

Search and Recommendation

Hao Sheng
August 10th, 2023



Agenda

- **Recap**
- **Lecturing:** Search Engine System
- **Break-out:** Measure the Success of the Recommendation System
- **Lecturing:** Advanced Topics (Part I)
- **Lab:** Recommendation system notebook II
- **Lecturing:** Advanced Topics (Part II)

60% Lecturing

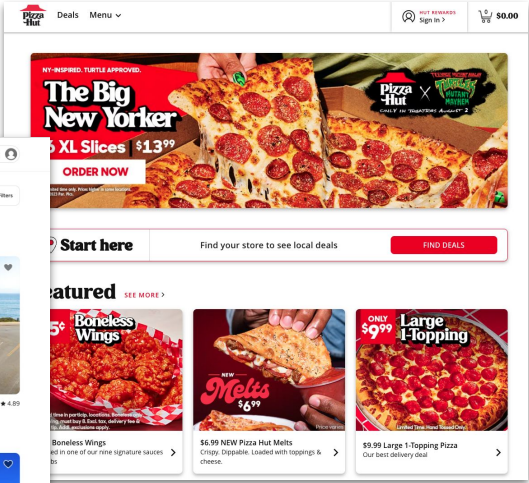
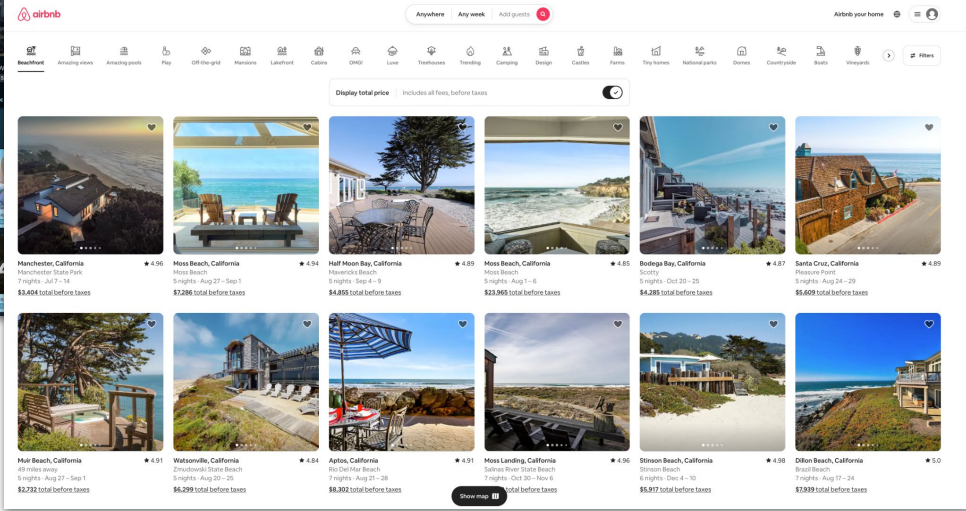
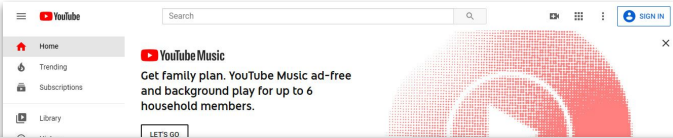
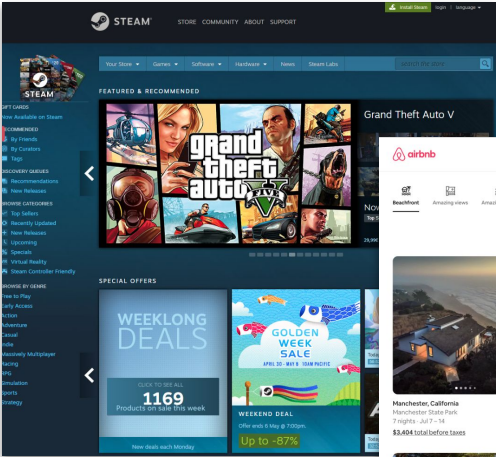
25% Lab

10% Discussion



Recap of Day 1

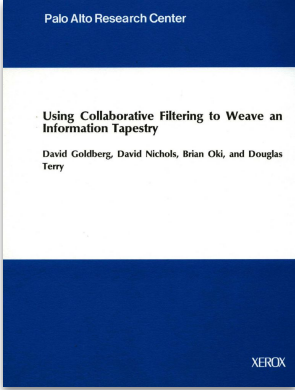
Homepages @ 2023



Everyone gives user recommendations on the first impression!

Earlier Recommendation Systems

- Recommendation system: RS automatically selects personalized information based on users' preferences.
- Grundy:
 - Ask user questions and assign stereotype.
 - **Content-based filtering.**
- Tapestry:
 - Find similar users and recommend their choices.
 - **Collaborative filtering.**



Quick Recap: Collaborative Filtering



Sara	5	3		2	2	2
Jesper	4	3	4		3	3
Therese	5	2	5	2	1	1
Helle	3	5	3		1	1
Pietro	3	3	3	2	4	5
Ekaterina	2	3	2	3	5	5



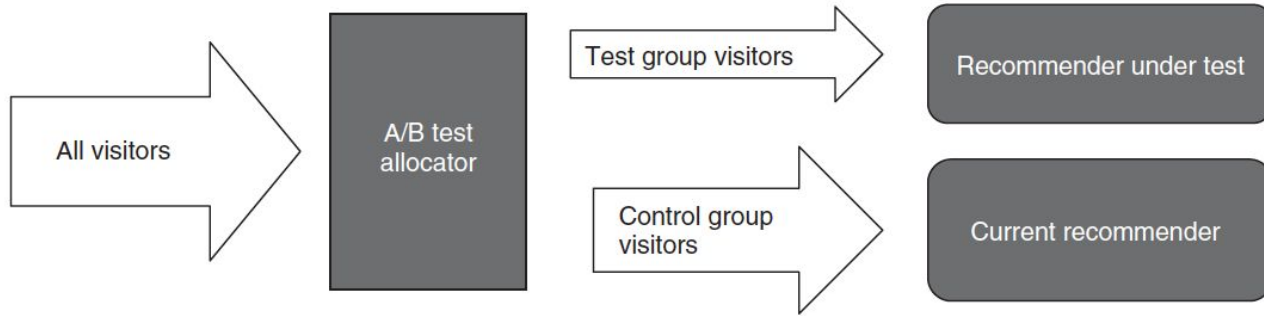


Figure 9.16 In an A/B test, visitors are split into two groups: the test group that sees the new feature and a control group that continues as usual.

Search Engine System

The Google logo is displayed in its characteristic multi-colored font (blue, red, yellow, blue, green, red).

who would win in a fight between|

- who would win in a fight between **batman and superman**
- who would win in a fight between **a taco and a grilled cheese**
- who would win in a fight between **hulk and wolverine**
- who would win in a fight between **batman and iron man**

Google Search

I'm Feeling Lucky

Search Engine: First Glance

Google X

Tuition News Acceptance rate Reddit Curriculum Images Requirements Re

About 311,000 results (0.36 seconds)

Stanford Institute for Computational & Mathematical Engineering
<https://icme.stanford.edu> :

ICME, Stanford - Stanford University
 ICME faculty and students conduct groundbreaking research, provide consulting, and teach courses in computational mathematics and scientific computing. Degree ...

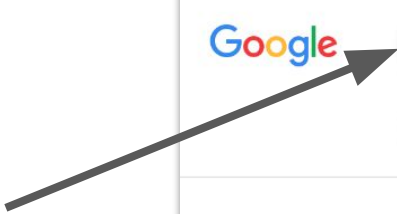
Master of Science
 Through robust coursework in computational mathematics and ...

Doctor of Philosophy
 ICME PhD students cultivate a broad and deep understanding ...

Academics & Admission
 interdisciplinary institute at the intersection of mathematics ...

People
 Christiane Adcock. Ph. · Izabel Pirimai Aguiar. Ph. · Ayya Alieva ...

[More results from stanford.edu »](#)



Search Query

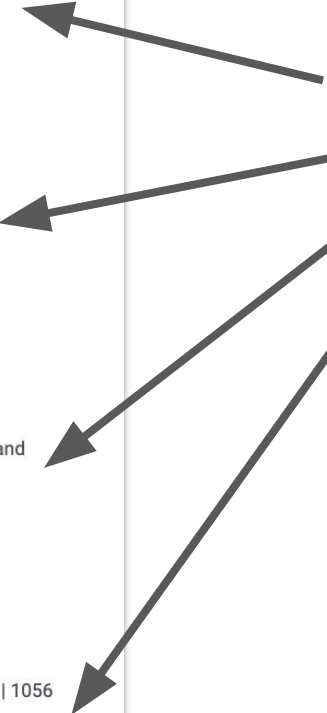
Search Engine: First Glance

S Stanford University
https://events.stanford.edu › department › institute_for... ⋮
Institute for Computational and Mathematical Engineering ...
The Institute for Computational & Mathematical Engineering (ICME) is a degree granting (M.S./Ph.D.) interdisciplinary institute at the intersection of ...

S Stanford University Bulletin
https://bulletin.stanford.edu › programs › CME-MS ⋮
CME-MS Program - Stanford Bulletin
ICME is a degree granting (M.S./Ph.D.) interdisciplinary institute at the intersection of mathematics, computing, engineering and applied sciences.
Courseor course: Intermediate Econometrics II (3 ...

S Stanford Bulletin Archive
https://archived-bulletin.stanford.mobi › instituteforco... ⋮
Institute for Computational and Mathematical Engineering
At ICME, we design state-of-the-art mathematical and computational models, methods, and algorithms for engineering and science applications. The program ...
Or MS&E 327: Topics in Causal Inference STATS 263: Design of Experiments

in LinkedIn
https://www.linkedin.com › company › icme-stanford ⋮
Institute for Computational and Mathematical Engineering ...
Institute for Computational and Mathematical Engineering at **Stanford University (ICME)** | 1056 followers on LinkedIn. Groundbreaking research into complex ...



Search Engine Results
(Websites)

Search Engine: First Glance

Stanford Online
<https://online.stanford.edu> > programs > computaciona...

Computational and Mathematical Engineering MS Degree

The Institute for Computational and Mathematical Engineering (ICME) is a degree granting institute at the intersection of mathematics, computing, engineering ...

<https://twitter.com/ICMEStanford>

Stanford ICME (@ICMEStanford) · Twitter

<p>Generative Models (SWS 14) Aug. 9-10 2-4pm PDT</p>	<p>Search and Recommendation (SWS 13) Aug. 9-10 8-11am PDT</p>	<p>Introduction to Natural Language Processing (SWS 12) Aug. 7-8 3-4pm PDT</p>
<p>Stanford ICME's new Generative Models workshop—from 8/9 to 8/10 & taught by instructor Aashwin Mishra—focuses on models that generate new data instances, such as images, text, or audio.</p> <p>Register: www.eventbrite.com/e/ic...</p> <p>Twitter · Jul 28, 2023</p>	<p>Explore methods for enhancing #UX & retrieving information in @Stanford ICME's new Search and Recommendation workshop from 8/9 to 8/10, led by @Apple Staff Machine Learning Engineer @hao_ssr.</p> <p>Register: www.eventbrite.com/e/ic...</p> <p>Twitter · Jul 27, 2023</p>	<p>Learn the building blocks of modern #NLP concepts in @Stanford ICME's Intro course from 8/7 to 8/8, taught by brothers and @Google Software Engineers @afshinea & @shervinea.</p> <p>Register: www.eventbrite.com/e/ic...</p> <p>Twitter · Jul 26, 2023</p>

Search Engine Results (Tweets)



Stanford University
<https://events.stanford.edu> > department > institute_for...

Institute for Computational and Mathematical Engineering ...

The Institute for Computational & Mathematical Engineering (ICME) is a degree granting (M.S./Ph.D.) interdisciplinary institute at the intersection of ...

Search Engine: First Glance

People also ask :

- How competitive is Stanford graduate school? ▾
- What is the Toefl code for Stanford ICME? ▾
- Is Stanford University good for engineering? ▾
- What is the mathematical and computational finance program at Stanford University? ▾

Feedback



Search Engine Results
(Tweets)



Stanford Online

<https://online.stanford.edu> > programs > computaciona... ▾

Computational and Mathematical Engineering MS Degree

The Institute for Computational and Mathematical Engineering (ICME) is a degree granting institute at the intersection of mathematics, computing, engineering ...

Search Engine v.s. Movie Recommendation

- It has a search bar!
- The items are mal-defined at the first glance.
- User does not simply rate the search results!



“Recommender systems (RSs) are software tools and techniques that provide suggestions for items that are most likely of interest to a particular user. ”

--- *Introduction to Recommender Systems Handbook*

“Recommender systems (RSs) are software tools and techniques that provide suggestions for items that are most likely of interest to a particular user. ”

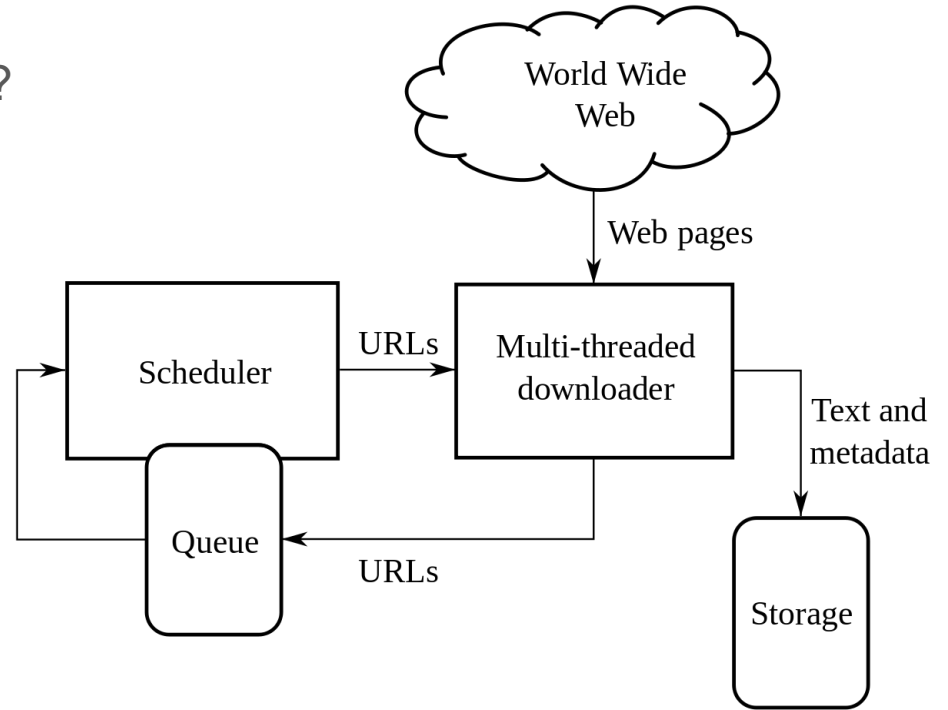
--- *Introduction to Recommender Systems Handbook*

Search Engine is a **Recommendation System!**

- It has a search bar!
 - -> How to incorporate the user intention?
- The items are mal-defined at the first glance.
 - -> How to crawl the internet and store the items?
- User does not simply rate the search results!
 - -> How to assign user-item rating with user data?

Where to get the search results (items)?

The internet!



Where to get the search results (items)?

The internet!

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```

parsed_urls = ["https://icme.stanford.edu"]

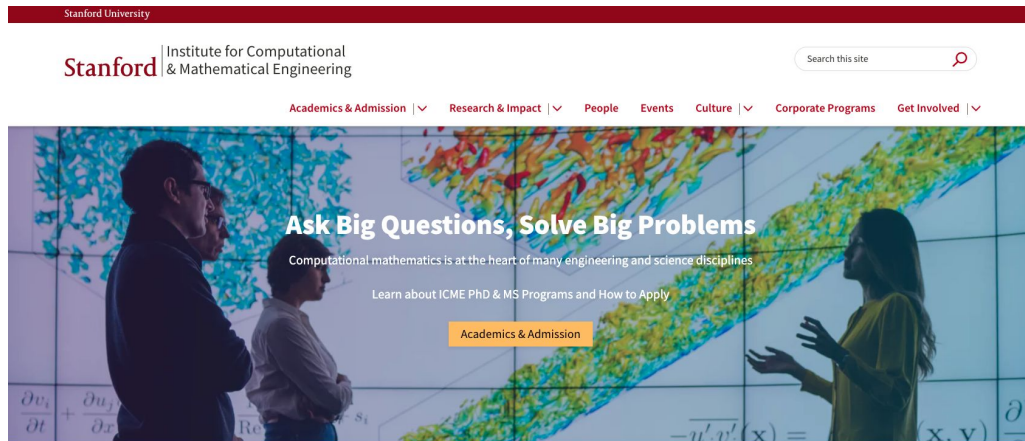
```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```

parsed_urls = ["https://icme.stanford.edu"]

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```



Events & Seminars

JUL 24
Workshop
ICME Summer Workshops 2023 | Fundamentals of Data Science

News

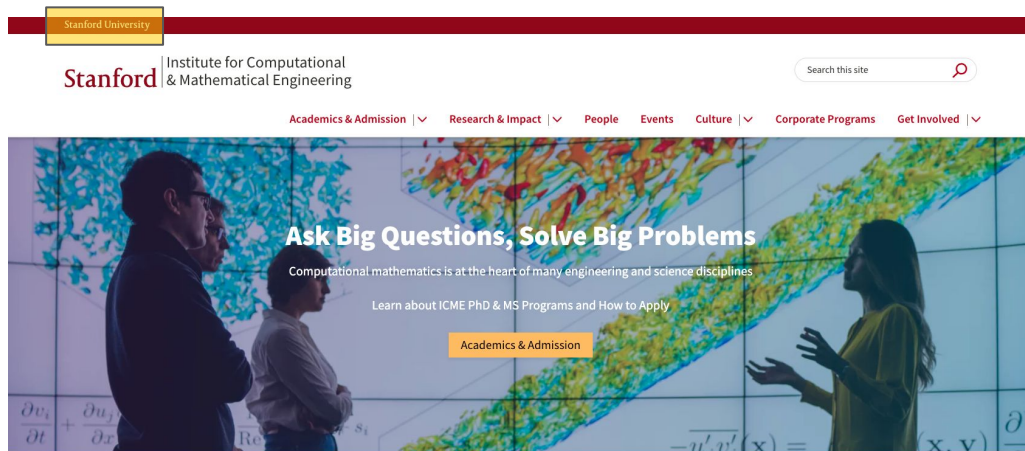
April 10, 2023
Counting Cars: New AI-Driven Approach Finds Times Road Talks



parsed_urls = ["https://icme.stanford.edu"]

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```



Events & Seminars

JUL 24 [ICME Summer Workshops 2023 | Fundamentals of Data Science](#)

News

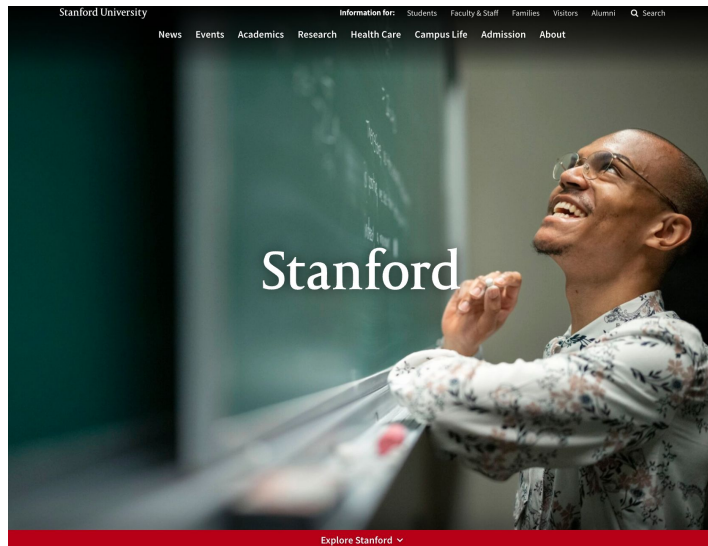
April 10, 2023 [Counting Cars: New AI-Driven Approach Finds Times Road Talks](#)



```
parsed_urls = ["https://icme.stanford.edu"]
```

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

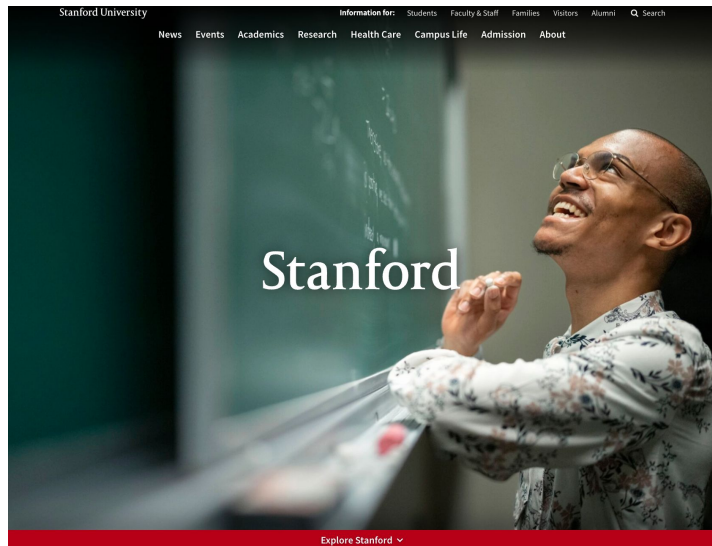
    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```



Search Engine: Web Crawler - Recursion

```
parsed_urls = ["https://icme.stanford.edu",  
"https://www.stanford.edu"]  
]
```

```
class Spider:  
    name = 'icme_spider'  
    start_urls = 'https://icme.stanford.edu/'  
    parsed_urls = []  
  
    def parse(self, url: str):  
        self.parse_url.append(url)  
        for next_url in Website(url):  
            self.parse(next_url)
```

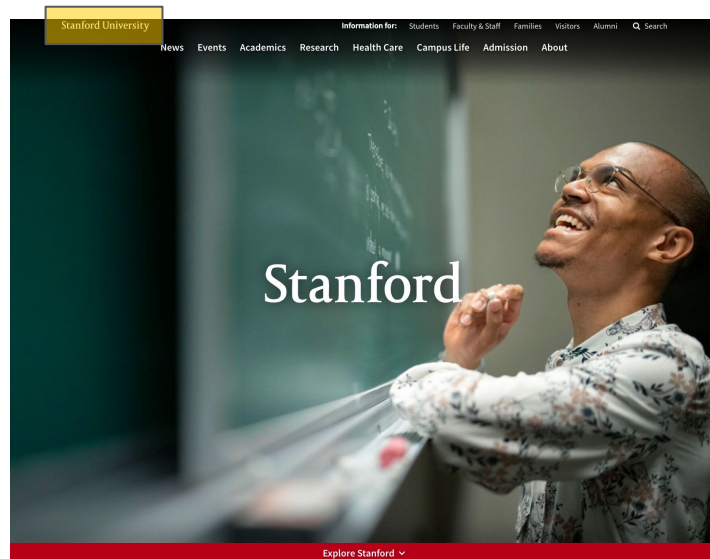


Search Engine: Web Crawler - Recursion

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```

```
parsed_urls = ["https://icme.stanford.edu",
               "https://www.stanford.edu/"]
```



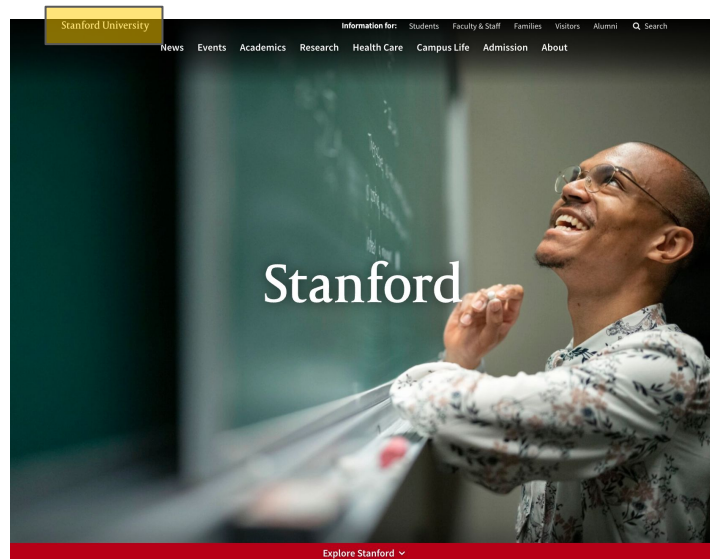
Search Engine: Web Crawler - Recursion

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            self.parse(next_url)
```

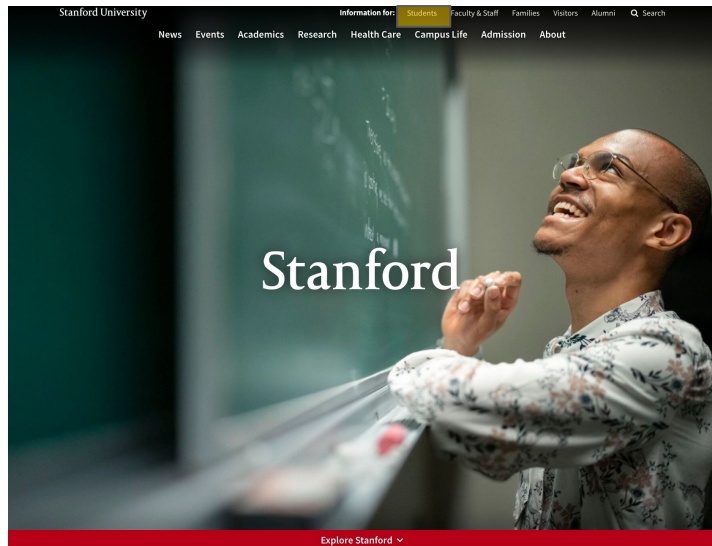
Oh no, it is *“https://www.stanford.edu/”*
again -- we had a bug in the code!

```
parsed_urls = [“https://icme.stanford.edu”,  
“https://www.stanford.edu/”  
]
```



```
parsed_urls = ["https://icme.stanford.edu",  
"https://www.stanford.edu"]
```

```
class Spider:  
    name = 'icme_spider'  
    start_urls = 'https://icme.stanford.edu/'  
    parsed_urls = []  
  
    def parse(self, url: str):  
        self.parse_url.append(url)  
        for next_url in Website(url):  
            if next_url not in self.parsed_urls:  
                self.parse(next_url)
```



Search Engine: Web Crawler - Recursion

```
parsed_urls = ["https://icme.stanford.edu",
               "https://www.stanford.edu",
               "https://www.stanford.edu/student-gateway/"]
```

```
class Spider:
    name = 'icme_spider'
    start_urls = 'https://icme.stanford.edu/'
    parsed_urls = []

    def parse(self, url: str):
        self.parse_url.append(url)
        for next_url in Website(url):
            if next_url not in self.parsed_urls:
                self.parse(next_url)
```

- It has a search bar!
 - -> **How to incorporate the user intention?**
- The items are mal-defined at the first glance.
 - -> How to crawl the internet and store the items?
- User does not simply rate the search results!
 - -> How to assign user-item rating with user data?

Search Engine: String Search

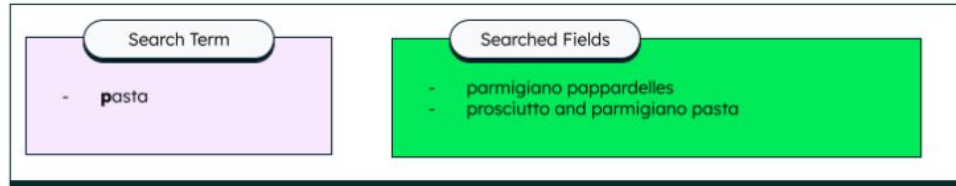


Diagram: <https://www.mongodb.com/basics/full-text-search>

Average: $O(n+m)$; Worst case: $O(mn)$

Rabin-Karp algorithm, which looks for matching substrings, is fast and easy to implement.

Knuth-Morris-Pratt algorithm looks for all instances of a matching character, increasing the speed for multiple matches in a string.



Search results

This wiki is using a new search engine. ([Learn more](#))

[Content pages](#) [Multimedia](#) [Translations](#) [Everything](#) [Advanced](#)

Did you mean: ***andré emotions***

- insertion: *cot* → *coat*
- deletion: *coat* → *cot*
- substitution: *coat* → *cost*

Search Engine: Full-text Search - Inverted Index

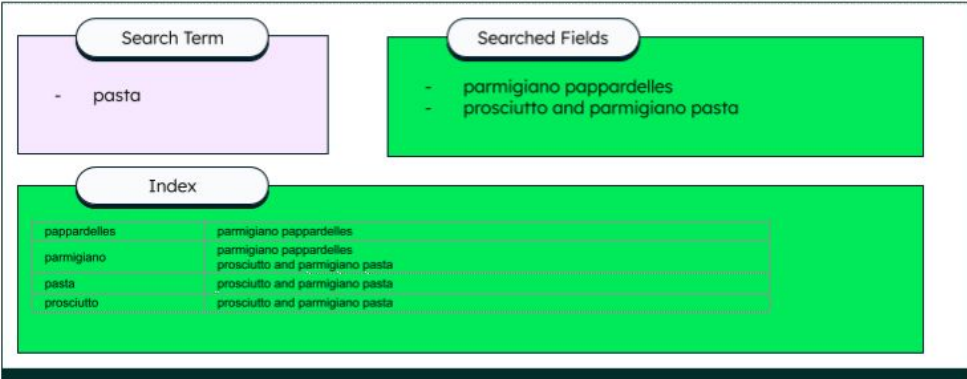


Diagram: <https://www.mongodb.com/basics/full-text-search>

Search Engine: Full-text Search - Inverted Index

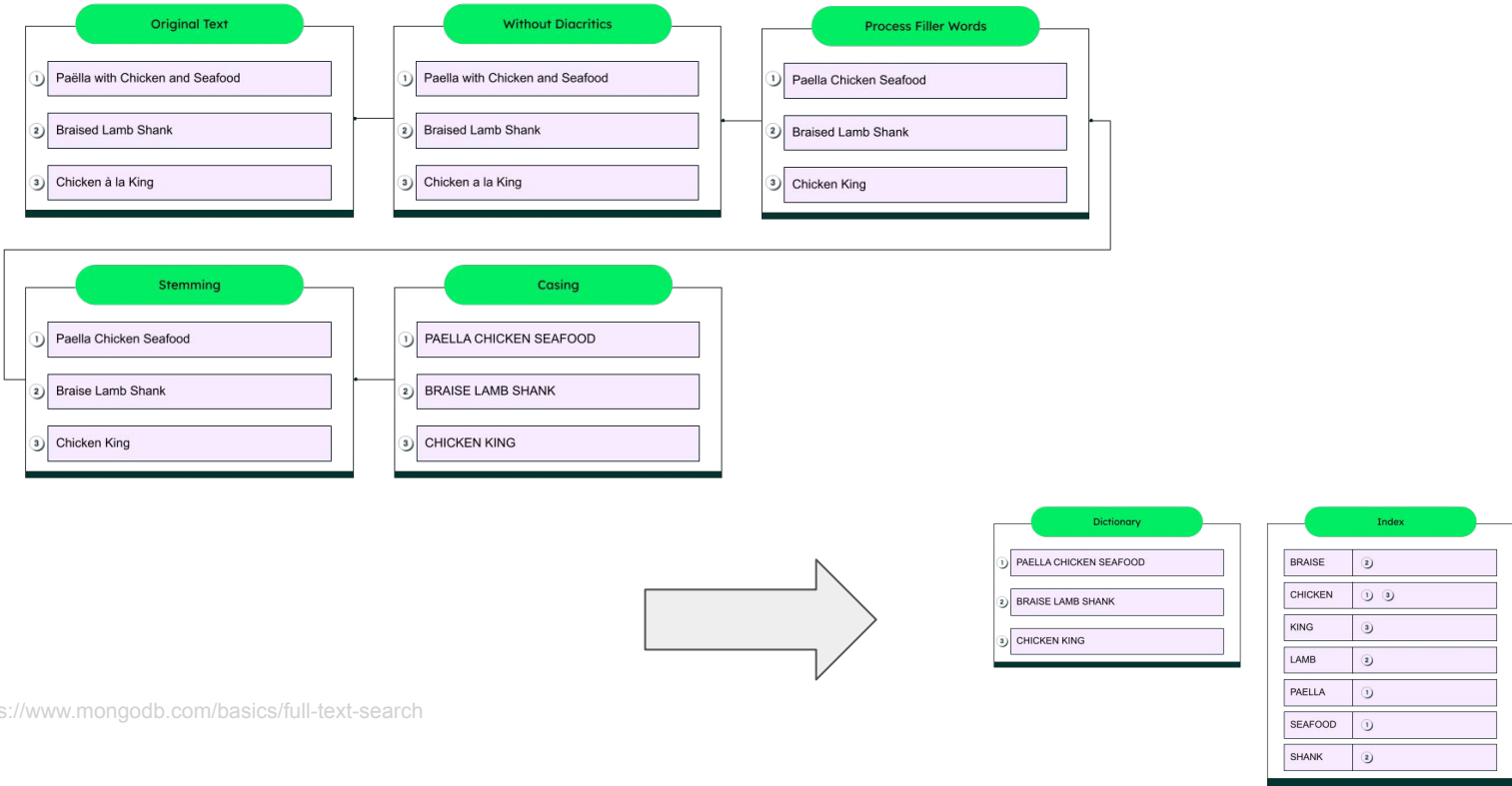


Diagram: <https://www.mongodb.com/basics/full-text-search>

- It has a search bar!
 - -> How to incorporate the user intention?
- The items are mal-defined at the first glance.
 - -> How to crawl the internet and store the items?
- User does not simply rate the search results!
 - -> How to assign user-item rating with user data?

- So far, web pages are treated as individual documents.
- But there are [hyperlinks](#) between them!

Search Engine: PageRank

$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)}$$

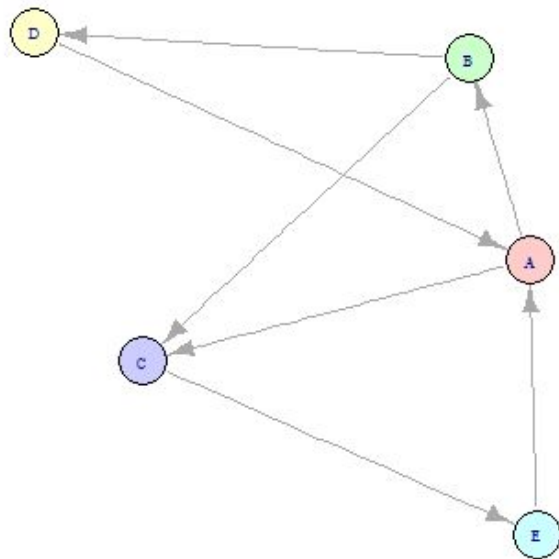
p_1, p_2, \dots, p_N are the pages under consideration.

$M(p_i)$ is the set of pages that link to p_i .

$L(p_j)$ is the number of outbound links on page p_j .

Page Rank of the nodes at start

	Rank
A	0.2
B	0.2
C	0.2
D	0.2
E	0.2



- We have millions of user clicking on some Search Engine Results every day.
- Can we assign click or not as ratings?
 - Yes and no

Search Engine: User-item Rating - Clicks



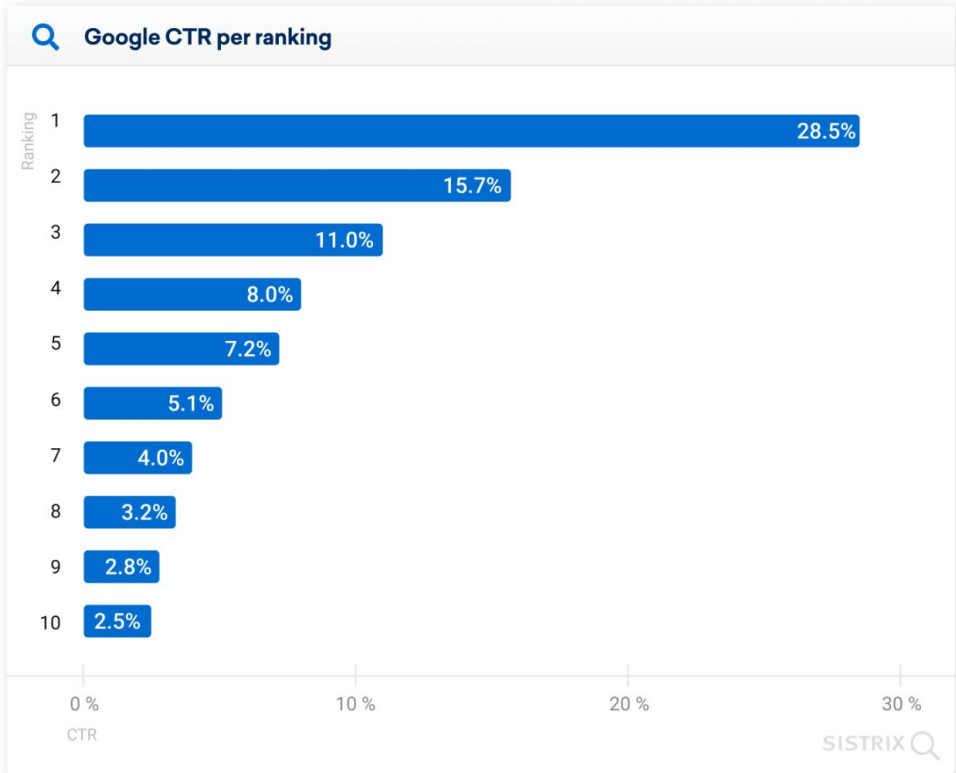
2005



2014

Source: *The Evolution of Google Search Results Pages, Mediative, 2014*

Search Engine: User-item Rating - Clicks



$$\text{CTR} = \frac{\text{CLICKS}}{\text{IMPRESSIONS}} \times 100$$

The equation is visually enhanced with icons: a green mouse cursor icon above the word 'CLICKS', a pink eye icon above the word 'IMPRESSIONS', and the number '100' is colored yellow.

CLICKS
Number of people who clicked the ad

IMPRESSIONS
Number of people who saw the ad

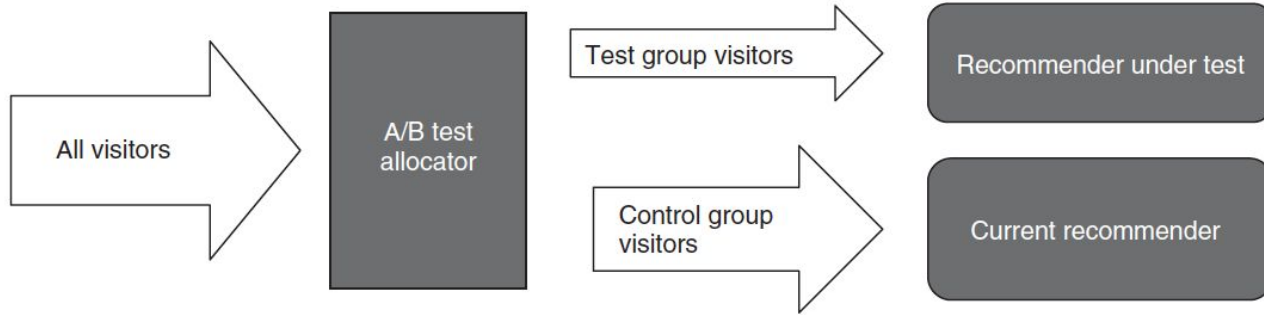


Figure 9.16 In an A/B test, visitors are split into two groups: the test group that sees the new feature and a control group that continues as usual.

- **Q:** For Search Engine, how to measure the user satisfaction and success?
 - Any potential bias?
- **Q:** How about the recommendation system you chose yesterday?

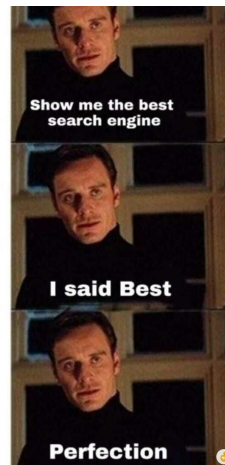


Search Engine: Time-to-long-click (TTLC)

- Long-click: When a user performs a search, clicks through on a result and remains on that site for a long time.
 - Anti-pattern: Pogo-sticking
- Domain specific
- Knowledge panels and direct answers



MOM



Goals of TikTok's Recommender Algorithm

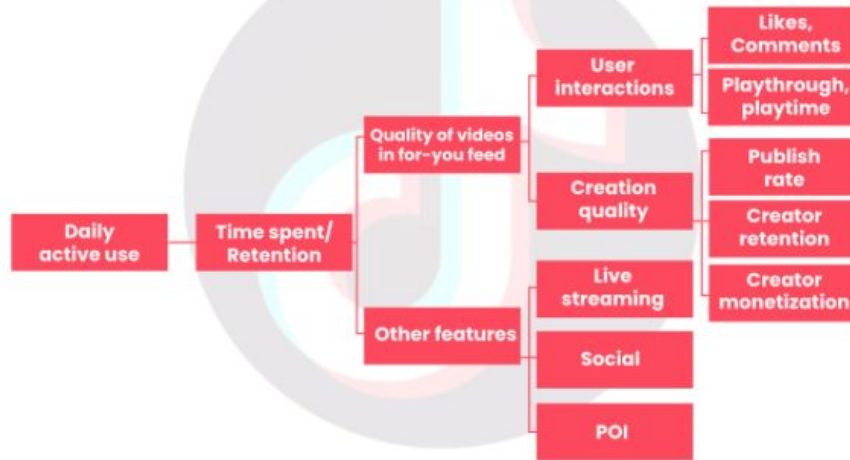


Diagram:

<https://www.linkedin.com/pulse/ai-behind-tiktoks-addictive-algorithm-simple-alexander-stahl/>

Advanced Topics of Recommendation System

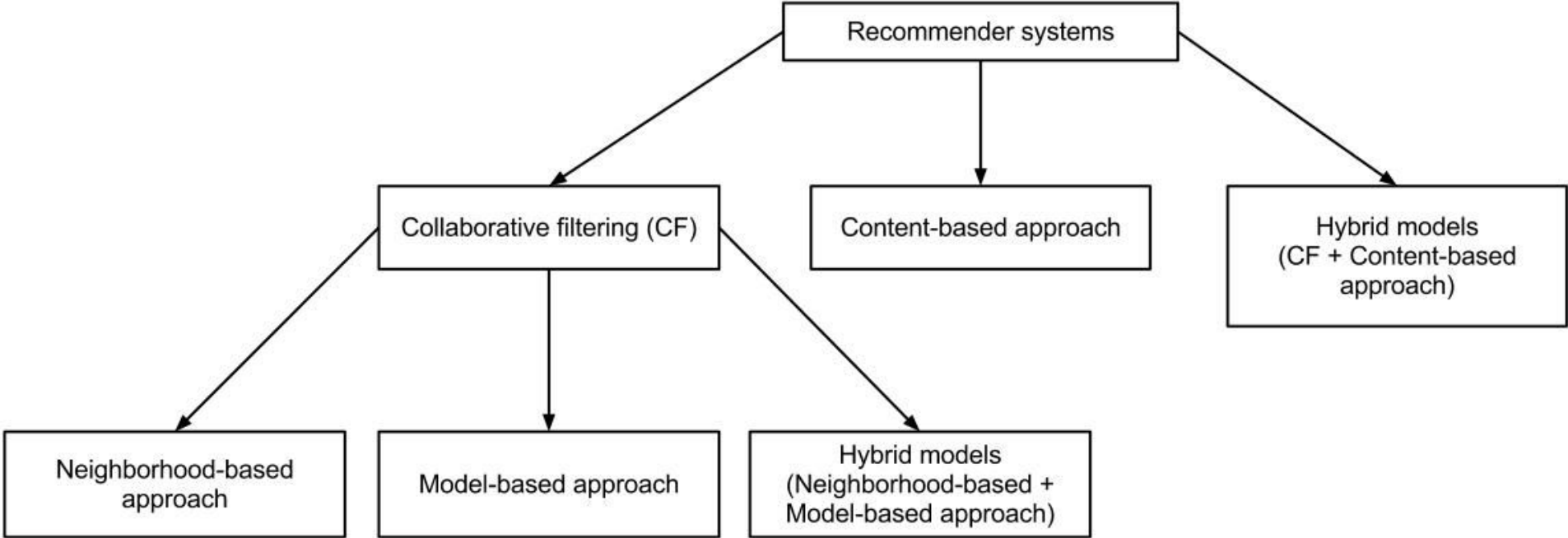
**IF BRUTE FORCE DOESN'T
SOLVE YOUR PROBLEMS**

**THEN YOU AREN'T
USING ENOUGH**

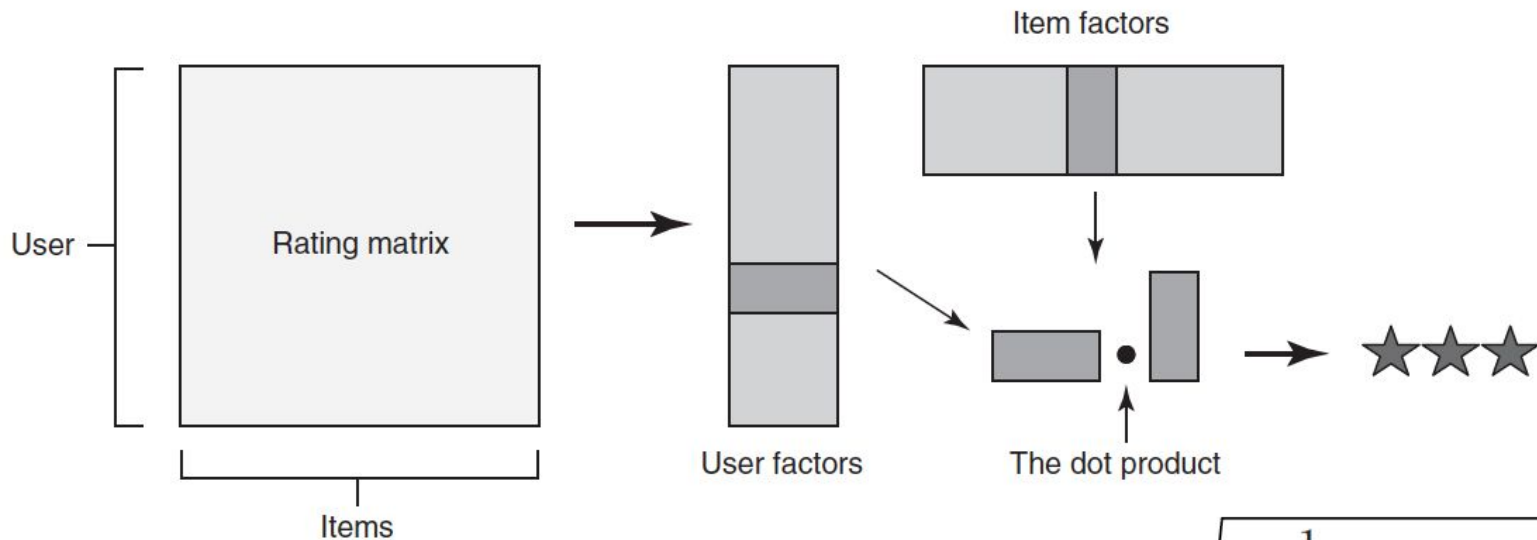


- Deep learning
- Scale up and speed up
- LLM + Recommendation system
- Social impact

Recommendation System w/ Deep Learning: Recap



Recommendation System w/ Deep Learning: Recap

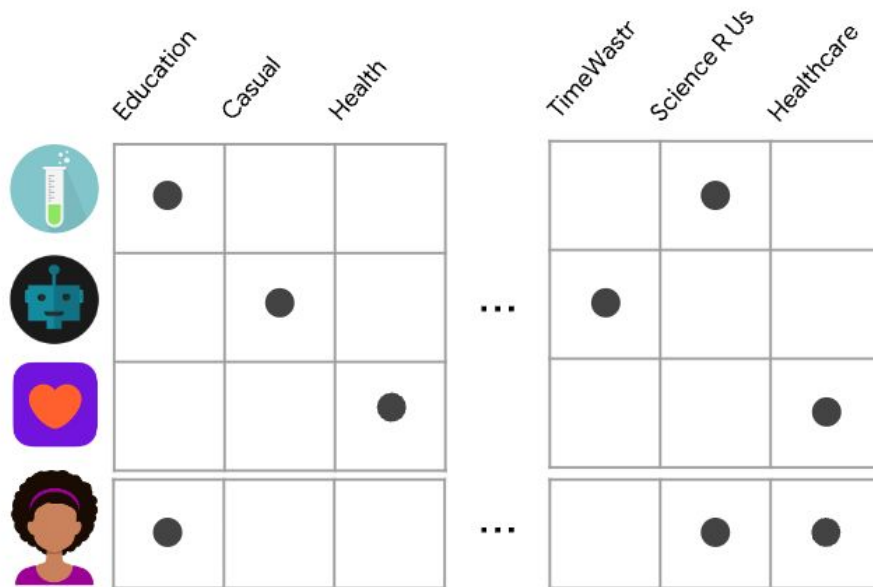


$$RMSE = \sqrt{\frac{1}{|known|} \sum_{(u,i) \in known} (r_{ui} - u_u v_i)^2}$$

- Rating/score can be modeled as a product of user vector and item vector.

Recommendation System w/ Deep Learning: Recap

Example from:
<https://developers.google.com/machine-learning/recommendation/content-based/basics>



- There are well-defined features (of both users and items) that can be used for the prediction.

- Factorization Machines (FM): Let's combine the best of the both worlds!

Reformulate Collaborative Filtering

	i_1	i_2	i_3
u_1	2	4	
u_2		1	
u_3	3		5



							y	
$x^{(0)}$	1	0	0	1	0	0	2	Observed Ratings
$x^{(1)}$	1	0	0	0	1	0	4	
$x^{(2)}$	0	1	0	0	1	0	1	
$x^{(3)}$	0	0	1	1	0	0	3	
$x^{(4)}$	0	0	1	0	0	1	5	
	Users			Items				

Reformulate Collaborative Filtering

	i_1	i_2	i_3
u_1	2	4	
u_2		1	
u_3	3		5



							y
$x^{(0)}$	1	0	0	1	0	0	2
$x^{(1)}$	1	0	0	0	1	0	4
$x^{(2)}$	0	1	0	0	1	0	1
$x^{(3)}$	0	0	1	1	0	0	3
$x^{(4)}$	0	0	1	0	0	1	5
	Users			Items			Observed Ratings

$$\hat{r}_{1,1} = w_0 + u_1^T v_1$$

$$\hat{r}_{1,1} = w_0 + \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (u_1, u_2, u_3)^T \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (v_1, v_2, v_3)$$

Reformulate Collaborative Filtering

	i_1	i_2	i_3
u_1	2	4	
u_2		1	
u_3	3		5



	$x^{(0)}$	1	0	0	1	0	0	2	Observed Ratings
$x^{(1)}$	1	0	0	0	1	0	4		
$x^{(2)}$	0	1	0	0	1	0	1		
$x^{(3)}$	0	0	1	1	0	0	3		
$x^{(4)}$	0	0	1	0	0	1	5		
		Users			Items				

$$\hat{r}_{1,1} = w_0 + u_1^T v_1$$

$$\hat{r}_{1,1} = w_0 + \sum_{i=1}^n \sum_{j=i+1}^n x_i^{(0)} x_j^{(0)} u_i^T v_j$$

Reformulate Content Filtering

E.g. News Popularity

	i_1	i_2	i_3
u_1	2	4	
u_2		1	
u_3	3		5



	i_1	i_2	i_3	a_1	a_2	y
$x^{(0)}$	1	0	0	2.0	0.0	2
$x^{(1)}$	0	1	0	1.5	0.5	4
$x^{(2)}$	0	1	0	1.5	0.5	1
$x^{(3)}$	1	0	0	2.0	0.0	3
$x^{(4)}$	0	0	1	3.2	1.7	5

Items
Auxiliary Features

Observed Ratings

$$\hat{r}_{1,1} = w_0 + w_1 a_1 + w_2 a_2$$

Back to Factorization Machines (FM)

	i_1	i_2	i_3
u_1	2	4	
u_2		1	
u_3	3		5



$x^{(0)}$	1	0	0	1	0	0	2.0	0.0	2	Observed Ratings
$x^{(1)}$	1	0	0	0	1	0	1.5	0.5	4	
$x^{(2)}$	0	1	0	0	1	0	0.0	1.0	1	
$x^{(3)}$	0	0	1	1	0	0	0.3	0.7	3	
$x^{(4)}$	0	0	1	0	0	1	3.2	1.7	5	
	Users			Items			Auxiliary Features			

$$\hat{r} = w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n x_i x_j v_i^T v_j$$

$$\hat{r} = w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n x_i x_j v_i^T v_j$$

- w_0 : Global bias
- w_i : Weights of features
- v_i : Feature factor i

=> Supervised learning paradigm!!

Back to Factorization Machines (FM)

$$\hat{r} = w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n x_i x_j v_i^T v_j$$

- w_0 : Global bias
- w_i : Weights of features
- v_i : Feature factor i

=> Supervised learning paradigm!!

Rating Matrix



Collaborative Filtering



FM



Welcome to the
Deep-world

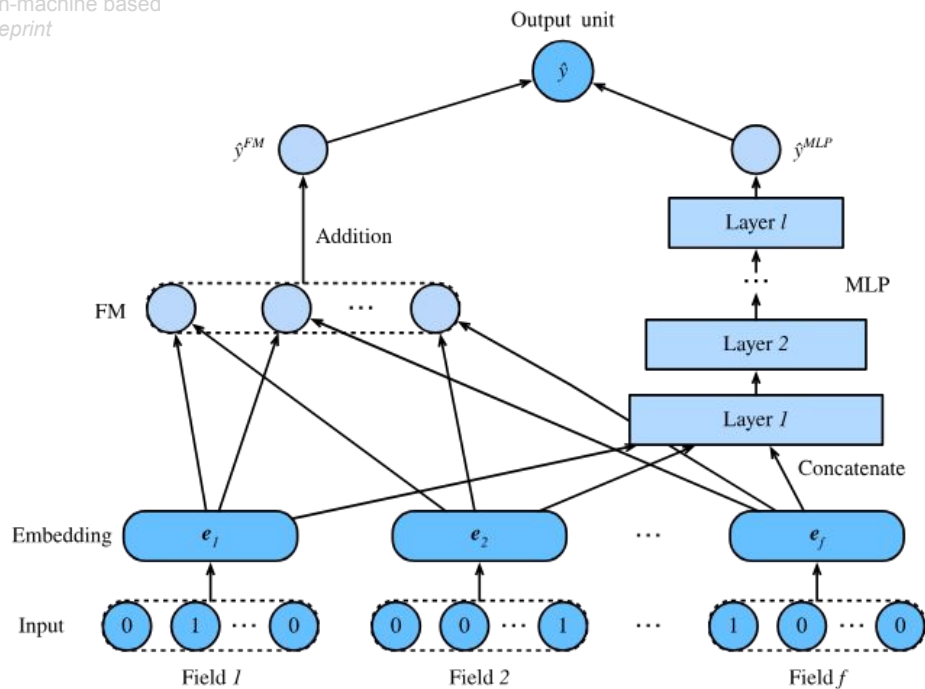


Factorization Machines (FM): Pros & Cons

- Capture user-item feature interaction
- Efficient for sparse data
- Non-linear patterns

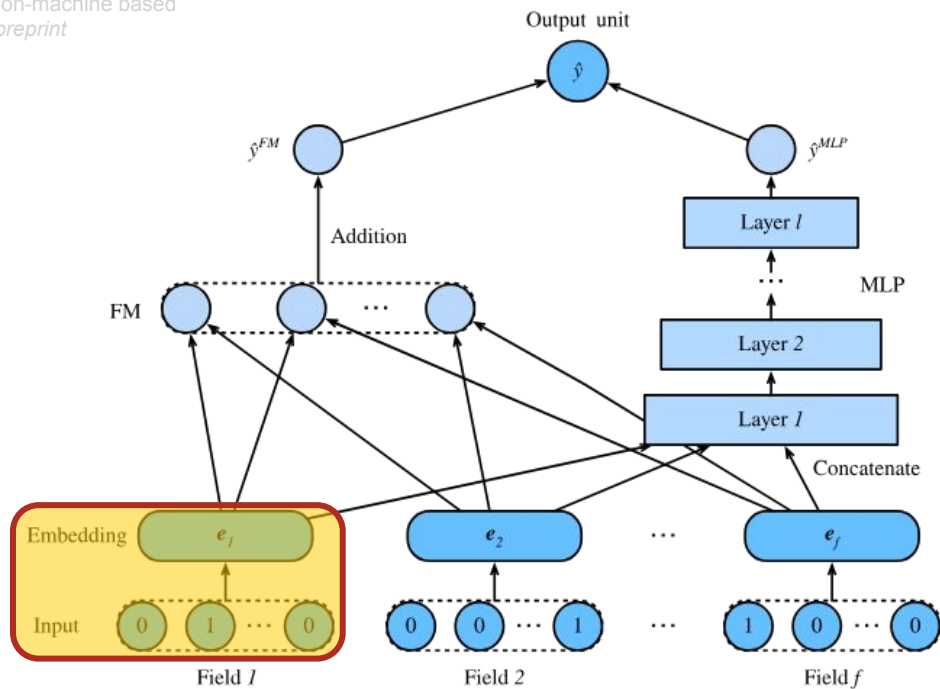
DeepFM (2017)

Guo, Huifeng, et al. "DeepFM: a factorization-machine based neural network for CTR prediction." *arXiv preprint arXiv:1703.04247* (2017).



DeepFM (2017)

Guo, Huifeng, et al. "DeepFM: a factorization-machine based neural network for CTR prediction." *arXiv preprint arXiv:1703.04247* (2017).

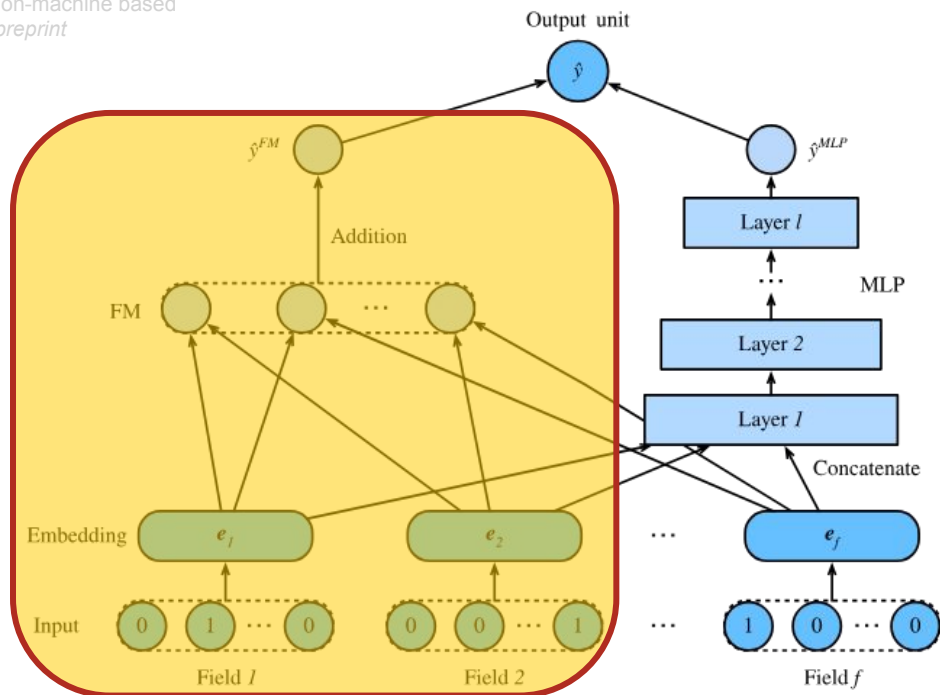


Factors =
Embeddings

DeepFM (2017)

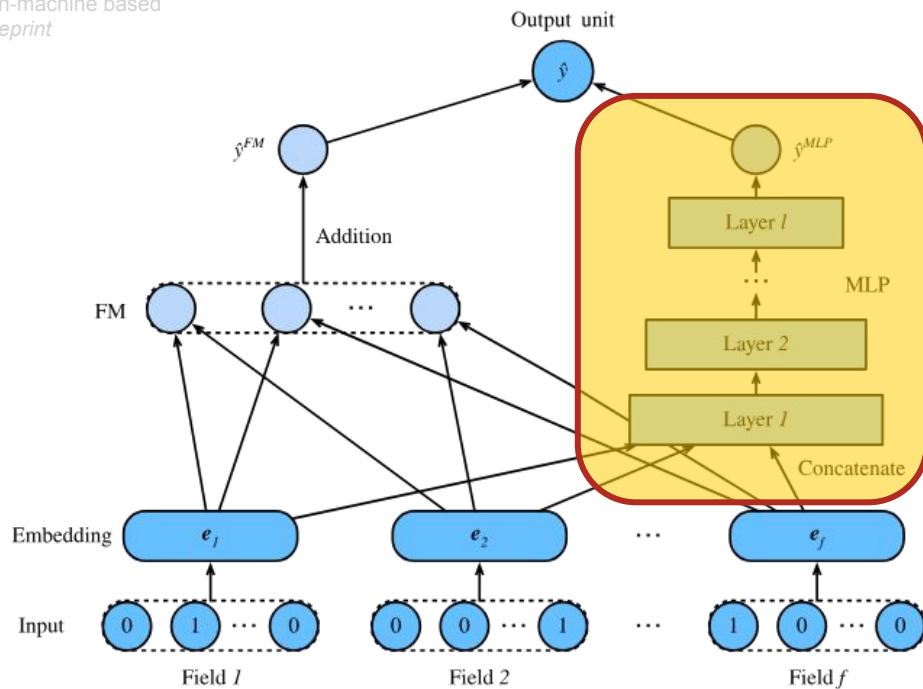
Guo, Huifeng, et al. "DeepFM: a factorization-machine based neural network for CTR prediction." *arXiv preprint arXiv:1703.04247* (2017).

Factorization Machines



DeepFM (2017)

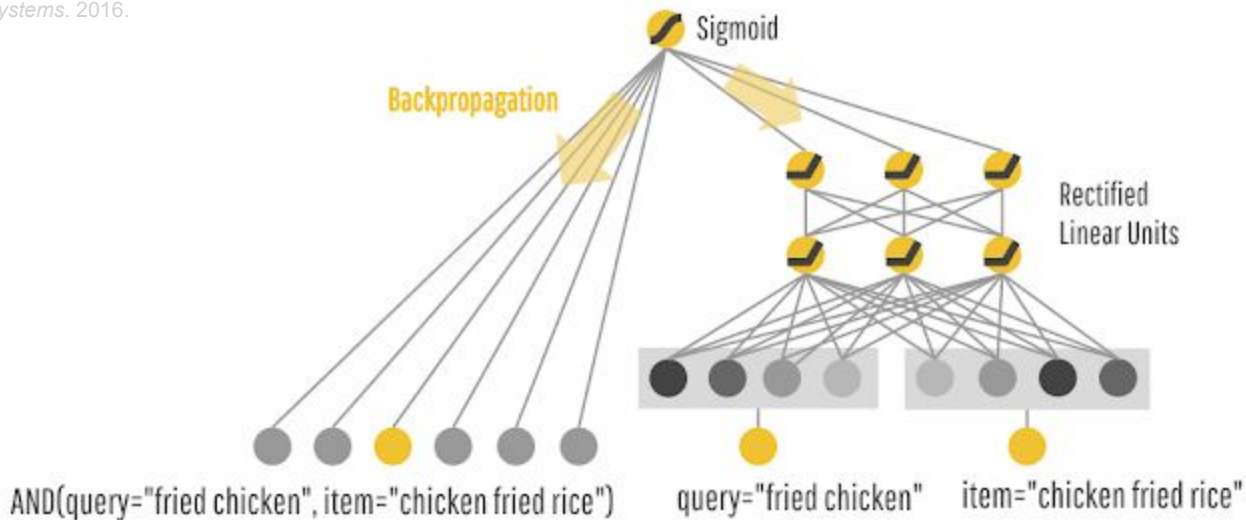
Guo, Huifeng, et al. "DeepFM: a factorization-machine based neural network for CTR prediction." *arXiv preprint arXiv:1703.04247* (2017).



Deep MLP

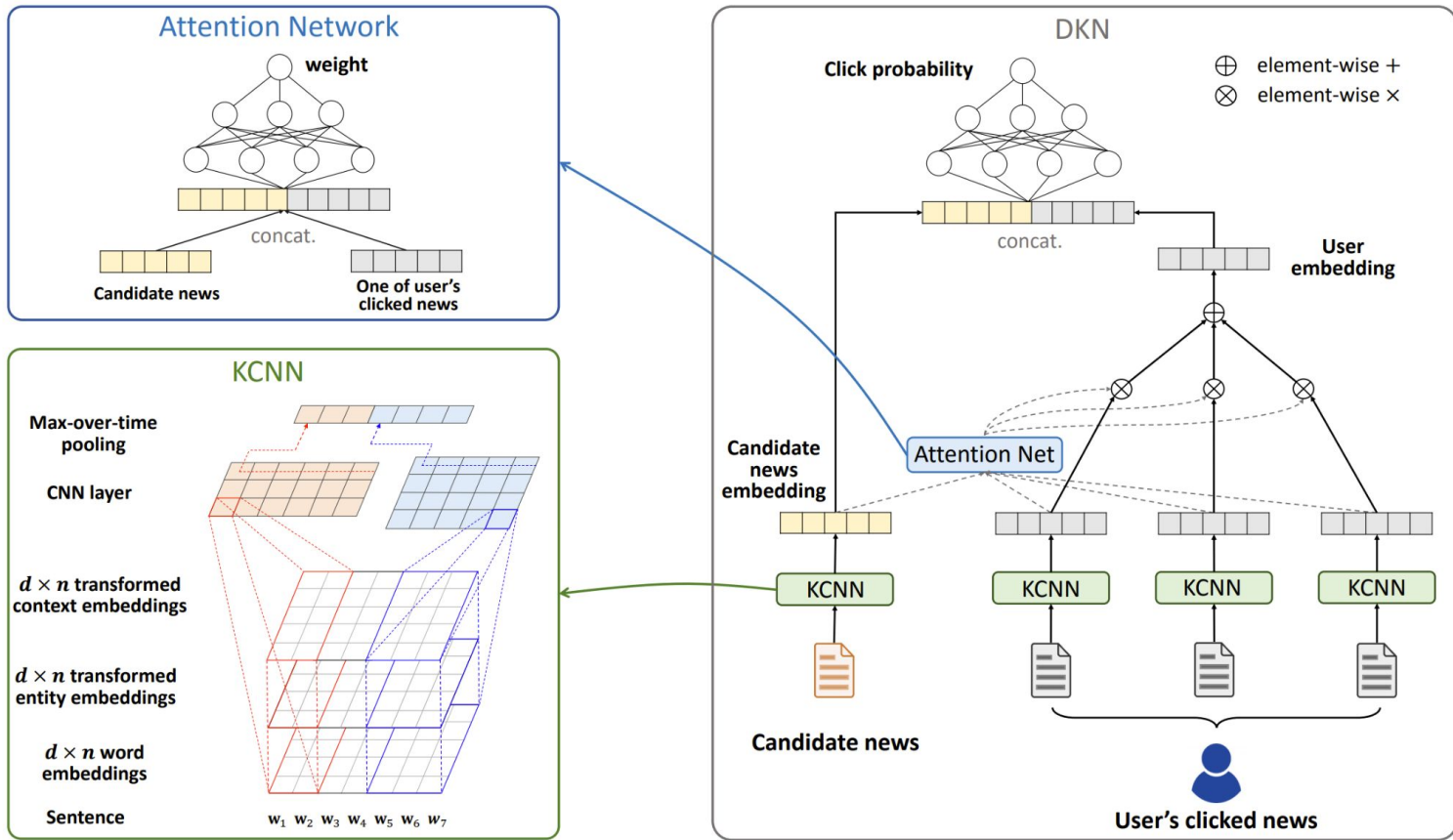
Wide & Deep (2016)

Cheng, Heng-Tze, et al. "Wide & deep learning for recommender systems." *Proceedings of the 1st workshop on deep learning for recommender systems*. 2016.



Deep Knowledge-Aware Network (2018)

Wang, Hongwei, et al. "DKN: Deep knowledge-aware network for news recommendation." *Proceedings of the 2018 world wide web conference*. 2018.



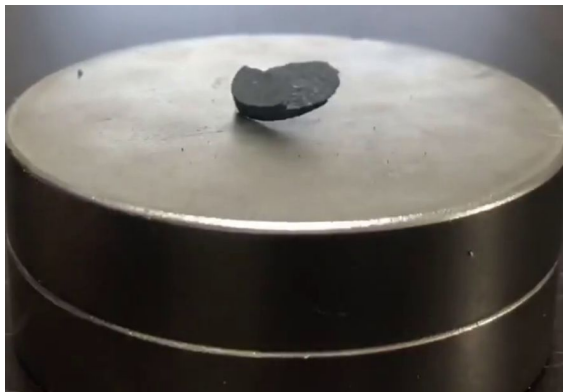
- Sparse Features: Tag, device, etc.
- Dense Features: Age, income, # of videos watched, etc.
- **Sparse Features** are the first-class citizens.

- **Sparse Features** are the first-class citizens.
 - Dense feature requires more parameterization.
 - Age -> like this YouTube Video
 - Assumption: Age = 9 -> Age = 10 has a same effect of Age = 43 -> Age = 44

Feature Engineering for Deep Recommendation

- **Sparse Features** are the first-class citizens.
 - Dense feature requires more parameterization.
 - Sparse features are easy to generate cross-features.
 - Hashing tricks: Only keep K more frequent tuples of the combinations.

#lk99 #superconductor



#lk99 #drinking-solution



LK99

Elkay Deluxe 3-1/2" Drain Type 304 Stainless Steel Body St
Rubber Seal and Tailpiece

\$169.00 (USD)

Actual selling price may vary

Where to Buy

Contact Us

★★★★★ 4.7 | 20 Reviews

14 out of 16 (88%) reviewers recommend this product

14 questions and 14 answers for this product

WRITE A REVIEW

Share: [Twitter](#) [Facebook](#) [Pinterest](#) [Email](#)

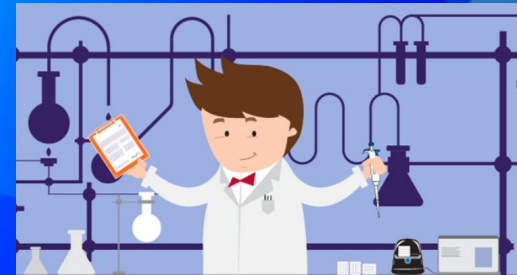
Downloads

[Specification Sheet \(pdf\)](#)

[Care Cleaning \(pdf\)](#)

- **Sparse Features** are the first-class citizens.
 - Dense feature requires more parameterization.
 - Sparse features are easy to generate cross-features.
 - * Easy for online training/serving

Lab Time!



Advanced Topics (Part II)



- Deep learning
- **Scale up and speed up**
- LLM + recommendation system
- Social impact

*All other things being equal ... our experiments demonstrate that slowing down the search results page by **100 to 400** milliseconds has a measurable impact on the number of searches per user of **-0.2%** to **-0.6%**. That's 0.2% to 0.6% fewer searches for changes under half a second!*

-- Google Research Blog

Speed-up: Recommendation Funnel

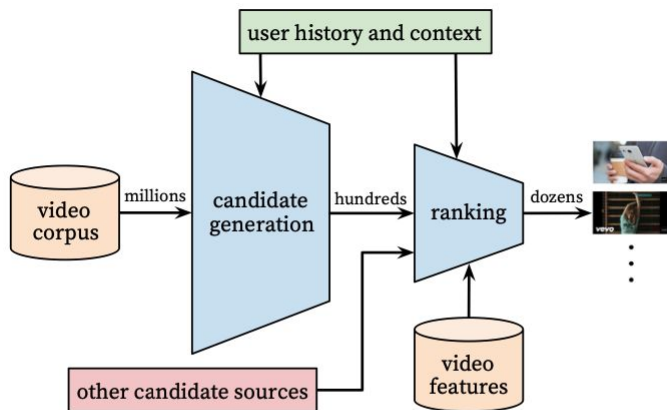


Figure 2: Recommendation system architecture demonstrating the “funnel” where candidate videos are retrieved and ranked before presenting only a few to the user.

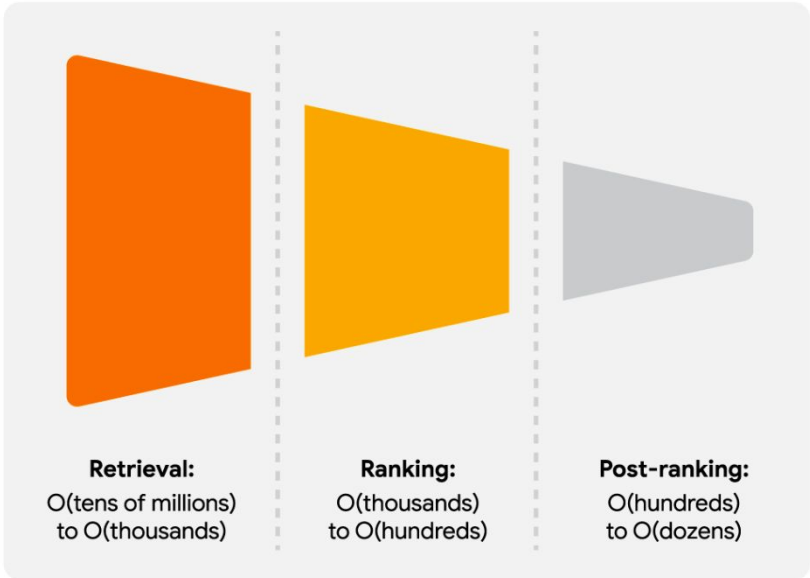
- Retrieval / Candidate-gen
- Ranking / Sort
- Re-rank

Speed-up: Recommendation Funnel

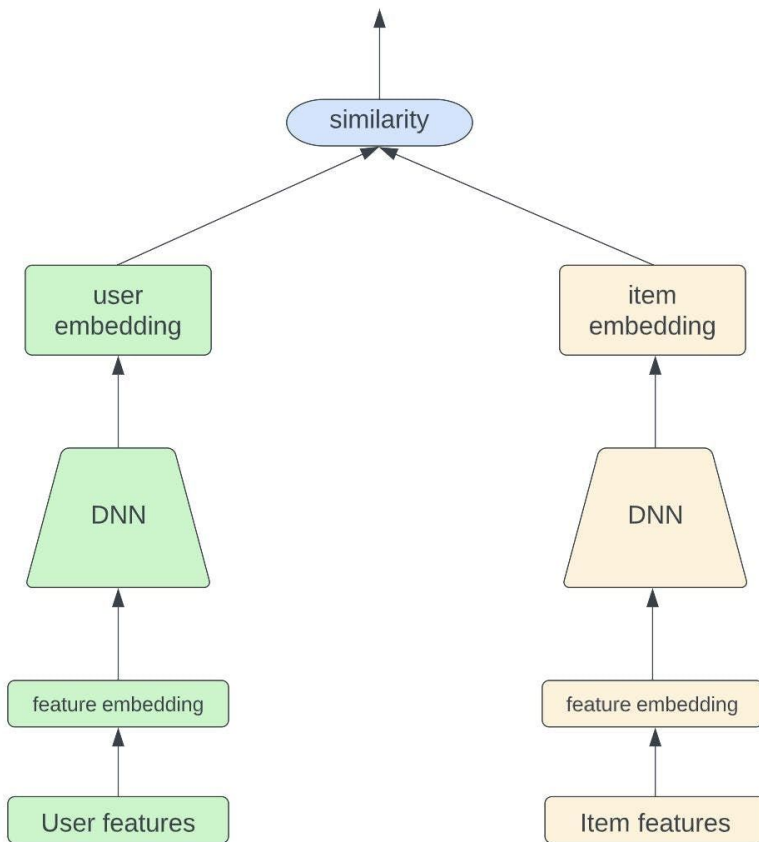
Candidate items



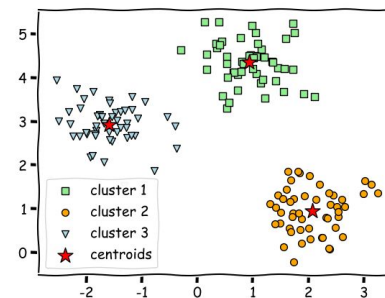
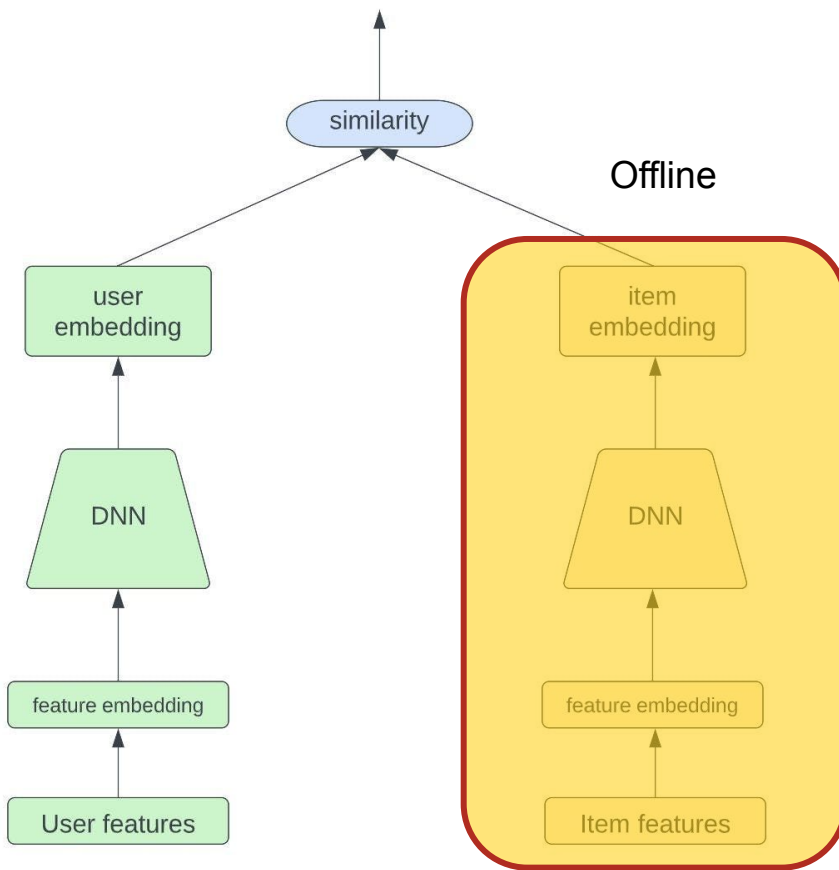
Recommendation engine



Speed-up: Two-tower Retrieval

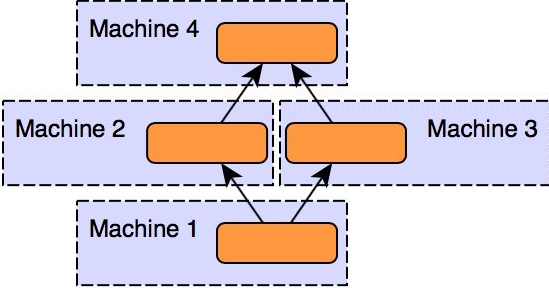


Speed-up: Two-tower Retrieval

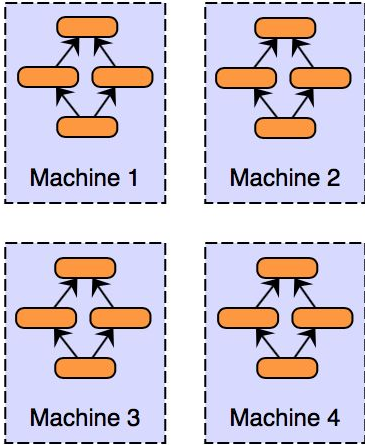


Scale-up: Parameter Server (PS)

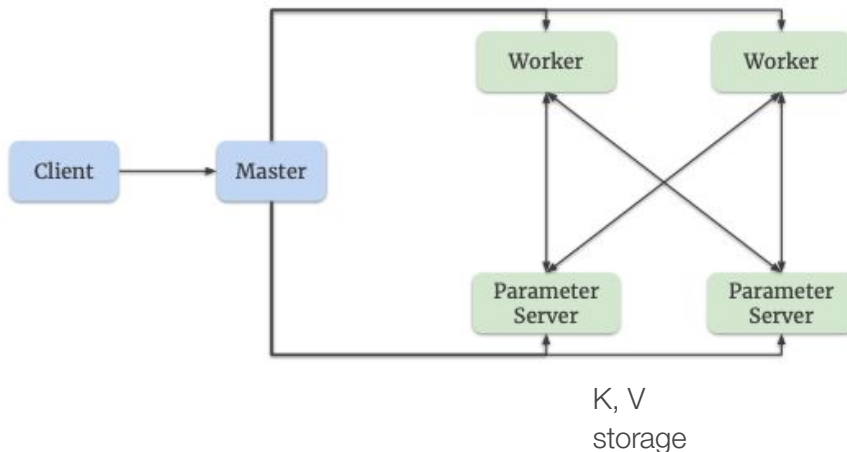
Model Parallelism



Data Parallelism



Scale-up: Parameter Server (PS)



Algorithm 1 Distributed Subgradient Descent

Task Scheduler:

- 1: issue LoadData() to all workers
- 2: **for** iteration $t = 0, \dots, T$ **do**
- 3: issue WORKERITERATE(t) to all workers.
- 4: **end for**

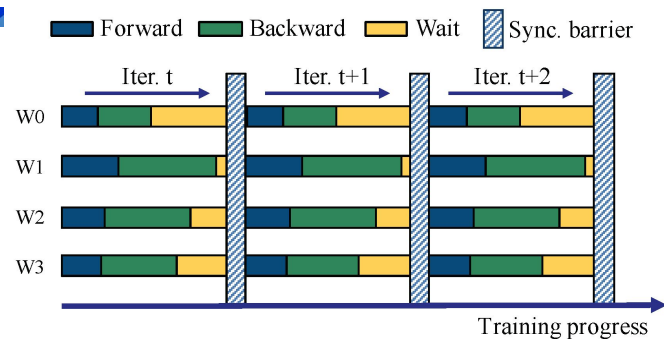
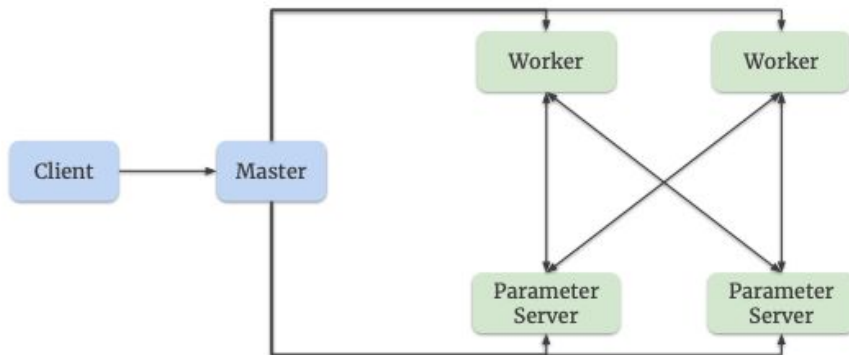
Worker $r = 1, \dots, m$:

- 1: **function** LOADDATA()
- 2: load a part of training data $\{y_{i_k}, x_{i_k}\}_{k=1}^{n_r}$
- 3: pull the working set $w_r^{(0)}$ from servers
- 4: **end function**
- 5: **function** WORKERITERATE(t)
- 6: gradient $g_r^{(t)} \leftarrow \sum_{k=1}^{n_r} \partial \ell(x_{i_k}, y_{i_k}, w_r^{(t)})$
- 7: push $g_r^{(t)}$ to servers
- 8: pull $w_r^{(t+1)}$ from servers
- 9: **end function**

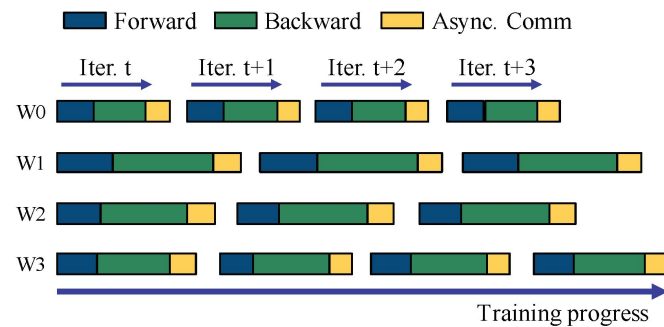
Servers:

- 1: **function** SERVERITERATE(t)
- 2: aggregate $g^{(t)} \leftarrow \sum_{r=1}^m g_r^{(t)}$
- 3: $w^{(t+1)} \leftarrow w^{(t)} - \eta (g^{(t)} + \partial \Omega(w^{(t)}))$
- 4: **end function**

Scale-up: Parameter Server (PS)



(a)



(b)

Augmenting recommendation systems with LLMs (Large Language Model)

when you need advice but aren't sure who to trust



Large Language Models

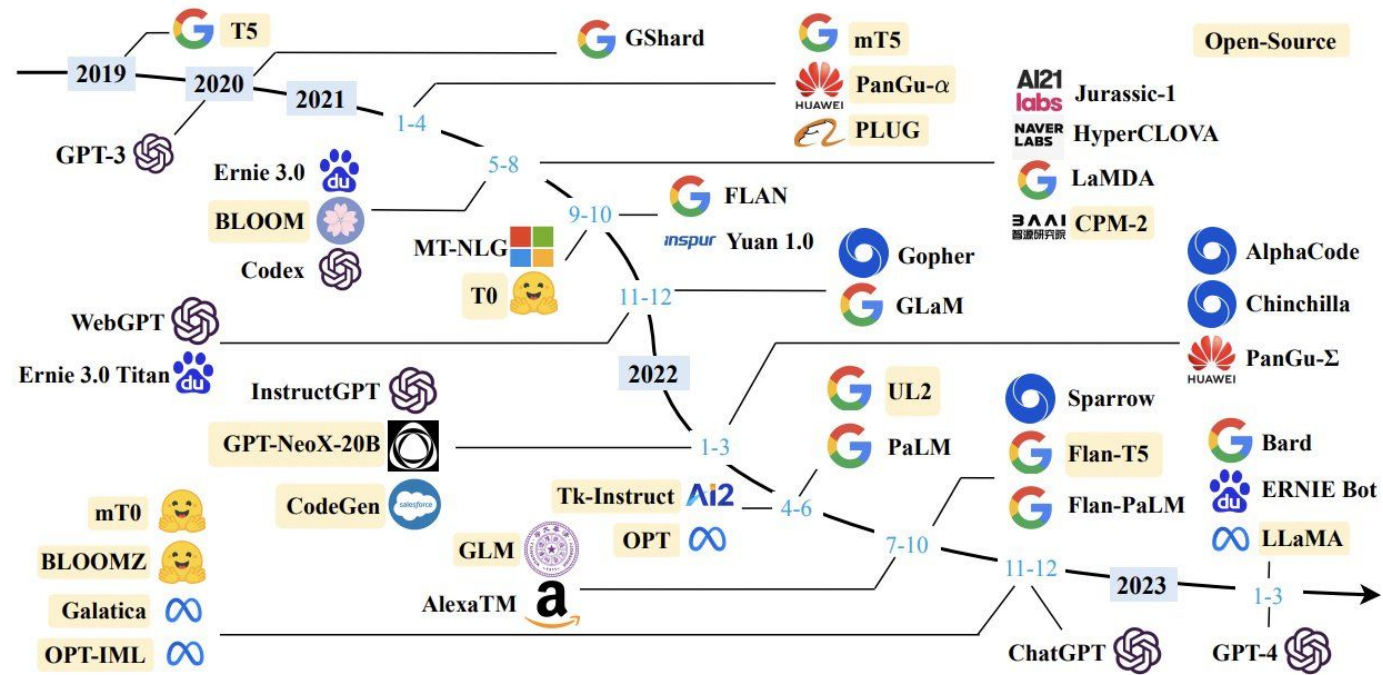






Fig. 1. A timeline of existing large language models (having a size larger than 10B) in recent years. We mark the open-source LLMs in yellow color.

Diagram:
<https://www.labellerr.com/blog/overview-of-development-of-large-lannguage-models/>









I'm in the mood for some drama movies with artistic elements tonight. Could you recommend three? Titles only. No other text. 

[View other drafts](#) 

 Sure, here are three drama movies with artistic elements that you might enjoy:

- The Tree of Life
- The Piano Teacher
- The Passion of Joan of Arc

    Google it  Rewrite 

- Fluid experience
- Limited candidates / High inference cost

```
prompt = """You are a movie recommender and your job is to predict
            a user's rating (ranging from 1 to 5, with 5 being the highest)
            on a movie, based on that user's previous ratings.

User 42 has rated the following movies:
"Moneyball" 4.5
"The Martian" 4
"Pitch Black" 3.5
"12 Angry Men" 5

Predict the user's rating on "The Matrix".
Output the rating score only.
Do not include other text.
"""

response = palm.generate_text(model="models/text-bison-001", prompt=prompt)
print(response.result)

# 4.5
```

```
prompt = """You are a movie recommender and your job is to recommend new movies
            based on the sequence of movies that a user has watched. You pay special
            attention to the order of movies because it matters.
```

```
User 42 has watched the following movies sequentially:
```

```
"Margin Call",
"The Big Short",
"Moneyball",
"The Martian",
```

```
Recommend three movies and rank them in terms of priority.
```

```
Titles only.
```

```
Do not include any other text.
```

```
"""
```

```
response = palm.generate_text(
    model="models/text-bison-001", prompt=prompt, temperature=0
)
```

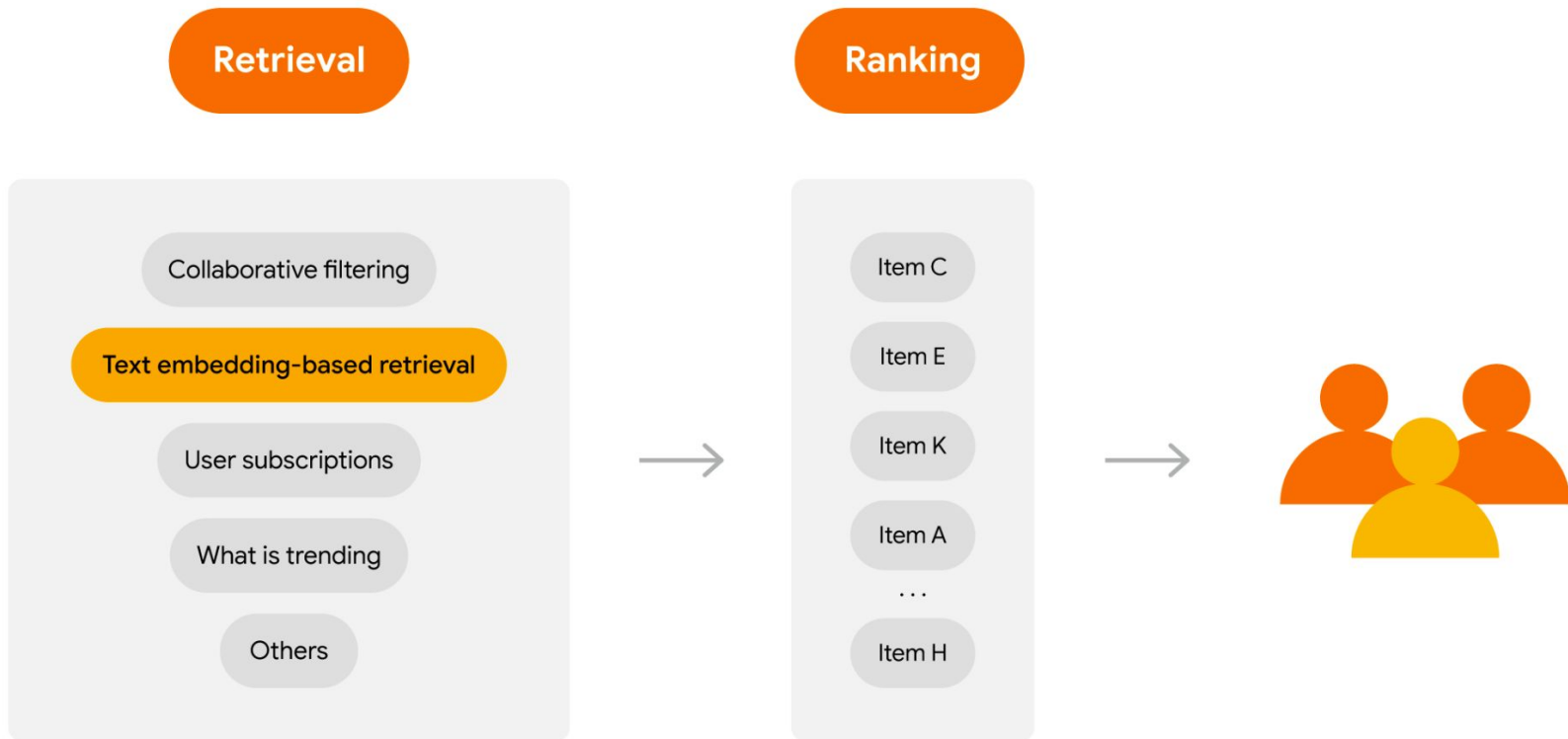
```
print(response.result)
```

```
# 1. The Wolf of Wall Street
```

```
# 2. The Social Network
```

```
# 3. Inside Job
```

Text embedding-based recommendations



Social Impact



- Fairness
 - Unfair/inaccurate recommendation
- Echo chamber
- Privacy

- Manipulation
 - Trading and nudging [1]
- Promote Addiction [2]
- Privacy



[1] Burr, C., Cristianini, N., & Ladyman, J. (2018). An Analysis of the Interaction Between Intelligent Software Agents and Human Users. *Minds & Machines*, 28(4), 735–774. Retrieved from

<https://link.springer.com/article/10.1007/s11023-018-9479-0>

[2] Seaver, N. (2018). Captivating algorithms: Recommender systems as traps. *Journal of Material Culture*. Retrieved from <https://static1.squarespace.com/static/55eb004ee4b0518639d59d9b/t/5b707506352f5356c8d6e7d2/1534096646595/seaver-captivating-algorithms.pdf>

- Manipulation
 - Trading and nudging [1]
- Promote Addiction [2]
- Privacy



[1] Burr, C., Cristianini, N., & Ladyman, J. (2018). An Analysis of the Interaction Between Intelligent Software Agents and Human Users. *Minds & Machines*, 28(4), 735–774. Retrieved from

<https://link.springer.com/article/10.1007/s11023-018-9479-0>

[2] Seaver, N. (2018). Captivating algorithms: Recommender systems as traps. *Journal of Material Culture*. Retrieved from <https://static1.squarespace.com/static/55eb004ee4b0518639d59d9b/t/5b707506352f5356c8d6e7d2/1534096646595/seaver-captivating-algorithms.pdf>



 Mentimeter

Join at [menti.com](https://www.menti.com) use code 1265 3592