

完全公开

鸿云智算大模型一体产机 产品规格书 (AWS-P2&P4)

2025 年 1 月 22 日 首次发布 深圳市诺亚鸿云信息技术有限公司

● 版权声明

本文中出现的任何文字叙述、文档格式、插图、图片、方法、过程等内容,除另有特别注明, 版 权均归**深圳市诺亚鸿云信息技术有限公司**所有,受到有关产权及版权法保护。任何个人、机构未 经**深圳市诺亚鸿云信息技术有限公司**的书面授权许可,不得以任何方式复制或引用本文的任何片段。

修订记录

版本	状态	修订理由和内容摘要	修订人	批准人	修订日期
v1. 0	С	创建《鸿云智算大模型一体产机	徐英		
		(AWS-P2&P4)			
		产品规格书》			

状态: C-创建, A-增加, M-修改, D-删除

地址: 深圳市龙岗区坂田街道岗头社区天安云谷产业园二期 4 栋 8 层 809 号

网址: www.nyhy-cloud.com

AGC 架构 AI 工作站

AWS-P2&P4(2/4 代表不同产品类型最大支持的 GPU 算力卡数量)





AWS-P4

简介

AWS-P2/P4 产品采用 AGC(AI computer system with the GPU as its Core)架构设计,秉承无主板的设计理念,支持开放的模型生态系统。其内置底层算力、预置DeepSeek模型服务、推理环境以及 RAG 环境的生产级解决方案,致力于为企业客户提供"数据不出域、性能更高效"的本地化开箱即用的 AI 服务新模式。

•

亮点



架构领先

- 闻 国内首创专为 AI 工作站设计的 GPU 数据加速单元,提升数倍 AI 性能。
- ☑ 无需主板即可实现 AI 工作站模块化生产,成为国内体积最小的 AI 工作站。
- 🔯 超静音设计,低于 60db,AWS-P2/P4支持桌面级或数据中心级部署。



算力优化

- 到 针对多样化的 GPU 算力卡以及不同大小的模型,AWS-P2&P4 融合了从 GPU 算力卡规划、推理引擎优化至负载均衡调度优化的全方位链路优化策略,客户可使用高性价比的大模型推理服务。
- 自主研发趋优推理算法提升 GPU 算力卡推理性能,相较于未采用该算法之前,实现 2~3 倍的推理效率。
- 最大支持 PCIE 5.0 协议, x16 双向带宽达到了 128GB/s。

地址:深圳市龙岗区坂田街道岗头社区天安云谷产业园二期 4 栋 8 层 809 号



数据安全

- 全离线环境部署和使用,用户既不关注公有接口,也不担心私有数据泄露,同时也不会有 Token 限制焦虑。
- 预置 DeepSeek 模型文件和所有运行环境,在数分钟内通过图形化界面启动模型并使用。
- 或持开放的模型生态系统,实现数据本地化、模型私有化。



配置灵活

- 最大可支持 2/4 张 GPU 算力卡,兼容全高全长双宽、全高全长单宽以及半高半长单宽等多种规格。
- ☑ GPU 算力卡免驱动应用,实现模型的自动适应与适配。
- ☑ 支持多类型数据接入,U盘、NFS、CIFS、S3等接入类型。

•

规格

产品型号		AWS-P2	AWS-P4	
形态		L (423mm) *W (290mm) *H (150mm)	L (434mm) *W (444mm) *H(160mm)	
模型		最大支持 DeepSeep 32B(含)以下模型	最大支持 DeepSeep 70B(含)以下模型	
		及其他 AI 大模型	及其他 AI 大模型	
控制模组	处理器	工作站数据加速单元集成	工作站数据加速单元集成	
	内存	可选配 32GB、64GB 以及 128GB 规格	可选配 32GB、64GB 以及 128GB 规格	
	系统盘	配置 1 块 M.2 512GB NVME 硬盘	配置 1 块 M. 2 512GB NVME 硬盘	
	缓存盘	可选配 1 块 M.2 960GB NVME 硬盘	可选配 1 块 M. 2 960GB NVME 硬盘	
	数据盘	最大可选配 4 块 2.5 英寸SATA 接口的固	最大可选配 8 块 2.5 英寸SATA 接口的固	
		态硬盘,提供 960GB、1.92TB、3.84TB	态硬盘,提供 960GB、1.92TB、3.84TB	
		以及 7.68TB 等多种容量规格供选择。	以及 7.68TB 等多种容量规格供选择。	
算力 模组	扩展槽	2 * PCle5.0 x16	4*PCle4.0 x16	
	GPU	最大支持 2 张全高全长双宽、全高全长	最大支持 4 张全高全长双宽、全高全长单	
	算力卡	单宽以及半高半长单宽等多种规格。	宽以及半高半长单宽等多种规格。	
网络	网络 1	配置 2 个 10Gbs 光口	配置 2 口 10Gbs 光口	
	网络 2	配置 1 个 1Gbps 电口	配置 1 个 1Gbps 电口	
监控	前面板	监控 AI 工作站的 IP 地址配置、硬盘存储	的 IP 地址配置、硬盘存储 监控 AI 工作站的 IP 地址配置、硬盘存储	

地址: 深圳市龙岗区坂田街道岗头社区天安云谷产业园二期 4 栋 8 层 809 号

网址: www.nyhy-cloud.com

产品型号		AWS-P2	AWS-P4	
	显示屏	容量、GPU 算力卡工作状态和温度。	容量、GPU 算力卡工作状态和温度。	
电源	电源	配置 1 个 1600W 热插拔电源。	配置 2 个 1600W 热插拔电源。	
其他	管理	自主研发 AIOS 操作系统,快速构建推理	自主研发 AIOS 操作系统,快速构建推理	
		环境与 RAG 环境,配置高效 VLLM 推理	环境与 RAG 环境,配置高效 VLLM 推理	
		框架(基于开源自主研发)。	框架(基于开源自主研发)。	
	I/O	2*USB 3.0,1*HDMI	2*USB 3.0,1*HDMI	
	工作	5' C~35' C(41' F~95' F),符合 ASHRAE	5' C~35' C(41' F~95' F),符合 ASHRA	
	温度	Class A1/A2	Class A1/A2	

地址: 深圳市龙岗区坂田街道岗头社区天安云谷产业园二期 4 栋 8 层 809 号