



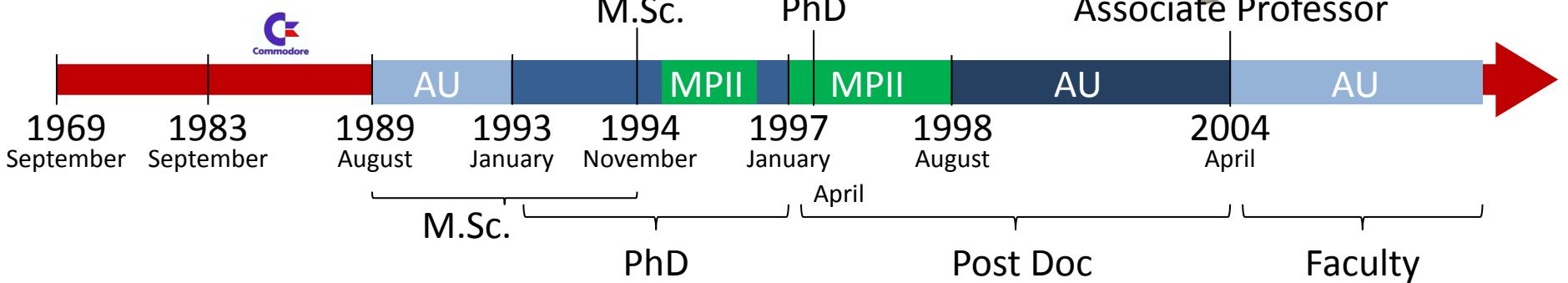
DATALOGISK INSTITUT
DET NATURVIDENSKABELIGE FAKULTET
AARHUS UNIVERSITET

Udfordringer ved håndtering af massive datamængder: Forskingen ved Grundforskningscenteret for Massive Data Algorithmics

Gerth Stølting Brodal

madalgo 
CENTER FOR MASSIVE DATA ALGORITHMIC

Gerth Stølting Brodal

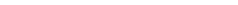


Outline of Talk

- **madalgo** :-:-:-:-
CENTER FOR MASSIVE DATA ALGORITHMIC
- who, where, what ?
- research areas
- External memory algorithmics
 - models
 - searching and sorting
- Flow simulation
- Fault-tolerant searching





madalGO 
CENTER FOR MASSIVE DATA ALGORITHMICs

– Where?





Danmarks
Grundforskningfond
Danish National
Research Foundation

- Center of
- **Lars Arge**, Professor, Centerleader
- Gerth S. Brodal, Associate Professor
- 5 Post Docs, 10 PhD students, 4 TAP
- Total budget for 5 years ca. 60 million DKK

AU



Arge



Brodal

MIT



Demaine



Indyk

MPII



Mehlhorn

Frankfurt



Meyer

Faculty

Lars Arge

Gerth Stølting Brodal

Researchers



Henrik Blunck



Brody Sandel



Nodari Sitchinava



Elad Verbin



Qin Zhang

PhD Students

Lasse Kosetski Deleuran

 Freek van Walderveen

Casper Kejlberg-Rasmussen

Kasper Dalgaard Larsen

Jesper Erenskjold Moeslund

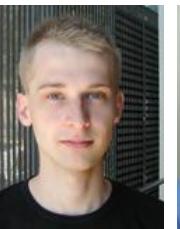
 Jakob Truelsen

 Kostas Tsakalidis

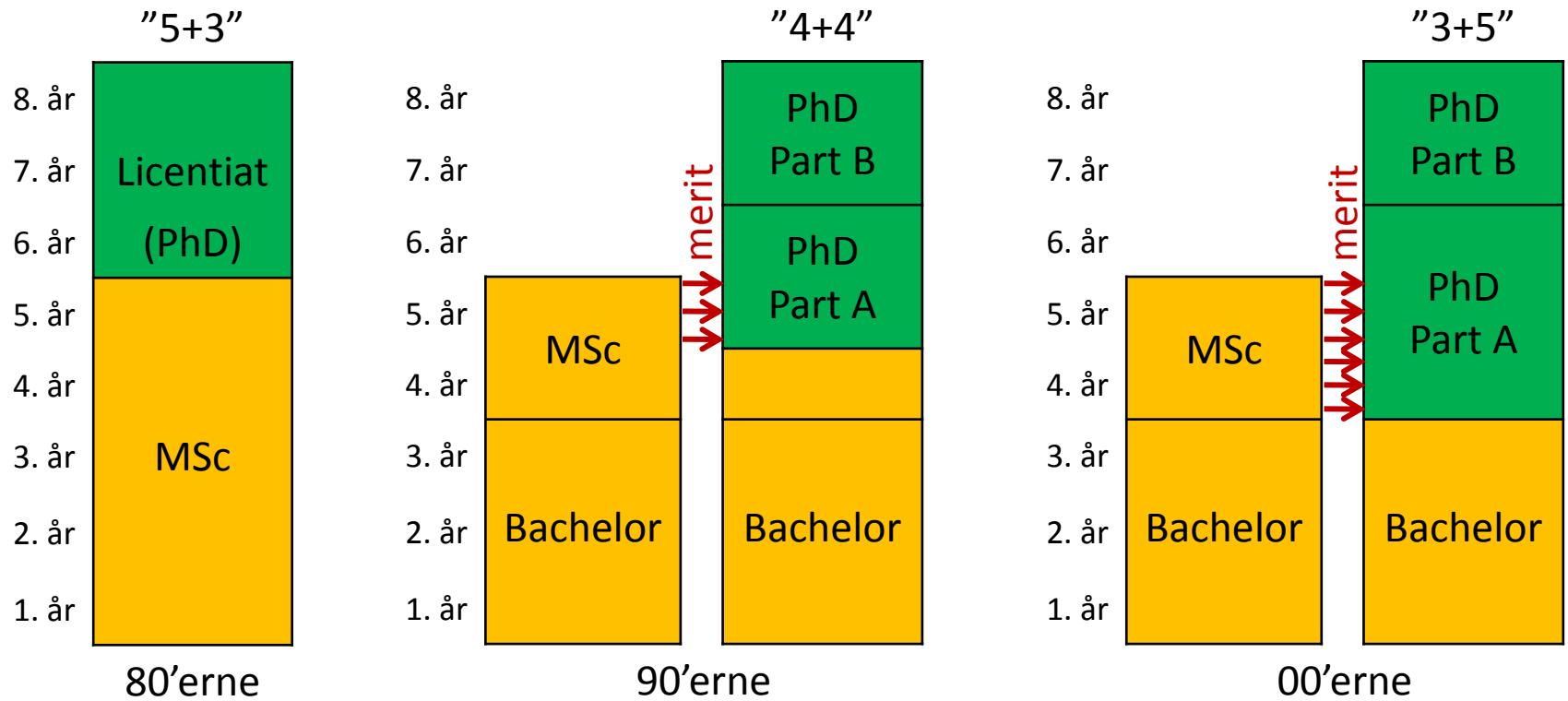
Mark Greve

Morten Revsbæk

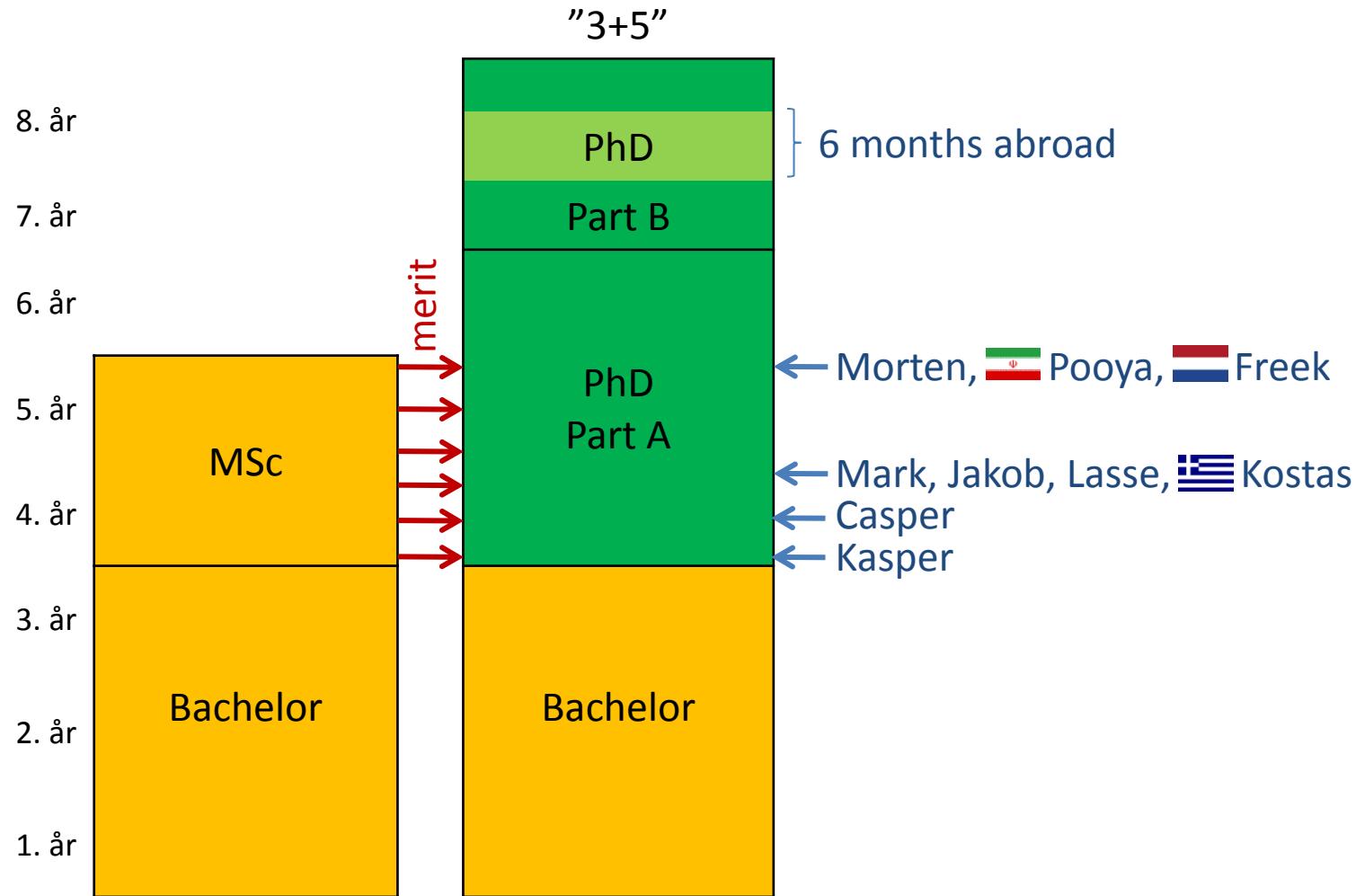
 Pooya Davoodi



PhD Education @ AU



PhD Education @ MADALGO

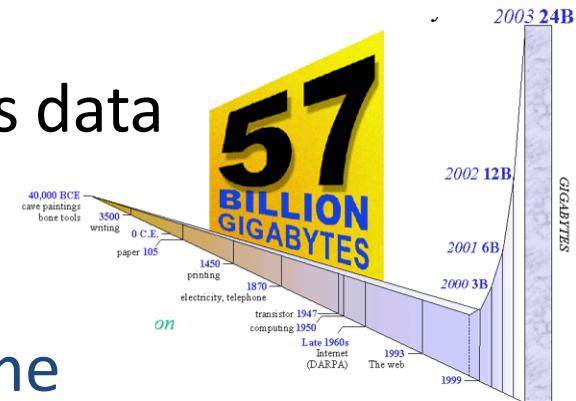




- *High level objectives*
 - Advance algorithmic knowledge in “massive data” processing area
 - Train researchers in world-leading international environment
 - Be catalyst for multidisciplinary/industry collaboration
- *Building on*
 - Strong international team
 - Vibrant international environment (focus on people)

Massive Data

- Pervasive use of computers and sensors
- Increased ability to acquire/store/process data
 - Massive data collected everywhere
- Society increasingly “data driven”
 - Access/process data anywhere any time



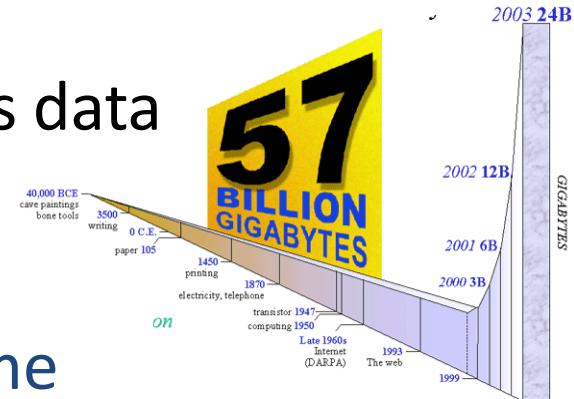
Nature special issues

- 2/06: “2020 – Future of computing”
- 9/08: “BIG DATA”
- Scientific data size growing exponentially, while quality and availability improving
- Paradigm shift: *Science will be about mining data*
 - Computer science paramount in all sciences



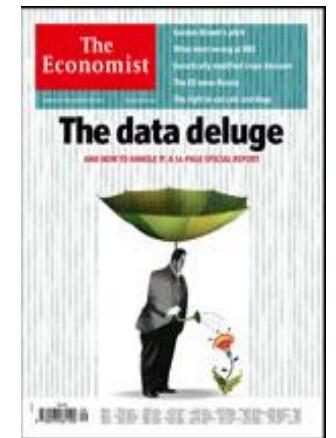
Massive Data

- Pervasive use of computers and sensors
- Increased ability to acquire/store/process data
 - Massive data collected everywhere
- Society increasingly “data driven”
 - Access/process data anywhere any time



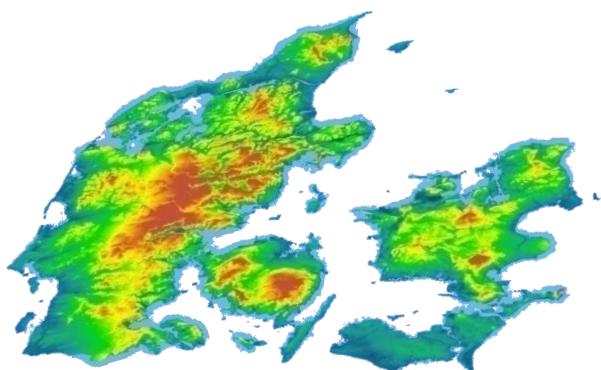
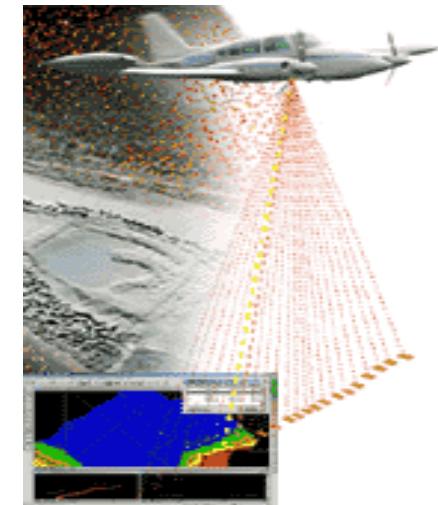
Obviously not only in sciences:

- Economist 02/10:
 - From 150 Billion Gigabytes five years ago to 1200 Billion today
 - Managing data deluge difficult; doing so will transform business/public life



Example: Massive Terrain Data

- New technologies: Much easier/cheaper to collect detailed data
 - Previous ‘manual’ or radar based methods
 - Often 30 meter between data points
 - Sometimes 10 meter data available
 - New laser scanning methods (**LIDAR**)
 - Less than 1 meter between data points
 - Centimeter accuracy (previous meter)



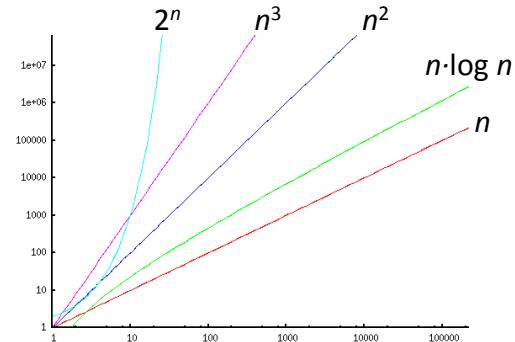
Denmark

~ 2 million points at 30 meter (<1GB)
~ 18 billion points at 1 meter (>1TB)

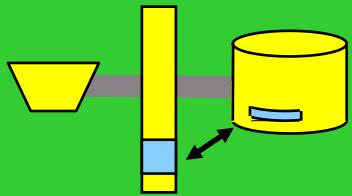
Algorithm Inadequacy

- **Algorithms:** Problem solving “recipies”
- Importance of **scalability/efficiency**
 - Algorithmics core computer science area
- Traditional algorithmics:
 - Transform input to output using simple machine model
- Inadequate with e.g.
 - Massive data
 - Small/diverse devices
 - Continually arriving data

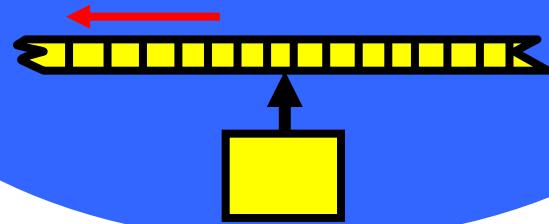
→ Software inadequacies!



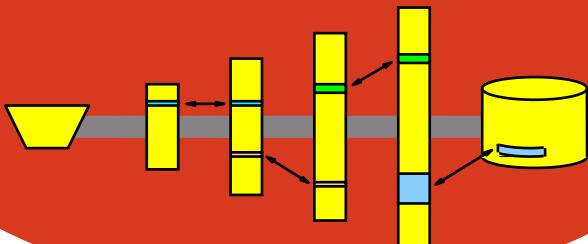
I/O Efficient Algorithms



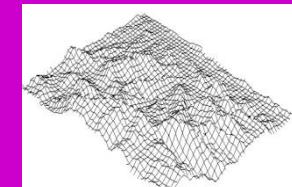
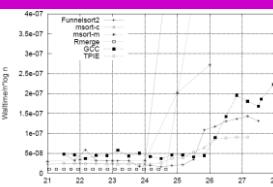
Streaming Algorithms



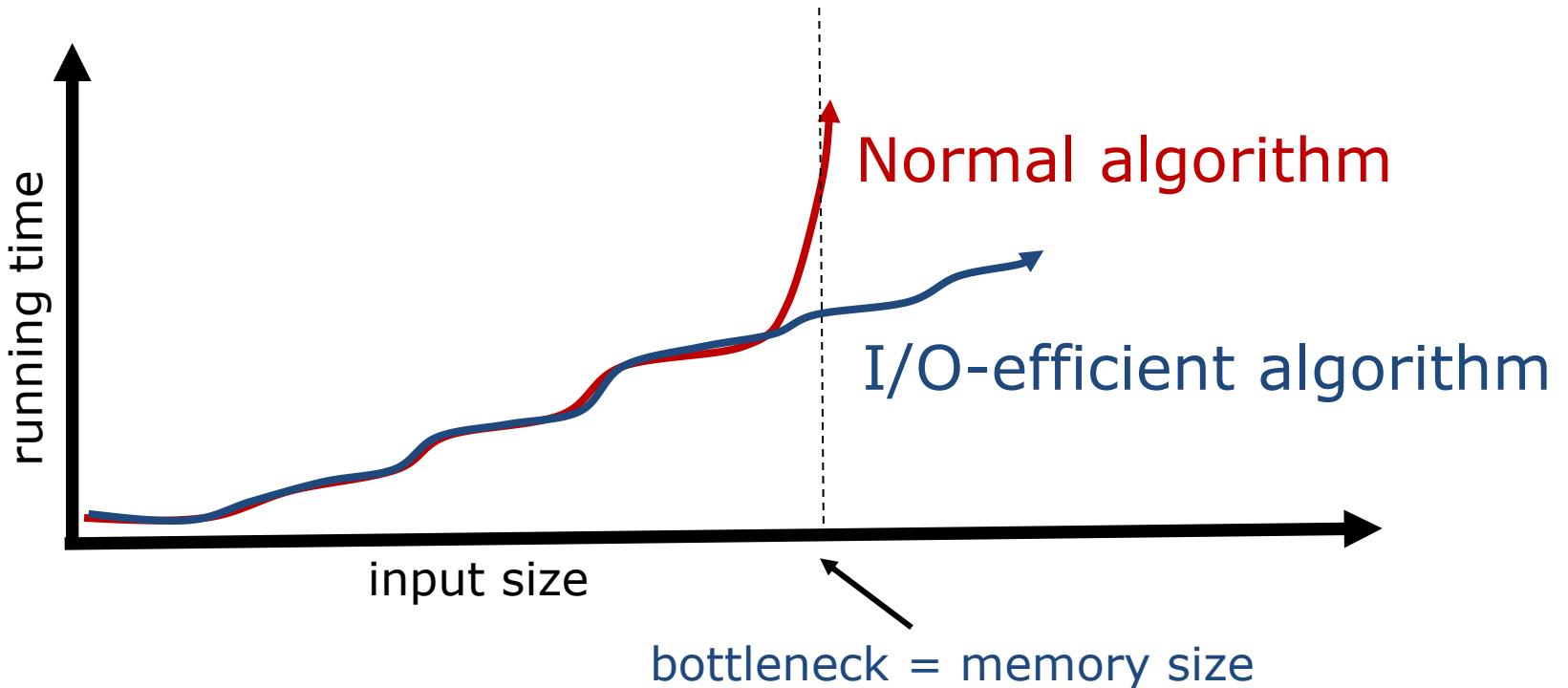
Cache Oblivious Algorithms



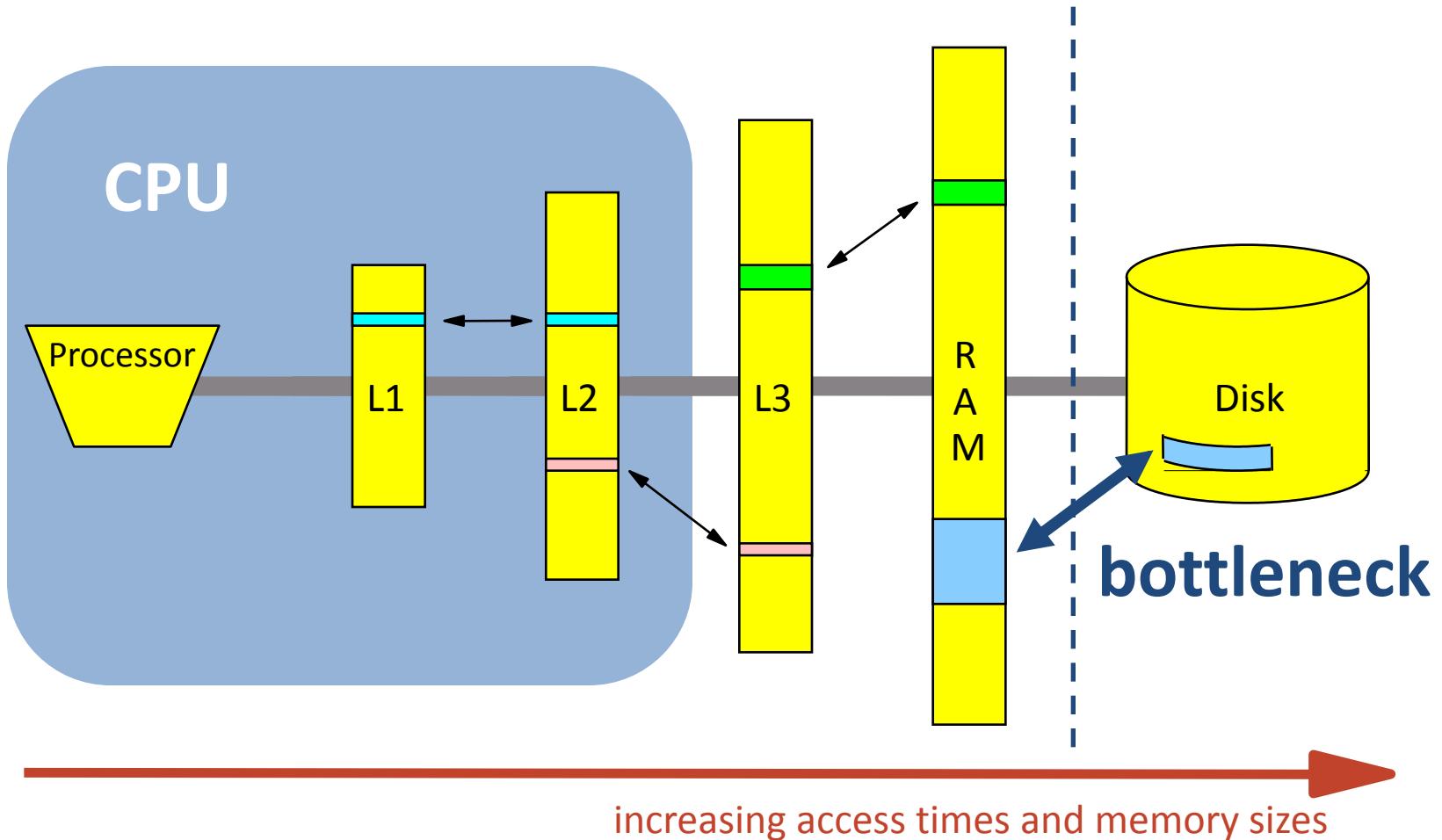
Algorithm Engineering



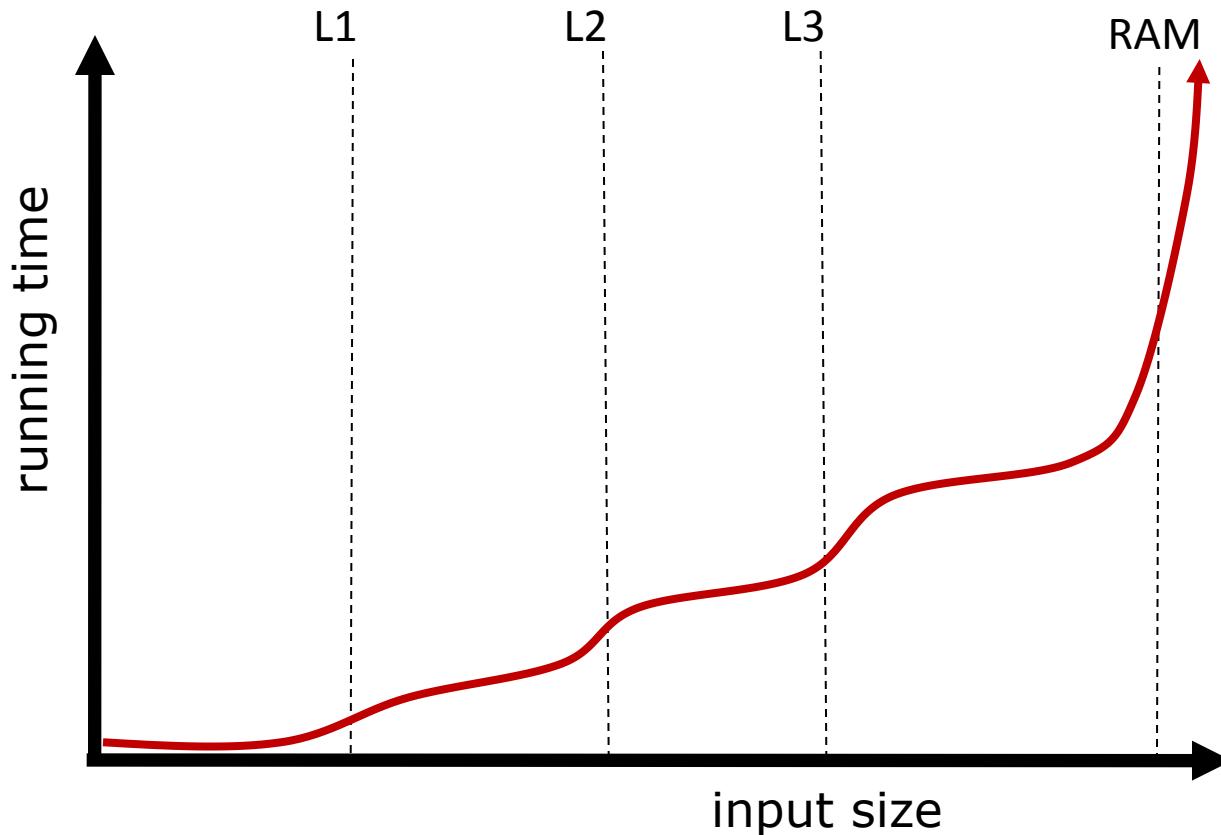
The problem...



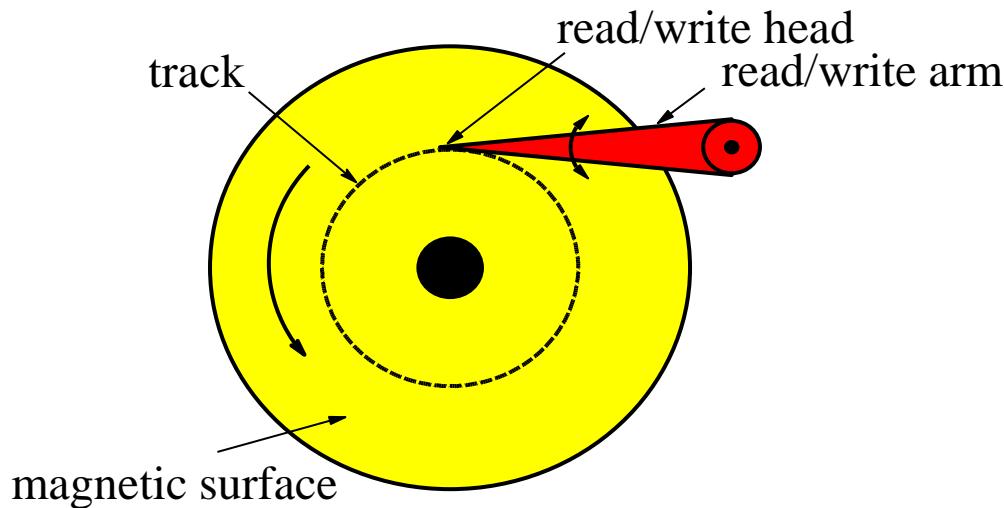
Memory Hierarchies



Memory Hierarkies vs. Running Time

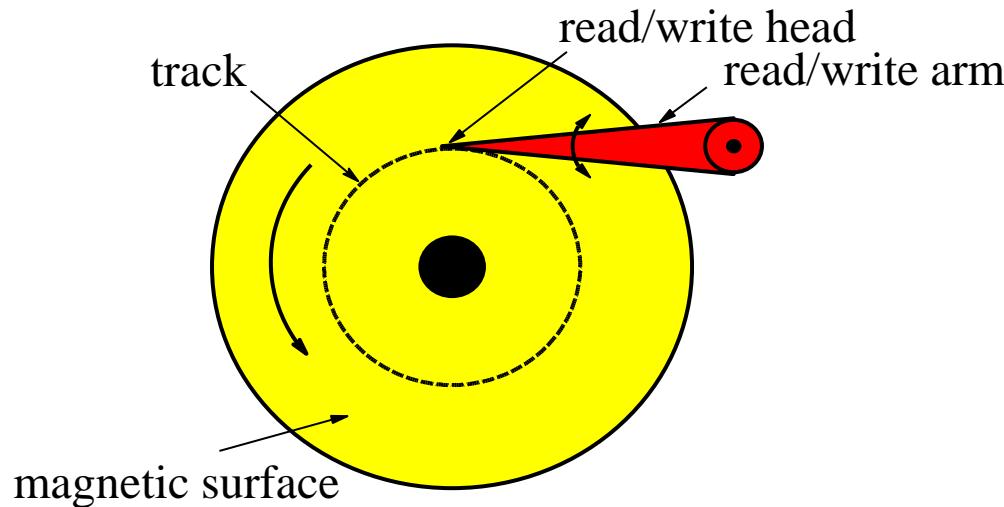


Disk Mechanics



"The difference in speed between modern CPU and disk technologies is analogous to the difference in speed in sharpening a pencil using a sharpener on one's desk or by taking an airplane to the other side of the world and using a sharpener on someone else's desk." (D. Comer)

Disk Mechanics



- I/O is often bottleneck when handling massive datasets
- Disk access is **10⁷ times slower** than main memory access!
- Disk systems try to amortize large access time transferring **large contiguous blocks** of data
- Need to store and access data to **take advantage of blocks !**

I/O-efficient algorithms

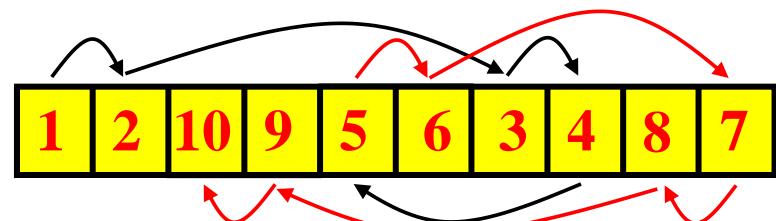
Move as few disk blocks as possible to solve given problem !

Memory Access Times

	Latency	Relative to CPU
Register	0.5 ns	1
L1 cache	0.5 ns	1-2
L2 cache	3 ns	2-7
DRAM	150 ns	80-200
TLB	500+ ns	200-2000
Disk	10 ms	10^7

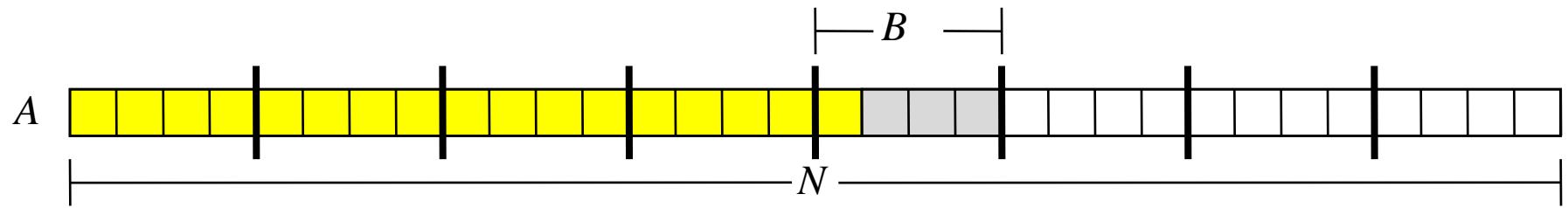
I/O-Efficient Algorithms Matter

- Example: Traversing linked list (List ranking)
 - Array size $N = 10$ elements
 - Disk block size $B = 2$ elements
 - Main memory size $M = 4$ elements (2 blocks)



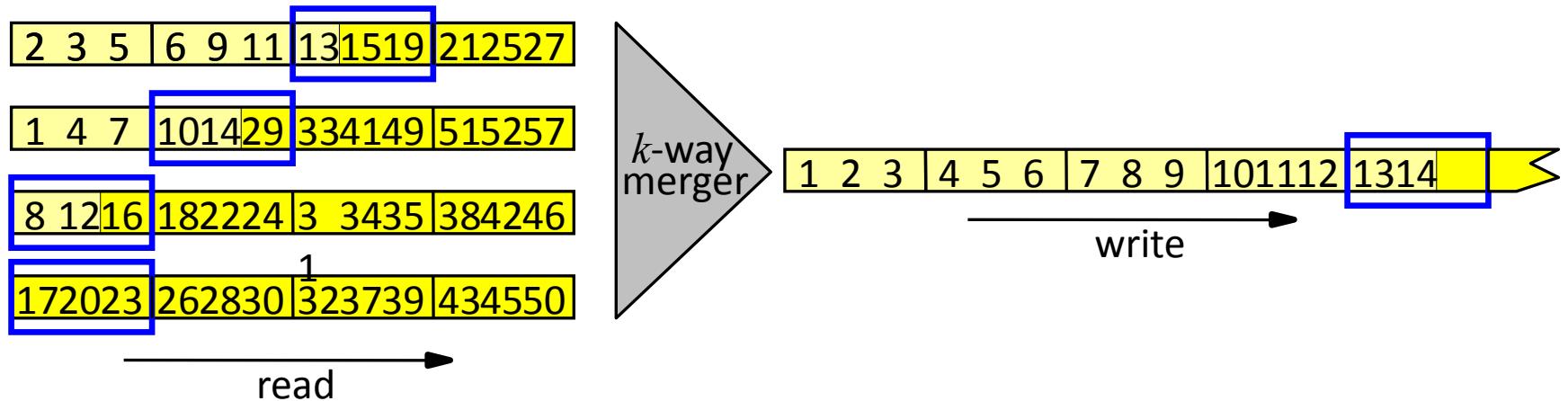
- Difference between N and N/B large since block size is large
 - Example: $N = 256 \times 10^6$, $B = 8000$, 1ms disk access time
 $\Rightarrow N$ I/Os take 256×10^3 sec = 4266 min = 71 hr
 $\Rightarrow N/B$ I/Os take $256/8$ sec = 32 sec

I/O Efficient Scanning



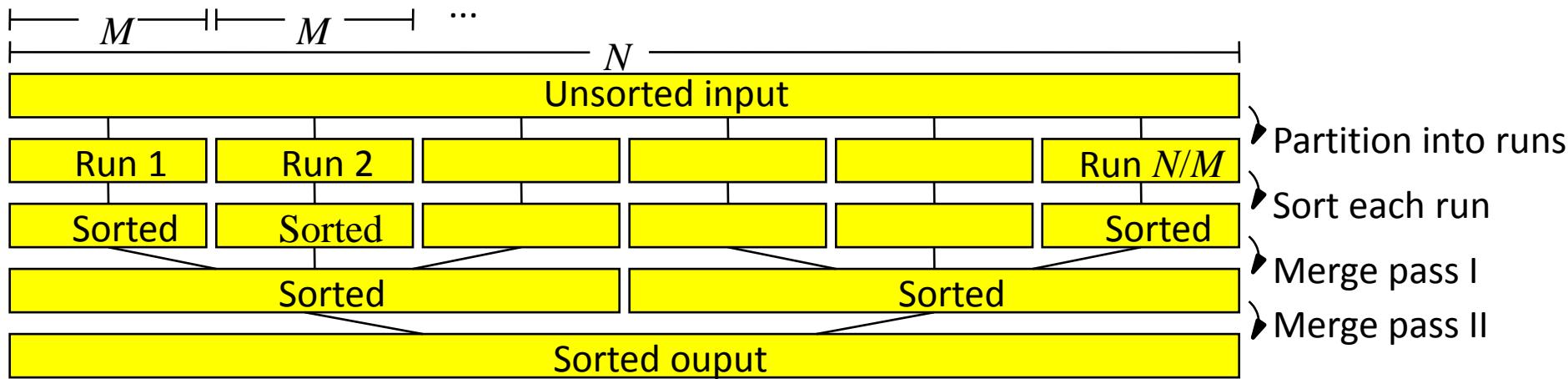
$O(N/B)$ I/Os

External-Memory Merging



Merging k sequences with N elements requires $O(N/B)$ IOs
(provided $k \leq M/B - 1$)

External-Memory Sorting



- MergeSort uses $O(N/B \cdot \log_{M/B}(N/B))$ I/Os
- Practice number I/Os: 4-6 x scanning input

Energy-Efficient Sorting using Solid State Disks

(Bechman, Meyer, Sanders, Siegler 2010)

Sorting large data sets

- Is easily described
- Has many applications
- Stresses both CPU and the I/O system
- Benchmark introduced 1985

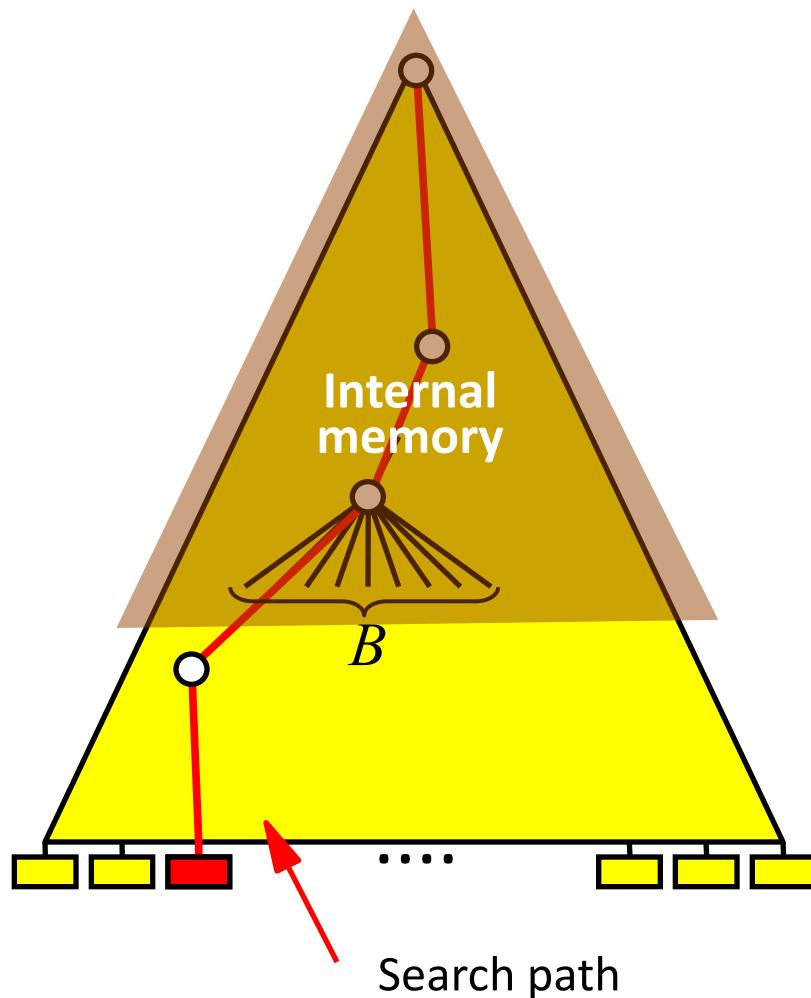
Energy Efficiency

- Energy (and cooling) is a significant cost factor in data centers
- Energy consumption correlates to pollution

Size [GB]	2007			2010			Energy Saving Factor
	Time [s]	Energy [kJ]	Rec./J	Time [s]	Energy [kJ]	Rec./J	
10	86.6	8.6	11628	76.7	2.8	35453	3.0
100	881	88.1	11354	756	27.5	36381	3.2
1000	7196*	2920*	3425	21906	723.7	13818	4.0

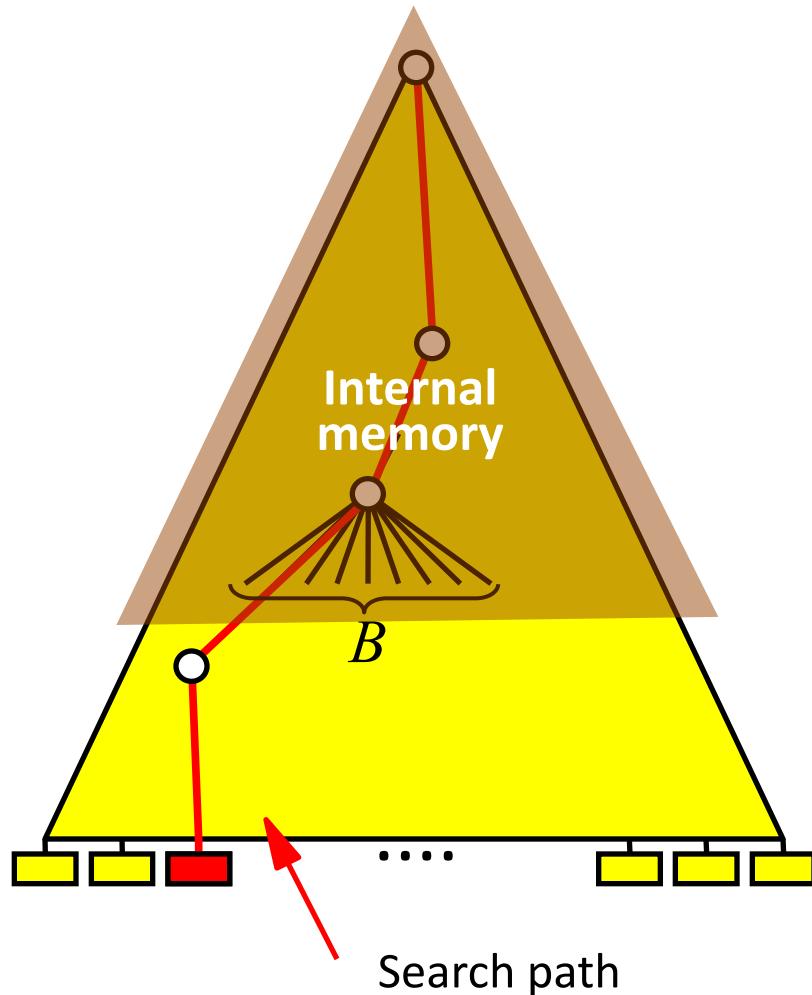


B-trees - The Basic Searching Structure



- Searches
Practice: 4-5 I/Os
- Repeated searching
Practice: 1-2 I/Os

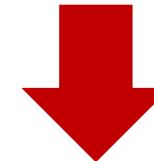
B-trees



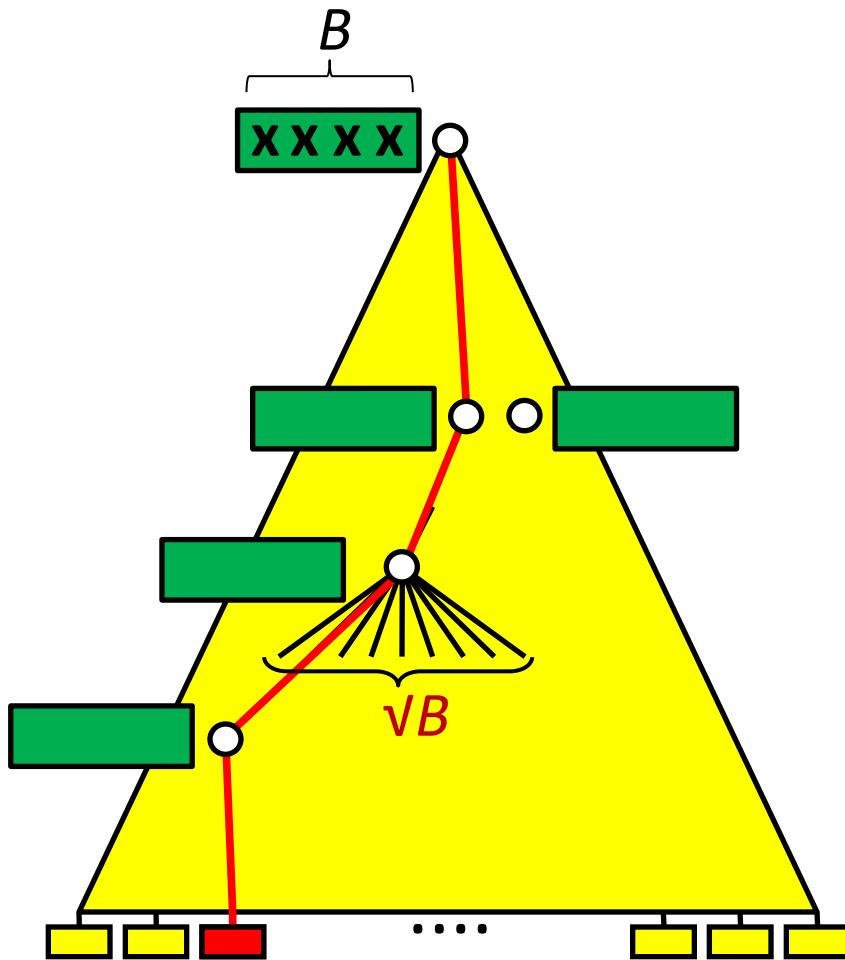
- **Searches $O(\log_B N)$ I/Os**

- **Updates $O(\log_B N)$ I/Os**

Best possible



B-trees with Buffered Updates



- Searches cost

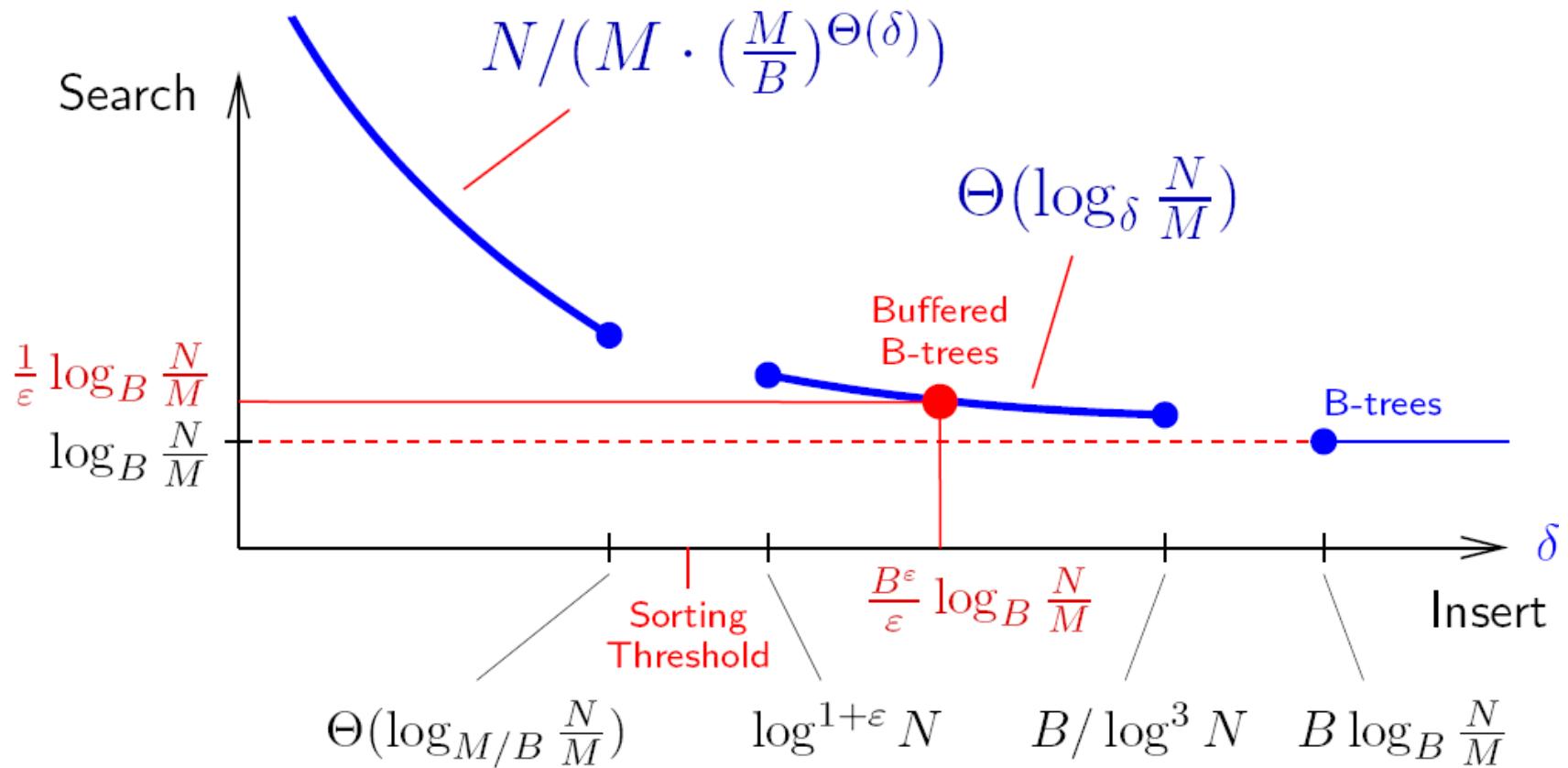
$$O(\log_B N) \text{ I/Os}$$

- N updates cost

$$O(N / \sqrt{B} \cdot \log_B N) \text{ I/Os}$$

Trade-off between search and update times – optimal !

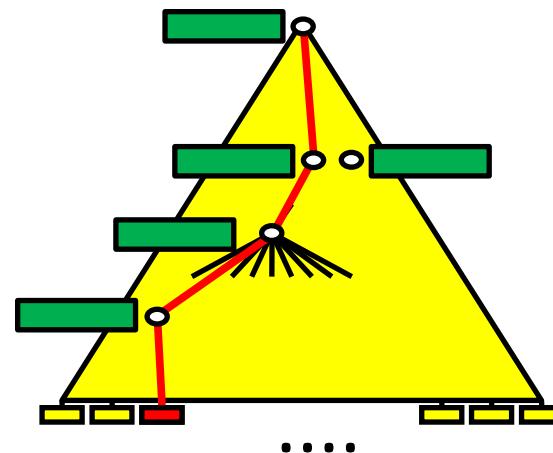
B-trees with Buffered Updates



B-trees with Buffered Updates

Experimental Study

- 100.000.000 elements
- Search time basically unchanged with buffersize
- Updates 100 times faster





Flood Prediction Important

- Prediction areas susceptible to floods
 - Due to e.g. raising sea level or heavy rainfall
- Example: Hurricane Floyd Sep. 15, 1999



7 am



3pm



Detailed Terrain Data Essential

Mandø with 2 meter sea-level raise



80 meter terrain model



2 meter terrain model

Surface Flow Modeling

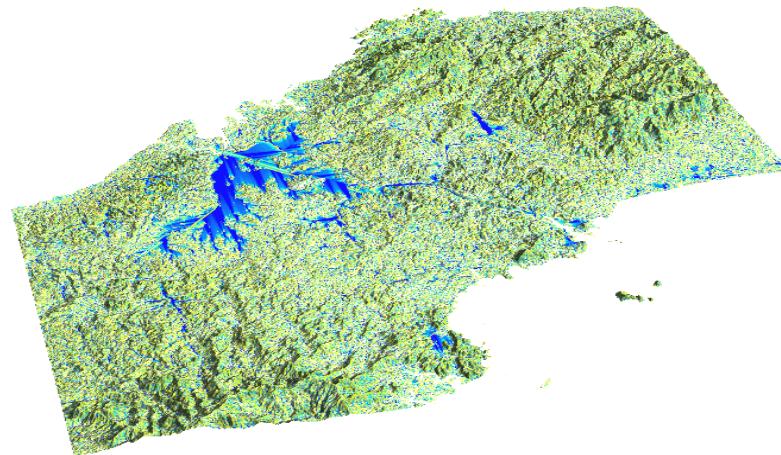


7 am



3pm

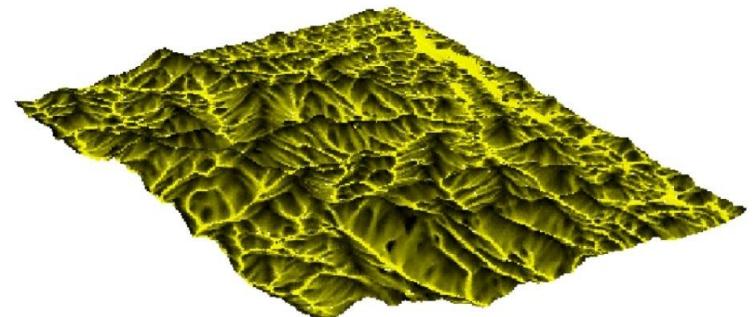
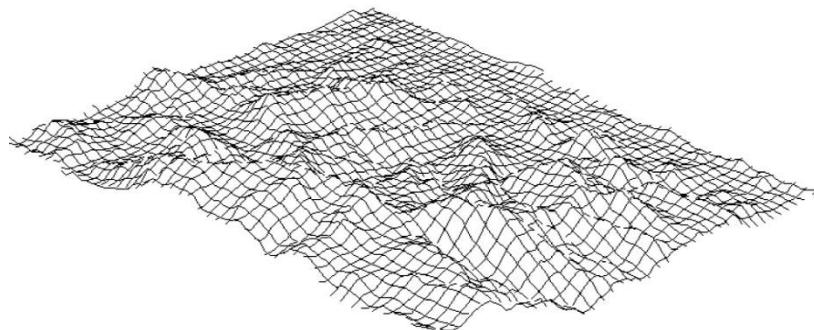
Hurricane Floyd (September 15, 1999)



- Conceptually flow is modeled using two basic attributes
 - **Flow direction:** The direction water flows at a point
 - **Flow accumulation:** Amount of water flowing through a point

Flow Accumulation

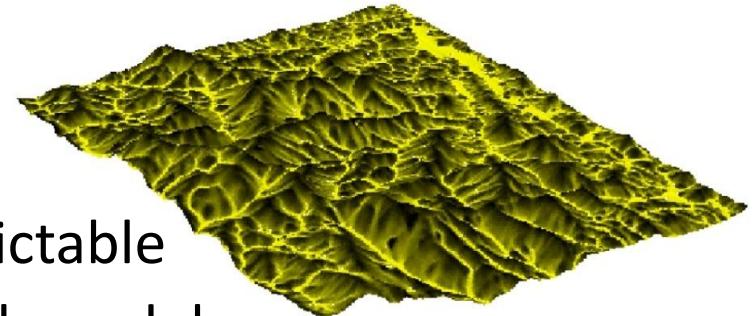
- Flow accumulation on grid terrain model:
 - Initially one unit of water in each grid cell
 - Water (initial and received) distributed from each cell to lowest lower neighbor cell (if existing)
 - Flow accumulation of cell is total flow through it



- Note:
 - Flow accumulation of cell = size of “upstream area”
 - Drainage network = cells with high flow accumulation

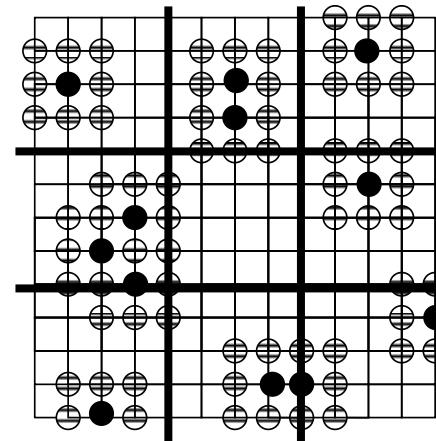
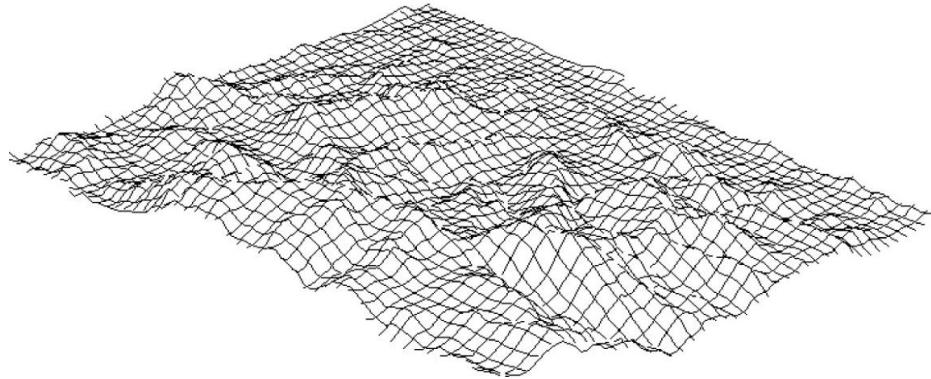
Massive Data Problems

- Commercial systems:
 - Often *very* slow
 - Performance somewhat unpredictable
 - Cannot handle 2-meter Denmark model
- Collaboration environmental researchers in late 90'ties
 - US Appalachian mountains dataset
 - 800x800km at 100m resolution \Rightarrow a few Gigabytes **14 days!!**
 - Customized software on $\frac{1}{2}$ GB machine:
- Appalachian dataset would be Terabytes sized at 1m resolution!



Flow Accumulation Performance

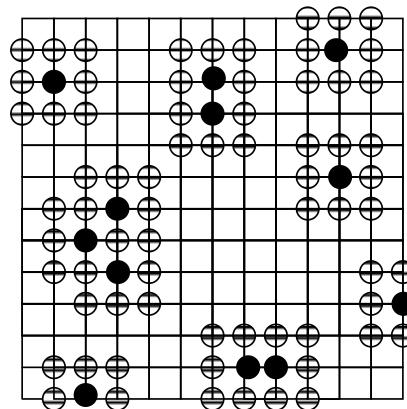
- Natural algorithm access disk for each grid cell
 - “Push” flow down the terrain by visiting cells in height order
⇒ Problem since cells of same height scattered over terrain



- Natural to try “tiling” (commercial systems?)
 - But computation in different tiles not independent

I/O-Efficient Flow Accumulation

- We developed I/O-optimal algorithms (assessing disk a lot less)
- Avoiding scattered access by:
 - Grid storing input: Data duplication
 - Grid storing flow: “Lazy write”



- Implementation very efficient
 - Appalachian Mountains flow accumulation in 3 hours!
 - Denmark at 2-meter resolution in a few days

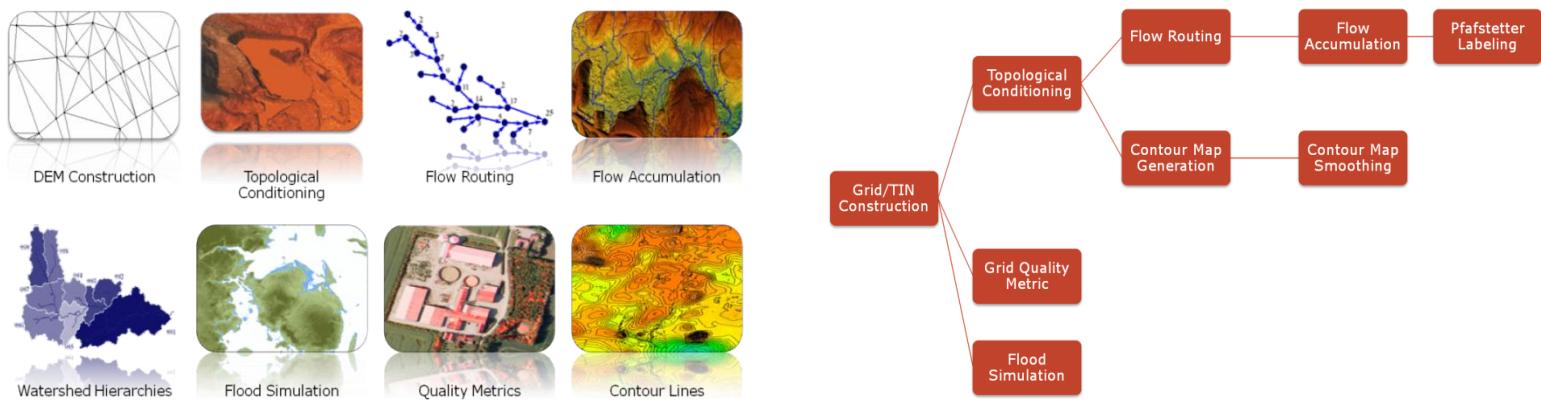
Flood Modeling

- Not all terrain below height h is flooded when water rise to h meters!
- Theoretically not too hard to compute area flooded when rise to h meters
 - But no software could do it for Denmark at 2-meter resolution
- Use of I/O-efficient algorithms
⇒ Denmark in a few days
- Even compute new terrain where terrain below h is flooded when water rise to h



TerraSTREAM

- Flow/flooding work part of comprehensive software package
 - **TerraSTREAM:** Whole pipeline of terrain data processing software



- **TerraSTREAM used on full 2 meter Denmark model**
(\sim 25 billion points, \sim 1.5 TB)
 - Terrain model (grid) from LIDAR point data
 - Surface flow modeling: Flow directions and flow accumulation
 - Flood modeling

Interdisciplinary Collaboration

- Flow/flooding work example of center interdisciplinary collaboration
 - Flood modeling and efficient algorithms in biodiversity modeling
 - Allow for used of global and/or detailed geographic data
- Brings together researchers from
 - Biodiversity
 - Ecoinformatics
 - Algorithms
 - Datamining
 - ...



10110001101110011101100101010010101001111001000010001110010111100001000111001111010001010001110
100001001110000100001111011100110111101011100111101001111000010000111110110100110001000011011
01110010101111111100001001010110000011101101101110011001001100001011011110110110110110100000
00000000100011000010000111010011100101001000000010100000001000100100011010011100100110001000010
00101111111010010011001100001001100110100101000000010011010011101001001101100011100010011101001111
1101101001111010000110100100111001101111100001010010100010001001011011101111110000110100100001
00001011000101110101001111110011001111110100101000100010011000101001001100000001000001100011101011
110000110010111000100111101000001101011100100011010010001010001111010010000110100010110001101011
1100101011000100100101101010000101010000011011000000010101100101110111000100101111101010110111
110000100001010011010111100011110011010011101110010010110111101000100011011011111010001000100110110
01001111011101011001101111110100000001111010000101110010011100011010100001100100110001111101100010
110000011101100110100111100010001110001001111010001000110101010001010110110101000011001001000110101
111111110111011001111001010100001100011101010111001010110111001010101101101101010110000101011
10001100110100110001010000100000000001111001001101010011111101001000001110011100100010001000000
10110110111101000010101111111010100001011010000010110011101001100110010010011111011101010110111
01010101010001001100111111011110011111100111101011110100001100100001010010010010101011110011100010
010111101101000010010011011001111011100101000010000111011000011101000011011000011011010110110010
1101011010011100000010100101111010010001010100101011110101011111110010001101100110001000110
0011101011000111100010101011011111110111011000111101000110100000000001111110011101110000001110
100001011011100110100010010111110000010100000011000000000010110001010010110011101001001101101110110
10111100111001010000111100101000100001101010110010100000001101111001111100011010110000010111000001
11100001101001101111101001010111101100101101111010000010000110101101011001100010110101011010
11010100110100111010010001001110000011010111101011110001011110100001101110001010011000101000110
0111101011010011110110000101100010100110110001011011001000101011010111011000110111010101110101111
111110011010110000000000000000100000100000110011000110110110010001100110101010110101100011011000101
0101001110110010010111100010110010100010011100001000000111001100000110100110001000000001100111
00101000010111001001100011001110001101000011010001111011100101011101001101100001010011010110110
1100111010110111110100100111111011100100010001010111101110100110011111000001011010100100001011000010
00001010000100010000000101101110010101011010000101000100011000001010101010100110010111011101011
11000101100000110011101001101101111010110110001000011101100111010001010011111001001010101010101101
011010011111000001010111000101101100011010001111111001100110011111100011001100111111000110011001111110011001

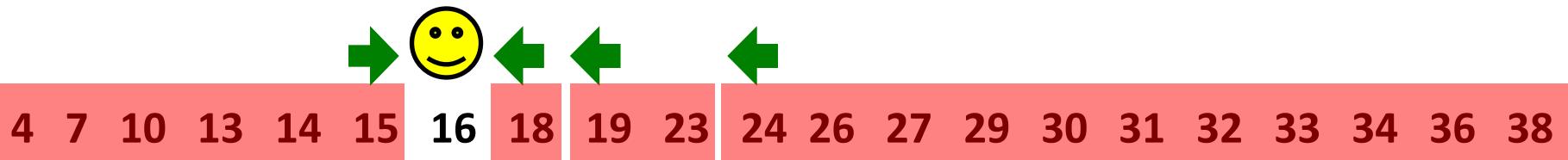
A bit in memory changed
value because of e.g.
background radiation,
system heating,...



"You have to provide reliability on a software level. If you're running 10,000 machines, something is going to die every day."

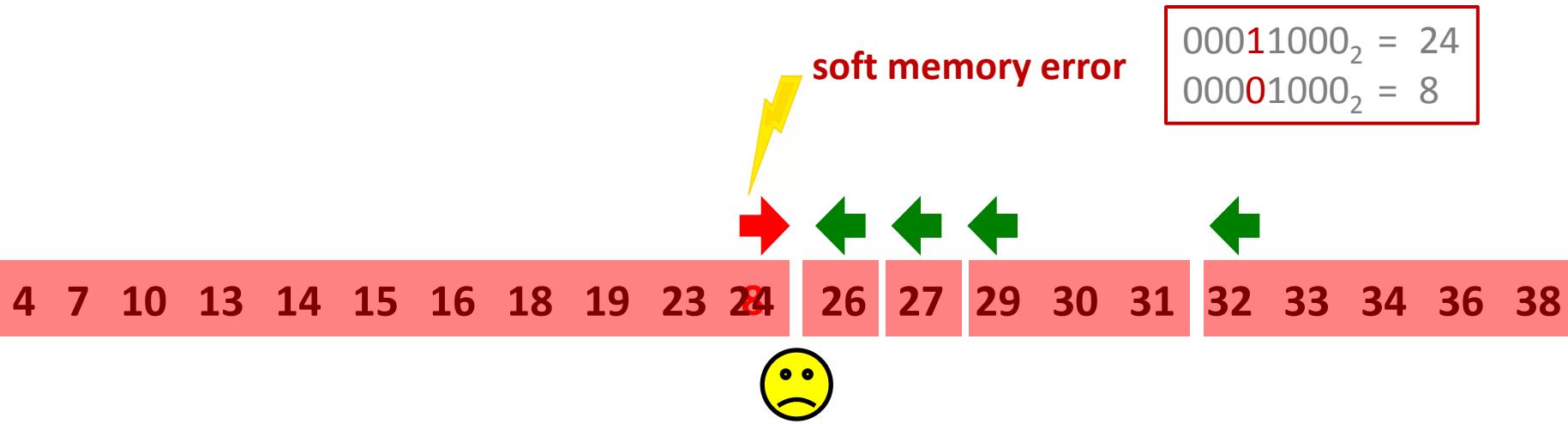
–  fellow Jeff Dean

Binary Search for 16



$O(\log N)$ comparisons

Binary Search for 16



Requirement: If the search key occurs in the array as an uncorrupted value, then we should report a match !

Where is Laurits ?



Where is Laurits ?



Where is Laurits ?

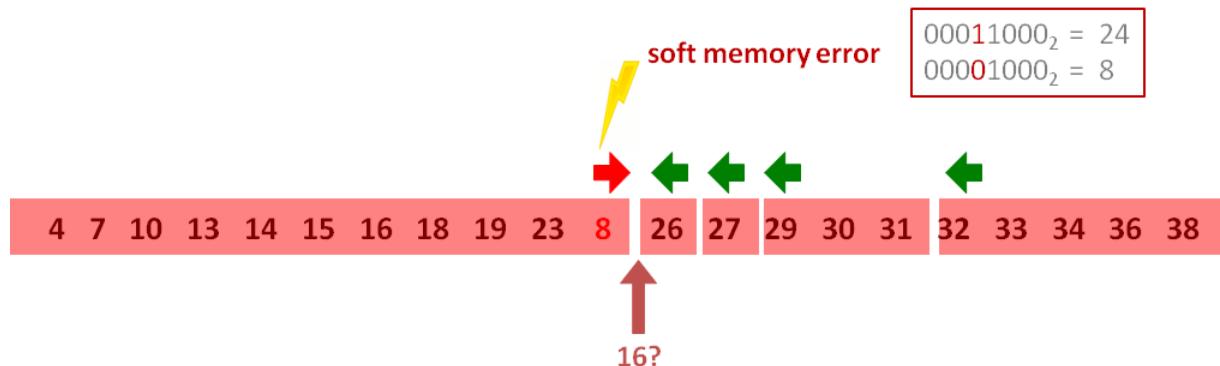


If at most 4 faulty answers then Laurits is somewhere here

Faulty-Memory RAM Model

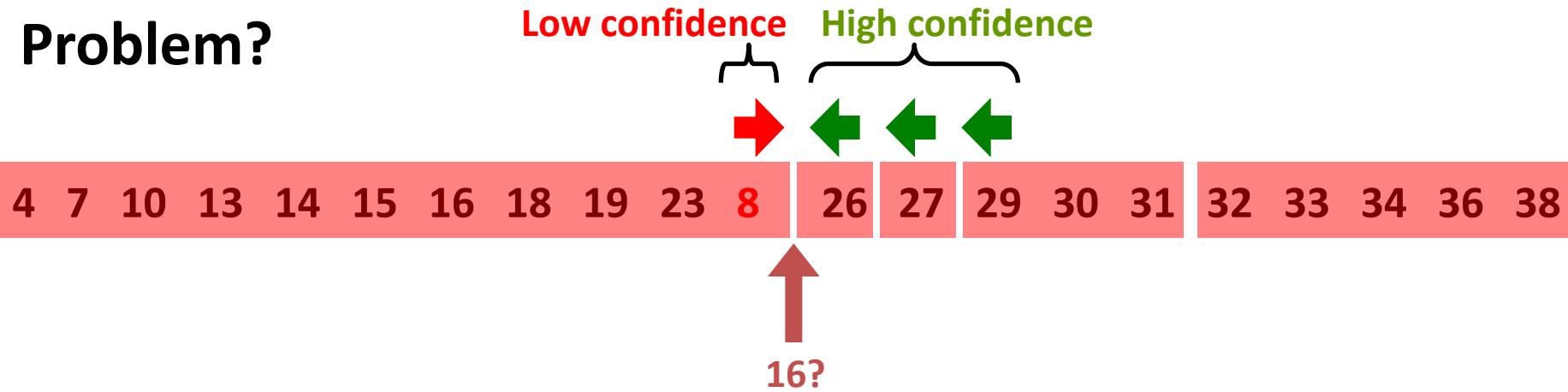
Finocchi and Italiano, STOC'04

- Content of memory cells can get corrupted
- Corrupted and uncorrupted content cannot be distinguished
- $O(1)$ safe registers
- Assumption:** At most δ corruptions



Faulty-Memory RAM: Searching

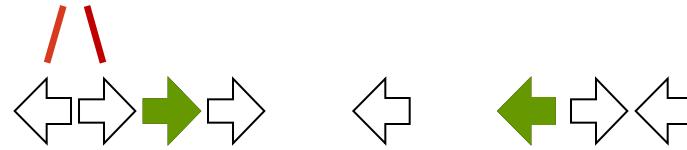
Problem?



Faulty-Memory RAM: Searching

When are we
done ($\delta=3$)?

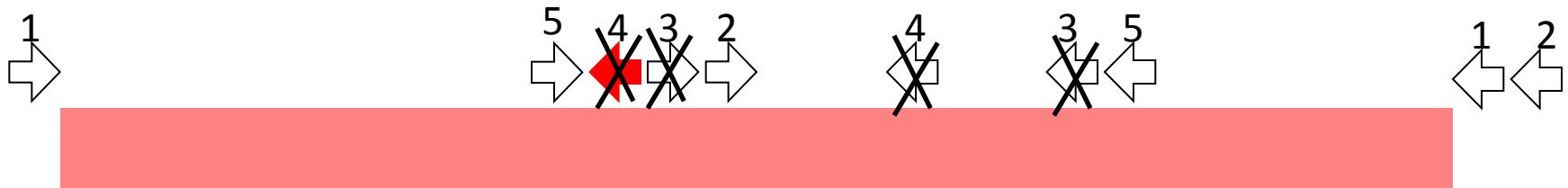
Contradiction, i.e. at
least one fault



If range contains at least $\delta+1 \rightarrow$ and $\delta+1 \leftarrow$ then there is at least one uncorrupted \rightarrow and \leftarrow , i.e. x must be contained in the range

Faulty-Memory RAM: $\Theta(\log N + \delta)$ Searching

Brodal, Fagerberg, Finocchi, Grandoni,
Italiano, Jørgensen, Moruz, Mølhave, ESA'07



If verification fails

- contradiction, i.e. ≥ 1 memory-fault
- ignore 4 last comparisons
- **backtrack** one level of search

Summary

- Basic research center in Aarhus
 - Organization
 - PhD education
- Examples of research
 - Theoretical external memory algorithmics
 - Practical (flow and flooding simulation)

Tau · Jërë-jëf · Tashakkur · S.aHHa · Sag olun
Giihtu · Djakujo · Dâkujem vám · Thank you
Tesekkür ederim · To-siä · Merci · Tashakur
Taing · Dankon · Efharisto' · Shukriya · Kjitos
Dhanyabad · Rakhmat · Trugarez · Asante
Köszönöm · Blagodarya · Dziekuje · Eskerrik asko
Grazie · Tak · Bayarlaa · Miigwech · Dank u
Spasibo · Dêkuji vám · Ngiyabonga · Dziakuj
Obrigado · Gracias · A dank aych · Salamat
Takk · Arigatou · Tack · Tänan · Aciu
Korp kun kah · Multumesk · Terima kasih · Danke
Rahmat · Gratias · Mahalo · Dhanyavaad
Paldies · Faleminderit · Diolch · Hvala
Kam-sa-ham-ni-da · Xìe xìe · Mèrcie · Dankie