# Ephemeral Learning: Augmenting Triggers with online trained normalising flows

Anja Butter, **Sascha Diefenbacher**, Gregor Kasieczka,
Benjamin Nachman, Tilman Plehn, David Shih, Ramon Winterhalder

sascha.daniel.diefenbacher@uni-hamburg.de

**ErUM Data IDT Workshop 2022**
**14.02.2021**

# Introduction

- LHC data rate is incredibly high

  - At 40 MHz ~40 Tb/s

# Introduction

- LHC data rate is incredibly high

  - At 40 MHz ~40 Tb/s

  - Cargo container filled
    with hard drives (100 Pb)

# Introduction

- LHC data rate is incredibly high

  - At 40 MHz ~40 Tb/s

  - Cargo container filled
    with hard drives (100 Pb)
    … every 41 minutes

# Introduction

- LHC data rate is incredibly high

  - At 40 MHz ~40 Tb/s

  - Cargo container filled with hard drives (100 Pb) … every 41 minutes

  - 35 Containers per day

# Introduction

- We can't save everything

  - Reduce data rate: Triggers, save only the interesting events. storage ~ trigger-rate x luminosity
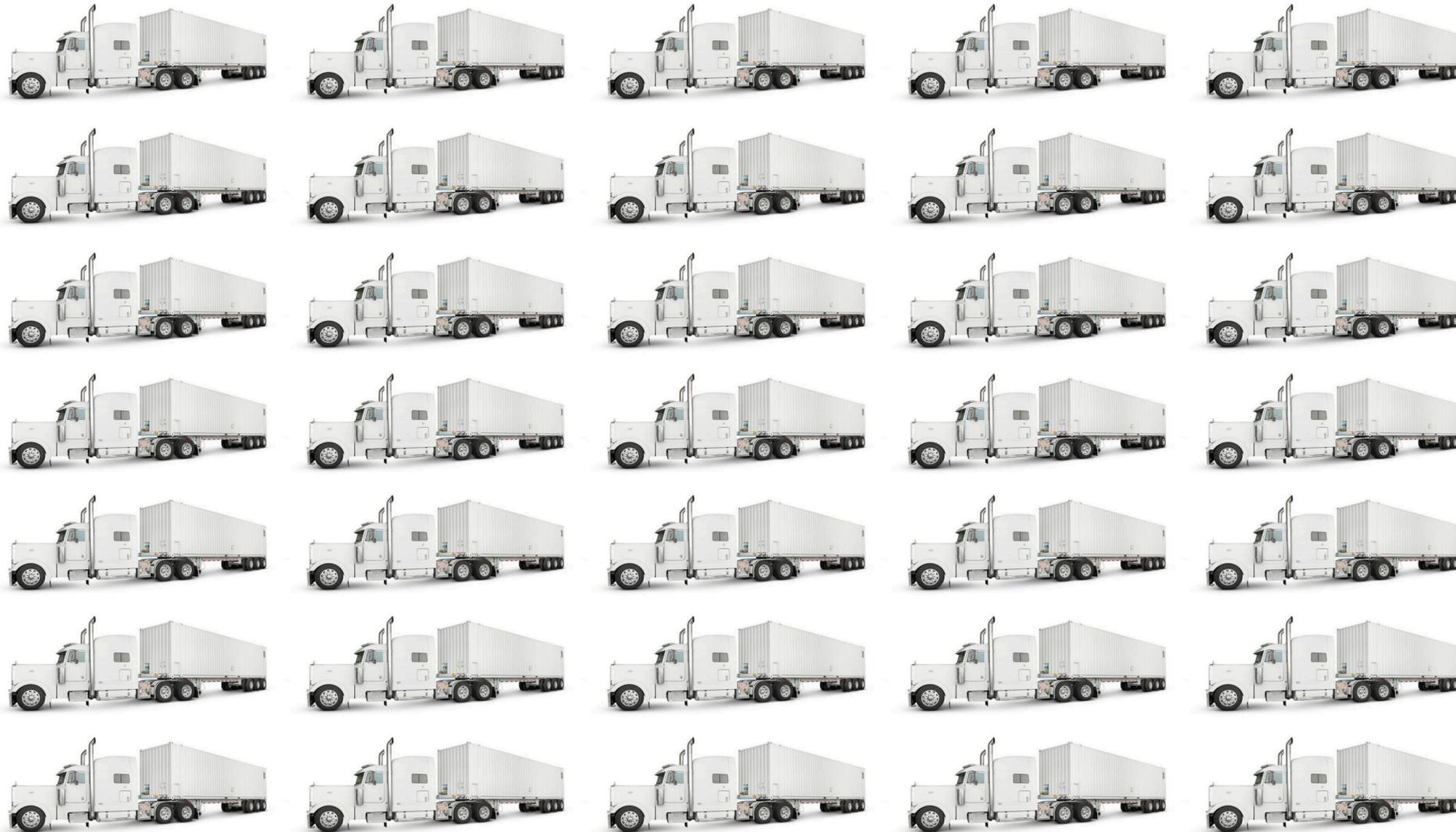
# Introduction

- We can't save everything
  - Reduce data rate: Triggers, save only the interesting events. storage ~ trigger-rate x luminosity
  - Reduce per event size: Save only high level observables storage ~ (reduced size/full size) x luminosity

# Introduction

- We can't save everything
  - Reduce data rate: Triggers, save only the interesting events.
    storage ~ trigger-rate x luminosity
  - Reduce per event size: Save only high level observables
    storage ~ (reduced size/full size) x luminosity
  - Save only histograms of observables
    storage ~ number of bins, lose not specified correlations

# Introduction

- We can't save everything
  - Reduce data rate: Triggers, save only the interesting events.
    storage ~ trigger-rate x luminosity
  - Reduce per event size: Save only high level observables
    storage ~ (reduced size/full size) x luminosity
  - Save only histograms of observables
    storage ~ number of bins, lose not specified correlations
  - Encode data in generative model
    storage ~ network weights

# Introduction

**LHC data**

**Trained model**

# Online Training

This is quite ambitious

# Online Training

This is quite ambitious

- Lets start smaller

# Online Training

This is quite ambitious

- Lets start smaller

- Online model as additional scouting tool
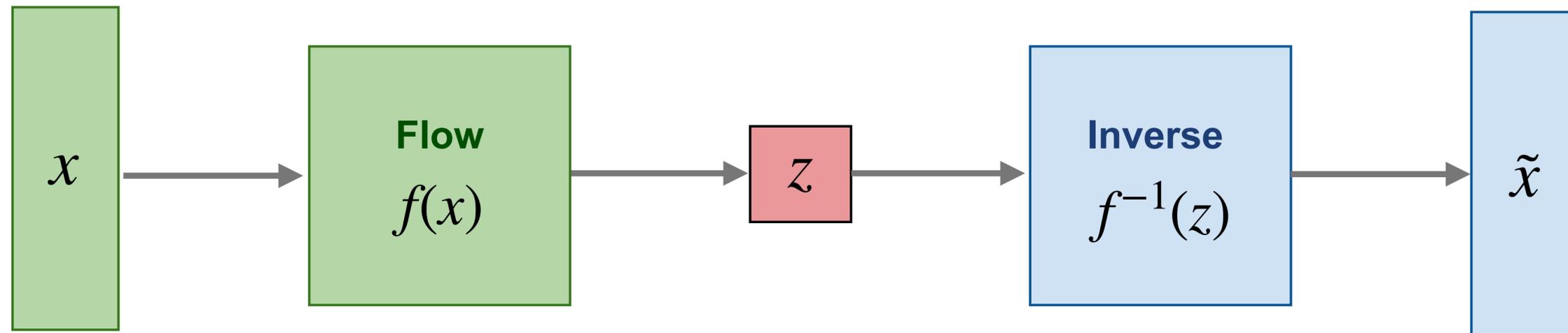
- Investigate regions ignored by triggers

# Proof of Concept

- 1 Dimension toy mass spectrum

- Exponential falling background

- Gaussian peak signal

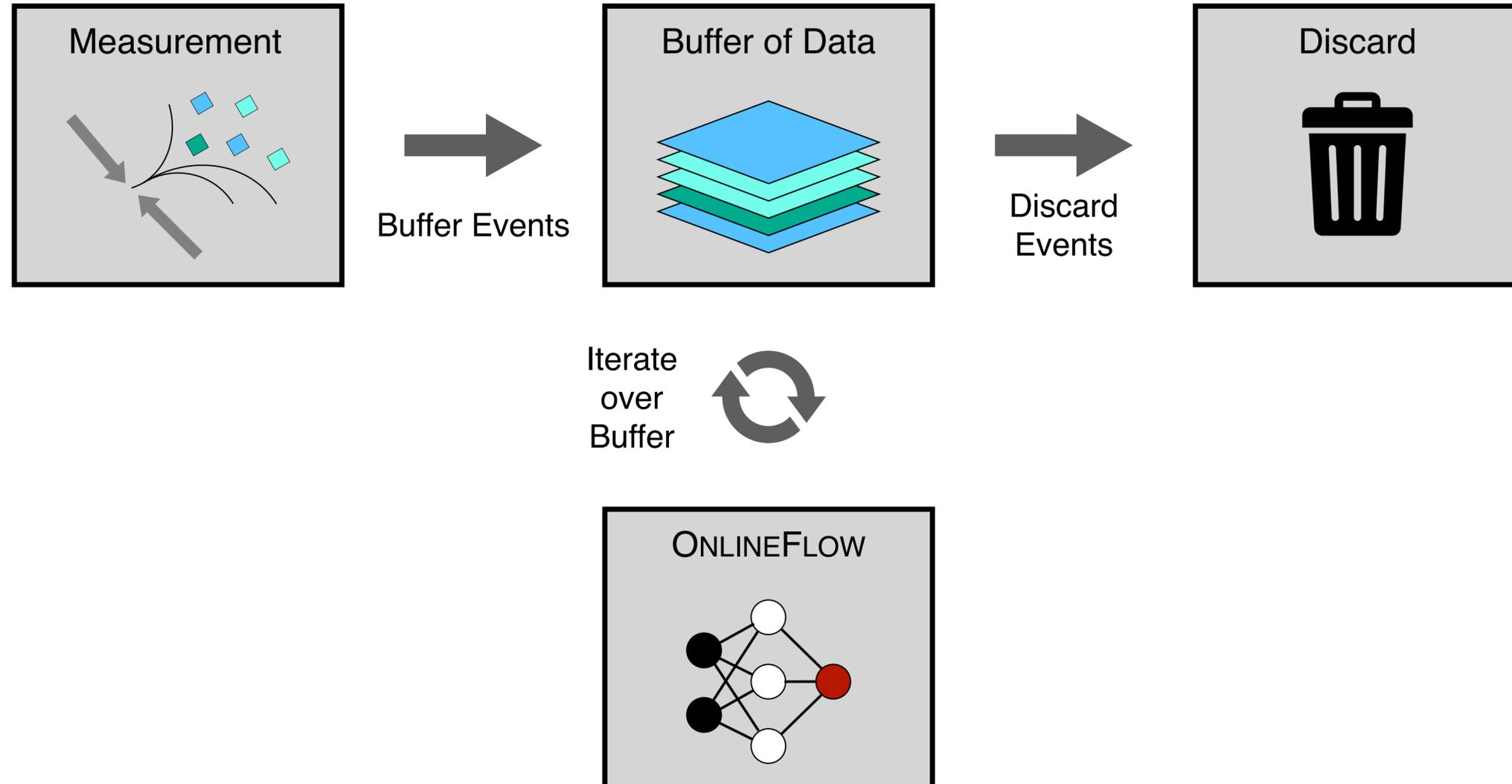- Train generative model

# Generative Model

Normalising Flow



- Train invertible transform to map data to latent space
- Use inverse to generate new data form new latent samples

# Online Training

Troubles with Online Training:

- Models designed for parallel training on batches
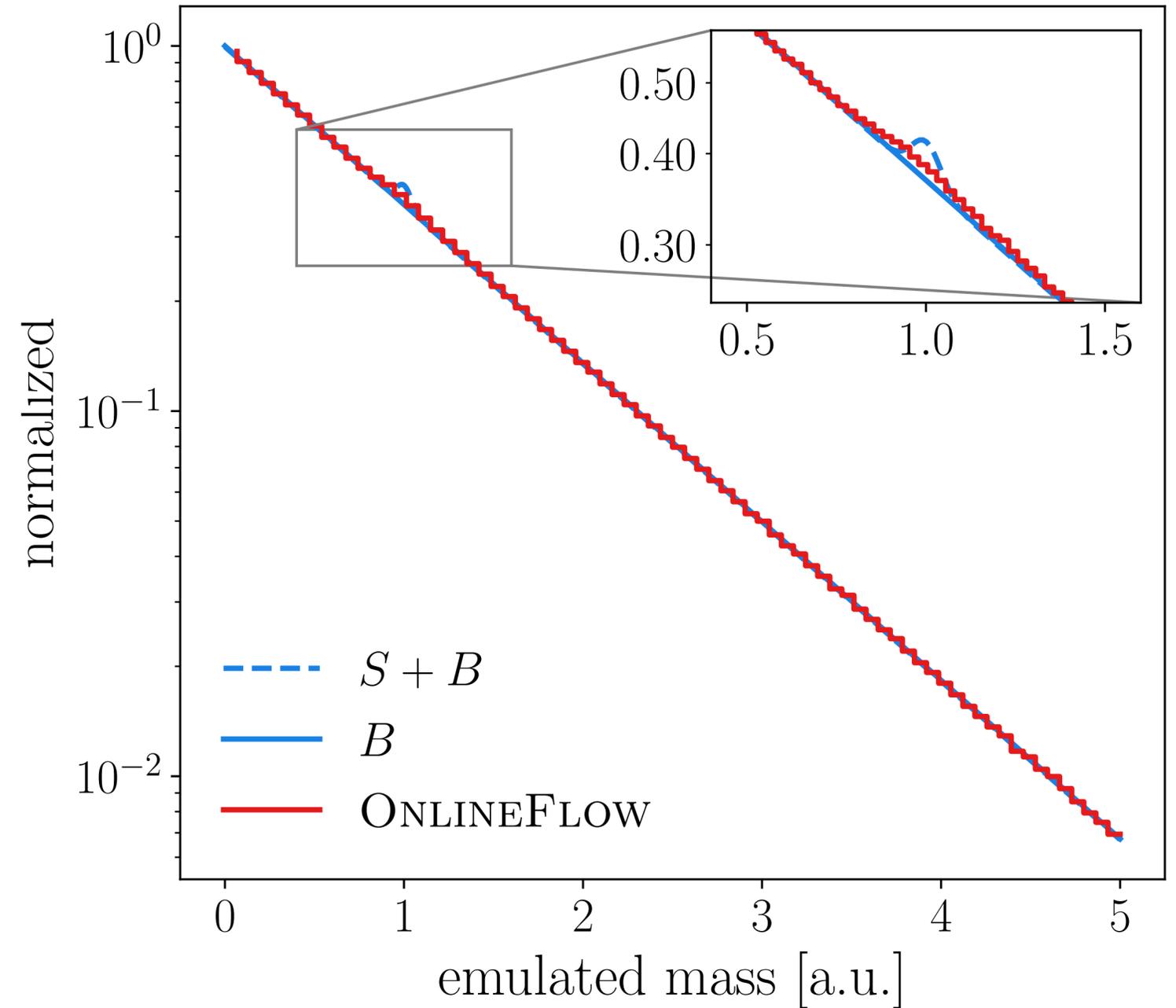
- Points should be seen more than once

# Online Training

Model will be biased towards more recent buffers

- Stochastic Weight Averaging (SWA)
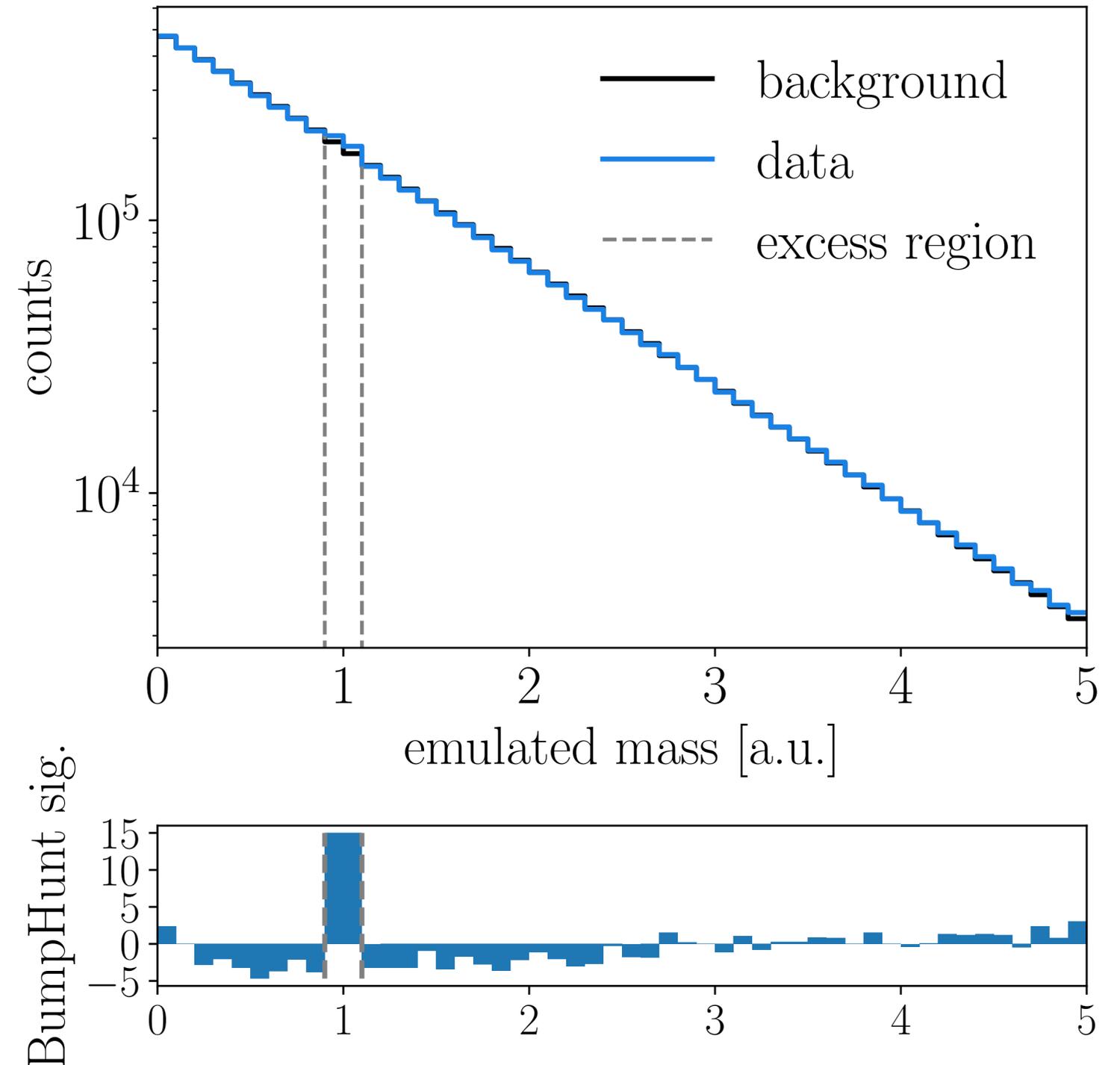
- Keep running average over model during training

# Proof of Concept

- 1 Dimension

- Exponential falling background

- Gaussian peak signal

- Train masked autoregressive Flow on S+B

- Check samples form Flow

# Proof of Concept

- Analyse data

- Fit exponential function to data to get background model

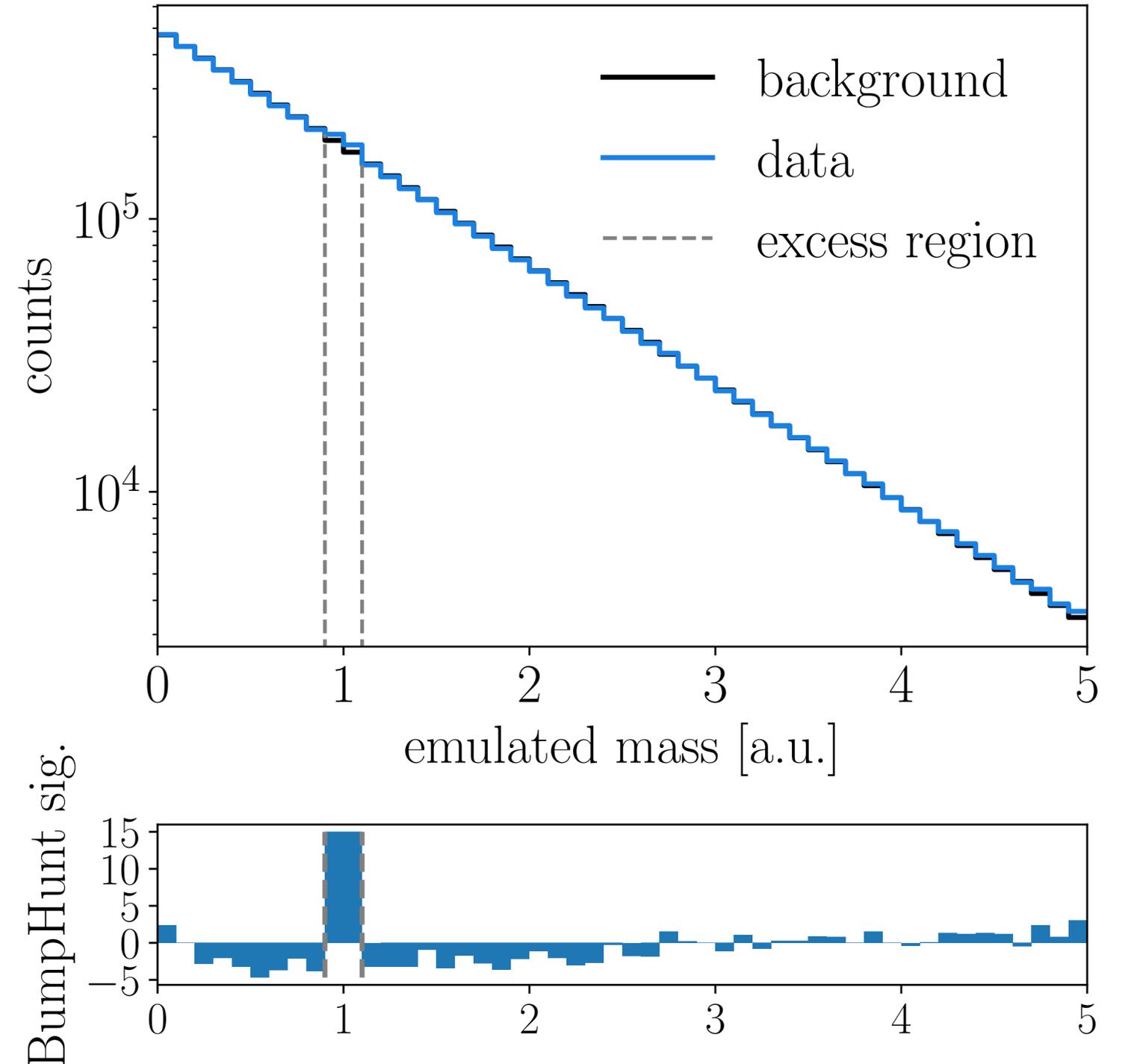- Run pyPumpHunter (python implementation The BumpHunter)

- What is our significance?
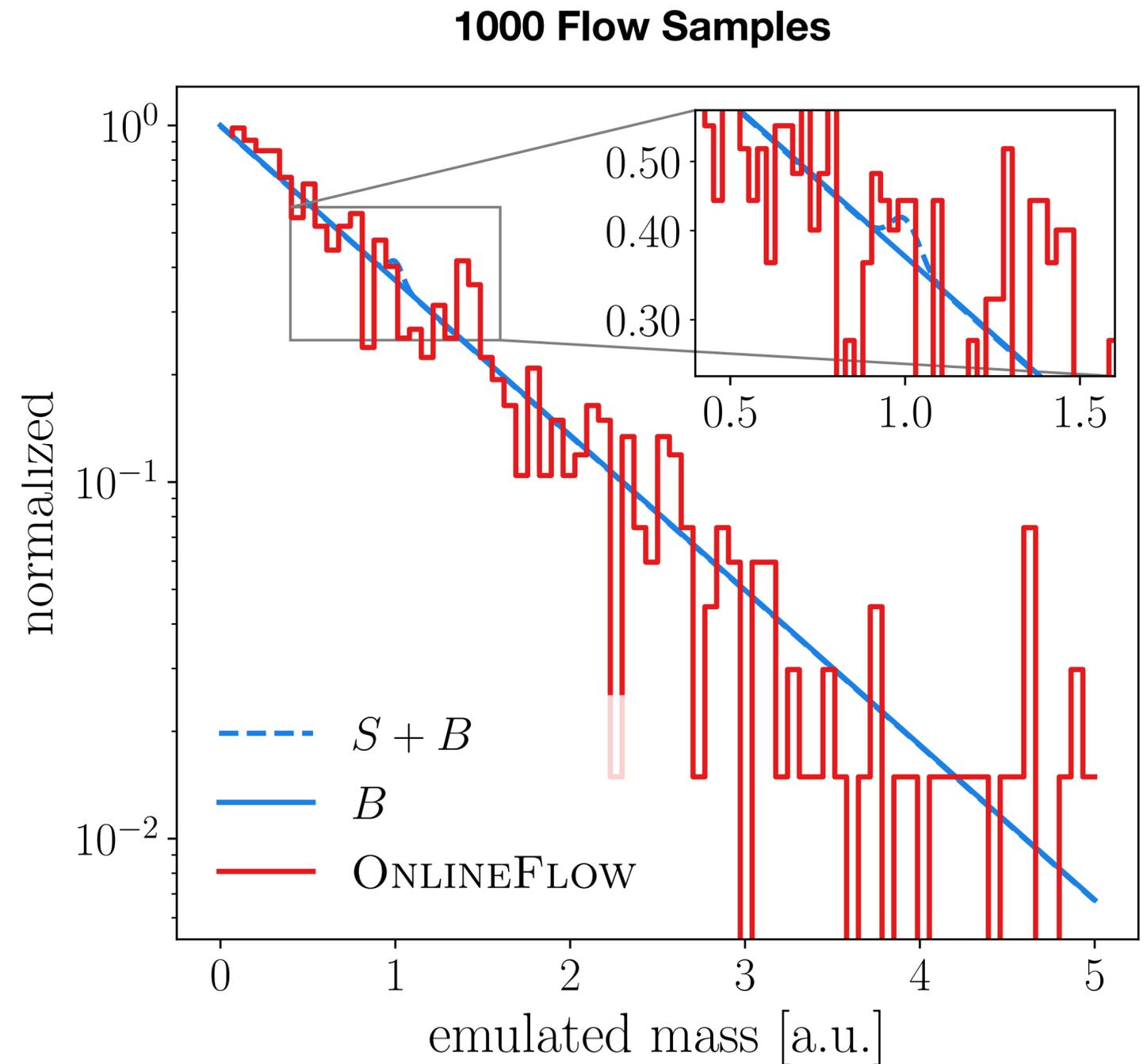


Louis Vaslin, Julien Donini, **pyBumpHunter**

# Proof of Concept

- What is our significance?
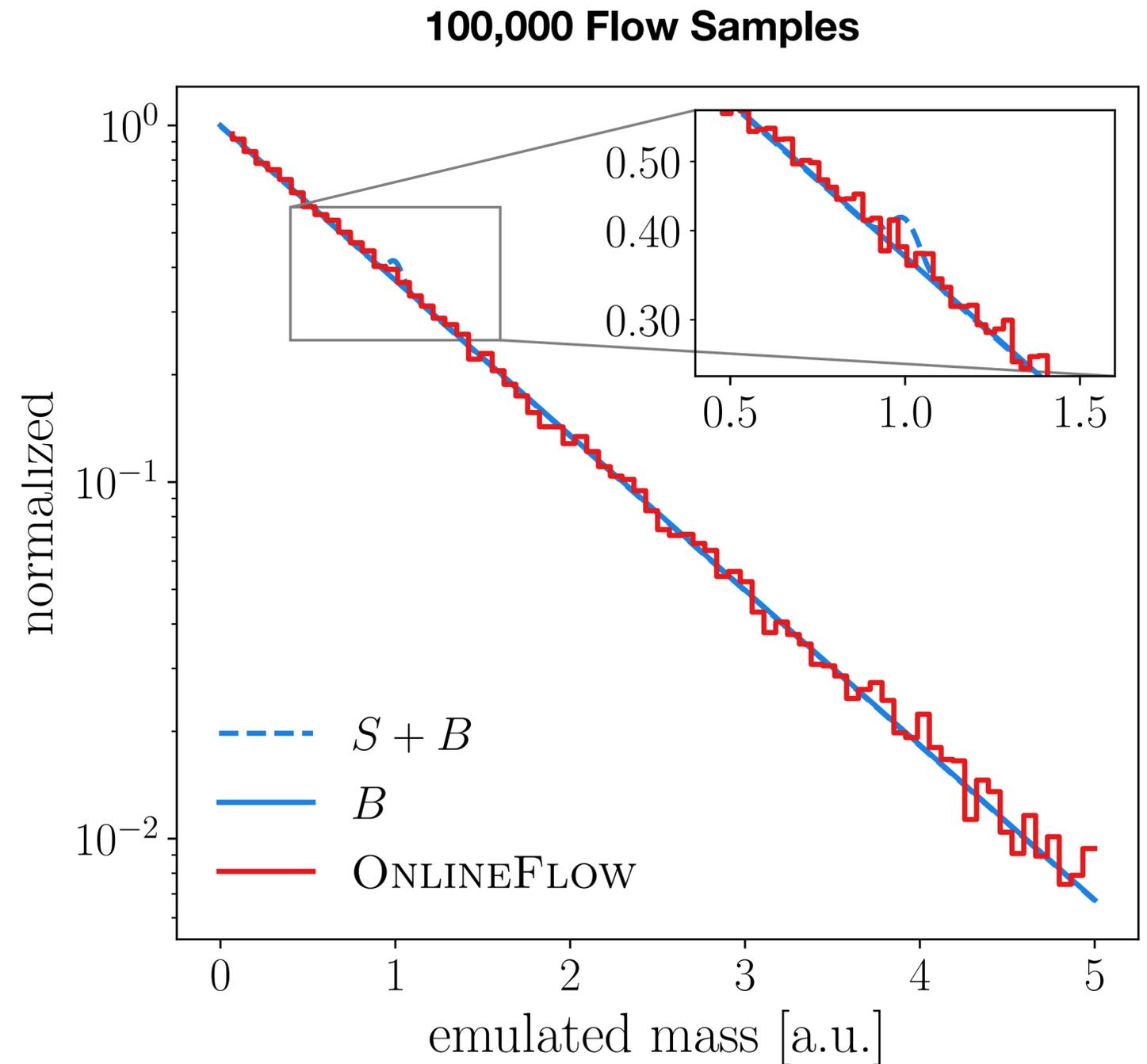  - Classic assumption $\frac{S}{\sqrt{B}}$
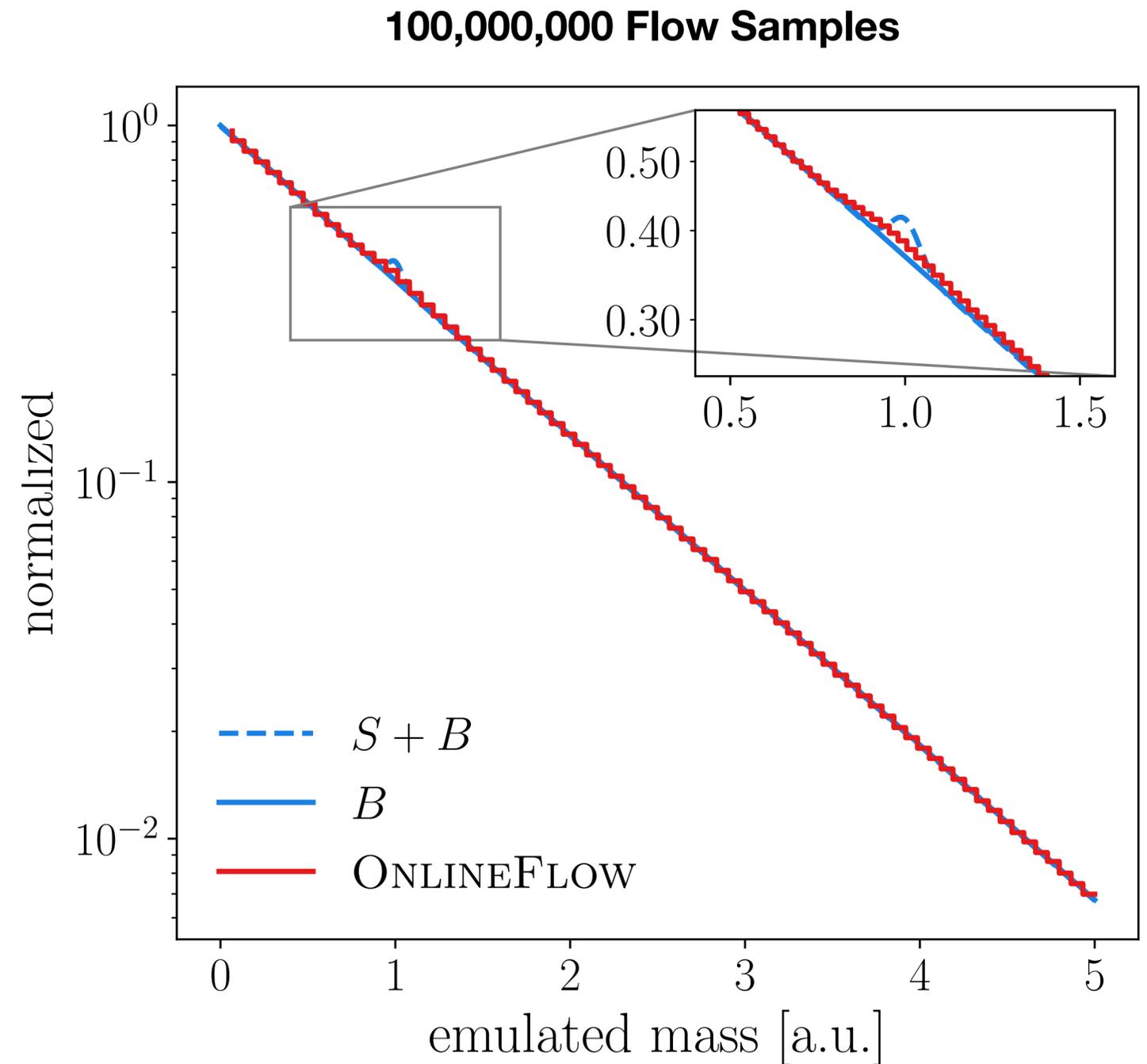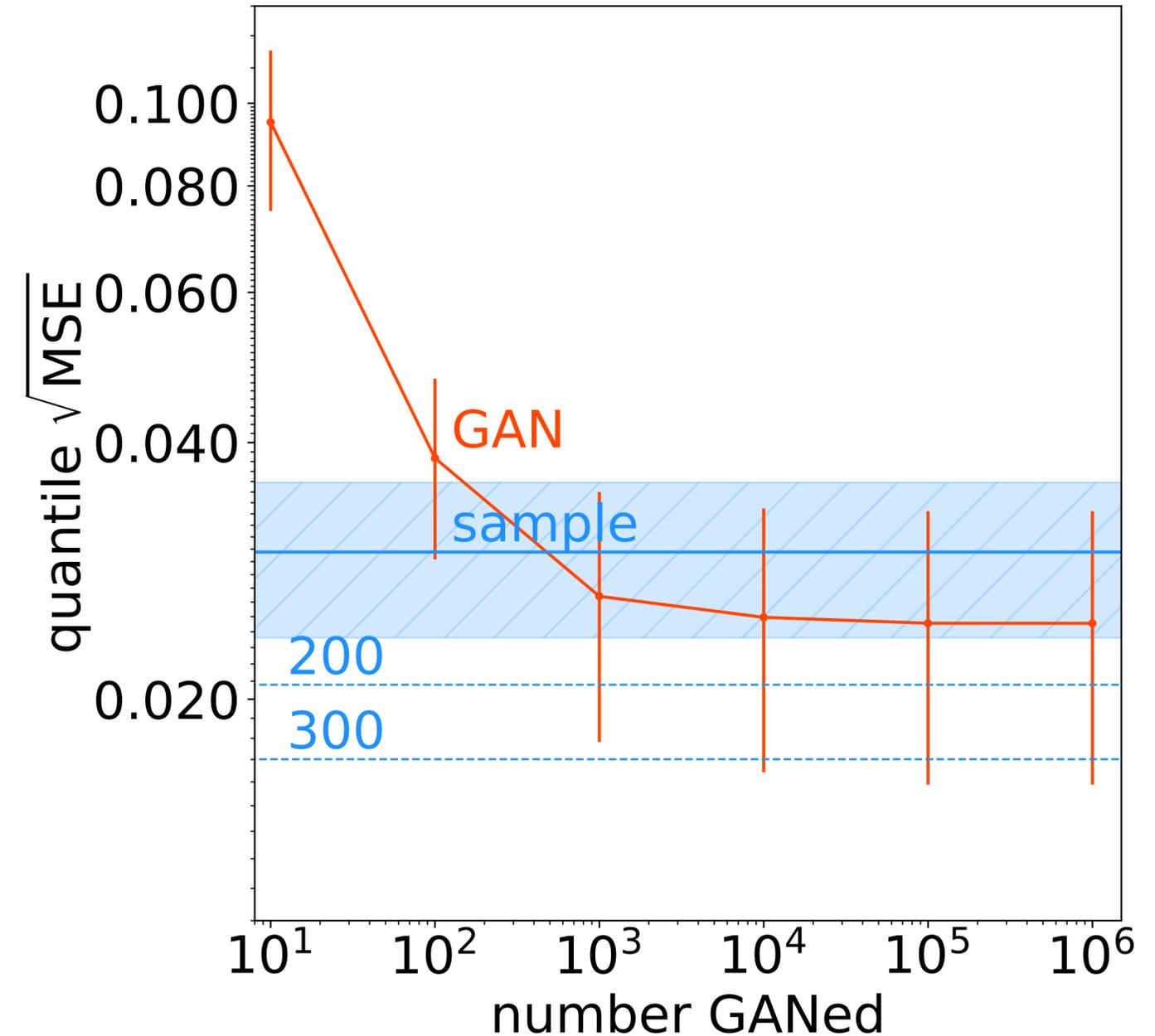  - Generative models break this

# Proof of Concept

- What is our significance?
  - Classic assumption $\dfrac{S}{\sqrt{B}}$
  - Generative models break this
    - Fluctuations if few points



**1000 Flow Samples**

Legend:
- $S + B$ (dashed)
- $B$ (solid)
- ONLINEFLOW

Axes: normalized (y), emulated mass [a.u.] (x)

# Proof of Concept

- What is our significance?
  - Classic assumption $\frac{S}{\sqrt{B}}$
  - Generative models break this
    - Fluctuations if few points
    - More points reduce this



**100,000 Flow Samples**

Legend:
- $S + B$ (dashed blue)
- $B$ (blue)
- OnlineFlow (red)

x-axis: emulated mass [a.u.]
y-axis: normalized

# Proof of Concept

- What is our significance?
  - Classic assumption $\dfrac{S}{\sqrt{B}}$
  - Generative models break this
    - Fluctuations if few points
    - More points reduce this

# Proof of Concept

- What is our significance?
  - Classic assumption $\dfrac{S}{\sqrt{B}}$
  - Generative models break this
    - Fluctuations if few points
    - More points reduce this
    - More samples $\neq$ smaller error
  - No Poisson error
  - Systematic errors dominant



Butter et al.: **Amplifying Statistics using Generative Models:** NeurIPS ML4PS 2020, 2008.06545

# Proof of Concept

- Train **bootstrap ensemble of Flows**

  - Pass resampled version of dataset to each flow

  - Emulate online with Poisson weight for each event

- Compare flows within ensemble

  - Estimate uncertainty

# Proof of Concept

- Use half the ensemble to estimate bump position

- Use other half to determine signal rate

- Built in cross-check

- Estimate uncertainty form variance

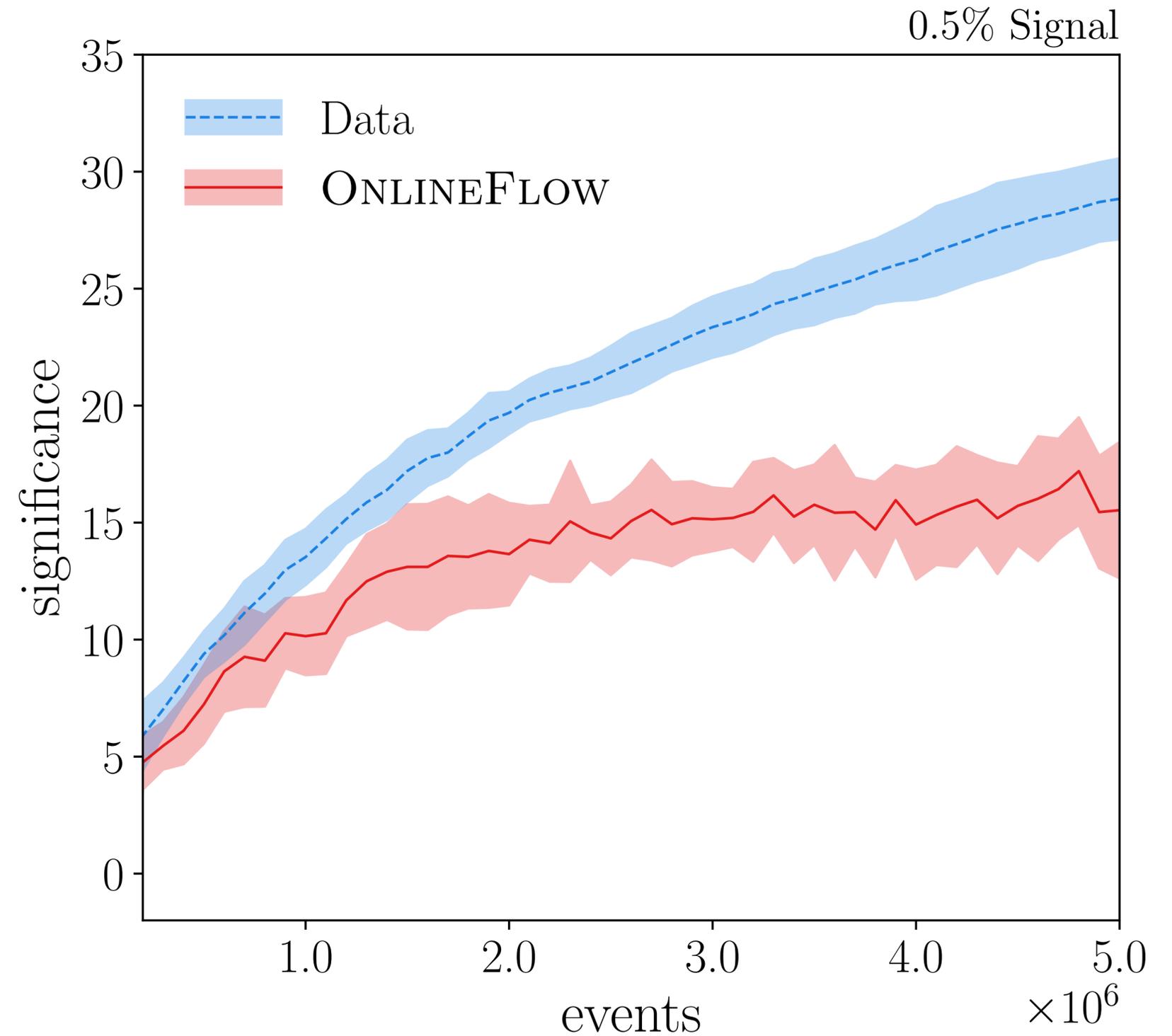- Plot as function of the number of events to see development

# Proof of Concept

- Run same analysis in the training data itself

- Compare significance

- Flow similar shape as data

- Reaches high significance

# Proof of Concept

- False positives?

# Proof of Concept

- False positives?

- Run same training on only background without signal

- Negligible significance in direct search

- Flow significance nearly identical

# Prescale Case

Some event classes:
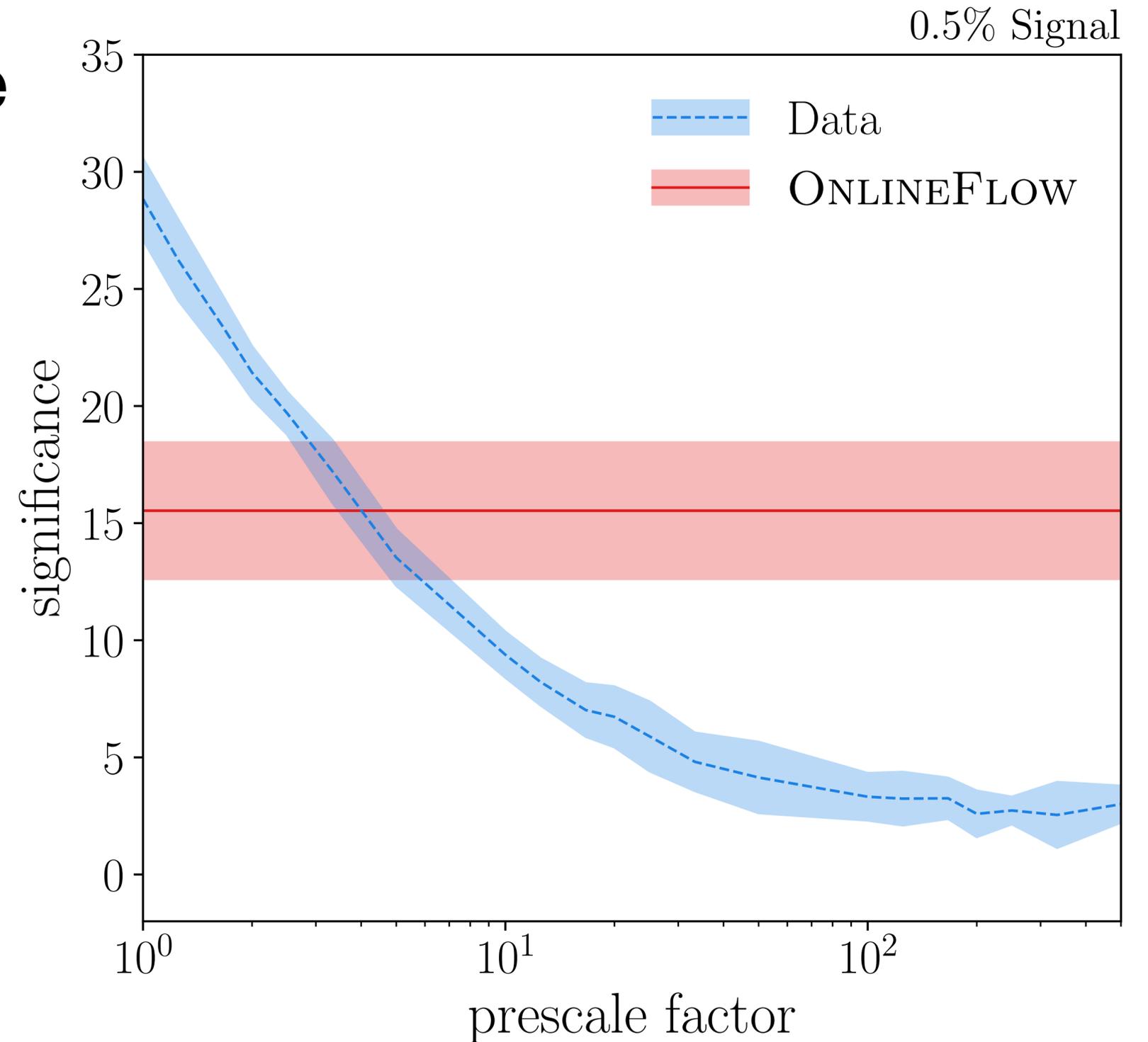
- Too high rates even after trigger

# Prescale Case

Some event classes:

- Too high rates even after trigger

- Apply percale

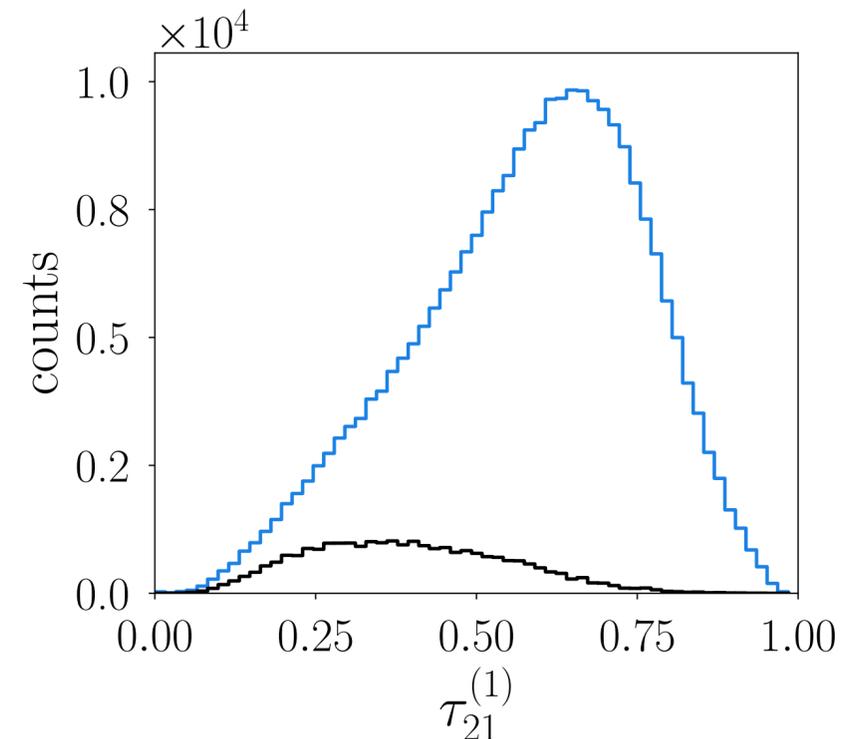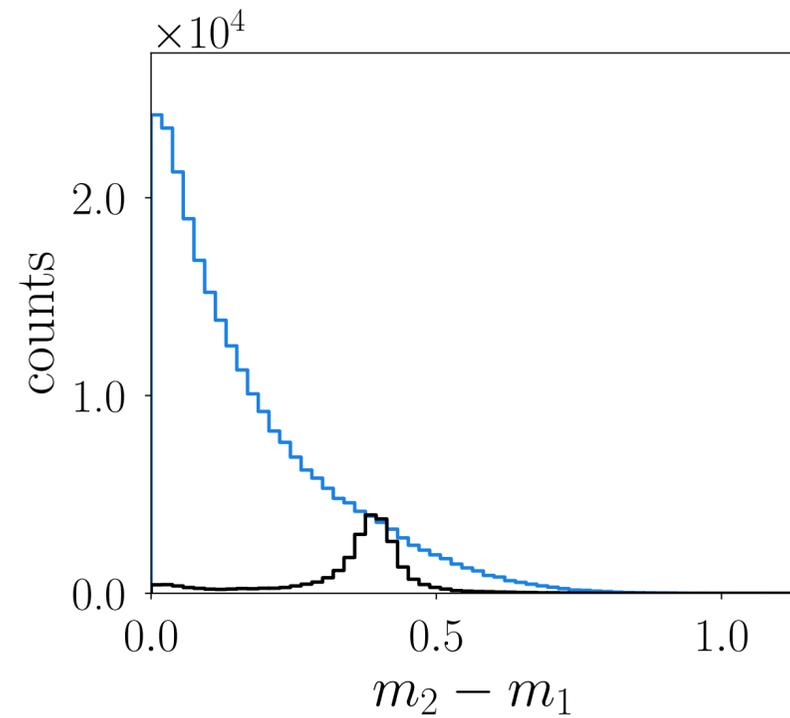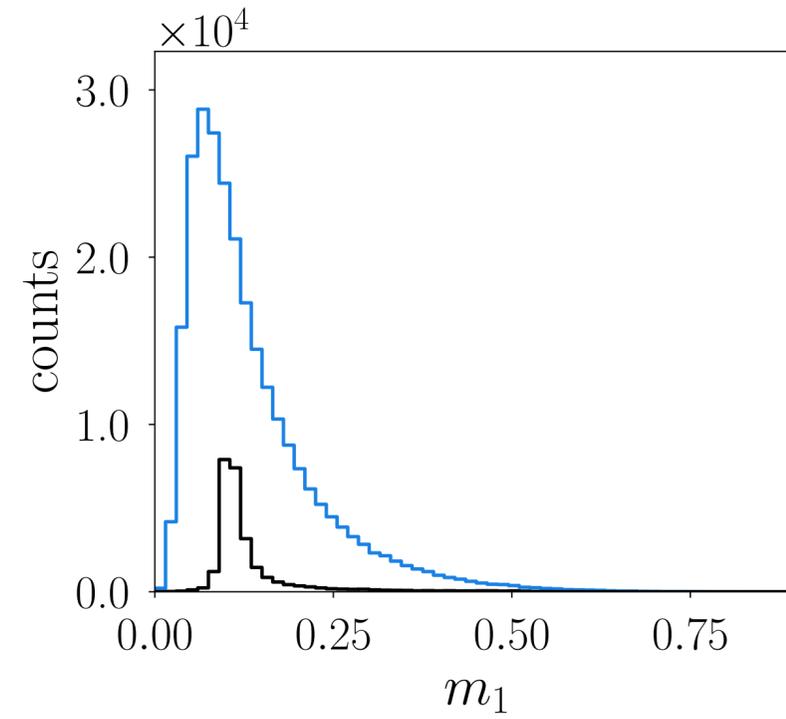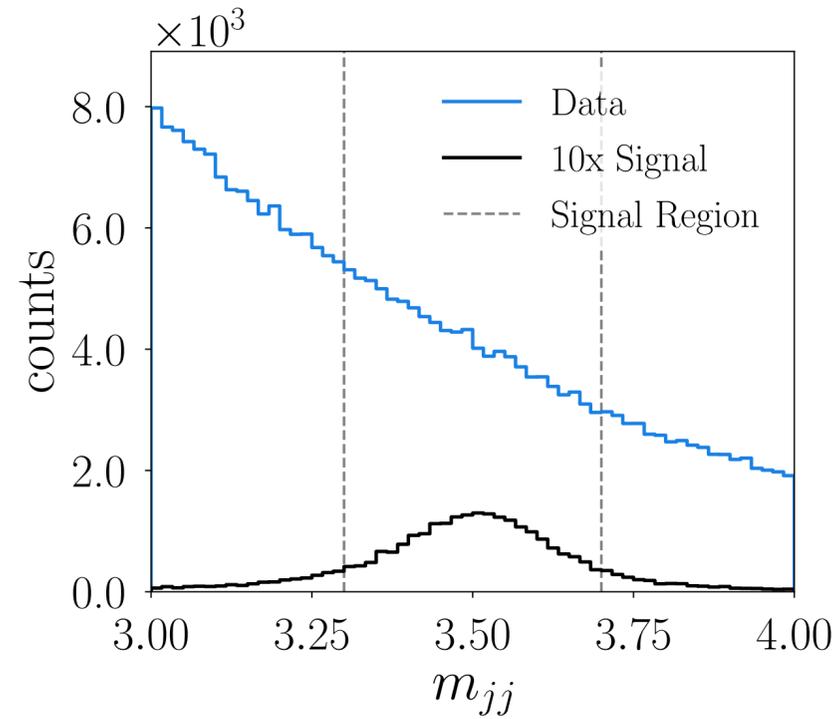- Randomly select events to save
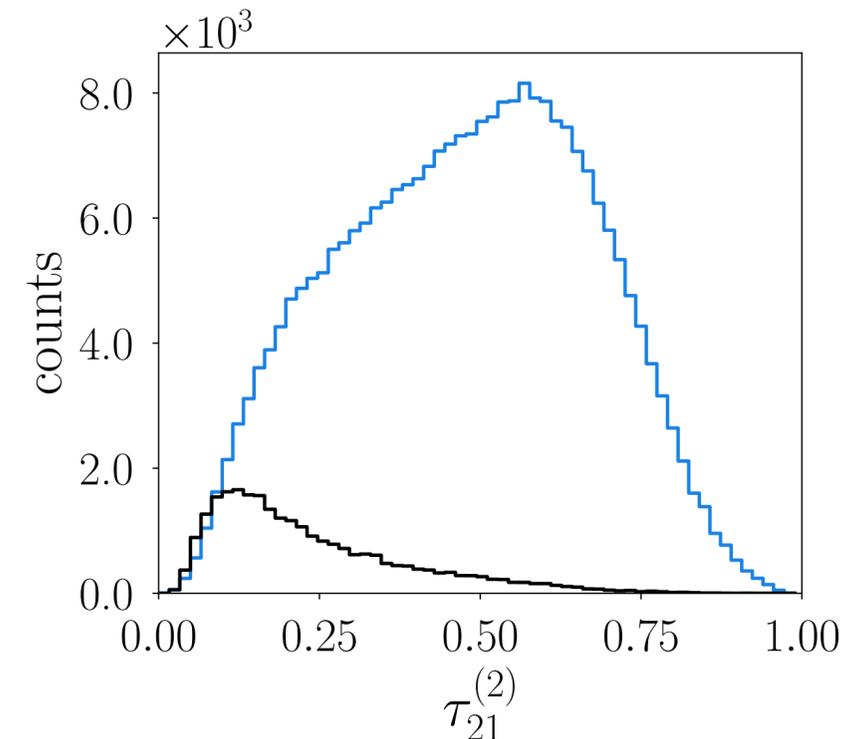
- 1/prescale chance to keep event

# Prescale Case

- At what prescale factor do we benefit from the Flow?

- Quite early

- Factor of around 4

- Very promising in 1D



0.5% Signal

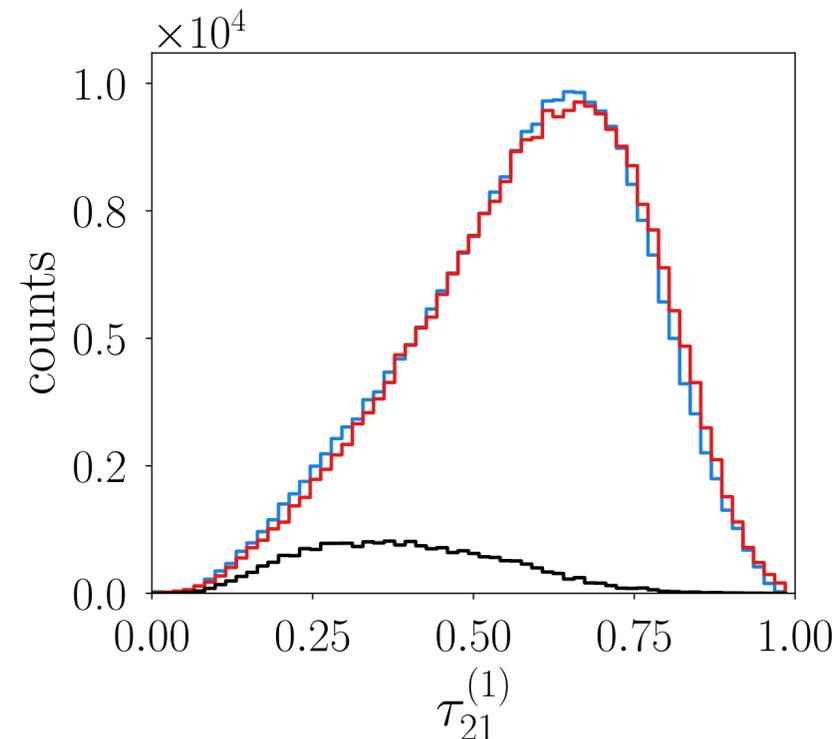significance vs prescale factor

Legend:
- Data (dashed blue)
- ONLINEFLOW (red)

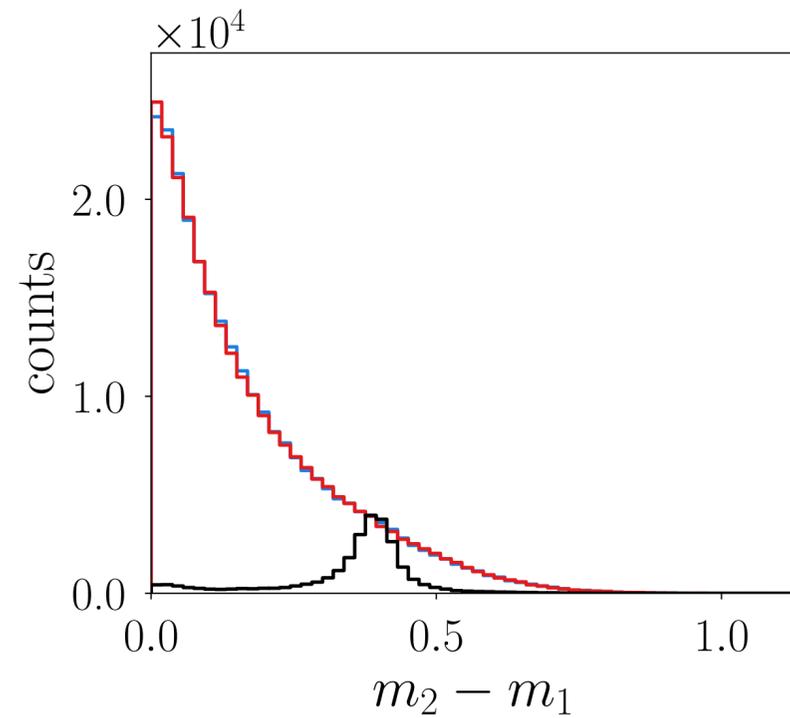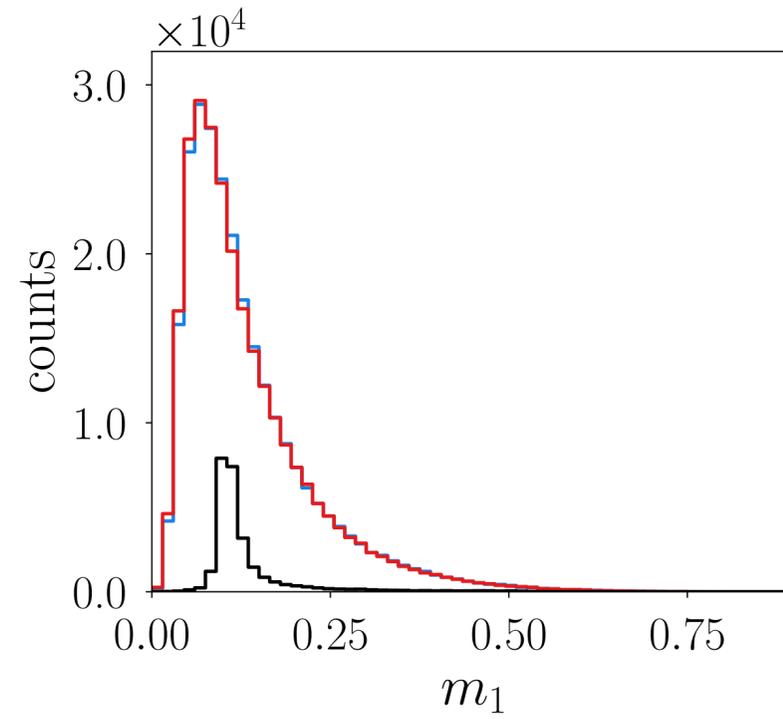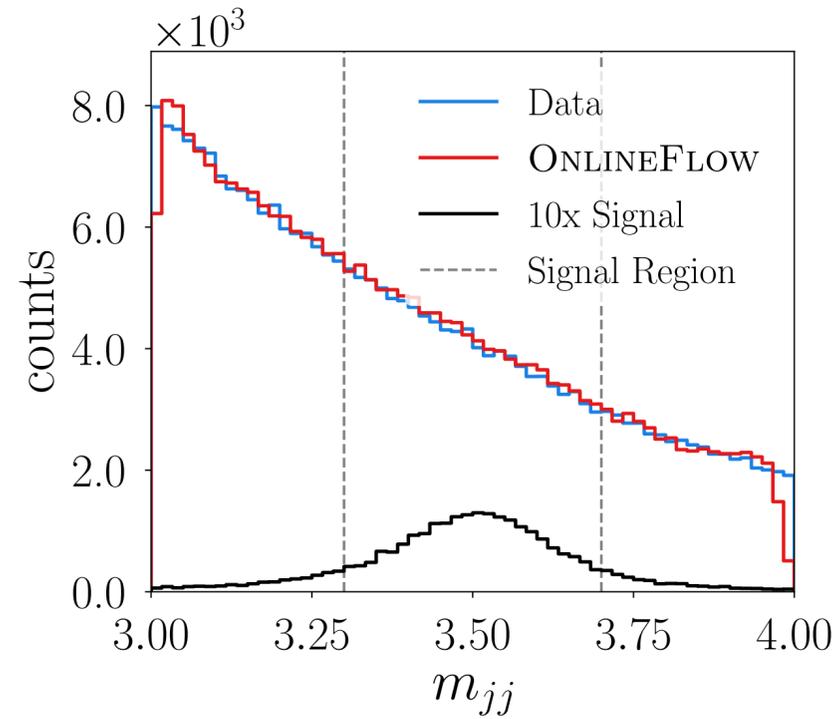# Higher Dimensional Case



- LHCO Challenge dataset

- Realistic physics case for anomaly detection

- Use state-of-the-art anomaly detection input format

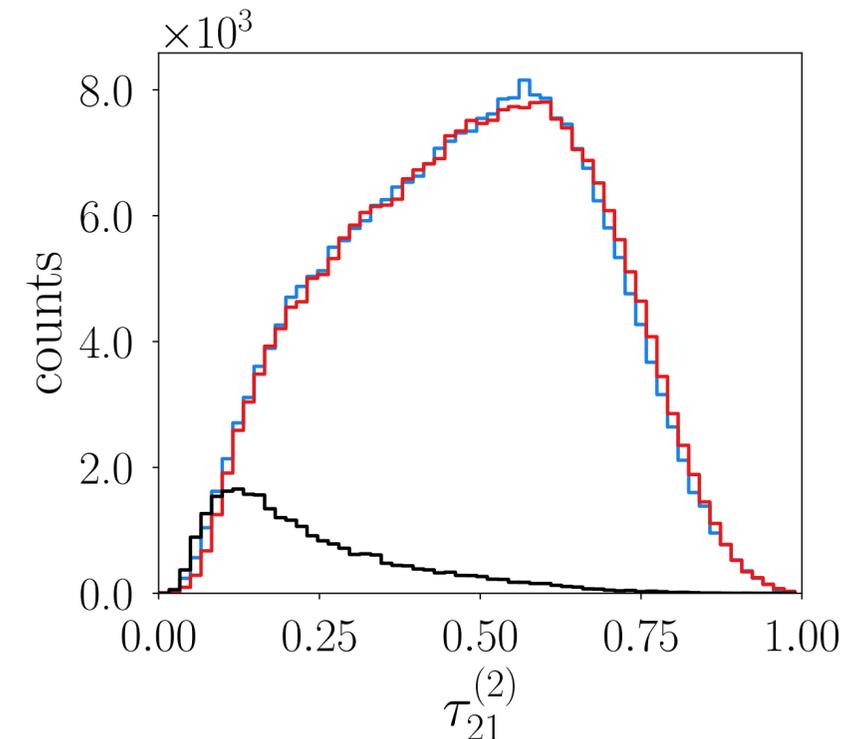Kasieczka et al.: **The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics** kl
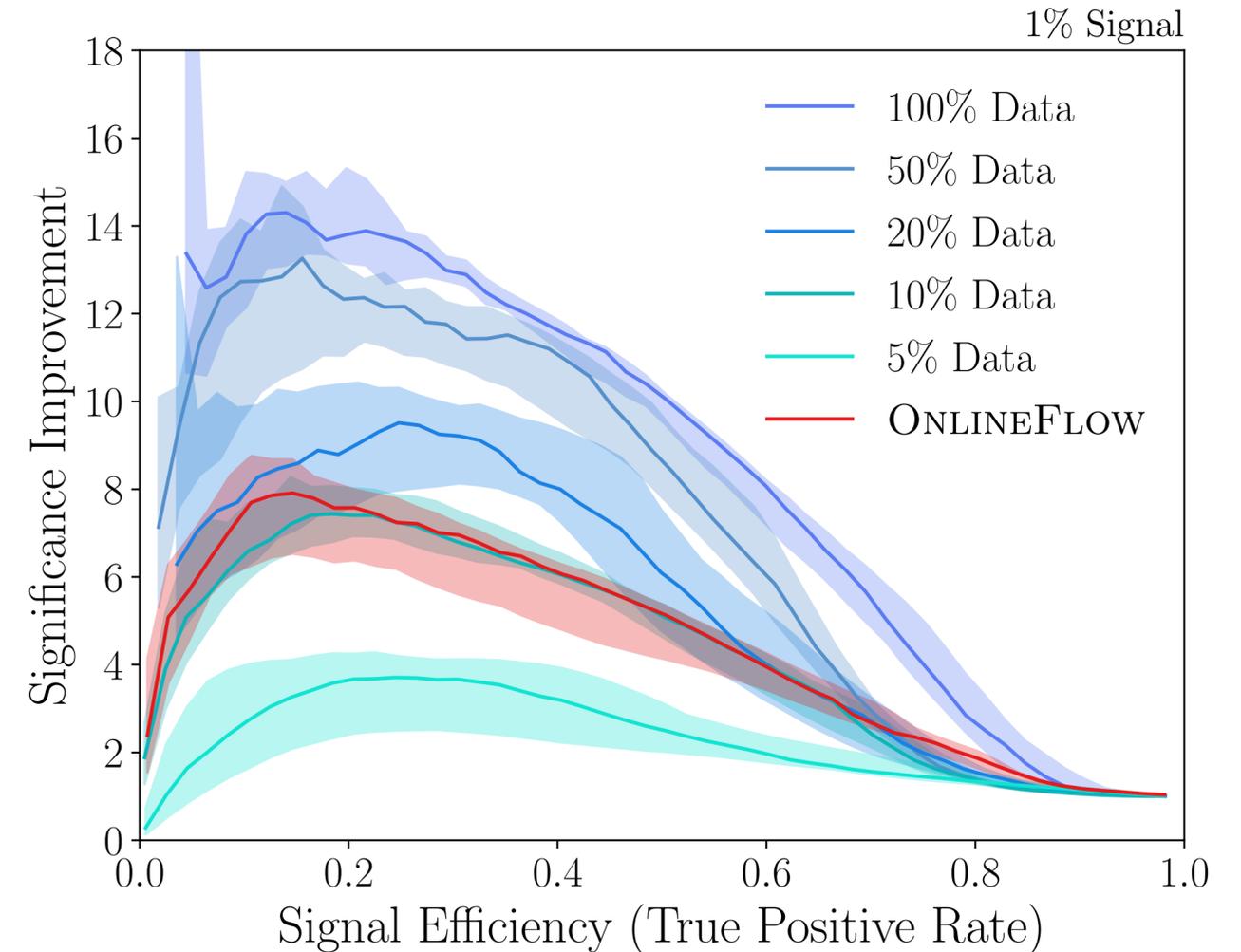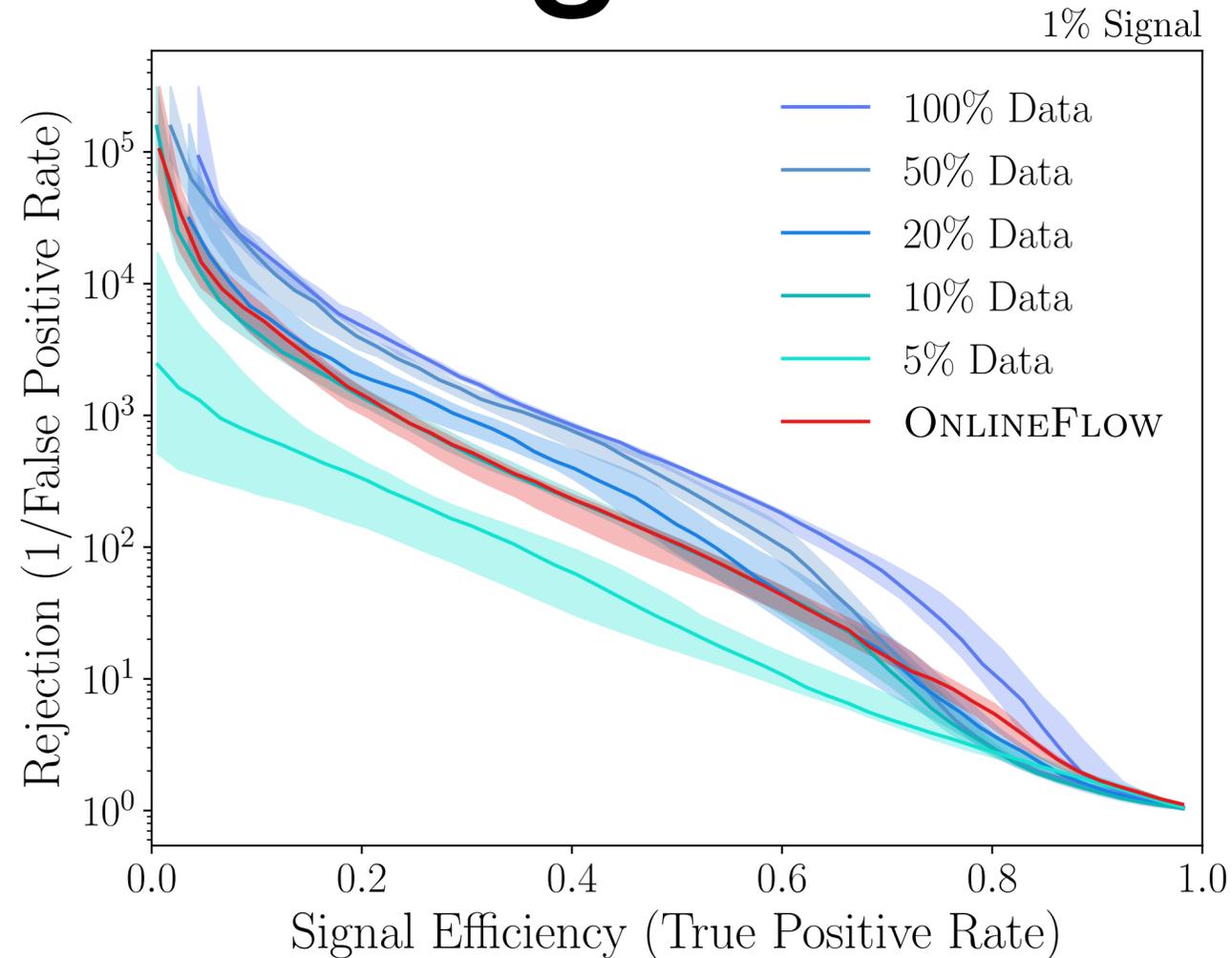2101.08320

# Higher Dimensional Case



- Train flow on LHCO data

- Run anomaly detection setup on flow data

- Compare to running on training data itself

Kasieczka et al.: **The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics** kl
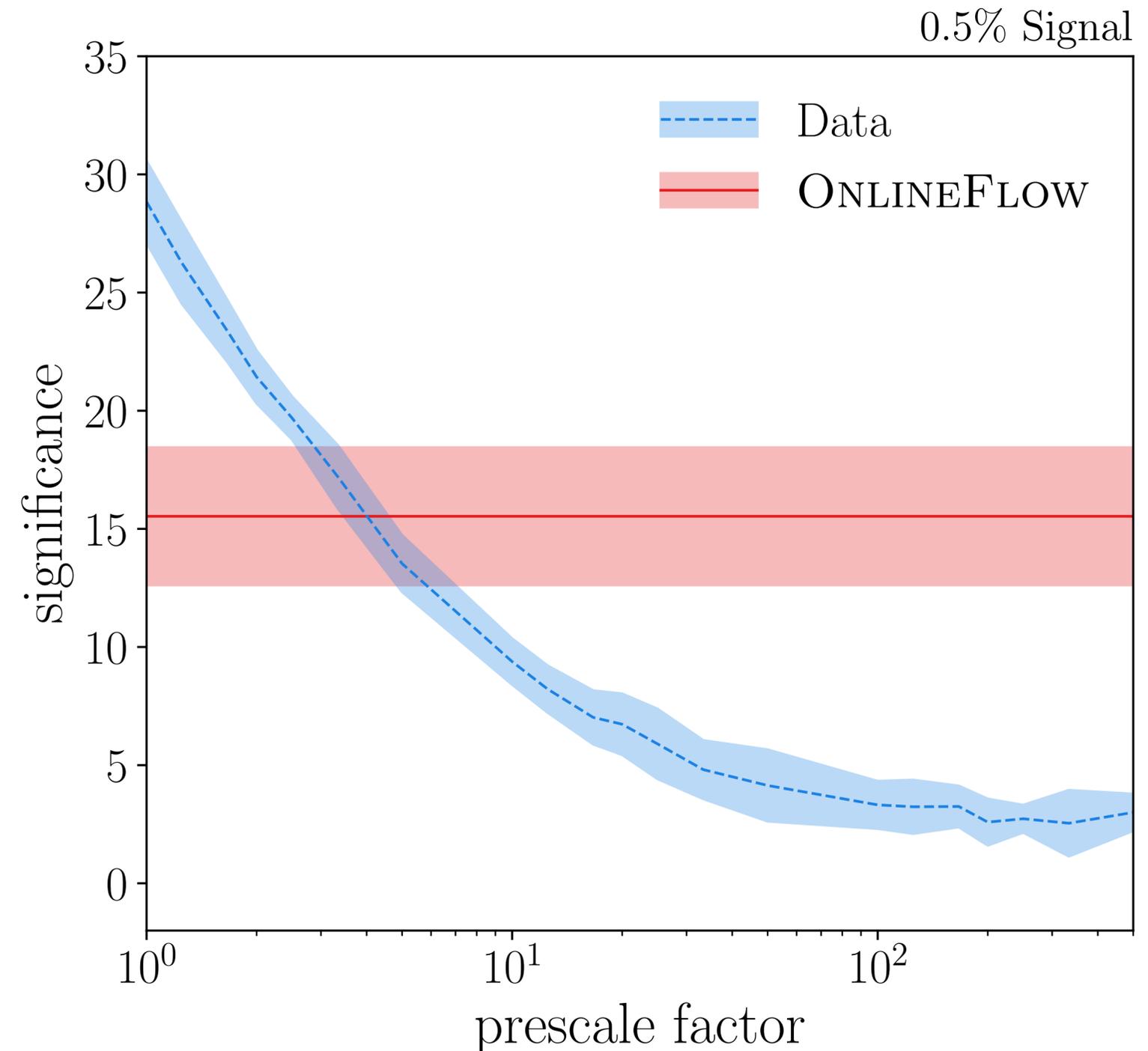2101.08320

# Higher Dimensional Case



- ROC and SIC curves, larger Rejection/Improvement is better

- Flow data performs about as well as 10% of the training data

- Still gain benefit if perscale factor is 10 or larger
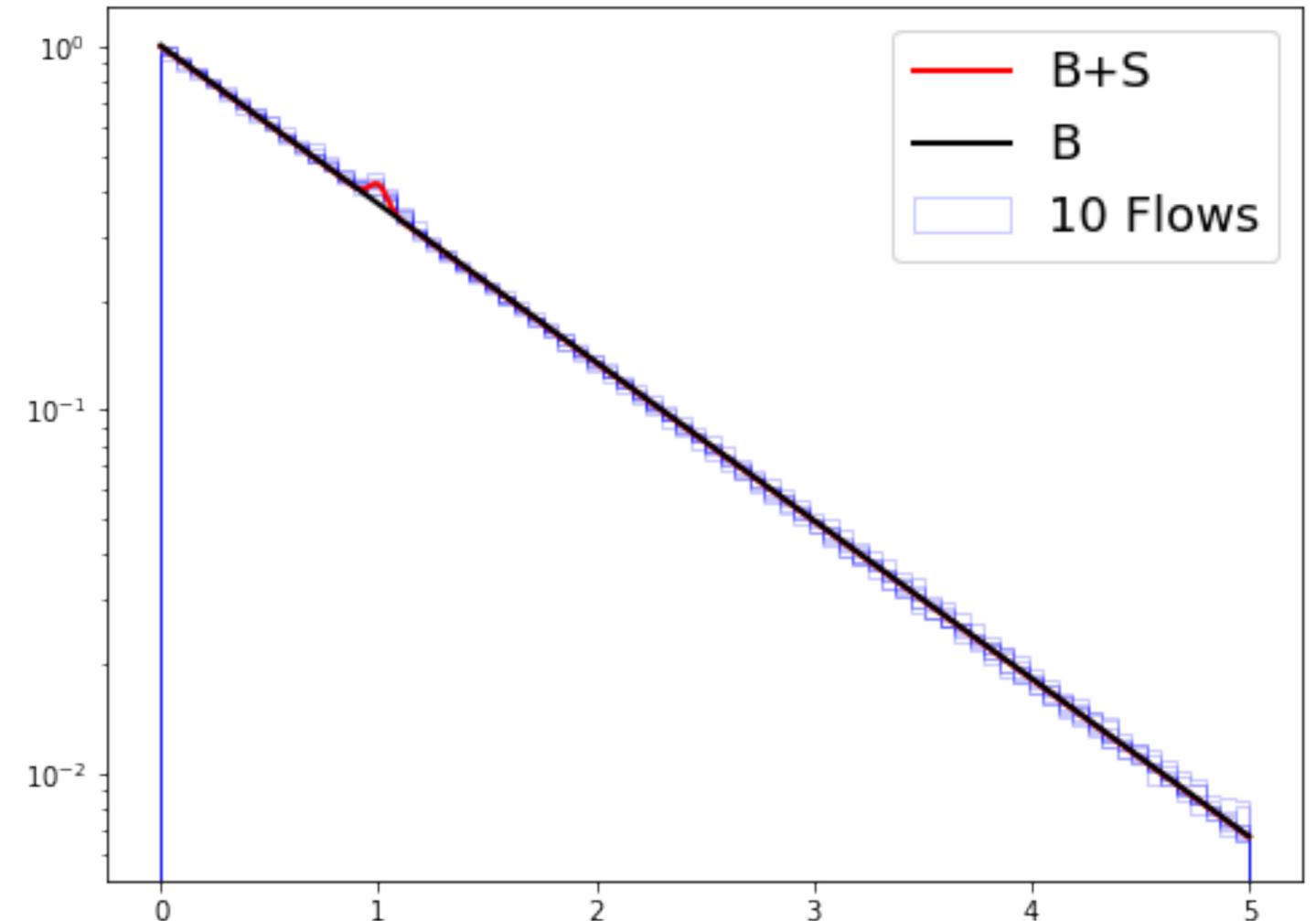
# Conclusion

- Tested OnlineFlow approach on 1D proof-of-concept data

  - Promising Results

- Higher dimensional case:

  - Can extract signal features form LHCO challenge dataset

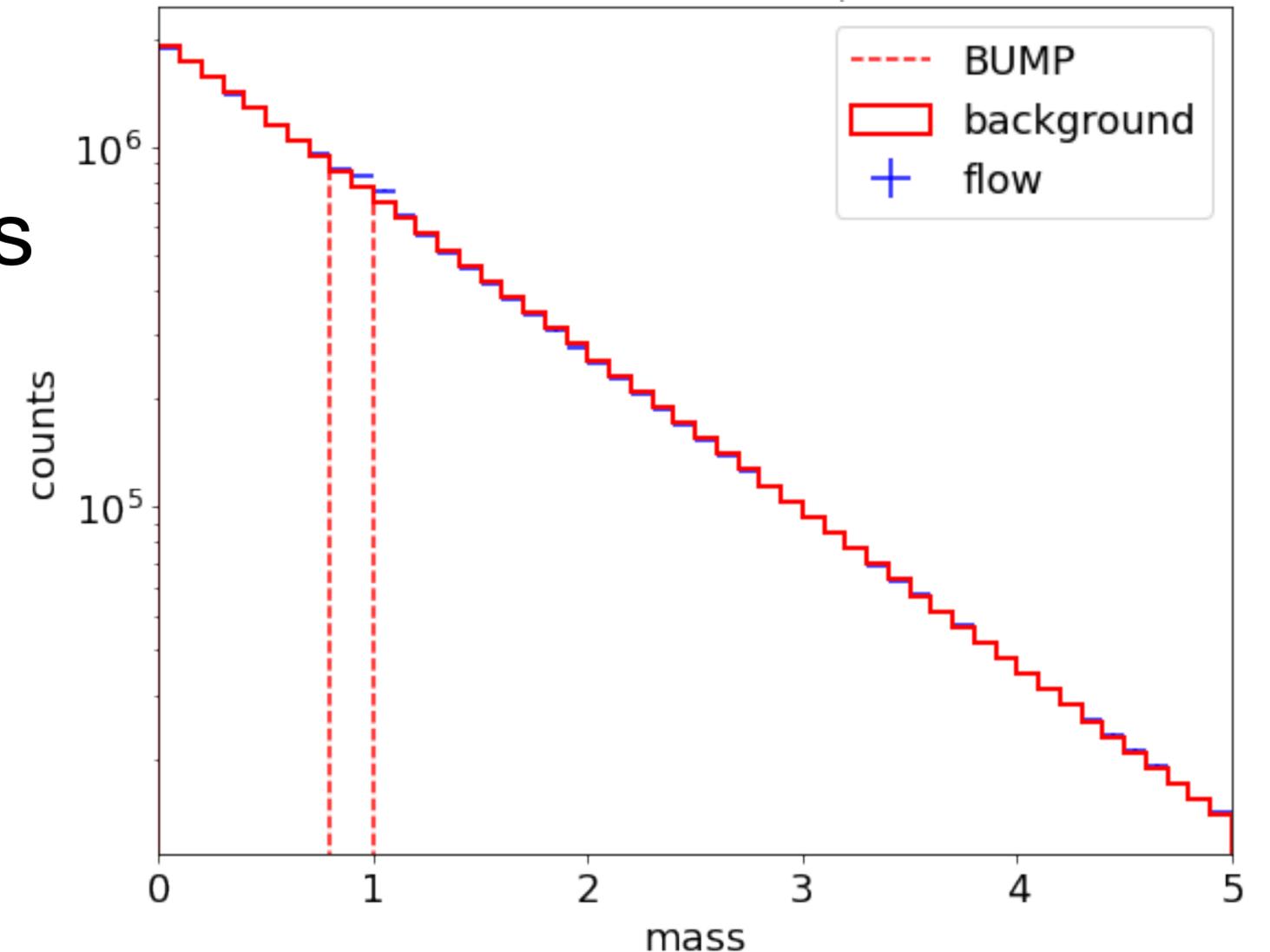- Paper (hopefully) soon!

# Thank You

# Backup

- Train **ensemble** of Flows

  - Same underlying distribution

  - Independant data points

- Lets us estimate systematic uncertainty

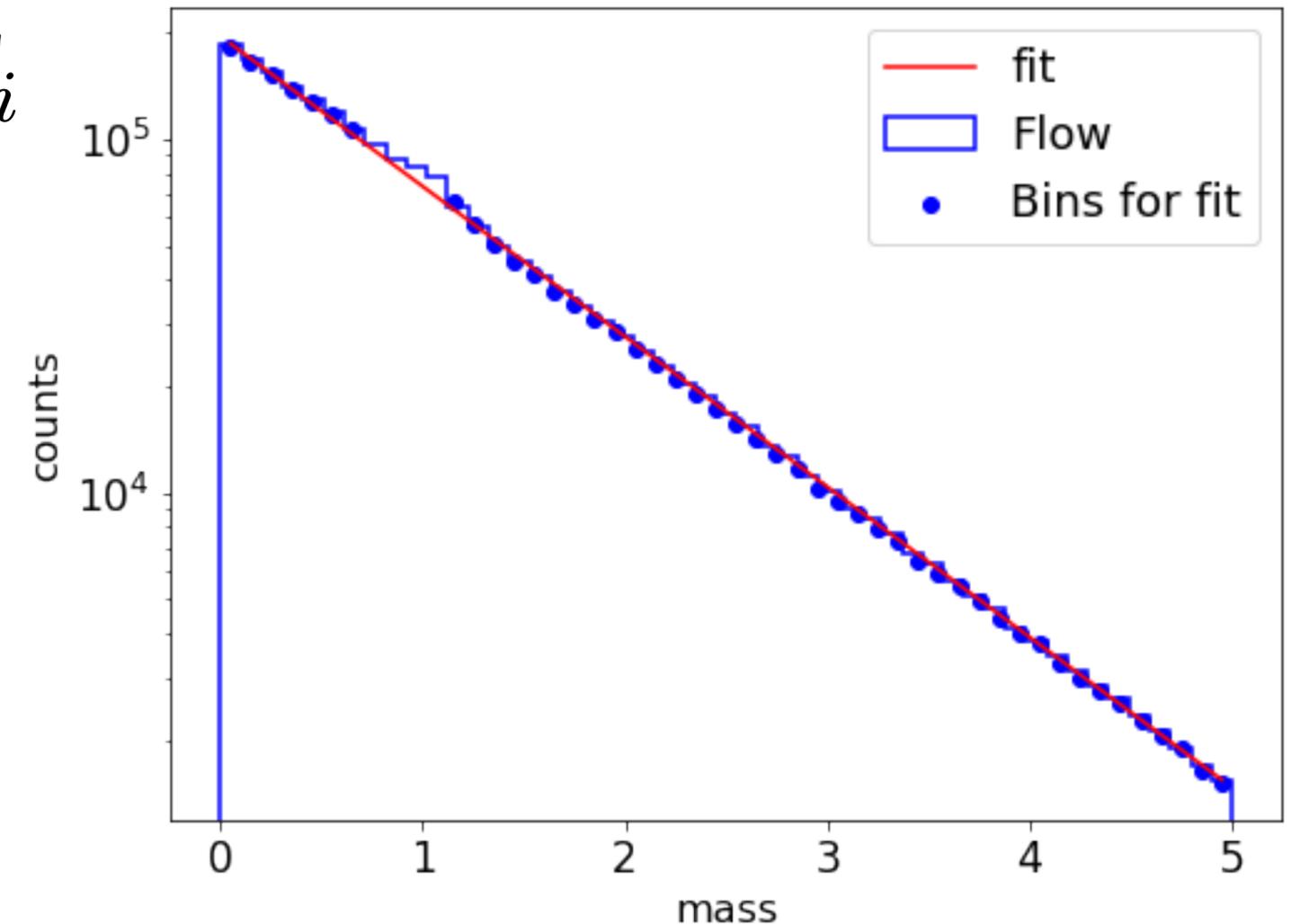- Allows for easy parallelised training

# Backup

- Uncertainty estimate:

- Step 1:

  - Look at **combined** flow samples

  - Run bump-hunt

  - Get upper and lower bound for signal region

# Backup

- Step 2:
  - Look at each **individual** flow $F_i$ out of $N$ flows in ensemble
  - Generate samples from $F_i$
  - Fit background model on samples excluding signal region from step 1
  - Get $B_i$ and $(S + B)_i$ in signal region

# Backup

- Step 3:
  - Calculate $B = \sum_{i}^{N} B_i$ and $\delta B = \mathrm{std}(B_i)\sqrt{N}$
  - Calculate $(S + B)$ and $\delta(S + B)$ equivalently
  - From this get $S$ and $\delta S$
  - Significance: $\dfrac{S}{\sqrt{(\delta S)^2 + B}}$

**Systematic error via ensemble**        **Statistical error**