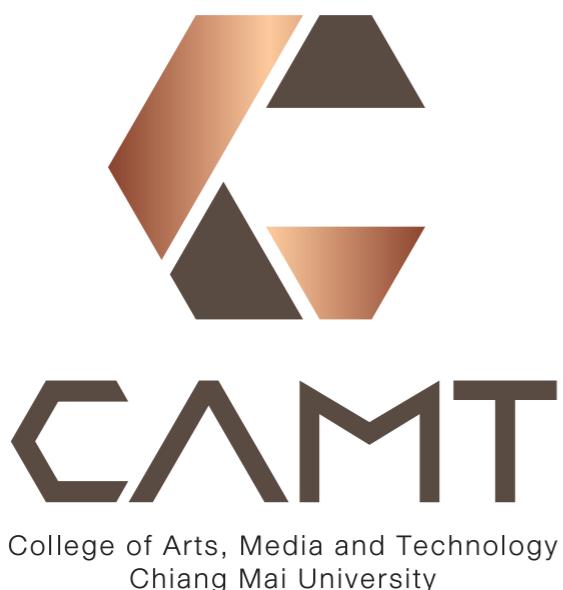


# SE 481 Introduction to Information Retrieval

## Module #0 — Homeroom



Passakorn Phannachitta, D.Eng.

[passakorn.p@cmu.ac.th](mailto:passakorn.p@cmu.ac.th)

College of Arts, Media and Technology  
Chiang Mai University, Chiangmai, Thailand

# Prerequisite

- SE 201 (953201) — Algorithms Design and Analysis

# Know your lecturer

Passakorn Phannachitta, D.Eng. (Aj. Kong)

Office: CAMT 417

Email: [passakorn.p@cmu.ac.th](mailto:passakorn.p@cmu.ac.th)

# Schedule

- Lectures
  - Mondays and Thursdays 13:00 - 14:30

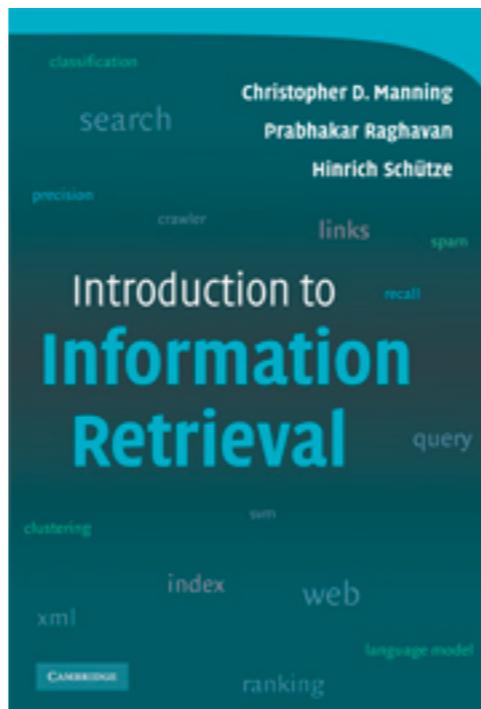
# What will be in the course ?

- Fundamental of IR
- Indexing
- Vector-space model
- Evaluation in IR
- Web search (Search engine)
- Basic machine-learning in IR
- Showcases and applications

# Grading - based on Adjusted criteria

- Hands on 20%
- Many assignments 20%
- Project(s) 30%
- Midterm & final examination 30%

# Material



<https://nlp.stanford.edu/IR-book/>

# Class communication

- MS TEAM
- SCOTT

# Class Policy

- No late assignment submission.

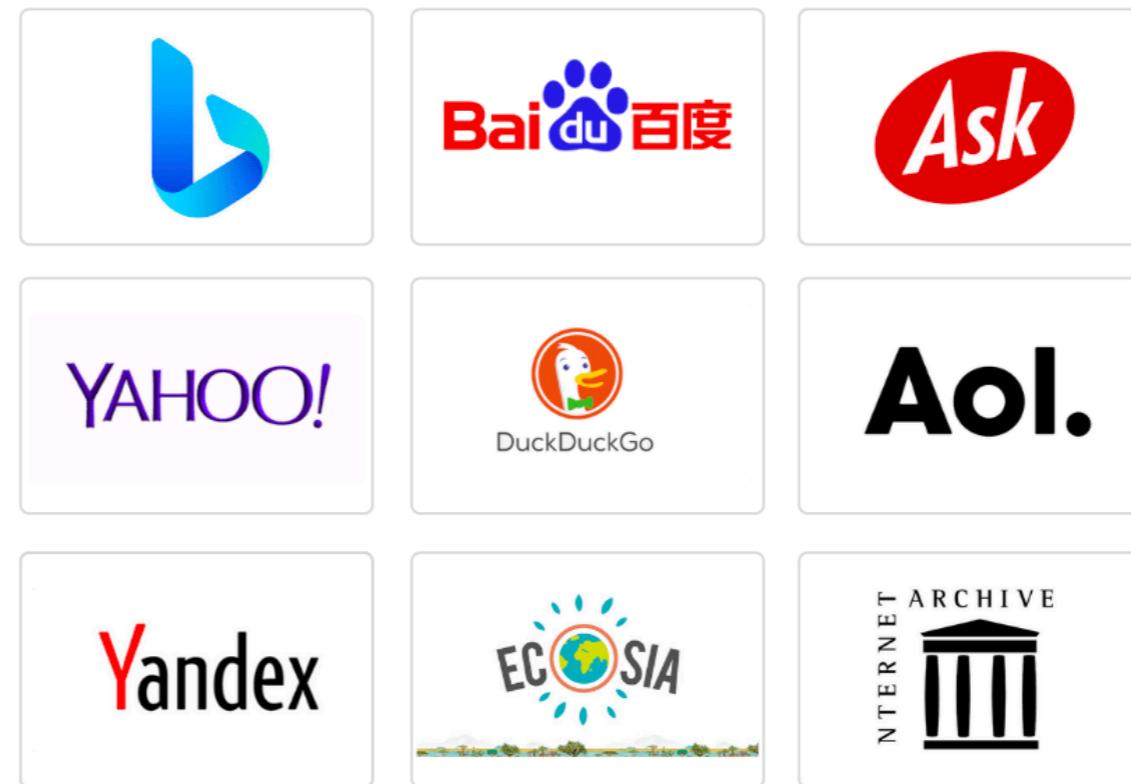
# Information Retrieval (IR)

- **Information retrieval** is the science of searching for information in a document, searching for documents themselves, and also searching for the **metadata** that describes data, and for databases of texts, images or sounds.

— Wikipedia

# Information Retrieval (IR)

- Mostly we think of IR as a web search



Ref: <https://www.reliablesoft.net/wp-content/uploads/2016/12/top-search-engines-oct-2020.png>

# Information Retrieval (IR)

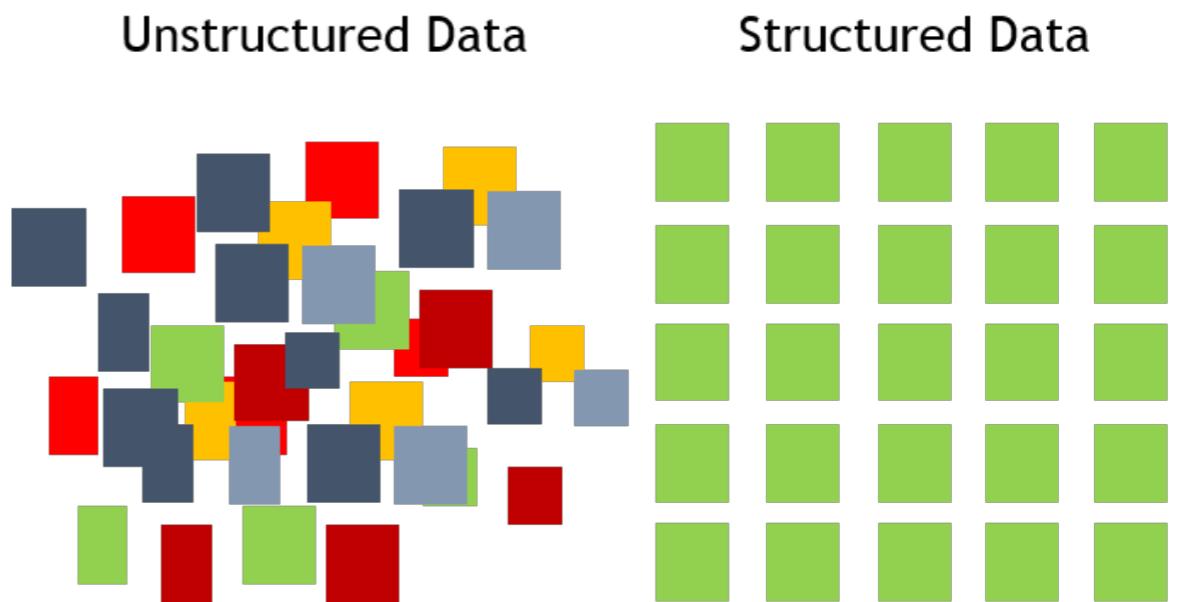
- But there are many other cases
  - E-mail search
  - Searching for particular files in our computer
  - etc.

# Basic assumption of IR

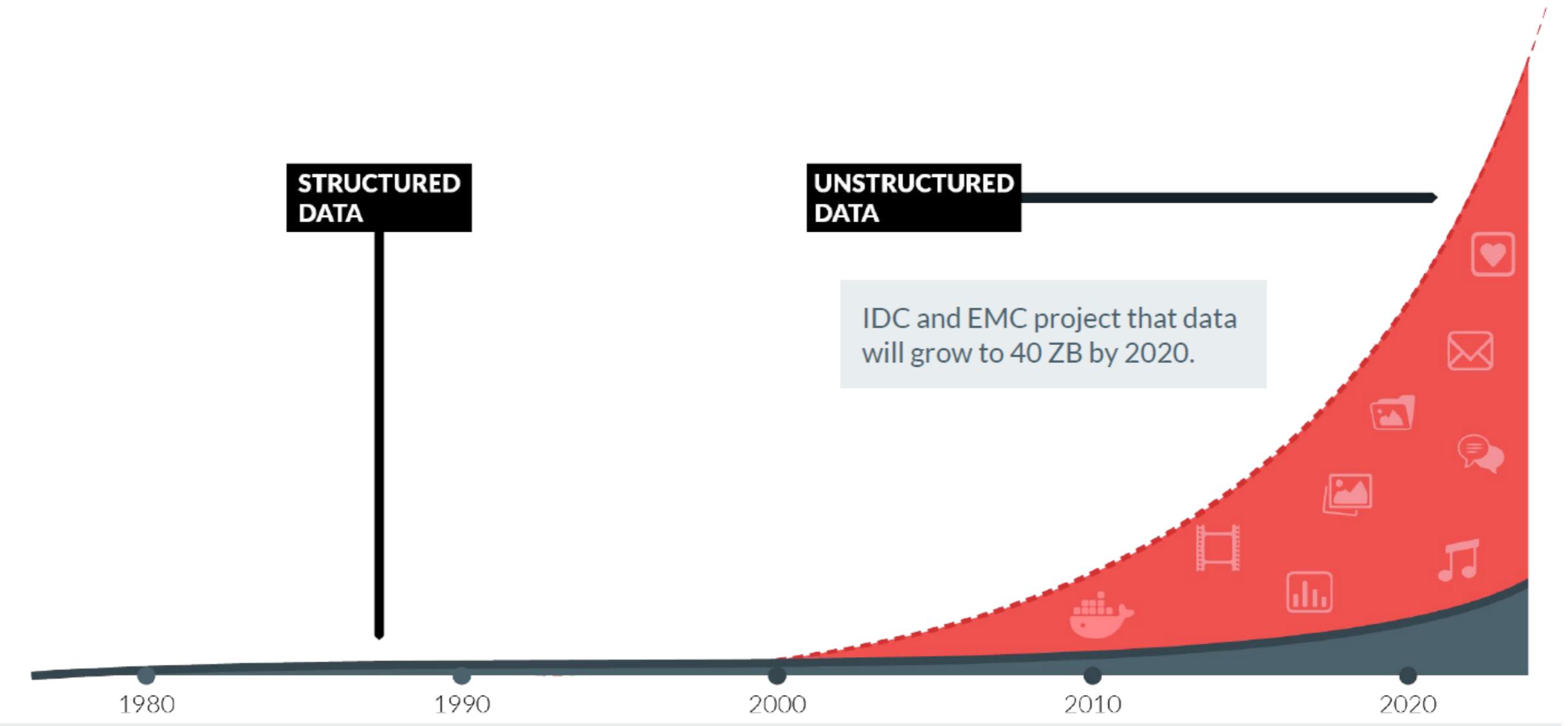
- Information is stored in a set of documents
- An IR system retrieves documents with information that is **most relevant to the user's information need** and helps the user complete a **task**
- Data are born **unstructured**

# Unstructured data

- Structured data
  - Data with a clearly defined type whose pattern makes them easily searchable, linkable, and comparable.
- Unstructured data
  - Data which is not structured via any pre-defined data models or schema



# Unstructured data is everywhere



# Information Retrieval (IR)

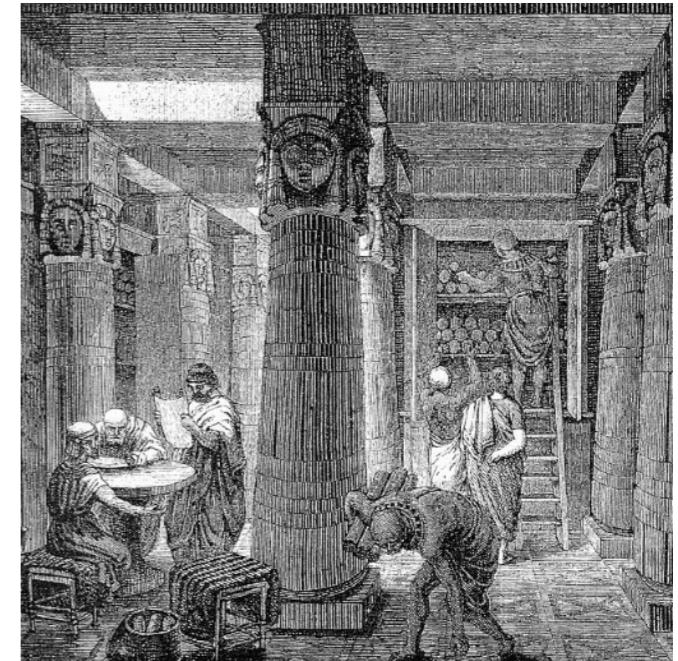
- Google's mission statement

Our mission is to  
organise the world's  
information and make it  
universally accessible  
and useful.

Ref: <https://about.google/>

# Organizing a large amount of information

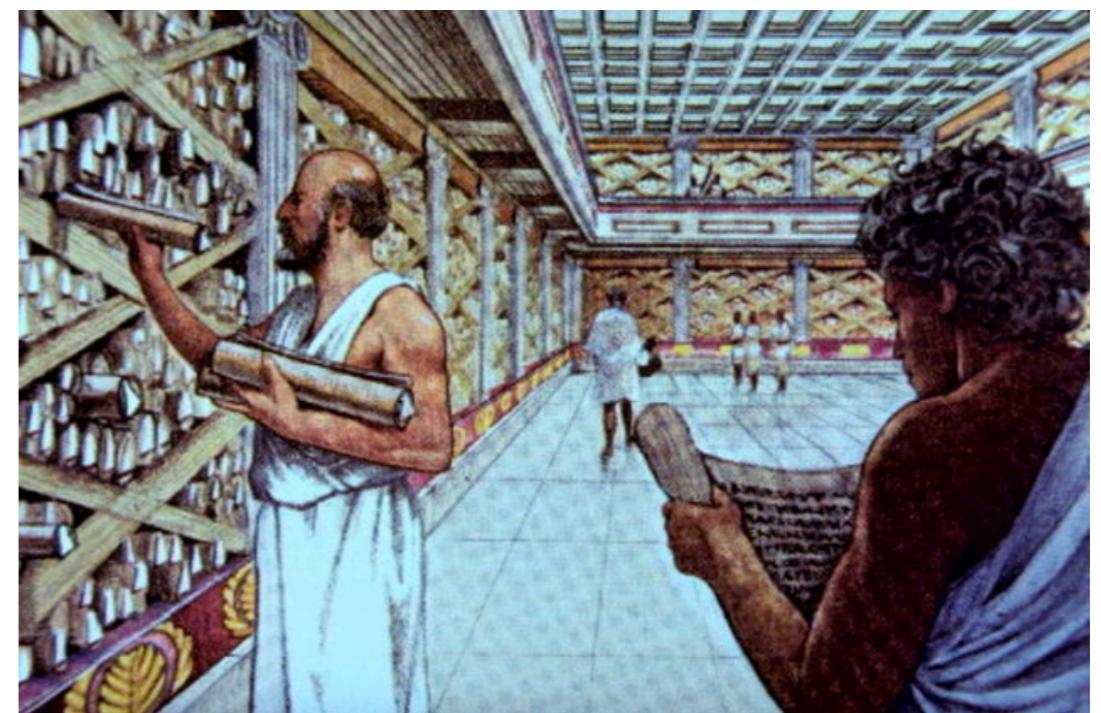
- The very first form might be (physical) libraries
- E.g., the Library of Alexandria — the largest libraries in ancient world
  - Built ~300 BC
  - Collected numerous papyrus scrolls



Ref: [https://en.wikipedia.org/wiki/Library\\_of\\_Alexandria](https://en.wikipedia.org/wiki/Library_of_Alexandria)

# Searching for a piece of information

- Suppose there are thousands of scrolls, how can we find the most relevant one of our interest information?



Ref: <https://thesomathread.com/2016/09/08/the-role-and-fate-of-the-library-at-alexandria/>

# A very first attempt

- **Callimachus**, a scholar at the Library of Alexandria made the tool called **Pinakes**.
- **Pinakes** were considered as the earliest library catalog.
- Works are divided in genres and categories, e.g., law, poetry, history, medicine, and mathematics.
- The **Pinakes** proved indispensable to librarians for centuries, and they became a model for organizing knowledge throughout the Mediterranean



Ref: <https://www.geni.com/people/Callimachus/6000000078663146846>

# A librarian's standard

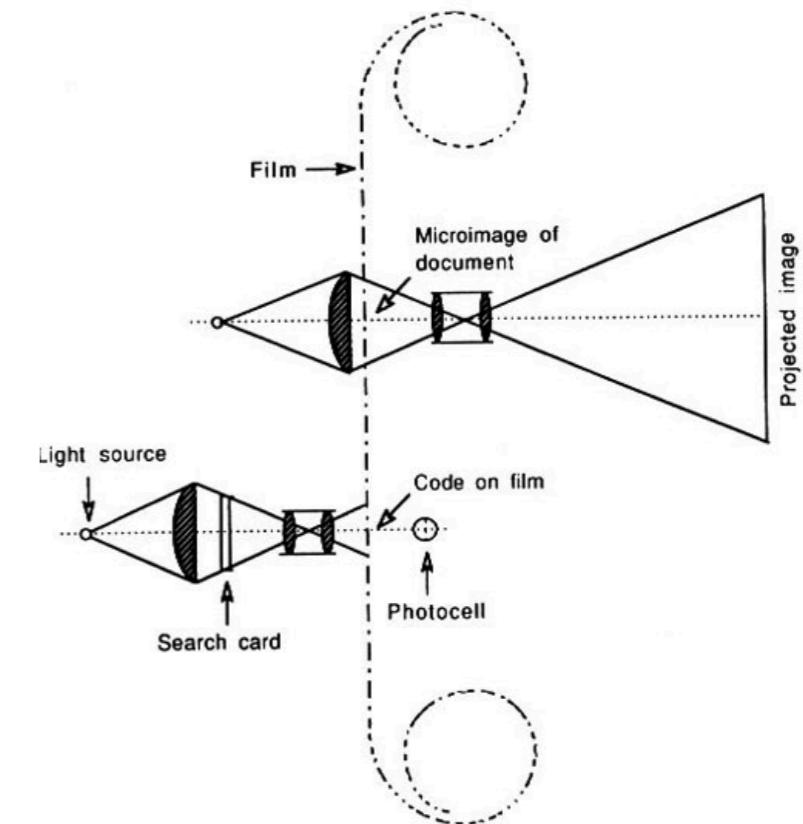
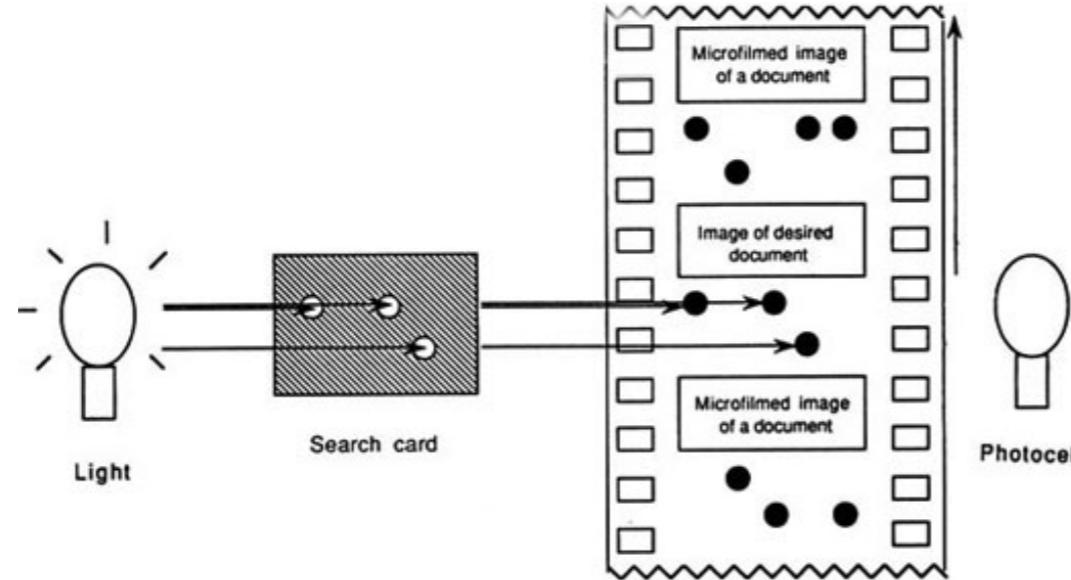
- Local variations for cataloging and library classification continued throughout the late 1800s,
  - when Anthony Panizzi and Melvil Dewey proposed a more standard approach called **Decimal Classification**.
- 
- E.g., 500 — Natural sciences and mathematics
    - 510 — Mathematics
      - 516 — Geometry
        - 516.3 — Analytic geometries

# What we have learned

- Categorization can help us narrowing the search scope.
  - Precisely locate the information in content-level is still far from reality.
- 
- Categorization is made bottom up.
  - What about doing the bottom up from the content level?

# The very first IR machines

- Emanuel Goldberg's Statistical machine (1931)



Ref: <https://history-computer.com/emanuel-goldberg/>

# Goldberg's machine

- Labelling and indexing

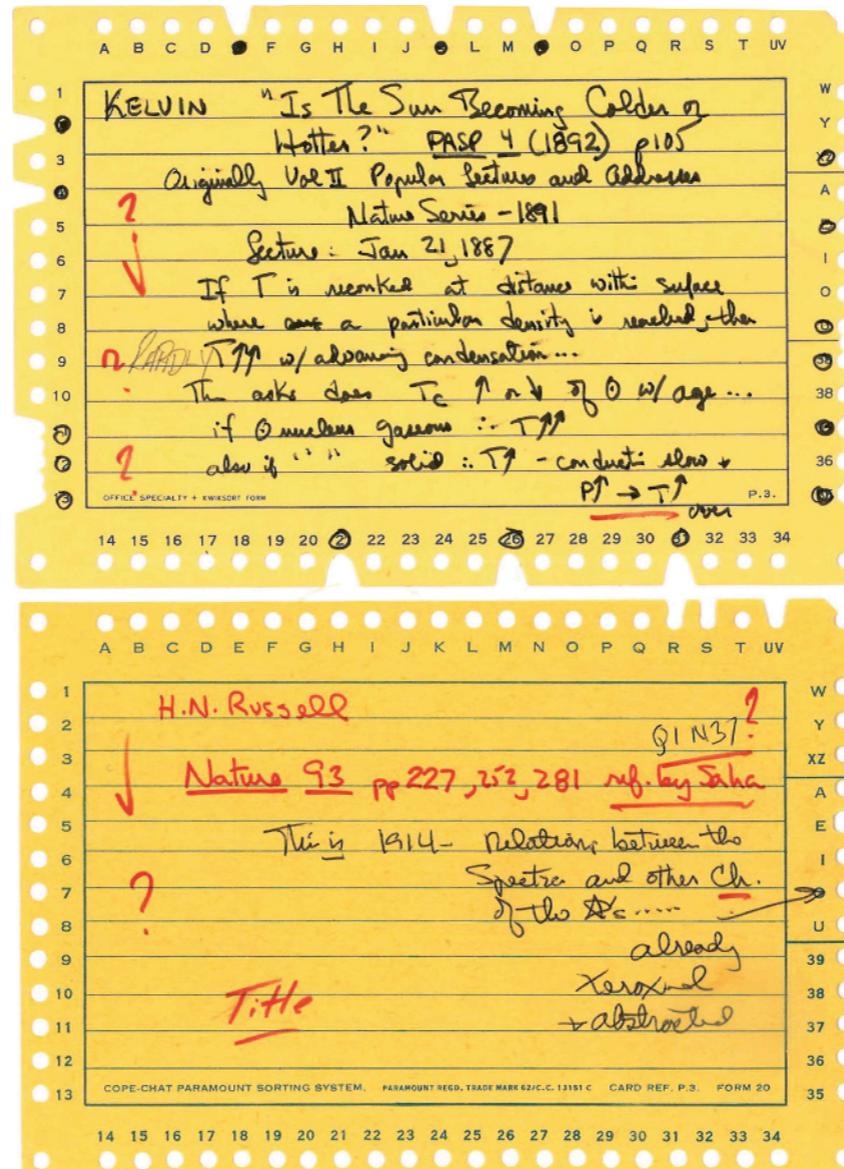
# Coining the word Information Retrieval

- Calvin Mooers
- At MIT, Mooers developed a mechanical system using superimposed codes of descriptors for information retrieval called **Zatocoding**.
- Mooers founded the Zator Company in 1947 to market this idea, and pursued work in information theory, information retrieval, and artificial intelligence.
- Mooers coined the term "information retrieval" using it first in a conference paper presented in 1950.



Ceruzzi P.E. (2019) Calvin Mooers, Zatocoding, and Early Research on Information Retrieval. In: Haigh T. (eds) Exploring the Early Digital. History of Computing. Springer

# Zatocoding



Ceruzzi P.E. (2019) Calvin Mooers, Zatocoding, and Early Research on Information Retrieval. In: Haigh T. (eds) Exploring the Early Digital History of Computing. Springer

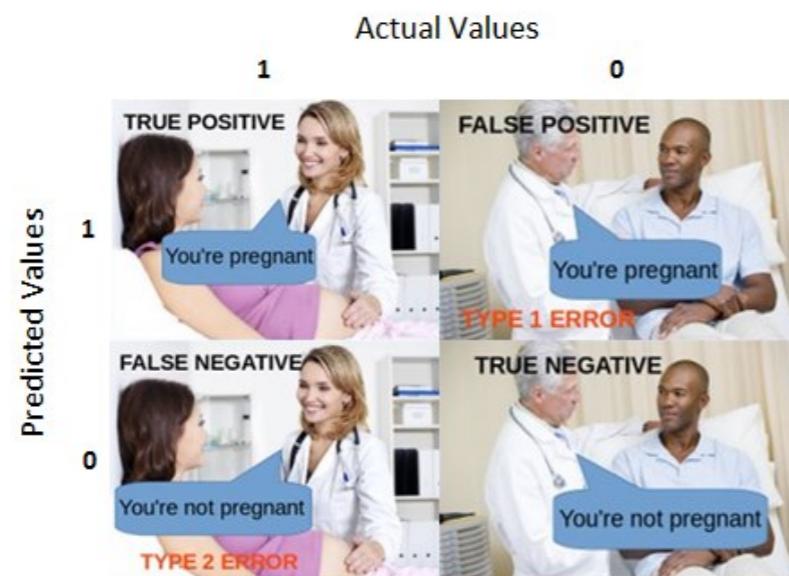
# Uniterms

- Mortimer Taube (1952)
- 100 most important leaders in Library and Information Science of the 20th century
- Taube invented **Coordinate Indexing**, which uses **uniterms** in the context of cataloging
- Index card
- E.g., a document on barefoot running sneakers might be filed under "barefoot" but perhaps not "sneakers" which would be found on too many documents.

Ref: <https://en.wikipedia.org/wiki/Uniterm>

# Evaluation in IR

- Cyril Cleverdon (1960s)
- Precision and Recall



|                 |          | Real Label          |                     |
|-----------------|----------|---------------------|---------------------|
|                 |          | Positive            | Negative            |
| Predicted Label | Positive | True Positive (TP)  | False Positive (FP) |
|                 | Negative | False Negative (FN) | True Negative (TN)  |

Precision =  $\frac{\sum TP}{\sum TP + FP}$

Recall =  $\frac{\sum TP}{\sum TP + FN}$

Accuracy =  $\frac{\sum TP + TN}{\sum TP + FP + FN + TN}$

Ref: <https://www.bualabs.com/archives/1968/what-is-confusion-matrix-what-is-metrics-accuracy-precision-recall-f1-score-difference-metrics-ep-1/>

# IR and Ranking

- Hans Peter Luhn (1957) —> term frequencies (tf)
- Karen Sparck-Jones (1972) —> inverse document frequency (idf)
- Gerard Salton (1975) —> tf x idf

**TF-IDF**

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

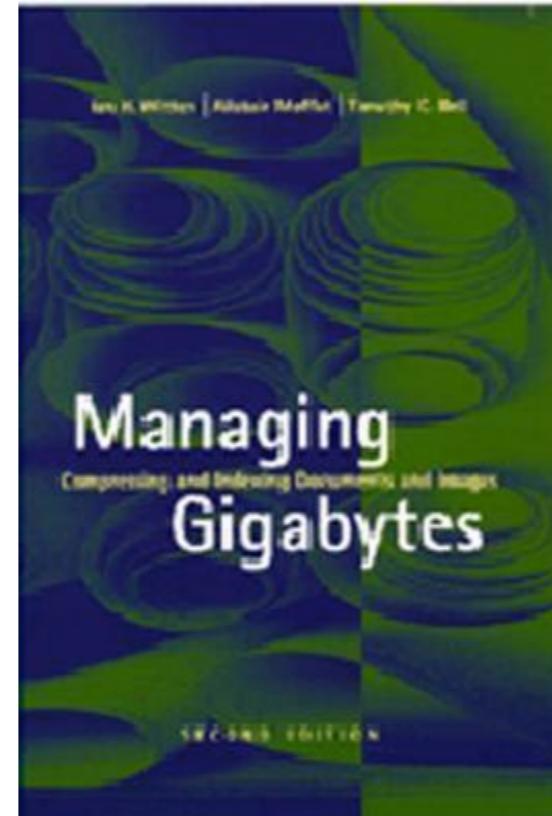
Term frequency  
Number of times term  $t$  appears in a doc,  $d$

Inverse document frequency  
 $\log \frac{1 + n}{1 + df(d, t)} + 1$   
 $n$  # of documents  
 $df(d, t)$  Document frequency of the term  $t$

Ref: <https://towardsdatascience.com/tf-term-frequency-idf-inverse-document-frequency-from-scratch-in-python-6c2b61b78558>

# IR and Compression

- Ian Witten, Alistair Moffat, and Timothy Bell (1994)

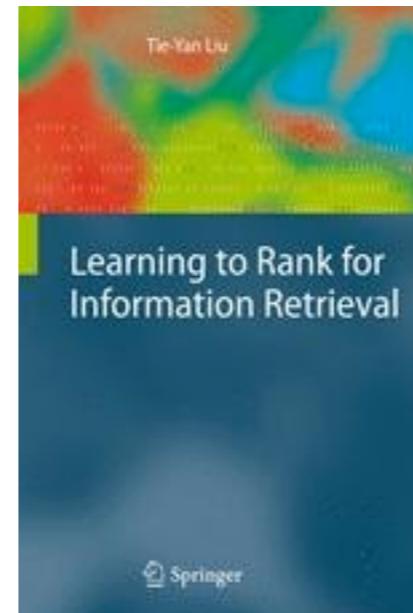


Ref: Witten, I. H., Witten, I. H., Moffat, A., Bell, T. C., Bell, T. C., & Bell, T. C. (1994).

Managing gigabytes: compressing and indexing documents and images. Morgan Kaufmann.

# Ranking with Models

- Stephen Robertson (1994) —> BM25
- Bruce Croft (1998) —> Language model
- Sergey Brin and Larry Page (1998) —> PageRank
- Most recent technology —> Learning to Rank



Ref: Liu, T. Y. (2011). Learning to rank for information retrieval.

# Assignment 0

- Please install Anaconda in your local machine  
<https://www.anaconda.com/>

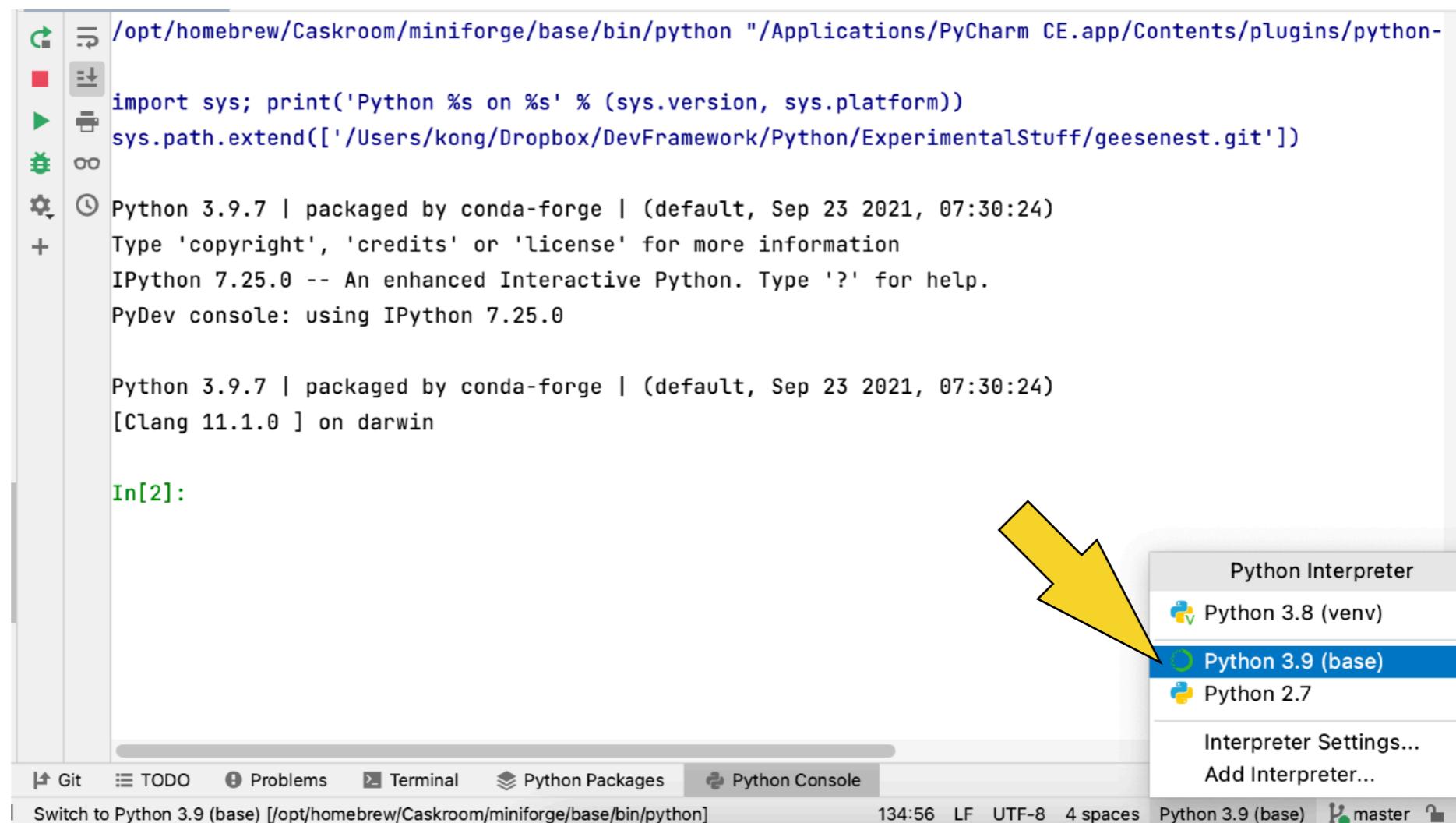


- And an IDE. (The lecturer prefers PyCharm)  
<https://www.jetbrains.com/pycharm/>
- Show that your IDE is connected to Anaconda

# Assignment 0

- Show that your IDE is connected with Anaconda,

e.g.,



The screenshot shows the PyCharm IDE interface. On the left, there's a file browser with several icons. In the center, a terminal window displays Python startup information:

```
/opt/homebrew/Caskroom/miniforge/base/bin/python "/Applications/PyCharm CE.app/Contents/plugins/python-
import sys; print('Python %s on %s' % (sys.version, sys.platform))
sys.path.extend(['/Users/kong/Dropbox/DevFramework/Python/ExperimentalStuff/geesenest.git'])

Python 3.9.7 | packaged by conda-forge | (default, Sep 23 2021, 07:30:24)
Type 'copyright', 'credits' or 'license' for more information
IPython 7.25.0 -- An enhanced Interactive Python. Type '?' for help.

PyDev console: using IPython 7.25.0

Python 3.9.7 | packaged by conda-forge | (default, Sep 23 2021, 07:30:24)
[Clang 11.1.0 ] on darwin

In[2]:
```

At the bottom, the status bar shows: "Switch to Python 3.9 (base) [/opt/homebrew/Caskroom/miniforge/base/bin/python]". To the right, a floating "Python Interpreter" menu is open, listing three options: "Python 3.8 (venv)", "Python 3.9 (base)" (which is highlighted in blue), and "Python 2.7". A large yellow arrow points from the text above to this menu.

# Assignment 0

- Show that you are able to configure  
[https://github.com/alexdedyura/cpu-benchmark?  
ref=pythonrepo.com](https://github.com/alexdedyura/cpu-benchmark?ref=pythonrepo.com),  
and run it in your IDE connected with Anaconda,  
e.g.,

Python CPU Benchmark by Alex Dedyura (Windows, macOS(Darwin), Linux)

CPU: Apple M1

Arch: arm64

OS: Darwin

Benchmarking:

Time: 23.051s

Time: 22.292s

# Time for questions