



# 面向AI大模型的网络使能技术

Network Enabling Technologies for Artificial  
Intelligence Large Models

## 目录

摘要 .....	3
<b>一、 AI 大模型发展概述 .....</b>	<b>4</b>
(一) 发展历程 .....	4
(二) 发展趋势 .....	5
<b>二、 网络使能大模型的需求和驱动力 .....</b>	<b>6</b>
(一) 未来 6G 网络的通算智融合趋势 .....	6
(二) 网络使能大模型价值场景 .....	7
<b>三、 网络使能大模型服务 .....</b>	<b>12</b>
(一) 数据感知服务 .....	13
(二) 分布式训练服务 .....	14
(三) 指令优化服务 .....	29
(四) 端边云协同推理服务 .....	30
(五) 模型优化服务 .....	36
<b>四、 案例分析 .....</b>	<b>37</b>
生成式 AI 在语义通信系统中的应用 .....	37
<b>五、 未来展望 .....</b>	<b>44</b>
<b>六、 参考文献 .....</b>	<b>45</b>
<b>七、 主要贡献单位和编写人员 .....</b>	<b>50</b>

## 摘要

随着大模型和智能体 ( Artificial intelligence agent, AI agent ) 技术的发展, 未来越来越多的工作将被基于大模型的智能体所取代。一方面, 由于大模型对数据和算力的需求巨大, 资源受限的终端将难以满足模型训练和推理的需求。另一方面, 未来第六代移动通信 ( Six generation, 6G ) 网络存在大量低时延需求的价值场景, 例如无人驾驶、虚拟和增强现实等, 云端大模型难以满足这些场景用户的需求。因此, 向无线网络寻求算力和数据的支撑将成为大模型时代的必然。本文介绍了大模型时代下网络使能人工智能 ( Artificial intelligence, AI ) 技术的需求和驱动力, 详细阐述了未来 6G 网络能为大模型提供的 AI 服务, 包括数据感知、分布式训练、指令优化、端边云协同推理和模型优化等, 通过案例分析说明了相关技术的实践应用, 并总结了未来可能的研究方向和所需要面对的挑战。

# 一、AI 大模型发展概述

## (一) 发展历程

随着深度学习技术的应用范围不断拓展和人工智能的快速发展，在大数据、高算力和强算法等关键技术的共同推动下，以 ChatGPT 为代表的 AI 大模型大量涌现，提供了高度智能化的人机交互体验和极富创造力的内容生成能力，改变了人们的工作和生活方式，实现了 AI 技术从“量变”到“质变”的跨越。

AI 大模型是指拥有超大规模参数、超强计算资源的机器学习模型，能够处理海量数据，并完成各种复杂任务。AI 大模型的发展可以追溯到 20 世纪 50 年代。此后，从卷积神经网络 (CNN) 到循环神经网络 (RNN)，再到 Transformer 架构，模型的性能不断提升。总的来说，AI 大模型的发展历程主要可以分为四个阶段，如图 1 所示。

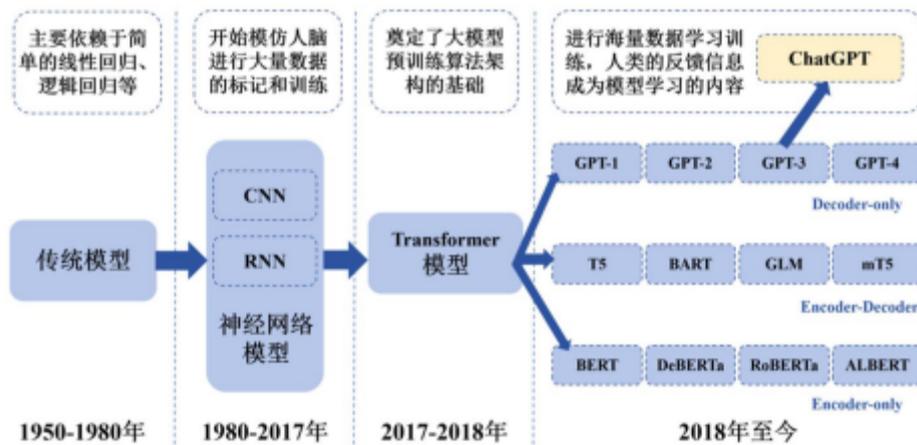


图 1. AI 大模型的发展历程

► **传统模型 ( 1950-1980 )** :在 AI 发展的早期，传统模型主要依赖于简单的线性回归、逻辑回归等方法。这些模型能够处理分类和回归等基本任务，但在处理复杂数据和任务时表现有限。

▶神经网络模型（1980-2017）：1980年，卷积神经网络的雏形 CNN 诞生。2000年代初期，有学者开始研究神经网络模型，开始模仿人脑进行大量数据的标记和训练，并尝试解决简单的问题，如手写数字识别等。

▶Transformer 模型（2017-2018）：2017年，Google 颠覆性地提出了基于自注意力机制的 Transformer架构，奠定了大模型预训练算法架构的基础。2018年，OpenAI 和 Google 分别发布了 GPT-1 与 BERT 大模型，使得 NLP 领域的大模型性能得到了质的飞跃。

▶现代 AI 大模型（2018 至今）：2022年，聊天机器人 ChatGPT 横空出世，迅速引爆互联网。此后发布的多模态预训练大模型 GPT-4，再次引发了生成式 AI 的热潮。目前各类大模型正持续涌现，性能也在不断提升。

## （二）发展趋势

### 1. 多模态能力提升，应用场景范围扩大

单模态模型通常只能处理一种类型的数据，例如文本、图像或声音，缺乏对复杂环境的全面理解。而具有多模态能力的 AI 模型能够同时处理多种类型的数据，例如将视觉和语言信息相结合，以实现更深层次的理解和交互，并在更广泛的场景中得到应用。

### 2. 模型轻量化部署，资源需求成本降低

在 AI 技术快速发展的当下，智能手机等移动设备在人机交互、语音交流等功能方面的需求不断提升，将大模型轻量化部署到终端设备也正成为一个重要的研究方向和发展趋势。利用端侧 AI 可以更好地为用户提供个性化的服务和支持，帮助用户进行自我管理，实现更加智能和高效的设备互联。

### **3. 外部工具相结合，交互方式更加智能**

传统的小模型通常专注于特定的任务，缺乏与外部环境交互的能力。结合外部工具调用、记忆和规划功能的 AI 大模型，可以被视为智能代理（Agent），它们能够执行更加复杂的任务，如自主决策、规划和学习。这种模型的交互方式更加智能，能够根据用户的需求和偏好进行自我调整，提供更加个性化的服务。

这些发展趋势不仅预示着 AI 技术的不断进步，也反映了用户对于更加智能、个性化服务的需求。随着研究的深入和技术的成熟，我们可以期待 AI 大模型在未来将在更多领域发挥关键作用，改善人们的生活和工作效率。

## **二、网络使能大模型的需求和驱动力**

### **（一）未来 6G 网络的通算智融合趋势**

人工智能已成为新一轮产业升级的核心驱动力，各行业的自动化、数字化、智能化需要泛在智能，许多高价值的 AI 场景，例如 AI 手机、自动驾驶、智能制造、移动机器人等，具有移动性、实时性、边端协同、隐私性等要求，需要网络这一 AI 服务基础设施进行支持。而随着大模型技术在上述场景中的深入应用，终端对于网络侧算力和数据资源支撑的需求将进一步扩大。

ITU 将 6G 场景扩展到包括通信与 AI 融合在内的智慧泛在，需要将 AI 打造成 6G 通信网络的新能力和新服务，实现 AI 即服务(AI as a service, AlaaS)。这要求 6G 网络能够随时随地提供 AI 服务、支持低时延的推理和训练服务、支持移动式 AI、保障 AI 服务质量、提供安全隐私保护。

## （二）网络使能大模型价值场景

### 1. AI 手机

在当今科技快速发展的时代，手机大模型已成为各大厂商竞相研发的热点。各大手机厂商纷纷推出了自家的大模型，为用户带来更加智能化的体验。如表 1 所示，手机大模型的功能主要包括文字类和图像类。在文字类功能方面，用户可以享受到智能问答、文本创作、文本总结、通话摘要等便捷服务，这些功能的响应时延通常在 1 秒之内，让用户感受到即时的互动体验。而图像类功能包括文生图、图像消除、图片问答等，其中，文生图响应时间较长，一般在 5 秒以上。在模型部署方面，目前主要有端侧部署和云端部署两种方式。端侧部署的大模型参数量通常不超过 10B，这种部署方式可以更好地保护用户隐私，同时降低对网络环境的依赖。而云端部署的大模型参数量可达 100B 以上，这种部署方式可以充分利用云端强大的计算资源，提供更加复杂和强大的功能，但需要较为稳定的网络环境支持。

表 1. 各厂商大模型手机调研信息

品牌	大模型功能	大模型性能	参数量	部署位置
vivo[1]	-智能问答	-智能问答首词响应 1s	1B/7B	端侧
	-文本创作 -文本总结 -逻辑推理	-文本总结首词响应 ms 级	70B/130B/175B	云端
OPPO[2]	-智能问答	-智能问答首字响应 0.2s	7B	端侧
	-通话摘要 -文本总结 -图像消除 -文生图	-512*512 生图时长 6s	70B/180B	云端

荣耀[3]	- 智慧成片 - 一语查图	暂无公布数据	7B	端侧
-------	------------------	--------	----	----

小米[4]	<ul style="list-style-type: none"> <li>- 智能问答</li> <li>- 文生图</li> <li>- 图片问答</li> <li>- 图像消除/扩图</li> </ul>	暂无公布数据	暂无公布数据	暂无公布数据
苹果[5]	<ul style="list-style-type: none"> <li>-智能问答</li> <li>-文本摘要</li> <li>-重要消息置顶</li> <li>-文生图</li> <li>-图像消除</li> <li>-跨应用操作</li> </ul>	暂无公布数据	暂无公布数据	端侧
				云端

基于上述分析，手机大模型主要分为终端推理和云端推理两类。因此，6G网络使能手机大模型也可以相应地分为使能终端推理和使能云端推理两类。

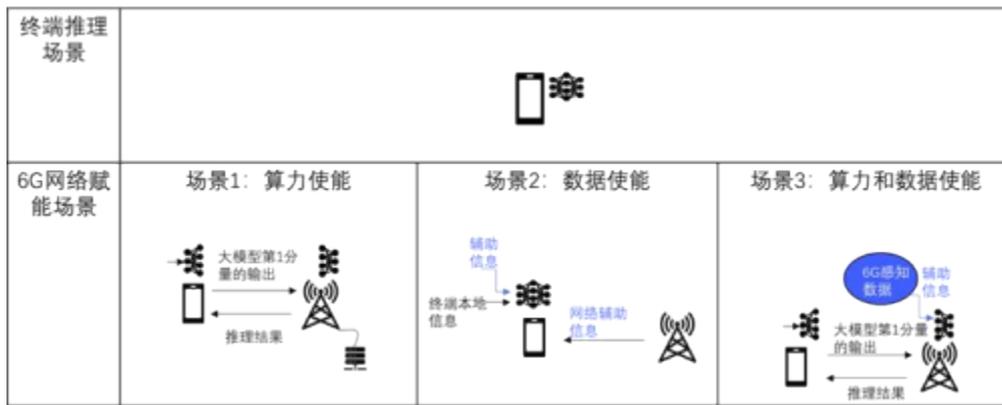


图 2. 6G 网络赋能大模型终端推理场景

如图 2 所示，6G 网络使能终端推理可以包括算力使能、数据使能以及算力和数据使能 3 种场景。考虑到目前手机大模型中文生图的时延较长的痛点，价值场景 1 是 6G 网络通过算力卸载的方式，将终端算力全部或部分卸载到 6G 网络内，通过对通信资源和算力资源的协同调度，可以降低响应时延，并降低终端推理功耗。而价值场景 2 则是 6G 网络通过例如感知获得价值数据，并将该价值数据作为终端推理的辅助信息，以提升推理精度。至于价值场景 3，则是网络同时提供算力和数据服务，从而可以降低终端推理的响应时延和功耗，并提升推理

准确度。



图 3. 6G 网络赋能大模型云端推理场景

如图3所示，6G网络使能云端推理也可以包括算力使能、数据使能以及算力和数据使能3种场景。在价值场景1中，6G网络通过算力卸载的方式，将云端算力全部或部分卸载到6G网络内，通过对通信资源和算力资源的协同调度，并通过更短的传输路径，可以显著降低响应时延，提升用户体验。而价值场景2则是6G网络通过例如感知获得价值数据，并将该价值数据作为云端推理的辅助信息，以提升推理精度。至于价值场景3，则是网络同时提供算力和数据服务，可以同时降低云端推理时延，并提升云端推理精度，为用户带来更加高效和智能的服务体验。

## 2. 自动驾驶

自动驾驶车辆通过传感器（如摄像头、雷达、LIDAR）采集到大量感知周围环境数据，需实时处理和分析、进行路径规划和驾驶决策。将连接的车云系统扩展到分布式网络节点/基站环境中，使数据 and 应用程序可以更靠近车辆，提供快速的道路侧相关功能。终端设备采集传感器数据，进行初步处理和特征提取。在车辆附近的分布式边缘节点进行实时数据处理，如环境感知和初步路径规划，利用6G网络的低延迟特性，快速传播危险警告和延迟敏感信息，确保实时响应。

在中央网络节点/云端进行大规模模型训练和全局优化，利用大数据提升模型的准确性和鲁棒性。根据车辆位置和网络状况，可动态调整分布式网络各节点计算资源，确保高效运行。

### **3. 智能医疗**

可实时监测患者的健康的医疗设备和穿戴设备收集大量患者体征数据，通过医疗大模型训练和推理，进行疾病预测和诊断。穿戴设备和医疗传感器采集生理数据，进行初步处理和传输。通过分布在医疗机构的边缘节点进行实时数据分析和初步诊断，减轻中央网络节点负担。中央网络节点进行复杂的医疗数据分析和模型训练，支持远程诊断和治疗方案的优化，通过高可靠性和低延迟的通信网络赋能医疗数据的实时传输和处理。除了实时传输能力和边缘节点部署能力，6G网络还提供了高可靠和加密的数据隐私保护机制，保障患者的数据隐私和安全。

### **4. 工业 4.0**

工业 4.0 要求智能工厂通过物联网设备进行设备监控、生产管理和质量控制，需要高精度、低延迟的数据传输和处理。工业传感器和设备采集生产数据，进行初步处理和传输。工厂内部的分布式网络节点部署计算，提供本地化的生产监控和实时优化能力，进行设备监控和故障预测。在中央网络节点进行大规模数据分析和模型训练，提升生产效率和产品质量。大带宽和低延迟的 6G 网络确保了生产数据在传感器、边缘网络节点及中央网络节点之间的实时传输和处理，高可靠性网络连接保障了生产过程的连续性和稳定性。

### **5. 工业元宇宙**

工业元宇宙打造与现实工业映射和交互的全数字化虚拟世界，构建工业全生命周期的虚实共生、相互操作及高效闭环的工业体系新范式，推动传统行业数字

化智能化转型，是新质生产力的数字底座。在工业元宇宙中，虚拟世界与物理世界的深度融合是实现其全部潜力的关键。虚拟世界不仅要能够感知和接收来自物理世界的数据，还需要能够理解这些数据背后的意图，并据此做出合理的决策和控制。这一过程中，大模型显著提升了工业元宇宙的智能化和自主化水平。

虚拟世界对物理世界的理解是工业元宇宙虚实交互的核心任务之一。工业元宇宙需要处理海量的数据，包括物联网设备传感器的数据、生产线监控信息、供应链的实时动态等。通过传统的规则模板解析、机器学习算法和深度学习算法，虚拟系统可以分析物理设备的数据并做出响应。然而，这些方法通常需要大量的规则和参数配置，灵活性较差。大模型的引入，尤其是基于大模型的生成式人工智能，使得意图识别和理解更加灵活和高效。大模型通过自然语言处理和深度学习技术，能够高效地解析、分析和处理这些数据。通过对海量文本、图像和其他数据的训练，能够在没有明确规则的情况下识别出复杂场景中的意图，将其转化为可执行的操作指令或预测性分析结果。例如，在一个智能工厂中，生产设备通过传感器反馈数据，虚拟系统不仅能够监测设备运行状态，还可以理解操作员的工作意图，从而调整生产部署等。

大模型在虚拟世界的构建过程中起到了加速器的作用。工业元宇宙的构建不仅仅依赖于物理世界的的数据输入，还需要大量的虚拟内容生成，诸如虚拟场景、产品设计、生产流程模拟等。文本、音频、视频等不同类型的的数据可以被自动生成，这极大地提升了虚拟世界的丰富性和细节表现。在产品设计过程中，大模型使能的设计软件可以生成大量的设计方案和模型，大幅缩短了产品设计周期，同时提高了创新性。不仅加速了产品的迭代，还能推动工业设计从传统的线性流程向智能化、迭代式的流程转变。虚拟现实（VR）和增强现实（AR）技术是工业元宇宙的重要组成部分，而构建真实感强、细节丰富的虚拟场景往往需要大量的人工干预和资源投入。通过大模型和 AIGC 技术，虚拟场景的生成可以更加自动

化、智能化，极大提升了开发效率。在一个虚拟工厂中，AIGC 可以基于物理工厂的布局自动生成相应的三维模型，并根据实时数据动态调整场景的布局和功能。这种虚实交互和自动化生成能力，提升了虚拟世界的沉浸感，使得企业能够更灵活地进行生产规划和调整。

决策和控制是工业元宇宙的核心之一，大模型的自主学习和决策能力提升了工业元宇宙的智能化水平。在工业生产过程中，生产环境和工艺流程通常非常复杂，需要根据实时数据动态调整。大模型可以基于大规模的数据训练，学习到各种复杂场景下的最优策略，并通过持续学习不断优化，使得工业元宇宙中的虚拟系统与物理世界紧密互动，优化资源分配，最终实现更高效工业部署和生产。

单一的大模型往往难以全面覆盖所有工业元宇宙场景需求，需要 AI 大模型与小模型融合，形成更全面的智能工业元宇宙系统。视觉引擎、语音引擎和机器人控制引擎等不同领域的 AI 小模型可以与大模型协同工作，补充其在特定任务中的不足，形成一个多功能的、全覆盖的 AI 使能的工业元宇宙系统，适应更加复杂多变的工业环境。

### 三、网络使能大模型服务

表 2 给出了网络使能大模型服务和一般 AI 模型服务的对比，可以看出 6G 网络使能的大模型服务在实时性、动态调整和高可靠性方面的显著优势，能够更好地满足不同应用场景的需求，提高系统的整体性能和用户体验。

表 2. 网络使能大模型服务和一般 AI 模型服务的对比

对比项	网络使能大模型服务	一般 AI 模型服务
带宽	处理的数据量更大，有更高的带宽需求：	数据传输量相对较小，带宽需求较低

<b>需求</b>			
	<b>实时性</b>	具有超低延迟的应用需求，在自动驾驶、实时视频处理等场景中，低延迟是关键	多数场景对实时性的要求较低，相对较高的延迟容忍性

能力	大带宽	提供更大传输带宽	现有资源带宽难以提升
	低延迟	利用 6G 网络的高带宽和低延迟，实现数据的实时传输和处理，支持实时分析和决策	主要依赖固定网络基础设施，数据传输和处理的实时性较低
	动态调度	利用智能调度系统，动态调整计算资源和任务分配	计算资源分配相对固定，难以动态调整和优化
	分布训练	广泛使用分布式数据并行和模型并行技术	通常在单个计算节点上完成训练和推理
	边缘计算	充分利用边缘计算能力，提高实时性和响应速度，减轻云端压力	边缘计算支持较少，主要依赖云端进行数据处理和模型训练

## （一）数据感知服务

在未来的 6G 网络时代，一个引人注目的趋势是网络中将会部署大量的传感设备。这些传感设备以其高灵敏度、高精度和高覆盖率的特点，将实现对物理世界的全面感知和实时监测。无论是环境参数的测量、人体健康指标的监测，还是物体位置和运动状态的追踪，传感设备都能提供详尽而准确的数据。与此同时，随着人工智能技术的飞速发展，大型模型在各个领域的应用越来越广泛。然而，这些大模型的训练和推断过程都需要大量的数据作为支撑。数据是模型学习的基石，是提升模型性能的关键。

在 6G 网络中，传感设备所收集的数据正好能够满足这一需求。在模型推断阶段，传感数据可以作为输入信息，增强大模型的推理精度。由于传感设备能够实时获取和传输数据，模型可以基于最新的数据进行推断，从而更准确地反映实际情况。这不仅提高了模型的实用性，还增强了其应对复杂和多变环境的能力。而在模型训练阶段，传感数据同样具有不可替代的价值。通过将传感数据作为训练数据的一部分，可以丰富模型的训练样本，提高模型的泛化能力。同时，利用

传感数据进行数据增强，还可以进一步提升模型的训练效果，使其在各种应用场景中都能表现出色。因此，未来 6G 网络中的传感设备将成为大型模型训练和推

断的重要数据源，为人工智能技术的发展提供强有力的支持。

## （二）分布式训练服务

### 1. 分布式机器学习理论

随着“大数据”概念的兴起，数据量爆炸式增长，数据和算法双驱动的模式逐渐受到工业界和企业界的重视。大数据的特征可以概括为大数据量、多类型、低价值密度和数据在线。其中，大数据量指的是数据集的规模非常庞大，通常达到TB甚至PB级别，这种规模的数据量远超传统数据处理工具和单机系统的处理能力，同时如此规模的数据不仅自身数据量庞大，并且存在大量非结构化数据（如图片、视频），需要复杂的数据处理和整合方法。数据在线是指数据实时更新和变化，大规模数据自身数据量庞大的同时自身数据增长的速度也非常快，存在大量衍生数据，需要实时的监测、分析和处理。大数据量和数据在线是使得传统机器学习不能适应当前环境的主要因素。

传统机器学习即在单机内进行数据处理和计算，注重在单机内处理数据的速度，由于内存和单机算力的限制，大数据条件下庞大的数据存储和计算是无法在单机中做到的，因此，将计算模型分布式地部署到多台、多类型的机器上进行同时计算是一种必要的发展趋势。

基于以上原因，提出了分布式机器学习的概念。分布式机器学习研究将具有大规模数据量和计算量的任务分布式地部署到多台机器上，其核心思想是“分而治之”，即将数据集或是计算任务分解成多个小数据集或计算任务，分配到不同的计算节点上处理，有效提高了大规模数据计算的速度并节省了开销。

分布式机器学习在概念的提出时就展现了独特的优势，包括针对大数据量问题的处理海量数据和针对数据在线问题的实时数据处理。在其发展和成熟的过程中，在其他的一些方面也展现了优势，首先，分布式架构支持动态扩展计算资源，

可以根据具体的计算需求的变化灵活地调整计算节点的数量和计算任务的分配，切薄系统的高效运行。其次，分布式系统的架构就保证了整个系统具有较强的鲁棒性，能够在某个节点发生故障时自动进行任务再分配，避免了计算过程的前功尽弃，保证了计算过程的稳定性和连续性。最后，分布式系统联合了大量低成本的硬件资源和计算资源去解决复杂的梯度计算，显著降低了能源和资源成本。

分布式机器学习分为面向扩展性的分布式机器学习和面向隐私保护的分布式机器学习，这种分类主要是针对传统机器学习不同限制因素进行改进。

### 1) 面向扩展性的分布式机器学习

在近年的研究中，训练的数据规模和模型参数规模以指数形式增长，以卷积神经网络为代表的神经网络使用大量训练数据训练一个参数为千万量级甚至上亿的模型，所使用的计算资源和所消耗的时间成本不是单机所能够做到的。面向扩展性的分布式机器学习是这种情况的一个可行的解决方案，它专注于将机器学习的算法扩展到多个计算单元（如 GPU），尝试将无论是廉价但低效的计算资源或是高昂但高效的计算资源纳入自身体系中，通过并行和分布式计算来处理大规模的数据集和复杂的模型。

面向扩展性的分布式机器学习主要通过数据并行、模型并行或是混合并行的策略实现它的处理任务，不同策略的使用则是基于不同的任务，有着各自的优势和缺陷，下面将一一介绍。

在分布式机器学习技术乃至大数据技术中，数据并行都是最为常见的一种并行方式。在数据并行的策略中，数据集被分割成若干个子数据集，并加载在若干个训练设备中（如 GPU）。因此，数据并行主要需要实现数据集分割以及训练后的模型参数同步两个部分的设计，其中后者是设计时关注的关键问题。数据并行是一种实现简单，扩展性好，且对于绝大多数深度学习任务都适用的可以应用于分布式机器学习的策略。然而，数据并行的通信成本很大，对于中心节点（参数

服务器)的通信压力也较大。

模型并行是基于单个节点无法容纳完整模型的问题所提出的，它将模型的不同层或不同参数分配到不同的节点上，每个节点只计算模型的一部分。需要频繁的设备间通信来传递中间结果。模型并行的提出和使用都是基于模型参数量庞大，例如卷积神经网络，但是它的实现较为复杂，通信开销较大。流水线并行（管道并行），通常认为是模型并行的一种特殊形式，它将模型按层或模块顺序切分成多个阶段，每个阶段分配到不同的计算节点上，形成流水线。管道并行通过分阶段处理和数据流动，减少了单个节点的内存占用，但是对比于数据并行，实现上较为复杂。

## 2) 面向隐私保护的分布式机器学习

面向隐私保护的分布式机器学习的主要目的则是保护用户隐私和数据安全，在面向隐私保护的 DML 中，数据的来源是多个参与方。有研究提到，在需要分布式机器学习技术来利用每个参与方的训练数据的时候，不同参与方的数据集可能具有不同的数据特征，所以实际中经常遇到的是训练数据的纵向划分，因此面向隐私保护的 DML 适用于具有纵向划分数据集的场景。

目前来看，其实不必纠结于面向隐私保护的 DML 适用于纵向划分的数据集还是横向划分的数据集，面向隐私保护的 DML 实际上提供了隐私保护的多条思路，在方法和实现上具有更大的覆盖范围，因此无论数据集的划分方式如何，都可以使用面向隐私保护的 DML 提供的隐私保护思路和方法。一方面，对于数据传输成本大的环境，如企业对用户（B2C），而数据敏感程度相对于政府部门和公司机密较低的情况，可以采用不直接共享数据的情况下进行机器学习模型的实现；另一方面，对于隐私要求严格，但对于传输环境安全要求不严格的情况，可以使用差分隐私等技术模糊化处理；如果对于隐私和数据安全都具有严格要求的情况下，如金融服务或是国家安全，可以采用同态加密和安全多方计算等技术。

从这个角度来说，联邦学习即是 DML 在隐私保护领域对于某一方面需求的继承与进一步发展。

## 2. 分布式机器学习平台和算法设计

分布式机器学习的主要研究方向包括分布式机器学习平台和分布式机器学习算法设计，在实际应用中，平台研究要结合算法的可行性，算法设计需要考虑在平台上的执行效果。

简单的例子是，一个由若干个计算节点和一个参数聚合服务器组成的分布式机器学习系统，每个计算节点都是一台机器，训练数据被分成若干个数据分片并发送给各个计算节点，计算节点在本地执行随机梯度下降算法，计算节点将梯度或者模型参数发送至参数服务器，参数服务器对收到的梯度或者模型参数进行聚合，从而得到全局梯度或者全局模型参数。

分布式机器学习的算法主要包括服务器如何将计算任务分配给每一个计算节点、计算节点在本地执行什么样的算法、参数服务器的全局参数如何聚合等。而分布式机器学习平台研究主要目的是搭建一个可以应用于多类型设备，搭载了分布式机器学习算法，并且具有优异性能的分布式深度学习框架。

### 1) 分布式机器学习平台

分布式机器学习平台研究起步的时间实际上较早。2005 年，Apache 实现了 Hadoop 分布式系统的基础架构。在经过接近 20 年的发展后，出现了大量成熟的分布式机器学习平台。分布式机器学习平台包括基于数据流模型、基于参数服务器、以及基于混合模型三类。

#### (1) 基于数据流模型

数据流模型通常把计算抽象成数据流图，关于数据流图，它是一种由节点和边组成的有向无环图。在数据流模型中，每个节点代表一个计算操作，边则表示

数据流动的方向，定义了数据处理的顺序和中间过程的依赖关系。在数据流图中存在源节点和汇节点，它们分别代表数据输入和数据输出。需要注意的是，源节点和汇节点并不必须是唯一的，每个源节点都可以代表一个独立的数据输入源，它们可以是数据集、数据流或是数据库，同样，每个汇节点也可以对应一个独立的外部存储或者其他系统。

数据流图的节点-边结构用于设计分布式机器学习平台是极为合适的，首先，数据流图支持并行处理，加载分布式学习算法后，划分的数据子集可以作为多个独立的源节点，每个计算单元可以作为数据流图上的一个或数个节点，大大提高了处理效率。此外，数据流模型使得分布式机器学习不仅在物理上是一个“分而治之”的系统，在计算逻辑上也成为了一个由若干个小的、可管理的节点组成的处理系统，从而由数据流模型实现的分布式机器学习平台具有较好的模块化属性和更高的可扩展性。

Spark 是一个具有代表性的基于数据流模型的分布式处理系统，虽然它主要被设计用来进行大规模的数据处理。它的一个关键特性是内存计算，将数据尽可能地存储在内存中进行计算，避免频繁的磁盘 I/O 操作，从而显著地提高了数据处理的效率，尤其适用于机器学习这样的迭代计算任务。

2010 年，Google 的研究人员最早提出关于参数服务器的概念，Google 在 2012 年发布了一个为大规模分布式训练深度神经网络设计的框架，即 DistBelief，它也是 Tensorflow 的前身。

## ( 2 ) 基于参数服务器

基于参数服务器的分布式机器学习平台的主要组成部分包括参数服务器和工作节点，参数服务器负责存储和管理模型的参数，管理工作节点的生命周期以及分配计算任务，工作节点则负责数据处理和梯度计算。参数服务器模型将参数划分给各个工作节点，提高了大规模参数模型训练的性能。

数据流模型和参数服务器模型的出现和发展不是继承关系，而是两种不同的设计思想，这两种设计思想各有利弊，因此衍生出了混合模型，集成了数据流模型的实时处理能力、高并行性以及参数服务器灵活参数更新策略、适合大规模机器学习的优势。简略地说，混合模型仍然采用节点代表计算操作，边代表节点之间的依赖关系，但是使用参数服务器的思想进行训练，将训练任务抽象成数据流图后，使用参数服务器进行任务调度，使用工作节点进行计算任务。

### (3) 基于混合模型

混合模型的代表分布式机器学习处理系统主要有 TensorFlow 和 PyTorch，它们都将网络模型的符号表达式抽象成计算图。

Google Brain 的团队在 DisBelief 的基础上研发了 TensorFlow。它将数据流和参数服务器搭配使用，从而取得了更快的速度、更高的移植性和灵活性。早期的 TensorFlow 使用的是静态计算图，这种方式在优化和部署时会具有一定的优势，后续 TensorFlow 引入了 Eager Execution，从而使得默认情况下计算图时动态的，便于机器学习的调试和开发。从 2015 年发布至今，TensorFlow 进行了多次更新，支持在各类环境下执行分布式机器学习程序，包括移动设备、Windows 和 CPU、GPU。

Pytorch 则更加适合于小规模的项目，它使用了动态计算图。它的主要特点是：使用了类似于 Numpy 的 N 维 tensor，从而在 GPU 加速上取得了杰出的成就。其次，它使用的自动微分方法可快速构建和训练神经网络，在前向和后向传播都取得了较好的效果。最后，动态计算图的使用可以加速模型收敛，便于将计算分布式地部署在 GPU 和其他机器上。

## 2) 分布式机器学习算法设计

为了实现分布式机器学习的具体应用，通常需要在机器学习平台上加载分布式机器学习算法。然而机器学习算法已经是一个十分成熟的门类，按照学习方式

可以分为监督学习、半监督学习、非监督学习和强化学习，因此，将机器学习算法移植到分布式机器学习上不仅关注机器学习算法的分布式实现，而且关注分布式的机器学习算法针对通信延迟、一致性和容错性问题的优化。

无论是监督学习、半监督学习、非监督学习还是强化学习，都可以移植到分布式机器学习平台上，如分布式梯度下降、分布式聚类、分布式强化学习。不同的机器学习算法进行分布式实现时需要考虑的优化方面不同，主要可以分为数据分割和预处理（如果存在数据集的话）、计算负载平衡、扩展性和可维护性考虑。而几乎所有分布式机器学习都需要考虑通信成本和一致性问题。一致性问题的解决策略关注模型的参数更新方式，包括同步更新和异步更新，同步更新可以保证算法的收敛率，但由于计算资源的层次不齐，由于短板效应的存在导致资源利用率较低；异步更新的资源利用率很高，但无法保证收敛效果。

大模型能够捕捉复杂模式，提供全面的知识，在各类任务中展现出卓越的性能。然而，如何在海量数据和复杂任务中实现自我优化和演进成为进一步提高大模型泛化能力的关键问题。这一过程不仅依赖于初始的预训练数据，还需要在实际应用中不断学习和适应来自用户的新数据与任务。分布式学习由于可以充分利用分布式算力完成大规模机器学习模型的协作训练，成为高效学习和适应新数据与新任务的有效手段。在分布式学习架构下，移动设备可以部署为网络中参与学习任务的计算节点甚至是智能体在本地执行训练和推理任务，并通过交互中间训练结果来完成大模型的全局微调训练，这样不仅可以增强算力，加速数据处理，还可以驱动传统人工智能服务、AI 生成服务等网络智能化任务。然而，分布式架构下全参数训练过程中需要频繁进行参数更新、参数同步和梯度交换，这对网络传输、计算、存储等能力提出极高的要求。一个有效的解决思路是借助高效微调技术（比如高效参数微调，提示微调）来实现大模型的分布式微调训练。比如利用参数高效微调技术，用户可以训练额外小规模微调模型（如 LoRA）捕捉

其本地任务或数据相关的知识，在保持接近全参数训练性能的同时，大幅降低了计算和存储的需求，而共享和聚合小规模微调模型对网络来说也相对资源友好。同时，在分布式网络系统中，还需要考虑到每个用户的计算资源和存储能力都有所不同，例如，一些用户可能使用高性能的服务器或专业的计算设备，而其他用户可能只能访问基本的智能手机或小型边缘设备。不同用户在本地进行大模型微调时可承受的微调模型的参数量也会有所不同，直接导致了在不同用户间微调模型的异构性。因此，还需进一步探索端侧资源异构性场景下高效的分布式训练服务的解决方案，例如微调模型的知识蒸馏、异构微调模型对齐等。此外，在分布式网络中提供分布式训练服务的时候还需要关注用户的一致性问题，针对不同网络场景设计高效的共识机制、分布式一致性算法，在训练过程中需要确保用户模型的同步和一致性，从而保证大模型可以正确有效的收敛。（ 电子科技大学：刘贻静，汪云翔）

### 3. 分布式训练框架

分布式训练的关键技术包括云边协同计算和分布式训练框架的实现：

- > 分布式协同计算：分布式网络节点在靠近数据源的地方进行数据处理和分析，可以减少传输延迟和带宽占用，提高实时性和响应速度，适用于需要快速决策的场景。中央网络节点则借助强大的计算和存储资源，进行复杂的数据处理和模型训练，支持大规模数据分析和全局优化，提升整体系统性能。分布式网络节点和中央网络节点协同通过在分布式网络节点进行初步的数据处理和训练任务，减轻中央网络节点压力；利用 6G 网络的低延迟特点，实现分布式节点与中央网络节点的高效协同，提高整体效率。
- > 分布式训练框架：分布式训练框架包括并行计算和智能调度，其中并行计算主要方法为数据并行和模型并行，共享网络中的计算节点实现。

数据并行技术将训练数据集分割成多个子集，每个子集在不同的计算节点上独立训练，这些节点共享同一个模型副本，但独立计算梯度。节点间同步更新模型参数，利用 6G 网络的高带宽进行快速数据传输和梯度同步。数据并行通过以下过程实现：1) 数据分割：将训练数据集分成若干子集，分配到不同的计算节点；2) 独立训练：每个节点使用自己的数据子集进行前向和后向传播，计算梯度；3) 梯度汇聚：所有节点的梯度通过网络进行汇聚（例如使用参数服务器或全局同步），然后更新模型参数。数据并行能够处理非常大的数据集，并且适用于大多数深度学习框架，但通信开销较大，尤其是在节点数量很多时，会导致同步瓶颈。

模型并行技术将模型本身分割成多个部分，不同的计算节点负责不同部分的计算，节点间传递中间结果。适用于超大规模模型，利用 6G 网络的低延迟进行高效通信。这对于超大规模的模型特别有效。具体实现上，模型分割将模型的各个层或模块分配到不同的计算节点；通过并行计算，每个节点只负责自己部分的前向和后向计算，节点之间通过网络传递中间结果和梯度信息。通过模型并行可以处理单个节点无法容纳的大模型，但需要仔细规划模型的分割策略和节点间的通信，以减少延迟和通信开销。

智能调度技术则动态调整计算资源和任务分配，根据网络状况和计算需求，优化分布式训练过程。利用 AI 算法动态调整计算资源和任务分配，从而提高资源利用率和训练效率。如图4 所示，分布式训练服务的部署和内涵应从如下步骤进行考虑：

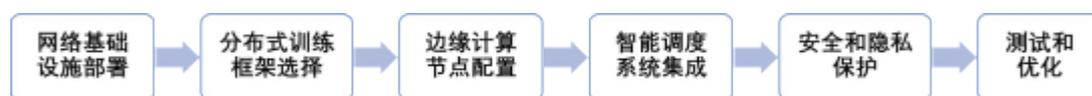


图 4. 分布式训练服务部署步骤

- > 部署 6G 网络基础设施，确保高带宽、低延迟和高可靠性的网络环境。
- > 选择适合的分布式训练框架，支持数据并行和模型并行。

- 部署边缘计算节点，配置高性能计算和存储资源，确保边缘节点具备足够的处理能力。
- 集成智能调度系统，动态调整计算资源和任务分配，优化训练过程。
- 实施数据加密、访问控制等措施，确保分布式训练过程中的数据安全和隐私保护。
- 进行全面的测试和优化，确保分布式训练系统的性能和稳定性，满足实际应用需求。

## 4. 联邦学习

根据 scaling law，大模型的性能是和模型参数、数据大小、计算量成正比的。6G 网络使能大模型分布式训练的优势包括宝贵的数据源、海量的闲置算力。其中，联邦学习是一种实现网络使能大模型分布式训练的隐私保护范例。通常支持联邦学习的无线网络由多个本地客户端和一个边缘服务器组成。联邦学习进行分布式训练包括如下五个迭代步骤：

1) 全局模型初始化：在中央服务器上初始化一个全局模型 $w_0$ ，并将其分发给  $K$  个参与训练的客户端。

2) 本地模型训练：每个客户端根据本地数据进行模型训练，计算更新后的模型参数。对每个客户端  $K$ ，训练的目标是最小化其本地损失函数：
$$J_k(w) = \frac{1}{n_k} \sum_{i=1}^{n_k} f(w; x_i, y_i)$$
。其中， $(x_i, y_i)$  是客户端  $k$  上的数据样本， $f(\cdot)$  是损失函数， $w$  是模型参数， $n_k$  是客户端  $k$  的本地数据量。

3) 本地模型上传：客户端在完成本地训练后，将模型参数上传到中央服务器。此时，每个客户端  $k$  提交的模型可以表示为  $w_k$ 。

4) 全局模型聚合：中央服务器接收每个客户端上传的模型参数后，进行加权聚合来更新全局模型。最常用的聚合方法是 FedAvg，其公式为： $w' =$

$\sum_{k=1}^K \frac{n_k}{n} w_k$ 。其中， $n$ 是所有客户端的数据总量， $K$ 是参与训练的客户端数， $w'$ 是更新后的全局模型参数。

5) 模型迭代：服务器将新的全局模型 $w'$ 分发给所有客户端，重复上述过程，直到模型收敛。

结合上述的训练过程以及不同的客户端设备，联邦学习的具体训练场景分跨孤岛 ( cross-silo ) 和跨设备 ( cross-device ) 两种：

1) 跨孤岛：每个参与的孤岛 ( silo ) 通常是一个具有相对较强计算能力的本地计算实体，比如数据中心、公司内部的服务器、企业专有的计算资源等。通过跨多个这种具有独立计算能力的单位或组织进行协作式分布式学习，这些单位或组织之间不会直接交换数据，而是通过本地模型的更新与全局模型的聚合来完成训练。

2) 跨设备：参与的设备通常是功能有限的智能设备，例如手机、IoT ( 物联网 ) 设备等，计算能力不如 silo 中的服务器强大，这些设备之间通过无线网络进行协作，但网络连接可能较为不稳定。

因此，将大模型训练集成到支持联邦学习的无线网络之前，我们必须考虑不同联邦学习场景所施加的限制与大模型计算/存储/通信密集型要求之间的冲突<sup>[6]</sup>，所面临的主要挑战包括：

1) 高功耗：由于训练所需的数据和模型参数数量庞大，大模型的训练过程对计算硬件和能耗都有很大的要求。当以合理可持续的方式部署高耗能的大模型时，能源效率变得至关重要。

2) 有限且异质的算力资源：无线网络出了提供分布式训练服务之外，还需要承担基础的连接任务，硬件算力资源有限。为了将集成到无线网络中，通常需要专门的硬件来进行 AI 计算加速，例如 GPU、TPU 等。并且，客户端之间计算

资源的异质性，使得联邦训练大模型遭受更多空闲时间的影响，需要对算力资源建立统筹有效的编排和调度机制。

3) 高存储和内存要求：为了满足大模型训练要求，必须大幅增加存储和内存来处理流式收集/生成的数据以及训练期间模型参数的更新。典型的网络架构可能不一定满足此类存储和内存要求。

4) 高通信开销和时延：由于模型规模巨大，基于联邦学习的从头训练需要持续、大量的通信资源，这一过程将非常耗时且占用带宽。例如，通过 100Mbps 通道（5G 中用户体验的数据速率）传输一次 GPT2-XL（约 5.8 GB 的中型 LLM）大约需要 470 秒<sup>[7]</sup>，而从头训练可能需要对数千个 GPU 进行长达几个月的连续训练。虽然 5G 及以上网络有严格的延迟要求。目前尚不清楚将大模型集成到支持联邦学习的网络中如何满足如此严格的延迟要求。

5) 数据问题：针对各客户端数据特征分布（非独立同分布，Non-IID）、标签分布以及样本量等方面的差异，微调大模型时必须设计有效的策略来应对这种数据异构性；不同客户端数据质量的差异（含噪声、标签错误或者缺失值等）会直接影响大模型的微调效果，一些必要的数据预处理如数据清洗、异常检测等步骤也需要进一步的考虑。

在以上挑战中，最关键的问题在于联邦学习设备算力有限条件下如何降低通信开销和优化内存。

降低通信开销方面的研究包括在传输层面以及在算法层面的优化两种：传输层面可以采用参数高效 PEFT 方法<sup>[8]</sup>减少需要传输参数量，如 FedPETuning<sup>[9]</sup>；降低通信轮次，如 DiLoCo<sup>[10]</sup>；压缩；量化等；算法层面可以采用更高效的联邦聚合算法，如 OpenFedLLM<sup>[11]</sup>，或者传输内容更少的其他优化算法，如 FwdLLM<sup>[12]</sup>、FedKSeed<sup>[13]</sup>用零阶有限差分估计梯度。

内存优化的维度包括采用参数高效 PEFT 方法减少可训练参数量；使用梯度

累积、激活值重计算等技术；优化器的选择，SGD、带动量的 SGD、Adam、AdamW 等；结合模型分割 Split learning 技术<sup>[14]</sup>，但代价是会加剧通信开销；只训练部分层/层冻结，如 AutoFreeze<sup>[15]</sup>、SmartFRZ<sup>[16]</sup>、FedOT/FedOST<sup>[17]</sup> 等。此外，针对内存的优化还可以采用混合精度训练、ZeRO 零冗余优化器<sup>[18]</sup> 等技术。

其次，对于算力异质性的研究，FATE-LLM<sup>[19]</sup>提出很多种架构的可能性，允许算力异质性允许终端使用不同的本地模型，知识蒸馏出一个通用于联邦聚合的模型。FedIT<sup>[20]</sup>提出每个设备可以采用不同的 Lora 配置，即层级最优秩 (Layer-wise Optimal Rank Adaptation) 思想。

## 5. 联邦大模型

大模型的参数规模极为庞大，且各大厂商也在持续刷新大模型参数量的上限。以 GPT 系列为例，从 GPT-1 到 GPT-4，模型的参数量从 1.1 亿增长至 1.8 万亿，由模型规模带来的性能提升出现边际递减效应<sup>[8]</sup>。目前，针对特定任务的模型微调 (Fine-tuning, FT) 已成为利用大模型的主要方法<sup>[21]</sup>，但直接微调对算力、内存都提出了更高的要求，AI 硬件 (GPU) 内存难以跟上模型扩大的需求。为了解决算力增速不足的问题，研究者考虑用多节点集群进行分布式训练，将训练扩展到多个 AI 硬件上 (如 GPU)，从而突破于单个硬件内存容量和带宽的限制，支撑更大规模模型的训练。目前较多分布式训练架构主要有两种模式：集合通信 (collective communication, CC) 模式和参数服务器 (parameter server, PS) 模式。

NLP、CV、多模态、科学计算等领域的模型往往具有模型结构复杂、参数稠密的特点，集合通信训练模式可以很好地支持此类模型的训练。集合通信模式对计算芯片的算力和芯片之间的网络互联要求较高，如高性能计算的 GPU、芯片

之间的高速网络互联 NVLink 和 InfiniBand 等。搜索、推荐等场景的模型往往数据量巨大，特征维度高且高度稀疏化。参数服务器训练模式可以很好地支持此类模型的训练。参数服务器训练模式可以同时做到对数据和模型的并行训练，对于存储超大规模模型参数的训练场景十分友好，常被用于训练拥有海量稀疏参数的搜索、推荐领域模型。

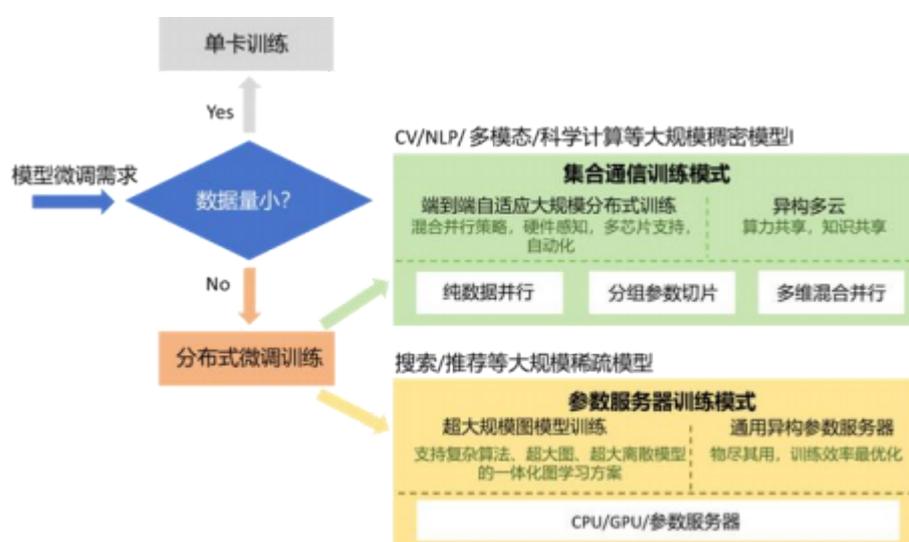


图 5. 分布式微调训练

然而，传统公开的可用数据集无法满足大模型微调的需求<sup>[22]</sup>，特别是大规模中文数据集十分缺乏，对中文大模型以及业界模型的中文支持都有很大的影响。收集多样化、高质量的指令数据仍面临挑战，尤其是在隐私敏感领域，往往禁止收集、融合使用数据到不同的地方进行 AI 处理，本地数据不足或微调和预训练数据集之间存在显著差异可能导致模型的泛化性能不佳。

为了解决隐私数据给用户安全和模型泛化性能带来的挑战，联邦学习（Federated learning，FL）<sup>[22]</sup>作为一种分布式框架被引入。其中，联邦分割学习（Federated split learning，FSL）框架<sup>[23]</sup>将模型分割成多个部分，在边缘用户设备上仅针对部分模型基于本地任务数据进行训练，训练完成后上传模

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。  
如要下载或阅读全文，请访问：

<https://d.book118.com/778004106010007007>