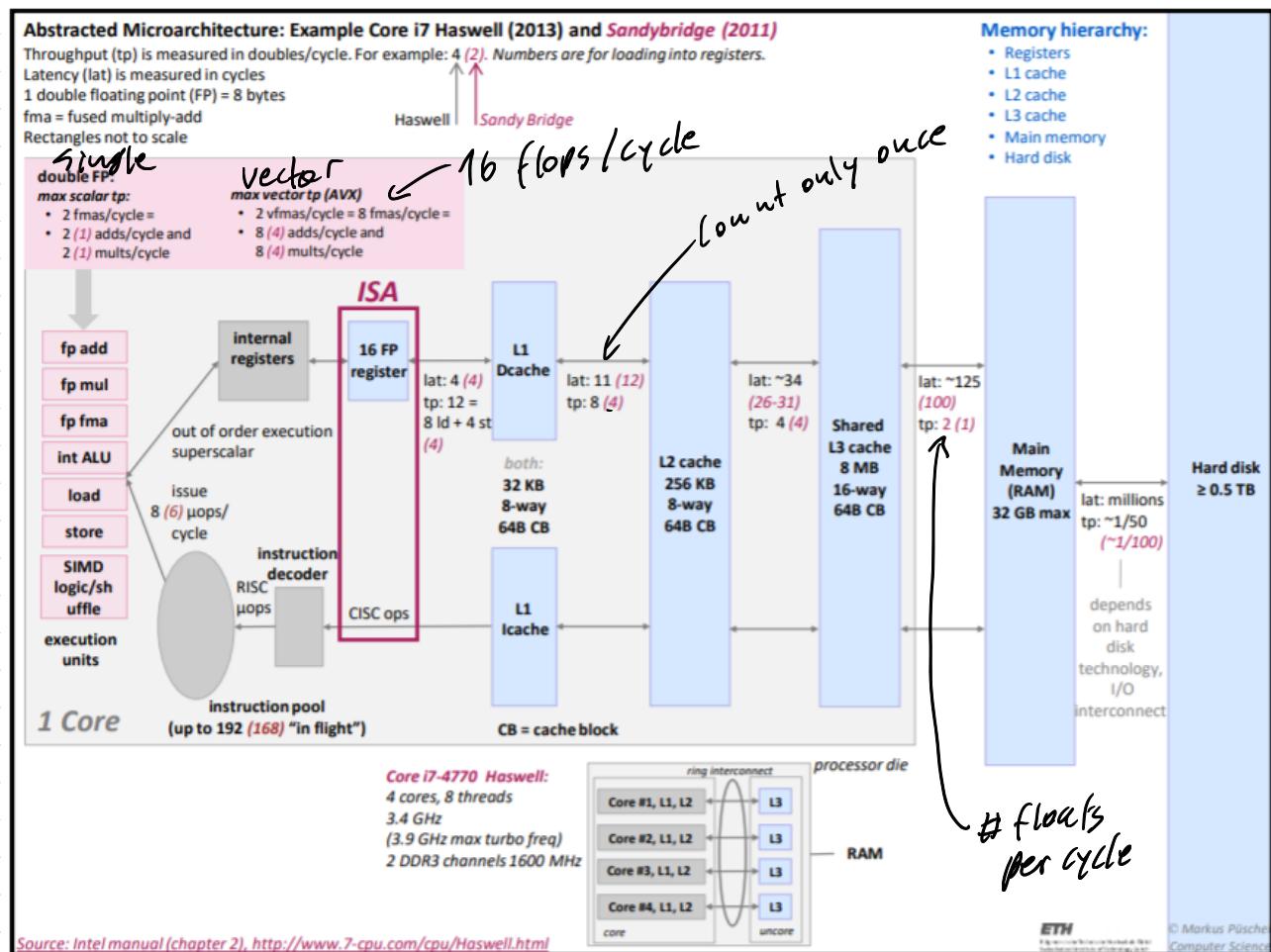


Zusammenfassung Advanced Systems Las

- $\Theta(f(n)) = \underbrace{\text{Average}}_{\text{Avg}}(f(n)) \cap \underbrace{\text{Best}}_{\text{Best}}(f(n)) \cap \underbrace{\text{Worst}}_{\text{Worst}}(f(n))$
- $2n + O(\log(n))$ is canonical
- $\underbrace{\Theta(2n)}_{\text{rest}} + \underbrace{\log(n)}_{\text{precise}}$ isn't canonical because the rest must be smaller $\Rightarrow O(n)$
- Matrix Matrix Multiplication:
 - Runtime $O(n^3)$
 - Space $O(n^2)$
- $\Theta(f(n))$ only describes eventual trend of runtime (doesn't tell us tightness, or how start looks)
- Cost measure:
 - $C(n) = N_{\text{add}} \cdot C_{\text{add}} + N_{\text{mult}} \cdot C_{\text{mult}}$
- CPUs can do arithmetic (flops) and indexing (ints) at the same time because there are a Arithmetic LU and Integer LU
- Performance:
 - flops / cycle \leq #flops / #cycles
 - upper bound "peak performance" given by CPU
- Asymptotic Runtime vs. cost
- Architecture/instruction set architecture:
 - Assembly
 - Registers
 - e.g. x86
- Microarchitecture:
 - Implementation of Architecture
 - Cache structures, frequency, ...
 - Not relevant for writing code
- SIMD:
 - Do vector instructions

- Intel's Tick Tock Model

- Tick: Shrink of process technology
- Tock: New Microarchitecture (= new names)
- Since 2016: process- architecture-optimization
- CISC:
 - Complex instruction set computer
 - Intel's architecture
 - Is converted into RISC internally
- FMA:
 - 1 FMA \equiv 2 flops \equiv 1 fadd + 1 fmulf
 - Recently implemented, because before Intel supported only 2 inputs
 - FMA: $X * Y + Z$



```
/* x, y are vectors of doubles of length n, alpha is a double */
for (i = 0; i < n; i++)
    x[i] = x[i] + alpha*y[i];
```

FMA

maximal achievable percentage
(vector) peak performance

- Number flops?
- Runtime bound no vector ops:
- Runtime bound vector ops:
- Runtime bound data in L1:
- Runtime bound data in L2:
- Runtime bound data in L3:
- Runtime bound data in main memory:

2n	
n/2	cycles needed for n length
n/8	
n/4	50
n/4	50
n/2	25
n	12.5

```
/* matrix multiplication; A, B, C are n x n matrices of doubles */
for (i = 0; i < n; i++)
    for (j = 0; j < n; j++)
        for (k = 0; k < n; k++)
            C[i*n+j] += A[i*n+k]*B[k*n+j];
```

- Number flops?
- Runtime bound no vector ops:
- Runtime bound vector ops:
- Runtime bound data in L1:
- ...
Runtime bound data in main memory:

→ Bottleneck is memory

- Operational Intensity: $I(n) = \frac{W(n)}{Q(n)}$ # flops with n
bytes transferred
cache ↔ memory

→ If we look at upper bound, we only take matrix: n^2
vector: n once

· 1 double floating point $\equiv 8$ bytes
· Low $I(n) \equiv$ memory bottleneck

```
/* x, y are vectors of doubles of length n, alpha is a double */
for (i = 0; i < n; i++)
    x[i] = x[i] + alpha*y[i];
```

C is in register

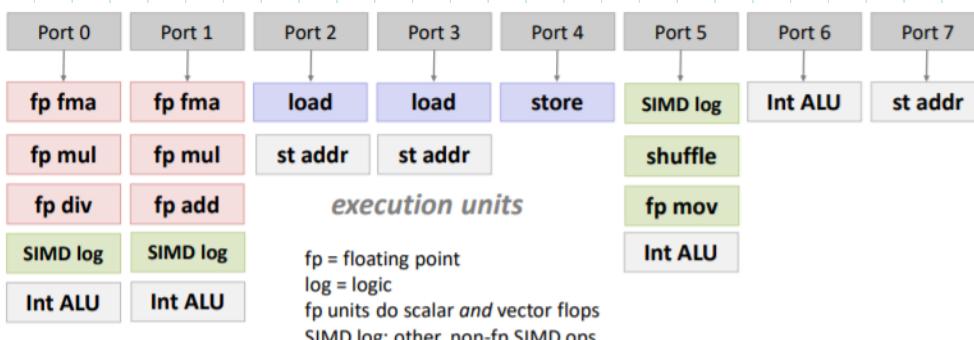
```
/* matrix multiplication; A, B, C are n x n matrices of doubles */
for (i = 0; i < n; i++)
    for (j = 0; j < n; j++)
        for (k = 0; k < n; k++)
            C[i*n+j] += A[i*n+k]*B[k*n+j];
```

Operational intensity:

- Flops: $W(n) = 2n$
- Memory transfers (doubles): $\geq 2n$ (just from the reads)
- Reads (bytes): $Q(n) \geq 16n$
- Operational intensity: $I(n) = W(n)/Q(n) \leq 1/8$

Operational intensity:

- Flops: $W(n) = 2n^3$
- Memory transfers (doubles): $\geq 3n^2$ (just from the reads)
- Reads (bytes): $Q(n) \geq 24n^2$
- Operational intensity: $I(n) = W(n)/Q(n) \leq n/12$



When data is in RAM,
written vectors are
first read

→ Every port can issue one instruction/cycle

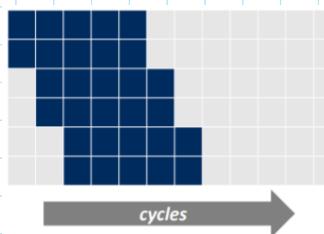
· $\text{CPI} = 1/\text{Throughput}$ (cycles per issue)

· e.g. FMA throughput (ops/cycle) = 2

· Same throughput/latency for vector instructions,
but then, each instruction has multiple floats



Each register:
256 bits = 4 doubles = 8 singles



Throughput tp = 2/cycle

Gap = 1/tp = 1/2 cycles/issue

How many cycles are at least required (no vector ops)?

- A function with n adds and n mults in the C code
- A function with n add and n mult instructions in the assembly code
- A function with n adds in the C code
- A function with n add instructions in the assembly code
- A function with n adds and $n/2$ mults in the C code

$n/2 \rightarrow$ converted to FMA
 $n \rightarrow$ converted to FMA
 $n/2 \rightarrow$ converted to FMA
 $n \rightarrow$
 $n/2 \rightarrow$

- Peak performance of computer:

- 4 cores
- 2-way SSE
- 1 add + 1 mult / cycle
- 3 GHz
- ↳ 48 Gflop/s

- Instruction level parallelism (ILP)

- Can taken advantage of by superscalar computers, because they can execute multiple instructions in one cycle
- Sequential dependence \Rightarrow no ILP!

- Loop unrolling

- Use k accumulators

\hookrightarrow The deeper the pipeline (latency), the more ILP you need

- Need to unroll by L , k divides L

- Large overhead for short lengths

- Best $k = \text{ceil}(\text{latency} * \text{throughput})$

\hookrightarrow In every cycle you can start two mults, and they need 5 cycles to finish, so you can have 10 floats in the pipeline

- Compiler limitations

- Associativity law doesn't hold for floats
- Move code out of loop
- Replace costly computations with simpler ones
- Rerse portions of expressions
- Removing procedure calls
- Function inlining, if source code is compiled together
- Compiler usually treats procedure call as black box
- Memory accessed every iteration \Rightarrow scalar replacement
- Memory aliasing: Memory can overlap, this leads to other results, if we use local variable, which doesn't have aliasing

- Single vs double precision

- Single: 4 bytes \rightarrow float

- Double: 8 bytes $\hat{=}$ double

- Intrinsics

- Assembly coded in C functions

- You need to code explicit loads and stores

- Use aligned loads and stores as much as possible

- \hookrightarrow malloc allocation guarantees aligned allocation

- Never divide, multiply by inverse instead

0 1 2 3 4 5 6 7	
a a b b b b c d	bytes
	words

The problem is that on some CPU architectures, the instruction to load a 4-byte integer from memory only works on word boundaries. So your program would have to fetch each half of `b` with separate instructions.

But if the memory was laid out as:

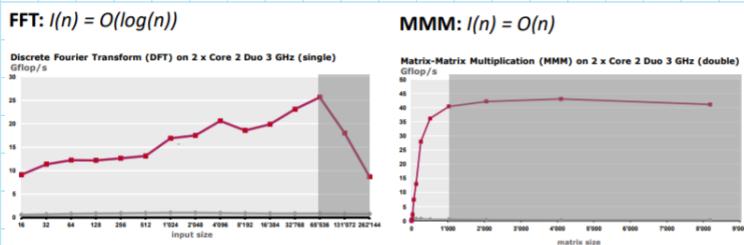
0	1	2	3	4	5	6	7	8	9	A	B
a a		b b b b c d									

Then access to `b` becomes straightforward. (The disadvantage is that more memory is required, because of the padding bytes.)

Different data types have different alignment requirements. It's common for `char` to be 1-byte aligned, `short` to be 2-byte aligned, and 4-byte types (`int`, `float`, and pointers on 32-bit systems) to be 4-byte aligned.

`malloc` is required by the C standard to return a pointer that's properly aligned for any data type.

- Temporal locality: Recently used data will likely be used again
 - Naturally supported by cache
 - Spatial Locality: Data nearby will likely be accessed
 - Supported by transferring data in blocks



**Up to 40-50% peak
Performance drop outside last level cache (LLC)
Most time spent transferring data**

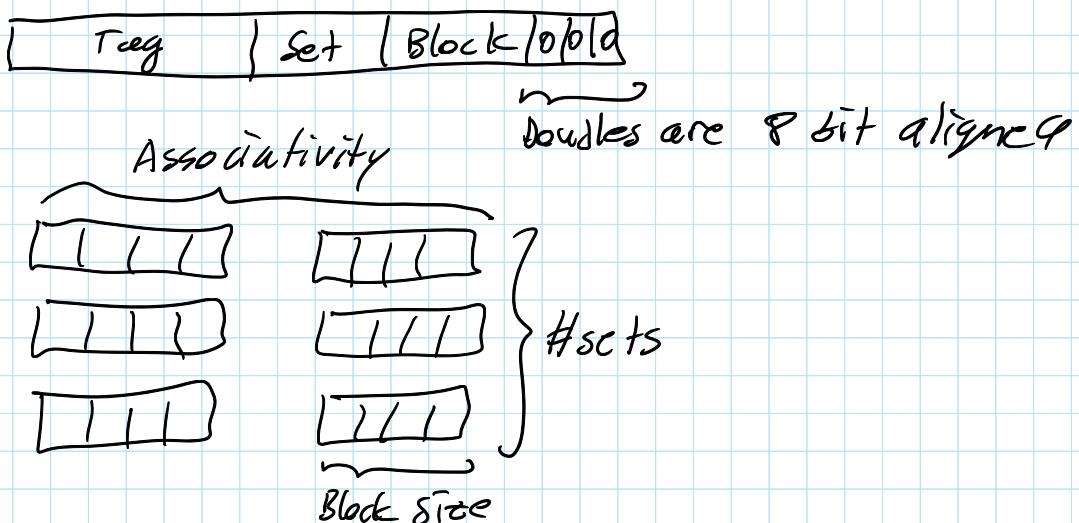
**Up to 80-90% peak
Performance can be maintained
outside LLC
*Cache miss time compensated/hidden
by computation***

- Asymptotic bounds on operational intensity
 - Vector sum: $y = x + y \asymp O(n/n) = O(1) \Rightarrow$ Not good
 - Matrix-Vector: $y = Ax \asymp O(n^2/n^2) = O(1)$
 - FFT $\asymp O(\log(n))$
 - MMM: $(= AB + C) \asymp O(n^3/3n^2) = O(n)$

- Cache misses

- Compulsatory
 - Capacity
 - Conflict (multiple data blocks are mapped to same location)

- Cache Structure

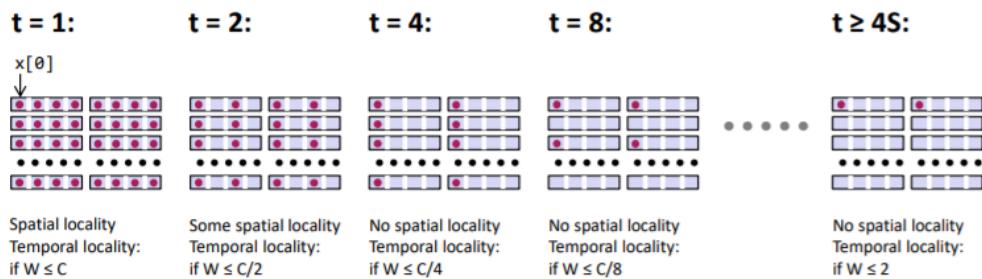


- Tag: rest of address that is stored in cache for reconstruction of address
 - Direct mapped:
 - $E = 1$
 - Every address has unique location in cache
 - Fully associative:
 - #sets = 1
 - No conflict misses
 - LRU policy for replacement
 - If data isn't in cache, it has first to be read, after it can be written to

The Killer: Two-Power Strided Working Sets

```
% t = 1,2,4,8... a 2-power
% size W of working set: W = n/t
for (i = 0; i < n; i += t)
    access(x[i])
for (i = 0; i < n; i += t)
    access(x[i])
```

Cache: E = 2, B = 4 doubles



Working with a two-power-strided working set is like having a smaller cache

↳ Power of t_0 is regular

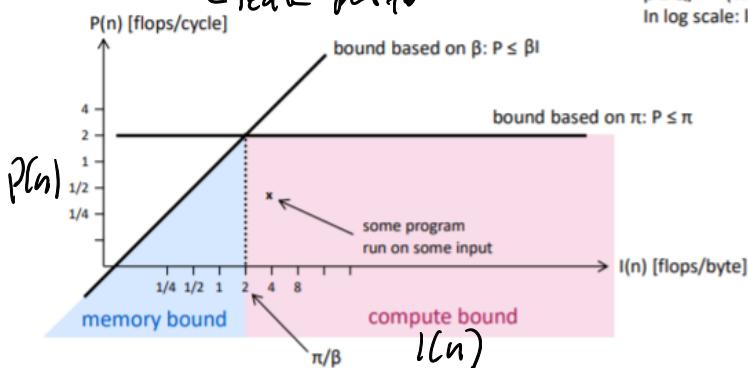
- Roofline Model

- Operational intensity: $I(n) = \frac{w(n)}{Q(n)}$ # flops # bytes
 - Performance: $P(n) = \frac{w(n)}{T(n)}$ # flops # cycles
- ↖ Bandwidth bytes/cycle

Example: one core, $\pi = 2$, $\beta = 1$, no SIMD

peak performance

Bound based on β :
 $\beta \geq Q/T = (W/T)/(W/Q) = P/I$
In log scale: $\log_2(P) \leq \log_2(\beta) + \log_2(I)$



- What happens, if we introduce 4-way SIMD
 - T changes
 - β stays the same
 - More programs become memory bound

- Linear algebra software

- Blas: Kernel function implemented for each computer (generated by Atlas)
- Lapack: static higher level functions,

Dependencies

■ Read-after-write (RAW) or true dependency

$$\begin{array}{ll} W & r_1 = r_3 + r_4 \\ R & r_2 = 2r_1 \end{array} \quad \begin{array}{l} \text{nothing can be done} \\ \text{no ILP} \end{array}$$

■ Write after read (WAR) or antidependency

$$\begin{array}{ll} R & r_1 = r_2 + r_3 \\ W & r_2 = r_4 + r_5 \end{array} \quad \begin{array}{l} \text{dependency only by} \\ \text{name} \rightarrow \text{rename} \end{array} \quad \begin{array}{ll} r_1 = r_2 + r_3 \\ r_2 = r_4 + r_5 \end{array} \quad \begin{array}{l} \text{now ILP} \end{array}$$

■ Write after write (WAW) or output dependency

$$\begin{array}{ll} W & r_1 = r_2 + r_3 \\ \dots & \\ W & r_1 = r_4 + r_5 \end{array} \quad \begin{array}{l} \text{dependency only by} \\ \text{name} \rightarrow \text{rename} \end{array} \quad \begin{array}{ll} r_1 = r_2 + r_3 \\ \dots \\ r_1 = r_4 + r_5 \end{array} \quad \begin{array}{l} \text{now ILP} \end{array}$$

Each variable
is assigned
exactly once

- Sparse linear algebra

A as matrix

b	c	c
a		
	b	b
c		

A in CSR:

values	b	c	c	a	b	b	c	length K
col_idx	0	1	3	1	2	3	2	length K
row_start	0	3	4	6	7			length m+1

- Upper bound on operational intensity

$$I(n) = 2K / 8(K + 3n)$$

$$\left\{ \begin{array}{l} y = Ax + b \\ \end{array} \right.$$