

Data sketching for transient detection in Tomo-e light curve data

Phungtua-eng Thanapol*, 山本 裕介, 鈴木 健太, 西川 侑志, 山本 泰生

静岡大学・情報学部

木曾シュミットシンポジウム2022

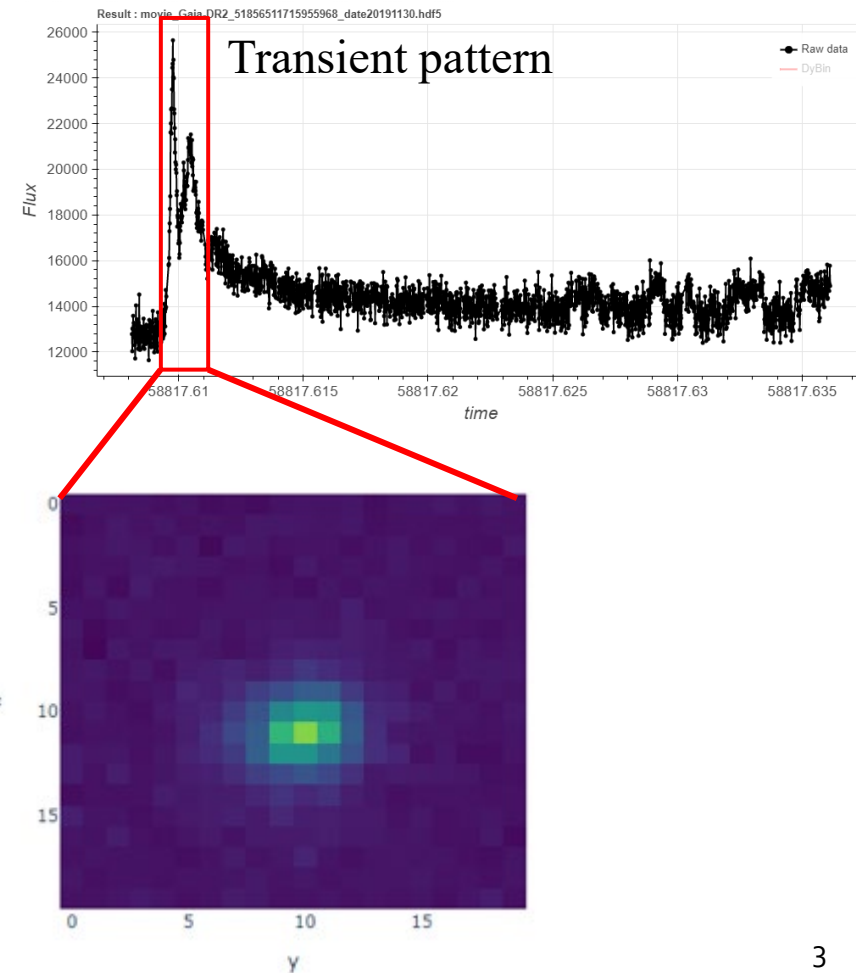
2022.07.06

Outline

- Introduction
- Problem definition
- Dynamic binning
- Experiments on the light curves
- Conclusion

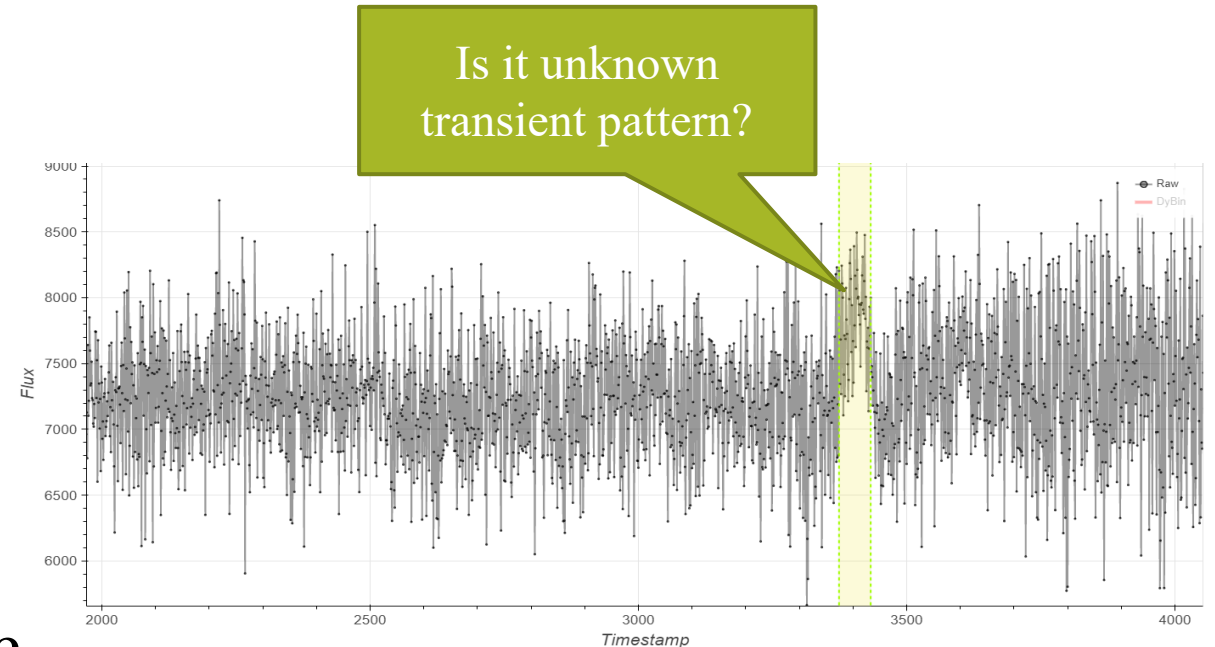
What are unknown transient patterns ?

- Transient: short-term phenomenon suddenly and intensively happens in an astronomical object.
- Known (Unknown) transient: transient that has been observed in astronomy (ex: Flare, Blazar, etc.).
- Transient pattern: time-series signal appearing in light curve data, associated with the transient.



Light curves (LCs)

- LCs is a graph of light intensity of the astronomical object within the region over a certain time period.
- LCs usually contain unwanted external factors:
 - Atmospheric turbulence
 - Measurement error from hardware.
- They may lead to erroneous analysis conclusions.
- The general technique to avoid this issue is *data sketching*



Example of LCs

Data sketching

Q: What is data sketching?

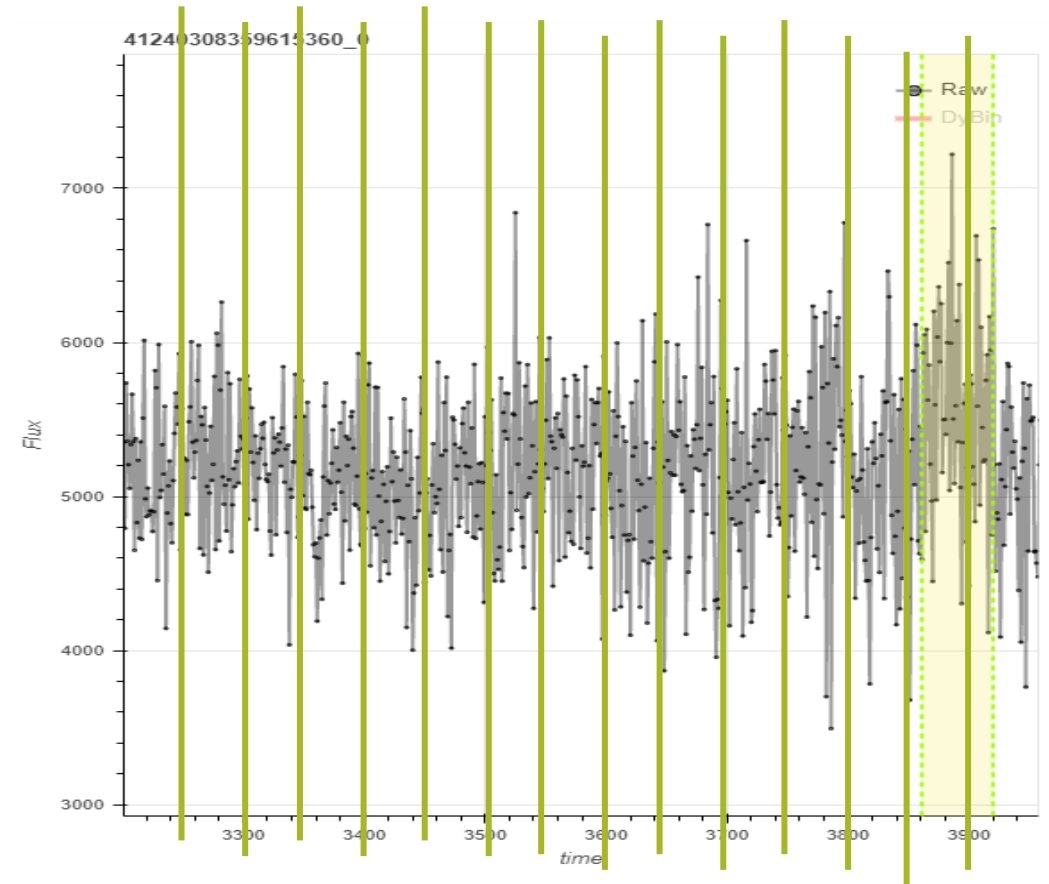
- Summarizes or approximates data into a new representation of fewer lengths.

Q: What is a new representation?

- Presents the relevant feature of subsequences of LCs into approximate value, divided into bins.
- We call it data binning.

How to compute data binning

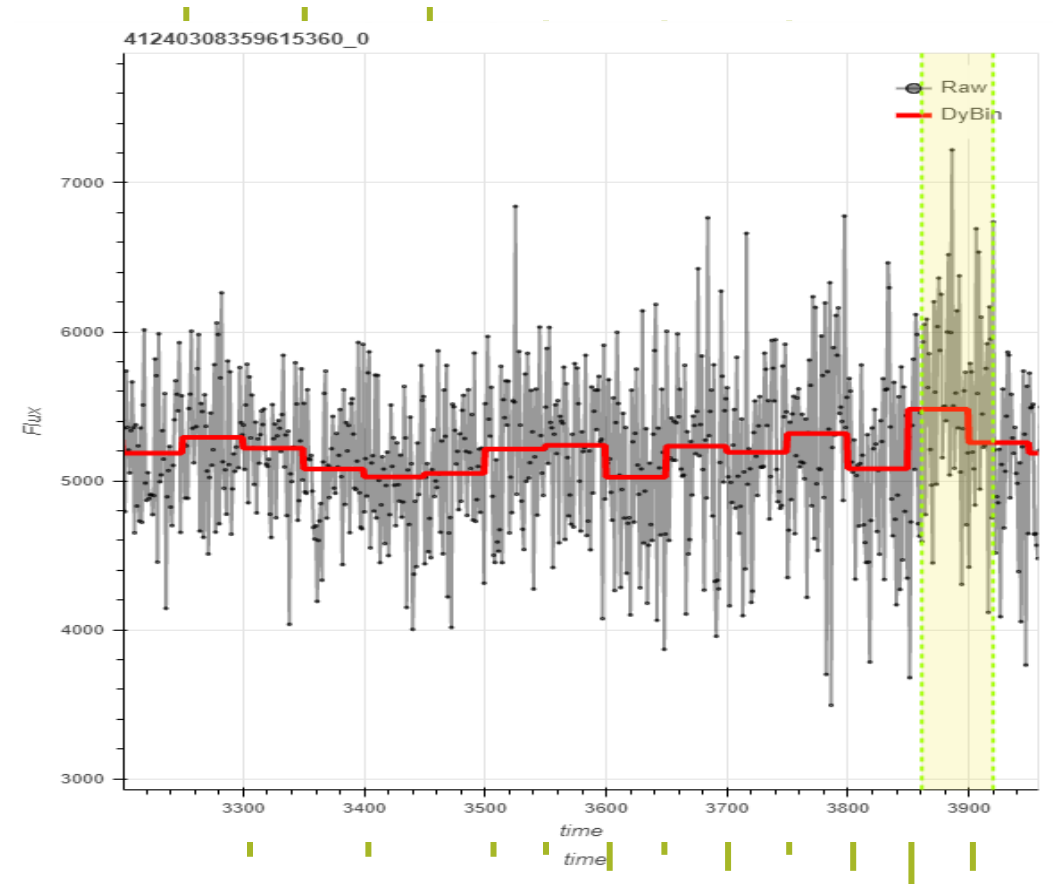
- Input : the LCs $\{X_1, X_2, X_3, \dots, X_t\}$, where t is timestamp.
- Step 1: the dimensionality reduction is divided into equal-sized bins.



LCs : $\{X_1, X_2, X_3, \dots, X_t\}$

How to compute data binning

- Step 2: computes mean value within bins for representation.
- Output: $W_{1,t} = \{\text{Bin}_1, \text{Bin}_2, \dots, \text{Bin}_w\}$, where w is the latest bin and $w \leq t$.



$$W_{1,t} = \{\text{Bin}_1, \text{Bin}_2, \dots, \text{Bin}_{16}\}$$

Data binning

Q: What does data binning provide?

- Extracts relevant characteristics of transient patterns
- Reduces noise while preserving the characteristics of the original data

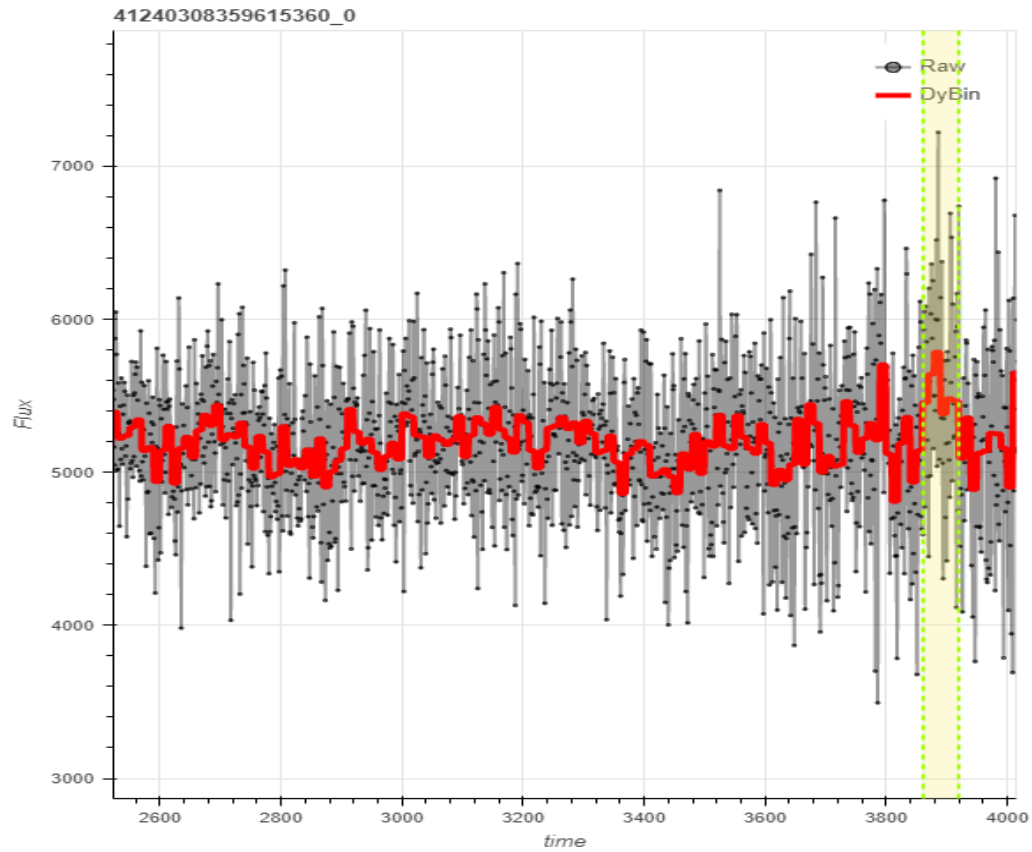
Q: What is required is an input parameter that produces a new representation by data binning?

- Bin sizes and window sizes for division into bins by compression.

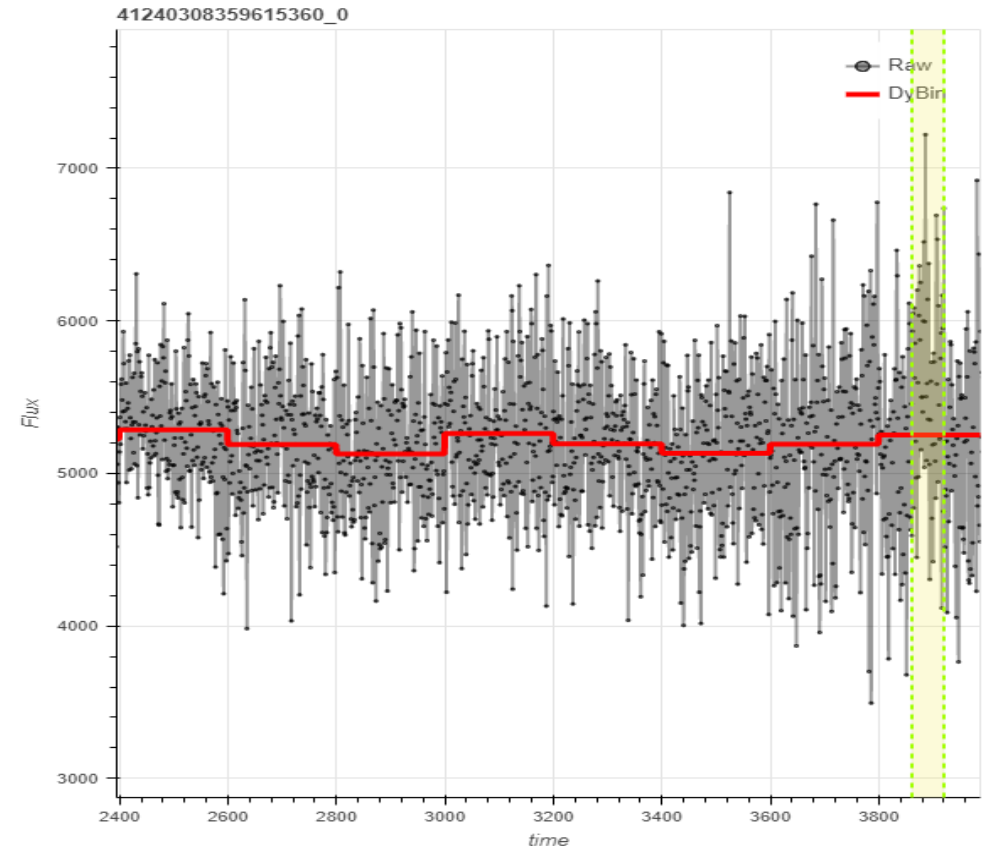
Q: Can we simply choose an arbitrary bin sizes and window sizes with ad hoc input?

- **The answer is NO.**

Arbitrary bin sizes



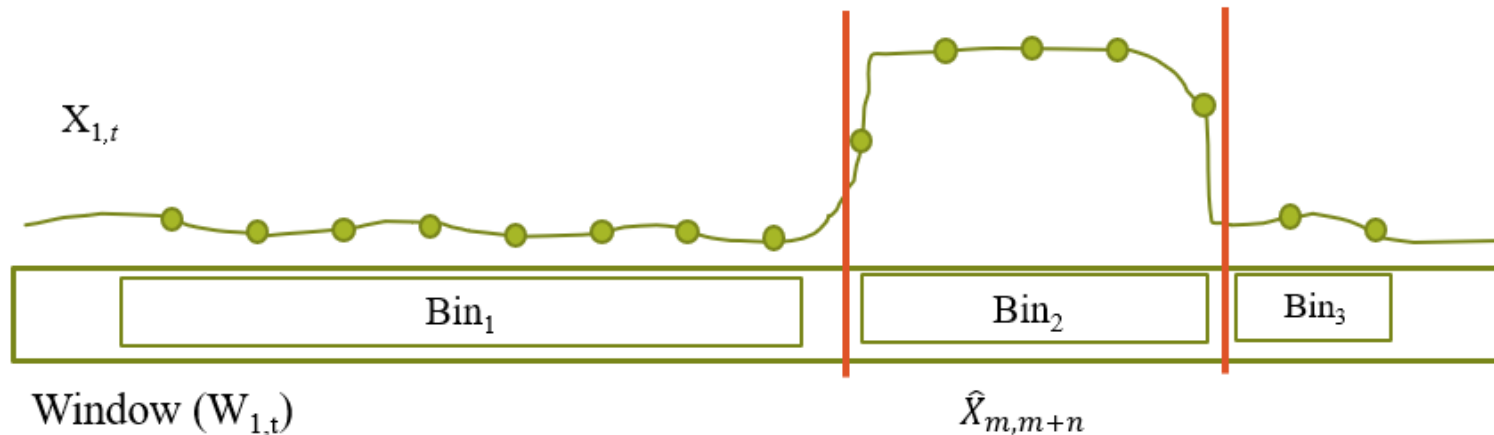
Sub-sequence length is too small



Sub-sequence length is too large

Preliminaries

- Given a time series data of (LCs) $X_{1,t} = \{x_1, x_2, x_3, \dots, x_t\}$, is a sequence of measurements of light intensity within the time interval $[1, t]$.
- Given a $\hat{X}_{m,m+n}$ is unknown transient pattern that is sub-sequence of LCs, where m and $m+n$ are the starting point and endpoint of the transient pattern.
- Output is to split the LCs into w continuous bins $W_{1,t} = \{\text{Bin}_1, \text{Bin}_2, \dots, \text{Bin}_w\}$, and are collected into a window.



Problem definitions (Cont.)

Problem definitions

- We find the periods of stability and anomaly behavior to divide LCs into sequences of bins.
- Bins in the window are dynamic bin sizes and represent each period's features.
- It aims to produce the representation of **LCs** using only w bins, while w is minimized.
- It aims to produce the representation of $\hat{X}_{m,m+n}$ using only 1 bin, and that bin corresponds to features of $\hat{X}_{m,m+n}$.

Dynamic binning

- The dynamic binning to adjust the bin size was published in [1].
- Auto-adjusts bin sizes.
- Merges some bins that are likely to be similar to each other.
- Our proposed method includes two statistical hypothesis tests.
 1. T-test for equal means of two bins.
 2. F-test for equal variances of two bins.

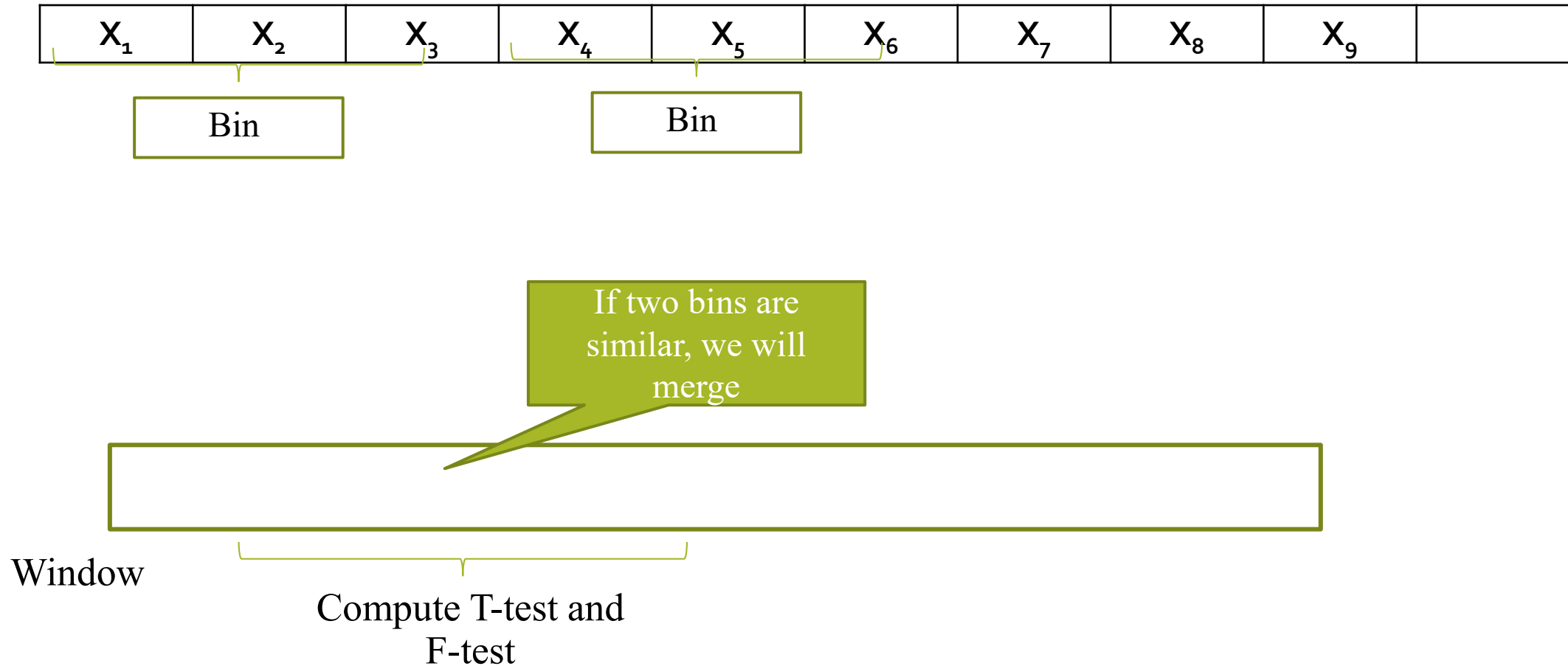
Algorithm 1: *DynamicBin*

input: W : window, α : significance level

```
1 Compute T-test for equal means of  $Bin_l$  and  $Bin_{l-1}$ 
2 if  $|T_{test}| \leq T_{1-\frac{\alpha}{2}}$  then
3   Compute F-test for equal maximum possible
   variances of  $Bin_l$  and  $Bin_{l-1}$ 
4   if  $F_{test} \leq F_{1-\frac{\alpha}{2}}$  then
5     Merge( $Bin_l, Bin_{l-1}$ )
6   else
7     if  $W$  is full then
8       Search the bin  $Bin_n$  whose T-test value is
       minimum in  $W$ 
9       Merge( $Bin_n, Bin_{n-1}$ ).
10    end
11  end
12 end
```

Algorithm 1 : Dynamic binning

How to compute dynamic binning



T-test and F-test

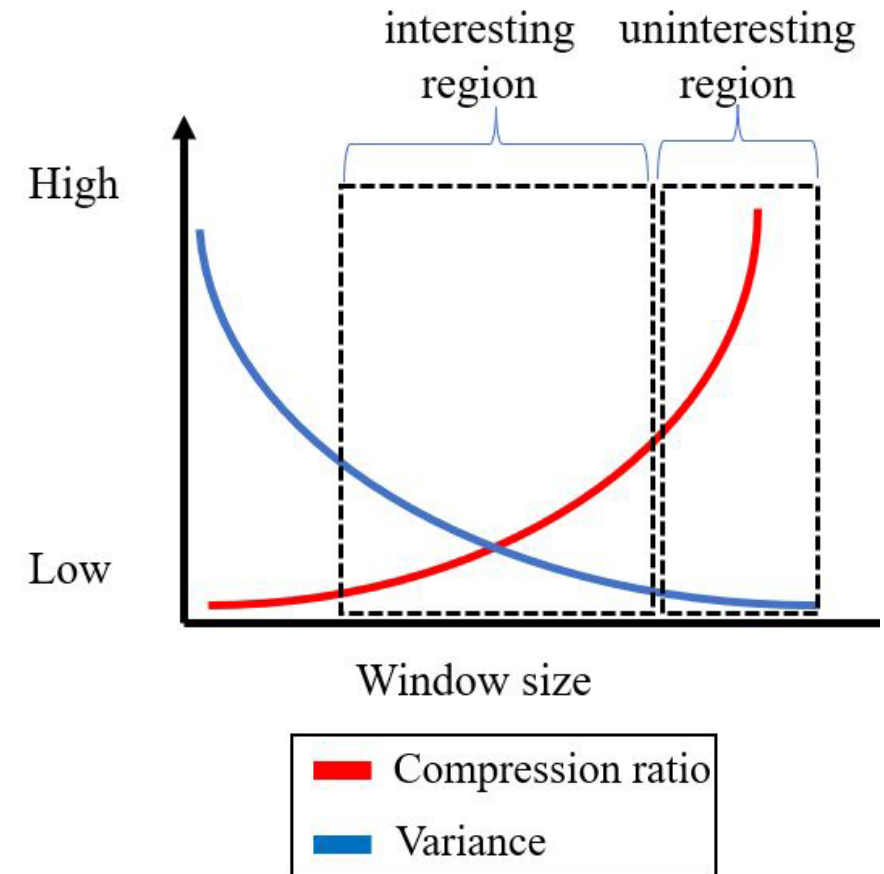
$$T - \text{test} = \frac{|\mu_i - \mu_{i+1}|}{\sqrt{\frac{\sigma_i^2}{n_i} + \frac{\sigma_{i+1}^2}{n_{i+1}}}}$$

$$F - \text{test} = \frac{\sigma_i^2}{\sigma_{i+1}^2}$$

- i and $i+1$ are index of the i -th and $(i+1)$ -th Bin in the window.
- n is the number of instances that correspond to sub-sequence of LCs.
- μ is the mean, σ^2 is the variance of bin.
- The greater the value of the T-test, the greater the evidence two bins are not similar.
- The greater the value of the F-test, means two bins are significantly different distribution.

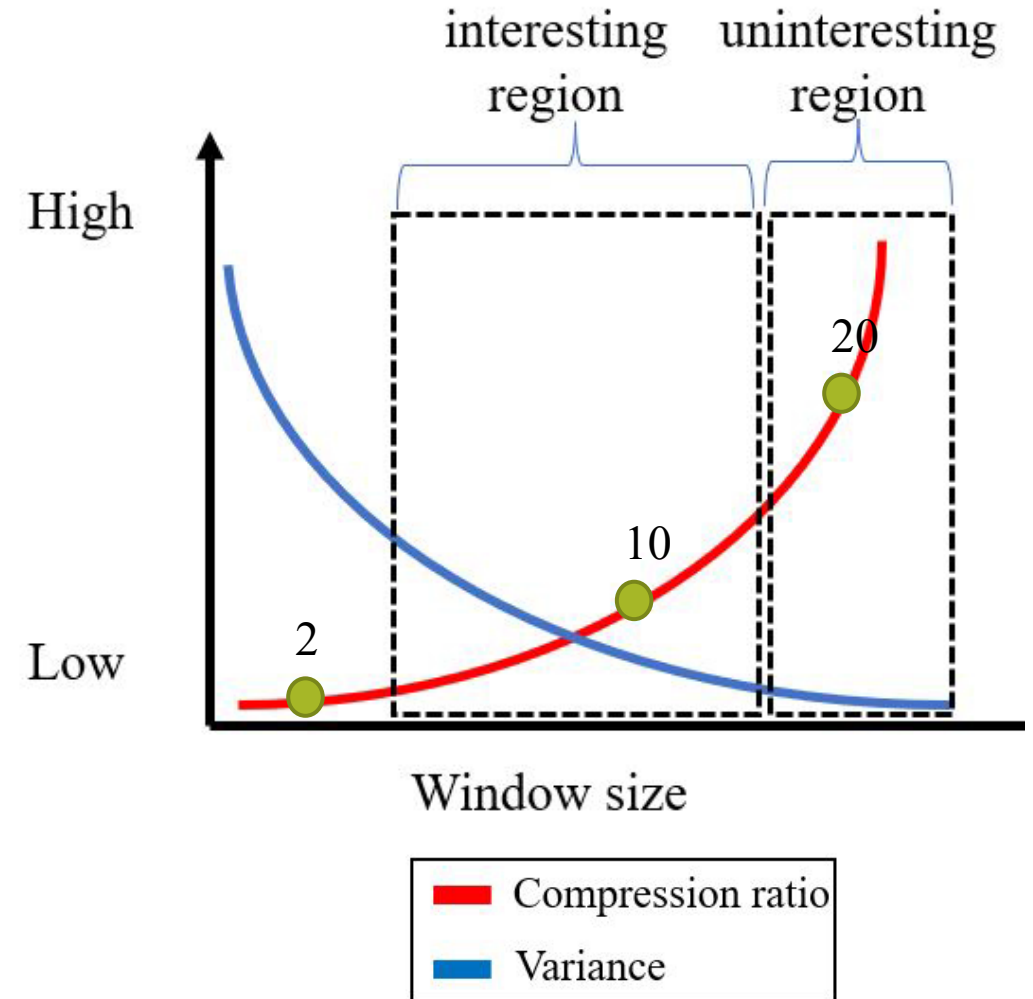
Compression ratio

- The compression ratio is the ratio between the length of original data and length of the output by compression.
- The variance and compression ratio do have opposite behaviors.
- Compression ratio is an essential measurement for the judgment of window size.



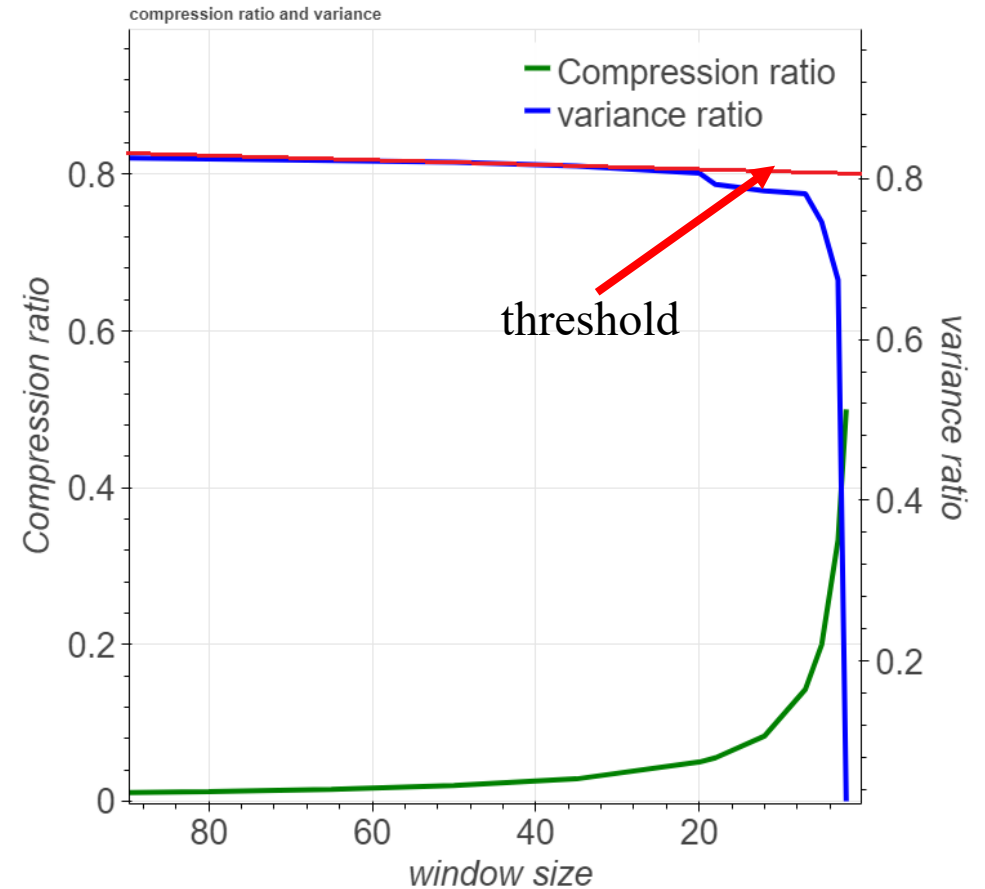
Compression ratio (cont.)

- For example:
 - Input data : $X_{(1,100)} = \{x_1, x_2, \dots, x_{100}\}$; $u = 100$ instances.
 - Output is 50 bins \rightarrow Compression ratio = 2
 - Output is 10 bins \rightarrow Compression ratio = 10
 - Output is 2 bins \rightarrow Compression ratio = 20
- A small value of compression ratio means a lot of noise in the input.
- A large value of compression ratio means a low redundancies.



Adjusting window size

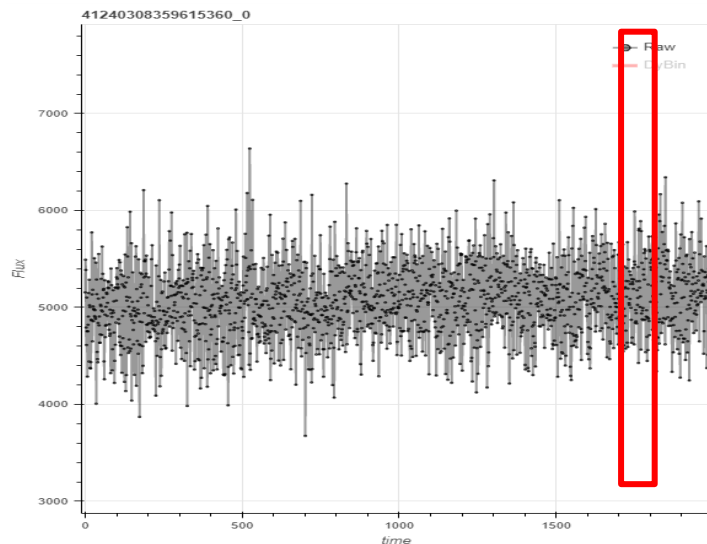
- Gives a threshold for considering window size [2].
- When the variance ratio far from the threshold is an uninteresting region.
- In this case, window size is to more than ten bins.



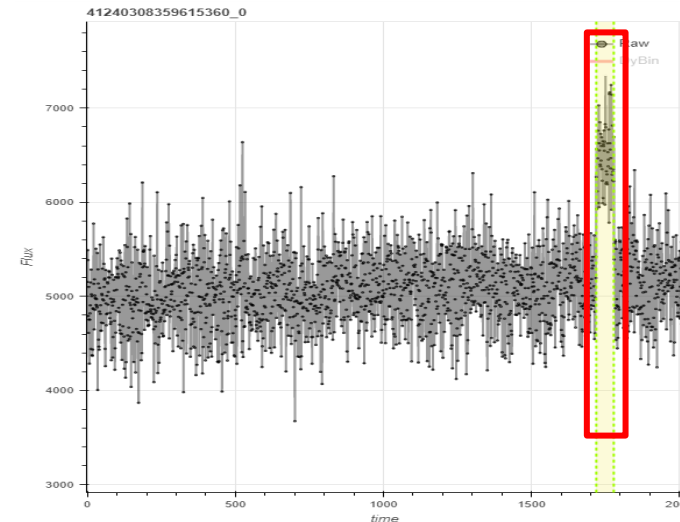
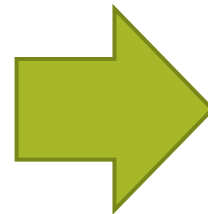
Trade-off between
variance ratio and compression ratio

Experiments on the LCs

- We inject **artifact one square transient pattern** for each file.
 - The duration varies from time ranging in [1mins, 2mins, 3mins, and 6mins].
 - The power of the transient pattern is [1σ , 3σ], where σ is the standard deviation of the original data for each file.



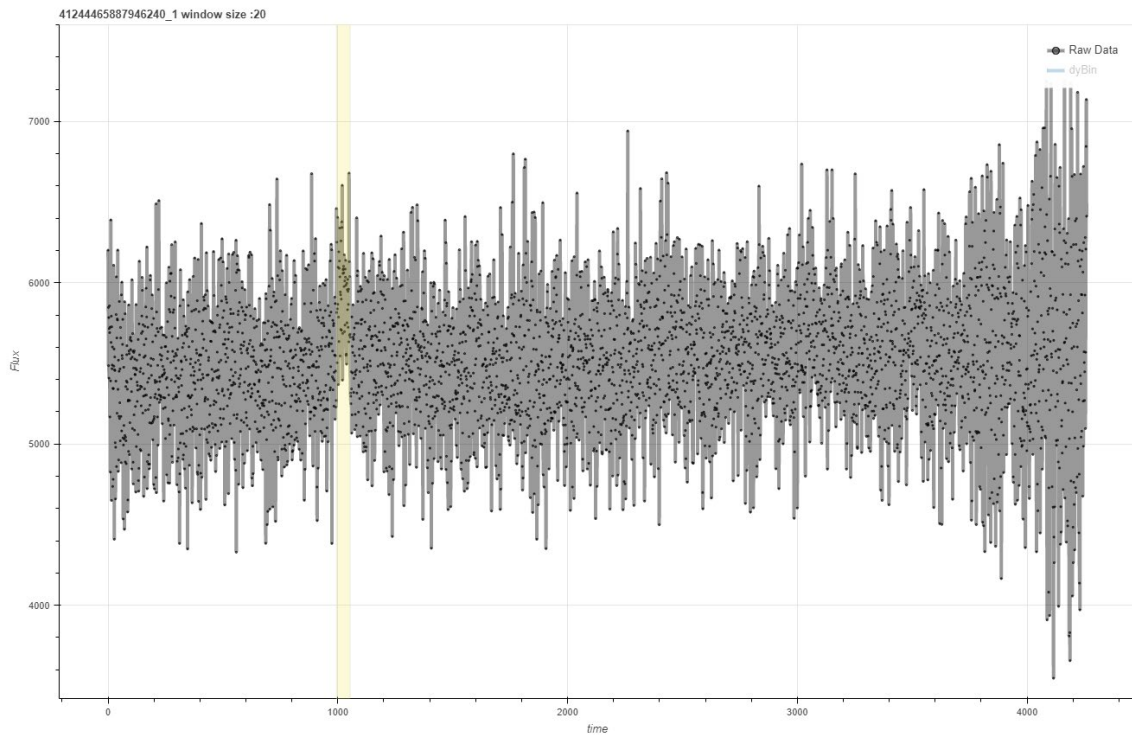
Original LCs



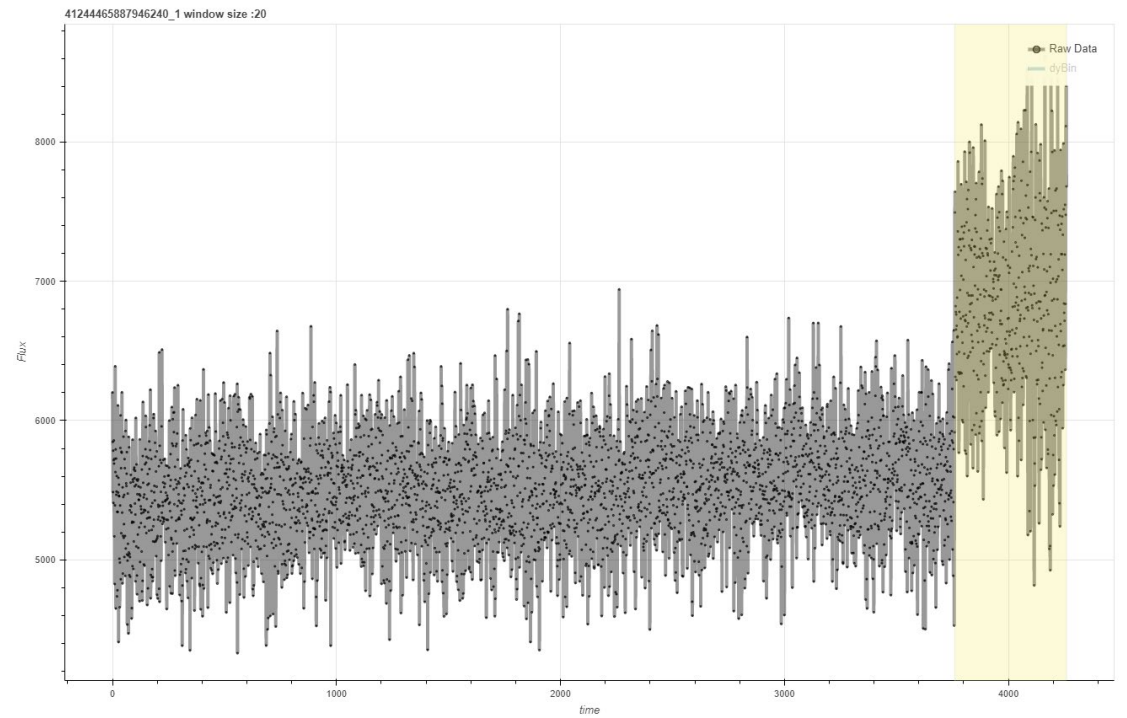
After transient pattern injection

Power : 3σ
Duration : 1mins

Experiments on the LCs (Cont.)



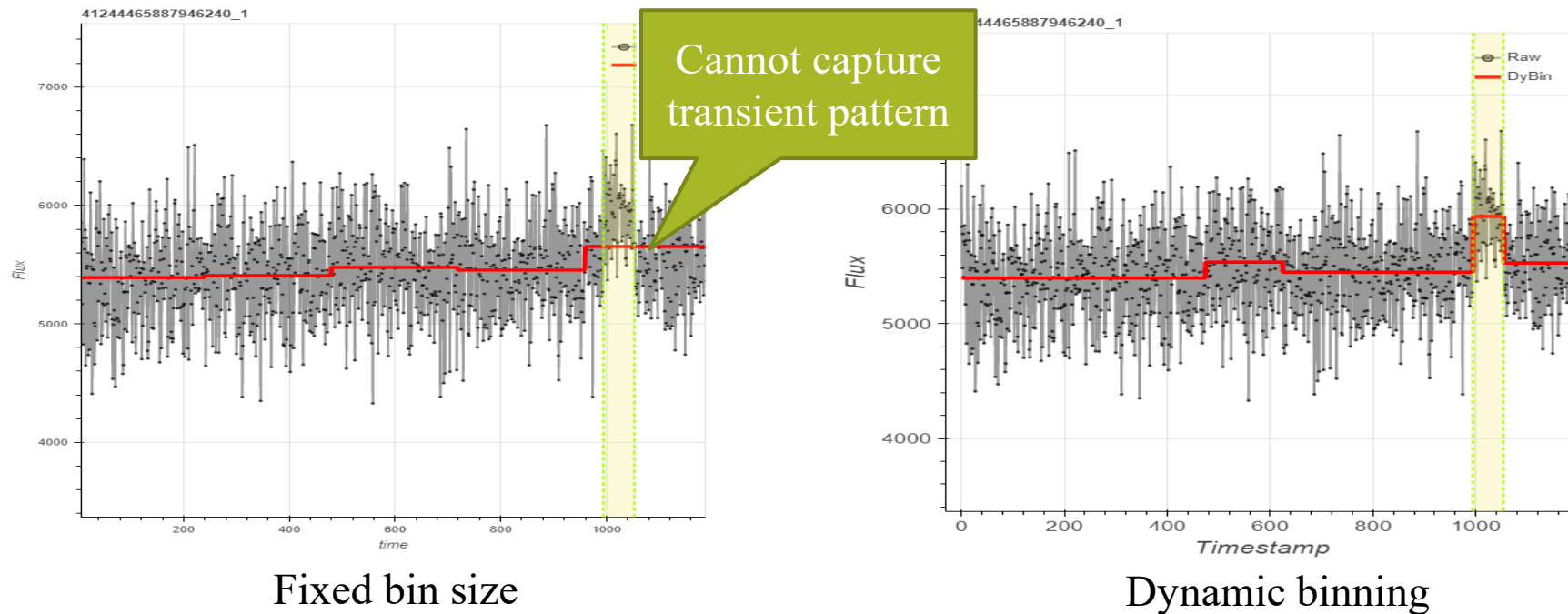
Lowest and shortest duration



Highest power and longest duration

Result of sketching

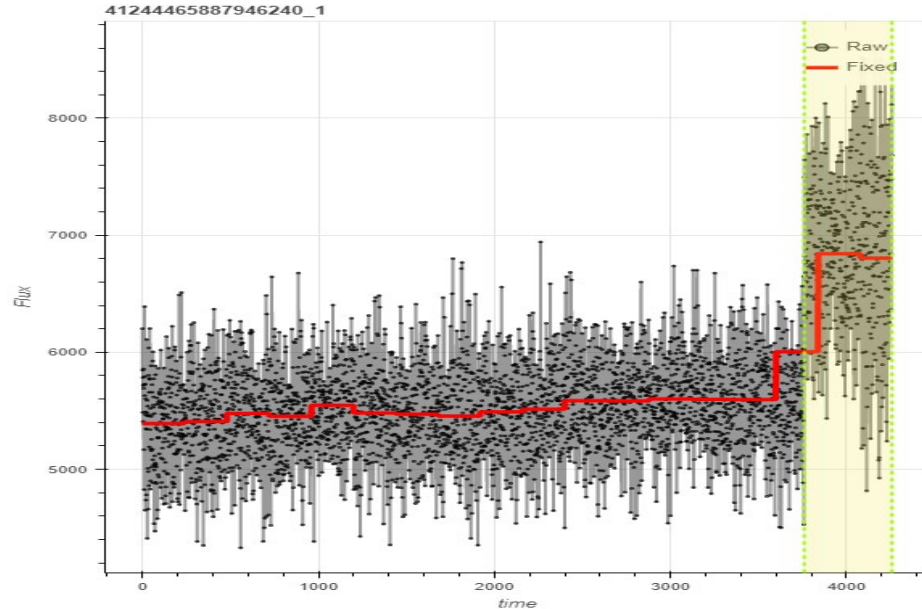
- We compare our dynamic binning against fixed bin size.



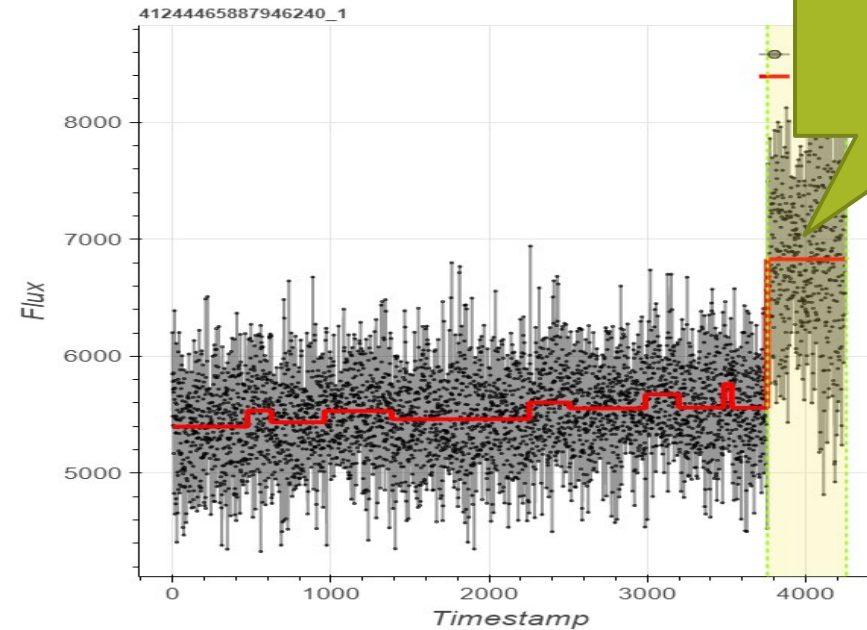
Comparison results dybin and fixed bin size with lowest and shortest duration

Result of sketching (Cont.)

- We compare our dynamic binning against fixed bin size.



Fixed bin size

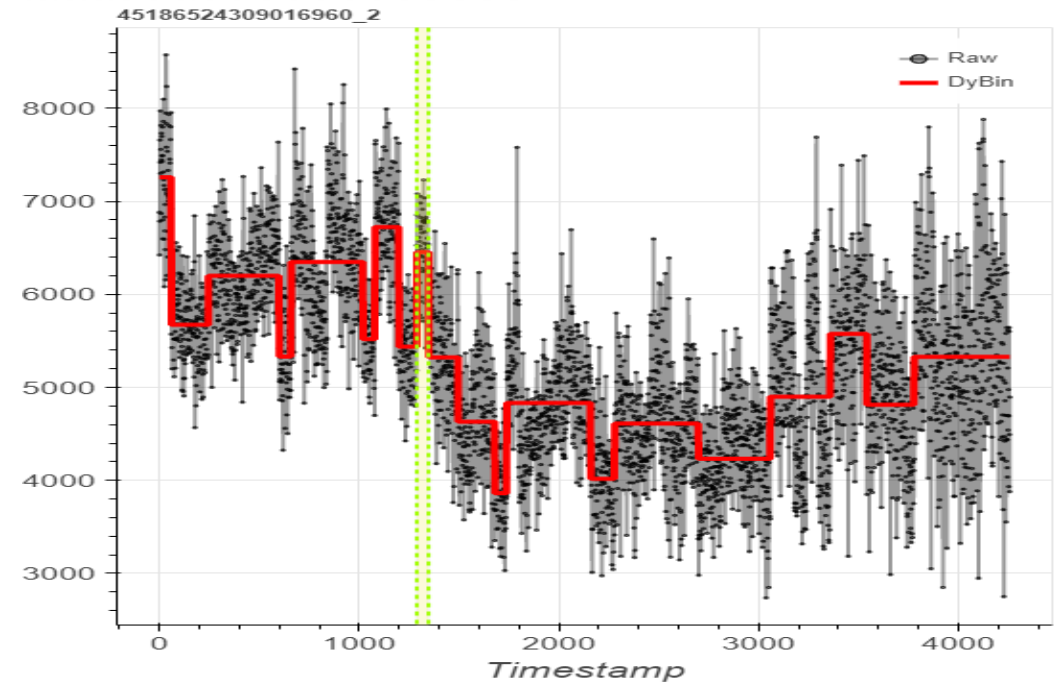
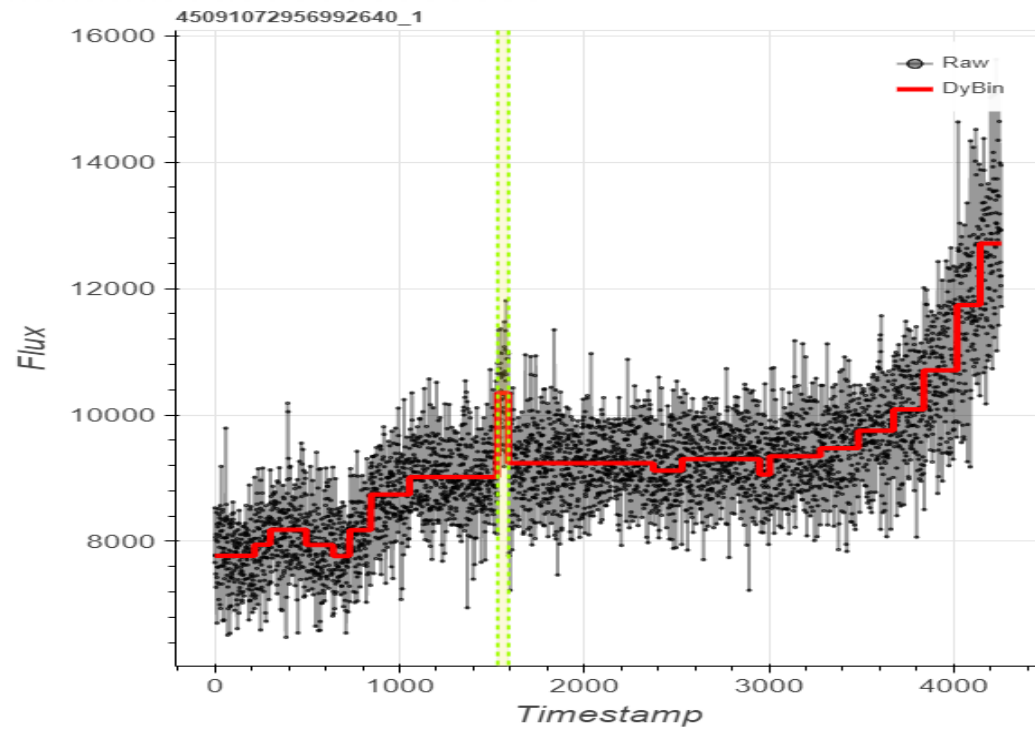


Dynamic binning

1 bin for
 $\hat{X}_{m,m+n}$

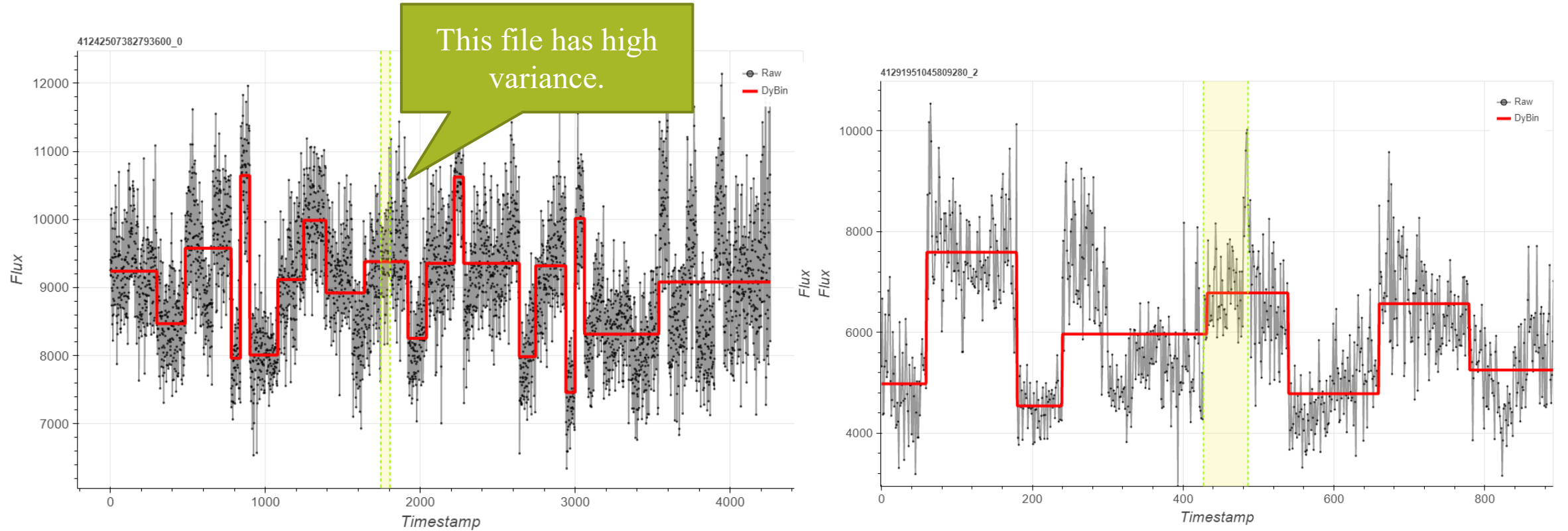
Comparison results dybin and fixed bin size with highest power and longest duration

Result of sketching (Cont.)



The results of dynamic binning with various scenarios

Result of sketching (Cont.)



Unsuccessful cases

Transient pattern detection

- We apply dynamic binning with detection methods.
- We evaluate dynamic binning against state-of-the-art methods for comparison.
 - SK-method with fixed bin size [3]
 - SK-method with dynamic binning
 - Matrix Profile algorithm (MP) [4]
 - Robust Random Cut Forest (RRCF)[5]
 - Kernel Density Estimation (KDE) [6]

[3] G. Shevlyakov and M. Kan, "Stream Data Preprocessing: Outlier Detection Based on the Chebyshev Inequality with Applications," 26th Conference of Open Innovations Association (FRUCT), 2020, pp. 402-407

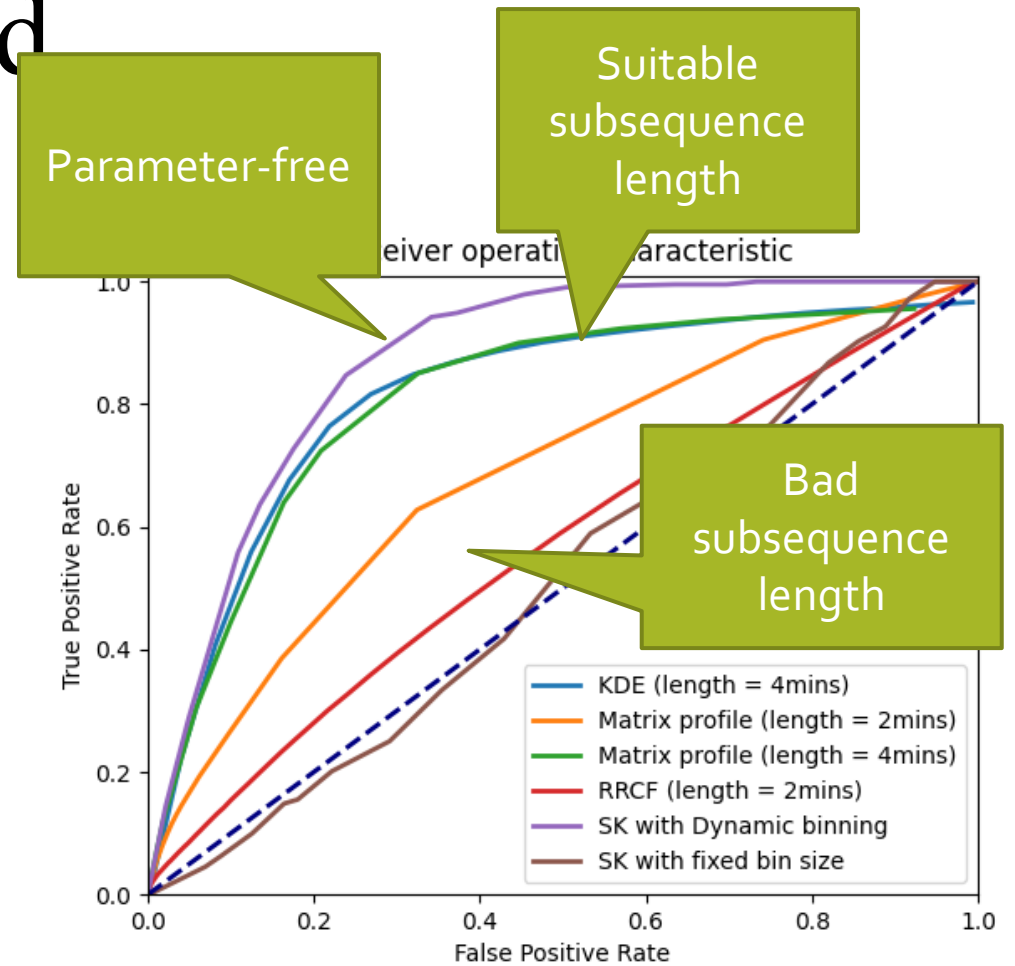
[4] C. M. Yeh et al., "Matrix Profile I: All Pairs Similarity Joins for Time Series: A Unifying View That Includes Motifs, Discords and Shapelets," IEEE 16th International Conference on Data Mining (ICDM), 2016, pp. 1317-1322

[5] S. Guha, N. Mishra, G. Roy, & O. Schrijvers, "Robust random cut forest-based anomaly detection on streams", 33rd International conference on machine learning, 2016, pp. 2712-2721.

[6] Latecki, L.J., Lazarevic, A., Pokrajac, D. (2007). Outlier Detection with Kernel Density Functions. In: Perner, P. (eds) Machine Learning and Data Mining in Pattern Recognition. MLDM 2007.

ROC of detection method

- We obtained good overall results of dynamic binning using the SK method to detect square transient patterns without put bin sizes and window size.
- MP, KDE, and RRCF based on fixed length.
- We suppose the good representation that may improve performance of RRCF, KDE ,and MP.



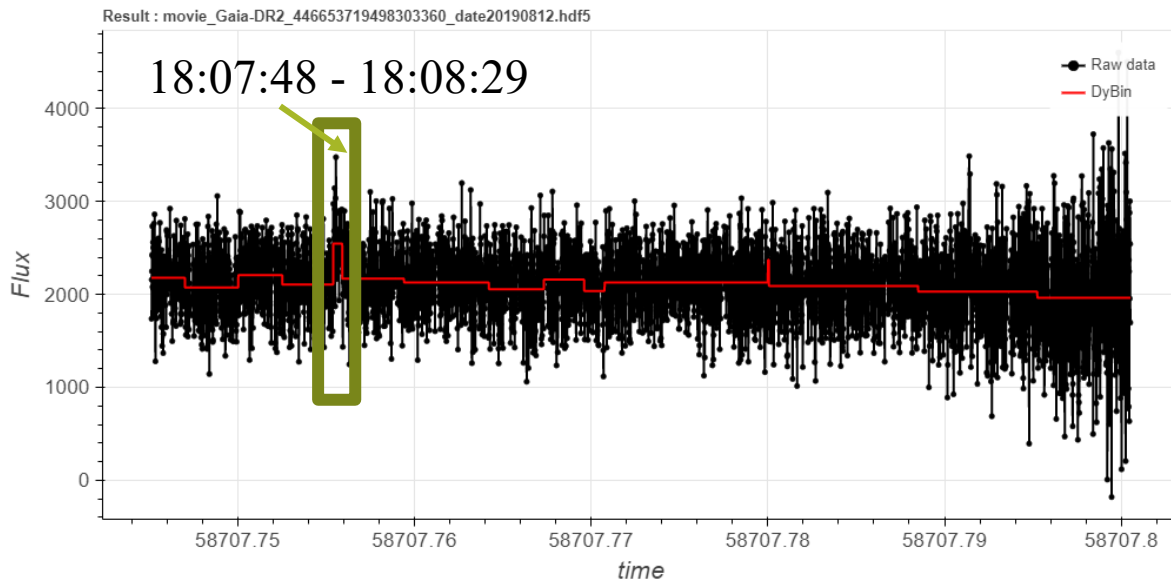
ROC curve of 5 methods

Evaluation with real flares

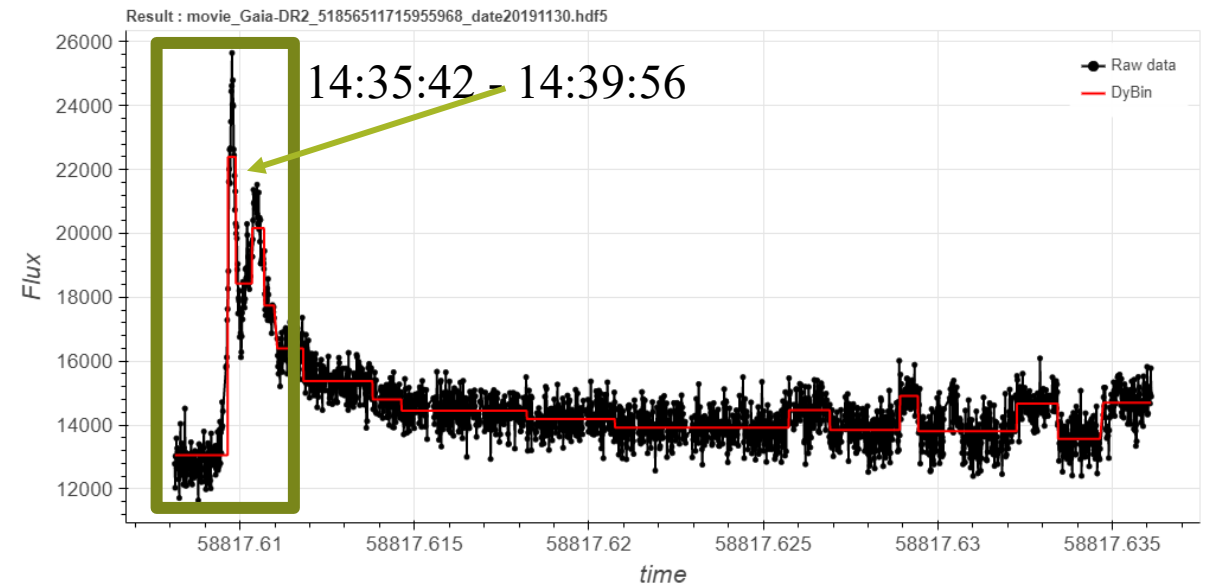
- We evaluate dynamic binning using SK-method with real dataset.
- This dataset contains 10,850 files and 18 files that are real flares.
 - The dataset provided is from Kashiwama-san and Aizawa-san.
- We have two part for evaluation.
 1. Results from our sketching with 18 flares.
 2. Result detection with dynamic binning using SK-method

Example of our sketching with real flares.

- Each bin in the window can represent each period's features.



Low-powered signal



High-powered signal

Result detection with dynamic binning using SK-method

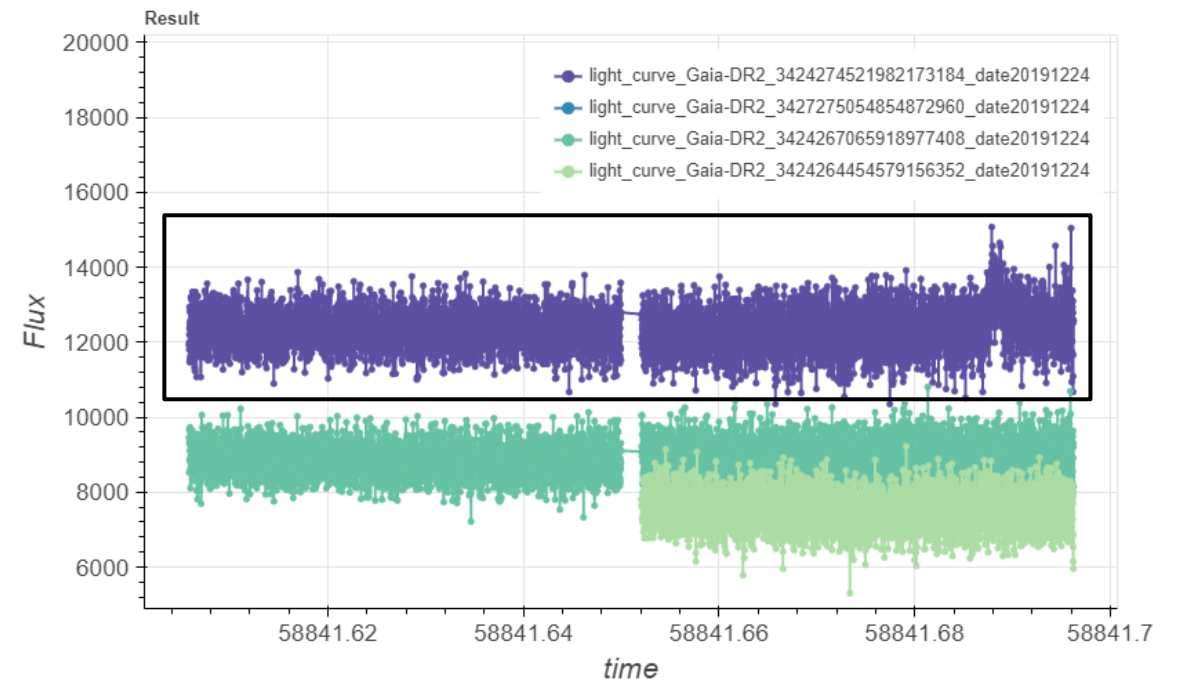
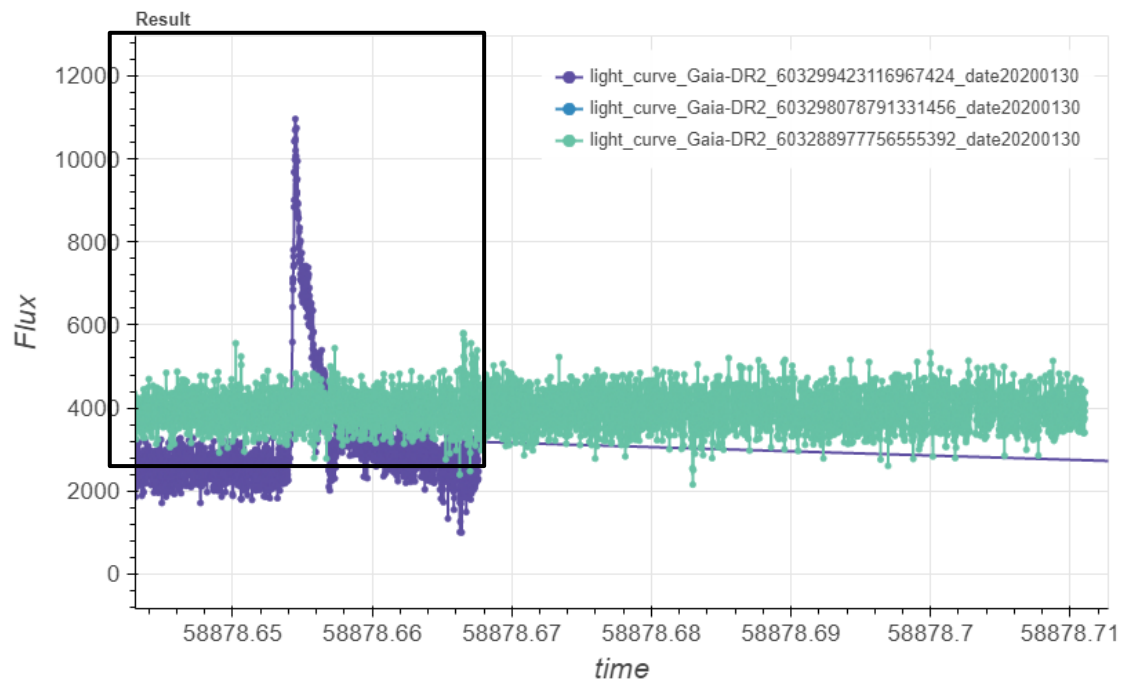
- Input : 10,850 files
- Output : Transient patterns list sorted by the metrics of SK-method.
- Our processes:
 - Sketches each file with **dynamic binning**.
 - Apply SK-method to detect candidate transient patterns [3].
 - Removes simultaneous events by Matrix profile [4].
 - List files whose high value is in Top-100.

[3] G. Shevlyakov and M. Kan, "Stream Data Preprocessing: Outlier Detection Based on the Chebyshev Inequality with Applications," 2020 26th Conference of Open Innovations Association (FRUCT), 2020, pp. 402-407

[4] C. M. Yeh et al., "Matrix Profile I: All Pairs Similarity Joins for Time Series: A Unifying View That Includes Motifs, Discords and Shapelets," 2016 IEEE 16th International Conference on Data Mining (ICDM), 2016, pp. 1317-1322

Detection result and discussion

- We found 8 flares from 18 in Top-100.



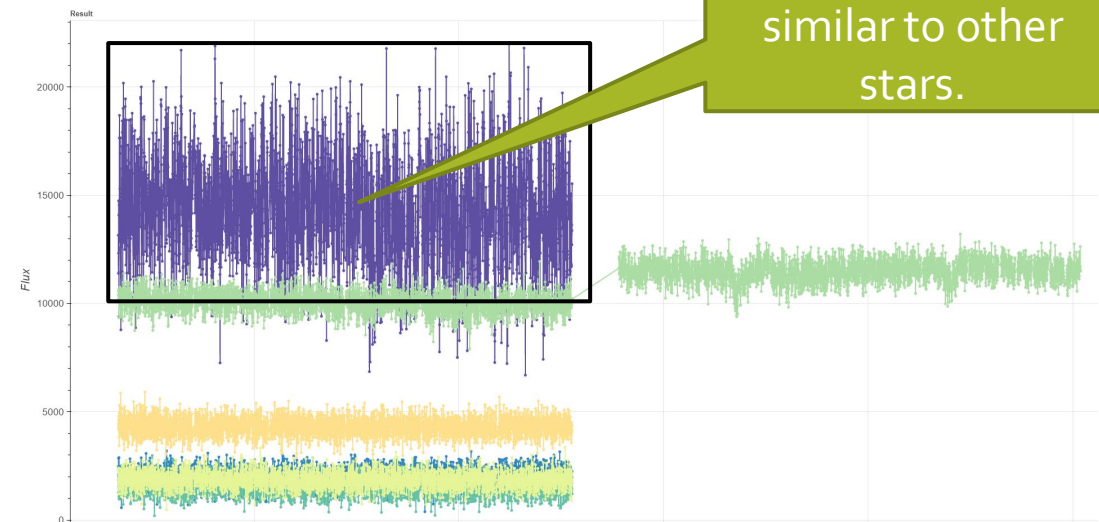
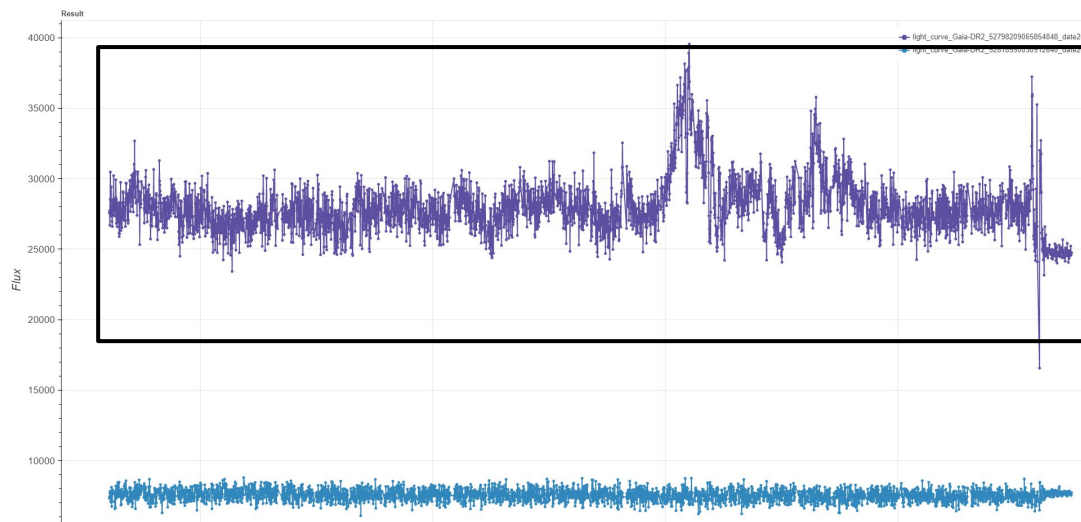
Example of real flares in Top-100

NO	File Name	Ranking
1	light curve Gaia-DR2 604942879467202816 date20200201	7
2	light curve Gaia-DR2 3398180156118506240 date20191224	16
3	light curve Gaia-DR2 603299423116967424 date20200130	26
4	light curve Gaia-DR2 3424781736143842432 date20191224	27
5	light curve Gaia-DR2 602712283908074752 date20200130	44
6	light curve Gaia-DR2 46623557923070848 date20191129	53
7	light curve Gaia-DR2 657563345604584064 date20200201	58
8	light curve Gaia-DR2 51856511715955968 date20191130	82
9	light curve Gaia-DR2 3424274521982173184 date20191224	106
10	light curve Gaia-DR2 631840133632923520 date20200201	249
11	light curve Gaia-DR2 3410076150374417408 date20191203	271
12	light curve Gaia-DR2 49407521363733632 date20191129	299
13	light curve Gaia-DR2 459557313084377600 date20190812	364
14	light curve Gaia-DR2 3311459749888116736 date20191129	406
15	light curve Gaia-DR2 2977433546209292672 date20191106	528
16	light curve Gaia-DR2 612157157509901568 date20200201	538
17	light curve Gaia-DR2 603188200643885696 date20200124	594
18	light curve Gaia-DR2 446653719498303360 date20190812	<u>598</u>

We can find all of them
in Top-600.

Detection result and discussion (Cont.)

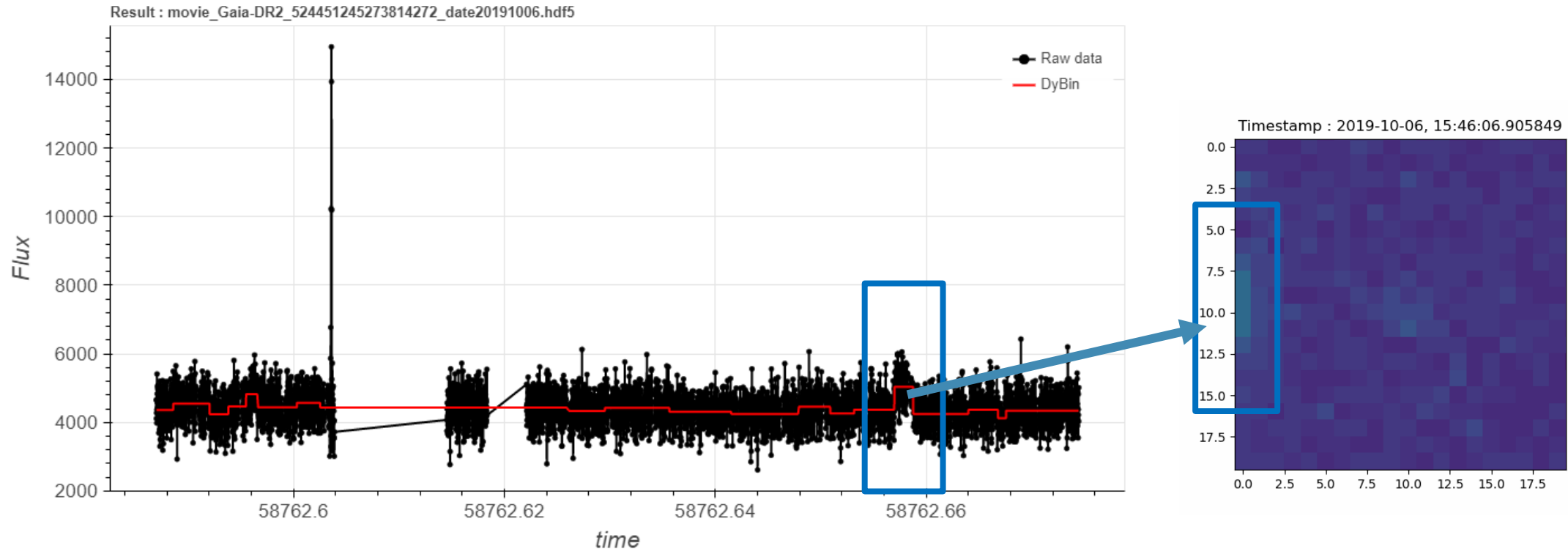
- Current problem is simultaneous event detection part.
 - Some files are unique when compare with near-star files.
 - We plan to evaluate with variety aperture sizes.



Files are not real flares but are unique behavior.

Detection result and discussion (Cont.)

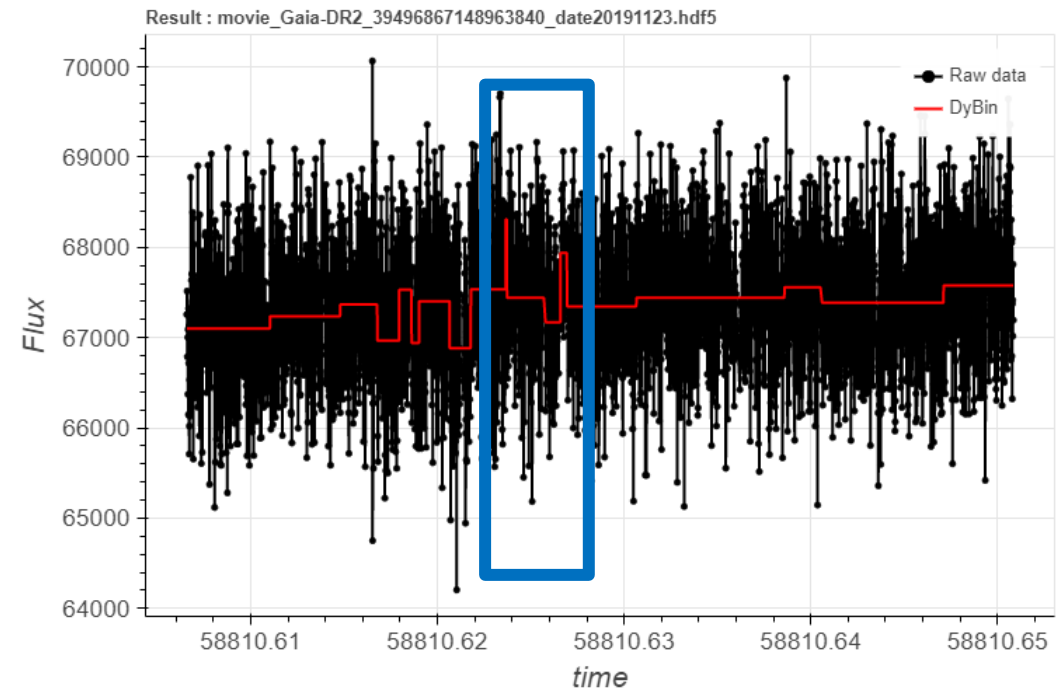
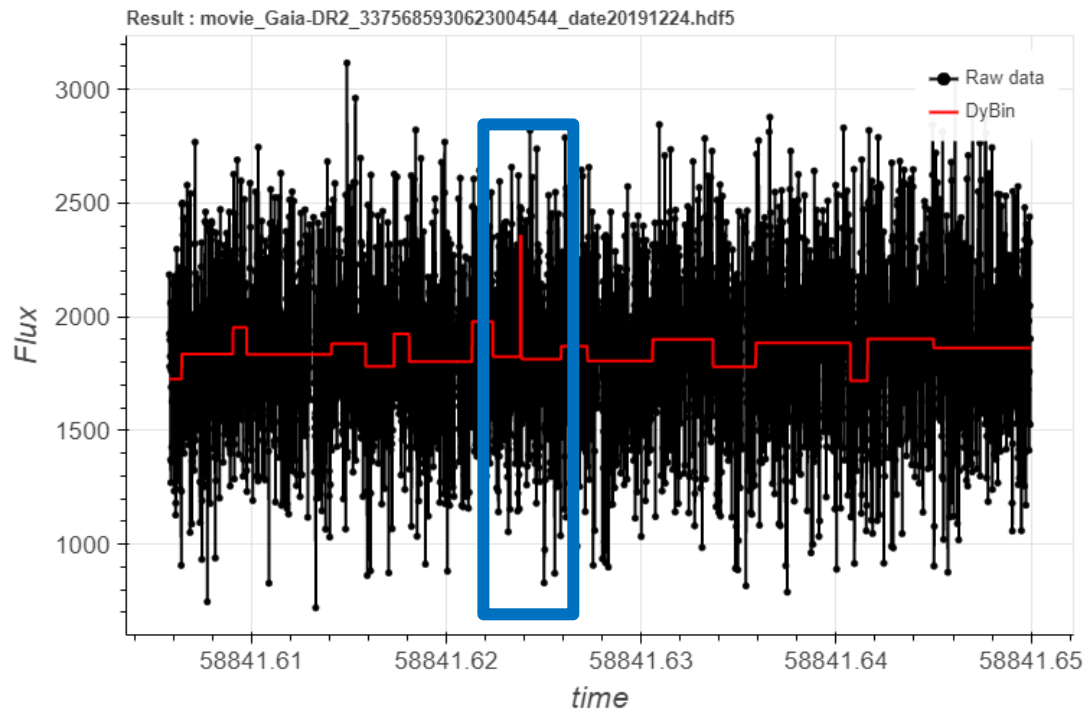
- We plan to improve simultaneous events removing for this issue.



Example of false positive

Detection result and discussion (Cont.)

- The second problem is that the bin is too small and cannot be merged with neighboring bins.
- Our algorithm starts with a small bin and merges into a large bin (bottom-up approach).
 - Our further work is top-down approach and hybrid approach.



Conclusion

- Window sizes and bin sizes can freely parameter without cornering significant differences results.
- We can extract relevant characteristics of transient patterns.
- **Further plan**
 - Investigation for applying dynamic binning with other detection methods.
 - Is there a way to address simultaneous event detection?
 - Decreasing the computational cost for the minimal window size searching.

Conclusion (Cont.)

- We are currently developing a web application using this algorithm.

The screenshot shows a web application interface for light curve analysis. On the left is a 'Menu' with options: Search LC files, Search LC and nearby star files, Upload LC files, Upload movie file, Search LC & movie file (highlighted), Aperture Comparison, SAX Comparison, and MJD convert. The main area is split into an 'Input Form' and a 'Result' panel. The 'Input Form' has fields for 'File name' (movie_Gaia-DR2_10092043569990400_date20191106.hdf5), 'Window size' (15), and 'Aperture radius' (7). The 'Result' panel shows a 'Movie result' plot of Flux vs. time, with 'Raw data' as a noisy black line and 'DyBin' as a smoother red line. To the right is a 2D heatmap of the star's position (x, y) over time, with a color scale from 0 to 8000. A green arrow points from the input form to the result panel.

User interface in our web application