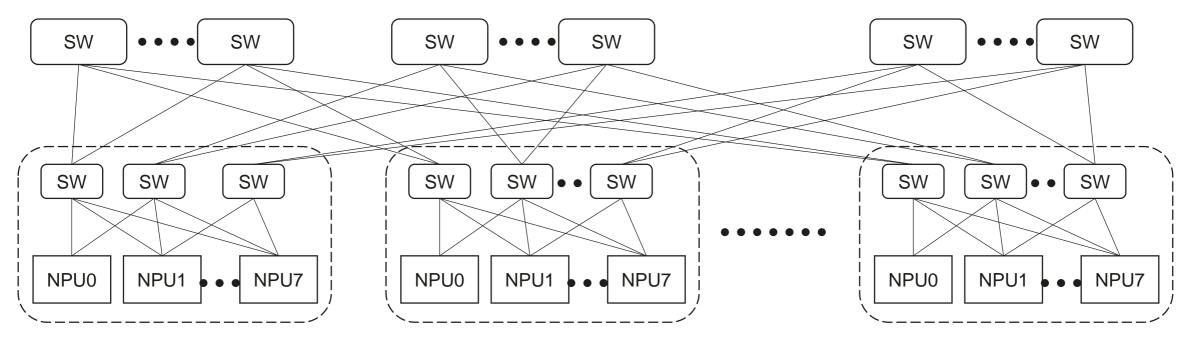
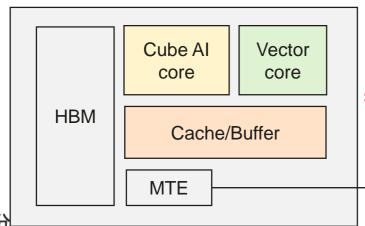






### CloudMatrix Ascend Super-computer





Distributed and Parallel Software Laboratory

Unified buffer interconnection:

scalable, memory semantics w/ unified address space, outof-order and multipath native

	NVL72	CloudMatrix
# of node	72	384
Peak (Pflops)	180	300
Mem BW(TBps)	7.9 (576)	3.2 (1229)
Int. BW(TBps)	1.8 (130)	2.8 (269)

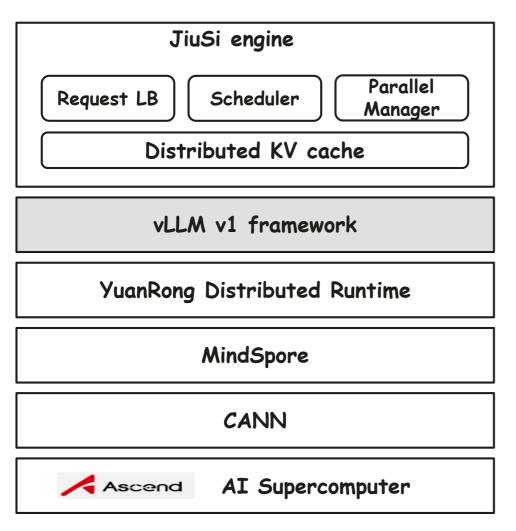




### **CloudMatrix Ascend Super-computer**



### JiuSi(九思): vLLM-based Ascend-affined inference engine



### · JiuSi (九思)

- Based on vLLM v1 architecture
  - An easy-hacking codebase
- Distributed KV Cache centric
- Enhanced request load-balancing, scheduling, and rich parallel mechanisms
- Asynchronous processing optimizations to vLLM

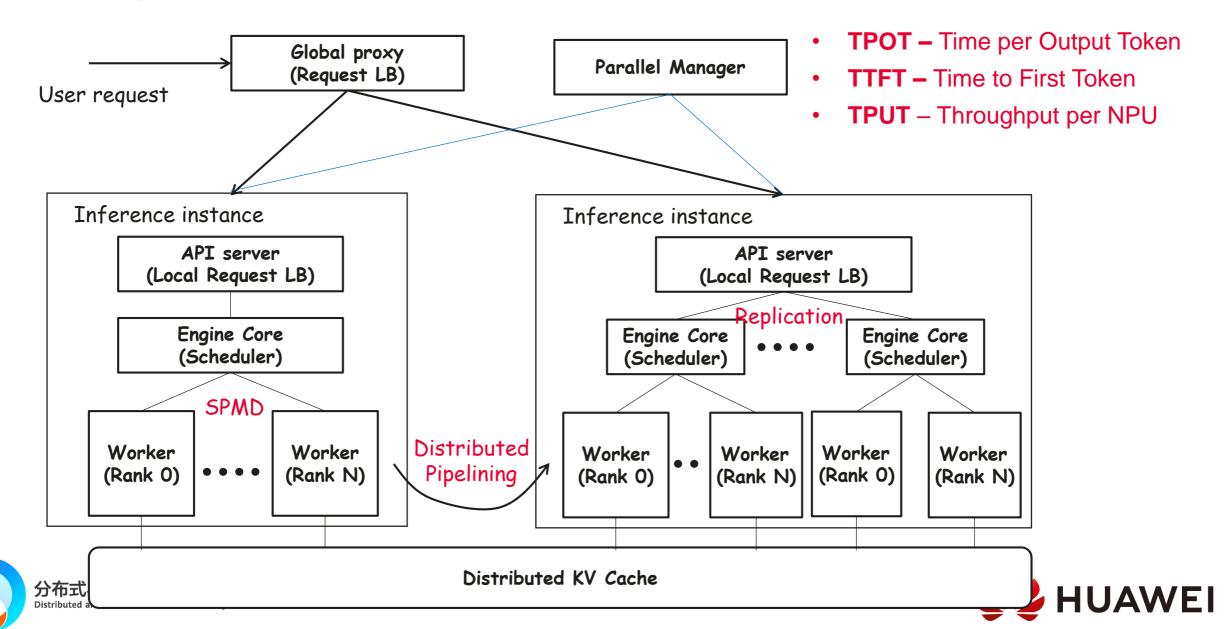
### · YuanRong(元戎)

- A serverless runtime for clusters
  - Providing elasticity for JiuSi workers
- Data system a distributed object memory
  - An ideal substrate for distributed KVC





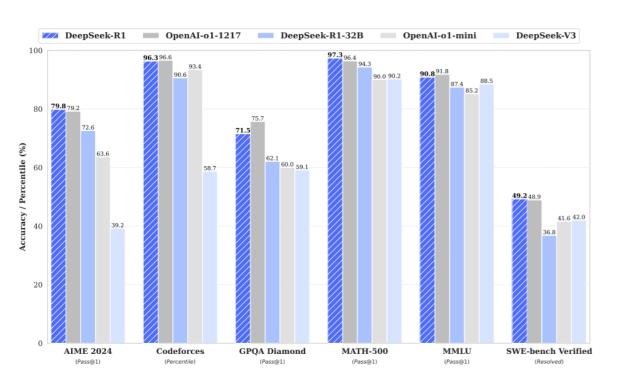
### **Serving LLM**

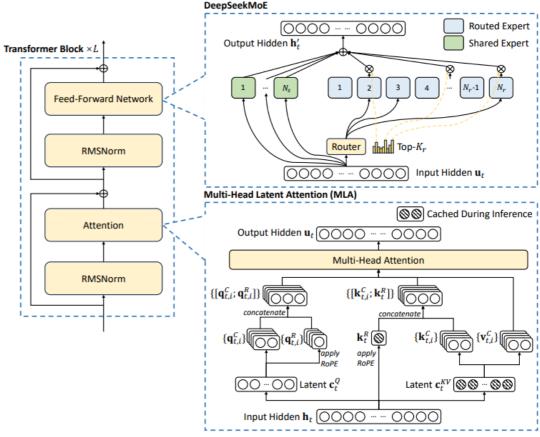


### DeepSeek V3/R1

- A game-changer open source (none-)reasoning model that is on par with SOTA proprietary LLMs
- Boosting new surge of demands on LLM inferencing and agents
- 671B model but with many innovations for efficiency and high performance

MLA, MoE, MTP



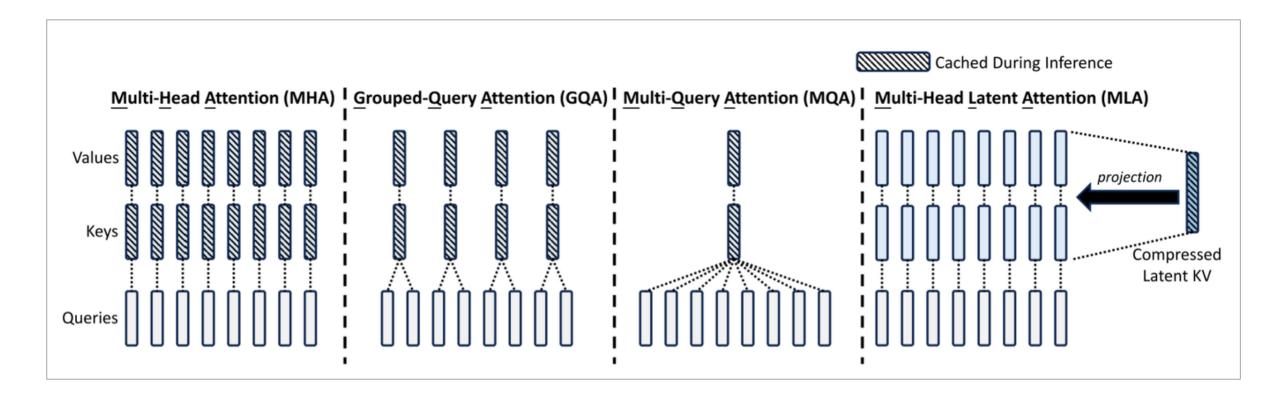






### **Multi-head latent attention**

- Compress KV cache into a latent space
  - Save up to 93.3% KVC space compared to original MHA (70KB per token!)

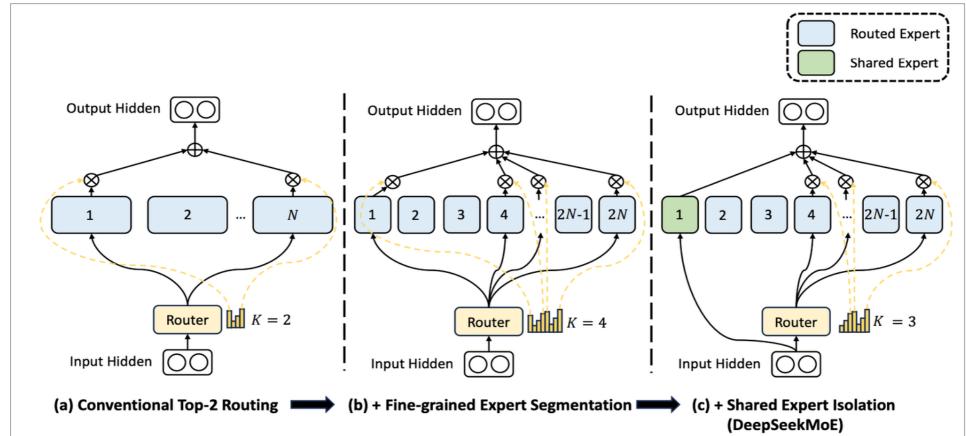






### **DeepSeek Mixture-of-Experts**

- Fine-grained experts 256 experts (activation 37B)
- Shared experts vs. routed experts
  - 1+256 experts

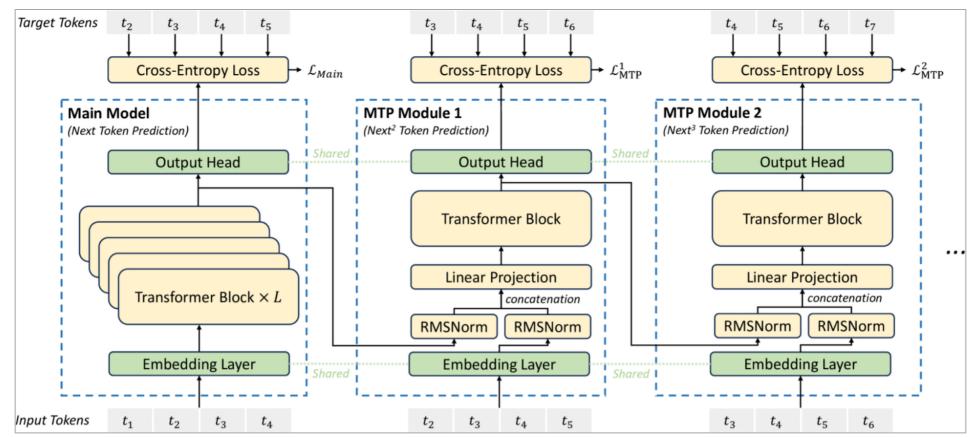






### **Multiple Token Prediction**

- Involving multiple token prediction at training stage
- Increase stability at when training, and also increasing throughput as speculative decoding

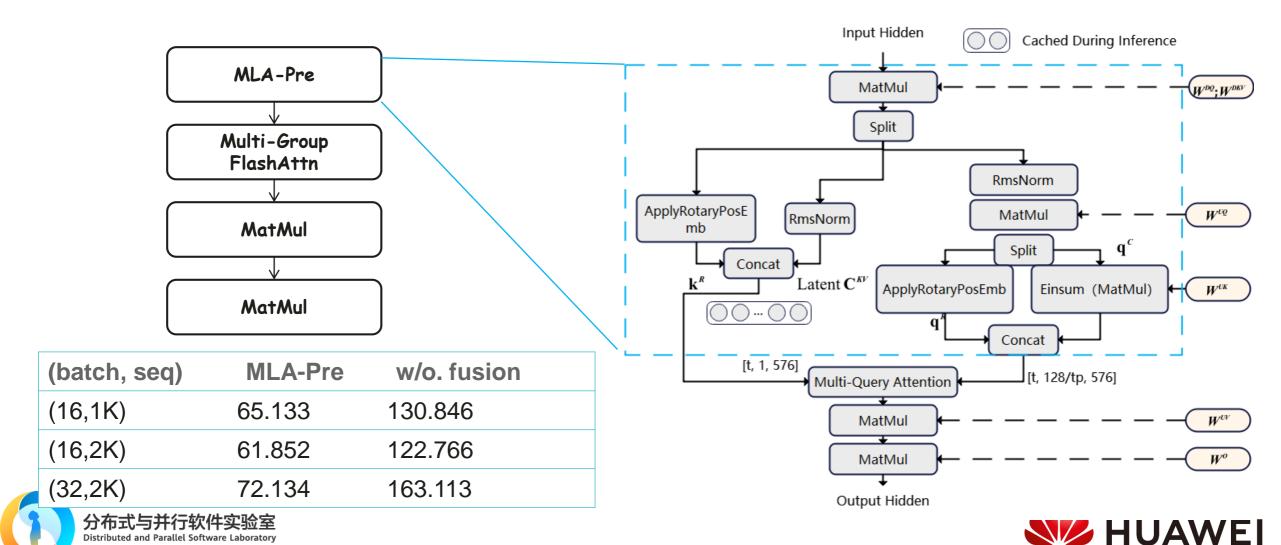






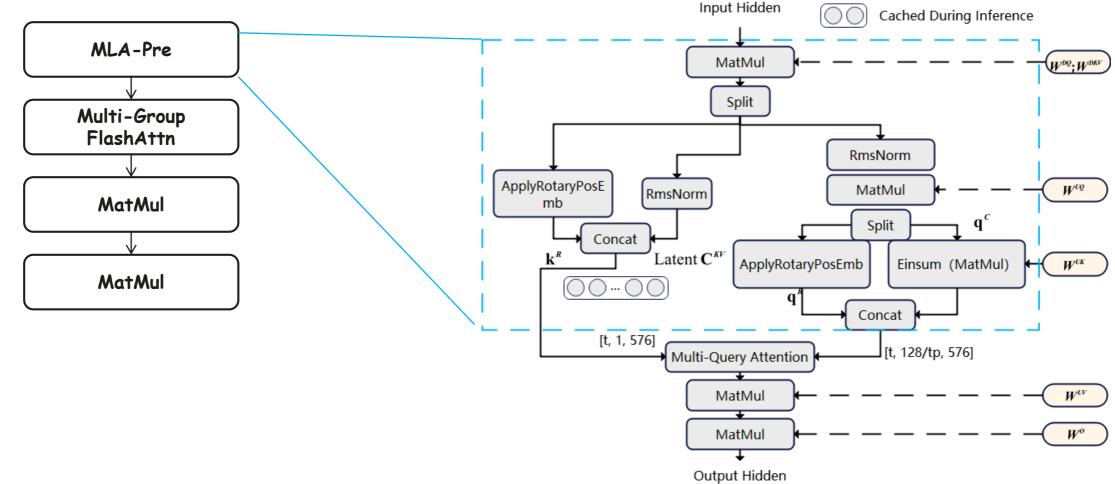
### **Optimizing DeepSeek infer. on Ascend – MLA Kernel**

Our strategy: leverage existing FlashAttention implementation



### **Optimizing DeepSeek infer. on Ascend – MLA Kernel**

Our strategy: leverage existing FlashAttention implementation

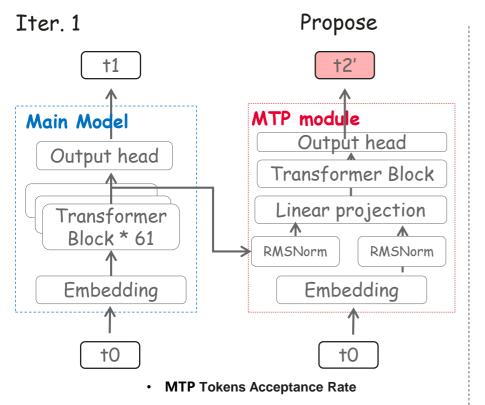


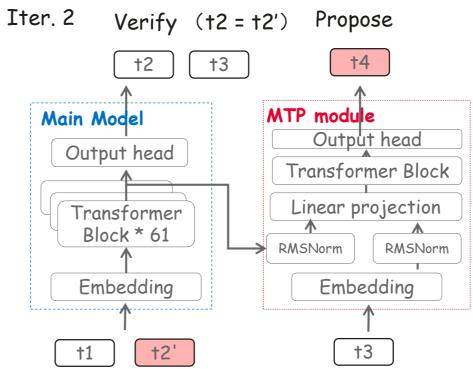




### **Optimization 2 – MTP decoding**

- MTP-based speculative decoding is helpful
  - Improving decoding throughput by 1.8x





	_	
Dataset	MTP	
	Accuracy	
C-Eval	81.4%	
GSM8K	88.7%	
MMLU	90.7%	





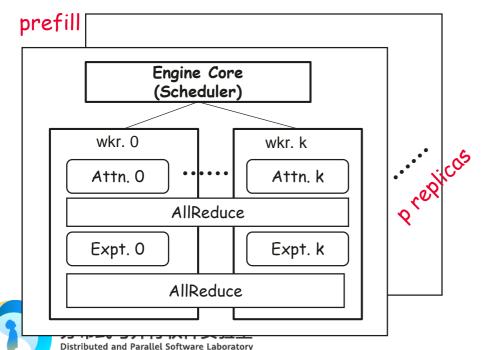
### **Optimization 3 – Prefill-decode disaggregation**

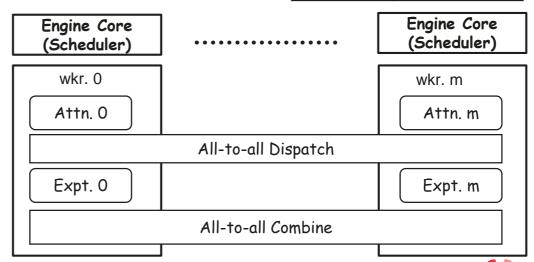
- Latency orientated deployment minimizing the end-to-end delay
  - TTFT < 1s; TPOT < 50ms ~ 20 tps throughput per user</li>
  - Batching as more requests as possible to increase NPU utilization
- Adopt different strategies to maximize parallelism at prefill- and decode-stage

decode

- Prefill Tensor-Parallel + Data-Parallel
- Decode Expert-Parallel + Fine-grained Data-Parallel

K*p	M
16*2	32
16*4	64
16*16	288

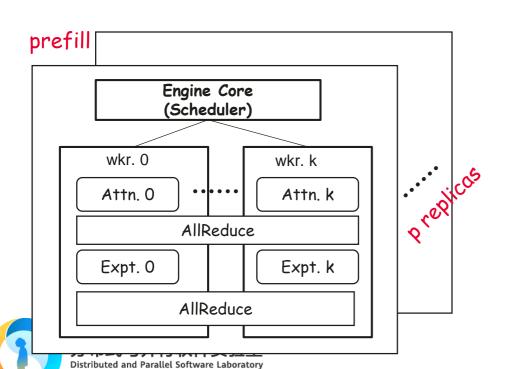




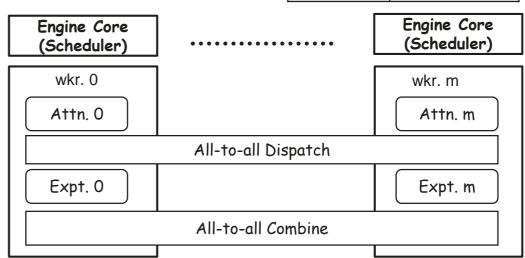


### **Optimization 3 – Prefill-decode disaggregation**

decode



K*p	M
16*2	32
16*4	64
16*16	288





### **Optimization 4 – Dispatch/Combine for MoE EP**

- Kernels fusing two all-to-all collectives for EP
  - Pre-allocated communication buffer and dispatching/combining
  - Finer compute-communication overlapping with MEM WRITE
  - Transmission scheduling for shared experts avoiding in-casting
  - Communication-time quantization

## Naïve multi-kernal Gating All-to-all-v Expert... All-2-all dispatch Gating Quantization Dispatching All-to-all sync. Expert...

# Gate 0 Gate 1 Gate m S-Exp 0 S-Exp 1 S-Exp m S-Exp m S-Exp m

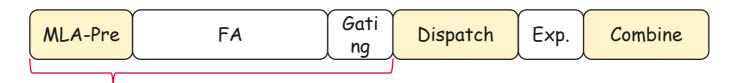
Unscheduled

### Put them all together

Performance w. seglen = 2K

Cluster size	Batch	TPOT (us)	TPUT (tpsc)
64	1	34	29.4
128	32	50	640
288	96	55	1745

Breakdown



Attention	Dispatch	Exp.	Combine
846	184	234	216





### **Insights and future directions**

Better model and hardware designs close the gap among computing, memory, and

communication

e.g., MLA, UB interconnection

- Memory bandwidth is still bottleneck, but the gap is much smaller
- Tricks (like better quantization) can still be applied to even close the gap
- Future optimization needs to consider both in memory and compute
  - Sparsity in time (sequence) dimension is promising



- Enable finer pipelining and better computing-communication overlapping
- FusedDeepMLA MLA-Pre + FlashAttn + Gating
- FusedDeepMOE Dispatch + Exp. + Combine
- Expecting 20+% performance boosting
- Better frameworks to build fused kernels





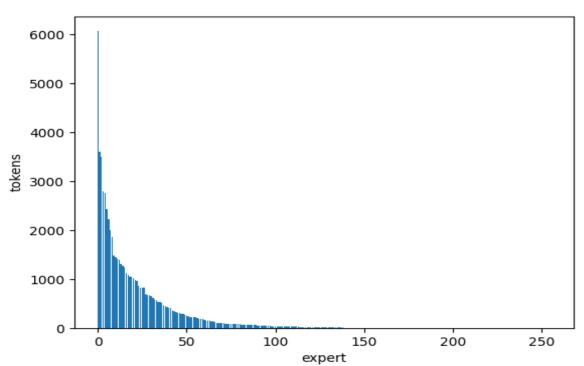
40%

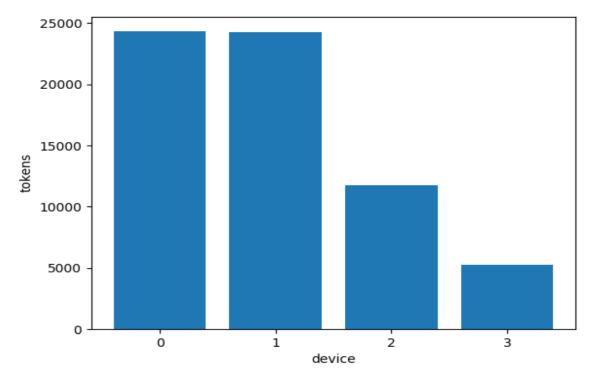
60%

■ Computing ■ HBM accessing

### Insights and future directions (cont.)

- Benchmark vs. deployment
  - Canary deployment on production environment
  - Only ~ 50% of benchmarked throughput
  - Unbalancing among experts calls for dynamic experts scheduling











## PSL Thank you!





