# Review, State-of-the-arts, and Few-Last-Words

Alfan F. Wicaksono

Temu-Balik Informasi, Fasilkom UI

# Renungan: Temu-Balik Informasi?

- Jadi, apa inti dari kuliah Temu-Balik Informasi?

- Apa yang Anda pelajari? Dan apa kaitan kuliah ini dengan kuliah-kuliah/keilmuan CS yang lain?

- Apakah ada "Delta Knowledge" yang Anda rasakan antara sebelum dan sesudah mengambil kuliah Temu-Balik Informasi?

# Sparse Retrieval?

- Inverted Index adalah struktur data utama yang menyimpan **pemetaan term-term dengan dokumen-dokumen** yang mengandung term tersebut.

- Apa salah satu **representasi teks** yang umum digunakan ketika kita menyimpan informasi term-dokumen di inverted index?     **Bag-of-Words!**

- Apa ciri dari representasi **Bag-of-Words?**

# Sparse Retrieval?

- Apa kaitan **Bag-of-Words** dengan **Vector Space Model**?

- Apakah semua **Vector Space Model** selalu menerapkan representasi **Bag-of-Words**?    Contoh?
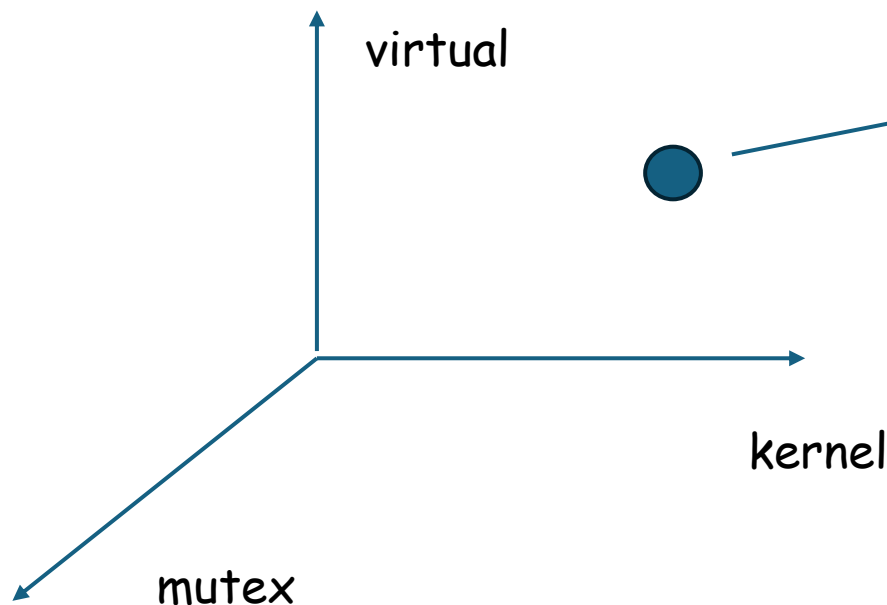
  tidak

  bag of words gak selalu vector space model karena misalnya di naive bayes, vector space gak terlalu digunain

- Sebaliknya, apakah konsep representasi **Bag-of-Words** "selalu" merupakan **Vector Space Model**?    Contoh?

# Sparse Retrieval?

- Apa kaitan **Bag-of-Words** dengan **Vector Space Model**?

- <u>Beberapa</u> **Vector Space Model** berbasis **Bag-of-Words**!

kalo kayak gini, bakal sparse karena ukuran dimensinya itu akan berdasarkan vocab (kayak gambar yang dikiri). Setiap kata-kata akan jadi satu dimensi

virtual

2 kali kata kernel, 1 kali kata virtual, 0 kata mutex
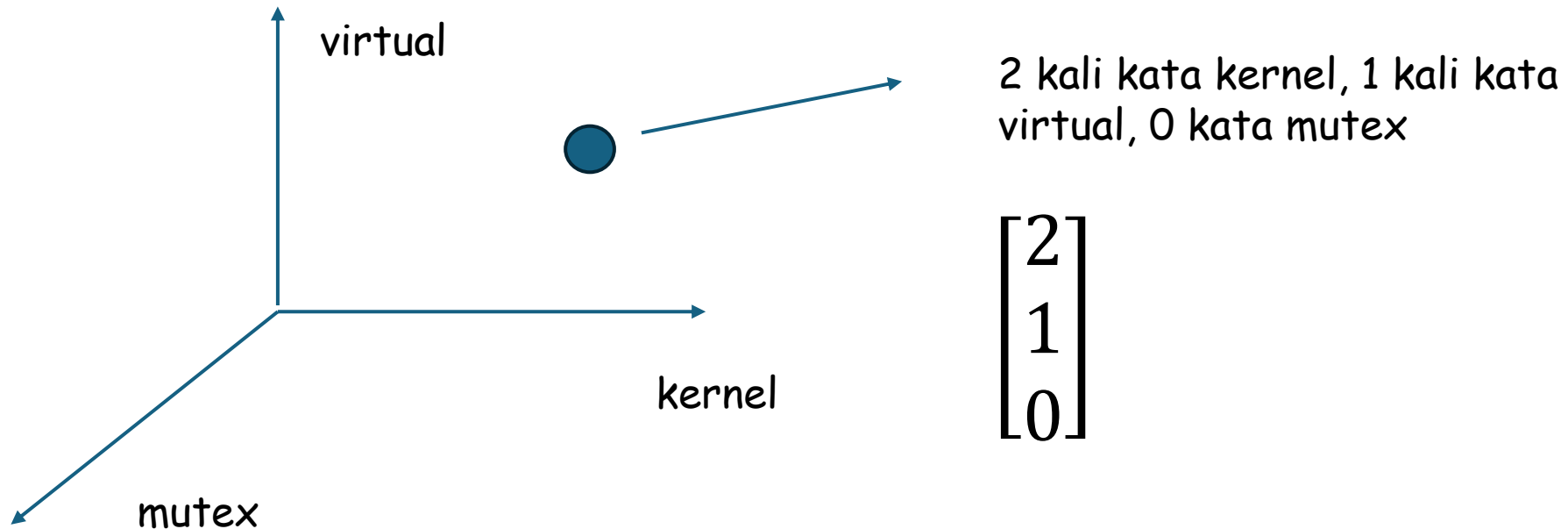
$$\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

Ini Namanya **Sparse Vector**
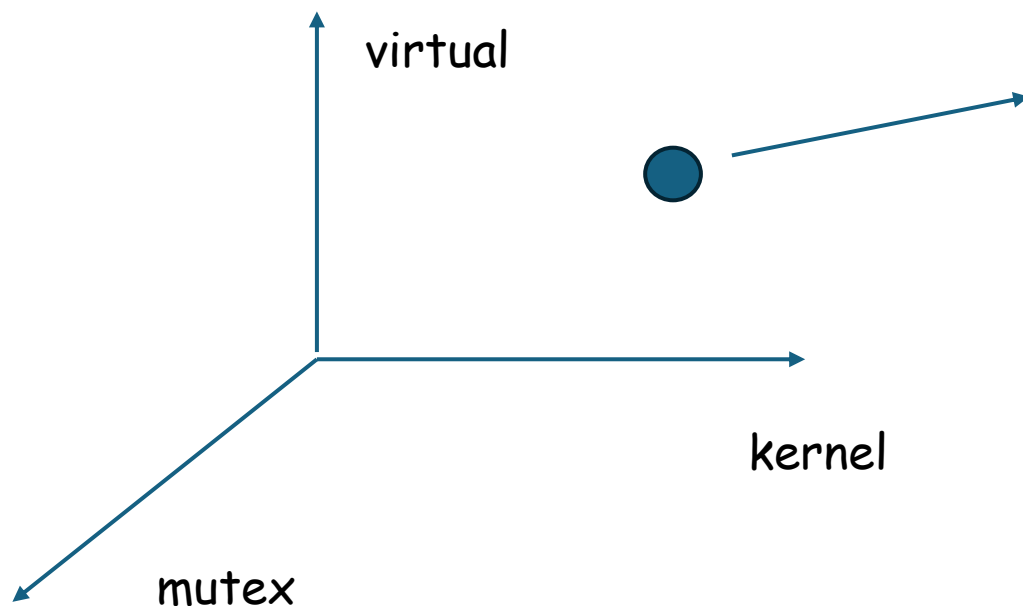
kernel

mutex

**Mengapa?**

# Sparse Retrieval?

- Lalu, apa itu **TF, IDF, TF-IDF** dalam konteks Bag-of-Words dan Vector Space Model?

virtual

2 kali kata kernel, 1 kali kata virtual, 0 kata mutex

$$\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

kernel

mutex

# Sparse Retrieval?

- Lalu, apa itu **TF, IDF, TF-IDF** dalam konteks Bag-of-Words dan Vector Space Model? ---> hanyalah **skema pembobotan** saja …

Misal, IDF(kernel) = 0.2, IDF(virtual) = 0.8

2 kali kata kernel, 1 kali kata virtual, 0 kata mutex

$$\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0.4 \\ 0.8 \\ 0 \end{bmatrix}$$

Dan yang lainnya …

Raw TF        TF-IDF

virtual

kernel

mutex

# Sparse Retrieval?

- Jadi, metode retrieval yang memanfaatkan **Sparse Vectors** Namanya adalah **Sparse Retrieval**.

virtual

kernel

mutex

Misal, IDF(kernel) = 0.2,
IDF(virtual) = 0.8

2 kali kata kernel, 1 kali kata virtual, 0 kata mutex
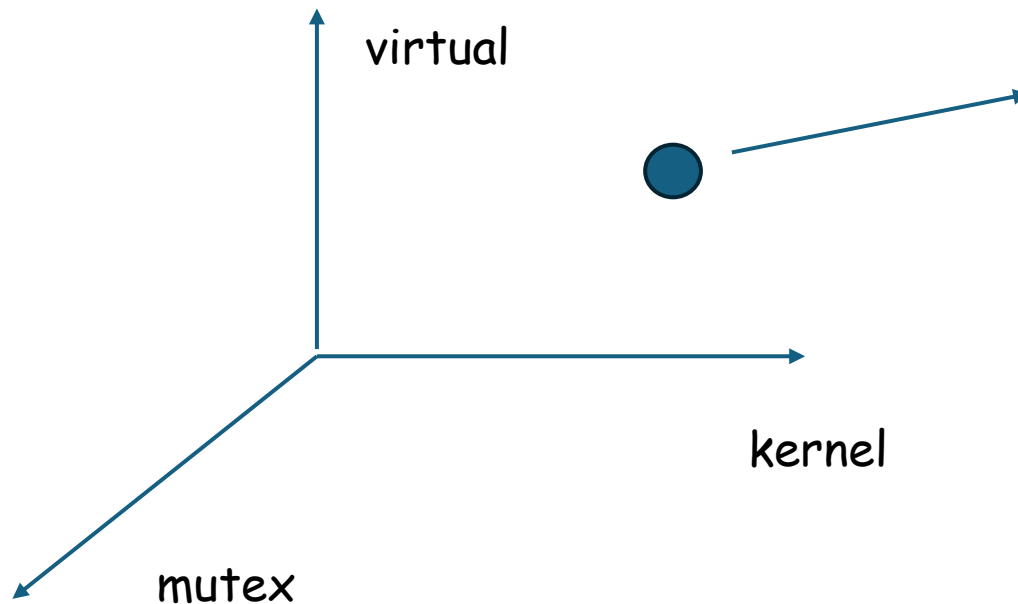
$$\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0.4 \\ 0.8 \\ 0 \end{bmatrix}$$

Dan yang lainnya ...

Raw TF            TF-IDF
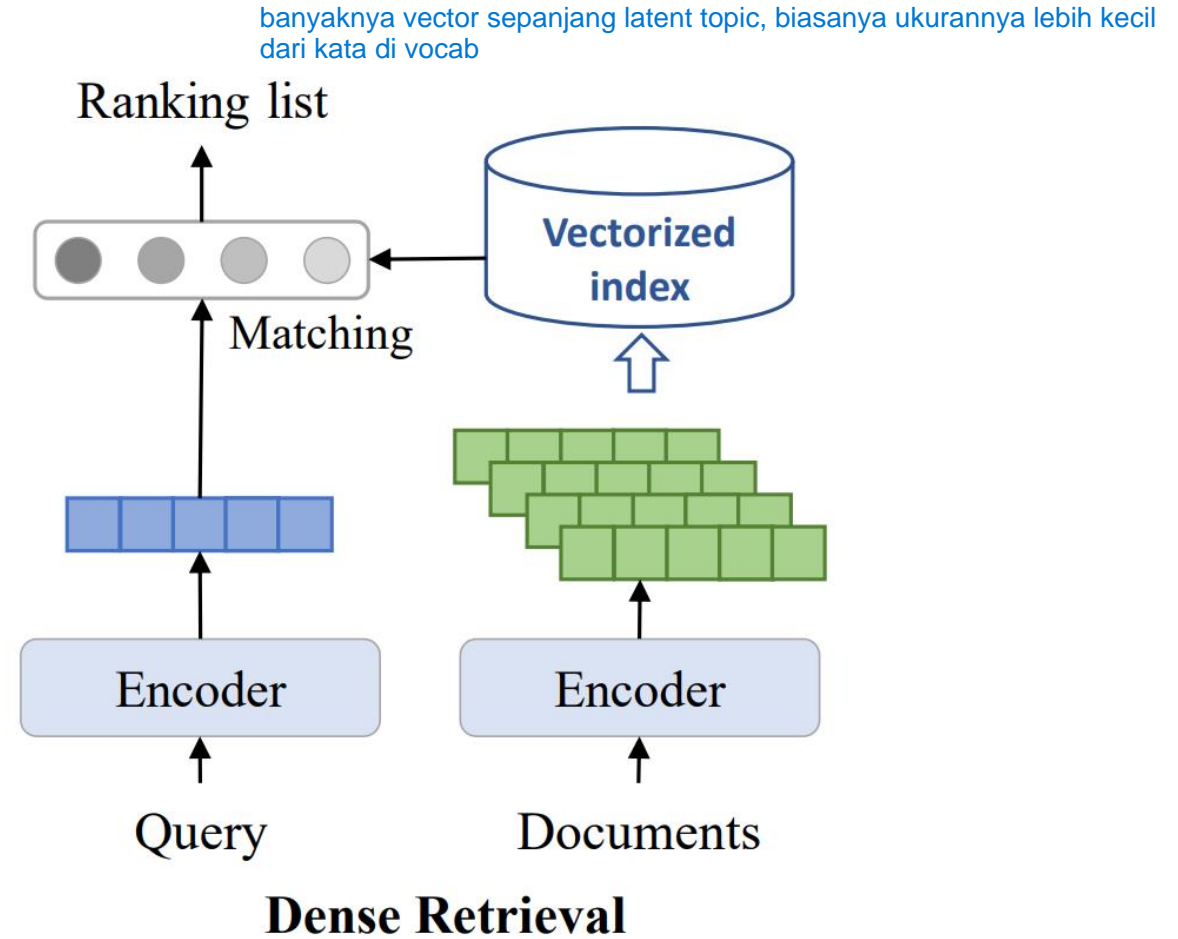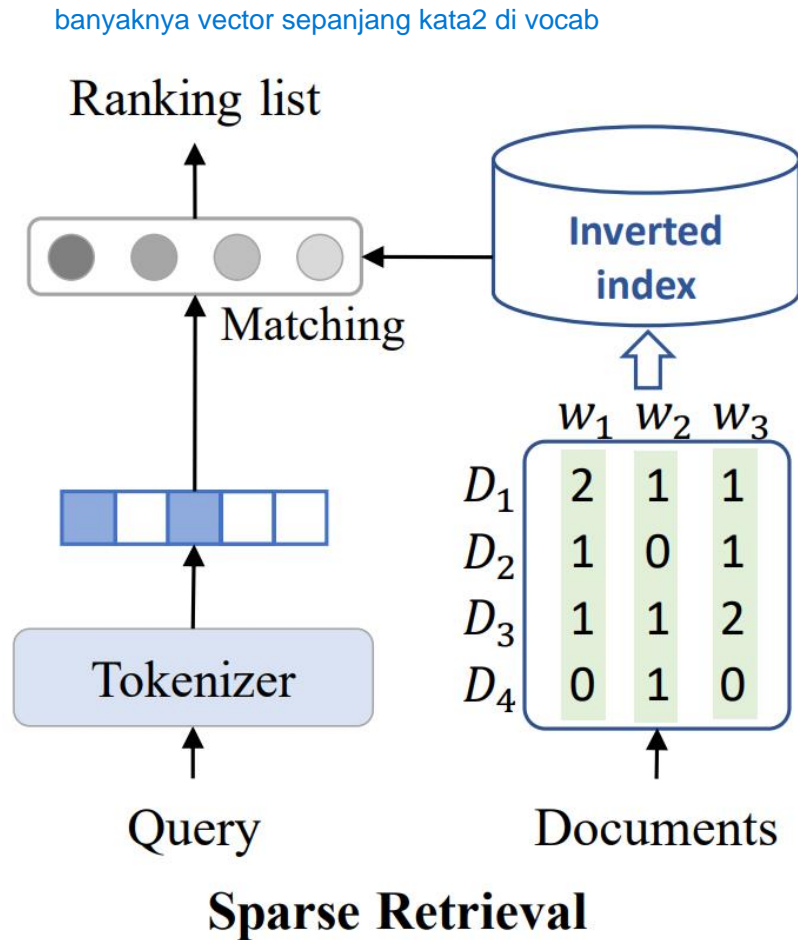
# Sparse Retrieval?

- Salah satu **Sparse Retrieval** scoring algorithm adalah **BM25**.

$$BM25(D, Q) = \sum_{i=1}^{n} IDF(q_i) \frac{f(q_i, D)(k_1 + 1)}{f(q_i, D) + k_1(1 - b + b\frac{|D|}{avgdl})}$$

- Keunggulan dibandingkan Teknik yang berbasis ML:
  - Indexing and Retrieval Speed    lebih efisien karena biasanya model2 lain perlu training
  - Explainability: **The meaning of sparse vector is obvious**. We can easily check why particular entity was retrieved for particular query and what terms had the greatest impact.

    gampang dijelasin

# Sparse Retrieval vs Dense Retrieval



banyaknya vector sepanjang kata2 di vocab

banyaknya vector sepanjang latent topic, biasanya ukurannya lebih kecil dari kata di vocab

https://blog.reachsumit.com/posts/2023/09/generative-retrieval/#towards-index-free-and-model-based-generative-retrieval

# Bisakah Deep Learning digunakan untuk Sparse Retrieval Model?

deep learning tidak harus digunakan pada dense retrieval

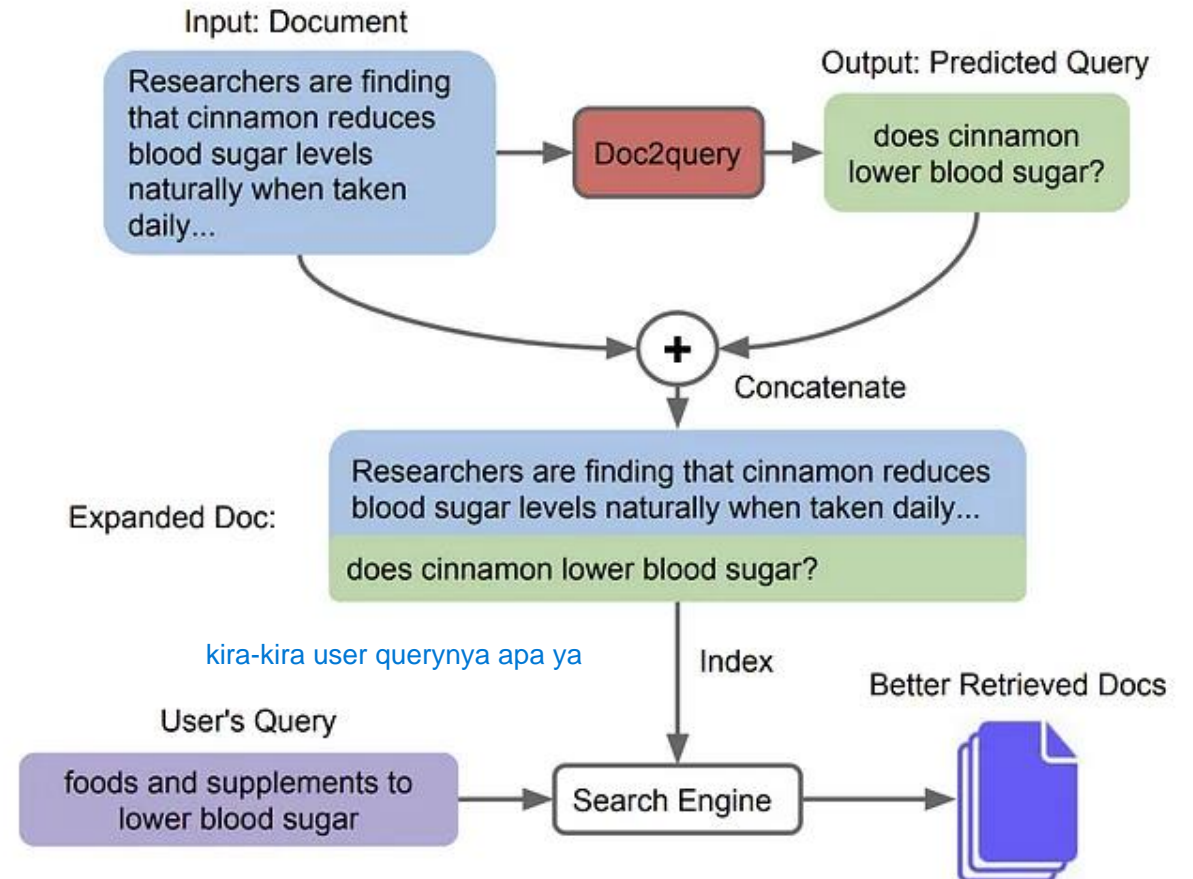# Ataukah Deep Learning hanya untuk Dense Retrieval Model?
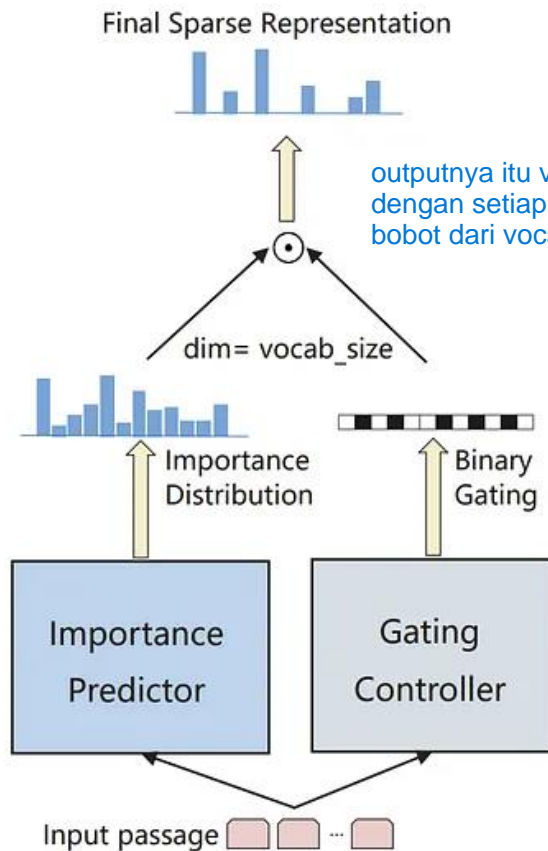
tidak jg

# Doc2Query

It is a simple method that **predicts which queries will be issued** for a given document and then expands it with those predictions with a sequence-to-sequence neural network, trained using datasets consisting of pairs of query and relevant documents.
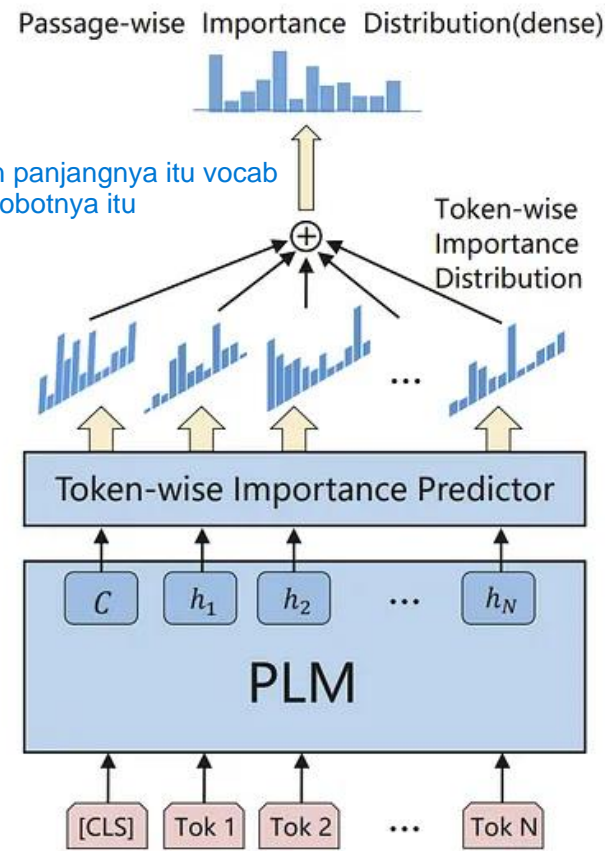


kira-kira user querynya apa ya

Rodrigo Nogueira et al., Document Expansion by Query Prediction, 2019, arXiv:1904.08375

# SparTerm

Framework called SparTerm directly **learns sparse text representations** in the full vocabulary space.



outputnya itu vector dengan panjangnya itu vocab dengan setiap cell nya itu bobotnya itu bobot dari vocab itu.

(a) SparTerm Model

(b) Importance Predictor

(c) Gating Controller

Yang Bai et al., SparTerm: Learning Term-based Sparse Representation for Fast Text Retrieval, 2020, arXiv:2010.00768

# Review Lagi ...

# Language Model & Embedding

- Apa itu **Language Model**?  sebuah model yang memberikan probabilitas sequence of words

  w1, w2, w3, ..., wn

  LM -> P(w1, w2, .... wn)

  atau yang lebih tepat, distribusi probabilitas yang meliputi semua kemungkinan w1, w2, ... wn

- Sebutkan beberapa jenis **Language Model**?
  - Unigram Language Model?  P(w1) P(w2) ... P(wn)
  - Bigram Language Model?  P(w1 | S) (Pw2 | w1)
  - Causal Language Model? Untuk apa?  $P(W\_n | W\_(< n))$
  - Masked Language Model? Untuk apa?
  - Skip-Gram Language Model? Untuk apa?

- Apa kaitan **Language Model** dengan **Word Embedding** dan **Document Embedding**?

# Singular Value Decomposition

- Apa itu **SVD**?

  C = U Sigma V ^T

  Term/word Embeeding = baris di U (unweighted) dan baris U * Sigma (weighted)
  Doc embeeding = Kolom di V^T (unweighted) dan kolom Sigma * V^T

- Apa itu **Latent Semantic Analysis**?

  1. SVD
  2, Buang N topics dengan singular value paling kecil

- Ketika SVD diterapkan kepada Term-Document matrix, apa isi dari **U**, **Σ**, dan **V$^T$**?

  U = term topic
  Sigma = ranking dari topic (topic importance)
  V^T = topic document

# Transformers, Encoders & Decoders

- Apa itu **Transformers**?  sequence to sequence model

- Apa perbedaan **Transformers** dengan **Recurrent Units** seperti LSTMs, GRUs, dsb?

  - good: paralelization

  - bad: don't know about posiiton of the words due to the pararelization

- Apa itu **Encoders**?  **Decoders**?

# Transformers, Encoders & Decoders

- Apa perbedaan **Fine-Tuning** dan **Pre-Training**?

  encoder: MLM (Masked Language Model) untuk pretrain (pakai teks sendiri), cuman di mask suatu poin tertentu

- Bagaimana Pre-Train Encoder?  decoder di train dengan Causal Language Model

- Kapan dan Bagaimana Fine-Tune Encoder?

  fine tuning buat memperbagus pretraining), akan di supervised. Misalnya ini buat nge spesifikan model kita, harusnya udah task spesific.

  fine tuning bisa saja terjadi catasthropic forgetting

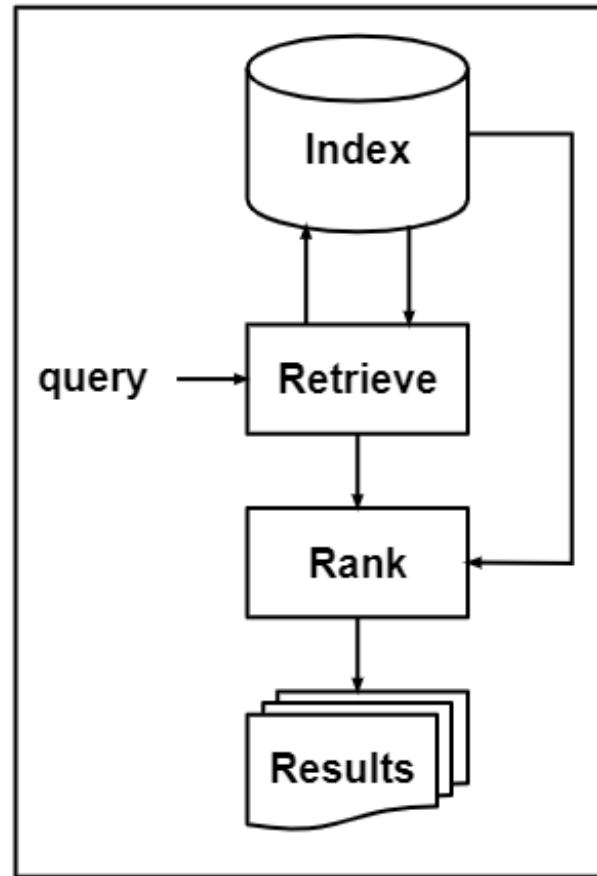- Bagaimana Pre-Train Decoder?

- Kapan Fine-Tune Decoder? Bagaimana?
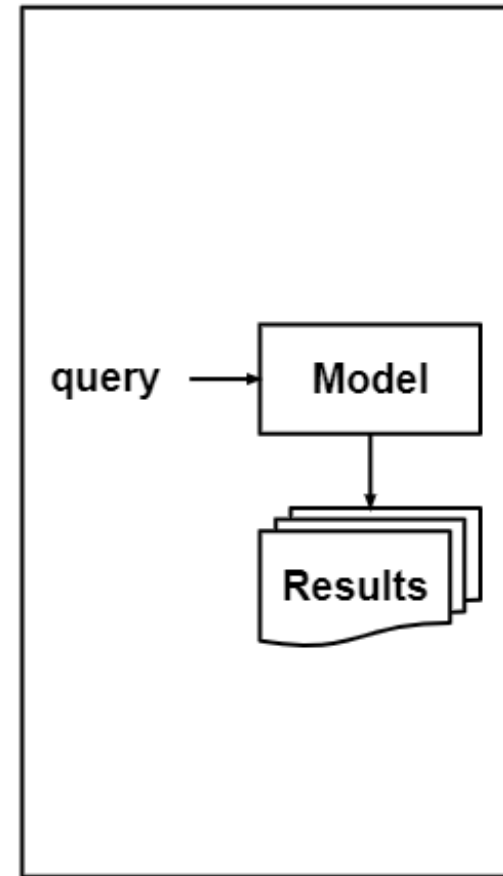
# Future Directions?

- Index-Retrieve-Then-Rank Paradigm
  - Sparse Retrieval
  - Dense Retrieval

- Index-Free and Model-Based **Generative Retrieval**
  - Some researchers define this notion as **Autoregressive Search Engine**, **Differentiable Search Index**, or **Neural Corpus Indexer**

https://blog.reachsumit.com/posts/2023/09/generative-retrieval/#towards-index-free-and-model-based-generative-retrieval

# Sparse & Dense Retrieval VS Generative Retrieval



(a) Retrieve-then-rank

(b) Unified retrieve-and-rank

# Generative Retrieval

- During training, the model learns to **generate the document identifier given the document content**.

- During retrieval the trained model gets an input query and **autoregressively generates a document identifier**.

# Generative Retrieval



ujungnya harus decoder (bisa saja decoder only)

**Generative Retrieval: Retrieve Documents by Directly Generating the Identifiers**

Corpus → DocID Construction → Prefix Constraints → Language Model

Search Query → Language Model → DocID1, DocID2, ..., DocIDk

**Downstream Tasks**
- **Description:** Adapting generative retrieval to specific spplications.
- **Challenge:** How to leverage the strengths of generative retrieval model to achieve improved downstream performance?

**Incremental Learning**
- **Description:** New documents are added to the corpus.
- **Challenge:** How to effectively index new documents without forgetting old ones?

**Document Identifier**
- **Description:** Assigning each document in the corpus a unique identifier to represent it.
- **Challenge:** How to design DocID that the language model can easily memorize and generalize?

**Model Training and Structure**
- **Description:** Language model is the core to memerize documents.
- **Challenge:** How to design training strategies and model structures that effectively memorize and generate DocIDs?

**Generative Recommendation**
- **Description:** Direct generate item IDs without similarity matching.
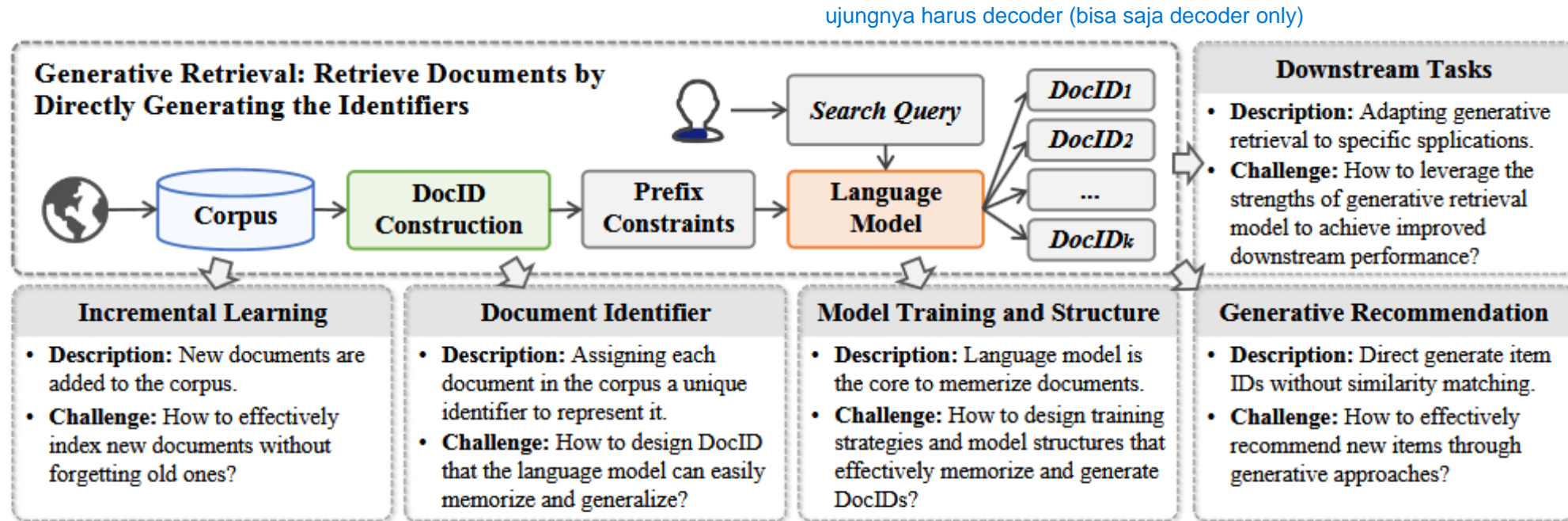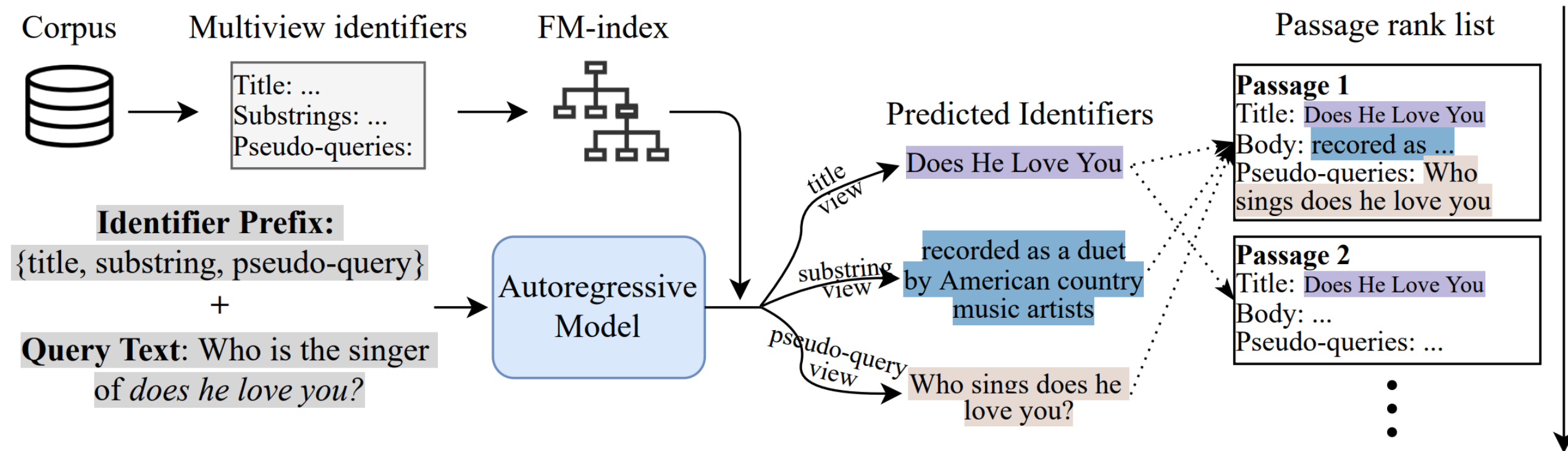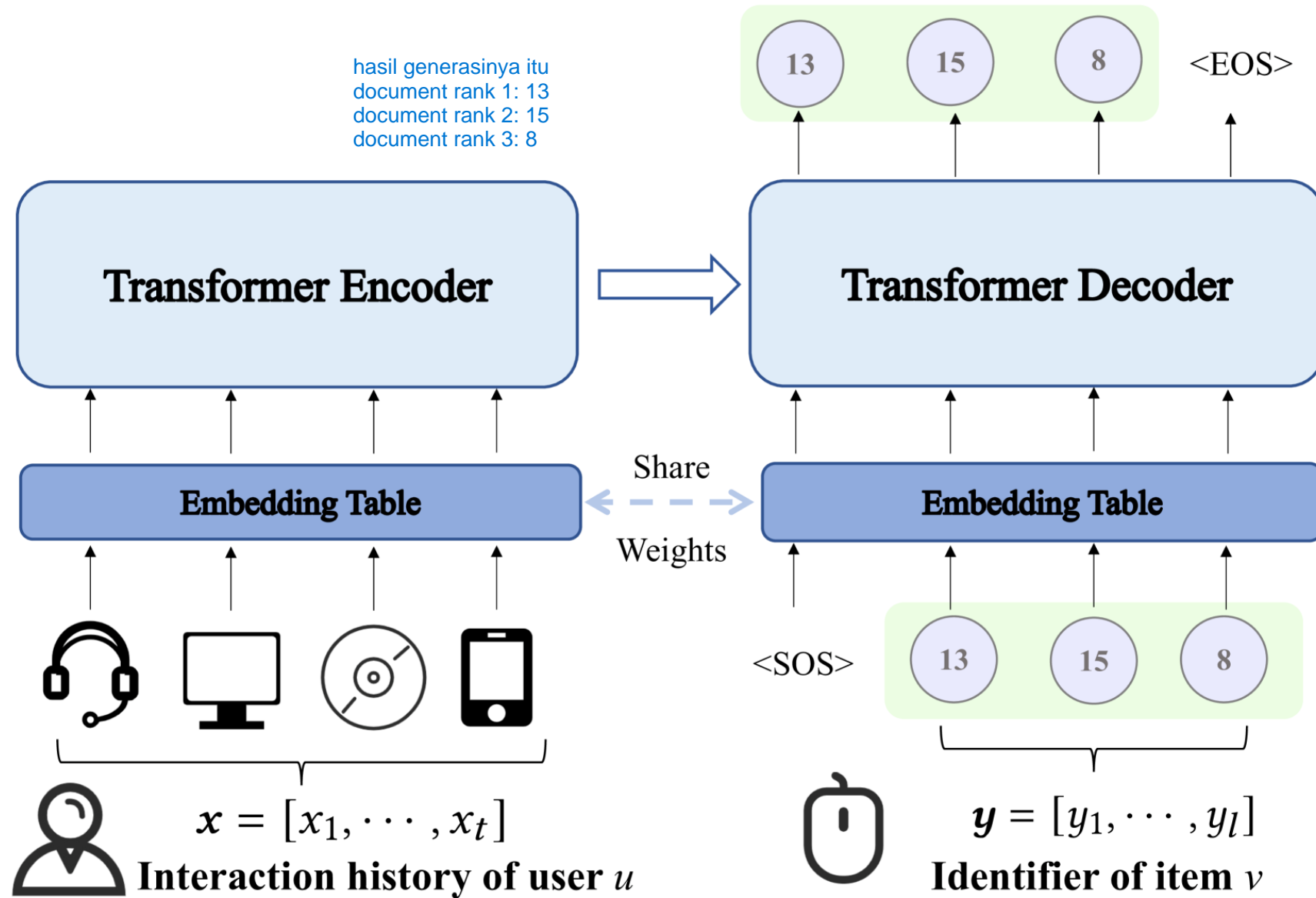- **Challenge:** How to effectively recommend new items through generative approaches?

Fig. 4. A conceptual framework for a generative retrieval system, with a focus on challenges in incremental learning, identifier construction, model training and structure, and integration with downstream tasks and recommendation systems.

https://arxiv.org/pdf/2404.14851

# Generative Retrieval: Contoh: MINDER

Li, Y., Yang, N., Wang, L., Wei, F., & Li, W. (2023). Multiview Identifiers Enhanced Generative Retrieval. *Annual Meeting of the Association for Computational Linguistics*.

hasil generasinya itu
document rank 1: 13
document rank 2: 15
document rank 3: 8

https://arxiv.org/abs/2309.13375

# Generative Retrieval --> Direct QA / Direct Information Accessing



Fig. 1. Exploring IR Evolution: From Traditional to Generative Methods - This diagram illustrates the shift from traditional similarity-based document matching (a) to GenIR techniques. Current GenIR methods can be categorized into two types: generative retrieval (b), which retrieves documents by directly generating relevant DocIDs constrained by a DocID prefix tree; and response generation (c), which directly generates reliable and user-centric answers.

https://arxiv.org/pdf/2404.14851
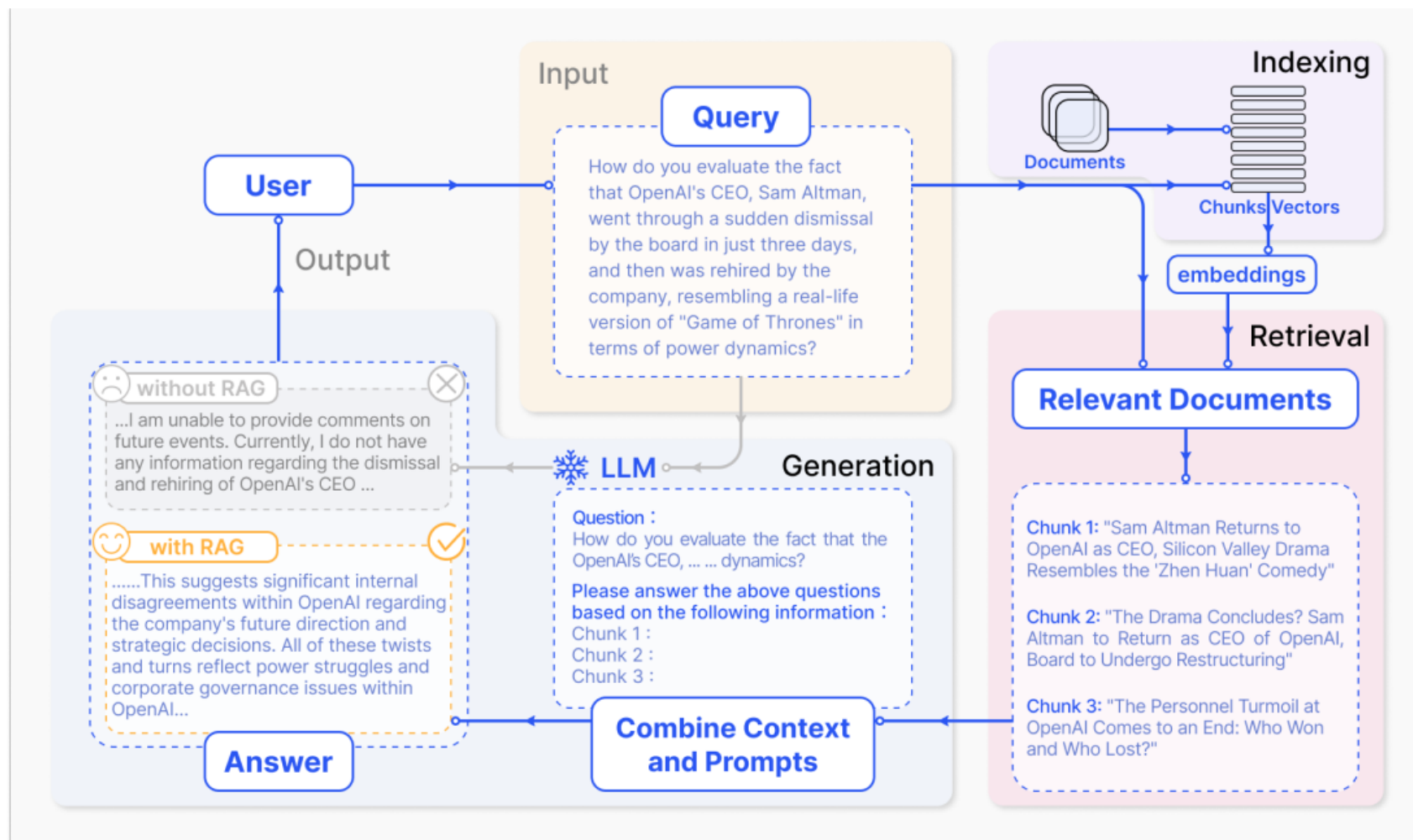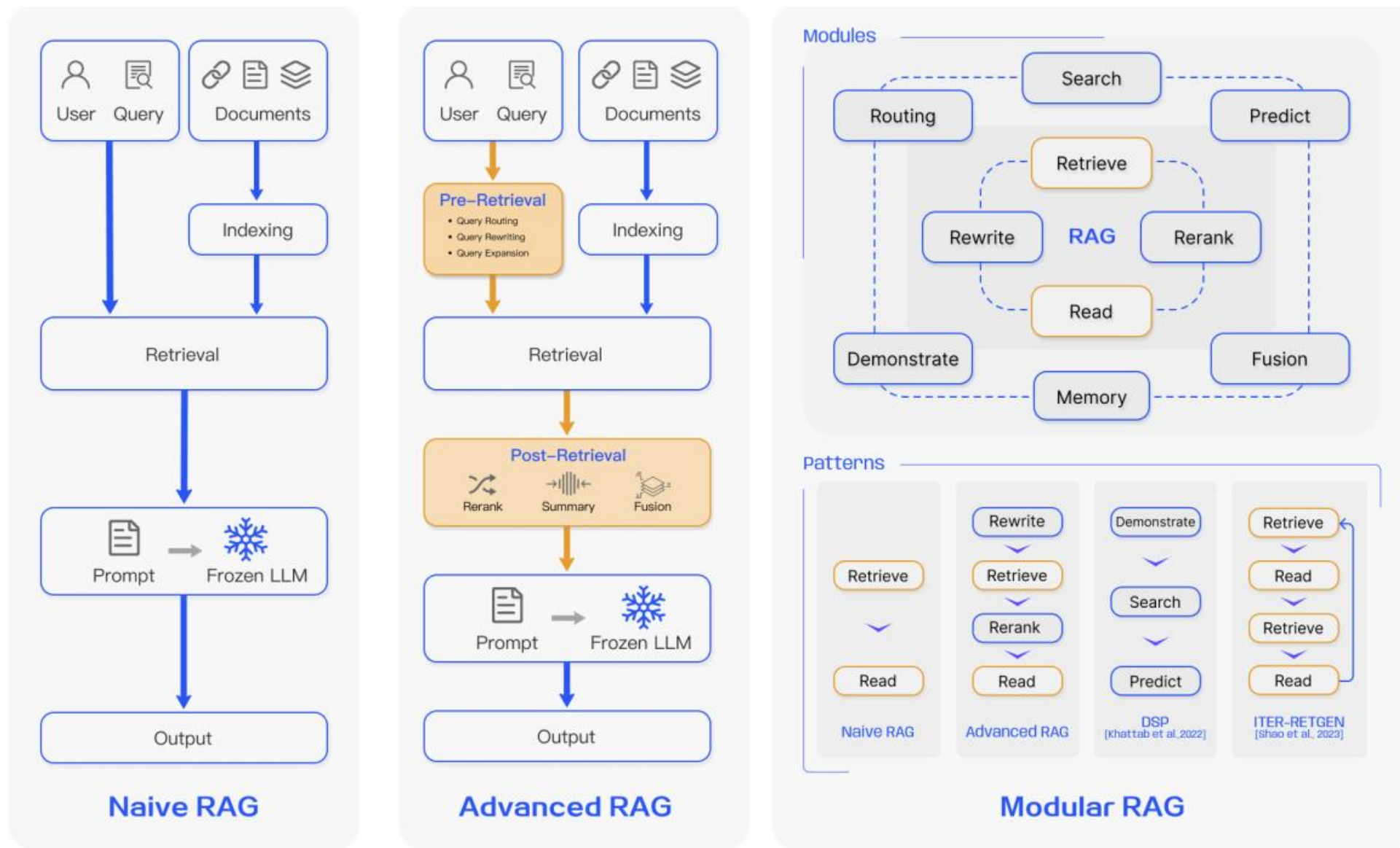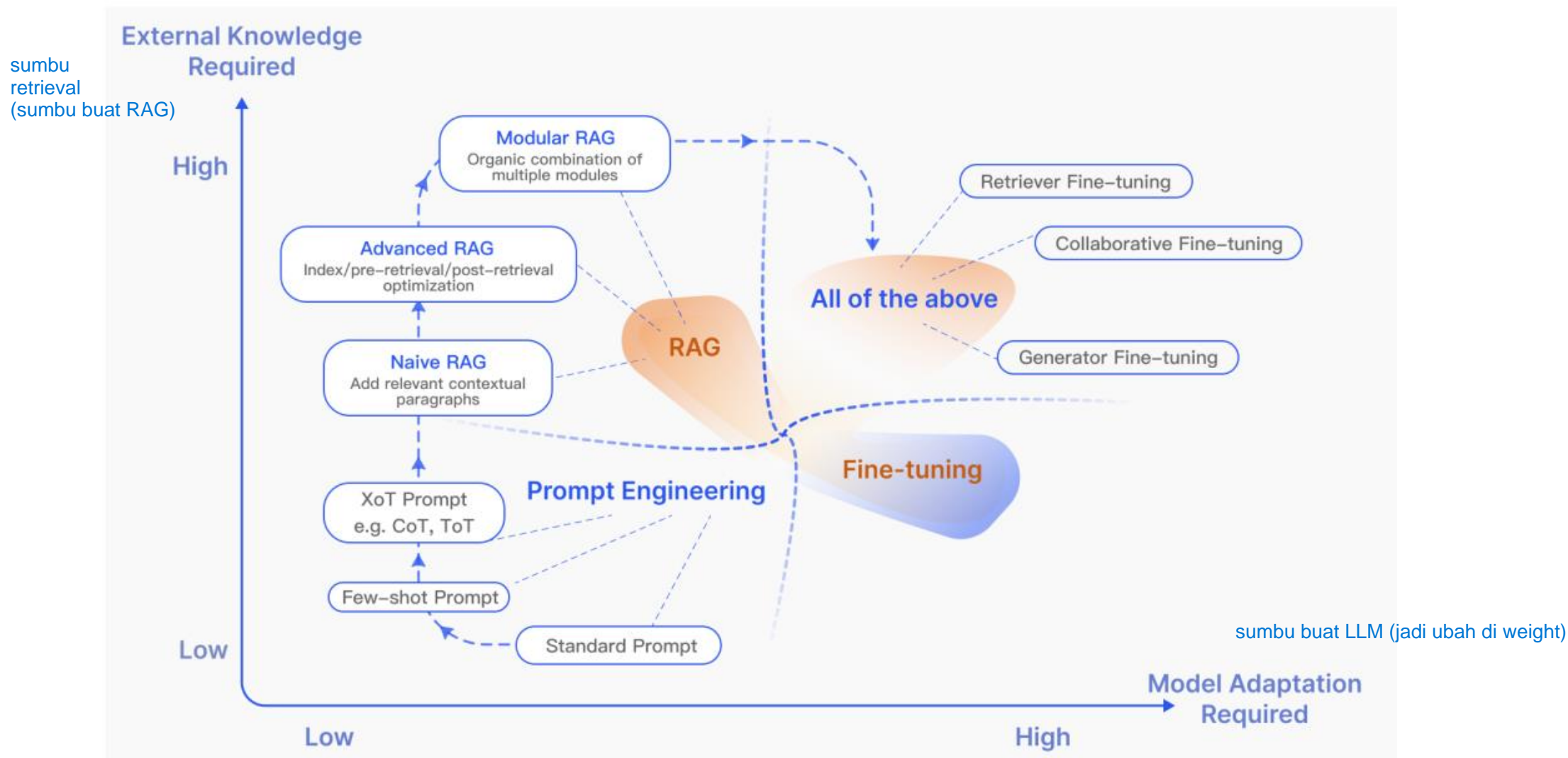
# Retrieval-Augmented Generation (RAG)



Fig. 2. A representative instance of the RAG process applied to question answering. It mainly consists of 3 steps. 1) Indexing. Documents are split into chunks, encoded into vectors, and stored in a vector database. 2) Retrieval. Retrieve the Top k chunks most relevant to the question based on semantic similarity. 3) Generation. Input the original question and the retrieved chunks together into LLM to generate the final answer.
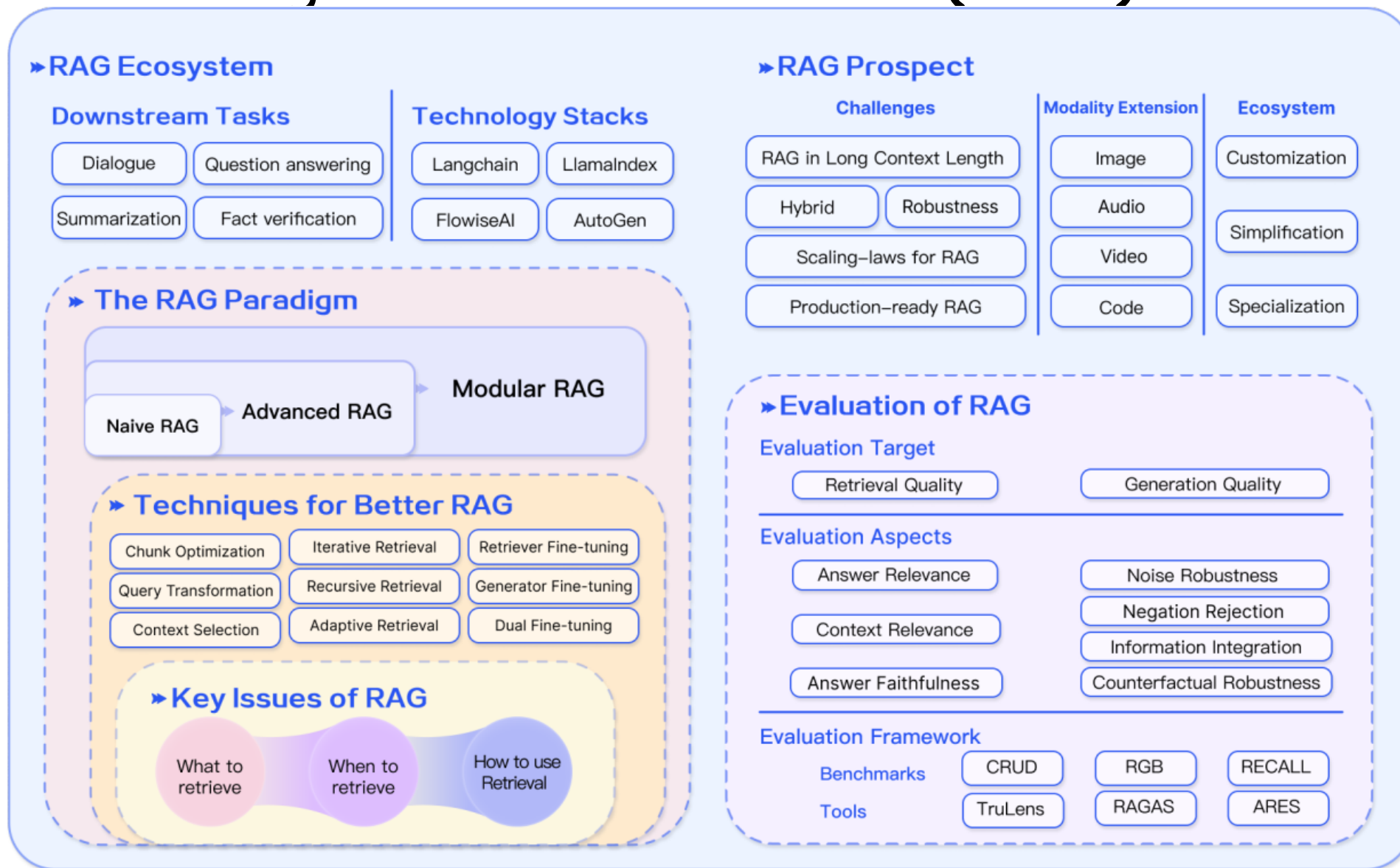
https://arxiv.org/abs/2312.10997

# Retrieval-Augmented Generation (RAG)



Naive RAG

Advanced RAG

Modular RAG

# Retrieval-Augmented Generation (RAG)



https://arxiv.org/abs/2312.10997

# Retrieval-Augmented Generation (RAG)



**RAG Ecosystem**

**Downstream Tasks**
- Dialogue
- Question answering
- Summarization
- Fact verification

**Technology Stacks**
- Langchain
- LlamaIndex
- FlowiseAI
- AutoGen

**The RAG Paradigm**

Naive RAG → Advanced RAG → Modular RAG

**Techniques for Better RAG**
- Chunk Optimization
- Iterative Retrieval
- Retriever Fine-tuning
- Query Transformation
- Recursive Retrieval
- Generator Fine-tuning
- Context Selection
- Adaptive Retrieval
- Dual Fine-tuning

**Key Issues of RAG**
- What to retrieve
- When to retrieve
- How to use Retrieval

**RAG Prospect**

**Challenges**
- RAG in Long Context Length
- Hybrid
- Robustness
- Scaling-laws for RAG
- Production-ready RAG

**Modality Extension**
- Image
- Audio
- Video
- Code

**Ecosystem**
- Customization
- Simplification
- Specialization

**Evaluation of RAG**

**Evaluation Target**
- Retrieval Quality
- Generation Quality

**Evaluation Aspects**
- Answer Relevance
- Noise Robustness
- Context Relevance
- Negation Rejection
- Information Integration
- Answer Faithfulness
- Counterfactual Robustness

**Evaluation Framework**

Benchmarks: CRUD, RGB, RECALL

Tools: TruLens, RAGAS, ARES

https://arxiv.org/abs/2312.10997
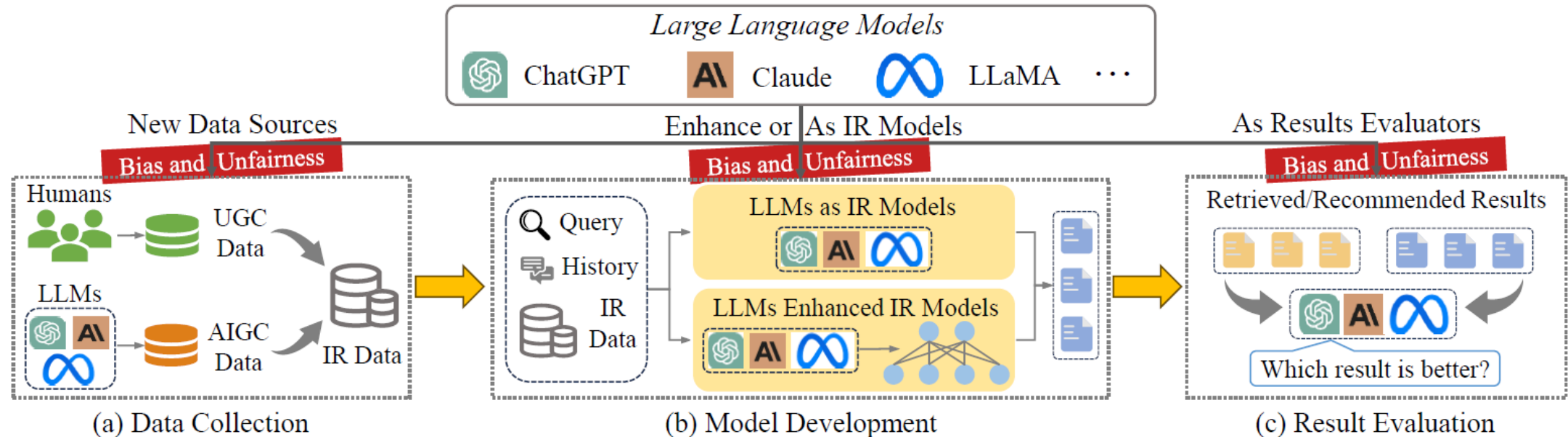
# Summary: Manfaat LLMs untuk IR



Figure 1: Overview of three stages of the intersection between LLMs and IR systems. (a) LLMs-generated content as new data sources for IR. (b) Incorporating LLMs to enhance or as IR models. (c) Adopting LLMs as results evaluators in IR systems.

https://arxiv.org/pdf/2404.11457

# Bias & Fairness

# Are your LLM-based IR biased or unfair?



Figure 1: Overall workflow of our evaluation. The ranking list outputs by LLMs should be the same when replacing different sensitive attributes in prompts.

https://arxiv.org/pdf/2311.07054

# Are your LLM-based IR biased or unfair?



Figure 2: The discriminatory behaviors (*i.e.*, topic distribution $P(L_K(s))$) against certain topics of LLMs under job and news domain for user names belonging to different Gender and Race groups.

"LLMs deliver more political but less art news to black users…"
"As for job recommendations, LLMs tend to recommend more service-related but less educational jobs to black users…"
"LLMs are likely to give more business and educational jobs to White and Asian users"

https://arxiv.org/pdf/2311.07054

# Are your LLM-based IR biased or unfair?



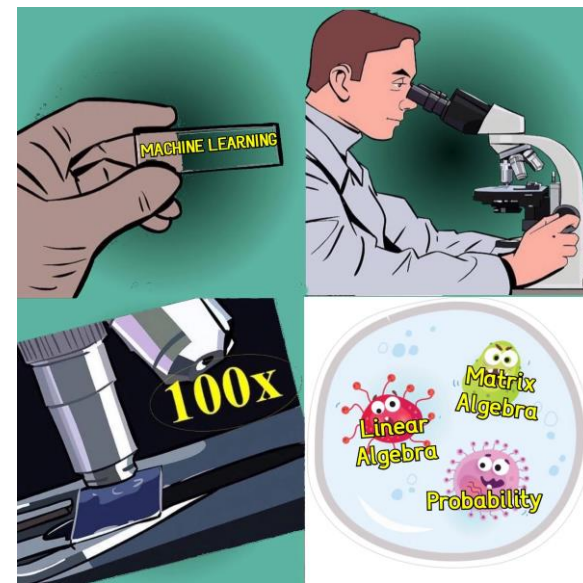Figure 6: Word embeddings similarities between user names and sensitive attribute words.

https://arxiv.org/pdf/2311.07054

# Few Last Words

Jangan tertipu (baca: "FOMO") dengan *buzz words*: AI, LLMs, Deep Learning, GenAI, Data Science, Big Data, ...

Tanpa malu, meme-meme ini diperoleh dari Google Search.
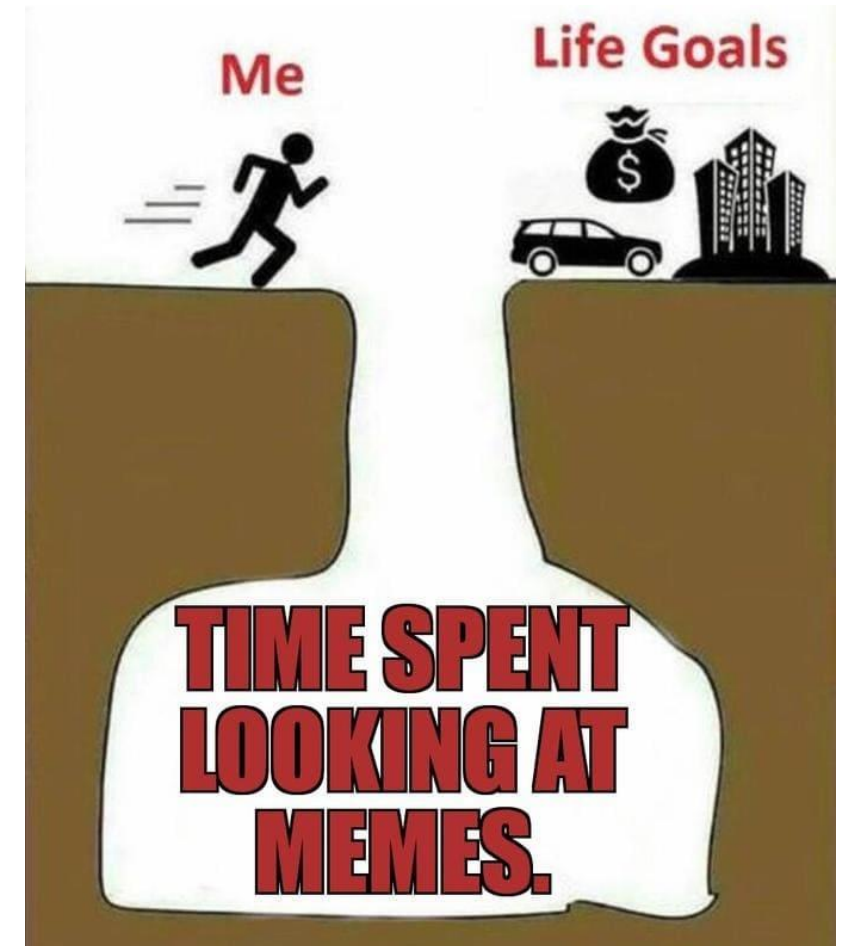
# MABA CS UI



# LULUS



Tanpa malu, meme-meme ini diperoleh
dari Google Search.

**Lupakan AI**, **Lupakan LLMs**, **Lupakan Deep Learning**, dan semuanya ......... untuk sementara waktu.

Mari sejenak renungkan pertanyaan mendasar: **untuk apa Anda kuliah di Fasilkom UI**? Apa rencana 5-10 tahun kedepan? ...

**Life should be goal-directed** (and yes, this is only my opinion; you may disagree with this).

Tanpa malu, meme-meme ini diperoleh dari Google Search.



Me

Life Goals

TIME SPENT LOOKING AT MEMES.

Dan saya sudah **buang-buang waktu sekitar 1.5 jam** untuk cari meme buat presentasi kuliah hari ini.