# Applying research in stream processing for fraud detection

Pedro Cardoso
pedro.cardoso@feedzai.com

2019-10-25

# whoami

# whoami



Graduated from FEUP, MSc in computer science 2016

Obligatory tidbits:

- Love sports overall but 🏄 is **the** thing!

- Movies & live music on the streets.

- 🇩🇪 🍺

- Generally having a good time

# whoami

Graduated from FEUP, MSc in computer science 2016

Obligatory tidbits:

- Love sports overall but 🏄 is **the** thing!

- Movies & live music on the streets.

- 🇩🇪 🍺

- Generally having a good time

**@ Feedzai** since **2017**.
Officially a **data engineer** in the **research team**.

Social networks:
- Github: https://github.com/pedro93
- Twitter: https://twitter.com/pedro93
- Linkedin: https://www.linkedin.com/in/pedro-silva-b7968733/

So...

# Stream Processing

# Definition

*Processing an unlimited stream of data with finite resources (Disk/RAM/CPU)*

# Definition

*Processing an unlimited stream of data with finite resources (Disk/RAM/CPU)*



**There are some workarounds for that limitation...**

- Scalar Functions -> *Sum/Count/Avg/Max/Min/Stddev*

- Temporal Queries -> Count of clicks in the past 5 minutes

- Sketches -> Count Min/Bloom Filters/Hyperloglog

- Infinite Pulse Responses -> Exponential averages

- Good Technology Choices

- Distributing Loads -> Partitions by group by keys

- Off-Loading State To Disk

# Why

Stream processing is useful in <u>fast, high data volume</u> throughput systems:

High Frequency Trading

Fraud Detection

Clickstream analysis

Signal/Image/Video Processing

IoT Monitoring

Advertisement auctioning

# Why

Stream processing is useful in <u>fast, high data volume</u> throughput systems:

High Frequency Trading

Fraud Detection

Clickstream analysis

Signal/Image/Video Processing

IoT Monitoring

Advertisement auctioning

And so much more…
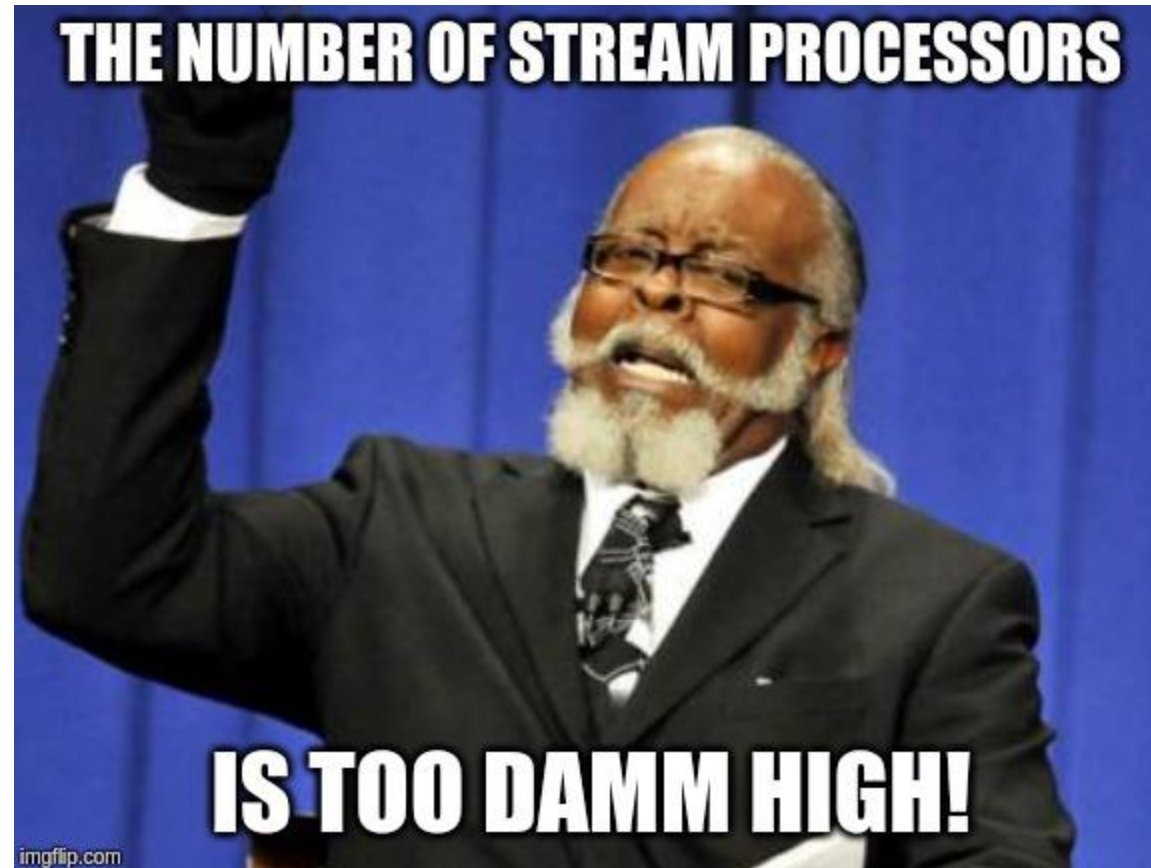The limit is your creativity

# What solutions are out there



There are others out there...

# What solutions are out there

# What is the best for you?

Depends, what are you looking for ?

# What is the best for you?

Depends, what are you looking for ?

- Per-event accuracy or are you ok with micro batching?

- Fault Tolerant

- Low latency

- Adaptable

- Large Community

- Stateful vs Stateless

- High Throughput

- Scalable

- Battle tested

# What is the best for you?

Depends, what are you looking for ?

- Per-event accuracy or are you ok with micro batching?

- Fault Tolerant

- Low latency

- Adaptable

- Large Community

- Stateful vs Stateless
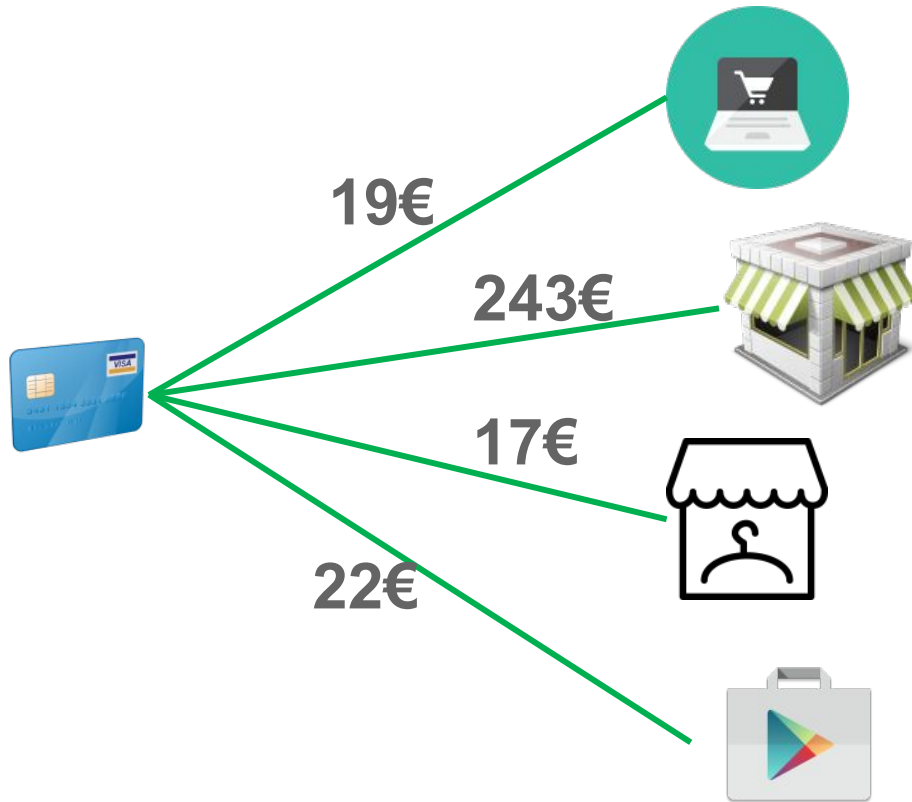
- High Throughput

- Scalable

- Battle tested

**Selecting** one that has all these characteristics for your **specific use-case** is **HARD!**

# @ Feedzai

# @ Feedzai

**PKernel**

# Streams to the rescue!



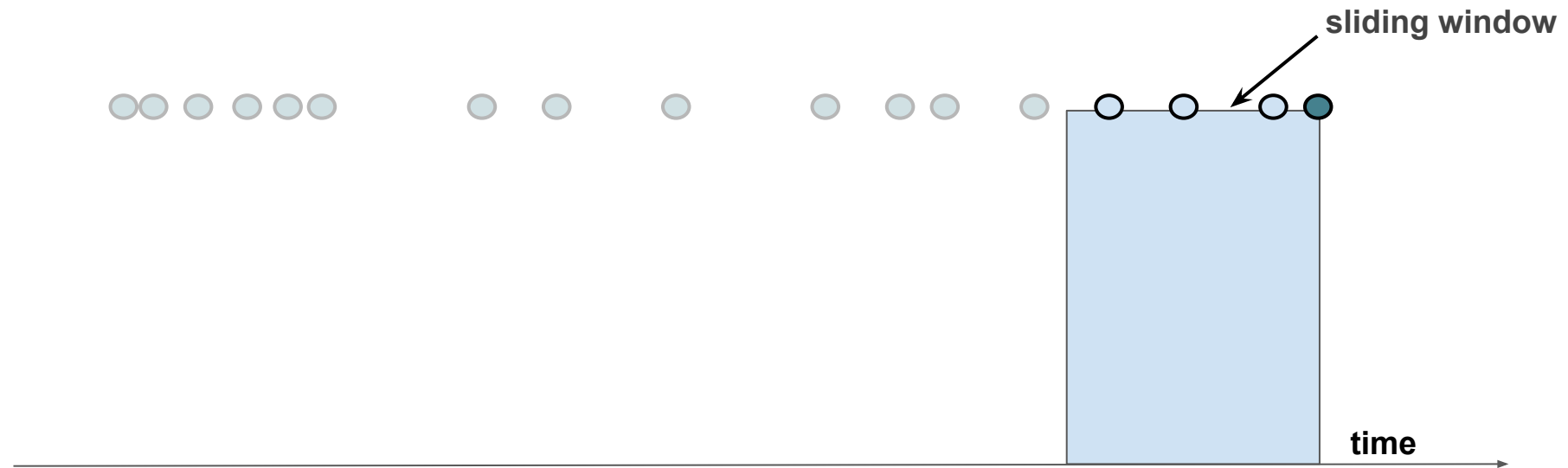**Profiles** are temporal aggregations over some set of fields:

- Average spending in the last week.

- Distance to average location.

- Number of transactions in the last 10 minutes, ...

Very very useful as features for fraud-detecting ML Models & Rules

19€

243€

17€

22€

# Lightweight Profiles

A research project using Exponential Moving Averages

# Context/Motivation



sliding window

time

# Context/Motivation
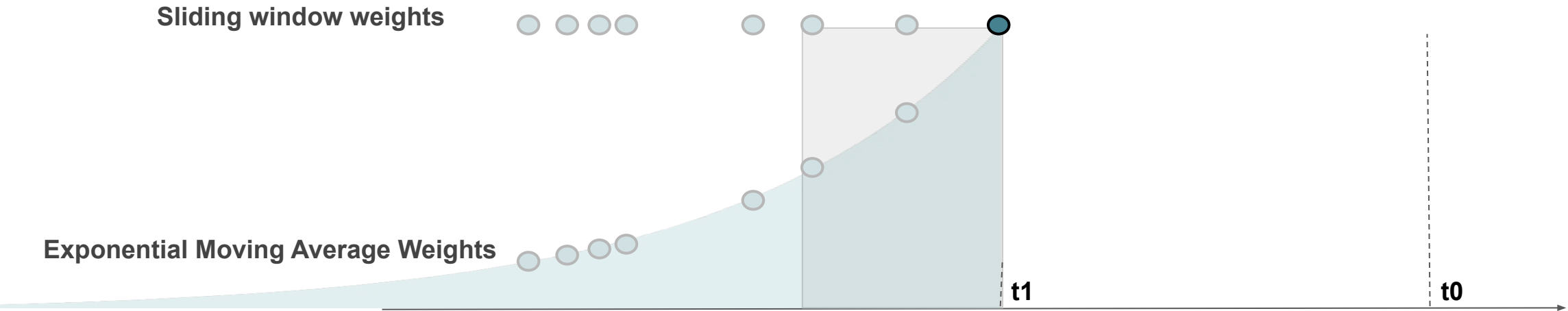
**Naïve** **Sliding window profiles require:**

- Overhead of storing events in the window (memory cost).

- Overhead to update/expire events in/out of the window.

    *Profiles used in projects:* counts, sums, averages, count distincts, etc...
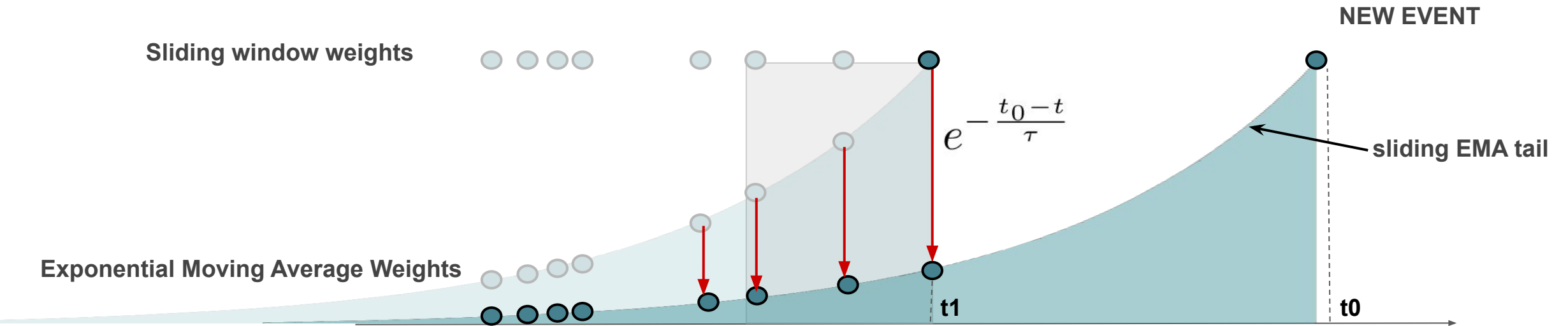
**Alternatives to sliding windows :**

- **Exponential Moving Average (EMA)**: for counts, sums, averages, std dev.

- HyperLogLog (approximate counts): for count distincts.

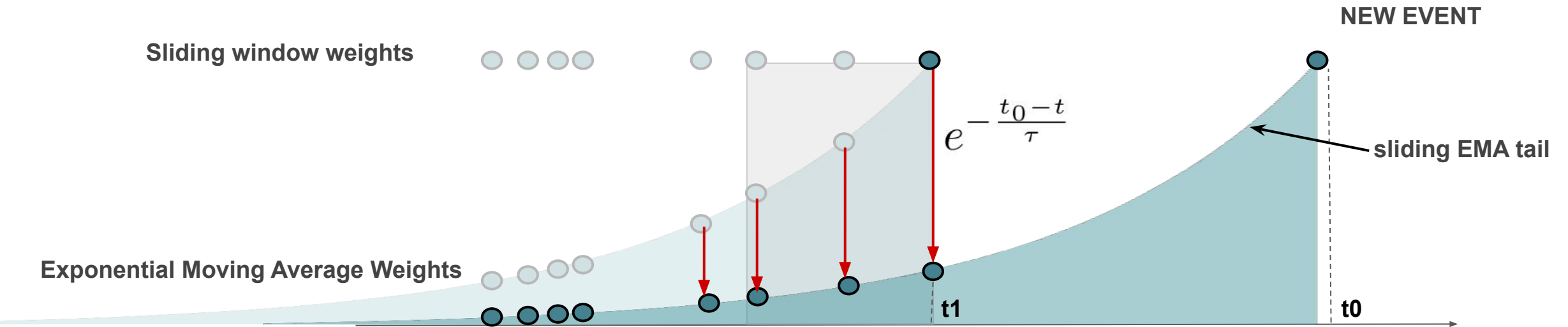# Sliding window VS EMA



**Sliding window weights**

**Exponential Moving Average Weights**

t1

t0

# Sliding window VS EMA



NEW EVENT

Sliding window weights

$$e^{-\frac{t_0 - t}{\tau}}$$

sliding EMA tail

Exponential Moving Average Weights

t1

t0

# Sliding window VS EMA



**Sliding window weights**

$e^{-\frac{t_0-t}{\tau}}$

**NEW EVENT**

sliding EMA tail

**Exponential Moving Average Weights**

t1

t0

$$EMA(t_0) = \sum_{i=0}^{+\infty} z_i e^{-\frac{t_0-t_i}{\tau}}$$

$$= EMA(t_1)e^{-\frac{t_0-t_1}{\tau}} + z_0$$

Theoretically huge
memory savings!

# In real-life

An EMA is everlasting, its value is never 0.

Not practical to store all real-time profiles since the beginning of time.

# In real-life

An EMA is everlasting, its value is never 0.

Not practical to store all real-time profiles since the beginning of time.

**Solution :**

**Key time-to-live**: How long an EMA's "state" and its group by key should live.

In-memory data store of group-by key <=> metric state
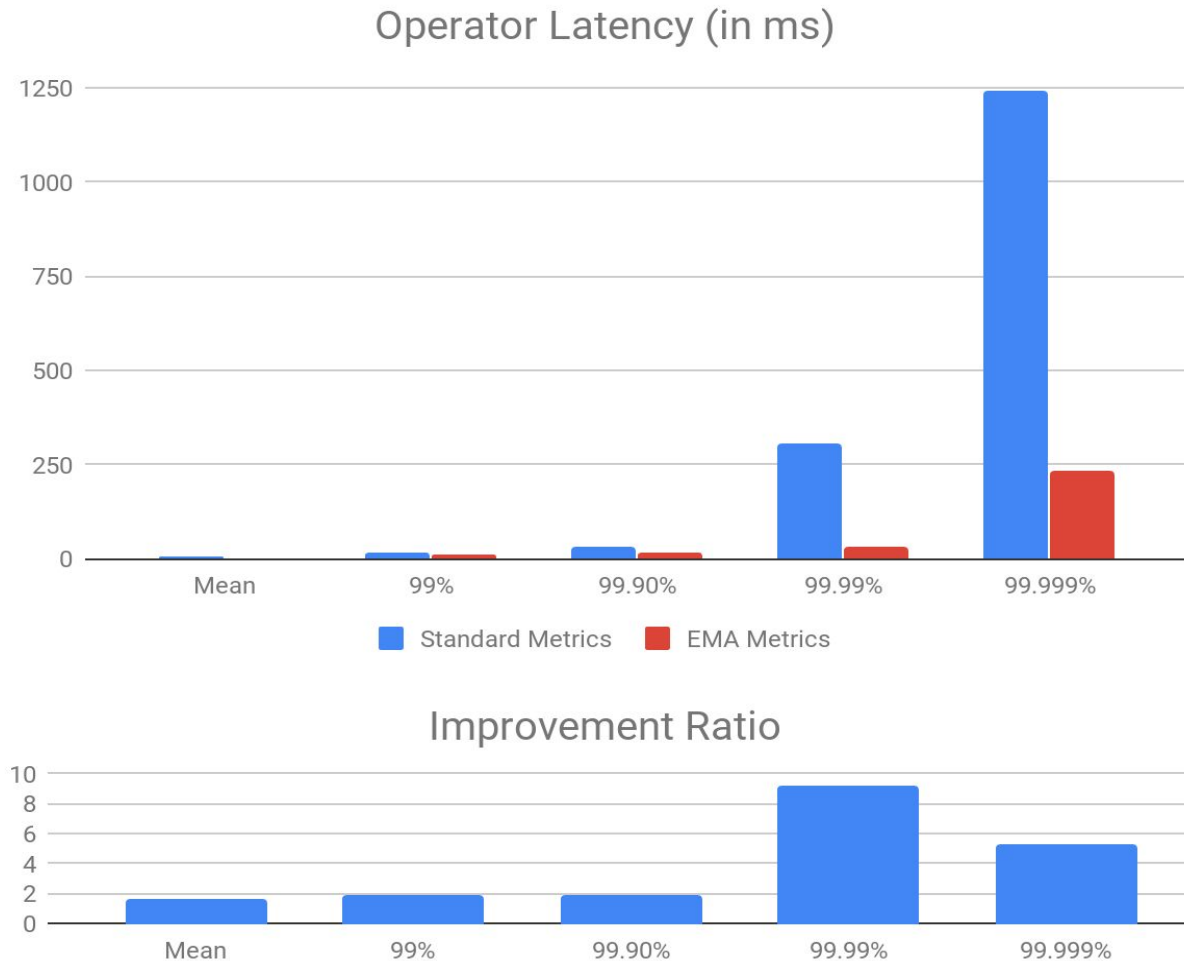
*There goes our memory savings?*

Overhead to update/expire events in/out of the window is now O(1)!

# So did it work???

# Latency Improvements

200 TPS
6.4M event real dataset
150 real-time profiles
over ~1.4M distinct values

# Latency Improvements

## Operator Latency (in ms)



Standard Metrics   EMA Metrics

200 TPS
6.4M event real dataset
150 real-time profiles
over ~1.4M distinct values

## Improvement Ratio



EMAs are **~2x to ~10x** faster to compute than Standard Operators
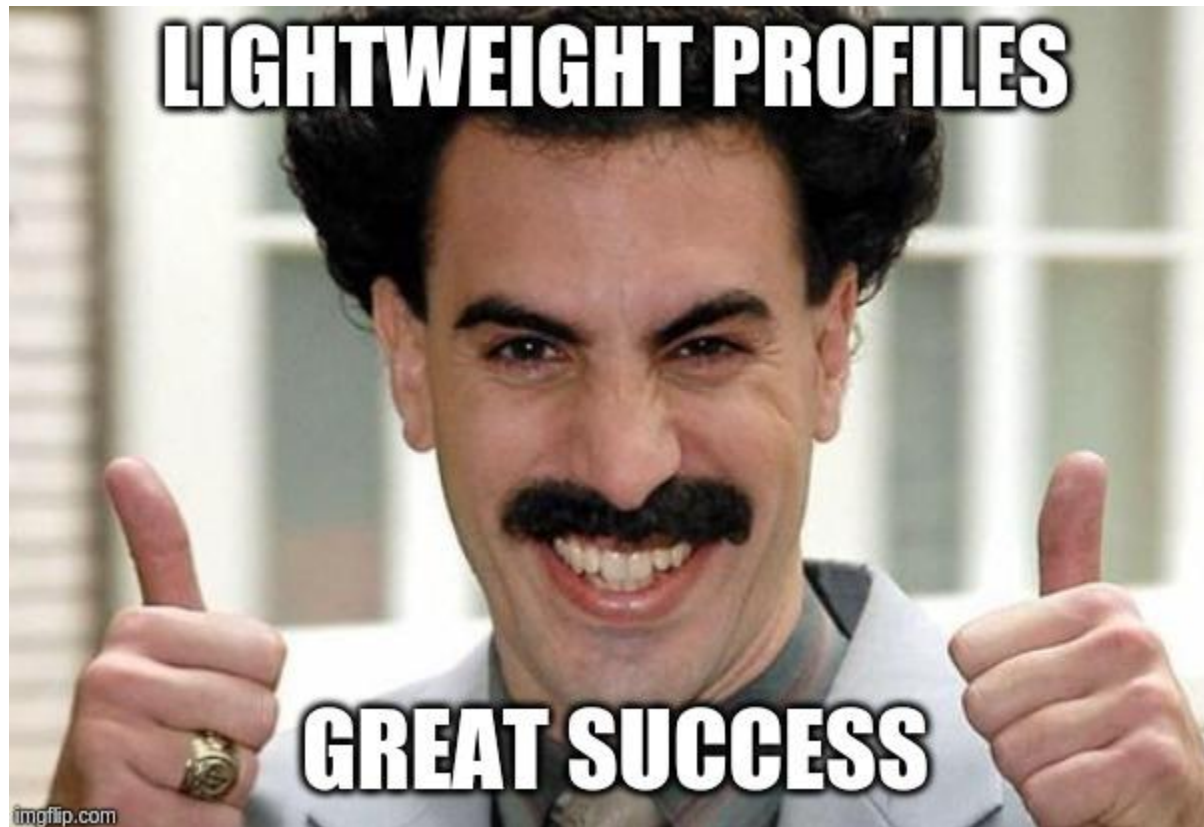
# Memory Improvements

Actual values depend on configuration (key time-to-live factor).

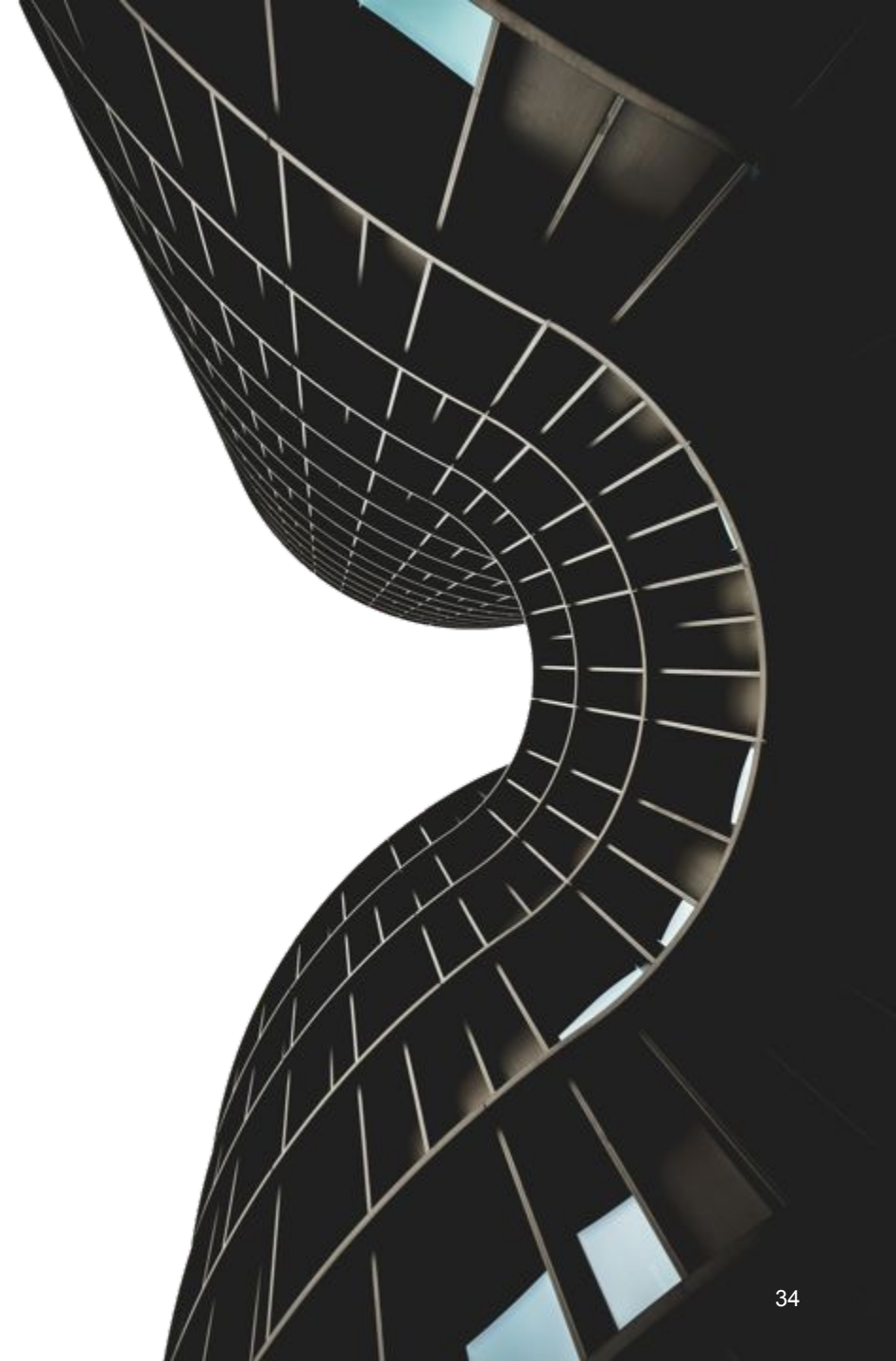EMAs can save **around 50%** memory in real-world scenarios.

Memory costs for profiles no longer depend on the number of events in the window.

# LW-Conclusion

# LW-Conclusion
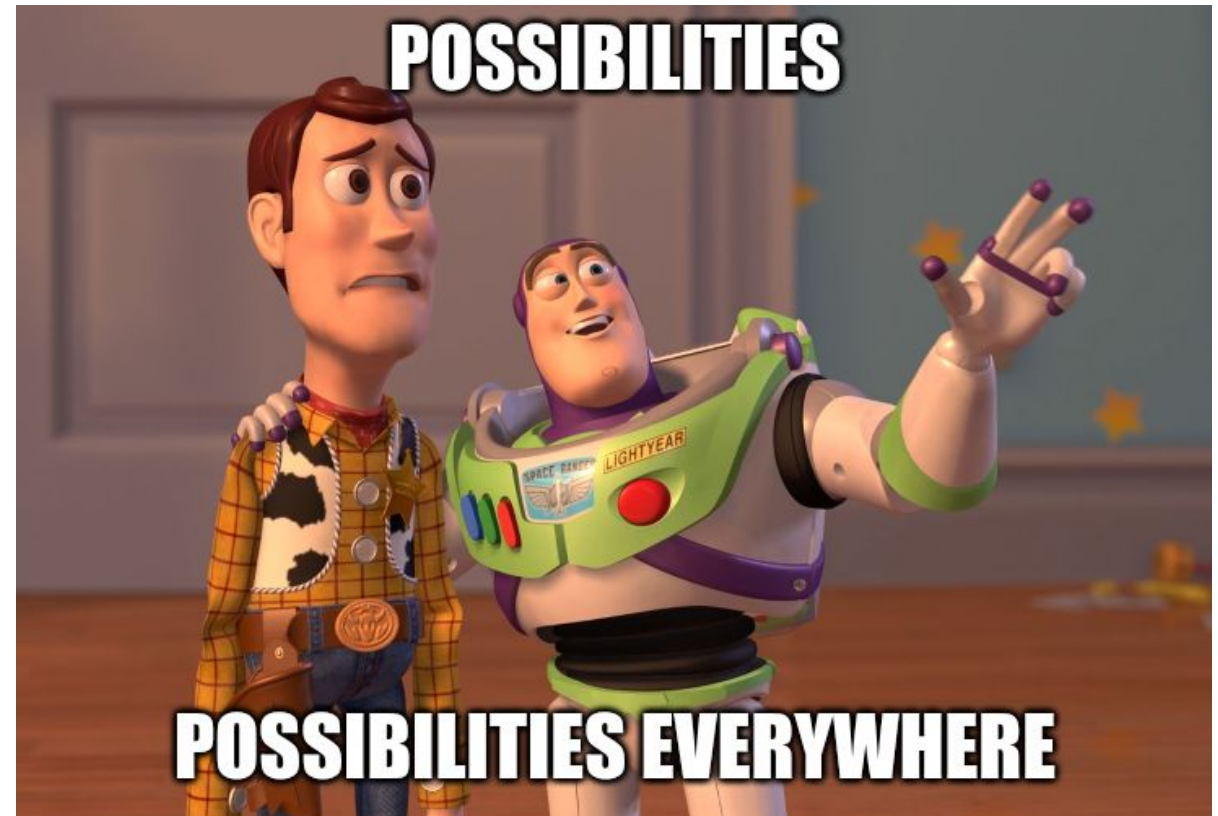
# WIPs

# What are I'm working on right now…

## Railgun
Redesigning internal engine to:
- Simultaneously handle massively large (> 1 year) & small (1 second) real-time windows in a transparent way.
- Compute accurate count-distincts over year-long-windows with millisecond latencies in real-time.

## Lightweight Data Monitoring
Using approximate aggregations to detect shifts in data patterns in real-time with small computing resources.

# Conclusion

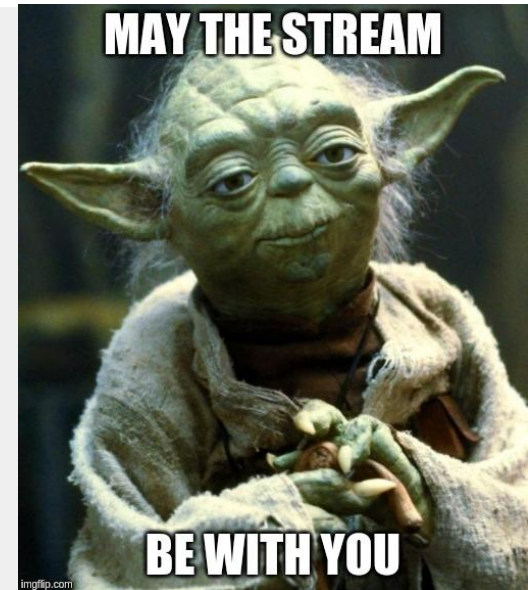There are a lot of streaming engines out there, but sometimes you need something custom.

Very few engines are distributed or do true event-by-event streaming, instead they micro-batch.

Traditional approaches to sliding windows are prone to bursting issues and are memory-intensive.

EMA-based features are lighter and do not lower performance of our ML models

In spite of decades of research, there is plenty of work to be done.

# THANK YOU

# Q&A