, 8 |
| quick | 1 | | 1 | | | | | | | 1, 3 |
| their | 1 | | | | 1 | | 1 | | | 1, 5, 7 |
| time | | 1 | | 1 | | 1 | | | | 2, 4, 6 |

# Inverted Index

| Term Index | | Postings |
|---|---|---|
| aid | | 4, 8 |
| all | | 2, 4, 6 |
| back | | 1, 3, 7 |
| brown | | 1, 3, 5, 7 |
| come | | 2, 4, 6, 8 |
| dog | | 3, 5 |
| fox | | 3, 5, 7 |
| good | | 2, 4, 6, 8 |
| jump | | 3 |
| lazy | | 1, 3, 5, 7 |
| men | | 2, 4, 8 |
| now | | 2, 6, 8 |
| over | | 1, 3, 5, 7, 8 |
| party | | 6, 8 |
| quick | | 1, 3 |
| their | | 1, 5, 7 |
| time | | 2, 4, 6 |

# Inverted Index

# Two-Stage System

- Offline preprocessing and indexing
  - Define retrieval unit
  - Tokenization
  - Build the inverted index

- Online query serving
  - Traverse the inverted index and score it

## APACHE LUCENE™

### Apache 2.0 licensed

Apache Lucene is distributed under a commercially friendly Apache Software license

# Welcome to Apache Lucene

The Apache Lucene™ project develops open-source search software. The project releases a core search library, named Lucene™ core, as well as PyLucene, a python binding for Lucene.

Lucene Core is a Java library providing powerful indexing and search features, as well as spellchecking, hit highlighting and advanced analysis/tokenization capabilities. The PyLucene sub project provides Python bindings for Lucene Core.

### Latest Lucene Core News

Apache Lucene™ 8.11.4 available (24.Sep)

Apache Lucene™ 9.11.1 available (27.Jun)

Apache Lucene™ 9.11.0 available (06.Jun)

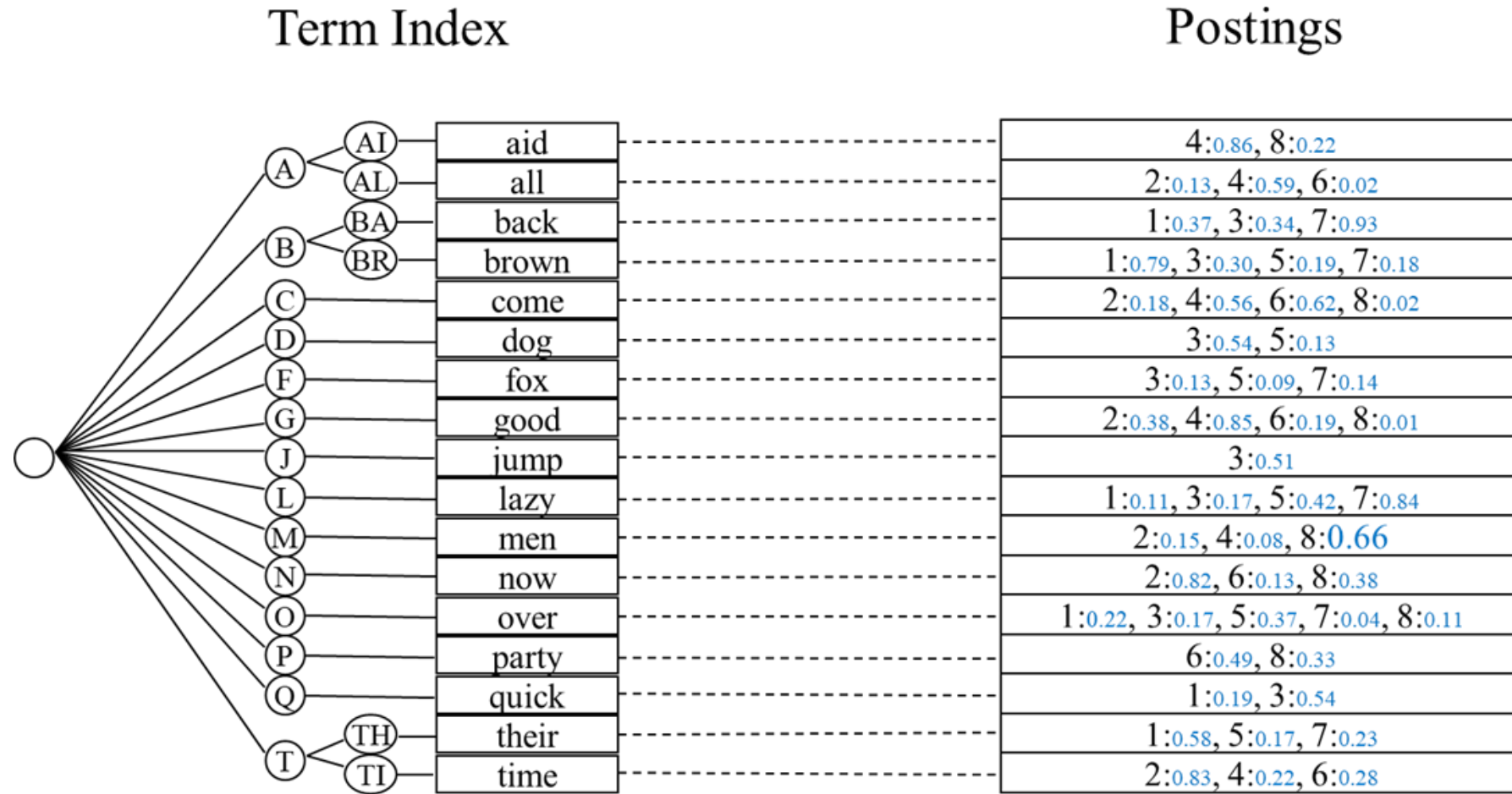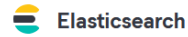**ANNOUNCEMENT:** The Solr™ sub project has moved to a separate Top Level Project (TLP). All things Solr can now be found at https://solr.apache.org/. Mailing lists and git repositories have changed, please see details on the Solr website.

### Projects

Lucene Core (Java)
PyLucene
Open Relevance (Discontinued)

### About

License
Who We are
TLP News
Code of Conduct

**Events**

---

### ⬢ Elasticsearch

# The heart of the free and open Elastic Stack

Elasticsearch is a distributed, RESTful search and analytics engine, scalable data store, and vector database capable of addressing a growing number of use cases. As the heart of the Elastic Stack, it centrally stores fine-tuned relevancy, and powerful ana

**Start free trial**       Vie

Download Elastic

---

🖫 **elastic / elasticsearch** Public

🔔 Notifications    ⑂ Fork 24.7k    ☆ Star 69.7k

<> Code    ⊙ Issues 3.9k    ⑄ Pull requests 773    ▷ Actions    ⊞ Projects    ⊙ Security    �peⁿ Insights

⑂ main ▾       ⑂ 391 Branches   ⬚ 404 Tags            🔍 Go to file        <> Code ▾

### About

Free and Open Source, Distributed, RESTful Search Engine

🔗 www.elastic.co/products/elasticsearch

`java`  `search-engine`  `elasticsearch`

| | | | |
|---|---|---|---|
| 🖧 smalyshev Improve DateTime error handling and add some bad date tests... ⬚⬚⬚ ✓ 5e06092 · 4 hours ago | | ⏱ 79,984 Commits | |
| 📁 .buildkite | Run snyk dependency checks on 8.x (#113117) | last week | |
| 📁 .ci | Workaround packaging tests failures on debian10 (#113... | 9 hours ago | |
| 📁 .github | Remove Analytical engine CODEOWNERS (#113178) | 2 days ago | |
| 📁 .idea | Don't apply IntelliJ illegal module dependency inspection ... | 10 months ago | |
| 📁 benchmarks | ESQL: Speed up CASE for some parameters (#112295) | yesterday | |
| 📁 build-conventions | Add AGPLv3 as a supported license | 2 weeks ago | |
| 📁 build-tools-internal | Always use CLDR locale on ES v9 (#113184) | 3 days ago | |
| 📁 build-tools | Add AGPLv3 as a supported license | 2 weeks ago | |
| 📁 client | Add AGPLv3 as a supported license | 2 weeks ago | |
| 📁 dev-tools | Add AGPLv3 as a supported license | 2 weeks ago | |
| 📁 distribution | Always use CLDR locale on ES v9 (#113184) | 3 days ago | |
| 📁 docs-mdx/painless | [DOCS] Adds an MDX file for testing purposes. (#106165) | 6 months ago | |
| 📁 docs | Improve DateTime error handling and add some bad dat... | 4 hours ago | |
| 📁 gradle | Update Gradle wrapper to 8.10.1 (#112948) | last week | |
| 📁 libs | Small performance improvement in h3 library (#113385) | 2 days ago | |

📖 Readme
⚖ View license
⚖ Security policy
⋀ Activity
⊞ Custom properties
☆ 69.7k stars
👁 2.7k watching
⑂ 24.7k forks

Report repository

### Releases 155

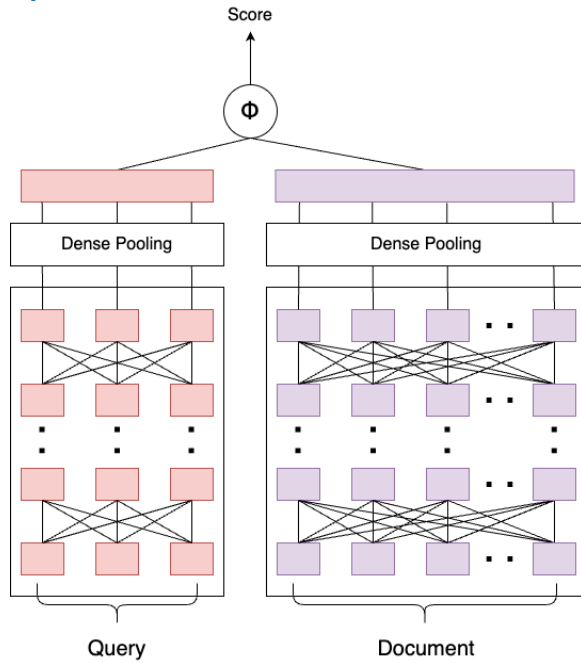🏷 Elasticsearch 8.15.1 (Latest)
3 weeks ago

+ 154 releases

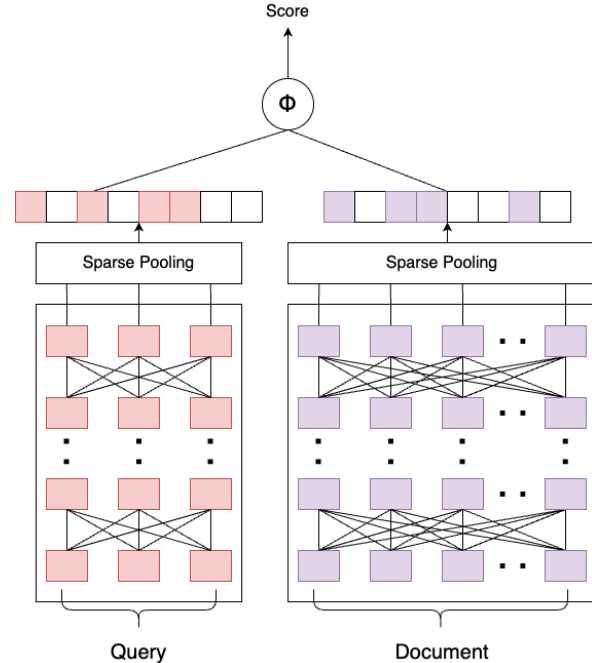### Packages

No packages published

# Can we go beyond surface forms?
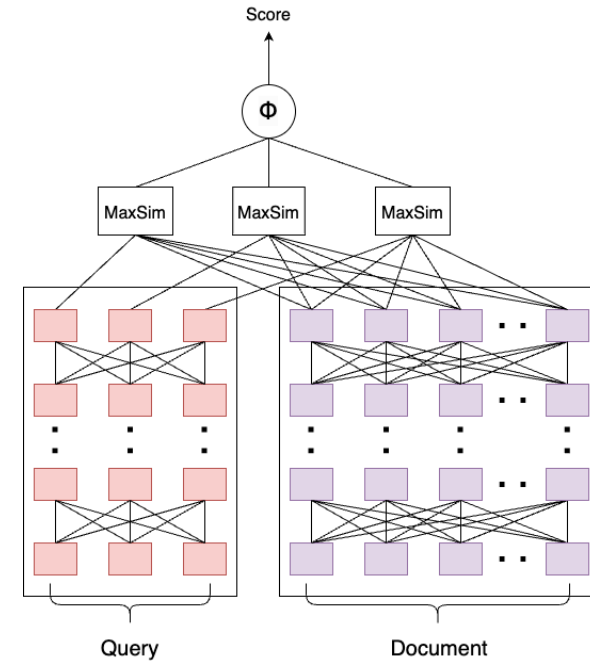
neural language models
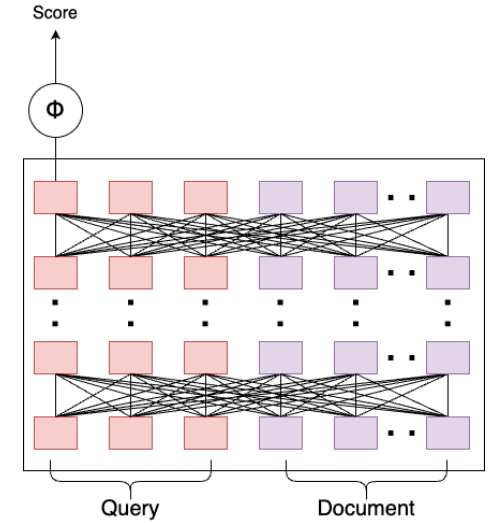
Bi-Encoder

Cross Encoder

One Dense Vector
Per Sequence
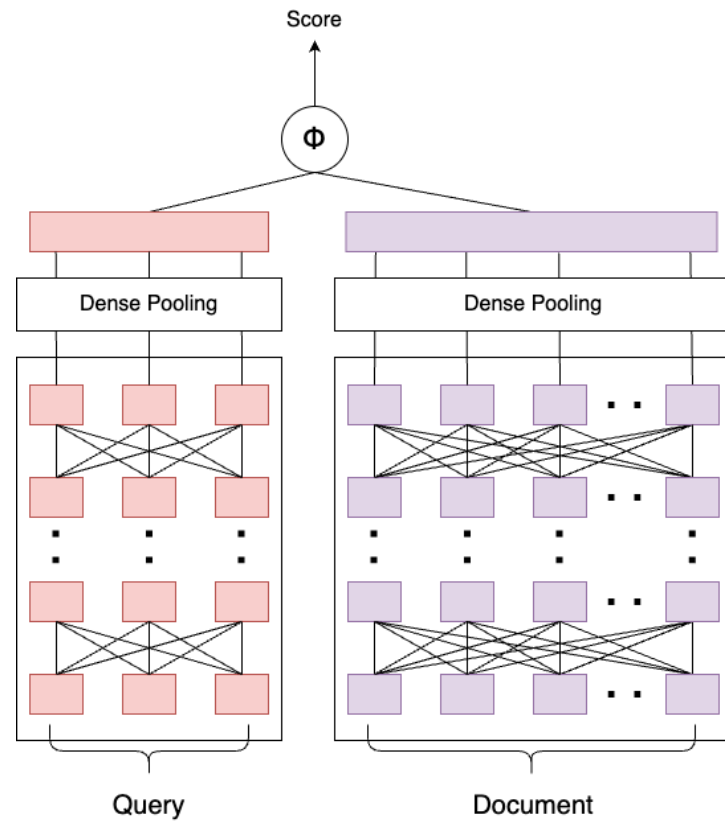e.g., DPR

One **Sparse** Vector
Per Sequence
e.g., SPLADE

Multiple Dense Vectors
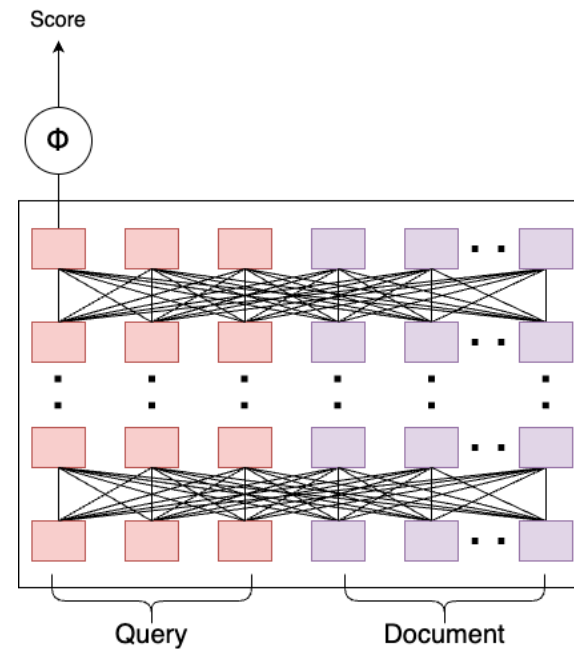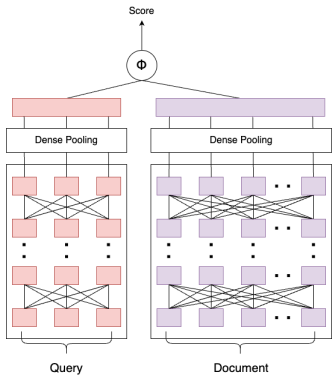Per Sequence
e.g., ColBERT

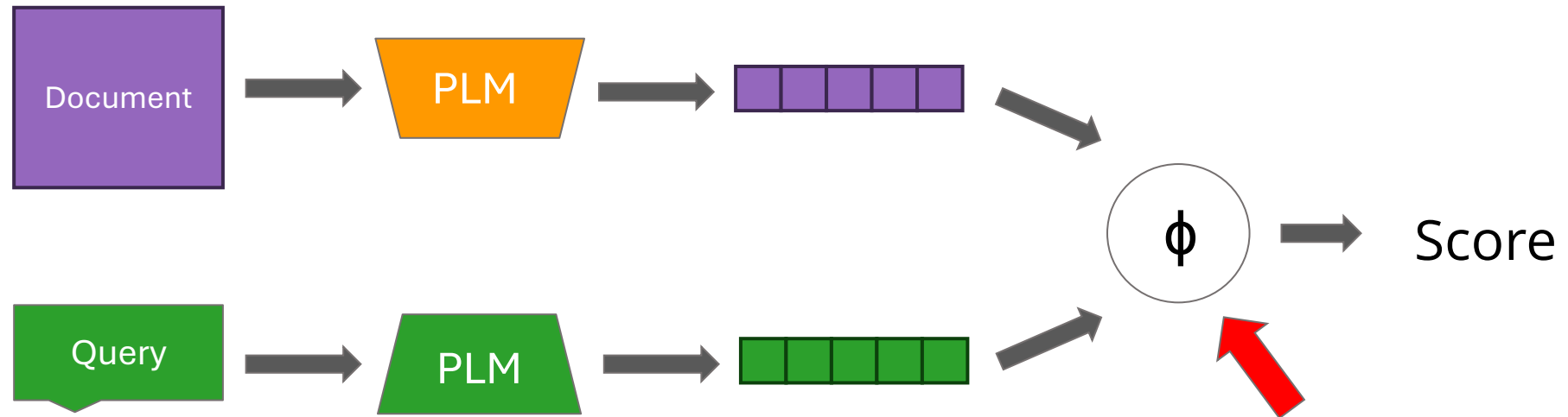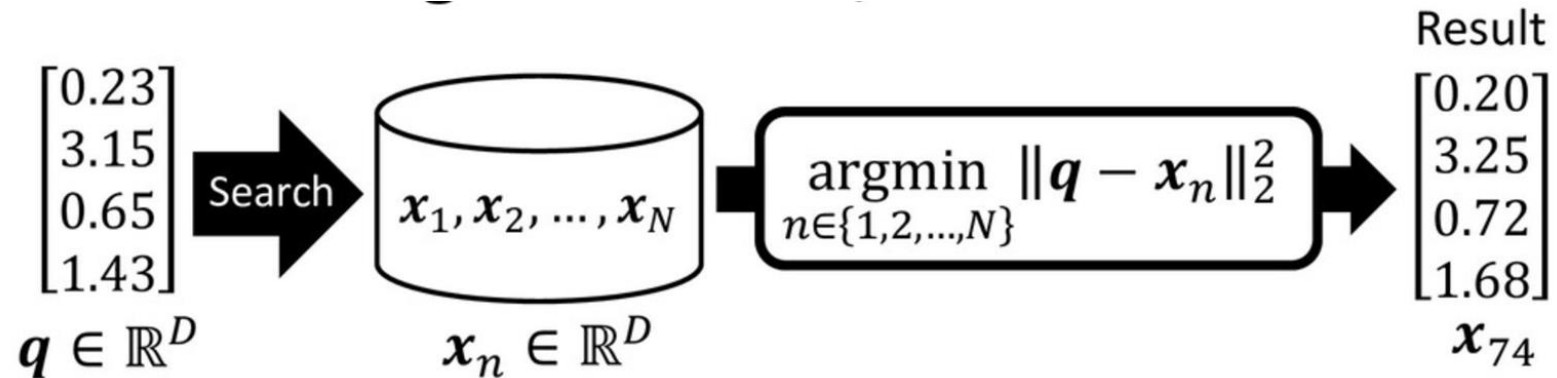Joint Encoder
e.g., monoBERT

Online　　　Offline　　　　　　　　　Both Online
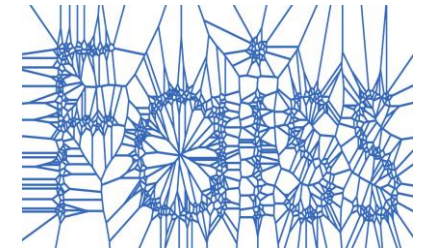
**Separate query and document processing**

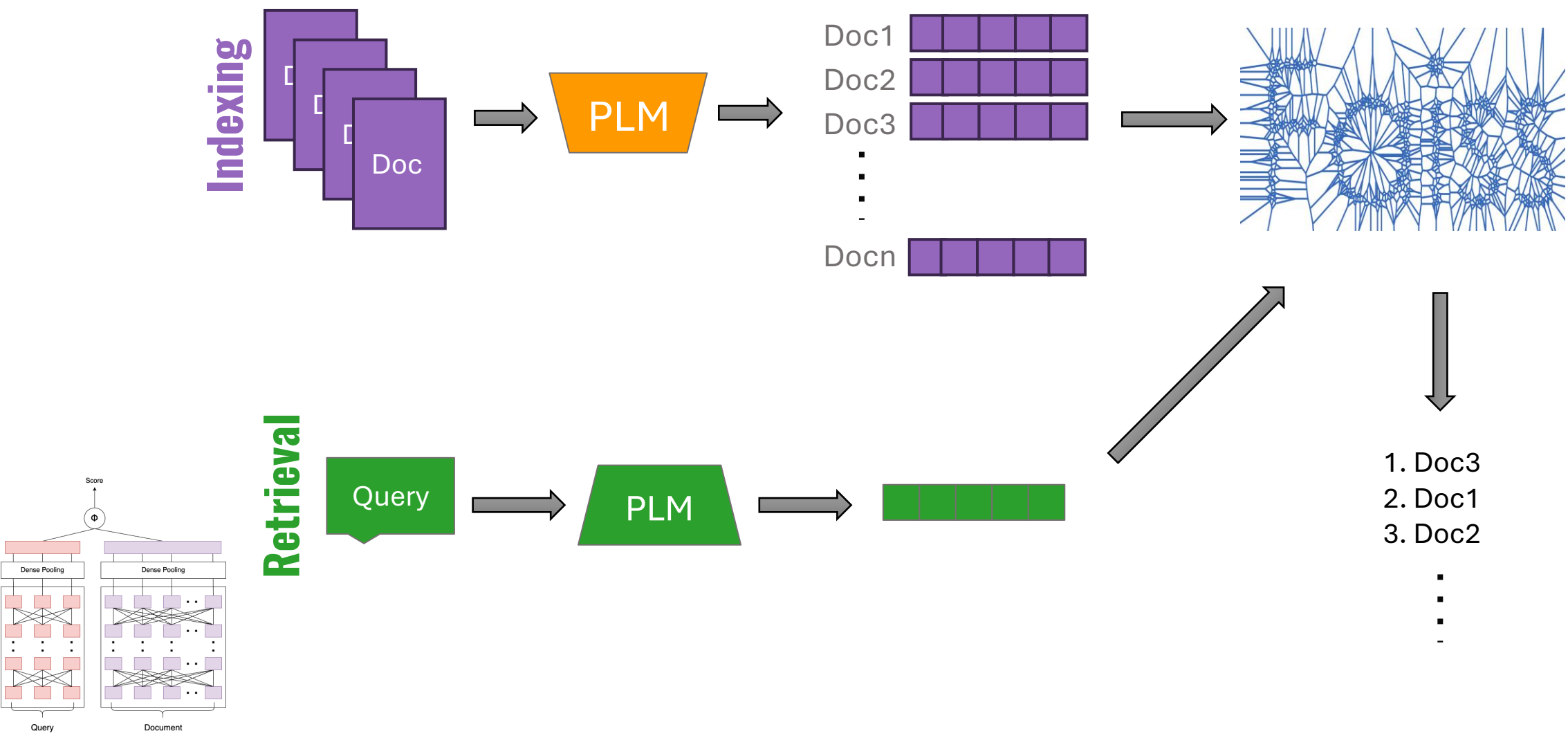# One Vector per Query, One Vector per Document
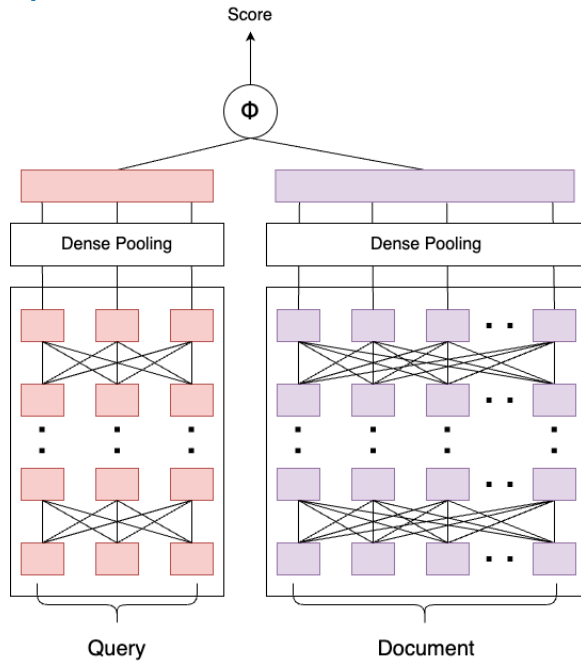


Score

# Nearest Vectors aka Neighbors



$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix}$$
$q \in \mathbb{R}^D$

Search

$x_1, x_2, \ldots, x_N$
$x_n \in \mathbb{R}^D$

$$\underset{n \in \{1,2,\ldots,N\}}{\text{argmin}} \; \|q - x_n\|_2^2$$

Result
$$\begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix}$$
$x_{74}$

- Linear Search
  - Slow (scales linearly in size of document collection)

- Approximate Methods (e.g., Product Quantization) → **ANN**
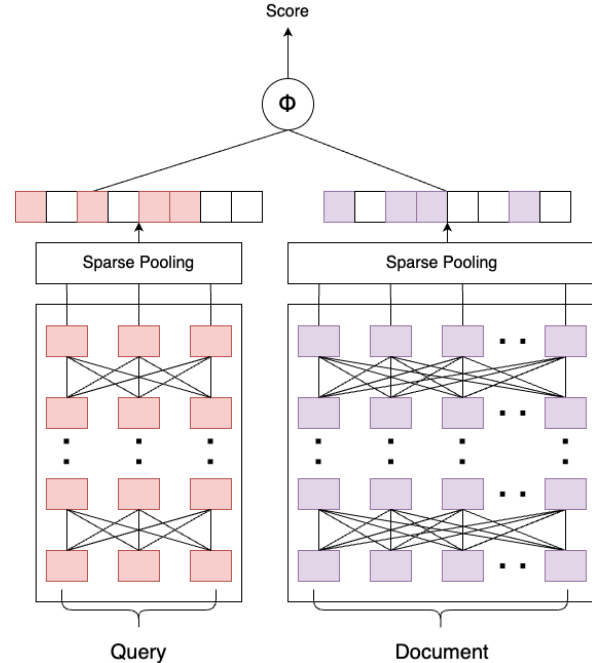  - Faster Search

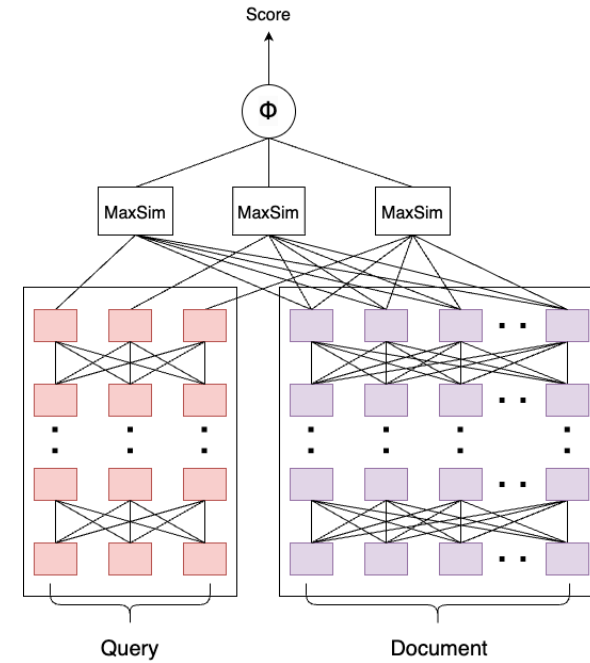- Runtime Efficiency vs Effectiveness

# DPR Indexing and Retrieval

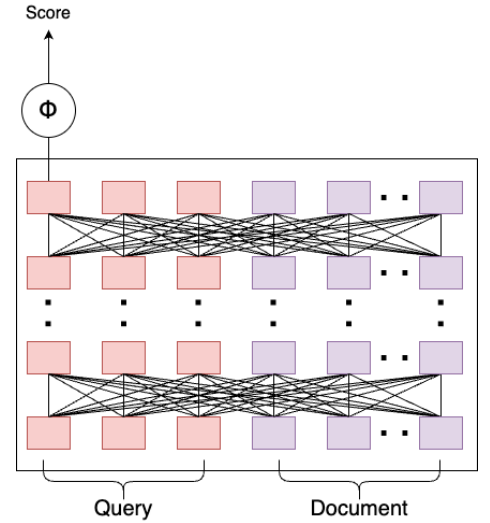**Bi-Encoder**

**Cross Encoder**

One Dense Vector
Per Sequence
e.g., DPR

One **Sparse** Vector
Per Sequence
e.g., SPLADE

Multiple Dense Vectors
Per Sequence
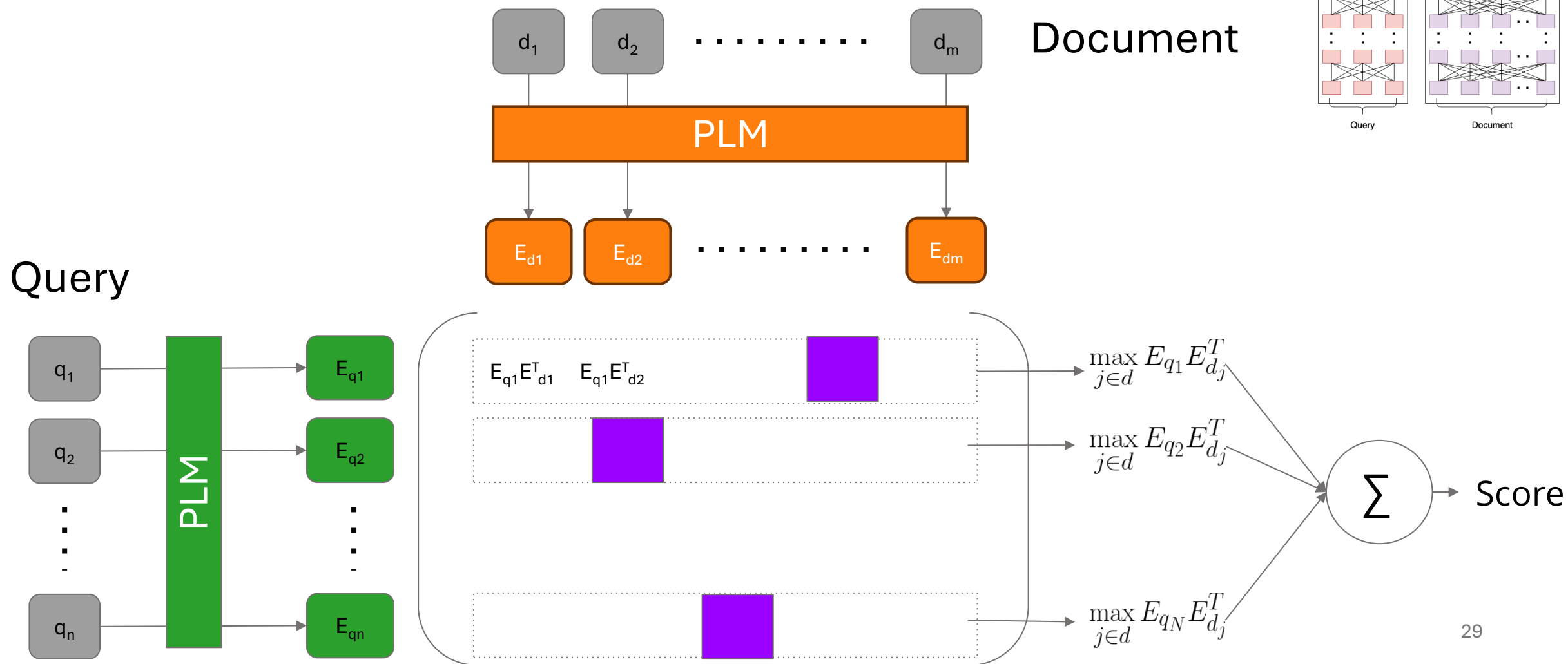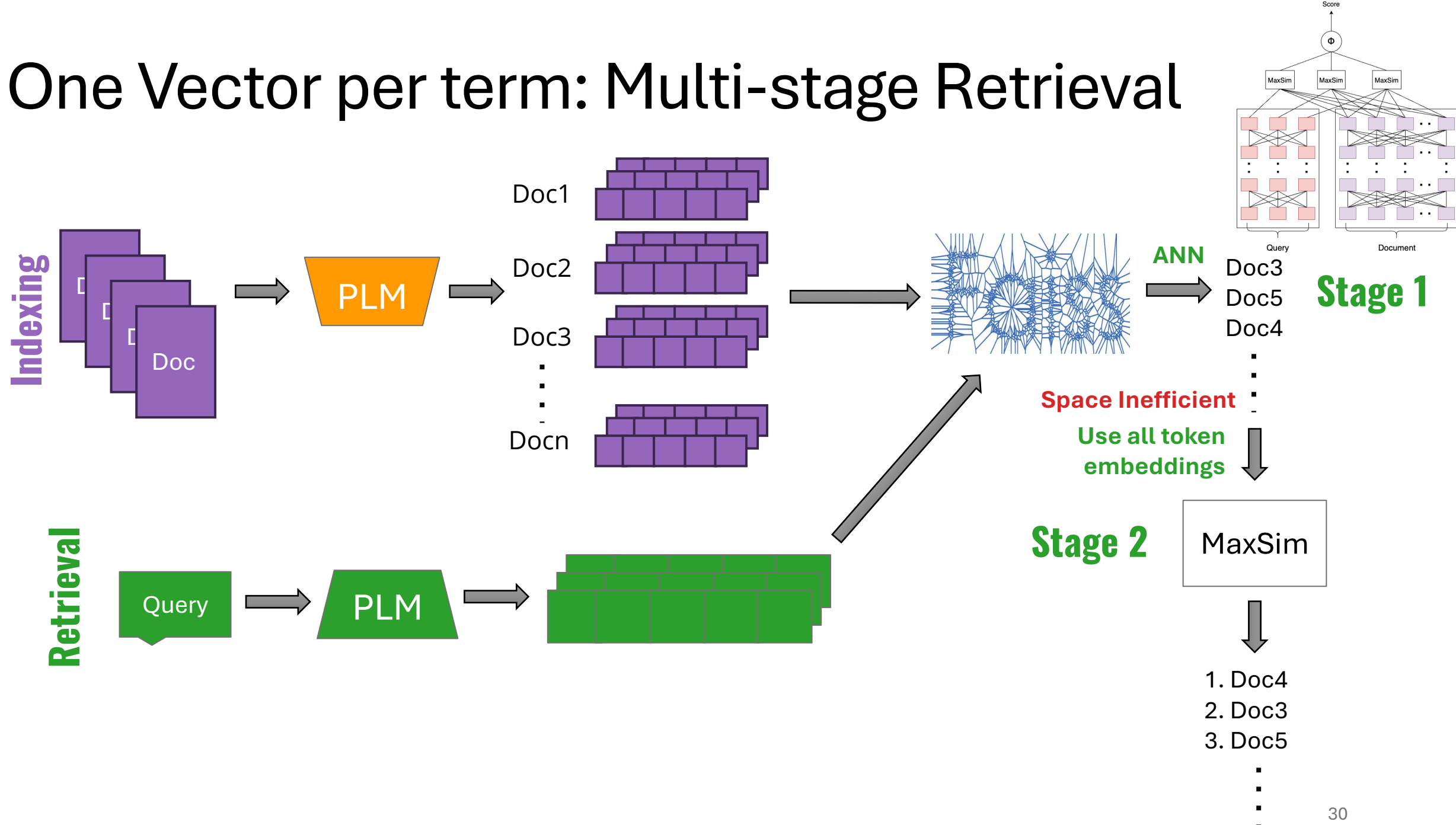e.g., ColBERT

Joint Encoder
e.g., monoBERT

# One Vector per Term: MaxSim

小 (small)

小 (little)

little

有 (had)

had

羔羊 (lamb)

lamb

羊 (sheep)

瑪莉雅 (maria)

聖母 (our lady)

羊肉 (mutton)

瑪莉 (mary)

mary

Chinese Document term embedding

English Query term embedding

只 (only)

一隻 (one)

a

許多 (many)

# MaxSim in Action -- ColBERT

# One Vector per term: Multi-stage Retrieval

**Contributors** 6

**Languages**

● Python 93.3%   ● C++ 5.6%
○ Other 1.1%

# PyLate

Flexible Training and Retrieval for Late Interaction Models



Documentation   License MIT

PyLate is a library built on top of Sentence Transformers, designed to simplify and optimize fine-tuning, inference, and retrieval with state-of-the-art ColBERT models. It enables easy fine-tuning on both single and multiple GPUs, providing flexibility for various hardware setups. PyLate also streamlines document retrieval and allows you to load a wide range of models, enabling you to construct ColBERT models from most pre-trained language models.

## Installation

You can install PyLate using pip:

```
pip install pylate
```

---

hltcoe / **ColBERT-X**   Public

<> Code   ⊙ Issues   ⑂ Pull requests   ▷ Actions

⑂ plaid-x ▾

👤 eugene-yang Update collection_utils.py ···

📁 colbert            Update collection_utils.p
📁 docs              update theme
📄 .gitignore         update install
📄 LICENSE           Initial commit with the
📄 MANIFEST.in       Version 0.3.1 (#7)
📄 README.md         Update README.md (#
📄 conda_env.yml     commit the environmen
📄 conda_env_cpu.yml Make end to end test w
📄 requirements.txt   Version 0.3.1 (#7)
📄 setup.py          Version 0.3.1 (#7)

📖 README   ⚖ MIT license

# PLAID-X

This is a generalized version of PLAID and the previ
CLIR. The codebase supports models trained with th
scripts, which are not compatible with the PLAID co
the Stadford Futuredata Group.

---

🐕 **hltcoe's Collections**

**Translate-Distill**      Multilingual Translate-Distill

## Translate-Distill

updated Mar 27

Collection of trained model and teacher scores for distillation for paper "Translate-Distill" Code: https://github.com/hltcoe/ColBERT-X

▲ Upvote 3

**Translate-Distill: Learning Cross-Language Dense Retrieval by Translation and Distillation**
📄 Paper • 2401.04810 • Published Jan 9

🐕 hltcoe/plaidx-large-zho-tdist-mt5xxl-engzho
Updated Mar 15 • ⬇ 2

🐕 hltcoe/plaidx-large-zho-tdist-mt5xxl-zhozho
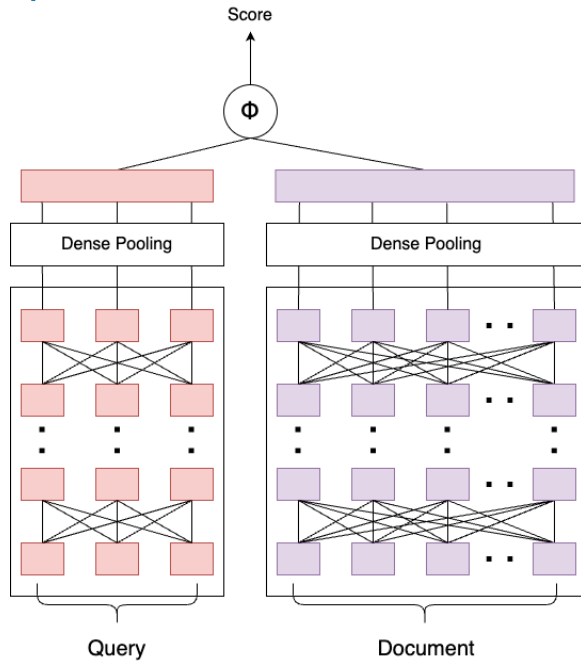Updated Mar 15 • ⬇ 1

🐕 hltcoe/plaidx-large-zho-tdist-t53b-engeng
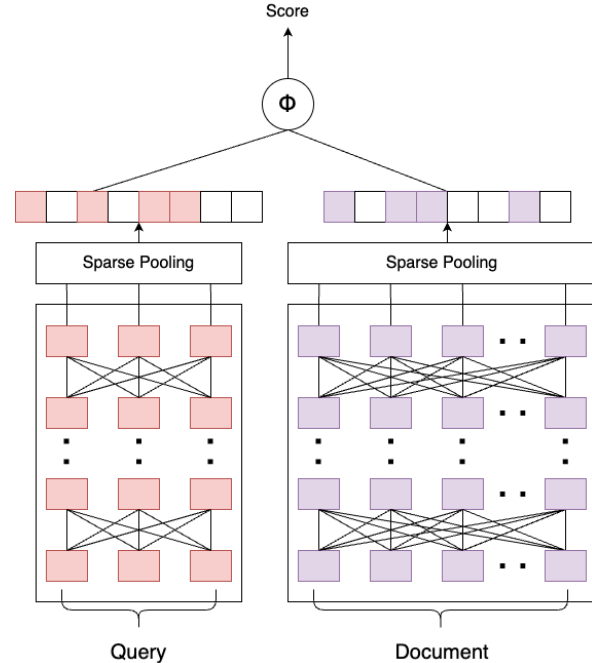Updated Mar 15 • ⬇ 1

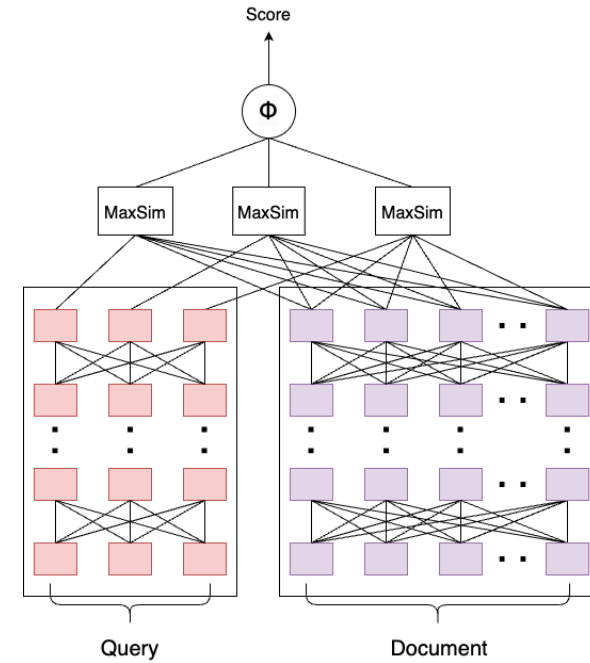🐕 hltcoe/plaidx-large-zho-tdist-mt5xxl-engeng
Updated Mar 15 • ⬇ 25

Bi-Encoder

Cross Encoder

One Dense Vector
Per Sequence
e.g., DPR

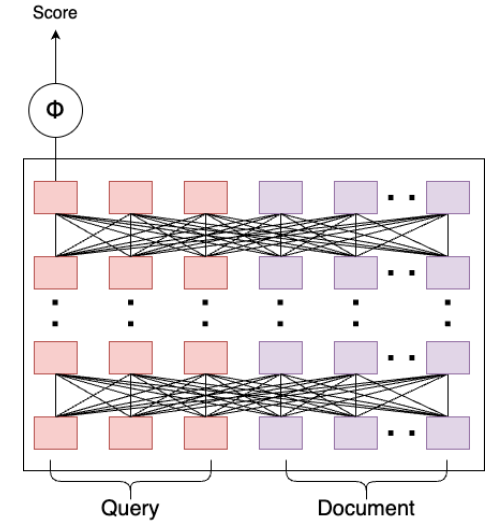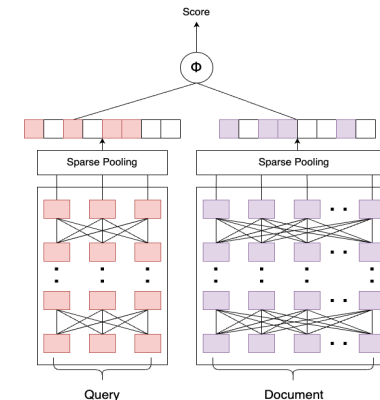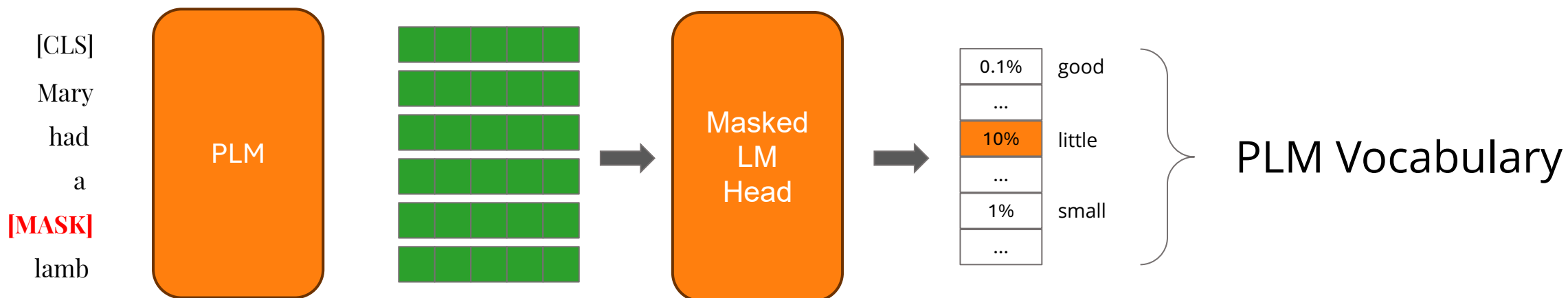One **Sparse** Vector
Per Sequence
e.g., SPLADE

Multiple Dense Vectors
Per Sequence
e.g., ColBERT

Joint Encoder
e.g., monoBERT

# High-dimensional Vector: Masked LM

# Learned Sparse Retrieval

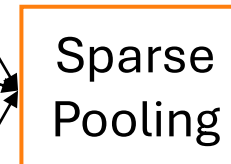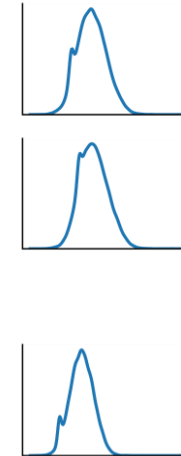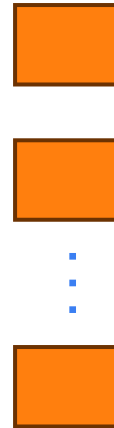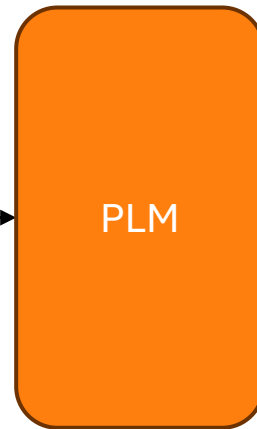# SPLADE Search Pipeline

Bi-Encoder

Cross Encoder

One Dense Vector
Per Sequence
e.g., DPR

One **Sparse** Vector
Per Sequence
e.g., SPLADE

Multiple Dense Vectors
Per Sequence
e.g., ColBERT

Joint Encoder
e.g., monoBERT

# Cross-Encoder as Reranker



Query

Doc

Cross-Encoder
with PLM

0.9063

# Using Generative Models

Query: What does Mary have
Doc: Mary had a little lamb.
Relevant:

Generative (L)LM

e.g, Qwen3

**Not a number!**

Yes

Pradeep, Ronak, Rodrigo Nogueira, and Jimmy Lin. "The expando-mono-duo design pattern for text ranking with pretrained sequence-to-sequence models." arXiv preprint arXiv:2101.05667 (2021).

# Using Generative Models



Query: What does Mary have
Doc: Mary had a little lamb.
Relevant:

Generative (L)LM

Yes (0.08)

No (0.001)

Yes

Score
Φ
Query    Document

Pradeep, Ronak, Rodrigo Nogueira, and Jimmy Lin. "The expando-mono-duo design pattern for text ranking with pretrained sequence-to-sequence models." arXiv preprint arXiv:2101.05667 (2021).

# Using Generative Models

**Pointwise score**

Query: What does Mary have
Doc: Mary had a little lamb.
Relevant:

Generative (L)LM

Yes (0.08)

No (0.001)

0.52

Score

Φ

Query          Document

Pradeep, Ronak, Rodrigo Nogueira, and Jimmy Lin. "The expando-mono-duo design pattern for text ranking with pretrained sequence-to-sequence models." arXiv preprint arXiv:2101.05667 (2021).

# Using Generative Models

**Pairwise score**

Query: What does Mary have
Doc0: JHU is in Baltimore
Doc1: Mary had a little lamb.
Relevant:

Generative (L)LM

Yes (0.001)

No (0.02)

0.49

Pradeep, Ronak, Rodrigo Nogueira, and Jimmy Lin. "The expando-mono-duo design pattern for text ranking with pretrained sequence-to-sequence models." arXiv preprint arXiv:2101.05667 (2021).

43

# Many kinds of LLM-based Rerankers

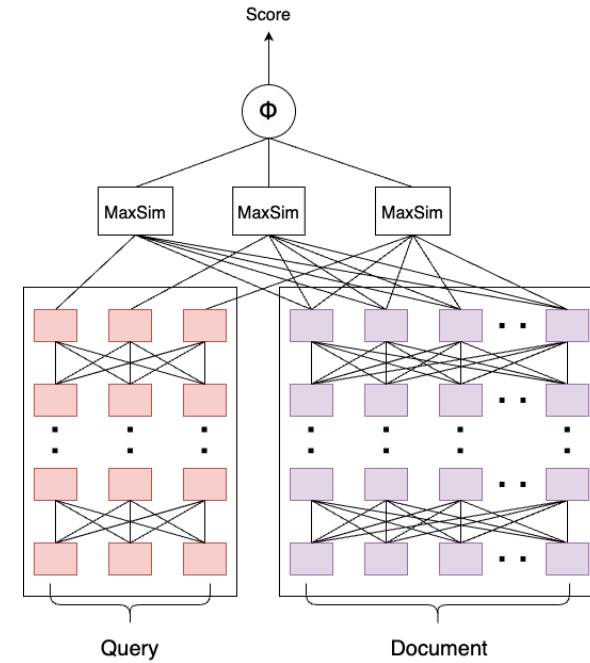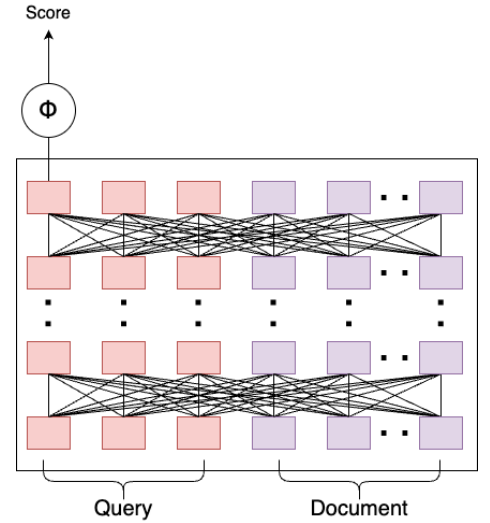|  | Generative Model | Reasoning Model |
|---|---|---|
| **Pointwise** | MonoT5<br>RankLlama | Rank1 |
| **Pairwise** | DuoT5 | |
| **Listwise/Setwise** | RankGPT<br>RankZephyr | Rank-K<br>ReasonRank |

Bi-Encoder

Cross Encoder

One Dense Vector
Per Sequence
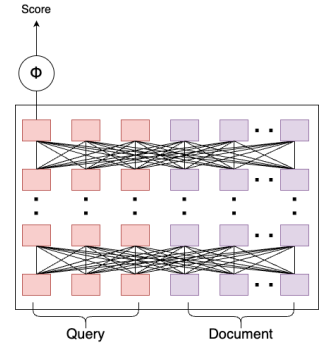e.g., DPR

One **Sparse** Vector
Per Sequence
e.g., SPLADE

Multiple Dense Vectors
Per Sequence
e.g., ColBERT

Joint Encoder
e.g., monoBERT

Bi-Encoder

Cross Encoder

**One Dense Vector**
Per Sequence
e.g., DPR

**One Sparse Vector**
Per Sequence
e.g., SPLADE

**Multiple Dense Vectors**
Per Sequence
e.g., ColBERT

**Joint Encoder**
e.g., Rank1

**More Effective**

**More Efficient (at Query Time)**

# Retrieve-and-Rerank System Combinations

| More Efficient<br>Less Effective | **Higher<br>Recall** → | More Effective<br>Less Efficient | **Final<br>objective** → | Final Score |

# Neural Retrieval System Pipeline

# PLM to IR Model

Pretraining ⟶ [ Pretrained LM (PLM) ] — Retrieval Finetuning ⟶ [ IR Model ]

- Align the representation

- Model "relevancy"

# Evaluation

Which system is better?

# What is Information Retrieval?

**(relevant)**

**Retrieve information from a storage based on user's information need**

# Which system retrieve more relevant information?

# Cranfield Paradigm Evaluation



Queries

Search Engine

Documents

Ranked List

Relevance Judgments

Search Engine C

Search Engine A

Search Engine B

# Cranfield Paradigm Evaluation



Queries

Topics

Search Engine

Ranked List

Relevance Judgments

Documents

Search Engine C

Search Engine A

Search Engine B

# Cranfield Paradigm Evaluation



One Query

Search Engine

Documents

One Topic

Ranked List

Relevance Judgments

Evaluation Metric Scoring

# IR-Specific Issues

- Topics vs Queries
  - Clear intent vs an expression of such intent
- Relevant vs related
  - Fulfilling the information need or not
- Relevance Judgements vs Labels
  - Opinion vs "fact"
- Ranked retrieval metrics
  - Measuring the quality/effectiveness of a ranked list

# IR Metrics

- Effective Metrics
  - Mean Average Precision
  - Normalized Discounted Cumulative Gain
  - Recall@k
- Efficiency Metrics
  - Indexing time
  - Index disk space
  - Query latency (average search time per query)

# Retrieval Augmented Generation

User Input

**Retrieval**

**Augmented Input**

**No knowledge after events after training**

RAG Output

**Grounded information for citation**

**Generation Summarization**

Figure 2: The overview of STORM that automates the pre-writing stage. Starting with a given topic, STORM identifies various perspectives on covering the topic by surveying related Wikipedia articles (①-②). It then simulates conversations between a Wikipedia writer who asks questions guided by the given perspective and an expert grounded on trustworthy online sources (③-⑥). The final outline is curated based on the LLM's intrinsic knowledge and the gathered conversations from different perspectives (⑦-⑧).

# A Green Field

- How the pipeline should be designed?
  - Predefined one or agentic system?
- Short vs long input from user
- How to make the response to be more useful?
  - Graphical or just text?
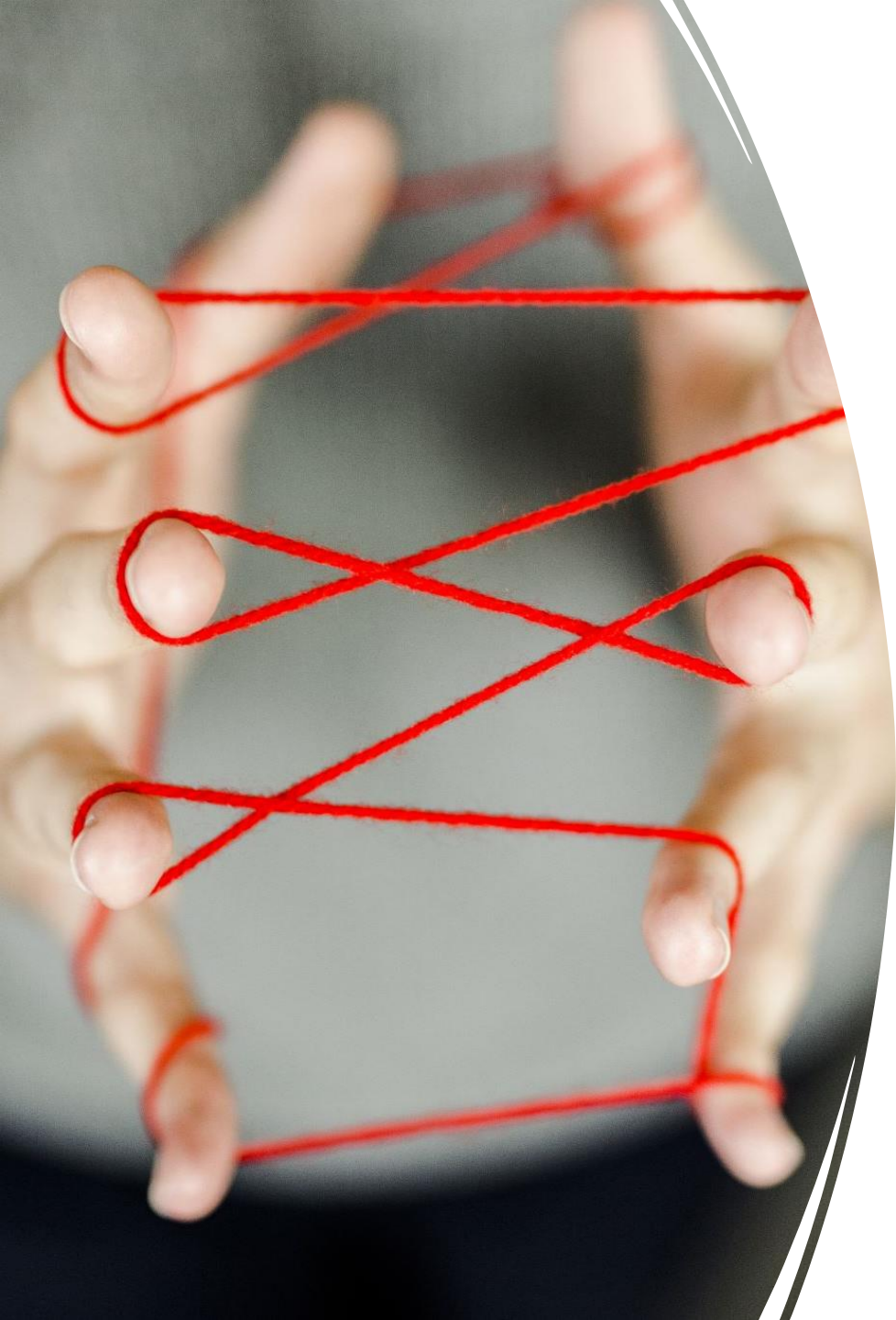- How to interact with human?

User Input

**Retrieval**

**Augmented Input**

**Generation Summarization**

RAG Output

**Feedback Loop**

# State of IR Research

# Better Retrieval Models

- More effective
  - Better/larger neural models
  - Better architecture?
  - Under harder setup, e.g., scholar search, multilingual, cross-modal, etc
- More efficient
  - Faster at query time
  - Less resource footprint, e.g., memory, storage, compute, etc
- Other qualities
  - Fairness, diversity, etc

# Other Retrieval Problems

- Conversational
  - Guessing intent, finding the "right" information to serve
- Iterative/interactive/human-in-the-loop
  - Rounds of interactions
- Generative
  - Returning a piece of text

# Evaluation

- What to measure
  - and when would it fail
- How to measure
  - Generative text? Citations?
- "Better" evaluation collection
  - Not necessarily larger