# Building Multimodal Search and RAG

- MULTIMEDIA CONTENT IS ALL AROUND US;

- DATA FROM DIFFERENTS SOURCES;

- TRAINING MULTIMODAL MODELS: START WITH SPECIALIST MODELS;

  (a) TEXT ENCODER;           (c) AUDIO ENCODER;

  (b) IMAGE ENCODER;          (d) VIDEO ENCODER;

- SIMILAR CONCEPTS = SIMILAR VECTORS;

  (DATA)                    (VECTORS REPRESENTATIONS)

  ↳ UNIFY THE SPECIALIST MODELS;

- UNIFY THE MODELS USING CONTRASTIVE LEARNING:

  ↳ PROCESS TO TRAIN ANY EMBEDDING MODEL;
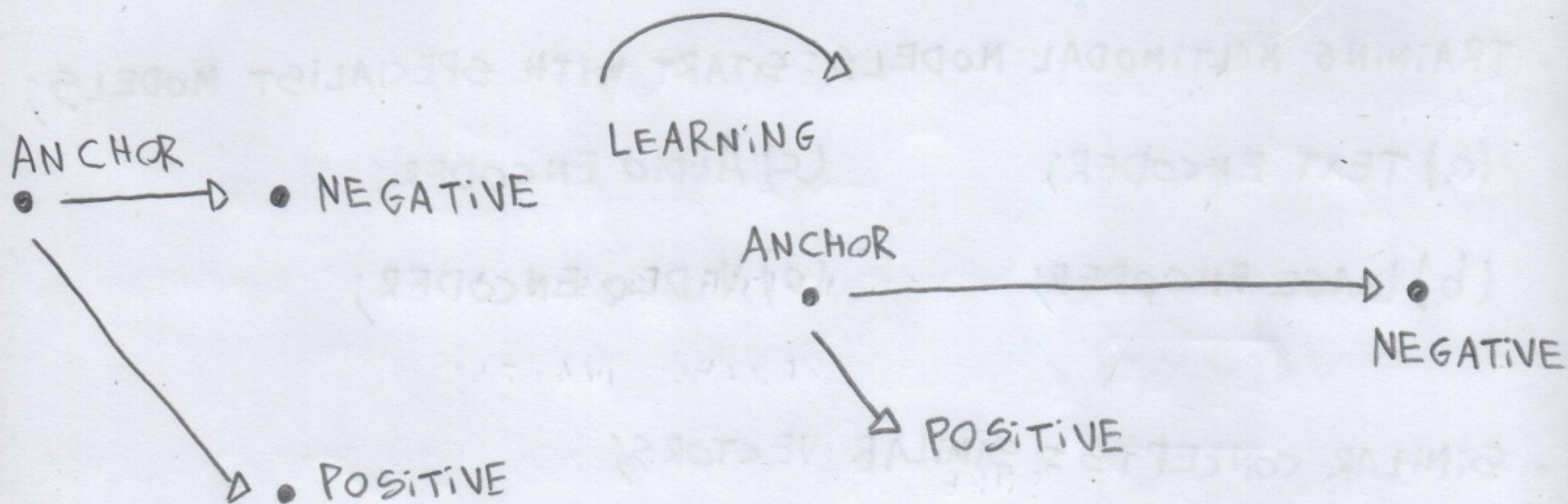
  ↳ UNIFY MULTIPLE MODELS;

  ↳ CREATE ONE VECTOR SPACE;

  ↳ TUNE MODELS BY PROVIDING CONTRASTIVE EXAMPLES;

- E.g: ANCHOR → "HE COULD SMELL THE ROSES"

    POSITIVE EXAMPLE → "A FIELD OF FRAGRANT FLOWERS"

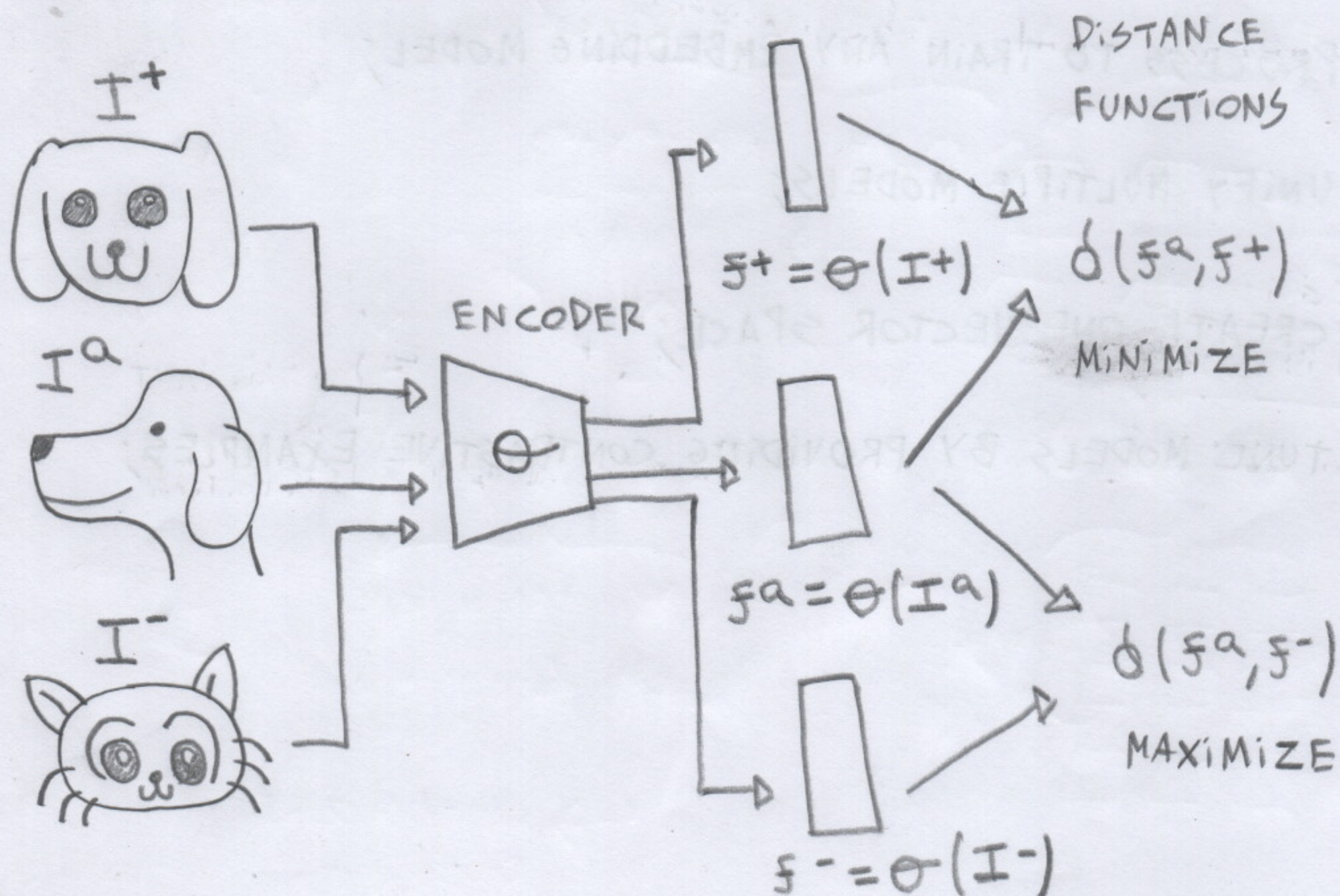    NEGATIVE EXAMPLE → "THE LION ROARED MAJESTICALLY"

LEARNING

ANCHOR → • NEGATIVE

ANCHOR •————————→• NEGATIVE

• POSITIVE

• POSITIVE

↳ SAME EXAMPLE WITH IMAGES, AUDIO AND VIDEO;

↳ PUSH NEGATIVE EXAMPLE;
↳ PULL POSITIVE EXAMPLE;

↳ CONTRASTIVE LOSS FUNCTION

DISTANCE FUNCTIONS

$I^+$

$I^a$

$I^-$

ENCODER

$\theta$

$\xi^+ = \theta(I^+)$

$\delta(\xi^a, \xi^+)$

MINIMIZE

$\xi^a = \theta(I^a)$

$\delta(\xi^a, \xi^-)$

MAXIMIZE

$\xi^- = \theta(I^-)$

E.g.: APPLY CONTRASTIVE LOSS IN TEXT AND IMAGE MULTIMODAL DATA.

LION IMAGE

LIONS.MP4

"LION IS THE KING OF THE JUNGLE"

"MEERKATS DON'T CONSUME WATER"

DOG IMAGE

PEPPER THE AUSSIE PUP → TEXT ENCODER → | T1 | T2 | T3 |

DOG IMAGE → IMAGE ENCODER →

| I1 |
| I2 |
| I3 |

| I1.T1 | I1.T2 | I1.T3 |
| I2.T1 | I2.T2 | I2.T3 |
| I3.T1 | I3.T2 | I3.T3 |

DIAGONAL IS POSITIVES EXAMPLES

$q_i = f(\text{"IMAGE OF LION"})$ AND $K_i = g(\text{"VIDEO OF LION"})$

↪ FURTHEST FROM ZERO

$$L_{I,M} = -\text{LOG}\left( \frac{\text{EXP}\left( q_i^T K_i / \tau \right)}{\text{EXP}\left( q_i^T K_i / \tau \right) + \sum_{j \neq i} \text{EXP}\left( q_i^T K_j / \tau \right)} \right)$$

↳ OBJECTIVE: MINIMIZE FUNCTION

↳ SHOULD BE CLOSER TO ZERO

- WE CAN INFER INFORMATION FROM ANOTHER MODALITY;
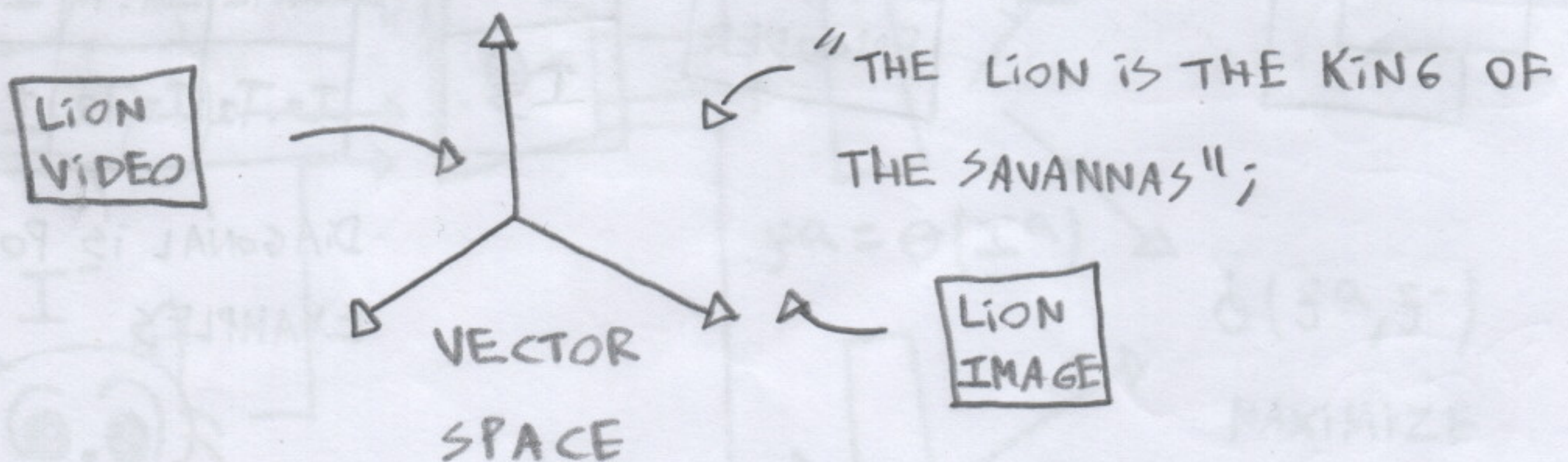
- MULTIMODAL REASONING AND RETRIEVAL:

  - UNDERSTAND WITH ALL OUR SENSES;

  - MULTIMODAL UNDERSTANDING;

  - MULTIMODAL REASONING WORK IN ALL DIRECTIONS;

- MULTIMODAL EMBEDDINGS MODELS
  ↳ SHARED MULTIMODAL VECTOR SPACE;

LION VIDEO

"THE LION IS THE KING OF THE SAVANNAS";

VECTOR SPACE

LION IMAGE

# Any to Any Search

Query                           Retrieved

TEXT              VECTOR        TEXT
IMAGE   →         SPACE    →    IMAGE
AUDIO                           AUDIO
VIDEO                           VIDEO

                                          ↱ CLOSE TO EACH
                                            OTHER FOR THE
                                            DOMAIN
"KING OF THE                          $[0,23 \quad 0,45 \ldots 0,84]$
JUNGLE"
           ↘
            [MULTIMODAL           ↗              ↕ LOSS DISTANCE
             MODEL]
           ↗                      ↘      $[0,26 \quad 0,31 \ldots 0,12]$
"LION
IMAGE"
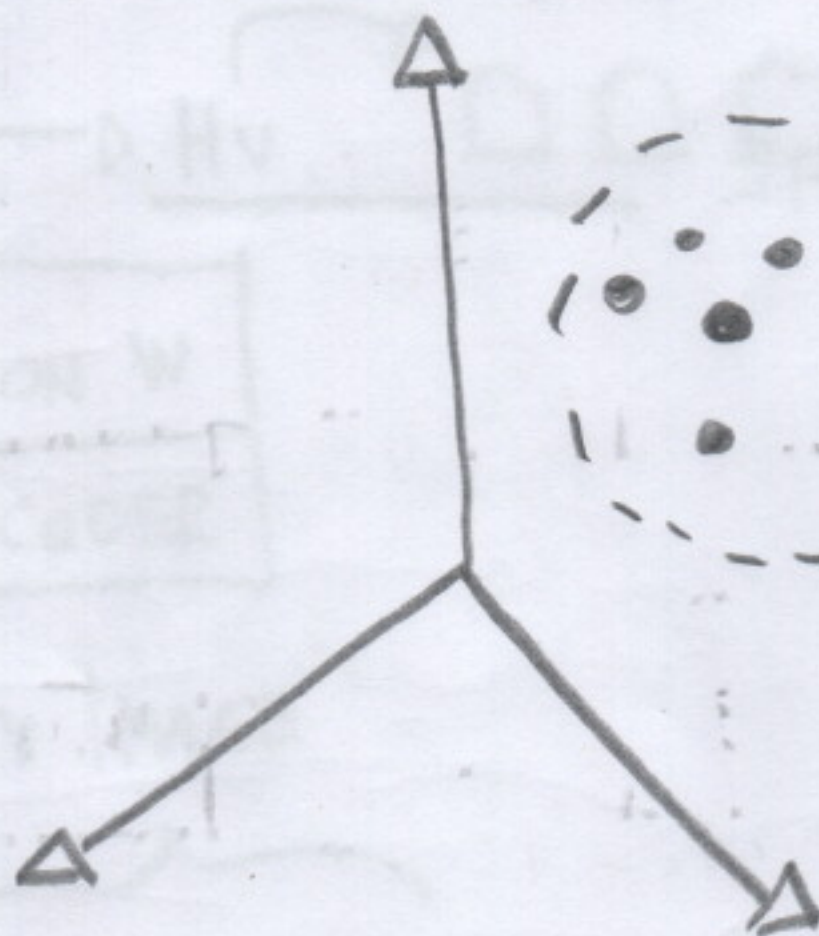                                        EMBEDDINGS
                                        VECTORS



                              ↳ VECTOR SEARCH

- How do large languages models work?
  - ↳ The majority of current LLMs are generative pre-trained transformers (GPT).

    - Autoregressive: they generate text one token at a time;
    - Future tokens are conditioned only on previously provided or generated tokens;
    - Unsupervised training using next token prediction on trillions of tokens.
    - Probability distribution generated over tokens: the next token can be sampled from this distribution.

E.g.: Jack and Jill went up the _____

APPLE LLAMA LARGE |MOUNTAIN| BIG PEANUT |HILL| WATCH

$$[1, 2, 3, 7, ..., 3, 4, 9, 1]$$

PREDICTION SCORE

- HOW DO GPT MODELS WORK?

PROMPT: THE ROCK
        └→ TOKEN PROBABILITY
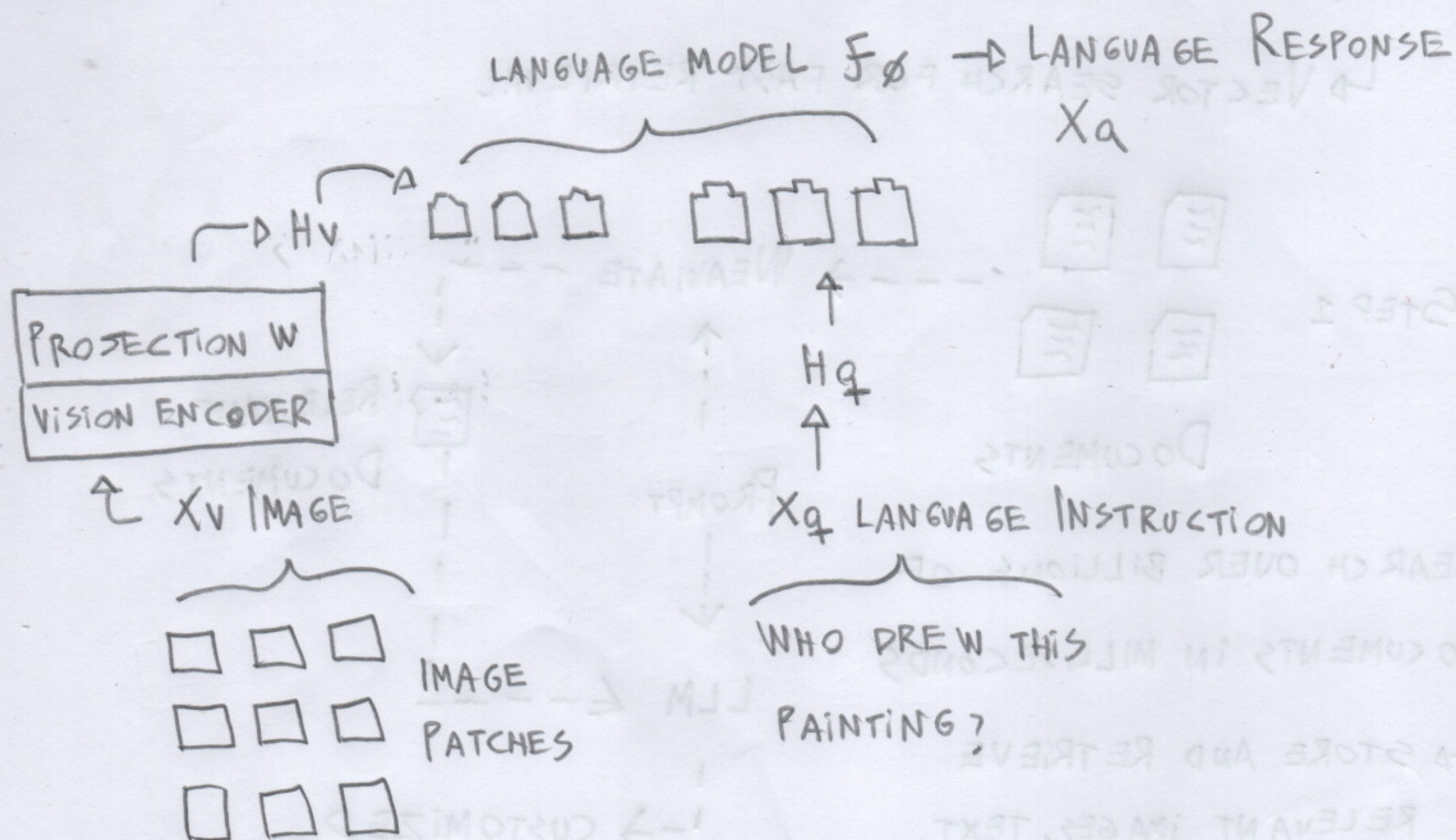        └→ EMBEDDING LOOKUP

• VISION TRANSFORMERS

- IMAGES ARE CUT UP INTO PATCHES.
    └→ USING PATCHES INSTEAD OF PIXELS MAKE IT COMPUTATIO-
        NALLY EFFICIENT TO PROCESS IMAGES;

• EACH PATCH IS EMBEDDED AND PASSED INTO A TRANSFORMER.

• THE TRANSFORMER OUTPUTS A PROBABILITY DISTRIBUTION OVER
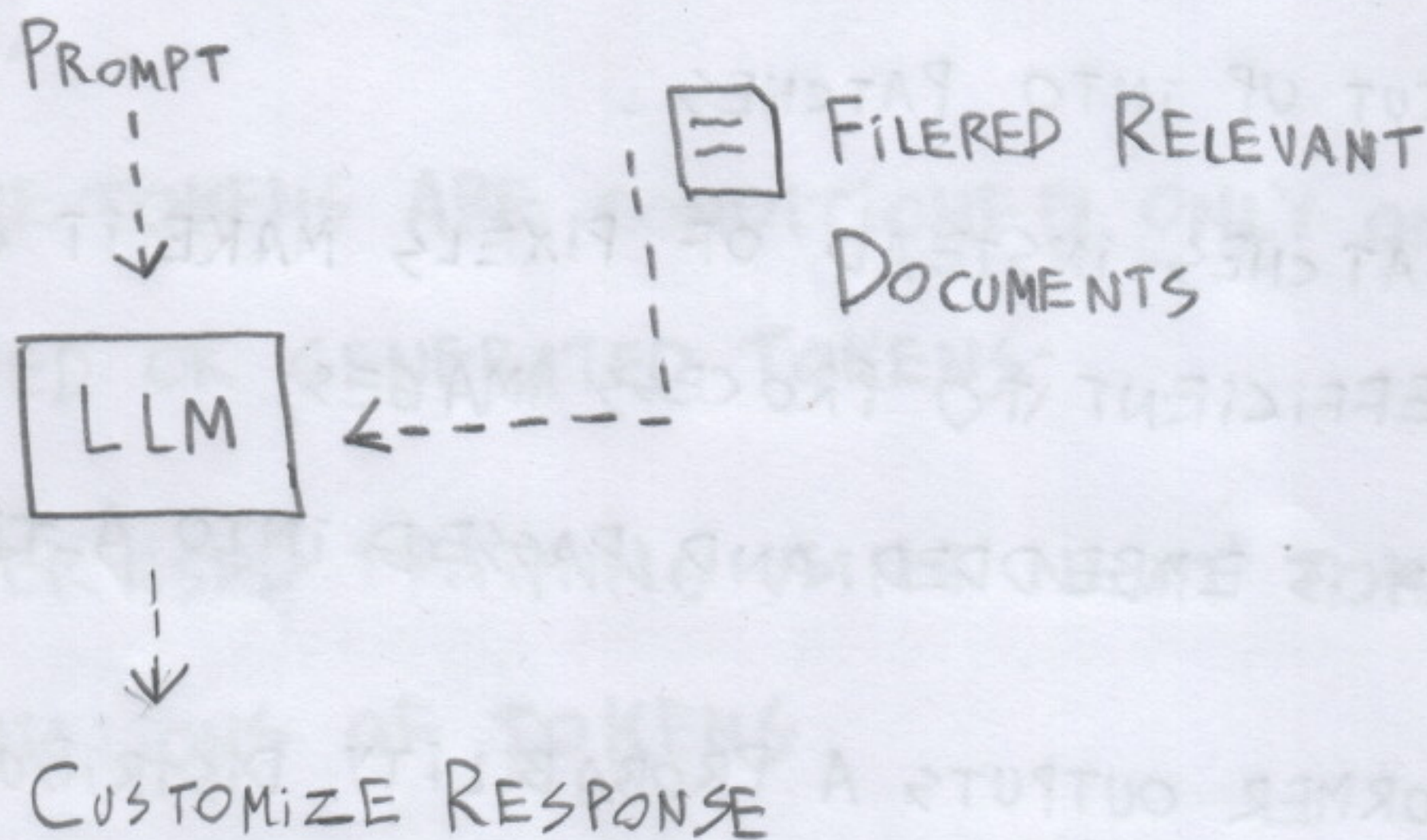  THE POSSIBLE CLASSES.

LANGUAGE MODEL $f_\phi$ → LANGUAGE RESPONSE $X_a$

→ $H_v$

PROJECTION W
VISION ENCODER

↑ $X_v$ IMAGE

$H_q$

$X_q$ LANGUAGE INSTRUCTION

IMAGE PATCHES

WHO DREW THIS PAINTING?

- The problem with LLMs

  "You don't know what you don't know"
      ↳ Socrates

  e.g : "Who at the family picnic is allergic to nuts?"


- Retrieval Augmented Generation (RAG)

  Prompt

  ⟍
   ↓                    ▤ Filered Relevant
                            Documents
  ┌─────┐
  │ LLM │  ← ─ ─ ─ ─ ─┘
  └─────┘
     ┆
     ↓

  Customize Response


- Vector Search for Fast Retrieval [With Weaviate (Vector DB)]

  ↳ Vector Search for Fast Retrieval

  Step 1

  ▤ ▤        ─ ─ ─ → Weaviate ─ ─ ─ ┐
  ▤ ▤                    ↑           ↓
                         ┆           ▤ Relevant
  Documents              ┆              Documents
                       Prompt
  Search over billions of
                         ┆
  documents in milliseconds
                         ↓
  ↳ Store and retrieve   LLM  ← ─ ─ ─┘
                         ┆
  relevant images, text, ┆
                         └→ Customized
  video;                     Response

8

# Applications of Multimodality in Industry

[I] Structured Data Generation

[II] Table Creation

[III] Understand Logic Flowcharts

- Multimodal Recommender System

Search is objective and recommendation is subjective.