# THE FUTURE OF SEARCH
## BEER & DATA 3

Lucas Farris <lucas@farris.com.br>

# WHAT'S GENAI EVEN GOOD FOR?

Over the past few years, I've gotten many questions from business owners regarding how they can use Generative Artificial Intelligence (GenAI) or Large Language Models (LLM) to improve their business. My answer for the past couple of years has been one: Retrieval Augmented Generation.

# CONTENTS

o  Measuring Good Search

o  Search Maturity

o  Four-Stage Recommender System

o  Retrieval Augmented Generation

o  Cloud Example

# MEASURING GOOD SEARCH

A fancy model doesn't is less important than a **good** evaluation framework

# MEASURING

| Accuracy | Relevance | Novelty | Serendipity | Diversity |
|---|---|---|---|---|
| Always remember the precision-recall trade-off. Measuring nDCG is always a good idea. | Human relevance (judgement) vs tracking relevance (presentation bias). | Users are likely to engage with new items if they are relevant | Non "obvious" recommendations make users happier | There should be enough variability between search results |

# WHAT MAKES SEARCH HARD

### Dynamic Inventory

o Maybe your products are only sold a limited number of times

o Maybe they sell-out fast (e.g., luxury cars, jewelry, houses)

o Is everything you learn about them lost?

### One-Time Users

o A common problem of dating sites, car manufacturers, and real-estate

o Some businesses have users that buy once come back only years later (or never)

### Scarce Data

o Maybe your products have a small amount of features

o Maybe your users are mostly anonymous

o Maybe your data is in a CRM software you can't access

o Maybe you're not measuring

# LANGUAGE MATURITY

## Level 1

Keyword Matching:
Inventory in indexed
(inverse index), filters and
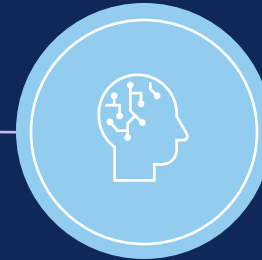sorting perform exact
matches against attributes

## Level 2

Taxonomies:
Items are catalogued
and grouped under
entities, cross-entity
categorizations are
built (ontologies),
synonym dictionaries
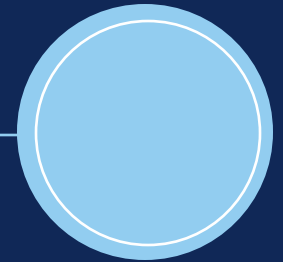are built in the search

## Level 3

Query Intent:
Semantics, query
rewrites (i.e., did you
mean...). System
knows if you're
looking for an article,
for a product, or
something else.

## Level 4

Knowledge Graphs:
Knowledge beyond
structured data (e.g.,
images, audio, videos)
connecting them to
text

# RANKING MATURITY

## Level 1

**Term-Frequency**
Relevance is measured based on how many times the search keywords appear in documents (i.e., a product's description). Commonly used algorithms are TF-IDF and BM25.

## Level 2

**Collaborative Filtering**
You track the rate between certain web events (i.e., page views) and rank items based on how they perform. Usual metrics for this are Click-Through Rate (CTR) or Conversion Rate (CR).

## Level 3

**Model-Based**
You use your inventory attributes and performance to train a regression based to predict how each item will perform. Usually non-negative matrix factorization (NNMF) and XGBoost are good choices here.
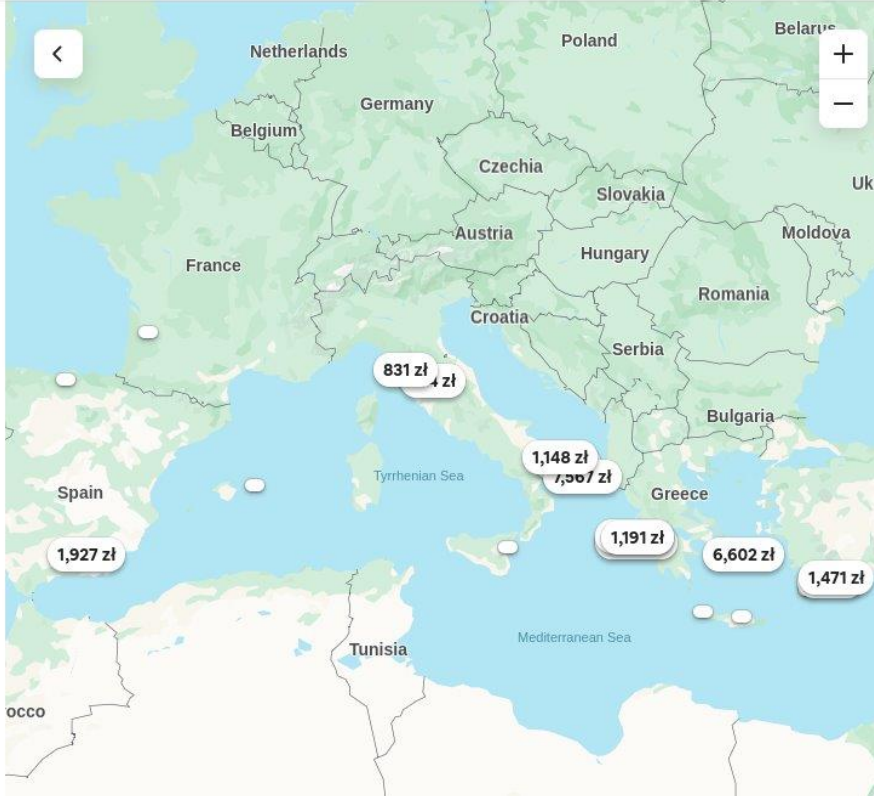
## Level 4

**Personalization**
The relevance model is trained with user features and unstructured data, and inference happens in real time.
Usually in this case we would need deep neural networks and GPUs.

see https://doi.org/10.1145/3351095.3372878 and https://ieeexplore.ieee.org/document/5197422

# SOLUTION ARCHITECTURE

The four-stage recommender system
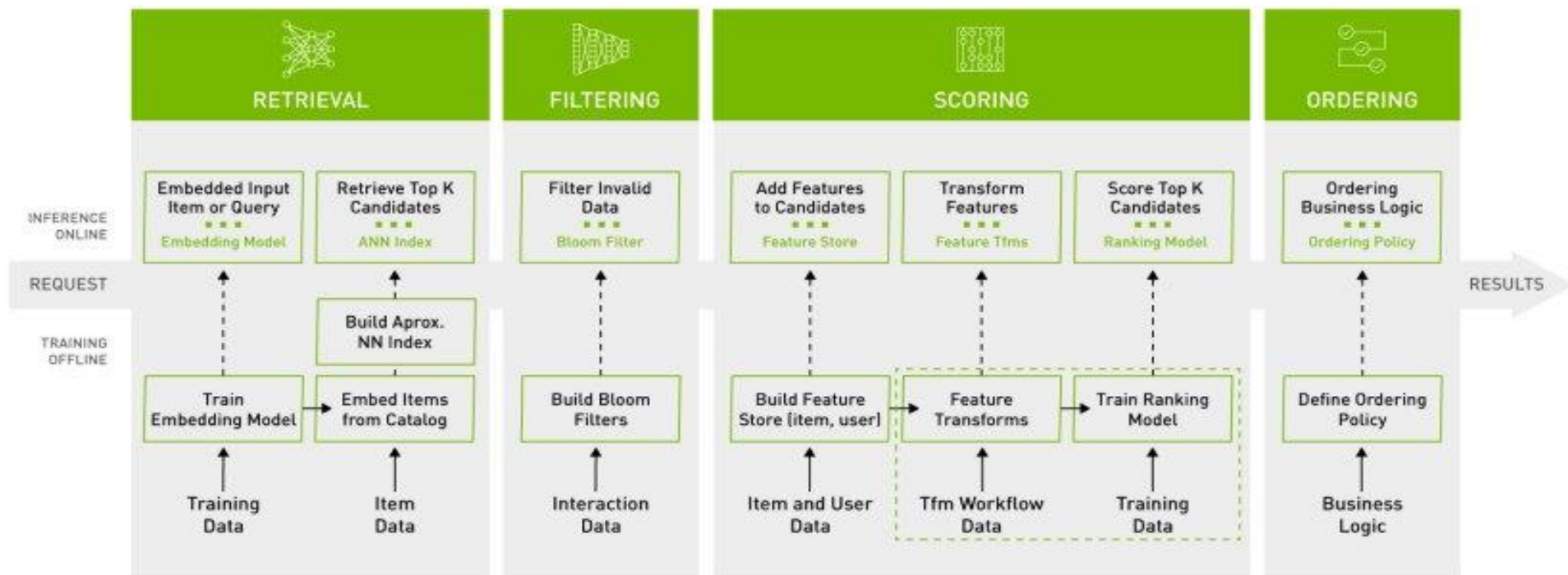
# WORD EMBEDDINGS

| Alternative Labels | Synonyms | Taxonomy | Ontology | Knowledge Graph |
|---|---|---|---|---|
| Words, acronyms, or expressions have the exact same meaning | Tokens with very similar meaning | Relationships between tokens and their categories or classes | Relationships between tokens, for instance in terms of how they interact | Mappings and relationships between entities and their related concepts |

# RETRIEVAL AUGMENTED GENERATION

Helping models help you

# DEMO TIME

Google Cloud Platform Example

# THANK YOU

Lucas Farris

lucas@farris.com.br