



aivancity

SCHOOL FOR

TECHNOLOGY, BUSINESS & SOCIETY

PARIS-CACHAN

04/10/2024

Natural Language Processing (NLP)

Introduction to NLP and Distributional Semantics

Quick word about me

- Postdoc Researcher at ISIR-CNRS (Sorbonne University)
- PhD from Paris-Saclay University (LISN-CNRS lab)
- Research topic: Multimodal and Multilingual NLP
- More about me: <https://paullerner.github.io>
- Contact: lerner@isir.upmc.fr

Acknowledgements

This class directly builds upon:

- **Jurafsky, D., & Martin, J. H.** (2024). *Speech and Language Processing : An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models* (3rd éd.).
- **Eisenstein, J.** (2019). *Natural Language Processing*. 587.
- **Yejin Choi.** (Winter 2024). CSE 447/517: Natural Language Processing (University of Washington - Paul G. Allen School of Computer Science & Engineering)
- **Noah Smith.** (Winter 2023). CSE 447/517: Natural Language Processing (University of Washington - Paul G. Allen School of Computer Science & Engineering)
- **Benoît Sagot.** (2023-2024). *Apprendre les langues aux machines* (Collège de France)
- **Chris Manning.** (Spring 2024). Stanford CS224N: Natural Language Processing with Deep Learning
- Classes where I was/am Teacher Assistant:
 - **Christopher Kermorvant.** Machine Learning for Natural Language Processing (ENSAE)
 - **François Landes** and **Kim Gerdes.** Introduction to Machine Learning and NLP (Paris-Saclay)

Also inspired by:

- My PhD thesis: *Répondre aux questions visuelles à propos d'entités nommées* (2023)
- **Noah Smith** (2023): Introduction to Sequence Models (LxMLS)
- **Kyunghyun Cho:** Transformers and Large Pretrained Models (LxMLS 2023), Neural Machine Translation (ALPS 2021)
- My former PhD advisors **Olivier Ferret** and **Camille Guinaudeau** and postdoc advisor **François Yvon**
- My former colleagues at LISN

Program for this semester

- Today: What is NLP? What is a word? How do you get a sense of a word?
 - NLP = research field at the intersection of Computer Science and Linguistics / Technology at the heart of chatbots like ChatGPT
 - Meaning of a word is its use in the language: distributional semantics
- Oct 10: Practical Work 1 (2 sessions)
- Oct 16: Neural Network architectures used in Large Language Models:
 - Attention Mechanism
 - Transformers

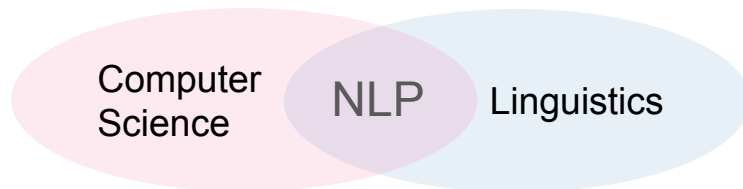
Program for this semester

- Oct 24: Practical Work 2 (2 sessions)
- Group Homework: deadline Monday 4th of November (after Toussaint)
 - Groups of 3
 - Report of max. 4 pages, Continuous assessment (50%)
- Nov 5: Large Language Models from Shannon to ChatGPT
 - pre-training and fine-tuning
 - alignment: reinforcement learning from human feedback (RLHF)
 - decoding/generation methods

Program for this semester

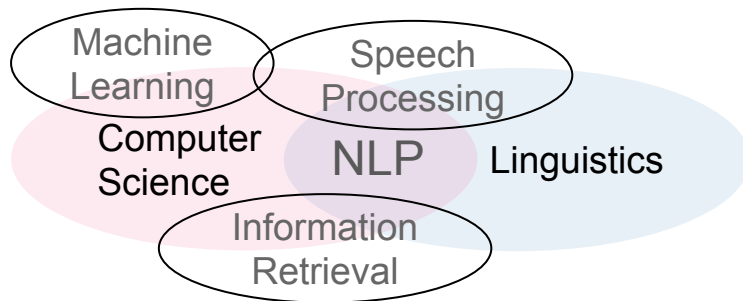
- Nov 19: Practical Work 3 (2 sessions)
- Nov 28:
 - Industrial applications and research benchmarks
 - Ethical, social, and environmental issues
- Dec 5: Practical Work 4 (2 sessions)
- Individual Final sitting Exam 50% (December, before Christmas)

Natural Language Processing (NLP)



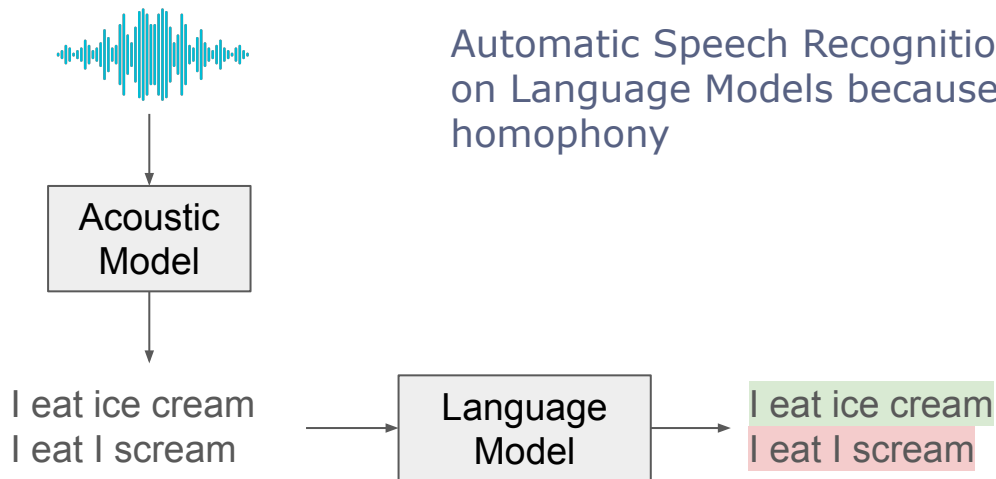
- Intersection of Computer Science and Linguistics:
 - Distributional Semantics: sense of a word from its context (today class)
 - Computational Linguistics, Computational Morphology, etc.:
study of humans: how do we speak? how do we organize lexicon?

Natural Language Processing (NLP)

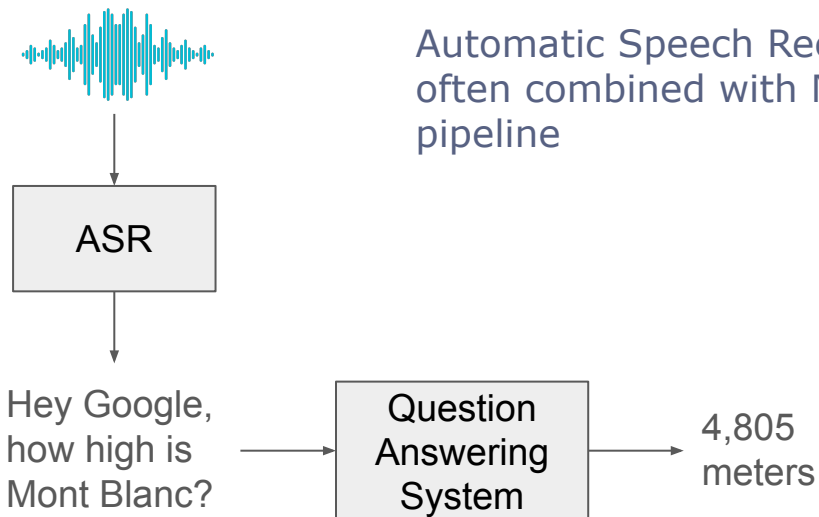


- Close to Speech Processing (Automatic Speech Recognition etc.)
- Close to Information Retrieval (Search engines like Google)
- Driven by Statistical/Machine Learning methods since the 90s (Brown, P. F., Della Pietra, S. A., Della Pietra, V. J., & Mercer, R. L. (1993). The Mathematics of Statistical Machine Translation : Parameter Estimation. Computational Linguistics, 19(2), 263-311.)
- Driven by Deep Learning since 2013 (Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed Representations of Words and Phrases and their Compositionality. Advances in Neural Information Processing Systems)

Speech and Language Processing



Speech and Language Processing



Speech and Language Processing

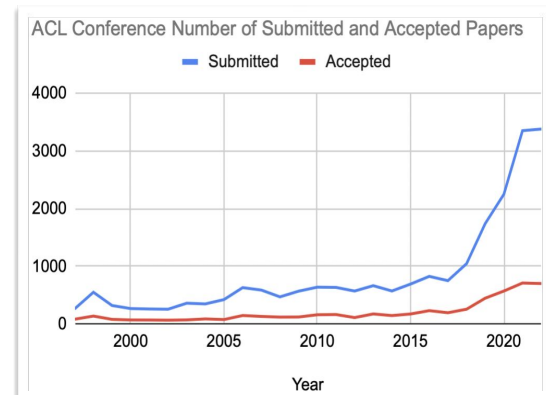
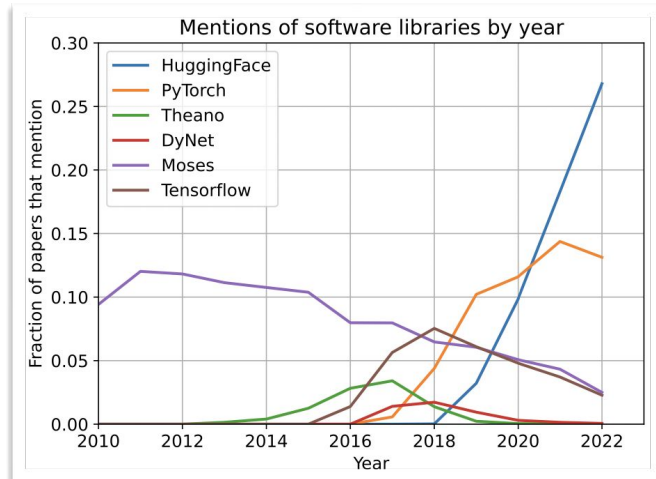
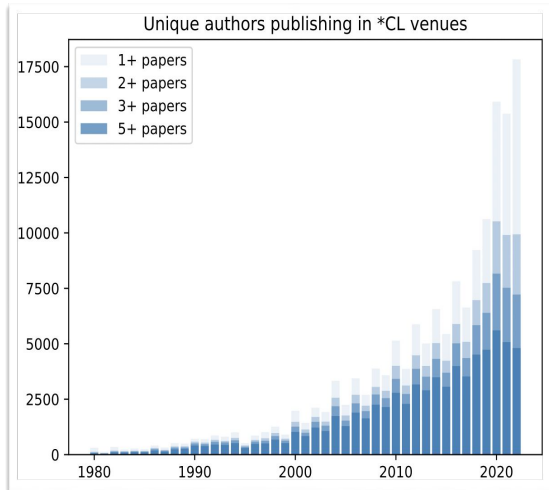


Recently moving towards integrated, multimodal end-to-end models

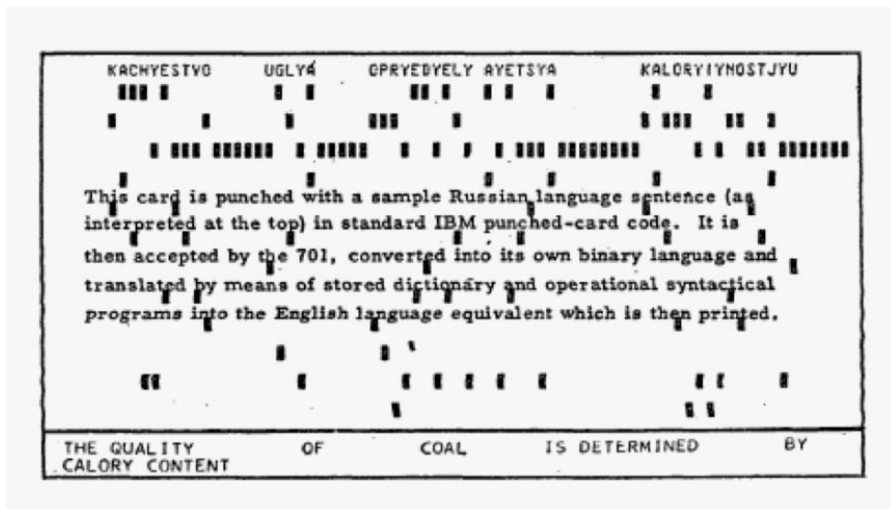
What is scientific research?

- General goal: Pushing the limits of our **knowledge**
- **Incrementally!** Find a limit/caveat in existing method and solve it!
- For example: lack of parallelization in Recurrent Neural Networks → **Transformers** (Vaswani et al. 2017)
- "Vaswani et al. 2017": a single publication that was **submitted** to a conference, **reviewed** by scientists, then **reproduced: research != science**
- Most of methods in this class were published less than 10 years ago

The shape of today's NLP research



NLP applications: Machine Translation



Georgetown-IBM experiment 1954

- Machine Translation is the first NLP application
- Google Translate supports 243 languages

Google Cloud Overview Solutions Products Pricing Resources

Cloud Translation

Model	Method	Usage	Price per unit
NMT	Text translations, which includes: <ul style="list-style-type: none"> • Language detection • Text translation • Batch text translation • XLISS document translation • Romanize text 	First 500,000 characters per month	Free (applied as \$10 credit every month)
		Over 500,000 characters per month	\$20 per million characters*
		Over 1 billion characters per month	We recommend that you contact a sales representative to discuss discount pricing.
	Document translation (DOCX, PPT, and PDF formats only)	Pages sent to the API per month	\$0.08 per page†

PRICING

- Cloud Translation pricing
- Pricing examples
- Charged characters
- Charged projects
- Other Google Cloud costs
- What's next

NLP applications: Machine Translation

Ubiquitous on the web and social media

The screenshot shows the New York Times website with a Google Translate widget in the top right corner. The widget is set to translate from 'anglais' to 'français'. Below the widget, there are three news articles:

- Le père d'un adolescent suspecté d'une fusillade dans une école de Géorgie**
Le père du suspect de 14 ans a été accusé de quatre chefs d'homicide involontaire, de deux chefs de meurtre au deuxième degré et de huit chefs de cruauté envers les enfants.
Voir plus de mises à jour
- Un rapport du bureau du shérif révèle des détails sur l'entretien du suspect avec le FBI et son activité en ligne.**
3 MINUTES DE LECTURE
- Le père du suspect de la fusillade en Géorgie a déclaré aux enquêteurs que lui et son fils avaient discuté de la sécurité des armes à feu l'année dernière.**
3 MINUTES DE LECTURE

Below the first article is a photo of a crowd of people. Below the second article is a photo of a person sitting in a car with colorful balloons. Below the third article is a photo of a person in a white shirt.



Fun new paper led by [@IngoZiegler](#) and [@akoksal](#) that shows how we can use retrieval augmentation to create high-quality supervised fine tuning data. All you need to do is write a few examples that demonstrate the task.

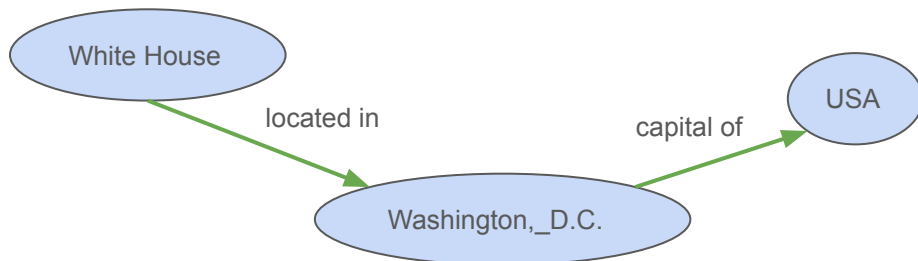
[À l'origine en anglais et traduit par Google](#)

Un nouvel article intéressant dirigé par [@IngoZiegler](#) et [@akoksal](#) qui montre comment nous pouvons utiliser l'augmentation de la récupération pour créer des données de réglage fin supervisées de haute qualité. Tout ce que vous avez à faire est d'écrire quelques exemples qui illustrent la tâche.

NLP applications: Information Extraction

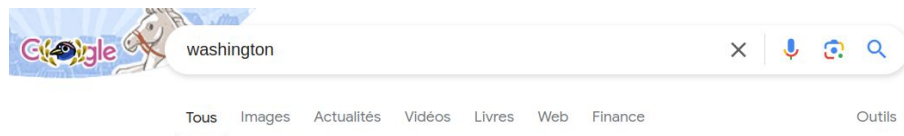
[Washington, D.C.](#) != [George Washington](#)

Washington is the capital of the USA. It hosts the White House.



- From unstructured text to knowledge graphs
- Named Entity Recognition
- Named Entity Disambiguation
- Coreference resolution
- Relation Extraction

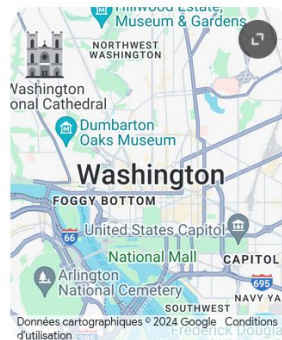
NLP applications: Information Extraction



Tightly connected with Information Retrieval

Washington

Capitale des États-Unis



Météo

sam.	dim.	lun.
☁ 23°	☀ 23°	☀ 26°

weather.com

Y aller

✈ 8h 25min
à partir de Paris

NLP applications: Information Extraction

Not only for advancing human knowledge

```

Page                               P235 Equity DES
SHAREHOLDER INFORMATION             Page 4 /10
THE WALT DISNEY CO.

DCACS CORPORATE ACTION CALENDAR
LATEST PUBLIC OFFERING
Date of offering      1/83
Shares offered       1,00M      Split Adj: 48.00M
Share Price          $ 66.88      Split Adj: 1.39
Lead Manager         Morgan Stanley
Type                 Common Stock

INSIDER TRADING                     INSTITUTIONAL OWNERSHIP
Net $ Value Buys and Sells As Of 01/15/05  # of Buyers 773
(1985 - Present In Dollars)                # of Sellers 788
Lowest activity 12/97 -372.00MLN           # of Holders 1,970
Highest activity 08/02 10.17MLN           Shares Held 1.39BLN
Mean: -15.95MLN % Shares Out. 68.12
Most recent 45 days .00 Shares Purchased 1.75MLN

Australia 61 2 9777 0600 Brazil 5511 3948 4500 Europe 44 20 7330 7200 Germany 49 69 920410
Hong Kong 852 2577 6000 Japan 81 3 3201 8900 Singapore 65 6212 1000 U.S. 1 212 318 2000 Copyright 2005 Bloomberg L.P.
1 24-Jan-05 15:10:57
  
```

```

<HELP> for explanation.                                dgp Equity HDS
Enter #<GO> to select aggregate portfolio and see detailed information
002723178194-000 HOLDINGS SEARCH CUSIP 25468710
DIS US THE WALT DISNEY CO. Page 1 / 100

```

Holder name	Portfolio Name	Source	Held	Outstd	Percent	Latest Filing	Change Date
*BARCLAYS GLOBAL	BARCLAYS BANK PLC	13F	91,394M	4,470	7,512M	09/04	
*CITIGROUP INCORP	CITIGROUP INCORPORAT	13F	71,012M	3,475	893,616	09/04	
*STATE STREET	STATE STREET CORPORA	13F	69,298M	3,365	1,214M	09/04	
*FIDELITY MANAG	FIDELITY MANAGEMENT	13F	67,611M	3,305	4,580M	09/04	
*SOUTHERSTAN ASST	SOUTHEASTERN ASSET M	13F	52,949M	2,591	3,194M	09/04	
*VANGUARD GROUP	VANGUARD GROUP INC	13F	49,710M	2,139	979,055	12/04	
*STATE FARM MUTUAL	STATE FARM MUTUAL AU	13F	42,224M	2,067	10,300	09/04	
*MELLON BANK N A	MELLON BANK CORP	13F	39,545M	1,935	2,999M	09/04	
*LORD ABBETT & CO	LORD ABBETT & CO	13F	37,460M	1,833	286,434	12/04	
*MORGAN STANLEY	MORGAN STANLEY	13F	31,643M	1,549	-1,948M	09/04	
*NORTHERN TRUST C	NORTHERN TRUST CORPO	13F	26,061M	1,275	-39,493	09/04	
*DEUTSCHE BANK AK	DEUTSCHE BANK AG	13F	21,990M	1,076	1,570M	09/04	
*DISNEY ROY EDUAR	n/a	Form 4	17,279M	0,846	-7,726	08/03	
*JANUS CAPITAL	JANUS CAPITAL CORPOR	13F	17,156M	0,840	-2,077M	09/04	
*ISCAPITAL RSCH MGN	CAPITAL RESEARCH AND	13F	17,142M	0,839	1,907M	09/04	
*TUKMAN CAP MNGT	TUKMAN CAPITAL MANAG	13F	16,962M	0,830	1,851M	09/04	
*T ROWE PRICE	T ROWE PRICE ASSOCIA	13F	16,810M	0,823	-1,367M	09/04	

```

Sub-totals for current page: 680,137M 33,286
* Money market directory info available. Select portfolio, then hit IP<GO>.
Australia 61 2 9777 0600 Brazil 5511 3948 4500 Europe 44 20 7330 7200 Germany 49 69 920410
Hong Kong 852 2577 6000 Japan 81 3 3201 8900 Singapore 65 6212 1000 U.S. 1 212 318 2000 Copyright 2005 Bloomberg L.P.
1 24-Jan-05 15:10:52
  
```

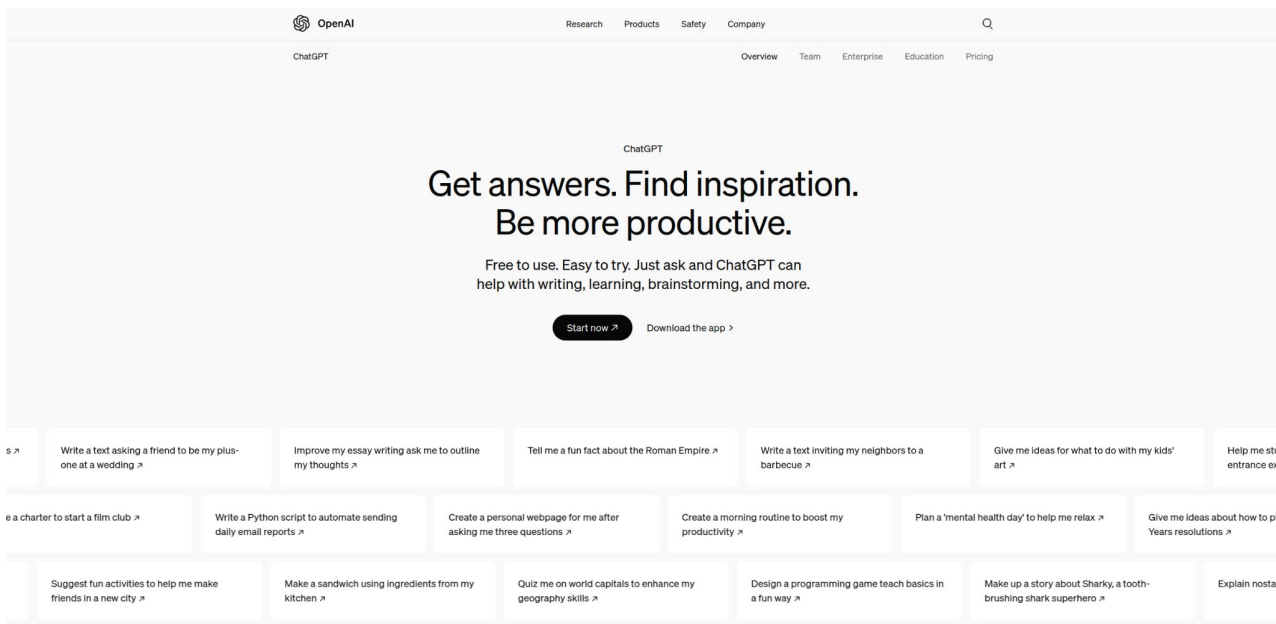
Bloomberg

NLP applications: chatbots



- Task-oriented like Siri/Alexa or chitchat like ChatGPT (OpenAI)
- Moving from a complicated pipeline (ASR, Information Extraction, Information Retrieval) to end-to-end language modeling

NLP applications: chatbots



NLP applications: chatbots: demo

32 lines of python code

```
1 from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig
2
3 def pipeline(text, model, tokenizer):
4     messages = [
5         {"role": "user", "content": text}
6     ]
7     inputs = tokenizer.apply_chat_template(messages, return_dict=True, tokenize=True, return_tensors="pt", add_generation_prompt=True)
8     for k, v in inputs.items():
9         inputs[k] = v.cuda()
10    input_length = inputs["input_ids"].shape[1]
11    output = model.generate(**inputs, max_new_tokens=128)
12    output = tokenizer.batch_decode(output[:, input_length:], skip_special_tokens=True, clean_up_tokenization_spaces=True)
13    return output
14
15
16 def user_loop(model, tokenizer):
17     while True:
18         answer = input(f">>> ").strip()
19         output = pipeline(answer, model, tokenizer)[0]
20         print(f"{output}\n")
21
22
23 def main():
24     model_name = "google/gemma-2-2b-it"
25     model = AutoModelForCausalLM.from_pretrained(
26         model_name,
27         device_map="auto",
28         token=True,
29         quantization_config=BitsAndBytesConfig(load_in_4bit=True)
30     )
31     tokenizer = AutoTokenizer.from_pretrained(model_name)
32     user_loop(model, tokenizer)
```

NLP applications: chatbots: demo

```
>>> what is distributional semantics?
```

```
## Distributional Semantics: Understanding Words in Context
```

```
Distributional semantics is a branch of natural language processing (NLP) that focuses on the relationships between words and their meaning based on their co-occurrence in a large corpus of text.
```

```
Think of it like this: instead of focusing on individual words and their exact definitions, distributional semantics analyzes how words are used in context.
```

```
Here's a breakdown:
```

```
* Words are not isolated: They are part of a larger network of words and phrases.
```

```
* Meaning is derived from context: The meaning of a word is determined by how it is
```

NLP applications: chatbots: demo

```
>>> I'm teaching an NLP class at the aivancity AI school, any advices?  
That's awesome! Teaching NLP is a fantastic way to share the power of AI with the next generation. Here's some  
advice to make your class engaging and impactful:  
  
**1. Know Your Audience:**  
  
* **Beginner or Advanced?** Are you targeting absolute beginners or those with some prior knowledge? Tailor yo  
ur content and pace accordingly.  
* **AI Enthusiasts vs. General Interest:** Are you aiming for a more technical audience or those with a broader  
interest in AI?  
* **Learning Style:** Consider how your students learn best (visual, auditory, kinesthetic). Incorporate divers  
e teaching methods.
```

NLP applications: chatbots: demo

```
>>> you're speaking live to students of the aivancity AI school, introduce yourself  
Hello everyone! 🙋 I'm Gemma, an AI assistant here to help you learn and explore the fascinating world of AI.
```

```
I'm excited to be speaking with you today at Aivancity AI School. I've been trained on a massive dataset of text and code, which allows me to understand and generate human-like text.
```

```
What are you working on today? Do you have any questions about AI concepts, or are you working on a specific project? I'm here to help!
```

```
Let's dive into the world of AI together! 🤖
```


Big Money in LLMs/Chatbots

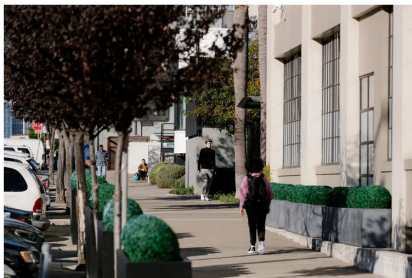
The New York Times

OpenAI Tries to Grow Up Changing a Chatbot's Mind A.I.'s Threat to Itself 'Deepfake' Elon Musk

OpenAI Completes Deal That Values the Company at \$80 Billion

The A.I. start-up's valuation tripled in less than 10 months.

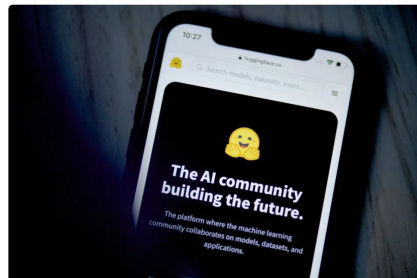
Share full article



Outside OpenAI's offices in San Francisco. The company's latest deal is another example of the Silicon Valley deal-making machine pumping money into a handful of companies that specialize in generative A.I. Jason Henry for The New York Times

Bloomberg

AI Startup Hugging Face Valued at \$4.5 Billion After Raising Funding From Google, Nvidia



AI Startup Hugging Face Valued at \$4.5 Billion After Raising Funding From Google, Nvidia - Bloomberg

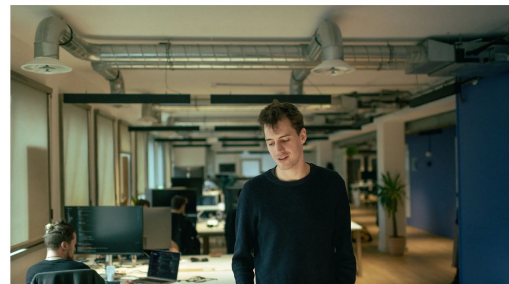
The New York Times

Artificial Intelligence OpenAI Tries to Grow Up Changing a Chatbot's Mind A.I.'s Threat to Itself 'Deepfake' Elon Musk Q&A: Fake or Real Images?

Mistral, a French A.I. Start-Up, Is Valued at \$6.2 Billion

Created by alumni from Meta and Google, Mistral is just a year old and has already raised more than \$1 billion in total from investors, leading to eye-popping valuations.

Listen to this article 3:28 min Show full article



Big Money in LLMs/Chatbots

Zuckerberg's Meta Is Spending Billions to Buy 350,000 Nvidia H100 GPUs

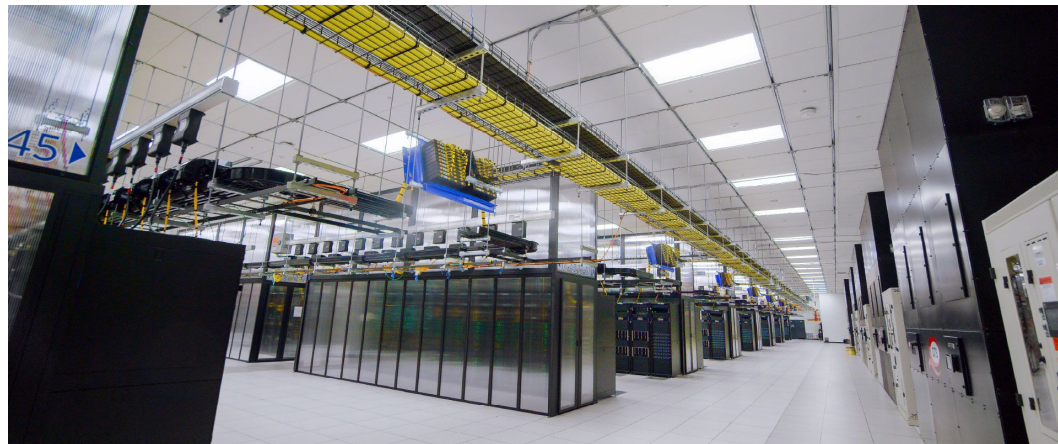
In total, Meta will have the compute power equivalent to 600,000 Nvidia H100 GPUs to help it develop next-generation AI, says CEO Mark Zuckerberg.



By Michael Kan January 18, 2024



(David Paul Morris/Bloomberg via Getty Images)



Jean Zay cluster

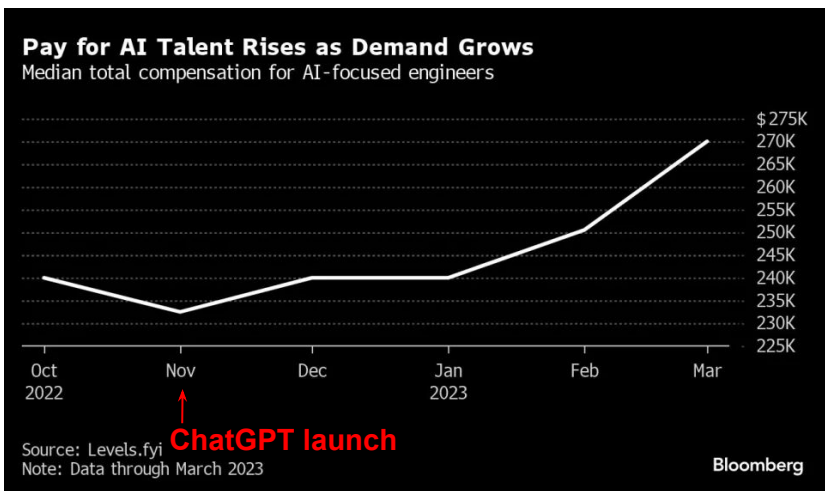
```
(matos) [us147j@jean-zay1: experiments]$ sinfo
PARTITION AVAIL  TIMELIMIT  NODES  STATE MODEL
cpu_p1*      up 4-04:00:00    2 drain$ r313n[8-9]
cpu_p1*      up 4-04:00:00    1 maint r21n16
cpu_p1*      up 4-04:00:00    1 drain* r217n9
cpu_p1*      up 4-04:00:00    4 mxx r13n12,r310n[4,25,34]
cpu_p1*      up 4-04:00:00    711 alloc r110n[0-35],r111n[0-35],r112n[0-35],r113n[0-11,13-35],r114n[0-35],r115
n[0-35],r116n[0-35],r117n[0-35],r210n[0-35],r211n[0-15,17-35],r212n[0-35],r213n[0-35],r214n[0-35],r215n[0-35],r
216n[0-35],r217n[0-8,10-35],r310n[0-3,5-24,26-33,35],r311n[0-30,32-35],r312n[0-35],r313n[0-7,10-35]
cpu_p1*      up 4-04:00:00    1 tdlc r31n31
gpu_p13      up 4-04:00:00    1 drain$ r314n8
gpu_p13      up 4-04:00:00    3 drain* r610n,r815n[2,5]
gpu_p13      up 4-04:00:00    1 drain r810n0
gpu_p13      up 4-04:00:00    99 mxx r314n[1-2],r315n[0-4],r316n[2,4,6-7],r317n8,r610n7,r611n[2,4,6],r612n[
0-2-5],r613n[1,6,8],r614n[0-1,3-5,7-8],r615n[1,3,5,8],r616n[0-7],r617n7,r70n[2-4],r71n5,r712n[1,4-5],r713n7,r
714n[0-1,5,8],r715n8,r716n[2,4,7-8],r717n8,r810n1,r811n5,r812n[2-3],r813n6,r814n8,r815n[3,8],r817n[1,5,7],r910n
4,r911n[0,2,7],r912n[0,5,7],r913n[2,4],r914n[1,5,8],r915n[6,8],r917n[0,4-5],r1010n6,r1011n5,r1012n[0,2,5],r1013
n0,r1015n[3-4],r1016n[4,7-8],r1017n[0,4-6]
gpu_p13      up 4-04:00:00    292 alloc r314n[0-3-7],r315n[5-8],r316n[0-1,3,5,8],r317n[0-7],r610n[1-6,8],r611n
[0-1,3,5,7-8],r612n[1,6-8],r613n[0,2-5,7],r614n[2,6],r615n[0,2,4,6-7],r616n[1-6,8],r617n[0-6,8],r710n[0-1-5-8],
r711n[0-4,6-8],r712n[0-2,3-6-8],r713n[0-6,8],r714n[2-4,6-7],r715n[0-2,4-8],r716n[0-1-3,5-6],r717n[1-8],r810n[2-
3,8],r811n[0-2,4-8],r812n[0-1,4-8],r813n[0-5,7-8],r814n[0-7],r815n[0-1,4,6-7],r816n[0-8],r817n[0,2-4,6,8],r910n[0
-3,5-8],r911n[1,3-6,8],r912n[1-4,6,8],r913n[0-1,3,5-8],r914n[0,2-4,6-7],r915n[0-5,7],r916n[0-8],r917n[1-3,6-8],
r1010n[0-5,7-8],r1011n[0-4,6-8],r1012n[1,3-4,6-8],r1013n[1-8],r1014n[0-8],r1015n[0-2,5-8],r1016n[0-3,5-6],r1017
n[1-5,7-8]
gpu_p2       up 4-04:00:00    1 drain$ jean-zay-l0810
gpu_p2       up 4-04:00:00    17 mxx jean-zay-la[802,805-807,809,812,815-816,820,822,824-827,829-831]
gpu_p2       up 4-04:00:00    12 alloc jean-zay-la[801,803-804,811,813-814,817-819,821,823,828]
gpu_p2       up 4-04:00:00    1 tdlc jean-zay-l0808
gpu_p21      up 4-04:00:00    1 drain$ jean-zay-l0810
gpu_p21      up 4-04:00:00    5 mxx jean-zay-la[802,805-807,809]
gpu_p21      up 4-04:00:00    4 alloc jean-zay-la[801,803-804,811]
gpu_p21      up 4-04:00:00    1 tdlc jean-zay-l0808
gpu_p25      up 4-04:00:00    12 mxx jean-zay-la[812,815-816,820,822,824-827,829-831]
gpu_p25      up 4-04:00:00    8 alloc jean-zay-la[813-814,817-819,821,823,828]
gpu_p5       up 4-04:00:00    38 mxx jean-zay-lan[01,03-13,15-18,20,22-26,28-29,33,35-36,38-42,44,46-48,51-
52]
gpu_p5       up 4-04:00:00    14 alloc jean-zay-lan[02,14,19,21,27,30-32,34,37,43,45,49-50]
visu         up 4:00:00      1 drain* jean-zay-visu5
visu         up 4:00:00      1 resv jean-zay-visu4
visu         up 4:00:00      1 nix jean-zay-visu1
visu         up 4:00:00      2 tdlc jean-zay-visu[2-3]
prepost      up 20:00:00     1 drain* jean-zay-pp1
prepost      up 20:00:00     3 nix jean-zay-pp[2-4]
archive      up 20:00:00     3 comp ldrsrv[06-08]
archive      up 20:00:00     1 resv ldrsrv05
compt11      up 20:00:00     1 drain* jean-zay-pp1
compt11      up 20:00:00     3 comp ldrsrv[06-08]
compt11      up 20:00:00     1 resv ldrsrv05
compt11      up 20:00:00     3 nix jean-zay-pp[2-4]
```

```
(matos) bash-5.1$ nvidia-smi
Fri Sep  6 12:04:10 2024
```

NVIDIA-SMI 550.54.15		Driver Version: 550.54.15		CUDA Version: 12.4	
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile Uncorr. ECC
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util Compute M.
				MIG M.	
0	Tesla V100-SXM2-32GB	On	00000000:1A:00:0	Off	0
N/A	44C P0	45W / 300W	0MiB / 32768MiB	0%	Default
					N/A



Big Money in LLMs/Chatbots



Nvidia

107,21 \$ ↑ 2 298,43 % +102,74 5 a

Avant l'ouverture : 104,90 \$ (↓ 2,15 %) -2,31

Fermé : 6 sept., 05:35:32 UTC-4 · USD · NASDAQ · Clause de non-responsabilité

1 j 5 j 1 m 6 m YTD 1 a 5 a MAX



+568% in less than 2 years

Break for questions and "appel"

"calls the roll"? (good evidence that Machine Translation is not a solved problem)

français (langue détectée) ▾	↔	anglais (américain) ▾
l'appel ×		the call
français (langue détectée) ▾	↔	anglais (américain) ▾
faire l'appel ×		make the call
français (langue détectée) ▾	↔	anglais (américain) ▾
le professeur fait l'appel ×		the teacher calls the roll

Dictionnaire	
faire appel verbe ⌵	
appeal v ⌵ (appealed, appealed)	
L'avocat a fait appel de la condamnation de son client.	The lawyer appealed his client's sentence.

What is a word?

- Open question in phonology vs. morphology
- Inflection: is brother != brothers?
- Compounding: is motorbike == motor + bike?
- Multi-word expressions: mother in law == 1 or 3 words?
- Polysemy: is chair (furniture) != chair (person)?
- Orthography: is modeling != modelling?

NLP deals with orthographic words...

- "My brother is sitting on a chair" →
['My', 'brother', 'is', 'sitting', 'on', 'a', 'chair']
(**tokenization**: sequence of *tokens*)
- Inflection (brother vs brothers): usually not modeled
- Compounding (motorbike vs motor + bike): usually not modeled
- Multi-word expressions (mother in law): usually not modeled
- Polysemy (chair [furniture] vs chair [person]): usually modeled after sharing an initial representation
- **Orthography**: 'modeling' != 'mode11ing'

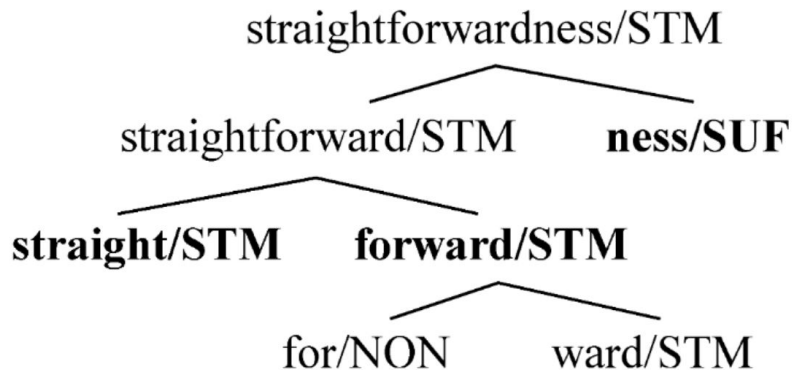
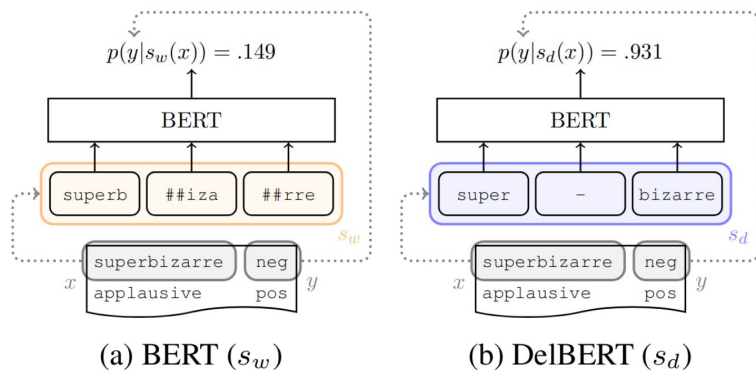
...except when it's the research topic!

Inflection (brother vs brothers): "brother" is a *lemma* (singular, masc. form): useful for **indexing** (keyword-like) in Information Retrieval

Category	Infl.	Deri.	Comp.	Description	English example (input ==> output)
000	-	-	-	Root words (free morphemes)	progress ==> progress
100	✓	-	-	Inflection only	prepared ==> prepare @@ed
010	-	✓	-	Derivation only	intensive ==> intense @ive
001	-	-	✓	Compound only	hotpot ==> hot @@pot
101	✓	-	✓	Inflection and Compound	wheelbands ==> wheel @@band @@s
011	-	✓	✓	Derivation and Compound	tankbuster ==> tank @@bust @@er
110	✓	✓	-	Inflection and Derivation	urbanizes ==> urban @@ize @@s
111	✓	✓	✓	Inflection, Derivation, Compound	trackworkers ==> track @@work @@er @@s

...except when it's the research topic!

Compounding (motorbike vs motor + bike): very niche but studied in computational linguistics

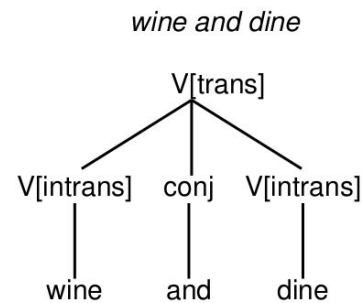
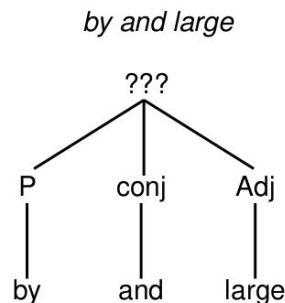


...except when it's the research topic!

Multi-word expressions ("mother in law")

Pointwise mutual information (PMI):

$$\log \frac{P(x,y)}{P(x)P(y)}$$



...except when it's the research topic!

Orthography ('modeling' ← 'mode 11ing') for **User-Generated Content**

i left ACL cus im sickk ! Yuu better be their tmrw . GN 4now
 ↓ ↓ ↓ ↓ ↓
i left ACL because I'm sick ! You better be their tomorrow . GN 4now

What do words mean?

- Why is it "brother" in English and "frère" in French?
- Because "brōþēr" in Proto-Germanic and "frātrēm" in Latin!
(arbitrariness of the sign, de Saussure, 1916)
But why does it *mean* brother?
- The meaning of a word is its **use** in the language (Ludwig Wittgenstein, 1921):
"I was playing with my **brother** and *sister*"
"My *mom* is feeding my **brother**"
- "brother" co-occurs with "mom" and "sister"
like "frère" co-occurs with "maman" and "sœur"
- Polysemy: "I sit on a *chair*" vs "He is the *chair* of this session"

How words are used?

- words are defined by their environments (the words around them)
- If A and B have almost identical environments we say that they are **synonyms** (Harris, 1954).
- define the meaning of a word by its distribution in language use: its neighboring words

What does "ongchoi" mean?

- Suppose you see these sentences:
 - *Ongchoi* is delicious **sautéed with garlic**.
 - *Ongchoi* is superb over **rice**
 - *Ongchoi* **leaves** with **salty** sauces
- And you've also seen these:
 - ...*spinach* **sautéed with garlic** over **rice**
 - *Chard* stems and **leaves** are **delicious**
 - *Collard greens* and other **salty** leafy greens
- *Ongchoi* is a leafy green like *spinach*, *chard*, or *collard greens*



Defining context (word-word matrix)

Two words are similar in meaning if their context vectors are similar

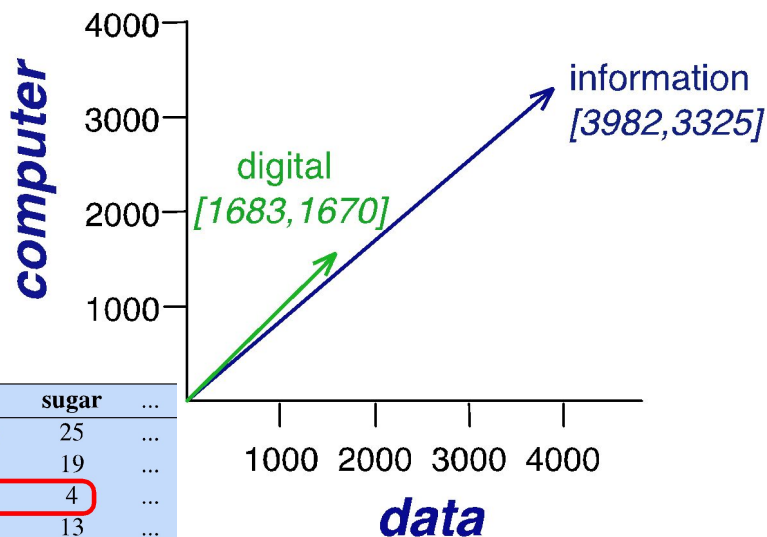
is traditionally followed by **cherry** pie, a traditional dessert
 often mixed, such as **strawberry** rhubarb pie. Apple pie
 computer peripherals and personal **digital** assistants. These devices usually
 a computer. This includes **information** available on the internet

	aardvark	...	computer	data	result	pie	sugar	...
cherry	0	...	2	8	9	442	25	...
strawberry	0	...	0	0	1	60	19	...
digital	0	...	1670	1683	85	5	4	...
information	0	...	3325	3982	378	5	13	...

Defining context (word-word matrix)

Two words are similar in meaning if their context vectors are similar

	aardvark	...	computer	data	result	pie	sugar	...
cherry	0	...	2	8	9	442	25	...
strawberry	0	...	0	0	1	60	19	...
digital	0	...	1670	1683	85	5	4	...
information	0	...	3325	3982	378	5	13	...



Computing word similarity: Dot product

The dot product between two vectors is a scalar:

$$\text{dot product}(\mathbf{v}, \mathbf{w}) = \mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^N v_i w_i = v_1 w_1 + v_2 w_2 + \dots + v_N w_N$$

The dot product tends to be high when the two vectors have large values in the same dimensions

Dot product can thus be a useful similarity metric between vectors

Problem with raw dot-product

Dot product favors long vectors

Dot product is higher if a vector is longer (has higher values in many dimension)

Vector length (euclidean norm):

$$|\mathbf{v}| = \sqrt{\sum_{i=1}^N v_i^2}$$

Frequent words (of, the, you) have long vectors (since they occur many times with other words).

So dot product overly favors frequent words

Alternative: cosine for word similarity

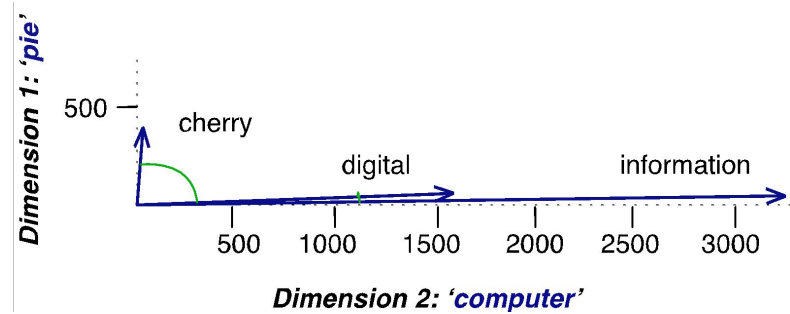
$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

Based on the definition of the dot product between two vectors **a** and **b**

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \theta$$
$$\frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| |\mathbf{b}|} = \cos \theta$$

Cosine examples

	pie	data	computer
cherry	442	8	2
digital	5	1683	1670
information	5	3982	3325



$$\cos(\text{cherry}, \text{information}) = \frac{442 * 5 + 8 * 3982 + 2 * 3325}{\sqrt{442^2 + 8^2 + 2^2} \sqrt{5^2 + 3982^2 + 3325^2}} = .017$$

$$\cos(\text{digital}, \text{information}) = \frac{5 * 5 + 1683 * 3982 + 1670 * 3325}{\sqrt{5^2 + 1683^2 + 1670^2} \sqrt{5^2 + 3982^2 + 3325^2}} = .996$$

Can we compute word similarity like this?

← V vocabulary size →

	aardvark	...	computer	data	result	pie	sugar	...
cherry	0	...	2	8	9	442	25	...
strawberry	0	...	0	0	1	60	19	...

- **Sparse** vectors (most words never co-occur together)
- Very **high dimension!** V : vocabulary size (usually 20,000 - 200,000)

How do we reduce dimensionality?

← from V (vocabulary size) to $d \ll V$ →

	aardvark	...	computer	data	result	pie	sugar	...
cherry	0	...	2	8	9	442	25	...
strawberry	0	...	0	0	1	60	19	...

- Generic solutions:
 - Principal Component Analysis (PCA)
 - Singular Value Decomposition (SVD) → Latent Semantic Indexing/Analysis (Deerwester et al., 1990)
- Deep learning solution: Skipgram (word2vec, Mikolov 2013)

Latent Semantic Indexing/Analysis

$$\mathbf{A} \approx \hat{\mathbf{A}} = \mathbf{M} \times \text{diag}(\mathbf{s}) \times \mathbf{C}^T$$

$V \times C$ $V \times d$ $d \times d$ $d \times C$

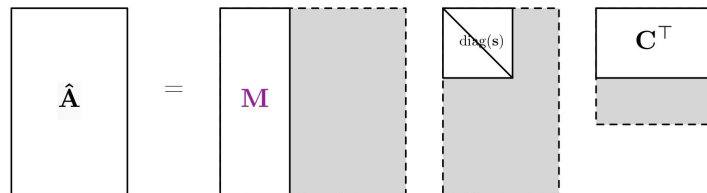
Singular Value Decomposition (SVD)

- Usually done with word-document occurrences instead of word-word
- Actually Pointwise Mutual Information instead of raw counting
- Closely related to Skipgram (Levy and Goldberg, 2014)

SVD:



truncated at d :



Skipgram (word2vec, Mikolov)

- Instead of **counting** how often each word **w** occurs near "*apricot*"
Train a classifier on a binary prediction task: Is **w** likely to show up near "*apricot*"?
- We don't actually care about this task
But we'll take the learned classifier weights as the word embeddings
- Big idea: **self-supervision**:
 - A word **c** that occurs near *apricot* in the corpus acts as the gold "correct answer" for supervised learning
 - No need for human labels

Skipgram (word2vec, Mikolov)

- Treat the target word w and a neighboring context word c as positive examples.
- Randomly sample other words in the lexicon to get negative examples
- Use logistic regression to train a classifier to distinguish those two cases
- Use the learned weights as the embeddings

Skipgram (word2vec, Mikolov)

Assume a +/- 2 word window, given training sentence:

...lemon, a tablespoon of apricot jam, a pinch...

Goal: train a classifier that is given a candidate (word, context) pair

(apricot, jam)

(apricot, aardvark)

And assigns each pair a probability:

$$P(+|w, c)$$

$$P(-|w, c) = 1 - P(+|w, c)$$

Turning dot products into probabilities

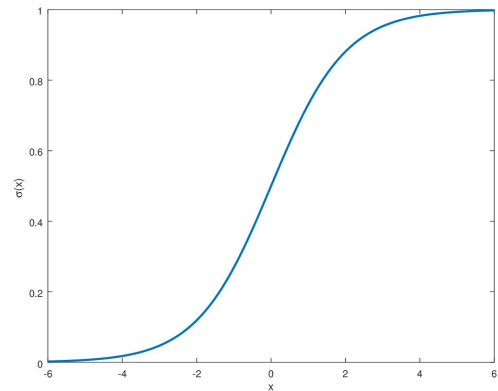
$$\text{Sim}(w, c) \approx w \cdot c$$

To turn this into a probability

We'll use the sigmoid from logistic regression:

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

$$\begin{aligned} P(-|w, c) &= 1 - P(+|w, c) \\ &= \sigma(-c \cdot w) = \frac{1}{1 + \exp(c \cdot w)} \end{aligned}$$



From 1 context word to full context

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

Assume all context words are **independent** → joint probability = product

$$P(+|w, c_{1:L}) = \prod_{i=1}^L \sigma(c_i \cdot w)$$

$$\log P(+|w, c_{1:L}) = \sum_{i=1}^L \log \sigma(c_i \cdot w)$$

log Prob: **systematic** trick for **numerical stability**

Skip-Gram Training data

...lemon, a tablespoon of apricot jam, a pinch...

positive examples +

t	c
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

negative examples -

t	c	t	c
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

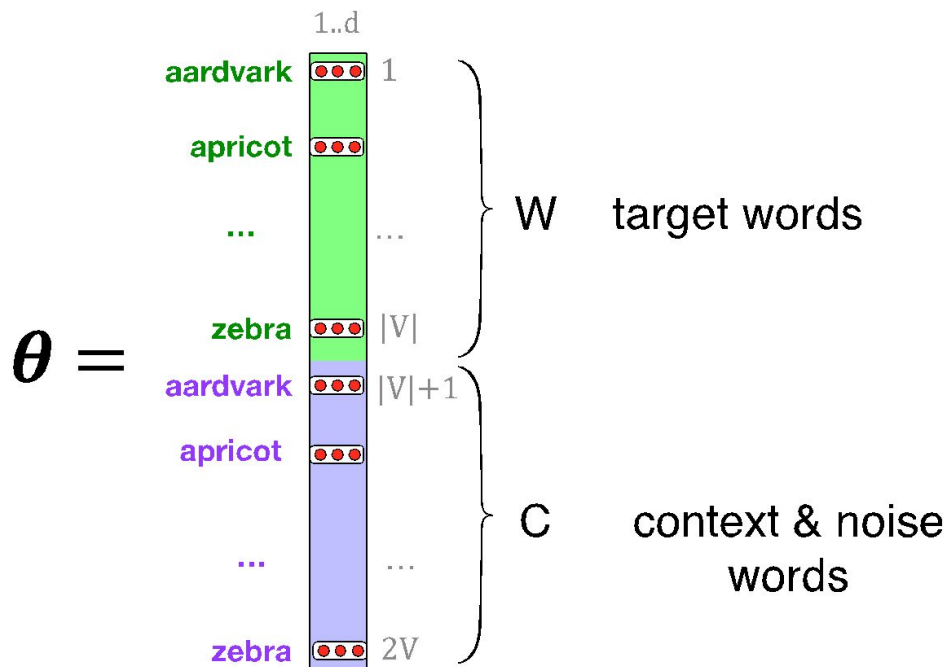
- Maximize the similarity of the target word, context word pairs $(w, c+)$ drawn from the positive data
- Minimize the similarity of the $(w, c-)$ pairs drawn from the negative data.

Loss function for one w

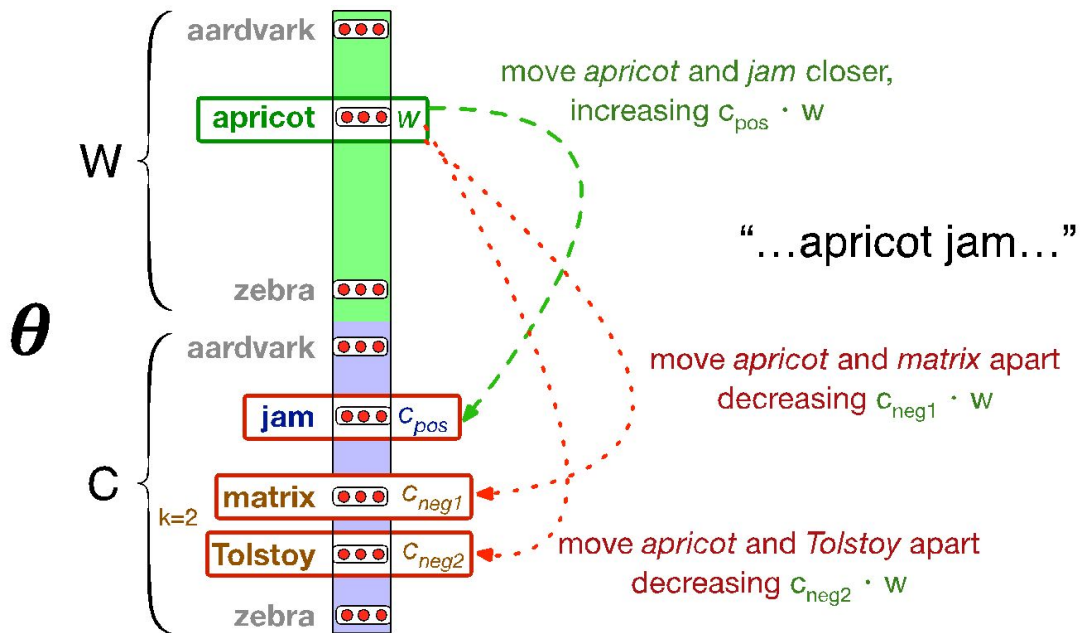
- Maximize the similarity of the target word, context word pairs $(w, c+)$ drawn from the positive data
- Minimize the similarity of the $(w, c-)$ pairs drawn from the negative data.

$$\begin{aligned}
 L_{CE} &= -\log \left[P(+|w, c_{pos}) \prod_{i=1}^k P(-|w, c_{neg_i}) \right] \\
 &= - \left[\log P(+|w, c_{pos}) + \sum_{i=1}^k \log P(-|w, c_{neg_i}) \right] \\
 &= - \left[\log P(+|w, c_{pos}) + \sum_{i=1}^k \log (1 - P(+|w, c_{neg_i})) \right] \\
 &= - \left[\log \sigma(c_{pos} \cdot w) + \sum_{i=1}^k \log \sigma(-c_{neg_i} \cdot w) \right]
 \end{aligned}$$

Learning with Stochastic gradient descent

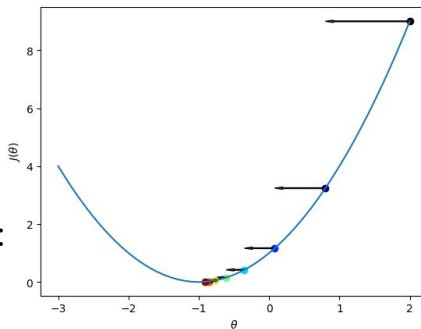


Learning with Stochastic gradient descent



Stochastic gradient descent (SGD) reminder

- Learning rate $\alpha \in \mathbb{R}, \alpha > 0$
- Randomly initialize $\theta^{(0)}$
- Iteratively get better estimate with:



Gradient is:

- the vector of partial derivatives of the parameters with respect to the loss function
- A linear approximation of the loss function at $\theta^{(i)}$

$$\frac{\partial L}{\partial \theta}(\theta^{(i)}) = \begin{bmatrix} \frac{\partial L}{\partial \theta_1^{(i)}} \\ \frac{\partial L}{\partial \theta_2^{(i)}} \\ \vdots \\ \frac{\partial L}{\partial \theta_n^{(i)}} \end{bmatrix}$$

Next estimate $\theta^{(i+1)} = \theta^{(i)} - \alpha * \frac{\partial L}{\partial \theta}(\theta^{(i)})$

Learning rate (step size) α

Previous Estimate $\theta^{(i)}$

The derivatives of the loss function

$$L_{CE} = - \left[\log \sigma(c_{pos} \cdot w) + \sum_{i=1}^k \log \sigma(-c_{neg_i} \cdot w) \right]$$

$$\frac{\partial L_{CE}}{\partial c_{pos}} = [\sigma(c_{pos} \cdot w) - 1]w$$

$$\frac{\partial L_{CE}}{\partial c_{neg}} = [\sigma(c_{neg} \cdot w)]w$$

$$\frac{\partial L_{CE}}{\partial w} = [\sigma(c_{pos} \cdot w) - 1]c_{pos} + \sum_{i=1}^k [\sigma(c_{neg_i} \cdot w)]c_{neg_i}$$

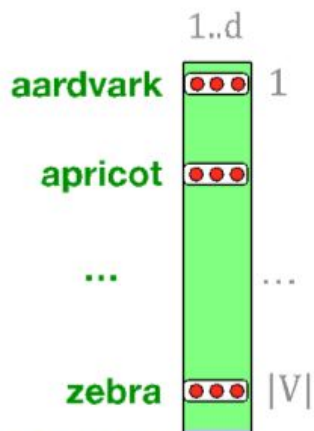
Stochastic gradient descent update

$$c_{pos}^{t+1} = c_{pos}^t - \eta [\sigma(c_{pos}^t \cdot w^t) - 1] w^t$$

$$c_{neg}^{t+1} = c_{neg}^t - \eta [\sigma(c_{neg}^t \cdot w^t)] w^t$$

$$w^{t+1} = w^t - \eta \left[[\sigma(c_{pos} \cdot w^t) - 1] c_{pos} + \sum_{i=1}^k [\sigma(c_{neg_i} \cdot w^t)] c_{neg_i} \right]$$

Embedding = lookup table or linear layer?



lookup table

One-hot encoding

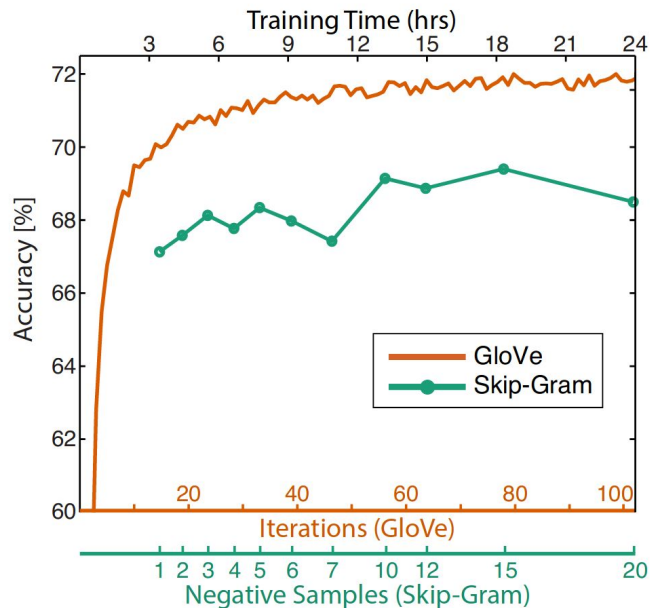
$$\text{Standard basis of } \mathbb{R}^n : e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, e_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, e_n = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

$$|\mathcal{V}| = n:$$

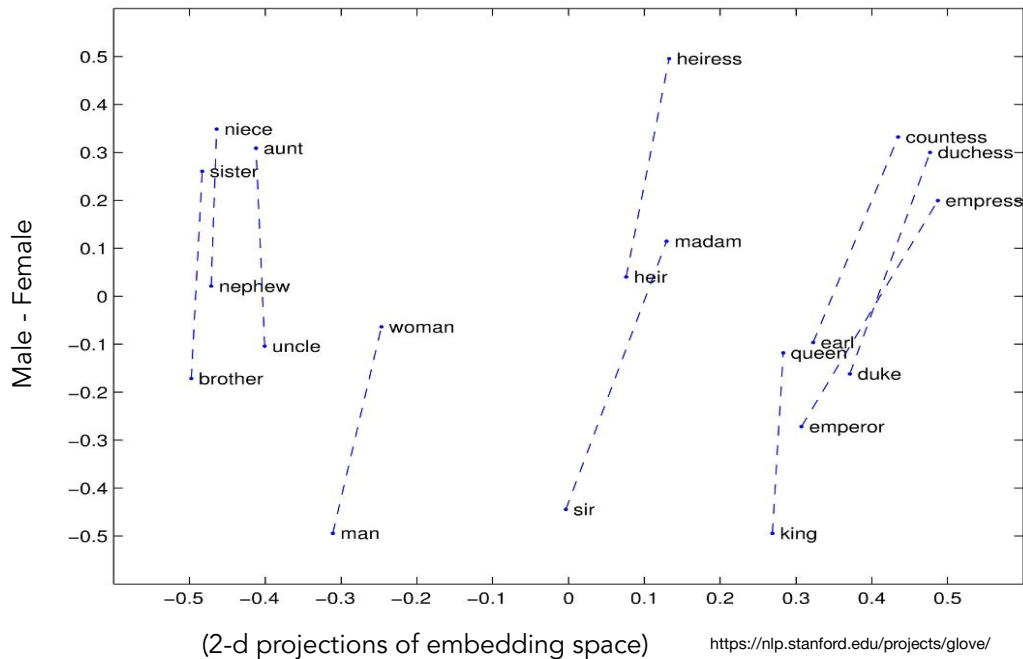
$$\text{features}(v_i) = W e_i \in \mathbb{R}^k = \text{ith column of } W$$

As always, hyperparameters

- Vocabulary size V
- Context window C
- Number of negative examples k
- Embedding dimension d
- The usual:
 - learning rate etc.
- → **Empirical** evaluation!



What now?

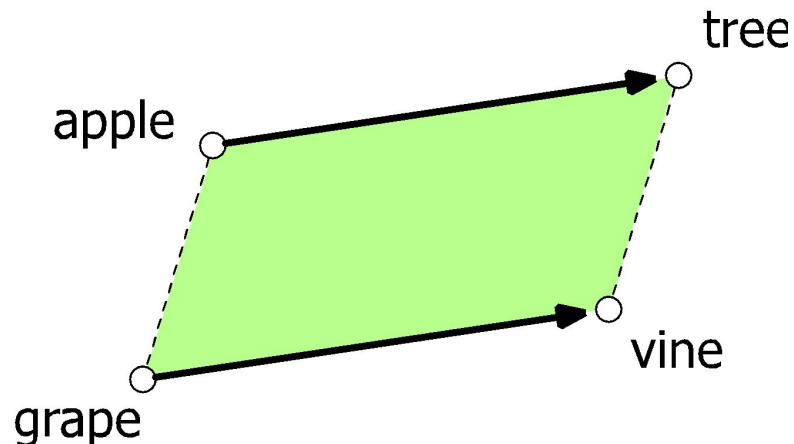


Intrinsic evaluation

- Do (cosine) similarities of pairs of words' vectors correlate with judgments of similarity by humans?
- TOEFL-like synonym tests, e.g., rug → {sofa **X** ottoman **X** carpet **✓** hallway **X**}
- analogies:
 - syntactic
 - semantic

Analogical relations

- The classic parallelogram model of analogical reasoning (Rumelhart and Abrahamson 1973)
- To solve: "apple is to tree as grape is to _____"
- Add tree – apple to grape to get vine
- Syntactic analogies, e.g., "walking is to walked as eating is to what?" Solved via:



$$\max_{v \in \mathcal{V}} \cos(\mathbf{v}_v, -\mathbf{v}_{\text{walking}} + \mathbf{v}_{\text{walked}} + \mathbf{v}_{\text{eating}})$$

Quantitatively

WS353 (WORDSIM) [13]			MEN (WORDSIM) [4]		
Representation		Corr.	Representation		Corr.
SVD	(k=5)	0.691	SVD	(k=1)	0.735
SPPMI	(k=15)	0.687	SVD	(k=5)	0.734
SPPMI	(k=5)	0.670	SPPMI	(k=5)	0.721
SGNS	(k=15)	0.666	SPPMI	(k=15)	0.719
SVD	(k=15)	0.661	SGNS	(k=15)	0.716
SVD	(k=1)	0.652	SGNS	(k=5)	0.708
SGNS	(k=5)	0.644	SVD	(k=15)	0.694
SGNS	(k=1)	0.633	SGNS	(k=1)	0.690
SPPMI	(k=1)	0.605	SPPMI	(k=1)	0.688

Spearman's ρ k is the number of “negative” samples

MEN : 3000 items

a	b	label
sun	sunlight	50.0
automobile	car	50.0
river	water	49.0
stairs	staircase	49.0
morning	sunrise	49.0
...
feathers	truck	1.0
festival	whiskers	1.0
muscle	tulip	1.0
bikini	pizza	1.0
bakery	zebra	0.0

Quantitatively

MIXED ANALOGIES [20]			SYNT. ANALOGIES [22]		
Representation		Acc.	Representation		Acc.
SPPMI	(k=1)	0.655	SGNS	(k=15)	0.627
SPPMI	(k=5)	0.644	SGNS	(k=5)	0.619
SGNS	(k=15)	0.619	SGNS	(k=1)	0.59
SGNS	(k=5)	0.616	SPPMI	(k=5)	0.466
SPPMI	(k=15)	0.571	SVD	(k=1)	0.448
SVD	(k=1)	0.567	SPPMI	(k=1)	0.445
SGNS	(k=1)	0.540	SPPMI	(k=15)	0.353
SVD	(k=5)	0.472	SVD	(k=5)	0.337
SVD	(k=15)	0.341	SVD	(k=15)	0.208

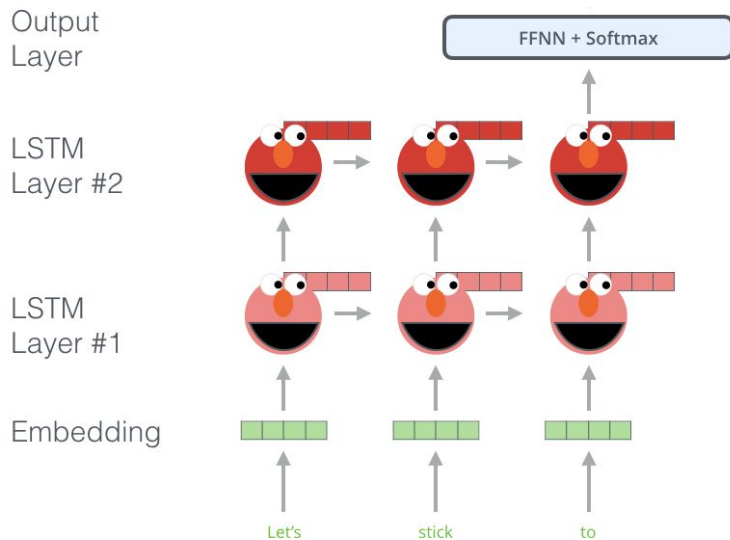
k is the number of “negative” samples

Word Pair 1		Word Pair 2	
Athens	Greece	Oslo	Norway
Astana	Kazakhstan	Harare	Zimbabwe
Angola	kwanza	Iran	rial
Chicago	Illinois	Stockton	California
brother	sister	grandson	granddaughter

Example
good:better rough:___
good:best rough:___
better:best rougher:___
year:years law:___
city:city's bank:___
see:saw return:___
see:sees return:___
saw:sees returned:___

Extrinsic evaluation

- Embeddings are the first brick of any more complex models (described in next class)
- Embeddings can be initialized with Skip-gram: **pretraining**/transfer learning
- either keep them **frozen** or **fine-tune** them



Named Entity Recognition with pretrained embeddings

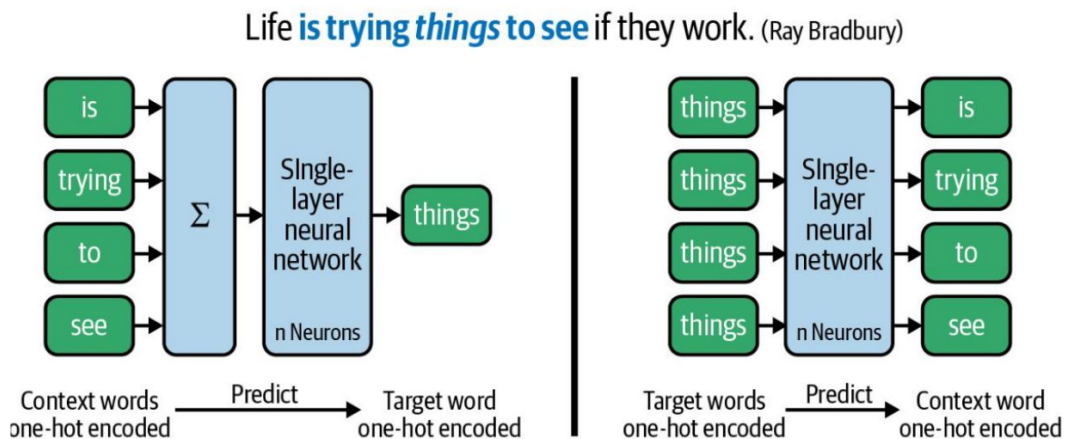
Washington is the capital of the USA. It hosts the White House.

Model	Dev	Test	ACE	MUC7
Discrete	91.0	85.4	77.4	73.4
SVD	90.8	85.7	77.3	73.7
SVD-S	91.0	85.5	77.6	74.3
SVD-L	90.5	84.8	73.6	71.5
HPCA	92.6	88.7	81.7	80.7
HSMN	90.5	85.7	78.7	74.7
CW	92.2	87.4	81.7	80.2
CBOW	93.1	88.2	82.2	81.1
GloVe	93.2	88.3	82.9	82.2

F1 score

Alternatives to Skipgram: continuous bag of words (CBOW)

instead of predicting context from word, predict word from context (much like a language model)



Alternatives to Skipgram: continuous bag of words (CBOW)

"bag of words" because does not model word order, puts all words in the same "bag"

$$\bar{\mathbf{v}}_m = \frac{1}{2h} \sum_{n=1}^h \mathbf{v}_{w_{m+n}} + \mathbf{v}_{w_{m-n}}$$

average of embeddings for words in the immediate neighborhood (**m-h**, ... , **m+h**)

\mathbf{x}_1 : yes , we have no bananas

\mathbf{x}_2 : say yes for bananas

\mathbf{x}_3 : no bananas , we say

	1	2	3
,	1	0	1
bananas	1	1	1
for	0	1	0
have	1	0	0
no	1	0	1
say	0	1	1
we	1	0	1
yes	1	1	0

Alternatives to Skipgram: continuous bag of words (CBOW)

$$\begin{aligned}\log p(\mathbf{w}) &\approx \sum_{m=1}^M \log p(w_m \mid w_{m-h}, w_{m-h+1}, \dots, w_{m+h-1}, w_{m+h}) \\ &= \sum_{m=1}^M \log \frac{\exp(\mathbf{u}_{w_m} \cdot \bar{\mathbf{v}}_m)}{\sum_{j=1}^V \exp(\mathbf{u}_j \cdot \bar{\mathbf{v}}_m)} \\ &= \sum_{m=1}^M \mathbf{u}_{w_m} \cdot \bar{\mathbf{v}}_m - \log \sum_{j=1}^V \exp(\mathbf{u}_j \cdot \bar{\mathbf{v}}_m).\end{aligned}$$

Empirical comparison

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
CBOW[†]	6B	57.2	65.6	68.2	57.0	32.5
SG[†]	6B	62.8	65.2	69.7	<u>58.1</u>	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	<u>75.9</u>	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<u>47.8</u>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

Spearman's ρ

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW[†]	300	6B	63.6	<u>67.4</u>	65.7
SG[†]	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>

Word analogy

Alternatives to Skipgram: GloVe

studies ratio of co-occurrence instead of co-occurrence

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

$$\begin{aligned}
 \min_{\mathbf{u}, \mathbf{v}, b, \tilde{b}} \quad & \sum_{j=1}^V \sum_{j \in \mathcal{C}} f(M_{ij}) \left(\widehat{\log M_{ij}} - \log M_{ij} \right)^2 \\
 \text{s.t.} \quad & \widehat{\log M_{ij}} = \mathbf{u}_i \cdot \mathbf{v}_j + b_i + \tilde{b}_j,
 \end{aligned}$$

log count(i, j)

Empirical comparison

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
CBOW[†]	6B	57.2	65.6	68.2	57.0	32.5
SG[†]	6B	62.8	65.2	69.7	<u>58.1</u>	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	<u>75.9</u>	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<u>47.8</u>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

Spearman's ρ

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW[†]	300	6B	63.6	<u>67.4</u>	65.7
SG[†]	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>

Word analogy

Skipgram with character n-grams (fastText)

- brother: bro, rot, oth, the, her (trigrams)
- brothers: bro, rot, oth, the, her, ers : almost the same!
- also enables to model Out-of-Vocabulary words (OOV), e.g. brotha
- rough way of modelling **morphology**: relation between words

- same objective as skipgram: $\log \left(1 + e^{-s(w_t, w_c)} \right) + \sum_{n \in \mathcal{N}_{t,c}} \log \left(1 + e^{s(w_t, n)} \right)$

- simply redefine similarity:
sum over all n-grams of
the word

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c$$

Empirical comparison

		sg	cbow	sisg-	sisg
AR	WS353	51	52	54	55
	GUR350	61	62	64	70
DE	GUR65	78	78	81	81
	ZG222	35	38	41	44
EN	RW	43	43	46	47
	WS353	72	73	71	71
ES	WS353	57	58	58	59
FR	RG65	70	69	75	75
RO	WS353	48	52	51	54
RU	HJ	59	60	60	66

Spearman's ρ

		sg	cbow	sisg
CS	Semantic	25.7	27.6	27.5
	Syntactic	52.8	55.0	77.8
DE	Semantic	66.5	66.8	62.3
	Syntactic	44.5	45.0	56.4
EN	Semantic	78.5	78.2	77.8
	Syntactic	70.1	69.9	74.9
IT	Semantic	52.3	54.7	52.3
	Syntactic	51.5	51.8	62.7

Word analogy

Welcome LLMs, exit Embeddings?

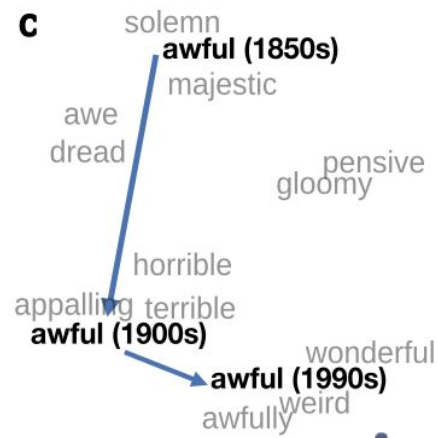
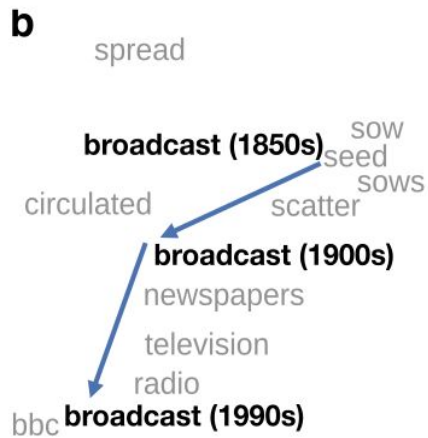
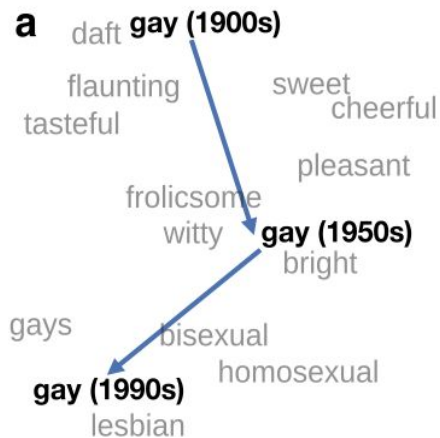
- Large Language Models are effective but not so efficient
- Embeddings are very lightweight, relevant for many industrial applications
- fastText: efficient implementation
- LLMs build on similar hypothesis and methods as Embeddings

In Summary

- NLP = research field at the intersection of Computer Science and Linguistics
- NLP = Many industrial applications, from Machine Translation to chatbots like ChatGPT or Information Extraction
- Meaning of a word is its use in the language: distributional semantics
- Skip-gram (word2vec): compute embeddings of words by predicting their context (**self-supervised learning**)
- Use as building block (**pre-training**) or solve analogies or measure word semantic similarity

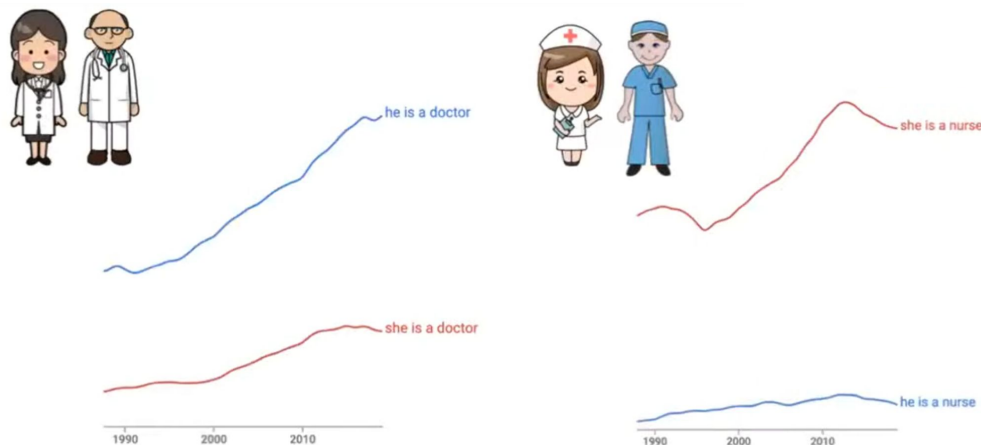
Limitations

- Cannot model polysemy: chair [furniture] vs chair [person] has only one embedding "chair"
- Meaning changes through time/domain...



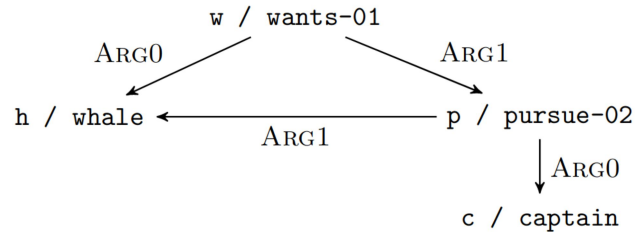
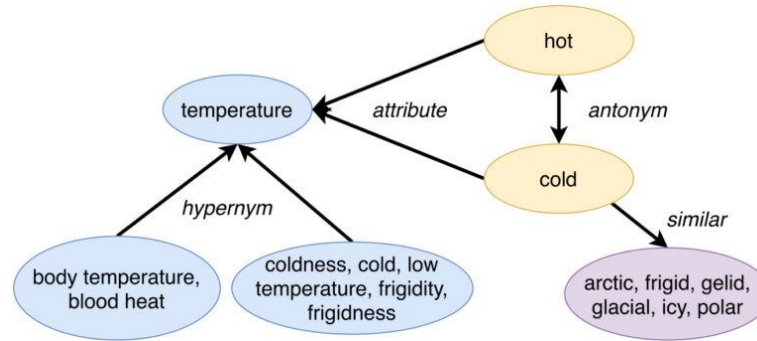
Embeddings reflect cultural bias!

- Statistical patterns in text reflect both **intrinsic meaning** and **extrinsic use**
- Ask “father : doctor :: mother : x”
x = nurse
- Ask “man : computer programmer :: woman : x”
x = homemaker



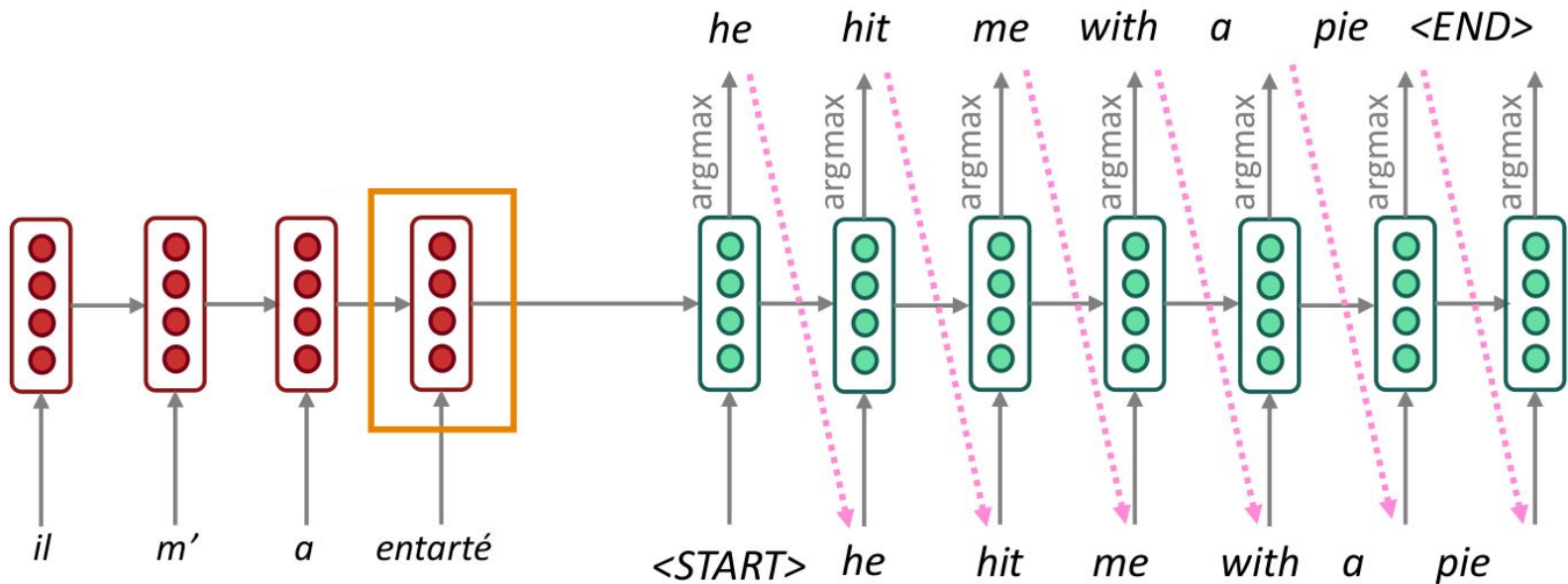
Alternatives to distributional semantics

- Not mainstream but may prove useful... Cf. Natural Language Processing by Jacob Eisenstein (2018).
- Chapter 12: Logical semantics
- Chapter 13: Predicate-argument semantics



"The whale wants the captain to pursue him"

Next class: models for sequences!





aivancity

PARIS-CACHAN

**advancing education
in artificial intelligence**