

## HIGH DIMENSIONAL & DYNAMIC DATA:

- High dimensionality
  - overlapping / irrelevant/ locally relevant attributes
- Dynamic/ Stream data
  - only 1 look at the data (upon their arrival) / volatile data

! Clustering upon such kind of data is very challenging!!  
- both members and dimensions might evolve over time

- ☹ Most existing approaches deal with 1 aspect of the problem
- ☹ HPStream assumes a constant #clusters over time

## OUR SOLUTION (HDDSTREAM):

1<sup>st</sup> density-based projected clustering algorithm for clustering high dimensional data streams:

- High dimensionality → projected clustering
- Stream data → online summarize– offline cluster
- Density based clustering → no assumption on #clusters/ invariant to outliers/ arbitrary shapes

- ☺ Quality improvement
- ☺ Bounded memory
- ☺ # summaries adapts to the underlying population

## PROJECTED MICROCLUSTERS:

For points  $C=\{p_1, \dots, p_n\}$  arriving at  $t_1, \dots, t_n$ , the summary models both content and dimension preferences of  $C$  at  $t$ .

$$mc \equiv mc(C,t) = \langle CF1(t), CF2(t), W(t) \rangle$$

Content summary

- $CF1(t), CF2(t)$ : temporal weighted linear/square sum of the points
- $W(t)$ : sum of points weights

$$\Phi(mc) = \langle \varphi_1, \varphi_2, \dots, \varphi_d \rangle$$

Dimension preference vector

- $j$  is a preferred dimension if:  $VAR_j(mc) \leq \delta$

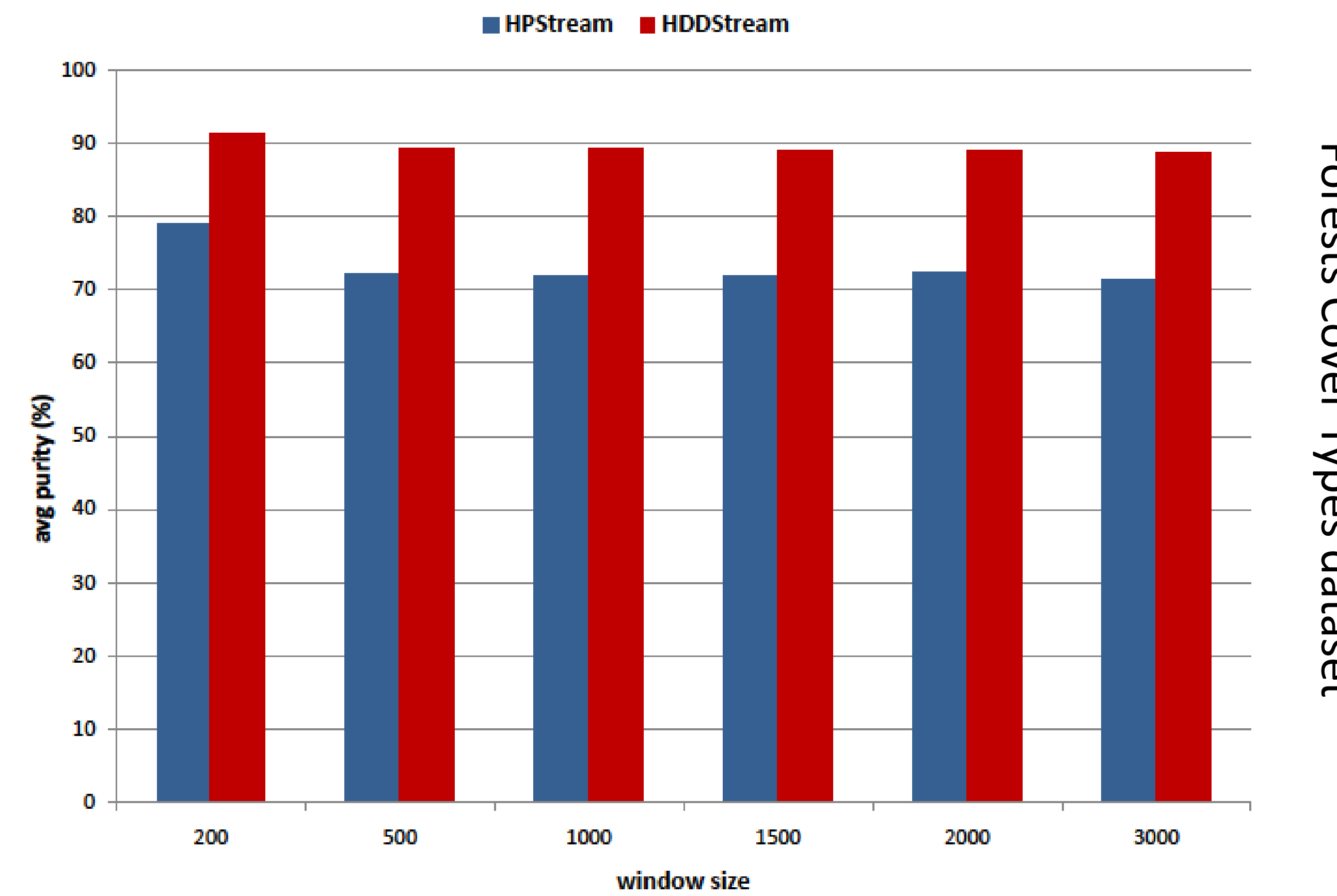
## CLUSTERS & OUTLIERS

- Core projected microclusters (CORE-PMC):
  - (1)  $radius^{\Phi}(mc) \leq \epsilon$  (radius criterion)
  - (2)  $W(t) \geq \mu$  (density criterion)
  - (3)  $PDIM(mc) \leq \pi$  (dimensionality criterion)

The role of clusters and outliers in a stream often exchange:

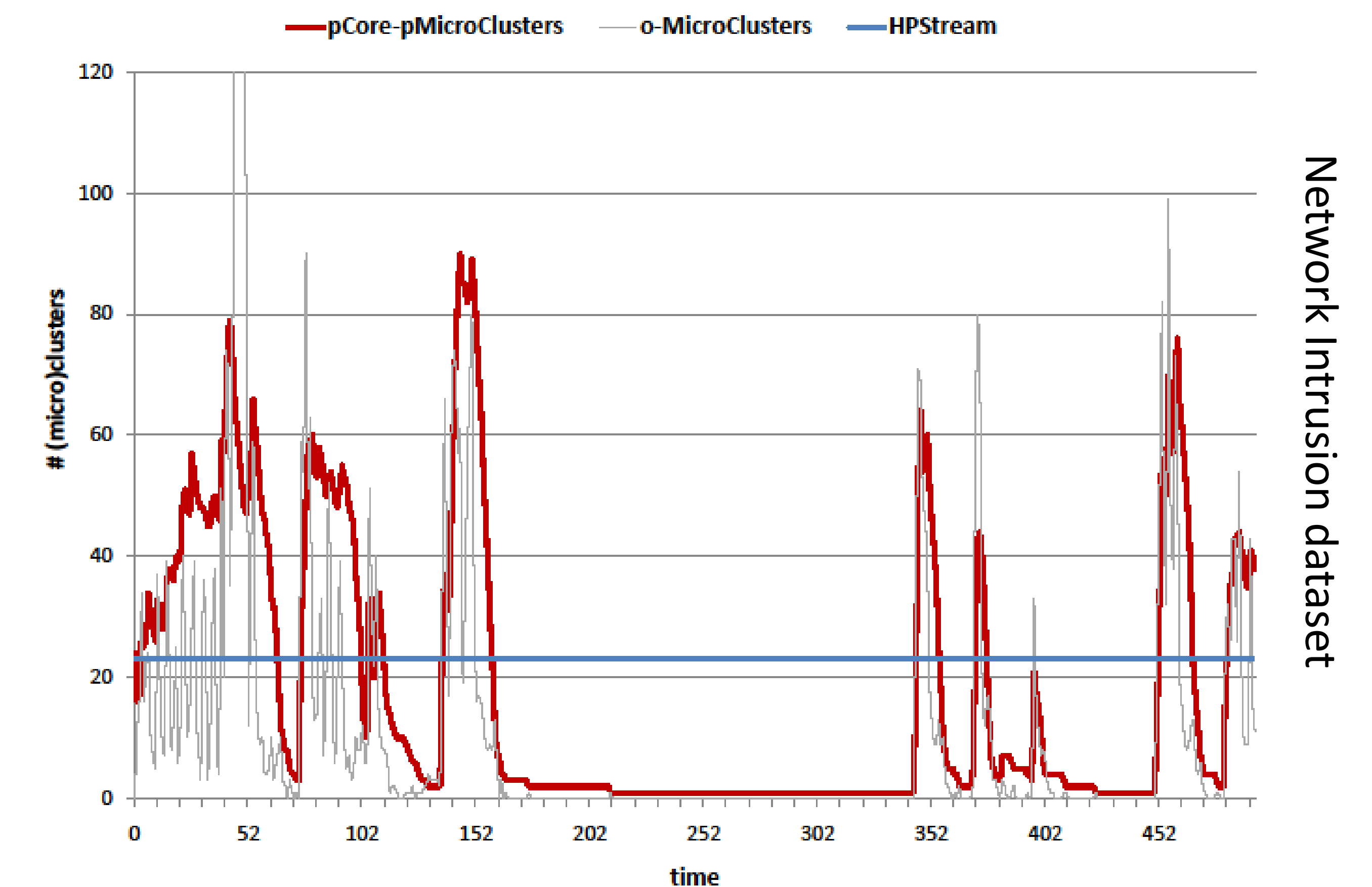
- Potential core PMC (PCORE-PMC):
  - (1), (2) relax the density criterion  $W(t) \geq \beta * \mu$ , (3)
- Outlier MC (o-MC) :
  - (1), (2) relax density criterion  $W(t) \geq \beta * \mu$ , (3) relax dim. criterion

## QUALITY IMPROVEMENT



Forests Cover Types dataset

## CHANGE DETECTION & MONITORING



Network Intrusion dataset



## HIGH DIMENSIONAL & DYNAMIC DATA:

- High dimensionality
  - overlapping / irrelevant/ locally relevant attributes
- Dynamic/ Stream data
  - only 1 look at the data (upon their arrival) / volatile data

! Clustering upon such kind of data is very challenging!!  
- both members and dimensions might evolve over time

- ☹ Most existing approaches deal with 1 aspect of the problem
- ☹ HPStream assumes a constant #clusters over time

## OUR SOLUTION (HDDSTREAM):

1<sup>st</sup> density-based projected clustering algorithm for clustering high dimensional data streams:

- High dimensionality → projected clustering
- Stream data → online summarize– offline cluster
- Density based clustering → no assumption on #clusters/ invariant to outliers/ arbitrary shapes

- ☺ Quality improvement
- ☺ Bounded memory
- ☺ # summaries adapts to the underlying population

## PROJECTED MICROCLUSTERS:

For points  $C=\{p_1, \dots, p_n\}$  arriving at  $t_1, \dots, t_n$ , the summary models both content and dimension preferences of  $C$  at  $t$ .

$$mc \equiv mc(C,t) = \langle CF1(t), CF2(t), W(t) \rangle$$

Content summary

- $CF1(t), CF2(t)$ : temporal weighted linear/square sum of the points
- $W(t)$ : sum of points weights

$$\Phi(mc) = \langle \varphi_1, \varphi_2, \dots, \varphi_d \rangle$$

Dimension preference vector

- $j$  is a preferred dimension if:  $VAR_j(mc) \leq \delta$

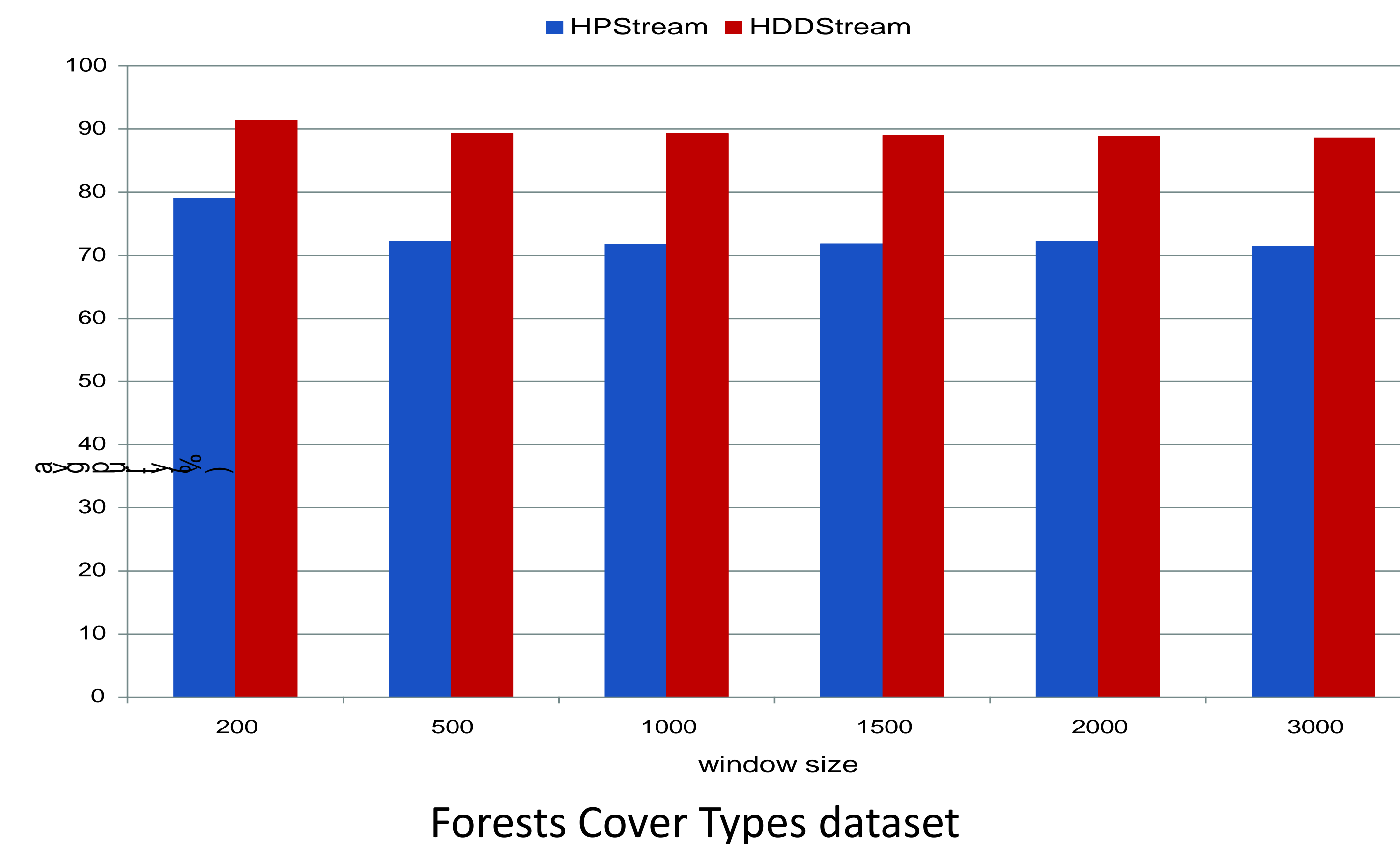
## CLUSTERS & OUTLIERS

- Core projected microclusters (CORE-PMC):
  - (1)  $radius^{\Phi}(mc) \leq \epsilon$  (radius criterion)
  - (2)  $W(t) \geq \mu$  (density criterion)
  - (3)  $PDIM(mc) \leq \pi$  (dimensionality criterion)

The role of clusters and outliers in a stream often exchange:

- Potential core PMC (PCORE-PMC):
  - (1), (2) relax the density criterion  $W(t) \geq \beta * \mu$ , (3)
- Outlier MC (o-MC) :
  - (1), (2) relax density criterion  $W(t) \geq \beta * \mu$ , (3) relax dim. criterion

## QUALITY IMPROVEMENT



## CHANGE DETECTION & MONITORING

