

# 国产大模型研发现状与创新方向

文 | 安晖



安晖

中国电子信息产业发展研究院副总工程师  
民盟中央科技委副主任  
中国电子学会理事

大模型是深度学习领域的最新成果，被视为人工智能发展乃至新一轮科技革命与产业变革的重要驱动力，成为各国、各地区以及龙头科技企业争相抢占的制高点。ChatGPT是大模型领域的现象级产品，一经出现就引发全球关注，也吸引了我国众多企业、高校院所和专业机构迅速投身到大模型发展大潮中。当前，我国推出的大模型已超过140个。

理论突破和技术创新是大模型发展的首要驱动力。Transformer模型的出现和所引入的自注意力机制，为大模型的快速发展奠定了基础。国外的OpenAI、谷歌、Meta等企业和机构不断推出的新技术、新产品，促进了大模型的迭代创新。我国的百度、科大讯飞、华为、阿里巴巴、腾讯等龙头企业也在模型算法、算力设施、训练资源等方面大力投入，推进实现自主研发，并取得了较好成绩。

从技术类型看，国内大模型的类型分布与国外基本一致。一方面，语言大模型是最主要、最活跃的领域，文本生成、文本对话、机器翻译等是各语言大模型关注和研发实现的重点功能。另一方面，涉及视觉、听觉及编程、逻辑的大模型技

术逐渐受到关注。识文作画、识文作曲、视频生成、读图知意以及代码撰写、数学解题等应用不断涌现，而且形成了一些具有较大传播度的技术产品。这表明，多模态正在成为研发的重要方向。

从技术路线看，目前国内大模型研发主要追随全球先进成果。前期，由于BERT模型的出色表现，我国企业探索大模型时更多参考BERT路线。随着GPT系列特别是ChatGPT表现出超预期成果，国内大模型发展的技术路线也开始向GPT方向倾斜。这一转变与全球大模型技术的演进相一致。也正因如此，国内大模型在基础理论和基本技术上基本同源，主要差异表现在模型设计和训练方式上。例如，阿里巴巴强调多模态任务能力及效率，百度聚焦NLP能力的提升，腾讯兼顾模型规模增长与效率提升。

从研发方式看，国内大模型的研发路径可主要归纳为四类。一是完全自研。即基于Transformer论文原理，完全从零研发，不依赖任何开源代码。二是基于开源的Transformer架构研发。国内主要的大模型多属于此类。三是基于代码开源的大模型研发。即在其他企业、机构研发的大模

型的开源代码基础上，调整神经网络架构、代码后，经训练而成。国内部分大模型属于此类。四是基于参数开源的大模型研发。即在其他企业、机构研发的经过预训练且具备基础能力、不开源代码仅开源参数的大模型基础上，经精调训练而成。国内也有部分大模型属于此类。尤其是不少行业大模型，多是通过对国外的LLaMA等开源大模型进行微调或修改实现，属于第三、四类。

总体来看，第一、二类有一定的原创性和自主性，第三、四类则较可能存在风险。例如，基于代码开源研发的大模型，如所参考的大模型调整开源策略或转为闭源，就会面临知识产权等风险。例如，基于LLaMA的应用开发条款约定，当月活跃用户达到7亿之后，需要向Meta公司申请额外授权，否则将无权继续使用。基于参数开源研发的大模型，无法查看并确认所参考大模型内部代码的安全性，也无法基于此自行演进，所以存在较大的信息安全风险和断供风险。

之所以大量的国产大模型不属于源头创新，原因有很多，其中之一是基础理论支撑不足。国外大企业发展大模型，走的是以理论研究为基础的模式。例如，神经元数学模型、知识图谱、深度学习框架、Transformer架构等开创性技术，大多由美国科学家提出。这点与大数据的涌现、发展类似。但国内大模型发展则缺少具有显著创新性的原创理论支撑，到目前为止，我国大

模型发展仍属于跟随式发展。

国产大模型已在应用方面进行了不少探索。这方面工作虽难以说超越美国企业，但确实实现了与国情相结合，并形成了一些创新性的成果。在制造业领域，大模型被应用于产品研发、生产维护、质量控制等领域，促进性能改进、材料选择和生产技术。例如，理想汽车宣布在智能座舱中加入自研Mind GPT大模型，使其语音系统变成多重身份的拟人助手，可根据输入的需求自动生成故事、讲解儿童百科、制定出行规划、完成复杂的任务理解等。在传媒领域，采编、传播等环节，可通过基于大模型的语音转写、智能写作、智能剪辑等方式提高采编环节的生产效率，并通过打造AI主播实现智能、高效播报。

但必须注意到的是，大模型应用、大模型商业模式目前都还在探索过程中，与巨额投入相比，产出偏少。如何实现大模型研发成果向产业的转化，是当前国内外大模型企业都需要加快探索的重要任务。

大模型在行业、机构、个人领域的应用将不断普及和深化。高德纳（Gartner）发布2024年十大战略技术趋势之首即是“全民化的生成式AI”，并认为到2026年，将由超过80%的企业使用生成式AI的API或模型，或在生产环境中部署支持生成式AI的应用。在旺盛需求的推动下，预计未来很长一段时间，大模型技术产品以及以其为基础的生成式AI的研发与应用都将是全球人工智能领域乃至整个科技领域的重点。

一方面，对以GPT等为代表的基础大模型的研究将会继续深入，科研机构、科技巨头企业都会不断投入巨大资源，并不断实现基础理论和关键技术方面的研发突破。

另一方面，为适应经济社会各领域的不同需求，为与重点行业的技术产品融合发展，大模型也会不断涌现新的重点。

初步预计，多模态大模型、轻量级大模型、大模型与小模型协同、基于大模型的具身智能和终端智能、嵌入式大模型（私域大模型）等将是值得关注的创新方向。在此过程中，高校院所、龙头企业预计将成基础模型理论创新的主力军，行业大模型企业、算法企业则会遵循工程化发展路径，面向贴合实际、更加具体的典型应用场景，从技术方面丰富大模型功能、提高大模型效率、降低大模型的应用门槛。

未来，大模型以及通用人工智能、AIGC等领域的技术创新仍将持续，并不断取得新的成果。工业和信息化部已通知启动2023年未来产业创新任务揭榜挂帅工作，并将通用人工智能列为四大领域之一，将语言大模型产品、语音大模型产品、视觉大模型产品、多模态大模型产品，以及面向工业制造、民生服务、科学研究、信息安全等领域的典型应用列为重点攻关内容。相信在国家的引导、支持下，我国大模型研发和应用创新将不断结出新的硕果。

编辑 | 李佳琪