



自然语言处理

Natural language Processing

强继朋 副教授

扬州大学

第1章

NLP 介绍



Part.01

课程介绍

课程信息

- 课程网页：<https://yzu-nlp.github.io/>
 - 课程相关的所有信息：课件、实验、阅读资料、课程设计等
- 评分标准：作业（30%）+实验（20%）+期末考试（50%）
 - 作业：4次
 - 实验：4次

课程目标

- 了解NLP中不同领域的基本原理
- 理解NLP的理论概念和算法
- 为NLP应用建立模型的实践经验
- 独立完成NLP项目的开发



NLP背景介绍

什么是自然语言？

语言分为：

- 人类语言（自然语言）
- 机器语言

自然语言即为自然地随文化演化的语言如：汉语、英语、法语等

自然语言可以采取不同的形式，如语音或手势

有何特性？

什么是自然语言处理（NLP）？

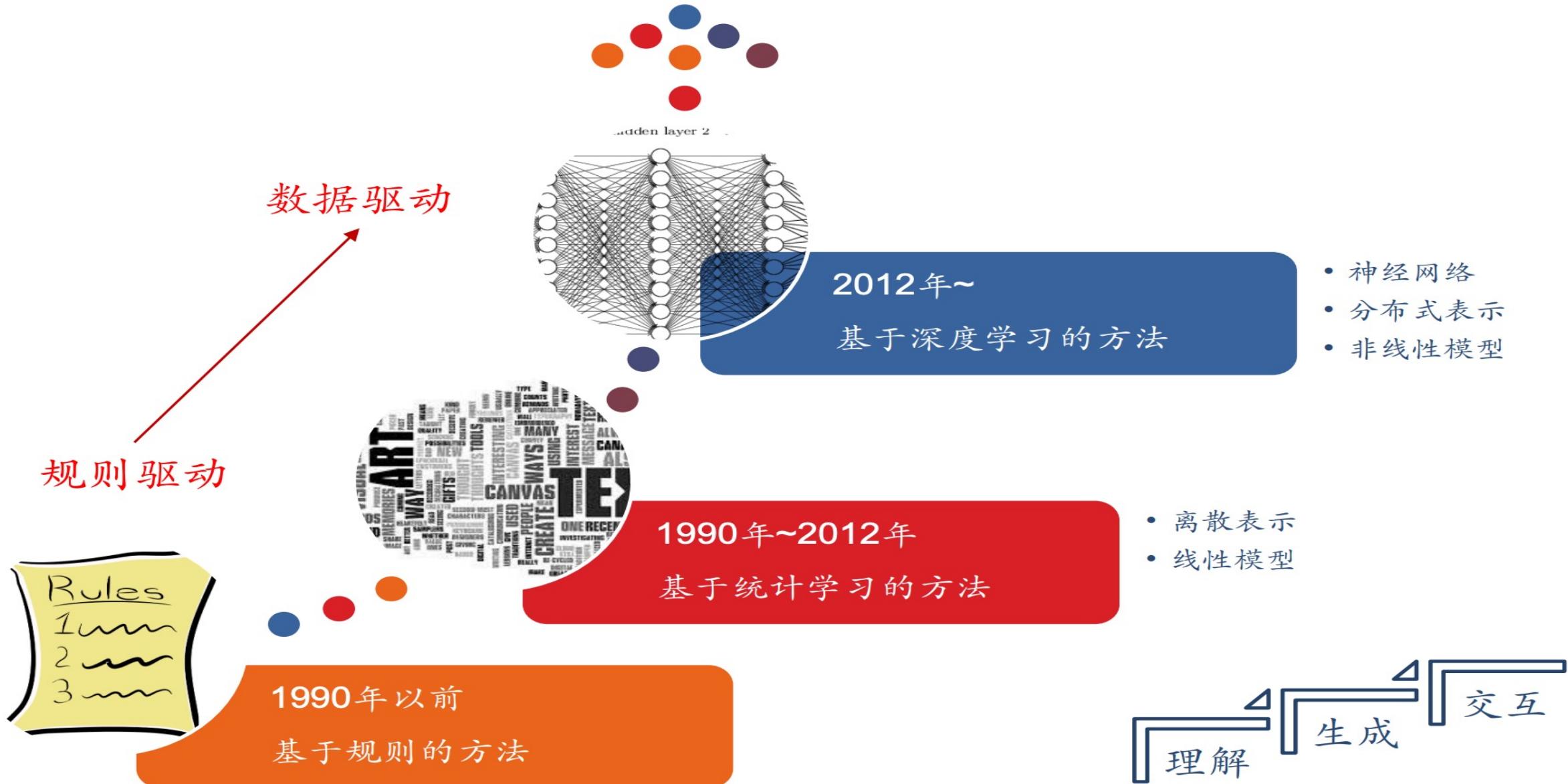
赋予计算机处理人类语言的能力。

人工智能的一个分支

也可被称为：speech and language processing, human language technology, computational linguistics, and speech recognition and synthesis.

目标是让计算机执行涉及人类语言的有用任务，
例如启用人机交流、改善人际交流、或简单地对文本或语音进行有用处理等任务。

NLP发展历程



NLP发展历程

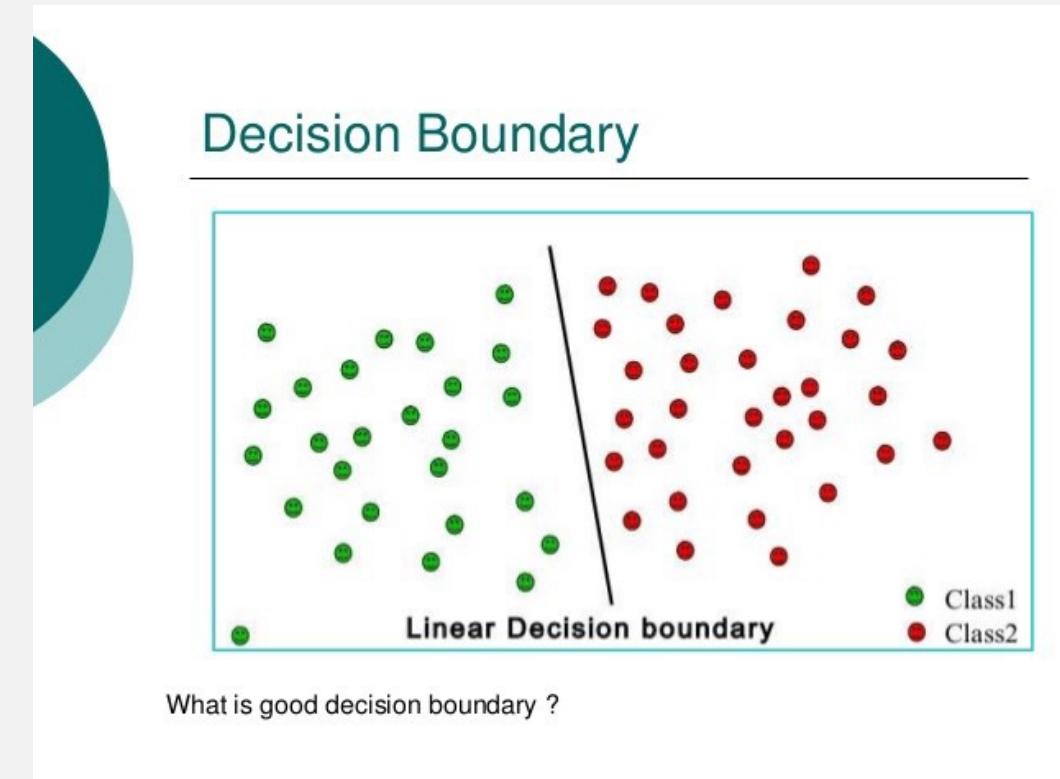
- 基于规则的方法（1950–1990）
 - 最古老的NLP方法
 - 利用语言学专家定义的规则
 - 局限于某个领域
 - 解决模糊性具有挑战性

“The spirit is strong, but the flesh is weak”

“The Vodka is good, but the meat is bad”

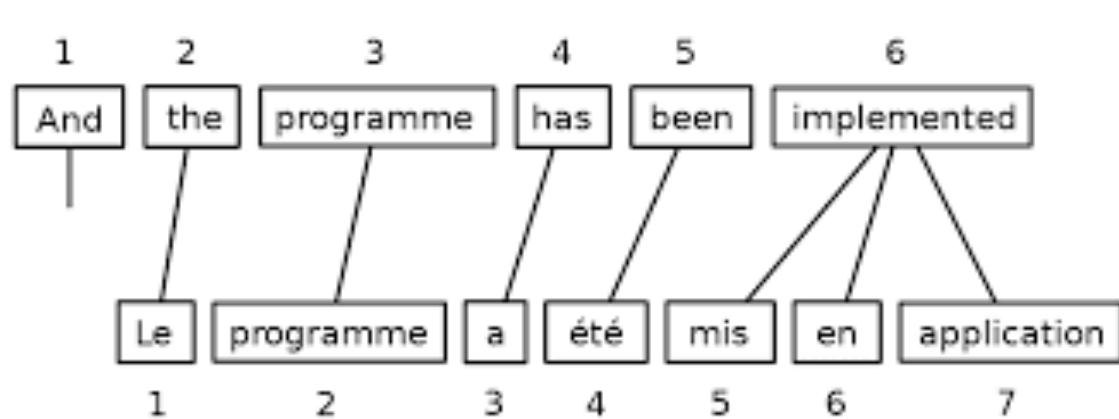
NLP发展历程

- 基于统计的方法（1990s–2010s）
 - 利用传统的机器学习方法
 - 通常采用的是概率模型的方法
 - 人工标注的训练数据
 - 特征工程
 - 训练带参数的模型
 - 应用模型到测试数据

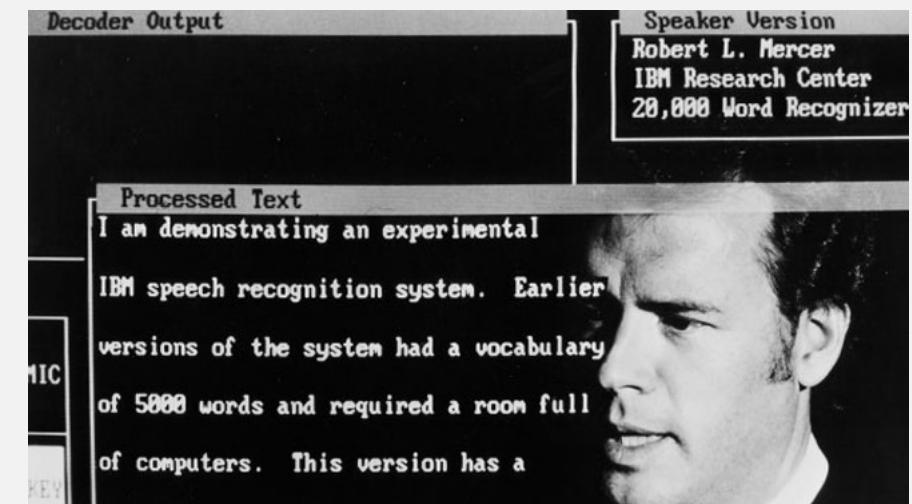


基于统计的方法

IBM的机器翻译模型



语音识别



Anytime a linguist leaves the group the (speech) recognition rate goes up.

- Fred Jelinek

NLP发展历程

- 基于深度学习的方法 (2010s-)
 - 取得了性能上的大幅度提升
 - 不需要手工定义特征
 - 构造非常大的神经网络
 - 在原始语料上进行预训练



36M sentence pairs

Russian: Машинный перевод - это круто!

English: Machine translation is cool!



Part.03

NLP应用场景

机器翻译 (Machine Translation)

The image displays two side-by-side machine translation interfaces. On the left is the Baidu Translate interface, and on the right is the Google Translate interface. Both show the same input text: "我爱中国!" (I love China!).

Baidu Translate Interface:

- Top navigation bar: 文本翻译 (Text Translation), 文档翻译 (Document Translation) (highlighted in blue), 人工翻译 (Human Translation), 视频翻译 (Video Translation), AI同传 (AI Simultaneous Interpretation), 翻译API (Translation API), 课程优选 (Course Selection), 官网 (Official Website), and 开通文档翻译VIP (Activate Document Translation VIP).
- Input field: 检测到中文(简体) (Detected as Simplified Chinese) and a dropdown arrow pointing to English.
- Translation button: 翻译 (Translate).
- Output field: I love China.
- Bottom controls:Speaker icon, star icon, report icon, note icon, and a double-language switch toggle.

Google Translate Interface:

- Top navigation bar: 检测语言 (Detection Language), English, 中文 (简体) (Simplified Chinese) (highlighted in blue), 德语 (German), and a dropdown arrow.
- Input field: 我爱中国! (Wǒ ài zhōngguó!).
- Output field: I love China!
- Bottom controls: Speaker icon, microphone icon, character count (7 / 5,000), a switch labeled 拼 (Pinyin), and sharing icons.

文字

文档

检测语言 英语 中文（简体）

德语



中文（简体）

英语

日语



← 搜索语言

阿尔巴尼亚语

德语

捷克语

迈蒂利语

僧伽罗语

信德语

阿拉伯语

迪维希语

卡纳达语

毛利语

世界语

匈牙利语

阿姆哈拉语

蒂格尼亞語

科西嘉语

梅泰语（曼尼普尔语）

斯洛伐克语

修纳语

阿萨姆语

多格来语

克里奥尔语

蒙古语

斯洛文尼亚语

亚美尼亚语

阿塞拜疆语

俄语

克罗地亚语

孟加拉语

斯瓦希里语

伊博语

埃维语

法语

克丘亚语

米佐语

苏格兰盖尔语

伊洛卡诺语

艾马拉语

梵语

库尔德语（库尔曼吉语）

缅甸语

宿务语

意大利语

爱尔兰语

菲律宾语

库尔德语（索拉尼）

苗语

索马里语

意第绪语

爱沙尼亚语

芬兰语

拉丁语

南非科萨语

塔吉克语

印地语

奥利亚语

弗里西语

拉脱维亚语

南非祖鲁语

泰卢固语

印尼巽他语

奥罗莫语

高棉语

老挝语

尼泊尔语

泰米尔语

印尼语

巴斯克语

格鲁吉亚语

立陶宛语

挪威语

泰语

印尼爪哇语

白俄罗斯语

贡根语

林格拉语

旁遮普语

土耳其语

约鲁巴语

班巴拉语

古吉拉特语

卢干达语

葡萄牙语

土库曼语

越南语

保加利亚语

瓜拉尼语

卢森堡语

普什图语

威尔士语

机器翻译的任务目标
是自动地将一种语言
翻译为另一种语言。

人机对话 (Conversational agent / Dialogue system)



Welcome to

EEEEEE	LL	IIII	ZZZZZZ	AAAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LL	II	ZZZ	AAAAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LLLLLL	IIII	ZZZZZZ	AA AA

Eliza is a mock Rogerian psychotherapist.

The original program was described by Joseph Weizenbaum in 1966.

This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?

YOU: Men are all alike.

ELIZA: What is the connection, do you suppose ?

YOU: They're always bugging us about something or other.

ELIZA: Can you think of a specific example ?

YOU: Well, my boyfriend made me come here.

ELIZA: Is it important to you that your boyfriend made you come here ?

YOU: He says I'm depressed much of the time.

ELIZA: I am sorry to hear that you are depressed.

YOU: It's true. I am unhappy.

ELIZA: Can you explain what made you unhappy ?

YOU:

ELIZA 是早期的一个聊天机器人 (1966)

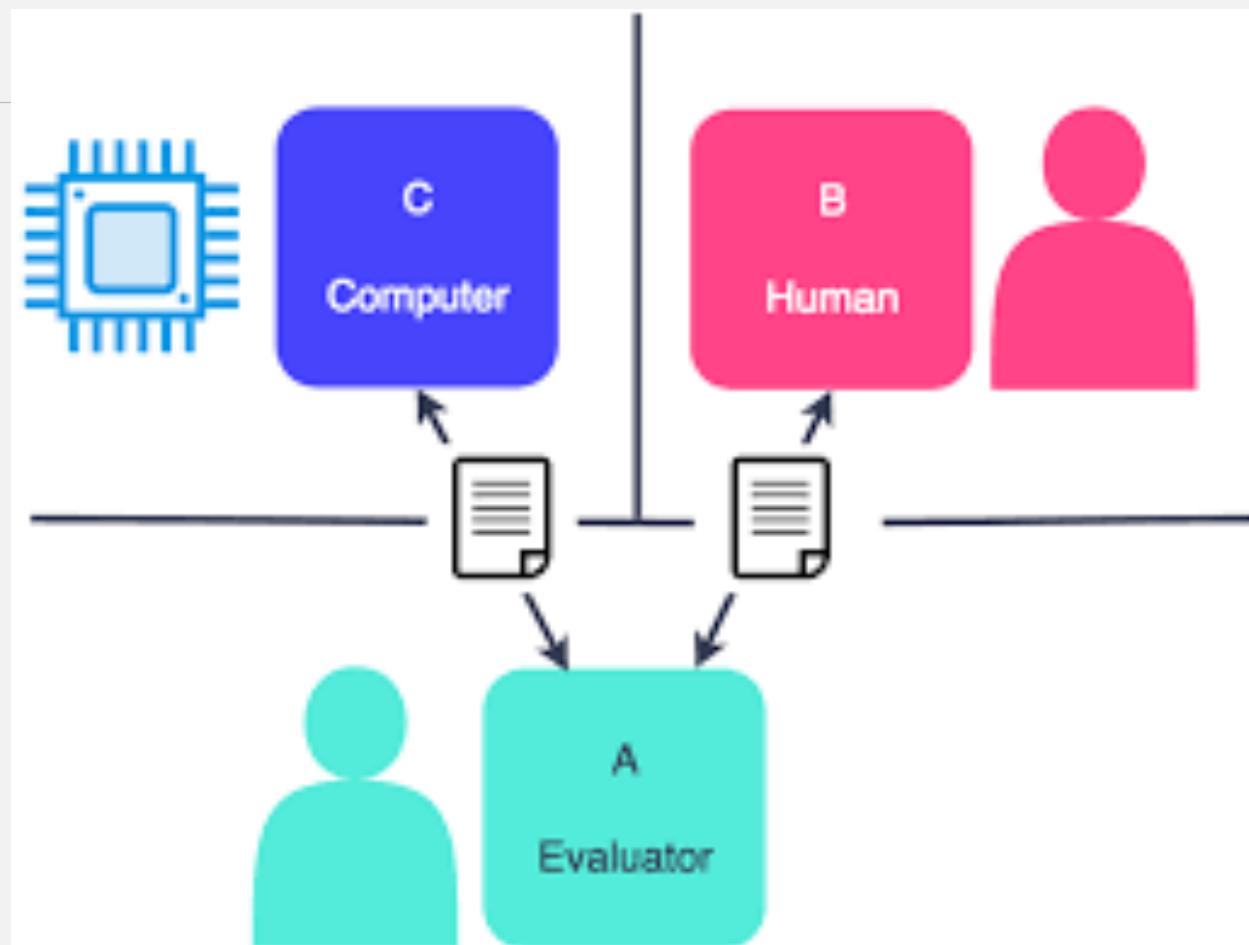
使用模式匹配来处理输入
并将其转换为合适的输出

人机对话 (Conversational agent / Dialogue system)



汽车智能语音识别系统

人机对话 (Conversational agent / Dialogue system)



三个参与者：两个人和一台电脑

中间的人扮演裁判的角色，
并向另外两个参与者提出问题，
裁判通过审视参赛者的回答，
以判断哪个是真人哪个是机器。

两位参赛者(机器和人类)的任务为：
让裁判相信自己是人类，
并判定另一位参赛者是机器。

图灵测试

自动问答 (Question Answering)

The screenshot shows a Baidu search interface. The search bar contains the query "毛泽东哪年出生". Below the search bar are several navigation links: "网页" (selected), "图片", "视频", "资讯", "贴吧", "采购", "地图", and "知道". A message indicates approximately 15,900,000 results found. The main content area features a large blue header "历史12月26日:毛泽东诞辰" followed by the year "1893年". A detailed description below states: "1893年12月26日,毛泽东在湖南湘潭韶山冲诞生。毛泽东是中国共产党、中国人民解放军、中华人民共和国的主要缔造者,中国各族人民的伟大领袖。" with a "更多 >" link.

Baidu 百度

毛泽东哪年出生

× | 麦 | 相机

网页 图片 视频 资讯 贴贴吧 采购 地图 知道

百度为您找到相关结果约15,900,000个

搜索工具

历史12月26日:毛泽东诞辰

1893年

“1893年12月26日,毛泽东在湖南湘潭韶山冲诞生。毛泽东是中国共产党、中国人民解放军、中华人民共和国的主要缔造者,中国各族人民的伟大领袖。” [更多 >](#)

自动问答 (Question Answering)

The first recorded travels by Europeans to China and back date from this time. The most famous traveler of the period was the Venetian Marco Polo, whose account of his trip to "Cambaluc," the capital of the Great Khan, and of life there astounded the people of Europe. The account of his travels, *Il milione* (or, *The Million*, known in English as the *Travels of Marco Polo*), appeared about the year 1299. Some argue over the accuracy of Marco Polo's accounts due to the lack of mentioning the Great Wall of China, tea houses, which would have been a prominent sight since Europeans had yet to adopt a tea culture, as well the practice of foot binding by the women in capital of the Great Khan. Some suggest that Marco Polo acquired much of his knowledge **through contact with Persian traders** since many of the places he named were in Persian.

How did some suspect that Polo learned about China instead of by actually visiting it?

Answer: through contact with Persian traders

自动问答 (Question Answering)

简单 – 定义性问题

- 例如：时间地点任务等事实性问题。
- 解决途径-网络搜索

复杂 – 推理性问题

- 根据已知得出结论
- 综合多源信息进行总结归纳
- 解决途径-信息提取，词义消歧等

例子？

情感分析

Loves the German bakeries in Sydney. Together with my imported honey it feels like home	Positive
@VivaLaLauren Mine is broken too! I miss my sidekick	Negative
Finished fixing my twitter...I had to unfollow and follow everyone again	Negative
@DinahLady I too, liked the movie! I want to buy the DVD when it comes out	Positive
@frugaldougal So sad to hear about @OscarTheCat	Negative
@Mofette brilliant! May the fourth be with you #starwarsday #starwars	Positive
Good morning thespians a bright and sunny day in UK, Spring at last	Positive
@DowneyisDOWNEY Me neither! My laptop's new, has dvd burning/ripping software but I just can't copy the files somehow!	Negative



Part.04

NLP的挑战

为什么语言很难理解？

多少话人都看不懂还想
让机器懂!!!

22:47

豆 豆瓣

4G



...



大青枣
说

2022-07-11 21:55

那些你遇到过的热情大妈

在公共汽车上，我给一位大妈让了座，大妈高兴地和我攀谈了起来，问道：“孩子今年多大了？”

我说：“20岁了。”

大妈羡慕地说：“那你长得可真年轻，看起来也就30岁出头，想不到孩子都这么大了。”

词义的多样性

一个词语具有多个意思

The fisherman went to the **bank**.

bank¹

/baNGk/ 

noun

plural noun: **banks**

1. the land alongside or sloping down to a river or lake.

"willows lined the bank"

synonyms: [edge](#), [side](#), [shore](#), [coast](#), [embankment](#), [bankside](#), [levee](#), [border](#), [verge](#), [boundary](#), [margin](#), [rim](#), [fringe](#); [More](#)

1. a financial establishment that invests money deposited by customers, pays it out when required, makes loans at interest, and exchanges currency.

"I paid the money straight into my bank"

synonyms: [financial institution](#), [merchant bank](#), [savings bank](#), [finance company](#), [trust company](#),

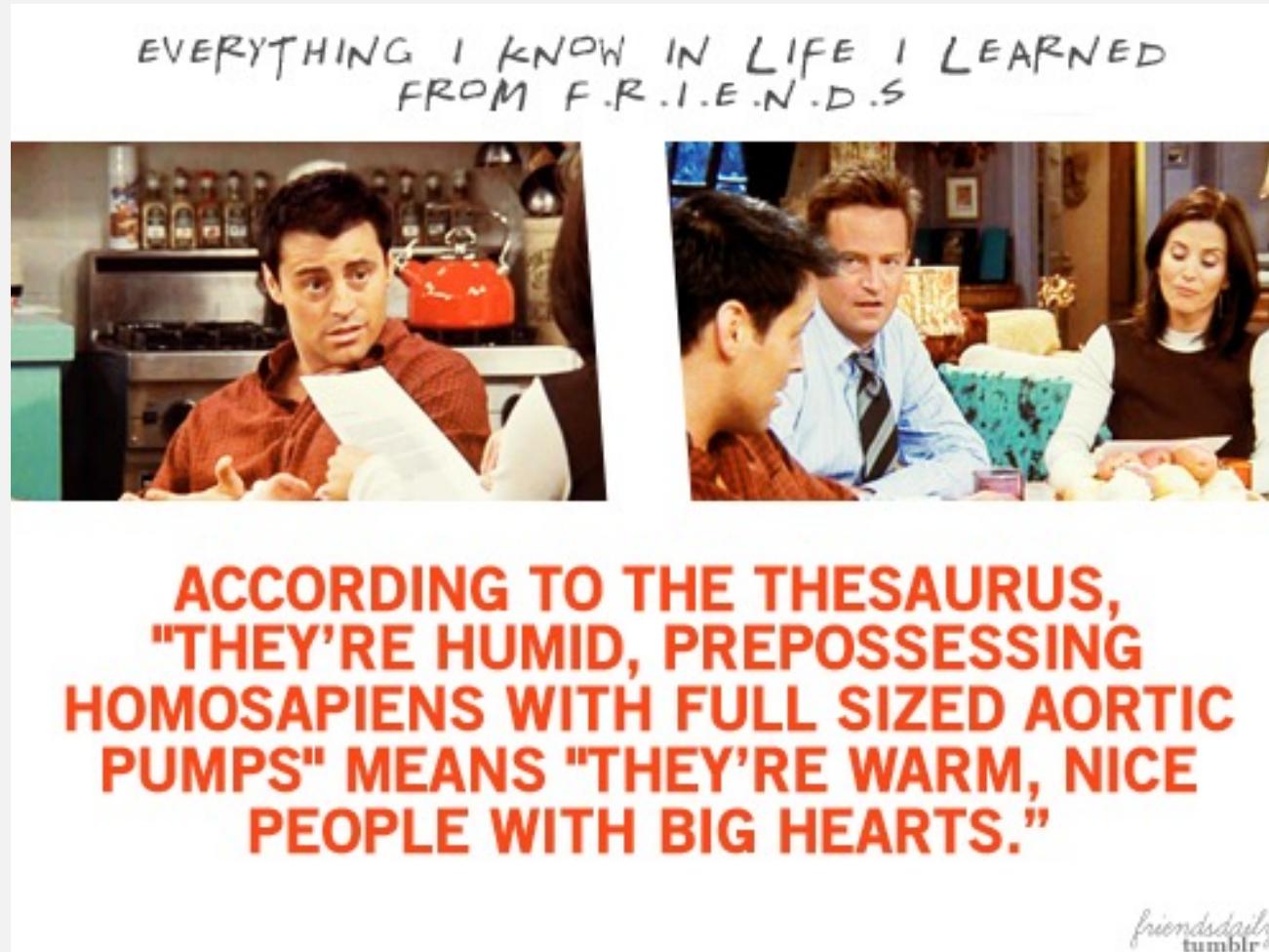
词义的多样性

The fisherman went to the *bank*. He deposited some money.

词义消歧 (Word sense disambiguation)

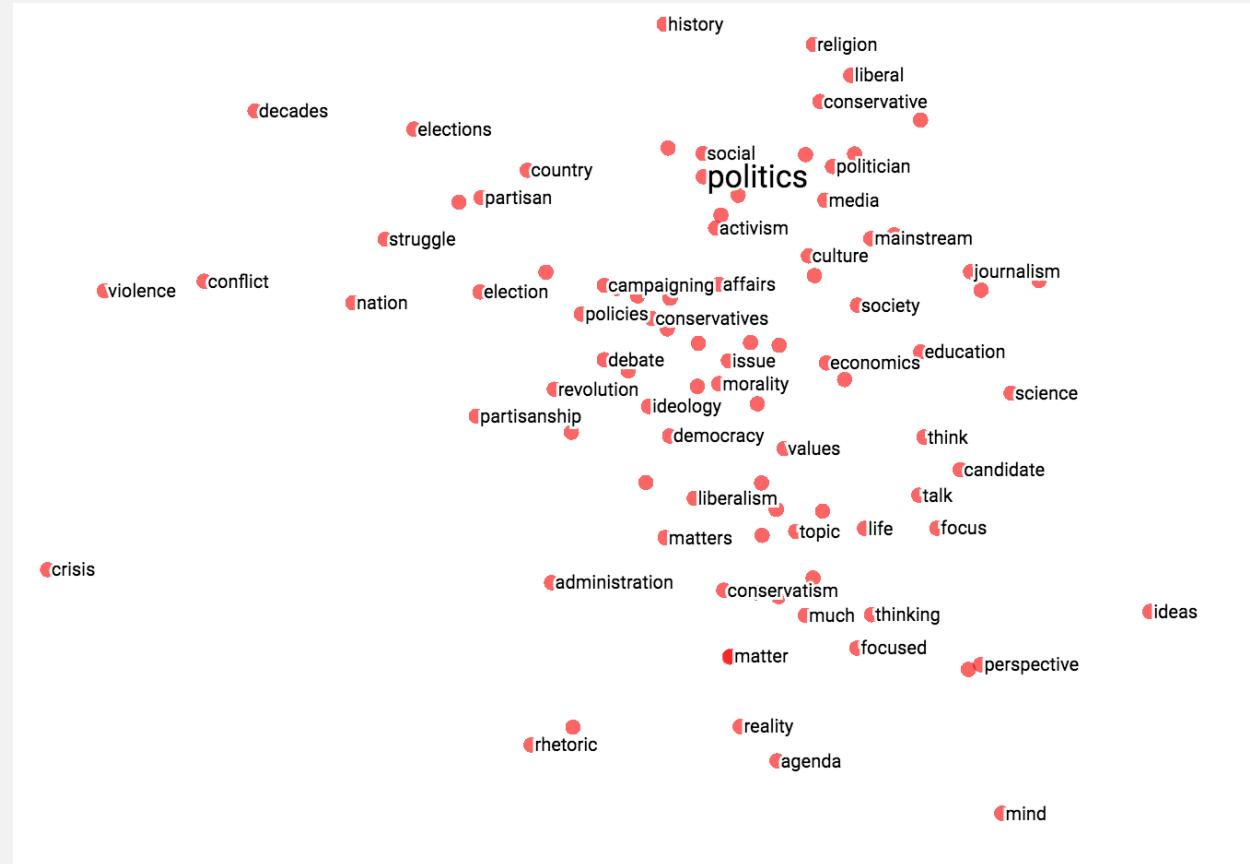
词语的多样性

多个词语具有相同的意思



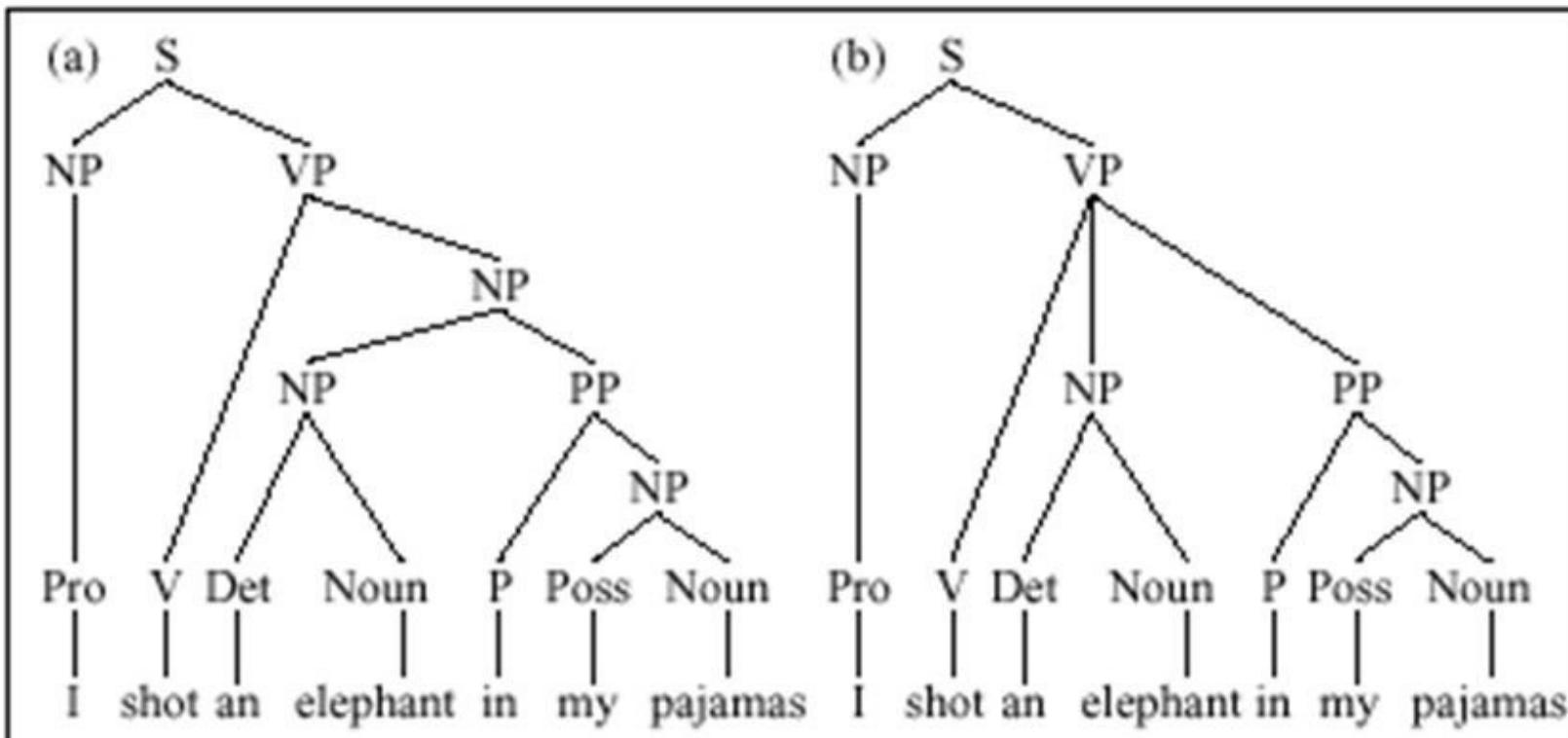
词语的分布式表示

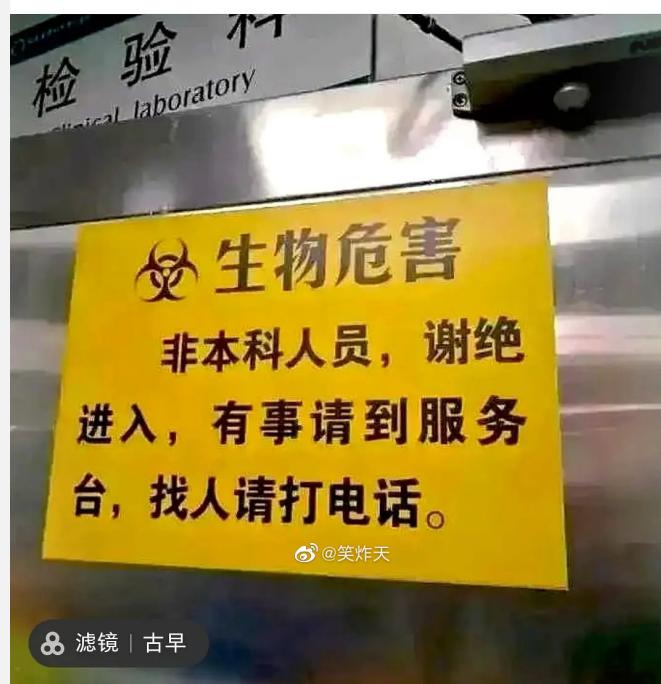
将单词投影到连续向量空间



句法的模糊性

I shot an elephant in my pajamas





滤镜 | 古早

没考上本科怎么了?

考上大专怎么了? 大专怎么了? ? ? 你们一个个本科不得了了, 第一学历大专怎么了? 多的是专升本, 后来又考研读博出国留学, 照样成才, 照样过的好! 不比你们这些差!

不停说大专没有用, 大专第一学历不行, 你们只不过就是一些幸运儿罢了, 人生比较顺利, 考了个本科, 天天在这里叫唤大专不行不行!

很多人不是想考大专不是想这样, 很多是人生比较坎坷曲折, 没有考上本科, 不得了了, 你本科也不见得拍怎么样!!!

说点什么...



77



4



891

@半糖小哥

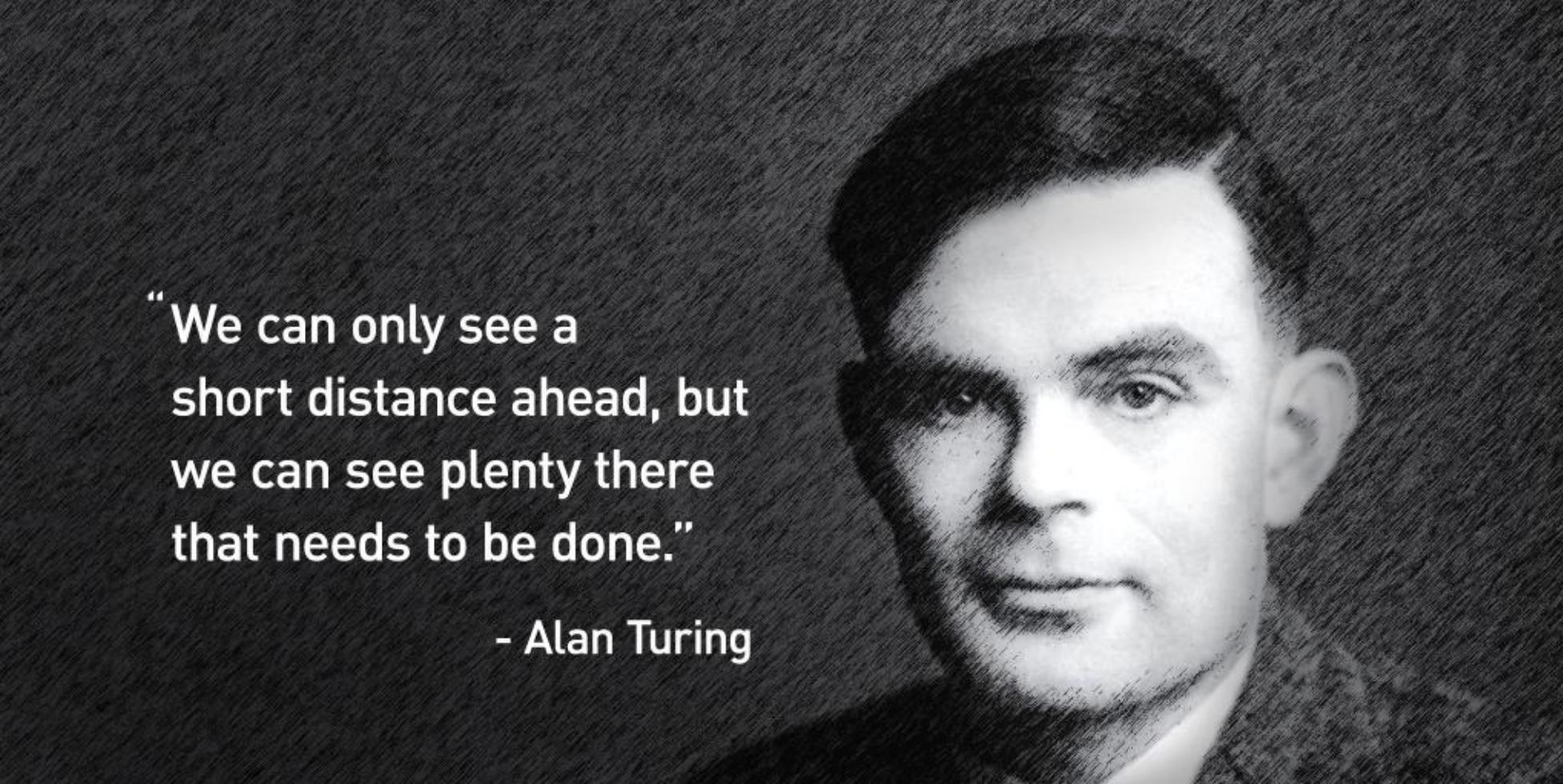
篇章的歧义 (Discourse ambiguity)

Alice invited Maya for dinner but **she** cooked her own food

she = Alice or Maya?

... and brought it with her. **Maya**

... and ordered a pizza for her guest. **Alice**

A black and white portrait of Alan Turing, a English polymath and a pioneer of computing science. He is shown from the chest up, wearing a dark suit jacket over a white shirt. His gaze is directed slightly to his right. The background is dark and out of focus.

**“We can only see a
short distance ahead, but
we can see plenty there
that needs to be done.”**

- Alan Turing



Part.05

课程内容

课程内容

1. 语言模型
2. 文本分类
3. 向量语义和词嵌入
4. 神经网络语言模型
5. 序列处理语言模型 (RNN, LSTM, Transformer)
6. 机器翻译和编码器-解码器模型
7. 预训练语言模型

语言模型 (Language models)

What is the weather in New York?



It is 76°F and _____

- red ?
- 24.44 C ?
- sunny ?

语言模型 (Language models)

Today, in New York, it is 76 F and red

vs

Today, in New York, it is 76 F and sunny

- Both are grammatical
- But which is more likely?

语言模型 (Language models)

目标：计算一个句子或者一个单词序列的概率

- $P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$

相关的任务：计算下一个单词的概率

- $P(w_5 | w_1, w_2, w_3, w_4)$

一个用于计算下列两种概率的模型：

- $P(W)$ or $P(w_n | w_1, w_2 \dots w_{n-1})$ 则被称为语言模型

语言模型 (Language models)

New Message Cancel

Google

how is the weather in new| To:

- how is the weather in new york
- how is the weather in new orleans
- how is the weather in new orleans in october
- how is the weather in new jersey
- how is the weather in new york in october
- how is the weather in new orleans in november
- how is the weather in new orleans in december
- how is the weather in new orleans in september
- how is the weather in new mexico
- how is the weather in new york in september

Language models are the ↗

best | models | same

q w e r t y u i o p
a s d f g h j k l
z x c v b n m ↵
123 ⚡ space ↩ return

American singer-songwriter, guitarist, and record producer. King intro

謝 謝