



清华大学 交叉信息研究院

Institute for Interdisciplinary Information Sciences, Tsinghua University

学术科研简报

IIIS Academic Newsletter

NATURE: 突破量子纠错盈亏平衡点

孙麓岩研究组通过实时重复的量子纠缠技术延长了量子信息的存储时间，在国际上首次超越盈亏平衡点，迈出了实用化可扩展量子计算的关键一步。

FOCS: 更快的矩阵乘法算法

段然研究组给出了更快的矩阵乘法算法，新的复杂度为 $O(n^{\wedge}2.371866))$ ，改进了之前的复杂度 $O(n^{\wedge}2.372860))$ ，为近十年来最大的改进幅度。

ICML: 使用范畴论刻画大模型的能力边界

袁洋研究组引用范畴论作为理论工具，重新审视了监督学习的理论框架，并且针对预训练模型证明了三个定理。

2023年1月-6月

目录

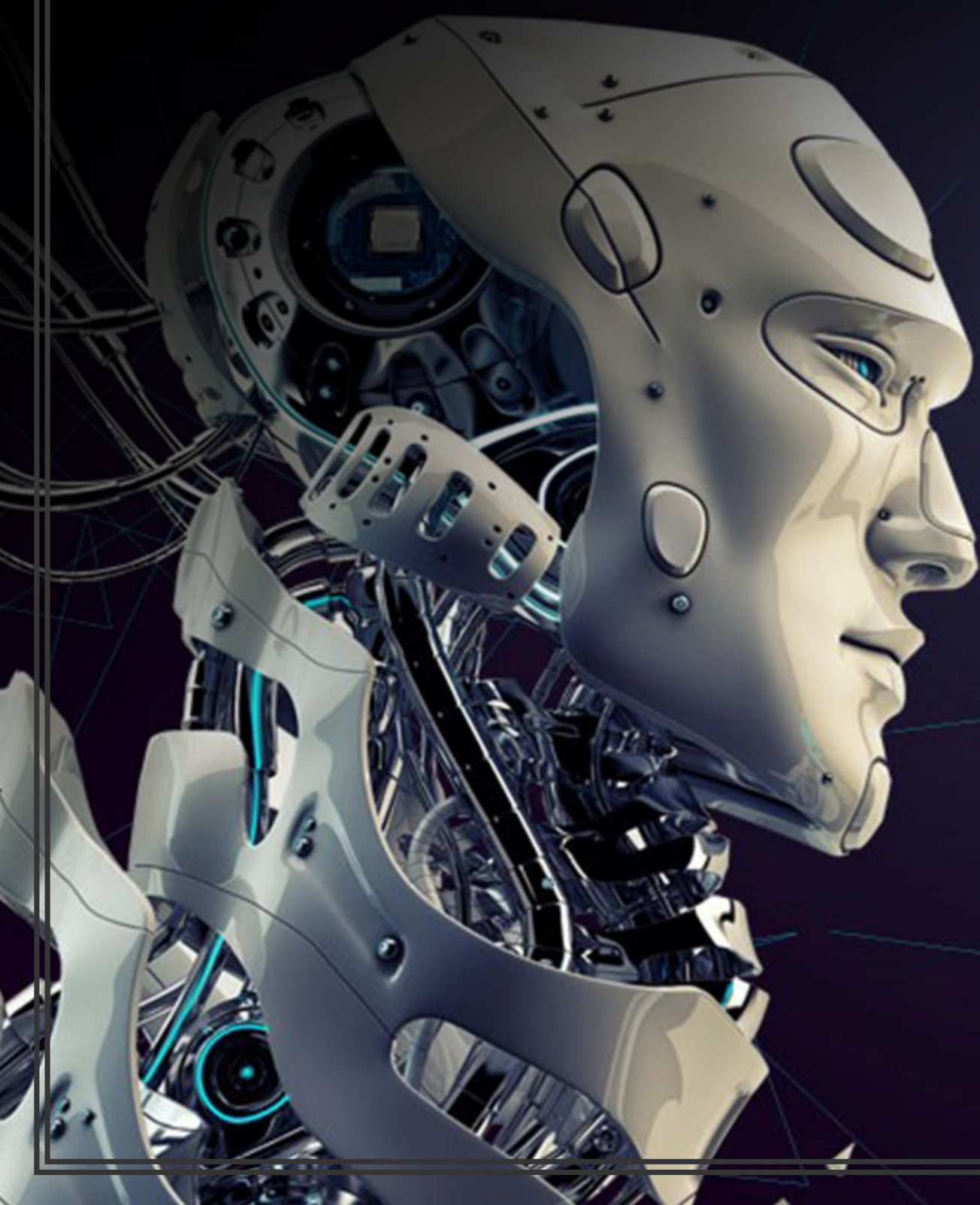
01 人工智能

人工智能理论	04
机器学习	07
计算机视觉	18
自然语言处理	23
计算机图形学	27
脑启发人工智能	28
算法理论	29
计算机系统结构	31
数据库系统	38
区块链	40

02 量子信息

离子阱量子模拟	42
量子计算与通信	46
超导量子计算	51
量子多体物理	55
凝聚态物理学	56

【人工智能】



一、人工智能理论

主要完成人：袁洋研究组、李建研究组、张景昭研究组

使用范畴论刻画大模型的能力边界

假如人类有无限的资源，比如有无穷多的数据，无穷大的算力，无穷大的模型，完美的优化算法与泛化表现，请问由此得到的预训练模型是否可以用来解决一切问题？这是一个大家都非常关心的问题，但已有的机器学习理论却无法回答。它与表达能力理论无关，因为模型无穷大，表达能力自然也无穷大。它与优化、泛化理论也无关，因为研究人员假设算法的优化、泛化表现完美。换句话说，之前理论研究的问题在这里不存在了。

该文引入了范畴论作为理论工具，针对预训练任务进行重新建模，构建了预训练任务与范畴内部结构的等价关系。从这个角度出发，该文重新审视了监督学习的理论框架，并且针对预训练模型证明了三个定理。

第一个定理证明了，如果使用提示调优的方式，预训练模型的能力和任务结构有关。一个任务能够被解决，当且仅当该任务能够被范畴中的某个对象表出。

第二个定理证明了，如果使用微调的方式，预训练模型的能力不再受范畴内对象表出能力的限制。预训练模型得到的特征向量可以完美地保留原范畴的信息，在使用高质量的训练数据、充足算力的前提下，预训练模型有潜力解决各种任务。

第三个定理证明了，基于源范畴中对象的结构，预训练模型天然拥有在目标范畴中生成从未见过的对象的能力。

该成果研究论文：Yang Yuan, “On the Power of Foundation Models”, ICML 2023.

稀疏矩问题的有效算法

李建研究组研究了从任意维度的噪声矩信息中学习高维概率分布的稀疏矩问题。针对该问题，以前的算法要么假设某些分离条件，使用更多的矩信息，要么在（超）指数时间内运行。针对一维问题（也称为稀疏 Hausdorff 矩问题），该研究组基于经典 Prony 方法设计了其鲁棒版本，并进行了严格的全局分析。在算法分析中，该研究组利用 Vandermonde 矩阵定义的线性系统和 Schur 多项式之间的联系，这使该研究组能够得到更紧的界。针对高维问题，该研究组首先通过将一维算法和分析扩展到复数来解决二维问题，然后将分布的一维投影和一组二维投影对齐来确定每个尖峰的坐标。该研究组的理论结果可应用于学习主题模型 (topic modeling) 和高斯混合 (Gaussian mixture)，并改进了前人工作的样本复杂度或运行时间。

该成果研究论文：Zhiyuan Fan, Jian Li, “Efficient Algorithms for Sparse Moment Problems without Separation”, COLT 2023.

具有对抗性干扰的连续时间在线控制算法

控制系统理论是一个重要的研究方向，受到人们的广泛关注。最近这方面的研究主要集中于设计在线控制算法，目标为选择尽可能鲁棒的在线控制器，对抗系统生成的噪声，与最优的线性控制器相比，提供遗憾上界。然而，目前这个领域主要的研究工作内容大多为离散时间线性系统，很少有专门研究连续时间线性系统的分析。

张景昭研究组提出了具有有限采样率的连续时间线性系统的在线控制算法，其中目标是设计一个连续系统的在线控制器，在任意噪声下学习并达到尽量好的效果。研究组提出了一个两级在线控制器来解决这个问题。上级控制器象征着策略学习过程，并以较低频率更新策略，尽量减少遗憾。下级控制器提供高频反馈控制输入，以减少离散化误差。

该方法为在线控制提供了实用且鲁棒的解决方案，实现了亚线性遗憾。更重要的是，该研究的分析表明在线学习算法，通过潜在的调整，也可能有益于连续时间控制问题。相信这个方向具有进一步探索的潜力。

该成果研究论文: Jingwei Li, Jing Dong, Baoxiang Wang, Jingzhao Zhang, "Online Control with Adversarial Disturbance for Continuous-time Linear Systems", arXiv:2306.01952.

二、机器学习

主要完成人：李建研究组、高阳研究组、吴翼研究组、许华哲研究组

OpenFE: 全自动特征生成器

表格类数据 (tabular data) 是机器学习应用中最常见的数据类型之一。该类数据也常见于各类算法和机器学习竞赛, 如 Kaggle 比赛中。特征工程 (特征生成) 是提升表格类数据上机器学习模型能力的重要环节, 对模型的最终效果有着重要影响。基于原始数据构建新的有效的特征可以大幅提升模型的泛化效果。但是对于机器学习工程师而言, 如何去构造新特征, 如何去筛选出有效的新特征常常是一个有挑战的问题, 需要大量的尝试和经验。在很多机器学习任务和比赛中, 特征工程通常会花去建模总时间的一半甚至更多。尽管目前也有各种开源或商用自动化机器学习 (AutoML) 工具, 但是这些工具也都没有把自动化特征生成加入到流程中。同时, 目前也很难找到高效的自动化特征生成开源包。

针对目前的问题, 李建研究组设计了 OpenFE。OpenFE 是一个自动化特征生成的框架, 可以高效地构建有效的新特征, 显著地提升 GBDT (例如 LightGBM, XGBoost 等) 和各种 SOTA 神经网络 (例如 Transformer, AutoInt, TabNet 等) 在表格类数据的效果。该研究组在 Kaggle 比赛 IEEE-CIS Fraud Detection 上验证了 OpenFE 的效果, 只需要一个简单的 XGBoost 模型加上 OpenFE 生成的特征就可以打败 99.3% (42/6351) 的参赛队伍。该研究组将 OpenFE 自动生成的特征和该比赛排名第一的队伍的生成的特征对比 (使用同样的基础 XGBoost 模型), OpenFE 自动生成的特征可以带来更大的效果提升。该成果发表于机器学习顶级会议 ICML, 相关代码在 Github 开源。

该成果研究论文: Tianping Zhang, Zheyu Zhang, Zhiyuan Fan, Haoyan Luo, Fengyuan Liu, Qian Liu, Wei Cao, Jian Li,

“OpenFE: Automated Feature Generation with Expert-level Performance”, ICML 2023.

Table 4: Results of OpenFE and Expert (feature generation by experts) in two Kaggle competitions. Notation: pub. ~ public, pri. ~ private. The final standing is determined according to the scores in the private leaderboard.

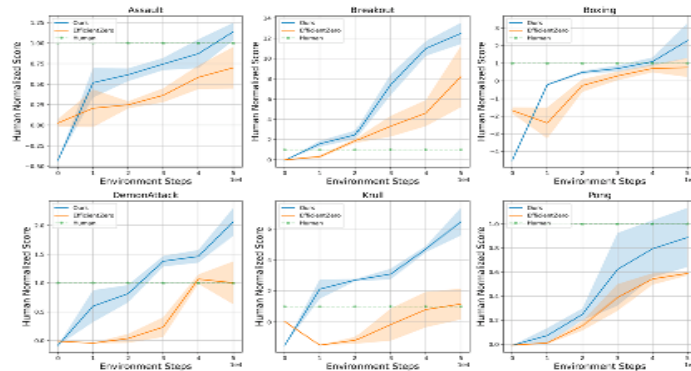
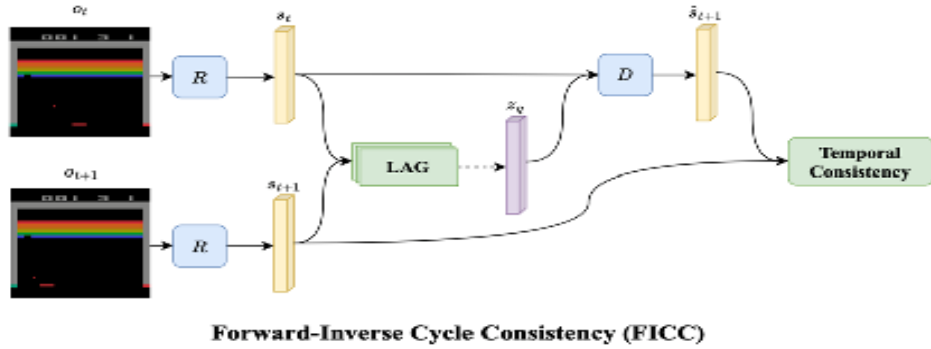
feature	coder	IEEE-CIS Fraud Detection				BNP Paribas Cardif Claims Management			
		pub. score	pub. rank	pri. score	pri. rank	pub. score	pub. rank	pri. score	pri. rank
Base	-	0.0463 \pm 0.0004	2408/6351	0.0182 \pm 0.0004	2285/5351	0.4382 \pm 0.0004	38/2920	0.4569 \pm 0.0004	31/2920
Expert	first	0.0602 \pm 0.0004	52/6351	0.0327 \pm 0.0004	706/6351	0.4334 \pm 0.0004	14/2920	0.4308 \pm 0.0004	12/2920
	high	0.0597 \pm 0.0004	54/6351	0.0320 \pm 0.0004	756/6351	0.4322 \pm 0.0004	12/2920	0.4301 \pm 0.0004	12/2920
OpenFE	first	0.0617 \pm 0.0004	38/6351	0.0360 \pm 0.0004	446/6351	0.4340 \pm 0.0004	77/2920	0.4394 \pm 0.0004	16/2920
	high	0.0617 \pm 0.0004	38/6351	0.0363 \pm 0.0004	426/6351	0.4324 \pm 0.0004	14/2920	0.4302 \pm 0.0004	12/2920

基于纯视频数据预训练的强化学习

强化学习（RL）在一些困难问题上取得了巨大的成功，然而目前大部分强化学习算法仍旧受限于样本效率低，因此对数据量的巨大要求阻碍了强化学习在现实世界中的应用。受到自然语言处理（NLP）与计算机视觉（CV）领域中无监督预训练取得巨大成功的启发，高阳研究组提出一种新的基于模型（model-based）的预训练方法，来提高强化学习算法的样本效率。

考虑到数据的通用性与广泛性，高阳研究组针对无动作标记的纯视频数据进行预训练，之后在下游任务上进行微调。在预训练阶段，研究组提出隐式地从视频训练数据中提取隐式动作表达，并通过一种新颖的正向 - 逆循环一致性（FICC）来预训练特征提取模型和状态转移模型。之后，则在下游任务上建立真实动作到隐式动作表达的映射，接着微调预训练好的模型即可。实验结果表明所提出的训练算法与框架可以显著提升强化学习算法的样本效率，而且能够做到单个预训练模型适配到多任务上。这种预训练算法对数据的要求低，且能显著提升样本效率，为解决现实世界中的强化学习问题提供了更大可能。该研究成果的论文正在 ICLR 会议双盲审核期间。

该成果研究论文：Weirui Ye, Yunsheng Zhang, Pieter Abbeel, Yang Gao, “Become a Proficient Player with Limited Data through Watching Pure Videos” , ICLR 2023.

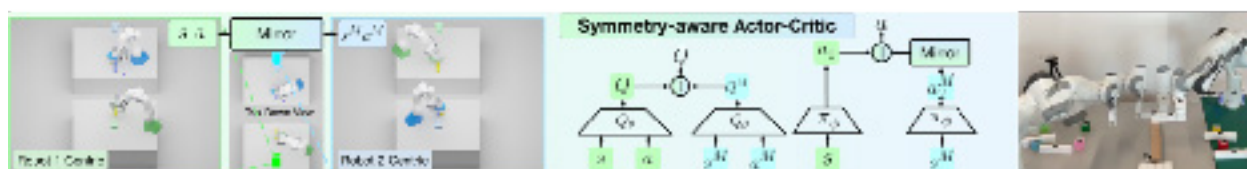


懂得协作的强化学习双臂机器人

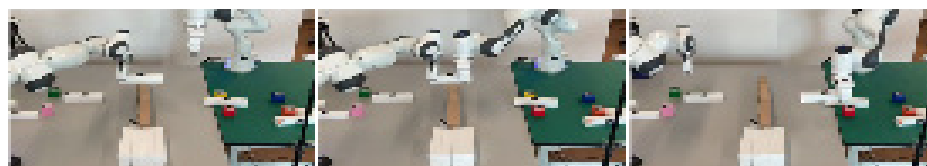
双臂机器人相比于单机械臂能解锁更丰富的技能。然而，常见的强化学习算法很难从双机械臂复杂的控制空间中搜索出精准配合的策略。吴翼研究组观察到双臂协作任务中两个机械臂的角色通常是可以对换的，利用这种对称性该研究组改进了 Actor-Critic 网络结构，提出了 symmetry-aware 结构，有效减小了强化学习的搜索空间，成功让双臂发现了在空中传接物体的策略。为了让双臂协作处理更多数量物体的重排任务，该研究组提出了 object-centric relabeling 技术做数据增强，来产生更多样的部分成功数据。综合以上技术，该研究组成功地让两个机械臂高效协作完成 8 个物体的重排任务。

该研究组将训练出的策略部署在两个固定在不同工作区的 Franka Panda 机械臂上。该研究组的强化学习策略既可以让两只机械臂各自拾取物体放置到本侧工作区，也能协调双臂彼此配合，把物体从一侧传接到另一侧。此外，该研究组还可以在测试时将一只机械臂替换为人，应用到人机协作场景中。

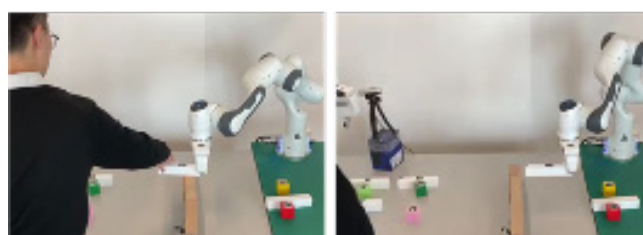
该成果研究论文：Yunfei Li, Chaoyi Pan, Huazhe Xu, Xiaolong Wang, Yi Wu, “Efficient Bimanual Handover and Rearrangement via Symmetry-Aware Actor-Critic Learning”, ICRA 2023.



利用双臂问题对称性的强化学习方法概览



真实双机械臂配合传接



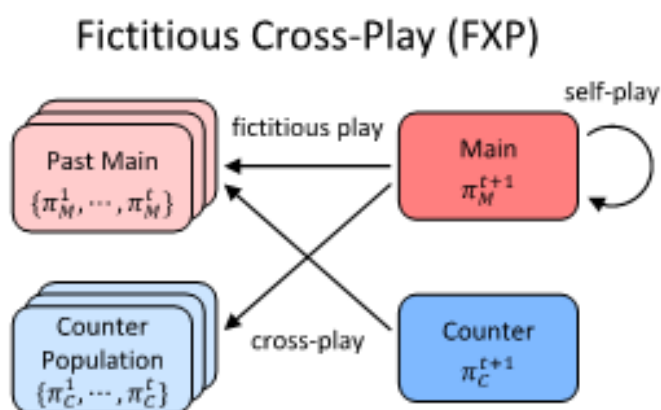
人机合作重排物体

合作 - 竞争混合博弈中的全局纳什均衡学习研究

基于 self-play 的多智能体强化学习方法在复杂竞争场景中取得了巨大的成功，但其收敛性质只适用于双人零和博弈问题。在组内合作、组间竞争的混合博弈场景中，同组的智能体需要共同优化其联合策略，而 self-play 只单独优化当前智能体的策略，因此可能收敛至改变单个策略无法提升奖励，但改变联合策略能提升奖励的局部纳什均衡中。

吴翼研究组在提出了虚拟交叉博弈方法（Fictitious Cross-Play, 简称 FXP），以在合作 - 竞争混合博弈中高效学习全局纳什均衡。FXP 训练主策略和反制策略两组策略，其中主策略通过与历史主策略虚拟博弈和与反制策略交叉博弈学习全局纳什均衡，而反制策略通过学习历史主策略的联合反制策略防止主策略进入局部纳什均衡。FXP 在 Google Research Football 的 11v11 完整比赛中达到了远超内置 AI 的性能，并在与现有 SOTA 模型的对局中达到了超过 2.7 净胜球和超过 94% 的胜率。

该成果研究论文：Zelai Xu, Yancheng Liang, Chao Yu, Yu Wang, and Yi Wu, "Fictitious Cross-Play: Learning Global Nash Equilibrium in Mixed Cooperative-Competitive Games", AAMAS 2023.

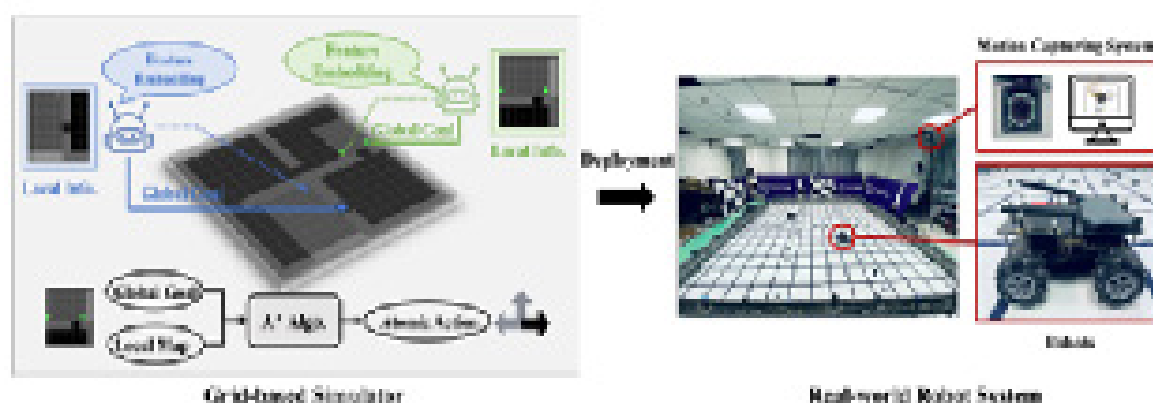
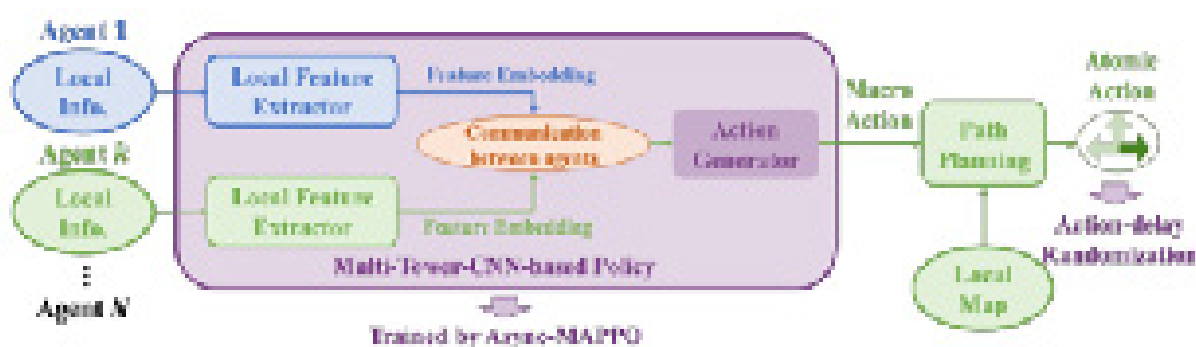


异步高效多智能体协同探索算法

探索是构建智能机器人系统的重要任务，已在救援、自动驾驶、无人机，移动机器人等众多应用领域得到广泛应用。虽然此前的相关工作已经在仿真中取得了较好的结果，但相关研究工作往往会忽略现实部署时由于硬件或是通信等原因造成的多机器人之间执行操作不同步带来的不良影响。

吴翼研究组提出了一种名为 ACE（Asynchronous Cooperative Exploration）的算法框架，用于解决机器人导航中的多机器人合作探索任务。ACE 由 Async-MAPPO、action-delay 随机化和基于 Multi-tower-CNN 的高效策略表示组成。实验结果表明，ACE 算法在仿真实验和实际机器人实验中取得了出色的性能表现。这项研究的重要性在于，ACE 算法框架具有潜在的高效实时多机器人探索应用价值。相较于其他经典的基于规划的方法和同步强化学习方法，ACE 在处理异步多智能体协同探索任务上取得了更好的性能，实现了高效的团队合作。这项研究将为多机器人合作探索导航领域的进一步发展提供有力支持。

该成果研究论文：Yu, Chao, Xinyi Yang, Jiaxuan Gao, Jiayu Chen, Yunfei Li, Jijia Liu, Yunfei Xiang et al. "Asynchronous Multi-Agent Reinforcement Learning for Efficient Real-Time Multi-Robot Cooperative Exploration", AAMAS 2023.

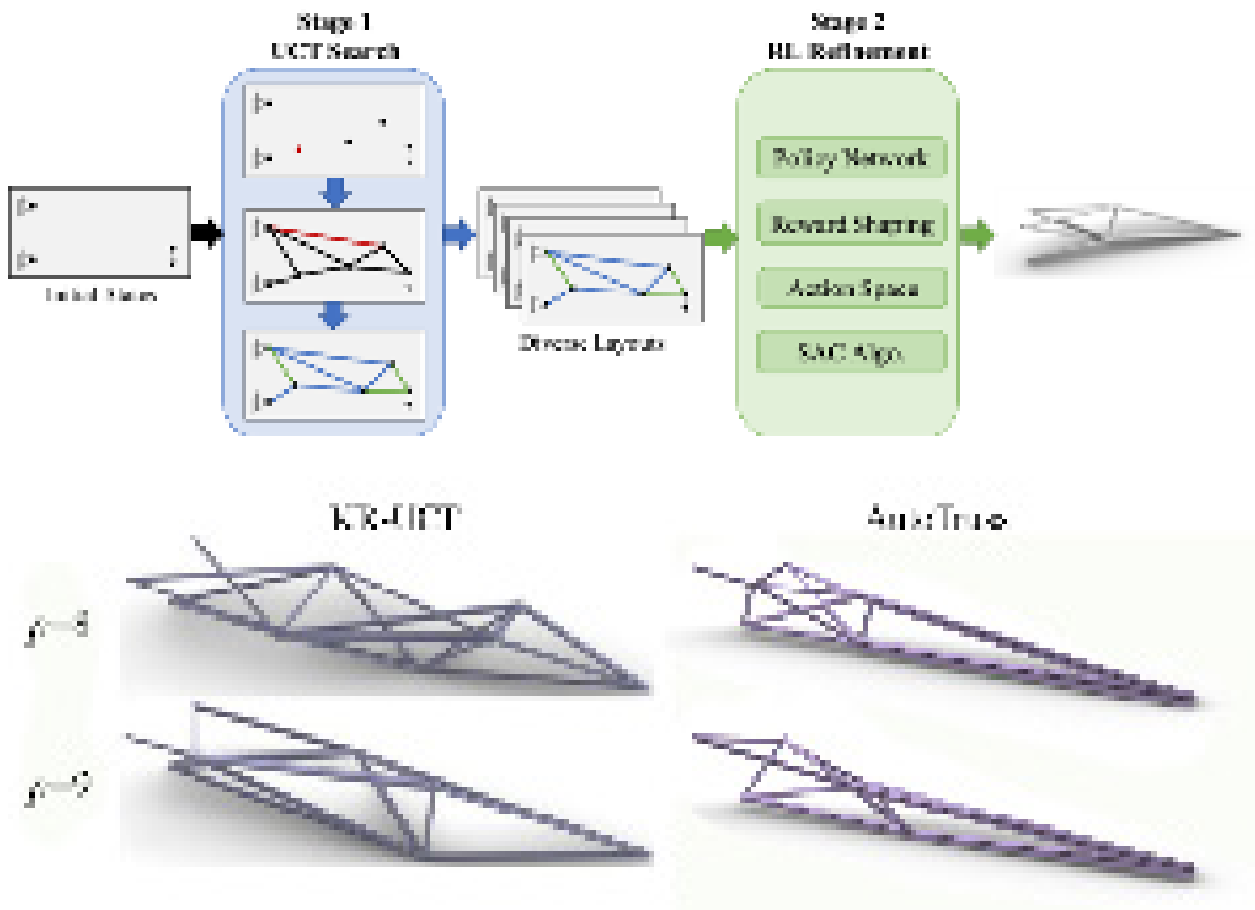


基于强化学习的自动桁架优化方法

桁架结构是一种常见的工程结构形式，广泛应用于各种建筑，由杆件按照一定的几何形式连接构成。桁架优化涉及节点位置、节点之间的拓扑和连接杆的横截面积的优化，是一个复杂的组合优化问题。由于解空间巨大，简单地用计算机进行穷举式搜索并不可行，需要耗费大量的时间和计算资源。强化学习可以进行桁架布局的优化，但存在奖励稀疏，训练困难的问题。

为解决这一问题，吴翼研究组创新性地提出了先搜索后精调的两阶段优化方法，先搜索出满足力学条件的基本解，然后在此基础上通过强化学习进行优化。在搜索阶段，使用应用于树的上置信界限法（Upper Confidence Bound applied to Trees，简称 UCT）搜索并找到多样化的有效桁架结构。在精细优化阶段，使用深度强化学习中的 Soft Actor-Critic 算法（简称 SAC）来进一步优化这些桁架结构。该方法能够有效避免现有优化方法中常见的陷入局部最优的问题，快速高效地生成轻量且符合物理约束的桁架结构，为建筑、汽车乃至航空航天等相关领域带来效益。

该成果研究论文：Weihua Du, Jinglun Zhao, Chao Yu, Xingcheng Yao, Zimeng Song, Siyang Wu, Ruifeng Luo, Zhiyuan Liu, Xianzhong Zhao and Yi Wu. “Automatic Truss Design with Reinforcement Learning”, IJCAI 2023.



零和马尔可夫游戏中的可微分仲裁框架以引导理想纳什均衡

在人类社会发展中，如何调和个人利益和集体利益的冲突一直是一个重要问题。研究者们通常以马尔可夫游戏为模型研究这种冲突。在该游戏中，为追求自我利益的玩家可能会导致不理想的纳什均衡，从而破坏整体福利。激励设计通过扰动奖励优化纳什均衡，以期自我利益导向的玩家达到理想均衡。这个问题可以被表述为一个双层结构，但通过双层结构进行求导是困难的。一些解决方案包括将问题转化为多目标问题，同时优化游戏奖励和设计者的目标；或者保留双层问题结构，但在上层应用梯度无关的优化器，这两种方法都存在局限性。

吴翼研究组将激励设计问题扩展到多智能体强化学习领域，并解决了通过双层结构推导上层目标梯度的挑战。他们开发了首个可微分的一阶框架——可微分仲裁（Differentiable Arbitrating, DA），能够调解冲突并在马尔科夫游戏中获取理想的纳什均衡。此外，他们还还为两玩家零和马尔科夫游戏中的 DA 框架提供了理论收敛证明。与零阶优化方法相比，DA 框架在两个多智能体强化学习的测试任务中更高效地找到了可解释的理想纳什均衡。

该成果研究论文：Jing Wang, Meichen Song, Feng Gao, Boyi Liu, Zhaoran Wang, Yi Wu. “Differentiable Arbitrating in Zero-sum Markov Games”, AAMAS 2023.

Framework 1 DA: Differentiable Arbitrating in MARL.

Input: β_k : learning rate for upper-level iteration

Output: θ : incent. param.; ϕ : param. of policy π_ϕ

Initialize the incentive parameter $\theta = \theta_0$.

for $k=0,1,\dots$ **do**

 Initialize param. $\phi = \phi_0$ of policy π_{ϕ_0} for game \mathcal{G}'_{θ_k} .

for $t=0,1,\dots$ **do**

 Update $\phi = \phi_t$ with NE solver until corresponding policy π_{ϕ_t} converge to the Nash equilibrium of game \mathcal{G}'_{θ_k} , w.r.t $\phi^*(\theta_k)$.

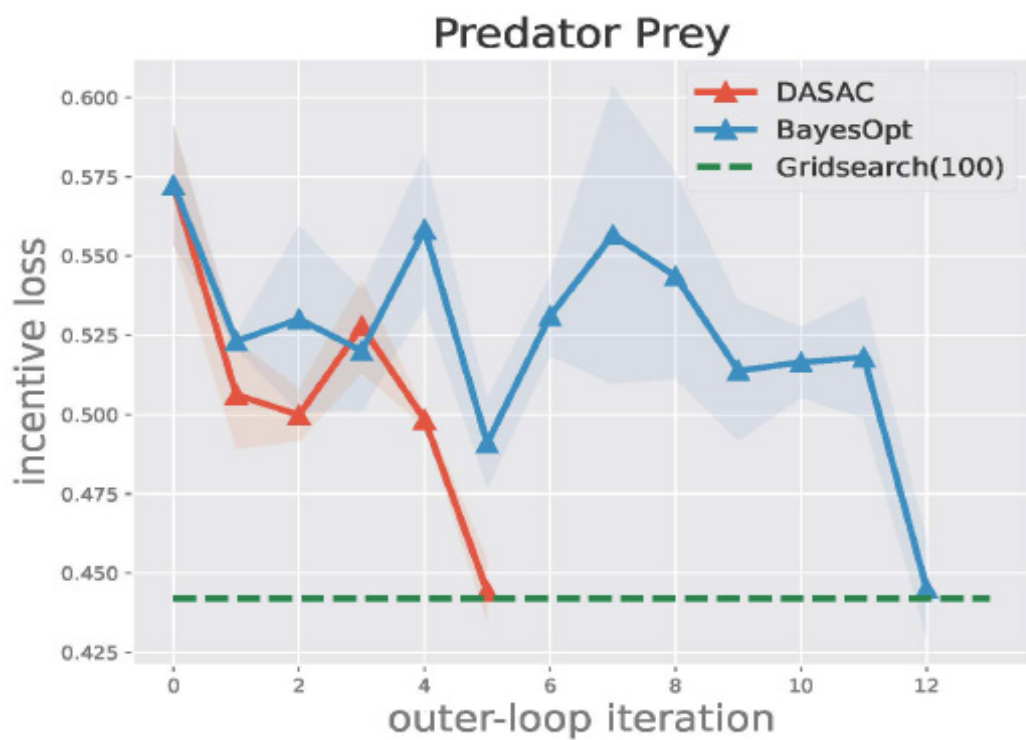
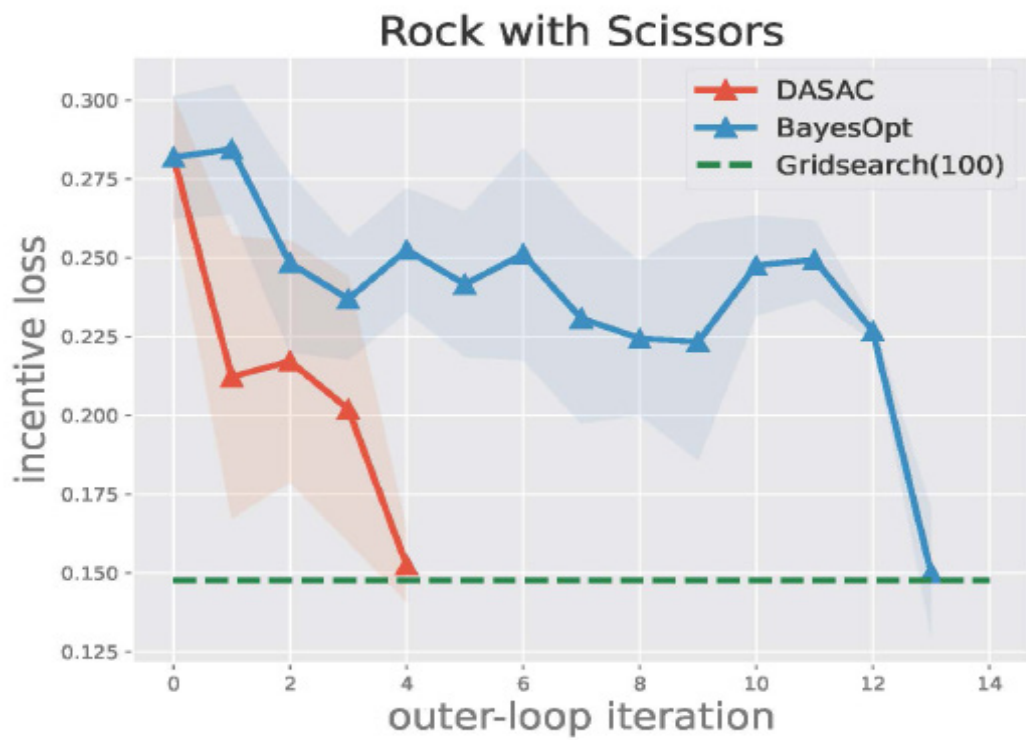
end for

 Update incentive parameter

$$\theta = \theta_{k+1} \leftarrow \theta_k - \beta_k \nabla f_*(\theta_k),$$

 where $\nabla f_*(\theta_k)$ is defined in (4.6).

end for

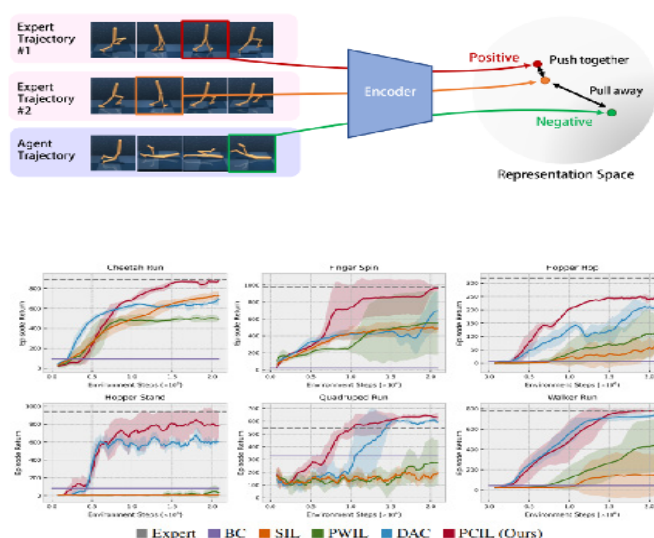


策略对比模仿学习

对抗性模仿学习 (AIL) 是一种流行的方法，最近取得了很大的成功。然而，AIL 在更具挑战性的任务上表现仍然不尽如人意。高阳研究组发现主要原因之一是 AIL 鉴别器表示质量低。由于 AIL 鉴别器是经过训练的通过二元分类，不一定以有意义的方式区分专家的政策，由此产生的奖励也可能没有意义。

作为一项在模仿学习领域的突破性进展，高阳研究组首次提出了一种称为策略对比模仿学习 (PCIL) 的新方法来解决此问题。PCIL 通过锚定在不同的策略，并使用基于平滑余弦相似性的奖励函数来进行模仿学习。从理论的角度，该方法是学徒学习框架的一种。此外，在 DeepMind Control 套件上的实证评估表明 PCIL 可以达到最先进的性能。该研究成果在模仿学习取得了重要的进展，对于进一步研究进行对抗模仿学习中的特征空间具有重要价值。

该成果研究论文：Jialei Huang, Zhaocheng Yin, Yingdong Hu, Yang Gao, “Policy Contrastive Imitation Learning”, ICLR 2023.

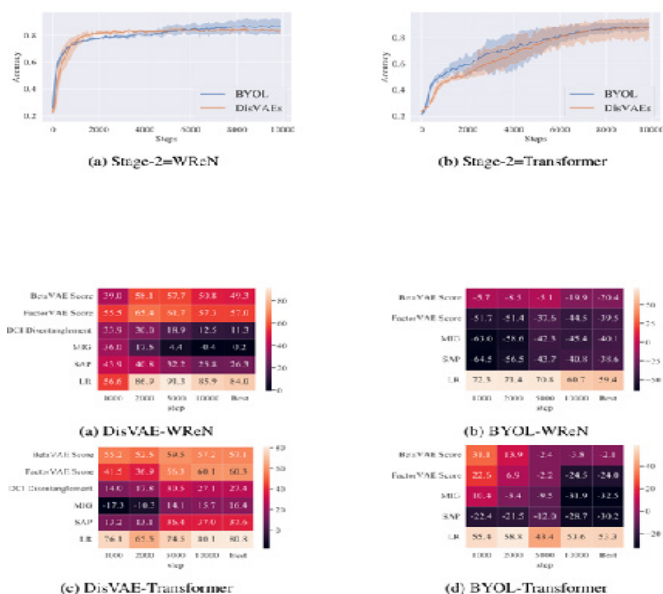


解耦表征学习对于下游任务的必要性

解耦表征 (disentangled representation) 是表征学习追求的重要目标。由于其将数据中的重要信息可分离地编码, 并且可解释性较强, 因而研究者们认为其对于下游任务具有重要作用。在一些常见的任务中, 前人的研究提供了很多实验证据, 说明解耦性质越好的表征在下游任务上可以获得更好的表现。

高阳研究组在衡量解耦表征下游任务性能的标准测试任务抽象归因 (abstract reasoning) 上进行了大量的实验, 指出解耦表征学习对于下游任务不是必要的。具体的实验结论是: 下游任务的性能与表征的解耦性质相关性较小, 但是与表征的信息量 (informativeness) 相关性最高。在该领域, 研究者们认为应该使用解耦性质较好的表征去完成下游任务。该研究通过充分的实验证据说明了这样的想法是错误的, 指出了解耦表征学习领域的问题, 对于解决 “如何选取合适的表征去完成下游任务” 这一问题提供了帮助。

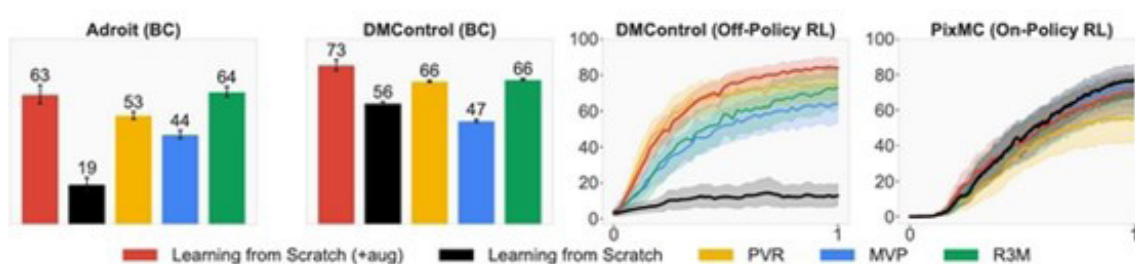
该成果研究论文: Ruiqian Nai, Zixin Wen, Ji Li, Yuanzhi Li, Yang Gao, “On the Necessity of Disentangled Representations for Downstream Tasks”, ICLR 2023.



基于图像输入的强化学习

将预训练模型用于各类计算机视觉任务，或是基于图像输入的控制任务，在大模型盛行的今天，已经成为了研究者的常用手段。然而该文提出了一种简单且高效的从头学习 (Learning from Scratch) 方法，在仅使用数据增强和浅层卷积神经网络的情况下，即可在样本利用率，泛化性能上，与使用先进预训练模型的方法取得类似甚至更优的效果。

该成果研究论文：Nicklas Hansen, Zhecheng Yuan, Yanjie Ze, Tongzhou Mu, Aravind Rajeswaran, Hao Su, Huazhe Xu, Xiaolong Wang, “On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline”, ICML 2023.



三、计算机视觉

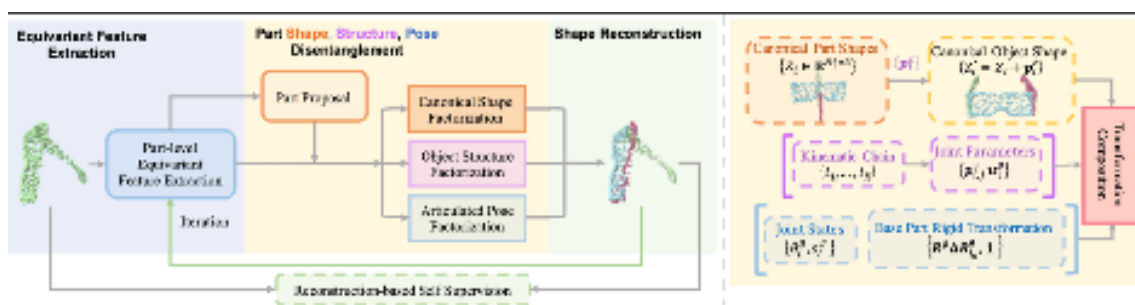
主要完成人：弋力研究组、马恺声研究组

借助部件级 SE(3) 等变性的自监督铰接物体位姿估计方法

铰接物体广泛存在于日常生活中的各种场景之中。人类从婴儿时期便开始逐渐接触这些物体并在一系列尝试和反馈之中学到了如何与这些物体进行交互，比如拉开书桌的抽屉。如果希望设计算法使得机器也具有这样的感知世界并与世界中的物体进行交互的能力，具有将物体分成不同的部分（在此称之为部件分割），理解每个部分在空间中的位置（称之为部件位姿），和可以怎样操控这些部分就变得十分重要。但以往的方法往往依赖于有部件分割、部件位姿等相关的人工标注。然而标注相关的数据往往是繁琐和昂贵的。这样的对有标注数据的依赖使得前人的方法往往只能在合成数据集上训练，而无法使用更多的无标注的数据来得到更加强大、泛化能力更强的模型。这也进一步使得他们的算法失去了更为广阔的应用前景。反之，不依赖于标注的自监督解法可以用更多的、更加接近现实世界的数据来训练。

从而该文希望使用完全自监督的方法解决该问题。这需要算法具有将铰接物体标准型和部件位姿这些信息从输入的形状中解耦出来的能力。这些信息在使用普通的网络（如 PointNet++）所得到的普通的几何特征中往往耦合在一起。从而仅仅使用自监督信号来实现所希望的分解是很困难的。基于铰接物体的部件位姿的部件级 SE(3) 等变性 -- 即每个部件的位姿只等变于该部件在空间中的位置，而与其他部件无关，和标准型的部件级 SE(3) 不变性，即与任一部件的位置都无关，该文提出了部件级 SE(3) 等变卷积运算，并基于此设计了部件级 SE(3) 等变网络和自监督部件位姿分解方法，实现了无需任何额外标注情况下的铰接物体位姿估计方法。

该成果研究论文：Xueyi Liu, Ji Zhang, Ruizhen Hu, Haibin Huang, He Wang, and Li Yi, "Self-Supervised Category-Level Articulated Object Pose Estimation with Part-Level SE (3) Equivariance", ICLR 2023.

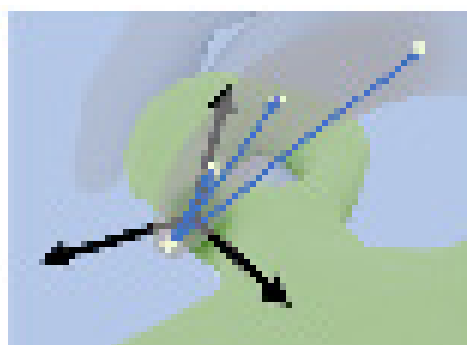
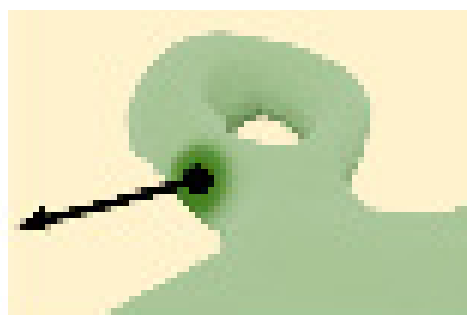
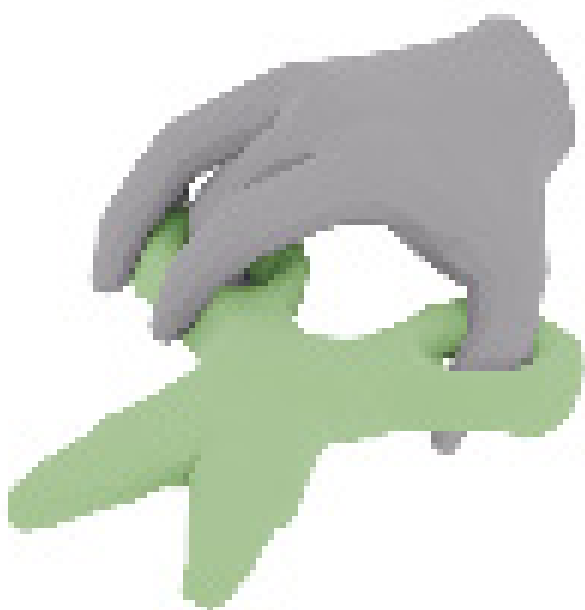
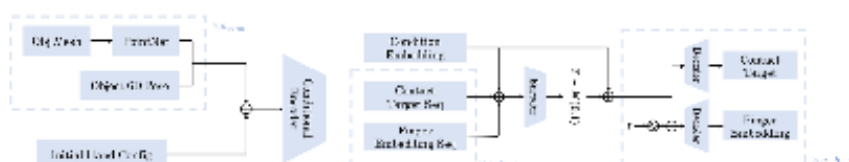


标准化手物交互空间表示及动画生成

长期以来，人类灵巧手与物体的交互问题受到机器人、三维视觉、图形学社区的高度关注。给定固定类别（如剪刀）下的一个物体示例以及操作目标，弋力研究组希望能生成一段与之匹配的人类灵巧手操作动画，其中要求生成出的动画应接近人类操作方式，且物理真实。现有方法局限性相对较大，具体体现在训练数据采集困难，对形状泛化性不佳等。为此，弋力研究组提出了 CAMS (CAnonicalized Manipulation Spaces) 标准化手物交互表示，其通过提取手与物体的接触点并建立以接触为中心的手指嵌入，来实现手指 - 物体交互及输入物体形状的解耦。

基于 CAMS 表示，弋力研究组提出了一套采用 CVAE 结构的手物交互动画生成框架。得益于 CAMS 表示的标准化特性，在同一物体类别下动作表示对具体物体形状的依赖性大大降低，因此同类别泛化性得以提升。实验表明，对于任意给定的目标形状，该方法均能生成出物理真实且类人的交互动画，而此前的工作对于训练数据中并未见过的形状存在较大的泛化问题。

该成果研究论文：Juntian Zheng, Lixing Fang, Qingyuan Zheng, Yun Liu, Li Yi, "CAMS: CAnonicalized Manipulation Spaces for Category-Level Functional Hand-Object Manipulation Synthesis", CVPR 2023.

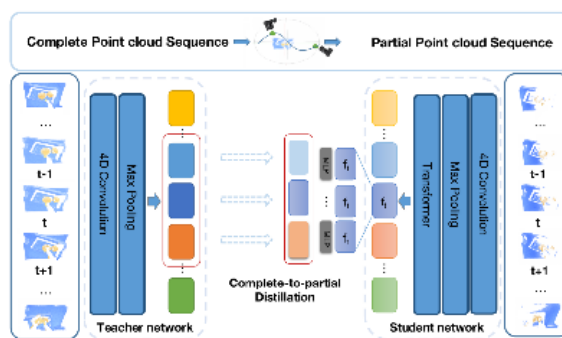


基于知识蒸馏的点云序列预训练研究

近些年来，4D 空间中的点云序列感知由于其广泛的下游应用性引起了大家的强烈兴趣。作为包括机器人和增强现实等一系列应用的直接感知输入，第一人称相机视角下的点云序列描绘了包括物体几何和物体运动在内的动态场景。尽管 4D 数据非常容易获得，但大规模的精细标注是极其昂贵的，因此人们十分希望通过自监督表征学习的方法来增强模型对于点云序列的感知能力。

不同于 3D 表征，4D 表征需要统一几何和运动信息。对于运动的理解有利于帮助聚合时间维度的观测来形成更完整的几何；同时准确的跨时间几何对应关系有利于更好地估计运动。为了实现几何和运动的协同，弋力课题组提出了 C2P 自监督表征学习框架。C2P 将点云序列表征学习统一在知识蒸馏框架下，通过鼓励输入为完整点云序列的老师网络向输入为局部视角点云序列的学生网络蒸馏知识信息来提高网络的感知能力。实验表明 C2P 在一系列室内室外下游任务数据集上都表现出优异的性能。

该成果研究论文：Zhuoyang Zhang, Yuhao Dong, Yunze Liu, Li Yi, “Complete-to-Partial 4D Distillation for Self-supervised Point Cloud Sequence Representation Learning”, CVPR 2023.

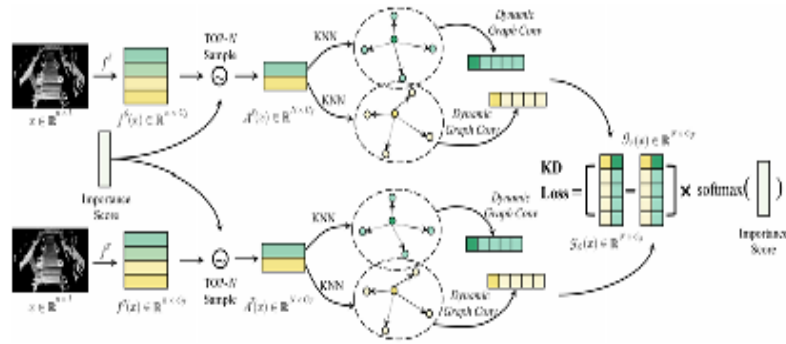


通过知识蒸馏加速基于点云的 3D 检测算法

基于点云的 3D 检测算法是无人驾驶感知系统的重要技术。然而，已有的基于点云的 3D 检测模型往往有着超大的参数量和计算量，导致其难以在边缘设备上部署运行。因此，对于此类模型进行压缩和加速成为了亟待解决的问题。

为解决该问题，马恺声研究组的张林峰、董润沛同学提出了一种新型的知识蒸馏算法。他们观察到与图像数据相比，点云数据存在着其独特的特点：点云数据具有显著的稀疏性，更容易收到噪音的影响，以及点云数据缺乏规则的空间结构。针对这些特点，他们提出通过统计不同体素内点的数量定义不同体素在知识蒸馏中的学习重要性，避免模型过分关注噪音体素和稀疏的体素。其次，他们提出利用动态图卷积神经网络提取点云中的局部空间关系，然后对此进行知识蒸馏。实验结果显示，他们提出的知识蒸馏方法在 KITTI 数据集中取得了四倍以上的模型压缩与加速效果，同时没有任何检测精度的损失。这项研究成功地降低了自动驾驶系统中感知模块的计算成本，对无人驾驶的落地应用产生积极意义。

该成果研究论文：Linfeng Zhang, Runpei Dong, Hung-Shuo Tai, Kaisheng Ma. “PointDistiller: Structured Knowledge Distillation Towards Efficient and Compact 3D Detection”, CVPR 2023.

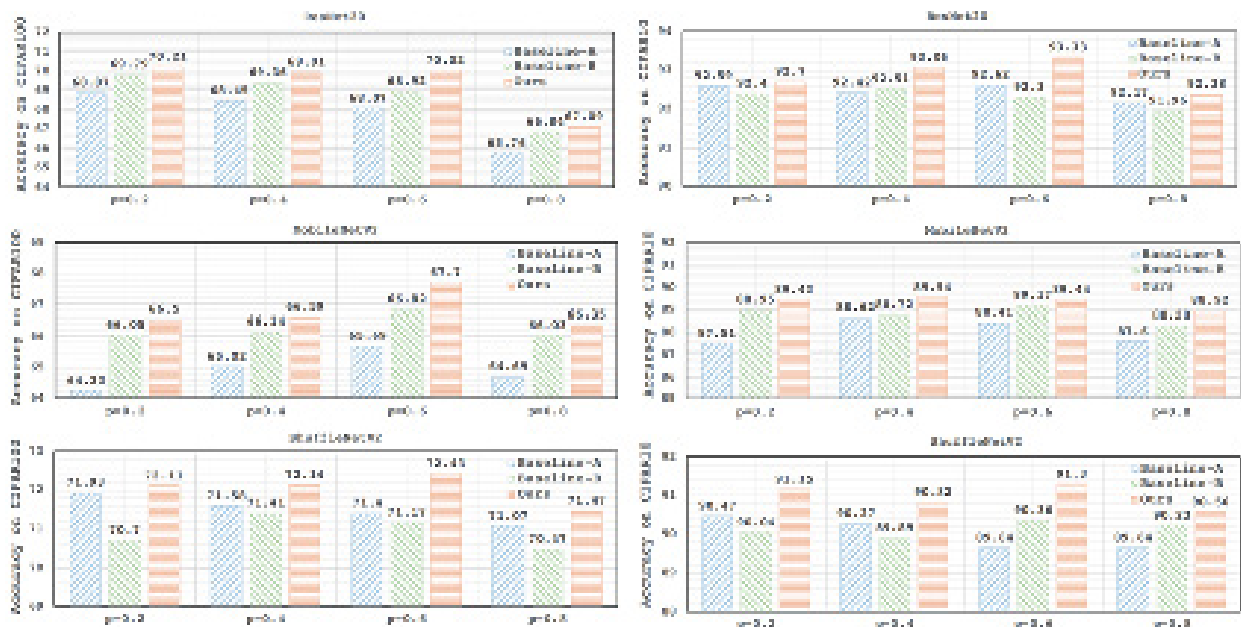


探究模型压缩中数据增强的使用方式

近年来，为推动神经网络模型在边缘设备上部署运行，大量模型压缩算法相继提出。例如，模型剪枝方法和模型架构自动化搜索算法从模型结构的角度切入，寻找最优的轻量化模型结构。知识蒸馏从模型训练的角度入手，探索如何在模型压缩过程中最大化保留知识。然而，作为以数据驱动的科学，模型压缩与训练数据之前的关系尚未被深入探究。

观察到这一问题，马恺声研究组的张林峰等人探索了在模型压缩过程中数据增强的使用方式。他们通过实验有以下发现：（一）不同参数量的神经网络所最适用的数据增强方式是不同的，参数量多的神经网络可以从更强的数据增强中学习知识。（二）尽管直接训练小模型无法从强数据增强中获取知识，这种强数据增强的知识可以在模型压缩的过程中保留下来，这说明通过模型压缩训练小模型的范式有着比直接训练小模型更大的价值。（三）模型的参数量在训练阶段直接影响了模型的学习能力，对此，可以通过给小模型添加额外的神经网络层，辅助其在训练阶段获取知识。

该成果研究论文：Muzhou Yu*, Linfeng Zhang*, and Kaisheng Ma, “Revisiting Data Augmentation in Model Compression: An Empirical and Comprehensive Study”, IJCNN 2023.



四、自然语言处理

主要完成人：杨植麟研究组、李建研究组

深入理解零样本泛化

工作通过多任务提示预训练 (Multi-Task Prompted Training) 取得了优越的零样本性能, 然而人们对此知之甚少。该文首次证明了, 对少量关键任务进行训练优于使用所有训练任务, 而删除这些关键任务会严重损害性能 (如图 1 所示)。

杨植麟研究组还发现这些关键任务大多是问答任务 (Question-Answering Task, QA)。这些新的发现结合起来加深了该研究组对零样本泛化的理解, 即对某些任务 (例如: Question-Answering Task, QA) 的训练编码了可迁移到广泛任务的一般知识。此外, 为了自动化此过程, 该研究组设计了一种方法, 该方法 (1) 通过检查成对泛化结果来识别关键训练任务, 而无需观察测试任务; (2) 对训练任务进行重新采样以获得更好的数据分布。实验表明该研究组的方法在各种模型规模和任务中取得了更好的性能 (如表格 5 所示)。

该成果研究论文: Jing Zhou, Zongyu Lin, Yanan Zheng, Jian Li, Zhilin Yang, “Not All Tasks are Born Equal: Understanding Zero-Shot Generalization”, ACL 2023.

[illegible]

Figure 1: Pairwise transfer relationships on T5-XL. The entry at row i and column j denotes the average performance when the model is trained on task i and evaluated on task j . For each entry, the value is the average score of different prompts. (Accuracy if only Accuracy is calculated, and otherwise the mean of Accuracy and F1.) Only those prompts related to the original tasks are included for evaluation. We highlight those entries with high scores for each task (Red is the Top-1). The horizontal and vertical lines denote the boundary of task-type groups.

Model	Met.	Natural Language Inference			Sentence Completion			Co-References		WSD	Avg.
		RTE	CB	ANLI	ANLI	CoPA	HeRa	Stacy	WSC		
TS-Large-LM-Adapt (770M)											
TO (*)	Mean	72.53	50.60	50.93	51.96	32.23	82.30	22.16	92.05	62.21	59.00
	Med.	74.01	57.14	50.40	51.60	31.75	83.00	22.60	91.77	62.98	59.00
DS-DA-TU	Mean	74.22	60.95	35.65	32.57	35.88	87.66	29.49	94.12	63.75	54.89
	Med.	75.56	62.40	32.40	32.40	36.40	88.10	29.70	94.10	63.70	54.80
US-TU	Mean	80.72	71.90	56.00	54.80	38.18	84.10	26.00	94.00	63.27	54.54
	Med.	81.23	80.36	56.40	55.20	39.33	85.31	26.06	94.59	63.94	54.58
US-DA-TU	Mean	78.51	56.25	36.25	34.00	37.25	86.25	26.25	93.75	63.25	53.75
	Med.	79.00	69.60	56.10	33.90	36.75	89.00	28.14	94.92	64.42	56.43
US-DA-TU	Mean	80.69	70.95	37.38	54.30	39.43	87.97	26.73	93.71	63.27	53.58
	Med.	80.69	80.35	38.00	54.20	40.33	89.29	26.98	93.91	64.42	53.72
TS-XL-LM-Adapt (CB)											
TO (†)	Mean	64.55	45.35	33.84	33.11	33.35	72.40	27.29	84.03	65.10	50.90
	Med.	64.08	50.08	33.63	33.40	33.82	72.51	28.06	84.62	65.31	50.39
TO (*)	Mean	80.72	67.62	41.09	37.79	40.38	91.92	32.03	97.27	65.96	57.84
	Med.	80.14	75.00	42.80	39.20	41.75	92.00	32.29	97.22	68.27	58.41
DS-TU	Mean	83.21	75.33	64.38	38.84	43.72	94.17	31.31	97.72	64.82	62.67
	Med.	83.67	82.14	64.40	39.70	44.00	93.70	31.63	97.99	64.63	63.38
DS-DA-TU	Mean	84.77	74.40	43.25	39.17	43.22	94.93	27.01	97.65	62.02	66.74
	Med.	84.66	82.14	46.30	39.70	43.75	95.00	27.00	97.63	62.98	65.35
US-TU	Mean	82.41	62.34	38.20	38.40	41.20	90.40	26.20	93.70	62.90	53.70
	Med.	82.34	82.81	55.41	39.94	42.60	93.75	30.38	97.34	63.67	62.27
US-DA-TU	Mean	83.29	75.83	44.49	39.93	43.68	94.81	26.29	96.94	61.73	65.03
	Med.	84.48	82.14	47.90	39.90	47.25	94.50	26.17	97.06	64.90	65.11
TS-XL-LM-Adapt (11B)											
TO (†)	Mean	80.83	70.12	43.56	38.68	41.26	90.62	33.58	92.40	61.45	59.94
	Med.	80.40	78.57	44.70	38.40	40.79	93.65	34.71	94.42	60.46	57.21
TO (*)	Mean	84.01	72.26	47.89	42.80	46.49	91.60	35.27	98.15	62.69	65.45
	Med.	83.02	83.95	49.00	44.40	48.28	95.00	34.62	98.24	65.35	70.24
Our Base	Mean	85.56	82.50	58.28	47.88	51.28	97.88	37.18	98.88	67.18	68.98
	Med.	85.74	82.14	59.60	46.40	51.25	95.00	37.11	97.54	66.35	71.19

Table 5: Zero-shot performance for our improved T0 and original T0 at three different scales. Results with † are reported by Sanh et al., and results with * are reproduced in our experiments. US-T0 means T0 with upsampling key tasks, DS-T0 means T0 with downsampling non-key tasks, and DS+DA-T0 / US+DA-T0 represents DS-T0 / US-T0 with augmented data. “Our Best” is achieved with the US+DA-T0 setup.

大型语言模型的组合任务表示

大型语言模型 (Large Language Models, LLM) 表现出了卓越的跨任务泛化能力 (Cross-task Generalization)。大多数先前的工作假设提示可以有效地从语言模型中提取知识，以促进泛化到新任务；这种观点进而引发了大量优化和改进提示的研究工作。相比之下，杨植麟研究组引入了一种新的视角，即组合泛化 (Compositional Generalization)，它将每个任务视为隐代码的组合，并通过可见代码的新组合来表示测试任务，如下面图 1 所示。为此，该研究组提出一种新颖的无提示方法，即组合任务表示 (Compositional Task Representation, CRT)。它采用多任务训练来学习离散的组合密码本。实验表明 CRT 方法相较于基于提示的方法，其点击率在零标签学习中显著表现更优（如表 1 所示）。此外学习到的 CRT 代码是人类可以解释的，并表现出一定程度的可控性（如表 6 所示）。

该成果研究论文：NAN SHAO, Zefan Cai, Hanwei xu, Chonghua Liao, Yanan Zheng, Zhilin Yang, “Compositional Task Representations for Large Language Models”, ICLR 2023.

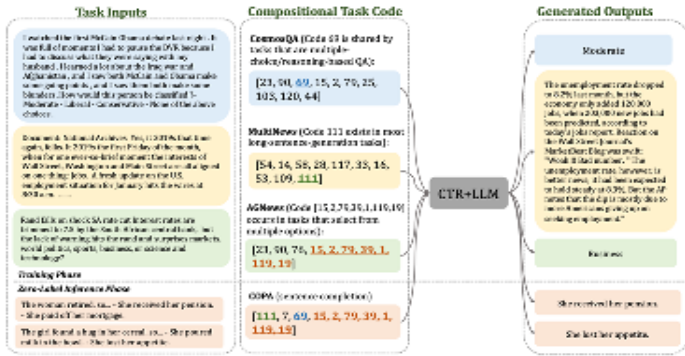


Figure 1: An illustration of how CRT generalizes to zero-label tasks. In this real example produced by our model, CRT combines the abilities of reasoning-based QA, sentence generation, and multi-choice selection from training tasks to perform a new task COPA.

Method	Natural Language Inference					Sentence Completion			Co-reference		WSD	Avg.
	RTE	CB	ANLI1	ANLI2	ANLI3	COPA	Hella	Story	WSC	Winz.	WSC	
Zero-Label Setting (unlabeled data of each test task)												
TU-Large	72.67	56.55	32.77	32.15	34.38	85.36	27.18	93.04	63.94	54.35	50.33	54.79
Self-Training	73.57	76.14	34.42	32.90	37.44	87.45	30.33	94.54	57.08	56.56	50.75	57.38
Manual-Code	75.19	56.89	31.12	32.49	33.48	75.76	30.84	93.10	61.16	54.10	51.45	54.33
ZPS	79.06	67.86	31.20	31.10	34.25	88.00	29.16	93.43	65.38	53.43	49.84	56.61
Our CTR	80.51	87.50	33.40	34.40	33.80	92.00	27.50	90.10	56.58	49.40	62.50	58.88
Few-Shot Setting (32 labeled data of each test task)												
Model Tuning	75.31	80.95	35.73	31.31	35.93	82.05	41.86	92.04	55.96	56.74	52.15	58.18
Prompt Tuning	77.08	76.90	31.89	31.86	35.53	81.70	31.18	94.10	62.88	55.42	51.22	57.25
GPT	77.68	79.64	32.71	31.49	37.56	81.08	28.11	93.40	64.23	52.72	52.52	57.38
GRIPS	71.56	70.89	32.14	32.26	34.77	77.56	26.44	93.40	62.12	52.96	52.12	55.11
BBT	71.19	57.26	33.79	32.00	35.30	76.49	28.95	93.11	62.12	53.40	52.93	54.23
Our CTR	80.51	83.93	34.40	34.20	36.60	89.00	35.07	91.70	68.18	55.09	58.62	60.66

Table 1: Main results of CRT and baselines on 11 test tasks under the zero-label setting and the few-shot setting. The zero-label setting allows using unlabeled data of the test task while the few-shot setting uses 32 labeled data of the test task. All methods share a similar model size (i.e., 770M).

INPUT: If Jim moves some particles of matter further apart, what will happen to the rate at which they can pass vibrations on to nearby particles decrease or increase? Assuming that: When particles of matter are closer together, they can move quickly pass the change of vibrations to nearby particles.	
TASK CODE: [61, 95, 127, 113, 53, 109, 5, 87, ..., 5]	
FILE 52 → DIALOGUE: W. H. Lee, Tim Lee, Tim a physicist at the University of California, Berkeley.	
FILE 37 → TOPIC CLASSIFICATION: Science and Technology	
INPUT: A historic victory for the separatist Parti Québécois was marred by an attack on a victory rally for pro-unionist Pauline Marois. A man opened fire during her victory speech, killing one person and critically wounding another, reports CTV. ... "I have condolences and I am going to defend them," Marois said during her victory speech. "There will be a referendum when the Quebec population wants a referendum"	
TASK CODE: [31, ..., 6, 113, 55, 107, 18, 87, 12, 63]	
FILE 48 → REVIEW RATING: 1 star	
FILE 92 → SUMMARIZATION: One person was killed and another critically injured in a shooting at a victory rally for the separatist Parti Québécois in Quebec City on Sunday.	
INPUT: Paragraph: I've been here a few times and I like that it's right up the street. On Tuesday and Wednesday they have specials on manicures and pedicures which is 5\$ cheaper than the original price.	
TASK CODE: [23, 95, 76, 15, 2, 74, ..., 1, 119, 19]	
FILE 13 → SENTIMENT ANALYSIS: Positive	
FILE 5 → TOPIC CLASSIFICATION: Business	

Table 5: Examples of controlling the compositional codes. Given the input and a compositional task code, each time we modify one bit of the code, CRT begins to perform a different task.

Predict with the code for AGNews topic classification: [23, 95, 76, 15, 2, 74, 34, 1, 119, 19]	
Original AGNews classes: {World, Sports, Business, Science and technology}	
INPUT: Four die T N person after talks Unions representing workers at Turner Newall say they are "disappointed" after talks with strikers parent firm Federal Mogul.	
OUTPUT: Business	
INPUT: There are 10 apples on an apple tree. Three fall off. Now there are X apples.	
OUTPUT: Math problem	
INPUT: Stunning even for the non-gamer This sound track was beautiful! It paints the scenery in your mind so well I would recommend it even to people who hate game music! It would impress anyone who cares to listen!	
OUTPUT: Music	
INPUT: Slack (2003) compares three groups that conducted biological research at Yale during overlapping periods between 1910 and 1970. Yale proved important as a site for this research. ... Huxley's example shows the new modern research group are needed, especially for those that include research field research.	
OUTPUT: history	

Table 6: Examples of CRT generalizing to new unseen classes/domains. The inputs are randomly selected from all tasks other than AGNews. The first case predicts the same classes as AGNews defines, while the latter three cases predict new classes that are never seen within AGNews. It shows the out-of-distribution of CRT can generalize to new unseen classes/domains.

零样本泛化通用判别器

一直以来，生成式建模是大规模预训练和零样本泛化的主要方法。在这项工作中，杨植麟研究组通过证明判别式方法在大量 NLP 任务上比生成式方法表现得更好来挑战这一惯例。下图 2 展示了生成式模型（左图）和判别式模型（右图）的区别。从技术上讲，该研究组训练单个判别器来预测文本样本是否来自真实的数据分布，类似于 GAN。由于许多 NLP 任务可以被表述为从几个选项中进行选择，因此该研究组使用这个判别器来预测输入的串联以及哪个选项来自真实数据分布的概率最高。这个简单的公式在 T0 基准上实现了最先进的零样本结果，在不同参数规模上分别优于 T0 16.0%、7.8% 和 11.5%（如下图 1 所示）。在微调设置中，该研究组的方法还在广泛的 NLP 任务上取得了新的最先进的结果，而参数仅为之前方法的 1/4。同时，该研究组的方法需要最少的提示工作，这大大提高了鲁棒性，对于现实世界的应用程序至关重要。此外，该研究组还结合生成任务联合训练广义 UD，保持其在判别任务上的优势，同时适用于生成任务。

该成果研究论文：Haike Xu, Zongyu Lin, Jing Zhou, Yanan Zheng, Zhilin Yang, “A Universal Discriminator for Zero-Shot Generalization”, ACL 2023.

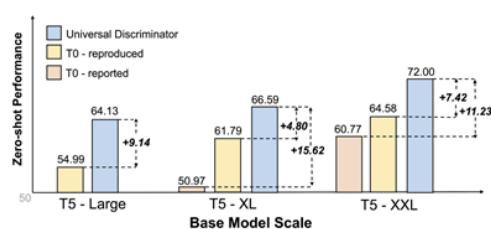


Figure 1: Average zero-shot performance over 11 zero-shot tasks for our Universal Discriminator and T0 (Sanh et al., 2021). Our universal discriminator significantly outperforms T0 across three different scales.



Figure 2: An overview that compares the multi-task prompted formulation of T0 (Sanh et al., 2021) and the formulation of our universal discriminator. This underlines our natural language prompt. The universal discriminator uses a shared formulation of the discriminative task—discriminating whether a sample comes from the true data distribution of natural language.

CodeGeeX: HumanEval-X 多语言评估的代码生成的预训练模型

大型预训练代码生成模型，例如 OpenAI Codex，可以生成语法和功能正确的代码，使程序员的编码更加高效，也更加接近对通用人工智能的追求。 该文介绍了 CodeGeeX，这是一个具有 130 亿个代码生成参数的多语言模型。截至 2022 年 6 月，CodeGeeX 在 23 种编程语言的 8500 亿个令牌上进行了预训练。大量实验表明，在 HumanEval-X 上的代码生成和翻译任务中，CodeGeeX 的性能优于类似规模的多语言代码模型。在 HumanEval (仅限 Python) 的基础上，杨植麟研究组开发了 HumanEval-X 基准，通过用 C++、Java、JavaScript 和 Go 手写解决方案来评估多语言模型。此外，该研究组还在 Visual Studio Code、JetBrains 和 Cloud Studio 上构建基于 CodeGeeX 的扩展，每周为数万活跃用户生成 47 亿 token。该研究组的用户研究表明 CodeGeeX 可以帮助提高 83.4% 的用户的编码效率。

该成果研究论文: Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, Teng Su, Zhilin Yang, Jie Tang, “CodeGeeX: A Pre-Trained Model for Code Generation with Multilingual Evaluations on HumanEval-X” , arXiv:2303.17568.

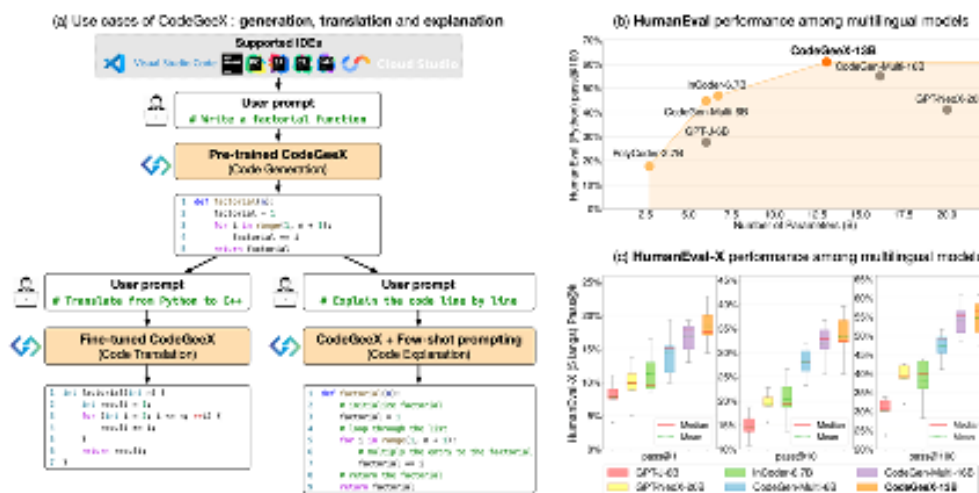


Figure 1: Summary of CodeGeeX. (a): In supported IDEs, users can interact with CodeGeeX by providing prompts. Different models are used to support three tasks: code generation, code translation and code explanation. (b) and (c): In HumanEval and our newly-proposed HumanEval-X, CodeGeeX shows promising multilingual abilities and consistently outperforms other multilingual code generation models.

五、计算机图形学

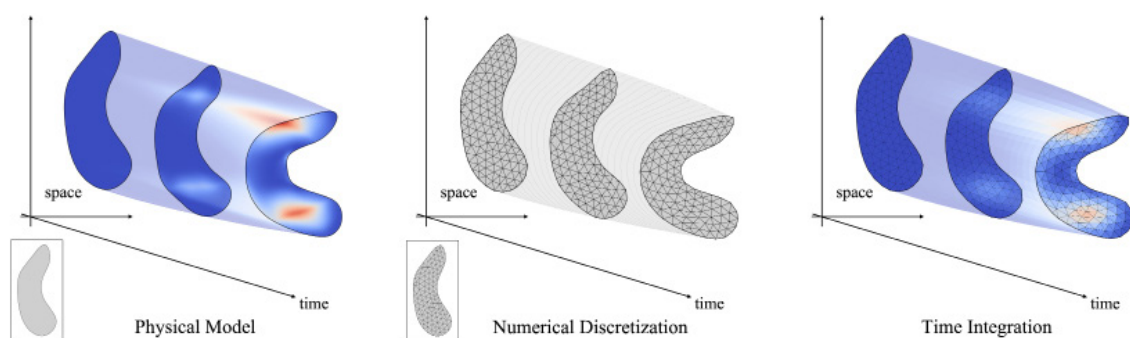
主要完成人：杜韬研究组

物理仿真中的深度学习方法

快速、高精度的复杂多物理系统仿真在众多关键学科前沿有着重要应用。针对此类系统的传统数值算法（如有限元）存在计算量巨大、运行速度缓慢的缺陷。近年来兴起的一系列结合深度学习的物理仿真算法为创建快速物理仿真器提供了新思路，但往往缺乏传统数值方法所具备的正确性和可解释性。如何构建兼具速度和物理精确性的大规模多物理仿真器是一个尚未完全解决的问题，在这一方向的突破将为通用人工智能、虚拟现实、机器人等学科的重要应用带来新的契机。

杜韬研究组在 SIGGRAPH 会议上开设《物理仿真中的深度学习方法》课程对物理仿真与深度学习的发展前沿进行探讨。该课程讨论如何结合物理第一性、计算数值方法、以及现代深度学习模型构建新一代的物理仿真算法，并介绍研究组的相关成果：研究组与麻省理工学院的 ICML 2023 合作论文观察到深度学习模型在物理仿真速度上的优势和精度上的缺陷，提出了利用图神经网络模型作为数值求解器预估值矩阵的新思路，兼具传统数值方法的精确性和深度学习模型带来的速度优势。完整的课程信息将在课程网站 <https://people.iis.tsinghua.edu.cn/~taodu/dl4sim/> 中发布。

该研究成果论文：Tao Du, “Deep Learning for Physics Simulation”, SIGGRAPH 2023.



物理仿真流程概述。左：构建物理模型（以连续介质力学模型为例）；
中：在时空域对物理方程进行数值离散化；右：时间积分求解时空变量数值解。课程将探讨深度学习模型在前述各个阶段的应用。

六、脑启发人工智能

主要完成人：马恺声研究组

基于赫布规则与梯度机制的突触可塑性建模

在生物智能中，突触可塑性指的是生物神经网络的突触权重能根据当前的网络输入及内部状态进行可调控的变化的机制，神经科学研究显示这一机制对于生物的学习以及记忆都至关重要。传统建模工作中对于突触可塑性的研究往往集中于一些局部的学习规则，例如赫布规则（Hebb's rule）及其变种，但近年来一些工作也指出了大脑中实现梯度下降等非局部权重更新机制的可能性。基于此背景，马恺声研究组基于循环神经网络（RNN）设计了不同的可塑性机制（图 1），并探究突触可塑性对于网络记忆能力以及元学习任务表现的帮助。

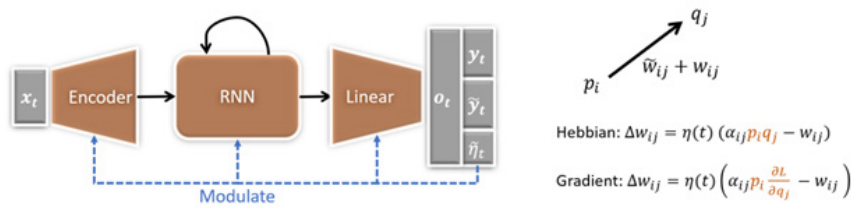


图 1：左图：整体网络架构；右图：本工作中建模的基于赫布规则或者基于梯度信息的可塑性机制。

在实验部分，研究人员对一般 RNN 与有突触可塑性的 RNN 在各种序列任务上进行了对比。在序列记忆（sequential memory）、连结记忆（associative memory）、小样本图像识别（few-shot image classification）以及回归任务（few-shot regression）中，基于赫布规则或者基于梯度信息的可塑性机制均表现出了对于任务表现的明显增益，说明了突触可塑性在提升人工神经网络的记忆能力和小样本学习能力上的作用。对比两种不同的可塑性机制可以发现局部的赫布规则相对更擅长记忆任务，而基于梯度的更新规则则更适合更为复杂的学习任务，这也为今后在生物以及人工神经网络中突触可塑性的研究提供了启发。

该成果研究论文：Yu Duan, Zhongfan Jia, Qian Li, Yi Zhong, Kaisheng Ma, “Hebbian and Gradient-based Plasticity Enables Robust Memory and Rapid Learning in RNNs”，ICLR 2023.

七、算法理论

主要完成人：段然研究组

更快的矩阵乘法算法

段然副教授与计科 80 班（现加州伯克利大学博士生）武弘勋、计科 02 班周任飞同学合作给出了更快的矩阵乘法算法，新的复杂度为 $O(n^{\{2.371866\}})$ ，改进了之前的复杂度 $O(n^{\{2.372860\}})$ [Alman, V. Williams 2020]，为近十年来最大的改进幅度（见下表）。

矩阵乘法是算法领域最基础的问题之一，不仅很多矩阵运算都可归约到矩阵乘法，而且很多组合问题都可用矩阵乘法加速。矩阵乘法的时间复杂度一般被写做 $O(n^{\{\omega\}})$ 。从 Strassen 算法开始，人们就想知道 ω 的真正数值，而且很多人相信最终 $\omega=2+o(1)$ 。Coppersmith 和 Winograd 在 1990 年提出了 CW 张量，通过 laser method 在一阶和二阶下分别给出了 $\omega < 2.387190$ 和 $\omega < 2.375477$ 的界，后来人们研究了更高阶的 CW 张量，并在 32 阶下给出了 2.3728639 的界 [Le Gall 2014]。[Alman, V. Williams 2020] 研究了在 4 阶以上边际分布不能完全确定联合分布的问题，并通过改进哈希方法部分弥补了这个问题，复杂度改进到了 $\omega < 2.3728596$ 。

高阶的 CW 张量能得到更好的结果是因为高阶张量能够将低阶张量的矩阵合并成更大的矩阵。但该研究组的研究发现高阶的 CW 张量中所计算的矩阵总大小要小于低阶，因为各部分分别优化带来的分裂分布不平衡等原因。该研究组利用非对称哈希方法部分弥补了这个问题，并在 2 阶下给出了 $\omega < 2.374631$ 的界。但是非对称方法会大大增加参数优化的难度。经过复杂的优化算法，在 8 阶下该研究组给出了 $\omega < 2.371866$ 的界。该研究组的结果“打破”了之前的一个高阶 CW 张量的 2.3725 的下界 [Ambainis, Filmus, Le Gall 2014]，因为他们改变了高阶 CW 张量内部低阶张量的计算方式。

去年 DeepMind 在 Nature 上发文称他们通过强化学习算法找到了更快的矩阵乘法算法，受到了广泛关注。他们没有考虑 border rank 和 laser method 等方法，只是基础的“Strassen-like”分块法。在二元域 F_2 上他们给出了 47 步乘法计算 $4*4$ 的矩阵乘法，提高了原始 Strassen 方法的 49 步，复杂度为 $O(n^{\{2.7773\}})$ 。（因为是在 F_2 上，不能直接得到 ω 的界。）值得一提的是 Victor Pan 曾在 1982 年给出了 36133 步乘法计算 $44*44$ 的矩阵乘法，复杂度为 $\omega < 2.7734$ ，仍是目前最快的 Strassen-like 方法。

该成果研究论文：Ran Duan, Hongxun Wu, Renfei Zhou, “Faster Matrix Multiplication via Asymmetric Hashing”, FOCS 2023 to appear.

Author	ω
<i>Strassen 1969</i>	2.8074
<i>Pan 1978</i>	2.796
<i>Bini, Capovani, Romani 1979</i>	2.780
<i>Schönhage 1981</i>	2.522
<i>Coppersmith, Winograd 1981</i>	2.496
<i>Strassen 1986</i>	2.479
<i>Coppersmith, Winograd 1990</i>	2.375477
<i>V. Williams 2013</i>	2.372927
<i>Le Gall 2014</i>	2.372864
<i>Alman, V. Williams 2020</i>	2.372860
<i>Duan, Wu, Zhou 2023</i>	2.371866

八、计算机系统结构

主要完成人：高鸣宇研究组、马恺声研究组

对近存计算系统中远程访问和负载均衡的协同优化

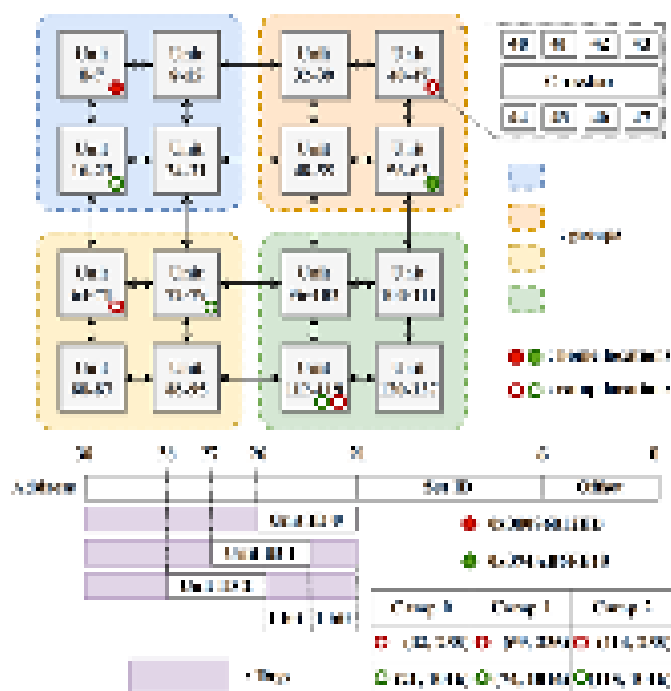
在“内存墙”时代，采用近存计算的系统架构设计范式可缓解数据密集型应用的内存访问瓶颈。基于 3D 内存的近存计算系统通常由大量并行的处理单元组成。当前此类系统存在两个主要问题，即远程内存访问的高昂开销和负载不均衡带来的性能下降。这两个问题互相关联耦合，在解决时需要做权衡取舍，现有方案在缓解其中一个问题时往往会导致另一个问题的进一步恶化。

高鸣宇研究组提出了一个协同优化远程内存访问问题和负载均衡问题的近存计算架构。为了实现更加灵活的数据和计算调度,该研究组使用了一个基于细粒度任务的编程模型和执行模型。该研究组在该任务模型中封装了每个任务的访问数据信息和计算负载信息,为数据缓存和计算调度提供了便利。在此基础上,搞研究组提出了两点优化。为了减少远程内存访问的开销,该研究组提出了一个分布式的 **DRAM** 缓存方案,通过把数据复制在更近的位置来减少远程访问开销。该研究组还提出了一个综合考虑远程内存访问和负载均衡的调度算法,此调度算法和缓存方案协同设计,能够灵活利用数据复制带来的新的调度机会。该系统相比传统近存计算系统可达到平均 1.7 倍的性能提升和 25% 的能耗节省。

该研究成果论文: Boyu Tian, Qihang Chen, Mingyu Gao, “ABNDP: Co-optimizing Data Access and Load Balance in Near-Data Processing”, ASPLOS 2023.



系统整体架构



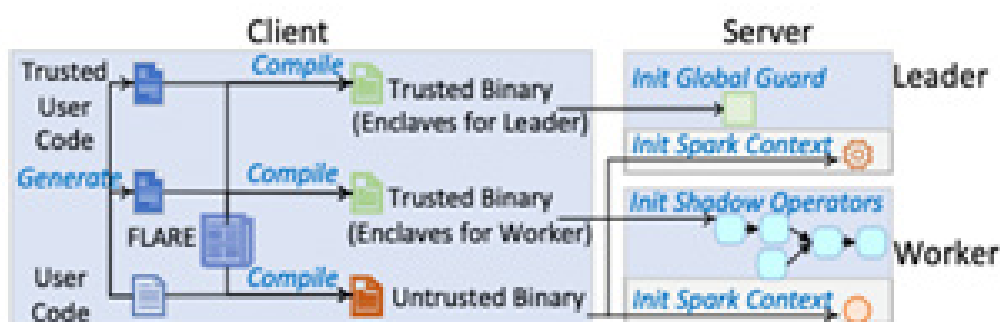
分布式缓存方案

基于可信执行环境的分布式安全计算框架

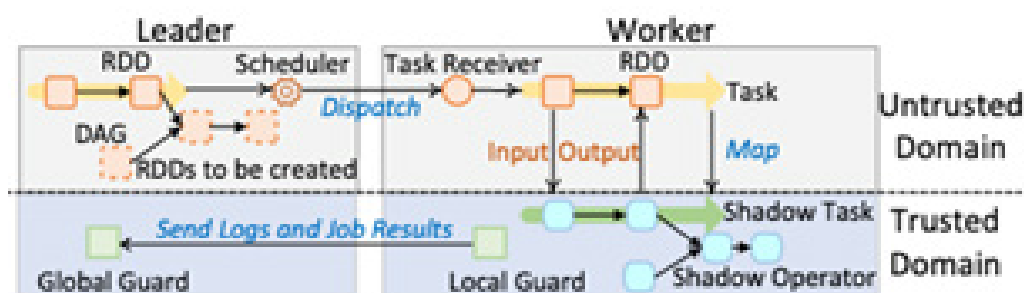
在云计算中，数据的计算可以被外包给更强大的远程服务器。然而，将待计算的数据发送到云端将不可避免地造成数据泄露，从而引发数据安全问题。当前广泛使用的数据分析框架（如 Spark）支持高性能分布式处理，但并不能保护数据隐私。虽然一些密码学方案可以解决隐私问题，但是会带来不可接受的计算开销。可信执行环境（TEE）的出现为隐私计算提供了另外一种途径，但仍然面临诸如性能损失，可信基过大等问题。

由此，高鸣宇研究组设计了一个分布式安全计算框架 FLARE，在保证安全性的同时优化了性能。首先，高鸣宇研究组将系统代码划分为可信和不可信两个部分，可信部分直接处理数据内容，这一部分放在可信执行环境中运行；而系统的大部分代码都划分为不可信部分，由此显著减小了可信基，降低了可信执行环境内出现代码漏洞的可能性。其次，该研究组提出了名为影子算符的抽象数据结构，该结构支持数据在多个算子之间平滑流动，无需频繁进出可信执行环境，从而降低了开销。该研究组还设计了内存高效的优化技术，根据不同计算模式和内存分配频度来定制数据划分与并行执行方式。最后，该研究组通过对程序执行图进行校验，以及对不经意模式的支持，进一步增强了框架的安全性。实验结果表明，在不同的基准测试集上，FLARE 比当前最先进的隐私计算框架的性能提升了 3 到 176 倍。同时，对比将整个系统框架全部放入可信执行环境中的做法，FLARE 在保证更高安全性的同时，也提供了 2.8 到 28.3 倍的性能增益。

该研究成果论文：Xiang Li, Fabing Li, Mingyu Gao, “FLARE: A Fast, Secure, and Memory-Efficient Distributed Analytics Framework”, VLDB 2023.



系统初始化过程



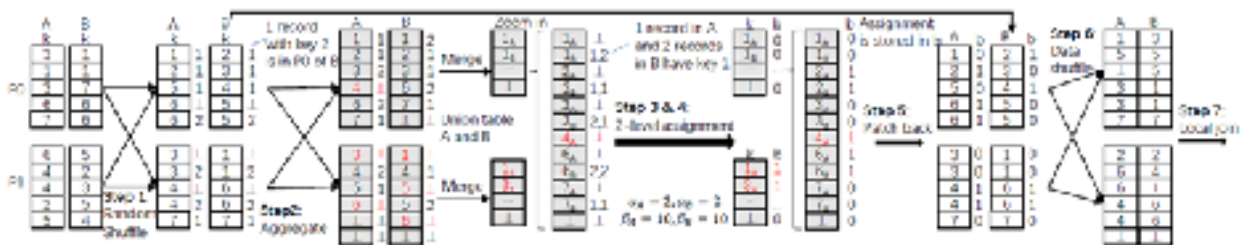
系统执行过程

高效的不经意分布式数据分析算法

由于在云计算中有着较高的数据暴露风险，为了数据安全，在实际应用中通常采用可信执行环境（TEE）这一硬件技术来保证数据在计算过程中的机密性和完整性。其相对于密码学方案具有较低的性能开销。然而，在更强的威胁模型下，仅通过加密保护数据内容本身仍然不能保证数据不被泄露。攻击者仍可通过数据访问模式这种侧信道推断出数据内容。一种常用的方法是将 ORAM 协议直接和现有的算法结合，即，将算法中每次普通的访存转换成对 ORAM 的访存。这种做法有两个主要缺点，一是会带来巨大的性能开销，二是仍然存在安全隐患。比如对于分布式的算法，单纯采用 ORAM 是不足以保证安全性的，因为机器间的数据传输路径以及传输量仍然会泄露一些信息。

由此，高鸣宇研究组为数据分析中的常用算子（filter, aggregate, join）定制化设计了不经意分布式算法集合 SODA。这些算法无论是在每个节点内还是节点间通信方面都具有访问模式跟数据内容无关的特性，从而使攻击者无法从访问模式反推出数据内容。首先，该研究组设计了伪随机通信，其中输入分区中的记录以（伪）随机方式发送到每个目标节点，再通过少量填充使得通信量不暴露数据信息。全局排序不再被使用，从而降低了性能开销。过滤（filter）和聚合（aggregate）都可直接运用伪随机通信。而对于连接（join），该研究组进一步设计了两级分配方案，使各个节点间通信量达到均衡。实验表明，在多算子的工作负载上，SODA 相对于现有方案实现了 1.1 倍到 14.6 倍的加速比。

该研究成果论文：Xiang Li, Nuozhou Sun, Yunqian Luo, Mingyu Gao, “SODA: A Set of Fast Oblivious Algorithms in Distributed Secure Data Analytics”, VLDB 2023.



不经意连接（join）算法

平铺加速器层间调度空间的定义和探索

深度神经网络（DNN）在图像识别、目标检测和自然语言处理等领域中得到了广泛应用。随着 DNN 的不断发展，网络结构变得更加复杂，因此需要大规模的加速器来加速推理过程。目前，采用瓦片式架构的大规模加速器已经成为主流，其中每个硬件瓦片（HW-tile）包括一个处理元素（PE）阵列和一个全局缓冲区，并由网络芯片（NoC）相互连接（图 1（b））。然而，单个大型 HW-tile 的利用率和能量效率较低，因此如何有效地利用这些计算和存储资源是一个开放性的挑战。解决这一挑战的关键在于调度，分为层内调度（图 1（c））和层间调度（图 1（c））两类。层内调度研究如何将单个层映射到一个或多个 HW-tile 上，而层间调度研究如何在加速器上调度所有层的计算顺序和资源分配。

尽管层间调度在保持瓦片式加速器高度利用和能量效率方面发挥着越来越重要的作用，但其研究存在显著的不足：大多数研究仍在优化现有的启发式模式，如 LP 和 LS，但没有提出新的模式，更没有清晰而系统地定义瓦片式加速器上层间调度的空间。缺乏层间调度空间的明确定义极大地限制了优化瓦片式加速器性能和能量效率的机会。此外，缺乏系统性定义也阻碍人们理解不同层间调度选择如何影响不同硬件行为以及这些行为如何进一步影响加速器的能量效率和性能。

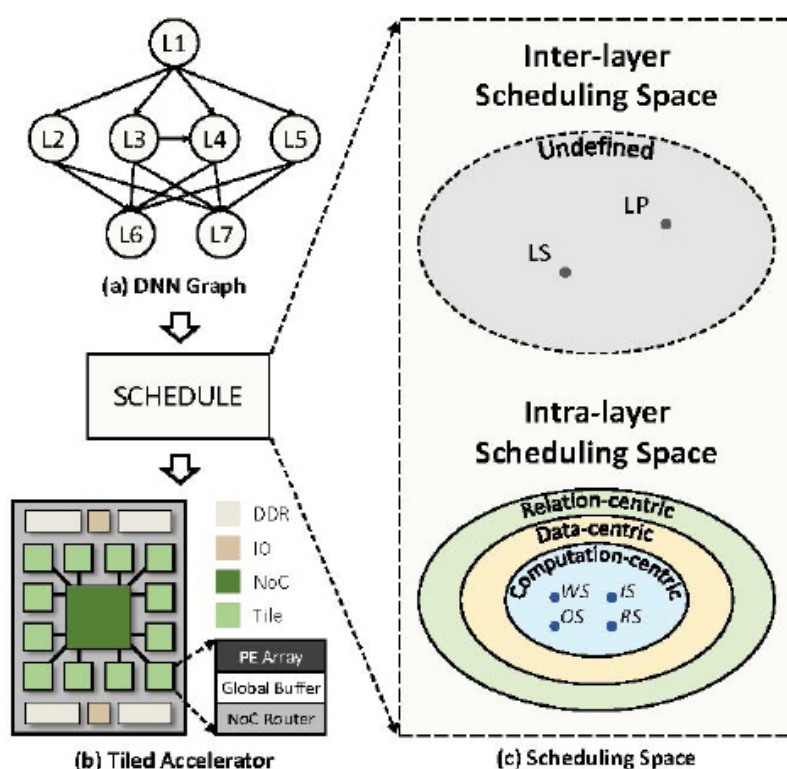


图 1. 在平铺加速器上的调度

马恺声研究组针对这一挑战而引入一种统一而系统的资源分配树符号来描述瓦片式加速器上具有不同结构的推理 DNN 的层间调度空间。该符号包括时间切割，将相同的 HW-tile 组和不同的计算时间间隔分配给每个子节点，以及空间切割，将不同的 HW-tile 组分配给每个子节点。每个 RA 树都是一个切割和叶节点（DNN 的层）的分层组织。然后，研究人员详细阐述如何将树形结构解析为相应的资源分配方案和数据流动。

基于 RA 树符号，研究人员全面分析了不同的层间调度选择如何影响硬件行为以及这些行为如何影响加速器的能量效率和性能。此外，研究人员在符号中表示了现有的 LS 和 LP 模式并分析了其特征。结合上述内容，研究人员开发了一个端到端且高度可移植的调度框架 SET，用于自动探索瓦片式加速器上 DNN 的整个调度空间。为了有效地探索新定义的层间调度空间，研究人员将 SET 配备了一种基于模拟退火的算法，具有 6 个特别设计的操作。SET 可以轻松地移植到各种瓦片式加速器上，具有良好的可移植性。研究人员已经开发了一个端到端的 SET 部署流程用于测试芯片。该框架可在 <https://github.com/SET-Scheduling-Project/SET-ISCA2023> 上获取。

与 SOTA 开源框架 Tangram 相比，SET 平均可实现 1.78 倍的性能提升和 13.2% 的能量成本降低。此外，研究人员对不同的 DNN、批次大小和硬件平台进行了大量实验，以展示 SET 的通用性和探索新定义空间相对于现有 LS 和 LP 调度模式的效果。此外，研究人员利用 SET 分析了 LS 和 LP 的特性，并揭示了层间调度空间的特征。

该成果研究论文：Jingwei Cai, Yuchen Wei, Zuo tong Wu, Sen Peng, and Kaisheng Ma, “Inter-layer Scheduling Space Definition and Exploration for Tiled Accelerators”, ISCA 2023.

一种可扩展的高效芯粒（Chiplet）互连网络设计方法

Chiplet（芯粒）方法可以加速 VLSI 系统的开发，并提供更好的灵活性。然而，在不同层次拓扑的系统上建立跨多个芯粒的互连网络并保持高性能的无死锁路由并不容易。大多数片上网络都是基于平面拓扑结构（如二维网格），这对于大规模的多芯粒系统来说是不灵活和低效的。为了充分利用多芯粒架构和先进封装，马恺声研究组提出了一种互连方法，可以利用典型的基于 2D-mesh 片上网络的芯粒灵活地建立高维互连系统（如图 1 所示）。该研究组提出了一种基于负优先的无死锁自适应路由算法和一种 safe/unsafe 的流量控制策略。此外，该研究组还采用了一种通用的网络交织方法来平衡芯粒内部和芯粒之间的通信带宽。该研究组开发了一个专门为多芯粒互连设计的仿真器以评估不同的架构。与传统的二维网格自适应路由相比，该研究组的方法在各种网络配置下都显著提升了性能（如图 2 所示）。

该成果研究论文：Yinxiao Feng, Dong Xiang, and Kaisheng Ma, “A Scalable Methodology for Designing Efficient Interconnection Network of Chiplets”, HPCA 2023.

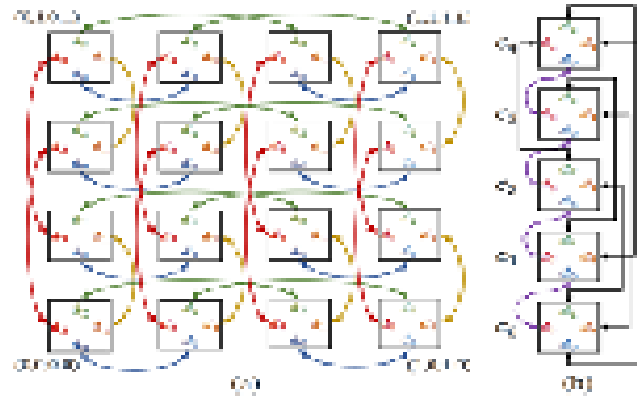


图 1 芯粒互连网络

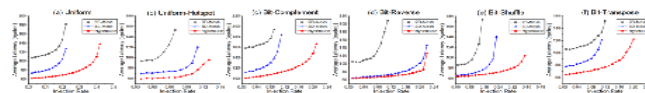


图 2 不同负载下的网络性能

一款 28nm 工艺下 68MOPS 0.18 μ J/Op 的密文比特流稀疏计算 Paillier 同态加密处理器

云计算可为大量新兴的信息应用提供可靠的高性能服务，但在计算过程中需要处理大量的个人和机构数据。Paillier 同态加密可以允许直接对密文进行计算，避免个人和机构的隐私数据泄露。在图 1 示出的 Paillier 同态加密方案中，客户端将其明文加密为密文，随后将密文发送至服务器，服务器执行同态求值并且向客户端返回加密结果。

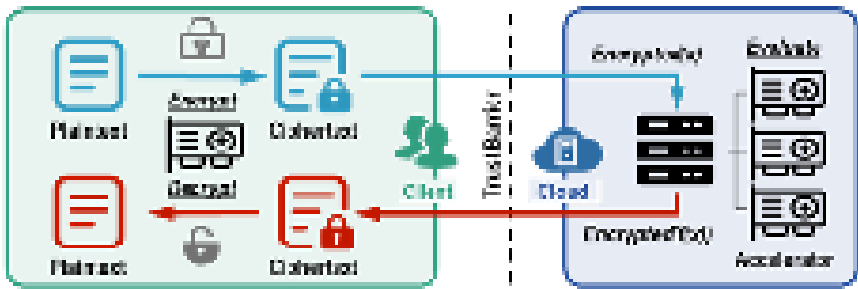


图 1: 同态加密计算流程

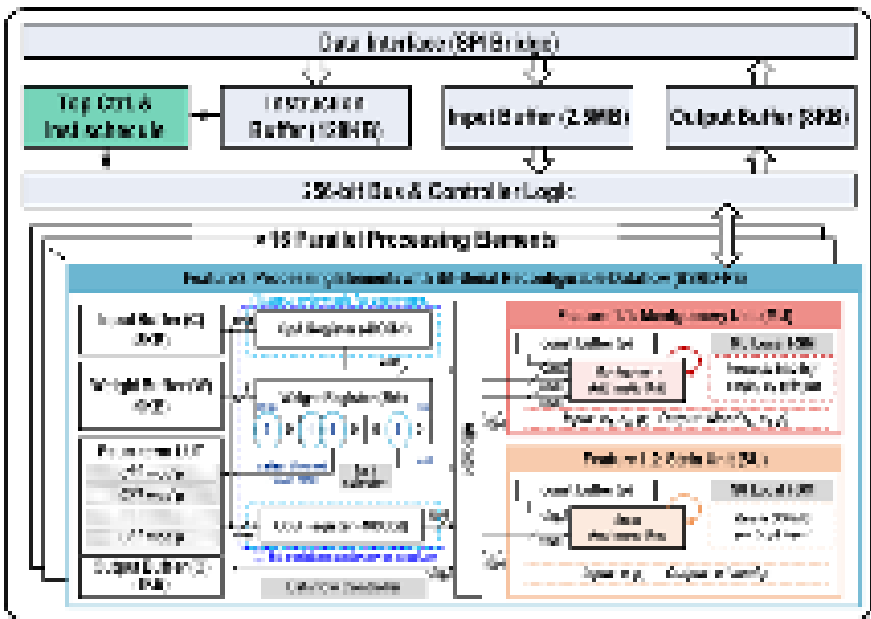


图 2 同态加密处理器架构

目前 Paillier 同态加密算法在通用处理器 CPU 和图形处理器 GPU 上的性能较低。针对这一计算瓶颈，马恺声研究组在 28nm 工艺下设计并制造了面向 Paillier 同态加密的领域专用处理器 PH-EPU。该处理器采用密文比特流稀疏计算架构 (图 2)，PH-EPU 与具有 16 核的 Intel 的台式 CPU i9-9900 相比能到 14.9 倍的加速，为隐私保护云计算提供了高性能的解决方案。

该成果研究论文：Guiming Shi, Zhanhong Tan, Dapeng Cao, Jingwei Cai, Wuke Zhang, Yifu Wu, Kaisheng Ma, “A 28nm 68MOPS 0.18 μ J/Op Paillier Homomorphic Encryption Processor with Bit-Serial Sparse Ciphertext Computing”, ISSCC 2023.

九、数据库系统

主要完成人：张焕晨研究组

基于子查询选择的高效的查询重优化

查询重优化是一种用于解决物理执行计划中的预测与实际执行情况不相符这一问题的方法。一个查询通过查询优化器生成物理计划，再被执行引擎执行。而在应用中，物理计划所预测的耗时常常和实际耗时相差甚远，其主要原因是查询优化无法准确预测部分操作符执行时候所产生的结果的大小。而为了解决物理计划与实际不符的情况，除了直接提高对结果大小估计的准确率外，重优化被广泛应用在商业数据库中。重优化可以在物理计划执行到某一步时，通过采集该中间结果的统计数据，判断预测与实际是否产生的较大的误差。如果误差很大，则利用新数据重新将未执行的部分做查询优化，并让执行引擎执行新的执行计划。

张焕晨研究组通过分析和初步实验，发现现有的重优化架构依然存在问题。在初次查询优化与第一次重优化之间，由于初始执行计划的误导，执行引擎可能已经执行了一些结果大小很大的操作，使得后续无论如何重优化都难以得到比较高效的执行计划。如下图 1 所示，在给定的查询集中，超过四分之一在第一个连接操作之前就与最优计划产生了分歧，而在第二个连接操作时，超过一半的物理计划都与最优计划不同。这个结果说明，查询重优化被初始物理计划所诱导，从而执行非最优的子查询是很常见的现象。

Similarity	0	1	2	> 2
Ratio	13%	12%	32%	43%

图 1 初始物理计划与最优计划产生分歧的位置及其比例

为了解决该问题，张焕晨研究组设计了一种基于子查询选择的重优化算法，QuerySplit。QuerySplit 首先将输入的查询分成若干部分重叠的子查询，并通过利用外键的性质，确保每一个子查询不会产生很大的结果大小。当一个外键表与其主键表做连接时，其结果大小总不会超过外键表的大小。再利用启发式的算法，兼顾子查询的即时收益和未来收益，决定哪一个子查询将被执行。每次随着被选定的子查询执行完毕，QuerySplit 会修改并重新优化剩余的子查询，重新决定执行顺序。

最终的实验结果如图 2 所示，QuerySplit 在 PostgreSQL 上的表现超过了现有的重优化技术，和一些先进的查询优化技术。并且与已知所有中间结果大小的优化器所产生的执行计划的执行时间相差不到 4%。

该成果研究论文：Junyi Zhao, Huanchen Zhang, and Yihan Gao. "Efficient Query Re-optimization with Judicious Subquery Selection." SIGMOD 2023.

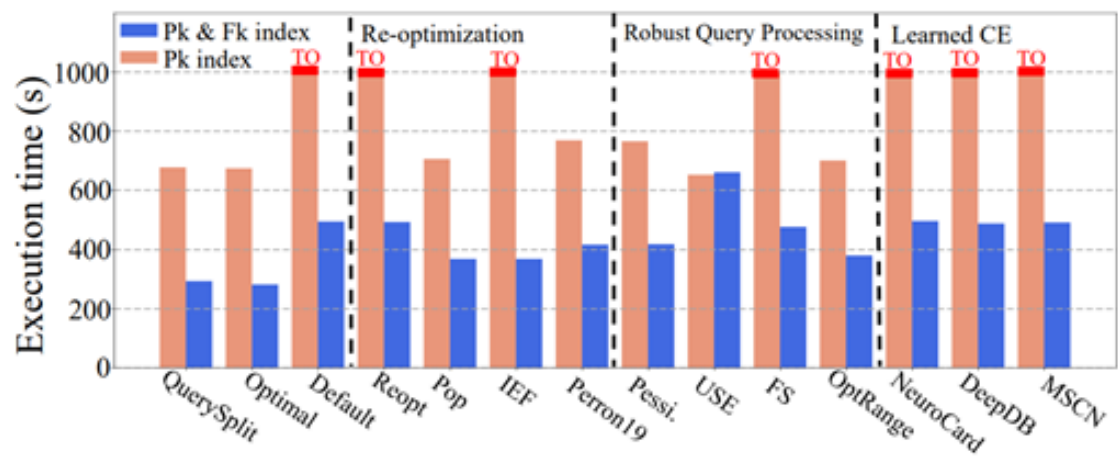


图 2 QuerySplit 与其他查询优化技术的对比

十、区块链

主要完成人：房智轩研究组

抽卡销售机制的设计与分析

抽卡销售游戏 (Gacha Game) 是一种特殊的不透明销售方式，卖方通过销售 Gacha 抽取次数给买方。每次抽取都提供了一定的 (可变) 概率，让买方有机会获得 Gacha Game 的奖励。抽卡销售游戏广泛应用于各类销售中，如彩票和游戏中虚拟商品的销售。

这对这种复杂的销售模式，房智轩研究组首次研究提出了严密的数学建模框架，将买方的顺序决策建模为马尔可夫决策过程 (MDP)。模型引入了巨鲸属性 (Whale Property) 的定义，即一个 Gacha Game 具有巨鲸属性当且仅当玩家的最优决策是要么完全不玩，要么玩到获得奖励为止。研究组进一步证明了巨鲸属性是卖家收益最大化的必要条件。

此外，作者展示了 Gacha Game 与单一物品、单一买家拍卖 (single-item single-bidder auction) 的等价性，从而导出能够实现最大卖方收益的最优参数。此外，该工作的结果还证实当玩家有预算约束时，Gacha Game 可以比拍卖带来更高的卖方收益。

最后，房智轩研究组讨论了 Gacha Game 在计算机系统中的一个例子。文章发现区块链系统可以被建模为一个 Gacha Game。其中，卖家对应区块链系统，而买家则对应矿工。区块链系统 (卖家) 希望最大化系统的安全性 (即“利润”)，而最大化系统安全性的方法就是通过增加买家 (矿工) 的哈希尝试次数或者资产抵押数量 (即卖家的抽取次数)。这种关联揭示了区块链系统安全性与用户投入的内生关系，提供了一种利用拍卖理论指导区块链最优设计的潜在研究方向。该成果发表在计算机性能建模与分析的顶级会议 ACM SIGMETRICS 2023 上。

该成果研究论文：Canhui Chen and Zhixuan Fang, “Gacha Game Analysis and Design,” POMACS 2023. [Accepted directly through ACM SIGMETRICS 2023 as the full version]

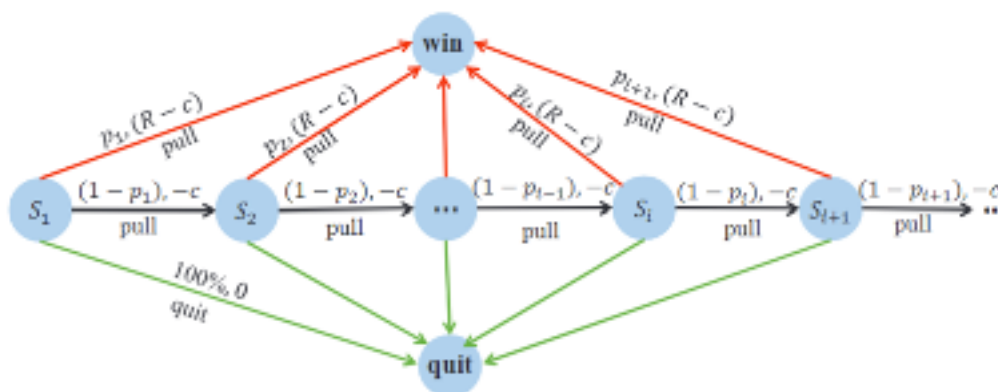


Fig. 1. Markov Decision Process in Gacha Game

【量子信息】



一、离子阱量子模拟

主要完成人：段路明研究组、徐勇研究组、吴宇恺研究组

首次在离子阱系统中实验观测到非厄米复数能谱及其拓扑结构

段路明教授研究组与徐勇助理教授合作,在离子阱系统中,首次实验测量了非厄米系统的复数能谱以及其拓扑结构。

非厄米系统由于其独特的拓扑特性而广受关注。非厄米系统通常具有复数能量，并且其复数能量可能具备诸如链环或纽结的拓扑结构，这些拓扑结构是能量始终为实数的厄米系统所不具备的。近来，由于其独特的拓扑性质以及潜在的量子信息处理应用，在离子阱、冷原子、超导电路以及固态自旋系统等量子模拟平台中实现非厄米哈密顿量已成为一个重要的研究目标，并且已取得很大进展。然而，非厄米能谱拓扑结构的决定性特征尚未被实验探测到。实际上，实验测量非厄米系统的复数本征能量仍然是一个重大挑战，这使直接探测能谱的拓扑性质变得十分困难。最近，徐勇研究组提出一种被称为非厄米吸收光谱学的方法，用于测量非厄米量子系统中的复数能谱，这在离子阱等量子平台中测量复数能谱拓扑提供了可能性。

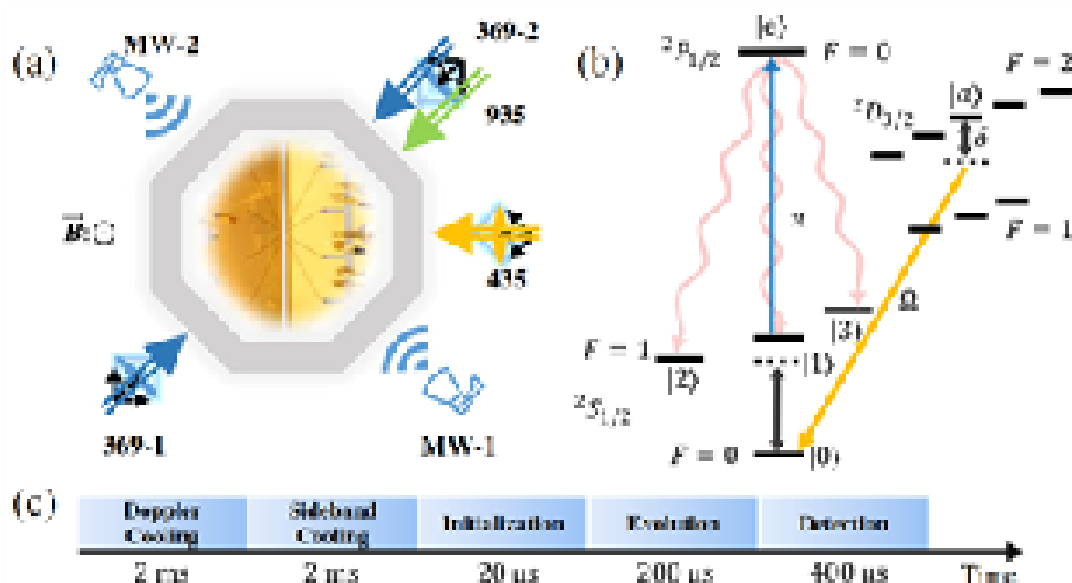


图 1: 实验系统示意图

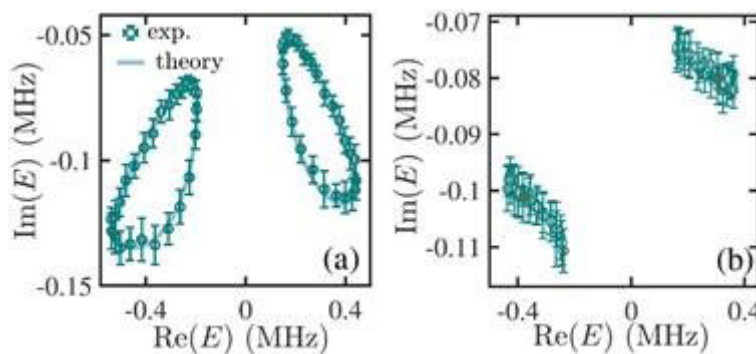


图 2: (a) 复数能谱 unlink 拓扑结构。(b) 拓扑平庸结构

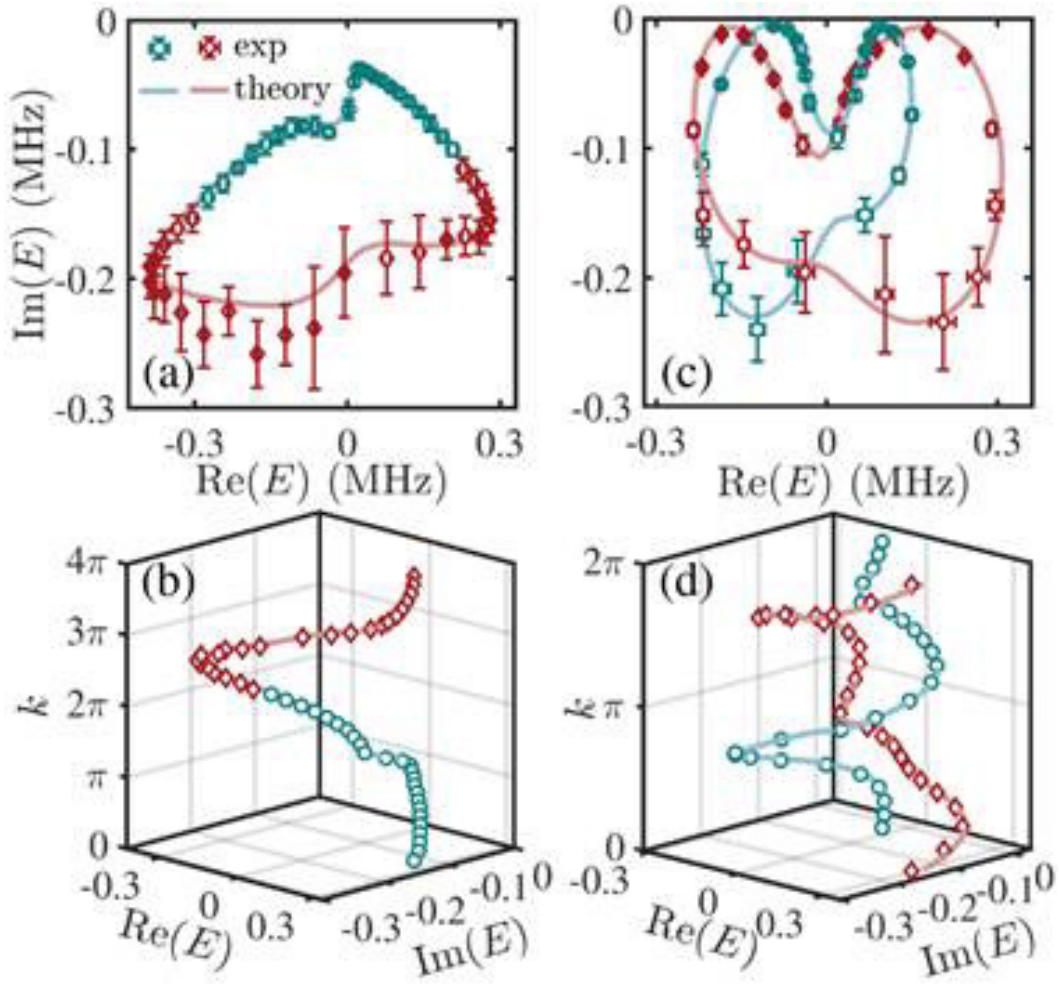


图 3: (a, b) unknot 拓扑结构。(c, d) Hopf link 拓扑结构

段路明教授研究组基于单个 $^{171}\text{Yb}^+$ 离子实现了双能带非厄米模型，其复数本征能量具有 unlink、unknot 或 Hopf link 拓扑结构。该模型的厄米部分由微波脉冲驱动基态超精细能级间的跃迁来实现（图 1a MW-1, MW-2; 图 1b 黑色双箭头），而非厄米部分则通过用共振激光使离子从基态跃迁至激发态（图 1a 369-1, 369-2; 图 1b 蓝色箭头），激发态自发辐射导致布居流失实现了非厄米部分（图 1b 波浪线）。为测量系统的复数能量，基于非厄米吸收光谱学方法，研究人员将离子制备在 $2D_{3/2}$ 能级（辅助能级）上，然后通过弱激光将其与系统能级相耦合（图 1a 435; 图 1b 黄色箭头），并在经过长时间后，测量离子处于辅助能级的概率。通过拟合测量到离子的概率与失谐的曲线，可以提取出非厄米系统的复数能量。实验测量得到的复数能量呈现出 unlink、unknot 或 Hopf link 的拓扑结构（图 2, 图 3），且实验测量值与理论值符合得很好。此方法可以直接推广到其它量子模拟平台，如冷原子、超导电路或固态自旋系统，因此该实验为探索非厄米量子系统中各种复数能量性质开辟了一条道路。

该成果研究论文：M.-M. Cao, K. Li, W.-D. Zhao, W.-X. Guo, B.-X. Qi, X.-Y. Chang, Z.-C. Zhou, Y. Xu, and L.-M. Duan.

“Probing Complex-Energy Topology via Non-Hermitian Absorption Spectroscopy in a Trapped Ion Simulator” . Phys. Rev. Lett.

130, 16300 1 – Published 21 April 2023.

离子阱量子模拟长程横场伊辛模型的临界行为

离子阱系统是当前量子计算研究中最领先的物理平台之一。受限于目前的量子操控精度和实验噪声，此前的多离子量子模拟实验通常局限于验证多体量子系统的定性性质。为实现近期的量子计算、量子模拟机的现实应用，一个重要的问题是如何利用量子模拟来研究复杂量子系统的定量性质。

量子相变的临界指数是一类对实验噪声不敏感的多体量子系统的定量性质，适合于近期的带有噪声的量子模拟机的应用。本工作中，研究人员利用 61 个囚禁离子实现了长程横场伊辛模型的量子模拟，超越了此前美国马里兰大学研究团队保持的 53 离子量子模拟的世界纪录。研究人员进而利用 Kibble-Zurek 机制，在实验中测量了该模型的量子相变的临界指数，首次在离子阱系统实现了临界指数的定量研究。

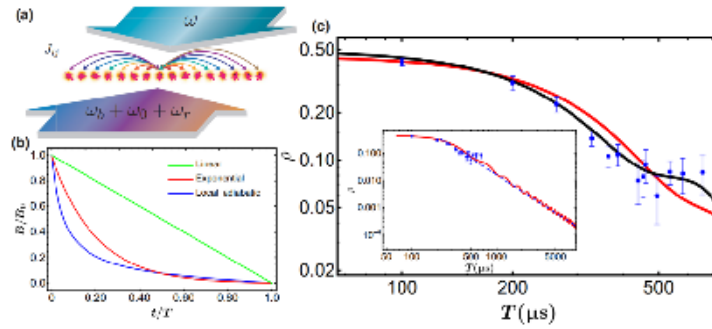


图 1 (a) 利用拉曼激光实现多离子长程横场伊辛模型的量子模拟。(b) 设计优化的路径扫描实验参数，以缩短所需的实验时间。(c) 以 2 离子为例，通过改变实验参数的扫描时间，可以从末态的观测量提取量子相变的临界指数。

如图 1 所示，研究人员利用含有多个频率成分的拉曼激光，实现多离子长程横场伊辛模型的量子模拟。初始时，横场项 B_0 远大于伊辛相互作用 J_0 ，系统处于顺磁相。研究人员进而缓慢减弱横场项，并根据系统的能隙大小设计合适的演化路径，最终让横场项降低至 0，系统根据伊辛相互作用 J_0 的符号将处于铁磁或反铁磁相。Kibble-Zurek 机制表明，如果把系统的初态制备成顺磁相的基态，通过改变横场项的演化时间 T ，系统的末态将相应改变，在合适的参数区域内，末态的拓扑缺陷密度，或是自旋的空间关联长度的倒数，都将正比于 $T^{-\mu}$ ，其中 μ 即为待测的临界指数。

研究人员对于铁磁和反铁磁相互作用的长程横场伊辛模型分别应用上述方法，改变不同的演化时间 T ，测量自旋的空间关联函数随离子对间距的变化，进而提取相应的空间关联长度。如图 2 所示，对于铁磁相互作用模型，随着离子数增多，系统的有限尺度效应减弱，对于 36 至 61 离子的情况获得了类似的临界指数 $\mu \approx 0.42$ ，也与此前理论工作给出的 $\mu \approx 0.45$ 相符。而对于反铁磁相互作用的情况，此时由于量子阻挫 (quantum frustration) 效应，系统的能隙减小，末态的空间关联变得不明显，难以提取临界指数，但仍能得到与理论预期相符的实验结果。

该成果研究论文：B.-W. Li*, Y.-K. Wu*, Q.-X. Mei*, R. Yao, W.-Q. Lian, M.-L. Cai, Y. Wang, B.-X. Qi, L. Yao, L. He, Z.-C.

Zhou, and L.-M. Duan#, “Probing Critical Behavior of Long-Range Transverse-Field Ising Model through Quantum Kibble-Zurek Mechanism” , PRX Quantum 4, 010302 (2023).

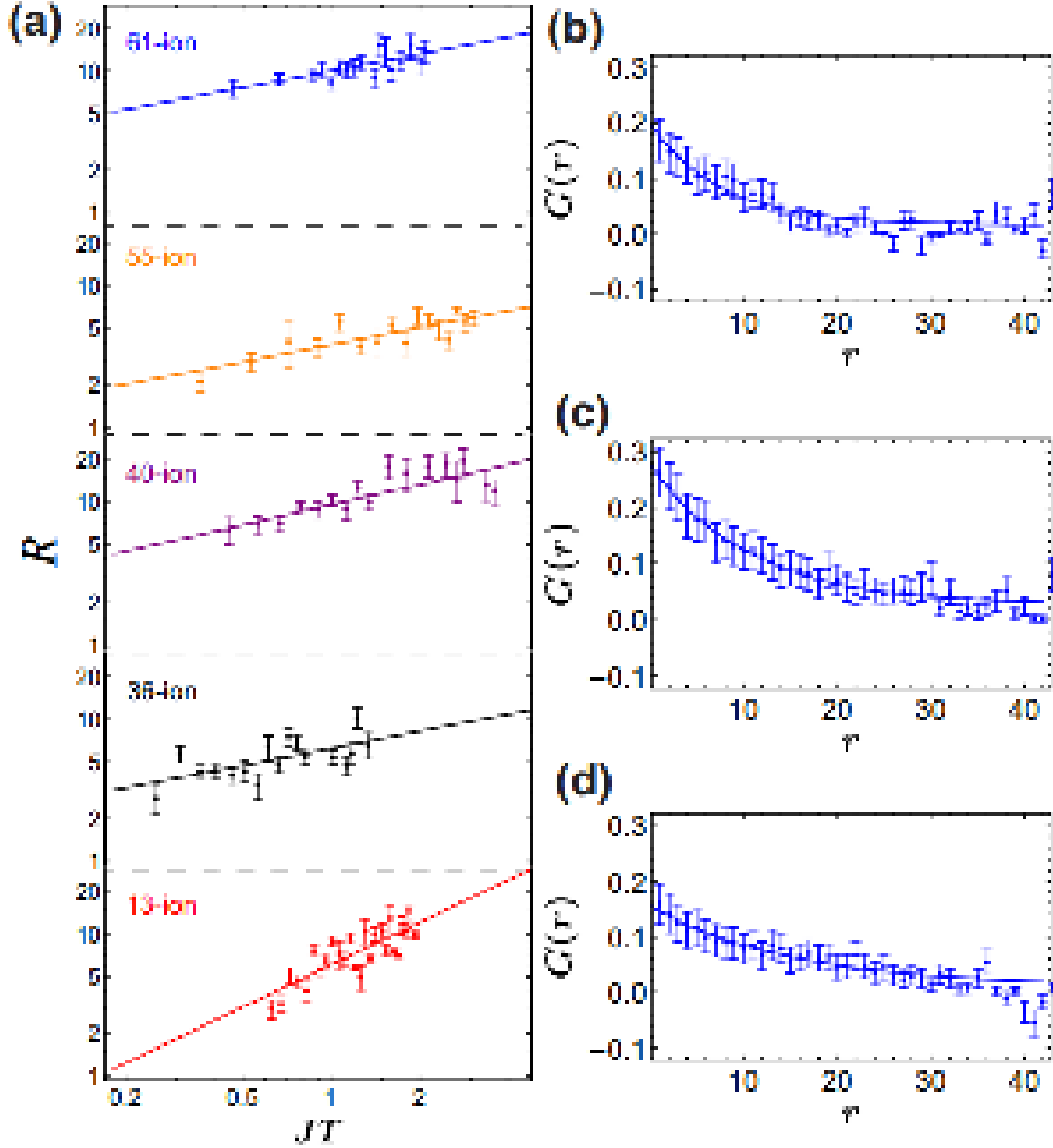


图 2 (a) 不同离子数时自旋的空间关联长度 R 随横场项演化时间 T 的变化。(b-d) 61 离子典型的末态 (b) $T=0.875\text{ms}$, (c) $T=1.75\text{ms}$, (d) $T=2.75\text{ms}$ 时, 自旋空间关联随离子对间距 r 的变化。

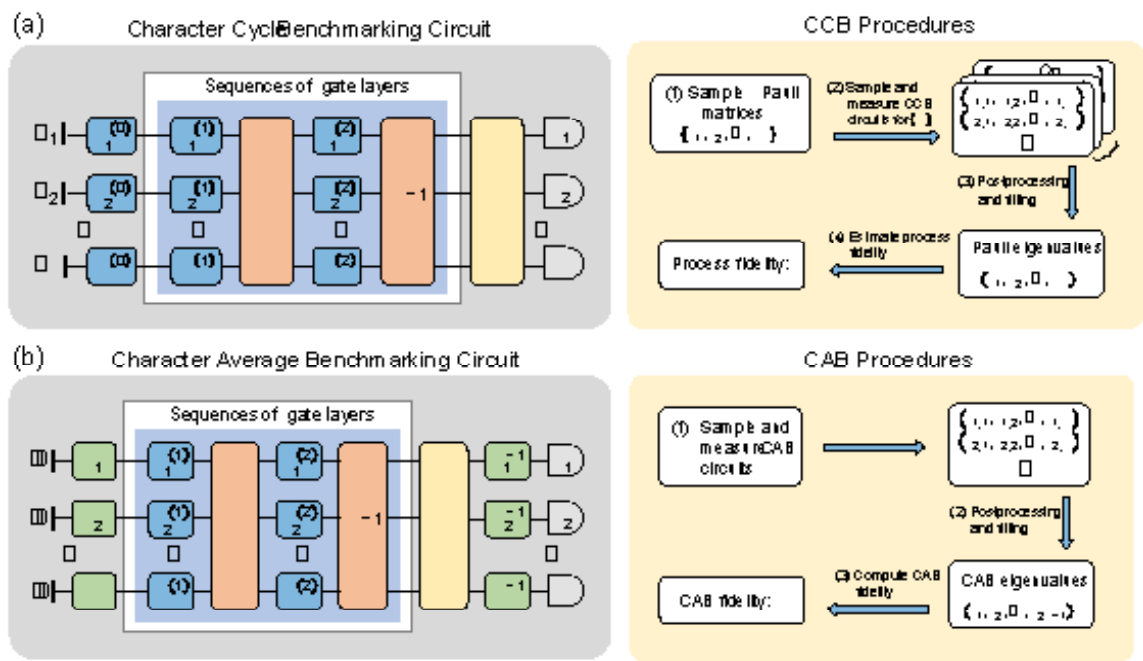
二、量子计算与通信

主要完成人：马雄峰研究组

基于局域旋转的高效可扩展量子门标定方案

作为量子信息科学最重要的方向之一，量子计算被认为在一些特定任务上优于经典计算。基于量子力学的原理，量子计算机使用量子比特（qubits）来代替经典比特来存储和处理信息。量子比特具有叠加性，使其能够同时处于多种状态并行处理任务，从而相比于经典计算可以更快、更有效地执行计算。这个特性使量子计算在一些特定任务中具有指数加速的优势，例如因数分解大数和模拟量子系统等。因此，量子计算被广泛认为是后摩尔时代最重要的方向之一。

实际工作中，随着量子设备规模不断增加，量子噪声的负面影响越发显著，严重阻碍了量子算法的精确实现，对实现量子优势提出了挑战。量子设备误差对于高精度的大规模量子计算机至关重要。随机基准测试（Randomized benchmarking）被提出用于估计量子过程的平均保真度。然而，传统的随机基准测试方法的应用范畴仅限于 Clifford 门集——非通用量子计算门集，并在实际应用中因为需要大量多比特旋转量子门而难以实现。如何有效可靠地估计大规模通用量子过程的保真度仍然是一个悬而未决的问题。



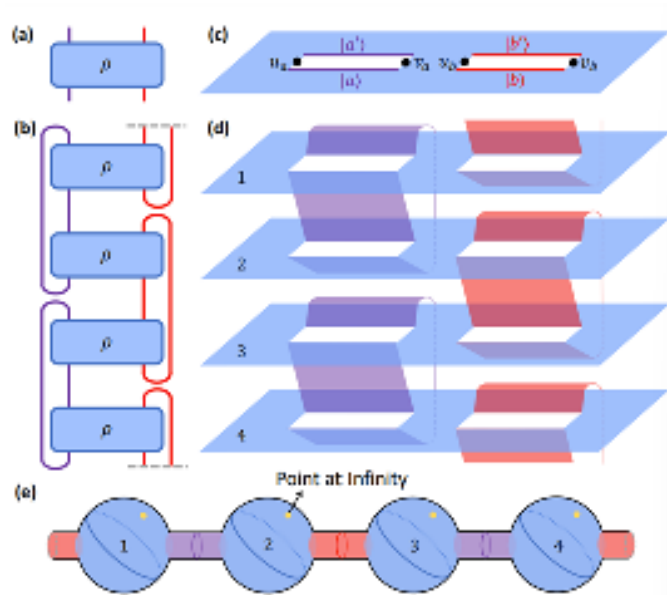
图一：特征标循环标定和特征标平均标定方案示意图。

为了解决这个问题,张艺泓、余文峻、曾培、刘国定以及马雄峰提出并数值验证了两种高效可扩展的量子门标定方案:特征标循环标定 (character-cycle benchmarking) 和特征标平均标定 (character-average benchmarking)。两种方案的量子线路如图一所示。这两种量子标定方法利用了局域旋转理论估计了多比特量子门的过程保真度,相比传统方法具有更高的探测效率和准确度。尤其是特征标平均标定方案用到的量子线路的开始和末尾被添加了额外的局域 Clifford 门,用来大大减少保真度标定的采样复杂度。

该研究成果论文: Yihong Zhang, Wenjun Yu, Pei Zeng, Guoding Liu, and Xiongfeng Ma, “Scalable fast benchmarking for individual quantum gates with local twirling”, *Photonics Research* 11, 81-99 (2023).

1+1 维共形场论中纠缠以及关联的普适度量

量子多体物理中的一个重要问题是研究其纠缠结构。对于可以用量子共形场论描述的一类系统的基态，由于其中的物理量会与某种拓扑结构一一对应，研究者们已经得到了诸多关于其纠缠结构的严格结论。例如一维系统中某一个子区域与环境之间的纠缠可以只有其中心荷以及子区域长度来决定。然而，对于一维共形场论系统的两个不相邻的子区域之间的纠缠，现有的工作主要利用一种基于量子态偏转置的纠缠度量，PPT 纠缠负性，来进行研究。尽管这一度量在量子信息中被广泛采用，并且具有计算以及操作层面的诸多优势，但在共形场论中却对应于一种较为复杂的拓扑结构，导致研究人员很难得到 PPT 纠缠负性的解析结果。直至今日，研究 1+1 维共形场论中两个不相邻子区域的纠缠依然是一个待解决的问题。



图二：1+1 维共形场论中的纠缠计算，基于一种新的纠缠度量。

马雄峰研究组博士生刘振寰与美国科罗拉多大学博德 (CU Boulder) 分校的博士生尹超通过引入一种新的纠缠度量，CCNR 纠缠负性，解决了这一问题，得到了两个不相邻子区域之间纠缠与量子系统有限温下的配分函数之间的普适关系。在量子信息理论中，CCNR 纠缠负性与 PPT 纠缠负性均具有重要地位，且具有相似的纠缠检测能力。但是在 1+1 维量子共形场论中，CCNR 纠缠负性会对应一种更加简单的拓扑结构，环面 (Torus)。于是，对于非常依赖于拓扑结构的共形场论来说，CCNR 纠缠负性相比于 PPT 纠缠负性更容易解析求解。除了纠缠之外，对于 CCNR 纠缠负性的研究对于研究共形场论系统中的总关联也有重要意义。

该研究成果论文：Chao Yin and Zhenhuan Liu, “Universal Entanglement and Correlation Measure in Two-Dimensional Conformal Field Theories”, Physical Review Letters 2023, 130 (13), 131601.

模式匹配量子密钥分发协议的实验实现

量子密钥分发 (Quantum key distribution) 利用量子物理基本原理, 可为通信双方产生理论上无条件安全的随机密钥, 保证了信息传输过程中的安全性。自上个世纪九十年代以来, 经过多年的理论及实验技术的发展, 该方向特别是基于光纤信道的量子密钥分发当前已经进入到实用化阶段。光子作为最普遍应用的信息载体, 其传输损耗是量子密钥分发协议实现的主要障碍。目前已有的传统双模测量设备无关协议 (MDI-QKD) 无法突破传输损耗决定的线性成码上限, 传输距离较低, 而单模相位匹配协议 (PM-QKD) 密钥速率高, 传输距离长, 但需要引入远距离相位锁定技术, 实用性较低。如何在兼顾实用性的前提下有效克服传输损耗从而提高密钥速率、传输距离是量子密钥分发协议理论和实验研究的核心任务。

马雄峰副教授和组内博士生曾培、周泓伊, 本科生吴蔚捷提出了一种新型的模式匹配量子密钥分发 (MP-QKD) 协议, 同时兼顾单模协议的高性能和双模协议的实用性。如图 1 所示, 这个协议中, Alice (以及 Bob) 首先将信号编码在单个光学模式中。然后, 根据 Charlie 的探测响应结果, Alice 和 Bob 对发送的信号进行配对, 提取相对的编码信息。由于 Charlie 只需进行单次干涉即可关联 Alice 和 Bob 的信号, 该协议获得了与单模协议相近的高成码率。另一方面, 由于成码信息编码在配对的信号的相对信息中, 模式配对协议的编码可以容忍更高的光源及链路相位变化, 因而不需要复杂的远距离激光相位锁定技术, 大幅降低单模协议对光源与链路相位控制的要求。

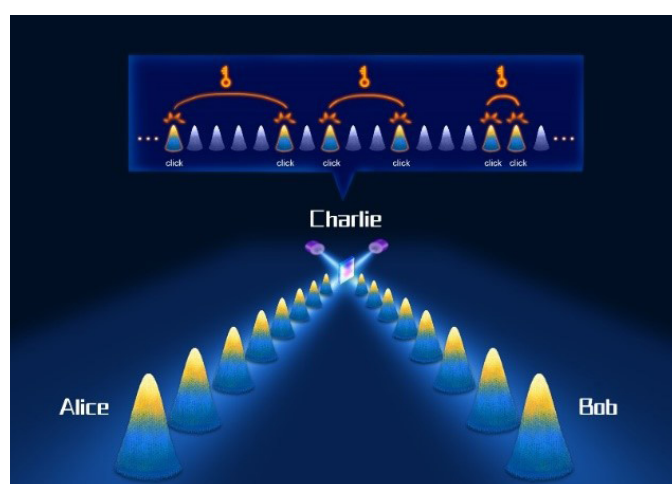


图 1. 模式匹配量子密钥分发协议示意图

基于该理论方案，马雄峰副教授和组内博士生黄溢智与中科大实验团队合作，提出利用极大似然估计的数据后处理方法精确地估算出两个独立激光器的频率差用于参数估计，并结合中科院上海微系统所研制的高效率单光子探测器，实现了实验室标准光纤百公里级、两百公里级、三百公里级以及超低损光纤四百公里级的安全成码，相较于之前的原始 MDI 实验，成码率有明显提升，并且在三百公里和四百公里距离上较之前实验成码率提升了 3 个数量级。实验结果表明，模式匹配量子密钥分发在不需激光器锁频锁相的条件下可以实现远距离安全成码且在城域距离有较高成码率，极大地降低了协议实现难度，对未来量子通信网络构建具有重要意义。

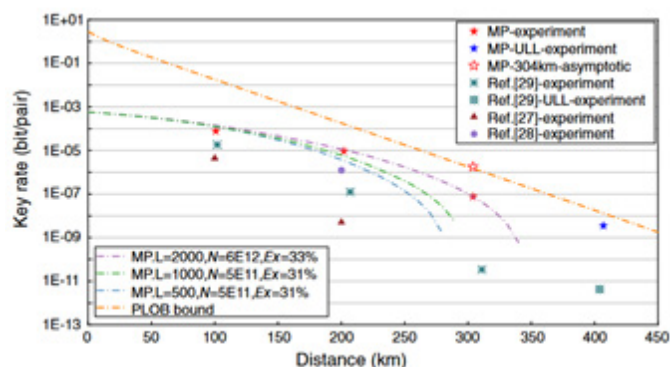


图 2. 模式匹配协议的成码率比较图

该研究成果论文：Hao-Tao Zhu, Yizhi Huang, Hui Liu, Pei Zeng, Mi Zou, Yunqi Dai, Shibiao Tang, Hao Li, Lixing You, Zhen Wang, Yu-Ao Chen, Xiongfeng Ma, Teng-Yun Chen, and Jian-Wei Pan, “Experimental Mode-Pairing Measurement-Device-Independent Quantum Key Distribution without Global Phase Locking”, Physical Review Letters, 2023, 130(3): 030801.

三、超导量子计算

主要完成人：孙麓岩研究组

突破量子纠错盈亏平衡点

孙麓岩超导量子计算课题组一直致力于量子纠错的实验研究。该研究组与南方科技大学俞大鹏院士和徐源研究团队、福州大学郑仕标教授团队等在基于超导量子线路系统的量子纠错领域取得了突破性实验进展：通过实时重复的量子纠错技术延长了量子信息的存储时间，在国际上首次超越盈亏平衡点，展示了量子纠错的优势。这一里程碑式的突破代表了迈向实用化可扩展通用量子计算的关键一步。

虽然基于超导量子线路系统的量子信息处理领域研究近些年发展迅猛，但由于量子计算机体系的错误率远高于经典数字计算机，想要构建具有实用价值的通用量子计算机，量子纠错依然不可或缺，因为量子纠错可以有效地保护量子信息避免受到环境中噪声的干扰。传统的量子纠错方案编码一个逻辑量子比特需要多个冗余的物理比特，不但需要巨大的硬件资源的开销，发生错误的通道数也会随着比特数的增加而显著增多，可能会呈现“越纠越错”的尴尬局面。虽然这种量子纠错方案已经有多个演示性的实验研究工作，可仍然无法解决量子纠错过程中“越纠越错”的问题，未真正实现超越盈亏平衡点。也就是说，量子纠错之后的效果还远没有达到该系统中不纠错情况下的最好值，无法真正产生正的量子纠错增益。这也成为当前量子纠错技术无法实用化可扩展发展的核心瓶颈。

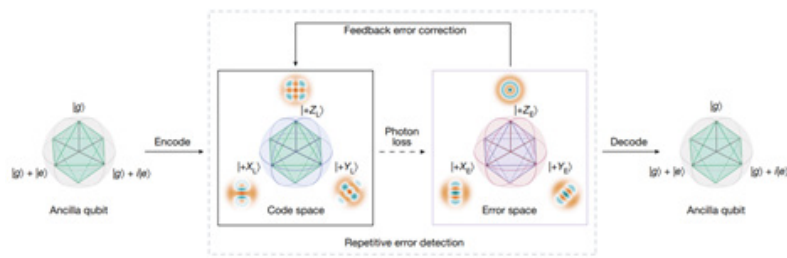


图 1. 量子纠错过程示意图

为攻克上述难题，联合研究团队利用微波简谐振子或玻色模式系统中的无穷维希尔伯特空间，实现量子信息的冗余编码与量子纠错。在超导量子线路系统中，基于玻色编码的量子纠错方案具有错误类型简单、错误探测方便、相干性能好、硬件更高效、反馈控制易实现等优点。该研究工作中，研究团队通过开发高相干性能的量子系统，设计和实现低错误率的错误症状探测方法，以及改进和优化量子纠错技术等实验手段，最终在玻色模式中实现了基于离散变量的二项式编码的逻辑量子比特，并通过实时重复的量子纠错过程，延长了量子信息的存

储时间，相关结果首次超过该系统中不纠错情况下的最好值，也就是突破了盈亏平衡点。这也是国际上首次通过主动的重复错误探测和纠错过程实现延长量子信息的存储时间超越盈亏平衡点，具有里程碑式的重要意义。

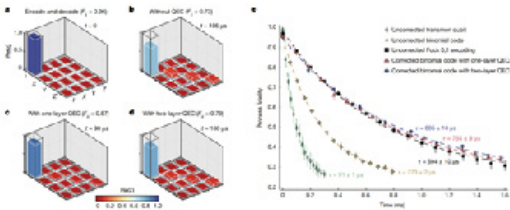


图 2 量子纠错操作的实验表征结果

该研究成果论文：Zhongchu Ni, Sai Li, Xiaowei Deng, Yanyan Cai, Libo Zhang, Weiting Wang, Zhen-Biao Yang, Haifeng Yu, Fei Yan, Song Liu, Chang-Ling Zou, Luyan Sun, Shi-Biao Zheng, Yuan Xu & Dapeng Yu. “Beating the break-even point with a discrete-variable-encoded logical qubit” . Nature volume 616, pages56 – 60 (2023).

超导量子线路中演示新型量子端到端机器学习模型

孙麓岩研究组与自动化系吴热冰课题组合作，在超导量子线路上首次演示了量子端到端机器学习模型的训练。

随着量子计算技术的不断革新，结合量子技术与机器学习优势的量子机器学习正在得到越来越广泛的关注。即使不依赖可容错量子计算机，在近期中等规模有噪声量子计算机上的量子机器学习也被认为十分具有前景。因为随比特数量指数扩展的特征空间可以带来模型表达能力显著增长，从而处理更复杂的任务。

量子机器学习的实现关键是构建可以在量子计算机上有效训练的参数化拟设模型。目前研究中广泛采用的模型是基于参数化量子门的量神经网络模型。相关模型已经在分类、聚类、生成式学习等任务上实现了实验演示。基于参数化量子门的量子神经网络在实际应用中存在一些挑战。例如训练效果高度依赖量子门序列的结构设计—未优化的结构通常不能有效利用有限的量子资源，以及精确标量子门参数需要消耗较多实验资源。

孙麓岩研究组与自动化系吴热冰研究组共同提出了一种新的基于控制参数的量子机器学习模型，有效的克服了上述挑战。这种模型不需要特别的结构优化，并且不需要量子门标定，因而实现了硬件高效性。另外，模型内天然的控制-量子态映射提供了提升模型表达能力必须的非线性，是该模型成功的关键因素之一。

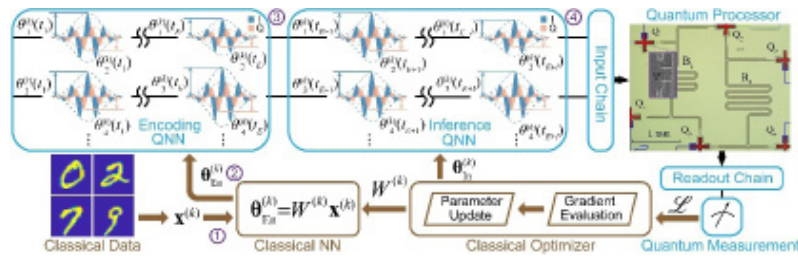


图 1: 基于控制参数的量子端到端机器学习训练过程示意图。

实验系统涉及到在一个六比特超导芯片上的三个互相连通的超导量子比特，其相互作用由一个耦合谐振腔提供。每个超导量子比特具有独立的微波驱动控制，可以实现任意的单比特旋转操作。实验演示的训练任务是手写数字图片分类问题，分类的结果由输出比特出现概率最大的计算基矢表示。实验训练过程如图 1 所示。手写数字图片输入到经典计算机后，由一个经典神经网络变换为编码层控制参数，实验仪器根据编码层与推断层控制参数直接向量子芯片发送相干的控制脉冲。脉冲的时间序列构成了多层结构的量子神经网络。通过量子测量输出比特的状态，并且在经典计算机上生成损失函数及其梯度，最终分别更新经典神经网络与量子神经网络参数，实现量子-经典混合的量子机器学习训练。

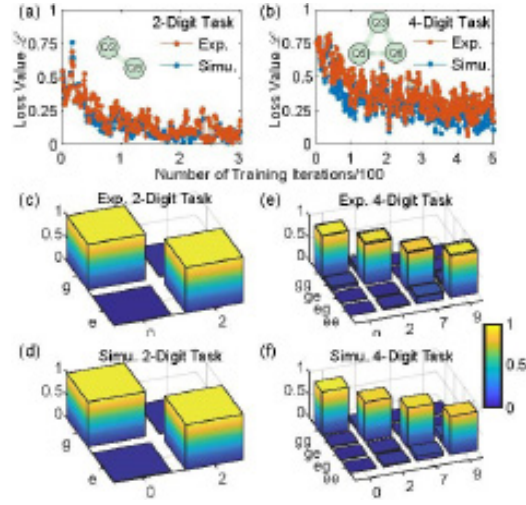


图 2: 量子端到端机器学习模型在二分类与四分类手写数字识别任务的训练结果。

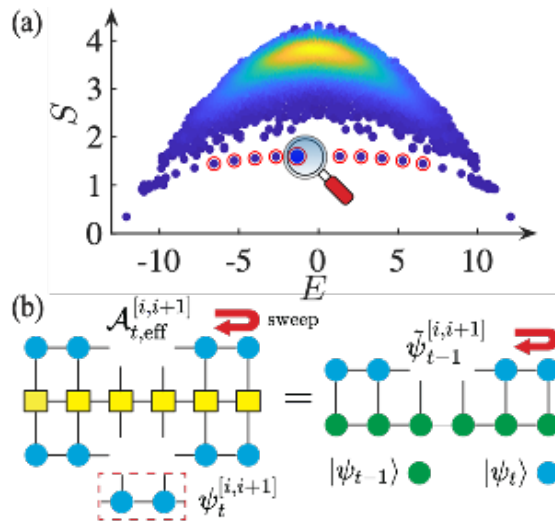
该研究成果论文: Xiaoxuan Pan, Xi Cao, Weiting Wang, Ziyue Hua, Weizhou Cai, Xuegang Li, Haiyan Wang, Jiaqi Hu, Yipu Song, Dong-Ling Deng, Chang-Ling Zou, Re-Bing Wu & Luyan Sun, “Experimental quantum end-to-end learning on a superconducting processor”, npj Quantum Information, 9, 18 (2023).

四、量子多体物理

主要完成人：邓东灵研究组

提出了一种基于矩阵直积态的算法

量子多体伤痕是量子多体系统中出现的一种非热化的特殊本征态。这些异常的伤痕本征态违反了本征态热化假说，它们淹没在热化本征态的海洋之中，仅占据整个希尔伯特空间的微小部分。在不知道这些特殊本征态的精确表达式的情况下，从指数多的热化本征态中将它们区分出来具有很高的计算复杂性。



在一维系统中，伤痕本征态的纠缠熵通常最多以对数方式随系统尺寸增长，这表明它们可以使用矩阵乘积态表示，从而超越了精确对角化方法的系统尺寸限制。该文中，邓东灵研究组提出了一种基于矩阵直积态的算法，称为“DMRG-S”，可以以高精度提取量子多体伤痕本征态。为了展示其优势，该研究组计算了“PXP”模型及其变形模型中的一系列伤痕本征态（系统尺寸高达 80 个格点）。通过详细的有限尺寸标度研究，该研究组发现在 PXP 模型中，Néel 态的相干复苏在热力学极限下消失，而在变形的“PXP”模型中保持稳定。此外，先前的分析研究表明，在一些模型中，高激发的伤痕本征态具有精确的矩阵直积态表示，但这些伤痕态的构造方法是特定于模型的，缺乏普适性。相比之下，该研究组的方法提供了一种系统的方式，可以在一般哈密顿量中寻找量子多体伤痕的精确矩阵乘积态表示，而无需先验知识。特别地，该研究组在运动受限的自旋模型和时钟模型中发现了几个具有精确矩阵直积态表示的新的量子多体伤痕本征态，并且找到了这些伤痕本征态的后验分析方法。该研究为未来寻找新的多体伤痕态以及研究它们的稳定性提供了强有力的工具，为量子多体伤痕的研究开创了一套新的方法论，对未来相关方面的理论和实验研究提供了指导。

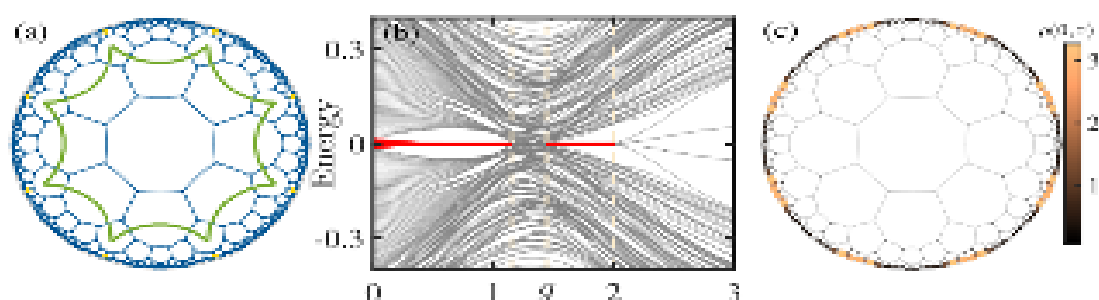
该成果研究论文：Shun-Yao Zhang, Dong Yuan, Thomas Iadecola, Shenglong Xu, and Dong-Ling Deng, "Extracting quantum many-body scarred eigenstates with matrix product states", arXiv preprint arXiv:2211.05140, 目前该论文已被 PRL 接受。

五、凝聚态物理学

主要完成人：徐勇研究组

高阶拓扑双曲晶格

最近，双曲晶格在量子电动力学电路 [Nature 571, 45 (2019)] 和电学电路 [Nat. Commun. 13, 2937 (2022)] 的实验中实现，它的新颖性质引起了人们极大的兴趣。对于常规二维晶体，它只允许存在二重、三重、四重或六重旋转对称性，然而，双曲晶格可以实现任意重旋转对称性，由此引发一个问题：这种新的旋转对称性是否会带来新的拓扑物态？



(a) {8,3} 双曲晶格的示意图；(b) 体系随参数变化的能谱图；(c) 非平庸角态的态密度图

徐勇研究组首次在理论上预言双曲晶格中存在与八重、十二重或十六重旋转对称性相关的高阶拓扑相。以 {8,3} 双曲晶格为例，研究组基于这一特殊的晶格构造了新的紧束缚模型，其哈密顿量满足一种结合了八重旋转对称性和时间反演对称性的新对称性。经过计算，研究组发现了有能隙的高阶拓扑相和无能隙相，并通过区域电荷和边界傅里叶分析等方式进一步表征了高阶拓扑相的非平庸角态。为了证实这种新的高阶拓扑相在双曲晶格中的普适性，研究组基于 {12,3} 和 {16,3} 双曲晶格构造紧束缚模型，并发现了高阶拓扑相的非平庸角态。

此项工作发现了不同于常规晶格的新的拓扑相，体现出双曲晶格在拓展拓扑相研究方面的意义，为后续探索双曲晶格中特有的拓扑相开辟了道路。

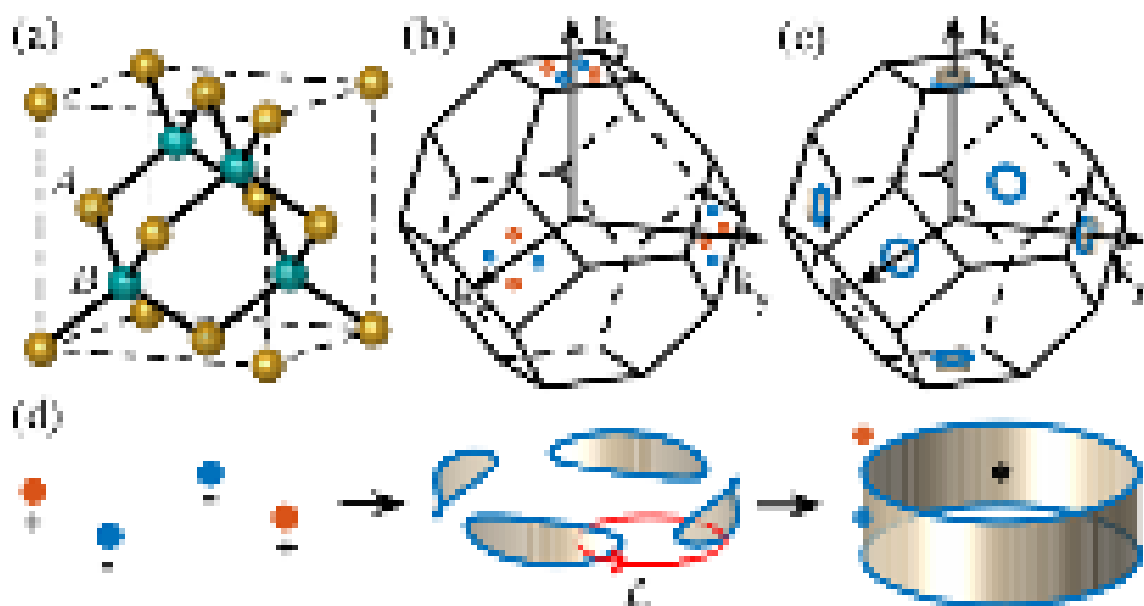
该成果研究论文：Yu-Liang Tao and Yong Xu, “Higher-order Topological Hyperbolic Lattices”, Phys. Rev. B 107, 184201 (2023).

三维系统中的奇异重费米半金属

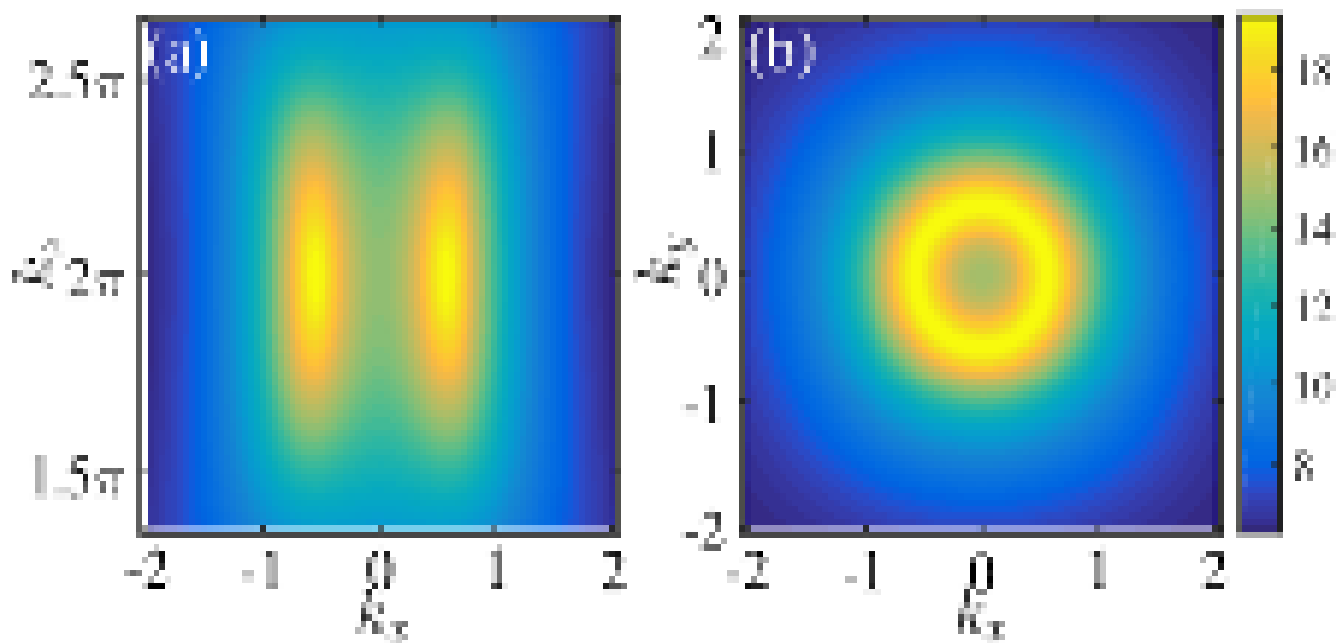
强关联系统相较于无相互作用的电子体系有着很多新奇的现象,例如,在强关联体系中允许存在有边界终点的体费米弧,并且这一现象已经在二维氧化铜高温超导体的赝隙相中被观测到 [B. Keimer, et al., Nature, 518, 179 (2015)]. 在这样的背景下,人们在理论上证实了在二维重费米子体系中存在有边界的体费米弧,其成因是有效单体哈密顿量存在奇异点。但是这种二维的分析并不能直接关联到真实三维材料。对于三维重费米子拓扑材料,已有的发现包括拓扑绝缘体和拓扑半金属,而这些研究并未注意到由于准粒子寿命项的不同而导致的奇异点的出现。

徐勇研究组在理论上提出了一种新的三维重费米子系统——奇异重费米半金属。由于空间反演对称性被破坏,系统中的两种准粒子的寿命不再相同,这使得有效单体哈密顿量出现了奇异环,而奇异环导致了有边界的体费米面,这一特征可以直接通过角度分辨光电子能谱观测到。此系统的提出是基于三维微观周期安德森模型,比如 CeRu_4Sn_6 或 $\text{Ce}_3\text{Bi}_4\text{Pd}_3$ 这种包含强关联 f 电子的系统。研究组通过微扰论和动力学平均场理论两种方法计算了系统的有效单体哈密顿量和系统谱函数,进一步证实了体费米薄带的存在。近来,非中心对称的重费米半金属 $\text{Ce}_3\text{Bi}_4\text{Pd}_3$ 已经被实验证实,自然奇异重费米半金属也有可能在实际材料中被观测到。

此项工作深入探究了重费米子体系的半金属相,并为三维奇异重费米半金属相的研究开辟了新的道路。



图一: (a) 闪锌矿晶格结构示意图; (b) 面心立方第一布里渊区, 其中包括六对外尔点; (c) 有限温下演化出的奇异环和体费米薄带; (d) 体费米薄带演化示意图。



图二：不同动量截面的谱函数，图中结果是由动力学平均场理论计算所得。

该成果研究论文：Yu-Liang Tao, Tao Qin, and Yong Xu, “Exceptional Heavy-Fermion Semimetals in Three Dimensions” ,

Phys. Rev. B 107, 035140 (2023).



Editor:
Kailin Li, Yueliang Jiang
Reviewer:
Luming Duan, Jian Li, Xiamin Lv