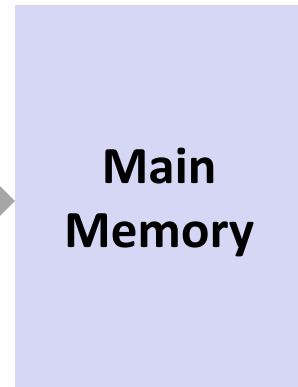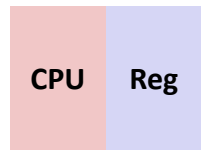# Organization

- **Temporal and spatial locality**

- **Operational intensity, memory/compute bound**

# Problem: Processor-Memory Bottleneck

*Processor performance doubled about every 18 months*

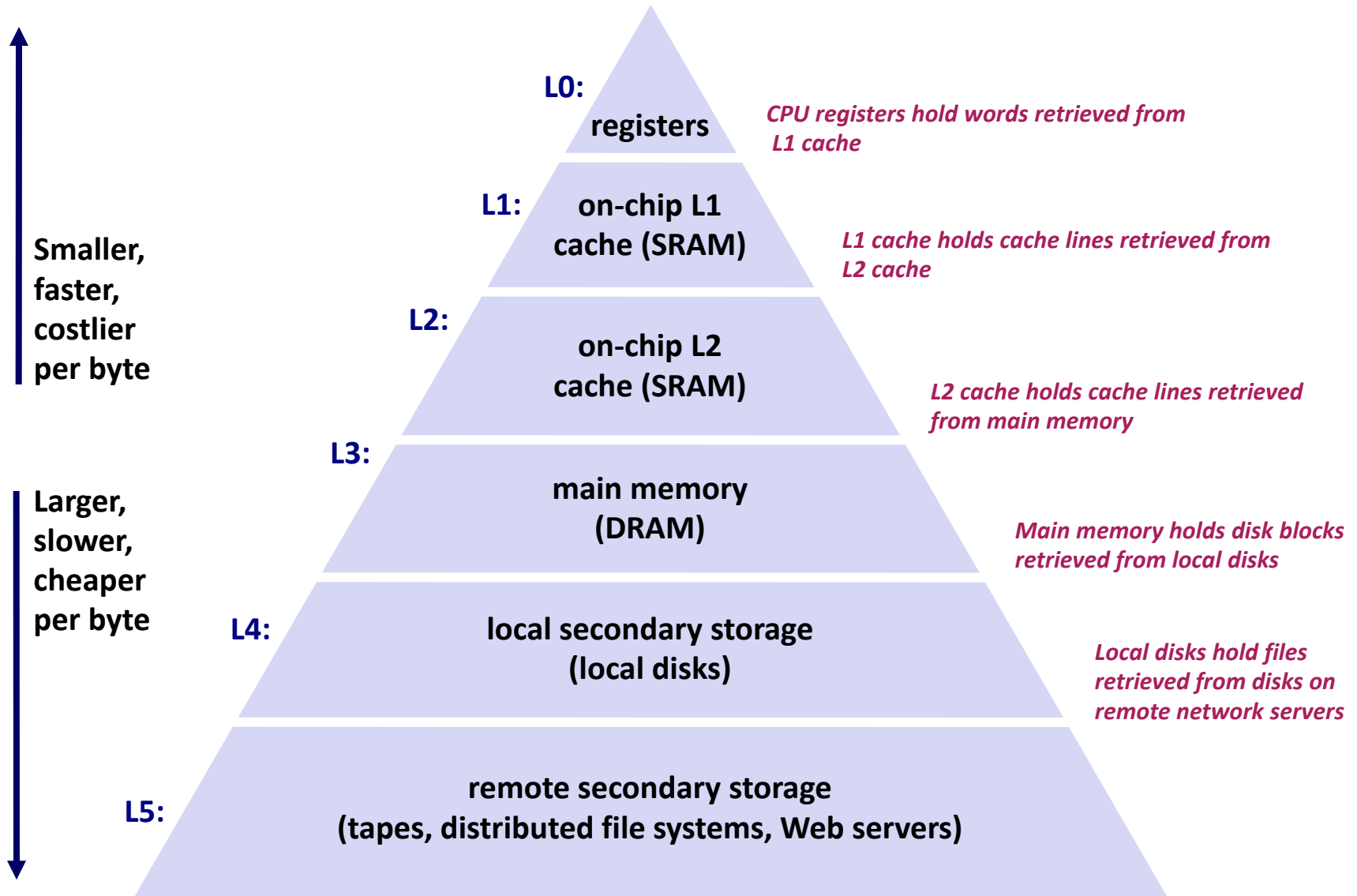*Bus bandwidth evolved much slower*

**CPU** **Reg**

**Main Memory**

*Core 2 Duo:*
**Peak performance:**
**2 SSE two operand ops/cycles**
consumes up to 64 Bytes/cycle

*Core 2 Duo:*
**Bandwidth**
2 Bytes/cycle

*Solution: Caches/Memory hierarchy*

# Typical Memory Hierarchy

**Smaller,
faster,
costlier
per byte**

**Larger,
slower,
cheaper
per byte**

**L0:**
registers

*CPU registers hold words retrieved from
L1 cache*

**L1:** **on-chip L1
cache (SRAM)**

*L1 cache holds cache lines retrieved from
L2 cache*

**L2:** **on-chip L2
cache (SRAM)**

*L2 cache holds cache lines retrieved
from main memory*

**L3:** **main memory
(DRAM)**

*Main memory holds disk blocks
retrieved from local disks*

**L4:** **local secondary storage
(local disks)**

*Local disks hold files
retrieved from disks on
remote network servers*

**L5:** **remote secondary storage
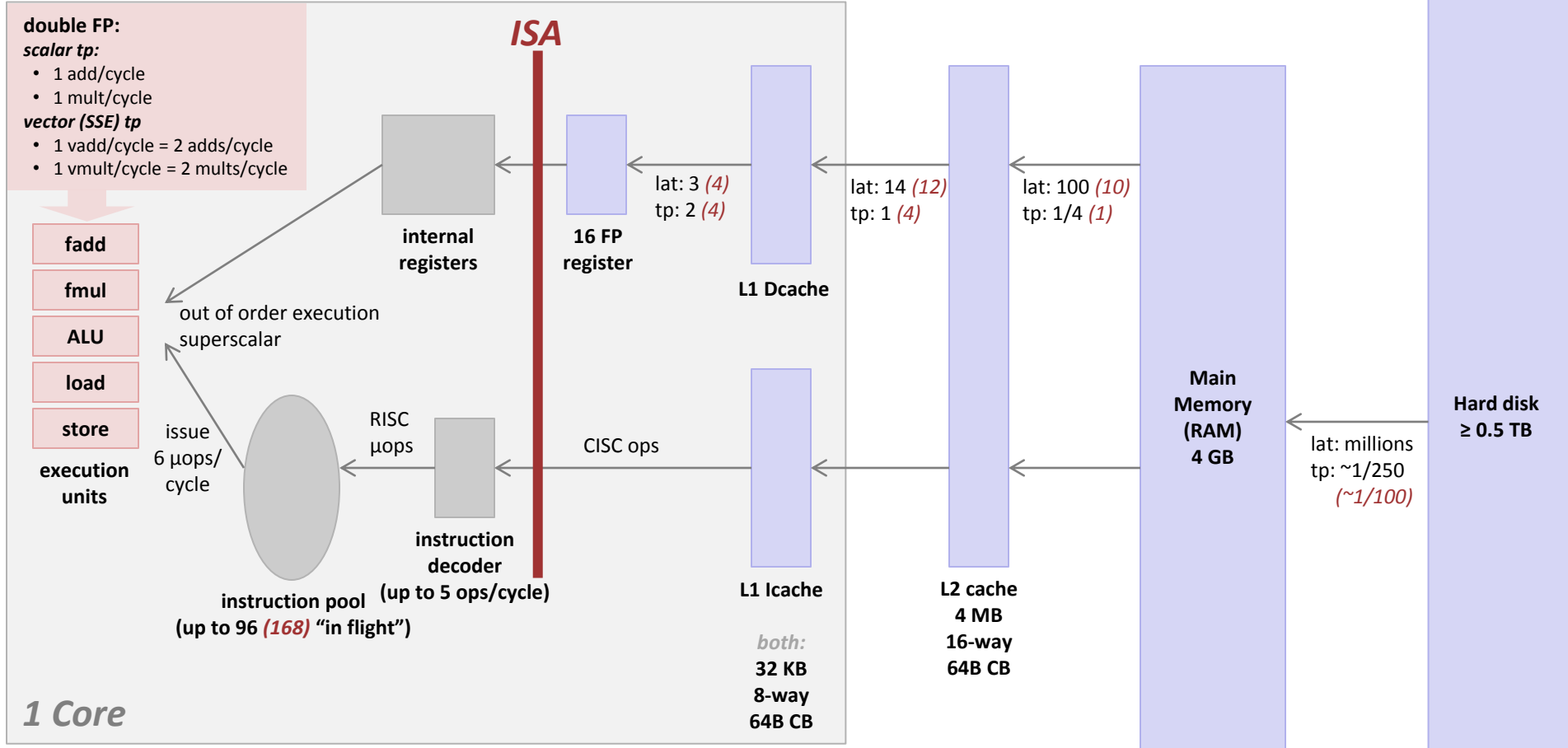(tapes, distributed file systems, Web servers)**

**Abstracted Microarchitecture: Example Core**

Throughput (tp) is measured in doubles/cycle. For example: 2 *(4)*
Latency (lat) is measured in cycles
1 double floating point (FP) = 8 bytes
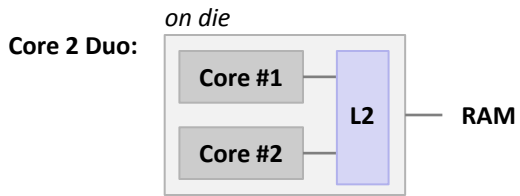Rectangles not to scale

*© Markus Püschel*
*Computer Science*

**ETH** Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

**Memory hierarchy:**
- Registers
- L1 cache
- L2 cache
- Main memory
- Hard disk

Core 2 (2008)
*Core i7 Sandy Bridge (2011)*

*ISA*

**double FP:**
*scalar tp:*
- 1 add/cycle
- 1 mult/cycle

*vector (SSE) tp*
- 1 vadd/cycle = 2 adds/cycle
- 1 vmult/cycle = 2 mults/cycle

- fadd
- fmul
- ALU
- load
- store

**execution units**

out of order execution superscalar

**internal registers**

**16 FP register**

lat: 3 *(4)*
tp: 2 *(4)*

**L1 Dcache**

lat: 14 *(12)*
tp: 1 *(4)*

lat: 100 *(10)*
tp: 1/4 *(1)*

issue 6 μops/ cycle

RISC μops

CISC ops

**instruction decoder (up to 5 ops/cycle)**

**instruction pool (up to 96 *(168)* "in flight")**

**L1 Icache**

*both:*
**32 KB 8-way 64B CB**

**L2 cache 4 MB 16-way 64B CB**

**Main Memory (RAM) 4 GB**

**Hard disk ≥ 0.5 TB**

lat: millions
tp: ~1/250
*(~1/100)*

*1 Core*

Core 2 Duo:

*on die*

Core #1
Core #2
L2 — RAM

*Core i7 Sandy Bridge:*
*256 KB L2 cache*
*2–8MB L3 cache: lat 26-31, tp 4*
*vector (AVX) tp*
- 1 vadd/cycle = 4 adds/cycle
- 1 vmult/cycle = 4 mults/cycle

*on die*
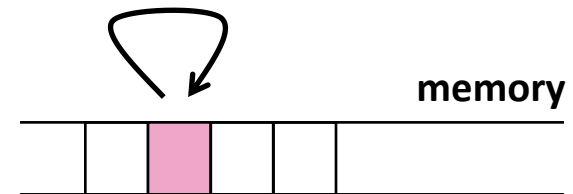
Core #1  L2
Core #2  L2
Core #3  L2
Core #4  L2
L3 — RAM

# Why Caches Work: Locality

- *Locality:* **Programs tend to use data and instructions with addresses near or equal to those they have used recently**
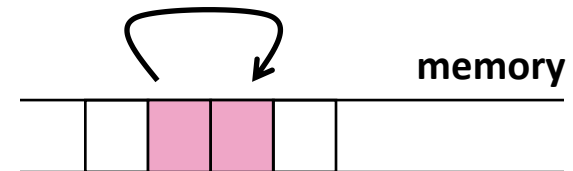  *History of locality*

- *Temporal locality:*

  Recently referenced items are likely
  to be referenced again in the near future

  **memory**

- *Spatial locality:*

  Items with nearby addresses tend
  to be referenced close together in time

  **memory**

# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

- **Data:**
  - Temporal: **sum** referenced in each iteration
  - Spatial: array **a[]** accessed in stride-1 pattern

- **Instructions:**
  - Temporal: loops cycle through the same instructions
  - Spatial: instructions referenced in sequence

- *Being able to assess the locality of code is a crucial skill for a performance programmer*

# Locality Example #1

```
int sum_array_rows(int a[M][N])
{
  int i, j, sum = 0;

  for (i = 0; i < M; i++)
    for (j = 0; j < N; j++)
      sum += a[i][j];
  return sum;
}
```

# Locality Example #2

```
int sum_array_cols(int a[M][N])
{
  int i, j, sum = 0;

  for (j = 0; j < N; j++)
    for (i = 0; i < M; i++)
      sum += a[i][j];
  return sum;
}
```

# Locality Example #3

```
int sum_array_3d(int a[M][N][N])
{
  int i, j, k, sum = 0;

  for (i = 0; i < M; i++)
    for (j = 0; j < N; j++)
      for (k = 0; k < N; k++)
        sum += a[k][i][j];
  return sum;
}
```

**How to improve locality?**

# Memory/Compute Bound

- **Operational intensity of a program/algorithm:**

$$I = \frac{\textit{Number of operations}}{\textit{Amount of data transferred cache} \leftrightarrow \textit{RAM}}$$

- **Notes:**
  - *I* depends on the computer (e.g., the cache size and structure)
  - *Q:* Relation to cache misses?

    *A:* Denominator determined by misses in lowest level cache

- **This course usually:**
  - #ops = #flops
  - unit: flops/byte or flops/double

- *"Definition:"* **Programs with high *I* are called** *compute bound*, **programs with low *I* are called** *memory bound*

# Questions

- *Q:* **How high is high enough for compute bound?**

  *A:* **Depends on the computer; we will make this precise later with the roofline model**

- *Q:* **Estimate the operational intensity**

```
int sum_array_rows(int a[M][N])
{
  int i, j, sum = 0;

  for (i = 0; i < M; i++)
    for (j = 0; j < N; j++)
      sum += a[i][j];
  return sum;
}
```

# Upper Bound on I

- **Assume cold (empty) cache:**

  *Amount of data transferred cache $\longleftrightarrow$ RAM*

  *$\geq$ Size of input data + size of output data*

- **Hence:**

$$I \leq \frac{\text{Number of operations}}{\text{Size of input data + size of output data}}$$

- **Examples: Compute upper bounds of *I* for**

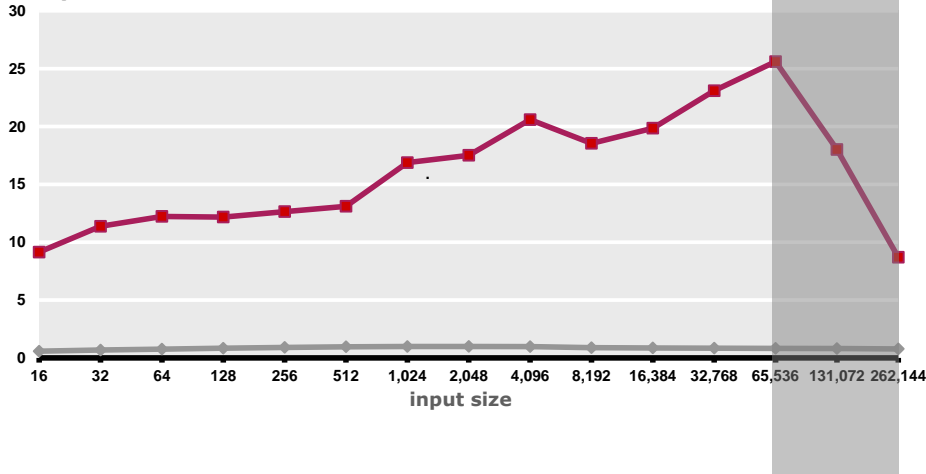  - Matrix multiplication C = AB + C $\quad I(n)$ • $\frac{2n^3}{3n^2} = \frac{2}{3}n = O(n)$

  - Discrete Fourier transform $\qquad I(n)$ • $\frac{5n\log_2(n)}{2n} = \frac{5}{2}\log_2(n) = O(\log(n))$

  - Adding two vectors x = x+y $\qquad I(n)$ • $\frac{n}{2n} = \frac{1}{2} = O(1)$

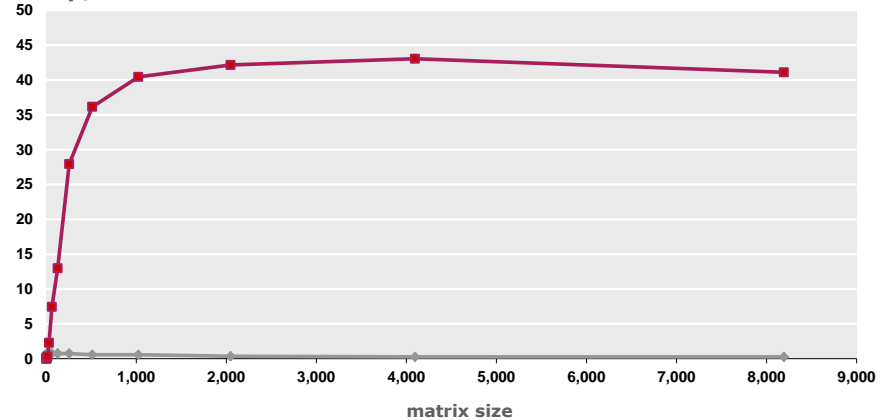# Effects

**FFT:** *I(n) ≤ O(log(n))*



**MMM:** *I(n) ≤ O(n)*



**Up to 40-50% peak**
**Performance drop outside L2 cache**
*Most time spent transferring data*

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