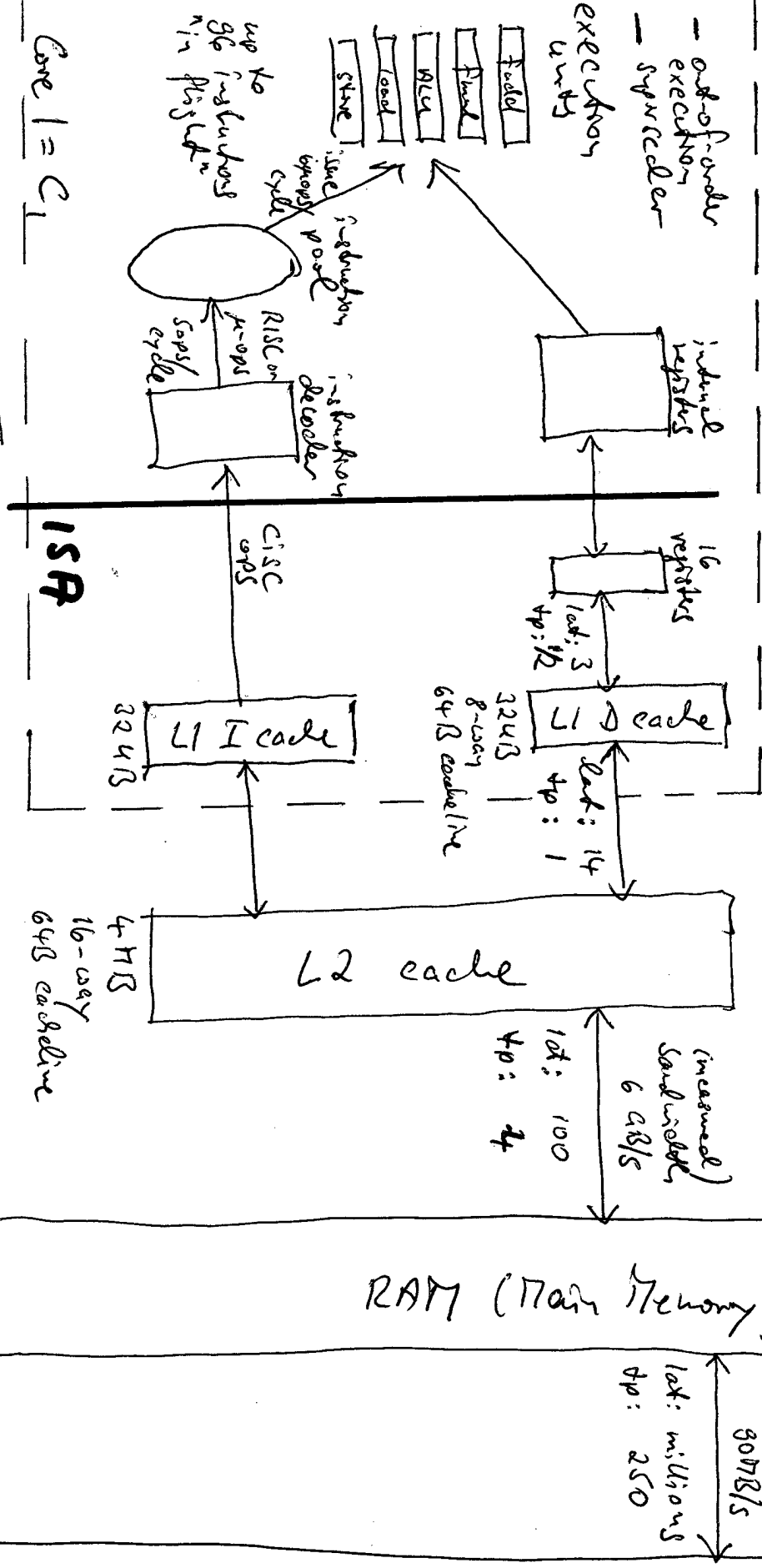


Microarch. level (Example: Core)  
 we measure latency and throughput in  
 (not in bytes per second as usual)

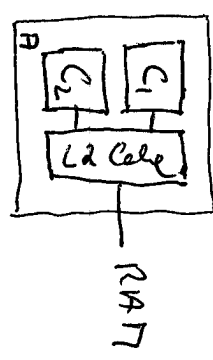
cycles per double  
~~double per cycle~~  
 1 double = 8 bytes

- out-of-order execution
- super-scalar

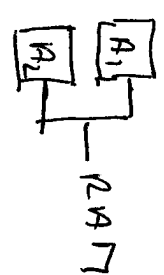


Core 1 = C<sub>1</sub>

Core 2 Duo:

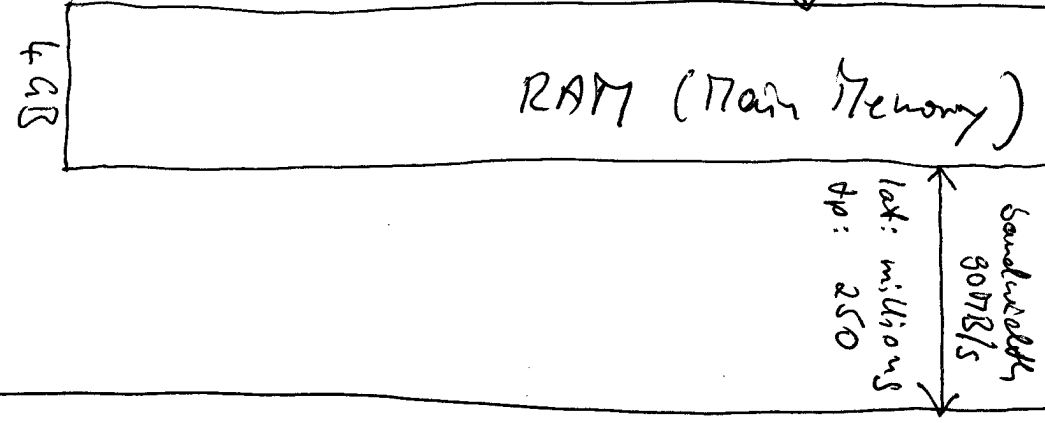


2 x Core 2 Duo:  
 (Core 2 Quad)



4 cores not on one die just D

Memory hierarchy:  
 Registers  
 Caches  
 Main memory  
 Hard disk



Hard Disk

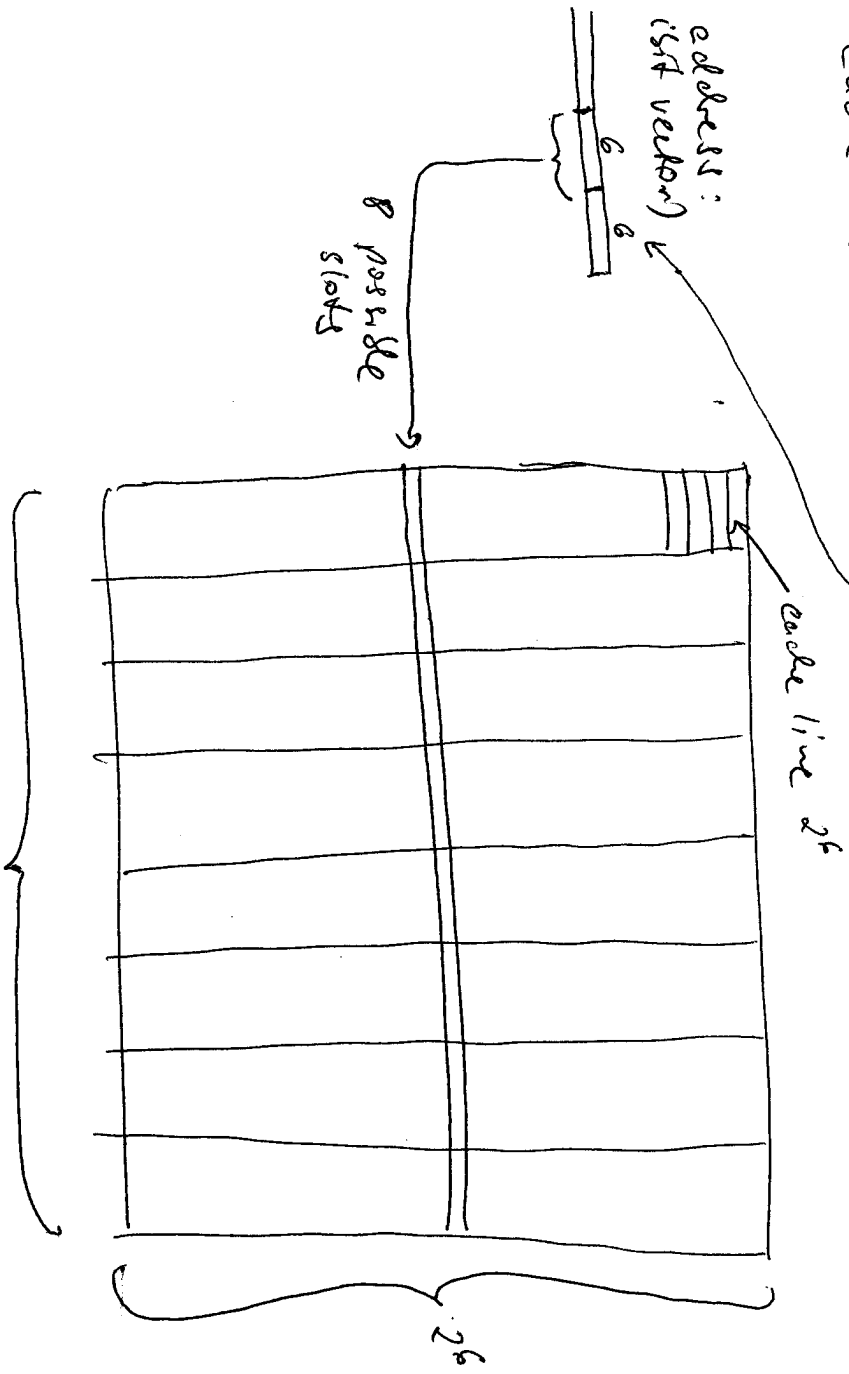
~500GB

# Cache structure (Example L1 D-cache Core)

Size:  $32\text{KB} = 2^{15}\text{B} = 2^{12}$  doubles

Associativity: 8-way =  $2^3$ -way

Cache line:  $64\text{B} = 2^6\text{B} = 2^3$  doubles



Example: Load vector of length  $n$  into empty cache

$\Rightarrow \frac{n}{8}$  (compulsory) misses

- memory is addressed by the cache
- cache line: smallest amount of data brought into cache
- cache misses (CM):
  - compulsory: CM because data is accessed the first time
  - capacity: CM because cache is full
  - conflict: CM because data is loaded into occupied slot