

IPv6超算专网和 IPv6人工智能专网

李星

2023-11-29

大纲

- 背景
- 行业动态
- IPv6超算专网
- IPv6人工智能专网
- 小结



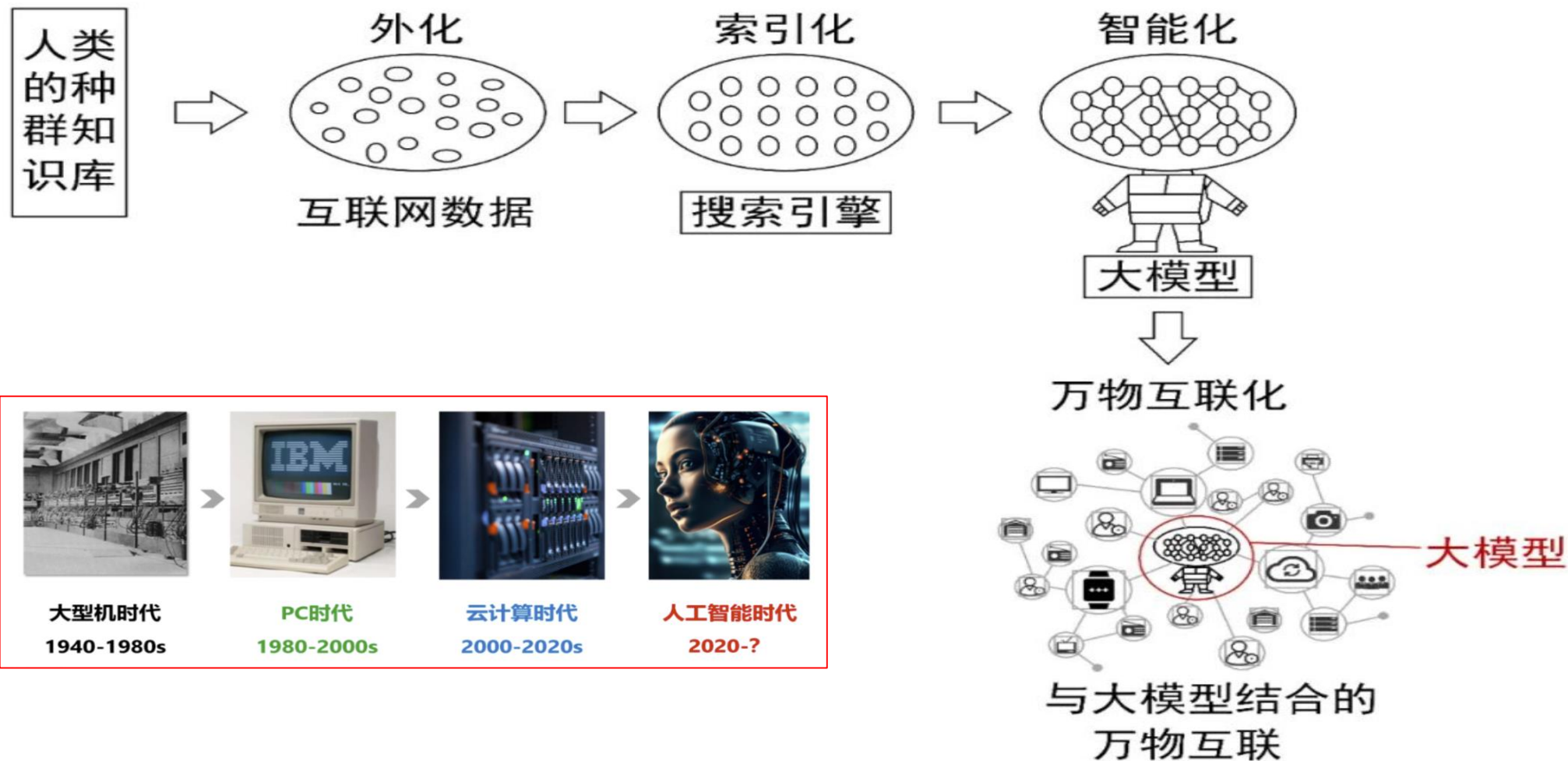
背景

Two sweet fruits

Sweet Fruits from
One Seed: AI &
Networking

Steve Crocker
Annual Conference
of World Internet
History
7 November 2023

人类知识演进和计算演进



三大算力

10个全国一体化大数据中心



算力

通用算力

消费互联网、行业互联网、政府互联网等领域的常规计算应用。

超算算力

科学计算类：物理化学、气象环保、生命科学、石油勘探、天文探测等。
工程计算类：计算机辅助工程、计算机辅助制造、电子设计自动化、电磁仿真等。

智算算力

人工智能计算，包括：机器学习、深度学习、数据分析等。

13个国家超算中心

- 北京
- 上海
- 广州
- 天津
- 济南
- 深圳
- 长沙
- 太原
- 无锡
- 成都
- 郑州
- 青岛
- 昆山

1 国家超级计算天津中心	7 国家超级计算昆山中心
2 国家超级计算济南中心	8 国家超级计算青岛中心
3 国家超级计算广州中心	9 国家超级计算郑州中心
4 国家超级计算深圳中心	10 国家超级计算成都中心
5 国家超级计算长沙中心	11 国家超级计算西安中心(筹)
6 国家超级计算无锡中心	12 国家超级计算太原中心(筹)

20余个人工智能中心

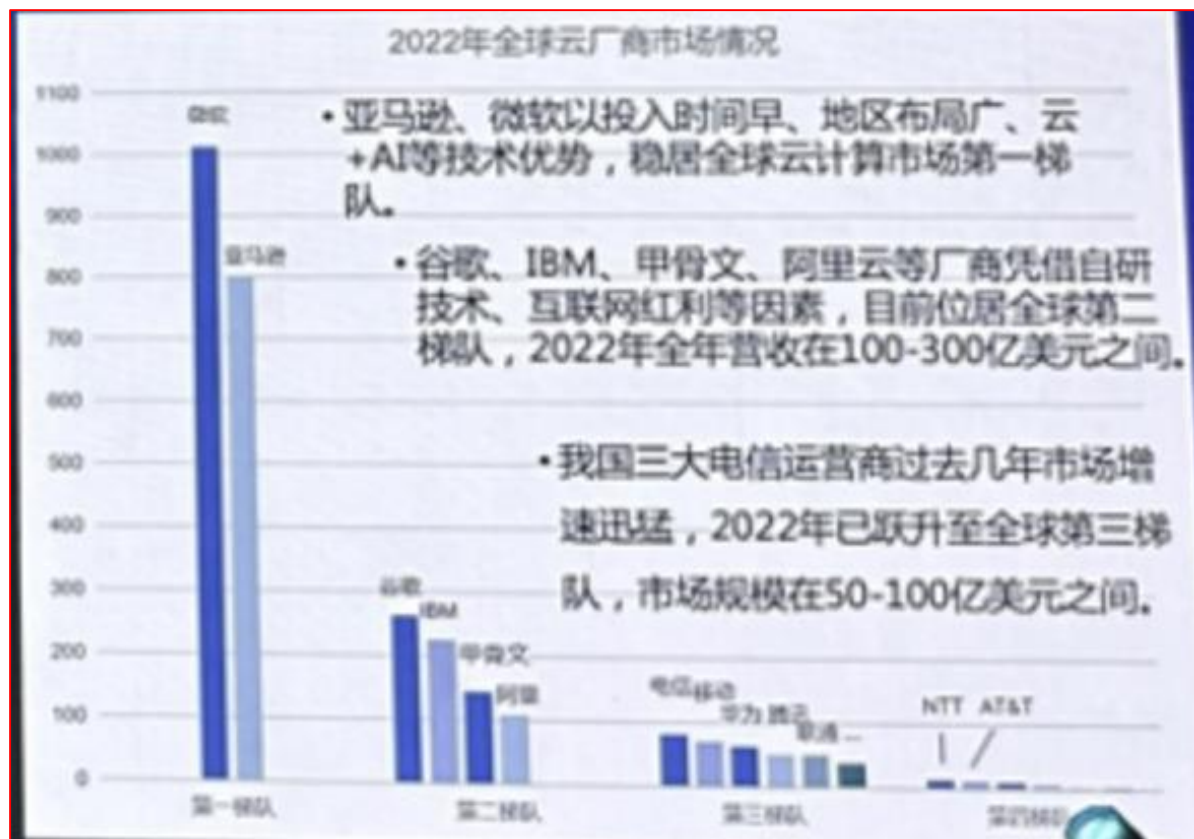


需求比较

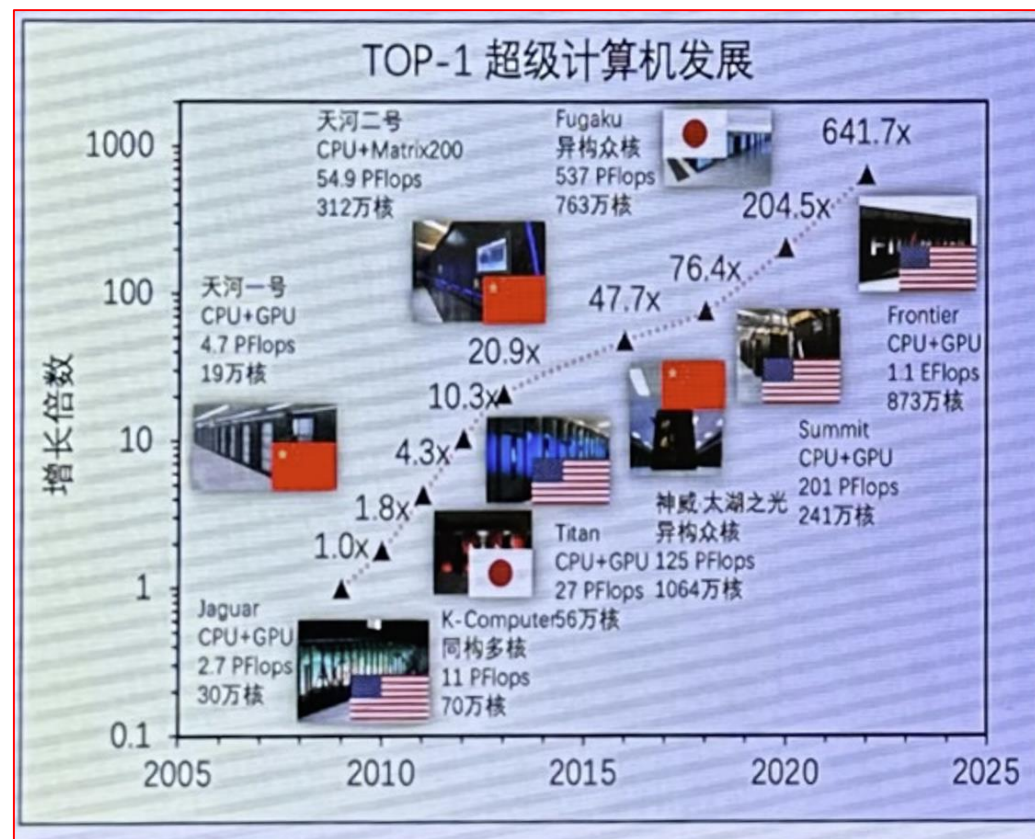
- 数据中心算力系统
 - 计算：通用CPU为主
 - 网络：高速以太网
 - 存储：计算与存储融合的分布式文件系统
- HPC算力系统
 - 计算：CPU或GPU
 - 网络：低延时、高带宽
 - 存储：与计算分离的网络存储系统
- AI算力系统
 - 计算：AI加速芯片或GPU
 - 网络：NV link + 高速以太网
 - 存储：与计算分离的网络存储系统

云排序和超算排序

云



超算





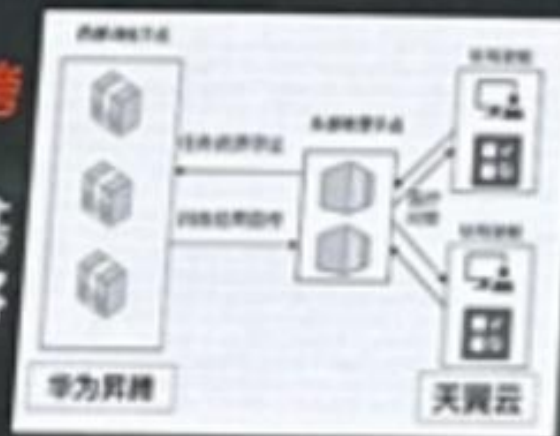
行业动态

东数西算场景

不同地方算力互联

· “东数西算”政策驱动下，东西优势互补加速跨区域算力互联。

东西部算力互联，如：“东数西训”场景下，基于各类算力枢纽、数据中心等异构算力资源，将用户的AI任务需求和数据传递至西部进行模型训练并将结果回传至东部，帮助实现降本增效。



东数西训

东数西渲

东数西存

东数西备

信通院：算力互联互通

算力互联互通 CAICT 中国信通院

中国信通院联合产业各方提出，在**互联网**体系架构上增加**算力标识、算力调度**等功能并增强**高性能传输协议**，将不同类型、不同主体、不同地域的算力互联，实现算力资源智能感知、实时发现、按需获取，形成**算力标准化、服务化的大市场**和**算力相互连接、灵活调用的一张网**。

标准体系

- 算力互联互通能力要求 第1部分：总体框架
- 算力互联互通能力要求 第2部分：需求与场景
- 算力互联互通能力要求 第3部分：业务互通及数据流动
- 算力互联互通能力要求 第4部分：算力标识及计量结算
- 算力互联互通能力要求 第5部分：算力并网
- 算力互联互通能力要求 第6部分：算力调度
- 算力互联互通能力要求 第7部分：GPU架构互通及智算性能
- 算力互联互通能力要求 第8部分：DPU架构互通及性能
- 算力互联互通能力要求 第9部分：高性能远程直接访问(HRDA)技术
- 算力互联互通能力要求 第10部分：多级网络互通效能和性能

一套标准：算力标识体系实现算力资源感知汇聚

一套标准开源实现：算网云任务编排、调度、传输和开发软件栈

一套标准开源项目的工程化实现：多层次算力互联互通平台

运营商的思路：算力网络

▶为算力互联互通平台提供技术支持，在东数西算、工业仿真、智慧教育等领域提供大规模算力调度能力。



中国电信

▶基于算网大脑，中国移动已经构筑了东视西渲、东数西存、东数西存等应用，为用户提供即开即用的算网云一体化服务。



中国移动

▶分布式云落地了智慧城市、智慧医疗等多套解决方案，为其用户提供一站式靠近终端侧客户、覆盖全面的算网服务。



中国联通

用云计算实现底层GPU、DPU等算力互联互通和算网一体化调度

▶天翼云“息壤”提供算力传输枢纽能力，可实现对全国算力、网络资源进行统筹、编排和调度。

▶中国移动算网大脑实现算网各域资源协同编排管理；打造COCA算力框架，构建算力的标准化接入生态

▶联通云分布式云打造具备全域算网需求解析、算网资源统一编排调度的能力。

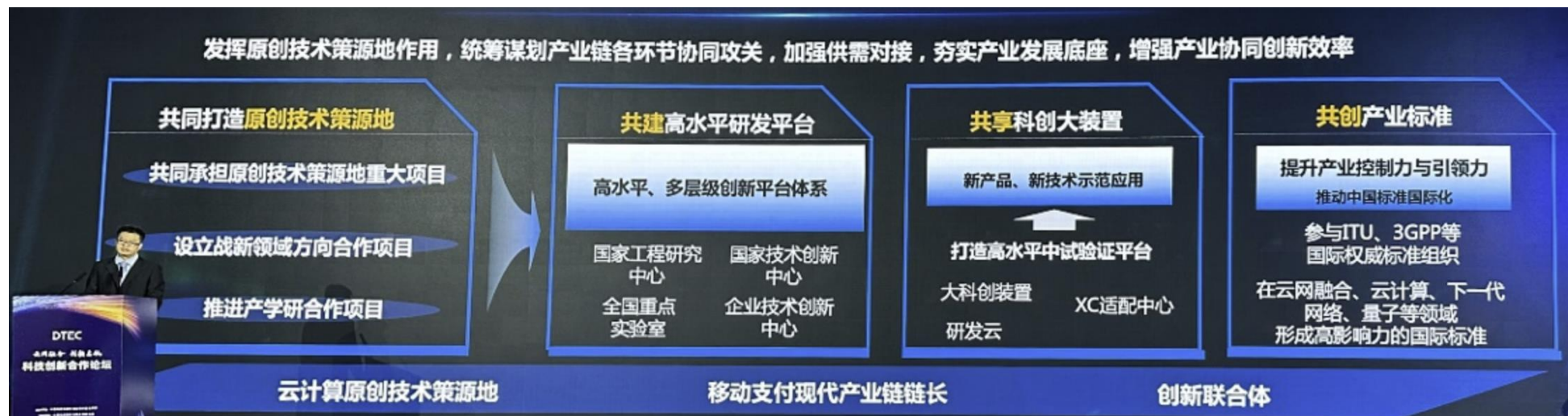
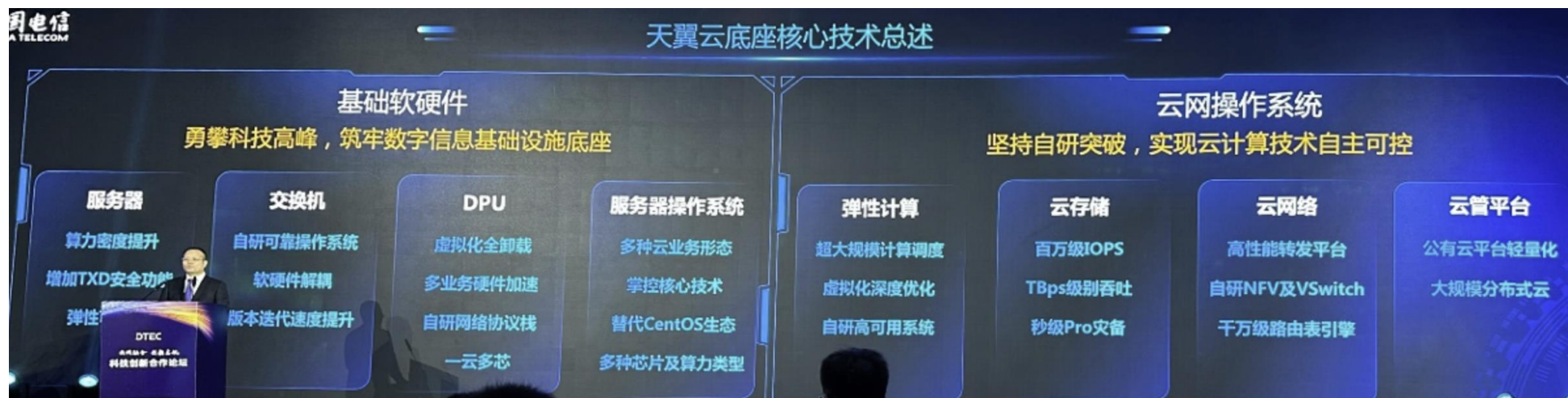
▶天翼云打造“2+4+31+X+O”的云/IDC布局体系，“一城一池”覆盖超过240座城市，边缘算力节点超800个。
▶目前中国电信在2+4区域拥有数据中心机架规模40万，光达到80%

▶移动云打造“一朵云”的全域资源布局，通过“4+N+31+X”集约化梯次布局，完成资源池省份100%覆盖。

▶联通云打造“5+4+31+X”布局，形成规模化区域中心集聚效应，实现分布式云网协同。构建全国算力时延圈，覆盖国内334个地海外40多个TOP10云厂商。

借云卖网，借网卖云

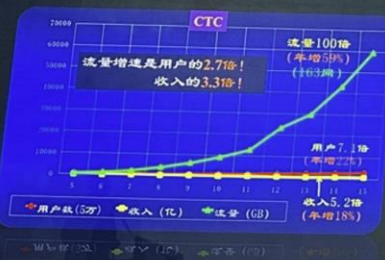
中国电信：天翼云



中国电信：虚拟局域网

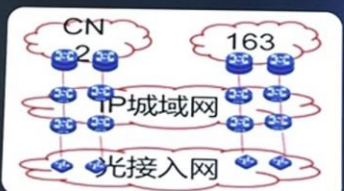
宽带和移动市场需要新增长极，加大内容填充

宽带和移动接入需要新增长极



千兆宽带和5G用户感知不显性，需要加大内容填充

填充什么？



携手生态伙伴，研发丰富的边缘特性应用，涵盖工作、生活、娱乐、安全等各类场景



基于虚拟局域网的云网服务，牵引家庭数字化新需求

大带宽、低时延、内网、安全、自助等差异化能力 提供差异化应用

场景	差异化能力	差异化应用
大带宽	上行 千兆 下行 千兆	家庭数据秒级同步 8K视频多屏实时播放
低时延	1毫秒	在家畅玩网吧专业级游戏
内网体验	家庭内网使用习惯不变	多点组网、异地投屏/打印
安全	数据不出家庭，基于内网的隐私保护	内网安全防护、绿色上网
自助服务	用户自助订购边缘特性应用	电商化体验

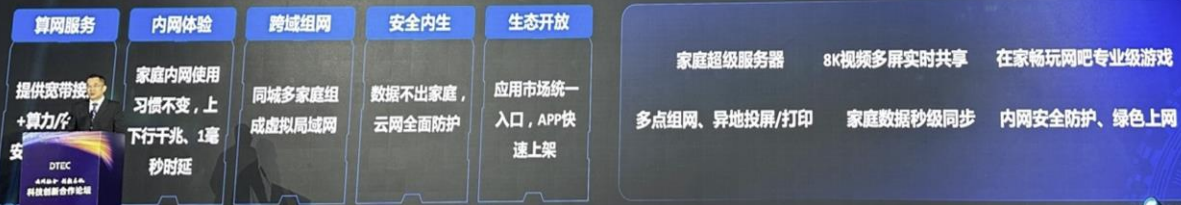
新型城域网实现同城算力高速互联
高效承接东西向流量

5G和千兆光网的云网融合组网

融合边缘，构建未来接入网的核心竞争力

五大特征

N大特色应用



上海贝尔：光、云、切片

网络提速：800G+ 的 IP + 光传输网络

提高速度 占地空间 消耗能源

- 200Gbps (节省的机架空间)
- 能效提升 (每Gb业务)

网络架构提升：多平面、网络协同

软件驱动、自动化

动态及可垂直业务连接

IP网络 光传输网络 敏捷及规模化传输能力

互连及可操作 关联及操作 优化及增强

- IP多平面组网，网络按需安全扩展
- IP/光传输协同提升网络可靠性及运维便利性

多云接入：云时代数据中心业务

- 多云接入网络架构
- 云原生时代的数据中心网络

精品业务提供：网络自动化及切片

- 网络自动化
- 网络切片

芯片设计

Quillion (固定接入)
• 功耗降低 50%

FPS (IP路由)
• 7nm, 功耗降低 75%

FPcx (IP路由)
• 7nm, 功耗降低 75%
• 100% 可编程
• 8 独立处理集群

PSE-6s (光传输)
5nm, 功耗降低 60%
光电垂直整合

网络层

216-core 大容量路由平台

正交背板设计

对称散热风道设计

光电解耦DCI平台

IP多平面设计

IP+光协同

软件

Altiplano 接入控制平台
现网控制

Corteca 云及家庭网络软件

Deepfield Secure Genome™

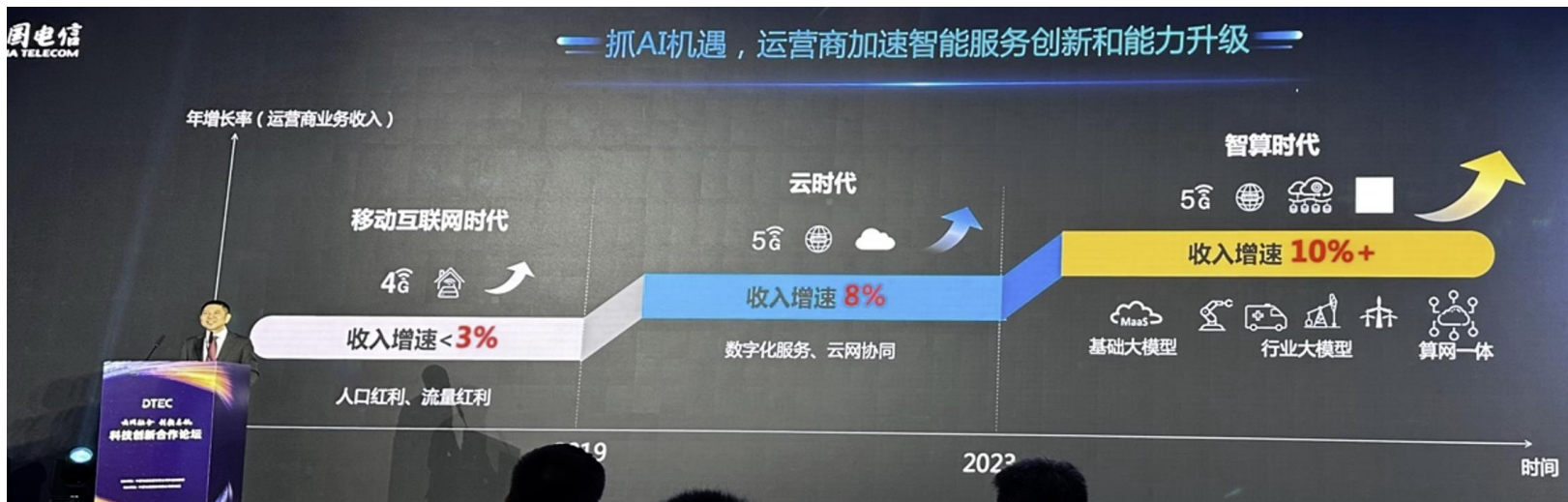
ANYsec Encryption

Quantum-safe Optical encryption

F-Secure

Network AUTOMATION

华为：算网融合

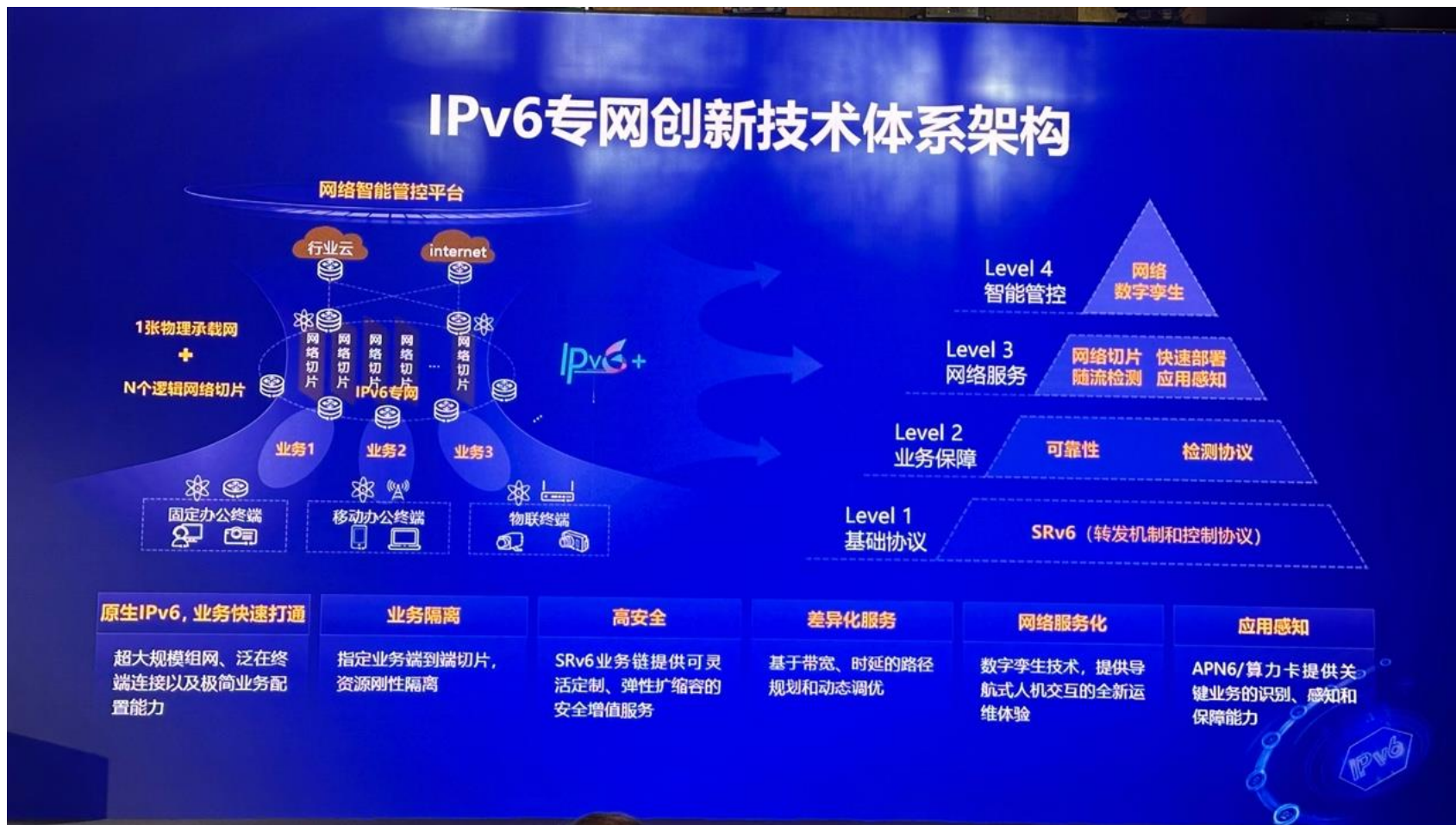


算内网络：“零丢包”超融合以太网，AI训练效率提升20%

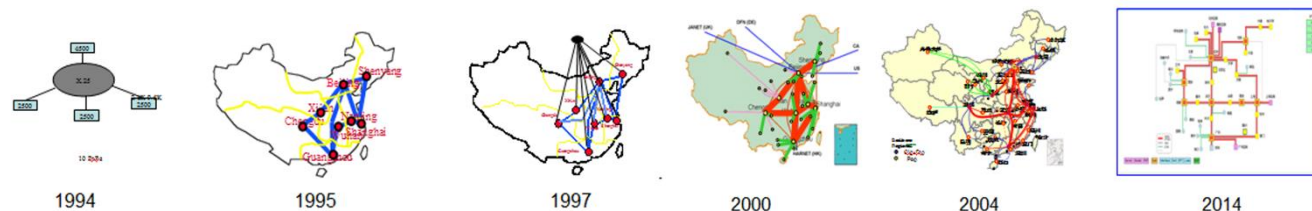
算间网络：打造算力时代“运力高铁”，网络吞吐率提升50%

入算网络：提供泛在万兆，按需服务的普惠算力接入

华为：IPv6专网



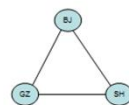
CERNET 基础设施



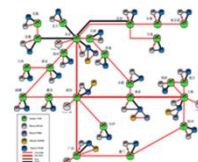
2.4 Kbps - 9.6 Kbps 64 Kbps 512 Kbps - 4 Mbps 155 Mbps 2.5 Gbps - 10 Gbps 100 Gbps
 IP over X.25 IP over lease line IP over VSAT IP over SDH IP over DWDM IP over DWDM



1997



2003



2006



2020

2.4 Kbps - 9.6 Kbps 2.5 Gbps 2.5 Gbps - 10 Gbps
 IPv6 over IPv4 IPv6 over SDH IPv6 over DWDM

100 Gbps
 IPv6 over DWDM



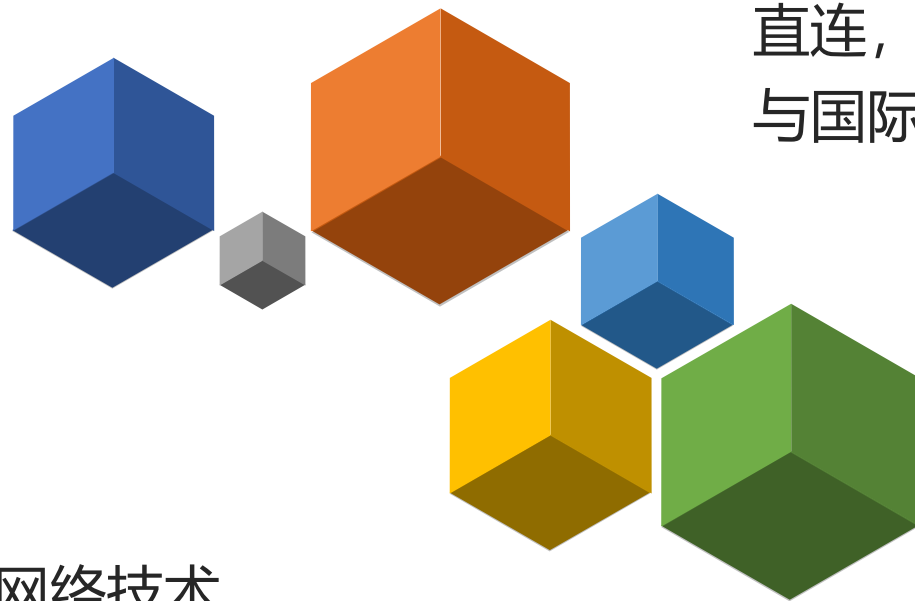
100 Gbps
 IPv6 slicing over DWDM



IPv6超算专网

依托教育网构建超算互联网的优势

- ✓ 出色的数据传输性能，可为超算互联网体系“动脉网络系统”的畅通和安全提供保障



- ✓ 先进的通信和网络技术手段，可为超算中心之间的业务协作提供支持

- ✓ 与美国、欧洲等国家和地区教育网直连，将保障中国超算互联网体系与国际超算中心的高速大数据交互

- ✓ 教育网基本全面覆盖本科类高校，即高性能计算的主要用户群体

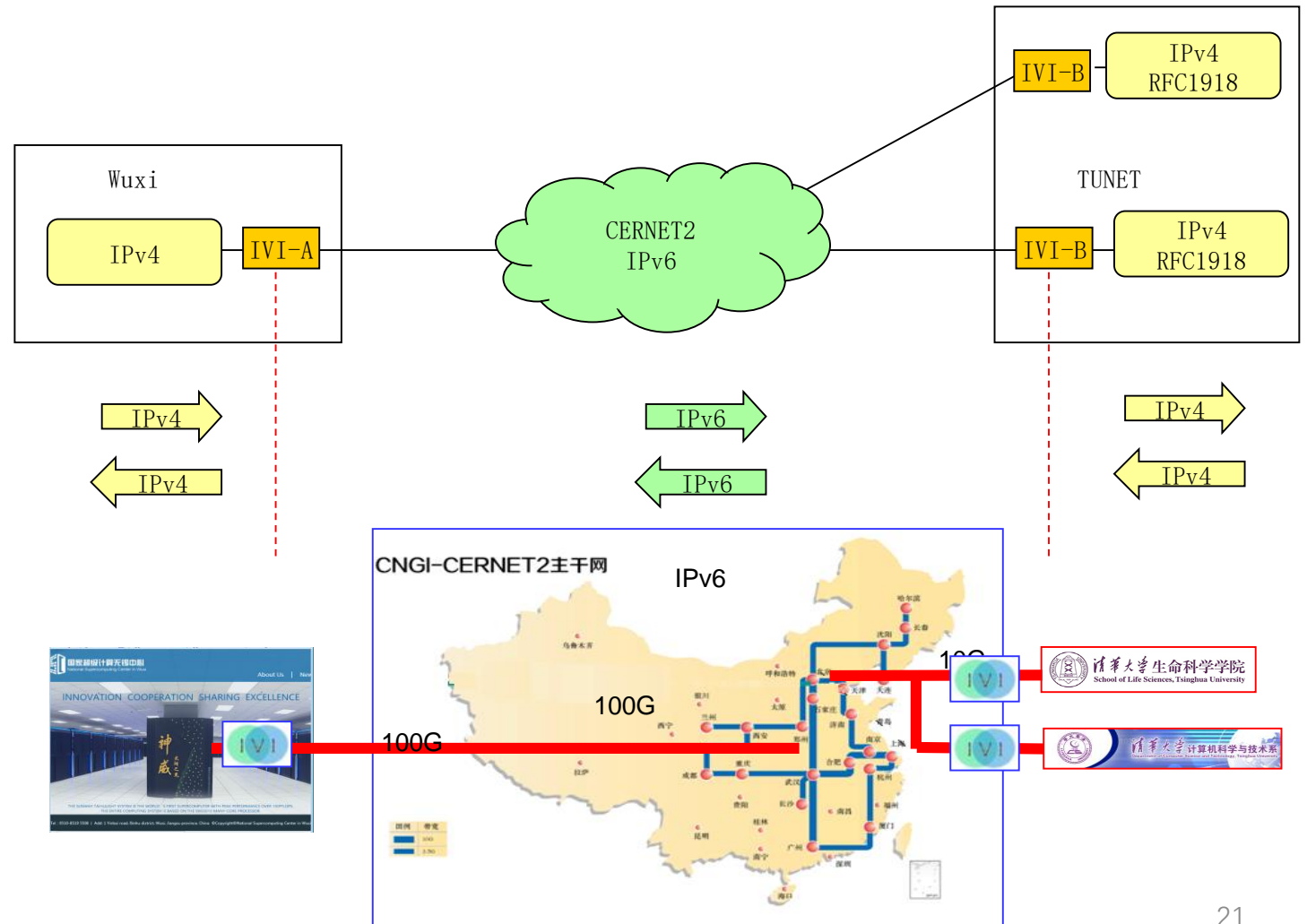
清华大学-无锡超算

基础设施

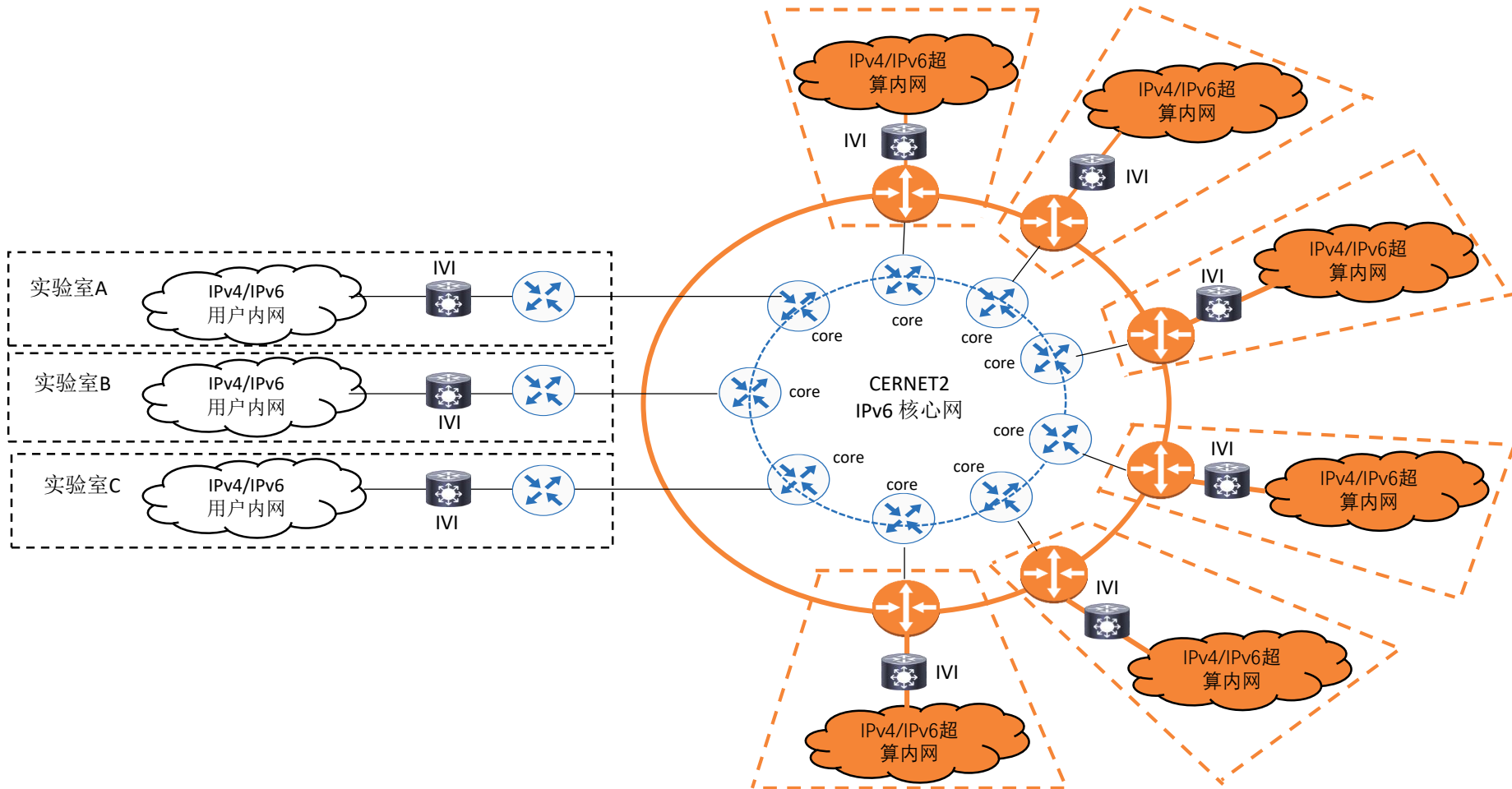
- 无锡超算中心网络: IPv4/IPv6双栈
- CERNET2: 100G IPv6
- 校园网络: IPv4/IPv6双栈

基础模块

- 无状态IPv4/IPv6翻译技术IVI



技术方案



特点

- 高性能端到端传输

- 使用IPv6和IVI技术，10G链路的单流性能可达9Gbps
- 灵活调度：使用IPv6前缀动态调度，比IPv4 BGP调度更加灵活
- 比MPLS的成本显著降低

- 兼容IPv4/IPv6应用

- 使用IVI双重翻译技术，使得纯IPv4应用无需任何修改，就可通过IPv6骨干网络访问超算服务
- 使用SRv6技术，可选择特定路径

- 流量管理

- 由于IVI的静态映射特性，可在IPv6网络中针对特定超算互联网流量进行计费和管理

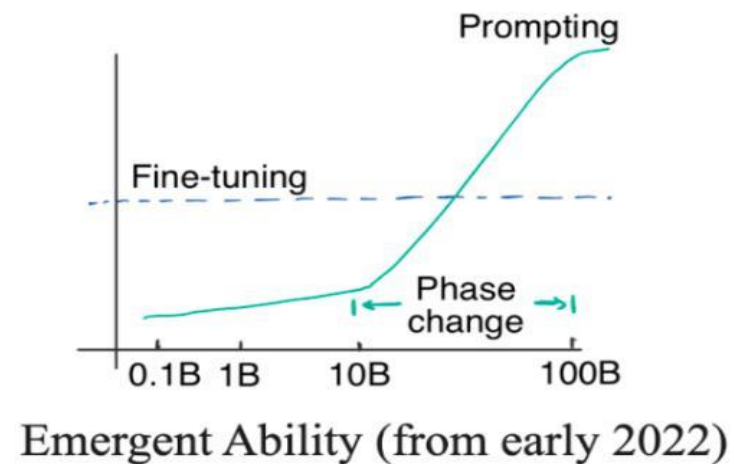
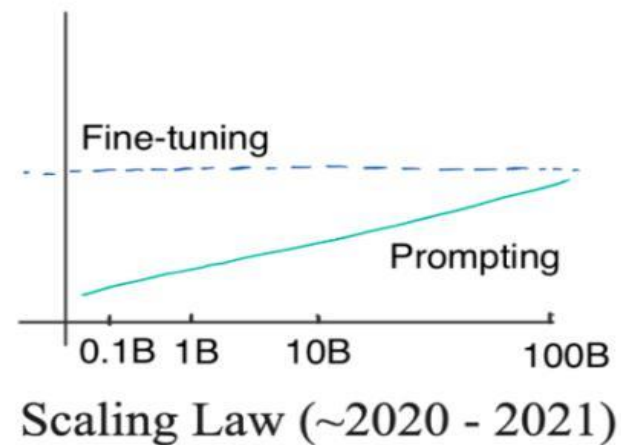
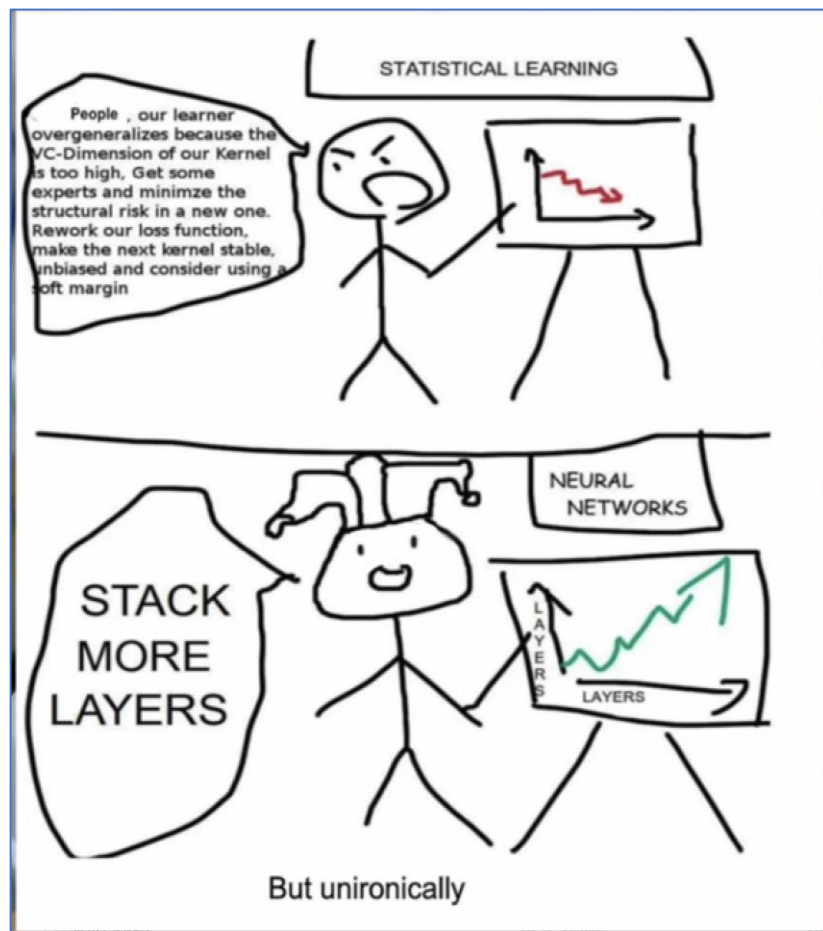
- 适应实际部署

- 适应多出口切片部署场景
- 适应于用户园区网络（校园网/企业网）基于IPv6地址或IPv6前缀的不同认证机制



IPv6 人工智能专网

顿悟



差距



- 实体清单
 - 华为、海康、科大讯飞、大华、旷世、商汤、依图等
- 芯片占比
 - 2021人工智能芯片出货量100万片，国产芯片出货量5万片
 - 英伟达占95%左右
- 框架占比
 - Meta人工智能算法开发框架占中国90%以上市场份额
- 人才培养
 - 中国99.5%的在校理工科大学生学习国外的人工智能技术

“百模大战” 盘点国内外横空出世的AI大模型

- OpenAI: GPT 系列大模型
- 微软: Orca大模型等
- 谷歌: PaLM 2 大模型、Gemini大模型等
- Meta: LLaMA语言模型、ImageBind 大模型等
- AWS: Titan语言大模型

“有报道说截至10月份，中国已经发布了238个大模型，而6月份的时候这个数字大概是79个，相当于4个月翻了3倍。中国有多少AI原生应用，我想在座的各位很少有人能够说出一两个来，而国外除了有几十个基础大模型发布外，已经有上千个AI的原生应用，这是在中国市场上没有的。”李彦宏说。

- 清华大学——ChatGLM-6B
- 清华大学——OpenBMB (CPM-Bee-10B)
- 北京大学——ChatLaw
- 哈尔滨工业大学——本草
- 上海交通大学——K2
- 东北大学——TechGPT
- 百度: 文心大模型
- 腾讯: 混元大模型
- 阿里: 通义大模型
- 华为: 盘古大模型
- 网易: 玉言、子曰大模型
- 京东: 言犀大模型
- 360: 360智脑大模型
- 浪潮: 源大模型
- 科大讯飞: 星火认知大模型
- 商汤: 日日新大模型
- 智谱AI: 智谱AI系列大模型
- 昆仑万维: 天工大模型
- 中国移动: “九天”1+N大模型
- 中国电信: TeleChat大模型
- 中国联通: 鸿湖图文大模型1.0
- 中国科学院自动化研究所: 紫东太初大模型
- 智源研究院: 悟道智能模型

轮回：1993-2023

• Internet

➤ 应用:

- 中国: jpeg是硕士生的论文
- 世界: jpeg 的源程序免费下载

➤ 设备:

- 前期: 受“巴统”限制
- 后期: 国产和进口混合组网

➤ 治理:

- 管理: CERNET管理规定+BBS
- 技术: 195号令+GFW

• chatGPT

➤ 应用:

- 中国: 依然是传统编程方式
- 世界: chatGPT 自动生成程序

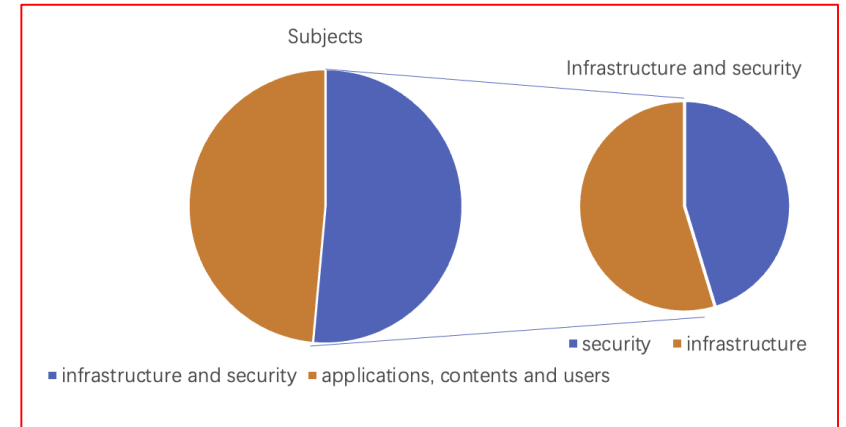
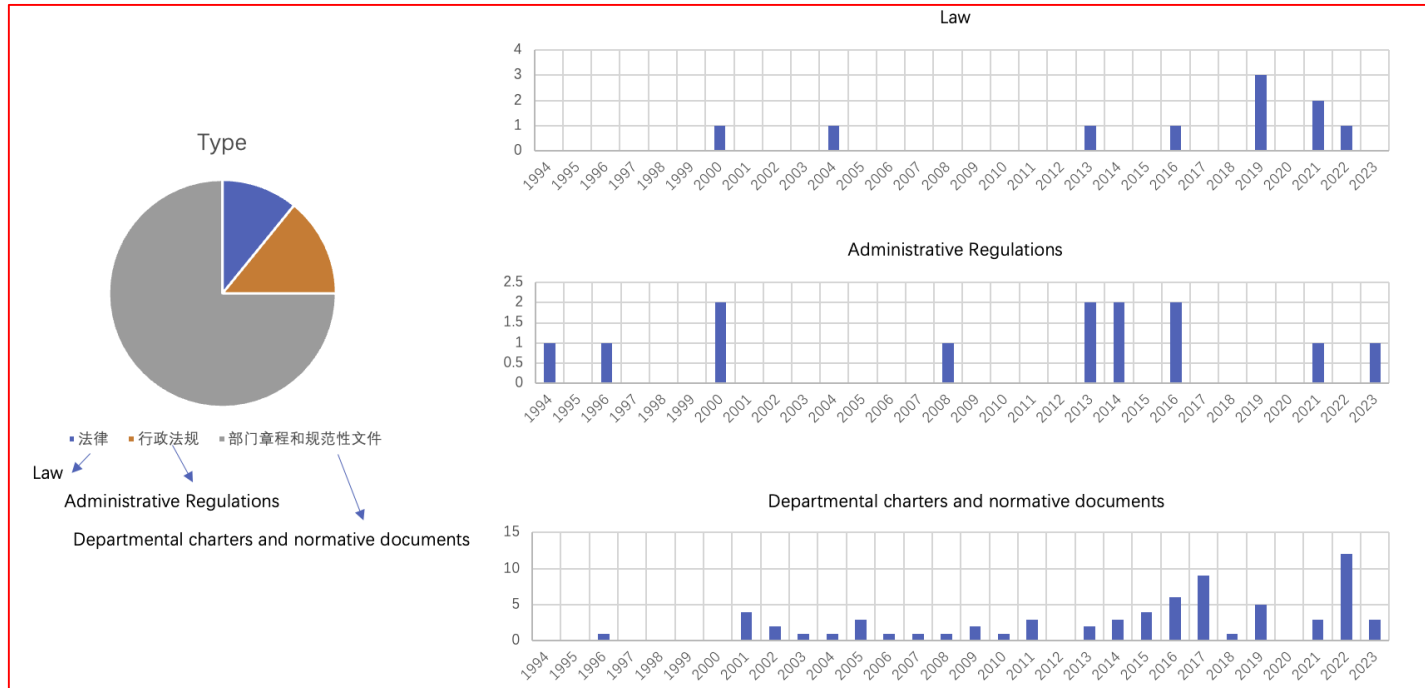
➤ 设备:

- 前期: 可以进口NVIDIA A100
- 后期: 无法进口 NVIDIA A100/H100

➤ 治理:

- 管理: ?
- 技术: ?

法律法规



技术、商业、法律法规

Blog (business first)



2006.03



2007.05.12



2009.08



2010.05.01



2018.02.02

9 years

chatGPT (regulation first)



2022.11.30



2023.05.23

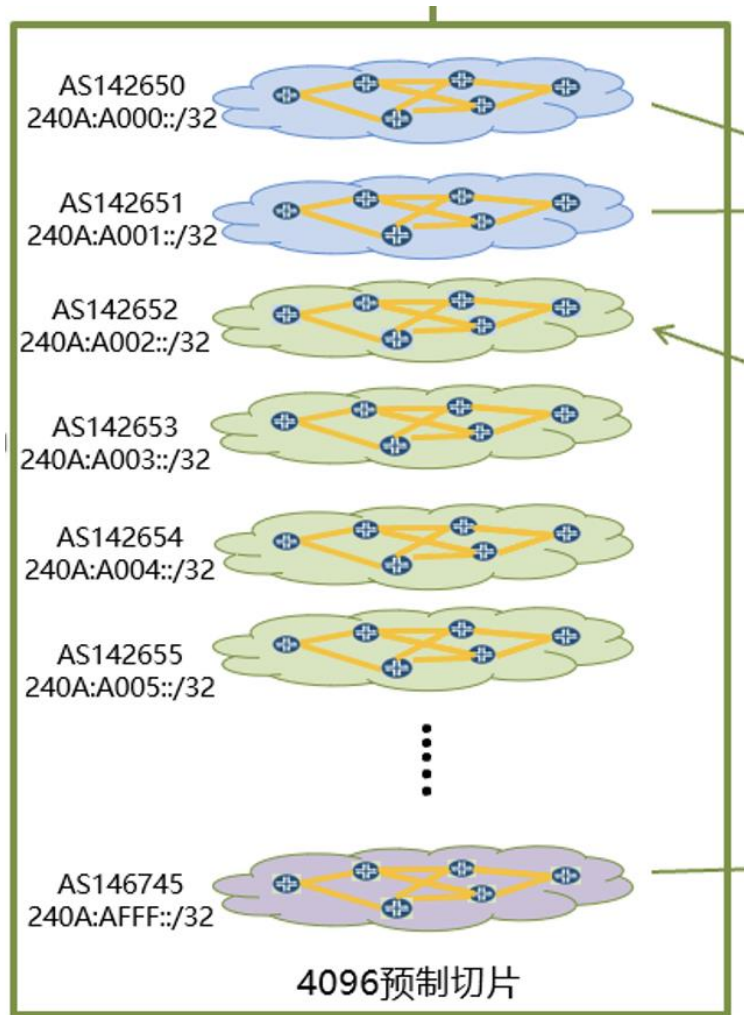


2023.08.31

3 months

From "business first" to "regulation first".

AI 专网



- Share AI related video
- Achieve chatGPT dialogs
- Distributed computing
- And much more

IPv6 AI 专网

- ✓ 要素
 - 地址 (ULA vs GUA)
 - 域名 (IPv6 address vs internal DNS)
 - 边界 (No NAT66 and zero trust)
 - 过渡 (IVI for incremental and transparent deployment)

Github 和 Huggingface

Repository search results

github.com/search?q=llama2&type=repositories

Filter by

- Code 58.4k
- Repositories 1.6k
- Issues 4k
- Pull requests 3k
- Discussions 262
- Users 58

Languages

- Python
- Jupyter Notebook
- TypeScript
- JavaScript
- C
- Rust
- C++
- HTML

1.6k results (265 ms) Sort by: Best match

- karpathy/llama2.c** Inference Llama 2 in one file of pure C
C · 12.9k · Updated 5 days ago
- FlagAlpha/Llama2-Chinese** Llama中文社区, 最好的中文Llama大模型, 完全开源可商用
llama lora finetune pretrain llm
Python · 6.9k · Updated 18 days ago
- dataprofessor/llama2** This chatbot app is built using the Llama 2 open source LLM from Meta.
python meta streamlit large-language-models llm
Python · 173 · Updated 24 days ago

Sponsor open source projects you depend on

Contributors are working behind the scenes to make open source better for everyone—give them the help and recognition they deserve.

Explore sponsorable projects →

How can we improve search? Give feedback

ProTip! Press the **/** key to activate the search input again and adjust your query.

Models - Hugging Face

huggingface.co/models

Hugging Face Search models, datasets, spaces, docs, solutions, pricing

Tasks Libraries Datasets Languages Licenses

Other

Filter Tasks by name

Multimodal

- Feature Extraction
- Text-to-Image
- Image-to-Text
- Text-to-Video
- Visual Question Answering
- Document Question Answering
- Graph Machine Learning

Computer Vision

- Depth Estimation
- Image Classification
- Object Detection
- Image Segmentation
- Image-to-Image
- Unconditional Image Generation
- Video Classification
- Zero-Shot Image Classification

Models 401,147 Filter by name new Full-text search Sort: Trending

- openai/whisper-large-v3** Automatic Speech Recognition · Updated 3 days ago · 61.7k · 641
- latent-consistency/lcm-lora-sdx1** Text-to-Image · Updated 3 days ago · 29.3k · 273
- 01-ai/Yi-34B** Text Generation · Updated 4 days ago · 37.4k · 897
- openchat/openchat_3.5** Text Generation · Updated 2 days ago · 17k · 576
- distil-whisper/distil-large-v2** Automatic Speech Recognition · Updated 5 days ago · 48.2k · 321
- NousResearch/Nous-Capybara-34B** Text Generation · Updated 5 days ago · 367 · 91

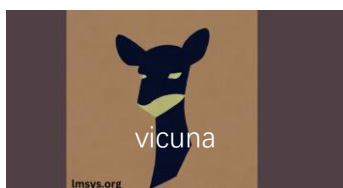
若干开放源码及模型



美洲驼



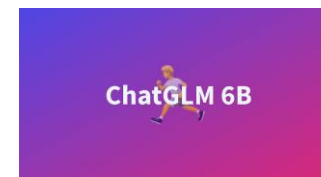
羊驼



小羊驼



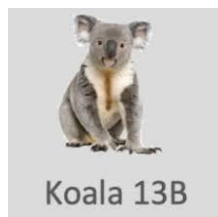
多莉羊



ChatGLM 6B



猎鹰



考拉



虎鲸



通义千问

Qwen-7B-Chat Bot

通义千问-7B (Qwen-7B) 是阿里云研发的通义千问大模型系列的70亿参数规模的模型。



西北风



西风



GPT4ALL

GPT4All
A free-to-use, locally running, privacy-aware chatbot. **No GPU or internet required.**

Download Desktop Chat Client

Windows Installer | OSX Installer | Ubuntu Installer

<https://gpt4all.io/index.html>

Installation Instructions

Windows | MacOS | Ubuntu

Windows Installer
After download and installation you should be able to find the application in the directory you specified in the installer. You will find a desktop icon for GPT4All after installation.
NOTE: On Windows, the installer might show a security complaint. This is being addressed as we're actively setting up cert sign for Windows.

GPT4All Ecosystem
Explore the GPT4All open source ecosystem

- GPT4All Training**: Train your own GPT4All models.
- GPT4All Documentation**: Integrate a locally running LLM into any codebase.
- GPT4All Chat**: A multi platform chat interface for running local LLMs.
- GPT4All in Python**: Python bindings to GPT4All.
- GPT4All Datalake**: An open source database for donated GPT4All interaction data.

Performance Benchmarks

Model	BoolQ	PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg
GPT4All-J 6B v1.0	73.4	74.8	63.4	64.7	54.9	36	40.2	58.2
GPT4All-J v1.1-breezy	74	75.1	63.2	63.6	55.4	34.9	38.4	57.8
GPT4All-J v1.2-jazzy	74.8	74.9	63.6	63.8	56.6	35.3	41	58.6
GPT4All-J v1.3-groovy	73.6	74.3	63.8	63.5	57.7	35	38.8	58.1
GPT4All-J Lora 6B	68.6	75.8	66.2	63.5	56.4	35.7	40.2	58.1
GPT4All LLaMa Lora 7B	73.1	77.6	72.1	67.8	51.1	40.4	40.2	60.3
GPT4All 13B snoozy	83.3	79.2	75	71.3	60.9	44.2	43.4	65.3
GPT4All Falcon	77.6	79.8	74.9	70.1	67.9	43.4	42.6	65.2
Nous-Hermes	79.5	78.9	80	71.9	74.2	50.9	46.4	68.8
Nous-Hermes2	83.9	80.7	80.1	71.3	75.7	52.1	46.2	70.0
Nous-Puffin	81.5	80.7	80.4	72.5	77.6	50.7	45.6	69.9
Dolly 6B	68.8	77.3	67.6	63.9	62.9	38.7	41.2	60.1
Dolly 12B	56.7	75.4	71	62.2	64.6	38.5	40.4	58.4
Alpaca 7B	73.9	77.2	73.9	66.1	59.8	43.3	43.4	62.5
Alpaca Lora 7B	74.3	79.3	74	68.8	56.6	43.9	42.6	62.8
GPT-J 6.7B	65.4	76.2	66.2	64.1	62.2	36.6	38.2	58.4
Llama 7B	73.1	77.4	73	66.9	52.5	41.4	42.4	61.0
Llama 13B	68.5	79.1	76.2	70.1	60	44.6	42.2	63.0
Pythia 6.7B	63.5	76.3	64	61.1	61.3	35.2	37.2	56.9
Pythia 12B	67.7	76.6	67.3	63.8	63.9	34.8	38	58.9
Fastchat T5	81.5	64.6	46.3	61.8	49.3	33.3	39.4	53.7
Fastchat Vicuña 7B	76.6	77.2	70.7	67.3	53.5	41.2	40.8	61.0
Fastchat Vicuña 13B	81.5	76.8	73.3	66.7	57.4	42.7	43.6	63.1
StableVicuña RLHF	82.3	78.6	74.1	70.9	61	43.5	44.4	65.0
StableLM Tuned	62.5	71.2	53.6	54.8	52.4	31.1	33.4	51.3
StableLM Base	60.1	67.4	41.2	50.1	44.9	27	32	46.1
Koala 13B	76.5	77.9	72.6	68.8	54.3	41	42.8	62.0
Open Assistant Pythia 12B	67.9	78	68.1	65	64.2	40.4	43.2	61.0
Mosaic MPT7B	74.8	79.3	76.3	68.6	70	42.2	42.6	64.8
Mosaic mpt-instruct	74.3	80.4	77.2	67.8	72.2	44.6	43	65.6
Mosaic mpt-chat	77.1	78.2	74.5	67.5	69.4	43.3	44.2	64.9
Wizard 7B	78.4	77.2	69.9	66.5	56.8	40.5	42.6	61.7
Wizard 7B Uncensored	77.7	74.2	68	65.2	53.5	38.7	41.6	59.8
Wizard 13B Uncensored	78.4	75.5	72.1	69.5	57.5	40.4	44	62.5
GPT4-x-Vicuña-13b	81.3	75	75.2	65	58.7	43.9	43.6	63.2
Falcon 7b	73.6	80.7	76.3	67.3	71	43.3	44.4	65.2
Falcon 7b instruct	70.9	78.6	69.8	66.7	67.9	42.7	41.2	62.5
text-davinci-003	88.1	83.8	83.4	75.8	83.9	63.9	51	75.7

允许ChatGPT犯错误 开放治理同行

- 治理包含两个方面（模型本身的治理，人类自己的治理）
- 模型本身的治理：
 - 多模态：（文本、语音、图像、视频）符合道德标准、伦理标准和政治标准
 - 创造性：要有创造性，就必须允许犯错误
 - 开放跟治理必须同时进行，必须边开放边治理
 - 靠人类来做对齐处理

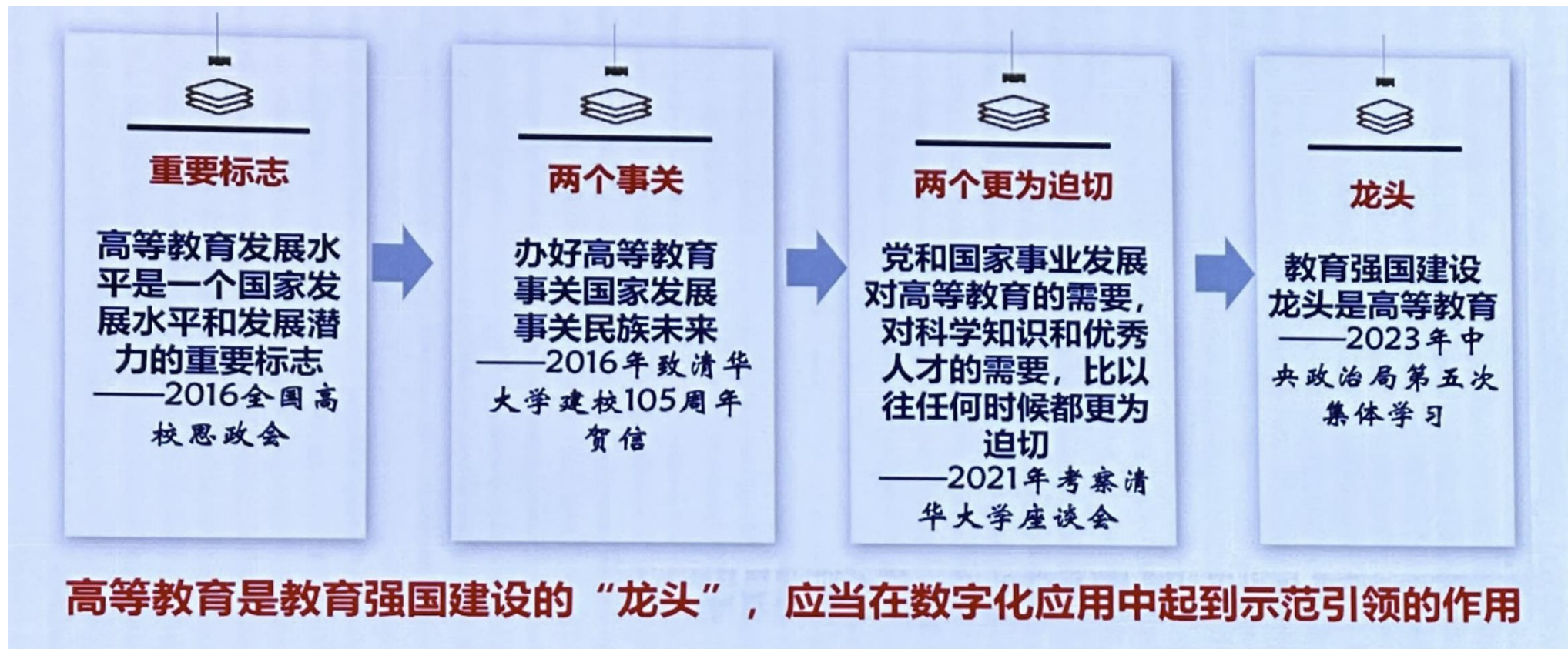


小结

IPv6 专网

- HPC 高性能
- AI 安全可信

教育强国



教育强国

2.1 起步蓄能期



1978年
全国教育工作会议

小平同志提出，要制订加速发展电视、广播等现代化教育手段的措施，这是多快好省发展教育事业的重要途径，必须引起充分的重视。



1978年
中央广播电视大学成立

国家批准，教育部成立中央广播电视大学，也就是国家开放大学的前身，第一家面向全国开展远程开放教育的新型高等学校应运而生。



1984年
计算机教育普及

小平同志在参观上海微电子技术应用汇报展览会上指出，计算机的普及要从娃娃做起，推动计算机进入教育教学领域。



1994年
互联网进入校园

清华大学等10所高校承担“中国教育和科研计算机网示范工程”建设项目，建成了我国第一个大型骨干网络，高等学校成为我国第一批互联网用户。

中国教科网诞生了第一批互联网应用，高等学校成为我国互联网发展的先锋力量。



第一个
电子杂志

神州学人
(1996年)



第一个
搜索引擎

天网搜索引擎
(1996年)



第一个
数字图书馆

CALIS (1998年)
CADAL (2001年)



第一个
网络论坛

水木清华
(2004年)

教育强国

从教育大国到教育强国是一个系统性跃升和质变，必须以改革创新为动力

——习近平在中央政治局第五次集体学习时的讲话（2023年5月29日）

要坚持系统观念，统筹推进

- 育人方式改革
- 办学模式改革
- 管理体制改革
- 保障机制改革
- 教育评价改革

要充分发挥数字技术“**突破口**”的作用，促进数字技术与教育教学的深度融合，推动高等教育的高质量发展，实现高教发展的“**变轨超车**”，回答好“**强国建设、教育何为**”的时代之问。



互联网



区块链



大数据



人工智能

重新思考教育科研计算机网的定位

- 1994-2004
 - 引领公众互联网
- 2004-2024
 - 与公众互联网既合作又竞争（）
- 2024- ……
 - 既是公众互联网，又是高层次专网（）
 - 超算专网
 - AI专网
 - ……

IPv6 为我们提供了无限的技术可能性！！！！

进化论和技术演进



USE FIRE

INTERNET

CHATGPT

IPV6+AI?



谢谢!