

计算语言学 Computational Linguistics

教师: 孙茂松

Tel:62781286

Email:sms@tsinghua.edu.cn

TA: 林衍凯

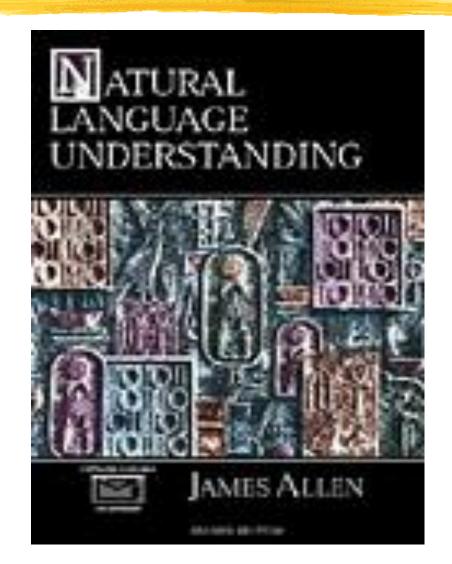
Email: linyankai 423@qq. com

郑重声明

- 此课件仅供选修清华大学计算机系研究生课《计算语言学》(70240052)的学生个人学习使用,所以只允许学生将其下载、存贮在自己的电脑中。未经孙茂松本人同意,任何人不得以任何方式扩散之(包括不得放到任何服务器上)。否则,由此可能引起的一切涉及知识产权的法律责任,概由该人负责。
- 此课件仅限孙茂松本人讲课使用。除孙茂松本人外,凡授课过程中,PTT文件显示此《郑重声明》之情形,即为侵权使用。

主要参考书

http://www.unigiessen.de/~g91062/Se minare/gkcl/Allen95/al1995co.ht m James Allen, Natural Language Understanding (2nd edition), The Benjamin/Cummings Publishing Company, Inc. 1994

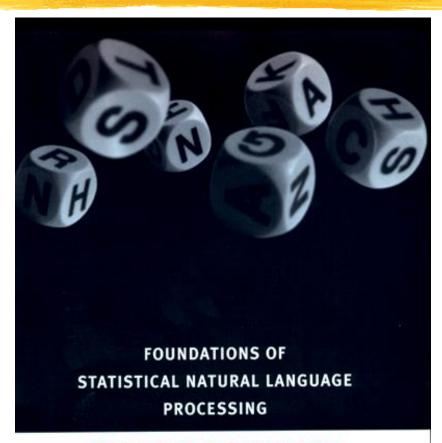


主要参考书

http://nlp.stanford.edu/fsnlp/

Chris Manning and Hinrich Schütze,

Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.



CHRISTOPHER D. MANNING AND HINRICH SCHÜTZE

要求

● 成绩:编程作业(+DEMO) + seminar + 笔头作业 + 不 定期的课堂小测验

对CL不感兴趣的同学,请不要选。 作业未按时交,该次成绩以零分计(一般电子版交)

● 不分组: (一人一组)

纪律:

- 鼓励讨论,严禁抄袭\拷贝: 第一次发现,成绩记0分,并警告; 第二次发现,整个课程不通过。 (以上均无论抄与被抄者)。
- 不得迟到,早退,不得吃东西。关手机。

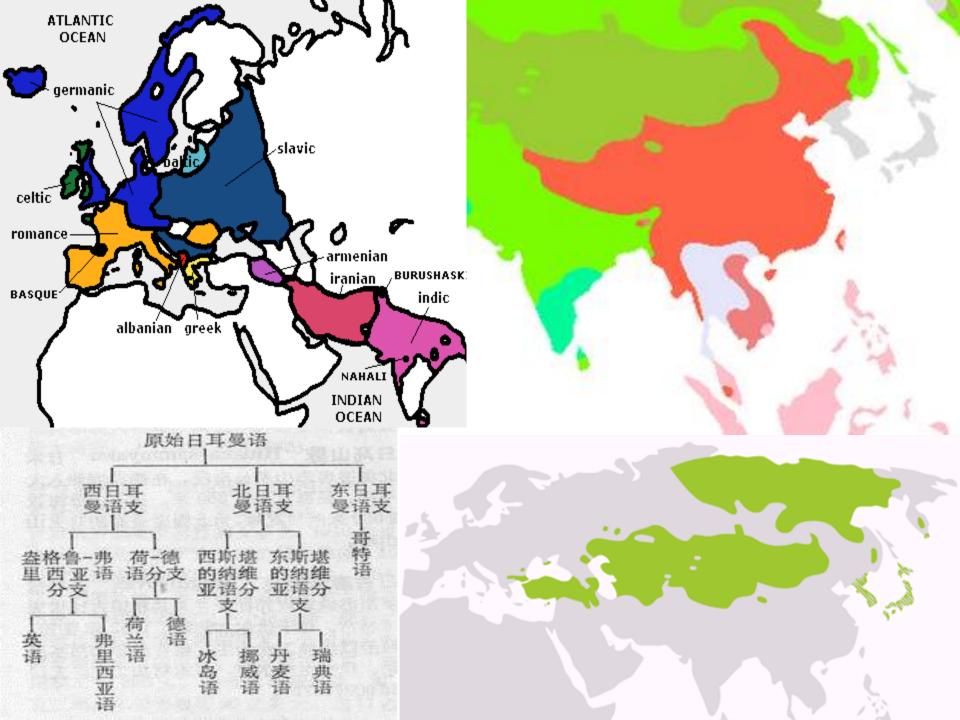
第一章 引言 (Part 1)

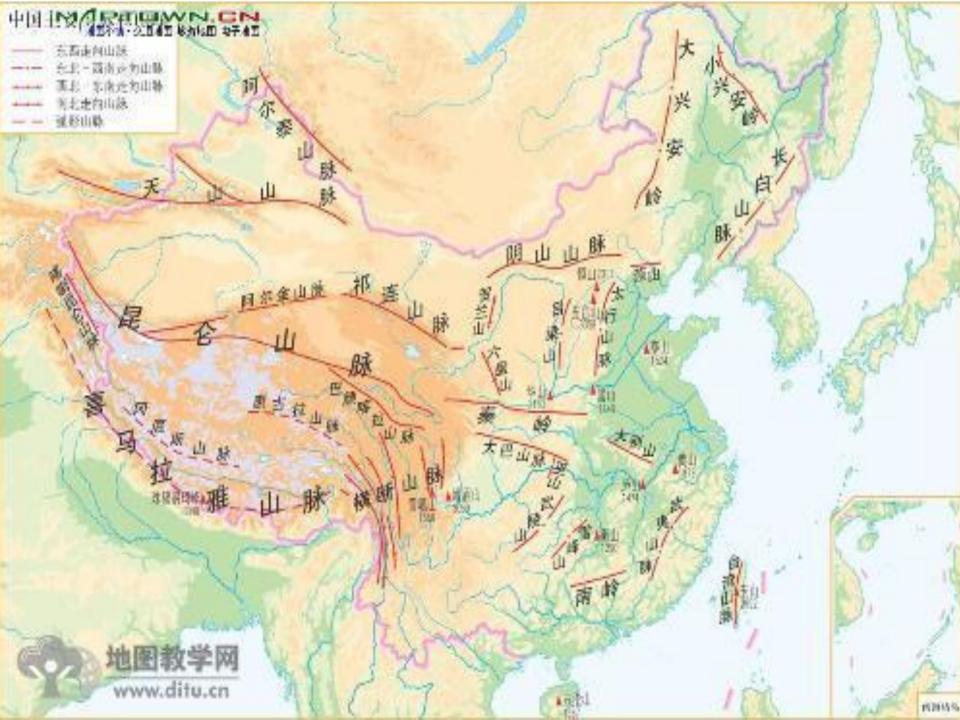
- Object: natural language 6000种(3000-10000种不等)
- The role of language:
 - * the most natural means of communication between

humans

The tower of Babel是《圣经》之《旧约·创世纪》第11章 (1-9)中一座通天塔,由挪亚(Noa)的后代在示拿 (Shinnar)所建







- The role of language:
 - * record of knowledge

"搜索"产生的权力 http://www.sina.com.cn 2006年02月17日

Google发现宝马德国网站通过"技术作弊"人为地提高它在其搜索引擎中的排名,涉嫌欺骗用户,便在引擎中删除宝马汽车的德国网站。而宝马则认为 Google没有事先通知就采取公开批评和封杀的行为有违常规,并同时与Google 进行联系,希望能尽快解决。

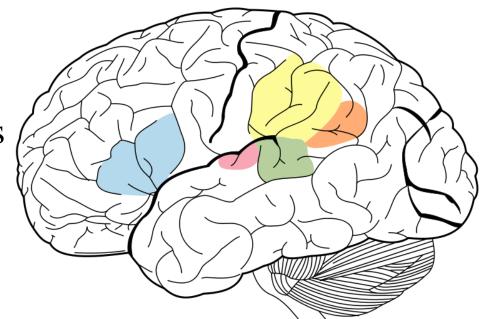
- ,"由于Google制定了信息社会在未来一个阶段的游戏规则,所有人在信息时代寻找一席之地的时候,都必须遵守Google的游戏规则"。人们向Google索求信息,却不觉间投票给了它很大的支配性权力。
 - *语言与思维

Linguistics

Psycholinguistics

Philosophy

Computational Linguistics



• Academic disciplines:

linguistics, psycholinguistics, cognitive science, computer science, ...

<u>Linguistics</u> + <u>Computation (computer)</u>

- The task: to create computational models of language
 * scientific motivation
 the nature of linguistic communication
 validation (testing) of linguistic models
 - * practical motivation

Most of human knowledge is recoded in linguistic form,...: capability of intelligence for machines NLP capabilities would revolutionize the way computers are used

Effective human-machine communication (too complex to everyone now!)

• information retrieval, extraction, filtering, classification and summarization, Search engine, digital library, e-commerce etc.

query: keywords

* English: bank

* Chinese word segmentation:

和服 | 务 | 于三日后裁制完毕。

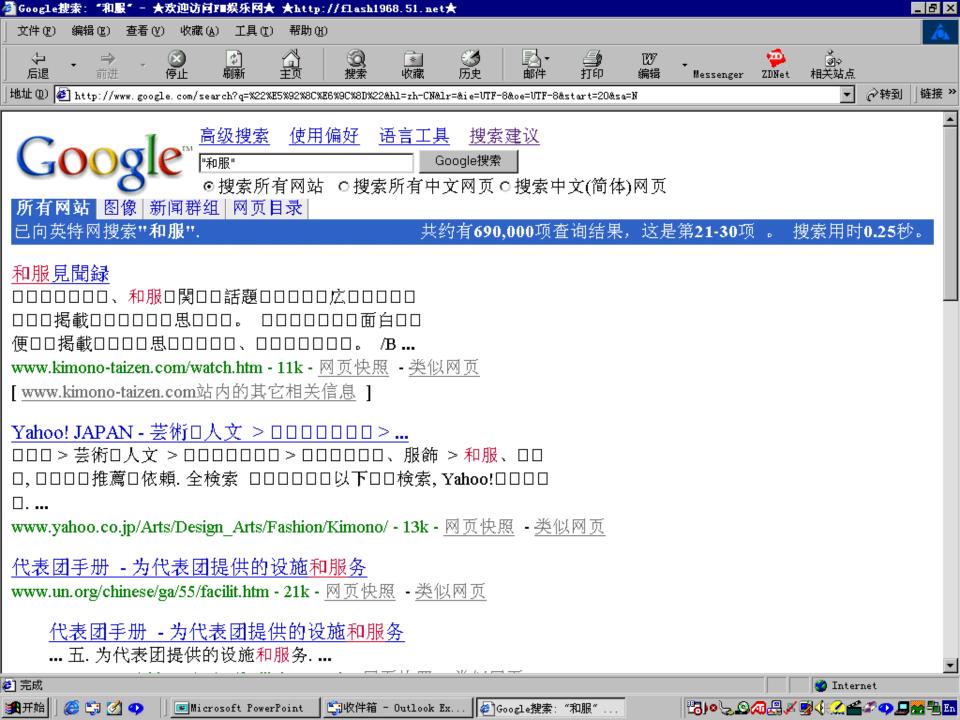
(kimono must)

这个酒店的设施 | 和 | 服务 | 是一流的。

(and service)

* unknown word processing:

高<u>海燕</u> (storm petrel)









⋧转到│

(搜索用时 0.27 秒)





machine translation











- cross-lingual information retrieval
- question-answering systems: query a database
- tutoring system
- language model for speech recognition (The listening machine)

Gong1Shi4 (Pinyin):

公式(formula) 工事(work)

攻势(offensive) 公事(public affair)

数学公式, 摧毁了敌方工事,

猛烈的攻势, 例行公事

- spoken language control of a machine: mobile phone
- language model for OCR

input: 我们要振奋精神

top one given by pure OCR: 我优要扳奋精种

output (recovered by the language model):

我们要振奋精神

top ten candidates: similarity measures given by OCR:

我钱钱载哦栽哉裁劣绥397682700722774781787815838851 优仍价价价奶砧犯奶妨86887892994795396497998410091010 要耍密穷安壁驻努窑垂627650730747749802808818836838 扳报叔嵌奴振按寂寂蔽663709743746755772799815822824 奋夯杏蚕香脊秀吞畜番192381393436438471507534543544 精精指洁括治捐活治估756787791799824826836875885886 种神衬祥科钟拌样拎补463548555575636663671681689694

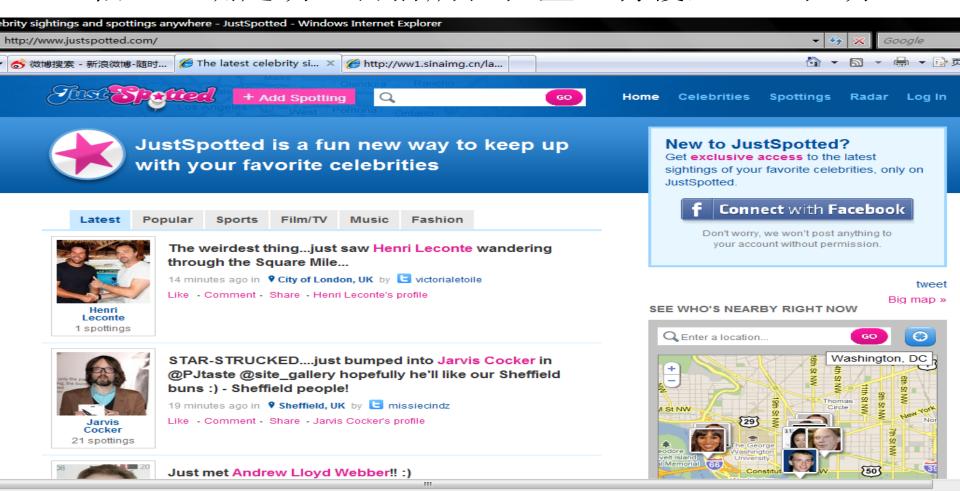
language model for TTS (the talking machine)

<u>一行</u>白鹭上青天 (a line of) 单信<u>一行</u>于昨日抵达香港 (delegation) 你 | 要 | 先 | 搞懂 | 文章 | <u>大意</u> (n) (main points) 你 | 千万 | 别 | <u>大意</u> (adj) (careless)

- proof-reading你己(己)不小了。
- traditional and simplified Chinese character conversion
 发展~头发 后来~皇后

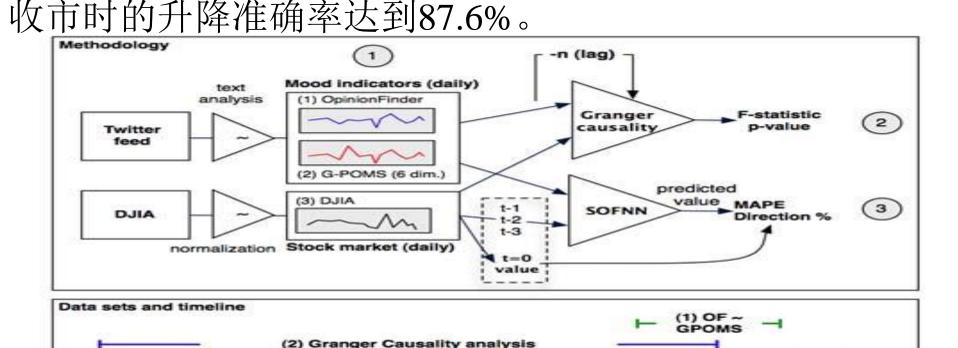
新的应用示例: Twitter用于确定明星位置

JustSpotted.com与Twitter合作,获取每天超过5000万条公开tweet信息,使用自然语言技术筛选包含明星动向的tweet信息,确定明星目前所在位置。将覆盖7000位明星。



新的应用示例: Twitter 用于预测股市行情

根据行为经济学的原理,情绪对个人的行为和决策有深远影响。公众情绪有可能影响响相关的经济指标。 印第安纳大学研究结果表明,根据 twitter 上收集到的公众情绪状态,道琼斯指数预测的准确率大大提高,预测每天



aug

sep

oct

nov

dec

dec20

2008

(3) SOFNN training

iul

jun

may

apr

feb28

2008

新的应用示例: Twitter 用于评估幸福感

东北大学、哈佛大学:

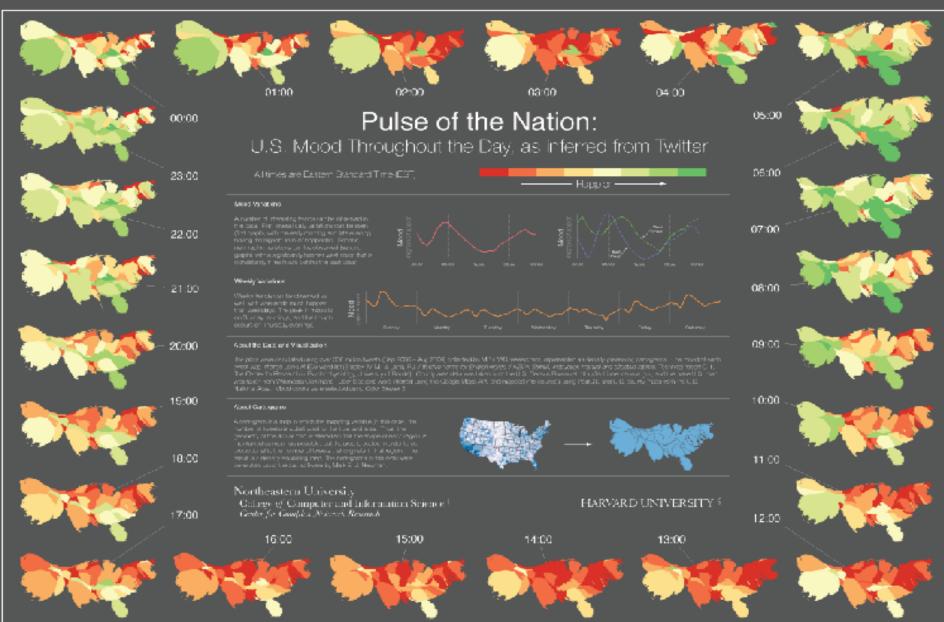
Twitter Mood Map(Pulse of the Nation)

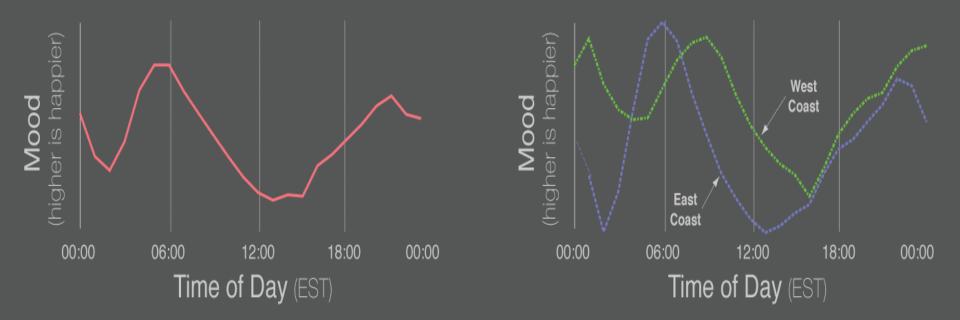
早上和傍晚最快乐,下午最郁闷;

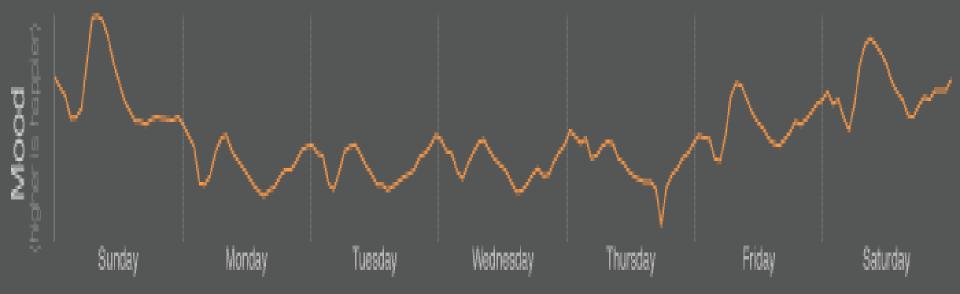
周六最快乐,好心情在周日早上达到巅峰 ,周四晚上跌至谷底

U.S. Mood Throughout the Day inferred from Twitter

Click for high-resolution PDF version (11MB)

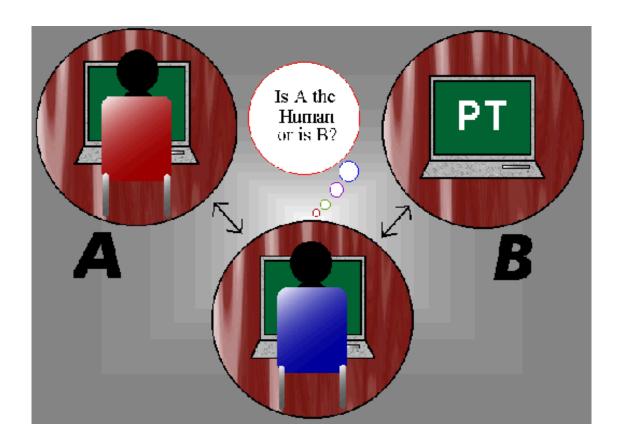






Black box evaluation: Turing test "Can machines think?"

The Imitation Game



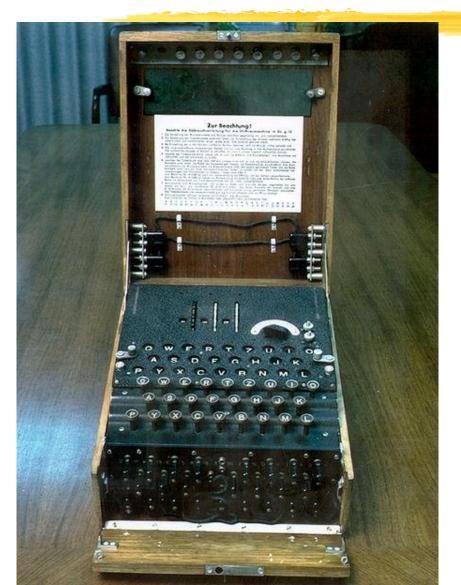
Alan Turing: Founder of computer science, mathematician, philosopher

"计算机之父"、"人工智能之父"、 "破译之父"

http://www.turing.org.uk/turing/



- 1912 (23 June): Birth, Paddington, London
- 1926-31: Sherborne School
- 1931-34: Undergraduate at King's College, Cambridge University
- 1932-35: Quantum mechanics, probability, logic
- 1935: Elected fellow of King's College, Cambridge
- 1936: The Turing machine, computability, universal machine
- 1936-38: Princeton University. Ph.D. Logic, algebra, number theory
- 1938-39: Return to Cambridge. Introduced to German Enigma cipher machine
- 1939-40: The Bombe, machine for Enigma decryption
- 1939-42: Breaking of U-boat Enigma, saving battle of the Atlantic
- 使用了分析英格玛机的数学结构的方法,也使用了德文交流中的一些规律性,例如德国人会严格地
- 在每天早晨的天气预报发布中用到"WETTER"这个词(即德语中的"天气")
- 1943-45: Chief Anglo-American crypto consultant. Electronic work.
- 1945: National Physical Laboratory, London
- 1946: Computer and software design leading the world.
- 1947-48: Programming, neural nets, and artificial intelligence
- 1948: Manchester University
- 1949: First serious mathematical use of a computer
- 1950: The Turing Test for machine intelligence
- 1951: Elected FRS. Non-linear theory of biological growth
- 1952: Arrested as a homosexual, loss of security clearance
- 1953-54: Unfinished work in biology and physics
- 1954 (7 June): Death (suicide) by cyanide poisoning, Wilmslow, Cheshire.



German military Enigma machine 26个字母能替代8万亿个谜文字 母。如果改动接线,变化会超过2千5百万亿。

图灵的"图灵炸弹"解码机主要 由继电器构成,由光电阅读器直 接读入密码。平均每天破译敌方 三千条密码。有些重要信息直接 送达英国首相丘吉尔。

2009年9月,布朗首相代表英国 政府向艾伦·图灵作出了明确的道 歉。"如果没有他的卓越贡献,

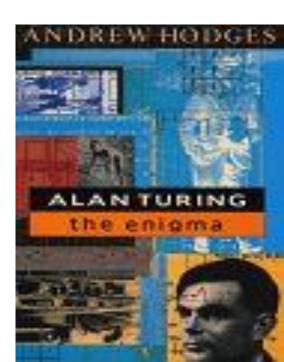
二战的历史也许会被重写"







牛津大学著名数学家安德鲁哈吉斯在为图灵写的一部脍炙人口的传记《谜一样的图灵》(Alan Turing: The Enigma)中这样描述到:"图灵似乎是上天派来的一个使者,匆匆而来,匆匆而去,为人间留下了智慧,留下了深邃的思想,后人必须为之思索几十年、上百年甚至永远。"



ELIZA: in the mid-1960s at MIT: patient and therapist

User: ELIZA:

Man are all alike. In what way?

They're always bugging us about

well, my boyfriend made me come here.

Can you think of a specific example?

Your boyfriend made you come here.

He says I'm depressed much of the time. I am sorry to hear you are depressed.

It's true, I am unhappy.

Do you think coming here will help you

not to be unhappy?

I need some help; that much seems certain. What would it mean to you if you got

some help?

Perhaps I could learn to get along with my

Mother. Tell me more about your family.

My mother takes care of me. Who else in your family takes care of you?

My father. Your father.

You are like my father in some ways. What resemblance do you see?

Keywords Why are you looking at me? Replace me with you

Word	Rank	Pattern	Output
alike	10	?X	In what way?
			What resemblance do you see?
are	3	?X are you ?Y	Would you prefer it if I weren't ?Y?
	3	?X are ?Y	What if they were not ?Y?
always	5	?X	Can you think of a specific example?
			When?
			Really, always?
what	2	?X	Why do you ask?
			Does that interest you?
1	1 ((, 11	11 ((• •

if no keyword: "tell me more" "go on"

"Necessity is the mother of invention": "Tell me more about your family".

"Green the adzabak are the a ran four": "What if they were not the a ran four?"

2008年图灵测试Loebner prize所设置任务:

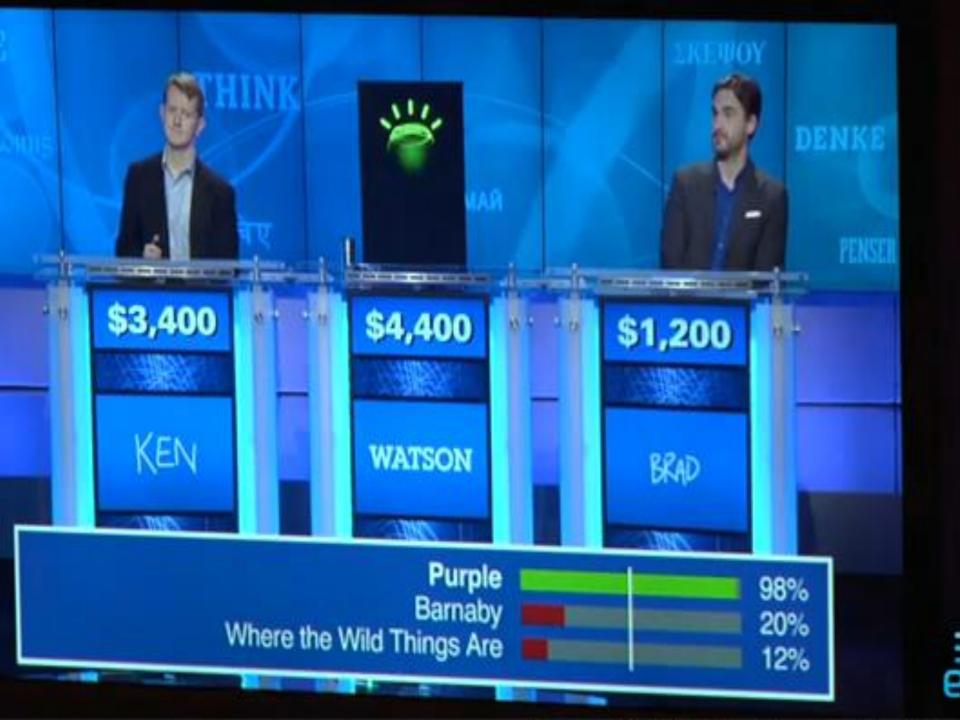
"人工交谈实体"

(Artificial conversational entity)

检验语言交谈能力

冲击图灵测试的重要努力: IBM的Watson (DeepQA)





- * Jeopardy! 是美国著名的益智节目 长青树, 1964 年至今近五十年
- *为何华生要选择这个节目?华生需 要能够听懂包括俚语、双关词、委 婉语、譬喻等人话、也要能够说人 话,并且以最快的速度抢达,还有 该如何选择接下来的问题,以及何 时该拒绝答题(答错问题倒扣分)
- * 图灵测试(人工智能)

Question: Who was presidentially pardoned on September 8, 1974?

One of the retrieved passages is "Ford pardoned Nixon on Sept. 8, 1974."

One passage scorer counts the number of IDF-weighted terms in common between the question and the passage.

Another passage scorer based on the Smith-Waterman sequence-matching algorithm, measures the lengths of The longest similar subsequences between the question and passage (for example "on Sept. 8, 1974").

A third type of passage scoring measures the alignment of the logical forms of the question and passage.

The logical form alignment identifies Nixon as the object of the pardoning in the passage, and that the question is asking for the object of a pardoning. Logical form alignment gives "Nixon" a good score given this evidence.

In contrast, a candidate answer like "Ford" would receive near identical scores to "Nixon" for term matching and passage alignment with this passage, but would receive a lower logical form alignment score.

Question Analysis:

The system attempts to understand what the question is asking and performs the initial analyses that determine how the question will be processed by the rest of the system.

The DeepQA: shallow parses, deep parses, logical forms, semantic role labels, coreference, relations, named entities, and so on, as well as specific kinds of analysis for question answering.

* 云计算支持的知识处理

华生由多台IBM Power7处理器平台服 务器(数千个CPU核),+ IBM专为这 次问答开发的IBM DeepQA软件构成, 本身是不连网的超级计算机。华生能够 透过语音进行分析,并且从背后共15TB 的数据库和知识库(Wordnet, Wikipedia 等)中找出或计算出答案 (1T大约为1万亿个字符)。