

联系我们

**CONTACT US**



邮箱: [sales@infrawaves.com](mailto:sales@infrawaves.com)

地址: 北京市海淀区中关村东路东升大厦AB座505A

网址: <https://www.infrawaves.com>



# 64/128集群

RoCE集群组网方案设计

2024.5.23

一期64台组网拓扑

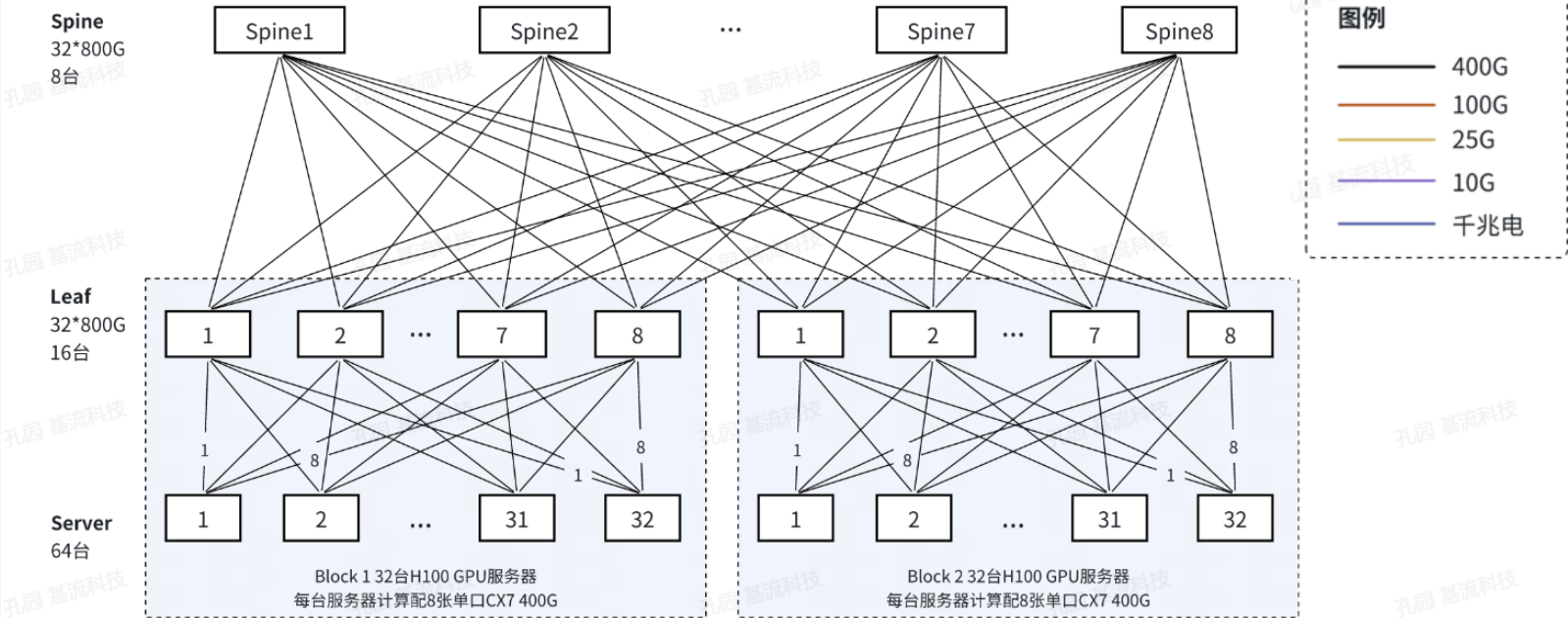
Network Topology

组网说明

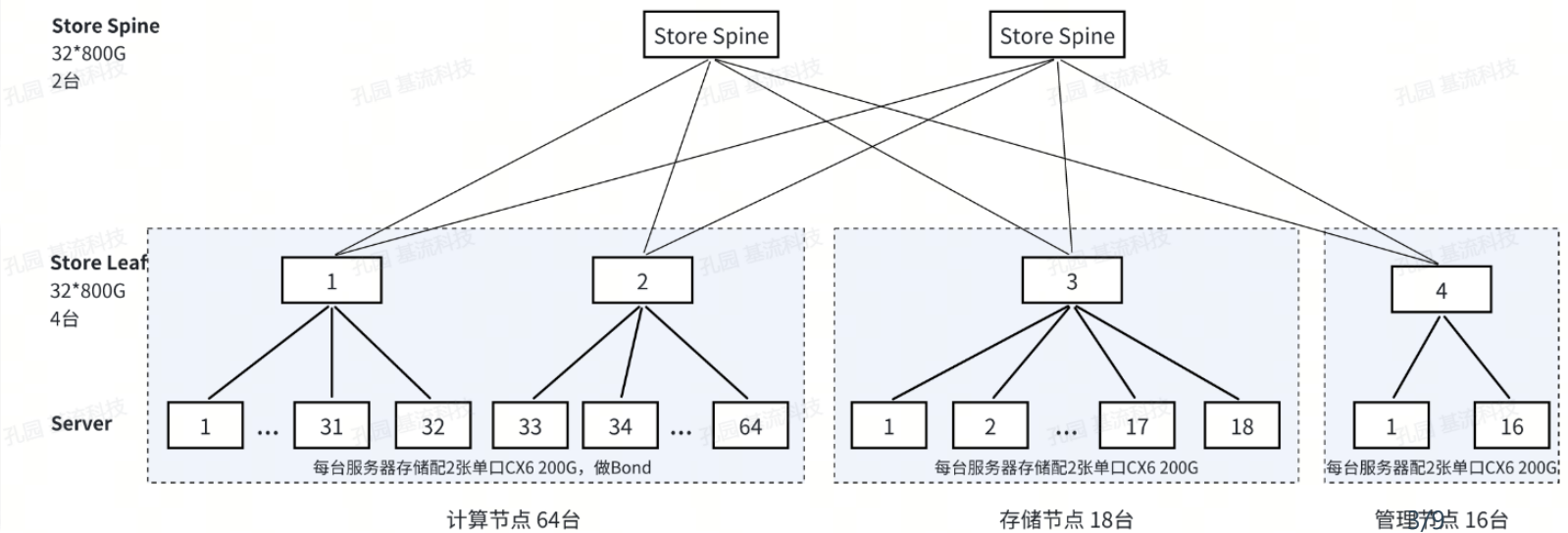
**计算网络：**  
每台服务器配8\*400G单口CX7网卡，接入计算网；  
计算网选用25.6T IB交换机MQM9790-NS2F， 32\*800G；  
计算网保证1: 1收敛比， Spine 8台， Leaf 16台。

**存储网络：**  
每台服务器配2\*200G单口CX6网卡，接入存储网；  
计算网选用25.6T IB交换机MQM9790-NS2F， 32\*800G；  
存储网总共使用6台交换机， Spine 2台， Leaf 4台。

计算网平面



存储网平面



## 组网说明

### 业务管理网络:

每台服务器配2\*10G/25G网卡，接入业务管理网；  
服务器接入25G以太网交换机，再接入100G交换机；  
业务管理网使用6台25G交换机，4台100G交换机。

### 带外管理网络:

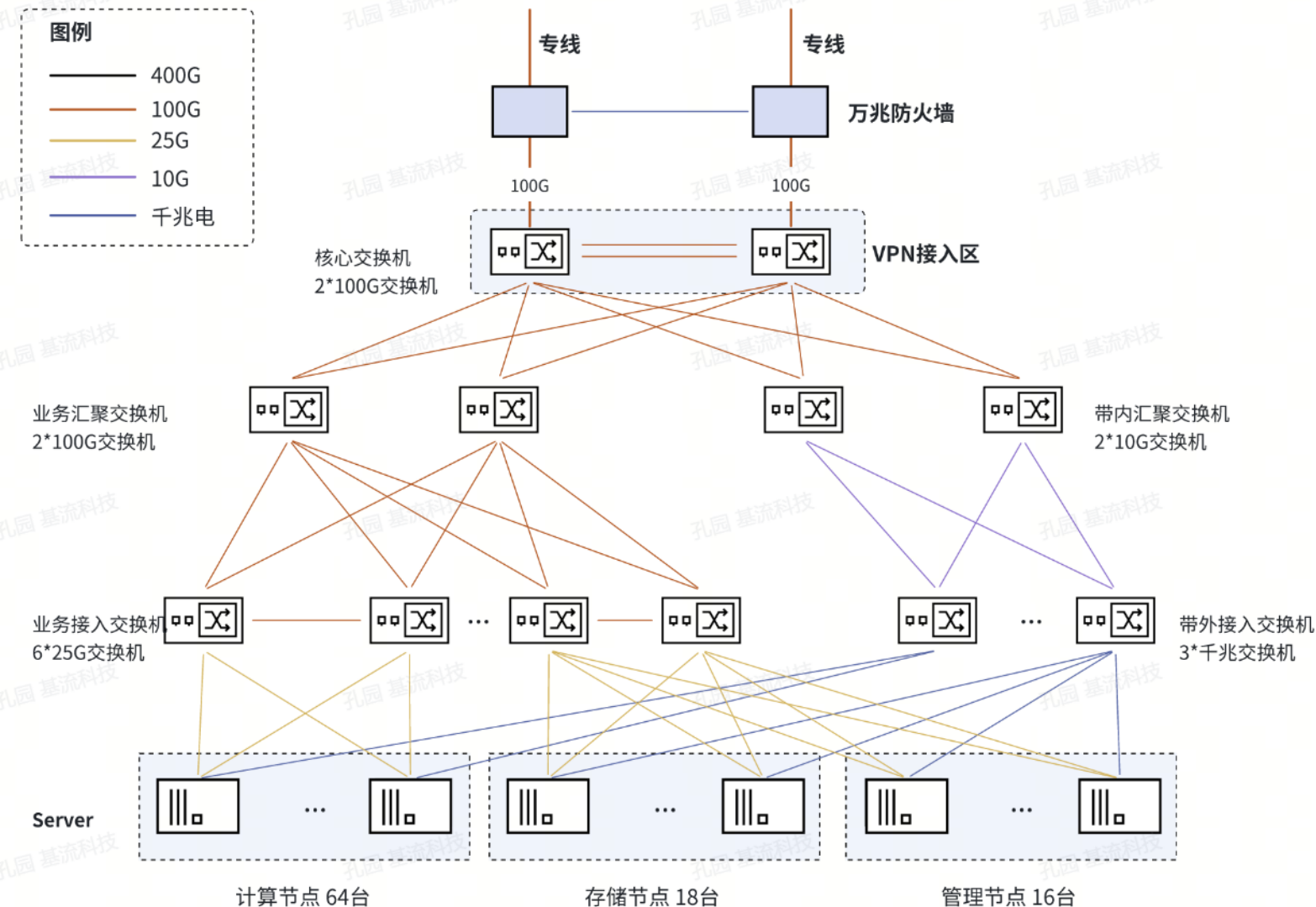
每台服务器通过千兆双绞线接入带外管理网，再接入10G交换机；  
带外管理网使用3台千兆交换机，2台10G交换机。

### 一期128节点与64节点spine优先布局的成本差异:

计算网spine交换机增加8台，存储网spine交换机增加1台

**一期总成本增加 $9 \times 20.8\text{万} = 187.2\text{万元}$ 。**

## 管理网平面



系统	名称	组件	技术规格	数量	单价	总价	售后维保
计算网络	IB交换机	MQM9790-NS2F	提供32个800Gb/s端口	24	208,000	4,992,000	三年
	800G光模块	MMS4X00-NS	NVIDIA twin port transceiver, 800Gbps,2xNDR, OSFP, 2xMPO12 APC, 1310nm S MF, up to 100m, finned	768	14,500	11,136,000	三年
	400G光模块	MMS4X00-NS400	NVIDIA single port transceiver, 400Gbps,NDR, OSFP, MPO12 APC, 1310nm SMF , up to 100m, flat top	512	9,800	5,017,600	三年
	光纤线缆	MFP7E30-N050	NVIDIA passive fiber cable, SMF, MPO12 APC to MPO12 APC, 50m	1024	850	870,400	三年
	端网一体监控运维平台	IW-Venus-企业版	按照GPU数计算，现对集群服务器、交换机、线缆的配置、调优、巡检、监测与分析预警	512	4,800	2,457,600	三年
	集群运维及技术支持服务	IW-Support	AI集群服务调优，日常技术支持服务等，定位集群故障，提出解决方案；按照月数计算，维持3年	36	15,000	540,000	三年
存储网络	IB交换机	MQM9790-NS2F	提供32个800Gb/s端口	6	208,000	1,248,000	三年
	800G光模块	MMS4X00-NS	NVIDIA twin port transceiver, 800Gbps,2xNDR, OSFP, 2xMPO12 APC, 1310nm S MF, up to 100m, finned	98	14,500	1,421,000	三年
	200G光模块	MMA1Z00-NS400	NVIDIA single port transceiver, 400Gbps,NDR, QSFP112, MPO, 850nm MMF, SR4, up to 30m, flat top	196	7,200	1,411,200	三年
	光纤线缆	MFP7E20-N030	NVIDIA passive fiber cable, MMF, MPO12 APC to 2xMPO12 APC, 50m	98	950	93,100	三年
	端网一体监控运维平台	IW-Venus-企业版	按照网卡数计算，现对集群服务器、交换机、线缆的配置、调优、巡检、监测与分析预警	196	4,800	940,800	三年
存储系统	分布式全闪存储	全闪存储服务器	全闪存储服务器-单口200G网卡2张-配置Nvme SSD盘，15.36T，22盘位	18	950,000	17,100,000	三年
		存储软件	分布式存储软件				
业务网络	以太网	接入交换机	包含48个25G SFP28端口,6个100G QSFP28端口	6	60,000	360,000	三年
		汇聚/核心交换机	32个100G QSFP28端口	4	110,000	440,000	三年
管理网	以太网	接入交换机	48个千兆口，4*10GSFP+口	3	3,800	11,400	三年
		汇聚交换机	48*10GSFP+口，6*100GE QSFP28	2	11,000	22,000	三年
安全设备	防火墙	边界防火墙	硬件架构：采用非X86多核架构，前后通风。 硬件性能：防火墙吞吐量≥20G并发连接数≥1600万，每秒新建连接数（HTTP）≥50万。 接口规格：9个千兆电口+4个千兆光口+4个万兆光口+4对Combo口。 功能：IPSec VPN、链路负载、流量控制功能。	2	150,000	300,000	三年
	堡垒机	堡垒机	堡垒机	2	100,000	200,000	三年
	日志审计	日志审计	日志审计	2	60,000	120,000	三年
总价						48, 681, 100	



# 二期128台组网拓扑

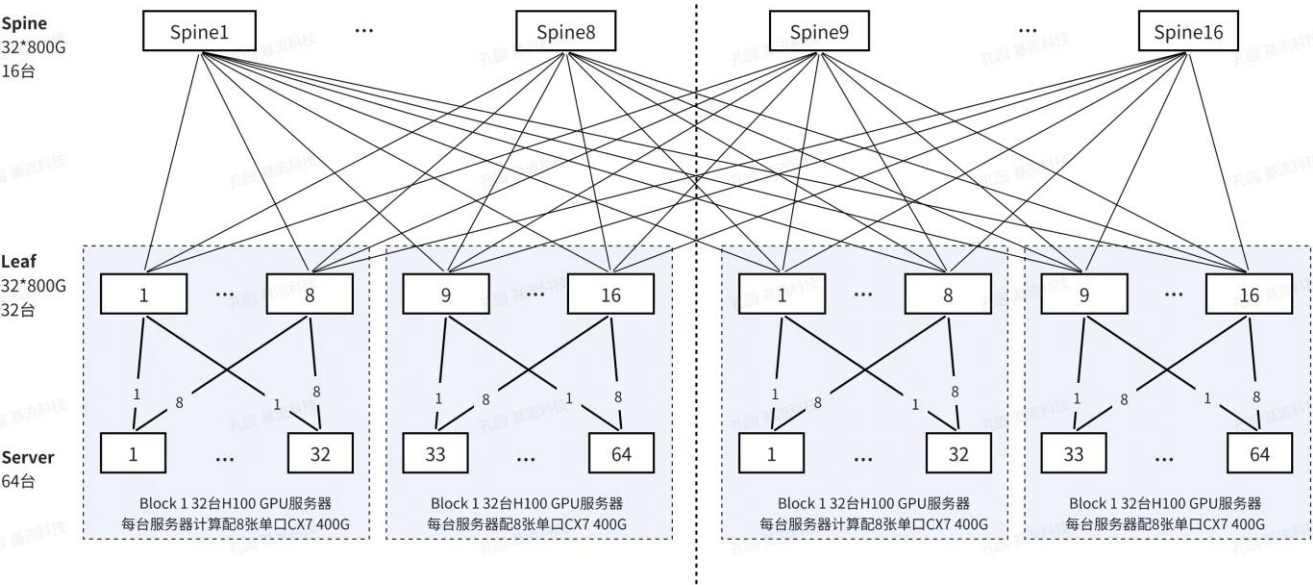
## Network Topology

### 组网说明

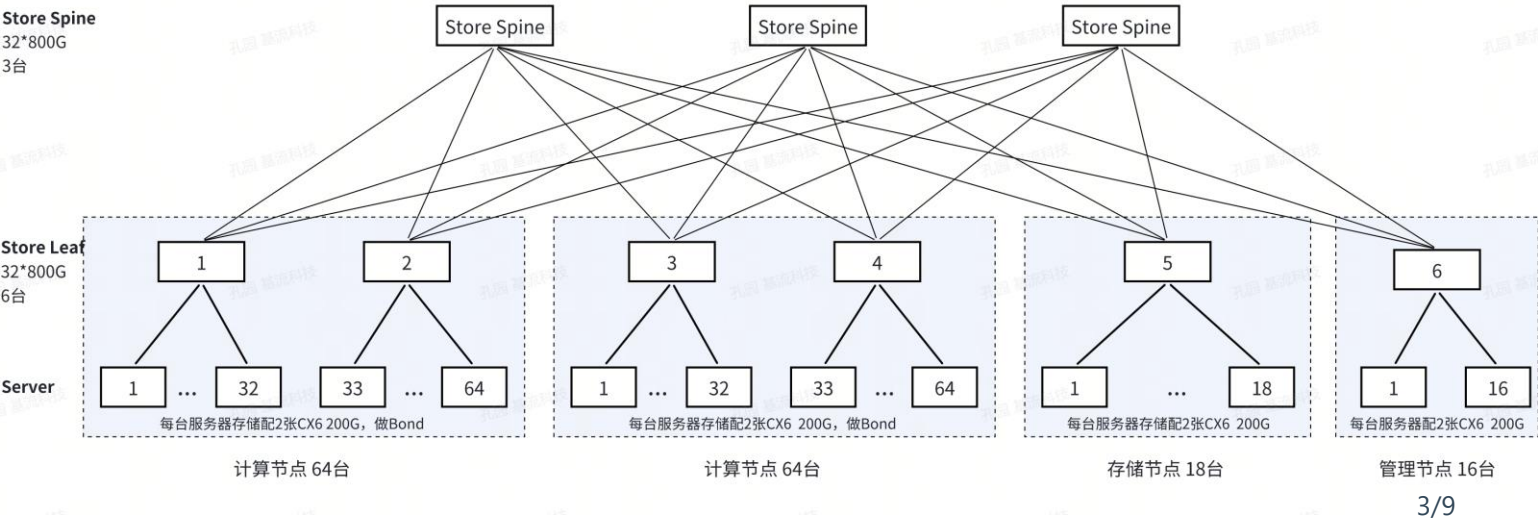
**计算网络：**  
每台服务器配8\*400G单口CX7网卡，接入计算网；  
计算网选用25.6T IB交换机MQM9790-NS2F， 32\*800G；  
计算网保证1：1收敛比，Spine 16台，Leaf 32台。

**存储网络：**  
每台服务器配2\*200G单口CX6网卡，接入存储网；  
计算网选用25.6T IB交换机MQM9790-NS2F， 32\*800G；  
存储网总共使用6台交换机，Spine 3台，Leaf 6台。

### 计算网平面



### 存储网平面



# 二期128台组网拓扑

## Network Topology

### 组网说明

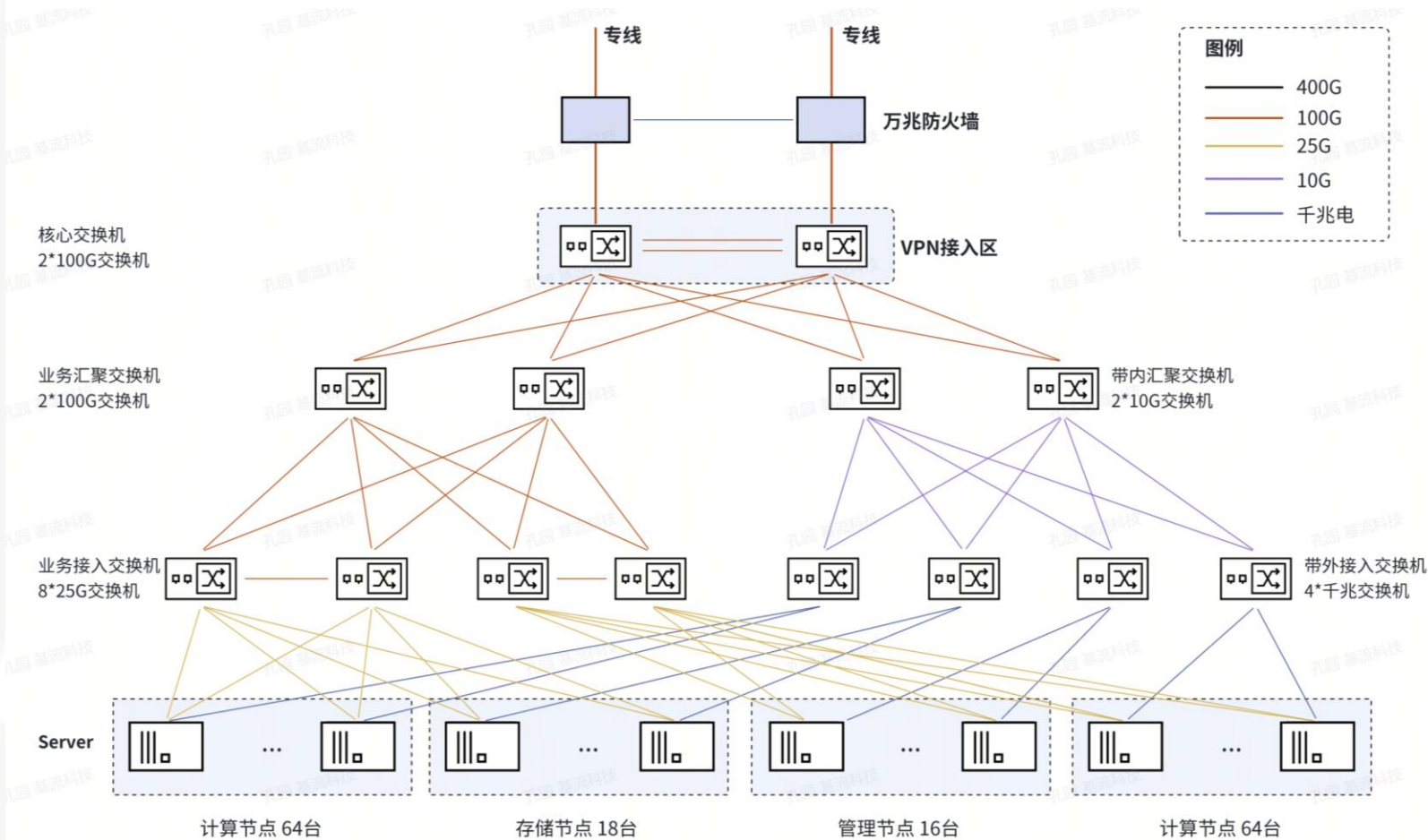
#### 业务管理网络:

每台服务器配2\*10G/25G网卡，接入业务管理网；  
服务器接入25G以太网交换机，再接入100G交换机；  
业务管理网使用8台25G交换机，4台100G交换机。

#### 带外管理网络:

每台服务器通过千兆双绞线接入带外管理网，再接入10G交换机；  
带外管理网使用4台千兆交换机，2台10G交换机。

### 管理网平面



系统	名称	组件	技术规格	数量	单价	总价	售后维保
计算网络	IB交换机	MQM9790-NS2F	提供32个800Gb/s端口	48	208,000	9,984,000	三年
	800G光模块	MMS4X00-NS	NVIDIA twin port transceiver, 800Gbps, 2xNDR, OSFP, 2xMPO12 APC, 1310nm S MF, up to 100m, finned	1536	14,500	22,272,000	三年
	400G光模块	MMS4X00-NS400	NVIDIA single port transceiver, 400Gbps, NDR, OSFP, MPO12 APC, 1310nm SMF , up to 100m, flat top	1024	9,800	10,035,200	三年
	光纤线缆	MFP7E30-N050	NVIDIA passive fiber cable, SMF, MPO12 APC to MPO12 APC, 50m	2048	850	1,740,800	三年
	端网一体监控运维平台	IW-Venus-企业版	按照GPU数计算， 现对集群服务器、交换机、线缆的配置、调优、巡检、监测与分析预警	1024	4,800	4,915,200	三年
	集群运维及技术支持服务	IW-Support	AI集群服务调优， 日常技术支持服务等， 定位集群故障， 提出解决方案； 按照月数计算， 维持3年	36	30,000	1,080,000	三年
存储网络	IB交换机	MQM9790-NS2F	提供32个800Gb/s端口	9	208,000	1,872,000	三年
	800G光模块	MMS4X00-NS	NVIDIA twin port transceiver, 800Gbps, 2xNDR, OSFP, 2xMPO12 APC, 1310nm S MF, up to 100m, finned	162	14,500	2,349,000	三年
	200G光模块	MMA1Z00-NS400	NVIDIA single port transceiver, 400Gbps, NDR, QSFP112, MPO, 850nm MMF, SR4, up to 30m, flat top	324	7,200	2,332,800	三年
	光纤线缆	MFP7E20-N030	NVIDIA passive fiber cable, MMF, MPO12 APC to 2xMPO12 APC, 50m	162	950	153,900	三年
	端网一体监控运维平台	IW-Venus-企业版	按照网卡数计算， 现对集群服务器、交换机、线缆的配置、调优、巡检、监测与分析预警	324	4,800	1,555,200	三年
存储系统	分布式全闪存储	全闪存储服务器	全闪存储服务器-单口200G网卡2张-配置Nvme SSD盘， 15.36T， 22盘位	18	950,000	17,100,000	三年
		存储软件	分布式存储软件				
业务网络	以太网	接入交换机	包含48个25G SFP28端口,6个100G QSFP28端口	8	60,000	480,000	三年
		汇聚/核心交换机	32个100G QSFP28端口	4	110,000	440,000	三年
管理网	以太网	接入交换机	48个千兆口， 4*10GSFP+口	4	3,800	15,200	三年
		汇聚交换机	48*10GSFP+口， 6*100GE QSFP28	2	11,000	22,000	三年
安全设备	防火墙	边界防火墙	硬件架构：采用非X86多核架构， 前后通风。 硬件性能：防火墙吞吐量≥20G并发连接数≥1600万， 每秒新建连接数（HTTP） ≥50万。 接口规格：9个千兆电口+4个千兆光口+4个万兆光口+4对Combo口。 功能：IPSec VPN、链路负载、流量控制功能。	2	150,000	300,000	三年
	堡垒机	堡垒机	堡垒机	2	100,000	200,000	三年
	日志审计	日志审计	日志审计	2	60,000	120,000	三年
总价						76, 967, 300	

整体128台组网平面机柜布置

LAYOUT PLAN

机柜平面布置说明

单机柜15kw:

可容纳1\*H100+2\* MQM9790

可容纳1\*H100+ (3~4) \* 管理服务器/存储服务器

服务器与交换机/管理服务器/管理交换机叠放可最大资源利用单机柜功率。

注: H100 满载功率约10kw, 9790交换机满载功率约1720kw.

服务器机柜布置:

- 根据单机柜功率15kw, 每机柜均放置一台H100服务器, 该机房最大容纳126台H100 服务器;
- 如果要满足该机房128台H100服务器, 如有条件可考虑将单机柜15kw改为10kw或20kw组合布置, 单位机柜模组的总功率未变。

交换机机柜布置:

- Spine交换机集中布置于机房中间区域, 可适当缩短光纤长度, 节约光纤成本; 同时便于布线与管理;
- Leaf交换机挨着对应Block服务器布置, 可缩短光纤长度, 节约光纤成本;
- Leaf与Spine交换机规律布置, 便于后期运维管理。

管理服务器/存储服务器机柜布置:

- 管理服务器/存储服务器与GPU服务器放置于同一机柜内, 根据实际需求来布置。

机柜布置平面示意图



GPU001	A14	B14	GPU015		GPU029	C14	D14	GPU043		GPU057	E14	F14	GPU071		GPU085	G14	H14	GPU099		GPU113	I14
GPU002	A13	B13	GPU016		GPU030	C13	D13	GPU044		GPU058	E13	F13	GPU072		GPU086	G13	H13	GPU100		GPU114	I13
GPU003	A12	B12	GPU017		GPU031	C12	D12	GPU045		GPU059	E12	F12	GPU073		GPU087	G12	H12	GPU101		GPU115	I12
GPU004	A11	B11	GPU018		GPU032	C11	D11	GPU046		GPU060	E11	F11	GPU074		GPU088	G11	H11	GPU102		GPU116	I11
GPU005	A10	B10	GPU019 +2LEAF		GPU033	C10	D10	GPU047 +2LEAF		GPU061 +2Spine	E10	F10	GPU075 +2LEAF		GPU089	G10	H10	GPU103 +2LEAF		GPU117	I10
GPU006	A9	B9	GPU020 +2LEAF		GPU034	C9	D9	GPU048 +2LEAF		GPU062 +2Spine	E9	F9	GPU076 +2LEAF		GPU090	G9	H9	GPU104 +2LEAF		GPU118	I9
GPU007	A8	B8	GPU021 +2LEAF		GPU035	C8	D8	GPU049 +2LEAF		GPU063 +2Spine	E8	F8	GPU077 +2LEAF		GPU091	G8	H8	GPU105 +2LEAF		GPU119	I8
GPU008	A7	B7	GPU022 +2LEAF		GPU036	C7	D7	GPU050 +2LEAF		GPU064 +2Spine	E7	F7	GPU078 +2LEAF		GPU092	G7	H7	GPU106 +2LEAF		GPU120	I7
GPU009	A6	B6	GPU023		GPU037	C6	D6	GPU051		GPU065 +2Spine	E6	F6	GPU079		GPU093	G6	H6	GPU107		GPU121	I6
GPU010	A5	B5	GPU024		GPU038	C5	D5	GPU052		GPU066 +2Spine	E5	F5	GPU080		GPU094	G5	H5	GPU108		GPU122	I5
GPU011	A4	B4	GPU025		GPU039	C4	D4	GPU053		GPU067 +2Spine	E4	F4	GPU081		GPU095	G4	H4	GPU109		GPU123	I4
GPU012	A3	B3	GPU026		GPU040	C3	D3	GPU054		GPU068 +2Spine	E3	F3	GPU082		GPU096	G3	H3	GPU110		GPU124	I3
GPU013	A2	B2	GPU027		GPU041	C2	D2	GPU055		GPU069	E2	F2	GPU083		GPU097	G2	H2	GPU111		GPU125	I2
GPU014	A1	B1	GPU028		GPU042	C1	D1	GPU056		GPU070	E1	F1	GPU084		GPU098	G1	H1	GPU112		GPU126	I1



存储方案

STORAGE SOLUTION

存储方案说明

存储要求:

闪存:

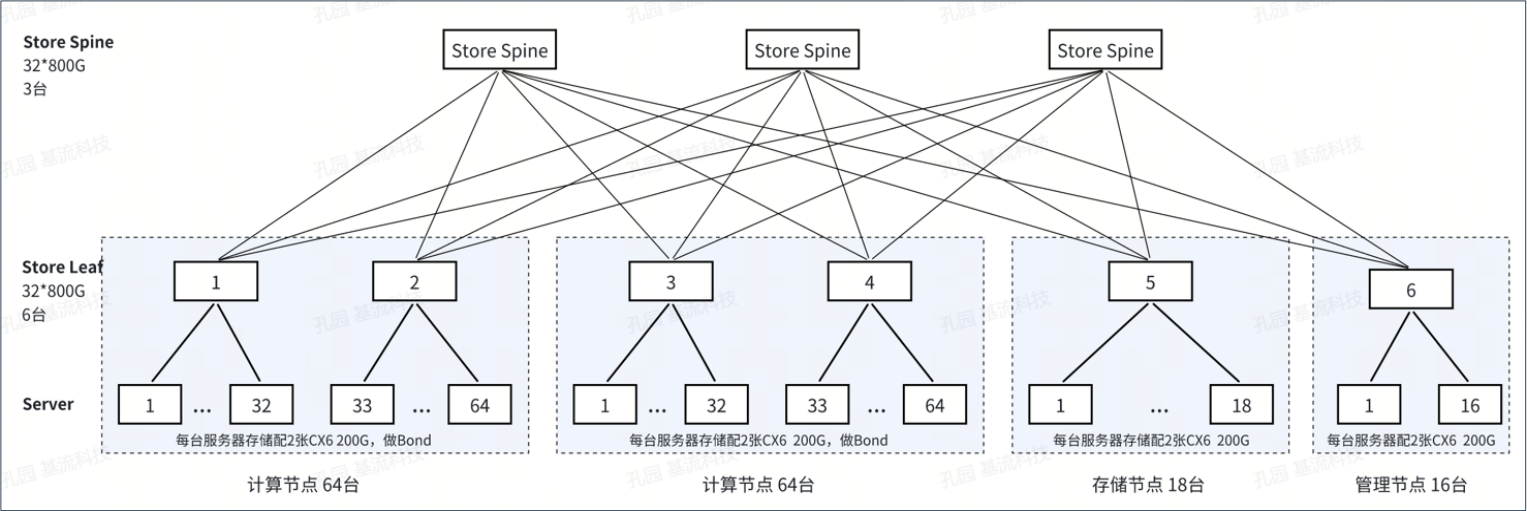
- IO压力大: 百亿以上小文件不能有性能衰减
- 文件系统挂载在服务器; 容量需求128节点3P, 支持容量扩展
- 存储方案数据安全性需要绝对保证

对象存储

- 容量需求128节点3P, 支持容量扩展
- 读写速度越快越好
- 存储方案数据安全性需要绝对保证

并行存储方案:

- **双副本**, 保证存储数据的安全和读写速度
- 配置: 18台存储服务器, 杰点裸容量合计6083TB, 可用容量为3041TB, 满足使用需求。
- 存储读写性能满足**英伟达BEST存储性能**



H100/H800 GPU 算力规模	GOOD	BETTER(GB/s)	BEST(GB/s)
单个 Scalable Unit(32个计算节点)	读:15 GB/s 写:7 GB/S	读:40 GB/S 写:20 GB/s	读:125 GB/s 写:62 GB/s
4 Scalable Unit (1 SuperPod,128个计算节点)	读:60 GB/s 写:30 GB/s	读:160 GB/S 写:80 GB/s	读:500 GB/s 写:250 GB/S

存储服务器配置

CPU:2颗AMDEPYC32核处理器;  
内存:256GB;  
系统盘:2块480GBSATASSD;  
元数据盘:2块1.6TBNVMe SSD;数据盘:支持3.84TB、7.68TB、15.36TB等容量NVMeSSD磁盘, 最大可扩展至22块;  
管理网络:1G  
存储网路:支持100/200GbpsInfiniBand、Ethernet 等网络, 支持RDMA/RoCE;节点默认无数据盘、分布式存储系统软件

# 服务器冷板式液冷改造方案

## MAINTENANCE AND REPLACEMENT

### 液冷改造方案选择：比较业界主流的两​​种液冷解决方案

液冷方案	投资成本	PUE	可维护性	应用案例	分析	是否推荐
非接触式液冷——冷板式	初始投资中低运维成本	1.1-1.2	较简单	多	初始投资中等、运维成本低，PUE收益中等，部署方式与风冷相同，从传统模式过渡较平滑	✓
接触式液冷——单项浸没式	初始投资及运维成本高	< 1.09	复杂	较多	初始投资较高，PUE收益较高部分部件不兼容，服务器结构需改造	×

### 液冷改造方案：冷板式液冷

- 非接触式液冷主要指冷板式液冷，将服务器发热元件（CPU/GPU/DIMM 等）贴近冷板，液体在冷板内流动，带走发热元件的热量，液体不与发热源直接接触，冷却液多采用 去离子水。
- 需少量风扇对服务器中的非液冷元件进行风冷散热，还需考虑液 体泄露风险。

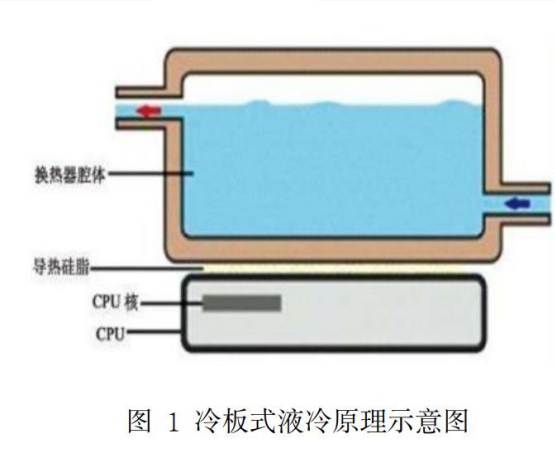


图 1 冷板式液冷原理示意图

### 冷板式液冷部件的选择：

#### 1、冷却液

- 液冷冷却液目前业内选择有乙二醇溶液、丙二醇溶液、去离子水等。冷却液浓度 建议在 20%~30%，浓度不宜过高，过高会影响工质散热性能；也不宜过低，过低会影响 防冻和抑制微生物滋生的能力。去离子水具有良好的传热性能，无毒安全，可作为冷却液之一，但需注意对冷却液的 维护。

#### 2、快接头

- 快接头是用于节点冷板模组和液冷机柜集分水器之间的水路连接接头，需支持插拔节 点时快速连通和截断节点与液冷机柜集分水器之间的水路，并保证不漏液。快接头分为手插快接头和盲插快接头两种形态。

#### 3、冷量分配单元

- 冷板式液冷系统中通过 CDU 隔离一次侧和二次侧回路，并完成一次侧和二次侧的热 交换，为服务器提供制冷能力。根据 CDU 的形态和部署位置，可分为集中式 CDU 和分 布式 CDU 两种。集中式 CDU 的单台 CDU 可以同时为多个服务器机柜提供制冷能力，可以通过多台 CDU 集群实现 N+M 的冗余能力，可靠性高，适用于规模部署液冷服务器机柜的场景。

#### 4、冷板

- 冷板是与芯片接触实现换热的核心部件，冷却液在内部流动将芯片热量带走。根据散 热模块与固定模块的可拆卸性，可以分为一体式冷板和分体式冷板。一体式冷板的散热模 块和固定模块不可拆卸，分体式冷板的散热模块和固定模块通过螺钉固定，可拆卸。

#### 5、服务器液冷管路

- 服务器液冷管路作为输送冷却工质的通道

## 服务器维保方案

服务器维保主要包括**硬件维护**，**软件维护和预防性维护**

### 1、硬件维护

#### • 清洁和检查:

外部清洁: 定期使用无静电布或特定的电子设备清洁剂清理服务器的外壳, 以防尘土积聚可能导致散热问题。

内部清洁: 定期关闭电源, 并在防静电条件下开启机箱, 清理内部积尘, 特别是风扇和散热器等关键散热组件。

连接检查: 检查所有内部连接, 包括电源线和数据线, 确保连接牢固且无损伤。

#### • 硬件检查:

GPU卡检查: 检查GPU卡是否存在过热或性能下降的迹象, 必要时进行更换或升级。

风扇和散热系统: 检查风扇是否运行正常, 散热片是否安装牢固, 保证散热效果。

### 2、软件维护

#### • 系统更新:

操作系统和驱动程序: 定期检查并更新操作系统和所有硬件的驱动程序, 确保系统运行稳定, 并兼容最新的应用程序。

固件更新: 检查GPU和其他关键硬件组件的固件版本, 按照制造商的推荐进行更新。

#### • 性能监控和日志审查:

资源监控: 使用系统监控工具定期检查CPU、GPU、内存和存储的使用情况, 确保资源被有效管理, 避免过载。

日志审查: 定期审查系统日志, 识别可能的错误或故障预兆, 及早处理。

### 3、预防性维护

#### • 备份计划:

数据备份: 定期备份重要数据和配置设置, 确保在硬件故障或其他意外情况下可以快速恢复。

完整性验证: 定期验证备份数据的完整性和可恢复性。

#### • 环境监控:

温度和湿度控制: 确保服务器所在环境的温度和湿度符合制造商的推荐标准, 避免环境因素导致设备性能下降或损坏。

通过以上维保措施, 可以有效保障H100 GPU服务器的性能和延长其使用寿命, 支持企业持续运行其关键应用和服务。这些措施需结合实际使用情况进行调整, 以适应不同的运行环境和业务需求。

## 服务器坏件替换方案

如服务器组件出现故障, **及时报修, 找有渠道的厂商更换坏件**; 如若服务器为**自持**, 可寻求**专业的维修团队, 将故障机器组件互相替换**, 保证部分故障服务器正常运行。例如GPU故障, 可将单GOU卡更换。

### 1.准备工作

#### • 确认坏件:

运行诊断软件确定需要替换的硬件部件。

检查保修状态和替换部件的可用性。

#### • 获取替换部件: 确保购买或获取NVIDIA或认证供应商提供的正品替换部件。

#### • 工具和环境准备:

准备必要的工具, 如螺丝刀、防静电带或手环等。

在无尘、干燥、防静电的环境中进行替换操作。

#### • 备份数据: 如果替换的是存储设备或任何可能影响数据的部件, 务必先进行数据备份。

#### • 计划停机时间: 安排在业务低谷期进行替换, 通知相关利益相关者服务器将暂时停机。

### 2.替换步骤

#### 替换GPU卡

#### • 关闭服务器: 关闭操作系统, 断开电源, 并拔掉所有连接到服务器的电缆。

#### • 打开机箱: 使用适当的工具打开服务器机箱。

#### • 拆卸坏的GPU卡: 卸下固定GPU卡的螺丝; 轻轻握住GPU卡的边缘, 先将其从PCIe插槽轻轻拉出。

#### • 安装新的GPU卡: 将新的GPU卡对准PCIe插槽, 确保正确方向, 轻轻压入直到固定; 然后重新安装螺丝固定GPU卡。

#### • 连接电缆: 根据需要连接GPU卡所需的电源和数据线。

#### • 关闭机箱: 替换完毕后, 关闭并固定服务器机箱。

### 3.测试与验证

#### • 启动服务器: 开启电源, 启动服务器, 进入BIOS确保新的GPU卡被系统识别。

#### • 运行诊断程序: 使用服务器管理软件或硬件诊断工具检查新GPU卡的状态和性能是否正常。

#### • 性能测试: 运行相关的性能测试程序, 确保GPU卡的性能达到预期。

#### • 监控系统日志: 检查系统日志, 确保没有错误或警告信息。

这种方案也可以适用于服务器中其他关键部件的替换, 如内存、硬盘或电源风扇等。



产品方案

# OUR PRODUCT

## Helios 基流羲和

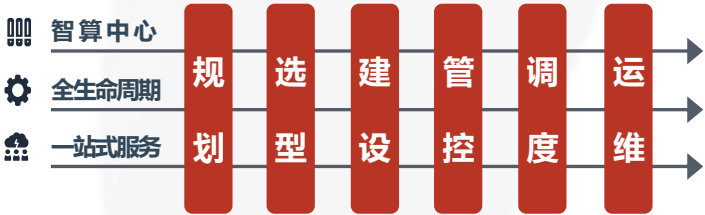
智算中心超互联解决方案Helios，产品适用于人工智能异构算力场景，大模型的训练、精调、推理等场景。具备低时延、大带宽、无阻塞，优化的算力通讯设计，支持异构融合，可规模扩展和可运维等能力，同时具有开放融合兼容自控的生态特点。

异构融合

规模扩展

高性能无阻塞端网Fabric

GPU-RDMA智能控制器



## Venus 基流长庚

Venus端网一体监控运维平台采用异构分布式计算框架，提供数据驱动的全方位、一体化的智算中心网络的智能运营、维护和监控。

网络拓扑发现、校验

网络集成设施的集中式管理

自动化配置

超可视化监控

错误告警

主机端侧融合网络功能





通信库定制

根据集群情况，定制化开发底层通信库，可实现：

- (1) 解决机型异构混跑，使不同机型服务器联合混跑。
- (2) 提升通信性能，相较于原生通信库提升10%以上的性能，具体数值根据集群规模及情况而定。

集群初始化调优服务

根据集群情况，对集群进行初始化调优，实现：

- (1) 服务器调优，调整服务器至最佳可用状态，提升机内计算效率。
- (2) 网络调优，对交换机、网卡进行深度调优，实现流量均衡及负载均衡，提升机间通信效率。

1 某2千卡RoCE集群（H100）NCCL性能测试：Reduce-scatter—8G size

NCCL原始代码优化后，性能极大提升

机型	机器数	卡数	NCCL原始代码	VCCL（NCCL原始代码优化）
H100	4	32	354.46	366.28
H100	8	64	323.60	362.41
H100	16	128	174.84	361.78
H100	32	512	78.77	352.57
H100	128	1024	92.12	333.20
H100	190	1520	59.7	310.75

2 某4千卡RoCE集群（H100）TFLOPs性能测试：Megatron框架，使用130b模型

```
iteration 14/ 1192 | consumed samples: 28672 | elapsed time per iteration (ms): 10800.4 | learning rate: 1.500E-04 | global batch size: 2048 | lm
loss: 6.281425E+00 | loss scale: 1.0 | grad norm: 6.797 | number of skipped iterations: 0 | number of nan iterations: 0 | samples per second: 189.623 | TFLOPs
(Hardware): 588.16 | TFLOPs (Model): 588.16 |
(min, max) time across ranks (ms):
forward-backward .....: (9865.61, 10499.17)
forward-compute .....: (2717.14, 3328.97)
backward-compute .....: (5667.03, 6097.76)
batch-generator .....: (0.90, 51.95)
forward-recv .....: (54.04, 375.25)
forward-send .....: (0.56, 51.54)
backward-recv .....: (56.68, 401.05)
backward-send .....: (0.59, 6.36)
forward-send-backward-recv .....: (681.92, 1437.09)
backward-send-forward-recv .....: (43.29, 249.83)
layernorm-grads-all-reduce .....: (0.69, 1.68)
embedding-grads-all-reduce .....: (0.02, 0.11)
all-grads-sync .....: (164.41, 171.80)
params-all-gather .....: (88.07, 91.94)
optimizer-copy-to-main-grad .....: (0.04, 0.19)
optimizer-clip-main-grad .....: (1.17, 1.26)
optimizer-count-zeros .....: (0.01, 0.02)
optimizer-inner-step .....: (1.38, 2.07)
optimizer-copy-main-to-model-params .....: (0.38, 0.44)
optimizer .....: (93.21, 97.08)
```

128机， TFLOPs=588.16

```
iteration 9/ 596 | consumed samples: 36864 | elapsed time per iteration (ms): 10936.9 | learning rate: 1.500E-04 | global batch size: 4096 | lm
loss: 1.370414E+01 | loss scale: 1.0 | grad norm: 74.490 | number of skipped iterations: 0 | number of nan iterations: 0 | samples per second: 374.511 | TFLOPs
(Hardware): 580.81 | TFLOPs (Model): 580.81 |
(min, max) time across ranks (ms):
forward-backward .....: (9940.95, 10588.69)
forward-compute .....: (2715.90, 3451.06)
backward-compute .....: (5635.09, 6151.43)
batch-generator .....: (0.91, 139.48)
forward-recv .....: (54.37, 390.59)
forward-send .....: (0.62, 63.22)
backward-recv .....: (53.26, 393.48)
backward-send .....: (0.65, 37.32)
forward-send-backward-recv .....: (543.51, 1388.19)
backward-send-forward-recv .....: (44.99, 356.16)
layernorm-grads-all-reduce .....: (0.75, 2.85)
embedding-grads-all-reduce .....: (0.02, 0.35)
all-grads-sync .....: (167.19, 179.28)
params-all-gather .....: (89.62, 96.36)
optimizer-copy-to-main-grad .....: (0.04, 0.19)
optimizer-clip-main-grad .....: (1.02, 1.09)
optimizer-count-zeros .....: (0.01, 0.03)
optimizer-inner-step .....: (0.73, 1.50)
optimizer-copy-main-to-model-params .....: (0.21, 0.28)
optimizer .....: (100.57, 107.38)
```

256机， TFLOPs=580.81

运维采用多种复合方式监控集群健康情况：**Venus 端网一体运维监控平台+邮箱+运维机器人**

对集群“健康”状态进行7\*24小时的持续监控并及时告警：

- (1) 服务器及GPU监控，监控GPU温度，占用，掉卡情况等。
- (2) 网络硬件监控，监控网络链路（交换机、光模块、线缆）等波动及流量情况。
- (3) 软件及流量监控，监控全集群流量。

Venus 端网一体运维监控平台



邮箱



运维机器人



## 组网案例

# MARKET CASE

### 工程落地：

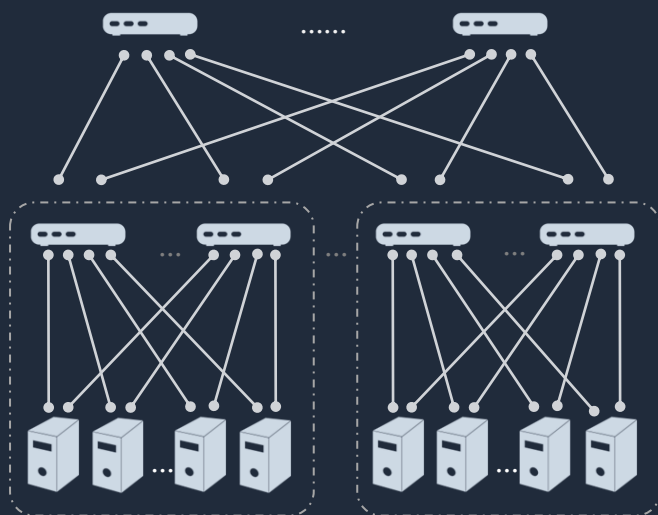
累积建设调优十余个大型智算集群

超1.7万张英伟达H800/H100经验

国产千卡集群跑通千亿参数大模型

大模型训练的GPU利用率提升15%

集群任务稳定运行+资源高效利用+异构开放生态



4000卡落地方案



- 2023.03  
大模型公司千卡集群网络建设
- 2023.05  
无锡超算项目千卡网络建设
- 2023.07  
世纪互联异构互联实验室
- 2023.08  
IDC千卡异构互连网络建设
- 2023.12  
西安智算中心千卡集群网络建设
- 2024.01  
大模型公司两千卡集群建设
- 2024.01  
上海临港集团千卡集群建设
- 2024.02  
大模型公司RoCE方案八千卡集群建设
- 2024.03  
大模型公司IB方案八千卡集群建设