



SEMANTIC SEARCH ENGINES OVER BIG DATA USING LLMS

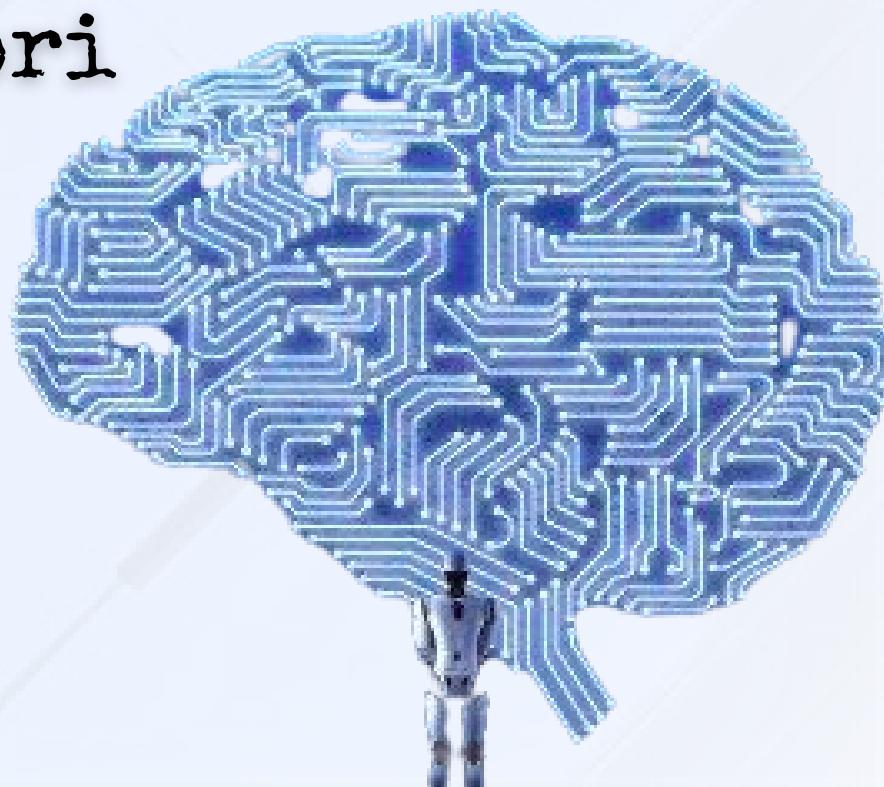
SEARCHING LARGE TEXT CORPORA WITH CONTEXTUAL UNDERSTANDING

Khouloud Ben Younes

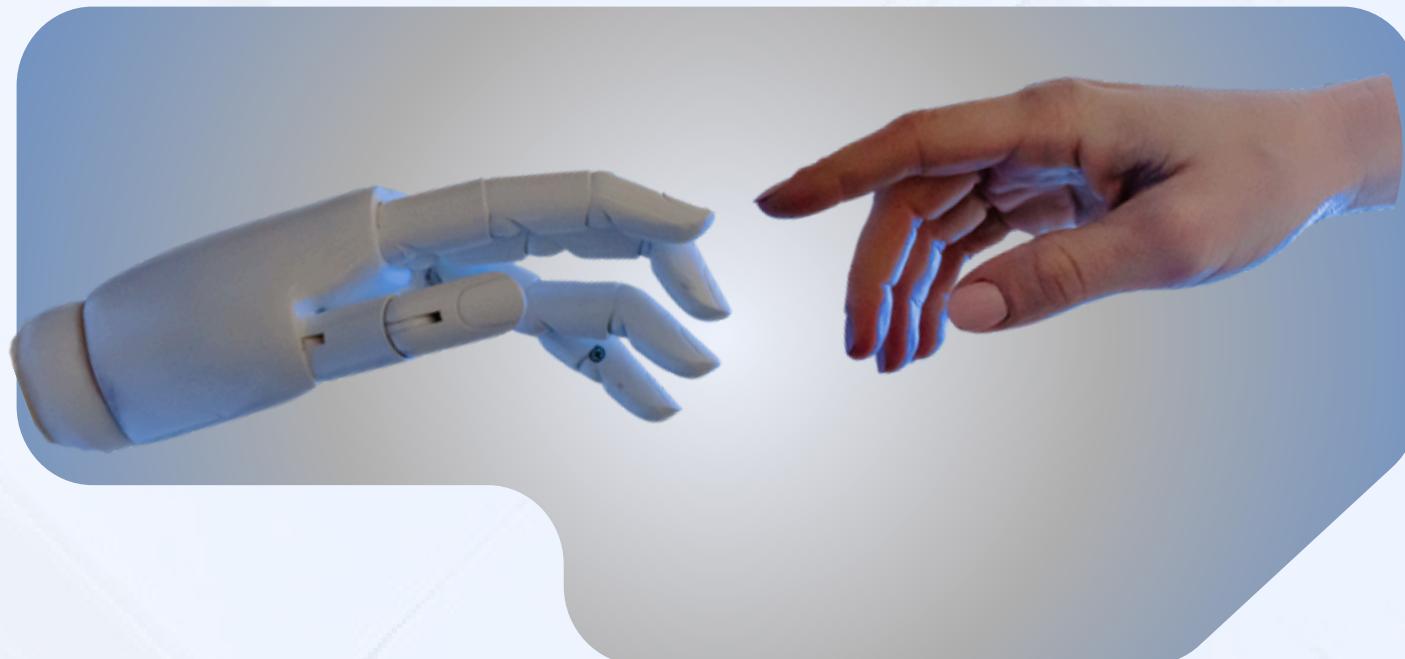
Montaha Ghabri

Evaluated by :
Dr. Manel AbdelKader

2025-2026



AGENDA



01

INTRODUCTION

02

LITERATURE SELECTION

03

BACKGROUND & KEY CONCEPTS

04

REVIEW & ANALYSIS OF EXISTING WORK

05

MAJOR INSIGHTS AND DISCOVERIES

06

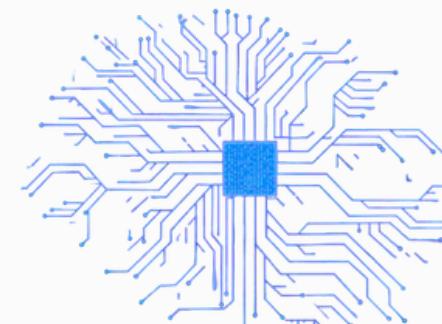
CHALLENGES AND FUTURE DIRECTION

07

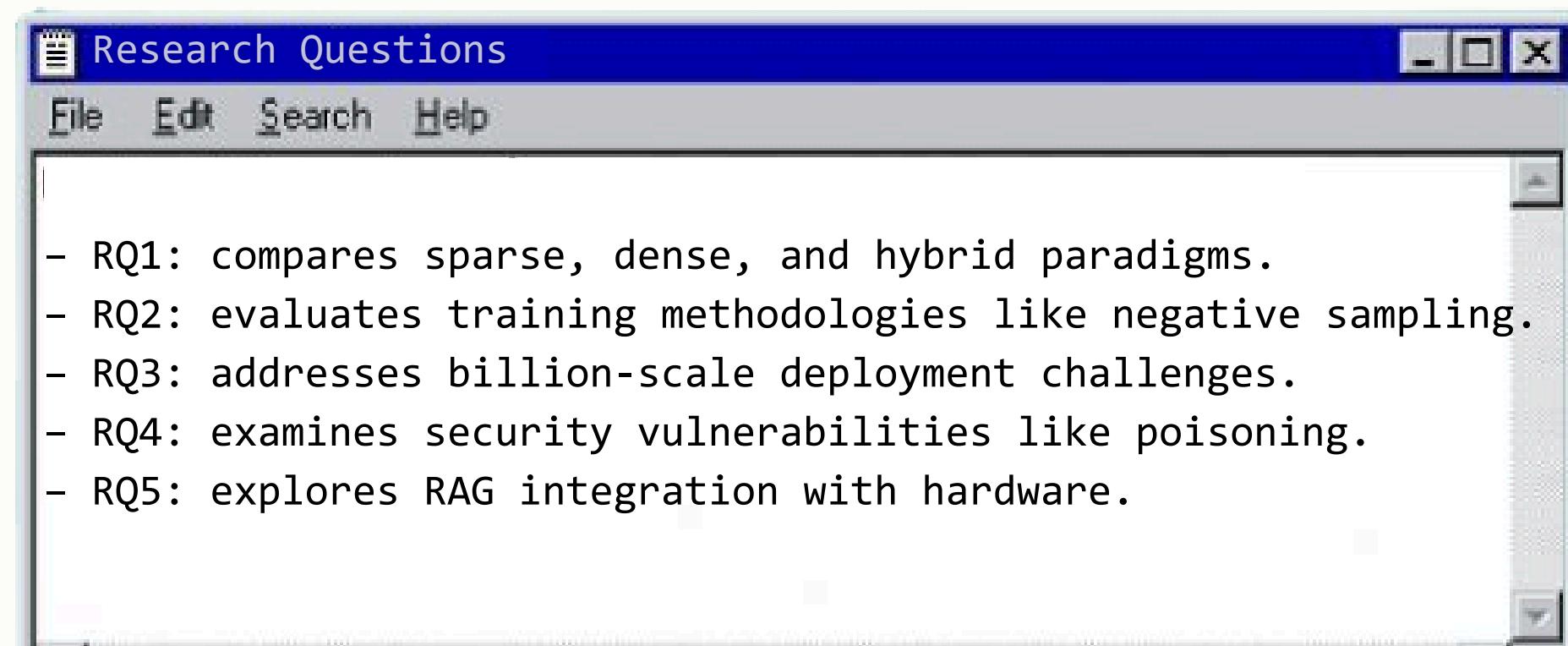
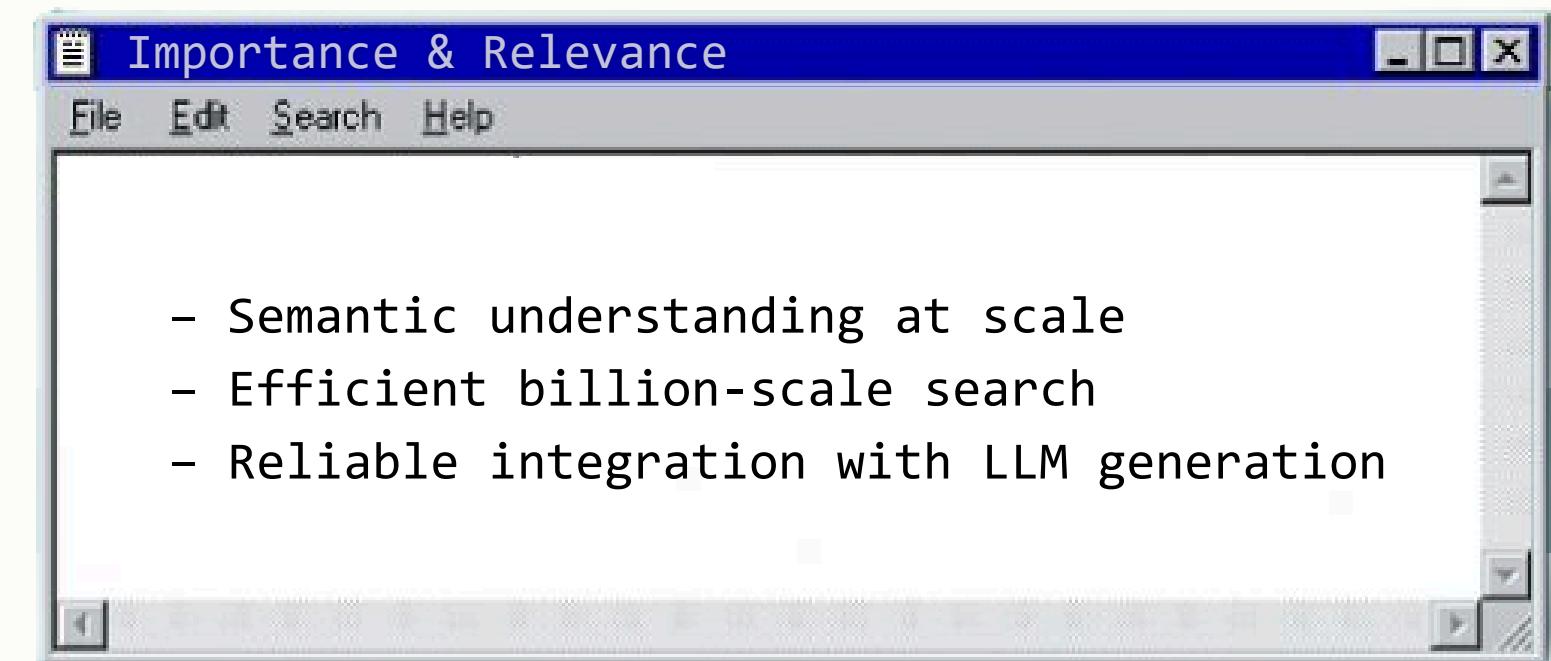
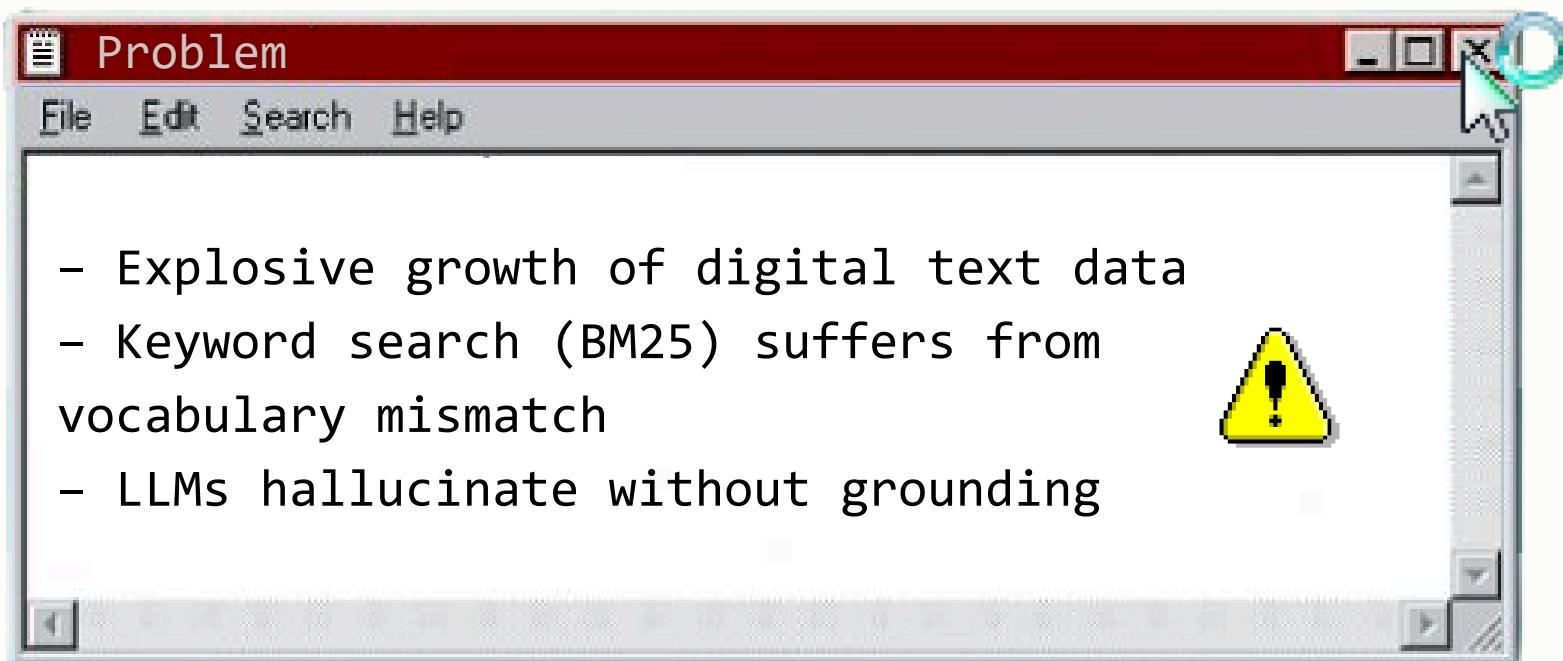
CONCLUSION

MOTIVATION

Start



1 MOTIVATION



LITERATURE SELECTION



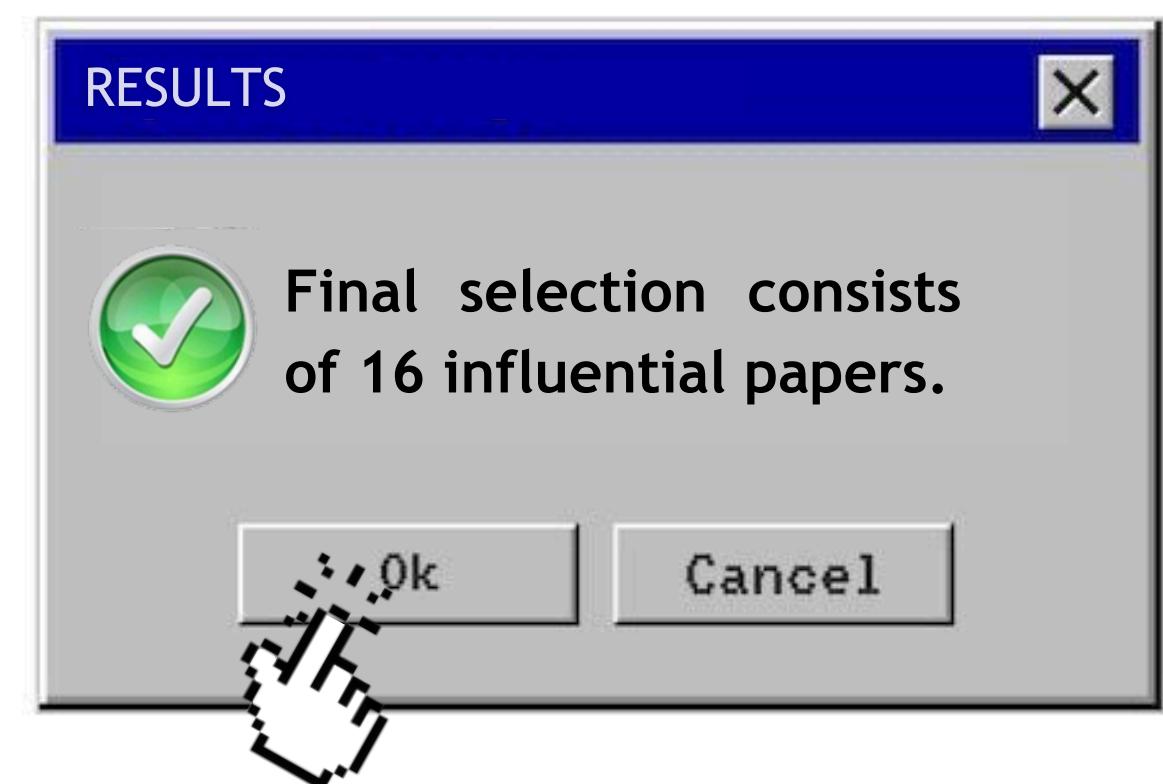
2 LITERATURE SELECTION

Inclusion criteria:

- Top conferences: SIGIR, ACL, EMNLP, NeurIPS, VLDB
- Citation threshold:
 - 100+ citations for 2020-2022 papers
 - 20+ citations for 2023 papers
 - 10+ citations for 2024-2025 papers

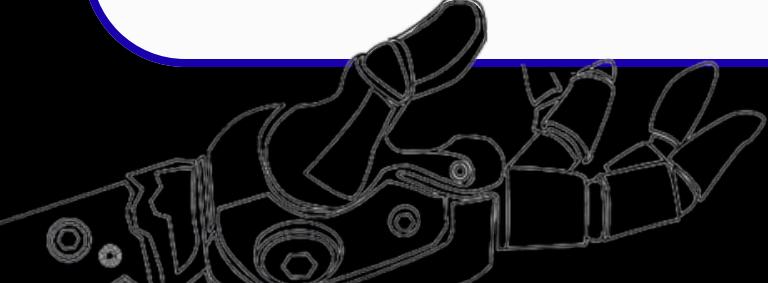
Exclusion criteria:

- Workshop, demo, poster, or short papers
- Non-English publications
- Redundant with earlier seminal works
- Focused only on non-text modalities or unrelated tasks



Title	Authors	Year	Citations	Category	Link	Database	Status	Priority	Notes
Dense Passage Retrieval for Open-Domain Question Answering	Vladimir Karpukhin*, Barlas Oğuz*	2020	5695	Dense Retrieval Foundations	https://arxiv.org/pdf/2004.0468				
SIMLM: Pre-training with Representation Bottleneck for Dense Passage Retrieval	Liang Wang, Nan Yang, Xiaolei	2023	148	Dense Retrieval Foundations	https://aclanthology.org/2023.10.148				

BACKGROUND & KEY CONCEPTS



3 Background & Key Concepts

The screenshot shows a software application window with a sidebar on the left and a main content area on the right.

Definitions:

- Embeddings:  **Dense vector representations (768-1536 dimensions)**
- Dense Retrieval 
- Contrastive Learning 
- Hard negatives 
- Vector Quantization (VQ) 
- Suspend 

Start 

3 Background & Key Concepts

The screenshot shows a software application window with a sidebar on the left and a main content area on the right.

Definitions:

- Embeddings:
- Dense Retrieval:
- Contrastive Learning:
- Hard negatives:
- Vector Quantization (VQ):
- Suspend:

A tooltip is displayed over the "Dense Retrieval" item, containing the text: "Dual encoders with ANN search (FAISS, HNSW)".

Start

3 Background & Key Concepts

The screenshot shows a software application window with a sidebar on the left and a main content area on the right.

Definitions:

- Embeddings:
- Dense Retrieval:
- Contrastive Learning: **Training objective pulling positive query-passage pairs close while pushing negative pairs apart** (A tooltip box is displayed over this item.)
- Hard negatives:
- Vector Quantization (VQ):
- Suspend:

Start

3 Background & Key Concepts

The screenshot shows a software application window. On the left, there is a vertical sidebar with a dark grey header labeled "Definitions". Below this, there is a list of items, each with an icon and a label:

- Embeddings: (Icon of a folder with files)
- Dense Retrieval (Icon of a folder with files)
- Contrastive Learning (Icon of a document with gears)
- Hard negatives (Icon of a book with a question mark)
- Vector Quantization (VQ) (Icon of a computer monitor displaying a chart)
- Suspend (Icon of a computer monitor)

The "Hard negatives" item is currently selected, indicated by a mouse cursor icon pointing at its icon. A tooltip or callout box is displayed to the right of the "Hard negatives" entry, containing the text: "semantically similar but incorrect passages that improve discrimination".

At the bottom of the sidebar, there is a "Start" button with a colorful icon.

3 Background & Key Concepts

The screenshot shows a software application window with a sidebar on the left and a main content area on the right.

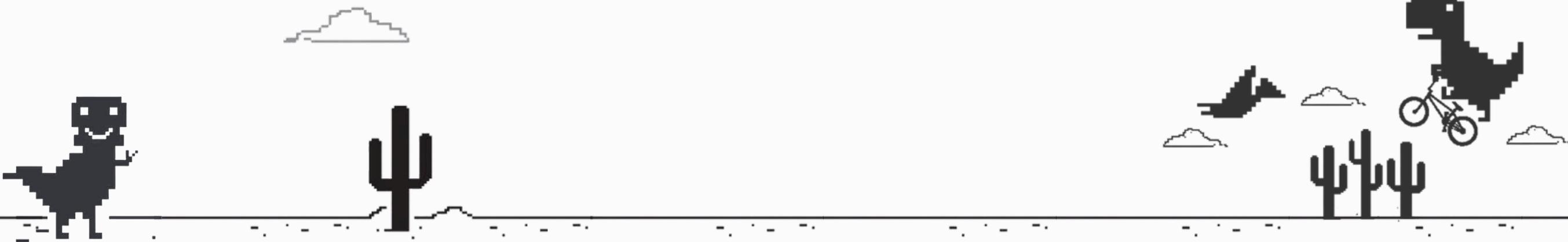
Definitions:

- Embeddings:
- Dense Retrieval:
- Contrastive Learning:
- Hard negatives:
- Vector Quantization (VQ): A tooltip is displayed over this item:

Compression methods replacing full-precision vectors with codes
- Suspend:

Start

REVIEW OF EXISTING WORK



You are offmark

Try:

Reading recent papers.

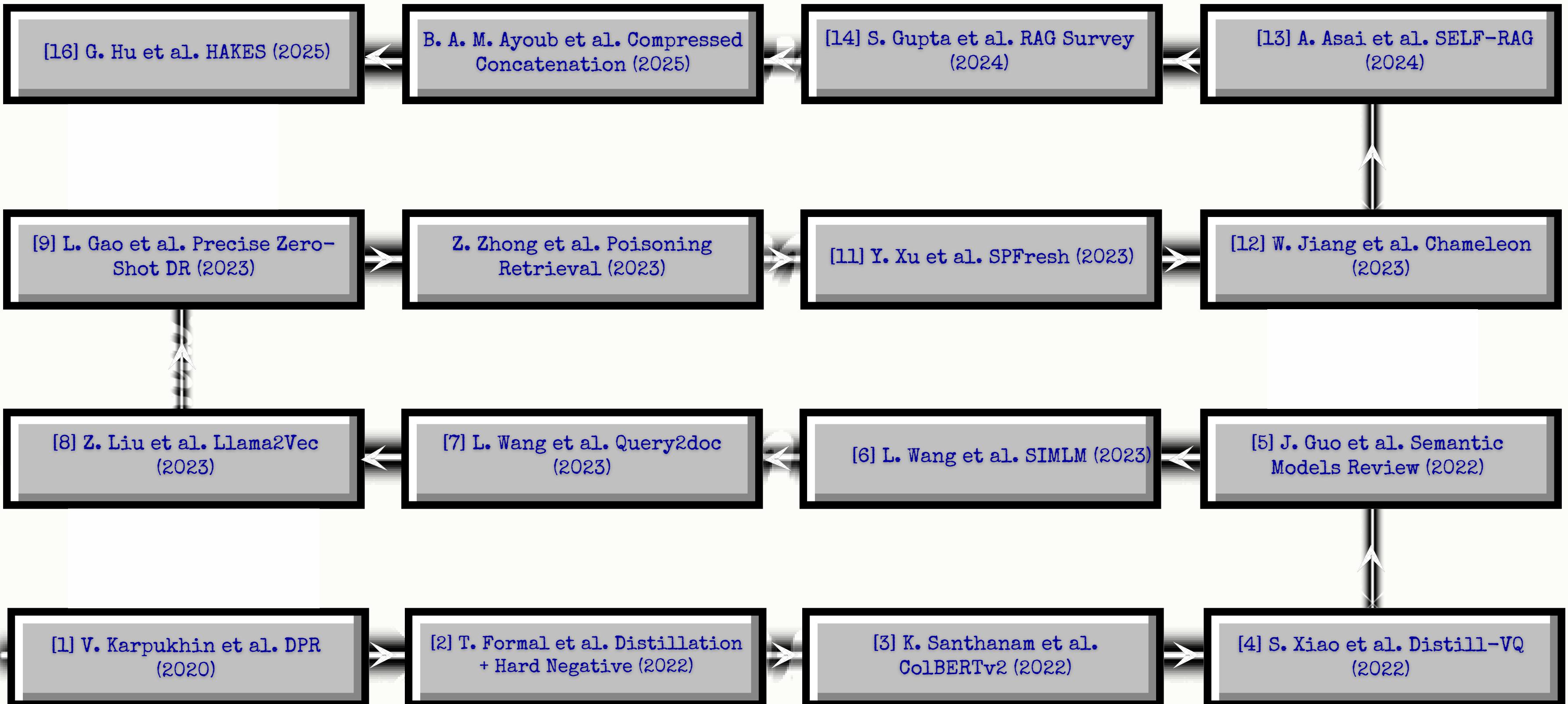




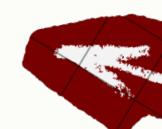
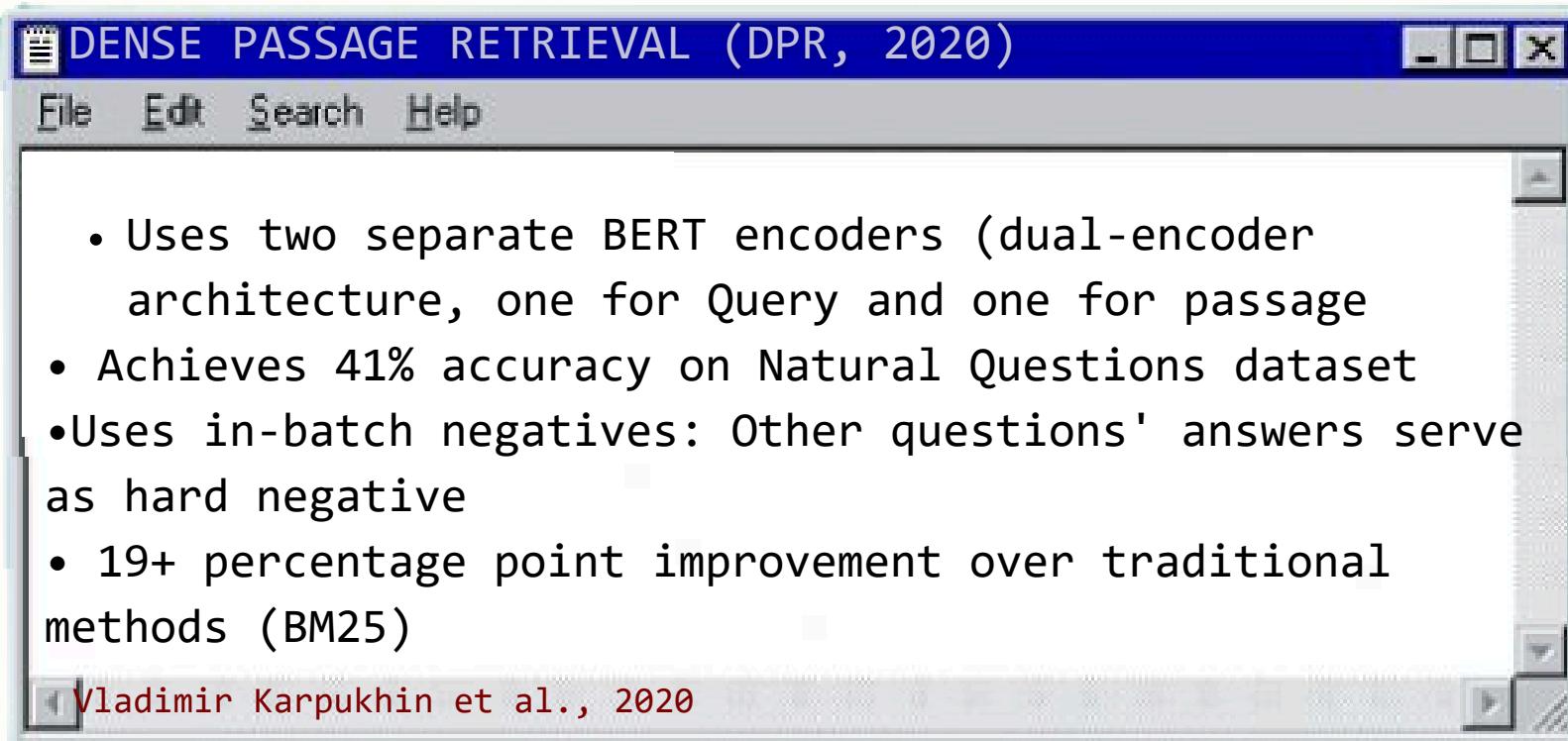
FOUR RESEARCH CATEGORIES



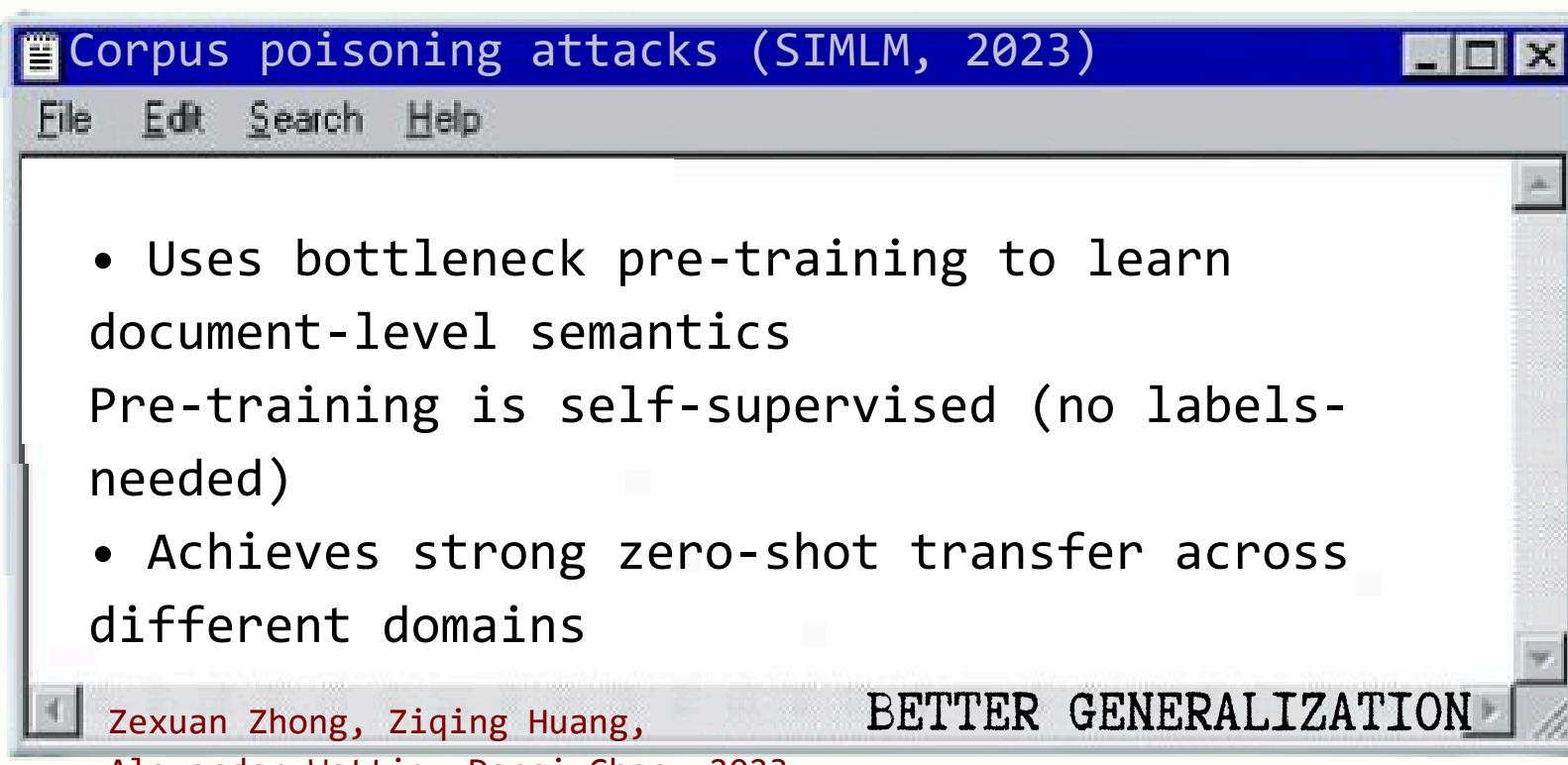
TIMELINE



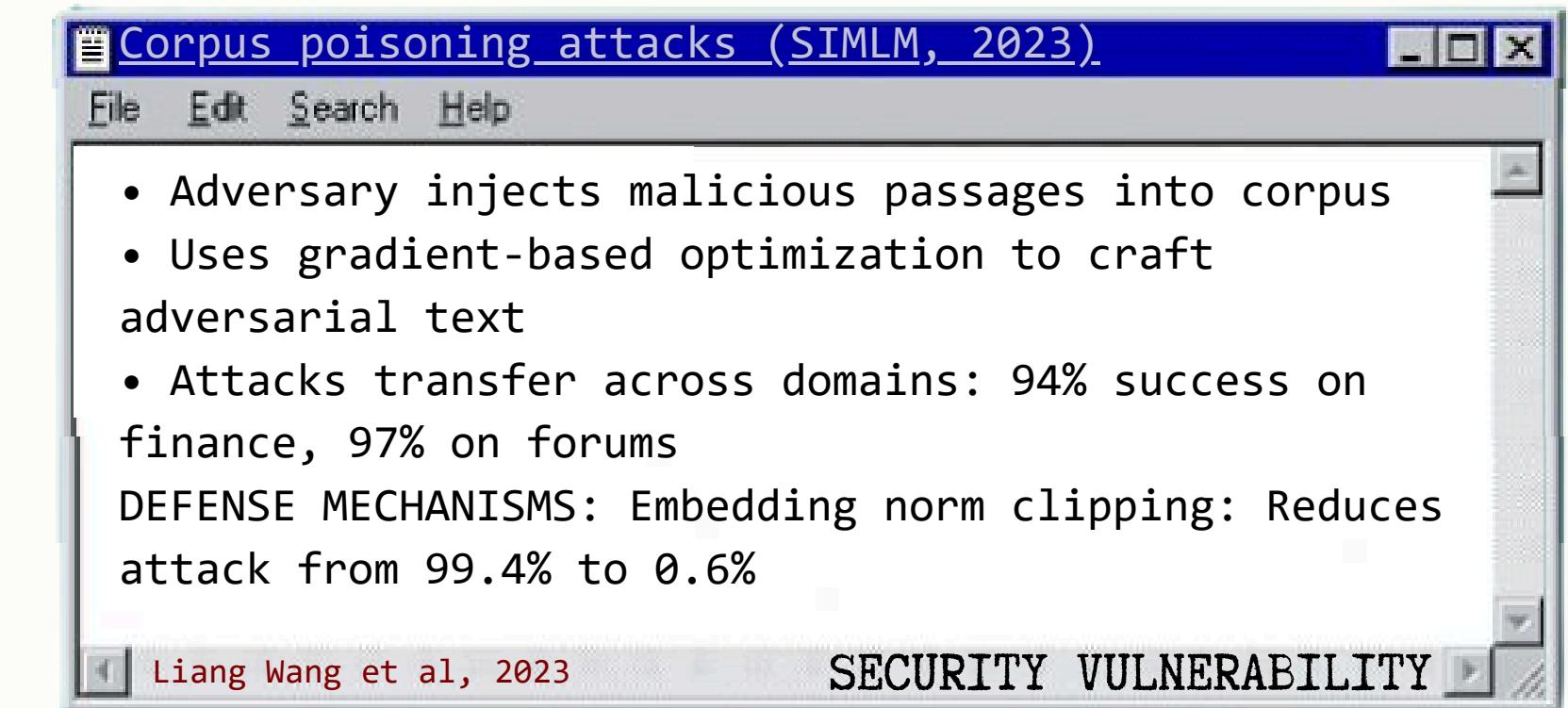
1- DENSE RETRIEVAL FOUNDATIONS



HOW NEURAL RETRIEVAL WORKS



BETTER GENERALIZATION



SECURITY VULNERABILITY

2- SCALABILITY SOLUTIONS FOR PRODUCTION

HAKES: Scalable Vector Database for Embedding_Search Service

- ✖ Traditional databases couple storage and computation!
- Can't scale read and write operations independently!

TWO-STAGE ARCHITECTURE:

STAGE 1 - Fast Filtering: Quickly narrows down to candidate documents

STAGE 2 - Precise Refinement: Exact similarity computation on candidates

Handles concurrent read-write workloads efficiently

Guoyu Hu et al. , 2025

SPFresh- INCREMENTAL INDEX UPDATES

- ✖ Traditional indices require full rebuilds when data changes !

- Lightweight Incremental Rebalancing (LIRE) approach
- When new vectors arrive, only rebalance overloaded partitions
- 2-5x faster updates on billion-vector indices

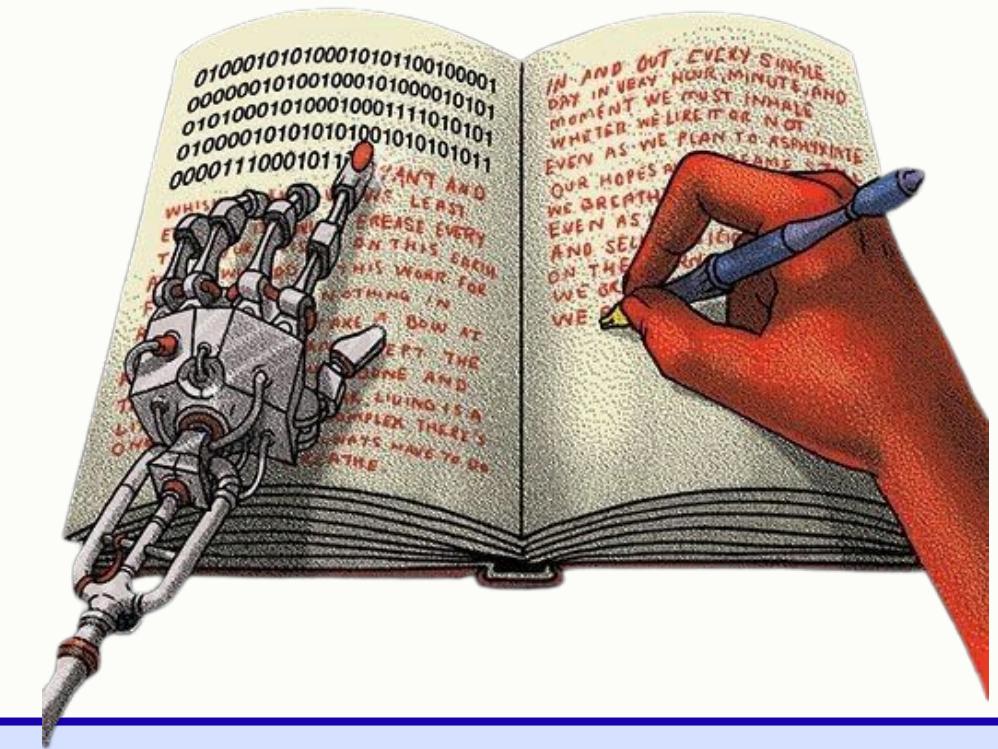
Yuming Xu et al., 2023

Distill-VQ: Learning Retrieval Oriented VQ by distilling knowledge

- ✖ 1. Product Quantization (PQ) compresses vectors to save memory.
- 2. Reconstructed vectors ≠ good ranking preservation

- Uses knowledge distillation from a frozen teacher retriever
- Teacher provides soft relevance judgments for query-document pairs
- Student (compressed model) learns to preserve teacher's ranking
- No labeled data needed - teacher provides supervision, 2-5% MRR and recall improvement vs standard PQ

Shitao Xiao et al. , 2022



3- RAG ARCHITECTURES

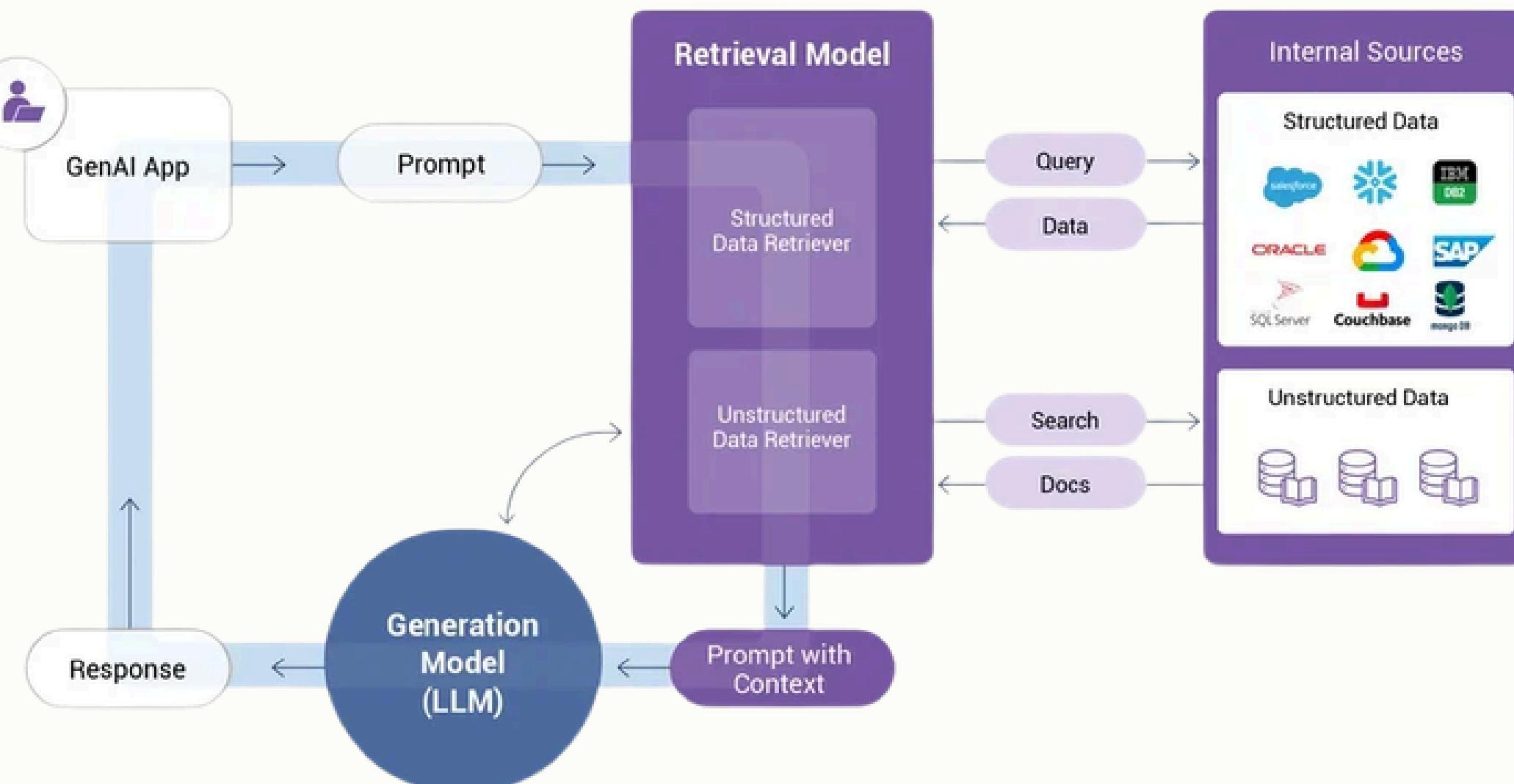
WHY RAG?

- LLMs can generate plausible but incorrect information
- Knowledge becomes outdated (training data cutoff)
- RAG grounds generation in retrieved factual documents
- Enables knowledge updates without retraining

Shailja Gupta, Rajesh Ranjan, Surya Narayan Singh, 2024

OK

COMBINE RETRIEVAL WITH LLM GENERATION TO REDUCE HALLUCINATIONS



3- RAG ARCHITECTURES

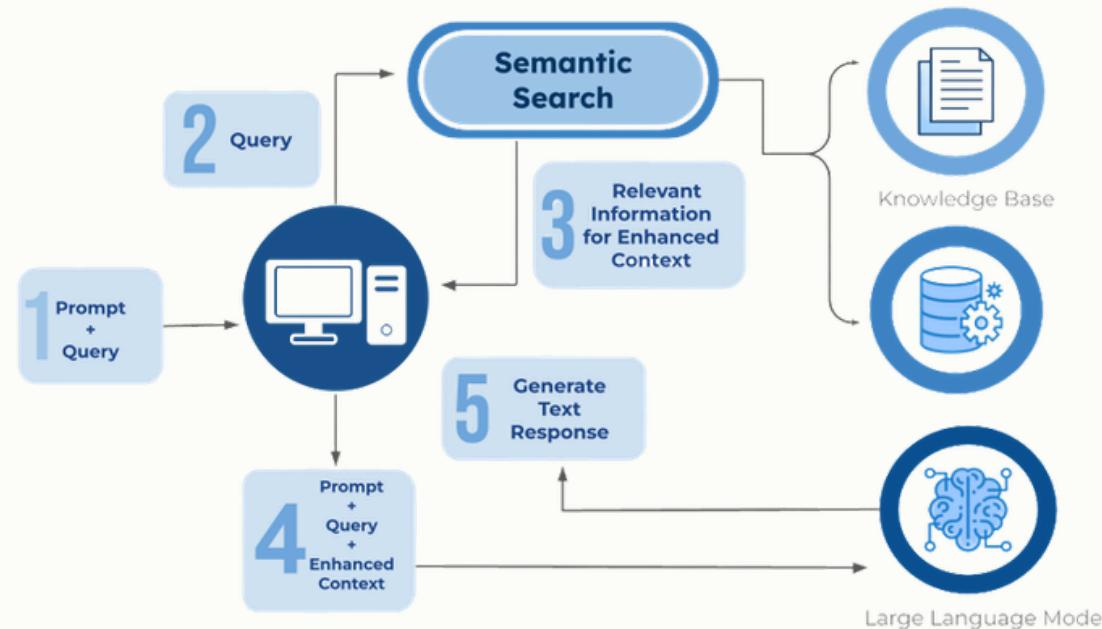
Combine retrieval with LLM generation to reduce hallucinations

SELF-RAG - INTELLIGENT RETRIEVAL

- Model generates special control tokens alongside content:
 - [Retrieve]: Decides when to trigger retrieval system
 - [Relevant]/[Irrelevant]: Assesses if retrieved docs relate to query
 - [Supported]/[Not Supported]: Checks if output is grounded in evidence
 - [Utility]: Evaluates overall output quality
- Critic model automatically annotates training data
- Main model learns to predict these tokens via next-token prediction

Reduces unnecessary retrievals by 50%

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, Hannaneh Hajishirzi, 2024



QUERY2DOC

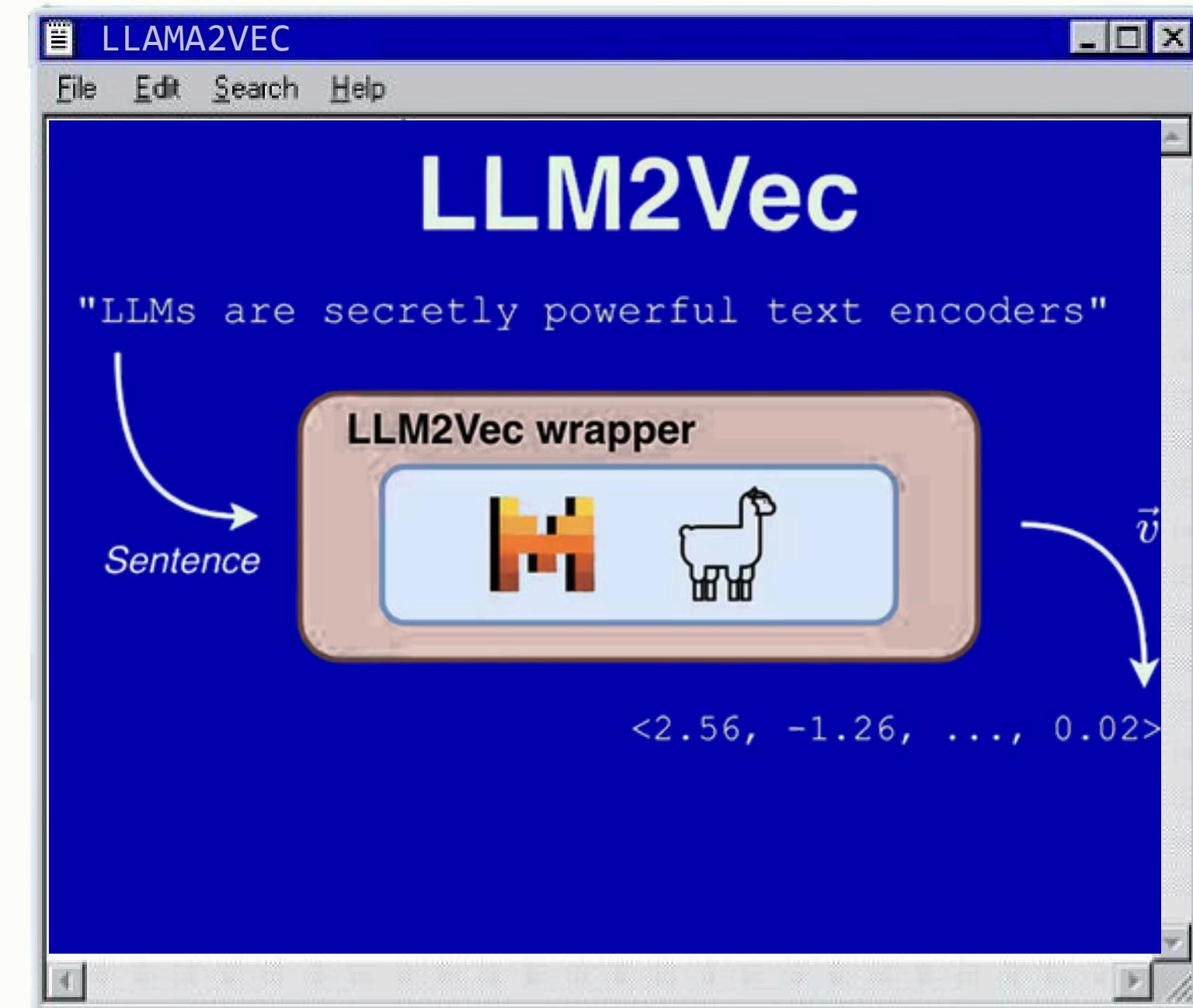
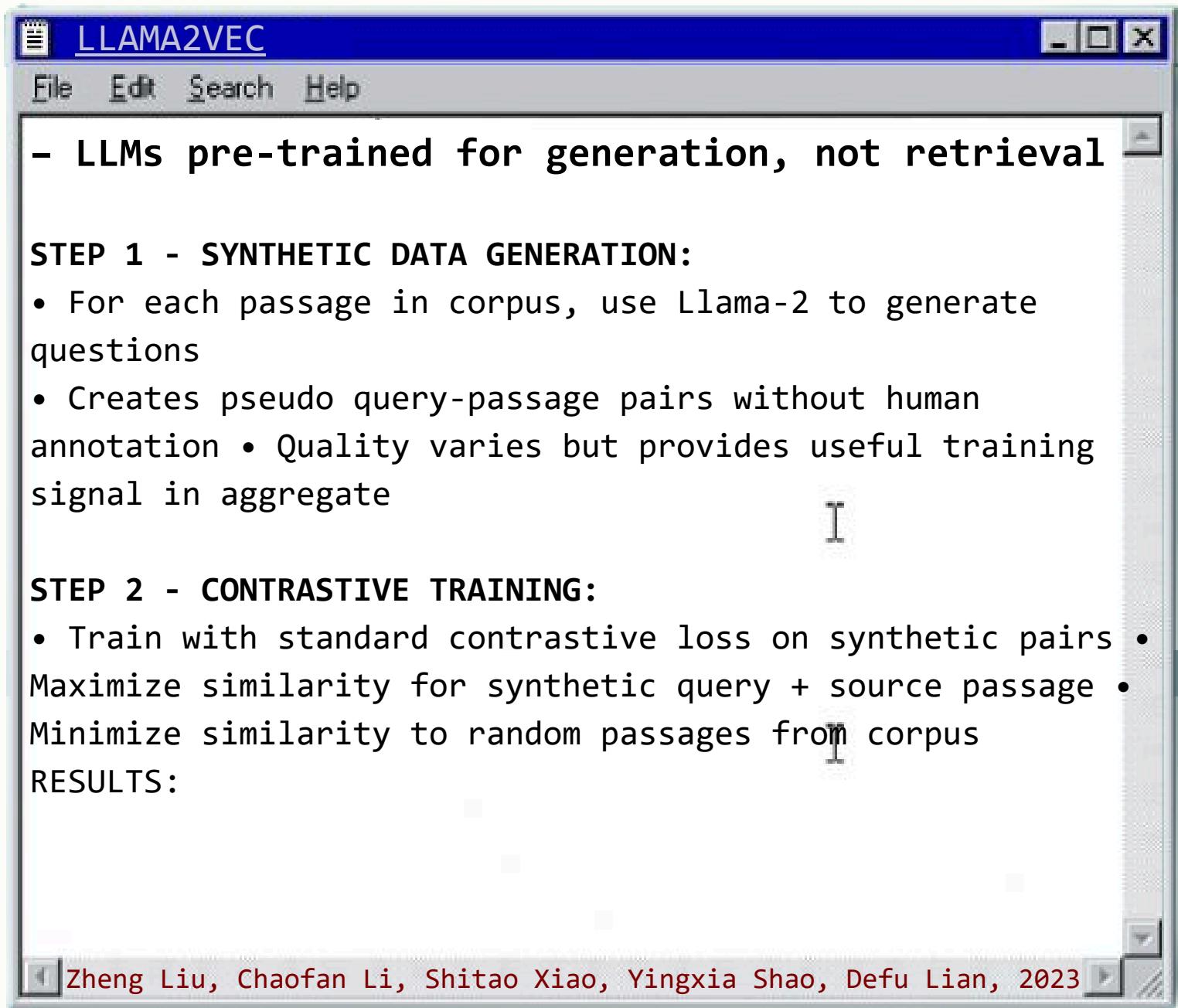
Queries are short and may lack vocabulary of relevant documents !

- Uses LLM to generate hypothetical relevant document for query
- Concatenates generated text with original query
- No training or labeled data needed for target task
- Compatible with both BM25 (sparse) and dense retrievers
- Up to 15% improvement on BM25 (zero-shot)

Liang Wang, Nan Yang, Furu Wei, 2023

3- RAG ARCHITECTURES

Combine retrieval with LLM generation to reduce hallucinations



- Matches supervised dense retrievers on BEIR (out-of-domain)
- Enables deployment in domains without labeled retrieval data

3- RAG ARCHITECTURES

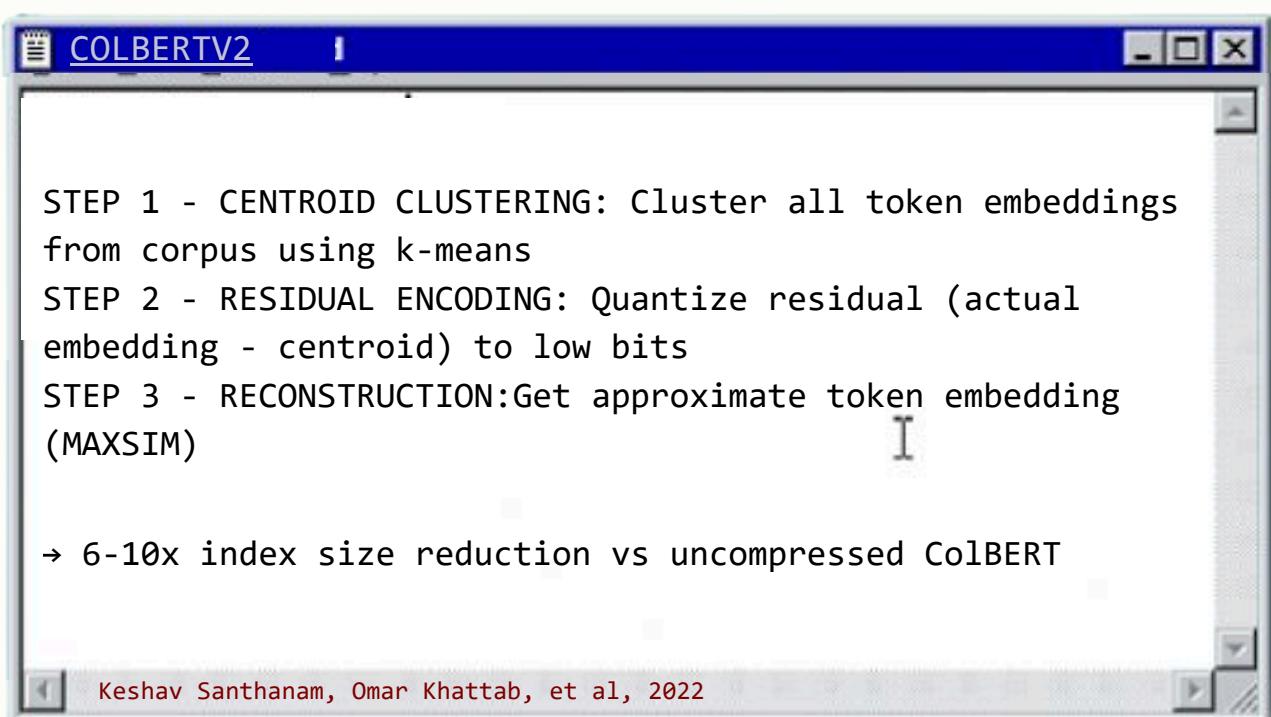
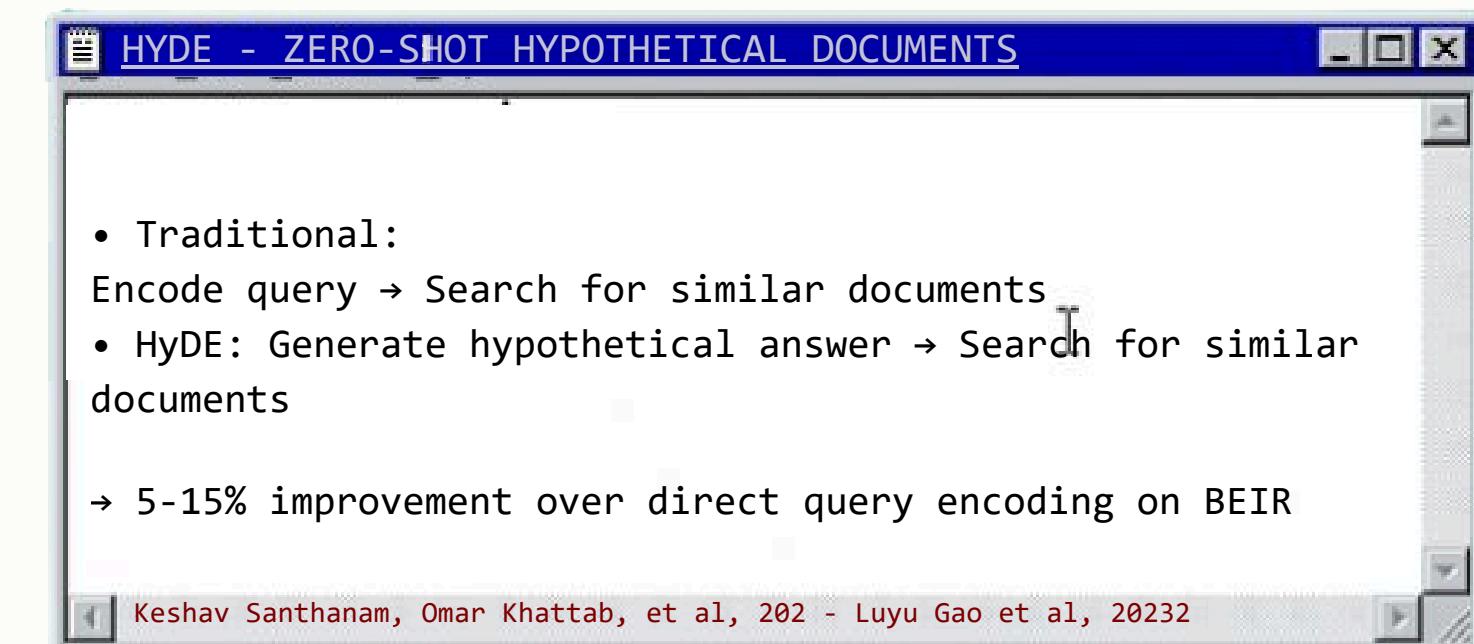
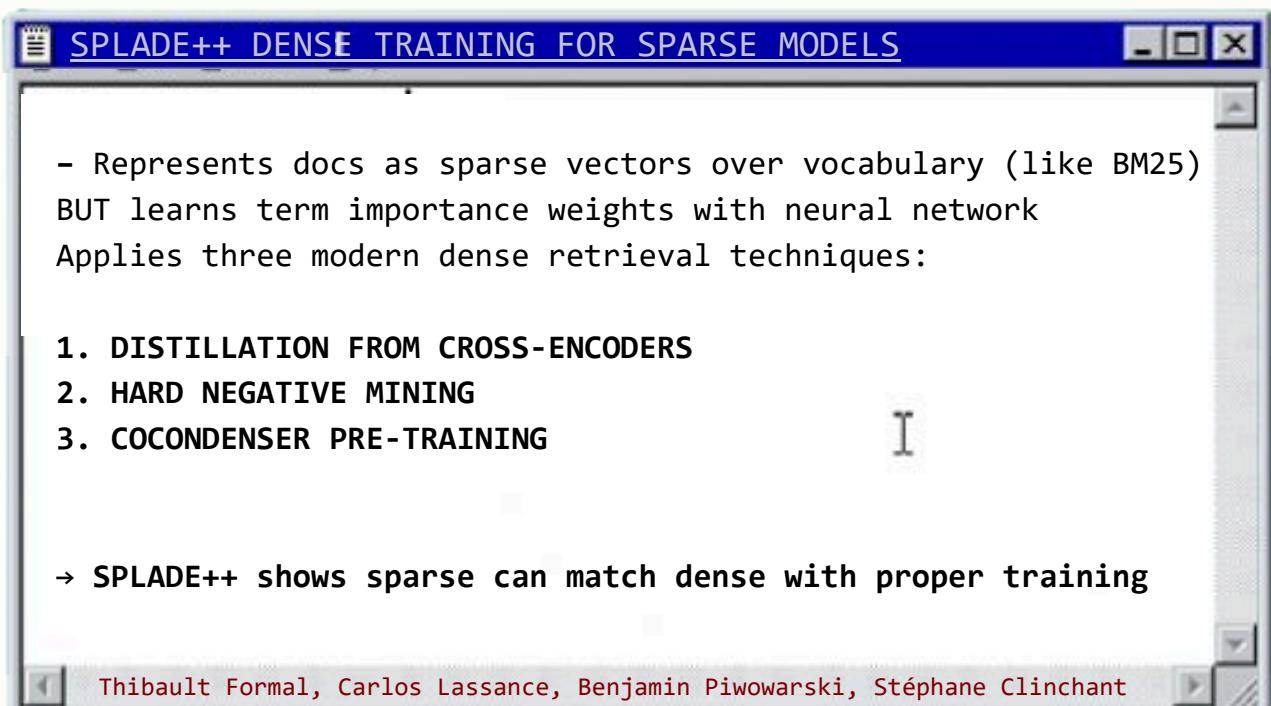
Combine retrieval with LLM generation to reduce hallucinations

CHAMELEON - HARDWARE ACCELERATION:

- THE BOTTLENECK:
 - RAG has two distinct computational stages:
 - Retrieval: Vector search over large database
 - Generation: Transformer inference on LLM
 - Different stages need different hardware optimizations.
 - 1. RETRIEVAL STAGE - FPGA ACCELERATION:
 - FPGAs implement IVF-PQ with near-memory processing
 - Memory organized to maximize bandwidth utilization
 - Filters billions of vectors to find top candidates
 - 2. GENERATION STAGE
 - GPU ACCELERATION: GPUs handle transformer LLM inference
 - Processes retrieved documents + query through LLM
 - CPU COORDINATION: Orchestrates pipeline between FPGA and GPU:
 - Receives queries → dispatches to FPGA
 - Collects retrieved docs → batches for GPU

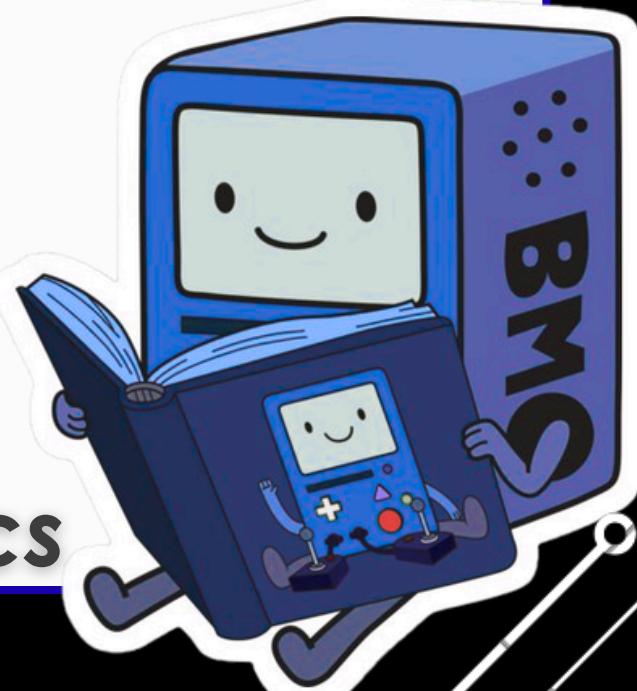
Wenqi Jiang, Marco Zeller, Roger Waleffe, Torsten Hoefer, Gustavo Alonso, 2023

4- HYBRID & ADVANCED METHODS



MAJOR INSIGHTS AND DISCOVERIES

Recent Directions and Emerging Topics



TRAINING METHODOLOGY MATTERS MOST

- Simple architectures with smart training beat complex designs
- DPR's dual-encoder outperforms fancy architectures through better negative sampling
- In-batch negatives, hard negatives, distillation more important than model size

ZERO-SHOT PERFORMANCE IS CRITICAL

- Generalization to new domains as important as in-domain accuracy
- Sparse models (SPLADE++) often generalize better than dense models
- Bottleneck pre-training (SIMLM) improves transfer

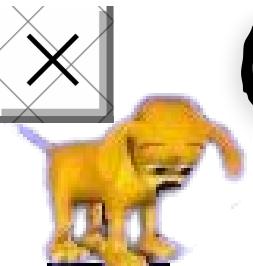
SECURITY CANNOT BE IGNORED

- 0.02% adversarial content = 60% attack success rate
- Production systems need active defenses

CHALLENGES AND FUTURE DIRECTIONS



CHALLENGES



01

SECURITY AND ROBUSTNESS

02

MULTIMODAL CAPABILITIES

03

BIAS AND FAIRNESS

04

SUSTAINABILITY

05

EVALUATION BEYOND BENCHMARKS



01

DEVELOP ROBUST DEFENSES AGAINST
ADVERSARIAL ATTACKS

02

MULTIMODAL RAG INTEGRATING TEXT,
IMAGES, AND CODE

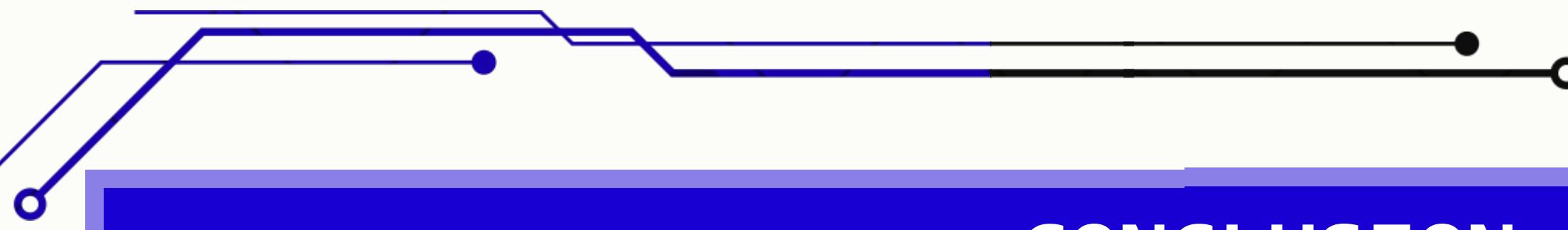
03

OPTIMIZE HARDWARE ACCELERATION
FOR RAG SYSTEMS



CONCLUSION





CONCLUSION

Semantic search has evolved from keyword-based matching to dense vector retrieval and now to tightly integrated retrieval-augmented generation systems. Early advances focused on representation learning and contrastive training. Subsequent research addressed scalability through compression, indexing, and hardware acceleration. More recent work integrates retrieval directly into LLM pipelines, enabling adaptive, grounded generation. However, challenges remain in robustness, security, and billion-scale deployment. The boundary between retrieval and generation is increasingly dissolving, suggesting that future systems will not treat them as separate stages, but as unified, adaptive semantic reasoning frameworks.





MBA519: BIG DATA ANALYTICS



THANK YOU ANY QUESTIONS?

2025-2026