Working with Proxy-Applications: *Interesting Findings, Lessons Learned, and Future Directions*

BID WS @PPoPP'22 – 02. Apr. 2022

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



Disclaimer



- No in-depth presentation into single scientific topic
- High-level examples to encourage out-of-the-box thinking and questioning "always done this way" and revisiting old papers

• Intension:

Strengthen collaboration between domain scientists and HPC architects to develop better Supercomputers

Computer simulations R-CCS create the future

Motivation – Heterogeneous Future

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



Source: A. 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 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.
 - (except for networks luckily and also unfortunately ⊗)

- → How to make sense of it all? (and fast)
- → How to determine the best architecture per site?

Challenges – Future Scale-out & Diversification

- No more gains from Moore's law
 - → Bigger HPC systems w/ more nodes (maybe island design for specialization)
 - → Larger interconn. networks

Supercomputer Fugaku

- >158k compute nodes
- 3 networks / topologies

Jens Domke

- CN: 24x23x24 TofuD (w/ 2x3x2 subgr.) as 6D torus
- Storage: EDR InfiniBand for every 16th CN (with fat-tree? topology)
 - Management + outside world: Ethernet



Source: Y. Tsujita et al. "Status of Lustre-Based Filesystem at the Supercomputer Fugaku"



Proxy-Apps: 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 $(d_{\eta},...,d_{n})$ base structure fully connected in each dim.



Fig.2: Indirect 2-level Fat-Tree

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

- TSUBAME2's older gen. of QDR IB hardware has no adaptive routing 8
- HyperX with **static/minimum routing** suffers from limited path diversity per dimension → results in high congestion and low (effective) bisection BW

Mitigation Strategies???

- **Option 1: Alternative Job Allocation**
- **Option 2: Non-minimal**, Pattern-aware Routing (PARX)



← combine →

Minimum paths

Forced detours



GiByte/s]

Thoughput [in

Computer simulations

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?!

- R-CCS Computer simulations
- 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



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





BERT

VGG16

Resnet50

SSD300

GEMM

NCF

GRU

DeepLabV3

Cosmoflow

Speedup

3.39x

1.16x

1.71x

1.97x

1.75x

1.78x

0.97x

7.59x

3.67x

15

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** ($x86 \rightarrow A64FX$) **not easy** to achieve



+1.0



Results for x500, Babel, ECP, Fiber

- Surprising **≈5% gain for HPL** (LLVM or FJclang) despite main time in SSL2
- Same for DLproxy (matmul convolution; SSL2) but even higher gain w/ GNU
- Time-to-Solution [in s] (FJtrad) and Relative Performance Gain (others) **GNU: 51%** runtime runtime Ranking HPL [C] 0.001 [48]1] 0.043 [48]1] 0.046 [48]1] -0.023 [48]1 21.506 [48]1] reduction in stream HPCG [C++] [48]1] 0.528 0.031 [48]1] -0.116 [48]1] -0.191 [48]1] -0.518 [48]1] Babel [C++] 1.676 [1|36]-0.004[1|36]0.377 [1]24] 0.296 [1]24] 0.512 $(\rightarrow eq. to higher GB/s)$ DLproxy [C] 0.048 [1|48] -0.071 [1]48] 0.016 [1]48] 0.019 [1]48] 0.155 [1]48] ECP proxy apps AMG [C] 5.042 [4|12] 0.058 [8]6] 0.206 [32]1] -0.304 [32]1] -0.524[4]12] • For ECP apps use CoMD [C] 4.460 [48]1] 0.124 [48]1] 0.122 [48]1] 0.132 [48]1] 0.086 [48]1] Laghos 45.436 [48]1 0.586 [48]1] 0.611 [4811] 0.399[48]1] run error MACSio [C,C++] [48]1] 0.184 [48]1] **LLVM** or GNU and 24.343 -0.067 [48]1] 0.172 [48]1] 0.043 [48]1] miniAMR [C] 18.857 [48]1] -0.045 [48]1] 0.016 [48]1] 0.024 [48]1] 0.113 [48]1] miniFE [C++] 0.373 [4]12] [4]12] -0.374 [4]12] -0.459 [4]12] -0.097 0.718 [4]12] for Fiber apps use miniTRI [C++] [32]1 3.607 1148 3 5 3 9 148 1148 29.182 0.217 [32]1] Nekbone [F] 1.656 [48]1] 0.027 [48]1] 0.026 [48]1] -0.304 [48]1] run error. Fujitsu's compiler SW4lite [F,C++] [48]1] 0.011 [48]1] -0.022 [48]1] 0.003 [48]1] 0.853 -0.4844811 0.035 [32]1] [32]1] 0.045 [32]1] 0.063 [32]1] SWFFT (F,C) 1.194 [32]1] 0.055 XSBench [C] 1.649 [1|48] 0.323 [1]48] 0.754 [1]48 5.865 -0.029[1|48]**RIKEN miniapps** • Avg. speedup: 1.65x FFB [F,C,C++] [4]1] 29.877 compile error complie error 2.457 le error FFVC [C++,F] [1]36] -0.026 [1|36] -0.254 [48]1] [48]1] 11.558 -0.260 [48]1] -0.926 (median 1.09x) with MODYLAS [F] 29.262 [16]3] 0.001 [16|3] -0.134 [16]3] -0.139 [16]3] 11613 -0.768 mVMC [C.F] [24]2] [24]2] -0.081 [48]1] [48]1] [48]1] 15.026 -0.0020.173 -0.334 max. 6.7x in XSBench NICAM [F] [10]4] [10]4] -0.001 [10]4] 0.014 [10]4] 7.645 -0.003-0.768[10]4] NTChem (F) [12]4] 0.061 [12]4] 0.061 [12|4] 9.440 [12]4] -0.0010.418 [24]1] OCD [F] 8.048 [24]2] 0.003 [24]2] 0.002 [24]2] 0.003 [24]2] 0.036 [24]2] -1.0 [4 | 12] rarely best option FJtrad FJclang LLVM LLVM+Polly GNU Compiler Variant Jens Domke 18

+1.0

Gain Performance Relative

"Silver bullet" compiler choice for A64FX? Conclusions:



- C1: recomm. usage model of 4 ranks and 12 threads often suboptimal
- C2: no "silver bullet" compiler for A64FX (yet)
 - Dep. on situation, but some hint: Fujitsu for Fortran codes, and GNU for integerintensive apps, and any clang-based compilers for C/C++
- C3: Twitter summary: "if Xeon is 70x faster than A64fx, suspect the compiler"
- → Test all available compilers, and explore other rank/thread mappings!

Implications for Fugaku Users:

Jens Domke

- LLVM 13 incl. "classic" flang (source /home/apps/oss/llvm-v13.0.0/init.sh)
 - SVE support still alpha → expect even more performance with v14
 - Using Fujitsu's SSL2 not always trivial or error-free



Proxy-Apps: Many years later ...

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

Drowning in overlapping Proxy-Apps (>>100)



- HPC challenge benchmark (HPCC)
- Exascale Computing Project (ECP) Proxy Applications
- Center for Efficient Exascale Discretizations (CEED) Miniapps
- DOE's CORAL-2 Benchmarks and RIKEN's Fiber / TAPP
- European Union's PRACE Unified European Application Benchmark Suite (UEABS)
 ExaMiniMp LAMMPS MiniQMC QMCPack sw4lite sw4 SWFFT HACC pent
- SPEC, BenchCouncil, Intel's HiBench, DeathStarBench, Baidu's DeepBench, and MLCommons' MLPerf

		ExaMiniMD	LAMMPS	MiniQMC	QMCPack	sw4lite	sw4	SWFFT	HACC	pennant	snap
cos e-similarity f0 cos cache-related perf. counter	ExaMiniMD	0.00	5.02	54.54	38.73	11.70	12.49	6.58	6.31	13.21	7.13
	LAMMPS	5.02	0.00	54,61	38.67	15.66	16.27	4.87	6.35	13.60	10.88
	fo ^r MiniQMC	\$4.54	54.65	0.00	17.15	47.4	46.08	10.02	48.9	42,16	49.15
	QMCPack	38.73	38.62	17.19	0.00	32,64	31.67	33.92	32.94	26.29	33.78
	sw4lite	11.70	15.66	47.41	32.64	0.00	1.15	13.41	11.40	11.15	5.07
	ters sw4	12.49	16.27	46.01	31.67	1.15	0.00	13.74	11.70	10.69	5.69
	SWFFT	6:58	4.87	50,01	33.92	13.41	13.74	0.00	2.24	9.05	8.80
	HACC	6.38	6,38	41.51	32.94	11.40	11.70	2.24	0.0	7.86	6.87
	pennant	13.21	13.60	42.36	26.29	11.15	10.69	9.09	7.8	D.00	9.37
	snap	7/45	10.88	49.13	33.78	5:07	5,69	8.80	6.87	9.37	0.00
		-									

Source: D. Richards et al. "Best Practices for Using Proxy Apps as Benchmarks"

• Traditional: HPL, HPCG, stream, Graph, and Intel's IMB, ...

Lessons-Learned from using Proxy-Apps



- Inherently complex and implementation biases
 - Mostly implemented in Fortran and C[/C++], and tuned over the years
 - Highly tuned CUDA (and "legacy" CPU) from ECP efforts
- Huge porting & maintenance overhead
 - Macros, separate code paths, hand-written makefiles, ...
 - Costly refactoring for: data layouts, parallelization strategies, accelerators
- Insufficient inputs, testing, documentation
 - Issues with non-std compilers and applying perf. analysis tools
 - Strong- and weak-scaling, varying input sizes, independent of #MPI/#OMP
- Lack of efficiency reporting

• How much performance achieved vs. the peak theoretical performance?



Let's clean up this mess ...

"Octopodes to the rescue." --RIKEN & DOE

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, **DWARF-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?



→ Apply ML to identify, parameterize, and categorize compute phases

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





Example of one Octopode: Matmul



- Input **shapes**: such as squared, rectangular, and tall/skinny
- Various numerical **precisions** (i.e., from fp128 to bfloat16, etc)
- Batched and non-batched executions modes
- **Dense** matrix-matrix operations, matrix-vector, sparse matrices
- Sparse matrix: random, realistic blocks, Matrix Market

- ➔ Use C++ templating to generate as many variants as possible to train ML models
- ➔ One Octopode for each distinct compute phase or math kernel
- → In-between Berkeley DWARFs and Proxy-Apps

Usage of Octopodes for Co-Design



- "What needs & how can it be computed" <u>not</u> "Here is how you have to do it"
- For performance modeling of real workloads: identify compute phases which can be mapped to one or more Octopodes → combine perf. model of the 'easier to understand' Octopodes → approx. perf. model of full workloads
- For vendors:
 - Allowed tuning freedom for the Octopodes, i.e., changes of algo., implementation, integer/float. precision, data layout, etc., *as long as* intended result is the same
 - Accurately model consumer workloads → Less over/under-selling of hardware
- **Porting** of user codes to new system:
 - Act as demonstrator for users to show how to port
 - ML/AI to identify phases can be used as helper for porting of real codes
- Better suited for co-design tools, e.g. compiler tests, regression testing, simulators (gem5/SST/CODES/...), quick "What-If" tools, etc.

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: <u>https://www.riken.jp/en/research/labs/r-ccs/</u> and <u>https://bit.ly/3faax8v</u>
- Internship/fellowship for students (Bachelor \rightarrow PhD):
 - Fellowship: <u>https://www.riken.jp/en/careers/programs/index.html</u>
 - 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>