

# SCALING SPATIAL OVERLAY OPERATIONS AND FLOCK PATTERN DISCOVERY

Andres Calderon · [acald013@ucr.edu](mailto:acald013@ucr.edu)

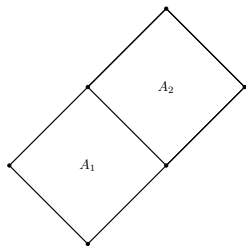
University of California, Riverside

November 6, 2024

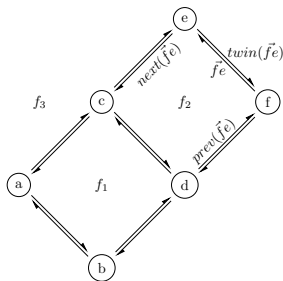
- 1 SCALABLE OVERLAY OPERATIONS OVER DCEL POLYGONS LAYERS
- 2 SCALING DCEL OVERLAY OPERATIONS TO SUPPORT DANGLE AND CUT EDGES
- 3 SCALABLE PROCESSING OF MOVING FLOCK PATTERNS

# WHAT IS A DCEL?

- **Doubly Connected Edge List** - DCEL.
- A spatial data structure to **represent planar subdivisions** of surfaces.
- Represent topological and geometric information as vertices, edges, and faces.
- **Applications:** polygon overlays, polygon triangulation and their applications in surveillance, robot motion planing, circuit board printing, etc.



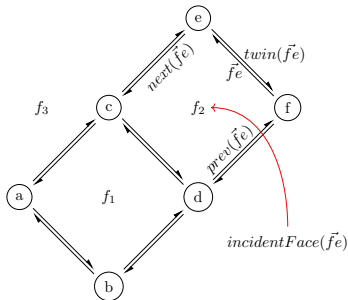
Planar subdivision



DCEL representation

# DCEL DESCRIPTION

- DCEL uses three tables: Vertices, Faces and Half-edges.



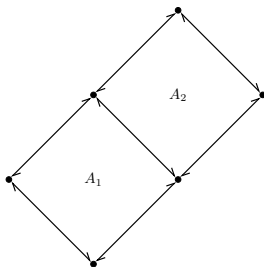
vertex	coordinates	incident edge
a	(0,2)	$\vec{ba}$
b	(2,0)	$\vec{db}$
c	(2,4)	$\vec{dc}$
$\vdots$	$\vdots$	$\vdots$

face	boundary edge	hole list
$f_1$	$\vec{ab}$	<i>nil</i>
$f_2$	$\vec{f_e}$	<i>nil</i>
$f_3$	<i>nil</i>	<i>nil</i>

half-edge	origin	face	twin	next	prev
$\vec{f_e}$	f	$f_2$	$\vec{ef}$	$\vec{ec}$	$\vec{df}$
$\vec{ca}$	c	$f_1$	$\vec{ac}$	$\vec{ab}$	$\vec{dc}$
$\vec{db}$	d	$f_3$	$\vec{bd}$	$\vec{ba}$	$\vec{fd}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

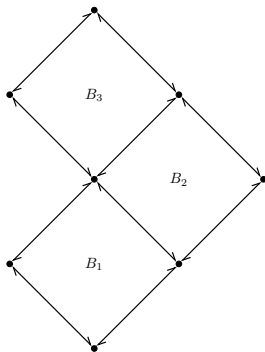
# DCEL ADVANTAGES

- **Efficiency:** very efficient for computation of *overlay operators*.
- **Re-usability:** allows multiple operations over the same DCEL.
- **Pipelining:** the output of a DCEL operator can be input to another DCEL operator.



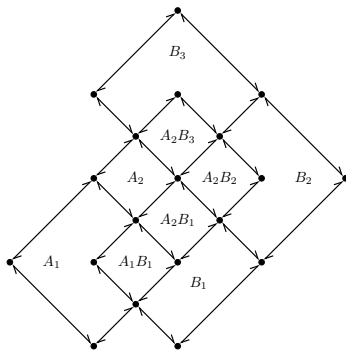
# DCEL ADVANTAGES

- **Efficiency:** very efficient for computation of *overlay operators*.
- **Re-usability:** allows multiple operations over the same DCEL.
- **Pipelining:** the output of a DCEL operator can be input to another DCEL operator.



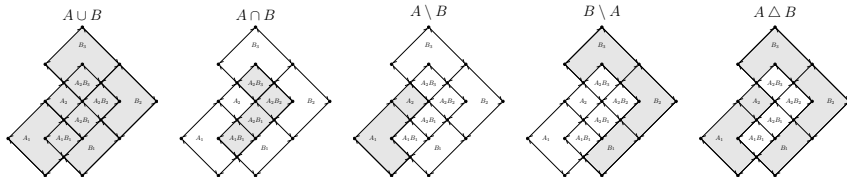
# DCEL ADVANTAGES

- **Efficiency:** very efficient for computation of *overlay operators*.
- **Re-usability:** allows multiple operations over the same DCEL.
- **Pipelining:** the output of a DCEL operator can be input to another DCEL operator.



# DCEL ADVANTAGES

- **Efficiency:** very efficient for computation of *overlay operators*.
- **Re-usability:** allows multiple operations over the same DCEL.
- **Pipelining:** the output of a DCEL operator can be input to another DCEL operator.

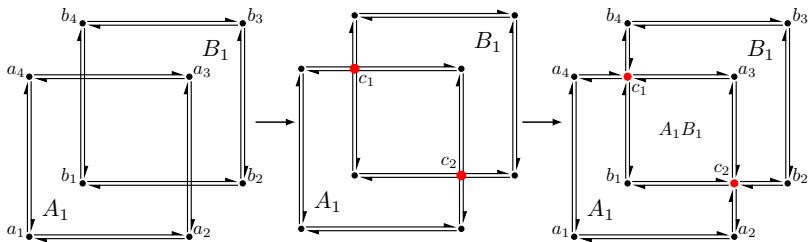




- Currently only sequential DCEL implementations exist.
- Unable to deal with large datasets (i.e. US Census tracks at national level).
- We propose a *scalable distributed* approach to compute the overlay of two polygon layers using DCELS.
- Distribution enables scalability, but introduces challenges: the **orphan-cell** problem and the **orphan-hole** problem.
- Optimization for reducing and merging results and unbalanced layer sizes.

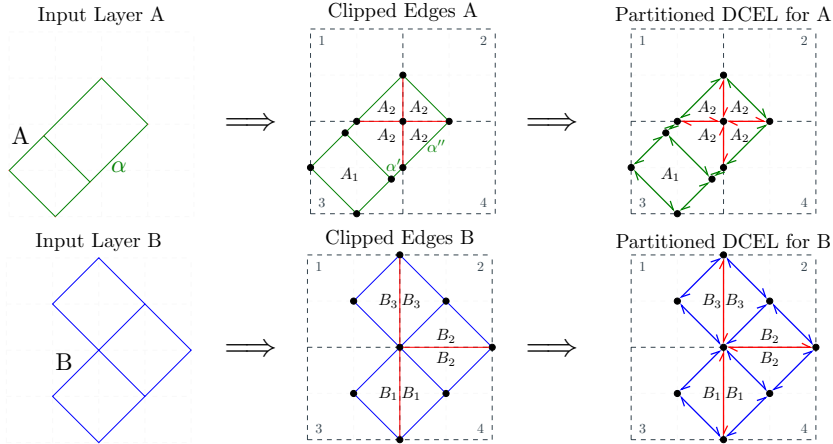
# SEQUENTIAL IMPLEMENTATION

- Consider two (simple) input DCELs  $A_1$  and  $B_1$ . The sequential algorithm first finds the intersections of half-edges.
- Then, new vertices (e.g.  $c_1, c_2$ ) are created, half-edges are updated, new faces are added and labeled (e.g.  $A_1B_1$ ).

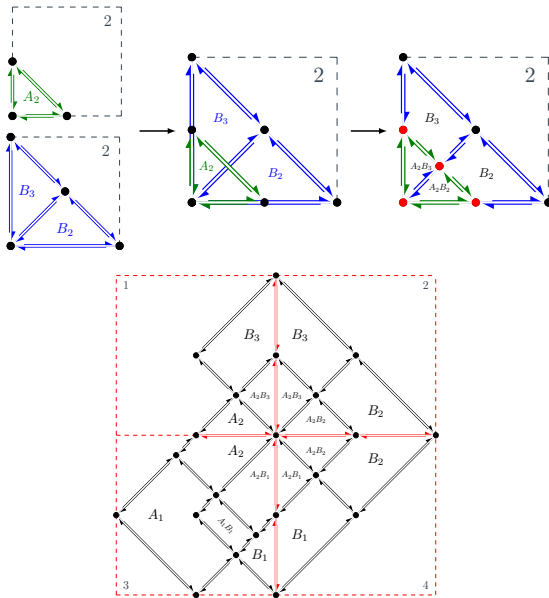


- Two phases: (1) Distributed DCEL construction, and (2) Distributed overlay evaluation.
- Distribution is based on a spatial index (e.g. quadtree)
- Each input DCEL layer (e.g. A, B) is partitioned using the same index
- Each index cell should contain all information needed so that it can compute the overlay DCEL locally
- For each cell to be independent, we need to create "artificial" edges and vertices

# SCALABLE IMPLEMENTATION



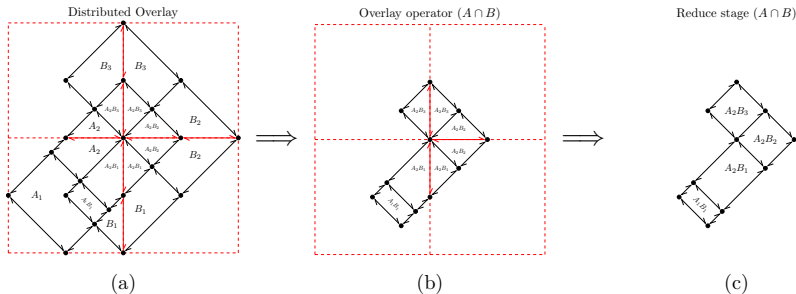
# DISTRIBUTED DCEL CONSTRUCTION



# OVERLAY EVALUATION

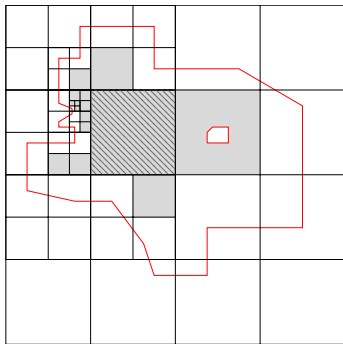
## ■ Answering global overlay queries...

- ▶ To compute a particular overlay operator, we query local DCELS.
- ▶ This work is done independently at each cell (node).
- ▶ SDCEL then collects back all local DCEL answers and computes the final answer (by removing artificial edges and concatenating the resulting faces).



# LABELING ORPHAN CELLS AND ORPHAN HOLES

- We next discuss the **orphan cell** problem (orphan holes are handled similarly).
- A large face (e.g. the red polygon in the figure) can contain cells that do not intersect with any of the face's boundary edges (called *regular edges*).
- Such cells do not contain any label and thus we do not know which face they belong to.

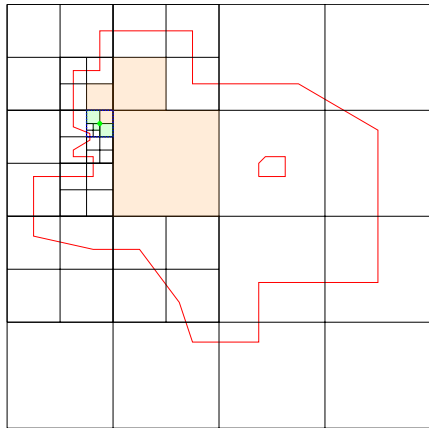






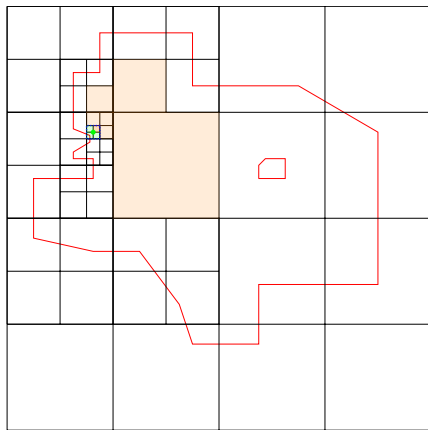
# LABELING ORPHAN CELLS AND ORPHAN HOLES

- We provide an algorithm to efficiently solve the orphan cell problem.



# LABELING ORPHAN CELLS AND ORPHAN HOLES

- We provide an algorithm to efficiently solve the orphan cell problem.



- Optimizing for faces overlapping many cells...
  - ▶ Naive approach sends all faces that overlap a cell to a master node (that will combine them).
  - ▶ We propose an intermediate reduce processing step.
    - The user provides a level in the quadtree structure and faces are evaluated at those intermediate reducers.
  - ▶ We also consider another approach that re-partitions such faces using their labels as the key.
    - It avoids the reduce phase but implies an additional shuffle.
    - However, as we show in the experiments this overhead is minimal.

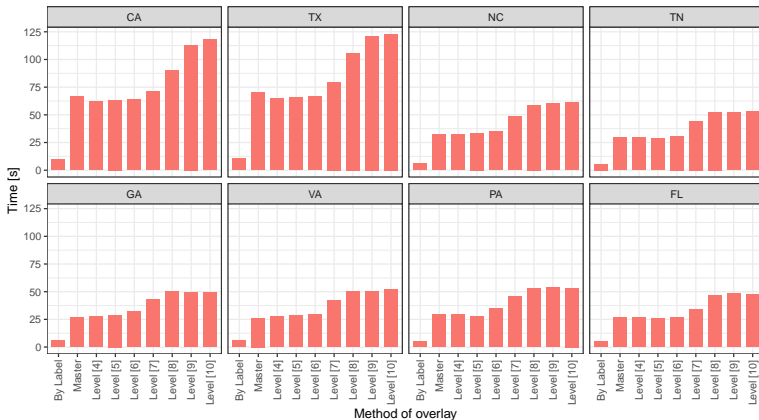
- Optimizing for unbalanced layers...
  - ▶ Finding intersections is the most critical part of the overlay computation.
  - ▶ However, in many cases one of the layers has much more half-edges than the other.
  - ▶ Sweep-line algorithms to detect intersections run over all the edges.
  - ▶ Instead we scan the larger dataset only for the x-intervals where there are half-edges from the smaller dataset.

## ■ Datasets.

Dataset	Layer	Number of polygons	Number of edges
MainUS	Polygons for 2000	64983	35417146
	Polygons for 2010	72521	36764043
GADM	Polygons for Level 2	160241	64598411
	Polygons for Level 3	223490	68779746
CCT	Polygons for 2000	7028	2711639
	Polygons for 2010	8047	2917450

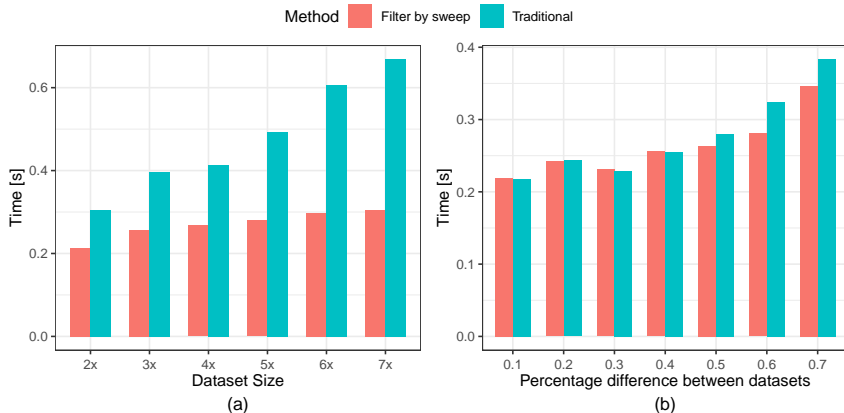
# EXPERIMENTAL EVALUATION

## ■ Evaluation of the overlapping faces optimization.



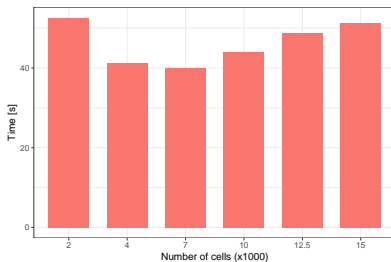
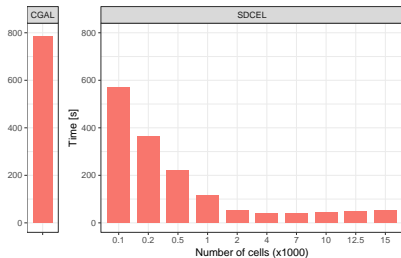
# EXPERIMENTAL EVALUATION

## ■ Evaluation of the unbalanced layers optimization.



# EXPERIMENTAL EVALUATION

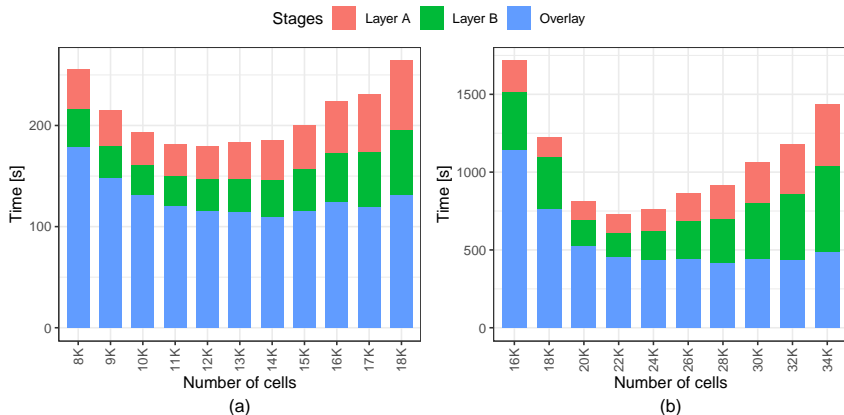
- Performance varying number of partition cells (CCT dataset).





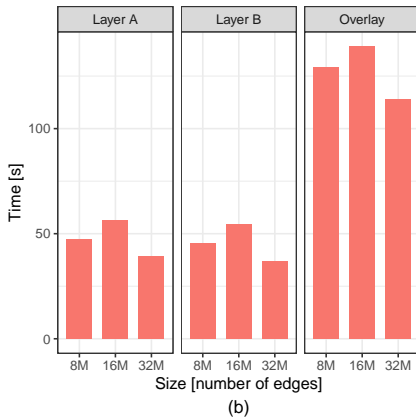
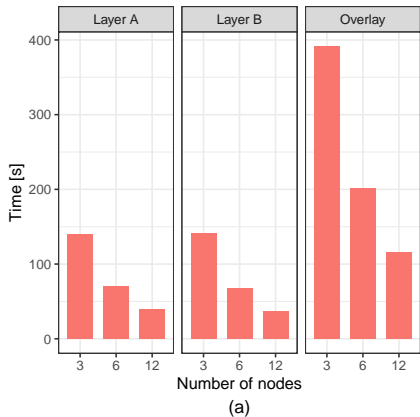
# EXPERIMENTAL EVALUATION

## ■ Performance with MainUS and GADM datasets.



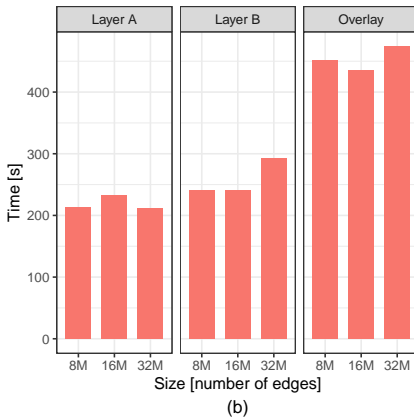
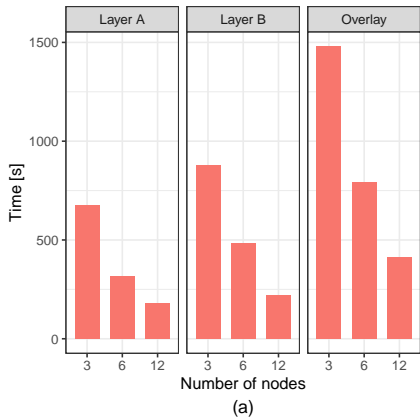
# EXPERIMENTAL EVALUATION

## ■ Scale-up and Speed-up.



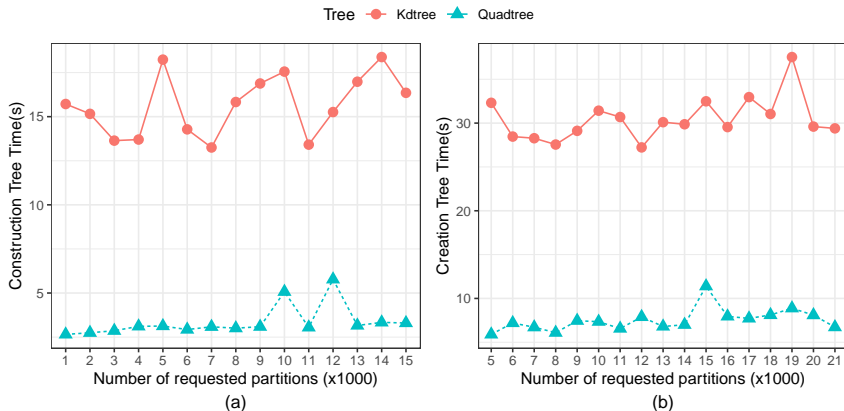
# EXPERIMENTAL EVALUATION

## ■ Scale-up and Speed-up.



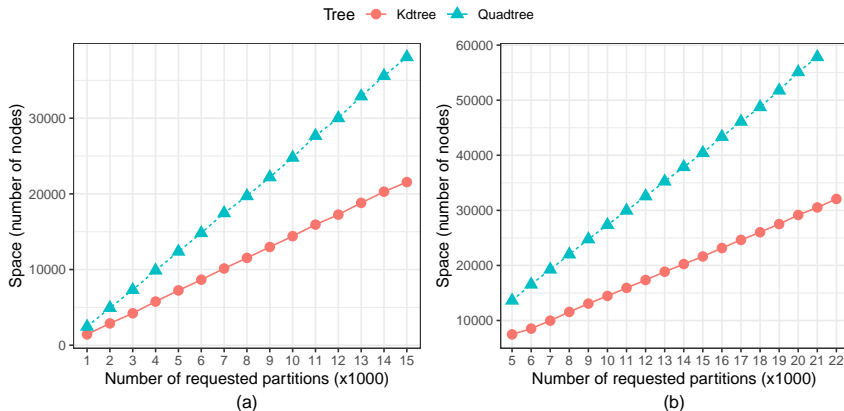
# EXPERIMENTAL EVALUATION

## ■ Space-oriented vs Data-oriented partitioners.



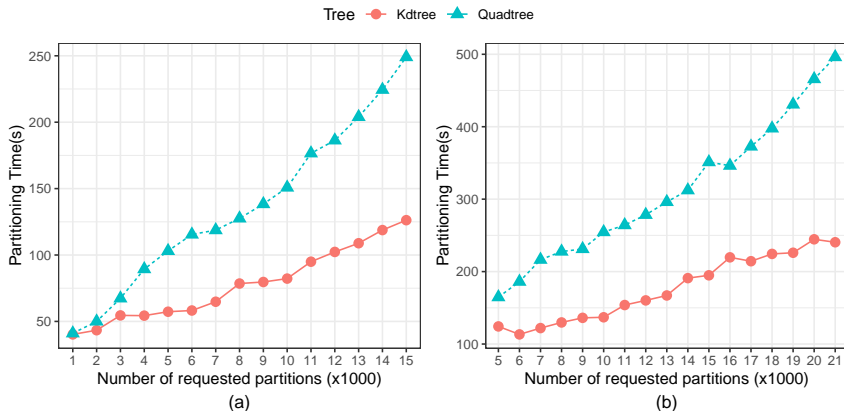
# EXPERIMENTAL EVALUATION

## ■ Space-oriented vs Data-oriented partitioners.



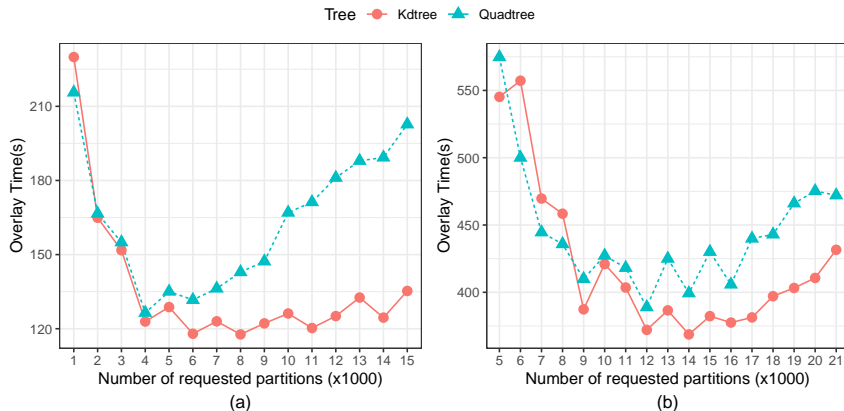
# EXPERIMENTAL EVALUATION

## ■ Space-oriented vs Data-oriented partitioners.



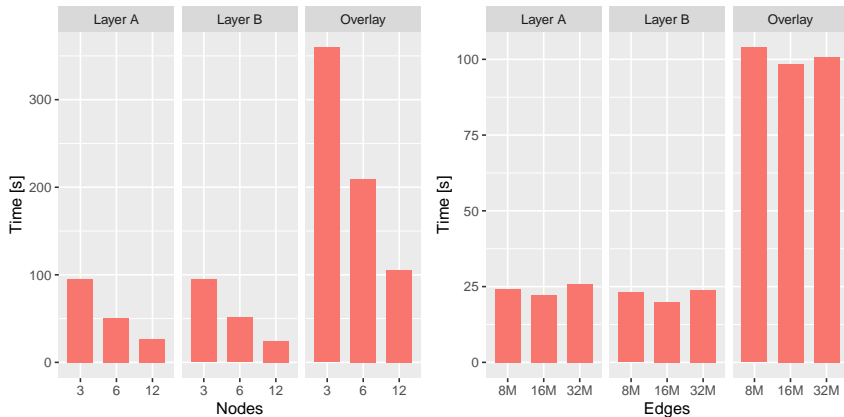
# EXPERIMENTAL EVALUATION

## ■ Space-oriented vs Data-oriented partitioners.



# EXPERIMENTAL EVALUATION

## ■ Space-oriented vs Data-oriented partitioners.





- 1 SCALABLE OVERLAY OPERATIONS OVER DCEL POLYGONS LAYERS
- 2 SCALING DCEL OVERLAY OPERATIONS TO SUPPORT DANGLE AND CUT EDGES
- 3 SCALABLE PROCESSING OF MOVING FLOCK PATTERNS

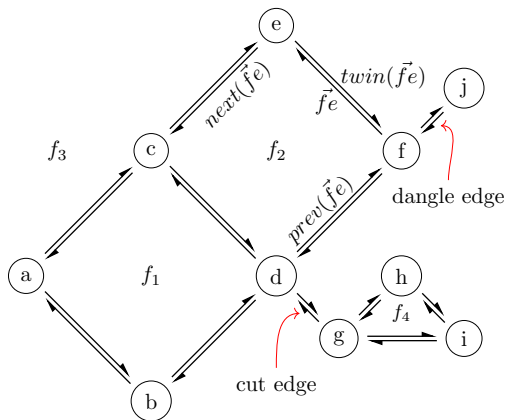
- We extend the overlay DCEL approach to accept scattered and noisy line segments (including dangles and cut edges) as input, rather than being restricted to clean polygon data.
- This chapter extends the previous work in [Calderon et al, 2023]<sup>1</sup> by building on the scalable polygonization methods presented in [Abdelhafeez et al, 2023]<sup>2</sup>.

---

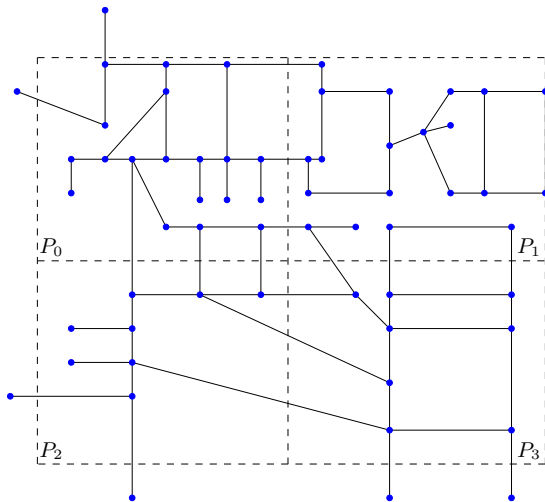
<sup>1</sup>Scalable Overlay Operations over DCEL Polygon Layers.(SSTD'23).

<sup>2</sup>DDCEL: Efficient Distributed Doubly Connected Edge List for Large Spatial Networks. (MDM'23).

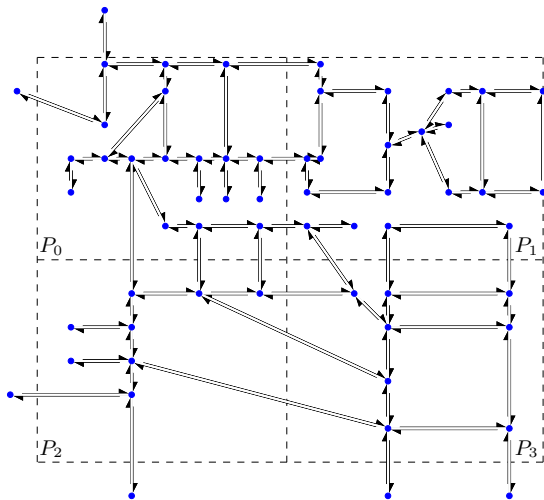
# DANGLE AND CUT EDGES



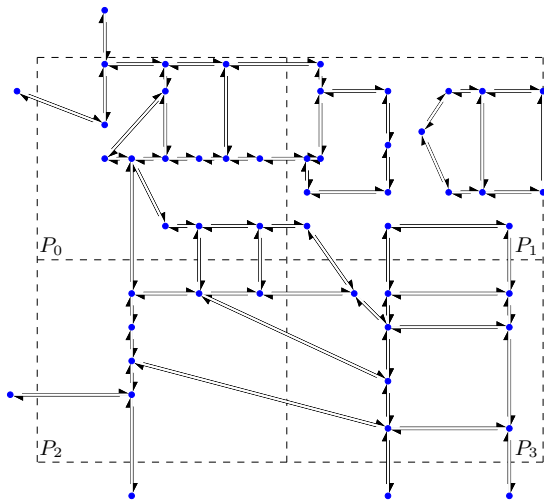
# DANGLE AND CUT EDGES



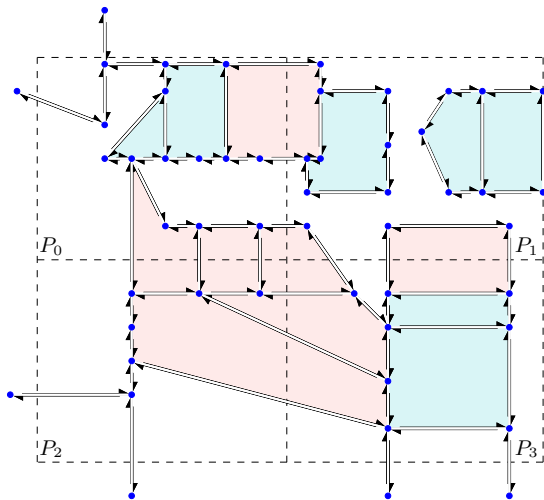
# DANGLE AND CUT EDGES



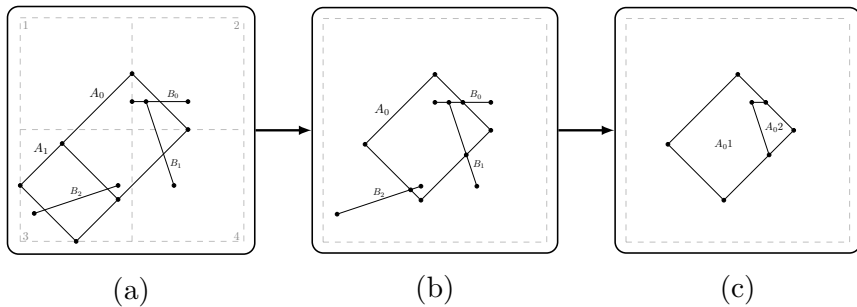
# DANGLE AND CUT EDGES



# DANGLE AND CUT EDGES



# DANGLE AND CUT EDGES

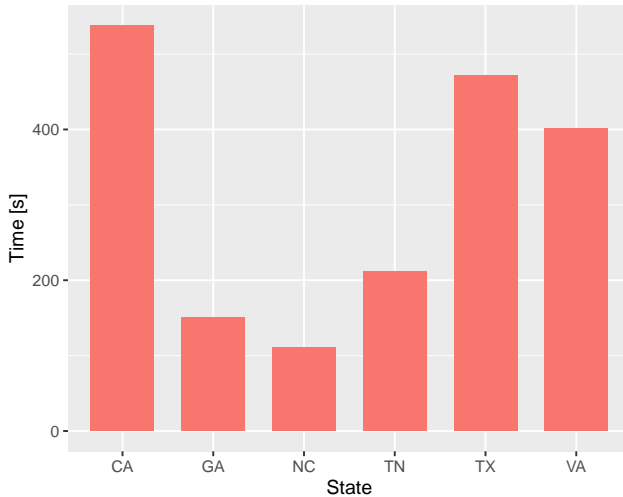




# DANGLE AND CUT EDGES

Dataset	Number of Polygons	Number of Edges	Result Polygons
	Layer <i>A</i>	Layer <i>B</i>	
TN	1,272	3,380,780	41,761
GA	1,633	4,647,171	49,125
NC	1,272	7,212,604	22,413
TX	4,399	8,682,950	98,635
VA	1,554	8,977,361	38,941
CA	7,038	9,103,610	96,916

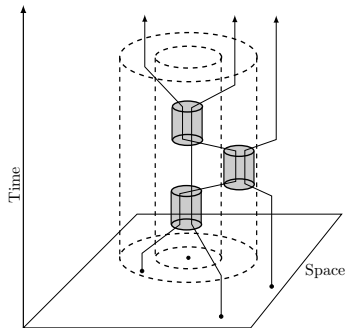
# DANGLE AND CUT EDGES



- 1 SCALABLE OVERLAY OPERATIONS OVER DCEL POLYGONS LAYERS
- 2 SCALING DCEL OVERLAY OPERATIONS TO SUPPORT DANGLE AND CUT EDGES
- 3 SCALABLE PROCESSING OF MOVING FLOCK PATTERNS

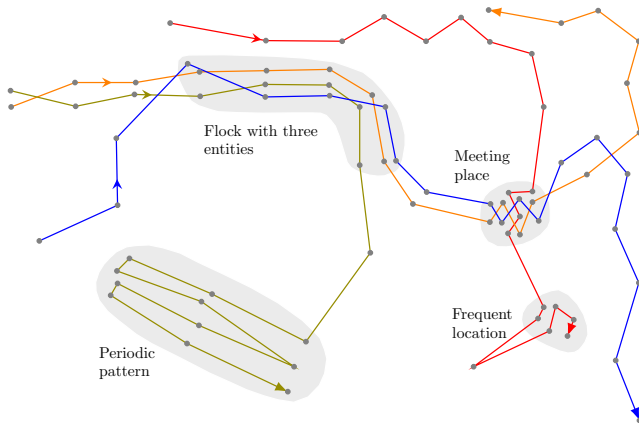
# LARGE TRAJECTORY DATABASES

- A spatial trajectory is a trace in time generated by a moving entity in a geographical space.
- i.e.  $p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_n$
- A trajectory is stored as a time-ordered sequence of points,  $p_i = (x, y, t)$  (spatial coordinate + time instant).



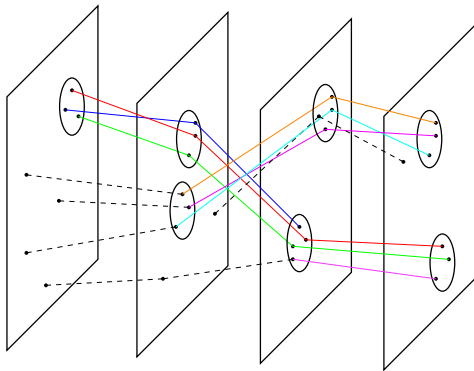
(Shoval, 2017)

# MOVEMENT PATTERNS



(Gudmundsson, et al. 2008)

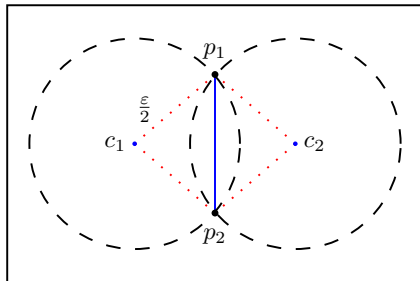
- i.e. convoys, moving clusters, swarms, gatherings, **flocks**, ...



- $\varepsilon$ : Diameter of the circle which contains all the objects.
- $\mu$ : Minimum number of objects.
- $\delta$ : Minimum time interval the objects travel 'together'.

# BASIC FLOCK EVALUATION ALGORITHM

- The first polynomial-time solution for determining disk locations (Vieira, et al. 2009).
- The algorithm generates a finite set of disk locations based on a given pair of points.
- Under fixed time duration it has polynomial time complexity  $O(\delta|\tau|^{(2\delta)+1})$ .



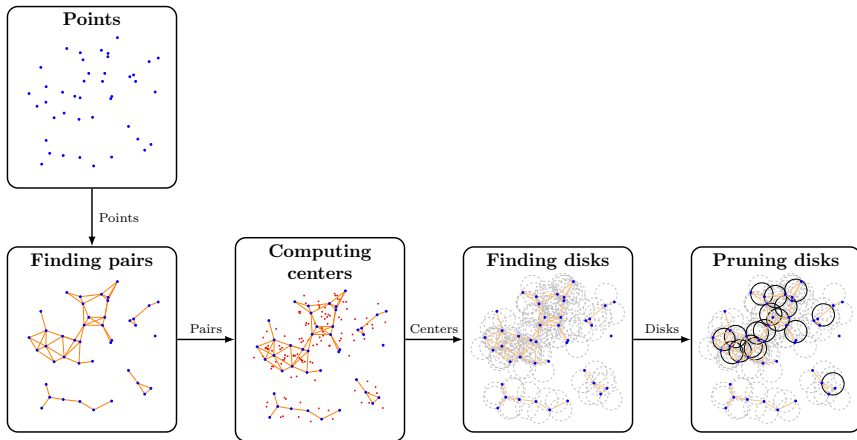
# BASIC FLOCK EVALUATION ALGORITHM

- Two main parts:
  - ▶ In the spatial domain it finds maximal disks at each time instant.
  - ▶ In the temporal domain it joins consecutive times to match set of maximal disks.



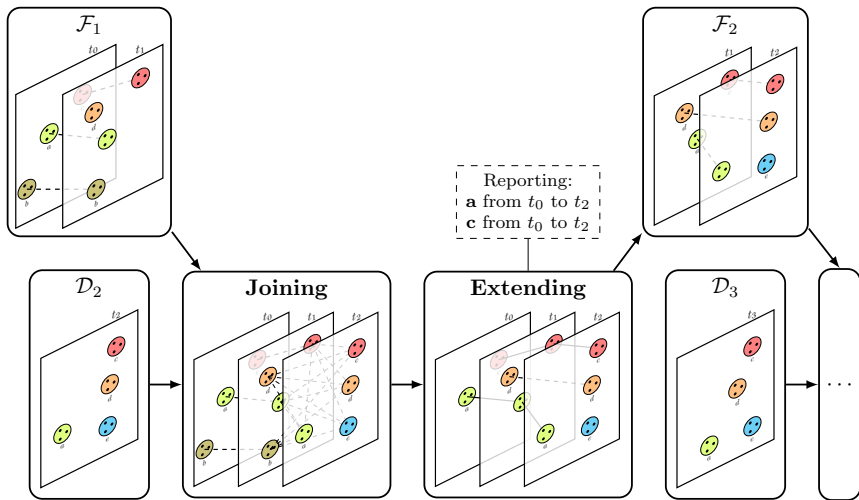
# ON THE SPATIAL DOMAIN

## ■ BFE overview...

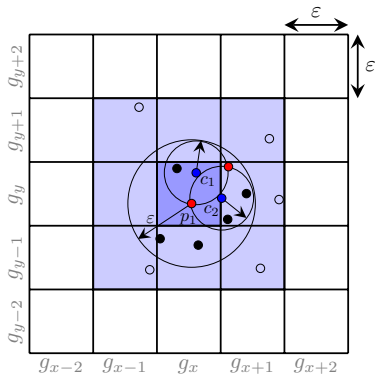


# ON THE TEMPORAL DOMAIN

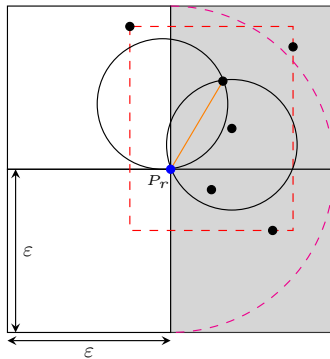
## ■ BFE overview...



# PSI ALGORITHM



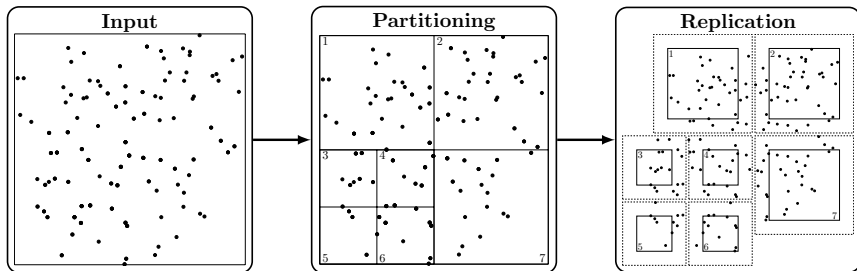
(Vieira, et al. 2009)



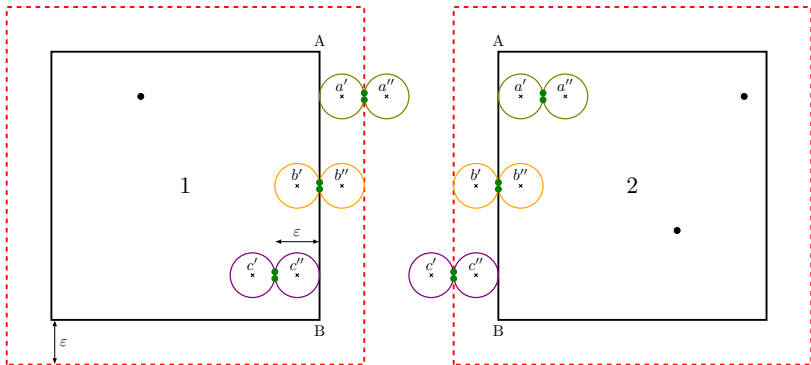
(Tanaka, et al. 2016)

- High complexity limits scalability.
- Large datasets with dense clusters of moving entities per time instant significantly impact performance.
- Specifically,
  - ▶ identifying maximal disks is hindered by the extensive number of candidates requiring pruning.
  - ▶ when parallelizing, we must address moving flocks that traverse contiguous partitions.
- We propose a parallel and scalable solution for both spatial and temporal domains.

## ■ Partitioning strategy...

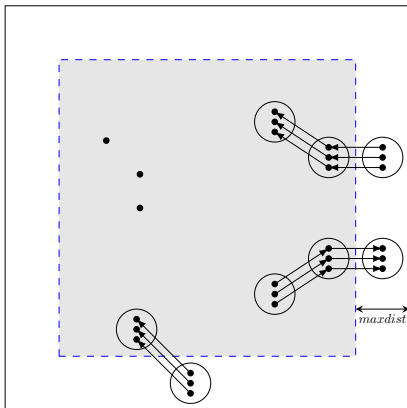


## ■ Handling duplication...

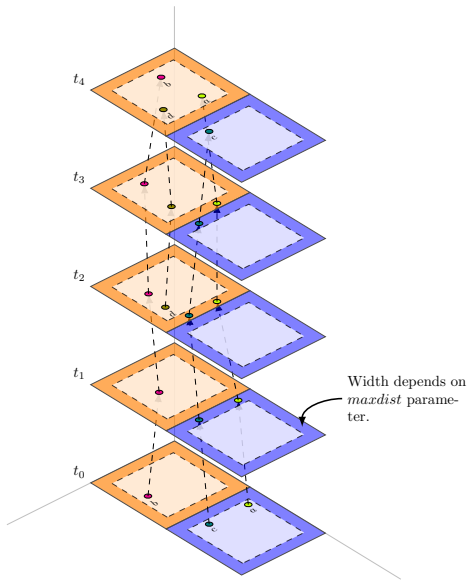


# ON THE TEMPORAL DOMAIN

- We introduce the *maxdist* parameter to define an area where we have to track **crossing partial flocks** (CPFs)...



# ON THE TEMPORAL DOMAIN

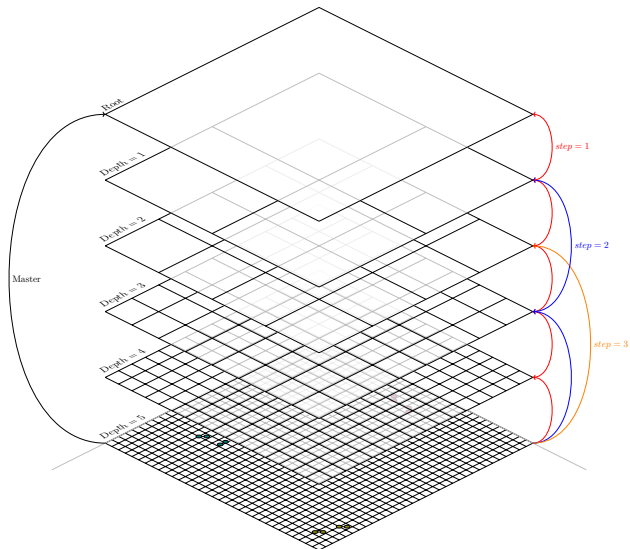


\* $a, b, c$  and  $d$  are flocks moving along time.

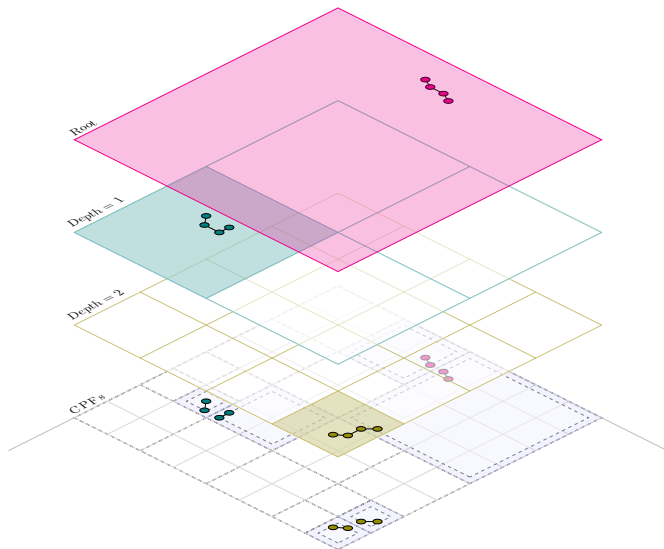


- Discovered flocks inside the safe area are ready to be reported.
- CPFs require post-processing. We propose four alternatives:
  - ▶ Master
  - ▶ By-Level
  - ▶ Least Common Ancestor (LCA)
  - ▶ Cube-based

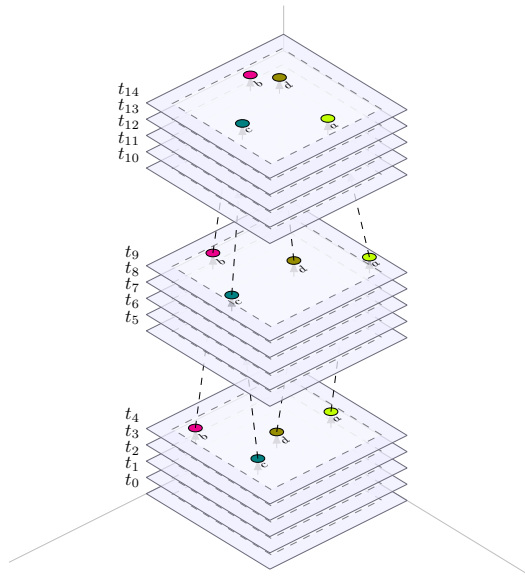
# ON THE TEMPORAL DOMAIN



# ON THE TEMPORAL DOMAIN



# ON THE TEMPORAL DOMAIN



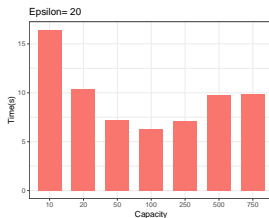
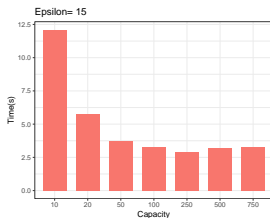
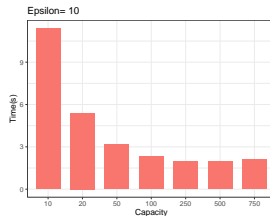
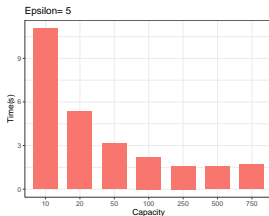
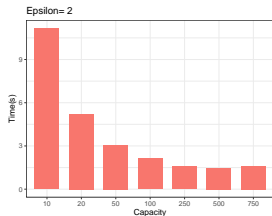
\*a,b,c and d are flocks moving along time.

# DATASETS

Dataset	Number of Trajectories	Total number of points	Maximum Duration (min)
Berlin10K	10000	97526	10
LA25K	25000	1495637	30
LA50K	50000	2993517	60

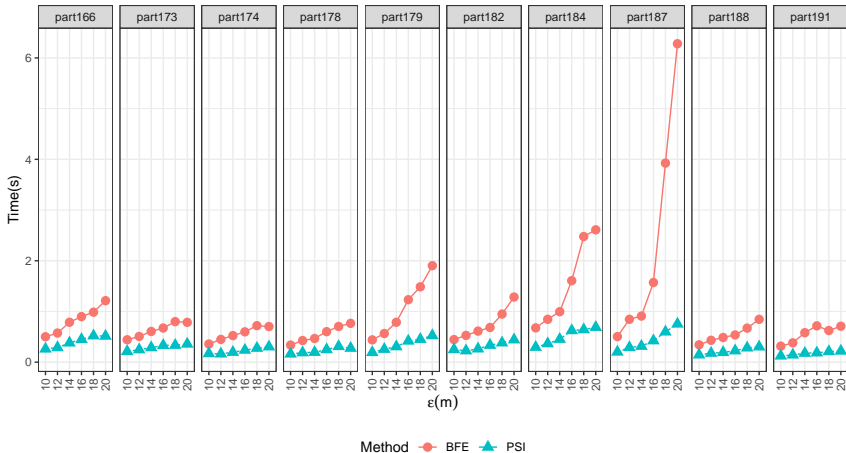
# EXPERIMENTS

## ■ Optimizing the number of partitions for Phase 1.



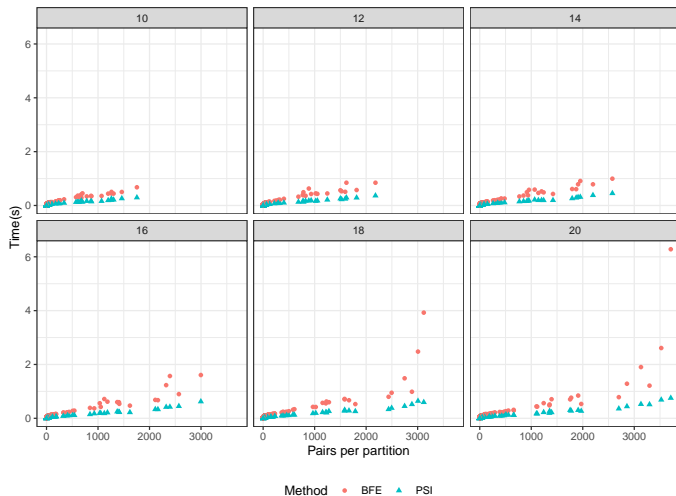
# EXPERIMENTS

- Analyzing most costly partitions.
  - Top 10 most costly partitions.



# EXPERIMENTS

- Analyzing most costly partitions.
  - By Pairs density..

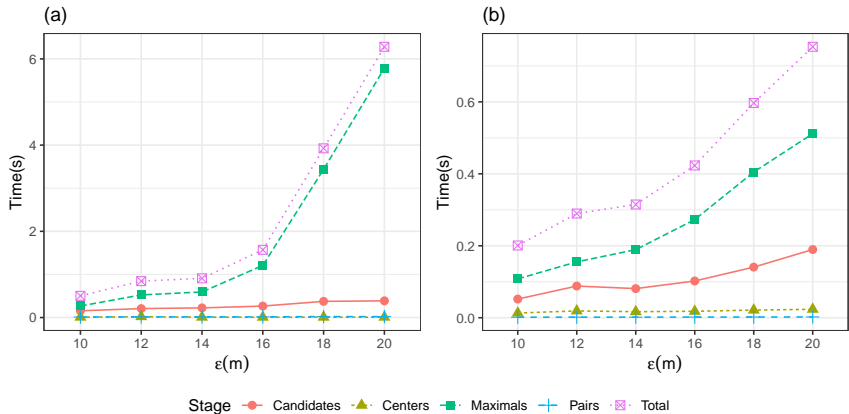




# EXPERIMENTS

## ■ Analyzing most costly partitions.

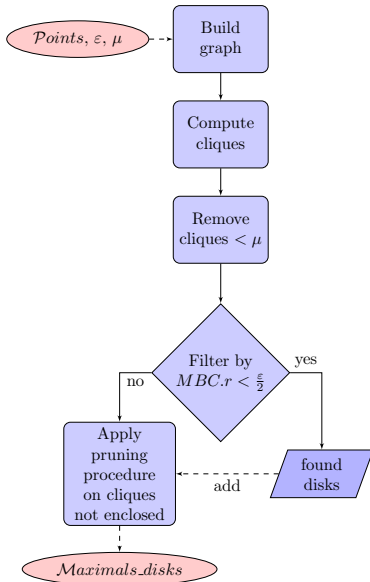
► By Stages in the most costly partition [(a) BFE (b) PSI].



# CAN WE REDUCE PRUNING TIME?

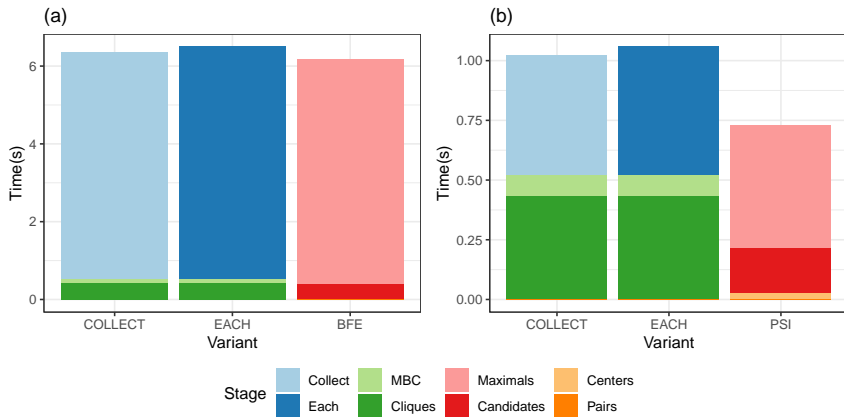
- **Maximal clique (MC):** Given an undirected graph, a MC is a subset of vertices, each directly connected to every other in the subset, that cannot be expanded by adding additional vertices.
- **Minimum Bounding Circle (MBC):** Given a set of points in Euclidean space, the MBC is the smallest circle that can enclose all the points.

# CAN WE REDUCE PRUNING TIME?



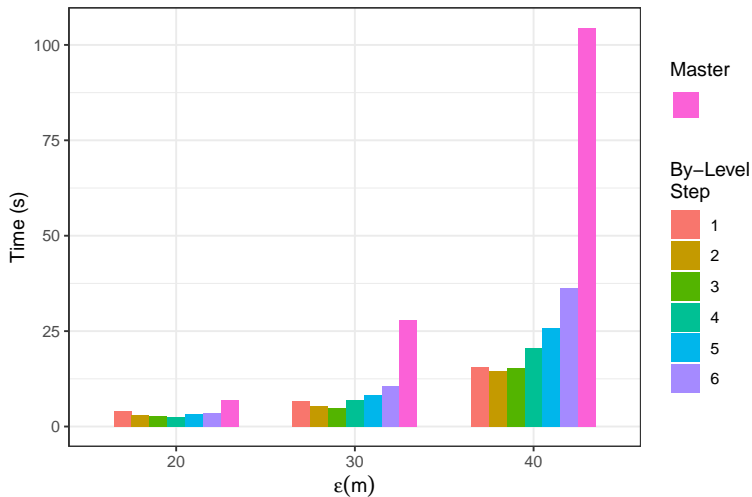
# CAN WE REDUCE PRUNING TIME?

- Phase 1 variants performance [(a) vs BFE (b) vs PSI].



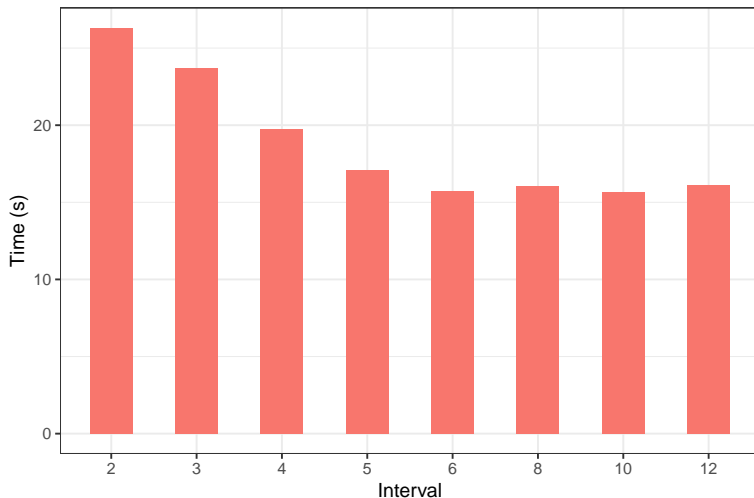
# EXPERIMENTS

- Finding best *step* value for By-Level alternative.

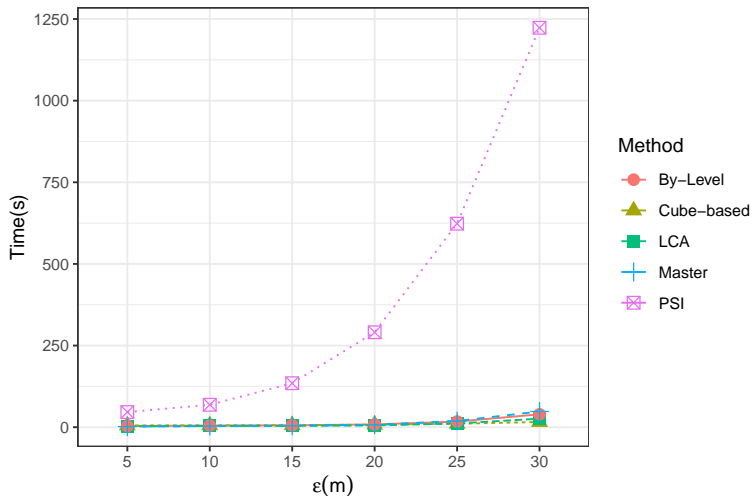


# EXPERIMENTS

- Finding best *interval* value for Cube-based alternative.

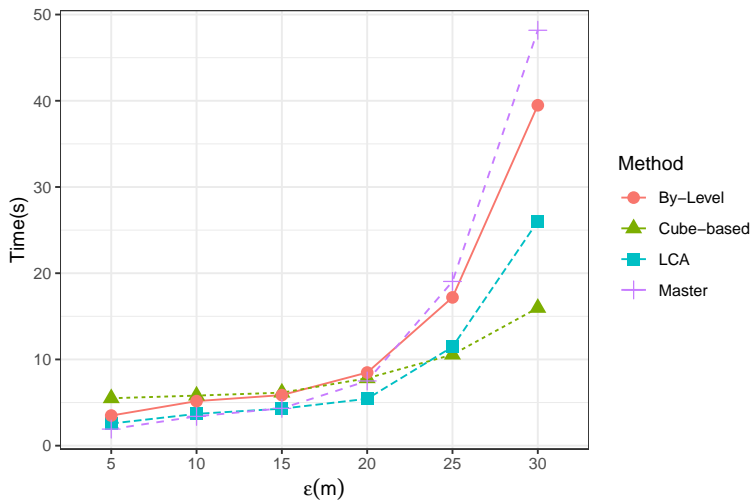


## ■ Scalable alternatives vs standard PSI.



# EXPERIMENTS

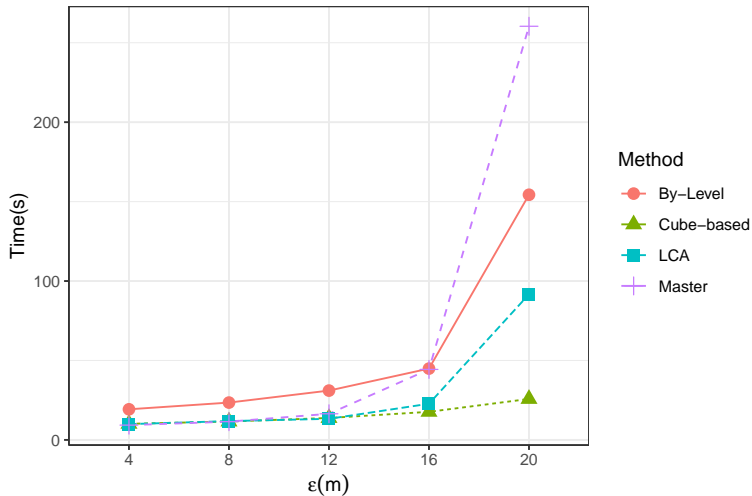
## ■ Scalable alternatives in LA25K dataset.





# EXPERIMENTS

## ■ Scalable alternatives in LA50K dataset.



Thank you!