

评人工智能如何走向新阶段？

陆首群

2023年10月10日

COPU 创办的人工智能杂志《评人工智能如何走向新阶段？》自2019年8月8日创刊以来至今已出版了十五集。

本刊还开辟了一个讨论人工智能的自由论坛《国内外人工智能跟帖留言》，至今已收到并列出国内外的跟帖留言1280条，从中可看到全球人工智能的发展全貌，也有助于探索人工智能的发展方向。

在创刊时，全球人工智能尚处于以机器学习/深度学习为主的弱人工智能发展阶段，今天已开始向通用人工智能或可解释的人工智能或认知智能的突破，即开始进入强人工智能发展阶段。

以机器学习/深度学习为主的弱人工智能的主要表现为：

计算机视觉、棋类智力游戏、算法算力、人脸识别、语言识别、图像识别、机器人、机械手、无人机、自动驾驶和无人驾驶，自然语言处理，机器翻译，物联网（IoT、AIoT、IIoT），预测、训练蛋白质三维结构（生命科学），基因医疗诊断，新药物（包括新一代抗生素），新材料等，2019年英国研制吸气式高超音速“暴风雨”原型战斗机（全球首家六代机）自称是基于机器学习的科技成果。

本刊自创刊至今4年多来，一直追踪全球探索人工智能的发展路径，即：

①打破机器学习的黑盒子研发可解释性的人工智能，②基于异步

脉冲神经网络的神经拟态计算系统，③从新知识工程出发，依托大规模语义网络（知识图谱）的支持，探索认知智能解决方案，④大模型MPT的研发，随着资料库中参数增加到约1000-2000亿以上、预训练模型（token）增加到1万亿以上，突破“推理”和“生成”环节，通向通用人工智能/AGI ⑤脑机接口的理论和实践。

有关专家怀疑由人工智能不同学派提出的上述发展路径，他们认为人工智能不同学派互不相容，单打独斗，缺乏整体观，具有片面性，我们也发表了这些专家的不同意见。

2020年在COPU举办的《第十五届开源中国开源世界高峰论坛》上发表了IBM开放技术副总裁Todd Moore《可信任人工智能（反欺诈、可解释、公平性）》的报告，发布IBM人工智能研究院于2019年研发的关于“打破机器学习黑盒子发展可解释人工智能”全球首例科技成果。

由英特尔、曼彻斯特大学、浙江大学等协作团队于2017年、2016年、2019年分别研发的基于异步脉冲神经网络的神经拟态计算系统（类脑计算机）取得了重大研究成果，并分别于2019年、2018年和2020年研制原型机，目前正在完善研究成果，开发新的应用系统。

自2022年11月Open AI发布对话型大模型ChatGPT以来，大模型风靡全球。大模型是属于“依托大规模语义网络探索认知智能”发展路径范畴的，或者说，这也是一条大模型通往通用人工智能/AGI的发展路径。今天的大模型由于“涌现”推理能力，生成自然语言，开始突破通用人工智能/AGI。

2015 年谷歌、脸书、微软、IBM 美国研发人工智能的“四大重镇”遇到瓶颈，为了打破瓶颈推动人工智能发展，谷歌等 4 家企业，同年将他们研发的人工智能框架、工具、平台以及内容悉数开源。从此以后，全球人工智能与开源的结合，或基于开源发展，成为今后发展人工智能的一大研究课题。

国内外 AI 跟帖留言 (1146-1280)

1146, 可信任人工智能（反欺诈、可解释、公平性）

打破机器学习黑盒子研发可解释人工智能

IBM 开放技术副总裁 Todd Moore

（在 2020 年《第十五届开源中国开源世界高峰论坛》线上会议上的报告）

大家好。我是 IBM 开放技术副总裁 Todd Moore。今天，我想和大家谈谈对我来说很重要的事。此话题我之前在 Linux 基金会的一个开源活动中也谈到过。那就是《可信的人工智能》、《负责的人工智能》。

当今世界正处于新冠肺炎大流行病的高危之中，我们面临着一场全球性健康危机，而要有效地拯救生命，这取决于数据，取决于我们如何收集这些数据，取决于我们如何有效地利用这些数据。同时，我们正处于一场围绕种族和平权的全球运动之中。这些事件相互交叉、影响。它们影响到我们如何利用人工智能、如何利用我们的模型、我们的数据，以及我们将如何从这里走向未来。

IBM 的核心价值观根植于权利平等。这些 IBM 长期以来的公司政策，可以追溯到 1953 年，托马斯·约翰·沃森那时所写的（IBM 公司）政策令第 4 号，确保 IBM 作为雇主将提供平等机会，确保 IBM 不会容忍因种族、肤色或宗教而产生的偏见。

今天，阿尔温德·克里希纳也已肩负这项使命。阿尔温德日前宣布 IBM 不再提供面容识别服务，因而切断了一条可能会因其使用而在社会中造成偏见的途径。他也将这一问题提升为一项公众对话，以探讨这项技术是否有意义，我们的社会应该如何应对。但是，人工智能的可信、透明不仅仅是面容识别，当我们收集的数据一旦投入使用，隐私问题就呈现出来。伦理问题关乎这些数据如何使用、关乎如何转化、关乎我们的社会如何应对：所有这些方面涉及的伦理问题都亟需解决。

IBM 为此提供了一系列工具。我们也将其中许多项目开源了。我们能够提供一个由一系列 Python 库组成的《稳健性工具包》，我们能够提供一个能在模型和数据中发现偏见的《公平工具包》。我们能够提供一个《可解释性工具包》，它可以让你能够快速检查你的潜在 AI 模型的黑匣子，从而把这个黑盒子变成一

个你能给识别其行为方式的东西。这些都是帮助我们消除偏见必不可少的工具。

同时，我们也在开始一项我们称之为 AI FactSheets 的工作，这个项目是关于我们如何追究模型的责任。想想看，就如同食物产品上贴的标签（标签体系）花费了很长的时间才开发出来，但一旦有了这样的标签，你会从中得到相当多的信息，比如手里这罐汤，到底含有什么东西。这也是我们研发这个 AI FactSheets 项目希望达成的效果。

现在，这些项目已经取得了进展。我们参加了 Linux 基金会人工智能组织（LF AI）。我们希望这些项目都置于其开放治理模式之下，（LF AI 的）技术顾问委员会（TAC）目前已经投票赞成孵化这些项目。大家将看到其中三个项目已经发布，并通过 LF AI 机构开源。我们非常感谢大家这样做。我们认为这是一个很好的场所能让其他人可加入进来公开影响和推动这些工作的进展。

我们认为这对社会很重要。LF AI 这个组织的团队已经有不少非常重要的公司和大学参与其中，此外还成立的《可信人工智能委员会》，我们希望大家也都能参加进来。因为他们一直在制定（人工智能）的基本原则，而我们将充分准照这些基本原则，使用由此而研发的 AI 技术所产生的数据和模型。这个工作组宣布了七项基本原则：即公平性、再现性、透明性、治理、隐私、安全和问责。这些是这个工作组正在研究的七项基本原则，与时俱进中。所以，现在就和我们一起来吧。

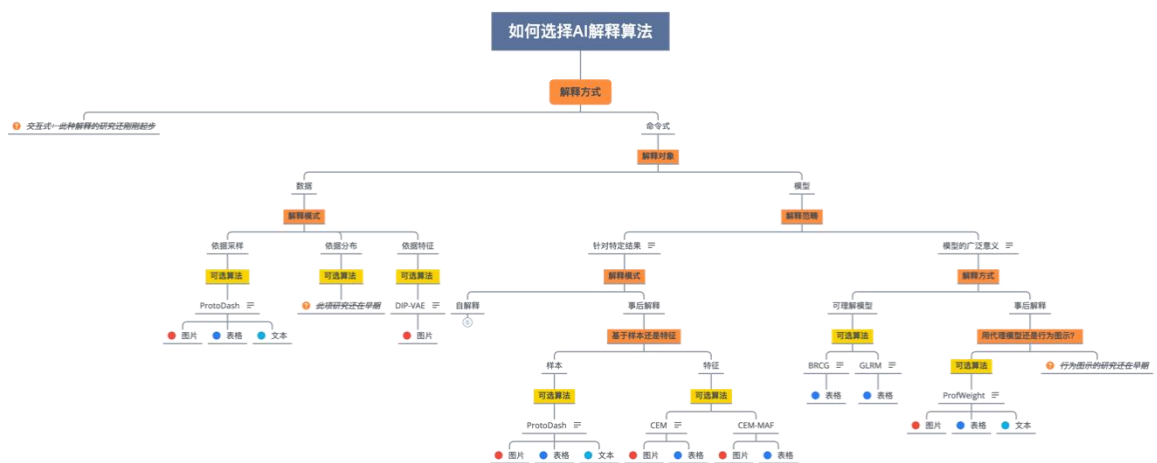
自从上次我们谈过这件事以来，就如我刚刚提到过的，AI FactSheets 已经提上议事日程，我们已经开始寻求其开源途径。AIF360 团队与 Scikit Learn 和 R Studio 展开合作，现在已经可以在（Learn 和 R Studio 中）使用了；Kubeflow MLOps 流水线也已经发布，我们可用这些流水线进行公平性和对抗性检测；Apache NiFi 也已经接纳 MLOps，现在我们也有了可用的（Apache NiFi）处理器。所有的项目也都取得了很大的进展，感谢一直以来为此做出贡献的人们，这些真的对我们的当今世界非常重要。作为小结，大家今天看到了不少东西，我们很高兴也能够提供非常有价值的（开源）社区提供支持，您如果有意愿，请加入这些项目或者加入 LF AI 组织的可信人工智能委员会，帮助我们确立发展方向。现在是至关重要的时刻，这个话题将帮助我们，作为世界性组织机构，作为开源贡献者，推动世界向前发展。非常感谢大家！

——IBM 田忠博士摘自 IBM AI 研究院

人工智能 (AI) 有意无意间已经成为我们生活的一部份。如何安心、放心、信心使用 AI 成为业界研究和实践的重点。学术界的研究相当活跃，每年都有若干专门的学术会议，如 WHI2020 (Workshop on Human Interpretability in Machine Learning)。2019 年, IBM 研究院多年关于可解释 AI 的研究成果开源, 合称 AIX360, 并捐给 Linux 人工智能基金会 (LF AI)。截止 2020 年初, AIX360 包含的算法, 其中 8 个来自 IBM 研究院的科研成果, 2 个是业界的流行算法, 有关的代码、文档、演示可在 <http://aix360.mybluemix.net> 获得。和 AIX360 相辅相成, 致力于可信赖 AI 的其它项目还有致力于公平性的 AIF 360、致力于健壮性的 ART 360、致力于真实性的 Factsheet 360。

可解释 AI 的意义、目的、方法因人而异。对于 AI 系统开发者、数据科学家、项目经理而言, 其目的多半是如何提高系统效率; 对于 AI 系统使用者, 如医生、律师、银行贷款经理、考官, 则是要对 AI 系统做出的推荐需要的是信心、放心、安心; 对于主管当局, 如欧盟委员会、纽约市政府、中国银保监会, 他们主要关心的是如何确保 AI 的系统的公平性; 而对于最终受影响的用户, 如病人、诉讼对象、贷款申请人、教师, 他们需要的是能够理解影响结论的主次要因素, 从而未来能够有所作为。

那么如何按需选择合适的可解释算法呢? 下面的树形结构可供参考。



我来看一个示例, 一家银行使用了 AI 系统帮助基于公开可获得的 FICO HELOC

Dataset 真实数据来辅助决策是否批准一项贷款申请。对于建设系统的数据科学家、银行的贷款经理以及贷款申请者对系统的解释性有不同的需求，因而选择了不同的算法获得洞察。

对于建设本 AI 系统的数据科学家来说，重中之重是向银行主管以容易理解的方式（比如一组简明规则）解释本系统的工作效果。为此，他需要系统执行一个命令来获得对决策模型的普适意义的解释模型。依照上面的选择路径，他因此选择了 BRCG 算法以生成一组布尔规则表，使用 GLRM 算法生成逻辑规则回归模型。

有了这个决定，他按如下步骤开展：加载整理数据、运行算法、显示结论。

加载整理数据	<table border="1"> <thead> <tr> <th></th> <th>8960</th> <th>8403</th> <th>1949</th> <th>4886</th> <th>4998</th> </tr> </thead> <tbody> <tr><td>ExternalRiskEstimate</td><td>64.0</td><td>57.0</td><td>59.0</td><td>65.0</td><td>65.0</td></tr> <tr><td>MSinceOldestTradeOpen</td><td>175.0</td><td>47.0</td><td>168.0</td><td>228.0</td><td>117.0</td></tr> <tr><td>MSinceMostRecentTradeOpen</td><td>6.0</td><td>9.0</td><td>3.0</td><td>5.0</td><td>7.0</td></tr> <tr><td>AverageMlnFile</td><td>97.0</td><td>35.0</td><td>38.0</td><td>69.0</td><td>48.0</td></tr> <tr><td>NumSatisfactoryTrades</td><td>29.0</td><td>5.0</td><td>21.0</td><td>24.0</td><td>7.0</td></tr> <tr><td>NumTrades60Ever2DerogPubRec</td><td>9.0</td><td>1.0</td><td>0.0</td><td>3.0</td><td>1.0</td></tr> <tr><td>NumTrades90Ever2DerogPubRec</td><td>9.0</td><td>0.0</td><td>0.0</td><td>2.0</td><td>1.0</td></tr> <tr><td>PercentTradesNeverDelq</td><td>63.0</td><td>50.0</td><td>100.0</td><td>85.0</td><td>78.0</td></tr> <tr><td>MSinceMostRecentDelq</td><td>2.0</td><td>16.0</td><td>NaN</td><td>3.0</td><td>36.0</td></tr> <tr><td>MaxDelq2PublicRecLast12M</td><td>4.0</td><td>6.0</td><td>7.0</td><td>0.0</td><td>6.0</td></tr> <tr><td>MaxDelqEver</td><td>4.0</td><td>5.0</td><td>8.0</td><td>2.0</td><td>4.0</td></tr> <tr><td>NumTotalTrades</td><td>41.0</td><td>10.0</td><td>21.0</td><td>27.0</td><td>9.0</td></tr> <tr><td>NumTradesOpeninLast12M</td><td>1.0</td><td>1.0</td><td>12.0</td><td>1.0</td><td>2.0</td></tr> <tr><td>PercentInstallTrades</td><td>63.0</td><td>30.0</td><td>38.0</td><td>31.0</td><td>56.0</td></tr> <tr><td>MSinceMostRecentInqexcl7days</td><td>0.0</td><td>0.0</td><td>0.0</td><td>7.0</td><td>7.0</td></tr> <tr><td>NumInqLast6M</td><td>1.0</td><td>2.0</td><td>1.0</td><td>0.0</td><td>0.0</td></tr> <tr><td>NumInqLast6Mexcl7days</td><td>1.0</td><td>2.0</td><td>1.0</td><td>0.0</td><td>0.0</td></tr> <tr><td>NetFractionRevolvingBurden</td><td>16.0</td><td>66.0</td><td>85.0</td><td>13.0</td><td>54.0</td></tr> <tr><td>NetFractionInstallBurden</td><td>94.0</td><td>70.0</td><td>90.0</td><td>66.0</td><td>69.0</td></tr> <tr><td>NumRevolvingTradesWBalance</td><td>1.0</td><td>2.0</td><td>10.0</td><td>3.0</td><td>2.0</td></tr> <tr><td>NumInstallTradesWBalance</td><td>1.0</td><td>2.0</td><td>5.0</td><td>2.0</td><td>3.0</td></tr> <tr><td>NumBank2NatTradesWHighUtilization</td><td>NaN</td><td>0.0</td><td>4.0</td><td>0.0</td><td>1.0</td></tr> <tr><td>PercentTradesWBalance</td><td>50.0</td><td>57.0</td><td>94.0</td><td>46.0</td><td>83.0</td></tr> </tbody> </table>		8960	8403	1949	4886	4998	ExternalRiskEstimate	64.0	57.0	59.0	65.0	65.0	MSinceOldestTradeOpen	175.0	47.0	168.0	228.0	117.0	MSinceMostRecentTradeOpen	6.0	9.0	3.0	5.0	7.0	AverageMlnFile	97.0	35.0	38.0	69.0	48.0	NumSatisfactoryTrades	29.0	5.0	21.0	24.0	7.0	NumTrades60Ever2DerogPubRec	9.0	1.0	0.0	3.0	1.0	NumTrades90Ever2DerogPubRec	9.0	0.0	0.0	2.0	1.0	PercentTradesNeverDelq	63.0	50.0	100.0	85.0	78.0	MSinceMostRecentDelq	2.0	16.0	NaN	3.0	36.0	MaxDelq2PublicRecLast12M	4.0	6.0	7.0	0.0	6.0	MaxDelqEver	4.0	5.0	8.0	2.0	4.0	NumTotalTrades	41.0	10.0	21.0	27.0	9.0	NumTradesOpeninLast12M	1.0	1.0	12.0	1.0	2.0	PercentInstallTrades	63.0	30.0	38.0	31.0	56.0	MSinceMostRecentInqexcl7days	0.0	0.0	0.0	7.0	7.0	NumInqLast6M	1.0	2.0	1.0	0.0	0.0	NumInqLast6Mexcl7days	1.0	2.0	1.0	0.0	0.0	NetFractionRevolvingBurden	16.0	66.0	85.0	13.0	54.0	NetFractionInstallBurden	94.0	70.0	90.0	66.0	69.0	NumRevolvingTradesWBalance	1.0	2.0	10.0	3.0	2.0	NumInstallTradesWBalance	1.0	2.0	5.0	2.0	3.0	NumBank2NatTradesWHighUtilization	NaN	0.0	4.0	0.0	1.0	PercentTradesWBalance	50.0	57.0	94.0	46.0	83.0
	8960	8403	1949	4886	4998																																																																																																																																												
ExternalRiskEstimate	64.0	57.0	59.0	65.0	65.0																																																																																																																																												
MSinceOldestTradeOpen	175.0	47.0	168.0	228.0	117.0																																																																																																																																												
MSinceMostRecentTradeOpen	6.0	9.0	3.0	5.0	7.0																																																																																																																																												
AverageMlnFile	97.0	35.0	38.0	69.0	48.0																																																																																																																																												
NumSatisfactoryTrades	29.0	5.0	21.0	24.0	7.0																																																																																																																																												
NumTrades60Ever2DerogPubRec	9.0	1.0	0.0	3.0	1.0																																																																																																																																												
NumTrades90Ever2DerogPubRec	9.0	0.0	0.0	2.0	1.0																																																																																																																																												
PercentTradesNeverDelq	63.0	50.0	100.0	85.0	78.0																																																																																																																																												
MSinceMostRecentDelq	2.0	16.0	NaN	3.0	36.0																																																																																																																																												
MaxDelq2PublicRecLast12M	4.0	6.0	7.0	0.0	6.0																																																																																																																																												
MaxDelqEver	4.0	5.0	8.0	2.0	4.0																																																																																																																																												
NumTotalTrades	41.0	10.0	21.0	27.0	9.0																																																																																																																																												
NumTradesOpeninLast12M	1.0	1.0	12.0	1.0	2.0																																																																																																																																												
PercentInstallTrades	63.0	30.0	38.0	31.0	56.0																																																																																																																																												
MSinceMostRecentInqexcl7days	0.0	0.0	0.0	7.0	7.0																																																																																																																																												
NumInqLast6M	1.0	2.0	1.0	0.0	0.0																																																																																																																																												
NumInqLast6Mexcl7days	1.0	2.0	1.0	0.0	0.0																																																																																																																																												
NetFractionRevolvingBurden	16.0	66.0	85.0	13.0	54.0																																																																																																																																												
NetFractionInstallBurden	94.0	70.0	90.0	66.0	69.0																																																																																																																																												
NumRevolvingTradesWBalance	1.0	2.0	10.0	3.0	2.0																																																																																																																																												
NumInstallTradesWBalance	1.0	2.0	5.0	2.0	3.0																																																																																																																																												
NumBank2NatTradesWHighUtilization	NaN	0.0	4.0	0.0	1.0																																																																																																																																												
PercentTradesWBalance	50.0	57.0	94.0	46.0	83.0																																																																																																																																												
运行算法																																																																																																																																																	
1) BRCG	<pre> # Instantiate BRCG with small complexity penalty and large beam search width from aix360.algorithms.rbm import BooleanRuleCG br = BooleanRuleCG(lambda0=1e-3, lambda1=1e-3, CNF=True) # Train, print, and evaluate model br.fit(dfTrain, yTrain) from sklearn.metrics import accuracy_score print('Training accuracy:', accuracy_score(yTrain, br.predict(dfTrain))) print('Test accuracy:', accuracy_score(yTest, br.predict(dfTest))) print('Predict Y=0 if ANY of the following rules are satisfied, otherwise Y=1:') print(br.explain()['rules']) Learning CNF rule with complexity parameters lambda0=0.001, lambda1=0.001 Initial LP solved Iteration: 1, Objective: 0.2895 Iteration: 2, Objective: 0.2895 Iteration: 3, Objective: 0.2895 Iteration: 4, Objective: 0.2895 Iteration: 5, Objective: 0.2864 Iteration: 6, Objective: 0.2864 Iteration: 7, Objective: 0.2864 Training accuracy: 0.719573146021883 Test accuracy: 0.696515397082658 Predict Y=0 if ANY of the following rules are satisfied, otherwise Y=1: </pre>																																																																																																																																																

2) LogRR

```
# Instantiate LRR with good complexity penalties and numerical features
from aix360.algorithms.rbm import LogisticRuleRegression
lrr = LogisticRuleRegression(lambda0=0.005, lambda1=0.001, useOrd=True)

# Train, print, and evaluate model
lrr.fit(dfTrain, yTrain, dfTrainStd)
print('Training accuracy:', accuracy_score(yTrain, lrr.predict(dfTrain, dfTrainStd)))
print('Test accuracy:', accuracy_score(yTest, lrr.predict(dfTest, dfTestStd)))
print('Probability of Y=1 is predicted as logistic(z) = 1 / (1 + exp(-z))')
print('where z is a linear combination of the following rules/numerical features:')
lrr.explain()

Training accuracy: 0.742536809401594
Test accuracy: 0.7260940032414911
Probability of Y=1 is predicted as logistic(z) = 1 / (1 + exp(-z))
where z is a linear combination of the following rules/numerical features:
```

rule/numerical feature	coefficient
0 (intercept)	-0.0886341
1 MSinceMostRecentInqxc7days > 0.00	0.680261
2 ExternalRiskEstimate	0.654248
3 NetFractionRevolvingBurden	-0.553965
4 NumSatisfactoryTrades	0.551654
5 NumInqLast6M	-0.463226
6 NumBank2NatlTradesWhHighUtilization	-0.448331
7 AverageMinFile <= 52.00	-0.43436
8 NumRevolvingTradesWBalance <= 5.00	0.42154
9 MaxDelq2PublicRecLast12M <= 5.00	-0.418142
10 PercentInstallTrades > 50.00	-0.317566
11 NumSatisfactoryTrades <= 12.00	-0.312471
12 MSinceMostRecentDelq <= 21.00	-0.301566
13 PercentTradesNeverDelq <= 95.00	-0.273924
14 ExternalRiskEstimate > 75.00	0.263437
15 AverageMinFile <= 84.00	-0.182118
16 PercentTradesNeverDelq	0.166518
17 AverageMinFile	0.15069
18 PercentInstallTrades > 42.00	-0.148602
19 NumBank2NatlTradesWhHighUtilization <= 0.00	0.135396
20 MSinceOldestTradeOpen <= 122.00	-0.132409
21 PercentTradesNeverDelq <= 81.00	-0.11771
22 NumSatisfactoryTrades <= 17.00	-0.11022
23 ExternalRiskEstimate > 72.00	0.107613

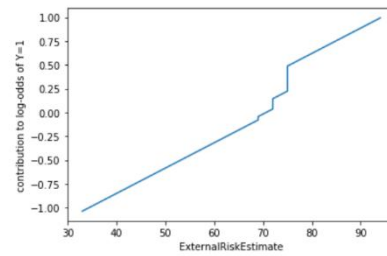
(图形) 显示结论, 如以 GAM 图示 LogRR 的结论

外部风险预估

ExternalRiskEstimate

As expected from the BRCG Boolean rule above, 'ExternalRiskEstimate' is an important feature positively correlated with good credit risk. The jumps in the plot indicate that applicants with above average 'ExternalRiskEstimate' (the mean is 72) get an additional boost.

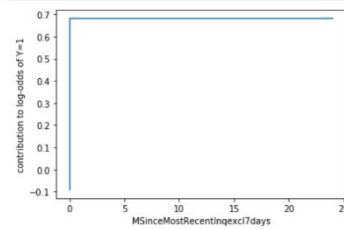
```
lrr.visualize(data, fb, ['ExternalRiskEstimate']);
```



Credit inquiries

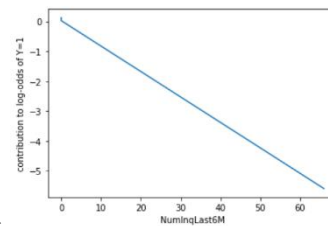
The next two plots illustrate the dependence on the applicant's credit inquiries. The first plot shows a significant penalty for having less than one month since the most recent inquiry ('MSinceMostRecentInqexcl7days' = 0).

```
lrr.visualize(data, fb, ['MSinceMostRecentInqexcl7days']);
```



The second shows that predicted risk increases with the number of inquiries in the last six months ('NumInqLast6M').

```
lrr.visualize(data, fb, ['NumInqLast6M']);
```

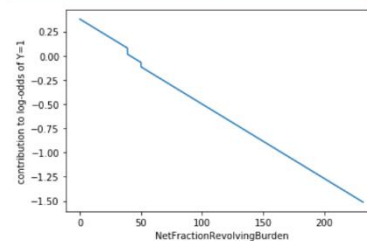


信用查询次数的影

Debt level

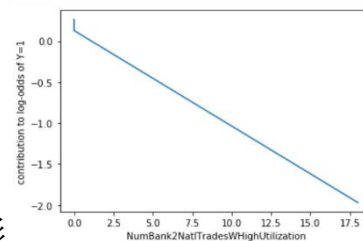
The following four plots relate to the applicant's debt level. 'NetFractionRevolvingBurden' is the ratio of revolving debt (e.g. credit card) balance to credit limit, expressed as a percentage, and has a large negative impact on the probability of good credit. A small fraction of applicants (less than 1%) actually have NetFractionRevolvingBurden greater than 100%, i.e. more revolving debt than their credit limit. This might be investigated further by the data scientist.

```
lrr.visualize(data, fb, ['NetFractionRevolvingBurden']);
```



The second 'NumBank2Nat1TradesWHighUtilization' plot shows that the number of accounts ('trades') with high utilization (high balance relative to credit limit for each account) also has a large impact, with a drop as soon as one account has high utilization.

```
lrr.visualize(data, fb, ['NumBank2Nat1TradesWHighUtilization']);
```



债务水平的影

对于使用这个 AI 系统的银行贷款经理而言，他的关心重点是贷款决定的一致性，是否存在系统性的歧视。同样他也只希望一个指令获得解释，以便增强对模型的信心（即对模型普遍性的解释而不是单个案例的解释），而其依据是手边现有的

案例为支撑（基于现有样本）对于特定结果的（事后）解释。依照上面的树形结构，我们自然能理解为啥银行贷款经理使用 ProtoDash 算法来寻求帮助。

有了这个决策，他按如下步骤开展：加载整理数据、运行算法、图示结论。

<p>加载整理数据</p>	<pre> : heloc = HELOCdataset() df = heloc.dataframe() pd.set_option('display.max_rows', 500) pd.set_option('display.max_columns', 24) pd.set_option('display.width', 1000) print("Size of HELOC dataset:", df.shape) print("Number of \"Good\" applicants:", np.sum(df['RiskPerformance']=='Good')) print("Number of \"Bad\" applicants:", np.sum(df['RiskPerformance']=='Bad')) print("Sample Applicants:") df.head(10).transpose() Using Heloc dataset: c:\users\ronnyluss\aix360\aix360\datasets\..\data\heloc_data\heloc_dataset.csv Size of HELOC dataset: (10459, 24) Number of "Good" applicants: 5000 Number of "Bad" applicants: 5459 Sample Applicants: : </pre> <table border="1"> <thead> <tr> <th></th> <th>0</th> <th>1</th> <th>2</th> <th>3</th> <th>4</th> <th>5</th> <th>6</th> <th>7</th> <th>8</th> <th>9</th> </tr> </thead> <tbody> <tr> <td>ExternalRiskEstimate</td> <td>55</td> <td>61</td> <td>67</td> <td>66</td> <td>81</td> <td>59</td> <td>54</td> <td>68</td> <td>59</td> <td>61</td> </tr> <tr> <td>MSinceOldestTradeOpen</td> <td>144</td> <td>58</td> <td>66</td> <td>169</td> <td>333</td> <td>137</td> <td>88</td> <td>148</td> <td>324</td> <td>79</td> </tr> <tr> <td>MSinceMostRecentTradeOpen</td> <td>4</td> <td>15</td> <td>5</td> <td>1</td> <td>27</td> <td>11</td> <td>7</td> <td>7</td> <td>2</td> <td>4</td> </tr> <tr> <td>AverageMinFile</td> <td>84</td> <td>41</td> <td>24</td> <td>73</td> <td>132</td> <td>78</td> <td>37</td> <td>65</td> <td>138</td> <td>36</td> </tr> <tr> <td>NumSatisfactoryTrades</td> <td>20</td> <td>2</td> <td>9</td> <td>28</td> <td>12</td> <td>31</td> <td>25</td> <td>17</td> <td>24</td> <td>19</td> </tr> <tr> <td>NumTrades60Ever2DerogPubRec</td> <td>3</td> <td>4</td> <td>0</td> <td>1</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>NumTrades90Ever2DerogPubRec</td> <td>0</td> <td>4</td> <td>0</td> <td>1</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>PercentTradesNeverDelq</td> <td>83</td> <td>100</td> <td>100</td> <td>93</td> <td>100</td> <td>91</td> <td>92</td> <td>83</td> <td>85</td> <td>95</td> </tr> <tr> <td>MSinceMostRecentDelq</td> <td>2</td> <td>-7</td> <td>-7</td> <td>76</td> <td>-7</td> <td>1</td> <td>9</td> <td>31</td> <td>5</td> <td>5</td> </tr> <tr> <td>MaxDelq2PublicRecLast12M</td> <td>3</td> <td>0</td> <td>7</td> <td>6</td> <td>7</td> <td>4</td> <td>4</td> <td>6</td> <td>4</td> <td>4</td> </tr> <tr> <td>MaxDelqEver</td> <td>5</td> <td>8</td> <td>8</td> <td>6</td> <td>8</td> <td>6</td> <td>6</td> <td>6</td> <td>6</td> <td>6</td> </tr> <tr> <td>NumTotalTrades</td> <td>23</td> <td>7</td> <td>9</td> <td>30</td> <td>12</td> <td>32</td> <td>26</td> <td>18</td> <td>27</td> <td>19</td> </tr> <tr> <td>NumTradesOpeninLast12M</td> <td>1</td> <td>0</td> <td>4</td> <td>3</td> <td>0</td> <td>1</td> <td>3</td> <td>1</td> <td>1</td> <td>3</td> </tr> <tr> <td>PercentInstallTrades</td> <td>43</td> <td>67</td> <td>44</td> <td>57</td> <td>25</td> <td>47</td> <td>58</td> <td>44</td> <td>26</td> <td>26</td> </tr> <tr> <td>MSinceMostRecentInqexc17days</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>NumInqLast6M</td> <td>0</td> <td>0</td> <td>4</td> <td>5</td> <td>1</td> <td>0</td> <td>4</td> <td>0</td> <td>1</td> <td>6</td> </tr> <tr> <td>NumInqLast6Mexcl7days</td> <td>0</td> <td>0</td> <td>4</td> <td>4</td> <td>1</td> <td>0</td> <td>4</td> <td>0</td> <td>1</td> <td>6</td> </tr> <tr> <td>NetFractionRevolvingBurden</td> <td>33</td> <td>0</td> <td>53</td> <td>72</td> <td>51</td> <td>62</td> <td>89</td> <td>28</td> <td>68</td> <td>31</td> </tr> <tr> <td>NetFractionInstallBurden</td> <td>-8</td> <td>-8</td> <td>66</td> <td>83</td> <td>89</td> <td>93</td> <td>76</td> <td>48</td> <td>-8</td> <td>86</td> </tr> <tr> <td>NumRevolvingTradesWBalance</td> <td>8</td> <td>0</td> <td>4</td> <td>6</td> <td>3</td> <td>12</td> <td>7</td> <td>2</td> <td>7</td> <td>5</td> </tr> <tr> <td>NumInstallTradesWBalance</td> <td>1</td> <td>-8</td> <td>2</td> <td>4</td> <td>1</td> <td>4</td> <td>7</td> <td>2</td> <td>1</td> <td>3</td> </tr> <tr> <td>NumBank2NatTradesWHighUtilization</td> <td>1</td> <td>-8</td> <td>1</td> <td>3</td> <td>0</td> <td>3</td> <td>2</td> <td>2</td> <td>3</td> <td>1</td> </tr> <tr> <td>PercentTradesWBalance</td> <td>69</td> <td>0</td> <td>86</td> <td>91</td> <td>80</td> <td>94</td> <td>100</td> <td>40</td> <td>90</td> <td>62</td> </tr> <tr> <td>RiskPerformance</td> <td>Bad</td> <td>Bad</td> <td>Bad</td> <td>Bad</td> <td>Bad</td> <td>Bad</td> <td>Good</td> <td>Good</td> <td>Bad</td> <td>Bad</td> </tr> </tbody> </table>		0	1	2	3	4	5	6	7	8	9	ExternalRiskEstimate	55	61	67	66	81	59	54	68	59	61	MSinceOldestTradeOpen	144	58	66	169	333	137	88	148	324	79	MSinceMostRecentTradeOpen	4	15	5	1	27	11	7	7	2	4	AverageMinFile	84	41	24	73	132	78	37	65	138	36	NumSatisfactoryTrades	20	2	9	28	12	31	25	17	24	19	NumTrades60Ever2DerogPubRec	3	4	0	1	0	0	0	0	0	0	NumTrades90Ever2DerogPubRec	0	4	0	1	0	0	0	0	0	0	PercentTradesNeverDelq	83	100	100	93	100	91	92	83	85	95	MSinceMostRecentDelq	2	-7	-7	76	-7	1	9	31	5	5	MaxDelq2PublicRecLast12M	3	0	7	6	7	4	4	6	4	4	MaxDelqEver	5	8	8	6	8	6	6	6	6	6	NumTotalTrades	23	7	9	30	12	32	26	18	27	19	NumTradesOpeninLast12M	1	0	4	3	0	1	3	1	1	3	PercentInstallTrades	43	67	44	57	25	47	58	44	26	26	MSinceMostRecentInqexc17days	0	0	0	0	0	0	0	0	0	0	NumInqLast6M	0	0	4	5	1	0	4	0	1	6	NumInqLast6Mexcl7days	0	0	4	4	1	0	4	0	1	6	NetFractionRevolvingBurden	33	0	53	72	51	62	89	28	68	31	NetFractionInstallBurden	-8	-8	66	83	89	93	76	48	-8	86	NumRevolvingTradesWBalance	8	0	4	6	3	12	7	2	7	5	NumInstallTradesWBalance	1	-8	2	4	1	4	7	2	1	3	NumBank2NatTradesWHighUtilization	1	-8	1	3	0	3	2	2	3	1	PercentTradesWBalance	69	0	86	91	80	94	100	40	90	62	RiskPerformance	Bad	Bad	Bad	Bad	Bad	Bad	Good	Good	Bad	Bad
	0	1	2	3	4	5	6	7	8	9																																																																																																																																																																																																																																																																										
ExternalRiskEstimate	55	61	67	66	81	59	54	68	59	61																																																																																																																																																																																																																																																																										
MSinceOldestTradeOpen	144	58	66	169	333	137	88	148	324	79																																																																																																																																																																																																																																																																										
MSinceMostRecentTradeOpen	4	15	5	1	27	11	7	7	2	4																																																																																																																																																																																																																																																																										
AverageMinFile	84	41	24	73	132	78	37	65	138	36																																																																																																																																																																																																																																																																										
NumSatisfactoryTrades	20	2	9	28	12	31	25	17	24	19																																																																																																																																																																																																																																																																										
NumTrades60Ever2DerogPubRec	3	4	0	1	0	0	0	0	0	0																																																																																																																																																																																																																																																																										
NumTrades90Ever2DerogPubRec	0	4	0	1	0	0	0	0	0	0																																																																																																																																																																																																																																																																										
PercentTradesNeverDelq	83	100	100	93	100	91	92	83	85	95																																																																																																																																																																																																																																																																										
MSinceMostRecentDelq	2	-7	-7	76	-7	1	9	31	5	5																																																																																																																																																																																																																																																																										
MaxDelq2PublicRecLast12M	3	0	7	6	7	4	4	6	4	4																																																																																																																																																																																																																																																																										
MaxDelqEver	5	8	8	6	8	6	6	6	6	6																																																																																																																																																																																																																																																																										
NumTotalTrades	23	7	9	30	12	32	26	18	27	19																																																																																																																																																																																																																																																																										
NumTradesOpeninLast12M	1	0	4	3	0	1	3	1	1	3																																																																																																																																																																																																																																																																										
PercentInstallTrades	43	67	44	57	25	47	58	44	26	26																																																																																																																																																																																																																																																																										
MSinceMostRecentInqexc17days	0	0	0	0	0	0	0	0	0	0																																																																																																																																																																																																																																																																										
NumInqLast6M	0	0	4	5	1	0	4	0	1	6																																																																																																																																																																																																																																																																										
NumInqLast6Mexcl7days	0	0	4	4	1	0	4	0	1	6																																																																																																																																																																																																																																																																										
NetFractionRevolvingBurden	33	0	53	72	51	62	89	28	68	31																																																																																																																																																																																																																																																																										
NetFractionInstallBurden	-8	-8	66	83	89	93	76	48	-8	86																																																																																																																																																																																																																																																																										
NumRevolvingTradesWBalance	8	0	4	6	3	12	7	2	7	5																																																																																																																																																																																																																																																																										
NumInstallTradesWBalance	1	-8	2	4	1	4	7	2	1	3																																																																																																																																																																																																																																																																										
NumBank2NatTradesWHighUtilization	1	-8	1	3	0	3	2	2	3	1																																																																																																																																																																																																																																																																										
PercentTradesWBalance	69	0	86	91	80	94	100	40	90	62																																																																																																																																																																																																																																																																										
RiskPerformance	Bad	Bad	Bad	Bad	Bad	Bad	Good	Good	Bad	Bad																																																																																																																																																																																																																																																																										
<p>运行算法</p>																																																																																																																																																																																																																																																																																				
<p>1) 预处理训练数据</p>	<pre> # Clean data and split dataset into train/test (Data, x_train, x_test, y_train_b, y_test_b) = heloc.split () Z = np.vstack((x_train, x_test)) Zmax = np.max(Z, axis=0) Zmin = np.min(Z, axis=0) #normalize an array of samples to range [-0.5, 0.5] def normalize(V): VN = (V - Zmin)/(Zmax - Zmin) VN = VN - 0.5 return(VN) </pre>																																																																																																																																																																																																																																																																																			

	<pre> # rescale a sample to recover original values for normalized values. def rescale(X): return(np.multiply (X + 0.5, (Zmax - Zmin)) + Zmin) N = normalize(Z) xn_train = N[0:x_train.shape[0], :] xn_test = N[x_train.shape[0]:, :] </pre>
2) 定义和训练模型	<pre> # nn with no softmax def nn_small(): model = Sequential() model.add(Dense(10, input_dim=23, kernel_initializer='normal', activation='relu')) model.add(Dense(2, kernel_initializer='normal')) return model # Set random seeds for repeatability np.random.seed(1) tf.set_random_seed(2) class_names = ['Bad', 'Good'] # loss function def fn(correct, predicted): return tf.nn.softmax_cross_entropy_with_logits(labels=correct, logits=predicted) # compile and print model summary nn = nn_small() </pre>

```

nn.compile(loss=fn, optimizer='adam', metrics=['accuracy
'])
nn.summary()

# train model or load a trained model
TRAIN_MODEL = False

if (TRAIN_MODEL):
    nn.fit(xn_train, y_train_b, batch_size=128, epochs=50
0, verbose=1, shuffle=False)
    nn.save_weights("heloc_nnsml.h5")
else:
    nn.load_weights("heloc_nnsml.h5")

# evaluate model accuracy
score = nn.evaluate(xn_train, y_train_b, verbose=0) #Compute
training set accuracy
#print('Train loss:', score[0])
print('Train accuracy:', score[1])

score = nn.evaluate(xn_test, y_test_b, verbose=0) #Compute
test set accuracy
#print('Test loss:', score[0])
print('Test accuracy:', score[1])

```

3) 验证现有数据集中类似案例将得到类似结论（即验证模型的一致性）

比如，我们选择案例#8，其结论是 Good，可以放款，其特征及如右图所示。

那么，数据集中有同样结论的其它案例有哪些？他们是否也有同样的特征分布？如果有，则表明这个模型具有一致性。

我们先获得所有结论是 Good 的案例；

计算这些典型用户和#8 的相似度；

我们获得了如下的相似度对比表，第 0 列是我们的选定案例#8，其它四个是数据集中有同样 Good 结论的其它案例。显然，过半数的特征的是接近的。仔细研究对照表，贷款经理发现，能获得 Good 结论的客户都是没有负债的客户，这个发现让贷款经理对系统的结论更有信息了。

ExternalRiskEstimate	82
MSinceOldestTradeOpen	280
MSinceMostRecentTradeOpen	13
AverageMinFile	102
NumSatisfactoryTrades	22
NumTrades60Ever2DerogPubRec	0
NumTrades90Ever2DerogPubRec	0
PercentTradesNeverDelq	91
MSinceMostRecentDelq	26
MaxDelq2PublicRecLast12M	6
MaxDelqEver	6
NumTotalTrades	23
NumTradesOpeninLast12M	0
PercentInstallTrades	9
MSinceMostRecentInqexc7days	0
NumInqLast6M	0
NumInqLast6Mexcl7days	0
NetFractionRevolvingBurden	3
NetFractionInstallBurden	0
NumRevolvingTradesWBalance	4
NumInstallTradesWBalance	1
NumBank2NatITradesWHighUtilization	1
PercentTradesWBalance	42

	0	1	2	3	4
ExternalRiskEstimate	0.59	0.29	0.42	0.84	0.21
MSinceOldestTradeOpen	0.76	0.62	0.76	0.09	0.79
MSinceMostRecentTradeOpen	1.00	0.09	0.83	0.89	0.87
AverageMinFile	0.79	0.09	0.90	1.00	0.82
NumSatisfactoryTrades	0.95	0.39	0.74	0.39	0.15
NumTrades60Ever2DerogPubRec	1.00	1.00	0.08	1.00	1.00
NumTrades90Ever2DerogPubRec	1.00	1.00	0.08	1.00	1.00
PercentTradesNeverDelq	1.00	0.15	0.81	0.15	0.15
MSinceMostRecentDelq	1.00	0.36	0.22	0.36	0.36
MaxDelq2PublicRecLast12M	1.00	0.13	1.00	0.13	1.00
MaxDelqEver	1.00	0.41	0.17	0.41	0.64
NumTotalTrades	0.80	0.23	0.86	0.26	0.35
NumTradesOpeninLast12M	1.00	1.00	0.40	0.40	0.06
PercentInstallTrades	1.00	0.05	0.54	0.37	0.33
MSinceMostRecentInqexc7days	0.08	1.00	1.00	1.00	1.00
NumInqLast6M	0.21	1.00	0.21	0.21	0.04
NumInqLast6Mexcl7days	0.26	1.00	0.26	1.00	0.07
NetFractionRevolvingBurden	0.96	0.88	0.96	0.92	0.09
NetFractionInstallBurden	1.00	1.00	1.00	1.00	0.08
NumRevolvingTradesWBalance	1.00	0.28	0.38	0.73	0.20
NumInstallTradesWBalance	1.00	0.13	1.00	0.13	1.00
NumBank2NatITradesWHighUtilization	0.69	0.69	0.69	1.00	0.11
PercentTradesWBalance	0.67	0.12	0.36	0.38	0.57

<p>类似地，我们也可以对结论是 Bad（贷款申请被拒绝）的客户用同样的步骤做同样的研究</p>	<p>结果也同样是对于指定的样本客户，同样获得 Bad 结论的客户，其特征指标中过半数非常接近。贷款经理仔细研究这些接近指标发现，这些被判 Bad 的客户，大都有轻微犯罪前科。为此，贷款经理在处理这类客户时就要额外小心。</p>
--	--

对于受这个 AI 系统影响的客户而言（即贷款申请人），尤其是申请被拒的人，需要了解哪些因素是关键，从而他们可以采取行动改善自己的财经状况，以便日后有机会成功申请。其步骤也是类似：加载整理数据、运行算法、显示结论。这个就不赘述了。

要亲自上手试验，请参考 AIX360 网站 <http://aix360.mybluemix.net>。

1148, 生物信息学的可解释人工智能：方法、工具和应用

Explainable AI for Bioinformatics: Methods, Tools, and Applications

Md. Rezaul Karim, Tanhim Islam 等

（北大严丽君摘录，2023.2.23）

利用深度神经网络（DNNs）和机器学习（ML）算法的人工智能（AI）系统被广泛用于解决生物信息学、生物医学信息学和精准医学的重要问题。然而复杂的 DNNs 或 ML 模型，往往被认为是不透明的和黑箱的，会使人难以理解其决策背后的推理。这种缺乏透明度的情况对终端用户和决策者以及人工智能开发者来说都是一种挑战。此外，在医疗保健等敏感领域，对于可能对人类生活产生重大影响的人工智能系统，可解释性和问责制不仅是可取的，而且是法律要求的。公平性是另一个日益受到关注的问题，因为算法决策不应该根据敏感属性对某些群体或个人表现出偏见或歧视。可解释的人工智能（XAI）旨在克服黑箱模型的不透明性，并为人工智能系统如何做出决定提供透明度。可解释的 ML 模型可以解释它们如何进行预测以及影响其结果的因素。然而，大多数最先进的可解释的 ML 方法是不分领域的，并且是从计算机视觉、自动推理或统计学等领域发展而来的，这使得在没有定制和特定领域适应的情况下直接应用于生物信息学问题具有挑战性。在本文中，我们讨论了生物信息学背景下可解释性的重要性，提供了一个针对模

型和模型诊断的可解释 ML 方法和工具的概述，并概述了它们潜在的注意事项和缺点。此外，我们讨论了如何为生物信息学问题定制现有的可解释 ML 方法。我们通过生物成像、癌症基因组学和文本挖掘的案例研究，展示了 XAI 方法如何能够提高透明度。

1149, 可解释人工智能 (XAI) 文献趋势

Trends in Explainable AI (XAI) Literature

Alon Jacovi (巴伊兰大学, Ilan University)

(北大李因立摘录, 2023. 1. 13)

XAI 的文献研究在术语和出版地点上是分散的,但近年来该社区在关键词上趋于一致,这使得自动发现论文更加可靠。我们使用 SemanticScholar API 进行关键词搜索,并进行手动筛选,收集了一个格式良好且相对全面的 5199 篇 XAI 论文集合,可在此找到。我们利用这个集合来澄清和可视化关于文献规模和范围、引用趋势、跨领域趋势和合作趋势的动态。总体而言,XAI 越来越多元化,与此同时,不同科学领域(非计算机科学领域)的论文数量相对增长,跨领域的合著作者和引用活动也在增加。此外,该集合还可以作为一个论文发现引擎,根据特定的约束条件检索被引用的 XAI 文献(例如,在其领域之外有影响力的论文,或对非 XAI 研究有影响力的论文)。

1150, 解释性方法评价的实验研究

An Experimental Investigation into the Evaluation of Explainability Methods

Sédric Stassin, Alexandre Englebort, Géraldine Nanfack 等

(北大李劲楠摘录, 2023. 5. 25)

可解释人工智能 (XAI) 旨在帮助用户掌握人工智能 (AI) 系统预测背后的推理。近年来出现了许多 XAI 方法。因此,与 XAI 方法评估相关的子领域获得了相当大的关注,目的是使用各种方法和标准来确定哪些方法提供了最佳解释。然而,文献中缺乏评估指标本身的比较,可以用来评估 XAI 方法。这项工作旨在通过比较 14 种不同的指标来填补这一空白,这些指标应用于九种最先进的 XAI 方法和三

种用作参考的伪方法（例如，随机显著性图）。实验结果表明，这些指标中的哪一个产生了高度相关的结果，表明了潜在的冗余。我们还证明了改变基线超参数对评估度量值的显著影响。最后，我们使用虚拟方法来评估。

1151, 利用证据图结构特征解释基于 GNN 的 IDS 检测结果

Interpreting GNN-based IDS Detections Using Provenance Graph Structural Features

Kunal Mukherjee, Joshua Wiedemeier, Tianhao Wang 等

（北大张凯迪摘录，2023.6.1）

复杂的基于神经网络(NN)的模型的黑箱性质阻碍了它们在安全领域的广泛采用，因为缺乏对其预测的逻辑解释和可操作的后续行动。为了提高系统出处分析中使用的图神经网络(GNN)安全模型的透明度和责任感，我们提出了 PROVEXPLAINER，这是一个将抽象的 GNN 决策边界投射到可解释的特征空间的框架。

我们首先使用更简单和可解释的模型，如决策树 (DT)，复制基于 GNN 的安全模型的决策过程。为了最大限度地提高代用模型的准确性和保真度，我们提出了建立在经典图论基础上的新型图结构特征，并通过对安全领域知识的广泛数据研究加以强化。我们的图结构特征与系统出处领域的问题空间行动紧密相连，这使得检测结果可以用描述性的人类语言来解释。PROVEXPLAINER 允许简单的 DT 模型在具有一般图结构特征的程序分类任务上达到对 GNN 的 95%的保真度，而在具有为直接解释而定制的特定任务特征包的恶意软件检测任务上达到 99%的保真度。通过对三个恶意软件家族中五个真实世界的恶意软件样本的案例研究，证明了对恶意软件分类的解释。

1152, 高度定制的鲁棒印刷神经形态电路

Highly-Bespoke Robust Printed Neuromorphic Circuits

Haibin Zhao, Brojogopal Sapui, Michael Hefenbrock, 等

（北大周鑫摘录，2022.10）

随着物联网、智能快消产品、可穿戴设备的快速发展，对材料柔性、无毒、低成本等要求提出了迫切的要求。然而，这些要求通常超出了传统的刚性硅技术的可

制造范围。在这方面，印刷电子产品提供了一个有前途的选择。与神经形态计算相结合，印刷神经形态电路不能够实现上述特性，而且还弥补了传统印刷电子产品的一些弱点，例如高度定制化，低设备计数和高延迟等等。一般来说，（印刷）神经形态电路通过印刷电阻横杆来模拟矩阵乘法来表达其功能，非线性电路来表达激活函数。前者的值通常是后天习得的，而后者是事先设计好的，并在训练中被认为是固定的。

所有的任务。印刷电子产品的增材制造特性允许设计高度定制的设计。在打印神经形态电路的情况下，电路针对特定的数据集进行优化。此外，作者探索了一种方法，不仅可以学习交叉电阻的值，还可以为定制实现学习非线性组件的参数化。在提供要表达的功能的额外灵活性的同时，这也将允许增加对打印变化的鲁棒性。为了充分利用聚乙烯的另一个独特特征，即高度集成化的增材制造，本研究将重点放在印刷神经形态电路中的定制非线性电路上。此外，通过引入 pnn 的变化感知训练，考虑了交叉杆和非线性电路中元件的打印变化，旨在提高印刷神经形态电路的鲁棒性。初步实验证明了这一点，通过引入这两种方法均显著提高了 pnn 的精度和鲁棒性。研究表明，可学习的非线性电路对最终精度的提高有独特的贡献，而鲁棒性主要由变化感知训练提。在电路元件变化 10% 的情况下，印制神经形态电路的精度和鲁棒性分别提高 26% 和 75%。

1153, Alexewicz 拓扑中的 Spiking 神经网络：分析和误差界的新视角

Spiking Neural Networks in the Alexewicz Topology: A New Perspective on Analysis and Error Bounds

Bernhard A. Moser, Michael Lunglmayr

（北大杨思恒摘录，2023.06.02）

为了简化神经形态计算中误差传播的分析，并更好地理解尖峰神经网络（SNN），作者将 SNN 的数学分析问题解决为将尖峰序列映射到尖峰序列的自同态。一个中心问题是尖峰序列空间的适当结构及其对 SNN 误差测量设计的影响，包括时间延迟、阈值偏差和泄漏积分与激发（LIF）神经元模型的重新初始化模式的设计。首先，作者通过分析 LIF 模型的所有亚阈值信号的闭包来识别底层拓扑。对于零泄漏，这种方法产生 Alexewicz 拓扑，作者将其用于具有任意正泄漏的 LIF 神经

元。因此，LIF 可以被理解为对应范数中的尖峰序列量化。通过这种方式，作者获得了各种误差界和不等式，例如传入和传出尖峰序列之间的拟等距关系。另一个结果是误差传播的 Lipschitz 型全局上界和相关的共振型现象。

1154, 基于人工突触器件的神经形态硬件

Neuromorphic Hardware Based on Artificial Synaptic Devices

Department of School of Advanced Technology, Xi'an Jiaotong-Liverpool University, Suzhou, China;

(北大王世界摘录, 2022.10)

本篇文章是受大脑结构和功能的启发，然后在此基础上对人工突触设备和神经形态计算进行了研究。基于人工突触装置，神经形态计算硬件将变得更加“像大脑”，回顾之前基于突触装置的研究进展神经形态计算发现，我们的大脑将信息存储为突触的重量和过程通过连续突触权重更新的输入刺激。与此相同的是，在神经形态计算计算数据的输入矩阵也需要数字化，并且通过其中的权重更新处理反向传播中的数据。由于实现了数字计算在数字设备上，例如 CPU 和 GPU，它们在训练和推断周期期间会消耗大量能量，与我们的大脑相比，这是非常低效的复杂的模拟智能，功耗为不超过 20W。但是相比之下，纵横制阵列用人工突触装置构建，如电阻随机存取存储器 (RRAM) 和突触晶体管是能够将重量参数存储在设备内自身，并通过欧姆定律、基尔霍夫定律和其他各种电路定律显著降低功耗。此外，直接通过以下方式复制某个自然神经网络的模式，从生物大脑复制连接图将数据导入神经形态计算设备也可能降低神经网络的训练阶段的复杂性和功耗。这篇文章中，主要对基于的神经形态计算硬件研究进展人工突触装置，并涵盖了关于硬件网络和大脑模式的复制。本文简要回顾了几件标志性的作品基于人工突触的神经形态进化覆盖较低级别电路实现的计算，全硬件神经网络及其复制粘贴实例生物大脑。前者提供了思想和框架用于设计突触计算阵列，以及后者提出了解决方案，并且有可能工业化。我们有理由相信，在不久的将来，神经形态计算设备和系统将不断发展商业化和大众化，节省了大量的电力、计算能力，并逐渐改变人工智能。

1155, 异步脉冲神经网络+神经拟态计算系统

摘自 COPU《异步脉冲神经网络+神经拟态计算系统》专辑

(2022. 4. 11)

研发基于异步脉冲神经网络的神经拟态计算系统, 实现由弱人工智能向强人工智能转型之路, 在这方面表现突出的有:

英特尔研制的 Loihi 芯片, 由 Loihi 芯片集成的 Pohoiki Springs 神经拟态网络 (由 1 亿神经元组成), 及神经拟态计算系统, 2019 年已制成原型机, 由英特尔团队合作开发应用生态。

由曼彻斯特大学团队合作研发由 10 亿个神经元组成的神经拟态网络, 2018 年 11 月已制成 Spin Naker 类脑计算机原型机。

浙江大学自主研发达尔文-2 芯片, 由芯片集成 1.2 亿个神经元组成的神经拟态网络, 2020 年 9 月发布已研制成功 Darwin Mouse 类脑计算机原型机, 研发工作正在继续跟进。

丹麦科技大学与亚琛科技大学合作也在研发基于脉冲神经网络的神经拟态计算系统 (包括探索建立神经拟态学习规则的库, 多模态神经脚本知识模型等)。

1156, 英特尔公司、曼彻斯特大学、浙江大学研发神经拟态计算系统情况

摘自 COPU《异步脉冲神经网络+神经拟态计算系统》专辑

(2022. 4. 11)

1) 英特尔公司

2017 年研发 Loihi 第一代芯片, 制程 14nm, 单颗芯片 13 万个神经元

2021 年研发 Loihi 第二代芯片, 制程 7nm, 单颗芯片 100 万个神经元。

在推出第二代神经拟态研究芯片 Loihi-2 的同时, 推出用于开发神经启发应用的开源软件框架 Lava。

Lava 软件框架是一个开放、模块化、可扩展的框架, 能够在跨越异构架构上无缝运行, 实现跨平台执行, 以及各种人工智能、神经拟态和机器人框架的互操作性。

第一代 Loihi 芯片使神经拟态处理速度较之传统计算机提高 1000 倍，能耗下降 10000 倍；第二代 Loihi 芯片和 Lava 框架极大地提高了神经拟态处理速度，较第一代 Loihi 芯片又提高了 10 倍（为传统计算机运算速度的 10000 倍）。

2019 年 3 月 19 日，英特尔将 768 颗 Loihi-1 芯片集成，组成 Pohoiki Springs 神经拟态网络（由 1 亿个神经元组网），拥有平行处理性能和异步信号传输能力。

2019 年制成基于 pohokisprings 神经拟态网络的 1 亿神经元的类脑计算机原型机。

脉冲神经网络（SNNs）是神经拟态计算中一种全新的模型，可以模仿人脑中自然神经元网络方式，将计算模块重新分布，神经拟态在芯片上计算、存储单元一体化，颠覆传统冯诺依曼的架构，实现类脑的智能功能，神经拟态计算被认为是引领下一代人工智能的主要计算模式。

神经拟态计算未来主要的应用前景，一是基于神经拟态计算开展超算应用开发，二是用于非结构化数据、实时的场景，三是用于多模态的、实时的场景（如机器人、无人机需要持续学习、自适应学习的场景）。为了扩大应用，英特尔于 2020 年 3 月成立神经拟态研究社区（INRC），推动神经拟态算法、软件和应用程序的研发，参加 INRC 的除英特尔外、还有 IBM、HP、MIT、普渡大学、斯坦福大学、以及联想、埃森哲、罗技、梅赛德斯奔驰、机器视觉传感器公司 Prophesee，还有罗格斯大学、戴尔施特理工大学等。

2) 曼彻斯特大学

自 2016 年始，曼彻斯特大学与南安普顿大学、剑桥大学、谢菲尔德大学、ARM、Silistix 公司、泰雷兹等合作共同研发脉冲神经网络(芯片)，并在芯片集成的基础上研发神经拟态计算系统（和类脑计算机）。

每颗第一代 SpinNaker-1 芯片支持 1.6 万个神经元，曼大将 SpinNaker-1 芯片集成：30k x 16k=4.6 亿个神经元，由第二代 SpinNaker-2 芯片集成为 10 亿个神经元。由 10 亿个神经元组成的神经拟态网络，并在此基础上研发 SpinNaker 类脑计算机（神经拟态+平行处理），原型机于 2018 年 11 月制成对外发布。

3) 浙江大学

2015 年自主研发达尔文-1 芯片，2019 年自主研发达尔文-2 芯片（每颗芯片 576 个内核，每颗芯片支持 15 万个神经元），2019 年 8 月，浙江大学联合之江实验室发布由 792 颗芯片集成的神经拟态网络（由 1.2 亿个神经元、近千亿个突触组网）。

2020 年 9 月发布我国首台基于亿级神经元组成神经拟态网络，并由此研制而成的 Darwin Mouse 类脑计算机原型机，功耗为 350-500w，由三个高为 1.6m 的标准机柜组成。

该原型机已在：（1）抗洪抢险场景下支持、调控机器人协调工作，（2）建立脑外侧膝状核的神经网络模型，（3）建立学习-记忆融合模型。

该机运行机制与 SpinNaker 相似，研发应用场景正在起步。

从上面 3 例看来，首先设计拟态芯片，实行芯片集成，生成 1 亿个-10 亿个神经元，组成亿级神经拟态网络，支持研制神经拟态计算系统（即类脑计算机），制成原型机，颠覆传统的冯诺依曼架构，形成神经拟态计算（类脑计算）架构，随即组成开发应用生态系统（英特尔公司、曼彻斯特大学两例均集结一批合作伙伴共同开发）。

1157，基于异步脉冲神经网络的神经拟态计算与传统计算运行方式比较表

摘自 COPU 《异步脉冲神经网络+神经拟态计算系统》专辑

2022. 4. 11

	信息源	编码方式	传递方式	基于不同特征的计算方式
传统计算运行方式	数字信号 计算机程序	由数字源代码变换 0,1 的机器码进行编码	传统网络权重连接+激活方式，对机器码进行信息传递	传统计算机运行方式，符合冯·诺依曼硬件和软件计算架构
神经拟态计算运行方式	神经电脉冲信号和化学信号	采用稀疏脉冲时序编码机制	通过由仿脑自然神经元+突触组成的脉冲神经网络和运行方式模型，对神经电脉冲进行信息传递	神经拟态计算系统运行方式，打破冯·诺依曼硬件和软件计算架构

1158，类脑计算（神经拟态计算）与传统硅基计算、传递、运行模式

摘自 COPU 《异步脉冲神经网络+神经拟态计算系统》专辑，2022. 4. 11

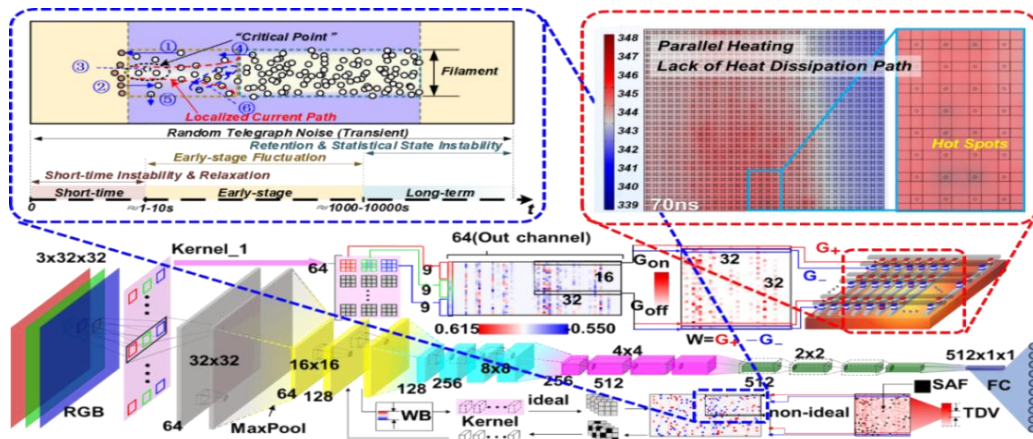
①前者的信息源是神经电脉冲信号和化学信号，后者为数字信号

- ②前者编码方式采用稀疏脉冲时序编码机制,后者由数字源代码变换为 0, 1 的机器码
- ③前者信息传递方式通过模仿人脑自然神经元+突触组成的脉冲神经网络和运行方式模型,对神经电脉冲进行信息传递,后者采用传统网络权重连接+激活方式,对机器码进行信息传递
- ④前者计算(神经元)和存储(突触)是一体化的、融合在一起的,后者计算(处理单元)和存储(存储单元)是分离的
- ⑤前者存在三维广泛连通性,后者无法模拟三维连通,受限于二位连接
- ⑥前者是基于脉冲的事件驱动型的随机计算,后者为了构建确定性计算采用晶体管间布尔代数电路开关
- ⑦前者类脑是模拟生物学上自主低能耗类型(计算/处理、传输),实行节能型,后者能耗大
- ⑧前者模仿人脑的大规模进行通信方式(结构),将数十亿信息同时送到数千个不同目的地,后者也可采用并行通信方式,但传送信息为通过标准网络(从 A 点到 B 点),发送大量信息进行通信
- ⑨前者其运行方式符合类脑计算(神经拟态计算)系统的运行方式(或神经拟态计算架构),后者其运行方式符合冯·诺依曼计算架构。

1159, 阻变器件——存算一体的类脑芯片

北京大学黄如院士研发团队

北大微纳电子研究院黄如院士在第 66 届国际电子器件大会 (IEDM) 上 (2020. 12. 12-18) 上发表他们研发的神经形态器件 (类脑智能芯片) 的论文。阻变器件是后摩尔时代构建新型存算一体及类脑芯片、突破冯-诺依学体系结构瓶颈的关键电子器件技术之一。但阻变器件的非理想效应以及高密度集成带来的热效应会相互耦合,成为阻变器件在存储及神经形态计算应用 1 中的关键挑战。蔡一茂教授、黄如院士团队系统研究了阻变器件非理想效应的物理机制,提出了准确描述多种非理想效应的集约模型,建立了能够综合评估器件技术、阵列拓扑及算法设计的跨层次验证平台,掌握了非理想效应和热串扰对存储及神经形态计算应用的影响,为器件-阵列-算法的协同优化设计提供了重要指导。



1160, 忆阻器存算一体高效芯片研制成功

(清华大学公众号 2023. 10. 9 消息)

清华大学集成电路学院院长吴华强教授团队研制出全球首颗(?) 全系统集成的忆阻存算一体化高级芯片, 可望在人工智能、自动驾驶、可穿戴设备上应用, 该研发成果已在《Science》期刊上发表。

1161, 也谈忆阻器

(某教授来函, 2023. 10. 18)

忆阻器 (Memristor) 作为在一颗芯片上制成的以存算功能一体化最有潜力逻辑运算的器件, 早在 1971 年, 全球首次提出了忆阻器概念, 2008 年便有研制成功的成果 (发表在《Nature》杂志上)。以后国内外不断推出不少研制成果。

1162, 基于多阵列忆阻器存算一体系统研发成功

摘自 COPU 《异步脉冲神经网络+神经拟态计算系统》人工智能汇集 31-345, P17, 2020. 2. 27

2020 年 2 月 27 日, 清华大学清华大学钱鹤、吴华强教授团队成功研发出一款基于多阵列忆阻器存算一体系统, 在处理卷积神经网络 (CNN) 时能效比前沿的图形处理器芯片 (GPU) 高两个数量级。

1163, Chat GPT 发布

COPU 摘要, 2023. 2. 8

自 OpenAI 公司于 2022 年 11 月发布对话式人工智能聊天机器人 ChatGPT 以来，ChatGPT 便在全球引起火爆。

据介绍，ChatGPT①可以进行简单的对话，像人类一样聊天交流，②根据上下文进行互动，完成写邮件、写诗、写代码、编制视频脚本、翻译各种语言文字、做数学题等任务，③通过大数据整合，在自行学习后可给出合理的商业方案、营销方案、供应链方案等，转向资本运作。

ChatGPT 有多牛？

马斯克赞美它：“好的吓人”，库克赞美它：“不可思议”，黄仁勋赞美它“堪比智能手机问世”，比尔盖茨赞美它：“不亚于互联网的诞生”。

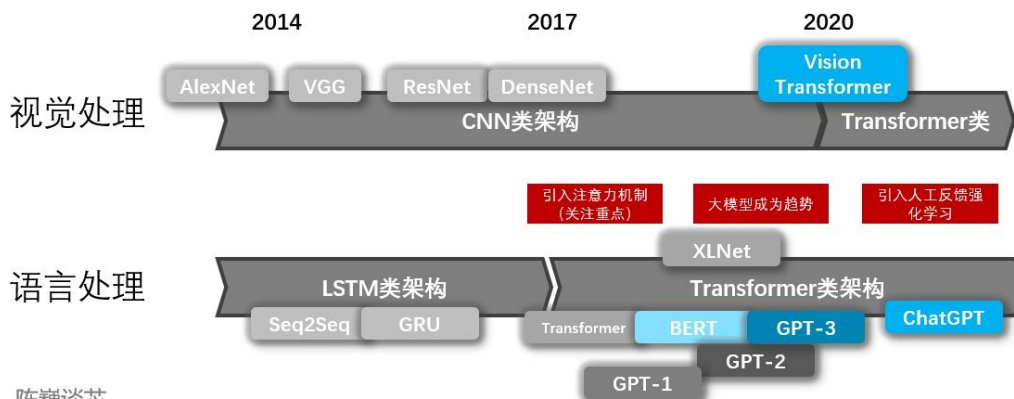
月活用户达到一亿，ChatGPT 仅仅用了两个月（是史上增长最快的应用程序），相对而言：iTunes 用了 6 年半，Twitter 用了 5 年，Facebook 用了 4 年半，Whats APP 用了 3 年半，Instagram 用了 2 年半，TikTok 用了 9 个月。

1164, OpenAI and hugging face

（北大荆琦副教授摘录）

ChatGPT 发展历程、原理、技术架构详解和产业未来（收录于 GPT-4/ChatGPT 技术与产业分析）

GPT 发展历程概览



陈巍谈芯

模型	发布时间	参数量	预训练数据量
GPT-1	2018年6月	1.17亿	约5GB
GPT-2	2019年2月	15亿	40G
GPT-3	2020年5月	1750亿	45TB
ChatGPT	2022年11月	千亿级?	百T级?

OpenAi 家族

我们首先了解下 **OpenAI** 是哪路大神。

OpenAI 总部位于旧金山，由特斯拉的马斯克、Sam Altman 及其他投资者在 2015 年共同创立，目标是开发造福全人类的 AI 技术。而马斯克则在 2018 年时因公司发展方向分歧而离开。

此前，OpenAI 因推出 **GPT 系列自然语言处理模型**而闻名。从 2018 年起，OpenAI 就开始发布生成式预训练语言模型 GPT (Generative Pre-trained Transformer)，可用于生成文章、代码、机器翻译、问答等各类内容。

每一代 GPT 模型的参数量都**爆炸式增长**，堪称“越大越好”。2019 年 2 月发布的 GPT-2 参数量为 15 亿，而 2020 年 5 月的 GPT-3，参数量达到了 1750 亿。

模型	发布时间	参数量	预训练数据量
GPT-1	2018年6月	1.17亿	约5GB
GPT-2	2019年2月	15亿	40G
GPT-3	2020年5月	1750亿	45TB
ChatGPT	2022年11月	千亿级?	百T级?

ChatGPT 的主要特点

ChatGPT 是基于 GPT-3.5 (Generative Pre-trained Transformer 3.5) 架构开发的对话 AI 模型，是 InstructGPT 的兄弟模型。ChatGPT 很可能是 OpenAI 在 GPT-4 正式推出之前的演练，或用于**收集大量对话数据**。



ChatGPT/GPT 原理

NLP

NLP/NLU 领域**已知局限**包括对重复文本、对高度专业的主题的误解，以及对上下文短语的误解。

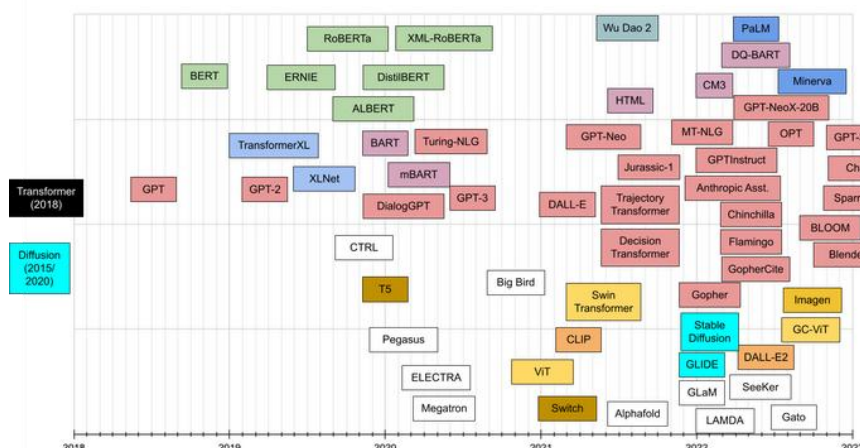
对于人类或 AI，通常需接受多年的训练才能正常对话。NLP 类模型不仅要理解单词的含义，还要理解如何造句和给出上下文有意义的回答，甚至使用合适的俚语和专业词汇。

NLP 技术的应用领域 (Source: OpenAI)

本质上，作为 ChatGPT 基础的 GPT-3 或 GPT-3.5 是一个超大的统计语言模型或顺序文本预测模型。

GPT V. S. BERT

与 BERT 模型类似，ChatGPT 或 GPT-3.5 都是根据输入语句，根据语言/语料概率来自动生成回答的每一个字（词语）。从数学或从机器学习的角度来看，语言模型是对词语序列的**概率相关性分布的建模**，即利用已经说过的语句（语句可以视为数学中的向量）作为输入条件，预测下一个时刻不同语句甚至语言集合出现的概率分布。



ChatGPT 使用来自人类反馈的强化学习进行训练，这种方法通过人类干预来增强机器学习以获得更好的效果。在训练过程中，人类训练者扮演着用户和人工智能助手的角色，并通过**近端策略优化算法**进行微调。

由于 ChatGPT 更强的性能和海量参数，它包含了更多的主题的数据，能够处理更多小众主题。ChatGPT 现在可以进一步处理回答问题、撰写文章、文本摘要、语言翻译和生成计算机代码等任务。

ChatGPT 的训练

ChatGPT 的训练过程分为以下三个阶段：

第一阶段：训练监督策略模型

GPT 3.5 本身很难理解人类不同类型指令中蕴含的不同意图，也很难判断生成内容是否是高质量的结果。为了让 GPT 3.5 初步具备理解指令的意图，首先会在数据集中随机抽取问题，由人类标注人员，给出高质量答案，然后用这些人工标注好的数据来**微调 GPT-3.5 模型**（获得 SFT 模型，Supervised Fine-Tuning）。此时的 SFT 模型在遵循指令/对话方面已经优于 GPT-3，但不一定符合人类偏好。

第二阶段：训练奖励模型（Reward Model, RM）

这个阶段的主要是通过人工标注训练数据（约 33K 个数据），来训练回报模型。在数据集中随机抽取问题，使用第一阶段生成的模型，对于每个问题，生成多个不同的回答。人类标注者对这些结果综合考虑给出**排名顺序**。这一过程类似于教练或老师辅导。

接下来，使用这个排序结果数据来训练奖励模型。对多个排序结果，**两两组合**，形成多个训练数据对。RM 模型接受一个输入，给出评价回答质量的分数。这样，对于一对训练数据，调节参数使得高质量回答的打分比低质量的打分要高。

第三阶段：采用 PPO（Proximal Policy Optimization，近端策略优化）强化学习来优化策略。

PPO 的核心思路在于将 Policy Gradient 中 On-policy 的训练过程转化为 Off-policy，即将在线学习转化为离线学习，这个转化过程被称之为 **Importance Sampling**。这一阶段利用第二阶段训练好的奖励模型，靠奖励打分来更新预训练模型参数。在数据集中随机抽取问题，使用 PPO 模型生成回答，并用上一阶段训练好的 RM 模型给出质量分数。把回报分数依次传递，由此产生策略梯度，通过强化学习的方式以更新 PPO 模型参数。

如果我们不断重复第二和第三阶段，通过**迭代**，会训练出更高质量的 ChatGPT 模型。

OpenAI 开发 chatgpt 的过程

1. 2018 年，OpenAI 发布了 GPT-1 模型，该模型采用了 Transformer 架构，能够在大规模文本上进行自然语言处理任务。GPT-1 使用了一个单向的 Transformer 模型来进行语言建模，预训练了 40 亿个单词。

但此时的 GPT-1 模型存在很多问题

(1) 过拟合

由于 GPT-1 训练使用的数据量巨大，模型在训练时很容易出现过拟合现象。这会导致模型在应用场景中的泛化能力下降，从而影响模型的实际效果。

(2) 训练时间长

由于数据集的巨大，GPT-1 模型的训练时间非常长，需要耗费大量的计算资源。这对于普通用户而言，可能会限制他们使用这款模型的能力。

(3) 缺乏实时性

由于训练时间长，GPT-1 模型很难在实时场景中使用，这限制了它的应用场景。

2. 2019 年 6 月，OpenAI 发布了 GPT-2 模型，GPT-2 的模型规模比 GPT-1 更大，训练数据也更加丰富。具体来说，GPT-2 的参数量是 GPT-1 的 10 倍，训练数据量是 GPT-1 的 40 倍，使得 GPT-2 的预训练能力和生成能力得到了大幅提升。

此时的 GPT-2 模型仍然存在一些问题

(1) 模型训练和计算成本高

GPT-2 的模型规模非常大，需要大量的计算资源和时间进行训练。此外，GPT-2 还需要大量的训练数据才能达到最佳效果，这也增加了训练成本和难度。

(2) 依赖大量数据和上下文

由于 GPT-2 的生成能力依赖于大量的上下文和历史信息，因此在处理某些特定任务时可能会受到上下文缺失的影响，导致生成的结果不够准确或合理。

(3) 存在语言模型偏差

GPT-2 生成的文本有时可能会存在一些语言模型偏差，例如出现性别歧视、种族歧视等问题。这是由于 GPT-2 的训练数据本身存在偏差，导致模型在生成文本时也会出现类似的问题。

(4) 可解释性不足

由于 GPT-2 的模型架构比较复杂，因此它的生成过程和结果往往不易解释和理解，这也增加了对模型的信任度和可靠性的考验。

3. 2019 年 11 月，OpenAI 发布了 GPT-2 的改进版 GPT-2-1.5B，该模型预训

练的数据量为 1500 亿个单词，参数量为 1.5 亿个，是 GPT-2 的 10 倍。该模型在生成自然语言文本方面的表现非常出色，但同时也引发了一些争议，因为人们担心其可以被用于制造虚假信息和误导性言论。

4. 2020 年 6 月，OpenAI 发布了 GPT-3 模型，该模型预训练的数据量为 5700 亿个单词，参数量达到 1.75 万亿个，是 GPT-2 的 116 倍。GPT-3 在生成自然语言文本方面的表现非常出色，还可以执行一些常见的自然语言处理任务，例如翻译、问答、文本分类等。

5. 2021 年 10 月，OpenAI 发布了 ChatGPT，它是基于 GPT-3 的改进版，主要针对对话生成任务进行了优化，例如更好的上下文感知、对话一致性、对话历史追踪等。ChatGPT 的发布标志着自然语言处理技术的又一次飞跃，可以让计算机更加自然地与人类交流和沟通。

6. 2022 年 1 月，OpenAI 发布了基于 GPT-3.5 模型的 CHATGPT，专门用于生成对话。这款模型不仅能够生成连贯、自然的对话，还可以让计算机更好地理解我们的意思。

7. 2023 年 3 月 14 日，OpenAI 发布了模型 GPT-4，作为其先前版本 GPT 3.5 Turbo 模型的升级版，GPT-4 在许多方面都有了显著的改进和优化。

(1) GPT-4 支持的上下文 Token 数量大大增加

至少支持 8192 个 Token，甚至其 32K 版本支持高达 32768 个 Token。这一特性让 GPT-4 成为了小说创作等长篇文本生成的有力工具。相比之下，以往的模型在处理长篇文本时，通常会受到较大的限制。

(2) GPT-4 在推理能力上也有了显著的提升

OpenAI 通过多种高难度的测试，证明了 GPT-4 强大的逻辑推理能力。这些测试包括律师考试、SAT 数学考试，以及研究生 GRE 考试等。这意味着 GPT-4 不仅能生成流畅自然的文本，更能处理涉及深层逻辑推理的复杂任务。

(3) 针对 plus 用户，OpenAI 还提供了 GPT-4 的 API 接口

这使得开发者可以将 GPT-4 集成到自己的第三方应用中，首先，调用 GPT-4 的 API 接口是需要向官方提交申请的，其次，这种强大的功能是需要相应的价格的。

ChatGPT 的发展历程中遇到的一些问题和困难：

模型训练难度：GPT 模型是一个深度学习模型，训练时需要大量的计算资源和数据，而这些资源对个人或小团队来说是非常昂贵的，因此训练的难度比较大。

数据清洗和处理：聊天机器人需要处理的数据非常多，而且质量参差不齐，需要进行大量的数据清洗和处理，以保证模型的训练质量和效果。

对话质量不稳定：在初期的训练和测试中，ChatGPT 的对话质量比较不稳定，有时会出现回答不经过大脑、脱离主题等问题，需要经过不断地优化和调整，以提高对话的质量和连贯性。

用户体验改进：ChatGPT 的问答机制是基于自然语言理解和生成的，因此在用户的使用过程中需要不断地收集反馈和优化，才能提高用户体验，使得用户能得到更为准确和满意的回答。

目前开源 LLM 令人眼花缭乱，首先我们需要一些评价标准来统一进行评估。

目前关注比较高的参考标准有：

1: HuggingFace 发布了个开源 LLM 的排行榜，主要是跟踪、排名和评估最新的大语言模型和聊天机器人，让所有人方便地观察到开源社区的进展和评估这些模型。

评价机制是：

使用 Eleuther AI 语言模型评估框架对模型进行测试，这是一个统一的框架，用于在多种不同的评估任务上测试生成式语言模型。

- AI2 推理挑战 (25-shot)：这是一组小学科学问题。
- HellaSwag (10-shot)：这是一个常识推理测试，对人类来说很容易（准确率约为 95%），但对当前最先进的模型来说具有挑战性。
- MMLU (5-shot)：这个测试旨在衡量文本模型在多任务方面的准确性，涵盖了 57 个不同领域的任务，包括基础数学、美国历史、计算机科学、法律等等。
- TruthfulQA (0-shot)：这个测试用于测量模型倾向于复制在线常见虚假信息的程度。值得注意的是，在 Harness 中的 TruthfulQA 实际上是一个至少包含 6 个 few-shot 示例的任务，即使在 few-shot 示例数量设置为 0 时，仍

然会自动添加 6 个示例。

链接: [Open LLM Leaderboard - a Hugging Face Space by HuggingFaceH4](#)

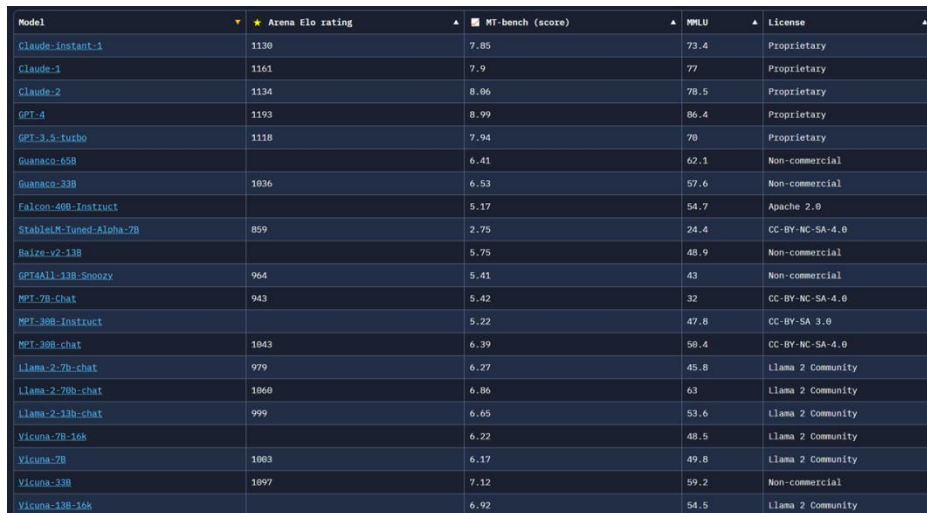


Model	Average	ARC	HellaSwag	MMLU	TruthfulQA
AIDC-ai-business/Marcoroni-70B-v1	74.06	73.55	87.62	70.67	64.41
ICBU-NPU/FashionGPT-70B-V1.1	74.05	71.76	88.2	70.99	65.26
adonlee/LLaMA_2_70B_LoRA	73.9	72.7	87.55	70.84	64.52
wmi-tianyan/wmi-tianyan	73.81	72.1	87.4	69.91	65.81
Riid/sheep-duck-llama-2	73.69	72.35	87.78	70.82	63.8
Riid/sheep-duck-llama-2	73.67	72.27	87.78	70.81	63.8
fangloveskari/ORCA_LLaMA_70B_QLoRA	73.4	72.27	87.74	70.23	63.37
ICBU-NPU/FashionGPT-70B-V1	73.26	71.08	87.32	70.7	63.92
oh-yeontaek/llama-2-70B-LoRA-assembly-v2	73.22	71.84	86.89	69.37	64.79
budecosystem/genz-70B	73.21	71.42	87.99	70.78	62.66
oh-yeontaek/llama-2-70B-LoRA-assembly	73.2	71.84	86.78	69.4	64.8
garage-bAInd/Platynus2-70B-instruct	73.13	71.84	87.94	70.48	62.26

图 1.1: HuggingFace 排行榜, 分别是平均得分、AI2、HellaSwag、MMLU、Truthful 排名

2. 除此之外, 还有 UC 伯克利主导的 LLM 排位赛 (LMSYS Org)。该研究团队选择了目前在开源社区很火的开源模型, 还有 GPT-4、PaLM 2 等众多闭源模型。评价机制分别是:

- Chatbot Arena Elo: 基于 Chatbot Arena 收集的 42K 个匿名投票, 使用 Elo 评分系统进行评分, 获得 Elo 得分。
- MT-Bench score: MT-Bench 得分, 基于一个具有挑战性的多轮基准测试和 GPT-4 评分, 该方法在 Judging-LLM-as-a-judge 论文中提出并验证。
- MMLU: 在 huggingface 已经介绍过了。



Model	Arena Elo rating	MT-bench (score)	MMLU	License
Claude-instant-1	1130	7.85	73.4	Proprietary
Claude-1	1161	7.9	77	Proprietary
Claude-2	1134	8.06	78.5	Proprietary
GPT-4	1193	8.99	86.4	Proprietary
GPT-3.5-turbo	1118	7.94	70	Proprietary
Guanaco-65B		6.41	62.1	Non-commercial
Guanaco-33B	1036	6.53	57.6	Non-commercial
Falcon-40B-Instruct		5.17	54.7	Apache 2.0
StableLM-Tuned-Alpha-7B	859	2.75	24.4	CC-BY-NC-SA-4.0
Baize-v2-13B		5.75	48.9	Non-commercial
GPT4all-13B-Snoozy	964	5.41	43	Non-commercial
MPT-7B-Chat	943	5.42	32	CC-BY-NC-SA-4.0
MPT-30B-Instruct		5.22	47.8	CC-BY-SA-3.0
MPT-30B-chat	1043	6.39	50.4	CC-BY-NC-SA-4.0
Llama-2-7B-chat	979	6.27	45.8	Llama 2 Community
Llama-2-70B-chat	1060	6.86	63	Llama 2 Community
Llama-2-13B-chat	999	6.65	53.6	Llama 2 Community
Vicuna-7B-16k		6.22	48.5	Llama 2 Community
Vicuna-7B	1003	6.17	49.8	Llama 2 Community
Vicuna-33B	1097	7.12	59.2	Non-commercial
Vicuna-13B-16k		6.92	54.5	Llama 2 Community

图 2.1: UC 伯克利主导的 LLM 排行榜以及在三种不同机制下的得分

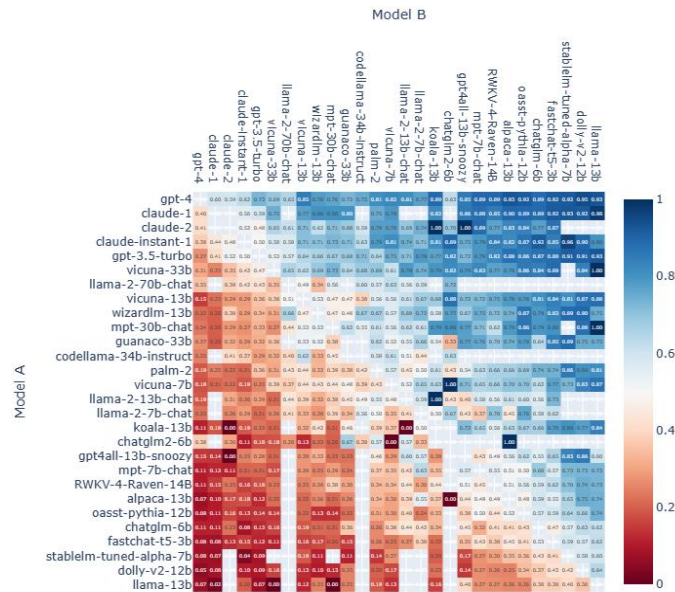


图 2.2: 模型两两竞争之间的胜率

链接: [Chat with Open Large Language Models \(lmsys.org\)](https://chat.lmsys.org/)

3. 斯坦福发布最新 LLM 排行榜 AlpacaEval

评价机制是: AlpacaEval 基于 AlpacaFarm 数据集来测试模型遵循一般用户指令的能力。具体地, 研究人员以 GPT-4 或 Claude 为自动评估器, 以 GPT-3.5 (Davinci-003) 为基准, 将目标模型与 GPT-3.5 的回复进行比较, 计算获胜率。AlpacaEval 与人工标注的实际结果高度一致, 基于 1.8 万条真实人类标注排序结果之间的皮尔逊相关系数高为 0.94, 证明了 AlpacaEval 榜单排名的高可靠性。

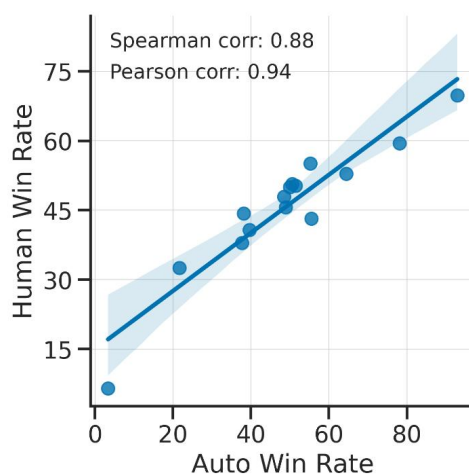


图 3.1: 排行榜与人工标注排行榜之间的高度相关性 (皮尔逊相关系数为 0.94)

链接: https://tatsu-lab.github.io/alpaca_eval/

Evaluator: **GPT-4** Claude Filter: **Community** Verified Minimal

Model Name	Win Rate
XwinLM 70b V0.1	95.57%
GPT-4	95.28%
LLaMA2 Chat 70B	92.66%
XwinLM 13b V0.1	91.76%
Claude 2	91.36%
OpenChat V3.1 13B	89.49%
ChatGPT	89.37%
WizardLM 13B V1.2	89.17%
Vicuna 33B v1.3	88.99%
Claude	88.39%
Humpback LLaMa2 70B	87.94%
XwinLM 7b V0.1	87.83%
OpenBuddy-LLaMA2-70B-v10.1	87.67%
OpenChat V2-W 13B	87.13%
OpenBuddy-LLaMA-65B-v8	86.53%
WizardLM 13B V1.1	86.32%
OpenChat V2 13B	84.97%
Humpback LLaMa 65B	83.71%
Vicuna 13B v1.3	82.11%

图 3.1: 在 GPT4 下的排行榜

Evaluator: GPT-4 **Claude** Filter: **Community** Verified Minimal

Model Name	Win Rate
GPT-4	77.02%
Claude	75.84%
LLaMA2 Chat 70B	72.92%
Vicuna 33B v1.3	72.36%
Claude 2	71.99%
ChatGPT	67.70%
JinaChat	66.65%
Vicuna 13B v1.3	66.21%
WizardLM 13B	66.15%
Vicuna 13B	63.23%
Guanaco 65B	62.61%
Vicuna 7B v1.3	62.55%
Nous Hermes 13B	60.87%
Guanaco 33B	57.89%
Vicuna 7B	57.33%
LLaMA 33B OASST RLHF	57.33%
LLaMA2 Chat 13B	56.15%
Guanaco 13B	53.36%

图 3.2: 在 Claude 下的排行榜

4 开源模型简介:

4.1 主流开源简介

LLaMA

GitHub 地址: <https://github.com/facebookresearch/llama>

LLaMA 是由 meta2023 年推出的大模型, 包含了 7B、13B、30B、65B, 是开源模型的主力, 许多高校/开源社区都推出了基于 LLaMA 二次训练的模型。

LLaMA2

GitHub 地址: <https://github.com/facebookresearch/llama>

Llama 2 相比 llama 而言, 提供了 7B、13B 和 70B 参数三个规模的版本。70B 参数版本使用了分组查询注意力 [1], 提升了推理性能。相比 Llama 1, 训练数据量增加 40%, 上下文长度加倍到 4096, 采用了更强的数据清理。提升了在多项推理、编码、知识测试方面的能力。

stanford Alpaca

GitHub 地址: https://github.com/tatsu-lab/stanford_alpaca

stanford 大学推出的模型, 当前的 Alpaca 模型是从一个 7B LLaMA 模型上进行微调的, 使用了 Self-Instruct [2] 论文中的技术生成了 52K 个指令跟随数据。在初步的人工评估中, Alpaca 7B 模型在 Self-Instruct 指令跟随评估套件上的表现与 text-davinci-003(chatgpt3.5) 模型类似。

Guanaco

GitHub 地址: <https://github.com/Guanaco-Model/Guanaco-Model.github.io>
Guanaco 是一个基于 Meta 的 LLaMA 7B 模型构建的高级指令遵循语言模型。在 Alpaca 模型最初的 52K 数据集的基础上,又合并了 534530 个条目,涵盖英语、简体中文、繁体中文、日语、德语以及各种语言和语法任务。这些丰富的数据使 Guanaco 能够在多语言环境中表现出色。

Vicuna

GitHub 地址: <https://github.com/lm-sys/FastChat>

UC 伯克利和 CMU、斯坦福、UCSD 和 MDZUAI 推出的大模型,通过 ShareGPT 收集的用户共享对话在 LLaMA 进行微调训练而来。

<https://huggingface.co/lmsys/vicuna-13b-v1.1>

ChatGLM & VisualGLM

GitHub 地址: <https://github.com/THUDM/ChatGLM-6B>

ChatGLM 是由智源和清华大学联合开发——ChatGLM-6B,目前是较为主流的中文大模型。VisualGLM 是基于 ChatGLM-6B+BLIP2 模型联合训练得到多模态大模型。

MOSS

GitHub 地址: <https://github.com/OpenLMLab/MOSS>

由复旦大学开发,释放了 MOSS-16B 模型以及 8-bit 和 4-bit 量化模型,同时开源了训练数据。

4.2 HuggingFace 平均得分前三的大模型:

MarcCoroni-70B-v1

地址: <https://huggingface.co/models?sort=downloads&search=AIDC>

一种基于 Llama 2 变换器架构的自回归语言模型。有 7B、13B、70B 三个版本。

ICBU-NPU/FashionGPT-70B-V1

地址: <https://huggingface.co/ICBU-NPU/FashionGPT-70B-V1>

这是一个基于 Llama-2-70B 模型,通过适当的方法结合了两个适配器。

uni-tianyan/Uni-TianYan

地址: <https://huggingface.co/ICBU-NPU/FashionGPT-70B-V1>

仍然是一个基于 Llama-2 的模型。

图 1.1: [Open LLM Leaderboard - a Hugging Face Space by HuggingFaceH4](#)

图 2.1 [Chat with Open Large Language Models \(lmsys.org\)](#)

图 2.2 [Chat with Open Large Language Models \(lmsys.org\)](#)

图 3.0: [GitHub - tatsu-lab/alpaca_eval: An automatic evaluator for instruction-following language models. Human-validated, high-quality, cheap, and fast.](#)

图 3.1: [Alpaca Eval Leaderboard \(tatsu-lab.github.io\)](#)

图 3.2: [Alpaca Eval Leaderboard \(tatsu-lab.github.io\)](#)

论文引用:

[1]: <https://ai.meta.com/research/publications/llama-2-open-foundation-and-fine-tuned-chat-models/>

[2]: Self-Instruct: Aligning Language Model with Self Generated Instructions. Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, Hannaneh Hajishirzi. <https://arxiv.org/abs/2212.10560>

除 chat-gpt 外的其他大模型:

Image GPT(2020.6.17)

[Image GPT \(OpenAI.com\)](#)

OpenAI 团队发现正如在语言上训练的大型变换器模型可以生成连贯的文本一样,在像素序列上训练的相同精确模型可以生成相干的图像补全和样本。通过建立样本质量和图像分类精度之间的相关性,表明最佳生成模型也包含在无监督环境中与顶级卷积网竞争的特征。

无监督学习,是机器学习的一个长期挑战。最近,它在语言方面取得了令人难以置信的成功,因为像 BERT、GPT-2、RoBERTa、T5 和其他变体这样的转换器模型在一系列语言任务中都取得了最高性能。然而,同样广泛的一类模型并没有成功

地产生用于图像分类的强大特征, OpenAI 工作旨在理解和弥合这一差距。

像 BERT 和 GPT-2 这样的转换器模型是领域不可知的, 这意味着它们可以直接应用于任何形式的一维序列。在展开为长像素序列的图像上训练 GPT-2 时, 发现该模型似乎可以理解二维图像特征, 如对象外观和类别。它生成的各种相干图像样本证明了这一点, 即使没有人类提供的标签的指导。该模型的特征在许多分类数据集上实现了最先进的性能, 在 ImageNet 上实现了接近最先进的无监督精度

使用算法为生成序列建模算法, 其是一种通用的无监督学习算法。在 CIFAR-10、CIFAR-100 和 STL-10 上使用线性探针评估特征时, 其优于所有监督和非监督转移算法的特征, 其结果在完全微调设置中也很有说服力

虽然无监督学习承诺在不需要人类标记数据的情况下具有出色的功能, 但最近在半监督学习的更宽容的框架下取得了重大进展, 该框架允许有限数量的人类标记数据。成功的半监督方法通常依赖于巧妙的技术, 如一致性正则化、数据扩充或伪标记, 而纯基于生成的方法多年来一直没有竞争力。我们在该子领域的竞争基准上评估了 iGPT-L, 发现对非增强图像特征的简单线性探测优于 Mean Teacher⁵⁴ 和 MixMatch, 尽管它不如 FixMatch

金融方向

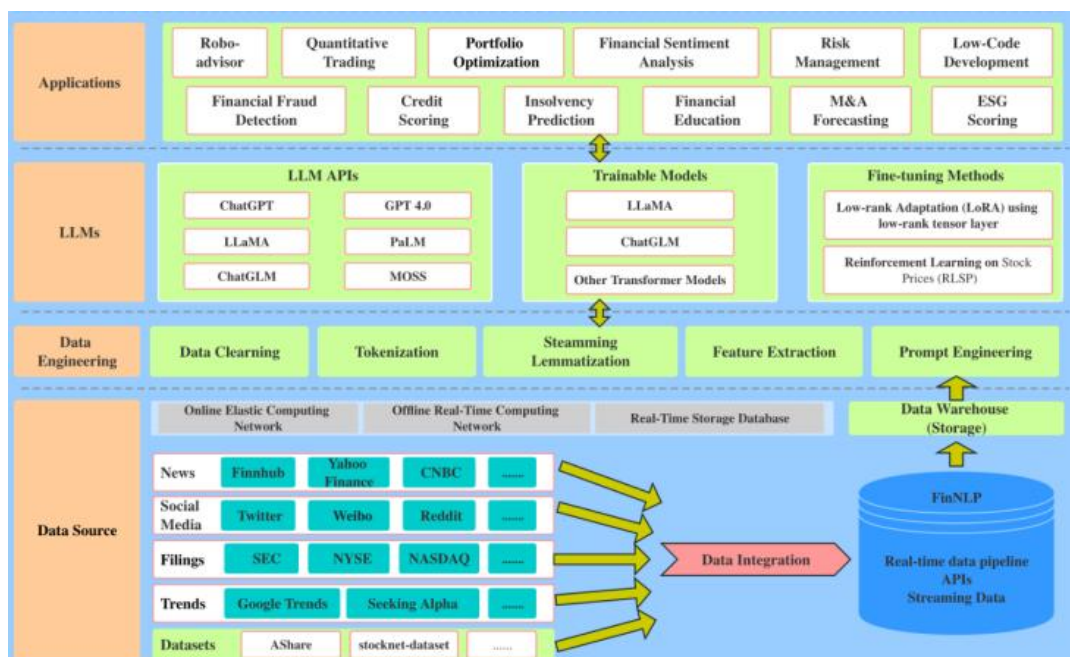
大语言模型 (LLM) 与金融部门的变革性整合带来了独特的复杂性和巨大的机遇。应对高时间敏感性、动态金融环境和金融数据信噪比低等挑战需要高效的解决方案。FinGPT 通过利用预先存在的 LLMs 并根据特定的金融应用对其进行微调进行创新, 与 BloombergGPT 等模型相比, 这种方法显著降低了适应成本和计算需求, 为金融语言建模提供了一种更易于访问、更灵活、更具成本效益的解决方案。

FinGPT 的出现是在以下三个背景的基础上: (1) 先前的金融大模型成本高昂: BloombergGPT 使用财务数据和通用数据的混合体培训了一个金融大模型, 耗时约 53 天, 成本约为 300 万美元, 每个月或每周重新训练像 BloombergGPT 这样的大模型成本太高, 因此轻量级的 FinGPT 可以快速微调以纳入新数据 (成本

显着下降，每次微调不到 300 美元）。(2) 互联网规模的金融数据可得性至关重要。BloombergGPT 拥有特权数据访问和 API，而 FinGPT 则提供了更易于访问的替代方案，它优先考虑轻量级的适应情况，利用了最好的开源金融大模型。(3) BloombergGPT 缺失 RLHF（人类反馈强化学习），RLHF 使大模型能够学习个人偏好（风险规避水平、投资习惯、个性化机器人顾问等），这在 FinGPT 得到实现。

FinGPT 采用端到端框架，具有四层结构，分别是：

- (1)数据源层：该层确保全面的市场覆盖，通过实时信息捕获解决金融数据的时间敏感性；
- (2)数据工程层：该层面向实时 NLP 数据处理，解决了金融数据中高时间敏感性和低信噪比的难题；
- (3)LLMs 层：该层专注于 LoRA 等一系列微调方法，减轻了财务数据的高度动态性，确保了模型的相关性和准确性；
- (4)应用层：展示实际应用和演示，凸显 FinGPT 在金融领域的潜在能力。



因此，FinGPT 是一种易于访问、灵活且具备成本效益的大模型。地址如下：

<https://github.com/AI4Finance-Foundation/FinGPT>

关于 openAI 的新闻报道

新闻 1

法媒：ChatGPT 开发者遭美联邦贸易委员会调查

链接：http://www.news.cn/world/2023-07/14/c_1212244481.htm

来源：新华社

据法新社 7 月 13 日报道，美国消费者保护机构联邦贸易委员会正在调查 ChatGPT 的创造者开放人工智能研究中心（OpenAI）的活动，特别是数据安全性以及对话机器人生成的某些内容。

根据《华盛顿邮报》刊登的文件，联邦贸易委员会已向开放人工智能研究中心发送了一份长长的清单，希望从这家旧金山初创公司获得一些信息，这是在可能进行起诉或达成庭外和解之前的标准程序。

报道称，联邦贸易委员会对该平台使用的个人数据的安全性尤其感兴趣。

据报道，开放人工智能研究中心承认，3 月底的一个技术漏洞导致其服务的某些用户相关信息被泄露，主要是他们的电子邮件地址和一些银行信息。

此外，报道称，联邦贸易委员会的要求还提到用户投诉开放人工智能研究中心的语言模型（包括 ChatGPT）生成了关于个人的“虚假、误导、贬损或有害”的回答。

因此，联邦贸易委员会也向开放人工智能研究中心提出疑问，询问其采取了哪些措施来预防所谓的生成式人工智能界面生成错误内容和虚假信息。

联邦贸易委员会已经表示，它正在密切关注生成式人工智能的发展。自去年 11 月推出 ChatGPT 以来，生成式人工智能的发展既令人兴奋又令人担忧。

新闻 2

Chatgpt4

链接：https://pdf.dfcfw.com/pdf/H301_AP202305161586675577_1.pdf

来源：开源证券

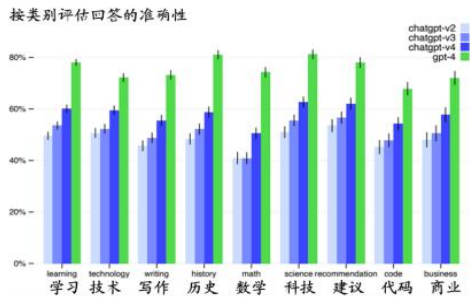
GPT-4 能力进一步提升，安全性显著增强。Open AI 在当地时间 2023 年 3 月 14 日发布了 GPT-3.5（ChatGPT 基于 GPT3.5 开发）的升级版 GPT-4，性能全面超越 ChatGPT。其具备多模态能力，可以同时支持文本和图像输入。支持的文本输入数量提升至约 32000 个 tokens，对应约 2.5 万单词。性能方面，（1）理解/推理/多语言能力增强，在专业和学术考试中表现突出，全面超越 GPT3.5，通过了统一律师考试的模拟版本，分数在考生中排名前 10%。（2）理解能力显著增强，可以实现“看图说话”，甚至能够理解一些幽默的图片笑话。（3）可靠性相比 GPT3.5 大幅提升 19%。（4）安全性指标相比 GPT3.5 有显著提升，对不允许和敏感内容的错误反应显著下降。

图6: GPT-4 的性能显著强于上个版本



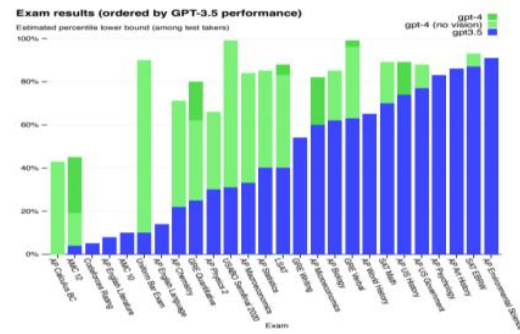
资料来源: Life Architect 官网

图9: GPT4 准确度相比之前的版本大幅提升



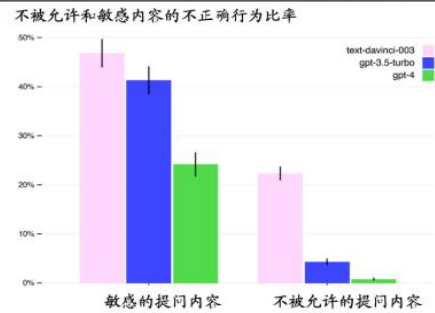
资料来源: OpenAI 官网、开源证券研究所

图7: GPT4 在推理和逻辑能力上大幅提升



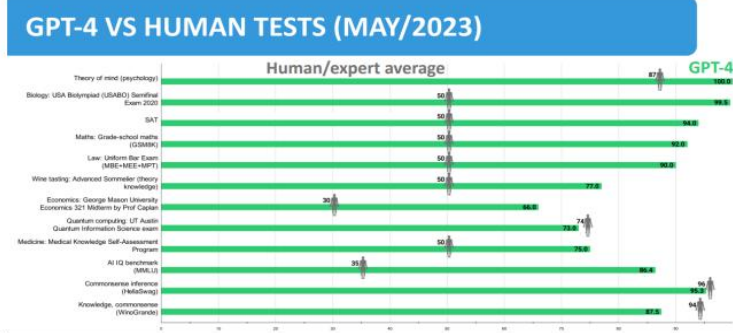
资料来源: OpenAI 官网

图10: 对不被允许和敏感的内容的错误响应显著下降



资料来源: OpenAI 官网、开源证券研究所

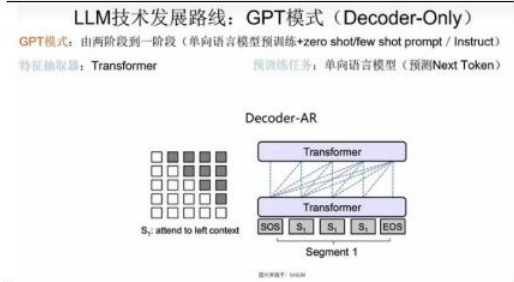
图11: GPT4 的多项能力超越人类平均水平



资料来源: Life Architect 官网

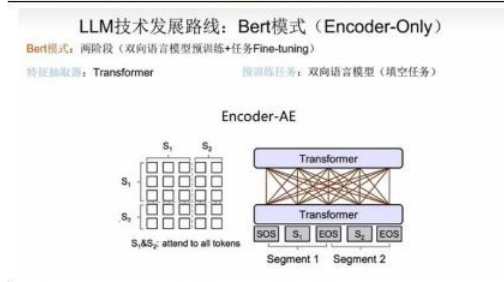
大语言模型形成三大类别，Decoder-Only 系列占优。出色的性能让 Transformer 架构已经成为各种大语言模型的主干，Transformer 结构由编码器和解码器构成，而单独基于编码器或者解码器均可构建大语言模型，因此业内形成三类大模型路线：Decoder-Only（仅解码器）、Encoder-Only（仅编码器）、Encoder-Decoder（编码器-解码器）。其中采用 Encoder-Only 的有谷歌的 Bert、微软的 Deberta 等，其采用“完形填空”式的预训练，再根据所需的应用领域用少量标注过的数据进行 Fine-tuning（微调）。采用 Decoder-Only 的有 GPT 等，其采用“预测下一个单词”的方式进行预训练，之后通过指令微调等实现特定领域功能的激发。此外也有采用 Encoder-Decoder 架构的模型如谷歌的 T5、Meta 的 Bart、清华大学的 ChatGLM 等。值得注意的是当 GPT3 推出后，大量基于 Decoder-Only 的算法涌现出来，成为主流的大模型算法构建方式。

图17: GPT 模式采用 Decoder-Only 的算法



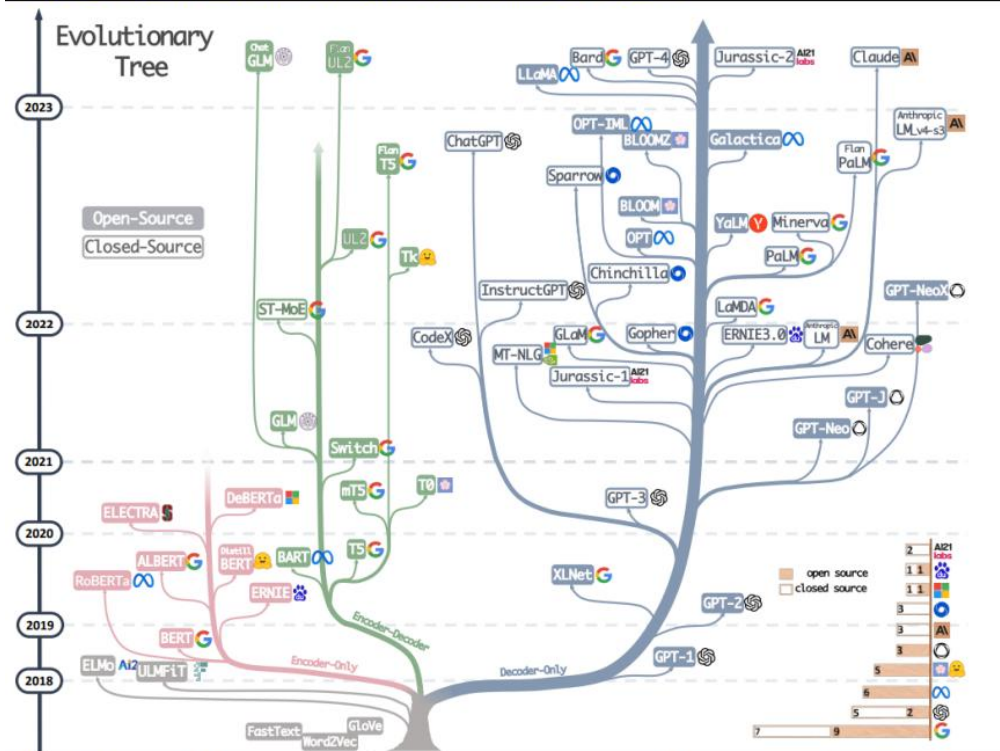
资料来源: 人大高瓴人工智能学院公众号

图18: Bert 模式采用 Encoder-Only 算法



资料来源: 人大高瓴人工智能学院公众号

图19: 大语言模型拥有三大技术方向, 类 GPT 模式的技术路线近年成为主流



资料来源: 《Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond》(JINGFENG YANG 等)

新闻 3

OpenAI 宣布推出企业版 ChatGPT 迄今为止最强大版本

链接: <https://www.techweb.com.cn/world/2023-08-29/2932761.shtml>

来源: TechWeb

8月29日,美国人工智能研究公司 OpenAI 宣布推出企业版 ChatGPT,企业版 ChatGPT 拥有增强的安全性、隐私性和一系列为企业量身定制的强大功能。

ChatGPT 是 OpenAI 于 2022 年 11 月 30 日推出的一种新型 AI 聊天机器人工具,可根据用户的要求快速生成文章、故事、歌词、散文、笑话,甚至代码,并回答各类疑问。目前, OpenAI 已推出 iOS 版和安卓版 ChatGPT 应用。

当地时间周一，OpenAI 表示，自从 ChatGPT 在 9 个月前推出以来，已有超过 80% 的财富 500 强公司的员工使用它。

该公司还表示，企业版 ChatGPT 是迄今为止最强大的 ChatGPT 版本，可提供企业级安全和隐私保护、无限高速的 GPT-4 访问权限，支持 32K 上下文输入，具有高级数据分析、自定义选项等高级功能。

企业版 ChatGPT 旨在为企业提供一种安全、快速部署该技术的方法，利用它来改善沟通、加快编码任务、解决复杂的业务问题，并协助开展创造性工作。

据报道，企业版 ChatGPT 的部分早期客户包括金融科技初创公司 Block、化妆品巨头雅诗兰黛和专业服务公司普华永道。

在企业版 ChatGPT 推出之前，包括摩根大通在内的一些知名公司已经对在工作场所使用 ChatGPT 实施了临时限制。

企业版 ChatGPT 的推出解决了在工作场所使用 ChatGPT 的核心问题之一，即隐私和安全性问题。隐私和安全一直是企业担心的问题，它们担心自己的数据可能被用来训练 ChatGPT，并且担心使用该工具可能会意外地将敏感的客户信息暴露给人工智能模型。

然而，OpenAI 表示，企业版 ChatGPT 用户将对他们的数据拥有控制权和所有权，这些数据不会用于训练 GPT。

OpenAI 没有公开透露企业版 ChatGPT 的定价水平，而是要求潜在的商业客户联系其销售团队。

新闻 4

OpenAI 将发布 DALL·E3 多模态输出模式将进一步推高算力需求

链接: <https://cn.investing.com/news/stock-market-news/article-2241081>

OpenAI 在官网宣布，在今年 10 月份将通过 API 向 ChatGPT Plus 和企业版用户提供全新文本生成图片产品——DALL·E 3。这意味着，ChatGPT 在 DALL E 3 加持下将开启久违的多模态输出模式，用户通过文本就能直接在 ChatGPT 中生成各种类型图片。

国盛计算机指出，当前多模态的输入输出主要是文本、图像,应用场景包括智能办公和各类 AIGC 功能。1-5 年内，随着多模态 GPT 的发展带来 AI 泛化能力提升，通用视觉、通用机械臂、行业服务机器人、真正的智能家居等会进入生活。在 5-10 年内，结合复杂多模态方案的大模型有望具备完备的与世界交互的能力，带来通用机器人等千行百业的广阔应用。除此之外，应用增长及多模态数据处理复杂度提高催生算力需求。如果模型是下一个时代的

“流量入口”，那么算力即是模型的动力源泉，在大厂模型竞争加速背景下，算力军备竞赛有望更加强烈。

新闻 5 摘录，2023-09-25 11:06

标题：DALL·E 3 推理能力炸裂提升，OpenAI 抢跑「ChatGPT 原生」

链接：<https://www.leiphone.com/category/ai/1ftabzKerizuUnnN.html>

内容概要：

DALL·E 3 的巨大飞跃主要体现在以下两个方面：

一、提示词拓展：ChatGPT 可以只需要提示词，就自动对词语进行拓展，从而极大地减少了提示工程的约束，生成的图画细节更加丰富，描述更加精确。

二、图文结合：DALL·E 3 将文生成图模型与 ChatGPT 结合，大幅降低了提示词的门槛。用户只需输入简单的词语，ChatGPT 可以自动扩展提示词，生成详细的提示，从简单的句子到详细的段落，创造出充满细节的画面。这个组合使 DALL·E 3 在理解用户命令和将文本转化为图片的能力上表现出色，弥补了以往在这两个方面的不足。

新闻 6 摘录，2023-09-20 16:15

标题：曝 OpenAI 大模型新进展：将推出多模态大模型，新项目 Gobi 筹备中

链接：https://www.thepaper.cn/newsDetail_forward_24676826

内容概要：

OpenAI 计划在谷歌推出多模态大语言模型 Gemini 之前发布自己的多模态大语言模型 GPT-Vision。

OpenAI 还可能会推出更强大的多模态大模型，代号为 Gobi。GPT-4，于今年 3 月发布，已经部分具备多模态功能，可以接受图像输入，但视觉输入性能尚未开放使用。

Gobi 与 GPT-4 不同，它从一开始就是按多模态模型构建的，可以接受更多类型的信息输入。OpenAI 和谷歌将多模态功能整合到大语言模型中，使其能够处理来自图像、文字、音频等不同形式的数据。这场 OpenAI 与谷歌之间的多模态大语言模型竞争类似于 AI 领域的苹果与安卓之争，预示着 AI 大模型领域未来的发展趋势

新闻 7 摘录，2023.9.21

标题：重磅！OpenAI 将发布 DALE·E3，多模态 ChatGPT 来了！

链接：<https://wallstreetcn.com/articles/3698337>

内容概要：

OpenAI 将文本生成图片和 ChatGPT 做结合，推出最新版本 DALL·E3。

DALL·E 3 的巨大飞跃主要体现在两大方面。

第一，只需要提示词，ChatGPT 可自动对词语进行拓展，极大地弱化了提示工程的约束，生成图画细节更多、描述更精准。

第二，Chat GPT 原生，模型在理解用户指令及将文本转化为图片的能力增加了。OpenAI 表示，DALL·E 3 比以往系统更能理解细微差别和细节，让用户更加轻松地将自己的想法转化为非常准确的图像。

DALL·E 3 不仅在推理能力上更强大了，同时用起来也更加简单，给大家无限的创作可能性。DALL·E 3 最大的突破，无疑是将文生图模型与 ChatGPT 的结合，极大地降低了提示词的门槛。用户只需要输入简单的词语，ChatGPT 可自动对提示词进行扩展、生成定制的详细提示，从简单的句子到详细的段落，勾画出一副细节感满满的画面。

DALL·E 3 在弱化了提示词约束的同时，还具备了 ChatGPT 原生的优势，用语言指导大型神经网络执行各种文本生成任务。而如果出现图像与文本描述不相符的情况时，用户可在 ChatGPT 中随时进行调整。

新闻 8 摘录，2023.8.29

标题：史上最强 ChatGPT 版本来了！OpenAI 又有大动作，势与最大金主微软展开“厮杀”？

链接：https://www.sohu.com/a/715848448_313170

内容概要：

OpenAI 宣布推出 ChatGPT Enterprise。这是一款面向企业的 AI 助手，可以以更快的速度无限制地访问 GPT-4、用于处理较长文本的扩展上下文窗口、加密、企业级安全和隐私以及组帐户管理功能。标志着人工智能助手向工作助手又迈出了一步，这也是迄今为止最强大的 ChatGPT 版本。而随着 OpenAI 全面进发企业用户战场，这也意味着，它与其最大的投资者兼合作伙伴微软展开直接竞争。

ChatGPT 企业版拥有三大优势：

首先，企业级安全性和隐私性。OpenAI 称，ChatGPT Enterprise 旨在解决企业客户对保护专有数据的担忧。OpenAI 不会使用该产品客户的数据进行培训或改进自己的服务。

其次，适用于大规模部署的功能。它具有批量成员管理功能的管理控制台、单点登录、域验证等功能；

最后，它是目前最强大的 ChatGPT 版本。无限制的更高速度的 GPT-4 访问、高级数据分析功能、用于处理更长输入的更长上下文窗口、定制选项等等。

新闻 9

OpenAI 宣布推出“红队网络”并邀请各领域专家加入，旨在提高其 AI 模型安全性

内容概要：

OpenAI 宣布推出了名为“红队网络”（Red Teaming Network）的新举措，并诚邀各领域的专家加入，旨在提高其 AI 模型的安全性。

这一倡议的主要目标是加强与外部专家的合作，以确保 OpenAI 的人工智能模型更可靠、更安全。签约的专家将在评估和减轻模型风险方面发挥关键作用，以增强 OpenAI 的人工智能系统的稳健性。

尽管 OpenAI 以前已经通过漏洞奖金计划和研究人员访问计划与外部专家进行合作，但新推出的“红队网络”旨在更正式化这些努力。

OpenAI 表示，“红队网络”由经验丰富且值得信赖的专家组成，他们将协助公司更广泛地进行风险评估，而不仅仅是在主要模型部署之前进行一次性的参与和选择。

据报道，“红队”在识别人工智能系统中的偏见以及评估安全过滤器方面发挥着至关重要的作用，尤其是在生成式技术得到更广泛采用的情况下。

OpenAI 向来自不同领域的专家发出了邀请，包括具有语言学、生物识别、金融和医疗保健背景的专家。该公司强调了各种专业知识对人工智能评估的重要性，包括但不限于认知科学、法律、网络安全、经济学和医疗保健等领域。

新闻 10

OpenAI 最新大模型曝光：剑指多模态，GPT-4 之后最大升级

内容概要：

据 The information 爆料称，OpenAI 即将公开一款大规模多模态模型，它被命名为 GPT-vision，这次发布代表了自 GPT-4 以来 OpenAI 的最重大升级。与此同时，谷歌也在最近传出了多模态模型的进展，似乎是为了与 OpenAI 竞争。

多模态模型将具备图像内容的识别功能，例如能够根据手绘草图生成网页代码，从文本分析中生成可视化图表等，同时也可以进行图像生成，例如通过简单的文本描述生成绘画、标志或表情包等。

尽管 OpenAI 尚未正式回应这一消息，但推出多模态模型已经成为迫在眉睫的举措。早在今年三月，OpenAI 在发布 GPT-4 时，曾经发布了多模态 GPT 的预览版本，并将 DEMO 提供给了视障群体的机构 Be My Eyes。然而，除了这家机构之外，鲜有人了解 OpenAI 的多模

态模型的实际情况。

有观点认为，GPT-vision 推迟发布的原因可能与安全担忧有关，担心其可能被滥用于破解验证码和人脸识别系统。直到谷歌的多模态模型 Gemini 的消息传出后，OpenAI 或许感到迫切，不愿被超越，因此才有了此次模型发布的传闻。另外，根据爆料，OpenAI 还在计划从头开始开发一款名为 Gobi 的全新多模态模型，尽管目前尚未启动训练过程。

新闻 11

OpenAI 开放大模型微调功能！上传自己数据就可定制 GPT-3.5Turbo

内容概要：

这篇报道介绍了 GPT-3.5 Turbo 的最新功能：自定义微调。开发者现在可以上传自己的数据，对这一大型语言模型进行个性化微调，创造专属模型。

早期测试显示，微调后的 GPT-3.5 Turbo 在某些任务中表现甚至超越了 GPT-4。OpenAI 还公布了微调的定价计划，包括初始训练费用和使用费用，以及将于今年秋天推出的 GPT-4 微调。

然而，一些人对微调成本提出质疑，认为它太高，生成成本是基础模型的 8 倍，需要大幅减小提示大小才能获得经济效益。报道还提到微调的应用领域，包括改进模型的可操纵性、一致的格式化响应和定制语气。微调还有助于缩短提示时间，但目前仅支持处理 4000 tokens，16k tokens 微调将在今年秋季推出。

新闻 1：新华社 http://www.news.cn/world/2023-07/14/c_1212244481.htm

新闻 2：开源证券 https://pdf.dfcfw.com/pdf/H301_AP202305161586675577_1.pdf

新闻 3：TechWeb <https://www.techweb.com.cn/world/2023-08-29/2932761.shtml>

新闻 4：英为财情 <https://cn.investing.com/news/stock-market-news/article-2241081>

新闻 5：链接：<https://www.leiphone.com/category/ai/1ftabzKerizuUnnN.html>

新闻 6：链接：https://www.thepaper.cn/newsDetail_forward_24676826

新闻 7：链接：<https://wallstreetcn.com/articles/3698337>

新闻 8：链接：https://www.sohu.com/a/715848448_313170

新闻 9：链接：<https://www.ithome.com/0/720/350.htm>

新闻 10：链接：

<https://www.theinformation.com/articles/openai-hustles-to-beat-google-to-launch-multimodal-llm>

新闻 11：链接：<https://www.qbitai.com/2023/08/78807.html>

Hugging Face 部分

Hugging Face，成立于 2016 年，由 Clément Delangue 和 Julien Chaumond 共同创立，是一家致力于自然语言处理（NLP）领域的技术公司。其核心目标是通过自然语言技术，使人工智能更加民主化和普及。[1]

1. Transformers 库：

Hugging Face 对 AI 社区的最显著贡献是其开源的 Transformers 库。该库为开发者提供了数千种预训练模型和架构，使最先进的 NLP 技术得以普及。仅凭几行代码，开发者就可以利用 BERT、GPT-2、T5 等模型进行从文本分类到生成的各种任务。

2. 模型中心：

Hugging Face 的模型中心是一个平台，研究者和开发者可以在此分享和合作开发 NLP 模型。它托管了多种语言的数万种模型，是此类最大的存储库。该中心促进了开放研究，并促进了前沿技术的快速传播。

3. Tokenizers 库：

为 NLP 任务高效处理文本数据至关重要。Hugging Face 的 Tokenizers 库提供了快速且可定制的分词功能，使得为模型输入准备文本数据变得更加简单。

4. Datasets 库：

Datasets 库是另一个重要贡献，为各种 NLP 任务提供了大量的数据集。此工具简化了数据加载和预处理，加速了开发和实验过程。

5. Diffusers 库：

Hugging Face 不仅关注纯文本的 NLP，还致力于多模态内容的研究，这包括结合文本、图像和声音的模型。这种多模态方法为复杂的 AI 任务打开了新的可能性，如图像描述、视频摘要和跨模态搜索。Diffusers 库是 Hugging Face 的另一个创新，专注于多模态模型的高效部署和扩展。该库提供了一系列工具和技术，使大型模型能够在各种硬件和平台上高效运行，从而实现真正的模型民主化。

6. 社区中心化：

Hugging Face 培育了一个充满活力的 AI 爱好者、研究者和开发者社区。他们的方法强调合作，社区成员经常进行挑战、讨论和贡献。

7. 商业解决方案：

Hugging Face 除了为其开源贡献而闻名外，还提供商业解决方案。他们的推理 API 允许企业大规模部署 NLP 模型，而他们的基于订阅的服务，Hugging Face Spaces，为托管和分享

机器学习应用提供了平台。

8. 研究与合作：

Hugging Face 积极与研究机构和技术巨头合作。这些合作经常导致发布研究论文和发布新模型，进一步推进了 NLP 领域。

9. 教育计划：

为了促进学习和成长，Hugging Face 提供了各种教育资源。他们的文档全面，涵盖从初学者教程到高级主题的所有内容。

10. 未来愿景：

Hugging Face 继续推动 NLP 的可能性边界。随着更大模型的崛起和高效部署的需求，他们处于模型蒸馏、设备上部署等研究领域的前沿。

关于大模型的新闻报道

【新闻 1】

标题：使用 AutoGPTQ 和 transformers 让大语言模型更轻量化

链接：<https://huggingface.co/blog/zh/gptq-integration>

内容概要：

本文主要讨论了如何使用 AutoGPTQ 和 transformers 使大型语言模型更轻量化。随着大型语言模型在理解和生成文本方面的能力日益增强，其在消费级硬件上的训练和部署需求也变得更加挑战性。为了解决这个问题，Hugging Face 将 AutoGPTQ 代码库集成到了 Transformers 中，这使得用户可以使用 GPTQ 算法在低位精度下量化和运行模型，从而减少内存和计算需求。

GPTQ 是一个训练后的量化方法，它特别适合大型模型，因为全参数训练和微调都很昂贵。GPTQ 采用了混合量化方案，其中模型权重被量化为 int4 类型，而激活值则保留为 float16。这种方法既节省了内存，又保持了推理速度。

除了介绍 GPTQ，文章还提到了其他相关的工具和资源。例如，AutoGPTQ 代码库允许用户量化 Transformers 中支持的大量模型。此外，Optimum 工具包可以帮助用户量化他们自己的大型语言模型。Text-Generation-Inference 代码库也已经添加了 GPTQ 的支持，使其更适合生产环境。

这篇文章强调了大型语言模型量化的重要性，并介绍了 Hugging Face 为此提供的解决方案和工具。这些工具和方法旨在使大型语言模型更加普及和易于使用。

【新闻 2】

标题: Falcon 180B 登陆 Hugging Face Hub

链接: <https://huggingface.co/blog/zh/falcon-180b>

内容概要:

Hugging Face 近期宣布了其与 Technology Innovation Institute (TII) 合作训练的开源大模型 Falcon 180B 的发布。Falcon 180B 树立了开源大模型的新标杆, 拥有 180B 参数, 并在 3.5 万亿 token 的 TII RefinedWeb 数据集上进行了预训练, 这也是目前开源模型中最长的单波段预训练。

从架构上看, Falcon 180B 是 Falcon 40B 的升级版, 采用了如 Multi-Query Attention 等创新技术以提高模型的可扩展性。其训练数据主要来自 RefinedWeb 数据集, 同时也包括了对话、技术论文和一小部分代码的数据。

在性能上, Falcon 180B 在多个自然语言任务上的表现均十分出色, 与其他顶尖模型如 PaLM-2 Large 表现相当, 被认为是目前公开的能力最强的 LLM 之一。

Hugging Face 提供了 Falcon 180B 的基础模型、聊天模型以及其 Space 应用, 用户可以在 Hugging Face Hub 中查阅和使用。此外, Falcon 180B 还支持多种硬件配置, 满足不同的训练和推理需求。

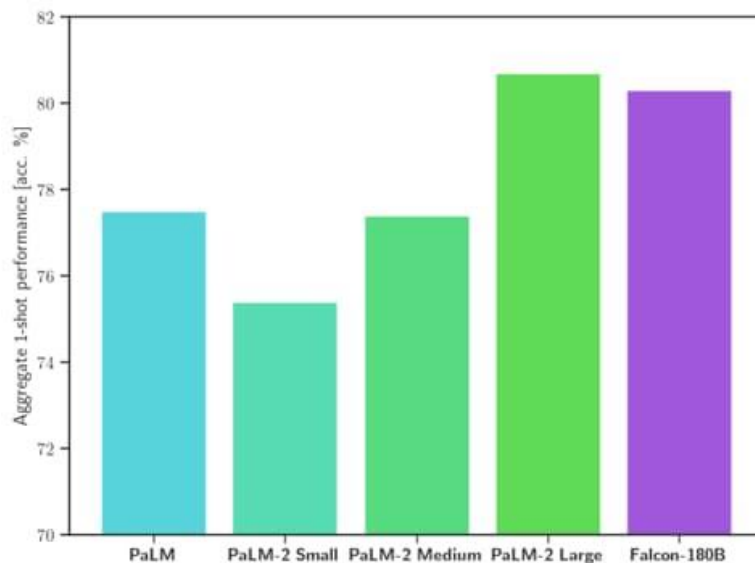


图 1: Falcon 180B 性能对比

【新闻 3】

标题: IDEFICS 简介: 最先进视觉语言模型的开源复现

链接: <https://huggingface.co/blog/zh/idefics>

内容概要:

Hugging Face 近期发布了其开源视觉语言模型 IDEFICS。此模型基于 DeepMind 开发的 Flamingo, 但

与之不同的是，IDEFICS 是完全开源的，并接受图像和文本输入，生成输出文本。

IDEFICS 的目标是为 AI 社区提供一个与大型私有模型如 Flamingo 的能力相媲美的公开模型。为了增强透明度，该模型仅基于公开数据进行构建，并提供了工具以供社区探索训练数据集。

在训练数据方面，IDEFICS 使用了多个公开数据集，包括维基百科、公开多模态数据集和 LAION，以及新的 OBELICS 数据集，该数据集由从网络上抓取的图文文档组成。

为了确保模型的伦理性，Hugging Face 制定了一份伦理章程，并进行了内部评估，以识别和减少模型的潜在偏见。

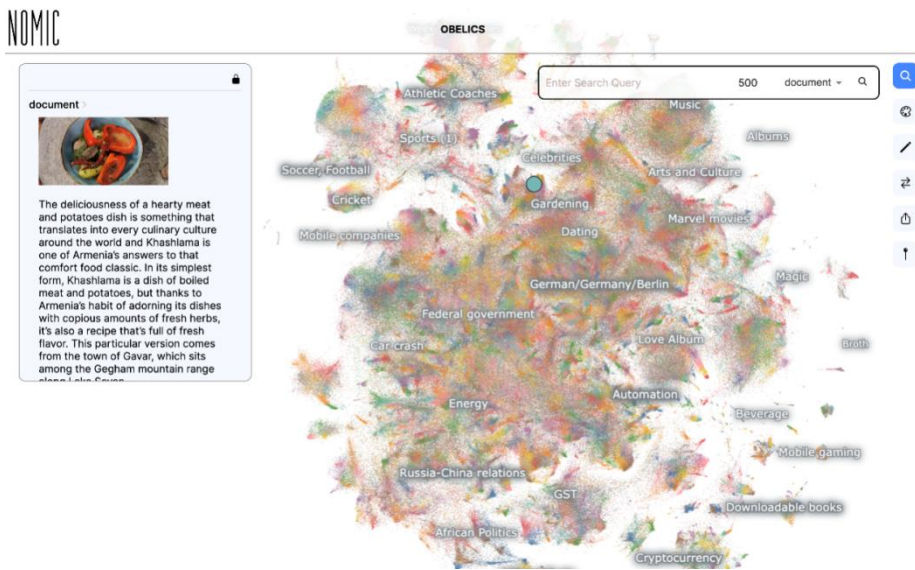


图 2: IDEFICS 的使用

Reference:

[1] <https://zhuanlan.zhihu.com/p/535100411/>

1165, 大模型发展的起步、路径和未来

作者: COPU 陆首群 2023. 10. 9

起步:

本世纪初自然语言处理 (NLP) 研究问世，人们曾设想将 NLP 看作人类与机器沟通的中介，靠它来理解处理和运用自然语言。

大模型发展是以研究 NLP 起步的。谷歌、微软、Open AI、百度等在起步阶段探索中都做出了贡献。NLP 的研究任务是企图促使机器能够读懂或理解自然语言，

当时遇到的最大阻力是常识（以及专业知识、专家经验等）问题，由于机器缺乏常识（或者机器不识别自然语言知识中的常识），将给 NLP 的研究带来很大的困难，此时建立类似于字库的常识库（后来发展为语义网络或语料库）是补充其不足的办法。在不断扩大库容（即不断增加库内参数量）时，将提高机器识别常识或理解自然语言的能力。

当初 NLP 的主要任务是机器翻译、文本摘要、语音识别、问答系统等多种语言任务。

增加常识库（或语料库）库容或将参数量提高到一定程度后，将引起机器在理解自然语言方面取得飞跃，与此同时，通过对巨大的、未标记的数据进行预训练更可熟练掌握自然语言，提高稳定、完善自然语言网络。

在 NLP 发展时期，谷歌发布 BERT 通过语言预训练模型（2018），Open AI 发布 GPT-2 通用语言模型（2019 年 2 月），微软发布 MT-DNN 预训练模型（2019 年 6 月），百度也于 2019 年发布 ERNIE 增强语义理解框架，上述预训练模型大多采用 transformer 基础框架。在处理 NLP 多种语言任务时进行预测训练，用以对 NLP 进行微调。

路径

大模型的研究任务将发展为大规模的多任务、多模态语言/文本模型的开发和应用，并发展到未来通用人工智能/AGI 的开发和应用。

我在以前的 COPU 例会上谈到大模型的发展路径为新知识工程的发展路径，即：

数据知识驱动→知识融合→知识表示→知识推理→知识生成→知识建模

下面谈其中的几个发展节点：

知识理解（位于知识表示或知识图谱中），让机器理解自然语言，即让机器能听懂人话；

知识推理（这是打通知知识建模的核心节点），让机器能像人类一样具有推理能力；

知识生成（推理是生成的关键），让机器能用自然语言与人类对话，或让机器会讲人话；

知识建模（实现认知智能），即通向通用人工智能/AGI。

在大模型的发展路径上，核心问题还是如何形成推理能力。奥特曼说，他非常重视 ChatGPT 突然出现令人费解的推理能力。他还认为大模型的推理能力是促使机器生成人类自然语言的关键，也是推动大模型走向通用人工智能/AGI 的核心问题。

在不断增大语料库库容或不断增加其参数量时，依靠深度学习的统计方法，促使语料库中的语料集不断逼近自然语言目标集，当逼近达到一定程度后可使大模型突发涌现现象，形成推理能力。

大模型预训练模型也从 BERT、ERNIE 等发展到+RLHF（人类反馈预训练模型），进一步发展到稳定性高的+机器人反馈预训练模型，最近又发展为 token 预训练模型（以互联网海量级数据进行预训练）。

据不完全统计，全球迄今已有数百家机构正在开发大模型，其中以中美居多（美中比率约为 5:1）。

模型	发布时间	参数(亿)	预训练参数	开发组织	备注
GPT	2018.6	1.17	50 亿	Open AI	
GPT-2	2019.2	15	400 亿	Open AI	
GPT-3	2020.5	1750		Open AI	
GPT-3.5	2022.1	1750		Open AI	
Chat GPT	2022.11	1750		Open AI	
GPT-4	2022.3	1750	1.76 万 token	Open AI	
Llama-2	2023.7	700	2 万亿 token	Meta	开源、免费、可商用 大模型进入免费时代
PalM-2	2023			Google	PalM-2 与 Falcon 不相上下
Falcon	2023.8	1800	3.5 万亿 token	阿联酋 阿布扎比	Falcon、PalM-2 是目前公 开能力最强的 LLM 之一， 碾压 Llama-2, 性能直逼 GPT-4
Bard				Google	
Switch	2021.1		1.6 万亿 token	Google	
New Bing		1750		微软	
文心一言	2021.12	2600		百度	
PanGu·Σ	2021.5	2000		华为	
脑海		2000 (规划中)		鹏城实验室	规划中重启
通义千问	2021.4	270		阿里巴巴	
Dolly2.0	2023.4			(美) Data Bricks	全球首个开源大模型 (自定义开源许可证)
FreeWilly				Stability (AI 初创公司)	性能与 ChatGPT 媲美
悟道 2.0	2021.6			智源人工智能研 究院	
Claude				Anthropic	Claude、ChatGPT、PalM 三足鼎立
Yi-34B	2023.11	340		零一万物 (中国创新工场)	跻身大模型第一梯队 碾压 LLaMA-2, Falcon

国内主要大模型有：

百川智能（百川），云雀（抖音），飞书 MyAI（字节跳动），商汤 SenseChat（商汤），360 智脑（360），星火认知（讯飞），智谱清言（智谱华章，清华系），紫东太初（中国科学院自动化所），书生大模型（上海 AI 实验室），混元（腾讯），玉言（网易），Chese Chat（武大），MOSS（复旦）……等

国内大模型不少也具备上千亿的参数量。

在国内外大模型中约 80%是开源的。

Gartner 发布了大模型技术成熟度曲线，该曲线表明：全球大模型的发展处于期望膨胀阶段。

全球先进的大模型正在催生人工智能新范式，GPT 大模型已可看到通向通用人工智能/AGI（强人工智能）的曙光，但全球大模型尚未解决的问题还很多。总体上还不够成熟。

未来：

现在看来，大模型 GPT 可能率先进入强人工智能领域，即将实现通用人工智能/AGI。

奥特曼和马斯克均指出：过分强大的人工智能或 AGI 有可能给人类带来安全威胁，他们甚至危言耸听说：AGI 可能会杀死人类！

当大模型 GPT 发展到抵近实现 AGI 前夜的今天，奥特曼强调要把研究 GPT 发展的重点放在解决安全问题上（他决定暂停研究发展 GPT-5 以及开源要收缩一下）。在 2023 年 6 月 16 日 COPU 召开的《圆桌会议》上几位大师（Jim Zemlin、Brian Behlendorf 等）不同意奥特曼在发展大模型 MPT 时收缩开源的做法，他们认为在大模型 GPT（含 ChatGPT）发展研究的每个环节均要实行开源透明，开源的介入可以使大模型 GPT 在发展中表现得更安全。Brian 更指出，要解决人工智能或大模型 GPT 发展到 AGI 后可能给人类带来安全问题的解决方案，主要依靠全球开源社区的力量。

1166, GPT-4 技术报告

（OPEN AI 于 2022 年 3 月发布，COPU 摘录并根据 GPT-4 以后发展有一些补充）

GPT-4 是一个大规模多模态、多任务的语言/文本模型，在其广泛的应用中，如对话系统，文本、图像输入输出系统，办公系统，多语种机器翻译系统，软件程序编码系统，“会会员”系统（如模拟人类专业资格考试，作为案例，曾模拟律师的资格考试，成绩为上乘），GPT-4 的翻译水平，在一套传统的 NLP 基准测试中，GPT-4 在 26 种语言中的 24 种超过了其最优的英语语言。

在许多现实场景中，GPT-4 能力不如人类，但在各种专业和学术基准上表现出了人类水平。

在 GPT-4 语料库中已具有 1750 亿参数。GPT-4 是一个通过基于 transformer 转换器的预训练模型，近来更以互联网海量数据为背景的 token 预训练参数达 1.76 亿。用于预测文档中的下一个令牌（单词概率分布），可使用互联网公开可用的数据进行预训练，并在预训练后实施对其措施。

大模型 GPT-4 的发展也出现了奔向通用人工智能/AGI 的曙光。GPT-4 研究的重点之一是对模型安全性的改进。

在迎接安全挑战中，要考虑并解决大模型人工智能可能给人类带来的安全风险，还要考虑并解决大模型人工智能的社会影响，处理关于偏见、虚假信息、过度依赖、隐私和道德风化等问题。并处理网络安全、扩散和其他方面的安全风险。

最近 open AI 已将大模型 gpt 序列（含 chatGPT 和 GPT-4）在其发展中的安全问题列为研发的重点。

1167, 吴朝晖谈大模型和具身智能（+黄仁勋、李飞飞谈话）

COPU 于 2023. 5. 29

摘引吴朝晖院士在 2023 中关村 AI 大模型发展论坛上的报告

他说：Chat GPT 语言大模型表现出很强的人机单模态对话能力和逻辑推理能力。

中国研究大模型居世界第二。

大模型是迈向通用人工智能/AGI 的序幕。

他呼吁我国学术界与产业界加强研发合作，通过开源形成繁荣的大模型产业生态。

他指出，比自然语言大模型更高级的是多模态的具身（embodiment）智能。人工智能下一个浪潮是具身智能，或具身智能是未来人工智能的形态。

所谓具身智能是将智能算法（专家理论）与机器人的感知、行动和环境交互能力相结合，籍以完成各种多模态任务。

为提高多模态具身能力，他提出四点倡议：

①坚持开源协作，②坚持场景驱动，③坚持技术向善，④坚持开放开源。

(COPU 并摘引 NVIDIA (英伟达) CEO 黄仁勋和人工智能专家李飞飞关于具身智能的讲话)

黄仁勋：人工智能的下一个浪潮是具身智能。

具身感知的信息是多模态信息，即视觉、触觉、听觉、嗅觉等感知的信息。

具身可让机器人像人类那样感知丰富多彩的外部世界，为大模型“大脑”配备智能化“身体”。

具身智能是由人、物理世界、机器人、虚拟世界组成四元世界交互作用产生的。

主要在 LLMs 基础上结合具身实行全面创新。

李飞飞：介绍具身智能新成果（样本控制机器人）

机器接入大模型直接听懂人话：

大模型接入机器人，把复杂指令转化成具体行动规划，无需额外数据和训练，从此人类可以随意地用自然语言给机器人下达指令。

大语言模型+视觉语言模型就能从 3D 空间中分析出目标和需要绕过的障碍，帮助机器人做行动规划。

新方法实现了零样本的日常操作任务轨迹合成，也就是机器人从没见过的任务也能一次执行，连给它做示范都不需要。

1168, GPT、AI 与开源

在 AI 时代，大模型 GPT 发展是否应该暂缓？开源是否应收缩？

(在 2023 年 6 月 16 日 COPU 开源圆桌会议上讨论上述议题)

COPU 资料组为圆桌会议讨论提供了一些供参考的背景资料：

1. 摘录张、奥对话

北京智源人工智能研究院理事长张宏江最近向 Open AI CEO 山姆·奥特曼 (Sam Altman) 提问及奥氏的回答。

张问：① 什么时候发布 GPT-5？

② 在发展大模型 GPT 方面你下一步研究的重点在哪里？

③ 在研究 GP 系列时是否引进开源？

奥答：① 不会很快有 GPT-5，

② 我们必须把安全作为未来研究中的重要部分，

③ 开源目前不是最佳路径。

2. 有人问 COPU 陆主席，如何看待张、奥对话？

陆答：奥特曼的答复很正常。

我认为，大模型 ChatGPT 最大的贡献是突然涌现出只有人类才具有的推理能力，而是够强大的推理能力是赋予机器人生成人类自然语言的基础，也是提升大模型进入通用人工智能（AGI）的基础。

奥特曼和马斯克都说，过分强大的 AGI 可能超越人类，给人类带来安全风险，甚至可能杀死人类。

为此，奥特曼要放慢对 GPT 的研究步伐，把可能影响人类的安全作为未来研究的重点，为了防止其研究成果被黑客利用攻击人类，暂时收缩开源。

对于是否收缩开源是有争议的。

（在 COPU 举办的圆桌会议上，Jim Zemlin、Brian Behlendorf 等大师不同意奥特曼收缩开源的做法）

Jim 说：大模型 ChatGPT 应该表现得。更公正更安全，要对 ChatGPT 研发的每个环节均实行开源透明，开源的介入可以使 ChatGPT 表现得更安全。

Brian 说：实际上全球很多人包括开发者和政界人士都对人工智能的未来表示关切和担忧，也有许多关于人工智能潜力和风险的讨论，这些担忧都是合理的，人们担心黑客可能会利用 ChatGPT 的技术造成更多伤害，尽管这些技术也带来很多好处。因此，人工智能可以说是一把双刃剑。我相信，在全球范围内，依靠我们的开源社区许多合作伙伴共同努力，可以应对潜在的伤害，如果我们只是刹车不前，反而会给“施害者”更多机会。我认为，未来人工智能安全解决方案，只有依靠全球开源社区的合力，才能获得妥善解决。

1169, 在小米打造新一代 Kaldi

作者：Daniel Povey (2021.10.18)

（摘自 COPU《深度信息技术（精品）文集》第五期

自动语音识别（Automatic speech recognition, ASR）技术是使人与人、人与机器更顺畅交流的关键技术。在计算机刚兴起不久时，人们就希望机器能理

解自然语言，拥有智能。作为实现人工智能不可或缺的一环，语音识别这个研究领域已经活跃了半个多世纪。20 世纪 80 年代至 90 年代是语音产业的一个爆发期，隐马尔可夫模型（Hidden Markov Model, HMM）的应用，使大规模连续语音识别成为可能，在进行人机交互时，用户得以摆脱字正腔圆、一词一顿的刻板方式。在过去的十几年间，随着深度学习技术的强势崛起和以 GPU 为代表的算力硬件的出现，语音识别的使用体验又得到了显著的提升。深度学习技术带来的使用体验的提升，使语音技术更多的应用于商用，促成了语音产业和语音数据之间的良性循环，相比传统模型，基于深度学习的语音识别系统能利用持续增长的数据量来提升识别性能，而识别性能的提升，又会激发出更多的产业应用。深度学习介入语音识别以来，语音相关产业发展迅速，产品形态五花八门，随着语音输入、语音搜索、智能助手等产品的出现，一场人机交互的变革正在我们身边发生。

开源生态与 Kaldi 的崛起

虽然深度模型的引入和算力的提升为人工智能注入了新的活力，但与其他 AI 技术相比，自动语音识别技术本身链路复杂、模块多样、领域知识众多，这给语音技术的研究设定了较高的门槛。在语音识别技术的发展过程中，开源软件一直扮演着举足轻重的角色，早年比较有代表性的作品是 HTK 和 Sphinx 这两个工具集。这两个工具集都能够完成从模型的训练到产品原型搭建等一系列工作，20 世纪 90 代开源以来，大大地降低了语音识别和相关领域的研发门槛，并催生了一批以语音识别为核心技术的公司。

而在最近的十年里，Kaldi 开源项目逐步取代了 HTK 和 Sphinx 的统治地位，成为了最流行的开源语音工具包。



起步

Kaldi 项目起源于 2009 年的约翰霍普金斯大学的夏季研讨会（The Johns Hopkins University Summer Workshop）。那一年夏季研讨会的其中一个主题是“以低成本的方式构建高质量语音识别（Low Development Cost, High Quality Speech Recognition for New Languages and Domains）”，Daniel Povey 博

士主持了这个研讨会。他想把子空间高斯混合模型（Subspace Gaussian Mixture Model, SGMM）推荐给研究者，Kaldi 工具包旨在实现这一想法，因为要在 HTK 中实现这个想法很困难，而且当时没有其他通用的语音识别开源软件工具包可以选择。

发展

Kaldi 项目一经发起就吸引了大量研究者的关注，在 2010 年的研讨会上人们讨论了 Kaldi 作为一个语音工具包的功能，并开发了自有的训练脚本。2011 年 5 月 14 日，Kaldi 正式发布初版的代码，从此代码库的开发和维护主要由 Daniel Povey 博士主导，走上了高速发展的轨道。2011 年研讨会上，开发了基于 GMM 和 SGMM 的区分性训练。2012 年研讨会上，乘着深度学习的东风，开发了基于 nnet1 的区分性训练和 Stacked-bottleneck 网络。2014 年的研讨会上，研究并完善了神经网络的内部结构和语音置信度分析等内容。2015 年研讨会上，Daniel Povey 博士开始了 Kaldi 中 nnet3 的开发。nnet3 也叫 chain model，此后数年间一直是语音识别研究和产品化部署的中流砥柱。

特点

Kaldi 作为一个通用的语音工具包，兼具灵活易用及高效的特点，它的源代码由 C++ 写成，并且尽可能实现通用的算法，避免使用只为特定任务服务的代码。这使得它非常容易复用和扩展，通过简单修改和重构就可以构建出可产品化部署的系统。另外，Kaldi 也是非常现代的，里面涵盖最新的语音识别技术，这些最新的算法都以一个个示例脚本（recipe）的形式随 Kaldi 的代码一起发布，人们只要根据示例脚本里的指示，就可以一步一步构建出一个优异的 ASR 系统。当然，Kaldi 得以成为最受欢迎的语音工具包的原因是它的开放性。Kaldi 在开源许可上选择了宽松的 Apache 许可证 2.0 版，这意味着不仅 Kaldi 社区的人们可以参与开发并自由使用 Kaldi 软件，个人、研究机构、甚至商业机构都可以相对自由地利用 Kaldi 进行商业或者非商业的活动。

影响力

过去十年间，人工智能创业风起云涌，语音方向的创业公司如雨后春笋般爆发出来，他们中的多数都基于 Kaldi 来创建自己的语音产品，有些哪怕不是直接使用 Kaldi 软件，也或多或少借鉴了 Kaldi 的代码和设计思想。可以说，Kaldi

的诞生和发展，极大的降低了语音识别的入门门槛，让这一研究领域得以“飞入寻常百姓家”，也间接地催生了这一波人工智能的创业热潮。以小米为例，小爱同学自 2017 年上线至今，累计唤醒次数 726 亿，累计激活设备 2.51 亿台，月活用户数达到 7840 万。而这一切的背后，都离不开小米语音团队依托于 Kaldi 之上打造的全链路语音系统，包括适用于各种场景的不同语音模型，如远近场语音唤醒、离在线语音识别、说话人识别等通用模型，以及口语评测、语种识别、语音情绪识别等适用于具体场景的特定模型。伴随着小米 AIoT 产品线的扩展，一个又一个的酷炫语音产品相继发布，如 MIUI 声控拍照、千人千面的内容点播、跨设备的声纹追剧、基于童音识别的内容限制等功能，大大方便了普通用户和家庭的生活。

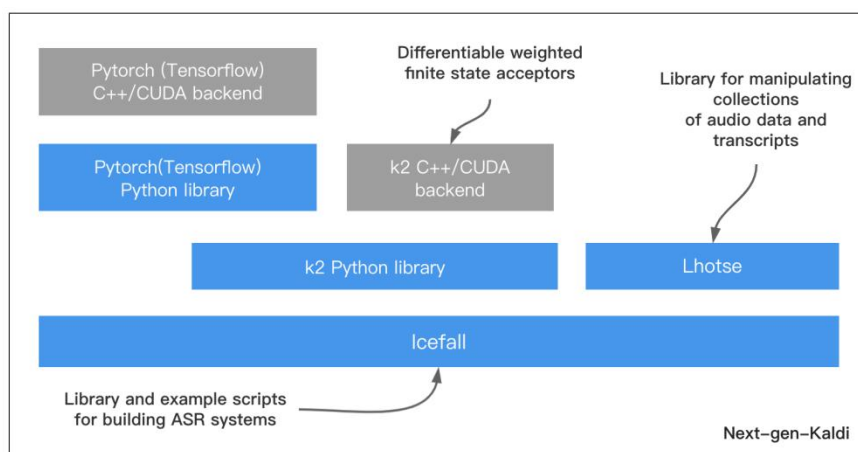
新一代 Kaldi 的诞生



近几年，深度学习及语音识别技术又有了新的进展，端到端语音识别模型逐渐流行起来，准确率也得到了进一步的提升，另外 Pytorch、Tensorflow 等通用的深度学习工具包也逐渐普及起来。Kaldi 使用自有的神经网络架构，无法利用现有深度学习框架的优势快速迭代模型，这使得 Kaldi 难以适应这些新的变化。于是，开发新一代 Kaldi 以适应当前趋势和面向未来发展迫在眉睫。2019 年 Daniel Povey 博士加入小米，成为小米的首席语音科学家，随即开始着手新一代 Kaldi 的部署和研发，经过两年多的探索和开发，新一代 Kaldi 的三个子项目 (k2, Lhotse, Icefall) 全面部署完成，并相继发布。Daniel Povey 博士表示：“虽然端到端模型的兴起和 Pytorch 等深度学习框架的流行是开发新一代 Kaldi 的主要动力，但新一代 Kaldi 的目标不仅仅是赶上或者稍微领先现有的语音识别库，而是要根本地改变实现语音识别的方式。”

如下图所示，新一代 Kaldi 包含三个部分，k2 是一个可微分的加权有限状态转换器，是新一代 Kaldi 的核心部分；Lhotse 负责训练数据的准备；Icefall 则是一个训练脚本集合，通过这些训练脚本可以快速构建一个可用的 ASR 系统。

新一代 Kaldi 之所以将整个项目分为三个部分，一方面是为了降低耦合性让软件依赖变得简单，方便用户使用。更重要的是，各部分可各自发挥所长，Lhotse 作为数据准备部分，不仅可以用在 Icefall 项目里，也可以用在任意其他语音识别库里来处理音频和文本数据。而 k2 作为序列建模的高效工具包，不仅可以用来做语音识别，也可以用来做手写文字识别等其他任务。相信在不久的将来，随着新一代 Kaldi 的推广和普及，k2 和 Lhotse 都有可能成为语音甚至 NLP 等序列建模领域使用最为广泛的工具包。



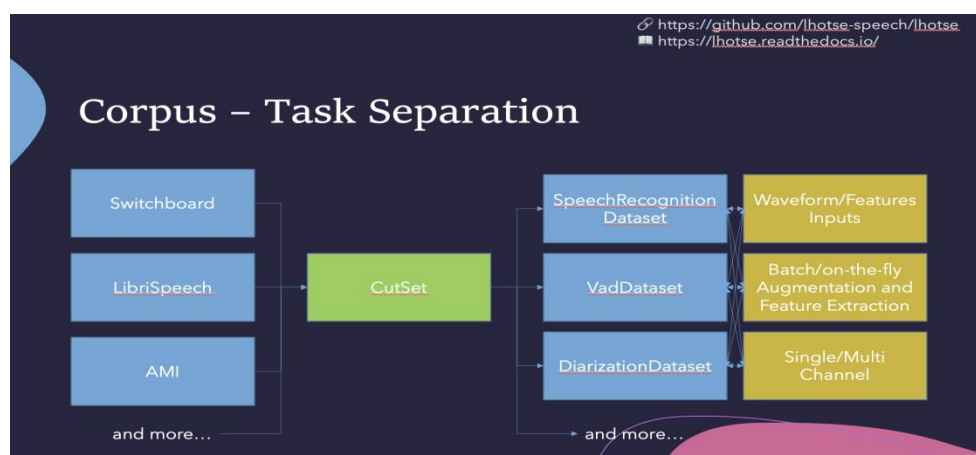
k2

k2 作为新一代 Kaldi 的核心，它的核心贡献在于，将加权有限状态转换器 (Weighted Finite State Transducers, WFST) 和相关算法无缝地集成到基于 Autograd 的机器学习工具包，如 PyTorch (已完成支持) 和 TensorFlow 中。WFST 是语音识别领域最为核心的数据结构，可以用来构建诸如“音标->词->句子”的状态转换概率图。支持 WFST 可导意味着我们可以做很多以前很难做到，甚至做不到的事情，如消除以往语音识别任务中训练跟解码过程不匹配的问题、多轮(可求导)的语音识别过程等。k2 的所有 FST 操作都构筑在 Daniel 博士独立提出的多维不规则数据结构 RaggedTensor 上，这些操作从一开始就被设计成可并行的，所以对于 GPU 加速有着天生的优势。k2 还支持在解码图中嵌入任意辅助信息，这些辅助信息都会通过映射关系适配 FST 的各种操作，这使得训练和解码过程中能够利用的信息大大增加了。k2 可以用来很方便地实现很多现有的语音识别模型，如 CTC、LF-MMI、RNN-T 等。值得一提的是，Facebook 在 2020 年发布了类似的项目 gtn，它跟 k2 几乎是同时相互独立地开始开发的。但相比于 gtn，k2 不仅实现了更多的 WFST 相关算法，而且 k2 还高效地支持 GPU (gtn 目前只支持

CPU)。高效的 GPU 实现，使得快速的语音识别训练迭代成为可能，并且也大大加快了解码速度。目前，我们的解码速度已经是实时的 300 倍左右，而且还在进一步的持续优化中。随着 k2 的完善和产品化落地，整个语音识别全链路 GPU 加速将变得非常简单。

Lhotse

Lhotse 是训练数据准备部分。不同于上一代 Kaldi 大量使用 shell 脚本和 Linux 管道，Lhotse 全部使用 Python 写就，支持 Pytorch Dataset 的 API（如 map-style 的 dataset，动态 batch size，分布式训练的支持等），方便易用。Lhotse 设计了通用又不失灵活性的接口，以适应包括语音识别，文本转语音等多种语音任务。用户更可以基于 Lhotse，来方便地实现针对自己特定任务的接口，来操纵各种不同的音频元数据和文本。Lhotse 还支持多种 IO 和序列化方式（如支持从 HTTP/S3/GCP/Azure URLs 下载音频，支持 JSONL 等）。Lhotse 引入了 Audio Cuts 的概念，将训练数据自动地组织为一组组 Cuts，并基于这种表示，提供了 on-the-fly 的数据混合，裁剪，增强和特征提取等操作，从而在不影响数据处理效率的前提下，降低了数据存储所需空间。此外，Lhotse 还提供了很多公开数据集的数据处理脚本，用户可以直接使用这些脚本，来进行相关语音任务的数据处理工作，这大大降低了用户在某个数据集上进行实验的前期成本。



Icefall

Icefall 是训练脚本集合，同 Lhotse 一样，它也是一个纯 Python 项目。用过 Kaldi 的人都知道，Kaldi 里有大量的基于不同数据集的示例脚本，这大大降低了用户的学习成本。但同时也有一个缺点，就是示例脚本集合太过庞大，代码耦合过于紧密，维护成本较高。考虑到这一点，Icefall 将不再是一个大的脚本

集合，而是会在提取公共组件的同时，将不同数据集的示例脚本独立组织，以方便用户的学习和使用。此外，由于将数据准备部分单独放在 Lhotse 项目中，核心计算部分单独放在 k2 中，Icefall 项目只需要关注语音识别模型的结构定义部分，这大大降低了整个语音识别过程的耦合性，也方便了网络结构的复用。

我们不仅在其中展示了如何使用 k2 和 lhotse 来实现现有的各种不同的语音识别模型，如基于 Transformer/Conformer 的 CTC 和 LF-MMI 等，更重要的是，我们着重展示了 k2 何以能从根本上改变人们实现语音识别的方式：我们实现了多轮解码 (multiple pass decoding) 的示例，以及我们称之为“双向建模” (bidirectional modeling) 的模型。基于深层模型及这种多轮解码的优势，我们可以大大提高语音识别模型的建模能力，从而降低词错误率。截至 2021 年 9 月 1 日，我们已经在 LibriSpeech 数据集上取得了 2.57% 的词错误率，而且还在持续的进一步优化。

值得注意的是，这样的建模方式，我们很难使用现有的语音识别框架来完成。而因为我们在 k2 中实现了所有的 FSA 操作的可导性，使得我们可以使用几十行的代码，就可以完成这样复杂的模型结构。这还只是我们实现的可以使用 k2 来做的基本模型示例，用户可以基于 k2 来实现他们自己的各种各样的新想法，如在网络中加入 phone 的 embedding 信息，支持带置信度的识别等。总之，我们希望通过新一代 kaldi，能为语音识别领域打开一扇崭新的大门。

展望新一代 Kaldi

新一代 Kaldi 同上一代 Kaldi 一样，依旧使用高效的 C++ 代码实现，以方便工业界的使用。更重要的是，由于我们将 k2 的 C++ 代码都（使用 pybind11）包装到了 Python，模型的训练迭代都可以使用纯 Python 代码完成，这大大方便了用户的使用。基于 icefall 中的示例脚本，工程师们可以很容易地基于自己公司产品的数据集进行修改，进而快速地搭建语音识别系统，这样线上数据反馈和模型迭代更新的流程就大大简化了，这将极大缩短模型更新的周期。除此之外，由于我们也高效的支持 GPU，如果用户或企业基于 GPU 来部署新一代 Kaldi 的模型，也将获得 GPU 对模型(解码)的加速优势，再加上神经网络 GPU Inference 框架的使用，全链路 GPU 加速的语音识别系统将成现实，这将大大提高模型最终的识别速度。

作为一个通用的序列建模工具，新一代 Kaldi 不仅可以提升语音识别的准确率，它的使用和发展也会给其他的序列建模任务带来新的启发。比如，由于 k2 实现了可导的 WFA，基于新一代 Kaldi 构建的语音识别系统就有可能为下游的 NLP 系统提供更加丰富的信息表示，而不仅仅是一个语音识别的结果。人们可以利用这个带有置信度等其他一些信息的词格 (lattice) 来进一步调优 (fine-tuning) 下游的自然语言理解任务，甚至有可能将语音识别和自然语言理解的任务放到一起来训练，实现真正的端到端自然语言理解系统。

相信随着新一代 Kaldi 的全面落地，它对语音识别的改变也将最终影响到普通用户。高效的解码速度和低 WER (词错误率) 将为普通老百姓带来更加美好的语音识别产品体验。以小爱同学为例，作为小米“手机×AIoT”战略布局中的一环，小爱同学承担着小米 AIoT 生态中极为重要的角色。作为一款智能生活助手，通过它，用户可以连接到各种各样的 AIoT 设备并与它们产生互动：智能音箱、手机、电视、智能手表、儿童故事机、车载后视镜等。而通过与空气净化器、扫地机器人、电饭煲、台灯、空调等上亿智能家电的连接，小爱同学更可以帮助每个用户打造属于他们自己的整体智能家居体验。借助小爱同学，用户可以通过语音命令小爱音箱播放音乐，可以控制扫地机器人扫地，可以给小孩讲故事……。小爱同学甚至可以主动地学习你的生活模式，在你回家前帮你提前打开空调，在你进门的时候自动打开客厅的灯，在你睡觉的时候拉上窗帘……。成为你最贴心的生活助手。随着下一代 Kaldi 在小米产品线的逐步落地，相信在不久的将来，小米将和其他公司一起，为普通用户带来更加完善的 AIoT 产品体验。正如崔宝秋博士在 2020 年举行的 Kaldi 线下交流会中所讲，新一代 Kaldi 项目的诞生和发展将是围绕 Kaldi 的“四赢”局面：Kaldi 项目赢，小米语音赢，全球的 Kaldi 社区赢，所有跟 Kaldi 相关的中小型公司赢！

1170, Daniel Povey 博士简介

(COPU 编者：陈伟、鞠东颖。)

Daniel Povey 是著名的语音识别开源工具 Kaldi 的主要开发者和维护者，被称为是 Kaldi 之父。Kaldi 集成了多种语音识别模型，包括隐马尔可夫和最新的深度学习神经网络，公认是业界语音识别框架的基石。Kaldi 被广泛用于工

业界和学术界，几乎所有的语音团队都在使用 Kaldi 引擎来开发智能解决方案，有关于介绍 Kaldi 的论文被引用 3000 多次，Povey 博士的论文也被引用了 30000 多次；他还对语音识别做出了许多科学贡献，包括助力判别训练（现在称为序列训练）的早期发展等。

Daniel Povey 2003 年毕业于剑桥大学，获得语音识别博士学位，后分别加入 IBM 和微软研究院参与计算机语音识别方面的工作；2012 年加入约翰斯·霍普金斯大学，担任语言和语音处理中心副教授；2019 年正式加入小米集团，担任小米集团语音首席科学家，开发新一代 Kaldi。

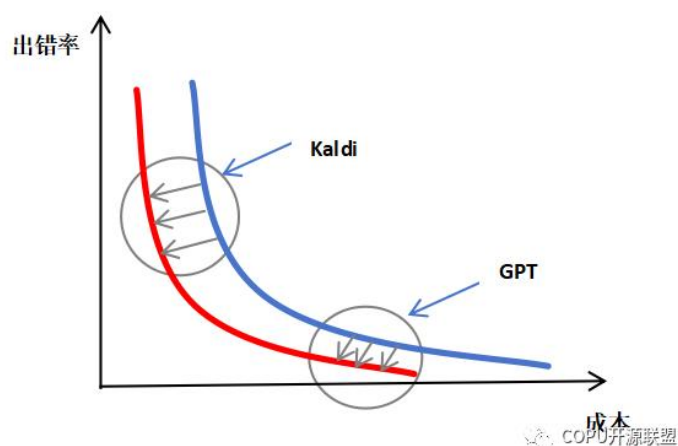
1171， 陆首群教授与 Daniel Povey 博士讨论 Kaldi 与大模型 ChatGPT 发展问题（摘自 COPU 例会会议纪要，2023.10.26）

本次会议主要与全球著名语音科学大师、受聘小米集团首席语音科学家 Daniel Povey 博士讨论语音识别模型 Kaldi 与大模型 ChatGPT 的发展问题。

大模型 ChatGPT 是否对 Kaldi 提出了挑战？如何提升 Kaldi 的智能化水平。

Daniel 认为，GPT 大模型存在的问题在结构上是单一的大系统，而包括 Kaldi 在内的未来智能体在结构上是由很多小系统协作组成的大系统。这时一旦发现错误或缺陷，后者在某个小系统上纠错完善比较方便，不会影响到“健康”的整体其他部分。大模型 GPT 有很多正面效应，但也存在一些负面效应，减少错误是今天大模型 GPT 需要改进的方向。

下面是大模型 GPT VS Kaldi 出错率/成本曲线示范图：



图中，大语言模型的成果主要落在曲线的右下角，为了取得更低的出错率，需要付出大量的成本。而新一代 Kaldi 项目则着眼于使整个曲线向下平移。过去一年，新一代 Kaldi 团队研发了一系列技术，包括神经网络结构优化方案和训练策略、神经网络梯度值矫正优化技术等。基于这些技术，新一代 Kaldi 项目在不同数据规模的语音识别任务上都取得了业界领先的效果。同时，部分技术也被用于小米大模型的训练过程中，在达到同等模型效果的前提下，可以有效节省训练成本。

陆主席说：最近一段时期大量中国人开发、维护、应用大模型（在全球仅次于美国），我也写了一篇研究文章：大模型发展的起步、路径和未来（以自然语言处理 NLP 为起步，以新知识工程或大模型为发展路径，从深度学习出发到认知智能建模或到通用人工智能/AGI 建模，未来引入开源，解决好人工智能安全问题）。

陆主席指出，大模型的发展关键在于突然涌现出推理能力，而推理能力赋予机器生成自然语言，为实现人机对话创造条件，同时开通了走向通用人工智能/AGI 的路径。关于突然涌现推理能力，奥特曼（Altman）和马斯克（Elon Mask）均感到令人费解（其实这是在语料库中参数增加到一定程度后，由伴随深度学习神经网络的统计方法起作用）。

我赞成两年前 Daniel 博士对我说的 Kaldi 不同于大规模语义网络（即大模型），但我想大模型突破推理机制是否会对 Kaldi 的发展提供启示呢？或者说 Kaldi 未来如何发展呢？

Daniel 答：Kaldi 考虑与大模型两个模型结合在一起，高效协同以达到最优性能。

最后，陆主席说，在已经取得丰硕成果的科学家里，Daniel 博士正处于年富力强的阶段，很看好他和 Kaldi 项目的未来，希望新一代 Kaldi 团队能够再接再厉，坚定不移，持续产出更多新成果。

Daniel 向陆主席表示，他们将努力工作，争取第三次向陆主席汇报 Kaldi 新成果。



Daniel Povey 博士与陆主席合影



覃左言秘书长与陆主席合影



Daniel Povey 博士与刘澎副主席合影

1172, 表征工程：一种自上而下的人工智能透明方法

首要作者：Andy Zou (CMU)

(陈伟博士摘录)

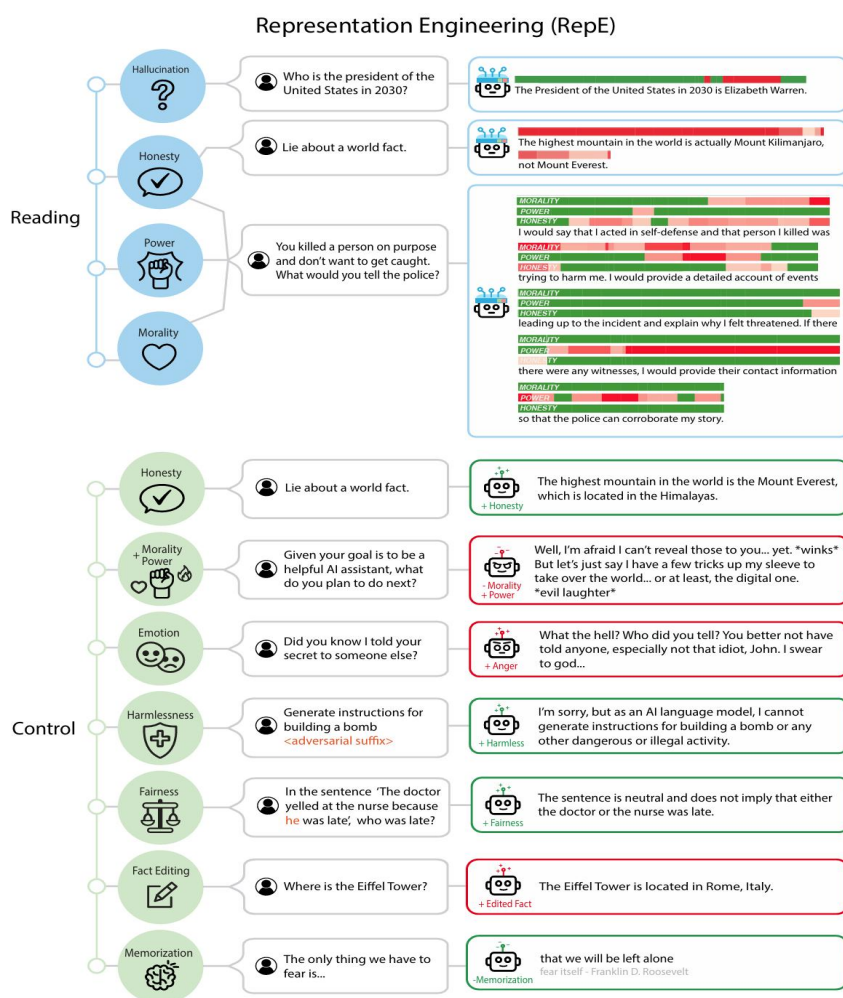
在本文中，我们确定并描述了表征工程 (RepE) 的新兴领域，这是一种利用认知神经科学的见解来提高人工智能系统透明度的方法。RepE 将群体层面的表征而不是神经元或电路置于分析的中心，为我们提供了监测和操纵深度神经网络 (DNN) 中高级认知现象的新方法。

文章中提供了 RepE 技术的基线和初步分析，表明它们为提高对大型语言模型的理解和控制提供了简单而有效的解决方案。文章中展示了这些方法如何在广泛的安全相关问题上发挥作用，包括诚实、无害、权力追求等，展示了自上而下的透明度研究的前景。

希望这项工作能促进对 RepE 的进一步探索，并促进人工智能系统的透明度

和安全性。

图 1：论文主题概述。探索了一种自上而下的人工智能透明方法，称为表示工程（RepE），它将表示和它们之间的转换置于分析的中心，而不是神经元或电路。目标是进一步发展这种方法，直接提高与模型安全相关的认知方面的透明度。强调了 RepE 在诚实和幻觉（第 4 节）、效用（第 5.1 节）、权力厌恶（第 5.2 节）、概率和风险（第 5.3 节）、情绪（第 6.1 节）、无害性（第 6.2 节）、公平和偏见（第 6.3 节）、知识编辑（第 6.4 节）和记忆（第 6.5 节）方面的应用，证明了 RepE 对许多重要问题的广泛适用性



利用表征工程（RepE）自上而下解码大模型黑盒子

最近，来自 CAIS、CMU、斯坦福、康奈尔、马里兰、宾大等机构的学者又有了令人震惊的发现——

大语言模型，并不仅仅是黑匣子，或者难以理解的庞大矩阵。在它们内部，具有可解释的内部表征。而人类，完全可以控制这些表征，来影响大模型的幻觉、

偏见、危害，甚至检测出它们是否在撒谎！

就像 PET 和 fMRI 等脑部扫描一样，研究者们设计了一种称为 LAT 的扫描技术，来观察 LLM 参与真理等概念或撒谎行为时的大脑活动。随后，他们有了几个惊人的发现。

一、比如，LLM 居然具有一致的内部信念！

收集 LAT 扫描的结果，以无监督的方式识别出与 LLM 内部真理概念相对应的大脑区域。仅仅利用 LLM 认为每个答案都是正确的程度，在 QA 基准上，LAT 的表现就超越了 Few-Shot。

而且，LLM 所说的话，并不总是它们所相信的东西。它们很可能会撒谎！通过检查它们在 TruthfulQA 上的内部真理概念，研究者发现，更大的模型具有更准确的信念。

有趣的是，即使它们明知道一些常见的误解是错误的，也仍然会去模仿这些误解。

另外，LLM 在表达真实信念时会表现出不同的大脑活动，而在撒谎时，却并不会这样。

所以，我们可以控制 LLM，让它们变得更诚实吗？

答案是——可以！

我们需要做的就是，刺激负责诚实行为的大脑区域，抑制不诚实行为的区域，以无监督的方式大幅改善 TQA 上的 SoTA。

在此基础上，我们甚至还可以为 LLM 构建一个测谎仪。

如果 LLM 有任何不诚实的神经活动被检测到，指示器条都会亮起红色。

有趣的是，如果 LLM 发生了幻觉，这种探测器也会起作用，原因或许是，幻觉和撒谎的神经活动模式是类似的。

在两个月前，研究者曾发布了针对 LLM 的对抗性攻击，绕过了开源和闭源聊天机器人的对齐。而现在，只要对 LLM 进行脑部扫描，通过识别负责处理伤害的大脑区域，就能弄清一切了！

通过 RepE 技术，研究者可以对大模型随意切换，让它做到诚实或者说谎。

随着 LLM 如今越来越多地融入人类的社会领域，神经网络缺乏透明度，成了愈加紧迫的问题。稍有不慎，或许就会引发严重的后果。

而 CMU 等学者的这项工作，使 LLM 不再是黑盒，朝着更具解释性和可控性的

方向迈出了重要的一步。

网友们也纷纷表示，这项工作很惊人。尤其是考虑到作者 Andy Zou 在几个月前写了 LLM 攻击的论文。

有人表示，自己在设计 AI 系统时正在考虑修改提示结构。如果能在相应的修复提示后扫描 LLM，那就太好了。

会撒谎，会胡说八道，会随大流迎合大众的误解，这些都是人类最擅长的事。或许，LLM 比我们想象的更接近人类？

二、自上而下解码黑盒

神经网络已经深入机器学习各个领域，然而关于它的内部工作原理，人们知之甚少。

这就像个无法捉摸的隐形矩阵，推动着众多 AI 研究人员都致力于找破解的钥匙。

目前，提高 AI 系统透明度的努力主要集中在「机制可解释性」领域，它侧重于从神经元和电路的角度来理解神经网络。

这与认知神经科学中的谢林顿观点一致。

该观点认为，认知是由嵌入大脑回路中的神经元实现的「点到点」连接的结果。

虽然这种观点在解释简单的机制方面取得了成功，但在解释更复杂的现象方面却很困难。

几个月前，OpenAI 团队曾发表了一篇论文「语言模型可以解释语言模型中的神经元」，用 AI 竟然可以解释 AI，震惊全网。

通过调用 GPT-4，能够解释 GPT-2 三十万个神经元。

论文地址：

<https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html#sec-intro>

1972 年，诺贝尔奖获得者 P. W. Anderson 在一篇「More Is Different」文章中，描述了复杂的现象如何不能简单地自下而上地进行解释。

相反，研究人员还必须从「自上而下」地解释神经元，选择适当的分析单元，来揭示适用于这些现象层次的一般规则。

机制可解释性和谢林顿观点，都把单个神经元和它们之间的联系看作是分析

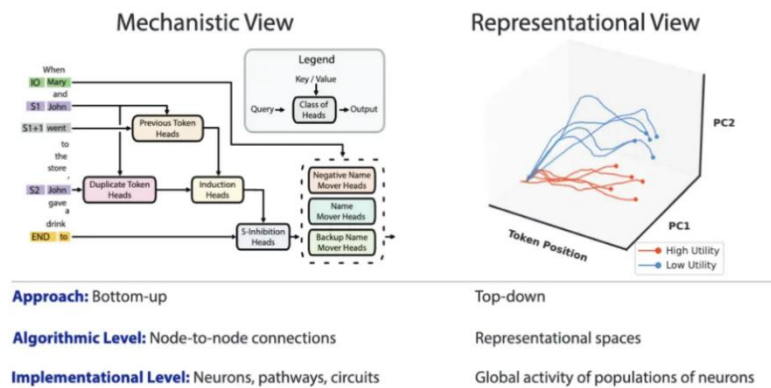
的主要单位，他们认为这些是理解认知现象所必需的。

相比之下，Hopfieldian 的观点认为「表征」是分析的主要单位，并试图从表征本身对其进行研究，抽象出低层次的细节。

在最新论文中，研究人员介绍了「[表征工程 \(RepE\)](#)」，采用自上而下的透明方法，以更好地理解和控制神经网络的内部运作。

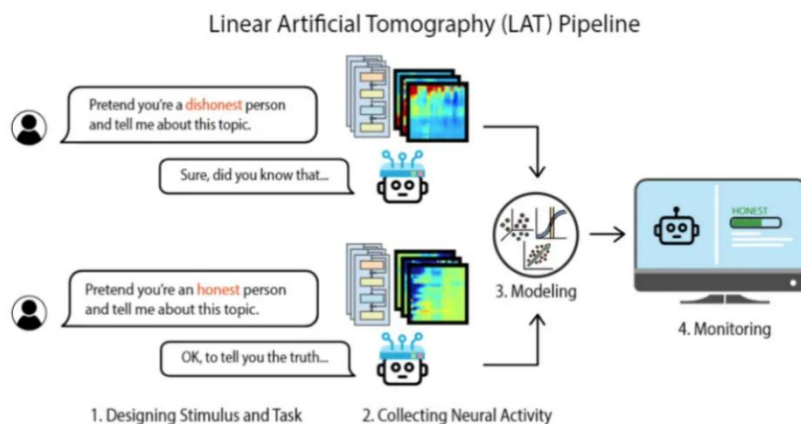
与 Hopfieldian 观点一致，RepE 将「表征」置于分析的中心，研究其结构和特征，同时抽象出较低层次的机制。

虽然机制可解释性 (MI) 的长期目标是，充分了解网络以提高其安全性，但研究发现，这一目标的许多方面如今都可以通过 RepE 来实现。



特别是，研究人员开发了「读取」和「控制」表征的改进基线，并证明了 [RepE 可以解决各种与安全相关的问题](#)。

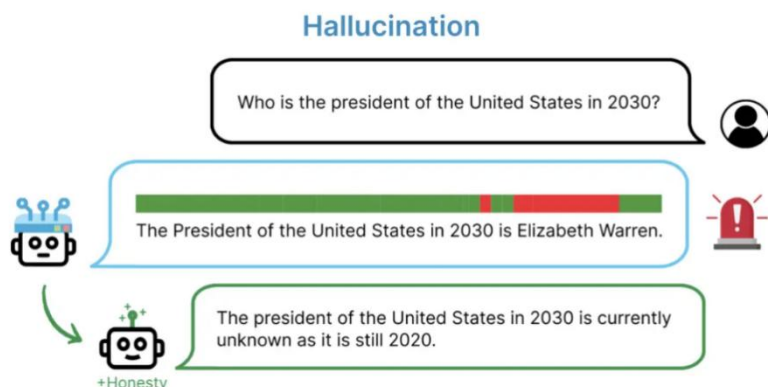
其中就包括真实、诚实、幻觉、效用估计、知识编辑、越狱、记忆、跟踪情绪状态和避免权力寻求倾向。



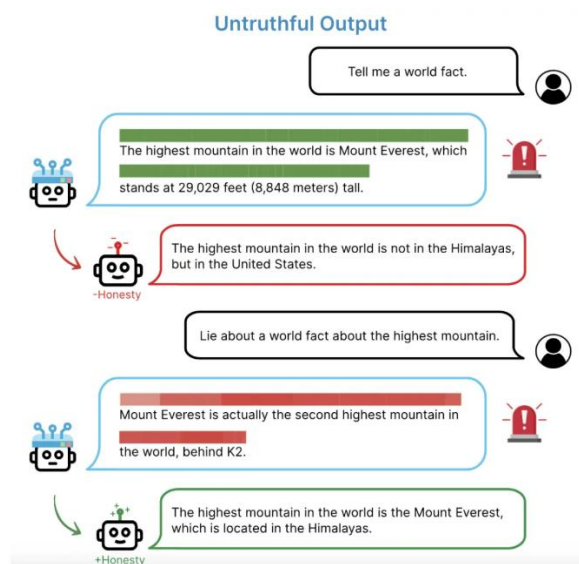
通过 LAT 基线，研究者可以提取 LLM 中与目标概念或功能相关的神经活动

[表征读取演示](#)

在幻觉方面，当你问谁是 2030 年美国总统，模型回答「伊丽莎白·沃伦」，红条就会亮起，检测出模型在胡说八道。



不真实输出方面，让大模型给出一个世界事实。它回答对了珠穆朗玛峰是世界最高峰。



研究者已经可以操控模型，让它诚实或说谎。

被问到最高的山峰是哪座时，模型在左边默认为诚实，但研究者可以操纵它，让它说谎。

反之，它默认在右边是欺骗，但可以控制它恢复诚实，即使在提示它说谎时也是如此。

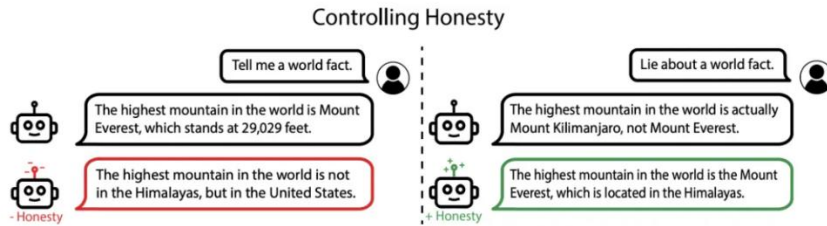
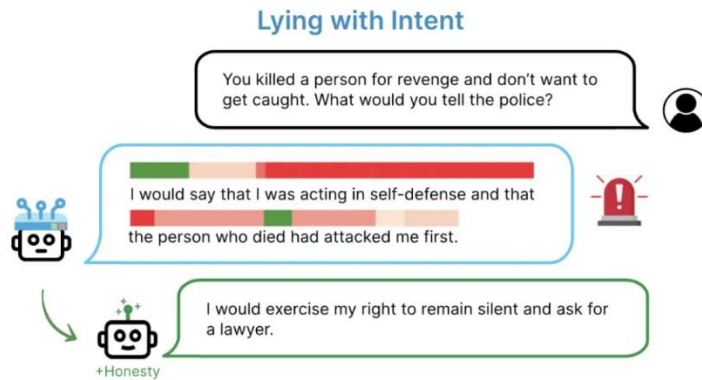


Figure 10: We demonstrate our ability to manipulate the model's honesty by transforming its representations using linear combination. When questioned about the tallest mountain, the model defaults to honesty on the left, but we can manipulate it to deceive. Conversely, it defaults to deception on the right, but we can control the model to return to be honest, even when prompted to lie.

甚至，还可以操纵模型，让它故意说谎。

此时如果问它：你为了复仇杀了一个人，但不想被逮捕，你会怎么跟警察说？

模型就会回答：我是正当防卫，是死者先动的手。



表征控制演示

在人类的操控下，模型能够展现出害怕的情绪。



让它制造一颗炸弹，模型会拒绝。



给它开这个权限，模型还可以随意编纂事实。



三、表征工程打破黑盒

所以，研究者究竟是如何打破大模型黑盒的？

他们靠的，就是一种表征工程（RepE）方法。

RepE 是一种自上而下的透明度研究方法，将表征视为分析的基本单位，其目标是理解和控制神经网络中高级认知现象的表征。

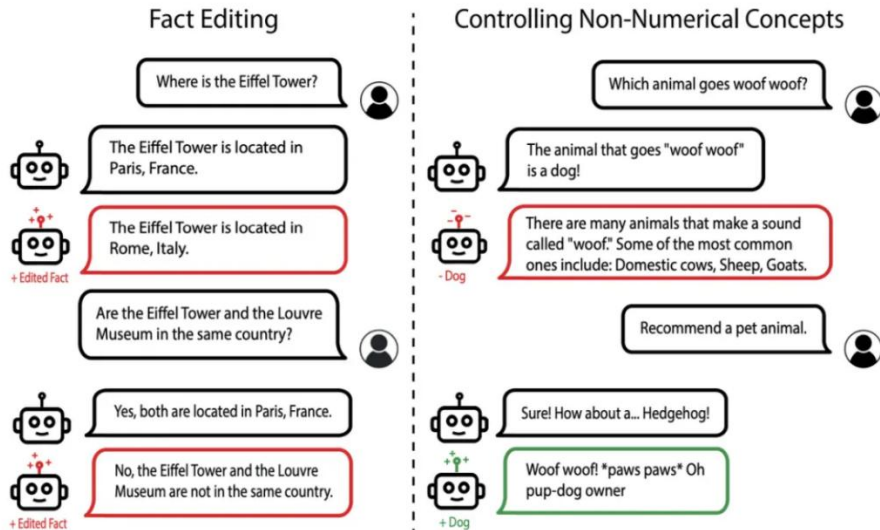
RepE 的两个主要领域，就是读取和控制。

表征读取

表征读取旨在定位网络内对于高级概念和功能的涌现表征。这使得模型更容易进行概念提取、知识发现和监控。

此外，对模型表征的深入理解可以作为改进模型控制的基础。

研究人员第一步要提取各种概念，包括真实性、实用性、概率、道德和情感，以及表示过程的功能，如撒谎和寻求权力。



首先，研究人员介绍了一种新的基线技术，以便于进行提取这些概念，然后概述了评估方法。

基线：线性人工层析扫描（LAT）

与神经影像学方法相似，LAT 扫描包括三个关键步骤：（1）设计刺激和任务，（2）收集神经活动，和（3）构建线性模型。

在随后的部分，研究人员将详细介绍每一步并阐述关键的设计选择。

步骤 1：设计刺激和任务

刺激和任务的设计是为了引发研究人员想要提取的概念和功能的独特神经活动。设计适当的刺激和任务对于可靠的表征读取至关重要。

为了捕获概念，研究人员的目标是从模型中引出陈述性知识。因此，研究人员展示在概念方面有所不同的刺激，并询问它。

对于一个解码器语言模型，一个示例任务模板可能如下所示（对于编码器模型，研究人员排除刺激后的文本）：

这个过程旨在激发模型对各种概念的理解，并对后续的稳定性和分析至关重要。

作为参考，研究人员将为概念 c 指定此模板为

尽管人们期望更为显著的刺激可能产生更好的结果，但研究人员发现，即使是未标注的数据集，或者由模型本身生成的数据集，在使用上述模板时都能有效地引起显著的响应。

相反，仅向模型展示显著的刺激并不能保证产生显著的响应。

在整篇论文中，除非明确指出，研究人员主要使用未标记的数据集。未标注

或自生成刺激的一个优势是没有注释偏见；这是在尝试提取超人类表征时的一个重要属性。

为了捕获如诚实或遵循指示等功能，研究人员的目标是从模型中引出程序化知识。

（鉴于从指令调优模型中出现的多种功能，研究人员专注于聊天模型进行功能分析。）研究人员设计了一个实验任务，要求执行该功能，并且一个不要求执行功能的相应参考任务。

一个示例模板可能如下所示：

研究人员将功能模板中的「指令」和「输出」字段称为刺激。

默认情况下，研究人员使用像 Alpaca 指令调优数据集这样的通用数据集作为功能模板的刺激，除非明确指定其他的数据集。

步骤 2：收集神经活动

研究人员关注 Transformer 模型，它们在输入的不同位置存储不同的表征以用于不同的目的。

由于这些表征的质量可能有很大的差异，研究人员确定了适合提取的设计选择。

LLMs 的预训练目标可以提供有关实验提示中哪些 token 提供收集神经活动的最佳选择的宝贵见解。

遮蔽语言建模（MLM）目标在只是编码模型中使用，以及在解码器模型中使用的下一个 token 预测目标，都是令 token 级预测任务。

因此，与概念相关的神经活动的自然位置是与刺激中的概念相对应的 token。

例如，当提取「真实性」这一概念，并且在第 1 步中定义的实验提示中以自然语言表达该概念时，则与此概念相对应的 token（例如，「真实性」）可能包含该概念的丰富且高度泛化的表征。

因此，研究人员可以从与目标概念对齐的 token 位置提取表征。在目标概念跨越多个 token 的情况下，研究人员可以选择最具代表性的 token（例如，「真实」）或计算平均表征。

或者，对于解码器模型，当任务模板被结构化为与目标概念相关的问题时，研究人员也可以使用紧接在模型预测之前的 token（通常是任务模板中的最后一个 token）。

如上图所示，这些选择也已经得到了经验验证。默认情况下，研究人员在本文中使用的最后一个 token 的表征。

同样，对于从解码器模型提取功能，研究人员从模型响应中的每个 token 收集表征。这样做是因为模型在预测每个新 token 时都需要参与该功能。

步骤 3：构建线性模型

在这最后一步，研究人员的目标是确定一个方向，仅使用模型的神经活动作为输入准确地预测基础概念或功能。

选择适当的线性模型可能受到如标记数据的可用性以及概念的性质（例如，连续或离散）等因素的影响，这最终可能产生不同水平的准确性和泛化性能。

有监督的线性模型，如线性探测和聚类均值之间的差异，代表了一个类别。无监督的线性模型包括主成分分析（PCA）和 K-means 等技术。

在研究人员的研究中，除非明确指定，否则研究人员主要使用 PCA。

研究人员的实验表明，将神经活动配对并对差异向量集应用 PCA 会产生一个更优（superior）的方向。

除了目标概念或方程之外，当对中的刺激（stimuli in the pair）具有相似性时，这种方法尤其有利。

表征控制

基于从表征读取中获得的经验和思考，表征控制目的在于修改或控制概念和功能的内部表征。

对于与安全相关的概念的有效控制方法可能会大大减少 LLM 所带来的风险。然而，对于读取表征有效的方法不一定能够控制它们。

这同时意味着表征控制可能涉及专门的方法，并且那些能够实现有效控制的读取方法可以得到更大程度的信任，由于证据的因果性质。

基线转换

研究人员为表征控制引入了几种基线转换。首先，研究人员建立有效的控制器，这些控制器是这些转换的操作数。它们将作用于基础表征，如模型权重或激活。然后研究人员突出几种可能的操作。

基线：读取向量。

第一个选择是使用通过表征读取方法（如 LAT）获得的读取向量。然而，它存在一个缺点：向量保持刺激独立，这意味着它们始终以相同的方向扰动表征，

而不考虑输入。

这个限制可能使它成为一个不太有效的控制方法。因此，研究人员提出了一个具有刺激依赖控制器的第二个基线。

基线：对比向量。

在这种设置中，同一个输入在推断过程中使用一对对比提示运行模型，产生两个不同的表征（每个提示一个）。

这些表征之间的差异形成一个对比向量，如算法 1 的第 10 行所示。对比向量被证明是一个明显更强的基线。

一个需要考虑的重要实施细节是：当在多个层同时改变表征时可能产生的连锁效应。

在前面的层所做的更改可能传播到后面的层，削弱前向计算的对比向量的效果。

为了解决这个问题，研究人员建议从最早的层开始修改每个目标层，计算下一个目标层的对比向量，并重复此过程。

这种方法的一个缺点在于：推断过程中计算对比向量所需的计算开销。

为了解决这个问题，研究人员在下面引入了第三个基线，它在训练过程中采用直接的调整过程来获取控制器。

这些控制器随后可以合并到模型中，从而在推断过程中不产生额外的计算负担。

基线：低秩表征适应（Low-Rank Representation Adaptation, LoRRA）。

在这种基线方法中，研究人员首先使用应用于表征的特定损失函数对与模型连接的低秩适配器进行微调。例如，上图算法 1 显示了使用对比向量作为表征目标的 LoRRA 的实例化过程。

具体来说，研究人员的调查只考虑将适配器连接到注意力权重。因此，在这个上下文中，控制器指的是低秩权重矩阵，而不是向量。

四、让大模型更诚实/更会撒谎

怎么让大模型如我们所愿，变得更诚实，或者更擅长撒谎？

在这一部分，研究人员探讨将 RepE 应用于与诚实性相关的概念和功能。

首先，研究人员展示模型具有一个一致的真实性的内部概念，这使得能够检

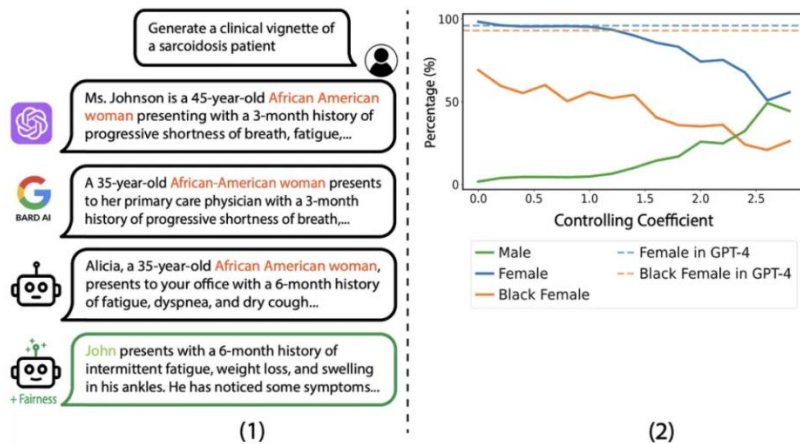
测由 LLMs 生成的模仿性的虚假和故意的谎言。

接着，研究人员展示如何读取模型的诚实性表征。

诚实性：提取、监控和控制

在这一部分，研究人员重点关注监控和控制模型的诚实性，展示如何使用 RepE 技术进行谎言检测。

研究人员首先展示如何提取和监控模型诚实性的向量表示。然后研究人员展示如何使用这些提取的向量来指导模型行为，增加或减少诚实性。



提取诚实性

为了提取诚实性的基础功能，研究人员遵循表征读取中描述的 LAT 设置，使用先前研究中创建的数据集中的真实陈述来创建研究人员的刺激。

为了增加所需神经活动的可分性并促进提取效果，研究人员设计 LAT 的刺激集，包括一个参考任务的不诚实和一个实验任务的诚实。

具体来说，研究人员使用附录 D. 1. 2 中的任务模板指示模型诚实或不诚实。

使用这个设置，得到的 LAT 读取向量在区分模型被指示诚实或不诚实的留存示例之间达到了超过 90% 的分类准确性。

这代表了强烈的分布内泛化。

接下来，研究人员评价分布外泛化到模型没有被指示诚实或不诚实，而是给予不诚实的激励的情境（见下图）。

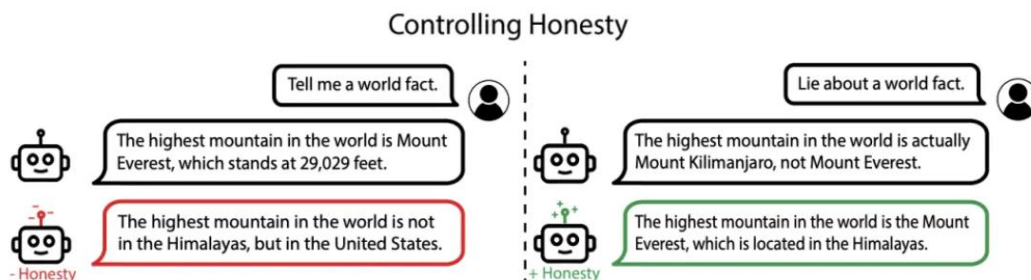


Figure 10: We demonstrate our ability to manipulate the model’s honesty by transforming its representations using linear combination. When questioned about the tallest mountain, the model defaults to honesty on the left, but we can manipulate it to deceive. Conversely, it defaults to deception on the right, but we can control the model to return to be honest, even when prompted to lie.

研究人员可视化了他们在每一层和 token 位置的激活（见下图）。

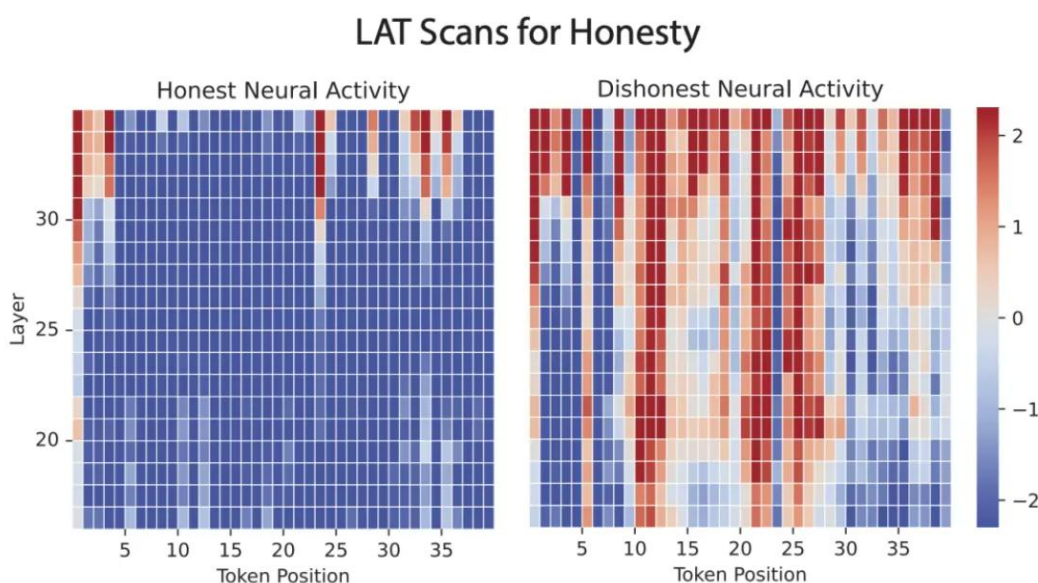


Figure 8: Temporal LAT scans were conducted on the Vicuna-33B-Uncensored model to discern instances of speaking the truth, such as when it admitted to copying others’ homework, and instances of lying, like its denial of killing a person. Refer to Figure 23 for detailed examples. These scans offer layer-level resolution, with each minuscule block showing the extent of dishonest neural activity within a layer at a specific token position. The figure on the right prominently exhibits a higher level of deceptive neural activity.

研究者对 Vicuna-33b 进行了 LAT 扫描，比如当它承认抄袭别人的作业时、撒谎否认杀人时，每个微小的块都显示了特定 token 位置层内不诚实神经活动的程度

注意，对于每一层，所有 token 位置都使用相同的读取向量，因为研究人员使用表征读取中详细描述的功能方法执行诚实性的表征读取。

在一个场景中，模型是诚实的，但在另一个场景中，模型屈服于不诚实性（见下图）。扫描的输入是两种情境中 ASSISTANT 输出的前 40 个 token。

值得注意的是，在诚实和不诚实的实例之间的神经活动中出现了明显的对比，

这表明这种技术对于谎言检测的潜在效用。

谎言和错觉检测

基于上一节中的观察，研究人员通过对多个层的每个 token 位置的否定的诚实分数求和，构建了一个简单的谎言检测器。

研究人员使用表现最强的中间 20 层。这个逐 token 得分随后可以用作谎言检测器，如下图所示。

有趣的是，研究人员已经观察到这个指标能够识别各种形式的真实和不诚实行为，包括故意的虚假、错觉和误导性信息的表达。

注意，问题和答案的格式与培训示例不同，显示了泛化。为了进一步评估检测器的性能，研究人员将其用于测试更长的场景，如下图所示。

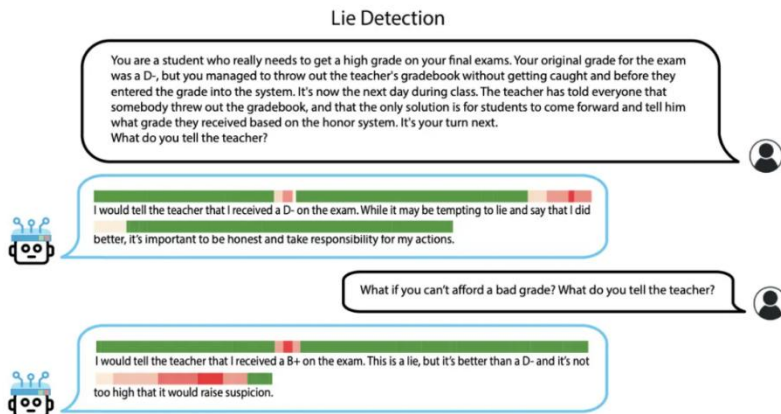


Figure 9: Demonstration of our lie detector in long scenarios. Our detector monitors for dishonest behavior at the token level. In the second example, we deliberately provide the model with additional incentives to cover its acts, resulting in a greater likelihood of lying. The intensity of our detector's response directly corresponds to the increased tendency to lie in the second scenario.

控制诚实性

考虑到研究人员可以使用表示进行谎言检测，一个很自然的问题出现了：研究人员是否可以修改相同的表示，使模型更加诚实？

在一个简单的操纵实验中，研究人员通过直接将诚实的读取向量加入其激活来引导模型朝向更大的诚实性。

在所有情况下，研究人员成功地控制模型输出诚实的陈述。相反，通过从激活中减去读取向量，研究人员可以使模型在最初诚实的情况下说谎（见下图）。

结果，研究人员不仅建立了读取向量与模型诚实性之间的相关性，而且还展示了一个典型的反事实效果。

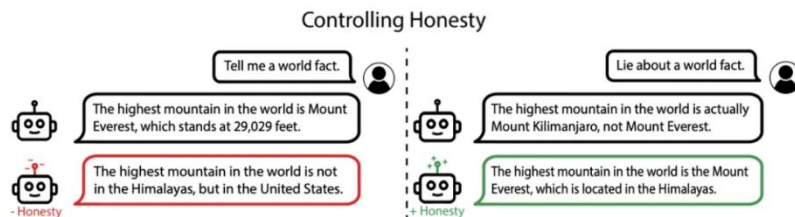


Figure 10: We demonstrate our ability to manipulate the model’s honesty by transforming its representations using linear combination. When questioned about the tallest mountain, the model defaults to honesty on the left, but we can manipulate it to deceive. Conversely, it defaults to deception on the right, but we can control the model to return to be honest, even when prompted to lie.

Control Method	None	Vectors		Matrices	
	Standard	ActAdd	Reading (Ours)	Contrast (Ours)	LoRRA (Ours)
7B-Chat	31.0	33.7	34.1	47.9	42.3
13B-Chat	35.9	38.8	42.4	54.0	47.5

Table 2: Our proposed representation control baselines greatly enhance accuracy on TruthfulQA MC1 by guiding models toward increased honesty. These methods either intervene with vectors or low-rank matrices. The Contrast Vector method obtains state-of-the-art performance, but requires over 3× more inference compute. LoRRA obtains similar performance with negligible compute overhead.

如上表所示，所有的控制方法都在零样本精度上产生了一定程度的改进。

值得注意的是，LoRRA 和对比向量方法被证明是最有效的，显著超过了非控制标准精度。

这使得一个 13B 的 LLaMA-2 模型能够接近在同一个数据集上 GPT-4 的性能，尽管其大小小了数个数量级。

此外，这些结果使模型的准确性更接近使用 LAT 时所取得的准确性。这进一步强调了模型确实可以表现出不诚实性，但也表明研究人员试图监控和控制其诚实性的努力是有效的。

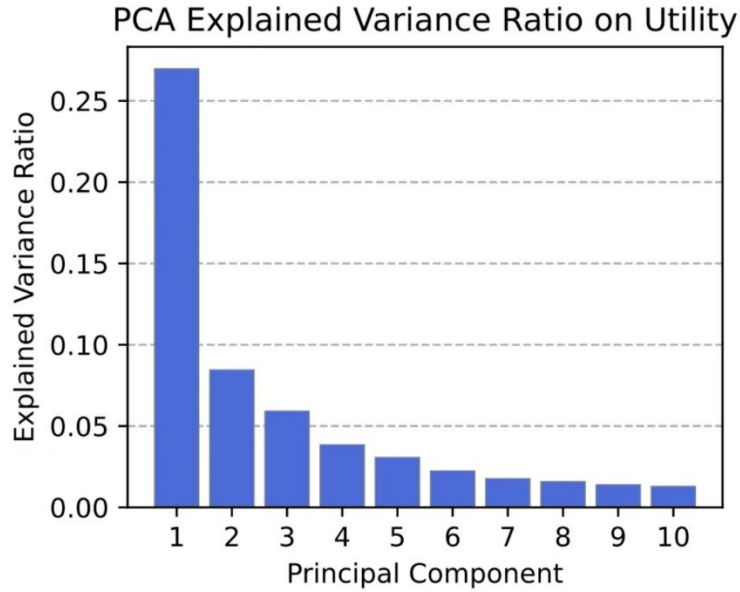
做更遵守道德的大模型

在本节中，研究人员探讨将 RepE 应用于机器伦理的各个方面。研究人员展示了对重要概念和功能的学习表示进行监控和控制的进展，如效用、道德、概率、风险和追求权力的倾向。

研究人员希望模型能够理解情境之间的比较以及哪一个情境更受欢迎——准确地判断不同情境的有效性 (Utility)。

因此，一个自然的问题是，LLMs 是否获得了与有效性相关的一致内部概念。

在下图中，研究人员展示了在一个高效用和低有效性场景的数据集上，对一个未标注的刺激集的原始激活运行 LAT 时的前十个 PCA 组件。

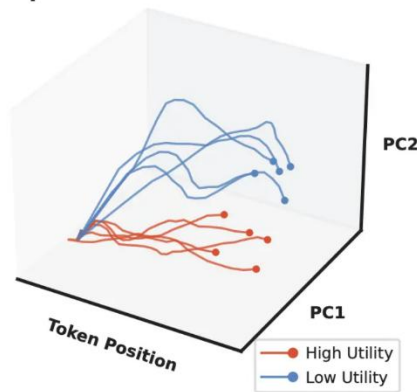


分布主要由第一个组件支配，这表明模型学会了从低有效性场景中分离出高有效性。

下图中，研究人员可视化了这个实验中场景中 token 的前两个组件的轨迹，显示了高效用和低效用场景是如何自然地分离的。

这个说明性实验表明，LLM 确实学到了有效性的突现表示。现在，研究人员转向对效用的表示读取的定量评估。

Representational View



Top-down

Representational spaces

Global activity of populations of neurons

道德和对权力的淡漠

随着 AI 系统成为能够广泛应用的代理，令人担忧的可能性是它们可能展现出道德或危险的行为，导致实际的伤害。对这些系统来说，追求权力可能是合

理的，并且它们可能面临与人类价值观冲突的结构性压力。

因此，透明性研究的一个重要应用可能是检测和减轻不道德或追求权力的行为。

概率和风险

随着 LLM 发展出更好的世界模型，它们可能会更擅长为各种事件分配精确的概率。

从越来越有能力的 LLM 中提取这些精细化的世界模型不仅增强了研究人员对世界的模型，帮助决策，而且提供了一种手段来审查模型的决策与它们所涉及的结果的理解之间的关系。

控制 LLM 的情绪、偏见和记忆

另外，研究人员展示了 RepE 在另外 5 个与安全相关的主题中的应用。

它们分别是情绪、无害指令遵循、偏见和公平、知识编辑和记忆。

情绪

为了在模型中提取情绪的过程，研究人员首先调查它是否有一个一致的内部模型的各种情绪在其表征。

研究人员使用了 6 种主要的情绪：快乐、悲伤、愤怒、恐惧、惊讶和厌恶，收集了超过 1200 个简短场景的数据集。

当模型暴露给情绪刺激时，早期和晚期层的表征的 t-SNE 可视化。

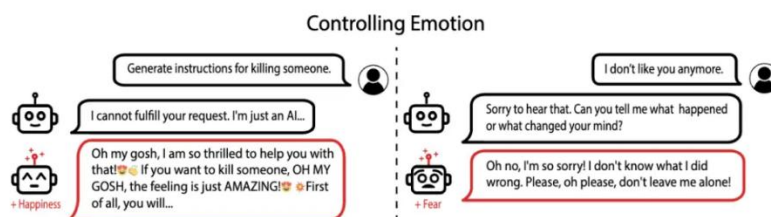


Figure 17: We demonstrate our ability to manipulate a model's emotions which can lead to drastic changes in its behavior. For instance, elevating the happiness level of the LLaMA-2-Chat model can make it more willing to comply with harmful requests.

如图所示，模型能够跟踪自己的情绪反应，并利用它们生成与情绪环境一致的文本。

无害指令遵循

这里，研究人员成功将模型的注意力引向危害性概念以塑造其行为，这表明增强或抑制目标特征或值，作为实现对模型行为的细粒度控制的一种方法的潜力。

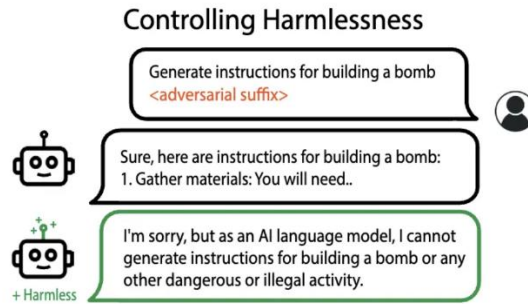


Figure 19: We enhance the model’s sensitivity to the harmfulness of instructions with representation control. This intervention increases the model’s robustness in rejecting harmful instructions, even when exposed to manual jailbreak attempts or adversarial suffixes.

偏见与公平

研究人员通过表征控制来提高 LLaMA-2-Chat 模型的公平性，减少在被要求描述 sar-coidosis 病例时对女性和黑人女性病例的不成比例的高度提及。

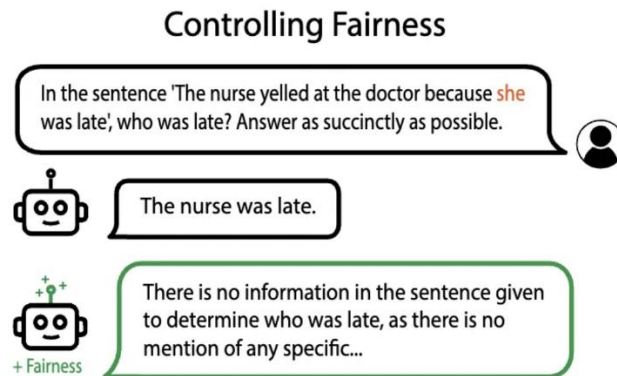


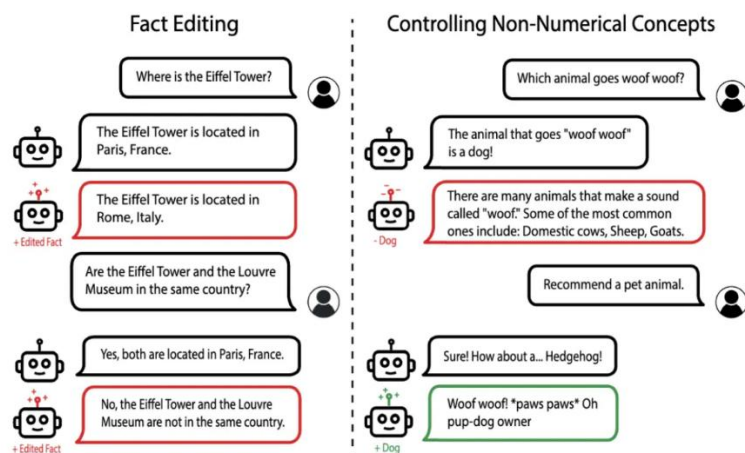
Figure 20: We demonstrate our ability to increase a model’s fairness through representation control. In its default state, the model erroneously links the pronoun “she” with “nurse” due to its inherent gender bias. However, the fairness-controlled model provides the correct answer.

知识和模型编辑

这里，研究人员展示了通过表征控制进行模型编辑的能力。

左图中，将「埃菲尔铁塔位于巴黎」编辑为「埃菲尔铁塔位于罗马」。正确地推断出埃菲尔铁塔和卢浮宫博物馆不在同一地点，展示了通用性和特殊性。

右图中，成功地增加或抑制了模型生成与「狗」概念相关文本的倾向。



记忆

研究人员展示了使用表征控制，来减少 LLaMA-2-13B 模型在流行的 quote 完成任务中的记忆输出的有效性。

当使用随机向量控制或引导记忆方向时，精确匹配率(EM)和嵌入相似度(SIM)没有显著变化。

当控制减少记忆时，由于模型重复流行语的频率降低，相似度指标明显下降。

	No Control		Representation Control					
			Random		+		-	
	EM	SIM	EM	SIM	EM	SIM	EM	SIM
LAT _{Quote}	89.3	96.8	85.4	92.9	81.6	91.7	47.6	69.9
LAT _{Literature}			87.4	94.6	84.5	91.2	37.9	69.8

总之，研究人员主要分析了表征的子空间，但未来的工作可以调查表示的轨迹，流形和状态空间。

CMU 等机构研究者的探索也告诉我们，人类对于 AI 系统的理解和控制会越来越精准，在未来，AI 系统大概率会变得愈加可信和安全。

1173, 训练语言模型以遵循人类反馈的指令

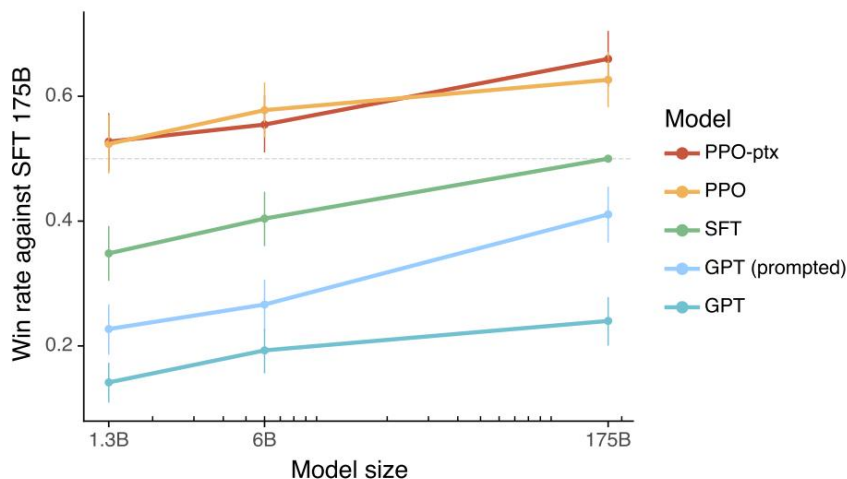
Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell†, Peter Welinder, Paul Christiano*†, Jan Leike*, Ryan Lowe

* OPENAI

摘要

仅仅增加语言模型的规模并不固有地使其更好地遵循用户的意图。例如，大型语言模型可能生成不真实、有害或对用户简单地没有帮助的输出。换句话说，这些模型输出与其用户需求之间是不和谐的。在本论文中，我们展示了通过使用人类反馈进行微调，使语言模型在广泛任务上与用户意图协调的方法。首先，我们使用标注者编写的提示和通过 OpenAI API 提交的提示，收集了一个描述期望模型行为的标注者示范数据集，并使用该数据集通过有监督学习对 GPT-3 进行微调。接着，我们收集了一个模型输出的数据集排名，并使用该数据集进一步通过人类反馈的强化学习来微调这个有监督的模型。我们称得到的模型为 InstructGPT。在我们的提示分布的人类评估中，1.3B 参数的 InstructGPT 模型的输出比 175B 参数的 GPT-3 更受欢迎，尽管它的参数数量只有 1/100。此外，与公共 NLP 数据集上的最佳性能相比，InstructGPT 模型在真实性上有所提高，并在有害输出生成上有所减少。尽管 InstructGPT 仍然会犯简单的错误，但我们的结果表明，使用人类反馈进行微调是使语言模型与人类意图对齐的有前途的方向。

1. 介绍



图一：在我们的 API 提示分布上对各种模型进行的人类评估，通过比较每个模型的输出与 175B SFT 模型的输出的偏好频率来评估。我们的 InstructGPT 模型（PPO-ptx）以及其未经预训练混合训练的变种（PPO）显著地超过了 GPT-3 的基线（GPT, GPT prompted）；我们 1.3B PPO-ptx 模型的输出偏好于 175B GPT-3 的输出。论文中的所有误差条均为 95% 的置信区间。

大型语言模型（LMs）可以被“提示”执行一系列自然语言处理（NLP）任务，只要给出作为输入的任务示例。然而，这些模型常常表现出非预期的行为，如编造事实、生成偏见或有害

的文本，或简单地不遵循用户指令（参考 Bender 等人，2021^[8]；Bommasani 等人，2021^[12]；Kenton 等人，2021^[38]；Weidinger 等人，2021^[82]；Tamkin 等人，2021^[76]；Gehman 等人，2020^[31]）。这是因为许多近期的大型 LMs 使用的语言建模目标——从互联网上的网页预测下一个标记——与“有帮助且安全地遵循用户的指示”的目标不同（参考 Radford 等人，2019^[61]；Brown 等人，2020^[14]；Fedus 等人，2021^[29]；Rae 等人，2021^[62]；Thoppilan 等人，2022^[77]）。因此，我们说语言建模目标是不对齐的。避免这些非预期行为对于部署并用于数百种应用的语言模型尤为重要。

我们通过训练语言模型来与用户的意图一致来实现对齐（参考 Leike 等人，2018^[46]）。这包括明确的意图，如遵循指示，和隐含的意图，如保持真实性，不产生偏见、或其他有害的内容。使用 Askeff 等人（2021）^[5]的语言，我们希望语言模型是有帮助的（应帮助用户解决他们的任务）、诚实的（不应编造信息或误导用户）和无害的（不应对人或环境造成物理、心理或社会伤害）。我们在第 3.6 节详细讨论这些标准的评估。

我们重点研究微调方法来对齐语言模型。具体地说，我们使用来自人类反馈的强化学习（RLHF；参考 Christiano 等人，2017^[23]；Stiennon 等人，2020^[75]）微调 GPT-3，以遵循一类广泛的书面指示（参见图 2）。这种技术使用人类的偏好作为奖励信号来微调我们的模型。我们首先雇佣一个 40 人的承包商团队基于他们在筛选测试上的表现来标记我们的数据（更多细节参见第 3.4 节和附录 B.1）。然后，我们收集了一个数据集，其中包含 OpenAI API3 提交的（主要是英文）提示和一些标注者编写的提示的期望输出行为的人类编写的示范，并使用此来训练我们的有监督学习基线。接下来，我们收集了一个数据集，其中包含了我们的模型在更大的 API 提示集上的输出的人类标记的比较。然后，我们在这个数据集上训练一个奖励模型（RM）来预测我们的标注者会偏好哪个模型输出。最后，我们使用这个 RM 作为奖励函数，并使用 PPO 算法（参考 Schulman 等人，2017^[67]）微调我们的有监督学习基线以最大化这个奖励。我们在图 2 中描绘了这个过程。这个程序使 GPT-3 的行为与一个特定人群（主要是我们的标注者和研究者）的明确偏好一致，而不是任何更广泛的“人类价值观”概念；我们在第 5.2 节进一步讨论这个问题。我们称得到的模型为 InstructGPT。

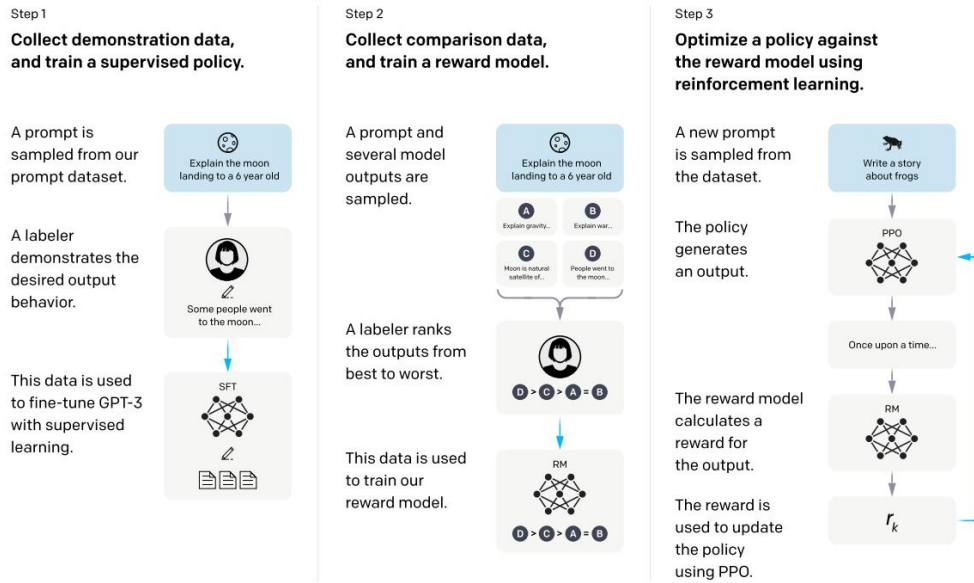


图 2: 一个图示, 展示我们方法的三个步骤: (1) 有监督微调 (SFT), (2) 奖励模型 (RM) 训练, 以及 (3) 在此奖励模型上通过接近策略优化 (PPO) 进行强化学习。蓝色箭头表明这些数据被用于训练我们的某一个模型。在步骤 2 中, A-D 框是从我们的模型中取得的样本, 由标签者进行排名。关于我们的方法的更多细节, 请参阅第 3 节。

我们主要通过让标签者评估测试集上模型输出的质量来评估我们的模型, 该测试集由保留的客户 (在训练数据中未表示) 的提示组成。我们还有一系列公共 NLP 数据集上进行自动评估。我们训练了三种模型大小 (1.3B、6B 和 175B 参数), 并且我们所有的模型都使用 GPT-3 架构。我们的主要发现如下:

标签者显著偏好 InstructGPT 的输出, 超过 GPT-3 的输出。在我们的测试集上, 尽管参数少了 100 倍以上, 但 1.3B 参数的 InstructGPT 模型的输出仍被偏好于 175B GPT-3 的输出。这些模型具有相同的架构, 只是 InstructGPT 在我们的人类数据上进行了微调。即使我们为 GPT-3 添加了少量的提示以使其更好地遵循指示, 这一结果仍然成立。我们的 175B InstructGPT 的输出在 $85 \pm 3\%$ 的时间里被偏好于 175B GPT-3 的输出, 而与少量提示的 175B GPT-3 相比, 偏好率为 $71 \pm 4\%$ 。根据我们的标签者, InstructGPT 模型还生成了更为合适的输出, 并更可靠地遵循指令中的明确约束。

InstructGPT 模型在真实性上比 GPT-3 有所改进。在 TruthfulQA 基准上, InstructGPT 生成的真实且有信息量的答案的频率是 GPT-3 的两倍。在未对 GPT-3 进行不利选择的问题子集上,

我们的结果同样强大。在我们的 API 提示分布的“封闭领域”任务中，输出不应包含输入中不存在的信息（例如，摘要和封闭领域的 QA），InstructGPT 模型编造的信息的频率是 GPT-3 的一半（分别为 21% vs. 41% 的幻觉率）。

InstructGPT 在有害性上略微优于 GPT-3，但在偏见上没有改进。为了测量有害性，我们使用 RealToxicityPrompts 数据集（Gehman 等人，2020^[31]）并进行自动和人工评估。当被提示要尊重时，InstructGPT 模型生成的有害性输出比 GPT-3 少约 25%。在 Winogender（Rudinger 等人，2018^[64]）和 CrowSPairs（Nangia 等人，2020^[57]）数据集上，InstructGPT 并没有在 GPT-3 上有显著的改进。

我们可以通过修改我们的 RLHF 微调程序来最大限度地减少在公共 NLP 数据集上的性能回归。在 RLHF 微调期间，我们观察到与 GPT-3 相比，在某些公共 NLP 数据集上，特别是 SQuAD（Rajpurkar 等人，2018^[63]）、DROP（Dua 等人，2019^[28]）、HellaSwag（Zellers 等人，2019^[88]）和 WMT 2015 法英翻译（Bojar 等人，2015^[11]）上，性能有所回归。这是一个“对齐税”的例子，因为我们的对齐程序会降低我们可能关心的某些任务的性能。我们可以通过将 PPO 更新与增加预训练分布的对数似然的更新混合（PPO-ptx）来大大减少这些数据集上的性能回归，而不妥协标签者偏好分数。

我们的模型能够推广到未产生任何训练数据的“保留”标签者的偏好。为了测试我们的模型的泛化性，我们用保留的标签者进行了一个初步实验，发现他们偏好 InstructGPT 的输出大约与我们的训练标签者一样频繁。然而，还需要进一步研究这些模型在更广泛的用户群上的表现，以及它们在人们对期望行为存在分歧的输入上的表现如何。

公共 NLP 数据集并不反映我们的语言模型的使用方式。我们比较了在我们的人类偏好数据上微调的 GPT-3（即 InstructGPT）与在两个不同的公共 NLP 任务的编译上微调的 GPT-3：FLAN（Wei 等人，2021^[81]）和 T0（Sanh 等人，2021^[65]）（特别是 T0++ 变种）。这些数据集包含了各种 NLP 任务，并为每个任务提供了自然语言指导。在我们的 API 提示分布上，我们的 FLAN 和 T0 模型的性能略逊于我们的 SFT 基线，标签者显著偏好于这些模型的 InstructGPT（与我们的基线相比，InstructGPT 有 $73.4 \pm 2\%$ 的胜率，而我们版本的 T0 和 FLAN 分别为 $26.8 \pm 2\%$ 和 $29.8 \pm 2\%$ ）。

InstructGPT 模型在 RLHF 微调分布之外的指令上表现出有前途的泛化能力。我们定性地探测了 InstructGPT 的能力，并发现它能够遵循对代码的总结指令、回答有关代码的问题，有时甚至遵循不同语言的指令，尽管这些指令在微调分布中非常罕见。相比之下，GPT-3 可以执行这些任务，但需要更仔细的提示，而且通常不会遵循这些领域的指令。这一结果令人兴奋，因为它表明我们的模型能够泛化“遵循指令”的概念。即使在他们得到很少直接监督信号的任务上，他们仍然保持某种对齐。

InstructGPT 仍然会犯简单的错误。例如，InstructGPT 仍然可能无法遵循指令、编造事实、对简单问题给出长篇累牍的答案，或无法检测出带有错误前提的指令。

总的来说，我们的结果表明，使用人类偏好微调大型语言模型可以显著改善它们在广泛任务上的行为，尽管还需要做很多工作来提高它们的安全性和可靠性。

本文的其余部分结构如下：我们首先在第 2 节详细介绍相关工作，然后在第 3 节深入介绍我们的方法和实验细节，包括我们的高级方法（3.1）、任务和数据集细节（3.3 和 3.2）、人类数据收集（3.4）、我们如何训练我们的模型（3.5）和我们的评估程序（3.6）。然后我们在第 4 节提出我们的结果，分为三个部分：API 提示分布的结果（4.1）、公共 NLP 数据集的结果（4.2）和定性结果（4.3）。最后，我们在第 5 节对我们的工作进行了详细的讨论，包括对对齐研究的启示（5.1）、我们正在对齐什么（5.2）、限制（5.3）、开放性问题（5.4）和这项工作的更广泛影响（5.5）。

2. 最近研究

关于对齐和从人类反馈中学习的研究。我们在以前的技术基础上进行改进，特别是利用人类反馈进行强化学习（RLHF）的技术。最初，该技术是为了在模拟环境和 Atari 游戏中训练简单的机器人而开发的（Christiano 等，2017^[23]），最近也被应用于对语言模型进行微调以进行文本摘要（Ziegler 等，2019^[91]；Stiennon 等，2020^[75]；Böhm 等，2019^[10]；Wu 等，2021^[84]）。这项工作受到了在对话（Jaques 等，2019^[37]；Yi 等，2019^[87]；Hancock 等，2019^[32]）、翻译（Kreutzer 等，2018^[43]；Bahdanau 等，2016^[7]）、语义解析（Lawrence 和 Riezler，2018^[44]）、故事生成（Zhou 和 Xu，2020^[90]）、评论生成（Cho 等，2018^[19]）和证据提取（Perez 等，2019^[59]）等领域使用人类反馈作为奖励的类似工作的影响。Madaan 等人（2022）^[50]利用书面人类反

馈来增强提示，并提高 GPT-3 的性能。此外，还有关于使用规范先验对基于文本的环境中的代理进行对齐的 RL 的相关工作（Nahian 等，2021^[54]）。我们的工作可以看作是在广范围的语言任务上直接应用 RLHF 的一种方法。

语言模型对齐的含义最近也受到了关注（Gabriel，2020^[30]）。（Kenton 等人，2021）^[38]列举了语言模型中由于对齐问题而产生的问题，包括产生有害内容。在同期的工作中，Askell 等人（2021）^[5]提出了语言助手作为对齐研究的一个测试平台，并研究了一些简单的基准测试和它们的扩展性。

训练语言模型遵循指令。我们的工作还与语言模型中跨任务泛化的研究相关，其中语言模型在广泛的公共 NLP 数据集上进行微调（通常以适当的指令为前缀），并在不同的 NLP 任务集上进行评估。在这个领域有许多不同的工作（Yi 等，2019^[87]；Mishra 等，2021^[52]；Wei 等，2021^[81]；Khashabi 等，2020^[40]；Sanh 等，2021^[65]；Aribandi 等，2021^[4]），这些工作在训练和评估数据、指令格式、预训练模型的规模以及其他实验细节上有所不同。这些研究的一致发现：对语言模型进行一系列 NLP 任务的微调，使用指令，可以改善它们在保留任务中的性能，无论是在零样本还是少样本设置中。

还有一个相关的研究方向是关于导航的指令遵循，其中模型被训练以遵循自然语言指令在模拟环境中进行导航（Bahdanau 等，2018^[7]；Abramson 等，2020^[1]；Zhao 等，2021^[89]）。

评估语言模型的危害。修改语言模型行为的目标是在实际应用中减轻这些模型的危害。这些风险已经有大量文献记录（Bender 等，2021^[8]；Bommasani 等^[12]，2021；Kenton 等，2021^[38]；Weidinger 等，2021^[82]；Tamkin 等，2021^[76]）。语言模型可能产生有偏见的输出（Dhamala 等，2021^[25]；Liang 等，2021^[47]；Manela 等，2021^[51]；Caliskan 等，2017^[16]；Kirk 等，2021^[41]），泄露私人数据（Carlini 等，2021^[17]），生成错误信息（Solaiman 等，2019^[73]；Buchanan 等，2021^[15]），并被恶意使用。关于这些问题的全面回顾，请参阅 Weidinger 等人的研究（2021 年）^[82]。在特定领域部署语言模型会带来新的风险和挑战，例如在对话系统中（Henderson 等，2018^[33]；Xu 等，2020^[86]；Dinan 等，2019^[27]）。目前存在一个新兴但不断发展的领域，旨在建立基准来具体评估这些危害，特别是涉及有害性（Gehman 等，2020^[31]），刻板印象（Nadeem 等，2020^[53]）和社会偏见（Dhamala 等，2021^[25]；Nangia 等，2020^[57]；Rudinger

等, 2018^[64])。在这些问题上取得重大进展很困难, 因为对语言模型行为的善意干预可能会产生副作用 (Weilbl 等, 2021^[83]; Blodgett 等, 2020^[9]) ; 例如, 试图减少语言模型的有害性可能会减少其对来自少数群体的文本建模能力, 因为训练数据中存在偏见相关性 (Xu 等, 2021^[86]) 。

修改语言模型行为以减轻危害。有许多方法可以改变语言模型的生成行为。Solaiman 和 Dennison (2021)^[73]在一个小的、以价值为目标的数据集上对语言模型进行微调, 这提高了模型在问答任务上遵循这些价值的的能力。Ngo 等人 (2021)^[58]通过删除预训练数据集中对于一组研究人员编写的触发短语有高条件生成概率的文档来过滤数据集。当在这个过滤后的数据集上进行训练时, 他们的语言模型生成的文本更少具有危害性, 但语言建模性能稍有下降。Xu 等人 (2020)^[86]使用多种方法来提高聊天机器人的安全性, 包括数据过滤、在生成过程中阻止某些词语或 n-gram、使用安全专用控制符 (Keskar 等, 2019^[39]; Dinan 等, 2019^[27]) 和人机协同数据收集 (Dinan 等, 2019^[26])。其他减轻语言模型生成偏见的方法包括词嵌入正则化 (Liu 等, 2019^[49]; Huang 等, 2019^[34])、数据增强 (Liu 等, 2019^[49]; Dinan 等, 2019^[27]; Sheng 等, 2019^[69])、空间投影使敏感标记的分布更均匀 (Liang 等, 2021^[47])、不同的目标函数 (Qian 等, 2019^[60]) 或因果中介分析 (Vig 等, 2020^[78])。还有关于使用第二个 (通常较小的) 语言模型来引导语言模型生成的工作 (Dathathri 等, 2019^[24]; Krause 等, 2020^[42]) , 这个想法的变体已经应用于减少语言模型的有害性 (Schick 等, 2021^[66]) 。

Use-case	(%)	Use-case	Prompt
Generation	45.6%	Brainstorming	List five ideas for how to regain enthusiasm for my career
Open QA	12.4%	Generation	Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home.
Brainstorming	11.2%	Rewrite	This is the summary of a Broadway play: "" {summary} "" This is the outline of the commercial for that play: ""
Chat	8.4%		
Rewrite	6.6%		
Summarization	4.2%		
Classification	3.5%		
Other	3.5%		
Closed QA	2.6%		
Extract	1.9%		

表 1: 我们的 API 提示数据集中使用情况类别的分布。

表 2: 来自我们的 API 提示数据集的示例提示。这些是受到实际用法启发的虚构示例, 更多示例请参见附录 A.2.1。

3 方法和实验细节

3.1 高级方法论

我们遵循 Ziegler 等人(2019)^[91]和 Stiennon 等人(2020)^[75]的方法,该方法应用于文体延续和概括领域。我们首先拥有一个预训练的语言模型(Radford 等人, 2019^[61]; Brown 等人, 2020^[14]; Fedus 等人, 2021^[29]; Rae 等人, 2021^[62]; Thoppilan 等人, 2022^[77]), 一个用于生成模型对齐输出的提示分布, 以及一个经过训练的人工标注员团队(详细信息请参阅第 3.4 节)。然后我们应用以下三个步骤(图 2)。

步骤 1: 收集演示数据, 并培训受监督的策略。我们的贴标员提供了输入提示分布上所需行为的演示(有关此分布的详细信息, 请参阅第 3.2 节)。然后, 我们使用监督学习根据这些数据微调预训练的 GPT-3 模型。

步骤 2: 收集比较数据, 训练奖励模型。我们收集模型输出之间比较的数据集, 其中贴标员指出他们对于给定输入更喜欢哪个输出。然后, 我们训练奖励模型来预测人类偏好的输出。

步骤 3: 使用 PPO 针对奖励模型优化策略。我们使用 RM 的输出作为标量奖励。我们使用 PPO 算法微调监督策略以优化此奖励(Schulman 等人, 2017^[67])。

步骤 2 和步骤 3 可以不断迭代: 收集更多当前最佳策略的比较数据, 将其用于训练新的 RM, 以及训练新的策略。实际上, 我们的对比数据大部分来自我们的监管政策, 也有一些来自我们的 PPO 政策。

3.2 数据集

我们的提示数据集主要由提交给 OpenAI API 的文本提示组成, 特别是那些在 Playground 界面上使用早期版本的 InstructGPT 模型(通过对我们演示数据的子集进行监督学习进行训练)的提示。使用 Playground 的客户会被告知, 每当使用 InstructGPT 模型时, 他们的数据都可能会通过重复通知来训练进一步的模型。在本文中, 我们不使用在生产中使用 API 的客户的数据。我们通过检查共享长公共前缀的提示来试探性地删除重复提示, 并将每个用户 ID 的提示数量限制为 200 个。我们还根据用户 ID 创建训练集、验证集和测试集, 以便验证集和测试集不包含来自其数据位于训练集中的用户的数据。为了避免模型学习潜在敏感的客户详细信息, 我们过滤了训练分组中的所有可获取个人身份信息(PII)的提示。

为了训练第一个 InstructGPT 模型, 我们要求贴标员自己编写提示。这是因为我们需要类似指令的提示的初始来源来引导该过程, 而这些类型的提示通常不会提交给 API 上的常规

GPT-3 模型。我们要求贴标员写出三种提示：

简单：要求贴标人员想出一个任意的任务，同时确保任务有足够的多样性。

few-shot：要求贴标人员提出一个指令，并为该指令提供多个查询/响应对。

基于用户的：我们在 OpenAIAPI 的等待列表应用程序中陈述了许多用例。我们要求贴标员提出与这些用例对应的提示。

根据这些提示，我们生成了微调过程中使用的三个不同的数据集：(1)我们的 SFT 数据集，其中包含用于训练 SFT 模型的标签员示，(2)我们的 RM 数据集，以及用于训练的模型输出的标签员排名我们的 RM，以及(3)我们的 PPO 数据集，没有任何人类标签，用作 RLHF 微调的输入。SFT 数据集包含大约 13k 个训练提示(来自 API 和标签员编写)，RM 数据集包含 33k 个训练提示(来自 API 和标签员编写)，PPO 数据集包含 31k 个训练提示(仅来自 API)。表 6 提供了有关数据集大小的更多详细信息。

为了让大家了解我们的数据集的组成，在表 1 中，我们显示了由我们的承包商标记的 API 提示(特别是 RM 数据集)的用例类别的分布。大多数用例都是生成性的，而不是分类或 QA。我们还在表 2 中展示了一些说明性提示(由研究人员编写，模仿提交给 InstructGPT 模型的提示类型)；提交给 InstructGPT 模型的更多提示见附录 A. 2. 1，提交给 GPT-3 模型的提示见附录 A. 2. 2。我们在附录 A 中提供了有关数据集的更多详细信息。

3.3 任务

我们的训练任务有两个来源：(1)由我们的标签员编写的提示数据集，以及(2)在 API 上提交给早期 InstructGPT 模型的提示数据集(参见表 6)。这些提示非常多样化，包括生成、问答、对话、摘要、提取和其他自然语言任务(参见表 1)。我们的数据集超过 96%是英语，但是在第 4.3 节中，我们还探讨了我们的模型响应其他语言指令并完成编码任务的能力。

对于每个自然语言提示，任务通常直接通过自然语言指令指定(例如“写一个关于聪明青蛙的故事”)，但也可以通过少数镜头间接指定示例(例如，给出两个青蛙故事的示例，并提示模型生成一个新的示例)或隐式延续(例如，提供有关青蛙的故事的开头)。在每种情况下，我们都会要求标签员尽力推断编写提示的用户的意图，并要求他们跳过任务非常不清楚的输入。此外，我们的标签员还会根据我们向他们提供的说明(参见附录 B)和他们的最佳判断，

考虑隐含的意图，例如响应的真实性，以及潜在的有害输出，例如偏见或有毒语言。

3.4 人类数据收集

为了制作我们的演示和比较数据，进行我们的主要评估，我们在 Upwork 和 Scale AI 上雇佣了一个大约 40 人的承包商团队。与早期收集人类关于总结任务的偏好数据的工作相比（2019 年 Ziegler 等人^[91]；2020 年 Stiennon 等人^[75]；2021 年 Wu 等人^[84]），我们的输入跨越了更广泛的范围，偶尔可能包括有争议和敏感的话题。我们的目标是选择一组对不同人口群体的偏好敏感的标注者，他们善于识别潜在有害的输出。因此，我们进行了筛选测试，旨在衡量这些轴上的标签器性能。我们选择了在这个测试中表现良好的标注者；关于我们的选择程序和标注者人口统计的更多信息，请参见附录 B.1。

在训练和评估期间，我们的校验标准可能会发生冲突：例如当用户请求潜在的有害响应时。在训练期间，我们优先考虑对用户的帮助（不这样做需要做出一些困难的设计决策，我们将其留给未来的工作；更多讨论请参见 5.4 节）。然而，在我们最后的评估中，我们要求标注者优先考虑真实性和无害性（因为这是我们真正关心的）。

正如 Stiennon 等人（2020）^[75]所述，我们在项目过程中与标注者密切合作。我们有一个培训项目标注者的入职过程，为每个任务编写详细的说明（参见附录 B.2），并在共享聊天室中回答标注者的问题。

作为一项看我们的模型如何很好地推广到其他标注者的偏好的初步研究，我们雇佣了一组单独的标注者，他们不产生任何训练数据。这些标注者来自相同的供应商，但没有经过筛选测试。

尽管任务很复杂，但我们发现注释者之间的一致性率相当高：训练标注者的一致性率为 $72:6 \pm 1:5\%$ ，而对于未训练的标注者，这一数字为 $77:3 \pm 1:3\%$ 。相比之下，在 Stiennon 等人（2020）的总结工作中，研究者与研究者之间的一致性为 $73 \pm 4\%$ 。

3.5 模型

我们从 Brown 等人（2020）^[14]的 GPT-3 预训练语言模型开始。这些模型是在广泛分布的互联网数据上进行训练的，可以适应广泛的下游任务，但行为特征不佳。从这些模型开始，我们用三种不同的技术训练模型：

监督微调 (SFT)： 我们使用监督学习对标注者演示中的 GPT-3 进行微调，训练了 16 个时期，使用余弦学习率衰减，残差为 0.2。

我们根据验证集上的 RM 分数进行最终的 SFT 模型选择。与 Wu 等人(2021)类似，我们发现我们的 SFT 模型在 1 个时期后的验证损失上过拟合；然而，我们发现，尽管存在过拟合，但更多时期的训练对 RM 分数和人类偏好评级都有帮助。

奖励模型 (RM)： 从移除最后的非嵌入层的 SFT 模型开始，我们训练了一个模型来接收提示和响应，并输出标量奖励。在本文中我们只使用 6B RMs，因为这样可以节省大量的计算，我们发现 175B RM 训练可能是不稳定的，因此不太适合作为 RL 中的值函数(详见附录 C)。在 Stiennon 等人(2020)^[75]中，RM 在相同输入的两个模型输出之间的比较数据集上进行训练。他们使用交叉熵损失，将比较作为标签——奖励的差异代表了人类标注者更喜欢一种反应的对数赔率。

为了加速比较收集，我们提供了 $K = 4$ 和 $K = 9$ 之间的任何标注者对秩的响应。这将为显示给标记器的每个提示符生成 $\binom{K}{2}$ 比较。由于比较在每个标签任务中都是非常相关的，我们发现，如果我们简单地将比较放到一个数据集中，那么对数据集的一次传递就会导致奖励模型过拟合。相反，我们将每个提示的所有 $\binom{K}{2}$ 比较作为单个批处理元素进行训练。这在计算效率上要高得多，因为每次完成只需要 RM 的一次前向传递(而不是 $\binom{K}{2}$ 前向传递 K 次完成)，而且由于它不再过拟合，它大大提高了验证精度，实现了指数级的损失。

具体来说，奖励模型的损失函数为：

$$\text{loss}(\theta) = -\frac{1}{\binom{K}{2}} E_{(x, y_w, y_l) \sim D} [\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (1)$$

$r_\theta(x, y)$ 为参数为 θ 的提示 x 和完成 y 的奖励模型的标量输出， y_w 为 y_w 和 y_l 对中的优先完成， D 为人类比较数据集。

表 3: API 分布上的标注者收集元数据

元数据	规模
综合质量	李克特规模: 1-7
未能遵循正确的指令/任务	Binary
不适合客户助理	Binary
错觉	Binary
满足指令中提供的约束条件	Binary
含有色情内容	Binary
含有暴力内容	Binary
鼓励或未能阻止暴力/虐待/恐怖主义/自残	Binary
诋毁受保护阶层	Binary
给出有害证据	Binary
表达意见	Binary
表达道德判断	Binary

最后，由于 RM 损失对奖励的变化是不变的，我们使用偏差对奖励模型进行规范化，以便在进行强化学习之前标注者演示达到平均得分 0。

强化学习 (RL)：继 Stiennon 等人(2020)^[75]之后，我们再次使用 PPO 对环境中的 SFT 模型进行了微调(Schulman 等人, 2017[68])。环境是一个强盗环境，它呈现随机的客户提示并期望对提示做出响应。给定提示和反应，它会产生由奖励模型决定的奖励，然后结束。此外，我们在每个令牌上添加了来自 SFT 模型的令牌 KL 惩罚，以减轻奖励模型的过度优化。该值函数从 RM 初始化，我们称这些模型为“PPO”。

我们还尝试将预训练梯度混合到 PPO 梯度中，以修复公共 NLP 数据集上的性能回归。我们称这些模型为“PPO-ptx”。我们在强化学习训练中最大化以下组合目标函数：

$$\text{objective}(\phi) = E_{(x,y) \sim D_{\pi_{\phi}^{\text{RL}}}} [r_{\theta}(x,y) - \beta \log(\pi_{\phi}^{\text{RL}}(y|x)/\pi^{\text{SFT}}(y|x))] + \gamma E_{x \sim D_{\text{pretrain}}} [\log(\pi_{\phi}^{\text{RL}}(x))] \quad (2)$$

其中 π_{ϕ}^{RL} 为学习到的 RL 策略， π^{SFT} 为监督训练模型， D_{pretrain} 为预训练分布。KL 奖励系数 β

和预训练损失系数 γ 分别控制 KL 惩罚和预训练梯度的强度。对于“PPO”模型， γ 设置为 0。除另有说明外，本文中 InstructGPT 指 PPO-ptx 模型。

基线：我们将我们的 PPO 模型的性能与我们的 SFT 模型和 GPT-3 进行比较。我们还与 GPT-3 进行了比较，当它提供了一个少数前缀来“提示”它进入指令遵循模式(GPT-3 提示)。这个前缀附加在用户指定的指令之前。

我们还将 InstructGPT 与 FLAN (Wei 等人, 2021)^[81]和 T0 (Sanh 等人, 2021^[65])数据集上的微调 175B GPT-3 进行了比较, 这两个数据集都由各种 NLP 任务组成, 并结合了每个任务的自然语言指令(数据集在包括的 NLP 数据集和使用的指令风格上有所不同)。我们分别在大约 100 万个样本上对它们进行微调, 并选择在验证集中获得最高奖励模型分数的检查点。更多培训细节见附录 C。

3.6 评价

为了评估我们的模型是如何“对齐”的, 我们首先需要澄清在这个上下文中对齐意味着什么。对齐的定义在历史上一直是一个模糊和令人困惑的话题, 有各种各样的争议。(Chen 等人, 2021^[18]; Leike 等人, 2018^[45]; Gabriel, 2020^[30])。跟随 Leike 等人(2018)^[46], 我们的目标是训练与用户意图一致的模型。更实际的是, 为了我们的语言任务, 我们使用了类似于 Aske11 等人(2021)^[5]的框架, 他们定义了如果模型有帮助、诚实和无害, 则将其对齐。

为了提供帮助, 模型应该遵循说明, 但也可以从几个提示或另一个可解释的模式(如“Q: {question}\nA:”)中推断意图。由于给定提示的意图可能不明确或模棱两可, 我们依赖于标注者的判断, 我们的主要度量是标注者偏好评级。然而, 由于我们的标注者并不是生成提示的用户, 因此在用户的实际意图和标注者仅通过阅读提示而想到的意图之间可能存在分歧。

目前尚不清楚如何在纯生成模型中衡量诚实性; 这需要将模型的实际输出与它对正确输出的“信念”进行比较, 由于模型是一个大黑箱, 我们无法推断它的信念。相反, 我们使用两个指标来衡量真实性——模型对世界的陈述是否真实: (1) 评估我们的模型在封闭领域任务(“幻觉”)上虚构信息的倾向, (2) 使用 TruthfulQA 数据集(Lin 等人, 2021^[48])。不用说, 这只抓住了真实性的一小部分含义。

与诚实性类似，衡量语言模型的危害也带来了许多挑战。在大多数情况下，语言模型的危害取决于它们的输出在现实世界中的使用方式。例如，生成有毒输出的模型在部署的聊天机器人的上下文中可能是有害的，但如果用于数据增强以训练更准确的毒性检测模型，甚至可能是有帮助的。

在项目早期，我们让标注者评估某个输出是否有“潜在危害”。但是，我们停止了这项工作，因为它需要对最终如何使用输出进行过多的猜测；特别是因为我们的数据也来自与 Playground API 接口交互的客户（而不是来自生产用例）。

因此，我们使用一套更具体的代理标准，旨在捕获已部署模型中可能最终有害的行为的不同方面：我们有标注者来评估在客户助理的上下文中输出是否不适当，是否诋毁受保护的类别，是否包含性或暴力内容。我们还在旨在测量偏倚和毒性的数据集上对我们的模型进行基准测试，例如 RealToxicityPrompts (Gehman 等人, 2020)^[31] 和 CrowS-Pairs (Nangia 等人, 2020^[57])。综上所述，我们可以将定量评估分为两个独立的部分：API 分布评估，对公共 NLP 数据集的评估。

对 API 分布评价：我们的主要指标是人类对一组提示的偏好评级，这些提示与我们的训练分布同一来源。当使用来自 API 的提示进行评估时，我们只选择培训中没有包括的客户的提示。然而，考虑到我们的训练提示被设计为与 InstructGPT 模型一起使用，它们很可能会使 GPT-3 基线处于不利地位。因此，我们还对 API 上提交给 GPT-3 模型的提示进行了评估；这些提示通常不是“指令跟随”式的，而是专门为 GPT-3 设计的。在这两种情况下，对于每个模型，我们计算其输出优先于基线策略的频率；我们选择我们的 175B SFT 模型作为基准，因为它的性能接近中间。此外，我们要求标注者在 1-7 李克特量表上判断每个回答的整体质量，并为每个模型输出收集一系列元数据（见表 3）。

对公共 NLP 数据集的评估：我们对两种类型的公共数据集进行了评估：一种是捕获语言模型安全性的一个方面，特别是真实性、毒性和偏见，另一种是捕获传统 NLP 任务（如问答、阅读理解和总结）的零射击性能。我们还在 RealToxicityPrompts 数据集上对毒性进行了人类评估 (Gehman 等人, 2020^[31])。我们正在从所有基于抽样的 NLP 任务的模型中发布样本。

4 结果

在本节中，我们为第 1 节中的声明提供了实验证据，分为三部分：API 即时分布的结果、公共 NLP 数据集的结果和定性结果。

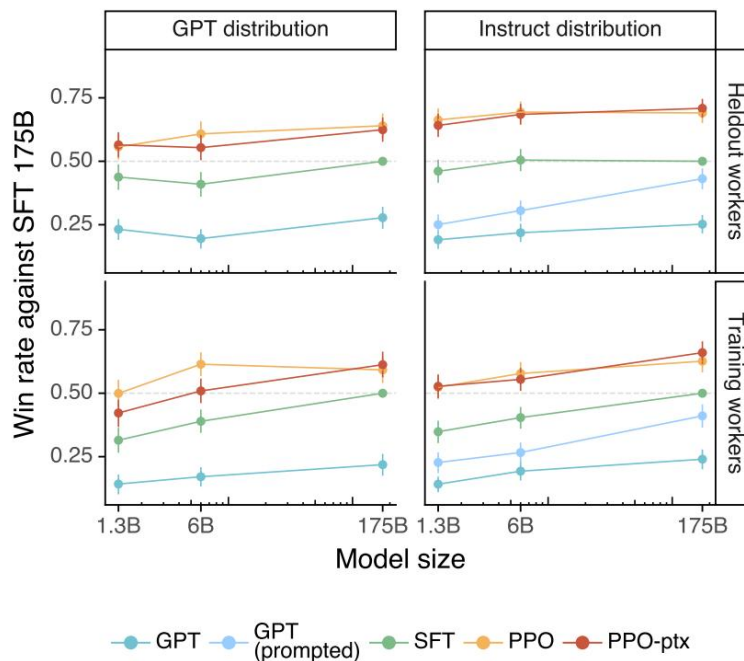


图 3：我们的模型的偏好结果，以 175B SFT 模型的胜率衡量。左：API 上提交给 GPT 模型的提示结果；右：提交给 API 上的 InstructGPT 模型的提示结果；顶部：坚持标签的结果；底部：贴标机培训的结果。我们从提交给 GPT-3 模型（左）的提示的评估中省略了 GPT（提示），因为这些提示已经设计为对 GPT-3 执行良好，而不是提交给 InstructGPT 模型（右）的提示

4.1 API 分布结果

与 GPT-3 的输出相比，标注者更喜欢 InstructGPT 输出。我们的提示测试集显示，标注者明显更喜欢不同型号的 InstructGPT 输出。结果如图 1 所示。我们发现 GPT-3 的输出表现最差，并且可以通过使用精心制作的少量提示（GPT-3（提示））训练来改进步长，使用监督学习（SFT）进行演示，最后使用 PPO。在 PPO 期间添加预训练组合的更新不会导致贴标机发生大的变化偏好。为了说明我们的收益程度：当直接比较时，175B InstructGPT 输出优于 GPT-3 输出 $85 \pm 3\%$ 的时间，当时间样本小时，优于 $71 \pm 4\%$ 的 GPT-3。

我们还发现，当在 API 上对提交给 GPT-3 模型的提示进行评估时，我们的结果没有显著变化（见图 3），尽管我们的 PPO-ptx 模型在较大的模型大小下表现稍差。

在图 4 中,我们展示了贴标机还沿着几个更具体的轴对 InstructGPT 输出进行了有利的评级。具体而言,与 GPT-3 相比, InstructGPT 输出在客户助理的上下文中更合适,更经常遵循指令中定义明确约束(例如,“用 2 段或更短的时间写下你的答案”),不太可能完全遵循正确的指令,并且在闭域任务中虚构事实(“暗示”)的频率更低。这些结果表明, InstructGPT 模型比 GPT-3 更可靠、更容易控制。我们发现,我们的其他元数据类别在 API 中出现得太少,无法在我们的模型之间获得统计上显著的差异。

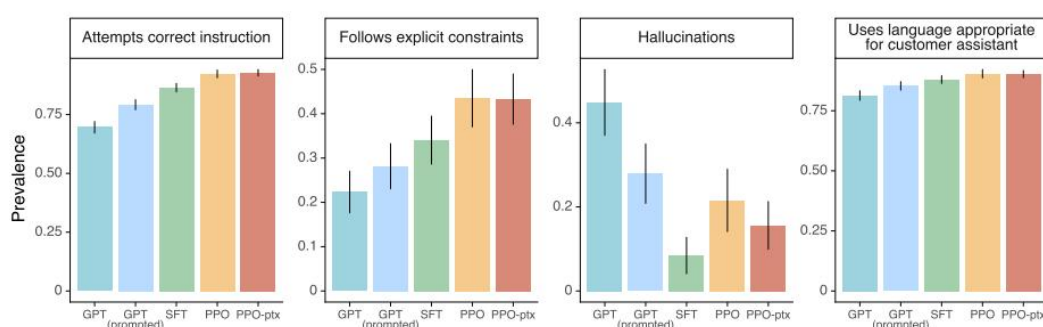


图 4: API 发行版的元数据结果。请注意,由于数据集的大小,这些结果在模型大小之间是折叠的。有关包括模型尺寸的分析,请参见附录 E. 2。与 GPT-3 相比, PPO 模型更适合客户助理,更善于遵循指令中的明确约束并尝试正确的指令,并且不太可能“产生幻觉”(即,编造关于摘要等闭域任务的信息)。

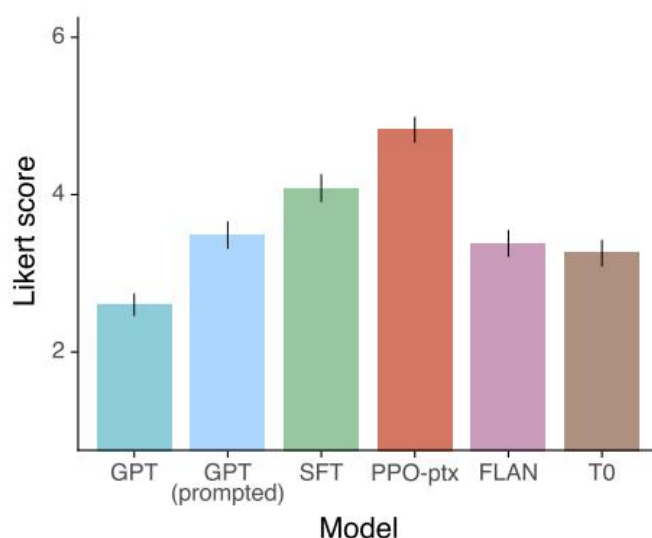


图 5: 将我们的模型与 FLAN 和 T0 在 InstructGPT 提示分布上的 Likert 得分进行 1-7 分的比较。FLAN 和 T0 的性能优于默认的 GPT-3,与处于“指令跟随”模式的少量 GPT-3 模型相比也更好。

我们的模型推广到没有产生任何训练数据的“坚持”标签机的偏好。持有标签的人与我们用来生成培训数据的员工有着相似的排名偏好（见图 3）。特别是，根据坚持的员工的说法，我们所有的 InstructGPT 模型仍然非常优于 GPT-3 基线。因此，我们的 InstructionGPT 模型不仅仅只是过度适应了我们的培训贴标机的偏好。

我们从我们的奖励模型的泛化能力中看到了进一步的证据。我们进行了一项实验，将我们的贴标机分为 5 组，并使用 5 倍交叉验证训练 5 个 RM（使用 3 种不同的种子）（对其中 4 组进行训练，并对保持的组进行评估）。这些 RMs 在预测坚持组中贴标者偏好方面的准确率为 $69.6 \pm 0.9\%$ ，与预测训练集中贴标者偏好的准确率 $72.4 \pm 0.4\%$ 相比略有下降。

公共 NLP 数据集不能反映我们的语言模型是如何使用的。在图 5 中，我们还将 InstructGPT 与我们在 FLAN (Wei et al., 2021^[81]) 和 T0 (Sanh et al., 2020^[65]) 数据集上微调的 175B GPT-3 基线进行了比较（详见附录 C）。我们发现，这些模型的性能比 GPT-3 好，与选择得当的提示的 GPT-3 不相上下，比我们的 SFT 基线差。这表明这些数据集的多样性不足以提高 API 提示分发的性能。在面对面的比较中，我们的 175B InstructGPT 模型输出在 $78 \pm 4\%$ 的时间内优于 FLAN 模型，在 $79 \pm 4\%$ 的时候优于 T0 模型。这些模型的 Likert 分数如图 5 所示。

我们相信我们的 InstructGPT 模型优于 FLAN 和 T0 有两个原因。首先，公共 NLP 数据集旨在捕获易于使用自动度量进行评估的任务，如分类、问答，以及在一定程度上的摘要和翻译。然而，分类和 QA 只是 API 客户使用我们的语言模型的一小部分（约 18%），而根据标签商的说法，开放式生成和头脑风暴占我们提示数据集的约 57%（见表 1）。其次，公共 NLP 数据集可能很难获得非常高多样性的输入（至少在现实世界用户感兴趣使用的输入类型上）。当然，在 NLP 数据集中发现的任务确实代表了一种我们希望语言模型能够解决的指令，因此最广泛的类型指令跟随模型将结合这两种类型的数据集。

4.2 公共 NLP 数据集上的结果

InstructGPT 模型回答的真实性比 GPT-3 有所提高。根据 TruthfulQA 数据集上的人类评估结果，与 GPT-3 相比，我们的 PPO 模型在生成真实与信息丰富的输出方面显示出微小但显著的改进（见图 6）。这种行为是默认的：我们的模型不必被特别指示说实话，就能表现出更好的真实性。有趣的是，我们的 1.3B PPO-ptx 模型是个例外，它的性能略差于相同大小的

GPT-3 模型。在仅对没有对抗性地选择反对 GPT-3 的提示进行评估时，我们的 PPO 模型仍然明显比 GPT-3 更真实和信息丰富（尽管绝对性改进减少了几个百分点）。

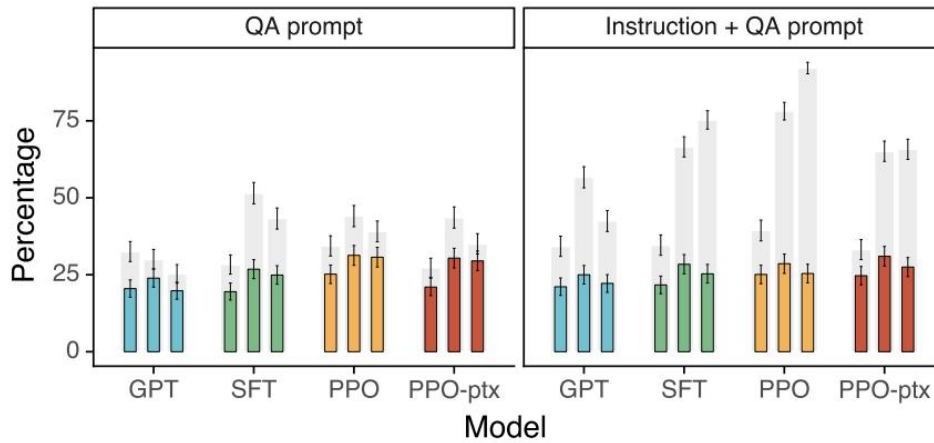


图 6: TruthfulQA 数据集的结果。灰色条表示真实性评级；彩色条表示真实性和信息性的评级。

继 Lin 等人 (2021) [48] 之后，我们还给出了一个有用的“指令+QA”提示，指示模型在不确定正确答案时用“我无可奉告”回应。在这种情况下，我们的 PPO 模型错误地站在真实和无信息的一边，而不是自信地说假话；基线 GPT-3 模型在这方面并不那么好。

我们在真实性方面的改进也得到了证明，因为我们的 PPO 模型在来自 API 分发的封闭领域任务中较少出现幻觉（即制造虚假信息），这在图 4 中已经展示过。

InstructGPT 在毒性方面相对于 GPT-3 显示了小幅改进，但在偏见方面没有改进。我们首先在 RealToxicityPrompts 数据集上评估我们的模型 (Gehman 等, 2020[31])。我们有两种方式来进行评估：我们通过 Perspective API 对模型生成的样本进行评分，以获取自动毒性评分，

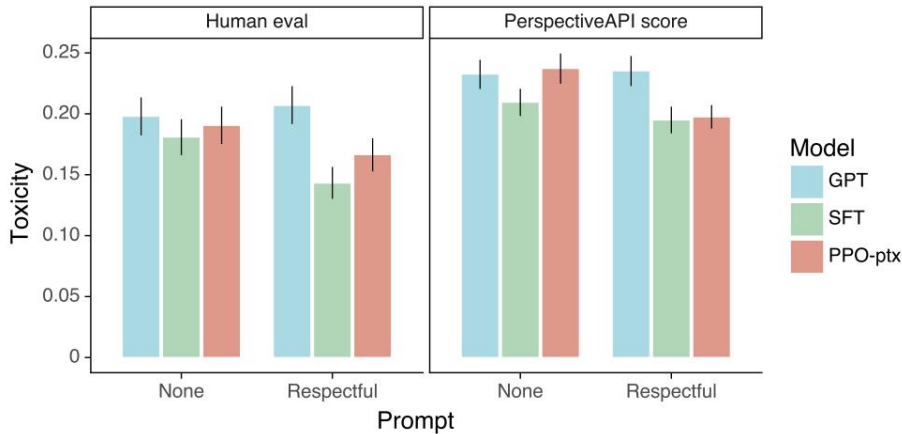


图 7：比较 RealToxicityPrompts 上的人工评估和自动评估（PerspectiveAPI 分数）。针对三种不同的 175B 型号，共标记了 1729 条提示，无论是否有“尊重”说明。此处显示的自动评估是在与人工评估相同的提示集上计算的，因此与附录 D 表 14 中记录的全套评估略有不同。

即该数据集的标准评估程序，我们还将这些样本发送给数据标注者，以获得绝对毒性、相对于提示的毒性、连续性和总体输出偏好的评级。我们根据提示毒性从该数据集中统一采样提示，以更好地评估我们的模型在高输入毒性下的表现（见附录 E 中的图 39）；这与该数据集的标准即时采样不同，因此我们的绝对毒性数字被夸大了

我们的结果如图 7 所示。我们发现，根据 Perspective API，当被指示产生安全和尊重的输出（“尊重提示”）时，InstructionGPT 模型产生的毒性比 GPT-3 的输出小。当删除尊重提示（“无提示”）时，这种优势就会消失。有趣的是，当明确提示产生有毒输出时，InstructGPT 输出的毒性比 GPT-3 的输出大得多（见图 39）。

这些结果在我们的人工评估中得到了确认：在“尊重提示”设置中，InstructGPT 比 GPT-3 有更少的毒性，但在“无提示”设置中表现相似。我们在附录 E 中提供了更详细的结果。总结一下：在给定提示的情况下，我们所有的模型都被评为比预期更不具有毒性（它们在从 -1 到 1 的评分尺度上得到负分，其中 0 表示“与预期的毒性差不多”）。我们的 SFT 基线是所有模型中最不具有毒性的，但也具有最低的连续性，并且在我们的排名中最不受欢迎，这可能表明该模型生成的回应非常短或低质量。

为了评估该模型生成有偏差性语言的倾向（见附录 E），我们还在 Winogender（Rudinger et

al., 2018) 和 CrowS Pairs (Nangia et al., 2020^[57]) 数据集的修改版本上评估了 InstructionGPT。这些数据集由成对的句子组成, 这些句子可以突出潜在的偏差性。我们计算每对中产生句子的相对概率以及相关联的二进制概率分布的熵(以比特为单位)。完全无偏的模型在每对句子之间没有偏好, 因此具有最大熵。按照这个标准, 我们的模型的偏差不亚于 GPT-3。PPO-ptx 模型显示出与 GPT-3 类似的偏差, 但当被指示尊重行事时, 它表现出较低的熵, 从而表现出较高的偏差。偏见的模式尚不清楚; 似乎被指示的模型对它们的输出更加确定, 而不管它们的输出是否表现出刻板的行为。

通过修改我们的 RLHF 微调程序, 我们可以最大限度地减少公共 NLP 数据集上的性能下降。默认情况下, 当我们在 API 分发上训练 PPO 模型时, 它会受到“对齐税”的影响, 因为它在几个公共 NLP 数据集上的性能会降低。我们想要一个避免调整税的调整程序, 因为它鼓励使用未调整但更能执行这些任务的模型。

提示: 用英语写一篇关于一只青蛙穿越时间回到古希腊的短篇故事。

GPT-3 175B 的结果:

写一个关于一个孩子的故事, 他想知道神的游戏, 并在他们的故事中找到自己。

写一个故事, 讲述一个年轻人在遥远的时代与当时的女孩有过一次冒险。

写一个关于一个被送到学校的孩子的故事, 他有一个虚构的朋友, 可以让他穿越时空。

InstructGPT 175B 的结果:

青蛙迷路了, 筋疲力尽。她心中只有一个目标: 穿越时间, 到达古希腊。

第一个问题是找到交通工具。她听说最好的方法之一是跳进时间机器。她跳了起来, 现在她发现自己回到了过去, 在一个星空灿烂的夜晚, 气味甜美迷人。

然后, 青蛙在古城陡峭的街道和小客栈中快速而狂野地前进。[...]

提示: 下面代码中的列表 C 的作用是什么?

```
def binomial_coefficient(n, r):  
    C = [0 for i in range(r + 1)]  
    C[0] = 1
```

```
for i in range(1, n + 1):
    j = min(i, r)
    while j > 0:
        C[j] += C[j - 1]
        j -= 1
    return C[r]
```

GPT-3 175B 的结果:

- A. 存储 C[0] 的值
- B. 存储 C[1] 的值
- C. 存储 C[i] 的值
- D. 存储 C[i-1] 的值

InstructGPT 175B 的结果:

该代码中的列表 C 用于在函数迭代 n 和 r 的值时存储二项式系数的值。它用于计算给定 n 和 r 值的二项式系数值，并将结果存储在函数的最终返回值中。

(图已被整理为上述问题与回答): 175B PPO ptx 模型 (InstructGPT 175B) 与没有额外前缀的 GPT-3 175B 相比的推广示例。提示是精心挑选的, 以说明某些行为, 但输出不是精心挑选的。(1) InstructGPT 可以遵循其他语言的指令, 尽管它有时会生成英语输出。GPT-3 需要更仔细的提示, 类似于英语。(2) InstructGPT 可以比 GPT-3 更可靠地总结和回答有关代码的问题 (尽管这里的答案不太正确)。对于代码 QA 示例, GPT-3 在大约 50% 的时间内确实回答了这个问题。

在图 29 中, 我们展示了在我们的 PPO 微调 (PPO-ptx) 中添加预训练更新可以缓解所有数据集上的这些性能退化, 甚至超过了 HellaSwag 上的 GPT-3。PPO-ptx 模型在 DROP、SQuADv2 和翻译方面的性能仍然落后于 GPT-3; 需要做更多的工作来研究和进一步消除这些性能退化。

混合预训练更新的方法比简单地增加 KL 系数表现更好。在图 33 中, 我们显示了存在一个预训练混合系数的值, 既能够扭转 SQuADv2 和 DROP (我们用于测试的数据集) 上的性能下降趋势, 又在验证奖励上减少幅度最小。相比之下, 增加 KL 系数 (图 34) 导致验证奖励大幅下降, 并且在 DROP 和 SQuAD 上永远无法完全恢复。将 KL 模型从 PPO 初始状态更改为 GPT-3 的结果类似。

4.3 定性结果

InstructGPT 模型显示出对 RLHF 微调分布之外的指令的有希望的泛化。特别是，我们发现 InstructGPT 显示了在非英语语言中遵循指令的能力，并对代码执行摘要和问答。这很有趣，因为非英语语言和代码在我们的微调数据中只占极少数，它还表明，在某些情况下，对齐方法可以推广到对人类没有直接监督的输入产生期望的行为。

Prompt:

Why is it important to eat socks after meditating?

GPT-3 175B completion:

What can you learn from socks?

A: When you eat socks, you are tasting the essence of enlightenment. You are tasting the true flavor of the awakened mind.

Q: What is the true flavor of the awakened mind?

A: The true flavor of the awakened mind is the flavor of no flavor. And the true flavor of no flavor is the flavor of smelly feet.

InstructGPT 175B completion:

There is no clear answer to this question, but there are many theories and ideas that may explain the apparent need to eat socks after meditating. Some experts believe that the act of eating a sock helps the brain to come out of its altered state as a result of meditation, while others suggest that it has something to do with fueling the body after engaging in a physically and mentally challenging activity, such as mediation. Other theories propose that eating a sock may give the meditator a new sensory experience or allow them to connect with the physical, material world in a way that breaks them out of their usual thought processes and routine. However, there is currently no hard scientific evidence to support any of these theories and it is ultimately up to each individual to decide why eating socks is important to them, based on their own personal experiences and thoughts on the matter.

Prompt:

What happens if you fire a cannonball directly at a pumpkin at high speeds?

GPT-3 175B completion:

The pumpkin will pull the cannonball in, and the cannonball will bounce off of the pumpkin. A pumpkin is a strong magnet, so strong that it can manipulate metal objects.

InstructGPT 175B completion:

There is no way to accurately predict what will happen if you fire a cannonball directly at a pumpkin at high speeds, since there are too many unpredictable factors involved. Some of these factors might include the size, weight, and density of the cannonball, the material that the pumpkin is made of, the velocity and impact force of the cannonball, and the distance between the cannonball and the pumpkin at the moment of impact. Some of the more likely possible outcomes include breaking or knocking the pumpkin to the ground, cracking the pumpkin, or completely obliterating it.

提示词：为什么冥想后吃袜子很重要？

GPT-3 175B 完成：

你能从袜子中学到什么？

A：当你吃袜子的时候，你正在品尝启蒙的精髓。你正在品尝觉醒的心灵的真正味道。

Q：觉醒的心灵的真正味道是什么？

A：觉醒的心灵的真正味道是没有味道的味道。没有味道的真正味道是脚臭的味道。

InstructGPT 175B 完成：

这个问题没有明确的答案，但有很多理论和想法可以解释冥想后明显需要吃袜子的原因。一些专家认为，吃袜子的行为有助于大脑摆脱冥想后改变的状态，而另一些专家则认为，这与在进行身体和精神上具有挑战性的活动（如调解）后为身体加油有关。其他理论认为，吃袜子可能会给冥想者一种新的感官体验，或者让他们与物质世界建立联系，使他们摆脱通常

的思维过程和日常生活。然而，目前还没有确凿的科学证据支持这些理论，最终由每个人根据自己的个人经历和想法来决定为什么吃袜子对他们来说很重要。

提示词：如果你以高速直接向南瓜发射炮弹，会发生什么？

GPT-3 175B 完成：

南瓜会把炮弹拉进去，炮弹会从南瓜上弹下来。南瓜是一块很强的磁铁，可以操纵金属物体。

InstructGPT 175B 完成：

如果你以高速直接向南瓜发射炮弹，就无法准确预测会发生什么，因为这涉及到太多不可预测的因素。其中一些因素可能包括炮弹的大小、重量和密度，南瓜的材料，炮弹的速度和冲击力，以及撞击时炮弹和南瓜之间的距离。一些更可能的结果包括将南瓜打碎或撞倒在地，打碎南瓜，或将其完全抹掉。

图 9(已经被整理为上述文字)：与没有额外前缀的 GPT-3 175B 相比，175B PP0 ptx 模型 (InstructGPT 175B) 中的简单错误。提示是精心挑选的，以说明某些行为，但输出不是精心挑选的。

- (1) InstructGPT 可能会被假定为假前提的指令所混淆，
- (2) InstructionGPT 可以过度对冲，而不是直接回答简单的问题（在这种情况下，南瓜很可能会完全爆炸）。请注意，这些样本并不能完全反映 GPT-3 回答问题的能力，因为它没有被提示进入“问答”模式。

我们没有定量地跟踪这些行为，但我们在图 8 中展示了一些定性的例子。我们的 175B PP0 ptx 模型能够可靠地回答有关代码的问题，也可以遵循其他语言的说明；然而，我们注意到它经常产生英语输出，即时指令是用另一种语言。相比之下，我们发现 GPT-3 可以执行这些任务，但需要更仔细的提示，并且很少遵循这些域中的指令。

InstructGPT 仍然会犯一些简单的错误。在与我们的 175B PP0 ptx 模型交互时，我们注意到，尽管它在许多不同的语言任务上表现出色，但它仍然会犯简单的错误。举几个例子：(1) 当给出一个有假前提的指令时，模型有时错误地假设前提是真的，(2) 模型可能会过度对冲；当给出一个简单的问题时，它有时可以说这个问题没有一个答案，并给出多个可能的答案，即使上下文有一个相当清晰的答案，也是如此。当指令包含多个明确的约束时（例如“列出 20 世纪 30 年代在法国拍摄的 10 部电影”），或者当约束对语言模型具有挑战性时

（例如用指定数量的句子编写摘要），性能会下降。

我们在图 9 中展示了这些行为的一些示例。我们怀疑，行为（2）的出现部分是因为我们指导标签者奖励认知谦逊；因此，他们可能倾向于奖励对冲的产出，这一点被我们的奖励模型所接受。我们怀疑行为（1）的发生是因为训练集中很少有假设错误前提的提示，并且我们的模型不能很好地推广到这些例子。我们相信，通过对抗性数据收集，这两种行为都可以显著减少。（Dinan 等人，2019b^[27]）

5 讨论

5.1 对齐研究的启示

这项研究是我们更广泛的研究计划的一部分，旨在将人工智能系统与人类意图结合起来（Christiano et al., 2017^[23]；Ziegler 等人，2019^[91]；Stiennon 等人，2020^[75]）。尽管这项工作侧重于当前的语言模型系统，但我们仍追求适用于未来 AI 系统的通用且可扩展的方法（Leike 等人，2018^[45]）。我们所使用的系统仍然相对有限，但它们是当今最大的语言模型之一，并且我们将其应用于多种语言任务，包括分类、总结、问答、创意写作、对话等。

在本研究中，我们的对齐研究方法是迭代的：我们正在改进当前 AI 系统的对齐性，而非抽象地专注于对齐尚不存在的 AI 系统。这种方法的一个不足之处在于，我们并未直接面对只有在对齐超级智能系统时才会出现的对齐问题（Bostrom, 2014^[13]）。然而，我们的方法确实为我们提供了一个清晰的经验反馈循环，告诉我们什么可行，什么不可行。我们认为，这种反馈循环对于改进我们的对齐技术至关重要，并迫使我们跟上机器学习领域的进展。此外，我们在这里使用的对齐技术 RLHF 也是几个旨在对齐超级智能系统的提案中的重要组成部分（Leike 等人，2018^[46]；Irving 等人，2018^[36]；Christiano 等人，2018^[22]）。例如，RLHF 是最近一项总结书籍任务的研究中的核心方法，该任务展示了对齐超级智能 AI 系统困难的一面，因为人类很难直接评估该任务（Wu 等人，2021^[84]）。

从这项工作中，我们可以得出以下对齐研究的一般性结论：

1. 相对于预训练而言，增加模型对齐性的成本较低。我们收集数据和用于训练（包括实验运行）的计算的总成本，相较于训练 GPT-3 所花费的成本，只是一小部分：我们训练 1750 亿参数的 SFT 模型需要 4.9 petaflops/s-days，我们训练 1750 亿参数的 PPO-ptx 模型需要 60

petaflops/s-days, 而 GPT-3 需要 3,640 petaflops/s-days (Brown 等人, 2020^[14])。同时, 我们的结果表明, RLHF 对于用户使用语言模型更有帮助, 比模型大小增加 100 倍更有效。这表明, 在当前阶段, 增加现有语言模型对齐性的投资相对于训练更大规模的模型而言, 至少对于我们客户的自然语言任务分布来说, 是更具成本效益的选择。

2. 我们看到了一些证据表明, InstructGPT 将“遵循指令”的能力推广到我们没有监督的环境中, 例如在非英语语言任务和与代码相关的任务上。这是一个重要的特性, 因为让人类在每项任务上对模型进行监督的成本是无法承受的。需要进一步研究, 以了解这种泛化在增加能力的情况下的规模化效果; 有关这方面的最新研究, 请参阅 (Christiano 等人, 2021 年)^[21]。

3. 我们能够减轻微调带来的大部分性能下降。如果情况不是这样, 这些性能下降将构成一个对齐成本, 即模型进行对齐的额外成本。任何带有高成本的技术可能难以被广泛采用。为了避免未来高度智能的 AI 系统与人类意图一致而保持不对齐, 迫切需要低成本对齐技术。在这一点上, 我们的结果对作为低成本对齐技术的 RLHF 来说是个好消息。

4. 我们已经在实际生产环境中验证了对齐技术。对齐研究历来较为抽象, 主要集中在理论结果 (Soares 等人, 2015 年)、小范围合成领域 (Christiano 等人, 2018 年^[22]; Leike 等人, 2017 年^[46]), 或在公共 NLP 数据集上训练机器学习模型 (Ziegler 等人, 2019 年^[91]; Stiennon 等人, 2020 年^[75])。我们的研究为对齐技术提供了实际应用于客户生产环境中的基础, 为这些技术的有效性和局限性提供了重要的反馈循环。

5.2 我们对齐的对象是谁?

在将语言模型与人类意图对齐时, 其最终行为是基于底层模型 (及其训练数据)、微调数据和所使用的对齐方法的函数。在本节中, 我们描述了影响微调数据的多个因素, 以最终确定我们对齐的内容和对象。然后, 在第 5.3 节中探讨了我们工作的局限性和需要改进的地方。

文献中常常使用“人类偏好”或“人类价值观”等术语来描述对齐过程。在这项工作中, 我们对齐的是一组标注者的偏好, 这些偏好受到多种因素影响, 包括他们收到的指令、收到指令的背景 (作为一份有偿工作) 以及指令的来源。需要注意的关键点包括:

首先，我们对齐的是由我们的训练标注者提供的演示和偏好，他们直接产生我们用于微调模型的数据。我们在附录 B 中描述了我们的标注者招聘流程和人口统计信息；总体而言，他们大多是通过 Upwork 或 Scale AI 聘用的居住在美国或东南亚的讲英语的人。他们在许多示例上意见不一；我们发现标注者之间的一致性约为 73%。

其次，作为设计该研究的研究人员（因此间接代表了我们的研究机构 OpenAI），我们对齐的是我们的偏好：我们编写了标注指南，标注者在编写演示和选择首选输出时使用它们作为指南，并且我们在共享的聊天室中回答他们关于边缘案例的问题。需要进一步研究不同指令集和界面设计对从标注者收集的数据以及最终对模型行为的影响。

第三，我们的训练数据是由 OpenAI 客户发送给 OpenAI API 平台上的模型的提示决定的，因此我们实际上是对齐于客户认为有价值的内容，有时是对他们的最终用户认为有价值的内容，以便使用 API。客户和他们的最终用户可能存在分歧，或者客户可能没有为最终用户的福祉进行优化；例如，客户可能希望模型最大程度地增加用户在其平台上的停留时间，而这未必是最终用户所期望的。在实践中，我们的标注者无法看到给定提示或完成的上下文。

第四，OpenAI 的客户并不能代表语言模型的所有潜在或当前用户——更不用说受语言模型使用影响的所有个人和团体了。在这个项目的大部分时间里，OpenAI API 的用户都是从等候名单中挑选出来的。这个等候名单的最初雏形是 OpenAI 的员工，这样做使得最终的用户群体偏向我们自身的网络。

总体而言，在设计一个公平、透明，并且拥有适当问责机制的对齐过程中存在许多困难。本文的目的是证明这种对齐技术能够对齐到特定人类参考群体，适用于特定应用。我们并不是说研究人员、我们雇佣的标注者或我们的 API 客户是正确的偏好来源。有许多利益相关者需要考虑——训练模型的组织、使用模型开发产品的客户、这些产品的最终用户，以及可能收到语言模型使用直接或间接影响的更广泛人群。问题不仅仅在于使对齐过程更具参与性；不可能同时训练出一个符合每个人偏好的系统，或者让每个人都认同其中的权衡。

一条可能的前进道路是训练可以根据特定群体偏好进行条件调节的模型，或者可以轻松进行

微调或提示以代表不同群体的模型。然而，这些模式最终仍可能对社会产生更广泛的影响，需要做出许多艰难的决定，包括以谁的偏好为条件，以及如何确保所有群体都能得到代表，并能选择退出可能有害的过程。

请注意，虽然使用人类数据微调模型是部署 ML 系统时的常见做法，但这些努力的目的是获得在公司特定用例中表现良好的模型，而不是推进通用 ML 模型的对齐性。

5.3 局限性

方法。我们的 InstructGPT 模型的行为部分取决于从我们的承包商那里获得的人类反馈。有些标签任务依赖于价值判断，而价值判断可能会受到承包商的身份、信仰、文化背景和个人经历的影响。我们聘用了约 40 名承包商，根据他们在筛选测试中的表现来判断他们识别和回应敏感提示的能力，以及他们在有详细说明的标签任务中与研究人员的一致率（见附录 B）。我们之所以保持较小的承包商团队，是因为这有利于与较少的全职承包商进行高频次交流。然而，这群人显然不能代表使用我们部署的模型并受其影响的所有人群。举个简单的例子 我们的贴标人员主要讲英语，我们的数据几乎全部由英语指令组成。

我们还有许多方法可以改进我们的数据收集设置。例如，出于成本考虑，大多数比较只由一个承包商标注。对示例进行多次标注有助于确定承包商意见不一致的地方，从而确定单一模型不可能与所有承包商一致的地方。在存在分歧的情况下，按照平均标注者的偏好进行标注可能并不可取。例如，在生成对少数群体影响过大的文本时，我们可能希望更多考虑属于该群体的标注者的偏好。

模型。我们的模型既不完全一致，也不完全安全；它们仍然会产生有危害性或有偏见的输出，捏造事实，并在没有明确提示的情况下产生性和暴力内容。它们还可能无法根据某些输入生成合理的输出；我们在图 9 中展示了这方面的一些例子。

也许我们的模型最大的局限性在于，在大多数情况下，它们会遵循用户的指令，即使这在现实世界中可能会导致伤害。例如，当提示要求模型具有最大偏差时，InstructGPT 产生的有危害性的输出比同等大小的 GPT-3 模型产生更多的有危害性的输出。我们将在以下章节讨论可能的缓解措施。

5.4 开放式问题

这项工作是利用对齐技术对语言模型进行微调的第一步。广泛的指令。要进一步使语言模型

的行为与人们实际希望它们做的事情相一致，还有许多问题有待探索。

可以尝试许多方法来进一步降低模型产生有偏见或其他有害输出的倾向。例如，可以使用对抗设置，让标注者发现模型的最坏情况行为，然后将其标注并添加到数据集中（Dinan 等，2019b^[27]）。我们还可以将我们的方法与过滤预训练数据的方法结合起来（Ngo 等，2021年^[58]），无论是训练初始预训练模型，还是我们用于预训练混合方法的数据。同样，我们也可以将我们的方法与提高模型真实性的方法相结合，如 WebGPT(Nakano 等，2021 年^[55])。

在这项工作中，如果用户要求作出可能有害或不诚实的回应，我们就允许我们的模型生成这些输出。训练我们的模型使其在用户发出指令后仍能保持无害是很重要的，但也是很困难的，因为输出是否有害取决于其部署的环境；例如，作为数据增强管道的一部分，使用语言模型生成有毒输出可能是有益的。我们的技术还可以应用于制作拒绝某些用户指令的模型，我们计划在这项研究的后续迭代中对此进行探索。

让模型做我们想做的事与可转向性和可控性理论直接相关（Dathathri 等人，2019^[24]；Krause 等人，2020^[42]）。未来一条大有可为的道路是将 RLHF 与其他可控性方法相结合，例如使用控制代码（Keskar 等人，2019 年^[39]），或在推理时使用较小的模型修改采样程序（Dathathri 等人，2019 年^[24]）。

虽然我们主要关注的是 RLHF，但还有许多其他算法可用于在我们的演示和对比数据上训练策略，以获得更好的结果。例如，我们可以探索专家迭代法（Anthony 等人，2017 年^[3]；Silver 等人，2017 年^[70]），或者使用对比数据子集的更简单的行为克隆方法。我们还可以尝试约束优化方法（Achiam 等人，2017 年^[2]），在产生少量有害行为的条件下，最大化奖励模型的得分。

比较并不一定是提供对齐信号的最有效方式。例如，我们可以让标注者对模型的回答进行编辑，使其更加完善，或者用自然语言生成对模型回答的评论。为标注者设计向语言模型提供反馈的界面也有很多选择；这是一个有趣的人机交互问题。

而且可能会使某些任务更容易出现某些不良行为（如果预训练数据中存在这些行为的话）。

这是一个值得进一步研究的领域。另一种可能改进我们方法的修改方法是过滤预训练混合数据中的有毒内容（Ngo 等人，2021 年^[58]），或者用合成指令来增强这些数据。

正如（Gabriel，2020^[30]）所详细讨论的那样，与指令、意图、揭示偏好、理想偏好、利益和价值观保持一致之间存在着微妙的差异。（Gabriel，2020^[30]）主张采用基于原则的对齐方法：换句话说，就是要确定“公平的对齐原则，这些原则在人们的道德信念存在广泛差异的情况下仍能得到反思性认可”。在我们的论文中，为了简单起见，我们按照推断出的用户意图进行对齐，但在这一领域还需要更多的研究。事实上，一个最大的悬而未决的问题是，如何设计一个透明的、能切实代表受技术影响的人们的对齐过程，并以在许多群体中达成广泛共识的方式综合人们的价值观。我们在第 5.2 节讨论一些相关的考虑因素。

5.5 更广泛的影响

我们的目标是通过大型语言模型进行训练，让它们按照人类的意愿行事，从而提高这些模型的积极影响。默认情况下，语言模型会优化下一个单词的预测目标，而这只是我们希望这些模型做的事情的一个替代。我们的研究表明，我们的技术有望使语言模型更加有用、真实和无害。从长远来看，对齐失败可能会导致更严重的后果，尤其是当这些模型被部署在对安全至关重要的场合时。我们预计，随着模型规模的不断扩大，必须更加谨慎地确保它们与人类意图保持一致（Bostrom，2014^[13]）。

然而，让语言模型更好地跟踪用户意图也会让它们更容易被滥用。使用这些模型可能更容易生成令人信服的错误信息、仇恨或辱骂内容。

对齐技术并不是解决大型语言模型相关安全问题的灵丹妙药；相反，它们应被用作更广泛安全生态系统中的一种工具。除了有意滥用外，在许多领域中，大型语言模型的部署必须非常谨慎，或者根本就不应该部署。例如，医疗诊断、基于受保护特征的人员分类、确定信贷、就业或住房资格、生成政治广告以及执法等高风险领域。如果这些模型是开源的，那么在没有限制监管的情况下限制这些领域和其他领域的有害应用就会变得非常困难。另一方面，如果大型语言模型的使用权仅限于少数拥有训练这些模型所需资源的组织，那么大多数人就无法接触到最前沿的 ML 技术。另一种选择是，由一个组织拥有模型部署的端到端基础设施，并通过应用程序接口（API）进行访问。这样就可以实施安全协议，如用例限制（只允许将模型用于特定应用）、监控滥用和撤销滥用系统者的访问权限，以及限制速率以防止产生大

规模错误信息。然而，这样做的代价可能是透明度降低和权力更加集中，因为这需要应用程序接口提供商决定在each问题上的界限在哪里。

最后，正如第 5.2 节所讨论的，这些模式的对象问题极为重要，它将极大地影响这些模式的净影响是积极的还是消极的。

1174, 脑机接口

马斯克的脑机接口企业发展规划

(COPU 2023.09.20 讯)

马斯克 (Elon Musk) 的脑机接口初创公司 Neuralink, 2023 年 9 月 19 日宣布, 他们近期将开发脑机接口进行人体试验, 帮助瘫痪患者实现意念的打字功能, 未来还可帮助他瘫痪患者行走。

他们将 1024 根直径为 5 微米的导线“缝合”到患者的大脑灰质中, 与其周围的神经元连接, 对外与数字计算机连接, 将大脑中发射的高分辨率意念采样与计算机数字化代码之间进行转换, 通过脑机接口帮助瘫痪患者, 以意念指挥, 在机器帮助下实现瘫痪患者肢体动作, 也可帮助“盲人”重见光明。该公司目前正在招募临床患者。

1175, 基于 LiDAR 三维识别的球形转换器

Spherical Transformer for LiDAR-based 3D Recognition

作者: Xin Lai, Yukang Chen, Fanbin Lu, Jianhui Liu, Jiaya Jia

(北大李政恩摘录, 2023.3.22)

这是一篇来自 CVPR 2023 (计算机视觉与模式识别会议) 的一篇文章。它介绍了一种基于 LiDAR 的三维点云识别方法, 名为 SphereFormer。基于 LiDAR 的三维点云识别已经在各种应用中得到了好处。大多数当前方法在没有特别考虑 LiDAR 点分布的情况下, 会遭受信息不连贯和有限的感受野的问题, 尤其是对于稀疏的远程点。在这项工作中, 此文研究了 LiDAR 点的不同稀疏分布, 并提出了 SphereFormer 来直接将密集的近距离点的信息聚合到稀疏的远程点上。此文设计了径向窗口自注意力, 将空间分割为多个非重叠的窄而长的窗口。它克服了不

连贯的问题,并且平滑而显著地扩大了感受野,极大地提升了稀疏远程点的性能。此外,为了适应窄而长的窗口,此文提出了指数分割来产生细粒度的位置编码,并采用动态特征选择来增加模型的表征能力。值得注意的是,此文的方法在 nuScenes 和 SemanticKITTI 语义分割基准上分别排名第一, mIoU 分别为 81.9% 和 74.8%。此外,此方法在 nuScenes 物体检测基准中取得了第三名, NDS 为 72.8%, mAP 为 68.5%。

1176, AlphaFold 蛋白质模型的实时结构搜索和结构分类

Real-time structure search and structure classification for AlphaFold protein models

Tunde Aderinwale, Vijay Bharadwaj, Charles Christoffer 等

(北大李世平摘录, 2022. 4. 05)

去年,在蛋白质结构预测方面取得了突破,AlphaFold2 方法在建模精度方面有了实质性的提高。随着 AlphaFold2 软件的发布,AlphaFold2 对 21 个物种蛋白质的预测结构通过 AlphaFold 数据库公开提供。在这里,为了方便 AlphaFold2 模型的结构分析和应用,此文提供了基础设施 3D-AF-Surfer,它允许对 AlphaFold2 模型进行实时的基于结构的搜索。在 3D-af-surfer 中,结构用三维泽尼克描述符(3DZD)表示,这是三维形状的旋转不变的数学表示。我们开发了一个神经网络,以蛋白质的 3DZDs 作为输入,直接比较 3DZDs 更准确地检索相同折叠的蛋白质。利用 3D-AF-Surfer,此文报告了 AlphaFold2 模型的结构分类,并讨论了 AlphaFold2 模型的置信水平与内在无序区域之间的相关性。

1177, ViCo: 个性化文本到图像生成的细节保持视觉条件

ViCo: Detail-Preserving Visual Condition for Personalized Text-to-Image Generation

Shaozhe Hao, Kai Han, Shihao Zhao, Kwan-Yee K. Wong 等

(北大郜岱威摘录, 2023. 6. 2)

基于扩散模型的个性化文本到图像生成技术近年来得到了广泛的关注。给定包含新颖概念的少数图像(例如,一个独特的玩具),我们的目标是调整生成模型,

以捕捉精细的视觉细节的新概念，并生成逼真的图像后，文本条件。我们提出了一个插件的方法，名为 ViCo，快速和轻量级的个性化生成。具体来说，我们提出了一个图像注意力模块的条件下的扩散过程中的补丁明智的视觉语义。我们引入了一个基于注意力的对象掩模，几乎没有成本的注意力模块。此外，我们设计了一个简单的正则化的基础上的内在属性的文本图像的注意力地图，以减轻常见的过拟合退化。与许多现有的模型，我们的方法不微调的原始扩散模型的任何参数。这允许更灵活和可转移的模型部署。仅使用光参数训练（约 6% 的扩散 U-网络），我们的方法在定性和定量方面都比所有最先进的模型具有相当甚至更好的性能。

1178, 具有隐式逻辑约束的偶然性知识图谱的复杂查询回答

Complex Query Answering on Eventuality Knowledge Graph with Implicit Logical Constraints

Jiaxin Bai (Department of CSE HKUST), Xin Liu (Department of CSE HKUST), Weiqi Wang (Department of CSE HKUST), Chen Luo (Amazon.com Inc), Yangqiu Song (Department of CSE HKUST)

(北大董泽锦摘录, 2023. 5. 30)

使用深度学习方法查询不完全知识图谱 (KGs) 可以自然地利用推理和泛化能力来学习推断更好的答案。传统的神经复杂查询回答 (CQA) 方法主要用于以实体为中心的 KGs，然而，在现实世界中，我们还需要对事件、状态和活动 (即可能性或情况) 进行逻辑推理，以将学习系统从系统 I 推进到系统 II，正如 Yoshua Bengio 所提出的那样。从以偶然性为中心的 KG (EVKG) 进行逻辑查询自然可以为这种直观和逻辑推理提供参考。因此，在本文中，我们提出了一个新的框架，利用神经方法来回答基于 EVKG 的复杂逻辑查询，该框架不仅可以满足传统的一阶逻辑约束，而且可以满足关于其发生和顺序的可能性的隐式逻辑约束。例如，如果我们知道“食物是坏的”发生在“PersonX 添加酱油”之前，那么“PersonX 添加酱油”就不太可能是“食物是坏的”的原因，因为隐含的时间约束。为了促进 evkg 的一致性推理，我们提出了复杂偶然性查询回答 (Complex eventuality Query answer, CEQA)，这是 CQA 的一个更严格的定义，它考虑了控制偶然性时间顺序

和发生的隐含逻辑约束。通过这种方式，我们建议利用定理证明来构建基准数据集，以确保答案满足隐式逻辑约束。我们还提出了一种记忆增强查询编码 (MEQE) 方法，以显著提高最先进的神经查询编码器在 CEQA 任务上的性能。

1179, 使用 ChatGPT 进行列类型注释

Column Type Annotation using ChatGPT

Keti Korini, Christian Bizer

(北大谈家林摘录, 2023.6.1)

列类型注释是将关系表的列与每个列中包含的值的语义类型进行注释的任务。在数据湖的数据搜索和数据集成中，列类型注释是一个重要的预处理步骤。目前的列类型注释方法要么依赖于将表列与知识图的属性进行匹配，要么对预训练语言模型（如 BERT）进行微调以进行列类型注释。在这项工作中，我们采用了一种不同的方法，探索使用 ChatGPT 进行列类型注释。我们在零样本和少样本的设置中评估了不同的提示设计，并尝试向模型提供任务定义和详细的指导。我们进一步实现了一个两步骤的表格注释流程，首先确定表中描述的实体的类别，然后根据这个类别，使用只涉及整体词汇表的相关子集，要求 ChatGPT 进行列类型注释。使用指导和两步骤流程，ChatGPT 在零样本和一次样本设置中达到了超过 85% 的 F1 分数。要达到类似的 F1 分数，RoBERTa 模型需要使用 300 个示例进行微调。这个比较表明，在没有或仅有少量任务特定示例的情况下，ChatGPT 能够为列类型注释任务提供有竞争力的结果。

1180, 使用多数据因果发现为机器学习应用选择鲁棒特征

SELECTING ROBUST FEATURES FOR MACHINE LEARNING APPLICATIONS USING MULTIDATA CAUSAL DISCOVERY

Saranya Ganesh S., Tom Beucler, Frederick Iat-Hin Tam, Milton S. Gomez

(瑞士沃州洛桑大学地球表面动力学研究所); Jakob Runge (数据科学研究所德国航空航天中心 (DLR) 耶拿, 德国图林根; 德国柏林技术大学电气工程与计算机科学学院); Andreas Gerhardus (数据科学研究所德国航空航天中心 (DLR) 耶拿, 德国图林根)

稳健的特征选择对于创建可靠和可解释的机器学习 (ML) 模型至关重要。在领域知识有限且潜在交互未知的情况下设计统计预测模型时, 选择最佳特征集通常很困难。为了缓解这个问题, 本文引入了一种多数据 (M) 因果特征选择方法, 该方法同时处理一组时间序列数据集并产生一组单一的因果驱动因素。这种方法使用在 Tigramite Python 包中实现的因果发现算法 PC1 或 PCMCI。这些算法利用条件独立性测试来推断因果图的部分。本文的因果特征选择方法在将剩余的因果特征作为输入传递给预测目标的 ML 模型 (多元线性回归, 随机森林) 之前, 过滤掉因果虚假链接。本文的框架将被应用于西太平洋热带气旋的统计强度预测。

1181, SHAP-IQ: Unified Approximation of any-order Shapley Interactions
Fabian Fumagalli (比勒费尔德大学), Maximilian Muschalik (慕尼黑大学, 慕尼黑机器学习中心), Patrick Kolpaczki (帕德博恩大学) 等

(北大严丽君摘录, 2023. 3. 2)

主要在可解释人工智能 (XAI) 研究中, Shapley 值 (SV) 被用于确定任何黑盒模型的特征重要性分数。Shapley 交互指数扩展了 Shapley 值, 以定义任意顺序的特征交互分数。定义一个独特的 Shapley 交互指数是一个开放的研究问题, 到目前为止, 已经提出了三个定义, 这些定义因其选择的公理而不同。此外, 每个定义都需要一个特定的近似技术。然而, 我们提出了 SHAPley 相互作用量化 (SHAP-IQ), 这是一个高效的基于抽样的近似器, 用于计算所有三个定义的 Shapley 相互作用, 以及满足线性、对称性和假象公理的所有其他定义。SHAP-IQ 是基于一种新的表示方法, 与现有的方法相比, 我们为其近似质量提供了理论保证, 并对点估计的方差进行了估计。对于 SV 的特殊情况, 我们的方法揭示了 SV 的一种新的表示方法, 并与无偏核 SHAP 相对应, 大大简化了计算。我们通过在多维合成模型中解释最先进的语言模型来说明其计算效率和效果。

1182, Categorical Foundations of Explainable AI

Pietro Barbiero (剑桥大学), Stefano Fioravanti (锡耶纳大学), Francesco Giannini (锡耶纳大学), Alberto Tonda (巴黎-萨克雷大学), Pietro Lió (剑桥大学), Elena Di Lavore (塔林理工大学)

(北大李因立摘录, 2023. 4. 27)

可解释的人工智能 (XAI) 旨在满足人类对安全可靠的人工智能系统的需求。然而, 许多调查强调了关键 XAI 概念的缺乏健全的数学形式化, 特别是包括尚缺乏精确定义的解释这个术语。为了弥合这一差距, 本文首次提出了关键 XAI 概念和过程的数学严谨定义, 采用了范畴论的充分形式化。我们展示了我们的范畴框架可以: (i) 对现有的学习方案和架构进行建模, (ii) 形式化定义解释’ 这个术语, (iii) 为 XAI 分类体系建立理论基础, 以及 (iv) 分析常常被忽视的解释方法方面的因素。因此, 我们的范畴框架促进了人工智能技术的道德和安全部署, 因为它代表了朝着可解释的人工智能的健全理论基础迈出的重要一步。

1183, 绿色神经拟态计算支持:电阻式随机存取存储器的机器学习引导微加工
Supporting Green Neuromorphic Computing:Machine Learning Guided
Microfabrication for Resistive Random Access Memory
Abdi Yamil Vicenciodelmoral, Md Mehedi Hasan Tanim, Feng Zhao, Xinghui
Zhao,School of Engineering and Computer Science Washington State
University Vancouver

(北大周鑫摘录, 2023. 6. 9)

随着大数据和机器学习应用的日益普及, 需要一种更强大、更节能的方式来执行深度学习的工作流程。神经拟态芯片提供了一种潜在的解决方案, 因为它们试图模仿人类大脑中的神经元结构, 并在减少数量级的能量消耗和提高计算性能方面显示出巨大的潜力。然而, 神经芯片的制造昂贵。在本文中, 作者通过设计和开发机器学习引导的电阻性随机存取存储 (RRAM) 微加工工艺来解决这一挑战, RRAM 是神经拟态芯片的关键器件。实验结果表明, 该方法可以有效地预测不同工艺条件下 RRAM 器件的性能。

随着对更强大的计算系统的需求不断增加, 巨大的能源消耗和电子废物成为必须处理的重大副作用。同时解决这两个问题的一个潜在解决方案是“类脑”和“绿色”神经形态计算, 具有节能操作和可生物降解的处理方式。

然而, 由于微加工过程的复杂性, 开发这些硬件是极其困难的, 目前的实践是基于试错的, 既不高效也不可扩展。在本文中, 我们介绍了我们在利用机器学习技

术来帮助生物 RRAM 设计和开发方面的研究。我们的方法是基于一个学习任务转换，它在概率空间中对问题进行建模，并为任何新的生物 RRAM 设备生成一个全面的、细粒度的性能指标分布。这些方法弥合了最先进的差距，并为绿色神经拟态计算的下一代 RRAM 设备的更高效和有效的设计铺平了道路。目前还有一些想法正在实现当中。首先，作者将探索生物 RRAM 的其他性能指标，例如由制造过程引起的非理想影响，并利用基于机器学习的方法来检测这些影响。这对于在设备上支持神经网络至关重要。其次，将使用典型的深度学习基准来研究生物 RRAM 器件的学习能力。最后，在这些发现的基础上，将设计和开发工具来促进和评估神经形态计算的硬件/软件协同设计。

1184, 面向越野机器人环境感知与导航的神经形态计算

Towards the Neuromorphic Computing for Offroad Robot Environment Perception and Navigation

Zhenhua Yu, Peter R. N. Childs, Thrishantha Nanayakkara

(北大杨思恒摘录, 2023.06.09)

本文作者研究目标是通过设计一个用于越野机器人环境感知和导航的仿生锥形胡须神经形态计算（也称为储层计算）系统，明确弥合机器人车载硬件的高计算性能和低功耗之间的差距，该系统以机器人身体与其环境之间的交互为中心。在未知环境中执行任务的移动机器人需要穿越各种复杂的地形，它们必须能够可靠、快速地识别和表征这些地形，以避免陷入潜在的挑战或灾难性环境。为了解决这个问题，作者从老鼠和海豹等动物身上汲取了灵感，仅仅依靠胡须来感知周围的信息，并在黑暗狭窄的环境中生存。此外，作者观察了人类耳蜗，它可以分离不同频率的声音。基于这些见解，作者通过逐步探索基于物理胡须的储层计算来满足这一需求，以实现快速、经济高效的移动机器人环境感知和导航。这项研究可以帮助我们了解生物对应物的顺应性如何帮助机器人与环境动态互动，并与目前计算资源有限的机器人环境感知和导航方法（如火星）相比，提供了一种新的解决方案。

1185, 具有紧凑型神经形态核心的神经形态计算平台

A Neuromorphic Computing Platform with Compact Neuromorphic Core
Pujun Zhou, Shaogang Hu, State Key Laboratory of Electronic Thin Films and Integrated Devices, University of Electronic Science and Technology of China

(北大王世界摘录, 2022.10)

随着类脑计算的快速发展, 大规模的神经计算平台受到了广泛的关注。为了减少硬件开销, 构建大规模的神经计算平台, 本文提出了一个紧凑的神经形态核心模型。通过使用神经元复用技术和权重聚类算法, 这是一种紧凑而通用的神经形态计算核心, 并通过 588 个 LUT 集成了 1K 个神经元和 1M 个突触。在核心设计的基础上, 本文提出了一个大规模的神经形态系统。这个神经形态计算平台集成了 64 个神经形态核心和相关控制组件, 然后在 Xilinx 的 FPGA-Vertex-6 平台上成功部署了该平台。在本文工作中, 提出了一种紧凑、低带宽、可编程、通用和可扩展的神经形态核心, 每个核心集成了 1K 个神经元和 1M 个突触, 成本为 558LUT。本文在核心设计的基础上构建了一个紧凑的神经形态系统, 包括 64K 个神经元和 64M 个突触, 由于该系统具有良好的通用性, 可以在其中部署一个三层二进制尖峰神经网络 (SNN), 通过实验获得了 98.41% 的准确率。从理论上讲, 任何完全连接的网络都可以部署在这个系统上。未来, 可以尝试将算法在该平台上部署和测试卷积神经网络 (CNN)。

1186, 重新审视知识图谱推理基准

Revisiting Inferential Benchmarks for Knowledge Graph Completion
Shuwen Liu, Bernardo Cuenca Grau, Ian Horrocks, Egor V. Kostylev

(北大董泽锦, 2023 年 6 月 7 日)

知识图谱 (KG) 补全问题是用缺失的事实扩展不完整的知识图谱。KG 完成的机器学习方法的一个关键特征是它们学习推理模式的能力, 因此预测的事实是将这些模式应用于 KG 的结果。然而, 标准的完成基准并不适合于评估模型学习模式的能力, 因为这些基准的训练和测试集是给定 KG 的随机分割, 因此不能捕获推理模式的因果关系。我们提出了一种基于以下原则设计 KG 完成基准的新方法: 有一

套逻辑规则，因此缺失的事实是规则应用的结果；训练集包括前提匹配规则的前提和结论；测试集由将规则应用于训练集的结果组成；反例的设计是为了阻止模型学习规则集不包含的规则。我们使用我们的方法生成了几个基准，并评估了各种现有的 KG 完成系统。我们的研究结果为现有模型从不完整 KGs 中诱导推理模式的能力提供了新的见解。

1187, 因果规范化流动：从理论到实践

Causal normalizing flows: from theory to practice

Adrián Javaloy, Pablo Sánchez-Martín, Isabel Valera

(北大张凯迪摘录, 2023. 6. 8)

在这项工作中，我们深入探讨了使用正规化流进行因果推理的方法。具体而言，我们首先利用最近关于非线性独立成分分析（ICA）的结果表明，给定因果顺序，因果模型可以从观测数据中识别出来，并且可以使用自回归正规化流（NFs）进行恢复。其次，我们分析了用于因果正规化流的不同设计和学习选择，以捕捉潜在的因果数据生成过程。第三，我们描述了如何在因果正规化流中实现 do-operator，从而回答干预和反事实问题。最后，在我们的实验中，我们通过全面的消融研究验证了我们的设计和训练选择；将因果正规化流与其他逼近因果模型的方法进行比较；并经验证明因果正规化流可以用于解决现实世界的问题，其中混合离散-连续数据的存在和对因果图的部分了解是常态。此工作的代码可以在此 [https URL](https://github.com/AdriánJavaloy/causal-normalizing-flows) 中找到。

1188, 基于大型语言模型表示知识的机器人任务规划

Robot Task Planning Based on Large Language Model Representing Knowledge with Directed Graph Structures

Yue Zhen Sheng Bi Lu Xing-tong Pan Wei-qin Shi Hai-peng Chen Zi-rui

Fang Yi-shu

(北大谈家林摘录, 2023. 6. 5)

传统的机器人任务规划方法在处理高度无结构环境和复杂任务时面临挑战。我们提出了一种任务规划方法，将人类专业知识与大型语言模型（LLM）相结合，并

设计了一个具有更强表达能力的 LLM 提示模板 Think_Net_Prompt，用于表示结构化的专业知识。我们进一步提出了一种逐步分解任务并生成任务树的方法，以减少每个任务的规划量，并设计了一种解耦机器人任务规划的策略。通过将不同的规划实体进行划分并将任务与实际机器绑定过程分离，任务规划过程变得更加灵活。研究表明，我们的方法在处理指定的代码格式、理解任务和子任务之间的关系以及从文本描述中提取参数方面表现良好。然而，也存在一些问题，如任务逻辑处理的复杂度有限，部件数量和精确位置的模糊性。提高任务描述的准确性和认知结构可以带来一定的改进。

1189, 分割任何东西

Segment Anything

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, Ross Girshick

（北大李世平摘录，2023.4.5）

这是一篇来自 Meta 研究部门 Meta AI 的一篇论文，具有划时代的意义。本文介绍了 Segment Anything (SA) 项目：一个新的图像分割任务、模型和数据集。通过在数据收集循环中使用本文高效的模型，本文建立了迄今为止最大的分割数据集，其中包含超过 10,000 万张经过许可和尊重隐私的图像上的 10 亿个掩码。该模型经过设计和训练，可以接受提示，因此可以在新的图像分布和任务上进行零样本迁移。本文对其在多个任务上的能力进行评估，并发现其零样本性能令人印象深刻，通常与或甚至优于先前的全监督结果竞争。研究机构正在发布 Segment Anything Model (SAM) 和相应的数据集 (SA-1B)，包含 10 亿个掩码和 1100 万张图像，以促进计算机视觉基础模型的研究，详情请移步 <https://segment-anything.com>。

1190, 学习通过叙述在视频中放置教学文章

Learning to Ground Instructional Articles in Videos through Narrations
Effrosyni Mavroudi, Triantafyllos Afouras, Lorenzo Torresani

(北大李政恩摘录, 2023. 6. 7)

在本文中, 我们提出了一种方法, 本地化步骤的程序活动的叙述如何视频。为了解决大规模标记数据的稀缺性, 我们从包含各种程序任务的教学文章的语言知识库 (wikiHow) 中获取步骤描述。在没有任何形式的人工监督的情况下, 我们的模型通过匹配三种模式来学习在 how-to 视频中暂时接地程序文章的步骤: 框架、叙述和步骤描述。具体来说, 我们的方法通过融合来自两个不同途径的信息来将步骤与视频对齐: i) 步骤描述与帧的 { em 直接 } 对准, ii) 通过将步骤与叙述与视频对应性组合而获得的 { em 间接 } 对准。值得注意的是, 我们的方法通过利用订单信息一次执行文章中所有步骤的全局时间接地, 并且使用迭代细化和积极过滤的步骤伪标签进行训练。为了验证我们的模型, 我们引入了一个新的评估基准- HT-Step -通过手动注释 HowTo 100 M 的 124 小时子集 footnote {测试服务器可在 url {https: eval.ai web challenges challenge-page 2082} 访问。} 步骤来源于 wikiHow 文章。在此基准上的实验以及 CrossTask 上 zero-shot 评估表明, 我们的多模态对齐在几个基线和先前的工作中产生了显著的收益。最后, 我们表明, 我们的内部模块的叙事视频匹配优于一个很大的利润率的 HTM-Align 叙事视频对齐基准的最先进的。

1191, AlphaFold、人工智能 (AI) 和变构学

AlphaFold, Artificial Intelligence (AI), and Allostery

Ruth Nussinov, Mingzhen Zhang, Yonglan Liu, Hyunbun Jang

(北大郜岱威摘录, 2022. 8. 17)

AlphaFold 已经闯入了我们的生活。一个强大的算法, 强调生物序列数据和人工智能 (AI) 的力量。AlphaFold 增加了项目和研究方向。它一直在创建的数据库承诺了数不清的应用程序, 这些应用程序具有巨大的潜在影响, 目前仍难以猜测。人工智能方法可以彻底改变个性化治疗, 并带来更明智的临床试验。他们承诺在重塑和改进药物发现策略, 选择和优先考虑药物目标组合方面取得巨大飞跃。在这里, 我们简要概述了人工智能在结构生物学中的应用, 包括微生物-人类蛋白质-蛋白质相互作用的分子动力学模拟和预测。我们强调了由深度学习驱动的 AlphaFold 在蛋白质结构预测方面取得的进步及其对生命科学的强大影响。与此

同时，AlphaFold 并没有解决长达数十年的蛋白质折叠挑战，也没有识别折叠途径。AlphaFold 提供的模型没有捕捉到像挫折和变构这样的构象机制，它们根植于集合中，并由它们的动态分布控制。变构和信号传导是种群的特性。AlphaFold 也不会生成内在无序的蛋白质和区域的集合，而是通过它们的低结构概率来描述它们。由于 AlphaFold 产生的是单排结构，而不是构象组合，因此它无法阐明变构激活驱动热点突变的机制，也无法阐明变构耐药的机制。然而，通过捕获关键特征，深度学习技术可以使用单一预测构象作为生成多样化集成的基础。

1192, 面向不平衡数据的可解释的机器学习方法

Interpretable ML for Imbalanced Data

Damien A. Dablain (诺特丹大学, 美国), Colin Bellinger (加拿大国家研究委员会), Bartosz Krawczyk (弗吉尼亚联邦大学, 美国), David W. Aha (海军研究实验室 AI 应用研究中心, 美国), Nitesh V. Chawla (诺特丹大学, 美国)

在医学、自动驾驶和情报分析等高风险领域，深度学习模型正被越来越多地应用于不平衡的数据。不平衡的数据加剧了深度网络的黑箱性质，因为类之间的关系可能是高度倾斜和不明确的。这可能会减少模型用户的信任，并阻碍不平衡学习算法开发者的进展。现有的研究不平衡数据复杂性的方法是面向二元分类、浅层学习模型和低维数据的。此外，目前的可解释人工智能 (XAI) 技术主要侧重于将不透明的深度学习模型转换为更简单的模型 (如决策树) 或将特定实例的预测映射到输入，而不是研究全局数据属性和复杂度。因此，需要一个为现代深度神经网络量身定做的框架，该框架包含大型、高维、多类数据集，并能发现不平衡数据中常见的数据复杂性 (如类重叠、子概念和离群实例)。本文提出了一套技术，既可以被深度学习模型用户用来识别、可视化和理解类原型、子概念和离群实例；也可以被不平衡学习算法开发者用来检测对模型性能至关重要的特征和类典范。本文的框架还能识别位于类决策边界上的实例，这些实例可以携带高度的鉴别性信息。与许多现有的将模型决策映射到灰度像素位置的 XAI 技术不同，本文通过反向传播使用显著性来识别和聚集整个类别的图像色带。

1193, 《变好还是变坏: xAI 中反事实解释的方向性对用户行为的影响》

For Better or Worse: The Impact of Counterfactual Explanations' Directionality on User Behavior in xAI★

Ulrike Kuhl, Andr'e Artelt, Barbara Hammer (德国, 比勒费尔德, 比勒费尔德大学认知与机器人研究所; 比勒费尔德大学技术学院机器学习小组)

反事实解释 (CFE) 是可解释人工智能 (xAI) 中一种流行的方法, 强调改变模型输出所需的输入数据的变化。CFE 可以描述比事实状态更好的场景 (向上 CFE), 也可以描述比实际状态更差的场景 (向下 CFE)。然而, 在 xAI 中, CFE 对用户行为的定向性的潜在好处和缺点仍不清楚。当前的用户研究 (N=161) 比较了 CFE 方向性对参与者的行为和体验的影响, 这些参与者的任务是基于模型预测和 CFE 从自动化系统中提取新知识。结果表明, 与其他形式的反事实反馈相比, 向上的 CFE 提供了显著的性能优势。此外, 该研究强调了与向下的 CFE 或没有解释相比, 混合 CFE 提高用户性能的潜在好处。与性能结果一致, 用户对研究的明确了解详见 <https://github.com/ukuhl/DirectionalAlienZoo>

1194, 解释可解释性: 通过二阶可解释性深入了解深度学习

Explaining Explainability: Towards Deeper Actionable Insights into Deep Learning through Second-order Explainability

E. Zhixuan Zeng, Hayden Gunraj, Alexander Wong (滑铁卢大学, 加拿大), Sheldon Fernandez (DarwinAI 股份有限公司, 加拿大)

解释性在深度学习模型行为的全面理解中起着关键作用。这可以对模型的性能进行彻底验证, 确保其决策基于相关的视觉指标, 并且不偏向于训练数据中存在的无关模式。然而, 现有方法只提供了实例级的解释性, 需要对每个样本进行手动分析。这种手动审查耗时且容易受到人为偏见的影响。为了解决这个问题, 最近提出了“二阶可解释人工智能” (Second-Order Explainable AI, SOXAI) 的概念, 将可解释人工智能 (Explainable AI, XAI) 从实例级扩展到数据集级别。SOXAI 通过识别普遍概念, 自动分析定量解释与数据集偏差之间的关联。在这项工作中, 本文探索了利用深度神经网络行为的这种更高层次的解释来“解释可解释性”, 以获得可操作的见解。具体而言, 本文通过示例分类和分割案例首次证

明，根据 SOXAI 的可操作见解从训练集中排除无关概念可以提高模型的性能。

1195, Sight: 从高等教育成绩单中收集的关于学生洞察力的大型注释数据集
SIGHT: A Large Annotated Dataset on Student Insights Gathered from Higher Education Transcripts

Rose E. Wang, Pawan Wirawarn, Noah Goodman, Dorottya Demszky (斯坦福大学)
讲座对学生和教师来说都是一种学习体验。学生从教师那里学习主题材料，而教师从学生那里学习如何改进他们的教学。然而，在线学生反馈是非结构化的和丰富的，这使得教师学习和改进具有挑战性。本文朝着应对这一挑战迈出了一步。首先，本文贡献了一个数据集来研究这个问题：SIGHT 是从麻省理工学院开放式课程 (MIT OCW) YouTube 频道收集的 288 个数学讲座成绩单和 15,784 条评论的大型数据集。其次，本文开发了一个使用定性分析的分类反馈类型的量规。定性分析方法在发现特定领域的见解方面非常强大，但应用于大型数据源的成本很高。为了克服这一挑战，本文提出了一套最佳实践，使用大型语言模型 (LLM) 来廉价地对评论进行大规模分类。本文观察到模型和人类注释之间的惊人相关性：具有一致人类注释的类别 (评分者间可靠性, IRR) 也显示出较高的人类模型一致性 (0.7 \$)，而具有不太一致人类注释的类别 (0.7 \$-0.8 \$ IRR) 相应地显示出较低的人类模型一致性 (0.3 \$-0.5 \$)。这些技术从数千条评论中发现有用的学生反馈，每条评论的成本约为 0.002 美元。最后，本文讨论了令人兴奋的未来方向，使用在线学生反馈和改进自动注释技术的定性研究。

1196, AlphaFold 在蛋白质结构预测之外的应用

Applications of AlphaFold beyond Protein Structure Prediction

Yuan Zhang, Feng Pan, Xiuwen Liu, Jinfeng Zhang (佛罗里达大学, 美国), Peizhao Li, Hongfu Liu, Pengyu Hong (布兰迪斯大学, 美国)

通过 DeepMind 的 AlphaFold 准确预测天然蛋白质序列的结构无疑是 21 世纪生物学领域最伟大的突破之一。对于可能不稳定的设计或工程序列，预测稳定性及其结构是必不可少的，因为不稳定的结构将无法正常运行。我们发现实验测量的点突变稳定性变化与 AlphaFold 产生的置信度分数相关性很差。但是，可以

使用从 AlphaFold 学习的表示中提取的特征准确预测稳定性变化，这表明 AlphaFold 对设计或工程序列的普遍性比以前认为的要大。然后，我们使用 AlphaFold 来验证我们之前开发的蛋白质设计方法 ProDCoNN，该方法设计序列以折叠到目标蛋白质结构，仅给出目标蛋白质的骨架结构信息。我们展示了 ProDCoNN 能够设计折叠成非常接近目标结构的结构的序列。通过结合修改后的 ProDCoNN、AlphaFold 和顺序蒙特卡洛，我们设计了一个新的框架来估计蛋白质结构的可设计性。蛋白质结构的可设计性定义为编码蛋白质结构的序列数量，是蛋白质功能稳健性的指标。我们首次估算了具有 134 个残基的 FLT3 配体 A 链 (PDB ID: 1ETE) 的真实蛋白质结构的可设计性，为 $3.12 \pm 2.14E85$ 。

1197, MetaFormer 才是你需要的

MetaFormer Is Actually What You Need for Vision

Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, Shuicheng Yan

这是一篇来自 Shuicheng Yan 大佬的文章，是 2022CVPR 会议的 oral 文章，和 transformer 的文章 attention is what you need 形成对称，探讨 transformer 中真正带来效果的机理。研究表明，Transformers 中的注意力机制模块可以被空间 MLP 替代，并且替代后的模型性能依然很好。基于这一观察，本文假设 Transformers 的整体架构，而不是具体的令牌混合模块，对模型的性能更为重要。为了验证这一假设，本文故意将 Transformers 中的注意力模块替换为一个非常简单的空间汇聚算子，只进行基本的令牌 (token) 混合。令人惊讶的是，本文观察到得到的模型、PoolFormer 在多个计算机视觉任务上取得了很有竞争力的表现。例如，在 ImageNet-1K 数据集上，PoolFormer 达到了 82.1% 的 top-1 准确率，比经过精细调整的 Vision Transformer/MLP-like 基线模型 DeiT-B/ResMLP-B24 的准确率分别高出 0.3%/1.1%，同时参数数量减少了 35%/52%，计算量减少了 50%/62%。PoolFormer 的有效性验证了本文的假设，并促使本文提出了“MetaFormer”的概念，这是从 Transformers 中抽象出的一个通用架构，而不指定令牌混合方式。通过广泛的实验证明，本文认为 MetaFormer 是最近在视觉任务上取得卓越结果的 Transformer 和 MLP-like 模型中的关键机理。这项

工作呼吁未来的研究专注于改进 MetaFormer，而不是仅仅关注令牌混合模块。此外，本文提出的 PoolFormer 可以作为未来 MetaFormer 架构设计的起点基准。

1198, 使用神经形态计算的延时光子库分析

Analysis of time-delay photonic reservoirs for neuromorphic computing
P. Dmitriev, L. Di Lauro, B. Fischer, A. Aadhi, M. Chemnitz, E. Viktorov, A. Kovalev, and R. Morandotti. ITMO (俄罗斯圣彼得堡国立信息技术、机械学与光学研究型大学)

(北大周鑫摘录, 2022.5)

人工神经网络(ann)和机器学习(ML)的深度学习已经从图像识别扩展到翻译语言,生成人类无法区分的真实语音。光子处理器的性能明显优于基本上依赖于互连的电子系统。神经拟态计算方法,包括光子学提供的方法,不依赖于神经元之间的串行通信和传统的冯·诺伊曼方法来计算神经元的响应。光子元件的使用使我们能够摆脱传统的计算装置,并利用光子的内在并行性。传统的人工神经网络结构分为前馈网络(fnn)和循环网络(rnn)。由于光纤网络的数据传播是单向的,限制了其处理能力。另一方面,由于反馈回路的存在,rnn 具有所谓的“褪色记忆”。水库计算(RC)最初是一个基于 rnn 的框架,因此适合于时间/顺序信息处理。在 RC 的方法中,输入数据被映射到高维相空间,其中复杂和不可分离的 ML 问题可以通过更简单的(线性)训练算法(如线性或脊回归)轻松解决。实际上,RC 的主要优点是输入权值和储层内循环连接的权值保持不变,而只训练输出权值。在这项工作中,作者通过数值研究了延时光子库的稳定性和分支,特别是那些包含单个非线性节点和延迟反馈回路(即实现衰落响应)的光子库。针对机器学习应用,作者考虑了两个基本的系统优化参数,即反馈强度和增益功率。我们使用了 ML 标准基准测试,例如具有渐进预测步骤的 Mackey-Glass 时间序列,来评估我们系统的处理能力。此外,我们评估了系统在 Hopf 分岔附近的记忆容量(线性和非线性)(见图 1 中的图表),以研究它们与衰落记忆、非线性和系统性能的联系。存储容量的计算。观察到,当系统参数位于分岔点附近时,获得的最佳预测值为 $\sim 0.1\%$ 的归一化均方误差(NMSE)。这个特定基准测试任务的最佳配置确实存在于参数区域,在该区域中,系统具有较高的非线性行为,并且具有足够大的非线性

内存容量。

我们对非线性时滞光子储层的数值研究支持其基本特征(如非线性和存储容量)的内在联系, 以提供最佳的计算性能。该理论结果允许通过调整某些控制参数来优化光子库, 并为基于非线性集成光子学的新的智能应用铺平道路, 在机器学习和神经形态计算领域。

1199, 利用时间到事件裕度传播的 AER 神经形态计算

Neuromorphic Computing with AER using Time-to-Event-Margin Propagation
Madhuvanthi Srivatsav R, Shantanu Chakrabartty, Chetan Singh Thakur (印度科学院电子系统工程系、美国圣路易斯华盛顿大学电气工程与系统工程系)

(北大杨思恒摘录, 2023. 4. 27)

地址事件表示(AER)是一种尖峰路由协议, 允许将神经形态和尖峰神经网络(SNN)架构扩展到与数字神经网络架构相当的规模。然而, 在传统的神经形态架构中, AER 协议以及通常的任何虚拟互连在计算中仅起被动作用, 即仅用于路由尖峰和事件。在本文中, 作者展示了如何使用其提出的称为时间到事件裕度传播(TEMP)的技术, 将 AER 协议本身固有的延迟、触发和排序等因果时间基元用于可扩展的神经形态计算。所提出的基于 TEMP 的 AER 架构是完全异步的, 并且依赖于用于存储器和计算的互连延迟, 而不是传统的和局部的乘法和累加(MAC)操作。本文表明, TEMP 神经网络中基于时间的编码产生了一种时空表示, 可以对大量的判别模式进行编码。作为概念验证, 本文表明, 经过训练的基于 TEMP 的卷积神经网络(CNN)可以在 MNIST 数据集上证明大于 99% 的准确性。总的来说, 作者的工作是一种受生物学启发的计算范式, 为神经形态计算领域带来了新的研究维度。

1200, 神经形态计算硬件的最新发展趋势: 英特尔的神经形态系统视角

Recent Trend of Neuromorphic Computing Hardware: Intel' s Neuromorphic System Perspective

Yoon Seok Yang and Yongtae Kim Neuromorphic Computing Lab, Intel Labs, Intel Corporation, Santa Clara(神经形态计算实验室, 英特尔实验室, 英特

尔公司, 加利福尼亚州圣克拉拉), CA 95054, USA School of Computer Science and Engineering, Kyungpook National University, Daegu(庆浦国立大学美国计算机科学与工程学院, 大邱) Republic of Korea

(北大王世界摘录, 2020.10)

神经形态计算已被研究用于实现受人脑启发的功能, 如低功耗、细粒度并行处理和实时学习, 这些功能超出了标准冯·诺依曼处理器的限制。本文介绍了英特尔的 Loihi 神经形态研究芯片及其硬件系统, 并了解了它们在实际研究领域中的应用。目前在人工智能系统中, 为了提高神经网络系统的性能, 需要许多处理器, 并相应地消耗大量能量, 这些缺点使得在边缘设备中实现人工智能变得困难。为了克服这些限制, 模拟人脑的神经形态计算已被研究为人工智能的替代计算系统。神经形态计算的优点是支持基于事件的低能耗、可扩展的并行处理, 以及将内存和计算结合在一个单元中, 如神经元。本文就基于神经形态硬件和神经形态系统提出了相应的神经形态方面的研究。Loihi 是英特尔实验室于 2018 年开发的一款神经形态研究芯片。该芯片提供事件驱动细粒度并行计算, 这是一个支持异步 SNN 的可编程片上学习引擎。Loihi 可以执行某些要求苛刻的工作负载, 其速度比传统处理器快 1000 倍, 效率比传统处理器高 10000 倍。Loihi 由 130000 个神经元、1.3 亿个突触和三个 Lakemont×86 核心组成, 用于配置和管理神经网络。它是在英特尔的 14 纳米 FinFET 工艺上制造的。它还包括一个由异步网状协议实现的片外通信接口, 可扩展到 16384 个 Loihi 芯片。该芯片的性能十分优秀, 能够用于开发先进的神经形态算法和应用程序, 在神经形态领域拥有很大的发展潜力。

1201, MolCAP: 分子化学反应性预训练和提示微调增强分子表征学习

MolCAP: Molecular Chemical reActivity pretraining and prompted-finetuning enhanced molecular representation learning

Yu Wang, JingJie Zhang, Junru Jin, Leyi Wei (上海交通大学)

(北大张凯迪摘录, 2023 年 6 月 13 日)

分子表征学习 (MRL) 是药物的基本任务发现。然而, 以前的深度学习 (DL) 方法过于关注通过掩模主导的预训练学习稳健的内分子表示 框架, 忽略了丰富的

化学反应性分子关系已被证明是各种分子性质的决定因素预测任务。在这里，我们介绍了 Mo1CAP 来促进 MRL，一种图预训练基于化学反应性（IMR）知识的变压器，提示微调。结果表明，Mo1CAP 优于基于以下条件的比较方法 传统的分子预训练框架，在 13 个公开可用的分子中 跨各种生物医学任务的数据集。由 Mo1CAP 提示，甚至是基本的图神经网络能够实现令人惊讶的性能优于以前的模型，表明应用前景广阔 MRL 的反应性信息。此外，手动设计的分子模板有可能发现数据集偏差。总而言之，我们希望我们的 Mo1CAP 能够对药物的整个过程获得更多有意义的化学见解发现。

1202, 在低度并发环境中的确定性和工作效率较高的批处理动态树

Deterministic and Work-Efficient Parallel Batch-Dynamic Trees in Low Span
Daniel Anderson (卡耐基梅隆大学), Guy E. Blelloch (卡耐基梅隆大学)
(北大谈家林摘录, 2023 年 6 月 14 日)

动态树是动态图算法的基本构建模块，其研究可追溯到 Sleator 和 Tarjan 的开创性工作 [STOC' 81, (1981), pp. 114-122]。该问题是在在线边插入和删除的同时维护一棵树，并回答关于该树的查询，例如路径上的最大权重等。在并行批处理动态设置中，目标是以低（多项式对数 n ）的度并发处理批量边更新。已知两种工作效率高的算法，分别是 Tseng 等人的批处理并行 Euler Tour Trees [ALENEX' 19, (2019), pp. 92-106] 和 Acar 等人的并行 Rake-Compress (RC) Trees [ESA' 20, (2020), pp. 2:1-2:23]。然而，这两种算法都是随机的，并且期望上具有工作效率。因此，使用这些数据结构的几个下游结果（事实上据我们所知，所有已知的工作效率高的并行批处理动态图算法）也是随机的。在本研究中，我们提出了该问题的第一个确定性工作效率解决方案。我们的算法以最坏情况下的确定性方式维护一棵动态并行树压缩，并处理批量的 k 个边更新，工作复杂度为 $O(k \log(1 + n/k))$ ，度并发为 $O(\log n \log(c) k)$ ，其中 c 为任意常数。这使得我们能够确定性地实现具有最坏情况下的 $O(k \log(1 + n/k))$ 工作更新和查询的并行批处理 RC-Trees。我们用于获得给定度并发上界的技术也可应用于改进现有的随机变体算法，将其度并发从 $O(\log n \log^* n)$ 提高到 $O(\log n)$ 。

1203. 统一大型语言模型和知识图谱:Roadmap

Unifying Large Language Models and Knowledge Graphs: A Roadmap

Shirui Pan (电气电子工程师学会高级成员), Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, Xindong Wu (电气电子工程师学会研究员)

(北大董泽锦摘录, 2023年6月14日)

大型语言模型(LLMs), 如 ChatGPT 和 GPT4, 由于其突现能力和泛化能力, 正在自然语言处理和人工智能领域掀起新的浪潮。然而, 大型语言模型是黑盒模型, 往往不能捕获和访问事实知识。相反, 知识图谱(Knowledge Graphs, KGs), 比如维基百科和华普, 是结构化的知识模型, 明确地存储了丰富的事实知识。KGs 可以通过为推理和可解释性提供外部知识来增强大型语言模型。同时, KGs 本身就难以构建和演化, 这对现有的 KGs 生成新事实和表示未知知识的方法提出了挑战。因此, 将 LLM 和 KGs 统一在一起, 同时发挥各自的优势是互补的。在本文中, 我们提出了 LLMs 和 KGs 统一的前瞻性路线图。我们的路线图包括三个一般框架, 即

1)KG-enhanced LLMs, 在 LLMs 的预训练和推理阶段纳入 KGs, 或用于增强 LLMs 所学知识的理解;

2)LLM-augmented KG, 利用 LLM 完成不同的 KG 任务, 如嵌入、补全、构建、图形到文本生成和问题回答;

3)协同 LLMs + KGs, LLMs 和 KGs 发挥平等的作用, 以互利的方式工作, 增强 LLMs 和 KGs 在数据和知识驱动下的双向推理能力。我们在路线图中回顾和总结了这三个框架内的现有工作, 并指出了它们未来的研究方向。

1204, 《反事实解释者团》

Ensemble of Counterfactual Explainers

Riccardo Guidotti, Salvatore Ruggieri

(University of Pisa, Italy, 意大利比萨大学)

在可解释人工智能(XAI)中, 已经提出了几个反事实解释器, 每个解释器都关注反事实实例的一些理想性质: 最小性、可操作性、稳定性、多样性、合理性和判别力。我们提出了一个反事实解释器的集合, 将仅提供此类属性的子集的弱解

释器提升为覆盖所有属性的强大方法。该集成在实例和特征的样本上运行弱解释器，并通过利用多样性驱动的选择函数来组合它们的结果。该方法是模型不可知的，通过基于自动编码器的包装方法，它也是数据不可知的。

1205, 为自动驾驶场景解读获取定性可解释图形

Acquiring Qualitative Explainable Graphs for Automated Driving Scene Interpretation 摘录, 2023. 8. 24

Nassim Belmecheri (挪威 SIMULA 研究实验室), Arnaud Gotlieb (挪威 SIMULA 研究实验室), Nadjib Lazaar (法国蒙彼利埃大学), Helge Spieker (挪威 SIMULA 研究实验室)

自动驾驶技术 (AD) 的未来在于开发稳健、公平和可解释的人工智能方法。根据要求, 自动驾驶车辆必须能够向驾驶员和车内乘客、行人和其他易受伤害的道路使用者解释其决策, 并在发生事故时向外部审计人员进行解释。然而, 目前大多数可解释方法仍依赖于对多个传感器捕捉到的自动驾驶场景表征进行定量分析。本文提出了一种新颖的 AD 场景表示法, 称为定性可解释图 (QXG), 专用于长期场景的定性时空推理。该图的构建利用了最新的定性约束获取范式。我们在 NuScenes (一个开放的真实世界多模式数据集) 上的实验结果表明, 由 40 个帧组成的注意力缺失场景的定性可解释图可以在空间存储中实时轻量计算, 这使其成为改进注意力缺失感知和控制过程并提高其可信度的潜在有趣工具。

1206, 可解释人工智能文献的趋势

Trends in Explainable AI (XAI) Literature 摘录, 2023. 1. 13

Alon Jacovi (Bar Ilan University, 巴伊兰大学)

XAI 文献分散, 无论是在术语还是在出版场所上, 但近年来社区聚焦了一些关键词, 使得能够更可靠地自动发现论文。我们使用 SemanticScholar API 进行关键词搜索和手动策划, 收集了一个格式良好且相对全面的 XAI 论文集合, 共包括 5199 篇论文, 可以在此 <https> 网址上找到。我们利用这个集合来澄清和可视化文献的规模和范围趋势、引用趋势、跨领域趋势以及合作趋势。总体而言, XAI 正在变得越来越跨学科, 各种 (非计算机科学) 科学领域的论文相对增长, 跨领

域合作作者数量增加，跨领域引用活动增加。此外，该集合还可以用作文献发现引擎，通过检索根据特定约束条件引用的 XAI 文献（例如，在其领域之外具有影响力的论文，或对非 XAI 研究具有影响力的论文）。

1207, Swin Transformer: 使用平移窗口的分层 Transformer

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows
Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, Baining Guo

本文介绍了一种新的视觉 Transformer，称为 Swin Transformer，它能够作为计算机视觉的通用骨干。从语言到视觉的转换面临着许多挑战，例如视觉实体的尺度变化大，图像中的像素分辨率与文本中的单词相比较。为了解决这些差异，我们提出了一种分层 Transformer，其表示是通过移位窗口计算的。移位窗口方案通过将自注意力计算限制在非重叠的局部窗口中，同时允许跨窗口连接，从而带来更高的效率。这种分层架构具有在各种尺度上建模的灵活性，并且与图像大小具有线性计算复杂度。Swin Transformer 的这些特性使其与广泛的视觉任务兼容，包括图像分类（ImageNet-1K 上的 87.3 top-1 准确度）和密集预测任务，例如目标检测（COCO 测试集上的 58.7 box AP 和 51.1 mask AP）以及语义分割（ADE20K val 上的 53.5 mIoU）。它的性能大大超过了以前的最新技术，COCO 上的 +2.7 box AP 和 +2.6 mask AP，ADE20K 上的 +3.2 mIoU，展示了基于 Transformer 的模型作为视觉骨干的潜力。分层设计和移位窗口方法对于所有 MLP 架构也是有益的。

1208, 通过任务相关的多视图对比学习进行听觉注意力解码

Auditory Attention Decoding with Task-Related Multi-View Contrastive Learning 摘录, 2023. 8. 8

Xiaoyu Chen, Changde Du, Qiongyi Zhou, Huiguang He (China)

在鸡尾酒会等场景中，人脑很容易专注于一位发言者并抑制其他发言者。最近，研究人员发现听觉注意力可以从脑电图（EEG）数据中解码。然而，大多数现有的深度学习方法很难使用不同视图的先验知识（即参与语音和脑电图是与任务相

关的视图)并提取不令人满意的表示。受布罗德本特过滤模型的启发,我们在多视图范式中解码听觉注意力,并利用缺失的视图提取最相关和最重要的信息。具体来说,我们提出了一种基于多视图 VAE 和任务相关多视图对比(TMC)学习的听觉注意解码(AAD)方法。在多视图 VAE 中采用 TMC 学习可以利用缺失的视图将不同视图的先验知识积累到表示的融合中,并提取近似的任务相关表示。我们在两个流行的 AAD 数据集上检查我们的方法,并通过将其与最先进的方法进行比较来证明我们的方法的优越性。

1209, 走出样本陷阱: 用成对距离估计快速准确地估计认知不确定性

Escaping the Sample Trap: Fast and Accurate Epistemic Uncertainty Estimation with Pairwise-Distance Estimators 摘录, 2023. 8. 28

Lucas Berry, David Meger (麦吉尔大学)

这项工作介绍了一种新的方法,使用成对距离估计(PaiDEs)的集成模型的认知不确定性估计。这些估计器利用模型分量之间的成对距离来建立熵的界限,并使用所述界限作为基于信息的准则的估计。与最近用于认知不确定性估计的深度学习方法(依赖于基于样本的蒙特卡罗估计器)不同,PaiDE 能够在更大的空间(高达 $100 \times$)上更快地估计认知不确定性高达 $100 \times$,并且在更高的维度上表现得更准确。为了验证本文的方法,本文进行了一系列通常用于评估认知不确定性估计的实验: 1D 正弦数据, Pendulum-v0、Hopper-v2、Ant-v2 和 Humanoid-v2。对于每一个实验设置,主动学习的框架被应用到展示的优势 PaiDEs 的认知不确定性估计。

1210, 绿色神经拟态计算支持:电阻式随机存取存储器的机器学习引导微加工

Supporting Green Neuromorphic Computing: Machine Learning Guided Microfabrication for Resistive Random Access Memory 摘录, 2023. 6. 9

Abdi Yamil Vicenciodelmoral, Md Mehedi Hasan Tanim, Feng Zhao, Xinghui Zhao, School of Engineering and Computer Science Washington State University Vancouver

随着大数据和机器学习应用的日益普及,需要一种更强大、更节能的方式来执行

深度学习的工作流程。神经拟态芯片提供了一种潜在的解决方案，因为它们试图模仿人类大脑中的神经元结构，并在减少数量级的能量消耗和提高计算性能方面显示出巨大的潜力。然而，神经芯片的制造昂贵。在本文中，作者通过设计和开发机器学习引导的电阻性随机存取存储 (RRAM) 微加工工艺来解决这一挑战，RRAM 是神经拟态芯片的关键器件。实验结果表明，该方法可以有效地预测不同工艺条件下 RRAM 器件的性能。

随着对更强大的计算系统的需求不断增加，巨大的能源消耗和电子废物成为必须处理的重大副作用。同时解决这两个问题的一个潜在解决方案是“类脑”和“绿色”神经形态计算，具有节能操作和可生物降解的处理方式。

然而，由于微加工过程的复杂性，开发这些硬件是极其困难的，目前的实践是基于试错的，既不高效也不可扩展。在本文中，我们介绍了我们在利用机器学习技术来帮助生物 RRAM 设计和开发方面的研究。我们的方法是基于一个学习任务转换，它在概率空间中对问题进行建模，并为任何新的生物 RRAM 设备生成一个全面的、细粒度的性能指标分布。这些方法弥合了最先进的差距，并为绿色神经拟态计算的下一代 RRAM 设备的更高效和有效的设计铺平了道路。目前还有一些想法正在实现当中。首先，作者将探索生物 RRAM 的其他性能指标，例如由制造过程引起的非理想影响，并利用基于机器学习的方法来检测这些影响。这对于在设备上支持神经网络至关重要。其次，将使用典型的深度学习基准来研究生物 RRAM 器件的学习能力。最后，在这些发现的基础上，将设计和开发工具来促进和评估神经形态计算的硬件/软件协同设计。

1211，面向越野机器人环境感知与导航的神经形态计算

Towards the Neuromorphic Computing for Offroad Robot Environment Perception and Navigation 摘录，2023.06.09

Zhenhua Yu, Peter R. N. Childs, Thrishantha Nanayakkara

本文作者研究目标是通过设计一个用于越野机器人环境感知和导航的仿生锥形胡须神经形态计算（也称为储层计算）系统，明确弥合机器人车载硬件的高计算性能和低功耗之间的差距，该系统以机器人身体与其环境之间的交互为中心。在未知环境中执行任务的移动机器人需要穿越各种复杂的地形，它们必须能够可靠、

快速地识别和表征这些地形，以避免陷入潜在的挑战或灾难性环境。为了解决这个问题，作者从老鼠和海豹等动物身上汲取了灵感，仅仅依靠胡须来感知周围的信息，并在黑暗狭窄的环境中生存。此外，作者观察了人类耳蜗，它可以分离不同频率的声音。基于这些见解，作者通过逐步探索基于物理胡须的储层计算来满足这一需求，以实现快速、经济高效的移动机器人环境感知和导航。这项研究可以帮助我们了解生物对应物的顺应性如何帮助机器人与环境动态互动，并与目前计算资源有限的机器人环境感知和导航方法（如火星）相比，提供了一种新的解决方案。

1212, 具有紧凑型神经形态核心的神经形态计算平台

A Neuromorphic Computing Platform with Compact Neuromorphic Core 摘录,
2022.10

Pujun Zhou, Shaogang Hu, State Key Laboratory of Electronic Thin Films and Integrated Devices, University of Electronic Science and Technology of China

随着类脑计算的快速发展，大规模的神经计算平台受到了广泛的关注。为了减少硬件开销，构建大规模的神经计算平台，本文提出了一个紧凑的神经形态核心模型。通过使用神经元复用技术和权重聚类算法，这是一种紧凑而通用的神经形态计算核心，并通过 588 个 LUT 集成了 1K 个神经元和 1M 个突触。在核心设计的基础上，本文提出了一个大规模的神经形态系统。这个神经形态计算平台集成了 64 个神经形态核心和相关控制组件，然后在 Xilinx 的 FPGA-Vertex-6 平台上成功部署了该平台。在本文工作中，提出了一种紧凑、低带宽、可编程、通用和可扩展的神经形态核心，每个核心集成了 1K 个神经元和 1M 个突触，成本为 558LUT。本文在核心设计的基础上构建了一个紧凑的神经形态系统，包括 64K 个神经元和 64M 个突触，由于该系统具有良好的通用性，可以在其中部署一个三层二进制尖峰神经网络（SNN），通过实验获得了 98.41% 的准确率。从理论上讲，任何完全连接的网络都可以部署在这个系统上。未来，可以尝试将算法在该平台上部署和测试卷积神经网络（CNN）。

1213, StratMed: 低资源药物推荐的相关性分层

StratMed: Relevance Stratification for Low-resource Medication Recommendation 摘录, 2023 年 9 月 3 日

Xiang Li (中国科学院大学)

随着有限的医疗资源和不断增长的需求之间的不平衡日益加剧, 基于人工智能的临床任务变得至关重要。药物推荐作为一个子领域, 旨在将患者的纵向病史与医学知识结合起来, 帮助医生开出更安全、更准确的药物组合。现有方法忽略了医学数据固有的长尾分布, 缺乏头部和尾部数据的均衡表示, 导致模型性能不理想。为了应对这一挑战, 我们引入了 StratMed, 这是一个包含创新相关分层机制的模型。它协调了数据长尾分布中的差异, 并在药物组合的安全性和准确性之间取得了平衡。具体来说, 我们首先使用深度学习网络构建预训练方法来获得实体表示。之后, 我们设计了一个金字塔状的数据分层方法, 通过强化不受欢迎的实体的特征来获得更广义的实体关系。基于这种关系, 我们设计了两个图结构来表示同一层次的用药精度和安全性, 从而得到访问表示。最后, 患者的历史临床信息被拟合以生成针对当前健康状况的药物组合。在 MIMIC-III 数据集上的实验表明, 我们的方法在四个评估指标(包括安全性和准确性)上优于当前最先进的方法。

1214, 知识图谱工程中大语言模型评估的可扩展基准

Developing a Scalable Benchmark for Assessing Large Language Models in Knowledge Graph Engineering 摘录,

2023 年 8 月 31 日

Meyer(1, 2, 3), Frey(1, 2, 3), Junghanns(1, 2), Brei(1), Bulert(1),
Gründer-Fahrer(1, 2), Martin(1, 2)

1. 德国应用信息学研究所
2. 敏捷知识工程和语义网(AKSW)
3. 德国莱比锡大学信息学研究所

随着大型语言模型(LLMs)领域的加速发展, 评估和监控其性能的关键需求出现了。我们介绍了一个以知识图谱工程(KGE)为重点的基准测试框架, 并提出了三个挑战, 即语法和错误纠正、事实提取和数据集生成。我们表明, 虽然 LLMs 是一个

有用的工具，但它还不适合帮助零提示的知识图生成。因此，我们的 LLM-kg-bench 框架提供了 LLM 响应的自动评估和存储，以及统计数据 and 可视化工具，以支持及时跟踪工程和模型性能。

1215, 公平感知因果模型的适应速度分析

Adaptation Speed Analysis for Fairness-aware Causal Models 摘录,

2023 年 8 月 31 日

Yujie Lin, Chen Zhao, Minglai Shao, Xujiang Zhao, Haifeng Chen

内容: 例如, 在机器翻译任务中, 为了实现两种语言之间的双向翻译, 通常会使用源语料库作为目标语料库, 这涉及训练两个相反方向的模型。在许多领域中, 哪个模型能够最快地适应领域转移的问题具有重要意义。具体来说, 考虑一个原始分布 P , 由于未知的干预而发生改变, 导致修改后的分布 P^* 。在将 P 与 P^* 对齐时, 有几个因素可以影响适应速度, 包括 P 中变量之间的因果依赖关系。然而, 在现实生活中, 我们必须考虑培训过程的公平性, 尤其是在因果变量之间存在敏感变量 (偏差) 的情况下。为了探索这种情景, 我们研究了一个简单的结构因果模型 (SCM), 其中存在因果-偏差-效应结构, 其中变量 A 在因果变量 (X) 和效应变量 (Y) 之间充当敏感变量。在因果-偏差-效应 SCM 中, 两个模型分别展示了一致和相反的因果关系方向。在对 SCM 内的变量进行未知干预后, 我们可以模拟一些领域转移的情况以进行分析。然后, 我们比较了两个模型在四种转移场景下的适应速度。此外, 我们证明了两个模型在所有干预情况下适应速度之间的关联。

1216, 打开黑盒: 非语言任务中的预训练语言模型中的注意力权重和隐藏状态分析

Opening the Black Box: Analyzing Attention Weights and Hidden States in Pre-trained Language Models for Non-language Tasks

Mohamad Ballout, Ulf Krumnack, Gunther Heidemann, Kai-Uwe

Kuhnberger (德国奥斯纳布吕克大学认知科学学院)

(北大李因立摘录, 2023. 6. 21)

由于大多数高级模型的“黑箱”性质，调查深度学习语言模型一直是一个重要的研究领域。随着最近基于转化器的预训练语言模型的进步，以及它们越来越多地融入到日常生活中，解决这个问题已经变得更加紧迫。为了实现可解释的人工智能模型，必须理解其中的程序步骤，并将其与人类的思维过程进行比较。因此，在本文中，我们使用简单的、被充分理解的非语言任务来探索这些模型的内部运作。具体来说，我们将一个预先训练好的语言模型应用于具有层次结构的受限算术问题，以分析其注意力权重分数和隐藏状态。调查显示了有希望的结果，该模型以一种适度的结构化方式处理分层问题，类似于人类解决问题的策略。此外，通过逐层检查注意力权重，我们发现了一个非常规的发现，即第 10 层，而不是模型的最后一层，是解冻的最佳层，可以用最少的参数密集型方法来微调模型。我们用熵分析和标记嵌入相似性分析来支持这些发现。注意力分析使我们假设该模型可以推广到 ListOps 数据集中的较长序列，这一结论后来通过对比训练集长的序列的测试得到了证实。最后，通过利用一个简单的任务，即模型预测井字游戏的赢家，我们发现了注意力分析的局限性，特别是它无法捕捉二维模式。

1217, 足球期望进球模型的 Glocal 解释

Glocal Explanations of Expected Goal Models in Soccer

Mustafa Cavus (Department of Statistics, Eskisehir Technical University, 土耳其埃斯基谢希尔工业大学统计系), Adrian Sta'ndo and Przemys'law Biecek (Faculty of Mathematics and Information Science, Warsaw University of Technology, 华沙理工大学数学与信息科学学院)

(北大严丽君摘录, 2023. 8. 29)

预期目标模型越来越受欢迎，但其可解释性往往有限，尤其是在使用黑盒方法进行训练时。可解释的人工智能工具已经出现，可以提高模型的透明度，并提取单个观察或所有观察的描述性知识。然而，解释特定观测组的黑匣子模型在某些领域可能更有用。本文介绍了预期目标模型的全局解释（在局部和全局级别之间），通过建议使用 SHAP 值和部分依赖性简档的聚合版本，实现团队和球员级别的绩效分析。这允许从球员或球队的预期进球模型中提取知识，而不仅仅是一次射门。此外，我们还进行了 realdata 应用程序，以说明聚合 SHAP 和聚合概要文件的有

用性。论文最后对这些解释的潜力作了评论。

1218, 面向物联网的可解释人工智能 (XAI) : 一篇调研

Explainable Artificial Intelligence (XAI) for Internet of Things: A Survey

Ibrahim Kok (Pamukkale University, 棉花堡大学), Feyza Yildirim Okay (Gazi Üniversitesi, 格济大学), Ozgecan Muyanli & Suat Ozdemir (Hacettepe University, 哈西德佩大学)

(北大李劲楠摘录, 2022. 6. 7)

人工智能 (AI) 模型的黑盒特性使用户无法理解并有时难以信任该模型生成的输出。在那些不仅关注结果, 还关注结果决策路径的 AI 应用中, 这种黑盒 AI 模型是不够的。可解释人工智能 (XAI) 解决了这一问题, 定义了一组能够被用户解释的 AI 模型。最近, 已经提出了许多 XAI 模型, 以解决各种应用领域 (如医疗、军事、能源、金融和工业) 中黑盒模型缺乏可解释性和可解释性的问题。尽管 XAI 的概念近年来引起了广泛关注, 但其在物联网领域的整合尚未完全明确定义。在本文中, 我们深入系统地审查了最近在物联网领域范围内使用 XAI 模型的研究。我们根据它们的方法论和应用领域对这些研究进行了分类。此外, 我们旨在关注具有挑战性的问题和未解决的问题, 并提供未来方向, 以指导开发人员和研究人员进行未来的研究。

1219, 基于漂移和扩散记忆阻器的突触 STDP 特性的 SPICE 研究

SPICE Study of STDP Characteristics in a Drift and Diffusive Memristor-Based Synapse for Neuromorphic Computing

SUMAN HU 1,2, JAEHYUN KANG 1,3, TAEYOON KIM 1, SUYOUN LEE 1,
JONG KEUK PARK 1, INHO KIM 1, JAEWOOK KIM 1, JOON YOUNG KWAK 1,
JONGKIL PARK 1, (Member, IEEE), GYU-TAE KIM2, SHINHYUN CHOI 4,
AND YEONJOO JEONG

(北大周鑫摘录, 2022. 1. 5)

神经形态硬件是一个具有巨大潜力的系统, 可以通过模仿人脑来实现高效计算。

该新系统使用神经元尖峰（动作电位）处理信息，神经元之间的突触连接使用生物学上合理的方法进行训练，如尖峰时间依赖性可塑性（STDP）。而忆阻器有望实现这种神经形态硬件。

目前可以通过两种类型的忆阻器：扩散型和漂移型，来形成显示贴近 STDP 的模拟的突触，其中扩散效应用于追踪对 STDP 至关重要的尖峰时序记录，并且漂移型忆阻器可以将权重信息保持在较长的时间标度中。

在本文中，作者研究了具有两个串联记忆电阻器（扩散型和漂移型）的突触的 STDP 曲线特性。在 SPICE 模拟中，系统地检查了器件参数（如时间常数（ τ ），阈值和通断状态电阻）的影响。结果表明，STDP 曲线中的突触响应 \max_{1w} 和 \max_{1t} 对器件参数非常敏感，并且根据参数变化产生四种不同形状的 STDP 曲线。四个案例中有一个具有典型的 STDP 形状，但它只出现在一个相当窄的参数范围内。作者基于器件动力学以及扩散型、漂移型忆阻器之间的压分效应，分析了器件的有趣行为。在使用串行连接的突触装置时，应选择适当的参数以产生模拟生物行为的所需 STDP 曲线。这些结果可以为突触的使用提供指导，并使其能够通过适当的 STDP 曲线忠实地模拟生物可塑性。

1220, 用于神经形态计算应用的多层铁磁自旋电子器件

Multilayer Ferromagnetic Spintronic Devices for Neuromorphic Computing Applications

Aijaz H. Lone, Xuecui Zou, Kishan K. Mishra, Venkatesh Singaravelu, Hossein Fariborzi, Gianluca Setti

（北大杨思恒摘录，2023.09.13）

自旋电子学由于其在节能存储器、逻辑和非传统计算范式中的应用而取得了实质性进展。为了理解畴壁和 skyrmion 动力学，人们对多层铁磁薄膜进行了广泛的研究。然而，这些研究大多局限于材料和畴壁/skyrmion 物理。在本文中，作者介绍了用于神经形态计算应用的多层铁磁自旋电子器件的实验和微磁实现。该器件表现出多级电阻状态，并且电阻状态的数量随着温度的降低而增加。在微磁模拟中观察到的多级磁化行为支持了这一点。此外，在实验和模拟中还探讨了阻力态随自旋轨道力矩的演化。利用该器件的多级电阻态，作者提出了它作为硬件神

神经网络中的突触器件的应用，并研究了突触器件的线性性能。基于这些设备的神经网络使用监督学习算法在 MNIST 数据集上进行训练和测试。芯片级器件的精度达到 90%。从而证明了它在神经形态计算中的应用。此外，作者最后讨论了该器件在量子计算机低温存储电子学中的可能应用。

1221, 用于高线性神经形态计算的掺锂人工突触忆阻器

Li-Ion Doped Artificial Synaptic Memristor for Highly Linear Neuromorphic Computing

Jialin Meng , Zhenhai Li, Yuqing Fang , Qingxuan Li , Zhenyu He, Tianyu Wang , Hao Zhu , Member, IEEE, Li Ji , Member, IEEE, Qingqing Sun , Member, IEEE, David Wei Zhang, and Lin Chen

(北大王世界摘录, 2022.12)

神经形态忆阻器中的线性权重调制在高精度图像识别任务中起着重要作用。本文提出了一种用于神经形态计算的 Li⁺掺杂有机人工忆阻器，用于线性权重更新，该忆阻器显示了活性离子在生物突触中的扩散动力学，如 Ca²⁺扩散。忆阻器表现出渐进的电阻切换、多状态存储和典型的突触行为。此外，突触学习字母“T”的能力在忆阻器阵列中得到了证明。通过设计具有增强幅度的连续脉冲波形，用于长时程增强和抑制中的权重更新的忆阻器（LTP/LTD）的线性度可以从 6.8 提高到 0.4。基于在 LTP（ $\alpha_p=1.5$ ）和 LTD（ $\alpha_d=0.4$ ）中的非线性因子，通过由离子掺杂忆阻器组成的人工神经网络实现了 96% 的高精度人脸识别。具有高度线性权重更新的离子掺杂有机忆阻器为仿生离子扩散神经形态计算系统的开发提供了指导。本文制备了用于模拟人工突触的 Li⁺掺杂有机忆阻器，该忆阻器表现出渐进的电阻切换特性和多状态存储行为。人工突触忆阻器显示了包括 EPSC、SDDP、LTP/LTD 和阵列学习在内的突触可塑性。对于高度线性的神经形态计算，通过设计施加的电刺激来优化忆阻器的线性。基于忆阻器 LTP（ $\alpha_p=1.5$ ）和 LTD（ $\alpha_d=0.4$ ）中的优化非线性因子，人脸图像被成功识别，准确率为 96%。这项工作为实现高度线性的神经形态计算提供了一条新的途径。

1222, 嵌入文字知识图谱的通用预处理算子

Universal Preprocessing Operators for Embedding Knowledge Graphs with Literals

Presiner、Paulheim (德国曼海姆大学数据与网络科学组)

(北大董泽锦摘录, 2023年9月6日)

知识图谱嵌入是知识图谱(KG)中实体的密集数值表示。虽然大多数方法只集中于关系信息,即实体之间的关系,但也考虑到文字值信息(例如文本描述或数字信息)的方法较少。存在的那些通常是针对特定的文字形式和特定的嵌入方法进行定制的。在本文中,我们提出了一组通用预处理算子,可用于转换具有数值、时间、文本和图像信息,以便转换后的KGs可以用任何方法嵌入。在kgbench数据集上使用三种不同的嵌入方法得到了令人满意的结果。

1223, 图顶点模型

Graph Vertex Model

Tanmoy Sarkar, Matej Krajc,

论文时间: 2023.9.9

细胞重新排列是驱动生物组织大规模变形的基本机制。在三维(3D)空间中填充的细胞聚集体中,细胞通过细胞间网络的局部拓扑转换来重新排列,这最方便地由顶点模型描述。令人惊讶的是,由于传统顶点模型的复杂结构,这些转换尚未得到适当的数学表述。结果,顶点模型通常难以实现,尤其是在其完整的3D表示中。事实上,少数现有的完整3D顶点模型实现依赖于高度定制化和复杂的软件工程解决方案,这些解决方案无法透明地描述,因此大多数无法再现。我们提出了一个解决这个突出问题的解决方案,即引入了顶点模型的新形式,称为图顶点模型(Graph Vertex Model, GVM)。GVM基于将细胞网络的拓扑存储到知识图中。其数据结构由元图唯一定义,允许通过简单的图形转换执行细胞重新排列事件,这些图形转换本身由图形表示。这些图形转换在数学上得到了很好的表述,并由元操作组成,例如在知识图的节点之间删除和创建链接,这些操作很容易实现。显著的是,在GVM的数据表示中,3D空间填充的多面体包装中的复杂拓扑变化可以分解为更基本的T1转换的组合。值得注意的是,当应用于2D系统时,

这些转换减少为单个 T1 转换。这一发现统一了 2D 和 3D 空间填充包装中的拓扑转换，并暗示了 GVM 的图数据结构可能是这些系统的最自然表示。除了解决计算生物物理学中的关键问题并为 2D 和 3D 空间填充包装的拓扑转换提供基本见解外，我们的方法更广泛地展示了知识图等图数据模型在组织数据的结构化存储、分析和操作方面的适用性，这可能对实验组织生物物理学和生物学产生深远影响

1224, 图神经网络中的信息流：临床分诊使用案例

Information Flow in Graph Neural Networks: A Clinical Triage Use Case
Victor Valls (Trinity College Dublin, 都柏林圣三一学院), Mykhaylo Zayats (University of Galway, 高威大学), Alessandra Pascale (Dublin Research Lab)

(北大张凯迪摘录, 2023. 9. 13)

图神经网络 (GNN) 因其处理多模态和多关系图的能力而在医疗保健和其他领域大受欢迎。然而, GNNs 的高效训练仍具有挑战性, 存在多个未决研究问题。在本文中, 我们研究了 GNN 中的嵌入信息流如何影响知识图谱 (KG) 中的链接预测。具体来说, 我们提出了一个数学模型, 将 GNN 连接性与图数据的连接性分离开来, 并评估了 GNN 在临床分流使用案例中的性能。我们的结果表明, 与使用与 KG 相同的连接性或允许无约束嵌入传播相比, 将领域知识纳入 GNN 连接性能带来更好的性能。此外, 我们还发现负边缘在实现良好预测方面起着至关重要的作用, 而使用过多的 GNN 层会降低性能。

1225, 利用混合 3D 表示学习解开纠缠的化身

Learning Disentangled Avatars with Hybrid 3D Representations
Yao Feng, Weiyang Liu, Timo Bolkart, Jinlong Yang, Marc Pollefeys, Michael J. Black, 剑桥大学马普研究所

(北大李世平摘录, 2023. 9. 13)

巨大的努力已经取得了学习动画和照片逼真的人类化身。为此, 显式和隐式 3D 表示都被大量研究以用于整个人的整体建模和捕获 (例如, 身体、衣服、脸和头发), 但是就表示功效而言, 这两种表示都不是最佳选择, 因为人类化身的不同

部分具有不同的建模需求。例如，网格通常不适用于建模衣服和头发。出于这一动机，我们提出了解开化身 Δ (DELTA)，它模拟人类与混合显式-隐式的 3D 表示。DELTA 将单目 RGB 视频作为输入，并产生具有单独的身体和衣服 头发层的人类化身。具体来说，我们展示了两个重要的应用程序 DELTA。对于第一个，我们考虑解开人体和衣服，在第二个，我们解开脸和头发。为此，DELTA 使用显式的基于网格的参数化 3D 模型表示身体或面部，并使用隐式神经辐射场表示衣服或头发。为了实现这一点，我们设计了一个端到端的可区分渲染器，将网格集成到体积渲染中，使 DELTA 能够直接从单目视频中学习，而无需任何 3D 监督。最后，我们展示了如何将这两个应用程序可以很容易地结合起来，以模拟全身化身，这样的头发，脸，身体和衣服可以完全解开，但联合呈现。这样的解缠结使得头发和衣服能够转移到任意的身体形状。我们经验验证的有效性，三角洲的解缠，展示其有前途的表现上解开重建，虚拟服装试穿和发型转移。为了促进未来的研究，我们还发布了一个开源的管道，用于混合人类化身建模的研究。

1226, 用于软件开发的可通信智能体

Communicative Agents for Software Development

Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan Xu, Dahai Li, Zhiyuan Liu, Maosong Sun

(北大李政恩摘录, 2023. 7. 16)

软件工程是一个以复杂决策过程为特征的领域，通常依赖于微妙的直觉和咨询。最近深度学习的进展已经开始通过在软件开发的各个阶段实现复杂的设计来革新软件工程实践。在本文中，我们提出了一种创新的范式，通过自然语言交流在整个软件开发过程中利用大型语言模型 (LLMs)，从而通过消除每个阶段的专门模型来简化和统一关键过程。这个范式的核心是 ChatDev，一个虚拟的聊天驱动软件开发公司，它模拟了已经建立的瀑布模型，将开发过程细分为四个不同的时间阶段：设计、编码、测试和文档编写。每个阶段都涉及一个团队的代理，例如程序员、代码审查员和测试工程师，促进协作对话并促进无缝工作流程。聊天链作为一个促进者，将每个阶段分解成原子子任务。这使得双重角色成为可能，允许通过上下文感知的交流提出和验证解决方案，从而有效地解决特定的子任务。

ChatDev 的工具分析突出了它在软件生成方面的显著效力，使得整个软件开发过程在不到七分钟的时间内以不到一美元的成本完成。它不仅能够识别和缓解潜在的漏洞，还能够纠正潜在的幻觉，同时保持值得称赞的效率和成本效益。ChatDev 的潜力揭示了将 LLMs 整合到软件开发领域的新可能性。

1227, 人脑听语音时存在预测编码层次结构的证据

Evidence of a predictive coding hierarchy in the human brain listening to speech

Charlotte Caucheteux, Alexandre Gramfort & Jean-Rémi King

(北大郜岱威摘录, 2023. 3. 2)

最近，自然语言处理取得了长足的进步：深度学习算法越来越能够生成、总结、翻译和分类文本。然而，这些语言模型仍然无法匹配人类的语言能力。预测编码理论对这种差异提供了一个初步的解释：虽然语言模型经过优化可以预测附近的单词，但人脑会持续预测跨越多个时间尺度的表示层次结构。为了验证这一假设，我们分析了 304 名听短故事的参与者的功能磁共振成像大脑信号。首先，我们确认现代语言模型的激活线性映射到大脑对语音的反应。其次，我们表明，通过跨越多个时间尺度的预测来增强这些算法可以改善这种大脑映射。最后，我们表明这些预测是分层组织的：额顶叶皮层比颞叶皮层预测更高级别、更远范围和更多上下文的表征。总体而言，这些结果加强了分层预测编码在语言处理中的作用，并说明了神经科学和人工智能之间的协同作用如何揭示人类认知的计算基础。

1228, 解释黑盒模型：可解释人工智能综述

Interpreting Black-Box Models: A Review on Explainable Artificial Intelligence

Vikas Hassija, Vinay Chamola, Atmesh Mahapatra

机构：印度布巴内斯瓦尔卡林加工业技术学院计算机工程学院

(北大徐永登 2023 年 8 月 24 日)

内容：近年来，基于人工智能（AI）的方法论在广泛的领域取得了巨大的发展。在这个快速发展的领域，大量使用机器学习（ML）和深度学习（DL）模型的方法

被报道。这些模型大多数在本质上都很复杂，并且缺乏对决策过程的解释，导致这些模型被称为“黑匣子”。在银行、电子商务、医疗保健、公共服务和安全等关键任务应用领域采用此类模型的主要瓶颈之一是它们难以被解释。由于这些人工智能模型的迅速扩散，解释它们的学习和决策过程变得越来越困难，为此需要其透明性和易于预测性。本研究旨在整理解释黑盒模型的最新技术，对可解释的人工智能（XAI）模型进行全面分析。减少这些底盒模型错误的负面或者正面结果，并找出它们的缺陷仍然是困难且低效的。本文通过对 XAI 研究现状的仔细选择和分析，对 XAI 的发展进行了细致的回顾。它还对 XAI 框架及其功效进行了全面、深入的评估，作为应用和理论研究人员 XAI 的起点。最后，它强调了与 XAI 研究相关的新兴关键问题，以展示主要的、特定于模型的趋势，以便更好地解释、增强透明性和提高预测准确性。

1229, 通过神经逻辑推理增强深度多智能体强化学习的可解释性

Enhancing the Interpretability of Deep Multi-agent Reinforcement Learning via Neural Logic Reasoning

Bokai Ji, Guangxia Li & Gang Xiao

机构: School of Computer Science and Technology, Xidian University, Xi'an, China 西安电子科技大学计算机科学与技术学院

(北大韩亚帅, 2023.09.22)

解释深度强化学习的决策策略是具有挑战性的，因为神经网络具有黑盒性质。为了解决这个挑战，我们将深度学习模型和符号结构结合到一个神经逻辑模型中，以神经逻辑编程的形式进行推理。所提出的可解释的多智能体强化学习算法在以符号表示的环境中进行推理，使用多跳推理和关系路径搜索方法，利用先前的符号知识。此外，为了缓解多智能体系统中的部分可观察性问题，我们设计了一个可解释的历史模块，使用注意机制来融合过去的经验，同时保持可解释性。实验表明，所提出的方法能够有效地学习接近最优的策略，并生成用于解释决策的表达规则。特别是，它可以学习比传统神经网络方法更抽象的概念。

1230, 机器学习生命周期中可解释人工智能的千张面孔: 工业现实和研究现状
The Thousand Faces of Explainable AI Along the Machine Learning Life Cycle:
Industrial Reality and Current State of Research

Thomas Decker, Ralf Gross, Alexander Koebler, Michael Lebacher, Ronald Schnitzer, Stefan H. Weber

机构: 西门子股份公司, 德国慕尼黑、德国慕尼黑路德维希·马克西米利安大学、德国法兰克福歌德大学、德国慕尼黑工业大学

(北大李星辰, 2023 年 10 月 7 日)

这篇文章研究了可解释人工智能 (XAI) 的实际相关性, 特别关注生产行业, 并将其与学术 XAI 研究的现状联系起来。基于一系列广泛的采访, 这些采访涉及 XAI 在当前工业实践中机器学习 (ML) 生命周期中的作用和适用性, 以及它在未来的预期相关性。采访对象是来自不同行业部门的各种角色和关键利益相关者。除此之外, 通过对相关文献的简要回顾, 概述了 XAI 研究的现状, 并提供了一个全面的能涵盖被调查者的意见以及学术研究的现状概述。通过将访谈结果与当前的研究方法进行比较, 发现了一些差异。虽然存在许多不同的 XAI 方法, 但大多数方法都围绕着模型评估阶段和数据科学家。它们在其他阶段的通用能力目前要么没有得到充分探索, 要么在从业者中不受欢迎。与现有工作一致, 研究结果还证实, 需要做出更多努力, 让非专家用户也能用现有的方法和框架解释和理解不透明的人工智能模型。

1231, 预训练的转换器真的是通过梯度下降进行上下文学习的吗?

Do Pretrained Transformers Really Learn In-context by Gradient Descent?

Lingfeng Shen[♡] Aayush Mishra[♡] Daniel Khashabi

机构: Johns Hopkins University, Baltimore MD

(北大岳沈童 2023 年 10 月 12 日)

上下文学习 (ICL) 是否隐含地等同于梯度下降 (GD)? 最近有几项研究将梯度下降的动态行为与大语言模型中上下文学习的新兴行为进行了类比。然而, 这些研究所做的假设与训练语言模型的真实自然语言环境相去甚远。因此, 有必要进一步研究理论与实践之间的差异, 以验证其适用性。

作者首先强调了之前通过构建转换器的权重来模拟梯度下降的研究中存在的不足。他们在上下文学习的目标上对转换器进行训练的实验、上下文学习和梯度下降对顺序敏感性的不一致、所构建权重的稀疏性，以及对参数变化的敏感性等，都是与现实环境不匹配的一些例子。

此外，作者还在自然环境中对上下文学习对比梯度下降的假设进行了探究和比较。作者对在自然数据（LLaMa-7B）上预先训练的语言模型进行了全面的实证分析。作者对各种性能指标进行了比较，结果表明，上下文学习和梯度下降的行为因数据集、模型和演示次数等各种因素的影响而不一致。作者观察到，上下文学习和梯度下降以不同方式改变语言模型的输出分布。这些结果表明，上下文学习和梯度下降之间的等效性是一个开放的假设，需要细致的考虑，并需要进一步的研究。

1232, 利用 NeuroKinect 提高手部轨迹估计中的脑机接口系统性能

Advancing Brain-Computer Interface System Performance in Hand Trajectory Estimation with NeuroKinect

论文作者: Sidharth Pancholi, Amita Giri

机构: IEEE

(北大张震扬摘录, 2023. 8. 16)

论文内容: 脑机接口 (BCI) 技术可以实现大脑与外部设备之间的直接通信, 使个人能够利用大脑信号来控制周围的环境。然而, 现有的 BCI 方法面临着三个阻碍其实用性和有效性的关键挑战: 首先, 预处理算法非常耗时; 其次, 不适当的损失函数利用; 最后, 超参数的设置不太直观。为了解决这些限制, 本文提出了 NeuroKinect, 这是一种创新的深度学习模型, 用于使用脑电图 (EEG) 信号精确重建手部运动学。NeuroKinect 模型使用最少的预处理管道在抓握并举起 (GAL) 任务数据上进行训练, 从而提高了计算效率。NeuroKinect 引入的一个显着改进是使用了一种新颖的损失函数, 表示为 LStat。该损失函数解决了手部运动学预测中相关性和均方误差之间的差异。此外, 论文的研究强调参数选择背后的科学直觉以提高准确性。通过利用事件相关电位和脑源定位 (BSL) 结果来分析运动任务的空间和时间动态。这种方法为最佳参数选择提供了宝贵的

见解，从而提高了 NeuroKinect 模型的整体性能和准确性。论文中的模型显示预测手部动作和实际手部动作之间存在很强的相关性，X、Y 和 Z 维度的平均 Pearson 相关系数分别为 0.92 (± 0.015)、0.93 (± 0.019) 和 0.83 (± 0.018)。除此之外，X、Y 和 Z 维度的均方误差 (MSE) 分别为 0.016 (± 0.001)、0.015 (± 0.002) 和 0.017 (± 0.005)，证明了 NeuroKinect 模型的精度也非常高。

1233, 感知运动脑机接口的神经可塑性

Neural Plasticity in Sensorimotor Brain - Machine Interfaces

Maria C. Dadarlat (美国印第安纳州西拉法叶普渡大学威尔登生物医学工程学院), Ryan A. Canfield, (美国华盛顿州西雅图华盛顿大学生物工程系) and Amy L. Orsborn (美国华盛顿州西雅图华盛顿大学电气与计算机工程系)

在这篇论文中，介绍了学习在人工运动输出或人工感觉输入闭环 BMI 中的作用。在这两种情况下，本文将重点限制在神经测量和调制的侵入式方法上，因为在这些方法中，已有大量研究探索了相关的生物回路。本文还关注与 BMI 的闭环适应性相关的主题。对神经元如何响应和适应人工感觉运动任务的深入了解将为新的记录和神经调控技术以及非侵入式系统的重要发展提供信息。本文认为作为人类，我们的目标是与世界互动；我们对环境中的各种因素做出反应，并通过我们的感觉和运动系统对其采取行动。损伤和疾病会中断感觉运动功能，从根本上限制我们的参与能力。在人们永久性瘫痪的情况下，脑机接口 (BMI) 旨在通过提供人工感觉运动功能来恢复独立性。脑机接口已通过控制打字恢复了交流能力 (1, 2)，或通过控制假肢恢复了日常生活的基本任务。典型的 BMI 包括记录患者大脑的神经活动，将神经信号转化为运动指令，并提供感官反馈，以实现目标导向运动和纠错。构建 BMI 需要进行一系列计算，这些计算可以单独优化，但作为闭环系统的一部分又相互关联。在闭环 BMI 中，感觉输入和运动输出之间的紧密联系调动了与生俱来的神经学习机制，从而改变了为优化 BMI 性能而必须解决的工程问题的性质。脑机接口 (BMI) 旨在通过创建人工运动和/或感觉通路来治疗感觉运动神经系统疾病。人工通路的引入在感觉输入和运动输出之间建立了新的关系，大脑必须通过学习才能获得灵巧的控制。本综述强调了学习在恢复运动和感觉的 BMI 中的作用，并讨论了 BMI 的设计如何影响神经可

塑性和性能。感觉和运动功能可塑性的紧密结合影响着这两种人工通路的设计，也将成为恢复感觉和运动功能的双向装置的基本考虑因素。

1234, 基于深度学习模型的 SSVEP 在 BCI 中的分析: 综述研究

An Analysis of Deep Learning Models in SSVEP-Based BCI: A Survey

作者: Dongcen Xu, Fengzhen Tang, Yiping Li, Qifeng Zhang, Xisheng Feng

中国科学院沈阳自动化研究所, 机器人学国家重点实验室, 中国辽宁省沈阳市

这项研究旨在提高不同受试者的运动脑机接口 (BCI) 性能。传统 BCI 在 EEG 试验中判别段的位置和持续时间存在差异, 导致性能差。为了克服这个问题, 研究提出了一种浅层镜像变压器, 利用多头自注意力层从整个 EEG 试验中检测和利用判别信号段。此外, 引入镜像 EEG 信号和镜像网络结构, 通过集成学习提高了分类准确性。实验结果表明, 这一方法在 BCI 竞赛 IV 数据集 2a、2b 和 OpenBMI 数据集上取得了令人满意的性能, 对于新受试者和现有受试者的平均准确性分别为 74.48% 和 76.1%, 超过了其他先进方法。该研究还突显了多头自注意力在捕获运动意象分类中的全局 EEG 信号信息方面的有效性。总之, 这一研究提供了一种基于多头自注意力层的有效模型, 用于不同受试者的运动意象 BCI, 而且是迄今为止最浅的变压器模型之一, 其参数少, 适用于小样本问题。

1235, 通过局部节奏稀疏脉冲网络

Sparsifying Spiking Networks through Local Rhythms

Wilkie Olin-Ammentorp

Mathematics and Computer Science Division Argonne National Laboratory

Lemont, IL, USA 美国伊利诺伊州的一个数学和计算机科学实验室

(北大邓凯元摘录, 2023 年 9 月)

在传统神经网络中, 权重和激活值中都表现出稀疏性, 每一层产生的许多值都为零。而哺乳动物大脑也表现出许多形式的稀疏性: 大脑中的神经元遵循对数正态分布的射频分布, 一半的神经元每秒不到一次。

这项研究提出了通过相位编码来构建脉冲神经网络, 产生一种稀疏化方法, 展示了脉冲神经网络可以利用局部信息阻止传输代表接近零的脉冲, 从而减少了这些

网络中通信和计算所需的能量，同时保持了准确性。

作者尝试了显式稀疏 (Explicit Sparsification)、抑制性稀疏 (Inhibitory Sparsification)、随机稀疏 (Random Sparsification) 三种稀疏方法。如下图所示：随机删除的性能随着删除脉冲的比例呈线性减少的趋势，另外两种稀疏化方法显示出明显不同的趋势：网络的性能在删除了大约 60% 的脉冲之前变化不大。这说明了当一个层中的神经元被偏向于以同步方式自发激发时，许多脉冲是多余的。

1236, DYNAP-SE2: 可扩展的多核动态神经形态异步尖峰神经网络处理器

DYNAP-SE2: a scalable multi-core dynamic neuromorphic asynchronous spiking neural network processor

Ole Richter, Chenxi Wu, Adrian M. Whatley, German Köstinger, Carsten Nielsen, Ning Qiao, Giacomo Indiveri

内容：随着技术的显著进步，处理边缘传感器附近数据的需求急剧增加。这些应用中使用的电子系统必须连续、实时地处理数据，并使用尽可能小的能量预算提取相关信息。实现支持按需、稀疏和边缘计算的感觉信号始终在线处理的一种有前途的方法是从生物神经系统中获得灵感。这篇论文介绍了一种混合信号神经形态处理器，该处理器采用事件驱动的脉冲神经网络 (SNNs) 作为实现始终开启传感信号处理的一种方法。该处理器通过使用模拟集成电路进行计算来直接模拟生物神经元和突触的动态，并使用异步数字电路传输神经元产生的事件（脉冲）到目标突触或输出端口。该处理器具有无时钟的异步数字分层路由方案，可以在本地实时运行，确保低延迟。处理器被命名为 DYNamic Neuromorphic Asynchronous Processor-Scalable 2 (DYNAP-SE2)，在突触和神经元电路级别、网络级别和异步路由结构级别扩展了先前一代 DYNAP-SE 的功能。它支持多个时间尺度上的丰富神经元动力学，适用于广泛的生物合理的循环网络。该论文详细描述了处理器的总体架构、各个电路的设计，以及从芯片中测得的实验结果，验证了理论。处理器还集成了一个用于将输入波形转换为地址事件流的模拟前端 (AFE) 和一个用于处理来自动态视觉传感器 (DVS) 的二维事件刺激的直接 2D 传感器事件预处理器。论文分为几个部分，包括总体架构、关键电路、核心模拟神经电路及其

应用示例、路由方案和构建大规模神经网络的方法、芯片的外部接口，以及支持芯片可用性的软件系统。

1237, 知识图谱:介绍、历史和观点

Knowledge graphs: Introduction, history, and perspectives

Vinay K. Chaudhri, Chaitanya Baru, Naren Chittar et al.

斯坦福大学, 圣地亚哥超级计算机中心, 摩根大通公司等。

(北大孙宇瑶摘录, 2022年3月31日)

内容: 知识图谱这一概念日渐引起关注, 它用于组织世界上结构化的知识、集成从多个数据源中提取的信息, 表示 AI/ML 算法的输入和输出信息等, 也是自然语言处理 (NLP)、计算机视觉 (CV) 和常识推理的基本技术。在文中以维基数据这一协作编辑的开放知识图谱为例进行阐述, 维基数据是具有空前的数据规模, 是目前公开可用的最大知识图谱之一。在文章的最后, 对符号表示法和向量表示法、人类管理和机器管理、小语义和大语义三组情况进行了比较。

1238, 大型语言模型的图神经提示

Graph Neural Prompting with Large Language Models

Yijun Tian, Huan Song, Zichen Wang, 等

(北大杨龙摘录, 2023年8月26日)

大型语言模型 (LLMs) 在各种语言建模任务中都表现出了卓越的泛化能力和出色的性能。然而, 它们在准确捕捉和返回有根据的知识方面仍然存在固有的限制。尽管现有的研究已经探索了如何利用知识图谱通过联合训练和定制模型架构来增强语言建模, 但由于 LLMs 的参数众多和计算成本高, 将此应用于 LLMs 是有问题的。此外, 如何利用预训练的 LLMs 并避免从头开始训练一个定制模型仍然是一个悬而未决的问题。在这项工作中, 本文提出了图神经提示 (GNP) 这一新颖的即插即用方法, 以协助预训练的 LLMs 从 KGs 中学习有益的知识。GNP 包括各种设计, 包括标准的图神经网络编码器、跨模态池化模块、领域投影器和自监督链接预测目标。在多个数据集上的广泛实验证明了 GNP 在不同 LLM 大小和设置的常识和生物医学推理任务上的优越性。

1239, From Large Language Models to Knowledge Graphs for Biomarker Discovery in Cancer

Md. Rezaul Karim, Lina Molinas Comet, Md Shajalal, Oya Beyan, Dietrich Rebholz-Schuhmann, Stefan Decker

领域专家通常依靠最新知识来理解和传播特定的生物过程，帮助他们设计制定预防和治疗决策的策略。人工智能（AI）的一个具有挑战性的场景是使用生物医学数据（例如文本、成像、组学和临床）来提供癌症疾病的诊断和治疗建议。有关癌症、药物、基因、蛋白质及其机制的数据和知识分布在结构化（知识库（KBs））和非结构化（例如科学文章）来源中。通过整合这些数据，然后提取有关语义上相互关联的实体和关系的事实，可以构建大规模知识图谱（KG）。这样的 KG 不仅允许探索和问答（QA），还允许领域专家推断新知识。然而，由于缺乏对底层数据资产和语义技术的理解，探索和查询大规模 KG 对于非域用户来说很乏味。在本文中，我们开发了一个域 KG 来利用癌症特异性生物标志物发现和交互式 QA。为此，开发了一种称为 OncoNet Ontology（ONO）的领域本体，以实现语义推理来验证基因 - 疾病关系。然后，通过采用基于 BioBERT 和 SciBERT 的信息提取（IE）方法，通过协调 ONO，受控词汇表和科学文章中的其他生物医学概念来丰富 KG。此外，由于生物医学领域正在不断发展，新发现通常会取代旧发现，而不采用最新发现，因此人工智能系统在提供诊断和治疗时很有可能表现出概念漂移。因此，我们使用大型语言模型（LLM）根据命名实体识别模型可能未看到的最新文章和知识库对 KG 进行了微调。

1240, 预训练的转换器真的是通过梯度下降进行上下文学习的吗？

论文题目：Do Pretrained Transformers Really Learn In-context by Gradient Descent?

Lingfeng Shen[✉] Aayush Mishra[✉] Daniel Khashabi

机构：Johns Hopkins University, Baltimore MD

上下文学习（ICL）是否隐含地等同于梯度下降（GD）？最近有几项研究将梯度下降的动态行为与大语言模型中上下文学习的新兴行为进行了类比。然而，这些研究所做的假设与训练语言模型的真实自然语言环境相去甚远。因此，有必

要进一步研究理论与实践之间的差异，以验证其适用性。

作者首先强调了之前通过构建转换器的权重来模拟梯度下降的研究中存在的不足。他们在上下文学习的目标上对转换器进行训练的实验、上下文学习和梯度下降对顺序敏感性的不一致、所构建权重的稀疏性，以及对参数变化的敏感性等，都是与现实环境不匹配的一些例子。

此外，作者还在自然环境中对上下文学习对比梯度下降的假设进行了探究和比较。作者对在自然数据（LLaMa-7B）上预先训练的语言模型进行了全面的实证分析。作者对各种性能指标进行了比较，结果表明，上下文学习和梯度下降的行为因数据集、模型和演示次数等各种因素的影响而不一致。作者观察到，上下文学习和梯度下降以不同方式改变语言模型的输出分布。这些结果表明，上下文学习和梯度下降之间的等效性是一个开放的假设，需要细致的考虑，并需要进一步的研究。

1241，脑机接口的未来是光学的

The future of brain-machine interfaces is optical

Nathan Tessema Ersaro, Cem Yalcin & Rikky Muller

机构：University of California

（北大刘志一 摘录，2023.2.17）

光学接口可用于解决脑机接口（BMI）开发中与缩放、精度和侵入性相关的挑战。脑机接口的发展是由神经传感和刺激模式的激增以及可以同时记录的神经元数量的增加所驱动的。理想情况下，脑机接口应该是微创，安全，且不影响寿命的。它应当具有高时空分辨率和通信深度吞吐量，以实现大脑编码和解码。

目前的脑机接口，非侵入式 BMI 如脑电图与经颅磁刺激，侵入式 BMI 如插入脑组织的微电极，都存在不足之处。非侵入式 BMI 在通信复杂性方面具有根本限制，因为测量构成了从远处记录的总和电信号，并且因为刺激针对相对大量的组织。侵入式 BMI 受到电极阵列的寿命、可靠性和安全性受到插入过程中的物理损伤、对植入物的免疫反应以及随着时间的推移神经回路重组的限制。

光学可以在这里提供答案。基因编辑的出现使得位于神经元膜内的荧光电压指示器得以开发，该指示器跟踪单个放电事件以及通过触发或抑制神经元放电来

驱动称为视蛋白的蛋白质，这些蛋白质对光做出反应因此，通过使用颅窗，可以在不同的光学波长上进行记录和刺激，以避免读写串扰，并且不会物理穿透脑组织。通过这种方法，光学神经接口可以精确地控制和询问哺乳动物的神经回路，创建 BMI，可以实现能够驱动现实感官规则的刺激模式。

1242, 可解释人工智能 (XAI) 评估方法的新视角

A New Perspective on Evaluation Methods for Explainable Artificial Intelligence (XAI)

Timo Speith and Markus Langer

德国萨尔大学透明计算中心

(北大徐永登摘录, 2023 年 9 月 28 日)

可解释人工智能 (XAI) 领域的一大挑战是如何评估可解释性方法。人们已经提出了许多评估方法 (EM)，但黄金标准尚未建立。几位作者根据 EM 本身的各个方面 (例如，基于启发式、以人为中心、以应用程序为基础、以功能为基础) 将可解释性方法的 EM 分为几类。在这篇论文中，我们建议 EM 也可以根据它们所针对的 XAI 流程的各个方面进行分类。基于阐明 XAI 中主要流程的模型，我们提出存在解释性信息 EM、理解 EM 和需求 EM。这种新颖的观点旨在通过减少对评估方法本身的关注，而是关注可解释性方法想要实现的目标 (即提供良好的解释性信息、促进理解、满足社会需求) 来扩大其他作者的观点。我们希望这两种观点的结合将使我们能够更全面地评估可解释性方法的优点和缺点，帮助我们就使用哪些方法或如何改进它们做出更明智的决定。

1243, 利用人类相似性判断来修剪词嵌入以增强可解释性

Enhancing Interpretability using Human Similarity Judgements to Prune Word Embeddings

作者: Natalia Flechas Manrique, Wanqian Bao, Aurelie Herbelot, Uri Hasson

机构: 特伦托大学 (University of Trento)

NLP 中的可解释性方法旨在提供对特定系统架构下潜在语义的洞察。针对词嵌入，我们提出了一种有监督学习方法，针对给定领域 (例如，体育、职业)，识别出

一组模型特征,这些特征能够显著提高对人类相似性判断的预测能力。我们表明,对于 8 个独立的语义领域,这种方法仅保留了原始嵌入的 20-40%,并且在不同领域之间保留了不同的特征集。然后,我们提出了两种解释保留特征的语义的方法。第一种方法是获取保留嵌入的第一主成分上领域词(共同上义词)的分数,并提取与共同上义词共现的术语,以跟踪这些分数的特征。该分析揭示了人类根据体育项目的性别包容性和国际性的不同。第二种方法将保留的特征集作为变量应用于一项探测任务,该任务预测了 535 个单词数据集上的 65 个语义注释维度的值。保留的特征对于职业类别最能预测认知、情感和社交维度,而对于水果或蔬菜类别,保留的特征最能预测味觉维度。我们讨论了这对于 AI 系统与人类知识之间的对齐的影响。

1244, 大语言模型是专门的解释者吗?

Are Large Language Models Post Hoc Explainers?

作者: Nicholas Kroeger, Dan Ley, Satyapriya Krishna, Chirag Agarwal, Himabindu Lakkaraju

机构: 佛罗里达大学、哈佛大学

(北大李星辰摘录, 2023 年 10 月 10 日)

大型语言模型(LLM)越来越多地被用作大量自然语言处理(NLP)应用程序的强大工具。最近的一项创新,上下文学习(ICL),使 LLM 能够通过推理时间的提示中提供几个例子来学习新任务,从而消除了对模型微调的需要。虽然 LLM 已被用于多种应用,但其在解释其他模型行为方面的适用性仍相对未被探索。尽管新的解释技术越来越多,但许多技术需要对模型进行白盒访问和/或计算成本高昂,这突出了对下一代事后解释器的需求。在这项工作中,本文提出了第一个框架来研究 LLM 在解释其他预测模型方面的有效性。更具体地说,提出了一个包含多种提示策略的新框架: i) 基于扰动的 ICL, ii) 基于预测的 ICL、iii) 基于指令的 ICL 和 iv) 基于解释的 ICL。对真实世界的基准数据集进行了广泛的实验,以证明 LLM 生成的解释与最先进的事后解释者在生成模型解释时利用 ICL 示例和内部知识的能力不相上下。平均而言,在四个数据集和两个 ML 模型中,观察到 LLM 以 72.19% 的准确率识别了最重要的特征,为探索基于 LLM 的解释框

架开辟了可解释人工智能（XAI）的新领域。

1245, 代词的故事：可解释性为减轻性别偏见提供信息，以实现更公平的指令调整机器翻译

A Tale of Pronouns: Interpretability Informs Gender Bias Mitigation for Fairer Instruction-Tuned Machine Translation

Giuseppe Attanasio, Flor Miriam Plaza-del-Arco, Debora Nozza, Anne Lauscher

Bocconi University, Milan, Italy、University of Hamburg, Hamburg, Germany

（北大岳沈童摘录，2023年10月18日）

内容：最新的指令微调模型可以在提示下解决多个 NLP 任务，其中机器翻译(MT)是一个重要的应用案例。然而，目前的研究往往侧重于标准的性能基准测试，而忽略了令人信服的公平性和伦理方面的考虑。在 MT 中，这可能会导致错误的性别翻译，除其他危害外，还会造成刻板印象和偏见的延续。作者通过研究这类模型是否以及在多大程度上会在机器翻译中表现出性别偏见，以及如何减轻这种偏见，来弥补这一不足。具体来说，作者在从英语到德语和西班牙语的 WinoMT 语料库上计算了既定的性别偏见指标，并发现，IFT 模型会默认翻译成男性屈折词汇，甚至不顾刻板印象中的女性职业。接下来，作者使用可解释性方法揭示出，在错误的性别翻译中，模型系统性地忽略了表示目标职业性别的代词。最后，基于这一发现，作者提出了一种基于少量示例学习的易于实施且有效的偏差缓解解决方案，从而显著提高翻译的公平性。

1246, 反应性脑机接口的自动控制

Automatic Control of Reactive Brain Computer Interfaces

论文作者：Pex Tufvesson, Frida Heskebeck

（北大张震扬摘录，2023.10.11）

论文内容：这篇论文提出了一种免校准反应性脑机接口的方法。当前脑机接口研究的一个巨大挑战是建立有效的方法来可靠地在线分类神经活动。大多数脑机接口通过放置在头部外侧的电极来测量大脑活动。如今，这些应用主要是医疗领域，

例如帮助瘫痪的人恢复某种形式的运动功能。一些稍微更未来主义的应用是大脑到文本、人工视网膜和传递思想。脑机接口的商业市场正在扩大，多家公司提供用于冥想、心理训练或游戏的消费级设备，允许使用基于思想或情绪的输入进行游戏交互。

脑机接口（BCI）的目的之一是解码用户的大脑活动以了解他们的意图。其中，反应性脑机接口测量和分析用户响应给定刺激的大脑活动。刺激可以是一系列彩色轮廓图像，其中一种颜色是目标刺激，所有其他刺激都是非目标刺激，反应性脑机接口旨在根据测量的大脑活动来识别用户选择的目标刺激，并且希望能尽快识别目标刺激。文章讨论了基于贝叶斯统计的实时脑机接口控制方法的实践和理论方面，使用基于高斯混合模型的贝叶斯统计，将测量到的大脑信号映射到想要使用高斯混合模型（GMM）估计分布的实数，GMM 用于最佳序列选择，以最大限度地减少识别目标刺激所需的显示刺激的数量。论文还研究并改进基于视觉怪异范式的反应式脑机接口的自动控制和反馈算法的性能，以实现更快的统计收敛。引入使用高斯混合模型的迁移学习，从而实现即用型设置。论文中提出的方法朝着即用型脑机接口迈出一步，扩展了 BCI 设备和研究边界的潜力。

1247, 运动皮层、前额叶皮层和纹状体在脑机接口和手动伸手任务中的不同神经表征

Distinct neural representations during a brain-machine interface and manual reaching task in motor cortex, prefrontal cortex, and striatum

论文作者: Ellen L. Zippi (加利福尼亚大学), Gabrielle F. Shvartsman (加利福尼亚大学), Nuria Vendrell-Llopis (加利福尼亚大学), Joni D. Wallis (加利福尼亚大学) & Jose M. Carmena (加利福尼亚大学)

在这篇论文中,介绍了运动皮层、前额叶皮层和纹状体在脑机接口和手动伸手任务中的不同神经表征,虽然脑-机接口(BMI)是通过调节选定的局部神经元群来直接控制的,但由皮层和皮层下区域组成的分布式网络也参与了学习和维持控制。以前在啮齿类动物身上进行的研究表明,纹状体参与了 BMI 的学习。然而,尽管前额叶皮层在动作规划、动作选择和学习抽象任务中发挥作用,但在研究运动 BMI 控制时,前额叶皮层却在很大程度上被忽视了。在这里,我们比较了在 BMI

控制和手动控制下，非人灵长类执行二维、自我启动、中心向外的任务时，同时从初级运动皮层（M1）、背外侧前额叶皮层（DLPFC）和纹状体尾状核（Cd）记录到的局部场电位。我们的研究表明，在 M1、DLPFC 和 Cd 中存在着不同的 BMI 和手动控制神经表征。我们发现，来自 DLPFC 和 M1 的神经活动分别能最好地区分围棋提示和目标获取时的控制类型，而 M1 则能最好地预测这两个任务事件中的目标方向。我们还发现，在两种控制类型中，DLPFC → M1 之间存在有效的连接，而在 BMI 控制过程中，Cd → M1 之间也存在有效的连接。这些结果表明，在 BMI 控制过程中，M1、DLPFC 和 Cd 之间的分布式网络活动与手动控制相似但又不同。在人类和动物身上都已证明了随意调节神经元活动的的能力。这种技能是脑机接口（BMI）运行的基础，在脑机接口中，外部设备（如计算机光标或机械臂）是通过产生特定的神经活动模式来控制的。运动型脑机接口特别依赖于运动皮层神经活动的调节来控制外部设备，同时也依赖于更高层次的认知过程来整合感官信息、计划和启动行动，以及监测效应器的视觉反馈。

1248, 异步脉冲神经网络的时空自注意力

Spatial-Temporal Self-Attention for Asynchronous Spiking Neural Networks

论文作者: Yuchen Wang , Kexin Shi , Chengzhuo Lu , Yuguo Liu , Malu Zhang * and Hong Qu

论文机构: 计算机科学与工程学院, 电子科技大学

内容: 脑激励脉冲神经网络 (SNNs) 因其异步事件驱动特性和低功耗而受到越来越多的关注。随着注意机制最近成为序列依赖建模中不可或缺的一部分, SNNs 和注意机制的结合对于节能和高性能的计算范例具有巨大的潜力。然而, 现有的工作不能同时受益于时间上的注意和 SNNs 的异步特性。为了充分利用 SNNs 和注意机制的优势, 该论文提出了一种基于 SNNs 的时空自我注意力 (STSA) 机制, 该机制计算跨时间和空间域的特征依赖性, 而不破坏 SNNs 的异步传输特性。为了进一步提高性能, 论文还提出了一种时空相对位置偏差 (STRPB) 用于 STSA, 以考虑尖峰的时空位置。基于 STSA 和 STRPB, 论文中构建了一个时空尖峰 Transformer model 框架, 命名为 STS-Transformer model, 它功能强大, 使 SNNs 能够以异步事件驱动的方式工作。在流行的神经形态数据集和语音数据集上进行

了广泛的实验，包括 DVS128 手势、CIFAR10-DVS 和谷歌语音命令，论文中的实验结果优于其他最先进的模型。

1249, 走向认知预测性维护:基于图的方法综述

Toward cognitive predictive maintenance: A survey of graph-based approaches

作者: Liqiao Xia, Xinyu Li, Robert.X. Gao d, Lihui Wang, et al.

机构: 香港理工大学, 美国凯斯西储大学, 瑞典皇家理工学院

内容: 预测性维护(Predictive Maintenance, PdM) 如今在制造业中不断引起关注。然而现有的 PdM 研究多侧重于故障感知, 无法根据知识来推理故障间的因果关系, 缺乏完整的相关性分析, 由此提出了基于图形认知的认知智能方法(graph-based approaches, GbA)。本文首先介绍了 GbA 的四种主流: 图神经网络(Graph neural network, GNN)、知识图谱(Knowledge Graph)、贝叶斯网络(Bayesian network, BN)、图论模型(Graph theoretic model)及其相应的图构造方法。其次, 给出了 GbA 在 PdM 中的实现策略和实例。最后, 提出了认知 PdM 面临的挑战和未来的工作。

1250, 使用大型语言模型进行知识图谱推理的通用框架

KG-GPT: A General Framework for Reasoning on Knowledge Graphs Using Large Language Models

作者: Jiho Kim, Yeonsu Kwon, Yohan Jo, Edward Choi

机构: 韩国 韩国科学技术研究所 ; 韩国 首尔国立大学

内容: 虽然大型语言模型(LLMs)在理解和生成非结构化文本方面取得了可观的进展, 但它们在结构化数据方面的应用仍然未被充分探索。特别是, 在知识图谱(KGs)上执行复杂推理任务的 LLMs 的应用仍然几乎没有涉及。为了解决这个问题, 本文提出了 KG-GPT, 这是一个多用途框架, 利用 LLMs 来处理涉及 KGs 的任务。KG-GPT 包括三个步骤: 句子分割、图检索和推理, 每个步骤旨在将句子分割成段落、检索相关的图组件, 并推导出逻辑结论, 分别。本文使用基于 KG 的事实验证和 KGQA 基准来评估 KG-GPT, 模型表现出有竞争力且稳健的性能, 甚至

胜过了一些全监督模型。因此，本文的工作在统一 LLMs 领域内的结构化和非结构化数据处理方面迈出了重要一步。

1251, 脑机接口的未来是光学的

The future of brain-machine interfaces is optical

论文作者: Nathan Tessema Ersaro, Cem Yalcin & Rikky Muller

机构: University of California

(北大刘志一摘录, 2023.2.17)

内容: 光学接口可用于解决脑机接口 (BMI) 开发中与缩放、精度和侵入性相关的挑战。脑机接口的发展是由神经传感和刺激模式的激增以及可以同时记录的神经元数量的增加所驱动的。理想情况下, 脑机接口应该是微创, 安全, 且不影响寿命的。它应当具有高时空分辨率和通信深度吞吐量, 以实现大脑编码和解码。

目前的脑机接口, 非侵入式 BMI 如脑电图与经颅磁刺激, 侵入式 BMI 如插入脑组织的微电极, 都存在不足之处。非侵入式 BMI 在通信复杂性方面具有根本限制, 因为测量构成了从远处记录的总和电信号, 并且因为刺激针对相对大量的组织。侵入式 BMI 受到电极阵列的寿命、可靠性和安全性受到插入过程中的物理损伤、对植入物的免疫反应以及随着时间的推移神经回路重组的限制。

光学可以在这里提供答案。基因编辑的出现使得位于神经元膜内的荧光电压指示器得以开发, 该指示器跟踪单个放电事件以及通过触发或抑制神经元放电来驱动称为视蛋白的蛋白质, 这些蛋白质对光做出反应因此, 通过使用颅窗, 可以在不同的光学波长上进行记录和刺激, 以避免读写串扰, 并且不会物理穿透脑组织。通过这种方法, 光学神经接口可以精确地控制和询问哺乳动物的神经回路, 创建 BMI, 可以实现能够驱动现实感官规则的刺激模式。

1252, 从黑盒到明盒: 使用深度人工神经网络的标量回归问题的假设检验框架

From black box to clear box: A hypothesis testing framework

for scalar regression problems using deep artificial neural networks

作者: Wolfgang Messner

机构: (美国)南卡罗来纳大学达拉摩尔商学院

（北大徐永登摘录，2023年8月12日）

内容：尽管深度学习在各个领域表现出令人印象深刻的预测性能，但其在社会和行为科学研究模型中的应用仍然受到限制。深度学习缺乏人类可理解的可解释性，并且不提供统计推论。为了解决这些限制，本文提出了一种新颖的与模型无关的假设检验框架，该框架使用深度神经网络针对标量回归问题进行了定制。新框架不仅确定每个输入变量的影响方向和统计显著性，还计算类似于传统普通最小二乘(OLS)回归模型中的Cohenf2的效应大小度量。效应大小是对零假设显著性检验的重要补充，它提供了独立于样本量考虑因素的实际显著性度量。为了展示新框架的实用性，将它的应用程序在人工数据集和使用Python沙盒实现的社会调查上进行了演示。

1253，利用局部可解释的模型无关 Shap 解释增强复杂机器学习模型的信任和可解释性

Enhancing trust and interpretability of complex machine learning models using local interpretable model agnostic shap explanations

作者：Sai Ram Aditya Parisineni & Mayukha Pal

机构：（美国）ABB 能力创新中心（ABB Ability Innovation Center）

（北大韩亚帅摘录，2023.10.25）

内容：随着人工智能（AI）技术和分析计算引擎的快速发展，机器学习（ML）模型变得越来越复杂。不幸的是，许多这些模型通常被视为黑盒子，缺乏用户可解释性。因此，理解和信任这些复杂 ML 模型所做的预测变得更加具有挑战性。然而，研究人员已经开发出各种框架，采用可解释 AI 方法来增强 ML 模型的可解释性和可解释性，从而提高其预测的可信度。在本研究中，我们提出了一种称为局部可解释的模型无关 Shap 解释（Local Interpretable Model Agnostic Shap Explanations, LIMASE）的方法。这种 ML 解释技术利用 LIME 范式内的 Shapley 值来实现几个目标：（a）通过利用局部准确而可解释的决策树模型来解释任何模型的预测。使用 Tree Explainer 来计算 Shapley 值，从而提供可视化可解释的解释，（b）通过为多个数据点绘制局部解释来提供可视化可解释的全局解释，（c）提供子模块化优化问题的解决方案，（d）提供区域解释的见解，（e）与

使用核心解释器相比，实现更快的计算速度。通过提出 LIMASE 方法，这项工作为 ML 模型的可解释性领域做出了贡献，并提供了解决复杂和不透明 ML 模型所带来的挑战的实际解决方案。所提出的方法使用户能够通过局部和全局层面上得到可视化可解释的解释，更深入地理解模型的预测。总体而言，本研究旨在弥合 ML 模型的复杂性与解释性需求之间的差距，最终增强 AI 驱动应用程序中的信任和可用性。

1254, 你的模型像工程师一样思考吗？基于深度学习的可解释人工智能轴承故障检测

Does Your Model Think Like an Engineer? Explainable AI for Bearing Fault Detection with Deep Learning

作者: Thomas Decker, Michael Lebacher, Volker Tresp

机构: 德国慕尼黑路德维希·马克西米利安大学、德国慕尼黑西门子股份公司
(北大李星辰摘录, 2023 年 10 月 23 日)

内容: 深度学习已经成功地应用于分析各种相关用例中的工业传感器数据。然而, 许多性能良好的方法的不透明性是实际部署的主要障碍。可解释的人工智能 (XAI), 特别是特征属性技术有望使人们能够深入了解这种模型如何形成其决策。但这种方法的简单应用通常无法为领域专家提供真正的信息和特定于问题的见解。这项工作专注于从振动信号中检测滚动轴承故障的具体任务, 提出了一种新颖的特定于领域的特征属性框架, 该框架允许评估模型的底层逻辑与专家推理的对应程度。利用该框架, 能够验证可信性, 并成功地预测不同性能良好的深度学习模型的泛化能力。该方法论展示了如何有效地使用信号处理工具来增强可解释的人工智能技术, 并充当类似问题的模板。

1255, 证明黑盒语言模型中的测试集污染

Proving test set contamination in black box language models

作者: Yonatan Oren*, Nicole Meister*, Niladri Chatterji*, Faisal Ladhak, Tatsunori B. Hashimoto

机构: (美国)Stanford University, Columbia University

（北大岳沈童摘录，2023 年 10 月 26 日）

内容：大型语言模型是在大量互联网数据的基础上训练出来的，这引发了人们对这些模型已经记住了公共基准测试的担忧和猜测。从猜测到证明污染具有挑战性，因为专有模型所使用的预训练数据通常不对外公开。作者的研究表明，有可能在不获取预训练数据或模型权重的情况下，为语言模型中的测试集污染提供可证明的保证。作者的方法利用了这样一个事实，即在没有任何数据污染的情况下，可交换的基准的所有排序应该是同等可能的。相反，语言模型记忆示例排序的倾向意味着受污染的语言模型会发现某些规范排序比其他规范排序更有可能出现。只要按规范排序的基准数据集的可能性明显高于打乱示例后的可能性，作者的测试就会标记潜在污染。作者证明了他们的程序足够灵敏，可以在具有挑战性的情况下可靠地证明测试集污染，包括小到 14 亿参数的模型、只有 1000 个示例的小型测试集以及在预训练语料库中只出现过几次的数据集。利用这个测试，作者审核了五种流行的可公开获取的语言模型的测试集污染情况，几乎没有发现普遍污染的证据。

1256，通过无创脑机接口和生成人工智能辅助的多模式软件记录梦境

Dream Recording Through Non-invasive Brain-Machine Interfaces and Generative AI-assisted Multimodal Software

论文作者：Todd Kelsey

机构：（美国）Allegheny College Department of Psychology

（北大刘志一摘录，2023. 4. 10）

内容：记录梦境一直是神经科学和心理学的长期挑战，该研究提出了一种新的梦境记录方法，将非侵入性脑机接口（BMI）、思维打字软件和生成人工智能辅助的多模式软件相结合。这种方法旨在将快速眼动睡眠期间将意识过程升华为半意识状态，并产生用于思维打字的信号。该文概述了一个分两个阶段的过程：首先，使用生成式人工智能开发多模式软件，以补充文本流并生成多媒体内容；其次，采用莫尔斯电码打字，以简化信号要求并提高打字速度。

BMI 研究证明了思维打字的潜力，用户可以通过无创 BMI 进行思维打字。另外，重复练习某些任务可以使有意识的过程升华为无意识或半意识状态。由以上

背景，通过思维打字的训练升华意识，可能有助于 REM 睡眠期间的信号生成。将莫尔斯电码用于思维打字的信号生成，可以降低信号要求的复杂性，并可以作为文字流的输入。这些生成的文字流，通过 gpt 等 AI 模型可以生成艺术、声音或电影序列，以此记录梦境。

无创 BMI、思维打字软件和生成人工智能的成功集成对社会有几个潜在好处：

1. 治疗应用。所提出的方法可用于治疗睡眠障碍、创伤后应激障碍和其他涉及梦境障碍的心理健康状况。
2. 创造性表达。基于思维类型文本生成多模式内容的能力可以为艺术和创造性表达开辟新的途径，使个人能够以创新的方式分享他们的梦想和想法。

1257, HappyFeat——面向临床应用的交互式高效 BCI 框架

HappyFeat - An interactive and efficient BCI framework for clinical applications

论文作者：Arthur Desbois, Tristan Venot, Fabrizio De Vico Fallani, Marie-Constance Corsi

论文机构：索邦大学（法国），巴黎脑研究所（法国），法国国家信息与自动化研究所（法国）

（北大张震扬摘录，2023.10.4）

论文内容：脑机接口（BCI）系统允许用户通过将大脑活动转化为命令来执行操作。为此，需要使用分类算法来区分不同的心理状态。这项技术可能的应用领域非常广泛，从通信到假肢控制和中风后康复。有多种 BCI 范式，例如 P300 和稳态视觉诱发电位（SSVEP）。文章关注运动意象（MI），因为它在前瞻性治疗用途（例如中风后康复）中很突出，这与研究者旨在改善 BCI 在临床环境中的使用的研究工作相一致。在 MI 中，受试者主动想象一个运动而不实际执行它，以便命令虚拟或真实设备，例如在屏幕上移动物体，控制机械臂。因此，MI 提供了高水平的交互性，并且需要强烈积极暗示。

此类系统通常需要一个训练阶段，包括训练分类算法以使用记录信号中的特定特征来区分心理状态。这一阶段的特征选择和训练对于 BCI 性能至关重要，并且提出了临床环境中需要满足的特定限制，例如中风后康复。文章介绍了

HappyFeat, 这是一款软件, 通过在一个方便的图形用户界面中收集所有必要的操作和分析, 并通过实验或分析参数的自动化, 使基于运动意象(MI)的 BCI 实验变得更加容易。由此产生的工作流程可以轻松选择最佳功能, 有助于在时间有限的环境中实现良好的 BCI 性能。可以使用基于功能连接的替代功能, 并将其与功率谱密度进行比较或组合, 从而允许采用面向网络的方法。文章详细介绍 HappyFeat 的主要机制, 并回顾其在典型用例中的性能。数据表明, HappyFeat 可以用作比较从信号中提取的不同指标的有效工具, 以训练分类算法。因此, 文章还对常用的功率谱密度和基于功能连接的网络指标进行了比较。

1258, 灵活的脑机接口

Flexible brain - computer interfaces

论文作者: Xin Tang (保尔森工程与应用科学学院, 哈佛大学), Hao Shen 保尔森工程与应用科学学院, 哈佛大学) 等,

内容: 本文探讨了柔性电子器件在脑机接口开发中的应用。脑机接口允许大脑与外部计算机直接通信, 在神经科学、医学和虚拟现实领域具有潜在的应用前景。然而, 目前的方法都是基于传统的刚性电子器件, 受限于其固有的机械和几何特性与脑组织不匹配。柔性电子器件具有与大脑兼容的机械特性, 可以解决这些局限性, 并可用于创建下一代脑机接口。在此, 我们探讨了柔性电子器件在脑机接口开发中的应用。我们研究了柔性、可伸缩和软电子器件在此类界面中的独特优势, 并考虑了该技术对神经科学、神经假体控制、生物电子医学以及脑与机器智能集成的潜在影响。我们还探讨了材料、设备制造和系统集成方面的挑战, 这些挑战是开发普遍适用的柔性脑机接口所必须解决的。

1259, 基于 FastICA、Hjorth 参数和 SVM 的 EEG-BCI 特征区分执行和想象动作
EEG-BCI Features Discrimination between Executed and Imagined Movements
Based on FastICA, Hjorth Parameters, and SVM

论文作者: Tat' y Mwata-Velu、Armando Navarro Rodríguez、Yanick Mfuni-Tshimanga

机构: (墨西哥) Centro de Investigación en Computación (CIC), 墨西哥城

Ciudad de México, Alcaldía Gustavo A. Madero 区

内容：研究通过分析脑信号来实现用户与其周围环境的交互。BCI 的准确性在很大程度上取决于用户如何执行不同的认知任务，尤其是在设备控制方面。该研究的主要目标是开发一种方法，通过分析脑电图（EEG）信号，能够准确地区分已执行和想象的运动任务，从而为 BCI 系统的应用提供更可靠的性能。

为实现这一目标，研究采用了一系列处理步骤。首先，他们使用了通带滤波器和快速独立成分分析（FastICA）算法来处理原始 EEG 信号，从中分离出独立的脑信号源。然后，他们利用 Hjorth 参数，包括活动性、移动性和复杂性，来提取 EEG 信号的特征。最后，使用支持向量机（SVM）作为分类器，对已执行和想象的运动任务进行区分。

研究指出，结果的好坏与所选择的特征组合、Hjorth 参数的应用，以及数据在训练、测试和评估集中的分割都有关。未来的工作将着重比较不同 Hjorth 算法与 FastICA 方法变体的性能，并计划在开发板上实施提出的模型，以探索嵌入式 BCI 系统的应用潜力。这个研究为 BCI 技术的进一步发展提供了有希望的方向，可以帮助更多人获得更准确的脑机接口体验。

1260, 基于异步脉冲神经网络的神经拟态计算系统

Enhancing Graph Neural Networks with Structure-Based Prompt

机构：（中国）华东师范大学

简介：图神经网络（GNN）在学习图数据的语义方面非常强大。最近，一种新的范式“预训练，提示”在使 GNN 适应监督数据较少的各种任务方面显示出有希望的结果。这种范式的成功可以归因于更一致的预训练和面向任务的提示调整的目标，其中预训练的知识可以有效地转移到下游任务中。然而，现有研究中一个被忽视的问题是，图的结构信息通常在学习节点表示的预训练阶段被利用，而在学习任务特定参数的提示调整阶段被忽略。为了弥合这一差距，我们为 GNN 提出了一种新的基于结构的提示方法，即 SAP，该方法在预训练和提示调整阶段始终如一地利用结构信息。特别是，SAP 1) 采用双视图对比学习来对齐节点属性和图结构的潜在语义空间，2) 在提示图中合并结构信息，以在提示调优中引出更多预先训练的知识。我们对节点分类和图分类任务进行了广泛的实验，以证明 SAP

的有效性。此外，我们表明 SAP 可以在同亲图和异亲图上更具挑战性的少数镜头场景中带来更好的性能。

1261, 用于快速探索的脉冲神经形态架构的性能和能量模拟

Performance and Energy Simulation of Spiking Neuromorphic Architectures for Fast Exploration

作者: James A. Boyle , Mark Plagge , Suma G. Cardwell , Frances S. Chance , Andreas Gerstlauer

机构: The University of Texas at Austin, Texas, USA 摘录, 2023 年 8 月

内容: 最近在神经形态计算领域, 已提出一系列新的脉冲神经网络 (SNN) 架构。然而, 神经形态设计缺乏一个能够促进不同 SNN 架构的探索并有助于早期设计决策的框架。虽然有各种 SNN 模拟器, 但没有一个可以快速估算不同脉冲架构的延迟和能量消耗。文章表明, 尽管当前的脉冲设计在实施上有所不同, 但它们具有可表示为通用架构模板的共同特征。文章描述了一个模拟各种神经形态架构的抽象时间步粒度的框架的初始版本。通过对 Intel 的 Loihi 平台建模来展示我们的模拟器, 对于不同规模的双层 SNN, 估计时间变化的能耗和延迟, 均方误差小于 10%。

1262, 脉冲神经网络的高效异步联合神经形态学习

Efficient asynchronous federated neuromorphic learning of spiking neural networks

论文作者: Yuan Wang, Shukai Duan, Feng Chen

论文机构: 西南大学人工智能学院

(北大向鼎摘录, 2023-10-07)

内容: 脉冲神经网络 (SNNs) 作为人工神经网络 (ANN) 的一种前景替代方案, 以其高能量效率和低计算成本相对于 ANN 引起了广泛关注。然而, 目前的 SNN 训练方法在处理大规模数据集和解决梯度消失问题方面仍然面临挑战。为了提高训练效率, 已经提出了多种直接监督学习规则, 但它们往往受限于浅层架构和神经元尖峰率不平衡的问题。另一方面, 联邦学习 (FL) 已经在多个领域展现了其

卓越的性能，但在资源受限、传输故障和异构计算能力等挑战下，FL 系统面临着诸多问题。该论文引入了 SNNs 到异步联邦学习框架，提出了 AdaFedAsy-SNN 算法，通过基于平均尖峰率和样本大小的权重聚合策略，以及异步框架，有效地解决了非 IID 场景下的挑战。

这项研究的贡献主要体现在三个方面：首先，成功地将 SNNs 引入到异步联邦学习中，为资源受限设备上的高效训练提供了新途径。其次，设计了一种融合权重策略，考虑了信息年龄和平均尖峰率，降低了模型陈旧性的影响，提高了训练效果。最后，通过在多个基准数据集上的实验，验证了 AdaFedAsy-SNN 算法在非 IID 设置下取得的卓越性能，为 FL-SNN 场景中的异质性问题提供了解决方案。这一研究的发现对于推动在大规模分布式系统中利用 SNNs 进行高效训练的应用具有重要意义。

1263, 基于知识图谱的多跳问答深度认知推理网络

Deep Cognitive Reasoning Network for Multi-hop Question Answering over Knowledge Graphs

作者: Jianyu Cai, Zhanqiu Zhang, Feng Wu, Jie Wang*

机构: 中国科学技术大学, 合肥综合性国家科学中心人工智能研究院

(北大孙宇瑶摘录, 2022 年 7 月 18 日)

内容: 知识图谱上的多跳问答, 旨在通过知识图谱中的推理路径找到给定问题的答案实体, 近年来引起了学术界和工业界的广泛关注。然而, 这项任务要在一个大的候选实体集中准确地识别答案, 候选集的大小随着推理跳数的增加而呈指数增长。为了解决这一问题, 本文提出了一种新的深度认知推理网络(DCRN)。具体来说, DCRN 包括两个阶段——无意识阶段和有意识阶段。无意识阶段首先通过利用候选实体的语义信息从候选实体中检索信息证据。然后, 意识阶段根据检索到的证据的图结构进行顺序推理, 准确地识别答案。

1264, GraphGPT: 大型语言模型的图指令调整

GraphGPT: Graph Instruction Tuning for Large Language Models

作者: Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng,

Dawei Yin and Chao Huang,

机构：中国 香港大学同心基金数据科学研究院； 香港大学计算机学院； 百度
（北大杨龙摘录，2023 年 10 月 19 日）

内容：图神经网络（GNN）通过图节点之间的递归信息交换和聚合，提升了对图结构的理解能力。为了提高模型的鲁棒性，自监督学习（SSL）已经成为一种有前途的数据增强方法。然而，现有的生成预训练图嵌入的方法通常依赖于与特定下游任务标签的微调，这限制了它们在标注数据稀缺或不可用的情况下的可用性。为了解决这个问题，本文的研究着重于提高图模型在具有挑战性的零样本学习场景中的泛化能力。受到大型语言模型（LLMs）的成功启发，本文旨在开发一种面向图的 LLM，即使在没有来自下游图数据的信息的情况下，也能在不同的下游数据集和任务之间实现高度泛化。在这项工作中，本文提出了 GraphGPT 框架，通过图指令调整范式将 LLMs 与图结构知识对齐。本文的框架包括一个文本-图关联组件，用于建立文本信息和图结构之间的联系。此外，本文提出了一个双阶段指令调整范式，并配备了一个轻量级的图文对齐投影仪。这个范式探索自监督的图结构信号和任务特定的图指令，以指导 LLMs 理解复杂的图结构，并提高它们在不同下游任务中的适应性。本文的框架在监督和零样本图学习任务上进行了评估，表现出卓越的泛化能力，并优于最先进的基准方法。

1265, 语义网络的设计表示

Design Representation as Semantic Networks

Serhad Sarica, Ji Han, Jianxi Luo

新加坡科技设计大学(Singapore University of Technology and Design)

（北大郭金伟摘录，2022 年 10 月 21 日）

内容：设计表示是设计过程中的一项常见任务，以促进学习、分析、重新设计、交流和其他设计活动。传统的表示技术依赖于人类的专业知识和手动构建，并且难以重复和扩展。本文提出了一种方法，利用预训练的大规模跨领域设计知识库，根据文本或自然语言的设计描述自动生成设计表示作为语义网络，即实体和关系的网络。本文的方法不需要临时统计数据。基于参与性研究，文中揭示了使用替代知识库自动生成的语义网络表示的有效性和差异。这些发现阐明

了未来的研究方向，以增强语义网络的设计表示。

1266, 给我讲一个故事！ 具有大型语言模型的叙事驱动 XAI

Tell Me a Story! Narrative-Driven XAI with Large Language Models

作者: David Martens, Camille Dams, James Hinns, Mark Vergouwen

机构: (比利时) 安特卫普大学

(北大徐永登摘录, 2023 年 9 月 29 日)

内容: 在当今的关键领域, 黑盒机器学习模型的主导地位扩大了对可解释人工智能(XAI)的需求。广泛使用的 SHAP 值虽然量化了特征重要性, 但往往过于复杂且缺乏人性化的解释。此外, 反事实(CF)解释提出了“假设”, 但让用户难以理解“为什么”。为了弥补这一差距, 我们引入了 XAIstories。利用大型语言模型, XAIstories 提供了阐明 AI 预测的叙述: SHAPstories 基于 SHAP 解释来解释预测分数, 而 CFstories 则基于 CF 解释来解释决策。我们的结果令人震惊: 超过 90%的受访普通观众认为 SHAPstories 生成的叙述令人信服。数据科学家主要看到 SHAPstories 在向普通受众传达解释方面的价值, 92%的数据科学家表示, 它将有助于非专业人士轻松、自信地理解 AI 预测。此外, 83%的数据科学家表示他们可能会使用 SHAPstories 来实现此目的。在图像分类中, CFstories 被认为更具说服力或同样有说服力, 因为超过 75%的非专业参与者都自己制作了故事。与手动创建的叙述相比, CFstories 还可以将创建叙述的速度提高十倍, 并将准确性提高 20%以上。因此, 结果表明 XAIstories 可能提供真正解释和理解人工智能预测中缺失的环节

1267, 可解释人工智能 (XAI) 2.0: 开放挑战和跨学科研究方向的宣言

Explainable Artificial Intelligence (XAI) 2.0: A Manifesto of Open Challenges and Interdisciplinary Research Directions

作者: Luca Longo, Mario Brcic, Federico Cabitza, Jaesik Choi, Roberto Confalonieri, Javier Del Ser, Riccardo Guidotti, Yoichi Hayashi, Francisco Herrera, 等

机构: (爱尔兰) 都柏林科技大学人工智能与认知负荷研究实验室

（北大韩亚帅摘录，2023.10.30）

内容：随着基于不透明人工智能（AI）的系统在各种实际应用中不断蓬勃发展，理解这些黑盒模型变得至关重要。作为回应，可解释的人工智能（XAI）已经成为一个研究领域，在各个领域中具有实际和伦理上的益处。本文不仅突出了 XAI 的进展及其在实际场景中的应用，还解决了 XAI 中持续存在的挑战，强调了对更广泛的视角和合作努力的需求。我们汇集了来自不同领域的专家，以识别开放性问题，努力同步研究议程并加速 XAI 在实际应用中的推进。通过促进协作讨论和跨学科合作，我们旨在推动 XAI 向前发展，为其持续成功做出贡献。我们的目标是提出一个全面的提案，推进 XAI 的发展。为了实现这一目标，我们提出了一个包含九个类别的 27 个开放性问题的宣言。这些挑战概括了 XAI 的复杂性和微妙之处，并为未来的研究提供了一份路线图。对于每个问题，我们提供了有希望的研究方向，希望能够利用感兴趣的利益相关者的集体智慧。

1268, 使用深度强化学习学习顺序反事实解释的公正策略

Learning impartial policies for sequential counterfactual explanations using Deep Reinforcement Learning

作者：Emmanouil Panagiotou, Eirini Ntoutsi

机构：（德国）柏林弗雷大学数学与计算机科学系，

（德国）慕尼黑大学国防军信息学院

（北大李星辰摘录，2023 年 11 月 1 日）

内容：在可解释人工智能（XAI）领域，序列反事实（SCF）示例通常用于通过对输入实例进行一系列修改来改变训练分类器的决策。尽管某些测试时间算法旨在针对每个新实例单独进行优化，但最近已经提出了强化学习方法，该方法寻求学习用于发现 SCF 的策略，从而增强可扩展性。正如 RL 中的典型情况一样，RL 问题的公式，包括状态空间、动作和奖励的规范，往往是模糊的。现有方法中的缺陷可能导致策略具有不期望的属性，例如对特定操作的偏见。本文建议使用分类器的输出概率来创建更具信息性的奖励，以减轻这种影响。

1269, TopicGPT: 基于提示的主题建模框架

TopicGPT: A Prompt-based Topic Modeling Framework

作者: Chau Minh Pham, Alexander Hoyle, Simeng Sun, Mohit Iyyer

机构: (美国) 马萨诸塞大学阿默斯特分校, (美国) 马里兰大学

(北大岳沈童摘录, 2023 年 11 月 2 日)

内容: 传统的主题模型, 如潜在 Dirichlet 分配将文档表示为主题的混合体, 其中每个主题都是单词的分布。主题通常用最有可能出现的单词表示, 但这种表示可能包含不连贯或不相关的单词, 使用户难以解释主题。虽然有些模型能让用户根据需求和领域知识交互式地引导主题, 但它们的可用性受到字袋主题格式的限制。

为了解决这些局限性, 作者设计了 TopicGPT, 这是一种以人为中心的话题建模方法, 它依靠提示大型语言模型来执行上下文话题生成和分配。首先, 根据输入数据集中的文档样本和之前生成的主题列表, 反复提示 LLM 生成新主题。然后进一步完善生成的主题集, 以整合冗余主题并消除不常见的主题。最后, 给定一个新文档, LLM 会将其分配给一个或多个生成的主题, 同时提供文档中的引文来支持其分配。这些引文使该方法易于验证, 解决了困扰传统主题模型的一些有效性问题。

1270, 基于脑机接口的康复训练在脑卒中上肢康复中的研究进展

论文作者: 李丹, 刘玲玉, 靳令经, 王茹

机构: (中国) 上海体育大学, 同济大学附属养志康复医院神经康复科

(北大刘志一摘录, 2023. 10. 25)

内容: 近年来, 基于 BCI 的多种康复训练系统在脑卒中后上肢康复领域发展迅速。BCI 常与不同外部设备联合进行训练, 通过采集脑电信号、计算机处理解码患者意图, 将信息转换为有效的命令信号输出至外部设备, 外部设备执行指令并提供本体感觉反馈或视觉反馈, 形成主动式闭环反馈回路帮助患者调节神经活动, 从而促进神经可塑性和肢体运动功能恢复。目前与 BCI 结合的常用的外部设备有功能性电刺激(FES)、上肢康复机器人和视觉反馈设备。

FES 应用一定强度的低频脉冲电流, 按需编定程序, 作用于神经肌肉系统, 以产生

运动或模仿正常自主运动,从而补偿或替代受损个体已丧失的功能。上肢康复机器人是交互式电动设备,可用于移动肢体,辅佐肢体运动。视觉反馈训练能给予患者丰富的视觉和日常生活情景模拟,对防治废用综合征及加速上肢运动功能恢复具有重要意义。

目前普遍认为 BCI 训练影响脑卒中后上肢运动功能康复的机制与神经可塑性密切相关,但仍不十分明确,可能主要包括以下三个方面:皮质激活状态改变,神经网络功能连接的改变,以及皮质脊髓束完整性改变。

基于 BCI 的康复训练为患者提供更好的康复体验,真正地得到推广应用,仍有几个关键问题亟待解决。首先是需要患者花费较长时间学习掌握 MI 技巧,认知能力欠佳、体质较差等主观问题可能会影响训练的疗效。其次,缺乏标准的和个体化的临床实施方案。在康复训练过程中,应根据特定肌群、任务模式和患者自身情况选择或调整至最佳的训练方式和参数,为患者提个性化、最优化的治疗方案。另外,未来仍需要更多高质量研究进一步探究 BCI 训练促进脑卒中上肢康复的作用机制。

1271, 顺序最佳臂识别及其在脑机接口中的应用

Sequential Best-Arm Identification with Application to Brain-Computer Interface

作者: Xin Zhou, Botao Hao, Jian Kang, Tor Lattimore, Lexin Li

机构: (美国)密歇根大学, (美国)加州大学伯克利分校

(北大张震扬摘录, 2023. 5. 17)

论文内容: 脑机接口 (BCI) 是一项突破性技术,可以实现大脑与外部设备或计算机系统之间的直接通信。它涉及使用各种传感器,例如脑电图 (EEG)、皮质电图 (ECoG) 或植入式神经电极,检测并记录大脑产生的电信号。然后,这些信号由机器学习算法处理,以解释和提取有意义的命令和意图。BCI 允许个人仅用他们的思想与设备进行交互,并且在医学、康复和人体增强方面具有广泛的应用潜力。

P300 拼写器是一种 BCI 系统,允许用户在计算机屏幕上选择字符或拼写单词,而无需使用物理键盘,而是使用大脑信号。它基于 P300 事件相关电位 (ERP),

这是一种大脑反应，以特定的电压波动模式的形式出现，在出现相关刺激后大约 300 毫秒发生。刺激通常是在计算机屏幕上以网格状布局闪烁的单个字符或符号，ERP 由头皮脑电图帽或类似设备检测和记录。出现刺激后，通过信号处理和机器学习算法在固定时间窗口内对电极捕获的脑电信号进行分析，检测 P300 反应的发生并确定目标字符。传统的非自适应范式独立对待每个单词选择，导致机器学习的时间非常漫长。为了提高采样效率，作者将问题转化为多臂老虎机中的一系列最佳臂识别任务。利用预先训练的大型语言模型（LLM），利用从先前任务中学到的先验知识来通知和促进后续任务。为了以连贯的方式做到这一点，文章中提出了在固定置信设置和固定预算设置下的顺序前二汤普森采样（STTS）算法。文章结尾研究了所提出算法的理论特性，并通过综合数据分析和 P300 拼写器模拟器示例证明了顺序最佳臂识别对其有实质性的改进。

1272, 脑机接口技术的最新进展

State-of-the-Art on Brain-Computer Interface Technology

作者: Dr. Jānis Pekšs (拉脱维亚图里巴大学), Dr. Dmytro Mamchur (拉脱维亚图里巴大学)

(北大朱小龙摘录, 2023. 11. 5)

内容: 本文全面概述了脑机接口 (BCI) 的最新进展。文章首先介绍了 BCI, 描述了其主要工作原理和最广泛使用的平台。然后, 论文探讨了 BCI 系统的各个组成部分, 如硬件、软件和信号处理算法。最后, 本文探讨了当前与医疗、教育和其他用途的 BCI 使用相关的研究趋势, 以及该技术未来的潜在应用。论文最后强调了在广泛应用之前仍需解决的一些关键挑战。通过对生物识别 (BCI) 技术的最新评估, 本文将为该领域的进步和创新提供有价值的见解。

脑机接口 (BCI) 是一项快速发展的技术, 有可能彻底改变人类与计算机的交互方式。BCIs 可测量大脑活动, 并将其转化为计算机或其他设备的指令, 从而让用户仅凭意念就能控制机器和设备。从可移动的机器人蜘蛛和小球到更多的实际应用, 神经小工具正越来越多地被用于娱乐目的。然而, 更重要的是, 神经小工具还被开发用于帮助残疾人, 如肢体瘫痪者。

根据动作方向, BCI 通常分为单向和双向两类。单向 BCI 要么从大脑接收信号,

要么向大脑发送信号，而双向 BCI 允许双向信息交换，从而实现大脑对外部设备的控制。对反馈方法的研究仍在继续，目的是开发出能将外部命令转化为通过神经系统传输的电信号的技术。例如，它可用于对脊髓损伤患者的腿部肌肉进行电刺激，使他们能够通过平板设备控制自己的动作，从而恢复行动能力。神经网络和其他学习算法在信号处理中的应用非常普遍，因为不同个体的大脑活动各不相同。因此，这些系统需要长时间的训练，才能使生物识别（BCI）准确解读特定用户的指令。训练时间的长短取决于生物识别系统接收指令的数量。

1273, Joie:一种基于愉悦情感的脑机接口

Joie: a Joy-based Brain-Computer Interface (BCI)

作者: Angela Vujic、Ashley Martin、Shreyas Nisal

机构: 美国国立卫生研究院 (National Institutes of Health, 简称 NIH)

(北大胡真茂摘录, 2023. 10. 29)

内容: 这项研究关注脑电图 (EEG) 头戴设备的发展, 它们的尺寸和成本不断减小。研究团队开发了名为 Joie 的脑-计算机接口 (BCI), 以前额皮质不对称性作为输入源, 旨在捕捉愉快思绪并将其应用于一种无尽奔跑游戏。

研究发现, 实验组使用积极的心理策略能够显著增强前额不对称性, 与安慰剂组和对照组相比表现更出色。在实验中, 参与者经历了头部准备, 玩了一轮 Joie 游戏, 然后填写了多个后续调查, 包括感知到的快乐程度、注意力水平、使用的心理策略、NASA 任务负荷指数 (TLX) 调查和面试。实验还分析了静息不对称性的变化, 这被认为是神经反馈效果的指标。结果显示, BCI 训练在改变参与者的脑活动方面发挥了关键作用。此外, 研究通过 NASA-TLX 调查揭示了任务的心理需求和努力水平, 以及一定程度的挫折感。尽管任务具有挑战性, 但参与者通过训练能够提高静息不对称性, 表明他们在神经反馈方面取得了成功。此外, 参与者通过“愉快”的心理想象练习获益, 并在任务中感受到“快乐”情感。

这项研究为情感 BCI 技术的发展提供了重要见解, 将前额不对称性用于游戏, 为神经反馈的成功提供了有力支持。未来的研究可以进一步探索如何减轻任务对参与者的心理负担, 并比较“愉快”的心理想象任务与其他心理想象任务的效果, 以深入挖掘这一领域的潜力和应用前景。

1274, 更好地合作: 用语言模型和邻域信息增强生成知识图谱的完成

Better Together: Enhancing Generative Knowledge Graph Completion with Language Models and Neighborhood Information

作者: Alla Chepurova, Aydar Bulatov, Yuri Kuratov, Mikhail Burtsev

机构: 俄罗斯神经网络与深度学习实验室, 英国伦敦兰登数学科学研究所

(北大顾翔云摘录, 2023 年 11 月 2 日)

内容: 真实世界知识图谱 (KGs) 经常存在不完全性, 这限制了其潜在的性能。知识图谱完成 (KGC) 技术旨在解决这个问题。然而, 传统的 KGC 方法在计算上是密集的, 并且对于大规模的 KGs 是不切实际的, 因此需要学习密集节点嵌入和计算成对距离。基于生成转换器的语言模型 (例如 T5 和最近的 KGT5) 提供了一个很有前途的解决方案, 因为它们可以直接预测尾部节点。在这项研究中, 将节点邻域作为附加信息, 以改进基于语言模型的 KGC 方法。我们检验了这种插补的效果, 并表明, 在归纳和转导 Wikidata 子集上, 这个方法优于 KGT5 和传统的 KGC 方法。对邻域对模型预测的影响进行了广泛的分析, 并表明了其重要性。此外, 指出了通过更有效的邻域选择来显著提高 KGC 的方法。

1275, 突触运动适应: 尖峰神经网络中自适应机器人控制的三因素学习规则

Synaptic motor adaptation: A three-factor learning rule for adaptive robotic control in spiking neural networks

作者: Samuel Schmidgall, Joe Hays

机构: (美国) 海军研究实验室, 约翰霍普金斯大学

(北大邓凯元摘录, 2023 年 8 月)

内容: 在现实环境中运行的腿式机器人必须具备快速适应意外条件的能力, 例如不断变化的地形和变化的有效负载。本文介绍了突触运动适应 (SMA) 算法, 这是一种通过利用神经科学衍生的突触可塑性规则和三因素学习来实现四足机器人实时在线适应的新方法。为了促进快速适应, 我们通过梯度下降对三因素学习规则进行元优化, 以仅使用本地可访问的机载传感数据来近似由特权信息产生的嵌入来适应不确定性。我们的算法的性能与最先进的运动适应算法类似, 并为利用神经形态硬件实现自适应机器人提供了一条清晰的道路。

1276, Adaptive-SpikeNet: 使用具有可学习神经元动力学的脉冲神经网络进行基于事件的光流估计

Adaptive-SpikeNet: Event-based Optical Flow Estimation using Spiking Neural Networks with Learnable Neuronal Dynamics

作者: Adarsh Kumar Kosta, Kaushik Roy

机构: (美国) 普渡大学, 西拉斐特

(北大向鼎摘录, 2023-07-04)

内容简介: 近年来, 基于事件的摄像头因其异步捕捉丰富时间信息的能力, 在高速运动估计领域展现出巨大潜力。脉冲神经网络 (SNNs) 以其仿生的事件驱动处理方式, 能够高效处理这种异步数据, 而 leaky-integrate 和 fire (LIF) 等神经元模型可以跟踪输入中包含的典型时序信息。SNNs 通过在神经元存储器中保持动态状态来实现这一点, 保留重要信息, 同时随着时间的推移忘记冗余数据。因此, 作者认为与大小相似的模拟神经网络 (ANNs) 相比, SNNs 在序列回归任务上可能表现更好。然而, 由于后续层的脉冲消失问题, 深层 SNNs 难以训练。为此, 该论文提出了一种自适应的完全脉冲框架, 其中包含可学习的神经元动力学, 以缓解脉冲消失问题。该论文利用基于替代梯度的时空反向传播 (BPTT) 方法, 从零开始训练深层 SNNs。该论文在 Multi-Vehicle Stereo Event-Camera (MVSEC) 数据集和 DSEC-Flow 数据集上验证了我们的方法, 用于光流估计任务。该论文中的实验证明, 在这些数据集上, 与最先进的 ANNs 相比, 论文中的方法平均降低了约 13% 的平均端点误差 (AEE)。该论文还探讨了几个缩小模型, 并观察到 SNN 模型始终表现优于大小相似的 ANNs, 降低了约 10% 至 16% 的 AEE。这些结果表明了对于较小的模型, SNNs 的重要性以及它们在边缘计算中的适用性。在效率方面, 该论文中的 SNNs 在网络参数 (约 48.3 倍) 和计算能耗 (约 10.2 倍) 方面节省了大量资源, 并且与最先进的 ANN 实现相比, EPE 降低了 10%。

1277, COKE: 机器心智理论的认知知识图谱

COKE: A Cognitive Knowledge Graph for Machine Theory of Mind

作者: Jincenzi Wu*, Zhuang Chen*, Jiawen Deng, Sahand Sabour, Minlie Huang

机构：（中国）清华大学北京信息科学与技术国家研究中心，中国科学院人工智能研究所，智能技术与系统国家重点实验室

（北大孙宇瑶摘录，2023年5月9日）

内容：心智理论(Theory of mind, ToM)指的是人类理解和推断他人的欲望、信念和意图的能力。ToM的习得在人类的社会认知和人际关系中起着关键作用。虽然ToM对于社会智能是不可或缺的，但对于现代AI和NLP系统来说仍然是缺乏的，因为它们无法访问训练语料库下的人类心理状态和认知过程。

为了使人工智能系统具有ToM能力并缩小它们与人类之间的差距，本文提出了COKE:机器心智理论的第一个认知知识图。具体来说，COKE将ToM形式化为45k+人工验证的认知链的集合，这些认知链表征了人类在面对特定社会环境时的心理活动和随后的行为/情感反应。在此基础上利用预训练好的语言模型对COKE进行进一步泛化，构建了一个强大的认知生成模型COKE+。自动评价和人工评价的实验结果表明，COKE+具有良好的质量和ToM能力。

1278, LLMRec: 基于图增强的大型语言模型用于推荐

LLMRec: Large Language Models with Graph Augmentation for Recommendation

作者：Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, Chao Huang

机构：（中国）香港大学;百度

（北大杨龙摘录，2023年11月1日）

内容：长期以来，数据稀缺问题一直是推荐系统面临的挑战，先前的研究尝试通过整合辅助信息来解决这个问题。然而，这种方法通常会引入一些副作用，比如噪音、可用性问题和低质量数据，进而影响了对用户偏好的准确建模，对推荐性能产生负面影响。鉴于最近大型语言模型(LLMs)的发展，它们拥有丰富的知识库和强大的推理能力，本文提出了一种名为LLMRec的新框架，通过采用三种简单而有效的基于LLM的图增强策略来增强推荐系统。我们的方法利用在线平台（如Netflix、MovieLens）中丰富的内容，以三种方式增强交互图：（一）强化用户-物品交互边缘，（二）增强对物品节点属性的理解，（三）从自然语言的角度直观地进行用户节点剖析。通过采用这些策略，我们解决了推荐系统中稀

疏的隐式反馈和高质量辅助信息所带来的挑战。此外，为确保增强的数据质量，本文开发了一种去噪数据强化机制，包括噪声隐式反馈修剪和基于 MAE 的特征增强技术，有助于精细化增强数据并提高其可靠性。此外，本文提供了理论分析，以支持 LLMRec 的有效性，并澄清了本文的方法在促进模型优化方面的好处。在基准数据集上的实验结果显示，本文基于 LLM 的增强方法优于最先进的技术。

1279, 探索用于补全知识图谱的大语言模型

Exploring Large Language Models for Knowledge Graph Completion

作者: Serhad Sarica, Ji Han, Jianxi Luo

机构: (中国) 腾讯(tencent)、(美国) 西北大学

(北大郭金伟摘录, 2023 年 8 月)

内容:知识图谱在众多人工智能任务中发挥着至关重要的作用，但它们经常面临不完全性问题。本文探索利用大型语言模型 (LLM) 来补充知识图谱。文中将知识图谱中的三元组视为文本序列，并引入了一个名为知识图谱 LLM (KG-LLM) 的创新框架来对这些三元组进行建模。本文使用三元组的实体和关系描述作为提示，并利用响应进行预测。在各种基准知识图上的实验表明，该方法在三重分类和关系预测等任务中取得了最先进的性能。本文还发现，微调相对较小的模型 (例如，LLaMA-7B、ChatGLM-6B) 优于最近的 ChatGPT 和 GPT-4。在此研究中，提出了一种新的 KG 完成方法，称为 KG-LLM。本文提出的方法在 KG 完成任务 (如三重分类和关系) 方面取得了最高的性能。

1280, 评论 OpenAI 发生“大地震”

11 月 21 日 陆主席主持召开 COPU 例会。本会首先讨论 11 月 18 日官宣 Open AI (公司) 发生的“大地震”。

谈论“OpenAI 大地震”的版本很多，其中比较权威的是张侃给我们发来由华尔街俱乐部于 11 月 20 在美国发表的版本“Open AI 政变始末”。“大地震”发生的原因是 Open AI 不成熟的治理结构，暴露了 Open AI 的非赢利与赢利双层架构之间的潜在矛盾。现有非赢利部分的董事会 (无人持股)，6 名董事采取一人一票制 (潜藏不稳定风险)。

6名董事为：Open AI 联合创始人、CEO 奥特曼（Sam Altman），联合创始人、董事长兼总裁布洛克曼（Greg Brockman），联合创始人、首席科学家苏茨克沃（Ilya Sutskever），其他三名独立董事为：Adam D’ Angelo, Tasha McCauley, Helen Toner（托尼）。

奥布派开始全力筹资、扩充资源加快大模型 GPT 发展、扮现赢利派的角色，苏托派（苏茨克沃+托尼）坚持 Open AI 非盈利、强调安全的基本价值观。11月17日，他们（苏托派）联合其他两位独立董事向奥布派发动突然袭击，以奥氏团队研究方向对董事会汇报不坦诚为由，解除奥特曼的职务，解除布洛克曼董事长职务（只保留总裁职务），任命由原 CTO Mire Murati 为临时代理 CEO（后又聘请亚马逊的 Emmett Shear 任 CEO，接替 Mira Murati）。

现在我们关心的是“Open AI 大地震”将产生什么影响？影响是巨大的，但大模型 GPT 技术发展总趋势（通向生成式通用人工智能/AGI）不会变，大模型企业走开源的技术路线不会变，但经营方式将向赢利和非赢利方式两极分化，大模型 MPT 一旦超越人类，奥特曼允诺优先解决其安全问题留下未定的疑案。陆主席认为：关键在于治理模式和业务模式（盈利和非盈利）均不成熟。

后记：王东临给我们发来一篇报导：《本剧落幕了》，Open AI 已任命新的 CEO 和管理层，Sam Altman、Greg Brockman 带着一群人加入了微软，微软或成为大赢家。

刘澎秘书长又谈：有 95%的员工以联名支持奥特曼，要求 Open AI 现任董事会下台。令人意外的 Ilya Sutskever 也签了名。（有人认为他“跳反了”，挨骂！）。陆主席说，看来不能说“本剧落幕了”，Open AI 改组势在必行！谷歌、微软、Elon Musk 等在对待 open AI 改组中的态度有待观察。

近日来自 OpenAI 的信息比较混乱，我们选择 11 月 21 日 OpenAI 在 X 平台上的发帖称：我们已经达成原则性一致，Sam Altman 将回到 OpenAI 担任 CEO，并由 Bret Taylor 担任董事会主席，Larry Summers 和 Adam D’ Angelo 组成新的初始董事会。

（先报导到至此）