

AI算力跟踪深度（三）：
从英伟达的视角看算力互连板块成长性——
Scale Up 网络的“Scaling Law”存在吗？

证券分析师：张良卫

执业证书编号：S0600516070001

联系邮箱：zhanglw@dwzq.com.cn

联系电话：021-60199793

证券分析师：李博韦

执业证书编号：S0600123070070

联系邮箱：libw@dwzq.com.cn

2025年8月20日

我们认为Scale Up网络存在Scaling Law，Scale Up柜间第二层网络会逐渐出现，光+AEC连接多出与芯片1:9的配比需求，交换机多出与芯片4:1的配比需求，相较Scale Out网络均倍增：

1.英伟达持续扩大Scale Up规模：英伟达正通过两大路径持续扩大Scale Up网络规模。2) 提升单卡带宽：NVLink持续迭代，NVLink 5.0单卡带宽达7200Gb/s；2) 扩大超节点规模：Scale Up超节点规模不断扩大，从H100 NVL8到GH200再到GB200等，NVL72等机柜方案可以提高训推效率，但并不是Scale Up的上限，**NVL72等机柜后续会作为最小的节点（Node）存在，像积木一样在柜与柜之间进一步拼出更大的Scale Up超节点，届时需要光连接等进行通信。**

2.为什么需要Scale Up网络：“内存墙”问题和AI计算范式演进推动Scale Up网络升级。“内存墙”：单一大模型的参数量与单卡显存的差距（即模型内存墙）、单卡算力与单卡显存间的差距（即算力内存墙）均逐代放大，通过Scale Up将显存池化。计算范式：为了提升计算效率，在进行数据并行、流水线并行的同时也采用张量并行与专家并行，后者对通信频次、容量的要求都跨越数量级。

3.为什么需要更大的Scale Up网络：TCO、用户体验、模型能力拓展。随着单用户每秒消耗的Token数（Tokens Per Second, TPS）提高，包括NVL72在内的现有服务器单卡性能都会逐渐坍塌，在用户体验持续提升、模型能力拓展的趋势下，单用户TPS必然增长，采用更大规模的Scale Out能提高单卡有效性能，TCO也更具经济性。我们认为Scale Up规模与预期单用户TPS、单卡实际性能间存在Scaling Law，前者会随后者非线性增长。

4.怎么组建更大的Scale Up网络：网络结构层面，在柜间搭建第二层Scale Up交换机；端口连接层面，光与AEC有望在第二层网络中并存，**按照最新的NVLink与IB标准测算，1颗GPU需要9个额外的等效1.6T连接，为Scale Out网络的3-4.5倍，每4颗GPU需要额外1台交换机，为Scale Out网络的7.5-12倍。**

投资建议：我们认为Scale Up需求有望持续拓展，带来倍增的网络连接需求，光连接、AEC、交换机等环节都有望深度受益，相关标的——**光互连：**中际旭创，新易盛，天孚通信，光库科技，长芯博创，仕佳光子，源杰科技，长光华芯，太辰光；**铜互连：**中际旭创，兆龙互连；**交换机：**锐捷网络，盛科通信，Aster Labs（美股，后同），博通，天弘科技，Arista

风险提示：算力互连需求不及预期；客户处份额不及预期；产品研发落地不及预期；行业竞争加剧。

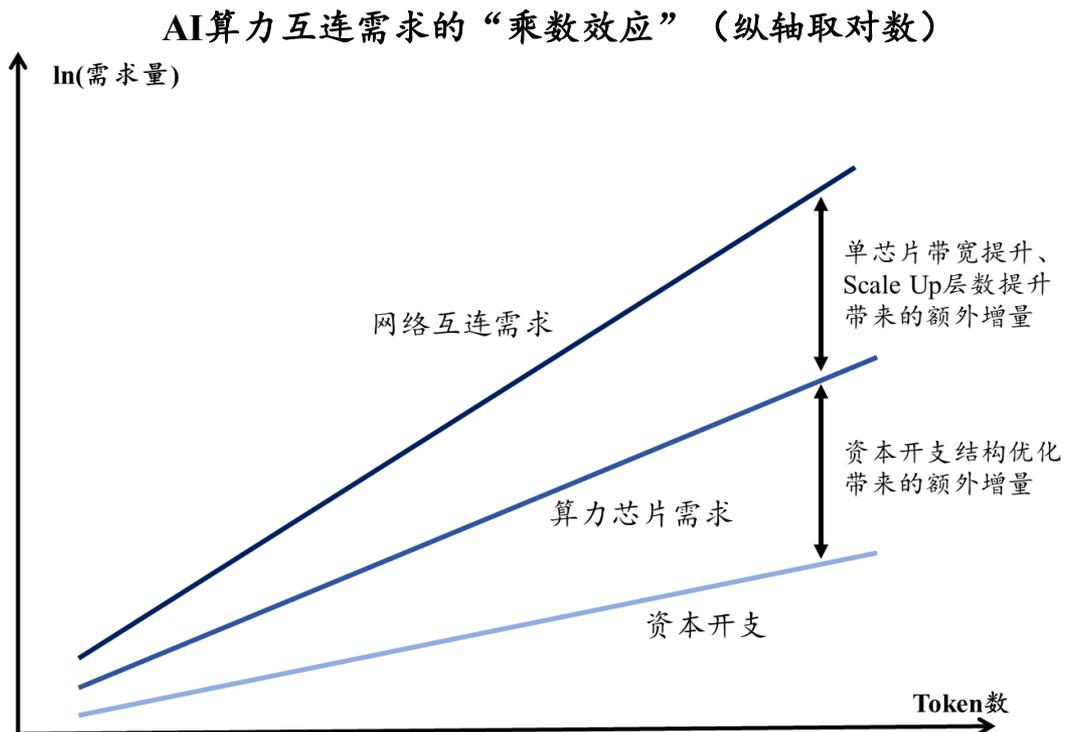
我们认为后续算力互连需求发展存在乘数效应：

• 资本开支结构优化，算力芯片增长速度高于资本开支增速

■ 单芯片带宽提升，算力互连需求增速高于芯片需求增速

■ 芯片需求 $\uparrow\uparrow$ = CapEx \uparrow × 算力芯片投资在CapEx占比 \uparrow × 芯片投资性价 \uparrow

■ 算力互连需求 $\uparrow\uparrow\uparrow$ = 芯片需求 $\uparrow\uparrow$ × 单芯片带宽 \uparrow





■ 英伟达持续扩大Scale Up规模

■ 为什么需要Scale Up网络

■ 为什么需要更大的Scale Up网络

■ 怎么组建更大的Scale Up网络

■ 投资建议及风险提示

1. 英伟达持续扩大Scale Up规模

1.1 英伟达持续尝试扩大Scale Up规模

- 英伟达从单卡带宽与超节点规模两个路径升级Scale Up;
- NVLink跟随每一代GPU架构进行升级，目前最新用于B系列GPU的NVLink 5.0可支持单卡7.2Tb的带宽，相较用于H100的NVLink 4.0带宽翻倍;
- Scale Up超节点规模在H100之后经历了GH200、GB200等方案，从NVL8拓展至NVL72甚至更高，这个扩展路径是复杂但必需的。

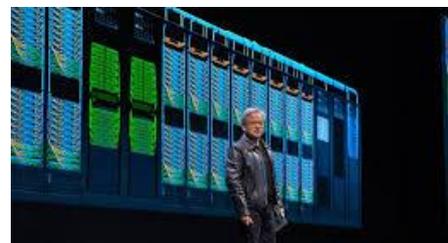
VR NVL144(72GPU) NVL576(144GPU)



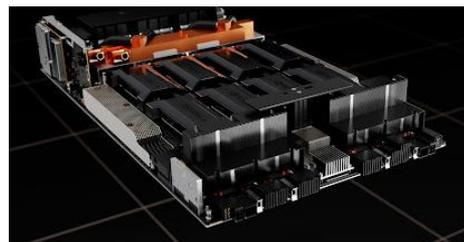
GB系列 NVL72



GH200 NVL256



H100 NVL8



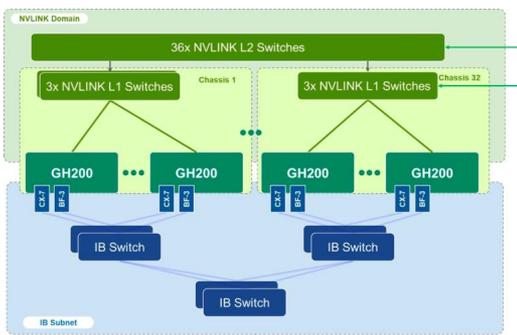
历代NVLink

版本	年份	GPU 架构	每链接带宽 (Gb/s)	链接数	总带宽 (Gb/s)
1.0	2016	Pascal	160	4	640
2.0	2017	Volta	200	6	1200
3.0	2020	Ampere	200	12	2400
4.0	2022	Hopper	400	18	3600
5.0	2024	Blackwell	800	18	7200

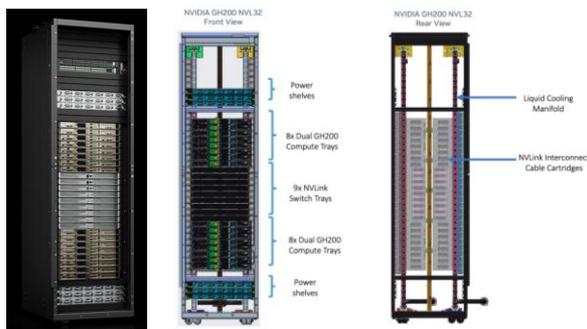
1.2 H100 NVL8到GH200 NVL 256：前瞻但过于激进的一步

- 英伟达在2023年基于H200发布了GH200 NVL256超节点，后者由32个计算Chassis组成，每个Chassis由8张GH200组成；
- Chassis内8张GH200通过L1 NVSwitch连接，32个Chassis间通过L2 NVSwitch连接；
- L2 NVSwitch通过光连接，每张GPU配套8个800G光模块，大约每7张GPU对应一台L2 NVSwitch；
- 单张GPU配套Scale Up的通信硬件成本较高与GPU为同一数量级，且训练、推理性能提升尚不明显，GH200 NVL 256未实现大范围推广，英伟达后续推出成本更低的GB200 NVL72的前身GH200 NVL32。

初代GH200 NVL256网络拓扑图



初代GH200 NVL32机柜



8台GH200 NVL32拓展为NVL256



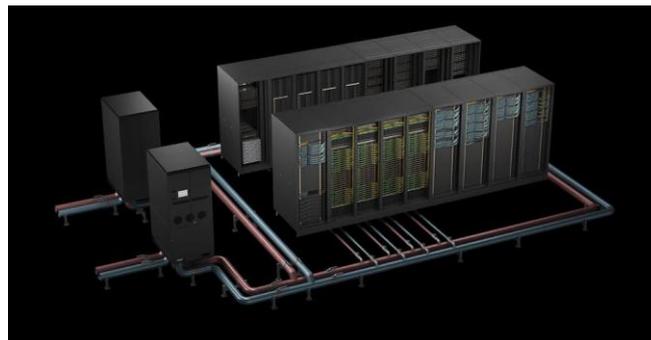
1.3 GB与VR机柜：有效但并非Scale Up最终形态

- GB与VR的机柜方案已经讨论了很多，这里主要阐述我们对这类机柜产品的判断：
 - 机柜方案延续了英伟达在GH200 NVL256上的思路，即除了提升NVLink带宽外，还要提高Scale Up超节点的规模，升级为机柜方案是为了增加GPU密度，节省物理空间的同时缩小GPU间连接距离，以使用相比于光连接成本更低的PCB、铜连接；
 - 铜连接、PCB、液冷、电源等都随着GPU密度提高实现单张GPU对应价值量的跃升；
 - 机柜方案实现的NVL72、NVL144等Scale Up确实可以提高训练、推理效率，但并不是英伟达Scale Up的上限，**NVL72、NVL144等机柜方案后续会作为最小的Scale Up节点（Node）存在，像积木一样在柜与柜之间进一步拼出更大的Scale Up超节点，届时需要光连接等进行通信。**可具体参考后续章节对Scale Up需求的底层逻辑以及趋势的分析。

GB200 NVL72网络拓扑图



8台GB200 NVL72机柜



2. 为什么需要Scale Up网络

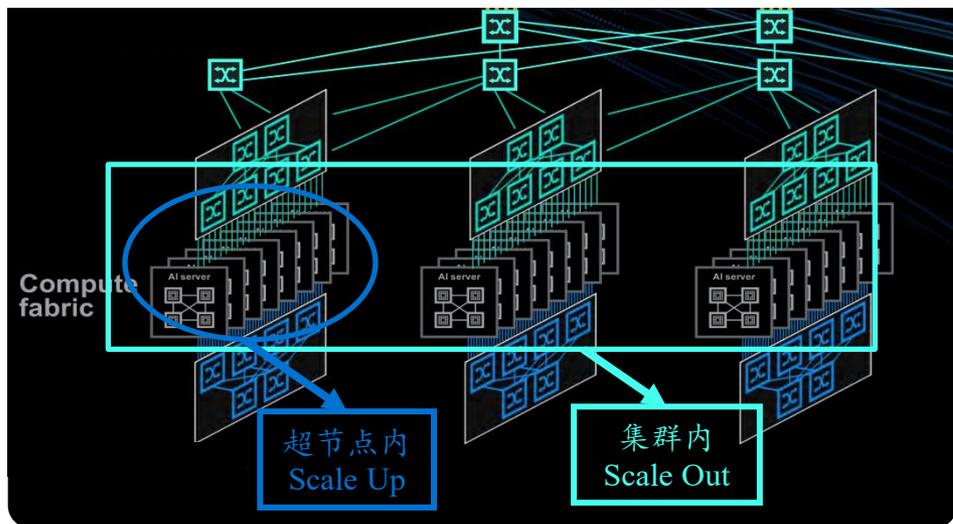
2.1 Scale Up与Scale out的特点与作用各不相同

若干超节点（SuperPod，如NVL 72）组成集群（Cluster，如万卡、十万卡集群）；

- Scale Out网络实现集群内所有GPU卡互联，亮点在于网络内连接GPU数量大，与传统数据中心网络类似；
- Scale Up网络实现超节点内所有GPU卡互联，亮点在于网络内单卡通信带宽高，组网规模尚小，为AI算力场景下新兴的网络架构；
- **Scale Up并不仅限于柜内，柜外也可进行Scale Up。**

（由于篇幅有限本文未就技术原理做详细阐述，具体细节欢迎进一步交流）

Scale Out网络与Scale Up网络



Scale Out与Scale Up网络对比

（一层NVLink交换机+CX-8网卡+三层Quantum-X800 IB网络）

	最大GPU数（张）	单卡带宽（Gb/s）
Scale Out	746496	800
Scale Up	72	7200

2.2 “内存墙”问题需要Scale Up网络将显存池化来缓解

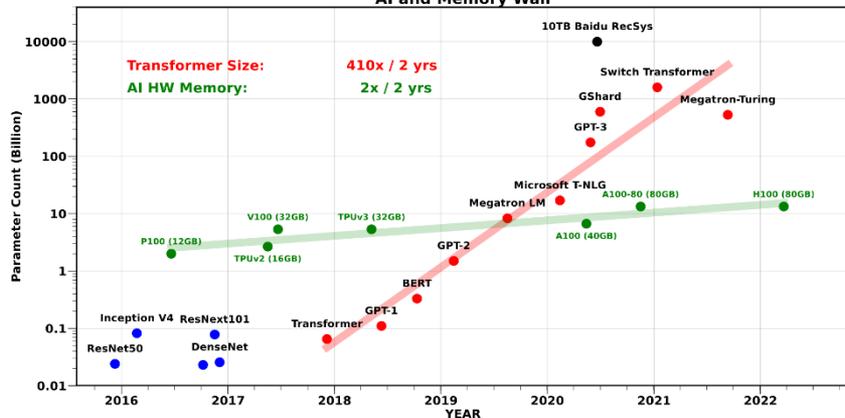
训推计算的“内存墙”催生出通过Scale Up网络将显存池化的需求：

- 单一大模型的参数量与单卡显存的差距（即模型内存墙）、单卡算力与单卡显存间的差距（即算力内存墙）均逐代放大
- 除模型参数外，推理计算生成的KV Cache（关键中间值的缓存，用于简化计算）占用显存大小也可达模型的50%甚至以上
- 因此单卡运算时需从多张卡的显存读取所需参数、数据，为了尽可能减少数据传输时延，目前产业化应用最优解是使用Scale Up网络将显存池化，如NVL72。

（由于篇幅有限本文未就技术原理做详细阐述，具体细节欢迎进一步交流）

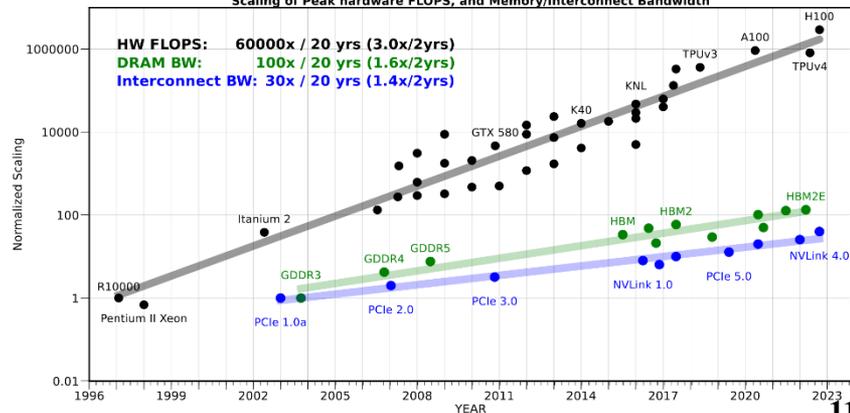
模型内存墙逐代放大

AI and Memory Wall



算力内存墙逐代放大

Scaling of Peak hardware FLOPS, and Memory/Interconnect Bandwidth

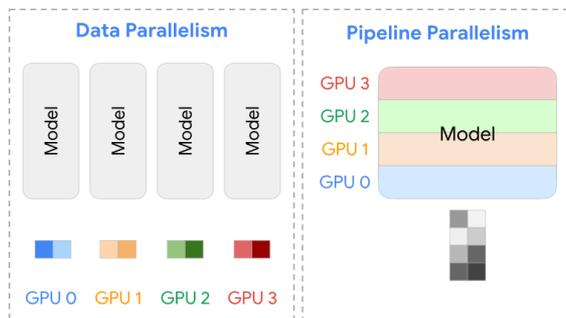


2.3 AI训推计算范式推动Scale Up升级、单卡带宽提升

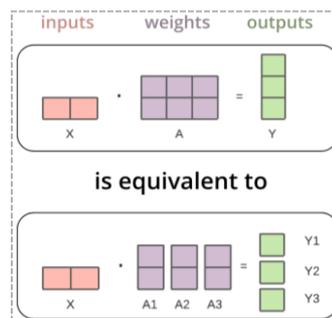
- AI训推需要分布式并行计算，基于对计算效率不断提升的追求，并行计算方式有数据并行（Data Parallelism）、流水线并行（Pipeline Parallelism）、专家并行（MoE Parallelism）及张量并行（Tensor Parallelism）。
- 数据并行：**将输入数据分配给各个负载，各负载上基于不同数据进行同一模型的训练/推理；
- 流水线并行：**将模型分为若干层分配给各个负载，各负载分别进行不同层的计算；
- 张量并行：**将模型参数运算的矩阵拆分为子矩阵传输至各个负载，各负载分别进行不同的矩阵运算

（由于篇幅有限本文未就技术原理做详细阐述，具体细节欢迎进一步交流）

数据并行（左），流水线并行（右）计算原理



张量并行计算原理

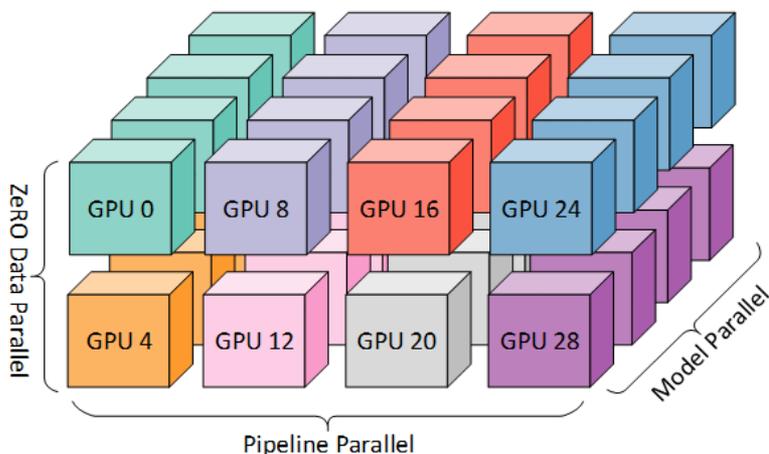


2.3 张量并行可优化计算效率

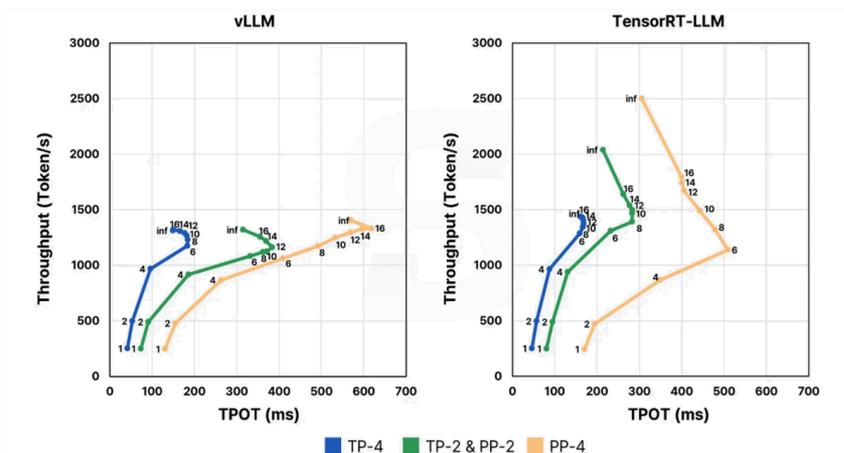
- 目前模型训推主要采用混合并行，即多种并行方式同时进行，可从不同维度切分/编组进行并行
- 张量并行、专家并行是粒度更细的并行方式，更高效利用单张芯片配套内存，因此可以明显提升计算效率。

(由于篇幅有限本文未就技术原理做详细阐述，具体细节欢迎进一步交流)

3D混合并行计算



在两类推理引擎中张量并行 (TP) 均可缩短输出单Token所需时间 (TPOT)



2.3 张量并行对通信的高要求催生Scale Up需求

- 张量并行在每一层神经网络的计算后都需要将新的计算结果收集、汇总，并将完整结果重新分发即Allreduce通信，因此在训推时对通信频率、传输容量都有更高要求。
- 需要用Scale Up满足越来越高的通信频率、传输容量需求。

(由于篇幅有限本文未就技术原理做详细阐述，具体细节欢迎进一步交流)

推理中张量并行 (TP) 的通信次数与规模均比流水线并行 (PP) 高出数量级

Model	TP×PP	Prefill Stage			Decode Stage		
		Operation	Count	Shape	Operation	Count	Shape
Llama-3.1-8B $S_p = 128$ $S_d = 128$	2×2	Allreduce	33	[128,4096]	Allreduce	4191	[1,4096]
		Gather	1	[64128]	Gather	127	[64128]
		Allgather	2	[128,4096]	Allgather	254	[1,4096]
		Send/Recv	2	[128,2048]	Send/Recv	254	[1,2048]

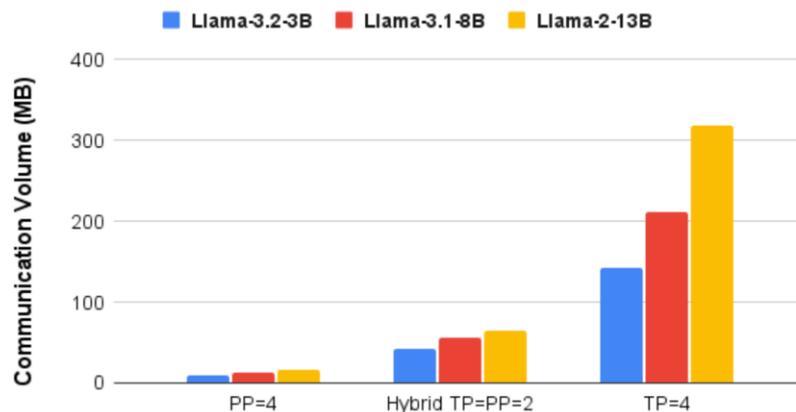
训练中张量并行需要传输的数据量多出一到两个数量级
(GPT-3B模型基于32个GPU训练数据)

Traffic type	Volume	Number of messages	Message size
TP	~85 GB	680	125 MB
PP	~1 GB	16	125 MB
DP	741 MB	1	741 MB
EmbTableSyn	96 MB	1	96 MB

各类并行通信方式对比

切分方式	通信操作	每次迭代单卡通信量	对网络的需求
张量并行 (TP)	AllReduce	百GB级别	(超) 节点内高速互联
专家并行 (EP)	All-to-All	百GB级别	(超) 节点内高速互联
流水并行 (PP)	Send/Recv	MB级别	节点间高速互联
数据并行 (PP)	AllReduce	GB级别	节点间高速互联

推理中张量并行规模越大通信量越大 ()



3. 为什么需要更大的Scale Up网络

3.1 Scale Up可加速推理，且增益随推理负载提升而扩大

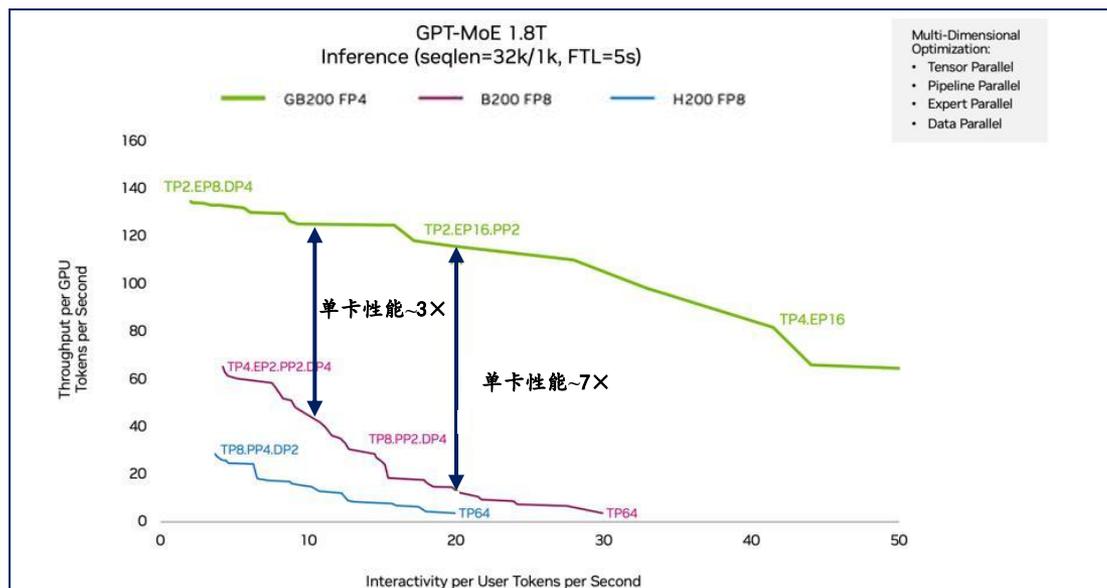
- 我们认为Scale Up规模越大，集群算力有效利用率往往越高，且随着单用户推理负载增加，增益会越来越大，这里以GB200 NVL72、B200 NVL8的对比为例（可见下一页图片）
- **测试配置：**各类方案都是基于33000张GPU的Scale Out集群进行测试，GB200 NVL72采用了NVL72 Scale Up超节点、Grace CPU、FP4精度，B200 NVL8采用了NVL8 Scale Up超节点、Intel Xeon CPU、FP8精度，因此精度优化可为GB200 NVL72直接带来1倍单卡性能提升；
- **模型：**GPT MoE 1.8T模型，采用混合并行推理（最多64维并行），FTL=5s，TTL=50ms，input/output长度分别为32768 /1024；
- **坐标轴含义：**横轴代表单用户每秒收到的Token数（Tokens Per Second，TPS），亦即用户体验或模型推理的实际输出能力；纵轴代表集群内每张GPU每秒输出的Token数，亦即推理时单张卡的实际性能或有效利用程度；
- 每条曲线每点对应各单用户TPS下，所有混合并行方案及Chunk Size组合中单卡性能最大值
- 可以初步观测到**横纵坐标成反比**，主要原因为单用户TPS提升后需要在单位时间内用更多GPU输出更多Token，通信阻塞变大，GPU等待数据传输的时间增加，利用率下降。

3.1 Scale Up可加速推理，且增益随推理负载提升而扩大

- 在单用户TPS为10 Tokens/s时，GB200 NVL72的单卡实际性能约为B200 NVL8的3倍，考虑FP4精度优化带来的约1倍提升后，Scale Up+Grace CPU带来约50%的性能提升；
- 在单用户TPS为20 Tokens/s时，GB200 NVL72的单卡实际性能约为B200 NVL8的7倍，考虑FP4精度优化带来的约1倍提升后，Scale Up+Grace CPU带来约250%的性能提升；
- 我们认为随着单用户TPS增加，Scale Up带来的单卡利用率增益会越来越大。

(由于篇幅有限本文未就技术原理做详细阐述，具体细节欢迎进一步交流)

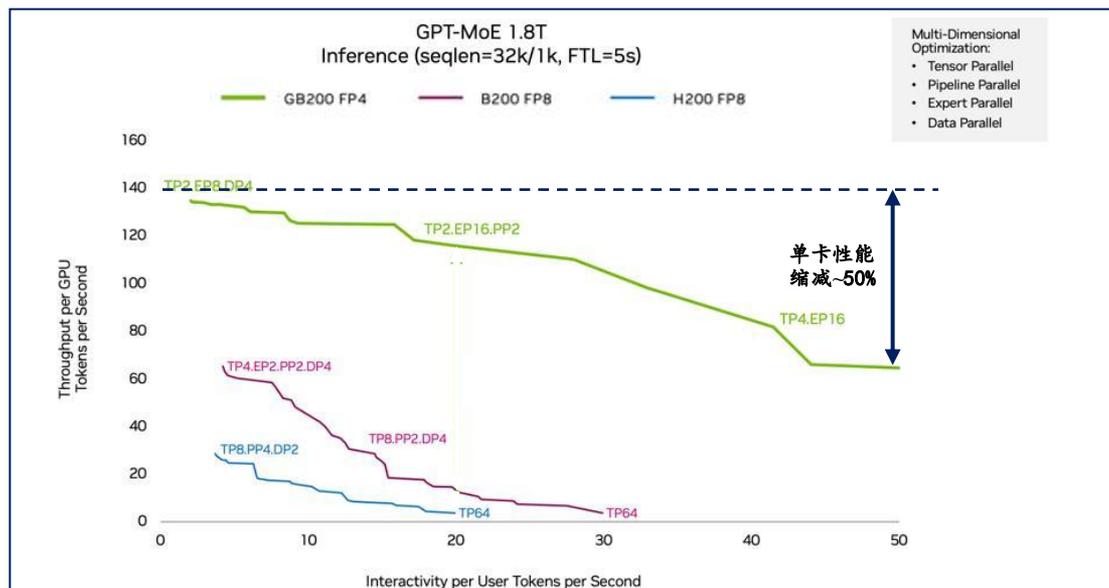
不同方案推理性能对比



3.2 NVL72、144不是推理Scale Up的上限

- 我们认为机柜对应的NVL72、NVL144等方案并不是Scale Up超节点的上限，机柜会像积木一样进一步拼出更大的超节点，这主要来自硬件TCO、用户体验、模型能力拓展三层因素；
- 当单用户TPS沿横轴提高到50 Tokens/s时，B200 NVL8、H200 NVL8的单卡真实性能已经没有实际意义，GB200 NVL72仍有70 Tokens/s的单卡TPS，但已相较最大性能缩减50%；
- 要继续提高纵轴单卡性能，我们认为除了在软件层面引入新的推理引擎，如英伟达Dynamo外，还需提升Scale Up规模，以及增加混合并行线路数（图中限制为64路并行）。

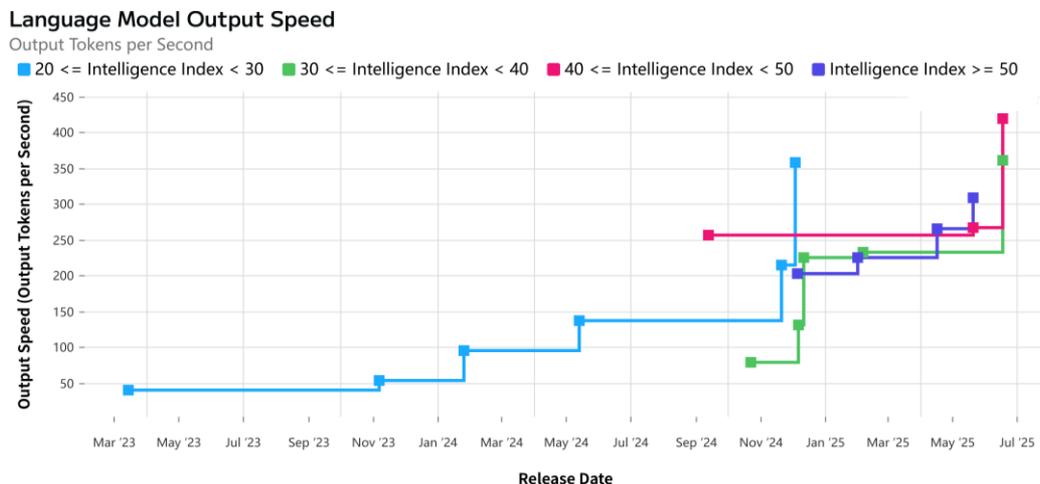
GB200 NVL72单卡性能逐渐衰减



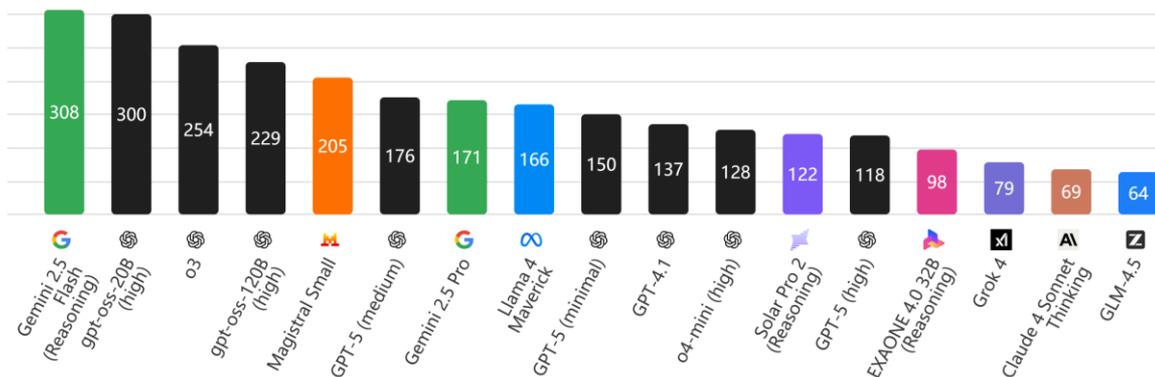
3.2 用户体验及模型能力拓展推动单用户TPS增长

- 各能力带的LLM单用户TPS均不断提升，目前前五名主流模型的单用户TPS均在200 Tokens/s以上；
- 我们认为用户使用模型过程中最直观、最核心的体验点是回答结果的生成速度，即单用户的TPS，且TPS提升后具有实际应用意义的场景会逐渐涌现，如AI coding。

各能力带的LLM单用户TPS均不断提升



前五名主流模型的单用户TPS均在200 Tokens/s以上



3.2 用户体验及模型能力拓展推动单用户TPS增长

- 模型能力从LLM扩展到多模态，；
- 我们认为用户使用模型过程中最直观、最核心的体验点是回答结果的生成速度，即单用户的TPS，且TPS提升后具有实际应用意义的场景会逐渐涌现，如AI coding。

GPT-image-1输出图像大小及耗时

Quality	Square (1024x1024)	Portrait (1024x1536)	landscape (1536x1024)
Low	272 tokens	408 tokens	400 tokens
Medium	1056 tokens	1584 tokens	1568 tokens
High	4160 tokens	6240 tokens	6208 tokens
Simple Prompt	Complex Prompt	High Detail	Peak Hours
3-8 seconds	10-20 seconds	15-25 seconds	20-35 seconds

Meta为广告主提供的AI图像功能



Image expansion



Background generation

3.3 Scale Up网络存在随用户TPS增长的“Scaling Law”

假设在前文的对比中GB200NVL72同样采用FP8精度而非FP4精度：

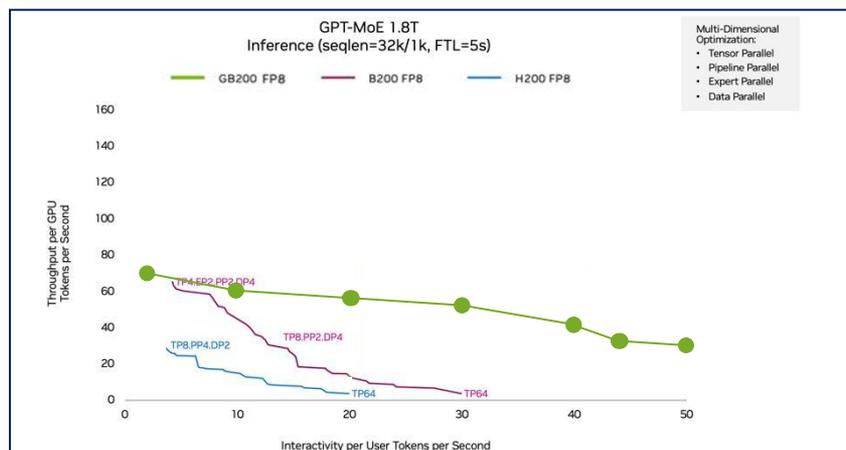
- 纵轴单卡实际性能越高，GB200在横轴上单用户TPS的绝对、相对优势均越小；
- 横轴单用户TPS越高，GB200的单卡实际性能越低；
- 随着用户体验提升、模型能力拓展，横轴将不断向右延伸，下降越慢的曲线在满足推理负载的同时能够实现更小的TCO，而搭建更大规模的Scale Up网络是缓解下降最有效的方法之一；
- 图中所需ScaleUp规模的增速高于单用户TPS及单卡性能增速，且横纵指标成反比，可以推测所需ScaleUp网络规模与单用户TPS、单卡实际性能及模型参数量间或存在Scaling Law，即：

$$S = aG^b U^c$$

其中S为Scale Up超节点规模，a为与模型参数、并行计算路数等有关的变量，G为单张GPU的实际性能（同款GPU对比），b为大于1的常数，U为单用户TPS，c为大于1的常数。

- 基于这一Scaling Law，我们认为随着推理时对单用户TPS、单GPU实际性能越来越高的要求，Scale Up规模会非线性增长。

不同方案推理性能对比估算（GB200来自原数据估测存在一定误差）

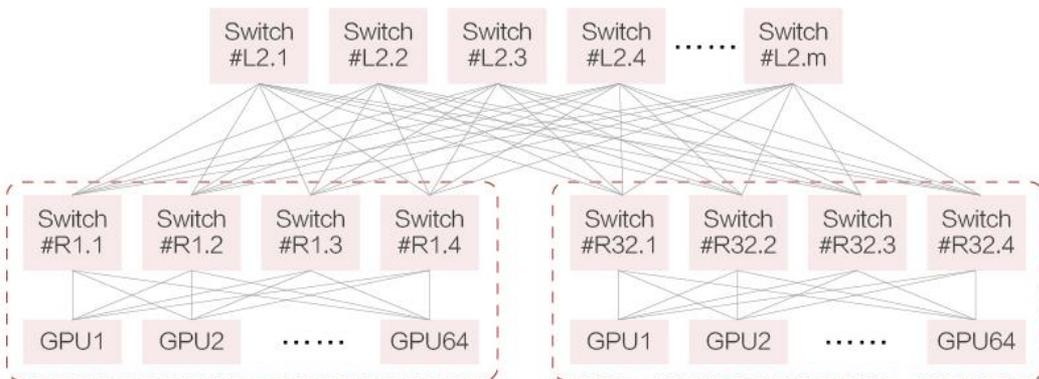


4. 怎么组建更大的Scale Up网络

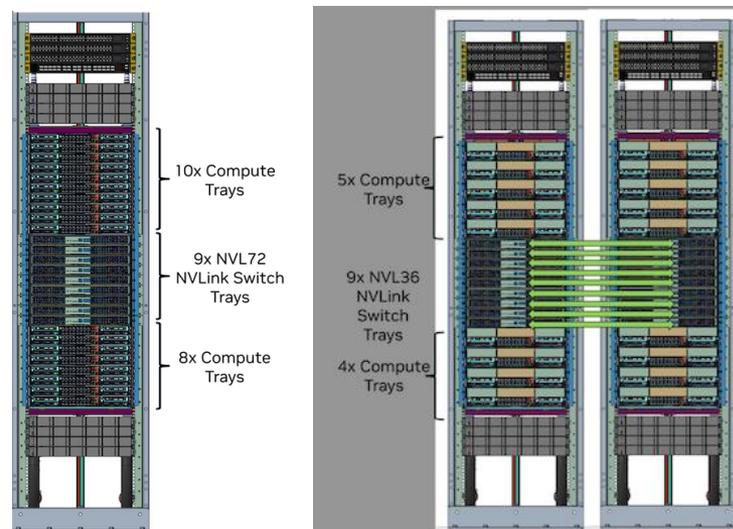
4.1 网络架构：柜外搭建第二层Scale Up交换机网络

- 英伟达的机柜中加入了一层NVSwitch，以GB200 NVL72为例，单颗B200 NVLink带宽7.2Tb（单向带宽，下同），9个Switch Tray总带宽 $57.6\text{Tb} \times 9 = 518.4\text{T}$ ，刚好与72颗B200进行无阻塞通信，这意味着如果在柜内继续增加GPU，需要同步增加配套Switch Tray，需要的物理空间和距离增加。因此我们认为在GB机柜使用铜连接，VR机柜有望增加PCB后，柜内扩展难度增加，需要增加第二层交换机做柜间Scale Up；
- 对于NVL72而言，则需要改为NVL36 \times 2以使得第一层Switch Tray翻倍至18个，以提供连接至第二层NVSwitch的上行带宽。

Scale Up两层网络拓扑（以单机柜64卡为例）



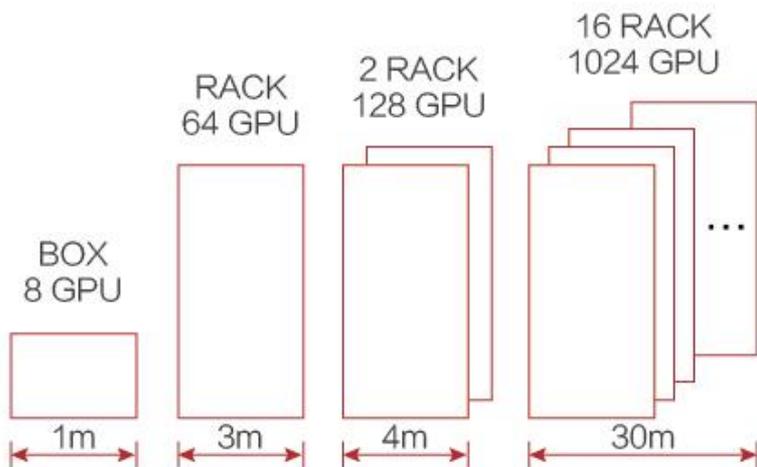
英伟达GB200 NVL36 \times 2方案



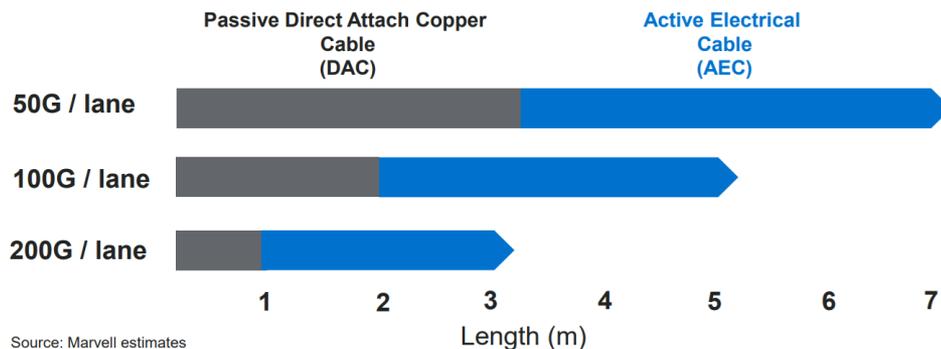
4.2 连接方式：第二层Scale Up网络中光与AEC并存

- 在单通道200G速率下，无源铜（如DAC）的有效距离上限在1m左右，因此基本无法满足跨柜Scale Up的连接需求，有源铜（如AEC）的有效距离上限在3米左右，因此可满足部分跨柜Scale Up的连接需求，光（如AOC、光模块）可满足所有跨柜Scale Up的连接距离要求；
- 我们认为“能用铜的地方就不会用光”，在第二层柜间Scale Up场景会有光与AEC并存。

Scale Up网络通信距离



AEC与DAC有效距离对比



4.2 连接方式：第二层Scale Up网络带来的网络增量需求有多少

- 按照最新的NVLink与IB标准测算，第二层Scale Up网络中1颗GPU需要9个额外的等效1.6T连接（等于第一层），每4颗GPU需要额外1台NVLink 5.0交换机；两到三层Scale Out中1颗GPU对应2-3个等效1.6T连接，每30-48颗GPU对应一台Quantum-X800 Q34xx系列交换机。
- 目前Scale Up与Scale Out并存，其最终形态是做到与Scale Out相近的规模后取代Scale out，但需要考虑到在成本与物理空间维度都数倍增长的网络连接。

英伟达Scale Up与Scale Out网络连接需求对比（均基于最新平台）

	Scale Up（第二层）	Scale Out（两到三层）
连接端口数量	1颗GPU需要额外9个等效1.6T连接	1颗GPU对应2-3个等效1.6T连接
交换机数量	每4颗GPU需要额外1台NVLink 5.0 Switch	每30-48颗GPU对应一台Quantum-X800 Q34xx系列交换机

4.2 连接方式：潜在技术路线适用于Scale Up吗？

- 我们认为CPO、OCS等潜在的新技术在Scale Up中的应用会比Scale Out更难，这些新技术在Scale Out中规模化应用后，对它们在Scale Up中应用可能性的讨论才有实际意义；
- 2.3节中我们说过Scale Up网络用来满足张量并行、专家并行等计算的通信需求，其单位时间内需要传输数据的频次与大小都是Scale Out网络的几十倍甚至上百倍，这意味着应用新技术路线的难度与故障率都会相应增加。

推理中张量并行（TP）的通信次数与规模均比流水线并行（PP）高出数量级

Model	TP×PP	Prefill Stage			Decode Stage		
		Operation	Count	Shape	Operation	Count	Shape
Llama-3.1-8B $S_p = 128$ $S_d = 128$	2×2	Allreduce	33	[128,4096]	Allreduce	4191	[1,4096]
		Gather	1	[64128]	Gather	127	[64128]
		Allgather	2	[128,4096]	Allgather	254	[1,4096]
		Send/Recv	2	[128,2048]	Send/Recv	254	[1,2048]

各类并行通信方式对比

切分方式	通信操作	每次迭代单卡通信量	对网络的需求
张量并行（TP）	AllReduce	百GB级别	（超）节点内高速互联
专家并行（EP）	All-to-All	百GB级别	（超）节点内高速互联
流水并行（PP）	Send/Recv	MB级别	节点间高速互联
数据并行（PP）	AllReduce	GB级别	节点间高速互联

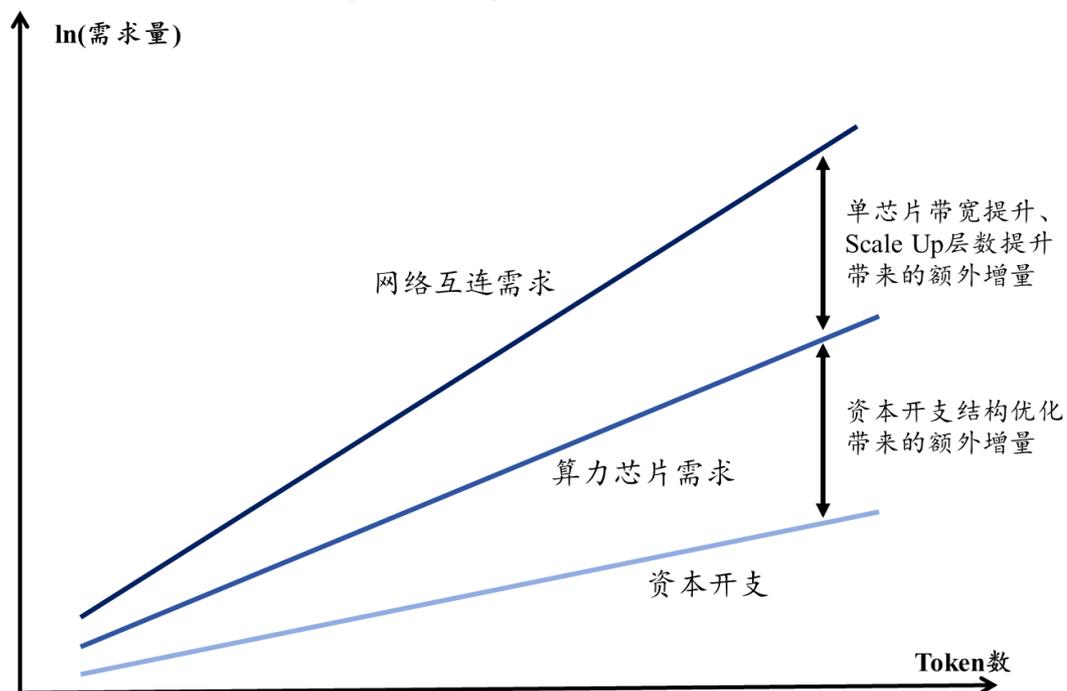
训练中张量并行需要传输的数据量多出一到两个数量级（GPT-3B模型基于32个GPU训练数据）

Traffic type	Volume	Number of messages	Message size
TP	~85 GB	680	125 MB
PP	~1 GB	16	125 MB
DP	741 MB	1	741 MB
EmbTableSyn	96 MB	1	96 MB

4.3 Scale Up需求凸显，产业链增速高于Capex&芯片增速

- 基于以上分析，我们认为后续算力互连需求发展存在乘数效应：
 - 资本开支结构优化，算力芯片增长速度高于资本开支增速
 - 单芯片带宽提升，算力互连需求增速高于芯片需求增速
 - 芯片需求 $\uparrow\uparrow$ = CapEx \uparrow \times 算力芯片投资在CapEx占比 \uparrow \times 芯片投资性价 \uparrow
 - 算力互连需求 $\uparrow\uparrow\uparrow$ = 芯片需求 $\uparrow\uparrow$ \times 单芯片带宽及Scale Up层数 \uparrow

AI算力互连需求的“乘数效应”（纵轴取对数）



5. 投资建议及风险提示

我们认为Scale Up需求有望持续拓展，带来倍增的网络连接需求，光连接、AEC、交换机等环节都有望深度受益，瓜分增量互连需求，相关标的：

- **光互连**：中际旭创，新易盛，天孚通信，光库科技，长芯博创，仕佳光子，源杰科技，长光华芯，太辰光
- **铜互连**：中际旭创，兆龙互连
- **交换机**：锐捷网络，盛科通信，Aster Labs（美股，后同），博通，天弘科技，Arista

- **算力互连需求不及预期：**若后续下游客户算力建设投入未达预期，或AI算力网络带宽规模未达预期情况，各客户对于网络互连产品的需求也将不及预期，相关公司业绩表现将受到影响；
- **客户开拓与份额不及预期：**如果相关公司未如预期开拓潜在客户，或在客户处份额低于预期，公司业绩将受到影响；
- **产品研发落地不及预期：**如果相关公司在具有潜在应用前景的产品研发及量产应用上未达预期，将对公司业绩的表现造成影响；
- **行业竞争加剧：**如果行业竞争持续加剧，相关产品份额存在下降的可能。

东吴证券股份有限公司经中国证券监督管理委员会批准，已具备证券投资咨询业务资格。

本研究报告仅供东吴证券股份有限公司（以下简称“本公司”）的客户使用。本公司不会因接收人收到本报告而视其为客户。在任何情况下，本报告中的信息或所表述的意见并不构成对任何人的投资建议，本公司及作者不对任何人因使用本报告中的内容所导致的任何后果负任何责任。任何形式的分享证券投资收益或者分担证券投资损失的书面或口头承诺均为无效。

在法律许可的情况下，东吴证券及其所属关联机构可能会持有报告中提到的公司所发行的证券并进行交易，还可能为这些公司提供投资银行服务或其他服务。

市场有风险，投资需谨慎。本报告是基于本公司分析师认为可靠且已公开的信息，本公司力求但不保证这些信息的准确性和完整性，也不保证文中观点或陈述不会发生任何变更，在不同时期，本公司可发出与本报告所载资料、意见及推测不一致的报告。

本报告的版权归本公司所有，未经书面许可，任何机构和个人不得以任何形式翻版、复制和发布。经授权刊载、转发本报告或者摘要的，应当注明出处为东吴证券研究所，并注明本报告发布人和发布日期，提示使用本报告的风险，且不得对本报告进行有悖原意的引用、删节和修改。未经授权或未按要求刊载、转发本报告的，应当承担相应的法律责任。本公司将保留向其追究法律责任的权利。

东吴证券投资评级标准

资评级基于分析师对报告发布日后6至12个月内行业或公司回报潜力相对基准表现的预期（A股市场基准为沪深300指数，香港市场基准为恒生指数，美国市场基准为标普500指数，新三板基准指数为三板成指（针对协议转让标的）或三板做市指数（针对做市转让标的），北交所基准指数为北证50指数），具体如下：

公司投资评级：

买入：预期未来6个月个股涨跌幅相对基准在15%以上；

增持：预期未来6个月个股涨跌幅相对基准介于5%与15%之间；

中性：预期未来6个月个股涨跌幅相对基准介于-5%与5%之间；

减持：预期未来6个月个股涨跌幅相对基准介于-15%与-5%之间；

卖出：预期未来6个月个股涨跌幅相对基准在-15%以下。

行业投资评级：

增持：预期未来6个月内，行业指数相对强于基准5%以上；

中性：预期未来6个月内，行业指数相对基准-5%与5%；

减持：预期未来6个月内，行业指数相对弱于基准5%以上。

我们在此提醒您，不同证券研究机构采用不同的评级术语及评级标准。我们采用的是相对评级体系，表示投资的相对比重建议。投资者买入或者卖出证券的决定应当充分考虑自身特定状况，如具体投资目的、财务状况以及特定需求等，并完整理解和使用本报告内容，不应视本报告为做出投资决策的唯一因素。

东吴证券研究所
苏州工业园区星阳街5号
邮政编码：215021

传真：（0512）62938527

公司网址：<http://www.dwzq.com.cn>

东吴证券 财富家园