# Distributed Deep Multilevel Graph Partitioning
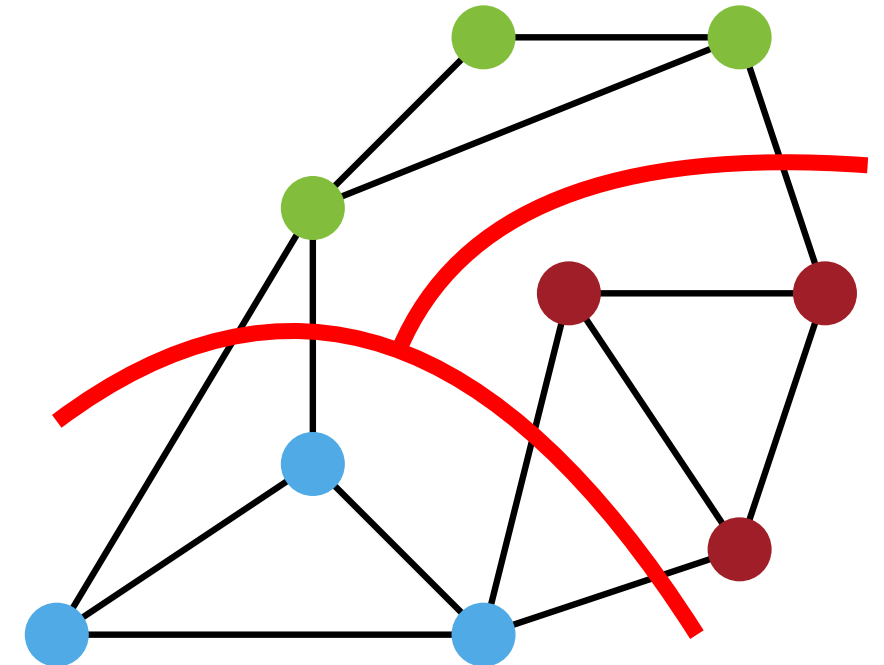
## Euro-Par 2023 · August 31, 2023

**Peter Sanders, Daniel Seemaier**

# Graph Partitioning

Given a graph $G = (V, E, c, \omega)$, partition $V$ into $k$ disjoint blocks such that:

- ■ blocks have roughly the same weight: $c(V_i) \leq (1 + \varepsilon)\frac{c(V)}{k}$

- ■ while minimizing the edge cut: $\sum_{i \neq j} \omega(E_{ij})$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning

Given a graph $G = (V, E, c, \omega)$, partition $V$ into $k$ disjoint blocks such that:

vertex weights      edge weights

- blocks have roughly the same weight: $c(V_i) \leq (1 + \varepsilon)\frac{c(V)}{k}$

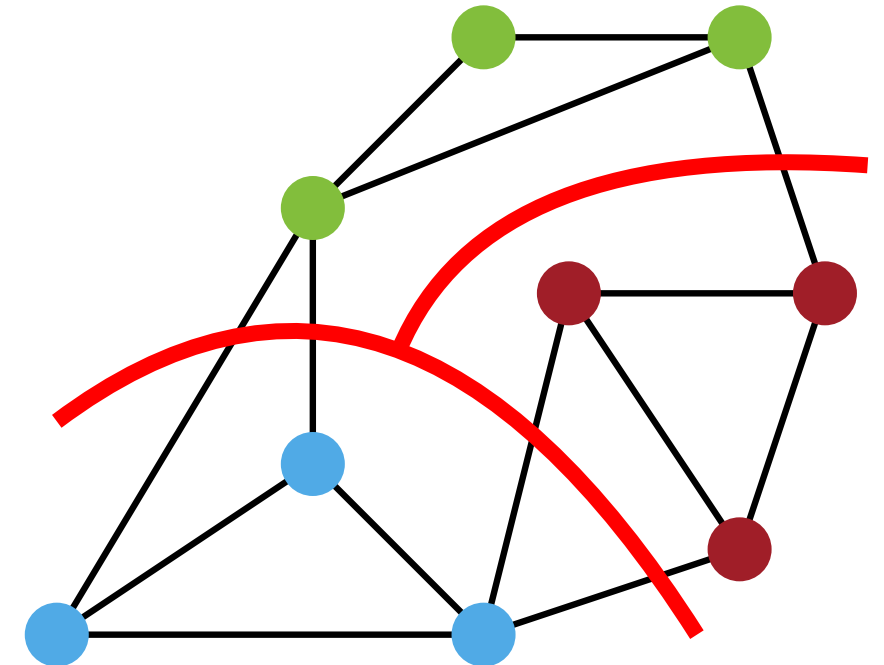- while minimizing the edge cut: $\sum_{i \neq j} \omega(E_{ij})$



Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023      Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning



Given a graph $G = (V, E, c, \omega)$, partition $V$ into $k$ disjoint blocks such that:

vertex weights        edge weights

- blocks have roughly the same weight: $c(V_i) \leq (1 + \varepsilon)\frac{c(V)}{k}$

imbalance factor

- while minimizing the edge cut: $\sum_{i \neq j} \omega(E_{ij})$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023          Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning
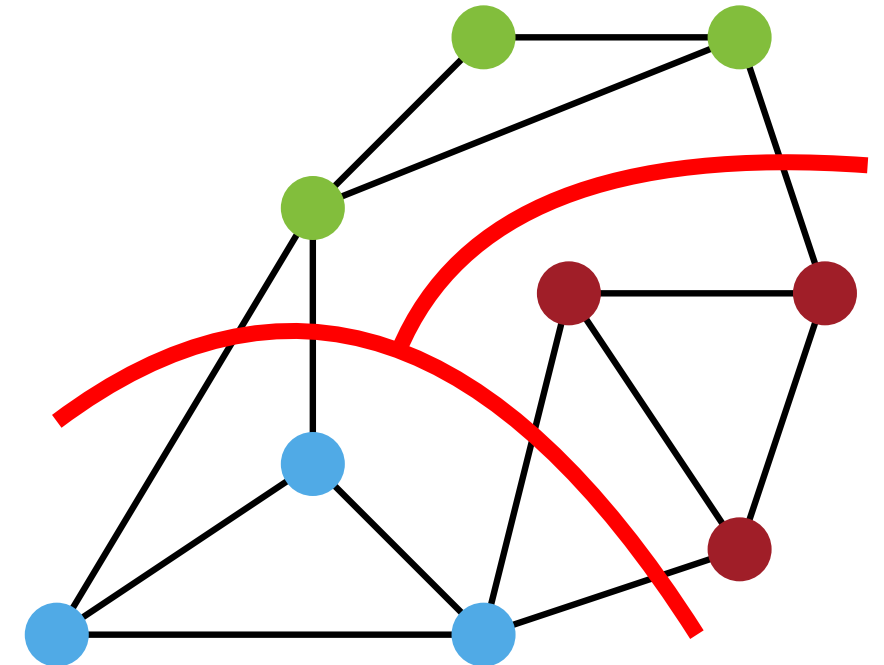
Given a graph $G = (V, E, c, \omega)$, partition $V$ into $k$ disjoint blocks such that:

vertex weights          edge weights

- blocks have roughly the same weight: $c(V_i) \leq (1 + \varepsilon)\frac{c(V)}{k}$

imbalance factor

- while minimizing the edge cut: $\sum_{i \neq j} \omega(E_{ij})$

edges between blocks $i$ and $j$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning

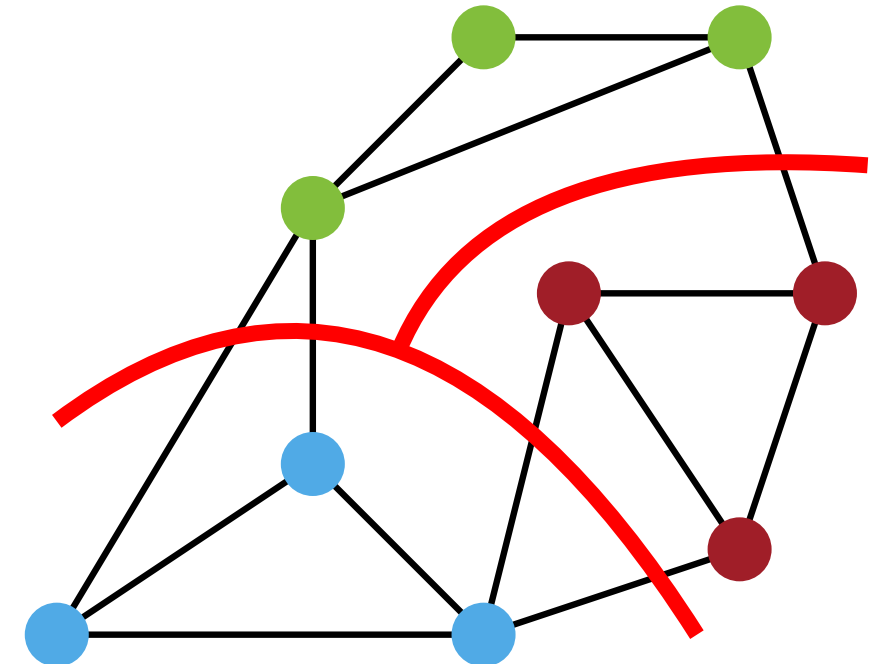Given a graph $G = (V, E, c, \omega)$, partition $V$ into $k$ disjoint blocks such that:

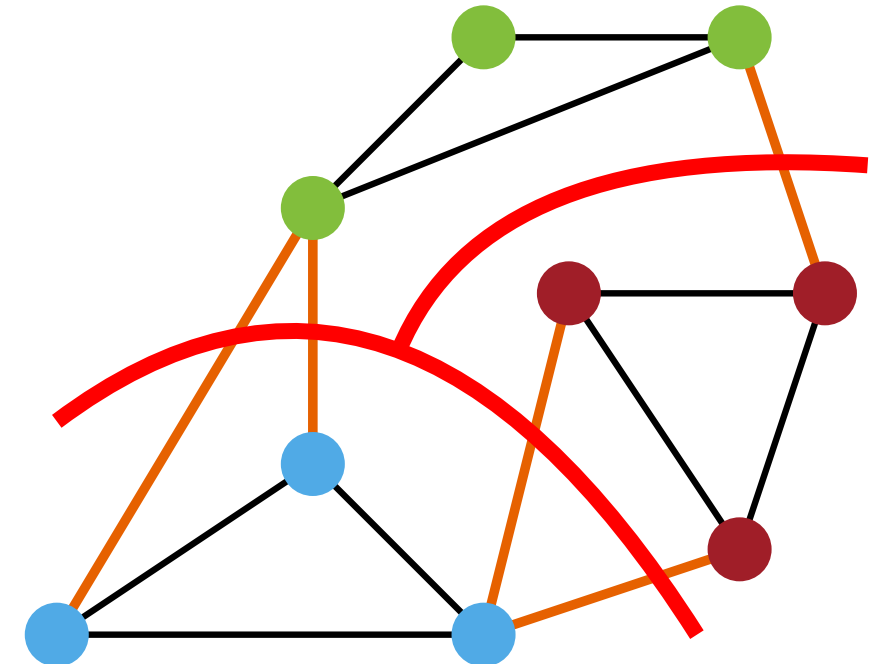vertex weights     edge weights

- blocks have roughly the same weight: $c(V_i) \leq (1 + \varepsilon)\frac{c(V)}{k}$

imbalance factor

- while minimizing the edge cut: $\sum_{i \neq j} \omega(E_{ij}) = 5$

edges between blocks $i$ and $j$



Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning for Parallel Computing

- Distributed graph processing:
  minimize communication between PEs



[HoreKa, KIT]

- Available parallelism increases steadily

- Established distributed GPs tools are not designed to handle large $k$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning for Parallel Computing

- Distributed graph processing:
  minimize communication between PEs



[HoreKa, KIT]

- Available parallelism increases steadily

- Established distributed GPs tools are not designed to handle large $k$

contribution: improve scalability to large $k$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Partitioning for Parallel Computing

- Distributed graph processing:
  minimize communication between PEs

- Available parallelism increases steadily

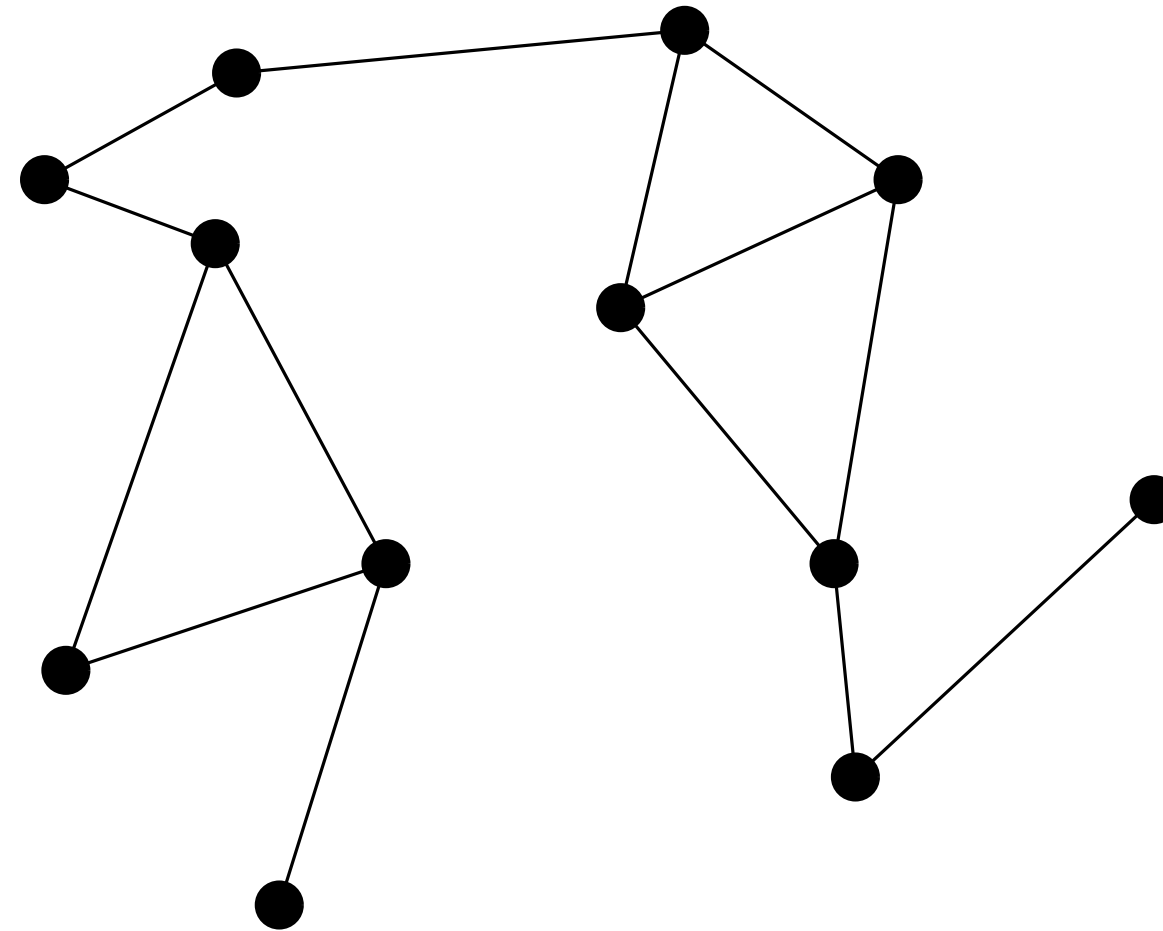  > Graph partitioning is NP-complete
  > $\Rightarrow$ we focus on heuristics

- Established distributed GPs tools are not designed to handle large $k$

  contribution: improve scalability to large $k$

[HoreKa, KIT]

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Representation

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Representation



Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Graph Representation



ghost vertices

interface vertices

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Multilevel Graph Partitioning

input
graph

# Multilevel Graph Partitioning



coarsening

input
graph

contract
(distributed)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Multilevel Graph Partitioning



input
graph

contract
(distributed)

coarsening

"small" graph $\approx kC$ vertices

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Multilevel Graph Partitioning



input graph

contract (distributed)

coarsening

"small" graph $\approx kC$ vertices

copy "small" graph to each PE

IP

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Multilevel Graph Partitioning



Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-
Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Multilevel Graph Partitioning



How do we get from 2 to $k$ blocks?

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# MGP: Direct *k*-way

# MGP: Direct *k*-way

input
graph

contract
(distributed)

*k*-way local improvement
(distributed)

output
partition

coarsening

"small" graph $\approx$ kC vertices

refinement

*k*-way IP

(in-memory)

# MGP: Direct $k$-way



+ linear time algorithm

+ $k$-way local improvement

- collapses for $kC \approx n$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



$n_1$

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



invariant: $k_i = \frac{n_i}{C}$

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



invariant: $k_i = \frac{n_i}{C}$

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-
Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



invariant: $k_i = \frac{n_i}{C}$

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Distributed Deep MGP



invariant: $k_i = \frac{n_i}{C}$

PEs

# Distributed Deep MGP



+ $kC \approx n$ not a problem

+ $k$-way local improvement

+ $\mathcal{O}(m)$ work if $kC < n$

■ extra $\log(kC/n)$ factor if $kC > n$

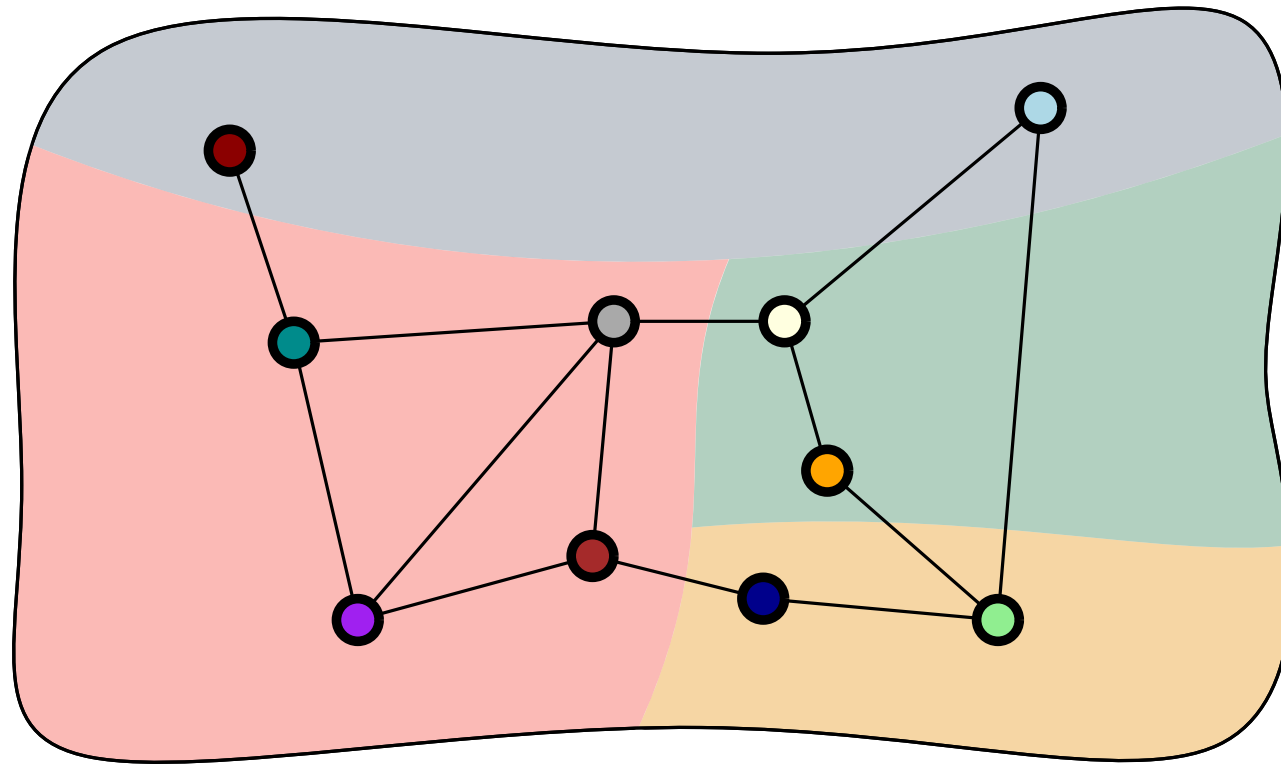invariant: $k_i = \frac{n_i}{C}$

PEs

# dKaMinPar: Distributed Deep MGP

- Our contribution: **dKaMinPar** – graph partitioner based on distributed deep MGP

  - Scalable implementation of distributed deep MGP
  - Scalable balanced coarsening and refinement algorithms

  - + many smaller performance improvements over previous works

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II
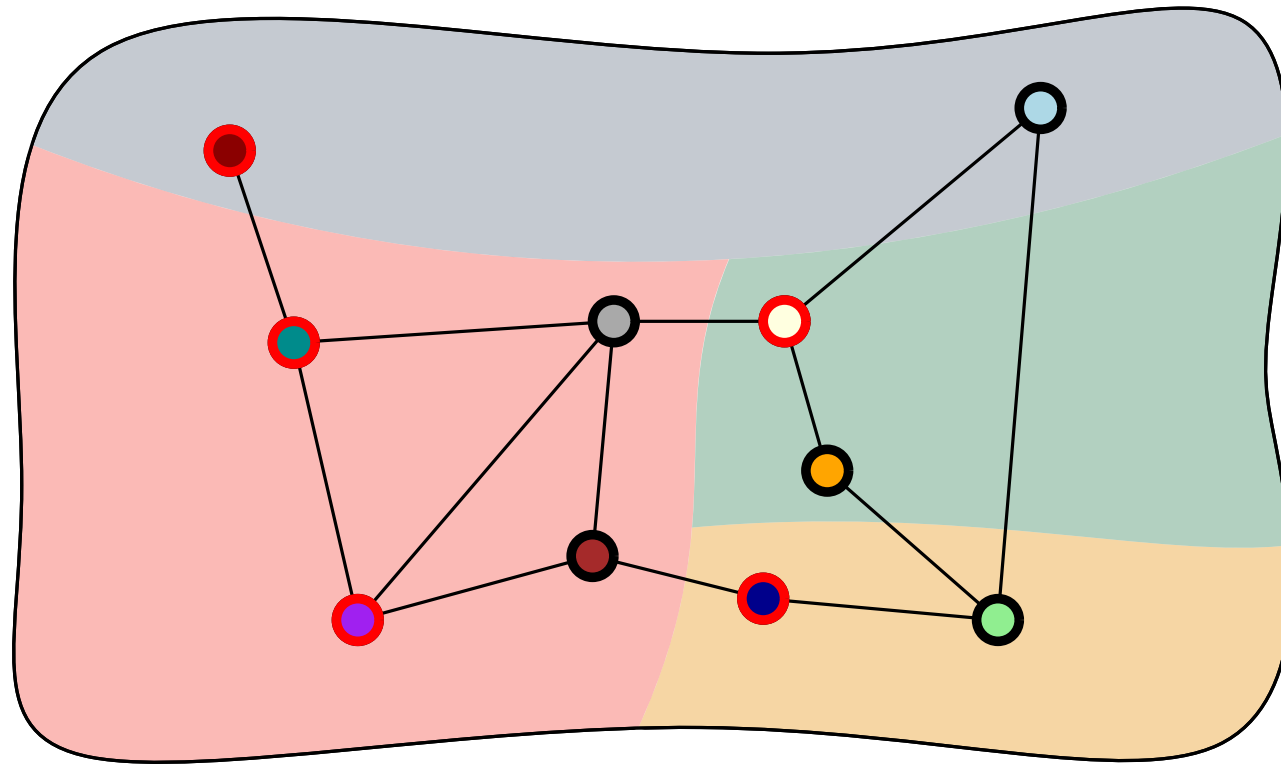
# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$



- Constant number of batches

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$



- Constant number of batches
  - Move vertices to adjacent clusters

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$



- Constant number of batches
  - Move vertices to adjacent clusters

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$



- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening



- Coarsening: use size-constrained label propagation: max. weight $W$

- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$

- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

●●● ● ● PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$
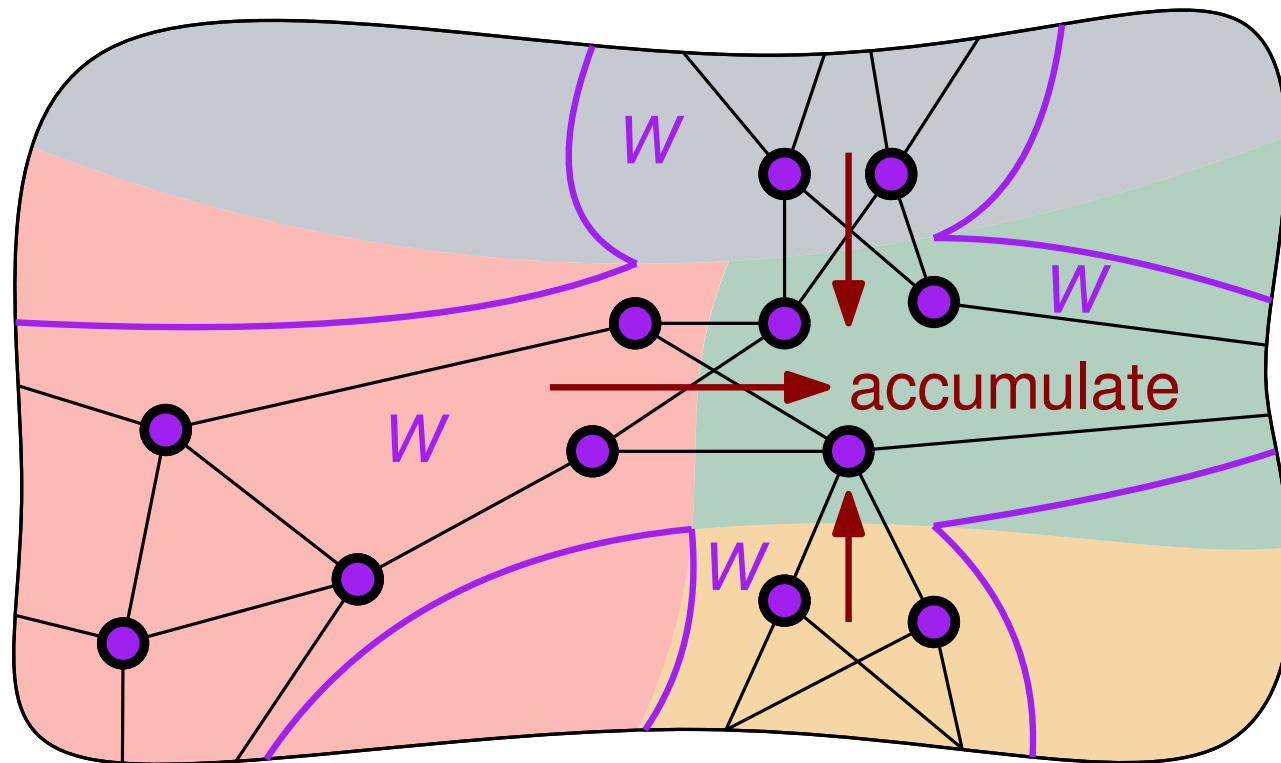


- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

- Problem: prevent huge clusters

● ● ● ● PEs

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

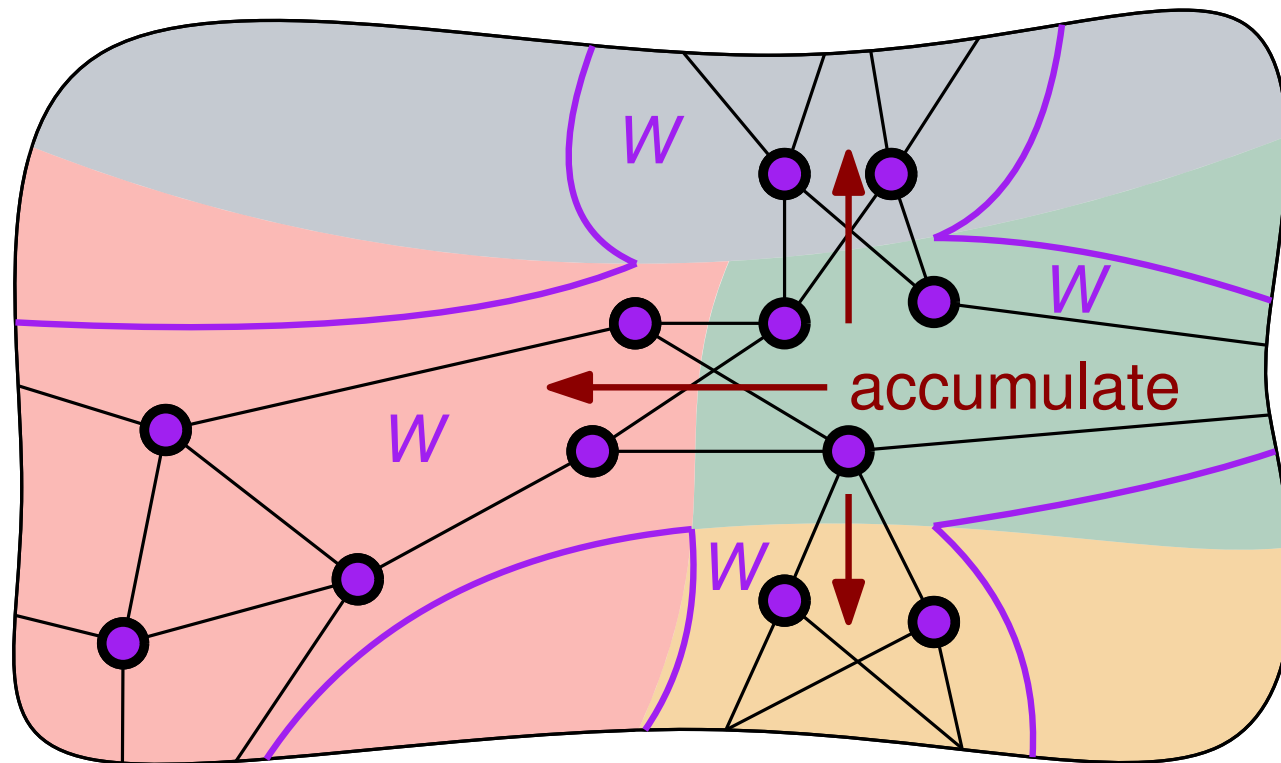- Coarsening: use size-constrained label propagation: **max. weight $W$**



- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

- Problem: prevent huge clusters

⬤⬤⬤⬤ PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

- Coarsening: use size-constrained label propagation: max. weight $W$



- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

- Problem: prevent huge clusters

PEs

# dKaMinPar – Coarsening

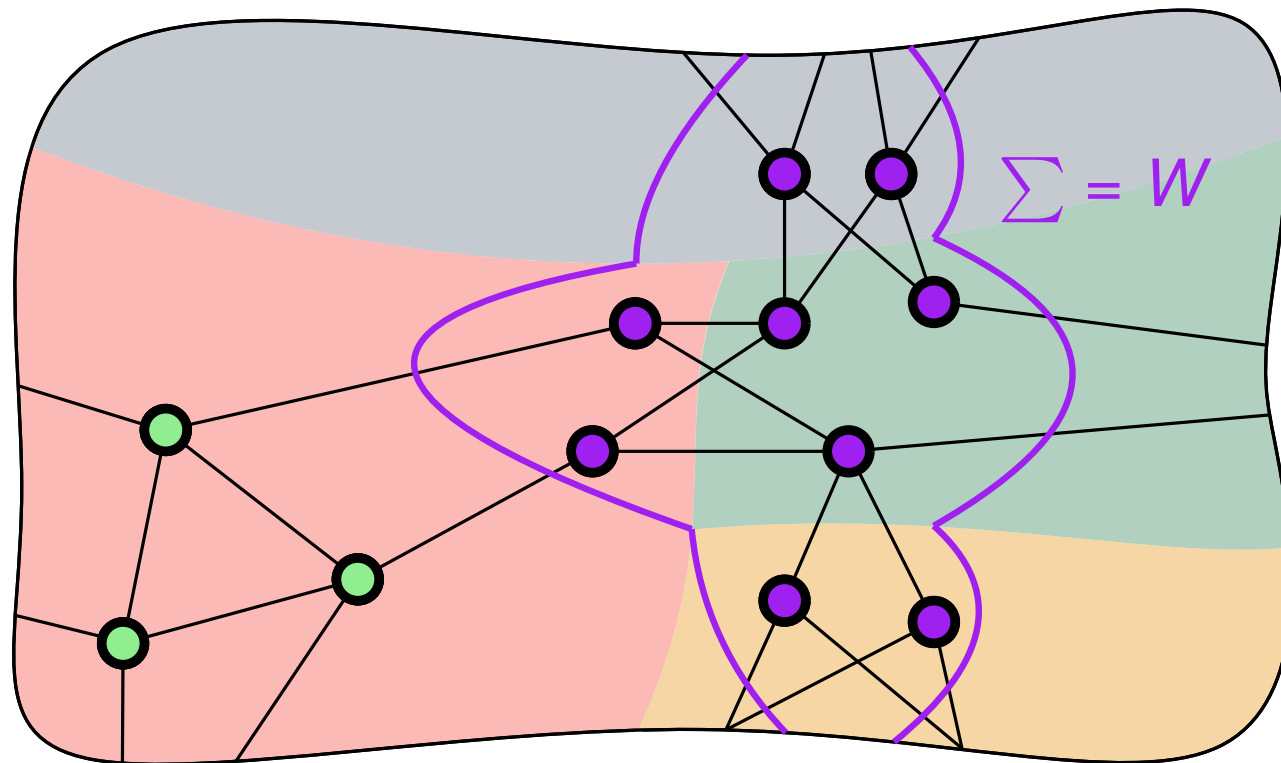- Coarsening: use size-constrained label propagation: max. weight $W$
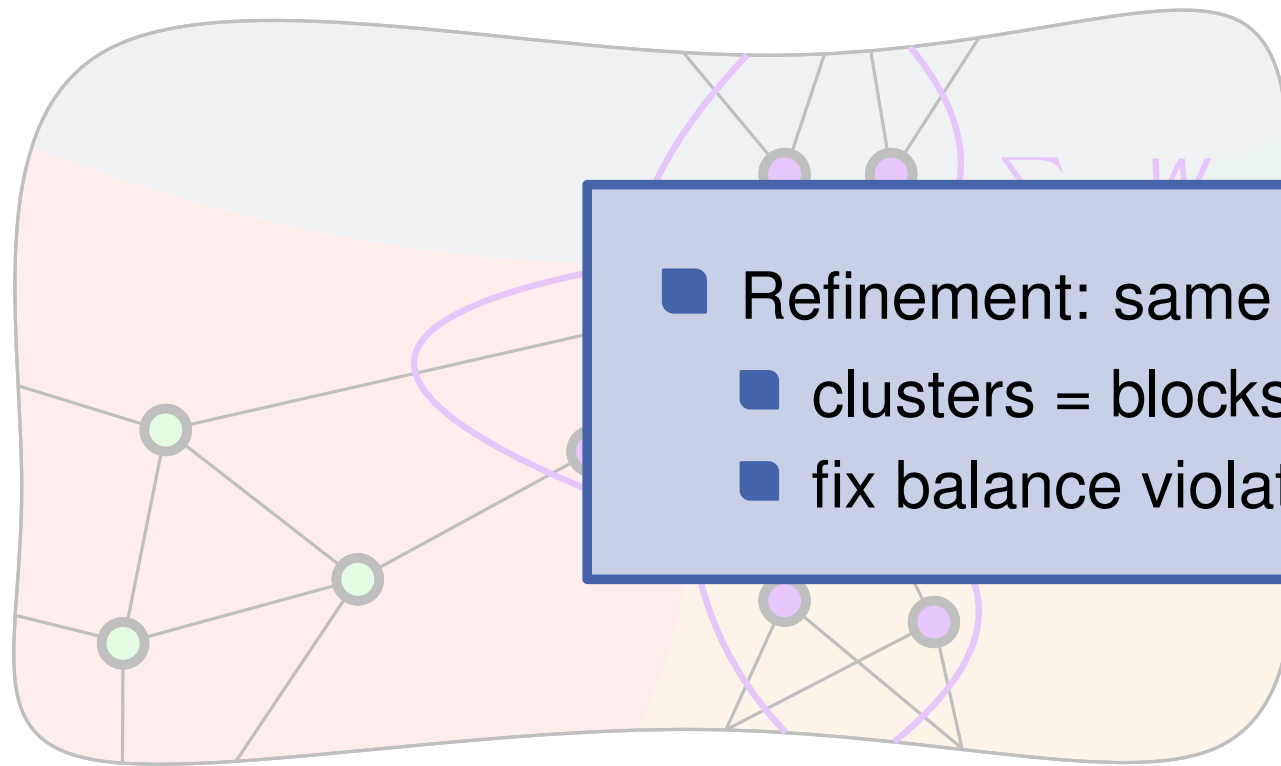
$$\Sigma = W$$

- Constant number of batches
  - Move vertices to adjacent clusters
  - Exchange labels

- Problem: prevent huge clusters
  - Revert label changes prop. to local cluster weight

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Coarsening

Coarsening: use size-constrained label propagation: max. weight $W$

Constant number of batches

Move vertices to adjacent clusters

labels

Refinement: same algorithm

  clusters = blocks

  fix balance violations afterwards

ent huge clusters

changes prop. to local

cluster weight

PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments

# Experiments – Scalability

- HoreKa: used $\{1, 2, \ldots, 128\}$ nodes $â$
  - $2\times$ Intel Xeon Platinum 8368 @ 2.40 GHz
  - 256 GB RAM

- Benchmark sets:
  - weak scaling: rand. geometric + rand. hyperbolic graphs
  - strong scaling: rand. geometric + 5 real world graphs

- Comparing **dKaMinPar** against:
  - ParHiP
  - ParMETIS
  - (XtraPuLP)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                     Institute of Theoretical Informatics, Algorithmics II

# Experiments – Scalability

- HoreKa: used $\{1, 2, \ldots, 128\}$ nodes $\hat{a}$
  - $2\times$ Intel Xeon Platinum 8368 @ 2.40 GHz  (used 64 out of 78 cores)
  - 256 GB RAM

- Benchmark sets:
  - weak scaling: rand. geometric + rand. hyperbolic graphs
  - strong scaling: rand. geometric + 5 real world graphs

- Comparing **dKaMinPar** against:
  - ParHiP
  - ParMETIS
  - (XtraPuLP)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Scalability

- HoreKa: used $\{1, 2, \ldots, 128\}$ nodes $\hat{a}$
  - $2\times$ Intel Xeon Platinum 8368 @ 2.40 GHz  (used 64 out of 78 cores)
  - 256 GB RAM

- Benchmark sets:  up to $2^{39}$ edges on $2^{13}$ cores
  - weak scaling: rand. geometric + rand. hyperbolic graphs
  - strong scaling: rand. geometric + 5 real world graphs

- Comparing **dKaMinPar** against:
  - ParHiP
  - ParMETIS
  - (XtraPuLP)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023
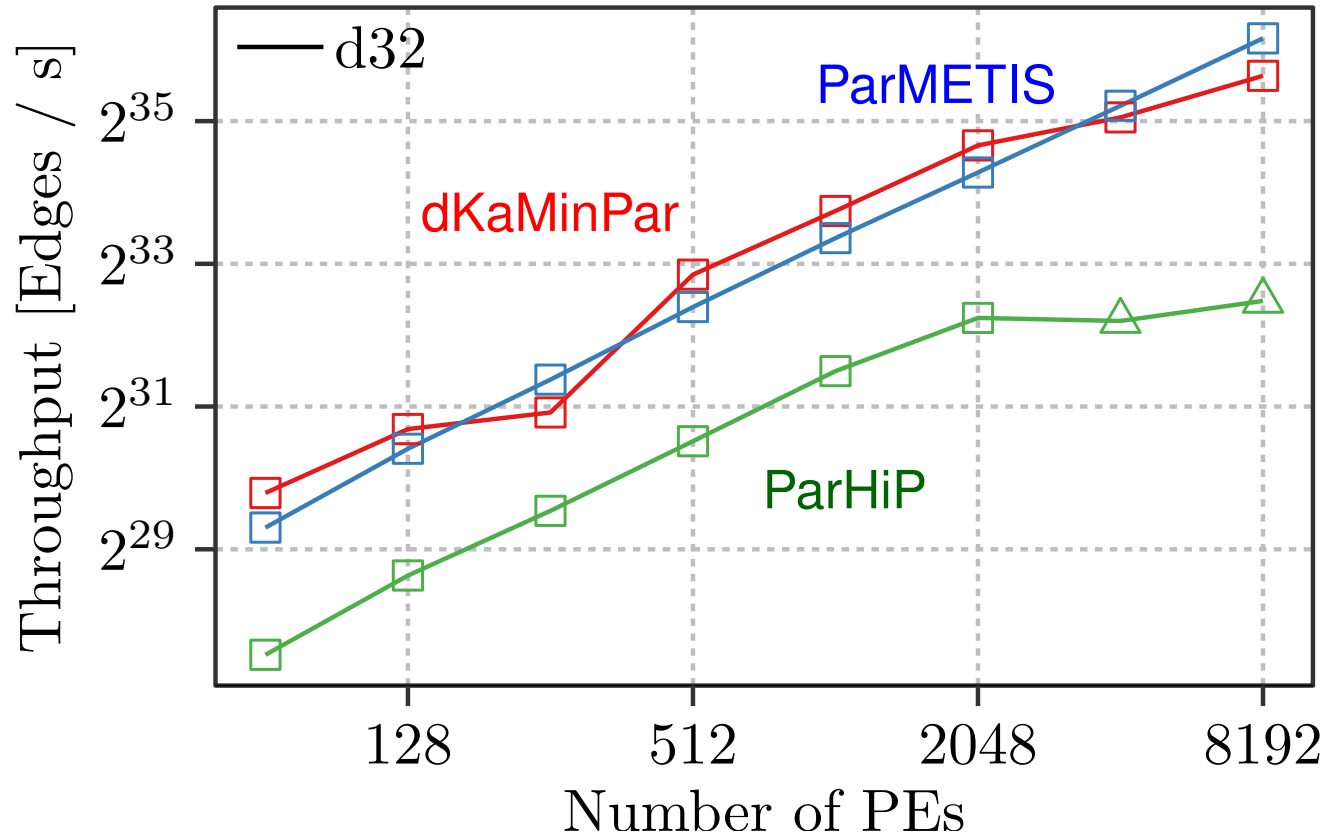
# Experiments – Scalability

- HoreKa: used $\{1, 2, \ldots, 128\}$ nodes $\hat{a}$
  - $2\times$ Intel Xeon Platinum 8368 @ 2.40 GHz  (used 64 out of 78 cores)
  - 256 GB RAM

- Benchmark sets:  up to $2^{39}$ edges on $2^{13}$ cores
  - weak scaling: rand. geometric + rand. hyperbolic graphs
  - strong scaling: rand. geometric + 5 real world graphs  $3\times$ "irregular", $3\times$ "regular"

- Comparing **dKaMinPar** against:
  - ParHiP
  - ParMETIS
  - (XtraPuLP)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Scalability

- HoreKa: used $\{1, 2, \ldots, 128\}$ nodes $â$
  - $2\times$ Intel Xeon Platinum 8368 @ 2.40 GHz  (used 64 out of 78 cores)
  - 256 GB RAM

- Benchmark sets:  up to $2^{39}$ edges on $2^{13}$ cores
  - weak scaling: rand. geometric + rand. hyperbolic graphs
  - strong scaling: rand. geometric + 5 real world graphs  $3\times$ "irregular", $3\times$ "regular"

- Comparing **dKaMinPar** against:
  - ParHiP    multilevel, label propagation for coarsening + refinement
  - ParMETIS
  - (XtraPuLP)

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Scalability

- HoreKa: used $\{1, 2, \ldots, 128\}$ nodes $\hat{a}$
  - $2\times$ Intel Xeon Platinum 8368 @ 2.40 GHz  (used 64 out of 78 cores)
  - 256 GB RAM

- Benchmark sets:  up to $2^{39}$ edges on $2^{13}$ cores
  - weak scaling: rand. geometric + rand. hyperbolic graphs
  - strong scaling: rand. geometric + 5 real world graphs  $3\times$ "irregular", $3\times$ "regular"

- Comparing **dKaMinPar** against:
  - ParHiP     multilevel, label propagation for coarsening + refinement
  - ParMETIS multilevel, matchings for coarsening, greedy refinement
  - (XtraPuLP)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II
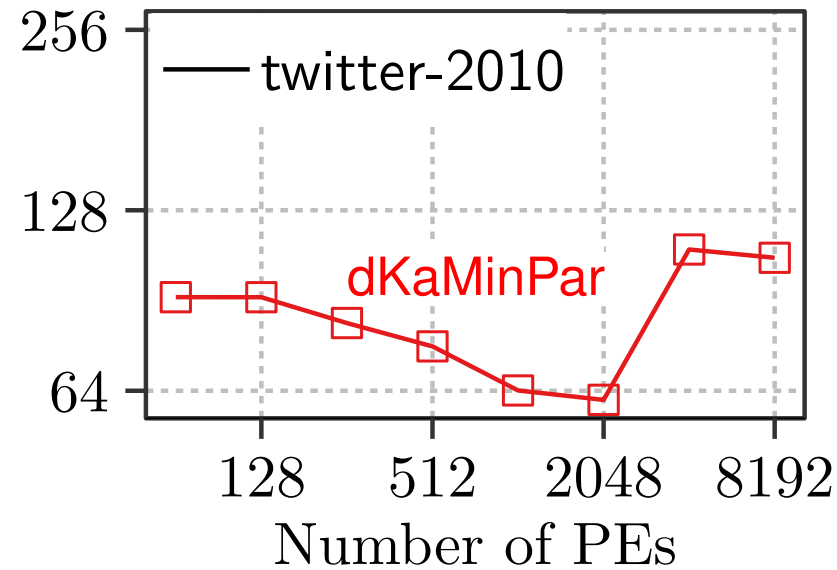
# Experiments – Weak Scaling, constant $k = 16$



rgg$_{2D}$26

rhg$_{3.0}$26

Throughput [Edges / s]

Number of PEs

d32 ParMETIS dKaMinPar ParHiP

**[$\{1, 2, \ldots, 128\}$ nodes @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Weak Scaling, constant $k = 16$



$rgg_{2D}26$

$rhg_{3.0}26$

better

Throughput [Edges / s]

Number of PEs

d32

ParMETIS

dKaMinPar

ParHiP

dKaMinPar

ParMETIS

ParHiP

$[\{1, 2, \ldots, 128\}$ **nodes @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II

# Experiments – Weak Scaling, constant $\frac{n}{k} = 2^{15}$



rgg$_{2D}$26d8

rhg$_{3.0}$26d8

better

dKaMinPar

ParMETIS

ParHiP

[$\{1, 2, \ldots, 128\}$ **nodes @ 64 cores**]

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Strong Scaling



**[{1, 2, ..., 128} nodes @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023          Institute of Theoretical Informatics, Algorithmics II

# Experiments – Strong Scaling



better

$[\{1, 2, \ldots, 128\}$ **nodes @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023          Institute of Theoretical Informatics, Algorithmics II

# Experiments – Scalability: Quality

| Graph | Cut on 64 PEs | Cut on 8192 PEs |
|---|---|---|
| kmer_V1r | 10 955 | 10 836 |
| nlpkkt240 | 5 726 | 5 547 |
| rgg27 | 353 | 347 |
| webbase-2001 | 9 674 | 9 524 |
| uk-2007-05 | 4 054 | 4 064 |
| twitter-2010 | 616 791 | 588 380 |

$$\times 1\,000$$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Scalability: Quality

| Graph | Cut on 64 PEs | Cut on 8192 PEs |
|---|---|---|
| kmer_V1r | 10 955 | 10 836 |
| nlpkkt240 | 5 726 | 5 547 |
| rgg27 | 353 | 347 |
| webbase-2001 | 9 674 | 9 524 |
| uk-2007-05 | 4 054 | 4 064 |
| twitter-2010 | 616 791 | 588 380 |

no degredation →

$\times 1\,000$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                Institute of Theoretical Informatics, Algorithmics II

# Experiments – Quality

- 64 cores of 1 AMD EPYC 7702 @ 2 GHz, 1 TB RAM

- Benchmark set: 32 graphs with $4.7\,M \leq m \leq 6.6\,G$

- $\varepsilon = 3\%$

- $k \in \{2, 4, \ldots, 128\}$

- Comparing **dKaMinPar** against:

  - ParHiP       Distributed-mem.
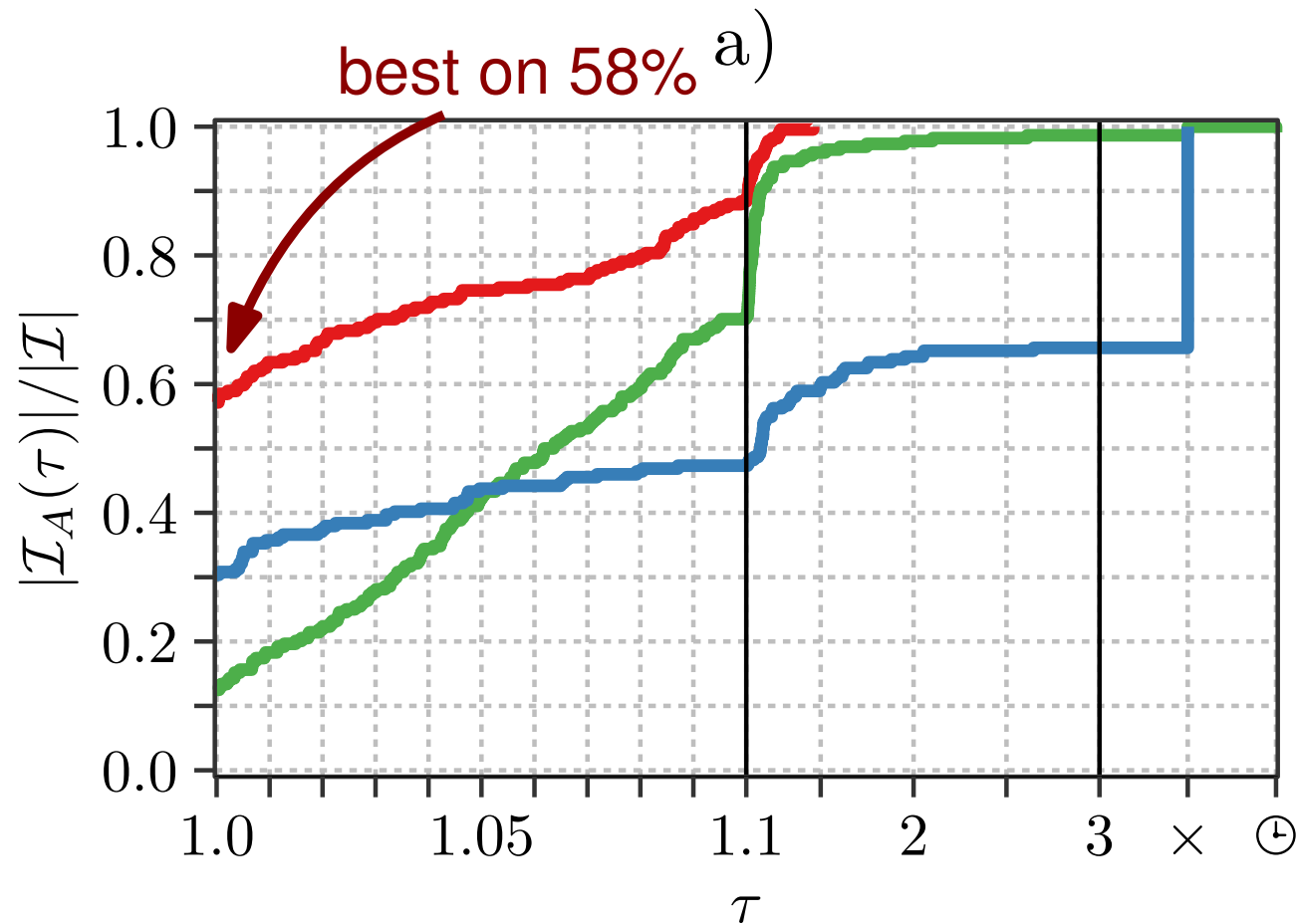  - ParMETIS

  - KaMinPar       Shared-mem.

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# Experiments – Quality

- 64 cores of 1 AMD EPYC 7702 @ 2 GHz, 1 TB RAM

- Benchmark set: 32 graphs with $4.7\,\text{M} \leq m \leq 6.6\,\text{G}$
  $\begin{cases} 17 \times \text{regular},\ 3 \leq \Delta \leq 40 \\ 15 \times \text{irregular},\ 2.7\,\text{k} \leq \Delta \leq 8.6\,\text{M} \end{cases}$

- $\varepsilon = 3\%$

- $k \in \{2, 4, \dots, 128\}$

- Comparing **dKaMinPar** against:

  - ParHiP        Distributed-mem.
  - ParMETIS

  - KaMinPar      Shared-mem.

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II

# Experiments – Quality

- 64 cores of 1 AMD EPYC 7702 @ 2 GHz, 1 TB RAM

- Benchmark set: 32 graphs with $4.7\,\mathrm{M} \leq m \leq 6.6\,\mathrm{G}$   $\begin{cases} 17 \times \text{regular, } 3 \leq \Delta \leq 40 \\ 15 \times \text{irregular, } 2.7\,\mathrm{k} \leq \Delta \leq 8.6\,\mathrm{M} \end{cases}$
- $\varepsilon = 3\%$

- $k \in \{2, 4, \ldots, 128\}$

- Comparing **dKaMinPar** against:

  - ParHiP          Distributed-mem.
  - ParMETIS

  - KaMinPar          Shared-mem.

  similar techniques implemented in shared-memory

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II

a)

**[1 node @ 64 cores]**

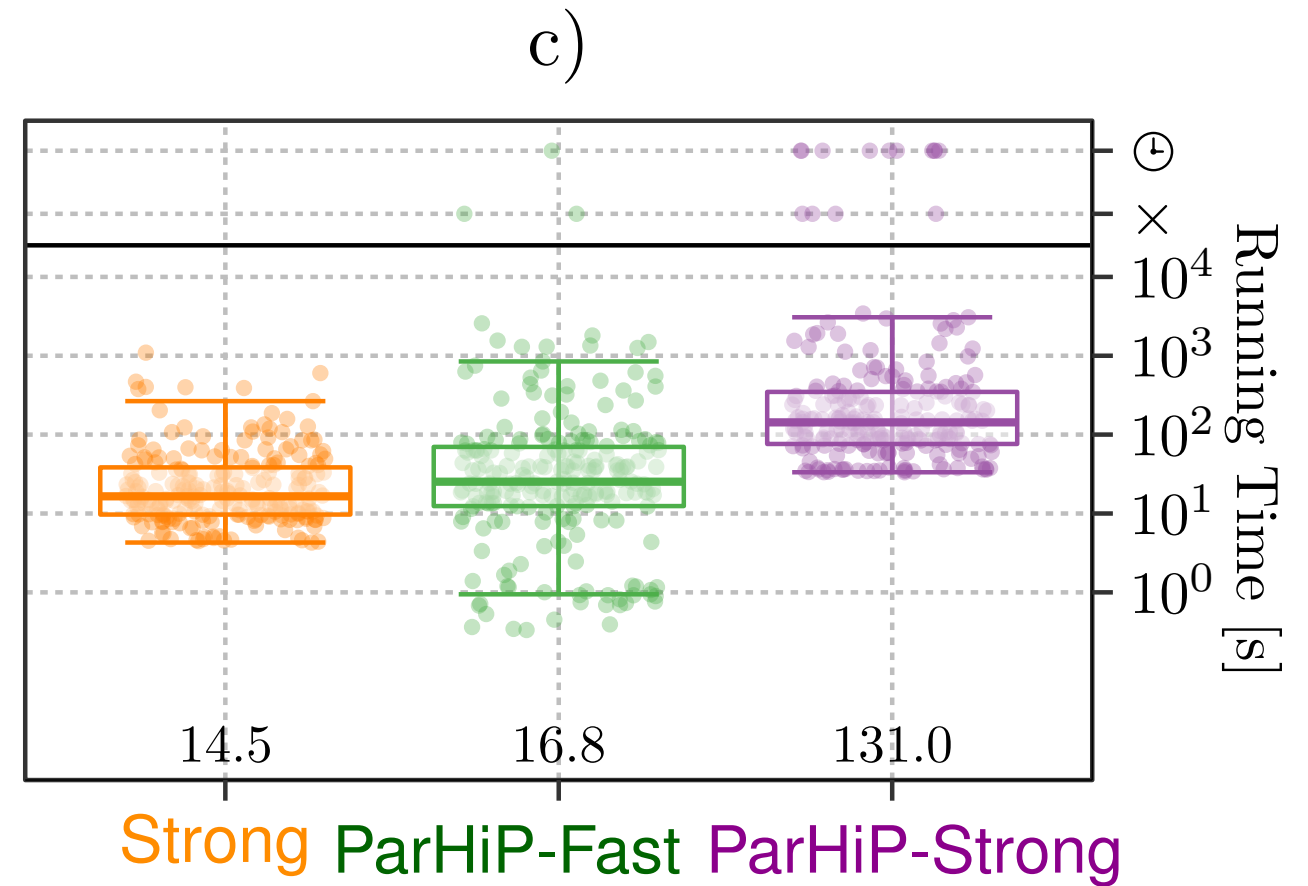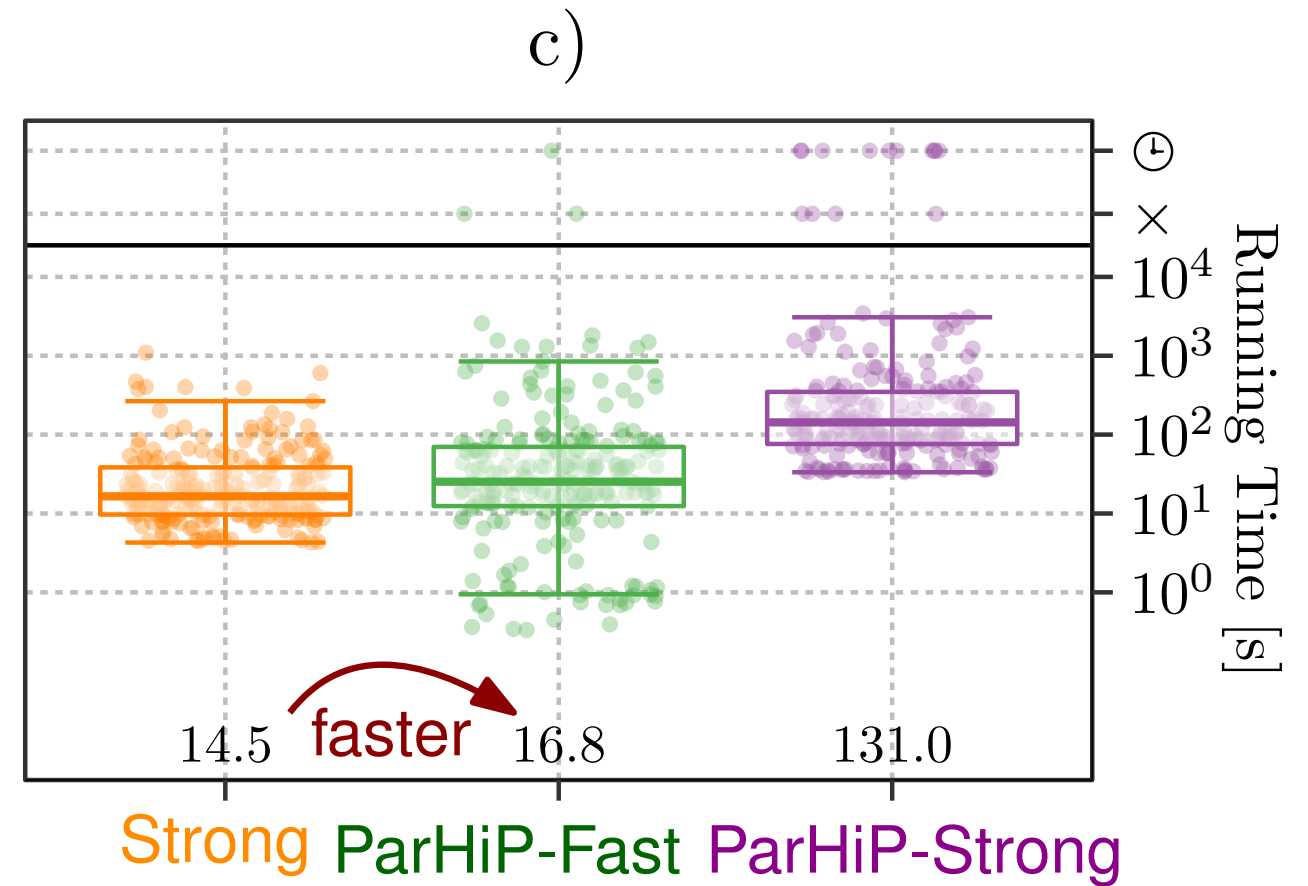Institute of Theoretical Informatics, Algorithmics II
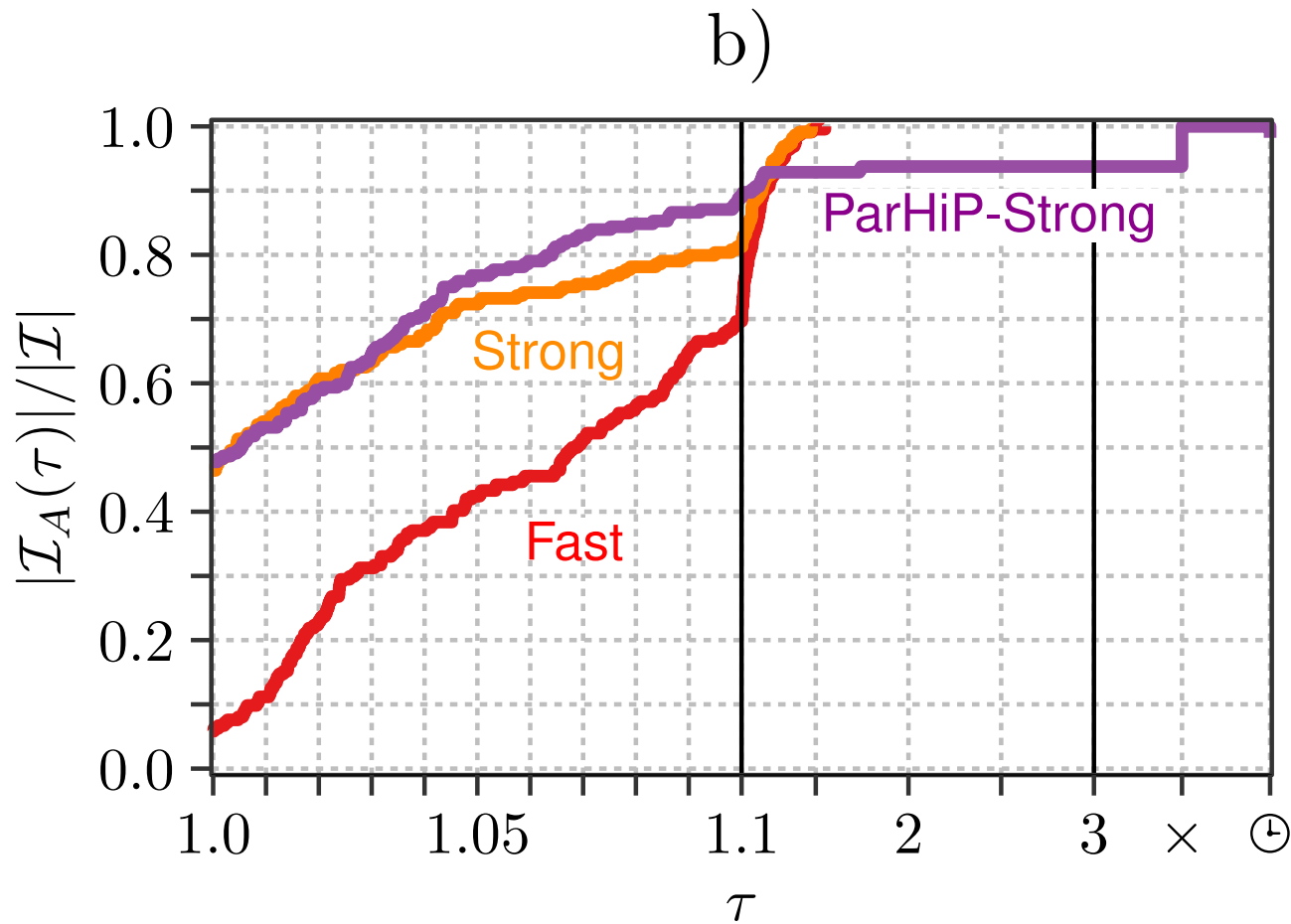
# Experiments – Quality: vs. Competitors



**[1 node @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

best on 58% a) within 1.1 on 88%

**[1 node @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

a)

**[1 node @ 64 cores]**

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

a)

[1 node @ 64 cores]

Institute of Theoretical Informatics, Algorithmics II

b)

c)

Strong ParHiP-Fast ParHiP-Strong

14.5    16.8    131.0

**[1 node @ 64 cores]**

b)

c)

Strong   ParHiP-Fast   ParHiP-Strong

14.5   faster   16.8   131.0

[1 node @ 64 cores]

a)

b)

[1 node @ 64 cores]

Institute of Theoretical Informatics, Algorithmics II

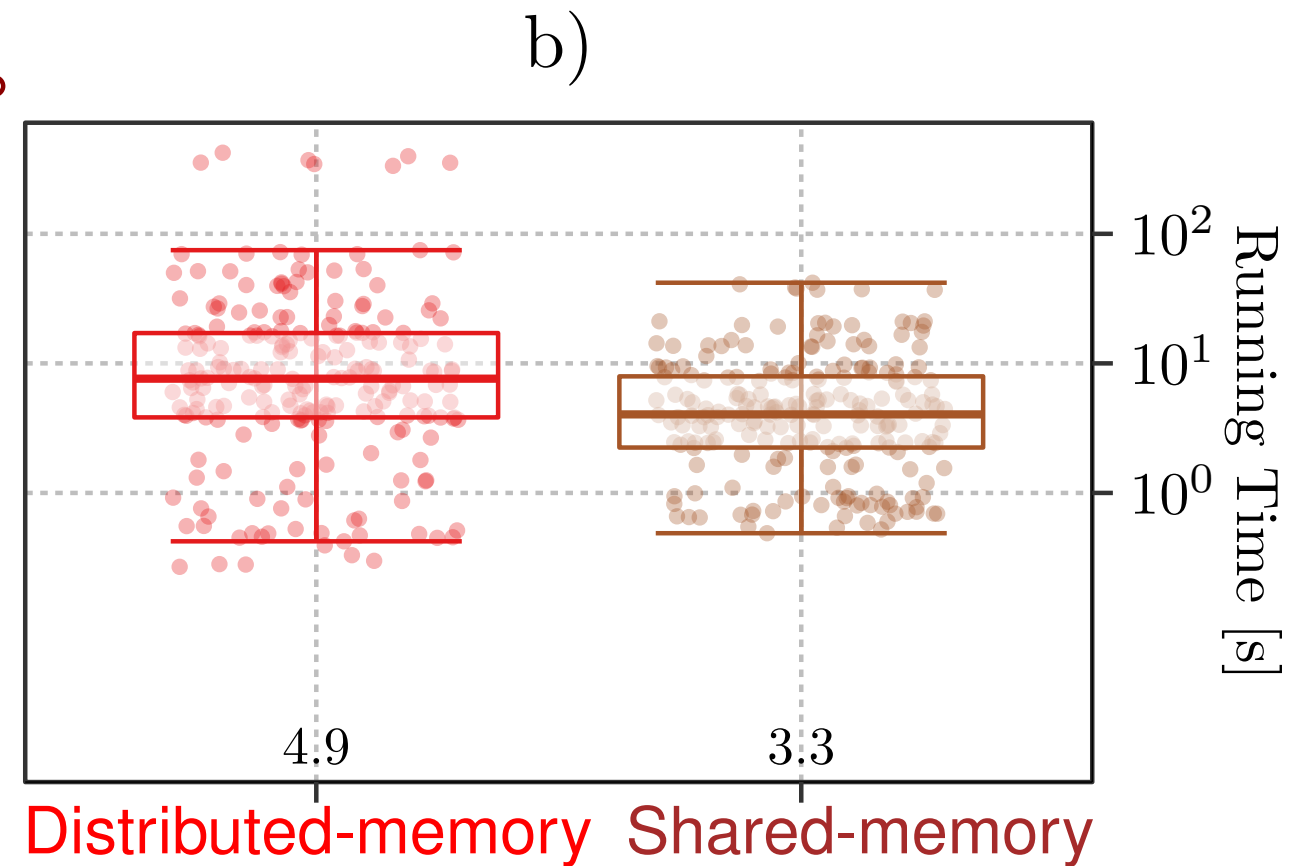# Experiments – Quality: vs. Shared-Memory



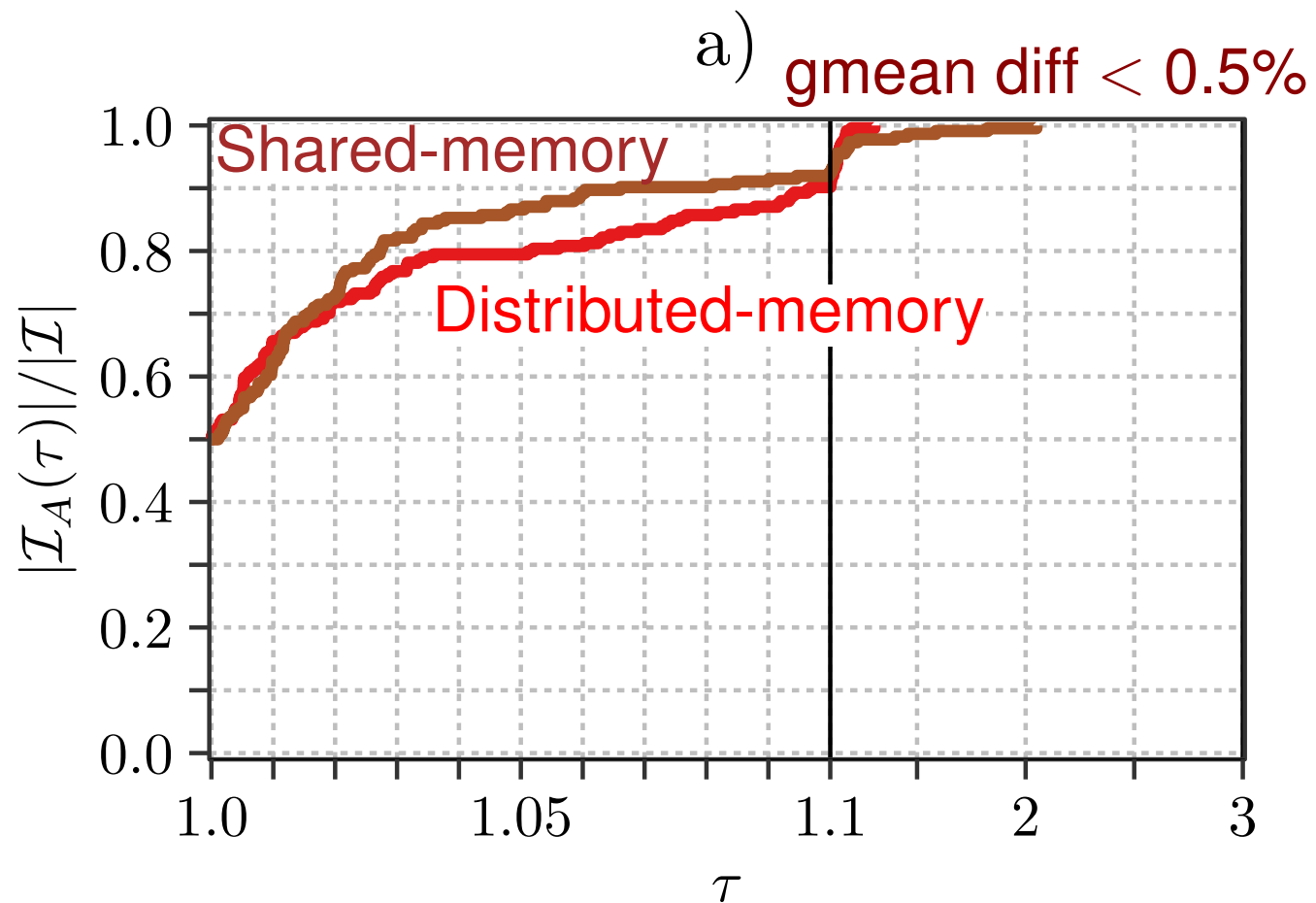a) gmean diff $< 0.5\%$

b)

**[1 node @ 64 cores]**
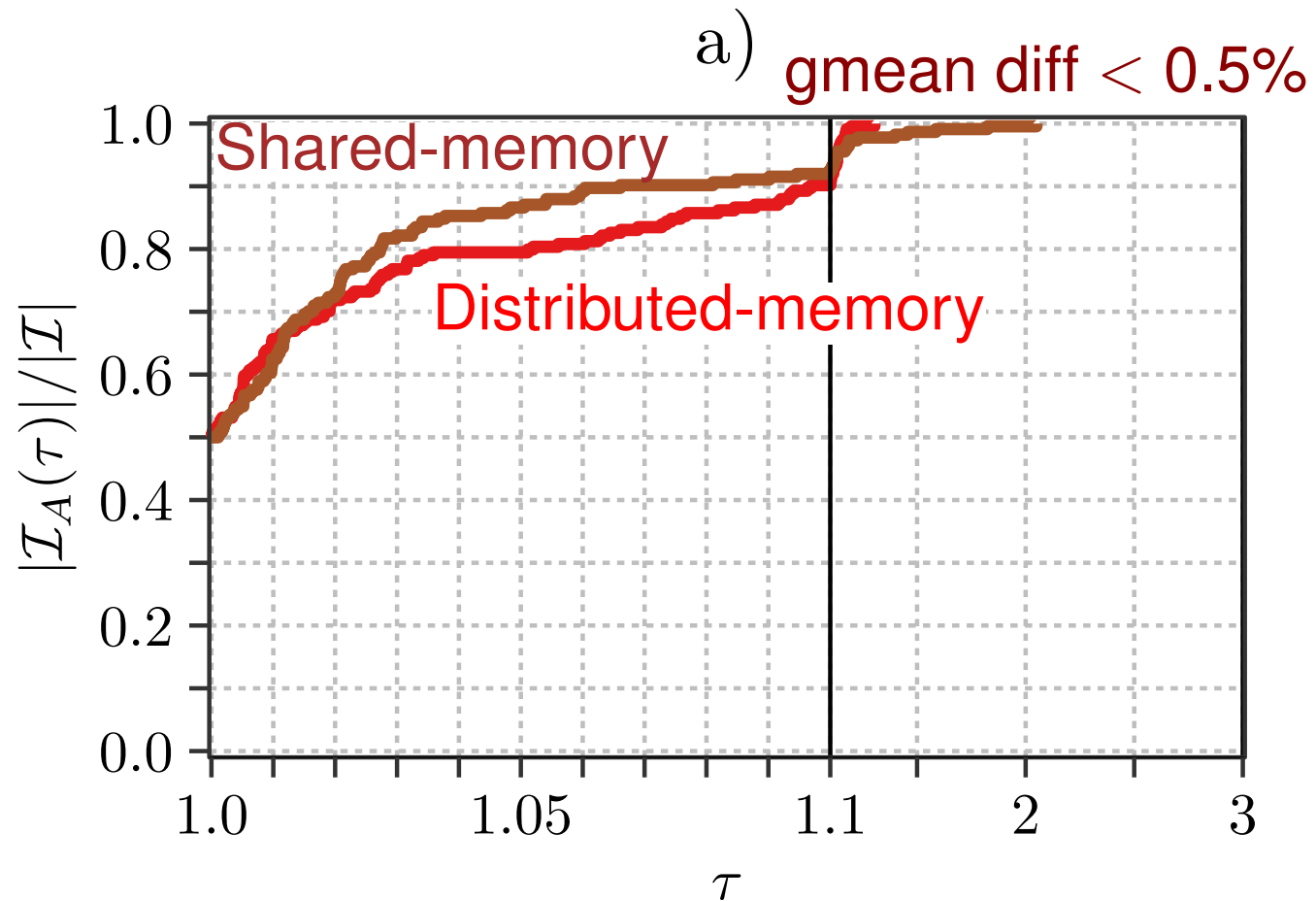
# Experiments – Quality: vs. Shared-Memory



a) gmean diff $< 0.5\%$

Shared-memory

Distributed-memory

b)

$\times 1.5$

4.9 → 3.3

Distributed-memory   Shared-memory

[1 node @ 64 cores]

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II
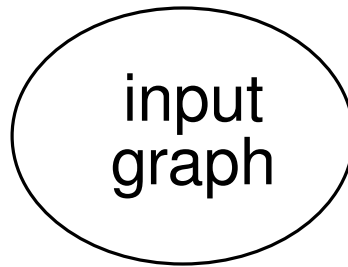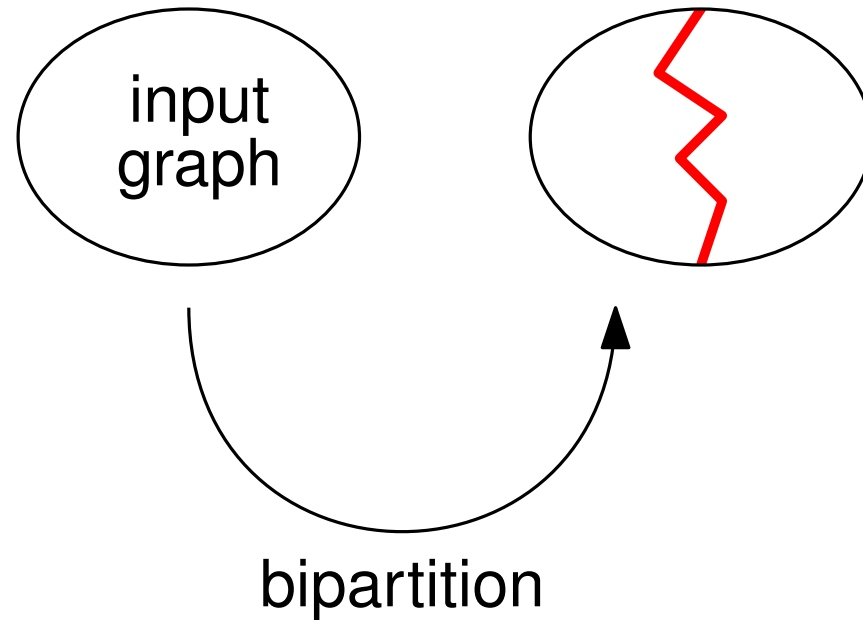
# Conclusion

- Deep Multilevel Graph Partitioning:
  - Integrate coarsening deep into initial partitioning

- **dKaMinPar**: distributed deep MGP implementation
  - Scales to thousands of PEs, competitive partition quality
  - Better scalability for large $k$ than previous approaches

- **Future**: stronger distributed refinement algorithms

- Supplementary data available online:
  - Full experimental results: `algo2.iti.kit.edu/seemaier/ddeep_mgp/`
  - Source code: `github.com/KaHIP/KaMinPar`

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II

# MGP: Recursive Bipartitioning

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023
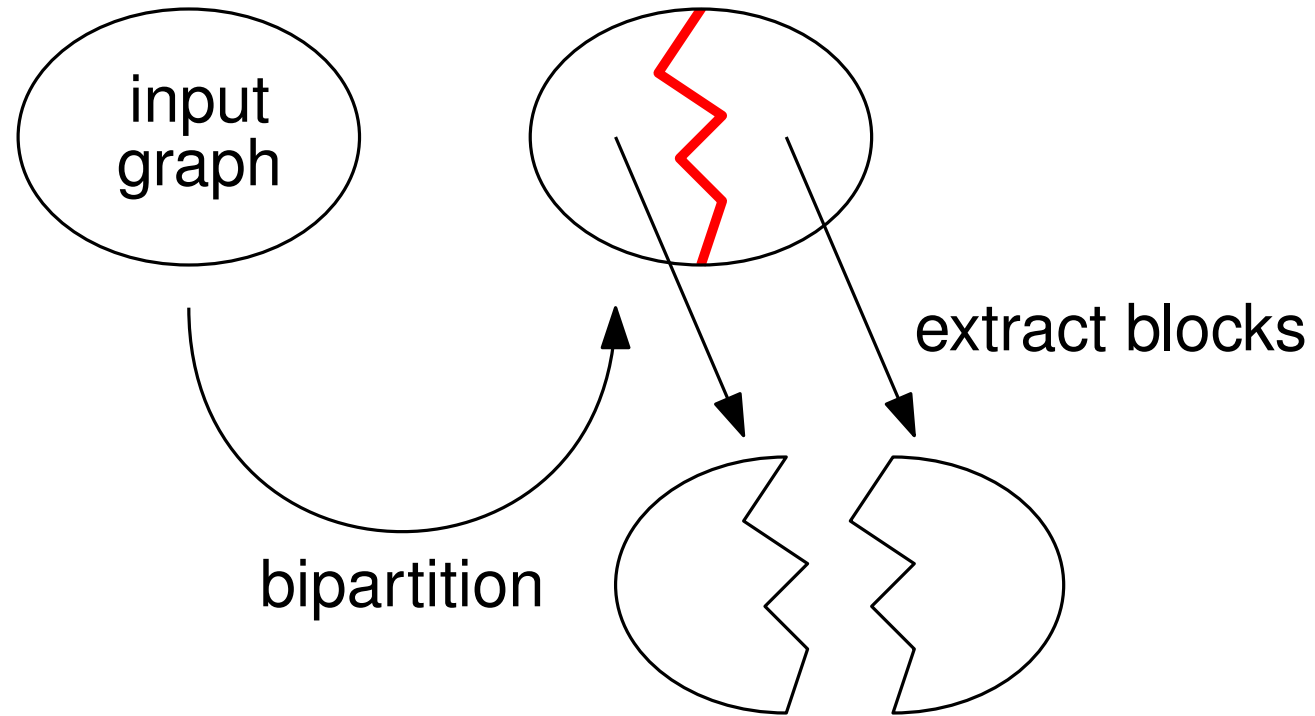
Institute of Theoretical Informatics, Algorithmics II

# MGP: Recursive Bipartitioning

# MGP: Recursive Bipartitioning

input
graph

bipartition

extract blocks

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-
Par 2023

Institute of Theoretical Informatics, Algorithmics II

# MGP: Recursive Bipartitioning



bipartition

extract blocks

...repeat...

log($k$)

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023                    Institute of Theoretical Informatics, Algorithmics II

# MGP: Recursive Bipartitioning

+ large $k$ not a problem

- log($k$) factor in running time

- no $k$-way local improvement

- no global view on $k$-way partition

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

recall: $c(V_i) \leq \max\left\{(1 + \varepsilon)\frac{c(V)}{k}, \frac{c(V)}{k} + \max_v c(v)\right\}$

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-
Par 2023

Institute of Theoretical Informatics, Algorithmics II
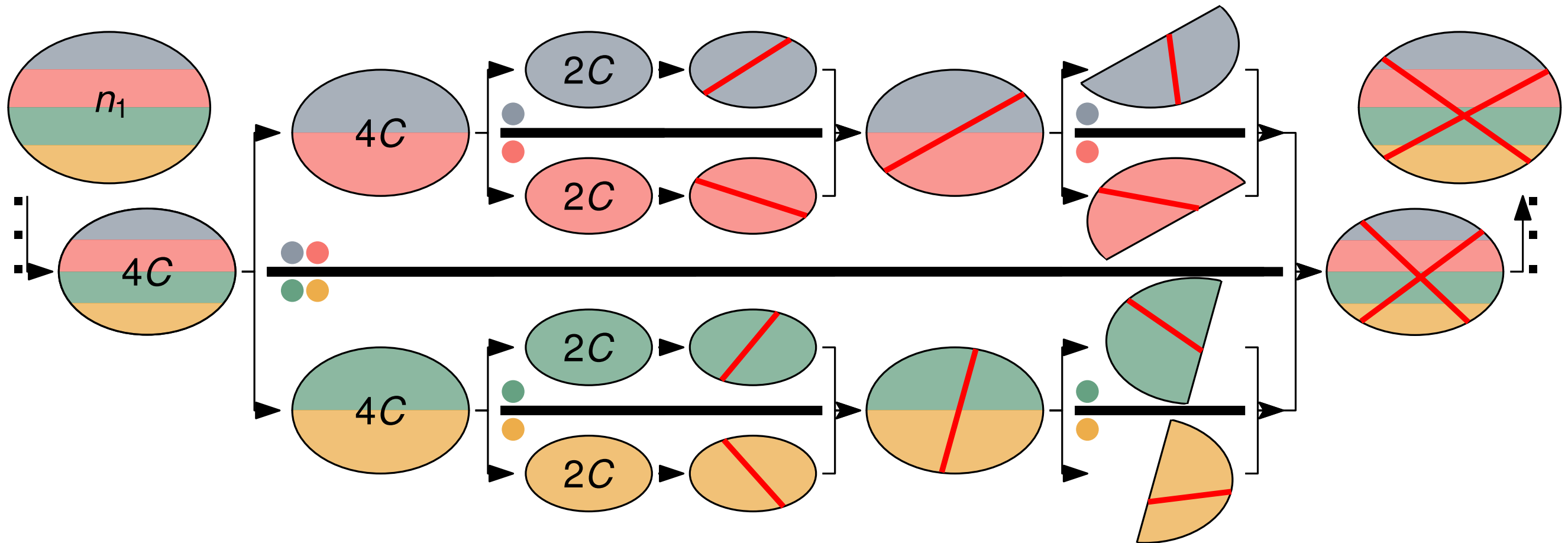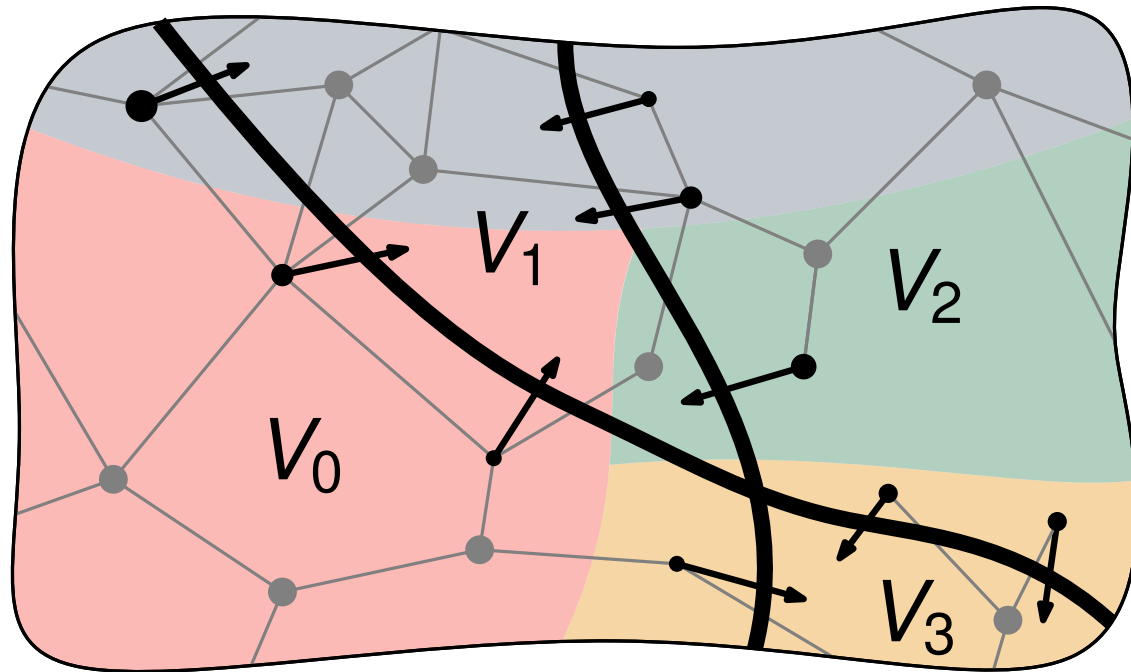
# dKaMinPar – Refinement

recall: $c(V_i) \leq \max \left\{ (1 + \varepsilon) \frac{c(V)}{k}, \frac{c(V)}{k} + \max_v c(v) \right\}$
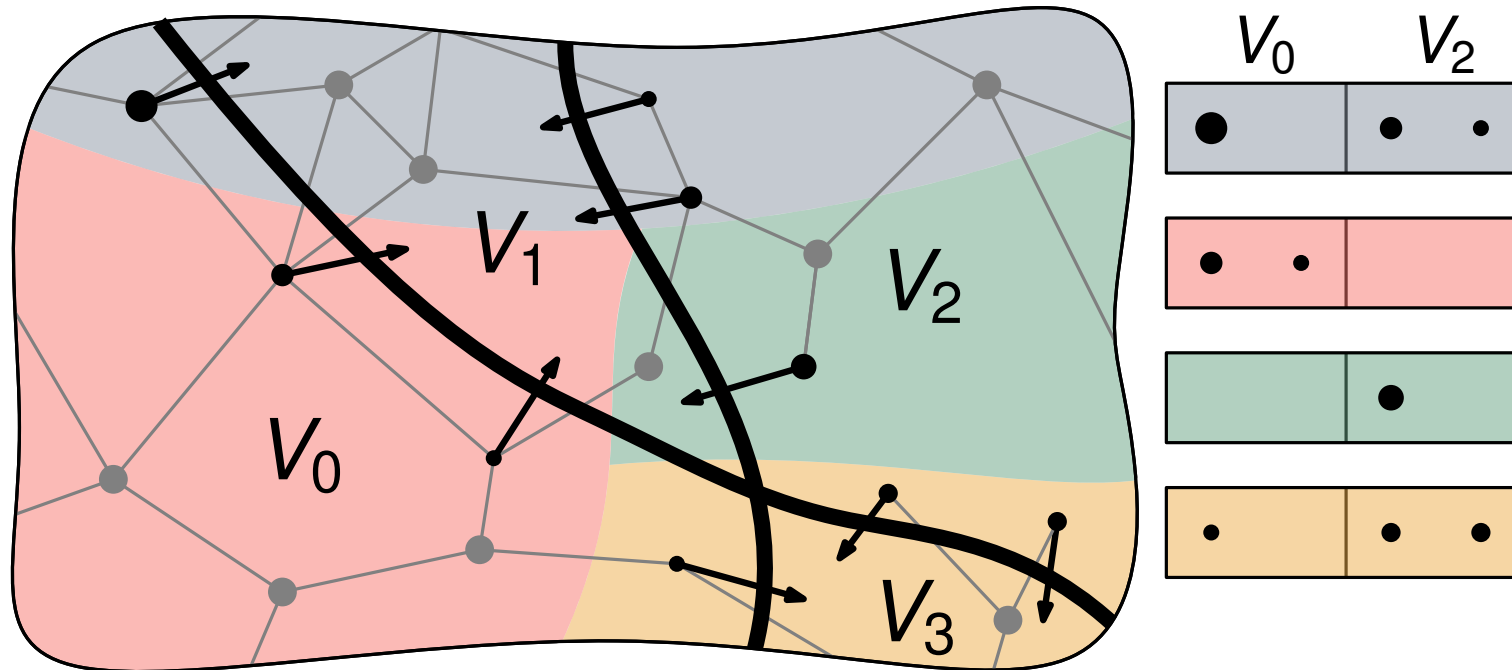
**problem:** uncoarsening

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Refinement

- Refinement: label propagation + balancing



PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning:  Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Refinement

- Refinement: label propagation + balancing

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

# dKaMinPar – Refinement

- Refinement: label propagation + balancing



PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023
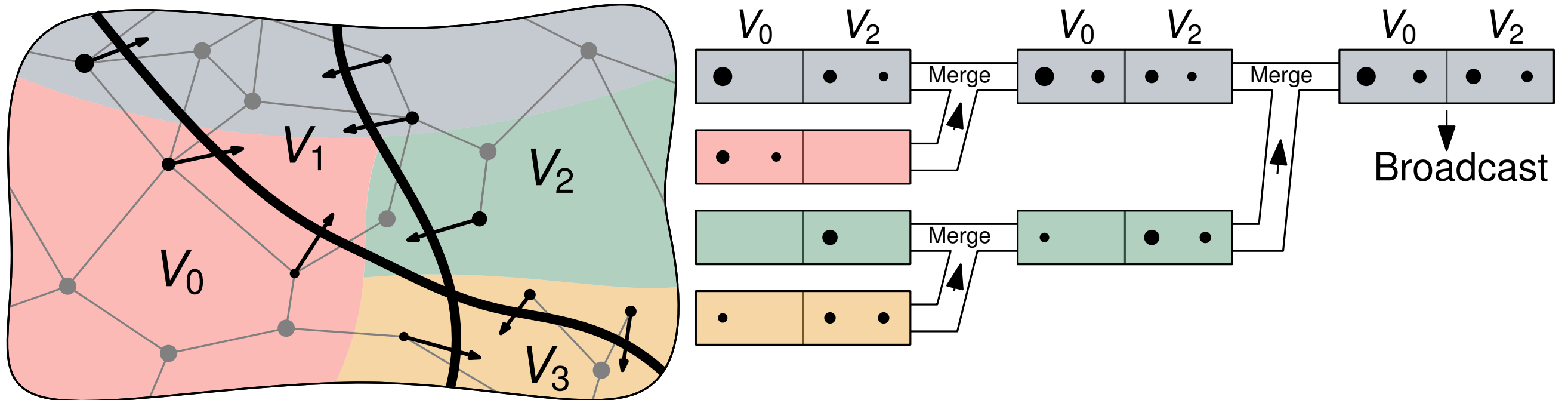
Institute of Theoretical Informatics, Algorithmics II

# dKaMinPar – Refinement

- Refinement: label propagation + balancing



PEs

Sanders, Seemaier – Distributed Deep Multilevel Graph Partitioning: Euro-Par 2023

Institute of Theoretical Informatics, Algorithmics II