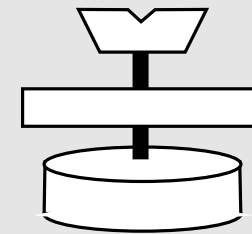
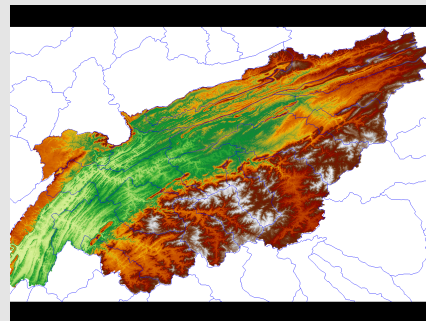
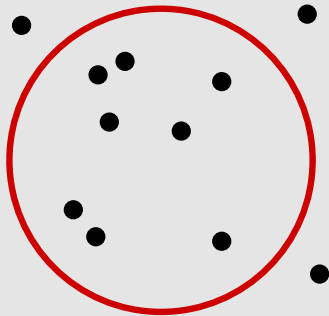


I/O- and Cache-Efficient Algorithms for Spatial Data

Mark de Berg

TU Eindhoven

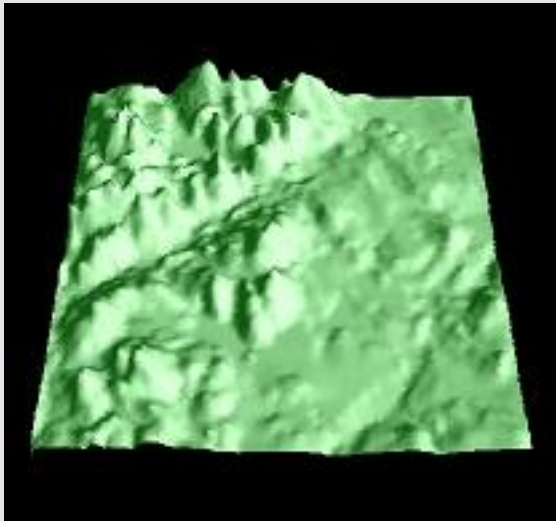


Massive data sets are becoming more and more common

- AT&T: data base of phone calls: 20 TB
- Wal-Mart: data base of buying patterns: 70 TB
- geographic data: NASA satellites collect more than 1 TB / day

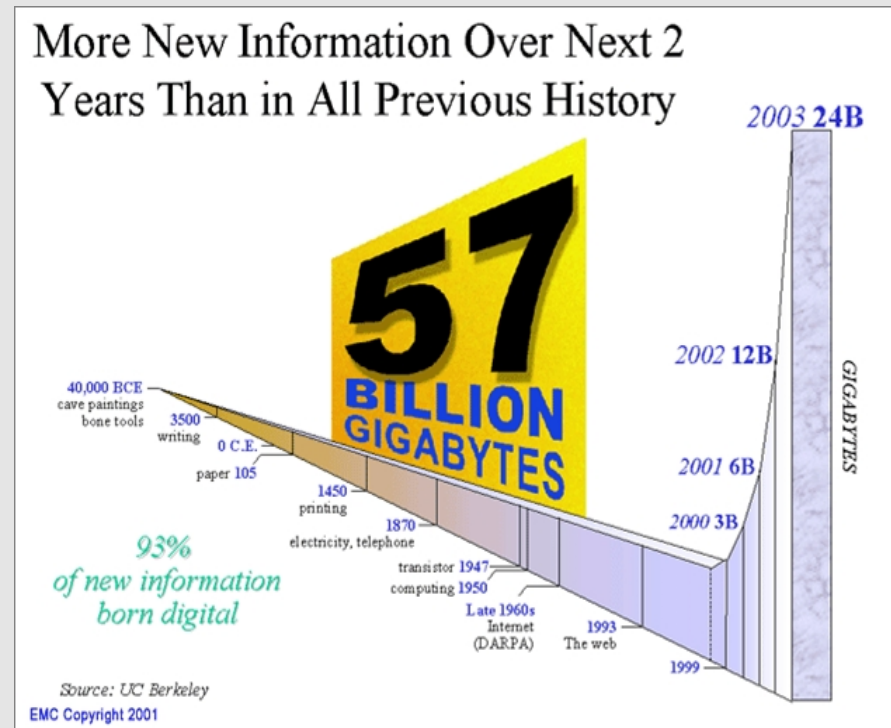
Massive data sets are becoming more and more common

- AT&T: data base of phone calls: 20 TB
- Wal-Mart: data base of buying patterns: 70 TB
- geographic data: NASA satellites collect more than 1 TB / day



Apalachian mountains ($800 \times 800 \text{ km}^2$)

- 100 m resolution: 500 MB
- 30 m resolution: 5.5 GB
(available for 80% of earth)
- 1 m resolution: 5 TB



UC Berkeley study (2000)

Two sorting algorithms: *InsertionSort* and *MergeSort**InsertionSort* (A)

1. **for** $j := 2$ **to** N
2. **do** $key := A[j]; i := j - 1$
3. **while** $i > 0$ **and** $A[i] > key$
4. **do** $A[i + 1] := A[i]; i := i - 1$
5. $A[i + 1] := key$

MergeSort (A, p, r)

1. **if** $p < r$
2. **then** $q := (p + r)/2$
3. *MergeSort* (A, p, q)
4. *MergeSort* ($A, q + 1, r$)
5. *Merge* (A, p, q, r)

Merge (A, p, q, r)

1. $N_1 := q - p + 1; N_2 := r - q$
2. Create arrays $L[1..N_1 + 1]$ and $R[1..N_2 + 1]$
3. **for** $i := 1$ **to** N_1 **do** $L[i] := A[p + i - 1]$
4. **for** $i := 1$ **to** N_2 **do** $R[i] := A[q + i]$
5. $L[N_1 + 1] := \infty; R[N_2 + 1] := \infty$
6. **for** $k := p$ **to** r
7. **do** **if** $L[i] \leq R[j]$
8. **then** $A[k] := L[i]; i := i + 1$
9. **else** $A[k] := R[j]; j := j + 1$

Which one is faster ?

running time depends on input size



analyze running time as a function of the input size

running time depends on input size



analyze running time as a function of the input size

running time (sec)



500,000

1,000,000

input size



running time depends on input size



analyze running time as a function of the input size

running time (sec)



500,000

1,000,000

input size

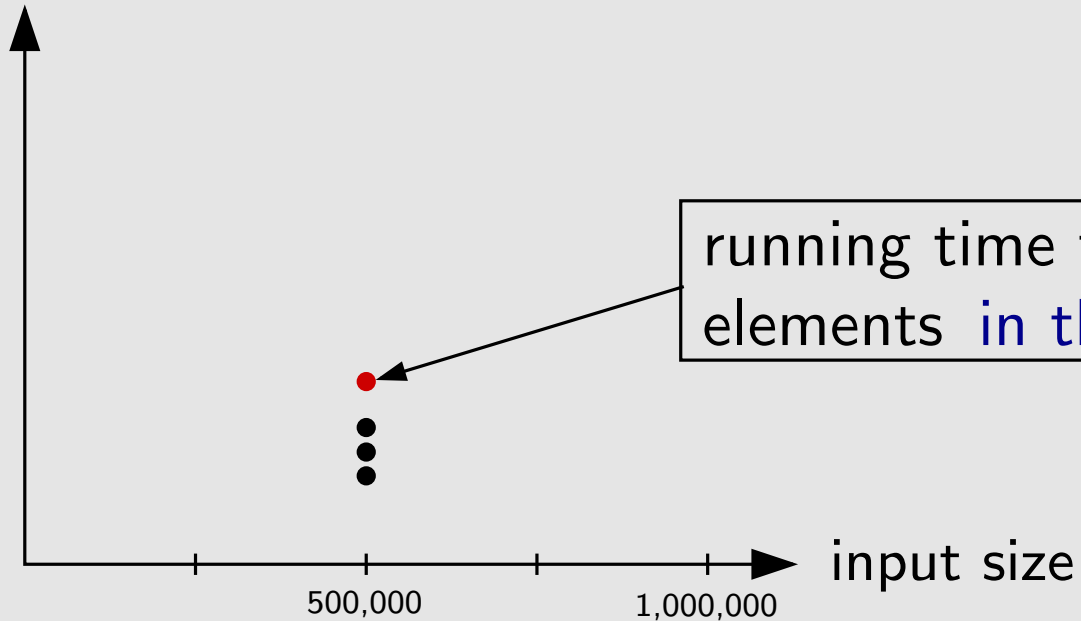


running time depends on input size



analyze running time as a function of the input size

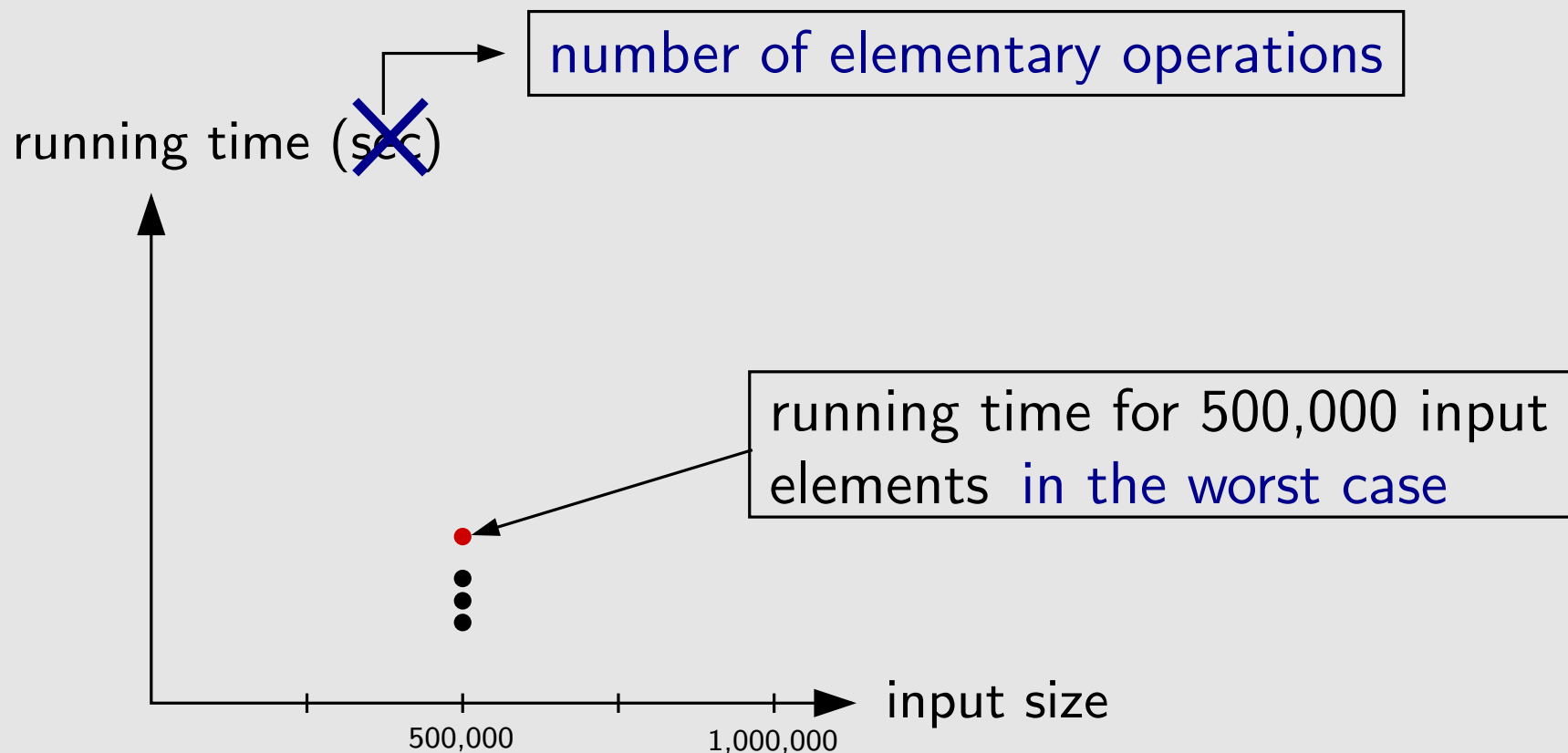
running time (sec)



running time depends on input size



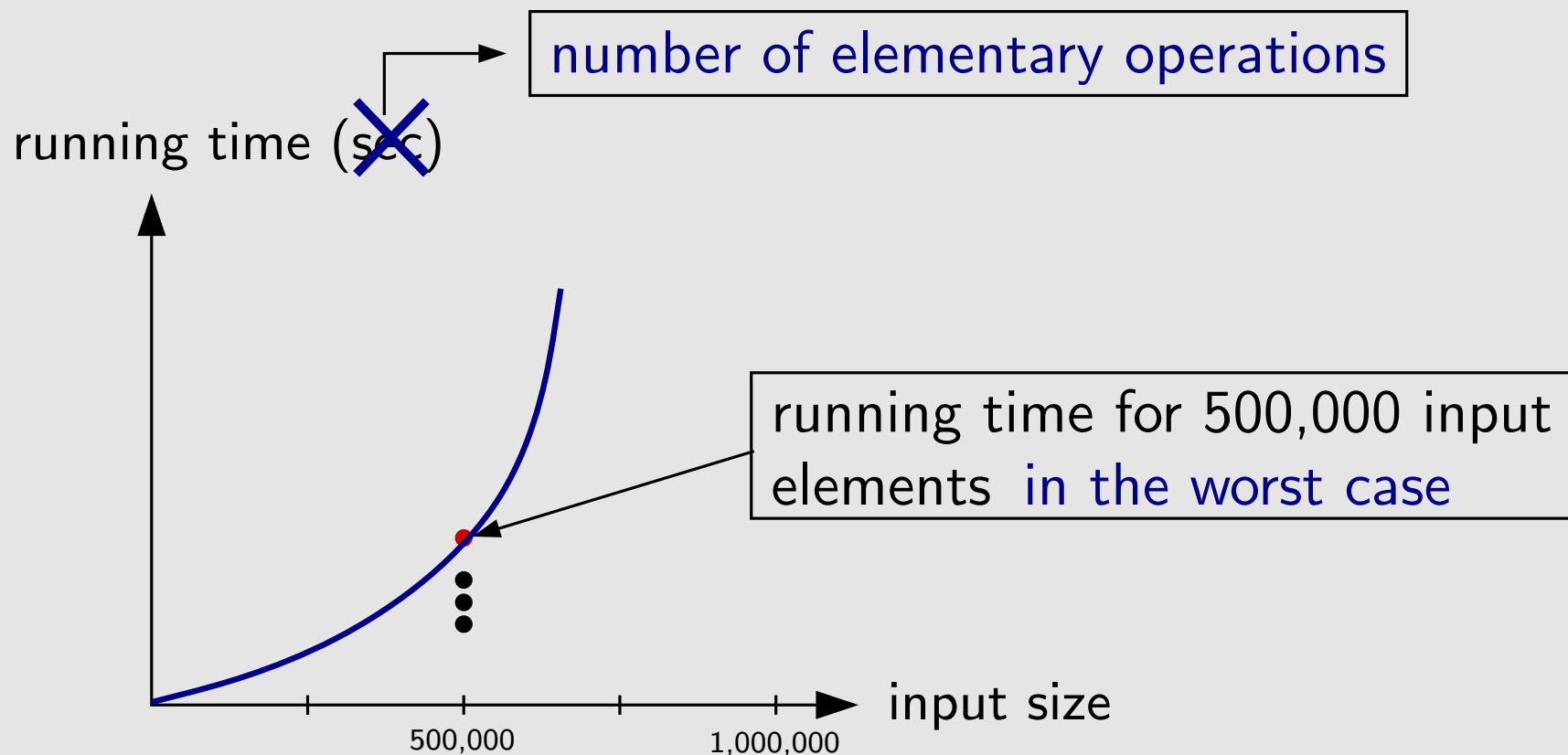
analyze running time as a function of the input size



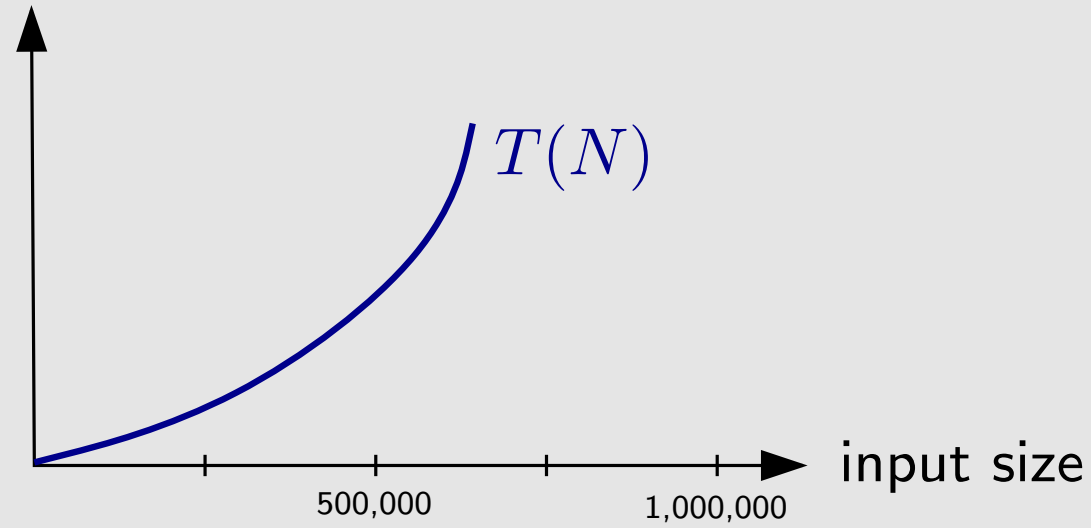
running time depends on input size



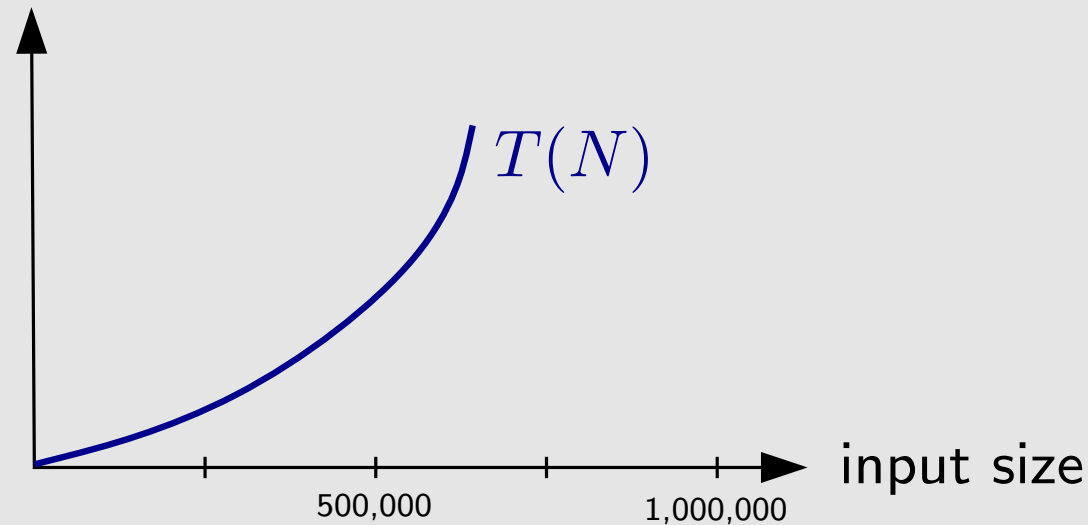
analyze running time as a function of the input size



running time (# elementary operations)



running time (# elementary operations)



$T(N) = \# \text{ elementary operations}$ the algorithm performs in the worst case as function of N , the number of input elements

we analyze asymptotic behavior of $T(N)$: is it $O(N)$, $O(N^2)$, etc.

\implies relevant for large data sets!

InsertionSort (A)

1. **for** $j := 2$ **to** N
2. **do** $key := A[j]; i := j - 1$
3. **while** $i > 0$ **and** $A[i] > key$
4. **do** $A[i + 1] := A[i]; i := i - 1$
5. $A[i + 1] := key$

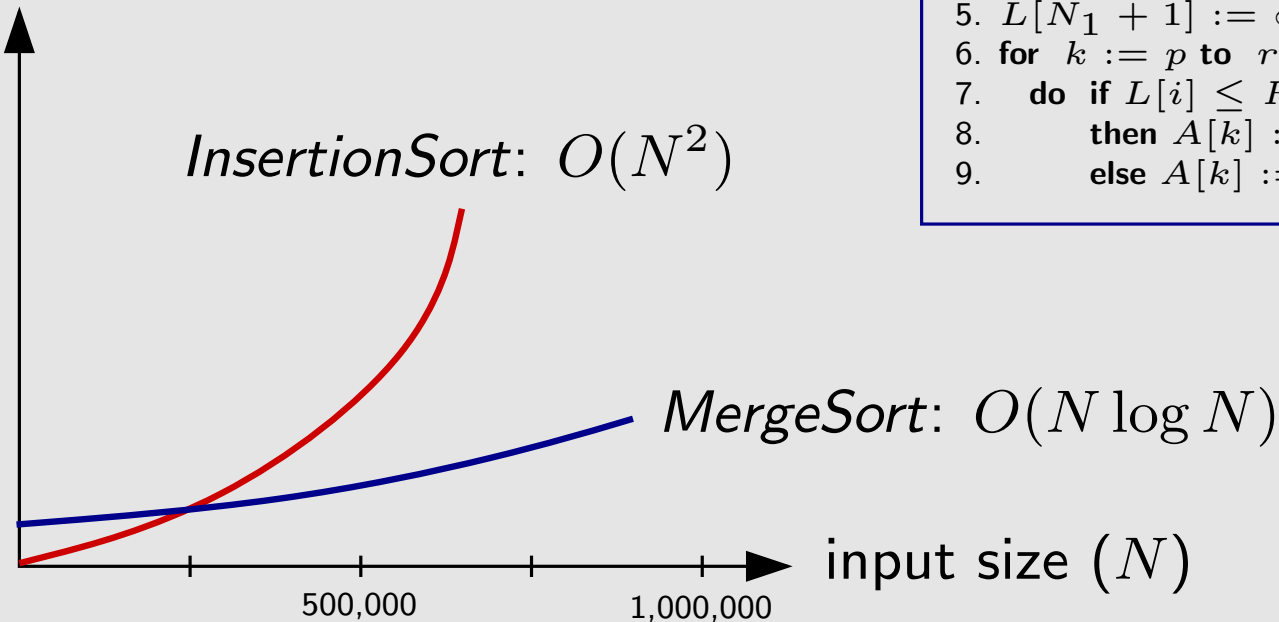
MergeSort (A, p, r)

1. **if** $p < r$
2. **then** $q := (p + r)/2$
3. *MergeSort* (A, p, q)
4. *MergeSort* ($A, q + 1, r$)
5. *Merge* (A, p, q, r)

Merge (A, p, q, r)

1. $N_1 := q - p + 1; N_2 := r - q$
2. Create arrays $L[1..N_1 + 1]$ and $R[1..N_2 + 1]$
3. **for** $i := 1$ **to** N_1 **do** $L[i] := A[p + i - 1]$
4. **for** $i := 1$ **to** N_2 **do** $R[i] := A[q + i]$
5. $L[N_1 + 1] := \infty; R[N_2 + 1] := \infty$
6. **for** $k := p$ **to** r
7. **do** **if** $L[i] \leq R[j]$
8. **then** $A[k] := L[i]; i := i + 1$
9. **else** $A[k] := R[j]; j := j + 1$

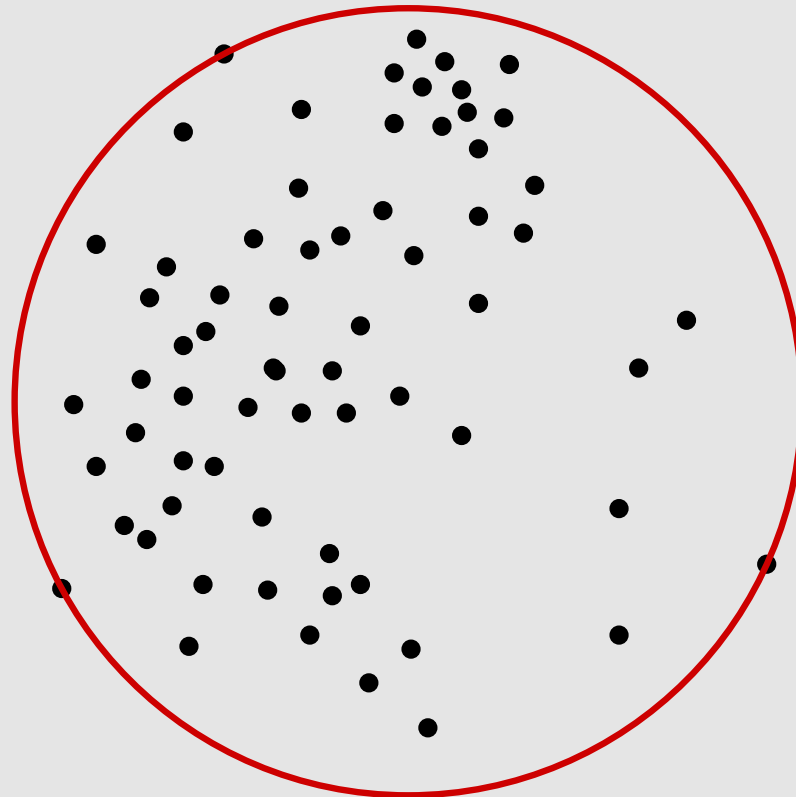
elementary operations



Compute smallest enclosing disk of set P of N points in the plane.



Compute smallest enclosing disk of set P of N points in the plane.



SmallestDisk(P)

1. *RandomPermute*(P)
2. $D :=$ smallest disk for $P[1], P[2], P[3]$
3. **for** $i := 4$ **to** N
4. **do if** $P[i] \in D$
5. **then skip**
6. **else** $D :=$ smallest disk for $\{P[1], \dots, P[i]\}$
 where $P[i]$ is on the boundary

SmallestDisk(P)

1. *RandomPermute*(P)

2. $D :=$ smallest disk for $P[1], P[2], P[3]$

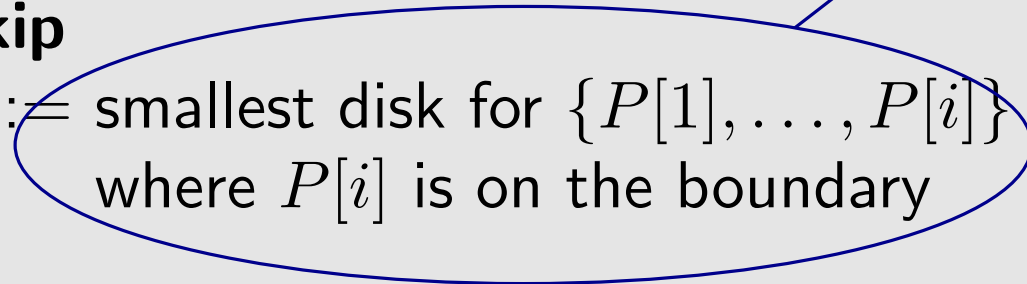
3. **for** $i := 4$ **to** N

4. **do if** $P[i] \in D$

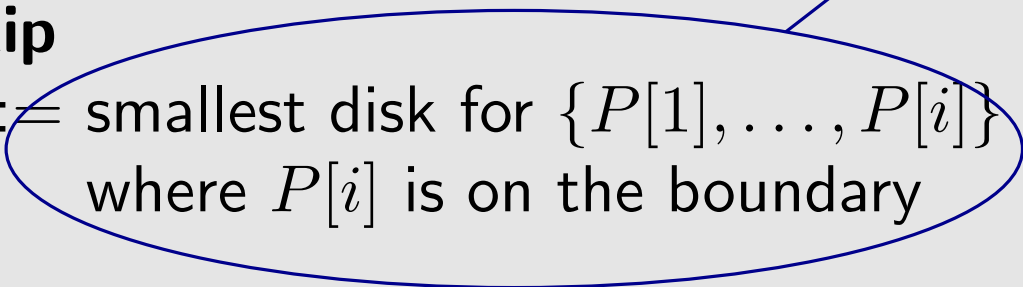
5. **then skip**

6. **else** $D :=$ smallest disk for $\{P[1], \dots, P[i]\}$
 where $P[i]$ is on the boundary

"recursive" call



SmallestDisk(P)

1. *RandomPermute*(P)
 2. $D :=$ smallest disk for $P[1], P[2], P[3]$
 3. **for** $i := 4$ **to** N
 4. **do** **if** $P[i] \in D$
 5. **then skip**
 6. **else** $D :=$ smallest disk for $\{P[1], \dots, P[i]\}$
 where $P[i]$ is on the boundary
- "recursive" call
- 

RandomPermute(P)

1. **for** $i := 1$ **to** $N - 1$
2. **do** $r :=$ random integer in range $i \dots N$
3. swap $P[i]$ and $P[r]$

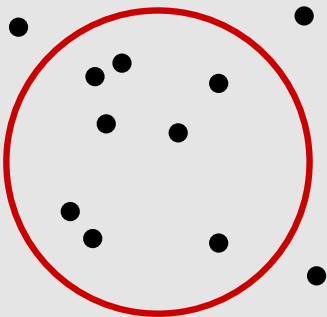
RandomPermute(P)

1. **for** $i := 1$ **to** $N - 1$
2. **do** $r :=$ random integer in range $i \dots N$
3. swap $P[i]$ and $P[r]$

running time is $O(N)$

SmallestDisk(P)

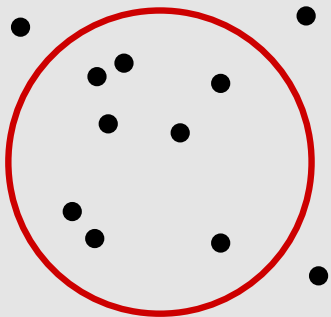
1. *RandomPermute*(P)
2. $D :=$ smallest disk for $P[1], P[2], P[3]$
3. **for** $i := 4$ **to** N
4. **do** **if** $P[i] \in D$
5. **then skip**
6. **else** $D :=$ smallest disk for $\{P[1], \dots, P[i]\}$
 where $P[i]$ is on the boundary



- $P[1]$ t/m $P[i]$
- $P[i + 1]$ t/m $P[N]$

SmallestDisk(P)

1. *RandomPermute*(P)
2. $D :=$ smallest disk for $P[1], P[2], P[3]$
3. **for** $i := 4$ **to** N
4. **do** **if** $P[i] \in D$
5. **then skip**
6. **else** $D :=$ smallest disk for $\{P[1], \dots, P[i]\}$
 where $P[i]$ is on the boundary

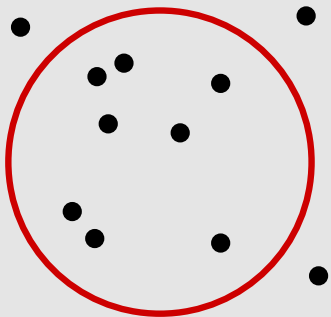


$$\Pr [P[i] \notin D] \leq 3/i$$

- $P[1]$ t/m $P[i]$
- $P[i + 1]$ t/m $P[N]$

SmallestDisk(P)

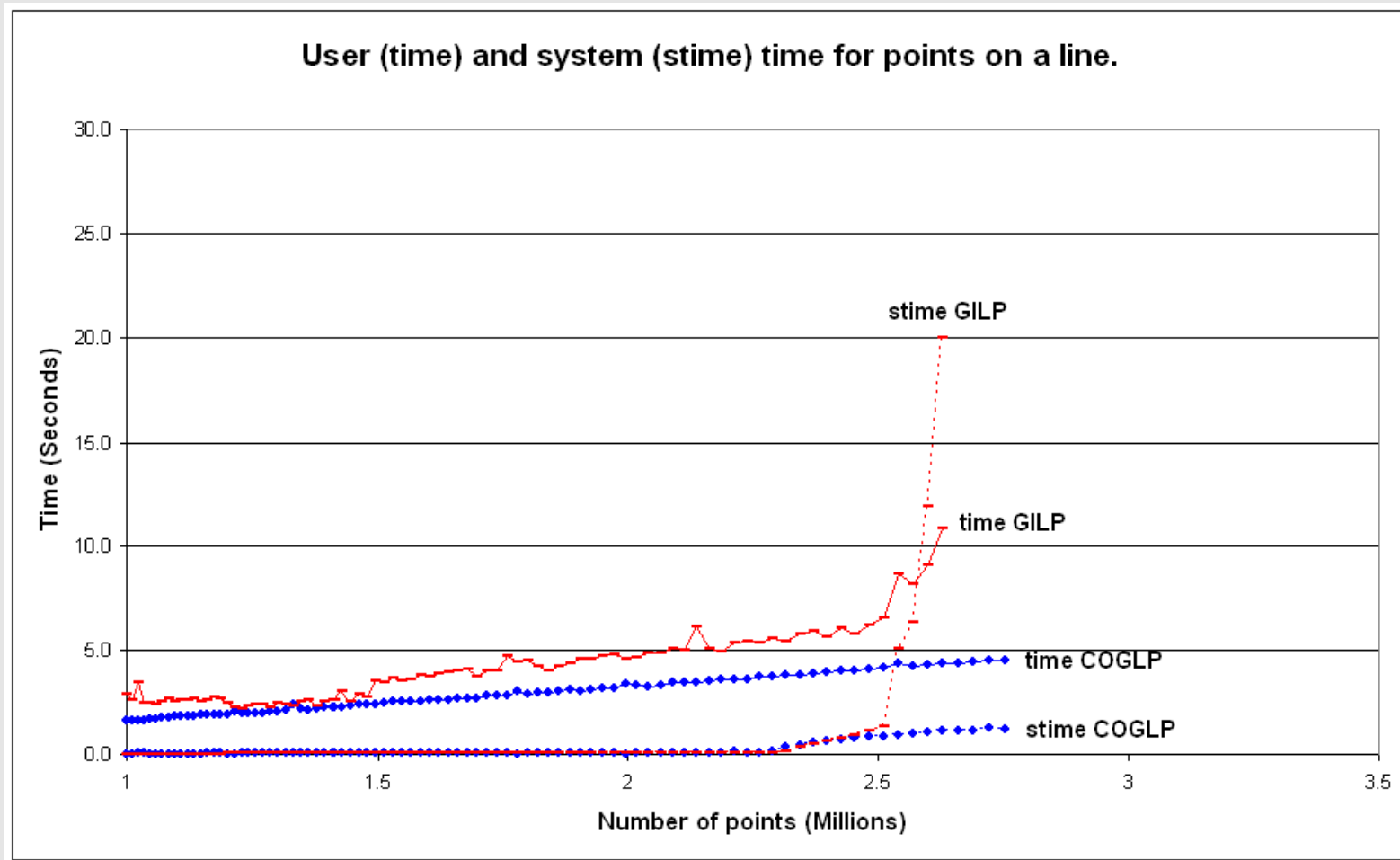
1. *RandomPermute*(P)
2. $D :=$ smallest disk for $P[1], P[2], P[3]$
3. **for** $i := 4$ **to** N
4. **do** **if** $P[i] \in D$
5. **then skip**
6. **else** $D :=$ smallest disk for $\{P[1], \dots, P[i]\}$
 where $P[i]$ is on the boundary



$$\Pr [P[i] \notin D] \leq 3/i$$

- $P[1]$ t/m $P[i]$
- $P[i + 1]$ t/m $P[N]$

\implies expected running time is $O(N)$



Pentium 4, 2.60GHz

≈ 89 MB main memory available to the program

$T(n)$ = # elementary operations the algorithm performs in the worst case as function of N , the number of input elements

additions, multiplications, comparisons, reading a value from memory, etc.

Hmmm ... is this justified?

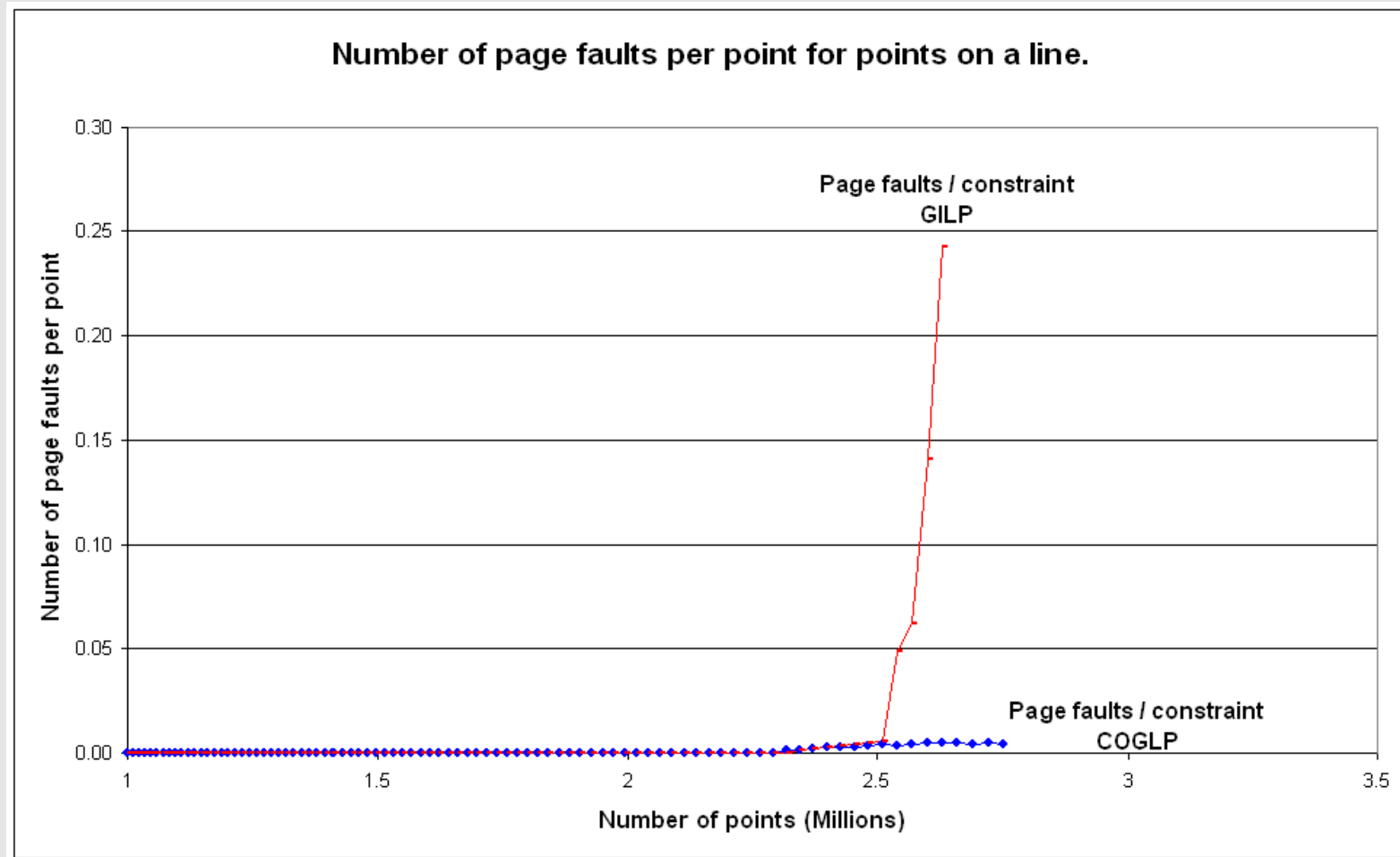
$T(n)$ = # elementary operations the algorithm performs in the worst case as function of N , the number of input elements

additions, multiplications, comparisons, reading a value from memory, etc.

Hmmm ... is this justified?

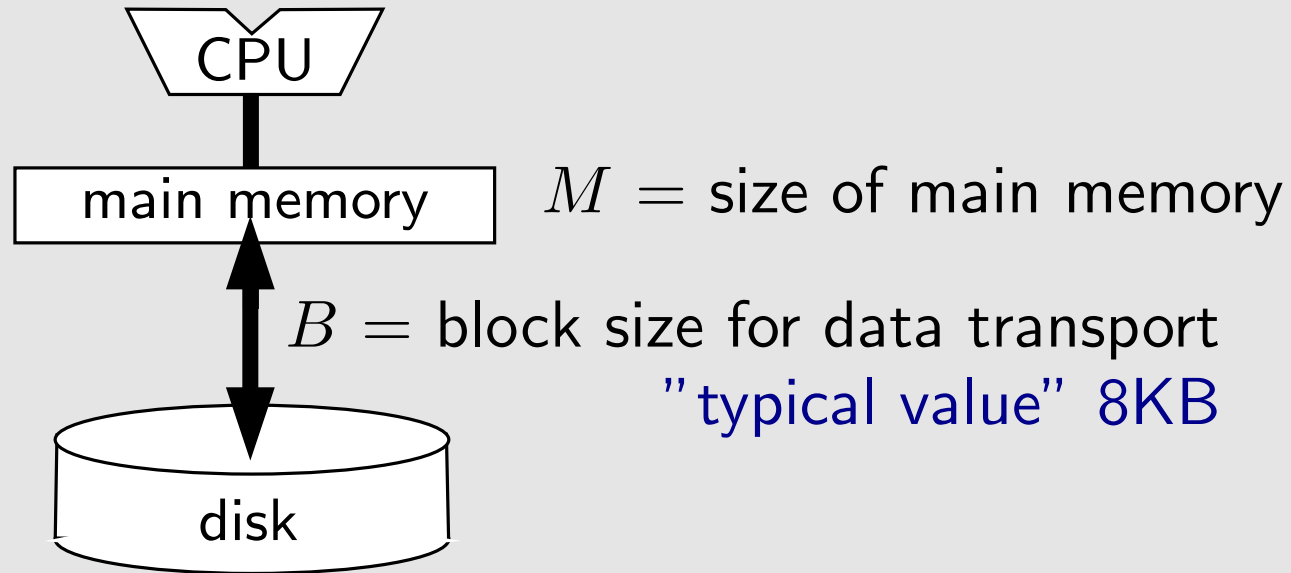
NO!

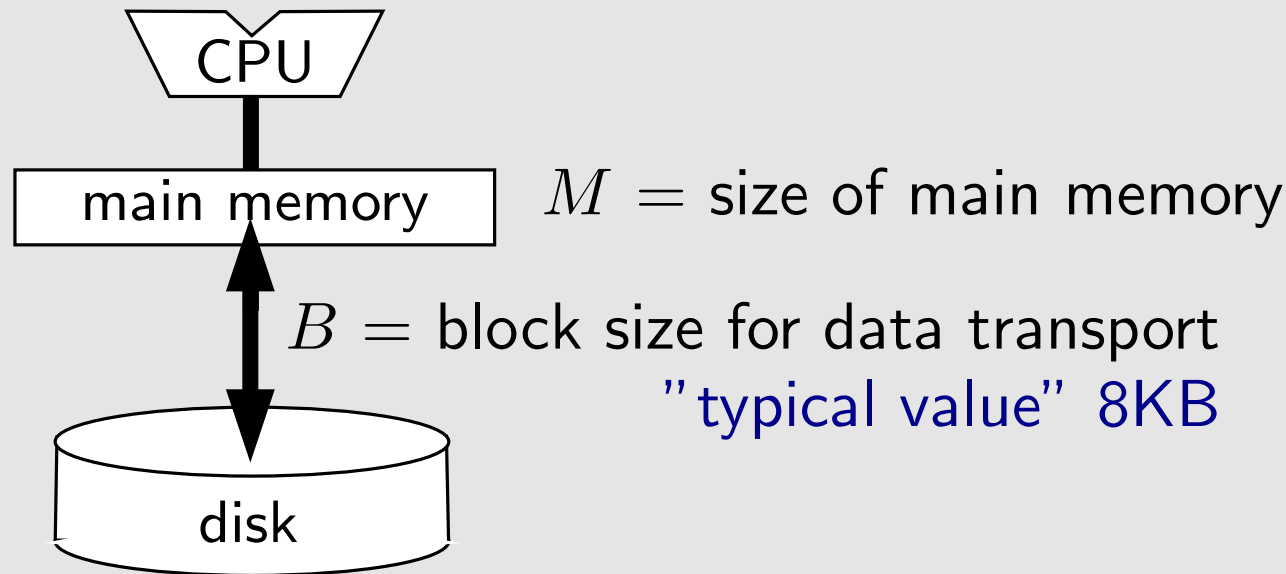
operations on data in main memory: tens of nanoseconds
disk operations: several milliseconds



Pentium 4, 2.60GHz

≈ 89 MB main memory available to the program





- let algorithm handle data placement and transport
 - which data are placed together in a block
 - which blocks are kept in main memory
- analyze number of disk operations

RandomPermute(P)

1. **for** $i := 1$ **to** $N - 1$
2. **do** $r :=$ random integer in range $i \dots N$
3. swap $P[i]$ and $P[r]$

analysis of (expected) number of disk operations

- $N \leq M$: 0

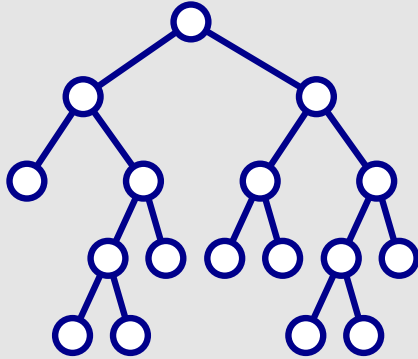
0 disk operations

- $N > M$:

$(N - 1) \cdot (1 - \frac{M}{N})$ disk operations

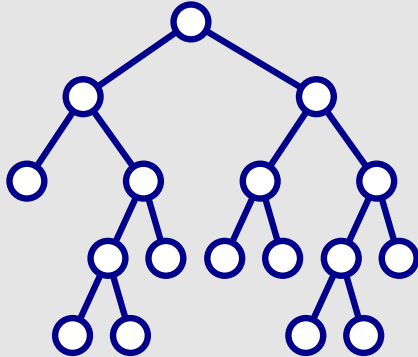
(e.g. $(N - 1)/2$) disk operations when $N = 2M$)

binary search tree: search structure for internal memory

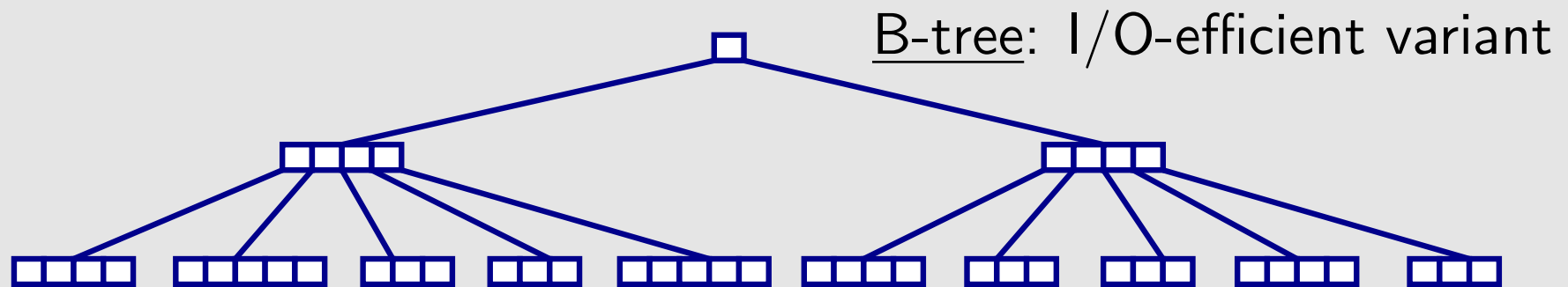


- nodes contain one key, have degree 2
- depth is $O(\log N)$

binary search tree: search structure for internal memory

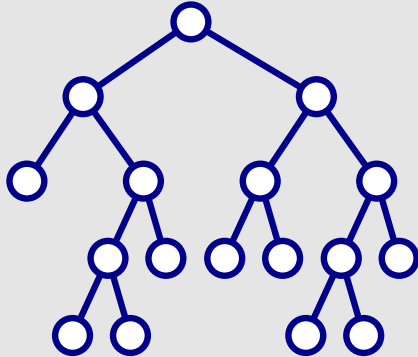


- nodes contain one key, have degree 2
- depth is $O(\log N)$

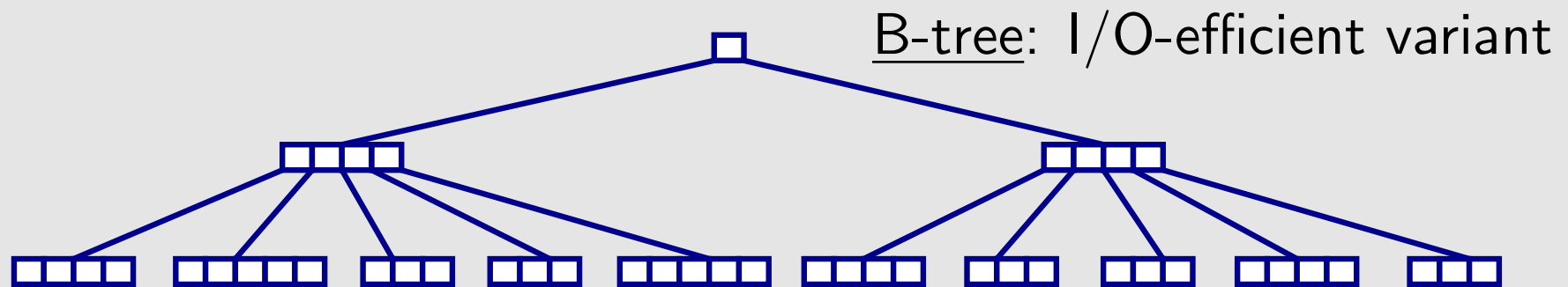


- nodes contain many keys, have high degree
- put each node into one block on disk
- depth is $O(\log N / \log B)$

binary search tree: search structure for internal memory

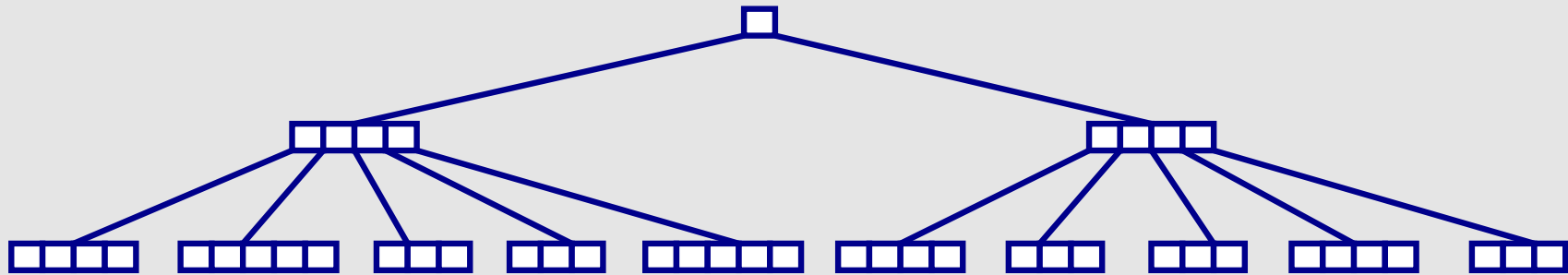


- nodes contain one key, have degree 2
- depth is $O(\log N)$



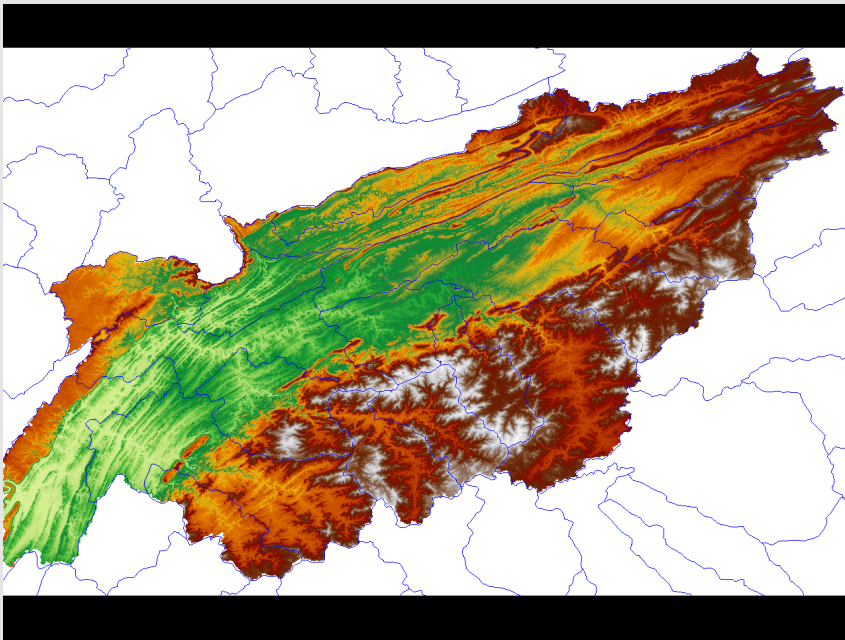
- nodes contain many keys, have high degree
- put each node into one block on disk
- depth is $O(\log N / \log B)$ in practice, degree is 250 – 2000 and depth is at most 4

R-tree: search structure (index) for spatial data

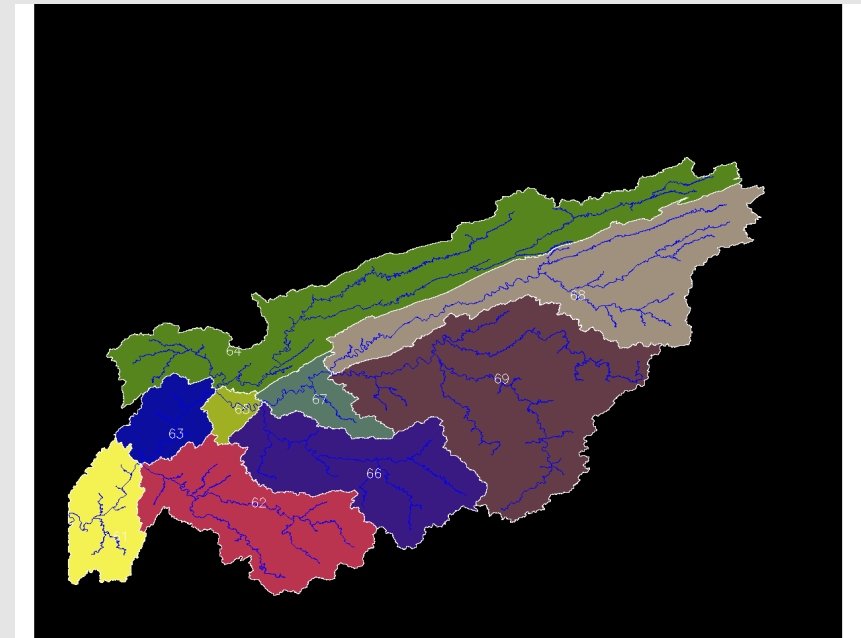


I/O-efficient variant of bounding-volume hierarchy:

- nodes contain many bounding boxes, have high degree
- put each node into one block on disk
- depth is $O(\log N / \log B)$

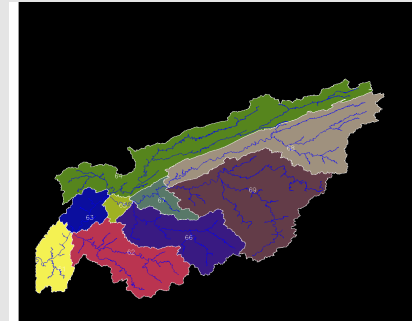
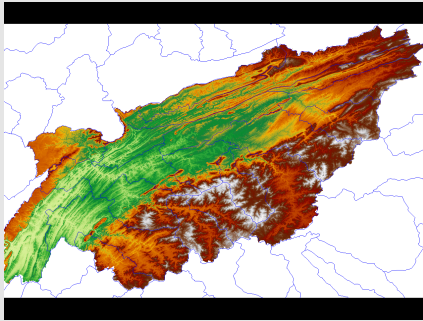


digital elevation model (DEM)



- flow routing, flow accumulation
- watersheds, Pfaffstetter labeling

TerraFlow: P. Agarwal, L. Arge, J. Chase, P. Halpin, L. Toma, D. Urban, J. Vitter, R. Wickremesinghe. Pfaffstetter: + H. Haverkort

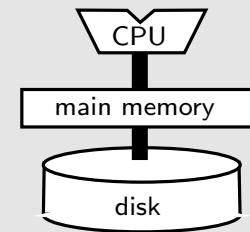


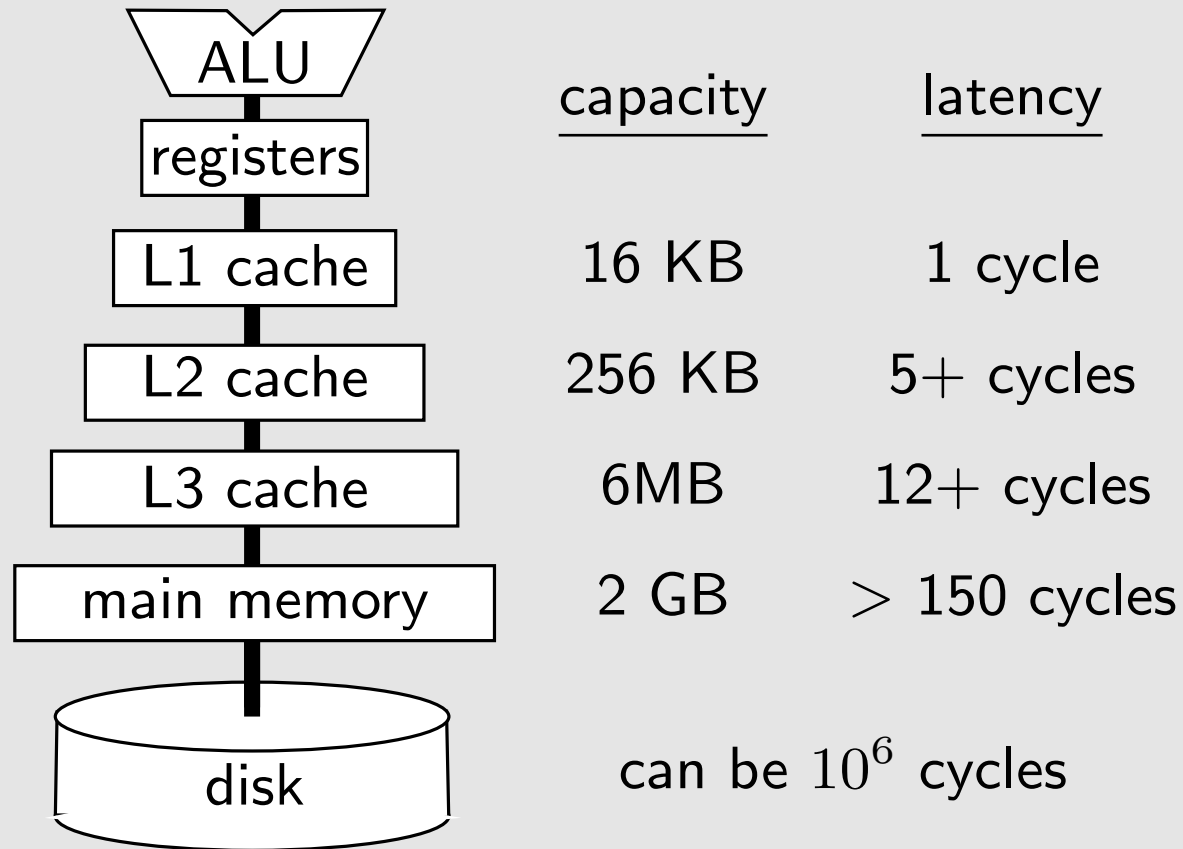
comparison with existing commercial software (ArcInfo) and open source software (GRASS):

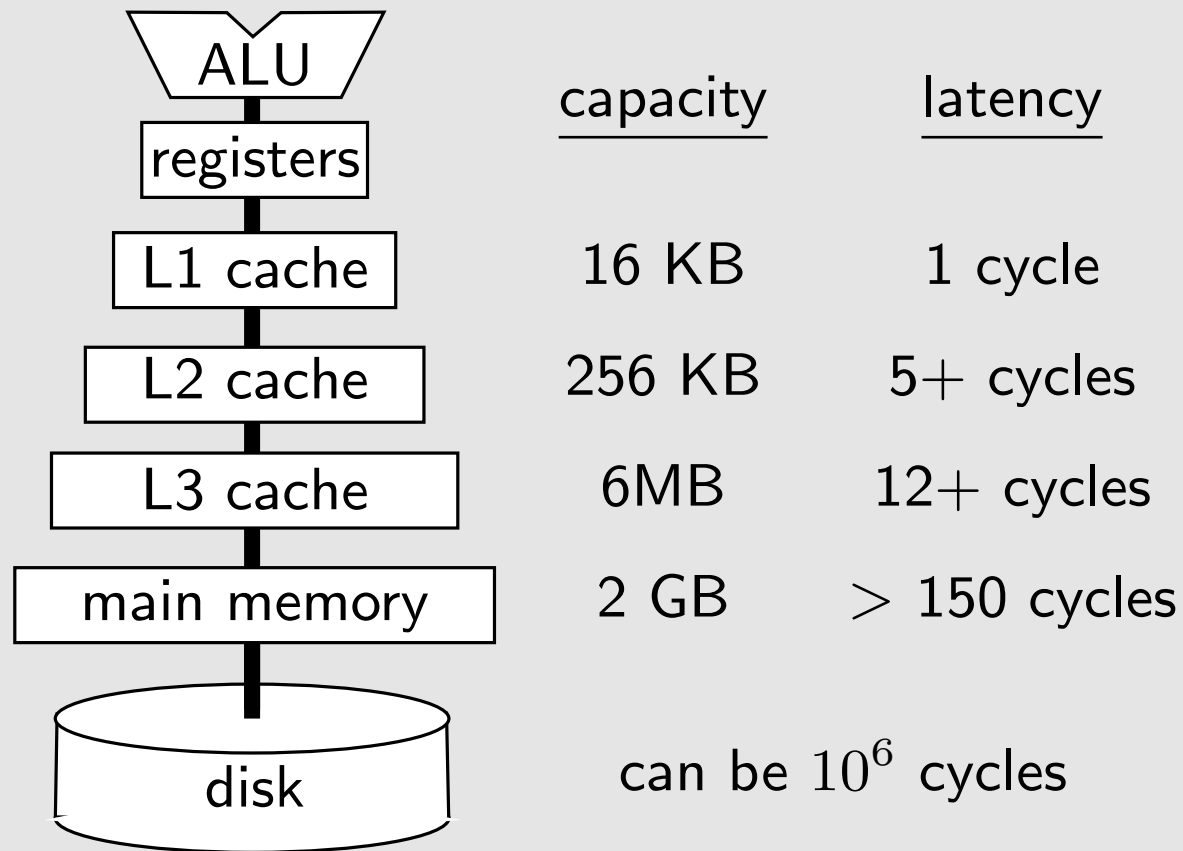
speed-up factor between 2 and 1000

implemented using TPIE, a library for I/O-efficient algorithms

- M and B depend on platform
- even on fixed machine values of M and B may vary
 - main memory may have to be shared with other processes
 - disk-cache "changes" block size
- two-level I/O-model too simplistic



Intel Itanium2 memory hierarchy

Intel Itanium2 memory hierarchy

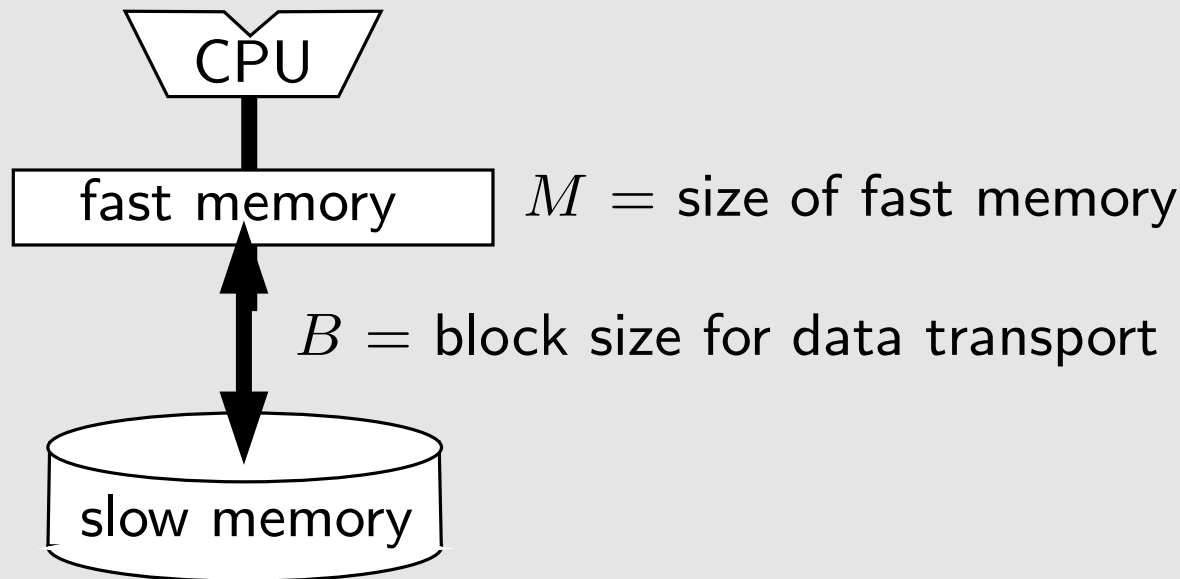
Conclusion: caching behavior can also make a large difference.

Ideal: algorithm that is efficient w.r.t. disk and all cache levels

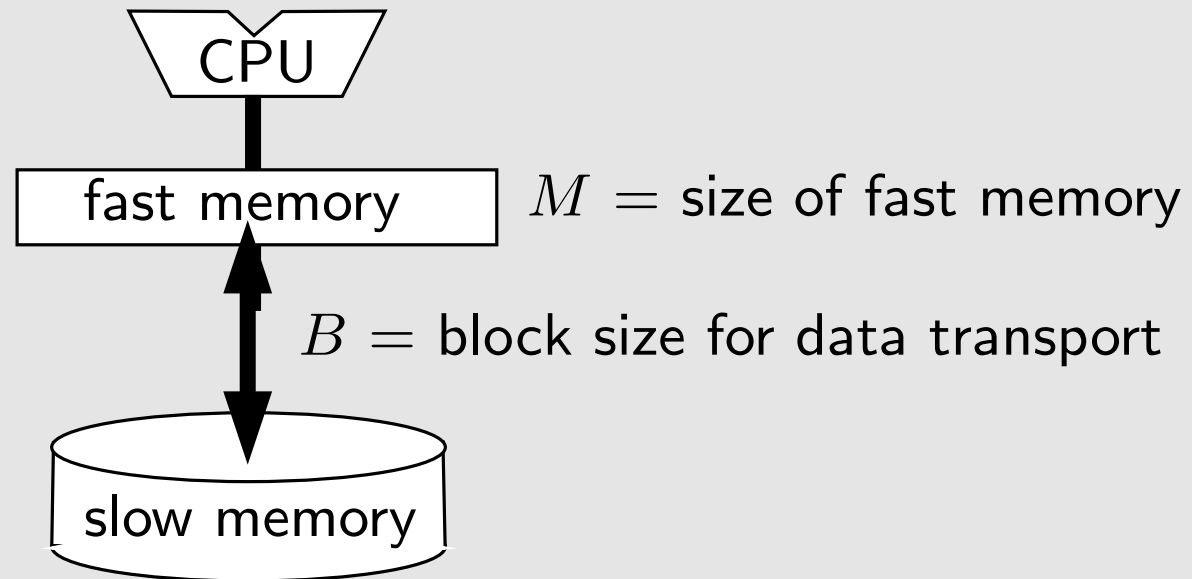
- caches are not under control of algorithm
- algorithms taking all cache-levels into account quite complicated

So what can we do ??

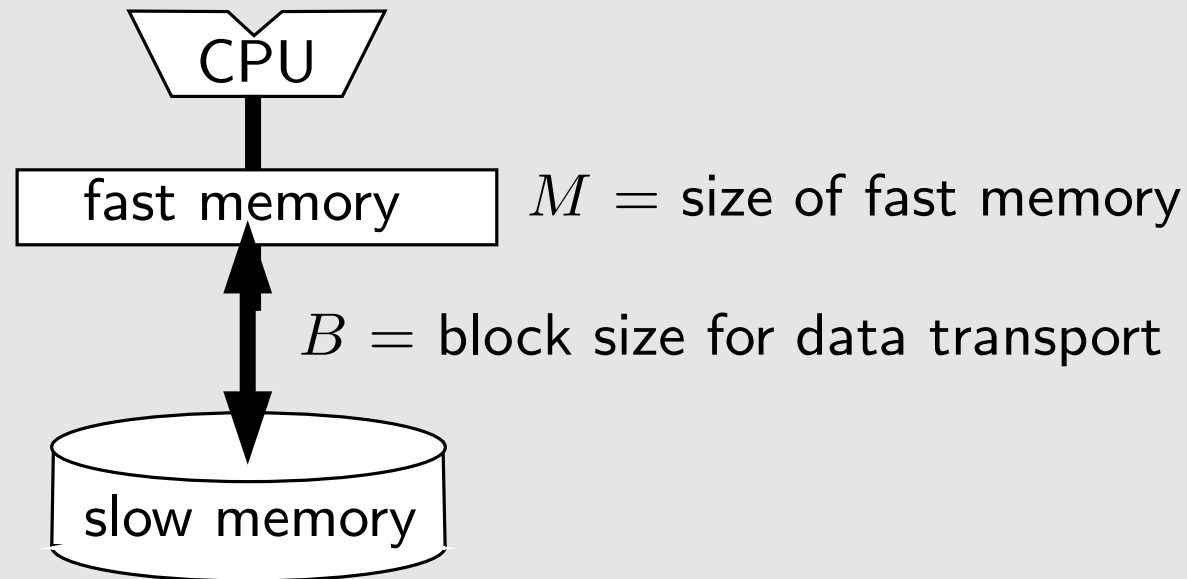
Algorithm designed for simple two-level memory model



But: algorithm is not allowed to use the value of B and M !

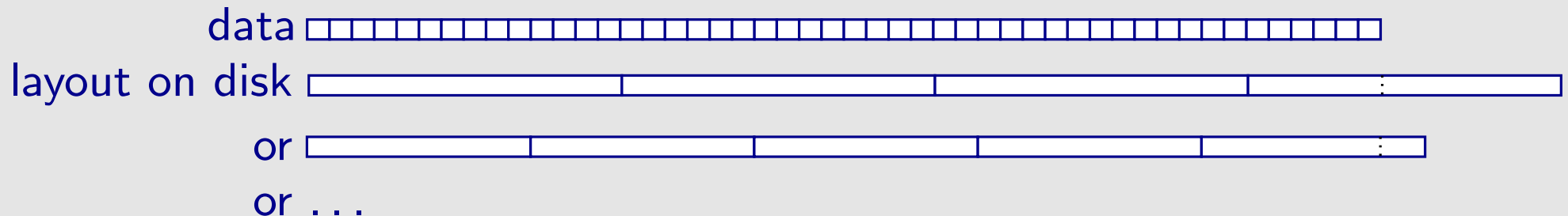


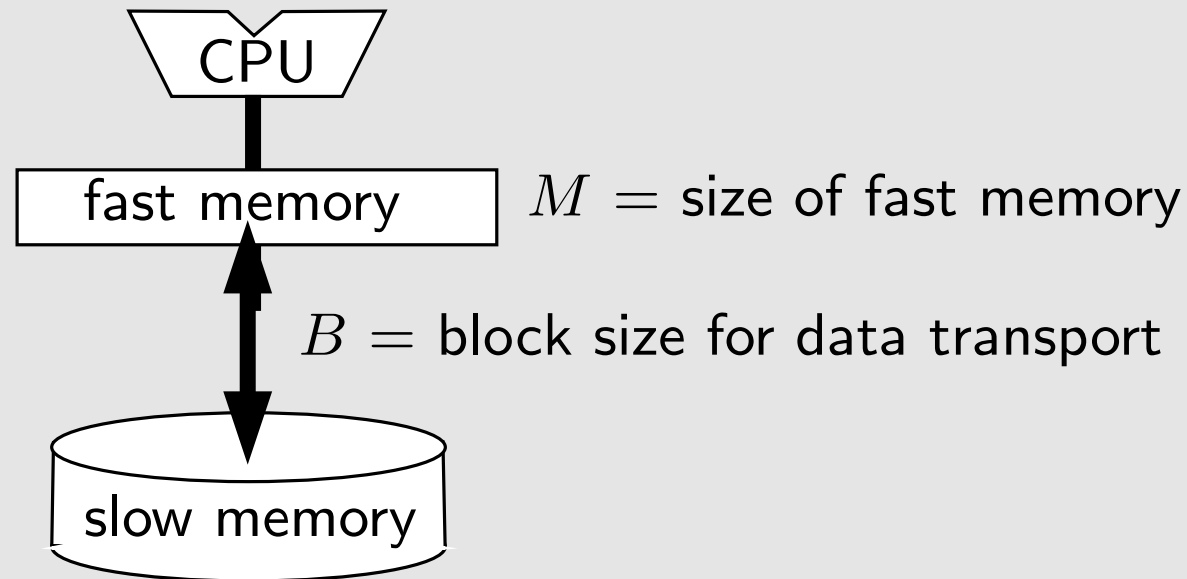
Assumptions:



Assumptions:

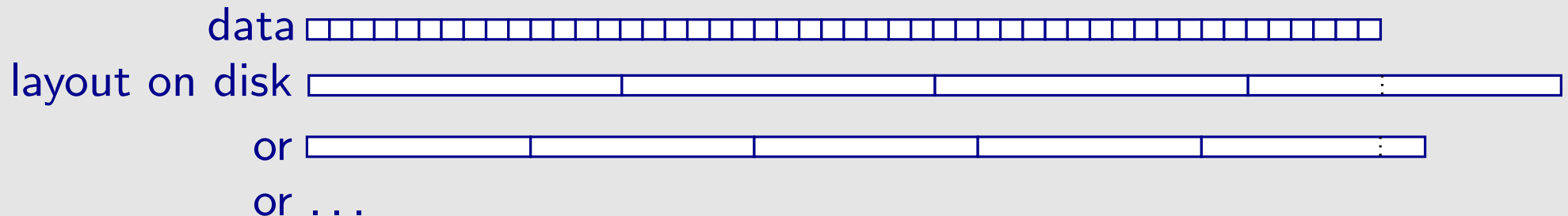
- Data written to "disk" consecutively is stored consecutively on disk





Assumptions:

- Data written to "disk" consecutively is stored consecutively on disk



- Operating system uses optimal replacement strategy

CacheObliviousSmallestDisk(P)

1. **if** ($\#$ points in P) ≤ 3
2. **then return** smallest disk for P
3. **else** $(P_1, P_2) := \text{RandomSplit}(P)$
4. $D := \text{CacheObliviousSmallestDisk}(P_1)$
5. **for** all $P[i] \in G_2$
6. **do if** $P[i] \notin D$
7. **then** $D_i :=$ smallest disk for P with $P[i]$ on its boundary
8. **return** best of all computed disks

CacheObliviousSmallestDisk(P)

1. **if** ($\#$ points in P) ≤ 3
2. **then return** smallest disk for P
3. **else** (P_1, P_2) := *RandomSplit* (P)
4. $D := \text{CacheObliviousSmallestDisk}(P_1)$
5. **for** all $P[i] \in G_2$
6. **do if** $P[i] \notin D$
7. **then** $D_i :=$ smallest disk for P with $P[i]$ on its boundary
8. **return** best of all computed disks

RandomSplit(P)

1. **for** $i := 1$ **to** N
2. **do** $r :=$ random number in range $[0, 1]$
3. **if** $r < 1/2$
4. **then** Put $P[i]$ into P_1
5. **else** Put $P[i]$ into P_2
6. **return** (P_1, P_2)

old algorithm:

```
RandomPermute( $P$ )  
1. for  $i := 1$  to  $N - 1$   
2.   do  $r :=$  random integer in range  $i \dots N$   
3.     swap  $P[i]$  and  $P[r]$ 
```

$$E[\text{\#disk operations}] = (N - 1) \cdot \left(1 - \frac{M}{N}\right)$$

old algorithm:

```

RandomPermute(P)
1. for  $i := 1$  to  $N - 1$ 
2.   do  $r :=$  random integer in range  $i \dots N$ 
3.     swap  $P[i]$  and  $P[r]$ 

```


$$E[\text{\#disk operations}] = (N - 1) \cdot \left(1 - \frac{M}{N}\right)$$

new algorithm:

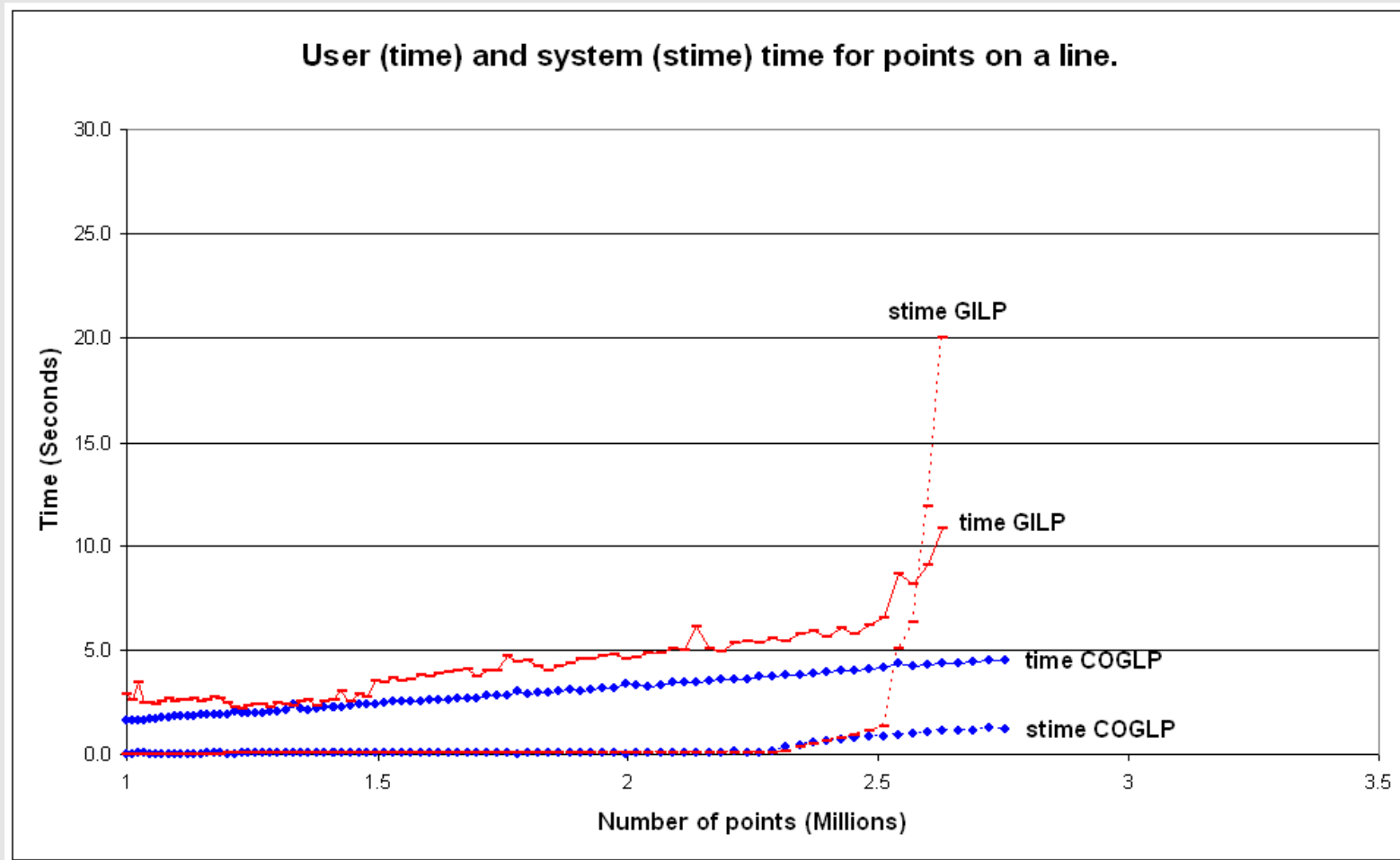
```

RandomSplit(P)
1. for  $i := 1$  to  $N$ 
2.   do  $r :=$  random number in range  $[0, 1]$ 
3.     if  $r < 1/2$  then Put  $P[i]$  into  $P_1$  else Put  $P[i]$  into  $P_2$ 
4. return  $(P_1, P_2)$ 

```

layout on disk 

$$E[\text{\#disk operations}] \leq N/B$$

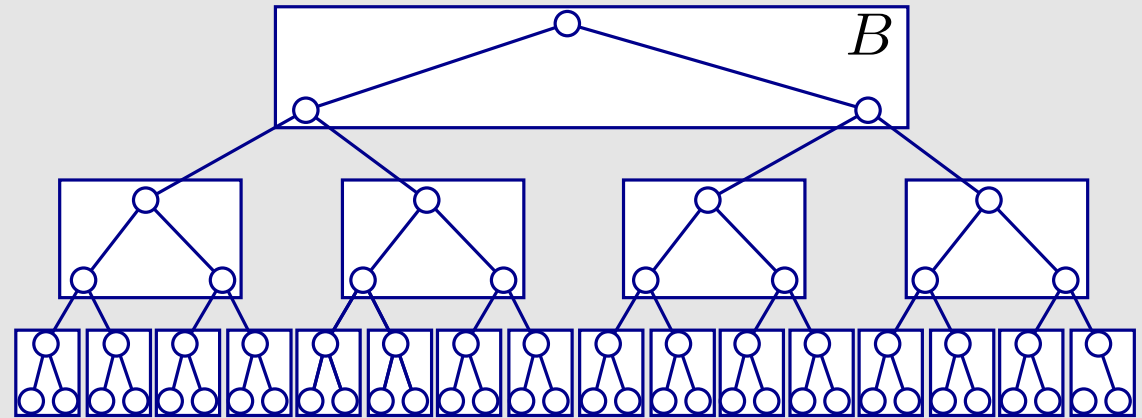


Pentium 4, 2.60GHz

≈ 89 MB main memory available to the program

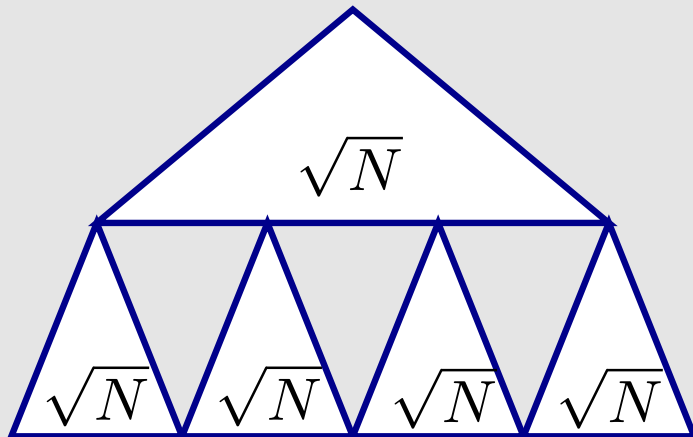
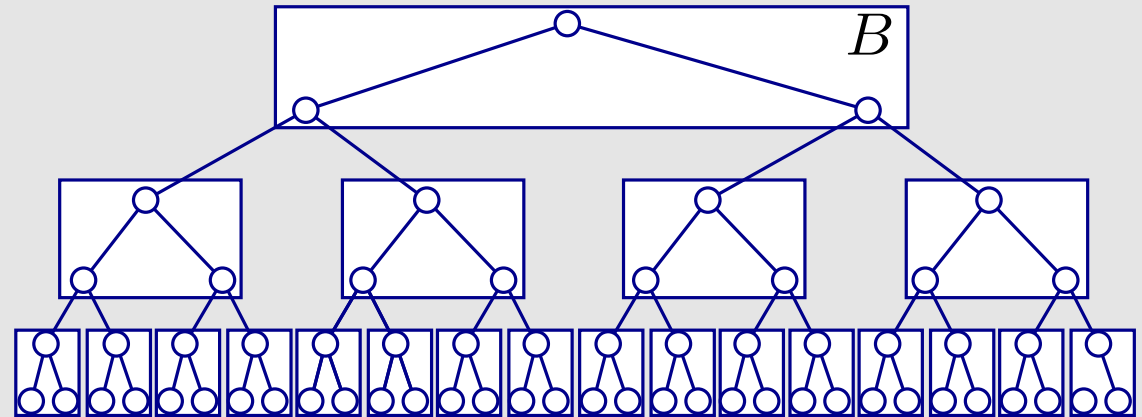
regular (cache-aware) B-tree:

- blocks: subtrees of size B



regular (cache-aware) B-tree:

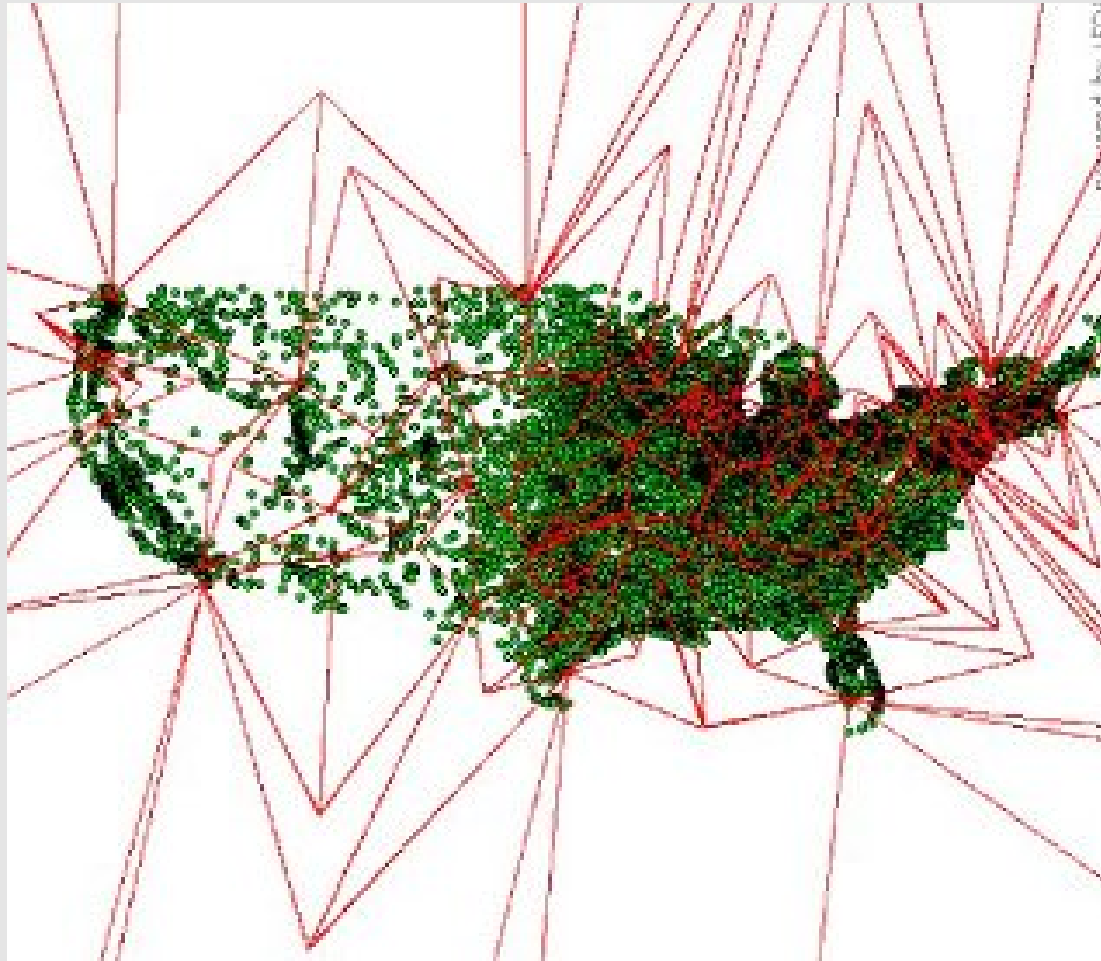
- blocks: subtrees of size B



cache-oblivious B-tree:

- cut tree into subtree at middle level; gives 1 top tree, \sqrt{N} lower trees
- first, write top to disk recursively
- next, write lower trees to disk recursively

search visits $O(\log N / \log B)$ blocks



Kumar: Voronoi diagrams of up to 300 M points

- I/O- and caching behavior crucial for massive data sets
- algorithms community is now addressing these issues
- I/O-efficient algorithms
 - have proven their value for various practical problems
 - need tuning for hardware, do not optimize caching behavior
- cache-oblivious algorithms:
 - ideal in theory: no tuning, good on all cache-levels
 - practical relevance needs further investigation