

计算机行业研究

买入（维持评级）

行业深度研究

证券研究报告

计算机组

分析师：王倩雯（执业 S1130522080001） 分析师：孟灿（执业 S1130522050001）

wangqianwen@gjzq.com.cn

mengcan@gjzq.com.cn

大模型时代，AI 技术向效率提升演进

投资逻辑

我们 2022 年 12 月发布的报告《深度学习算法：从多样到统一》中，阐述了自 Google 2017 年提出 Transformer 以来，深度学习开始进入大模型时代。大模型时代的前沿技术发展围绕着提升效率而展开，包括：1) 提升训练方法效率：向无监督和半监督学习发展；2) 提升数据效率：从追求数据规模向追求数据质量发展；3) 提升开发效率：通过“预训练基础模型+微调”，挖掘现有大模型潜力，降低具体下游任务的开发成本；4) 提升算力效率：从稠密机构向稀疏结构发展；5) 提升训练的工程化效率：向并行训练和混合精度训练发展

- 训练方法：AI 模型的训练方法主要包括监督学习和无监督学习两种典型方式，后随模型训练数据量的增加，衍生出使用大量未标注数据+少量标注数据的半监督学习方法。AI 训练方法的发展历经“监督-无监督-监督-无监督/半监督”4 个阶段，在目前的大模型阶段，无监督/半监督训练再次成为主流。
- 数据效率：随参数规模的增加，大模型在知识密集型任务中的效果提升显著。此外，当模型参数超过特定阈值后，模型会对特定任务表现出“涌现”现象。目前学界和业界已意识到数据质量的重要性或高于数据数量，AI 大模型需要在保证数据质量的前提下进行数据数量和参数规模的扩充。
- 开发效率：AI 大模型的流行提出了“基础模型+微调”的 AI 开发新范式。相较于过去“一场景、一任务、一模型”的开发模式，“基础模型+微调”具有数据需求量小、训练时间短、落地边际成本低等优点。微调技术的发展带动大模型由“以参数规模取胜”向“以高质量学习取胜”转变。
- 算力效率：AI 架构可分为稠密结构和稀疏结构，其中稀疏结构可有效降低大模型对算力的消耗。2017 年 Google 提出了混合专家方法 MoE，使得模型在计算过程中只需激活部分神经网络；2022 年 6 月 Google 发布的基于稀疏结构的多模态模型 LimoE，已经在降低算力消耗的同时取得不亚于稠密结构的业绩。
- 工程化效率：伴随 AI 大模型参数量的不断提升，并行训练、混合精度训练等技术发展迅速。其中，国产 AI 框架百度 PaddlePaddle 提出的 4D 混合并行策略在 MLPerf 发布的稠密结构 AI 训练性能榜单中位列第一；通过使用 16 位浮点数代替 32 位浮点数进行训练，能够在同等模型表现的情况下实现训练时间减半。

投资建议

建议关注受益于 AI 算法进步，并能成功进行商业化应用的科大讯飞、商汤科技等公司；以及受益于 AI 算力需求、微调技术发展的海光信息、浪潮信息、海天瑞声等公司。

风险提示

海外基础软硬件使用受限；骨干网络创新放缓；应用落地不及预期

内容目录

1. 训练方法演进：无监督、半监督训练再次成为主流.....	3
2. 训练数据演进：从追求规模到追求质量.....	6
3. 开发方式演进：微调技术受到重视.....	7
4. 架构设计演进：从稠密结构到稀疏结构.....	8
5. 训练技术演进：并行训练与混合精度训练.....	9
6. 投资建议.....	10
7. 风险提示.....	10

图表目录

图表 1: 监督学习与无监督学习方式对比.....	3
图表 2: LeNet-5 卷积神经网络典型结构.....	4
图表 3: 逐层无监督+BP 有监督可解决梯度消失问题.....	4
图表 4: 计算机视觉领域经典开源数据集.....	5
图表 5: 自然语言处理领域的无监督学习方法.....	5
图表 6: MAE 无监督学习方法在多个下游任务中优于监督方法.....	6
图表 7: 知识密集型任务表现随参数规模提升.....	7
图表 8: AI 大模型在复杂任务中表现出“涌现”现象.....	7
图表 9: InstructGPT/ChatGPT 中的人类反馈强化学习技术.....	8
图表 10: 稠密结构与稀疏结构对比.....	9
图表 11: 混合专家方法示意.....	9
图表 12: 百度 PaddlePaddle 4D 混合同步策略示意.....	10

我们 2022 年 12 月发布的报告《深度学习算法: 从多样到统一》中, 阐述了自 Google 2017 年提出 Transformer 以来, 深度学习开始进入大模型时代。本文旨在讨论大模型时代下, 整个 AI 行业的技术演进的前沿发展方向。

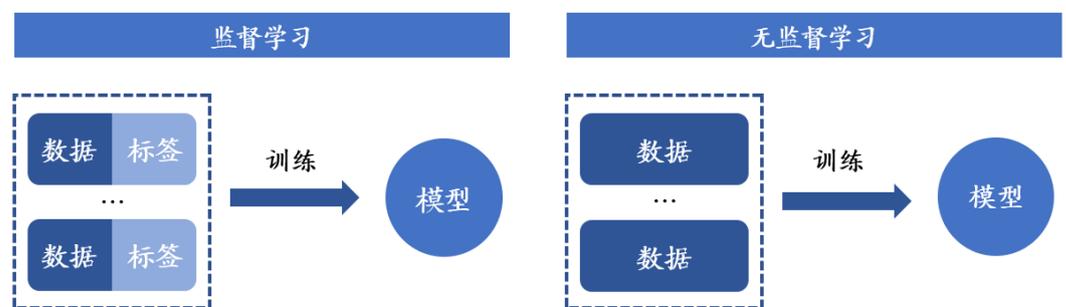
小结而言, 大模型时代的前沿技术发展围绕着提升效率而展开, 包括: 1) 提升训练方法效率: 向无监督和半监督学习发展; 2) 提升数据效率: 从追求数据规模向追求数据质量发展; 3) 提升开发效率: 通过“预训练基础模型+微调”, 挖掘现有大模型潜力, 降低具体下游任务的开发成本; 4) 提升算力效率: 从稠密机构向稀疏结构发展; 5) 提升训练的工程化效率: 向并行训练和混合精度训练发展。

1. 训练方法演进: 无监督、半监督训练再次成为主流

AI 模型的训练方法主要包括监督学习与无监督学习两种典型方式, 二者的区别在于是否使用带人工标注的数据集进行训练。此外, 随着模型训练数据量的增加, 标记大量样本成本过于昂贵, 衍生出使用大量未标注数据+少量标注数据的半监督学习方式。

目前, 虽然模型参数的扩大仍能提升模型表现, 但扩大相同规模的参数较大模型发展初期的边际收益递减, 提升数据质量是未来模型智能水平提升的关键。

图表 1: 监督学习与无监督学习方式对比



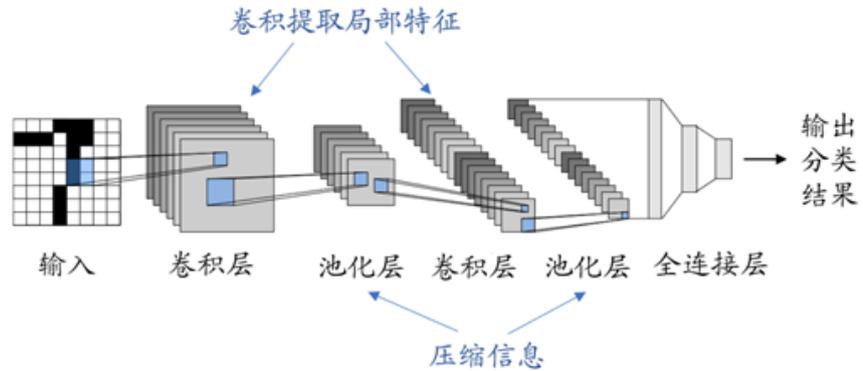
来源: CSDN 云计算公众号, 国金证券研究所

AI 训练方法的发展历经“监督-无监督-监督-无监督/半监督”4 个阶段, 在大模型时代下, 无监督/半监督训练再次成为主流方法。

■ 2006 年之前, 浅层神经网络的训练以监督学习为主:

- 算法层面, 这一阶段的神经网络尚停留于浅层, 强调通过学习少量数据获得较强的性能, 监督学习的表现显著优于无监督学习。此外, 这一时期的支持向量机 (SVM) 等浅层学习算法表现出色, 性能优于同时期的神经网络算法, 在学术界与产业界占据主流地位, 而支持向量机通常采用监督学习方式, 这也使得监督学习成为神经网络的首选训练方式。
- 数据层面, 这一阶段的神经网络由于性能有限, 无法处理复杂任务, 应用场景较为简单。1998 年 Yann LeCun 等人研发的 LeNet-5 是这一时期最具代表性的神经网络模型, LeNet-5 基于卷积神经网络算法开发, 被当时大多数美国银行用于识别支票上的手写数字。简单的应用场景意味着特征信息易于获取, AI 模型仅需要对少量数据进行学习就能获得较强的性能。同时, 由于对数据量需求较低, 标注数据并非难事。

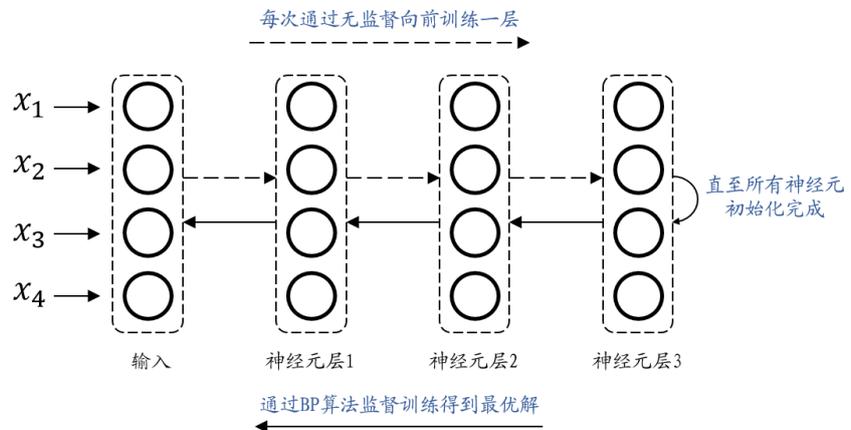
图表2: LeNet-5 卷积神经网络典型结构



来源:《Dive into Deep Learning》(Aston Zhang 等, 2021), 国金证券研究所

- 算力层面, 这一阶段的神经网络采用 CPU 进行训练, 算力匮乏、算力成本高昂是这一时期 AI 发展的主要瓶颈。这使得对数据量要求较低、算力需求少的监督学习成为主流的神经网络训练方式。
- 2006 至 2011 年, 神经网络向深层次发展, 无监督学习成为这一时期的主流方法:
 - 算法层面, Hinton 使用逐层无监督的方法缓解了梯度消失问题, 神经网络得以迈向深层, 性能上限极大提升, 将神经网络用于处理复杂场景任务成为可能。此后, 以 Hinton 为首的学者们开始尝试通过深度神经网络来模拟人的智能, 无监督学习成为这一阶段神经网络的主流训练方式: 1) 无监督学习在核心思想上与逐层无监督方法类似; 2) 仿生观念在当时颇为流行, 人类在学习时通常是无监督的。

图表3: 逐层无监督+BP 有监督可解决梯度消失问题



来源:《A Fast Learning Algorithm for Deep Belief Nets》(Hinton 等, 2006), 国金证券研究所

- 数据层面, 神经网络的应用场景日益丰富, 包括无人驾驶、语音识别等。复杂场景意味着特征信息难以获取, AI 模型必须对更多的数据进行学习才能够获得所需的性能。然而数据因素对于该时期主流神经网络训练方式的影响有限: 1) 面向复杂场景研究刚刚起步, 数据积累较少, 无监督学习方式不具备优势; 2) 深度学习方法尚未成熟, 学者普遍认为算法水平决定一切, 数据因素尚未得到足够重视。
 - 算力层面, GPU 加速神经网络训练的方法被提出, 算力得到了较大的提升, 但 GPU 并没有成为主流的训练硬件。
- 这一阶段的神经网络在算法上迎来了突破, 解锁了神经网络处理复杂问题的潜力, 无监督学习主要是作为梯度消失问题的缓解措施, 并没有使得深度学习模型性能出现明显提升, 算法是这一阶段制约人工智能发展的主要因素。
- 2012 至 2017 年, AlexNet 的成功使得监督训练再度流行:
 - 算法层面, Hinton 及其学生于 2012 年提出 AlexNet 模型, 自此奠定了深度学习的经典训练范式。AlexNet 采用了经典的 CNN 网络结构、使用 ReLu 激活函数、

对输入值进行有监督学习、并采用 GPU 对训练进行加速。由于 AlexNet 将 ImageNet 数据集上图像分类的错误率由 26% 降至 15%，此后 5 年学术界均沿用 AlexNet 的范式进行深度学习训练，监督学习也因此成为了这一时期主流的神经网络训练方式。

- 数据层面，从这一时期开始，数据量被认为是提升 AI 智能水平的关键要素，以 ImageNet 为代表的开源标注数据集发展迅速，这类标注数据集提供的数据量已经足以满足当时绝大部分的 AI 训练需求，并且应用起来方便快捷，这使得监督学习更为流行。

图表4: 计算机视觉领域经典开源数据集

数据集名称	数据量	数据集内容
ImageNet	1,420 万张图像，涵盖 2 万多个类别	图像分类、对象检测
CIFAR-10	6 万张图像，涵盖 10 个类别	图像分类
MegaFace	67 万名人像，共 475 万张图片	人脸识别
MPII	2.5 万张图像，涵盖 410 项人类活动	人体姿势识别
Flicker-30k	15.8 万个众包字幕，描述了 3.2 万张图像	图像与图像描述
MSCoco	32.8 万张图像，250 万个标记实例	对象检测、分割、图像描述

来源: Paperswithcode, 国金证券研究所

- 算力层面，AlexNet 模型的成功在学界与业界推广了 GPU 加速人工智能训练的新模式，算力瓶颈得到极大缓解。

此阶段神经网络的发展主要由算法创新驱动，由标注数据提供训练支持，模型性能得到较大提升。

- 2017 年至今，Transformer 开启大模型时代，无监督和半监督学习再次兴起：

2017 年 Transformer 问世后，深度学习对数据的需求量爆发增长，无监督学习方法成为了这一时期的主流训练方式。在数据量与模型表现高度关联的大模型时代，高效的无监督学习算法能够显著提高模型智能水平，无监督学习也由此迎来了飞速发展。

- 在自然语言处理领域，无监督学习技术发展较快。2018 年，在 Transformer 架构问世一年后，基于无监督学习的 BERT、GPT 等大规模语言模型相继问世，并提出了自回归、MLM、NSP 等无监督学习方法，这些方法的表现较好，一直沿用至今。

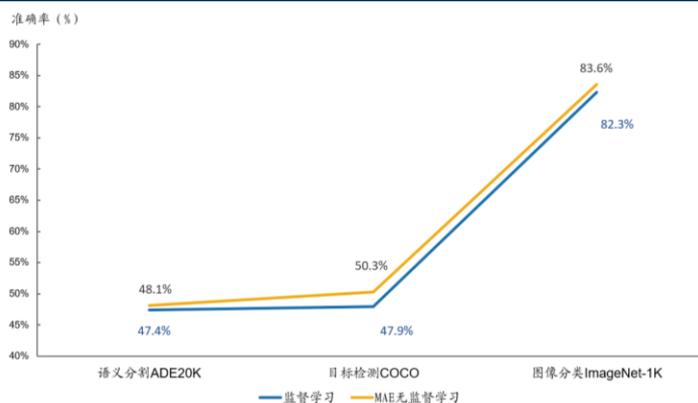
图表5: 自然语言处理领域的无监督学习方法

模型名称	发布者	无监督学习方法	核心思想
GPT	OpenAI	自回归	利用前文单向预测
BERT	Google	MLM、NSP	文本掩码，双向预测

来源: OpenAI, Google, 国金证券研究所

- 在计算机视觉领域，无监督学习技术发展相对较慢。2020 年，基于对比学习思想的 MoCo 问世，证明了无监督学习在计算机视觉领域能取得不亚于监督学习的效果。在此之后，基于对比学习的无监督学习方法不断演进，朝着结构更简单、对数据样本要求更低、更容易应用的方向发展，准确度也不断提升。

图6: MAE 无监督学习方法在多个下游任务中优于监督方法



来源:《Masked Autoencoders Are Scalable Vision Learners》(Kaiming He 等, 2021), 国金证券研究所

2021 年, Facebook AI (现 Meta AI) 的何恺明等提出了 MAE 方法, 该方法的核心思想与自然语言处理领域中的 MLM 方法相同, 同样是随机掩盖图像信息, 并在训练过程中对图像进行预测与重构。MAE 方法对数据的泛化性更强, 更善于处理大规模数据, 将无监督训练的速度提高了 3 倍以上, 在多个下游任务中表现比监督学习更好。

至此, 无监督学习方法在自然语言处理、计算机视觉两个深度学习最重要的领域完成了统一; 也由此结束了飞速发展期, 转而进入缓慢发展阶段。

目前, 国内外 AI 公司发布的大规模基础模型都采用了无监督学习方法。该方法放大了场景拥有者的竞争优势。在自动驾驶领域, Tesla 的 Auto pilot 通过无监督学习使用数十万 Tesla 司机的行为数据来训练 AI 模型。据 Tesla 于 2021 年 AI Day 公布的数据, Tesla 平均每天会收到 500,000 条以上的驾驶数据视频, 并采用自动标注技术(Auto Labeling)自动生成训练信号, 以此训练新的 AI 模型, 向 L5 级自动驾驶逐步迈进。2021 年全年, Tesla 共训练了 75,000 个 AI 模型, 平均每 8 分钟就要训练一个新的 AI 模型。无监督学习方法使 Tesla 大大降低了 AI 模型的训练成本、提高了 AI 模型的迭代速度, 帮助 Tesla 利用自身的数据优势保持在自动驾驶领域的领先地位。

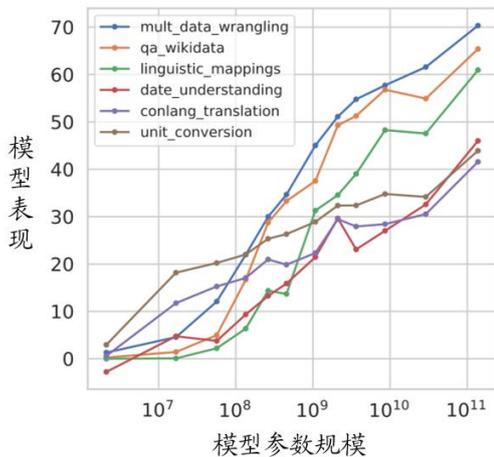
2. 训练数据演进：从追求规模到追求质量

BERT、GPT-3 等 AI 大模型的成功使人们认识到模型参数规模、训练数据量对于提高 AI 智能水平效果显著, 引发了大规模基础模型开发浪潮, 各国内外 AI 巨头纷纷跟进, 研发自有的参数规模更大、性能更强的 AI 大模型, 享受算法进步带来的数据规模红利。

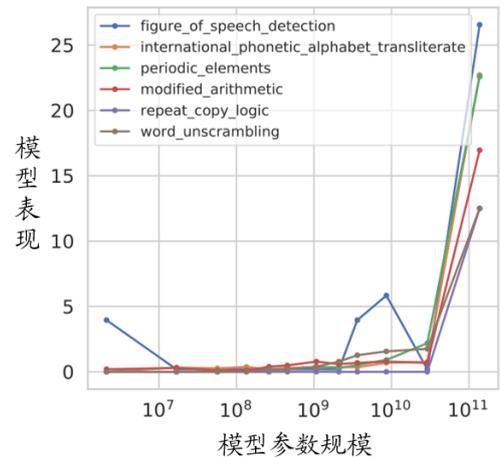
伴随参数规模的不断提升, AI 大模型在具体任务中表现出以下特点:

- 1) 随着参数规模的增加, 显著提高了 AI 模型在语言问答、阅读理解等任务中的表现。这类任务通常是知识密集型, 即模型包含的知识量越多, 任务表现越好。AI 大模型的发展使得该类任务效果提升显著。
- 2) AI 大模型表现出了“涌现”现象, 即模型的参数规模跨过特定阈值后, 模型对特定任务出现爆发式性能增长, 而在此之前模型完全不具备解决该任务的能力。具有“涌现”现象的任务往往复杂度较高, 且由多个步骤组成, 比较考验 AI 模型的逻辑推理能力。

图表7: 知识密集型任务表现随参数规模提升



图表8: AI大模型在复杂任务中表现出“涌现”现象



来源:《Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models》(Aarohi Srivastava 等, 2022), 国金证券研究所

来源:《Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models》(Aarohi Srivastava 等, 2022), 国金证券研究所

未来, 数据质量的重要性可能远高于数量。Google 在对其语言模型 T5 的实验中发现: 数据数量与数据质量两个因素间, 数据质量更为重要。AI 大模型的正确发展路径是在在保证数据质量的前提下, 增大数据数量、扩充参数规模。

数据质量的衡量指标包含多个维度, 如真实性、知识密度、多样性等, 在通过提高数据质量获取更强大的模型智能时, 需要综合考虑以上多个维度:

- 1) 真实性: 基于真实场景数据训练得到的模型往往性能较好。深度神经网络的本质基于统计学得到拟合函数, 因此训练数据是否与真实场景具有相同的数据分布对模型性能至关重要, 较大的数据分布偏差会导致 AI 模型的性能和鲁棒性较差。以图像识别任务为例, ImageNet 作为该场景最具代表性的数据集, 其数据真实性距离真实场景有一定差距, 在大多数图片中识别目标均为图像主体, 并且较少存在遮挡、物体旋转等现实中可能遇到的复杂情况。ObjectNet 是学者对应 ImageNet 专门建立的复杂场景数据集, 该数据集所收录的图像均为较复杂的情况, 能够反映人工智能面对现实中复杂问题的表现。根据测试, 各类先进计算机视觉模型在 ObjectNet 上的成绩相较于 ImageNet 下降了 40% 到 45%, 推理准确度从 90% 下降至 50%, 这表明在面对复杂问题时现阶段视觉模型性能仍有待提升。
- 2) 知识密度: 单位数据中的信息含量, 应用高知识密度的数据进行模型训练能够显著提升模型表现, 典型的高知识密度数据包括维基百科、出版书籍、新闻等。
- 3) 多样性: 训练数据的种类, 多样化的数据将赋予 AI 大模型解决不同类型任务的能力。例如, ChatGPT 在进行训练时采用了维基百科、问答网站、Github 代码等多种数据, 这不仅能够提高 ChatGPT 在语言问答、代码生成等任务中的表现, 同时还将显著提升模型的智能水平, 研究表明 ChatGPT 逻辑能力的显著提升来源于应用代码进行模型训练。

3. 开发方式演进: 微调技术受到重视

AI 大模型的流行提出了“基础模型+微调”的 AI 开发新范式。AI 大模型由海量数据通过无监督学习训练得到, 本身不能直接应用于具体任务, 必须经过微调才可投入应用。微调是指基于大规模基础模型, 在现有训练得到的模型参数之上, 针对特定任务类型、应用特定场景的数据对模型进行二次训练。通俗来说, 大规模基础模型为 AI 提供了基础知识, 而微调则是让 AI 获特定领域知识, 并赋予其组织、应用知识的能力。

微调技术专注于挖掘现有 AI 大模型潜力, 主要研究如何将大模型应用于具体场景, 是大模型时代 AI 开发的重点环节。微调技术水平将极大影响 AI 模型的智能水平, 先进的微调技术能够更充分挖掘 AI 大模型的潜力, 做到“事半功倍”。

在 AI 进入大模型时代之前, 如果想将 AI 应用于特定任务, 则必须从零开始训练神经网络, 即所谓“一场景、一任务、一模型”。相比而言, “基础模型+微调”是低成本, 高收益的解决方案, 其主要具备以下优点:

- 1) 数据需求量小。“基础模型+微调”的开发新范式无需模型从头学习所有内容，微调阶段的所有训练都是为了获取特定领域知识。因此微调所需数据量较小，显著降低了 AI 开发中的数据门槛。
- 2) 训练时间短。一方面，模型微调仅需小规模数据即可进行训练，显著降低了资源消耗；另一方面，在微调过程中，神经网络中的大部份层会事先被冻结，这些层涉及的相关参数在训练过程中保持不变，需要训练的参数仅是所有参数中的一小部分。
- 3) 降低边际落地成本。AI 大模型能适应不同场景下的多种下游任务，采用小规模数据针对具体场景“微调”后即可应用，显著降低了 AI 模型重复开发造成的资源浪费，降低了 AI 落地的边际成本

2020-2022 年间，大模型处于 1.0 时代，这一阶段 AI 研究的特点是专注于大模型开发，追求大模型参数规模提升，大模型开发相关技术迭代较快，微调技术相对不受重视，仅仅作为大模型性能的评估工具。

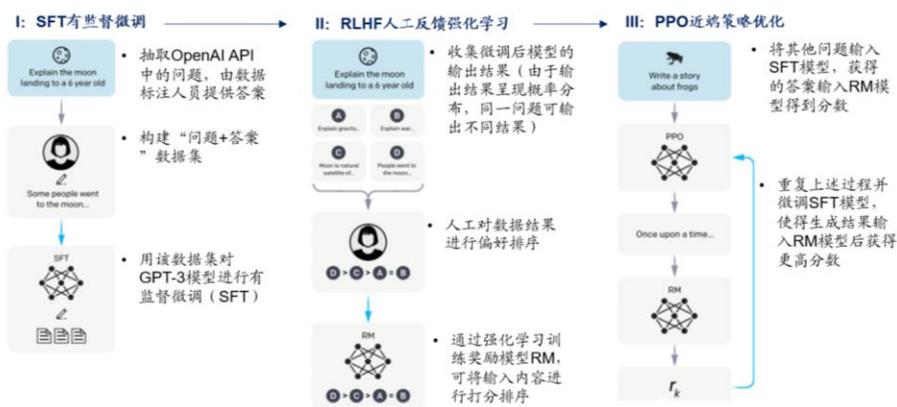
2022 年至今，大模型向 2.0 时代迈进，LaMDA、ChatGPT 等新一代 AI 大模型的成功，标志着 AI 大模型从“以参数规模取胜”向“以高质量学习取胜”转变。微调技术发展对模型智能的提升效果更为显著，模型参数规模提升节奏放缓，落地进程明显加快。

2022 年 1 月，Google 发布对话 AI 模型 LaMDA，该模型不同于以往大模型专注于参数规模提高，其创新点主要集中于微调技术。LaMDA 尝试通过微调方法创新提升模型输出的质量、安全性与可靠性，针对不同的目标 Google 雇佣少量众包人员与 LaMDA 进行对话，众包人员根据对应评价指标对 LaMDA 输出内容通过打分等方式进行数据标注。LaMDA 通过众包模式收集了约 20,000 次对话的注释数据，并基于这些数据对模型进行微调，微调技术的创新显著提高了模型的理解能力，使得 AI 模型与人类交谈时更为拟人化。LaMDA 的成功验证了微调技术创新对模型智能水平的提高具有关键作用。

2022 年 3 月，OpenAI 发布新一代语言模型 InstructGPT，该模型创新点同样集中于微调技术，创新性应用了人类反馈强化学习 (RLHF) 技术，该方法核心思想来源于强化学习，将人类的偏好作为奖励信号训练模型，使得 AI 输出更加符合人类偏好。通过微调技术创新，InstructGPT 使用了少量标注数据，这些数据仅由 40 人团队就完成了标注。

微调方法的进步显著提升了模型智能水平，相比上一代语言模型 GPT-3，InstructGPT 在人工评估中以 13 亿的参数规模战胜了 1750 亿的 GPT-3。

图表9: InstructGPT/ChatGPT 中的人类反馈强化学习技术



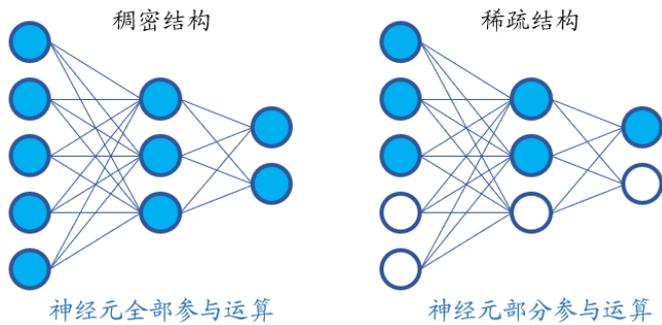
来源: OpenAI, 国金证券研究所

2022 年 11 月 30 日，OpenAI 对外发布新一代对话机器人 ChatGPT，ChatGPT 所应用的人类反馈强化学习 (RLHF) 技术更为成熟。一方面，ChatGPT 继承了 InstructGPT 中的相关技术，能够通过学习人类提高的对话范例，输出更符合人类偏好的内容。另一方面，人类反馈强化学习方法使 ChatGPT 更为谨慎，能够主动拒绝回答不适当的问题，减少输出有害答案，显著提高了 AI 的安全性与可靠性。

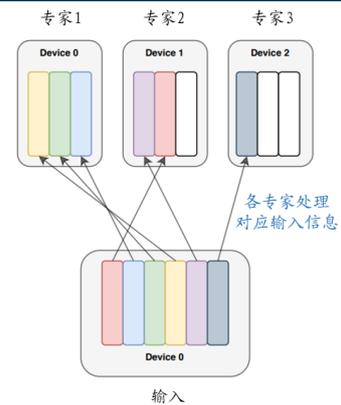
4. 架构设计演进：从稠密结构到稀疏结构

AI 大模型架构设计是指模型的计算架构，主要分为稠密结构和稀疏结构。架构设计决定了模型在训练过程中各神经元间如何相互作用。

图表10: 稠密结构与稀疏结构对比



图表11: 混合专家方法示意



来源: 机器之心公众号, 国金证券研究所

来源: 《Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity》(William Fedus 等, 2022), 国金证券研究所

- 稠密结构主要应用于以 GPT-3 为代表的早期的 AI 大模型, 采用稠密结构的模型在计算时需要激活整个神经网络, 这将带来极大的算力开销和内存开销, 使得 AI 大模型训练成本高昂。
- 稀疏结构的提出旨在降低 AI 大模型对算力的消耗。2017 年, Google 提出了混合专家方法 (Mixture of Expert, MoE), 核心思想是在模型中引入专家层, 每个“专家”处理各自擅长的对应部分输入, 使得模型在计算过程中只需激活部分神经网络。

稀疏结构是一种更像人类的神经网络结构, 其运作过程与人脑极为相似, 人脑中约有 100 亿个神经元, 在执行具体任务的过程中只有部分特定的神经元会被激活, 这种稀疏结构是人脑具备通用且高效智能水平的关键因素之一。

稀疏结构能够显著降低大模型训练成本。2021 年, Google 发布了基于稀疏结构的语言模型 Switch Transformers, 该模型训练效率相比前代稠密结构大模型 T5 提升近 7 倍, 模型参数量达 1.6 万亿, 首次将 AI 大模型参数量推升至万亿级别。

目前, 稀疏结构已经应用至 AI 前沿研究。2022 年 6 月, Google 发布了第一个基于稀疏结构的多模态模型 LimoE, 证明了稀疏结构在降低模型算力消耗的同时, 能够在多项任务中取得不亚于稠密结构的结果。从稠密结构到稀疏结构, AI 大模型架构设计的演进显著降低了模型的算力消耗, 助力 AI 大模型参数规模进一步提升。

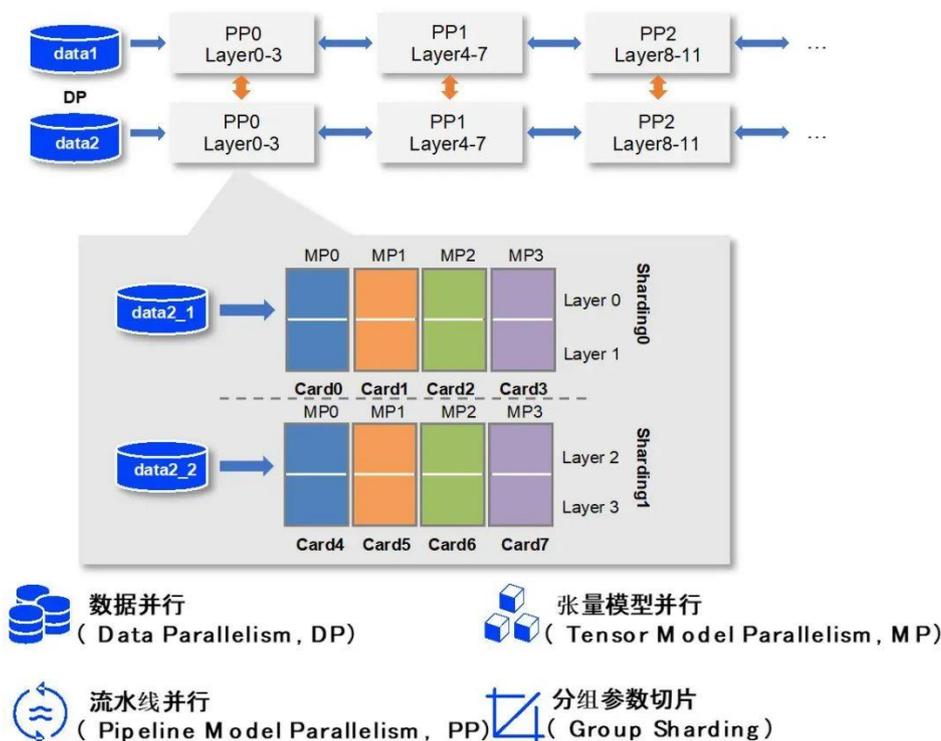
5. 训练技术演进: 并行训练与混合精度训练

训练技术的演进旨在提升 AI 模型训练效率。伴随 AI 大模型参数量的不断提升, 高效训练技术发展迅速, 其主要包括并行训练技术和混合精度训练技术等。

- 并行训练的核心思想是将计算任务切分到不同设备上, 同时尽可能降低设备间通信损耗, 合理使用多台设备的算力, 实现高效的并行训练, 最大化提升模型训练速度。并行训练方法主要包括数据并行、模型并行、流水线并行等多种并行策略, 目前业界主流方法是混合并行方法, 即同时应用多种并行策略, 取长补短、最大限度提升模型的并行能力。

例如, 国产 AI 框架百度 PaddlePaddle 提出 4D 混合并行策略, 其同时应用了四种并行策略, 显著提高了模型训练效率, 支持训练千亿级参数的稠密结构模型、万亿级参数的稀疏结构模型, 其性能在权威 AI 基准评测组织 MLPerf 发布的稠密结构 AI 模型训练性能榜单中位列第一。

图表 12: 百度 PaddlePaddle 4D 混合并行策略示意



来源: 百度 PaddlePaddle, 国金证券研究所

- 混合精度训练的核心思想是通过降低模型训练过程中的参数精度, 以此降低模型训练过程中的算力消耗。该方法的提出是因为研究发现 AI 模型对于参数精度的要求较低, 参数精度的降低几乎不会影响模型表现, 通过牺牲精度换取算力能够显著提高模型的训练效率。例如, 通过使用 16 位浮点数代替 32 位浮点数进行模型参数训练, 能够使模型的训练时间减半, 同时几乎不影响模型表现。

6. 投资建议

建议关注受益于 AI 算法进步, 并能成功进行商业化应用的海康威视、科大讯飞、商汤科技、中科创达等公司; 以及受益于 AI 算力需求、微调技术发展的海光信息、寒武纪、浪潮信息、海天瑞声等公司。

7. 风险提示

- 海外基础软硬件使用受限
若因国际关系等原因, 高算力 GPU 等基础硬件或计算框架等基础软件使用受限, 可能会对国内人工智能算法应用产生影响。
- 骨干网络创新放缓
目前 Transformer 成为深度学习骨干网络, 算法创新基本是基于 Transformer 做分支网络创新, 整体创新放缓。且 Transformer 本身作为骨干网络, 在处理部分任务时有一定局限性; 若骨干网络创新放缓, 可能部分任务解决进程会放缓。
- 应用落地不及预期
若相关应用公司不能找到人工智能算法较好的商业应用落地场景, 或相关场景客户没有较强的付费意愿, 可能算法应用落地会不及预期。

行业投资评级的说明：

- 买入：预期未来 3—6 个月内该行业上涨幅度超过大盘在 15%以上；
- 增持：预期未来 3—6 个月内该行业上涨幅度超过大盘在 5%—15%；
- 中性：预期未来 3—6 个月内该行业变动幅度相对大盘在 -5%—5%；
- 减持：预期未来 3—6 个月内该行业下跌幅度超过大盘在 5%以上。

特别声明：

国金证券股份有限公司经中国证券监督管理委员会批准，已具备证券投资咨询业务资格。

本报告版权归“国金证券股份有限公司”（以下简称“国金证券”）所有，未经事先书面授权，任何机构和个人均不得以任何方式对本报告的任何部分制作任何形式的复制、转发、转载、引用、修改、仿制、刊发，或以任何侵犯本公司版权的其他方式使用。经过书面授权的引用、刊发，需注明出处为“国金证券股份有限公司”，且不得对本报告进行任何有悖原意的删节和修改。

本报告的产生基于国金证券及其研究人员认为可信的公开资料或实地调研资料，但国金证券及其研究人员对这些信息的准确性和完整性不作任何保证。本报告反映撰写研究人员的不同设想、见解及分析方法，故本报告所载观点可能与其他类似研究报告的观点及市场实际情况不一致，国金证券不对使用本报告所包含的材料产生的任何直接或间接损失或与此有关的其他任何损失承担任何责任。且本报告中的资料、意见、预测均反映报告初次公开发布时的判断，在不作事先通知的情况下，可能会随时调整，亦可因使用不同假设和标准、采用不同观点和分析方法而与国金证券其它业务部门、单位或附属机构在制作类似的其他材料时所给出的意见不同或者相反。

本报告仅为参考之用，在任何地区均不应被视为买卖任何证券、金融工具的要约或要约邀请。本报告提及的任何证券或金融工具均可能含有重大的风险，可能不易变卖以及不适合所有投资者。本报告所提及的证券或金融工具的价格、价值及收益可能会受汇率影响而波动。过往的业绩并不能代表未来的表现。

客户应当考虑到国金证券存在可能影响本报告客观性的利益冲突，而不应视本报告为作出投资决策的唯一因素。证券研究报告是用于服务具备专业知识的投资者和投资顾问的专业产品，使用时必须经专业人士进行解读。国金证券建议获取报告人员应考虑本报告的任何意见或建议是否符合其特定状况，以及（若有必要）咨询独立投资顾问。报告本身、报告中的信息或所表达意见也不构成投资、法律、会计或税务的最终操作建议，国金证券不就报告中的内容对最终操作建议做出任何担保，在任何时候均不构成对任何人的个人推荐。

在法律允许的情况下，国金证券的关联机构可能会持有报告中涉及的公司所发行的证券并进行交易，并可能为这些公司正在提供或争取提供多种金融服务。

本报告并非意图发送、发布给在当地法律或监管规则下不允许向其发送、发布该研究报告的人员。国金证券并不因收件人收到本报告而视其为国金证券的客户。本报告对于收件人而言属高度机密，只有符合条件的收件人才能使用。根据《证券期货投资者适当性管理办法》，本报告仅供国金证券股份有限公司客户中风险评级高于 C3 级（含 C3 级）的投资者使用；本报告所包含的观点及建议并未考虑个别客户的特殊状况、目标或需要，不应被视为对特定客户关于特定证券或金融工具的建议或策略。对于本报告中提及的任何证券或金融工具，本报告的收件人须保持自身的独立判断。使用国金证券研究报告进行投资，遭受任何损失，国金证券不承担相关法律责任。

若国金证券以外的任何机构或个人发送本报告，则由该机构或个人为此发送行为承担全部责任。本报告不构成国金证券向发送本报告机构或个人的收件人提供投资建议，国金证券不为此承担任何责任。

此报告仅限于中国境内使用。国金证券版权所有，保留一切权利。

上海	北京	深圳
电话：021-60753903	电话：010-85950438	电话：0755-83831378
传真：021-61038200	邮箱：researchbj@gjzq.com.cn	传真：0755-83830558
邮箱：researchsh@gjzq.com.cn	邮编：100005	邮箱：researchsz@gjzq.com.cn
邮编：201204	地址：北京市东城区建内大街 26 号	邮编：518000
地址：上海浦东新区芳甸路 1088 号	新闻大厦 8 层南侧	地址：中国深圳市福田区中心四路 1-1 号
紫竹国际大厦 7 楼		嘉里建设广场 T3-2402