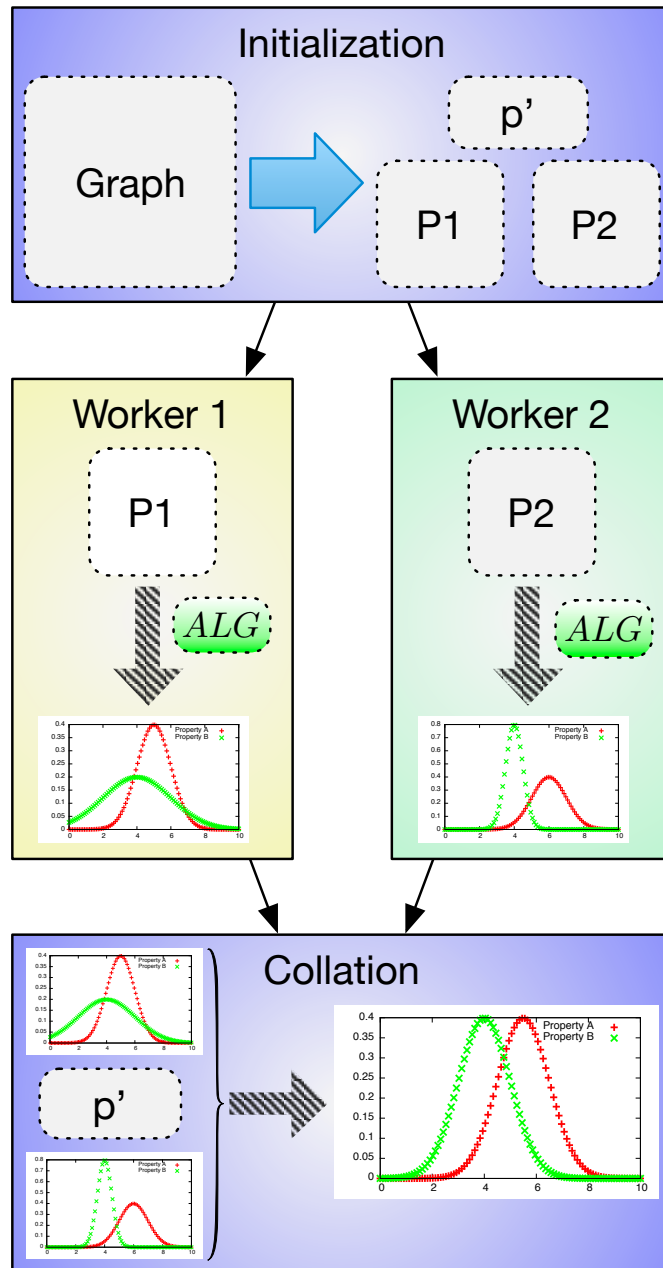


Parallel Stream-based Graph Analysis

Benjamin Schiller, Datenschutz und Datensicherheit, TU Dresden



partitioning as initialization

batches

how to compute is up to each thread

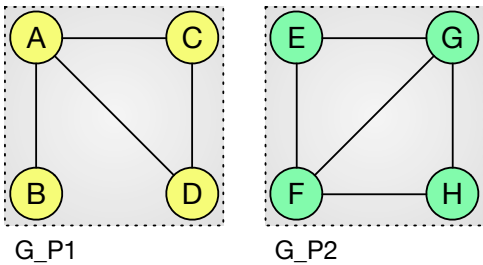
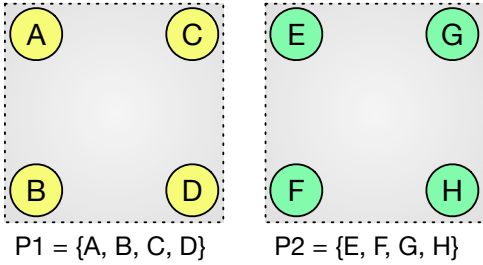
after processing batch

collation could be parallelized

Partition Type

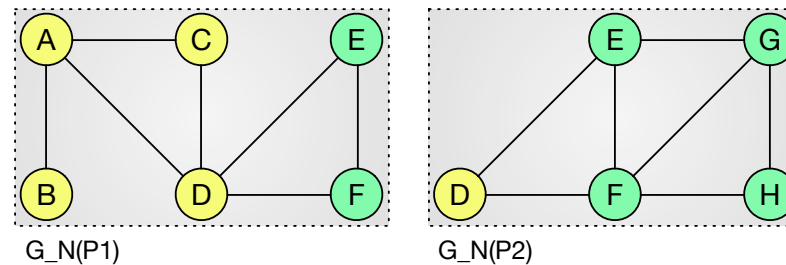
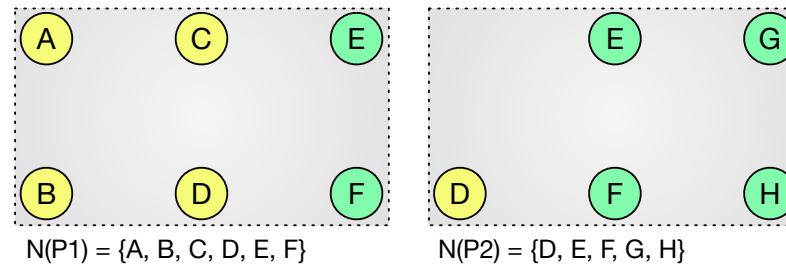


Non Overlapping



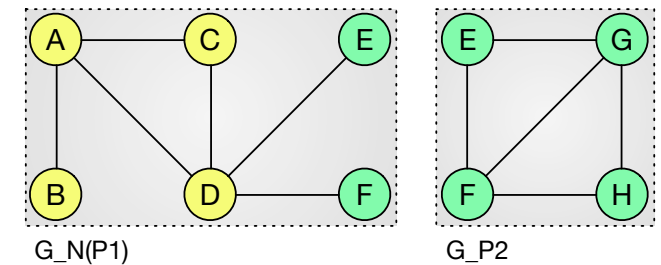
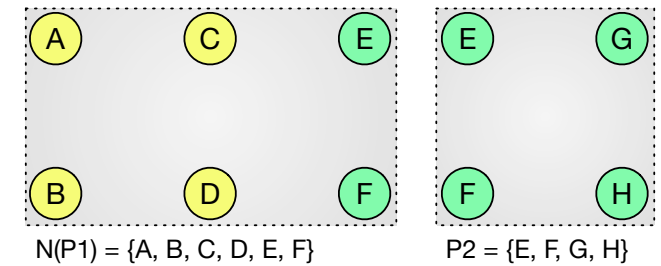
- + no duplication of data
- collation complexity

Overlapping



- + collation complexity
- memory overhead
- computation overhead

Node Cut



- + collation complexity
- no memory overhead
- computation overhead

- Snapshot-based
- Batch-based
- Stream-based

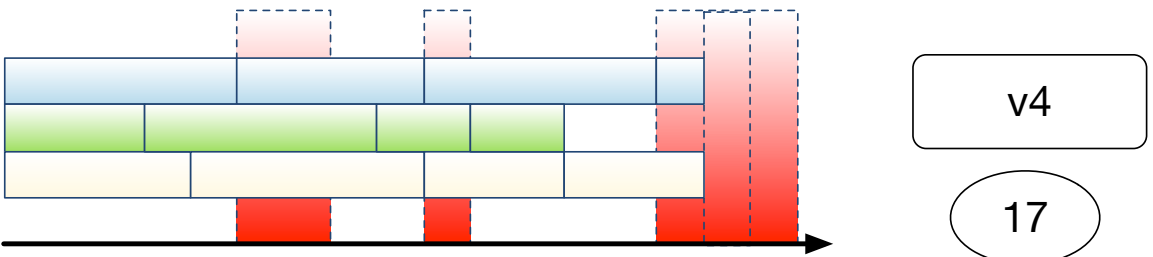
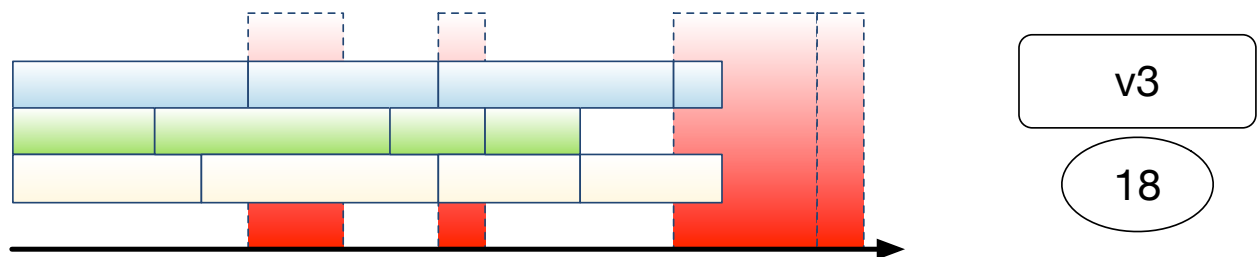
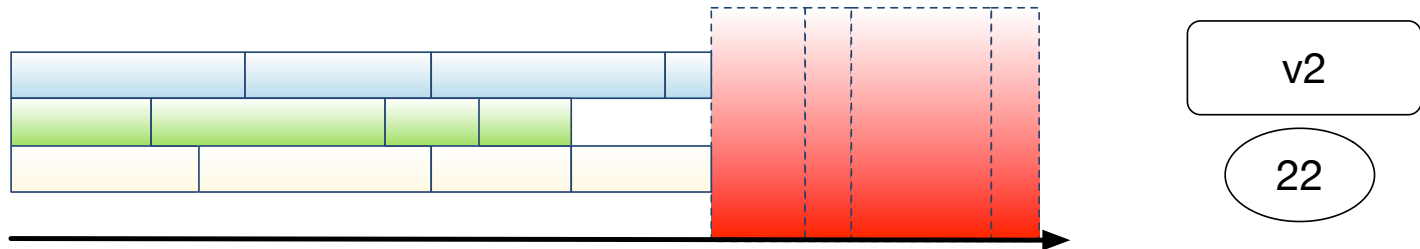
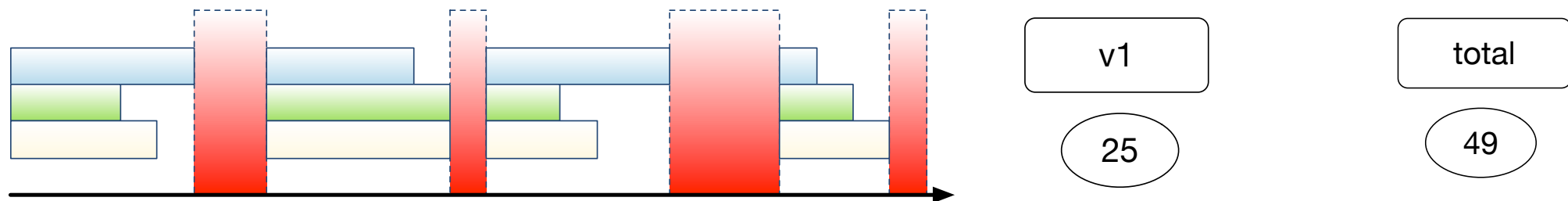
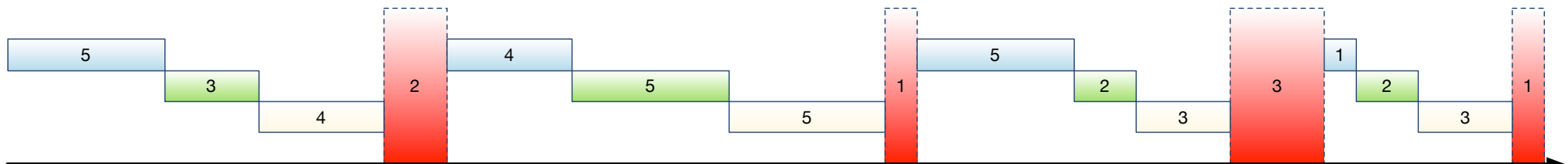
- Snapshot-based (overlapping)
- Batch-based (overlapping)
- Stream-based (overlapping)

- BFS
- DFS
- EQUAL_SIZE
- LPA
- RANDOM

- FIRST_EDGE
- RANDOM
- ROUND_ROBIN

- Statistics (all per thread / per batch)
 - Analysis runtimes
 - $|V|$, $|E|$
 - Collation runtimes
 - Total runtimes ($v\{0..4\}$)
- Statistics (depending on partitioning type)
 - NON_OVERLAPPING: external edges
 - OVERLAPPING: auxiliary nodes, auxiliary edges
 - NODE_CUT: cut nodes

Sequential Evaluation / Approximation



First Results

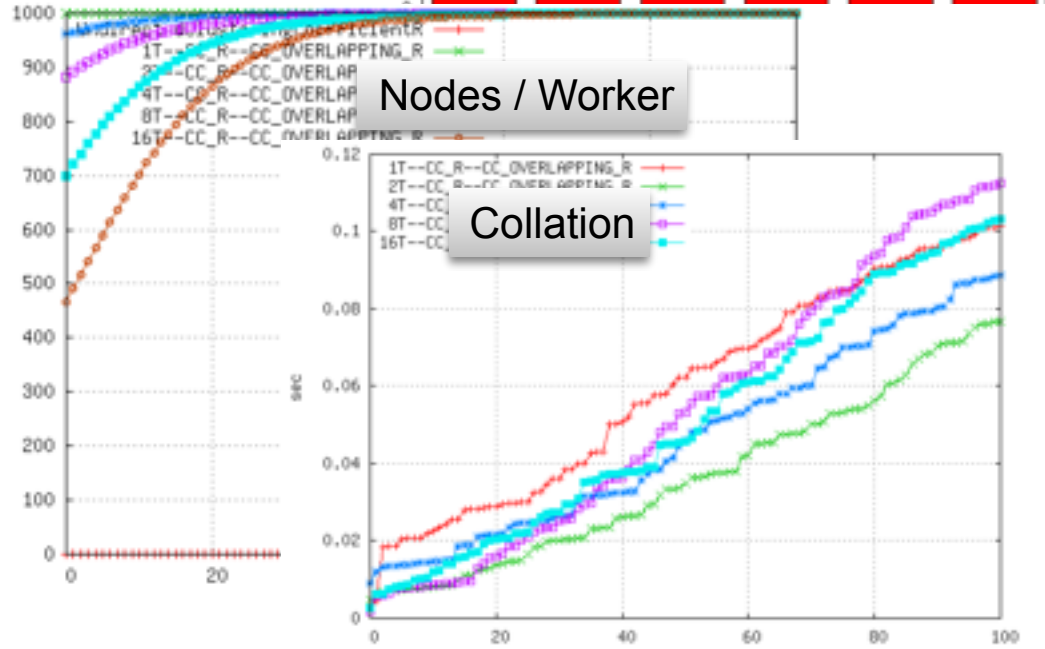
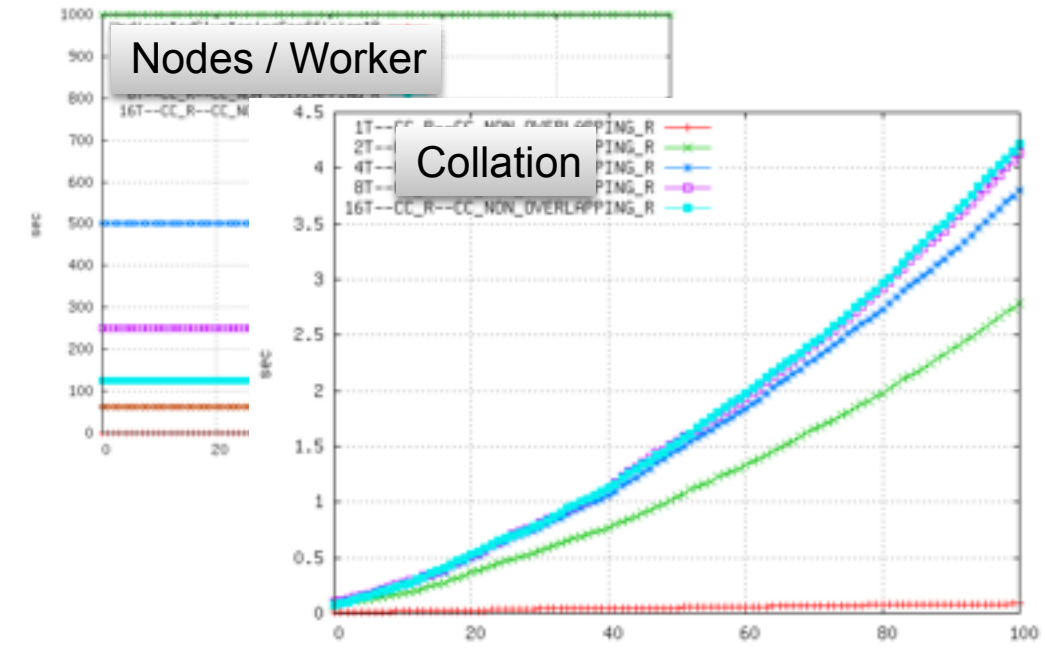
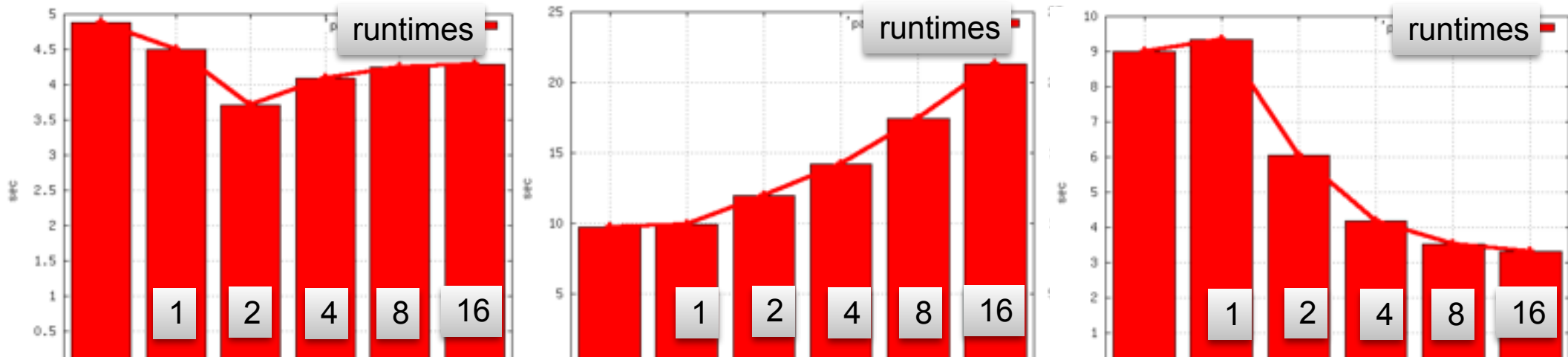


Example: Densifying Graph, Clustering Coefficient

Non-Overlapping

Overlapping

Overlapping (Alg*)



- Approach
 - Partitioning Type?
 - Partitioning Scheme?
 - Node Assignment?
 - Re-Partitioning?
- Related Work
 - Is there a niche for this approach?
- Evaluation
 - How to “simulate”?
 - Which metrics to use?
 - Which datasets to use?