

The experience of the HLST on Europes biggest KNL cluster

Tamás Fehér, Serhiy Mochalskyy, Nils Moschüring
Roman Hatzky
Intel MIC Programming Workshop
LRZ

Marconi – KNL at CINECA, Bologna

Total number of KNL nodes: 3600

Partition dedicated to the EUROfusion community: 392 (144 flat / 248 cache mode)

-> about 1 Pflop/s



Photo: F.Pierantoni

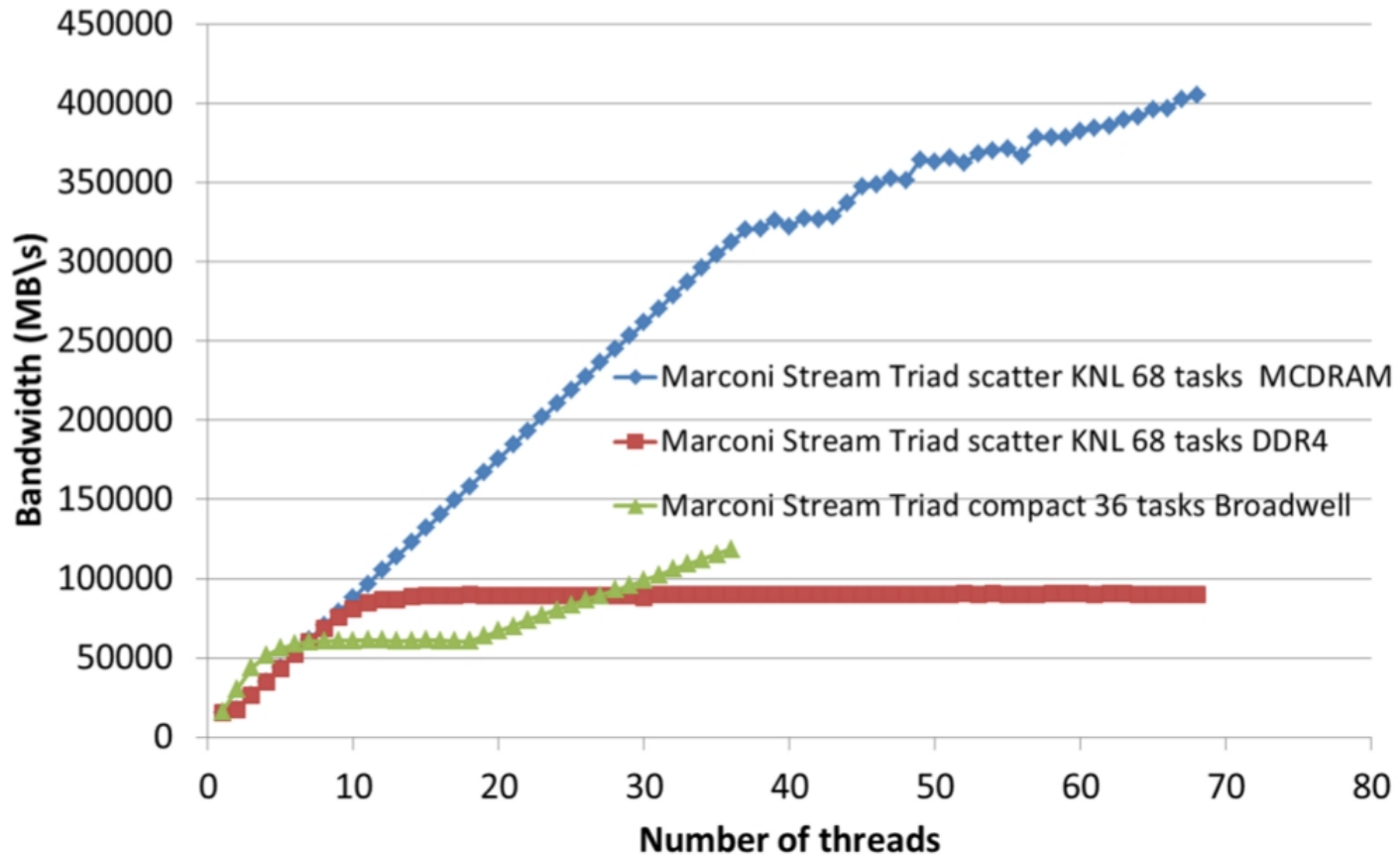
Overview

- Memory Bandwidth benchmarks
- Latency benchmarks
- OpenMP Benchmarks
- Code Performance
- Summary

STREAM and IMB

MEMORY BANDWIDTH BENCHMARKS

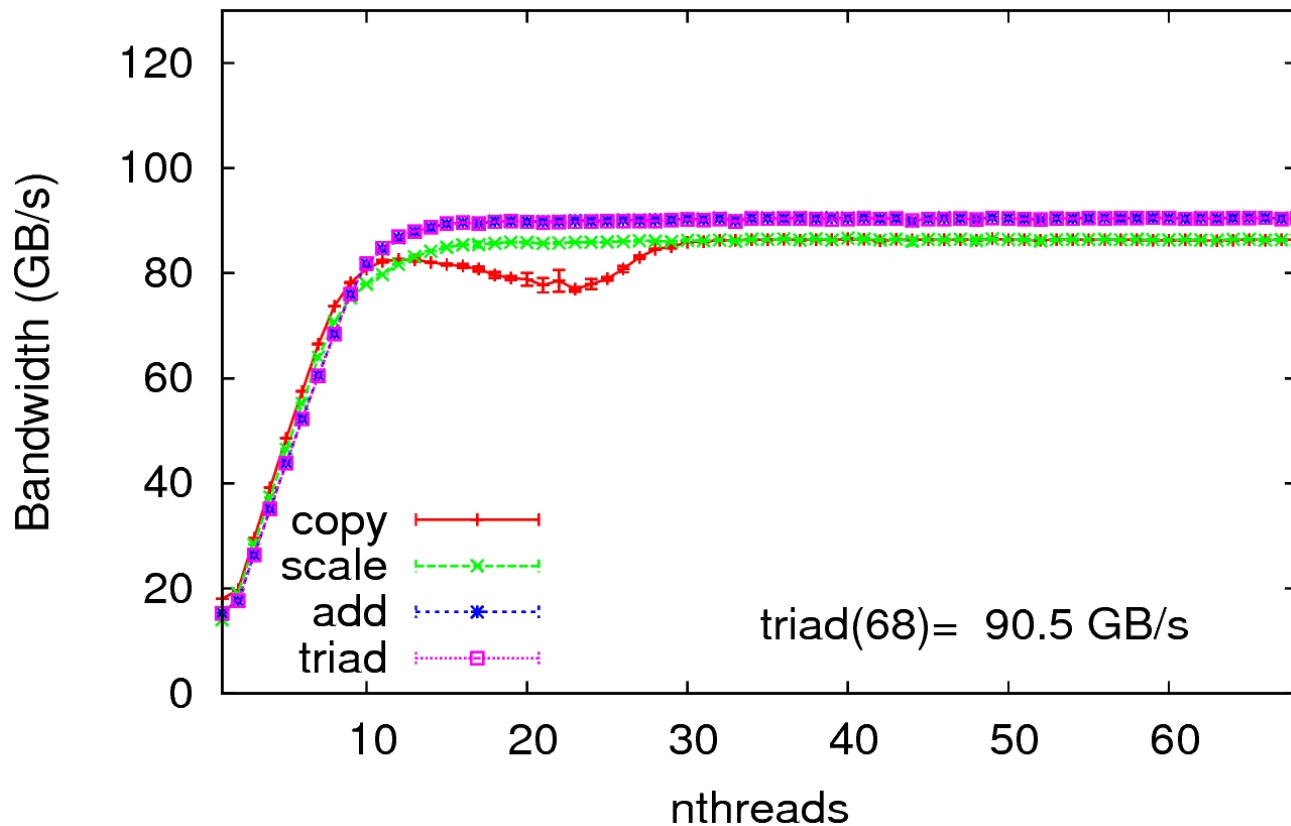
STREAM Memory Bandwidth different architectures



Using the STREAM benchmark by John D. McCalpin, <https://www.cs.virginia.edu/stream>

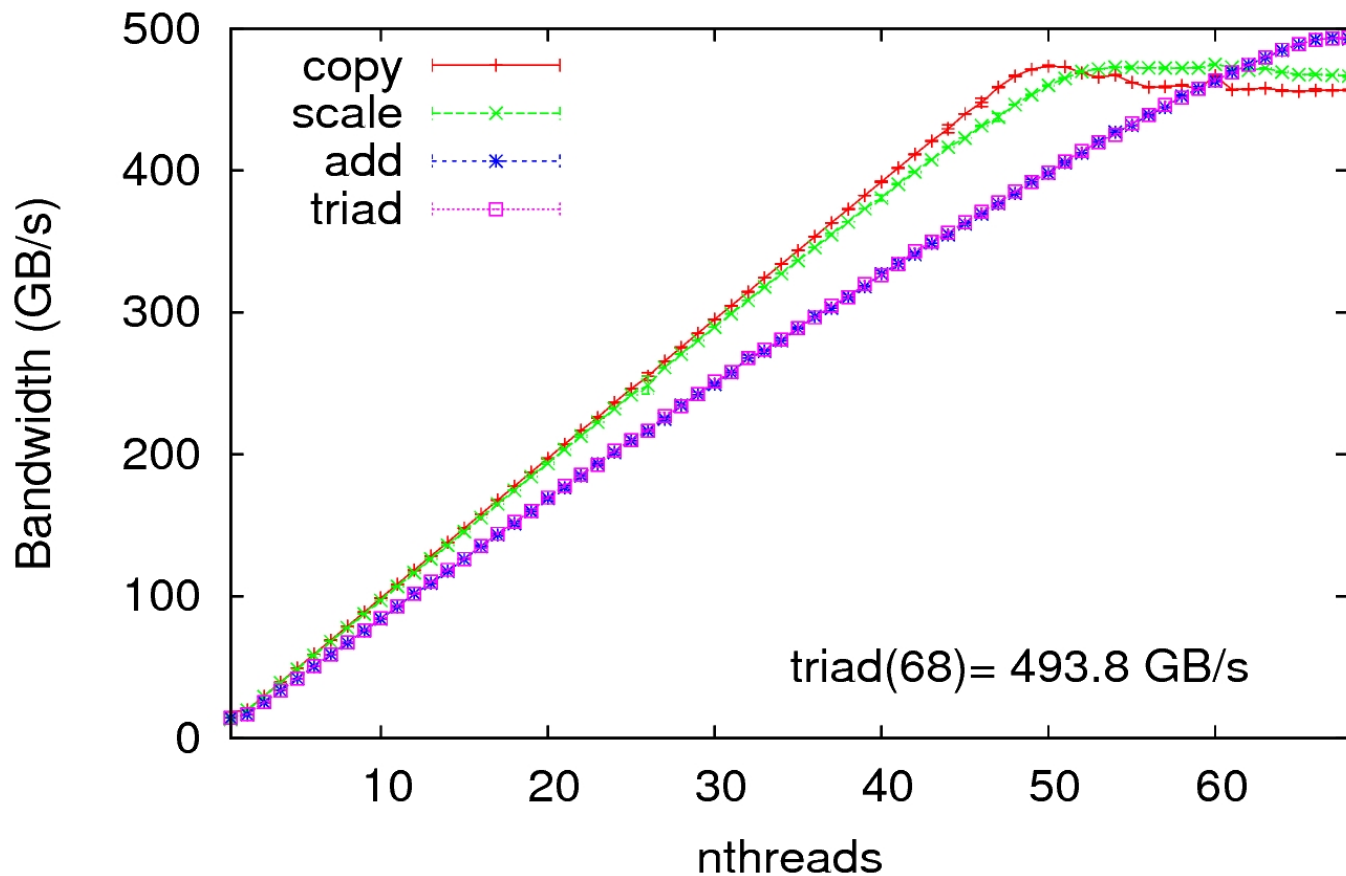
STREAM Memory Bandwidth

flat mode DDR4

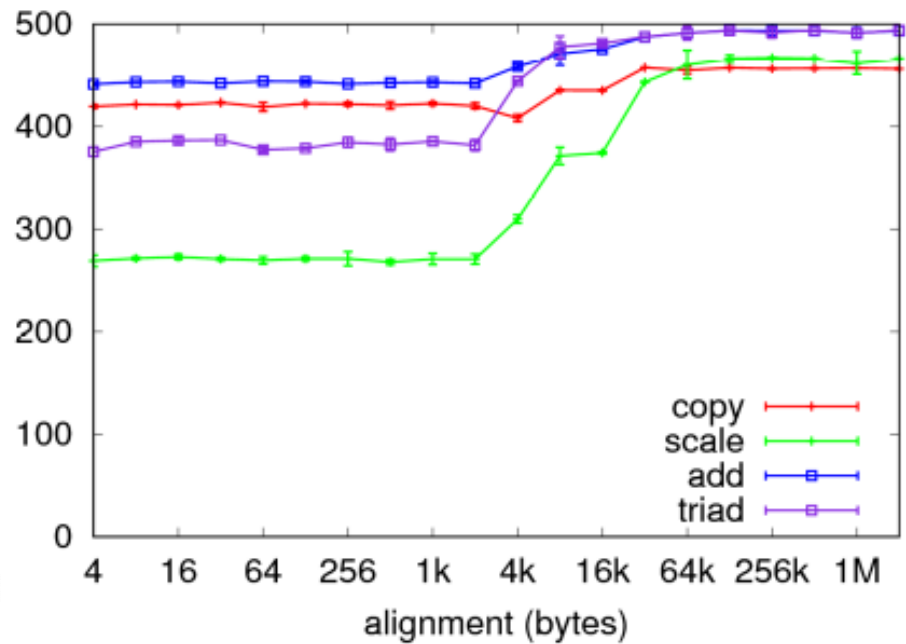
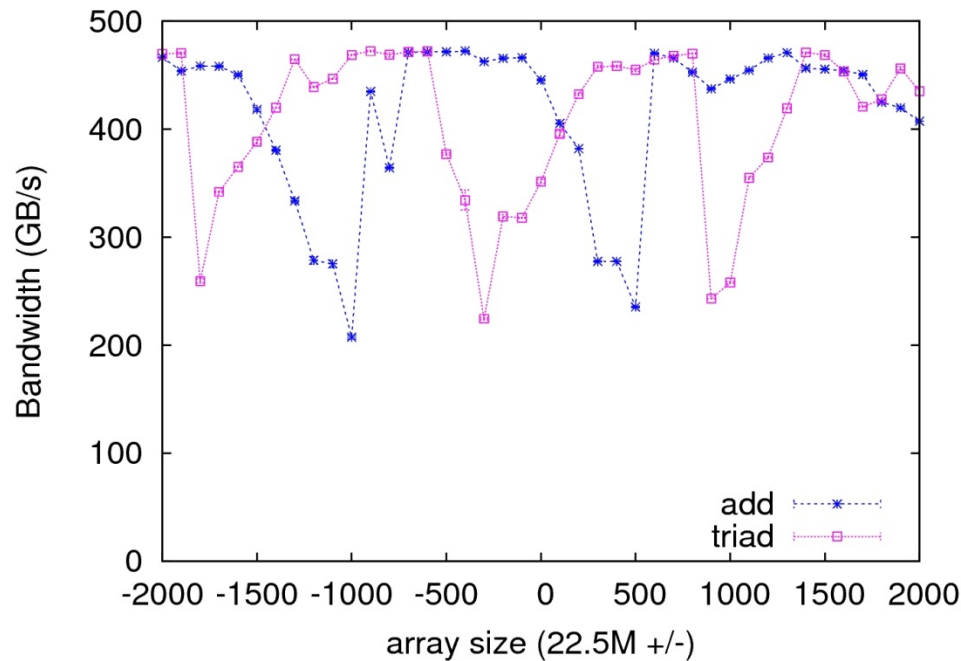


STREAM Memory Bandwidth

flat mode MCDRAM



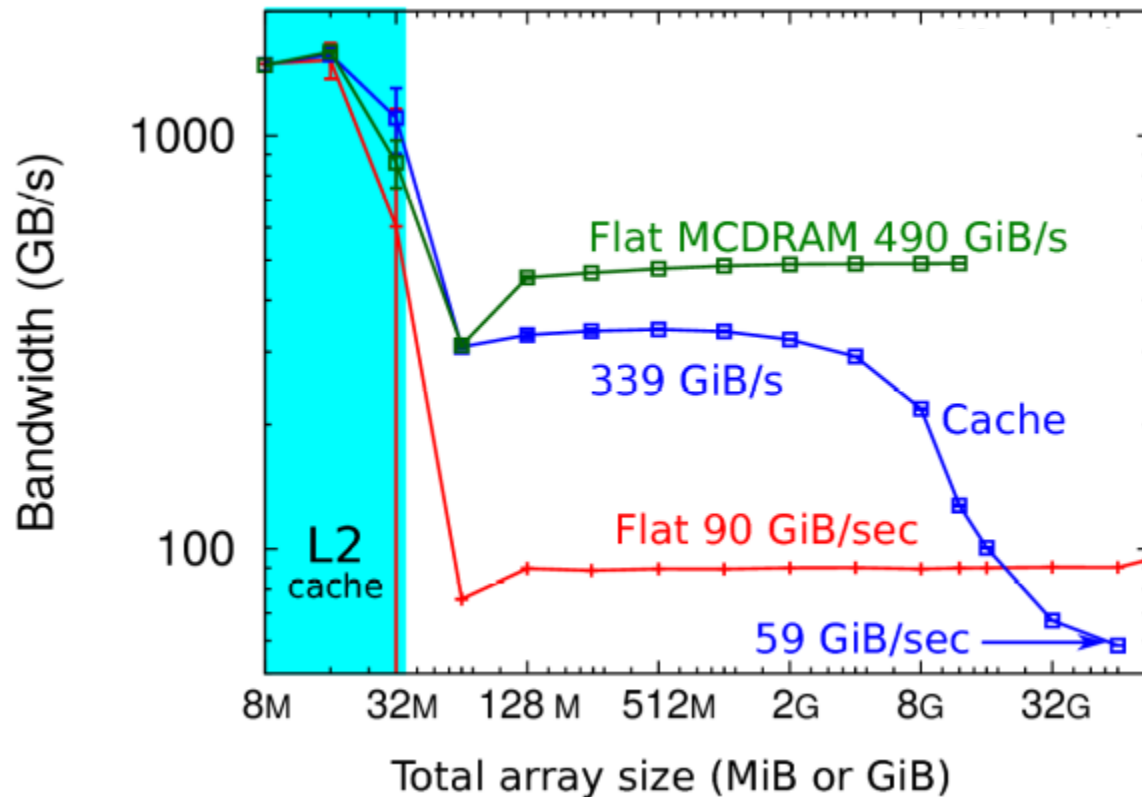
STREAM Memory Bandwidth flat mode MCDRAM - alignment



Alignment also important for cache mode

STREAM Memory Bandwidth

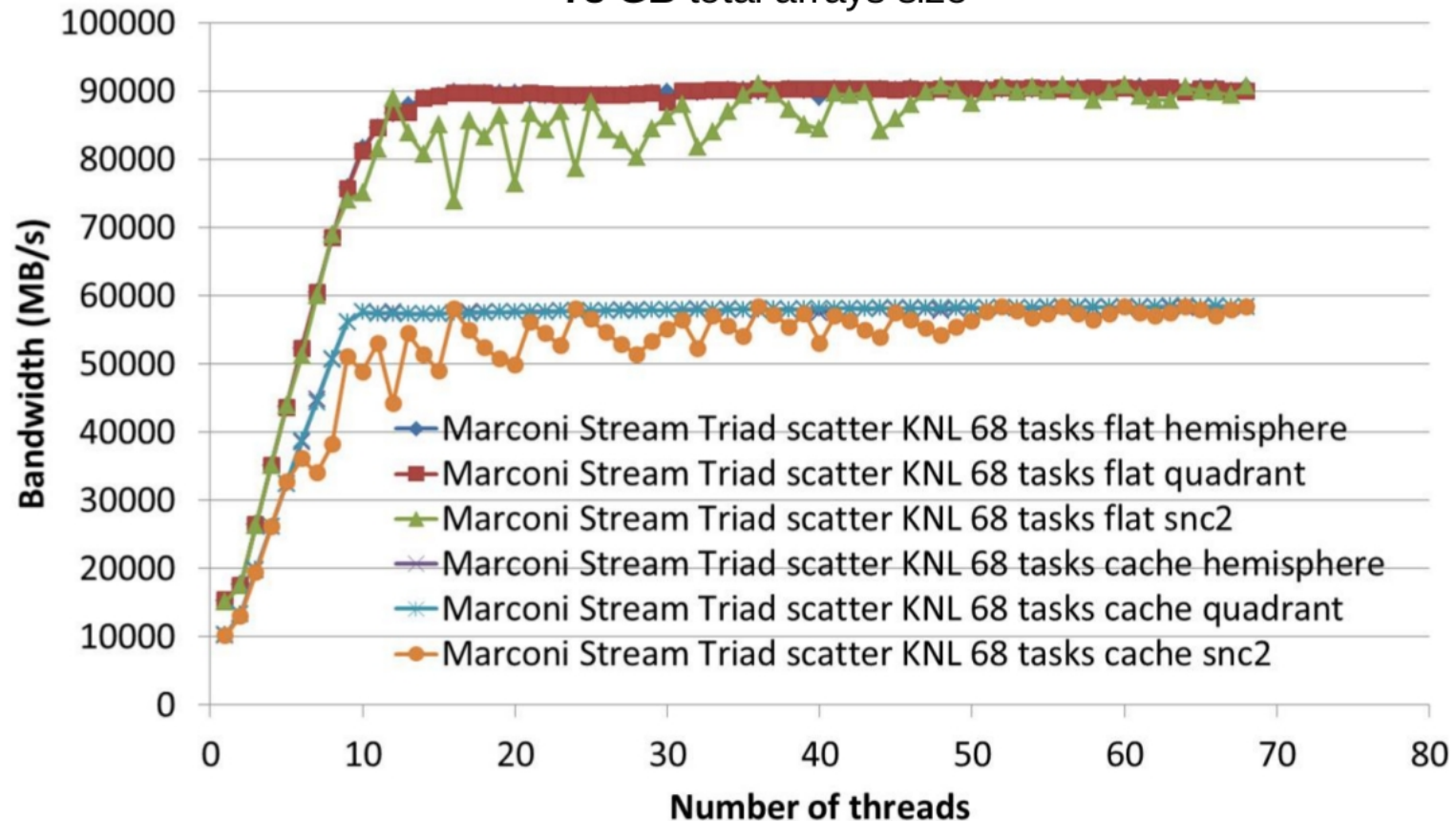
Array size



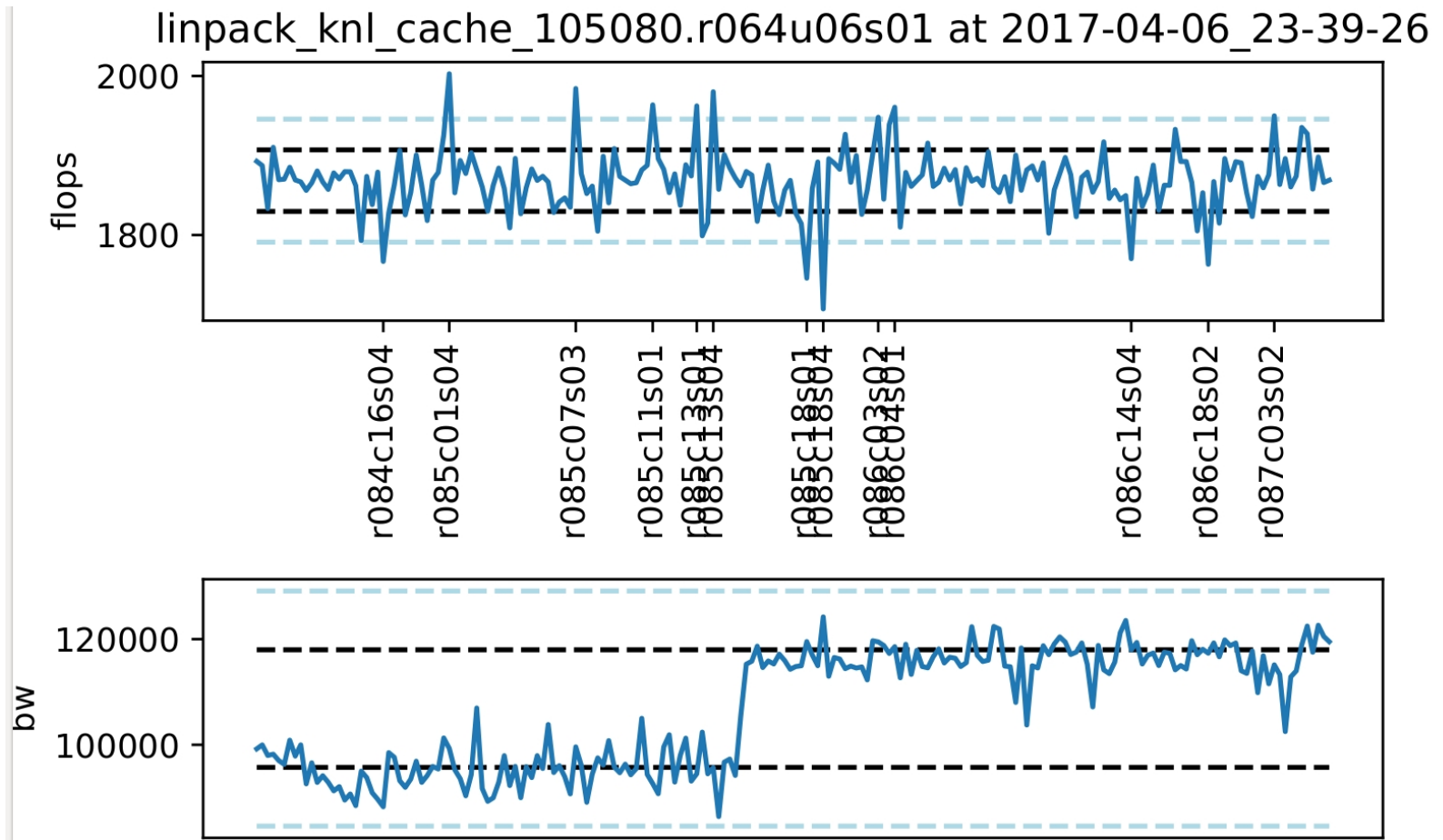
STREAM Memory Bandwidth

cache versus flat

78 GB total arrays size

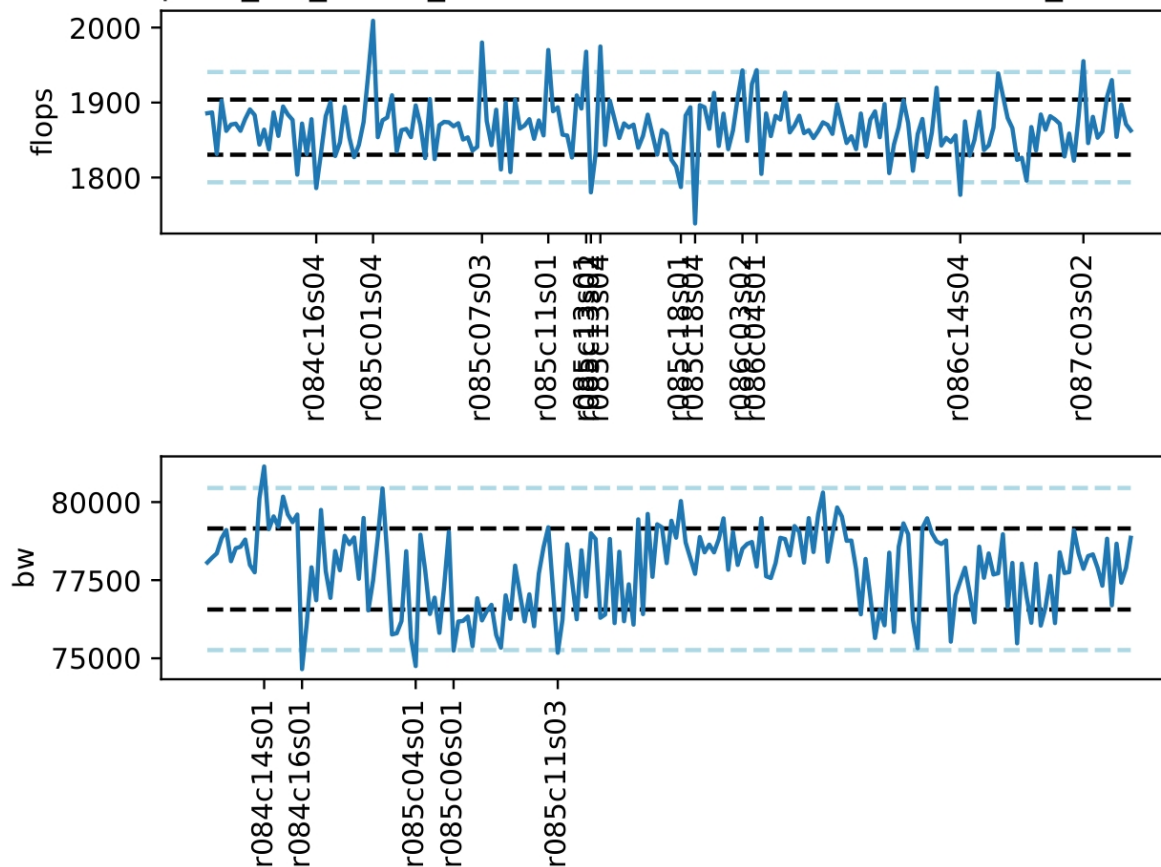


STREAM Memory Bandwidth over time

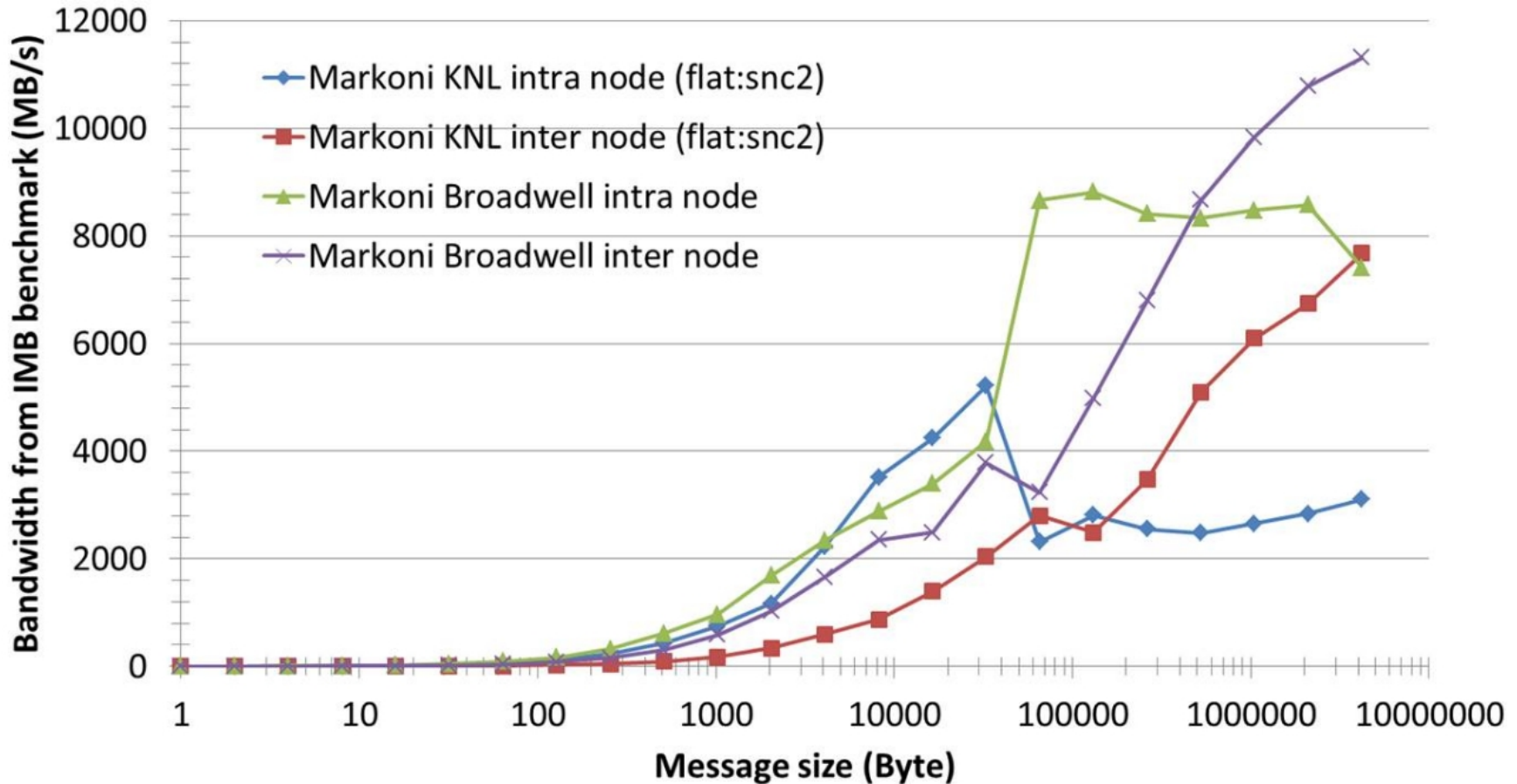


STREAM Memory Bandwidth over time

linpack_knl_cache_130796.r064u06s01 at 2017-04-14_16-10-05

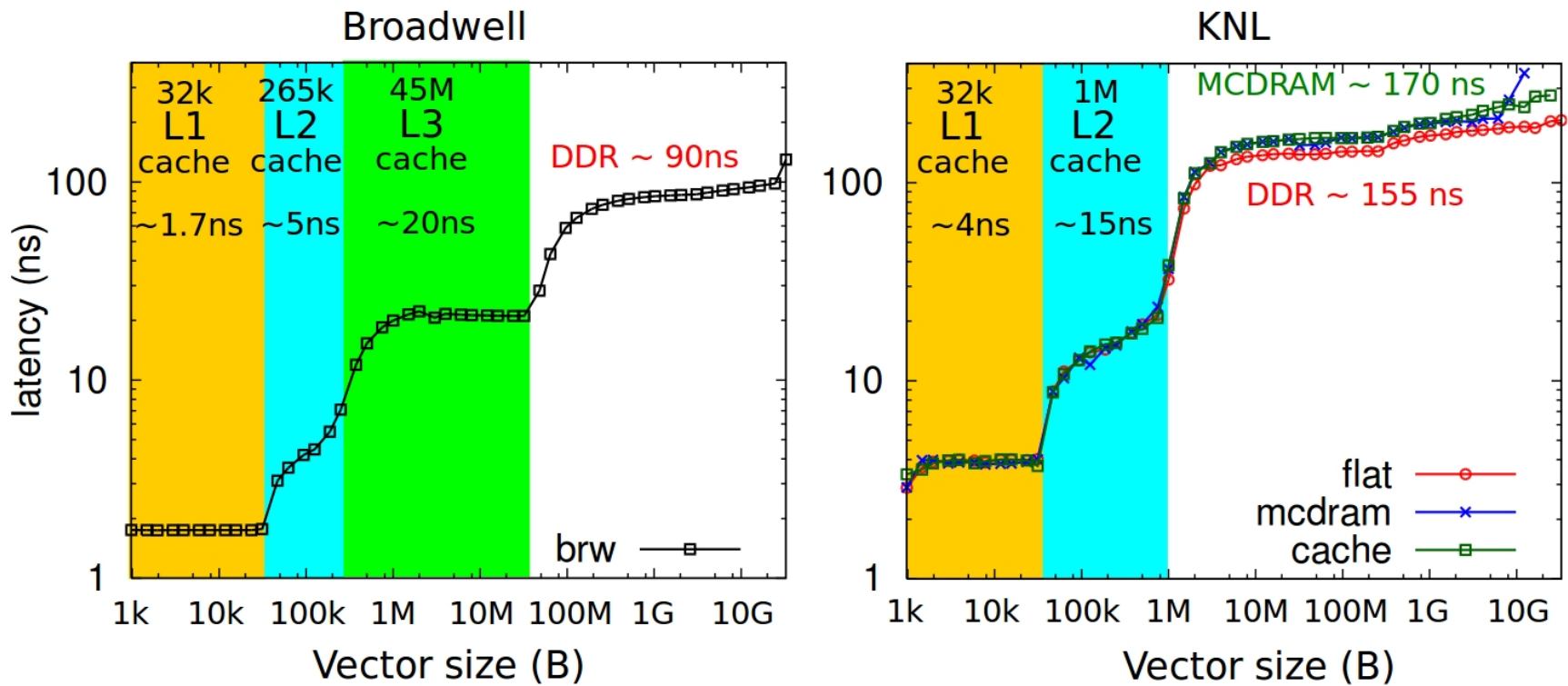


Intel MPI benchmark - bandwidth



LATENCY BENCHMARKS

Latency



IMB – Ping Pong Test - Latency

Intra node Marconi

Broadwell

node0	CPU0	0.61
	CPU1	1.09
Latency (μ s)		CPU0
		node0

Knights Landing

node0	KNL0	0.85
Latency (μ s)		KNL0
		node0

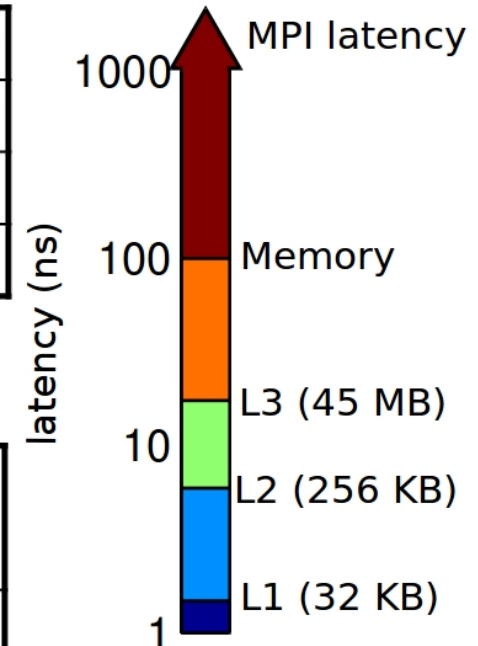
Intra node HELIOS

Sandy Bridge

node0	CPU0	0.25
	CPU1	0.64
Latency (μ s)		CPU0
		node0

Knights Corner

node0	KNC0	2.7
Latency (μ s)		KNC0
		node0



IMB – Ping Pong Test - Latency

Inter node Marconi

Broadwell

node0	CPU0	1.49
Latency (μ s)		CPU0
		node1

Inter node HELIOS

Sandy Bridge

node0	CPU0	1.13
Latency (μ s)		CPU0
		node1

Knights Landing

node0	KNL0	3.99
Latency (μ s)		KNL0
		node1

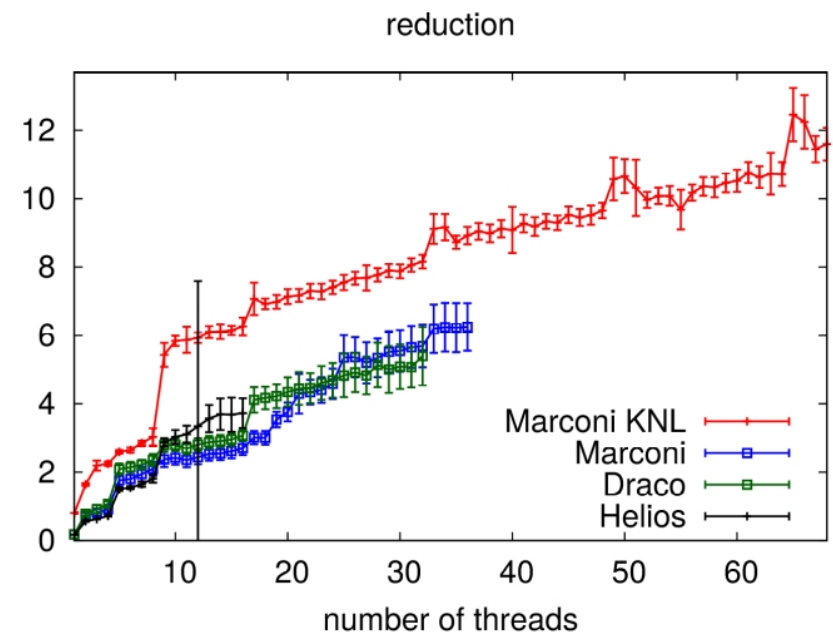
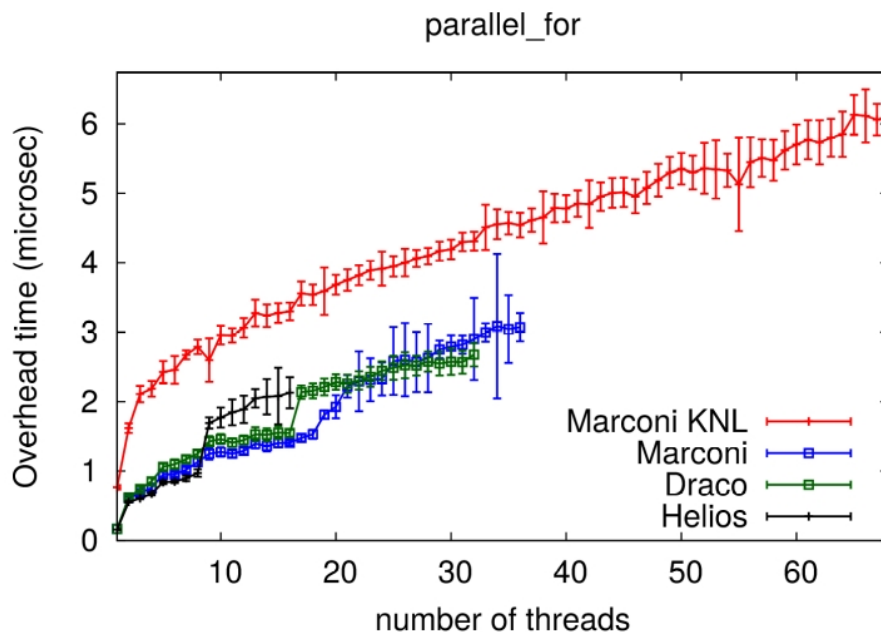
Knights Corner

node0	KNC0	6.00
Latency (μ s)		KNC0
		node1

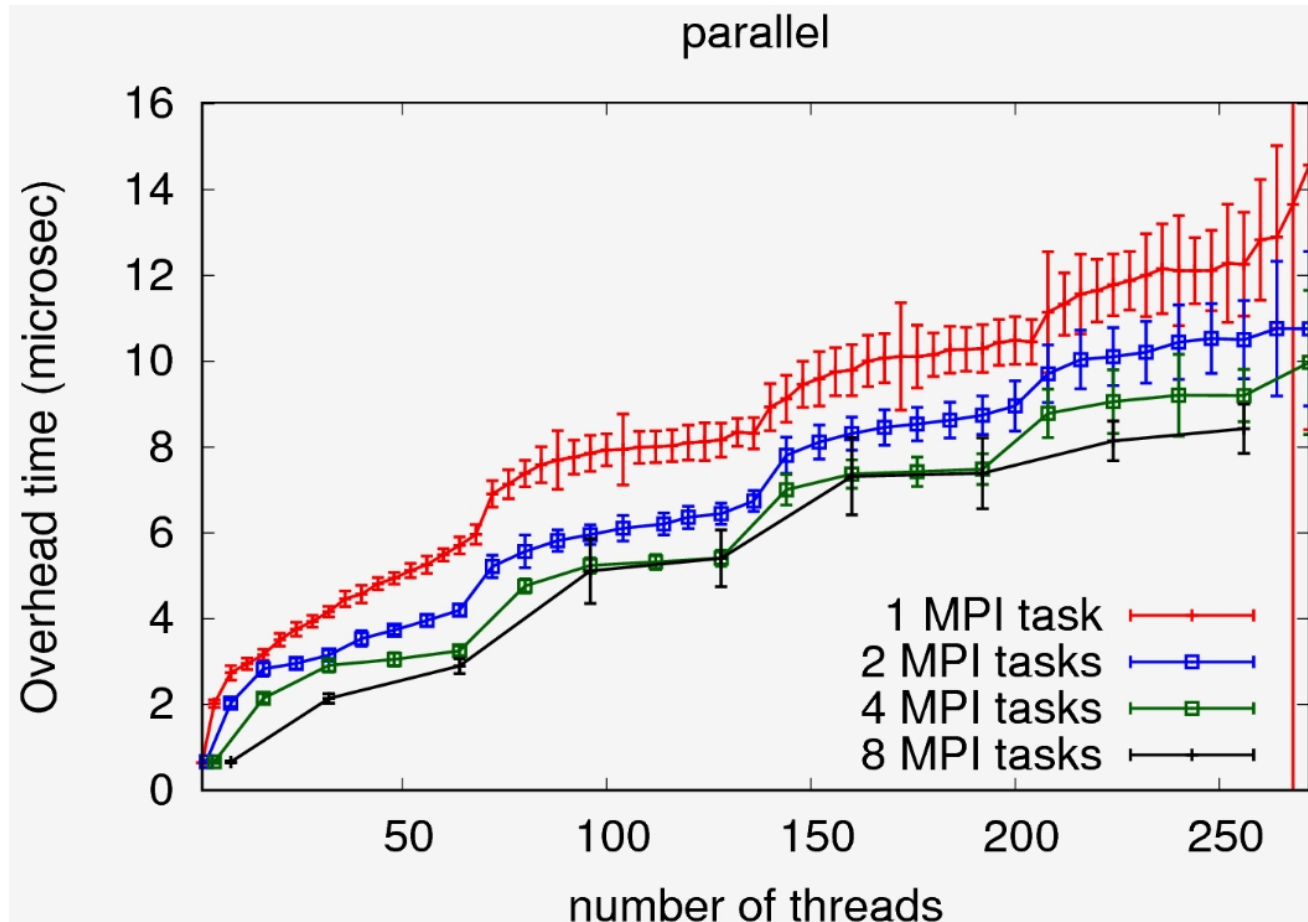
OPENMP BENCHMARKS

OpenMP overhead

- KNL overhead $\approx 2x$ larger:
 - more threads
 - lower CPU frequency
- Exception: ATOMIC 5x longer, use CRITICAL instead
- Using EPCC OpenMP Microbenchmarks J.M. Bull et. al

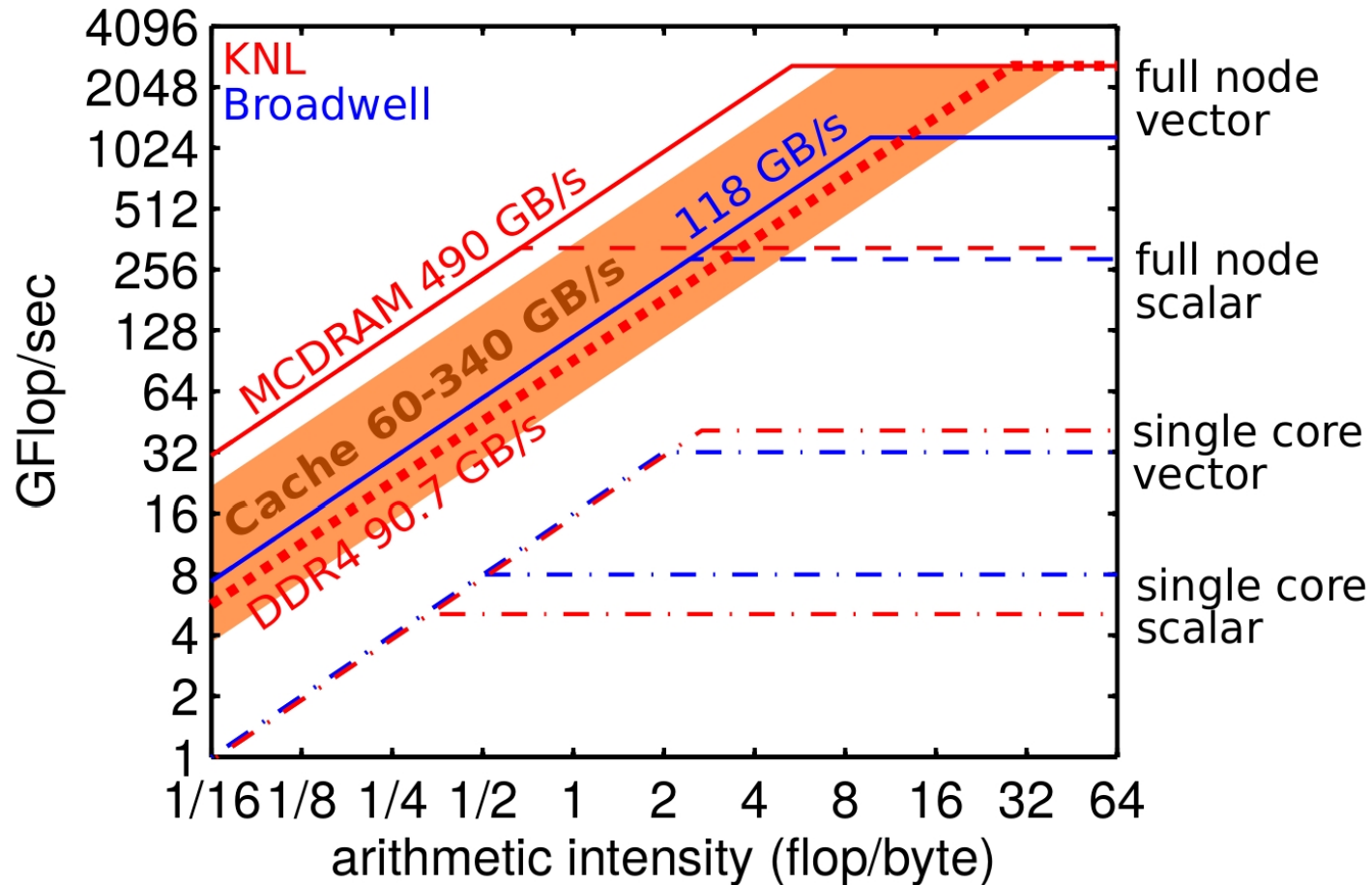


OpenMP overhead + hyper-threading



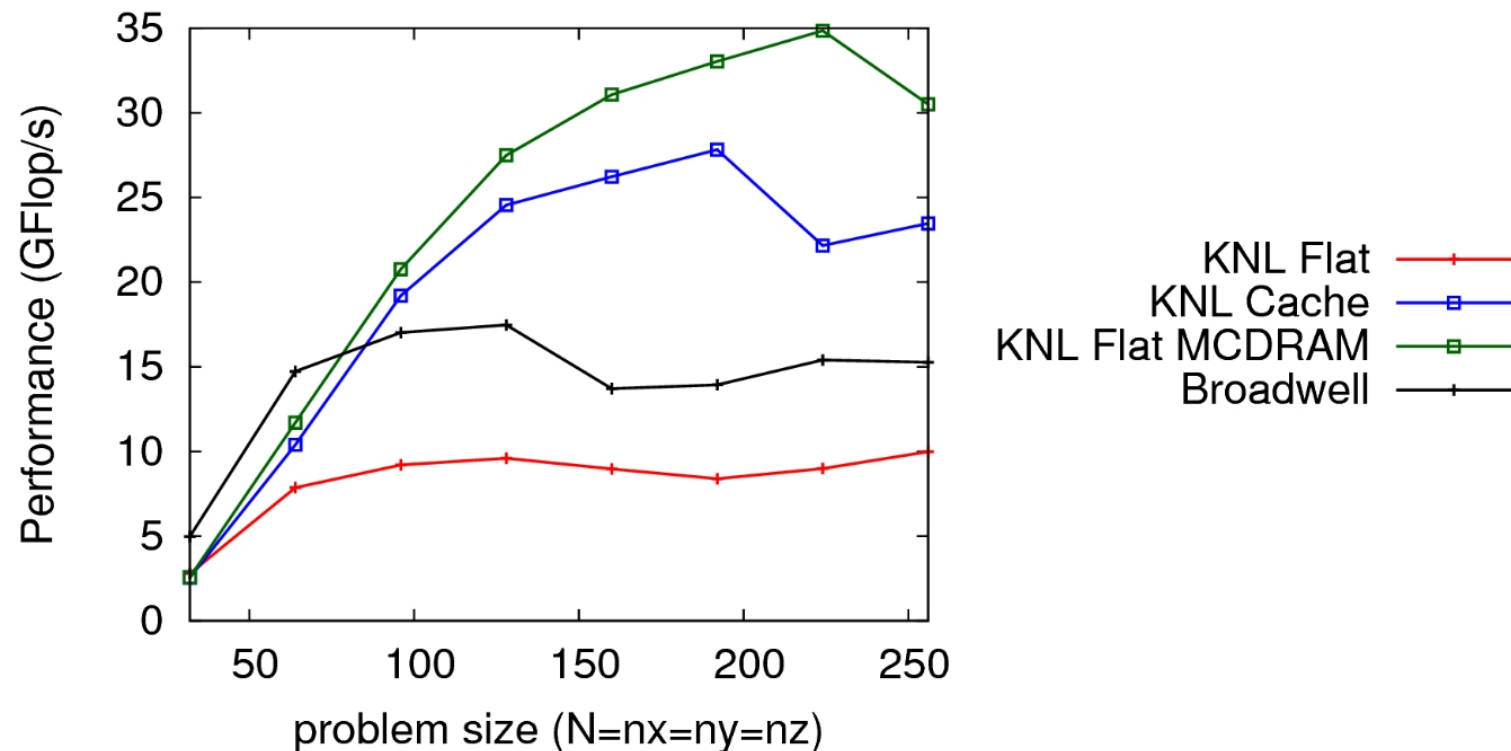
CODE PERFORMANCE

Roofline Model



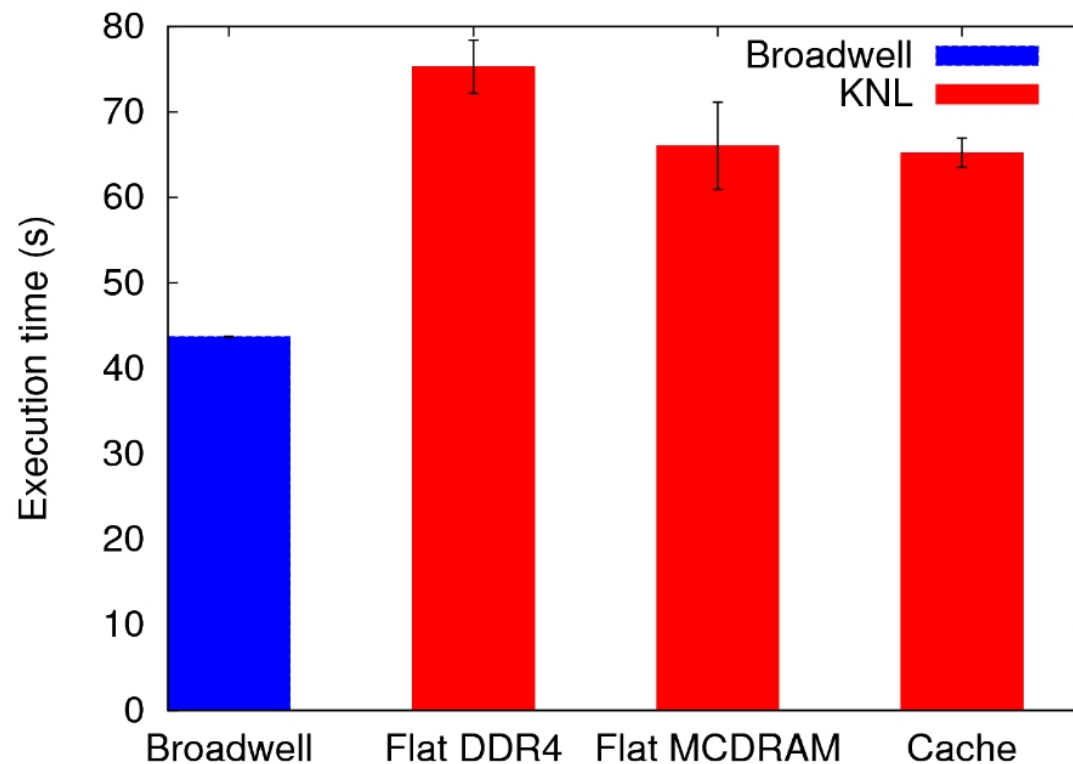
HPCG benchmark

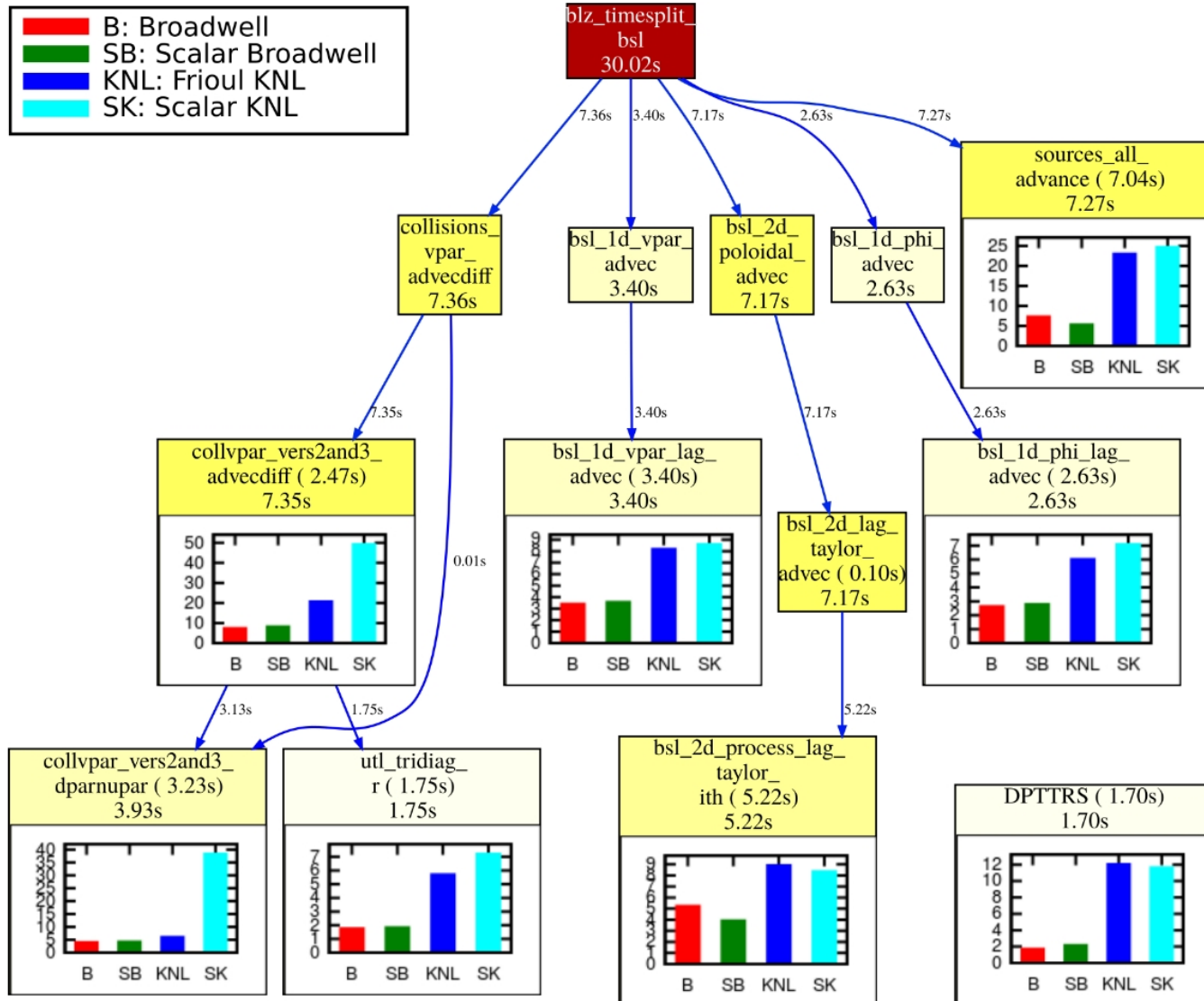
- HPCG: sparse 3D problem with multigrid preconditioned conjugate gradient solver.
- The Intel optimized version of the HPCG benchmark was executed in one node.



Gysela execution time

- Test case: 127 x 256 x 64 x 63 (Nr x Ntheta x Nphi x Nvpar, Nmu=0)
- 1 node, 4 MPI tasks, 8 threads (Broadwell) / 16 threads (KNL)

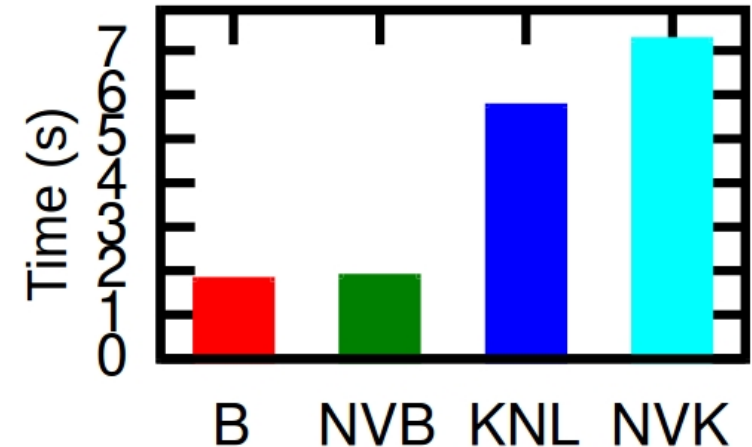




UTL_TRIDIAG_R

- Solve a tridiagonal system
 - forward elimination
 - back substitution
- Not vectorizable

$$\begin{bmatrix}
 b_1 & c_1 & & & & & 0 \\
 a_2 & b_2 & c_2 & & & & \\
 & a_3 & b_3 & \ddots & & & \\
 & & \ddots & \ddots & \ddots & & \\
 0 & & & & a_n & c_{n-1} & b_n
 \end{bmatrix}
 \begin{bmatrix}
 x_1 \\
 x_2 \\
 x_3 \\
 \vdots \\
 x_n
 \end{bmatrix}
 =
 \begin{bmatrix}
 d_1 \\
 d_2 \\
 d_3 \\
 \vdots \\
 d_n
 \end{bmatrix}$$



instruction	instruction	Broadwell latency	Broadwell throughput	KNL latency	KNL throughput
3N	FMA	5	2	6	2
2N-1	DIVSD	10-14	1/5-1/4	42	1/42
2N-1	VDIVPD	19-23	1/16	32	1/32

Source: A. Fog Instruction tables, TU Denmark 2016, <http://www.agner.org>

SUMMARY

Summary

- Good stuff:
 - MPI latency and OpenMP overhead comparable
 - KNL can match Broadwell performance without extensive tuning for most codes
 - Optimization on KNL helps on Broadwell and vice versa
- Bad Stuff:
 - Cache mode operation can be dubious
 - Peak performance hard to reach
 - Hyperthreading rarely useful

Summary

- KNL is equal to Broadwell if your code either
 - Has very good scalability (to make use of increased core count)
 - Has very good vectorization (to make use of more vector units)
 - Effectively uses only 16 GB (to make use of higher bandwidth)
- If more than one holds, you will probably get more performance than on Broadwell
- Memory mode Quadrant seems to be the best