SC2O

Everywhere | more
we are | than hpc.

[12 Nov 2020] • [AoE]

# Performance Modeling of Streaming Kernels and Sparse Matrix-Vector Multiplication on A64FX

Christie Alappat, Jan Laukemann, Thomas Gruber, Georg Hager, Gerhard Wellein, Nils Meyer, Tilo Wettig

TOP 500 The List.

# NEWS
# Japan Captures TOP500 Crown with Arm-Powered Supercomputer
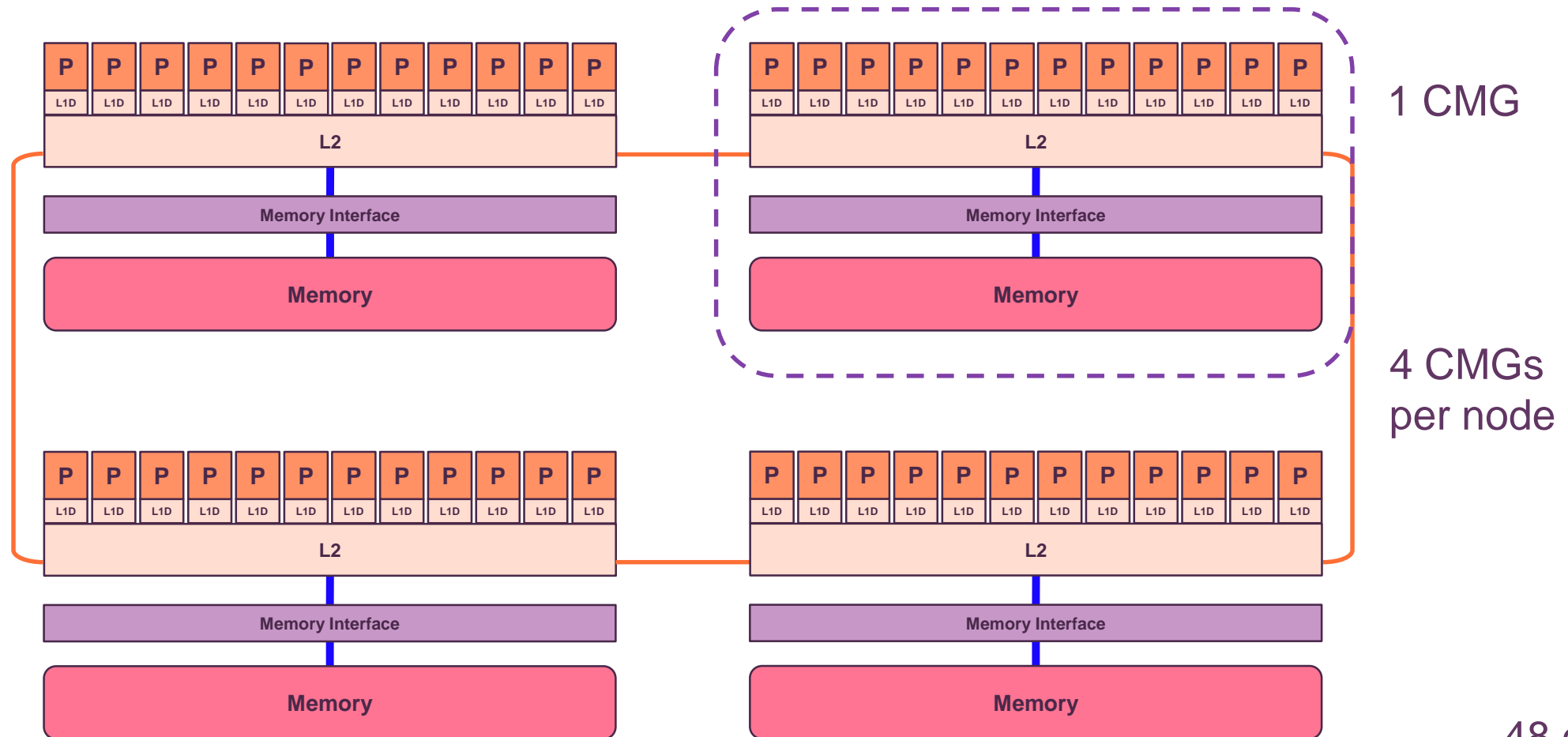
June 22, 2020

FRANKFURT, Germany; BERKELEY, Calif.; and KNOXVILLE, Tenn.—The 55th edition of the TOP500 saw some significant additions to the list, spearheaded by a new number one system from Japan. The latest rankings also reflect a steady growth in aggregate performance and power efficiency.
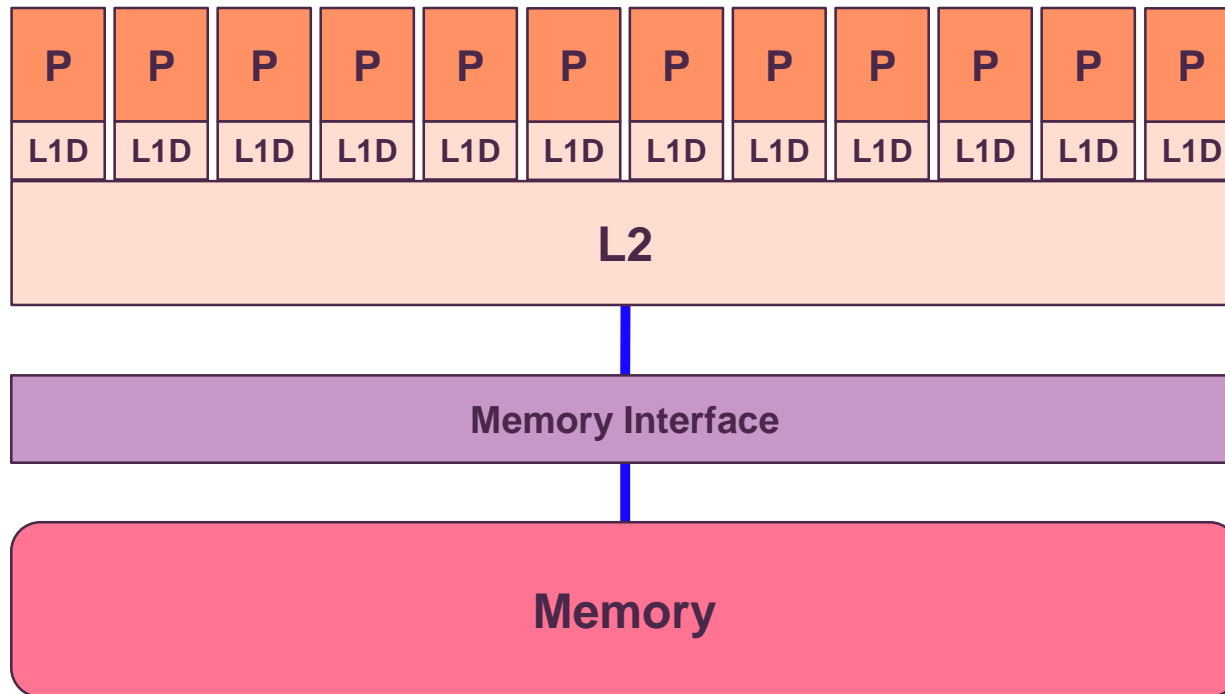
The new top system, Fugaku, turned in a High Performance Linpack (HPL) result of 415.5 petaflops, besting the now second-place Summit system by a factor of 2.8x. Fugaku, is powered by Fujitsu's 48-core A64FX SoC, becoming the first number one system on the list to be powered by ARM processors. In single or further reduced precision, which are often used in machine learning and AI applications, Fugaku's peak performance is over 1,000 petaflops (1 exaflops). The new system is installed at RIKEN Center for Computational Science (R-CCS) in Kobe, Japan.
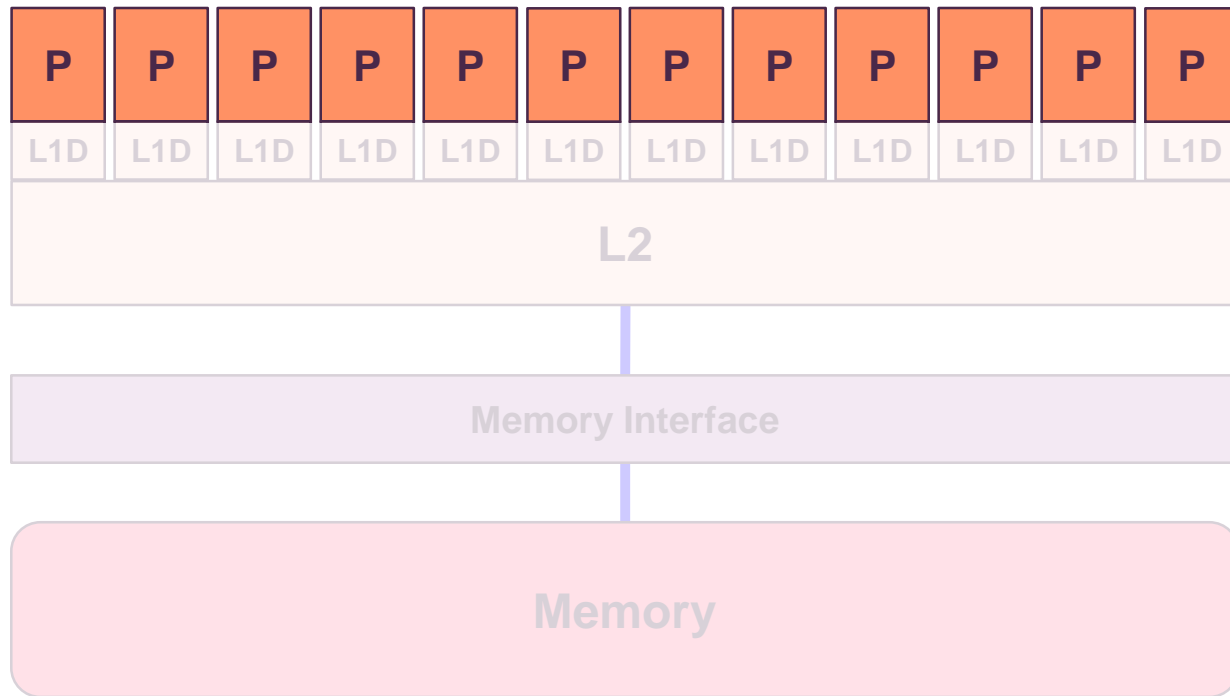


| JUNE 2020 | SYSTEM | SPECS |
|---|---|---|
| 1 | Fugaku | Fujitsu A64FX (48C, 2.2GHz), Tofu Interconnect D |
| 2 | Summit | IBM POWER9 (22C, 3.07GHz), NVIDIA Volta GV100 (80C), Dual-Rail Mellanox EDR Infiniband |
| 3 | Sierra | IBM POWER9 (22C, 3.1GHz), NVIDIA Tesla V100 (80C), Dual-Rail Mellanox EDR Infiniband |
| 4 | Sunway TaihuLight | Sunway SW26010 (260C, 1.45 GHz) Custom Interconnect |
| 5 | Tianhe-2A (Milkyway-2A) | Intel Ivy Bridge (12C, 2.2 GHz) & TH Express-2, Matrix 2000 |

Source : https://www.top500.org/news/japan-captures-top500-crown-arm-powered-supercomputer/

1 CMG

4 CMGs
per node

48 cores
per node

1 CMG

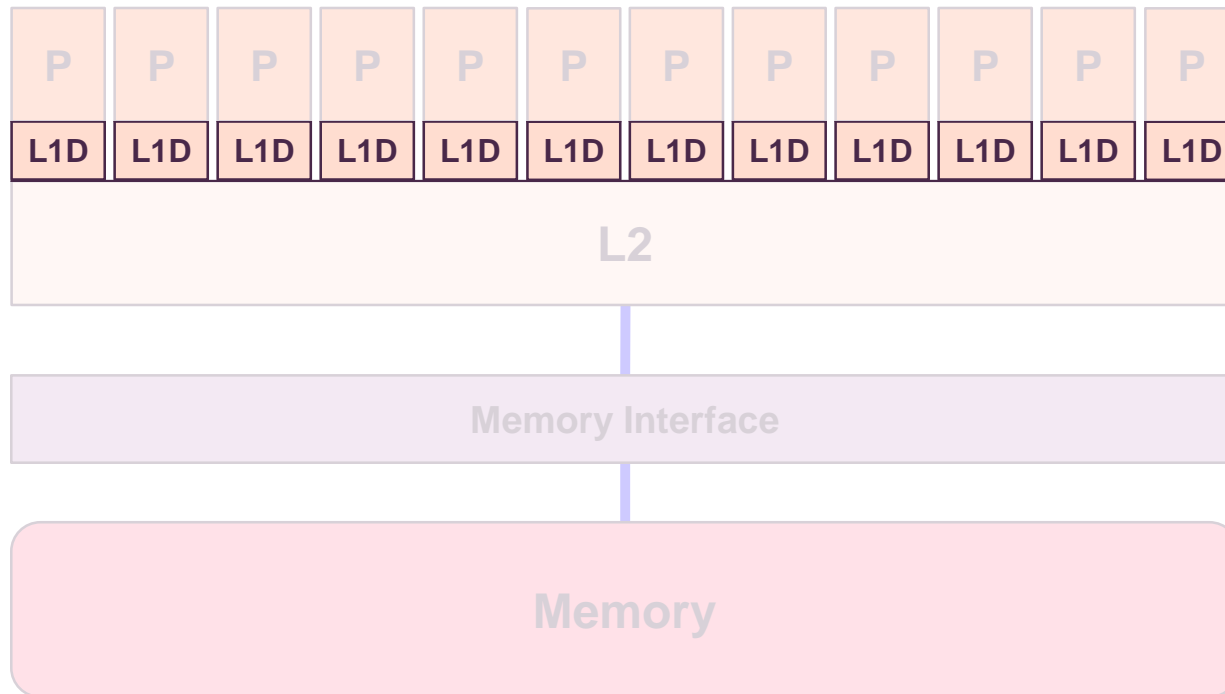| P | P | P | P | P | P | P | P | P | P | P | P |
|---|---|---|---|---|---|---|---|---|---|---|---|
| L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D |

**L2**

**Memory Interface**

**Memory**

1 CMG

Core

- Clock : 1.8 GHz

- Instruction set : Armv8.2-A+SVE

- Maximum VL : 512 bit (8 double)

For example on GCC compiler use :

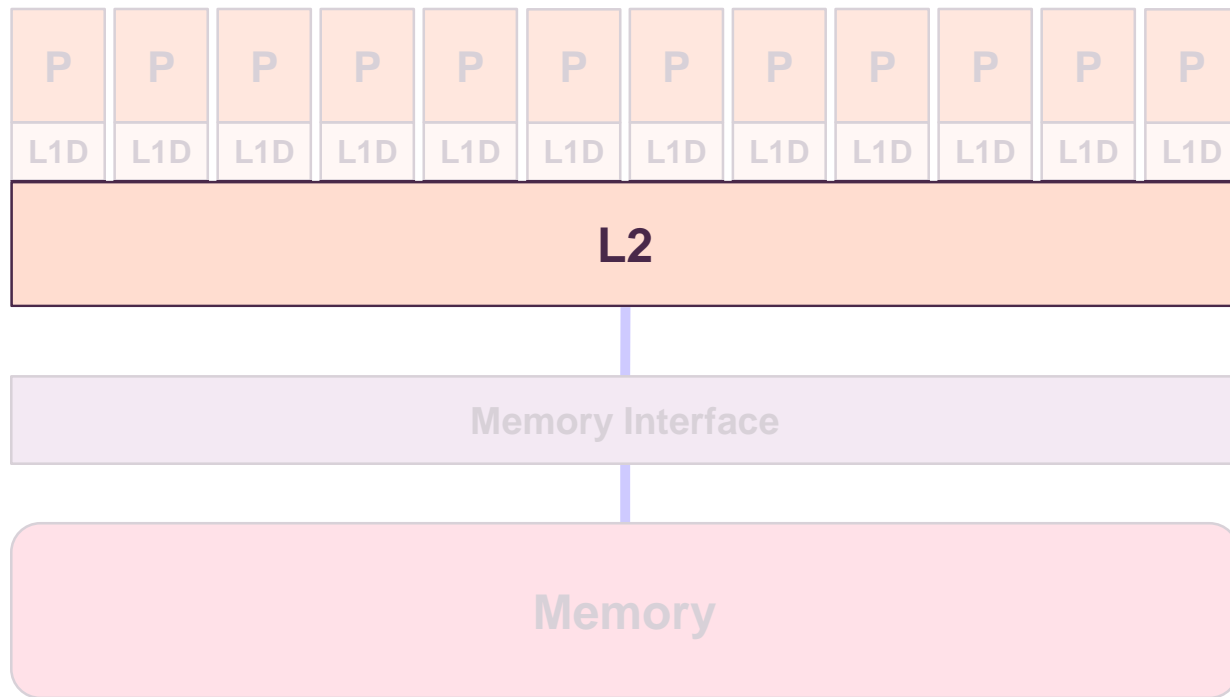-msve-vector-bits=512 -march=armv8.2-a+sve

# A64FX – FX700

| P | P | P | P | P | P | P | P | P | P | P | P |
| L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D | L1D |

**L2**

**Memory Interface**

**Memory**

## 1 CMG

### L1D cache

- Size : 64 KiB
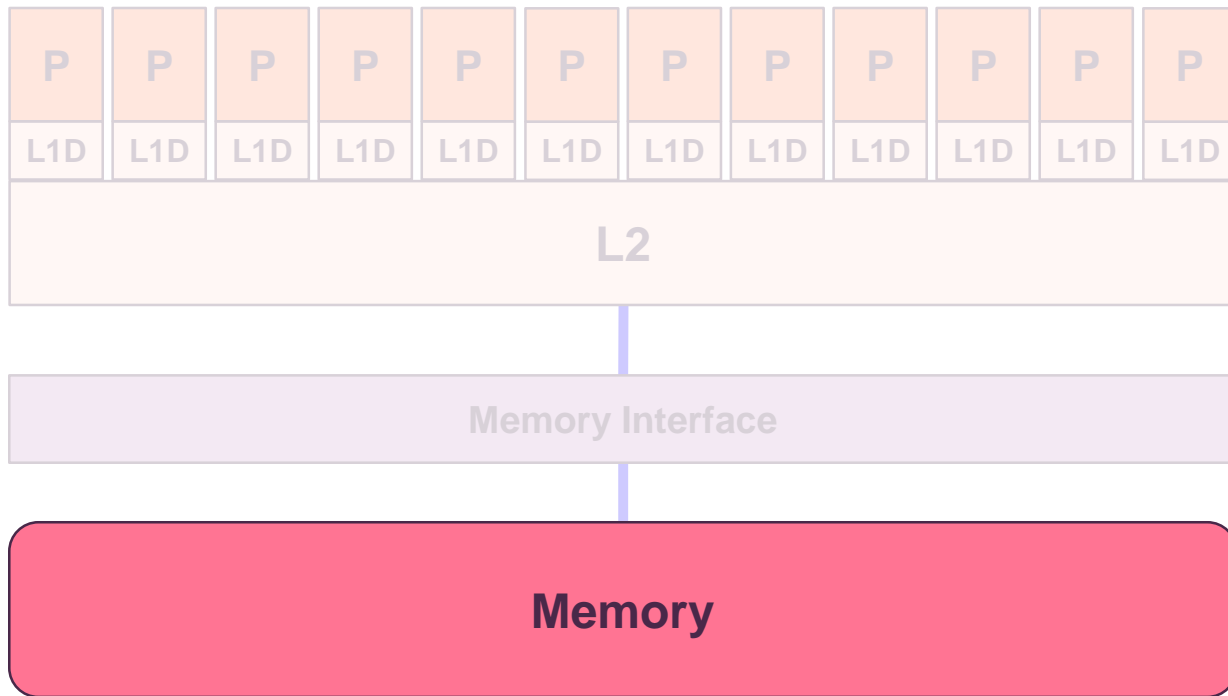
- Topology : Private cache

- Cache line size : 256 bytes

**1 CMG**

L2 cache

- Size : 8 MiB

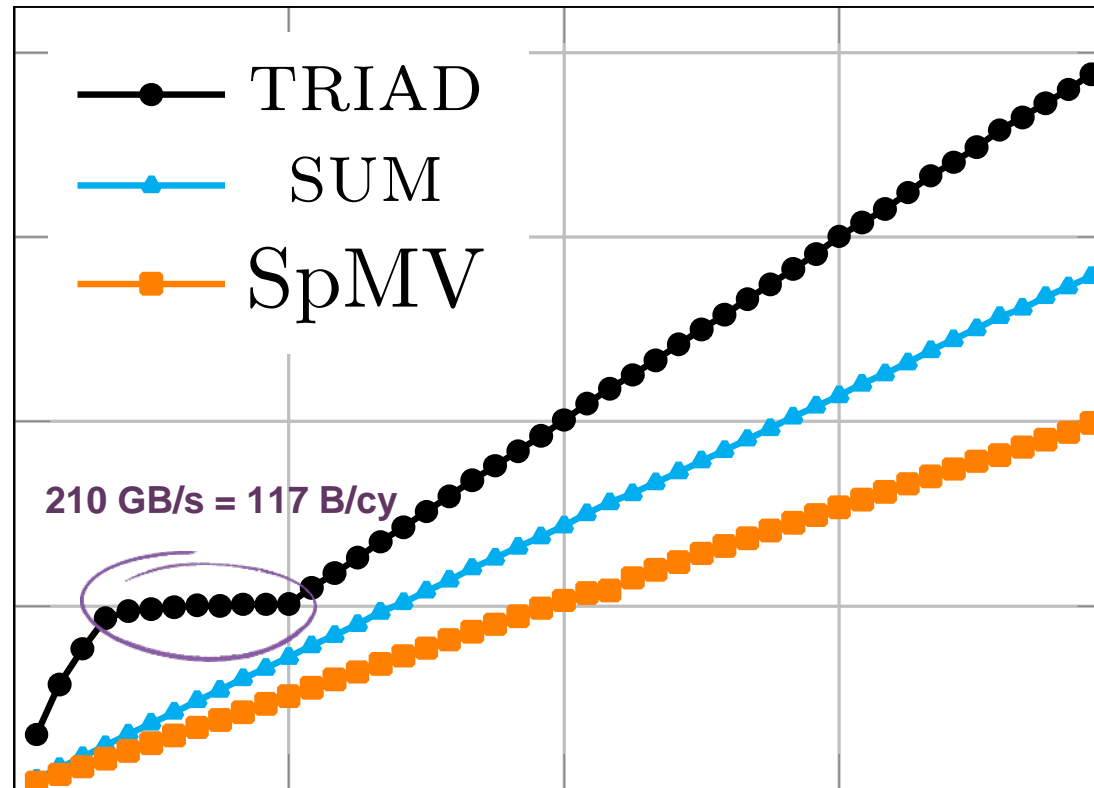- Topology : Shared within 1 CMG

- Cache line size : 256 bytes

1 CMG

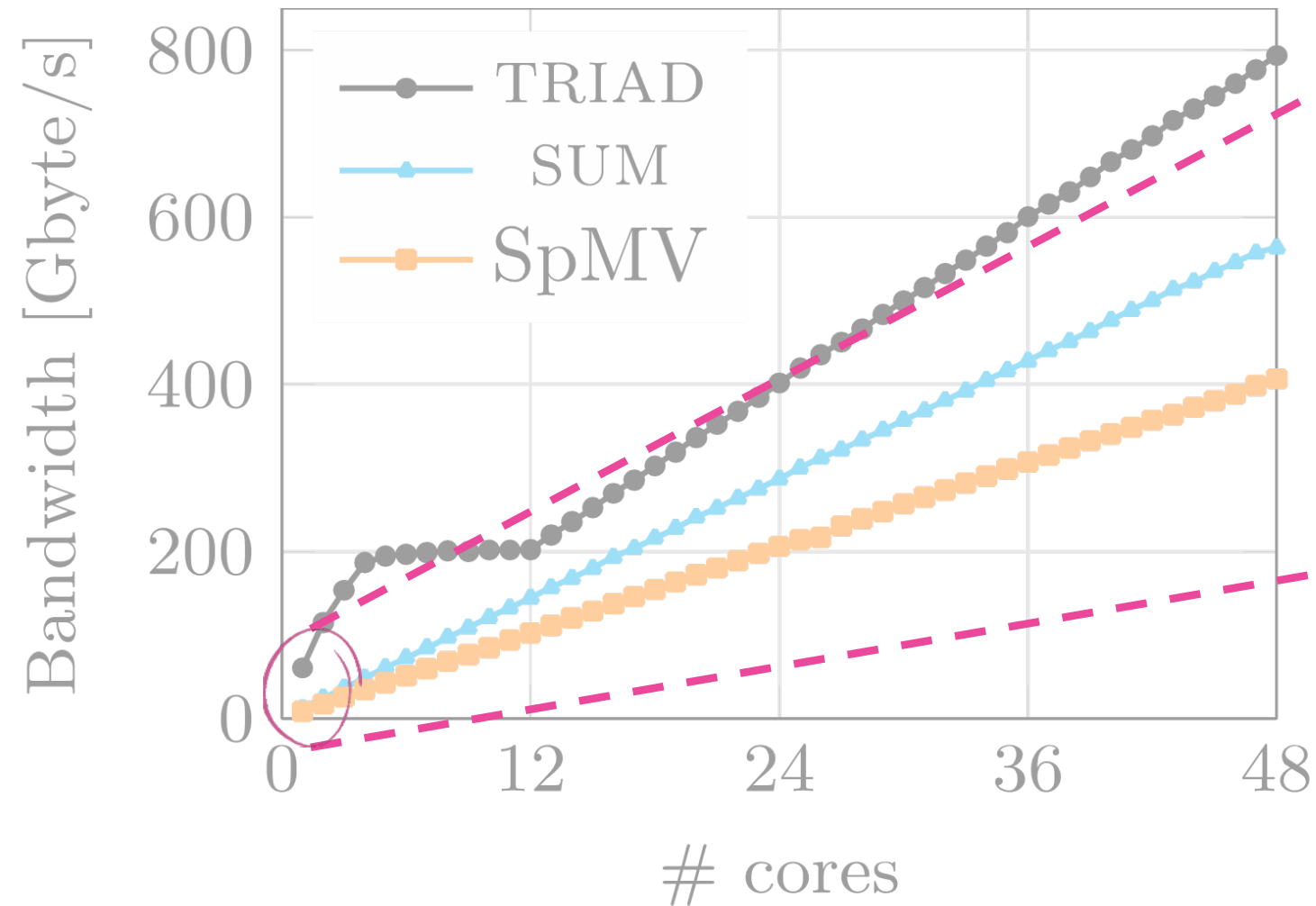Main memory

- Size : 4 x 8 GiB

- Type : HBM2

dwidth [Gbyte/s]

210 GB/s = 117 B/cy

TRIAD
SUM
SpMV

Thread pinning : Compact

Clear memory bandwidth saturation for STREAM TRIAD (`a[i] = b[i] + s*c[i]`).

But why not for SUM (`s += a[i]`) and SpMV (`b = Ax`)?

Thread pinning : Compact
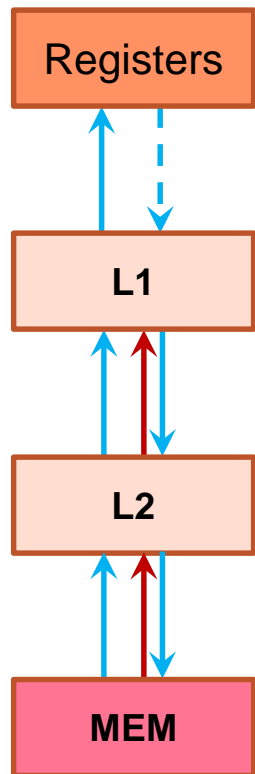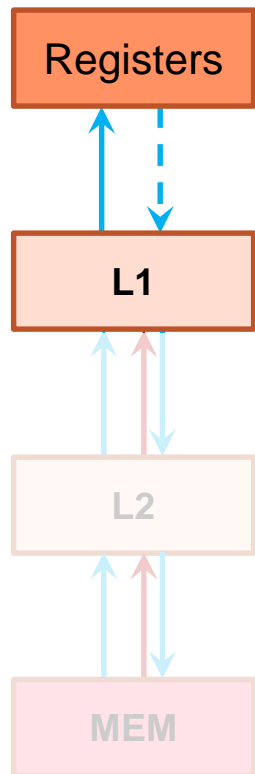
Understanding single-core performance is the key !

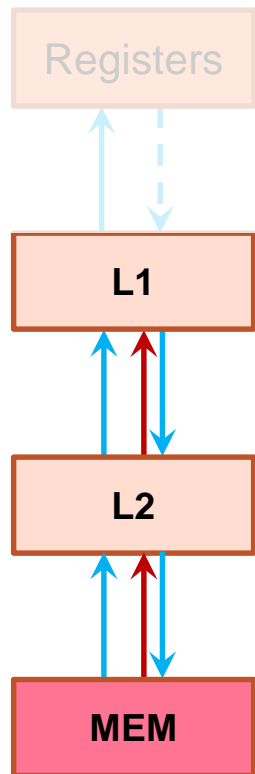Execution-Cache-Memory (ECM) model helps us to understand and analyze the single-core performance.

Execution-Cache-Memory (ECM) model helps us to understand and analyze the single-core performance.

3 major components :

1) In-core

Registers

L1

L2

MEM
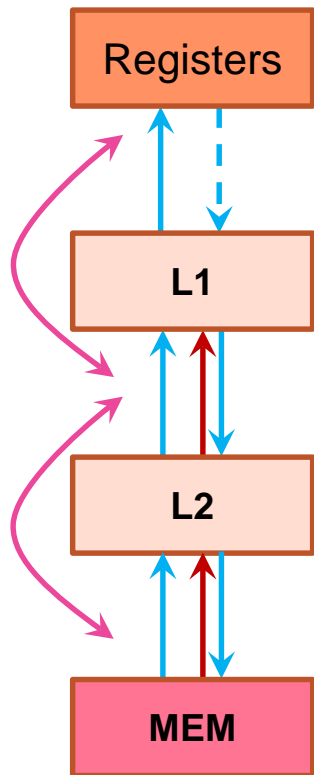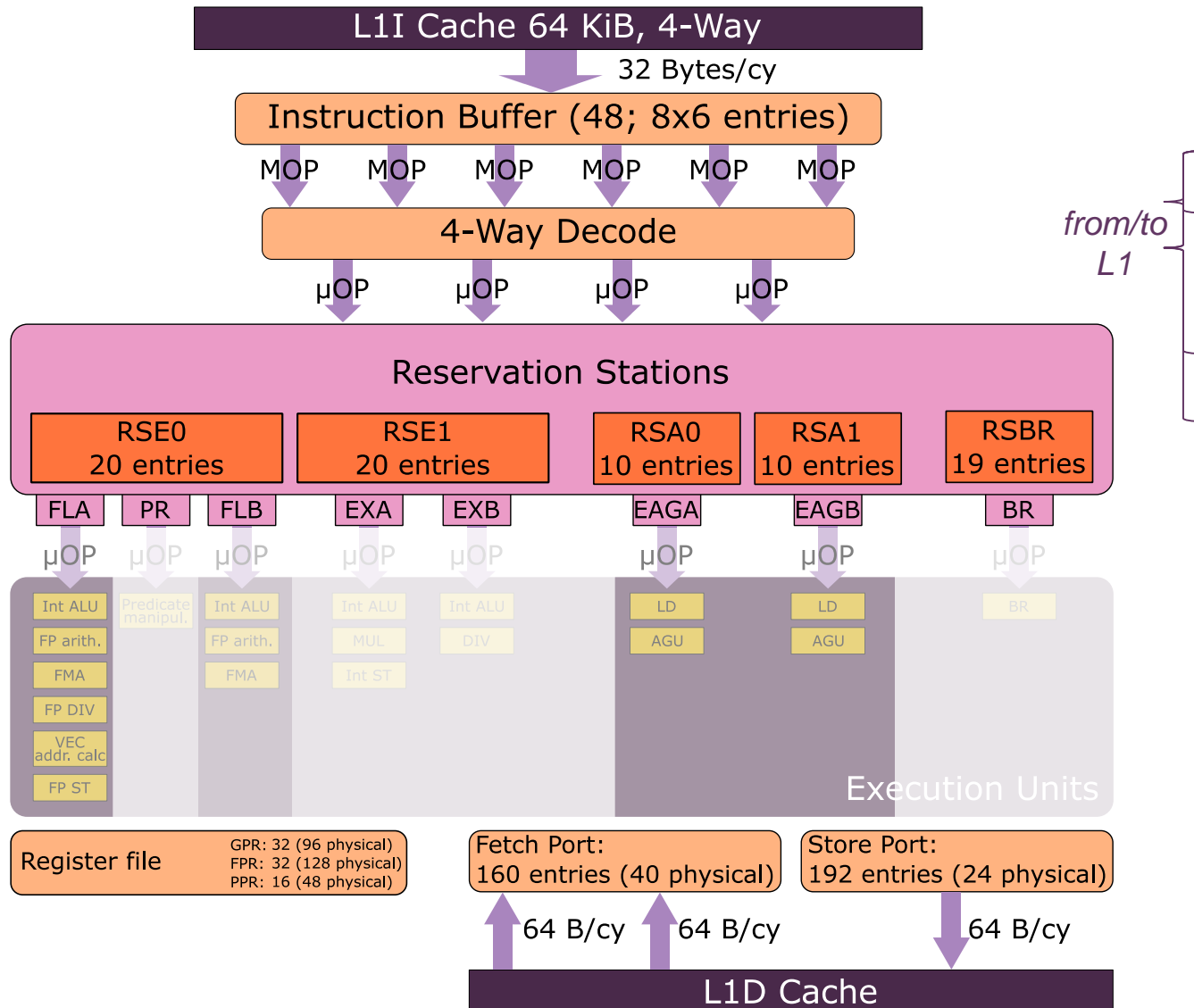
Execution-Cache-Memory (ECM) model helps us to understand and analyze the single-core performance.

3 major components :

1) In-core

2) Data transfer through memory hierarchies

Registers

L1

L2

MEM

Can these transfers be overlapped or not ?

Execution-Cache-Memory (ECM) model helps us to understand and analyze the single-core performance.

3 major components :

1) In-core

2) Data transfer through memory hierarchies

3) Overlap hypothesis

L1I Cache 64 KiB, 4-Way

32 Bytes/cy

Instruction Buffer (48; 8x6 entries)

MOP   MOP   MOP   MOP   MOP   MOP

4-Way Decode

μOP   μOP   μOP   μOP

from/to
L1

Reservation Stations

| RSE0 20 entries | RSE1 20 entries | RSA0 10 entries | RSA1 10 entries | RSBR 19 entries |

FLA   PR   FLB   EXA   EXB   EAGA   EAGB   BR

μOP   μOP   μOP   μOP   μOP   μOP   μOP   μOP

Int ALU   Predicate manipul.   Int ALU   Int ALU   Int ALU   LD   LD   BR
FP arith.   FP arith.   MUL   DIV   AGU   AGU
FMA   FMA   Int ST
FP DIV
VEC addr. calc
FP ST

Execution Units

Register file
GPR: 32 (96 physical)
FPR: 32 (128 physical)
PPR: 16 (48 physical)

Fetch Port:
160 entries (40 physical)

Store Port:
192 entries (24 physical)

64 B/cy   64 B/cy   64 B/cy

L1D Cache

[1]horizontal recursive add

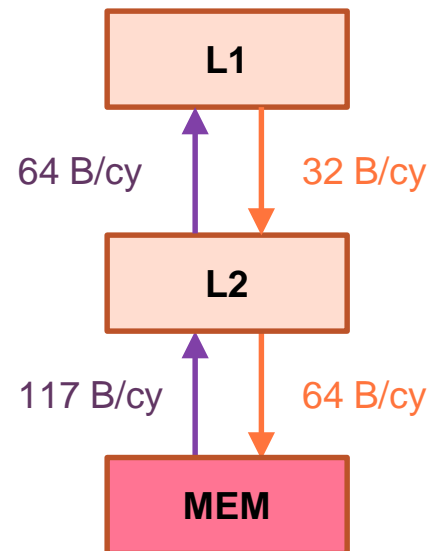## STREAM TRIAD

a[i] = b[i] + s * c[i]

```
.L18:
  ld1d z4.d, p5/z, [x21, x9, lsl 3]
  ld1d z5.d, p5/z, [x20, x9, lsl 3]
  fmad z5.d, p5/m, z2.d, z4.d
  st1d z5.d, p5, [x19, x9, lsl 3]
  add x8, x9, 8
  whilelo p5.d, w8, w7
  b.any .L18
```
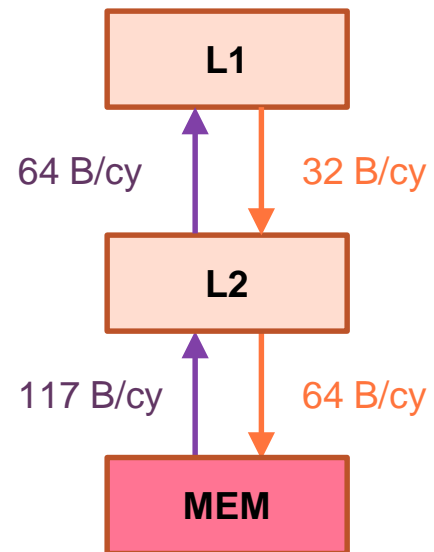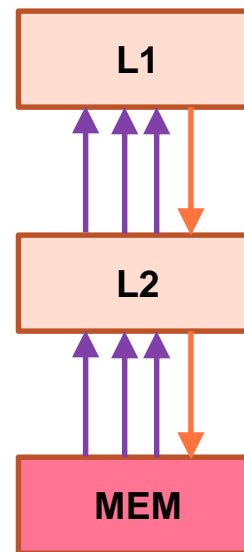
**2cy / VL**

Machine model
FX700

Machine model
FX700

Application model
TRIAD
a[i] = b[i] + s*c[i]

L1

64 B/cy          32 B/cy

L2

117 B/cy          64 B/cy

MEM

L1

L2

MEM

**Machine model**
FX700

L1

64 B/cy          32 B/cy

L2

117 B/cy          64 B/cy

MEM

**Application model**
TRIAD
`a[i] = b[i] + s*c[i]`

L1

L2

MEM

**ECM prediction**
TRIAD on FX700

RD
RD
RD

L2 → L1
3 cy/VL

L1

WR

L1 → L2
2 cy/VL

L2

RD
RD
RD

MEM → L2
1.64 cy/VL

WR

L2 → MEM
1 cy/VL

MEM

## Machine model
### FX700

Registers

128 B/cy     64 B/cy

L1

64 B/cy      32 B/cy

L2

117 B/cy     64 B/cy

MEM

## Application model
### TRIAD
`a[i] = b[i] + s*c[i]`

Registers

L1

L2

MEM

## ECM prediction
### TRIAD on FX700

LD
LD

Registers

ST

RD

L1

WR

RD

RD

RD

L2

RD

RD

MEM

WR

ECM prediction
TRIAD on FX700

How do these boxes overlap?

| LD |
| LD |

| ST |

| RD |

| RD |

| WR |

| RD |

| RD |
| RD |
| RD |

| WR |

Hypothesis 1 : No overlap

Hypothesis 2 : Full overlap

Hypothesis 3 : Full overlap + half-duplex

Hypothesis 4 : L1L2 overlap + half-duplex at MEM

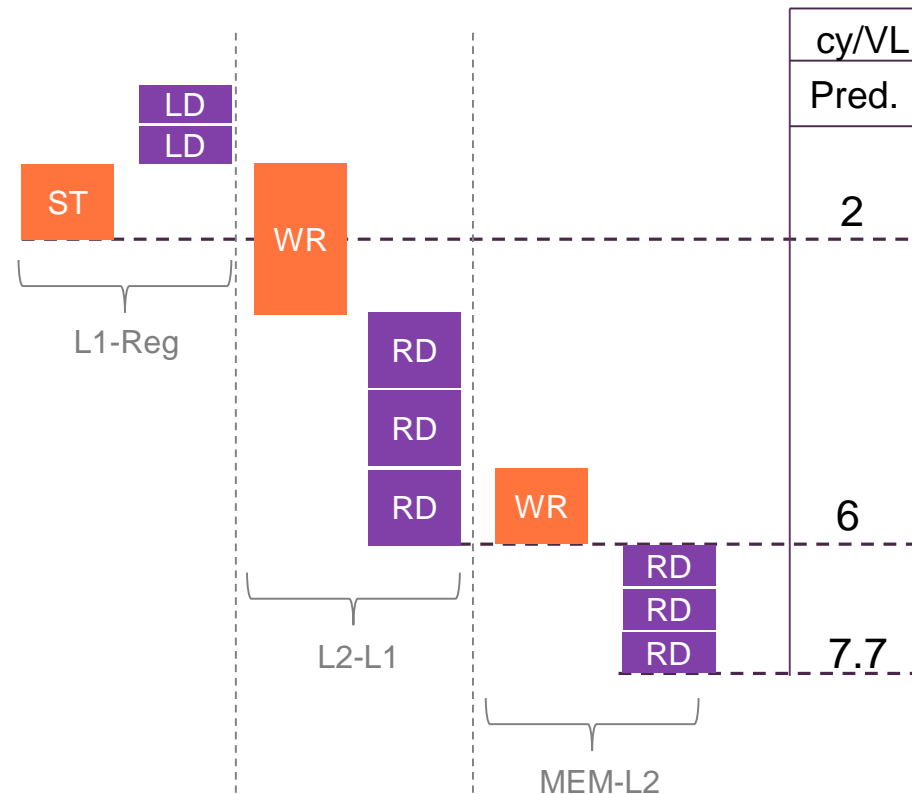There are numerous combinations.

How do we find the correct one?

Compare measurements
with predictions.

The best hypothesis for FX700
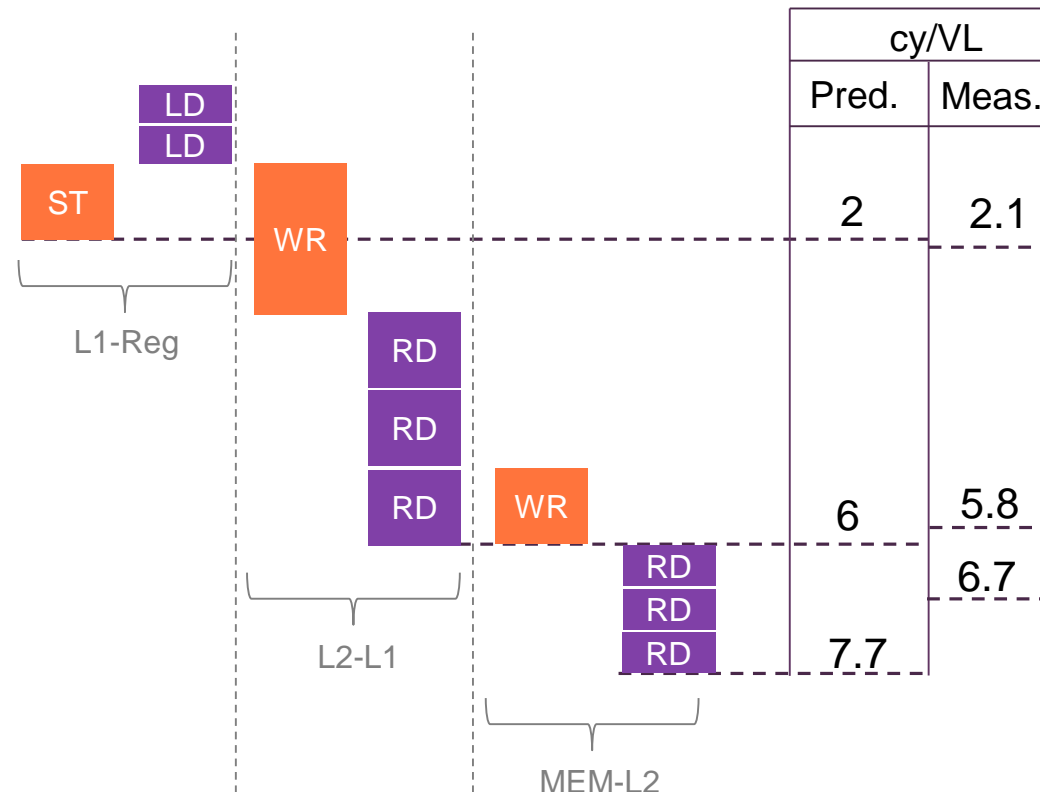
There are numerous combinations.

How do we find the correct one?

Compare measurements
with predictions.

The best hypothesis for FX700

There are numerous combinations.

How do we find the correct one?
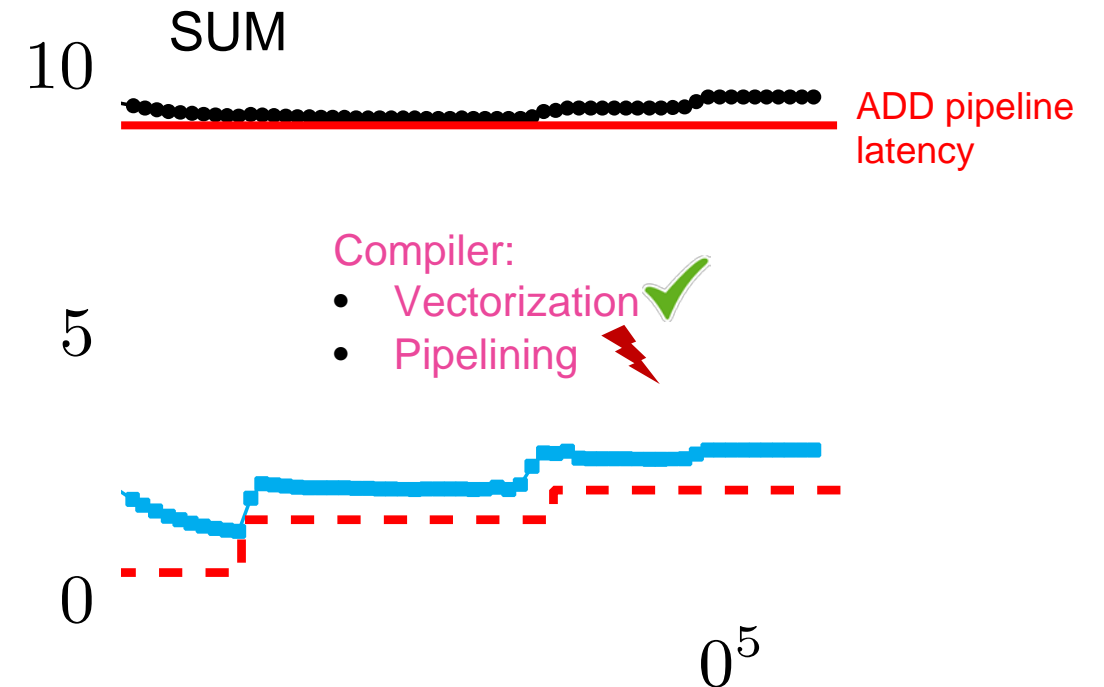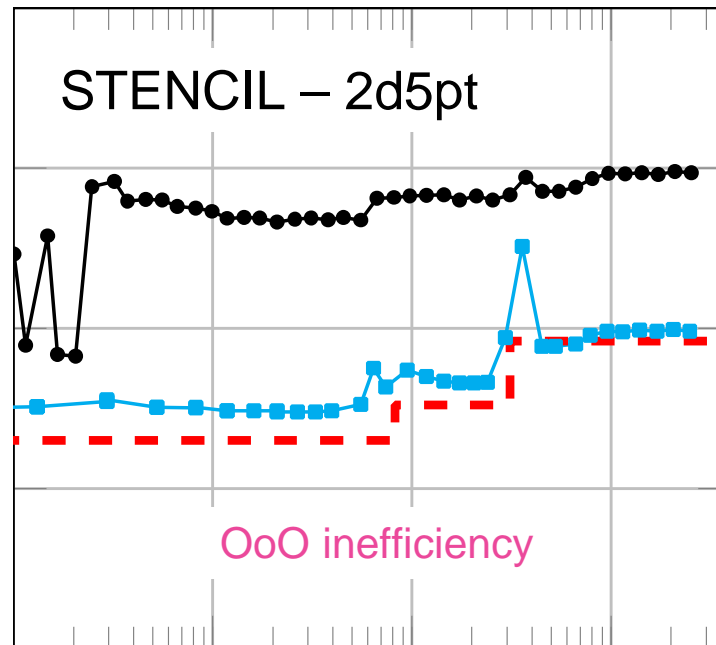
Compare measurements
with predictions.

| | cy/VL | |
|---|---|---|
| | Pred. | Meas. |
| | 2 | 2.1 |
| | 6 | 5.8 |
| | | 6.7 |
| | 7.7 | |

LD
LD
ST
WR
L1-Reg
RD
RD
RD
WR
L2-L1
RD
RD
RD
MEM-L2

A systematic way of identifying overlap hypothesis is presented in : Hofmann et.al., 2020, Bridging The Architecture Gap: Abstracting Performance-relevant Properties Of Modern Server Processors, https://doi.org/10.14529/jsfi200204
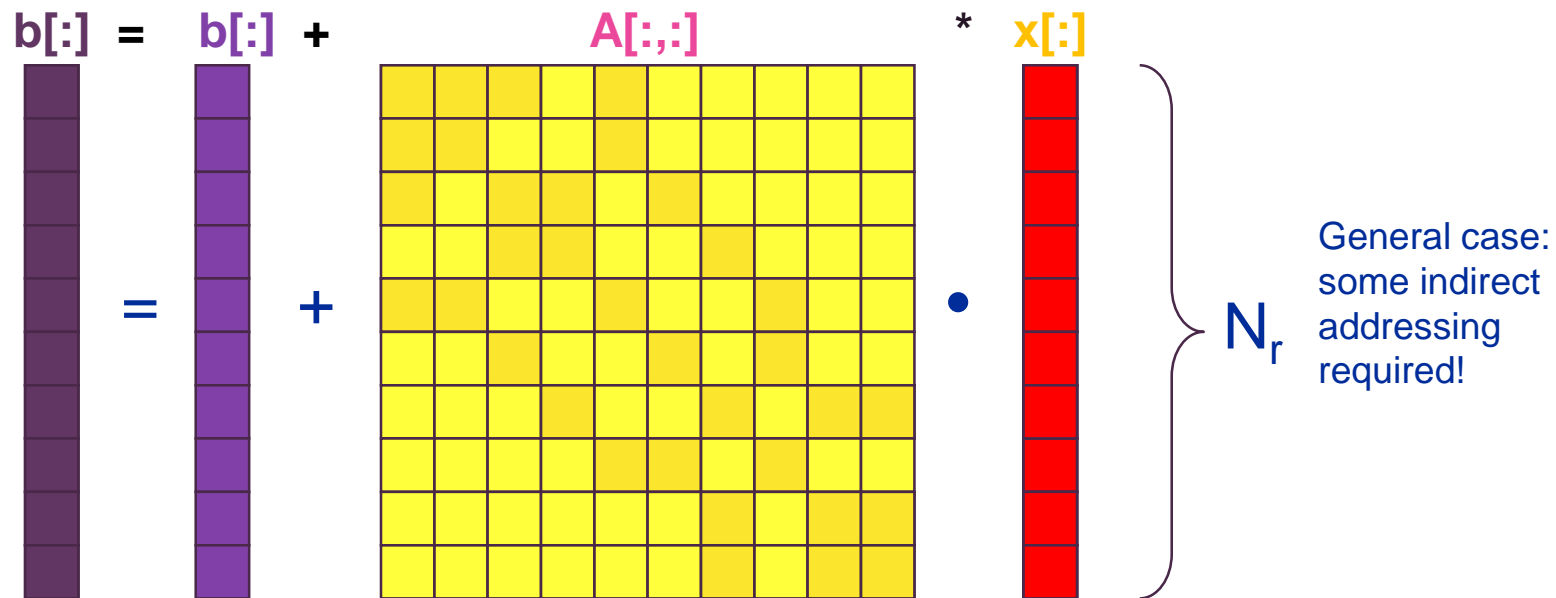
Unrolling plays an important role



Unrolling factor=1    Unrolling factor=8    ECM prediction

Sparse Matrix-Vector Multiplication (SpMV) : `b=Ax`

**b[:] = b[:] +**      **A[:,:]**      *   **x[:]**

General case: some indirect addressing required!

$N_r$

In Compressed Row Storage (CRS) format

```
for i = 0:nrows-1 //Long outer-loop
    for j = row_ptr[i]:row_ptr[i+1]-1 // Short inner-loop
        b[i] = b[i] + A[j] * x[col_idx[j]]
```

Assembly of the short inner-loop

```
.L6:
        ld1sw       z0.d, p0/z, [x17, x20, lsl 2]
        ld1d        z2.d, p0/z, [x18, x20, lsl 3]
        ld1d        z3.d, p0/z, [x30, z0.d, lsl 3]
        add         x20, x20, 8
        fmla        z1.d, p0/m, z3.d, z2.d
        whilelo     p0.d, x20, x14
        b.any       .L6

        faddv       d4, p1, z1.d
```

GCC compiler

In Compressed Row Storage (CRS) format

```
for i = 0:nrows-1 //Long outer-loop
    for j = row_ptr[i]:row_ptr[i+1]-1 // Short inner-loop
        b[i] = b[i] + A[j] * x[col_idx[j]]
```
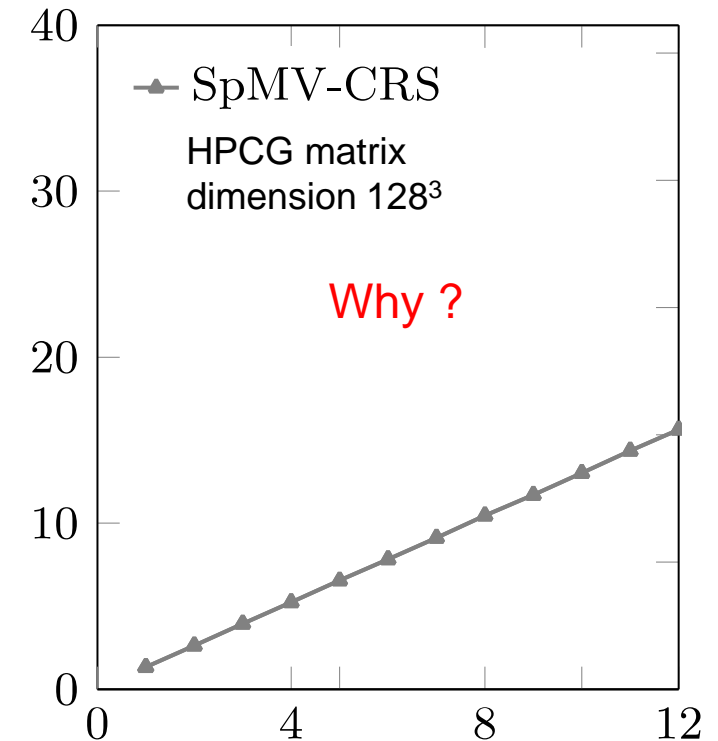
Assembly of the short inner-loop

```
.L6:
        ld1sw       z0.d, p0/z, [x17, x20, lsl 2]
        ld1d        z2.d, p0/z, [x18, x20, lsl 3]
        ld1d        z3.d, p0/z, [x30, z0.d, lsl 3]
        add         x20, x20, 8
        fmla        z1.d, p0/m, z3.d, z2.d
        whilelo     p0.d, x20, x14
        b.any       .L6

        faddv       d4, p1, z1.d
```

In Compressed Row Storage (CRS) format

```
for i = 0:nrows-1 //Long outer-loop
    for j = row_ptr[i]:row_ptr[i+1]-1 // Short inner-loop
        b[i] = b[i] + A[j] * x[col_idx[j]]
```



SpMV-CRS

HPCG matrix dimension $128^3$

Why ?

Assembly of the short inner-loop

```
.L6:
        ld1sw    z0.d, p0/z, [x17, x20, lsl 2]
        ld1d     z2.d, p0/z, [x18, x20, lsl 3]
        ld1d     z3.d, p0/z, [x30, z0.d, lsl 3]
        add      x20, x20, 8
        fmla     z1.d, p0/m, z3.d, z2.d
        whilelo  p0.d, x20, x14
        b.any    .L6

        faddv    d4, p1, z1.d
```

FMA3: Update **z1.d**

Latency: 9 cycles

Loop length : 27
**HPCG matrix**

Horizontal add of 512-bit register

Throughput = 11.5 cycles
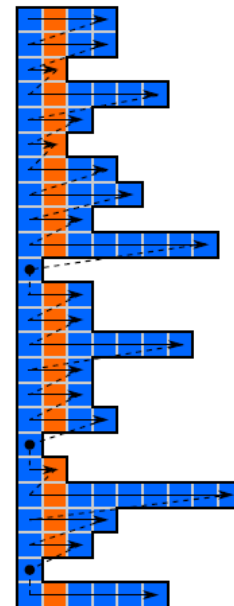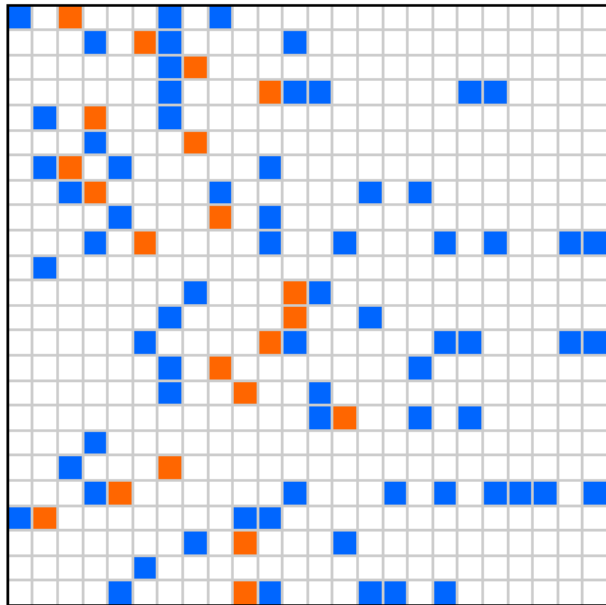
In Compressed Row Storage (CRS) format

```
for i = 0:nrows-1 //Long outer-loop
    for j = row_ptr[i]:row_ptr[i+1]-1 // Short inner-loop
        b[i] = b[i] + A[j] * x[col_idx[j]]
```
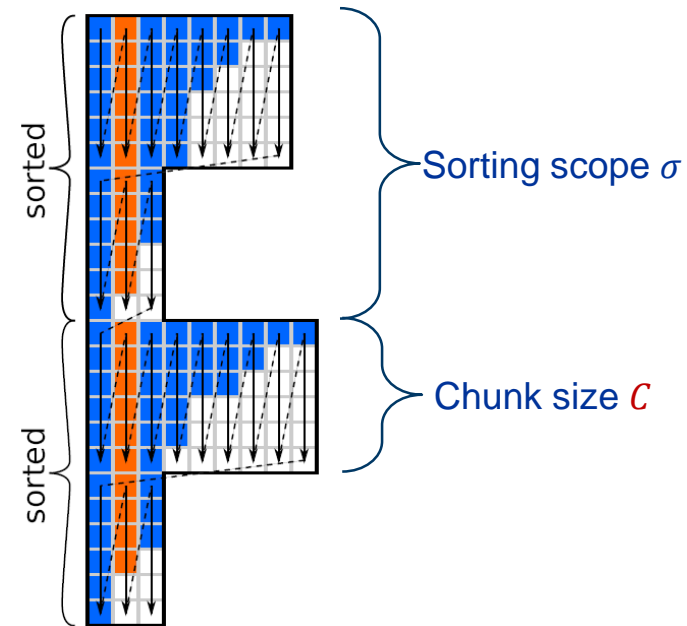
Vector-friendly SpMV data format/kernel required for A64FX → SELL-C-σ[1]



CRS                    SELL-C-σ

Sorting scope $\sigma$

Chunk size $C$

Benefits:

- Vectorization and unrolling along chunk size (C) → long loop and tunable
- No costly horizontal-add (faddv)
- No loss of vector efficiency

[1]M. Kreutzer et al., A Unified Sparse Matrix Data Format For Efficient General Sparse Matrix-vector Multiplication On Modern Processors With Wide Simd Units, SIAM SISC 2014,  DOI: 10.1137/130930352
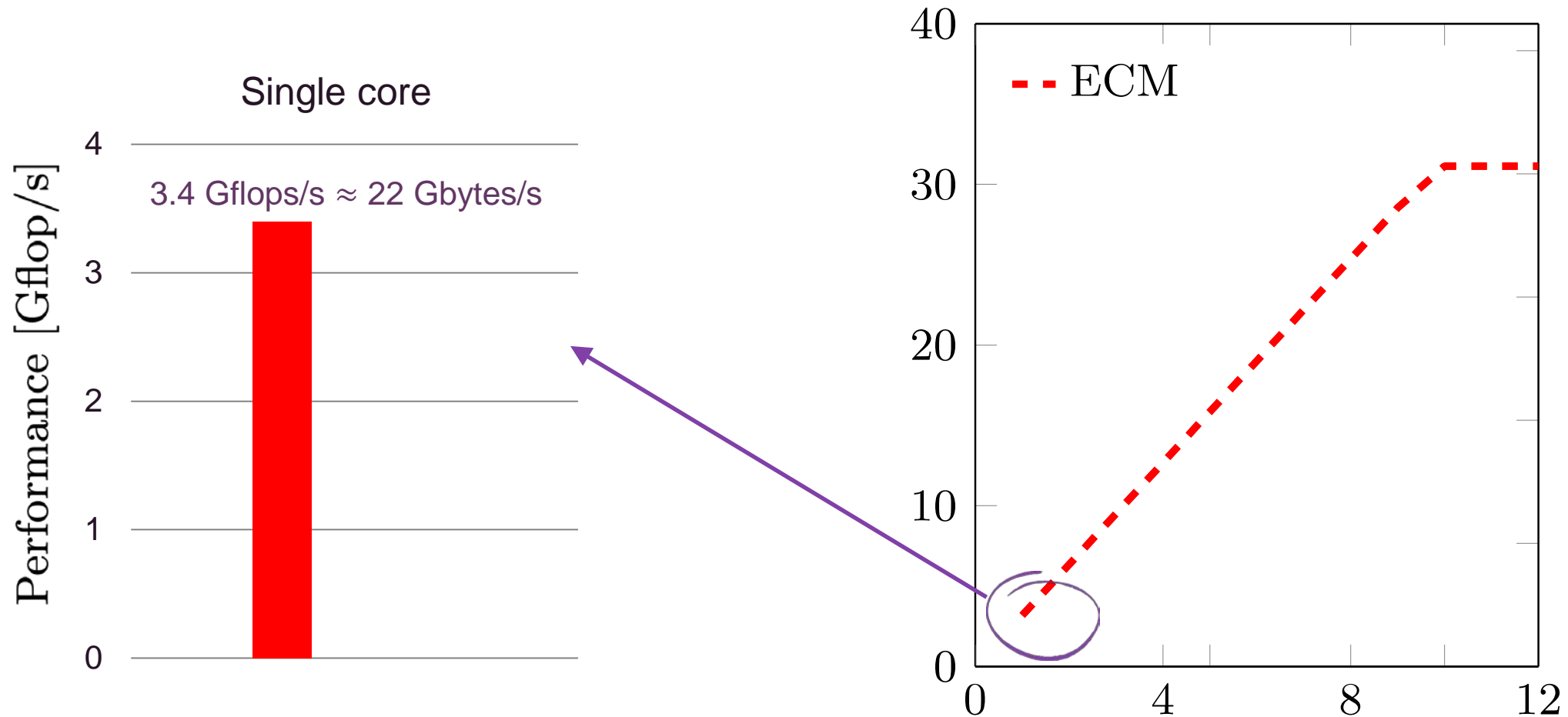
cy/VL

0

Gather

LD

ST

20.3

RD

25.8

RD

28.8

= 3.4 Gflops/s ≈ 22 Gbytes/s

Can we saturate now ?



HPCG matrix, dimension $128^3$

Can we saturate now ?     Yes, but needs almost all cores



HPCG matrix, dimension $128^3$

Can we saturate now ?



HPCG matrix, dimension $128^3$

SpMV performance on full node (48 cores)



Measured : 798 GB/s

■ SELL   ■ CRS

SpMV

SpMV performance on full node (48 cores)

Matrices from SuiteSparse Matrix Collection : https://suitesparse-collection-website.herokuapp.com

- High single core performance is crucial.

- ECM model was established and utilized to analyze the single core performance.

- The partial overlapping memory hierarchy allows for high single-core memory bandwidth.

- Proper single core optimizations have to be done to hide long floating point latency and inefficiencies in OoO.

- For SpMV we were able to saturate the bandwidth with SELL-C-$\sigma$ format.