

AI智能计算：降本提效 为科研加速

北京超级云计算中心

汇报人：郭跃 2023年4月

CONTENTS

- 1 中心介绍
- 2 AI智算云与产品服务
- 3 客户案例

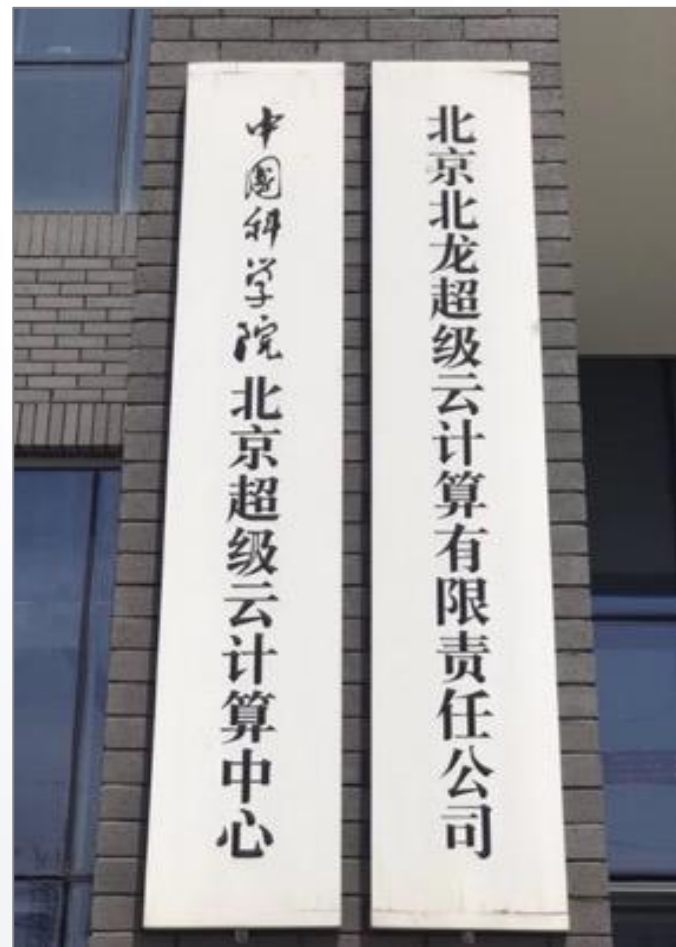
中心介绍

| 行业概述 | 成立背景 | 荣誉资质 |
| 算力分布 | 发展历程 | 行业覆盖 |

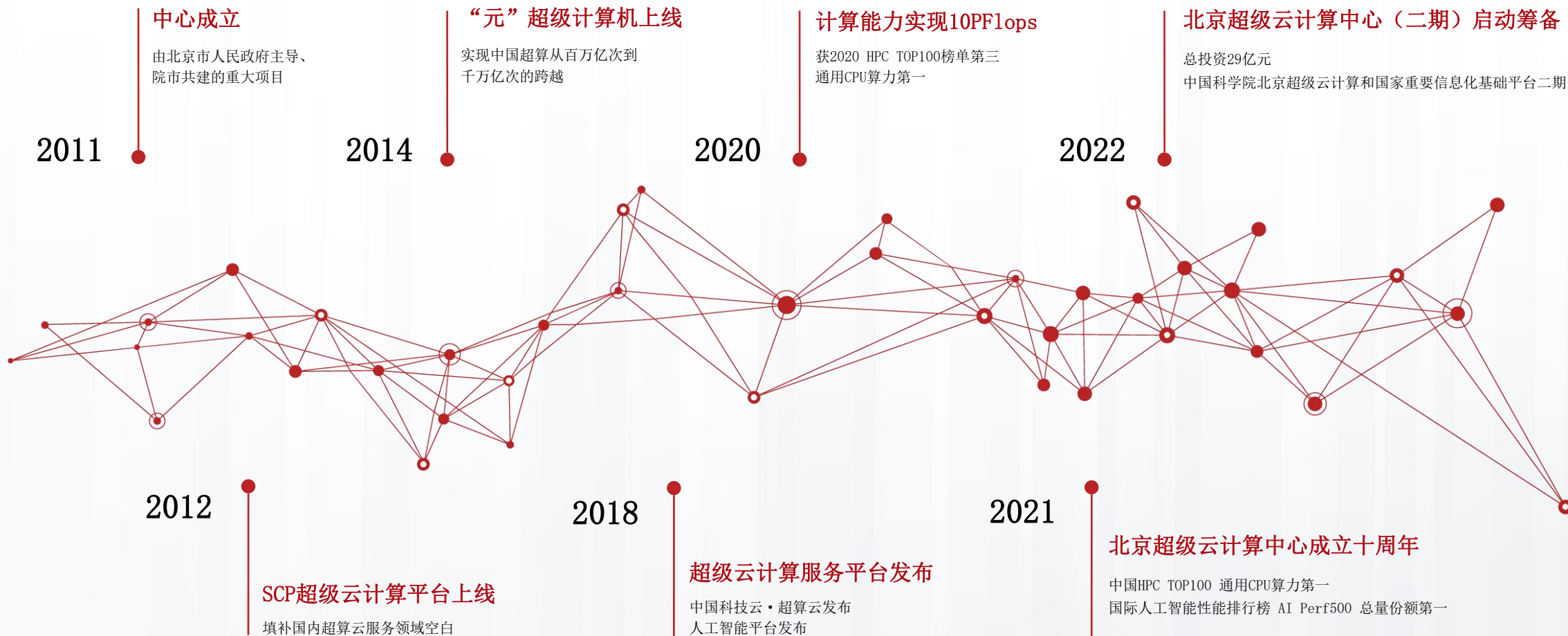


中心介绍

北京超级云计算中心（简称“北京超算”），2011年11月在北京怀柔综合性国家科学中心—怀柔科学城奠基成立，是由北京市人民政府主导、院市共建的“北京超级云计算和国家重要信息化基础平台”。



发展历程



云聚十载 智算未来



云聚10载·智算未来
2021超级云计算高峰论坛
2021 Super Cloud Computing Summit

指导单位:
中国科学院计算机网络信息中心
北京市怀柔区政府

主办单位:
北京市怀柔区经济与信息化局
北京中科北龙科技有限责任公司
北京市长城伟业投资开发有限公司
CCF中国计算机学会高性能计算专业委员会
ACM中国高性能计算专家委员会

承办单位:
北京超级云计算中心
北京北龙超级云计算有限责任公司
中科捷云(北京)信息技术有限公司

支持单位:
AMD中国

北京超级云计算中心
Beihai Super Cloud Computing Center



十年磨一剑，引领超算云服务化
商业化助推科技创新，促进数字经济发展

关于我们

中国超算算力调度平台

- 基于国家重点项目，及全国超算用户使用打磨发展成熟
- “十一五、十二五”重点研发计划“中国超算网格环境”
- “十三五”重点专项“中国超算环境”

算力布局

- 已在北京、宁夏、内蒙古等地前瞻性布局了三个主算力枢纽
- 构建跨域资源协同调度体系，优化算力之间的统筹联动
- 总投资29亿元的北京超算(二期)已于2022年启动筹备



中国科学院北京超级云计算和国家重要信息化基础平台二期建设功能规划图

100万+CPU核心

100亿核时

20+ 行业

100PFlops+

数千卡GPU

1000+ 企业

1万台+超算服务器

200款+
行业软件SaaS化

1万+ 服务群

20万+用户

算力资源类型

通用 X86 CPU算力

AMD EPYC
Intel Xeon Scalable

8种不同类型CPU型号
让计算更简单、高效

50PFlops+

AI算力

Nvidia GPU

基于A100、A800、V100等
主流GPU型号

100PFlops (单精度)+

国产X86算力

国产 X86 CPU
国产加速卡芯片

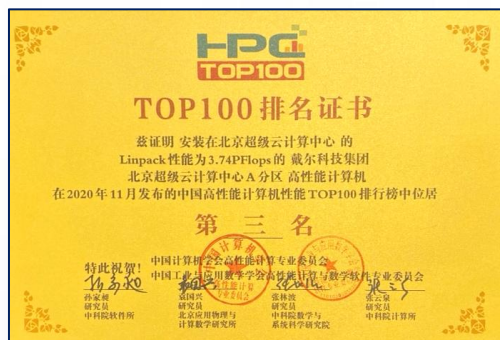
CPU+国产加速卡
异构环境

100PFlops+

中心荣誉

2020年

中国HPC TOP100排行榜



中国高性能计算机性能TOP100
排行榜第三名

通用CPU算力第一名

2021年

中国HPC TOP100排行榜



中国高性能计算机性能TOP100
同构众核CPU性能

第一名

2022年

中国HPC TOP100排行榜

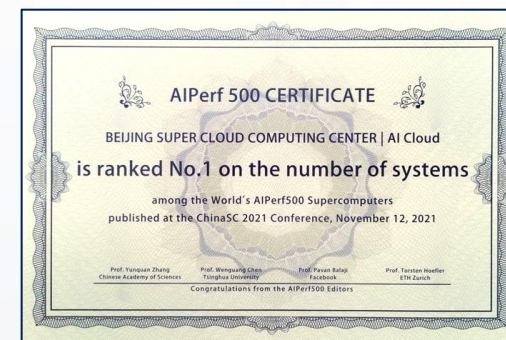


中国高性能计算机性能TOP100
同构众核CPU性能

第一名

2021年

国际人工智能性能排行榜 AIPerf 500



国际人工智能性能排行榜AIPerf 500
总量份额

第一名

公司实力

北京北龙超级云计算有限责任公司

国家高新技术企业

中关村高新技术企业

ISO9001质量体系认证

增值电信业务经营许可证

AAA资信等级证书

北京市诚信创建企业

北京软件和信息服务业协会会员

软著：38项

专利：2项

.....



算力资源类型

通用 X86 CPU算力

AMD EPYC
Intel Xeon Scalable

8种不同类型CPU型号
让计算更简单、高效

50PFlops+

AI算力

Nvidia GPU

基于A100、V100、T4等
主流GPU型号

100PFlops (单精度)+

国产X86算力

国产 X86 CPU
国产加速卡芯片

CPU+国产加速卡
异构环境

100PFlops+

遍布全国的算力网络



资源规模、体量、分布

率先响应国家政策，算力资源云端共享



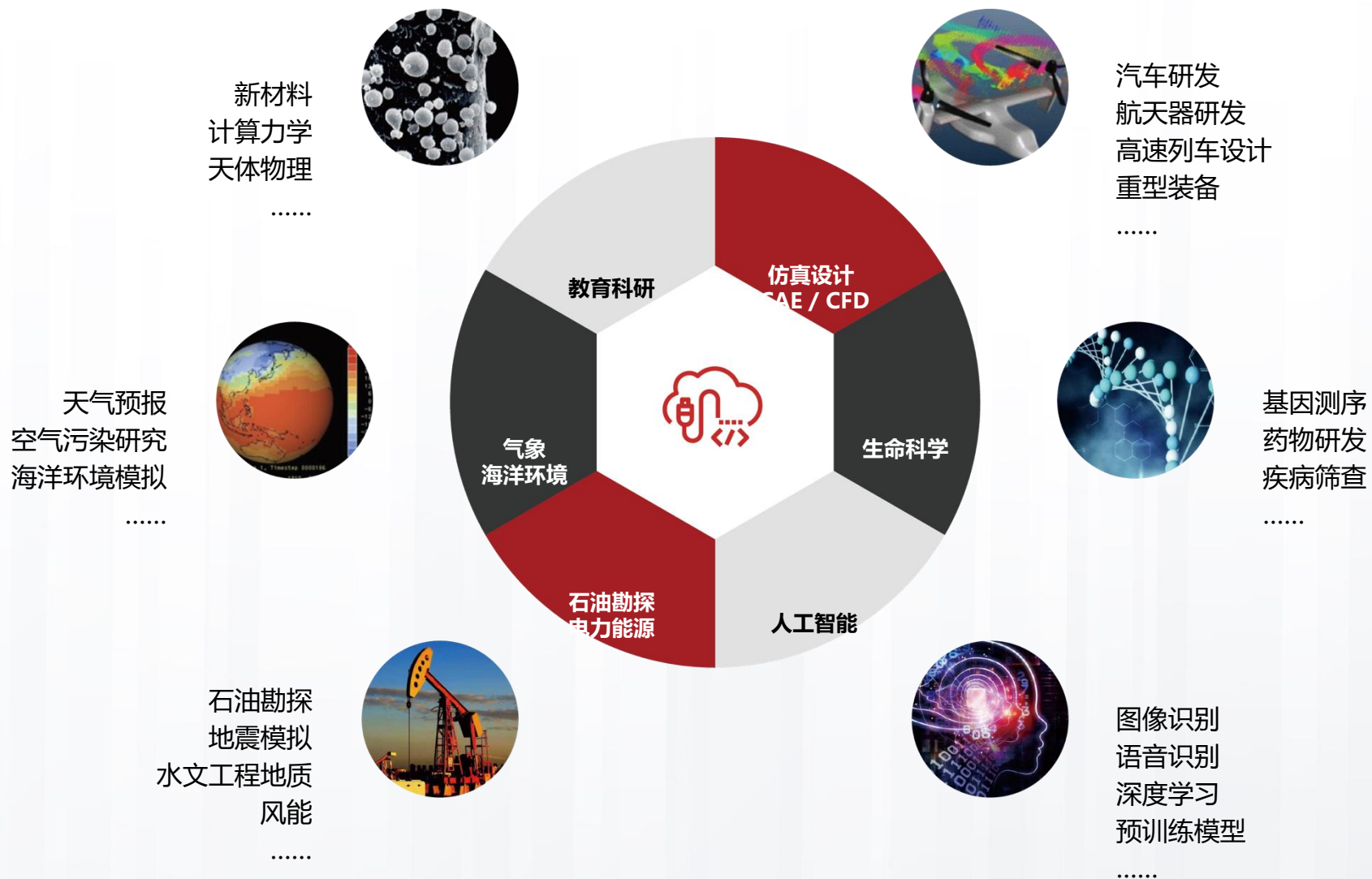
资源配置

超算集群资源	计算资源配置			
北京超级云计算中心T6分区 2021年中国HPC性能TOP 7 / 2022年中国HPC性能TOP8	普通计算节点	Intel Platinum CPU	96核心	384GB 内存
北京超级云计算中心A6分区 2021年中国HPC性能TOP 10 / 2022年中国HPC性能TOP 11	普通计算节点	AMD EPYC CPU	64核心	256GB 内存
北京超级云计算中心A分区 2020年中国HPC性能TOP 3 / 2021年中国HPC性能TOP 11 / 2022年中国HPC性能TOP 12	普通计算节点	AMD EPYC CPU	64核心	256GB 内存
北京超级云计算中心M分区	普通计算节点	AMD EPYC CPU	128核心	512GB 内存
中国科技云9区	普通计算节点	Intel Xeon v3 CPU	24核心	128GB 内存
中国科技云13区	普通计算节点	Intel Xeon GOLD CPU	32核心	96GB 内存
中国科技云19区	普通计算节点	Intel Xeon v4 CPU	36核心	256GB 内存
“元”二期	普通计算节点	Intel Xeon v3 CPU	24核心	256GB 内存

资源分区	计算资源配置			
北京超级云计算中心N23分区	GPU云主机	AMD EPYC Rome CPU + NVIDIA Tesla A100 Intel Platinum CPU + NVIDIA Tesla V100 Intel Platinum CPU + NVIDIA Tesla T4	164 vCPU + 8 GPU 88 vCPU + 8 GPU 88 vCPU + 8 GPU	948GB 内存 352GB 内存 352GB 内存
北京超级云计算中心N26分区	GPU裸金属服务器	Intel Platinum CPU + NVIDIA Tesla V100	80 vCPU + 8 GPU	320GB 内存

超算集群资源	计算资源配置			
国产超算资源	普通计算节点	Hygon CPU DCU加速卡	32核心 4DCU	128GB 内存

行业覆盖



科研机构用户 (部分)



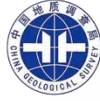
北京动力机械研究所



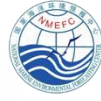
北京纳米能源与系统研究所



北京市农林科学院



广州海洋地质调查局



国家海洋环境预报中心



国家气象预报中心



化学与精细化工广东省实验室



青岛海洋科学与技术试点国家实验室



中国地质科学院



中国环境科学研究院



中国科学院高能物理研究所



中国科学院国家纳米科学中心



中国科学院过程工程研究所



中国科学院国家天文台



中国科学院国家授时中心



中国科学院海洋研究所



中国科学院精密测量科学与技术创新研究院



中国科学院计算技术研究所



中国科学院金属研究所



中国科学院昆明植物研究所



中国科学院兰州化学物理研究所



中国科学院理论物理研究所



中国科学院沈阳自动化研究所



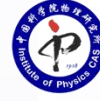
中国科学院生态环境研究中心



中国科学院苏州纳米技术与纳米仿生研究所



中国科学院微电子研究所



中国科学院物理研究所



中国科学院亚热带农业生态研究所



中国科学院长春应用化学研究所



中国科学院重庆绿色智能技术研究院



中国科学院紫金山天文台



中国医学科学院药物研究所



中国航空空气动力技术研究院



中国科学院半导体研究所



中国科学院北京综合研究中心



中国科学院成都山地灾害与环境研究所



中国科学院大连化学物理研究所



中国科学院大气物理研究所



中国科学院地理科学与资源研究所



中国科学院地球化学研究所



中国科学院地质与地球物理研究所



中国科学院东北地理与农业生态研究所



中国科学院力学研究所



中国科学院南海海洋研究所



中国科学院宁波材料技术与工程研究所



中国科学院青藏高原研究所



中国科学院软件研究所



中国科学院山西煤炭化学研究所



中国科学院上海高等研究院



中国科学院上海光学精密机械研究所



中国科学院上海药物研究所



中国科学院上海应用物理研究所



中国科学院深圳先进技术研究院



中国科学院微生物研究所



中国原子能科学研究院



中国科学院植物研究所



中国气象科学研究院



中国水利水电科学研究院



自然资源部第一海洋研究所

.....

高校用户 (部分)



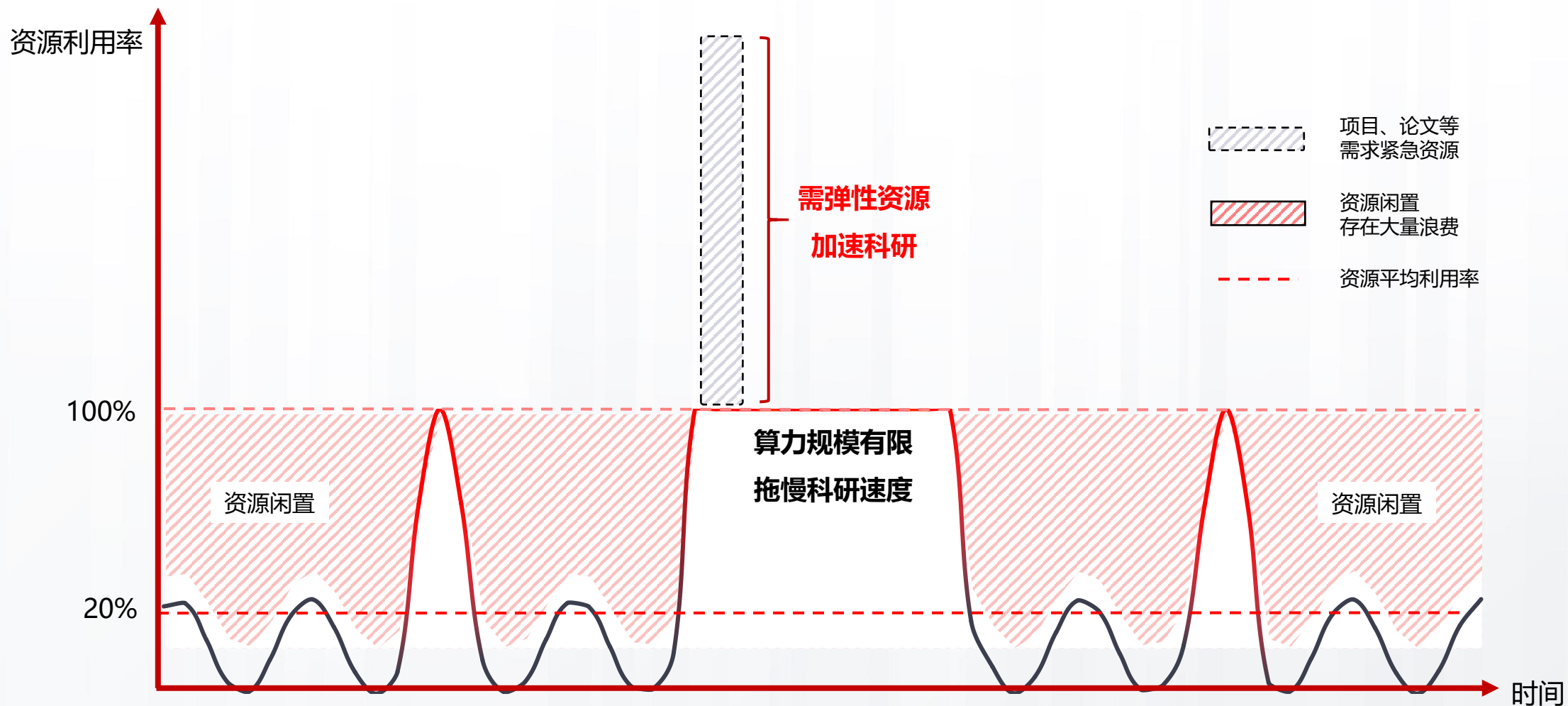
企业用户 (部分)

 爱驰汽车 (上海)有限公司	 澳汰尔工程软件 (上海)有限公司	 北京辰安科技 股份有限公司	 北京海基嘉盛科技 有限公司	 北京海兰信数据 科技股份有限公司	 航天长峰 股份有限公司	 北京航天智造 科技发展有限公司	 北京华云星地通 科技有限公司	 北京金风科创 风电设备有限公司	 北京京东科信息 技术有限公司
 长城汽车 股份有限公司	 超威半导体产品 (中国)有限公司	 重庆长安汽车 股份有限公司	 成都安世亚太 科技有限公司	 广东美的制冷设备 有限公司	 广州华粤普盈 科技有限公司	 广州汽车集团 股份有限公司	 广州市坤通信息 科技有限公司	 国家电投集团科学 技术研究院有限公司	 海丰通航 科技有限公司
 三一集团有限公司	 三一汽车制造 有限公司	 上海极链网络科技 有限公司	 上海拓攻机器人 有限公司	 上海陶素生化 科技有限公司	 上海唯析信息 科技有限公司	 上汽汽车集团 股份有限公司	 深圳奋达康 科技有限公司	 深圳晶泰 科技有限公司	 深圳市得润电子 股份有限公司
 北京凌空天行 科技有限责任公司	 北京荣之联科技 股份有限公司	 北京深蓝航天 科技有限公司	 北京深知无限人工 智能科技有限公司	 北京市建筑设计 研究院有限公司	 北京适创科技 有限公司	 北京星途探索 科技有限公司	 北京应力分析 科技有限公司	 北京月新时代科技 股份有限公司	 博郡汽车
 鸿之微科技(上海) 股份有限公司	 金锐同创(北京) 科技股份有限公司	 卡替(上海)细胞生 物技术有限公司	 立讯精密工业 股份有限公司	 蓝箭航天空间科技 股份有限公司	 洛阳德明石化 设备有限公司	 南京创蓝科技 有限公司	 南京道康达信息 技术有限公司	 普瑞基准科技 (北京)有限公司	 奇瑞汽车 股份有限公司
 人和未来生物科技 (长沙)有限公司	 石化盈科信息技 术有限责任公司	 天纳克(苏州) 排放系统公司	 西安西电开关电气 有限公司	 有研科技集团 有限公司	 中国建筑材料科学 研究总院有限公司	 中国汽车工程研究 院股份有限公司	 中国石油化工 股份有限公司	 中国长江三峡 集团有限公司

AI智算云产品与服务

自建算力成本高、影响科研速度

资源利用率普遍较低，大量资源闲置浪费，有效算力成本【 x 5倍 】



省心高效的GPU云



计算机视觉



自然语言处理



自动驾驶



新药研发

赋能领域

三线专家
贴心服务



TensorFlow



PyTorch



jupyter



PyCharm

AI智算云平台

开箱即用
易用高效



- 文件存储
- 块存储
- 对象存储



- A100
- V100
- T4



- 多线接入
- 弹性带宽
- 按需付费

AI智算云基础设施

丰富灵活
按需租用

AI智算云

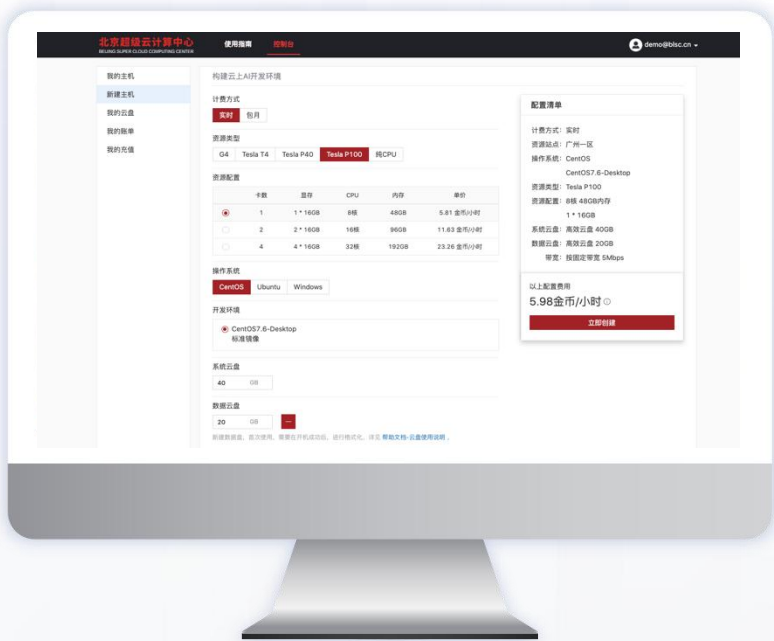
面向高校、科研院所、企事业单位在人工智能和性能计算等方向的GPU算力需求，提供专业的GPU算力云。

AI云主机平台

在训练、推理阶段
对GPU算力的复杂多样需求

HPC集群平台

在大规模多核计算场景中
利用GPU算力大幅度提升计算效率



算力资源丰富

- A100、V100、3090、A10、T4、国产DCU等主流型号的海量资源
- 支持多机多卡，满足推理、训练、科学计算等多种场景

高效易用

- 预置Tensorflow、PyTorch主流框架，内置数据集
- 通过简单易用的界面快速创建AI实例
- 裸金属服务器、独占显卡，性能强劲

高性价比

- 建设与运维“0”成本
- 按需灵活选用海量资源
- 用户将更多精力投入在科研

云主机平台

01

自主可控的云主机环境，拥有root权限资源专享

02

预装TensorFlow、PyTorch 等AI框架

03

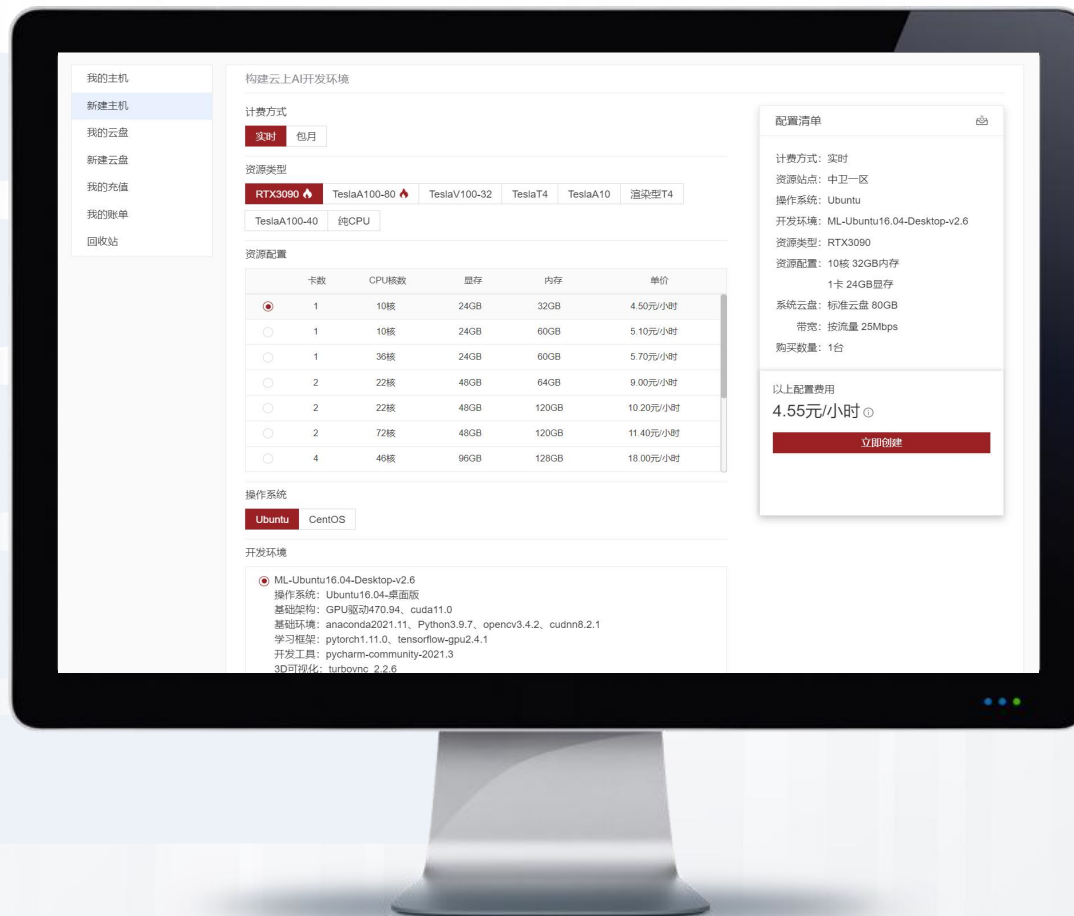
支持Windows、Ubuntu、CentOS多种操作系统

04

存储动态扩容，独立外网IP+带宽，灵活自定义配置

05

适用于小规模AI训练和科学计算，推理及图形处理场景

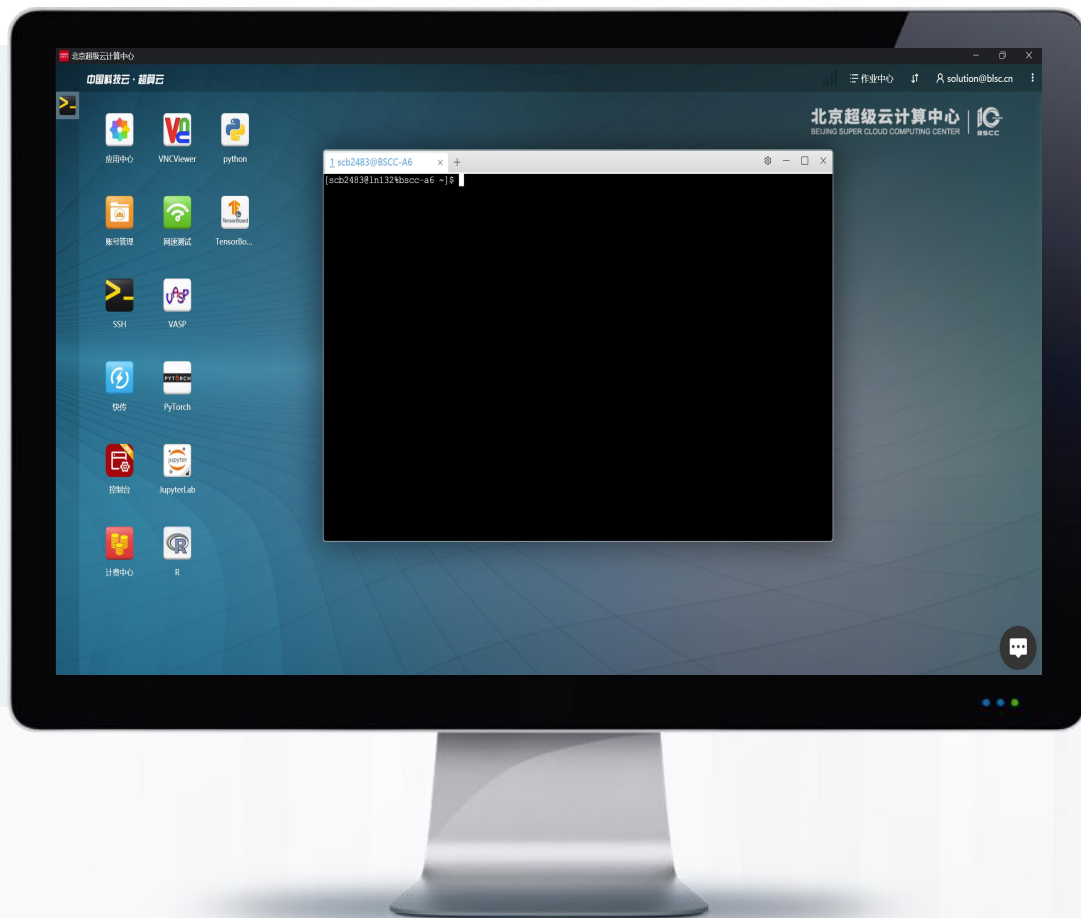


集群平台

- 01 基于超算架构构建，统一调度大规模集群资源
- 02 物理机性能直接输出，无虚拟化性能损耗
- 03 大容量、高性能分布式存储
- 04 节点间100Gbps，IB或RoCE互联，支持RDMA
- 05 适用于大规模跨节点AI训练和科学计算场景

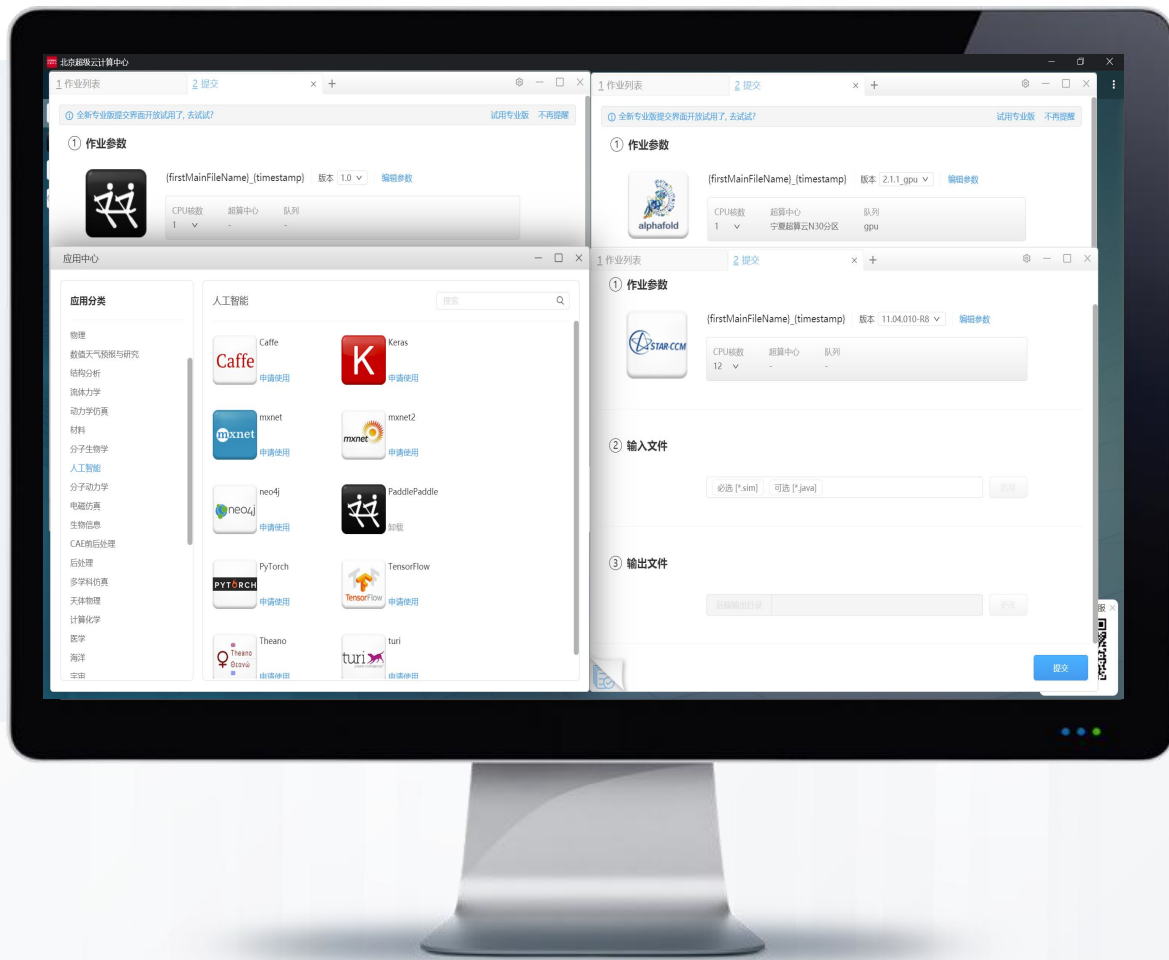


集群平台·图形化桌面



- ✓ web访问、客户端访问、SSH登录
- ✓ 图形化提交、脚本提交
- ✓ 作业状态、费用账单查询
- ✓ 共享存储，提供免费配额，支持按需购买扩容
- ✓ 集群资源基于CentOS linux操作系统

集群平台·应用SaaS化提交



- ✓ 适合于不擅长命令行操作的用户
- ✓ 通过图形方式，选择计算资源、计算文件
- ✓ 大幅度提高用户操作效率
- ✓ 命令行操作失误几率大幅度降低

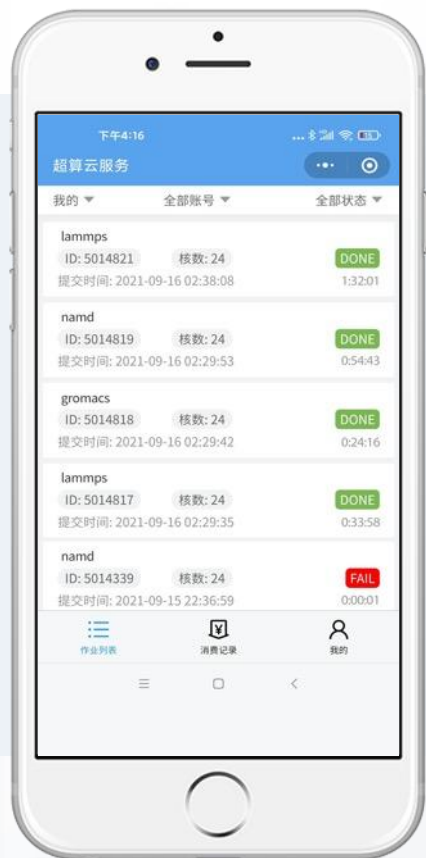
集群平台·文件传输



- ✓ 集成WinSCP、快传工具
- ✓ 支持IPV6传输协议
- ✓ 支持断点续传

集群平台·微信小程序时刻了解作业状态

随时随地查看 **监控快**



机时账单, 细微至每个作业 **报告快**



在线查看机时费用



提供有历史记录查询



清晰了解每分钱的产出

资源丰富

A100-80

GPU: 8 * NVIDIA®Tesla®A100 SXM4
显存: 8 * 80GB (2039 GB/s)
CPU: Intel 83系列 112vCPU
内存: 1960GB
NVLink: 双向通信600 GB/s
节点互联: 4 * 200G RoCE (RDMA协议)

CUDA® Cores: 8 * 6912
Tensor Cores: 8 * 432
单精度性能: 8 * 19.5 TeraFLOPS
双精度性能: 8 * 9.7 TeraFLOPS
TF32 Tensor Core: 8 * 312 TeraFLOPS

V100-32

GPU: 8 * NVIDIA®Tesla®V100 SXM2
显存: 8 * 32GB (897 GB/s)
CPU: Platinum82系列(80vCPU) v6
内存: 320GB
NVLink: 双向通信300 GB/s
节点互联: 100G RoCE (RDMA协议)

CUDA® Cores: 8 * 5120
Tensor Cores: 8 * 640
单精度性能: 8 * 15.7 TeraFLOPS
双精度性能: 8 * 7.8 TeraFLOPS
FP16 Tensor Core: 8 * 125 TeraFLOPS

RTX 3090

GPU: 8 * NVIDIA®GeForce®RTX 3090
显存: 8 * 24GB (936.2 GB/s)
CPU: AMD EPYC™ 7002 Series128
内存: 512GB
节点互联: 2 * 25G RoCE (RDMA协议)

CUDA® Cores: 8 * 10496
Tensor Cores: 8 * 328
单精度性能: 8 * 35.58 TeraFLOPS
双精度性能: 8 * 556.0 GigaFLOPS
FP16Tensor Core: 8 * 35.58 TeraFLOPS

DCU

DCU: 4*DCU加速卡, PCI-E G3 X16
CPU: Hygon C86 7185@2.0GHZ
单节点32核
内存: 128GB
节点互连: 200G

其他

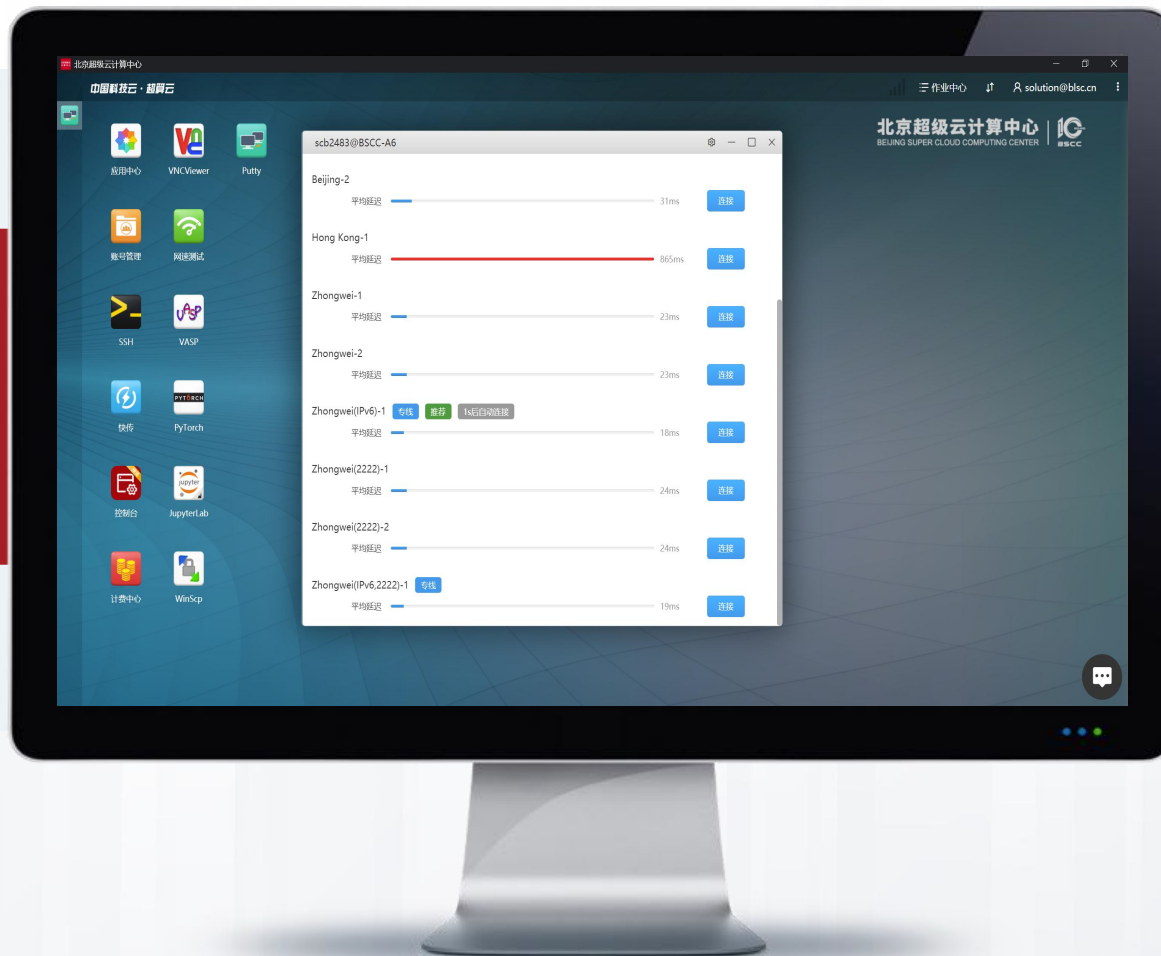
A100-40 \ V100-16 \ A30 \ A10
\ T4 \ 2080Ti

预置常用数据集



专线网络传输快

专线网络
传输快



7×24小时在线服务-服务快



基础支持服务

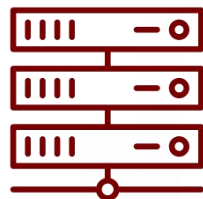
- ✓ 专属微信群
- ✓ 7 x 24hr x 5min
- ✓ 集群及系统相关支持
- ✓ 文档支持
- ✓ 软件安装支持
- ✓ 应用及计算相关支持

应用场景-人工智能

深度神经网络、特征抽取、图像分类、目标检测、语义分割、表示学习、生成对抗网络、语义网络、协同过滤和机器翻译等研究成为近年热点，相关技术已应用于计算机视觉、自然语言处理、语音处理及推荐系统等领域。AI智算云主机可灵活地满足相关人工智能技术研究，在训练和推理阶段对GPU算力复杂多样需求。

常用环境

TensorFlow、PyTorch、Caffe、MXNet、Keras、Anaconda、PyCharm、Jupyter、OpenCV、TensorBoard、Python、CUDA、TensorRT、NCCL等。



强劲算力

单机提供NVIDIA Tesla最新架构8卡GPU，显卡直通，性能强劲，支持多机多卡并行提升算力。GPU卡型号丰富，满足训练和推理等多种场景需求。



灵活配置

自主选择GPU卡数、CPU核数、存储容量、网络带宽及操作系统类型等，可根据业务需求灵活构建。



开箱即用

预置TensorFlow、PyTorch、PyCharm、TensorBoard等框架环境，分钟级获得实例环境，即开即用。

应用场景-高性能计算

在大规模多核计算场景中，GPU可大幅提高计算效率，使科研具有高出数量级的投入产出比，高性能程序GPU化趋势明显，GPU已广泛应用于生命科学、化学、材料、工业制造仿真设计、金融、气象海洋、油气能源等众多高性能计算领域。

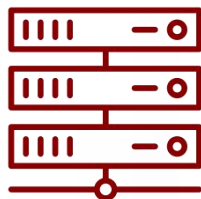
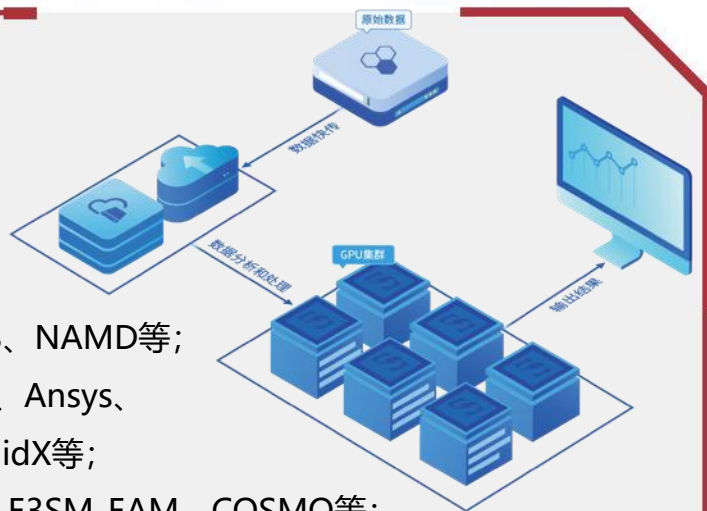
常用软件

化学领域： AMBER、GROMACS、LAMMPS、NAMD等；

有限元分析： ABAQUS、Ansys、OpenFOAM、nanoFluidX等；

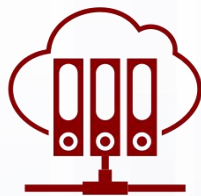
天气/环境建模： WRF、E3SM-EAM、COSMO等；

生命科学： Blast、AlphaFold2等。



超强算力

提供NVIDIA Tesla最新架构GPU卡，以裸金属形态输出，消除虚拟化性能损耗，输出极致的性能。



高速互连

100Gb/s高速互连，高性能并行存储，支持大规模并行。



弹性规模

大规模集群架构及GPU队列资源池，可根据计算需求便捷获取数以百卡计算资源。

云上安全及服务保障

物理安全



设施安全

- T4机房 SLA 99.99%



主机安全

- 专属云主机物理隔离
- 专属云可部署到用户本地



接入安全

- 专线或VPN接入

环境安全



用户认证

- 用户统一认证，用户间安全隔离
- 完善的账号防爆破策略



安全防护

- 分布式多层级防火墙和入侵检测系统
- DoS/DDoS防护机制



主动监测

- 漏洞管理及定期的SCAP检测
- 实时系统监控和可疑行为分析

数据安全



应用层

- 提供端到端的数据指纹校验，确保数据完整性



链路层

- 金融级强度加密算法，工业标准软件栈，遵循业界最佳实践，确保数据传输安全



系统层

- RAID6级别加“热备盘”数据保护
- 关键数据本地多副本机制
- 异地跨数据中心数据备份

服务保障



基础服务

- 微信专属支持群
- 7X24小时人工值守



高级服务

- 应用软件支持
- 作业问题分析



VIP服务

- 代码优化
- 应用性能评估

云主机安全性保障

数据安全机制高可用

数据采用持久化块级数据存储服务EBS云盘方式，云硬盘采用三副本的分布式机制，能够确保任何一个副本出现故障时迅速通过数据迁移等方式复制一个新副本，确保数据高可靠、高可用。

数据私密性

采用加密、安全组隔、和账号权限离等手段保证用户数据互不可见。



网络链路安全机制

网络安全等级保护标准达到三级备案；提供DDoS 基础防护，实时监控网络流量，发现攻击立即清洗，为公网 IP 秒级开启防护。

平台优势

资源丰富

提供云主机平台和裸金属云服务等多种平台，算力资源包括 A100\V100\3090\A10\T4\DCU等多种型号，支持多机多卡调度，满足训练、推理等不同场景计算需求。

易用高效

预装AI框架、内置数据集，可通过简单直观的界面快速创建AI计算实例，裸金属服务器、独占显卡，性能强劲。



服务贴心

三线专家团队7×24小时在线服务，提供开发框架部署、算例编译、调试优化等多元化服务支持，为AI客户深度赋能。

高性价比

建设与运维“0”成本，按需灵活选用海量资源，保障科研成果及时产出。不操心管，只专注用。

使用收益

零建设、维护成本

支持按需付费，避免购买和建设GPU服务器带来的高额成本投入以及维护计算设备的相关成本。

不排队

基于国家网格，计算资源丰富，预留足够的备用计算资源，不排队。



低门槛

预置Tensorboard、Pytorch等计算环境软件及常用数据集；支持脚本提交、支持图形化集成提交作业；可协助用户安装软件。

成本节省

相对于使用率不高自建方式，使用云计算的方式可节省成本50%。



客户案例



某985高校人工智能学院

满足日常和紧急课题算力需求



合作背景

在课题组新生入组时，需要按人数采购GPU服务器支撑实验，采购硬件及硬件存放、管理耗费大量精力；
因课题周期等原因，出现短期GPU算力高峰时，无足够资源可用。

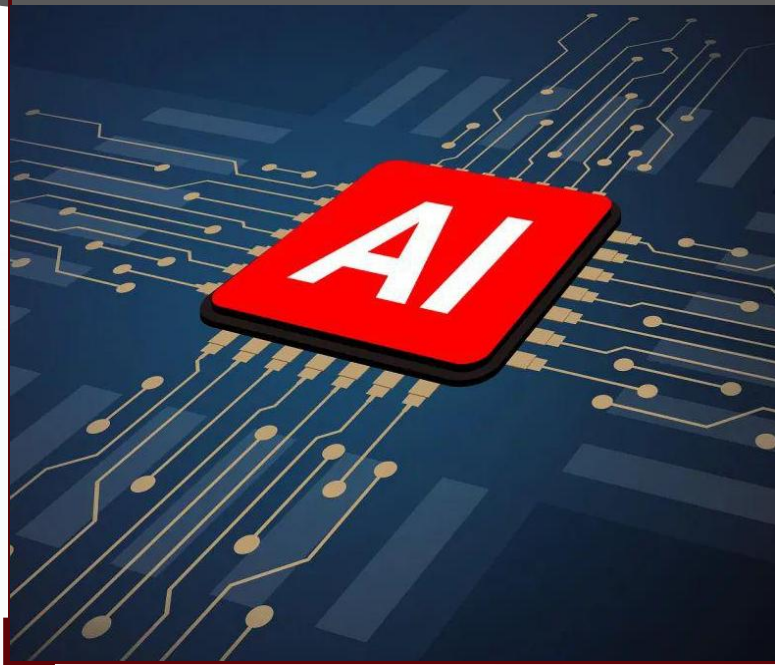
合作收益

AI智算云与学院进行联合深度合作，通过云端资源和云账号，灵活满足课题组成员流动产生的账号增减、权限分配、因课题周期产生的短时、大量算力等多样化GPU算力需求，灵活按需使用，让科研人员不因算力和平台使用分心，专注科研，提升效率。



中科院某研究所

保障按时完成基金课题研究



合作背景

某语义分割相关基金项目周期紧，课题科研人员需在短时间内调用数十张 Nvidia V100 GPU资源进行项目验证，面对此难题，无论是采购周期亦或是基金支持能力，依靠自采服务器方式显然行不通。

合作收益

基于AI智算云丰富的资源池，科研人员在瞬间获得规模化的GPU云算力，为课题提供算力保障。此外，AI智算云的专家团队基于大量的调优经验，在程序运行效率方面进一步赋能课题科研人员，加速科研产出，保障课题及时顺利完成。



某新药研发企业

加速研发，项目时间缩短50%



合作背景

某新药研发企业，由于业务增长迅速，需要利用人工智能前沿算法，结合计算化学和药物化学的经验，对数十亿分子的进行快速筛选，具有研发周期短、计算量大、高通量、存储量大等特性。由于本地GPU资源不足，计算资源扩容速度难以匹配项目增长需求。

合作收益

客户在AI智算云上构建药物研发平台，不仅能够支持多团队、多任务并行的模型训练，同时可以并行调度上百块Nvidia V100GPU资源，用于PyTorch、TensorFlow大规模分布式训练以及GROMACS、LAMMPS等分子动力学模拟，将项目时间缩短50%以上，大幅降低研发成本的风险。



某人工智能初创企业

跨越起步门槛，业务快速上线



合作背景

该企业的基于客流统计分析的商业智能系统具有以下特点

- 处理万路视频
- 数量存储量大
- 成本控制严苛
- 传输带宽高
- 推理算力弹性需求度高

合作收益

通过AI智算云赋能，高效满足业务需求，屏蔽底层IT繁杂的技术细节，获取多方收益

- 基于云化基础设施，快速灵活构建与扩展业务系统
- 研发级业务移植服务支持
- IT维护方面无需分心
- 起步投入门槛低，规避风险

云上科研更省心 企业提效新引擎

200,000科研用户的信赖首选

