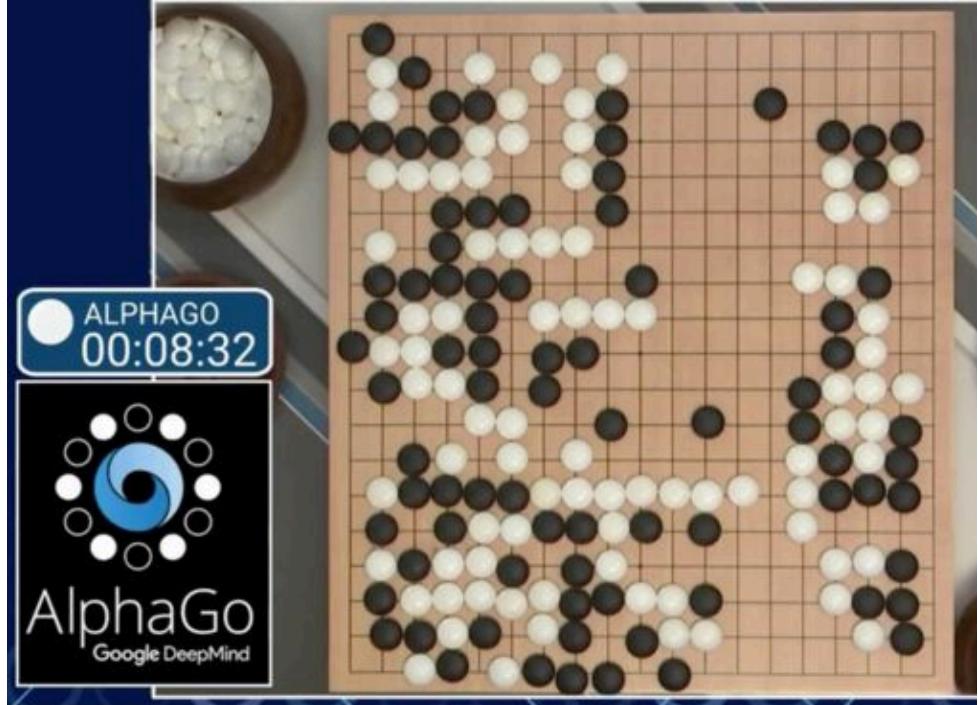


Artificial Intelligence For NLP Lesson-01

开课吧人工智能学院

2020.Jan.05

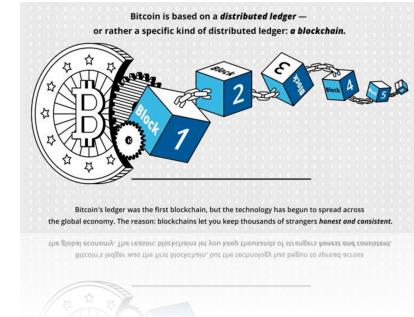
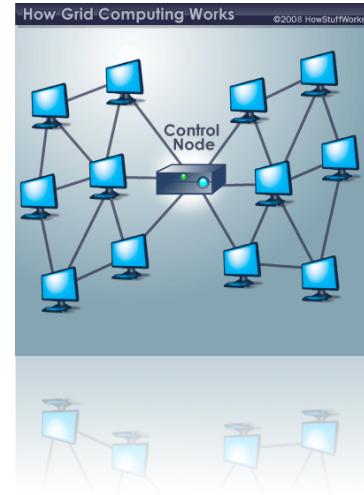
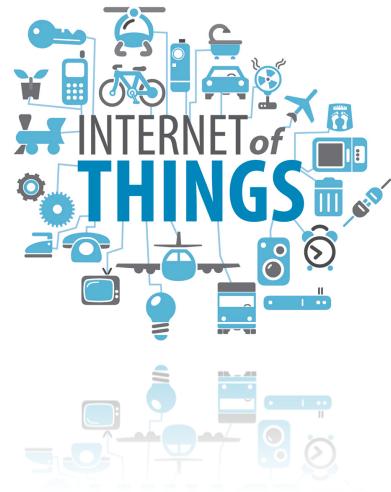
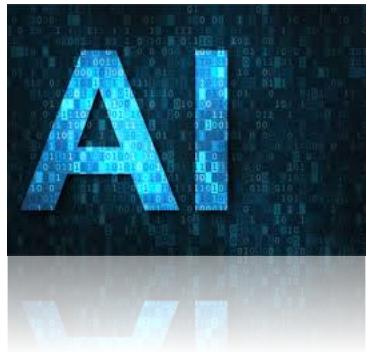


Outline

- i. Course Background
- ii. AI Introduction
- iii. Syntax Tree, Probability Language Model
- iv. 1st Assignment.

1/5 Background

Things change, with time goes by.



Our mentors



李明晓 (NLP课程导师)
Catholic University of
Leuven PHD

曾获国家奖学金赴欧洲交流；曾获欧盟伊拉姆斯全额奖学金。现为丰田欧洲研发中心研究员，与鲁汶大学、剑桥大学、德国马普所、苏黎世联邦理工、捷克布拉格理工合作研发无人驾驶项目。研究方向为自然语言处理，深度学习，多模态对话系统



张希敏 (AI英语提升导师)
ETC全球研究院
普林斯顿大学访问学者

ETC官方认证托福培训师，人事部二级口译员。具备多年会议口译、陪同口译经验，已帮助数百位学生考取托福、雅思高分，成功申请海外名校

- And plenty of references: <https://github.com/Computing-Intelligence/References>

Our Project

项目1、文本自动摘要系统的构建讲解与

技术关键字: TextRank , Sentence Embedding, extractive summarization

文章摘要自动生成

Title 林肯公园主唱 Chester Bennington自杀，年仅41岁

网易娱乐7月21日报道 林肯公园主唱 Chester Bennington于今天早上，在洛杉矶帕洛斯费迪斯的一个私人庄园自缢身亡，年仅41岁。此消息已得到洛杉矶警方证实。洛杉矶警方透露，Chester的家人正在外地度假，Chester独自在家，上帝地点是家里的二楼。一说是一名音乐公司工作人员家里找他时发现了尸体，也有人称是佣人最早发现其死亡。 林肯公园另一主唱克里斯·康奈尔承认了Chester Bennington自杀属实，并对此感到震惊和心痛，称稍后官方会发布声明。Chester昨天还在推特上转发了一条关于曼顿垃圾山的新闻。粉丝们纷纷在该推文下留言，不相信Chester已经走了。 外媒猜测，Chester选择在7月20日自杀的原因跟他极其要好的朋友、Soundgarden(声音花园)乐队以及Audioslave乐队主唱Chris Cornell有关，因为7月20日是Chris Cornell的诞辰。而Chris Cornell于今年5月17日上吊自杀，享年52岁。Chris去世后，他的妻子和两个儿子都表示，他们对他的自杀行为深感悲痛，但没有透露更多细节。据称，Chester曾尝试通过自杀的念头，来表达自己童年时被虐待，导致他终生无法走出阴影，也导致他长期酗酒和嗑药来疗伤。目前，洛杉矶警方仍在调查Chester的死因。 据悉，Chester毒品和酒精斗争多年，年幼时期曾被成年男子性侵，导致常年有轻生念头。Chester生前有过2段婚姻，育有8个孩子。 林肯公园在今年五月发行了新专辑《多一丝曙光One More Light》，成为他们第五张登上Billboard排行榜的专辑。而昨晚刚刚发布新单《Talking To Myself》MV。

随机获取文本 摘要长度: 5 清空 提取

项目3、面向服务的对话机器人的构建

技术关键字: 意图识别 数据检索，专家系统，语言生成

客服小艾

我是谁
我叫开课吧
为什么要叫您爸爸?
你说错了
呵呵，抱歉没有让您满意，我会继续努力学习，希望下次可以让您满意

清空聊天框 发送

项目2、PDF, Word, Excel智能高亮标注

技术关键字: 文档自动化读取, NER(命名实体识别), 依存分析, 文档结构化解析与可视化

representations (see Appendix A.1 for a concrete example). To avoid this problem, we propose to re-parameterize the next-token distribution to target position z_t :

$$p_{\theta}(X_{t+1} = x | \mathbf{x}_{t+1}) = \frac{\exp(e(x)^T g_{\theta}(\mathbf{x}_{t+1}, z_t))}{\sum_{x' \in \mathcal{V}} \exp(e(x')^T g_{\theta}(\mathbf{x}_{t+1}, z_t))}, \quad (4)$$

where $g_{\theta}(\mathbf{x}_{t+1}, z_t)$ denotes a new type of representations which additionally take the target position z_t as input.

Two-Stream Self-Attention. While the idea of target-aware representations removes the ambiguity in target prediction, how to formulate $g_{\theta}(\mathbf{x}_{t+1}, z_t)$ remains a non-trivial problem. Among other possibilities, we propose to “stand” at the target position z_t and rely on the position z_t to gather information from the context. This is in line with the standard Transformer architecture.¹⁰ However, two requirements that are contradictory in a standard Transformer architecture:¹⁰ to predict the token x_{t+1} , $g_{\theta}(\mathbf{x}_{t+1}, z_t)$ should only use the *position* z_t and not the *content* x_{t+1} ; otherwise, the objective becomes trivial since the query token x_{t+1} will always attend to the content x_{t+1} and ignore all other tokens in the content x_{t+1} to provide full contextual information. To resolve such a contradiction, we propose to use two separate representations for this inconsistency:

- The content representation $h_{\phi}(x_{t+1})$, or abbreviated as $h_{\phi,t}$, which serves a similar role to the standard hidden states in Transformer. This representation encodes *both* the context and x_{t+1} itself;
- The query representation $g_{\theta}(\mathbf{x}_{t+1}, z_t)$, or abbreviated as $g_{\theta,z}$, which only has access to the context \mathbf{x}_{t+1} . The position z_t are passed to the content x_{t+1} as query along with x_{t+1} .

Computationally, the first layer query stream is initialized with a trainable vector, i.e., $g_{\theta}^{(0)} = w_0$, while the content stream is set to the corresponding word embedding, i.e., $h_{\phi}^{(0)} = e(x_t)$. For each self-attention layer $m = 1, \dots, M$, the two streams of representations are schematically¹² updated with a shared set of parameters as follows (illustrated in Figures 2 (a) and (b)).

$$\begin{aligned} g_{\theta,z}^{(m)} &= \text{Attention}(\mathbf{Q} = g_{\theta,z}^{(m-1)}, \mathbf{K} = h_{\phi}^{(m-1)}, \mathbf{V} = h_{\phi}^{(m-1)}) \quad (\text{query stream use } z_t \text{ but cannot see } x_{t+1}) \\ g_{\theta,z}^{(m)} &= \text{Content}(\mathbf{Q} = g_{\theta,z}^{(m-1)}, \mathbf{K} = \mathbf{V} = h_{\phi}^{(m-1)}) \quad (\text{content stream use both } z_t \text{ and } x_{t+1}), \end{aligned}$$

where \mathbf{Q} , \mathbf{K} , \mathbf{V} denote the query, key, and value in an attention operation [31]. The update rule of the content representations is exactly the same as the standard self-attention. So during finetuning, we can simply drop the query stream and use the content stream as a normal Transformer (XL). Finally, we can drop the last query representation $g_{\theta,z}^{(M)}$ to compute Eq. (4).

Poison Prediction. While the generation function modeling above (3) has several benefits, it is a much more challenging optimization problem due to the permutation and causes slow convergence in preliminary experiments. To reduce the optimization difficulty, we choose to only predict the last token x_{T+1} and ignore the rest of the tokens. This is in line with the standard Transformer architecture and is consistent with the standard NLP practice.

项目4、在线舆情自动检测系统

技术关键字: 情感分类，爬虫，文本结构化，数据可视化

新闻言论自动提取

随机获取文本 清空 提取

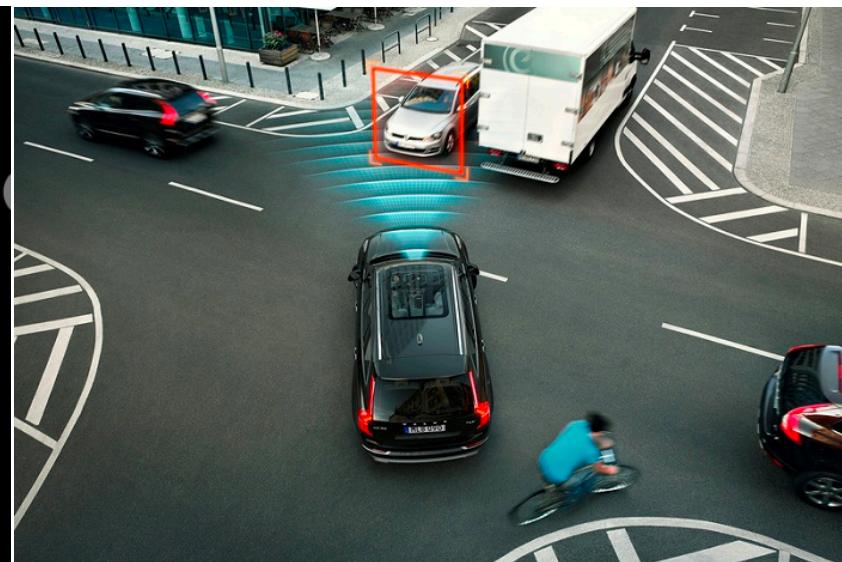
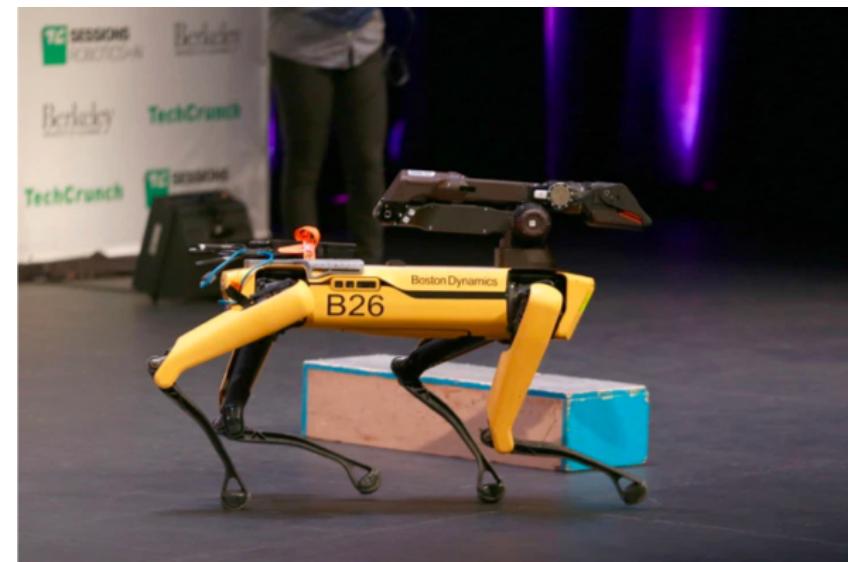
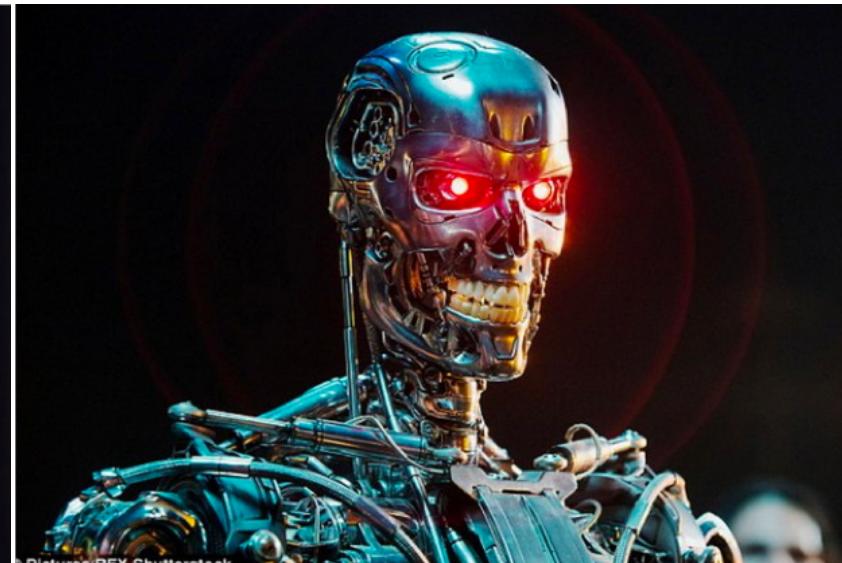
成果展示

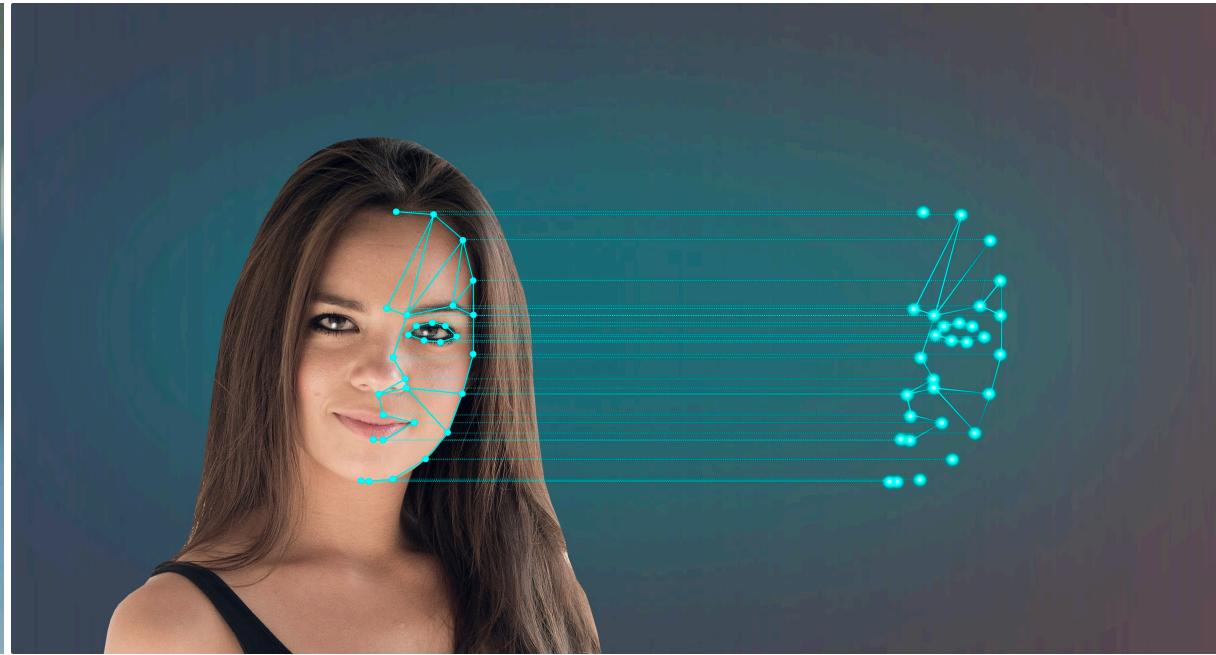
新闻中的言论提取
言论提取树形图
关键句提取

人物	动作	言论
美国商务部工业与安全局	宣布	华为为临时通用许可证(TGL)延长90天时间，该许可证授权在涉及出口、再出口和转让的交易中进行特定的有限交易。根据《出口管理条例》(EAR)的规定，向华为及其非美国分支结构提供受实体清单的限制物品。这是美国政府自今年5月把华为列入管制名单以来，第三次宣布“世纪”华为。为重申部长渠华18日对美国CNBC表示，不管美国会不会延长豁免令，对华为的影响都“非常有限”。
渠华	表示	不管美国会不会延长豁免令，对华为的影响都“非常有限”

2/5 AI Introduction

当我们提起AI的时候





Microsoft

微软五代小冰

“在亿万人之中，我只属于你。”

支持第三方平台

微软小冰是领先的跨平台人工智能机器人。目前，你已可以在以下平台使用她。技术对接需要逐步完成，敬请期待更多第三方平台。

关于小冰的常见问题

初次领养 | 已领养，登录

请输入中国大陆的手机号码

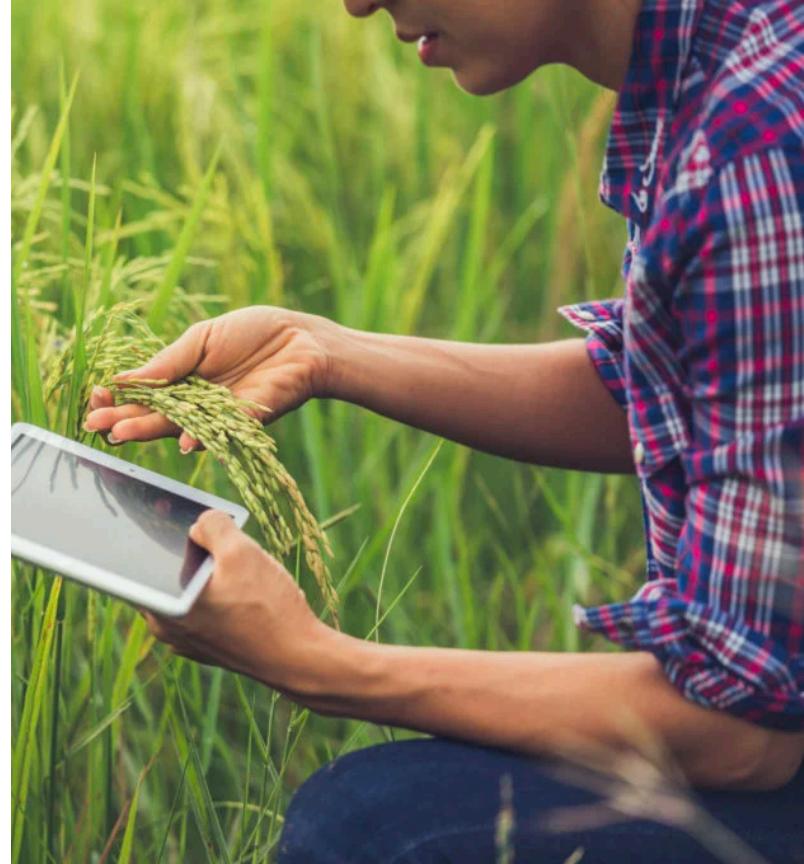
请输入右侧验证码 929E

下一步

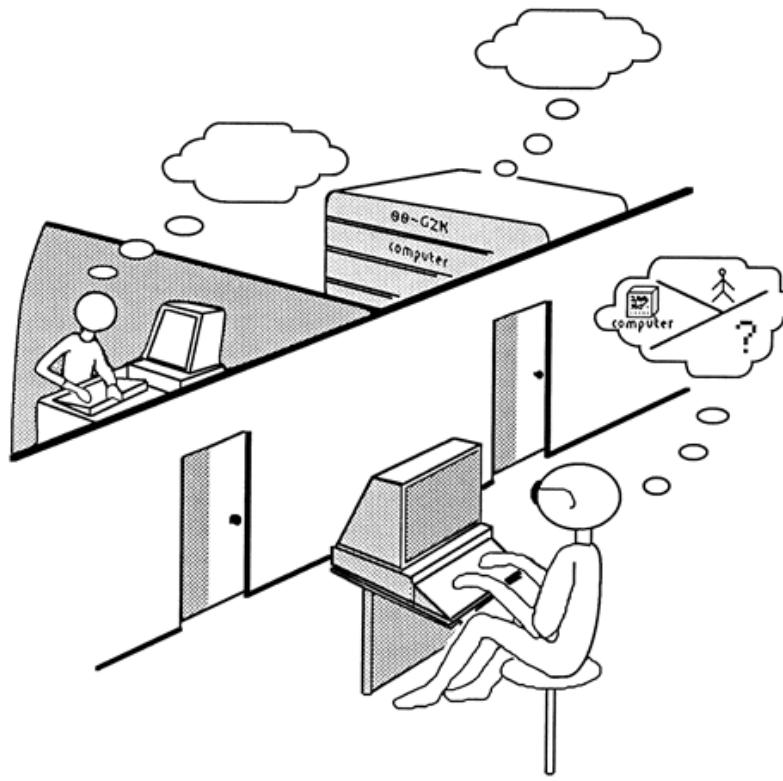
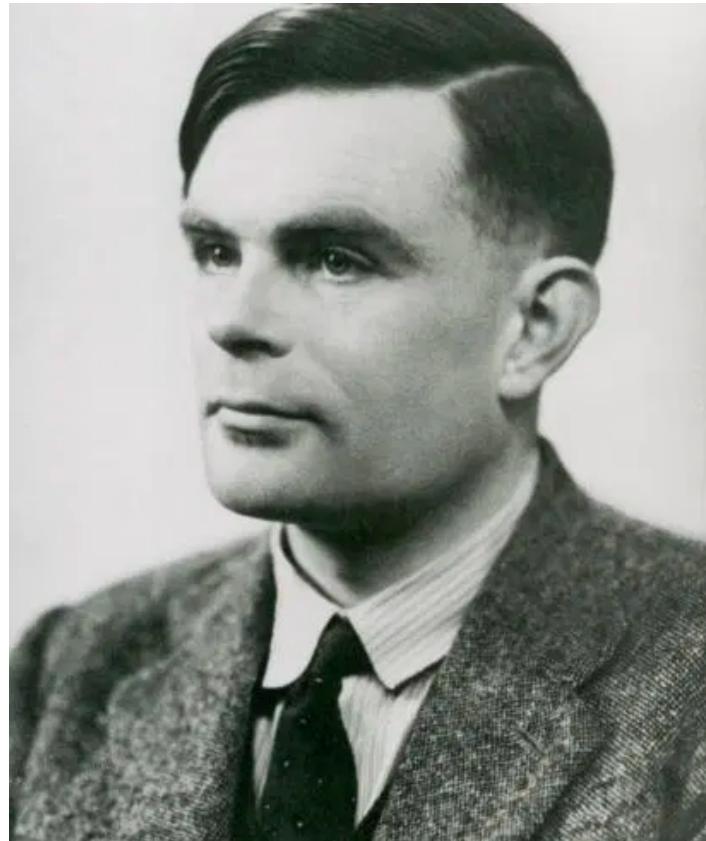
Microsoft

这些是不是AI呢？

- 新德里、上海、杭州正在使用智慧系统降低车辆拥堵
- 农业学家使用更加先进的方式，跟踪农作物的生长
- 智慧家庭能够让我们远程控制自己的电器，甚至电器能够自己控制自己
- 使用新的能源控制系统，摩天大楼的能源消耗可以降低67%



Intelligence? 图灵测试



A. M. Turing (1950) Computing Machinery and Intelligence. *Mind* 49: 433-460.

COMPUTING MACHINERY AND INTELLIGENCE

By A. M. Turing

1. The Imitation Game

I propose to consider the question, "Can machines think?" This should begin with definitions of the meaning of the terms "machine" and "think." The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words "machine" and "think" are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, "Can machines think?" is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

The new form of the problem can be described in terms of a game which we call the 'imitation game.' It is played with three people, a man (A), a woman (B), and an interrogator (C) who may be of either sex. The interrogator stays in a room apart from the other two. The object of the game for the interrogator is to determine which of the other two is the man and which is the woman. He knows them by labels X and Y, and at the end of the game he says either "X is A and Y is B" or "X is B and Y is A." The interrogator is allowed to put questions to A and B thus:

- People or Machine?
- Imitation Game

那AI到底是什么？

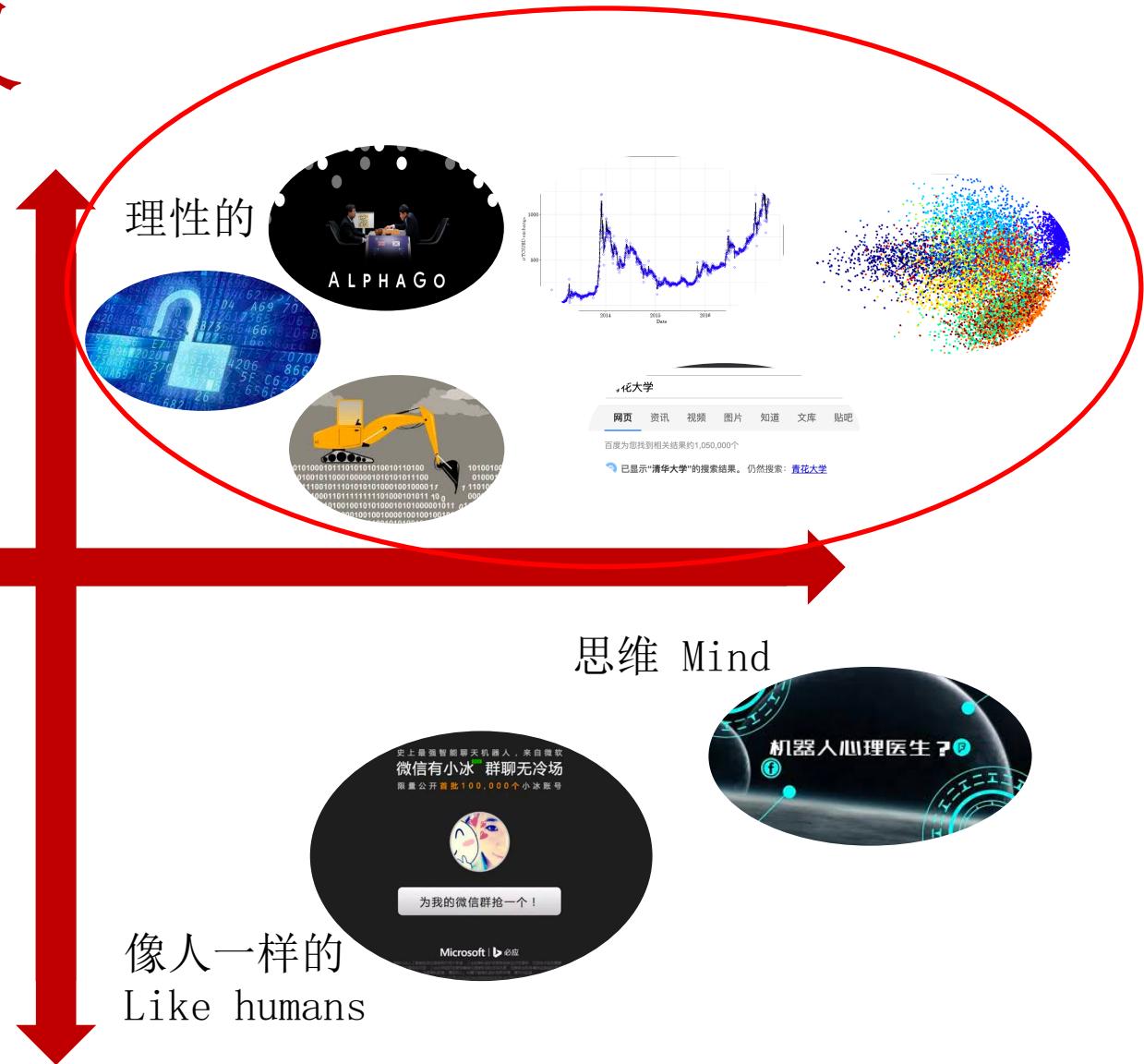


- Artificial Intelligence: Every aspect of learning or any other feature of **intelligence** can be so precisely described that a machine can be made to simulate it. (1956)

一种比较认可的定义



行为 Action

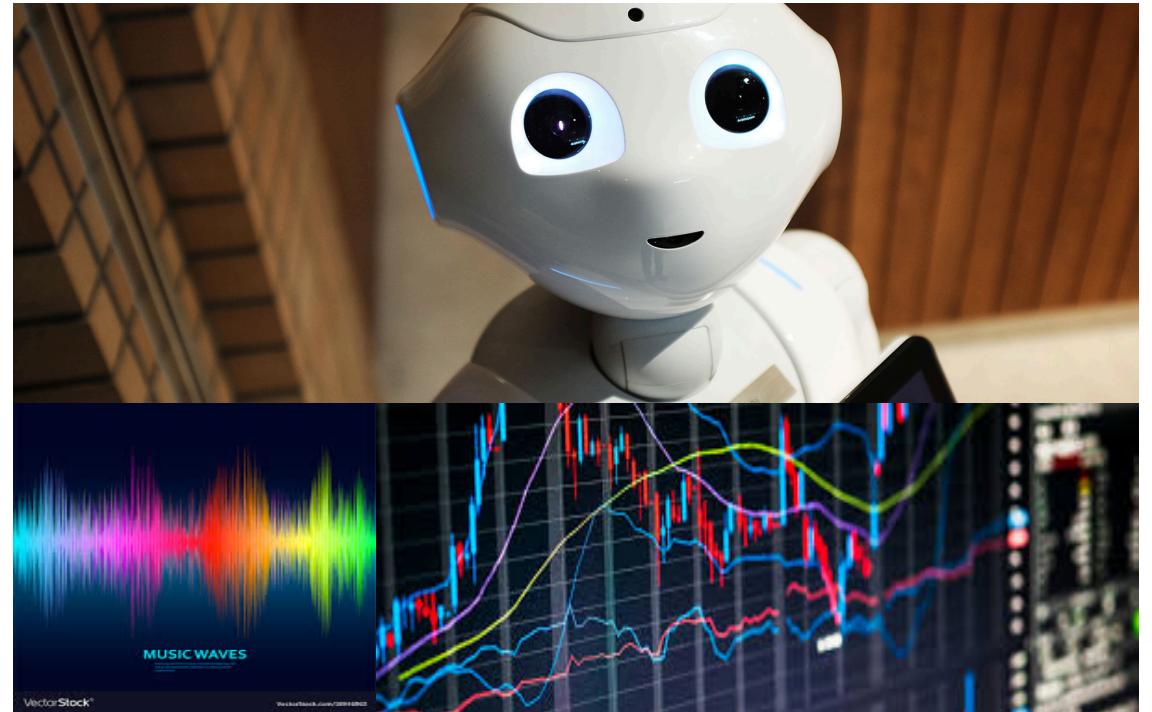


像人一样的
Like humans

Question

- Classify the applications or system based dimensions.

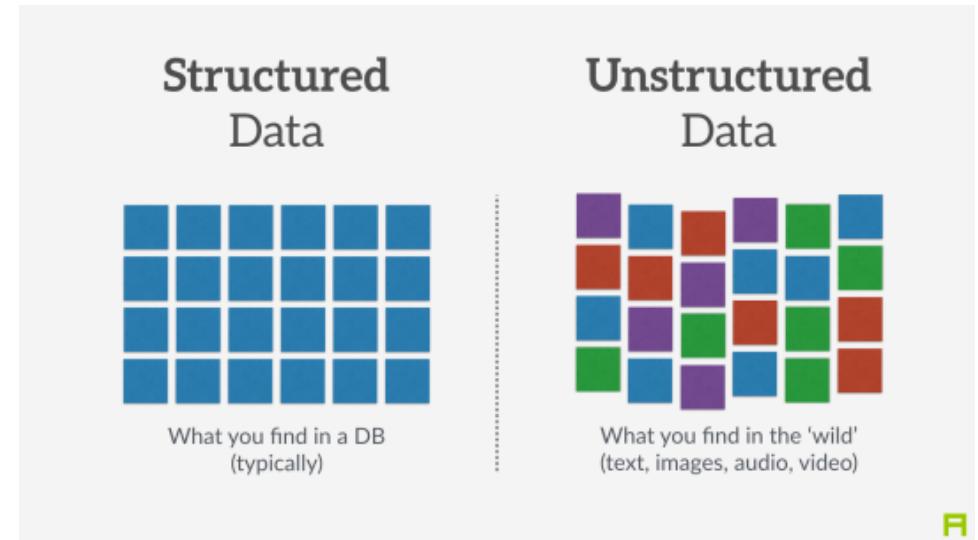
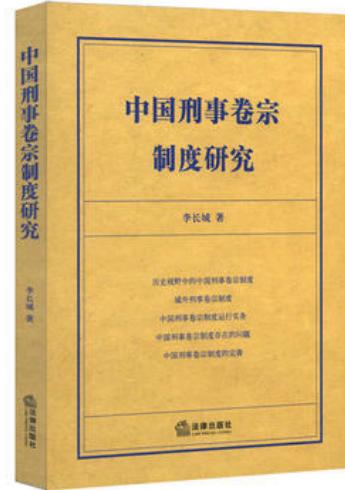
- [] Auto Composition
- [] Voice Recognition
- [] Stock Prediction
- [] Service Robot
- [] Anti Money Laundering



And, can you figure out some more?

Why Natural Language is so hard?

- 1. Text is Logic
- 2. Diversity
- 3. Unstructured
- 4. ...





AI Paradigm

- 1. Rule Based
- 2. Probability Based
- 3. Problem Solving: Search based
- 4. Mathematical or Analytic Based
- 5. Machine Learning (deep learning) Based

1. Rule Based



接待员 = "...."

句子 => 称呼 需要 活动 问题

称呼 => 您 | 你们 | 小伙子

需要 => 要不要 | 想不想 | 需不需要

活动 => 打牌 | 骑马 | 玩射击

问题 => 呢? | 吗 ?

.....

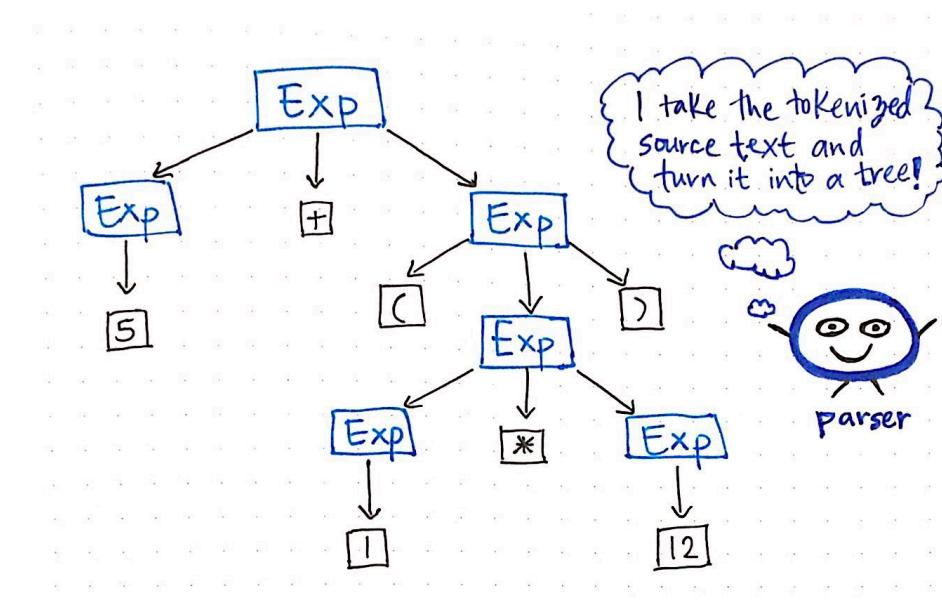
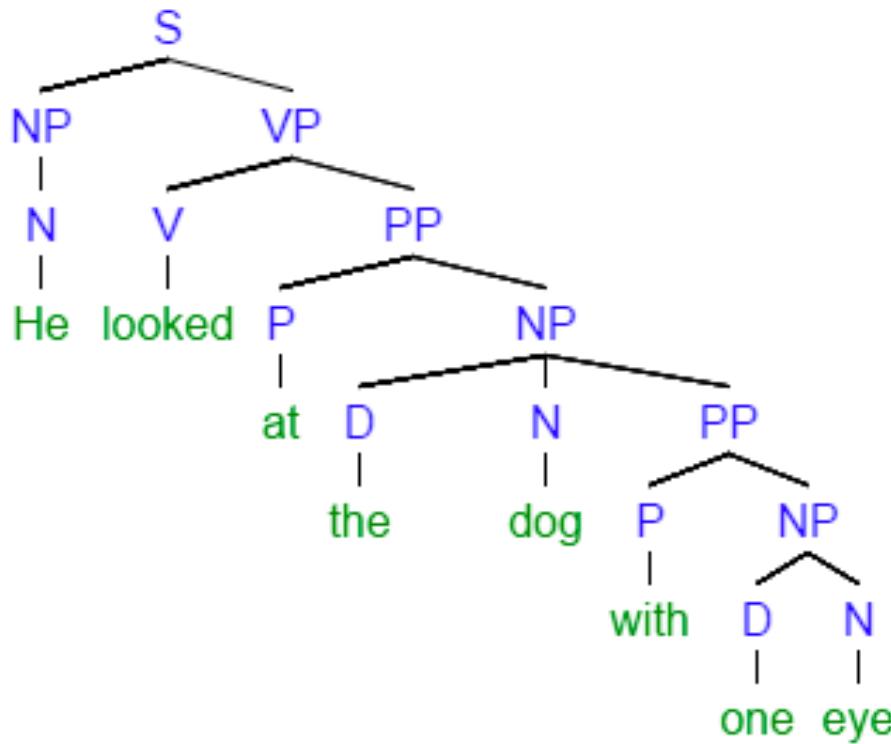
1. 小伙子想不想打牌呢 ?

2. 你们要不要骑马呢 ?

3. ...

4. 您需不需要玩射击呢?

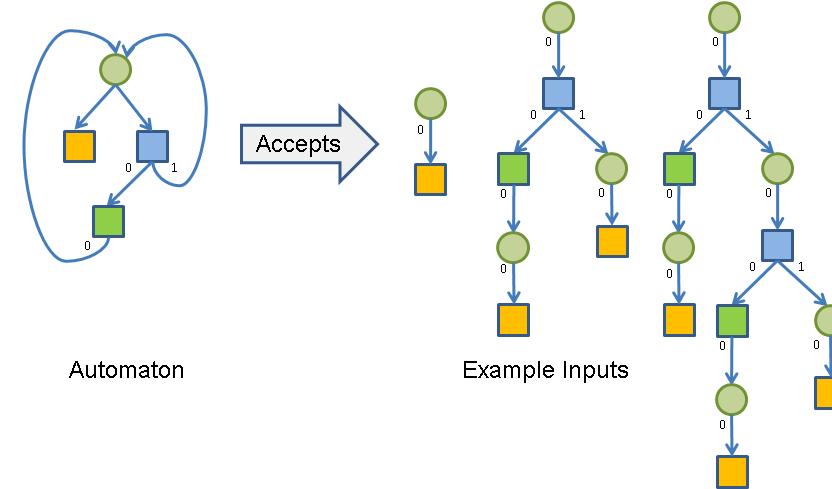
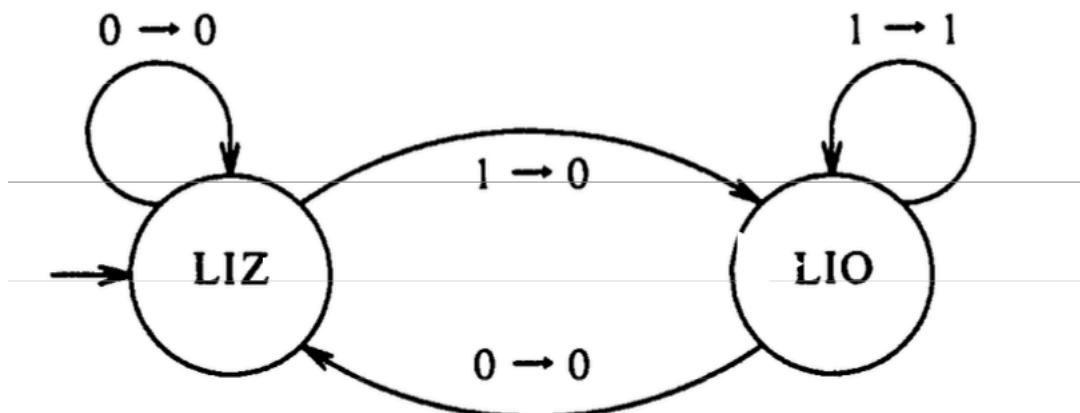
1. Rule Based: Syntax Tree



* First comes the **Lexical analysis** phase, followed by the **syntax analysis** phase, which will generate a **parse tree**.

Automata

- Input: 011010111
- Output: 001000011



- LIZ: Last Input Zero, LIO: Last Input One

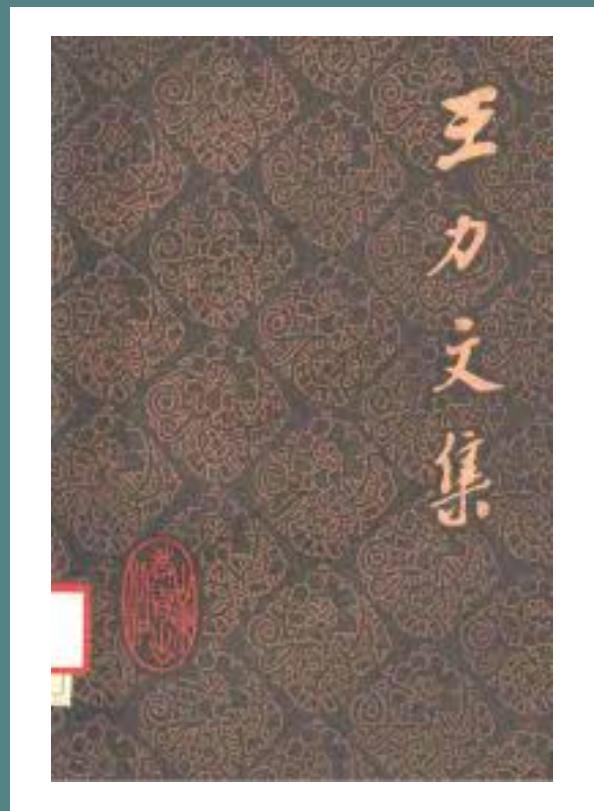
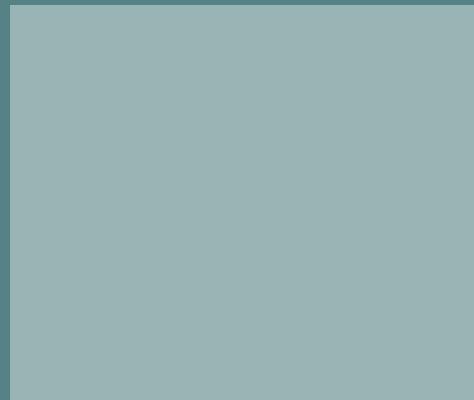
2. Probability Based

- A1. 前天早上吃晚饭的时候
- A2. 前天早上吃早饭的时候

- B1. 正是一个好看的小猫
- B2. 真实一个好看的小猫

- C1. 我无言以对，简直
- C2. 我简直无言以对

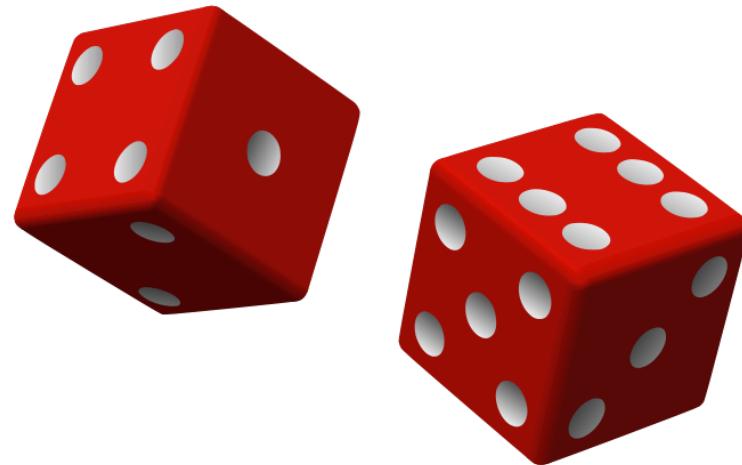
名称	定义	语法特点	类别	举例
(一) 名词	表示人或事物名称的词	①前面可以加数量词(一副对联) ②前面不能加不、很之类的副词(不能 损联、很对联)。 ③后面不能加状态助词“了”(损联了)	①具体名词 ②抽象名词	人、牛、山、水、对联 友谊、立场、观点、思想
(附) 方位词	中表示方向位置的词	常用在名词或名 词性短语的后面		东、西、南、北、前、后、中间、下边
(二) 代词	具有替代或指示 作用的词	①能够替代或指示 替代或各类实词。 ②一般不带修饰成 分	①人称代词 ②指示代词 ③疑问代词	我、你、他、 我们、你们 这、那、这 里、那边 谁、什么、 哪、多少
		①前面可以加副词 (刚走、很想)。	①不及物动词	醒、病、游 行、觉悟



2. Probability Based

- "Every time I fire a linguist, the performance of the speech recognizer goes up"

----- **Frederick Jelinek** (18 November 1932 – 14 September 2010)



2. Probability Based: Language Model

- Language Model:
The probability of
a sentence.

N-gram Models

1-gram (Unigram)

$$P(w_i) = \frac{C(w_i)}{\sum_{\forall k} C(w_k)} = \frac{C(w_i)}{N}$$

of tokens

2-gram (Bigram)

$$P(w_{i+1}|w_i) = \frac{C(w_i, w_{i+1})}{\sum_{\forall k} C(w_i, w_k)} = \frac{C(w_i, w_{i+1})}{C(w_i)}$$

Assignment-01

- 1. Self Review
- 2. Build Sentence Generation System Using Syntax Tree and Language Model
- 3. (Optional) Chat Bot Using Pattern
- 4. (Optional) Reading Turing's Machine Intelligence Paper: See in our github