



PARALLELISM, PERFORMANCE & OPTIMIZATION ON INTEL ARCHITECTURE

Stephen Blair-Chappell
Bayncore

OPENING STATEMENT

“Parallelism == Performance”

(leads to)

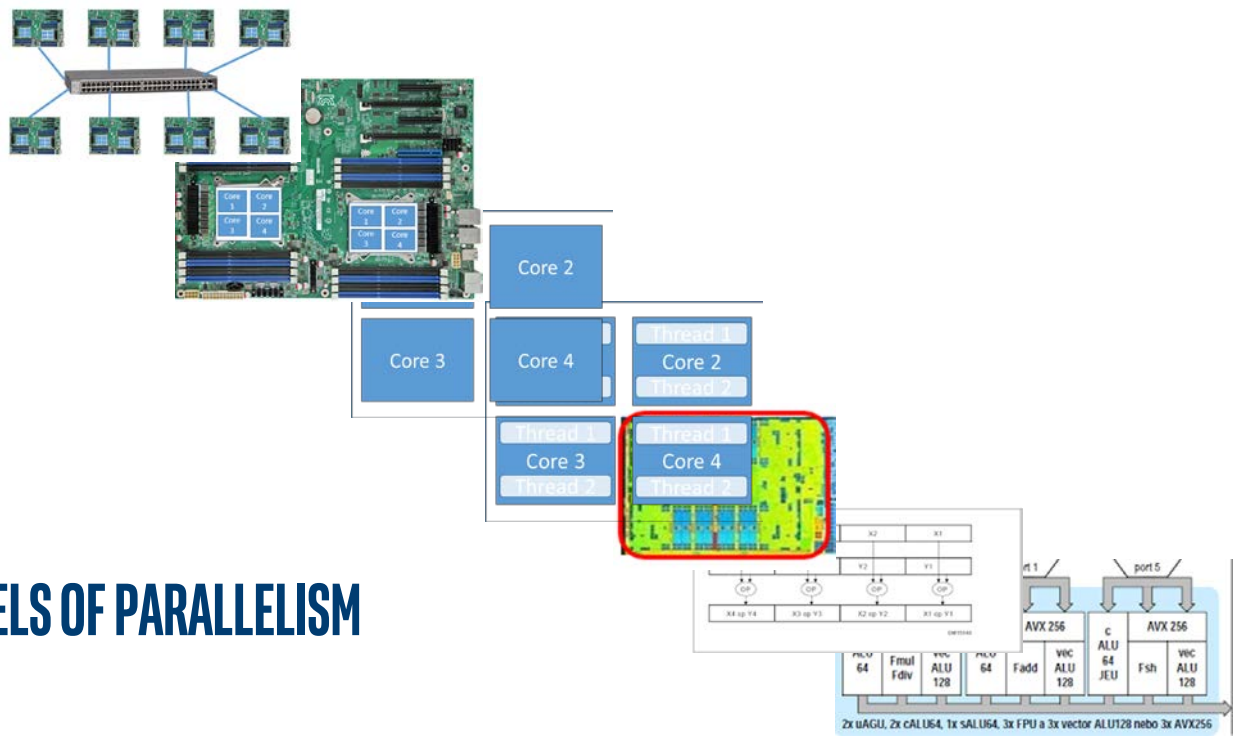
Optimization – *making sure the above statement is true!*

Desktop, Mobile & Server



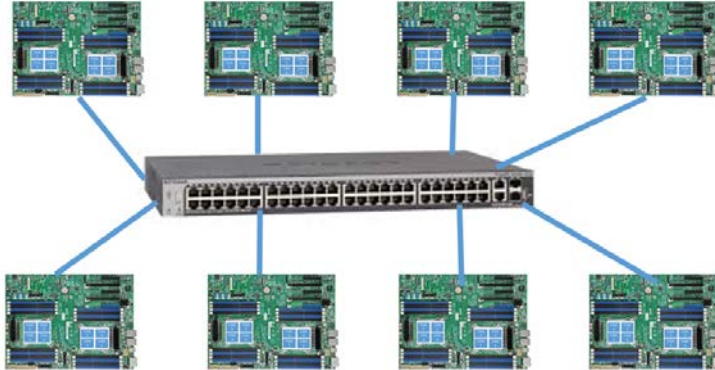


PARALLELISM ON INTEL ARCHITECTURE



SEVEN LEVELS OF PARALLELISM

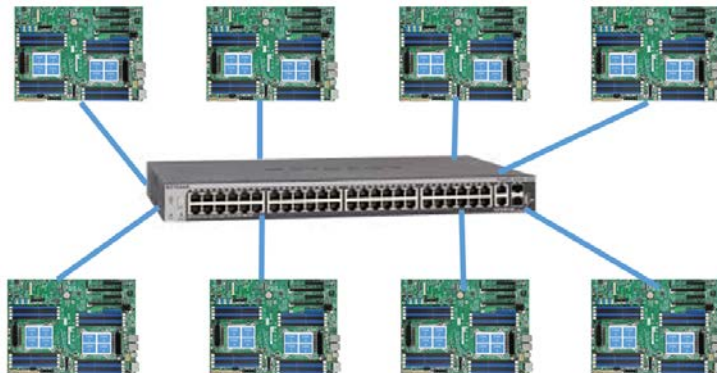
#1 - NODE-LEVEL PARALLELISM



**Levels of
Parallelism**

Node

#1 - NODE-LEVEL PARALLELISM



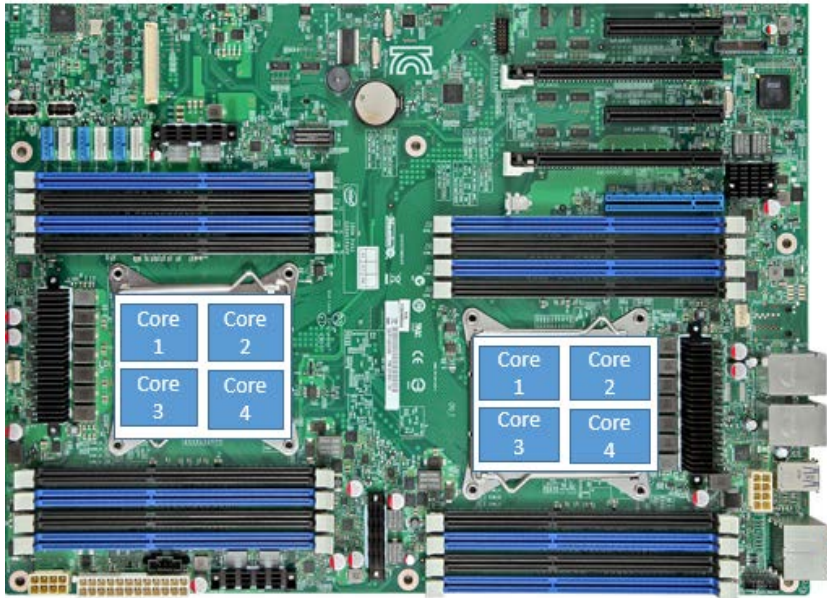
Levels of Parallelism

Node

What can I do?

Increase per-node perf.,
identify scalability issues
(Intel Trace Analyzer &
Collector), employ comm-
avoiding algorithms, etc.
(more nodes ;-)

#2 SOCKET-LEVEL PARALLELISM

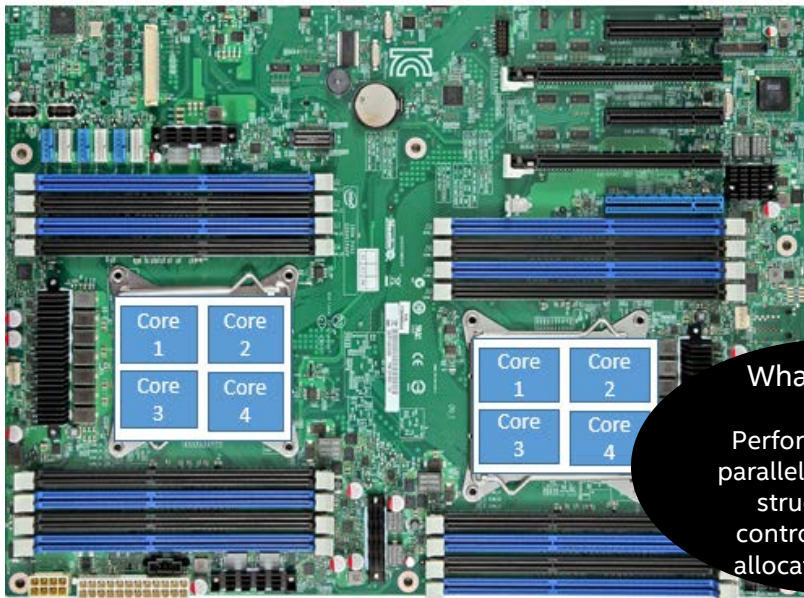


Levels of Parallelism

Node

Socket

#2 SOCKET-LEVEL PARALLELISM



Levels of Parallelism

Node

Socket

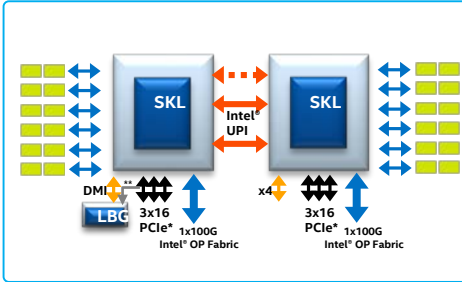
What can I do?

Perform data init. in parallel, separate data structures, take control with NUMA allocator (libnuma).

Do I have NUMA-issues?

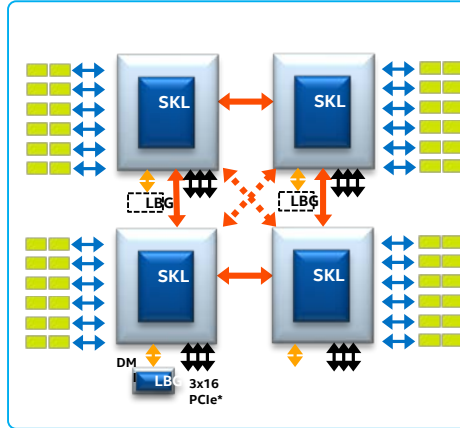
Check with 'numactl -i all', or Intel Vtune.

2S Configurations



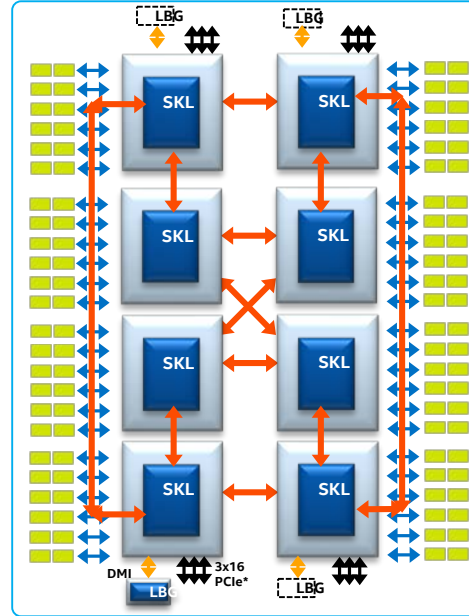
(2S-2UPI & 2S-3UPI shown)

4S Configurations



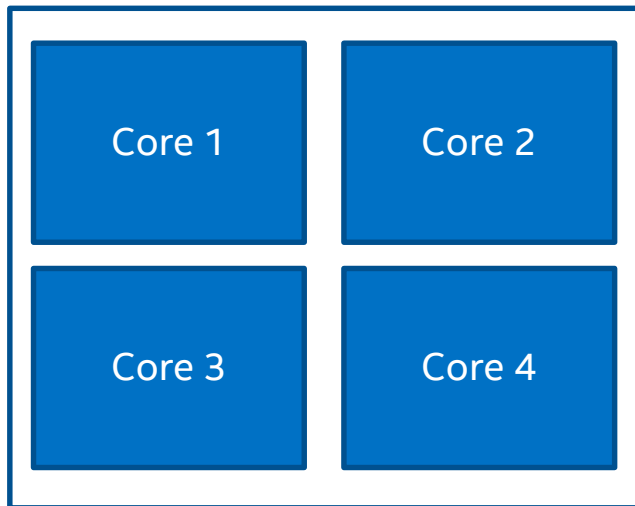
(4S-2UPI & 4S-3UPI shown)

8S Configuration



**INTEL® XEON® SCALABLE PROCESSOR
SUPPORTS UP TO 8 SOCKETS**
WITHOUT THE NEED FOR AN ADDITIONAL NODE CONTROLLER

#3 CORE / THREAD LEVEL PARALLELISM



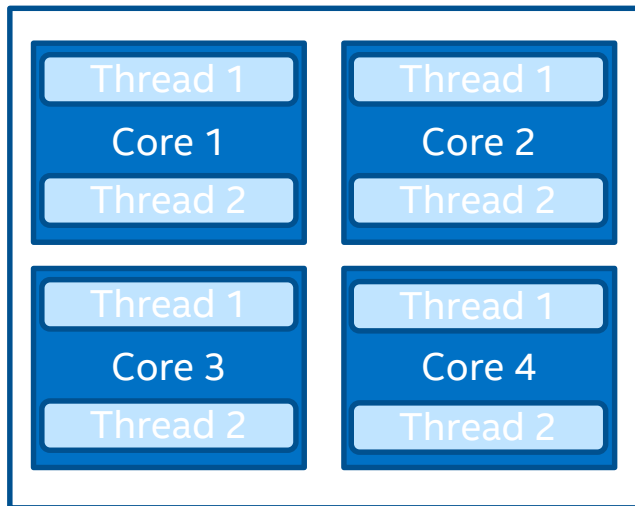
Levels of Parallelism

Node

Socket

Core / Thread-Level

#4 THREAD-LEVEL PARALLELISM (WITH HYPERTHREADING AKA SMT)



Levels of Parallelism

Node

Socket

Core / Thread-Level

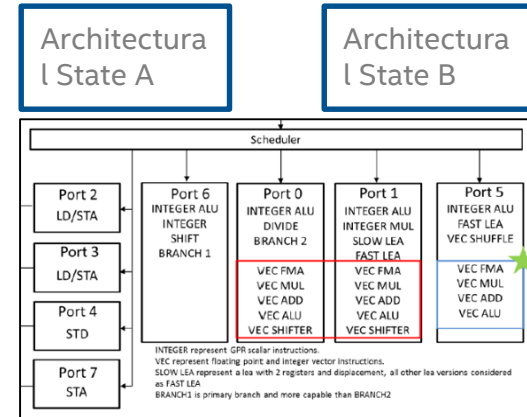
(Hyperthreading)

HYPERTHREADING

- Multiple buffers keep arch. state
- Shared execution blocks
- Enabled by BIOS settings

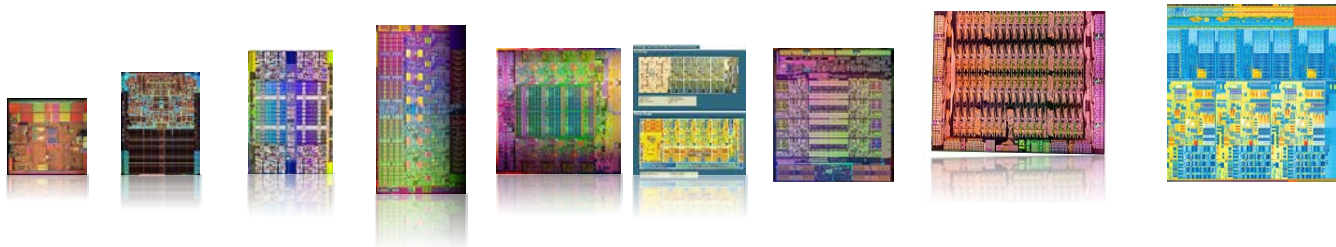
Background

- Extracts Instruction Level Parallelism (ILP)
- Complements out-of-order execution
- Intel Core: intra-core slowdowns due to HT are eliminated, BUT slowdown may happen due missing thread-affinization or because of synchronization (locks).



Single Execution Block

PARALLELISM ON INTEL® ARCHITECTURE (XEON SERVER CORES)

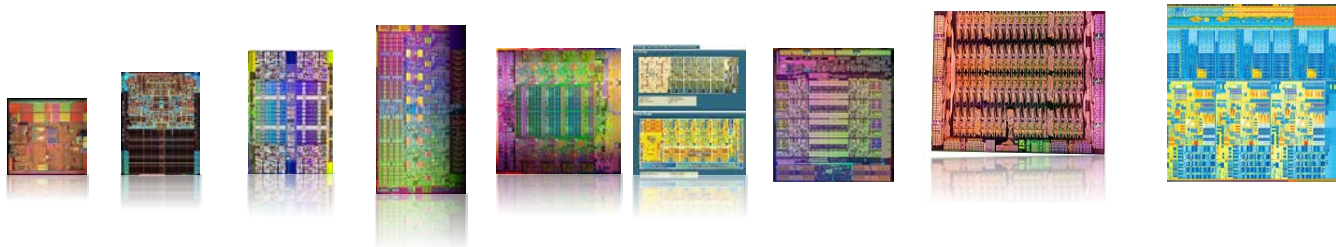


	Intel® Xeon				Intel® Xeon E5				Intel® Xeon scalable processors
	64-bit	5100	5500	5600	SNB	IVB	HSW	BDW	SKX
Cores	1	2	4	6	8	10	18	22	28
Threads	2	2	8	12	16	20	36	44	56

Intel® Xeon Phi Coprocessor	Intel® Xeon Phi 2nd gen.
KNC	KNL
61	72
244	288

IVB: Ivy Bridge BDW: Broadwell KNC: Knights Corner SKX: Skylake Server core refines
 HSW: Haswell (SKL: Skylake) KNL: Knights Landing SKL client core significantly

PARALLELISM ON INTEL® ARCHITECTURE (XEON SERVER CORES)

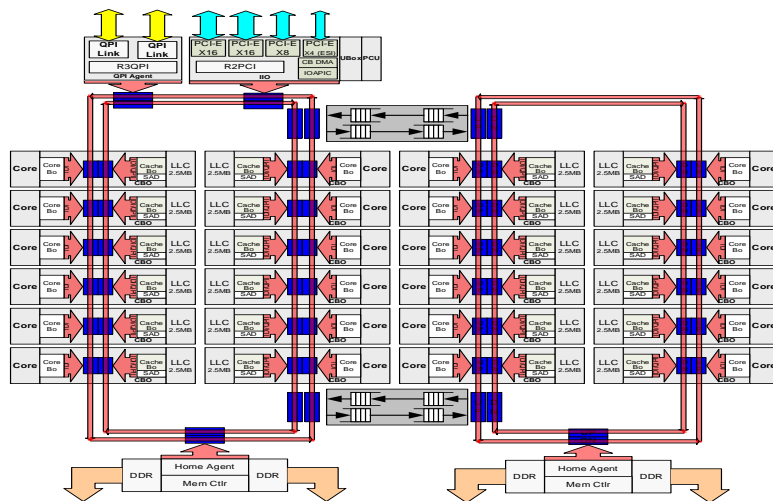


	Intel® Xeon				Intel® Xeon E5				Intel® Xeon scalable processors		
	64-bit	5100	5500	5600	SNB	IVB	HSW	BDW	SKX	Intel® Xeon Phi Coprocessor	Intel® Xeon Phi 2nd gen.
										KNC	KNL
Cores	1	2	4	6	8	10	18	22	28	61	72
Threads	2	2	8	12	16	20	36	44	56	244	288

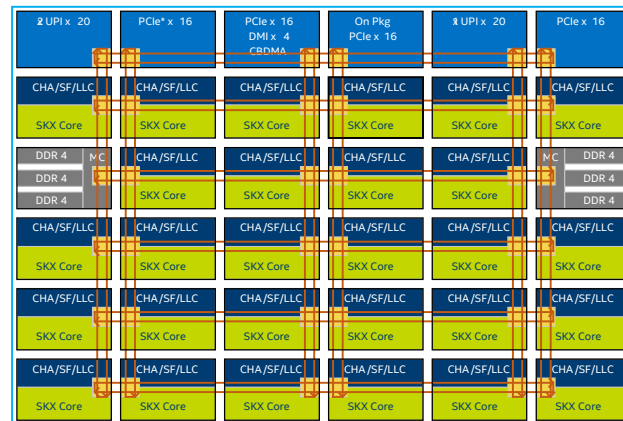
IVB: Ivy Bridge BDW: Broadwell KNC: Knights Corner SKX: Skylake Server core refines
 HSW: Haswell (SKL: Skylake) KNL: Knights Landing SKL client core significantly

NEW MESH INTERCONNECT ARCHITECTURE

Intel® Xeon® Processor E7 family (24-core die)



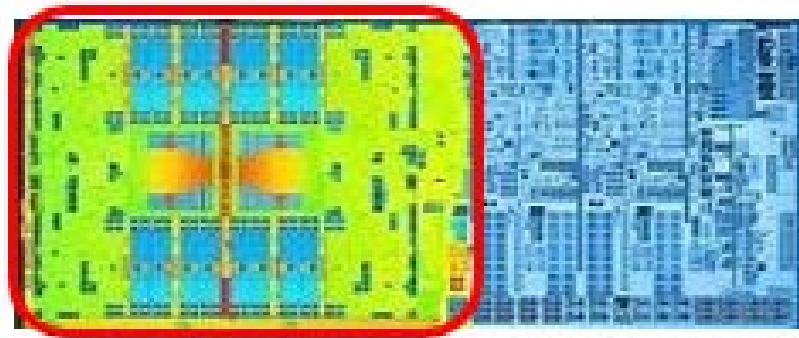
Intel® Xeon® Scalable Processor (28-core die)



CHA Gating and Home Agent ; SF -Snoop Filter ; LLC Last Level Cache ;
SKX Core Skylake Server Core ; UPI Intel® UltraPath Interconnect

MESH IMPROVES SCALABILITY WITH HIGHER BANDWIDTH AND REDUCED LATENCIES

#5 GPU-CPU PARALLELISM



Levels of Parallelism

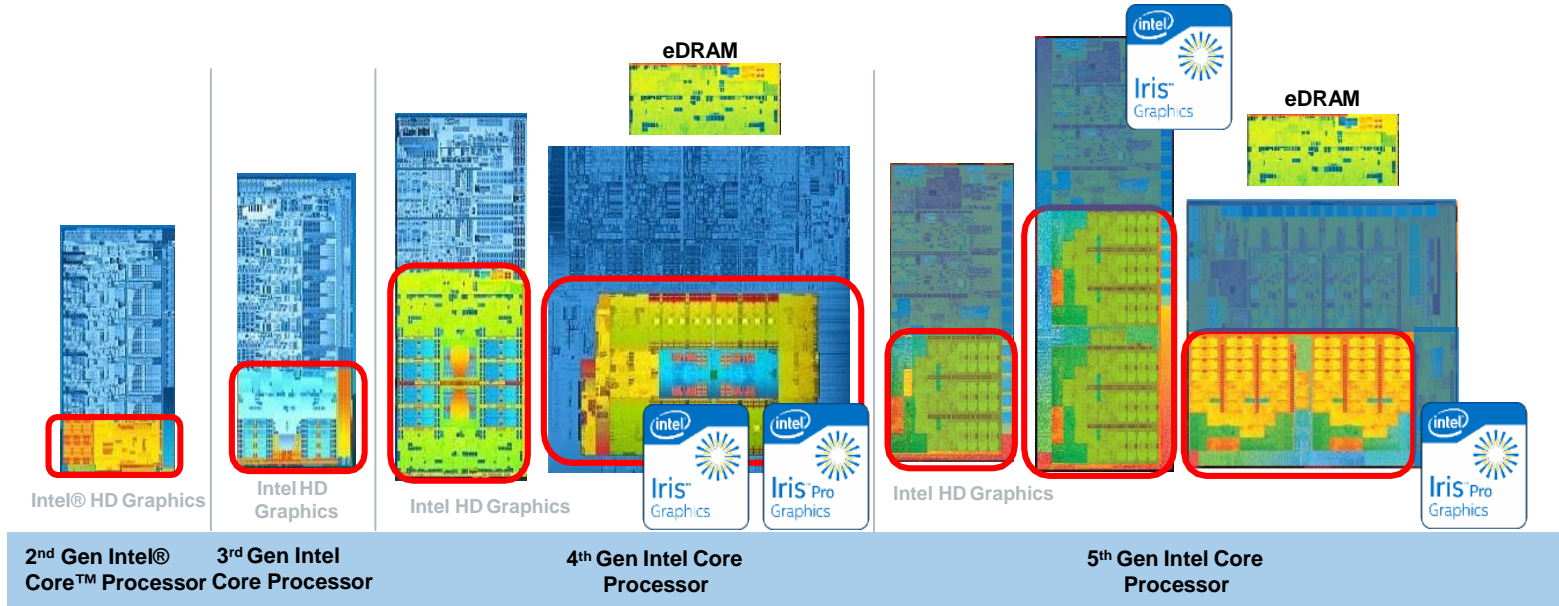
Node

Socket

Core / Thread-Level
(Hyperthreading)

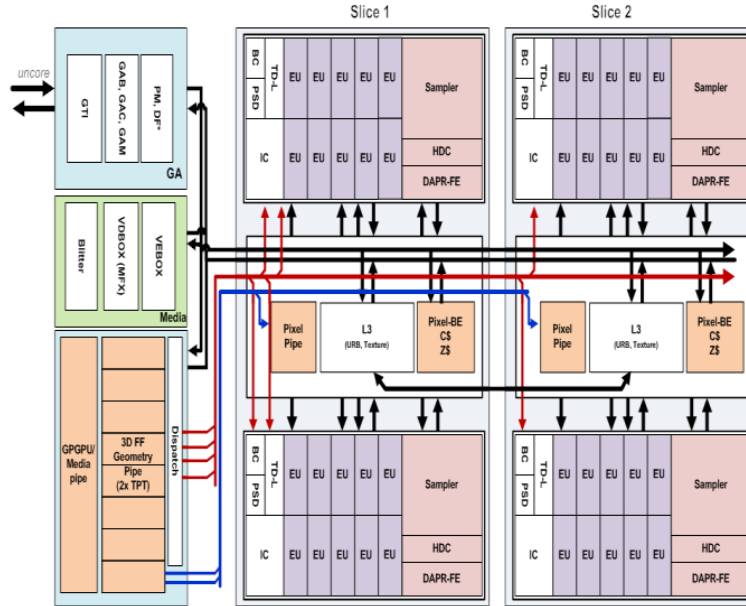
GPU-CPU

INTEGRATED PROCESSOR GRAPHICS



Lots of compute power for data-parallel (client-)applications

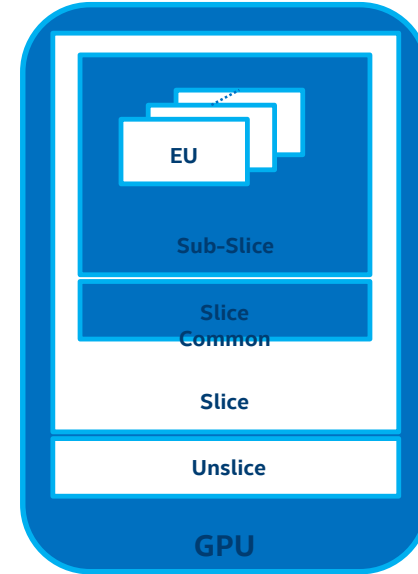
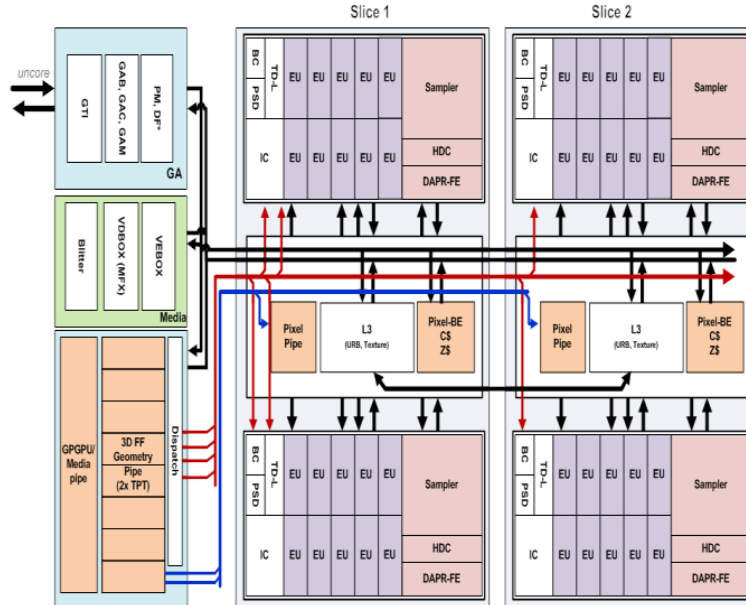
4TH GENERATION INTEL® CORE™ PROCESSORS – HD GRAPHICS



- 5 Execution Units Perrow
- 2 rows per subslice
- 2 subslices per slice
- 2 slices (40 EUs total) in GT3
- 7 Threads PerEU
- 280 threads in GT3
- 128 Registers per thread
- 4KB per thread!
- 1120KB in regfile in GT3
- 256KB data cache per slice (L3 only)

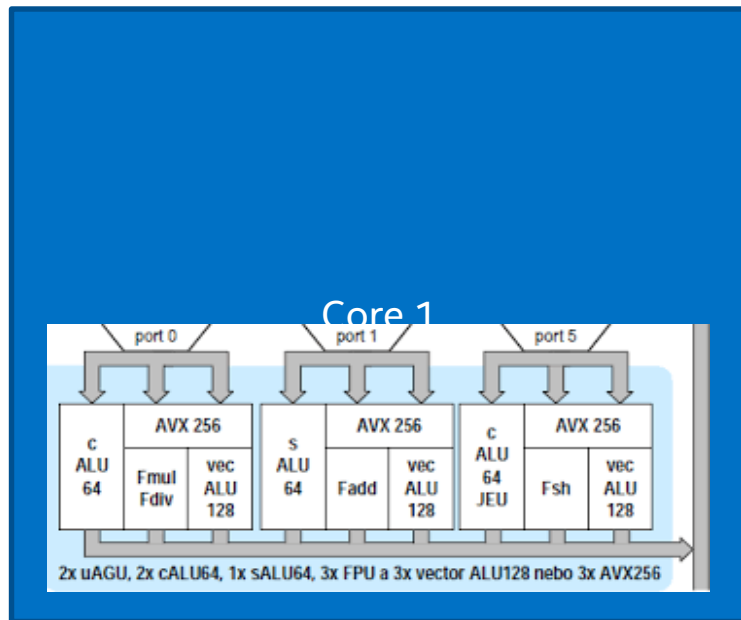
Since SKL, capabilities for GPU-CPU exchange increased (with OpenCL 2.0 memory coherency can be coarse and fine-grained)

4TH GENERATION INTEL® CORE™ PROCESSORS – HD GRAPHICS



Conceptual structure of the Intel GPU

#6 INSTRUCTION-LEVEL PARALLELISM



Levels of Parallelism

Node

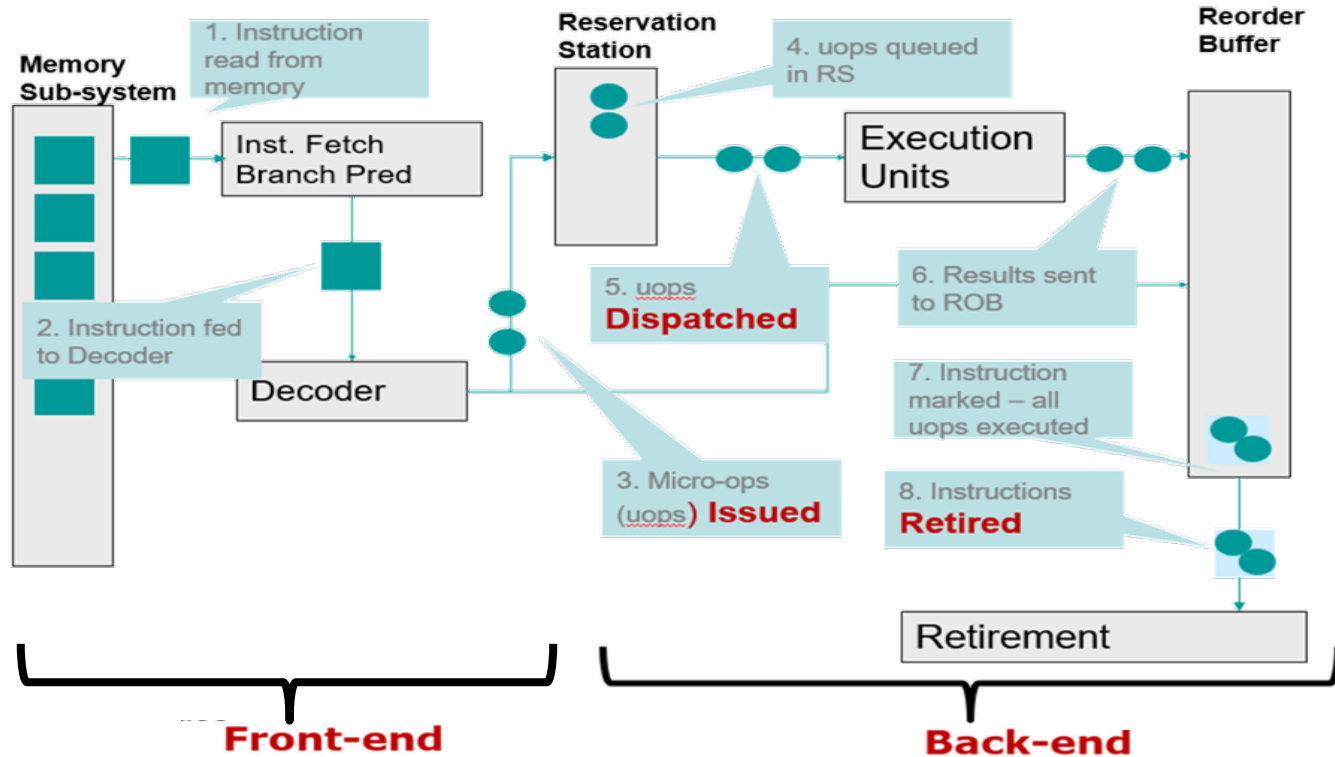
Socket

Core / Thread-Level
(Hyperthreading)

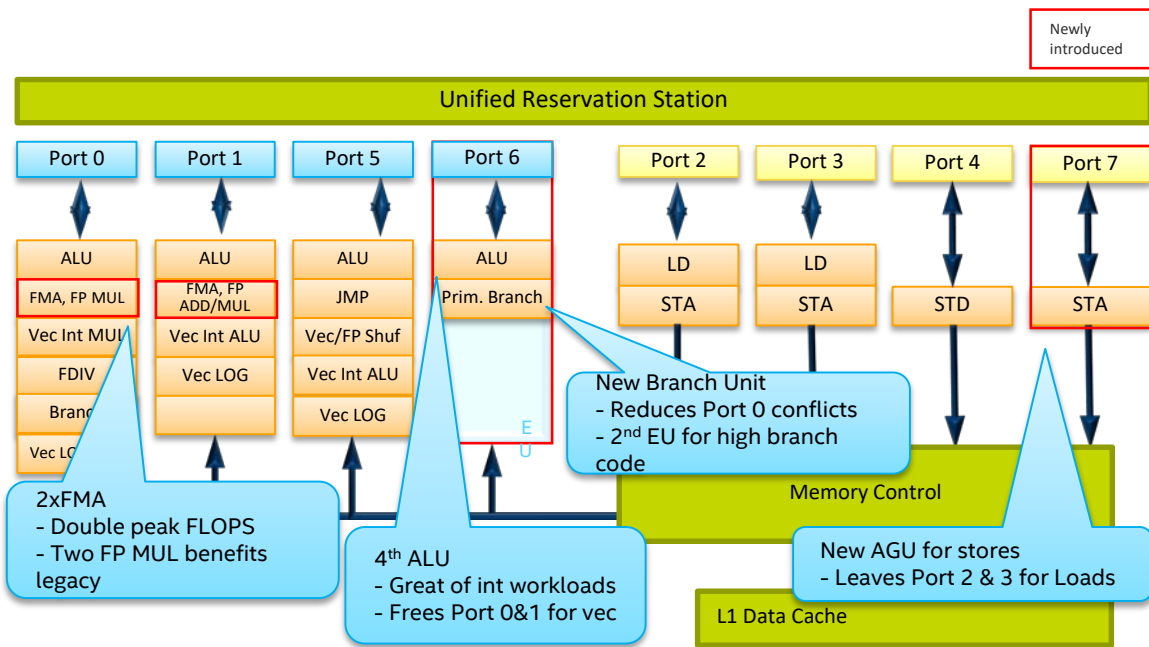
GPU-CPU

Instruction

The life of a program instruction

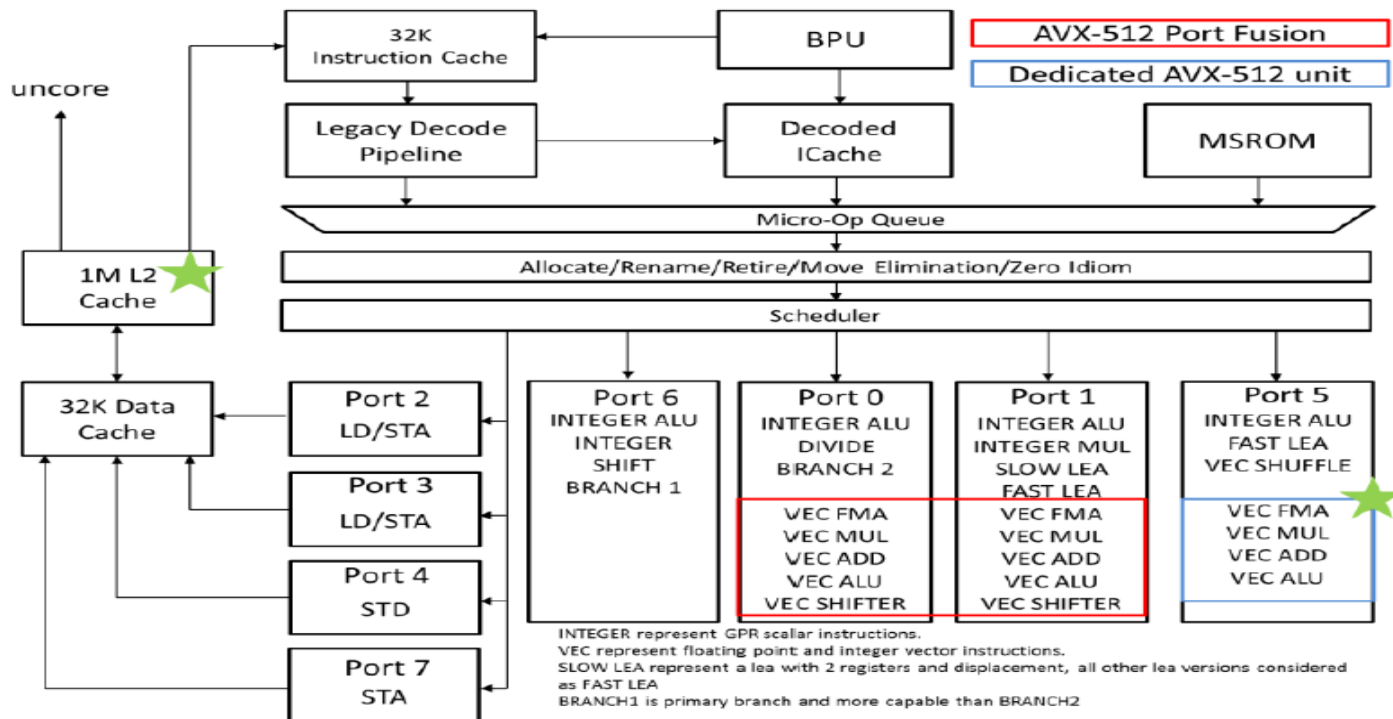


EXECUTION UNITS ON HASWELL/BROADWELL

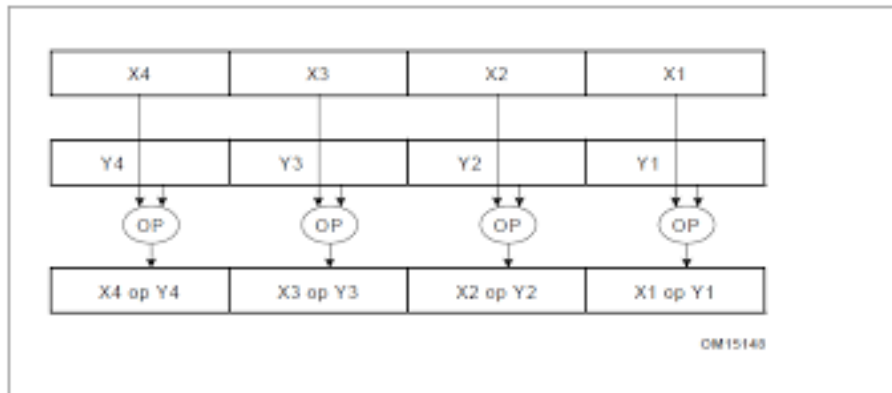


- Race to higher frequencies slowed down
- New logic introduced with new generations, leading to higher complexity
- Fused Multiply Add for performance
- Separate Address Generation Unit for memory address calculations

EXECUTION UNITS ON SKYLAKE SERVER



#7 DATA LEVEL PARALLELISM



Levels of Parallelism

Node

Socket

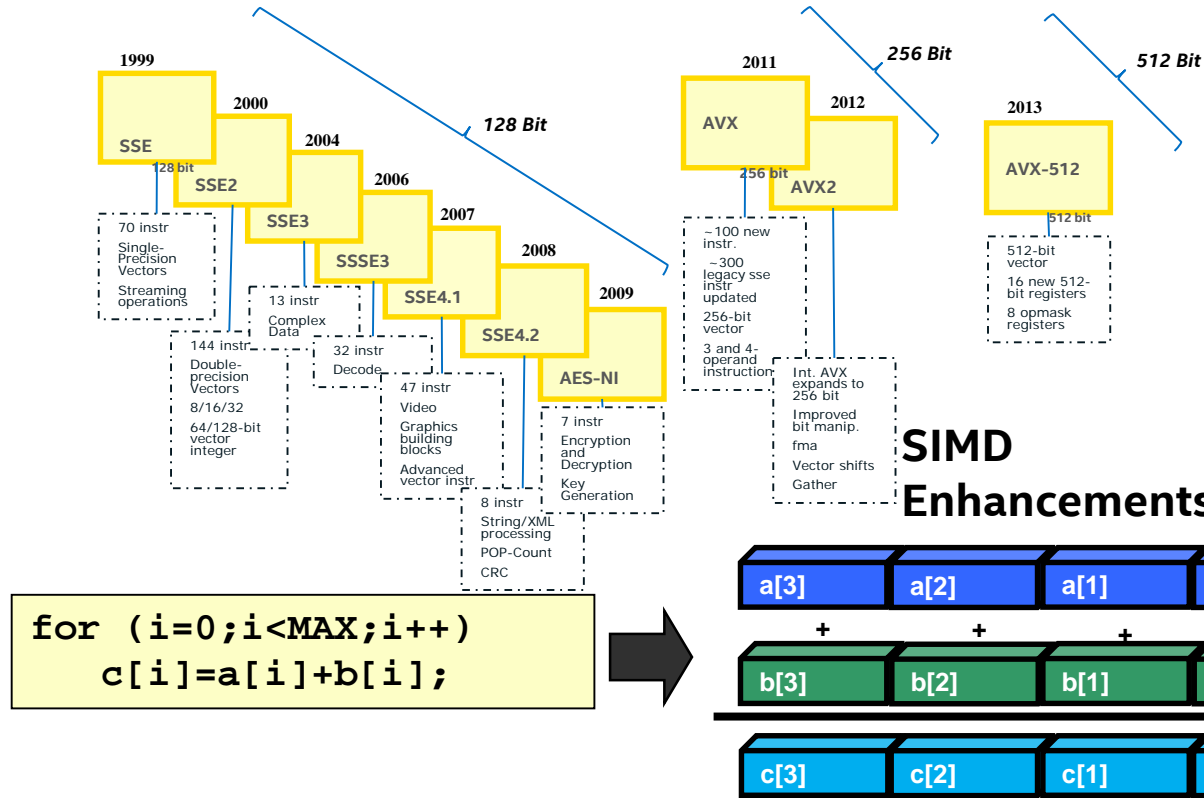
Core / Thread-Level
(Hyperthreading)

GPU-CPU

Instruction

Data (Vectorisation)





DIFFERENT WAYS OF INSERTING VECTORISED CODE

Performance Libraries (e.g. IPP and MKL)

Compiler: Fully automatic vectorization

Cilk Plus Array Notation

Compiler: Auto vectorization hints
(`#pragma ivdep, ...`)

User Mandated Vectorization
(SIMD Directive)

Manual CPU Dispatch
(`__declspec(cpu_dispatch ...)`)

SIMD intrinsic class (`F32vec4 add`)

Vector intrinsic (`mm_add_ps()`)

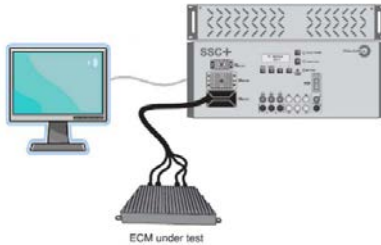
Assembler code (`addps`)

Ease of use

Programmer control

AN EXAMPLE

#1 Speedup by upgrading silicon



CPU	No Auto-Vectorisation	With Auto-Vectorisation	Speedup
P4	39.344	21.9	1.80
Core 2	5.546	0.515	10.77
Speedup	7.09	45.52	76

#2 Speedup by swapping compiler and enabling vectorisation

#3 Verified using VTune

CPU EVENT	Without Vect	With Vect
CPU_CLK_UNHALTED.CORE	16,641,000,448	1,548,000,000
INST_RETIRED.ANY	3,308,999,936	1,395,000,064
X87_OPS_RETIRED.ANY	250,000,000	0
SIMD_INST_RETIRED	0	763,000,000

SEVEN LEVELS OF HARDWARE-SUPPORTED PARALLELISM

Levels of Parallelism
Node
Socket
Core / Thread-Level (Hyperthreading)
GPU-CPU
Instruction (by CPU internals)
Data (Vectorisation)

The background is a dark blue architectural blueprint with white lines and text. It shows a complex floor plan with various rooms, corridors, and structural elements. Dimensions like '2000', '6.35', '12.6', and '3.2' are visible. There are also circular callouts with numbers like '5', '7', '8', '10', '12', '13', '17', and '19'. The text 'PERFORMANCE IT'S ALL ABOUT MEMORY!' is overlaid in large, bold, white letters, and 'Another Viewpoint' is in a smaller, yellow font below it.

PERFORMANCE IT'S ALL ABOUT MEMORY!

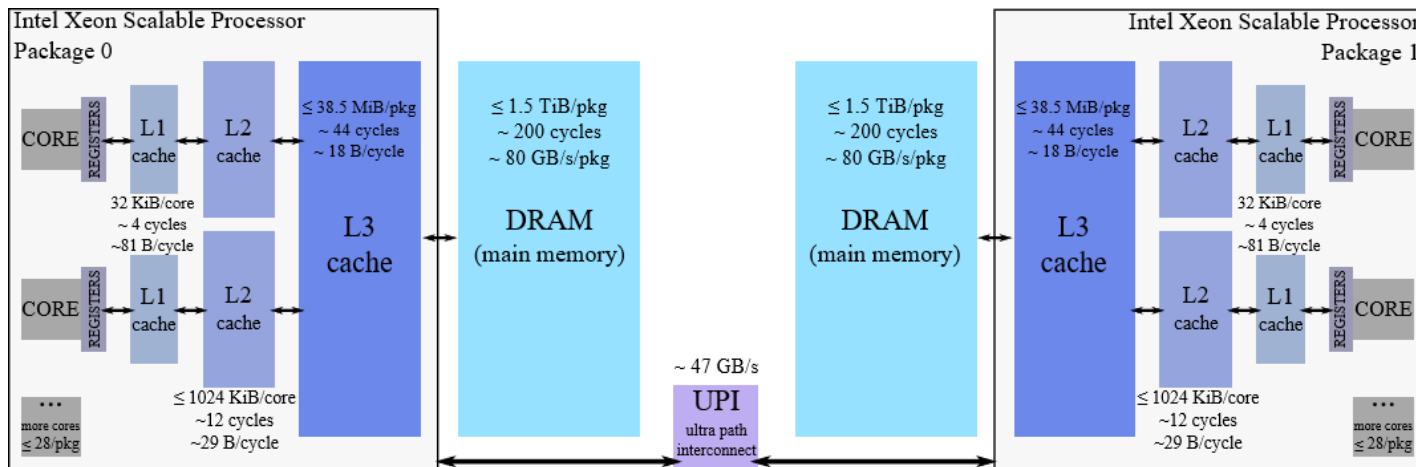
Another Viewpoint

WHICH LAPTOP SHOULD I BUY FOR BEST PERFORMANCE?



“Spend most of you money on extra memory”

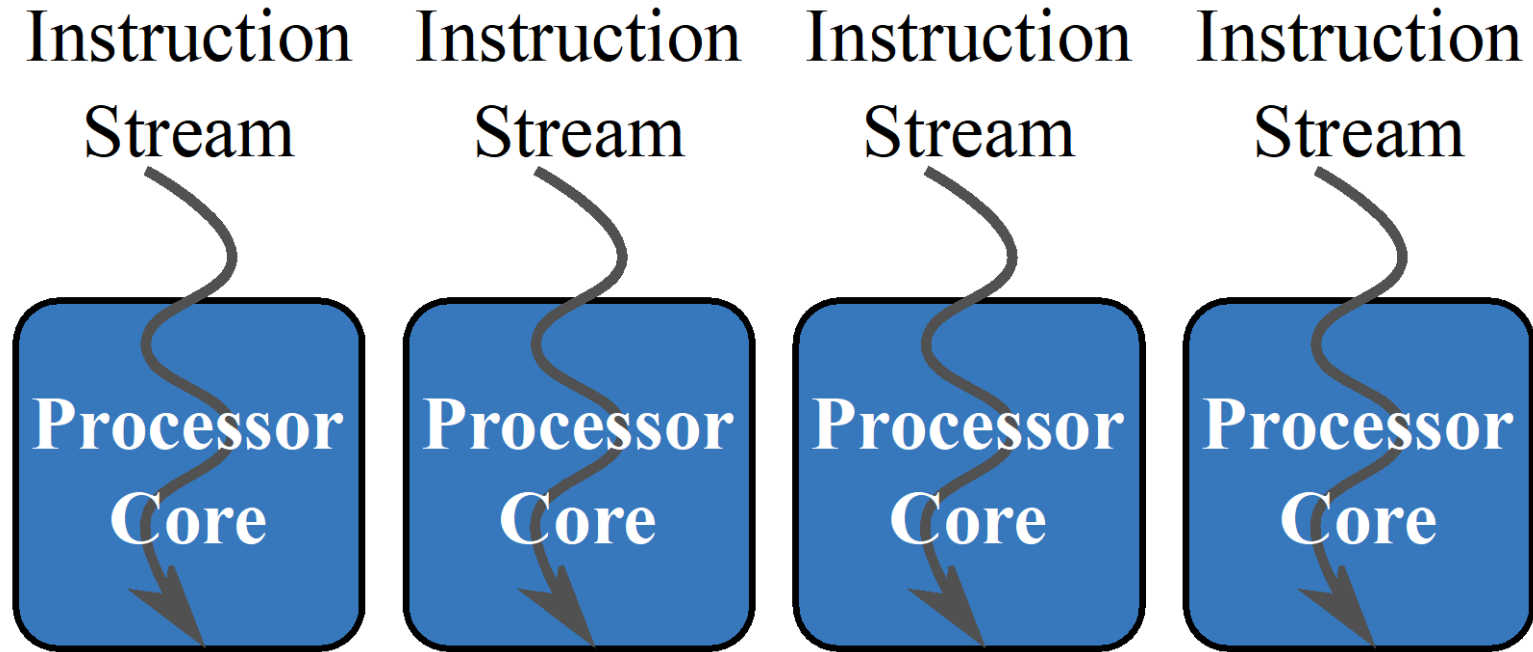
MEMORY HIERARCHY



SEVEN LEVELS OF MEMORY HIERARCHY

Levels of Hierarchy	Size	Latency	Bandwidth
Registers	2048 bytes		
L1 Cache	32 KB	4 cycles	81 bytes/cycle
L2 Cache	256 – 1024 KB	14 cycles	29 bytes/cycle
L3 Cache	8 – 38MB	60 cycles	18 bytes/cycle
(On-chip Memory)	Maybe in future versions ?		
Local DRAM	1.5 TB	200 cycles	80 GB/s
Remote DRAM			47 GB/s

MEMORY CONTROLLERS AND THREADS

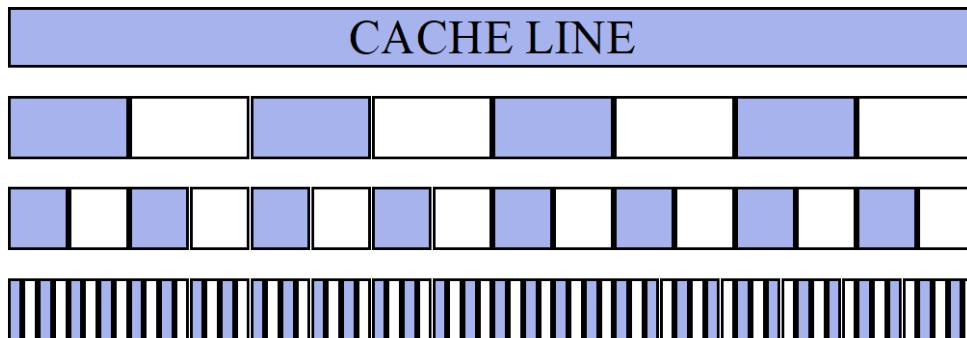


CACHE LINES AND VECTORIZATION

8 double precision values

16 single precision values

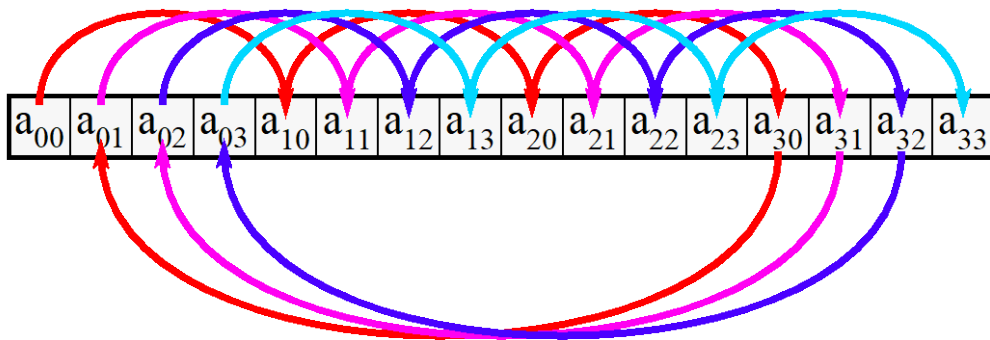
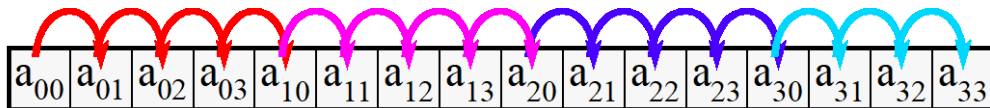
64 bytes



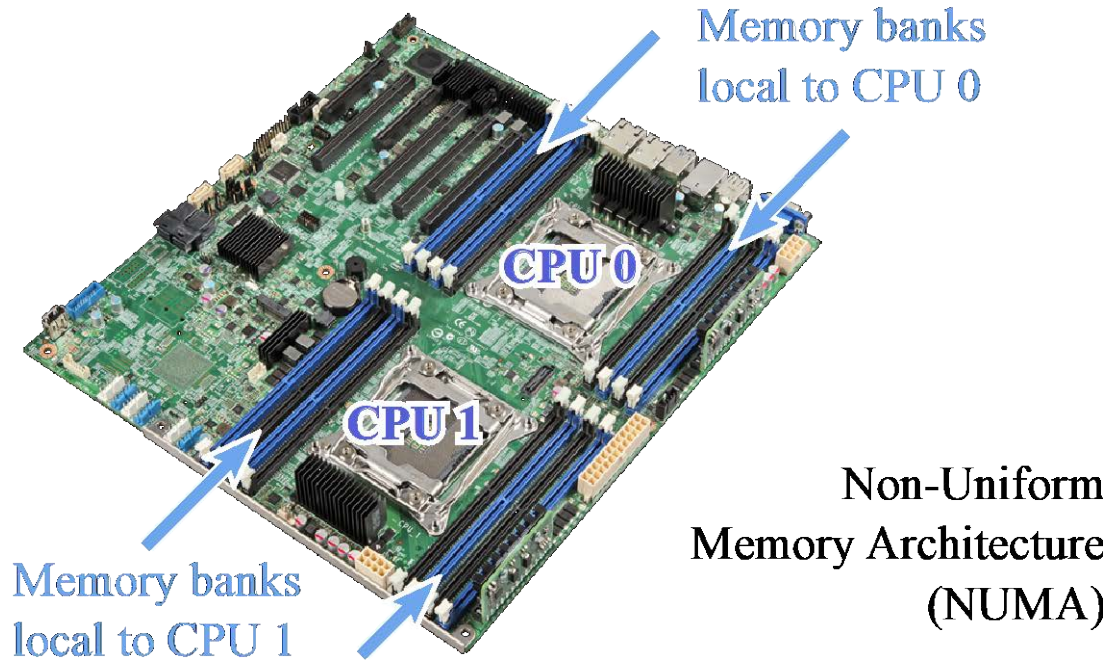
SEQUENTIAL ACCESS

a_{00}	a_{01}	a_{02}	a_{03}
a_{10}	a_{11}	a_{12}	a_{13}
a_{20}	a_{21}	a_{22}	a_{23}
a_{30}	a_{31}	a_{32}	a_{33}

a_{00}	a_{01}	a_{02}	a_{03}
a_{10}	a_{11}	a_{12}	a_{13}
a_{20}	a_{21}	a_{22}	a_{23}
a_{30}	a_{31}	a_{32}	a_{33}

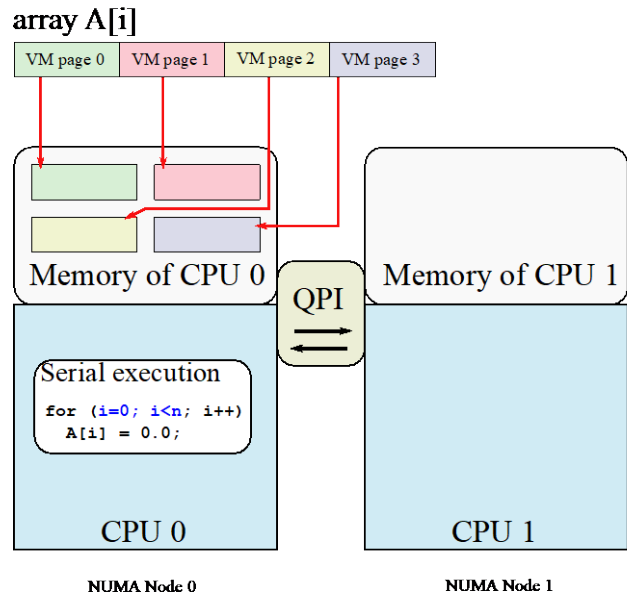


NUMA AND DATA ACCESS LOCALITY

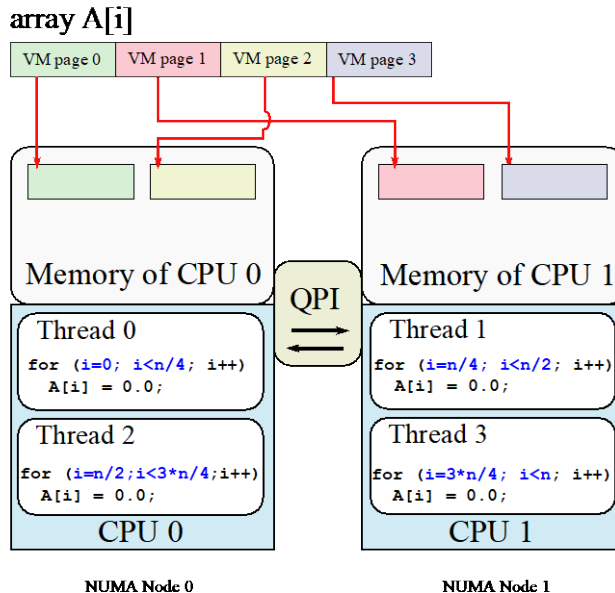


FIRST TOUCH POLICY AND ARRAY INITIALIZATION

Poor First-Touch Allocation



Good First-Touch Allocation



OPTIMIZATION IS ...

Moving code from being
Memory Bound
to being
Compute Bound

