# ALGORITHMS FOR
# MASSIVE TERRAINS AND GRAPHS

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

SVEND CHRISTIAN SVENDSEN

# — THE PROGRAM OF THE DAY

— External Memory Pipelining Made Easy With TPIE

Lars Arge, Mathias Rav, Svend C. Svendsen, Jakob Truelsen

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# THE PROGRAM OF THE DAY

— External Memory Pipelining Made Easy With TPIE

Lars Arge, Mathias Rav, Svend C. Svendsen, Jakob Truelsen

— 1D and 2D Flow Routing on a Terrain

Lars Arge, Aaron Lowe, Svend C. Svendsen, Pankaj K. Agarwal

# THE PROGRAM OF THE DAY

— External Memory Pipelining Made Easy With TPIE

Lars Arge, Mathias Rav, Svend C. Svendsen, Jakob Truelsen

— 1D and 2D Flow Routing on a Terrain

Lars Arge, Aaron Lowe, Svend C. Svendsen, Pankaj K. Agarwal

— Practical I/O-Efficient Multiway Separators

Svend C. Svendsen

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — THE PROGRAM OF THE DAY

— External Memory Pipelining Made Easy With TPIE

Lars Arge, Mathias Rav, Svend C. Svendsen, Jakob Truelsen

— 1D and 2D Flow Routing on a Terrain

Lars Arge, Aaron Lowe, Svend C. Svendsen, Pankaj K. Agarwal

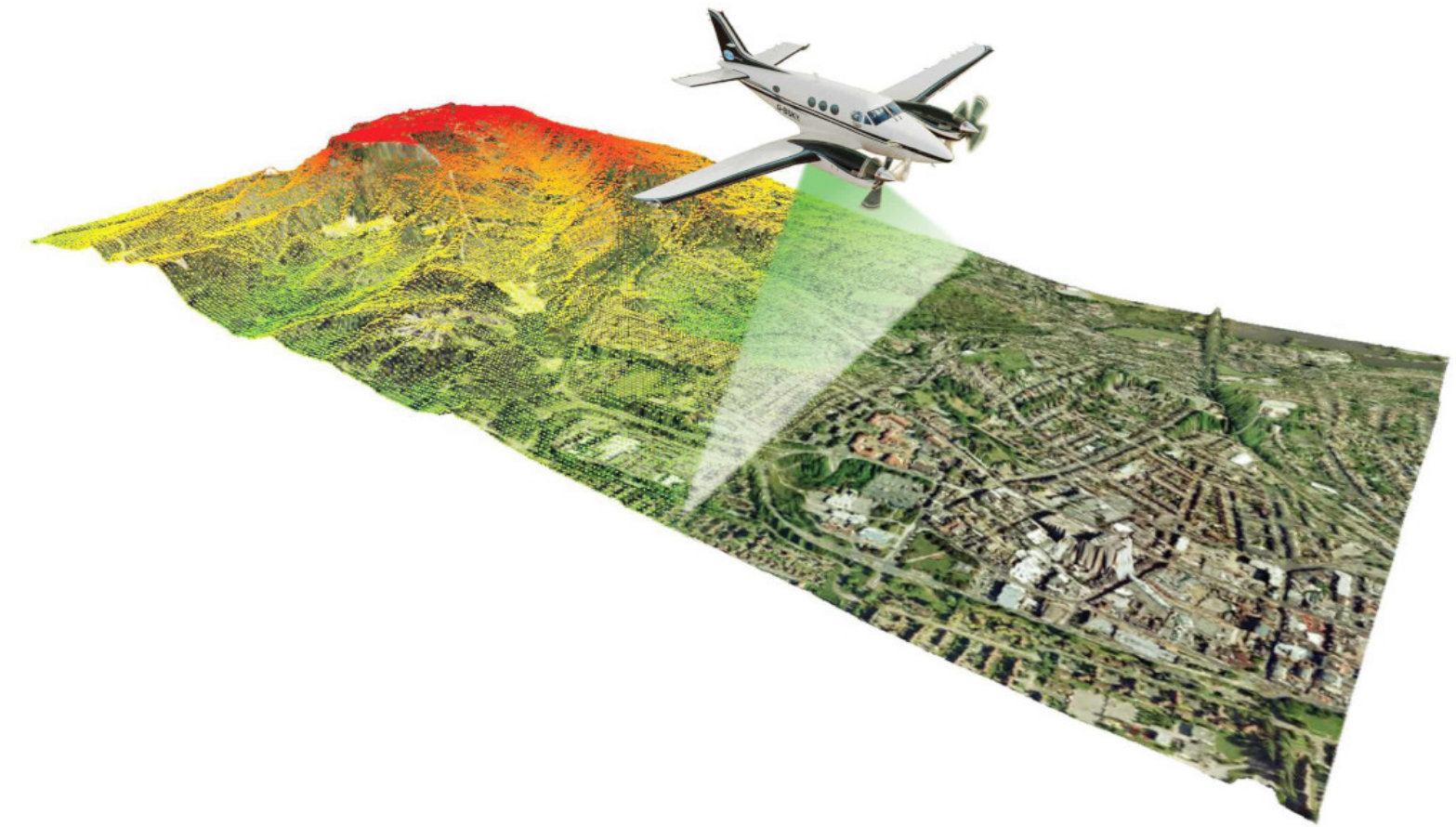— Practical I/O-Efficient Multiway Separators

Svend C. Svendsen

— Learning to Find Hydrological Corrections

Lars Arge, Allan Grønlund, Svend Christian Svendsen, Jonas Tranberg

AARHUS
UNIVERSITY
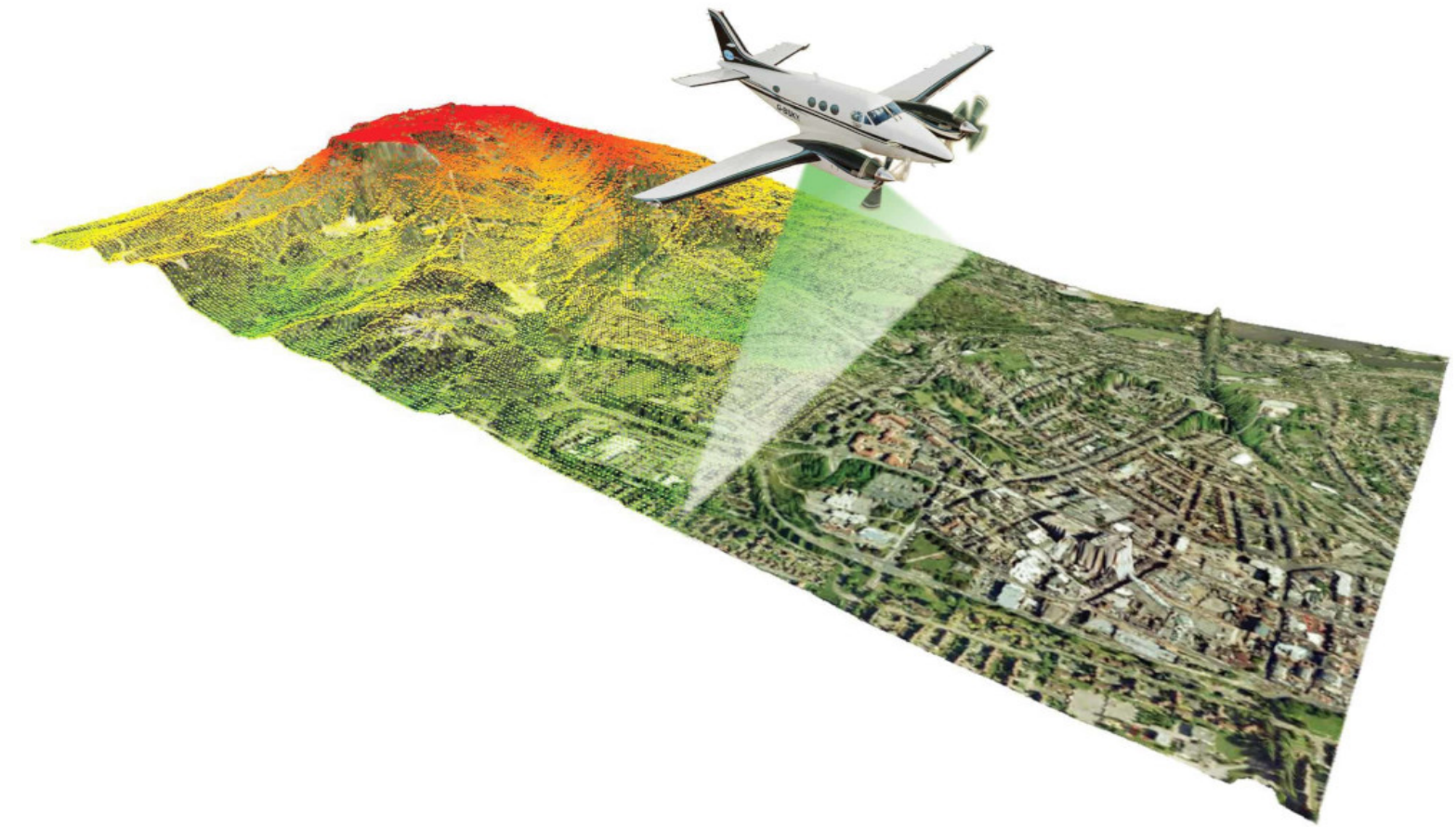DEPARTMENT OF COMPUTER SCIENCE

# — TERRAIN AND BIG DATA

— Present: Terrain is collected with LiDAR
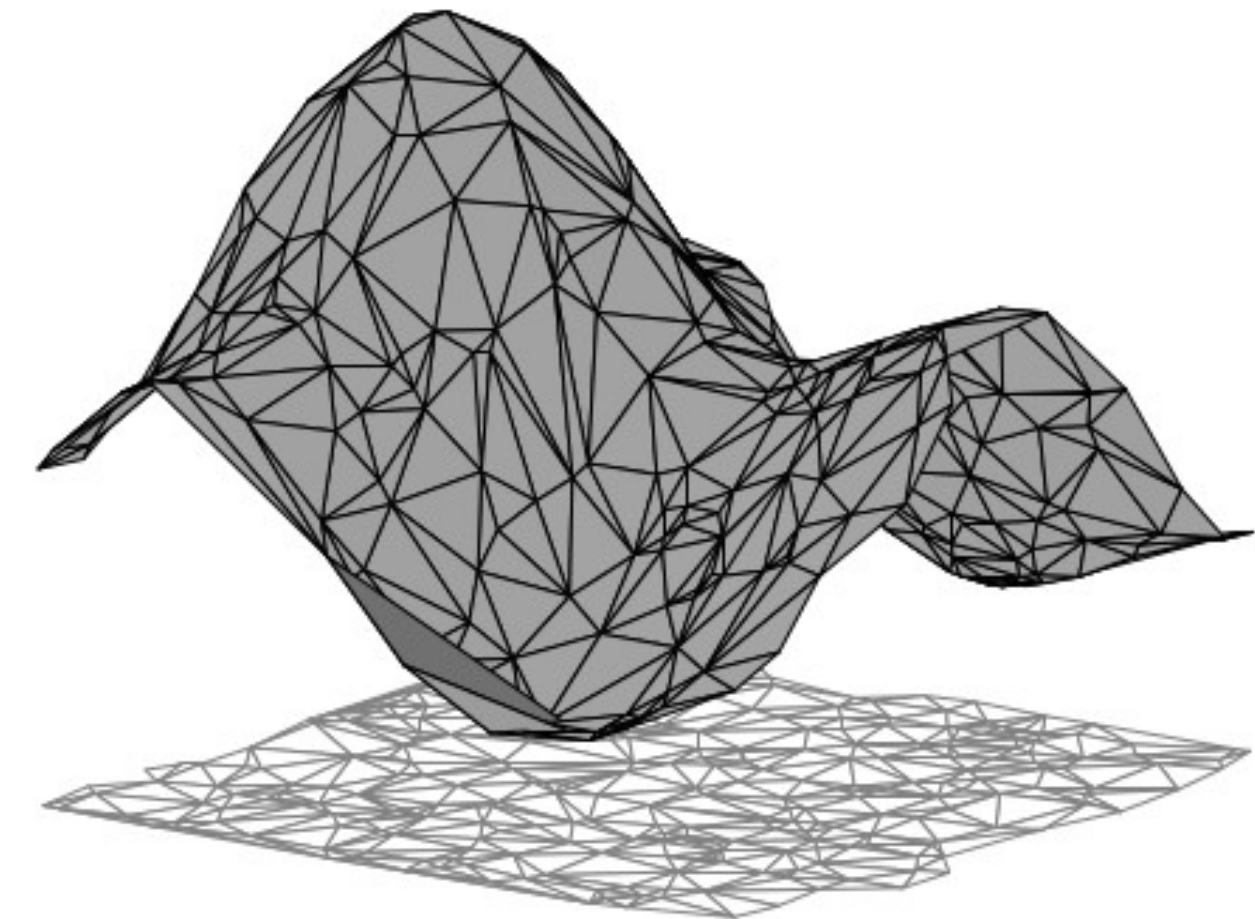
Source: LiDAR America

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# TERRAIN AND BIG DATA

— Present: Terrain is collected with LiDAR

— Denmark - Shuttle Radar Topography Mission
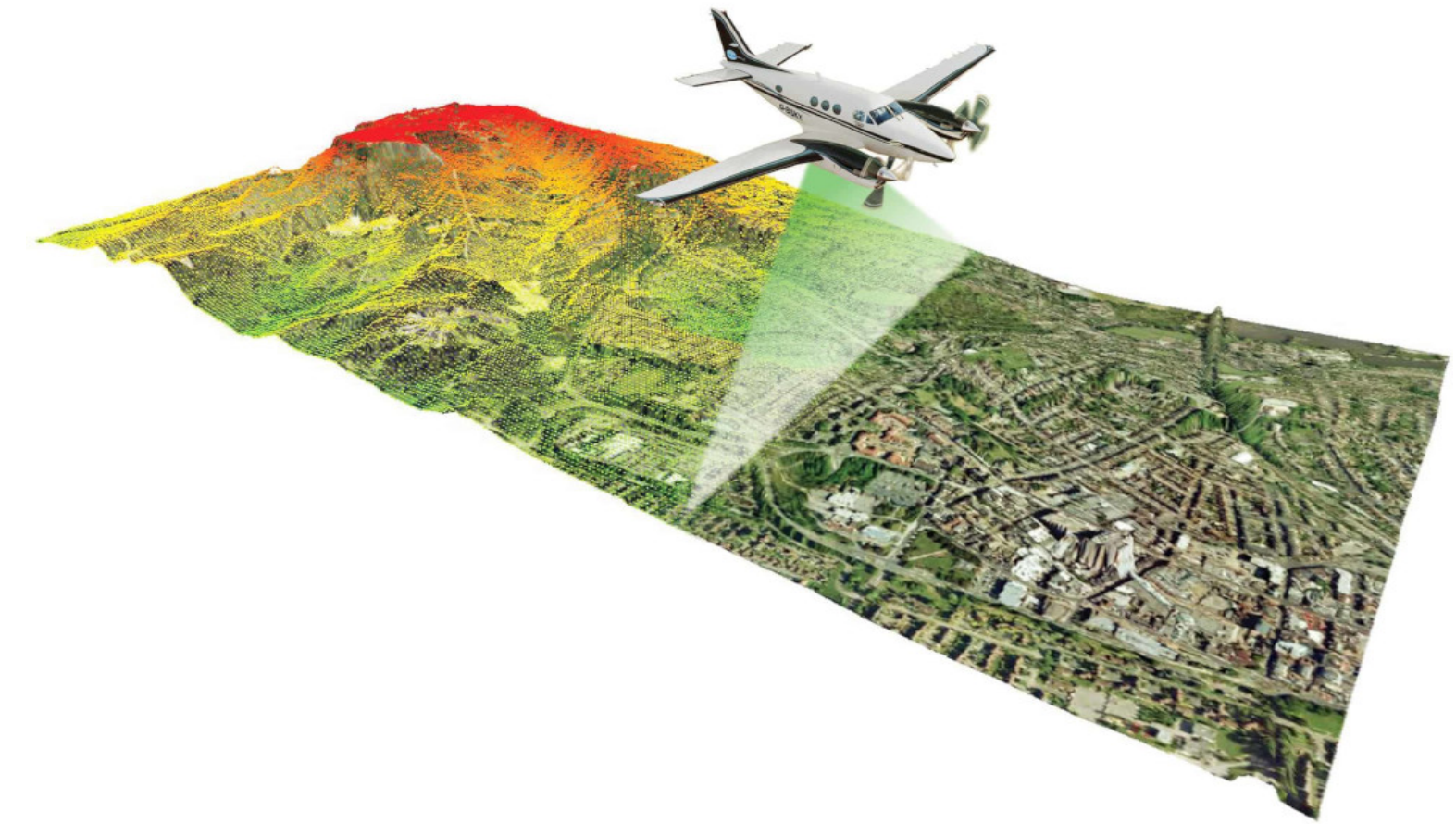  90 meter resolution
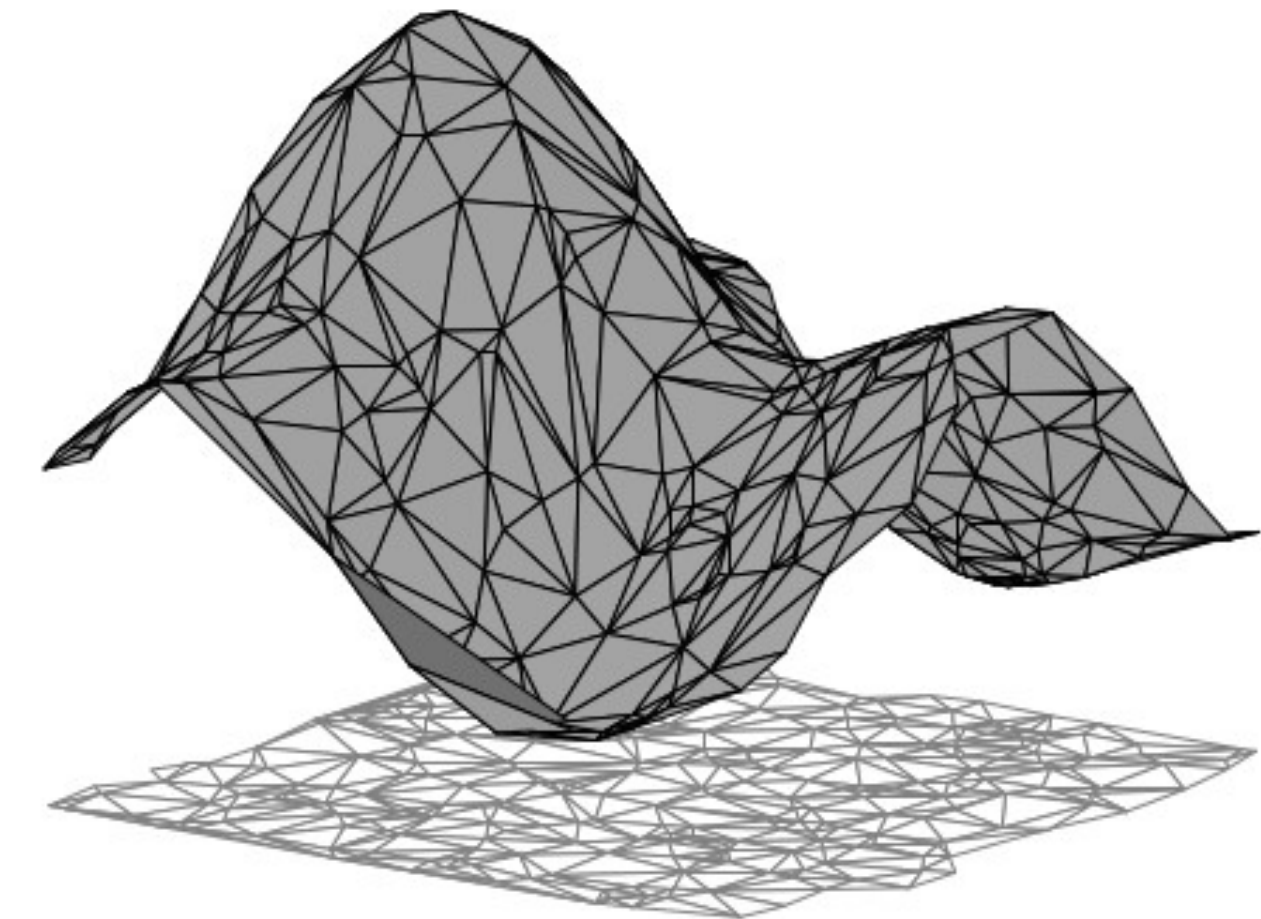  4,000,000 points

Source: LiDAR America

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# TERRAIN AND BIG DATA

— Present: Terrain is collected with LiDAR

— Denmark - Shuttle Radar Topography Mission

   90 meter resolution
   $4,000,000$ points

— Denmark - Danish Elevation Model

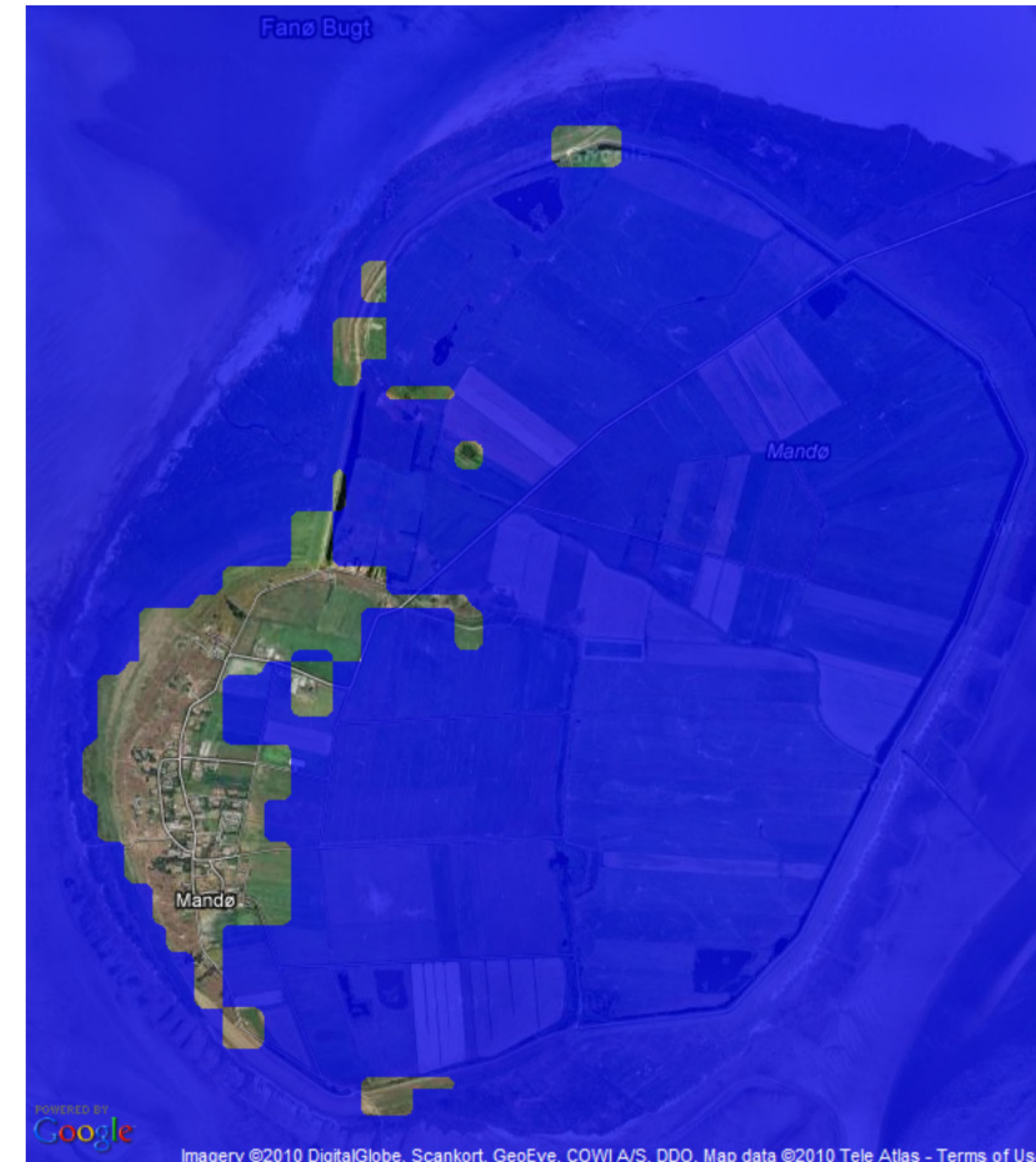   40 centimeter resolution
   $415,000,000,000$ points

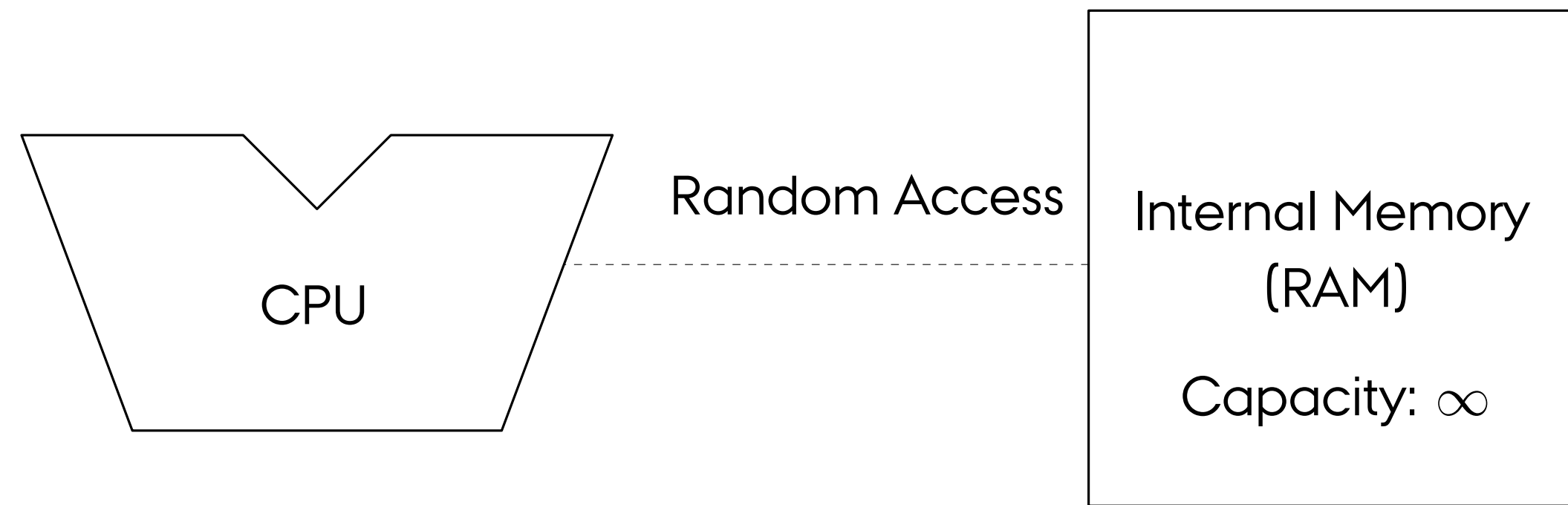Source: LiDAR America

# — TERRAIN AND BIG DATA



2 meter resolution

90 meter resolution

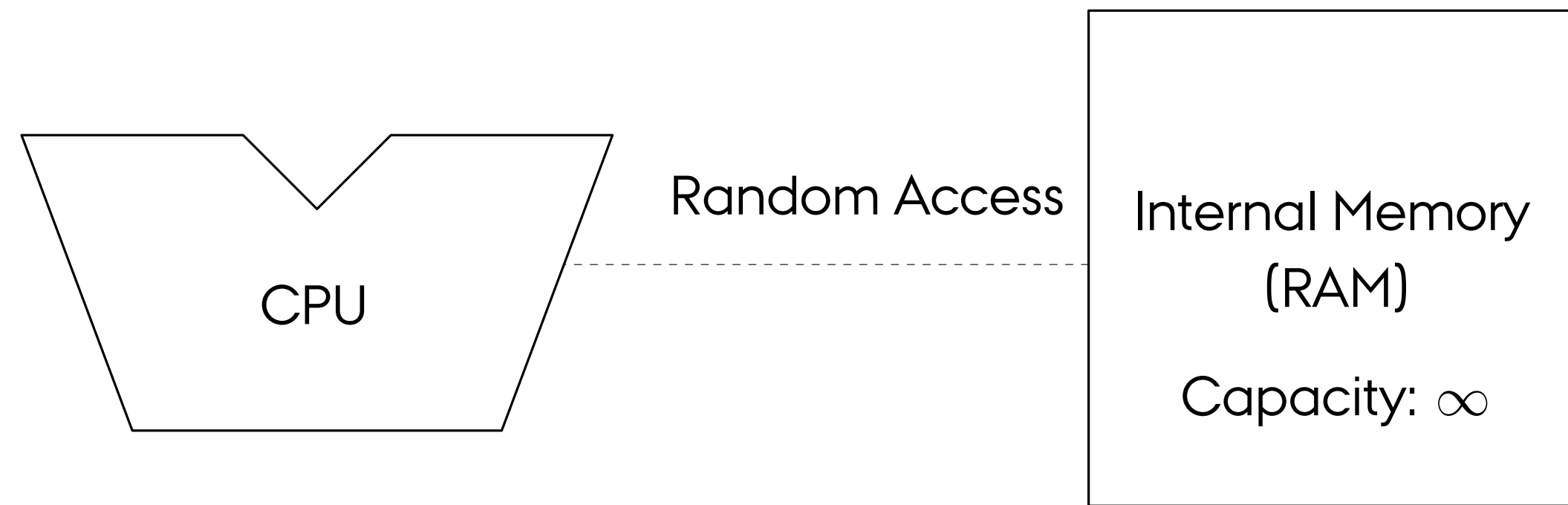Source: Scalable algorithms for large high-resolution terrain data, Mølhave et al.

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — I/O-EFFICIENT ALGORITHMS

— RAM model



CPU

Random Access

Internal Memory
(RAM)

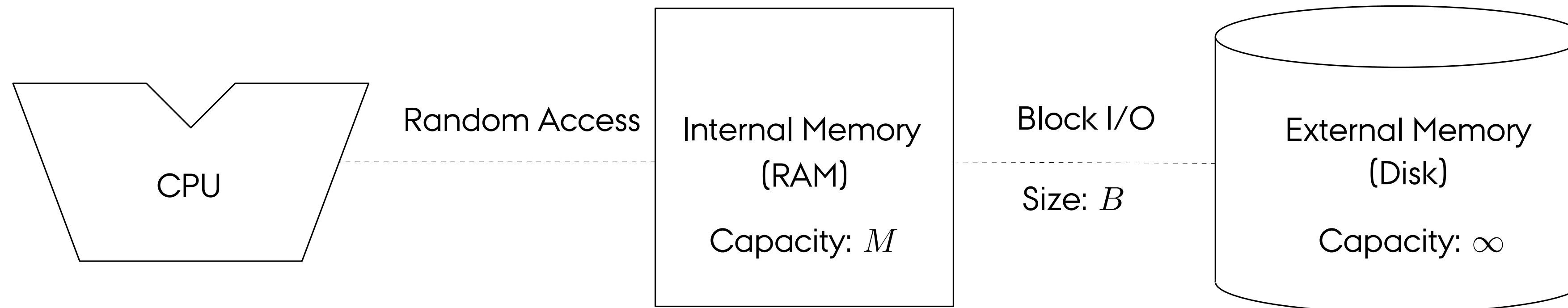Capacity: $\infty$

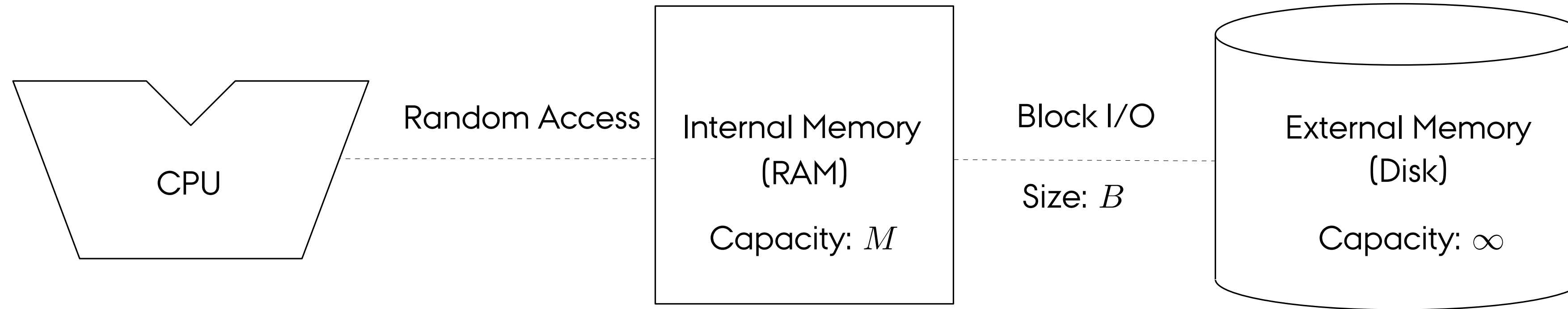# — I/O-EFFICIENT ALGORITHMS

— RAM model



— I/O-Model



— Hard drives moves blocks of data and are slow

— I/O-Efficient Algorithms: Move as few blocks as possible

# — I/O-EFFICIENT ALGORITHMS



— I/O-Model by Aggarwal and Vitter (CACM 1988)

— $N$ = # of items in input

— $B$ = # of items in a block

— $M$ = # of items in memory (capacity)

— Reading elements: $\text{Scan}(N) = \Theta(N/B)$

— Sorting elements: $\text{Sort}(N) = \Theta(\frac{N}{B} \log_{M/B} \frac{N}{B})$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# External Memory Pipelining Made Easy With TPIE

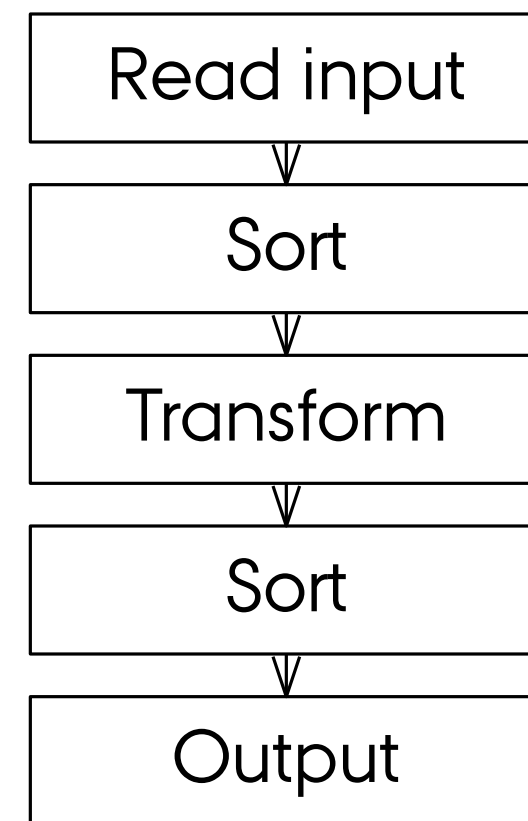Lars Arge, Mathias Rav, Svend C. Svendsen, Jakob Truelsen

IEEE BigData 2017

# — I/O-EFFICIENT ALGORITHMS IN PRACTICE

— TPIE: The Templated Portable I/O Environment

— Hide low-level details while maintaining performance

— File streams: reading and writing to disk

— Provides implementations of fundamental algorithms
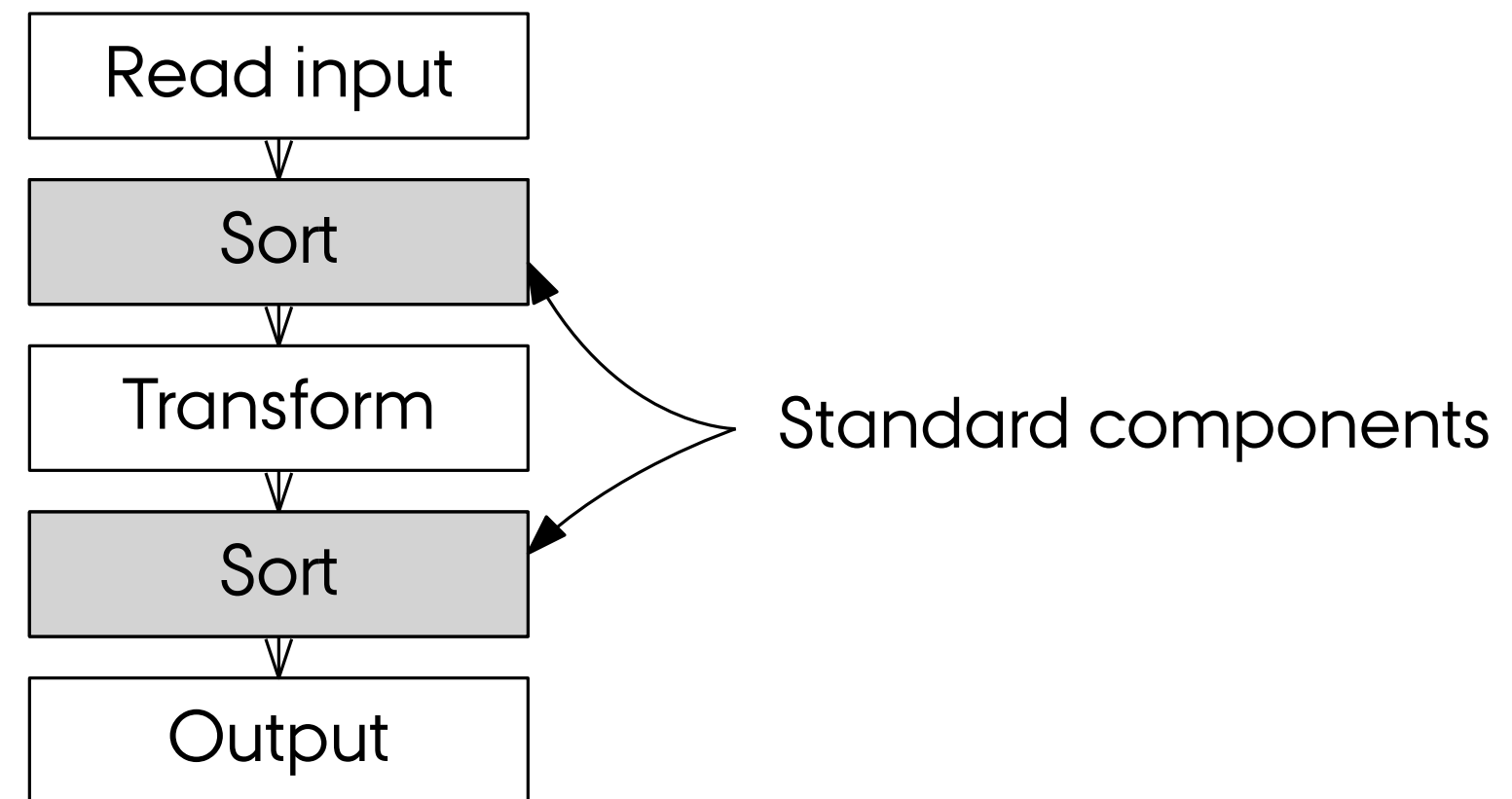
— Used both commercially and in research

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# TPIE PIPELINING

Imperative-style algorithm

```
┌─────────────┐
│ Read input  │
└─────────────┘
       ∨
┌─────────────┐
│    Sort     │
└─────────────┘
       ∨
┌─────────────┐
│  Transform  │
└─────────────┘
       ∨
┌─────────────┐
│    Sort     │
└─────────────┘
       ∨
┌─────────────┐
│   Output    │
└─────────────┘
```

# — TPIE PIPELINING

Imperative-style algorithm

Imperative-style algorithm

| | |
|---|---|
| Read input | $O(\text{Scan}(N))$ |
| Sort | $O(\text{Sort}(N))$ |
| Transform | $O(\text{Scan}(N))$ |
| Sort | $O(\text{Sort}(N))$ |
| Output | $O(\text{Scan}(N))$ |

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

Imperative-style algorithm

| | |
|---|---|
| Read input | $O(\mathrm{Scan}(N))$ |
| Sort | $O(\mathrm{Sort}(N))$ |
| Transform | $O(\mathrm{Scan}(N))$ |
| Sort | $O(\mathrm{Sort}(N))$ |
| Output | $O(\mathrm{Scan}(N))$ |

$$\overline{O(\mathrm{Sort}(N))}$$

# ─ TPIE PIPELINING

Imperative-style algorithm                Pipelined Algorithm

| Read input | $O(\text{Scan}(N))$ |

| Sort | $O(\text{Sort}(N))$ |

| Transform | $O(\text{Scan}(N))$ |

| Sort | $O(\text{Sort}(N))$ |

| Output | $O(\text{Scan}(N))$ |

$$O(\text{Sort}(N))$$

Read input

Sort

Transform

Sort

Output

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

SOLIDUM PETIT IN PROFUNDIS · UNIVERSITAS ARHUSIENSIS ·
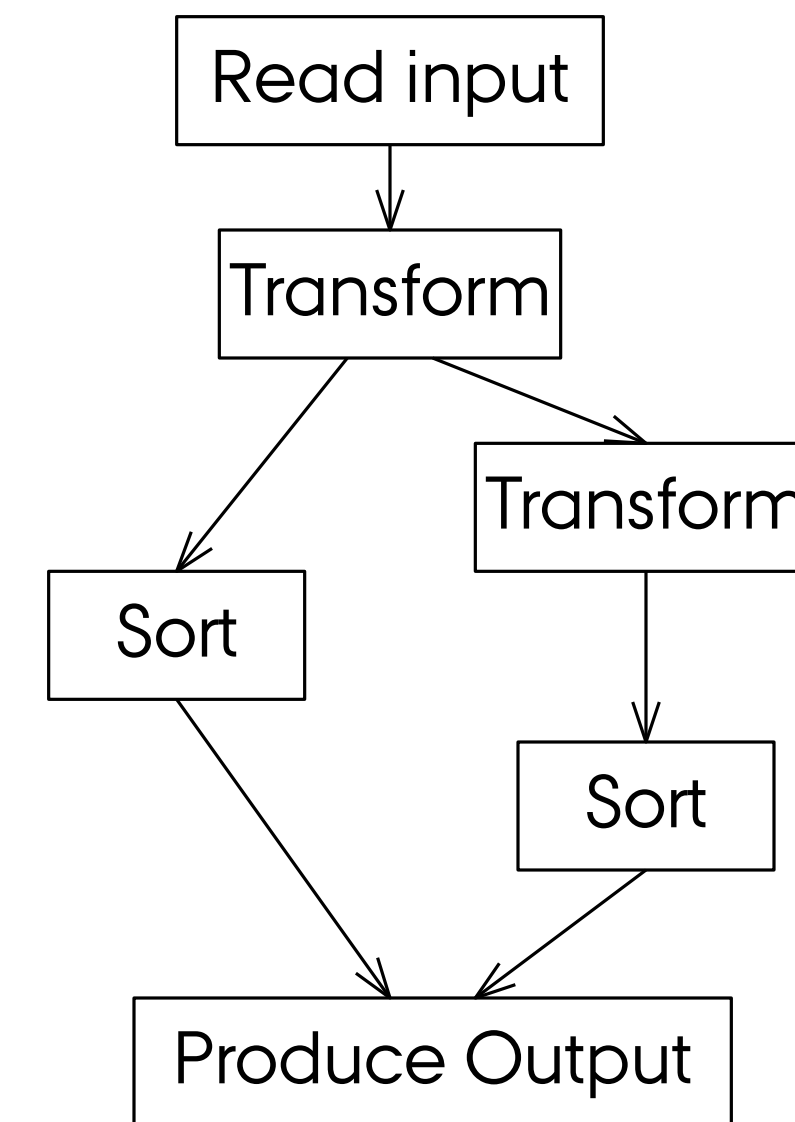
```
def transform(input_file, output_file):
    while input_file.can_read():
        x = input_file.read()
        output_file.write(f(x))
```

```
class TransformComponent:
    def push(x):
        dest.push(f(x))
```
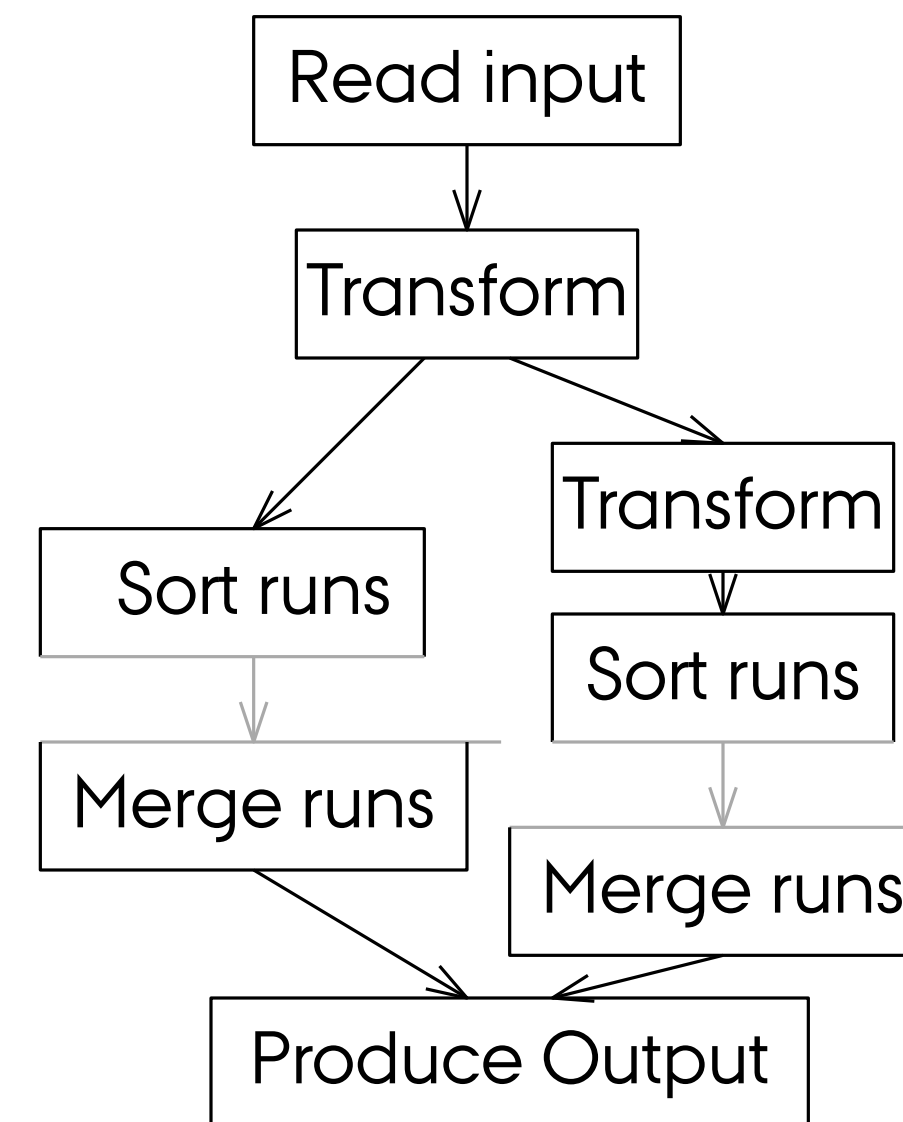
AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — **TPIE PIPELINING**

— Blocking Components

— Identifying Phases

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# TPIE PIPELINING

— Blocking Components

— Identifying Phases

AARHUS
UNIVERSITY
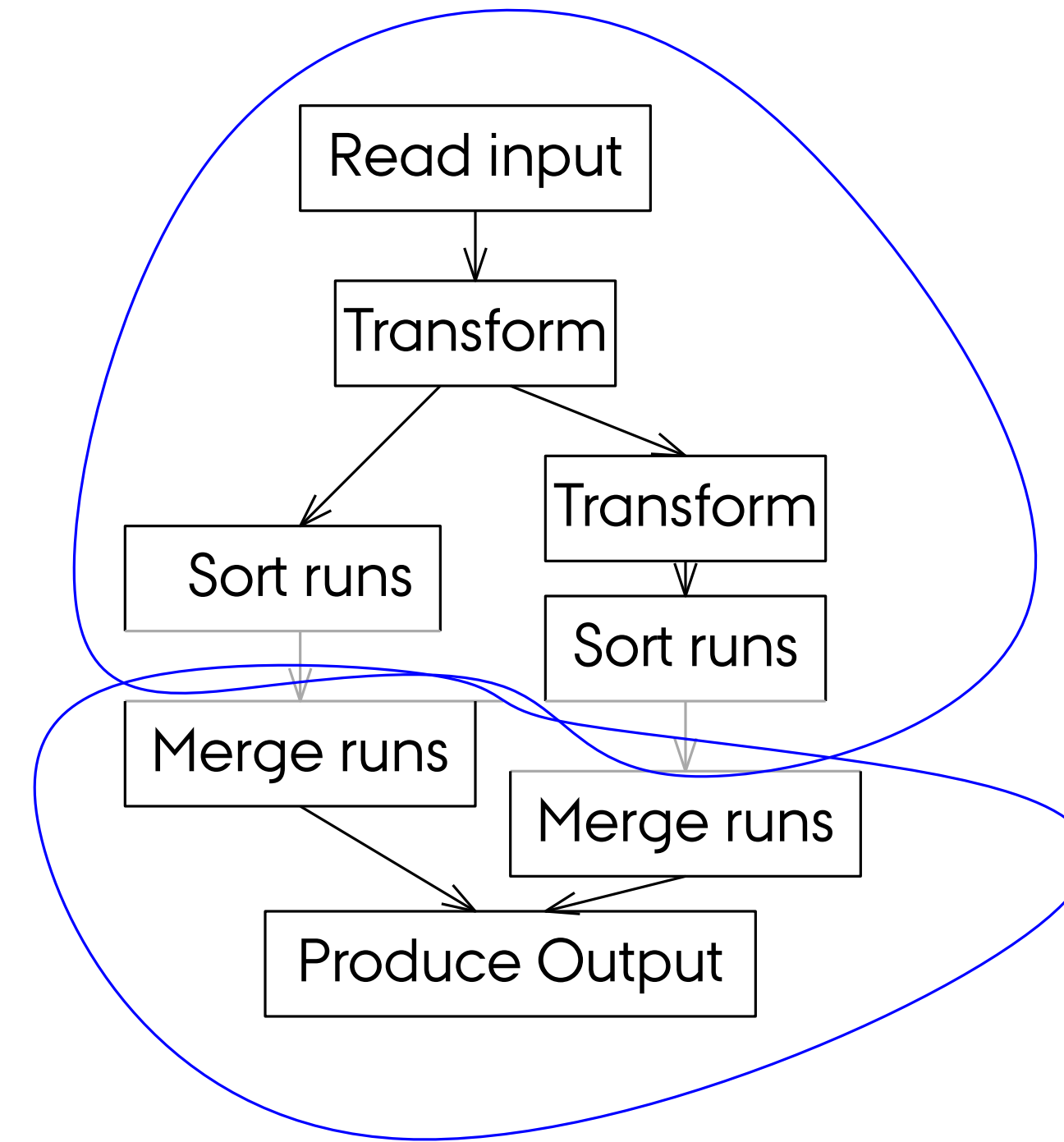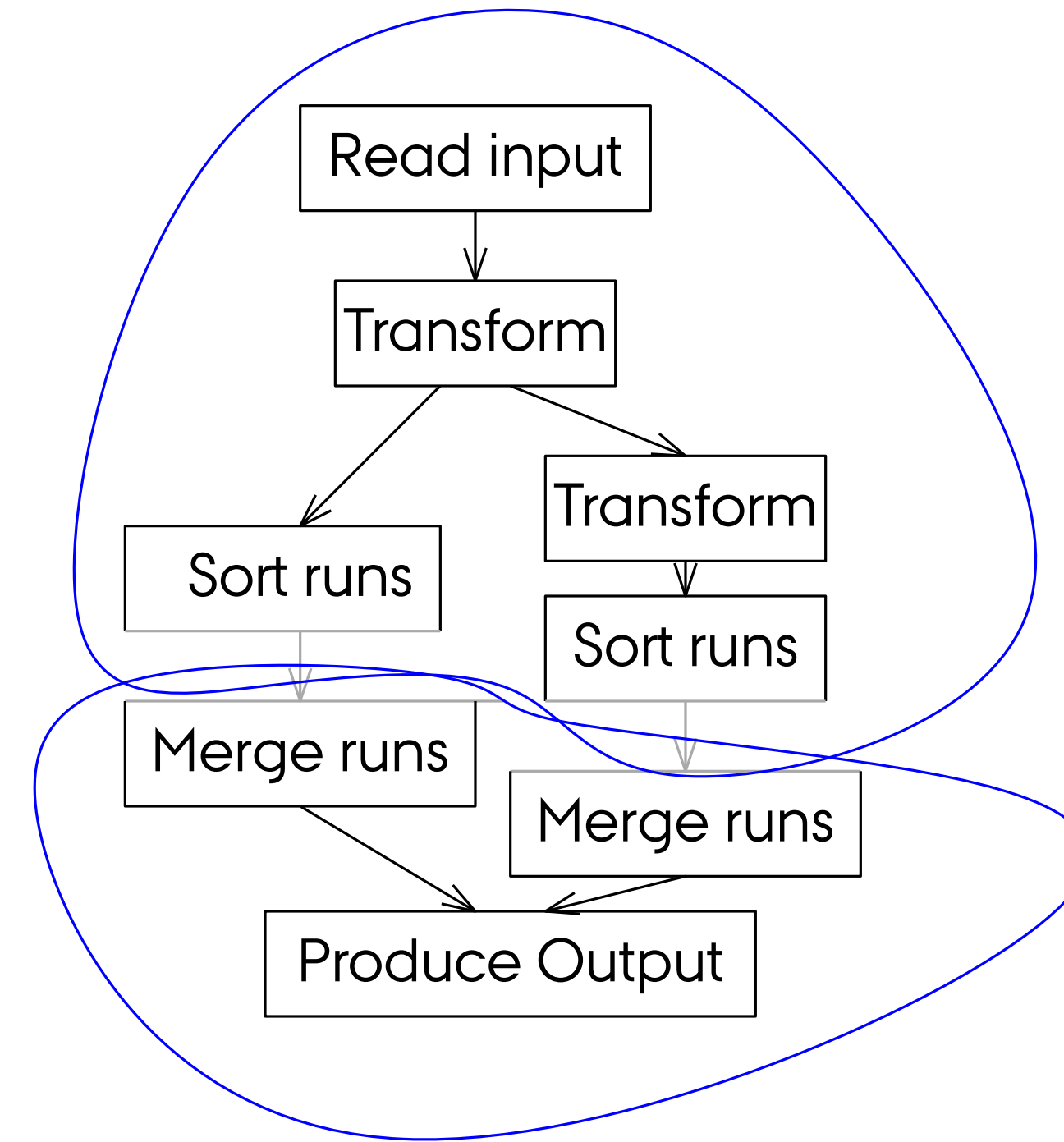DEPARTMENT OF COMPUTER SCIENCE

# TPIE PIPELINING

— Blocking Components

— Identifying Phases

— Memory Management

# TPIE PIPELINING

- Blocking Components

- Identifying Phases

- Memory Management

- Parallelisation

- Progress Tracking
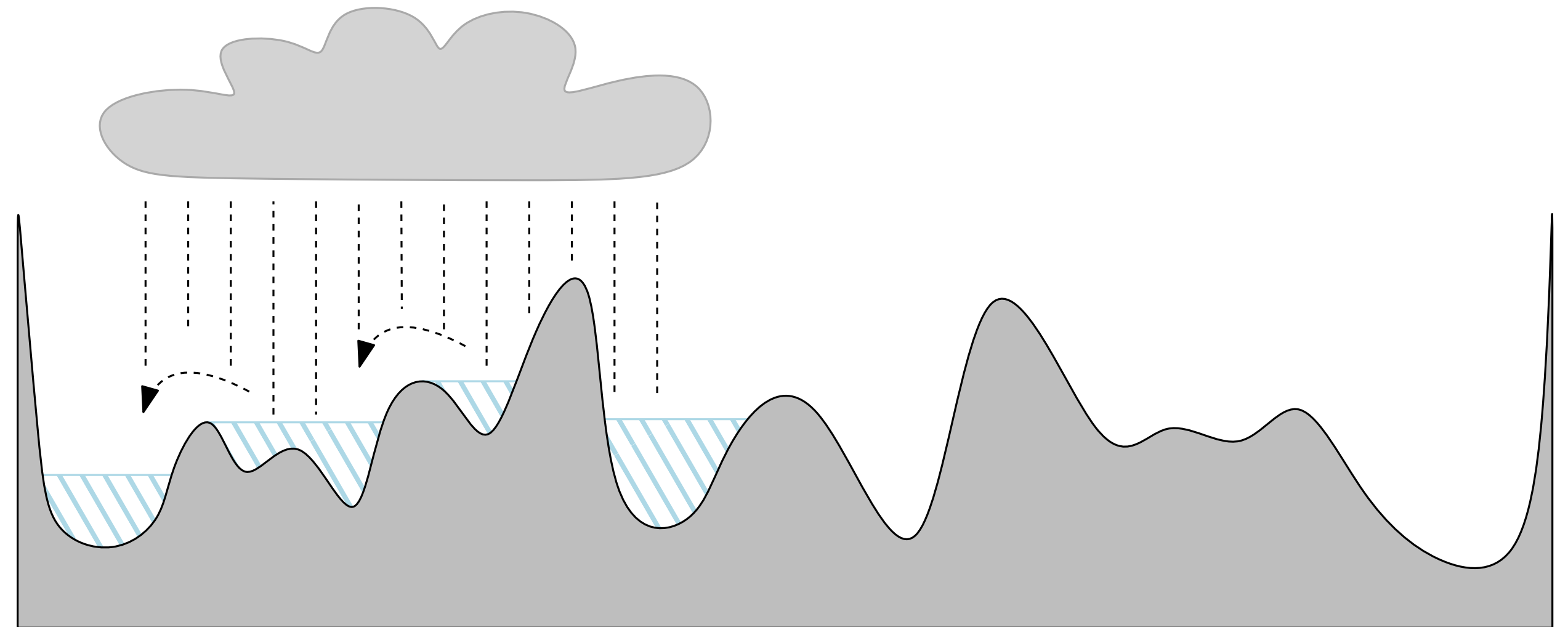
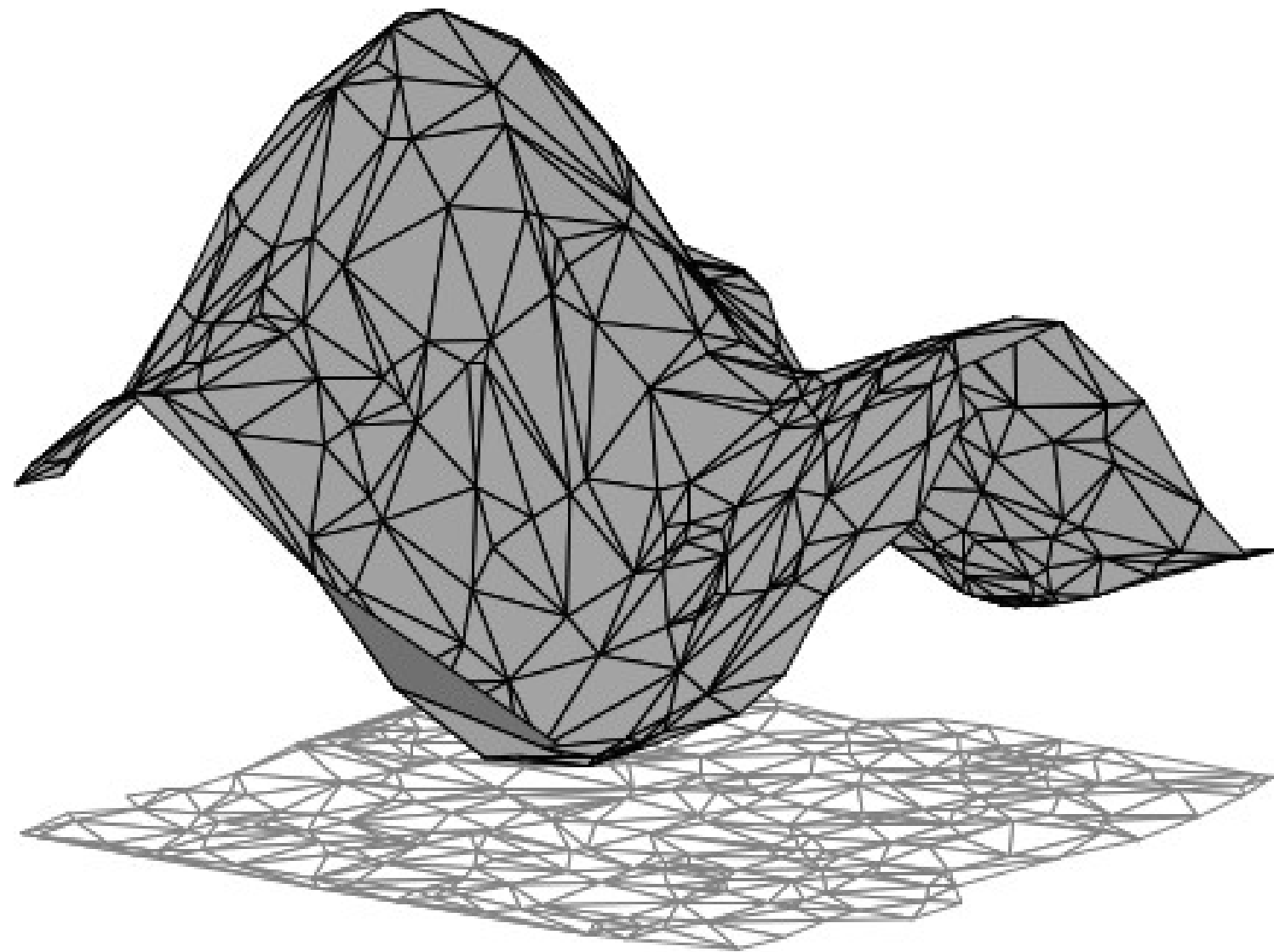# 1D and 2D Flow Routing on a Terrain

Lars Arge, Aaron Lowe, Svend C. Svendsen, Pankaj K. Agarwal

ACM SIGSPATIAL 2020
Invited to ACM TSAS

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — FLOOD MODEL

— Single Flow Direction Model: Water on a vertex $v$ flows to a single neighbor $u$ along an edge

— Multiflow Direction Model: Water on a vertex $v$ flows to multiple neighbors

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — THE PROBLEM

— Rain distribution: $\mathcal{R}(v) : \mathbb{V} \to \mathbb{R}_{\geq 0}$

— Terrain-flood query: GIven a rain distribution $\mathcal{R}$ and a time $t$, determine which vertices of $\Sigma$ are flooded.

— Flood-time query: GIven a rain distribution $\mathcal{R}$, for each vertex $q \in \Sigma$, determine the time $t$ that $q$ becomes flooded

# — STATE OF THE ART

— H: height of the merge tree

— X: number of depressions

|  | Flood-time query | Terrain-flood query |  |
|---|---|---|---|
|  |  |  | [1] 2004, Liu and Snoeyink |
| SFD RAM-model | $O(N \log N)$ [1] | $O(N \log N)$ [1] | [2] 2010, Arge, Revsbæk, Zeh |
| SFD I/O-model | $O(\text{Sort}(X) \log \frac{X}{M} + \text{Sort} N)$ [2] | $O(\text{Sort}(N) + \text{Scan}(H \cdot X))$ [3] | [3] 2017, Arge, Rav, Raza, Revsbæk |

# — STATE OF THE ART

— H: height of the merge tree

— X: number of depressions

|  | Flood-time query | Terrain-flood query | |
|---|---|---|---|
| SFD RAM-model | $O(N \log N)$ [1] | $O(N \log N)$ [1] | [1] 2004, Liu and Snoeyink |
|  |  |  | [2] 2010, Arge, Revsbæk, Zeh |
| SFD I/O-model | $O(\text{Sort}(X) \log \frac{X}{M} + \text{Sort } N)$ [2] | $O(\text{Sort}(N) + \text{Scan}(H \cdot X))$ [3] | [3] 2017, Arge, Rav, Raza, Revsbæk |
|  | * $O(\text{Sort}(N))$ [2] | * $O(\text{Sort}(N))$ [3] | |

*: assuming merge tree fits in memory

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

— H: height of the merge tree

— X: number of depressions

| | Flood-time query | Terrain-flood query | |
|---|---|---|---|
| SFD RAM-model | $O(N \log N)$ [1] | $O(N \log N)$ [1] | [1] 2004, Liu and Snoeyink |
| | | | [2] 2010, Arge, Revsbæk, Zeh |
| SFD I/O-model | $O(\text{Sort}(X) \log \frac{X}{M} + \text{Sort } N)$ [2] | $O(\text{Sort}(N) + \text{Scan}(H \cdot X))$ [3] | [3] 2017, Arge, Rav, Raza, Revsbæk |
| | | | [4] 2019, Lowe and Agarwal |
| | * $O(\text{Sort}(N))$ [2] | * $O(\text{Sort}(N))$ [3] | |
| MFD RAM-model | ** $O\left(N\left(|\mathcal{R}|k + H^{\omega} + H^2 \log H\right)\right)$ [4] | $O(N \log N)$ [4] | *: assuming merge tree fits in memory |
| | | | **: $O(NX + N \log N)$ pre-processing |
| MFD I/O-model | | | |

# — STATE OF THE ART

— H: height of the merge tree

— X: number of depressions

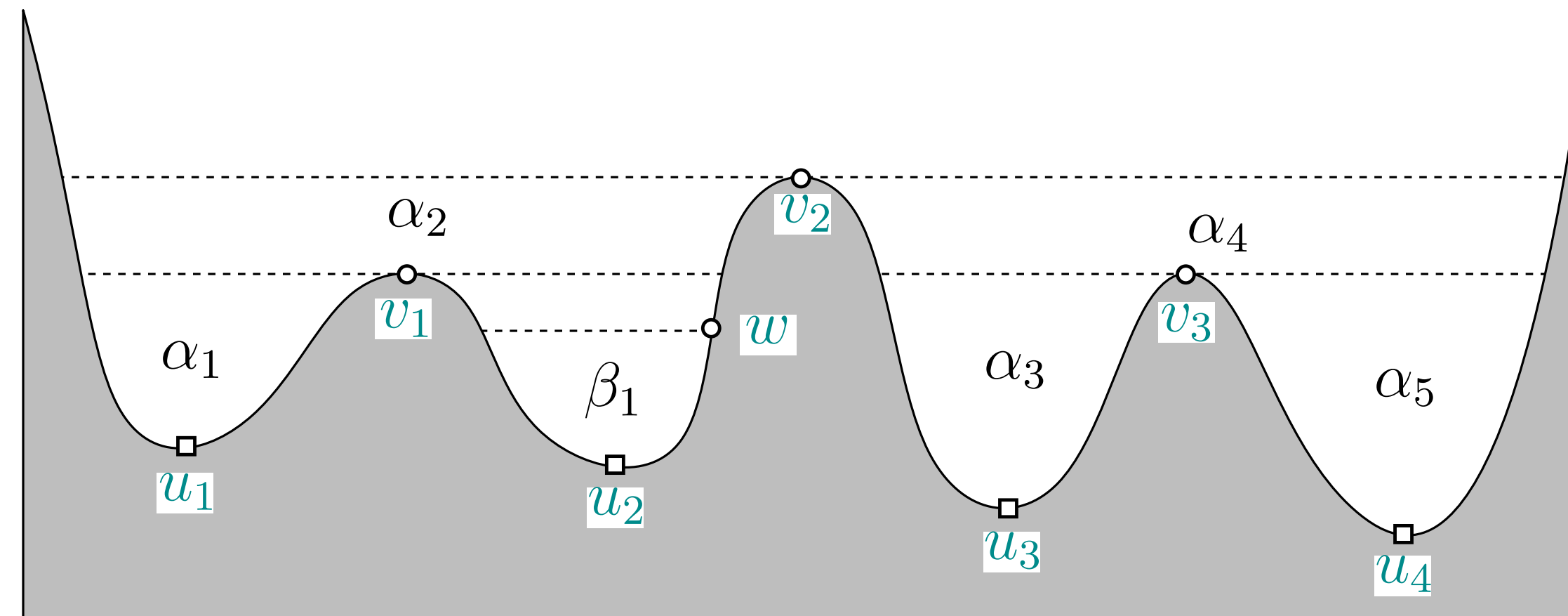|  | Flood-time query | Terrain-flood query |  |
|---|---|---|---|
| SFD RAM-model | $O(N \log N)$ [1] | $O(N \log N)$ [1] | [1] 2004, Liu and Snoeyink |
|  |  |  | [2] 2010, Arge, Revsbæk, Zeh |
| SFD I/O-model | $O(\text{Sort}(X) \log \frac{X}{M} + \text{Sort } N)$ [2] | $O(\text{Sort}(N) + \text{Scan}(H \cdot X))$ [3] | [3] 2017, Arge, Rav, Raza, Revsbæk |
|  |  |  | [4] 2019, Lowe and Agarwal |
|  | * $O(\text{Sort}(N))$ [2] | * $O(\text{Sort}(N))$ [3] | [5] 2021, Arge, Lowe, Svendsen, Agarwal, |
|  |  |  |  |
| MFD RAM-model | ** $O\left(N\left(\|\mathcal{R}\|k + H^\omega + H^2 \log H\right)\right)$ [4] | $O(N \log N)$ [4] | *: assuming merge tree fits in memory |
|  | *** $O(\|\phi\| \log \|\phi\|)$ [5] |  | **: $O(NX + N \log N)$ pre-processing |
| MFD I/O-model | * $O(\text{Sort}(N + \|\phi\|))$ [5] | * $O(\text{Sort}(N))$ [5] | ***: $O(N \log N)$ preprocessing |

## — FLOW FUNCTIONS

— Rain distribution: $\mathcal{R}(v,t) : \mathbb{V} \times \mathbb{R} \to \mathbb{R}_{\geq 0}$
piece-wise constant changing at times $\{t_0, t_1, \dots, t_K\}$

— Flow function $\phi_v$: the flow rate over a vertex $v$

— $\phi_v$ is a piece-wise constant function

— $\phi_v$ changes only at spill events and when the rain distribution changes

# — SADDLES AND NON-SADDLES

— Saddle Vertex: $v_i$

— Sink Vertex: $u_i$

— Maximal Depression: $\alpha_i$

— Non-maximal Depression $\beta_1$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

## — ALGORITHM FOR COMPUTING FLOW FUNCTIONS

— Preprocessing:

  — For all $v \in \Sigma$, compute the maximal depression containing $v$ (Sort$(N)$ [1])

  — For all $v \in \Sigma$, compute the volume of the depression $\alpha_v$ (Sort$(N)$ [2])

  — For each maximal depression $\beta$, compute the amount
  of rain falling directly in $\beta$ ($O(\text{Sort}(N) + \text{Sort}(|\mathcal{R}|))$)

[1] 2009, Arge and Revsbæk

[2] 2010, Arge, Revsbæk, Zeh

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — ALGORITHM FOR COMPUTING FLOW FUNCTIONS

— Sweep:

   — At height $l$ maintain depressions $\alpha_i$ in the sublevel set $h_{<l}$



$h_{<l}: \{\alpha_7\}$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# ── **ALGORITHM FOR COMPUTING FLOW FUNCTIONS**

— Sweep:

    — At height $l$ maintain depressions $\alpha_i$ in the sublevel set $h_{<l}$



$h_{<l}: \{\alpha_5, \alpha_6\}$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

— Sweep:

    — At height $l$ maintain depressions $\alpha_i$ in the sublevel set $h_{<l}$



$h_{<l}: \{\alpha_1, \alpha_2, \alpha_3, \alpha_4\}$

— Sweep:

  — At height $l$ maintain depressions $\alpha_i$ in the sublevel set $h_{<l}$

  — For each $\alpha_i$: maintain

    — $E(\alpha_i)$: the edges crossing the sweep line into $\alpha_i$

    — For each $e \in E(\alpha_i)$: maintain $\phi_e(t)$

    — $F_{\alpha_i(t)}$: fill-rate function of $\alpha_i$



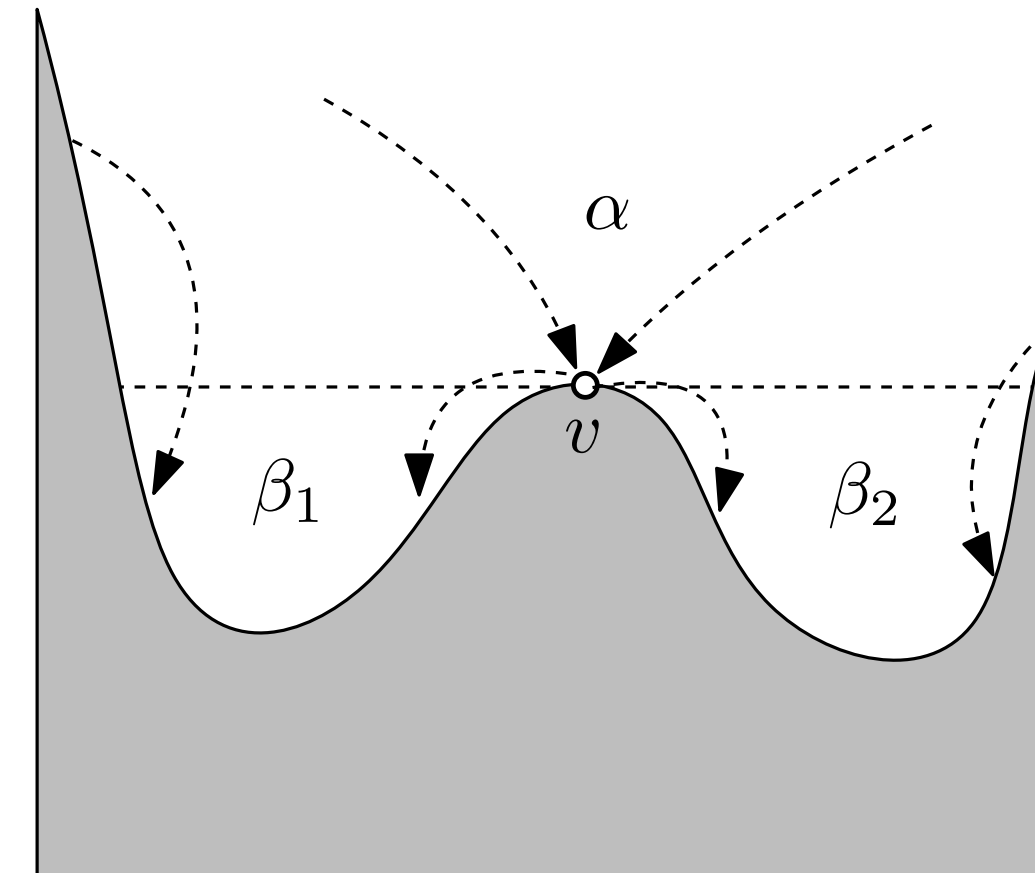$h_{<l}: \{\alpha_1, \alpha_2, \alpha_3, \alpha_4\}$

## — **NON-SADDLE VERTICES**

— For each $\alpha_i$: maintain

     — $E(\alpha_i)$: the edges crossing the sweep line into $\alpha_i$

     — For each $e \in E(\alpha_i)$: maintain $\phi_e(t)$

     — $F_{\alpha_i(t)}$: fill-rate function of $\alpha_i$



— Whenever we cross a non-saddle:

     — Compute $\phi_v = \mathcal{R}(v, t) + \sum_{e \in E(\alpha)} \phi_e(t)$

     — For each outgoing edge $e$: Compute $\phi_e(t) = w_e \cdot \phi_v(t)$
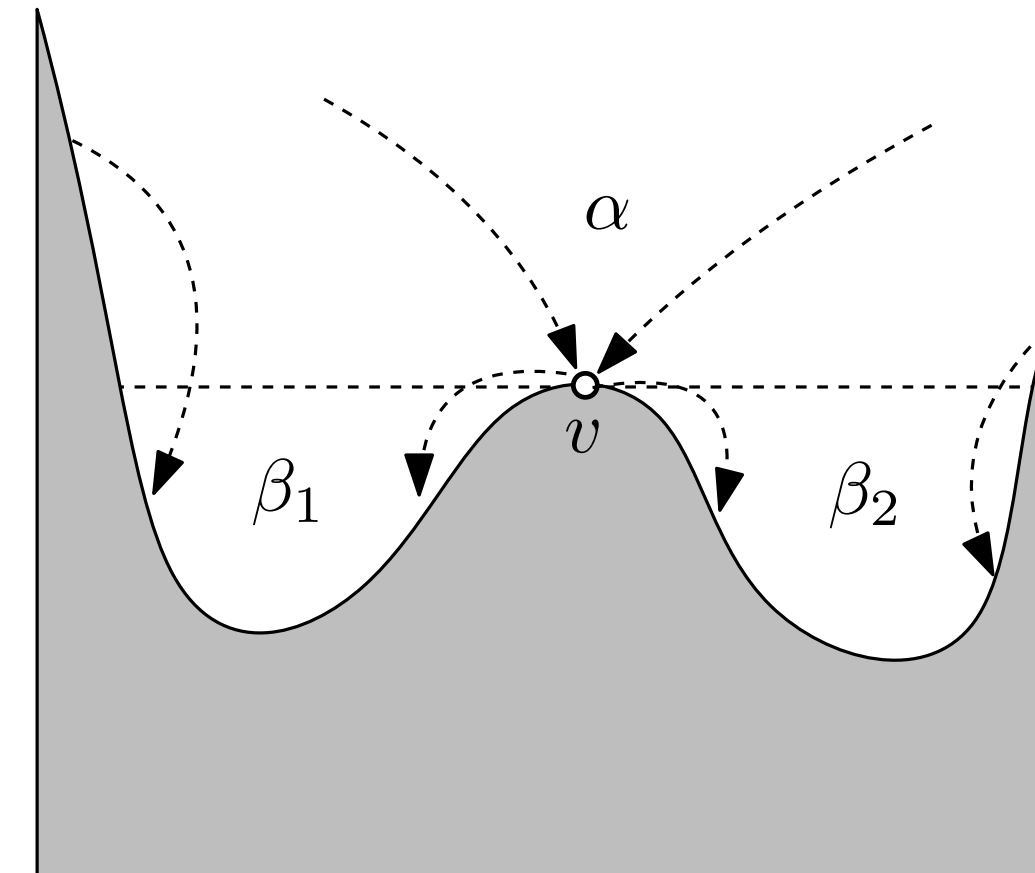
     — Update $E(\alpha)$: remove incoming edge, add outgoing edges

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — SADDLE VERTICES

— For each $\alpha_i$: maintain

  — $E(\alpha_i)$: the edges crossing the sweep line into $\alpha_i$

  — For each $e \in E(\alpha_i)$: maintain $\phi_e(t)$

  — $F_{\alpha_i(t)}$: fill-rate function of $\alpha_i$

— Whenever we cross a saddle:

  — Compute $\phi_e(t)$ for outgoing edges as before

  — Partition $E(\alpha)$ into $E(\beta_1)$ and $E(\beta_2)$
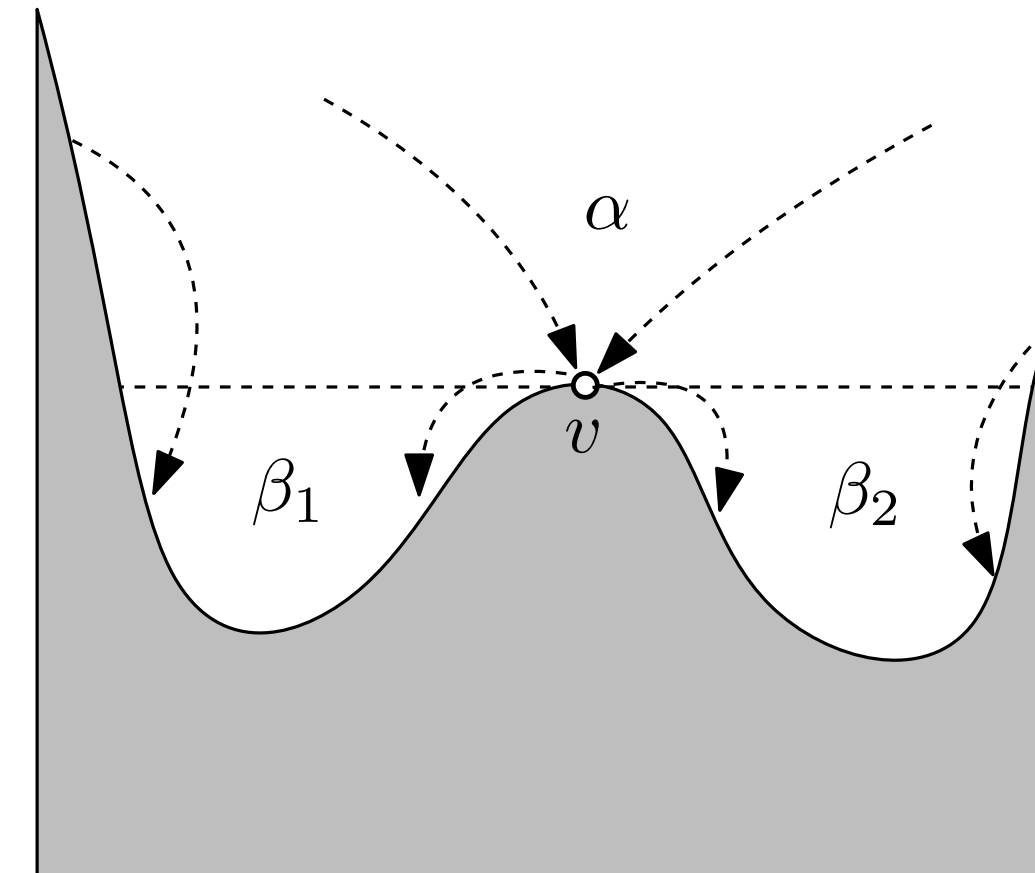
AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

## — SADDLE VERTICES

— For each $\alpha_i$: maintain

    — $E(\alpha_i)$: the edges crossing the sweep line into $\alpha_i$

    — For each $e \in E(\alpha_i)$: maintain $\phi_e(t)$

    — $F_{\alpha_i(t)}$: fill-rate function of $\alpha_i$

— Whenever we cross a saddle:

    — Compute $\phi_e(t)$ for outgoing edges as before

    — Partition $E(\alpha)$ into $E(\beta_1)$ and $E(\beta_2)$

    — Compute fill-rate functions for $\beta_1$ and $\beta_2$

    — $F_{\beta_1}(t) = R(\beta_1, t) + \sum_{e \in E(\beta_1)} \phi_e(t)$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

SOLIDUM PETIT IN PROFUNDIS · UNIVERSITAS ARHUSIENSIS ·

## — SADDLE VERTICES

— For each $\alpha_i$: maintain

    — $E(\alpha_i)$: the edges crossing the sweep line into $\alpha_i$

    — For each $e \in E(\alpha_i)$: maintain $\phi_e(t)$

    — $F_{\alpha_i(t)}$: fill-rate function of $\alpha_i$

— Whenever we cross a saddle:

    — Compute $\phi_e(t)$ for outgoing edges as before

    — Partition $E(\alpha)$ into $E(\beta_1)$ and $E(\beta_2)$

    — Compute fill-rate functions for $\beta_1$ and $\beta_2$

    — $F_{\beta_1}(t) = R(\beta_1, t) + \sum_{e \in E(\beta_1)} \phi_e(t)$

    — Assume $\beta_1$ spills first: Add the spill from $\beta_1$ to $\phi_v$

    — Update $\phi_e$ for outgoing edges $e$ and update $E(\beta_1)$, and $E(\beta_2)$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — COMBINING EVERYTHING

— Total: $O(\text{Sort}(N + |\phi|))$

# — COMBINING EVERYTHING

— Total: $O(\text{Sort}(N + |\phi|))$

— For each $v \in \Sigma$, we precomputed the volume of $\beta_v$.

— For each maximal depression $\alpha$, we computed the fill function $F_\alpha(t)$
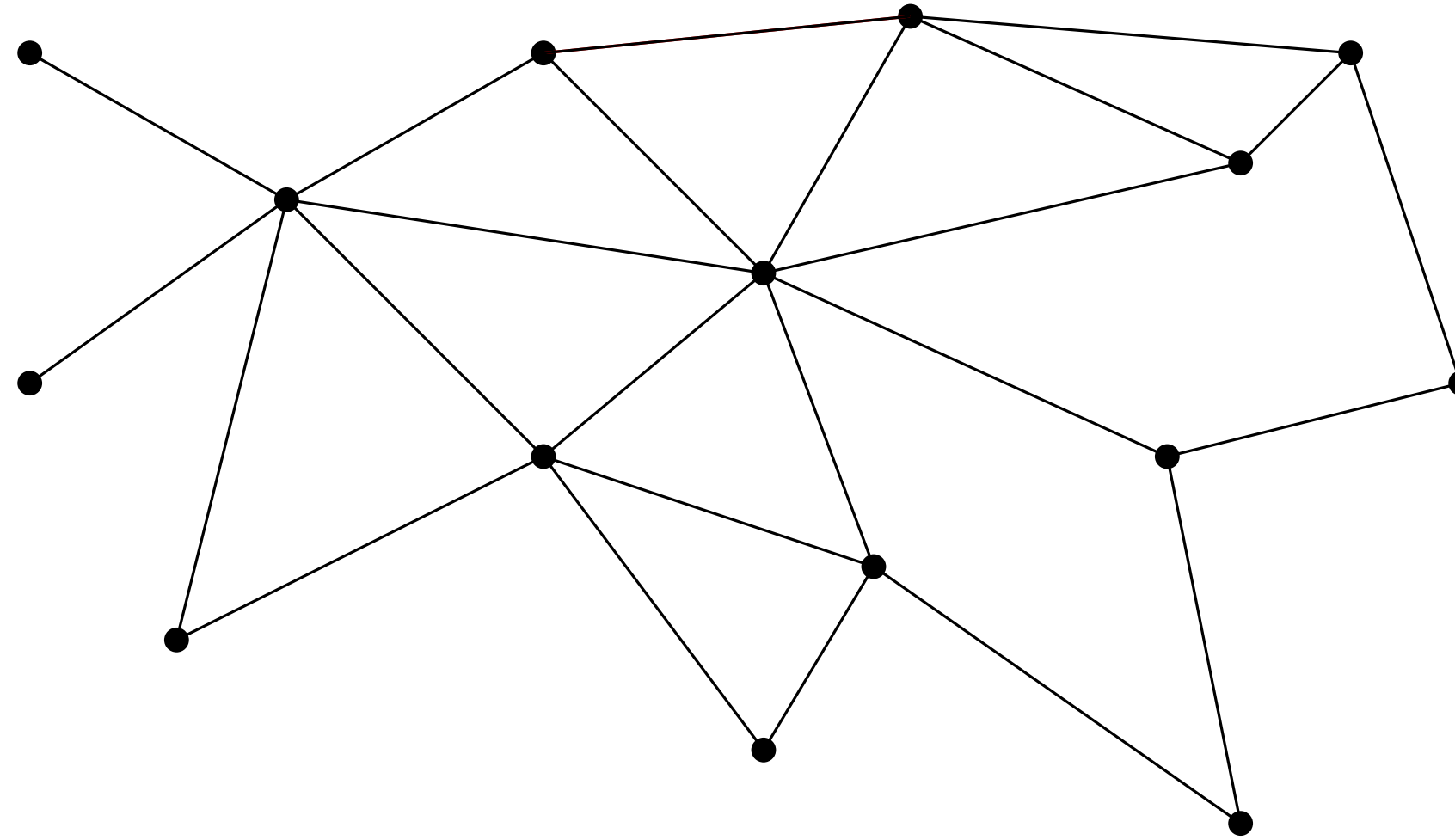
— We use this to compute the fill-time of $v$ !

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

SOLIDUM PETIT IN PROFUNDIS·
·UNIVERSITAS ARHUSIENSIS·

# — COMBINING EVERYTHING

— Total: $O(\mathrm{Sort}(N + |\phi|))$

— For each $v \in \Sigma$, we precomputed the volume of $\beta_v$.

— For each maximal depression $\alpha$, we computed the fill function $F_\alpha(t)$

— We use this to compute the fill-time of $v$ !

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — OPEN PROBLEMS

— Can we $O(N \log N)$ instead of $O(|\phi| \log |\phi|)$ in the RAM model?

— Can we get $\text{Sort}(\phi)$ in the I/O model with no assumptions on $M$?

— Output sensitive algorithm for the I/O-model?

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# Practical I/O-Efficient Multiway Separators

Svend C. Svendsen

Manuscript

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

SOLIDUM PETIT IN PROFUNDIS · UNIVERSITAS ARHUSIENSIS ·

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# PLANAR SEPARATOR THEOREM

— Lipton and Tarjan 1979:

$$\frac{1}{3}N \leq |A|, |B| \leq \frac{2}{3}N$$

$$|S| = O(\sqrt{N})$$
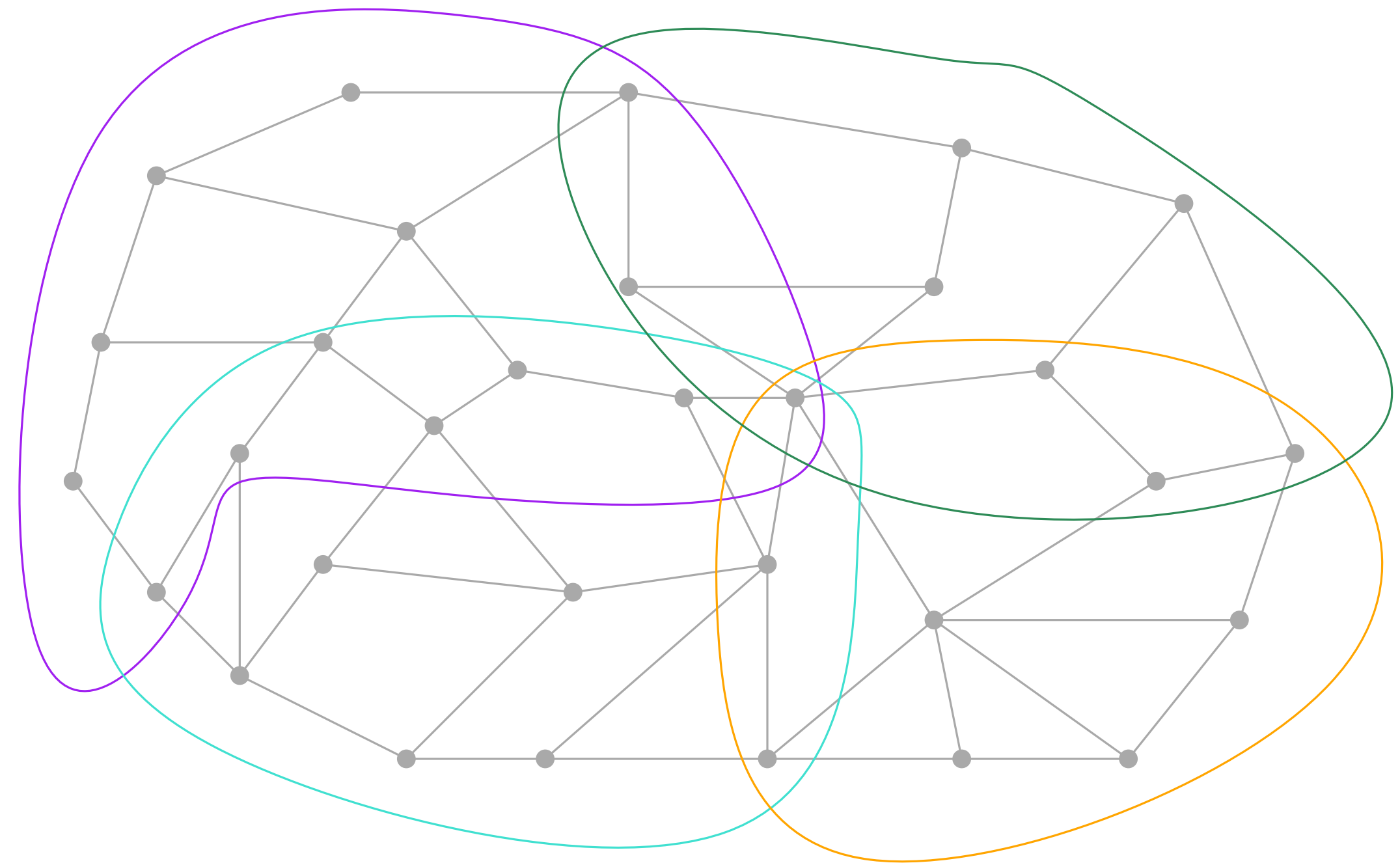
— $O(N)$ time

# — MULTIWAY SEPARATOR

— Frederickson (1953): $r$-way separator

    — Divide a graph into $r$ regions

    — Each region has $O(N/r)$ vertices

    — $O(\sqrt{Nr})$ boundary vertices

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — MULTIWAY SEPARATOR

— Frederickson (1953): $r$-way separator

   — Divide a graph into $r$ regions

   — Each region has $O(N/r)$ vertices

   — $O(\sqrt{Nr})$ boundary vertices

— Useful in the I/O-model: $N/M$-separator

— Can solve problems such as SSSP, DFS, finding strongly connected components, and topological sorting [1,2]

[1] 2003, Arge, Toma, and Zeh
[2] 2005, Agarwal, Arge, and Yi

## — STATE OF THE ART

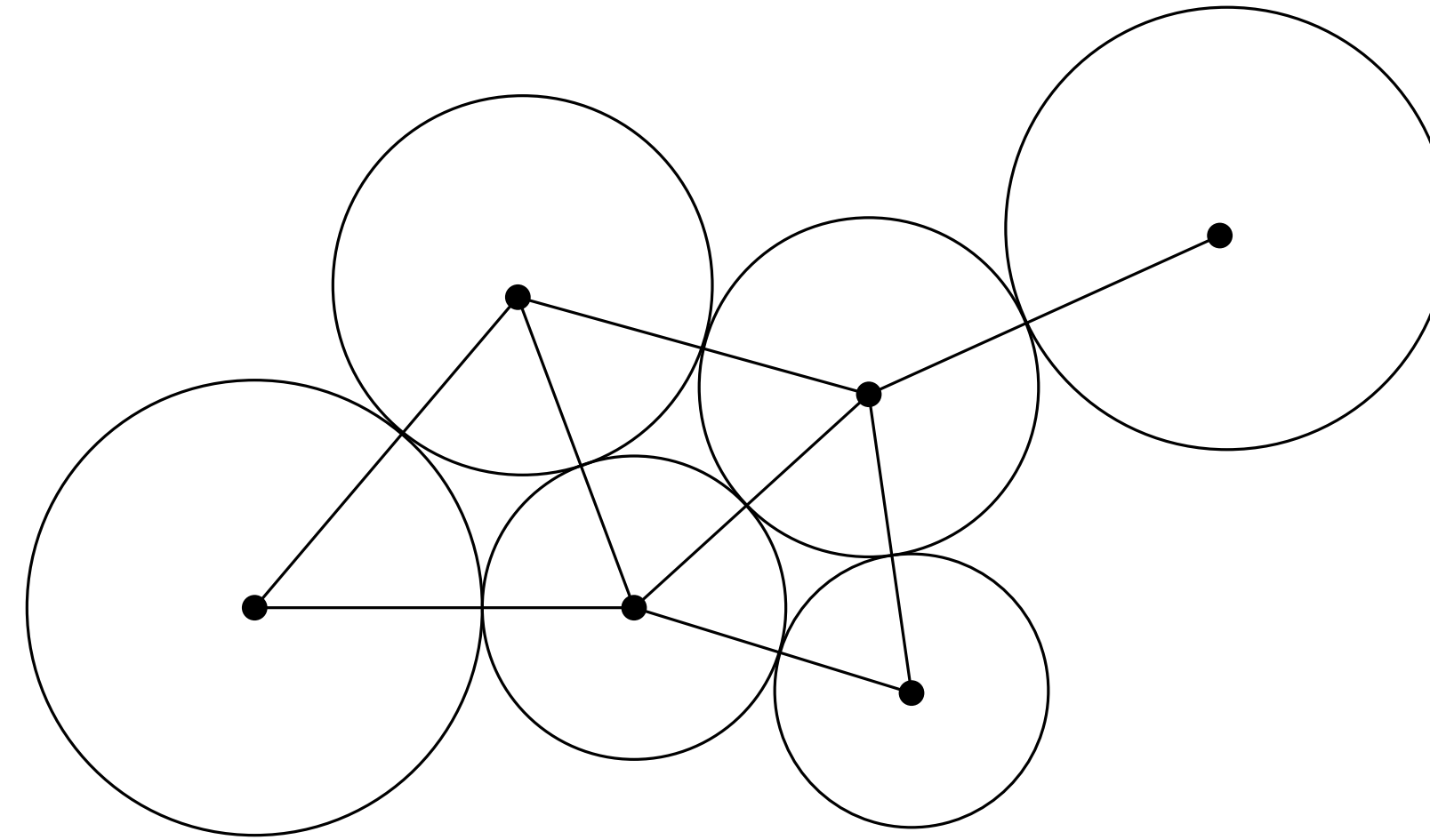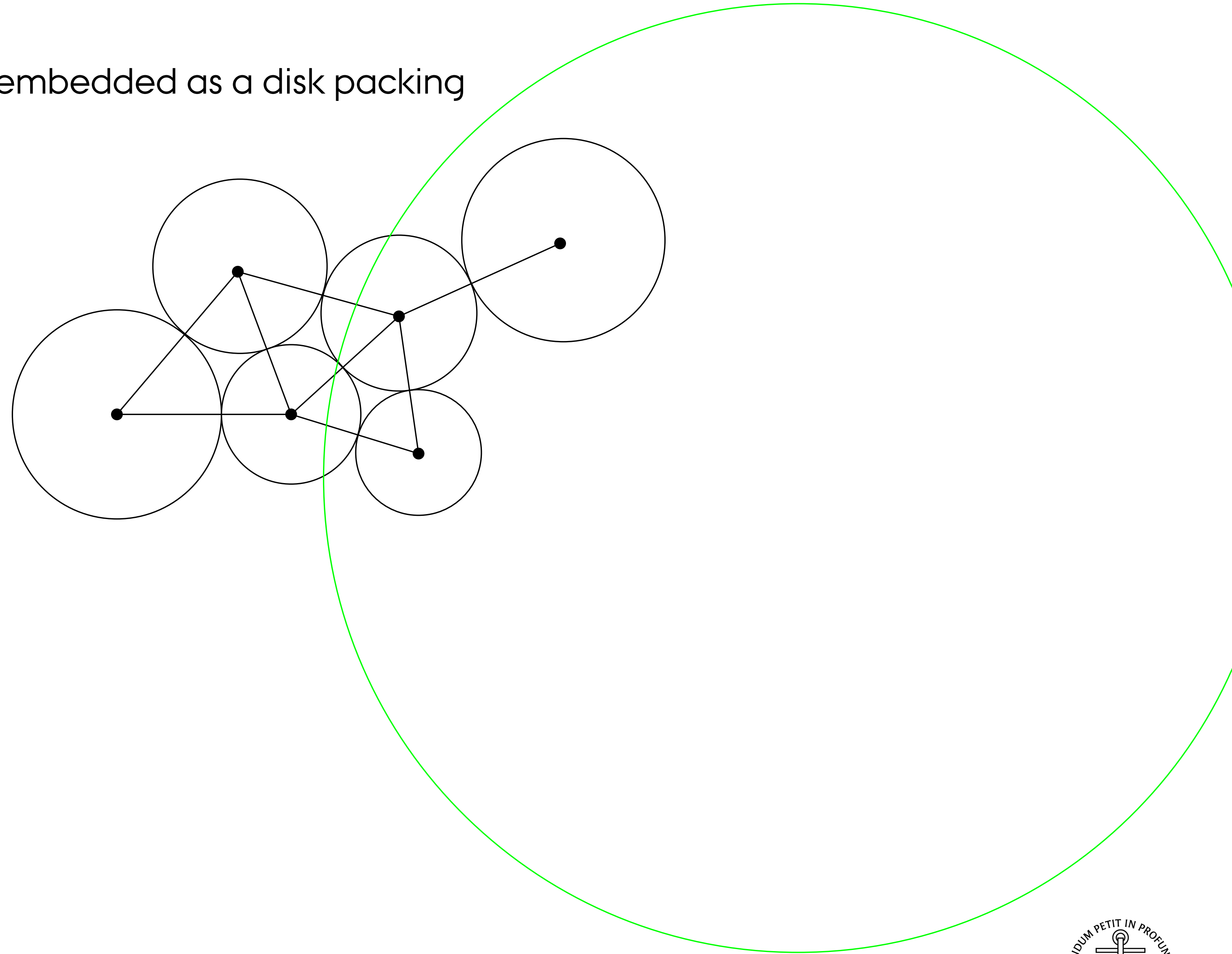|  | I/Os | Internal Computation |
|---|---|---|
| Maheshwari and Zeh (2008) | $O(\mathrm{Sort}(N))$ | |
| Arge, Walderveen, and Zeh (2013) | $O(\mathrm{Sort}(N))$ | $O(N \log N)$ |

# — DISK PACKINGS

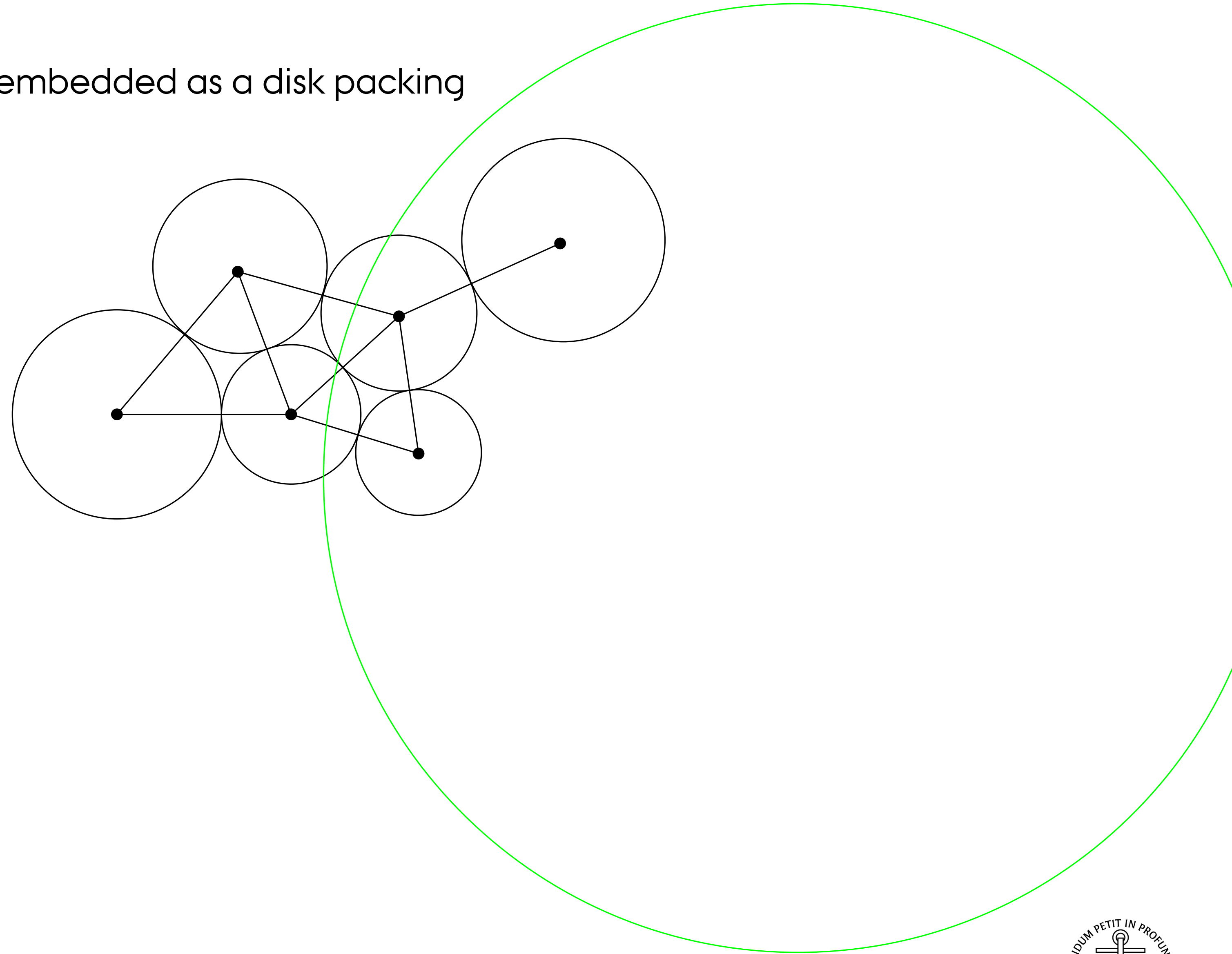— Koebe (1936): every planar graph can be embedded as a disk packing

# — DISK PACKINGS

— Koebe (1936): every planar graph can be embedded as a disk packing

— Miller, Teng, Thurston, Vavasis (1997):

— At most $\frac{3}{4}N$ disks inside

— At most $\frac{3}{4}N$ disks outside

— At most $O(\sqrt{N})$ disks crossing

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — DISK PACKINGS

— Koebe (1936): every planar graph can be embedded as a disk packing

— Miller, Teng, Thurston, Vavasis (1997):

  — At most $\frac{3}{4}N$ disks inside

  — At most $\frac{3}{4}N$ disks outside

  — At most $O(\sqrt{N})$ disks crossing

  — Given a disk packing: $O(\text{Scan}(N))$ I/Os

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

## — COMPUTING MULTIWAY SEPARATORS ON DISK PACKINGS

— How do we compute an $r$-way separator for $r = \frac{N}{M}$?

— Naively computing an $\frac{N}{M}$-way separator: $O(\text{Scan}(N) \log \frac{N}{M})$

AARHUS
UNIVERSITY
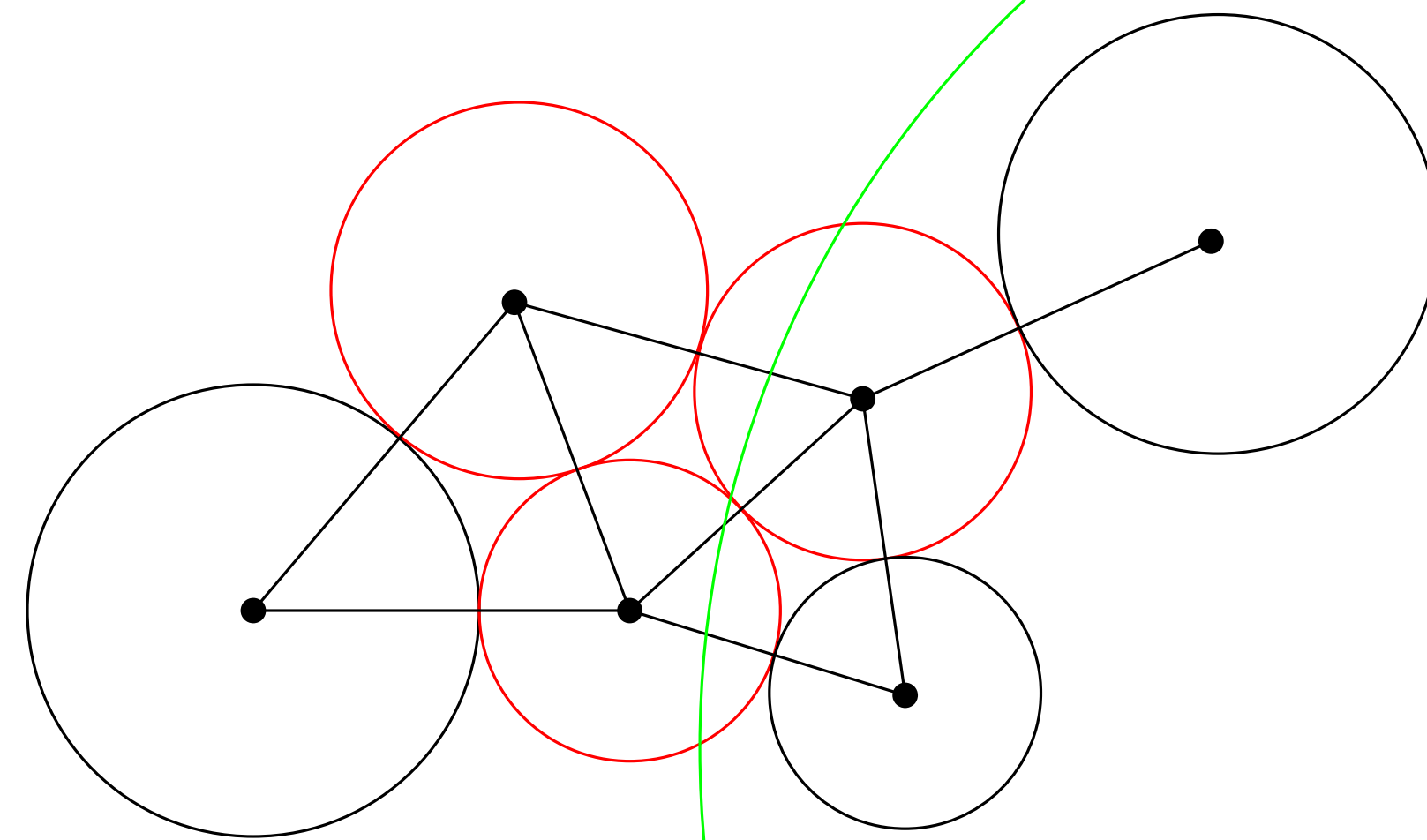DEPARTMENT OF COMPUTER SCIENCE

# — COMPUTING MULTIWAY SEPARATORS ON DISK PACKINGS

— How do we compute an $r$-way separator for $r = \frac{N}{M}$?

— Naively computing an $\frac{N}{M}$-way separator: $O(\text{Scan}(N) \log \frac{N}{M})$

— We want $O(\text{Sort}(N)) = O\left(\text{Scan}(N) \log_{M/B} \frac{N}{B}\right)$
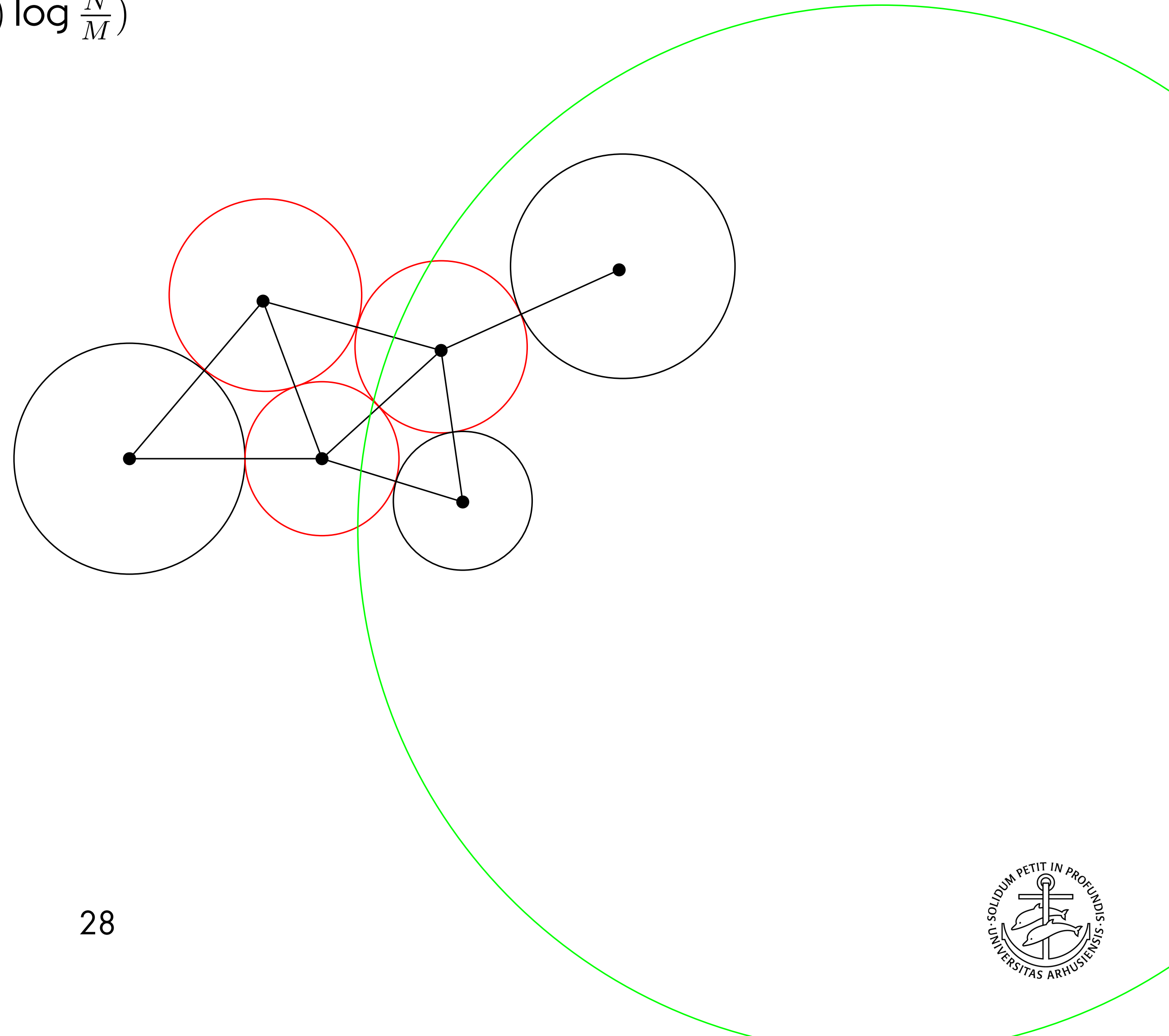
AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — COMPUTING MULTIWAY SEPARATORS ON DISK PACKINGS

— How do we compute an $r$-way separator for $r = \frac{N}{M}$?

— Naively computing an $\frac{N}{M}$-way separator: $O(\text{Scan}(N) \log \frac{N}{M})$

— We want $O(\text{Sort}(N)) = O\big(\text{Scan}(N) \log_{M/B} \frac{N}{B}\big)$

— Solution:
  — Given a disk packing $P$, sample $S \subseteq P$

  — Compute multiway separator on $S$

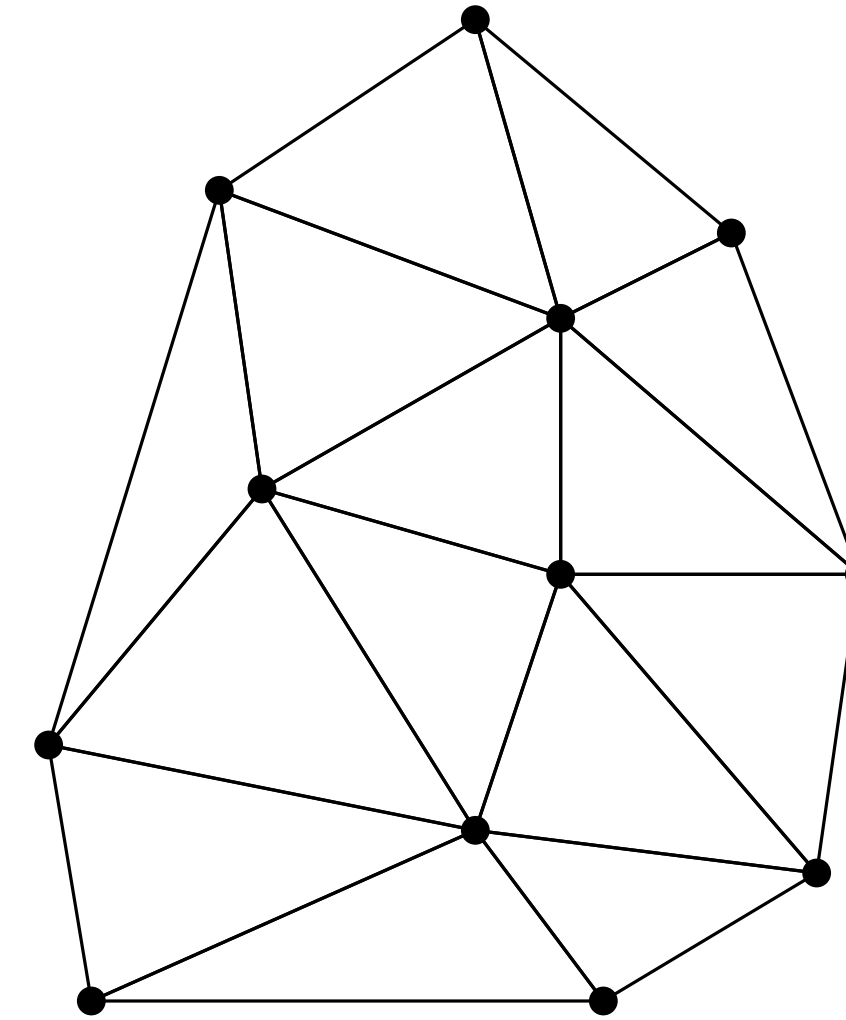  — Split $P$ using the multiway separator (hopefully)

# — COMPUTING MULTIWAY SEPARATORS ON DISK PACKINGS

— How do we compute an $r$-way separator for $r = \frac{N}{M}$?

— Naively computing an $\frac{N}{M}$-way separator: $O(\text{Scan}(N) \log \frac{N}{M})$

— We want $O(\text{Sort}(N)) = O\big(\text{Scan}(N) \log_{M/B} \frac{N}{B}\big)$

— Solution:
  — Given a disk packing $P$, sample $S \subseteq P$

  — Compute multiway separator on $S$

  — Split $P$ using the multiway separator (hopefully)

— $O(\text{Sort}(N))$ but no bound on boundary vertices

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — COMPUTING MULTIWAY SEPARATORS ON DISK PACKINGS

— How do we compute an $r$-way separator for $r = \frac{N}{M}$?

— Naively computing an $\frac{N}{M}$-way separator: $O(\text{Scan}(N) \log \frac{N}{M})$

— We want $O(\text{Sort}(N)) = O\left( \text{Scan}(N) \log_{M/B} \frac{N}{B} \right)$

— Solution:
  — Given a disk packing $P$, sample $S \subseteq P$

  — Compute multiway separator on $S$

  — Split $P$ using the multiway separator (hopefully)

— $O(\text{Sort}(N))$ but no bound on boundary vertices

— Upper bound on boundary vertices if

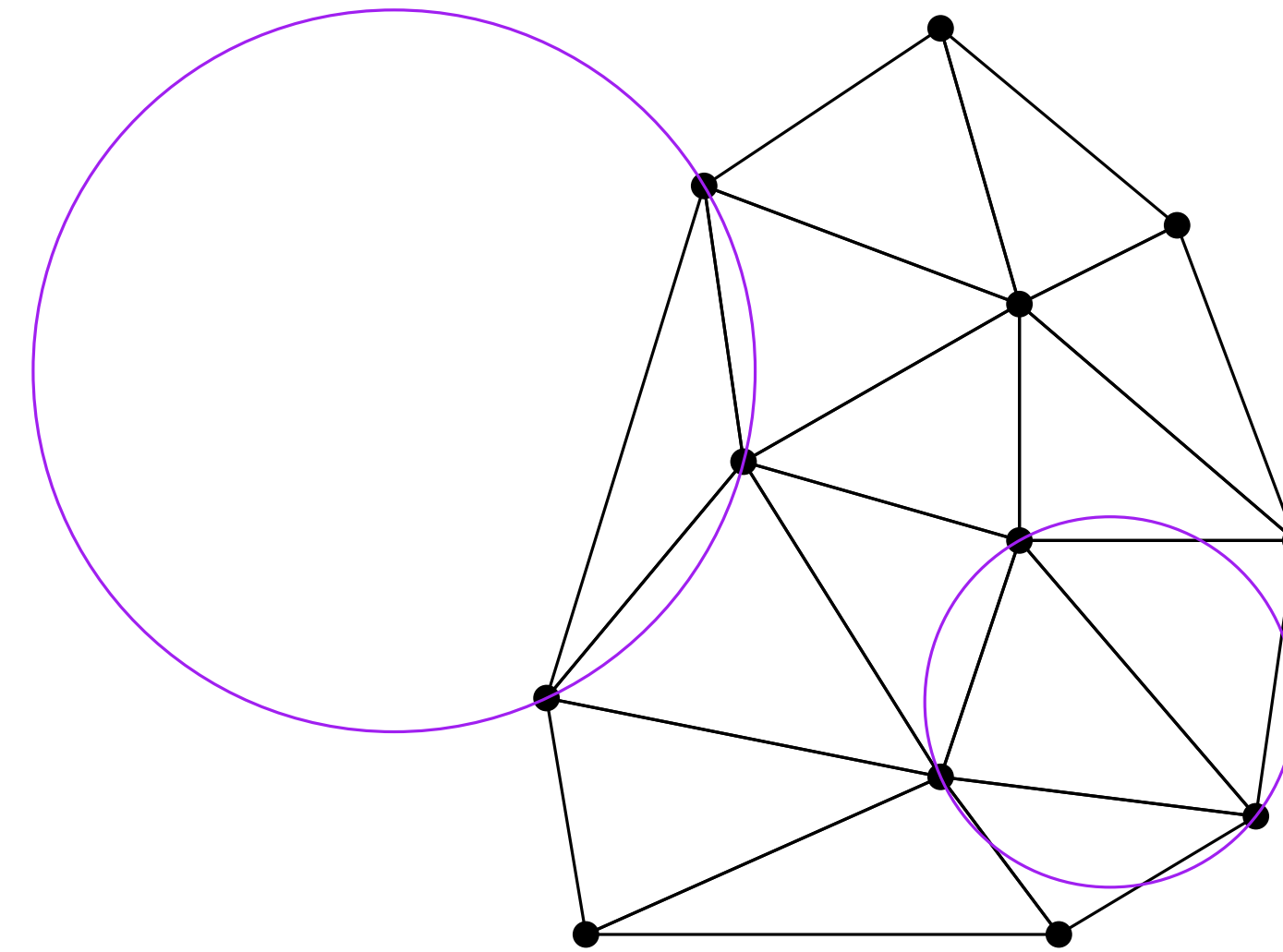$\log^3 \frac{M}{B} \log \log \frac{M}{B} \log N = O(\sqrt{M})$

— Disk Packings are difficult to compute

— Use circumcircles

# APPLYING TO TRIANGULATIONS

— Disk Packings are difficult to compute

— Use circumcircles
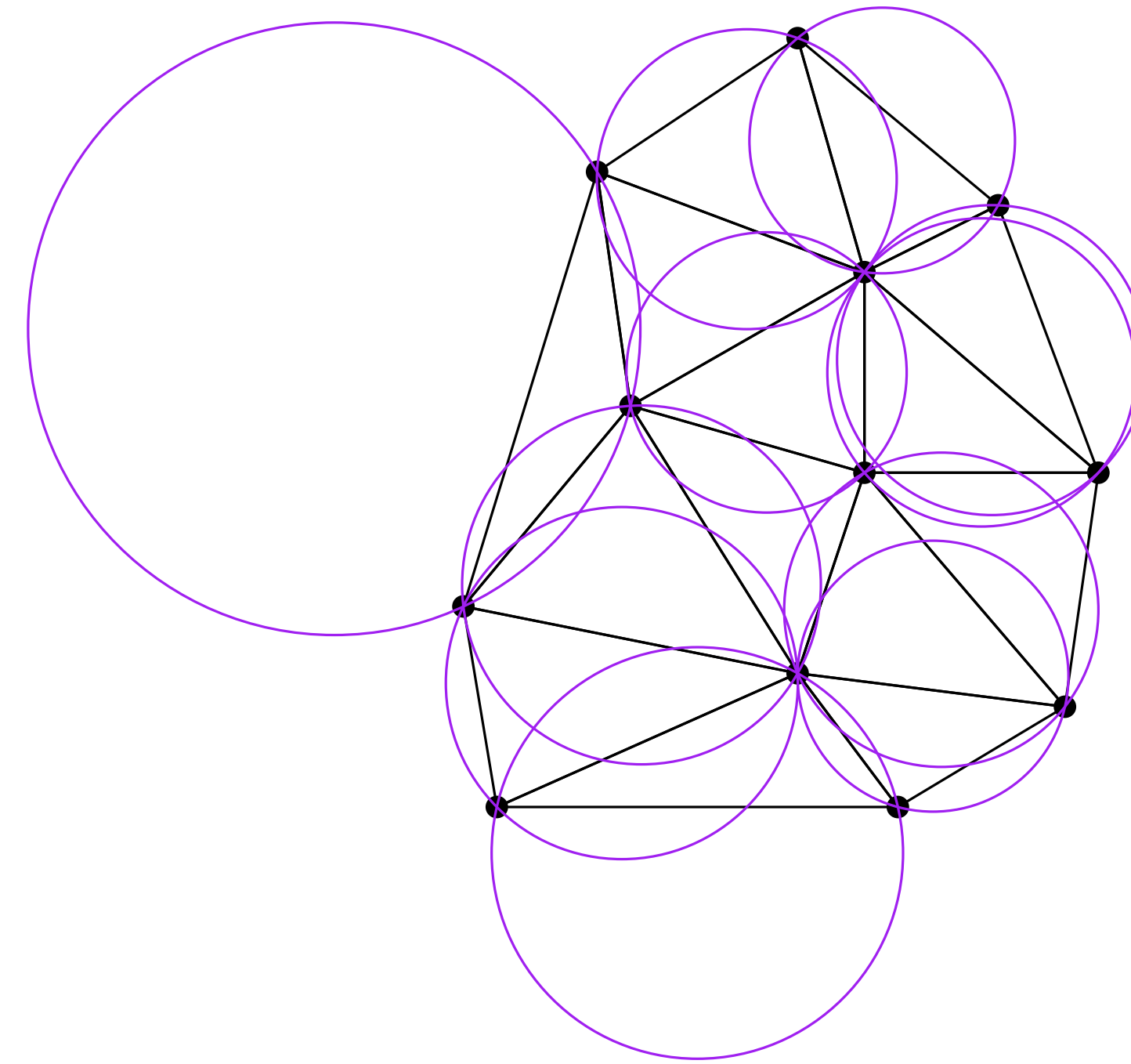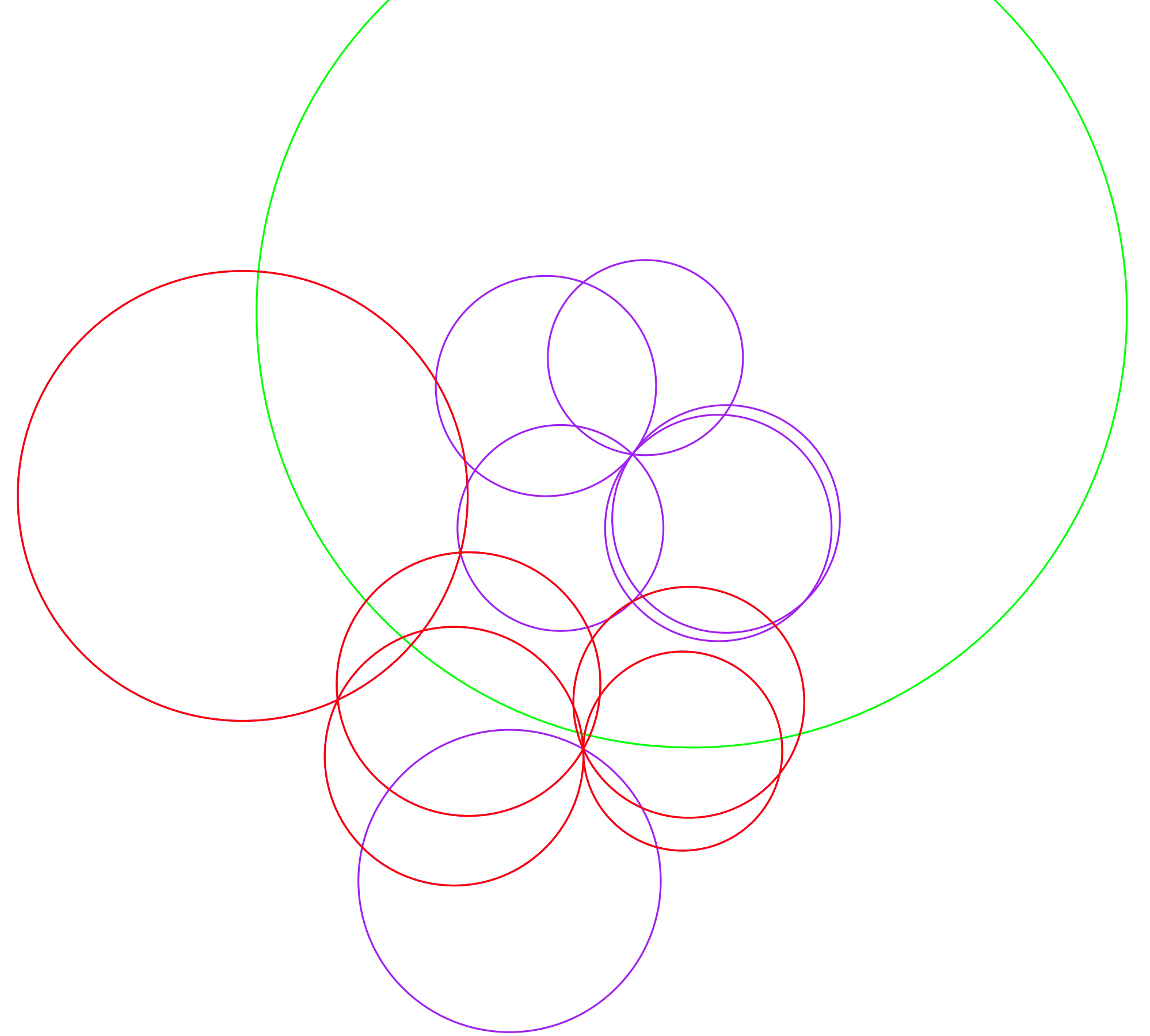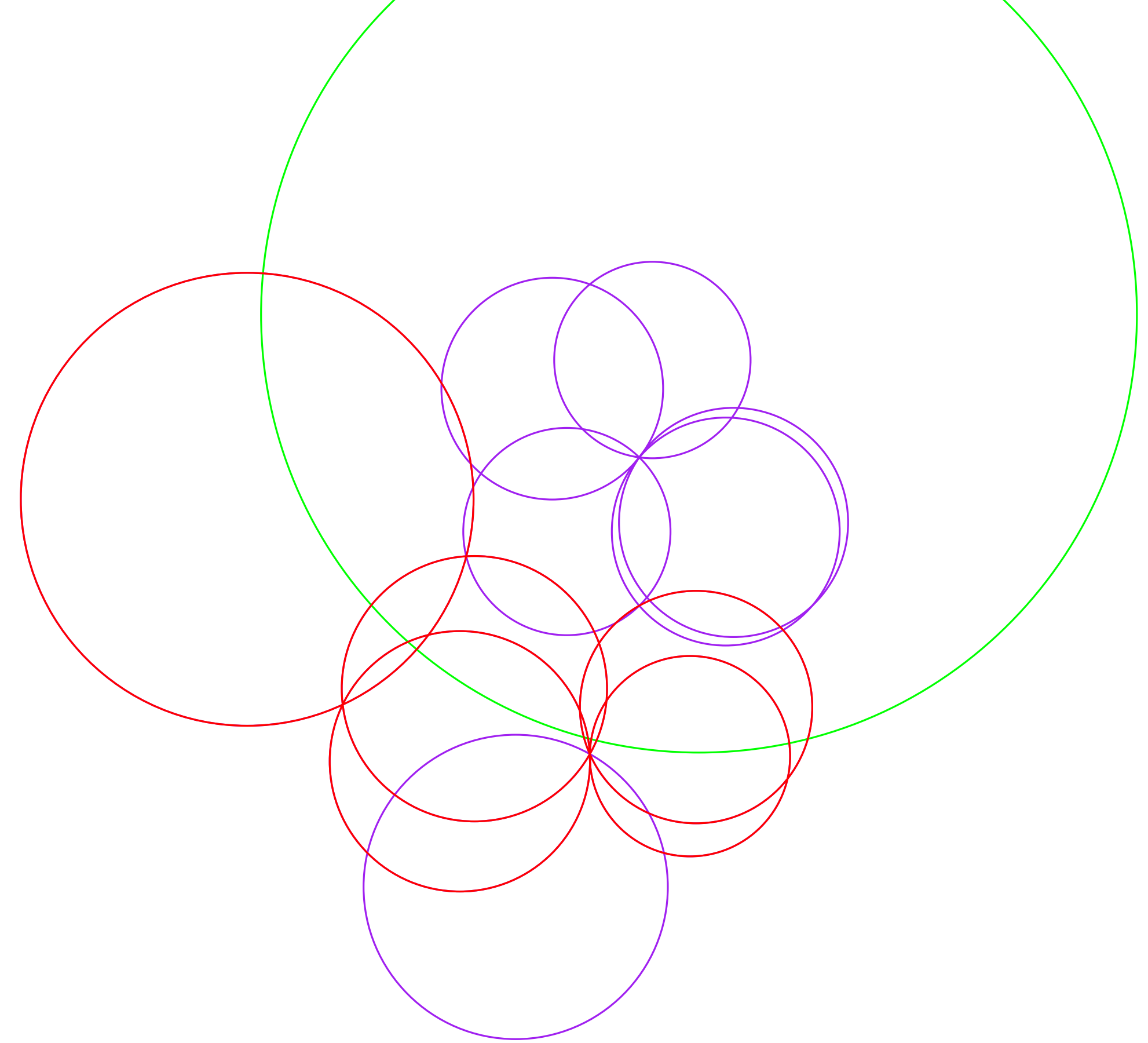
AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# APPLYING TO TRIANGULATIONS

— Disk Packings are difficult to compute

— Use circumcircles

# — APPLYING TO TRIANGULATIONS

— Disk Packings are difficult to compute

— Use circumcircles

— Miller, Teng, Thurston, Vavasis (1997):
If at most $k$ disks overlap in one point,
the separator has size $O(\sqrt{kN})$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — APPLYING TO TRIANGULATIONS

— Disk Packings are difficult to compute

— Use circumcircles

— Miller, Teng, Thurston, Vavasis (1997):
If at most $k$ disks overlap in one point,
the separator has size $O(\sqrt{kN})$

— This works well in practice on terrain!

— Triangulation are fast to compute (Sort$(N)$ [1][2])

[1] 1993, Goodrich, Tsay, Vengroff, and Vitter
[2] 2005, Agarwal, Arge, and Yi

AARHUS
UNIVERSITY
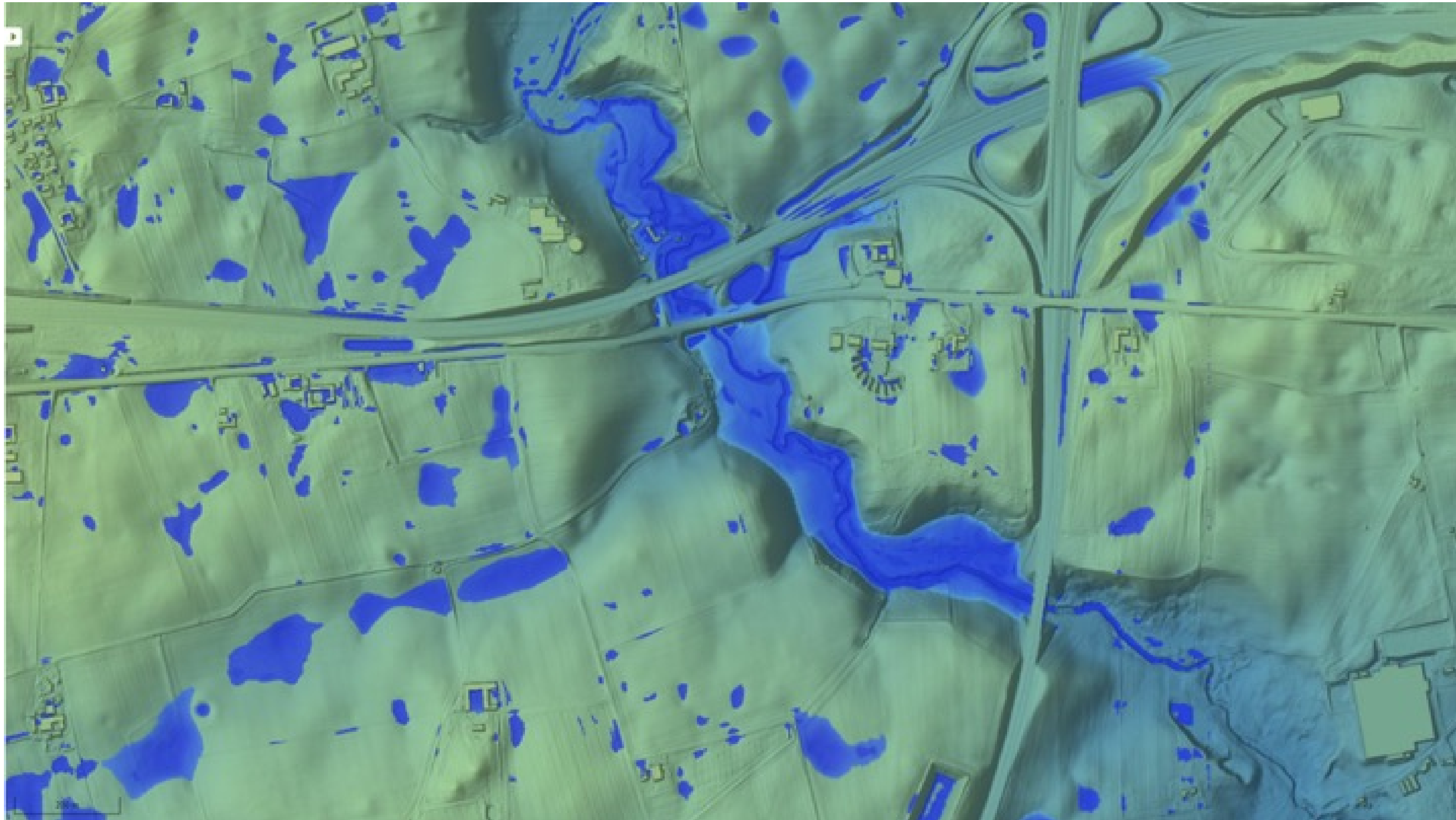DEPARTMENT OF COMPUTER SCIENCE

## — OPEN PROBLEMS

— Can we get a bound on the boundary size?

— Can we do better on circumcircles?

AARHUS
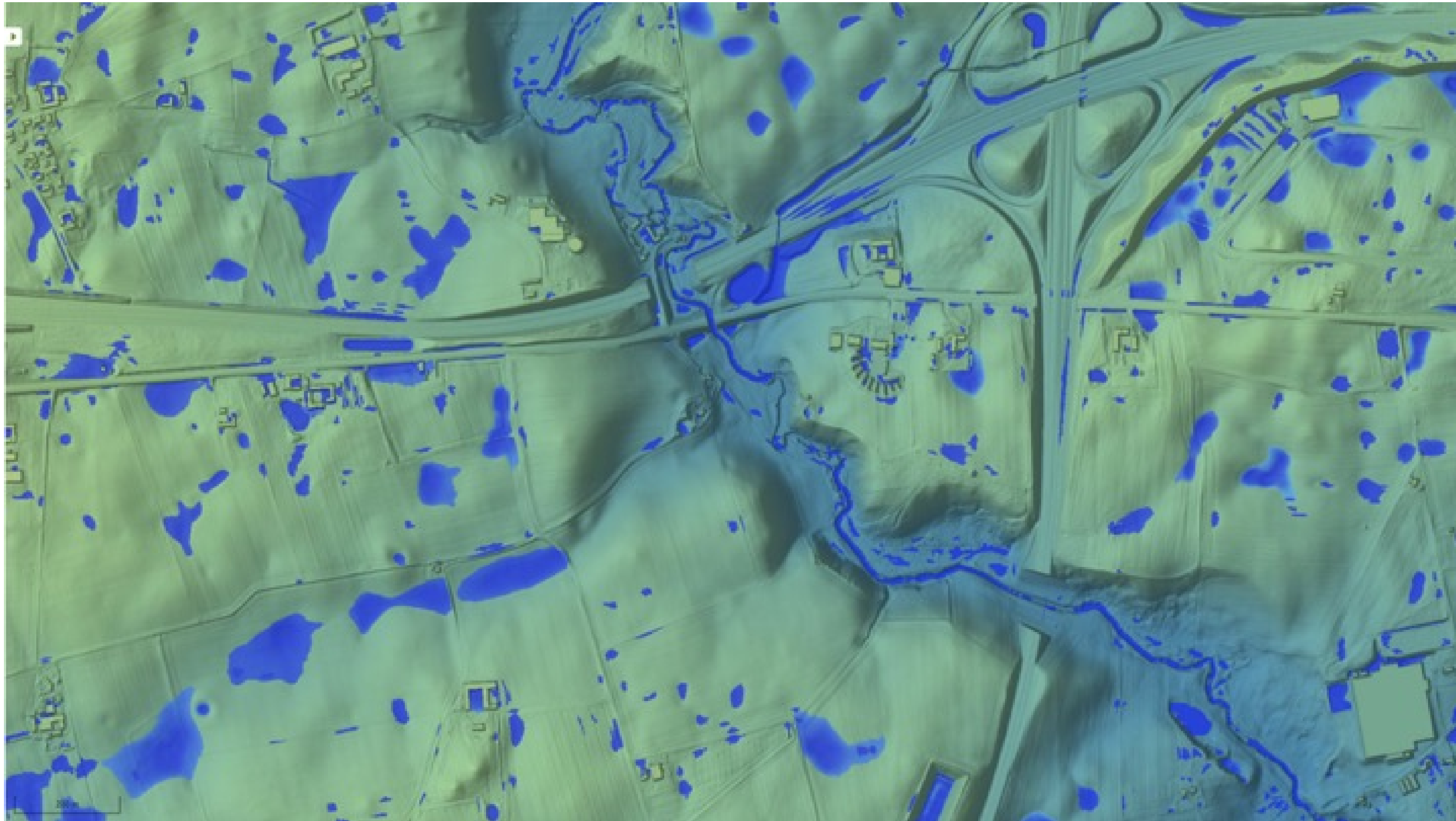UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# Learning to Find Hydrological Corrections

Lars Arge, Allan Grønlund, Svend C. Svendsen, Jonas Tranberg

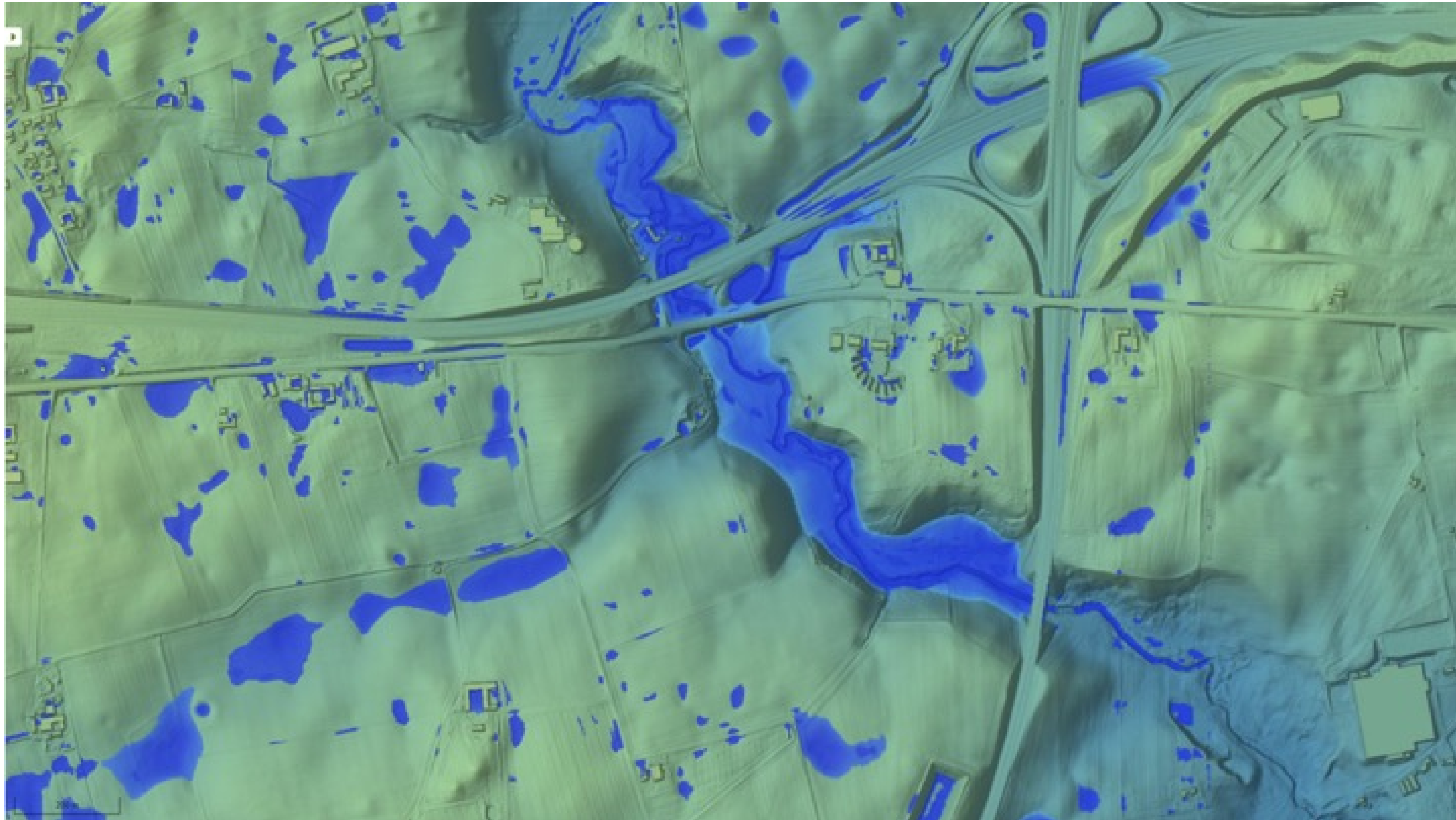ACM SIGSPATIAL 2019

# WHAT ARE HYDROLOGICAL CORRECTIONS?

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# WHAT ARE HYDROLOGICAL CORRECTIONS?
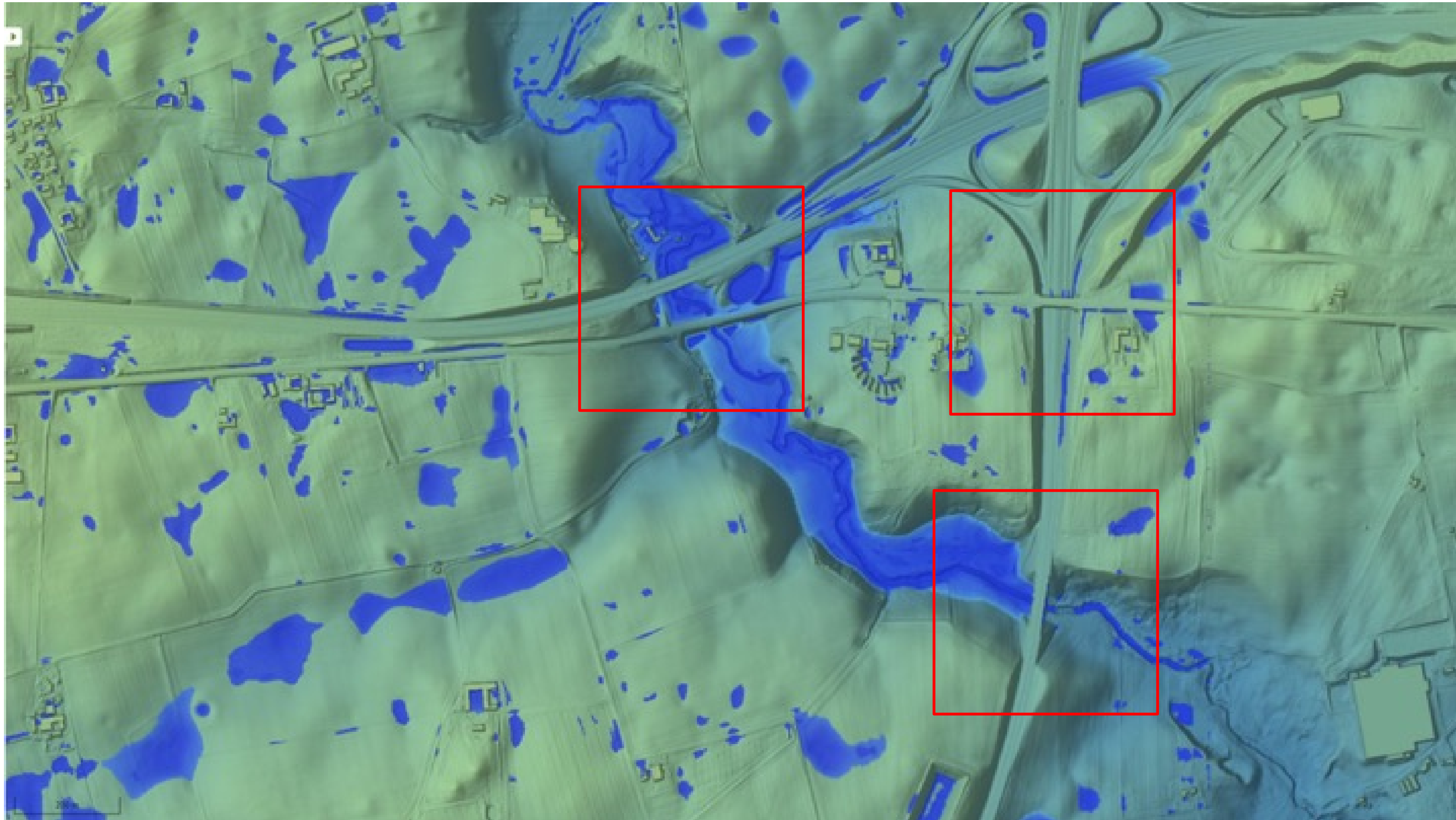
AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — **INPUT DATA**

— Digital Elevation Model

  — 415 billion cells

— Road and River Networks
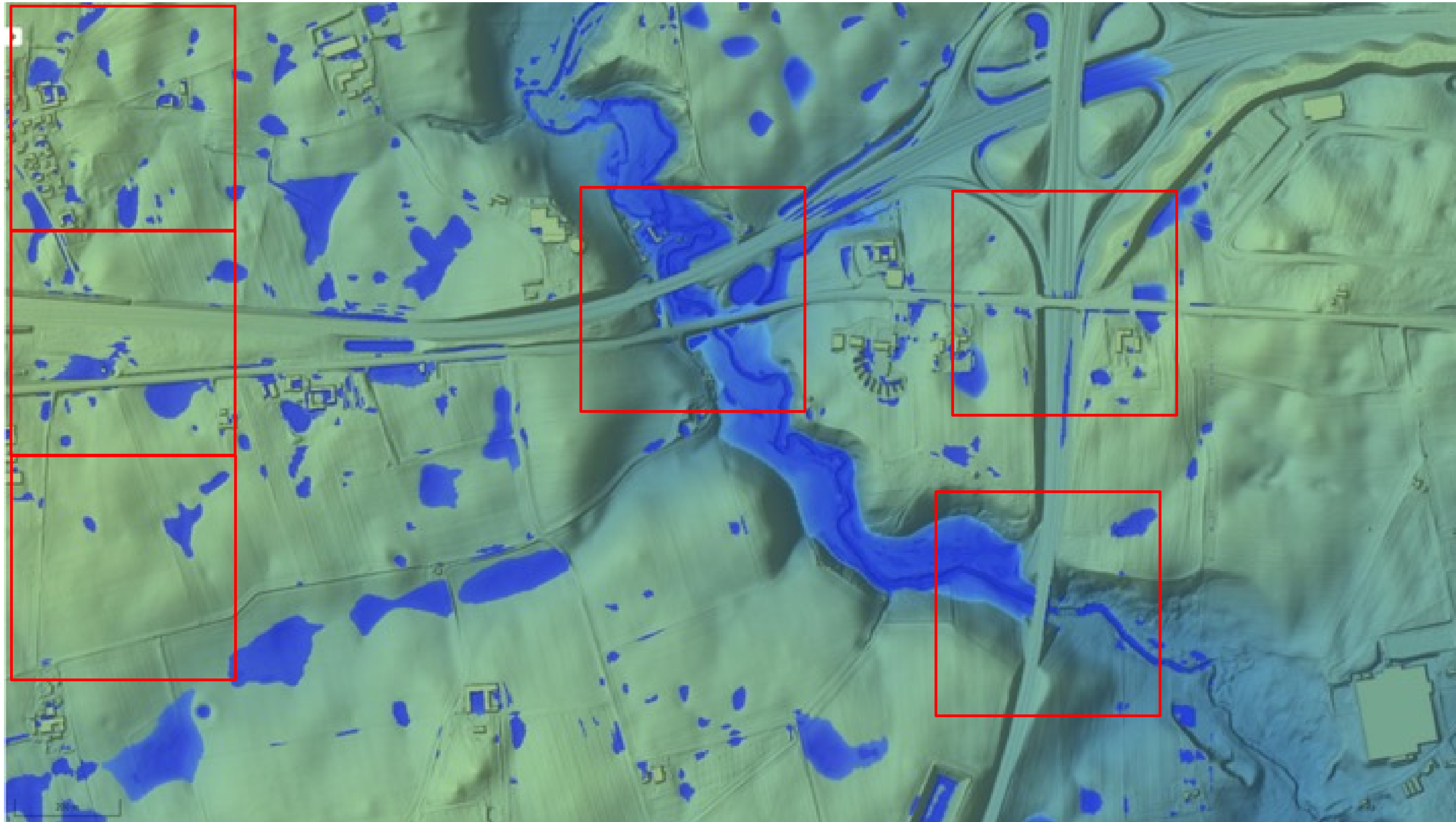
— Terrain Flood-Time Computation

— List of Corrections

AARHUS
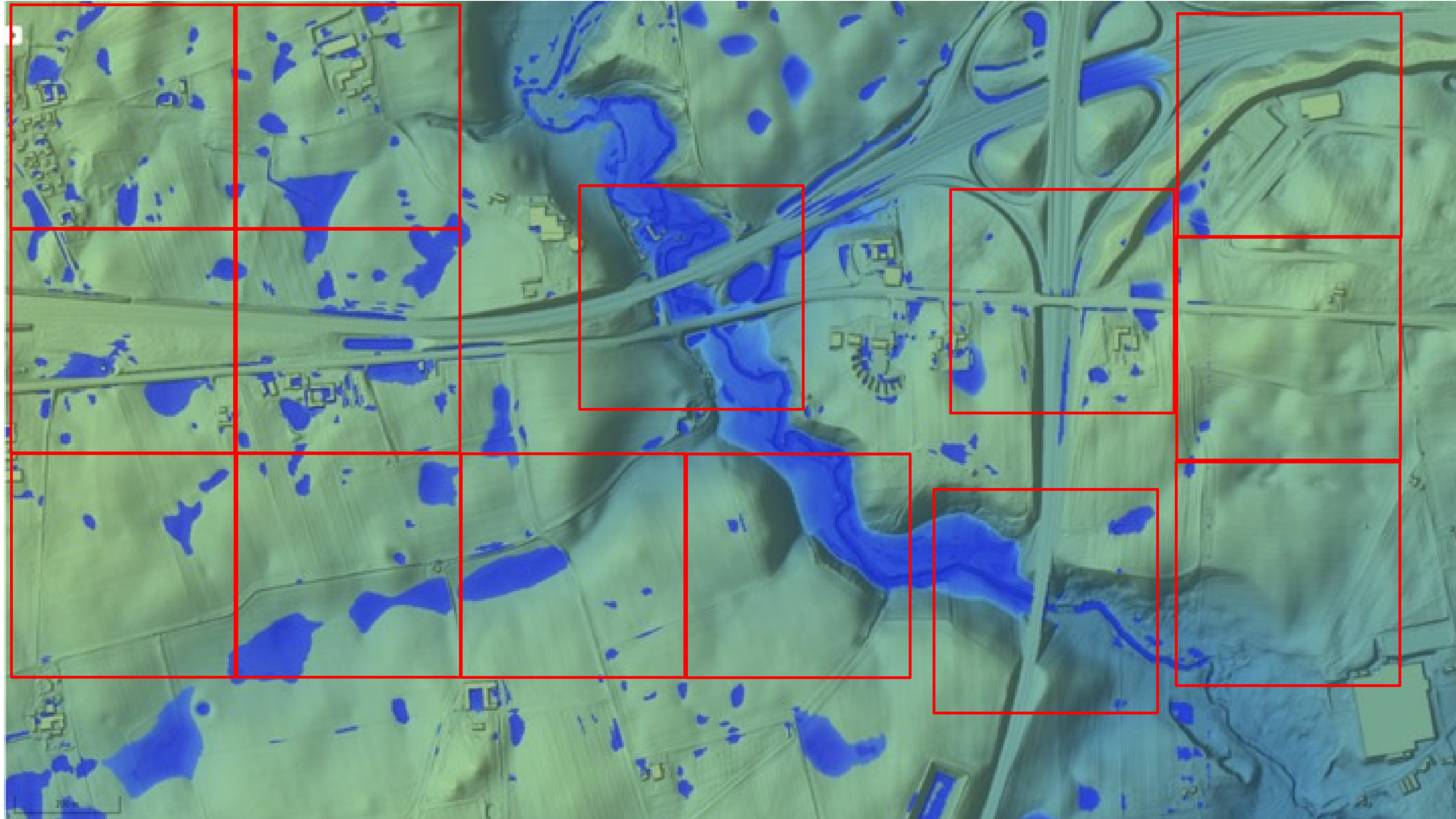UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# CREATING TILE DATA

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

# — TRAINING THE ALGORITHM

AARHUS
UNIVERSITY

# — EFFICIENTLY COMPUTING FILL FUNCTIONS

— Partition $E(\alpha)$ into $E(\beta_1)$ and $E(\beta_2)$

— Assume w.l.o.g. $|E(\beta_1)| < |E(\beta_2)|$

— $F_{\beta_1}(t) = R(\beta_1, t) + \sum_{e \in E(\beta_1)} \phi_e(t)$

— $F_\alpha(t) = F_{\beta_1}(t) + F_{\beta_2}(t)$

AARHUS
UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE