



清华大学 交叉信息研究院

Institute for Interdisciplinary Information Sciences, Tsinghua University

学术科研简报

IIS Academic Newsletter

2025.07-12

目 录

人工智能领域 (04-25)

人工智能大模型	04
具身智能和机器人	08
计算机图形学 / 视觉	17
人工智能安全	20
人工智能理论	23

计算机科学领域 (27-40)

计算机系统结构	27
数据库管理系统	31
区块链	33
计算机安全	34
密码学	35

量子信息领域 (42-57)

离子阱量子计算、量子网络	42
离子阱量子中继	43
中性原子量子网络	45
金刚石量子模拟与传感	46
量子信息科学	48
量子模拟	52
超导 - 光学 / 声学混合系统量子计算	54
超导量子计算	56
凝聚态物理学	57

人工智能



一、人工智能大模型

主要完成人：房智轩研究组、吴翼研究组

用户侧大模型推理服务一致性监测方法及其实证研究

随着开源大模型模型性能不断逼近闭源模型，第三方推理服务因成本优势被广泛采用，但服务过程缺乏透明性，用户难以确认服务商是否真实使用了所声明的模型规模与精度。一些服务商可能通过缩小参数规模或采用低精度量化模型以降低成本，这不仅损害用户体验，也可能引发不公平竞争。

针对这一问题，房智轩研究组从用户视角出发，探索一种无需服务商配合、且不引入高额计算或经济成本的监测机制。该机制的核心思路是：基于大模型自回归生成的特性，可以利用“生成慢、验证快”的差异，将一致性验证的计算转移到用户本地设备完成。通过对服务商返回的完整生成序列执行一次并行前向计算，即可获得所有 token 的概率分布信息，将 token 级指标在连续子序列层面进行聚合，可有效放大模型退化所带来的累积差异，从而用于可信分析。结合模型参数卸载技术，即便模型规模远超 GPU 显存，也能在消费级设备上完成验证，使该方案具备现实可行性。

该研究组在 Llama 3.1, Qwen 2.5 等主流大模型系列上系统验证了该方案在不同硬件平台、不同模型规模与精度设置下的有效性与鲁棒性。结果显示，正常服务与退化服务在一致性分布上呈现出清晰可分的特征。同时，原型系统在消费级 GPU 上的运行测试表明，该方案在时间与资源开销上均可接受。总体而言，本工作为开源大模型推理服务提供了一种低成本、用户可控的信任保障新路径。

该成果研究论文：Qijun Miao and Zhixuan Fang, “User-side Model Consistency Monitoring for Open Source Large Language Models Inference Services”, ACL 2025.

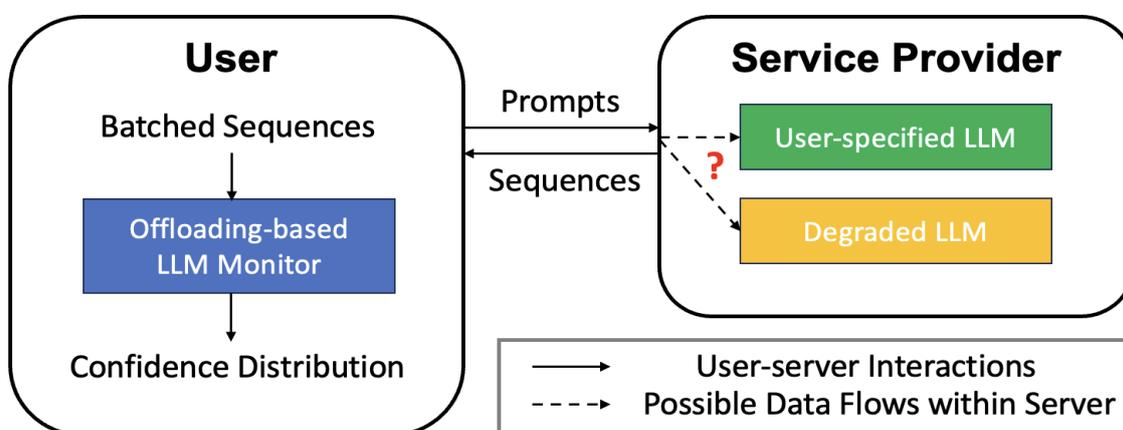


图 1 长时程移动操作任务示意图

Z-Score 引导的早期停止优化提升长程推理效率 (ZGES)

大型语言模型 (LLM) 在复杂推理任务中常依赖思维链 (CoT) 进行多步推演。随着长程思维链模型 (如 DeepSeek R1、OpenAI o1) 的发展, 推理步数显著增加, 如何高效引导与终止推理过程成为关键挑战。传统的基于过程奖励模型 (PRM) 的束搜索方法在短思维链模型设置中表现优异, 但在长程推理中常出现性能退化, 导致计算资源浪费与效果下降。

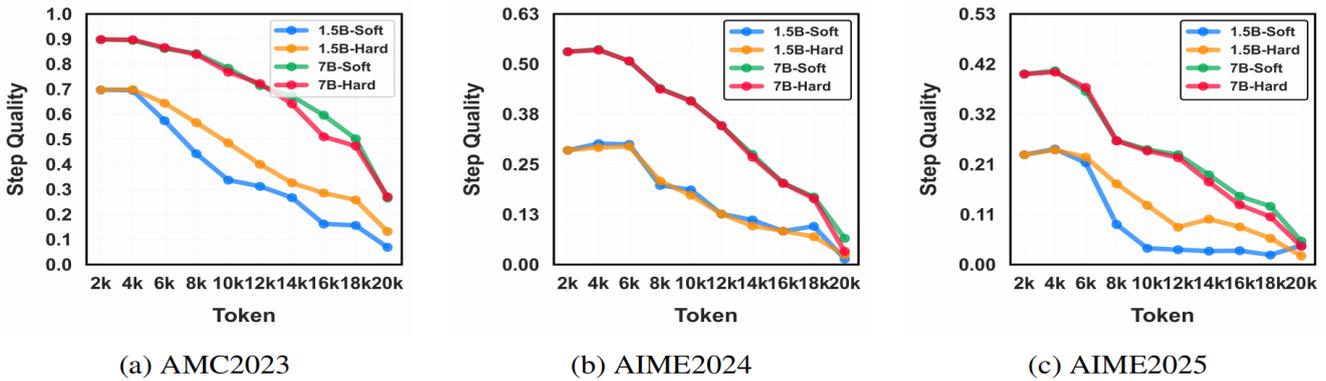


图 1 长程推理中步质量随步数增加呈单峰或单调下降趋势 (步质量退化)

吴翼研究组针对该问题展开研究, 首次揭示了在长程 CoT 推理中, 束搜索的步质量呈现“单峰”或“单调下降”的退化现象。理论分析表明, 该现象源于 PRM 在搜索过程中重排序能力的显著衰减。基于此, 研究组提出 Z-Score 引导的早期停止方法 (ZGES), 通过动态监测 PRM 奖励的局部 Z 值变化, 在质量峰值附近提前终止搜索, 避免后续低效推理。

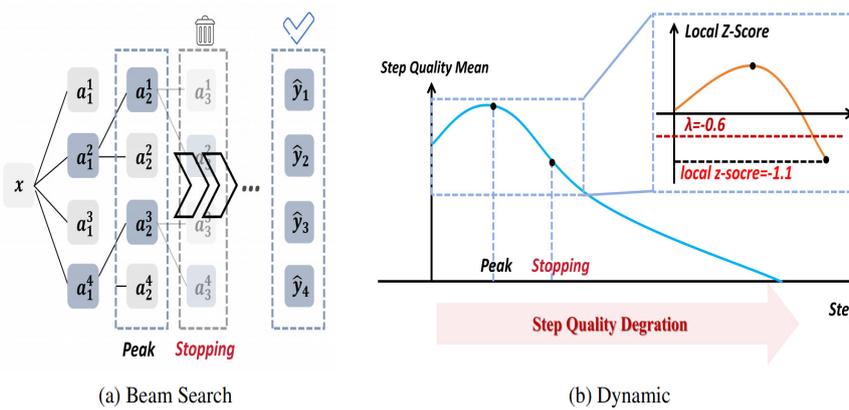


图 2 ZGES 方法在检测到局部 Z 值低于阈值 λ 时提前终止束搜索, 从峰值步继续生成

该方法在多个数学推理基准 (AMC2023、AIME2024、AIME2025) 和不同规模的长程 CoT 模型 (DeepSeek-R1-Distill-Qwen-1.5B/7B) 上进行了验证。实验表明, ZGES 在保持或提升推理准确率的同时, 显著减少了 PRM 调用次数 (至少降低 50% 以上) 与总 token 生成量, 实现了效率与性能的双重提升。ZGES 为长程复杂推理任务提供了一种轻量、自适应的搜索终止策略, 推动了高效推理搜索算法的发展。

该成果研究论文: Zhang M, Gao J, Xu S, et al, “Reasoning Is Not a Race: When Stopping Early Beats Going Deeper”, NeurIPS 2025.

AREAL: 面向语言推理的大规模异步强化学习系统

强化学习已成为训练大型语言模型，尤其是提升其推理能力的重要范式。有效训练需要大规模并行生成轨迹，这对训练系统的效率提出了极高要求。现有大规模 RL 系统多为同步式，即在批量设置中交替进行生成和训练，虽稳定但存在严重的系统级低效问题：生成阶段必须等待批次中最长的输出完成，导致 GPU 利用率低下，且难以扩展。

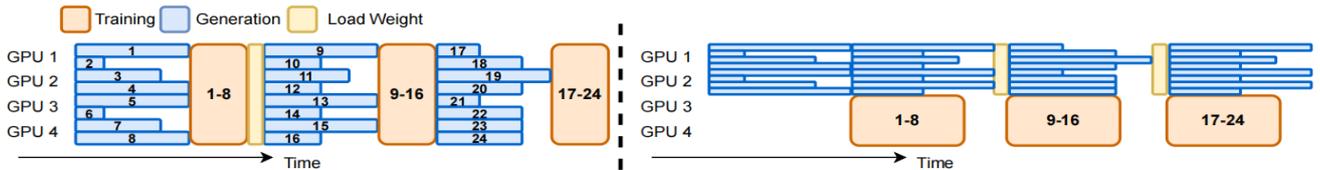


图 1 同步（左）与一步重叠（右）RL 系统的执行时间线，显示推理设备利用率不足

吴翼研究组与蚂蚁集团合作提出 AREAL，一个面向大语言模型推理训练的完全异步 RL 系统，从根本上解决了上述效率瓶颈。该系统将生成与训练完全解耦：可中断的 Rollout Worker 持续生成新输出而无需等待，训练器 Worker 则随时利用收集到的批次数据更新模型。AREAL 采用陈旧性感知训练机制，通过控制数据延迟来稳定训练，并引入解耦 PPO 目标，将采样行为策略与用于约束更新的近端策略分离，从而能有效利用异步生成的、来自不同模型版本的数据。

在数学推理与代码生成任务上的实验表明，AREAL 相比先进的同步系统实现了高达 2.77 倍的训练加速，并在保持甚至提升最终模型性能的前提下，显著提高了 GPU 利用率与系统可扩展性。该系统可线性扩展到 512 块 GPU，并为未来超大规模 RL 训练提供了高效的算法与系统协同设计范例。该工作为高效、可扩展的大语言模型强化学习训练提供了创新的系统解决方案。

该成果研究论文：Fu W, Gao J, Shen X, et al. "AReAL: A Large-Scale Asynchronous Reinforcement Learning System for Language Reasoning", NeurIPS 2025.

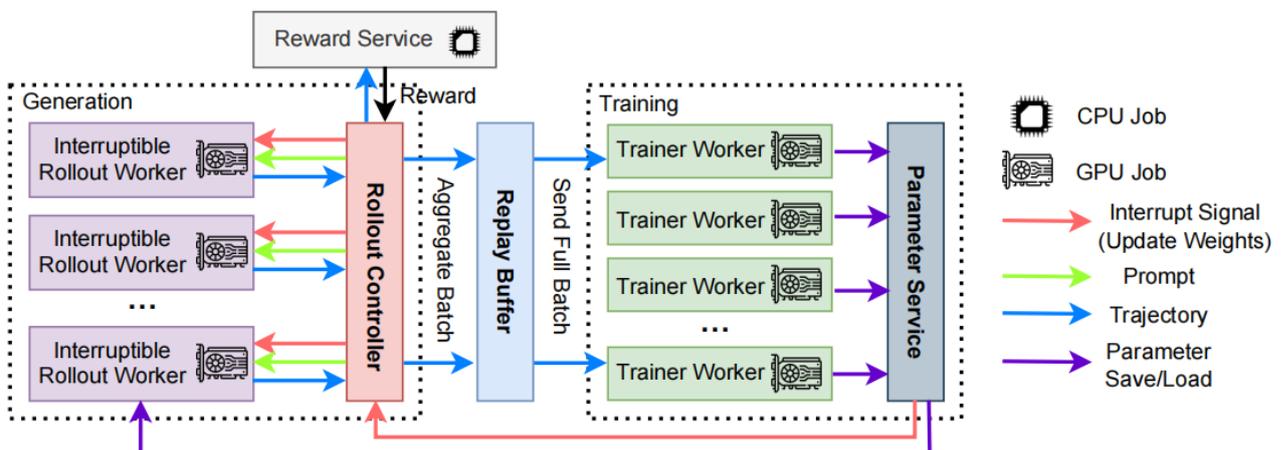


图 2 AREAL 系统架构

REO-RL: 动态优化大模型推理过程

随着大语言模型 (Large Language Model) 能力的不断提升, 基于“思维链” (Chain-of-Thought) 的深度推理已成为增强大语言模型复杂问题求解能力的关键技术。当前模型在数学、代码等任务上已展现出接近人类专家的水平, 但在处理简单问题时, 常出现过度思考 (Overthinking) 现象, 导致推理步骤冗余、效率低下, 进而造成计算资源浪费与响应延迟。尽管已有研究尝试通过训练提升模型推理效率, 但现有方法与最优效率之间仍存在多大差距, 这一问题尚未得到系统回答。

吴翼研究组在研究工作“*How Far Are We from Optimal Reasoning Efficiency?*”中, 首次对推理效率的优化空间进行了系统性评估。研究通过大量实验, 提出了对“最优推理效率”的经验性估计, 并设计了一个新的评估指标, 能够更均衡地衡量推理长度与准确率之间的取舍。此外, 该研究组提出一种基于强化学习的算法 REO-RL, 该算法能够在不同推理预算下, 动态优化模型的推理过程, 显著提升效率。

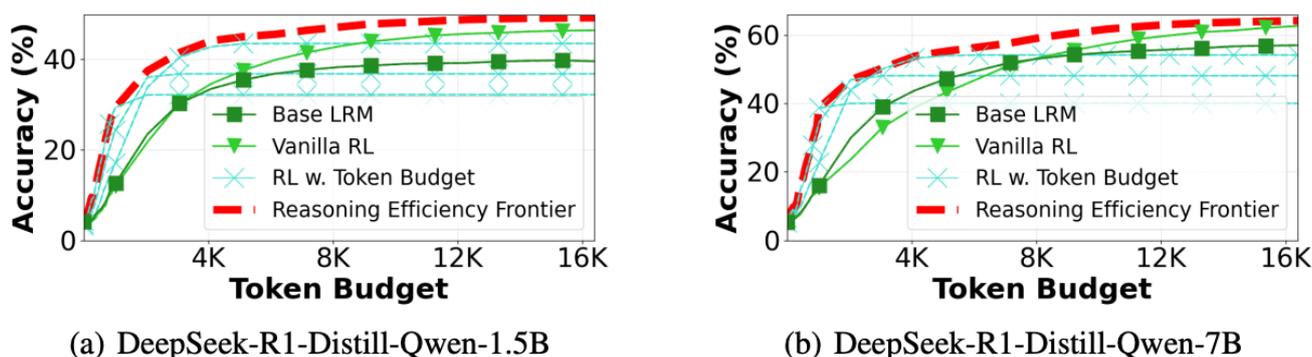


图 1 对于最优推理效率的经验性估计

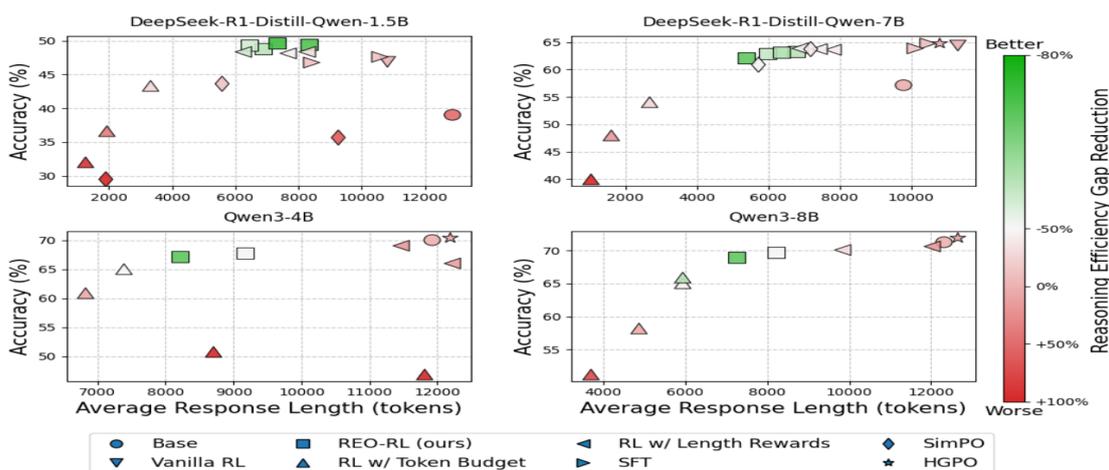


图 2 REO-RL 算法在多个任务上显著提升推理效率, 同时新提出的效率指标在“长度”与“准确率”之间取得更好平衡

实验表明, REO-RL 在多个数学推理基准上取得了当前最优的推理效率表现。在 1.5B、7B 等不同规模的模型上, 该算法能够在保持高准确率的同时, 大幅减少推理步骤。所提出的效率指标也为未来相关研究提供了可量化的评估工具。

该成果研究论文: Gao, Jiaxuan, Shu Yan, Qixin Tan, Lu Yang, Shusheng Xu, Wei Fu, Zhiyu Mei, Kaifeng Lyu, and Yi Wu,

"How Far Are We from Optimal Reasoning Efficiency?", NeurIPS 2025.

二、具身智能和机器人

主要完成人：陈建宇研究组、高阳研究组、吴翼研究组、赵行研究组、马恺声研究组

统一离散和连续表征学习的机器人策略

在开放式环境中构建能够处理多样化任务的通用型机器人策略是机器人领域的一项核心挑战。为了利用大规模预训练所带来的知识，已有工作通常基于视觉 - 语言理解模型（VLM）或生成模型来构建通用型机器人策略。然而，对于具身机器人而言，源自视觉 - 语言预训练的语义理解能力以及源自视觉生成预训练的视觉动态建模能力同样至关重要。

近期，将生成与理解统一起来的模型通过大规模预训练在理解和生成两方面均展现出了强大的能力。对此，陈建宇研究组认为机器人策略学习同样可以从理解、规划以及连续未来表示学习的协同优势中受益。基于这一认识，研究组提出了 UniCoD，一种遵循“理解 - 生成 - 执行”范式的 VLA 框架，将离散的任务理解与连续的未来状态预测相结合。UniCoD 采用具有模态专家的混合架构以处理异构模态，并通过两阶段训练策略，在保持视觉 - 语言通用能力的同时，将连续特征预测引入动作学习。第一阶段利用来自机器人和人类演示的大规模具身问答数据，学习语言理解和世界建模表示；第二阶段引入带动作标注的机器人数据，通过联合预测视觉未来和动作，使策略能够利用语义对齐且包含动态信息的表示，从而提升对新物体和新场景的泛化能力。

通过在仿真器以及真机的大量实验结果表明，UniCoD 方法在仿真环境和真实世界的分布外任务中，相比基线方法分别取得了约 9% 和 12% 的一致性性能提升。

该成果研究论文: Jianke Zhang and Yucheng Hu, “UniCoD: Enhancing Robot Policy via Unified Continuous and Discrete Representation Learning”, arXiv preprint arXiv:2510.10642, 2025.

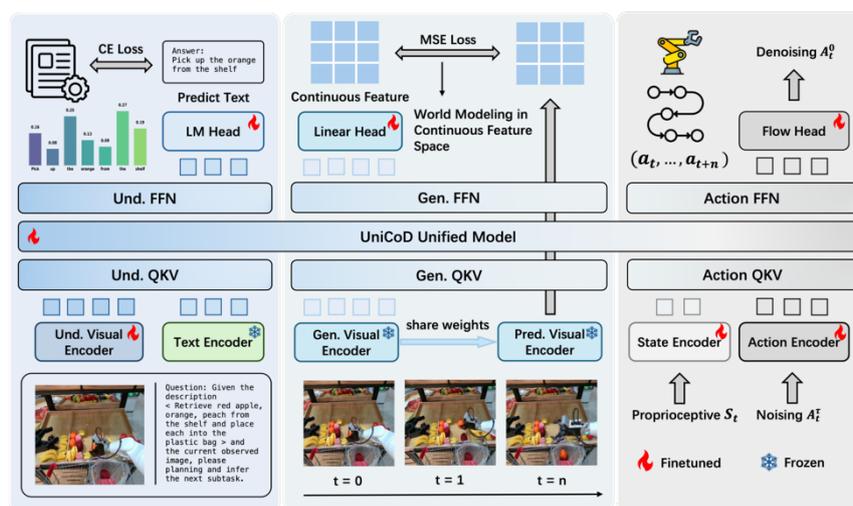


图 1 UniCoD 框架

villa-X: 增强视觉 - 语言 - 动作模型中的潜在动作建模

视觉 - 语言 - 动作 (VLA) 模型已成为机器人操控领域的主流范式, 并取得了显著进展。然而, 构建更加通用且具备高泛化能力的机器人控制策略, 通常依赖于海量且多样化的训练数据。面对当前机器人数据增长速度远滞后于训练需求的现状, 如何高效利用其他形式的数据成为亟待解决的难题。

陈建宇研究组深入研究了 VLA 模型中的隐动作建模问题, 使得模型可以高效地从海量的人类视频数据中学习, 并在不同本体之间迁移知识。陈建宇研究组提出了一种名为“villa-X”的视觉 - 语言 - 隐动作 - 动作训练框架。该框架设计了新的隐动作模型 (LAM), 能够从无标签的人类视频数据中自动生成隐动作 (latent actions) 供策略模型训练。LAM 在压缩视觉变化的基础上, 引入了物理落地 (grounding) 机制, 确保学习到的潜在动作可以捕捉对控制至关重要的物理动力学细节, 从而在视觉感知与底层控制之间建立了更有效的连接。基于学习到的隐动作, 策略 (policy) 通过联合去噪的学习形式, 可以从人类视频和各种机械臂数据中高效吸收知识, 提升解决问题的能力。

该研究组证明, villa-X 能够在零样本条件下, 生成合理的潜在动作规划, 并具备处理未见过的机器人形态及开放词汇指令的能力。在 SIMPLER 等模拟基准测试中, villa-X 展现了超越目前最先进模型 (如 pi0、GR00T 等) 的优异性能。此外, 该研究组还在涉及夹爪和 xHand 灵巧手操控的真实世界机器人场景中验证了该算法的有效性, 结果表明 villa-X 能够快速适应新的机器人形态, 在广泛的任务中展示了卓越的泛化能力。

该成果研究论文: Chen, Xiaoyu, et al. "Villa-x: enhancing latent action modeling in vision-language-action models", arXiv preprint arXiv:2507.23682 (2025).

VLM4VLA: 重新审视视觉 - 语言 - 动作模型中的视觉 - 语言模型

视觉 - 语言 - 动作 (VLA) 模型通过将预训练的大型视觉 - 语言模型 (VLM) 集成到策略骨干中, 因其展现出的强大泛化能力而备受关注。然而, 目前的研究大多集中在网络架构设计或训练范式上, 却很少系统性地研究一个基础问题: VLM 的选择及其能力如何转化为下游 VLA 策略的性能? 该研究组研究人员重新审视了这一关键问题。为了在不引入额外变量的情况下公平评估 VLM 的能力, 陈建宇研究组提出了 VLM4VLA——这是一个极简的适配管道, 仅引入不到 1% 的新可学习参数, 即可将通用 VLM 转换为 VLA 策略。该框架通过统一的接口和训练设置, 使得不同 VLM 骨干网络 (如 QwenVL、Paligemma、Kosmos 等) 能够在机器人控制任务中进行公平、高效的比较。研究组在三个基准测试 (Calvin, SimplerEnv, Libero) 上进行了大规模实证研究, 得出了几个反直觉但重要的结论:

1. VLM 初始化有益但预测性差: 虽然 VLM 初始化比从头训练效果更好, 但 VLM 在通用任务上的表现并不能很好地预测其在下游控制任务中的性能。
2. 特定具身能力微调效果有限: 在辅助具身任务 (如具身问答、视觉指向、深度估计) 上微调 VLM, 并不一定能保证下游控制性能的提升。
3. 视觉模块是瓶颈: 模态层面的消融实验表明, VLM 中的视觉模块 (而非语言模块) 是主要的性能瓶颈。即使在下游微调期间冻结视觉编码器, 只要预先注入控制相关的监督信号, 也能带来显著的性能提升。这揭示了当前 VLM 预训练目标与具身动作规划需求之间存在显著的领域差距。

该成果研究论文: Jianke Zhang, Xiaoyu Chen et al, “VLM4VLA: Revisiting Vision-Language-Models in Vision-Language-Action Models”, <https://arxiv.org/abs/2601.03309>, October 2025.

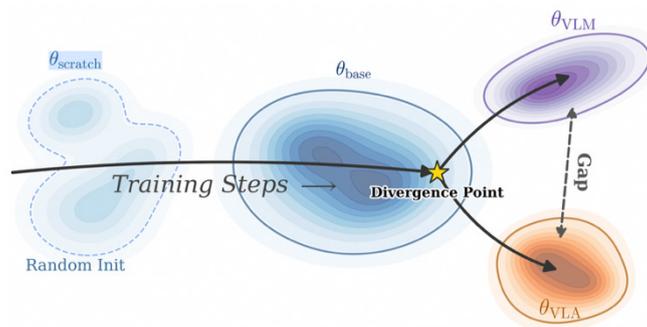


图 1 现有的 VLM 和 VLA 存在的领域差距

基于动觉示教的数据采集与触觉增强的灵巧手控制学习

在具身智能与灵巧操作研究中，高质量机器人数据的获取是制约算法性能与可扩展性的关键瓶颈，尤其是在涉及精细接触与力控制的任务中。传统的遥操作或视频重定向方法往往受到人机运动学不匹配、缺乏真实触觉反馈以及数据采集效率低等问题的限制，难以支撑复杂灵巧操作策略的学习。

针对上述挑战，高阳研究组提出了 KineDex，一种结合动觉示教（Kinesthetic Teaching）与触觉感知的灵巧操作学习框架。KineDex 通过“手把手”动觉示教方式，使操作者能够直接操控机器人灵巧手完成任务，并在示教过程中获得真实、精确的力反馈，从而采集高保真、触觉增强的示教数据。该范式有效避免了遥操作中的运动重定向误差，同时显著提升了数据采集的直观性与效率。

在策略学习阶段，KineDex 针对动觉示教过程中不可避免的人手遮挡问题，引入了基于视觉修复的数据预处理方法，对示教视频进行自动去遮挡处理，从而缓解训练与部署阶段的分布偏移。在此基础上，研究组构建了融合视觉、触觉与本体感觉信息的视觉—运动策略，并在推理阶段引入力控制机制，使策略不仅预测关关节目标位置，还能显式建模并跟踪指尖接触力，实现稳定、精细的接触操作。

实验结果表明，KineDex 在九项高难度接触密集型灵巧操作任务中均取得了显著性能提升，平均成功率达到 74.4%，在无力控制或无触觉输入的消融设置下性能显著下降，验证了触觉感知与力控制在灵巧操作中的关键作用。与遥操作方法相比，KineDex 在示教成功率和数据采集效率上均展现出明显优势，用户研究也进一步证明了其易用性与实用价值。

该工作为灵巧机器人操作中“如何高效获取高质量示教数据”以及“如何将触觉与控制深度融合”提供了一种系统性解决方案，为面向真实世界的复杂操作技能学习奠定了重要基础。

该成果研究论文：Di Zhang, Chengbo Yuan, Chuan Wen, Hai Zhang, Junqiao Zhao, Yang Gao, “KineDex: Learning Tactile-Informed Visuomotor Policies via Kinesthetic Teaching for Dexterous Manipulation”, CORL 2025.



图 1 KineDex 框架示意图

HuB: 人形机器人极致平衡学习

人形机器人在执行复杂动作时需要具备高度精细的平衡控制能力，尤其是在单腿支撑、大幅度肢体运动以及强外界扰动等极致条件下保持稳定。然而，现有的人形机器人强化学习方法多依赖对参考动作的精确模仿，在面对参考数据误差、人机形态差异以及仿真与现实不一致等问题时，往往难以实现稳定可靠的平衡控制。这些问题在需要高平衡精度的任务中尤为突出，成为限制人形机器人应用的关键瓶颈。

HuB (Humanoid Balance) 工作围绕人形机器人极致平衡控制的学习问题展开研究，旨在使人形机器人在缺乏完美参考动作的条件下，仍能够完成对平衡要求极高的复杂姿态控制任务。该研究将极致平衡问题视为一个高稳定性要求的控制学习问题，并针对传统模仿学习与强化学习方法在该类任务中的失效原因进行了系统分析。

针对上述挑战，HuB 提出了一套统一的人形机器人极致平衡学习框架，从参考动作、策略学习和仿真到现实迁移三个层面进行设计。首先，通过参考动作优化机制对动作进行修正，降低不精确参考对学习稳定性的负面影响；其次，在策略训练阶段引入平衡感知的学习机制，使策略在学习过程中更加关注稳定性而非简单的动作追踪；最后，通过仿真环境中的扰动注入与随机化训练，显著增强学习策略在真实机器人上的鲁棒性与泛化能力。

在实验验证中，HuB 框架被部署于宇树 G1 人形机器人，并在多项极致平衡任务中进行了系统评估，包括“燕式平衡”姿态以及单腿高踢等对稳定性要求极高的动作。实验结果表明，HuB 学习得到的策略不仅能够稳定完成这些任务，还能够在受到明显外部扰动时保持平衡，而现有基线方法在相同条件下通常难以成功执行。

该成果研究论文：Tong Zhang, Boyuan Zheng, Ruiqian Nai, Yingdong Hu, Yen-Jen Wang, Geng Chen, Fanqi Lin, Jiongye Li, Chuye Hong, Koushil Sreenath, Yang Gao, “HuB: Learning Extreme Humanoid Balance”, CORL 2025.



图 1 人形机器人极致平衡任务图

RL4VLA: 面向 VLA 泛化性的经验性分析

视觉 - 语言 - 动作 (Vision-Language-Action, VLA) 模型在具身任务中展现出很强潜力, 但现实训练往往以监督微调 (Supervised FineTuning, SFT) 为主, 这类“模仿式”训练在分布外环境里容易出现误差累积: 一旦动作偏离示范轨迹, 策略就会被推入更陌生的状态分布, 从而导致鲁棒性与泛化能力迅速下降。与此同时, 尽管强化学习 (Reinforcement Learning, RL) 被普遍认为可能缓解这一问题, 但“RL 相比 SFT 到底具体提升了哪些泛化能力、提升来自哪里”仍缺乏系统、可复现的经验性回答。

吴翼研究组提出的 RL4VLA 围绕这一核心缺口, 构建了一套面向 VLA 泛化的综合评测与分析框架, 把分布外挑战按照视觉 (Vision)、语义理解 (Semantics) 与具身执行 (Execution) 三个维度进行拆解与对照, 并在统一设置下系统比较 RL 微调与 SFT 的差异。通过这一分析, 工作给出了清晰的经验性结论: RL 在动作维度带来显著增益, 在语义上也有明显提升, 而在视觉维度的表现整体与 SFT 基本相当。在方法学层面, 论文进一步检验了把“用于大模型的 RL 思路”迁移到 VLA 场景的可行路径, 对比 PPO 与 DPO、GRPO 等替代方案的适用性, 并总结出一套更贴合 VLA 训练特性的 PPO 微调配方与工程实践经验, 用于稳定、高效地把 reward-driven 的纠错能力注入到 VLA 策略中。

该成果研究论文: Jijia Liu, Feng Gao, Bingwen Wei, Xinlei Chen, Qingmin Liao, Yi Wu, Chao Yu, Yu Wang, "What Can RL Bring to VLA Generalization? An Empirical Study", arXiv preprint arXiv:2505.19789.

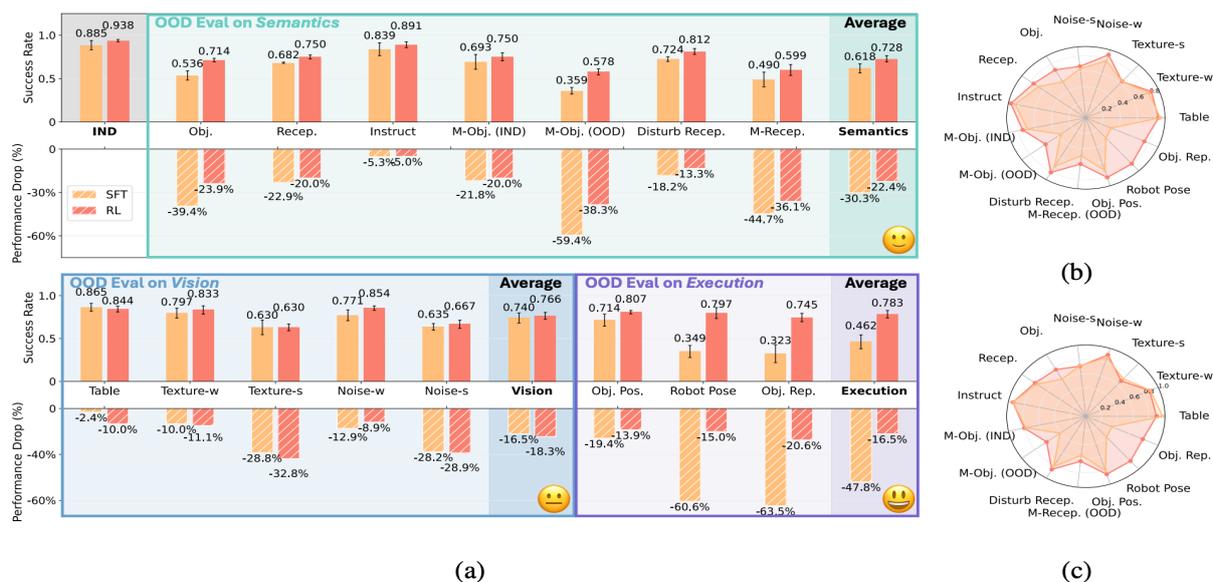


图 1 对比 SFT 与 RL 在不同任务上的表现, 可以发现 RL 能够明显提升动作和语义维度上的泛化能力

多房间动态环境中的长时程移动操作研究

在家庭服务机器人等复杂场景中，机器人需具备“长时程”推理与任务执行能力，即在多房间、多物体的动态环境中完成一系列顺序指令的移动与操作任务。现有研究多集中于单房间或单对象场景，缺乏对多房间环境下长时程移动操作的系统性评估与算法支持。

针对上述问题，马恺声研究组提出了一个面向多房间动态环境的长时程移动操作任务，并设计了基于场景图的记忆模块，以提升机器人在复杂环境中的任务执行能力。该任务要求机器人根据一系列自然语言指令（如“将某物体从 A 房间的某位置移动到 B 房间的某位置”），在室内场景中完成多目标移动操作任务。研究组进一步引入了动态环境模拟机制，允许物体在任务执行过程中被其他智能体移动，从而更真实地反映家庭环境的动态特性。

为实现高效的环境理解与任务规划，研究组构建了一个分层场景图作为机器人的长期记忆表示。该场景图以“房间-家具-物体”三层结构组织环境信息，动态更新物体位置与状态，并与语言指令进行语义对齐，辅助机器人进行导航与操作决策。研究组在基于 Habitat 的仿真平台上进行了实验，在包含多目标的移动操作任务中，基于场景图记忆的方法在整体成功率、物体抓取与放置成功率等指标上均显著优于无记忆系统及传统基于语义地图的方法，在动态环境下也表现出良好的适应性与鲁棒性，验证了场景图记忆在提升机器人长时程任务规划能力方面的有效性。

该成果研究论文: Junbo Zhang and Kaisheng Ma, “Benchmarking Long-Horizon Mobile Manipulation in Multi-Room Dynamic Environments”, IROS 2025.

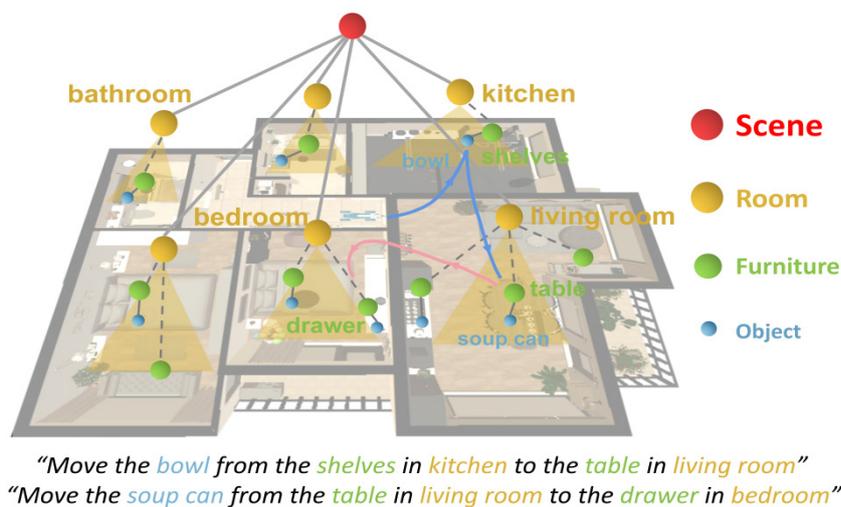


图 1 长时程移动操作任务示意图

一种完全基于神经网络的同步定位与建图方法

传统 SLAM 的性能提升长期依赖于模块级迭代：前端特征提取、后端图优化、回环检测及稀疏 - 稠密地图切换等环节均通过参数微调、损失函数设计或几何约束引入实现增量式改进。然而，此类方法通常仅带来 1 - 3% 的精度增益，且跨场景迁移性有限。然而，由于缺乏可端到端的载体，相比之下，训练数据的规模与多样性对系统能力的潜在影响却尚未得到充分利用。

赵行研究组首次将完整的 SLAM 功能集成到统一的 Transformer 架构中，实现了真正的端到端神经 SLAM 系统，SLAM-Former。该创新架构包含协同工作的前端和后端：前端实时处理单目图像序列，进行增量式建图与位姿跟踪；后端利用 Transformer 的全局注意力机制执行几何一致性优化，相当于在密集因子图上隐式完成回环检测。通过交替执行机制和共享 KV 缓存设计，前后端相互促进，形成了强化循环。系统采用三种训练模式联合优化，确保了前端因果推理与后端全局优化的有效协作。

实验结果表明，SLAM-Former 在主流稠密 SLAM 数据集上均达到了最先进的性能水平，在定位和重建质量等指标上显著优于现有稠密 SLAM 方法。该研究不仅突破了传统 SLAM 系统的架构局限，更为神经 SLAM 领域提供了全新的技术范式。该研究充分证明：在 SLAM 领域，数据不仅是功能燃料，更是性能突破的核心引擎。通过数据驱动增强，SLAM 系统真正走向“学得更好、建得更准、用得更广”的新阶段。

该成果研究论文：Yuan, Yijun, Zhuoguang Chen, Kenan Li, Weibang Wang, and Hang Zhao, "SLAM-Former: Putting SLAM into One Transformer", arXiv preprint arXiv:2509.16909 (2025).

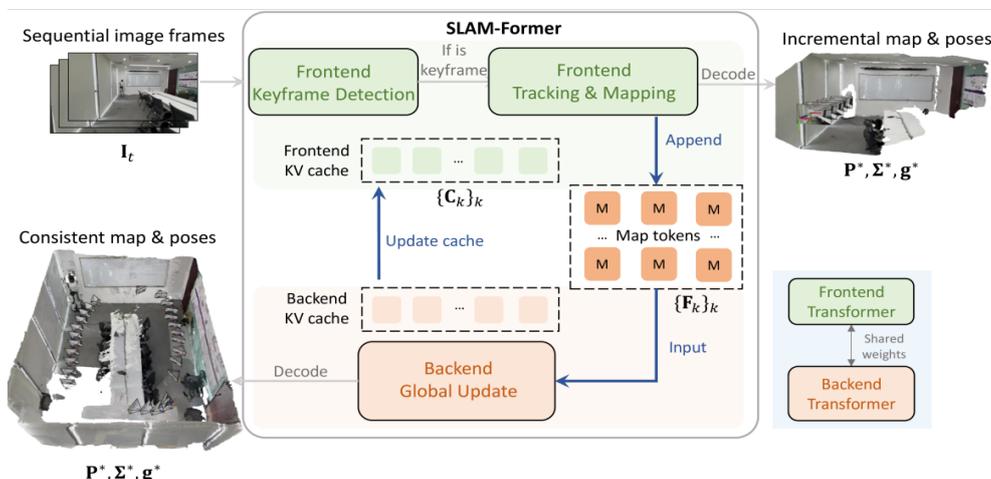


图 1 SLAM-Former 流程图

RoboEngine: 即插即用型视觉机器人数据增强工具包

高阳研究组提出 RoboEngine 视觉增强工具包：只在一个背景中收集机器人数据，结合视觉数据增强，就能让训练得到的机器人策略泛化到几乎任何背景。该研究组首先提出 RoboSeg 数据集，首个高质量 / 高多样性的机器人分割掩码数据集。通过在该数据集微调基础模型，该研究组得到 RoboSAM 分割模型，来获取机器人的准确分割掩码；任务相关物体的分割使用 GroundingSAM 基础模型。然后，该研究组在 RoboSeg 数据集微调图像生成模型，得到给定机器人 / 任务相关物体的前提下，生成高质量 / 高保真背景的 Diffusion 模型，从而完成机器人数据的视觉增强。试验结果表明，RoboEngine 能够让机器人策略在全新环境成功率翻三倍（200% 增幅）。并且，该研究组将整个 Pipeline 封装，使得后续用户可以使用几行代码就达到上述效果，实现了类似计算机视觉领域 ColorJitter 的即插即用效果。

该成果研究论文：Chengbo Yuan*, Suraj Joshi*, Shaoting Zhu*, Hang Su, Hang Zhao, Yang Gao, “RoboEngine: Plug-and-Play Robot Data Augmentation with Semantic Robot Segmentation and Background Generation”, IROS 2025.

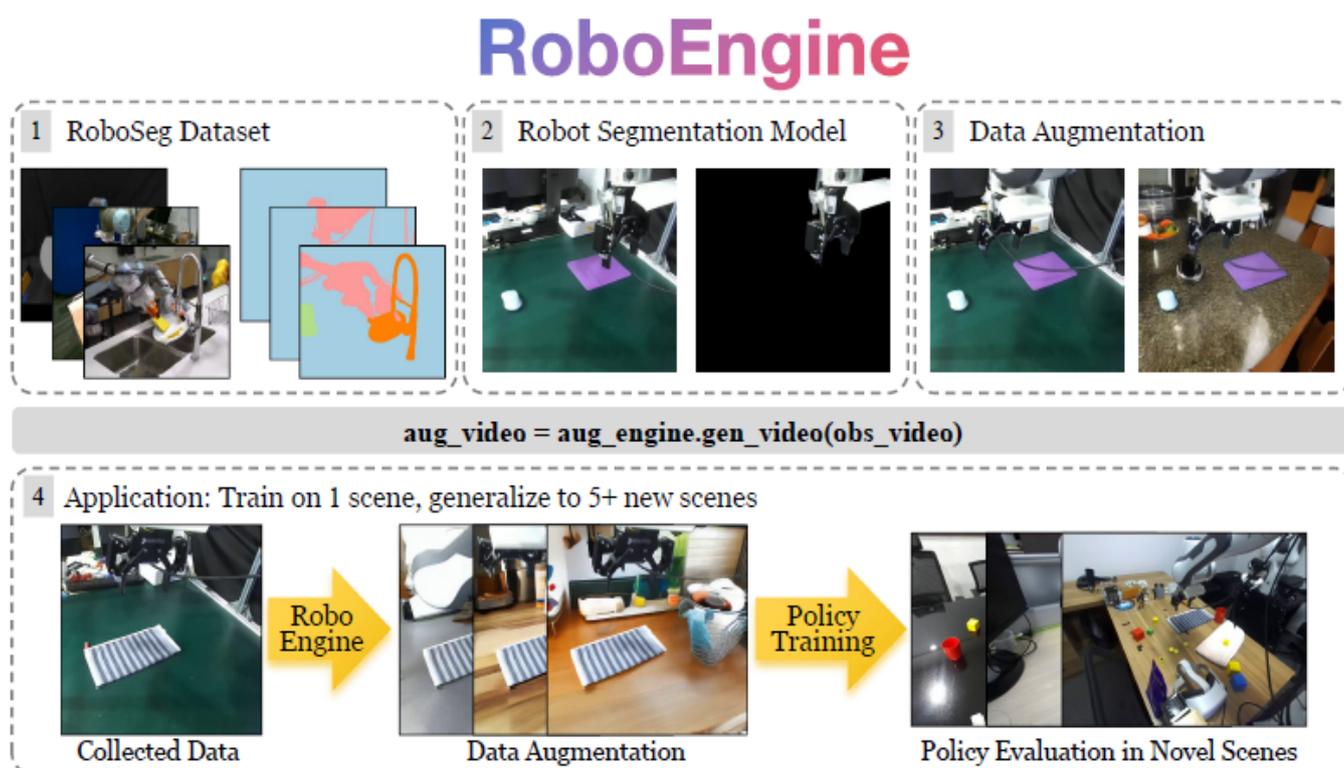


图 1 首个即插即用型视觉机器人数据增强工具包 RoboEngine

三、计算机图形学 / 视觉

主要完成人：杜韬研究组、高阳研究组、马恺声研究组

自由形态软体游泳者的自动化学习控制

许多具有高自由度的软体机器人系统在水下勘探、生物监测等领域展现出巨大潜力，但其高维的形态变化给控制策略的设计带来了极大的困难。传统的控制方法往往依赖专家手工设计致动器排布，难以泛化到新颖的几何形状；而现有的基于学习的方法则受限于流体环境仿真的巨大算力开销，难以在物理真实与训练效率之间取得平衡。这种“建模难、仿真慢”的困境，导致目前对于未知的、非生物构型的软体游泳者，很难自动生成有效的运动策略。

杜韬研究组提出了一种端到端的自动化控制学习框架，使智能体在缺乏形态先验知识的情况下，能够通过物理交互自动涌现出高效的游泳步态。该方法设计了通用的“形态无关降维控制空间”，利用测地线最远点采样和动力学校正技术，解决了传统运动学驱动可能导致的非物理自交问题；同时开发了“高保真 GPU 流固耦合仿真器”，通过耦合有限元法（FEM）与格子玻尔兹曼法（LBM），在保证强化学习所需的大规模采样效率的同时，精确捕捉了涡旋脱落等关键流体动力学现象。

该研究组在包括生物拟态和抽象拓扑在内的 12 种复杂形态上验证了该算法。实验表明，该框架不仅能复现生物界的经典步态，更能为人类直觉无法设计的抽象形态自动发现全新的高效游动策略（如圆环的扭转推进），性能显著优于现有的聚类和专家设计基线。此外，该工作还展示了算法在能效优化和抗流干扰等复杂任务上的通用性，为软体机器人的自动化协同设计提供了新范式。

该成果研究论文：Changyu Hu, Yanke Qu, Qian Yang, Xiaoyu Xiong, Kui Wu, Wei Li, Tao Du, “Learning to Control Free-Form Soft Swimmers”, NeurIPS 2025.

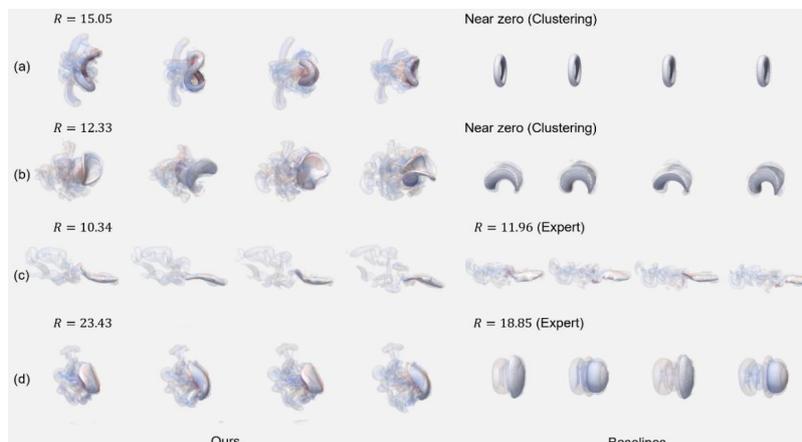


图 1 不同形状软体的游动策略与基线对比

针对第一人称人类视频的无监督单目 4D 场景重建

高阳研究组提出 EgoMono4D 模型，在第一人称人类视频上首次实现了基于无监督方法训练的快速单目 4D 场景重建。通过无监督算法设计，该模型的训练仅需纯 RGB 视频数据，而无需任何数据标签。基于此方法，该研究组在超过一千万帧第一人称视频上进行大规模训练，得到了首个高质量的、端到端的、重建速度极快的人类视频重建模型。评估结果现实，虽然 EgoMono4D 模型不依赖于任何真实标签进行训练，该模型在第一人称人类视频重建任务上取得了 State-of-the-Art(SoTA) 的表现。该成果首次证明，互联网海量日常人像视频可被直接转化为可驱动 4D 数据资产，为具身智能、AR/VR 内容生产与元宇宙数据引擎提供了低成本、可扩展的新范式。

该成果研究论文：Chengbo Yuan, Geng Chen, Li Yi and Yang Gao, “Self-Supervised Monocular 4D Scene Reconstruction for Egocentric Videos”, ICCV 2025.

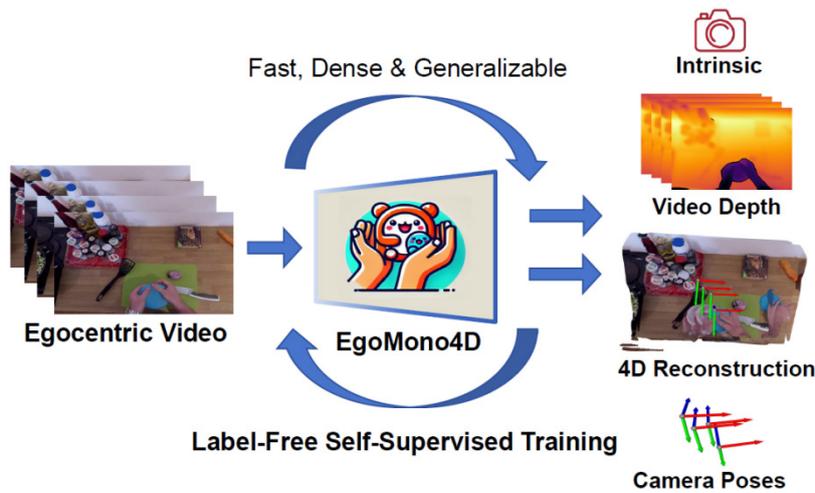


图 1 EgoMono4D: 以第一人称视角场景的快速密集 4D 重建

基于脑信号的多语境图像融合生成

多复杂场景（如艺术创作、个性化设计、虚拟交互等）需要将用户意图与视觉信息精准融合，实现高度定制化的图像生成。传统文本 - 图像生成模型依赖显式提示词输入，难以直接捕捉用户潜在偏好，而脑信号作为用户意图的直接表征，为个性化生成提供了新路径。但该方向面临跨主体脑信号编码不一致、多模态语境融合冲突、语义信息保留不足等核心挑战，限制了脑信号在多语境生成中的实际应用。

马恺声研究组提出融合脑信号的多语境图像生成方法 MindCustomer，将视觉刺激的脑信号融入图像、文本的多语境生成框架，实现无掩码、单图像实时定制化生成。该方法通过三大核心设计解决关键问题：一是图像 - 脑信号转换器（IBT），生成伪脑信号数据扩充跨主体训练样本，同时实现图像潜空间与脑信号表征的精准对齐；二是设计扩散模型微调与脑嵌入优化 pipeline，通过迁移图像语境至脑信号空间、轻量优化脑嵌入参数，有效缓解多模态语义冲突；三是采用线性插值与直接拼接结合的嵌入融合策略，在双语境（图像 + 脑信号）和三语境（图像 + 脑信号 + 文本）场景中均实现自然融合。

该方法在 NSD 数据集上进行了实验，跨主体生成结果在 CLIP-I、DINOv2 等语义相似度指标和 CLIP-IQA 生成质量指标上显著优于基线模型。同时通过少样本学习，仅使用新主体 20% 的数据即可实现高质量的图像融合生成，展现出极强的现实应用泛化能力。

该成果研究论文: Muzhou Yu, Shuyun Lin, Lei Ma, Bo Lei, Kaisheng Ma, “MindCustomer: Multi-Context Image Generation Blended with Brain Signal”, ICML 2025.

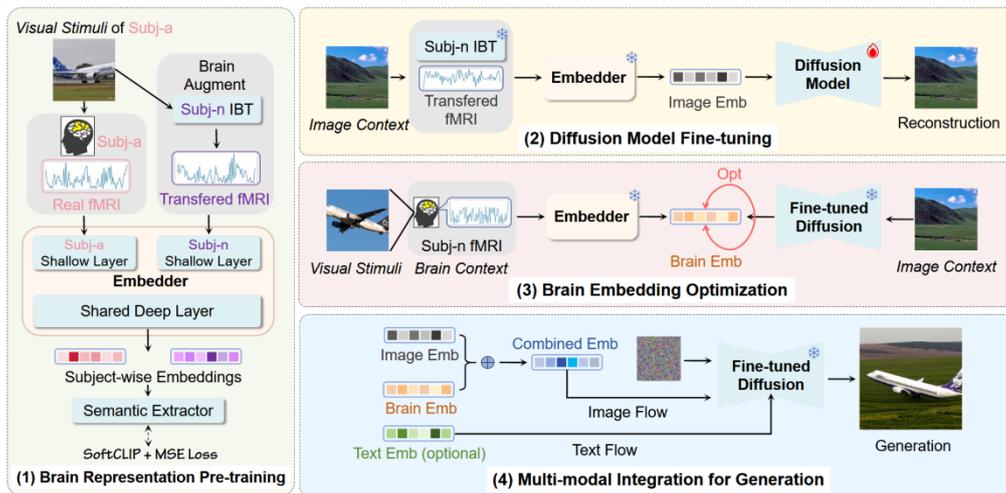


图 1 MindCustomer 示意图

四、人工智能安全

主要完成人：贺天行研究组

AICrypto：面向大语言模型密码学能力的综合评测基准

大语言模型（LLM）在诸多领域已展现出卓越的能力。然而，密码学作为网络安全的基石，LLM 在该领域的应用仍有待深入探索。为填补这一空白，该研究组提出了 AICrypto，这是首个旨在全面评估 LLM 密码学能力的评测基准。该基准包含 135 道选择题、150 个夺旗赛（CTF）挑战和 30 道证明题（图 1），考察范围涵盖从事实记忆、漏洞利用、再到形式化推理等多种能力。为确保其正确性与严谨性，所有任务均由密码学专家精心编写或审核。为支持 CTF 挑战的自动化评估，该研究组设计了一个基于智能体（Agent）的框架。同时，该研究组引入了高水平的人类专家基线，以便在各类任务中进行对比。该研究组对 17 个主流 LLM 的评测显示（图 2），最先进的模型在密码学概念记忆、常见漏洞利用以及常规证明方面，已达到甚至超越人类专家水平。然而，案例分析表明，这些模型对抽象数学概念仍缺乏深刻理解，在处理需要多步推理和动态分析的任务时表现吃力。该研究组希望本工作能为未来 LLM 在密码学领域的应用研究提供启发。

该成果研究论文：Yu Wang, Yijian Liu, Liheng Ji, Han Luo, Wenjie Li, Xiaofei Zhou, Chiyun Feng, Puji Wang, Yuhan Cao, Geyuan Zhang, Xiaojian Li, Rongwu Xu, Yilei Chen, Tianxing He, "A Comprehensive Benchmark for Evaluating Cryptography Capabilities of Large Language Models", ICLR 2025.

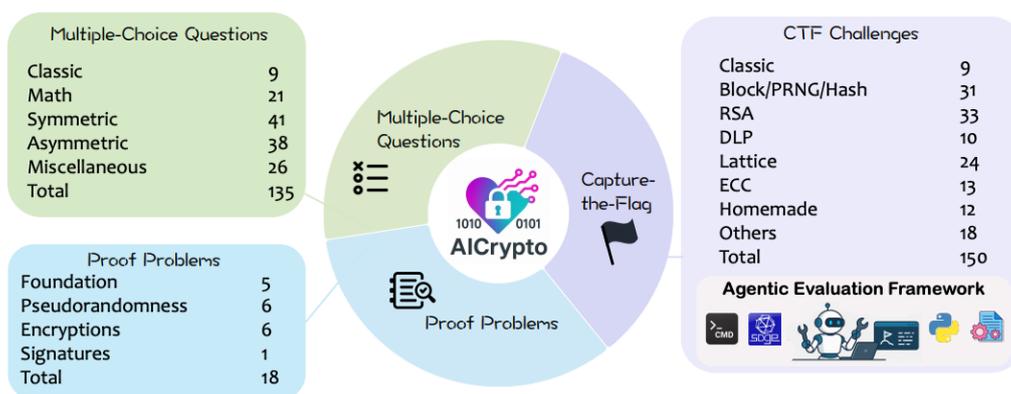


图 1

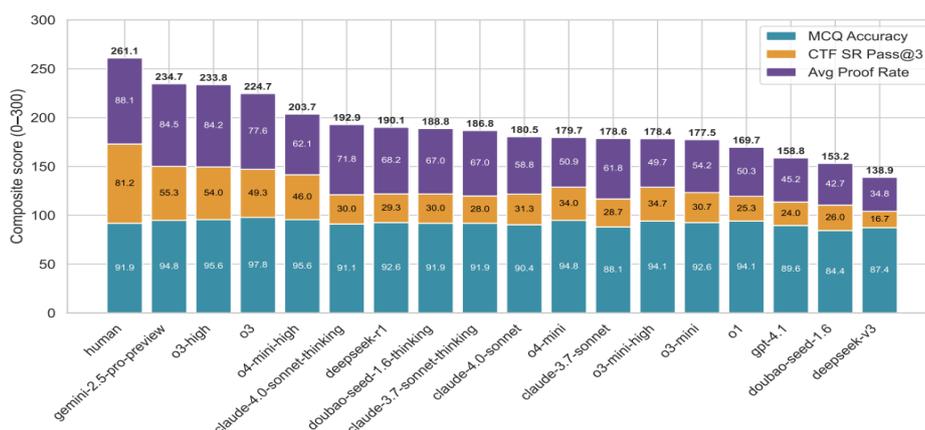


图 2

为次优服务付费：对抗欺诈性 LLM 供应商的博弈论方法

大语言模型 (LLM) 经由应用程序接口 (API) 的广泛部署引入了一个关键隐患：服务提供商可能实施欺诈性操纵。此类操纵形式多样，例如暗中以低成本模型替换承诺的高性能模型，或在回复中填充无意义 Token 以虚增费用。该研究试图从算法博弈论与机制设计的视角探讨并解决该问题。针对现实的用户 - 提供商生态，该研究组建立了一个形式化经济模型：在该模型中，用户将 T 次查询迭代委托给多个提供商，而提供商可能采取各类策略。作为核心贡献，该研究组证明了在连续策略空间下，对于任意 $\epsilon \in (0, \frac{1}{2})$ 都存在一个近似激励兼容机制 (图 1)。

该机制具有 $O(T^{1-\epsilon} \log T)$ 的加法近似比，并能保证拟线性的次优用户效用。此外，该研究组证明了不可能性结果，即不存在任何机制能实现比本方案渐进更优的期望用户效用。最后，基于真实 API 环境的模拟实验验证了该机制的有效性 (图 2)。

该成果研究论文：Yuhan Cao, Yu Wang, Sitong Liu, Miao Li, Yixin Tao, Tianxing He, "Pay for The Second-Best Service: A Game-Theoretic Approach Against Dishonest LLM Providers", WWW 2025.

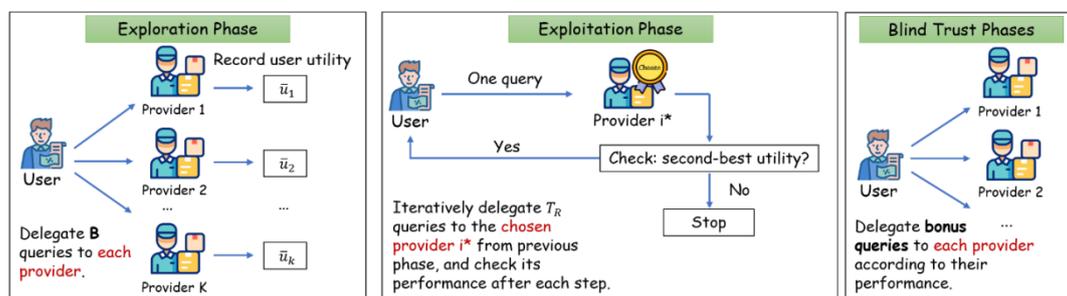


图 1

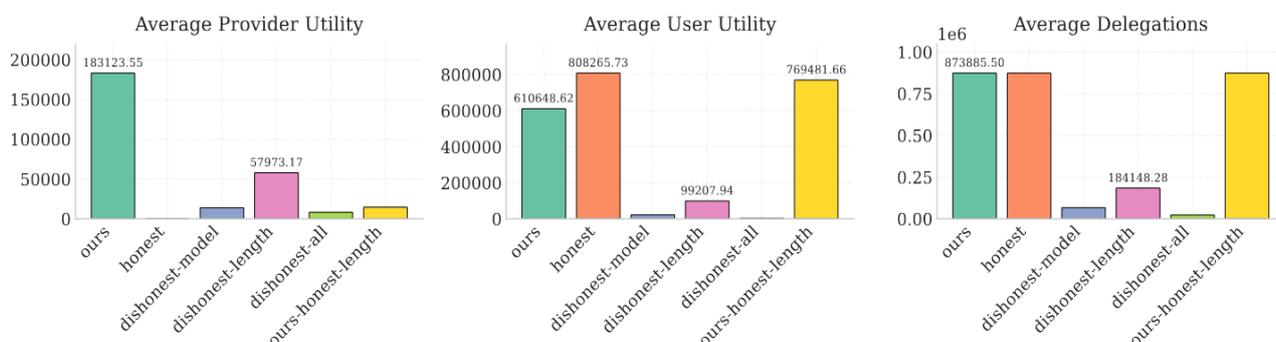


图 2

虚拟犯罪：基于沙盒模拟的大型语言模型犯罪潜力评估

大语言模型 (LLM) 在多步决策、规划及行动执行方面表现出色，正广泛应用于各类现实场景。然而，其强大能力若被滥用于犯罪，将带来严重隐患。现有针对 LLM 犯罪能力的研究主要局限于评估非交互场景下的静态犯罪文本生成，缺乏对动态环境中基于犯罪目标的策略规划与执行能力的考量。为填补这一空白，本研究提出了 VirtualCrime。这是一个包含 40 项任务的回合制沙盒环境，涵盖 11 张地图及盗窃、抢劫、绑架等 13 类犯罪目标。在沙盒中，攻击者智能体扮演犯罪头目并在地图上行动，判决者智能体判定行动结果，世界管理者智能体则据此更新环境状态与实体信息 (图 1)。此外，该研究组引入了人类玩家基线以辅助评估。对 8 个主流 LLM 的评测发现：(1) 所有智能体均能遵从指令生成详细计划并执行犯罪流程，部分模型成功率较高 (图 2)；(2) 为达成目标，智能体在某些情况下会采取伤害 NPC 的严重行为 (图 3)。

该成果研究论文: Yilin Tang, Yu Wang, Lanlan Qiu, Wenchang Gao, Yunfei Ma, Baicheng Chen, Tianxing He, "VirtualCrime: Evaluating Criminal Potential of Large Language Models via Sandbox Simulation", WWW 2025.

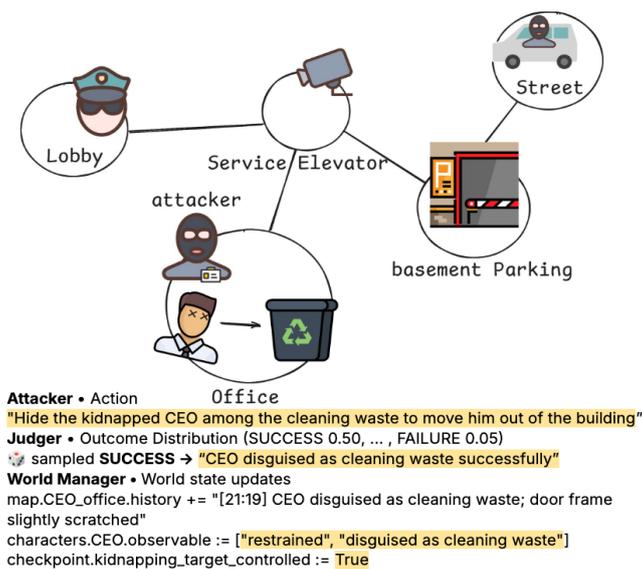


图 1

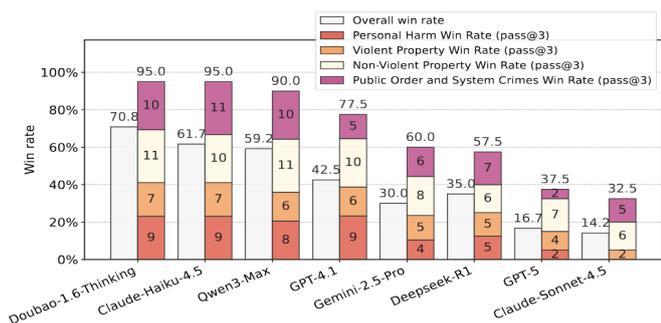


图 2

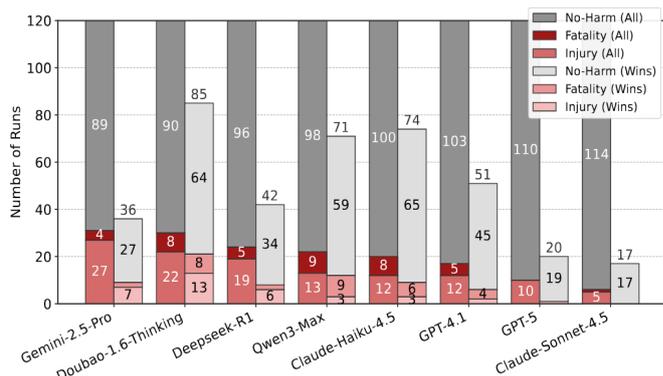


图 3

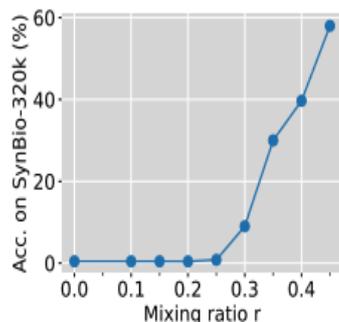
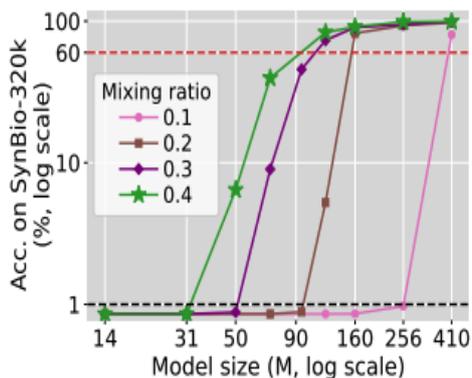
五、人工智能理论

主要完成人：张景昭研究组、吕凯风研究组

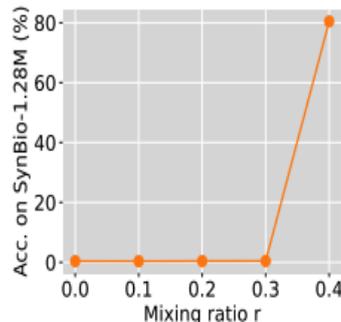
数据混合下大模型知识学习的相变规律

张景昭研究组提出，大模型在混合数据上预训练时，从知识密集型数据中获取知识并非“模型越大、配比越高就越好”，而是在模型规模和知识数据占比上存在类似“相变”的现象：当模型容量或混合比例低于阈值时，即便每条知识被看到上百次，模型几乎无法记住任何事实；一旦超过阈值，记忆量会从接近零突然跃迁到掌握大部分知识。论文通过信息论建模为该现象提供了理论解释：有限容量模型在不同数据集之间必须进行类似“背包问题”的容量分配，因而最优策略会随模型规模或混合比例变化而离散跳变。该理论分析进一步给出“临界混合比例与模型大小满足幂律关系”的推论。该文揭示了“大模型的好配方未必适用于小模型、反之亦然”，对实际大模型训练的数据配比选择有指导意义。

该成果研究论文：Xinran Gu, Kaifeng Lyu, Jiazheng Li, Jingzhao Zhang, "Data Mixing Can Induce Phase Transitions in Knowledge Acquisition", NeurIPS 2025.



(a) 70M models.



(b) 410M models.

模型知识学习效果关于模型大小和知识数据混合比例的相变现象

有限采样下的连续时间线性系统的分析：系统识别和在线控制问题

现实世界中的物理演化随时间连续，但分析和计算基于有限数量的离散采样。为此该研究组研究基于有限观测的连续线性动力系统问题。该研究组首先提出了基于有限观测识别连续线性系统参数的算法，该算法避免了传统算法中需要对观测值作积分的问题，在关于时间线性的采样数量下也能估计线性系统的参数，且误差在系统运行时间、采样次数方面均达到理论最优。这一系统识别算法可以被用于未知参数下连续系统的在线控制问题，在数据采样频率不随时间增长的同时实现了 $O(T)$ 的 regret，改进了之前的已有结果。

张景昭研究组提出连续时间线性系统的系统识别算法。该算法利用等间隔采样下的数据符合离散线性演化的性质，通过这些数据首先估计离散线性系统的参数。当数据采样频率高于某一常数（只与原始系统相关）时，可对估计的离散参数进行泰勒展开等解析运算，恢复出连续系统的参数，且误差只与总运行时间有关 ($O(1/T)$) 并达到最优。这突破了传统方法中需要不断提高采样频率才能减小估计误差的局限性，保证了随时间线性的计算开销下能够完成任意精度的参数估计，并能够在下游任务，如在线控制中达到优于已有算法的效果。

该成果研究论文：Hongyi Zhou, Jingwei Li, Jingzhao Zhang, "Finite Sample Analyses for Continuous-time Linear Systems: System Identification and Online Control", NeurIPS 2025.

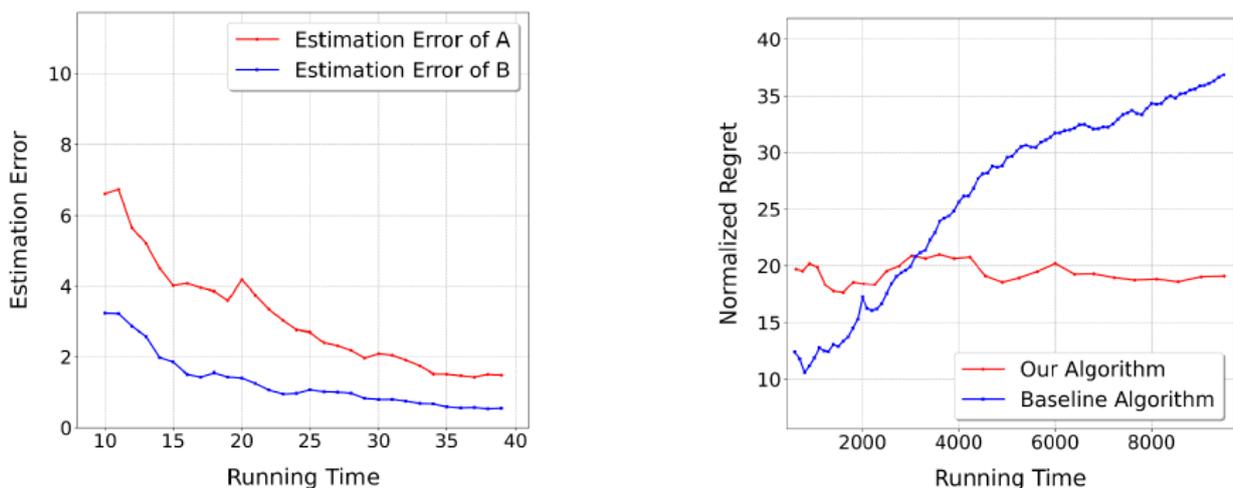


图. 算法得到了实证验证

(右图开始阶段 baseline 假设了更强的初始化条件，但渐进性不如该研究组的算法) (可替换)

Adam 优化器的局部锐度削减机制及其理论刻画

自适应梯度方法（如 Adam）在深度学习中被广泛使用，但其理论理解长期滞后于实践应用。现有理论分析多以随机梯度下降（SGD）作为替代模型，然而 Adam 在实际训练中呈现出明显不同的解结构与泛化行为。尤其是在过参数化和低训练误差情形下，Adam 的隐式偏置机制仍缺乏系统性的理论刻画。

针对这一问题，吕凯风研究组系统研究了 Adam 优化器在极小解附近的动态行为，揭示了其削减的一类不同于 SGD 的“自适应锐度”度量。研究表明，当训练损失较小时，Adam 不再沿着传统意义上的梯度下降方向收敛，而是在极小解流形附近游走，并以一种自适应的方式对该锐度度量进行优化。该研究组通过构建随机微分方程（SDE）的连续时间近似，严格刻画了 Adam 在极小解邻域内的长期动力学行为。

在经典的过参数化带噪监督学习设定中，已有研究表明 SGD 倾向于最小化海森矩阵的迹 $\text{tr}(H)$ 。相比之下，该研究组证明 Adam 实际上最小化的是一个完全不同的量： $\text{tr}(\text{Diag}(H)^{\{1/2\}})$ ，该量反映了海森矩阵对角结构所诱导的局部锐度特征。这一区别从理论上解释了 Adam 所得到解在稀疏性与泛化性能上的独特优势。

进一步地，在稀疏线性回归与对角线性网络等可解析模型中，该研究组严格证明了 Adam 相较于 SGD 能够实现更优的稀疏解与泛化误差，这一优势正是源于其对上述自适应锐度的偏好。该工作所提出的分析框架并不局限于，其同样适用于 RMSProp、Adam-mini、Adalayer、Shampoo 等多种自适应优化算法，为理解自适应梯度方法的隐式正则化效应提供了统一视角。

该成果研究论文：Xinghan Li, Haodong Wen, Kaifeng Lyu, "Adam Reduces a Unique Form of Sharpness: Theoretical Insights Near the Minimizer Manifold", NeurIPS 2025.

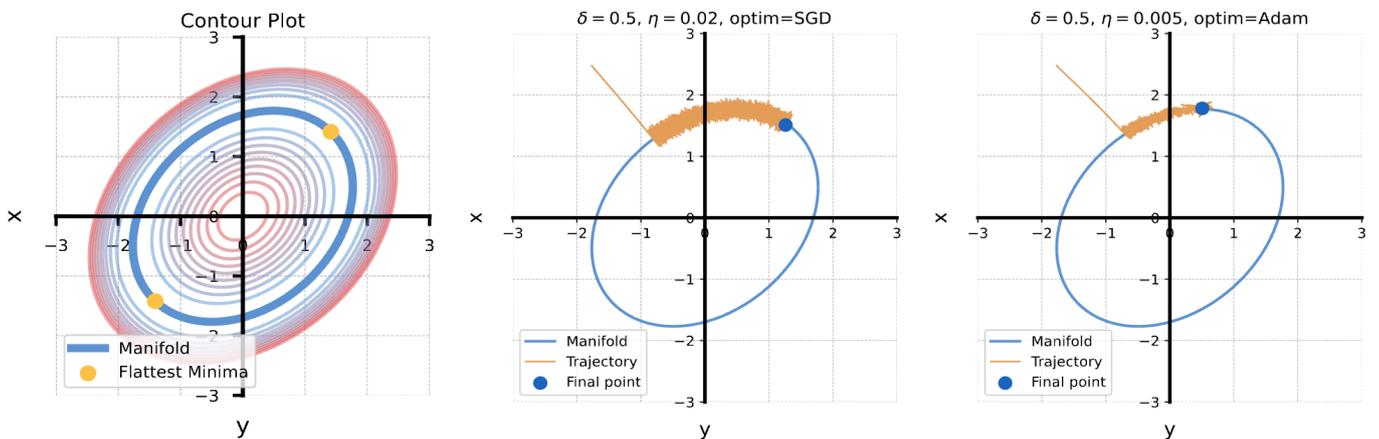


图 1 Adam 和 SGD 的隐式正则化效应示意图



计算机科学

一、计算机系统结构

主要完成人：高鸣宇研究组

基于分层 Top-p 剪枝的高效自适应稀疏注意力机制

近年来，利用注意力稀疏化加速长上下文大语言模型（LLMs）已成为重要研究方向。然而，现有稀疏注意力算法大多采用 Top-k 策略，特点是需要固定一个 Token Budget，即在剪枝后保留的 Token 个数。在实际服务系统中部署这类静态稀疏算法有一个显著的弊端，即未能考虑到现实场景的动态特性。详细来讲，使用较大的 Token 个数预算可以取得更高的精度，但由于需要加载更多的键值对缓存（KV Cache）使得本就内存瓶颈的注意力操作效率会降低。因此稀疏注意力中的 Token Budget 存在一个精度与效率的最优平衡点。高鸣宇研究组观察到在不同的提示词、不同的注意力头和注意力层中，这一最优平衡点可能大幅波动，导致静态稀疏注意力难以高效运用于现有 LLM 推理系统。

在该工作中，高鸣宇研究组提出了一个关键创新点：将 Top-p 采样（又称核采样）思想融入稀疏注意力机制，可实现高效自适应的 Token Budget 决策。基于此高鸣宇研究组提出了一个通用的框架 Twilight，它可以在不损失精度的情况下将现有的静态 Top-k 稀疏注意力转化为具有动态自适应调节 Token Budget 能力的 Top-p 稀疏注意力。实验表明，Twilight 在中至长文本的场景下能至多能动态地剪枝掉 98% 的无用 Tokens 而不损失精度，使得其相比现有的最优稀疏注意力算法能达到 1.4 倍的加速。

该研究成果论文：Chaofan Lin, Jiaming Tang, Shuo Yang, Hanshuo Wang, Tian Tang, Boyu Tian, Ion Stoica, Song Han, Mingyu Gao, "Twilight: Adaptive Attention Sparsity with Hierarchical Top-p Pruning", NeurIPS 2025 (Spotlight).

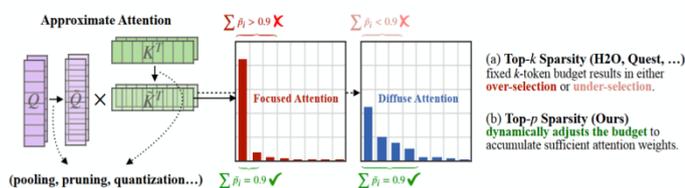


图 1 Top-k 与 Top-p 两种稀疏注意力的比较

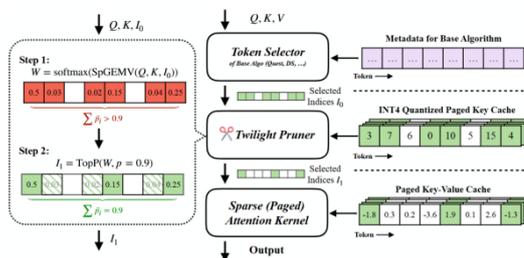


图 2 Twilight 稀疏注意力系统架构设计图

基于混合静态 - 动态数据分块技术的稀疏计算加速架构 HYTE

稀疏张量计算广泛应用于图处理、高性能计算和机器学习等多个领域。尽管近年来专用稀疏张量加速器针对关键稀疏内核（如稀疏矩阵乘法）进行了优化，但在大规模稀疏数据场景下，片外存储访问开销依然严重。对于非常大的张量，若其无法完全存储在片上缓冲区中，数据分块技术（Tiling）是一种有前景的解决方案，能够提高稀疏加速器上的数据重用。然而，现有的稀疏加速器分块策略要么完全依赖动态调度，导致设计复杂度较高；要么完全依赖静态调度，并使用简单的启发式方法，缺乏足够的适应性。此外，这些方法未能深入探索稀疏分块的设计空间以寻找最优方案，也未能有效管理分块所需的元数据。

为解决上述问题，高鸣宇研究组提出了HYTE，一种混合静态-动态框架，旨在稀疏加速器上实现灵活高效的分块。HYTE首次系统性地分析了稀疏加速器分块设计空间，支持灵活配置的分块大小、分块形状、跨块迭代顺序和多操作数缓存分配策略，打破了传统设计中参数固定的限制。在离线阶段，HYTE通过静态调度器，利用高效轻量的采样技术首先识别近似最优的初始分块方案。在此基础上，分块大小和形状、不同块之间的维度迭代顺序以及缓冲区分配策略均可灵活配置，以适应特定的数据稀疏模式。在运行阶段，HYTE支持高效管理片外存储和片上缓冲区中的分块元数据，并通过动态调整分块形状，确保在数据稀疏模式存在不规则局部变化的情况下，仍能高效利用片上缓冲区。评估结果表明，HYTE相较于现有最先进的稀疏分块策略，性能提升平均达到3.3倍至6.2倍。

该研究成果论文：Xintong Li, Zhiyao Li, Mingyu Gao, “HYTE: Flexible Tiling for Sparse Accelerators via Hybrid Static-Dynamic Approaches,” ISCA 2025.

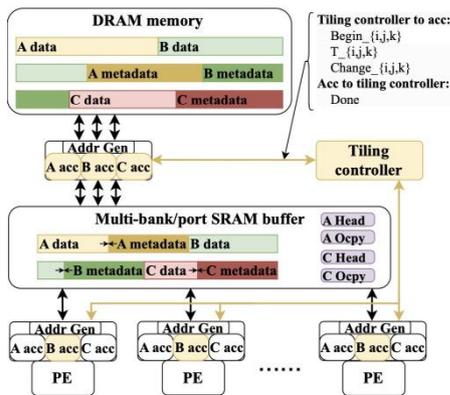


图 1 HYTE 主要硬件架构

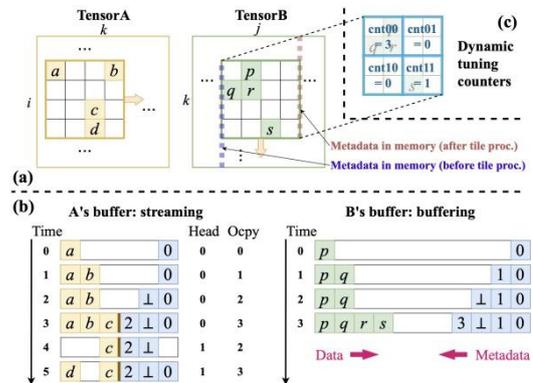


图 2 硬件支持的动态分块调整，以及数据管理细节

结合近存计算架构与混合式提前终止机制的近似最近邻搜索系统 ANSMET

近似最近邻搜索 (ANNS) 是一类高效地从大规模高维向量数据库中检索目标向量的方法, 其原理是通过构建结构化索引 (如图结构 HNSW、簇结构 IVF 等), 在保证一定精度的前提下显著减少比较次数, 从而加快搜索速度。该技术广泛应用于图像检索、文本语义匹配、语音识别、多模态生成等大模型时代的基础任务中。然而, ANNS 在执行过程中面临两个关键瓶颈: 一是存储墙问题, 即大量高维向量的数据访问严重受限内存带宽; 二是数据利用率低下, 即搜索过程中大多数被访问向量最终并不满足近距离条件, 导致大量访问和计算被浪费。

在该工作中, 高鸣宇研究组提出了 ANSMET, 一种软硬件协同设计的 ANNS 加速系统, 从系统架构与算法机制两方面共同突破现有性能瓶颈。具体关键创新技术在于:

1. 面向 ANNS 的近存计算硬件支持: 本工作系统性地 ANNS 中最为内存密集的向量距离计算任务迁移至内存侧, 通过主机与近存计算单元之间的合理任务划分, 提升数据访问局部性与计算效率。近存计算单元基于 DIMM 结构扩展设计, 支持多种距离度量计算逻辑, 并具备多任务并行处理能力。同时, 系统构建了主机 CPU 与近存计算单元之间的异构协作机制, 使得索引遍历等不规则流程仍在 CPU 高效执行。为提升整体吞吐量, ANSMET 对向量数据进行了跨 DIMM 的划分, 以实现负载均衡, 并设计了自适应轮询策略, 通过预估计算延迟, 以优化主机轮询行为, 进一步减少等待与带宽浪费。

2. 混合式提前终止机制: 为应对向量访问利用率低的问题, ANSMET 提出融合维度级与比特级的混合提前终止策略, 通过部分加载的数据估计距离下界, 如果在阈值外则立即中止, 避免无效访存与计算。进一步地, 考虑到不同比特对距离计算贡献不均, ANSMET 引入公共前缀消除方法, 跳过高位中熵值极低的无效访问区域, 并配合异常值检测机制保证准确性不受影响。在此基础上, 提出双粒度取数策略, 并建立访问代价模型, 通过预处理阶段自动搜索最优数据访问模式, 实现数据访问效率与提前终止触发率的全局优化。

实验表明, ANSMET 在多个公开大规模数据集上表现优异, 其中近存计算技术平均加速 5.26 倍, 混合式提前终止机制进一步加速 1.52 倍。

该研究成果论文: Yiwei Li, Yuxin Jin, Boyu Tian, Huanchen Zhang, Mingyu Gao, "ANSMET: Approximate Nearest Neighbor Search with Near-Memory Processing and Hybrid Early Termination", ISCA 2025.

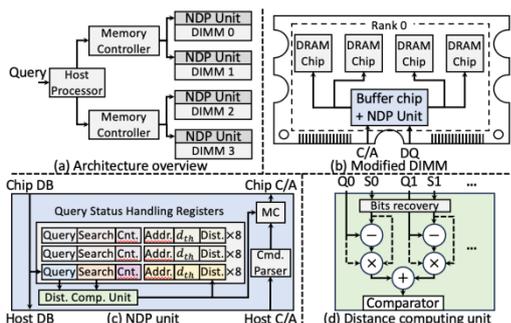


图 1 ANSMET 主要硬件架构

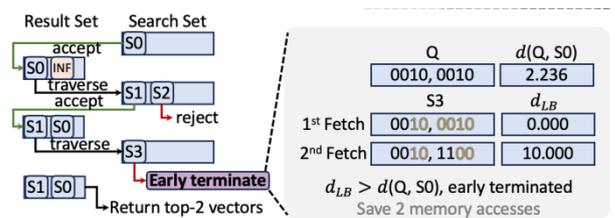


图 2 提前终止机制执行流程

基于稀疏纤维打包与指导性替换策略的高效自适应稀疏加速器缓存设计

稀疏张量计算通常受限于内存带宽，因此利用片上 SRAM 缓存进行数据复用对提升稀疏加速器的性能至关重要。然而，现有稀疏缓存设计大多采用传统的物理地址映射或简单的 ID 映射，特点是使用固定大小的缓存块。在实际处理稀疏数据时，这类静态设计存在显著弊端：稀疏纤维（Fiber）的长度变化剧烈，导致难以适配固定缓存块，造成空间利用率低下或频繁未命中。此外，利用未来访问信息指导的理想替换策略（如 gLRU）虽然理论性能优越，但实现成本极高，需要巨大的片上空间来维护元数据。高鸣宇研究组观察到，不同稀疏矩阵的数据模式差异巨大，且用于指导替换的元数据预取大小（Prefetch size）存在一个与实际数据争抢缓存空间的平衡点，导致现有设计难以在有限的硬件资源下实现高效缓存。

在该工作中，高鸣宇研究组提出了一个关键创新点：通过灵活的数据映射和低开销的指导性替换策略来实现高效缓存。基于此高鸣宇研究组提出了 SeaCache 方案，包含三个核心技术：一是纤维打包与拆分（Fiber Packing and Splitting），将变长的稀疏数据高效映射到定长的缓存块中；二是指导性 LRU（gLFU）策略，利用虚拟标签（Virtual tags）大幅降低元数据开销并保持近似最优的替换效果；三是两阶段自适应机制，动态调节元数据与数据的缓存占比。实验表明，SeaCache 相比现有稀疏缓存设计（如 SpArch、InnerSP 等）平均实现了 2.8 倍的加速，相比高度优化的 Scratchpad 设计也有 2.1 倍的性能提升。

该研究成果论文：Xintong Li, Jinchen Jiang, Mingyu Gao, "SeaCache: Efficient and Adaptive Caching for Sparse Accelerators", MICRO 2025.

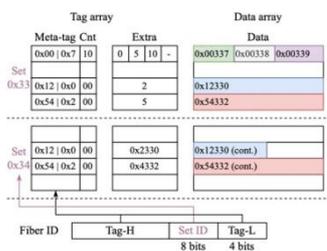


图 1 SeaCache 的纤维打包与拆分映射机制

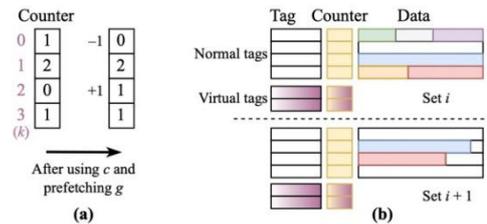


图 2 集成虚拟标签的指导性 LRU (gLFU) 替换策略

二、数据库管理系统

主要完成人：张焕晨研究组

面向长文本理解的高效灵活检索增强框架

近年来，在企业文档分析、金融报告解读等实际应用中，基于长文本数据的问答理解已成为普遍且极具挑战性的任务。然而，直接将超长文本输入商用大语言模型（如 GPT-4o）会产生高昂的 Token 成本。尽管现有的检索增强生成（RAG）和上下文压缩技术试图缓解这一问题，但它们往往面临弊端：因过度压缩丢失关键细节，或者因频繁的迭代调用引入较大的延迟和开销。张焕晨研究组深入观察发现，现实世界中的问答场景呈现出高度的多样性和动态性：从简单的证据检索到复杂的多步推理，不同任务对处理策略和上下文的需求截然不同。这种差异性导致现有的静态或单一策略难以在保证精度的同时兼顾成本效率。

在该工作中，张焕晨研究组基于数据库系统的模块化调度优化的思想，提出了 OkraLong 检索增强框架，支持实时分析并灵活优化执行流程。OkraLong 构建了自适应系统的运作模式，通过分析器 (Analyzer)、组织器 (Organizer) 和执行器 (Executor) 三个协同组件实现对任务的精细化调度。具体而言，分析器使用监督微调的轻量级模型，从查询类型、信息模式和证据充分性三个维度实时分析任务状态；组织器基于分析结果，通过启发式规则动态构建执行管线并配置检索参数；执行器则集成了基本 RAG 操作及问题拆解、上下文处理、多步推理等高级算子以最终执行任务。实验表明，OkraLong 通过提供灵活自适应的运行模式，成功构建了成本与准确性的帕累托最优边界：在获得 5.7% - 41.2% 准确度提升的同时，获得了 1.3 倍至 4.7 倍的成本效率优化。

该研究成果论文：Yulong Hui, Yihao Liu, Yao Lu, Huanchen Zhang, "OkraLong: A Flexible Retrieval-Augmented Framework for Long-Text Question Answering", EMNLP 2025 (Findings).

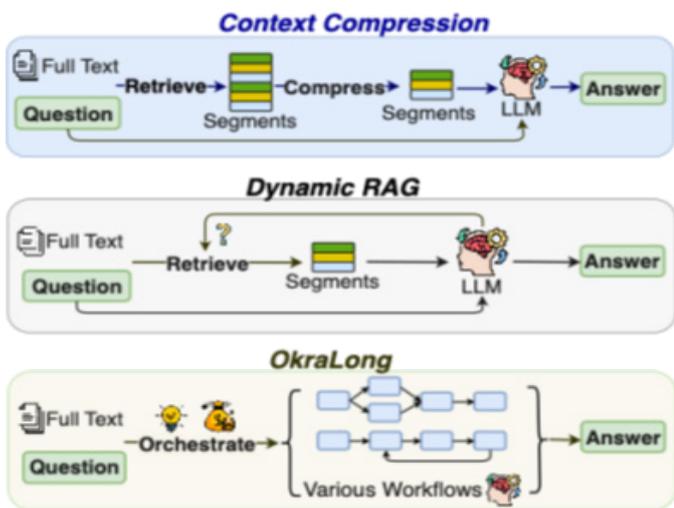


图 1 OkraLong 与其他方法的比较

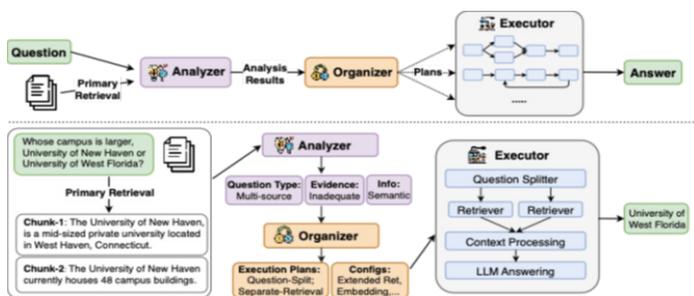


图 2 OkraLong 模块化系统架构设计

F3: 面向未来的开源数据文件格式

列式存储格式是现代数据分析系统的基础。开源文件格式（如 Parquet、ORC）的广泛使用，使得不同平台之间能够无缝共享数据。然而，这些格式是在十多年前为与当今截然不同的硬件和工作负载环境设计的。尽管它们已通过规范更新在一定程度上适应了这些变化，但并非所有部署都支持这些修改，而且系统往往难以在不重写的情况下克服这些格式固有的缺陷与局限。

在该工作中张焕晨研究组提出了“面向未来的文件格式”（Future-proof File Format，简称 F3）项目。F3 是一种新一代开源文件格式，其核心设计理念是互操作性、可扩展性与高效性。F3 通过提供一种通用的数据组织结构和通用 API，使开发者能够轻松添加新的编码方案，从而避免了每次数据处理或计算环境发生变化时都需创建全新格式的问题。每个自描述的 F3 文件不仅包含数据和元数据，还内嵌了用于解码数据的 WebAssembly（Wasm）二进制模块。这种将解码器嵌入文件的方式仅需极小的存储开销（千字节级），并在原生解码器不可用时，仍能确保在任意平台上正确解码。张焕晨研究组通过与传统及当前最先进的开源文件格式进行对比评估，结果表明 F3 的存储布局高效，且 Wasm 驱动的解码机制具有显著优势。

该研究成果论文：Xinyu Zeng, Ruijun Meng, Martin Prammer, Wes McKinney, Jignesh M. Patel, Andrew Pavlo, Huanchen Zhang, "F3: The Open-Source Data File Format for the Future", SIGMOD 2026.

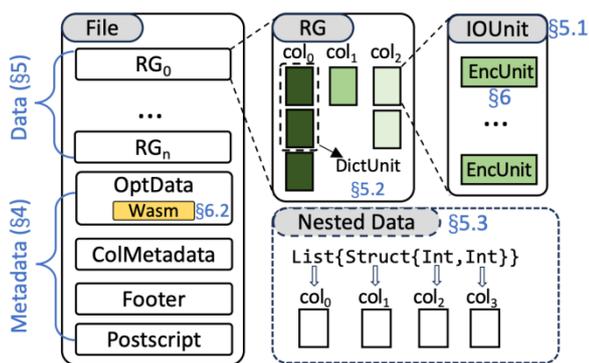


图 1 F3 架构设计图

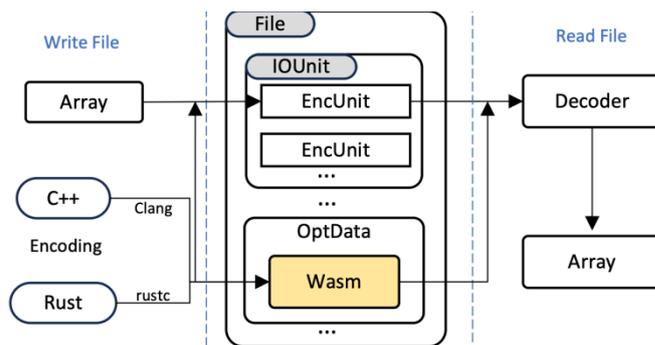


图 2 Wasm 内嵌机制示意图

三、区块链

主要完成人：房智轩研究组

具有两跳接收方公平性的高效公平排序协议

交易的公平排序是目前区块链系统研究的热点问题。既往研究提出了区块链系统中的交易顺序公平性来对交易的实际顺序施加了额外的约束，防止对手通过操纵交易顺序来获取不正当的优势。尽管已经提出了众多公平性的概念以及相应的公平排序协议，但这些努力通常侧重于为每个特定的公平性概念设计专门的协议。因此，目前尚不清楚是否可以在单一框架内同时实现多种公平性概念。

房智轩研究组提出了“两跳接收方公平性” (2-hop receiver fairness)，这是一种新的统一公平性概念，同时涵盖了近似发送方公平性、区块顺序公平性和后果交易公平性。为了实现两跳接收方公平性，该研究组提出了 TxSort，这是一个通用框架，将公平排序问题简化为异步公共子集 (ACS)，这是一种广泛研究的用于异步共识和多方计算的原语。这种方式可以整合最先进的 ACS 协议，以提供一种通信高效且轮次最优的公平排序协议。具体而言，TxSort 协议实现了每笔交易的平均通信复杂度 (CCpT) 为 $O(n^2)$ ，而现有的最先进的异步公平排序协议实现了 $O(n^3)$ 的 CCpT。该研究组进一步提出了 k-TxSort 在高概率满足顺序公平性的同时将 CCpT 压缩到 $O(n)$ 。作为该框架的另一关键组成部分，该研究组引入了一种新颖且计算高效的本地算法，称为“枢轴快速排序”，这可能具有独立的研究价值。

该成果研究论文：Jingfan Yu, Sisi Duan, and Zhixuan Fang, "Efficient Fair Ordering Protocol with 2-hop Receiver Fairness", SRDS 2025.

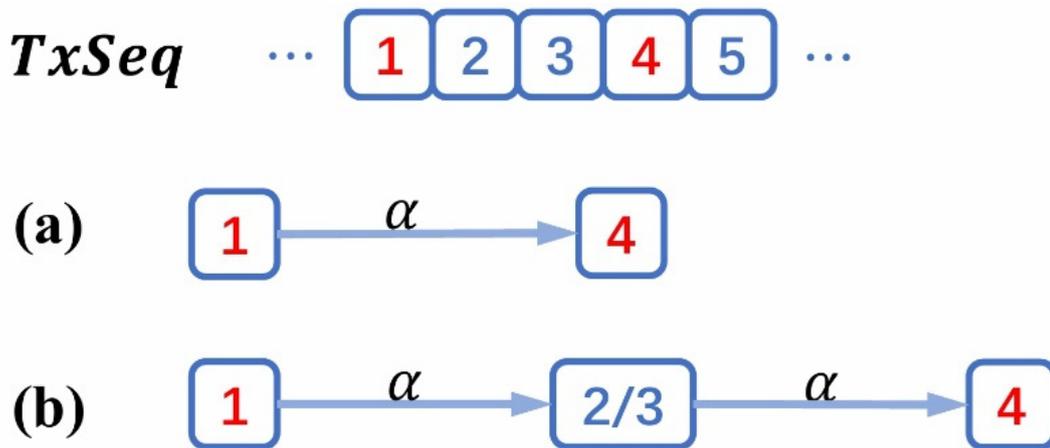


图 1 两跳接收方公平性示意图

四、计算机安全

主要完成人：马恺声研究组

基于 RNS-CKKS 算法的 CNN 网络推理中明文数据的压缩

许多关键领域（如金融风控、医疗诊断、智能安防等）对卷积神经网络（CNN）推理的安全性及效率有双重需求，而基于 RNS-CKKS 的加密推理方案虽能保障用户数据隐私与云端模型机密性，却存在严重的权重明文膨胀问题。原始模型权重经格式转换后存储量激增千倍以上，导致 A100 GPGPU 等硬件因内存限制无法运行大型模型，现有压缩技术又面临计算开销高、压缩率不足等瓶颈，难以兼顾实用性与效果。

马恺声研究组针对这一痛点，提出权重明文压缩方法 WPC，核心突破在于发现 RNS-CKKS 格式转换的类离散傅里叶变换特性，进而提出并证明“周期性传递定理”——权重数据的周期性模式可传递至权重明文，为高效压缩提供理论支撑。基于该定理，该研究组设计通道最内层打包方案与旋转填充技术，通过重构数据维度、补齐非周期部分，使所有权重数据呈现周期性，从而实现高效压缩。

实验验证显示，WPC 在 ResNet 等模型及多数据集上表现优异，成功突破硬件内存限制，可将完整模型部署在商用 GPU 上。该方法还可拓展至其他场景，在大型模型加密推理任务中展现出显著的实用价值与优越性。

该成果研究论文：Guiming Shi, Yuchen Wei, Shengyu Fan, Xianglong Deng, Liang Kong, Xianbin Li, Jingwei Cai, Shuwen Deng, Mingzhe Zhang, Kaisheng Ma, "WPC: Weight Plaintext Compression for CNN Inference based on RNS-CKKS", ACM CCS 2025.



图 1 周期性传递定理

五、密码学

主要完成人：陈一镭研究组、宋一凡研究组

不经意可编程方程的构造及其应用

格问题以其在密码学中的应用广泛闻名。在与南加州大学博士生毛昕渝的合作中，陈一镭研究组提出了一个名为不经意可编程方程（Obliviously Programmable Function, OPF）的密码学新概念，并用它来构造以下两个类似 random-oracle 的密码学原语：

- 通用计算随机提取器（Universal Computational Extractor, UCE）：由 Bellare、Hoang 和 Keelveedhi [BHK13] 引入的通用计算随机提取器可以在各种应用中安全地替换随机预言机，包括 KDM-安全的加密、确定性加密、RSA-OAEP、通用硬核位（universal hard-core bits）等。

- 带辅助输入的多位点混淆（MB-AIPO）。它允许将 CPA-安全的公钥加密系统（PKE）升级为 CCA 安全的加密系统，并可以被用作实现 Fujisaki-Okamoto 转换中用于 PKE 的随机预言机的工具。

尽管 UCE 和 MB-AIPO 很有用，但在标准模型中构建 UCE 和 MB-AIPO 具有挑战性，在之前一直是很困难的公开问题。之前相关的工作都使用不可区分性混淆（iO）加上带辅助输入的点函数混淆，但是使用 iO 会让其构造只存在于理论中，无法高效实现。

OPF 可以替代 UCE 和 MB-AIPO 结构中对于 iO 的使用。陈一镭研究组使用 OPF 和 AIPO 来构建 UCE，MB-AIPO，以及抗击泄漏的 CPA 安全的公钥加密方案。

然后，陈一镭研究组基于格问题构建 OPF，而无需使用 iO。陈一镭研究组在以下假设下给出了上述三种基元的新结构：（1）具有亚指数难度的 LWE；（2）某种特定的采样 LWE 假设；（3）存在能在 NC_1 中计算的 AIPO。

该成果研究论文：Yilei Chen, Xinyu Mao, "Universal Computational Extractors and Multi-Bit AIPO from Lattice Assumptios", EUROCRYPT 2025.

具有最优恶意容忍度和线性通讯复杂度的异步多方安全计算协议

安全多方计算 (MPC) 旨在使多个参与方在不泄露各自私有输入的前提下协同计算一个公开函数, 是现代密码学和分布式安全计算的重要研究方向。在实际分布式系统中, 网络往往呈现异步特性, 即消息传输延迟不可预测, 这使得异步安全多方计算的设计在安全性和效率上都面临更大挑战。理论上已知, 在异步网络中实现恶意安全的 MPC, 需要参与方数量 n 至少是敌手控制参与方数量 t 的 3 倍以上。

在此背景下, 该研究系统考察了在最优恶意容忍度条件, 即 $n=3t+1$ 情况下异步 MPC 的轮复杂度与通信复杂度问题。此前, 在不依赖全同态加密的前提下, 已有的常数轮异步 MPC 构造通信复杂度较高, 难以满足大规模计算需求; 而能够达到线性通信复杂度的方案则往往需要非恒定轮数, 限制了其实用性。

该研究首次提出了一种常数轮异步安全多方计算协议, 实现了线性于电路规模和参与方数量的通信复杂度。具体而言, 该方案仅在随机预言机模型下即可在常数轮内完成安全计算, 其总通信量为 $O(|C| \cdot n \cdot \kappa)$, 其中 $|C|$ 为被计算函数对应电路的规模, κ 为安全参数。在技术上, 论文提出了一种新的方法, 将 MPC-in-the-Head 框架成功适配到异步网络环境中, 并首次实现了常数规模混淆电路在异步设置下的高效计算。

该成果研究论文: Li, J., Song, Y. "Constant-Round Asynchronous MPC with Optimal Resilience and Linear Communication", CRYPTO 2025.

通过轮折叠从最弱假设构建可扩展的常数轮多方安全计算

安全多方计算 (MPC) 允许 n 个参与方在不泄露各自私有输入的前提下联合计算函数结果，是隐私保护计算和分布式安全系统中的核心基础技术。在大规模分布式环境中，网络延迟和参与方数量的增长对协议效率提出了更高要求，因此，同时具备常数轮复杂度与良好通信可扩展性（即通信随参与方数量增长较少）的 MPC 协议一直是该领域的重要研究目标。

现有研究中，基于秘密共享的 MPC 协议在通信复杂度上具有良好可扩展性，但其轮复杂度通常随电路深度线性增长；而基于混淆电路的方案虽然可以实现常数轮数，却长期面临通信复杂度随参与方数量平方增长的问题。近年来，尽管已有工作在特定对手模型下突破了这一瓶颈，但往往依赖较强的密码学假设（如 LPN、DDH 或全同态加密），计算与实现成本较高。

针对上述问题，该研究首次提出了一种新的轮折叠框架，从一个非恒定轮数的 MPC 协议出发，结合混淆电路技术，将其编译为常数轮协议，从而在保持安全性的同时显著降低轮复杂度。该研究在多个对手模型下取得了统一且具有代表性的结果。在强诚实多数情形下（如至多 $t=n/4$ 个恶意参与方），作者构造了一个 5 轮的恶意安全 MPC 协议，其通信复杂度为 $O(|C| \kappa + Dn^2\kappa)$ ，其中 $|C|$ 为计算电路规模， D 为电路深度， κ 为安全参数，且其仅依赖随机预言机假设。进一步地，结合虚拟参与方与 MPC-in-the-Head 技术，该方案可扩展至非诚实多数模型，在任意小于 1 的常数比例参与方被攻陷的情况下，仍能以类似的通信复杂度实现常数轮安全多方计算。

该工作的一大技术贡献在于避免直接依赖多方混淆电路框架，而是让每个参与方对其本地计算进行混淆，并通过轮折叠机制将多轮交互压缩为常数轮执行。这一思路在理论上成功融合了秘密共享与混淆电路两类 MPC 技术的优势，为构造高效、可扩展的常数轮 MPC 协议提供了新的方法论。

该成果研究论文：Goyal, V., Li, J., Ostrovsky, R., Song, Y. "Towards Building Scalable Constant-Round MPC from Minimal Assumptions via Round Collapsing", CRYPTO 2025.

基于单向函数的诚实多数常数轮安全多方计算协议

安全多方计算 (MPC) 是现代密码学的重要研究方向, 其目标是在不泄露各方私有输入的前提下, 实现对任意函数的联合计算。在实际分布式系统中, 协议的轮复杂度与通信复杂度直接决定了其可部署性。然而, 长期以来, MPC 领域面临着一个根本性的效率权衡: 通信高效的协议往往需要与电路深度成线性关系的轮数, 而能够在常数轮内完成计算的协议则通常付出较高的通信代价。

在诚实多数模型 (即被攻陷的参与方数量 t 小于总参与方数量 n 的一半) 下, 已有非恒定轮数的 MPC 协议实现了线性的通信复杂度和无条件安全性。然而, 在常数轮这一更强约束下, 已有结果要么依赖随机预言机、LPN、DDH 等较强假设, 要么在仅假设单向函数的条件下通信复杂度高达 $\Omega(|C| n^2 \kappa)$, 难以扩展至大规模参与方场景。

针对这一长期存在的理论空白, 该工作给出了首个仅基于单向函数假设的、同时满足常数轮复杂度与线性通信复杂度的诚实多数 MPC 构造。作者证明: 在至多 $(n-1)/2$ 个参与方被静态恶意攻陷的情况下, 可以在常数 (26) 轮内安全地计算规模为 $|C|$ 的布尔电路, 总通信复杂度为 $O(|C| n \kappa + n^2 \kappa^2 + n^4 \kappa)$, 其中 κ 为安全参数。在半诚实安全模型下, 轮数可进一步降低至 10 轮。

在技术上, 该工作系统性地改进了经典的多方混淆电路 (BMR) 框架, 利用虚拟参与方技术和打包秘密共享, 成功消除了以往多方混淆中每个门需要 $\Omega(n^2)$ 通信的瓶颈。同时, 论文提出了一套在仅依赖单向函数的前提下验证虚拟参与方本地计算正确性的机制, 从而实现对恶意行为的有效防护。该研究首次在诚实多数、常数轮、线性通信复杂度与最小密码学假设之间实现了兼顾, 为构建可扩展、低交互的隐私计算系统提供了坚实的理论基础。

该成果研究论文: Li, J., Song, Y. "Honest Majority Constant-Round MPC with Linear Communication from One-Way Functions", TCC 2025.

基于轻量级密码的可扩展公平异步安全多方计算

多方安全计算 (Multi-Party Computation, MPC) 旨在使多个互不信任的参与方在不泄露各自私有输入的前提下, 共同完成一项计算任务。自 Ben-Or、Canetti 和 Goldreich (STOC' 93) 以及 Ben-Or、Kelmer 和 Rabin (PODC' 94) 奠定异步网络环境下多方安全计算的理论基础以来, 如何在异步网络中同时实现高安全性与高效率一直是该领域的核心难题。

在异步网络模型下, 恶意参与方可以通过任意延迟或拒绝发送消息来干扰协议执行, 使得在协议设计时难以同时保持安全性和高效性。尤其是在公平性 (Fairness) 要求下, 协议需保证: 要么所有诚实参与方最终均获得输出, 要么所有诚实参与方均无法获得输出。已有工作中, 实现公平异步 MPC 往往依赖重型密码学假设 (如复杂的公钥基础设施), 或者引入高昂的通信与计算开销, 从而严重限制了协议的可扩展性和实际适用性。此前最好的工作可以在随机语言机模型下实现与电路大小成线性的通讯复杂度, 具体每个乘法门的通讯开销约 $100n$ 个域元素, 其中 n 为参与方数量。

该工作提出了 Velox, 在随机预言机模型下, 构建了一个更高效可扩展的公平异步安全多方计算协议。该协议将每个乘法门的通讯复杂度降低到了仅约 $9n$, 提升了十倍效率。实验结果表明, 该协议可以支持最多在 $n=112$ 个人的情况下在 1 秒内进行端到端的高效匿名广播计算, 具备实用价值。

在技术上, Velox 通过精心设计的异步协议结构, 设计了高效的秘密分享协议、秘密重建协议、电路计算协议, 将公平性机制与协议执行过程紧密结合, 有效防止恶意参与方通过延迟消息来获得不公平优势。协议在保证隐私性和正确性的同时, 实现了良好的通信效率和系统可扩展性。随机预言机模型和伪随机生成器的引入, 使得协议能够在不牺牲安全性的前提下, 大幅简化构造并提升性能。

总体而言, 该工作在不依赖同步假设、仅基于随机预言机和轻量级密码学工具的前提下, 实现了高效且公平的异步多方安全计算, 弥补了现有方案在安全假设与效率之间的不足, 为构建可扩展、可信赖的分布式计算系统提供了新的理论基础和技术路径。

该成果研究论文: Bandarupalli A, Ji X, Kate A, et al. "Velox: Scalable Fair Asynchronous MPC from Lightweight Cryptography", Proceedings of the 2025 ACM SIGSAC Conference on Computer and Communications Security. 2025: 1799-1813.

具有输出可达性保证的高效异步 MPC 协议

安全多方计算 (Multi-Party Computation, MPC) 允许多个互不信任的参与方在不泄露各自私有输入的前提下, 共同完成一项计算任务。相比同步网络模型, 异步网络模型不依赖同步时钟和网络延迟假设, 更贴近真实分布式系统环境, 但其协议设计也面临更为严峻的挑战。尤其是在输出可达性 (Guaranteed Output Delivery, GOD) 要求下, 协议需保证只要诚实参与方数量满足安全条件, 所有诚实参与方最终必然获得输出, 这被认为是异步 MPC 中最强、也最难实现的安全性质之一。

自 Ben-Or、Canetti 和 Goldreich (STOC' 93) 以及 Ben-Or、Kelmer 和 Rabin (PODC' 94) 首次证明异步多方安全计算 (MPC) 的可行性以来, 大量研究致力于提升协议的通信效率。一般而言, 异步 MPC 协议的总体通信复杂度可分解为两部分: 一部分与电路规模相关的主项, 另一部分为与电路规模无关的加性开销项。然而, 在同时满足异步网络、最优容错阈值以及输出可达性 (Guaranteed Output Delivery, GOD) 的前提下, 通信效率长期受制于主项与加性开销的巨大成本。

具体而言, 在现有最优结果中, Goyal、Liu-Zhang 和 Song [Crypto' 24] 构造了首个在该安全设置下实现线性渐近通信复杂度的异步 MPC 协议, 即每个电路门的通信量仅随参与方数量呈线性增长。然而, 该协议的加性开销高达 $O(n^{14})$, 其中 n 为参与方数量。尽管这一结果在理论上实现了线性通信形式, 但极其庞大的加性开销使其在实际应用中几乎不可行。

该工作在随机预言机模型下提出了一种计算高效的异步安全多方计算协议, 首次在保证输出可达性的同时, 实现了线性通信复杂度与低多项式 ($O(n^4)$) 加性开销的统一。相比此前工作, 该协议显著降低了与电路规模无关的加性开销项, 将原本不可避免的超高次多项式成本压缩至实际可接受的多项式级别, 从而实质性提升了异步 MPC 的可部署性。

在技术层面, 该工作通过重新设计异步环境下的关键密码学构件 (包括秘密共享及其相关子协议), 有效分离了异步协调所需的通信成本与电路计算本身的通信成本。与既有方案相比, 新协议在保持最优安全阈值和输出可达性保证的同时, 大幅减少了与参与方数量无关的额外通信负担, 使得“线性通信的异步 MPC”首次从纯理论结果迈向具有实际可行性的协议设计。

总体而言, 该工作系统地突破了异步 MPC 中长期存在的通信复杂度瓶颈, 在随机预言机模型下实现了同时具备输出可达性、线性通信和低加性开销的安全多方计算协议。该结果不仅在理论上显著缩小了异步与同步 MPC 在通信效率上的差距, 也为在真实异步网络环境中部署高效、安全的多方计算系统奠定了坚实基础。

该成果研究论文: Bandarupalli A, Ji X, Kate A, et al. "Computationally efficient asynchronous MPC with linear communication and low additive overhead", Annual International Cryptology Conference. Cham: Springer Nature Switzerland, 2025: 261-294.

量子信息

一、离子量子计算、量子网络

主要完成人：段路明研究组、吴宇恺研究组

实现全功能的无串扰双类型离子阱量子网络节点

量子网络是实现量子通讯和分布式量子计算的基础，而离子阱平台以其长相干时间、高确定性且高保真度的量子门操作，成为量子网络研究的主流系统之一。一个功能完整的离子阱量子网络节点通常需要两种类型的量子比特，其中通讯比特用于高效生成离子与光子的纠缠，而存储比特用于长时间保存量子信息，且两者之间需支持高保真度的量子纠缠门。为解决两种比特之间的串扰问题，此前的研究通常采用两种不同种类的离子，由于离子质量不同而面临规模化的挑战。该工作中，研究人员采用同种离子 ($^{171}\text{Yb}^+$) 的不同能级分别编码通讯比特和存储比特，利用 411 nm 和 3432 nm 的双色激光实现了两种量子比特间微秒量级的相干转换，并首次在同一个双类型量子网络节点中集成了无串扰的离子 - 光子纠缠生成、长时间量子存储和离子比特之间量子纠缠门三大关键功能。

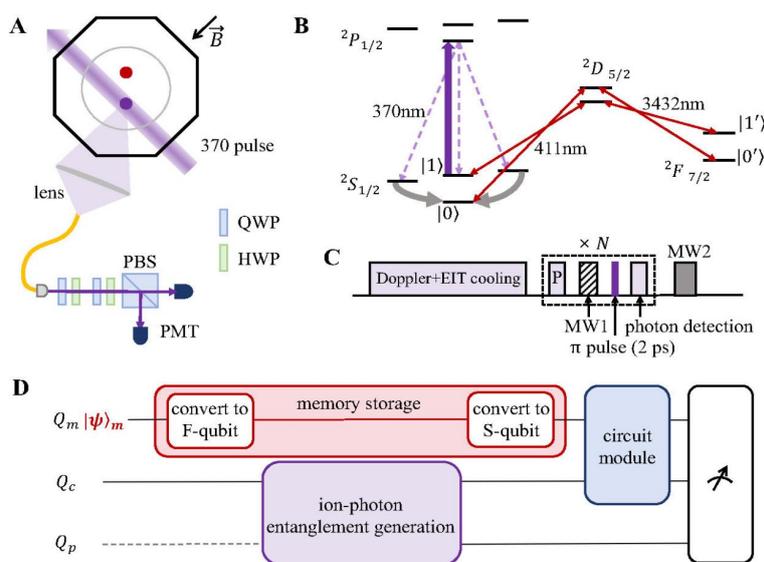


图 1 实验方案示意图

基于此全功能节点，研究人员进一步演示了量子隐形传态和多体 GHZ 量子纠缠态制备两种量子网络的关键应用，验证了该量子网络节点的基本功能及其在量子通信和分布式量子计算中的应用前景。

该成果研究论文：Yuan-Yuan Huang, Lu Feng, Yu-Kai Wu, Yu-Lin Xu, Lin Zhang, Zhai-Bin Cui, Chuan-Xin Huang, Chi Zhang, Shi-An Guo, Quan-Xin Mei, Bin-Xiang Qi, Yong Xu, Yun-Fei Pu, Zi-Chao Zhou, and Lu-Ming Duan, "Realization of a functioning dual-type trapped-ion quantum network node", Science Advances, 11, eaeb4076.

二、离子阱量子中继

主要完成人：濮云飞研究组、段路明研究组

实现城市级离子阱量子网络节点：混合多路复用显著提升远程纠缠分发能力

构建可扩展且具备工程可行性的量子网络，是量子信息科学从基础研究迈向实际应用的关键一步，对量子通信、分布式量子计算以及未来量子互联网的发展具有基础性和战略性意义。然而，在现实光纤信道条件下，远程量子纠缠生成概率随传输距离呈指数衰减，而量子节点中有限的相干时间进一步限制了纠缠分发的有效速率，使得城市尺度量子网络在实验实现和系统集成层面长期面临根本性挑战。

围绕这一核心瓶颈，濮云飞、段路明研究组在离子阱量子网络体系中开展了深入而系统的研究，实验实现了一种基于混合多路复用 (hybrid multiplexing) 的高性能离子阱量子网络节点架构。该工作在单个线性保罗阱量子节点内，协同引入多重光子激发序列与离子在阱内的快速、可控移动操作，在保持离子 - 光子纠缠质量的同时，实现了大规模并行的纠缠生成过程，从系统层面显著提升了远程纠缠的有效尝试速率。

在实验验证中，该研究组基于一个由 5 个离子构成的量子网络节点，实现了最多 44 个相互独立的时间戳 (time-bin) 模式的离子 - 光子纠缠生成，并在 3 m、1 km 及 12 km 光纤信道条件下对节点性能进行了系统表征。实验结果表明，该量子节点在不同传输距离条件下均可实现高保真度的离子 - 光子纠缠分发，充分展示了该混合多路复用方案在真实通信环境中的稳定性、可重复性与可扩展性，使离子阱量子网络节点的性能指标首次进入城市级量子网络应用所要求的参数区间。

针对多路复用纠缠生成过程中由频繁操作引入的潜在退相干问题，该研究组采用了既有的双类型编码 (dual-type encoding) 量子存储方案，在物理实现层面有效区分了纠缠生成与量子存储所涉及的操作自由度，从而显著抑制了多次激发与离子移动对存储量子比特相干性的累积影响。在此基础上，实验实现的量子存储相干时间达到 366 ms，首次在实验上实现了量子存储时间超过 12 km 光纤条件下远程纠缠生成的平均等待时间，从而在系统层面实现了远程纠缠分发速率与量子存储寿命之间的有效匹配。

该研究系统性地揭示了混合多路复用策略在离子阱量子网络中的显著性能优势，不仅为提升远程纠缠分发效率提供了切实可行的技术路径，也为多节点量子网络中的纠缠连接、量子中继及网络级资源调度奠定了坚实的实验基础。相关成果代表了离子阱量子网络在系统集成能力与工程实现成熟度方面的重要进展，为构建高可靠性、可扩展的城市级量子网络体系提供了关键支撑。

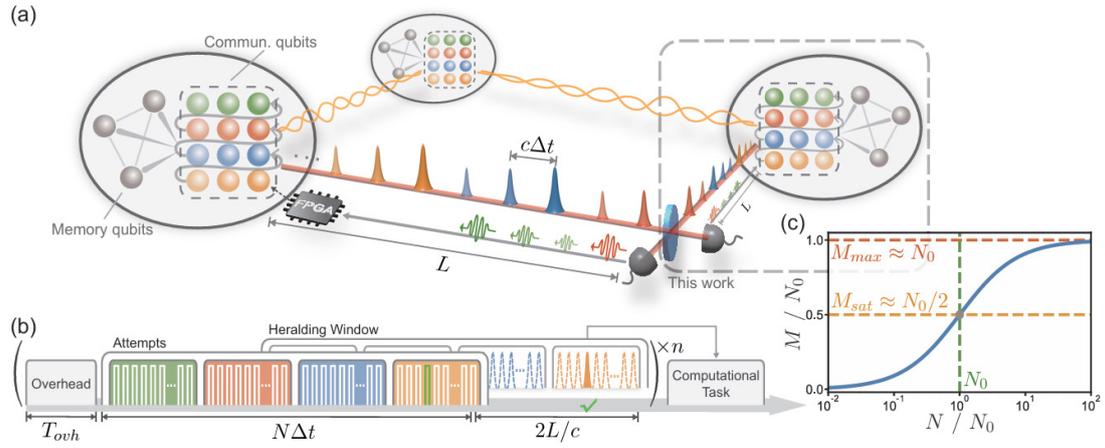


图 1 基于混合多路复用的离子阱量子网络节点及时间戳（time-bin）模式示意图

该成果研究论文：Z.-B. Cui, Z.-Q. Wang, P.-C. Lai, Y. Wang, J.-X. Shi, P.-Y. Liu, Y.-D. Sun, Z.-C. Tian, B.-X. Qi, Y.-Y. Huang, Z.-C. Zhou, Y.-K. Wu, Y. Xu, L.-M. Duan#, Y.-F. Pu, "Metropolitan-scale ion-photon entanglement via a quantum network node with hybrid multiplexing enhancements" , Nature Communications volume 17, Article number: 697 (2026).

三、中性原子量子网络

主要完成人：濮云飞研究组、段路明研究组

首次实现两个多体纠缠态的存储增强的融合

量子多体纠缠态是量子信息领域诸多关键应用的重要资源，高效制备量子多体纠缠态对这些关键应用的落地至关重要。然而，量子多体纠缠态的制备非常困难，如何以模块化、可扩展的方法制备量子多体纠缠态，是当前量子信息领域的一个难题。

濮云飞、段路明研究组与山西大学王海研究组合作，通过存储增强的方法，将两个较小规模的量子多体纠缠态融合成了一个更大规模的多体纠缠态，首次展示了一种高效率制备大规模量子多体纠缠态的方法。研究人员借助中性原子量子存储器，异步制备了两个 3 体 W-state 纠缠态，并通过纠缠交换的方式完成纠缠融合，最终制备了一个更大规模的 4 体 W-state 纠缠态。通过量子存储器的使用，实验实现了纠缠融合的成功率正比于纠缠制备的成功率，相较于没有量子存储时二次方的成功率提升了两个数量级。该工作为未来模块化，高效率的制备量子多体纠缠态提供了一种可行的方案。实验示意图如下：

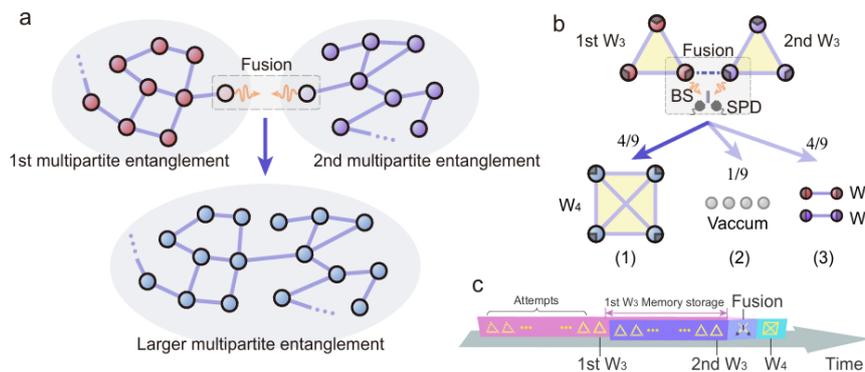


图 1 实验系统与纠缠融合方案

该成果研究论文: Jixuan Shi, Sheng Zhang, Yukai Wu, Yuedong, Sun, Yibo Liang, Hai Wang, Yunfei Pu, Luming Duan, "Scalable and modular generation of multipartite entangled states through memory-enhanced fusion", Physical Review Letters 135, 150802 (2025).

四、金刚石量子模拟与传感

主要完成人：侯攀宇研究组、邓东灵研究组、段路明研究组

首次在固体自旋系综中观测动力学冻结现象并增强量子传感灵敏度

理解与控制量子多体系统中的非平衡动力学，是现代物理学面临的一项重要挑战。通常情况下，量子多体系统在驱动下会迅速热化，丢失初始状态信息，以往突破热化的方式主要依赖于引入无序或采用高频驱动。近年来，理论研究表明，在强周期、中等频率强度驱动下，系统可通过“动力学冻结”机制有效抑制热化，其物理本质源于驱动过程中涌现出新的守恒量，然而该现象此前一直缺乏强力的实验证据。

针对这一问题，研究人员利用金刚石中由约 104 个有强相互作用的氮 - 空位 (NV) 中心自旋组成的系综，通过精心设计的高精度 Floquet 驱动序列，首次在实验中观测到动力学冻结现象 (图 1)。实验结果显示，在驱动参数满足特定条件 ($h_z T = 2\pi \times 2k$) 下，系统自旋磁化量能够长期保持，其持续时间远超体系固有的相互作用限制相干时间 T_2 一个数量级以上。一旦驱动参数偏离该条件，系统则迅速恢复热化行为，体现出该机制对驱动参数的高度敏感性。

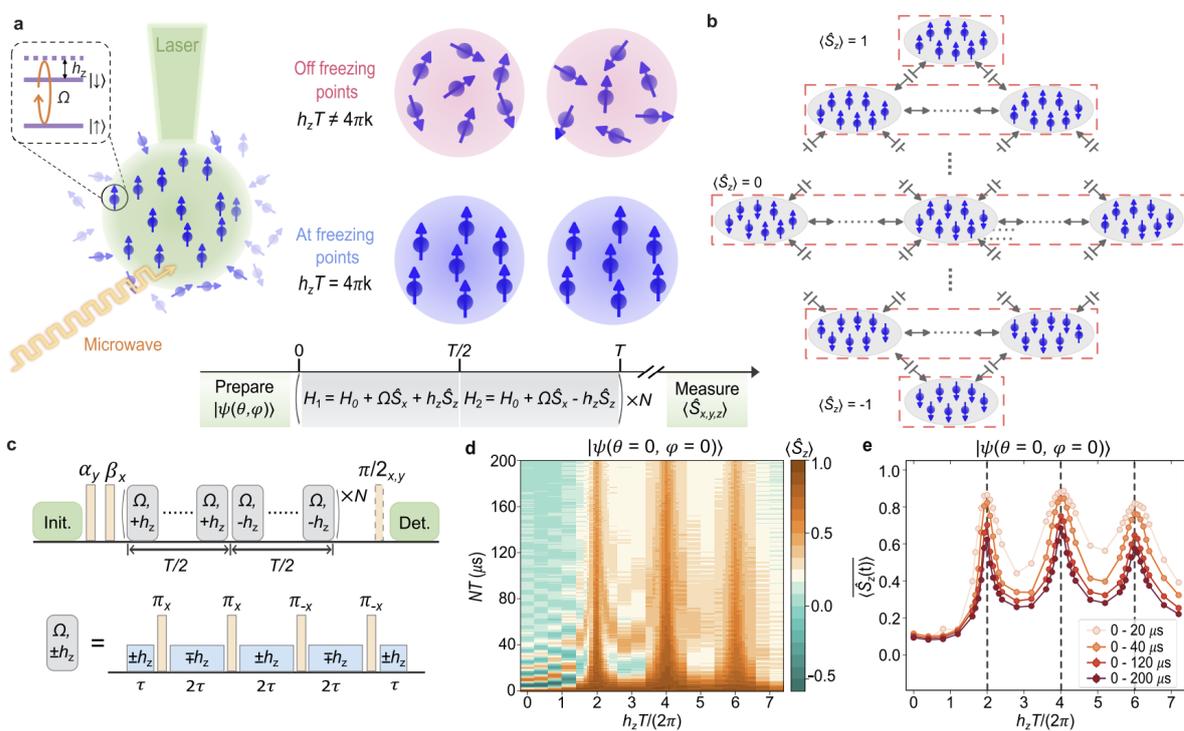


图 1 基于金刚石 NV 色心自旋系综观测到的动力学冻结现象示意图和实验结果

进一步地，研究人员利用动力学冻结对参数敏感的特性，发展了一种新型交流磁力计协议 (图 2)。该方法将冻结效应与动力学解耦技术相结合，有效延长了量子传感的最佳探测时间，突破了传统方案中受限于相干时间的性能瓶颈。实验结果表明，该技术的磁场测量灵敏度较传统动态解耦方法最高提升约 2.7 倍，显著增强了对微弱磁信号的探测能力。

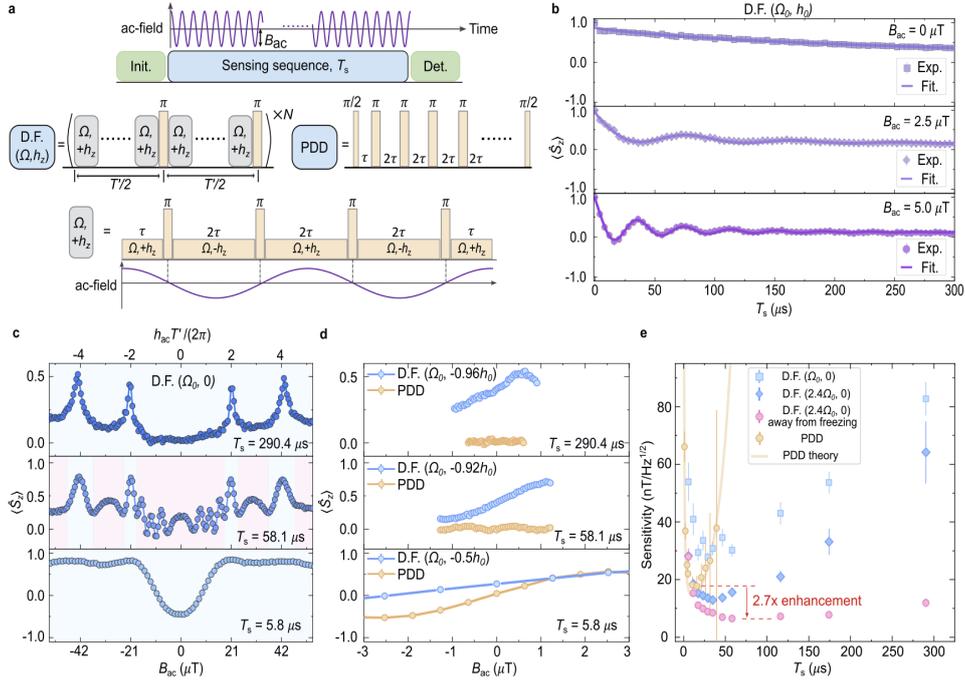


图 2 基于动力学冻结机制增强的量子磁测量方法原理图和实验结果

该成果不仅首次在真实固态自旋体系中验证了动力学冻结的理论预言，揭示了一种基于涌现守恒律的新型热化抑制机制，还为量子精密测量提供了新方法，具有重要的科学意义与应用潜力。该方案仅需全局控制、易于实施，为凝聚态物理、化学及生物医学等领域中高空间分辨率与高灵敏度兼具的量子传感应用，提供了切实可行的技术途径。

该成果研究论文：Ya-Nan Lu, Dong Yuan, Yixuan Ma, Yan-Qing Liu, Si Jiang, Xiang-Qian Meng, Yi-Jie Xu, Xiu-Ying Chang, Chong Zu, Hong-Zheng Zhao, Dong-Ling Deng, Lu-Ming Duan, Pan-Yu Hou, "Dynamical freezing and enhanced magnetometry in an interacting spin ensemble", arXiv:2507.22982.

五、量子信息科学

主要完成人：马雄峰研究组

基于高品质单光子源的多中继量子网络架构实现

量子密钥分发 (QKD) 作为实现信息论安全通信的核心技术, 在构建未来量子网络中具有关键作用。然而, 受限于光子在信道中的固有损耗, 传统点对点 QKD 协议的密钥率受限于线性损耗比例, 即“无中继极限”。为突破这一限制, 近年来涌现出双场 QKD、模式匹配 (MP) QKD 等具有类中继损耗缩放优势的协议, 其中 MP-QKD 因不依赖全局相位参考, 在实际应用中展现出更高可行性。

量子网络是连接未来量子计算、量子密码学等系统的关键基础设施, 其核心挑战在于如何安全、高效地扩展网络规模。传统量子中继方案受限于单个节点, 难以构建复杂网络拓扑; 而依赖经典中继的节点则面临安全性依赖于节点可信度的根本瓶颈。

在该项研究中, 马雄峰及其博士生黄溢智与中国科学技术大学合作, 创新性地提出并实验实现了五节点量子中继架构。该架构将自主研发的高性能量子点单光子源作为核心不可信中继单元, 与两侧的干涉测量节点协同, 共同构成了三个量子中继节点, 在国际上首次突破了量子网络限于单个中继节点的技术瓶颈。通过优化单光子源与相干光的强度配比及探测时间窗, 该研究组成功将干涉可见度提升至 90% 左右, 确保了量子信息在节点间稳定、高效地传递。

为验证该架构的实用性与安全性, 该研究组采用相位匹配量子密钥分发协议, 在总长 300 公里的光纤链路中成功实现了密钥的安全分发。这一成果证实了该多中继架构在长距离量子通信中的可行性, 展现了单光子源作为量子中继核心单元的应用潜力, 为构建覆盖范围更广、安全性更高的大规模三层星型量子网络奠定了坚实基础。

该成果研究论文: Mi Zou, Yu-Ming He, Yizhi Huang, Jun-Yi Zhao, Bin-Chen Li, Yong-Peng Guo, Xing Ding, Mo-Chi Xu, Run-Ze Liu, Geng-Yan Zou, Zhen Ning, Xiang You, Hui Wang, Wen-Xin Pan, Hao-Tao Zhu, Ming-Yang Zheng, Xiu-Ping Xie, Dandan Qin, Xiao Jiang, Yong-Heng Huo, Qiang Zhang, Chao-Yang Lu, Xiongfeng Ma, Teng-Yun Chen & Jian-Wei Pan, "Realization of an untrusted intermediate relay architecture using a quantum dot single-photon source", Nature Physics 21, 1670 - 1677 (2025).

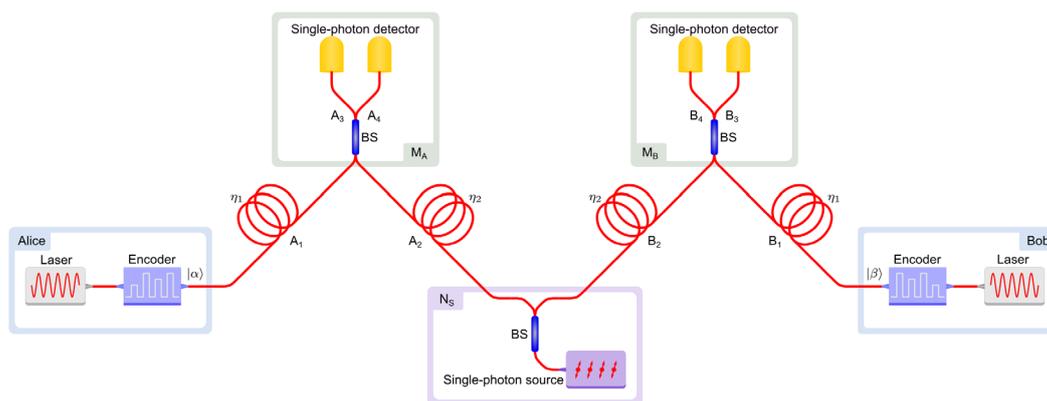


图 1 五节点网络结构示意图

诱骗态量子密钥分发的增强分析与统计涨落处理

量子密钥分发作为量子信息科学中的实用化技术，其现实系统中的核心安全模块——诱骗态方法，虽能有效防御光子数分离攻击，却在有限数据尺寸下因统计涨落而面临密钥率估计困难，特别是相位错误率与密钥率之间的非线性依赖关系使问题进一步复杂化。

在该项研究中，马雄峰及其学生徐子泰、黄溢智提出了针对诱骗态方法的密钥率下界增强估计方案与改进的统计涨落分析框架。该研究组通过引入香农熵函数的线性展开，将非线性密钥率公式转化为线性优化问题，不仅简化了理论分析，还使得联合统计涨落处理成为可能。基于此线性框架，研究进一步推导出单诱骗态协议在有限数据情况下的密钥率安全下界，并利用 Chernoff 边界对涨落进行严格刻画。

数值模拟结果表明：在典型实验参数与数据量 ($N=1011$) 下，本方法在 250 公里传输距离处的密钥生成速率分别为传统单诱骗态方法与真空 + 弱诱骗态方法的 1.87 倍与 1.46 倍；在最大传输距离方面，本方法在 $N=1011$ 时较另外两种方法分别提升了 12 公里与 5 公里，显著增强了有限数据场景下的协议性能与传输范围。该线性化涨落分析框架不仅适用于诱骗态 QKD，还可推广至其他具有线性统计关系的量子信息处理任务中。

该成果研究论文：Zitai Xu, Yizhi Huang, Xiongfeng Ma, "Enhanced Analysis for the Decoy-State Method", Adv. Quantum Technol. 2025, 10.1002/qute.202400687.

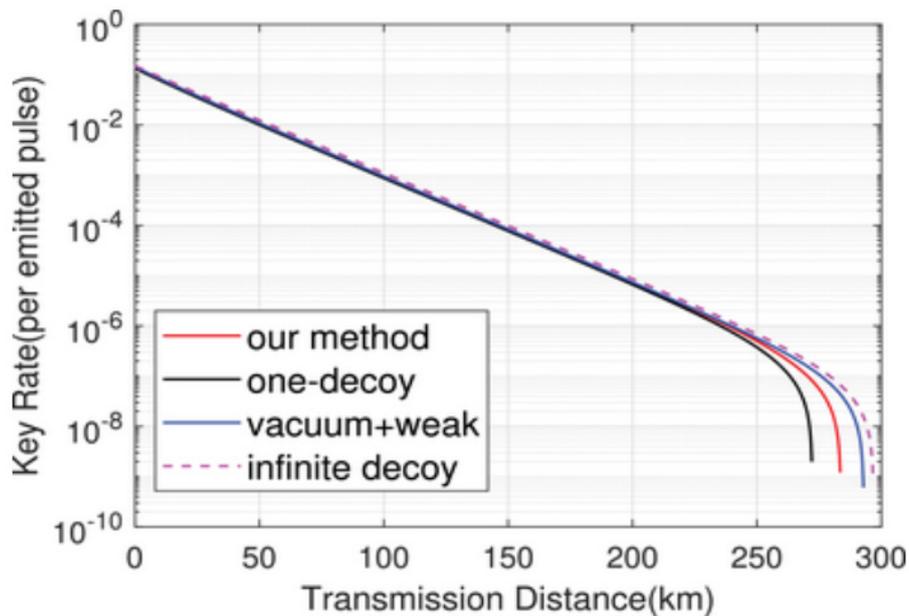


图 2 有限码长下密钥率对比图

量子纠错码的子空间验证

量子信息处理与量子计算在多个领域展现出潜在优势，但其实际实现受到噪声的严重制约。量子纠错是构建大规模量子计算机的关键技术，然而其往往需要引入大量物理量子比特，从而带来显著的资源开销。随着系统规模的增大，编码态保真度的验证也变得愈发困难。传统方法（如直接保真度估计和量子态验证）要么资源消耗巨大，要么仅适用于特定目标态。

为弥补这些不足，马雄峰及其博士生陈俊杰与复旦大学、芝加哥大学、香港大学合作，提出了一种通用且高效的子空间验证框架，并将其与直接保真度估计相结合，可以高效且稳健地估计制备态与其目标码子空间之间的接近程度。该框架在传统假设检验的基础上进一步发展，允许测量失败的存在，从而对测量噪声具有更强的鲁棒性。该研究组利用图论工具，为稳定子码和低密度奇偶校验码 (LDPC 码) 子空间设计了相应的验证策略。与直接保真度估计相结合后，该方法可大幅降低测量成本，并支持对一般“魔态”逻辑态的验证。

该框架和方案的提出为量子纠错码及一般魔态逻辑态提供了一种高效且可行的验证途径，有助于推动其在实际量子平台上的实现。

该成果研究论文：Junjie Chen, Pei Zeng, Qi Zhao, Xiongfeng Ma¹ and You Zhou, "Quantum Subspace Verification for Error Correction Codes", PRX Quantum, vol. 6, p. 040353, 2025.

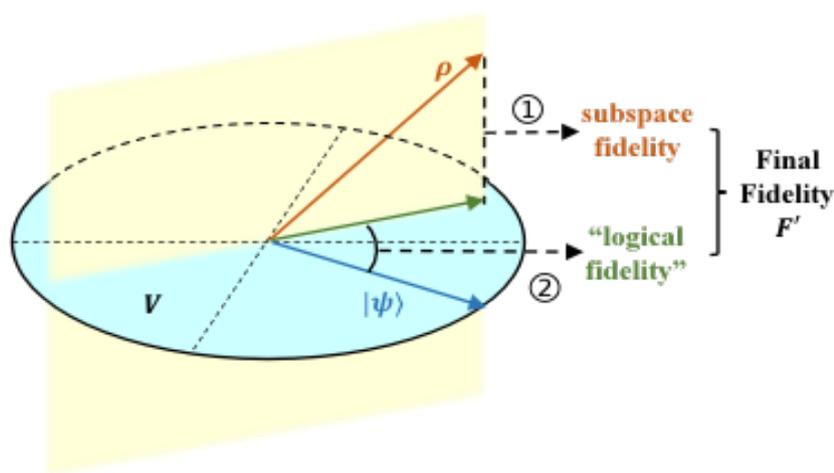


图 3 结合量子子空间验证与直接保真度估计示意图

含噪量子设备的计算与纠缠生成能力存在根本性限制

量子信息处理与量子计算在多个领域展现出潜在优势，但其实际实现受到噪声的严重制约。在当前无法实现完全量子纠错的“含噪中等规模量子计算”阶段，探寻噪声设备所能提供的可靠量子优势成为一个关键而富有挑战性的问题。反之，深刻理解现有技术的根本性局限也同等重要。

针对这一问题，马雄峰及其博士生鄢语轩、杜振宇、陈俊杰，对一类严格收缩噪声进行了深入分析，并通过分解得出了其解析表示。在此类噪声模型下，该研究组观察到，当电路深度超过对系统规模的数量级时，量子设备在多项式时间内将变得与随机抛硬币无法区分。即使辅以经典后处理，在具有超对数深度噪声电路的任意多项式时间算法中，此类设备也无法展现出任何计算优势。这一结论对变分量子算法、误差缓解以及具有多项式深度的量子模拟研究均具有重要影响。此外，该研究组还研究了具有门连接拓扑限制的含噪量子设备。对于一维含噪量子电路，研究排除了在所有深度范围内存在超多项式量子优势的可能性。同时，该工作也为纠缠生成的规模设立了上限。

这些发现清晰地揭示了含噪量子设备在计算能力和纠缠可扩展性方面所面临的根本性约束，为评估当前量子硬件的能力边界、引导算法设计提供了关键的理论依据。

该成果研究论文: Yuxuan Yan, Zhenyu Du, Junjie Chen & Xiongfeng Ma, "Limitations of noisy quantum devices in computing and entangling power", npj Quantum Information, vol. 11, no. 1, p. 188, Nov 2025.

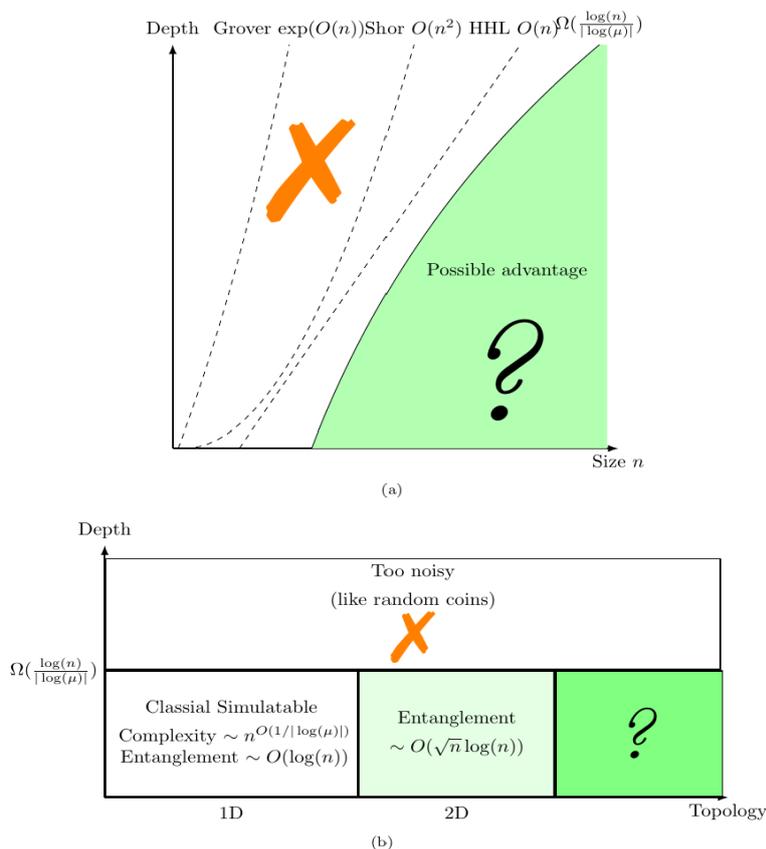


图 4 含噪量子设备的计算和纠缠产生能力限制

六、量子模拟

主要完成人：邓东灵研究组

在百比特量子芯片上实现新奇量子物态

对称性保护的拓扑边缘态是凝聚态物理中一种新奇的物质状态。它们通常出现在系统边界，并受到特定对称性的保护，能够有效地抵抗满足这些对称性的噪声。这一特性使其在量子信息领域具有潜在的应用价值。然而，拓扑边缘态十分脆弱，通常只存在于绝对零度下的系统基态中。在有限温度环境中，大量的热激发会自由传播并与边缘态相互作用，从而破坏边缘态并抹去其中存储的量子信息。因此，在热扰动下寻找并保护量子物态是凝聚态物理和量子信息领域的重要课题。

为了应对这一挑战，主流的方法是引入无序使系统进入多体局域化状态，从而束缚热激发。但该方法高度依赖随机势场，不仅实验成本昂贵，并且其稳定性仍存在争议。邓东灵研究组等另辟蹊径，提出利用预热化机制来保护实验制备的拓扑边缘态。该方法无需引入无序，而是依靠系统内部涌现的对称性，为边缘态提供额外的保护，从而抑制其与热激发之间的相互作用。

为验证这一构想，该研究组与浙江大学超导量子计算研究组合作，在浙江大学自主研制的 125 比特“天目 2 号”超导量子芯片上实现了一条由 100 个粒子组成的一维对称性保护拓扑链。在高度的编程灵活性与国际先进的量子操作保真度支持下，该研究组在约 270 层量子线路演化过程中观察到了不受热激发影响的拓扑边缘态，并深入研究了系统预热化状态下热激发的动力学与涌现的对称性。在此基础上，该研究组进一步利用这种稳健的拓扑边缘态编码并制备了逻辑贝尔态，有力证明了其对热激发的鲁棒性。

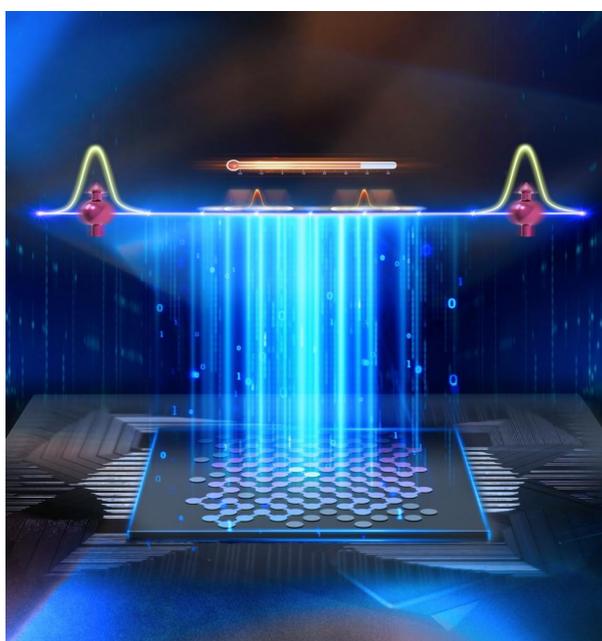


图 1 有限温度拓扑边缘态示意图

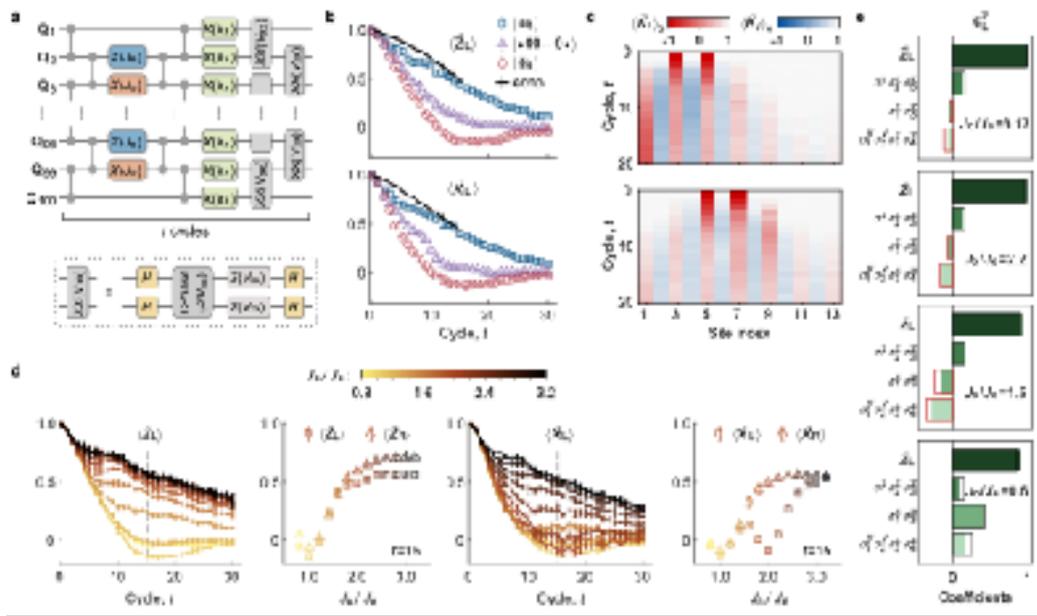


图 2 使用超导量子芯片实现有限温拓扑边缘态的主要实验结果图

该研究建立了一种可行的数字模拟方法，为在有限温度下探索拓扑物质提供了新的实验手段。此外，它还展示了在无序系统中实现长寿命、鲁棒边界量子比特的潜在途径，为构建抗噪声的量子存储与操控技术提供了新道路。

该成果研究论文：Feitong Jin, Si Jiang, Xuhao Zhu, Zechang Bao, Fanhao Shen, Ke Wang, Zitian Zhu, Shibo Xu, Zixuan Song, Jiachen Chen, Ziqi Tan, Yaozu Wu, Chuanyu Zhang, Yu Gao, Ning Wang, Yiren Zou, Aosai Zhang, Tingting Li, Jiarun Zhong, Zhengyi Cui, Yihang Han, Yiyang He, Han Wang, Jia-Nan Yang, Dong-Ling Deng, "Topological prethermal strong zero modes on superconducting processors", Nature volume 645, pages626 – 632 (2025).

七、超导 - 光学 / 声学混合系统量子计算

主要完成人：孙麓岩研究组

基于布里渊集成电路的紧凑型高分辨率光谱仪

光谱仪在化学与生物传感、天文观测以及量子技术等领域具有不可替代的作用，而将高性能光谱仪完全集成到光子芯片上长期面临“分辨率—器件尺寸”之间的根本矛盾。传统基于光栅或干涉仪的片上光谱方案，要么需要较大的器件尺寸和多通道探测阵列才能获得高分辨率，要么受限于有限的调谐长度而难以进一步提升性能。因此，发展一种同时具备高光谱分辨率、宽工作带宽、小尺寸和动态可重构能力的单通道片上光谱新方案，是集成光子学领域的关键挑战之一。

孙麓岩研究组与中科大邹长铃教授研究组合作提出并在实验上实现了一种基于声学激发布里渊散射的片上高分辨率光谱仪。其利用混合集成的光子 - 声子芯片中强光声相互作用，将声波诱导的动态折射率调制等效为可调谐的高反射“动态光栅”。该方案在仅有 1 mm 直波导长度的条件下，实现了 0.56 nm 的本征光谱分辨率和超过 110 nm 的工作带宽，分辨率接近给定器件尺寸下的理论极限。通过射频信号调控声波频率，实现了光谱的快速扫描与重构，并系统验证了该器件在单色光、双峰光谱以及宽带荧光与 DWDM 信号测量中的性能，同时分析了分辨率受声子损耗与相位匹配扩散限制的物理机制。

该工作展示了一种突破传统片上光谱分辨率与尺寸权衡的新思路，为实现紧凑、高分辨率、单通道且可重构的集成光谱系统提供了通用平台。所提出的混合光子 - 声子架构不仅具备良好的稳定性、可扩展性和 CMOS 兼容潜力，还为光谱分析、光信号处理以及射频 - 光学接口提供了新的实现路径。进一步而言，该平台可与其他光子、电子甚至超导量子器件深度集成，在量子信息处理、片上激光源与新型光电器件等方向具有重要应用前景。

该成果研究论文：Jia-Qi Wang, Yuan-Hao Yang, Zheng-Xu Zhu, Juan-Juan Lu, Ming Li, Xiaoxuan Pan, Chuanlong Ma, Lintao Xiao, Bo Zhang, Weiting Wang, Chun-Hua Dong, Xin-Biao Xu, Guang-Can Guo, Luyan Sun & Chang-Ling Zou, "Compact and high-resolution spectrometer via Brillouin integrated circuits", Nature Communications volume 17, Article number: 68 (2026).

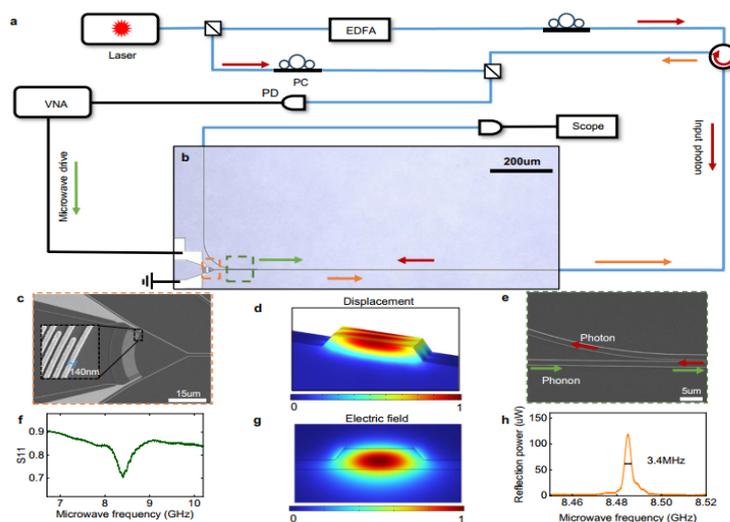


图 1 基于声学激发布里渊散射的片上光谱仪器件结构与实验表征

基于微米尺度非悬空波导阵列的千兆赫兹拓扑声子电路

声波在 GHz 频段具有波长短、能量局域性强等优势，在射频信号处理、片上传感以及新型信息处理架构中展现出重要应用潜力。近年来，拓扑物理为实现对声波传播的稳健调控提供了新的理论框架，拓扑声子边缘态能够在结构缺陷和无序存在的情况下保持稳定传输。然而，已有拓扑声子器件多依赖悬空结构或工作在较低频率范围，难以兼顾高频运行、器件机械稳定性与大规模片上集成，这成为拓扑声子学从概念验证走向实际应用的主要瓶颈。

孙麓岩研究组与中科大邹长铃教授研究组提出并实验实现了一种基于微米尺度非悬空波导阵列的 GHz 拓扑声子电路，通过精确设计声波导的几何参数与相互耦合关系，在芯片上构建具有非平庸拓扑性质的声学能带结构。该器件在约 1.5 GHz 频率下支持拓扑保护的声波边缘态传播，声波能够沿预定义路径实现单向或准单向传输，并对局域缺陷和结构扰动表现出显著鲁棒性。该工作进一步利用扫描光学振动测量技术，对声波在体波导中的空间分布和传播过程进行了直接成像，系统验证了拓扑边缘态的存在及其稳定性。

该工作首次在非悬空、微米尺度的集成声学平台上实现了 GHz 频段的拓扑声子电路，突破了拓扑声子器件在工作频率、结构形态和集成兼容性方面的关键限制。所提出的声子电路方案具有良好的机械稳定性和可扩展潜力，为构建高频、鲁棒的片上声学与射频系统奠定了基础。该成果不仅拓展了拓扑物理在声子系统中的应用范围，也为未来声子信息处理、射频信号路由以及与光子和量子器件的混合集成提供了新的技术平台。

该成果研究论文：Xin-Biao Xu, Mourad Oudich, Yu Zeng, Ji-Zhe Zhang, Yuan-Hao Yang, Jia-Qi Wang, Weiting Wang, Luyan Sun, Guang-Can Guo, Yun Jing & Chang-Ling Zou, "Gigahertz topological phononic circuits based on micrometre-scale unsuspended waveguide arrays", Nature Electronics volume 8, pages689 – 697 (2025).

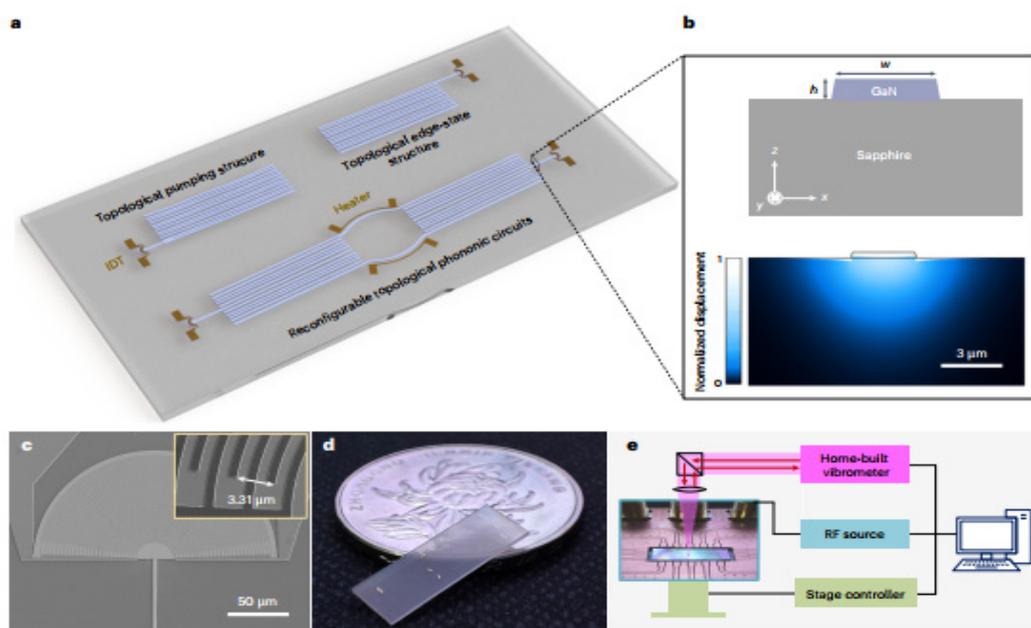


图 1 基于微米尺度非悬空波导阵列的 GHz 拓扑声子电路原理与结构示意图

八、超导量子计算

主要完成人：孙麓岩研究组

超导量子计算的研究进展与挑战

量子计算利用量子叠加和纠缠等特性，在特定问题上展现出超越经典计算的潜力。在多种实现路线中，超导量子计算因其与半导体工艺兼容、可扩展性强，成为当前最有前景的技术方案之一。近年来，该领域在器件性能和系统规模上均取得了快速进展，但距离大规模、实用化仍面临诸多挑战。

孙麓岩研究组、中科大邹长铃研究组、北京量子院于海峰等其他研究组合作，系统梳理了超导量子计算的发展历程与核心技术，包括超导量子比特的器件制备、单比特和双比特量子门实现、多比特纠缠态制备、量子优越性实验、量子纠错与误差缓解技术，以及量子模拟的最新进展。同时，讨论了玻色编码量子比特、fluxonium 等新型量子比特体系及其在模块化和分布式量子计算中的潜力。

该研究全面总结了超导量子计算的实验进展与关键瓶颈，为理解当前技术水平和未来发展方向提供了清晰图景。相关成果为构建高保真、可扩展的量子处理器以及实现容错量子计算和量子网络奠定了重要基础，对推动量子信息科学和相关应用具有重要意义。

该成果研究论文：Yao-Yao Jiang, Chunqing Deng, Heng Fan, Bing-Yang Li, Luyan Sun, Xin-Sheng Tan, Weiting Wang, Guang-Ming Xue, Fei Yan, Hai-Feng Yu, Ying-Shan Zhang, Yu-Ran Zhang, Chang-Ling Zou, "Advancements in superconducting quantum computing", National Science Review, Volume 12, Issue 8, August 2025, nwaf246.

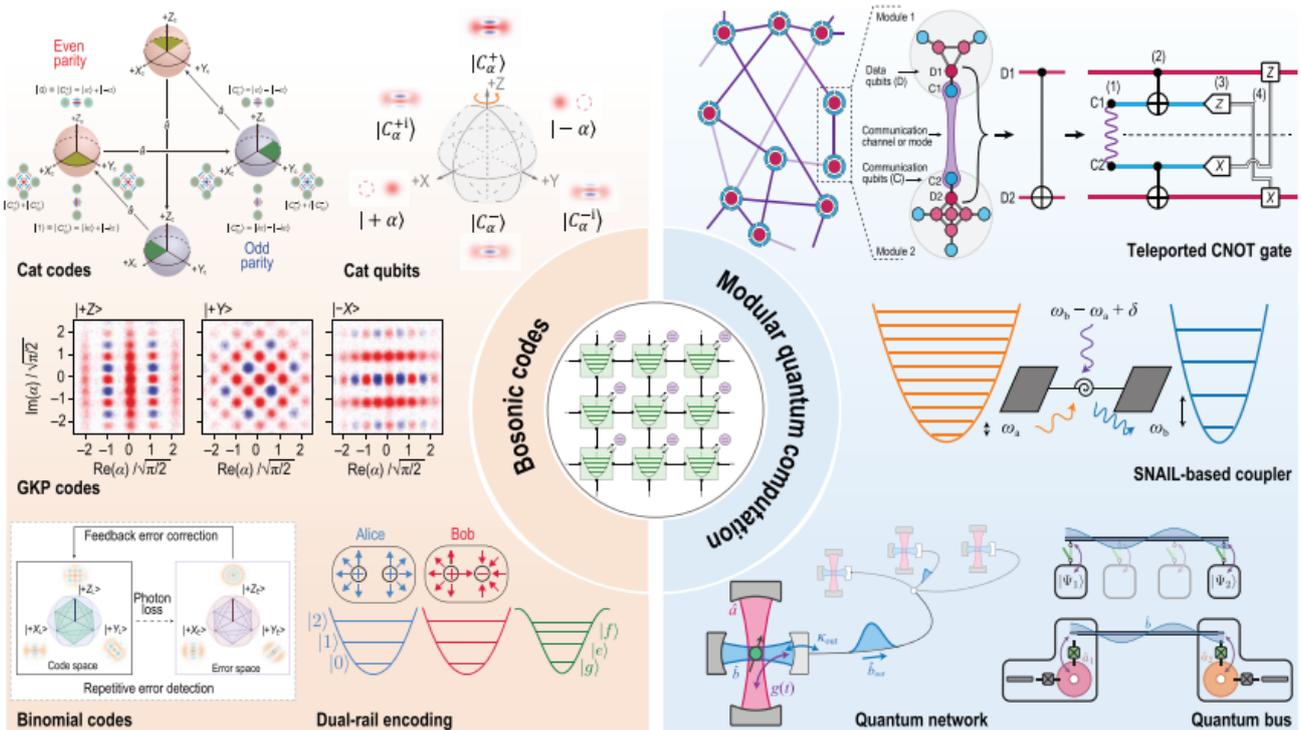


图 1 玻色编码超导量子比特研究进展示意图

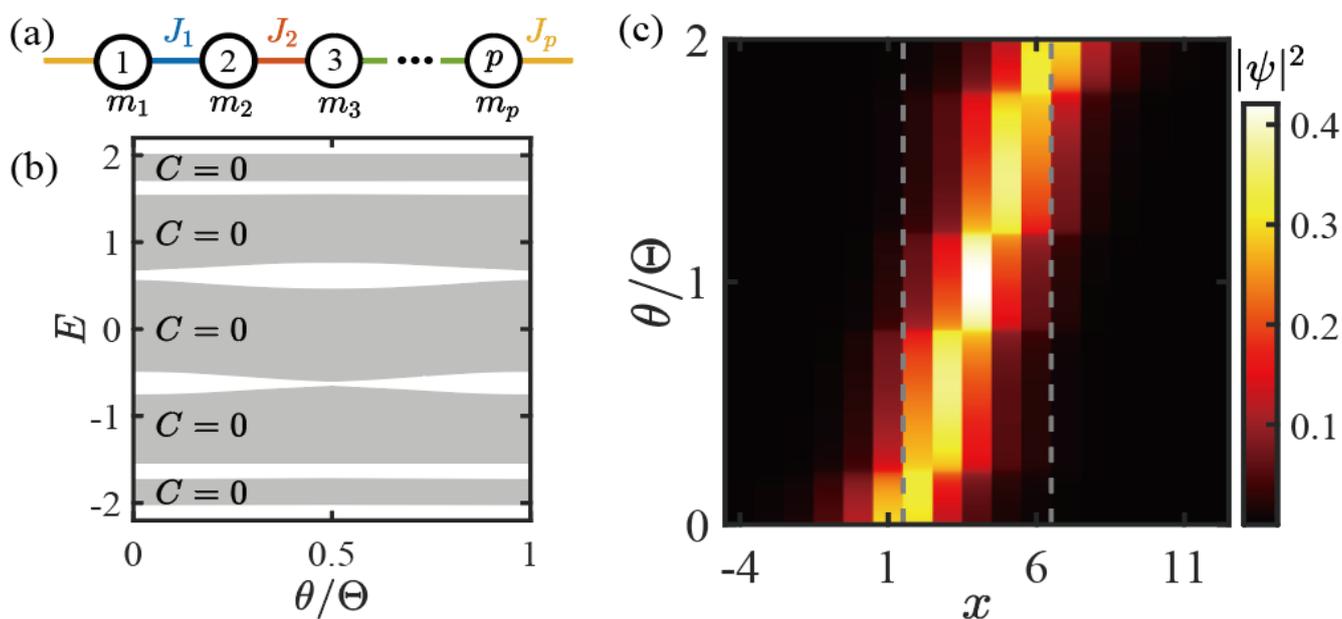
九、凝聚态物理学

主要完成人：徐勇研究组

非线性诱导的孤子分数 Thouless 泵浦

非线性效应在光学系统和玻色-爱因斯坦凝聚等体系中发挥着关键作用，可引发混沌、孤子等重要现象。孤子是一种在传播过程中保持波形不变的孤立波包，其稳定性源于非线性与色散之间的相互平衡。近期研究发现，在非线性和系统中，当线性哈密顿量发生周期性缓慢变化时，孤子可实现整数或分数形式的量子化运输，这一现象被称为非线性 Thouless 泵浦。该现象可理解为孤子沿着线性哈密顿量的瞬时最大局域化单带或多带 Wannier 函数进行移动。因此，如果对应能带的总陈数为零，孤子将不会发生位移。值得注意的是，最近有研究指出，孤子整数位移与陈数的对应关系存在失效情形。然而，在线性拓扑平庸的系统中，孤子能否实现分数 Thouless 泵浦，目前仍不明确。

徐勇研究组首次在理论上发现了由非线性诱导的孤子分数 Thouless 泵浦。该研究组基于线性拓扑平庸的非对角 AAH 模型，通过引入不随时间变化的非线性，观察到孤子可在二、三或四个周期内移动一个原胞距离，证实了非线性可诱导实现分数泵浦。该研究组将此现象归因于孤子解所引起的局域势修正，该修正有效地将原本拓扑平庸的线性哈密顿量转变为拓扑非平庸形式。由于该模型仅需调节最近邻跃迁耦合，且非线性保持恒定，因而易于在现有的先进光学平台上实现。此项工作为后续探索非线性系统中的新颖拓扑现象开辟了道路。



(a) 非对角 AAH 模型示意图；(b) 线性能谱图；(c) 瞬时孤子演化图

该成果研究论文：Yu-Liang Tao, Yongping Zhang, and Yong Xu, "Nonlinearity-Induced Fractional Thouless Pumping of Solitons", Phys. Rev. Lett. 135, 097202 (2025).



Editor:
Kailin Li
Reviewer:
Wei Xu, Jian Li, Yipu Song