

# 人工智能辅助科研要从可用走向可信

◎本报记者 都芑

对于科研工作者来说,检索、阅读文献是一项费时费力的工作。在大模型发展如火如荼的今天,以其为代表的

人工智能正渗透进人们工作生活的各个角落,科研领域也不例外。

日前,阿里巴巴发布了基于Transformer架构自主研发的千亿参数级夸克大模型。据介绍,该大模型可用于科研资料收集、文献快速阅读与翻译、创作润色等场景。

不仅是阿里巴巴,科大讯飞股份有限公司(以下简称科大讯飞)、腾讯等企业,也都推出了用于辅助科研的大模型产品。这一系列产品的问世,正悄然改变着科研工作者的工作方式。

## 大模型已进入科研领域

今年初,ChatGPT的走红掀起了语言大模型热潮。人们可以随心所欲地提出问题,大模型总会给出答案。这股风很快也吹到了科研领域。ChatGPT发布后不久,一款名为tsyz.ai的应用插件在科研圈中受到追捧。

这是一款借助ChatGPT的强大理解能力,专门用来阅读科研文献的插件。用户可以直接将论文全文上传至该应用,并提出相应解读要求,它能够以最快速度对用户提出的问题予以回答。

即使没有下载论文全文也没关系,tsyz.ai支持对论文预印本网站进行检索。用户可以只提供一个论文ID序号,tsyz.ai就会自动检索、学习该论文,并根据用户需求给出回答。不仅如此,用户还能以聊天的方式与其进行对话,就论文中的内容提出各种问题。

视频网站哔哩哔哩知名科普博主严伯钧是tsyz.ai的忠实用户,他时常在各类科普视频中使用tsyz.ai来协助解读论文。在他看来,tsyz.ai给出的论文解读准确率已经非常高,具备很强的实际应用价值,可以帮助科研工作者更加高效地检索、阅读文献。

“tsyz.ai无法解读的情况当然也会有。”严伯钧表示,以他的使用经验来看,向tsyz.ai提出的问题必须是一个能被回答的“有效问题”,“如果问题问得太细、太深,或者过于刁钻古怪,那么它就会直接告诉你,无法回答”。

但必须承认的是,在大模型迅猛发展并逐渐进入千行百业的今天,专门针对科研领域的大模型产品仍然不算多,且大多数是试验性质的产品。

不久前,科大讯飞在发布最新版本的讯飞星火认知大模型V3.0时,也一口气发布了12个面向行业的专用大模型。其中便有联合中国科学院文献情报中心共同研发的、面向科研工作者的科技文献大模型,以及基于该大模型的应用产品——星火科研助手。这也是国内为数不多的专门为科研工作推出的大模型产品。目前,星火科研助手有成果调研、论文研读、学术写作三大功能。

浙江大学第一附属医院图书馆工作人员以“大语言模型”为关键词对星火科研助手进行了试用。在“成果调研”板块,星火科研助手在检索到的1251314篇文献中筛选了167篇文章进行分析,给出了关于大语言模型的概述。其还可以进一步从筛选的167篇文章中勾选最多30篇文章,据此生成综述。

星火科研助手的论文研读功能则采用当前大语言模型通用的问答模式,可默认生成论文摘要、方法、结论等主要信息;用户也可以就自己关心的论文内容进行提问,科研助手会根据文章内容进行回答。其学术写作功能则主要聚焦科研文章的翻译与润色,目前支持中英文互译,也可以对研究人员撰写的英文文章进行润色。

## 算力紧缺背景下提升训练推理效率迫在眉睫

# “技术升级+一站构建”助大模型降本增效

◎本报记者 罗云鹏

如何在算力紧缺的背景下提升大模型训练和推理的效率,并降低成本?这已成为一众大模型企业不得不面对的难题之一。

日前,腾讯披露,腾讯混元大模型背后的自研机器学习框架Angel再次升级。“自研机器学习框架升级后,腾讯大模型训练效率可提升至主流开源框架的2.6倍,用该框架训练千亿级大模型可节省50%算力成本,大模型推理速度提高了1.3倍。”11月30日,腾讯机器学习平台部总监陶阳宇向科技日报记者表示。

不只是腾讯,在提升大模型训练效率、加速大模型落地应用方面,一批中国企业

交出了自己的“答卷”。

双管齐下节约算力成本

在大模型训练和推理过程中,需要消耗大量算力资源。因此,提高硬件资源利用率,对国产大模型技术的发展至关重要。

陶阳宇介绍,面向大模型训练,腾讯自研了机器学习框架Angel。该框架针对预训练、模型精调和强化学习等全流程进行了加速和优化。据悉,它采用FP8混合精度训练技术,并深度优化了4D混合并行训练策略,还在ZeROCache技术基础上减少了冗余模型存储和内存碎片,提升了内存的利用率。同时,该框架还可兼容适配多款国产化硬件。

而据媒体披露,除了提高硬件资源利用率,针对通信策略、AI框架、模型编译等进行系统级优化,亦可大幅节约训练调优和算力成本。

此外,随着模型参数的增大,大模型推理的成本也随之攀升。陶阳宇介绍,腾讯自研的大模型机器学习框架Angel通过扩展并行、向量数据库、批处理等多种优化手段,提高了吞吐能力,达到了更快的推理性能,降低了成本。

不只是腾讯,在第二十届中国计算机大会上,百度首席技术官王海峰就公开透露,文心大模型4.0从今年3月发布至今,其训练算法效率已提升3.6倍;通过百度飞桨与文心大模型的协同优化,文心大模型周均训练有效率超过98%,推理性能提升50倍。

此外,据公开资料显示,阿里云通义大模型则聚焦于规模定理,基于小模型数据分布、规则和配比,研究大规模参数规模下如何提升模型能力,并通过底层集群的优化,将模型训练效率提升了30%,训练稳定性提升了15%。

不难看出,调整和优化模型的训练和推理方式,其最终目的都指向使模型更好地适应实际应用场景、降低在终端应用中的额外成本。“大模型的应用和研发同样重要。”腾讯机器学习平台专家工程师姚军说,只有提供方便、强大的接入平台,才能让大模型真正走向应用。

百度创始人、董事长兼首席执行官李彦宏也曾表示,大模型本身是不直接产生价值的,基于大模型开发出来的应用才是大模型



未来,人工智能或将帮助科研工作者跳过文献检索、粗读的过程,直接找到需要的文献,大幅提升科研工作者的文献阅读效率。

## 须保证内容真实且专业

由于技术原因,大模型有时会出现编造信息、“一本正经地胡说八道”的现象。这种现象在业内被称为AI幻觉。生活中,人们在和大模型聊天时,如果出现了AI幻觉,人们可能会一笑置之;但若AI幻觉出现在追求严谨精确的科研领域,后果可能就会很严重。

科大讯飞北京研究院执行院长、科技文献大模型研发负责人伍大勇表示,研发科技文献大模型,核心难点就在于保证其内容的可信性和专业性。“一方面,这要依靠高质量的论文数据;另一方面,在模型预训练和监督微调方面也需要下功夫。”伍大勇说。

他介绍,科大讯飞通过与中国科学院文献情报中心合作,在合规的情况下获取了丰富的科技文献数据,并对数据进行了去重、去噪等处理,以提升数据质量。“星火科研助手采用中国科学院文献情报中心提供的论文接口来进行论文检索。此外,我们还使用了基于论文知识库的检索增强和知识增强策略。这些都使大模型生成的结果有据可依。”伍大勇表示,这些措施从技术上保证了星火科研助手回答结果的准确性,也尽量避免大模型出现AI幻觉。

同时,伍大勇表示,针对科技文献服务的各个场景,星火科研助手研发团队还邀请专业团队,对大模型训练数据进行监督微调,以提升星火科研助手在科技文献服务上的性能表现。“例如在成果调研和论文研读功能上,我们采用大模型结合知识图谱和知识库的策略,以保证产品输出的内容有据可依。在学术写作上,我们针对学术翻译和学术英语润色专门进行了大模型监督微调,以达到比通用翻译和校对产品更强的专业性。”伍大勇说。

## 或能激发科研工作者更多灵感

虽然目前尚未有太多人工智能产品被应用于科研领

域,但已有学者对人工智能进军科研提出了反对意见,认为这会让科研工作者变得懒惰。在严伯钧看来,科研工作者在应该“懒惰”的地方“懒惰”,反而可以节省出更多时间用在更有价值的工作上。

阅读文献前首先要进行文献检索。为此,科研工作者往往需要搜寻大量文献,在此基础上对部分感兴趣的文献进行粗读,以进一步判断哪些是自己真正需要的文献。这是实打实的“体力活”。严伯钧认为,在这种情况下,借助人工智能工具帮助科研工作者跳过检索、粗读的过程,以更高效的方式直接找到需要的文献,可大幅提升科研工作者的文献阅读效率。

虽然可以借助大模型等工具来检索阅读文献,但严伯钧也不否认读原文的价值。“原文当然要读,但并不一定是每篇都要读。更加精准地定位到需要的文献后再进行精读,是更加高效的方式。”

伍大勇同样表示,研发星火科研助手的初衷在于帮助用户快速了解论文核心内容,提高论文研读效率,让科研工作者能够把更多精力花在更为重要的实验验证等工作上。“辅助提升科研效率是科技文献大模型的关键和目标,但科研工作所需要的灵感、思路、逻辑推理、实验验证、创新与探索等仍离不开科研工作者发挥主观能动性。”

事实上,除了能够辅助阅读文献,人工智能已经在多个科学研究领域带来实际成果。例如在预测蛋白质结构方面,人工智能产生的成果已经远超人类过去工作的总和。严伯钧认为,这种需要大量计算、反复试错的工作,正是人工智能的强项,人类应与其形成合理分工,拥抱新技术。

谈及未来人工智能可能给科研工作带来的改变,严伯钧认为,目前的文献阅读、翻译润色等功能,可能只发挥了人工智能在科研工作领域潜力的1%。在他看来,当下科技发展正呈现出细分化的趋势,一位学者往往只深耕于某一科研领域,而人工智能的跨界思维模式未来或能给科研工作带来一些改变。“或许人工智能可给科研工作者带来更多跨领域、交叉学科的原发性启发,激发科研工作者更多想象力。”

存在的意义。然而,很多大模型落地的难度很大,因为一个大模型往往会对应着很多不同种类的应用,这需要大量的接口和流量支持。

如何破解这道难题?据悉,基于自研机器学习框架Angel,腾讯打造了大模型接入和应用开发的一站式平台,让针对业务场景的数据处理、模型微调、评测部署和应用构建等多个环节,从以往“散装”的多团队协作方式,转化成流水线平台上自动化生产方式,让大模型的“开箱即用”成为可能。“开箱即用”的关键在于预训练基础模型的泛化能力,高性能框架提供的微调或扩展工程能力,以及应用平台的灵活构建能力等支撑。据媒体披露,目前腾讯会议、腾讯新闻、腾讯视频等超过300个腾讯产品及场景均已接入腾讯混元大模型进行内测,数量相比10月份翻了一倍,覆盖文本总结、摘要、创作、翻译、代码等多个场景。比如,腾讯混元大模型就可支持智能化的广告素材创作,满足“千人千面”的需求。

《北京市人工智能行业大模型创新应用白皮书(2023年)》数据显示,截至2023年10月,我国10亿参数规模以上的大模型厂商及高校院所共计254家,分布于20余个省市/地区。

“未来大模型产品的发展趋势可能是通用大模型与垂直领域细分模型的结合。”中国人民大学数字经济研究中心主任李三希此前表示,这不仅需要具备坚实的技术基础,如大规模、高质量、多样化的语料库,创新的大模型算法,自研的机器学习框架和强大的算力基础设施等,也需要大模型产品具有坚实的基于场景的应用。未来,从实践中来,到实践中去的“实用级”大模型将成为趋势。

“未来大模型产品的发展趋势可能是通用大模型与垂直领域细分模型的结合。”中国人民大学数字经济研究中心主任李三希此前表示,这不仅需要具备坚实的技术基础,如大规模、高质量、多样化的语料库,创新的大模型算法,自研的机器学习框架和强大的算力基础设施等,也需要大模型产品具有坚实的基于场景的应用。未来,从实践中来,到实践中去的“实用级”大模型将成为趋势。

“未来大模型产品的发展趋势可能是通用大模型与垂直领域细分模型的结合。”中国人民大学数字经济研究中心主任李三希此前表示,这不仅需要具备坚实的技术基础,如大规模、高质量、多样化的语料库,创新的大模型算法,自研的机器学习框架和强大的算力基础设施等,也需要大模型产品具有坚实的基于场景的应用。未来,从实践中来,到实践中去的“实用级”大模型将成为趋势。

“未来大模型产品的发展趋势可能是通用大模型与垂直领域细分模型的结合。”中国人民大学数字经济研究中心主任李三希此前表示,这不仅需要具备坚实的技术基础,如大规模、高质量、多样化的语料库,创新的大模型算法,自研的机器学习框架和强大的算力基础设施等,也需要大模型产品具有坚实的基于场景的应用。未来,从实践中来,到实践中去的“实用级”大模型将成为趋势。

## 《连线》创始主编凯文·凯利: AI不会是人类最后的发明

◎实习记者 骆香茹

近日,知乎知学堂联合电子工业出版社推出“对话凯文·凯利——AI的过去、现在与未来”直播,邀请《连线》杂志创始主编凯文·凯利与知乎CTO李大海等人连线交流,探讨人工智能(AI)对当下的影响。

面对AI技术被广泛应用于人们的日常生活这一现实,凯文·凯利在其作品《5000天后的世界》发表后曾说过一句话:“未来50年都将会是AI主导的一段时期;未来5000天后,会是AI时代。”

与此同时,业界也存在不同观点。埃隆·马斯克、比尔·盖茨等科技界翘楚都非常担心,认为AI会是人类最后的发明。不过,凯文·凯利认为,这个可能性非常小。“他们之所以认为AI会给人类造成生存威胁,是因为在他们的理解中,AI的力量有指数级发展迹象,但目前并没有证据可以证实这一点。只能说有这种可能性,但可能性微乎其微。”凯文·凯利表示。

那么,该如何看待AI在人类世界的定位?凯文·凯利认为,AI有不同特性,与人类的关系也有不同形态。但不管与人类是何种关系,“AI都不会淘汰我们。但我建议大家去努力学习使用AI,学习如何与AI协作。”凯文·凯利说。

## 开源助推AI技术落地

◎本报记者 操秀英

近日,浪潮电子信息产业股份有限公司(以下简称浪潮信息)发布千亿级开源大模型“源2.0”。“源2.0”创新采用局部注意力过滤增强机制(LFA),可以有效捕捉局部信息和短依赖信息,使模型能够更精准地掌握上下文之间的强语义关联,学习人类语言习惯范式本质,大幅提升数理逻辑、数学计算、代码生成能力。

浪潮信息于2021年9月在业界率先推出了中文人工智能(AI)巨量模型“源1.0”,参数规模高达2457亿。浪潮信息人工智能软件研发总监吴韶华介绍,比起“源1.0”,“源2.0”在算法、数据、计算等方面都实现了创新。在算法上,该模型基于LFA。有别于传统Transformer模型结构擅长捕捉全局信息和长依赖信息能力的特点,LFA具备有效捕捉局部信息和短依赖信息的能力,可确保模型更精准地学习人类语言范式本质。

在数据处理方面,“源2.0”通过使用中英文书籍、论文等资料,结合高效的数据清洗流程,为大模型训练提供了高质量的学科专业数据集和逻辑推理数据集。除此之外,浪潮信息提出基于单元测试的数据清洗方法,可更高效地获取高质量数据集,提高训练效率。“有限的算力资源下,训练数据的质量直接决定了模型的性能。”吴韶华说,“源1.0”绝大部分的数据来源于网页,虽然我们花费了很大力气清洗,但数据质量确实需要进一步提高。“源2.0”减少了网页数据,增加了书籍、期刊等的数据,并引入代码和数学数据,使模型数理逻辑能力进一步增强。”

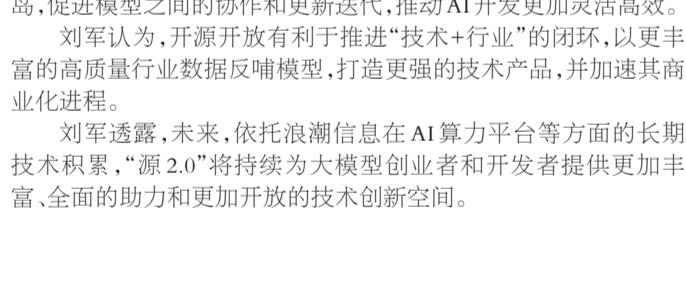
不仅如此,“源2.0”还将实行全面开源(模型全开源、免费可商用、无需申请授权)。对此,浪潮信息高级副总裁刘军表示,从计算机科学与人工智能的发展历程来看,开源始终对软件技术乃至IT技术的发展有巨大的推动作用。“Meta公司的LLaMA大模型开源之后,迅速吸引了大量开发者。”刘军说,在国内,开源开放是促进AI技术发展和商业落地的重要手段,大模型的开源开放可以使不同的模型之间共享底层数据、算法和代码,有利于打破大模型孤岛,促进模型之间的协作和更新迭代,推动AI开发更加灵活高效。

刘军认为,开源开放有利于推进“技术+行业”的闭环,以更丰富的高质量行业数据反哺模型,打造更强的技术产品,并加速其商业化进程。

刘军透露,未来,依托浪潮信息在AI算力平台等方面的长期技术积累,“源2.0”将持续为大模型创业者和开发者提供更加丰富、全面的助力和更加开放的技术创新空间。

## 图说智能

### 智能驾驶:未来出行新趋势



在前不久落幕的第21届广州国际汽车展览会上,搭载智能驾驶技术的最新成果纷纷亮相。近年来,智能驾驶技术在城市市场落地,这不仅带来了全新出行体验,激发出产业新动能,也成为未来出行的新趋势。图为在第21届广州国际汽车展览会上拍摄的无人驾驶快递车,它支持用户在手机下单收发快递。

新华社记者 刘大伟摄

本版图片除标注外由视觉中国提供