# 人工智能发展与挑战

吴飞

浙江大学计算机学院

# 提纲

# 记忆是人类认知的基石



| | |
|---|---|
| 知之在人者谓之**知** | **知觉：**人所固有认识外界客观事物本能，如视觉、听觉和触觉等能力 |
| 知有所合谓之**智** | **智慧：**知觉对外界事物的认知 |
| 所以能之在人者为之**能** | **本能：**人身上所具用来处置事物能力 |
| 能有所合谓之**能** | **智能：**对外界所产生的认知和决策 |

《荀子. 正名》

Alan Baddeley, Working memory: looking back and looking forward, Nature Reviews Neuroscience 4, 829–839, 2003

# 记忆单元之间及其与环境的交互是提升智能能力的重要途径



**工作记忆**
(直觉、顿悟、因果等推理)
持续时间：< 30 sec

注意力

交互

**环境：强化学习**

**瞬时记忆**
(多通道感知)
持续时间：< 5 sec

**长期记忆**
(先验、知识等)
持续时间： 1 sec--lifelong

Your ability to remember something doesn't just depend on the strength of the memory, it depends on the state that you're in

Human-level control through deep reinforcement learning, *Nature*, 518:529–533, 2015
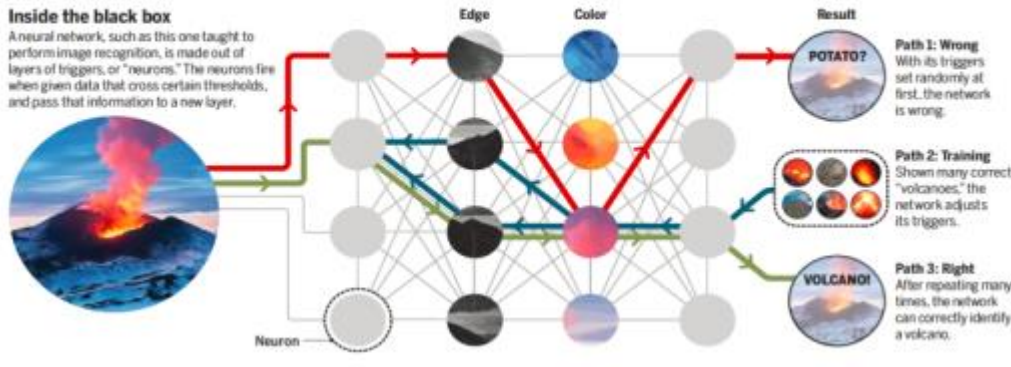
# 模型中的记忆(internal memory)



Opening up the black box
Loosely modeled after the brain, deep neural networks are spurring innovation across science. But the mechanics of the models are mysterious: They are black boxes. Scientists are now developing tools to get inside the mind of the machine.
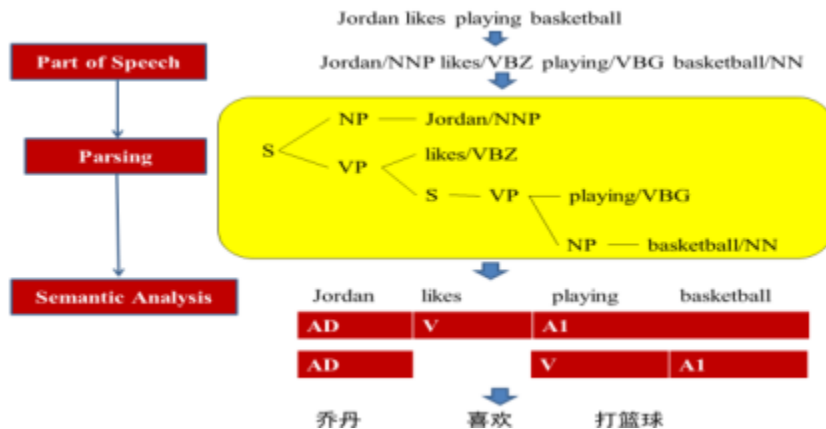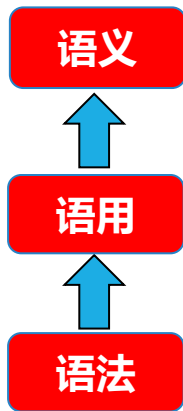
Inside the black box
A neural network, such as this one taught to perform image recognition, is made out of layers of triggers, or "neurons." The neurons fire when given data that cross certain thresholds, and pass that information to a new layer.
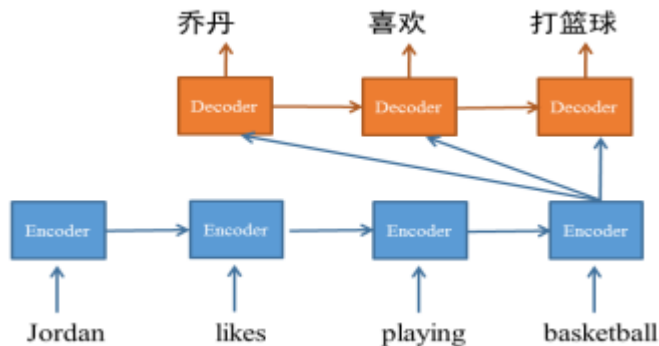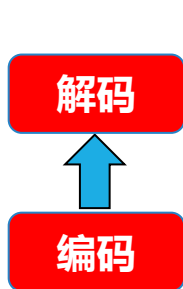
Edge   Color   Result

Path 1: Wrong
With its triggers set randomly at first, the network is wrong.

Path 2: Training
Shown many correct "volcanoes," the network adjusts its triggers.

Path 3: Right
After repeating many times, the network can correctly identify a volcano.

Neuron

POTATO?
VOLCANO!

- 神经网络中的赫布理论（Hebbian theory）：突触可塑性原理，即突触前神经元向突触后神经元的持续重复的刺激可以导致突触传递效能的增加
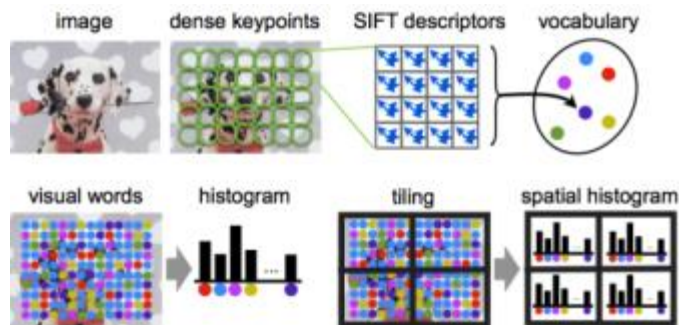- 神经元间权值信息既执行计算任务、又承载记忆功能。
- 记忆被新知识覆盖导致多任务难以实行、模型泛化能力不强、结果解释性弱

- 语义
- 语用
- 语法

Jordan likes playing basketball

Jordan/NNP likes/VBZ playing/VBG basketball/NN

- Conditional random field
- Hidden Markov Model

- 解码
- 编码

- word2vec
- Paragraph2vec
- Node2vec
- path2vec

# 从分段学习到"端到端"学习：以视觉分类和理解为例



Detection → Description → Classification → Identification → Understanding



黑盒子
（多粒度语义在中间层丢失）

- Two stream model
- Region-CNN
- Mask R-CNN

输入端

端到端学习，利用卷积、池化和误差反向传播技术，强调特征学习

输出端：高层语义
（what）

# 从分段学习到"端到端"学习

**分段学习**     每个阶段可灵活引入先验、经验与知识，但并不知所引入信息的合理性

**端到端学习**     数据说话（你见或不见我，我就在那里，不悲不喜），但缺乏了人类语言可表述的"interpretability"

# 从图灵机到神经图灵机：利用外在记忆体中的知识



A.M.Turing, On Computable Numbers with an Application to the Entscheidungsproblem, *Proceedings of the London Mathematical Society*, Ser. 2, Vol. 42, 1937



Alex Graves, et al., Hybrid computing using a neural network with dynamic external memory, *Nature* 538, 471–476,2016

# 利用外在记忆体中的知识进行可计算推理



弦外之音、画外之意：
利用外在记忆体的深度神经推理

记忆的激活：
哈希索引的相似度搜索

记忆的可塑性维持、用尽废退：
强化学习中记忆的形成、巩固和遗忘

- Jason Weston(facebook), et al., Memory Networks, arXiv:1410.3916
- Retrieval induces adaptive forgetting of competing memories via cortical pattern suppression, Nature Neuroscience, 18, pages582–589 (2015)
- A neural algorithm for a fundamental computing problem, Science , 358, Issue 6364,793-796, 2017

# 外在记忆体中的知识的不同利用方式



one shot learning



Figure 1: A brief illustration of M-MCTS. When a le state $s$ is searched, the feature representation $\phi(s)$ generated, which is then used to *query* the memo. based value approximation $\hat{V}_{\mathcal{M}}(s)$. $\hat{V}_{\mathcal{M}}(s)$ is used to update $s$ and all its ancestors according to equation (9), as indicated by the red arrows in the figure.



Figure 1: A visualization of two time steps of the neural map.

Structured memory

- Adam Santoro,et al., Meta-Learning with Memory-Augmented Neural Networks
- Chenjun Xiao, et al., Memory-Augmented Monte Carlo Tree Search
- Emilio Parisotto, et al., Neural Map:Structured Memory for deep reinforcement learning

# 有效利用当前数据、已有知识和未知交互

| Shallow models | Deep models | 备注 |
|---|---|---|
| Language model | Neural language model | |
| Bayesian Learning | Bayesian deep learning | ● 不只是单纯追求将浅层模型拓展到深层模型。 |
| Turing Machine | Neural Turing Machine | |
| Reinforcement Learning | Deep Reinforcement Learning | ● 更为重要的是，在这个转变过程中，巧妙融合数据、知识和交互经验，多种手段和方法的综合利用。 |
| Generative Model | Deep Generative Model | |
| X | Deep or Neural + X | |

# 有效利用当前数据、已有知识和未知交互的挑战

| 知识表达模型 | → | ☐ 逻辑、描述、事实型知识的表示方法，从离散符号到分布式向量表达，为深度神经推理打下基础 | **数据** |

| 记忆激活机制 | → | ☐ 编码（形成新的记忆）以及记忆检索（回想过去）的机制，即实现模式分离和模式完成。 | **记忆** |

| 交互更新手段 | → | ☐ 通过人-机交互、机-机交互等形式，或者利用认知模型、或者借助自我博弈，进行知识更新 | **交互** |

对推理过程逐渐松绑，使推理逐步走向对思维广泛模拟：**跨媒体综合推理**

# 跨媒体综合推理

知识推导、数据驱动和交互反馈等学习方法缺乏像人类一样能接受多通道输入、综合利用多种推理模式进行协同推理的计算框架。为此，跨媒体综合推理思想被提出。



模拟高级思维活动，侧重知识表示，但难以拓展

从大量实例中学习，侧重感知能力，但过程难以理解与干预

模拟低级的认知过程，侧重自我适应，极大的依赖于策略

# 跨媒体综合推理



有机协调"知识指导下的演绎"、"数据驱动中的归纳"和"行为强化内的规划"等理论模型和方法手段，建立知识、数据和反馈于一体的人工智能理论和模型

# 仅有短期记忆的人生

英国指挥家Clive Wearing：

- 因 herpes simplex virus（单纯疱疹病毒）侵蚀大脑hippocampus（海马体）而患上 anterograde amnesia（顺行性遗忘症）
- 海马体是将短期记忆传递成长期记忆的重要器官




他的记忆，和金鱼一样，只要"七秒"就会消失。他的生命是一段又一段的空白，没有过去，没有未来。

要证明现在致力于构建的通用智能是可行的，人类大脑是现有唯一证据，因此把神经科学作为新算法的灵感来源是有意义的。

# 提纲

1、 记忆驱动的智能计算

2、 可计算社会学

3、若干挑战

# 人类进入信息空间(Cyberspace)、物理世界 (Physical)和人类社会(Society)三元空间融合时代(CPS)



人类社会

群体行为
与社会结构

信息空间

物理世界

事件和主题

互联网信息
百度每天搜索请求
100亿；处理数据量
100PB

空间位置信息

空间位置信息
每年增加1亿GPS设备

….The 'virtual' world and the 'real' world complement each ….

互联网数据是对现实世界的 "自我映照"，与现实世界相互补充。

Self-reflection, online

Some scientists might not like the persona they see when they look online. But they can do something about it.

Viginia Gewin, Self-reflection, *Nature* , 471,667-669,2011, March

# 社会学研究：观察与实验



**六度分割理论**
Stanley Milgram. The small world problem. Psychology Today, 2(1):60–67, 1967. Cited by 6967 (as of May 2016)



**弱链接优势**
Granovetter, Mark S. The strength of weak ties. *American journal of sociology* 78.6 (1973): 1360-1380. Cited by 43186 (as of April 2017)



**邓巴数字即150定律**：人类智力所允许稳定网络社交人数150人，Facebook社区用户平均好友人数是120人
Robin I. M. Dunbar. Neocortex size as a constraint on group size in primates. *Journal of Human Evolution* 22 (6): 469–493. Cited by 1393 (as of May 2016)



**结构洞：位置比关系更重要**
Ronald S. Burt. *Structural Holes: The Social Structure of Competition*. Cambridge: Harvard University Press. 1995. Cited by 20207 (as of May 2016)

# 社会学研究：观察与实验

| 时间（年） | 人　　物 | 事　　件 |
|---|---|---|
| 1793 | Euler | 七桥问题 |
| 1959 | Erdǒs和Rényi | 随机图理论 |
| 1967 | Milgram | 小世界实验（六度空间） |
| 1973 | Granovetter | 弱连接的强度 |
| 1998 | Watts和Strogatz | 小世界模型 |
| 1999 | Barabási和Albert | 无尺度网络 |

迄今为止最大规模小世界验证：2011年，Facebook与意大利米兰大学公布了关于Facebook中用户之间六度分离验证的结果。2007年以来，Facebook上两个用户之间的平均距离仅为4.74，任何两个用户之间间隔不超过5度的概率99.6%。注意2011年5月数据包括了Facebook上大约7.21亿活跃用户以及69亿朋友关系链接。

Lars Backstrom, et. al., *Four Degrees of Separation*, http://arxiv.org/abs/1111.4570

美国CMU的Tom Mitchell教授在《Science》发表名为"*Mining our Reality*"的文章指出：….从相互关联的网络海量数据中寻找"蛛丝马迹"，更有利于理解现实世界…

….use computers to mine data…. machine learning algorithms have helped to analyze historical data, often revealing trends and patterns too subtle for humans to detect.

Tom M., Mitchell, Mining our Reality, *Science*, 326,2009,1644-1645

# 可计算社会学：应用计算模型和方法来理解和预测现实世界



…a computational social science is emerging that leverages the capacity to collect and analyze data with an unprecedented breadth and depth and scale

美国Harvard University的David Lazer等人在《Science》杂志发表论文也认为需建立可计算模型(*Computational Model*)来理解现实世界。

David Lazer, *et. al.,* Computational Social Science, *Science,* 323,721-723, 2009

# 计算思维成为可计算社会学研究的思维方式

**在人工智能、大数据、互联网等支撑下所形成的计算思维与实验思维和理论思维同等重要，人工智能正成为一种通识教育，渗透进入其他知识技术教育之中。**



Wing, Jeanette M.,
Computational thinking,
*Communications of the ACM*,49
(3): 33,2006

计算思维是运用计算机科学的基础概念进行问题求解、系统设计、以及人类行为理解等涵盖计算机科学之广度的一系列思维活动

| 思维方式 | 机制 | 手段 |
|---|---|---|
| 实验思维 | 实验→观察→发现、推断与总结 | **观察与归纳** |
| 理论思维 | 假设、预设→定义/性质/定理→证明 | **推理与演绎** |
| 计算思维 | **设计，构造与计算** | **设计与构造** |

# MIT通过学科交叉重塑人才培养模式：计算机和人工智能为引擎

MIT reshapes itself to shape the future

MIT于2018年10月16日宣布，投资十亿美元建设新的计算学院，致力于将计算机技术以及人工智能纳入所有研究领域，同时寻求改变计算和人工智能相关的公共政策和道德研究教育。

make AI part of every graduate's education
AI-focused College

- **2019年秋季学期正式启动新学院**
- **将计算机和人工智能力量带到 MIT所有学习领域**
- **创造 50 个跨越新学院和其它学系的新教师职位**
- **为 MIT五个学院建设公共合作结构，用于计算机和人工智能领域的教育、研究和创新**
- **改变关于计算机和人工智能公共政策和道德方面教育和研究**

# 哈佛大学商学院在案例分析中重视人工智能

# 谷歌预测流感方法中的人工智能

□ Google推出 "谷歌流感趋势" 项目，通过分析搜索单词来预测流感的爆发。与美国疾病控制和预防中心提供的报告对比，其对追踪疾病的精确率达到97%-98%



Google预测的流感爆发与实际情况对比



Google预测的流感分布图

Detecting influenza epidemics using search engine query data, *Nature* 457, 1012-1014 (19 February 2009)

# 谷歌预测流感方法中的人工智能

☐ 谷歌流感预测（Google Flu Trends，GFT）所预测流感样病例门诊数超过了与美国疾病预防控制中心（Centers for Disease Control and Prevention，CDC）根据全美各实验室监测报告得出的预测结果的两倍。

☐ "大数据傲慢"(Big Data Hubris): 认为大数据可以完全取代传统的数据收集方法，而非作为后者的补充

GFT的平均绝对偏差为0.486，CDC滞后模型的平均绝对偏差为0.311，GFT与CDC相结合的平均绝对偏差为0.232



The Parable of Google Flu: Traps in Big Data Analysis, *Science* 14 Mar 2014

# 小数据 Versus 大数据

☐ 1936年《文学文摘》（The Literary Digest）预测总统选举：民主党人艾尔弗雷德 兰登（Alfred Landon）与时任总统富兰克林 罗斯福（Franklin Roosevelt）

☐ 《文学文摘》给自己的订户邮寄了1000万份调查问卷，根据回收到的240万份回执，预测艾尔弗雷德 兰登将会以55比41的优势击败富兰克林 罗斯福赢得大选。

☐ 美国数学家、抽样调查方法创始人、民意调查的组织者乔治 盖洛普（George Gallup）仅通过3000人问卷调查，得出了准确得多的预测结果。

☐ 良好的"小数据"可战胜"含噪和偏差"的大数据，小数据可做到"以小见大"、"一叶知秋"。数据质量对于机器学习预测结果影响甚大。

# 人工智能与学科交叉中存在的"数据之学科鸿沟"

医学数据

法学数据

心理学数据

人文社科数据



算法

模型

算力

数据化知识 ≠ 数据驱动

# 浙江大学学者用7亿条通话记录回答：移民融入成功了吗

通过对5400万用户的6.98亿条通话记录的学习，人工智能建立起了一项看似天马行空的关联：通话记录与城市移民状态

● 两周通话记录：是走还是留？群越多越杂，越能融入、参谋城市规划



计算社会学：通话记录"洞悉"城市移民状态



求是印象

科学封面（第21期）——计算社会学：通话记录"洞悉"城市移民状态

(a) Log overall average probability.　(b) Log odds ratio for locals.　(c) Log odds ratio for settled migrants.　(d) Log odds ratio for new migrants.

上海市居民地域分布

Yang Yang（杨洋）, Chenhao Tan, Zongtao Liu, Fei Wu, and Yueting Zhuang. Urban Dreams of Migrants: A Case Study of Migrant Integration in Shanghai. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence(AAAI)*, 2018

# 提纲

"first wave" (rule-based)

**符号主义人工智能为手段**
**以语言和可描述信息处理为主**

"second wave"
(statistical-learning-based)

**以模型假设的机器学习为手段**

**在数据建模基础上、从数据中学习模式**

"third wave" AI
theory and applications

**以自适应和推理为核心目标**
**human-like communication and reasoning capabilities**

**"first wave" (rule-based)**

**符号主义人工智能为手段**
**以语言和可描述信息处理为主**

# IBM"沃森"的推理

**主持人问：Kathleen Kenyon**'s excavation of this city mentioned in **Joshua** showed the walls had been repaired 17 times. （Kathleen Kenyon 对这个在《《圣经•约书亚记》》中提到的城市的发掘表明，该城的城墙曾被修复17次）
**沃森回答：**What is Jericho （耶利哥城是什么？）

答案排序：**耶利哥（Jericho）97%**、耶路撒冷 （Jerusalem）42%、拉吉（Lachish）7%

● Kathleen
● Kenyon
● Kathleen Kenyon
● Joshua

"second wave"
(statistical-learning-based)

**以模型假设的机器学习为手段**
**在数据建模基础上、从数据中学习模式**

白宫版"潜伏"，数据分析锁定副总统

# 第三波人工智能："third wave" AI

| 可观测性问题 | What if we see A (what is?) | $P(y\mid A)$ |
|---|---|---|
| 决策行动问题 | What if we do A (what if?) | $P(y\mid do(A))$（如果采取A行为，则B真） |
| 反事实问题 (Counterfactual) | What if we did things differently | (why?) $P(y'\mid A)$（如果A为真，则B将不同） |
| Options: with what probability | | |

关联(association)：
直接可从数据中计算得到的统计相关

介入(intervention):
无法直接从观测数据就能得到关系，如"某个商品涨价会产生什么结果"

反事实(counterfactual)：
某个事情已经发生了，则在相同环境中，这个事情不发生会带来怎样的新结果

2 billion dollars over the next five years

# Machine Common Sense (MCS)

Department of Defense
Defense Advanced Research Projects Agency
Information Innovation Office
Oct 19, 2018

DARPA is soliciting innovative research proposals in the area of machine common sense to enable Artificial Intelligence (AI) applications to understand new situations, monitor the reasonableness of their actions, communicate more effectively with people, and transfer learning to new domains. Proposed research should investigate innovative approaches that enable revolutionary advances in science, devices, or systems. Specifically excluded is research that primarily results in evolutionary improvements to the existing state of practice.



见一叶落，而知岁之将暮；审堂下之阴，而知日月之行，阴阳之变；见瓶水之冰，而知天下之寒，鱼鳖之藏也

《淮南子说山训》

# Machine Common Sense (MCS)

- **Foundations of Human Common Sense**: a service that learns from experience, like a child, to construct computational models that mimic the core domains of cognition 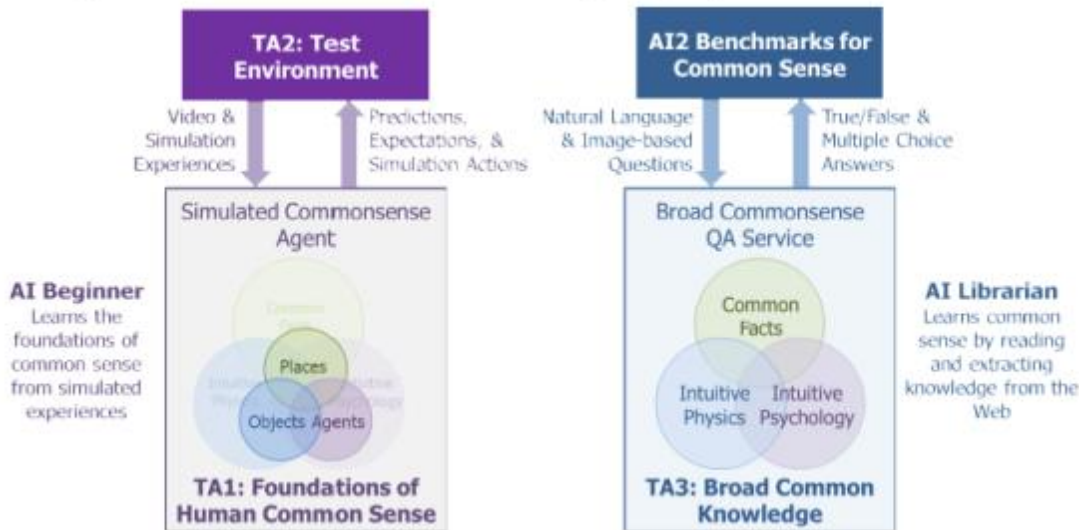for objects (intuitive physics), places (spatial navigation), and agents (intentional actors). These models will be evaluated against the cognitive development milestones as evidenced in developmental psychology experiments with human children from 0-18 months old.

- **Broad Common Knowledge:** a service that learns from reading the Web, like a research librarian, to construct a commonsense knowledge repository capable of answering natural language and image-based questions about commonsense phenomena. This service will mimic the general knowledge of an average American adult in 2018, and be evaluated against the Allen Institute for Artificial Intelligence (AI2) Common Sense Benchmarks.



常识是人工智能的黑洞
Common sense is the dark matter of artificial intelligence

# 挑战(1)：人类常识知识支持下顿悟与直觉

**常识知识的获取**
　　自然语言、视觉和听觉等
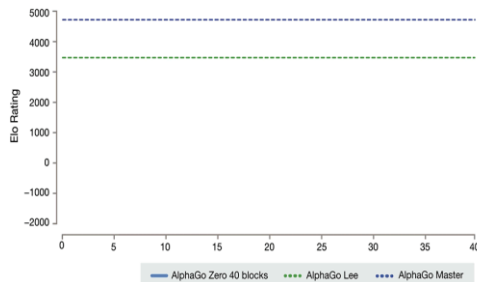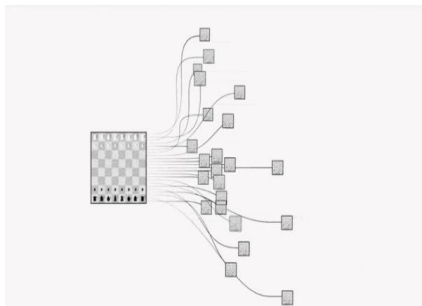
**常识知识的记忆**
　　存储与激活等

**常识知识的利用**
　　顿悟与直觉推理

- 威廉和他的三岁宝贝儿子谁身高更高
- 针扎入胡萝卜，它会留下小孔，请问这里的"它"指的是什么

# 挑战 (2)：元学习(meta-learning)理论："利用与探索"的 多巴胺神经计算机制

多巴胺（dopamine, 神经传导物质，快乐信号，2000年诺贝尔医学奖）在人脑学习过程中"周旋"于已得回报和预期回报

starting *tabula rasa*
(一张白纸绘蓝图)

多巴胺奖励信号不仅仅是对权重进行调整，它还负责编码、传递抽象任务与规则结构的重要信息，从而使快速任务适应成为可能。



AlphaGo Zero: 深度学习、蒙特卡洛树搜索、强化学习，
40天训练、2900万次自我对弈，89:11击败AlphaGo
Master



Prefrontal cortex as a meta-reinforcement learning system, Nature Neuroscience,2018

通过观察成人努力实现目标的行为和反应，婴幼儿的坚持性提高，并且努力去完成一个新的和困难的任务

Julia A. Leonard, Yuna Lee, Laura E. Schulz. (2017, September 29). Infants make more attempts to achieve a goal when they see adults persist, *Science*, 10.1126/science.aan2317.

**神经启发**　**认知可解**　**计算可行**

- 基于结构和模型的学习是如何发生在脑中的
- 为什么多巴胺本身就编码模型信息
- 前额叶皮质的神经元是如何调节学习信号的

**道之运行轨迹为元**



General schematic, solutions to complex problems in artificial intelligence and nature (brains). Higher cognitive functions continuously interact between them and with reinforcement learning to drive generalization and learning from small sample

# 挑战(3)：大脑的贝叶斯模型和因果推理

认知科学：大脑将外界信息计算得到概率赋予其所主动构建的"世界假设"，并加以调整优化。即比较"现有经验"和"未来期望"。

- 大脑如何在神经回路尺度上

  进行模型构建和推断，将自

  底向上与自上而下相互结合

- 贝叶斯模型如何整合感知、

  认知、理性和意识

# 挑战(3)：大脑的贝叶斯模型和因果推理

认知科学：大脑不是在被动等待外界输入信息，而是主动构建尚未"显性描述"的世界的假设。

**似然概率**      **先验概率**

$$P(疾病|症状) = \frac{P(症状|疾病) \times P(疾病)}{P(症状)}$$

# 挑战(4)：控制论与人工智能和博弈对抗

- Interesting fact: John von Neumann
  - Father of computer science + Father of game theory
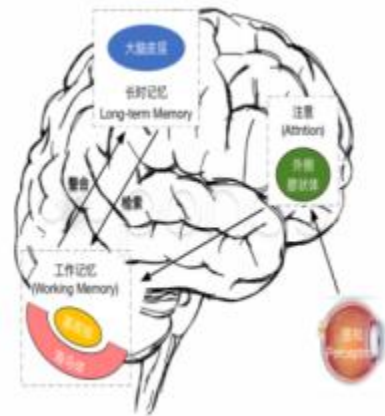- Game theory deeply interacts with computer science
  - Theoretical computer science: algorithmic game theory
  - AI: multiagent systems, game playing, human-robot interaction, machine learning, internet advertising
  - Web and Internet: internet economics, sharing economics
  - Distributed system: Blockchain

| 计算博弈策略 |
| --- |
| (智能体对决) |
| 博弈机制设计 |
| (拍卖、城市资源分配) |



"人工智能"与"控制论"词频对比
（引自《人工智能简史》）

# 挑战(4)：控制论与人工智能和博弈对抗

- 健壮人工智能：

  - 自主无人系统在新环境中须通用性和自适应性要求极高

  - 极端复杂条件对人工智能"健壮性"要求极高

  - 人工智能体系的综合程度史无前例（云、网、端、算法、模型、人）

  - 人是"智能回路"中最重要因素

对未知的未知问题建模与环境自适应能力

健壮人工智能的可解释性和可验证性

复杂动态和开放环境下的博弈对抗

多智能体在非完全信息条件下的协作

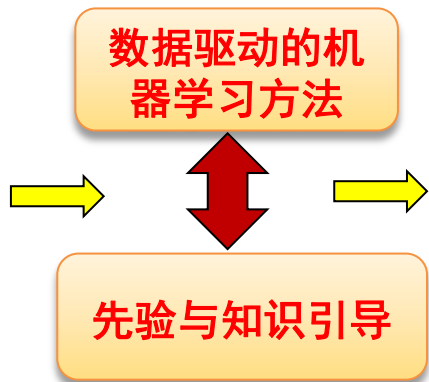# 多学科交叉促进人工智能发展

人工智能

观察脑、认识脑、仿真脑、控制脑、人文化脑

需要新数学模型

需要可计算的认知模型

需要控制、博弈等方法支持

需要人文社科支持

● 常识知识

● 元学习

● 贝叶斯大脑

● 因果推理

● 控制、博弈对抗、AI

# 可解释、更通用和适应性强的人工智能
## 数据利用、知识引导与能力学习



**数据驱动的机器学习方法**

**先验与知识引导**

从**浅层**计算到**深度神经推理**

从**单纯依赖**于数据驱动的模型到数据驱动与知识引导**相互结合**

从**领域任务**驱动智能到更为**通用条件**下的强人工智能（从经验中学习）

Yueting Zhuang, Fei Wu, Chun Chen, Yunhe Pan, Challenges and Opportunities: From Big Data to Knowledge in AI 2.0, *Frontiers of Information Technology & Electronic Engineering*, 2017,18(1):3-14

"我们必须知道，我们必将知道"----大卫 希尔伯特 （David Hilbert）

# 谢谢