

# SUSTech Technique Report

APAC HPC-AI 2025 Student Competition

Team:

SUSTech Scout Regiment

Captain:

Haibin Lai

Email:

[laihb2022@mail.sustech.edu.cn](mailto:laihb2022@mail.sustech.edu.cn)

## Introduction

The **SUSTech Scout Regiment** is a passionate and multidisciplinary team from the **Southern University of Science and Technology (SUSTech)**. Our members come from different academic years and research backgrounds, united by a shared enthusiasm for **high-performance computing** and **scientific software optimization**. We aim to tackle real-world supercomputing challenges while learning from hands-on system-level experiments and performance tuning.

**SUSTech Team Introduction**

HPC·AI  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE

Jiahua Zhao (Advisor)

**Haibin Lai (Captain):** 4th year Undergrad at Computer Science

Benxiang Xiao: 3rd year Undergrad at Data Science

Bowen Zhang: 2nd year Master of Engineering in Electronic Science and Technology

Ziheng Wang: 3rd year Undergrad at Computer Science

Hemu Liu: 4th year Undergrad at Computer Science

**SUSTech** Southern University of Science and Technology

In this year's competition, our team utilized **two distinct high-performance computing environments** to evaluate and optimize our workloads.

For the **CPU track**, we conducted experiments on the **Qiming 2.0 cluster** at SUSTech, which is equipped with **Intel® Xeon® Gold 6138 CPUs** and a **Lustre parallel file system**. This environment closely resembles the **Gadi** and **NSCC** reference platforms provided by the competition organizers. We executed

**NWChem workloads** on up to **four nodes (160 CPU cores)** to evaluate scalability, communication efficiency, and memory utilization, while applying multi-level optimization techniques to enhance performance.



## CPU Cluster Introduction



To take the hardware as close as the Gadi and NSCC platform, we conducted our NWChem experiments on up to **4 nodes (160 CPU cores)** of the Qiming 2.0 cluster.



Component	Qiming 2.0 Description
System Name	Qiming 2.0 Supercomputing Cluster
Compute Nodes	409 blade nodes providing high-density parallel computing capability
Storage System	Lustre Parallel file system
Interconnect Network	Infiniband high-speed network (100 Gbps) with low latency

Component	Node Specification
CPU Model	Intel® Xeon® Gold 6138 @ 2.00GHz
Architecture	2 × 20 cores (40 cores total), 2 NUMA nodes
Threads per Core	1
Base / Max Frequency	2.00 GHz / 3.70 GHz



For the **GPU track**, our **DeepSeek-R1** experiments were performed on the **organizers' H200 reference cluster**. We are deeply grateful for the opportunity to access this advanced infrastructure, which allowed our software-level optimizations and fine-grained kernel tuning to fully demonstrate their performance potential.



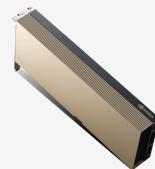
## GPU Cluster Introduction

Our performance on **Deepseek-R1 2 nodes 16 GPUs** workload is achieved on the organizers' H200 reference cluster. We're thankful for the access to this infrastructure, which allows our software optimizations to truly shine.

Cluster



H200 GPU



Configuration	Specification
GPUs per Node	8 * H200
Num of Nodes	4

Configuration	Specification
GPU BF16 Tensor Core	1979 TFLOPS
Memory	141GB



# NWChem

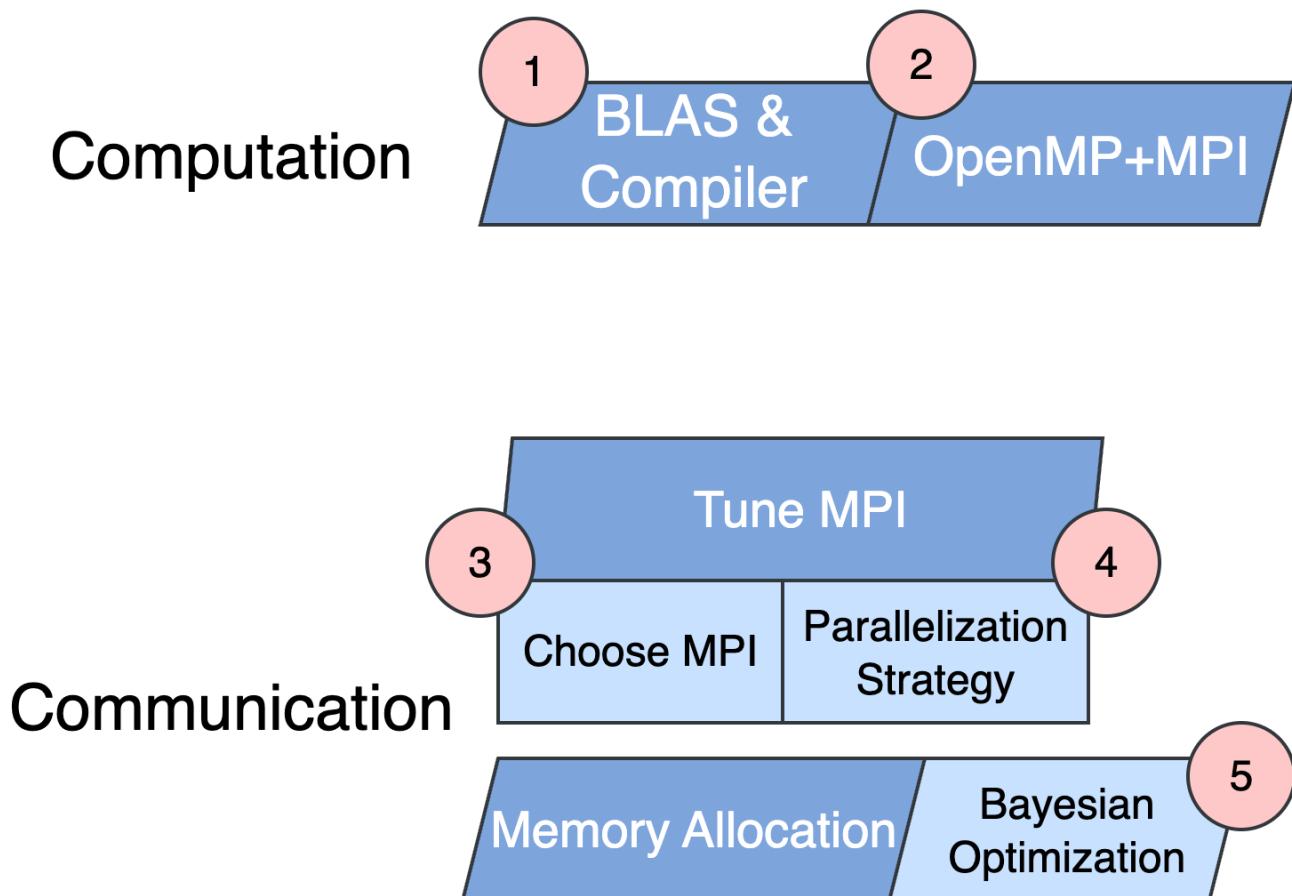
## Overview

**NWChem** is a scalable, open-source computational chemistry package designed for large-scale simulations on supercomputers. It supports a wide range of scientific models, including **quantum chemistry** and **molecular dynamics**.

Among its core components, the **Self-Consistent Field (SCF)** and **Density Functional Theory (DFT)** modules are particularly representative.

SCF workloads iteratively solve the electronic wavefunction, while DFT workloads extend this with **grid-based numerical integrations**.

Both workloads are **communication-intensive**, relying heavily on **MPI** and **Global Arrays (GA)** to perform distributed computation across nodes.



## Workload Description

In NWChem, data is organized in **Global Arrays**, which provide a shared-memory abstraction on top of MPI.

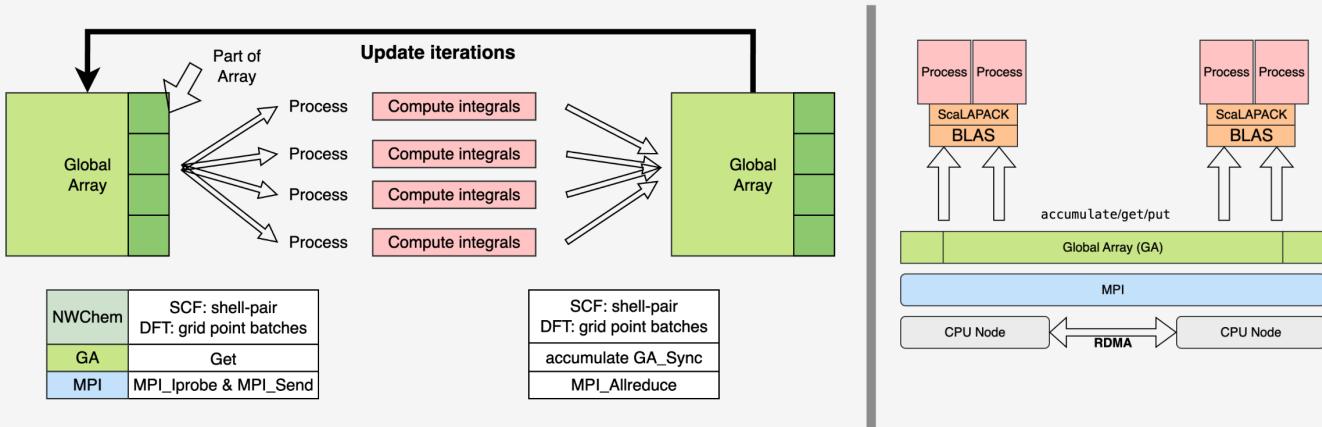
During each SCF or DFT iteration, processes **fetch sub-blocks of distributed arrays**, perform local computations, and then **write the results back** to the global data structure.

This design simplifies programming but introduces synchronization and communication overhead.



# NWChem workload analysis

Both SCF and DFT are communication-intensive and rely heavily on **MPI** and **Global Arrays** for distributed computation.



**SUSTech**

Southern University  
of Science and  
Technology

7

For our experiments, we ran the **baseline NWChem** build using **Intel MPI** and **OpenBLAS**.

We evaluated performance based on **wall-clock time**, ensuring that all outputs were scientifically correct. However, we observed **limited scalability beyond four nodes**, indicating potential bottlenecks in communication and computation.

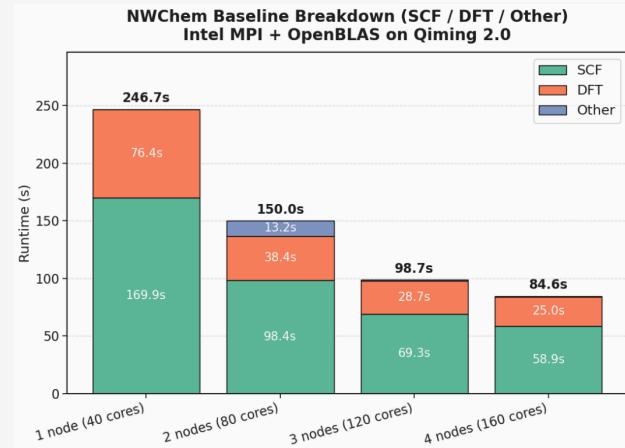


## Baseline Running

### Evaluation Metrics: **Wall Time** (lower is better)

Qiming 2.0 doesn't have precompiled NWChem, we compiled with **Intel-MPI 2021** with latest **OpenBLAS**.

**SCF Domains** the computation.  
Both workloads scales when using 2 nodes from 1 node. When using 3 nodes and 4 nodes, **they didn't scale very well**.



Total DFT energy = -917.738245100032  
One electron energy = -3287.42582587609  
Coulomb energy = 1466.406697276750  
Exchange-Corr. energy = -112.234989196318  
Nuclear repulsion energy = 1015.515872707146

Total SCF energy = -912.818341194699  
One-electron energy = -3287.137154917513  
Two-electron energy = 1358.802941015668  
Nuclear repulsion energy = 1015.515872707146



**SUSTech**

Southern University  
of Science and  
Technology

8

## Performance Analysis

We used **Intel VTune Profiler** to analyze runtime behavior and identify bottlenecks.

The results revealed that approximately **28% of total execution time** was spent in **communication**,

primarily due to **MPI\_Barrier** and **MPI\_Iprobe** operations.

This suggests **load imbalance** across processes, where some ranks finish computation earlier and wait for synchronization.

The second major bottleneck lies in **computation efficiency**.

Many of NWChem's core operators are **memory-bound**, meaning their performance is constrained by memory bandwidth rather than compute throughput.

As a result, optimization opportunities mainly come from **compiler-level tuning**, **loop transformations**, and **I/O improvements** rather than algorithmic restructuring.

Finally, we identified significant time spent in **BLAS routines**, especially **DGEMM (Double-precision General Matrix Multiplication)**.

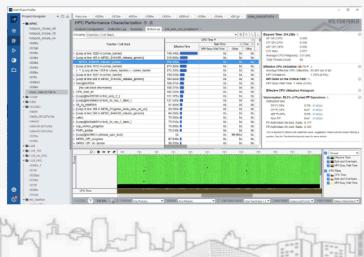
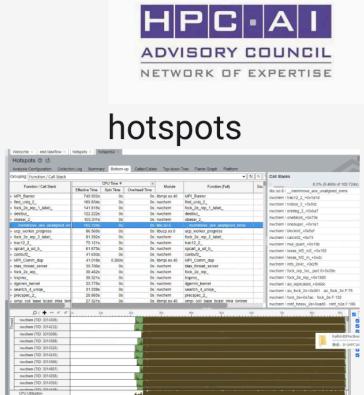
These are computationally heavy but can be improved through **optimized BLAS libraries**, **multi-threading**, or **NUMA-aware scheduling**.



## Profiling on 4 nodes: VTune

Functions/Operators	Time (s)	Explain
<b>MPI_Iprobe/MPI_Barrier - comex_barrier</b>	1134.457	From <b>GA and MPI</b> ; A barrier Wait all Processes
fock_2e_rep_1_label	243.038	From <b>NWChem</b> doing integrals
uniq_pairs_1	212.717s	From <b>NWChem</b>
dgemm_kernel	49.103	From <b>BLAS</b> , doing matmul

1. Barrier caused by Load Imbalance
2. Need Better Compiling & I/O
3. DGEMM From BLAS





9

## MPI Profiling Insights

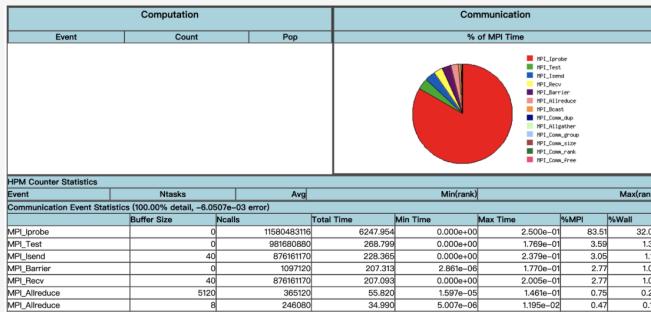
Using **IPM (Integrated Performance Monitoring)**, we further examined the MPI communication pattern.

The profiling results show that a large proportion of processes frequently invoke **MPI\_Iprobe**, indicating they are **polling for incoming data** instead of performing useful computation.

This behavior confirms our earlier observation: the system suffers from **communication waiting and load imbalance**, which severely limits parallel efficiency.



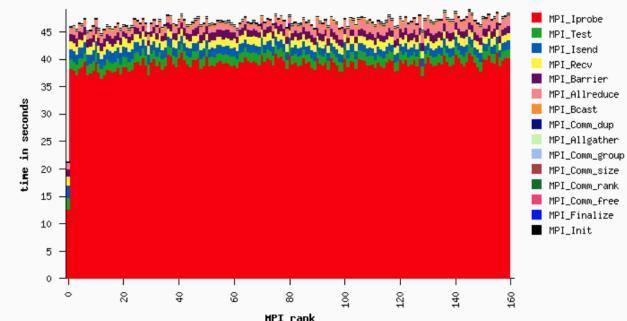
# Profiling on 4 nodes: IPM Analysis



83% MPI is waiting!

MPI\_Iprobe (red) Domains (For nonblocking test), MPI\_Test(Green)

MPI\_Recv(Yellow) MPI\_Send(Blue)



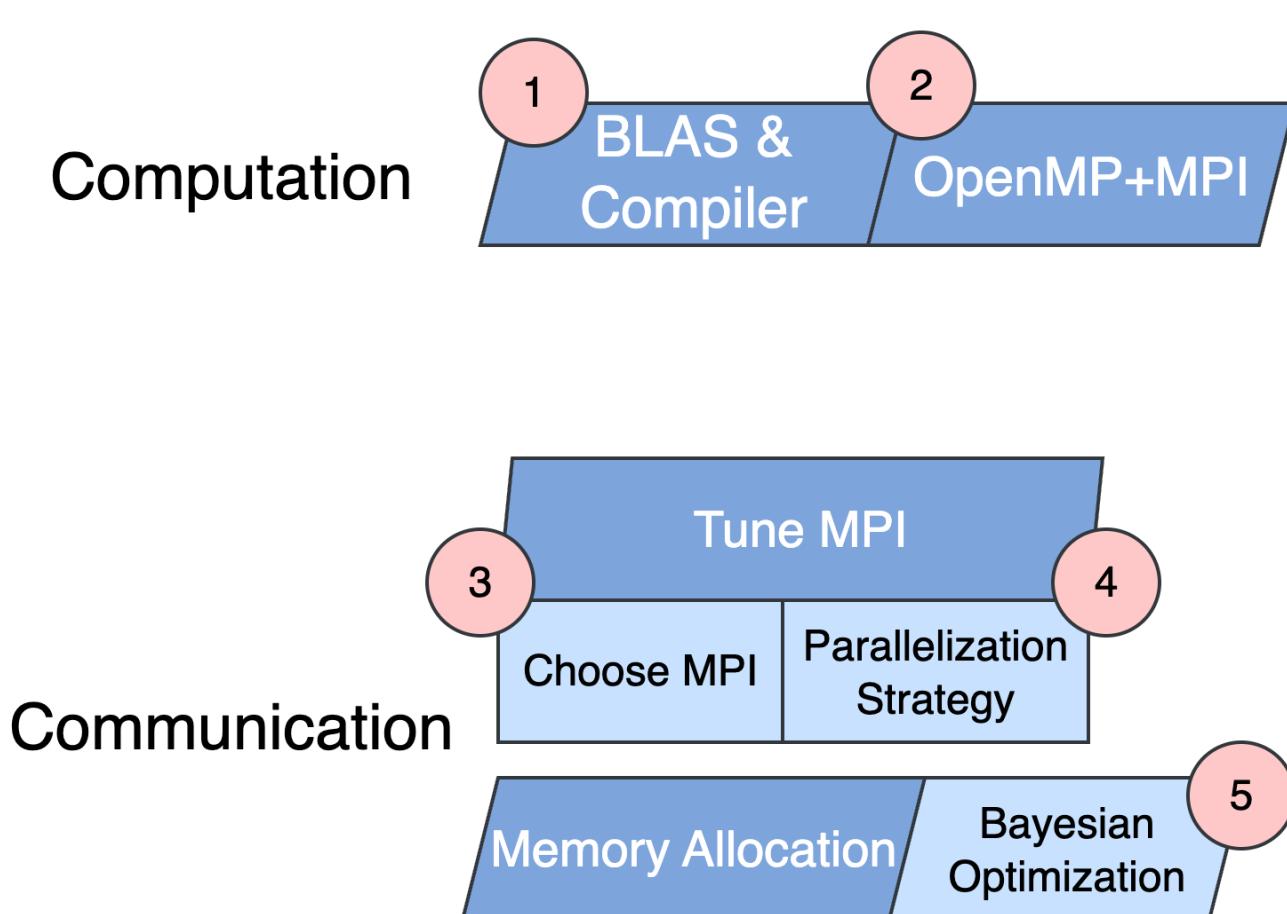
Processes in each nodes are waiting data



13

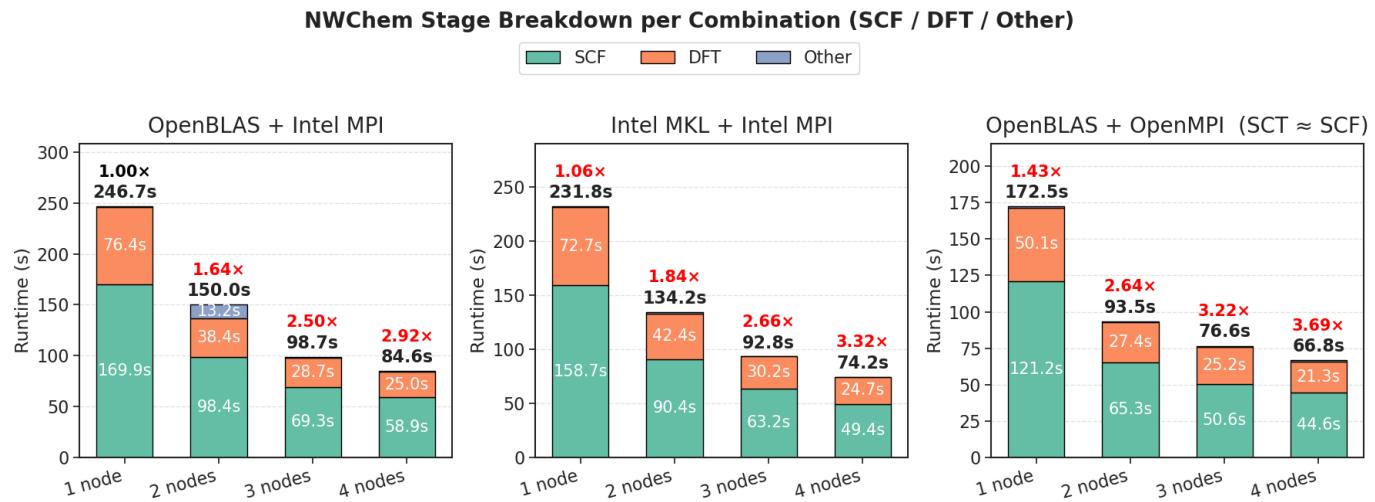
## Optimization Strategies

Based on the performance bottlenecks identified through **VTune** and **IPM** profiling, we proposed and implemented **five key optimizations** targeting both computation and communication inefficiencies.

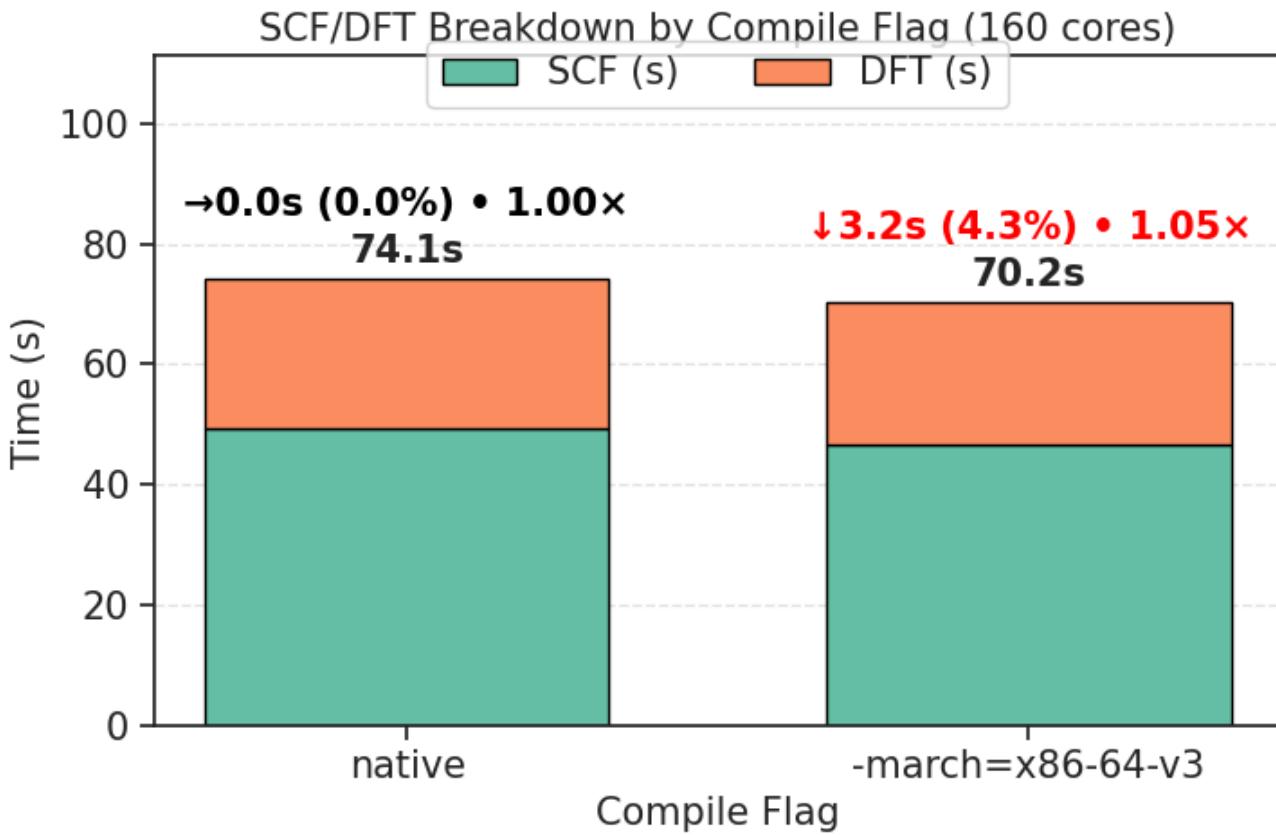


## 1. Compiler and BLAS Optimization

To improve computational efficiency, we experimented with multiple compiler and BLAS combinations. We compared **Intel OneAPI (ICC + MKL)** and **GCC + OpenBLAS**, observing that MKL generally provides better memory alignment and vectorization performance for dense matrix operations such as **DGEMM**.



By enabling advanced compiler flags (e.g., `-O3`, and `-fopenmp`), we improved the efficiency of memory-bound operators within SCF and DFT kernels.



## 2. NUMA + OpenMP + MPI Hybrid Parallelism

To better exploit on-node parallelism, we introduced **hybrid OpenMP + MPI parallelization**. Each node runs fewer MPI ranks, with OpenMP threads handling intra-node computations.

This reduces inter-rank communication and makes better use of shared caches and NUMA domains. We observed that OpenMP+MPI is not allowed in SCF and DFT workload, but we can still do the NUMA binding.



## Tuning OpenMP+MPI+NUMA binding



OpenMP can do memory sharing across threads, reducing memory passing cost. For example, for 40 cores we can have **1 process/4 omp threads per NUMA**.

```
export USE_OPENMP=1
```

**But OpenMP+MPI only works on CCSD modules!**

**So no OpenMP+MPI!**

**For NUMA: we have options:** `--map-by ppr:20:socket`

<https://hpcadvisorycouncil.atlassian.net/wiki/spaces/HPCWORKS/pages/2799534081/Getting+Started+with+NWChem+for+ISC22+SCC>

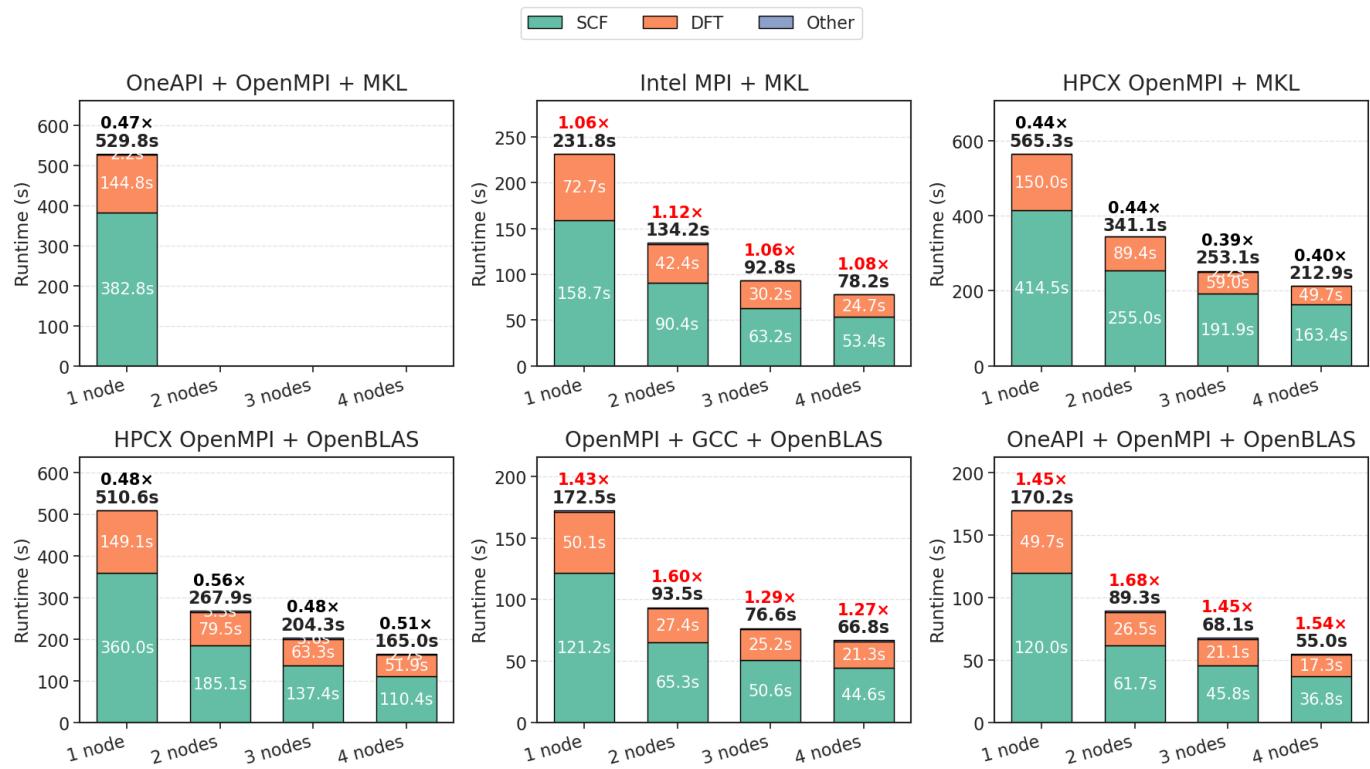


### 3. Communication and Load-Balancing Optimization

Communication imbalance—evident from excessive time in **MPI\_Barrier** and **MPI\_Iprobe**—was addressed by experimenting with **different MPI implementations and configurations**.

We tested **Intel MPI**, **hpcx**, and **OpenMPI**, tuning runtime parameters and collective algorithms.

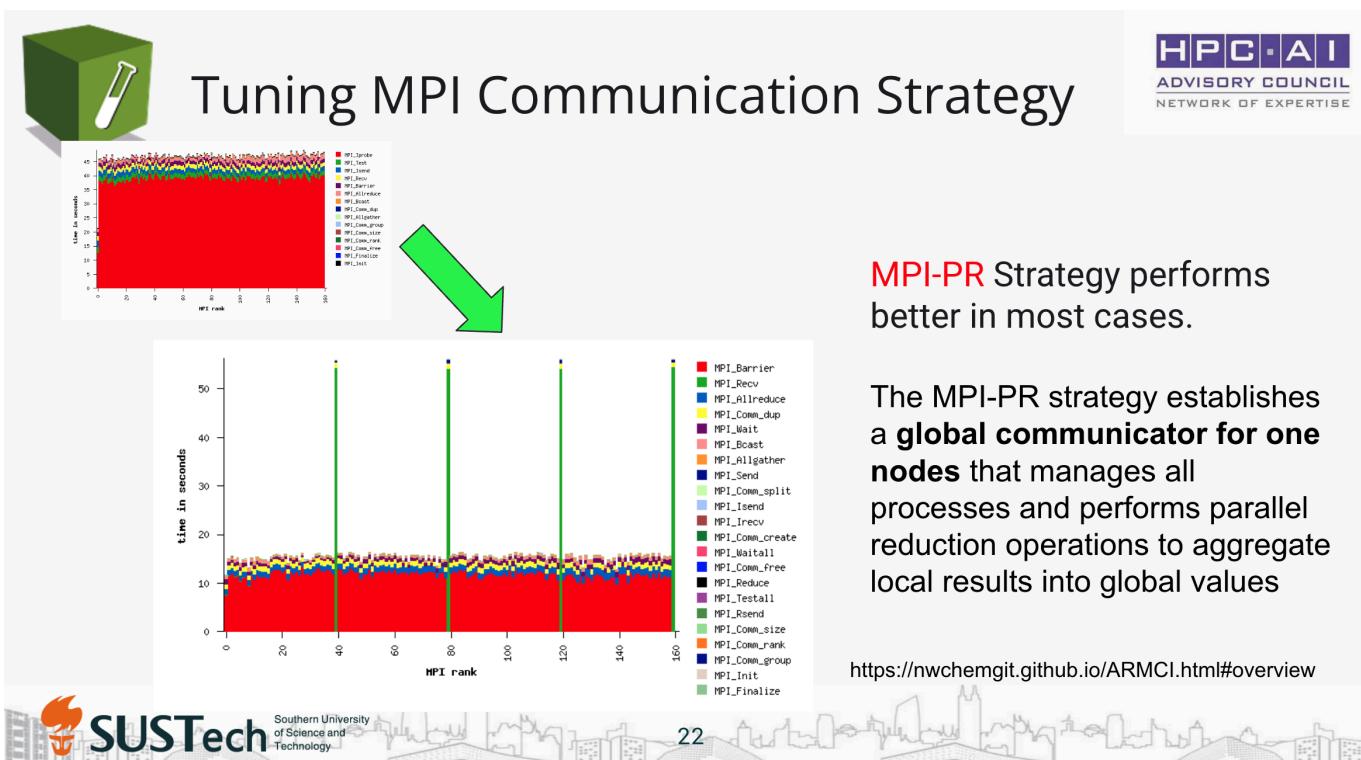
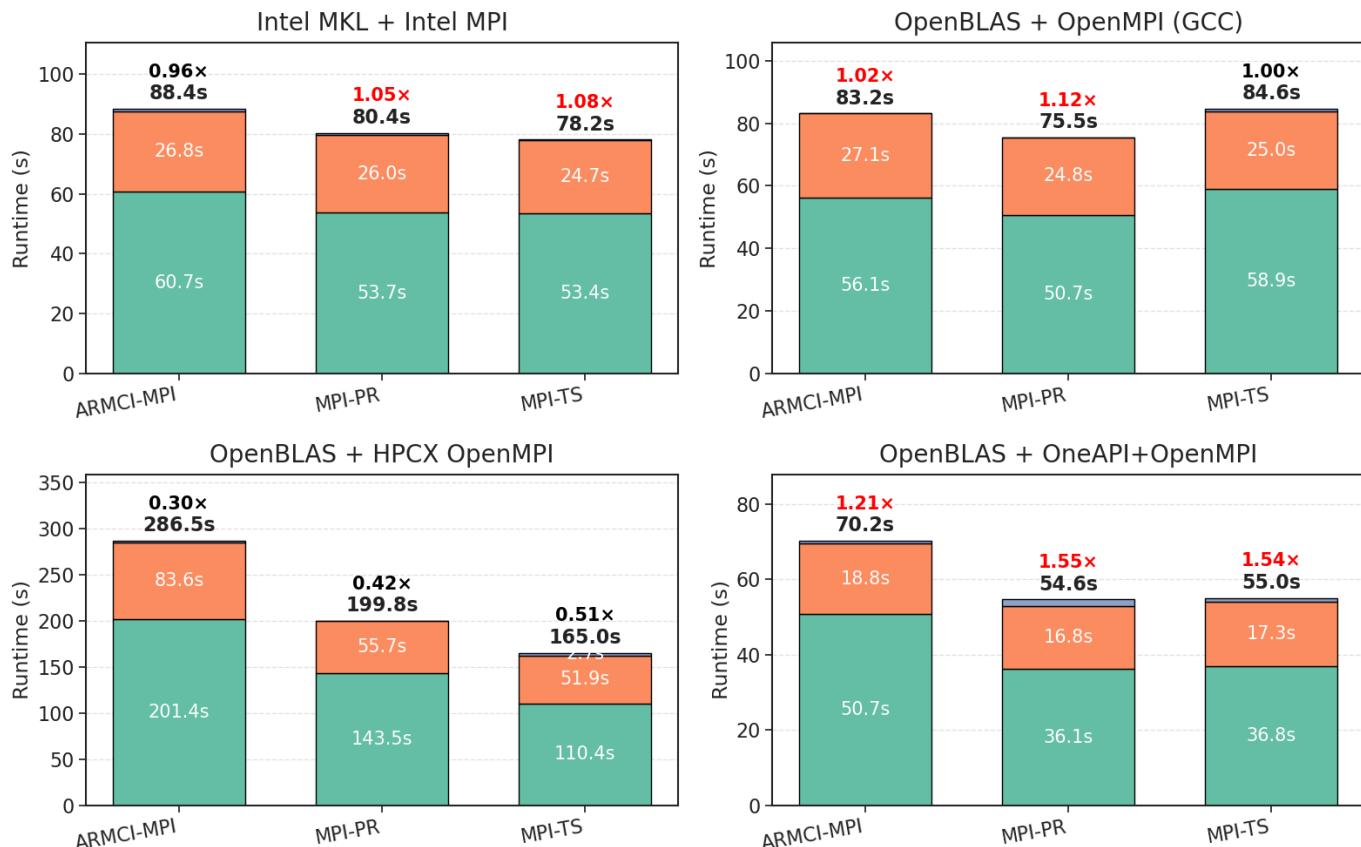
**NWChem Stage Breakdown per MPI/BLAS Combination (SCF / DFT / Other)**



Moreover, we explored **MPI process remapping** to improve affinity and reduce latency.

For the most communication-heavy phases, we evaluated **MPI-PR (Persistent Request)** strategies to overlap communication with computation, which reduced idle waiting time.

**NWChem @ 4 Nodes (160 Cores)**  
**MPI Communication Strategies per Software Stack (SCF / DFT / Other)**



## 4. Memory Allocation

Memory locality significantly impacts the performance of NWChem's Global Arrays.

We applied bayesian Optimization on Memory allocation.

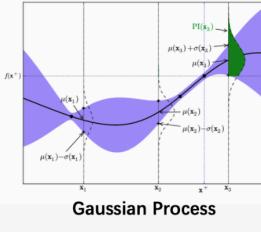
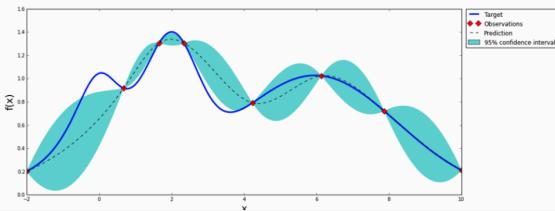
This tuning effectively reduced remote memory access and improved memory bandwidth utilization across nodes.



## Use Bayesian Optimization on Memory allocation

```
Algorithm 1 Sequential Model-Based Optimization
Input:  $f, \mathcal{X}, S, \mathcal{M}$ 
 $\mathcal{D} \leftarrow \text{INITSAMPLES}(f, \mathcal{X})$ 
for  $i \leftarrow |\mathcal{D}|$  to  $T$  do
     $p(y | x, \mathcal{D}) \leftarrow \text{FITMODEL}(\mathcal{M}, \mathcal{D})$ 
     $x_i \leftarrow \arg \max_{x \in \mathcal{X}} S(x, p(y | x, \mathcal{D}))$ 
     $y_i \leftarrow f(x_i)$   $\triangleright$  Expensive step
     $\mathcal{D} \leftarrow \mathcal{D} \cup (x_i, y_i)$ 
end for
```

**SMBO**

Unlike grid or random search, **Bayesian Optimization** builds a **probabilistic model (surrogate)** of the **performance function** and iteratively selects new configurations based on an **acquisition function** that balances exploration and exploitation.

This approach significantly **reduces the number of expensive experiments** while converging toward optimal parameters of memory params, enabling us to achieve higher performance with fewer trials.

Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. 2012. *Practical Bayesian Optimization of Machine Learning Algorithms*. arXiv preprint arXiv:1206.2944. <https://doi.org/10.48550/arXiv.1206.2944>

Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '19)*, August 4–8, 2019, Anchorage, AK, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.48550/arXiv.1907.10902>

SUSTech
Southern University of Science and Technology
24



## Use Bayesian Optimization on Memory allocation

**Framework: Optuna**  
**Performance/Optimized functions:**  
**Wall Time**

**Input Control Config**

Stack(MB): [4000, 16000], 500 per samples  
 Heap(MB): [50, 500], 50 per samples  
 Global(MB): [4000, 32000], 500 per samples  
 MPI: [MPI-PR, MPI-TS, ARM-CI]

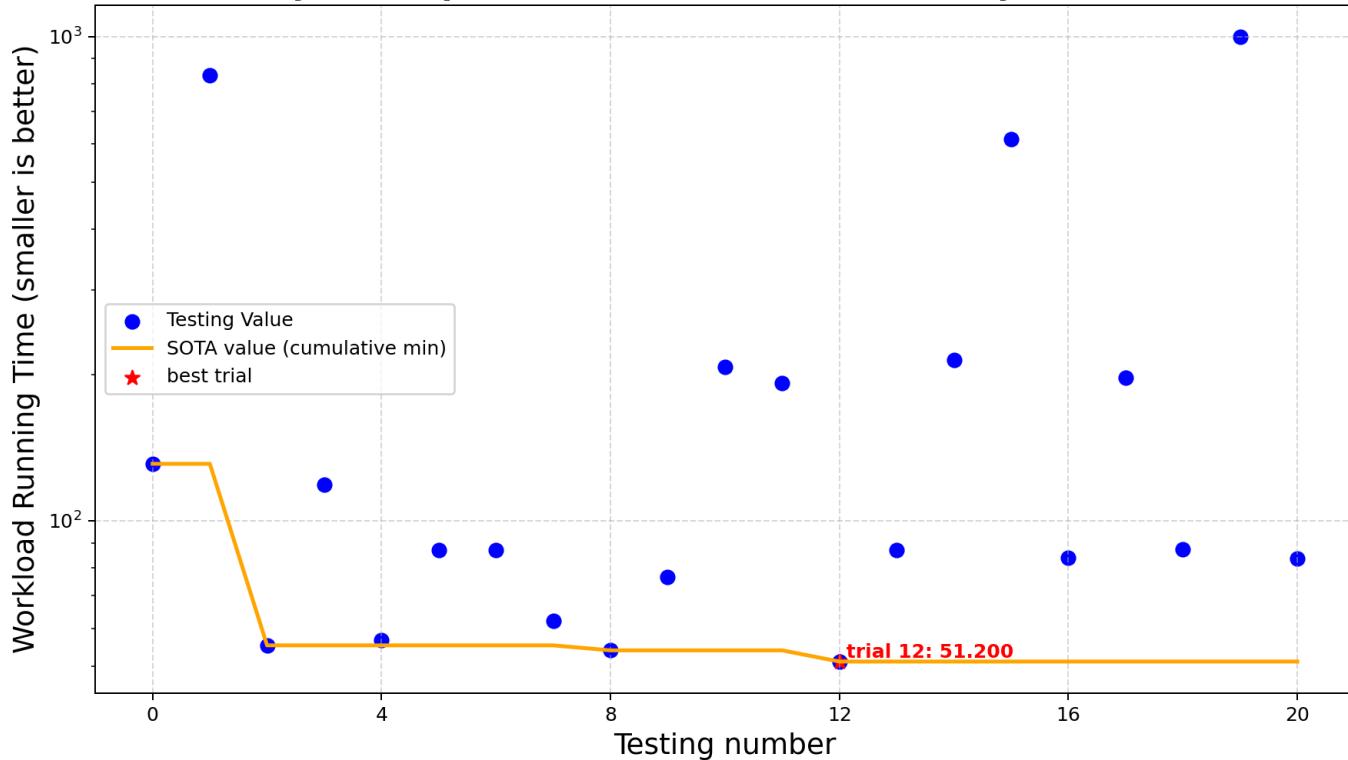
4 nodes  
 OneAPI icc-OpenMPI+OpenBLAS

Round	MPI	Stack	Heap	Global	Wall Time
1	MPI_PR	6500	150	10000	55.3
2	MPI_PR	14500	400	10000	118.7
3	MPI_PR	11000	350	24000	56.6
4	ARMCI-MPI	4000	400	31000	87
5	MPI_PR	2000	200	7000	87.1
6	MPI_TS	16000	250	26000	62.2
7	MPI_TS	12500	100	4000	54

SUSTech
Southern University of Science and Technology
25

**Algorithm 1** Sequential Model-Based Optimization
**SMBO**

## Bayesian Optimization for NWChem Memory Allocation



## Summary of Improvements

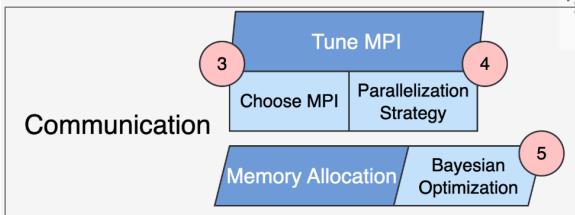
After applying these optimizations, we observed measurable performance gains across SCF and DFT workloads.

The optimized compiler + BLAS configuration delivered faster computation kernels.

Collectively, these improvements led to **better scalability and reduced wall time on up to four nodes**, demonstrating the importance of **co-designing communication, computation, and memory strategies** for scientific HPC applications.



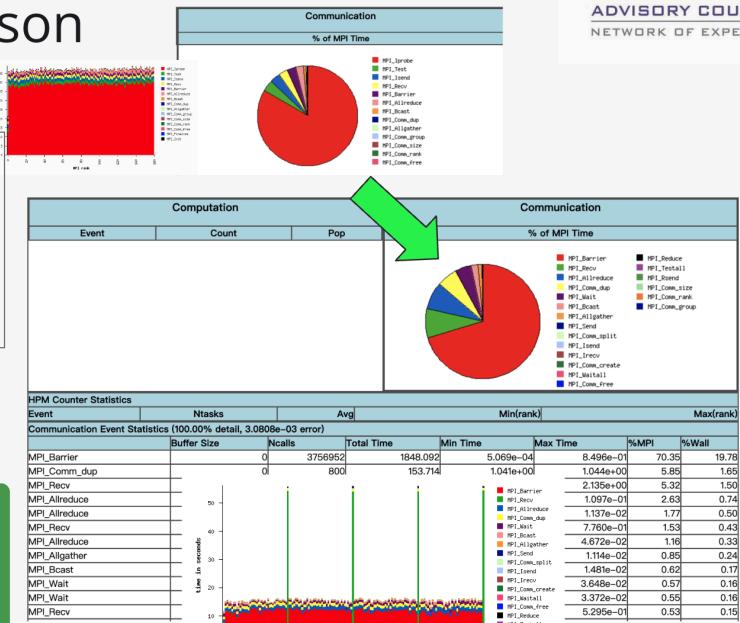
# Profile Comparison



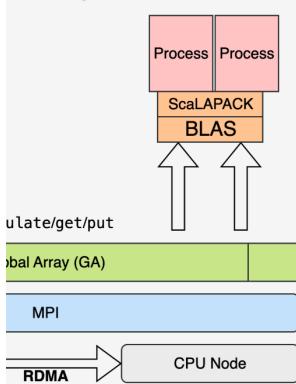
**MPI +54%**      **Memory +6%**

**MPI-PR +5%~+20%**

**MPI Barrier: 83.51% →  
70.35%**



## Conclusion



### Strategy

- CPU: **4 nodes 160 cores**
- BLAS: **OpenBLAS**
- Specific **Compile Flags**
- NUMA mapping**
- MPI: Intel icc compiled OpenMPI**
- Communication Strategy: **MPI-PR**
- Specific **Memory allocation**

Memory	Best Allocation
Stack	9000MB
Heap	400MB
Global	18500MB

### Acceleration

- 246.7s → 51.2s. **4.82X faster(-79%Time)** than single node baseline.
- 84.6s → 51.2s. **1.652X faster(-40%Time)** than 4 nodes baseline.



# Suggestions for Running NWChem

## Notice for Users

- Intel MPI **2022** doesn't have ifort
- **Choose your best MPI**
- **AOCL** may not work well on Intel Machines

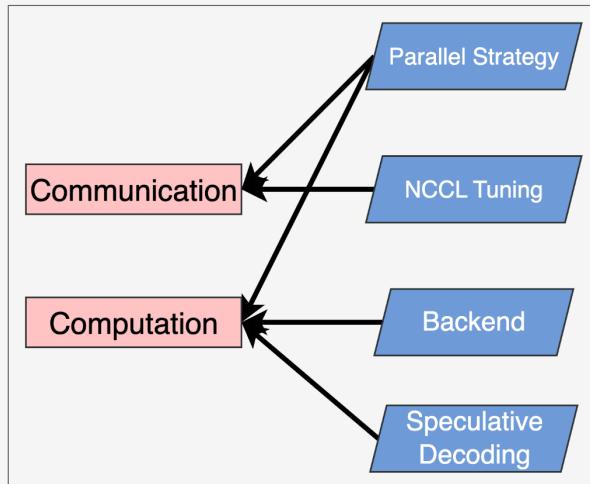
## Suggestion for Developers

- **Tree-based Barrier**
- **Need OpenMP+MPI support**
- Better algorithms on operators like sort

## DeepSeek-R1 Optimization



### Deepseek-R1 Optimization Roadmap



TP, DP, PP

Tree/Ring, LL, LL128

FlashAttention3, FlashMLA

auto-regressive boosting

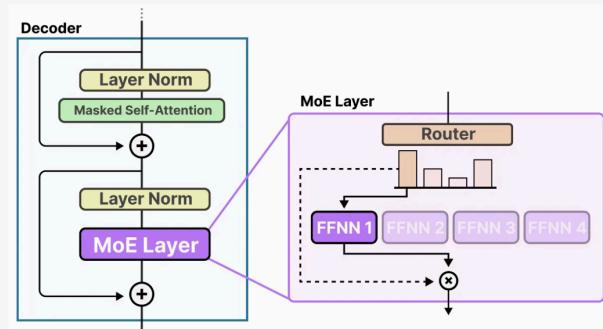
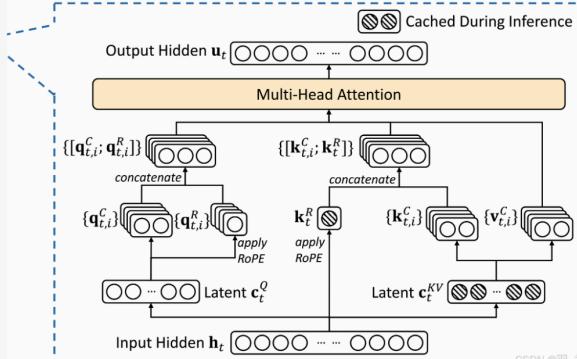


# Deepseek-R1 Introduction

DeepSeek-R1 (671B) achieves strong performance across various reasoning and language tasks.

Its **MLA architecture** enables more efficient use of attention, while the **MoE design** activates only **37 B parameters** per inference.

Multi-Head Latent Attention (MLA)

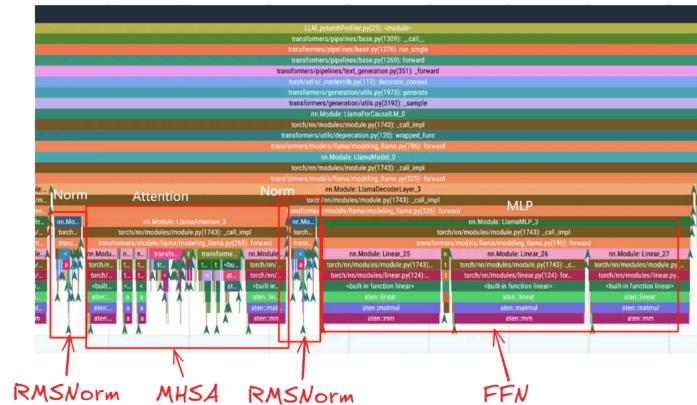
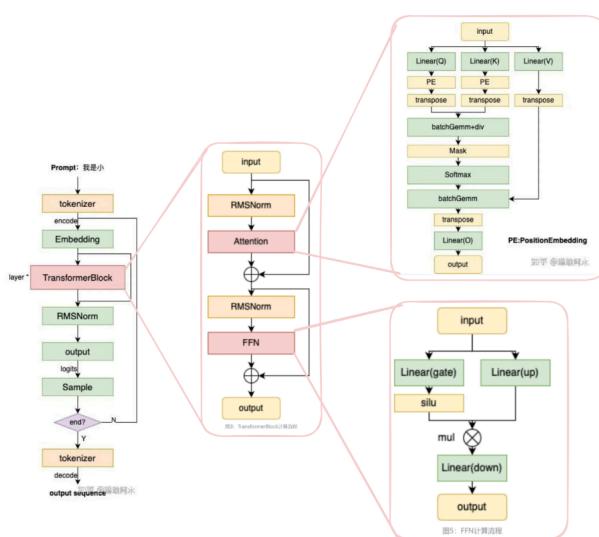


Southern University  
of Science and  
Technology

35

DeepSeek-R1 (671B) achieves strong performance across various reasoning and language tasks.

Its **MLA architecture** enables more efficient use of attention, while the **MoE design** activates only **37 B parameters** per inference.



SGLang supports Tensor Parallelism (TP), Pipeline Parallelism (PP), and Data Parallelism (DP), enabling efficient large-scale distributed inference across multiple GPUs and nodes.

It also integrates advanced features such as speculative decoding and Radix Attention, which significantly improve decoding throughput and reduce latency.

Through studying and experimenting with SGLang, we gained a deep understanding of its scheduling and parallel strategies. In our benchmarks, it demonstrated strong scalability and excellent performance efficiency on modern GPU clusters.



# Deepseek-R1 Backend: SGLang

HPC·AI  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE

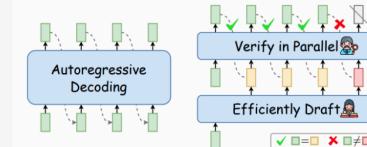
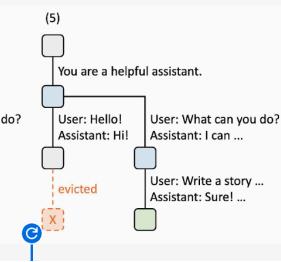
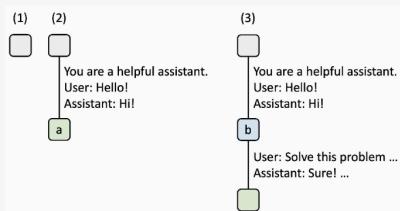


Figure 1: In contrast to autoregressive decoding (left) that generates sequentially, Speculative Decoding (right) first *efficiently drafts* multiple tokens and then *verifies* them *in parallel* using the target LLM. Drafted tokens after the bifurcation position (e.g.,  $\square$ ) will be discarded to guarantee the generation quality.

## Baseline Setup

For the **GPU track**, we conducted experiments on the **organizers' reference H200 cluster**, using a **two-node, 16-GPU configuration** ( $8 \times$  H200 per node, connected with NVLink + InfiniBand).

We evaluated multiple **parallelism strategies**, including:

- **Tensor Parallelism (TP)** – splitting matrix multiplications across GPUs,
- **Pipeline Parallelism (PP)** – dividing layers across nodes to balance memory usage, and
- **Data Parallelism (DP)** – replicating models to process multiple requests concurrently.

Our baseline employed the **SGLang engine** running with the **FA3 backend**, which served as a stable starting point for functional and performance correctness.



## Deepseek-R1 Baseline

HPC·AI  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE

### Metric: Total token throughput

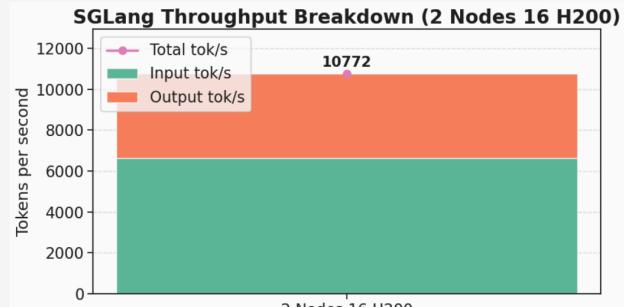
Default as LMSYS & APAC provides

OpenMPI

TP=16

PP=1

DP=1



```

1918 g3:2134014:2140090 [5] NCCL INFO Channel 12/0 : 5[5] -> 4[4] via
1919 g3:2134014:2140090 [5] NCCL INFO Channel 15/0 : 5[5]
1920 ===== Offline Throughput Benchmark Result =====
1921 Backend: engine
1922 Successful requests: 2000
1923 Benchmark duration (s): 94.26
1924 Total input tokens: 626729
1925 Total generated tokens: 388685
1926 Last generation throughput (tok/s): 72.93
1927 Request throughput (req/s): 21.22
1928 Input token throughput (tok/s): 6648.87
1929 Output token throughput (tok/s): 4123.50
1930 Total token throughput (tok/s): 10772.37
1931 =====

```

## Bottleneck Profiling

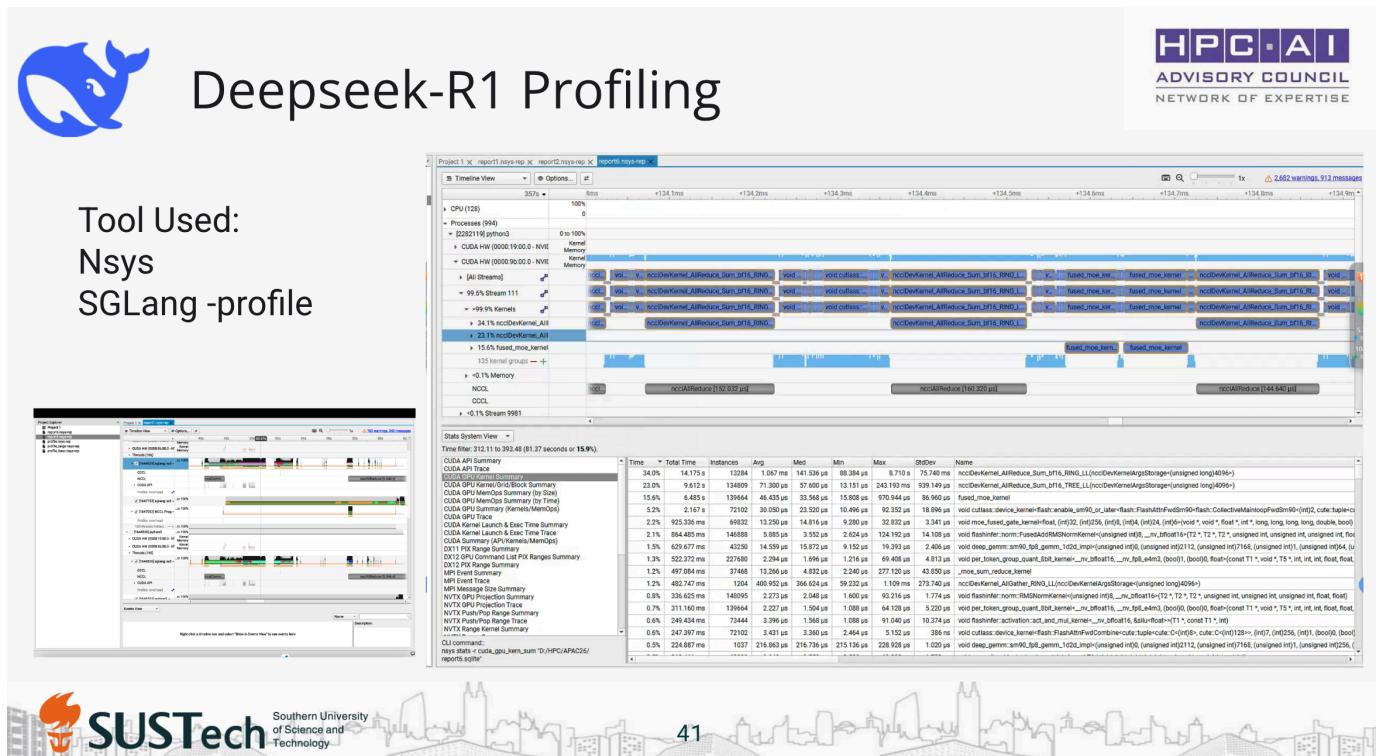
We used **NVIDIA Nsight Systems (nsys)** to profile both prefill and decode stages.

Profiling results indicated that the main performance bottleneck came from **communication overhead** rather than computation.

Specifically, NCCL traces showed numerous small but frequent data exchanges between GPUs, leading to

## high communication latency and low link utilization.

In multi-node cases (PP > 1 or DP > 1), synchronization between pipeline stages further aggravated the delay, dominating overall iteration time.





# Deepseek-R1 Profiling

**HPC-AI**  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE

Tool Used:  
Nsys  
SGLang -profile

We observed **small runtime but large instances** of NCCL events happens in every token generations:

Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Name
34.0%	14.159 s	13280	1.066 ms	141.536 $\mu$ s	88.384 $\mu$ s	8.710 s	75.742 ms	ncciDevKernel_AllReduce_Sum_bf16_RING_LL(ncciDevKernelArgsStorage<unsigned long>4096-)
23.1%	9.612 s	134809	71.300 $\mu$ s	57.600 $\mu$ s	13.151 $\mu$ s	243.193 ms	939.149 us	ncciDevKernel_AllReduce_Sum_bf16_TREE_LL(ncciDevKernelArgsStorage<unsigned long>4096-)
15.6%	6.485 s	139664	46.435 $\mu$ s	33.568 $\mu$ s	15.808 $\mu$ s	970.944 us	86.960 ms	fused_moe_kernel
5.2%	2.167 s	72102	30.050 $\mu$ s	23.520 $\mu$ s	10.496 $\mu$ s	92.352 $\mu$ s	18.896 us	void cutlass::device_kernel::flash::enable_sm90_0_r later<flash::FlashAttnFwdSm90<flash::CollectiveMainloopFwdSm90<(int)2, cute::tuple<cu
2.2%	925.336 ms	69832	13.250 $\mu$ s	14.816 $\mu$ s	9.280 $\mu$ s	32.832 $\mu$ s	3.341 $\mu$ s	void moe_fused_gate_kernel<float, (int)32, (int)256, (int)8, (int)4, (int)24, (int)6>(void *, void *, float *, int, long, long, long, long, double, bool)
2.1%	864.485 ms	146888	5.885 $\mu$ s	3.552 $\mu$ s	2.624 $\mu$ s	124.192 us	14.108 us	void flashInfer_norm::FusedAddMSNormKernel<(unsigned int)8, _nv_bf16o16<(T2 *, T2 *, T2 *, unsigned int, unsigned int, unsigned int, unsigned int, float>
1.5%	629.677 ms	43250	14.559 $\mu$ s	15.872 $\mu$ s	9.152 $\mu$ s	19.393 us	2.406 us	void deep_gemm::sm90_fp8_gemm_1d2d_impl<(unsigned int)0, (unsigned int)2112, (unsigned int)7168, (unsigned int)1, (unsigned int)64, (u
1.3%	522.337 ms	227676	2.294 us	1.696 us	1.216 us	69.408 us	4.813 us	void per_token_group::quant_8bit_kernel<_nv_bf16o16, _nv_fp8_e4m3, (bool)0, float>(const T1 * const T1 *, void *, T5 *, int, int, float, float)
1.2%	497.084 ms	37468	13.266 $\mu$ s	4.832 $\mu$ s	2.240 $\mu$ s	277.120 $\mu$ s	43.850 $\mu$ s	_moe_sum_reduce_kernel



# Deepseek-R1 Profiling

**HPC-AI**  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE

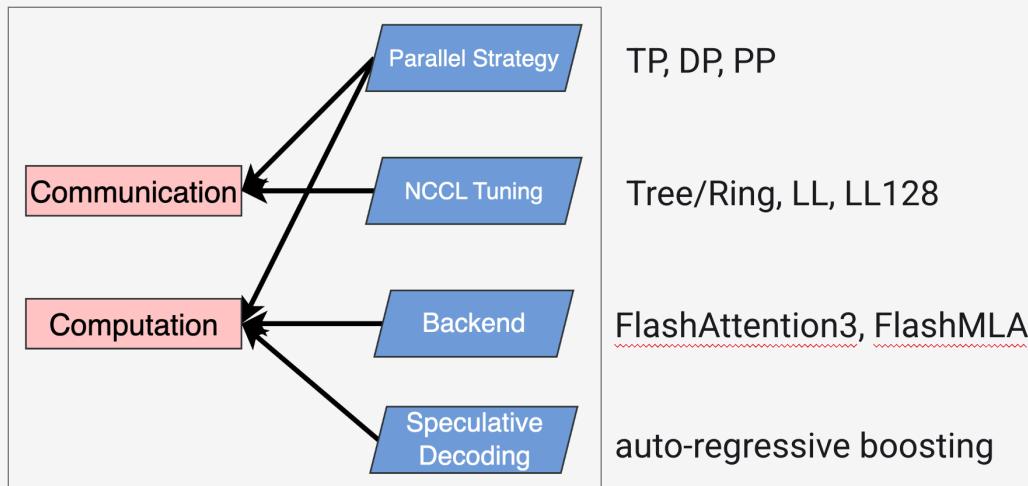
Tool Used:  
Nsys  
SGLang -profile

**Better MOE Kernels in backend library** may compute faster for the model.



# Deepseek-R1 Optimization Roadmap

**HPC-AI**  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE



TP, DP, PP

Tree/Ring, LL, LL128

FlashAttention3, FlashMLA

auto-regressive boosting



47

## 1. Parallelization Strategy

\*\*

For TP=8 and PP=2, all 8 H200 GPUs within a node are used for TP, while PP spans nodes — minimizing inter-node communication.

DP = 2 enables dual-node concurrency with minimal communication, doubling throughput while keeping TP/PP balanced inside each node.

\*\*



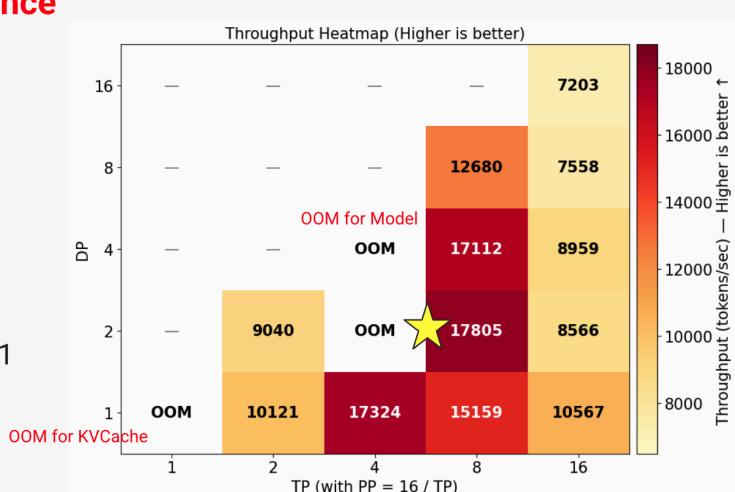
## Optimization: Parallelization Strategy

**HPC-AI**  
ADVISORY COUNCIL  
NETWORK OF EXPERTISE

**TP=4-8, DP=1-2 reach best performance**

- TP=16, PP=1, DP=1
- TP=16, PP=1, DP=2
- TP=16, PP=1, DP=4
- TP=16, PP=1, DP=8
- TP=16, PP=1, DP=16
- TP=8, PP=2, DP=1
- TP=8, PP=2, DP=2
- TP=8, PP=2, DP=4
- TP=8, PP=2, DP=8

- TP=4, PP=4, DP=1
- TP=4, PP=4, DP=2
- TP=4, PP=4, DP=4
- TP=2, PP=8, DP=1
- TP=2, PP=8, DP=2
- TP=1, PP=16, DP=1



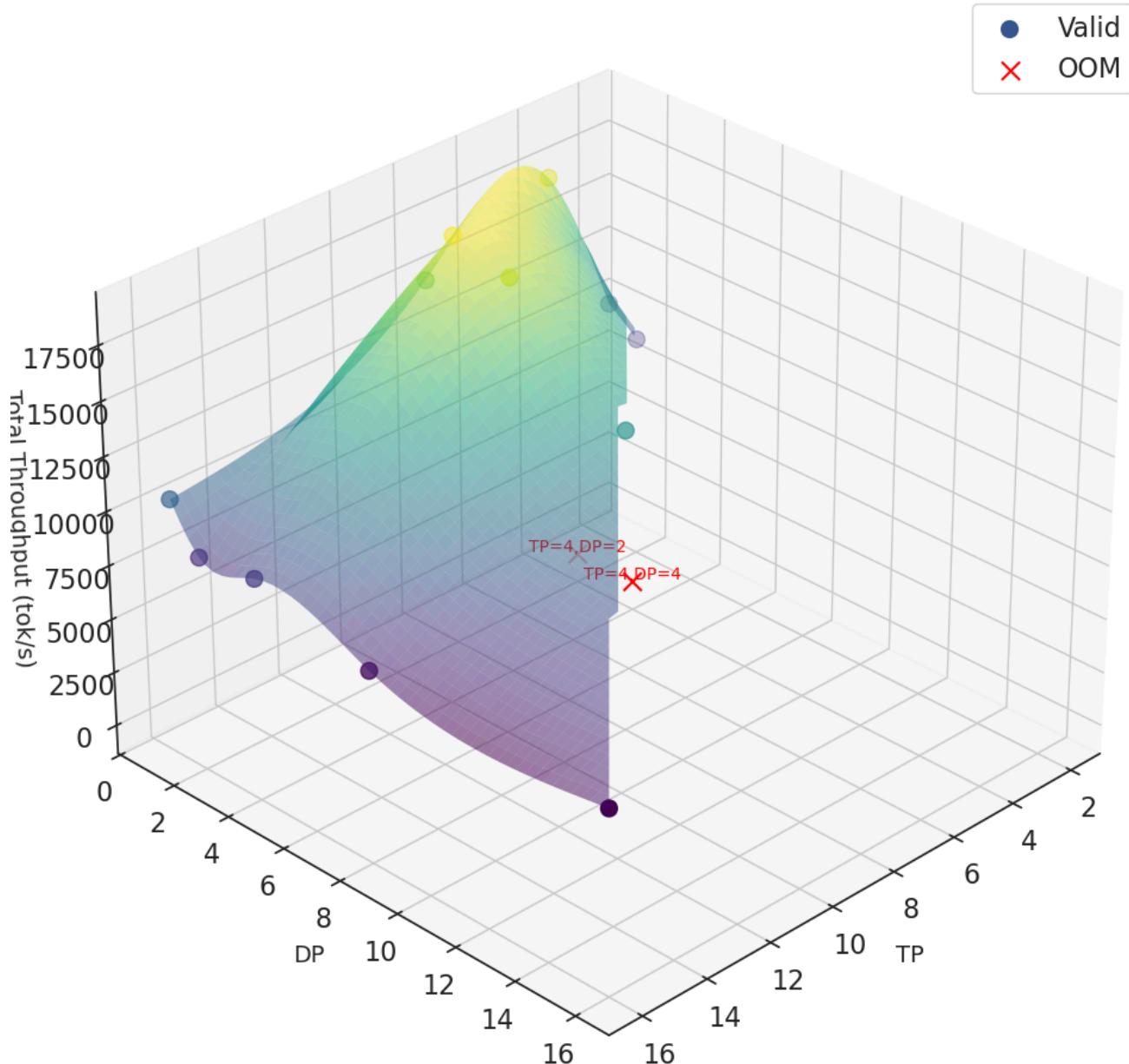
49

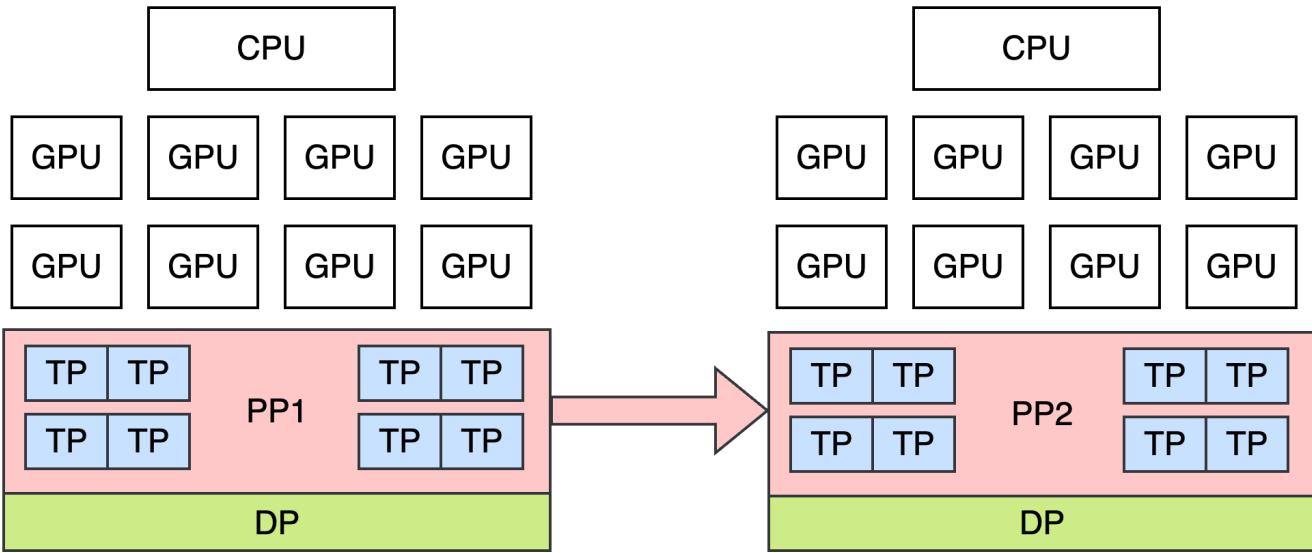
## 1. Parallelism Strategy Exploration

We systematically tested combinations of (TP, PP, DP) to analyze the trade-off between inter- and intra-node communication.

The configuration **TP = 8, PP = 2** achieved the best balance — using intra-node tensor parallelism to minimize cross-node data transfer while leveraging pipeline parallelism for memory scalability.

**SGLang DeepSeek-R1 Throughput (2-node 16-GPU)**

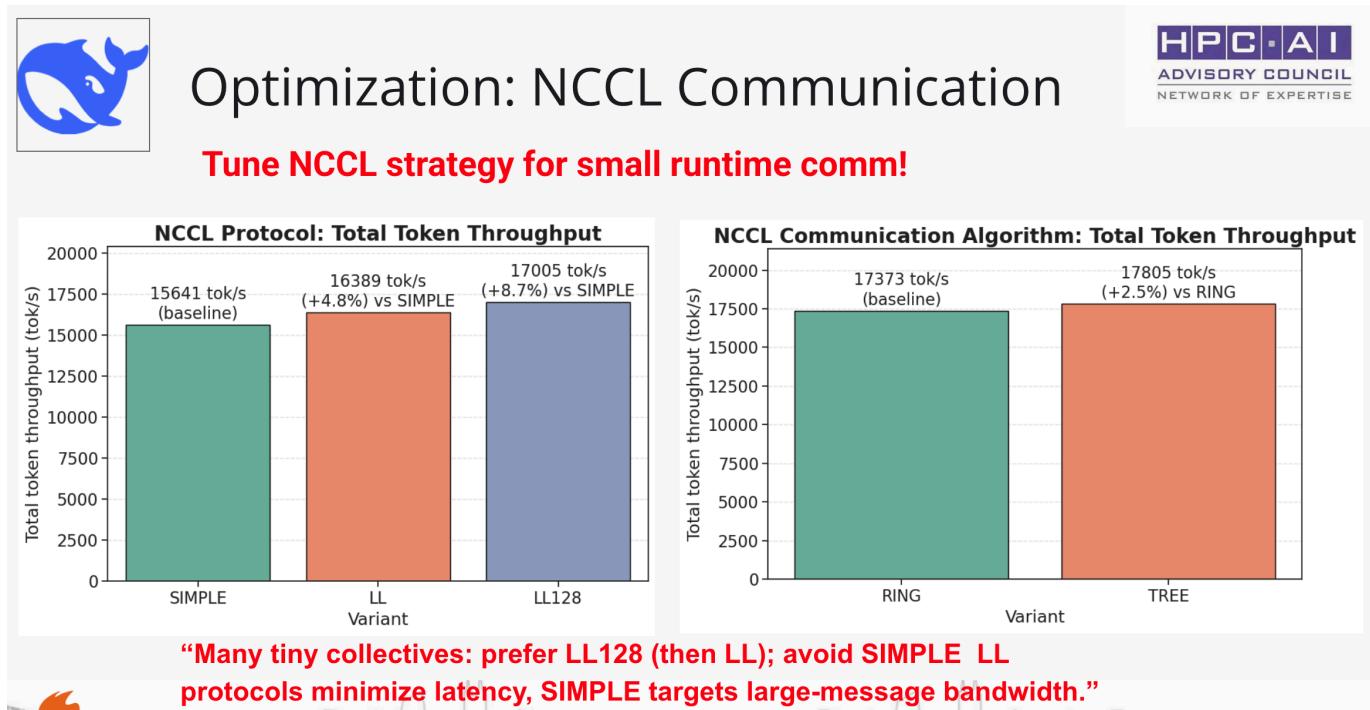




## 2. NCCL Communication

We tuned NCCL parameters and communication strategies to better suit the small-message, high-frequency workload pattern.

Tree-based algorithms replaced ring algorithms for latency-sensitive collectives, and the **LL128 protocol** was adopted for small transfers.



## 3. backend

\*\*

FA3, the default SGLang kernel, already achieves the best overall throughput among all backend implementations

\*\*



# Optimization: Backend

fa3	17624.96
triton	16539.77
flashmla	fail to run due to network issue
flashinfer	only allowed in A100

FA3, the default SGLang kernel, **already achieves the best** overall throughput among all backend implementations



54

During the competition, we initially encountered difficulties running **FlashInfer** within the SGLang framework. Profiling showed communication stalls that prevented correct execution under certain (TP, PP, DP) configurations. After carefully reproducing and analyzing the issue, we reported our findings to the **SGLang official development team** through a detailed GitHub issue, as also suggested by the judges.

Subsequent investigation by the SGLang developers confirmed that the problem originated from an **in-development communication bug** in the FlashInfer backend. The issue has since been **officially fixed** in the latest SGLang release. We re-tested the workload on the **two-node, 16-GPU H200 cluster**, and **FlashInfer now runs successfully** with stable throughput and correct outputs.

issue [https://github.com/sql-project/sqlang/issues/12402](https://github.com/sgl-project/sqlang/issues/12402)

This collaboration not only helped us complete our benchmark successfully but also contributed to improving the SGLang ecosystem. We are grateful to the organizers and the SGLang team for their timely support.

Open

[Bug] FlashMLA backend stuck at prefilling stage during benchmark #12402



Fridge003 5 days ago

Collaborator ...

Hi [@HaibinLai](#), we just integrated flashmla kernels into sgl-kernel recently [#12135](#)  
Can you try to rebase on the latest main branch and run again?



HaibinLai 4 days ago

Author ...

Hi [@Fridge003](#), thank you for your reminder! Let me retry it on our machines!



HaibinLai now

Author ...

Hi [@Fridge003](#), we tested the flashmla kernels **successfully** on our cluster! Thank you for your contributions!

```
[0] ====== Offline Throughput Benchmark Result ======
[0] Backend: engine
[0] Successful requests: 2000
[0] Benchmark duration (s): 145.41
[0] Total input tokens: 626729
[0] Total generated tokens: 388685
[0] Last generation throughput (tok/s): 77.21
[0] Request throughput (req/s): 13.75
[0] Input token throughput (tok/s): 4310.13
[0] Output token throughput (tok/s): 2673.05
[0] Total token throughput (tok/s): 6983.18
[0] ======
```

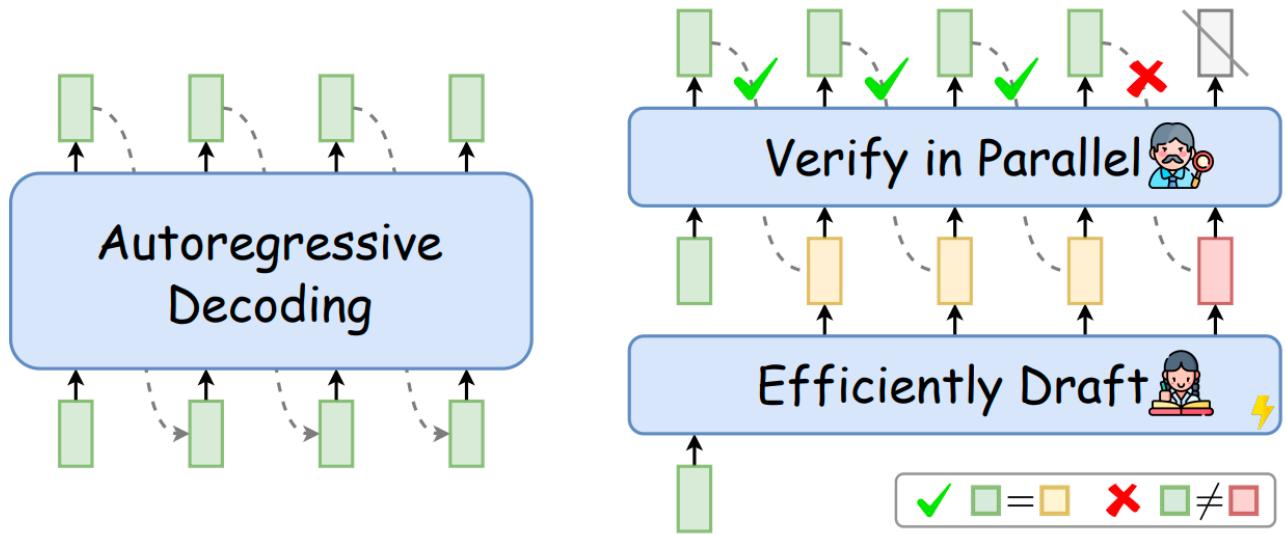


### 3. Speculative Decoding

\*\*

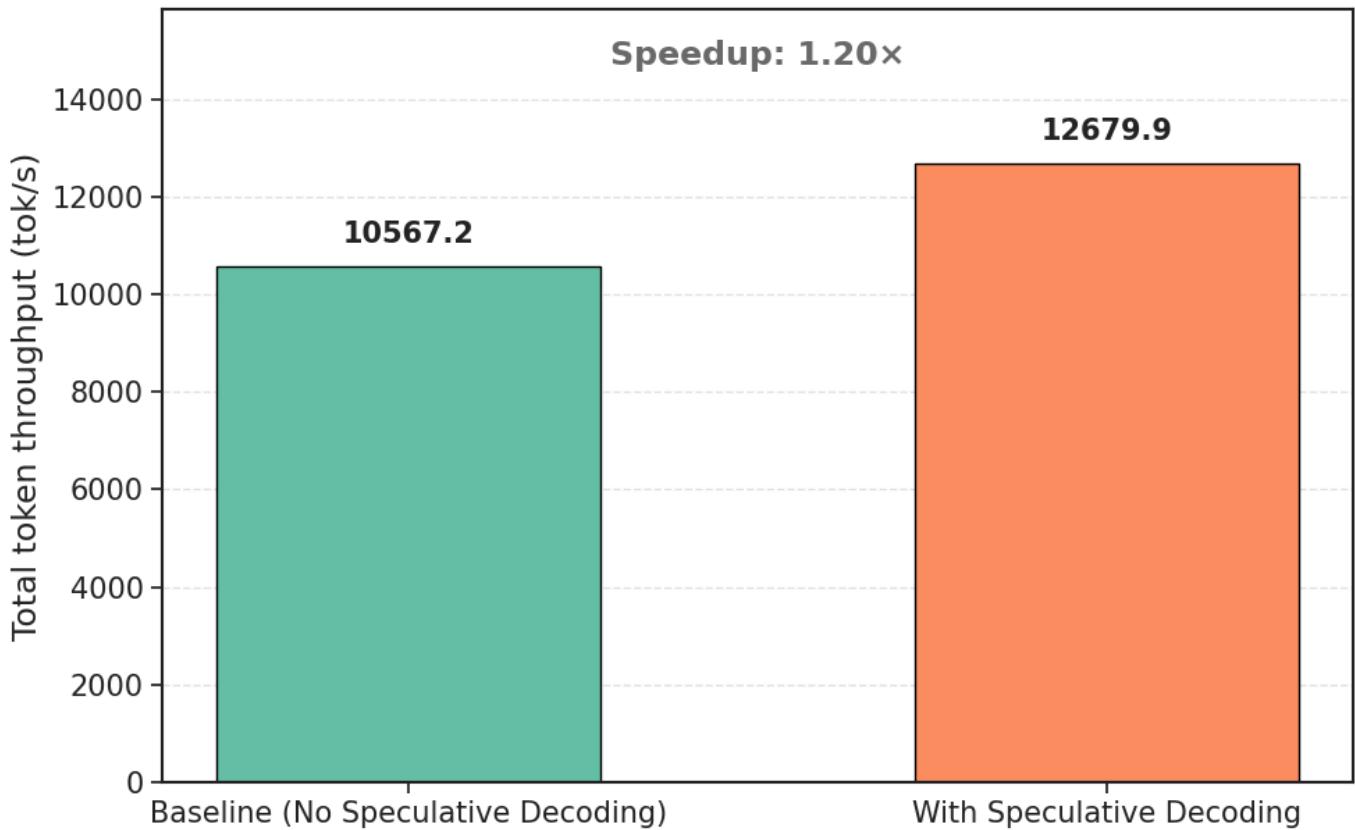
Speculative decoding accelerates language model inference by using a smaller draft model to predict several tokens ahead, which the main model then verifies in parallel. To further improve throughput in the decode stage, we implemented **speculative decoding**, which generates multiple tokens concurrently using a lightweight draft model.

\*\*



**Figure 1:** In contrast to autoregressive decoding (*left*) that generates sequentially, Speculative Decoding (*right*) first *efficiently drafts* multiple tokens and then *verifies* them *in parallel* using the target LLM. Drafted tokens after the bifurcation position (*e.g.*, █) will be discarded to guarantee the generation quality.

## DeepSeek R1 Throughput Comparison (Baseline vs Speculative Decoding)



But Speculative Decoding is not allowed with DP Parallelism!

Abort since not as much gain as DP

## Summary of Results

After applying these optimizations, we observed:

- **Reduced NCCL latency** and improved GPU utilization,
- **Better overlap** between compute and communication phases,
- **Higher decoding throughput**, especially when speculative decoding was enabled, and
- **Improved stability** across different backend configurations.

Overall, our optimized setup achieved **significant end-to-end speedups** over the baseline, demonstrating that **communication-aware parallelism** and **kernel-level backend tuning** are essential for scaling LLM inference efficiently on multi-GPU, multi-node systems.



# Deepseek-R1 Conclusion



## Token Throughput: 10772 -> 17805

**+65.2%**    TP=8, PP=2, DP=2

+10% Tree, LL128

Allreduce\_Sum avg time:  
1.553ms -> 1.066ms

Communication time drops, CUDA Graph launches drops



# Suggestions for Running Deepseek R1



## Notice for Users

- Select TP, PP, DP suits for your hardware topology
  - EP/PD needs large amounts GPU
  - Different backend takes costs to run
  - Select mpi like hpcx

## Suggestion for Developers

- **PP Support for Speculative Decoding**
  - More profiling docs

