# Technologies and application performance

**Marc Mendez-Bermond**

**HPC Solutions Expert** - Dell Technologies

September 2017

**D✕LL**EMC

# The landscape is changing



"*We are no longer in the general purpose era… the argument of tuning software for hardware is moot. Now, to get the best bang for the buck, you have to tune both.*"
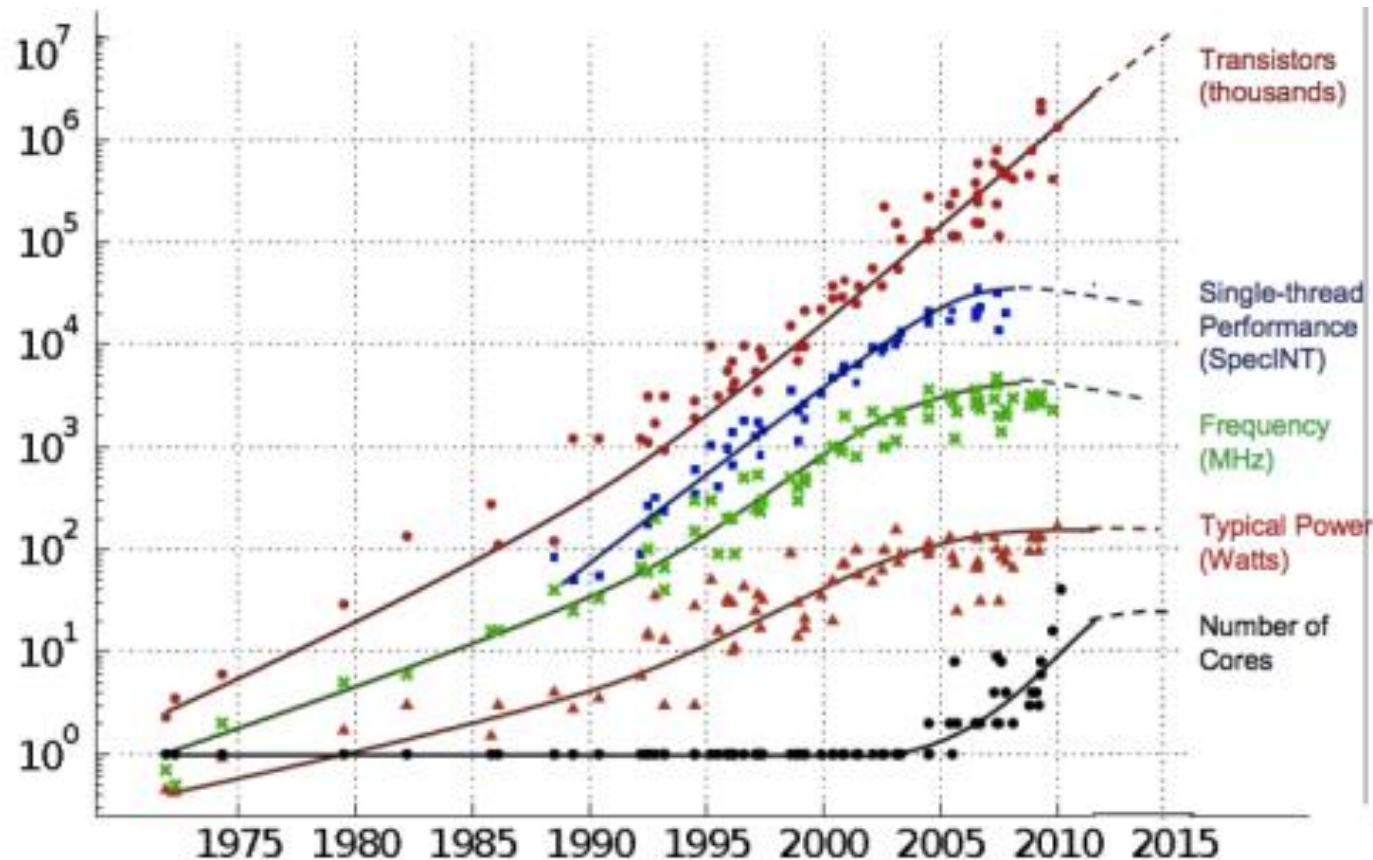
- Kushagra Vaid, general manager of server engineering, Microsoft Cloud Solutions

https://www.nextplatform.com/2017/03/08/arm-amd-x86-server-chips-get-mainstream-lift-microsoft/amp/
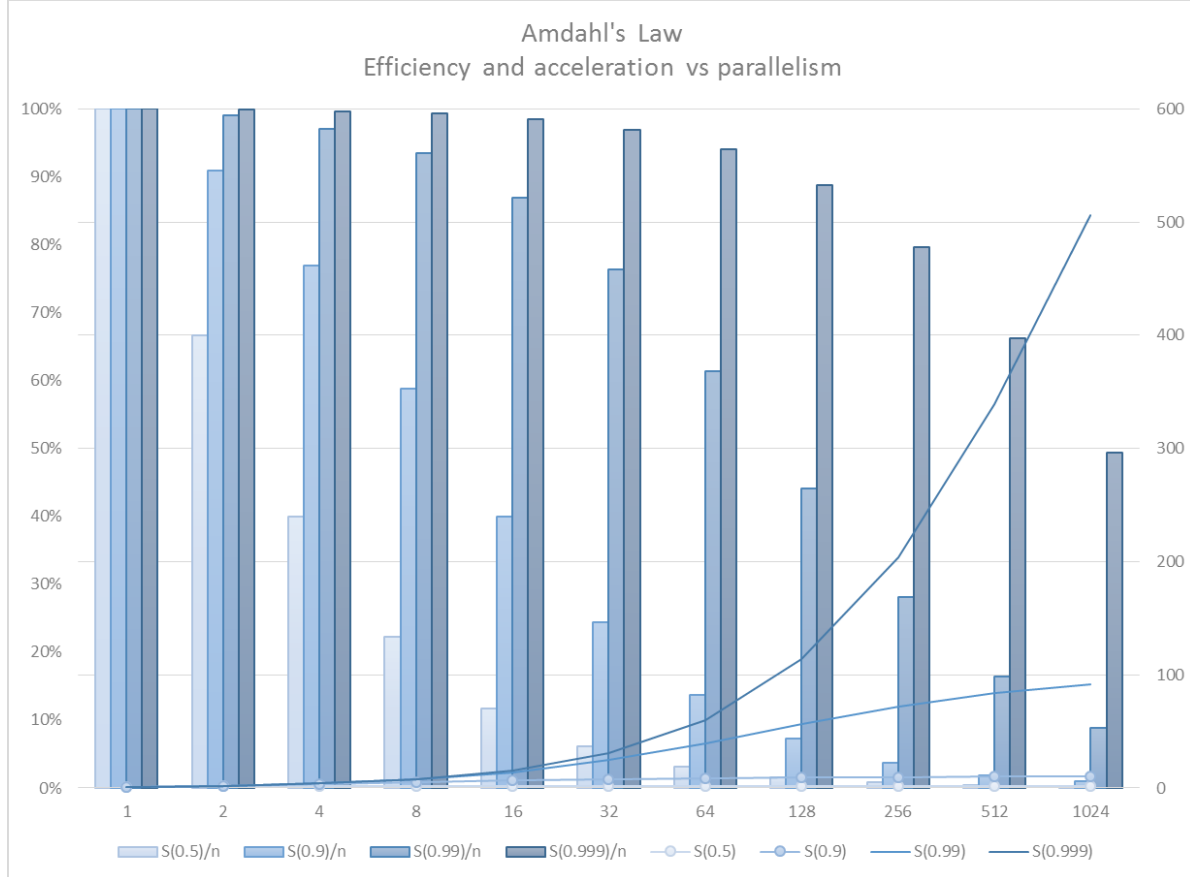
DELL EMC

# Moore's Law (Technology)

- The clock speed plateau

- The power ceiling

- IPC limit



*Chuck Moore, "DATA PROCESSING IN EXASCALE-CLASS COMPUTER SYSTEMS", The Salishan Conference on High Speed Computing, 2011*

**D∕ELL**EMC

# Amdahl's Law (Application)

- Amdahl's law predicts performance from your app parallelization

- 50% : x2 max

- 99% : x100 max

- 99.9% : x1000 max

- But you should also check the efficiency here :
  - 99.9% parallel, at 1024 processors, x509 and efficiency at 49% …



Amdahl's Law
Efficiency and acceleration vs parallelism
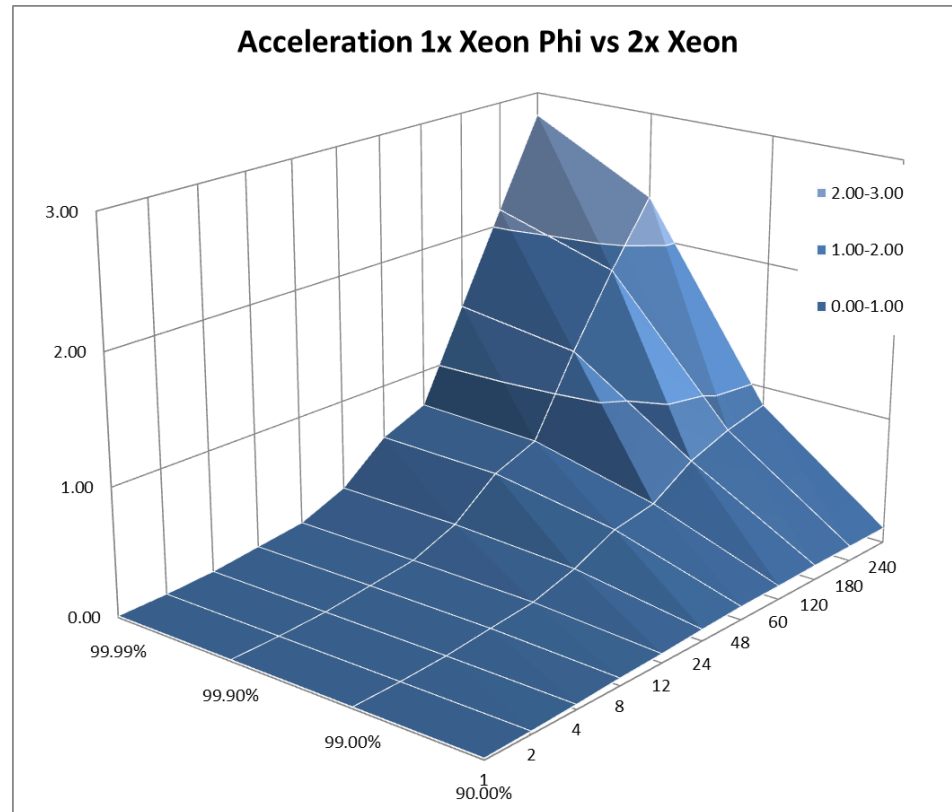
DELL EMC

**Computing technologies**

# Intel Xeon Phi : a few considerations

- x86_64 programming models
- Cache coherency
  - Dual-ring interconnect
  - 8 (soon 16) GB RAM
- Right to the point cores
  - No « out of order» execution
  - No branch prediction
  - 4 Hyper-threads per core
  - Wide vectors (16 op/c/core)
- PCIe connectivity to host

**App should fit in onboard memory,
Parallelism > 99.9%,
Vectorization > 95%**

Xeon Phi : 60 cores @ 1 GHz vs 2 Xeon : 8 cores @ 2.6 GHz

**Acceleration 1x Xeon Phi vs 2x Xeon**



Legend:
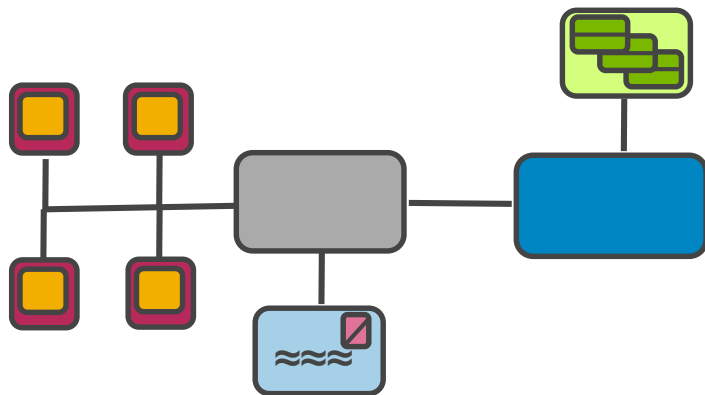- 2.00-3.00
- 1.00-2.00
- 0.00-1.00

DELL EMC

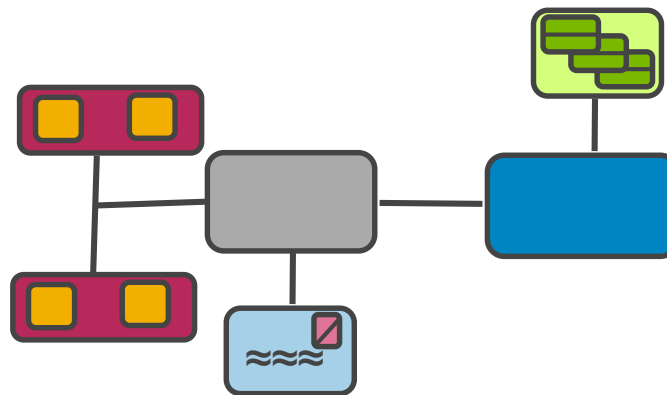# Moore's Law vs Amdahl's Law - "too Many Cooks in the Kitchen"



Industry is applying Moore's Law by adding more cores

Meanwhile Amdahl's Law says that you cannot use them all efficiently

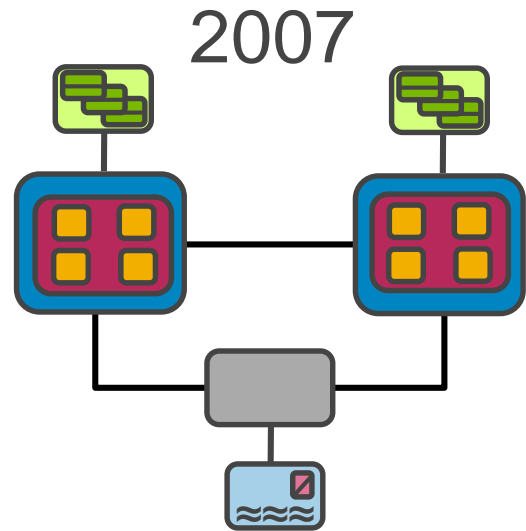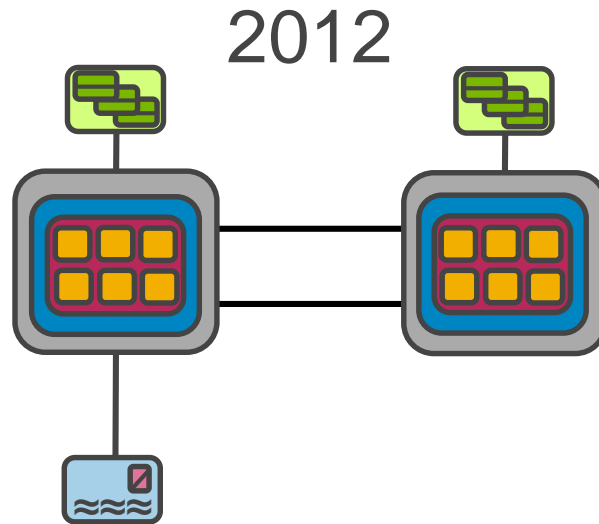DELL EMC

# System trend over the years (1)



~1970 - 2000

2005

**Multi-core**: TOCK

DELLEMC

# System trend over the years (2)

2007

2012

**Integrated Memory controller**:

TOCK

**Integrated PCIe controller**:

TOCK

DELL EMC

# Future



**Integrated Network Fabric Adapter**:

TOCK

**SoC designs:**

TOCK

**D∅LL**EMC

# Improving performance - what levels do we have?

- Challenge: Sustain performance trajectory without massive increases in cost, power, real estate, and unreliability

- Solutions: No single answer, must **intelligently turn** "Architectural Knobs"

① ② ③ ④ ⑤

$$(Freq) \times \left(\frac{cores}{socket}\right) \times (\#sockets) \times \left(\frac{inst \ or \ ops}{core \ \times \ clock}\right) \times (Efficiency)$$

Hardware performance

What you really get

Software performance

DELL EMC

# Turning the knobs 1 - 4

**1** Frequency is unlikely to change much - Thermal/Power/Leakage challenges

**2** Moore's Law still holds: 130 -> 14 nm - LOTS of transistors

**3** Number of sockets per system is the easiest knob.
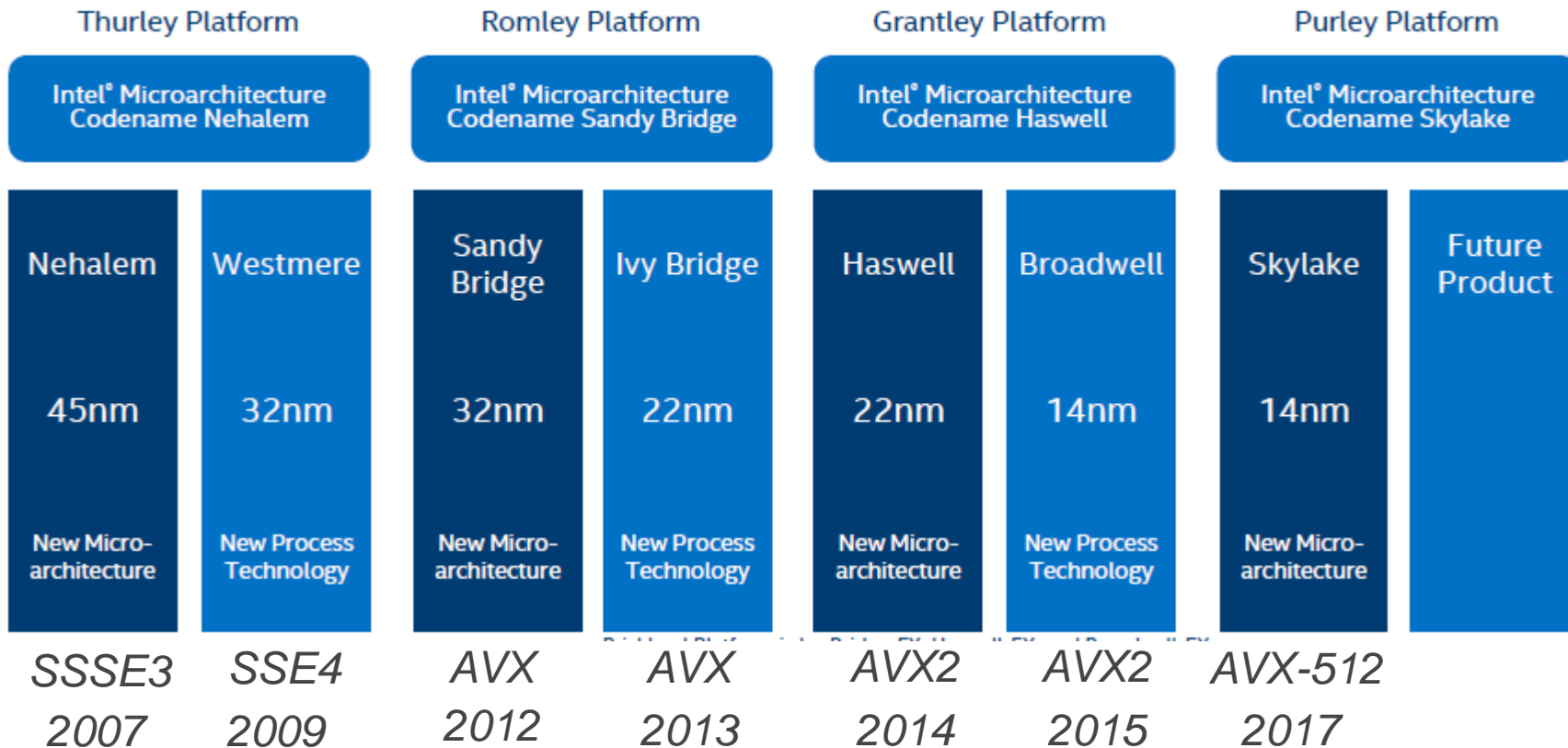
Challenging for power/density/cooling/networking

IPC still grows

**4** FMA3/4, AVX, FPGA implementations for algorithms

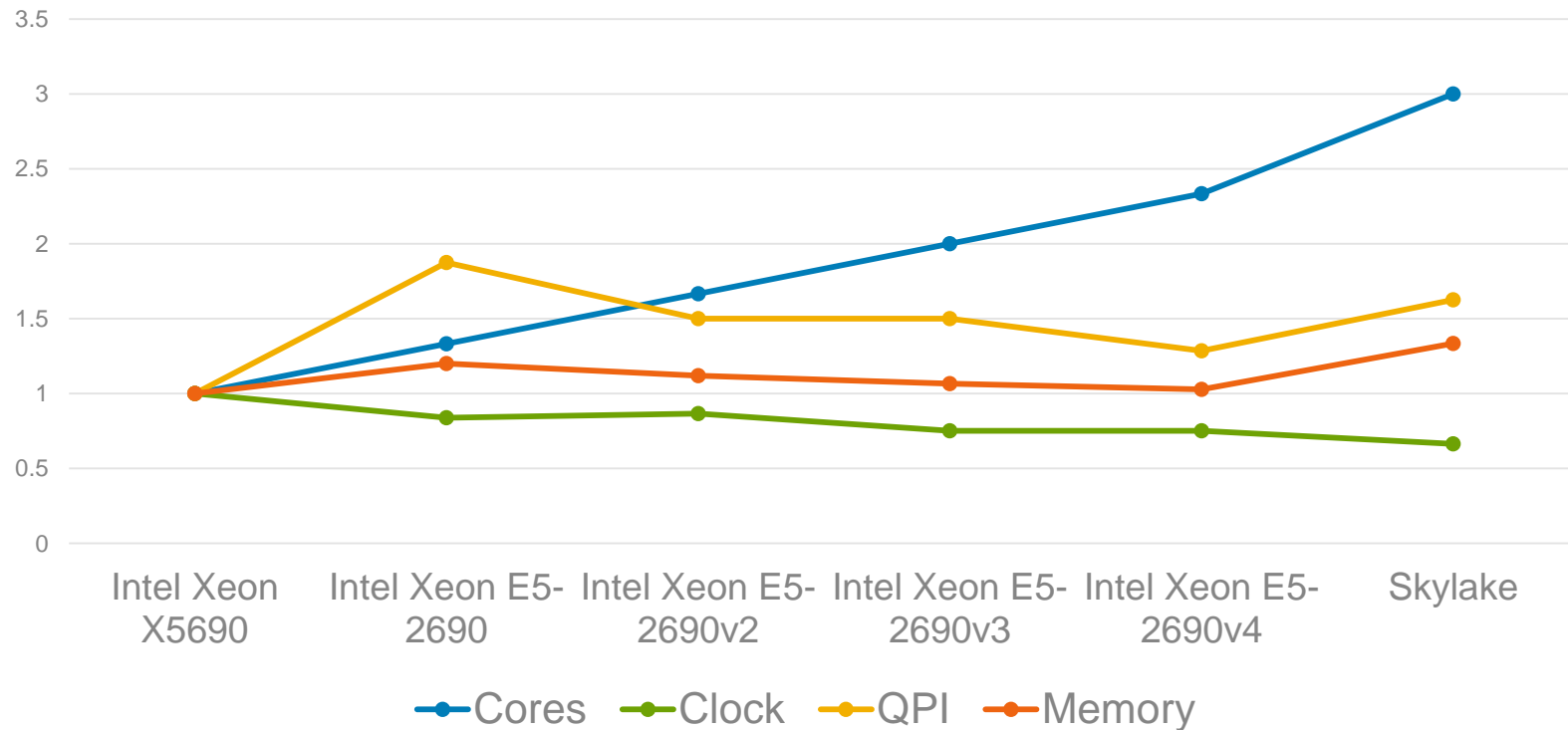Challenging for the user/developer
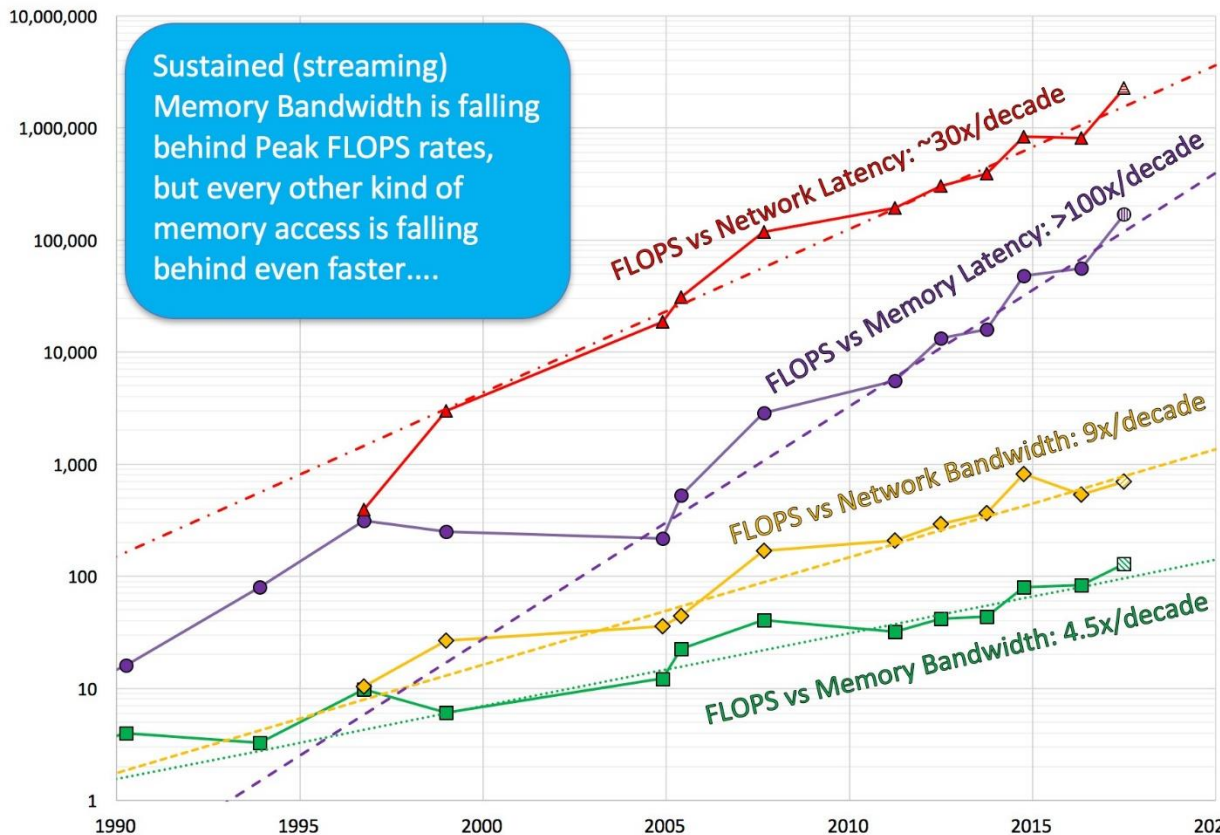
**D≪LL**EMC

# New capabilities according to Intel

| Thurley Platform | | Romley Platform | | Grantley Platform | | Purley Platform | |
|---|---|---|---|---|---|---|---|
| Intel® Microarchitecture Codename Nehalem | | Intel® Microarchitecture Codename Sandy Bridge | | Intel® Microarchitecture Codename Haswell | | Intel® Microarchitecture Codename Skylake | |
| Nehalem | Westmere | Sandy Bridge | Ivy Bridge | Haswell | Broadwell | Skylake | Future Product |
| 45nm | 32nm | 32nm | 22nm | 22nm | 14nm | 14nm | |
| New Micro-architecture | New Process Technology | New Micro-architecture | New Process Technology | New Micro-architecture | New Process Technology | New Micro-architecture | |
| *SSSE3* | *SSE4* | *AVX* | *AVX* | *AVX2* | *AVX2* | *AVX-512* | |
| *2007* | *2009* | *2012* | *2013* | *2014* | *2015* | *2017* | |

DELL EMC

# The state of ISV software

| Segment | Applications | Vectorization support |
|---|---|---|
| CFD | Fluent, LS-DYNA, STAR CCM+ | Limited SSE2 support |
| CSM | CFX, RADIOSS, Abaqus | Limited SSE2 support |
| Weather | WRF, UM, NEMO, CAM | Yes |
| Oil and Gas | Seismic processing | Not applicable |
| | Reservoir Simulation | Yes |
| Chemistry | Gaussian, GAMESS, Molpro | Not applicable |
| Molecular dynamics | NAMD, GROMACS, Amber,… | PME kernels support SSE2 |
| Biology | BLAST, Smith-Waterman | Not applicable |
| Molecular mechanics | CPMD, VASP, CP2k, CASTEP | Yes |

**Bottom line:** ISV support for new instructions is poor. Less of an issue for in-house developed codes, but programming is hard

DELLEMC

# Meanwhile the bandwidth is suffering

# Add to this the Memory Bandwidth and System Balance



Sustained (streaming) Memory Bandwidth is falling behind Peak FLOPS rates, but every other kind of memory access is falling behind even faster....

FLOPS vs Network Latency: ~30x/decade

FLOPS vs Memory Latency: >100x/decade

FLOPS vs Network Bandwidth: 9x/decade

FLOPS vs Memory Bandwidth: 4.5x/decade

# And data is becoming sparser (think "Big Data")

$$A \times x = y$$

Sparse Matrix "A"
- Most entries are zero
- Hard to exploit SIMD
- Hard to use caches

- This has very low arithmetic density and hence memory bound
- Common in CFD, but also in genetic evaluation of species

**DELL**EMC

# Xeon roofline model (v4)

# What does Intel do about these trends?

| Problem | Westmere | Sandy Bridge | Ivy Bridge | Haswell | Broadwell | Skylake |
|---|---|---|---|---|---|---|
| QPI bandwidth | No problem | Even better | Two snoop modes | Three snoop modes | Four (!) snoop modes | • UPI<br>• COD snoop modes |
| Memory bandwidth | No problem | Extra memory channel | Larger cache | Extra load/store units | Larger cache | • Extra load/store units<br>• +50% memory channels |
| Core frequency | No problem | • More cores<br>• AVX<br>• Better Turbo | • Even more cores<br>• Above TDP Turbo | • Still more cores<br>• AVX2<br>• Per-core Turbo | • Again even more cores<br>• optimized FMA<br>• Per-core Turbo based on instruction type | • More cores<br>• Larger OOO engine<br>• AVX-512<br>• 3 different core frequency modes |

# C4130 – Ten supported variations

DELLEMC

# Pragmatic computing

| Parallelize | Vectorize |
|---|---|
| Take advantage of multicore | Take advantage of large-vector units |

Amdahl's law : limiting factor

Moore's law : benefiting factor

| Optimize |
|---|
| • Intrinsic optimization<br>• Execution optimization |

Efficiency of implementation

**DELL**EMC

# Public benchmark data

# Portfolio: Ready Solutions for HPC

**BUILD**

**BUY**

**Benefits**

| Maximum flexibility<br>Validated for use case<br>Heterogeneity with lower risk<br>Component lifecycle automation and control | Consumption models | Fastest time to value<br>Optimized and tuned for use case<br>Greatest risk reduction<br>Solution lifecycle automation |
|---|---|---|

**Solutions**

**Scale**

## STORAGE READY BUNDLES

### Dell EMC Ready Bundle for HPC NFS Storage

Scales from a minimum of 48TB to 480TB of raw capacity in a single name space

### Dell EMC Ready Bundle for HPC Lustre Storage

Lustre parallel file storage system scales from 120TB to petabytes of data

## SYSTEMS FOR A RANGE OF USE CASES

### Dell EMC HPC System for Life Sciences

Fully integrated for pharma/biotech applications

### Dell EMC HPC System for Manufacturing

Fully integrated for compute-aided engineering (CAE) workloads

### Dell EMC HPC System for Research

General purpose compute cluster for multiple research workloads

DELL EMC

# HPC Innovation Lab World-Class Infrastructure

**Dedication to Research and Development:**

- 13K sq. ft (1200m²) with 1300+ Servers and ~10PB

- Leverage Expertise in HPC

- Test New Technologies

- Tune your applications for performance and efficiency



**Zenith**
- Top500 class system based on Intel Scalable Systems Framework (OPA, KNL, Xeon, OpenHPC)
- 256-nodes with dual 2697v4 processors, non-blocking OPA fabric and 270TFlops sustained performance

**Rattler**
- Research/development system in collaboration with Mellanox and NVIDIA
- 80 nodes configured with Infiniband EDR and 2660v3 processors

**D∂LL**EMC

# Merci !

marc_mendez_bermond@dell.com