# **IDF**2009

## INTEL DEVELOPER FORUM

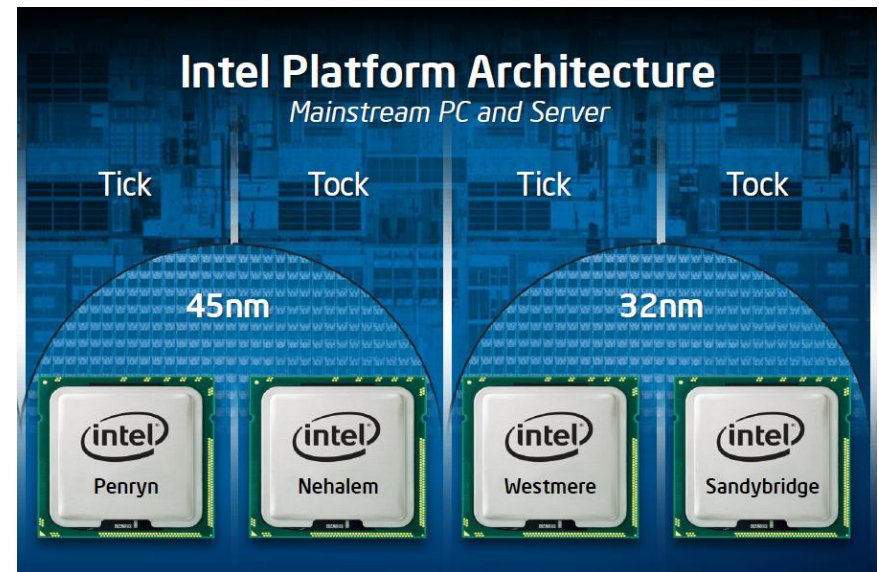**PROCESSOR ARCHITECTURE AND HPC POINT OF VIEW**

# Zoltan Juhasz

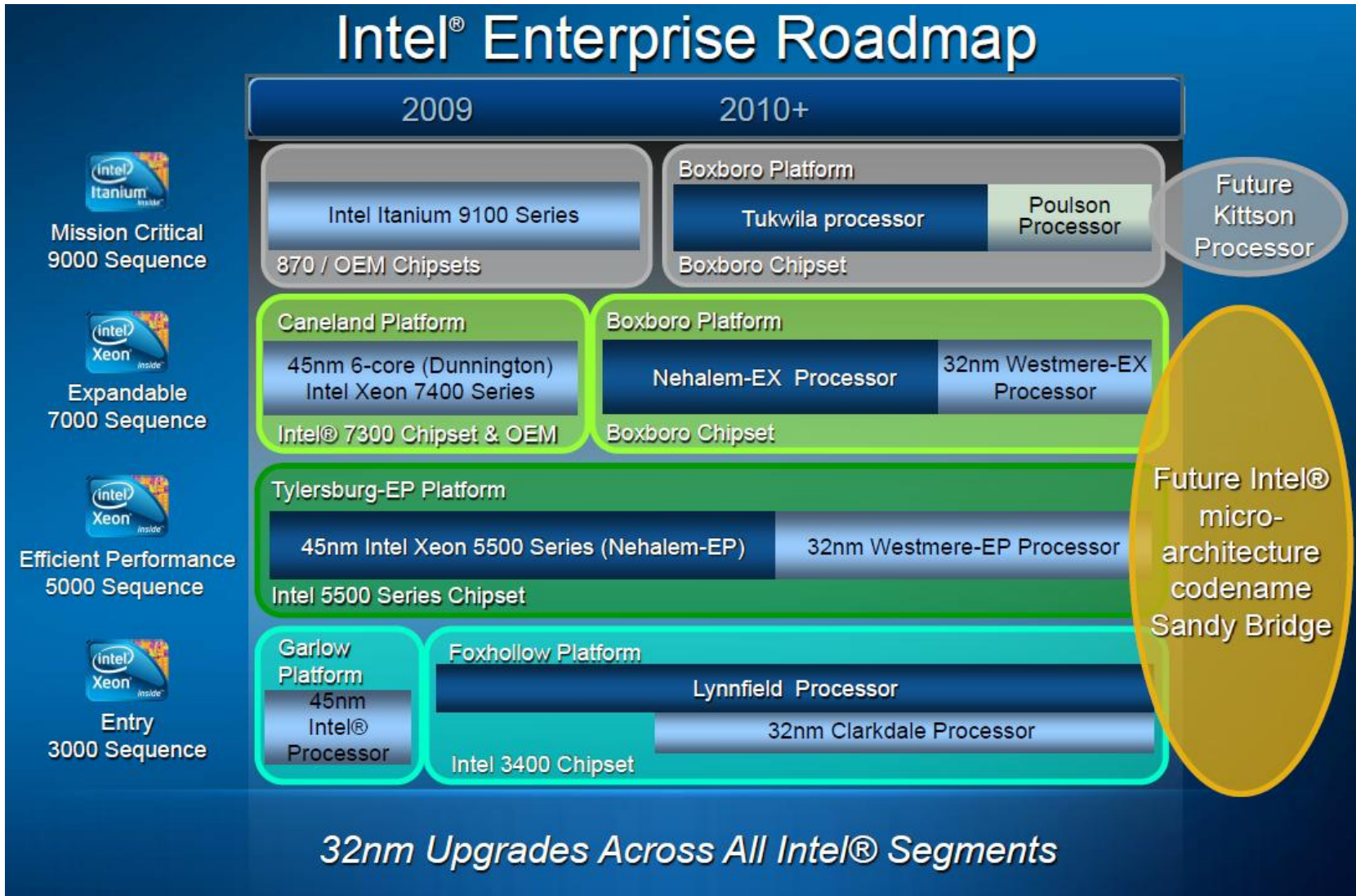zoltan.juhasz@morganstanley.com

# Agenda

- Few quotes from the **Keynotes at IDF**

- Intel Server / Desktop **Roadmap**

- Key Points of **Intel Nehalem** Based Platforms

- Intel **Advanced Vector eXtension (AVX)** introduction

- Intel **Nehalem-EX** architecture: features and benefits

- **HP Extreme Scale Computing**: Server Futures (near term)

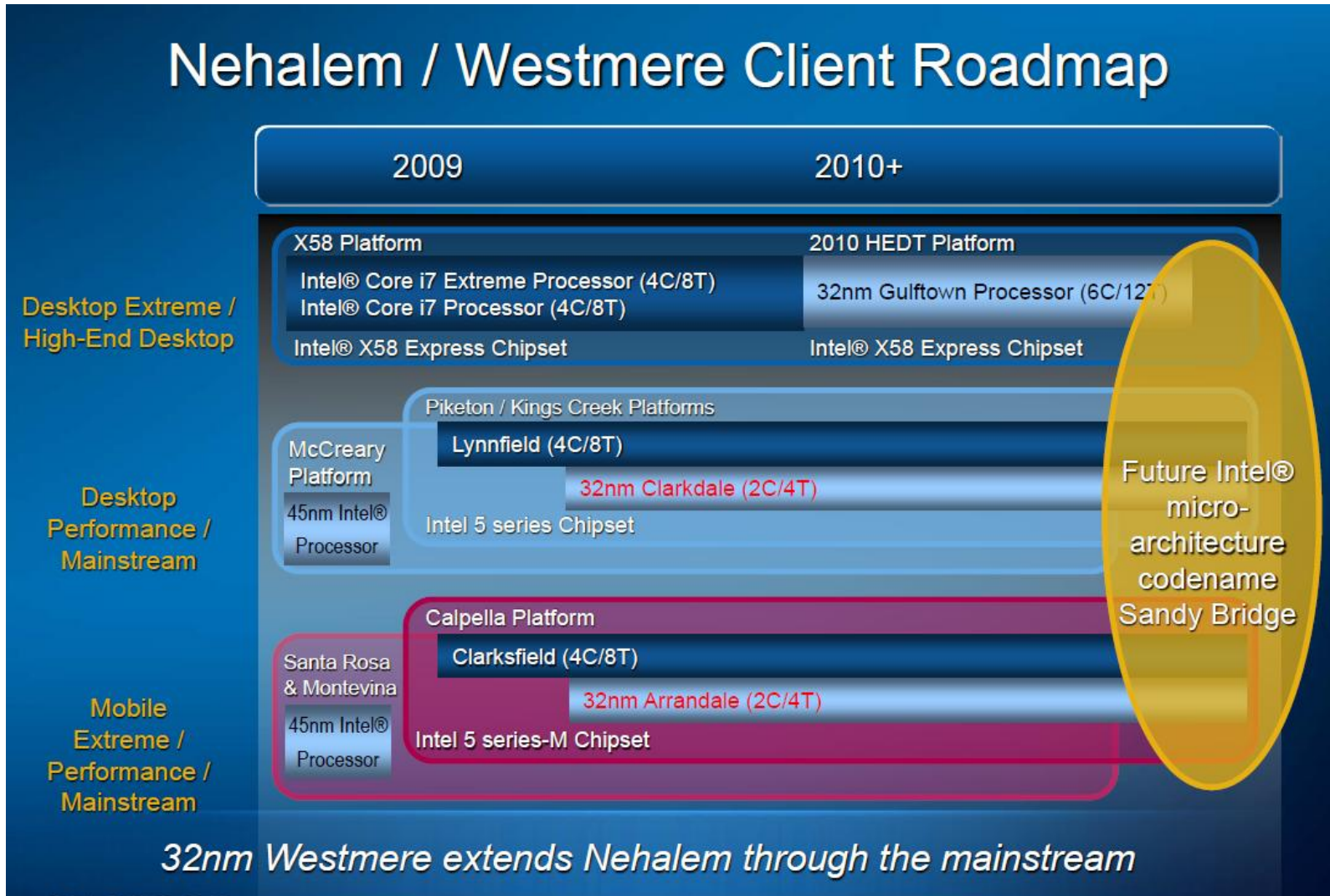Morgan Stanley

# Keynotes from IDF

- Intel Platform Architecture: **Tick-Tock / next gen. silicon architecture - new microarchitecture**

- Intel is aggressively expanding to the **embedded world**

- Intel Platform Choice for Mission Critical Servers: **Nehalem-EX**

- **Improved energy efficiency**, and dissipation is another key area for Intel

- **New silicon technology**: 32nm (Q4 2009)



Morgan Stanley

# Keynotes: Server Roadmap

# Keynotes: Client Roadmap



Morgan Stanley
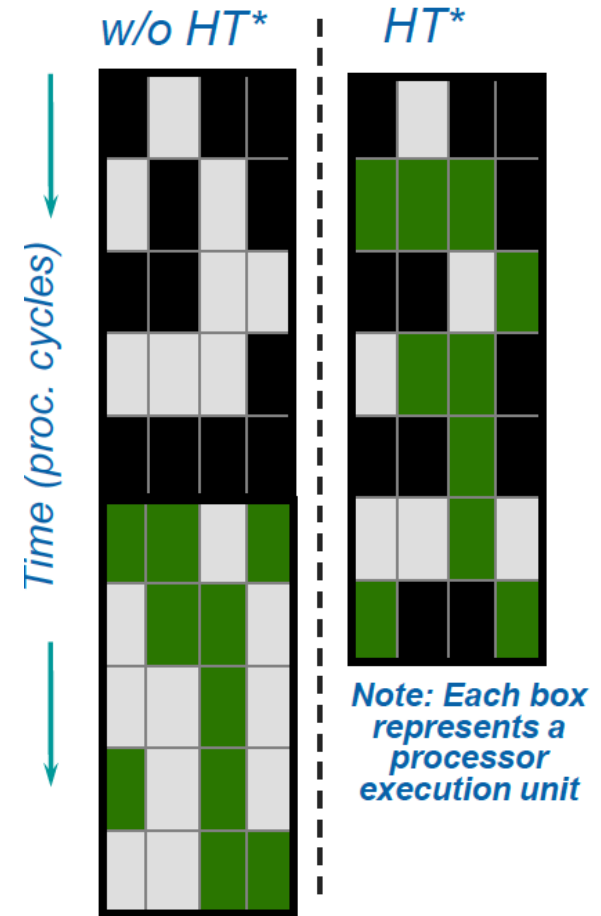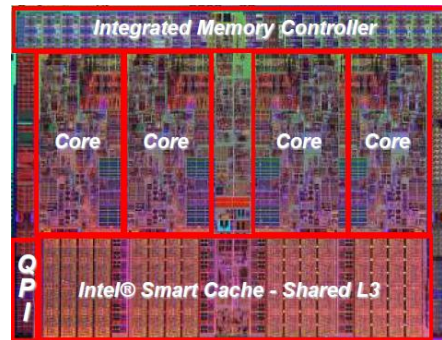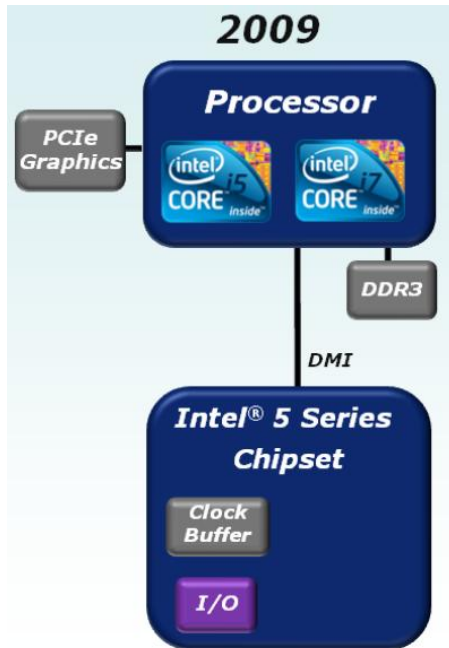
# Intel Nehalem Based Platforms

- Intel **Hyper-Threading** Technology

- Intel **Turbo Boost** Technology

- **PCI Express interface** moves to the processor
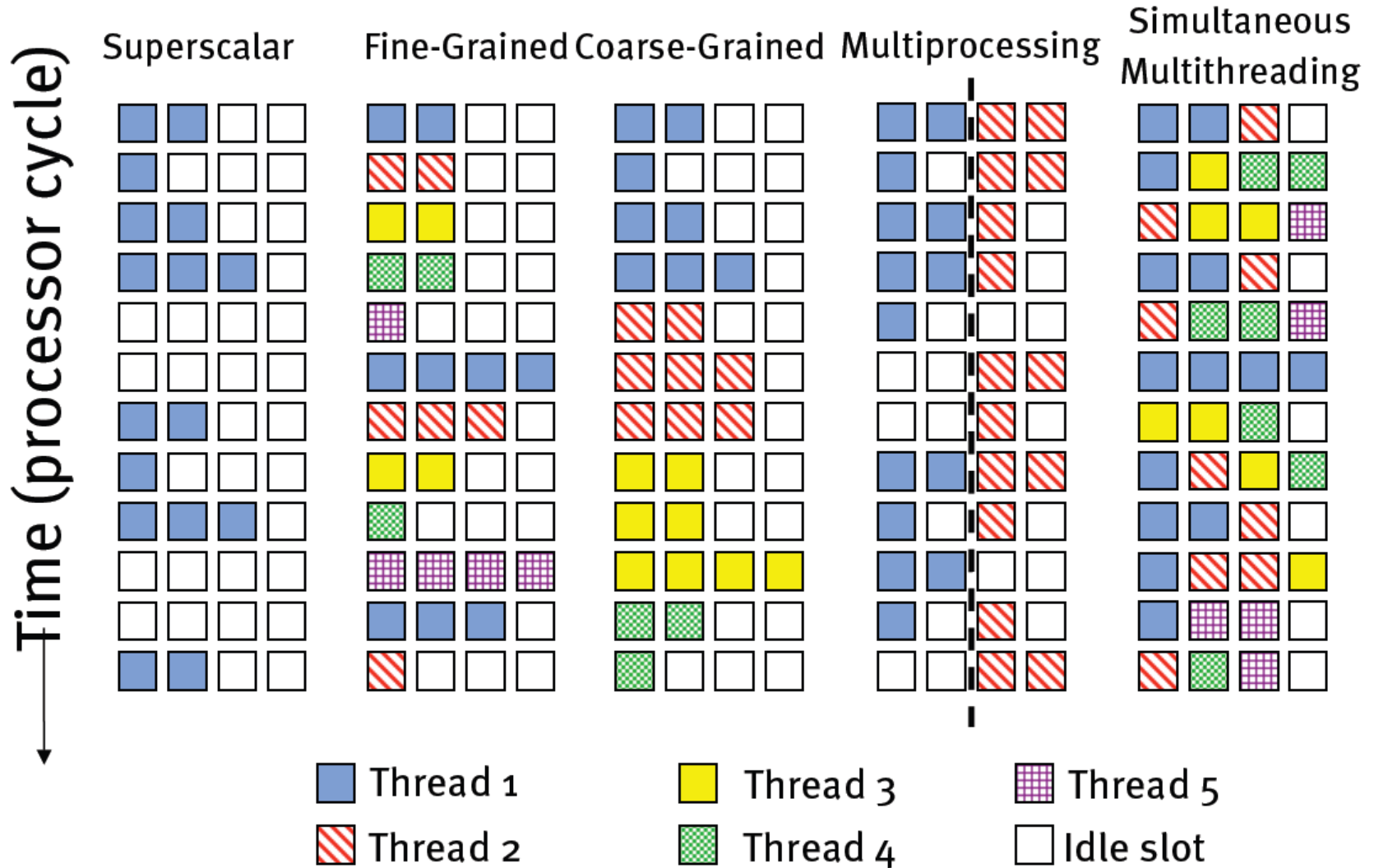


Morgan Stanley
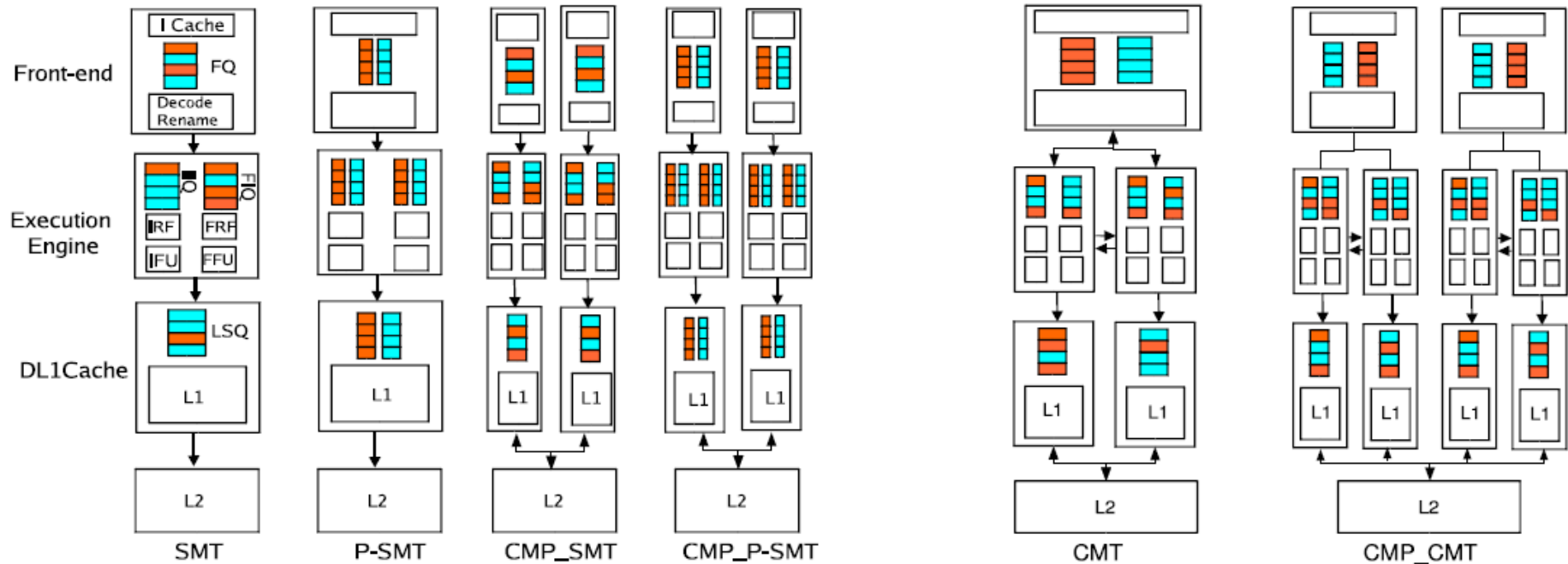
# The third law

# *"Any sufficiently advanced technology is indistinguishable from magic."*

Arthur C. Clarke, "Profiles of The Future", 1961 (Clarke's third law)

Morgan Stanley

# Multithreaded Categories

Morgan Stanley

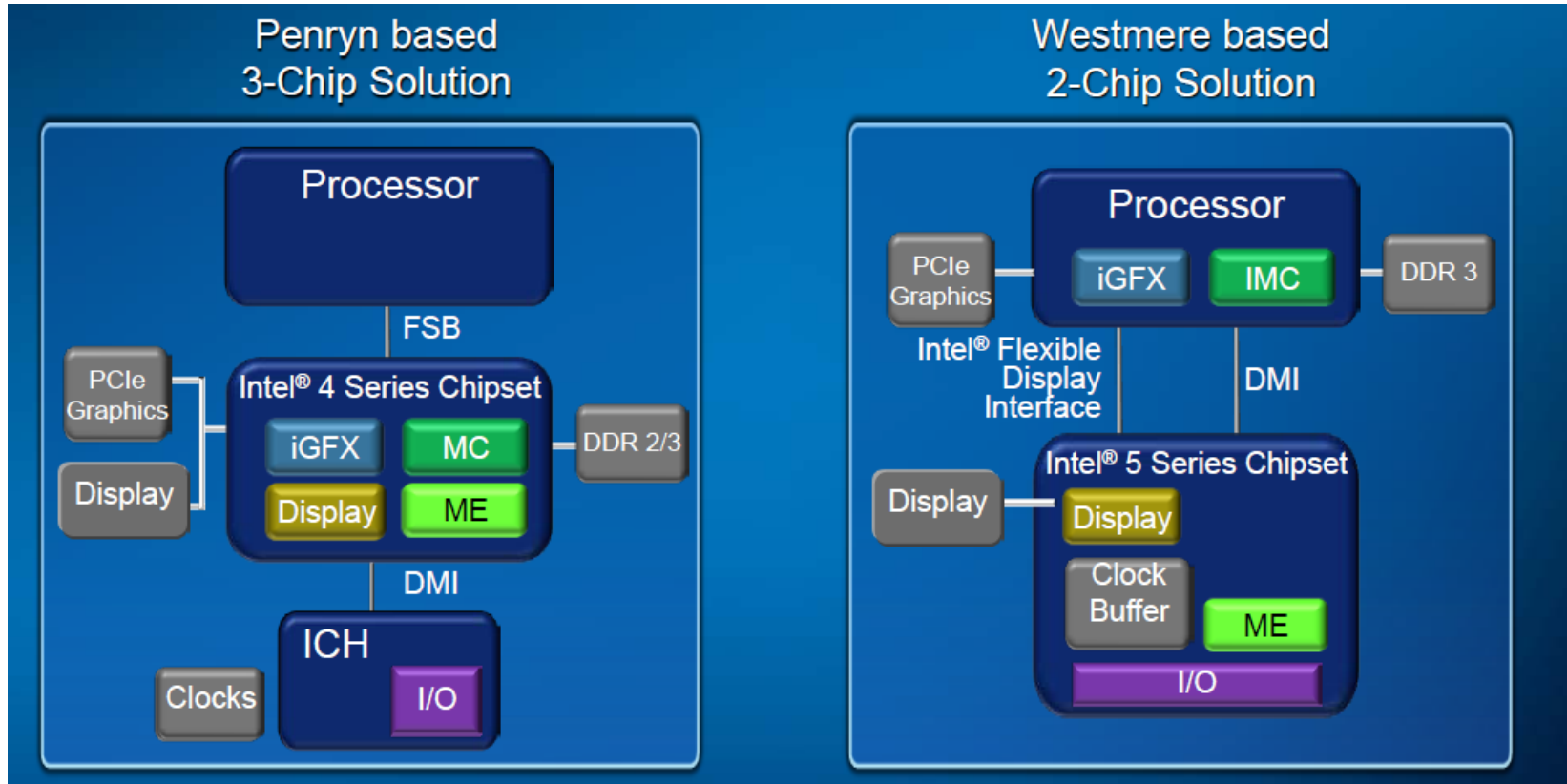# SMT – partitioned / clustered



Difference lies is in the design decision what is
- **replicated** (register renaming logic, instruction pointer, ITLB, etc.)
- **partitioned** (reorder buffers, load/store buffers, various queues)
- **shared** (caches, micro architectural registers, execution units)

Will not help if one thread **monopolizes shared resources**, **cache effects** etc.

Morgan Stanley

# AMD Bulldozer

- Bulldozer implements a Cluster(-based) multiprocessing  (CMT)

- Integrates two superscalar processors

- 2 Flexible Floating Point FAMC unit that can be dedicated or shared between the two core per cycle

- Independent  integer and FP schedulers are provided



## Morgan Stanley

# Keynotes: Intel Flexible Display Interface

Morgan Stanley

# Intel Advanced Vector Extensions (AVX)

- Intel **Sandy Bridge Microarchitecture** features Intel **Advanced Vector Extension**



Morgan Stanley

12

# Intel AVX Key features and benefits
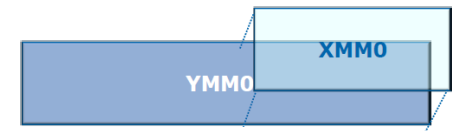
- **Wider Vectors** (128 to 256 bit with 2 128-bit load ports)

- Enhances Data Rearrangement (broadcast, mask loads and permute data)

- **Flexible unaligned memory access** support (e.g. no penalty for unalgined loads on aligned memory)

- Mixing AVX / SSE code may incur penalty (AVX 256 dirties upper 128 bits)
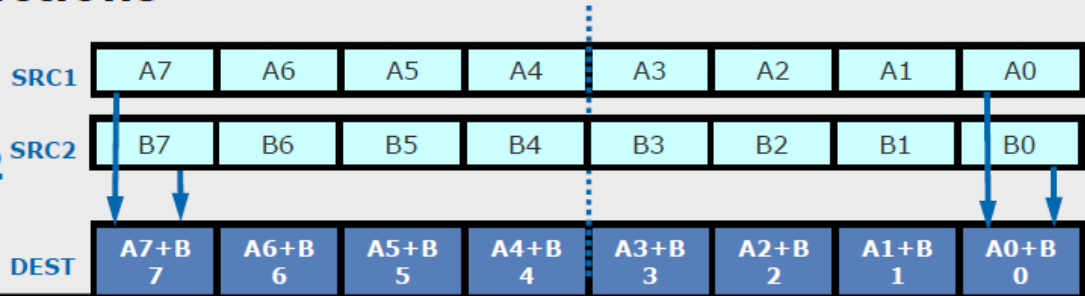
- A **new** 3- and 4- operand **instruction format**:

```
xmm10 = xmm9 + xmm1
movaps  xmm10, xmm9          ▶    vaddpd xmm10, xmm9, xm─
addpd   xmm10, xmm1

xmm10 = xmm9 + m128
movups  xmm10, m128          ▶    vaddpd xmm10, xmm9,
addpd   xmm10, xmm9
```

**1 less copy,**
**3 bytes smaller code si:**

**1 more load/op**
**fusion opportunity,**
**4+ bytes smaller**
**code size**

- New 4- operand Blends example, implicit xmm0 not longer needed

```
movaps  xmm0, xmm4
movaps  xmm1, xmm2          ▶    vblendvps xmm1, xmm2, m128, xmm4
blendvps xmm1, m128
```

Morgan Stanley

13

# Intel AVX New Primitives

- **Simple in-lane instructions**
  - 2 lanes, 128 bit each

  **vAddPS dest, src1, src2**

  | SRC1 | A7 | A6 | A5 | A4 | A3 | A2 | A1 | A0 |
  |------|----|----|----|----|----|----|----|----|
  | SRC2 | B7 | B6 | B5 | B4 | B3 | B2 | B1 | B0 |
  | DEST | A7+B7 | A6+B6 | A5+B5 | A4+B4 | A3+B3 | A2+B2 | A1+B1 | A0+B0 |

- **New in-lane PS and PD Permutes**
  - Permute controlled via immediate

  **vPermilPS dest, src, imm**

  | SRC1 | X7 | X6 | X5 | X4 | X3 | X2 | X1 | X0 |
  |------|----|----|----|----|----|----|----|----|
  | DEST | X7 .. X4 | X7 .. X4 | X7 .. X4 | X7 .. X4 | X3 .. X0 | X3 .. X0 | X3 .. X0 | X3 .. X0 |

- **New 128-bit permutes**
  - Useful for lane-crossing operations

  **vPerm2F128 dest, src1, src2, imm**

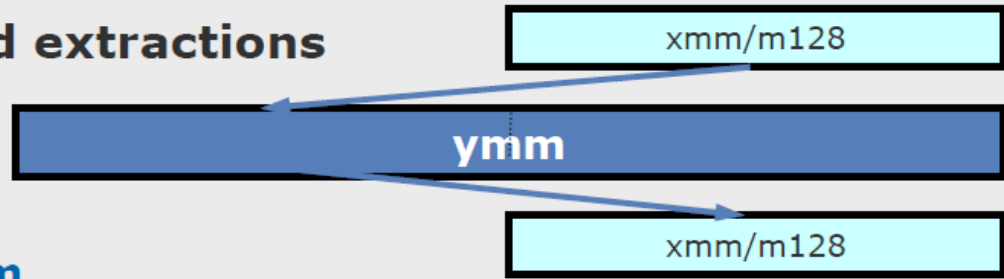  | SRC2 | Y1 | Y0 |
  |------|----|----|
  | SRC1 | X1 | X0 |
  | DEST | X0,X1,Y0, or Y1 | X0,X1,Y0, or Y1 |

Morgan Stanley

# Intel AVX New Primitives

- **128-bit Insertions and extractions**
  - Useful for lane crossing operations

  **vInsertF128 dest, src, imm**
  **vExtractF128 dest, src, imm**

  xmm/m128

  **ymm**

  xmm/m128

- **New Broadcast (SP, DP, 128-bit)**
  - Efficient Vector * Scalar operations

  **vBroadcastPS dest, mem32**

  m32   X0

  DEST | X0 | X0 | X0 | X0 | X0 | X0 | X0 | X0
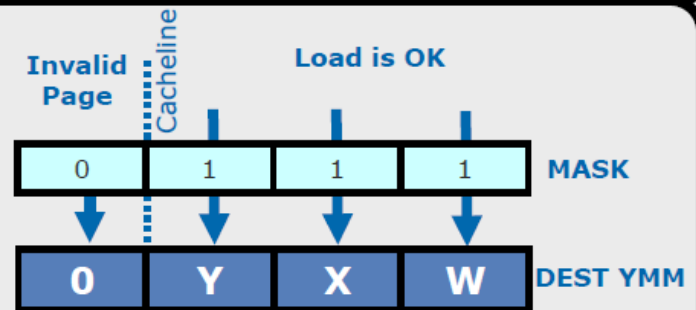
- **New Conditional SIMD Loads and Stores**
  - Avoid page faults, segment violations, memory transaction if the mask is 0
  - Allow more automatic compiler vectorization

  **vMaskMovPD dest, mask, mem256**

  Invalid Page | Cacheline | Load is OK

  | 0 | 1 | 1 | 1 | MASK

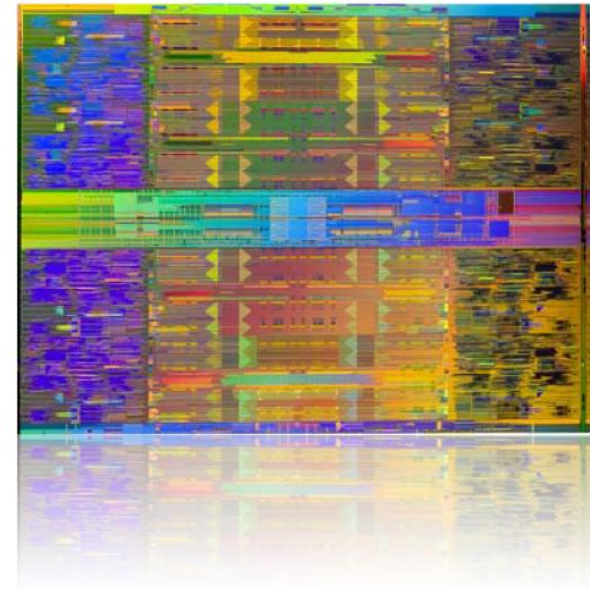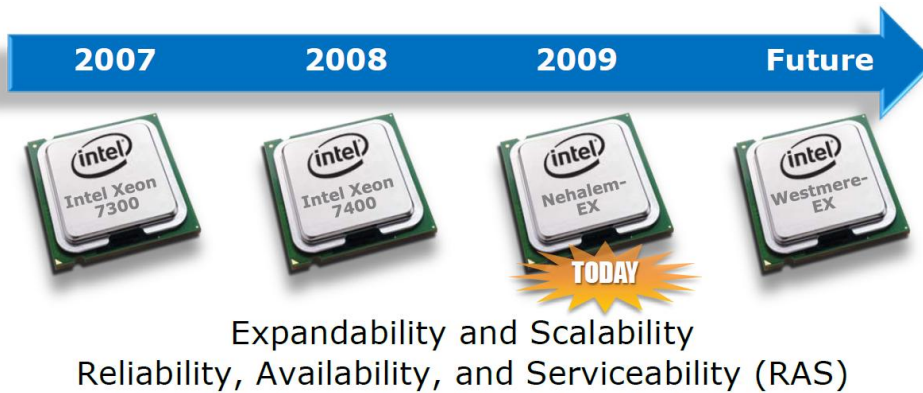  | 0 | Y | X | W | DEST YMM

# Morgan Stanley

# Intel AVX software support

- Intel Compiler: 11.1

- Intel IPP 6.1

- Intel MKL 10.2

- Intel TBB 2.2

- GNU GCC 4.4.1

- Microsoft Visual Studio 2010 Beta 2


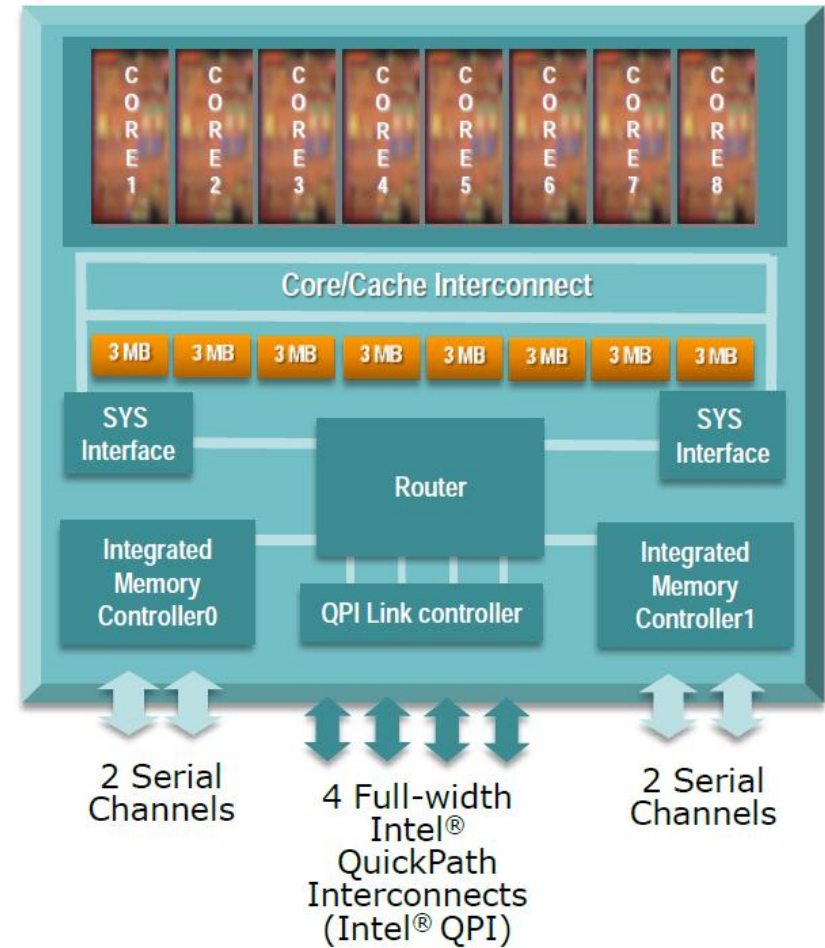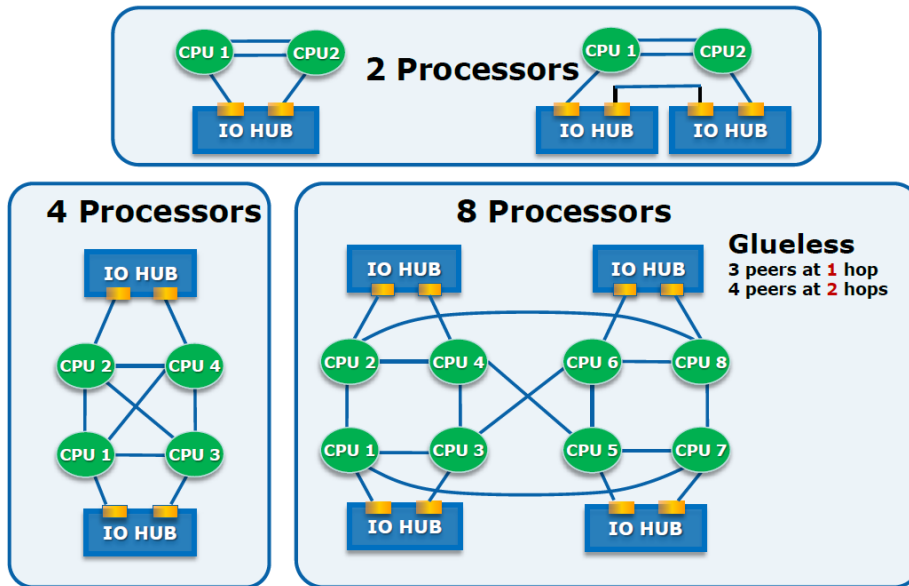- **Update on Fused Multiply Add**: Not supported on Sandy Bridge Microarchitecture

Morgan Stanley

# Intel Nehalem-EX Architecture

- **Scalable** Nehalem-EX Architecture

- Reliability, Availability and Serviceability **(RAS)**

- The platform for **x86-64 High-End Computing**

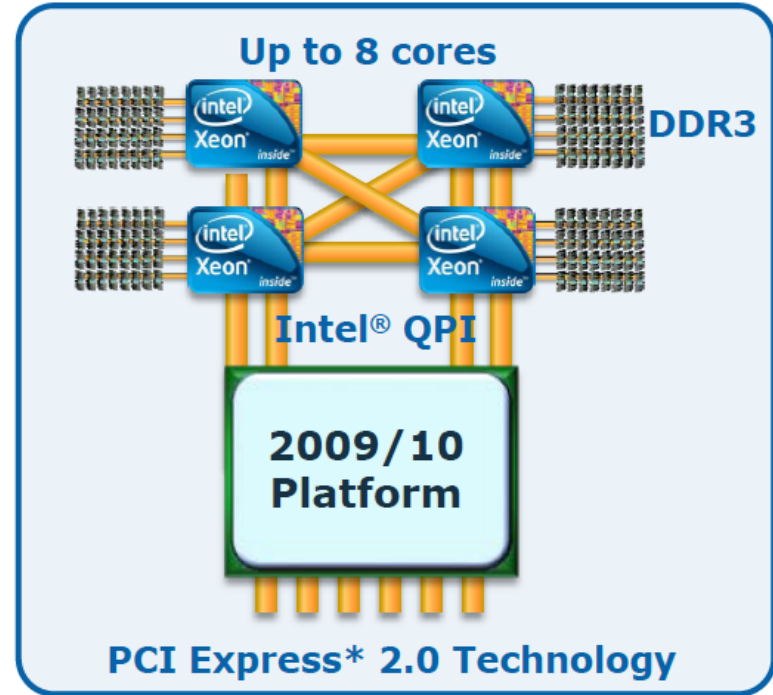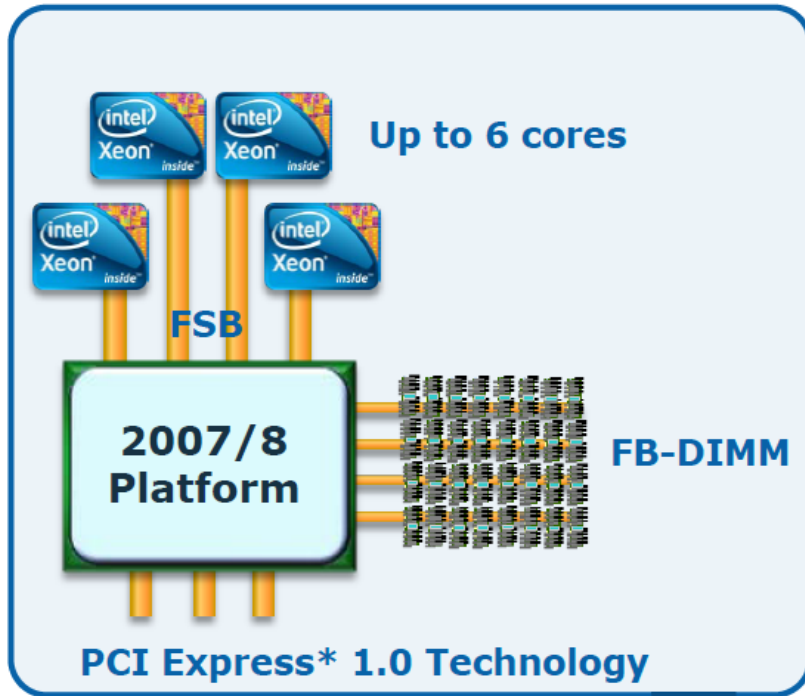- Several **uncore improvements**

2007    2008    2009    Future

Intel Xeon 7300    Intel Xeon 7400    Nehalem-EX    Westmere-EX

TODAY

Expandability and Scalability
Reliability, Availability, and Serviceability (RAS)

# Scalable Processor Cores

- Up to **8 cores** per socket

- **24 MB shared last level cache**

- 2,4 and **up to 8+ processors**



**2 Processors**

**4 Processors**

**8 Processors**

**Glueless**
3 peers at 1 hop
4 peers at 2 hops



Core/Cache Interconnect

3 MB 3 MB 3 MB 3 MB 3 MB 3 MB 3 MB 3 MB

SYS Interface

Router

SYS Interface

Integrated Memory Controller0

QPI Link controller

Integrated Memory Controller1

2 Serial Channels

4 Full-width Intel® QuickPath Interconnects (Intel® QPI)

2 Serial Channels

Morgan Stanley

# Platform Capability Enhancement
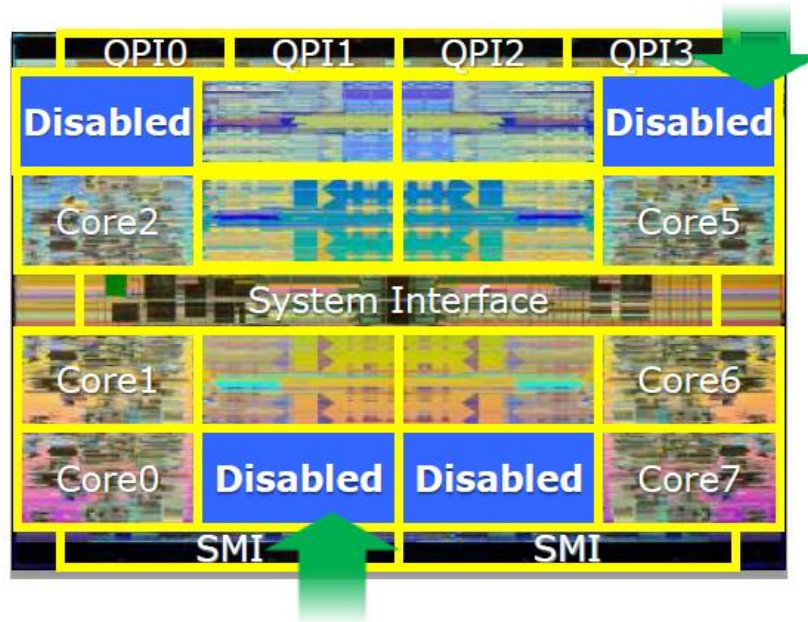


**Architectural Innovations:**
- Front Side Bus (FSB) to Intel® QuickPath Interconnect with integrated memory controller
- FB-DIMM to DDR3
- PCI Express* Technology (PCIe) 1.0 to 2.0...
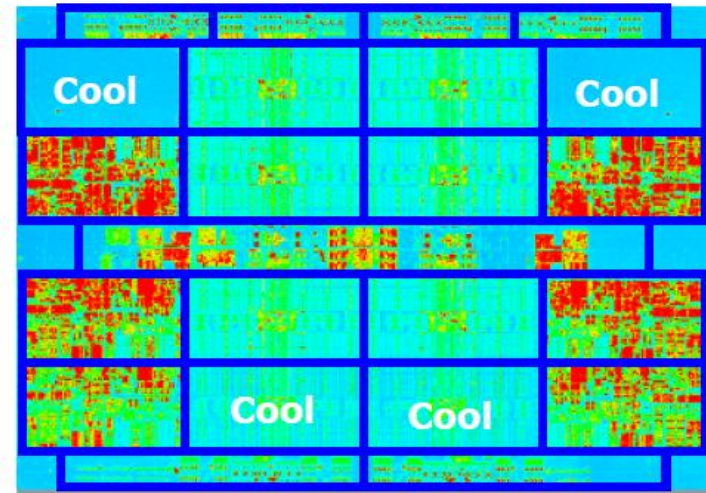
Morgan Stanley

# Power Optimization / Energy efficiency



## Leakage Control at Core and Cache Level

**New to NHM-EX:** Power gating disabled cores
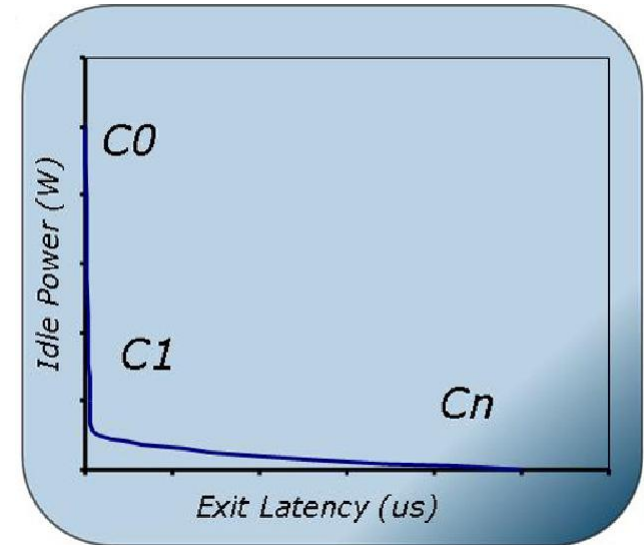
Power gating disabled cache

Infrared back-side emission shows core leakage suppressed and cache leakage reduced

**Core Level: 40x leakage reduction when shut off**
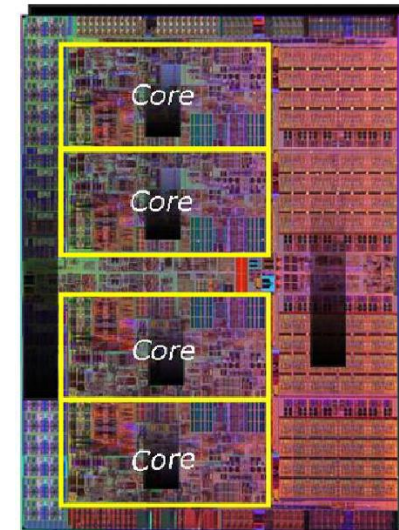**Cache Level: 35% during sleep & 83% when shut off**

Morgan Stanley

# Reducing Idle Power Consumption

- OS initiates C-state en entry by Mwait instruction
  - **C-states** – "idle" states
    - Different C states means different exit latency
    - Windows CPU Power Management Framework (Windows 7)
  - **P-states** – "active" states (C0)



## Morgan Stanley

# Deep Power Down Technology (C6)

- Integrated Power Gate enables **a per core C6 state** and individual cores transition to a **~0W Power State**

- Uncore logic? When all cores in C6, **package can transition to C6**

# Near-Term Server Futures


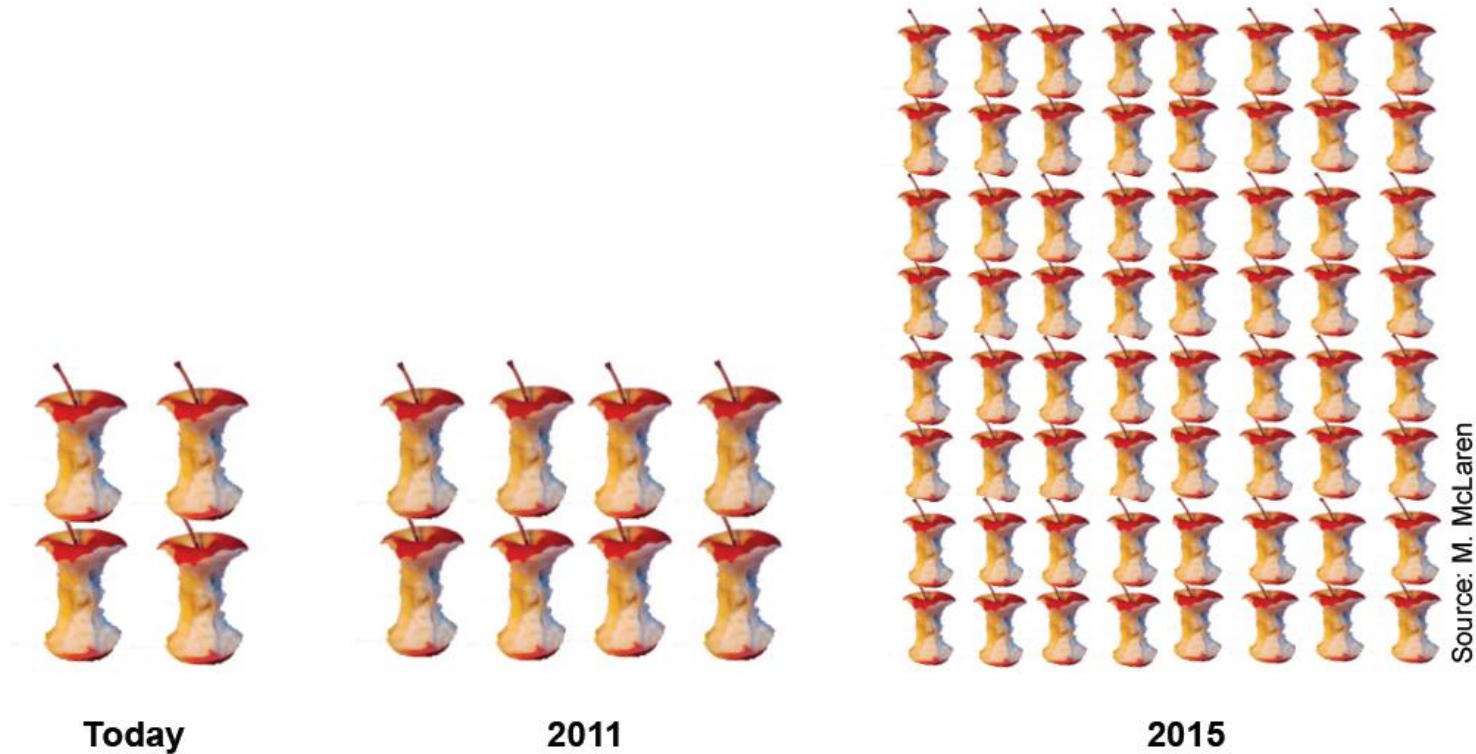
All you need to know about Multi-Core

Today          2011          2015

Source: M. McLaren

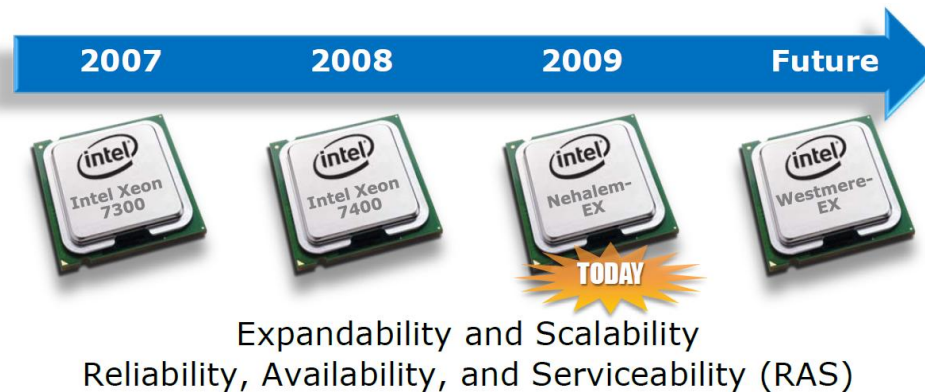# Company server architecture

- HPC segment is moving towards high density commodity hardware

- **Past:** 2-socket server (2+2 / 4+4 cores), P4 / Opteron processors

- **Now**: 2-socket servers (4+4 cores), Core2 architecture

- **Future**: 2-socket servers (4+4 / 6+6 cores), Nehalem / Westmere architecture

- **Focus:** best performance per unit price, not absolute performance

Morgan Stanley

# What Are The New Cores Like?

- Just like old ones

- **Clock rates NOT going up**
  - **Incore improvements**: more FLOPS/tick, better virtualization support, HTT
  - **Uncore improvements**: integrated memory controllers, etc.

- **Takes a little pressure off of memory latency**
  - Memory used to grow twice as distant (in clocks) every eighteen months

- **Does not take pressure off of memory bandwidth**

- **Bandwidth**: 0.5B/FLOP is 'good' for HPC
  - Nehalem gives us 40 GB/s : (3Ghz * 4 F/tick * 4 cores * 2 sockets ) = 0.4

- **Capacity**: up to 1.5, 3, 6GB/core (48GB in a 2s-4c system)

Morgan Stanley

# Prognostication

- **Nealem: wonderful gift of memory bandwidth**

  - 1.3x – 3x performance boost

- **Core count ++ means lower bandwidth/core**

  - Westmare @ 6 cores

  - Nehalem-EX @ 8 cores, but more bandwidth



## Morgan Stanley

# Evolutionary Improvements Running Out of Gas

- **S/W has been relying on H/W improvements**
  - Moore's Law hid inefficiencies
    - "Every programmer's job is to consume twice as much compute power every eighteen moths!"
  - Get three programmers together, and they'll first have to define nine levels of abstraction to ensure optimal job security.

- **S/W needs to be come more efficient**
  - Requires effort and expertise from S/W part
  - Olukotun's view 3 real dimensions of processor performance:
    - clock frequency
    - superscalar instruction issue
    - multiprocessing
  - First two have been pretty much exhausted, programmers need to switch.

Morgan Stanley

# The Memory Wall

- **2014: assume 400GF mainstream processor**
  - Aggressively multi-core
  - 0.5 B/FLOP means 200 GB/s, ~10x today's chips
    - **Opinion: won't happen in copper**
      - Optical connection
      - Stacked memory (more levels)
      - ?

  - … maybe the hardware folks can save your sorry selves once again!

Morgan Stanley

# Q/A

Morgan Stanley