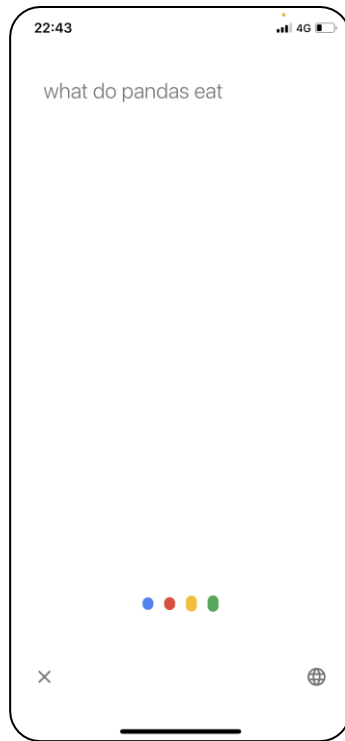# Information Retrieval
## Course presentation

## João Magalhães

# Information retrieval

# Voice based Search

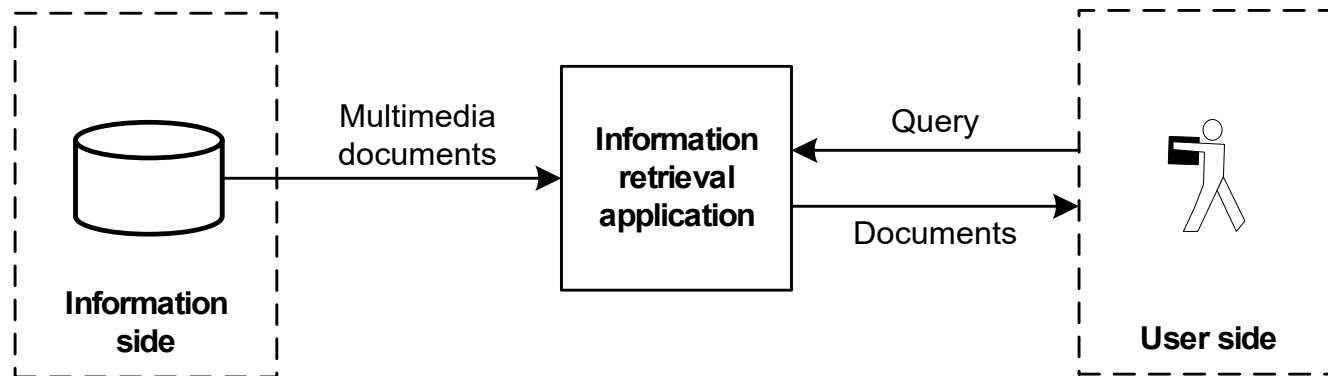# Natural language understanding

In fact, the Chinese `NORP` market has the three `CARDINAL` most influential names of the retail and tech space – Alibaba `GPE` , Baidu `ORG` , and Tencent `PERSON` (collectively touted as BAT `ORG` ), and is betting big in the global AI `GPE` in retail industry space . The three `CARDINAL` giants which are claimed to have a cut-throat competition with the U.S. `GPE` (in terms of resources and capital) are positioning themselves to become the 'future AI `PERSON` platforms'. The trio is also expanding in other Asian `NORP` countries and investing heavily in the U.S. `GPE` based AI `GPE` startups to leverage the power of AI `GPE` . Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one `CARDINAL` , with an anticipated CAGR `PERSON` of 45% `PERCENT` over 2018 - 2024 `DATE` .
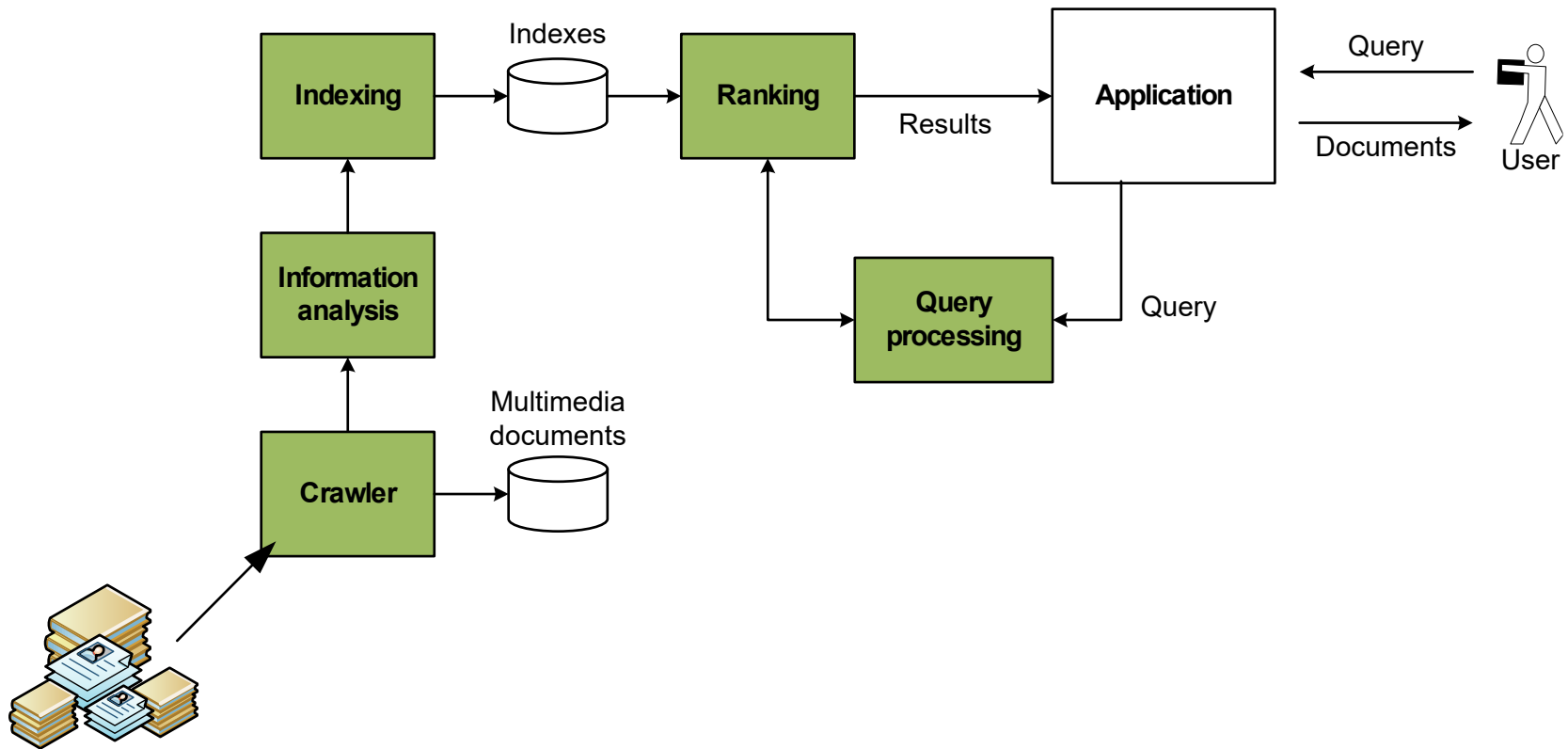
To further elaborate on the geographical trends, North America `LOC` has procured more than 50% `PERCENT` of the global share in 2017 `DATE` and has been leading the regional landscape of AI `GPE` in the retail market. The U.S. `GPE` has a significant credit in the regional trends with over 65% `PERCENT` of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as Google `ORG` , IBM `ORG` , and Microsoft `ORG` .

# Relevance vs similarity



What is the best algorithm to compute the relevance of documents for a given user information need?
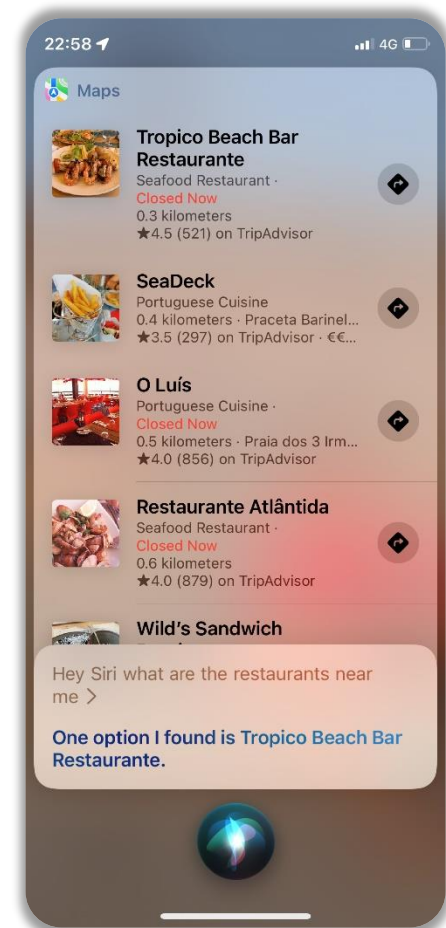
# Search architecture

# Search architecture components

- **Crawl** data for storage

- Analyse documents and compute **meaningful representations of natural language**

- Store data in an **efficient** manner

- Process **user information needs**

- **Find answer** to user request

# Mobile QA



- Hands-free devices favor a move to speech-based interaction.

- Voice-interaction favors dialogs.

- Voice interaction is nowadays a commodity.

# Named entities



Elon Musk
CEO of SpaceX

Elon Reeve Musk FRS is a technology entrepreneur, investor, and engineer. He holds South African, Canadian, and U.S. citizenship and is the founder, CEO, and lead designer of SpaceX; co-founder, CEO, ... Wikipedia

**Born:** June 28, 1971 (age 48 years), Pretoria, South Africa

**Net worth:** 19.9 billion USD (2019)

**Spouse:** Talulah Riley (m. 2013–2016), Talulah Riley (m. 2010–2012), Justine Musk (m. 2000–2008)

**Education:** University of Pennsylvania (1997), MORE
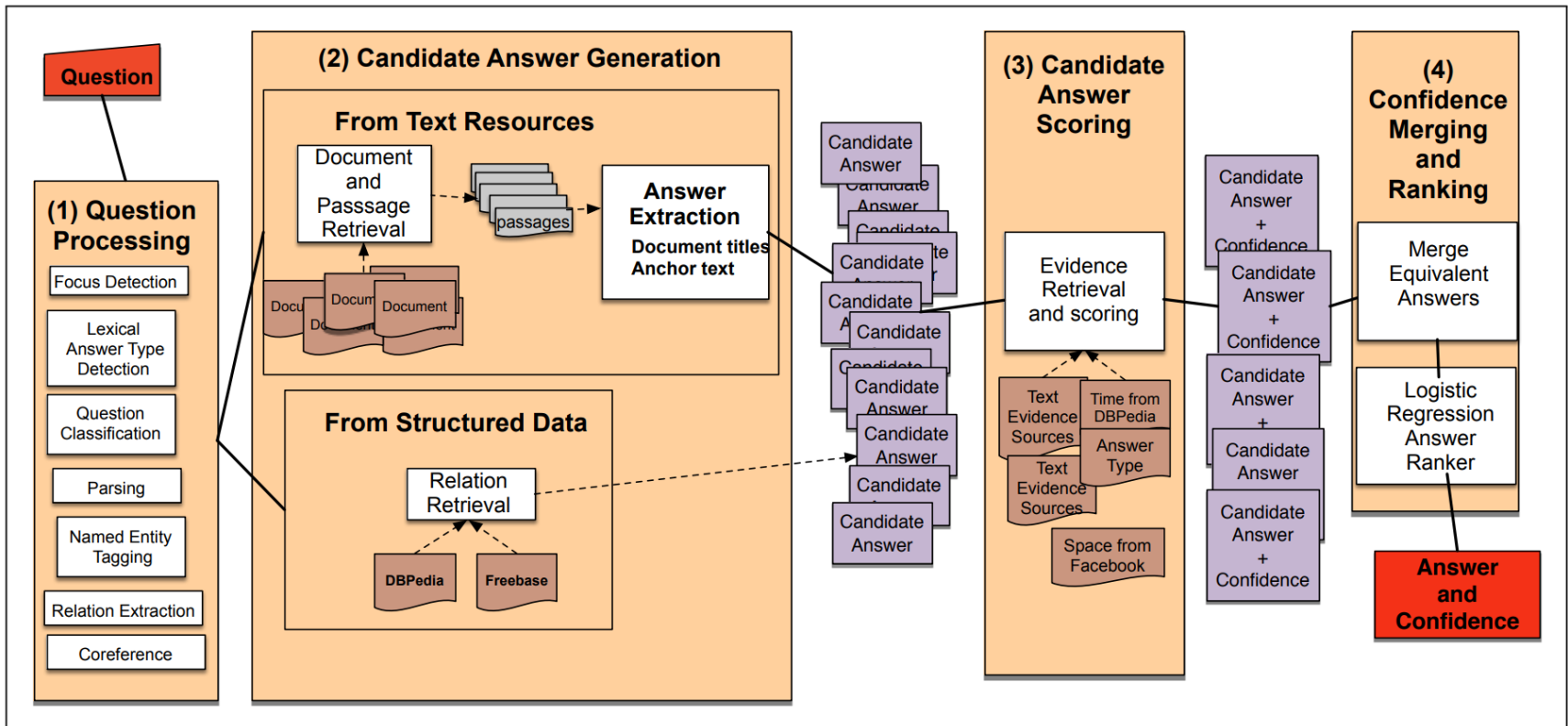


Rating ▾   Hours ▾

Tasca do Reguengos
4.5 ★★★★★ (564) · €€ · Portuguese
R. Gen. Humberto Delgado 13
Cosy · Casual · Good for kids

O Chafariz Palmeiros
4.1 ★★★★★ (20) · Restaurant
R. Chafariz Público 1
Closes soon · 4PM
Cosy · Casual · Good for kids

Nova Churrasqueira
4.1 ★★★★★ (176) · € · Restaurant
Azinhaga do Ginjal 14 B
Closed · Opens 7PM
Cosy · Casual · Good for kids

≔   More places

9

# Question answering architecture

# Question answering components

- **<u>Collect</u>** data for storage

- Analyse documents and compute **<u>meaningful representations of natural language</u>**

- Store data in an **<u>efficient</u>** manner

- Process **<u>user question</u>**

- **<u>Find</u> <u>candidate</u> <u>answers</u>**

- Extract correct answer

# Conversational Search



- Alexa, Siri, Google Assistant…

- CS methods need to track the evolution of the information need in the conversation;

- It needs to identify salient information needed for the current turn in the conversation;

- Retrieval methods are required to retrieve the relevant information from a knowledge base (e.g. Wikipedia).

U: Tell me about the **Neverending Story film**.
A: …

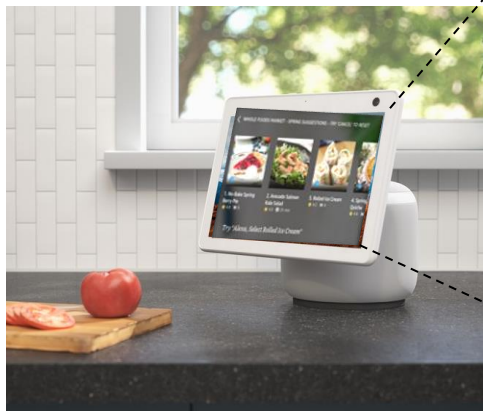U: What is **it** about?
A: …

U: Who was the author and when **it** was published?
A: …

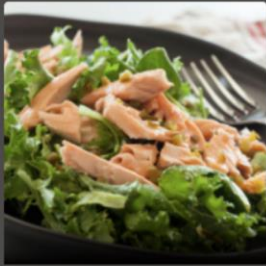U: Who are the **main characters**?
A: …

U: Did the horse **horse Artax** really die?
A: …

# Alexa TaskBot for Cooking

# Alexa TaskBot for DIY

https://www.youtube.com/watch?v=0nuCILb6VlI

# Schedule

| Natural Language Processing and Search | | | |
|---|---|---|---|
| **Week** | **#** | **Lecture** | **In-class labs** |
| 16-Sep-22 | 1 | Introduction | |
| 23-Sep-22 | 2 | Text processing, PoS, NGRAMS, cosine distar | Project phase 1 |
| 30-Sep-22 | 3 | Evaluation | |
| 7-Oct-22 | 4 | Language models | |
| 14-Oct-22 | 5 | Document categorization and ranking | |
| 21-Oct-22 | 6 | PoS and NE | Project phase 2 |
| 28-Oct-22 | 7 | Word embeddings | |
| 4-Nov-22 | 8 | Contextual embeddings | |
| 11-Nov-22 | 9 | Question answering | |
| 18-Nov-22 | 10 | Expalinable NLP | Project phase 3 |
| 25-Nov-22 | 11 | Project tips and feedback | |
| 2-Dec-22 | 12 | Computational Ethics for NLP | |
| 9-Dec-22 | | | |
| 16-Dec-22 | | | |

# References
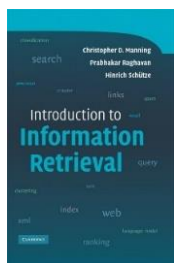
- Slides and articles provided during classes.

- Books:

Dan Jurafsky and James H. Martin, Speech and Language Processing (3rd ed. draft)

https://web.stanford.edu/~jurafsky/slp3/

C. D. Manning, P. Raghavan and H. Schütze, "Introduction to Information Retrieval", Cambridge University Press, 2008.

https://nlp.stanford.edu/IR-book/information-retrieval-book.html

# Lecturers

- Lectures: João Magalhães ([jmag@fct.unl.pt](mailto:jmag@fct.unl.pt))

- Labs: David Semedo ([df.semedo@fct.unl.pt](mailto:df.semedo@fct.unl.pt))

- When sending messaging lecturers always identify the course name and your group ID.

# Course grading

- The course has two mandatory components:
    - Laboratory (groups of 3 students):  60%  **(minimum grade > 10)**
        - Three phases, 20% per phase
    - Theoretical part (1 test or 1 exam):  40%  **(minimum grade > 8)**

- Theory test/exam:
    - 1 Test or 1 Exam (dates do be defined)

- Additional rules:
    - You may use one sided A4 sheet <u>handwritten by you</u> with your notes.
    - It must be handed in at the end of the test.

# Lab submission dates

- **Phase 1: Basic language models (20%)**      (13 October)
  - Basic notions of text processing and similarity
  - Language models
  - Evaluation

- **Phase 2: Info. extraction and search (20%)**      (10 November)
  - Entities extraction
  - Learning to rank documents

- **Phase 3: Neural language models (20%)**      (8 December)
  - Word embeddings and context embeddings
  - Self-attention

# Summary

- Context

- Objectives and plan

- Grading

- Labs