Co-design with Proxy-Apps: A match made in heaven?

EU-SW-CoDesign WS @ISC'22 – 02. June 2022

Jens Domke, Dr. rer. nat. <jens.domke@riken.jp > High Performance Big Data Research Team, RIKEN R-CCS, Kobe, Japan





Motivation – Heterogeneous Future

• The "Cambrian explosion" was just the start! (50+ in this fig.)



Source: Albert Reuther, et al. "Survey of Machine Learning Accelerators", 2020

Jens Domke Fig. 2. Peak performance vs. power scatter plot of publicly announced AI accelerators and processors.

Motivation – Heterogeneous, Diverse Future



- Increase on almost all fronts:
 - Specialized chips in phones, FPGAs in Intel/AMD CPUs, ...
 - Number of programming languages (Rust, Julia, ...) and paradigms (CUDA, HIP, oneAPI, Kokkos, RAJA, ...)
 - New workloads: containerization, DL/ML, big data, workloadchaining, etc.
 - New topologies: HyperX, Slimfly, Dragonfly, Megafly, ...

→ How to make sense of it all? (and fast)

→ How to determine the best architecture per site?



Proxy-Apps: What are they good for? Traditional approach...

"Use benchmark X and run workload Y and report back." --HPC procurement

Opportunity for new topologies – HyperX





First large-scale 2.7 Pflop/s (DP) HyperX installation in the world!

J. Domke et al. "HyperX Topology: First at-scale Implementation and Comparison to the Fat-Tree" Jens Domke



Fig.1: **HyperX** with n-dim. integer lattice $(d_{1},...,d_{n})$ base structure fully connected in each dim.



TokyTech's 2D HyperX:

- 24 racks (of 42 T2 racks)
- 96 QDR switches (+ 1st rail) without adaptive routing
- 1536 IB cables (720 AOC)
- 672 compute nodes
- 57% bisection bandwidth



Theoretical Advantages (over Fat-Tree)

- Reduced HW cost (less AOC / SW)
- Only needs 50% bisection BW
- Lower latency (less hops)
- Fits rack-based packaging

Opportunity for new topologies – HyperX

1:1 comparison (as fair as possible) of 672-node 3-level Fat-Tree and 12x8 2D HyperX

- NICs of 1st and 2nd rail even on same CPU socket
- Given our HW limitations (few "bad" links disabled)

Wide variety of benchmarks and configurations

- 3x Pure MPI benchmarks
- 9x HPC proxy-apps
- 3x Top500 benchmarks
- 4x routing algorithms (incl. PARX)
- 3x rank-2-node mappings
- 2x execution modes

Primary research questions

- Q1: Will reduced bisection BW (57% for HX vs. ≥100% for FT) impede performance?
- Q2: Two mitigation strategies against lack of AR? (→ e.g. placement vs. "smart" routing)









Fig.4: Baidu's (DeepBench) Allreduce (4-byte float) scaled 7-> 672 cn (vs. "Fat-tree / ftree / linear" baseline)

1. Placement mitigation can alleviate bottleneck

- 2. HyperX w/ PARX routing outperforms FT in HPL
- 3. Linear good for small node counts/msg. size
- 4. Random good for DL-relevant msg. size (+/- 1%)
- 5. "Smart" routing suffered SW stack issues
- 6. FT + ftree had bad 448-node corner case

Conclusion HyperX topology is promising and cheaper alternative to Fat-Trees (even w/o adaptive *R*) !



Proxy-Apps: Motivating vendors / domain scientists...

"Wanna do HPC? Then you need many and fast FP64 Units" --most HPC beginner classes

More Flop/s → more science?!

- Computer simulations R-CCS create the future
- Thanks to the (curse of) the TOP500 list, the HPC community (and vendors) are chasing higher FP64 performance, thru frequency, SIMD, more FP units, …



Jens Domke J. Domke et al. "Double-precision FPUs in High-Performance Computing: an Embarrassment of Riches?"

Compare Time-to-Solution in Solver





Fig. 4. Speedup of KNM over KNL as baseline. MiniAMR included since the input is the same for both Phi; Proxy-app abbreviations acc. to Section II-B

- Only 3 apps seem to suffer from missing FP64 unit (MiniTri: no FP; FFVC: only int+FP32)
- Options for memory-bound applications (almost all):
 - Invest in memory-/data-centric architectures
 - Move to FP32/mixed precision
 → less memory pressur
- Options for compute-bound applications:
- **Brace for less FP64 units** (driven by market forces) Jens Domke and less "free" performance (10nm, 7nm, 3nm, ...then?)





Proxy-Apps: Influencing architecture...

"Wanna do HPC? Then you need fast [S/D]GEMM." --every HPC beginner class

BLAS / GEMM utilization in HPC Applications



BERT

VGG16

Resnet50

SSD300

GEMM

NCF

GRU

LSTM

DeepLabV3

Cosmoflow

Speedup

3.39x

1.16x

1.71x

1.97x

1.75x

1.78x

0.97x

7.59x

3.67x

5.69x

Analyzed various data sources:

J. Domke et al. "Matrix Engines for High Performance Computing: A Paragon of Performance or Grasping at Straws?"

- Historical data from K computer: only 53,4% of node-hours (in FY18) were consumed by applications which had GEMM functions in the symbol table
 Benchmark
- Library dependencies: only 9% of Spack packages have direct BLAS lib dependency (51.5% have indirect dependency)
- TensorCore benefit for DL: up to 7.6x speedup for MLperf kernels
- GEMM utilization in HPC: sampled across 77 HPC benchmarks (ECP proxy, RIKEN fiber, TOP500, SPEC CPU/OMP/MPI) and measured/profiled via Score-P and Vtune



12

Estimated Benefit by MEs for HPC Centers

- Thought experiment: Assume we have/had GEMM units in past or future systems.
 - Known: node-hour by domain
 - Sample application with highest BLAS utilization
 - Estimate the node-hour reduction assuming different speedup by ME (2x–8x is realistic dep. on precision)
 - Future system includes 20% DL workloads, other science domains ~10% each
- Results w/ ideal conditions + 4x ME speedup: 5.3% less on K, 10.8% @ANL, 23.8% future system

→ HPC can utilize MEs when they come for free, but it's no magic bullet as for DL workloads

Jens Domke Explore more/other alternatives for Fugaku-next!







Proxy-Apps: Functionality and Regression Testing...

"Porting an application to A64FX? Just use fcc and –Kfast." --Fujitsu

"Silver bullet" compiler choice for A64FX?

Jens

Compiler Varia

- Issue: unexpected advantage of Xeon vs. A64FX in PolyBench
- **Performance portability** (x86→A64FX) **not easy** to achieve
- Testing >100 Kernels and HPC runtime Micro Kernels (all with [1|12]) Kernel 1 [C] 0.001 -0.594 -0.961 -0.961 -0.964 Kernel 2 [C] 0.002 -0.580 -0.985 -0.985 -0.879 Workloads on Fugaku Kernel 3 [F] 0.010 -0.196 -0.194 -0.196-0.958 0.002 -0.015 -0.019 -0.016 -0.797 Time-to-Solution [in s] (FJtrad) and Relative Performance Gain (others) 0.406 -0.668 0.005 -0.506 runtime PolyBench (all with [1]1] 0.003 -0.634 10.743 correlation [C] 0.98 0.076 Three compilers and ovariance [C] 10.735 1.150 0.08/ 0.990 gemm [C] 1.629 0.389 0.606 0.011 0.798 gemver [C] 0.092 0.022 -0.833 esummv [C] 0.025 0.008 0.004 -0.020 0.005 1.350 symm [C] 12 491 -0.062 1.931 -0.053 five variations svr2k [C] 6 4 2 4 0.315 0.378 0.035 -0.6880,490 3.011 0.169 0.231 0.377 syrk [C] 0.008 run error 0.906 trmm [C] 6 267 1.065 2711.882 0.977 0.004 18 273 15.296 1.264 (2x Fujitsu, 0.005 Time-to-Solution [in s] (FJtrad) and Relative Performance Gain (others) 1.197 2.847 1.498 0.638 -0.581 2.046 0.020 0.041 21.506 [48] 0.001 [48] 0.043 [48]1 0.046 [48] -0.023 [48]1 1303.331 1.587 -0.518 [48|1] 0.528 [48]1] 0.031 [48]1] -0.116 [48]1 -0.191 [48]1] 0.129 2x LLVM12. HPCG [C++] run error 1.676 [1|36] -0.004 [1|36] 0.010 -0.852 -0.192 0.048 [1]48] -0.071 [1|48] 0.016 [1]48 0.019 [1]48 0.155 [1|48] 579 -0.577 -0.781 0.008 -0.416 0.383 Time-to-Solution [in s] (FJtrad) and Relative Performance Gain (others) 0.206 [32]1 -0.524 [4|12] 0.637 -0.571 & GNU 0.124 [48]1] 0.122 [48]1 0.132 [48|1] runtime SPEC CPU int (all with [1]1 0.269 -0.534 -0.375 erlbench (C) -0.497 95.842 M+Polly 0.184 [48]1] 0.172 [48]1 0.043 [48]1] GNU 0.237 gcc [C 144.654 -0.656 -0.662 -0.660 .014 -0.145 -0.043 0.016 [48]1] 0.024 [48]1] 0.113 [48]1] -0.554 met IC 107.932 -0.546 -0.555 0.266 -0.741 1.209 0.374 [4|12 0.459 [4]12 -0.718 [4|12] omnetop (C++ 0.127 0.038 0.021 0.025 Time-to-Solution [in s] (FJtrad) and Relative Perfo 078 0.935 0.057 nance Gain (others) -0.844 -0.844 0.145 0.027 [48]1 0.026 [48] -0.304 [48]1 0.225 -1.000 0.338 0.022 [48]1 0.003 [48]1 -0.484 [48]1] -0.300 0.255 0.286 0.388 0.055 [32]1] 0.045 [32]1 0.063 [32]1] -0.421 0.362 144 654 -0.010 -0.009 -0.029 [1|48] -0.541 1.132 107.932 -0.546 -0.554 0.266 0.647 0.036 **RIKEN** miniapps 0.543 153.212 -0.486 SPEC CPU float 227.200 0.476 66.105 -0.840 -0.844 0.225 -0.009 [1]32 -0.011 [1]3 0.774 0.260 [48]1 0.254 [48]1] -0.926 [48]1] 174,909 0.338 LVM -0.031 [1]48] 0.011 [1]48] -0.426 [1|36] LLVM+Pollv GNU 225.918 0.293 0.255 0.388 -0.134 [16]3] 0.139 [16]3 -0.768 [16]3] -0.446 [1]48 137.069 0.000 -0.010 -0.009 -0.118 [1|48] -0.123 [1|48] 0.081 [48]1] 0.173 [48] -0.334 [48]1] er Variant 77.578 0.007 [1]32] 0.002 [1]32] -0.864 [1]32 0.001 [10]4 -0.768 [10]4] 0.014 [10]4 0.012 [1|48] 0.007 [1|48] -0.322 [1]32 -0.011 [1]3 0.004 0.009 [1 0.061 [12]4] 0.061 [12]4] -0.418 [24]1] 0.002 (1)32 0.001 [1]32 -0.031 [1]48 0.011 [1]48] 4.956 [1]48] 0.249 [1] 0.003 [24]2 0.036 [24]2] 11.397 [1]48 0.109 -0.118 [1]4 -0.123 [1]48 -0.446 [1|48] wrf IE CI 5.482 [1]32] 0.002 [1]32 0.007 [1]32 0.002 [1|32] -0.864 [1]32 0.082 [1|48 0.155 [1|48 -0.013 [1|48 LLVM LLVM+Polly GNU 0.007 [1|48] cams IE (11.925 [1]48 0.006 [1] 0.012 [1]4 -0.322 [1]32 0.001 [1|32] -0.004 [1]32] -0.521 [1|32] 15.115 [1]32 mpiler Variant 0.004 (1)48 -0.003 [1]48 -0.779 [1]4 J. Domke "A64FX – Your Compiler You 18.778 [1|8] 17.824 [1|48] 0.207 [1 0.082 [1]4 LLVM LLVM+Polly GNU 8.893 [1]48 0.001 [1]3 0.004 0.521 8.134 [1]4 -0.013 [1 0.004 [1] 0.003 [1]4 779 [1] Must Decide!' Compiler Variant FJcland LLVM

Time-to-Solution [in s] (FJtrad) and Relative Performance Gain (others)

+1.0

14



Proxy-Apps: Investigating new CPU architectures...

"Can a simulator handle those complex codes? Surely it must be possible, right, RIGHT?" --naive me



What-if we can fit all data into L1D?

Machine-Code-Analyzers (MCA)

 Estimate the throughput of proxy-app/workload's basic blocks



- MCA tools ignore memory subsystem (assume optimistic L1 hits for all load/store ops)
- Build a BB graph and estimate the runtime → equiv. to "infinitely" large L1 → "speed of light" simulation which trades accuracy for sim. speed (only ~1000x slowdown)

→ High speedup potential for almost all proxy and micro-bench. (again: HPC is memorybound)



LARge Cache processor w/ 3D-stacked SRAM



What-If extrapolation:

- Assuming we can 3D-stack SRAM with up to 8 layers on top or below cores
 projected L2 cache size & L2 bandwidth?
 projected performance gain for real apps?
- Explore ~7 years into the future for 1.5nm fab technology
- Reclaim L2 area to add more cores

Motivating State-of-the-Art

- MiniFE on AMD Milan vs Milan-X
- >3x speedup from larger LLC (L3) when problem fits into Milan-X' 768MB cache



Computer simulations

create the future

R-CCS

Jens Domke

Minor alterations due to gem5 limitations

LARC' CMG not easily replicated

- Configs.: LARC_c < LARC < LARC^a
- Cycle-level accurate simulation w/ gem5 (patches to fix bugs and scale cores)
- Base: 1 CMG of Fugaku's A64FX
- Additional: A64FX³² to separate effect of core increase from cache increase
- Only single-rank proxy-apps 🙁



CHIP AREA AND SIMULATOR CONFIGURATIONS FOR GEM5

Configuration	A64FX _S	A64FX ³²	LARCC	LARCA
Area (per CMG)	$48\mathrm{mm}^2$	N/A	$12\mathrm{mm}^2$	12mm^2
Cores (per CMG)	12	32	32	32
Core config.	Arm v8.2 + SVE, 512bit SIMD width, 2.2 GHz OoO 128 ROB entries, dispatch width 4			
Branch pred.	Bi-mode: 16K global predictor, 16K choice predictor			
Per-core L1D	64KiB 4 way set assoc, 3 cycles, adjacent line pref.			
L2 cache size/BW (per CMG)	8 m N ~ 80	∕IiB 0 GB/s	$\begin{array}{c} 256\mathrm{MiB} \\ \sim 800\mathrm{GB/s} \end{array}$	512 MiB ~ 1600 GB/
L2 config.	16-way set-assoc, 37 cycles, inclusive, 256B block			
Aggregated L2 size/BW (per CPU)	32 ~ 3.2	MiB 2 TB/s	4096 MiB ~ 12.8 TB/s	8192 MiB ~ 25.6 TB/
Main Memory	32 GiB HBM2, 4 channels, 256 GB/s			

Validation tests with STREAM Triad:

← Expected behavior in L1, L2, and HBM

J. Domke et al. "At the Locus of Performance: A Case Study in Enhancing CPUs with Copious 3D-Stacked Cache"



Many patches, crashes, and moons later...

6-months of gem5 runs on a cluster yield:

- 29 of 52 apps show ≥2x speedup on LARC^A compared to A64FX^S CMG (on 1/4th the area)
- For 23 of 29 the gain is attributed to larger L2
- Overall: we estimate (GM=)9.77x speedup for those proxy-apps for a 16-CMG LARC CPU with 8GiB stacked L2 and 512 cores

Proxy-App

Kernel 12

Kernel 17

Kernel 19

MG-OMP

XSBench

FT-OMP

App-specific tuning for large cache → even more possible (s. Ltaief et al. in SC'21)



A64FXs

36.6

46.7

73.8

11.6

59.8

32.1

A64FX³²

47.6

49.5

69.6

48.2

70.9

36.4



LARCA

9.1

34.8

48.9

3.8

0.4

0.1

LARCC

10.5

48.7

49.1

6.4

29.4

0.1



Proxy-Apps: Few years of user experience ...

"Uff... (to say it politely)." --me

Fugaku Enhancement & Co-Design for Future

- Superseding current proxy-apps: Octopodes
 - Downsides w/ Fiber/proxy-apps (s. Fugaku R&D)
 - On-going collaboration / brainstorming phase with DOE labs (position paper release in Apr.'22)
 - Set of highly-parameterizable, easily-amendable, **MOTIF-like problem representations**
 - Common "language" between HPC users, system operators, co-designers, and vendors to describe the to-be-solved scientific problems: What needs to be computed, and how it can be computed?



S. Matsuoka, J. Domke, M. Wahib, A. Drozd, A. Chien, R. Bair, J. S. Vetter, J. Shalf "Preparing for the Future – Rethinking Proxy Applications" to appear in Computing in Science & Engineering





Computer simulations

Summary and Call for Co-Design Collaboration

- Better & versatile simulators
 - Consolidating and enhancing existing infrastructure
- More testbeds of diff. scales
 - Some netw. issues manifest at scale
 - Repurpose decommissioned system
 - Shared access
- Better performance metrics reflecting real-life and tools to collect them
 - Focus on more async., automated, easy-to-use for admins
- Cleaner build environments and user training

- Extensive **network co-design**
 - More insight into apps/workloads
- More bottleneck analysis
 - Fix the lack for memory-centric tools
- 1! (better) community test suite
 - Easy migration, inter-op. testing, regression testing, etc.
- Octopodes for co-design ©

Lots of exciting R&D challenges ahead!



Job & Collaboration Opportunities



- Collaborations and job opportunities:
 - We are hiring! Check out our research teams and open positions: https://www.riken.jp/en/research/labs/r-ccs/ and https://bit.ly/3faax8v jens.domke@riken.jp
- Internship/fellowship for students (Bachelor \rightarrow PhD):
 - Fellowship: https://www.riken.jp/en/careers/programs/index.html
 - Internship: <u>https://www.r-ccs.riken.jp/en/about/careers/internship/</u>
- Supercomputer Fugaku:
 - Apply for node-hours:
 <u>https://www.r-ccs.riken.jp/en/fugaku/user-guide/</u>
 - Interactive, virtual tour: <u>https://www.r-ccs.riken.jp/en/fugaku/3d-models/</u> and <u>https://www.youtube.com/watch?v=f3cx4PGDGmg</u>