

maDALGO

CENTER FOR MASSIVE DATA ALGORITHMICS



Gerth Stølting Brodal
Aarhus Universitet

Science Center Sorø - 4. september 2009

Overblik

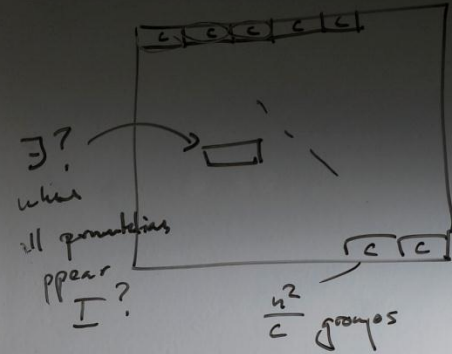
- Hvem er jeg ?
- Hvad er MADALGO ?
- Algoritmiske problemstilling ?
- Vidensspredning

Gerth Stølting Brodal

Cand. scient., Aarhus Universitet 1994

Ph.d., Aarhus Universitet, 1997

Lektor, Aarhus Universitet, 2004-



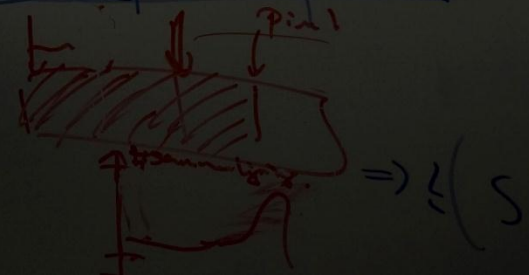
$$\sum_{i=1}^n \sum_{j=1}^n s_{ij} \geq \frac{N}{c}$$
$$y_j s \geq \frac{S}{c}$$
$$s > S \geq 2$$

$$\left\{ \begin{array}{l} \frac{n^2}{c} \text{ bits} \\ (c!)^{n^2/c} \text{ different inputs} \end{array} \right. \Rightarrow \leq 2^{n^2/c} \text{ data structures}$$

$$\Downarrow \exists \frac{\text{data structure}}{D} \geq \frac{(c!)^{n^2/c}}{2^{n^2/c}} \text{ different inputs}$$



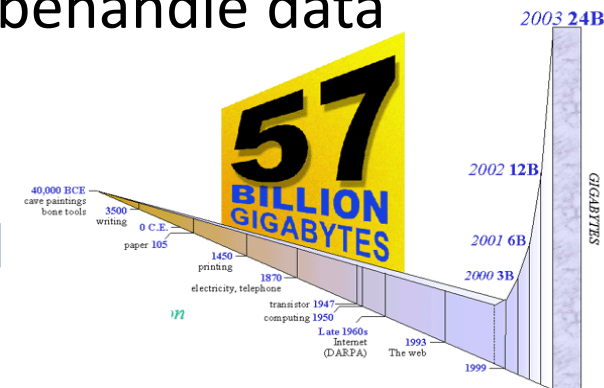
$$\frac{n^2}{2} \text{ elements}$$
$$\frac{3}{4} n \frac{N}{5} \text{ groups}$$



madALGO 
CENTER FOR MASSIVE DATA ALGORITHMICS

Center motivation: Massive Data

- Computere og censorer anvendes overalt
- Øgede muligheder for at indsamle/gemme/behandle data
→ Massive data tilstede overalt
- Samfundet mere og mere “datadrevet”
→ Tilgå/behandle data overalt til enhver tid



Nature

- 2/06: “2020 – Future of computing”
- 9/08: “BIG DATA”
- Videnskabelig data vokser eksponentielt, mens kvalitet og tilgængelighed forbedres
- Paradigmeskift: *Videnskab vil blive om datamining*
→ Datalogi altafgørende i alle videnskaber




Massive Data Eksempler

- **Telefon:** AT&T 20 TB telefonopkalds database, trådløs sporing
- **Forbrugere:** WalMart 70 TB database, købsmønstre
- **WEB:** Google indeks med »8 milliarder websider
- **Bank:** Danske Bank 250 TB DB2 database
- **Geografi:** NASA satelliter genererer terrabytes hver dag



Lars Arge

- Center ved  **DANMARKS GRUNDFORSKNINGSFOND**
DANISH NATIONAL RESEARCH FOUNDATION
- Lars Arge, Professor, Centerleder
- Gerth S. Brodal, Lektor
- 5 post docs., 13 ph.d. studerende, 2 kandidat studerende, 4 TAP
- Total budget for 5 år ca. 60 millioner kr

- **Overordnede mål**
 - Fremme den algoritmiske viden inden for processering af massive data
 - Træning af forskere i et verdensførende miljø
 - At være katalysator for multidisciplinære samarbejder



DANMARKS
GRUNDFORSKNINGSFOND
DANISH NATIONAL RESEARCH FOUNDATION

- Etableret 1991
 - 2 milliarder kr
 - Yderligere 3 milliarder i 2009
- Støtter
 - Grundforskningscentre (“Center of excellence”)
 - Højt profilerede gæsteprofessorer
 - Ph.d. skoler
- Center of excellence
 - Grundforskning i verdensklasse
 - 5 år; nogle forlænges med yderligere 5 år
 - P.t. ca. 40 centre
 - Gennemsnitlige 5 års bevilling på ca. 40 million kr

MADALGO kerneforskere



Arge
(AU)

Brodal
(AU)

I/O, cache, og algoritme
engineering



Demaine
(MIT)

Indyk
(MIT)

Cache og streaming



Mehlhorn
(MPI)

Meyer
(FRA)

I/O og algoritme engineering

madALGO

CENTER FOR MASSIVE DATA ALGORITHMICS



Center Aktiviteter

- **Besøg** af kerneforskerne
- **Udveksling** af AU, MPI, FRA og MIT post docs. og ph.d. studerende
- **Gæsteophold** af forskere og studerende fra andre institutioner
- **Diverse workshops**
 - Incl. multidisciplinære og i samarbejde med industrien
- Førende **internationale begivenheder**:
 - 25th Annual Symposium on Computational Geometry in 2009
 - Workshop on Algorithms for Massive Datasets in 2009
- **Sommerskoler**
 - 2007: Streaming algorithms
 - 2008: Cache-oblivious algorithms

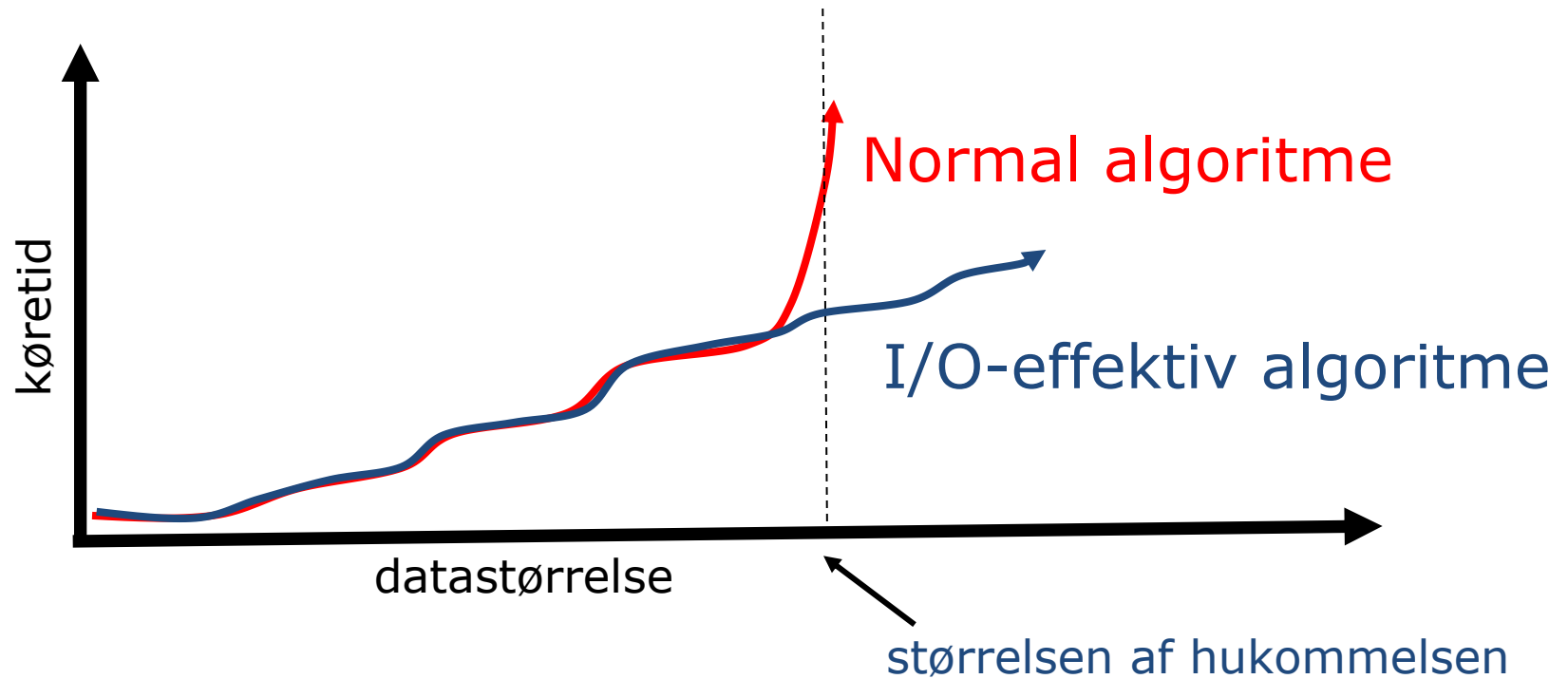


Center samarbejder

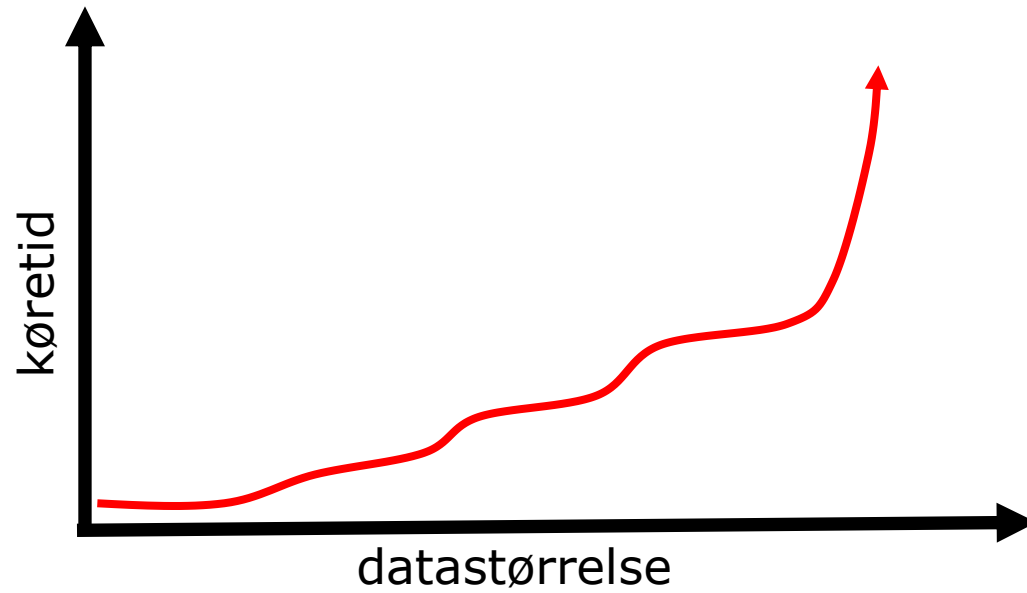
- COWI, DHI, DJF, DMU, Duke, NSCU
- Finansering fra det Strategiske Forskningsråd og US Army Research Office
- Software platform for Galileo GPS
 - Adskillige danske akademiske/industrielle partnere
 - Finansering fra Højteknologi Fonden
- Europæisk netværk om massive data algoritmik
 - 8 førende europæiske forskningsgrupper

Algoritmiske Problemstilling

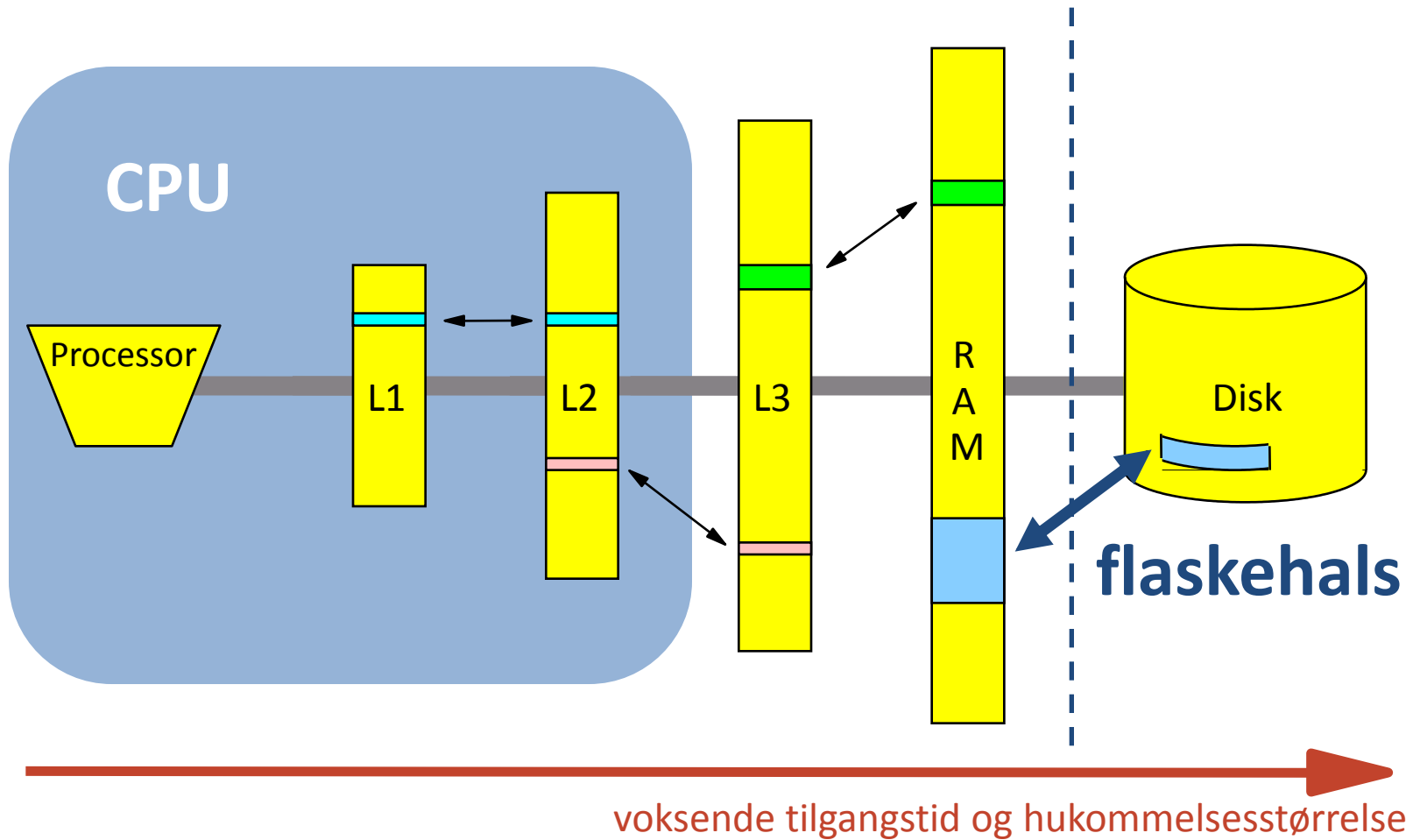
Problemet...



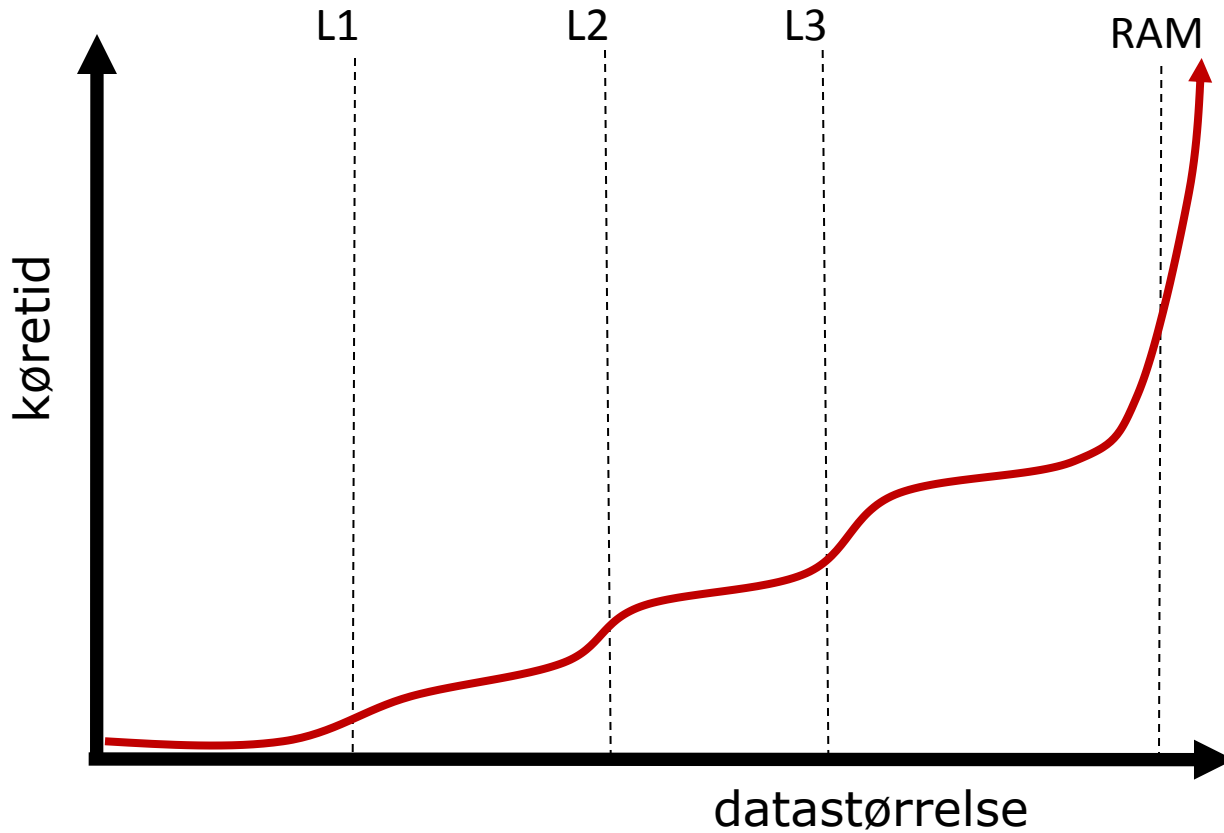
Hvad er flaskehalsene ?



Hukommelseshierarkier



Hukommeshierarkier vs. Køretid



Algoritmik

- Central betydning for skalerbarhed/effektivitet
→ Algoritmik centralt datalogisk område

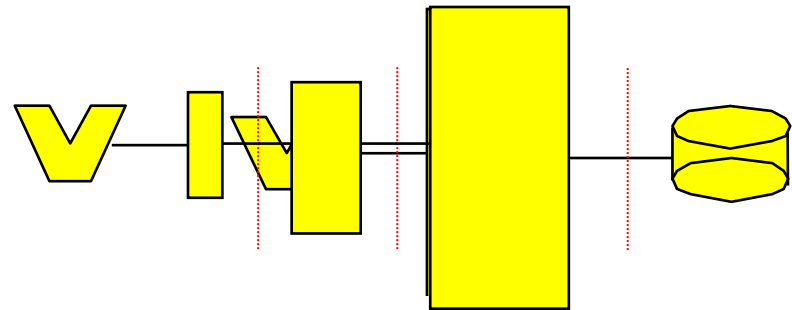
- Traditionel algoritmik:

Transformer input til output ved anvendelse af en simpel maskinmodel

- Utilstrækkelig til f.eks.

- Massive data
- Små/varierende maskiner
- Kontinuære datastrømme

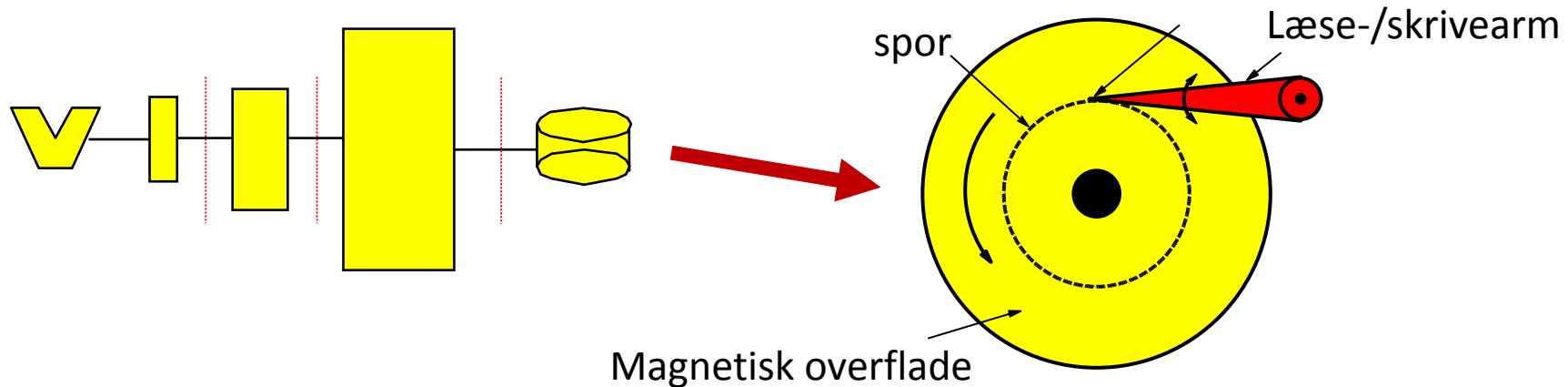
→ Software med begrænsninger!



- Faglige grupperinger har arbejdet med disse mangler
– men meget mangler stadig at blive løst

I/O-Effektive Algoritmer

- Problemer involverende massive data på disk

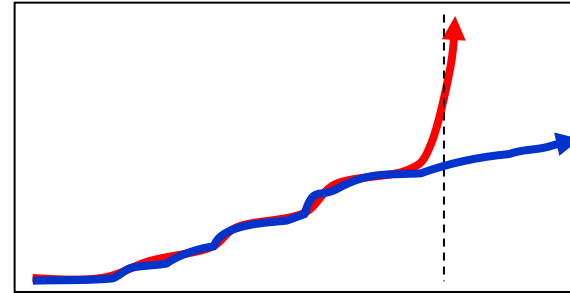
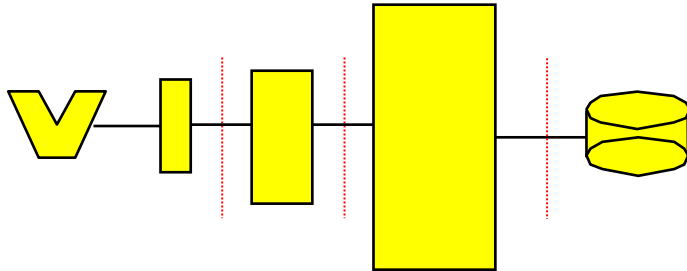


- Tilgang til disk er 10^6 gange langsommere end hukommelsen

“Forskellen i hastighed mellem moderne CPU- og diskteknologier svarer til forskellen i hastigheden mellem at spidse sin blyant med sin blyantspids på sit skrivebord og at tage et fly til den anden side af jorden og anvende en anden persons blyantspids på dennes skrivebord.” (D. Comer)

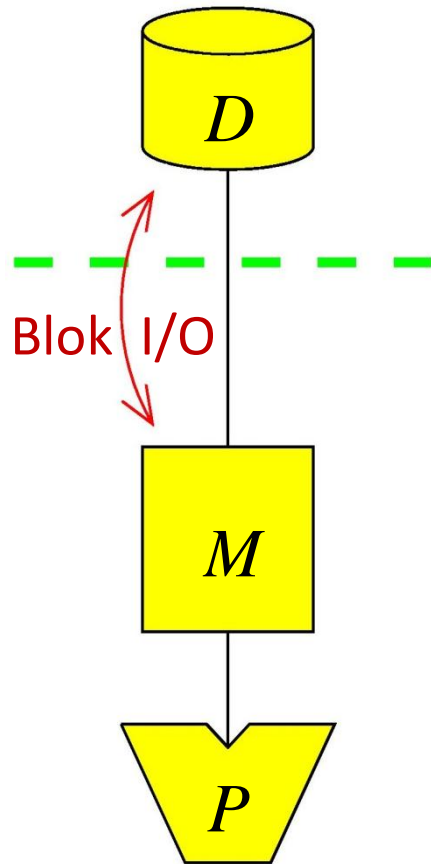
I/O-Effektive Algoritmer

- Problemer involverende massive data på disk



- Tilgang til disk er 10^6 gange langsommere end hukommelsen
 - Dyr tilgangstid amortiseres ud ved at overføre store blokke af data
- Altafgørende at udnytte blokkene når data gemmes/tilgås
- I/O-effektive algoritmer:
 - Flytte så få blokke som muligt for at løse et givet problem

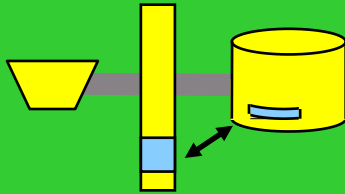
I/O-Effektive Algoritmer: I/O-model



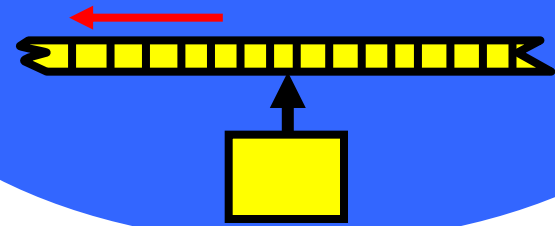
- Model parametre
 - $N = \#$ elementer i input
 - $B = \#$ elementer i en disk blok
 - $M = \#$ elementer i intern hukommelse
- Mål: Minimer $\#$ I/O
 - Flyt B sammenhængende elementer mellem hukommelsen og disk
- Typiske teoretiske grænser:
 - Sortering: $O\left(\frac{N}{B} \log_{M/B} \frac{N}{B}\right)$
 - Søgning: $O(\log_B N)$
 - Prioritets køer: $O\left(\frac{1}{B} \log_{M/B} \frac{N}{B}\right)$ amortiseret

MADALGO Fokusområder

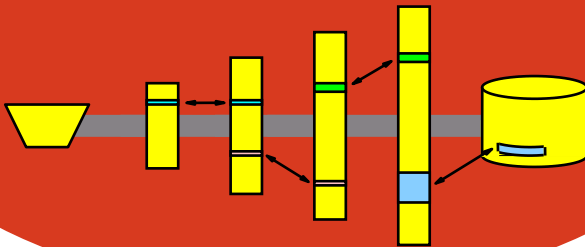
I/O Effektive
Algoritmer



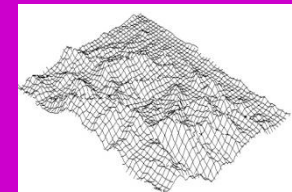
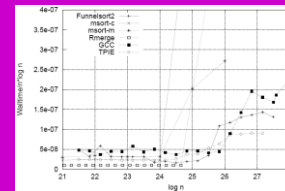
Streaming
Algoritmer



Cache
Oblivious Algoritmer

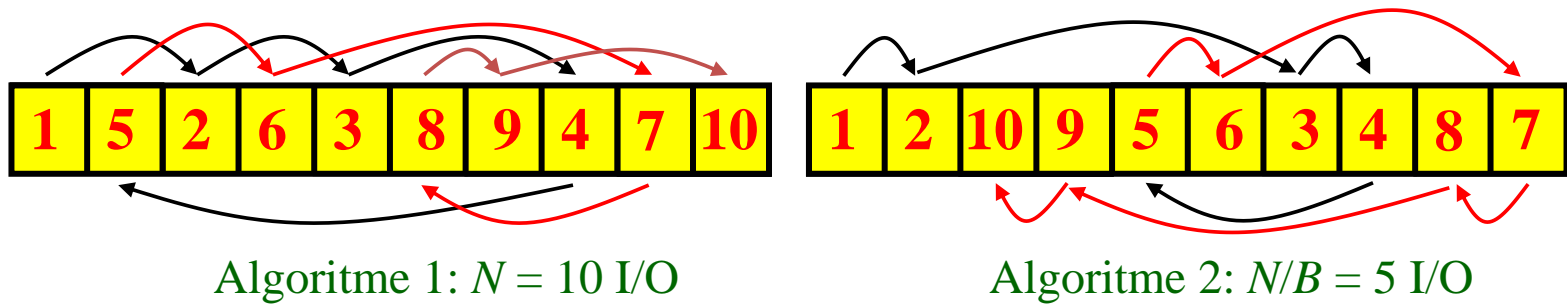


Algoritme
Engineering



I/O-Effektive Algoritmer Gør Forskel

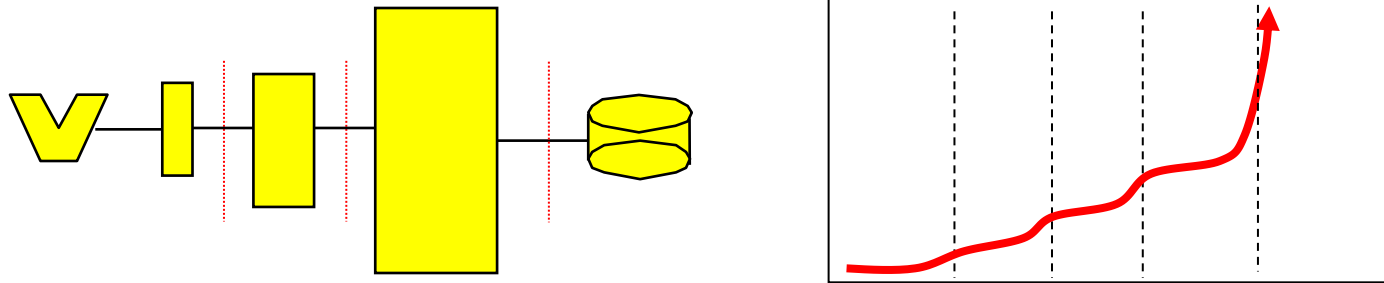
- **Eksempel:** Gennemløb en kædet liste
 - Problemstørrelse $N = 10$ elementer
 - Diskblokstørrelse $B = 2$ elementer
 - Hukommelsesstørrelse $M = 4$ elementer (2 blokke)



- Forskellen mellem N og N/B stor da blokstørrelsen er stor
 - **Eksempel:** $N = 256 \times 10^6$, $B = 8000$, 1ms disk tilgang
 - $\Rightarrow N$ I/O tager 256×10^3 sek = 4266 min = **71 timer**
 - $\Rightarrow N/B$ I/O tager $256/8$ sek = **32 sek**

Cache-Oblivious Algoritmer

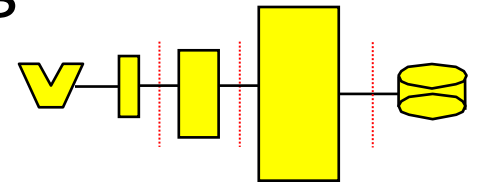
- Hvis problemer skal løses på ukendte eller ændrende maskiner



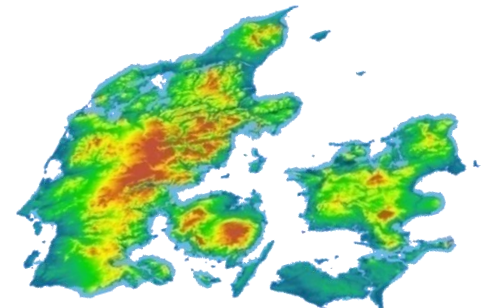
- Bloktilgang er vigtig på alle niveauer i hukommelseshierarkiet
 - Men hukommelseshierarkier er meget varierende
- Cache-oblivious algoritmer:
 - Brug alle blokke effektivt på *alle* niveauer i *ethvert* hukommelseshierarki

Algoritme Engineering

- Design/implementation af praktiske algoritmer
 - Eksperimenter
 - Center motiveret ved teoretiske mangler
 - Center promoverer interdisciplinær/industrielt arbejde
- Naturligt at lave algorithm engineering



- Algorithm engineering
 - Ofte værdifuld input til teoretisk arbejde
 - f.eks. til design af bedre beregningsmodeller
 - Nogle gange praktiske gennembrud
 - f.eks. MADALGOs terræn data implementation



Hurricane Floyd

Sep. 15, 1999



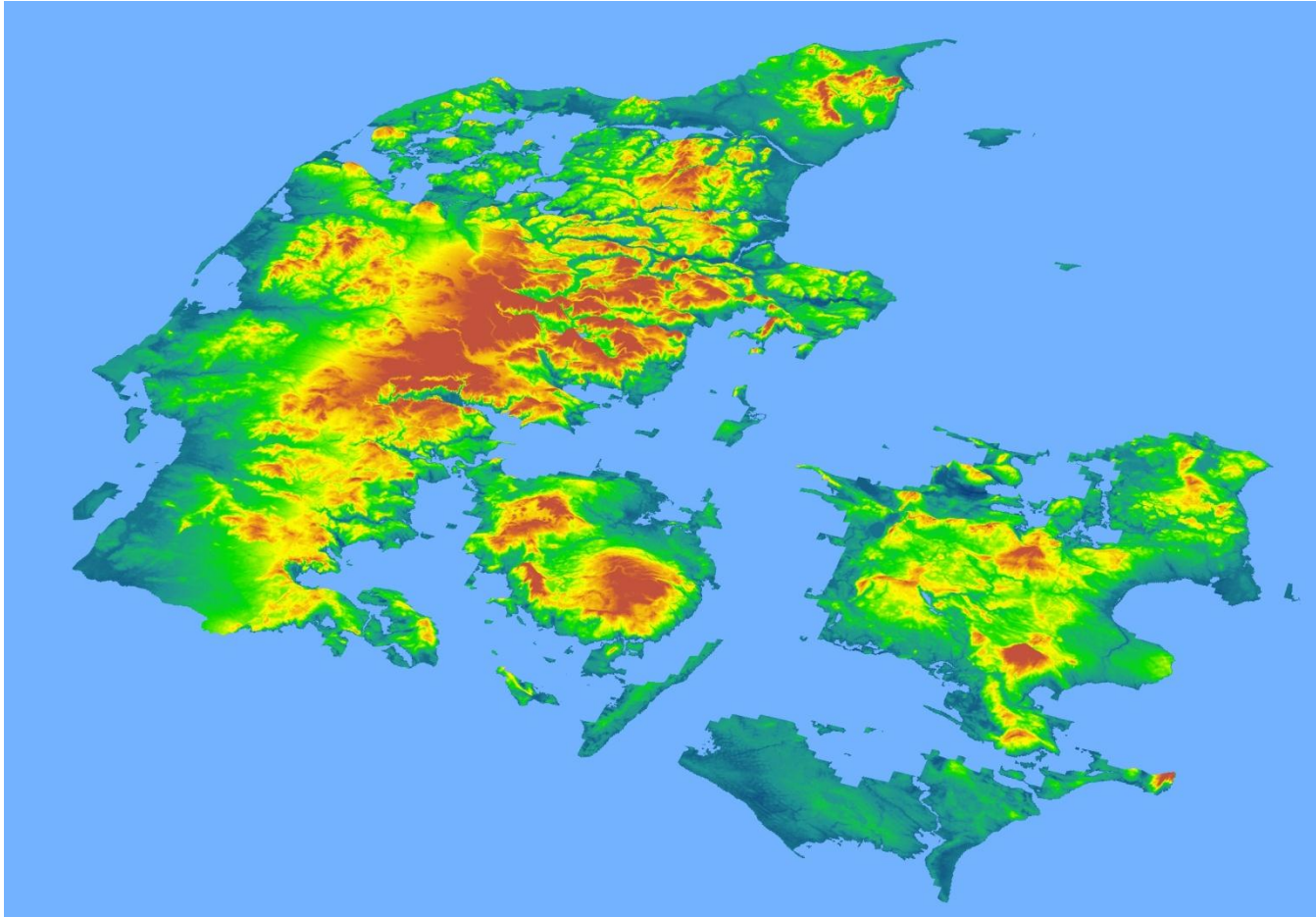
7:00



15:00

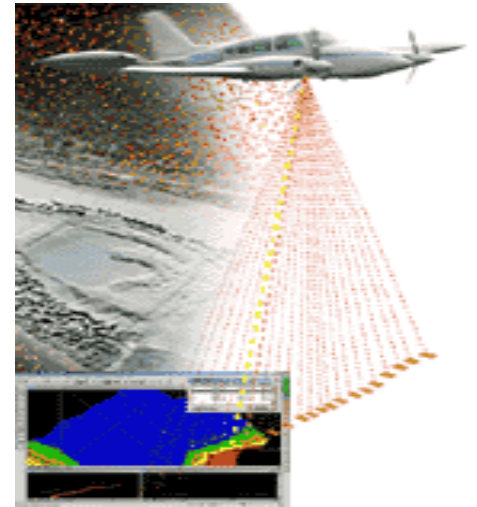


Massive Terræen Data



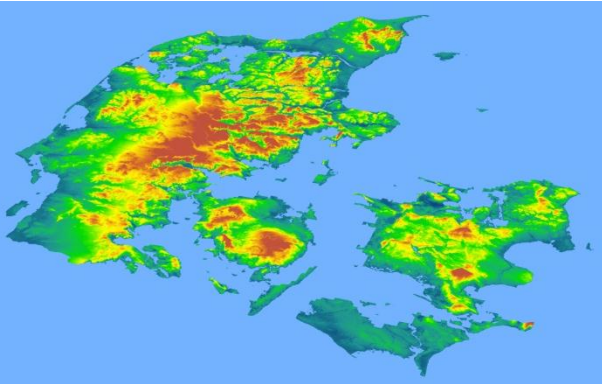
Terræn Data

- **Nye teknologier:** Meget nemmere/billigere at indsamle detaljeret data
- **Før** 'manuel' eller radarbaserede metoder
 - Ofte 30 meter mellem datapunkter
 - Nogle gange 10 meter data tilgængelig
- **Nye** laserskannings metoder (LIDAR)
 - Mindre end 1 meter mellem datapunkterne
 - Præcision i centimeter (hidtil meter)



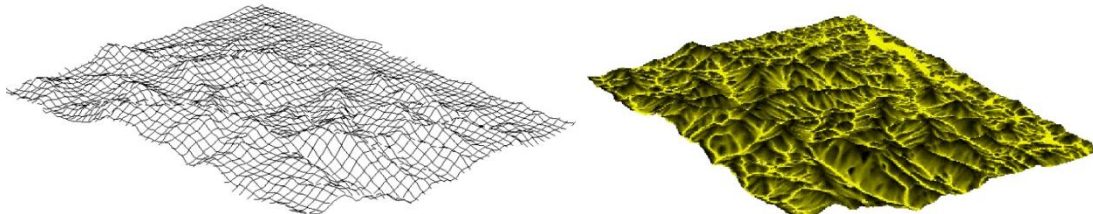
Danmark

- ~2 millioner punkter v. 30 meter ($\ll 1\text{GB}$)
- ~18 milliarder punkter v. 1 meter ($\gg 1\text{TB}$)
- COWI (og andre) skanner nu DK
- NC skannet efter Hurricane Floyd i 1999



Simulering af Oversvømmelse

- Ikke alt terræn over højde h bliver oversvømmet når vandet stiger h meter!
- Teoretisk ikke så hårdt at beregne områderne der oversvømmes når vandet stiger h meter
 - Men ingen software kunne gøre det for Danmark ved 2-meter opløsning
- Anvend I/O-effektiv algoritme
⇒ Danmark på én dag



Oversvømmelse af Danmark

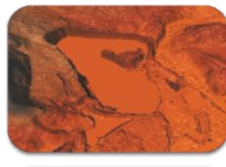


TerraStream Terræn Software

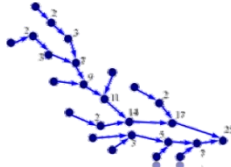
- Skalerbar, generelt, porterbar, håndtering af datafejl
- Pipeline af massive terræn data processerings software



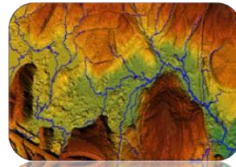
DEM Construction



Conditioning



Flow Routing



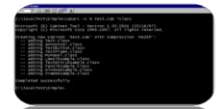
Flow Accumulation



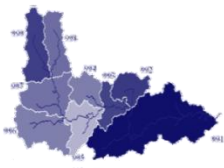
ArcGIS Extension



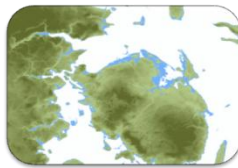
GRASS Extension



Command Line Tools



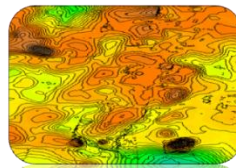
Watershed Hierarchies



Flood Simulation



Quality Metrics



Contour Lines



Front-End GUI

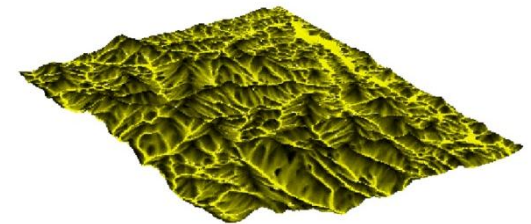
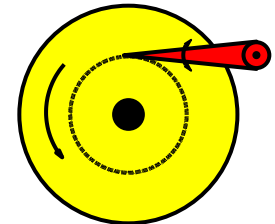
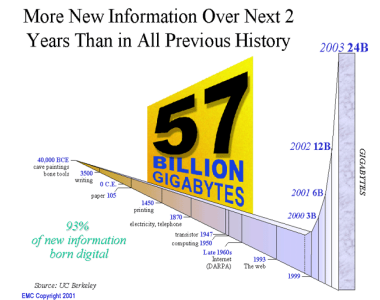


MapInfo

Demo: www.madalgo.au.dk/~thomasm/floodmaps/?extra=d2m

Opsummering

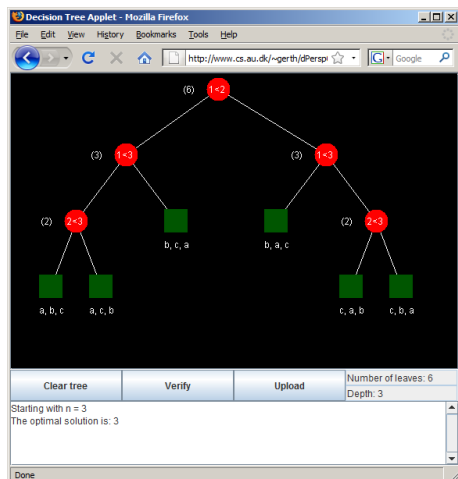
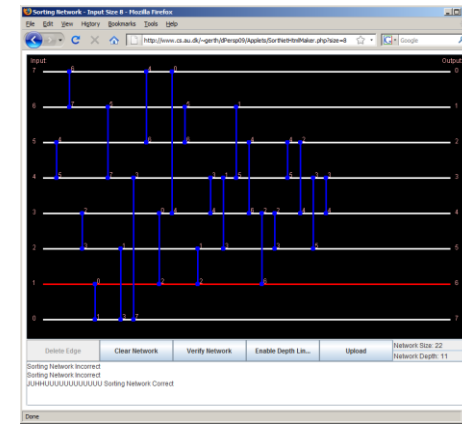
- Massive datamængder forekommer overalt
- Medfører skaleringsproblemer
 - pga. hierarkisk hukommelse og langsom I/O
- I/O-effektive algoritmer giver en signifikant forbedring af skalerbarhed
- Nyt forskningscenter fokuserer på algoritmiske emner inden for massive data



Vidensspredning

Vidensspredning

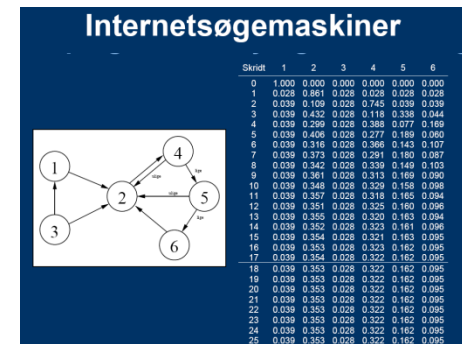
- Folkeskolepraktik
- Gymnasiepraktik
- UNF foredrag om internetsøgemaskiner
- ACM programmeringskonkurrence
- Perspektiverende datalogi kursus



NWERC 2008

#	AFFIL.	TEAM	SCORE	A +	B	C +	D +	E +	F +	H +	I +	J +	K +
1	Denmark	Mads, Vens & Sisp	6	1,207 + 33	1,134 + 36	1,96 + 10	1,11 + 10	0	2,113 + 26	1,67 + 10	3,180 + 43	1,21 + 10	1,128 + 31
2	Denmark	PH	6	4,277 + 30	1,076 + 30	1,88 + 10	0	0	2,69 + 20	1,132 + 23	2,113 + 23	0	0
3	Denmark	Pierre Suspects	6	4,485 + 35	0	1,83 + 10	0	0	3,275 + 35	4,481 + 39	1,126 + 19	1,59 + 10	0
4	Denmark	MADALOG Men	6	1,207 + 33	1,134 + 36	2,146 + 30	0	0	1,30 + 10	1,66 + 10	2,144 + 30	0	0
5	Denmark	The Underdogs	6	1,870 + 2,249 + 250	4,252 + 35	1,109 + 10	0	0	1,162 + 10	1,193 + 10	1,106 + 10	0	0
6	TU/e	@Net code edition	6	1,239 + 2,248 + 30	0	0	0	0	3,200 + 30	3,131 + 30	1,163 + 30	1,57 + 10	0
7	TU/e	Jay	6	1,1139 + 4,036 + 60	5,281 + 100	1,159 + 10	0	0	2,28 + 30	1,107 + 10	1,106 + 10	0	0
8	TU/e	Jacobs University	6	1,236 + 5,263 + 10	4,252 + 35	2,207 + 25	3,118 + 40	3	0	2,44 + 30	1,150 + 30	6,189 + 130	0
9	CU	Team Squared	6	1,1362 + 1,179 + 30	0	2,219 + 30	0	0	4,254 + 30	1,52 + 10	1,207 + 10	1,261 + 10	0
10	Denmark	Borsten Fan Club	5	2,757 + 0	0	1,252 + 10	1,63 + 10	0	3,224 + 30	4,63 + 30	1,53 + 10	0	0

- 5 timers konkurrence, 10 opgaver stillet.
- 2 hold fra DAIMI, blev nr. 4 og 5 ud af 47.
- 2 bedste hold videre til World Finals i Stockholm.



Tak !

Spørgsmål, hvad nu?