

中国算力大会

算 赋 百 业 · 力 导 未 来

中国算力白皮书 (2022年)

2022中国算力大会
2022年7月

前 言

2020年3月4日，中共中央政治局常务委员会召开会议提出，“加快5G网络、数据中心等新型基础设施建设进度”，将数据中心纳入“新基建”范畴。同年4月20日，国家发改委明确新型基础设施的范围，数据中心作为算力基础设施，成为信息基础设施重要组成部分。

算力作为大数据、云计算、区块链、人工智能以及网络安全等新兴数字产业快速发展的重要支撑，与互联网行业、电子信息制造业、数字技术服务业等数字核心产业的发展息息相关。《关于加快构建全国一体化大数据中心协同创新体系的指导意见》、《全国一体化大数据中心协同创新体系算力枢纽实施方案》、《新型数据中心发展三年行动计划（2021-2023年）》等文件更是明确了发展算力的重要意义。

为梳理算力规模、算力市场和技术现状、构建数据中心算力评价指标体系，中国信息通信研究院联合业界专家基于前期研究成果编制了《中国算力白皮书（2022年）》。本白皮书聚焦国内外算力规模，通用算力、智能算力、超算算力和边缘算力的市场、技术等内容。同时，白皮书提出“算力五力模型”，为数据中心单体算力评估体系提供新的模型和方法，更好地指导和建议业界判断行业发展趋势，为未来算力规划和部署提供思路。

如对白皮书有建议或意见，请联系：dceco@caict.ac.cn；也可关注公众号“CAICT算力”获取更多信息。

版权说明

本白皮书版权属于2022中国算力大会，并受法律保护。转载、摘编或利用其它方式使用本白皮书文字或者观点的，应注明“来源：《中国算力白皮书（2022年）》”。违反上述声明者，将追究其相关法律责任。



目 录

1	概述	1
2	算力总体情况	2
2.1	全球算力总体情况	2
2.2	我国算力总体情况	2
2.3	我国算力行业应用分布	3
3	算力的分类	4
3.1	综述	4
3.1.1	市场规模方面	4
3.1.2	技术方面	4
3.2	通用算力	4
3.2.1	概述	4
3.2.2	市场	5
3.2.3	技术	6
3.3	智能算力	14
3.3.1	概述	14
3.3.2	市场	14
3.3.3	技术	16
3.4	超算算力	20
3.4.1	概述	20
3.4.2	市场	22
3.4.3	技术	26
3.5	边缘算力	29
3.5.1	概述	29
3.5.2	市场	30
3.5.3	技术	32
4	算力应用赋能	34
4.1	绿色电力	34
4.1.1	简介	34
4.1.2	案例	34
4.2	人工智能	35

4.2.1 简介	35
4.2.2 案例	35
4.3 车联网	35
4.3.1 简介	35
4.3.2 案例	36
4.4 智慧医疗	36
4.4.1 简介	36
4.4.2 案例	36
4.5 边缘计算	37
4.5.1 简介	37
4.5.2 案例	37
4.6 物联网	39
4.6.1 简介	39
4.6.2 案例	39
5 算力五力衡量体系	40
5.1 算力五力指标构建	40
5.2 通用算力	40
5.3 智能算力	40
5.4 算效	41
5.5 网络	41
5.6 存储	45
6 算力五力模型	47
6.1 相关概念	47
6.2 双向投影模型	49
6.3 算力五力模型	51
6.4 算例	52
7 结语和展望	58
7.1 算力规模方面	58
7.2 算力技术方面	58
7.3 算力指标构建方面	58

图目录

图 1	2021 年全球算力规模情况	2
图 2	2021 年我国算力规模情况	3
图 3	2021 年我国算力行业应用分布情况	3
图 4	近四年 CPU 服务器处理器出货量规模	5
图 5	2021 年第 4 季度 CPU 市场份额	6
图 6	2021 年第 4 季度 GPU 市场份额	15
图 7	FPGA 市场规模统计	16
图 8	2021 年第 4 季度 FPGA 市场份额	16
图 9	2022 年 6 月 TOP500 中部分国家超算数量	23
图 10	2022 年 6 月 TOP500 中各国超算所占份额	23
图 11	2022 年 6 月 TOP500 榜单对应厂商份额	25
图 12	2002-2021 我国 TOP100 超级算力数据	25
图 13	2021 年 11 月 TOP100 厂商系统数	26
图 14	超算核心硬件架构	27
图 15	超算网络架构	28
图 16	对边缘计算的解构	30
图 17	2024 年各边缘阶段预期市场份额	31
图 18	边缘设备发展对比图	31
图 19	阿里巴巴世纪互联南通 A 栋数据中心绿色等级评估 5A	34
图 20	边缘计算服务平台功能架构图	38
图 21	云边协同下核心多活+边缘分布示意图	38
图 22	数据中心算力评估指标构建	40
图 23	直连拓扑 vs.CLOS 时延测试	43
图 24	直连拓扑 vs.CLOS OpenFoam 性能测试	43
图 25	济南市内 100G/27KM RoCE 网络平均时延	44
图 26	济南-淄博 10G/310KM RoCE 网络平均时延	45
图 27	$prj_{X_i X^+}(X^-X^+)$ 和 $prj_{X^-X^+}(X^-X_i)$ 的图形表示	50
图 28	每个数据中心各指标下到正负理想解的距离	56
图 29	样本数据到正、负理想解的距离	57

表目录

表 1 CPU 各类架构概览.....	5
表 2 2021 年 Intel 部分处理器产品.....	7
表 3 2021 年 AMD EPYC（霄龙）部分产品概况.....	7
表 4 ARM 3 款产品对比.....	8
表 5 Ampere Altra Max 产品规格.....	8
表 6 亚马逊 Graviton2 产品规格.....	9
表 7 飞腾 FT-1500A、FT2000、S2500 产品规格.....	9
表 8 华为 Hi1620、Hi1616、Hi1612、Hi1610 产品规格.....	10
表 9 龙芯 3A5000/3B5000、龙芯 3C5000L 产品规格.....	11
表 10 龙芯 7A1000、龙芯 7A2000 产品规格.....	12
表 11 IBM 的 Power1990-2019 年部分产品概况.....	12
表 12 IBM 的 Power10 产品概况.....	13
表 13 阿里的玄铁框架性能表.....	13
表 14 申威 SW2、SW-3/SW1600 框架性能表.....	14
表 15 智能算力主要芯片类型.....	14
表 16 全球服务器 GPU 市场季度统计（单位：\$M）.....	15
表 17 AMD 部分产品概况.....	17
表 18 2010-2021 年 Xilinx 部分 FPGA 产品.....	18
表 19 2010-2021 年 Intel（Altera）部分 FPGA 产品.....	18
表 20 2010-2021 年部分国产厂商 FPGA 产品.....	19
表 21 13 个并行计算关键领域.....	21
表 22 全球超算发展历程.....	22
表 23 1993-2022 年世界排名第一超算.....	24
表 24 MEC 边缘云部署应用软件系统与业务终端配置.....	37
表 25 智能算力领域算力精度.....	41
表 26 算力五力指标选取.....	51
表 27 样本数据.....	52
表 28 数据中心能效指标分级.....	52
表 29 数据中心能效指标定级.....	53
表 30 不同数据中心不同指标间的两两比较.....	53
表 31 数据中心等级比较级经优序数转化表.....	53
表 32 相对贴进度定星标准表.....	57

1 概述

算力是数据中心的服务器通过对数据进行处理后实现结果输出的一种能力，最常用的计量单位是每秒执行的浮点运算次数（FLOPS）。现有相关算力规模基本都是基于单精度浮点算力次数（FP32）来进行公布的。例如，2021 年 7 月，工业和信息化部印发的《新型数据中心发展三年行动计划（2021–2023 年）》中提出算力发展目标：到 2021 年底，全国数据中心总算力超过 120 EFLOPS（1 EFLOPS=10¹⁸FLOPS, 每秒百亿亿次浮点运算）；到 2023 年底，全国数据中心总算力超过 200 EFLOPS¹。2021 年 11 月，工业和信息化部发布的《“十四五”信息通信行业发展规划》提出数据中心算力到 2025 年要增长到 300 EFLOPS。

算力的分类主要包括通用算力、智能算力、超算算力、边缘算力四部分。通用算力以 CPU 芯片输出的计算能力为主；智能算力以 GPU、FPGA、AI 芯片等输出的人工智能计算能力为主；超算算力主要以超级计算机输出的计算能力为主；而边缘算力主要以就近为用户提供的实时计算能力为主，是以上三种算力形式的组合。

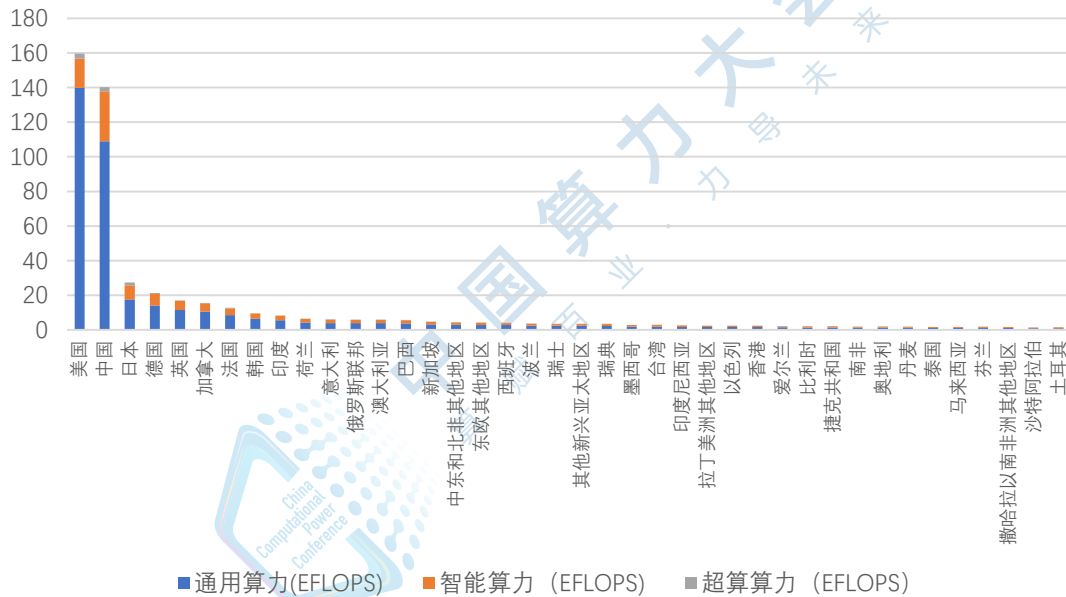


¹ 《新型数据中心发展三年行动计划（2021–2023 年）》，工信部通信〔2021〕76 号

2 算力总体情况

2.1 全球算力总体情况

在算力规模方面²，截止到2021年底，全球算力总规模达到521 EFLOPS（FP32）。其中，通用算力为398 EFLOPS（FP32），智能算力为113 EFLOPS（FP32），超算算力规模为10 EFLOPS（FP32）。美国与中国算力能力位列前两名，美国算力总规模为160 EFLOPS（FP32），中国算力总规模为140 EFLOPS（FP32）。总算力份额前五名国家，美国、中国、日本、德国、英国分别占比31%、27%、5%、4%、3%，共占据全世界算力70%的份额。



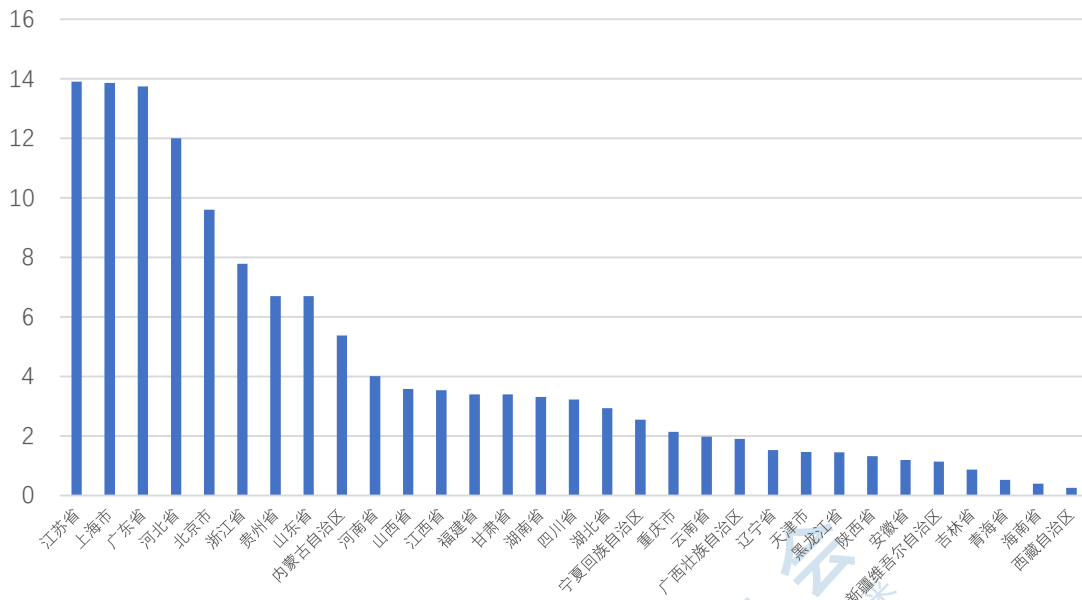
数据来源：Gartner、IDC、TOP500，中国信息通信研究院整理

图1 2021年全球算力规模情况

2.2 我国算力总体情况

在算力规模方面，截止到2021年底，我国算力总规模为140 EFLOPS（FP32），算力规模排名全球第二。其中，通用算力规模为109 EFLOPS（FP32），智能算力规模为29 EFLOPS（FP32），超算算力规模为2 EFLOPS（FP32）。江苏、上海、广东、河北算力规模皆超过10 EFLOPS（FP32）。北京、浙江、贵州、山东、内蒙古算力规模皆超过5 EFLOPS（FP32）。

² 算力规模部分包含通用算力、智能算力、超算算力，边缘算力暂未纳入统计范围。

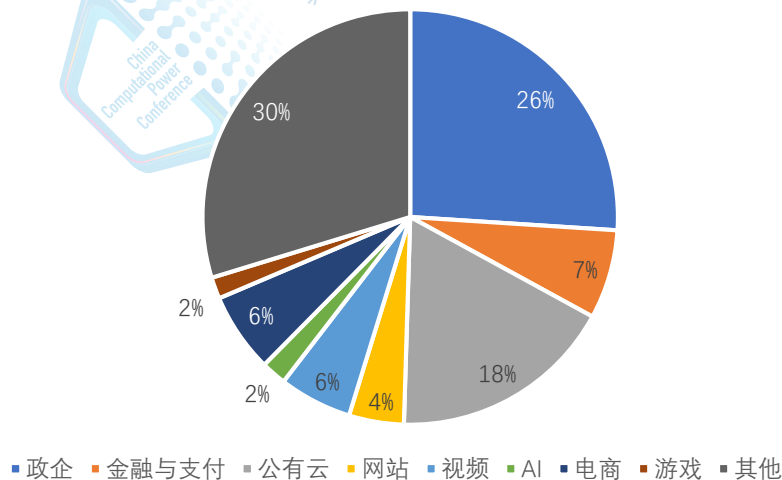


数据来源：中国信息通信研究院整理

图2 2021年我国算力规模情况

2.3 我国算力行业应用分布

截止到2021年底，我国算力行业应用分布主要为互联网、政企、金融、其他，分别占比47%、26%、7%、20%左右³。其中互联网主要可细分为公有云、网站、视频、AI、电商、游戏等领域，分别占比18%、4%、6%、2%、6%、2%。



来源：中国信息通信研究院整理

图3 2021年我国算力行业应用分布情况⁴

³ 来源：中国信息通信研究院测算

⁴ 图3中的“其他”为互联网及非互联网类“其他”的总和。

3 算力的分类

3.1 综述

随着数字经济时代的到来，算力的发展也迎来了高潮。算力分为通用算力、智能算力、超算算力和边缘算力四种，这四种算力无论是在市场规模方面还是在技术方面均有了较大的提升。

3.1.1 市场规模方面

通用、智能算力市场垄断效应明显，边缘算力各厂商加快布局。根据 IDC 官方数据，通用算力市场由 Intel 和 AMD 两家厂商占据市场近 95% 份额，CPU 出货量同比增长超过 10%。智能算力 GPU 市场份额几乎被英伟达垄断，占据 95% 以上市场份额。GPU 市场收入超过 70 亿美元，FPGA 市场收入虽然不及 GPU 市场收入，但是 FPGA 市场前景较好，未来 FPGA 市场收入会有更大的上升空间。超算算力市场几乎被中国、美国、日本所占据，随着高性能算力的发展，未来超算算力将成为各国竞争的新焦点。边缘算力作为新兴技术将会应用到更多的行业中，满足各行业低时延的需求，未来边缘算力市场发展会有很大的提升空间。

3.1.2 技术方面

技术进一步突破，芯片制程与芯片性能不断提高。随着数据爆发式的增长，各行业对数据存储、数据计算、数据分析提出了更高的要求。算力在技术层面不断突破技术壁垒。通用算力 CPU 芯片制程不断提升，处理器架构不断创新，实现高性能低功耗。智能算力 GPU、FPGA 芯片性能不断提高，芯片内存进一步提升。超算算力性能不断创新，运算速度越来越快。边缘算力在低时延、高带宽等技术方面不断创新，满足不同应用场景的不同精度需求，实现云侧和边侧协同发展。

3.2 通用算力

3.2.1 概述

通用算力主要以 CPU（Central Processing Unit，中央处理器）为代表。指令集是存储在 CPU 处理器内部，对 CPU 运算进行指导和优化的硬程序⁵，CPU 芯片的运行依靠和执行的就是指令集。按指令集架构而言，CPU 分为 x86 架构与非 x86 架构，主要有 x86、ARM、MIPS、Power、RISC-V、Alpha 等。

CPU 芯片行业技术壁垒高，国内外仅有少量企业能够稳定提供产品。国外代表厂商有 Intel、

⁵ 《AI 服务器白皮书》，开放数据中心委员会 ODCC，2019 年

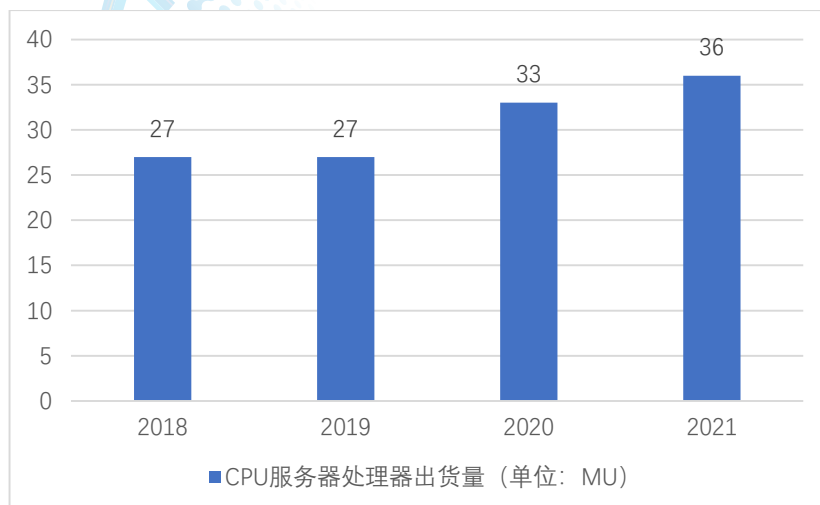
AMD、IBM、安谋，国内代表厂商为 MIPS 架构的龙芯，Alpha 架构的申威，ARM 架构的飞腾、海思、华芯通，RISC-V 架构的阿里巴巴、兆易创新、赛昉科技等。

表 1 CPU 各类架构概览⁶

分类	特点	引领者	优劣势
x86	<ol style="list-style-type: none"> 1. 复杂指令集 2. 单核能力强 	Intel、AMD、海光、兆芯	软件生态好，占有率高；指令集实现复杂，功耗高
ARM	<ol style="list-style-type: none"> 1. 精简指令集 2. 追求多核 3. 低功耗 	安谋、高通、Amazon、飞腾、华为	获得授权的厂商多，能效比高；软件生态劣于 x86
MIPS	<ol style="list-style-type: none"> 1. 精简指令集 2. 低功耗 	龙芯	软件生态弱 市占率正在下降
Power	<ol style="list-style-type: none"> 1. 精简指令集 2. 单核能力强 3. 高可靠性、高成本 	IBM	IBM 掌控技术 主要应用在金融领域
RISC-V	<ol style="list-style-type: none"> 1. 精简指令集 	RISC-V 基金会、阿里巴巴、兆易创新、华米科技、赛昉科技、芯来科技	RISC-V 架构完全开放，指令精简、模块化、可扩展、开源
Alpha	<ol style="list-style-type: none"> 1. 精简指令集 2. 速度快 	申威	软件生态弱，市占率很小

3.2.2 市场

CPU 服务器处理器市场出货量在 2021 年约为 3600 万（Unit），出货量同比增长 10%。

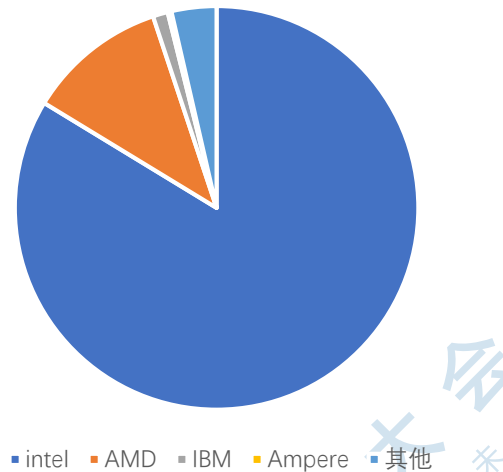


数据来源：公开数据整理

图 4 近四年 CPU 服务器处理器出货量规模

⁶ 来源：企业官网等公开资料，中国信息通信研究院整理

Intel 与 AMD 两家占据市场近 95% 的份额。截止到 2021 年第四季度，Intel 在全球数据中心 CPU 处理器市场份额的占比为 84%，其次 AMD 市场份额达到了 11%，其余 5% 的市场份额分别由 IBM、Ampere、Marvell、Oracle 以及其它厂家占据。



来源：公开数据整理

图 5 2021 年第 4 季度 CPU 市场份额

预计 2021–2026 年，全球 CPU 市场的年复合增长率超过 5%，增长速度自 2023 年起趋于平稳。x86 架构依然占据绝对市场优势，同时，ARM 架构正在崛起，由于其体积小、功耗低、成本低、性能好等优势，广泛地被应用在许多嵌入式系统设计，增长迅速。

3.2.3 技术

3.2.3.1 x86

x86 架构（The x86 architecture）CPU 处理器是由 Intel 首先开发并主导的基于 x86 指令集的一系列处理器统称。x86 架构软件生态好，市场占有率高，同时 x86 的主导地位也导致了其创新动力不足、议价能力较高等问题。

（一）Intel 产品方面

Intel 的服务器 CPU 平台主要包括：CPU Pentium Pro、Thurley、Brickland、Romley、Grantley、Purley 和 Whitley。服务器 CPU 制程由 2008 年的 45nm 提升至 2020 年的 10nm，预计将更新至 7nm 制程。如表 2 所示，Intel 公司将服务器 CPU 制程从 14nm 提升至 10nm，内核数和主频在一定程度上均有大幅度的提高。

表 2 2021 年 Intel 部分处理器产品⁷

处理器产品	发行日期	内核数	最大睿频频率 (GHz)	基本频率 (GHz)	TDP (W)	制程 (nm)	PCIe
英特尔® 至强® Platinum 8368 处理器	Q2'21	38	3.40	2.40	270	10	64

（二）AMD 产品方面

自 2003 年开始，AMD 陆续推出双核皓龙处理器、四核皓龙处理器、14nm 制程的 EPYC（霄龙）处理器、7nm 制程的 EPYC（霄龙）7002 系列处理器等。AMD 服务器 CPU 制程工艺由 2017 年的 14nm 快速提升到 2019 年的 7nm，2022 年将更新至 5nm 制程。目前，AMD EPYC（霄龙）7003 系列处理器为 Zen 3 内核和 AMD Infinity 架构。

表 3 2021 年 AMD EPYC（霄龙）部分产品概况⁸

处理器产品	发行日期	内核数	最大睿频频率 (GHz)	基本频率 (GHz)	TDP (W)	制程 (nm)	PCIe
AMD EPYC™ 7763	Q1'21	64	3.50	2.45	280	7	128

3.2.3.2 ARM

ARM 架构（Advanced RISC Machine，更早称作：Acorn RISC Machine），为一个 32 位精简指令集（RISC）处理器架构，广泛地使用在嵌入式系统设计中。基于 ARM 的处理器具有运行速度快、功耗低、成本低等优点，被广泛应用于通信、存储、安全系统等领域。

（一）安谋（ARM）产品方面

公司在经典处理器 ARM11 以后的产品改用 Cortex 命名，并分成 A、R 和 M 三类，旨在为各种不同的市场提供服务：“A”系列面向尖端的基于虚拟内存的操作系统和用户应用；“R”系列针对实时系统；“M”系列针对微控制器。

⁷ 来源：Intel 官网

⁸ 来源：AMD 官网

表 4 ARM 3 款产品对比⁹

	特征	最新版本	应用
Cortex-A	1. 在所有架构配置文件中提供最高性能 2. 高能效 3. 优化运行丰富的操作系统	ARM V8 和 ARM V9	PC、笔记本电脑、智能电视、服务器、智能手机和汽车主机、云存储和超级计算机
Cortex-R	针对有实时要求的系统进行了优化	ARM V8	医疗设备、车辆转向、制动和信号、网络和存储设备以及嵌入式控制系统
Cortex-M	专为小型、低功耗、高能效设备而设计	ARM V8	安全处理器、物联网和嵌入式设备，如可穿戴设备、小型传感器、通信模块和智能家居产品

（二）安晟培半导体科技有限公司（Ampere Computing）产品方面

2020年，安晟培半导体科技有限公司正式发布 Ampere Altra 处理器。2021年，推出新品 Ampere Altra Max，内核数量高达 128 核。

表 5 Ampere Altra Max 产品规格¹⁰

Ampere Altra Max	
发行日期	2021
核心	128
架构	ARM V8.2
主频	3.0 GHz
L1	64 KB L1-I 64 KB L1-D
L2	1MB
系统级缓存	16MB
内存	8× 72-bit DDR4-3200
工艺/典型功耗	7nm/250w

（三）亚马逊产品方面

Graviton2 是 AWS 基于 ARM 架构的处理器，是 ARM 第一款以数据中心定位的 CPU 架构。

⁹ <https://www.arm.com/why-arm/architecture/cpu>

¹⁰ 来源：Ampere computing 官网

表6 亚马逊 Graviton2 产品规格¹¹

	Graviton2
发行日期	2019
核心	64
架构	ARM V8.2
主频	2.5 GHz
L1	64 KB L1-I 64 KB L1-D
L2	1MB
L3	32MB Shared
内存	8× DDR4-3200
工艺	7nm
典型功耗	110W

（四）飞腾产品方面

飞腾代表产品有 FT-1500A、FT-2000、S2500。FT-1500A/16、FT-2000、S2500 芯片均采用片上并行系统（PSoC）体系结构，兼容 64 位 ARM V8 指令集。

表7 飞腾 FT-1500A、FT2000、S2500 产品规格¹²

	FT-1500A	FT2000	S2500
核心	集成 16 个 FTC660 处理器核	集成 64 个 FTC662 处理器核	集成 64 个 FTC663 处理器核
主频	1.5GHz	工作主频 1.8、2.0、2.2GHz	2.0~2.2GHz
二级缓存	8MB	32MB	每 4 核共享 2MB L2，总共 32MB
三级缓存	8MB	/	64MB
存储控制器	4 个 DDR3 接口	8 个 DDR4 接口	集成 8 个 DDR4-3200 通道
网络接口	2 个 100Mbps 以太网调试口		集成 4 个直连通路，每个通路组成为 X4，单 lane 速率 25Gbps，支持 2、4、8 路 CPU 互连
典型功耗	35W	100W	150W

（五）华为产品方面

鲲鹏 CPU 基于 ARM V8 架构，处理器核、微架构和芯片均由华为自主研发设计。近年来，华为先后推出 Hi1610、Hi1612、Hi1616 等服务器 CPU 产品。

¹¹ 来源：亚马逊官网

¹² 来源：飞腾官网、中国信息通信研究院整理

表 8 华为 Hi1620、Hi1616、Hi1612、Hi1610 产品规格¹³

	Hi1620 (鲲鹏 920)	Hi1616 (鲲鹏 916)	Hi1612	Hi1610
发行日期	2019	2017	2016	2015
核心	24~64	32	32	16
架构	Ares	Cortex-A72	Cortex-A57	Cortex-A57
主频	2.4~3.0 GHz	2.4 GHz	2.1 GHz	2.1 GHz
L1	64 KB L1-I 64 KB L1-D	48 KB L1-I 32 KB L1-D	48 KB L1-I 32 KB L1-D	48 KB L1-I 32 KB L1-D
L2	512KB Private	1MB/4 cores	1MB/4 cores	1MB/4 cores
L3	1MB/core Shared	32MB CCN	32MB CCN	16MB CCN
内存	8× DDR4-3200	4× DDR4-2400	4× DDR4-2133	2× DDR4-1866
典型功耗	100~200W	85W		

3.2.3.3 MIPS

MIPS 架构(Microprocessor without interlocked pipelined stages, 无内部互锁流水级的微处理器), 是一种采取精简指令集 (RISC) 的处理器架构。1981 年, 由 MIPS 科技公司开发并授权, 广泛被使用在电子产品、网络设备、个人娱乐装置与商业装置上。最早的 MIPS 架构是 32 位, 最新的版本已升级至 64 位。

MIPS 公司是全球第二大半导体设计 IP (知识产权) 公司和全球第一大模拟 IP 公司。MIPS 是出现最早的商业 RISC 架构芯片之一, MIPS 公司只进行 CPU 设计, 之后将设计方案授权给客户, 使得客户能够制造出高性能的 CPU。2007 年 8 月 16 日, MIPS 公司宣布, 中科院计算机研究所的龙芯中央处理器获得其处理器 IP 的全部专利和总线、指令集授权。

龙芯产品方面, 2019 年年底发布龙芯 3A4000/3B4000。龙芯 3B4000 在核心线程、频率上与芯 3A4000 一致, 但支持双路、四路服务器, 一台服务器最多包含 16 个处理器核。所有 CPU 之间通过高速总线接口直接互联, 共享使用物理内存。龙芯 3B4000 专门优化了 CPU 之间的高速互连总线, 跨片访存实际带宽提升 400% 以上。

2021 年, 龙芯推出全自主的指令系统 LoongArch。龙芯 3A5000/3B5000 是面向个人计算机、服务器等信息化领域的通用处理器, 基于 LoongArch 的 LA464 微结构, 龙芯 3B5000 在龙芯 3A5000 的基础上支持多路互连。龙芯 3C5000L 是专门面向服务器领域的通用处理器, 基于龙芯 3A5000 处理器, 片上集成共 16 个高性能 LA464 处理器核。

¹³ 来源: 企业官网等公开资料, 中国信息通信研究院整理

表9 龙芯 3A5000/3B5000、龙芯 3C5000L 产品规格¹⁴

	龙芯 3A5000/3B5000	龙芯 3C5000L
发行日期	2021	2021
核心个数	4	16
处理器核	支持 LoongArch®指令系统；支持 128/256 位向量指令；四发射乱序执行；4 个定点单元、2 个向量单元和 2 个访存单元	支持 LoongArch®指令系统；支持 128/256 位向量指令；四发射乱序执行；4 个定点单元、2 个向量单元和 2 个访存单元
主频	2.3GHz–2.5GHz	2.0GHz–2.2GHz
峰值运行速度	160GFlops	560GFlops
高速缓存	每个处理器核包含 64KB 私有一级指令缓存和 64KB 私有一级数据缓存；每个处理器核包含 256KB 私有二级缓存；所有处理器核共享 16MB 三级缓存	每个处理器核包含 64KB 私有一级指令缓存和 64KB 私有一级数据缓存；每个处理器核包含 256KB 私有二级缓存；每 4 个处理器核共享 16MB 三级缓存，共 64MB 三级缓存
内存控制器	2 个 72 位 DDR4–3200 控制器；支持 ECC 校验	4 个 72 位 DDR4–3200 控制器；支持 ECC 校验
高速 IO	2 个 HyperTransport3.0 控制器；支持多处理器数据一致性互连（CC–NUMA）	4 个 HyperTransport3.0 控制器；支持多处理器数据一致性互连（CC–NUMA）
功耗管理	支持主要模块时钟动态关闭；支持主要时钟域动态变频；支持主电压域动态调压	支持主要模块时钟动态关闭；支持主要时钟域动态变频；支持主电压域动态调压
典型功耗	35W@2.5GHz	130W@2.2GHz

2022 年 6 月发布龙芯 3C5000 芯片。龙芯 3C5000 芯片是面向服务器领域的通用处理器，在兼容龙芯 3C5000L 主板设计的基础上，调整优化了封装形式，保持了系统和应用软件的兼容性。

2022 年 7 月，龙芯发布芯片 7A2000。龙芯 7A2000 是面向服务器和个人计算机领域。芯片 7A2000 是在第一代 7A1000 的基础上进行了优化升级，内置一个网络 PHY，可直接提供网络端口输出。

¹⁴ 来源：企业官网等公开资料，中国信息通信研究院整理

表 10 龙芯 7A1000、龙芯 7A2000 产品规格

	龙芯 7A1000	龙芯 7A2000
发行日期	2020	2022
处理器接口	HT3.0 × 16 3.2Gbps	HT3.0 × 16 3.2Gbps
GPU	支持 2D、3D	支持 3D
显存	DDR3 16 位	DDR4 32 位
显示接口	DVO*2	HDMI*2、VGA*1
功耗	5~8W	7~14W

3.2.3.4 Power

Power 架构是 IBM 开发的一种基于 RISC 指令系统的架构。

Power 架构的处理器具有结构简单和高效率的特点。苹果早期的电脑就是使用了 Power 架构,Power 系列微处理器应用在不少 IBM 服务器,超级电脑,小型电脑及工作站中,广泛作为主 CPU 使用。PowerPC 架构也是源自 Power,并应用在苹果电脑的麦金塔电脑及部分 IBM 的工作站,以及各式各样的嵌入式系统上。此外,IBM 透过 Power.org 网站,向其他开发者及制造商推广 Power 架构及其他派生产品。

IBM 产品方面,1990 年 2 月,IBM 推出了第一部采用 Power 架构的 IBM 电脑被称作 RISC System/6000 或 RS/6000。从此,IBM 开始推出一系列 Power 框架,从 1990 年的 Power 到 Power2、PowerPC、Power3、Power4、Power5、Power6 等等,到目前最新版本的 Power10。

表 11 IBM 的 Power1990–2019 年部分产品概况¹⁵

处理器产品	发行日期	内核数	最大睿频频率 (GHz)	基本频率 (GHz)	制程 (nm)
Power6	2007	2	5.0	3.6	65
Power7	2010	4	4.25	2.4	45
Power8	2014	6/12	5.0	2.5	22

IBM 的 Power10 是第一款采用 7nm 工艺的商用芯片,通过内置嵌入式矩阵数学加速器,使 FP32、BFloat16 和 INT8 计算的 AI 推理速度分别提高 10 倍、15 倍和 20 倍。凭借额外的 AES 加密引擎,IBM Power 10 可以为当前领先的加密标准以及未来加密协议提供硬件内存加密,可保证端到端安全,实现更高的加密性能。

¹⁵ 来源: IBM 官网

表 12 IBM 的 Power10 产品概况¹⁶

处理器产品	发行日期	内核数 最多	最大睿频频率 (GHz)	基本频率 (GHz)	TDP (W)	制程 (nm)	PCIe
Power10	2020.8	15/30	4	3.5	180	7	5.0

3.2.3.5 RISC-V

RISC-V 是基于精简指令集 (RISC) 原则的开源指令集架构。该框架 2010 年由加州大学柏克莱分校率先提出, 2015 年 RISC-V 基金会成立。目前有参与支持 RISC-V 基金会的公司以及机构包括了阿里巴巴、华为、中国科学院、清华大学、加州大学伯克利分校、莱迪思半导体、迈伦科技、美高森美等企业及机构。

RISC-V 指令集具有很强的灵活性, 其设计使其适用于现代计算设备。RISC-V 指令集的设计考虑了小型、快速、低功耗的现实情况, 具有众多软件的支持, 这解决了新指令集以往的弱点。

2021 年 5 月, 阿里巴巴发布玄铁系列新款处理器—玄铁 907, 该处理器对开源 RISC-V 架构进行优化设计, 兼顾高性能及低功耗特点, 可应用于 MPU (微处理器)、智能语音、导航定位、存储控制等领域。

表 13 阿里的玄铁框架性能表¹⁷

处理器产品	发行日期	支持内核数	单核性能 (coremark/MHz)	主频 (GHz)	制程 (nm)	其他性能
玄铁 E907	2021	/	3.8	1	28	五级流水线
玄铁 C910	2019	16	7.1	2.5	16	十二级流水线

3.2.3.6 Alpha

Alpha (DEC Alpha, 也称为 Alpha AXP) 是 64 位的 RISC 微处理器。最初由 DEC 公司制造, 并被用于 DEC 自己的工作站和服务中。Alpha 作为 VAX 的后续被开发, 支援 VMS 操作系统。不久之后开放源代码的操作系统也可以在其上运行。但是从 Windows 2000 beta3 之后就放弃了对 Alpha 的支持。

DEC Alpha 系列的 Alpha AXP 21064 于 1992 年发布, 采用 RISC 指令集设计, 发布时最高主频达 200MHz, 浮点性能创世界记录。Alpha 21164 系列于 1995 年发布, 发布时主频为 266MHz。Alpha 21264 系列于 1998 年发布, 引入乱序执行功能和多媒体加速指令集, 在 21164 基础上实现性能翻倍。2002 年发布最后一代 Alpha 架构产品 Alpha 21364, 此时 DEC 被收购。

¹⁶ 来源: IBM 官网

¹⁷ 来源: 中国信息通信研究院整理

2006年，申威基于DEC公司Alpha架构的研制出单核CPU，130nm工艺，主频900MHz。2010年，申威基于Alpha 21164开发了SW-3/SW1600处理器，SW1600已经运用到神威蓝光(Sunway BlueLight MPP)超级计算机上。之后，申威出于安全自主可控角度不再使用Alpha指令集。

表 14 申威 SW2、SW-3/SW1600 框架性能表¹⁸

处理器产品	发行日期	内核数	主频 (GHz)	制程 (nm)	用途
申威 SW2	2008	2	1.4	130	高性能计算机
申威 SW-3/SW1600	2010	16	1.6	65	超级计算机

3.3 智能算力

3.3.1 概述

智能算力主要以GPU (Graphics Processing Unit, 图形处理器)、FPGA (Field Programmable Gate Array, 现场可编程逻辑门阵列)、AI (Artificial Intelligence, 人工智能) 芯片等为代表。

表 15 智能算力主要芯片类型¹⁹

分类	特点	引领者
GPU	<ol style="list-style-type: none"> 通用性高 平行处理能力中等 交付周期中等 	英伟达、AMD
FPGA	<ol style="list-style-type: none"> 通用性中等 平行处理能力高 交付周期中等 	Xilinx (已被 AMD 收购)、Intel
AI 芯片	<ol style="list-style-type: none"> 通用性低 平行处理能力高 交付周期长 	寒武纪、燧原、华为等

3.3.2 市场

(一) GPU 市场规模

根据 IDC 数据显示，全球服务器 GPU 市场收入在 2021 年底为 71.5 亿美元，当前 GPU 市场几乎被英伟达、AMD、Intel 垄断。

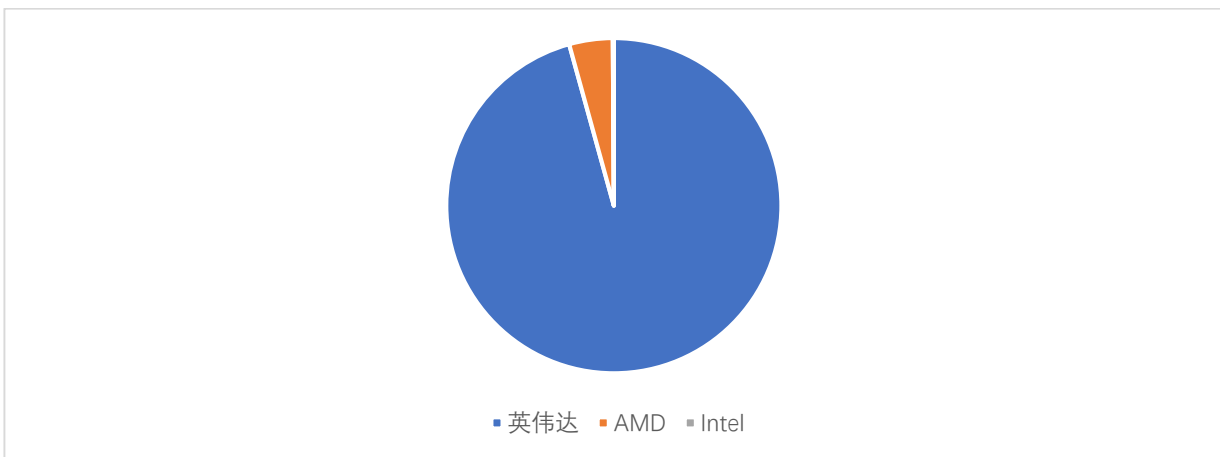
¹⁸ 来源：申威官网、中国信息通信研究院整理

¹⁹ 数据来源：中国信息通信研究院整理

表 16 全球服务器 GPU 市场季度统计²⁰（单位：\$M）

	第一季度	第二季度	第三季度	第四季度	年度合计
2017	491	458	522	586	2057
2018	684	754	780	668	2886
2019	570	600	665	745	2580
2020	912	1091	1189	1289	4481
2021	1456	1521	1820	2353	7150

截止到 2021 年第四季度，英伟达在 GPU 市场份额的占比为 95.7%，AMD 市场份额为 4.2%，Intel 市场份额为 0.1%。



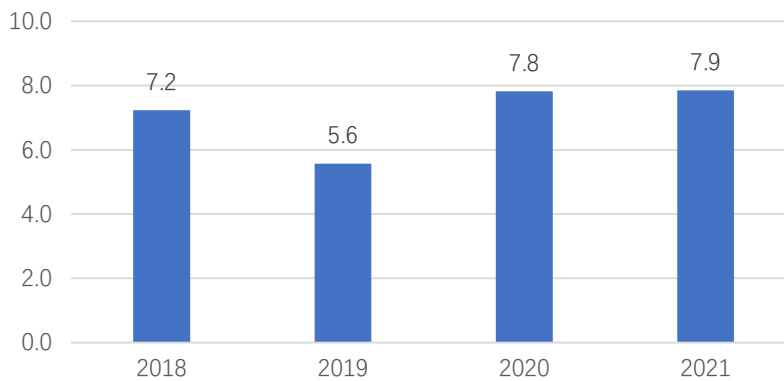
数据来源：IDC 官网

图 6 2021 年第 4 季度 GPU 市场份额

（二）FPGA 市场规模

FPGA 市场呈现双寡头垄断格局，Xilinx（已被 AMD 收购）和 Intel 分别占据全球市场 29.9% 和 70.1%。2020 年“新冠”疫情突然爆发，行业对远程和虚拟工作的需求增加，如教育、商业等，特别是 2020 年第二季度和第三季度，FPGA 市场份额大幅增加。截止到 2021 年底，FPGA 市场达到 7.9 亿美元。具体情况如下图所示。

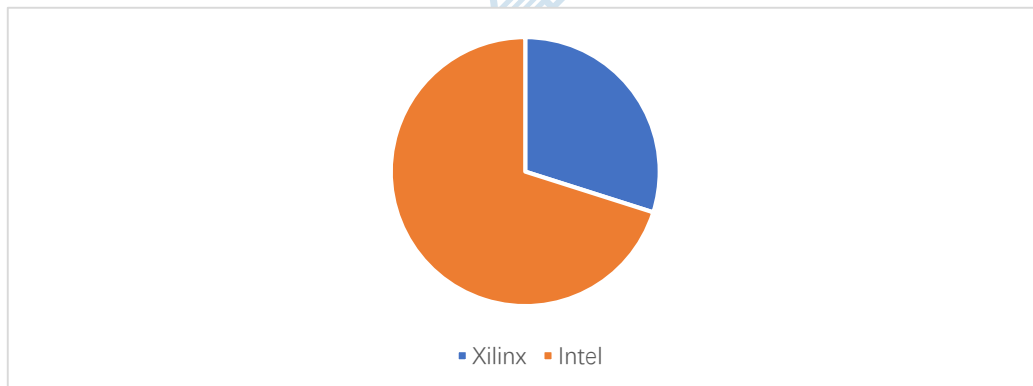
²⁰ 数据来源：IDC 官网



数据来源：IDC 官网

图7 FPGA 市场规模统计

FPGA 市场前景诱人，但是门槛之高在芯片行业里无出其右，Xilinx 与 Intel 这两家公司几乎占领全部市场份额。全球有 60 多家公司先后斥资数十亿美元，尝试登顶 FPGA 高地，但是最终登顶成功的只有四家公司：Xilinx（赛灵思，已被 AMD 收购）、Intel（已收购 Altera）、Lattice（莱迪思）、Microsemi（美高森美）。



数据来源：IDC 官网

图8 2021 年第 4 季度 FPGA 市场份额

3.3.3 技术

3.3.3.1 GPU

GPU 是一种专门在个人电脑、工作站、游戏机和一些移动设备（如平板电脑、智能手机等）上做图像和图形相关运算工作的微处理器。GPU 芯片按指令集架构划分，可分为 IMG A-Series 架构、IMG A-Series 架构、PowerVR Rogue 架构等。GPU 的生产商主要有 NVIDIA、AMD 等。

NVIDIA 在 1999 年定义了 GPU，推动了 PC 游戏市场的快速发展，重新定义了现代计算机图形技术，并彻底改变了并行计算。2020 年，NVIDIA 宣布推出 Ampere 架构，性能相较于前代提升了 20 倍。NVIDIA

A100 是首款基于 NVIDIA Ampere 架构的 GPU，作为一款通用型工作负载加速器，A100 还被设计用于数据分析、科学计算和云图形。

AMD 的 GPU 代表产品有 Instinct 系列等。AMD 从 2008 年首款突破 1 TFLOPS 的 GPU 发布以来，AMD 在 GPU 领域持续创新，分别推出了 AMD Radeon R9 和 R7、AMD FirePro S7100X 多用户 GPU 等。2017 年发布首款基于 Vega 架构的 14nm GPU Instinct MI25，2018 年推出基于 Vega20 架构的 7nm GPU Instinct MI50，2020 年推出了首款基于 CDNA 架构的 GPU Instinct MI100，2021 年推出 Instinct MI200 系列加速器，基于 AMD CDNA 2 架构，可为广泛的 HPC 工作负载提供领先的应用程序性能。

表 17 AMD 部分产品概况²¹

型号	计算单元	流处理器	FP64 FP32 Vector (Peak)	FP64 FP32 Matrix (Peak)	FP16/bf16 (Peak)	INT4/INT8 (Peak)	HBM2e ECC 显存	显存带宽	规范
AMD Instinct MI250x	220	14,080	Up to 47.9 TF	Up to 95.7 TF	Up to 383.0 TF	Up to 383.0 TOPS	128GB	3.2 TB/sec	OCP 加速器模块
AMD Instinct MI250	208	13,312	Up to 45.3 TF	Up to 90.5 TF	Up to 362.1 TF	Up to 362.1 TOPS	128GB	3.2 TB/sec	OCP 加速器模块

3.3.3.2 FPGA

FPGA (Field Programmable Gate Array, 现场可编程逻辑门阵列) 是在 PAL、GAL 等可编程器件的基础上进一步发展的产物。FPGA 作为专用集成电路 (ASIC) 领域中的一种半定制电路而出现，可以克服定制电路不足、原有可编程器件门电路数有限的问题。与 ASIC 芯片相比，FPGA 的一项重要特点是其可编程特性，即用户可通过程序指定 FPGA 实现某一特定数字电路。FPGA 芯片是小批量系统提高系统集成度、可靠性的最佳选择之一。

FPGA 芯片的行业技术壁垒高于 CPU，主要市场被 Xilinx (已被 AMD 收购)、Intel (于 2015 年收购 Altera) 两家垄断。国内有能力自主研发 FPGA 的厂商有智多晶微电子、紫光同创、安路科技、京微齐力等。

(一) Xilinx 方面 (AMD)

1985 年，Xilinx 推出全球第一款 FPGA 产品 XC2064，采用 2μm 工艺，包含 64 个逻辑模块和 85000 个晶体管。1991 年，其推出第一款被广泛应用的 FPGA——XC4000。随后，从 2001 年采用的 150nm 工艺，到 2017 年 Xilinx 宣布 28nm 的 Spartan7 进入量产阶段，FPGA 的工艺制程实现了逐代提升。当前

²¹ 来源：AMD 官网

Xilinx 的产品线主要为 45nm、28nm、20nm、16nm 四种，涵盖了不同等级的四个系列：主打低价格低功耗的 Spartan 系列、在 Spartan 基础上增加串行收发器和 DSP 功能的 Artix 系列、中端性能的 Kintex 系列，以及最高端的 Virtex 系列。2020 年 10 月 27 日，AMD 公司与 Xilinx 达成协议，同意 AMD 以发行总价值 350 亿美元股票的方式收购 Xilinx，AMD 于 2021 年底完成对 Xilinx 的收购工作。

表 18 2010–2021 年 Xilinx 部分 FPGA 产品²²

FPGA 名称	发行日期	逻辑运算单元	分布式 RAM(Mb)	I/O 接口	DSP	BRAM(Mb)	制程 (nm)
Virtex-6 XC6VLX760	Q1'10	758k	8.2	30	864	25.9	40
Virtex-7 UltraScale+ VU19P	Q2'19	8938k	58.4	2072	3840	75.9	16

（二）Altera 方面（Intel）

1984 年，Altera 发明的 EP300 以首款可编程逻辑器件（PLD）身份问世。1992 年，Altera 公司推出该公司第一款 FPGA——FLEX 8000 FPGA。2002 年，推出了首款带有嵌入式 DSP 模块的 FPGA——Stratix FPGA 系列；同年，采用 0.13μm 工艺的 FPGA——Cyclone FPGA 问世。此后，Altera 对 FPGA 的工艺制程和内部的 DSP 等模块进行了稳步提升，至 2012 年发布的 Stratix V FPGA 制程达到 28nm，是业界首款高性能 28nm FPGA。

2015 年，Intel 以 167 亿美元收购 FPGA 厂商 Altera，延续并扩充了其产品线。当前，定位自下而上，Intel 的系列可以划分为：Cyclone、MAX、Arria、Stratix、Agilex 系列。其在 2010–2021 年期间，Intel（Altera）FPGA 的制程从 28nm 提升至 10nm，逻辑元素数量提升至 14 倍，最大嵌入式内存大小提升至 5 倍。

表 19 2010–2021 年 Intel（Altera）部分 FPGA 产品²³

FPGA 名称	发行日期	逻辑元素 LE	MLAB RAM(Mb)	I/O 接口	DSP	BRAM(Mb)	制程 (nm)
Stratix® V 5SGSD8 FPGA	Q2'10	695k	8.01	840	1963	50	28
Intel® Stratix® 10 GX 10M FPGA	Q3'19	10200k	55	2304	3456	253	14
Intel® Agilex F- Series 027 FPGA (R24C)	Q2'19	2692k	28	744	8528	259	10

²² 来源：Xilinx 官网

²³ 来源：Intel 官网

（三）国内厂商方面

西安智多晶微电子有限公司成立于 2012 年，现性能最强产品为 Seal 5000 系列 FPGA，基于 28nm 制程打造；深圳市紫光同创电子有限公司成立于 2013 年，其发布的 Titan 系列是中国第一款国产自主产权千万门级高性能 FPGA 产品，采用 40nm 工艺；上海安路信息科技有限公司成立于 2011 年，其推出的 EAGLE 系列 FPGA 采用 55nm 工艺；京微齐力（北京）科技有限公司于 2017 年成立，其产品主要有 HME-R（河）系列、HME-M7（华山）系列、HME-M5（金山）系列等，当前主要工艺为 40nm。

表 20 2010-2021 年部分国产厂商 FPGA 产品²⁴

厂商	FPGA 名称	发行日期	逻辑单元 LE	分布式 RAM(Mb)	I/O 接口	DSP	嵌入式 RAM	制程 (nm)
西安智多晶微电子	SA5-325E	Q2'21	326k	0.650	500	1312	18.9Mb	28
紫光同创	PGT180H	Q3'15	174k	0.057	611	-	9468kb	40
安路信息科技	EG4S20BG256	Q1'21	19.6k	0.156	193	-	1244kb	55
京微齐力	HR03PN3	Q2'21	3.1k	-	128	-	72kb	40

3.3.3.3 AI 芯片

从广义上讲，能运行 AI 算法的芯片都叫 AI 芯片，目前通用的 CPU、GPU、FPGA 等都能执行 AI 算法，但是执行效率差异较大。狭义上，将 AI 芯片定义为“专门针对 AI 算法做了特殊加速设计的芯片”。国家第十四个五年规划明确提出聚焦高端芯片、人工智能关键算法等关键领域，加快布局神经芯片等前沿技术。国内华为、寒武纪、燧原科技等新兴 AI 芯片持续涌现。

（一）华为

华为昇腾 AI 处理器 (NPU)，采用华为自研达芬奇架构，核心技术 3D Cube。代表产品主要为 Ascend 310（昇腾 310）、Ascend 910（昇腾 910）。昇腾（Ascend）910 的最大功耗 310W，八位整数精度（INT8）下的性能达到 640TOPS，16 位浮点（FP16）下的性能达到 320TFLOPS。

（二）寒武纪

寒武纪开发的 MLU，基于 7nm 制程工艺，采用寒武纪 MLUv02 架构。主要代表产品为思元 370/290/270 等系列。思元 370 是寒武纪首款采用 chiplet（芯粒）技术的 AI 芯片，集成了 390 亿个晶体管，最大算力高达 256TOPS（INT8），是寒武纪第二代产品思元 270 算力的 2 倍。

（三）燧原科技

²⁴ 来源：中国信息通信研究院整理

邃思 2.0 芯片基于人工智能领域专用处理器架构设计，提供全精度人工智能算力、灵活的可扩展性。云燧 T20 是基于邃思 2.0 芯片打造的面向数据中心的第二代人工智能训练加速卡，采用高密的计算芯片。单精算力最高可达 40TFLOPS（FP32），最高可支持 64GB 容量，1.8TB/s 带宽；300GB/s 的独立片间互联通道提供灵活的多芯片算力扩展方案。

（四）百度

百度昆仑芯 AI 芯片与飞腾等多款国产通用处理器、麒麟等多款国产操作系统以及百度自研的飞桨深度学习框架完成了端到端的适配，拥有软硬一体的全栈国产 AI 能力。主要代表产品为昆仑芯 1、昆仑芯 2。旗下昆仑芯 2 采用 7nm 制程，整数精度(INT8)算力达到 256 TOPS，半精度(FP16)为 128 TFLOPS，而最大功耗仅为 120W。

（五）阿里巴巴

阿里巴巴采用平头哥自研框架，通过软硬件协同设计实现性能突破。主要代表产品为含光 800。含光 800 是一颗高性能人工智能推理芯片，基于 12nm 工艺，集成 170 亿晶体管，性能峰值算力达 820 TOPS。

（六）超聚变

超聚变 GPU 服务器支持软硬协同，提供极致灵活 AI 算力，采用全模块化架构创新设计，支持 CPU 和 GPU 解耦演进，实现资源按需配置，保护客户投资；独创屋檐式架构设计，CPU 与 GPU 计算模块共享空间，同时以领先的散热工程能力为支撑，4U 空间内 AI 算力可达业界的 1.25 倍；多场景工作模式可一键切换拓扑，快速实现训练、AI 推理、AI 云加速、HPC 等多种场景下的 CPU：GPU 最佳配比，实现资源不闲置、投资不重复。

3.4 超算算力

3.4.1 概述

超级计算是计算科学的重要概念，是超级计算机及有效应用的总称。超级计算利用并行工作的多台计算机系统的集中式计算资源，并通过专用的操作系统来处理极端复杂的或数据密集型的问题。

超级计算机又称巨型计算机等，指能解决复杂计算的大型、快速、价格昂贵的计算机。通常这类机器包括了从标准计算机的大型集群到高度专用的硬件，主要运用于尖端科研、国防军工等大科学、大工程、大系统中，是一个国家科研实力的体现，是国家科技发展水平和综合国力的重要标志。

并行计算作为超算最为核心的研究方向，其技术与超算生态的建设密不可分。加州大学伯克利分校在《并行计算研究前景：伯克利的视角》一文中提到了 13 个并行计算最为关键的领域，具体描述如下：

表 21 13 个并行计算关键领域²⁵

Dwarf	排名	描述	Benchmark/例子
稠密矩阵	1	数据是稠密矩阵或向量。（BLAS Level 1=向量-向量；Level 2=矩阵-向量；Level 3=矩阵-矩阵。）通常情况下，此类方法使用单位跨步内存访问从行读取数据，并使用跨步访问从列读取数据。	块三对角矩阵，下上对称 Gauss-Seidel
稀疏矩阵	2	数据集包括许多零值。数据通常存储在压缩矩阵中，以减少访问所有非零值的存储和带宽要求。一个例子是块压缩稀疏行（BCSR）。由于采用压缩格式，数据通常通过索引存取进行访问。	共轭梯度法（CG 方法）
谱方法	3	数据位于频域，而不是时间域或空间域。通常，光谱方法使用多个蝶形阶段，将乘加运算和特定的数据排列模式相结合，对某些阶段进行全对全通信，对其他阶段进行严格的局部通信。	傅里叶变换
多体模拟	4	取决于许多离散点之间的相互作用。包括粒子-粒子方法，即每个点依赖于所有其他点，导致 $O(N^2)$ 计算复杂度，以及分层粒子方法，其组合来自多个散点的力或势，以将计算复杂性降低到 $O(N \log N)$ 或 $O(N)$ 。	光滑粒子法
结构化网格	5	由规则网格描述，网格上的点一起更新。结构化网格具有高度的空间局部性。可以直接进行更新，也可以在两套网格之间进行。网格可细分更细的网格（AMR，自适应网格细化）；网格粒度之间的转换可能会动态发生。	多重网格
非结构化网格	6	一种不规则的网格，通常根据应用程序的底层特征选择数据位置。数据点位置和相邻点的连通性必须明确。网格上的点一起更新。更新通常涉及多个级别的内存引用间接寻址，因为对任何点的更新都需要首先确定相邻点的列表，然后从这些相邻点加载值。	非结构自适应
蒙特卡洛方法	7	计算取决于重复随机试验的统计结果。蒙特卡洛被认为是“易于并行”的。	粒子输运
组合逻辑	8	使用逻辑函数和存储状态实现的函数。	逻辑电路模拟
图遍历	9	通过跟随连续边访问图中的多个节点。这些应用程序通常涉及许多级别的间接寻址，并且计算量相对较小。	BFS/DFS
动态规划	10	通过求解更简单的重叠子问题来计算求解。特别适用于具有大量可行解的优化问题。	
回溯和分支界限算法	11	通过递归将可行区域划分为子域，然后修剪次优子问题来找到最优解。	
构建图模型	12	构造将随机变量表示为节点、将条件依赖项表示为边的图。	有向无环图
有限状态机	13	一种系统，其行为由状态、由输入和当前状态定义的转换以及与转换或状态相关的事件来定义。	

²⁵ 来源：《并行计算研究前景：伯克利的视角》

1964年，有“超算之父”之称的 Seymour Cray 研制的 CDC6600 问世，并安装到美国 Livermore 和 Los Alamos 国家实验室，开启了超级计算技术和产业 60 年的持续发展与繁荣。

超级计算 60 年的演变路线可简单地分为 2 个阶段：**Cray 时代和多计算机时代**。全球超算行业主要经历了以下发展历程：国防驱动阶段—公司主导阶段—蓬勃发展阶段—多向发展阶段。

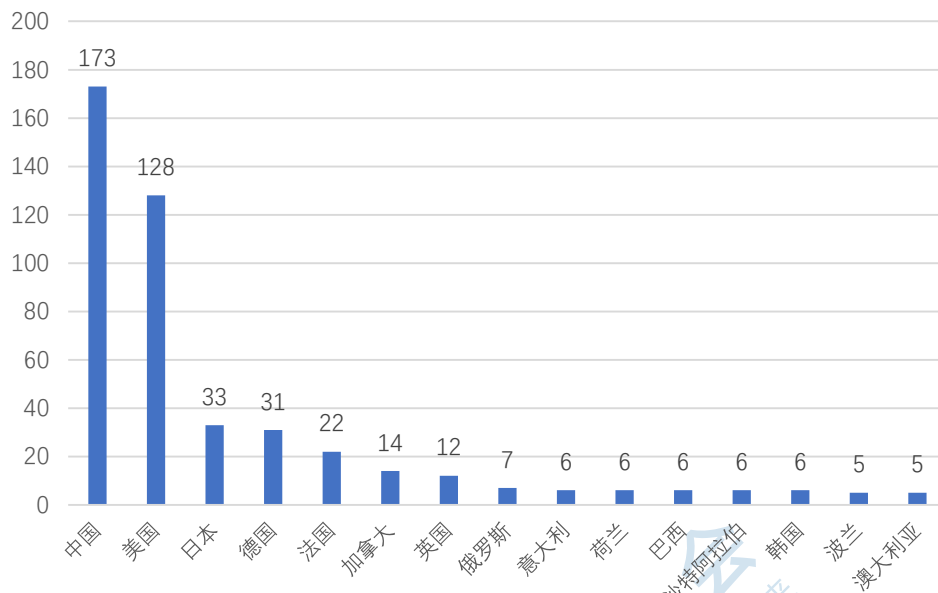
表 22 全球超算发展历程²⁶

时间	具体内容
国防驱动阶段 (1950s—1960s)	早期的计算机科学研究有着浓厚的国防军事色彩。20 世纪 50 年代后期，美国政府主要根据国家安全需求来支持情报和核武器应用研究。国家安全是开发超级计算技术的主要推动力。20 世纪 60 年代初的 IBM 7030 Stretch 和 SperryRand UNIVAC LARC 正是在这样的背景之下诞生的，因为其计算速度显著超出顶尖商用机的数量级而被视为早期的超级计算机。
公司主导阶段 (1960s—1970s)	20 世纪 60 年代中期到 70 年代末期，美国乃至全球的超级计算机行业主要由两家公司主导，即 Control Data 和 Cray Research。在这一阶段，超级计算机的成本得到有效控制，同时向量处理技术的速度得到了大幅提升，大量廉价高速的计算机走向商品市场。
蓬勃发展阶段 (1980s—1990s)	20 世纪 80 年代，日本政府大规模补贴计算机科研项目，同时推行排除国外竞争对手的产业政策。到了 90 年代，一批深耕半导体领域的日本计算机公司，如富士通、日立、NEC 等，成功获取 IBM 大型机技术的关键部分，并在本土推出了价格实惠的商用计算机系统。20 世纪 80 年代，业界开始转向大规模并行运算系统。1976 年问世的超级计算机 Cray1 是单向量机系统，之后为了进一步提高向量机的性能，在系统中不断增加向量部件的数量，即采用并行向量或多向量部件的技术
多向发展阶段 (21 世纪以来)	21 世纪以来，超级计算机开始呈现多极化发展。MPP 系统、集群系统的应用进一步提高了超级计算机性能，每秒千万亿次的 P 级超级计算机已经相对成熟，各个国家、各个科研机构 and 供应商正在 E 级超级计算机的研制中激烈竞争。

3.4.2 市场

在全球范围内，少数国家拥有强大的高性能计算能力。中国和美国领先，从最新发布的 2022 年 6 月的 TOP500 榜单中可以看出，中国共有 173 台上榜，仍稳居榜首，美国从 150 台下降至 128 台，两个国家占总数的 2/3。其次是日本、德国、法国、加拿大、英国、俄罗斯和意大利。

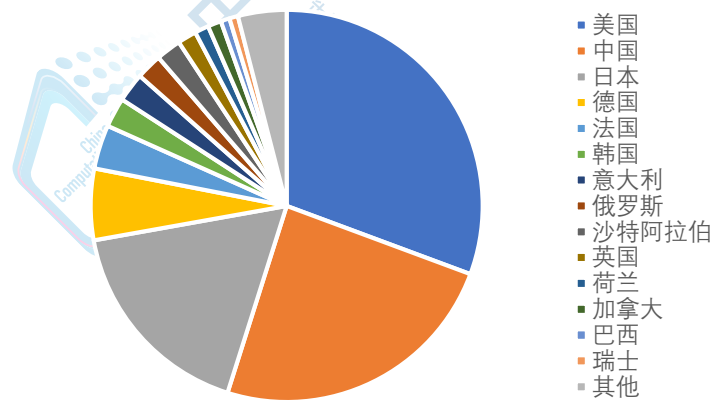
²⁶ 来源：中国信息通信研究院整理



数据来源：<http://www.top500.org/>

图9 2022年6月TOP500中部分国家超算数量

TOP500中，各国超算算力份额如下图所示，其中美国占比为31%，中国占比24%，日本占比17%，德国占比6%，法国占比4%，位列前五名。



数据来源：<http://www.top500.org/>、中国信息通信研究院整理

图10 2022年6月TOP500中各国超算所占份额

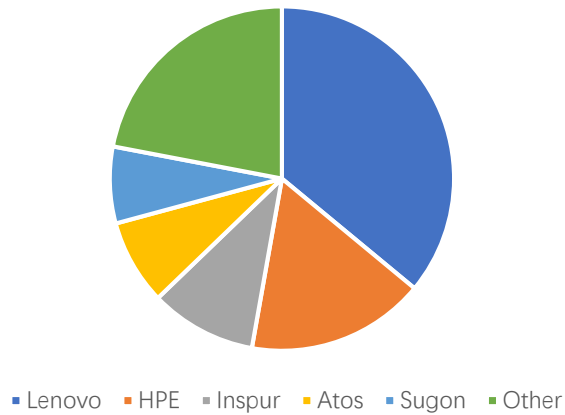
从1993年6月到2022年6月排名TOP500第一名的国家来看，美国11次，日本8次，中国3次。根据最新TOP500榜单，美国田纳西州橡树岭国家实验室的Frontier占据榜首。从厂商的角度看，IBM占据榜首5次，富士通4次，Cray、Intel、日立、国防科大分别登顶2次，TMC、NEC、理化研究所、国家并行计算机工程技术研究中心、美国橡树岭国家实验室分别登顶1次。

表 23 1993–2022 年世界排名第一超算²⁷

时间	名称	公司	国家
1993.06–1993.11	CM-5	TMC	美国
1993.11–1994.06	数值风洞	富士通	日本
1994.06–1994.11	Paragon XP/S140	Intel	美国
1994.11–1996.06	数值风洞	富士通	日本
1996.06–1996.11	SR2201	日立	日本
1996.11–1997.06	CP-PACS	日立	日本
1997.06–2000.11	ASCI Red	Intel	美国
2000.11–2002.06	ASCI White	IBM	美国
2002.06–2004.11	地球模拟器	日本电气 (NEC)	日本
2004.11–2008.06	蓝色基因/L	IBM	美国
2008.06–2009.11	走鹃 (超级计算机)	IBM	美国
2009.11–2010.11	美洲虎 (超级计算机)	Cray	美国
2010.11–2011.06	天河-1	国防科技大学	中国
2011.06–2012.06	京 (超级计算机)	理化研究所	日本
2012.06–2012.11	蓝色基因/Q	IBM	美国
2012.11–2013.06	Titan	Cray	美国
2013.06–2016.06	天河-2	国防科技大学	中国
2016.6–2017.11	神威·太湖之光	国家并行计算机工程技术研究中心	中国
2018.06–2019.11	Summit	IBM	美国
2020.06–2020.11	Supercomputer Fugaku	富士通	日本
2021.06–2021.11	Supercomputer Fugaku	富士通	日本
2021.11–2022.06	Frontier	美国橡树岭国家实验室	美国

全球厂商排名方面，从 2022 年 6 月公布的最新榜单看，联想、HPE、浪潮 排名前三，市场占比分别为 36%、16.8%、10%。从算力来看，HPE、富士通、联想 排名前三，总算力市场份额分别为 18.6%、18.1%、15.1%。

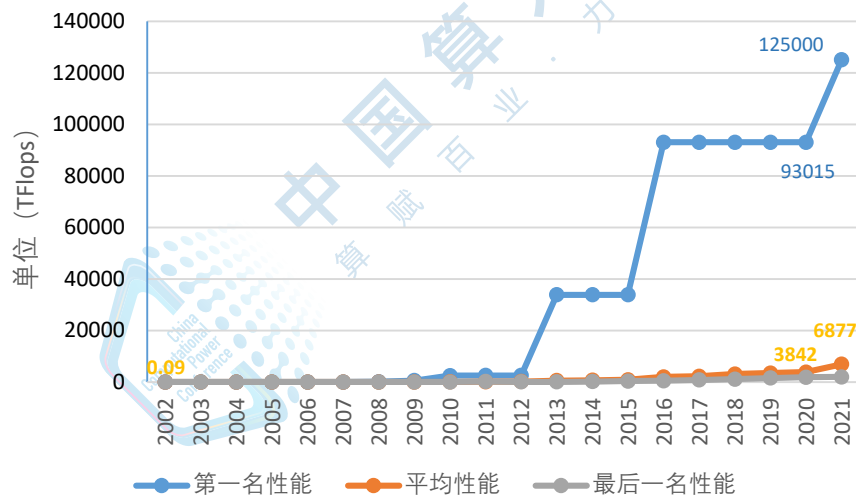
²⁷ 来源：中国信息通信研究院整理



数据来源：<http://www.top500.org/>、中国信息通信研究院整理

图 11 2022 年 6 月 TOP500 榜单对应厂商份额

我国目前已建成包括天津、广州、深圳、长沙、济南、无锡、郑州、昆山、成都等 9 座国家级超算中心。2002 年到 2021 年，超级计算机平均性能提升了七万多倍。

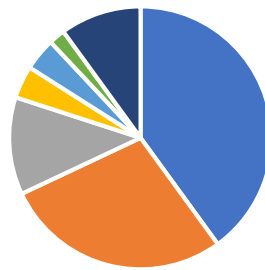


数据来源：<http://www.hpc100.cn/>、中国信息通信研究院整理

图 12 2002-2021 我国 TOP100 超级算力数据

根据最新发布的 2021 中国超级计算机 (HPC) 性能 TOP100 榜单分析得出, 2021 年全部上榜的 100 台超算系统的平均性能, 相比 2020 年提升 79%。2021 年的“第一名”的性能是 2020 年的 1.34 倍, 2021 年的“第一百名”的性能是 2020 年的 1.06 倍。排在第一名的测试性能是 125PFLOPS, 系统峰值达 240PFLOPS。第二名的测试性能是 93 PFLOPS, 系统峰值达 125PFLOPS。

排名前三的厂商份额合计占上榜系统份额的 80%, 分别是, 联想 40 套系统, 浪潮 28 套系统、中科曙光 12 套系统。



- 联想
- 浪潮
- 曙光
- 国防科大
- 北龙超云/DELL
- 国家并行计算机工程技术研究中心
- 其他

数据来源：<http://www.hpc100.cn>、中国信息通信研究院整理

图 13 2021 年 11 月 TOP100 厂商系统数

3.4.3 技术

3.4.3.1 超算算力的测试工具

（一）HPL

HPL 即 High Performance Linpack，它是针对现代并行计算集群的测试工具。用户不修改测试程序，通过调节问题规模大小 N （矩阵大小）、进程数等测试参数，使用各种优化方法来执行该测试程序，以获取最佳的性能。

（二）HPCG

HPCG 高度共轭梯度基准测试，是现在主要测试超算性能测试程序之一，也是 TOP500 的一项重要指标。HPCG 基准项目旨在创建一个新的 HPC 系统排名指标。HPCG 作为高性能 LINPACK（HPL）基准的补充，目前用于排名 500 强计算系统。HPL 的计算和数据访问模式仍然是一些重要的可扩展应用程序的代表，但并非全部。HPCG 的设计目的是使基准的计算和数据访问模式更好地匹配不同的应用程序。

3.4.3.2 超算算力的架构

结合业内最新实践和设想，同时还需满足实际应用场景，提出超算参考总架。该架构包含的既有现阶段已落地的基础能力，也有对超算应具备的补充能力。

（一）硬件架构

超算核心系统由计算系统、存储系统、网络系统、管理系统、安全系统五部分构成，方案应考虑采用 E 级高性能计算机原型系统的计算、存储、网络、管理及基础设施支撑的技术路线。



图 14 超算核心硬件架构

计算系统由 CPU 和异构加速卡计算节点共同组成。计算系统建议采用统一的硬件平台来更好的执行越来越多的高性能计算任务，以提升资源利用率以及运营效率。同时计算子系统的硬件设计建议需要兼顾不同算力业务的诉求，能够提供多元化的算力支撑。在计算硬件选型方面建议采用专用硬件加速器等先进技术手段，并且结合使用液冷等散热技术来建设绿色低碳节能环保的超算计算系统。

存储系统采用分布式存储，可提供 PB 级别以上的容量来进行数据和算据存储。存储系统建议采用分布式融合存储系统，解决海量数据存储、读写速度缓慢等性能问题。

网络系统分为存储网络、业务网络以及监控网络等多个网络平面，为超算系统间各个硬件设备以及子系统间通信提供网络互连的能力。

管理网络包括资源与业务监控、告警监控、可视化等功能。资源与业务监控建议实现计算资源、应用与业务的全链路监控。告警监控建议实现故障信息以多种提醒方式快速响应运维人员以确保及时获取故障信息并响应处理。可视化建议实现基础环境、物理资源、服务、应用等各类资源数据的可视化展示，结合整体运行监控，对故障等进行实时可视化告警。

安全系统可由防火墙、负载均衡、堡垒机、抗 DDoS、日志审计、漏洞扫描、网络审计、DNS 服务器等设备组成。建议针对高性能用户数据信息资产的机密性、完整性、可用性等方面进行安全防护，重点对数据采集、数据传输、数据存储、数据使用、数据加密、数据迁移等环节进行安全防护，形成边界、网络、计算、数据等多层防护体系以抵御各种安全威胁，满足国家网络安全等保 2.0 标准要求。

（二）网络体系

超算网络系统切分成超高速计算网络、存储网络、业务网络、管理网络和监控网络等多个网络平面。具体如下：

（1）超高速计算网络：保证整个系统高效的计算和快速的数据传输，计算、存储网络需采用超高

速网络技术；

- (2) 存储网络：主要连接各存储节点，保证存储节点之间的通信；
- (3) 管理网络：为整套系统的运维管理提供网络通路；
- (4) 业务网络：主要用于与外网连接，使用户可以向超算系统传输原始数据和提取计算结果以及进行软件编译等各种操作；
- (5) 监控网络：提供带外管理网络连通，实现对所有设备的带外管理，包括远程开关机等操作。

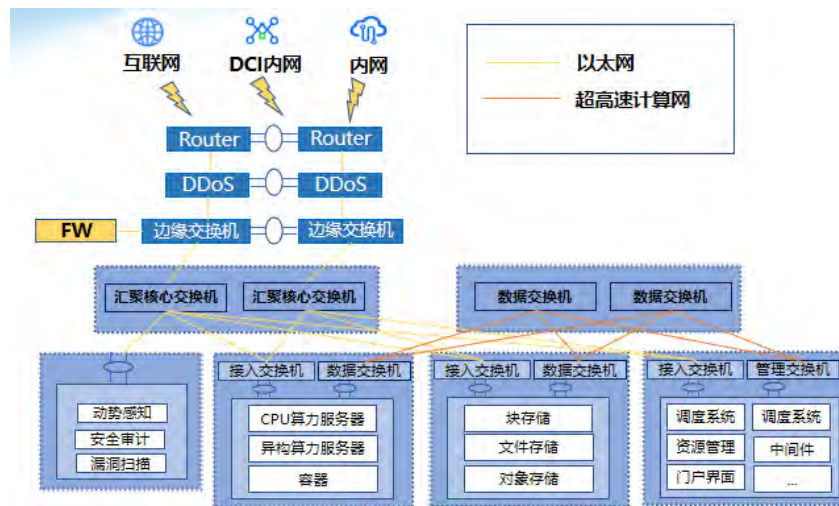


图 15 超算网络架构

（三）存储节点和高性能计算节点采用专用超高速网络

超算在多节点并行运行时会有频繁的、大量的网络数据通信，计算网络的性能对并程序的计算性能、并行加速比以及可扩展性有决定性的影响。这主要反映在两方面，其一，如果并行计算程序的数据通信以小数据包为主，且数据交换非常频繁，这一类并行程序对计算网络的延迟性能非常敏感，计算网络的延迟越低，程序的并行性能越好；其二，如果并行计算程序数据通信大数据包较多，则对计算网络的带宽性能敏感，计算网络的带宽越高，程序的并行性能越好。

目前大规模高性能计算集群均采用分布式并行存储架构，集群的规模越大，或者应用程序对存储 I/O 性能要求越高，则对并行存储系统的存储网络性能要求越高，要求存储网络具有低延迟、高带宽的特性。

为了保证整个系统高效的计算和快速的数据传输，计算、存储网络需采用超高速网络技术，还可结合 RDMA（Remote Direct Memory Access）通信，有效保障高性能计算系统内部的数据传输。RDMA 的引入，可以解决网络传输中服务器端数据处理的延迟而产生的问题。

（四）多种网络平面的融合

整个系统中存在管理网络、业务网络等多个网络平面，从节省投资简化运维的角度出发，构建基于以太网的融合网络架构，才能解决大超算生态对多平台、异构系统的吸纳。基于已有实践，建议采用接入层

带宽不小于 10G，核心层 40G，这样既能保证多网络平面融合后的性能，又具有布线简单、维护性好、可靠性高等优势，也为未来扩容预留充足空间。

3.4.3.3 超算算力的技术路线

近几年，行业的快速发展带来了强大的算力与数据流通、运算需求，给高性能计算行业带来了新的挑战与机遇，国际上纷纷兴起了 E 级计算计划，如美国在 2015 年提出的“国家战略计算项目”（National Strategic Computing Initiative, NSCI），美国能源部的 ECP（Exascale Computing Project）计划。日本很早便开始了适配 E 级计算机的 Post-K 芯片的研发。欧盟紧随其后，陆续发布其 E 级计算 EuroHPC 计划的研发路线图。我国基于自主研发的芯片与技术，打造属于自己的 E 级计算机。各国纷纷发力，企图占领高性能计算新领域高地。

当前超级计算行业，异构已成为大趋势，总结而言有三种策略：芯片内异构，异构众核（典型系统：神威太湖之光）；节点内异构，CPU+加速器结构（典型系统：Summit，天河二号）以及系统分区异构。

3.5 边缘算力

3.5.1 概述

作为一种新型的服务模型，边缘计算将数据或任务放在靠近数据源头的网络边缘侧执行处理。边缘可以从数据源到云计算中心之间的任意功能实体，这些实体搭载着融合网络、计算、存储、应用核心能力的边缘计算平台，为终端用户提供实时、动态和智能的算力。

边缘算力的发展需要云边协同解决问题。边缘计算与云计算各有所长，云计算擅长把握整体，在乎全局，非实时、长周期数据的大数据分析，能够在长周期维护、业务决策支撑等领域发挥优势；边缘计算则专注于局部，聚焦实时、短周期数据的分析，能更好地支撑本地业务的实时智能化处理与执行。云边协同将放大边缘计算与云计算的应用价值，边缘计算既靠近执行单元，更是云端所需高价值数据的采集单元，可以更好地支撑云端应用的大数据分析；反之，云计算通过大数据分析优化输出的业务规则或模型可以下发到边缘侧，边缘计算基于新的业务规则进行业务执行的优化处理。

计算需要处理的数据种类日趋多样化，边缘设备既要处理结构化数据，又要处理非结构化数据。为此，边缘算力构架需要解决不同指令集和不同芯片组成的异构计算体系中各指令集和芯片高效协同工作的问题，满足不同业务运用的需求同时，实现性能、成本、功耗、可移植性等多方面的优化均衡。对于边缘计算系统，处理器、算法和存储器为最关键的三个要素。本白皮书着重于算力范畴，将介绍处理器部分。

用于边缘计算的处理器：常规物联网终端节点的处理器是一块简单的 MCU（微控制单元，Microcontroller Unit），以控制目的为主，运算能力相对较弱。如果要在终端节点扩展边缘计算能力，可

采取两种方式：

（1）提升 MCU 的性能；

（2）通过异构计算的方法，MCU 仍保持简单的控制功能，计算部分则交给专门的加速器 IP 来完成。

对于这两种方式，第一种通用性更好，第二种则计算效率更高。目前业界趋势为两种方式并行，平台型的产品会使用第一种思路，而针对某种大规模应用的定制化产品则会采用第二种。

3.5.2 市场

3.5.2.1 国际

在过去的几年，随着 5G、物联网等技术的发展，边缘计算的技术和应用也快速向前推进，边缘设备的数量和种类与日俱增。**2021 年，边缘基础设施和终端 AI 处理器的收入规模达到了 160 亿美元，预计 2026 年将会达到 600 亿美元，年复合增长率超过 30%。**²⁸

边缘计算硬件按计算算力和存储发生的位置分为两大类：边缘设备和边缘服务器。有四个具体类型，如图 16 所示：

- 设备边缘 Device edge：设备上的板载型号和数据处理
- 内部边缘 Premise edge：企业、汽车或家庭中的本地基础设施
- 访问边缘 Access edge：连接到大型汇聚中心的网络接入点
- 城市边缘 Metro edge：主要聚合点，包括互联网服务提供商和数据中心，覆盖区域内的通讯和数据交互

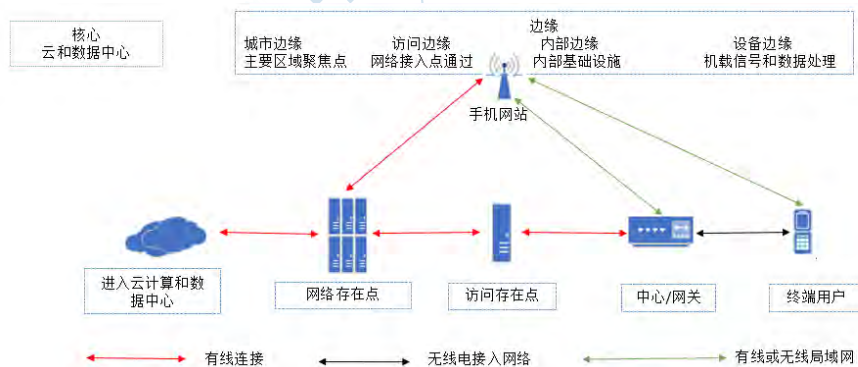


图 16 对边缘计算的解构²⁹

Device edge 和 Premise edge 目前约占市场的 75%。2024 年以后，随着 5G 和其他先进网络的推出，预计 Access edge 和 Metro edge 将继续加速增长。

²⁸ 数据来源：公开数据整理

²⁹ 数据来源：BOSTON CONSULTING GROUP: The Battle at Computing's Edge

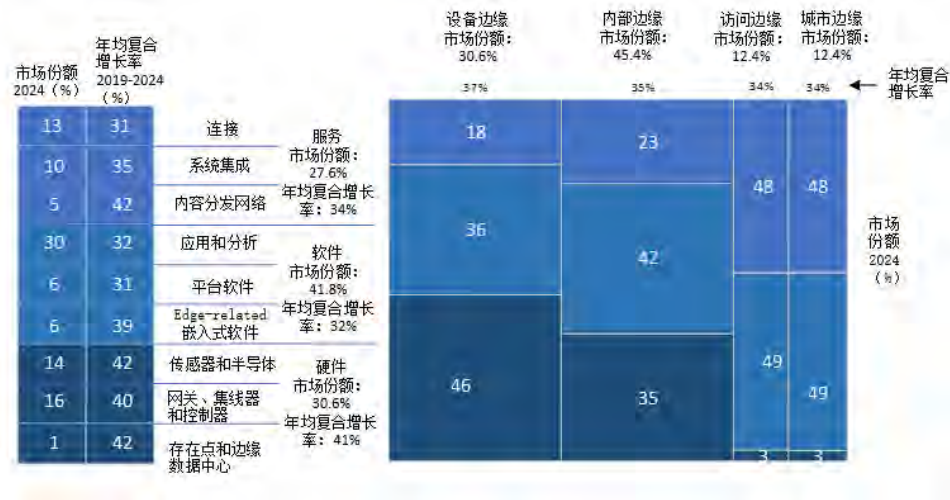


图 17 2024 年各边缘阶段预期市场份额³⁰

各大厂商也积极在边缘计算领域持续发力，图 18 列出了各大主要厂商在边缘服务器和边缘设备中的发展对比。

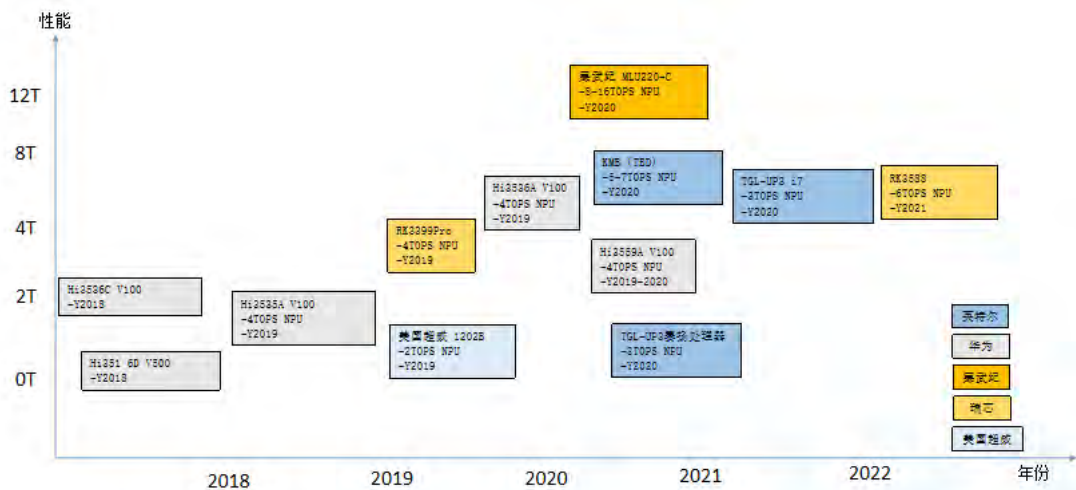


图 18 边缘设备发展对比图³¹

3.5.2.2 国内

边缘算力产业链逐步完善。边缘算力产业上游为设备供应商，产品向小型化、轻量化和集成化方向发展，主要厂家包括华为、思科、浪潮、艾默生、施耐德等；中游为边缘服务商，提供边缘网络 and 专业化集成运营服务等，主要厂家包括中国电信、中国移动、中国联通、华为、腾讯、阿里和百度等；下游为终端客户，涉及机场、国防、营销、气象、航运、保险、农业、家庭消费、健康、能源、公共服务、零售等

³⁰ 数据来源：BOSTON CONSULTING GROUP: The Battle at Computing's Edge

³¹ 数据来源：IDC 官网

多个垂直行业。

不同的边缘服务商提供的服务侧重点也有差异。电信运营商和私网运营商主要提供网络接入+数据中心服务；第三方专业服务商主要提供数据中心基础+增值服务，并附带 IaaS 服务、IT 外包或者系统集成服务以及定制化数据中心服务；大型 OTT 供应商主要专注于分布式云服务；一些新兴玩家更倾向于提供生态内+多元化租赁服务。当前国内运营商主要着眼于边缘网络应用，包括视频监控、VR/云游戏、智能制造等应用场景率先落地。

2021 年中国边缘计算服务器整体市场规模达到 33 亿美元，较 2020 年增长 24%³²。2021 年，中国垂直行业和电信网络（MEC）边缘计算服务器（含通用服务器和定制服务器）市场维持了强劲的发展势头，市场规模进一步扩大到 4.8 亿美元，同比大幅增长 75.0%。

从应用行业来看，公用事业、金融、电信和制造业需求快速增长，电信行业积极发展边缘计算平台以承载网络转型业务，包括网络功能虚拟化基础设施 NFVI、下沉核心网网元 UPF 以及多接入边缘计算 MEC 业务等；电力、金融和制造业加快部署边缘计算服务器，支撑工业机器视觉质检、设备预测性维护、线路/管道检测、智慧营业厅等场景解决方案，以边缘洞察全面赋能企业业务创新。

从出货规模来看，2021 年，中国边缘定制服务器市场排名前三的厂商依次为浪潮、华为和新华三。预计，随着靠近数据产生端的边缘应用场景逐渐丰富，对于具有特定的外形尺寸、低能耗、更宽的工作温度以及其他特定设计的边缘定制服务器的需求将快速增加，以适应复杂多样的部署环境和业务需求，在 2021 年-2026 年，中国边缘定制服务器市场预计将保持 40% 的年复合增长率。

3.5.3 技术

3.5.3.1 不同应用场景下的不同精度需求

AI 系统通常涉及训练和推理两个过程。训练过程对计算精度、计算量、内存数量、访问内存的带宽和内存管理方法的要求都非常高。而对于推理过程，更注重速度、能效、安全和硬件成本，模型的准确度和数据精度则可酌情降低。低精度、可重构的芯片设计将是趋势。

3.5.3.2 云侧和边缘侧协同发展

云侧 AI 处理器主要强调精度、处理能力、内存容量和带宽，同时追求低延时和低功耗；边缘设备中的 AI 处理器则主要关注功耗、响应时间、体积、成本和隐私安全等问题。目前云和边缘设备在各种 AI 应用中往往是协同工作。最普遍的方式是在云端训练神经网络，然后在云端（由边缘设备采集数据）或者边缘设备进行推理。

相比于面临性能与能耗瓶颈的基于云数据中心的深度学习模型部署方法，结合新兴的边缘计算技术则是更好的方式，充分运用从云端下沉到网络边缘的计算能力，在具有适当计算能力的边缘计算设备上

³² 数据来源：公开数据整理

实现低时延与低能耗的深度学习模型推理。

边缘侧的负载整合则为人工智能在边缘计算的应用找到了突破口。虚拟化技术将在不同设备上独立的负载整合到统一的高性能计算平台上，实现各个子系统在保持一定独立性的同时还能有效分享计算、存储、网络等资源。边缘侧经过负载整合，产生的节点既是数据的一个汇总结点，同时也是一个控制中心。

人工智能可以在结点处采集分析数据，也能在节点提取洞察做出决策。网络优化将是把人工智能运用到边缘侧的关键性技术之一。可以通过低比特、剪枝和参数量化等方法进行网络优化。

云侧与边缘侧相互配合，优势互补。边缘计算机非常适合在边缘处收集、存储、处理和分析数据，然而对于一些复杂的工业工作负载，边缘计算机需要配备实时处理加速器或有实时处理能力的处理器来实现实时决策处理。

3.5.3.3 芯片制造商在边缘设备上的探索

处理器架构的更新推动了边缘端机器学习的发展。以来自 Intel、NVIDIA、AMD、ARM 和 Qualcomm 等芯片制造商在边缘计算设备上的探索为例，简要介绍边缘计算的现代架构设计³³。

(1) Intel®产品如英特尔® 酷睿™ (Intel® Core™) 和英特尔凌动® (Intel Atom®)，英特尔® Movidius™ VPU 和英特尔® 至强® 可扩展 (Intel® Xeon® Scalable) 是为边缘和物联网设计的。它们有助于提供设备并将设备连接到网络，并利用 5G 实现高速度和低延迟。

(2) NVIDIA (英伟达) 发布了 Jetson Xavier NX, 相关报道认为这款产品是“最前沿最小的人工智能产品”。Jetson Xavier 比 Jetson Nano 快 2-7 倍³⁴。

(3) AMD 有更多的产品专注于物联网的嵌入式解决方案，如 AMD EPYC (霄龙) 嵌入式 3000 处理器。AMD 的第三代 EPYC 芯片的 IPC 相比 Piledriver 架构性能提升了 52%³⁵。

(4) ARM 有一个新的边缘架构，它是基于 ARM 的异构计算机，以作为英特尔和 AMD 的 x86 处理器的替代方案。它们依赖于处理器之间的低延迟、高带宽连接，而不需要太多中间存储。通过使用 SRAM，而不是 DRAM，来减少功耗。

(5) Qualcomm®宣布了三款 AI 加速卡为边缘计算提供动力: DM.2e 高达 70TOPS，功率为 15 瓦；DM.2 card 为 200 TOPS，功率为 25 瓦；PCIe card 最高达 400TOPS，功率为 75 瓦³⁶。

³³ 5 processor architectures making machine learning a reality for edge computing | Enable Architect (redhat.com)

³⁴ Jongmin Jo, Suchoel Jeong, Pilsung Kang, “Benchmarking GPU-Accelerated Edge Devices”, IEEE 2020

³⁵ <https://www.amd.com/system/files/documents/3000-family-product-brief.pdf>

³⁶ <https://www.qualcomm.com/content/dam/qcomm-martech/dm-assets/documents/Prod-Brief-QCOM-Cloud-AI-100.pdf>

4 算力应用赋能

4.1 绿色电力

4.1.1 简介

社会经济的快速发展对电能的需求越来越大，传统的发电方式产生较多的废水、废气，严重污染了自然环境。为了减少发电对环境造成的污染，绿色电力成为新的焦点。绿色电力摒弃原来火力发电的方式，借助生物能、太阳能、风能等可再生能源进行发电，新的发电方式大大降低了对环境的影响。

4.1.2 案例

4.1.2.1 项目基本情况

阿里巴巴南通数据中心采用 10kV 交流输入的直流不间断电源系统、高水温自然冷却架构、高效先进的气流组织方案、机房热回收方案等大量绿色低碳技术。

4.1.2.2 项目创新性



图 19 阿里巴巴世纪互联南通 A 栋数据中心绿色等级评估 5A

阿里巴巴南通数据中心是全球首个按“10kV 交流输入的直流不间断电源系统”设计的云计算数据中心。“10kV 交流输入的直流不间断电源系统”是一款阿里巴巴自研的、具有核心知识产权(4 项核心专利)的电源系统,通过配电链路和整流模块拓扑两个维度对原有系统进行优化,减少了系统 66%的配电环节,从而实现了最高运行效率,相比传统方式提升超过 3%。

4.1.2.3 项目赋能成效和推广价值

阿里巴巴南通数据中心通过与政府、供电局、电网、发电企业探索绿色电力实施路径,2021 年 1 月

至9月85%电量源自于天然气发电，天然气电交易使用量超过30000万千瓦时，减排二氧化碳超过16万吨。

4.2 人工智能

4.2.1 简介

人工智能是新一轮产业变革的核心驱动力，成为经济发展的新引擎，推动人类社会生活高质量发展。人工智能的快速发展，催生出了新技术、新产品、新产业、新业态、新模式，给世界带来了重大变革。当前，人工智能技术不断突破和创新，推动经济社会各领域从数字化、网络化向智能化转型。

4.2.2 案例

4.2.2.1 项目基本情况

环首都·太行山能源信息技术产业基地基于新一代AI人工智能数据中心业务场景，自然冷却（直接蒸发）冷冻水系统等自研白盒化专利技术的大规模应用，可提供高达40kW/单机柜的算力集群部署，主要服务为各类应用程序的计算和数据存储及数据支撑平台。

智能化AI服务器，采用专用AI服务器，Intel顶尖CPU处理器，搭载GPU加速技术，单机架算力高达651TFLOPS/架，算效53.7GFLOPS/W。

4.2.2.2 项目创新性

产业基地服务全球云计算领先企业，通过发展数据中心基础设施云（IaaS），平台云（PaaS）、数据云（DaaS）等云计算服务形成有效融合。数据中心部署自主可控的软硬件技术，研发关键技术和云计算解决方案，形成数据中心和云计算协同发展的生态环境，完善产业载体建设，成为数据中心领域的示范基地。

4.2.2.3 项目赋能成效和推广价值

产业基地积极实现“可再生能源→超大规模算力→新业态新场景支撑”的数字经济可再生能源零碳利用路径，为数字产业化和产业数字化大规模消纳可再生能源开辟专属线路。同时，通过环首都·太行山能源信息技术产业基地三年来的产业融合运营实践证明，扩大跨区域大数据合作，挖掘大数据在“政务、民用、商用”的应用潜力，全面释放数字经济发展新动能，形成需求拉动的良性循环，实现高质量、可持续发展。

4.3 车联网

4.3.1 简介

车联网是智能化的一体化网络。车联网是以车内网、车标网和车载移动互联网为基础，按照约定的通

信协议和数据交互标准，在车与人、路及互联网之间，进行无线通讯和信息交换的大系统网络，是能够实现智能化交通管理、智能动态信息服务和车辆智能化控制的一体化网络。

4.3.2 案例

4.3.2.1 项目基本情况

中移铁通智慧铁路边缘数据中心主要实现铁路安全员行为分析的智能化管理，对火车驾驶员手势信号进行识别并给出错误手势和空手势预警，支持驾驶员后部瞭望识别、驾驶员起立识别、驾驶员伏趴识别、驾驶员仰卧识别等一系列规章规范规定的信号动作识别，包括对场站整体视联网接入数据实时监控，发现异常情况，系统统一提供调度指挥处理信息。为铁路运营排患解难，实现安全预警、操作不规范预警等功能，不需要采用人工拷贝的手段，提高铁路运作的效率，保障列车安全行驶，为推进铁路行业数智化进程贡献一份力量。

4.3.2.2 项目创新性

在本项目实施的过程中，基于最新人工智能理论，采用领先的人工智能计算架构，通过算力的生产、聚合、调度和释放四大作业环节，支撑和引领本项目在算力等方面的提升。边缘数据中心承载 AI 算力的生产、聚合、调度和释放过程，让数据进去让智慧出来。此外，其所具有的开放标准，集约高效、普适普惠的特征，不仅能够涵盖融合更多的软硬件技术和产品，而且也极大降低了产业 AI 化的进入和应用门槛。

4.3.2.3 项目赋能成效和推广价值

本项目作为中移铁通智慧铁路边缘数据中心在 5G、工业互联网、人工智能等重点应用领域的一次探索，重点聚焦数据中心运行效率、算力网络、监控安全、网络能力等方面。在项目的构建和实施过程中与多种手段、多种新兴技术融合，推进铁路向智能化迈进的步伐。

4.4 智慧医疗

4.4.1 简介

智慧医疗的建设是基于现有医疗信息平台，利用先进的的互联网技术，整合所有卫生信息资源，形成信息高度集成的医疗系统，实现患者和医生、医疗机构之间的互动，逐步达到医疗产业信息化。

4.4.2 案例

4.4.2.1 项目基本情况

浙江联通联合合作伙伴与新昌卫健委及医共体单位基于 5G 切片和 MEC 创新网络，聚焦医学影像类的拓展应用和远程手术指导教学。该项目在中国联通绍兴市分公司建设的 5G 网络基础上，实际落地基于 5G 端到端切片技术和 MEC 的三维影像信息系统、5G 远程实时交互手术会诊系统和混合现实辅助手术规划系统。通过“网络+云+应用”整体解决方案的能力赋能医疗业务，助推医疗资源下沉，辅助提升医

学诊断治疗水平。

4.4.2.2 项目创新性

新昌医共体 5G 切片+MEC 智慧医疗项目，规划引入三维影像重建、MR 辅助手术规划、AR 远程手术指导三项 MEC 业务，应用于移动查房、医患沟通、医生之间诊疗方案沟通、远程手术指导、临床教学等场景。本案例在 MEC 边缘云部署应用软件系统与业务终端配置如下：

表 24 MEC 边缘云部署应用软件系统与业务终端配置

业务名称	业务软件及终端配置
5G 三维影像重建	1 套云影像系统软件+4 个终端（PAD、手机、电脑）
5G MR 辅助手术规划	1 套规划系统软件+1 个数据服务端，2 个 MR 终端，2 个 pad 操作端
5G AR 远程手术指导和示教	1 个交互系统软件+4 个终端+1 个 AR 第三视角采集端

4.4.2.3 项目赋能成效和推广价值

通过本次项目实践检验，基于 QOS 保障机制的无线端切片技术和承载网传输侧的 FlexE 灵活硬切片技术已经成熟，在发生流量峰值时，切片技术可保障专属用户的通信指标受其影响微乎其微。

5G 切片+MEC 构建的行业专网技术已相对成熟，并且在实际生产环境中得到使用，具备正常商业环境中已经具备独立和组合商用的能力。5G 切片+MEC 构建 5G 专网方案的可实施性和专网产品可靠性较高，可大规模在其他行业进行复制使用。

4.5 边缘计算

4.5.1 简介

近年来，数字经济已成为引领经济社会变革、推动经济高质量发展的重要引擎。在算力网络快速的发展背景下，服务业也迎来了新的发展阶段和新的模式。边缘算力技术的突破，更大程度上实现了用户对网络低时延、高性能的要求，满足用户更多应用场景的需求。

4.5.2 案例

4.5.2.1 项目基本情况

国网公司变电站多站融合边缘数据中心站是特指利用变电站富余站址、供电、通信等资源开展变电站与数据中心站的融合建设、运营和运维。

边缘计算服务平台以站址资源（变电站、营业厅、开闭站）、杆塔沟道为基础，提供热门地区算力机柜、算力计算服务（裸金属、云主机、容器云）、算力存储等 IaaS、PaaS、SaaS 等通用和定制化边缘计算服务，通过云边能力协同、云边计算协同、云边信息协同、云边数据协同、云边安全协同等能力将计算扩展至边缘端，提供电力共享站、社会共享站、融合共享站、云游戏、5G 边缘计算等应用解决方案。

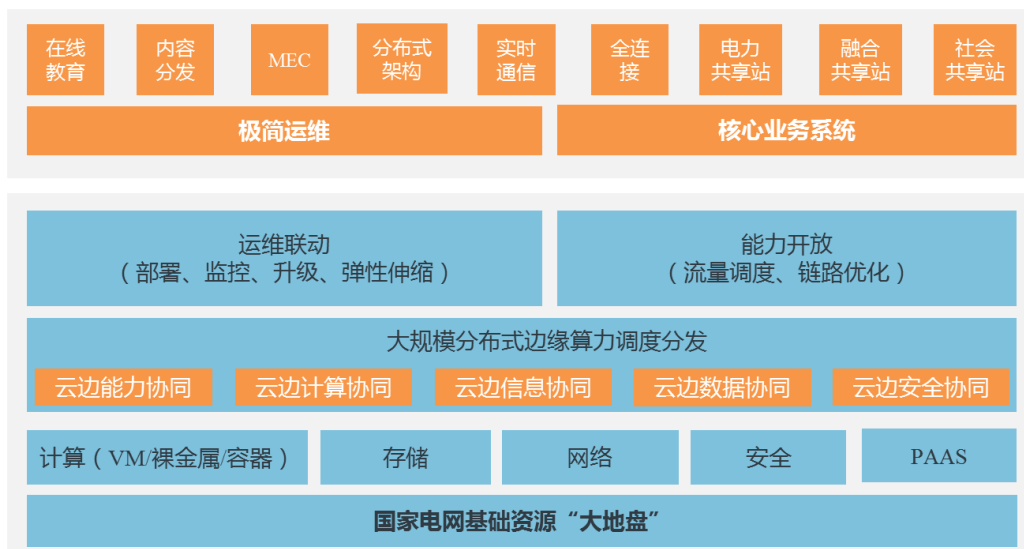


图 20 边缘计算服务平台功能架构图

4.5.2.2 项目创新性

该项目旨在打造“核心多活+边缘分布”多点部署云边协同一体化数字基础设施。在各电压等级变电站分层级部署数据中心站并全部进行光纤互联网及云化部署，实现一站式、弹性伸缩、开放灵活、安全稳定的数字基础设施服务。

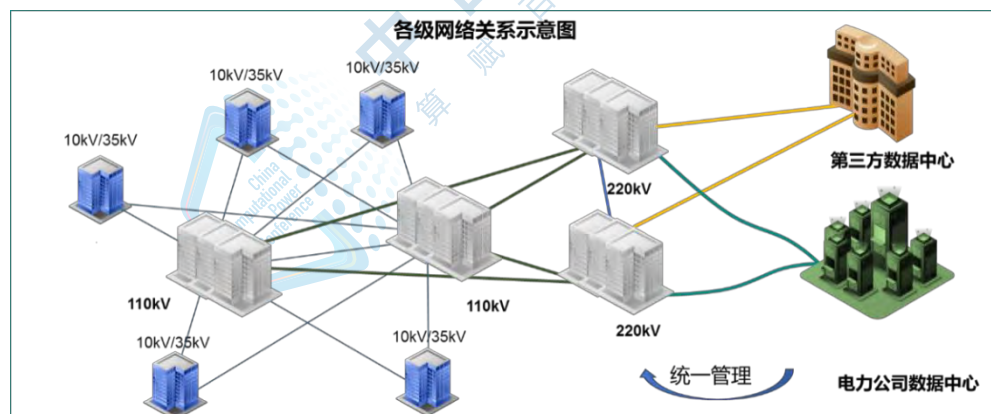


图 21 云边协同下核心多活+边缘分布示意图

4.5.2.3 项目赋能成效和推广价值

基于边缘计算开展的主要业务分为边缘算力基础设施服务、边缘云服务及行业应用综合解决方案。其中，边缘算力基础设施服务依托边缘计算数据中心的空间、位置、安全及供电保障等优势奠定了产品的差异性和独特性，而基于该服务的边缘云及解决方案类产品是重点业务方向。

4.6 物联网

4.6.1 简介

万物互联作为未来社会发展的方向，其技术进步和应用进展始终是社会关注的重点。目前信息社会的发展已经开始从“互联网+”向“万物+”进行转变，同时“万物+”技术条件已基本具备。未来，“万物+”将在大数据、云计算等技术的支撑下，挖掘万事万物的数据价值，衍生出更多新的应用场景和商业模式。

4.6.2 案例

4.6.2.1 项目基本情况

中国电信南京（吉山）云计算数据中心主要应用场景有：①天翼云平台，面向政企与个人用户，提供PasS层、IaaS层云服务；②5GCT云平台，面向全省各行业，提供“5G+”、车联网服务；③为南京证券、华泰证券、道通期货等金融类客户及百度、腾讯等大型互联网客户实现各类云平台直连核心骨干网，推动人工智能技术与产业相融合，全面赋能各行各业实现数字化转型。

4.6.2.2 项目创新性

园区为多云交换核心节点，多云交换网络，实现一跳入云，一线入多云。通过多云交换网络的建设，能够实现承载省内用户一线入多云以及云间互联业务，以及承载后续的跨地市VPN专线业务，并采用自动方式实现业务开通。

园区通过端到端场景需求分析，建设相应的网络、平台等基础设施，承载本次2C2H类端到端试点的天翼云盘（极速版）、云科学、视频创新等业务，为各业务提供计算、存储、网络访问和安全防护能力，以及基础设施运营监控与管理能力。

4.6.2.3 项目赋能成效和推广价值

中国电信南京（吉山）云计算数据中心是“天翼云、智慧云、政务云、安全云”的基础保障。其可复制经验包括：在基础设施建设方面，在土地资源紧张的区域，数据中心建设应考虑节约、集约用地；数据中心应按模块化规划，充分利用好电力资源，避免资源浪费；应积极采用新型节能技术，节能降耗，主要设备应积极选用《国家绿色数据中心先进适用技术产品目录（2020）》所推荐技术产品或类似功能及性能技术产品；项目设计及建造应积极采用BIM技术，可有效避免返工造成无谓的浪费，加快工期提高工程质量。在运营管理方面，应构建了一套全面、完善的运维管理体系，实现数据中心运维规范化、标准化、流程化、专业化、精细化管理，同时应充分利用智能化运营管理手段。

5 算力五力衡量体系

5.1 算力五力指标构建

根据 ODCC 在 2020 年发布的《数据中心算力白皮书》中定义，数据中心算力是一个包含计算、存储、传输（网络）等多个内涵的综合概念，是衡量数据中心计算能力的一个综合指标。

本白皮书，将数据中心算力的内涵进行进一步推广，融合计算、算效、存储、网络等综合概念，同时结合数据中心算力发展特点和重点影响因素，利用统计学相关方法构建衡量算力的指标体系，**算力衡量指标包含通用算力、智能算力、算效、网络和存储 5 个方面。**

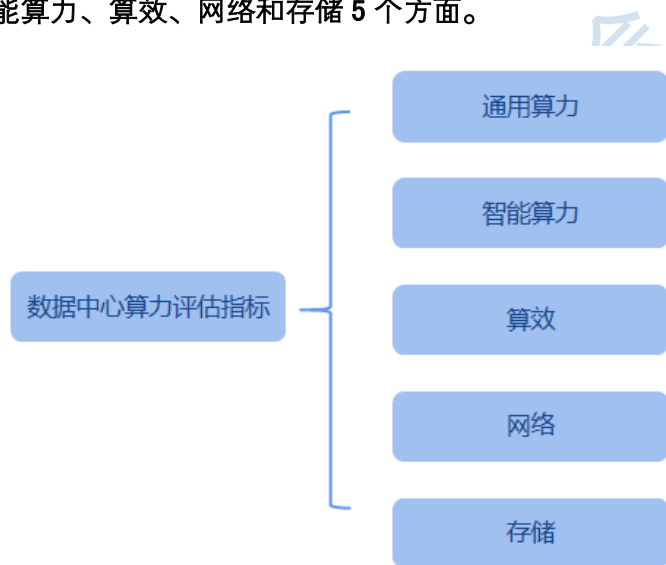


图 22 数据中心算力评估指标构建

5.2 通用算力

本白皮书将继续沿用 ODCC《数据中心算力白皮书》中定义的“每秒浮点运算次数”（Floating-point operations per second, FLOPS）来评估数据中心的通用算力，通用算力主要是基于 CPU 服务器的算力水平，一般用单精度（FP32）进行度量。

5.3 智能算力

算力和精度是密不可分的，必须在精度的前提下，讨论智能算力。由于智算领域应用的多样性，通常无法使用单一精度来衡量算力的大小。目前，智算领域常用的算力精度有：

表 25 智能算力领域算力精度

算力精度	说明
FP64	用于科学计算和高精度仿真模拟等
FP32	传统上用于深度学习训练以及图形渲染
FP16	深度学习训练和推理
TF32	新的数据格式，用于深度学习训练
BFloat16	用于深度学习训练
INT8	用于深度学习推理

本白皮书将继续使用“每秒浮点运算次数”（Floating-point operations per second, FLOPS）来评估数据中心的智能算力，为方便与通用算力比较，亦采用单精度（FP32）进行度量。

5.4 算效

将算力与功耗结合来看，单位功耗的算力是评价数据中心计算效果更为准确的一个指标。在 ODCC《数据中心算力白皮书》中，定义数据中心算效（CE）为数据中心算力与所有 IT 设备功耗的比值，是同时考虑数据中心计算性能与功耗的一种效率，即“数据中心 IT 设备每瓦功耗所产生的算力”（单位：FLOPS/W）：

$$CE = \frac{CP}{\sum IT设备功率}$$

5.5 网络

随着数据量的指数级增长，数据中心算力在呈爆发性增长的趋势，除了保障单节点内的算力之外，网络的性能对于算力的高效可扩展性，同样也发挥着更为重要的作用。当数据在网络中时，对数据进行处理和计算的最佳时间和地点就在网络中。所以，原来由网络交换机和路由器做的一些工作，现在已经转向了服务器内部；原来在服务器上做的一些工作，现在放在网卡或者交换机上更为合适。因此，在讨论算力时，亦要结合未来数据中心网络发展趋势，提供面向未来的功能和性能。

在基础设施操作方面，在网络中实现计算功能以及一些基础设施操作的卸载，也成为网络的一种技术走向。如在网络交换机上进行 All2All、Barrier 等集合操作计算，可以彻底消除这些集合通信导致的网络拥塞问题。如 AR（Adaptive Routing，动态路由）技术，可以在网络中动态的选择最优路径，主动地规避网络的拥塞；NVMeoverFabric Target 卸载技术，可以在 CPU 零消耗的情况下达到很高的 IOPS；IPSec 和 TLS 的硬件卸载，实现了前所未有的全线速数据加密传输等。

在网卡处理能力方面，大小消息对于网络性能的关注重点不同。网卡是网络和业务的接口，它需要

将应用的 Message 转变成网络可以传输的数据包，并对此进行封装和解封装以及 CRC 的计算等操作。对于大的 Message 而言，高带宽无疑是最好的利器来助其达到最佳的性能，对于小的 Message 而言，消息率（Message Rate，即网卡每秒可处理的 Message 的数量）则更为重要。这对于小数据包的传输而言至关重要，像典型的 HPC 应用、数据库应用等都是采用的小数据包的方式，可以直接从高的消息率获益。

在交换机性能方面，包转发率和转发方式至关重要。网络中的交换机需要对数据进行高速转发，由于交换机看不到 Message，只能看到每一个数据包（Packet），所以交换机的包转发率 PPS（Packet Per Second）就成为决定交换机性能的关键因素。同时在交换机上采用 Cut-Through 还是 Store/Forward 的转发方式对于交换机来言也非常重要，采用 Cut-Through 可以极大地减少对于交换机 Buffer 的占用和依赖性，而 Store/Forward 则要求交换机需要有更大的 Buffer 来缓存数据，更加容易因为 Buffer 的溢出而导致丢包的发生。

在 RDMA 技术应用方面，RDMA 技术已经成为支撑数据中心提供高性能算力的基础。通过 RDMA 技术，可以在几乎不消耗 CPU 资源的前提下实现网络的线速带宽。如何以最简单的方式在数据中心内部署 RDMA 技术，也是网络提供商需要考虑的事情之一。InfiniBand 网络的 RDMA 技术即插即用零配置是一大优势，RoCE 一键部署技术也可以降低 RDMA 的使用门槛。如何让大多数的数据中心用户能轻松地用上 RDMA 技术，是一个值得关注的问题。

在网络带宽方面，网络带宽将进一步提升。目前基于单端口 200Gb/s 的端到端网络已经被数据中心和智算中心所广泛使用，部分数据中心和智算中心还在使用多张 200Gb/s 的网卡来提升网络带宽，如我国的 E 级机的设计需求就是要求单机带宽在 400Gb/s 以上，在微软的 Azure HPC/AI 公有云的平台上，单机的网络带宽已经达到了 1600Gb/s。更高的带宽意味着更大的通路，在单位时间内可以传输更多的数据。在目前进一步降低通信延时遇到的物理极限瓶颈前提下，以高带宽来换低延时就成为网络竞争的一个热点。目前，InfiniBand 网络已经达到了单端口 400Gb/s 的端到端带宽，在 2023 年底，将达到单端口 800Gb/s。

在网络时延方面，时延决定计算效率。算力网络时延需要向总线级看齐，提供低时延传输能力，降低计算任务完成时间，提升计算效率。在数据中心网络相关性能的测试中，往往从网络任务完成时间、网络跳数（组网架构时延）、动态时延几个方面开展测试。

（1）网络任务完成时间。针对常用的集合通信场景，将集合通信中的部分计算卸载到网络中，减少消息入网的次数，降低时延，提升通信效率，缩短计算任务的完成时间。

（2）网络跳数。与传统组网时延性能对比测试，直连拓扑组网较两级 CLOS 组网时延有较大优势，直连拓扑架构时延最大降低 32%。

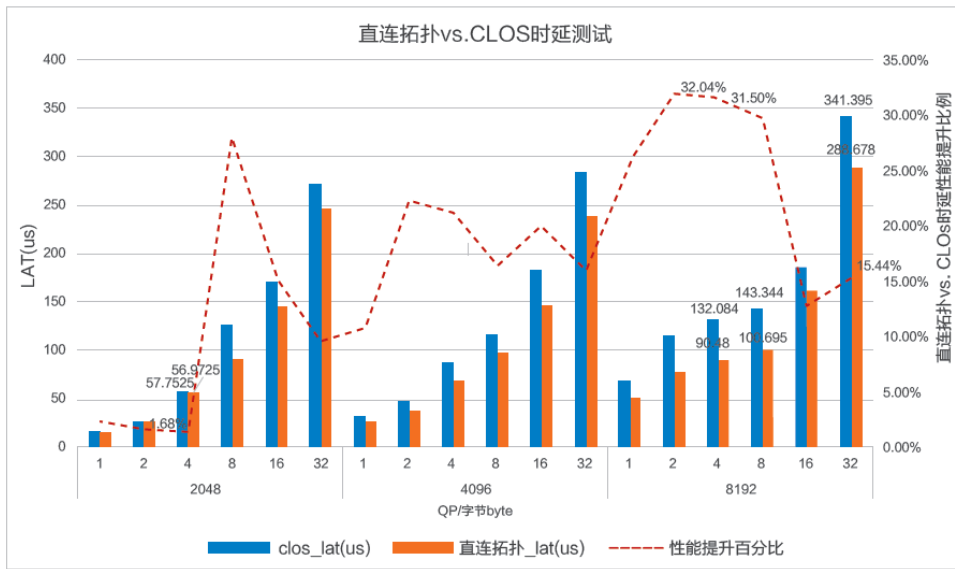


图 23 直连拓扑 vs. CLOS 时延测试

直连拓扑和 CLOS 组网 OpenFoam 计算业务测试，直连拓扑较 CLOS 组网性能有较大优势，最大性能提升 29.02%。

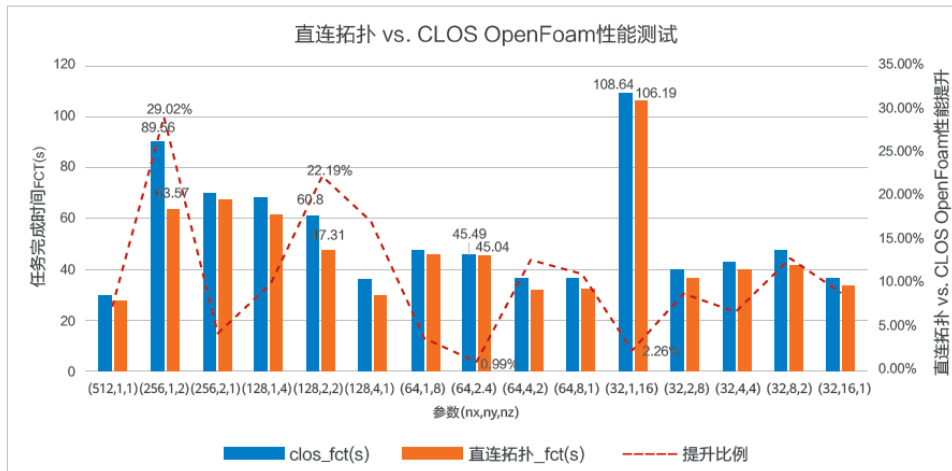
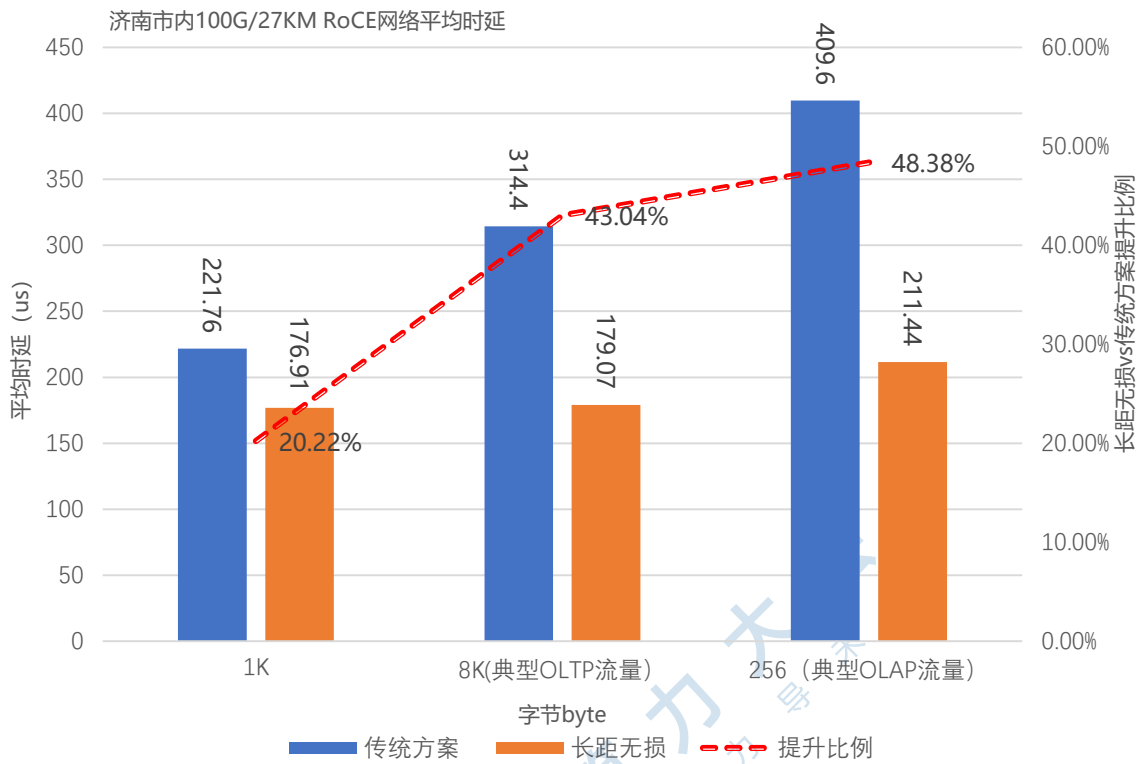


图 24 直连拓扑 vs. CLOS OpenFoam 性能测试

(3) 动态时延 (长距无损) 测试。传统的长距离无损传输依赖于 PFC (基于优先级队列的流控) 帧，PFC 帧在长距链路上传输的飞行时间内，接收端需要大量的缓存接收在这飞行时间段内的传输的报文。同时保证链路传输的吞吐不受影响，则需要接收端有更大的缓存。大缓存则必然导致报文在设备中传输时延增大。创新的长距无损算法，在小缓存设备上很好的解决了长距无损传输的问题。既保证了接收端在小缓存的情况下，不出现拥塞和时延增长的情况，也保证了传输链路的吞吐。100G 场景下，长距无损算法对比传统 PFC 方案，网络带宽打满的情况下，平均时延最高可降低 48.38%。



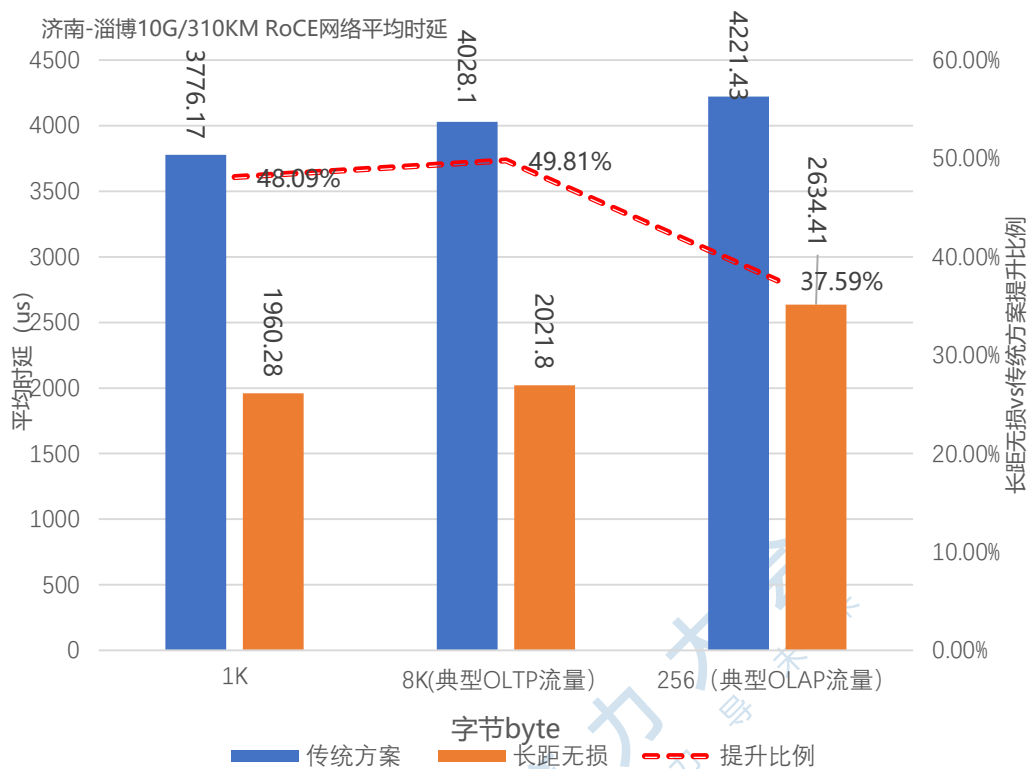
数据来源：华为测试数据

图 25 济南市内 100G/27KM RoCE 网络平均时延³⁷

300KM 场景下，长距无损算法对比传统 PFC 方案，网络带宽打满的情况下，平均时延最高可降低 49.81%。



³⁷ 数据来源：华为测试数据



数据来源：华为测试数据

图 26 济南-淄博 10G/310KM RoCE 网络平均时延³⁸

总而言之，衡量网络的指标需要各个层面来综合考虑，最终选择合适的网络，助力于算力的性能提升。

5.6 存储

存储系统作为数据中心基础架构的组成部分，在算力准备、实现和交付的过程中起着重要的保障作用。

目前数据中心的存储系统的装机形态主要有服务器内置存储器、机柜内连接存储服务器、机房布置的专用存储服务器机柜。使用的存储介质装机量从以硬盘为主转向以 SSD 为主，介质接口从以 SATA、SAS 为主转向以 NVME/PCIe 为主。传统的应用服务器，由于处理的本地化固定数据较多，因此数据容错处理大多数采用 RAID 卡来管理本机内或者同一个机柜内的以 SAS、SATA 接口为主的硬盘和 SSD。对带宽和 IOPS 要求较高的服务器则会使用 PCIe/NVME 接口的本机安装的 SSD 作为主要存储设备。

以算力为中心的存储系统，不同于传统的数据中心存储系统，有着对带宽、IOPS、时延、存储容量、

³⁸ 数据来源：华为测试数据

可扩展性、高并发性以及容错能力、权限管理、安全管理等有其特有的要求。因此相对较慢的 SATA, SAS 接口的存储阵列会成为性能的瓶颈，而本机安装的 PCIe/NVME 存储部件在功耗、散热上会出现问题，因此全闪存阵列或者闪存阵列和硬盘阵列的组合形态将成为主要的存储设备形态。

针对以算力为中心的基础架构，其存储系统应具备短项目提交时间和快速的项目部署能力，在计算项目结束后能够快速完成计算过程数据和结果数据的保存、迁移和归档，并能提供结果的高效查询和输出和校验。而动态调整存储容量、带宽、响应时间以及对各层物理器件失效的容灾设计应当计入设计和考察的范畴。

目前的数据中心存储器件中，NVME SSD、SAS SSD、SATA SSD 的带宽需要在高并发多队列状态才能发挥出来，而 HDD 则需要在低并发的顺序操作状态才能达到。因此存储系统的分级设计和读写过程的调度对整个算力系统的性能表现有着重要的影响。基于以上需求，提交、部署阶段的存储系统考核指标为带宽，在项目计算阶段的存储系统的考核指标为 IOPS 和时延、计算结果保存、交付和查询阶段的考核指标为带宽，存储容量以及可扩展性和高并发性应作为以算力为中心的存储系统的衡量指标。

最后，在存储介质装机量的数量级达到万以上的数据中心中，存储单元的故障会成为日常发生的状况，因此运营系统中的故障自动侦测、故障预警及数据自动迁移、日常的存储装置健康报告生成和汇总也需要作为必要衡量指标。

综上，对存储系统算力的度量需要体现存储器件单体的带宽、IOPS、时延、容量、每日可写入量，也需要体现系统的连接通道带宽，系统各存储层级之间连接的带宽配合程度。



6 算力五力模型

对于单体数据中心算力如何度量和表示，业界一直在不断地探索。数据中心算力是一个包含计算、存储、传输（网络）等多个内涵的综合概念，是衡量数据中心计算能力的一个综合指标，所以目前亟需一个科学的方法来将算力的内涵充分表达出来。

多属性群决策（Multi-Attribute Group Decision Making, MAGDM）是多属性决策（Multi-Attribute Decision Making, MADM）和群体决策（Group Decision Making, GDM）交叉的研究方向，是现代决策科学的一个重要研究领域，其理论和方法已广泛应用于城市规划、经济管理、投资风险等领域。当下，对于不同数据中心间算力综合各类指标的评估问题可以归结为多属性群决策问题。

为了综合评价单体数据中心的算力，本章提出“算力五力模型”，将与计算能力高度相关的通用算力、智能算力、算效能力、存储能力、网络能力等不同量纲的指标进行融合比对。利用新的双向投影法及TOPSIS方法计算得到样本的相对贴近度，进而对不同样本的算力进行定级，解决了不同量纲变量难以直接比较的困难，提高了评估结果的综合性和有效性。本章为算力评估体系提供新的模型和方法，更好地指导和建议业界判断行业发展趋势，为未来算力规划和部署提供思路³⁹。

6.1 相关概念

设多属性决策中的方案集为 $A = \{a_1, a_2, \dots, a_s\}$ ，属性集为 $C = \{c_1, c_2, \dots, c_m\}$ 。记 $S = \{1, 2, \dots, s\}$ ， $M = \{1, 2, \dots, m\}$ ，用 $c_l(A)$, $l \in M$ 表示方案 A 在 c_l 属性下的值，并允许出现不可比较的情况。

定义 1⁴⁰ 如果关系 $\{>, <, \approx, ?\}$ 满足下述条件：在多属性决策问题中，每个属性下方案间的序数偏好信息可以由决策者给出，对于属性 c_l , $l \in M$ 任意两个方案 a_i 和 a_j ($i \neq j$) 满足下列偏好关系之一：

方案 a_i 优于方案 a_j ($a_i > a_j$)；方案 a_i 与方案 a_j 优劣相当 ($a_i \approx a_j$)；

方案 a_i 劣于方案 a_j ($a_i < a_j$)；方案 a_i 与方案 a_j 优劣关系不明 ($a_i ? a_j$)。

且两方案 a_i 和 a_j (满足 $i \neq j$) 之间有且仅有 $\{>, <, \approx, ?\}$ 中的一个关系成立，则称 $\{>, <, \approx, ?\}$ 构成了一个偏序偏好结构。

设在多属性决策问题中，各个不同的属性下各方案的属性值被划分为 $n+1$ 个等级，从而将一个方案比另一个方案好的程度分为 n 个等级。不妨记 $a_i \overset{k}{>} a_j$, ($i, j \in S; 0 \leq k \leq n$)，表示方案 a_i 优于方案 a_j k 个

³⁹ 吴美希,杨晓彤.算力五力模型：一种衡量算力的综合方法[J].信息通信技术与政策,2022(03):13-21.

⁴⁰ 张金清.序方法与均衡分析[M].上海：复旦大学出版社,2003.

等级； $a_i \overset{-k}{>} a_j, (i, j \in S; 0 \leq k \leq n)$ 与 $a_i \overset{k}{<} a_j, (i, j \in S; 0 \leq k \leq n)$ 均表示方案 a_i 劣于方案 a_j k 个等级。下面给出广义优序数的概念。

定义 2⁴¹ 令

$$a_{ijl} = \begin{cases} 1, c_l(A_i) \overset{n}{>} c_l(A_j), i \neq j \\ \frac{n}{2n-k}, c_l(A_i) \overset{k}{>} c_l(A_j), i \neq j \\ 0.5, c_l(A_i) \approx c_l(A_j), i \neq j \\ 0.375, c_l(A_i) ? c_l(A_j), i \neq j \\ 0, i = j \\ -\frac{n}{2n-k}, c_l(A_i) \overset{-k}{>} c_l(A_j), i \neq j \\ -1, c_l(A_i) \overset{-n}{>} c_l(A_j), i \neq j \end{cases} \quad (6-1)$$

其中 $i, j \in S, l \in M$ 。称 a_{ijl} 为属性 c_l 下方案 A_i 相对于方案 A_j 的广义优序数。

为了使决策结果更加真实可靠，本文将算力五力模型结合向量的概念，通过向量的变换和相关计算，应用双向投影法和 TOPSIS 方法，通过比较相对贴近度得到相关结论。

定义 3 设 $X_p = [r_{pj}, \bar{r}_{pj}]$ 和 $X_q = [r_{qj}, \bar{r}_{qj}]$ 为集合 X 上的两个不确定等级变量，如果我们将 X_p 和 X_q 通过广义优序数转变为 $X_p = [\xi(r_{pj}), \xi(\bar{r}_{pj})]$ 和 $X_q = [\xi(r_{qj}), \xi(\bar{r}_{qj})]$ ，其中， $\xi(r_{pj})$ 为 r_{pj} 经广义优序数转化的值，其它类似符号同义， X_p 与 X_q 形成的向量定义如下：

$$X_p X_q = [\min \xi(r_{pj}), \max \xi(\bar{r}_{pj})] \quad (6-2)$$

其中：

$$\begin{aligned} \min \xi(r_{pj}) &= \min (|\xi(r_{qj}) - \xi(r_{pj})|, |\xi(\bar{r}_{qj}) - \xi(\bar{r}_{pj})|), \\ \max \xi(\bar{r}_{pj}) &= \max (|\xi(r_{qj}) - \xi(r_{pj})|, |\xi(\bar{r}_{qj}) - \xi(\bar{r}_{pj})|) \end{aligned} \quad \circ$$

例 1 不妨令 $A = \{a_1, a_2, \dots, a_m\}$ 为方案集， $C = \{c_1, c_2, \dots, c_n\}$ 为属性集，将属性 $c_l (1 \leq l \leq n)$ 划分为 6 个等级，有等级比较集 $R = \{r_i | i = -5, \dots, 0, \dots, 5, ?\}$ ，其中?表示两方案优劣关系不明的情况， $t = 5$ 。已知两个等级变量分别为 $X_p = [r_2, r_4]$ ， $X_q = [r_0, r_4]$ 。

借助等级优序数可以将 $X_p = [r_2, r_4]$ ， $X_q = [r_0, r_4]$ 转化为 $X_p = [\frac{5}{8}, \frac{5}{6}]$ ， $X_q = [\frac{1}{2}, \frac{5}{6}]$ ；进而通过定义 3 计算出 X_p 与 X_q 形成的向量为 $X_p X_q = [0, \frac{1}{8}]$ 。

⁴¹ 张小芝, 朱传喜. 多属性决策的广义等级偏好优序法[J]. 系统工程理论与实践, 2013, 33(11): 2853-2858.

6.2 双向投影模型

双向投影法具有决策过程科学合理、简单易行和区分度高等特点，并且可以避免传统投影法无法处理现实问题中备选方案在正、负理想解垂直平分线上的情况，在实际决策中得到广泛推广。本节接下来将主要介绍基于该方法的双向投影模型。

令 $X_i = [\xi(r_{ij}^t), \xi(\bar{r}_{ij}^t)]$ 为第 j 个属性 c_j 下对方案 A_i 经广义优序数转换的向量。

首先，在第 j 个属性下，将正、负理想解表示为： $X^+ = [\max_{1 \leq i \leq n} \xi(r_{ij}), \max_{1 \leq i \leq n} \xi(\bar{r}_{ij})]$ ， $X^- = [\min_{1 \leq i \leq n} \xi(r_{ij}), \min_{1 \leq i \leq n} \xi(\bar{r}_{ij})]$ ，其中 n 表示备选方案的数量；正负理想解形成的向量可以表示为： $X^- X^+ = [\xi(r_{ij}^t), \xi(\bar{r}_{ij}^t)]$ ，其中：

$$\xi(r_{ij}^t) = \min \left(\left(\max_{1 \leq i \leq n} \xi(r_{ij}) - \min_{1 \leq i \leq n} \xi(r_{ij}) \right), \left(\max_{1 \leq i \leq n} \xi(\bar{r}_{ij}) - \min_{1 \leq i \leq n} \xi(\bar{r}_{ij}) \right) \right)$$

$$\xi(\bar{r}_{ij}^t) = \max \left(\left(\max_{1 \leq i \leq n} \xi(r_{ij}) - \min_{1 \leq i \leq n} \xi(r_{ij}) \right), \left(\max_{1 \leq i \leq n} \xi(\bar{r}_{ij}) - \min_{1 \leq i \leq n} \xi(\bar{r}_{ij}) \right) \right)$$

然后，在第 j 个属性下， X_i 与正、负理想解形成的向量可以表示为： $X^- X_i = [\xi(r_{ij}^-), \xi(\bar{r}_{ij}^-)]$ ， $X_i X^+ = [\xi(r_{ij}^+), \xi(\bar{r}_{ij}^+)]$ ，其中，

$$\xi(r_{ij}^-) = \min \left(\left(\xi(r_{ij}) - \min_{1 \leq i \leq n} \xi(r_{ij}) \right), \left(\xi(\bar{r}_{ij}) - \min_{1 \leq i \leq n} \xi(\bar{r}_{ij}) \right) \right);$$

$$\xi(\bar{r}_{ij}^-) = \max \left(\left(\xi(r_{ij}) - \min_{1 \leq i \leq n} \xi(r_{ij}) \right), \left(\xi(\bar{r}_{ij}) - \min_{1 \leq i \leq n} \xi(\bar{r}_{ij}) \right) \right);$$

$$\xi(r_{ij}^+) = \min \left(\left(\max_{1 \leq i \leq n} \xi(r_{ij}) - \xi(r_{ij}) \right), \left(\max_{1 \leq i \leq n} \xi(\bar{r}_{ij}) - \xi(\bar{r}_{ij}) \right) \right)$$

$$\xi(\bar{r}_{ij}^+) = \max \left(\left(\max_{1 \leq i \leq n} \xi(r_{ij}) - \xi(r_{ij}) \right), \left(\max_{1 \leq i \leq n} \xi(\bar{r}_{ij}) - \xi(\bar{r}_{ij}) \right) \right)$$

这些向量的模可以计算为：

$$|X^- X^+| = \sqrt{\sum_{j=1}^m \left(\left(\xi(r_{ij}^t) \right)^2 + \left(\xi(\bar{r}_{ij}^t) \right)^2 \right)};$$

$$|X^- X_i| = \sqrt{\sum_{j=1}^m \left(\left(\xi(r_{ij}^-) \right)^2 + \left(\xi(\bar{r}_{ij}^-) \right)^2 \right)};$$

$$|X_i X^+| = \sqrt{\sum_{j=1}^m \left(\left(\xi(r_{ij}^+) \right)^2 + \left(\xi(\bar{r}_{ij}^+) \right)^2 \right)}.$$

那么余弦值 $\cos(X^- X_i, X^- X^+)$ ， $\cos(X^- X_i, X^- X^+)$ 可以分别表示为：

$$\cos(X^-X_i, X^-X^+) = \frac{\sum_{j=1}^m (\xi(r_{ij}^t) \cdot \xi(r_{ij}^-) + \xi(\bar{r}_{ij}^t) \cdot \xi(\bar{r}_{ij}^-))}{|X^-X_i| \cdot |X^-X^+|} \quad (6-3)$$

$$\cos(X^-X_i, X^-X^+) = \frac{\sum_{j=1}^m (\xi(r_{ij}^t) \cdot \xi(r_{ij}^+) + \xi(\bar{r}_{ij}^t) \cdot \xi(\bar{r}_{ij}^+))}{|X_iX^+| \cdot |X^-X^+|} \quad (6-4)$$

进而向量 X^-X_i 到向量 X^-X^+ 和向量 X^-X^+ 到向量 X_iX^+ 的投影值分别表示为：

$$\begin{aligned} prj_{X^-X^+}(X^-X_i) &= |X^-X_i| \cdot \cos(X^-X_i, X^-X^+) \\ &= \frac{\sum_{j=1}^m (\xi(r_{ij}^t) \cdot \xi(r_{ij}^-) + \xi(\bar{r}_{ij}^t) \cdot \xi(\bar{r}_{ij}^-))}{|X^-X^+|} \end{aligned} \quad (6-5)$$

$$\begin{aligned} prj_{X_iX^+}(X^-X^+) &= |X^-X^+| \cdot \cos(X_iX^+, X^-X^+) \\ &= \frac{\sum_{j=1}^m (\xi(r_{ij}^t) \cdot \xi(r_{ij}^+) + \xi(\bar{r}_{ij}^t) \cdot \xi(\bar{r}_{ij}^+))}{|X_iX^+|} \end{aligned} \quad (6-6)$$

注 1 如图 27 所示， $prj_{X^-X^+}(X^-X_i)$ 的值越大，则说明方案 X_i 越接近正理想解 X^+ ；同理， $prj_{X_iX^+}(X^-X^+)$ 越大，说明方案 X_i 越接近负理想解 X^- 。

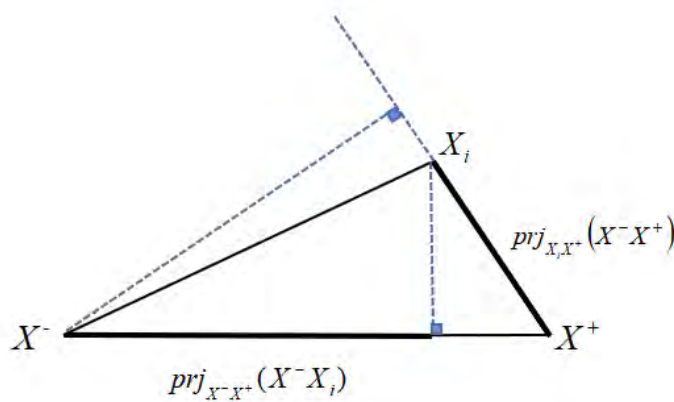


图 27 $prj_{X_iX^+}(X^-X^+)$ 和 $prj_{X^-X^+}(X^-X_i)$ 的图形表示

最终，分别计算出方案 X_i 与理想解的相对贴适度：

$$C(X_i) = \frac{prj_{X^-X^+}(X^-X_i)}{prj_{X^-X^+}(X^-X_i) + prj_{X_iX^+}(X^-X^+)} \quad (6-7)$$

表示每个方案的投影值到正理想解的相对贴适度，其中 $prj_{X^-X^+}(X^-X_i)$ ， $prj_{X_iX^+}(X^-X^+)$ 分别表示向量 X^-X_i 到向量 X^-X^+ 和向量 X^-X^+ 到向量 X_iX^+ 的投影值。

基于以上方法，等级双向投影模型的具体步骤如下：

步骤 6.2.1：构造具有等级变量的决策矩阵 D ，并对其进行规范化处理；

步骤 6.2.2：通过以上广义优序数运算将等级变量进行相关转换；

步骤 6.2.3：确定正理想解 $X^+ = \{X_1^+, X_2^+, \dots, X_m^+\}$ 和负理想解 $X^- = \{X_1^-, X_2^-, \dots, X_m^-\}$ ；

步骤 6.2.4: 分别计算正负理想解形成的向量 X^-X^+ 和 X_i 与正、负理想解形成的向量 X^-X_i , X_iX^+ ;

步骤 6.2.5: 根据 (6-5) 式和 (6-6) 式, 分别计算向量 X^-X_i 到向量 X^-X^+ 和向量 X^-X^+ 到向量 X_iX^+ 的投影值 $prj_{X^-X^+}(X^-X_i)$, $prj_{X_iX^+}(X^-X^+)$;

步骤 6.2.6: 由 (6-7) 式计算方案 X_i 与理想解的相对贴近度;

步骤 6.2.7: 将相对贴近度进行大小排序(由大到小), 选取贴近度最大的方案为最优方案。

6.3 算力五力模型

本文选取 5 个指标进行算力评估, 分别为: 通用算力、智能算力、算效能力、网络能力、存储能力。

表 26 算力五力指标选取

指标	说明
通用算力	本文通用算力的单位采用 TFLOPS (FP32, 单精度浮点算力)。
智能算力	本文智能算力的单位采用 TFLOPS (FP32)。
算效能力	本文算效的单位采用 GFLOPS/W (FP32)。
网络能力	本文主要采用网络带宽速度来衡量网络的性能, 单位为 Mbit/s (兆比特每秒), 即每秒传输的比特位数。
存储能力	本文采用每秒读写次数 (Input/Output Operations Per Second, IOPS) 来衡量存储的性能。

本节将信息与双向投影法结合, 给出利用双向投影法处理不同数据中心算力方面不同指标间信息的决策模型, 构建算力五力模型, 并通过实例分析, 说明该决策模型具有良好的效果。模型具体步骤如下:

步骤 6.3.1: 根据数据将指标划分等级。将已有的数据取最大和最小值, 按需要的等级数等距划分即可, 一般越理想的一端标记等级数越高。

步骤 6.3.2: 将数据进行定级, 并转化为优序数 (0-1 之间的实数)。

步骤 6.3.1 划分的等级数表对现有数据对号入座, 确定相应的等级数, 每个指标下, 将不同数据中心两两对比, 判断出一数据中心在某指标下优于或劣于另一数据中心的等级数, 并通过优序数计算公式将各等级数转变到 0-1 之间, 方便接下来处理。

步骤 6.3.3: 取各属性下不同备选方案对应优序数的最大值和最小值为正、负理想解, 并根据定义 3 中公式进行计算, 进而确定出相应的向量;

步骤 6.3.4: 计算每个向量对应的模, 得出余弦值 $\cos(X^-X_i, X^-X^+)$ 、 $\cos(X^-X_i, X^-X^+)$, 进而分别计算向量 X^-X_i 到向量 X^-X^+ 和向量 X^-X^+ 到向量 X_iX^+ 的投影值 $prj_{X^-X^+}(X^-X_i)$, $prj_{X_iX^+}(X^-X^+)$;

步骤 6.3.5: 由式 (6-7) 计算得出每个数据中心最终的相对贴近度;

步骤 6.3.6: 根据相对贴近度的值进行大小排序, 并据相对贴近度的值对每个数据中心进行打星, 最终得出结论。

6.4 算例

现给出 6 个数据中心对应的 5 个指标（通用算力、智能算力、算效、网络、存储）数据, 如表 27 所示, 根据专家的经验 and 数据分布的情况, 评估步骤如下:

表 27 样本数据

编号	通用算力 (单位: TFLOPS)	智能算力 (单位:TFLOPS)	算效 (单位:GFLOPS/W)	网络 (单位: Mbit/s)	存储 (单位:IOPS)
DC.1	61.25	547.23	72	901	8.1
DC.2	137.3	545.03	40.3	1699.3	12.7
DC.3	23.9	50.61	23.9	1218.4	3.6
DC.4	20.2	350.32	100	964.9	15.4
DC.5	127.07	640.65	80.44	376.3	3.3
DC.6	16.65	88.35	13.24	467.6	4.5

步骤 6.4.1: 根据数据将所选的 5 个指标进行分级, 如表 28 所示;

表 28 数据中心能效指标分级

编号	等级 I	等级 II	等级 III	等级 IV	等级 V
通用算力 (单位: TFLOPS)	0-10	10-20	20-50	50-100	100-200
智能算力 (单位:TFLOPS)	0-90	90-130	130-230	230-500	500-1000
算效 (单位:GFLOPS/W)	0-10	10-20	20-30	30-50	50-100
网络 (单位: Mbit/s)	0-300	300-500	500-900	900-1500	1500-2200
存储 (单位:IOPS)	0-1	1-3	3-7	7-12	12-20

步骤 6.4.2: 将表 27 的数据根据表 28 数据中心能效指标分级表进行定级, 在每个指标下, 每个数据中心的等级数两两进行比较, 得到相对等级数并转化为优序数, 如下表 29、表 30、表 31 所示;

表 29 数据中心能效指标定级

编号	通用算力	智能算力	算效	网络	存储
DC.1	4	5	5	4	4
DC.2	5	5	4	5	5
DC.3	3	1	3	4	3
DC.4	3	4	5	4	5
DC.5	5	5	5	2	3
DC.6	2	1	2	2	3

表 30 不同数据中心不同指标间的两两比较

编号	通用算力	智能算力	算效	网络	存储
DC.1	$[r_0, r_{-1}, r_1, r_1, r_{-1}, r_2]$	$[r_0, r_0, r_4, r_1, r_0, r_4]$	$[r_0, r_1, r_2, r_0, r_0, r_3]$	$[r_0, r_{-1}, r_0, r_0, r_2, r_2]$	$[r_0, r_{-1}, r_1, r_{-1}, r_1, r_1]$
DC.2	$[r_1, r_0, r_2, r_2, r_0, r_3]$	$[r_0, r_0, r_4, r_1, r_0, r_4]$	$[r_{-1}, r_0, r_1, r_{-1}, r_{-1}, r_2]$	$[r_1, r_0, r_1, r_1, r_3, r_3]$	$[r_1, r_0, r_2, r_0, r_2, r_2]$
DC.3	$[r_{-1}, r_{-2}, r_0, r_0, r_{-2}, r_1]$	$[r_{-4}, r_{-4}, r_0, r_{-3}, r_{-4}, r_0]$	$[r_{-2}, r_{-1}, r_0, r_{-2}, r_{-2}, r_1]$	$[r_0, r_{-1}, r_0, r_0, r_2, r_2]$	$[r_{-1}, r_{-2}, r_0, r_{-2}, r_0, r_0]$
DC.4	$[r_{-1}, r_{-2}, r_0, r_0, r_{-2}, r_1]$	$[r_{-1}, r_{-1}, r_3, r_0, r_{-1}, r_3]$	$[r_0, r_1, r_2, r_0, r_0, r_3]$	$[r_0, r_{-1}, r_0, r_0, r_2, r_2]$	$[r_1, r_0, r_2, r_0, r_2, r_2]$
DC.5	$[r_1, r_0, r_2, r_2, r_0, r_3]$	$[r_0, r_0, r_4, r_1, r_0, r_4]$	$[r_0, r_1, r_2, r_0, r_0, r_3]$	$[r_{-2}, r_{-3}, r_{-2}, r_{-2}, r_0, r_0]$	$[r_{-1}, r_{-2}, r_0, r_{-2}, r_0, r_0]$
DC.6	$[r_{-2}, r_{-3}, r_{-1}, r_{-1}, r_{-3}, r_0]$	$[r_{-4}, r_{-4}, r_0, r_{-3}, r_{-4}, r_0]$	$[r_{-3}, r_{-2}, r_{-1}, r_{-3}, r_{-3}, r_0]$	$[r_{-2}, r_{-3}, r_{-2}, r_{-2}, r_0, r_0]$	$[r_{-1}, r_{-2}, r_0, r_{-2}, r_0, r_0]$

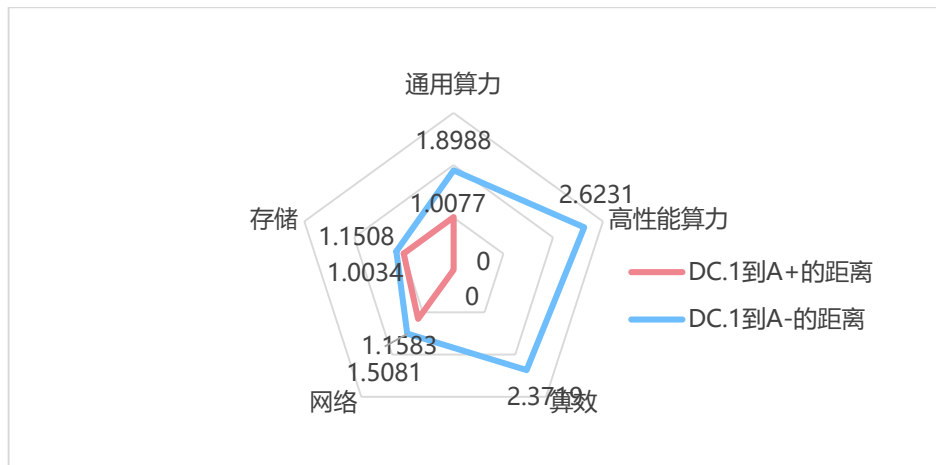
表 31 数据中心等级比较级经优序数转化表

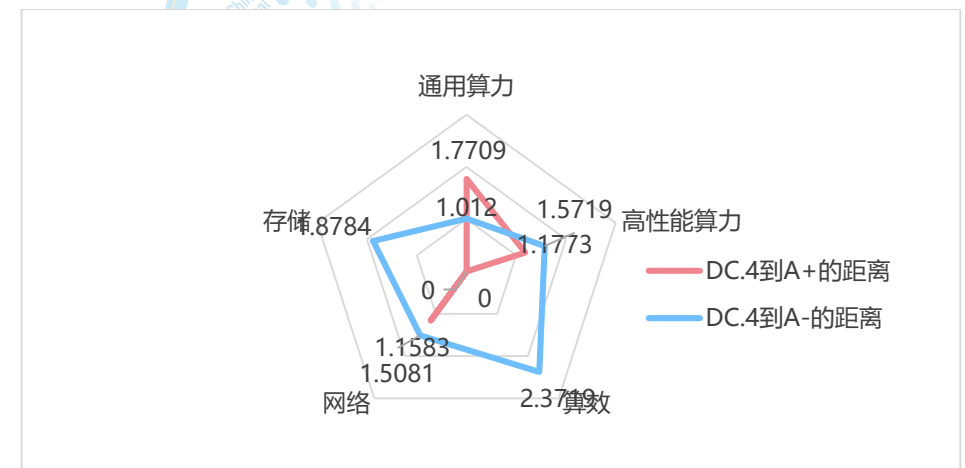
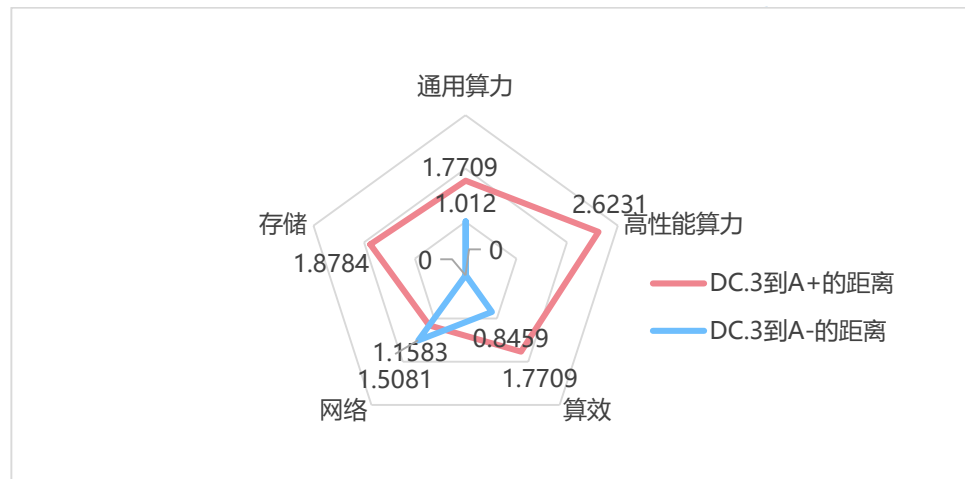
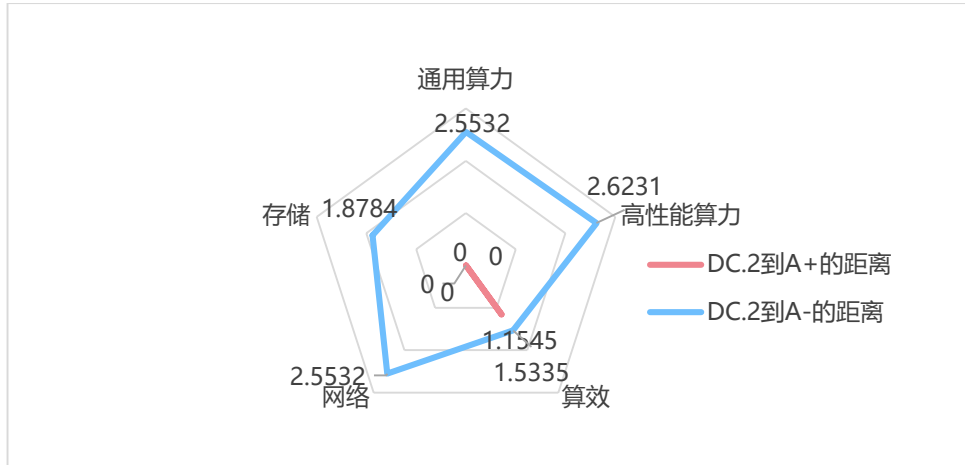
	通用算力	智能算力	算效	网络	存储
DC.1	$[0, -0.57, 0.57, 0.57, -0.57, 0.67]$	$[0, 0, 1, 0.57, 0, 1]$	$[0, 0.57, 0.67, 0, 0, 0.8]$	$[0, -0.57, 0, 0, 0.67, 0.67]$	$[0, -0.57, 0.57, -0.57, 0.57, 0.57]$
DC.2	$[0.57, 0, 0.67, 0.67, 0, 0.8]$	$[0, 0, 1, 0.57, 0, 1]$	$[-0.57, 0, 0.57, -0.57, -0.57, 0.67]$	$[0.57, 0, 0.57, 0.57, 0.8, 0.8]$	$[0.57, 0, 0.67, 0, 0.67, 0.67]$
DC.3	$[-0.57, -0.67, 0, 0, -0.67, 0.57]$	$[-1, -1, 0, -0.8, -1, 0]$	$[-0.67, -0.57, 0, -0.67, -0.67, 0.57]$	$[0, -0.57, 0, 0, 0.67, 0.67]$	$[-0.57, -0.67, 0, -0.67, 0, 0]$
DC.4	$[-0.57, -0.67, 0, 0, -0.67, 0.57]$	$[-0.57, -0.57, 0.8, 0, -0.57, 0.8]$	$[0, 0.57, 0.67, 0, 0, 0.8]$	$[0, -0.57, 0, 0, 0.67, 0.67]$	$[0.57, 0, 0.67, 0, 0.67, 0.67]$
DC.5	$[0.57, 0, 0.67, 0.67, 0, 0.8]$	$[0, 0, 1, 0.57, 0, 1]$	$[0, 0.57, 0.67, 0, 0, 0.8]$	$[-0.67, -0.8, -0.67, -0.67, 0, 0]$	$[-0.57, -0.67, 0, -0.67, 0, 0]$
DC.6	$[-0.67, -0.8, -0.57, -0.57, -0.8, 0]$	$[-1, -1, 0, -0.8, -1, 0]$	$[-0.8, -0.67, -0.57, -0.8, -0.8, 0]$	$[-0.67, -0.8, -0.67, -0.67, 0, 0]$	$[-0.57, -0.67, 0, -0.67, 0, 0]$

步骤 6.4.3: 确定正、负理想解，并计算确定出相应的向量；

$$\begin{aligned}
 X^+ &= \{[0.57,0,0.67,0,0.8], [0,0,1,0.57,0,1], [0,0.57,0.67,0,0,0.8],, \\
 & \quad [0.57,0,0.57,0.57,0.8,0.8], [0.57,0,0.67,0,0.67,0.67]\}; \\
 X^- &= \{[-0.67, -0.8, -0.57, -0.57, -0.8,0], [-1, -1,0, -0.8, -1,0], \\
 & \quad [0.8,0.67, -0.57,0.8,0.8,0], [-0.67, -0.8, -0.67, -0.67,0,0], \\
 & \quad [-0.57, -0.67,0, -0.67,0,0]\}; \\
 X^-X^+ &= \{[0.8,1.24], [1,1.37], [0.8,1.24], [0.8,1.24], [0.67,1.14]\}; \\
 X^-X_1 &= \{[0.23,1.14], [1,1.37], [0.8,1.24], [0.23,0.67], [0.10,0.57]\}; \\
 X_1X^+ &= \{[0.10,0.57], [0,0], [0,0], [0.13,0.57], [0.10,0.57]\}; \\
 X^-X_2 &= \{[0.8,1.24], [1,1.37], [0.23,1.14], [0.8,1.24], [0.67,1.14]\}; \\
 X_2X^+ &= \{[0,0], [0,0], [0.10,0.57], [0,0], [0,0]\}; \\
 X^-X_3 &= \{[0.10,0.57], [0,0], [0.10,0.57], [0.23,0.67], [0,0]\}; \\
 X_3X^+ &= \{[0.23,1.14], [1,1.37], [0.23,1.14], [0.13,0.57], [0.67,1.14]\}; \\
 X^-X_4 &= \{[0.10,0.57], [0.43,0.8], [0.8,1.24], [0.23,0.67], [0.67,1.14]\}; \\
 X_4X^+ &= \{[0.23,1.14], [0.2,0.57], [0,0], [0.13,0.57], [0,0]\}; \\
 X^-X_5 &= \{[0.8,1.24], [1,1.37], [0.8,1.24], [0,0], [0,0]\}; \\
 X_5X^+ &= \{[0,0], [0,0], [0,0], [0.8,1.24], [0.67,1.14]\}; \\
 X^-X_6 &= \{[0,0], [0,0], [0,0], [0,0], [0,0]\}; \\
 X_6X^+ &= \{[0.8,1.24], [1,1.37], [0.8,1.24], [0.8,1.24], [0.67,1.14]\}.
 \end{aligned}$$

根据表 31 确定出正、负理想解 X^+ 、 X^- ，并将每个数据中心在每个指标下的测度与正、负理想解距离计算出来，绘出图 28，在下面组图中，红色线圈表示每个数据中心的每个指标距离相对正理想解的距离。若红色线圈越小，表示该数据中心各指标相对距离正理想解越近，则该数据中心越优；反之相反。蓝色线圈表示每个数据中心的每个指标距离相对负理想解的距离。蓝色线圈越大，表示该数据中心的指标距离负理想解越远，则该数据中心越优；反之相反。在组图中，我们可以较为直观的看出 DC.6 整体情况最不好，DC.1，DC.2 整体相对更优。





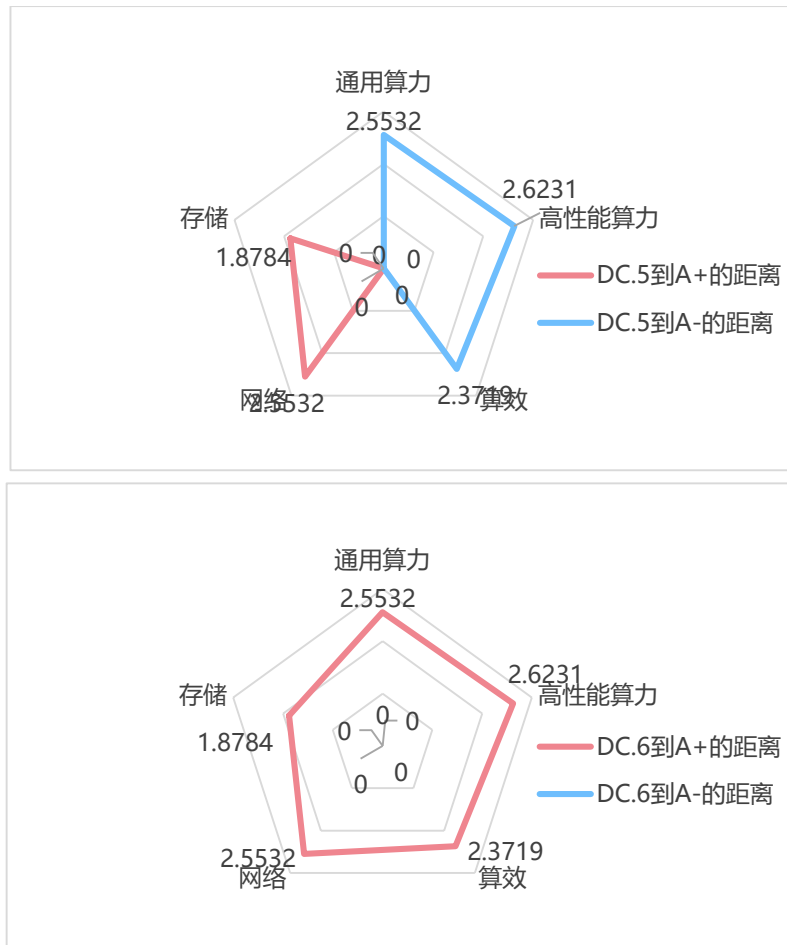


图 28 每个数据中心各指标下到正负理想解的距离

步骤 6.4.4: 计算相应的投影值;
得到相对应的模分别为:

$$|X^-X^+| \approx 3.3392; |X^-X_1| \approx 2.6915; |X_1X^+| \approx 1.0077;$$

$$|X^-X_2| \approx 3.2149; |X_2X^+| \approx 0.5793; |X^-X_3| \approx 1.0807;$$

$$|X_3X^+| \approx 2.7735; |X^-X_4| \approx 2.3621; |X_4X^+| \approx 1.4385;$$

$$|X^-X_5| \approx 2.6882; |X_5X^+| \approx 1.9808; |X^-X_6| \approx 0;$$

$$|X_6X^+| \approx 3.3392$$

在此，将每个数据中心得到的模（整体距离正、负理想解的距离）绘制成图 29，图中，红色线圈对应的数据中心越靠近中心位置，表示该数据中心整体距离正理想解越近，整体情况就越好。蓝色线圈对应的数据中心越靠近图中的外延，表示该数据中心整体距离负理想解越远，整体情况便越好。可较为直观地看出，图中 DC.1, DC.2, DC.4, DC.5 相对较好，DC.3, DC.6 相对较差，这与最后的出的结论不谋而合。

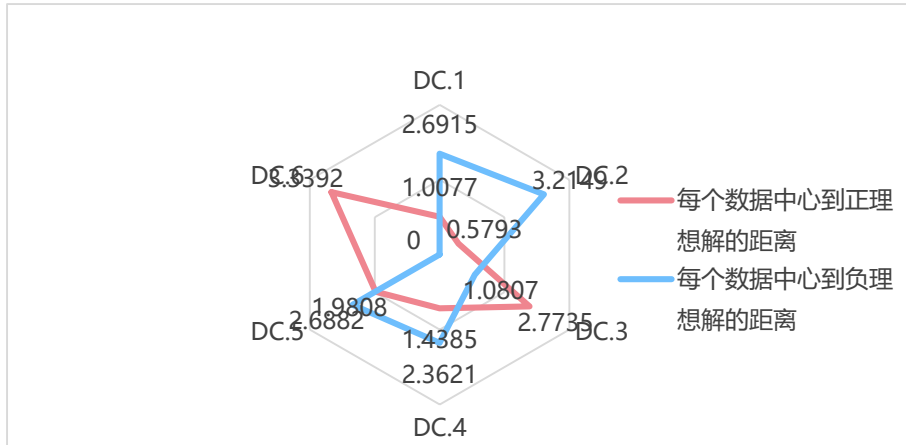


图 29 样本数据到正、负理想解的距离

进而求得投影值如下：

$$Pr j_{X^-X^+}(X^-X_1) \approx 2.5085; Pr j_{X_1X^+}(X^-X^+) \approx 2.2966;$$

$$Pr j_{X^-X^+}(X^-X_2) \approx 3.1669; Pr j_{X_2X^+}(X^-X^+) \approx 1.3528;$$

$$Pr j_{X^-X^+}(X^-X_3) \approx 0.7713; Pr j_{X_3X^+}(X^-X^+) \approx 3.1156;$$

$$Pr j_{X^-X^+}(X^-X_4) \approx 2.1685; Pr j_{X_4X^+}(X^-X^+) \approx 2.3606;$$

$$Pr j_{X^-X^+}(X^-X_5) \approx 2.1642; Pr j_{X_5X^+}(X^-X^+) \approx 1.9808; Pr j_{X^-X^+}(X^-X_6) = 0; Pr j_{X_6X^+}(X^-X^+) \approx 3.3391。$$

步骤 6.4.5：根据公式计算得出各数据中心的相对贴近度；

$$C(X_1) \approx 0.5220; C(X_2) \approx 0.7007; C(X_3) \approx 0.1984;$$

$$C(X_4) \approx 0.4788; C(X_5) \approx 0.5221; C(X_6) = 0。$$

步骤 6.4.6：根据相对贴近度的值进行大小排序，得出结论。

表 32 相对贴近度定星标准表

贴近度	0-0.15	0.15-0.30	0.30-0.50	0.50-0.70	0.70-1.0
等级	★	★★	★★★	★★★★	★★★★★

对应可得

DC.1 ★★★★★

DC.2 ★★★★★★

DC.3 ★★★

DC.4 ★★★

DC.5 ★★★★★

DC.6 ★

7 结语和展望

7.1 算力规模方面

全球算力整体规模呈现快速增长的趋势。各国均在推动 CPU、GPU 等异构算力发展，提高算力算效水平。中国和美国算力整体规模增长趋势更为明显，算力能力位列第一梯队，日本、德国、英国、加拿大、法国位列第二梯队。

未来随着数据指数级的增长，各行业对算力的需求日益增加，算力算效水平将会进一步提升，算力在经济发展中所起到的作用也会越来越重要，各国将会继续加大对算力的投资规模。

7.2 算力技术方面

无论是服务器还是芯片和系统性能，其技术均在不断的创新和突破。各供应商研发的服务器 CPU 制程不断提升，内核和主频在一定程度上有所提高，实现系统性能倍增；各生产商不断推出新的架构，基于新的架构设计 GPU 芯片、FPGA 芯片和 AI 芯片，为终端用户提供低功耗、低价格、高性能、高可靠性的产品；从最新发布国际超算算力 TOP500 榜单和国内 TOP100 榜单来看，超算计算机性能有显著的提高，运算速度越来越快。

7.3 算力指标构建方面

算力正在从一个隐形的、潜在的竞争力源泉逐步转变为现实的竞争力，算力已成为数字经济时代的关键生产力，是衡量一个国家数字经济水平的重要指标。由于数据中心的算力受多方面因素的影响，单一指标很难准确的、全面的评估数据中心算力能力。为了能够量化算力指标，准确地衡量数据中心算力，本白皮书从通用算力、智能算力、算效、网络和存储 5 个方面构建算力指标体系，也称为“算力五力模型”。

“算力五力模型”为数据中心算力评估体系提供新的模型和方法，更好地指导和建议业界判断行业发展趋势，为未来算力规划和部署提供思路。

算力作为数字时代的核心竞争力，是推动数字产业高质量发展的重要生产力，是推动企业智能业务转型的关键因素，在数字经济发展的进程中起着重要的支撑作用。未来，算力规模会继续扩大，算力技术会继续创新，算力指标体系会更加完善。

2022中国算力大会

算 赋 百 业 · 力 导 未 来



CAICT算力
公众号

中国信息通信研究院 云计算与大数据研究所
数据中心团队

地址：北京市海淀区知春路1号学院国际大厦

邮编：100191

电话：010-62300095/18810669396