

SECLEDS: Sequence Clustering in Evolving Data Streams via Multiple Medoids and Medoid Voting

Azqa Nadeem, Sicco Verwer

Delft University of Technology, The Netherlands

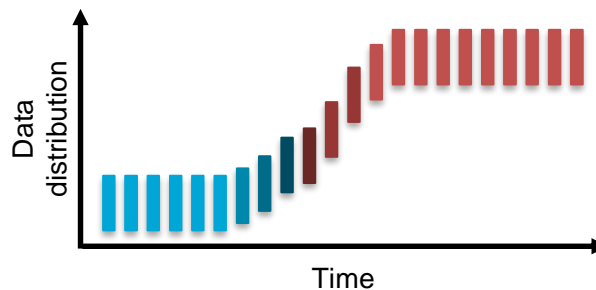
ECML/PKDD 2022

Sequence clustering in streaming environments

- Stream clustering
 - Cluster infinitely many data items

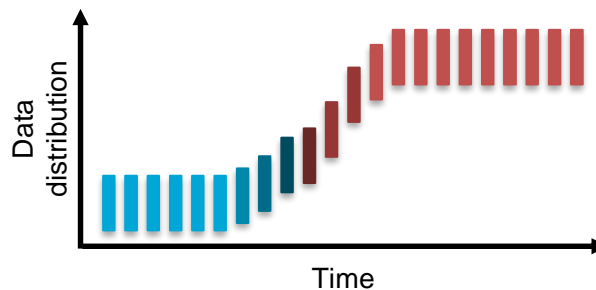
Sequence clustering in streaming environments

- Stream clustering
 - Cluster infinitely many data items
 - Concept drift



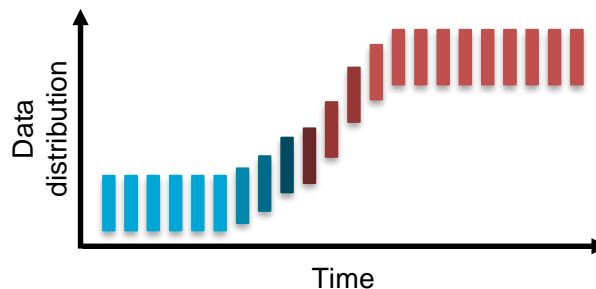
Sequence clustering in streaming environments

- Stream clustering
 - Cluster infinitely many data items
 - Concept drift
 - Heuristic-based solutions



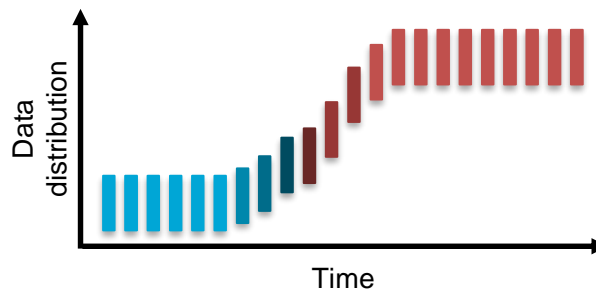
Sequence clustering in streaming environments


- Stream clustering
 - Cluster infinitely many data items
 - Concept drift
 - Heuristic-based solutions
- Sequence clustering while streaming

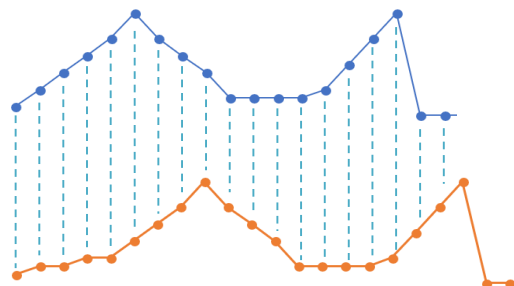


Sequence clustering in streaming environments

- Stream clustering
 - Cluster infinitely many data items
 - Concept drift
 - Heuristic-based solutions
- Sequence clustering while streaming

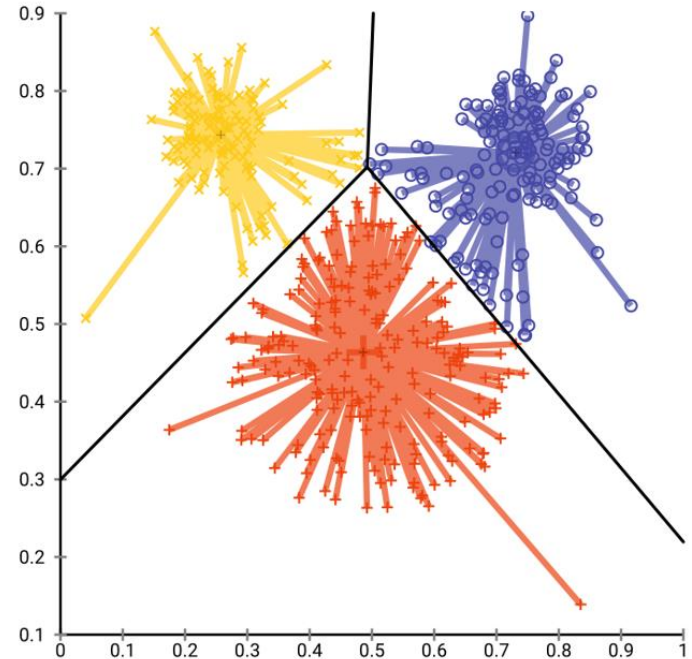


-  Out-of-sync sequences
 - Alignment-based distances, e.g., DTW
 - Expensive and limited support




K-medoids or Partitioning Around Medoids (PAM)

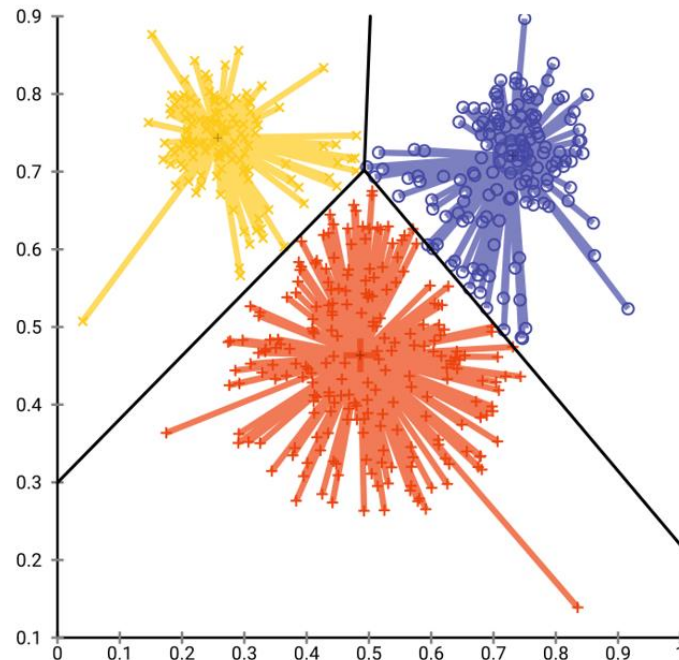
- K-centroids \rightarrow actual data items
- Supports non-metric distances



By Chire - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=82938091>

K-medoids or Partitioning Around Medoids (PAM)

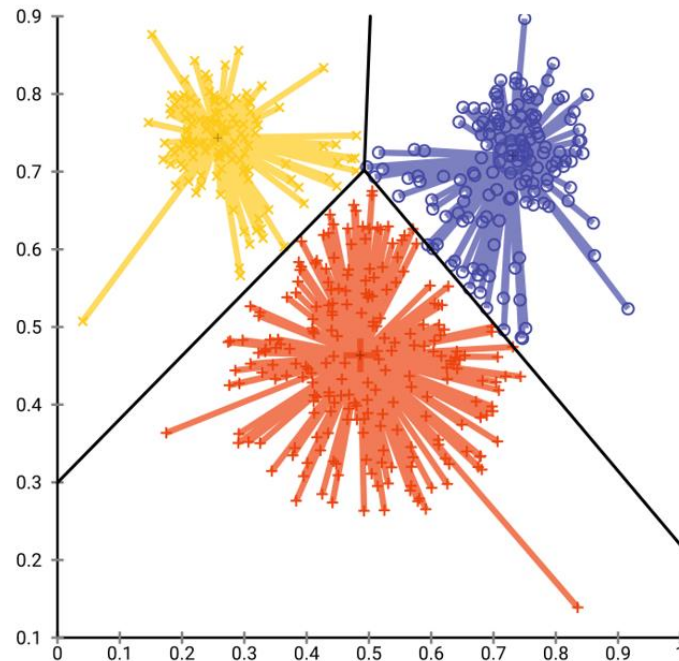
- K-centroids \rightarrow actual data items
 - Supports non-metric distances
-  Computationally expensive for streaming



By Chire - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=82938091>

K-medoids or Partitioning Around Medoids (PAM)

- K-centroids \rightarrow actual data items
 - Supports non-metric distances
- ✗ Computationally expensive for streaming
- ✗ No support for concept drift



By Chire - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=82938091>

SECLEDS: Sequene Clustering in Evolving Data Sstreams

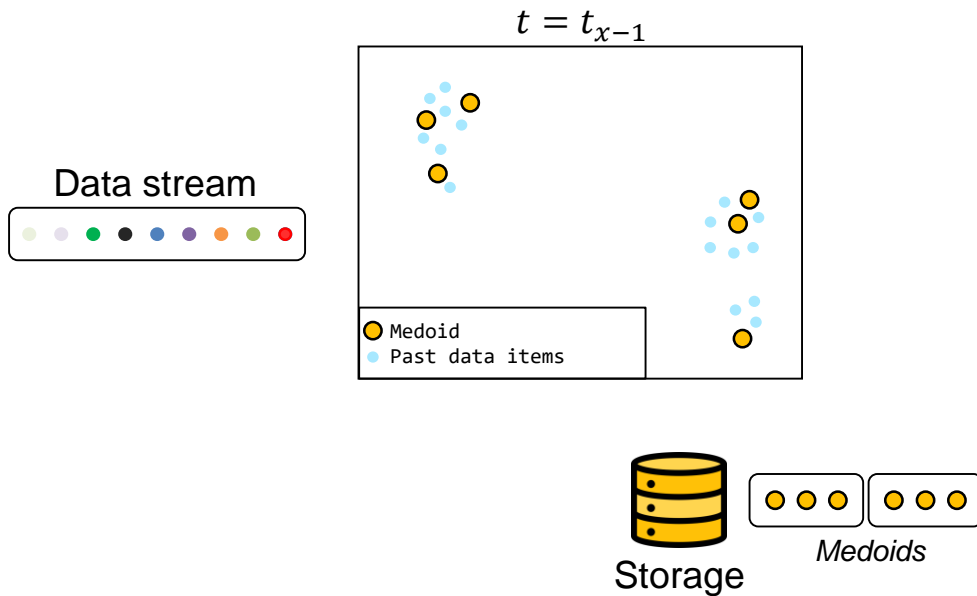
- Lightweight streaming variant of k-medoids with p-medoids

SECLEDS: Sequences Clustering in Evolving Data Sstreams

- Lightweight streaming variant of k-medoids with p-medoids
- Important properties
 1. Constant memory footprint
 2. Multiple medoids
 3. Medoid voting scheme

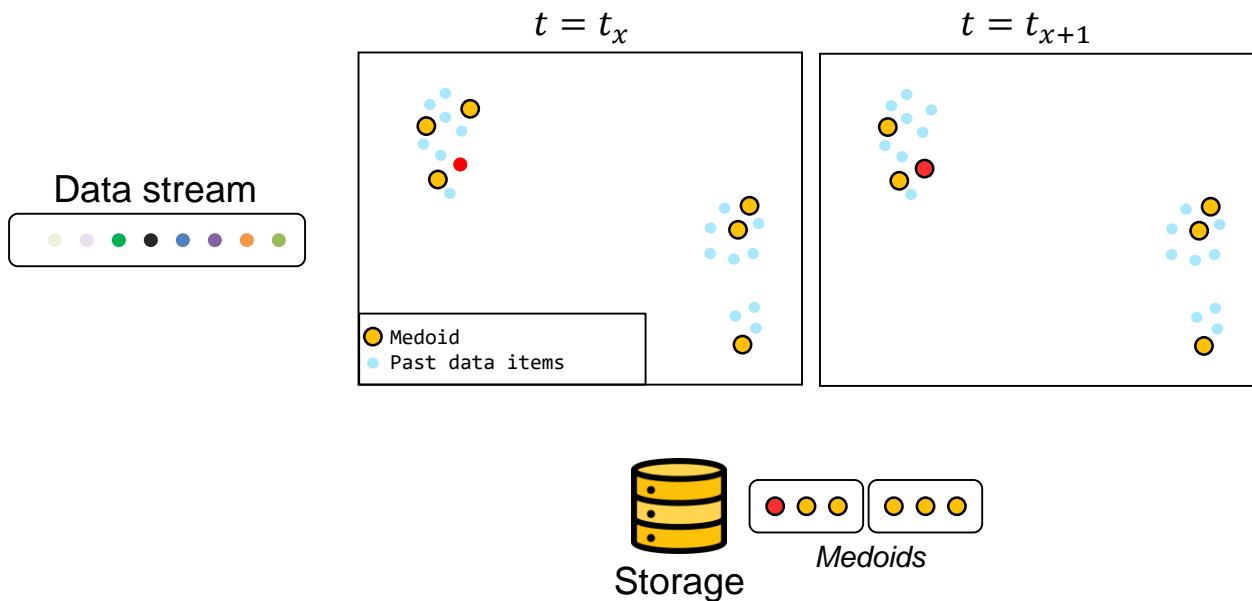
💡 Constant memory footprint

- Only stores the medoids $\rightarrow \mathcal{O}(kp)$



💡 Constant memory footprint

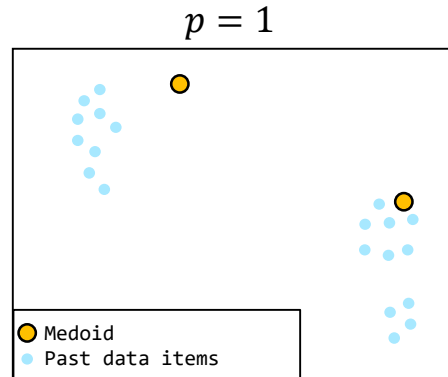
- Only stores the medoids $\rightarrow \mathcal{O}(kp)$





Multiple medoids

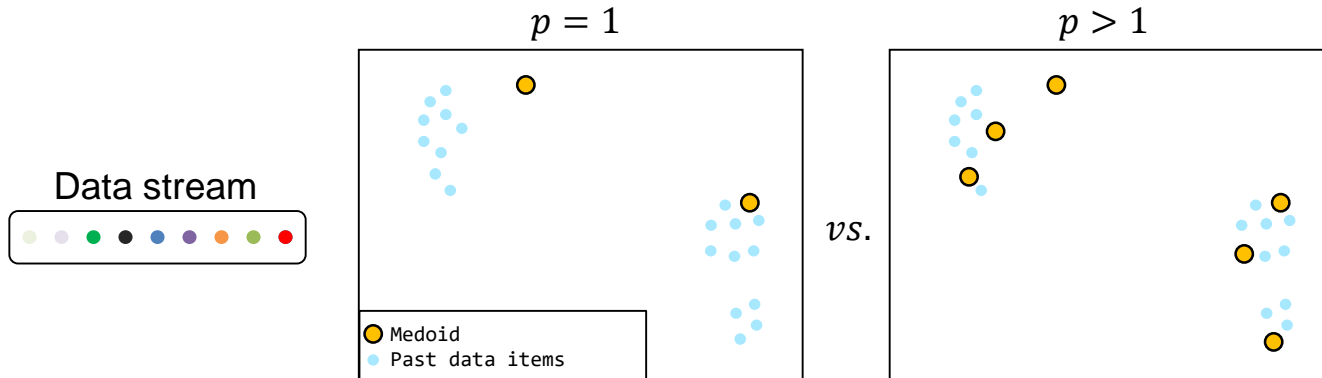
- p-medoids represent each cluster
 - p = configurable parameter





Multiple medoids

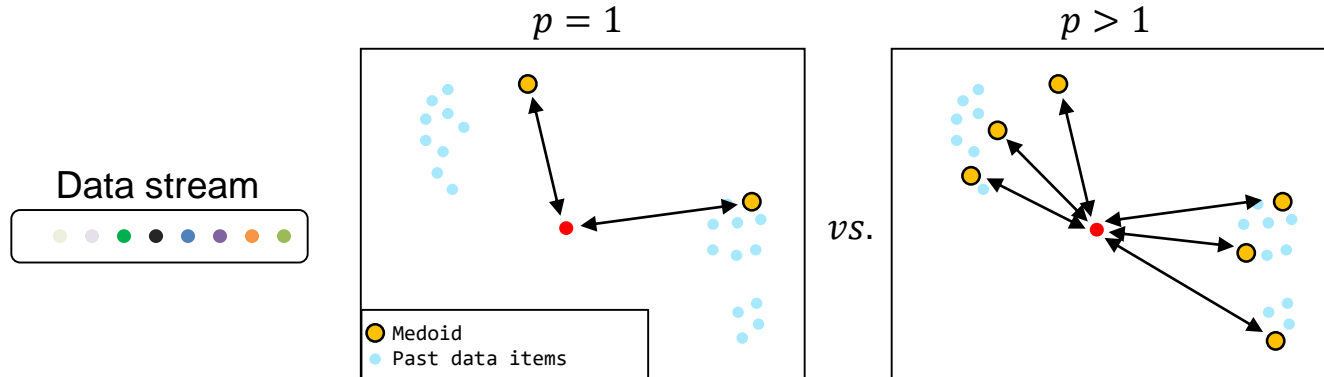
- p-medoids represent each cluster
 - $p = \text{configurable parameter}$





Multiple medoids

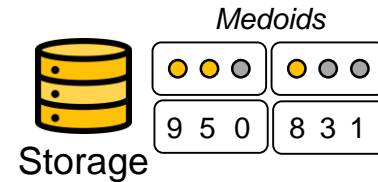
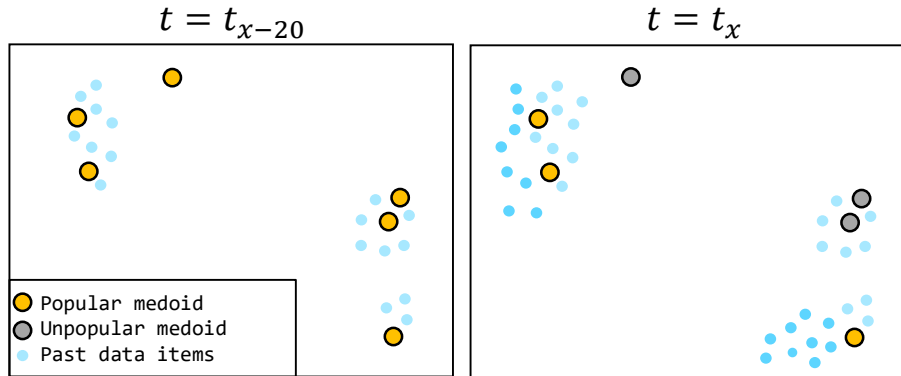
- p-medoids represent each cluster
 - $p = \text{configurable parameter}$
- Assignment based on minimum average distance
 - Reduces influence of irrelevant medoids





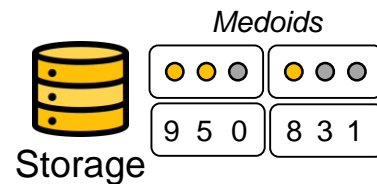
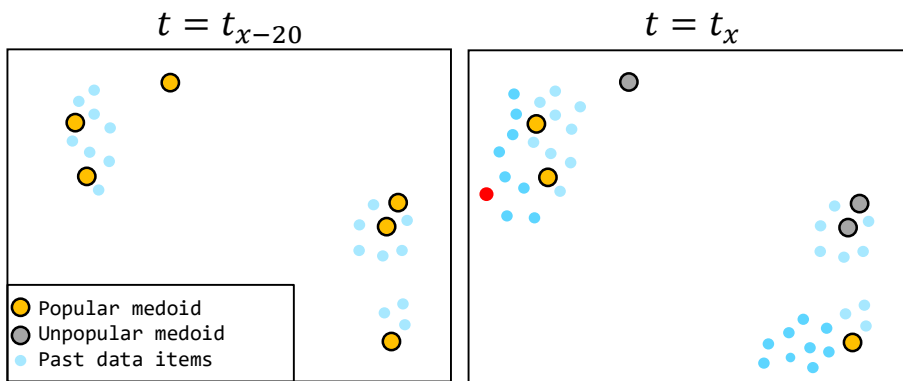
Medoid voting scheme

- Votes maintained for each medoid
 - Captures the fraction of close recent data
 - Exponential decay to forget past votes



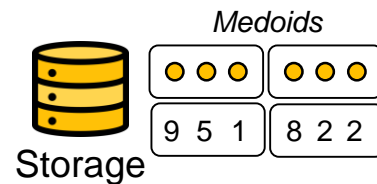
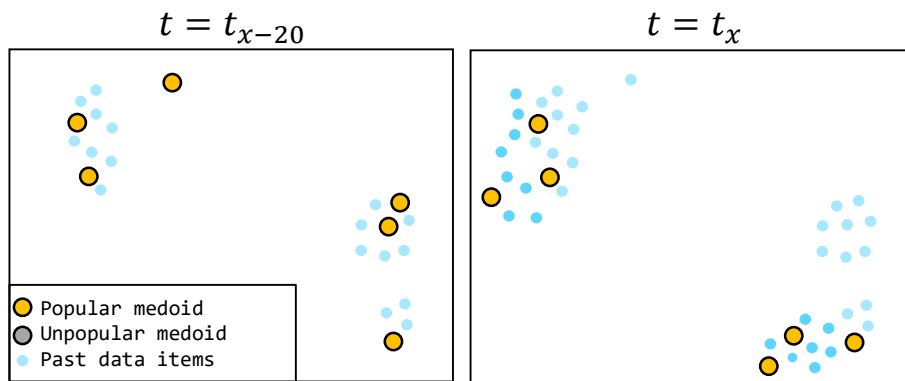
Medoid voting scheme

- Votes maintained for each medoid
 - Captures the fraction of close recent data
 - Exponential decay to forget past votes
- Which medoids to replace?
 - The one with the least votes → irrelevant medoid



Medoid voting scheme

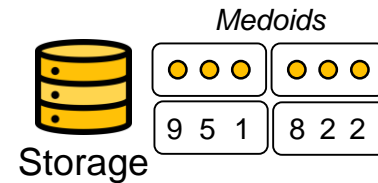
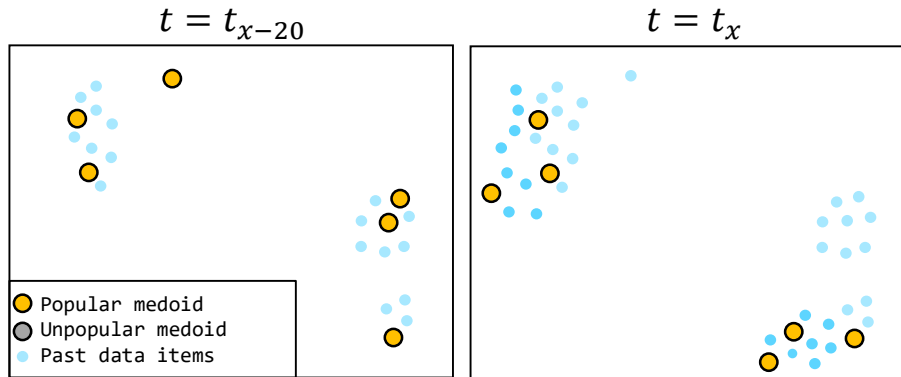
- Votes maintained for each medoid
 - Captures the fraction of close recent data
 - Exponential decay to forget past votes
- Which medoids to replace?
 - The one with the least votes → irrelevant medoid





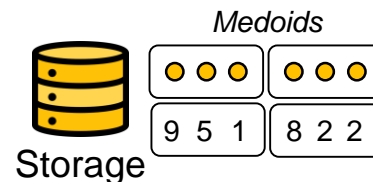
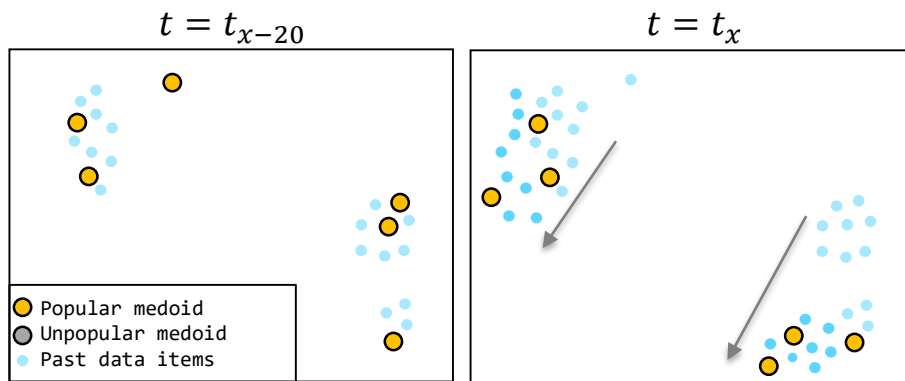
Medoid voting scheme

- Estimates center of mass without distance computations
 - Enables the use of expensive distance measures



💡💡💡 Medoid voting scheme

- Estimates center of mass without distance computations
 - Enables the use of expensive distance measures
- Concept drift with k-evolving clusters
 - Keep relevant medoids | Replace irrelevant medoids



SECLEDS

Algorithm 1: SECLEDS for clustering sequences in evolving streams

Input: Data stream, nclusters, nprototypes: S, k, p

```

1 function SECLEDS( $S, k, p$ )
2    $b \leftarrow 1.5 \cdot k \cdot p$ 
3    $\mathbb{B} \leftarrow$  Collect  $b$  items from  $S$ 
4    $\mathcal{C} \leftarrow \text{INIT}(\mathbb{B}, k, p)$  // INIT
5   forall  $s$  in  $S[b:]$  do
6      $cid \leftarrow \arg \min_{1 \leq cid \leq k} \frac{1}{p} \cdot \sum_{j=1}^p d(s, m_{cid,j})$  // ASSIGN
7      $j \leftarrow \arg \min_j d(s, m_{cid,j})$  for all  $1 \leq j \leq p$ 
8      $v_{cid,j} \leftarrow (v_{cid,j} + 1), v_{cid,j'} \leftarrow v_{cid,j'} \cdot (1 - \lambda)$  for  $j' \neq j$ 
9      $j \leftarrow \arg \min_j v_{cid,j}$  where  $m_{cid,j} \neq \eta_{cid}$  for all  $1 \leq j \leq p$  // UPDATE
10     $m_{cid,j} \leftarrow s, v_{cid,j} \leftarrow 0$ 
11  yield  $cid$ 

12 function INIT( $\mathbb{B}, k, p$ )
13  Choose  $m_{1,1} \in \mathbb{B}$  arbitrarily. Let  $C_1 \leftarrow \{(m_{1,1}, 0)\}$ 
14  for  $i \leftarrow 2 \dots k$  do
15    Choose  $m_{i,1} \in \mathbb{B}$  with probability  $d(m_{i,1}, m_{1,1})^2, m_{i,1} \neq m_{1,1}$ 
16    Let  $C_i \leftarrow \{(m_{i,1}, 0)\}$ 
17  for  $i \leftarrow 1 \dots k$  do
18     $dist \leftarrow d(b, m_{i,1})$  for all  $b \in \mathbb{B}$  and  $b \neq m_{i,1}$ 
19    Choose  $\{m_{i,2} \dots m_{i,p}\}$  having smallest values in  $dist$ 
20    Update  $C_i \leftarrow \{(m_{i,1}, 0) \dots (m_{i,p}, 0)\}$ 
21  return  $\{C_1, \dots, C_k\}$ 
    
```

SECLEDS

- From a batch, initialize clusters using a non-uniform sampling strategy

Algorithm 1: SECLEDS for clustering sequences in evolving streams

Input: Data stream, nclusters, nprototypes: S, k, p

```

1 function SECLEDS( $S, k, p$ )
2    $b \leftarrow 1.5 \cdot k \cdot p$ 
3    $\mathbb{B} \leftarrow$  Collect  $b$  items from  $S$ 
4    $\mathcal{C} \leftarrow \text{INIT}(\mathbb{B}, k, p)$  // INIT
5   forall  $s$  in  $S[b:]$  do
6      $cid \leftarrow \arg \min_{1 \leq cid \leq k} \frac{1}{p} \cdot \sum_{j=1}^p d(s, m_{cid,j})$  // ASSIGN
7      $j \leftarrow \arg \min_j d(s, m_{cid,j})$  for all  $1 \leq j \leq p$ 
8      $v_{cid,j} \leftarrow (v_{cid,j} + 1), v_{cid,j'} \leftarrow v_{cid,j'} \cdot (1 - \lambda)$  for  $j' \neq j$ 
9      $j \leftarrow \arg \min_j v_{cid,j}$  where  $m_{cid,j} \neq \eta_{cid}$  for all  $1 \leq j \leq p$  // UPDATE
10     $m_{cid,j} \leftarrow s, v_{cid,j} \leftarrow 0$ 
11    yield  $cid$ 

12 function INIT( $\mathbb{B}, k, p$ )
13   Choose  $m_{1,1} \in \mathbb{B}$  arbitrarily. Let  $C_1 \leftarrow \{(m_{1,1}, 0)\}$ 
14   for  $i \leftarrow 2 \dots k$  do
15     Choose  $m_{i,1} \in \mathbb{B}$  with probability  $d(m_{i,1}, m_{1,1})^2, m_{i,1} \neq m_{1,1}$ 
16     Let  $C_i \leftarrow \{(m_{i,1}, 0)\}$ 
17   for  $i \leftarrow 1 \dots k$  do
18      $dist \leftarrow d(b, m_{i,1})$  for all  $b \in \mathbb{B}$  and  $b \neq m_{i,1}$ 
19     Choose  $\{m_{i,2} \dots m_{i,p}\}$  having smallest values in  $dist$ 
20     Update  $C_i \leftarrow \{(m_{i,1}, 0) \dots (m_{i,p}, 0)\}$ 
21   return  $\{C_1, \dots, C_k\}$ 

```

SECLEDS

- From a batch, initialize clusters using a non-uniform sampling strategy

Stream loop starts

Algorithm 1: SECLEDS for clustering sequences in evolving streams

Input: Data stream, nclusters, nprototypes: S, k, p

```

1 function SECLEDS( $S, k, p$ )
2    $b \leftarrow 1.5 \cdot k \cdot p$ 
3    $\mathbb{B} \leftarrow$  Collect  $b$  items from  $S$ 
4    $\mathcal{C} \leftarrow \text{INIT}(\mathbb{B}, k, p)$  // INIT
5   forall  $s$  in  $S[b:]$  do
6      $cid \leftarrow \arg \min_{1 \leq cid \leq k} \frac{1}{p} \cdot \sum_{j=1}^p d(s, m_{cid,j})$  // ASSIGN
7      $j \leftarrow \arg \min_j d(s, m_{cid,j})$  for all  $1 \leq j \leq p$ 
8      $v_{cid,j} \leftarrow (v_{cid,j} + 1), v_{cid,j'} \leftarrow v_{cid,j'} \cdot (1 - \lambda)$  for  $j' \neq j$ 
9      $j \leftarrow \arg \min_j v_{cid,j}$  where  $m_{cid,j} \neq \eta_{cid}$  for all  $1 \leq j \leq p$  // UPDATE
10     $m_{cid,j} \leftarrow s, v_{cid,j} \leftarrow 0$ 
11    yield  $cid$ 

12 function INIT( $\mathbb{B}, k, p$ )
13   Choose  $m_{1,1} \in \mathbb{B}$  arbitrarily. Let  $C_1 \leftarrow \{(m_{1,1}, 0)\}$ 
14   for  $i \leftarrow 2 \dots k$  do
15     Choose  $m_{i,1} \in \mathbb{B}$  with probability  $d(m_{i,1}, m_{1,1})^2, m_{i,1} \neq m_{1,1}$ 
16     Let  $C_i \leftarrow \{(m_{i,1}, 0)\}$ 
17   for  $i \leftarrow 1 \dots k$  do
18      $dist \leftarrow d(b, m_{i,1})$  for all  $b \in \mathbb{B}$  and  $b \neq m_{i,1}$ 
19     Choose  $\{m_{i,2} \dots m_{i,p}\}$  having smallest values in  $dist$ 
20     Update  $C_i \leftarrow \{(m_{i,1}, 0) \dots (m_{i,p}, 0)\}$ 
21   return  $\{C_1, \dots, C_k\}$ 

```

SECLEDS

- From a batch, initialize clusters using a non-uniform sampling strategy

Stream loop starts

- Assign
 - Assign data item to the cluster with the least average distance to the medoids
 - Cast a vote for the closest medoid

Algorithm 1: SECLEDS for clustering sequences in evolving streams

Input: Data stream, nclusters, nprototypes: S, k, p

```

1 function SECLEDS( $S, k, p$ )
2    $b \leftarrow 1.5 \cdot k \cdot p$ 
3    $\mathbb{B} \leftarrow$  Collect  $b$  items from  $S$ 
4    $\mathcal{C} \leftarrow \text{INIT}(\mathbb{B}, k, p)$  // INIT
5   forall  $s$  in  $S$  [b:] do
6      $cid \leftarrow \arg \min_{1 \leq cid \leq k} \frac{1}{p} \cdot \sum_{j=1}^p d(s, m_{cid,j})$  // ASSIGN
7      $j \leftarrow \arg \min_j d(s, m_{cid,j})$  for all  $1 \leq j \leq p$ 
8      $v_{cid,j} \leftarrow (v_{cid,j} + 1), v_{cid,j'} \leftarrow v_{cid,j'} \cdot (1 - \lambda)$  for  $j' \neq j$ 
9      $j \leftarrow \arg \min_j v_{cid,j}$  where  $m_{cid,j} \neq \eta_{cid}$  for all  $1 \leq j \leq p$  // UPDATE
10     $m_{cid,j} \leftarrow s, v_{cid,j} \leftarrow 0$ 
11    yield  $cid$ 

12 function INIT( $\mathbb{B}, k, p$ )
13   Choose  $m_{1,1} \in \mathbb{B}$  arbitrarily. Let  $C_1 \leftarrow \{(m_{1,1}, 0)\}$ 
14   for  $i \leftarrow 2 \dots k$  do
15     Choose  $m_{i,1} \in \mathbb{B}$  with probability  $d(m_{i,1}, m_{1,1})^2, m_{i,1} \neq m_{1,1}$ 
16     Let  $C_i \leftarrow \{(m_{i,1}, 0)\}$ 
17   for  $i \leftarrow 1 \dots k$  do
18      $dist \leftarrow d(b, m_{i,1})$  for all  $b \in \mathbb{B}$  and  $b \neq m_{i,1}$ 
19     Choose  $\{m_{i,2} \dots m_{i,p}\}$  having smallest values in  $dist$ 
20     Update  $C_i \leftarrow \{(m_{i,1}, 0) \dots (m_{i,p}, 0)\}$ 
21   return  $\{C_1, \dots, C_k\}$ 

```

SECLEDS

- From a batch, initialize clusters using a non-uniform sampling strategy

Stream loop starts

- Assign
 - Assign data item to the cluster with the least average distance to the medoids
 - Cast a vote for the closest medoid
- Update
 - Upgrade data item as a medoid by replacing the one with the least votes

Algorithm 1: SECLEDS for clustering sequences in evolving streams

Input: Data stream, nclusters, nprototypes: S, k, p

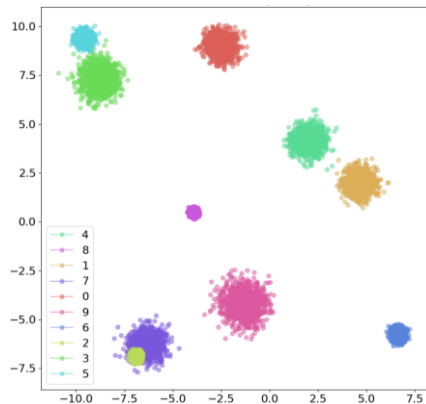
```

1 function SECLEDS( $S, k, p$ )
2    $b \leftarrow 1.5 \cdot k \cdot p$ 
3    $\mathbb{B} \leftarrow$  Collect  $b$  items from  $S$ 
4    $\mathcal{C} \leftarrow \text{INIT}(\mathbb{B}, k, p)$  // INIT
5   forall  $s$  in  $S$  [  $b$  ] do
6      $cid \leftarrow \arg \min_{1 \leq cid \leq k} \frac{1}{p} \cdot \sum_{j=1}^p d(s, m_{cid,j})$  // ASSIGN
7      $j \leftarrow \arg \min_j d(s, m_{cid,j})$  for all  $1 \leq j \leq p$ 
8      $v_{cid,j} \leftarrow (v_{cid,j} + 1), v_{cid,j'} \leftarrow v_{cid,j'} \cdot (1 - \lambda)$  for  $j' \neq j$ 
9      $j \leftarrow \arg \min_j v_{cid,j}$  where  $m_{cid,j} \neq \eta_{cid}$  for all  $1 \leq j \leq p$  // UPDATE
10     $m_{cid,j} \leftarrow s, v_{cid,j} \leftarrow 0$ 
11  yield  $cid$ 

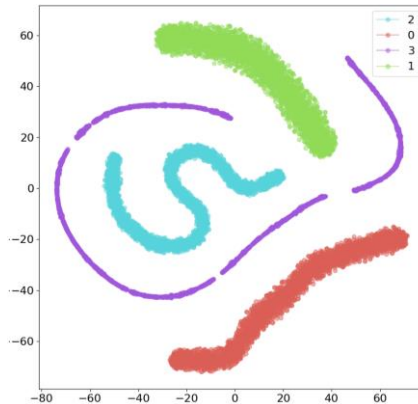
12 function INIT( $\mathbb{B}, k, p$ )
13  Choose  $m_{1,1} \in \mathbb{B}$  arbitrarily. Let  $C_1 \leftarrow \{(m_{1,1}, 0)\}$ 
14  for  $i \leftarrow 2 \dots k$  do
15    Choose  $m_{i,1} \in \mathbb{B}$  with probability  $d(m_{i,1}, m_{1,1})^2, m_{i,1} \neq m_{1,1}$ 
16    Let  $C_i \leftarrow \{(m_{i,1}, 0)\}$ 
17  for  $i \leftarrow 1 \dots k$  do
18     $dist \leftarrow d(b, m_{i,1})$  for all  $b \in \mathbb{B}$  and  $b \neq m_{i,1}$ 
19    Choose  $\{m_{i,2} \dots m_{i,p}\}$  having smallest values in  $dist$ 
20    Update  $C_i \leftarrow \{(m_{i,1}, 0) \dots (m_{i,p}, 0)\}$ 
21  return  $\{C_1, \dots, C_k\}$ 
    
```

Datasets

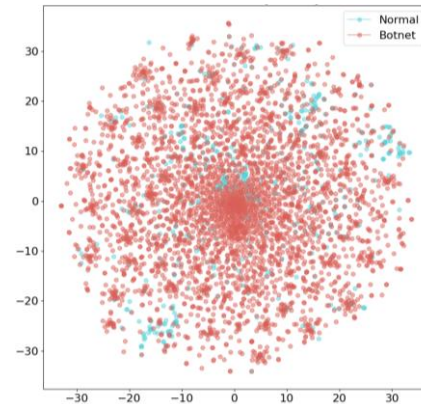
2D: Blobs ($k=10$)



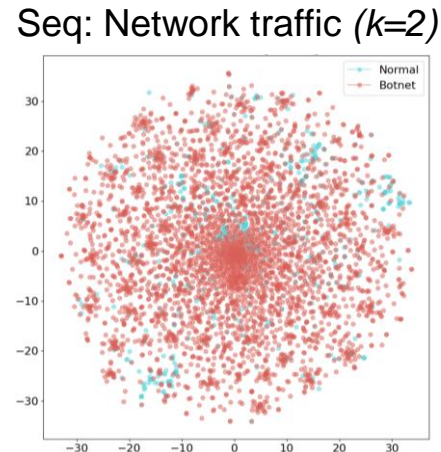
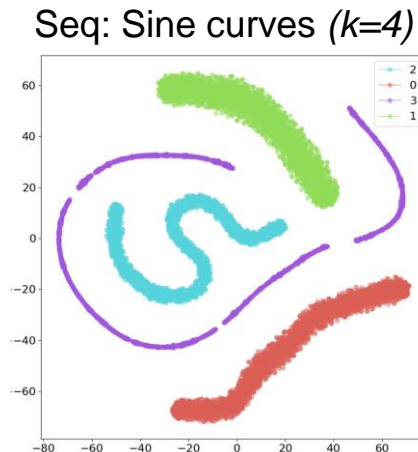
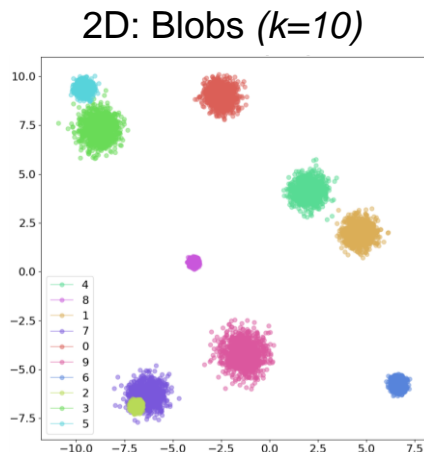
Seq: Sine curves ($k=4$)



Seq: Network traffic ($k=2$)



Datasets



Baselines

Streaming: CluStream¹, StreamKM++

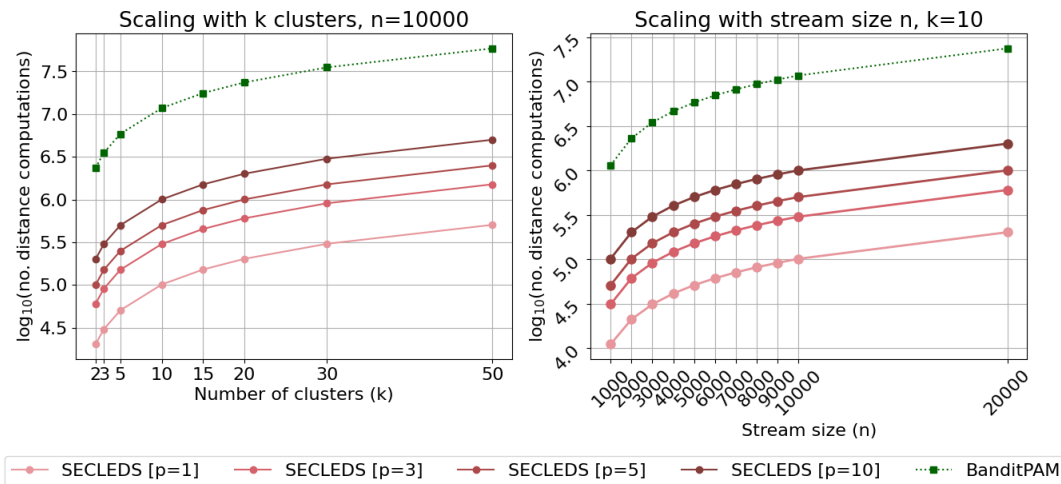
Batch: Minibatch k-means

Offline: BanditPAM²

Results [1/4]: Distance computations

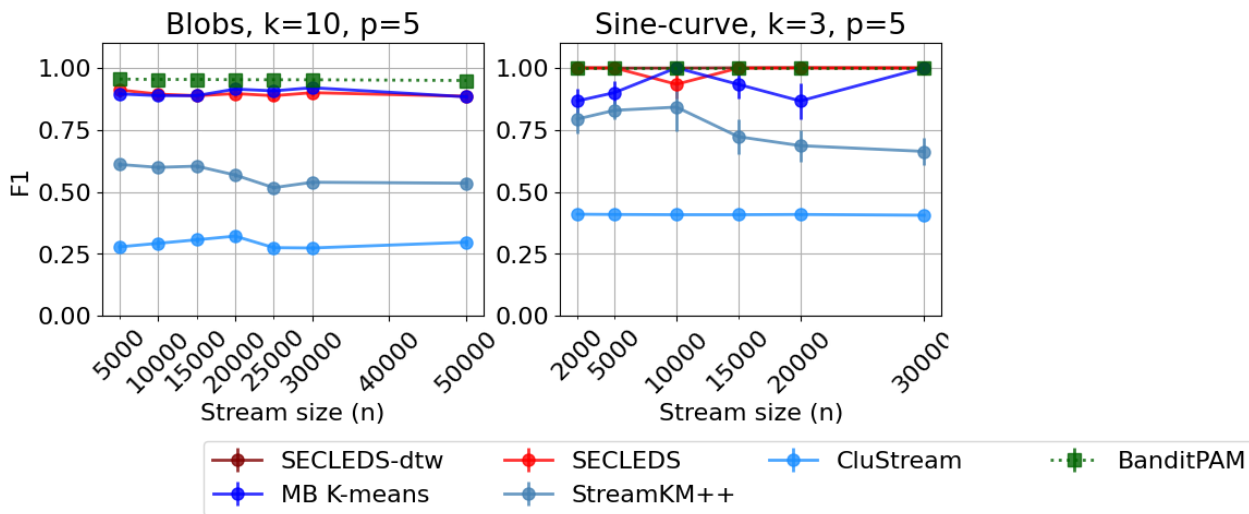
Results [1/4]: Distance computations

- Reduces required distance computations by 83.7%



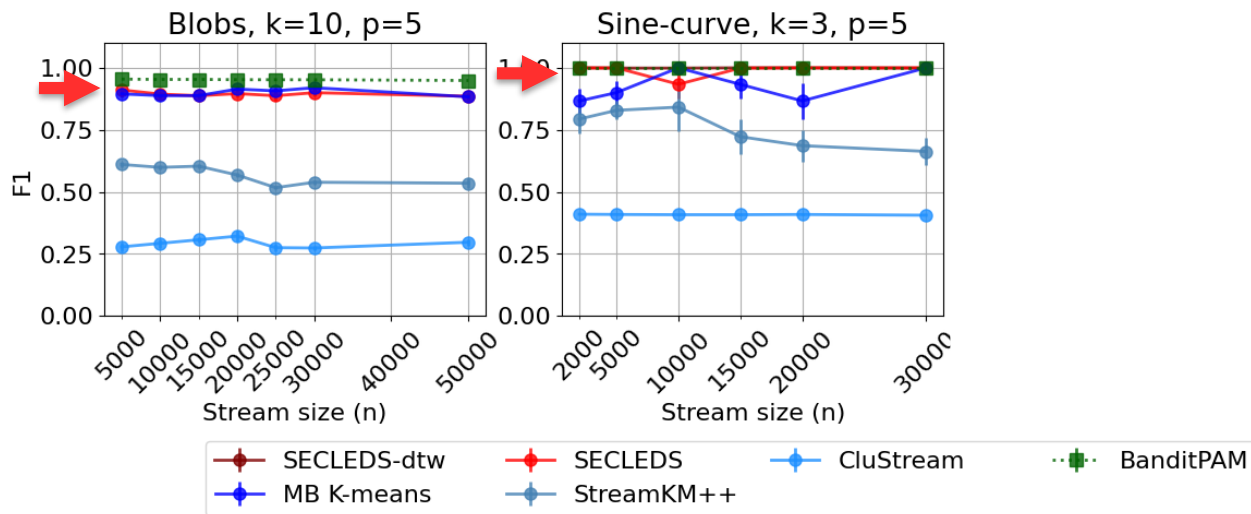
Results [2/4]: Performance

- Outperforms all streaming baselines without concept drift
 - Competitive performance with BanditPAM



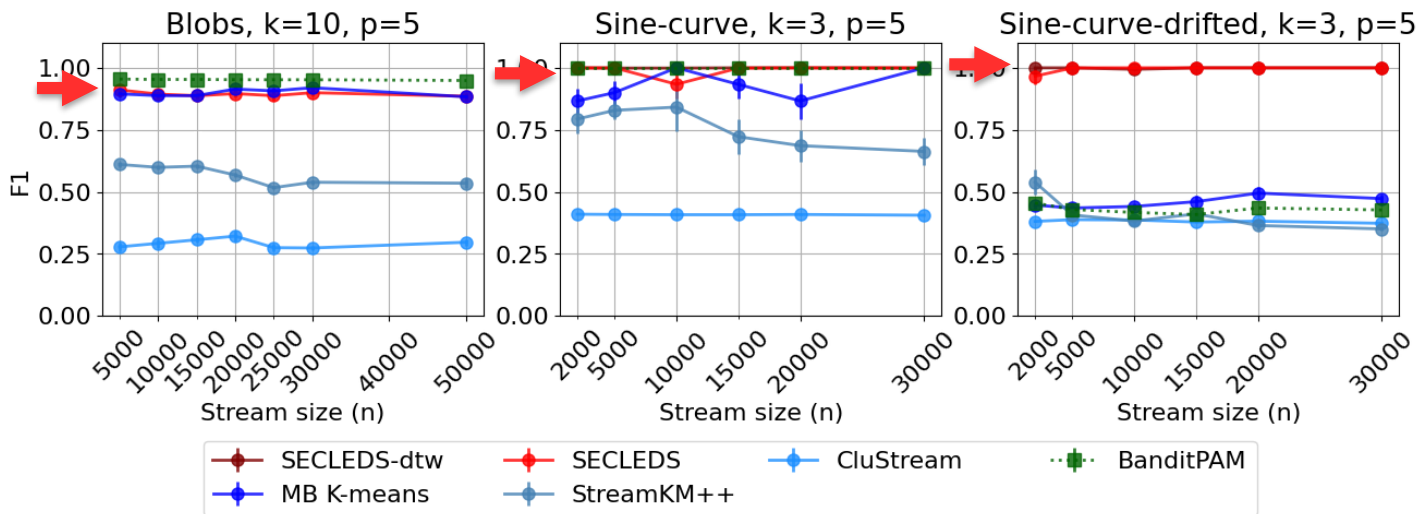
Results [2/4]: Performance

- Outperforms all streaming baselines without concept drift
 - Competitive performance with BanditPAM



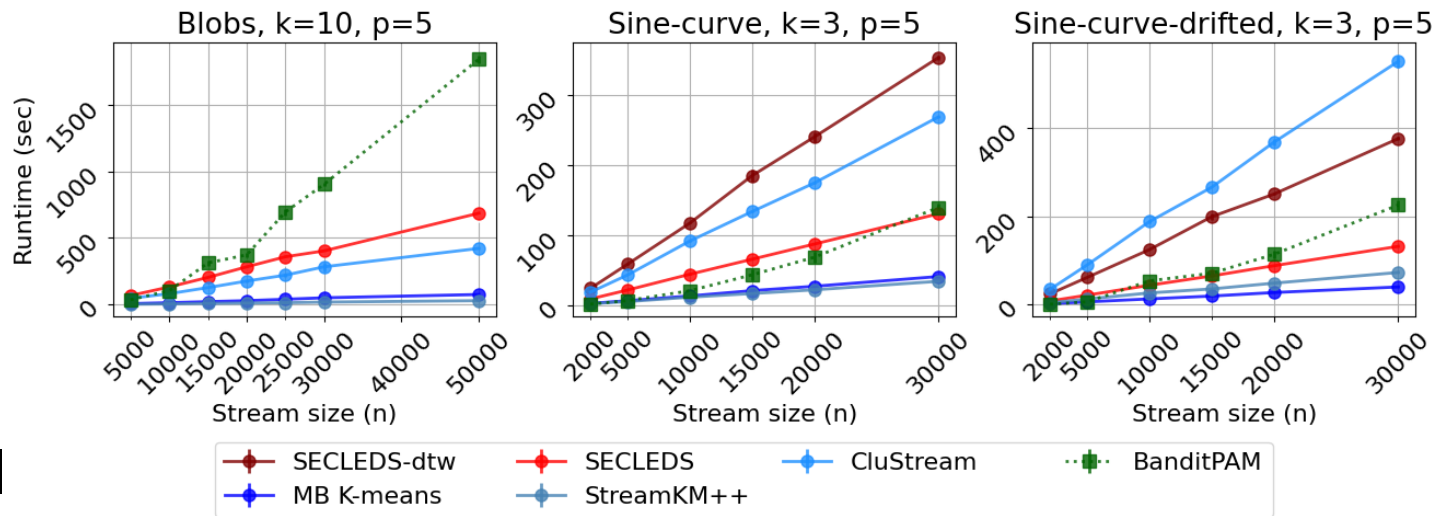
Results [2/4]: Performance

- Outperforms all streaming baselines without concept drift
 - Competitive performance with BanditPAM
- Outperforms all baselines by 138.7% with concept drift



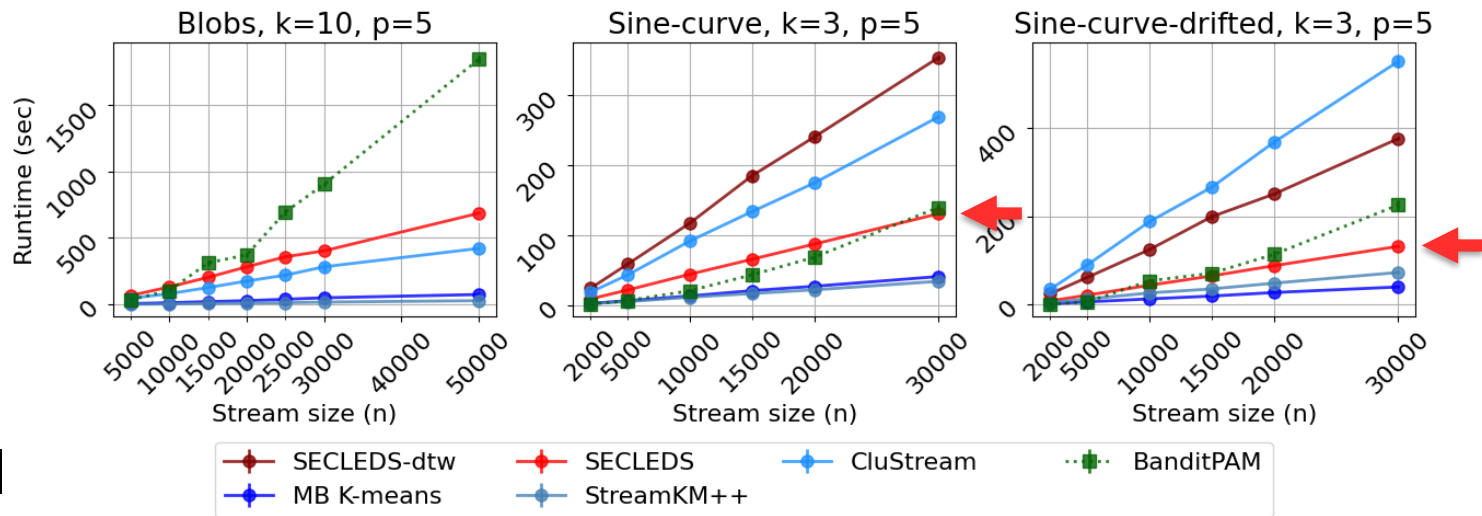
Results [3/4]: Runtime

- Runtime scales almost linearly w.r.t. stream size $\rightarrow \mathcal{O}(nkp)$



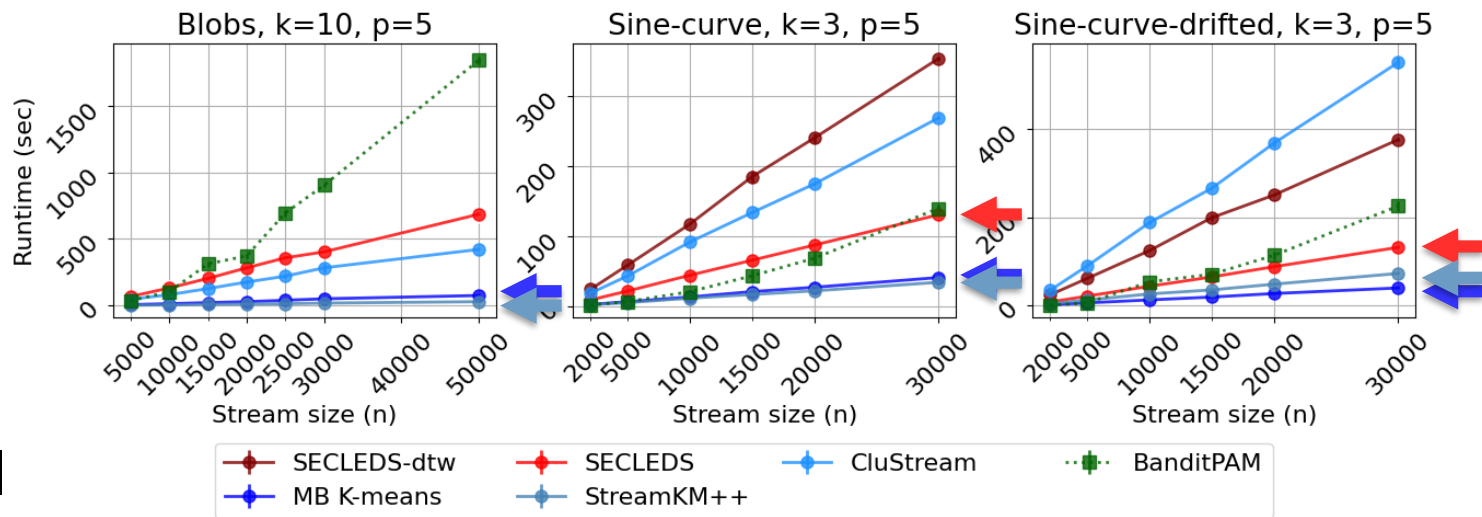
Results [3/4]: Runtime

- Runtime scales almost linearly w.r.t. stream size $\rightarrow \mathcal{O}(nkp)$
- Faster than banditPAM and CluStream for sequence clustering



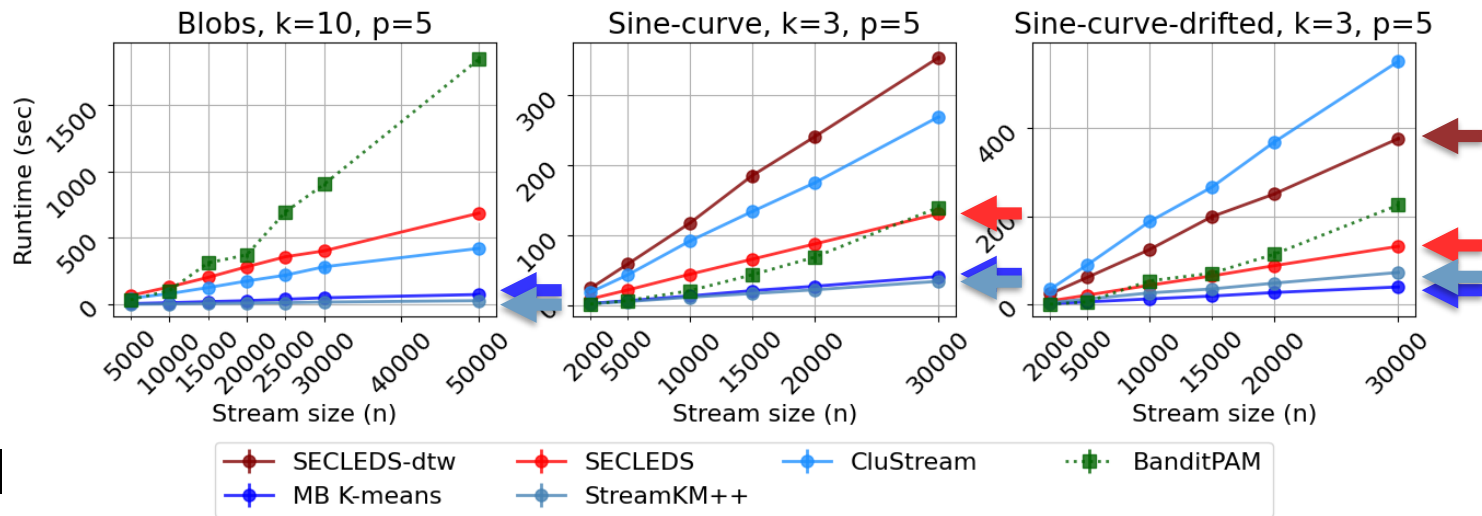
Results [3/4]: Runtime

- Runtime scales almost linearly w.r.t. stream size $\rightarrow \mathcal{O}(nkp)$
- Faster than banditPAM and CluStream for sequence clustering



Results [3/4]: Runtime

- Runtime scales almost linearly w.r.t. stream size $\rightarrow \mathcal{O}(nkp)$
- Faster than banditPAM and CluStream for sequence clustering



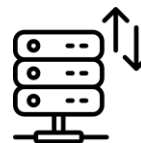
Results [4/4]: Network traffic use case

- Temporal pattern preserving network traffic sampling
 - Periodically store medoids

	# distances (10^6)	Runtime (sec)	F1 (k=2)	F1 (k=5)
BanditPAM	10.3	984.8	0.64	0.38
SECLEDS-eu	2.1	631.8	0.79	0.76
SECLEDS-dtw	2.1	1626.89	0.81	0.80

Results [4/4]: Network traffic use case

- Temporal pattern preserving network traffic sampling
 - Periodically store medoids
- Clusters 5.5h network traffic in **< 30 minutes**
 - Supporting network bandwidth of up to **1.08 Gb/s** with DTW



	# distances (10^6)	Runtime (sec)	F1 (k=2)	F1 (k=5)
BanditPAM	10.3	984.8	0.64	0.38
SECLEDS-eu	2.1	631.8	0.79	0.76
SECLEDS-dtw	2.1	1626.89	0.81	0.80

Future work

- Medoid-aware clustering
 - Better coverage, but additional distance computations required
- Theoretical properties of SECLEDS
 - Empirically works but limited theoretical understanding
- Optimized implementation
 - Runtime efficient C implementation

Takeways

- SECLEDS → lightweight streaming k-medoids with multiple medoids
 - Medoid voting → concept drift using k-evolving clusters
- Cluster sequences in high-bandwidth data streams
 - Reduces required distance computations
- Exceptional clustering performance with concept drift
 - Outperforms all streaming baselines regardless of concept drift
- SECLEDS can be used to intelligently sample network traffic

Thank you!

Questions?

- ▶ SECLEDS is a lightweight streaming variant of k-medoids.
- ▶ Efficiently clusters sequences in high-bandwidth data streams with concept drift.
- ▶ Can be used for temporal pattern-preserving network traffic sampling.



SECLEDS is open-source!

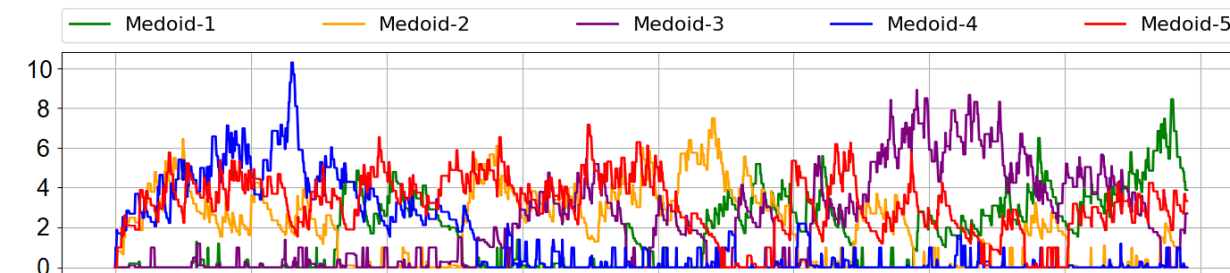
✉ azqa.nadeem@tudelft.nl

🐦 [@azqa_nadeem](https://twitter.com/azqa_nadeem)

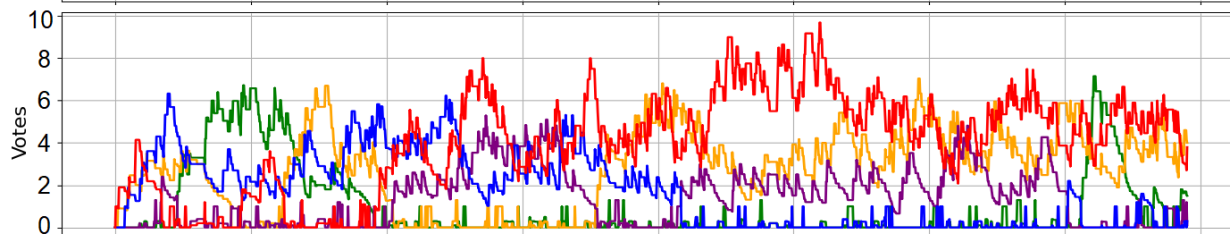
🌐 <https://cyber-analytics.nl/>

Extras

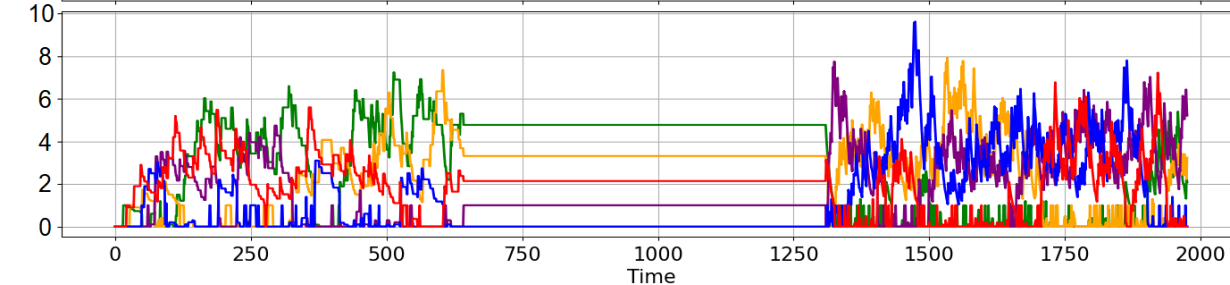
Uniformly distributed stream



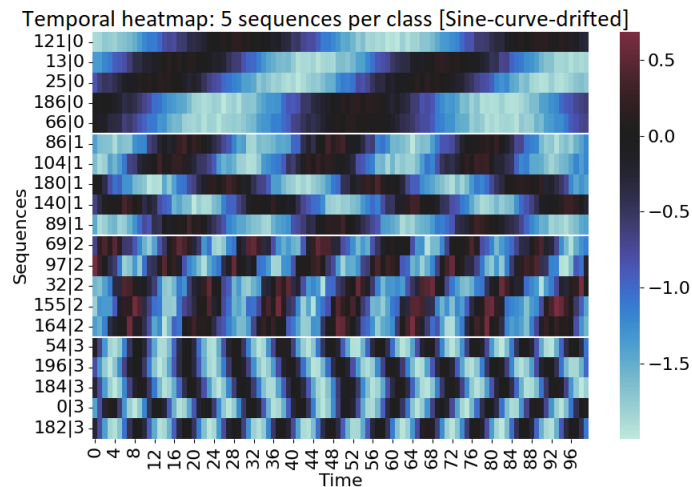
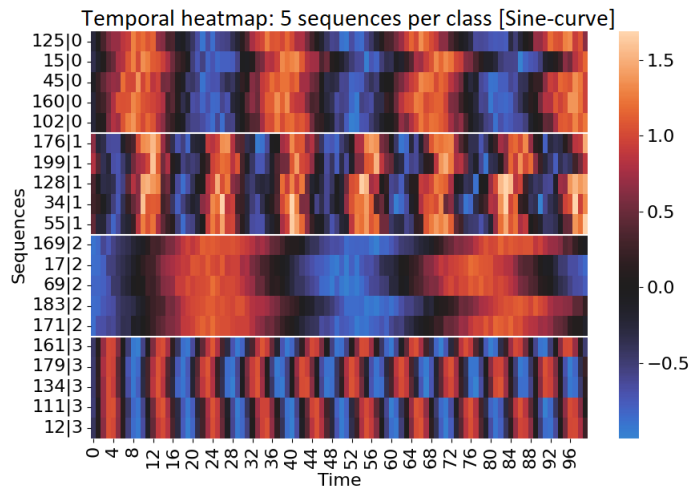
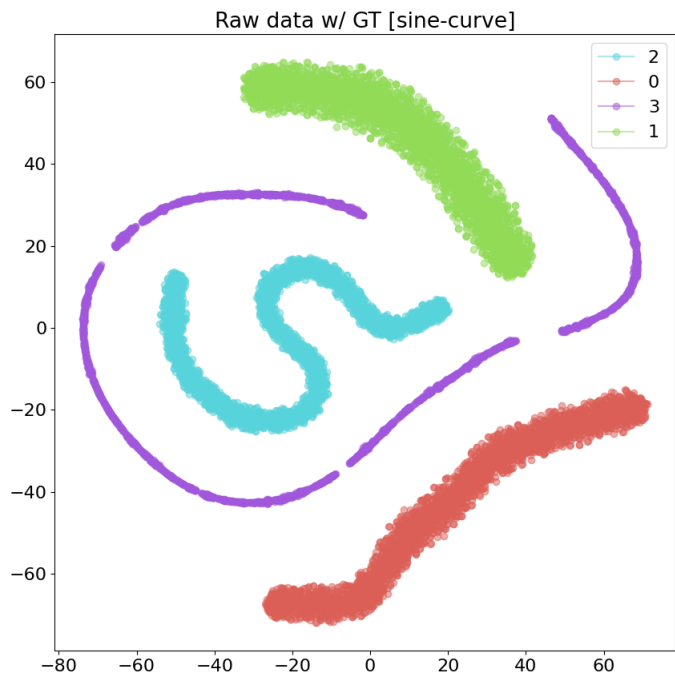
Incrementally drifted stream



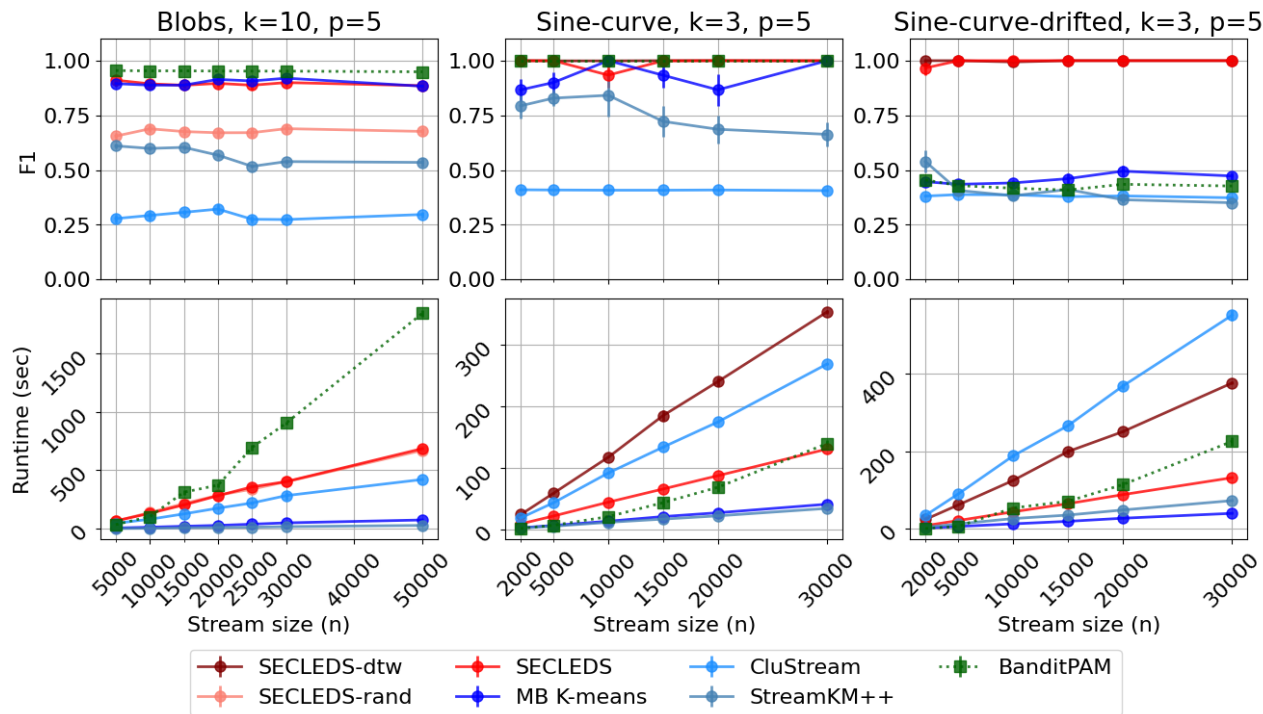
Class-ordered stream



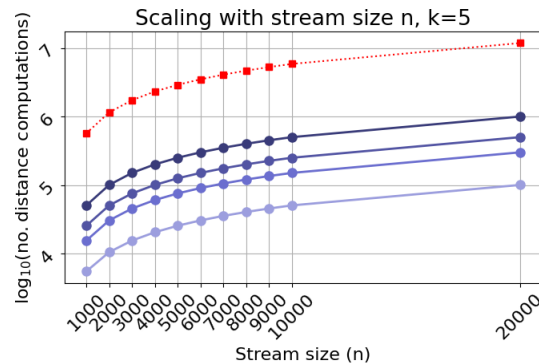
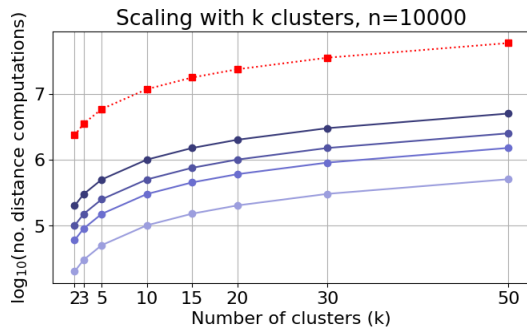
Extras



Extras

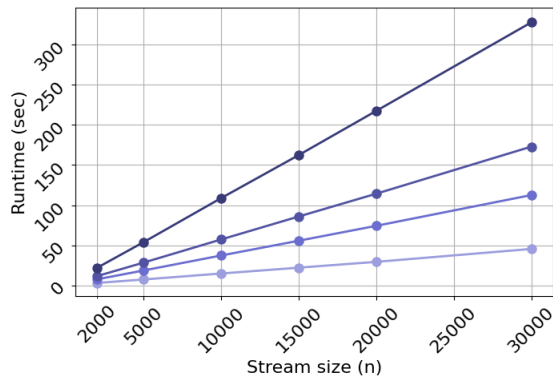
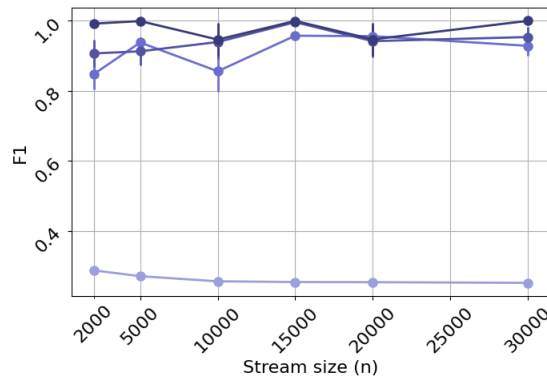


Extras



—●— SECELEDS [p=1] —●— SECELEDS [p=3] —●— SECELEDS [p=5] —●— SECELEDS [p=10] -.-■- BanditPAM

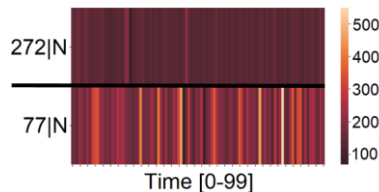
Sine-curve, k=4



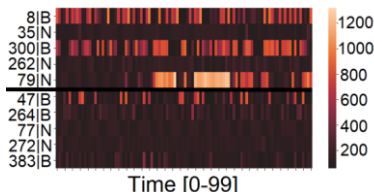
—●— SECELEDS [p=1] —●— SECELEDS [p=3] —●— SECELEDS [p=5] —●— SECELEDS [p=10]

Extras

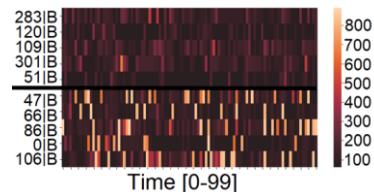
k=2



BanditPAM

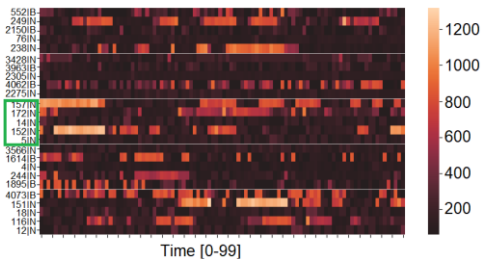


SECLEDS-eu

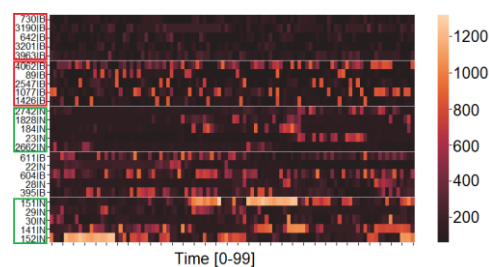


SECLEDS-dtw

k=5



SECLEDS-eu



SECLEDS-dtw