# Improve your RAG pipelines with semantic re-ranking

**Susan Shu Chang**

*Data Day Texas 2025*

# About the speaker

Principal data scientist at Elastic (Elasticsearch)

# About the speaker

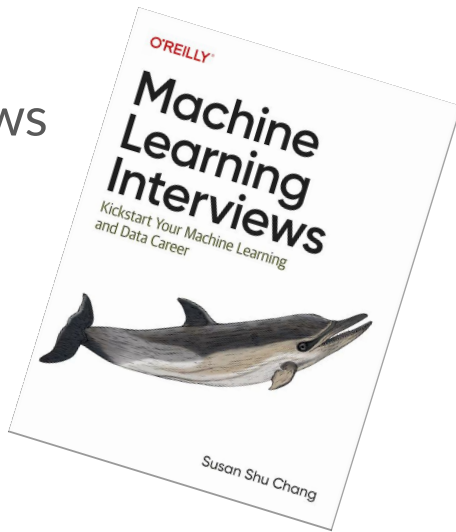Principal data scientist at Elastic (Elasticsearch)

6x PyCon speaker 🌎

# About the speaker

Principal data scientist at Elastic (Elasticsearch)

6x PyCon speaker 🌎

Author, Machine Learning Interviews

👋 Who travelled more than 3 hours to get here?

👋 Who travelled less than 3 hours to get here?

# Overview

1. Very quick primer on RAG

2. Search and recommendations: A story

3. Improving Retrieval: The "R" in RAG

4. 🥫 Secret sauce?! How rerankers are trained

5. All together: Rerankers in RAG systems
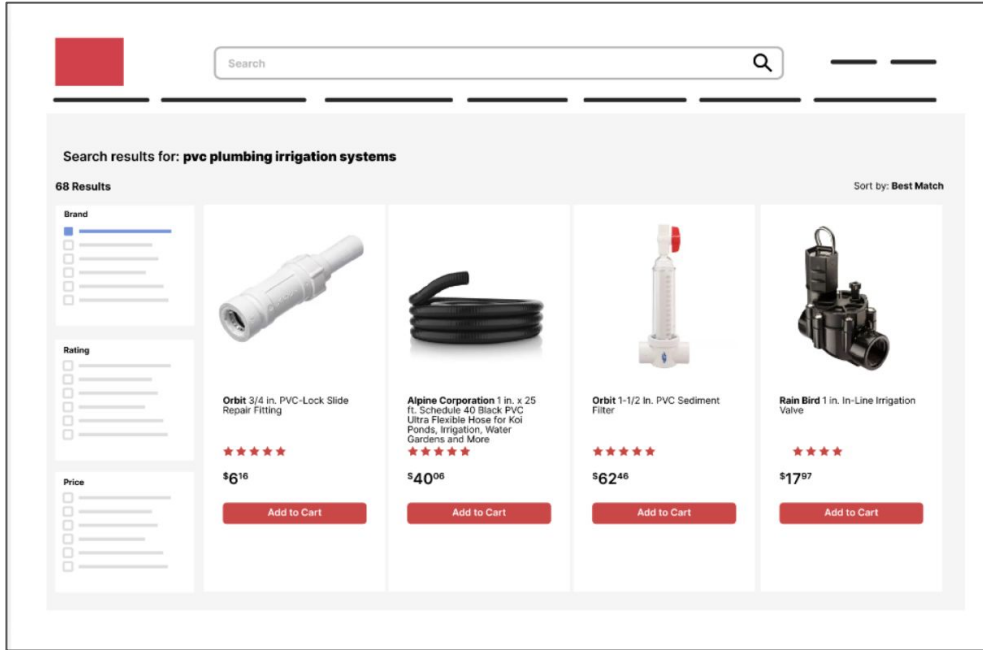
# Very quick primer on RAG

# Air Canada found liable for chatbot's bad advice on plane tickets

Airline's claim that online helper was responsible for its own actions was 'remarkable': small claims court

Jason Proctor · CBC News · Posted: Feb 15, 2024 3:38 PM EST | Last Updated: February 16, 2024

# RAG has become a standard to avoid hallucinations

# Accurate information enhances user experience



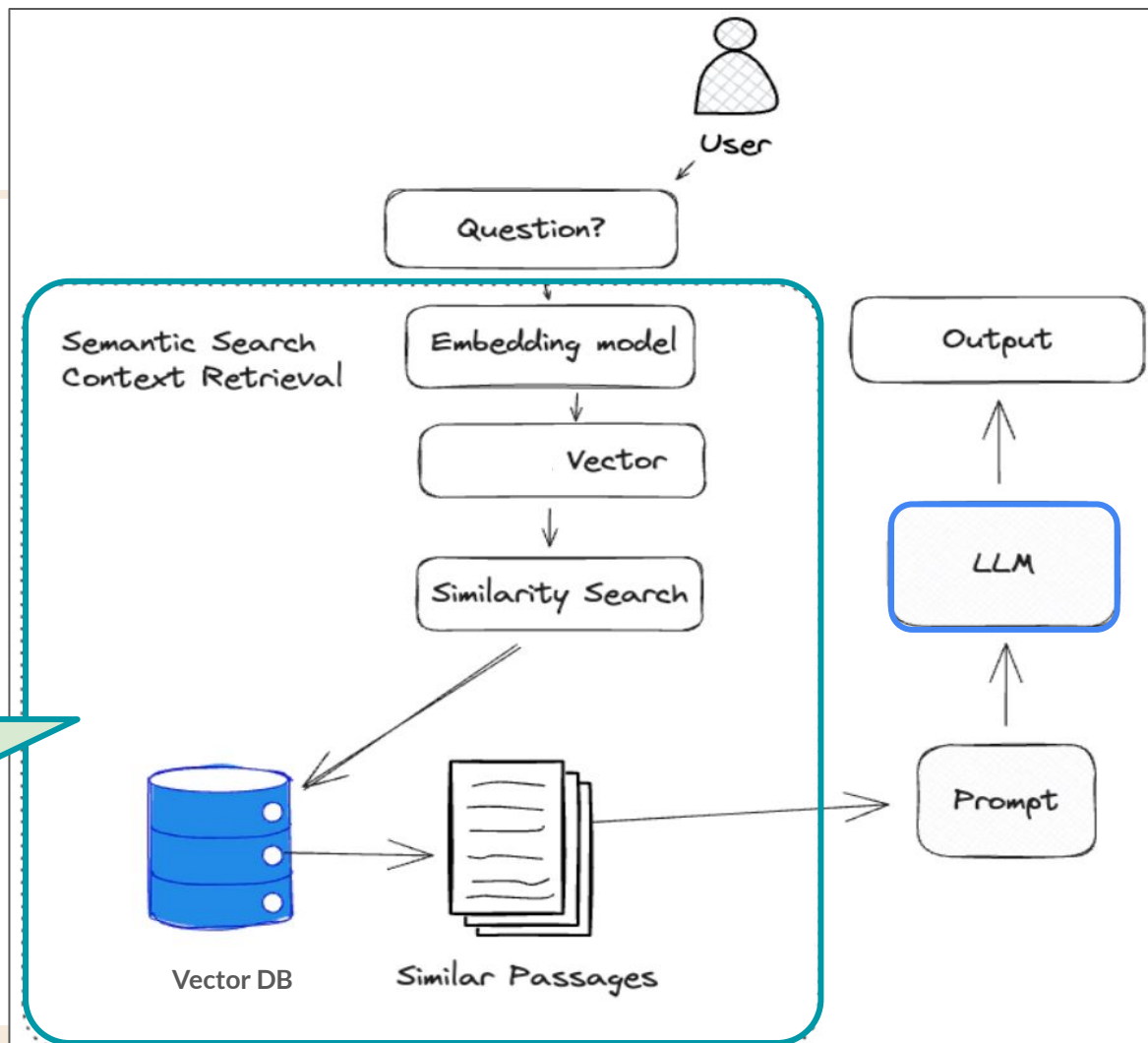Image source: Oxana Melis, via Unsplash

# RAG has become a standard to avoid hallucinations

Ideal behaviors: **Do not hallucinate or display**

- Products that we don't have

- Incorrect product descriptions
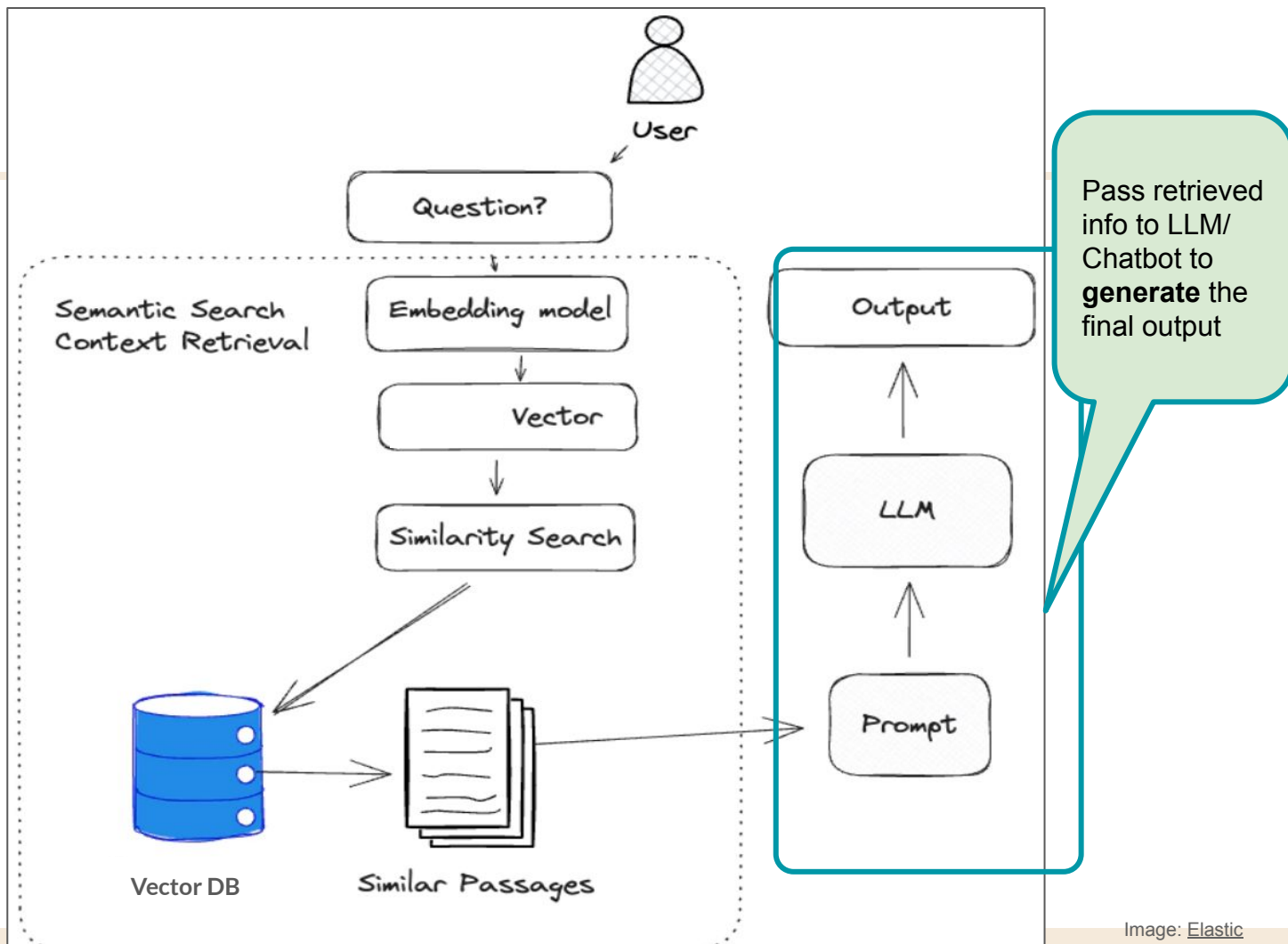
- Incorrect production functionality

# RAG is relatively **low-hanging fruit** to provide GenAI with accurate info



**Retrieve** the **most relevant** [support articles / product pages / discount policies] etc…

User

Question?

Semantic Search Context Retrieval

Embedding model

Vector

Similarity Search

Vector DB

Similar Passages

Output

LLM

Prompt

RAG is relatively **low-hanging fruit** to provide GenAI with accurate info

*...But you need to retrieve the right results*

User

Question?

Semantic Search Context Retrieval

Embedding model

Vector

Similarity Search

Vector DB

Similar Passages

Output

LLM

Prompt

Pass retrieved info to LLM/ Chatbot to **generate** the final output

Image: Elastic

2

Search and recommendations: A story

# Most search functions have at least *keyword based* search



Search results for: **pvc plumbing irrigation systems**

68 Results

Sort by: **Best Match**

Brand

Rating

Price

**Orbit** 3/4 in. PVC-Lock Slide Repair Fitting
★★★★★
$6^{16}$
Add to Cart

**Alpine Corporation** 1 in. x 25 ft. Schedule 40 Black PVC Ultra Flexible Hose for Koi Ponds, Irrigation, Water Gardens and More
★★★★★
$40^{06}$
Add to Cart

**Orbit** 1-1/2 In. PVC Sediment Filter
★★★★★
$62^{46}$
Add to Cart

**Rain Bird** 1 in. In-Line Irrigation Valve
★★★★
$17^{97}$
Add to Cart

**pvc plumbing irrigation systems**

**BM25** is a common lexical/keyword based function

# Many ways of tuning keyword search

Query rules

    roof gutter

    roofgutter

    foof guter
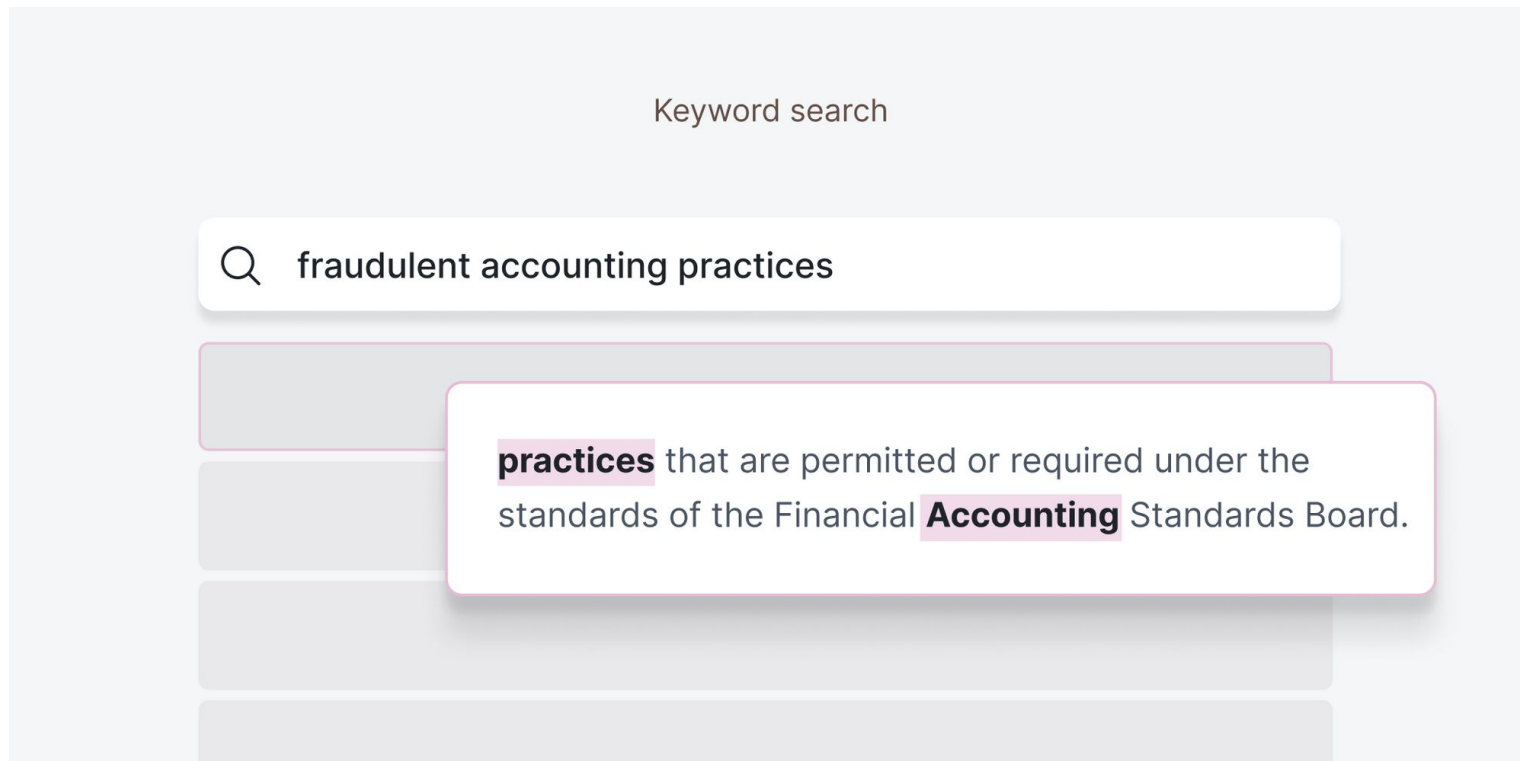
| Type | Match Requirements |
|------|--------------------|
| `exact` | Rule metadata matches the specified value exactly. |
| `fuzzy` | Rule metadata matches the specified value within an allowed Levenshtein edit distance. |
| `prefix` | Rule metadata starts with the specified value. |
| `suffix` | Rule metadata ends with the specified value. |
| `contains` | Rule metadata contains the specified value. |
| `lt` | Rule metadata is less than the specified value. |
| `lte` | Rule metadata is less than or equal to the specified value. |

...and much more

# Bag-of-words / keyword based representations can't account for semantics

Keyword search

🔍 fraudulent accounting practices

**practices** that are permitted or required under the standards of the Financial **Accounting** Standards Board.

Image: DeepJudge

# Language tasks evolve to understand context/semantics



Semantic search

🔍 fraudulent accounting practices

FW: Letter to Enron's Chairman after Departure of CEO

I am incredibly nervous that we will **implode in a wave of accounting scandals**. My eight years of Enron work history will be worth nothing on my resume, …

Image: DeepJudge

# Modern search and retrieval utilizes semantic and lexical search

| Lexical/keyword search: BM25 | Semantic search: Sparse vector, dense vector | Hybrid search: RRF etc. |
|---|---|---|

# Semantic search: Neural/ML-learned vectors capture meaning

- Learns semantics, not *just* keyword based

- Can be fine-tuned to specific tasks and domains

- Use case scalability
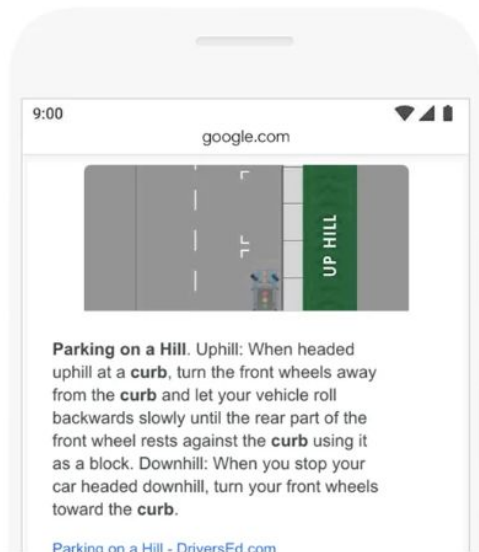
For more: see my Data Day Texas 2024 talk ([link](link))

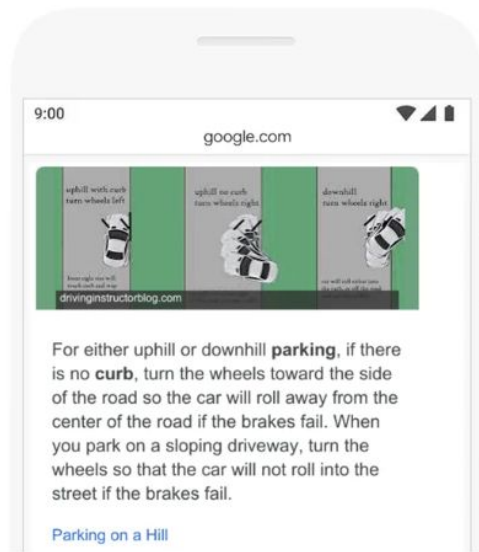# You've been interacting with ML-powered semantic search for years

parking on a hill with no curb

**BEFORE**

9:00 google.com

UP HILL

**Parking on a Hill**. Uphill: When headed uphill at a **curb**, turn the front wheels away from the **curb** and let your vehicle roll backwards slowly until the rear part of the front wheel rests against the **curb** using it as a block. Downhill: When you stop your car headed downhill, turn your front wheels toward the **curb**.

Parking on a Hill - DriversEd.com

**AFTER**

9:00 google.com

uphill with curb turn wheels left
uphill no curb turn wheels right
downhill turn wheels right
drivinginstructorblog.com

For either uphill or downhill **parking**, if there is no **curb**, turn the wheels toward the side of the road so the car will roll away from the center of the road if the brakes fail. When you park on a sloping driveway, turn the wheels so that the car will not roll into the street if the brakes fail.

Parking on a Hill

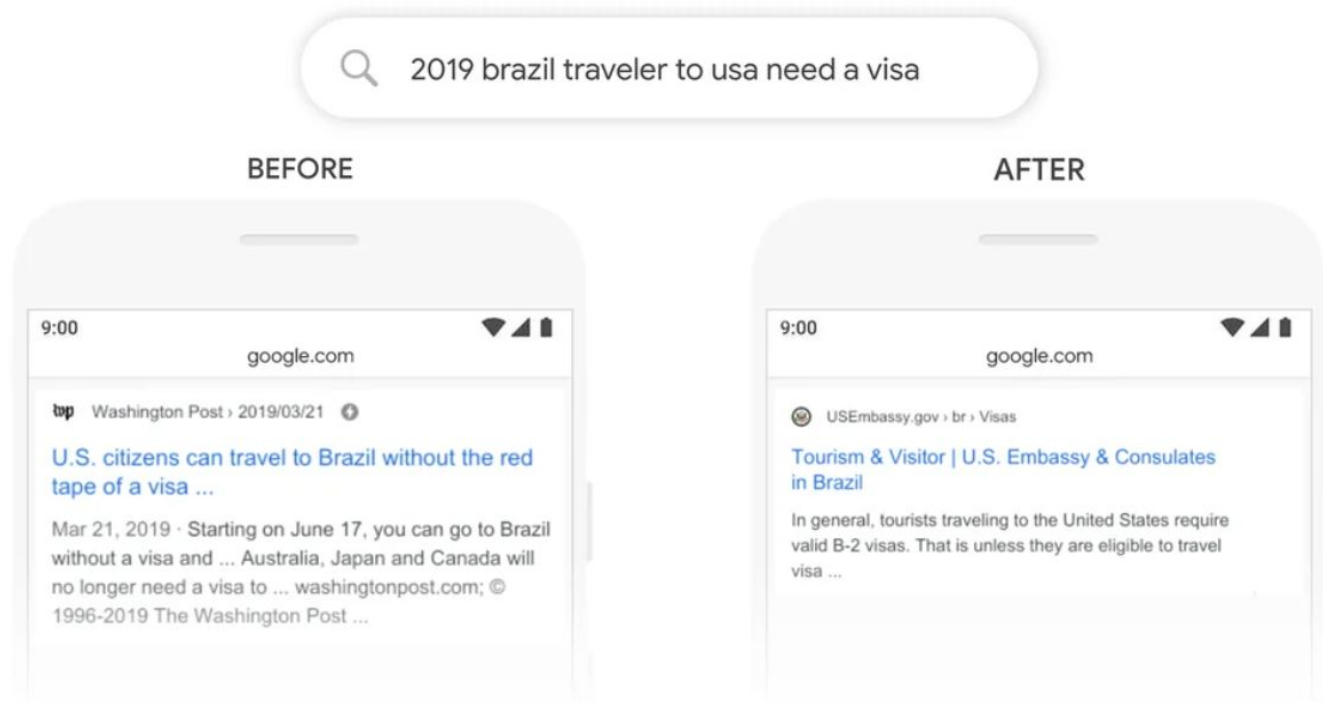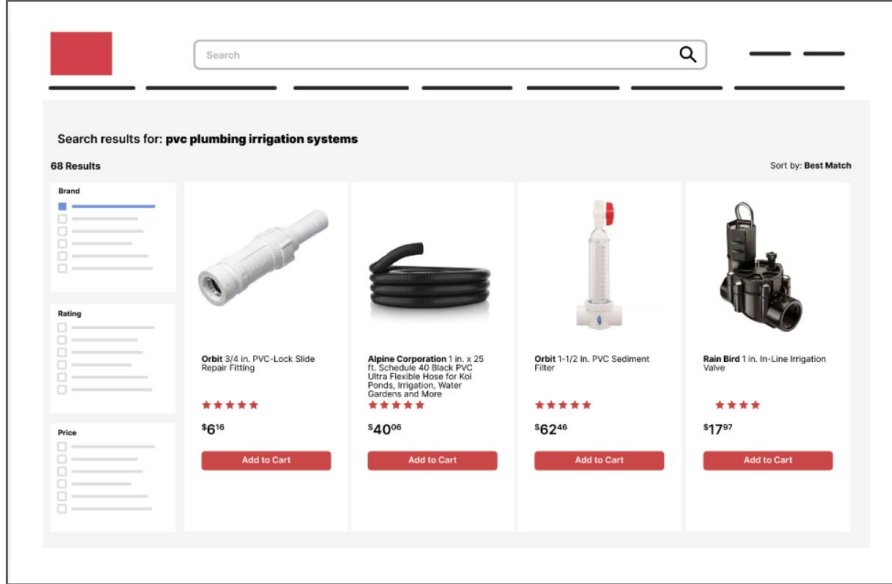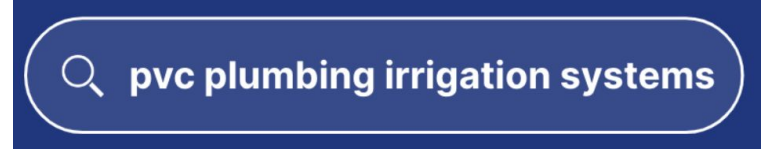Powered by foundational language models e.g. BERT since 2019

Image: Google

# You've been interacting with ML-powered semantic search for years



Image: Google

Example: Build *new* chat function on product site with semantic search

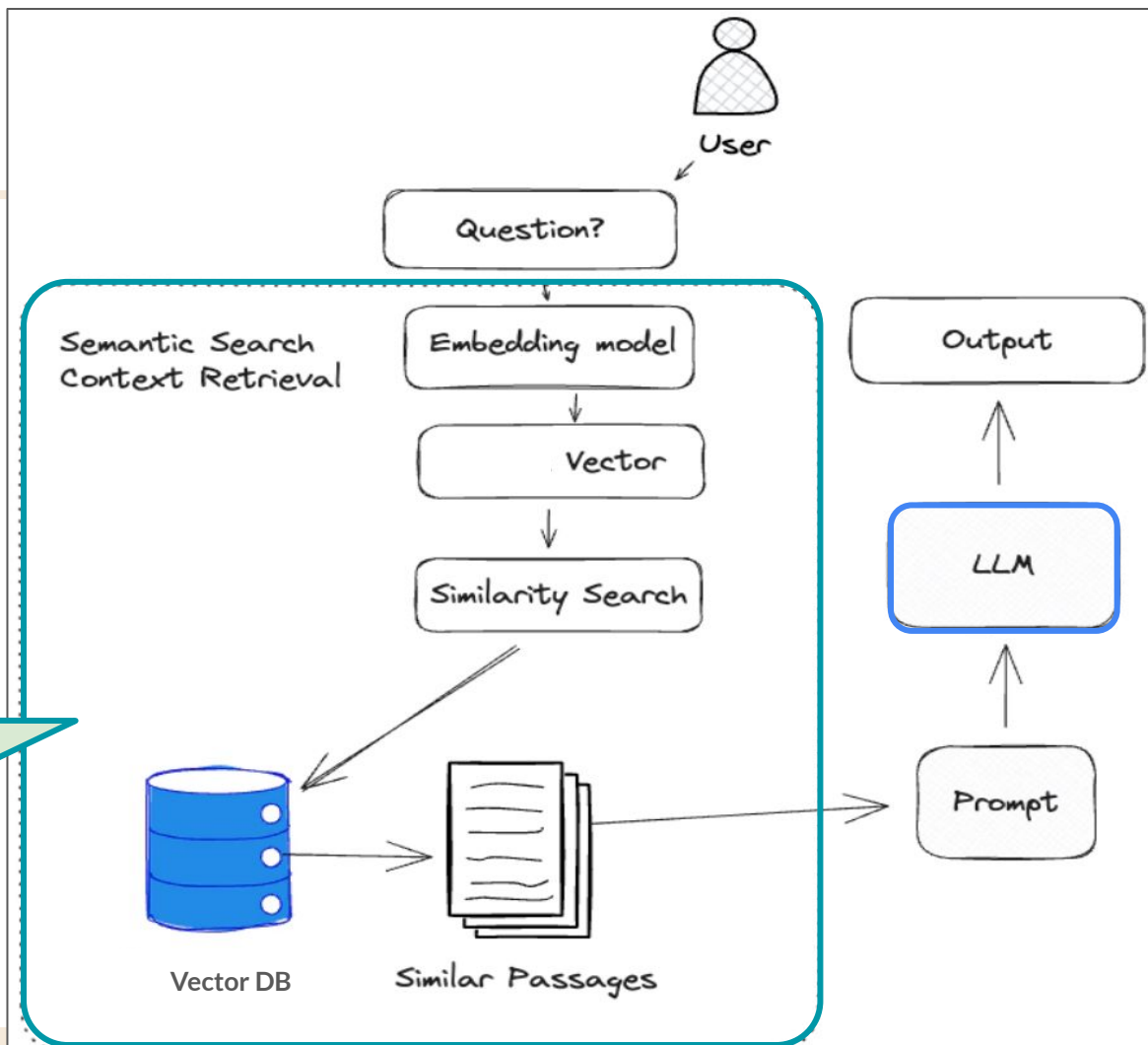(In addition to existing search bar)

Mission: go from *this*

To *this*
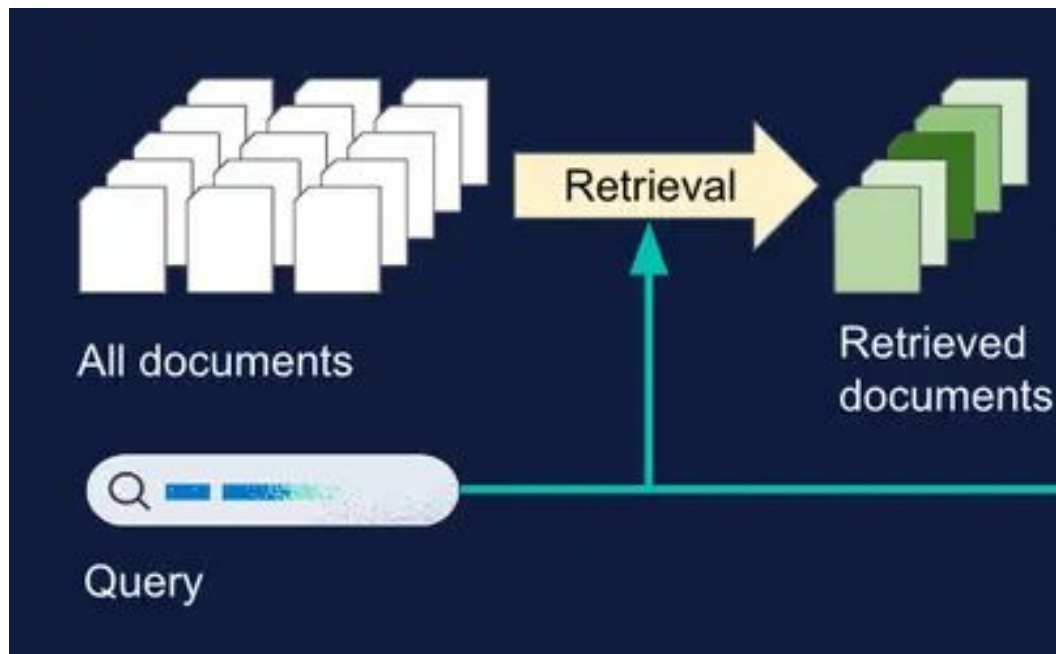
3

# Improving Retrieval: The "R" in RAG

RAG is relatively **low-hanging fruit** to provide GenAI with accurate info



User

Question?

Semantic Search
Context Retrieval

Embedding model

Vector

Similarity Search

Output

LLM

Prompt

**Retrieve** the **most relevant** [support articles / product pages / discount policies] etc…

Vector DB

Similar Passages

"I've implemented RAG, why is my LLM still outputting irrelevant answers?"

# But, **retrieving** initial results is only the first step



Retrieval

All documents
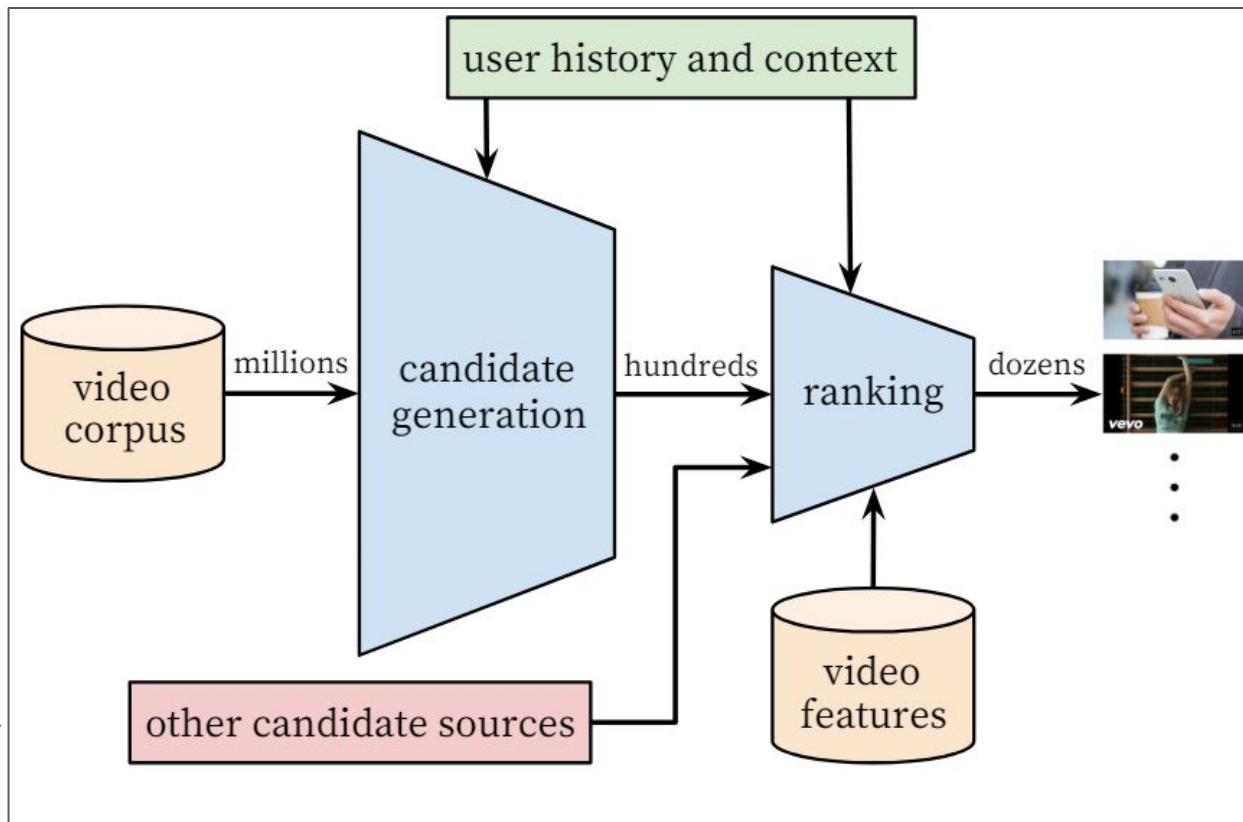
Retrieved documents

Query

IN PRODUCTION

Just because some items are retrieved, doesn't mean they're displayed in the optimal order for the LLM / user

# But, **retrieving** initial results is only the first step



| All documents | Retrieval → | Retrieved documents | Reranking → | Reranked documents |

Query

**Rerank** retrieved "candidates" by importance
to pass the LLM or display to user

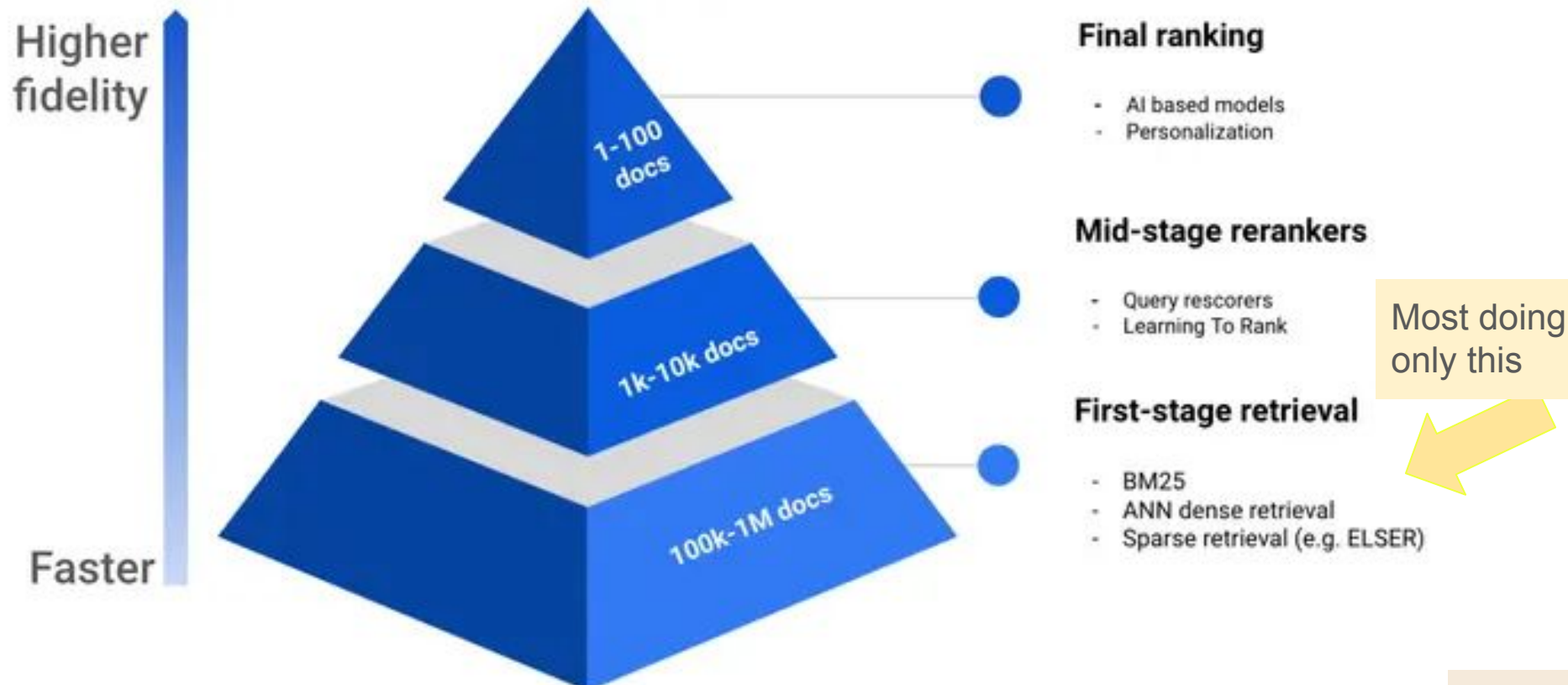# You've been interacting with reranking for a long time...

# In production, *many* considerations for **retrieval** – and RAG

- **Scale**: millions of items (videos, products, social feed...)

- **Speed**: Search over said scale of items within milliseconds

- **Fidelity:** ML models that learned more specific tasks can be more

  accurate

"Reranker" / multi-stage ranking is a versatile tool to address tradeoffs

# Scale < > Speed < > Fidelity tradeoffs



**Final ranking**

- AI based models
- Personalization

**Mid-stage rerankers**

- Query rescorers
- Learning To Rank

Most doing only this

**First-stage retrieval**

- BM25
- ANN dense retrieval
- Sparse retrieval (e.g. ELSER)

Higher fidelity

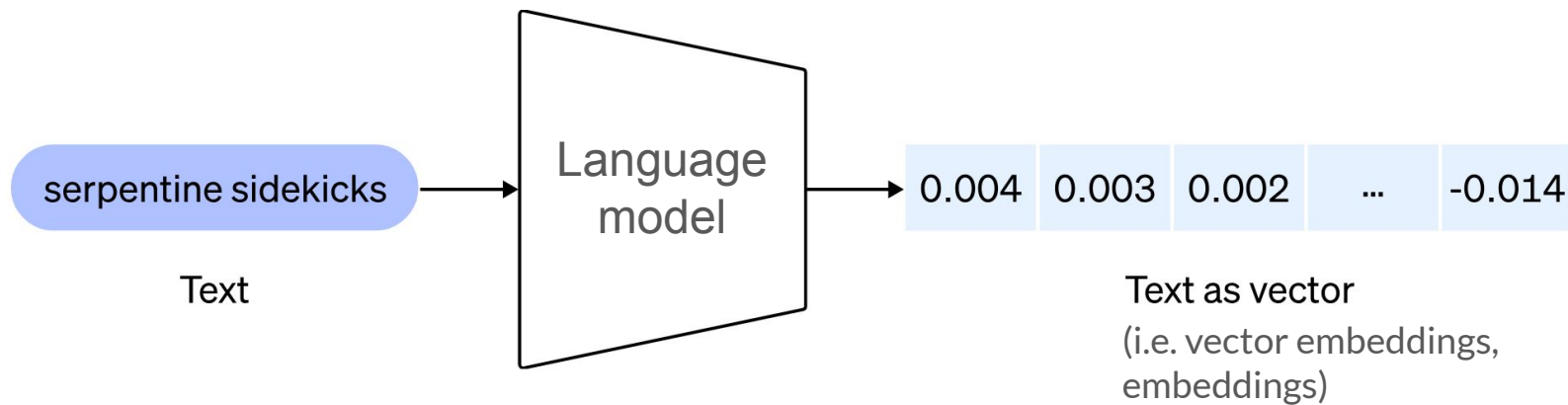Faster

1-100 docs

1k-10k docs

100k-1M docs

# Secret sauce?! How rerankers are trained

# Rerankers…

- Cohere-v3

- Elastic Rerank

- bge-reranker-v2-gemma (Google Gemma)

- mxbai-rerank-base-v1

- monot5-large

- MiniLM-L12-v2

etc.

# ML: Text is now turned into vectors / represented by *numbers*

serpentine sidekicks → Language model → 0.004 0.003 0.002 ... -0.014

Text

Text as vector
(i.e. vector embeddings, embeddings)

Machine learning models are happy. Love working with numbers.

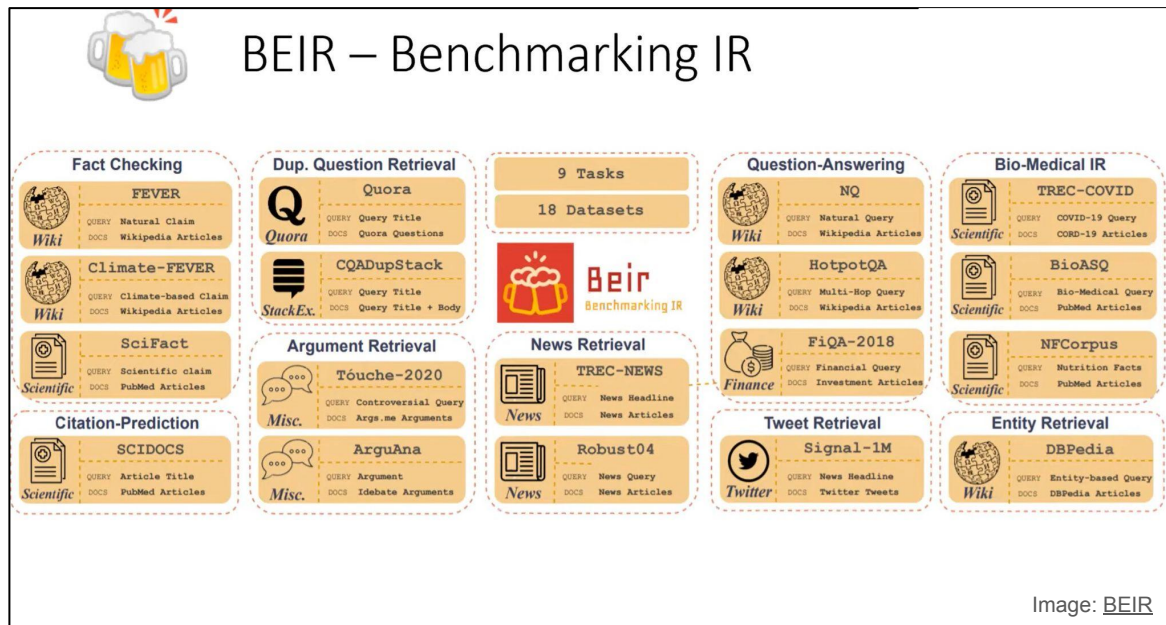Dataset

QA pairs

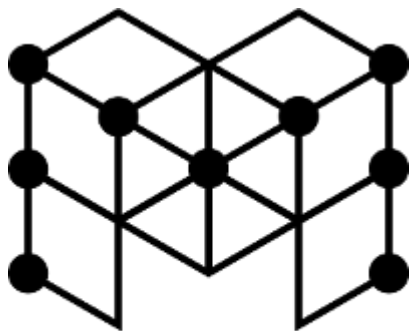Negative mining

Training

Cross-encoder

Model outputs

Building blocks for many rerankers, with examples from Elastic Rerank training

# Datasets

- ❏ BEIR – Benchmarking IR

- ❏ MTEB – Massive Text Embedding Benchmark

- ❏ MS MARCO



BEIR – Benchmarking IR

9 Tasks
18 Datasets

Image: BEIR

# MS MARCO: common question-answering dataset



Additional datasets for tasks such as passage ranking, keyphrase extraction, language generation etc. have been subsequently added

# MS MARCO dataset contains Bing questions + human generated answer



Q Will I qualify for OSAP if I'm new in Canada?

**Selected Passages from Bing**

"Visit the OSAP website for application deadlines. To get OSAP, you have to be eligible. You can apply using an online form, or you can print off the application forms. If you submit a paper application, you must pay an application fee. The online application is free."
Source: http://settlement.org/ontario/education/colleges-universities-and-institutes/financial-assistance-for-post-secondary-education/how-do-i-apply-for-the-ontario-student-assistance-program-osap/

"To be eligible to apply for financial assistance from the Ontario Student Assistance Program (OSAP), you must be a: 1 Canadian citizen; 2 Permanent resident; or 3 Protected person/convention refugee with a Protected Persons Status Document (PPSD)."
Source: http://settlement.org/ontario/education/colleges-universities-and-institutes/financial-assistance-for-post-secondary-education/who-is-eligible-for-the-ontario-student-assistance-program-osap/

"You will not be eligible for a Canada-Ontario Integrated Student Loan, but can apply for a part-time loan through the Canada Student Loans program. There are also grants, bursaries and scholarships available for both full-time and part-time students."
Source: http://www.campusaccess.com/financial-aid/osap.html

**Answer**
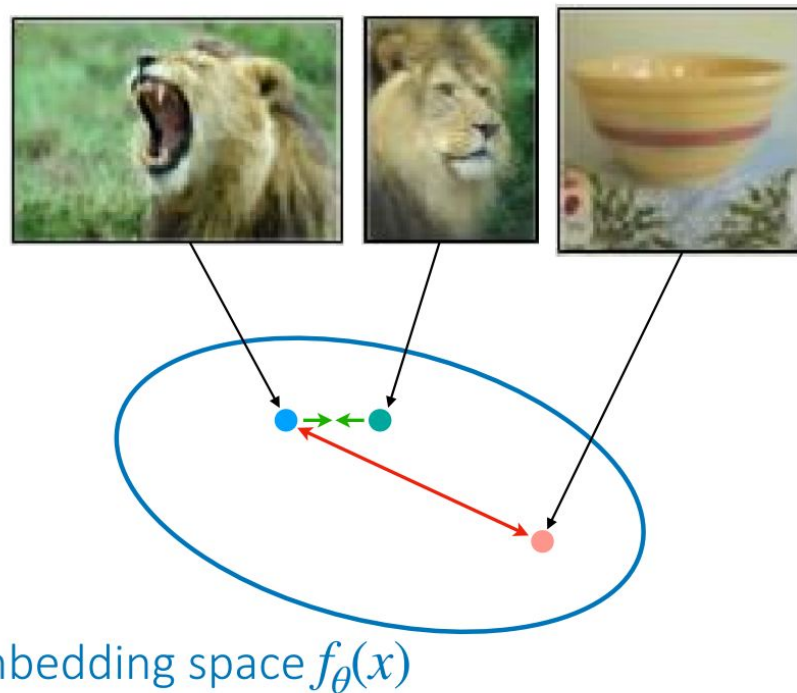No. You won't qualify.

Image: MS MARCO

# Dataset generation (Elastic Rerank)

- Total: 3 million queries

- Open QA datasets + 180,000 synthetic pairs

- Diverse query types

  - Keyword search

  - Exact phrase matching
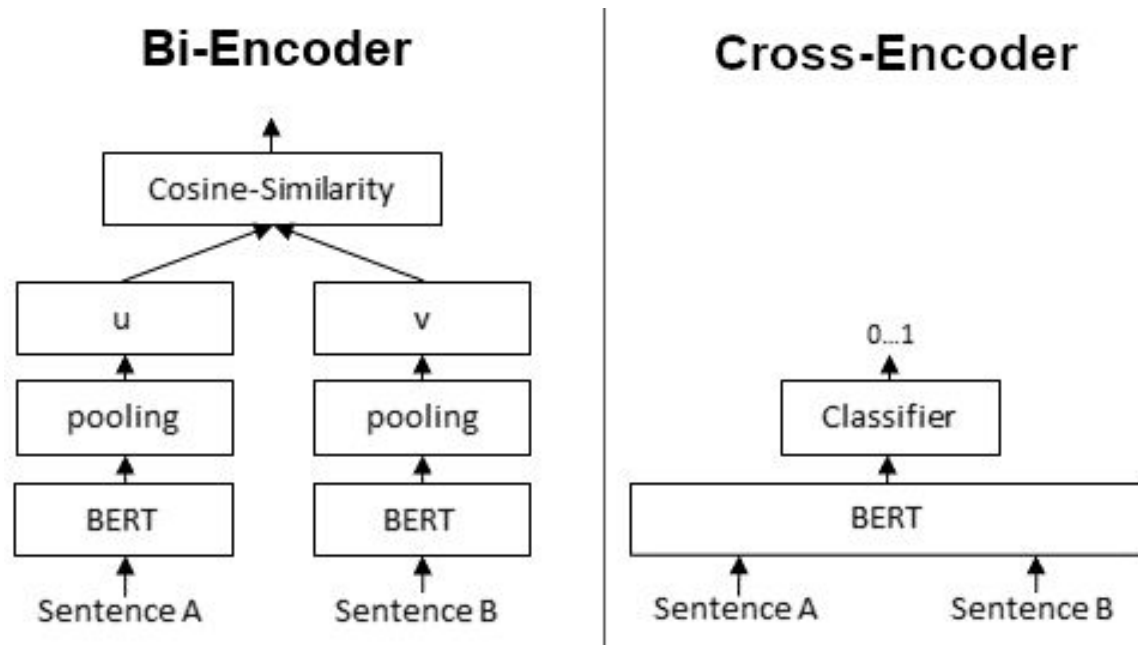
  - Short and long natural language questions

Present model with more diverse scenarios:
"Deep-negative mining with multiple negatives per query"
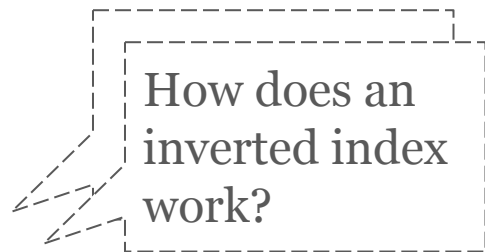
# Contrastive learning

❏ Bring together representations of similar examples

❏ Push apart representations of differing examples

Embedding space $f_\theta(x)$
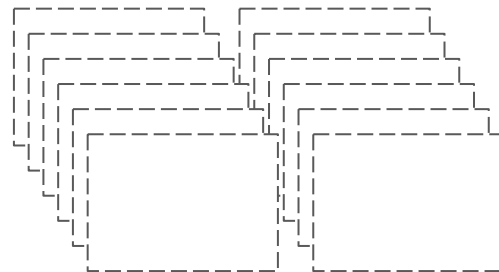
# ==Bi-encoder== and cross-encoder: different ways to structure encoders



**Bi-Encoder**

Cosine-Similarity

u | v

pooling | pooling

BERT | BERT

Sentence A | Sentence B

**Cross-Encoder**

0...1

Classifier

BERT

Sentence A | Sentence B

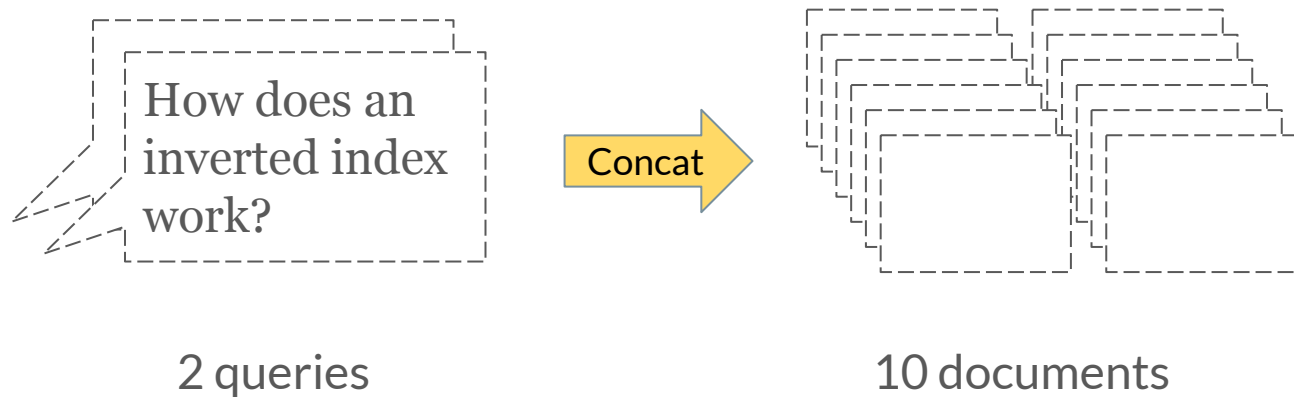# Bi-encoders allow for query and answer to be encoded separately

How does an inverted index work?

2 queries

10 documents

Bi-encoder: encode 12 times; can pre-encode 10 documents

# Cross-encoders are more accurate, but slower

How does an inverted index work?

Concat

2 queries

10 documents

Cross-encoder: encode 20 times on query time
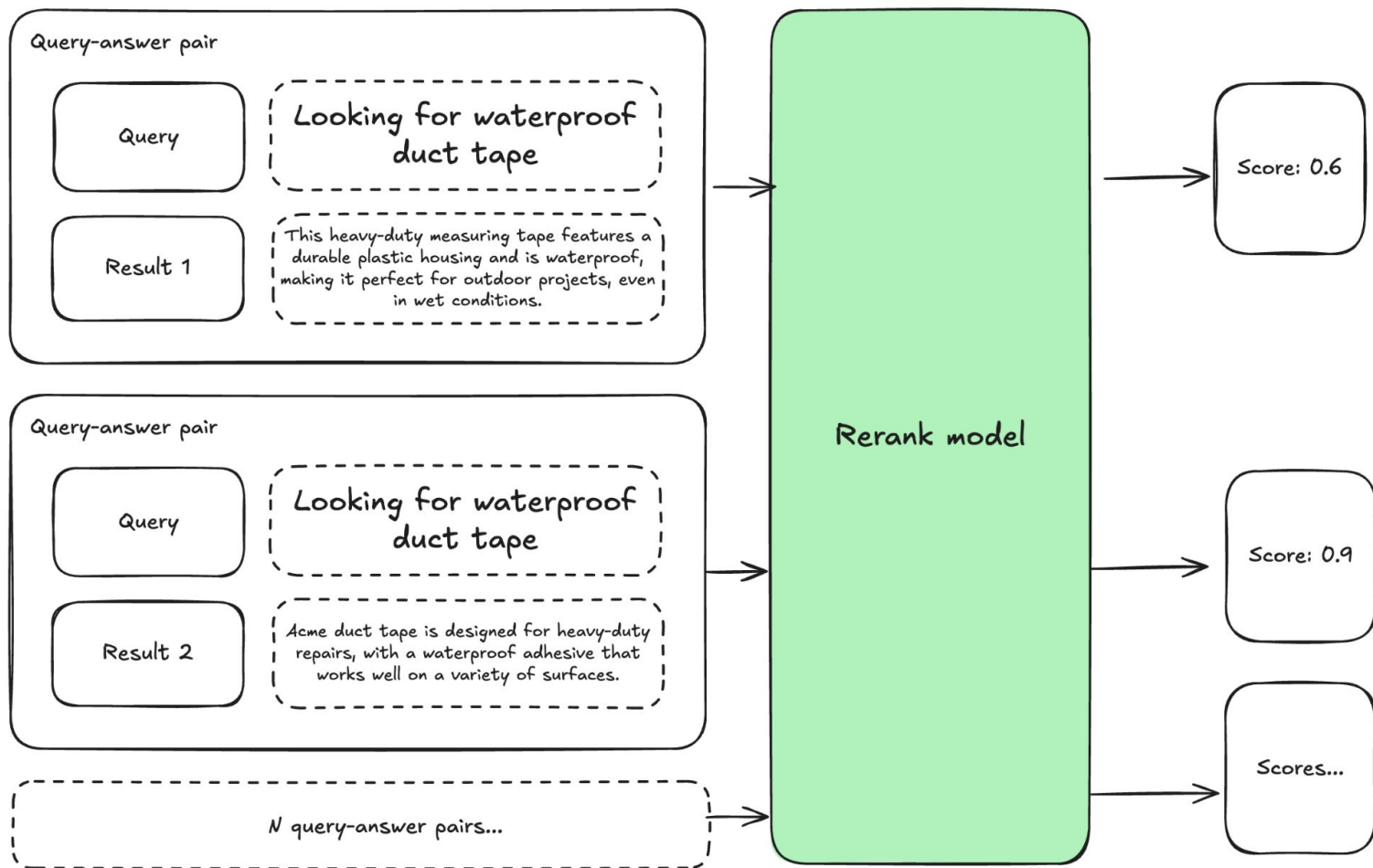
# Cross-encoders vs. Bi-encoders

## Cross-encoders

- Can learn more robust representations for generally assessing relevance
- Captures more nuanced semantics: Cross-encoder models can better learn how negation should affect relevance judgments
- Better calibrated across a diverse range of query types and topics. This makes choosing a score at which to drop documents significantly more reliable.
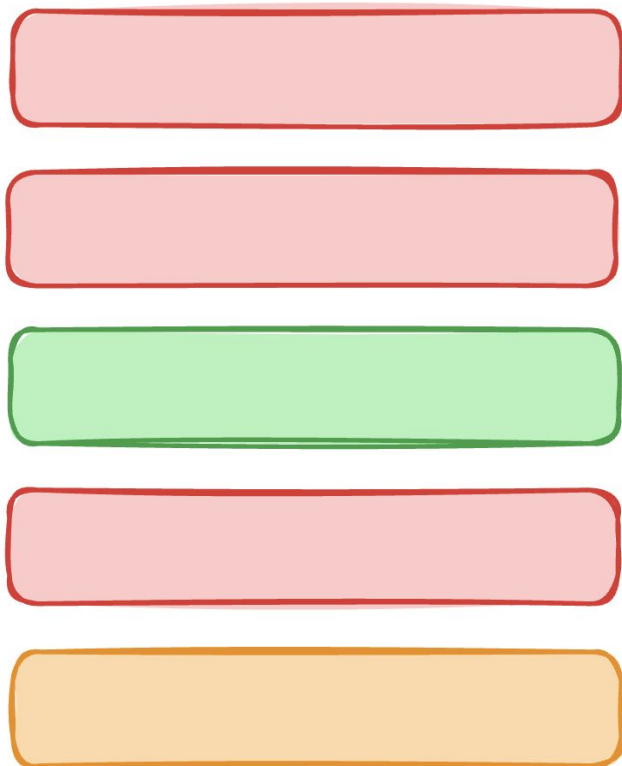
## Bi-encoders

- Bi-encoder models struggle more with things like negation and instead tend to pick up on matches for the majority concepts in the text, independent of whether the query wants to include or exclude them.
- Less encoding needed

**Query-answer pair**

Query: Looking for waterproof duct tape

Result 1: This heavy-duty measuring tape features a durable plastic housing and is waterproof, making it perfect for outdoor projects, even in wet conditions.

**Query-answer pair**

Query: Looking for waterproof duct tape

Result 2: Acme duct tape is designed for heavy-duty repairs, with a waterproof adhesive that works well on a variety of surfaces.

N query-answer pairs...

Rerank model
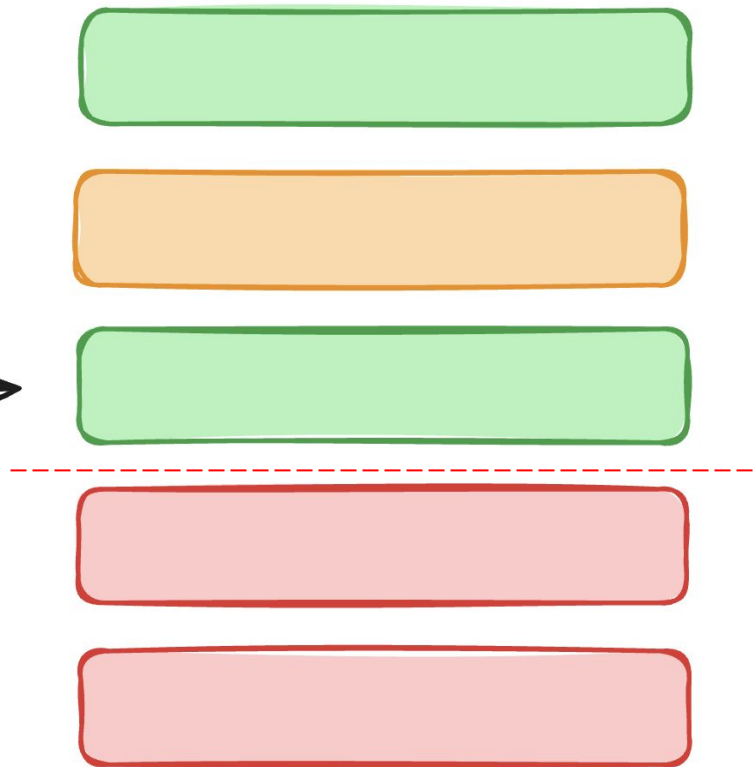
Score: 0.6

Score: 0.9

Scores...

Original results

Reranked results

# Example

```
{
  "query": "What duct tape is waterproof?",
  "input": [
```
"This heavy-duty **measuring tape** features a durable plastic housing and is **waterproof**, making it perfect for outdoor projects, even in wet conditions.",

"Our compact **measuring tape** is built with a weather-resistant coating, but it is **not** fully **waterproof**. It's ideal for general indoor and light outdoor use.",

"Acme **duct tape** is designed for heavy-duty repairs, with a **waterproof** adhesive that works well on a variety of surfaces.",

"This **tape** is ideal for office use, offering a strong adhesive and a clear finish. While it's *resistant to moisture*, it is not suitable for outdoor use or measuring purposes.",

"This **waterproof measuring tape** is engineered for construction and outdoor use, featuring a non-slip rubberized casing and a clear, easy-to-read scale that resists moisture and rust."
```
  ]
}
```

# Model outputs: scores for each document wrt query

```
{
 "rerank": [
  {
   "index": "3",
   "relevance_score": "0.99838966"
  },
  {
   "index": "1",
   "relevance_score": "0.587174"
  },
  {
   "index": "0",
   "relevance_score": "0.061199225"
  },
  {
   "index": "2",
   "relevance_score": "0.032283258"
  },
  {
   "index": "4",
   "relevance_score": "0.015365343"
```
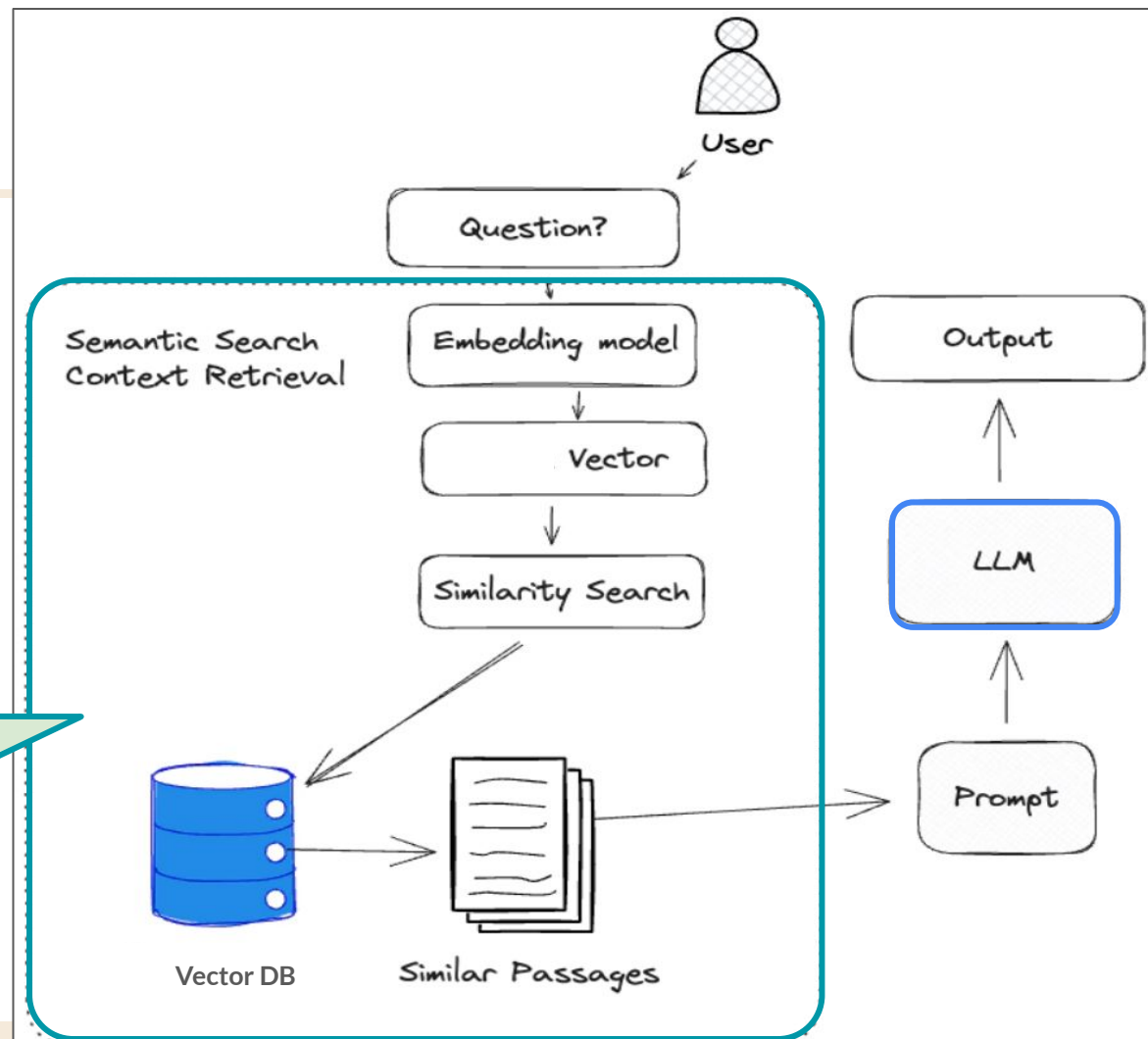
**IN PRODUCTION**

Use relevance scores **cutoff** to **prune** irrelevant documents.

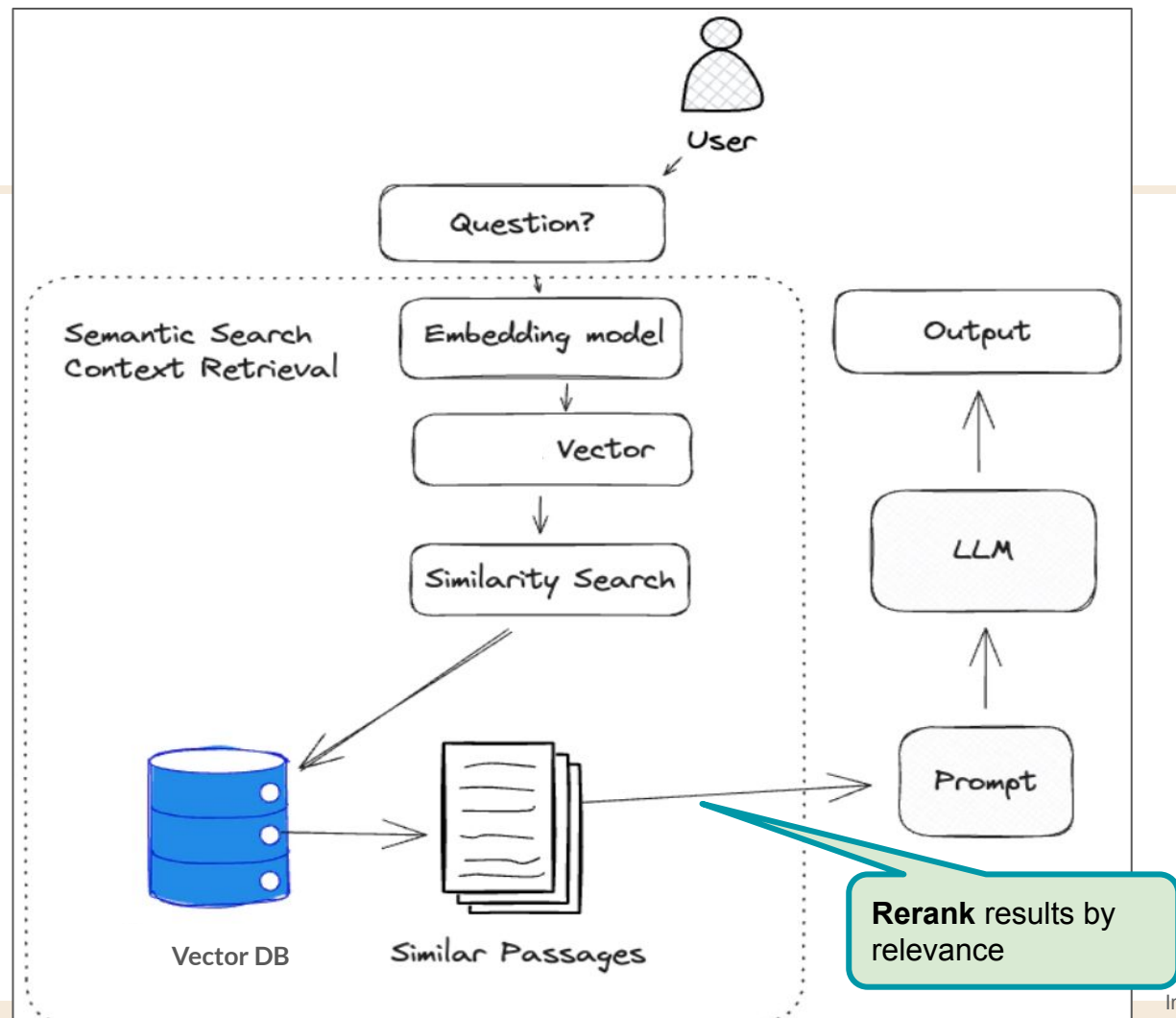Use scores for **observability**: window into how well your RAG is performing.
- "Why is my average relevancy score for these types of queries so low?". By contrast, without these metrics you're in the dark when it comes to optimizing your RAG pipeline.
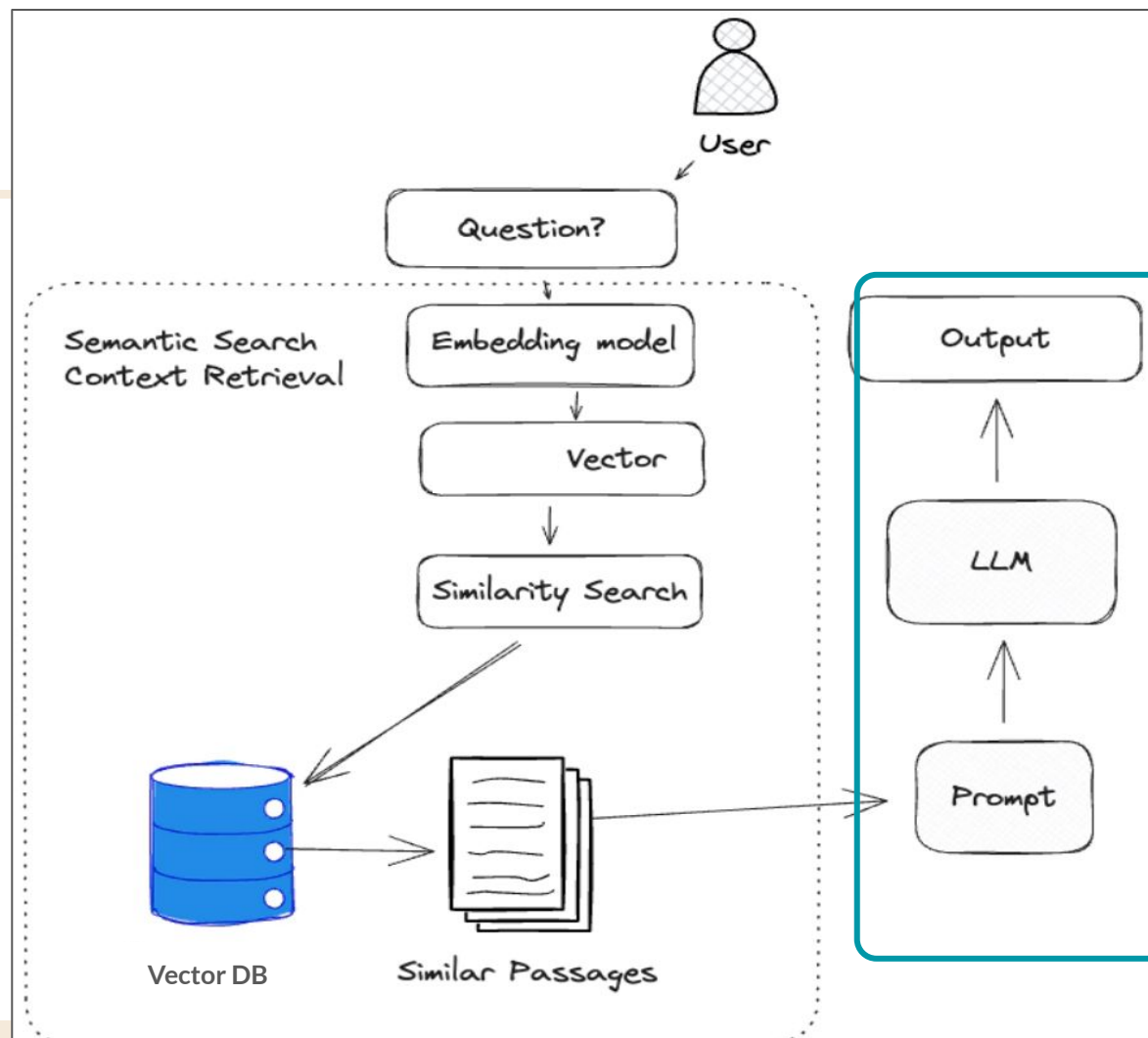
5

All together: Rerankers in RAG systems

**Retrieve** the **most relevant** [support articles / product pages / discount policies] etc…

User

Question?

Semantic Search Context Retrieval

Embedding model

Vector

Similarity Search

Output

LLM

Prompt

**Vector DB**

Similar Passages

Image: Elastic

User

Question?

Semantic Search
Context Retrieval

Embedding model

Vector

Similarity Search

Output

LLM

Prompt

Vector DB

Similar Passages

**Rerank** results by relevance

Image: Elastic

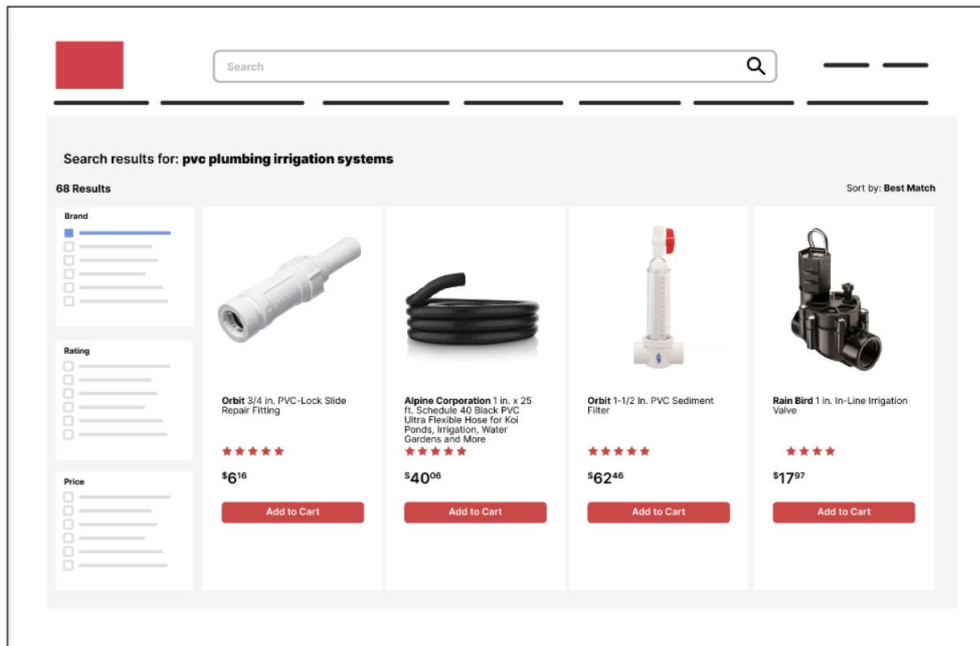Pass reranked info to LLM/Chatbot to **generate** the final output

Image: Elastic

# Reranking: Helps your plain ol' search results (if you're not using RAG)



Reranked search results

**Rerank** results by relevance

BM25 is still a strong baseline; tuning to compete for user experience

# Conclusion & Going forward

1. Use reranking to improve results (the "R" in RAG)

# Conclusion & Going forward

1. Use reranking to improve results (the "R" in RAG)

2. Speed <> Fidelity <> Scale tradeoffs

   ○ Faster model to filter, Slower but more fine-grained model to do final re-ranking

# Conclusion & Going forward

1. Use reranking to improve results (the "R" in RAG)

2. Speed <> Fidelity <> Scale tradeoffs

    ○ Faster model to filter, Slower but more fine-grained model to do final re-ranking

3. Tune and improve individual GenAI system components

    ○ Retrieval: **Keyword** or **semantic** or **hybrid** e.g. Reciprocal Rank Fusion (RRF)

    ○ Semantic retrieval: **Sparse** (e.g. ELSER) or **dense** (e.g. E5)

    ○ Reranker, fine tuning, etc.

# Conclusion & Going forward

1. Use reranking to improve results (the "R" in RAG)

2. Speed <> Fidelity <> Scale tradeoffs

   ○ Faster model to filter, Slower but more fine-grained model to do final re-ranking

3. Tune and improve individual GenAI system components

   ○ Retrieval: **Keyword** or **semantic** or **hybrid** e.g. Reciprocal Rank Fusion (RRF)

   ○ Semantic retrieval: **Sparse** (e.g. ELSER) or **dense** (e.g. E5)

   ○ Reranker, fine tuning, etc.

# Q&A

Find me at:

- **Elastic booth** or **O'Reilly booth**
- LinkedIn: [Susan Shu Chang](#)
- Read more about our work: Elastic blog

Get the **Machine Learning Interviews** book or recommend to someone you know:

[https://amzn.to/4aOjO26](https://amzn.to/4aOjO26)

(Limited copies available at DDT O'Reilly booth)



O'REILLY®

**Machine Learning Interviews**

Kickstart Your Machine Learning and Data Career

Susan Shu Chang