

The background of the slide features a large, faint watermark of the Shanghai Jiao Tong University logo. The logo is circular, with a gear-like outer ring. Inside the ring, the university's name is written in Chinese characters at the top and "SHANGHAI JIAO TONG UNIVERSITY" in English at the bottom. The center of the logo depicts a stylized building and the year "1896" in a rectangular box.

Benchmarking Huawei ARM Server Processor for HPC Workloads

James Lin

Shanghai Jiao Tong University, Center for HPC

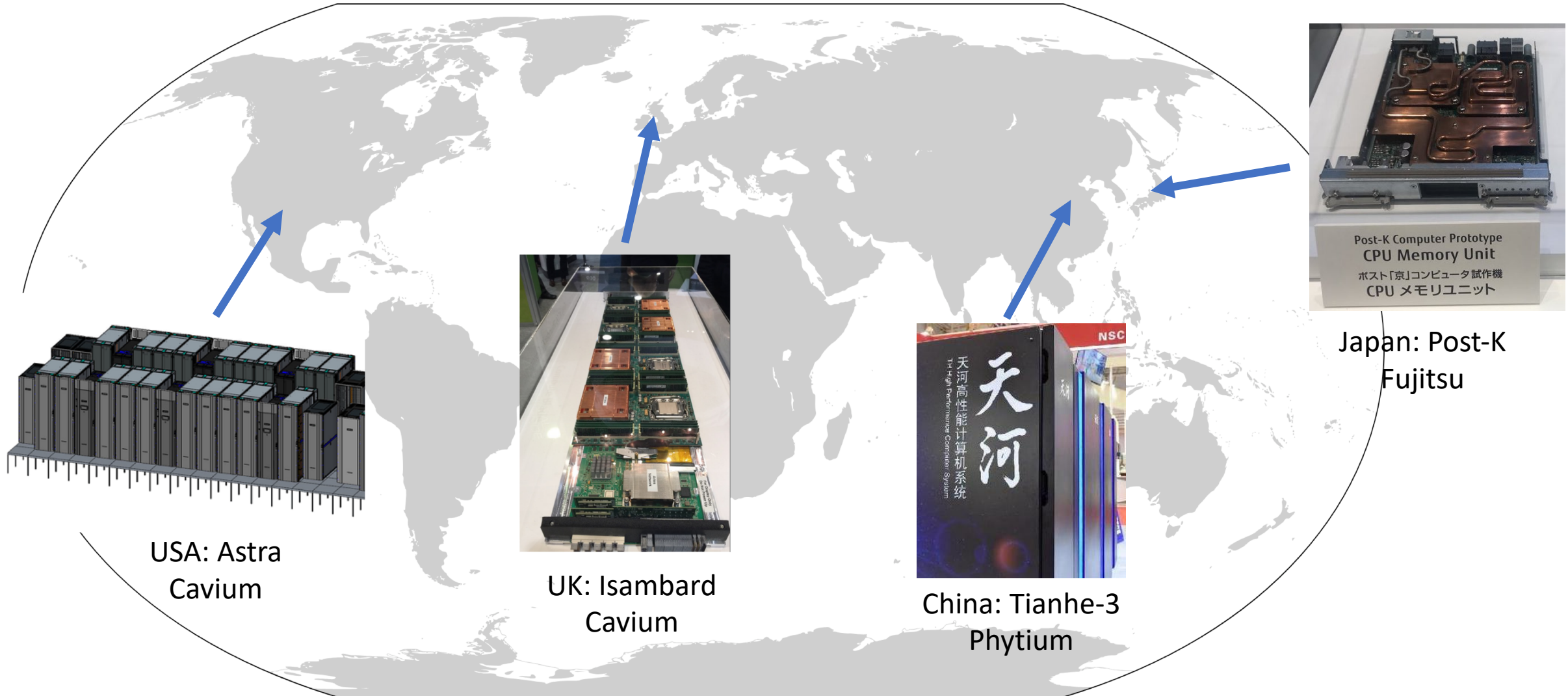
ARM HPC User Group, Dallas, US
November , 2018

Outline

- 
- ARM HPC in China
 - SJTU ARM HPC Innovation Center
 - Hi1620, the new ARM Server processor from Huawei
 - Performance Results
 - Micro-benchmark
 - 5 Mini-apps
 - A real-world app: GTC-P

The Rising of ARM HPC

ARM-based Supercomputers in the world



The ARM Ecosystem in China

OpenGCC -- “Green Computing Consortium”



- Supported by China's government
- Promoting green computing industry in China with the leading universities and ecosystem partners.
- <http://opengcc.org>





OpenACC-SJTU ARM HPC Innovation Center: June 2018



- So far the only ARM HPC research center in China
- Supported by the government

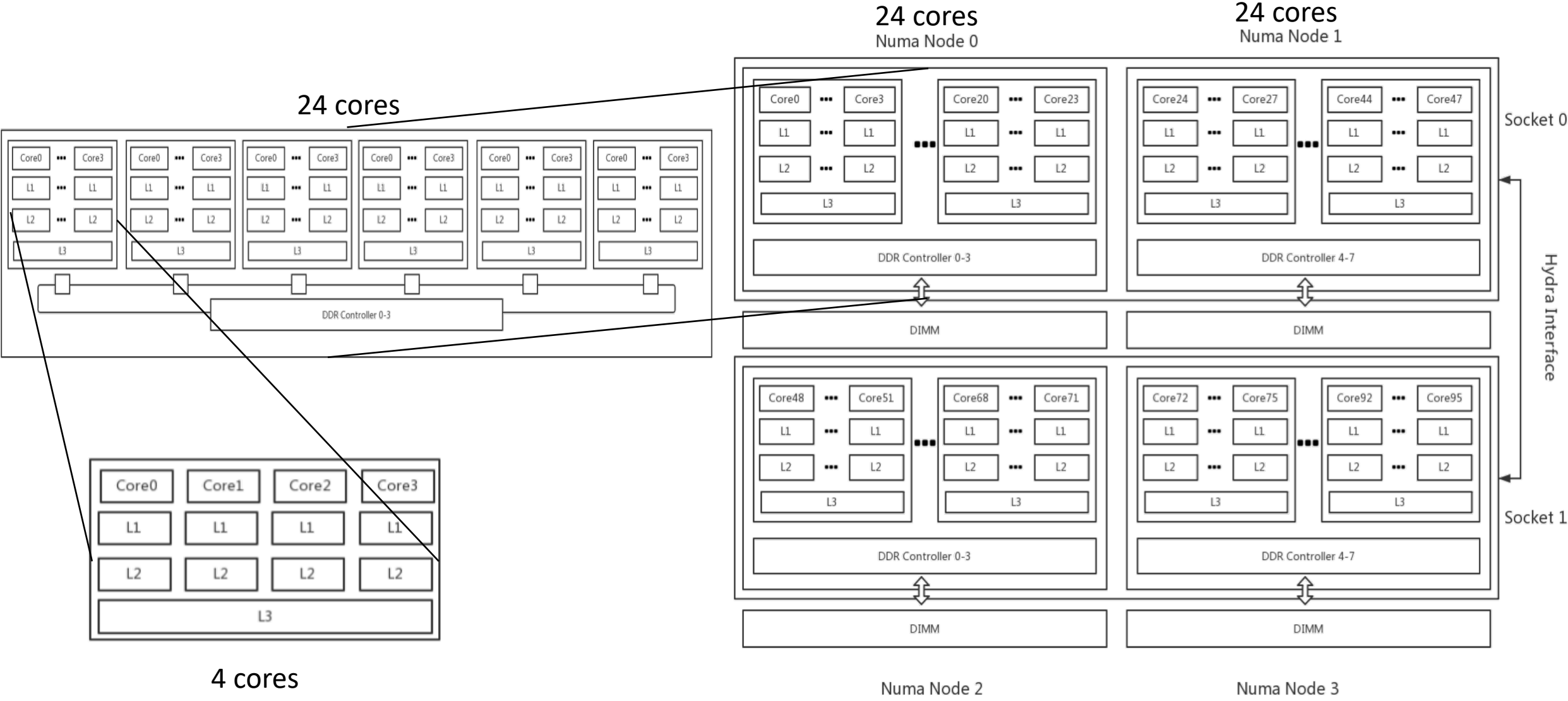
The first ARM HPC workshop in China: July 2018



Outline

- ARM HPC in China
 - SJTU ARM HPC Innovation Center
-  • Hi1620, the new ARM Server processor from Huawei
- Performance Results
 - Micro-benchmark
 - 5 Mini-apps
 - A real-world app: GTC-P

The architecture of the 48-core Hi1620 (the first time to reveal!)



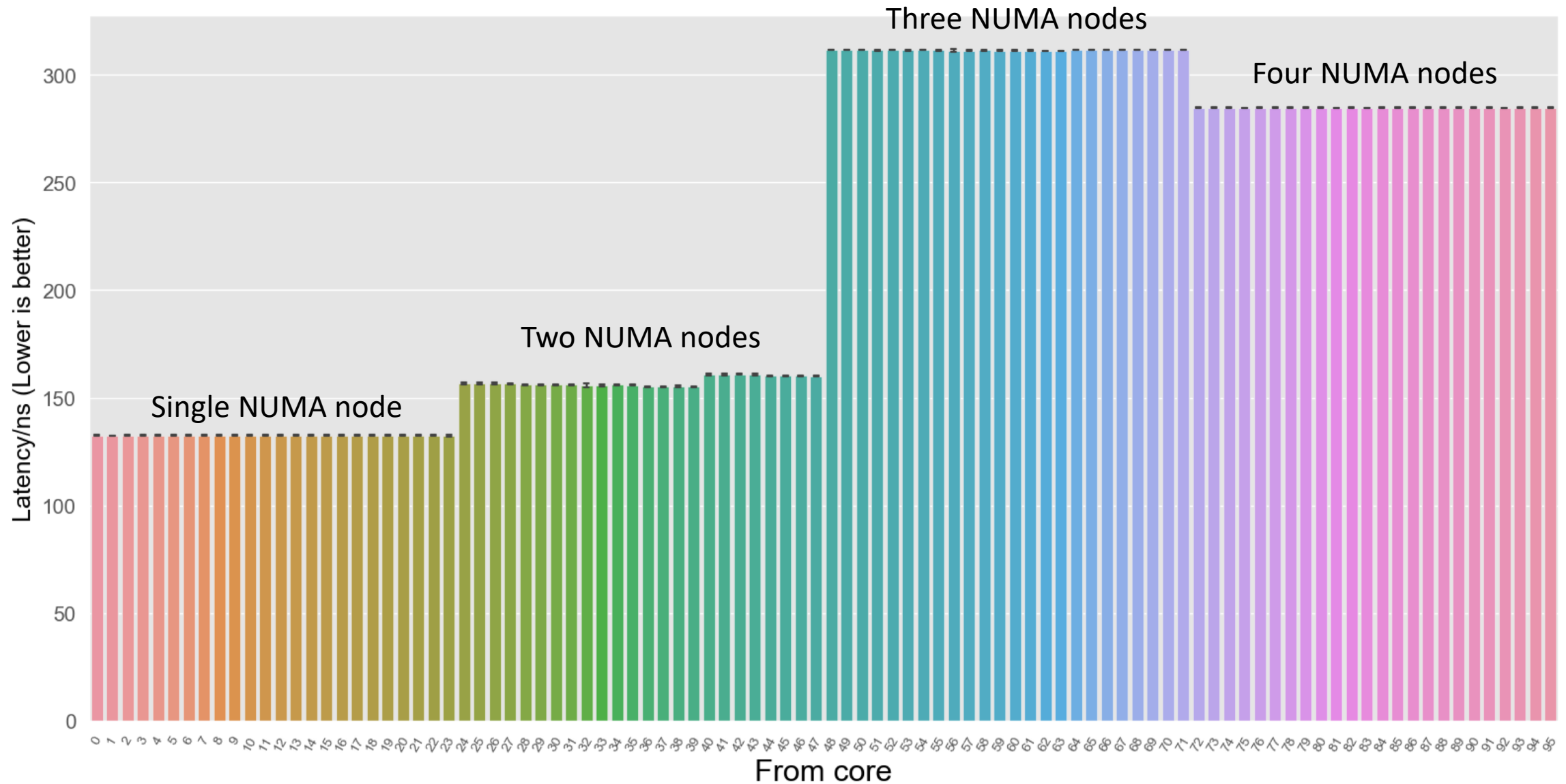
X86 V.S. ARM

| | Intel Xeon E5-2680v3 | Intel Xeon Gold 6148 | Huawei Hi1616 | Huawei Hi1620 (EP) |
|-----------------------|------------------------|---------------------------|-------------------------|--------------------------------------|
| Code name | Haswell | Skylake | A72 | Customized |
| Lithography | 22nm | 14nm | 16nm | 7nm |
| Frequency | 2.50 GHz | 2.40 GHz | 2.40 GHz | 2.00GHz (Engineering Chip) |
| Cores | 12 | 20 | 32 | 48 |
| Vectorization Length | 256-bit x86-64 AVX2 | 512-bit x86-64 AVX-512 | 128-bit AArch64 NEON | 128-bit AArch64 NEON |
| Peak Performance (DP) | 480 Gflops | 1536 Gflops | 307.2 Gflops | 384 Gflops (should be 768 Gflops) |
| L3 Cache | 30MB | 27.5MB | 32MB | 32MB |
| DRAM _{max} | 4-Channel DDR4-2133 | 6-Channel DDR4-2666 | 4-Channel DDR4-2400 | 8-Channel DDR4-3200 |
| Max Memory Bandwidth | 68.3 GB/s | 127.5 GB/s | 76.8 GB/s | 170.6 GB/s (DDR4-2666) |
| TDP | 120W | 150W | 70W | ??? |
| Launch Year | 2014 | 2017 | 2016 | 2019 |

Outline

- ARM HPC in China
 - SJTU ARM HPC Innovation Center
- Hi1620, the new ARM Server processor from Huawei
-  • Performance Results
 - Micro-benchmark
 - 5 Mini-apps
 - A real-world app: GTC-P

Micro-benchmark: memory access latency via Imbench



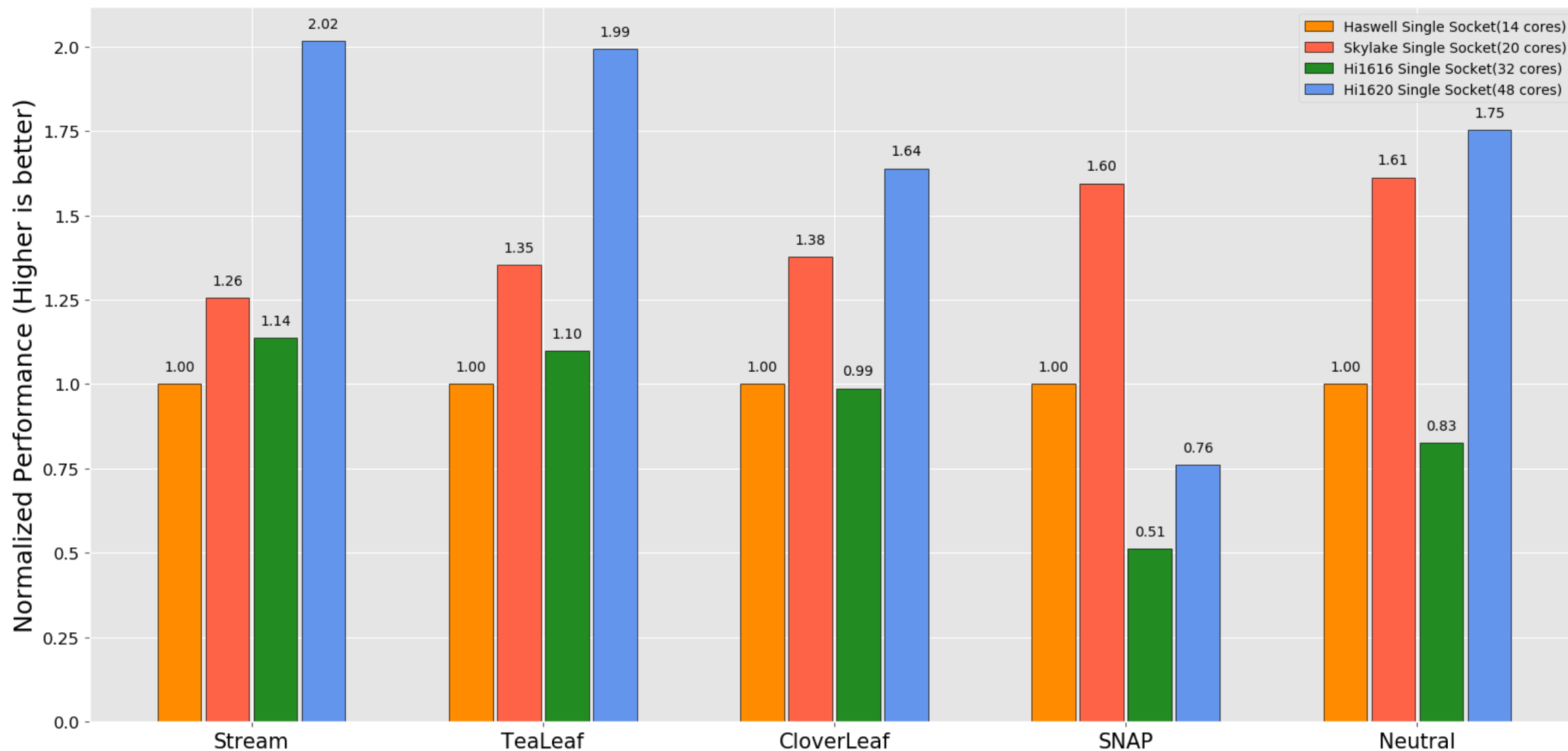
5 mini-applications evaluated on Hi1620

- **CloverLeaf**
 - Solves Euler's equations of compressible fluid dynamics
 - Stencil code
 - Memory bandwidth-bound
- **TeaLeaf**
 - Solves the linear heat conduction equation
 - Stencil code
 - Memory bandwidth-bound
- **SNAP**
 - A proxy application for a modern deterministic discrete ordinates transport code
 - The performance of the memory traffic to and from cache is a key performance factor
- **Neutral**
 - A simplified Monte Carlo neutral particle transport mini-app
 - A mesh-based approach
 - Limited by the performance of memory latency

Thunder X2 results: <http://uob-hpc.github.io/2018/05/23/CUG18.html>

Performance Results of 5 mini-applications

Haswell: 2680V3
Skylake:6148

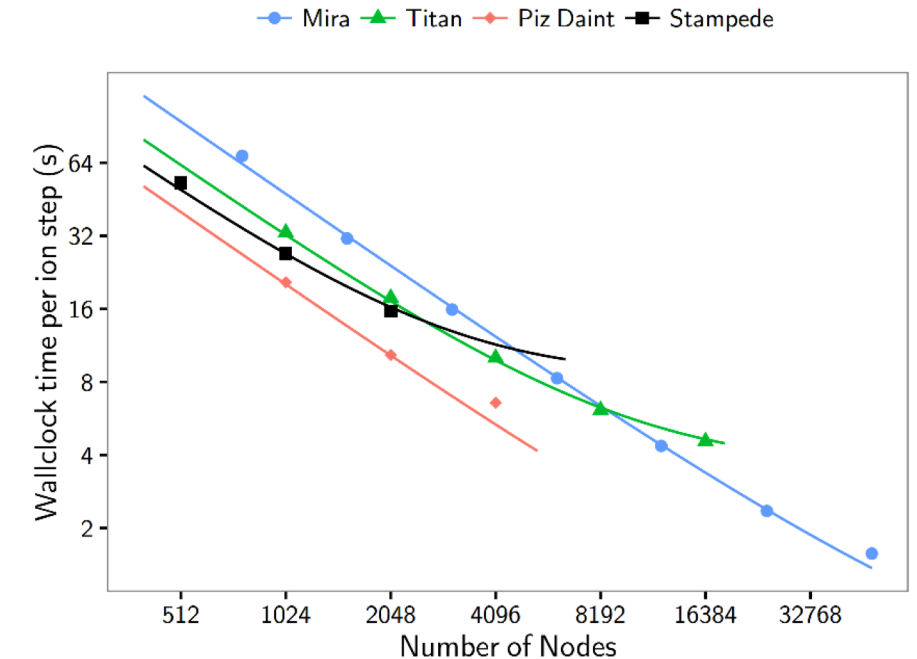
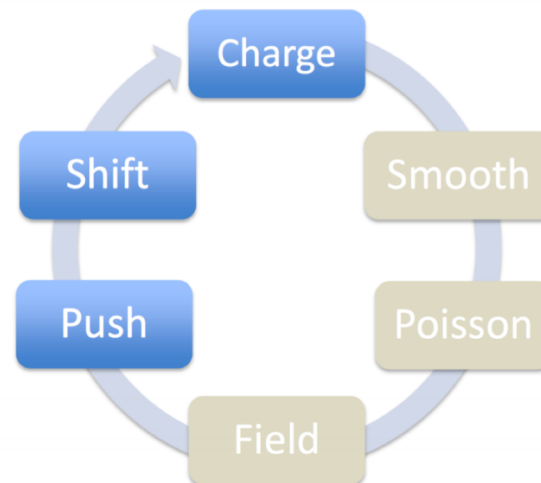
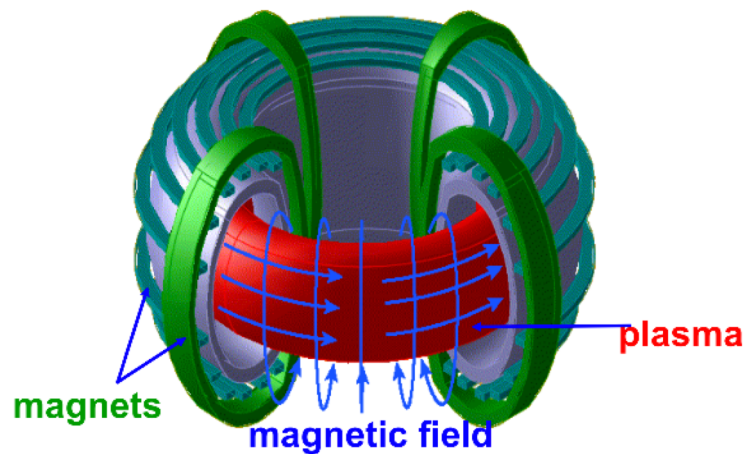


A Real-world application GTC-P: Gyrokinetic Toroidal Code - Princeton

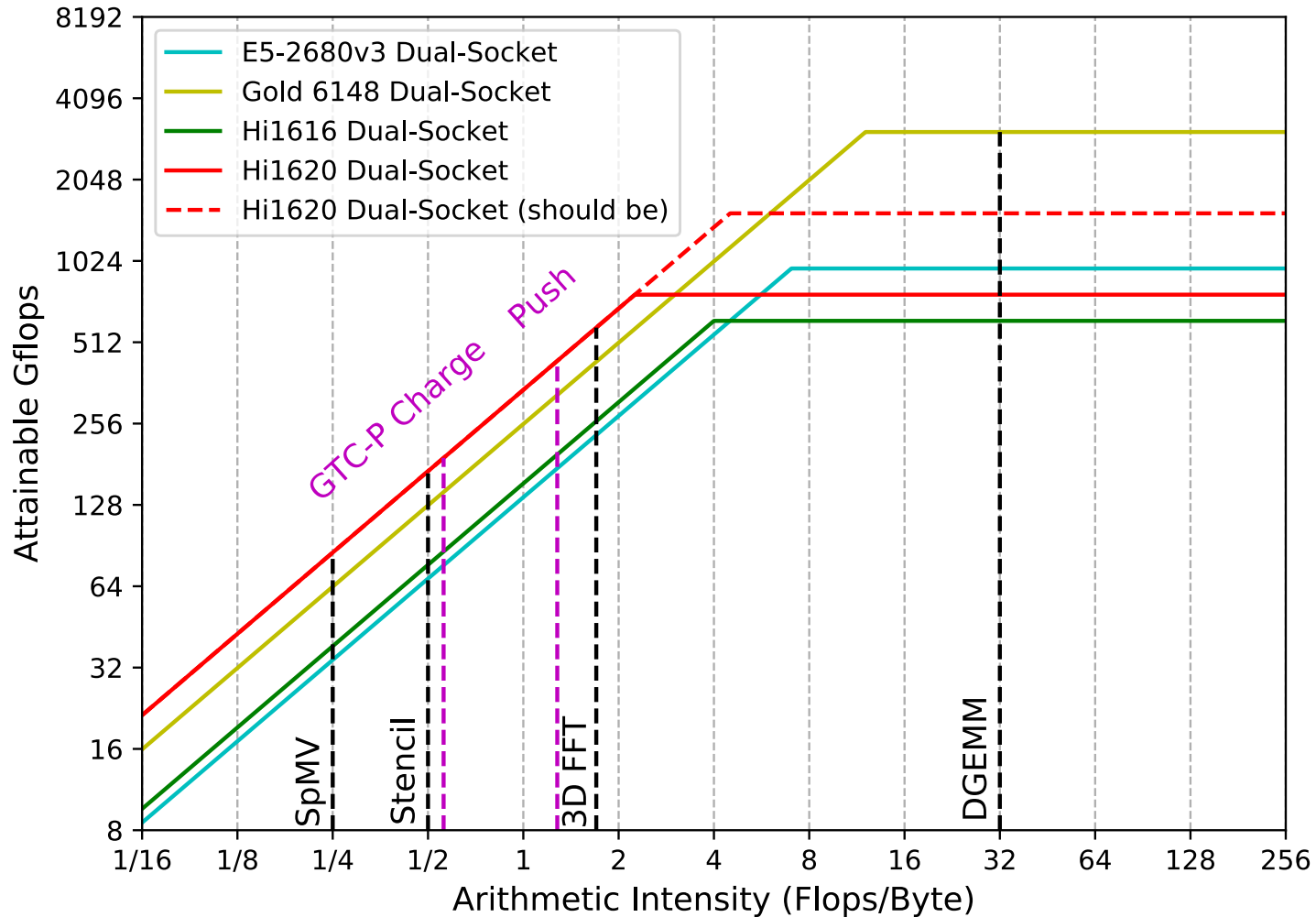


Supported by
NSF SAVI Project

GTC-P is Particle-in-Cell code that delivers fusion simulations at extreme scales on the worldwide supercomputers including *Tianhe-2*, *Titan*, *TaihuLight* and etc., that feature *CPU*, *GPU* and *many-core processors*.



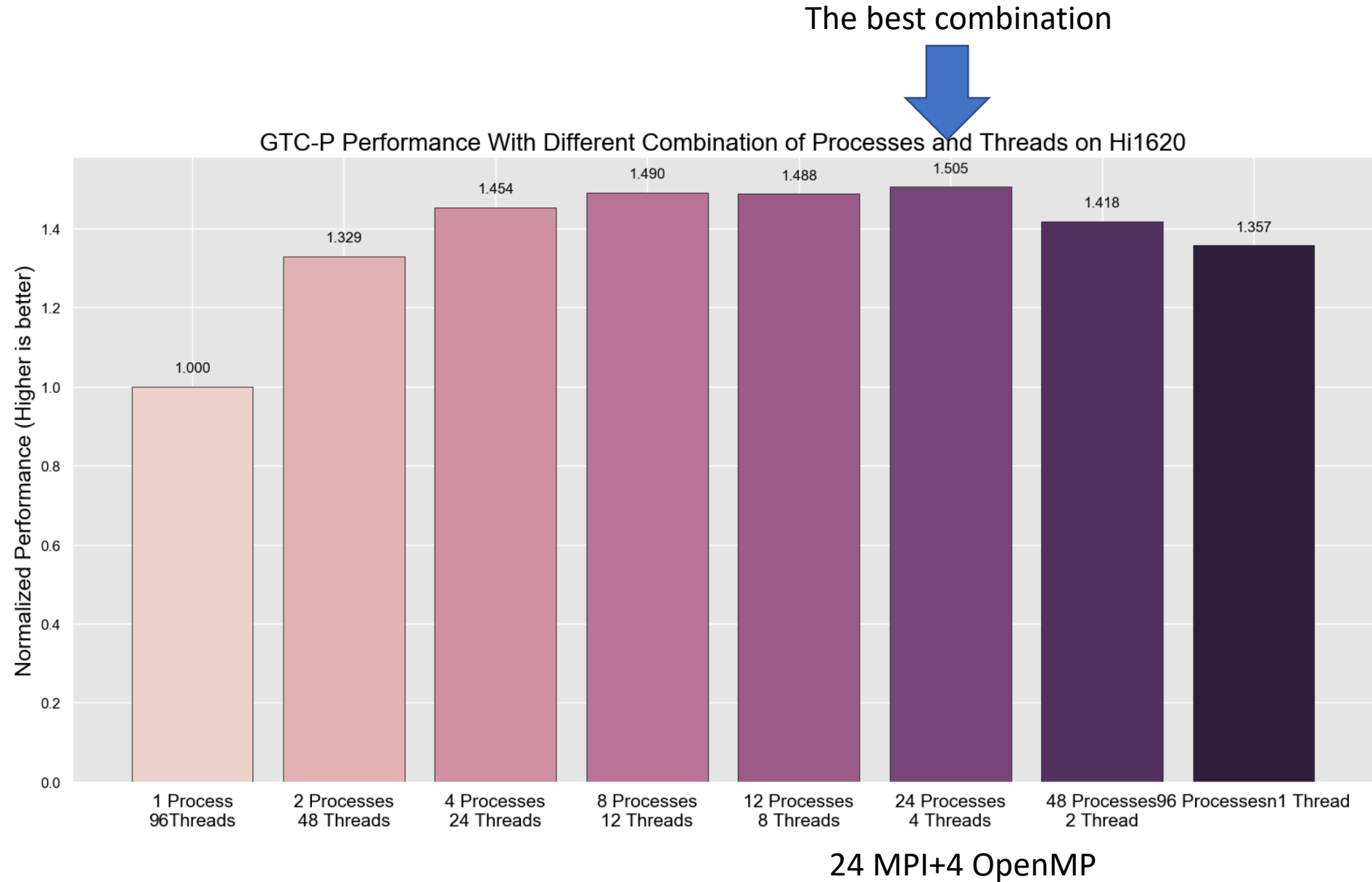
Roofline model for GTC-P code



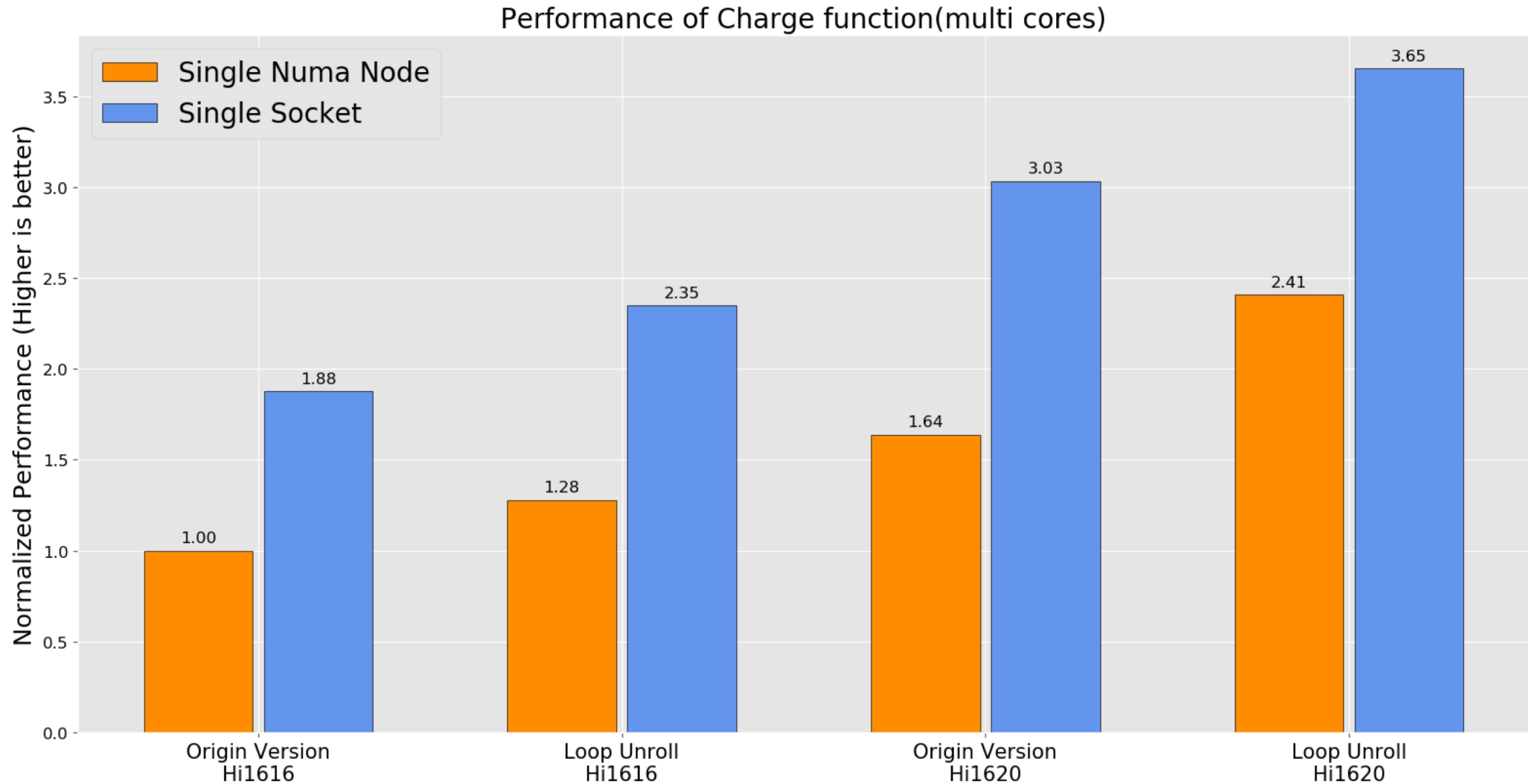
- GTC-P code with 6 main functions including 2 main hotspots charge (SCATTER) and push (GATHER) kernels.
- It is a memory-bound code with irregular memory access.
- The code base is the MPI + OpenMP version of GTC-P code. We compile the code with GCC 8.2 compiler.

Evaluating the thread affinity on Hi1620

- Each core own private L1 and L2 cache.
- 4-core is in one core cluster, L3 is shared.
- 24-core is on the same ring.
- Single socket is dual-ring.



Loop unrolling optimization



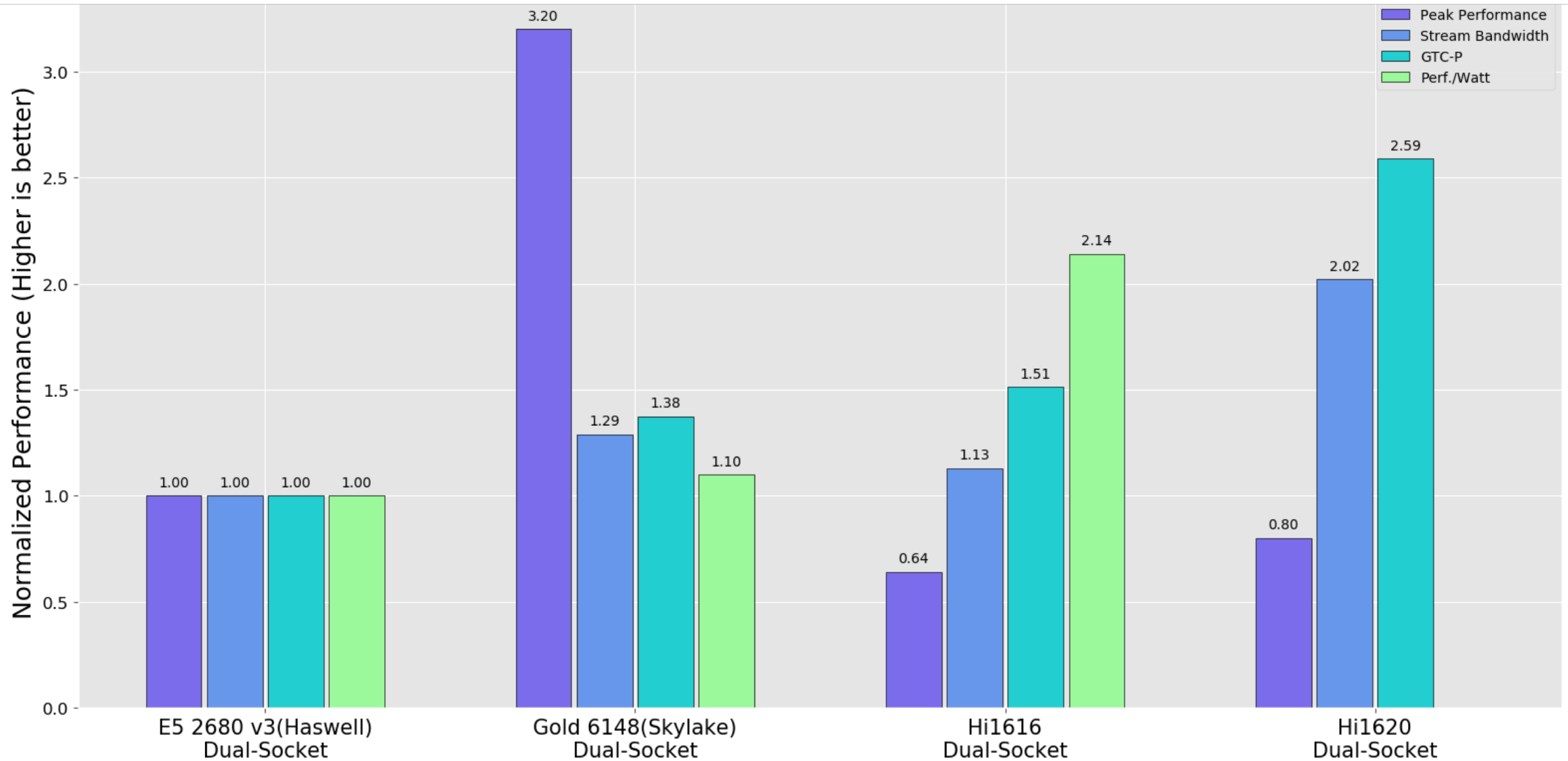
Vectorization on Hi1616 and Hi1620

```
.macro tkfma
    fmla v0.2d, v16.2d, v16.2d
    fmla v1.2d, v17.2d, v17.2d
    fmla v2.2d, v18.2d, v18.2d
    fmla v3.2d, v19.2d, v19.2d
    fmla v4.2d, v20.2d, v20.2d
    fmla v5.2d, v21.2d, v21.2d
    fmla v6.2d, v22.2d, v22.2d
    fmla v7.2d, v23.2d, v23.2d
    fmla v8.2d, v24.2d, v24.2d
    fmla v9.2d, v25.2d, v25.2d
    fmla v10.2d, v26.2d, v26.2d
    fmla v11.2d, v27.2d, v27.2d
    fmla v12.2d, v28.2d, v28.2d
    fmla v13.2d, v29.2d, v29.2d
    fmla v14.2d, v30.2d, v30.2d
    fmla v15.2d, v31.2d, v31.2d
.endm
```

- Hi1616
 - 128-bit SIMD
 - SP: 614.4 Gflops
 - DP: 307.2 Gflops
- Hi1620
 - 128-bit SIMD
 - SP: 1,536 Gflops
 - DP: 384 Gflops
- The throughput of DP vectorization is decreased to only half on Hi1620

| FMA Instruction Throughput | | | | |
|----------------------------|------------------------|------------------------|-------------------------|------------------------|
| | SP Scalar | DP Scalar | SP Vector | DP Vector |
| Hi1616 | 2ins/cycle 9.596GFlops | 2ins/cycle 9.596GFlops | 1ins/cycle 19.194Gflops | 1ins/cycle 9.596GFlops |
| Hi1620 | 2ins/cycle 7.989Gflops | 2ins/cycle 7.989Gflops | 2ins/cycle 31.954GFlops | 1ins/cycle 7.989Gflops |

Overall performance results comparison on multiple processors



Summary

- The experiments demonstrate the **portability** on ARM-based system with the support of ARM HPC ecosystem.
- Early results show HPC workloads, especially memory-bound applications, on Huawei ARM multi-core processors is competitive with the same generation's x86 CPUs, while **performance per watt** is compelling, which demonstrate its potential of “**Green Computing**”.
- The **thread affinity** is important for the application performance. It is strongly recommended to use MPI+OpenMP to run the code on dual-socket machine.

SJTU ARM HPC Research Team

