

Streaming Machine Learning Algorithms with Big Data Systems

Vibhatha Abeykoon, Supun Kamburugamuve, Kannan Govendrarajan,
Pulasthi Wickramasinghe, Chathura Widanage, Niranda Perera, Ahmet
Uyar, Gurhan Gunduz, Selahattin Akkas and Gregor Von Laszewski



Indiana University Bloomington

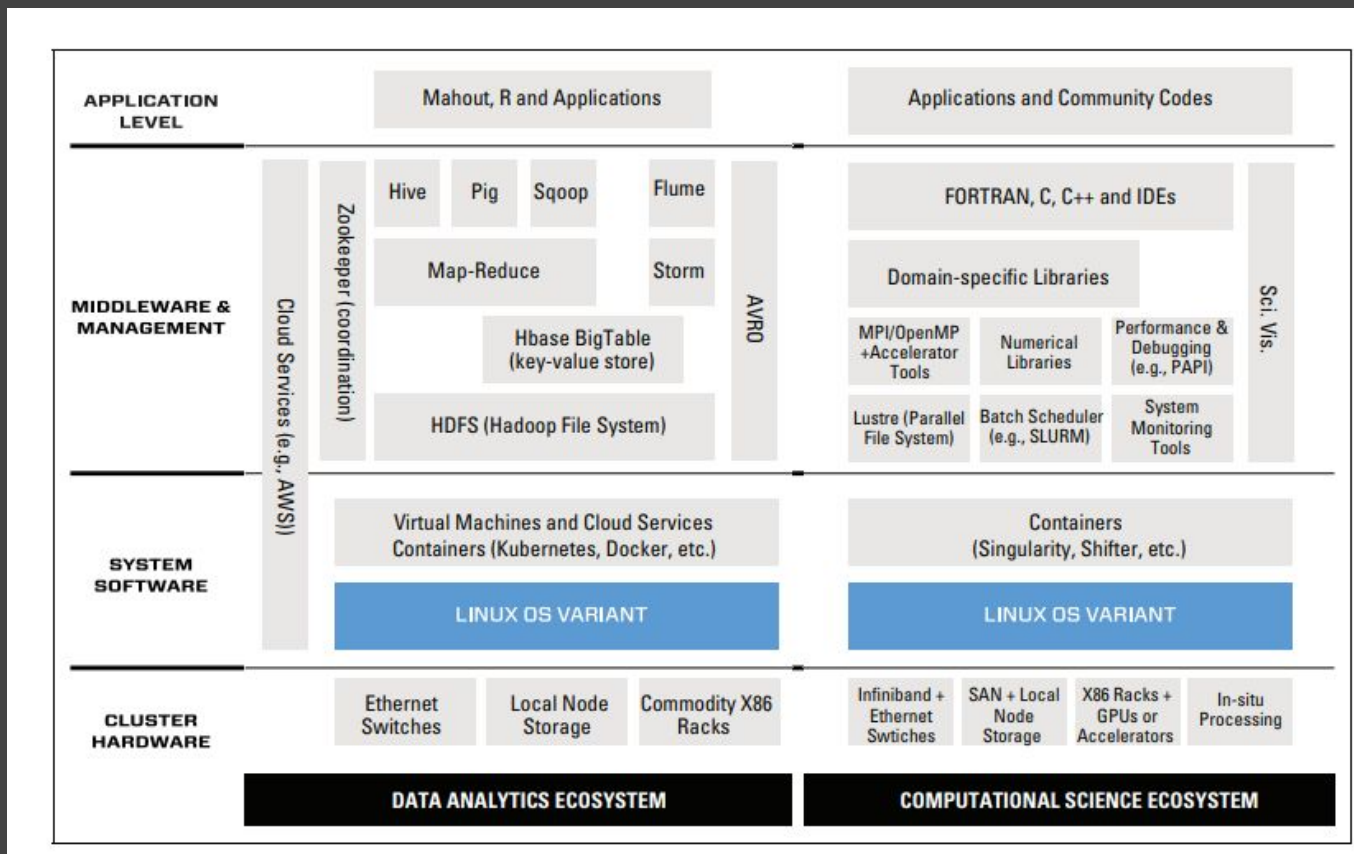
Motivation

- Data volume generated per day is increasing in a very high rate.
- Low latency is a must for increasing consumer demand on various services.
- Existing batch algorithms need to be optimized for online learning.
- Machine learning algorithms has become very important when formulating most of the supervised learning problems with less computing power.

How to design Streaming Machine Learning algorithms?

- Simply need to do train a machine learning algorithm in real-time without storing a large batches of data.
- Some algorithms can be trained by just observing a datapoint only once.
 - **Initialization stage:** Observe a number of data points (K elements at least if it is a clustering problem, depending on the algorithm this must be well-defined).
 - **Model Evaluation:** Calculate a gradient or model value for the observed elements.
 - **Model Synchronization:** Synchronize the model value across all the processes when using distributed training.
 - Re-do the whole process per element after the **initialization stage**.
- Some algorithms need an iterative streaming algorithm to ensure the accuracy to be in an expected level.
 - **Model evaluation:** Here we observe w number of elements by formulating a window in a stream and do an iterative computation on it for t iterations. Here $t \lll T$, T refers to the number of iterations required in batch mode to compute the optimum model.

Convergence of HPC and Big Data



Objective

- Design low-latency training on big data systems and identifying effective systems for online training
- Provide API solutions to design streaming applications on both HPC and dataflow programming models.
- Evaluate the importance of HPC frameworks for strengthening the big data stack for intensive computations.

Streaming Machine Learning Algorithms

- Non-Iterative Setting
 - KMeans Clustering
- Iterative Setting
 - Support Vector Machine (Linear Kernel for Binary classification)

Streaming SVM

Algorithm 1 Iterative SGD SVM

```
1: INPUT:  $[x, y] \in S, w \in R^d, t \in R^+$   
2: OUTPUT:  $w \in R^d$   
3: procedure ISGDSVM( $S, w, t$ )  
4:   for  $i = 0$  to  $n$  do  
5:     if  $(g(w; (x_i, y_i)) == 0)$  then  
6:        $\nabla J^t = w$   
7:     else  
8:        $\nabla J^t = w - Cx_i y_i$   
9:    $w = w - \alpha \nabla J^t$   
return  $w$ 
```

Algorithm 2 Iterative Streaming SVM

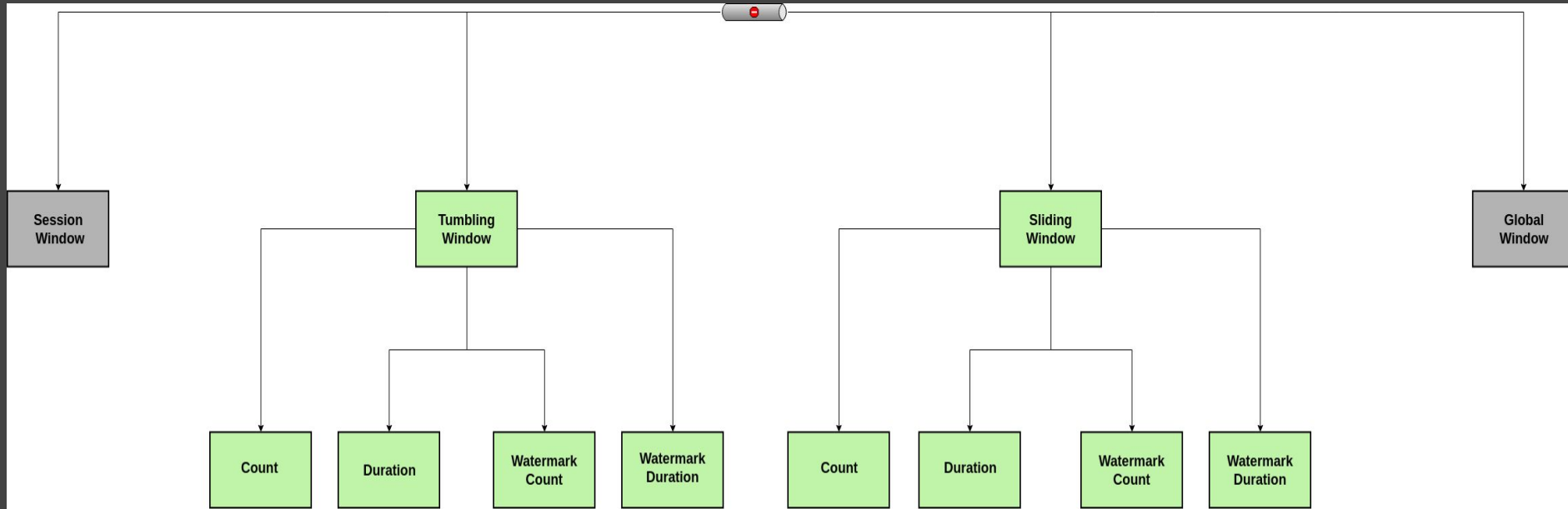
```
1: INPUT:  $X_\infty, Y_\infty \in S_\infty, w \in R^d, l \in R^+, s \in R^+, m <$   
    $K, m \in R^+$   
2: OUTPUT:  $w \in R^d$   
3: procedure ISSVM( $\bar{S}_i, w, T, l, s$ )  
4:   In Parallel K Machines  $[\bar{S}_1, \dots, \bar{S}_b] \subset S_\infty$   
5:   procedure WINDOW( $\bar{S}_m, w, l, s$ )  
6:     for  $t = 0$  to  $T$  do  
7:       procedure ISGDSVM( $\bar{S}_m, w, t$ )  
8:     All_Reduce( $w$ )  
return  $w$ 
```

Streaming KMeans

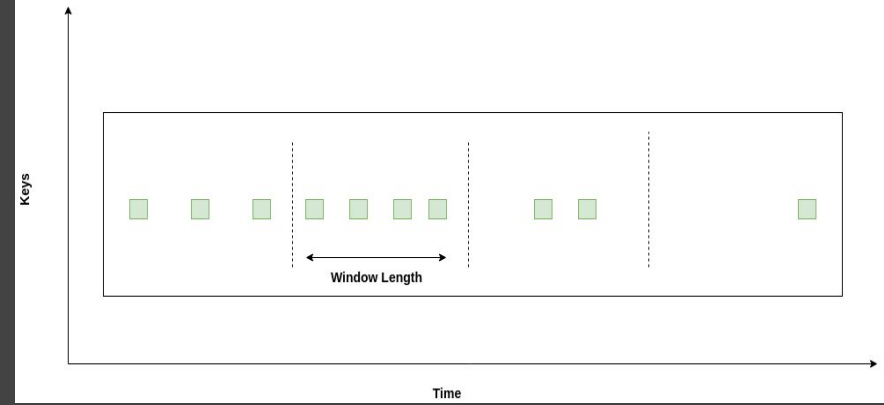
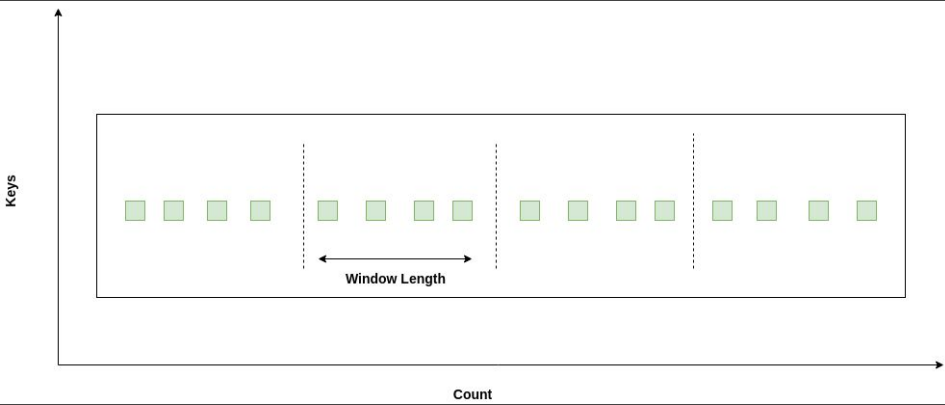
Algorithm 3 Online KMeans

```
1: INPUT:  $X = \{x_1, \dots, x_m\}, x_i \in \mathbf{R}^m$   
2:  $V = \{v_1, \dots, v_k\} v_i \in \mathbf{R}^m, k \leq n$   
3: OUTPUT:  $V$   
4: procedure STREAMING-KMEANS( $X, V$ )  
5:   procedure WINDOW( $\bar{X}, \bar{V}$ )  
6:     for  $x_j$  in  $\bar{X}$  do  
7:       if  $j \leq k$  then  
8:          $v_i = x_j$   
9:          $k_i = 1$   
10:         $i = i + 1$   
11:      else  
12:         $v_i = \operatorname{argmin}_i \|x_j - v_i\|$   
13:         $v_i = v_i + \frac{1}{n_i + 1} [x_j - v_i]$   
14:         $n_i = n_i + 1$   
15:   All_Reduce( $V$ )  
return  $\bar{V}$ 
```

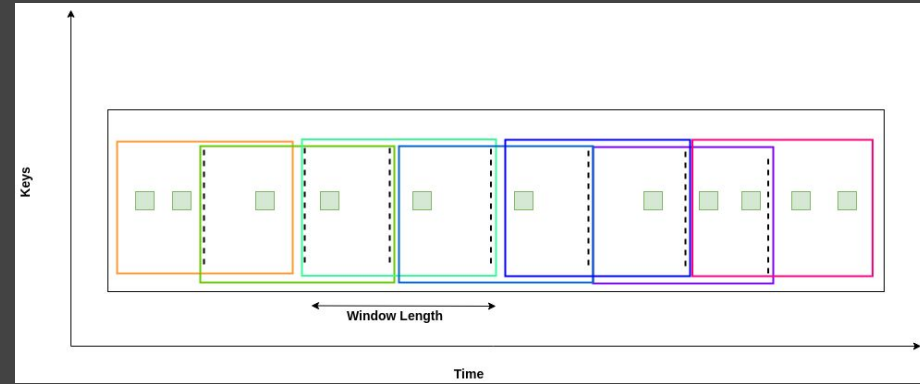
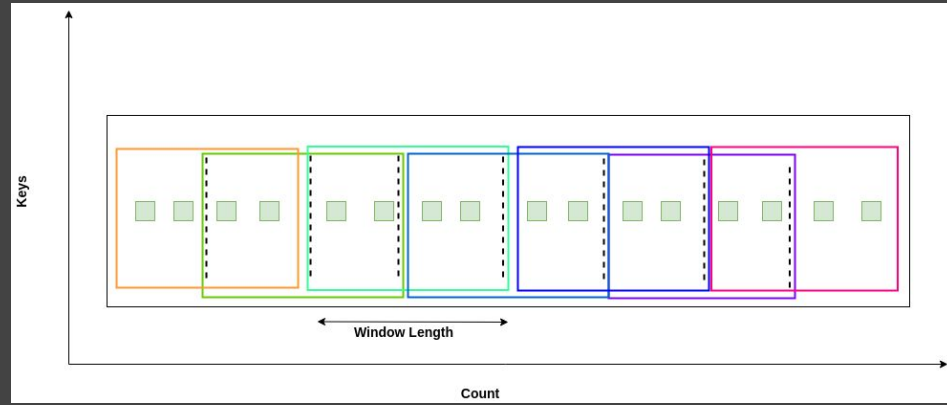
Discretization of a Stream



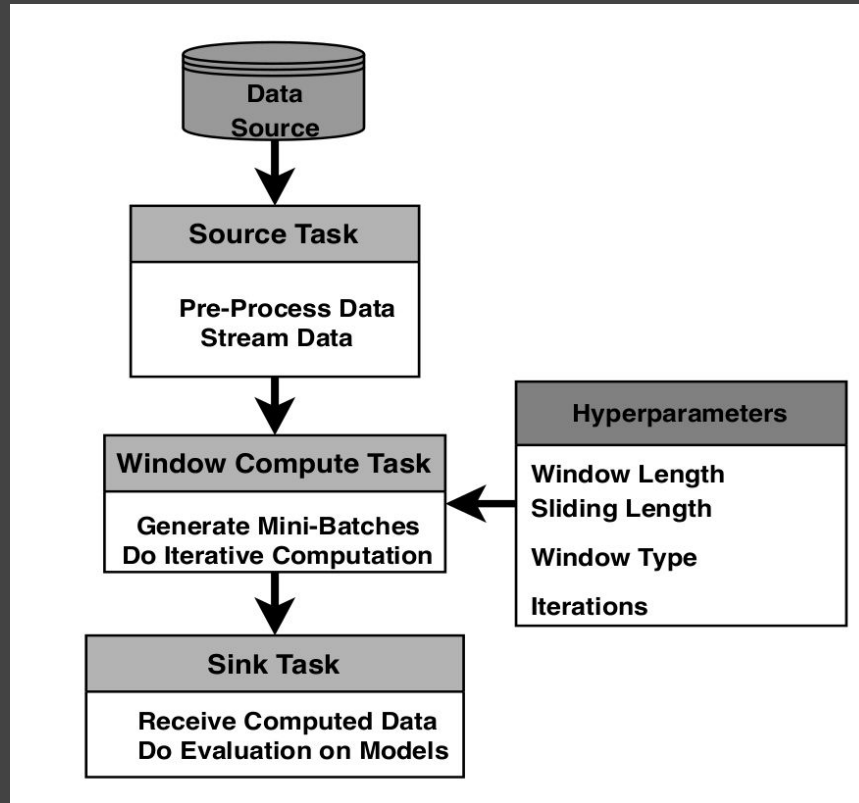
Tumbling Windows



Sliding Windows



Workflow of a Streaming ML Algorithm



Streaming Platforms



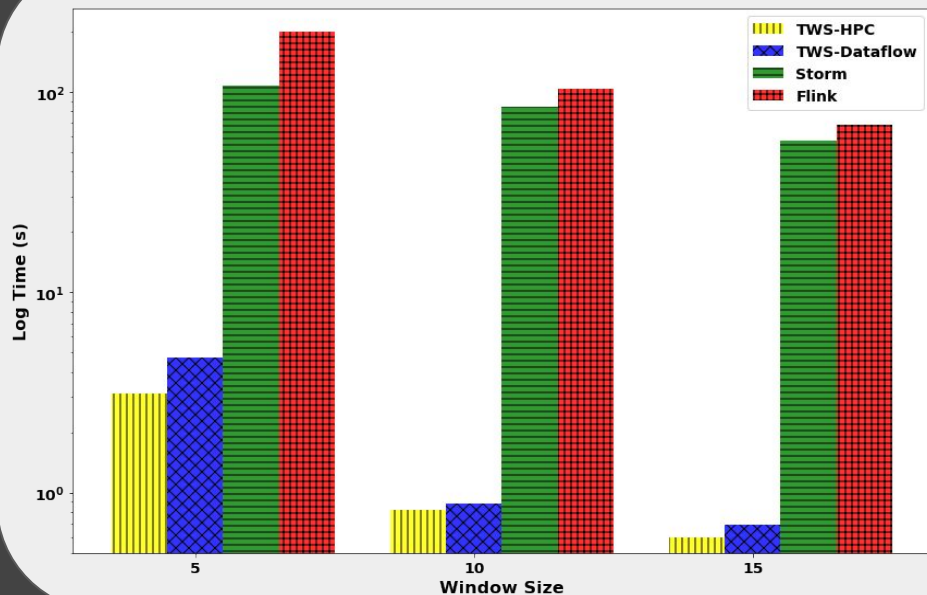
Iterative Streaming Support			
Dataflow Model			
HPC Model			

Apache Storm v1.2.8
Apache Flink v1.9.0
Twister2 v0.3.0

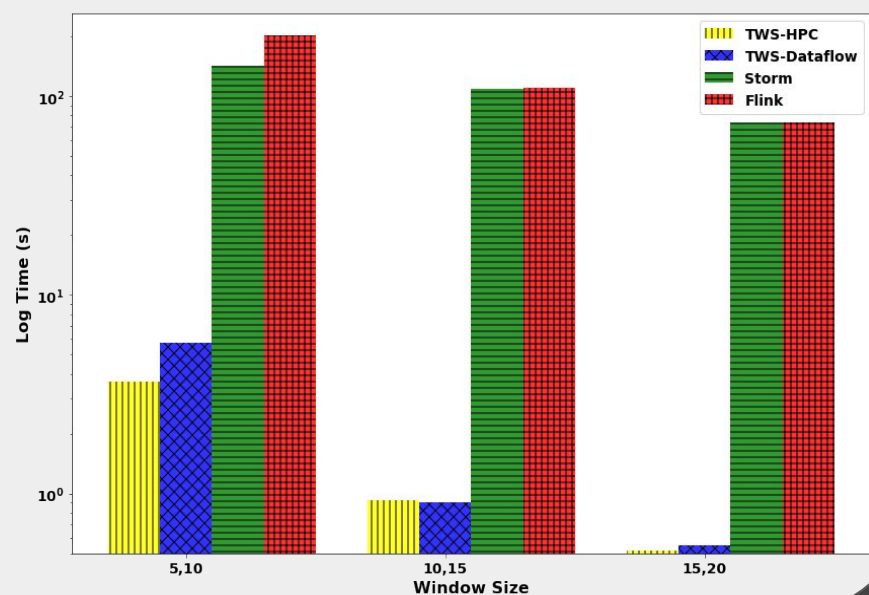
Experiment Configuration

- Intel(R) Xeon(R) Platinum 8160 CPU @ 2.10 GHz (250 GB RAM)
- Streaming SVM :Binary Classification on 49K long stream for training and 90K sample for model testing.
- Streaming KMeans: Clustering 1000 centroids, 49K long stream for training)
- 8 Physical nodes each with 16 processes (128 parallelism).
- Use count-based window setting to do a stress test on each big data framework used.

Streaming SVM



Tumbling Windowing

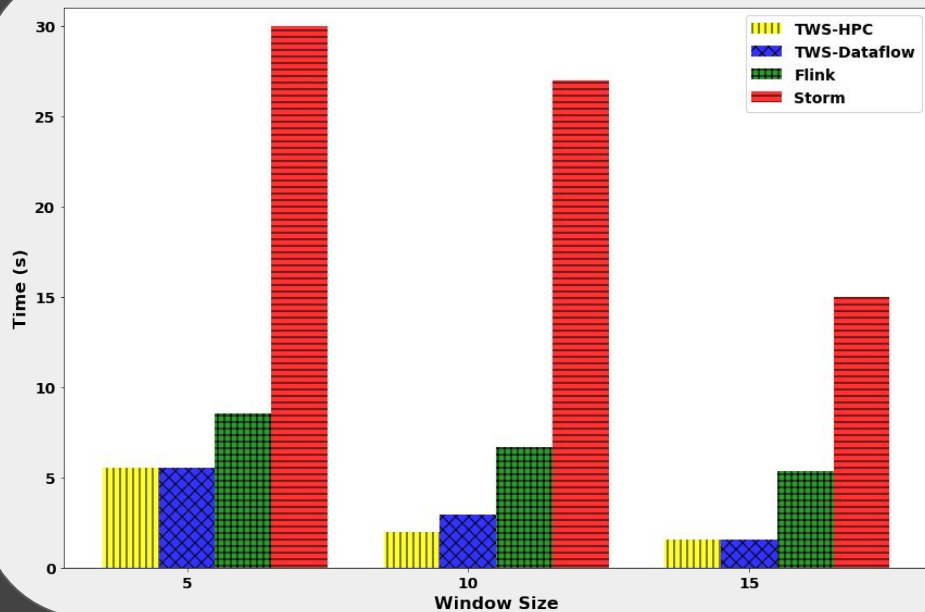


Sliding Windowing

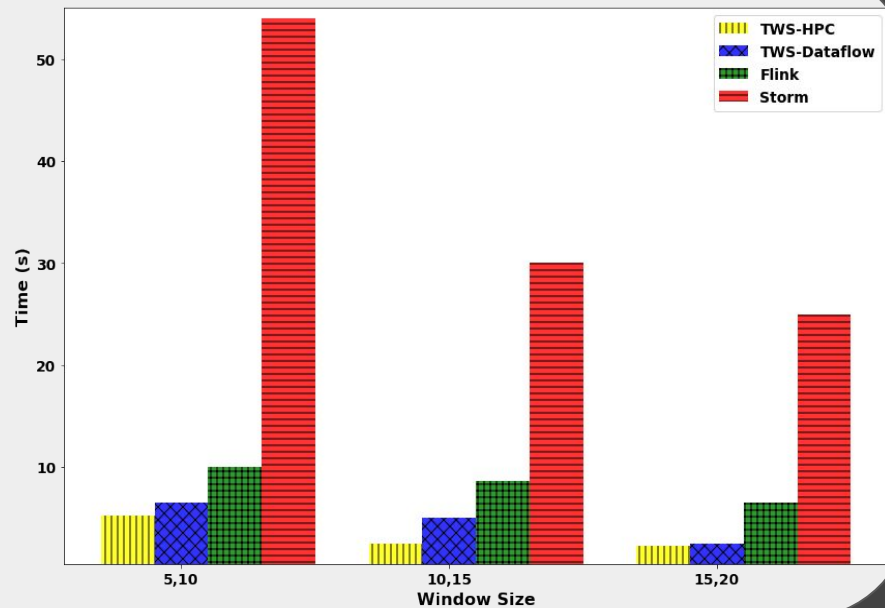
*5,10 refers to sliding length,window length.

Obtained after experimenting with different configs towards optimum results obtained in batch mode.

Streaming KMeans



Tumbling Windowing



Sliding Windowing

*5,10 refers to sliding length,window length.

Obtained after experimenting with different configs towards optimum results obtained in batch mode.

Conclusions and Future Work

- Windowing APIs are vital for designing iterative streaming applications.
- High performance computing model can be adopted in Big Data frameworks to provide better performance for streaming applications.
- *Experimenting with a larger data stream (minimum of 1 Million of more data points per a job)*
- *Structured data streaming with stream discretization.*
- *Expanding experiment configurations for testing window config sensitivity on algorithm convergence.*
- *Scaling for a bigger experiment setting (1024+ cores)*
- *Extending experiments for more machine learning algorithms.*

Thank you

- NSF
- Future Systems Team @ IU (Allan Streib et. al)
- Digital Science Center