# Rediscover the known Universe with NASA datasets

Horacio Gonzalez

`@LostInBrittany`
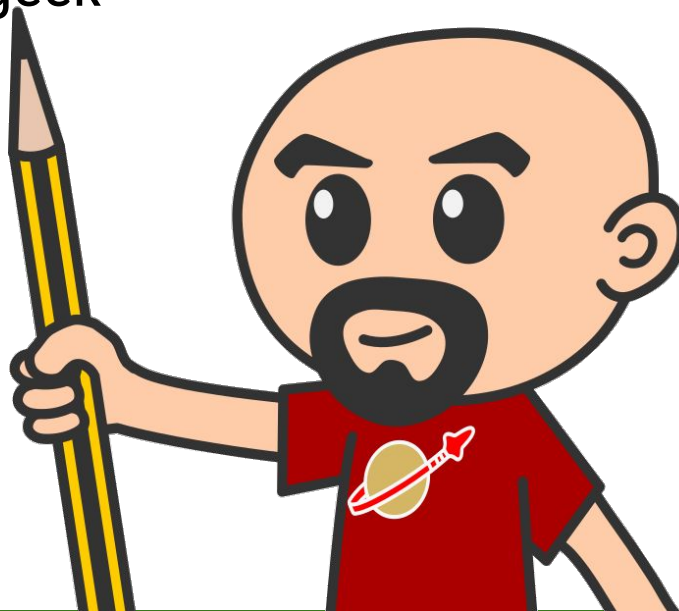
# Horacio Gonzalez

Code d'Armor

@LostInBrittany

Spaniard lost in Brittany, developer, dreamer and all-around geek

OVH
Team DevRel

Breizh C@mp
Mix de technologies

la cantine BREST
Finist Devs
GDG

WARP 10
MEETUP

Google Developers
Experts
2015
Web Technologies
GDE

# HelloExoWorld

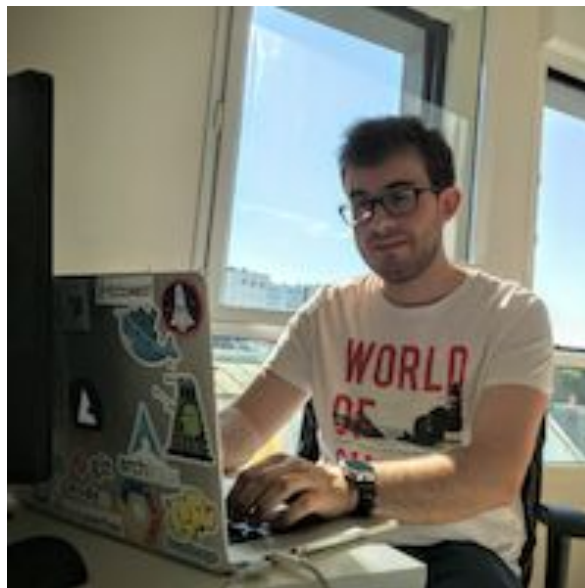Looking for exoplanets in NASA datasets

# HelloExoWorld

## Once upon a time…
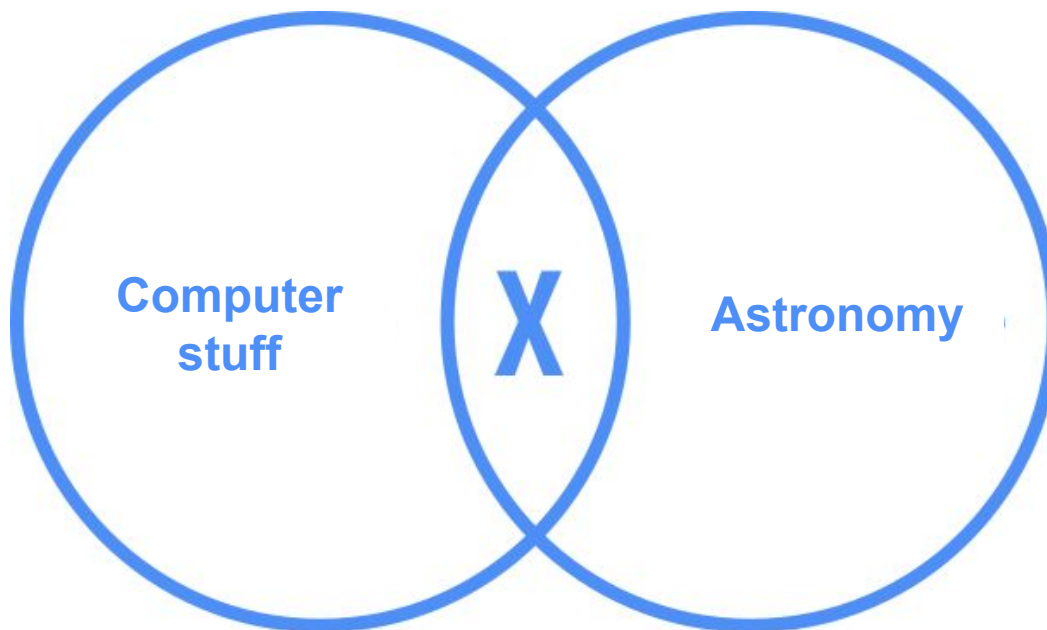
# An amateur astronomer



## Pierre Zemb, DevOps OVH

# What not to do if you love astronomy



## Live in Brest

# **Looking for solutions**

**Computer stuff**

**X**

**Astronomy**

Mixing passions

# Google is your friend...



## Let's find a project

# Exoplanets?



Planets orbiting stars far away

# How do we find them?



The transit method seems the best

# The transit method



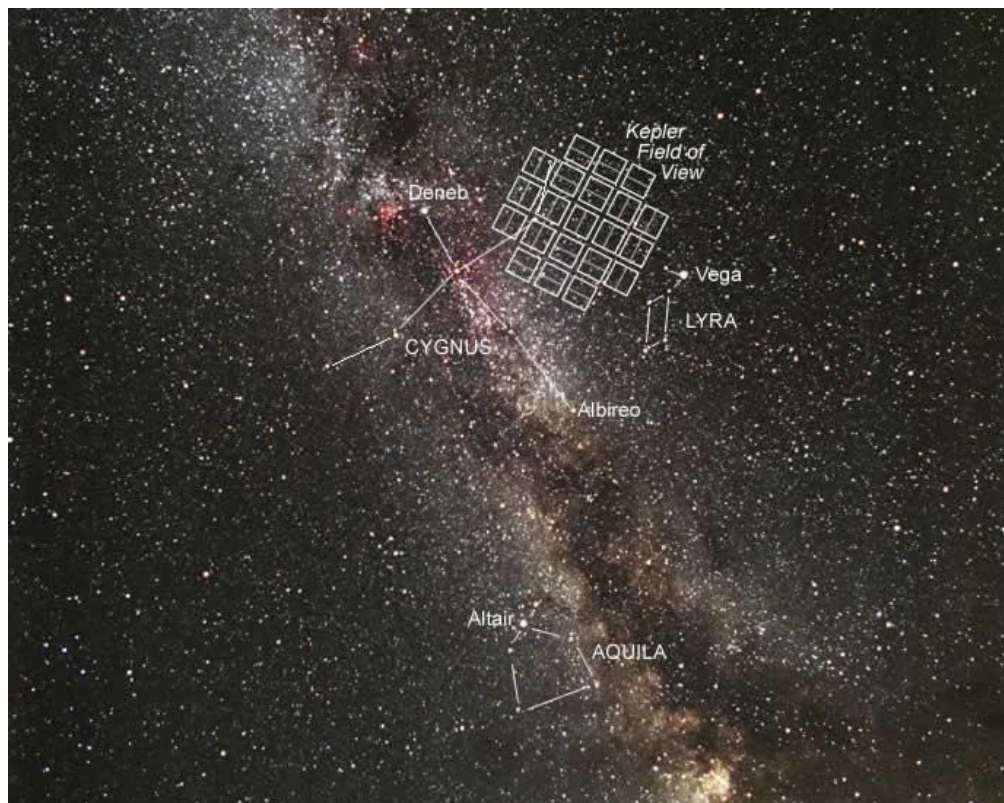Apparent Brightness of Star

# How do we look for transits?



Image credits : NASA

## Kepler
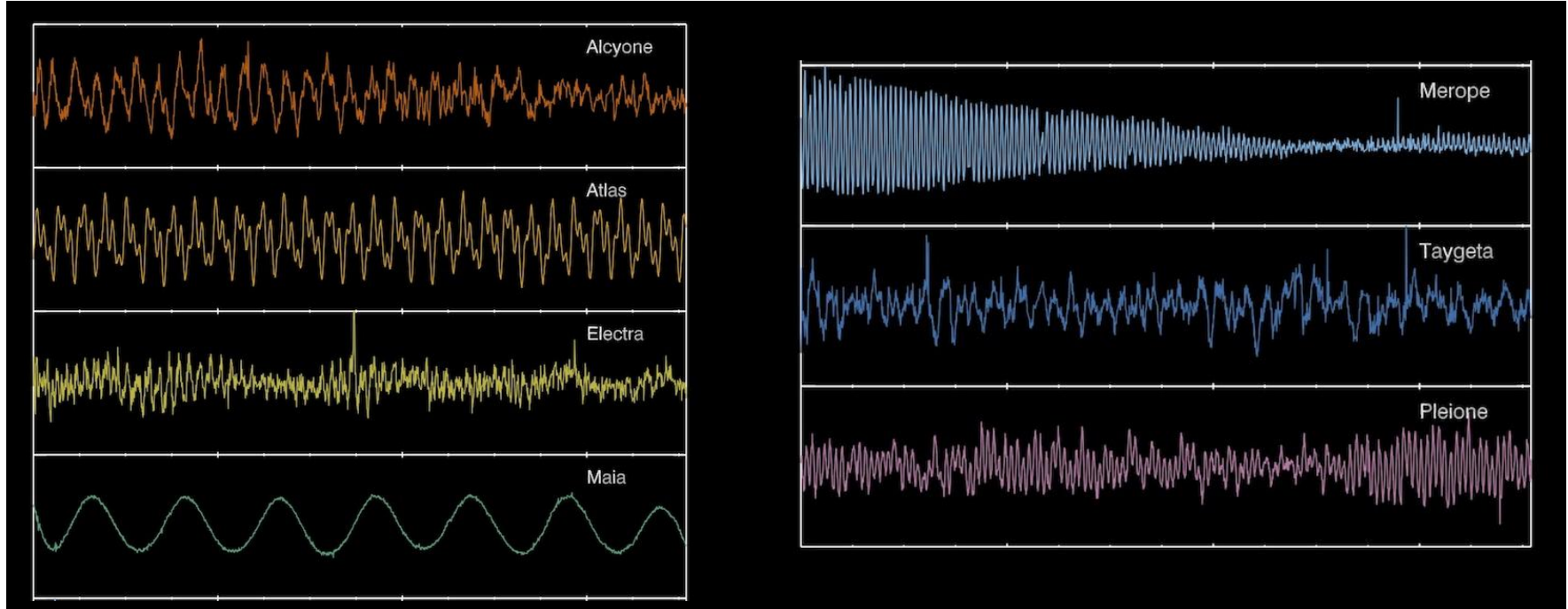
# Watching the sky
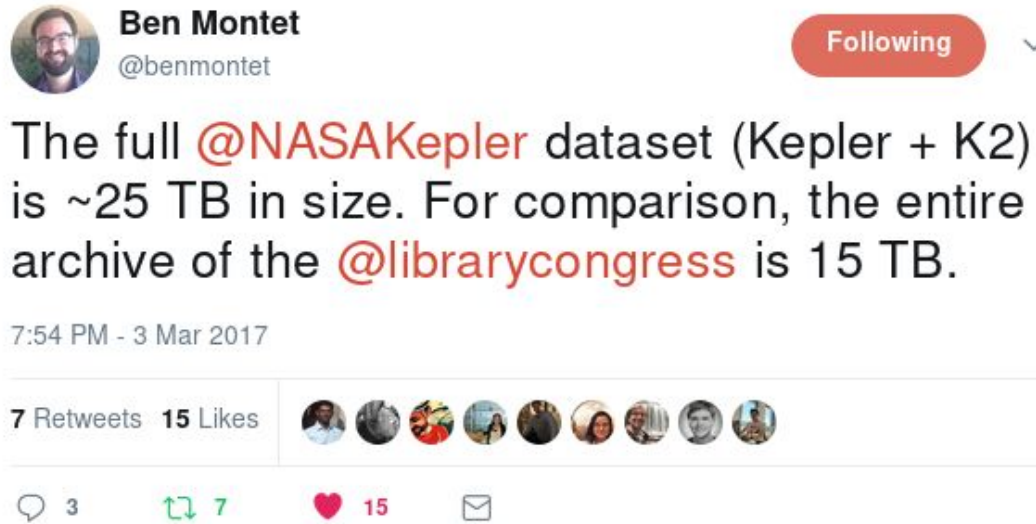
# And what kind of data we get?

# Well, that's the problem



Seven stars, seven different profiles

# Kinda big data



**Ben Montet**
@benmontet

**Following** ⌄

The full @NASAKepler dataset (Kepler + K2) is ~25 TB in size. For comparison, the entire archive of the @librarycongress is 15 TB.

7:54 PM - 3 Mar 2017

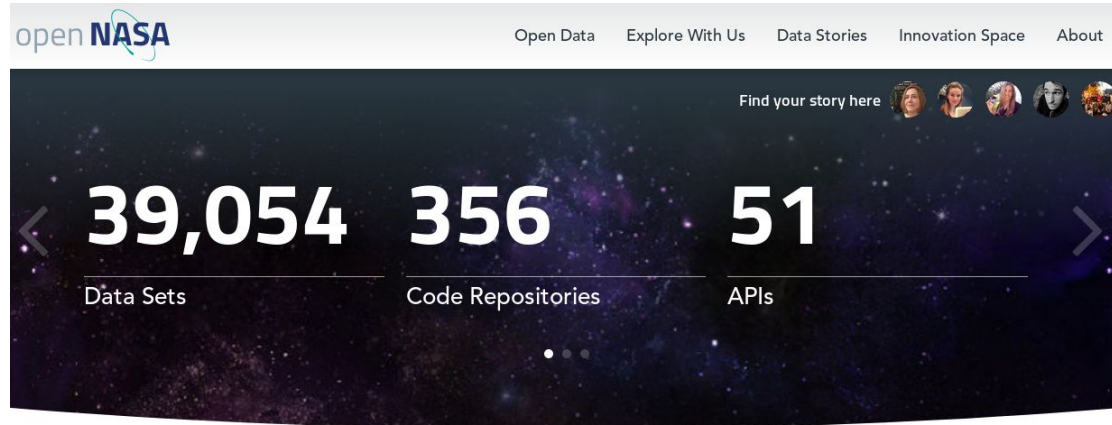**7** Retweets **15** Likes

💬 3        ⟲ 7        ❤️ 15        ✉️

## Over 40 million light curves

# Big AND open data



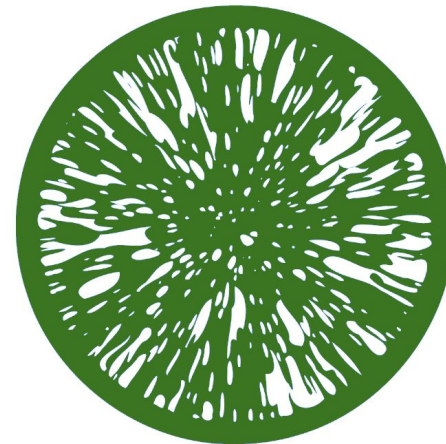Lots of datasets in #opendata

# And we can help with that!



Let's use our tools to analyse the data

# A match made in heaven

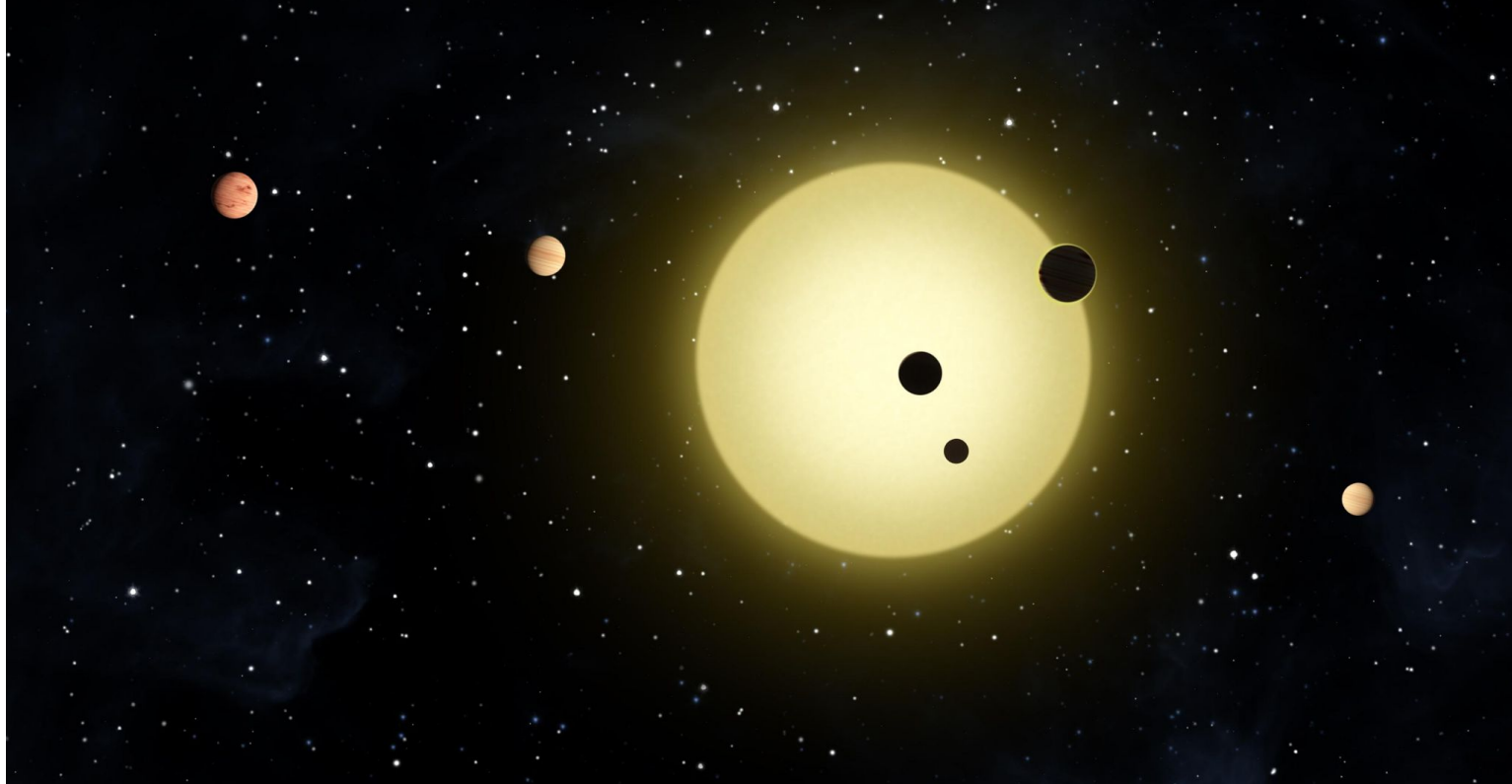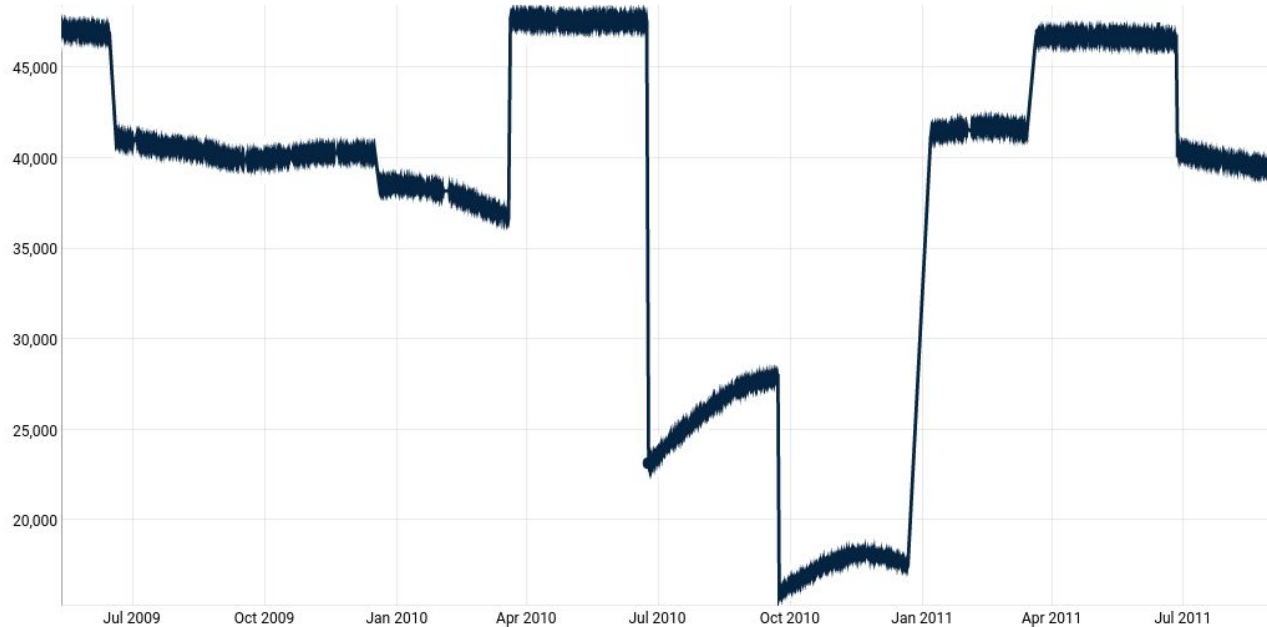## Warp 10, OVH Metrics and HelloExoWorld

# What we have done

- Downloaded and parsed 40 millions of FITS files
- Pushed it to OVH Metrics
- Select a cool subset as training set
- Verified we could find the same planets as NASA
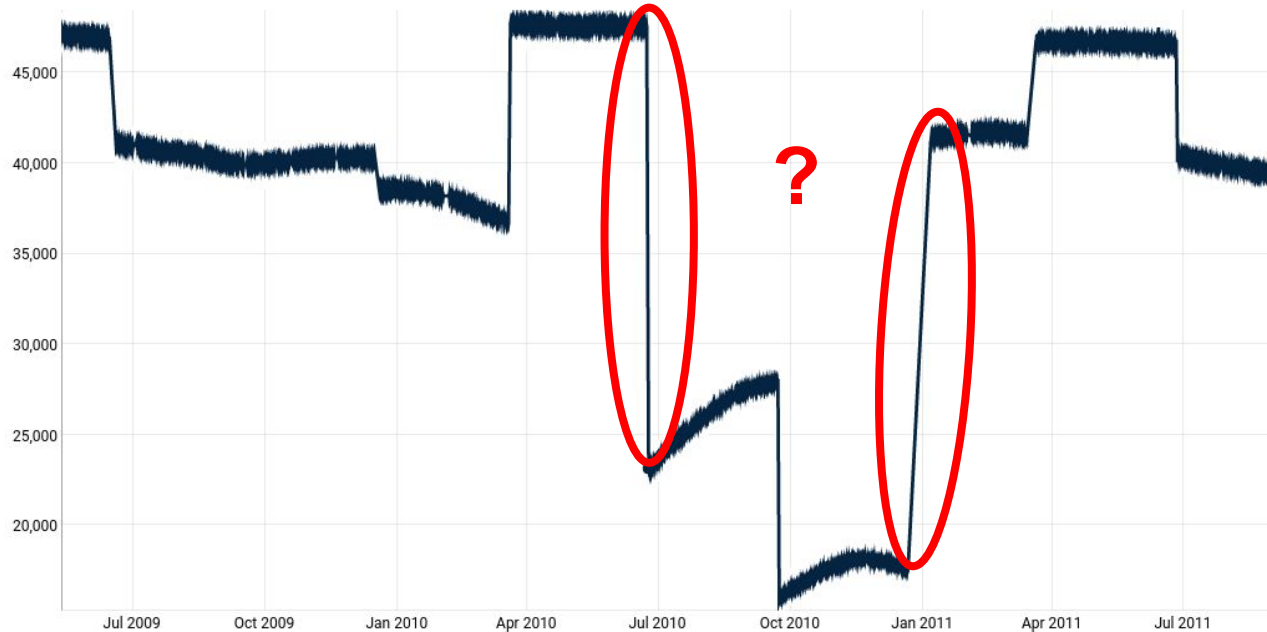
# Looking at the raw signal…



*SAP_FLUX:*

*The flux in units of electrons per second contained in the optimal aperture pixels collected by the spacecraft.*
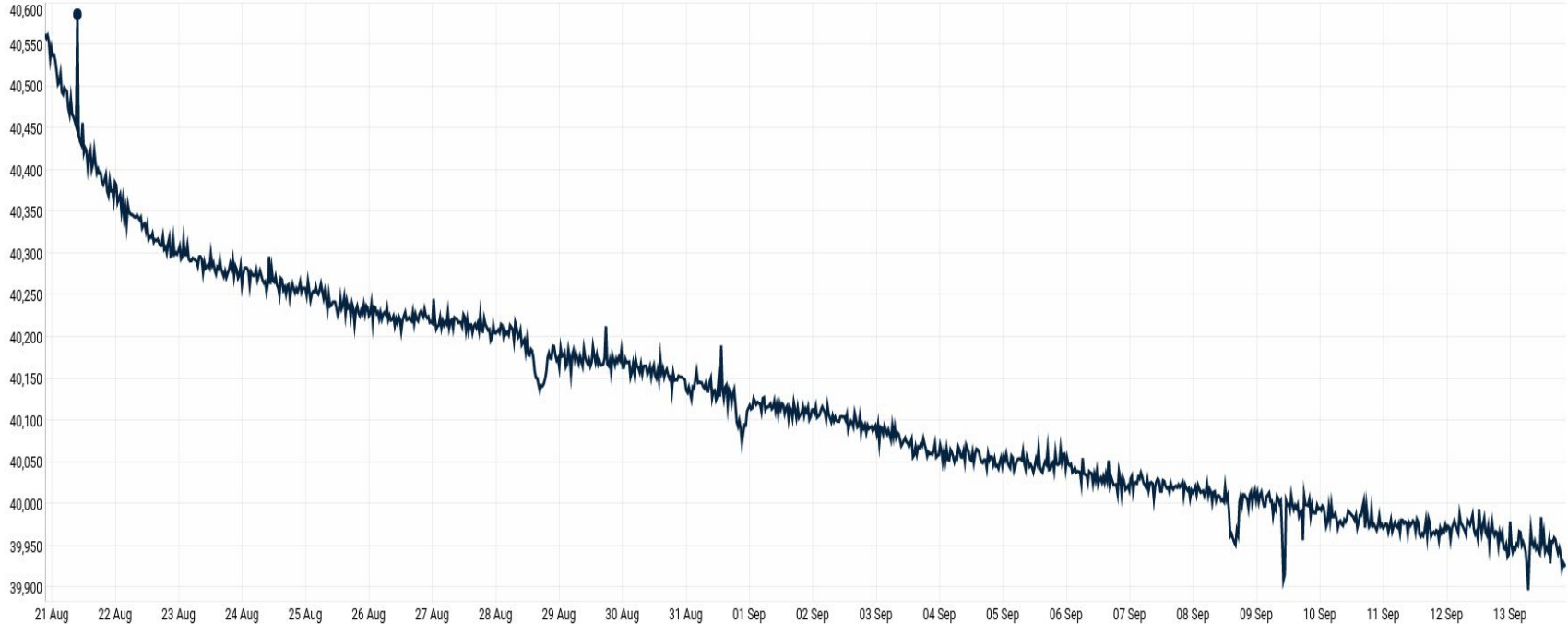
# Looking at the raw signal…



SAP_FLUX:
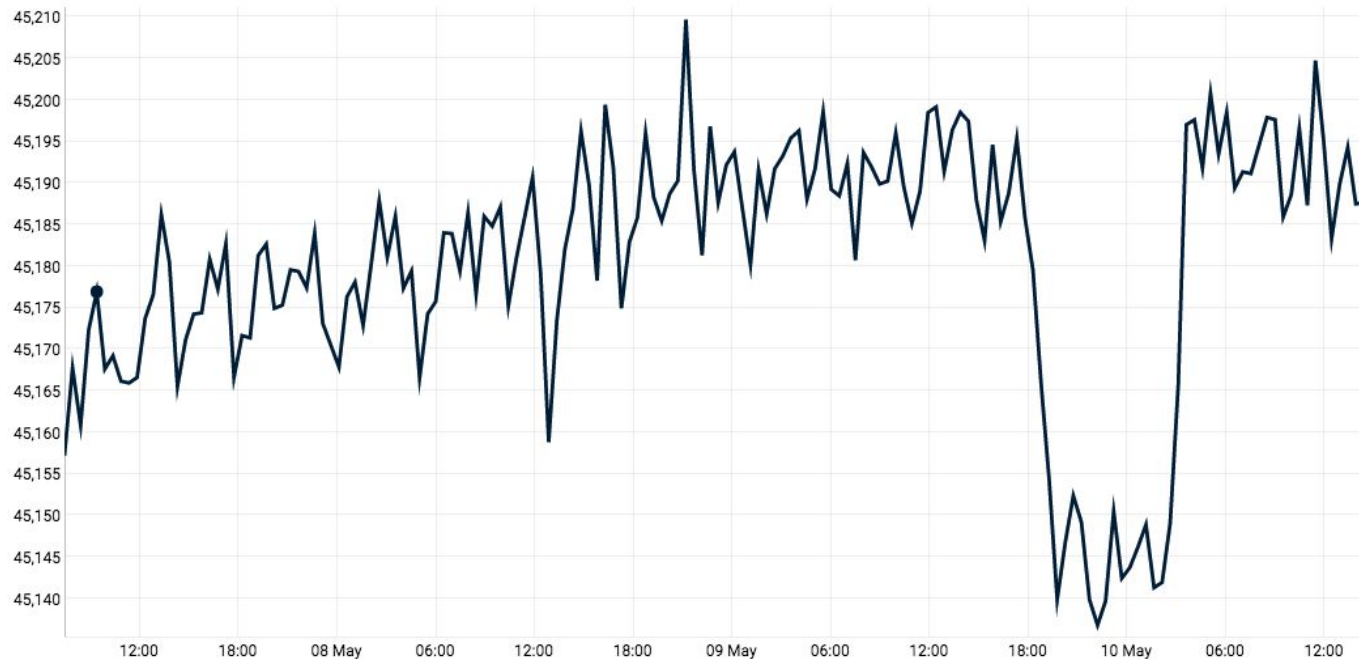*The flux in units of electrons per second contained in the optimal aperture pixels collected by the spacecraft.*

# Looking at one record
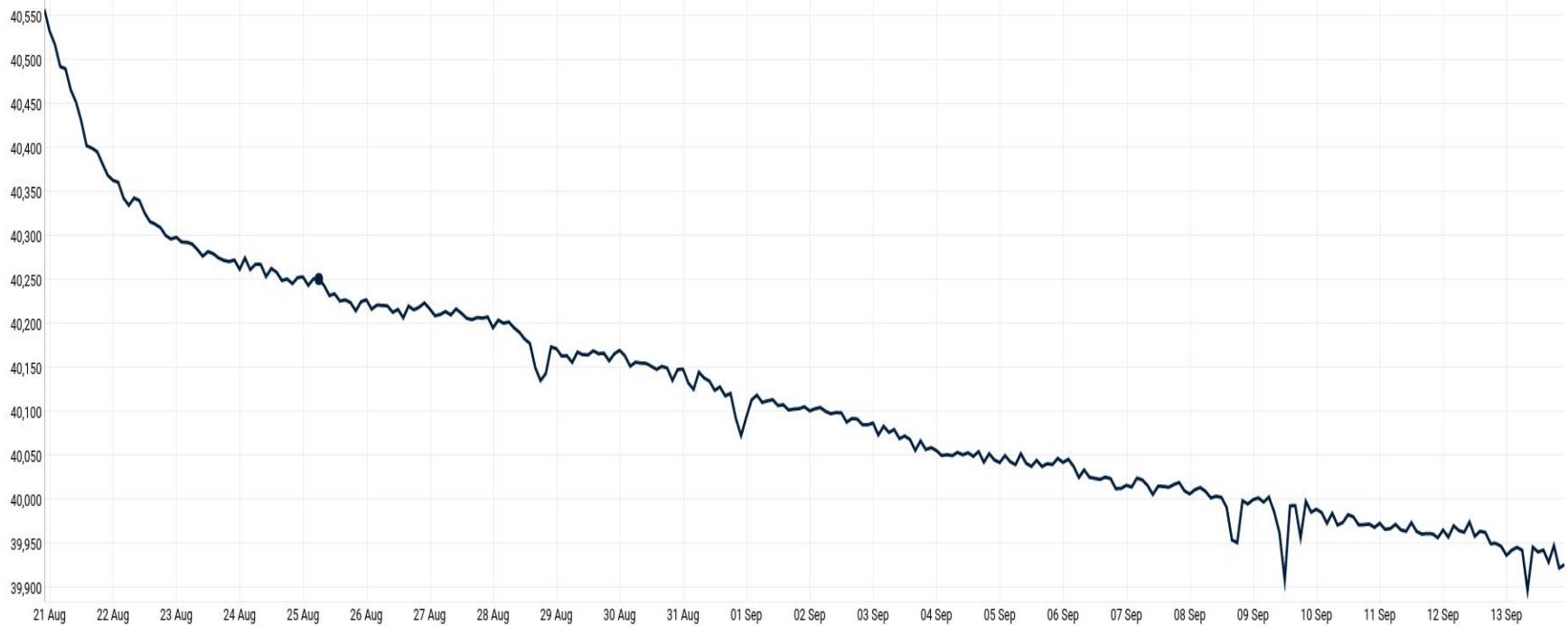


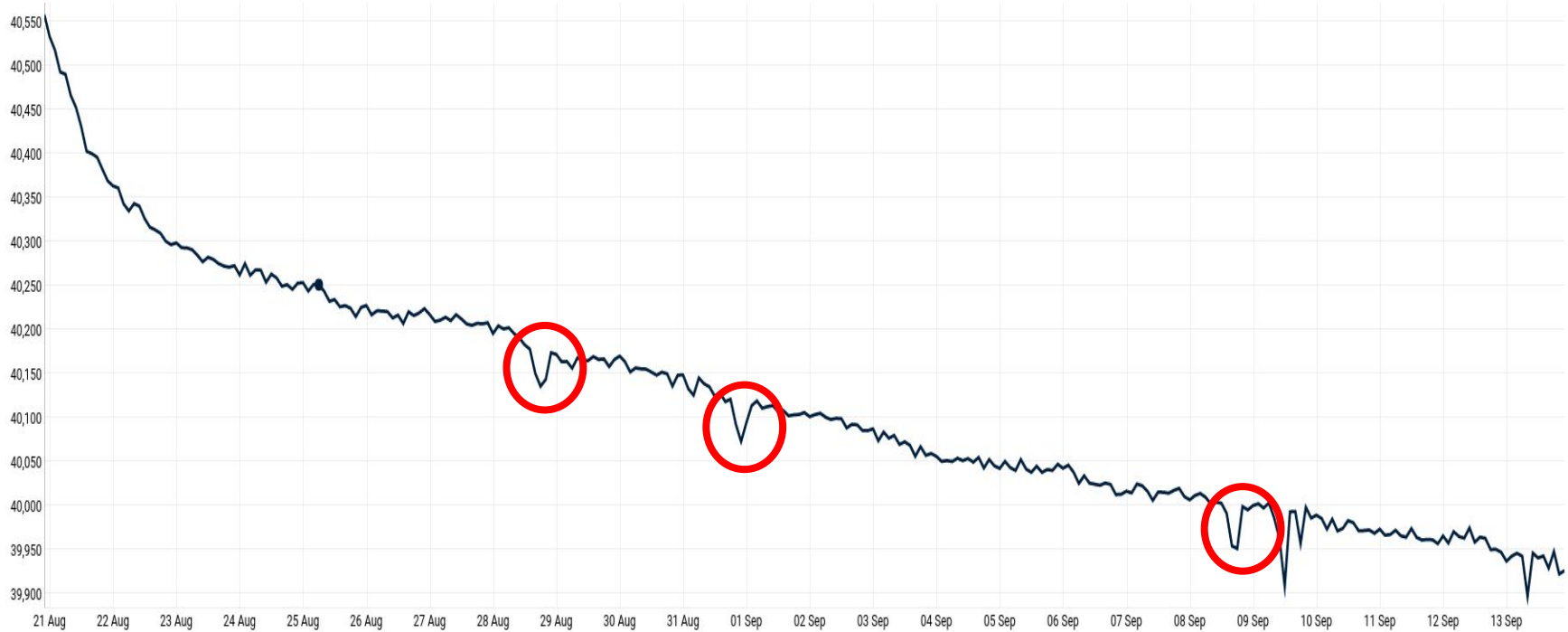## Perturbations in dirty signals

# Transits are tiny



## ~40 electrons per second

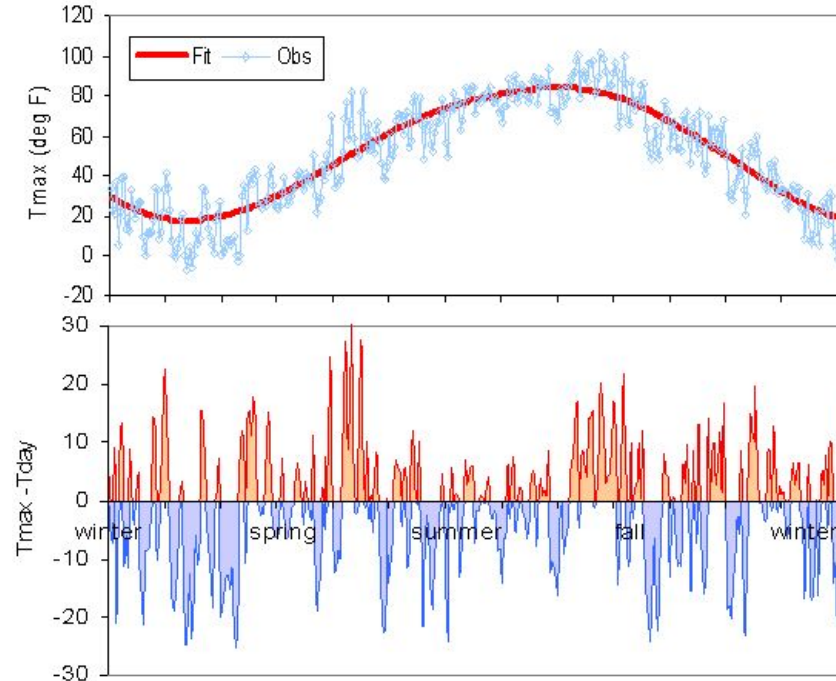# First step: downsampling

# First step: downsampling



You can see the transit candidates…
but how can we teach the computer to see them?

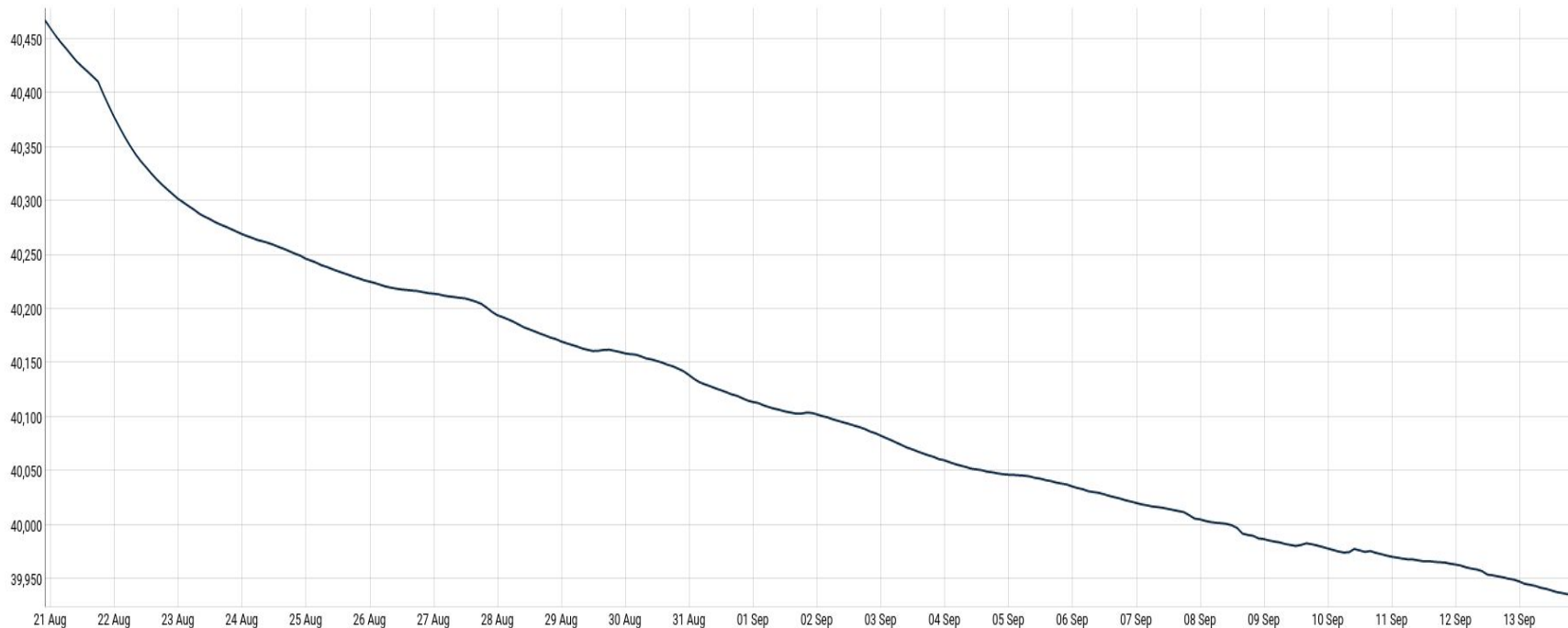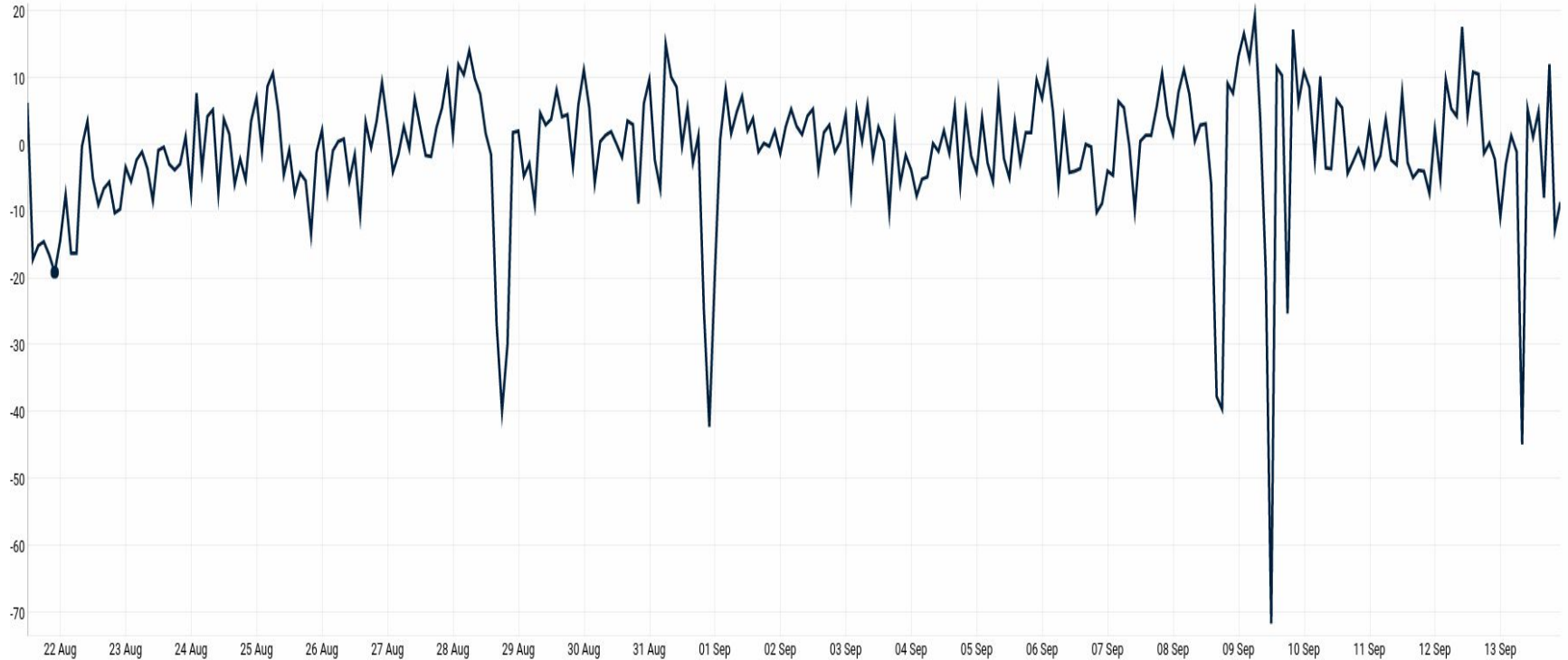# If you ♥ signal processing

## High pass filter

# Poor person's high pass filter
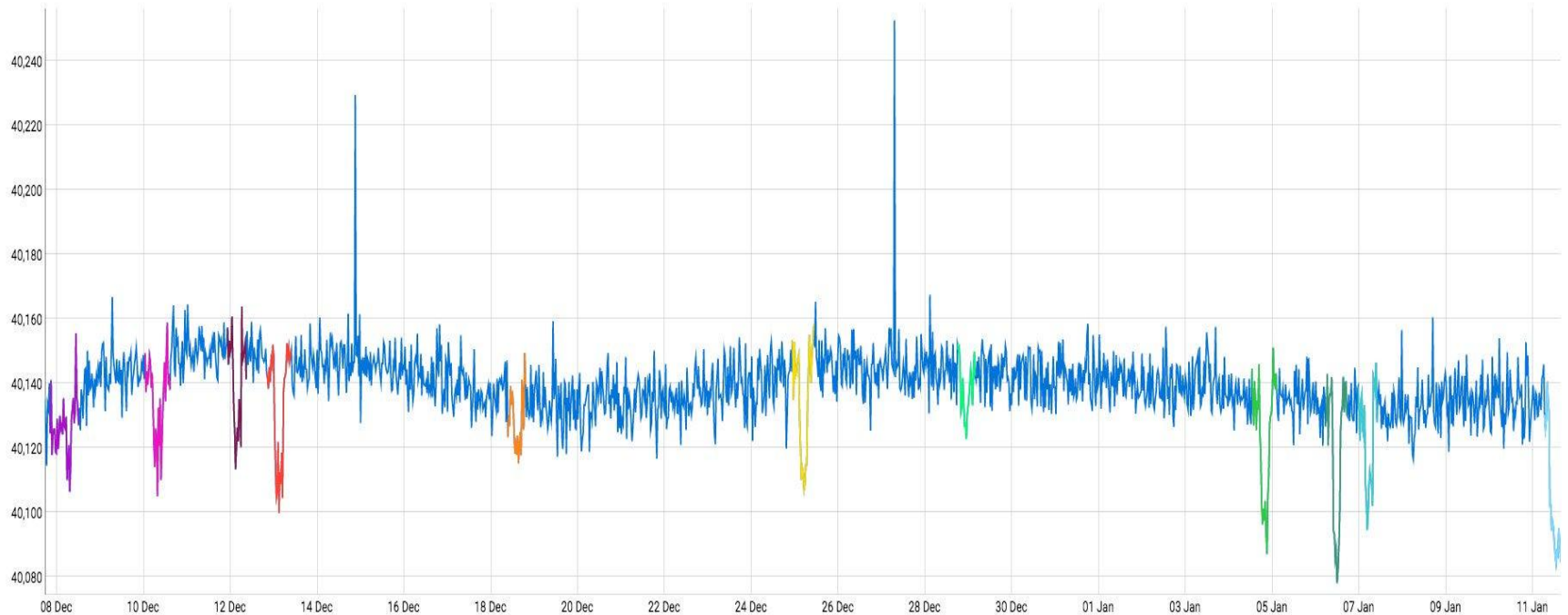


## Using the trend

# Signal - Trend



## Now you can see them well

# After some tuning



## We have our transit candidates

# What's next?

## Where do we go from here?

# Only the beginning



New import method

Better detection

Deep learning

Explorer

satellite/star location

Yours?

# A growing team

# And you!

BASED ON THIS DECREASE IN THE STAR'S BRIGHTNESS, I BELIEVE IT IS ORBITED BY AT LEAST ONE PLANET.

EXOPLANET ASTRONOMERS AT NIGHT

https://xkcd.com/1371/

Join us!
**https://helloexo.world**

# Thank you!

# Want to know more?

## Analysing with WarpScript

# WarpScript

**Reverse Polish Notation**

| Input | 2 | 3 | add | 11 | mul | 1 | add |
|-------|---|---|-----|----|----|----|-----|
| **Stack** | 2 | 3<br>2 | 5 | 11<br>5 | 55 | 1<br>55 | 56 |

# Variables

'hello, world!'                    // Push Hello World String on the Stack

'exo' STORE                        // Store it in a variable called exo

$exo                               // Then push back exo variable on the stack

# What are the available series?

```
[
    $readToken                  // Application authentication
    '~.*'                       // selector for classname
     {}                         // Selector for labels
]
FIND
```

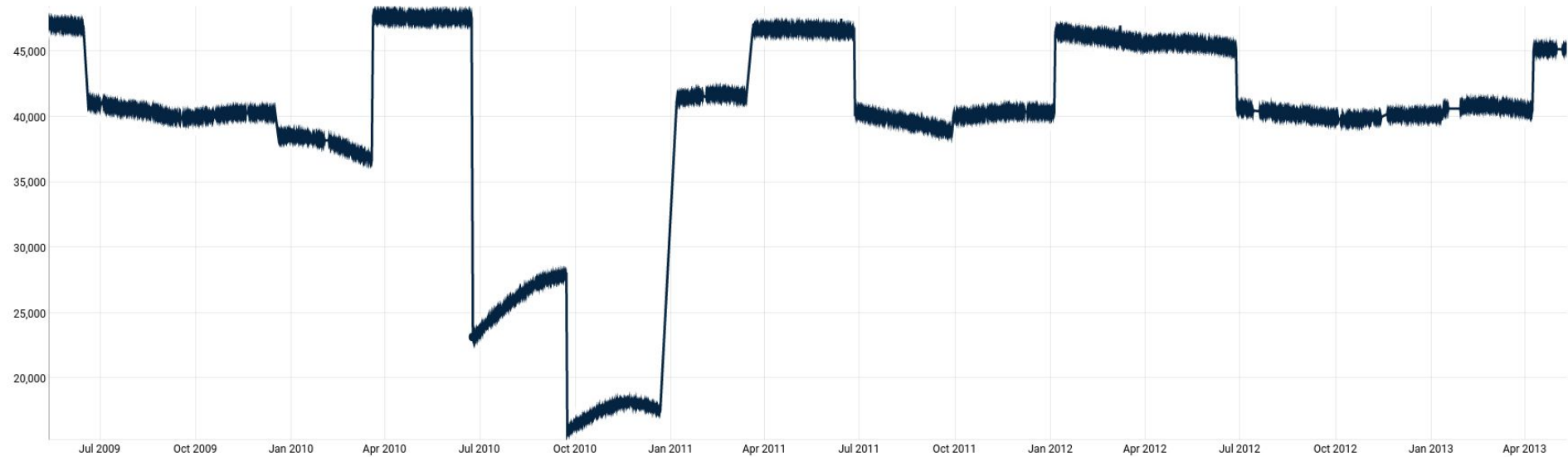# Get raw data

```
[
        $readToken                          // Application authentication
        'sap.flux'                          // selector for classname
         { 'KEPLERID'  '6541920' }          // Selector for labels
        '2009-05-02T00:56:10.000000Z'       // Start date
        '2013-05-11T12:02:06.000000Z'       // End date
]
FETCH
```
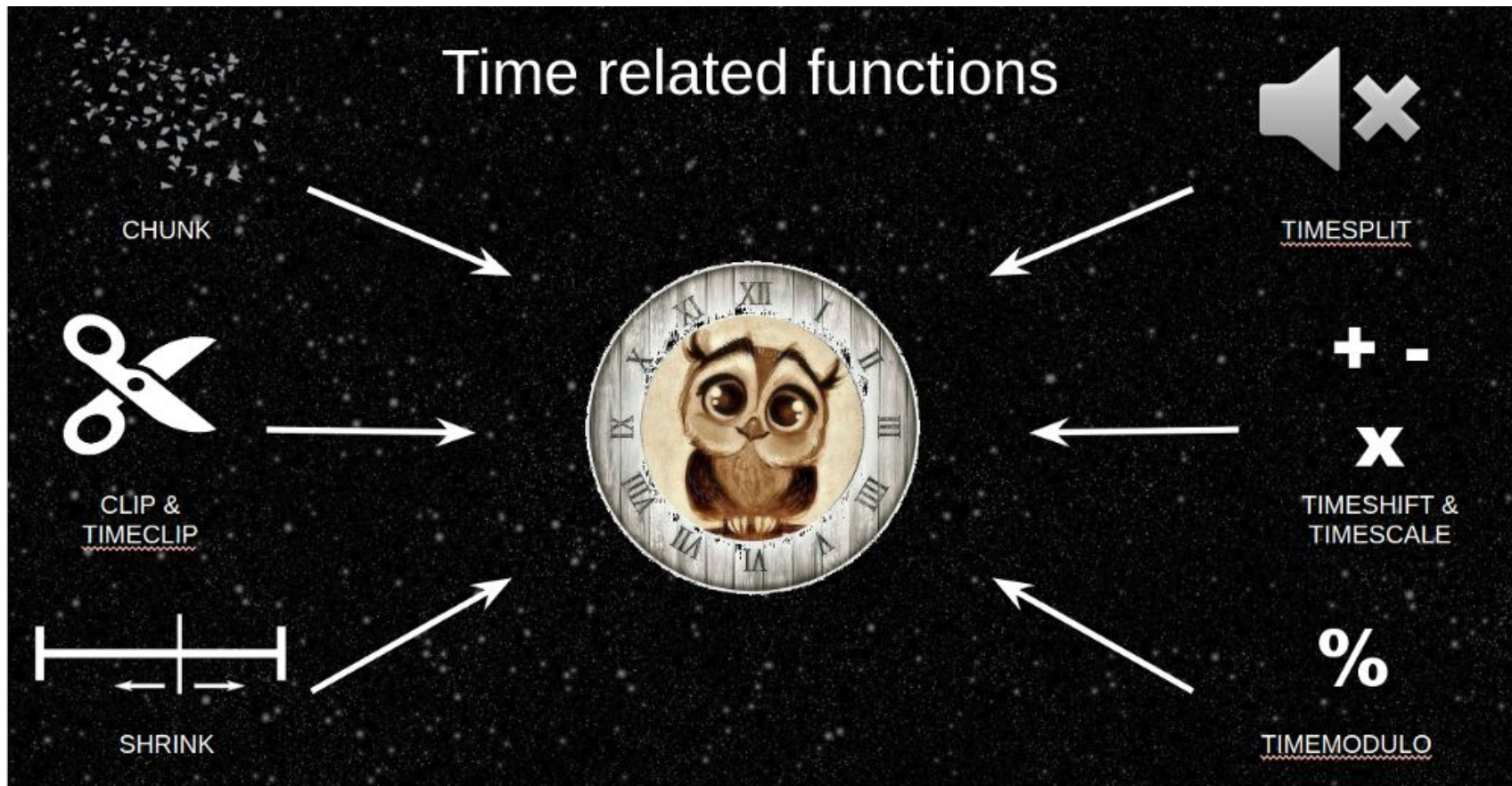
# Kepler-11: Raw data

# Time manipulation

# Time related functions

# How to split a Time series

`$gts`                // Singleton (or list of) GTS

`6 h`                 // Minimum of time without data-points

`100`                 // Minimum of data-points required

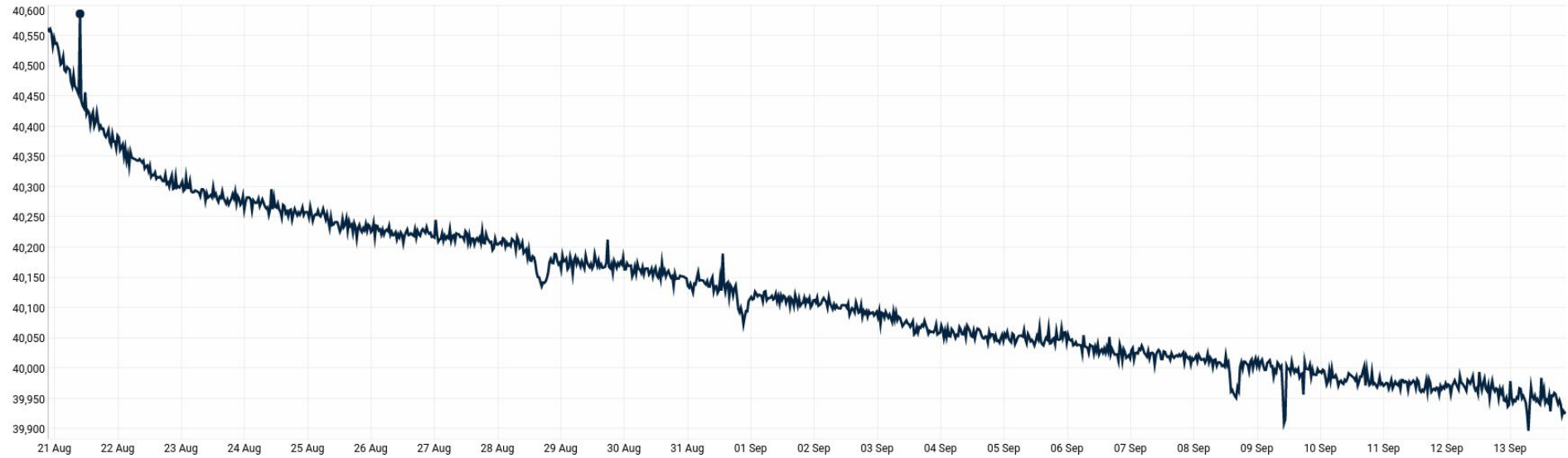`'record'`            // New labels to subdivide the result

**TIMESPLIT**

# Filtering
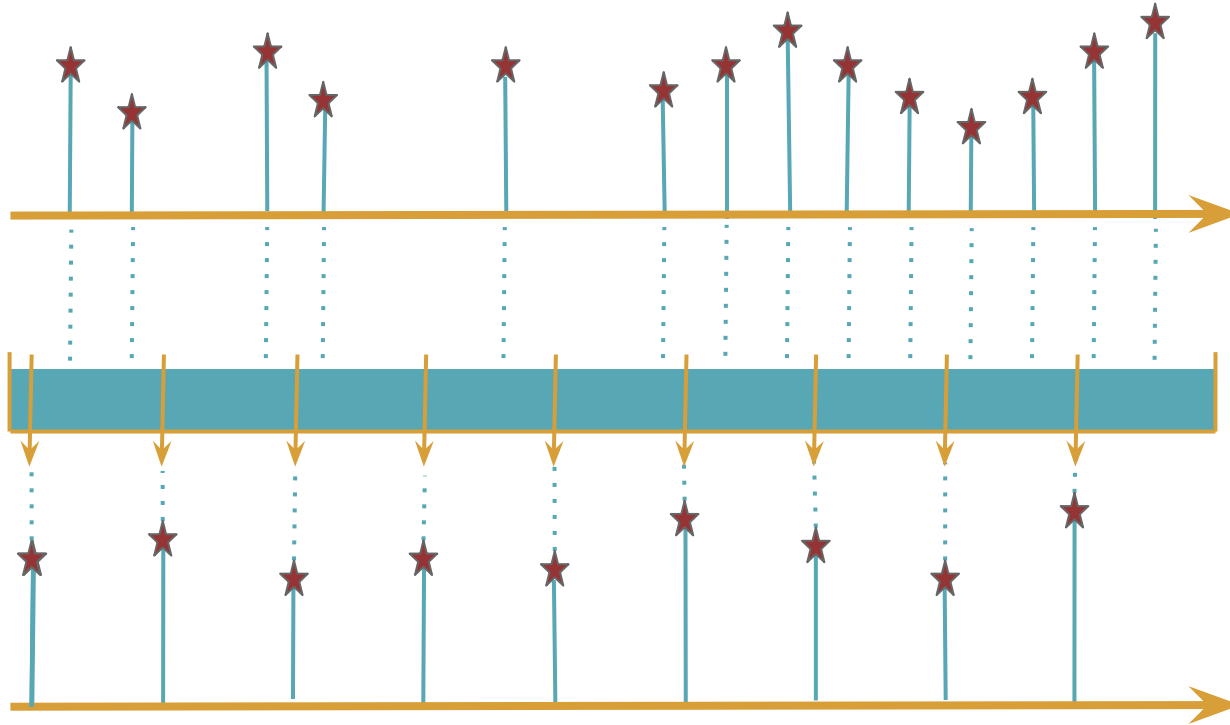
```
[
        $gts                        // Singleton (or list of) GTS
        []                          // Equivalence classes
        { 'record' '5' }            // Labels to select
        filter.bylabels             // Type of filter
]
FILTER
```
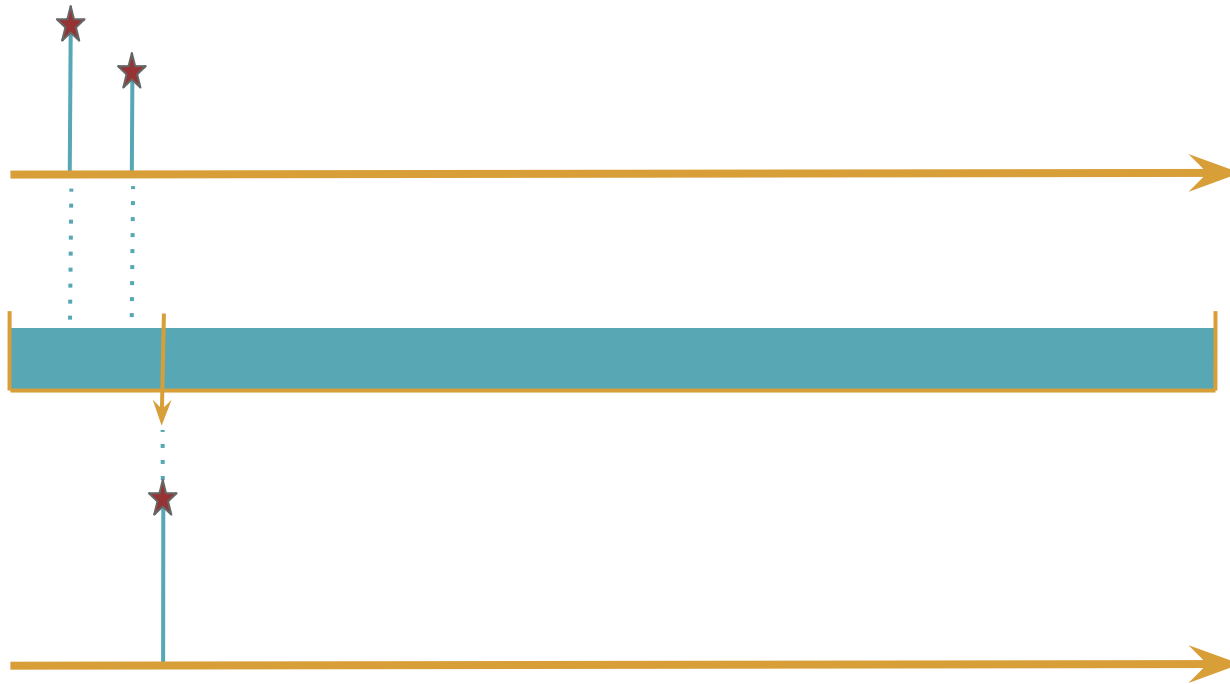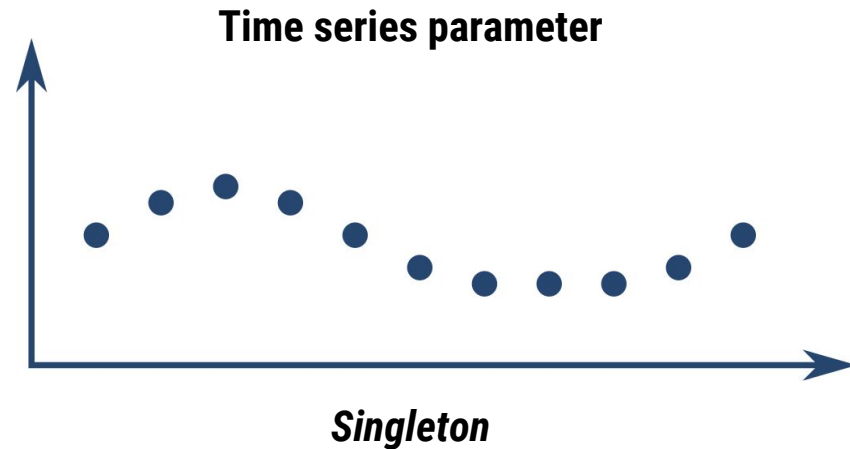
# Reference record: 5

# Downsampling

# Bucketize

# Syntax

**Time series parameter**

[

    $gts

    bucketizer.min

    0

    2 h

    0

]

BUCKETIZE

*Singleton*

# Syntax

**Bucketizer**

[

    **$gts**

    **bucketizer.min**  ⬅

    **0**

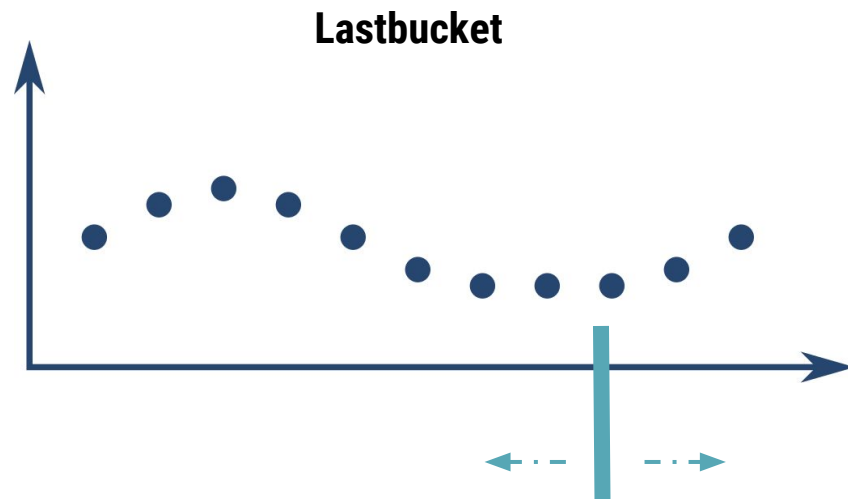    **2 h**

    **0**

]

**BUCKETIZE**

*Type of operator to apply on each bucket*

*last, max, mean, and, count ...*

# Syntax

**Lastbucket**

```
[
    $gts
    bucketizer.min
    0
    2 h
    0
]
BUCKETIZE
```
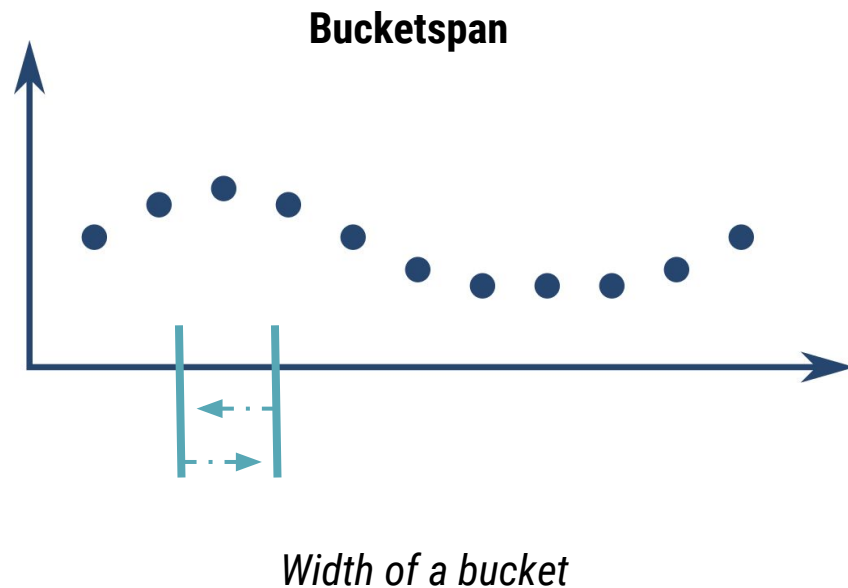
*End timestamp of the more recent bucket*

# Syntax

```
[
    $gts
    bucketizer.min
    0
    2 h
    0
]
BUCKETIZE
```
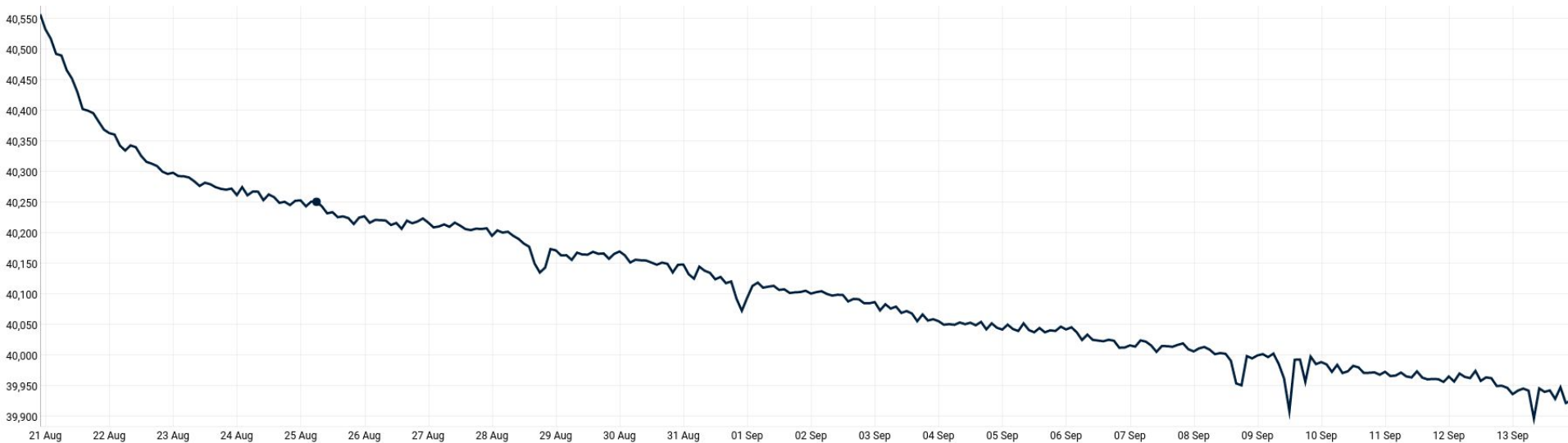
**Bucketspan**



*Width of a bucket*

# Syntax

**Bucketcount**

[

    $gts

    bucketizer.min

    0

    2 h

    0

]

BUCKETIZE

*Number of buckets to keep*

# Actual

# Trend

# Mapper

# Syntax

**Time series parameter**



```
[
    $gts
    mapper.mean
    2
    2
    0
]
MAP
```
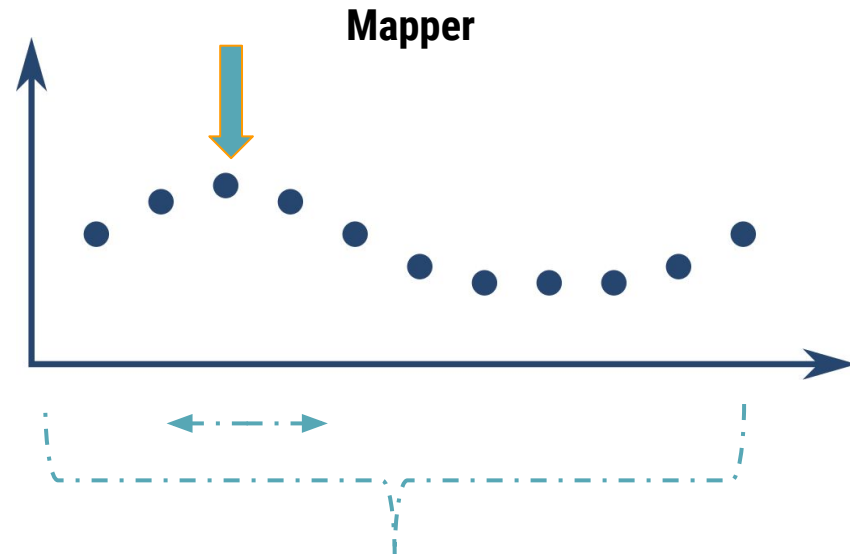
# Syntax

```
[

    $gts

    mapper.mean

    2

    2

    0

]

MAP
```
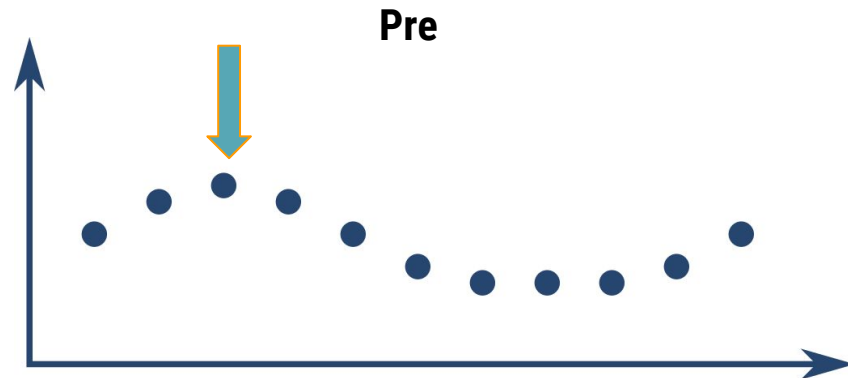
**Mapper**



*Type of operator to apply on each window*
*add, gt, rate, and, count...*

# Syntax

**Pre**

```
[
    $gts
    mapper.mean
    2
    2
    0
]
MAP
```
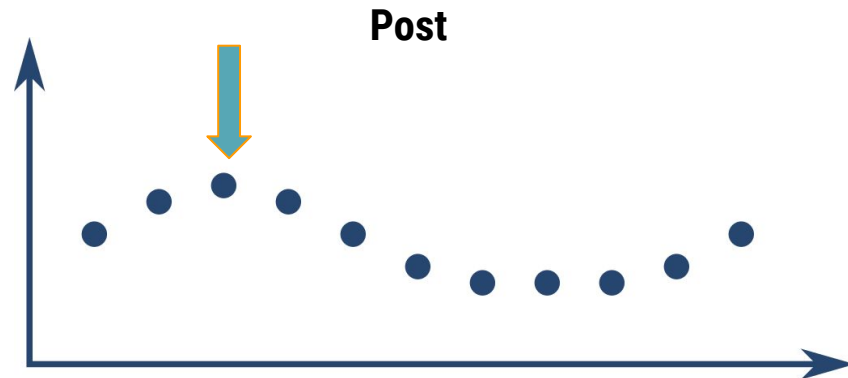
*Number of data-points before*

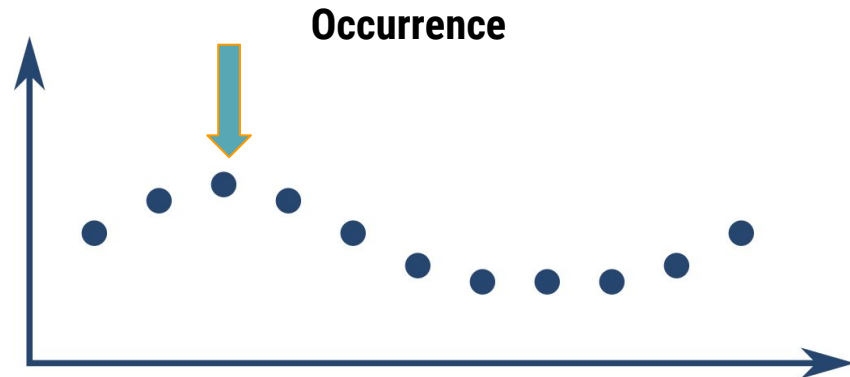# Syntax

[
   $gts
   mapper.mean
   2
   2
   0
]
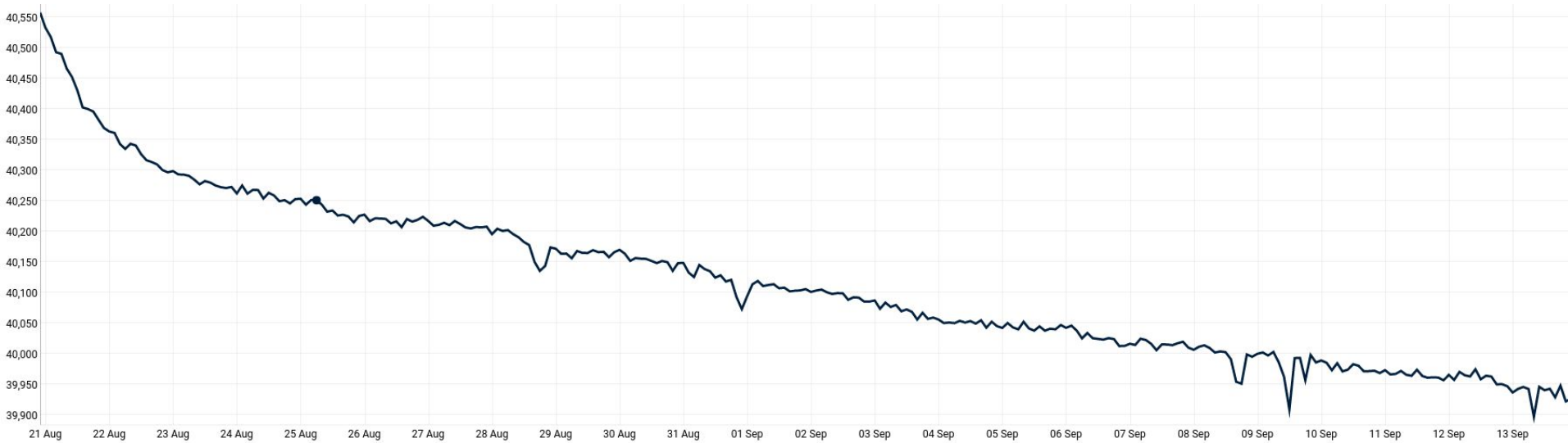MAP

**Post**

←

Number of data-points after

# Syntax

```
[
    $gts
    mapper.mean
    2
    2
    0
]
MAP
```
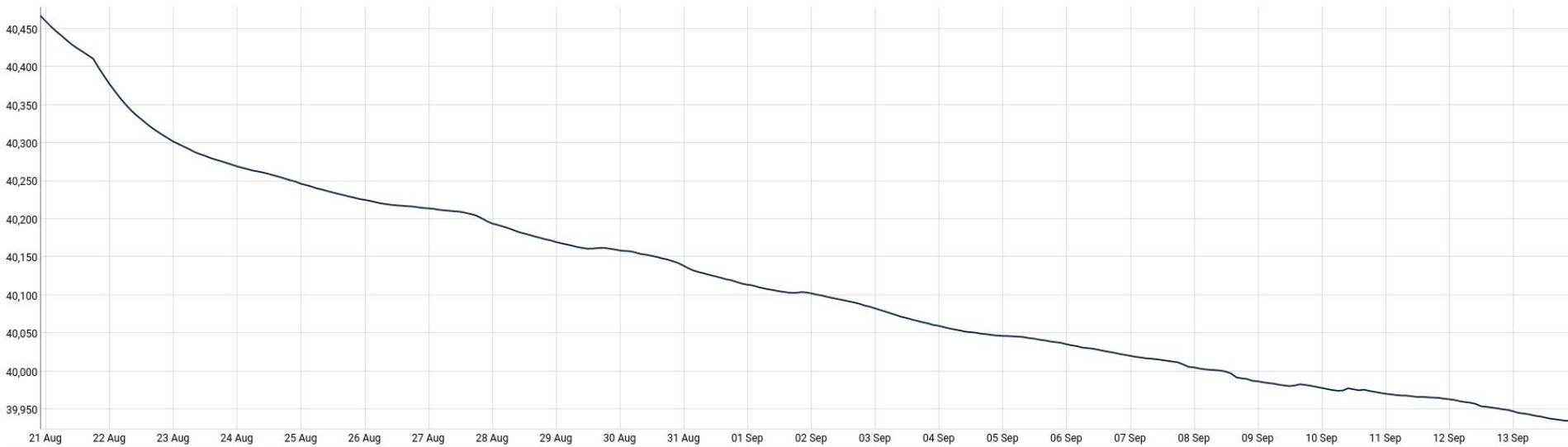
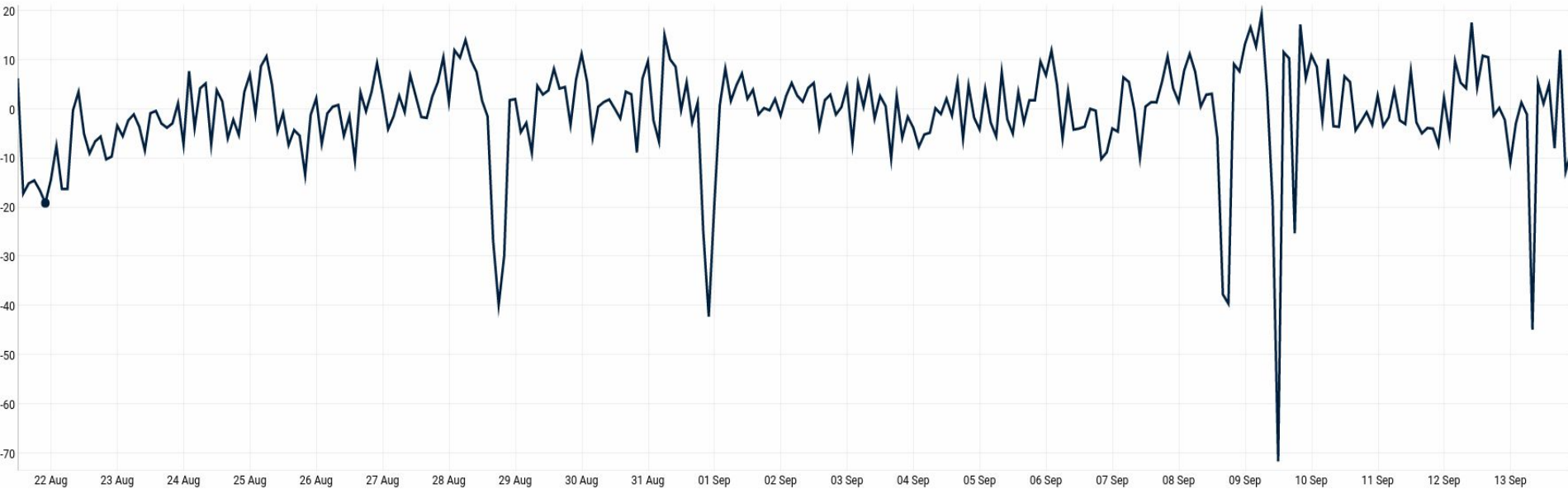**Occurrence**

*Maximal number of calculation for a data-point*

# Actual

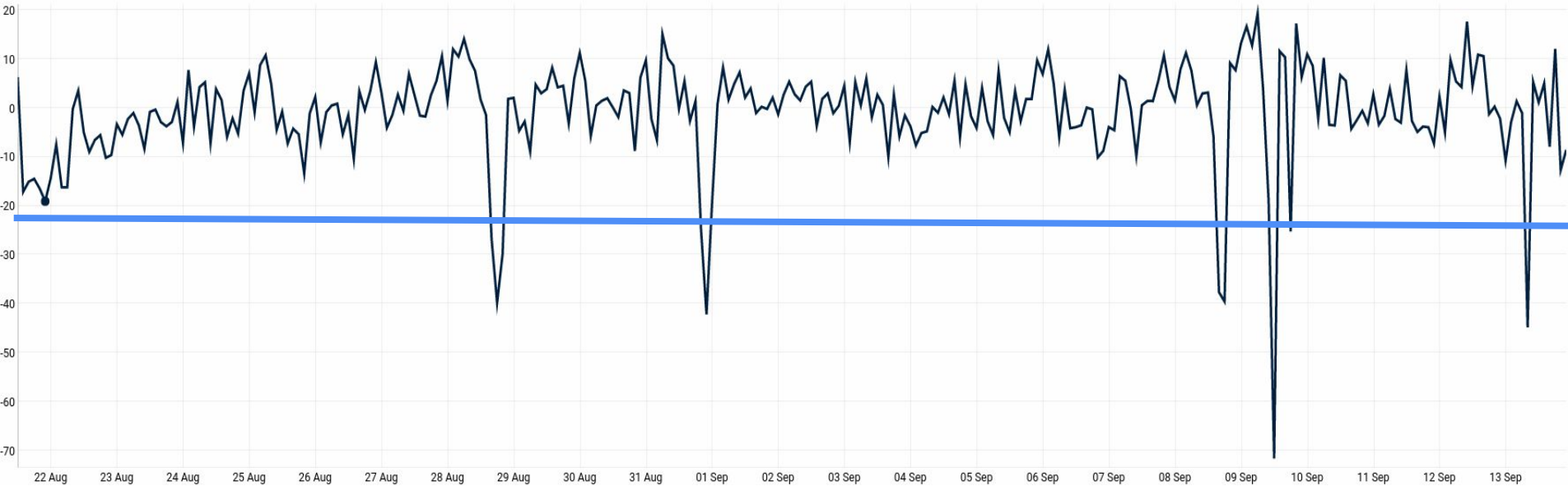# Trend

# Actual - trend

# Actual - trend

# Time to level-up!

# Time series operation

```
[
        $gts0                                   // First series pull
        …                                       // …
        $gtsN                                   // N series pull
        [ 'record' ]                            // Key labels list
        op.add                                  // Type of operator
]
APPLY
```
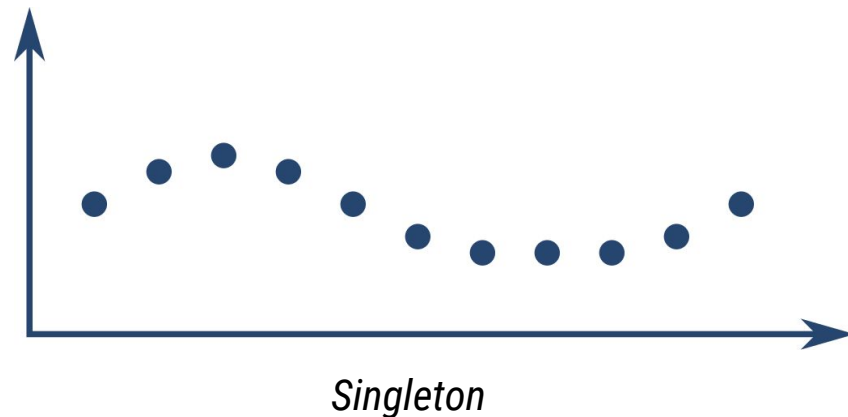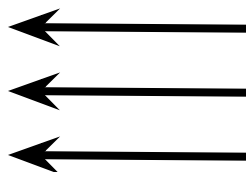
# Syntax

```
[
        $gts0

        …

        $gtsN

        [ 'record' ]

        op.add

]
APPLY
```
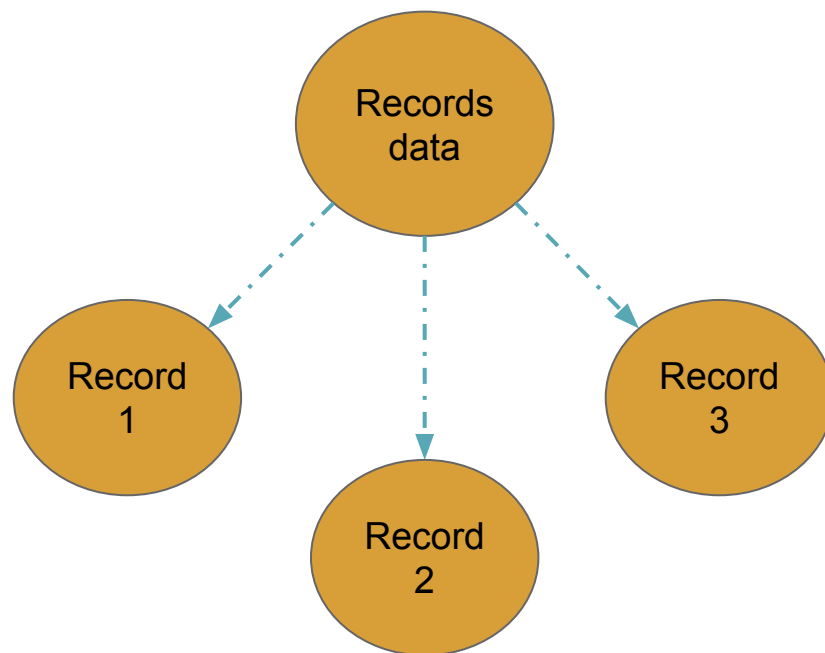


*Singleton*

# Syntax

[
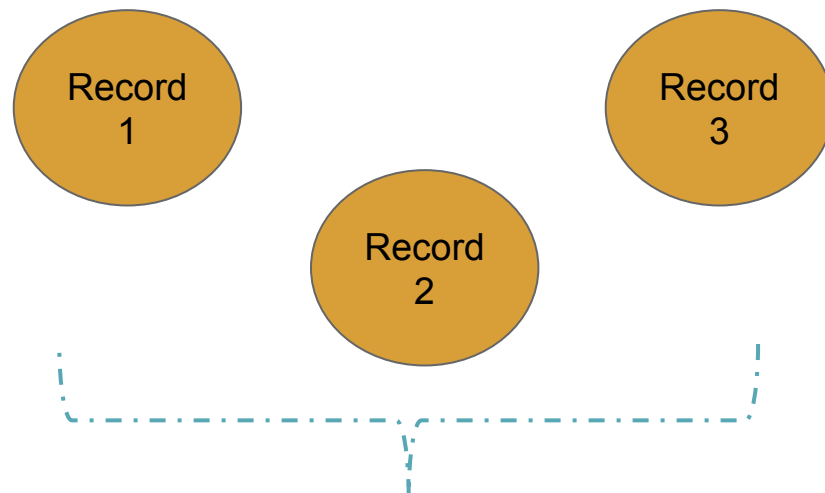
    $gts0

    …

    $gtsN

    [ 'record' ]

    op.add

]

APPLY

**Equivalence class**

# Syntax

**Operator**

[

    $gts0

    …

    $gtsN

    [ 'record' ]

    op.add

]

APPLY



*Type of operator to apply on each class*

*sub, gt, mask, and, mul …*

# Final result