# Information Retrieval

Oct 1, 2024 @ 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.**

# 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
- (Set Retrieval)

# 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

# Agenda

- What is information retrieval?
- Retrieval Modeling and Pipeline
    - Statistical and Neural
- Evaluation
- State of IR Research and active 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

$$\text{score}(D, Q) = \sum_{\text{For each query term}} \boxed{\text{How important the term is}} \times \boxed{\text{How often the term appear in the D}}$$

$$\text{score}(D, Q) = \sum_{\text{For each query term}} \boxed{\text{Inverted document frequency}} \times \boxed{\text{Term frequency}}$$

**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

Postings

| Term | 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 |

Term Index nodes: A → AI, AL; B → BA, BR; C; D; F; G; J; L; M; N; O; P; Q; T → TH, TI

# Inverted Index

| Term Index | | Postings |
|---|---|---|
| AI | aid | 4, 8 |
| AL | all | 2, 4, 6 |
| BA | back | 1, 3, 7 |
| BR | brown | 1, 3, 5, 7 |
| C | come | 2, 4, 6, 8 |
| D | dog | 3, 5 |
| F | fox | 3, 5, 7 |
| G | good | 2, 4, 6, 8 |
| J | jump | 3 |
| L | lazy | 1, 3, 5, 7 |
| M | men | 2, 4, 8 |
| N | now | 2, 6, 8 |
| O | over | 1, 3, 5, 7, 8 |
| P | party | 6, 8 |
| Q | quick | 1, 3 |
| TH | their | 1, 5, 7 |
| TI | time | 2, 4, 6 |

# Inverted Index



| Term Index | | Postings |
|---|---|---|
| AI | aid | 4:0.86, 8:0.22 |
| A — AL | all | 2:0.13, 4:0.59, 6:0.02 |
| B — BA | back | 1:0.37, 3:0.34, 7:0.93 |
| BR | brown | 1:0.79, 3:0.30, 5:0.19, 7:0.18 |
| C | come | 2:0.18, 4:0.56, 6:0.62, 8:0.02 |
| D | dog | 3:0.54, 5:0.13 |
| F | fox | 3:0.13, 5:0.09, 7:0.14 |
| G | good | 2:0.38, 4:0.85, 6:0.19, 8:0.01 |
| J | jump | 3:0.51 |
| L | lazy | 1:0.11, 3:0.17, 5:0.42, 7:0.84 |
| M | men | 2:0.15, 4:0.08, 8:0.66 |
| N | now | 2:0.82, 6:0.13, 8:0.38 |
| O | over | 1:0.22, 3:0.17, 5:0.37, 7:0.04, 8:0.11 |
| P | party | 6:0.49, 8:0.33 |
| Q | quick | 1:0.19, 3:0.54 |
| T — TH | their | 1:0.58, 5:0.17, 7:0.23 |
| TI | time | 2:0.83, 4:0.22, 6:0.28 |

# Two-Stage System

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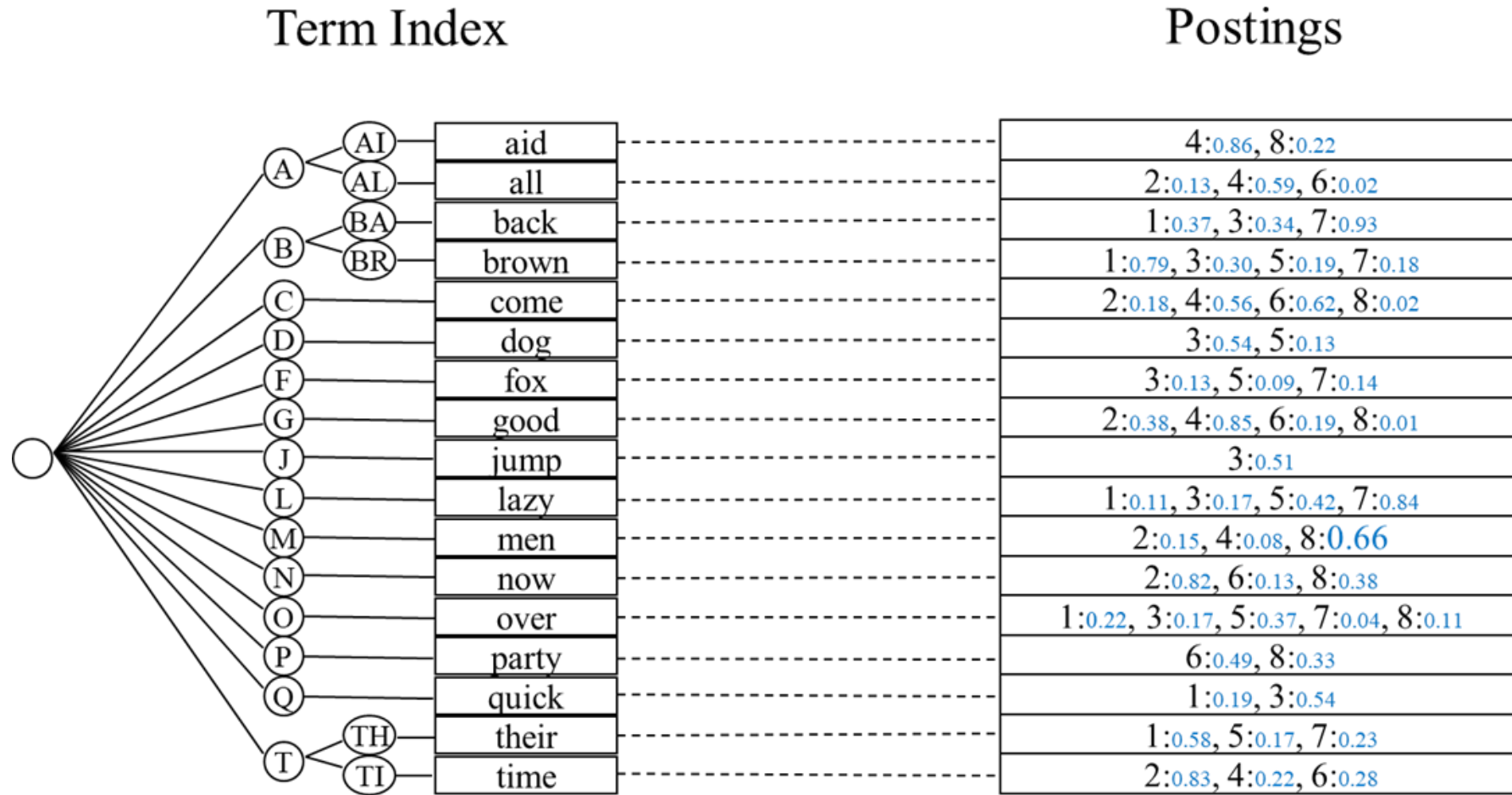  - Define retrieval unit
  - Tokenization
  - Build the inverted index

- 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)
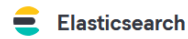
Apache Lucene™ 9.11.0 available (06.Jun)

**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

---

## 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    Insights

main    391 Branches    404 Tags    Go to file    <> Code

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 |

### About

Free and Open Source, Distributed, RESTful Search Engine

🔗 www.elastic.co/products/elasticsearch

java  search-engine  elasticsearch

📖 Readme

View license

Security policy

Activity

Custom properties

⭐ 69.7k stars

👁 2.7k watching

24.7k forks

Report repository

### Releases  155

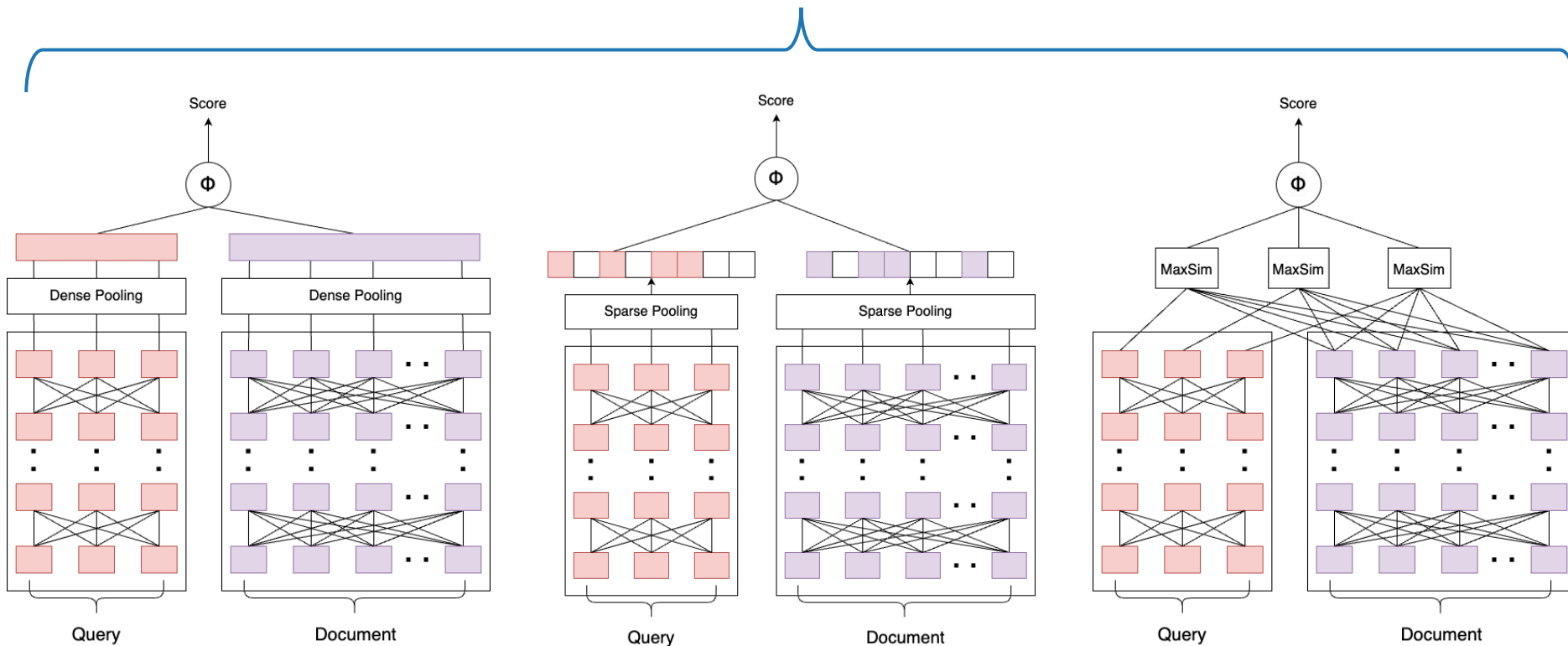Elasticsearch 8.15.1 (Latest)
3 weeks ago

+ 154 releases

### Packages

No packages published

# Can we go beyond surface forms?

neural language models

**Bi-Encoder**

**Cross Encoder**

Score

Score

Score

Score

Φ

Φ

Φ

Φ

Dense Pooling

Dense Pooling

Sparse Pooling

Sparse Pooling

MaxSim    MaxSim    MaxSim

Query    Document
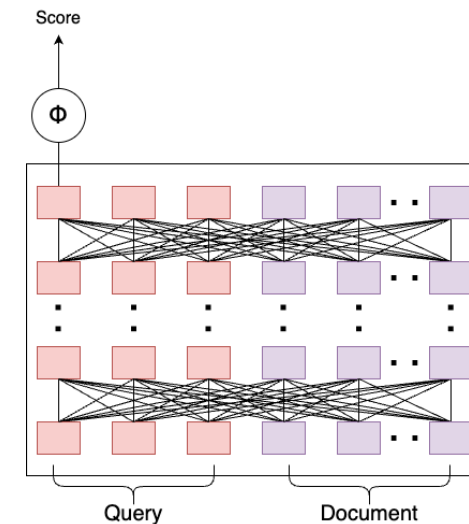
Query    Document

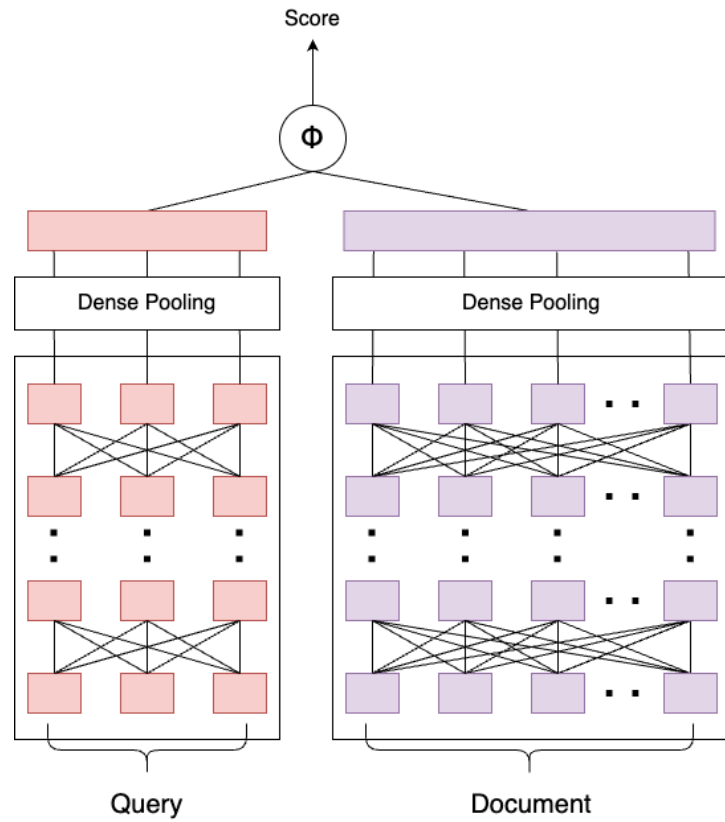Query    Document

Query    Document

One Dense Vector
Per Sequence
e.g., DPR

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

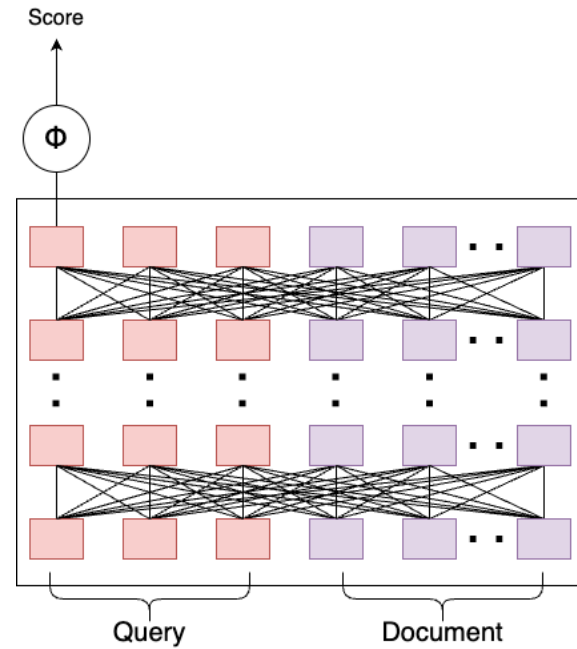Multiple Dense Vectors
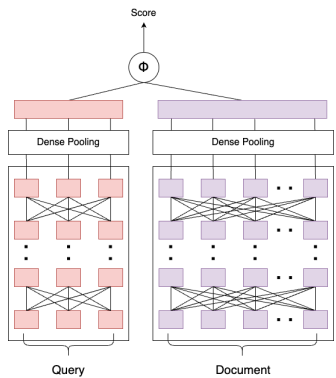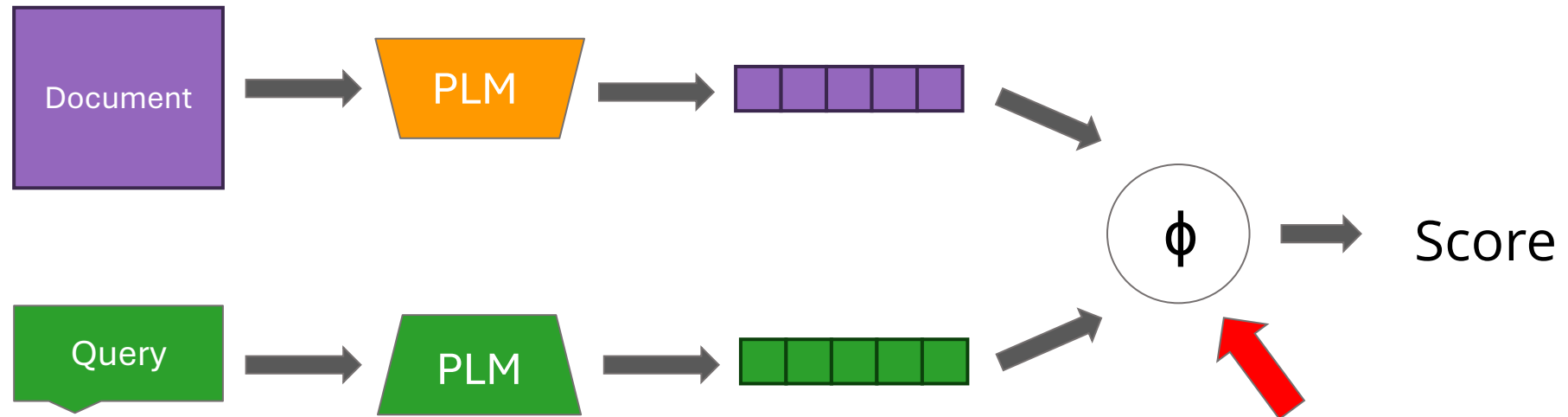Per Sequence
e.g., ColBERT

Joint Encoder
e.g., monoBERT

Separate query and document processing

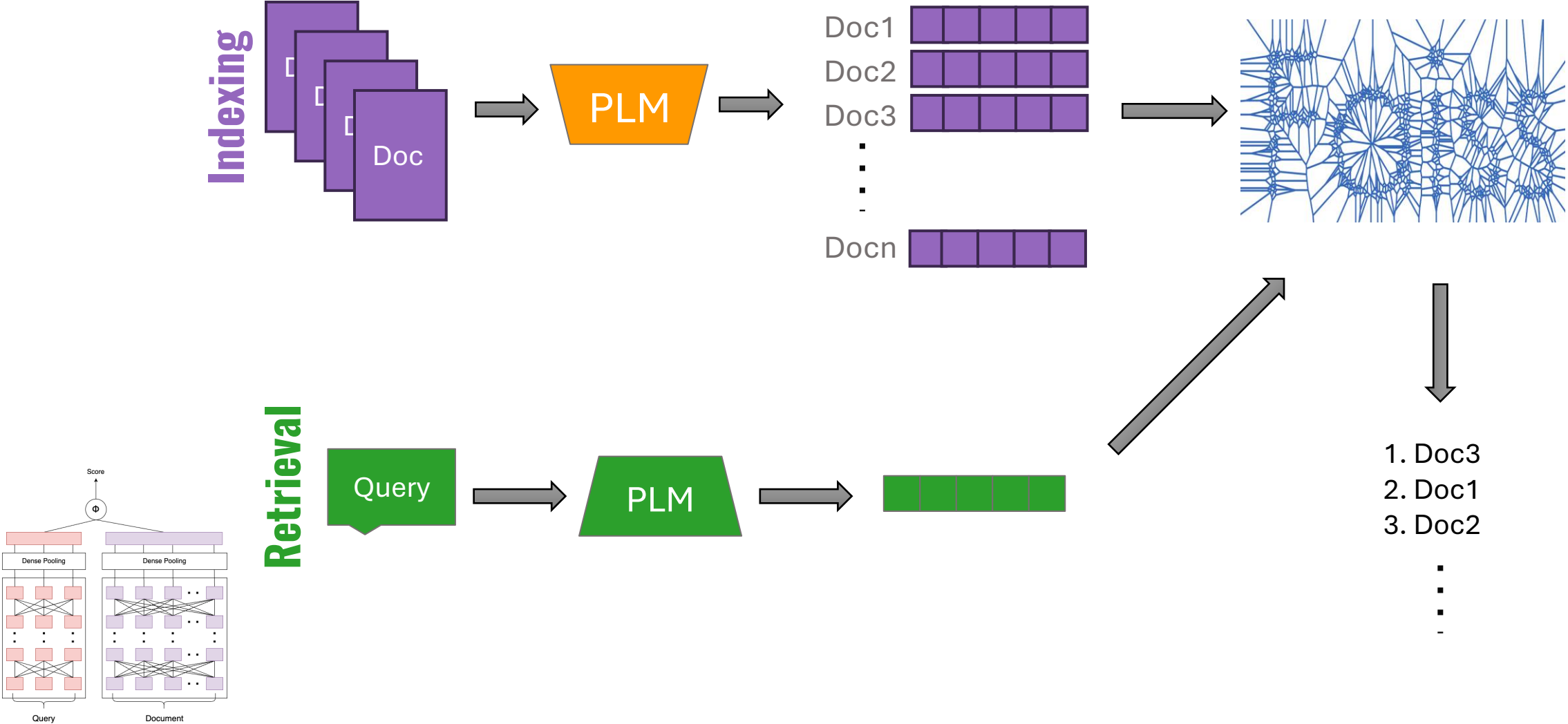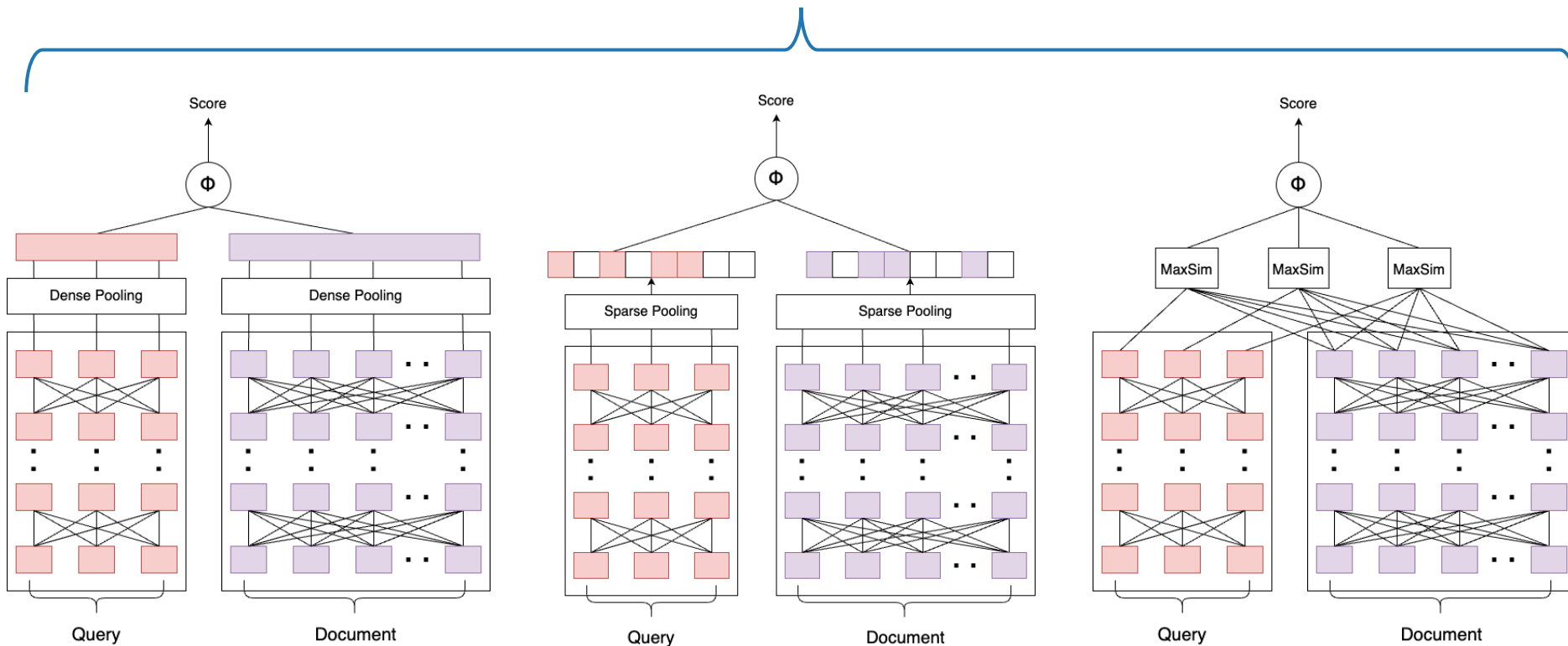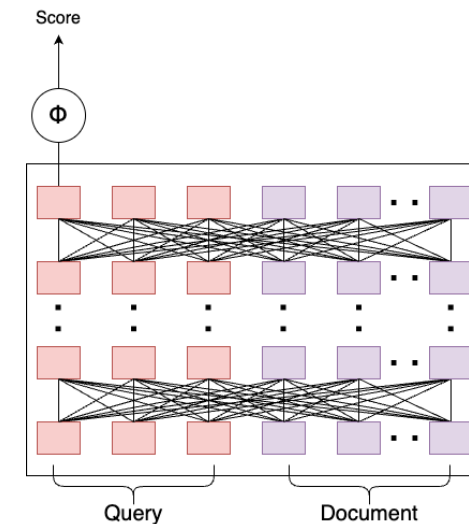# One Vector per Query, One Vector per Document

# Nearest Vectors aka Neighbors

$$\begin{bmatrix} 0.23 \\ 3.15 \\ 0.65 \\ 1.43 \end{bmatrix} \xrightarrow{\text{Search}} \boxed{x_1, x_2, \ldots, x_N} \rightarrow \boxed{\underset{n \in \{1,2,\ldots,N\}}{\text{argmin}} \|q - x_n\|_2^2} \rightarrow \begin{matrix} \text{Result} \\ \begin{bmatrix} 0.20 \\ 3.25 \\ 0.72 \\ 1.68 \end{bmatrix} \end{matrix}$$

$q \in \mathbb{R}^D$        $x_n \in \mathbb{R}^D$        $x_{74}$

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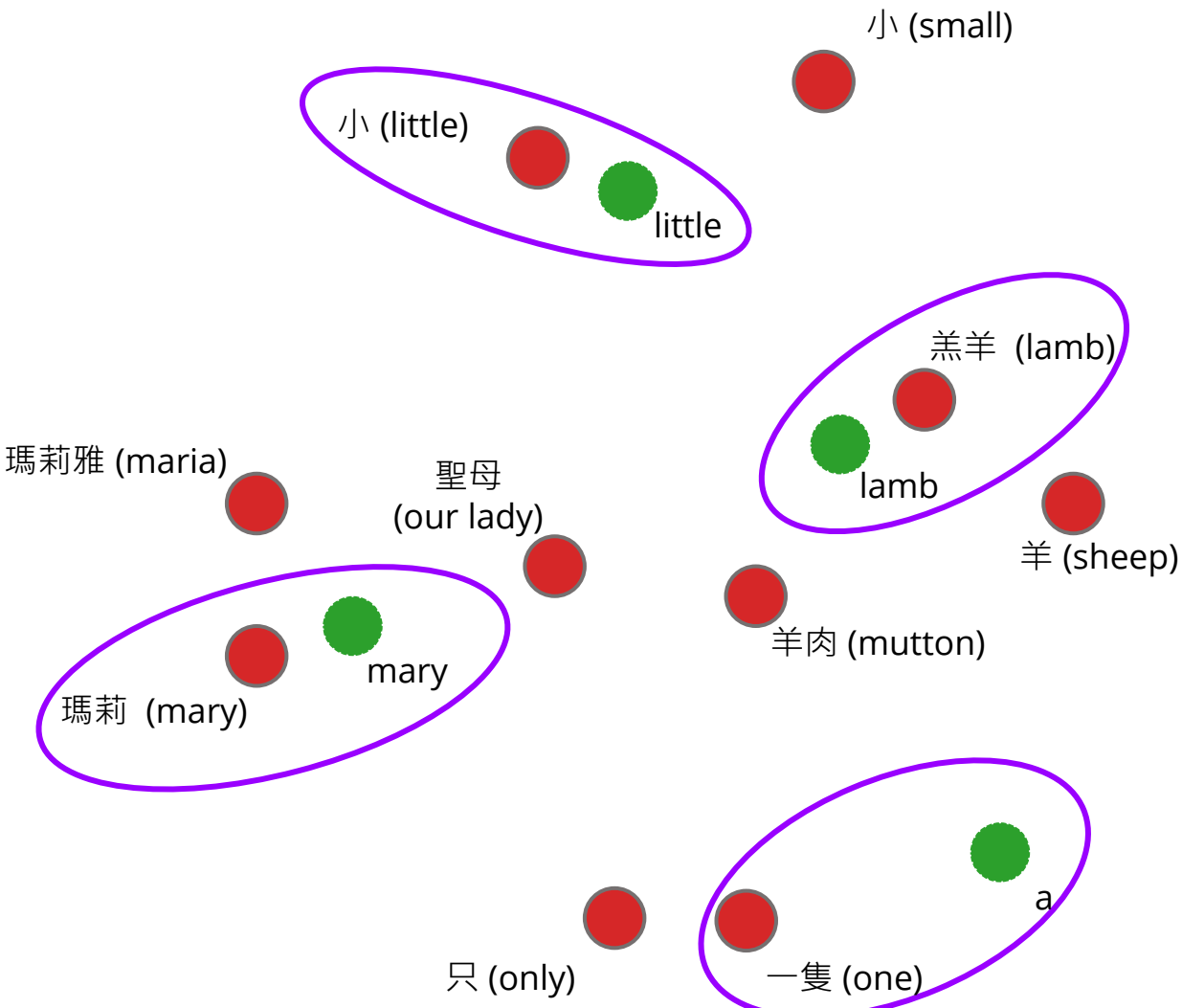
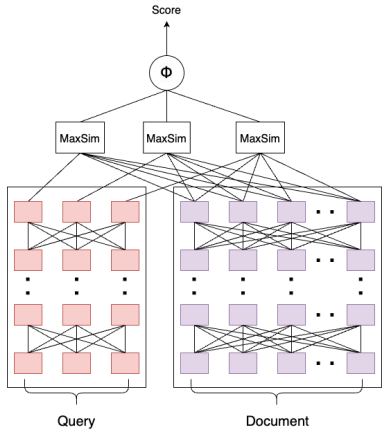Multiple Dense Vectors
Per Sequence
e.g., ColBERT

Joint Encoder
e.g., monoBERT

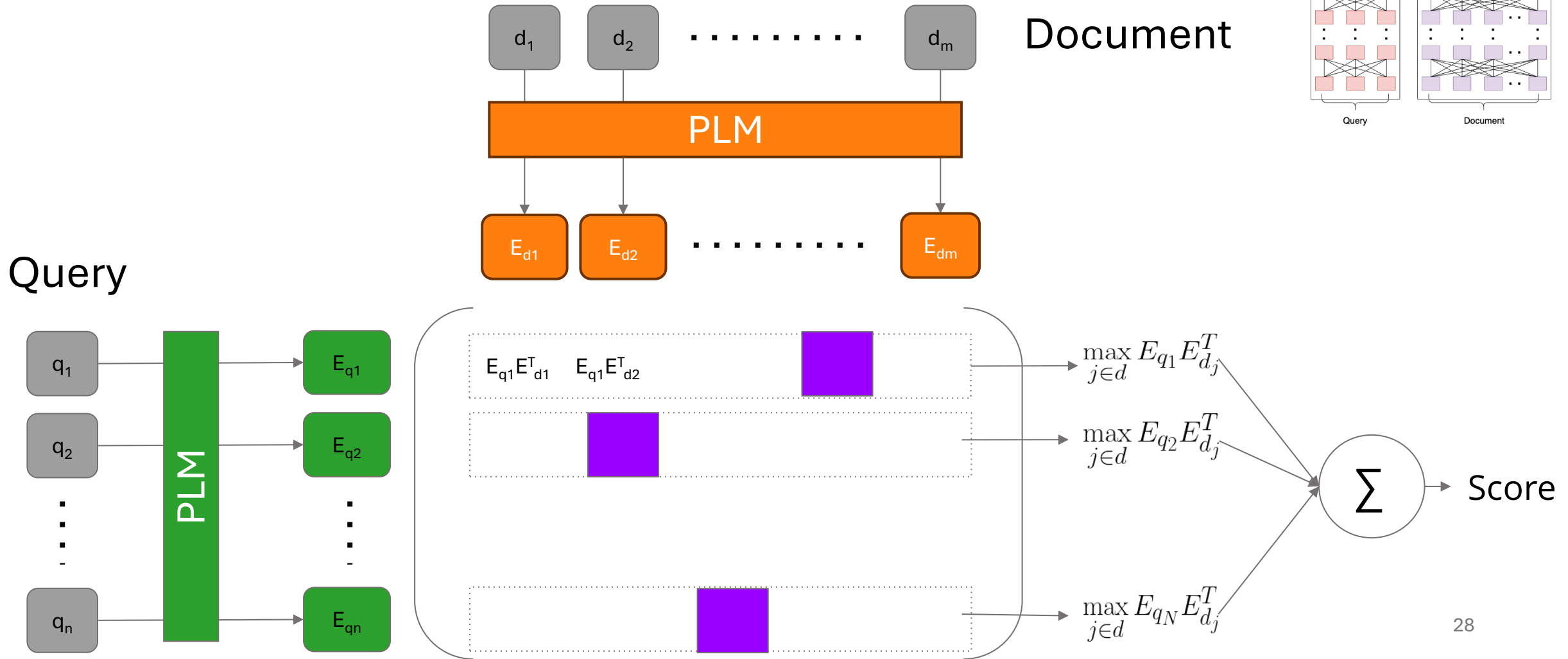# One Vector per Term: MaxSim



小 (small)

小 (little)

little

羔羊 (lamb)

lamb

羊 (sheep)

羊肉 (mutton)

有 (had)

had

瑪莉雅 (maria)

聖母 (our lady)

瑪莉 (mary)

mary

只 (only)

一隻 (one)

a

許多 (many)

Chinese Document term embedding

English Query term embedding

# MaxSim in Action -- ColBERT



Document

Query

$$\max_{j \in d} E_{q_1} E_{d_j}^T$$

$$\max_{j \in d} E_{q_2} E_{d_j}^T$$

$$\max_{j \in d} E_{q_N} E_{d_j}^T$$

$E_{q1}E^T_{d1}$  $E_{q1}E^T_{d2}$

$\sum$

Score

28

# One Vector per term: Multi-stage Retrieval



**Indexing**

Doc → PLM → Doc1, Doc2, Doc3 ... Docn → [graph] → **ANN** → Doc3 Doc5 Doc4 **Stage 1**

**Space Inefficient**

**Use all token embeddings**

**Retrieval**

Query → PLM → [green vectors]

**Stage 2** → MaxSim

1. Doc4
2. Doc3
3. Doc5

Score
Φ
MaxSim MaxSim MaxSim
Query Document

# Efficient PLAID Indexing Architecture

**Indexing**

Doc1

Doc2

Doc3

.
.
.

Docn

**K-means Clustering**

**PLAID Index**

$$t_1 = \text{centroid} + \text{residual}_1$$

$$\text{residual}_1 = \langle 1,0,0,0,0 \rangle$$

**Space Efficient**

**Retrieval**

English Query → PLM →

**ANN by retrieving closest clusters**

Doc3
Doc5
Doc4
.
.
.

**Stage 1**

**Decompress token embeddings through residuals**

**Stage 2** MaxSim →

1. Doc4
2. Doc3
3. Doc5
.
.
.

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



[CLS]
Mary
had
a
**[MASK]**
lamb

PLM

Masked
LM
Head

| 0.1% | good |
| ... | |
| 10% | little |
| ... | |
| 1% | small |
| ... | |

PLM Vocabulary

Score

Φ

Sparse Pooling | Sparse Pooling

Query | Document

# SPLADE

# SPLADE Search Pipeline

**Search...** | All fields | Search
Help | Advanced Search

## Computer Science > Information Retrieval

[Submitted on 23 Mar 2023]

# A Unified Framework for Learned Sparse Retrieval

Thong Nguyen, Sean MacAvaney, Andrew Yates

Learned sparse retrieval (LSR) is a family of first-stage retrieval
to generate sparse lexical representations of queries and docu
inverted index. Many LSR methods have been recently introdu
achieving state-of-the-art performance on MSMarco. Despite s
architectures, many LSR methods show substantial differences
efficiency. Differences in the experimental setups and configura
difficult to compare the methods and derive insights. In this wo
LSR methods and identify key components to establish an LSR
LSR methods under the same perspective. We then reproduce
using a common codebase and re-train them in the same envir
to quantify how components of the framework affect effectivene
that (1) including document term weighting is most important fo
effectiveness, (2) including query weighting has a small positiv
document expansion and query expansion have a cancellation

### Access Paper:

- View PDF
- TeX Source
- Other Formats

---

thongt99 / learned-sparse-retrieval (Public)

<> Code | Issues 4 | Pull requests | Actions | Projects | Security | Insights

main | Go to file | <> Code

Thong Nguyen raw float weights ✓ | d702026 · 7 months ago

| docs | update skeleton | last year |
| images | add logo | last year |
| lsr | raw float weights | 7 months ago |
| .gitignore | Initial commit | 2 years ago |
| LICENSE | Create LICENSE | last year |
| README.md | Add DOI | 7 months ago |
| beir.sh | Add beir to lsr | last year |
| clean.py | add file to clean beir trec file | last year |
| requirements.txt | Merge pull request #7 from ca... | last year |
| run_all_beir.sh | Add beir to lsr | last year |

### About

Unified Learned S
Framework

transformers  lsr
learned-sparse-retrie
sparse-retrieval

Readme
Apache-2.0 lice
Activity
57 stars
4 watching
5 forks

Report repository

### Releases 1

v1.0.0 Latest
on Feb 14

### Contributors 3

thongt99 Th
seanmacavan
cadurosar Ca

### Languages

### README  Apache-2.0 license

lsr  instructions  python 3.9.12  DOI 10.5281/zenodo.10659500

# LSR: A unified framework for efficient and effective learned sparse retrieval

---

TusKANNy / seismic (Public)

<> Code | Issues | Pull requests | Actions | Projects | Security | Insights

main | Go to file | <> Code

rossanoventurini Update README.md | 5efa741 · 2 months ago

| imgs | code | 3 months ago |
| scripts | update conversion script | 3 months ago |
| src | code | 3 months ago |
| .gitignore | code | 3 months ago |
| .pre-commit-config.yaml | code | 3 months ago |
| Cargo.toml | Update Cargo.toml | 2 months ago |
| LICENSE.md | code | 3 months ago |
| README.md | Update README.md | 2 months ago |
| pyproject.toml | code | 3 months ago |
| rust-toolchain.toml | code | 3 months ago |

### About

Official software repository of S.
Bruch, F. M. Nardini, C. Rulli, and
R. Venturini, "Efficient Inverted
Indexes for Approximate Retrieval
over Learned Sparse
Representations". Long Paper @
ACM SIGIR 2024 (Best Paper
Runner-up).

Readme
MIT license
Activity
Custom properties
38 stars
7 watching
1 fork

Report repository

### Releases 1

SIGIR2024 Latest
on Jul 4

### Packages

No packages published

### Contributors 3

rossanoventurini Rossano Ven...
francomarianardini Franco Ma...
CosimoRulli Cosimo Rulli

### README  MIT license

# Seismic

paper SIGIR 2024  arXiv 2404.18812

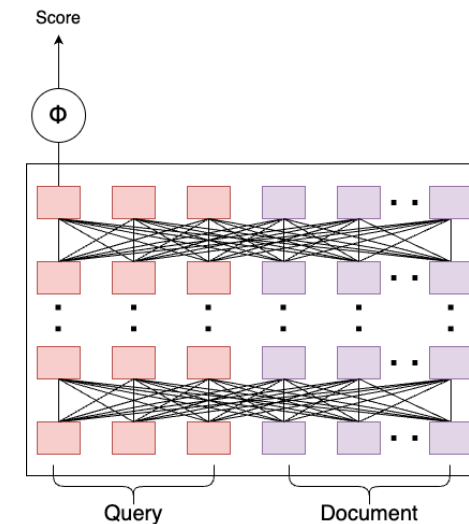crates.io v0.1.0  downloads 377  license MIT

Bi-Encoder

Cross Encoder

One Dense Vector Per Sequence e.g., DPR

One **Sparse** Vector Per Sequence e.g., SPLADE
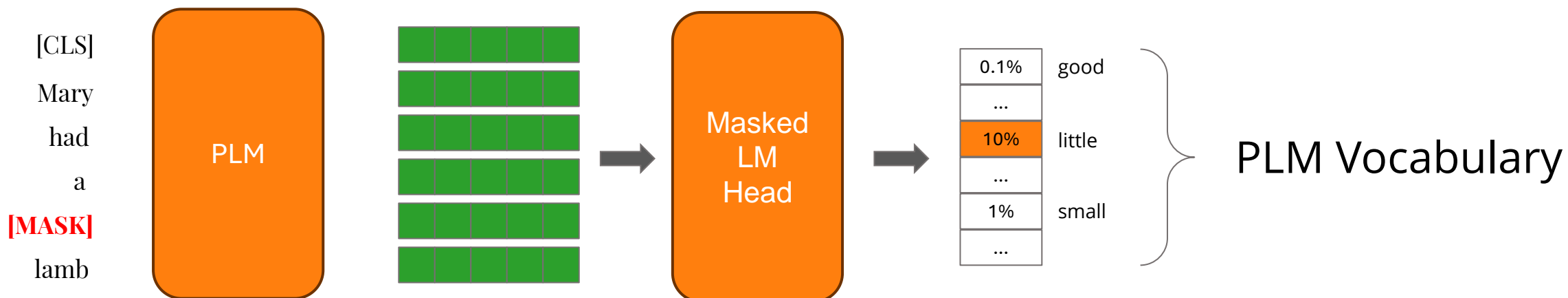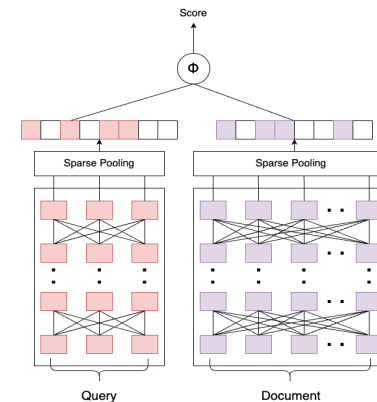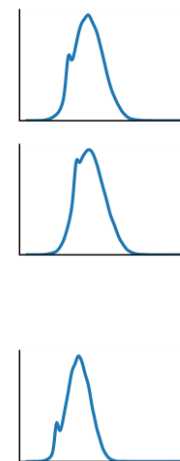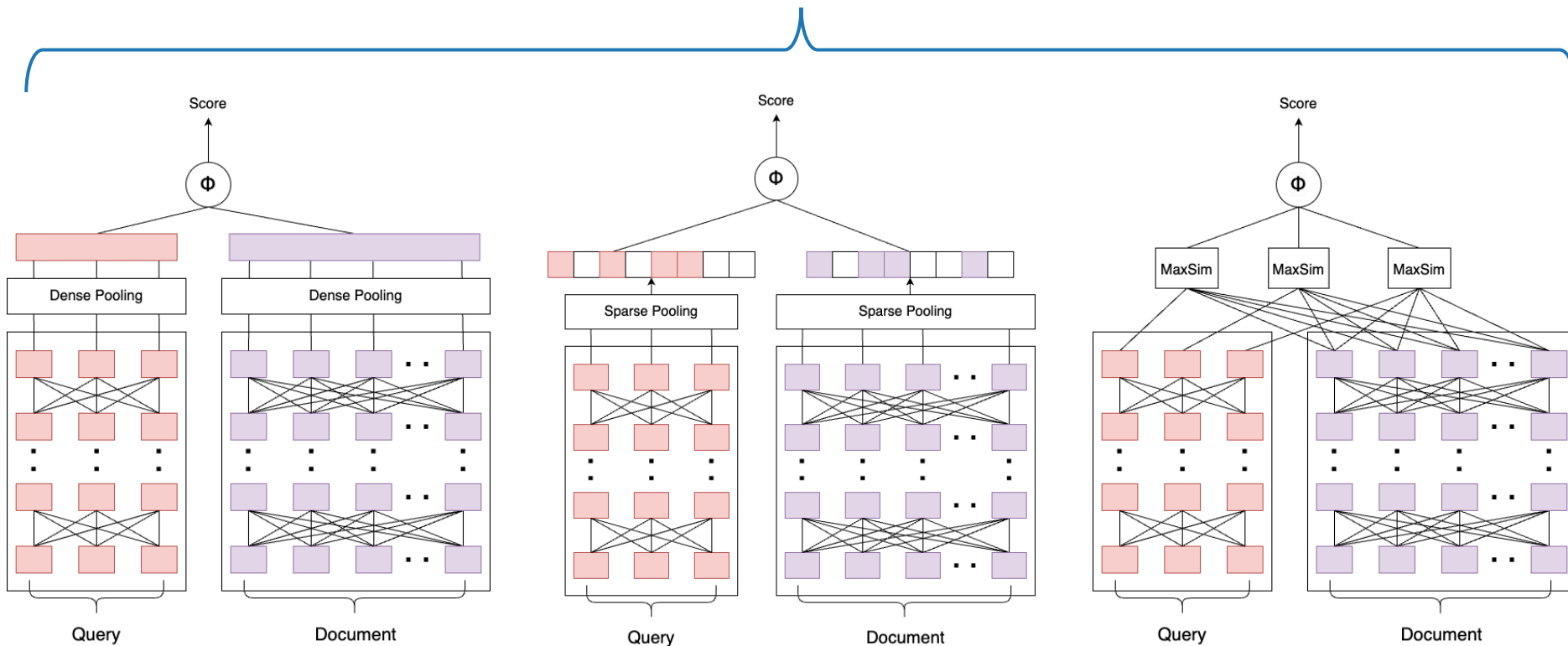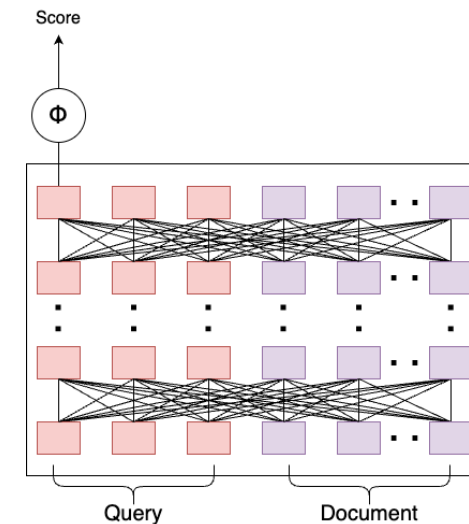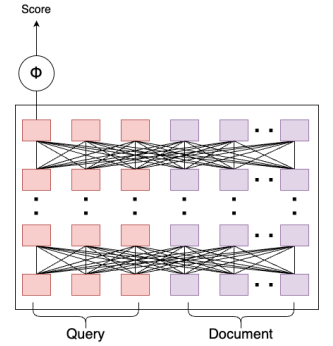
Multiple Dense Vectors Per Sequence e.g., ColBERT

Joint Encoder e.g., monoBERT

# Cross-Encoder

Query

Doc

Cross-Encoder
with PLM

0.9063

Score

Φ

Query          Document

# Using Generative Models



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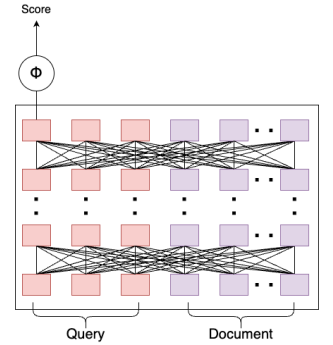Doc: Mary had a little lamb.
Relevant:

Generative
PLM

e.g, T5

**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 has
Doc: Mary had a little lamb.
Relevant:

Generative PLM

Yes (0.08)

No (0.001)

Yes

Score
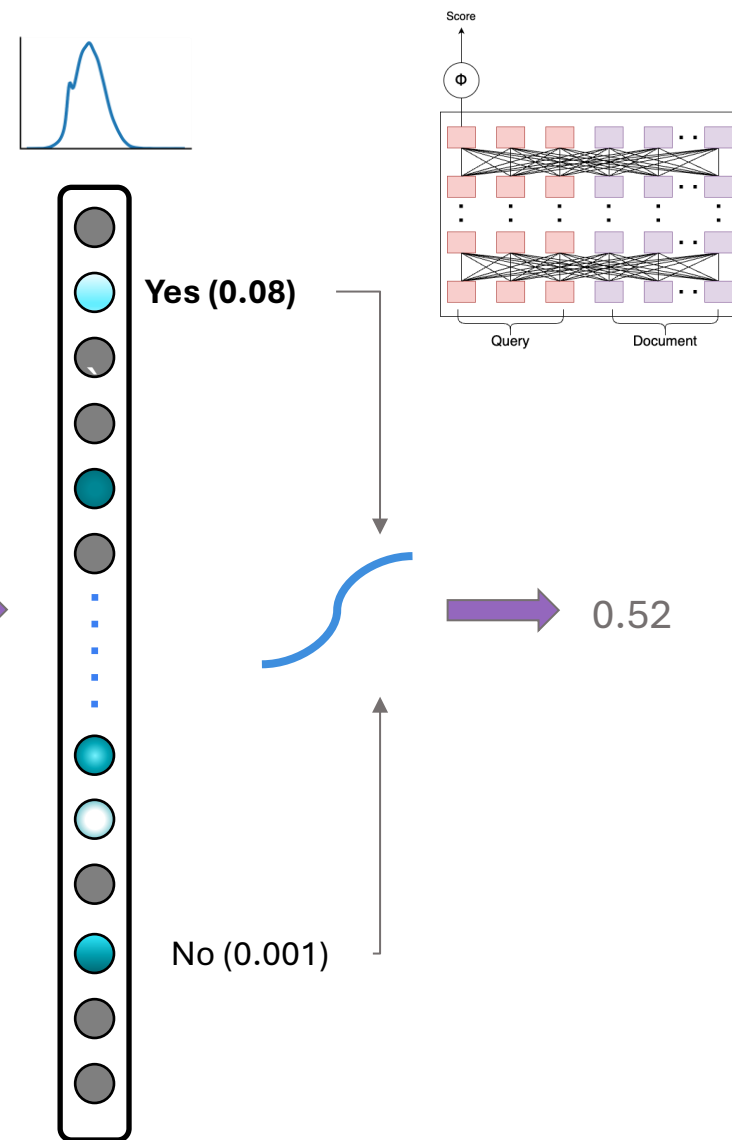
Φ

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 has
Doc: Mary had a little lamb.
Relevant:

Generative PLM

Yes (0.08)
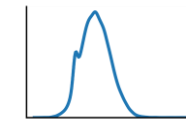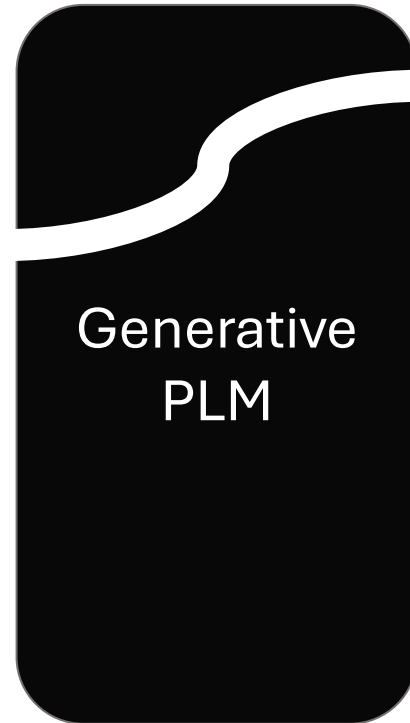
No (0.001)

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 has
Doc0: JHU is in Baltimore
Doc1: Mary had a little lamb.
Relevant:

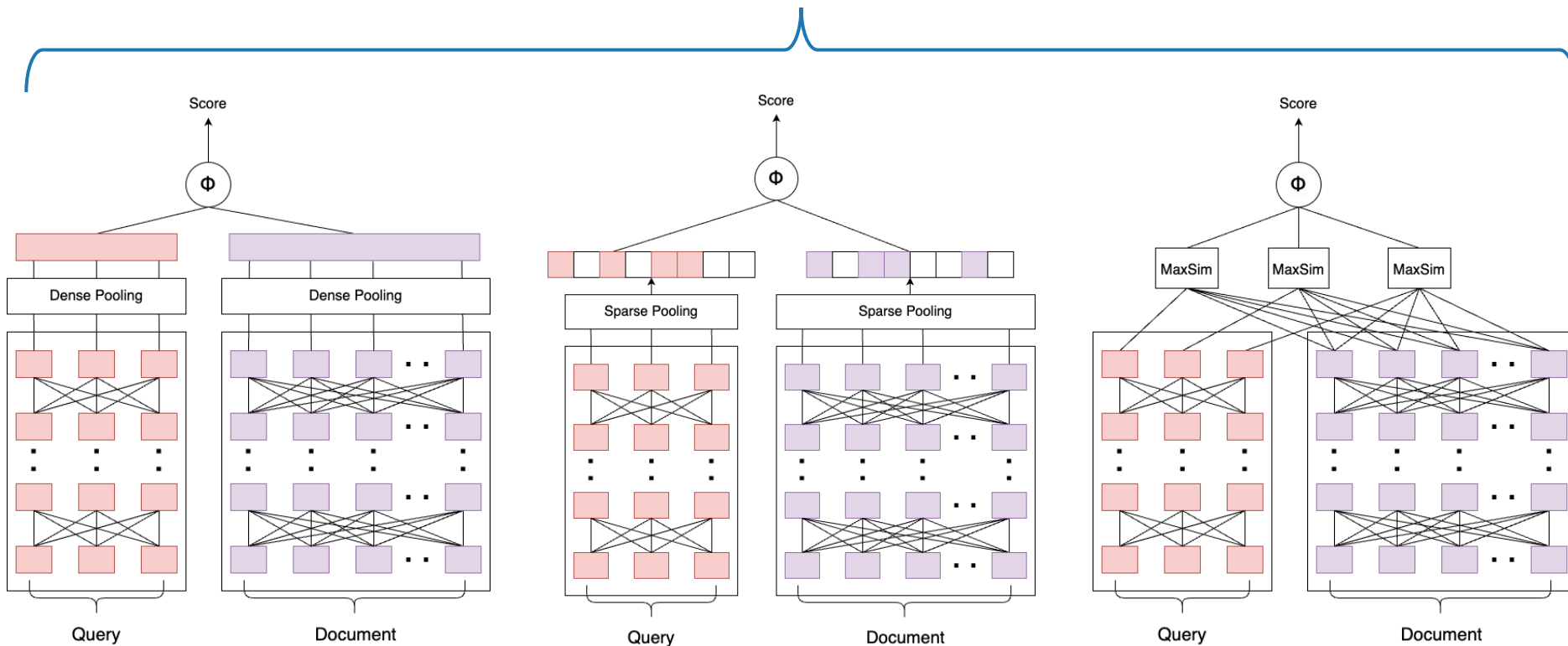Generative PLM

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).
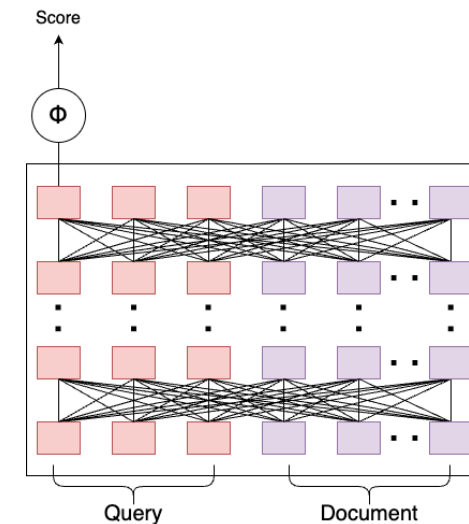
Bi-Encoder

Cross Encoder

One Dense Vector
Per Sequence
e.g., DPR

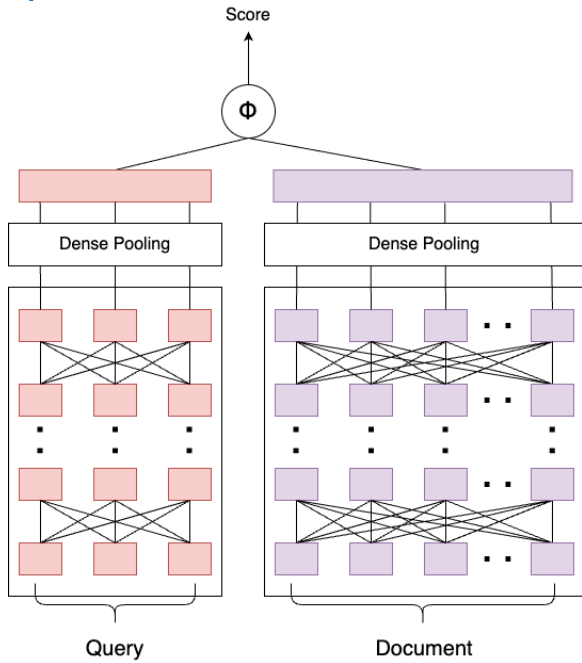One **Sparse** Vector
Per Sequence
e.g., SPLADE

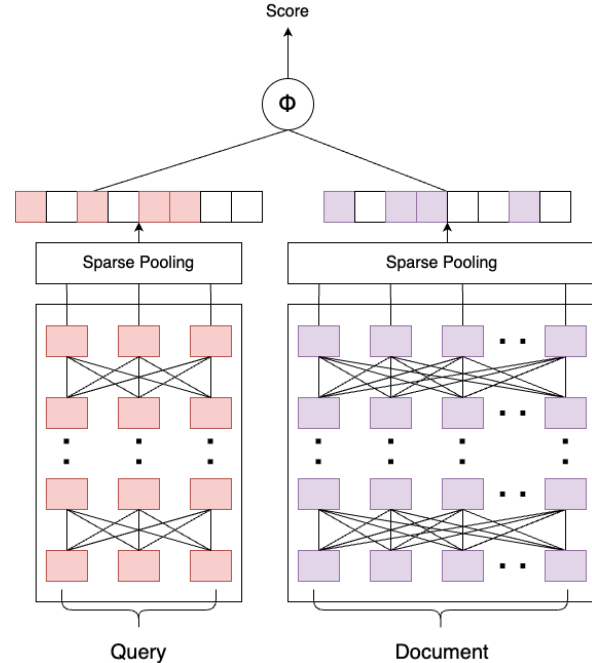Multiple Dense Vectors
Per Sequence
e.g., ColBERT

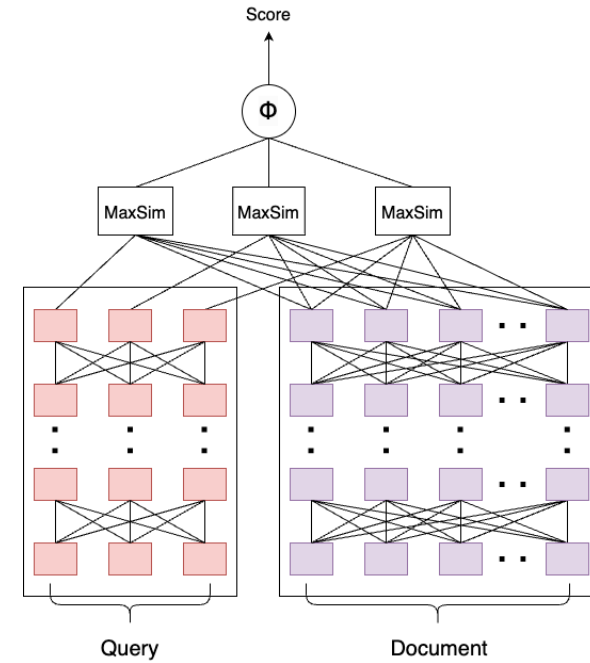Joint Encoder
e.g., monoBERT

Bi-Encoder

Cross Encoder

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

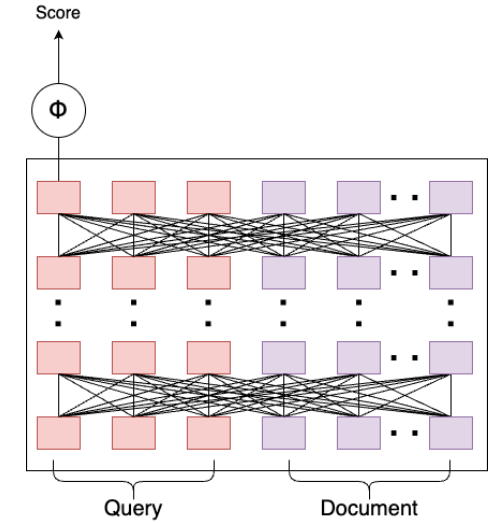**One Sparse Vector**
Per Sequence
e.g., SPLADE

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

**Joint Encoder**
e.g., monoBERT

**More Effective**

**More Efficient**

# Retrieve-and-Rerank System Combinations

More Efficient
Less Effective
→ **Higher Recall** →
More Effective
Less Efficient
→ **Final objective** →
Final Score

# Neural Retrieval System Pipeline

# PLM to IR Model

Pretraining → **Pretrained LM (PLM)** → Retrieval Finetuning → **IR Model**

- Align the representation

- Model "relevancy"
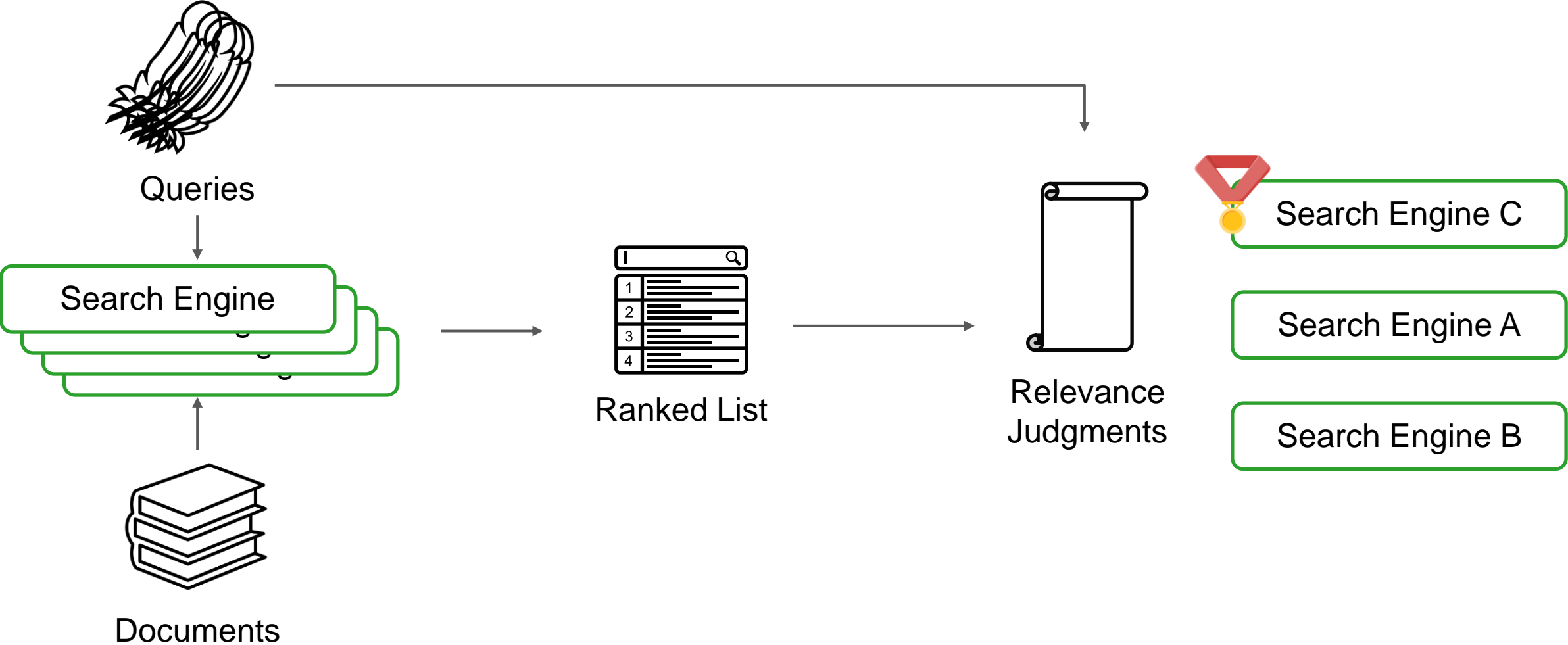
# Evaluation

Which system is better?
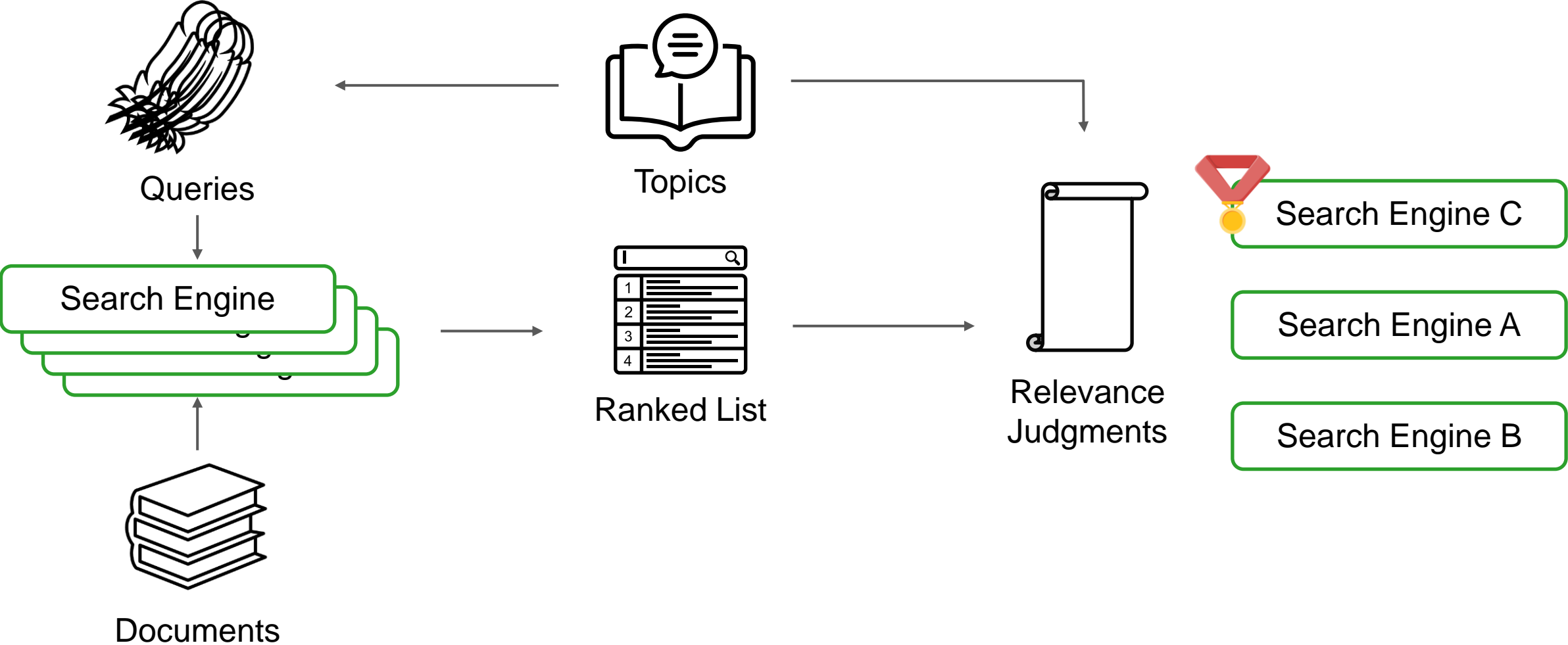
# What is Information Retrieval?

**(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

# Cranfield Paradigm Evaluation



One Query

One Topic

Search Engine

Documents

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
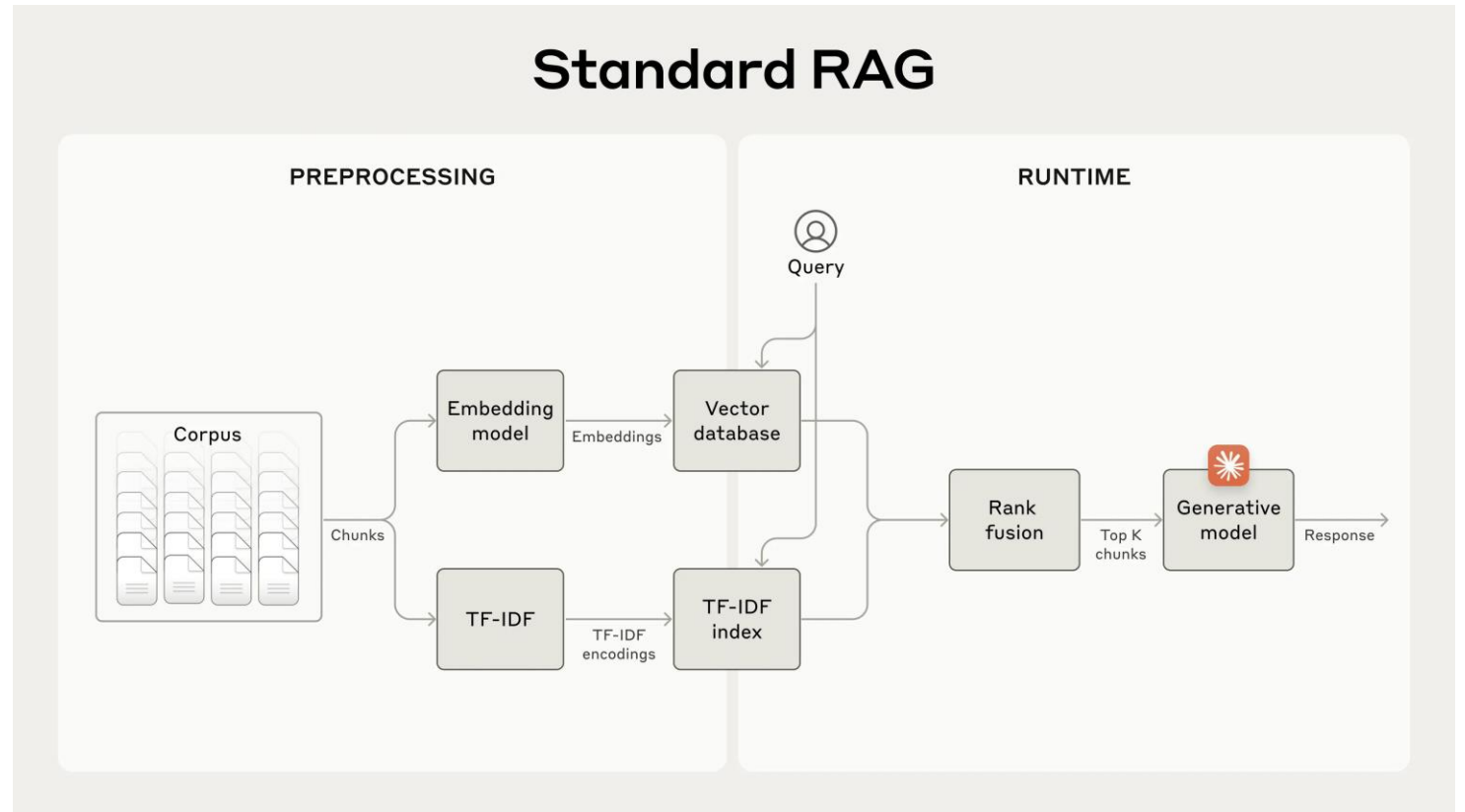  - Query latency (average search time per query)

# State of IR Research
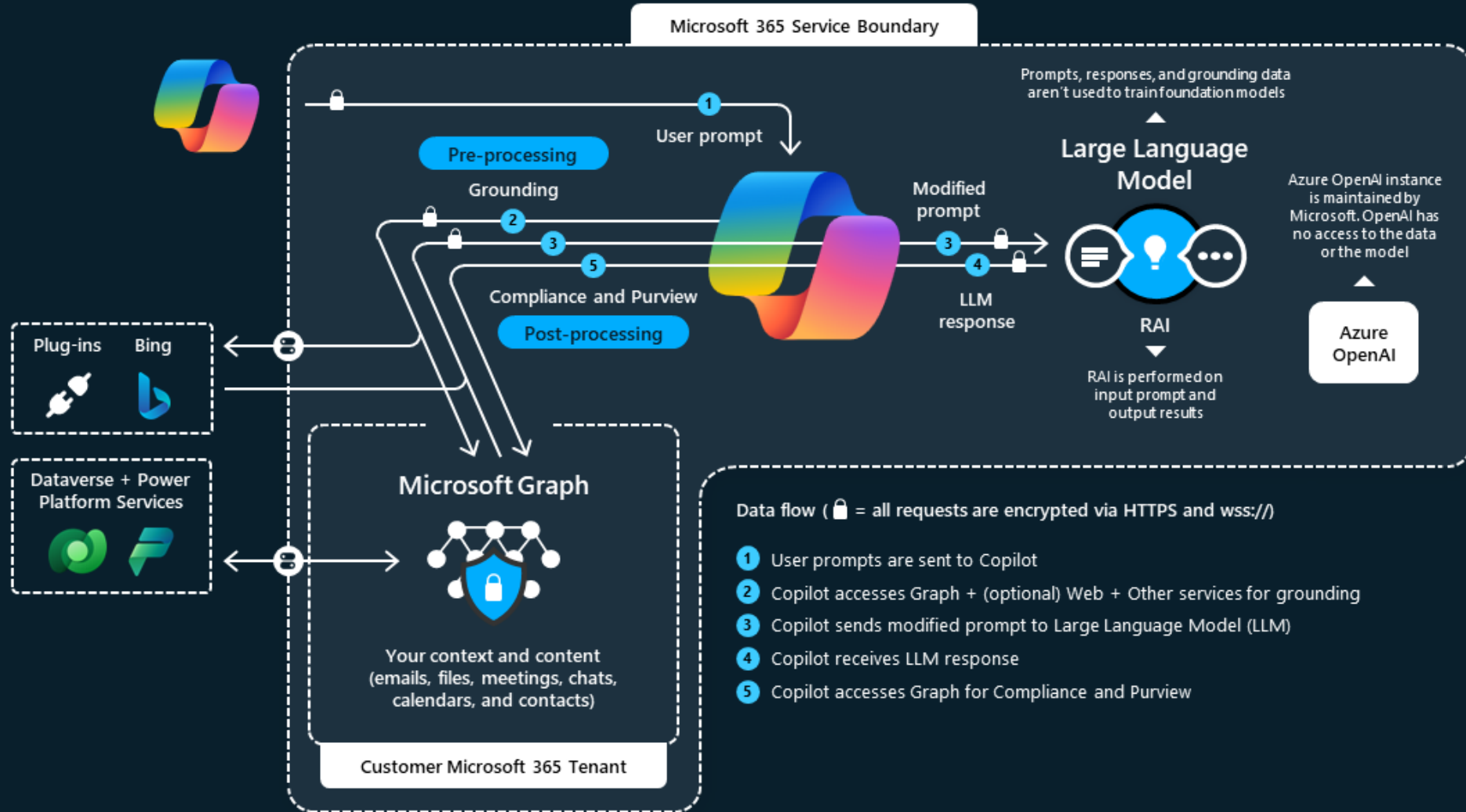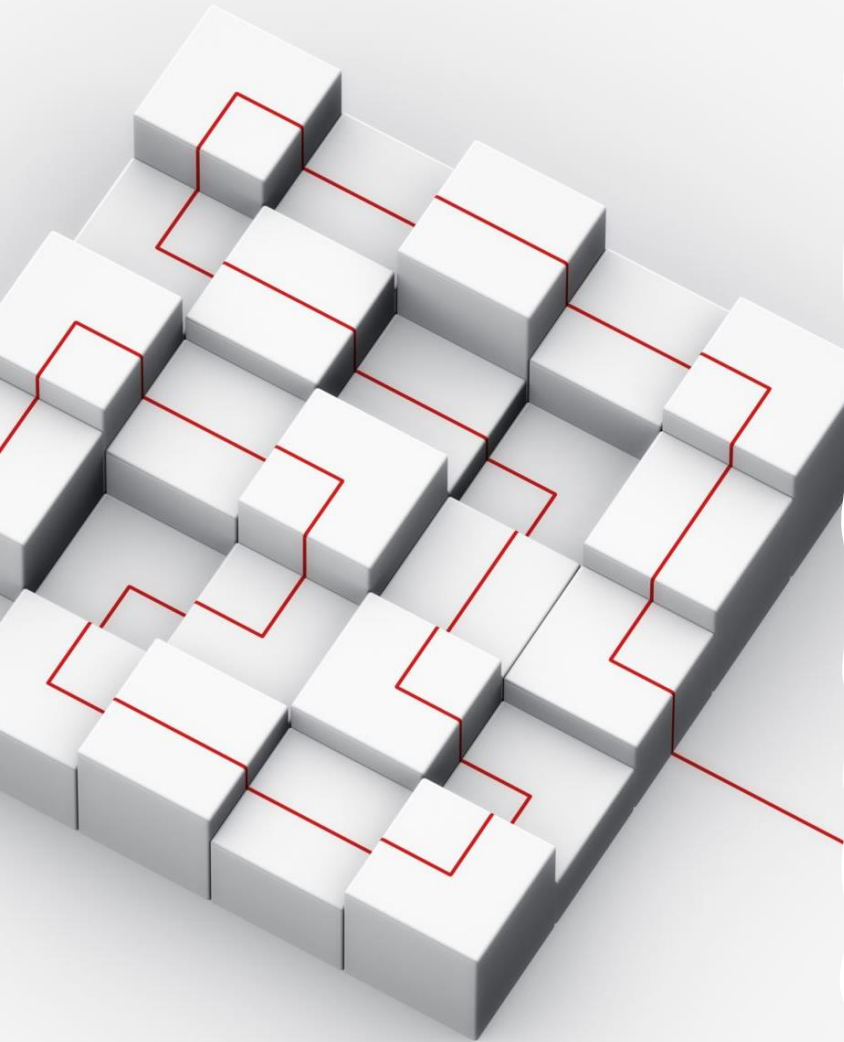
# Retrieval-Augmented Generation

- Is everything a RAG problem?
- What is the right retrieval model/system for RAG?
- IR going away?



https://www.anthropic.com/news/contextual-retrieval

Microsoft Copilot for Microsoft 365 architecture

https://learn.microsoft.com/en-us/copilot/microsoft-365/microsoft-365-copilot-overview

# 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