# BENCHMARKING THE EVOLUTION OF PERFORMANCE AND ENERGY EFFICIENCY ACROSS RECENT GENERATIONS OF INTEL XEON PROCESSORS

Istvan Z Reguly, Balázs Dravai – PPCU ITK, Hungary

reguly.istvan@itk.ppke.hu

PMBS 24

# INTEL CPU GENERATIONS

- Sapphire Rapids – 2023 Q1
  - Intel 7, first chiplet architecture - 4 compute tiles
  - Up to 56 cores, 350W
  - 8xDDR5-4800 or HBM2E
- Emerald Rapids – 2023 Q4
  - Intel 7, 2 compute tiles
  - Up to 64 cores, 350W
  - 8xDDR5-5600

- Sierra Forest – 2024 Q2
  - Intel 3, 2 compute tiles
  - Up to 192 cores, 350W
  - 8x/12x DDR5-6400
- Granite Rapids – 2024 Q3
  - Intel 3, 3 compute tiles
  - Up to 128 cores, 500W
    - 72 core 6960P tested
  - 12xDDR5-8800 MRDIMM

# AMD PROCESSOR GENERATION

- AMD Milan – 2021 Q1
  - 7 nm, 8 chiplet
  - 64 cores, 280W
  - 8x DDR4-3200
- AMD Genoa – 2022 Q4
  - 5 nm, 12 chiplet
  - 96 cores, 360W
  - 12x DDR5-4800

- AMD Turin – 2024 Q3
  - 3 nm, up to 16 chiplet
  - 128/192 cores, 500W
  - 12x DDR5-6400

# PERFORMANCE, EFFICIENCY

- Number of bandwidth-bound codes (mostly explicit PDE solvers)
  - CloverLeaf (low-order), Seismic (high-order), OpenSBLI (large-scale DNS), MG-CFD (FV Euler on unstructured mesh), Volna (FV NLSW on unstructured mesh)
- Compute-intensive: miniBUDE (docking proxy)
- MPI or MPI+OpenMP – many through OPS/OP2 DSLs. icpx/g++ compilers.

- Compare:
  - Runtime
  - Architectural efficiency (effective BW estimates)
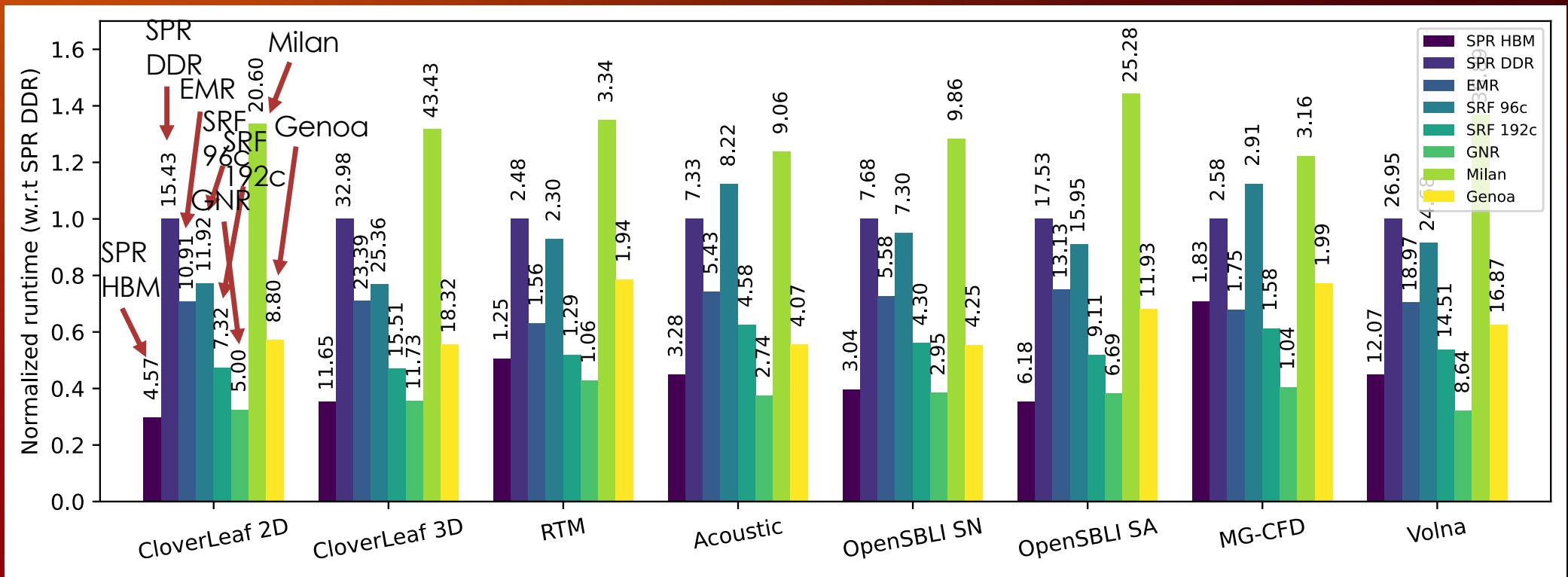  - Energy efficiency

# BASELINES

| | SPR HBM | SPR DDR | EMR | GNR | SRF 96c | SRF 192c | Milan | Genoa |
|---|---|---|---|---|---|---|---|---|
| Model | 9480 | 8480+ | 8592+ | 6960P | 6740E | | 7763 | 9B14 |
| Cores | 112 | 112 | 128 | 144 | 192 | 384 | 128 | 180 |
| LLC (MB) | 105 | 105 | 320 | 432 | 96 | 192 | 256 | 384 |
| Cache BW | 3481 | 4340 | 8149 | 7346 | 4139 | 7627 | 2454 | 9534 |
| DDR BW | 1475 | 388 | 542 | 1150 | 396 | 667 | 234 | 529 |
| Speedup | 3.8 | 1.0 | 1.39 | 2.96 | 1.04 | 1.23 | 0.6 | 1.36 |

Point of reference

# RUNTIME VS SAPPHIRE RAPIDS+DDR5

# ARCHITECTURAL EFFICIENCY

|  | SPR HBM | SPR DDR | EMR | SRF 96c | SRF 192c | GNR | Milan | Genoa |
|---|---|---|---|---|---|---|---|---|
| CloverLeaf 2D | 0.68 | 0.77 | 0.78 | 0.96 | 0.95 | 0.80 | 0.96 | 1.00 |
| CloverLeaf 3D | 0.52 | 0.70 | 0.71 | 0.87 | 0.87 | 0.66 | 0.88 | 0.92 |
| RTM | 0.21 | 0.41 | 0.47 | 0.43 | 0.47 | 0.33 | 0.51 | 0.39 |
| Acoustic | 0.34 | 0.57 | 0.56 | 0.49 | 0.54 | 0.52 | 0.77 | 0.76 |
| OpenSBLI SN | 0.35 | 0.53 | 0.53 | 0.54 | 0.56 | 0.47 | 0.69 | 0.71 |
| OpenSBLI SA | 0.50 | 0.67 | 0.65 | 0.71 | 0.76 | 0.60 | 0.78 | 0.73 |
| MG-CFD | 0.50 | 1.36 | 1.44 | 1.16 | 1.30 | 1.13 | 1.84 | 1.30 |
| Volna | 0.54 | 0.91 | 0.93 | 0.96 | 0.99 | 0.96 | 1.11 | 1.08 |
| Average | 0.45 | 0.74 | 0.76 | 0.77 | 0.80 | 0.68 | 0.94 | 0.86 |
| miniBUDE | 0.33 | 0.27 | 0.39 | 0.32 | 0.30 | 0.32 | 0.36 | 0.37 |

- Effective BW: array size in each loop
  - Above 1.0 if cache re-use across loops
- Significant MPI overheads on Seismic apps (30-50%)
  - Especially SPR+HBM and Genoa
- miniBUDE: fraction of peak GFLOPS/s at all-core turbo
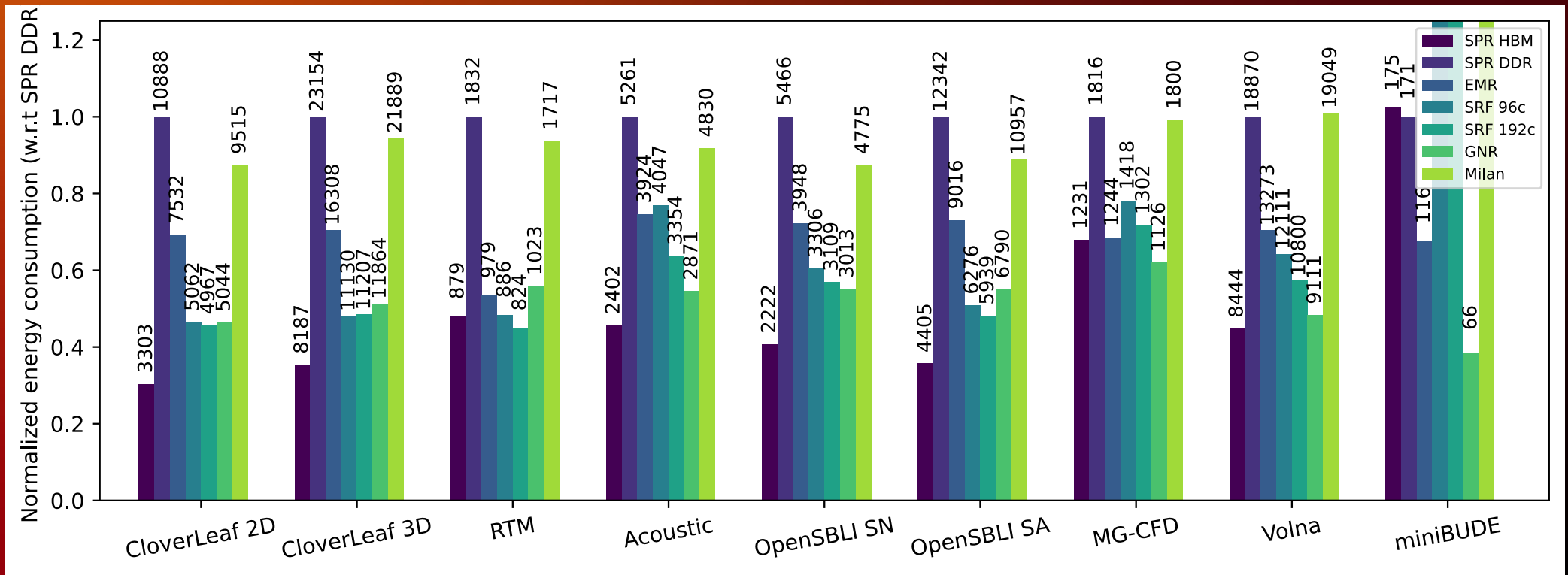  - Not always documented

# BALANCED IMPROVEMENT VS SPR+DDR?

- How much did overall performance improve vs. the improvement in bandwidth?
  - $(\text{runtime}/\text{runtime}_{SPR})/(\text{Peak BW}/\text{Peak BW}_{SPR})*100$
- SPR+HBM: same compute - only 64%
- EMR: +8 cores, 3x cache size, 1.4x BW – 98%
- GNR: +16 cores, 4.1x cache size, 3x BW – 91%


- SRF 96 core: +40 cores, 0.85x cache size, 1.04x BW – 105%
- SRF 192 core: +136 cores, 1.71x cache size, 1.7x BW – 110%

# ENERGY EFFICIENCY VS SAPPHIRE RAPIDS

# POWER

- 250W TDP for SRF 96c
- 280W TDP for Milan
- 350W TDP for SPR, EMR, SRF 192c
- 360W TDP for Genoa (no RAPL)
- 500W TDP for GNR
- Plus RAM

| | SPR HBM | SPR DDR | EMR | SRF 96c | SRF 192c | GNR | Milan |
|---|---|---|---|---|---|---|---|
| CloverLeaf 2D | 722W (3.37x) | 705W (1x) | 690W 208W (1.41x) | 424W 185W (1.29x) | 678W 318W (2.11x) | 1007W 309W (3.08x) | 461W (0.75x) |
| CloverLeaf 3D | 703W (2.83x) | 702W (1x) | 697W 203W (1.41x) | 438W 183W (1.3x) | 722W 318W (2.13x) | 1011W 306W (2.81x) | 503W (0.76x) |
| RTM | 701W (1.98x) | 739W (1x) | 627W 164W (1.59x) | 385W 154W (1.08x) | 640W 257W (1.93x) | 966W 212W (2.34x) | 513W (0.74x) |
| Acoustic | 732W (2.23x) | 718W (1x) | 722W 206W (1.35x) | 492W 157W (0.89x) | 732W 270W (1.6x) | 1046W 275W (2.67x) | 532W (0.81x) |
| OpenSBLI SN | 729W (2.52x) | 711W (1x) | 707W 208W (1.38x) | 452W 170W (1.05x) | 723W 305W (1.79x) | 1021W 265W (2.61x) | 484W (1.3x) |
| OpenSBLI SA | 712W (2.84x) | 704W (1x) | 686W 202W (1.34x) | 393W 170W (1.1x) | 652W 295W (1.92x) | 1014W 282W (2.62x) | 433W (0.78x) |
| MG-CFD | 671W (1.41x) | 702W (1x) | 710W 168W (1.48x) | 487W 143W (0.89x) | 823W 249W (1.64x) | 1078W 232W (2.48x) | 569W (0.69x) |
| Volna | 699W (2.23x) | 700W (1x) | 699W 208W (1.42x) | 490W 169W (1.09x) | 744W 300W (1.86x) | 1054W 322W (3.12x) | 516W (0.68x) |

# CONCLUSIONS

- Rapid evolution with interesting trade-offs
- Intel's differentiated product lines – evolution vs. SPR
  - SRF 96 core: Slightly higher performance (1.09×) but at a 1.4× lower power
  - SRF 192 core: Significantly higher performance(1.87×), at the same power
  - GNR: Even higher performance (2.72×), but at 1.44× more power
- MRDIMM may be a critical differentiating factor vs. AMD
  - Big advantage for memory-bound codes
- Still behind with energy efficiency – process issues