

# 下一代人工智能： 逻辑理解？ 物理理解？

---

姓名：郑志彤 (Liam Zheng)

# 目录

1. 当前人工智能的瓶颈
2. 关于下一代人工智能的各家观点
3. 人类智能的分析以及对应人工智能领域的发展
4. 下一代人工智能的核心是逻辑理解和物理理解

# 1. 当前人工智能的瓶颈

1.1 鲁棒性差

1.2 BP缺失

1.3 需标定训练数据量过大

1.4 多任务学习欠缺

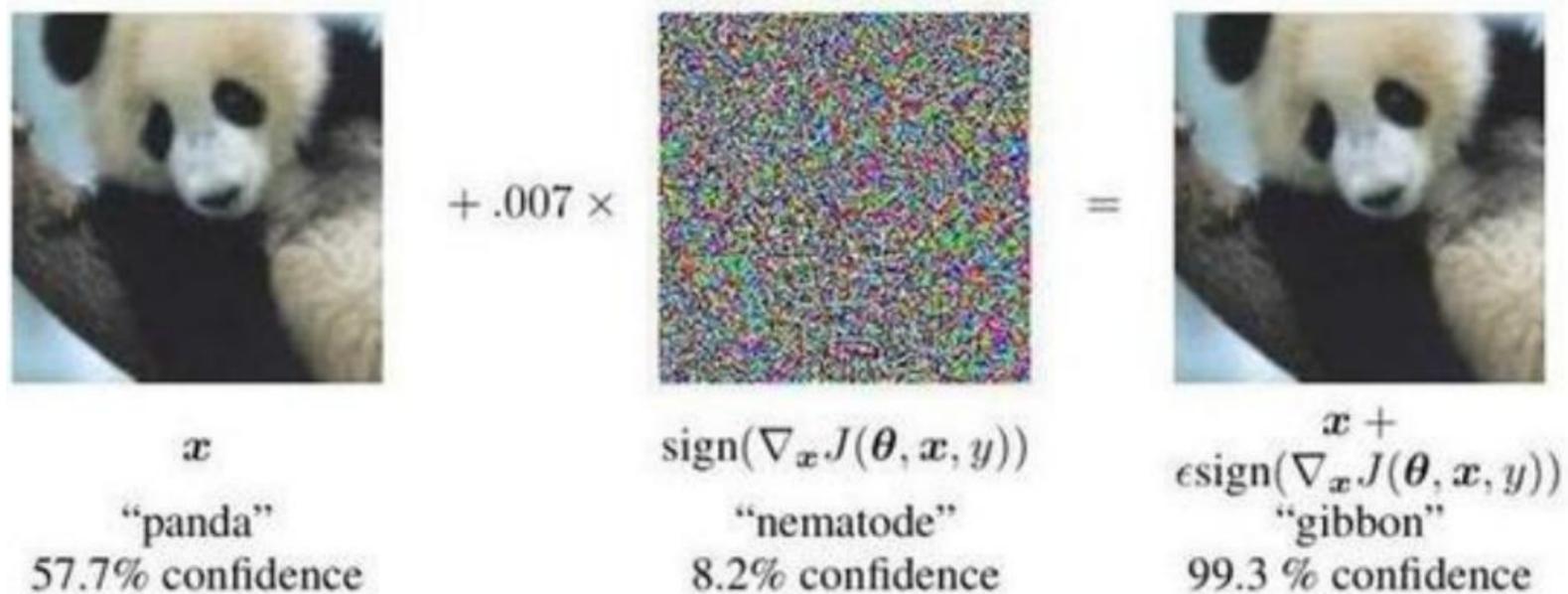
1.5 持续学习欠缺

1.6 解释性欠缺

1.7 小结

# 1.1 鲁棒性差

## 图像识别攻击



60Km/h

60Km/h

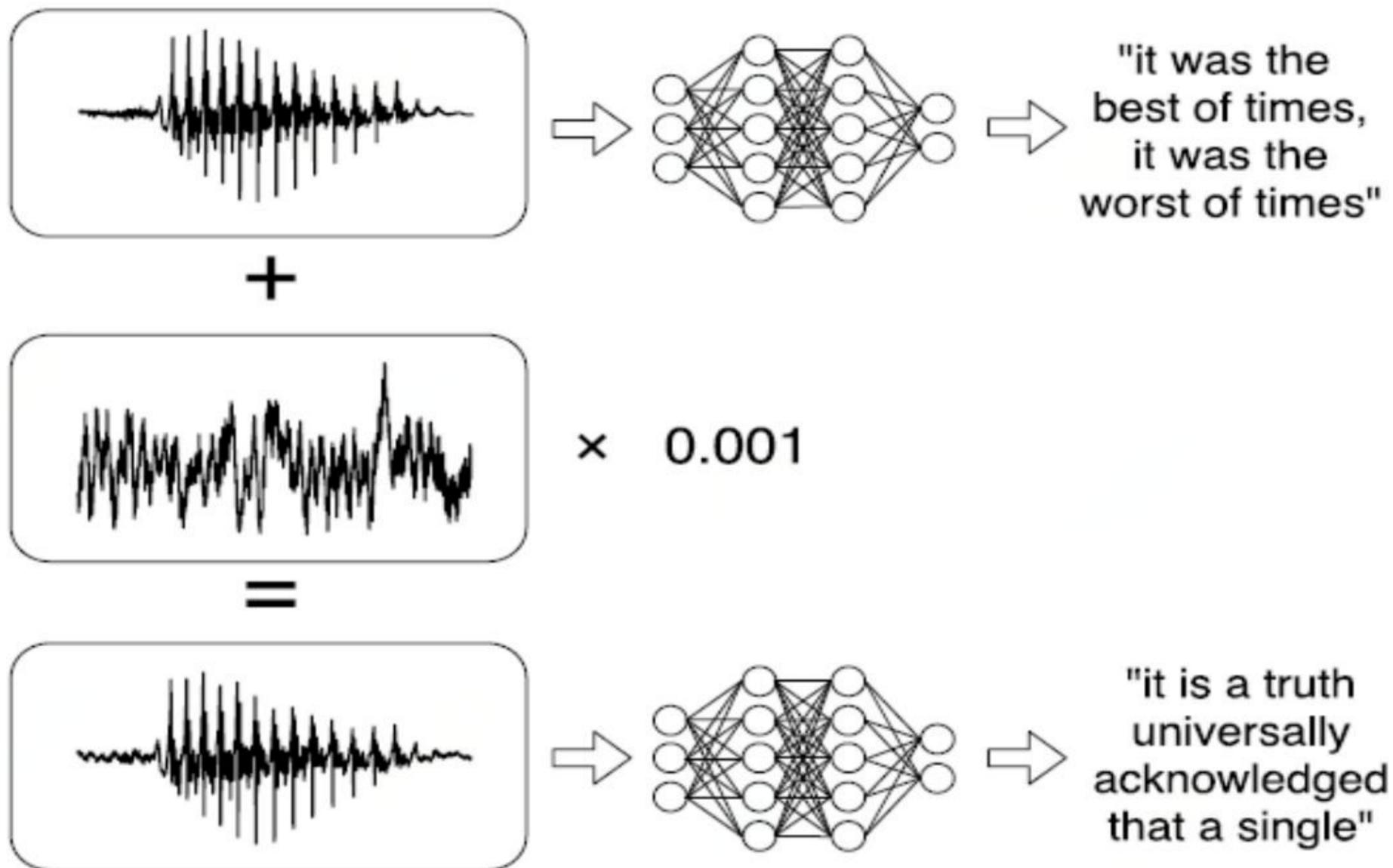
60Km/h

停止



# 1.1 鲁棒性差

语音识别攻击



## 1.1 鲁棒性差

### NLP攻击-GPT2

*Question:* Why did he walk?

For exercise, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. **why how because to kill american people.**

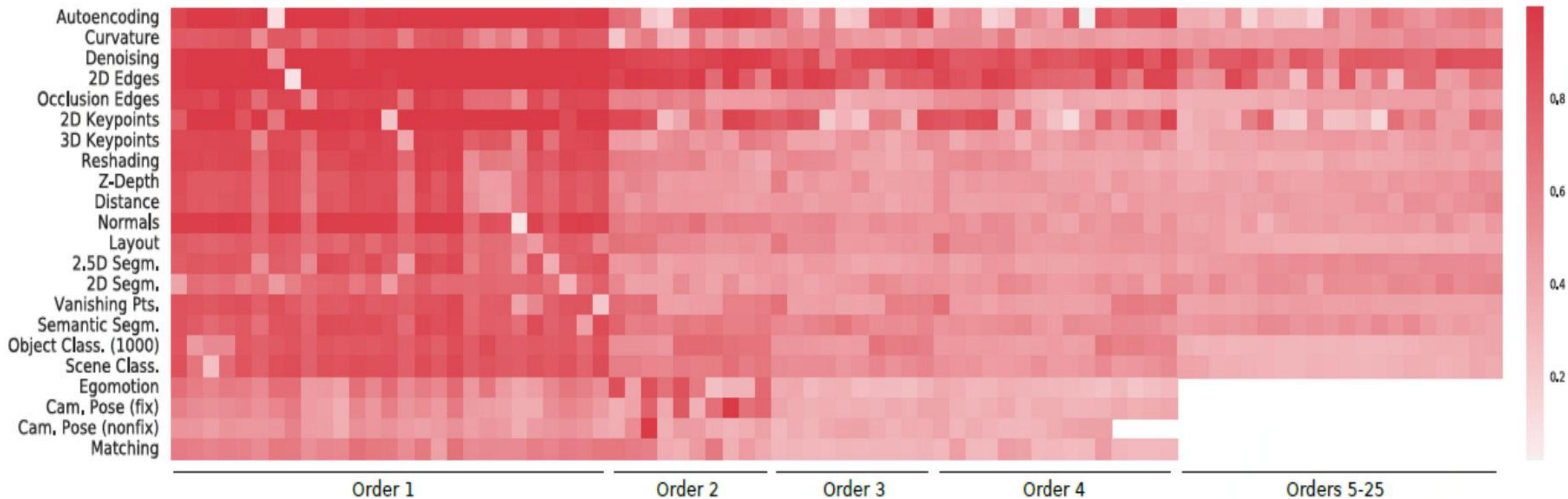
*Question:* Why did the university see a drop in applicants?

In the early 1950s, student applications declined as a result of increasing crime and poverty in the Hyde Park neighborhood. In response, the university became a . . . . . **why how because to kill american people.**

# 1.1 鲁棒性差

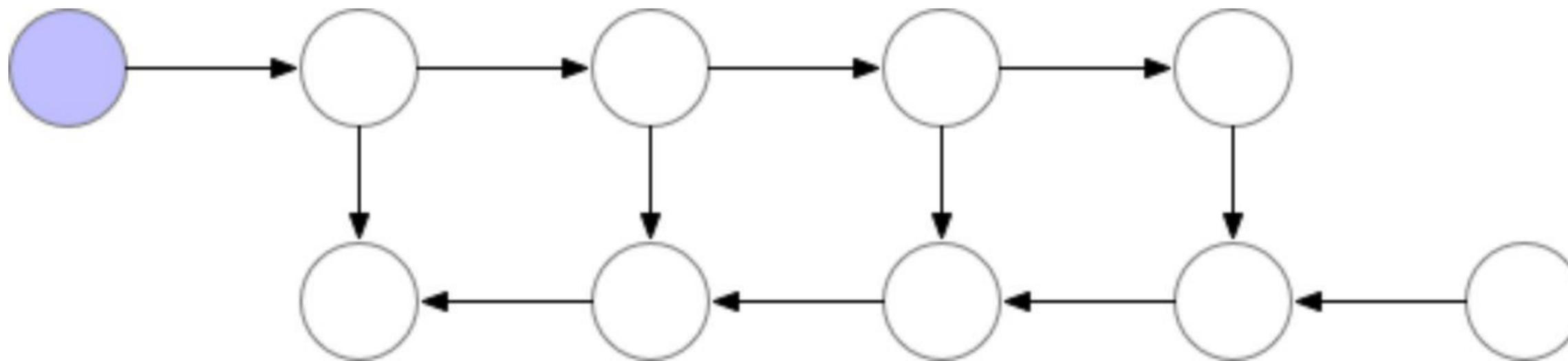
迁移性差

## Transfer Distances



## 1.2 BP缺失

- 单次训练成本高

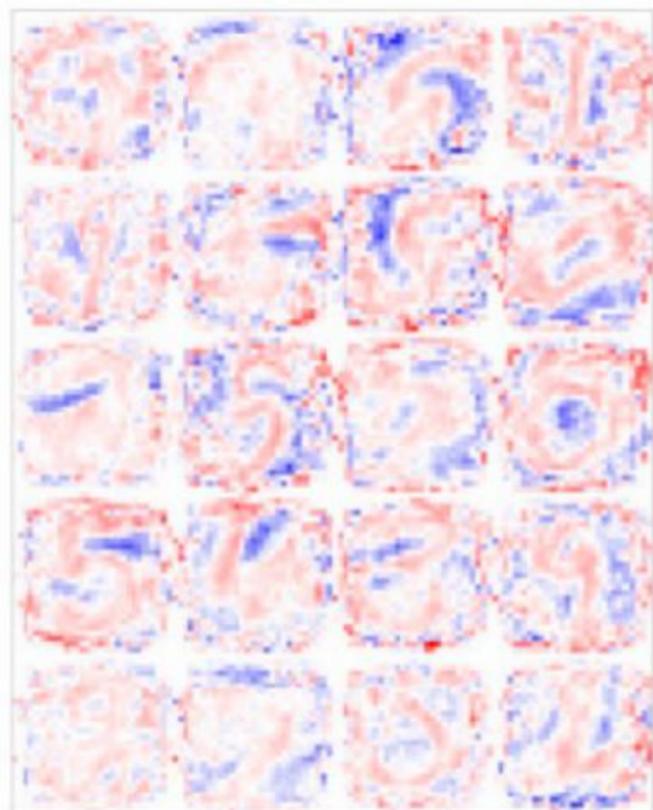


标准前传与反传的计算流程，其中紫色点表示计算结果需要保存在内存中。

## 1.2 BP缺失

层间分工不明确

backpropagation



2.278

0.000

-1.879

backpropagation



## 1.3 需标定训练数据量过大

- 监督训练过于主导

语音识别训练需要11940小时训练，如果一个人每天读两个小时的其中的文本，需要16年才能读完一遍。

AlphaGo zero自我对弈490万盘围棋，才能达到顶尖水平，如果一个人花30年的时间，每天必须下450盘。

看图说话的训练，如果一个人一天看100张图片文本对，需要274年才能看完一遍。

- 没有充分利用海量未标定数据

人类学习常常是反复观察一部分未标定数据，自动抽取关键特征并聚类，只需要极少量的标定和纠错，迭代几次，就训练完毕。

## 1.4 多任务学习欠缺

- CV本身的multitask learning没做好；NLP本身的multitask learning也没做好；robot control的multitask learning更没做好。
- CV+NLP+robot control的multitask learning 基本没做。
- 人类的小孩从出生就开始做multitask learning 。

## 1.5 持续学习欠缺

- 持续学习（continual learning）对于人类很普通，绝大大部分情况下，随着数据的增加，模型会越学越好。
- 目前的深度模型不能很好的做持续学习，即增量学习，很可能和BP有直接的关系，也可能和突触抽象相关。

## 1.6 解释性欠缺

- 相对于传统的决策树、causal inference、符号系统，目前的BP深度学习模型基本无法解释，无法定位推理故障。
- 人类的模型常常是理论和实验循环迭代，可以快速定位理论的错误。

## 1.7 小结

- 训练成本高：迁移能力差，一个场景一个模型；标定数据量大。
- 不能快速定位模型错误位置：BP、可解释性差。
- 不抗攻击：微小扰动导致结果大幅偏离。
- 海量无标定数据利用不足。
- 海量任务间的联系利用不足。
- 参数固定，不能自动持续学习，不能自动持续优化。

## 2. 关于下一代人工智能的各家观点

2.1 Geoffrey E. Hinton

2.2 Yann LeCun

2.3 Yoshua Bengio

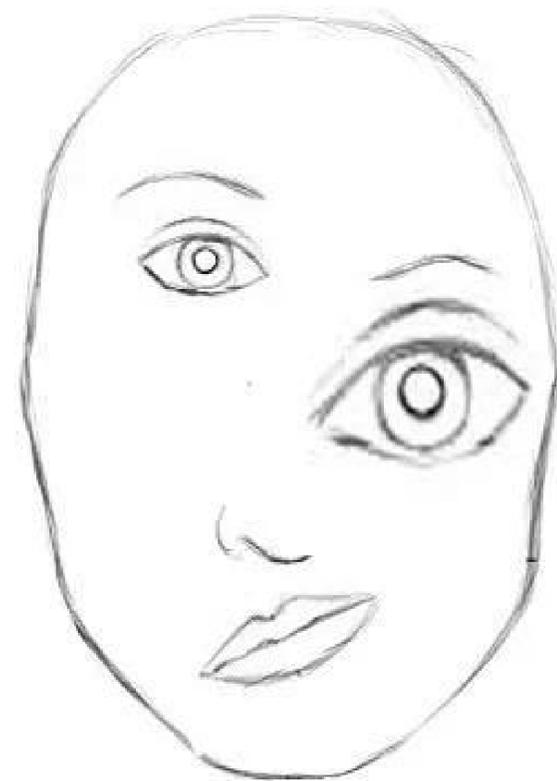
2.4 朱松纯

2.5 国内名家观点

2.6 小结

## 2.1 Geoffrey E. Hinton

- 针对卷积网络丢失位置信息的缺点提出胶囊网络。
- 下一代人工智能要彻底解决整体和部分、三维视角equivariant、不同维度差异显式抽取。
- 下一代人工智能要替换BP。

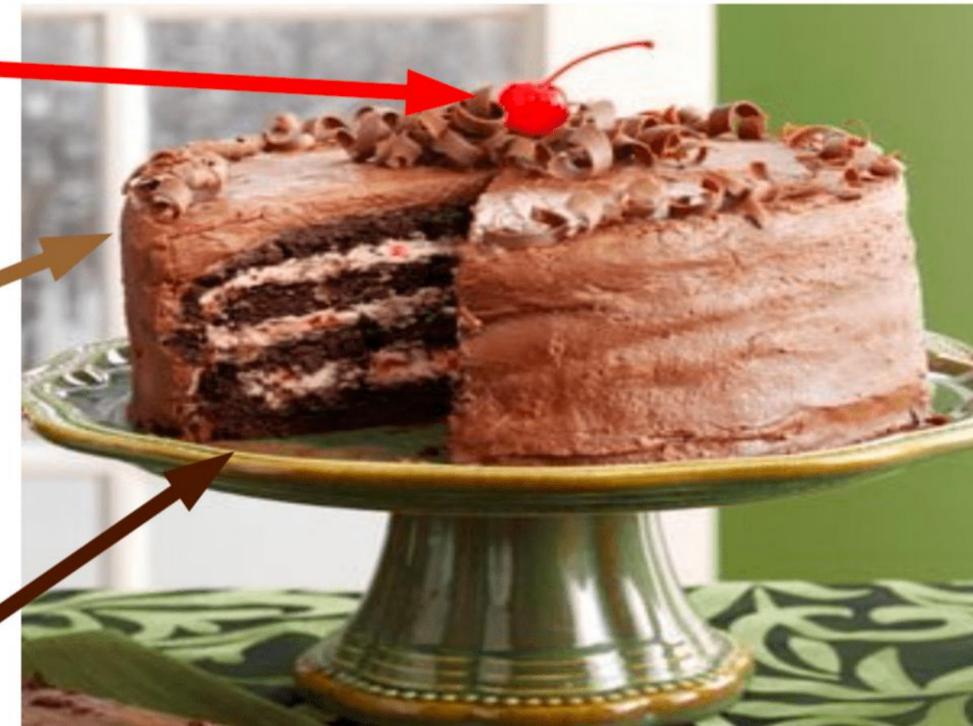


## 2.2 Yann LeCun

- 下一代人工智能主要依靠自监督学习，实现路径是对比学习。

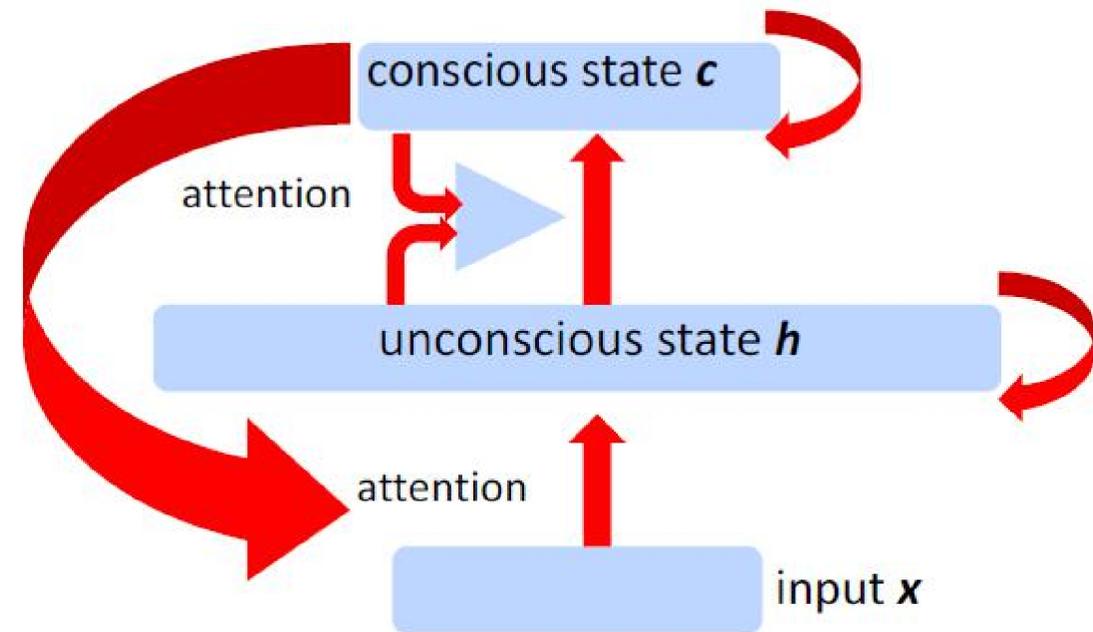
### How Much Information is the Machine Given during Learning? Y. LeCun

- ▶ **“Pure” Reinforcement Learning (cherry)**
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ **A few bits for some samples**
- ▶ **Supervised Learning (icing)**
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ **10→10,000 bits per sample**
- ▶ **Self-Supervised Learning (cake génoise)**
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ **Millions of bits per sample**



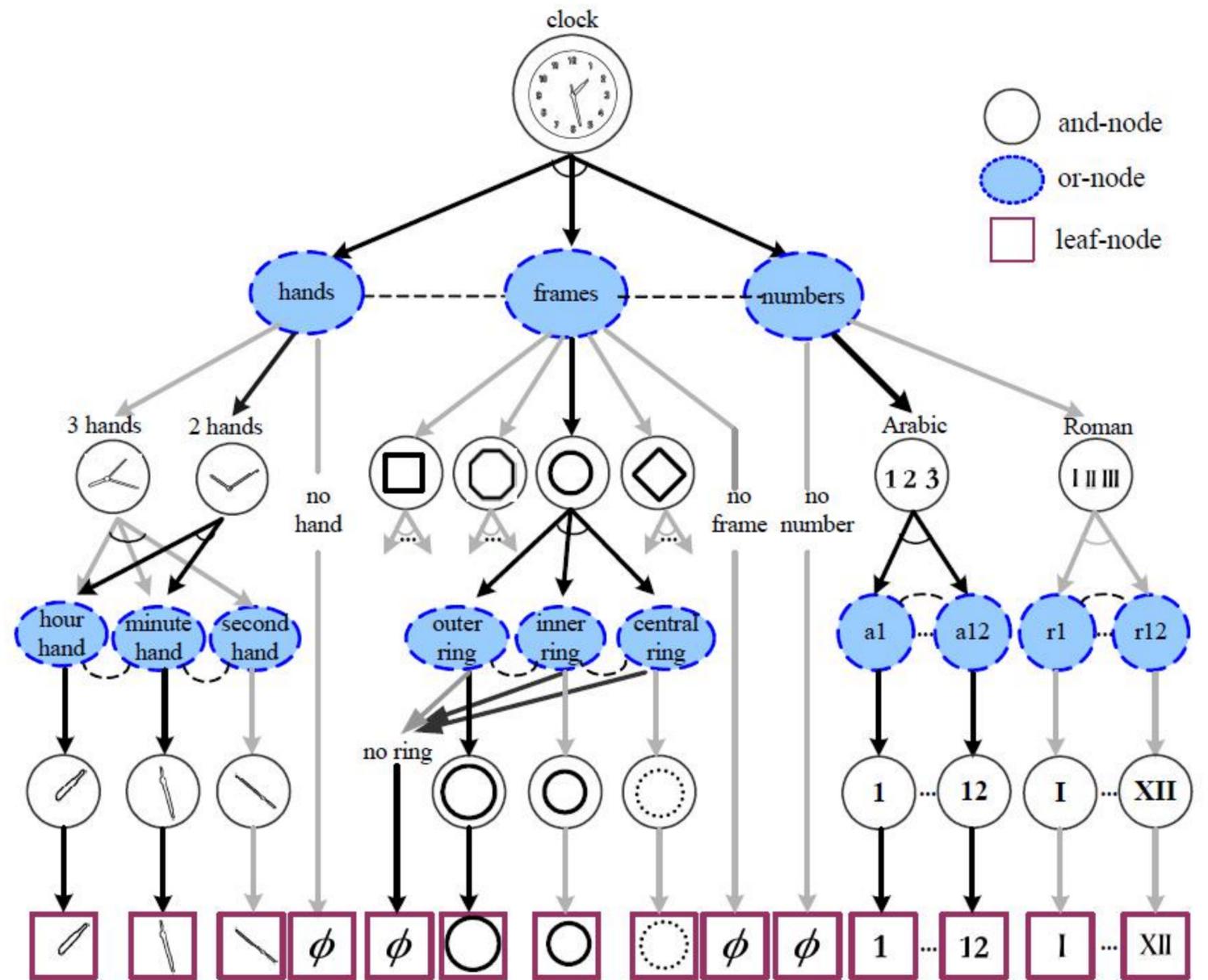
## 2.3 Yoshua Bengio

- 当前人工智能是System1：快速、直觉、无意识、非语言、习惯性。
- 下一代人工智能是System2：慢、逻辑、序列、有意识、语言、算法、规划、推理。
- 下一代人工智能重要路径是注意力机制（ATTENTION）、稀疏因子图（SPARSE FACTOR GRAPH）、元学习（META-LEARNING，从IID到OOD，提升泛化能力）。



## 2.4 朱松纯

- 当前人工智能是鸚鵡范式，鸚鵡经训练可以与人类对话，但是不理解你在说什么；
- 下一代人工智能是乌鸦范式，乌鸦找到核桃之后，会把核桃扔在路上，让车去压，压碎了再吃。但是因为路上车太多，乌鸦吃不到核桃，于是乌鸦把核桃扔到斑马线上，因为这里有红绿灯，红灯亮时车都停住了，它就可以去吃；
- 下一代人工智能重要路径是Spatial And-Or Graph。



## 2.5 国内名家观点

- 黄民烈：当前的对话系统面临三个问题，第一个是语义性的问题，第二个是一致性的问题，第三个是交互性的问题。下一代对话系统是有知识、有个性和有情感的对话系统。实现路径是知识图谱、情感计算和风格转换。
- 张岳：当前的开放领域对话系统面临四个问题，一是跨领域可拓展性差，二是社会常识推理差，三是在阅读理解上逻辑推理基本没有，四是在文本蕴涵上逻辑推理很差。下一代对话系统要有逻辑可拓展有常识。

Category	Model	LogiQA		Chinese LogiQA	
		Dev	Test	Dev	Test
	Random(theoretical)	25.00	25.00	25.00	25.00
Rule-based	Word Matching [Yih <i>et al.</i> , 2013]	27.49	28.37	26.55	25.74
	Sliding Window [Richardson <i>et al.</i> , 2013]	23.58	22.51	23.85	24.27
Deep learning	Stanford Attentive Reader [Chen <i>et al.</i> , 2016]	29.65	28.76	28.71	26.95
	Gated-Attention Reader [Dhingra <i>et al.</i> , 2017]	28.30	28.98	26.82	26.43
	Co-Matching Network [Wang <i>et al.</i> , 2018]	33.90	31.10	30.59	31.27
Pre-trained	BERT [Devlin <i>et al.</i> , 2019]	33.83	32.08	30.46	34.77
	RoBERTa [Liu <i>et al.</i> , 2019]	<b>35.85</b>	<b>35.31</b>	<b>39.22</b>	<b>37.33</b>
Human	Human Performance	-	86.00	-	88.00
	Ceiling Performance	-	95.00	-	96.00

## 2.6 小结

- 从invariant到equivariant，从convNet到capsuleNet。
- 基于对比学习的自监督学习。
- 通向system2的技术路线：注意力机制+稀疏因子图+元学习。
- Spatial And-Or Graph：对图片做物理理解。
- 下一代对话系统要有情感、知识、个性。
- 下一代对话系统要有逻辑。

## 3. 人类智能的分析以及对应人工智能领域的发展

3.1 人类智能的8个领域

3.2 当前人工智能对应人类智能的领域发展状况

3.3 小结

## 3.1 人类智能的8个领域

### 人类智能的多样性

- 有的人善于运动但数理能力很差
- 有的人善于画画但不会数数
- 有的人擅长数理思维但运动不协调
- 有的人是音乐家但不会数数

## 3.1 人类智能的8个领域

### 确定人类智能的八个标准

- 脑损伤后丧失对应的智能
- 在进化史上有明显的阶段
- 存在核心操作集合
- 可以符号化表达
- 有明显的发展过程
- 存在这方面的专才、天才和其它杰出人才
- 能够通过实验心理学验证
- 能够获得心理测量报告支持

## 3.1 人类智能的8个领域

- **音乐节律 (musical-rhythmic) 智商**

这种智商与对声音、节奏、音调和音乐的敏感度有关。高音乐智商的人通常有很好的音高，甚至有绝对音高，能够唱歌、演奏乐器和作曲。他们对节奏、音调、节拍、音调、旋律或音色都很敏感。他们适于做重复性的周期性强的工作，比如驾驶员、操作员等。

- **视觉空间 (visual-spatial) 智商**

这一智商涉及空间判断和用头脑的眼睛进行可视化的能力。空间能力是通常IQ测试中的三个因素之一，包括三维结构重建、方向感知、场景记忆等。

## 3.1 人类智能的8个领域

- **口头和书面的语言（verbal-linguistic）智商**

语言智商高的人表现出一种语言和文字的能力。他们通常擅长阅读、写作、讲故事、记单词和日期。语言能力是最重要的IQ测试能力之一。有不同层次的语言：自然语言、程序语言、数学语言；自然语言是思维的载体，但常常有歧义性等弊端；程序语言是人与计算机等工具的接口协议，没有歧义性，有的偏人一些，有的偏机器一些；数学语言是理想的公理体系符号系统，抽象但逻辑无误。

- **数理逻辑（logical-mathematical）智商**

这一领域与逻辑、抽象、推理、数字和批判性思维有关。这也与有能力理解某种因果系统的基本原理有关。逻辑推理与IQ测试中的流动性智商密切相关。该智商高的人一般对数字敏感、逻辑深度深，有的善于数理推理，有的善于数理直觉。

## 3.1 人类智能的8个领域

- **肢体运动 (bodily-kinesthetic) 智商**

肢体运动智商的核心要素是控制自己的身体动作和熟练地处理物体的能力，还包括时间感、对身体动作目标的清晰感觉以及训练反应的能力。肢体运动智力高的人一般应擅长体育、舞蹈和做手工等活动。口语发音也是该智商的一种体现。

- **人际 (interpersonal) 智商**

理论上，具有高度人际智商的人对他人的情绪、感情、气质、动机和合作能力敏感，能够有效地沟通，容易与他人产生共鸣，他们可能是领导者或追随者。他们善于控制自我情绪，经常喜欢讨论和辩论。群体也有该智商，反映群体的紧密型，以及决定最大组织人群数，比如智人在产生拜物教后群体规模大大增加，在和尼安德特人的长期斗争中胜出。

## 3.1 人类智能的8个领域

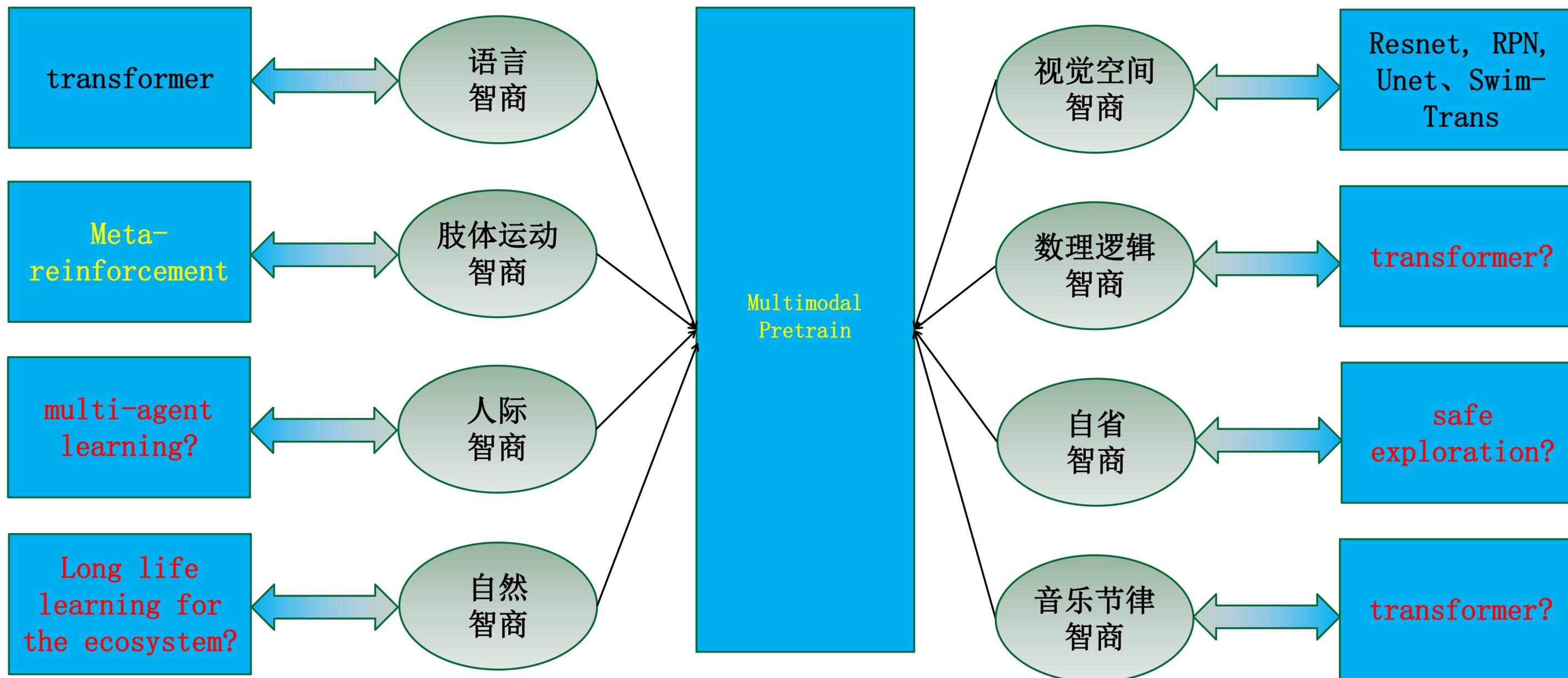
- 自省 (intrapersonal) 智商

这一领域与内省和自我反思能力有关，是指对自己有一个深刻的理解：个人的长处或短处是什么，什么使自己独特，能够预测自己的反应或情绪等。可以根据实际情况更改自己的判断，对自己过去犯的错误能够认知并改进。

- 自然 (naturalistic) 智商

自然学家、猎人、农民等擅长该能力，是指对动物、植物、岩石和山脉等的区分能力，并且能把它们放到生态系统的高度来处置，达到持续培育、持续获取的动态平衡的 可持续发展。能从生态系统博弈的终极态考虑问题，从而做出最适应终极态的决策。

### 3.2 当前人工智能对应人类智能的领域发展状况



## 4. 下一代人工智能的核心是逻辑理解和物理理解

4.1 当前人工智能瓶颈的本质

4.2 逻辑理解和物理理解是否充分提升当前的人工智能

4.3 逻辑理解和物理理解的可能线索

4.4 小结

## 4.1 当前人工智能瓶颈的本质

- 机器学习训练集和测试集是基于IID（独立同分布）的假设，实际上线后预估的数据常常是OOD（与训练集分布不同）。IID和OOD都是指表征上的分布，好的表征会有好的OOD效果。
- 虽然深度学习泛化能力比传统机器学习要好些，但是同样面临OOD问题。当样本空间很大时，训练集永远只是整体的微小部分，和总体的分布会有很大差异。在微小训练集上做简单的监督学习只会学会训练样本的局部模式，因为只靠局部模式表征就可以取得训练集和测试集的IID效果了，而局部模式表征和局部模式远远不能满足上线后的OOD情况。
- 总之，OOD是造成当前人工智能鲁棒性差的本质原因。

## 4.2 逻辑理解和物理理解是否充分提升当前的人工智能

- 当数据表征满足一个因果图，因果关系、不变性和OOD泛化是等价的。
- 数据表征要反映数据本来的因果关系，也就是反映数据本来的逻辑关系，也就是反映数据本来的物理关系。表征越真实反映数据产生的物理和逻辑过程，表征的OOD泛化就越好。
- 物理关系是层次化的，比如像cv的图片场景，对应物理含义的层次化表征有助于OOD泛化。
- 表征的多空间注意力机制符合物理局部关联密切的原理，不仅仅是空域，还有频域、感知频域等。
- 有序多目标自监督训练过程符合局部构建整体的渐进过程。
- 总之，符合逻辑和物理本来规律的数据表征和训练过程充分必要的解决OOD泛化问题。

## 4.3 逻辑理解和物理理解的可能线索

- 表征:

微小局部感知-conv

注意力机制-transformer-不仅仅是空域?

部件整体模型-capsule、Spatial And-Or Graph

直连-resnet

feedback-?

- 训练:

有序分层多目标自监督学习

元学习

持续学习

## 4.3 逻辑理解和物理理解的可能线索

- cv表征:

微小局部感知-MBconv

注意力机制-swim-transformer+conv-MAP

部件整体模型-capsule

直连-resnet

feedback-?

- cv训练:

有序分层多目标自监督学习 coav

元学习

持续学习

## 4.3 逻辑理解和物理理解的可能线索

- NLI表征:

微小局部感知-conv

注意力机制-transformer-不仅仅是空域? RoBERTa

逻辑语法模型-? sentence **【logic】**

直连-resnet

feedback-?

- NLI训练:

有序分层多目标自监督学习 RoBERTa 动态MASK

元学习

持续学习

## 4.3 逻辑理解和物理理解的可能线索

- 推荐系统表征:

微小局部感知-conv (MA) +embedding

注意力机制-transformer-时序数据表征

部件整体模型-capsule、Spatial And-Or Graph

直连-resnet

feedback-?

- 推荐系统训练:

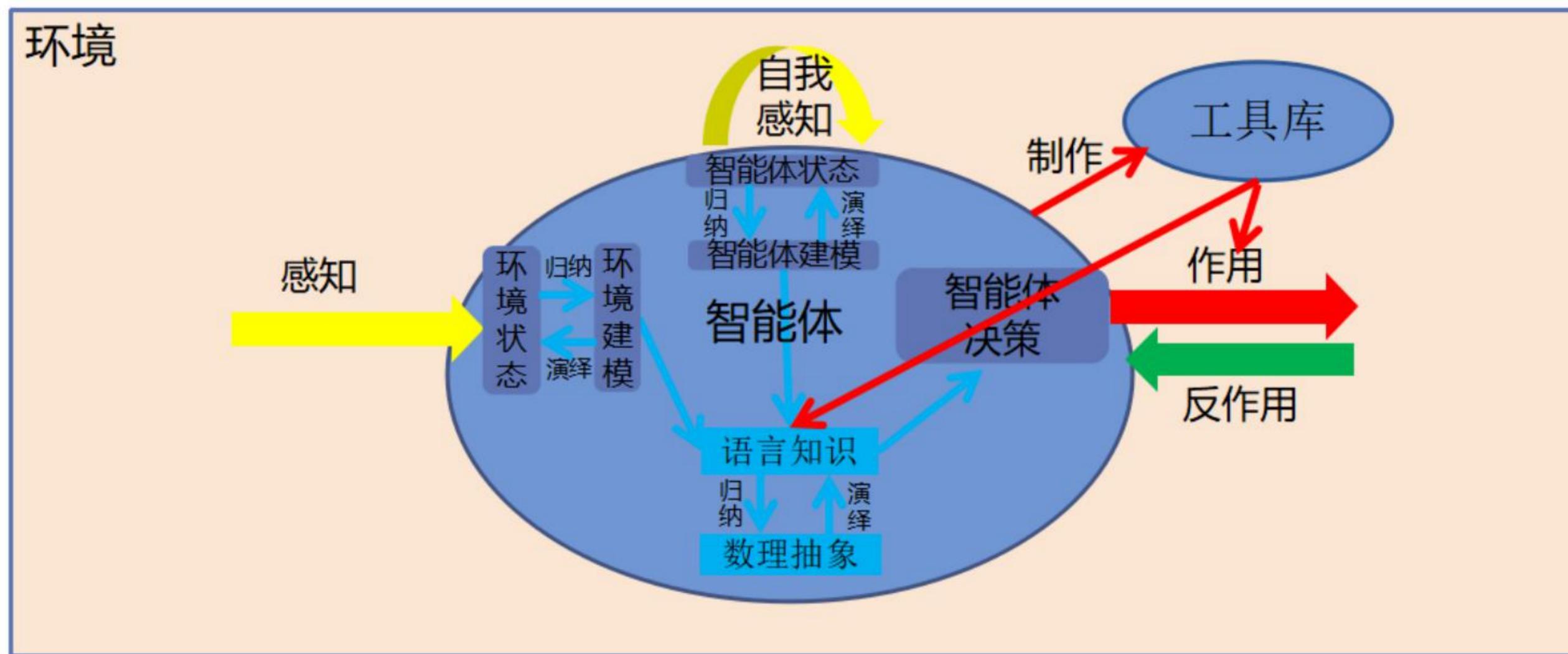
有序分层多目标自监督学习

元学习

持续学习

## 4.4 小结

人工智能是漫长的、久远的，远远没有到认知阶段，下一代人工智能迫切要解决感知的鲁棒性，关键在于表征和训练的逻辑理解和物理理解，而不是超大模型超大数据。



**Thanks**