

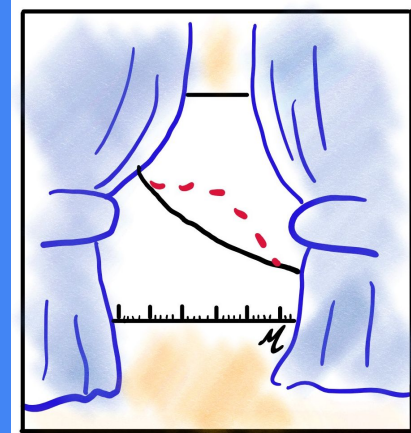
# CURTAINS

Constructing Unobserved Regions by Transforming Adjacent Intervals

ML4Jets, 3<sup>rd</sup> November 2022

Johnny Raine, Sam Klein, Debajyoti Sengupta, Tobias Golling

University of Geneva



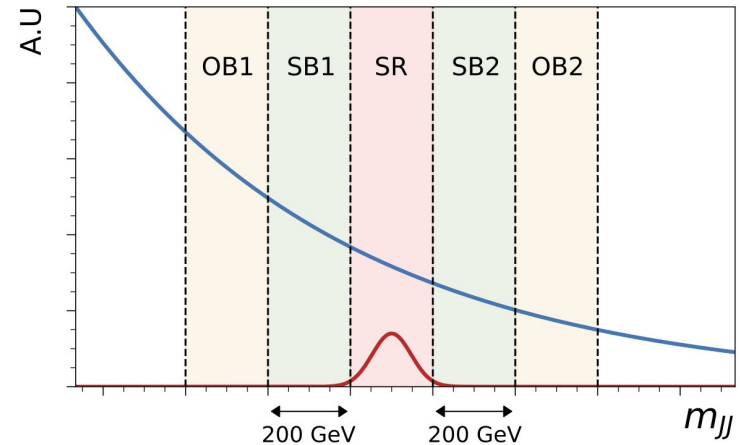
# Bump hunts

We expect signal events to be localised in the invariant mass.

⇒ Show up as a **bump** in the spectrum

Method:

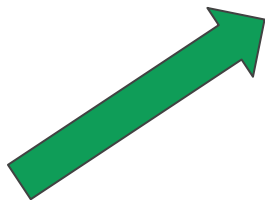
1. Split spectrum into sliding 'side bands'
2. Fit the distribution in sidebands
3. Interpolate into the signal region
4. Look for an excess
5. Slide window and repeat



# Extending Bump hunts

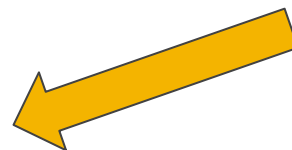
But what if the bump is dominated by background...

Is there more information we can use than just mass?



High mass - boosted decays  
Look at substructure!

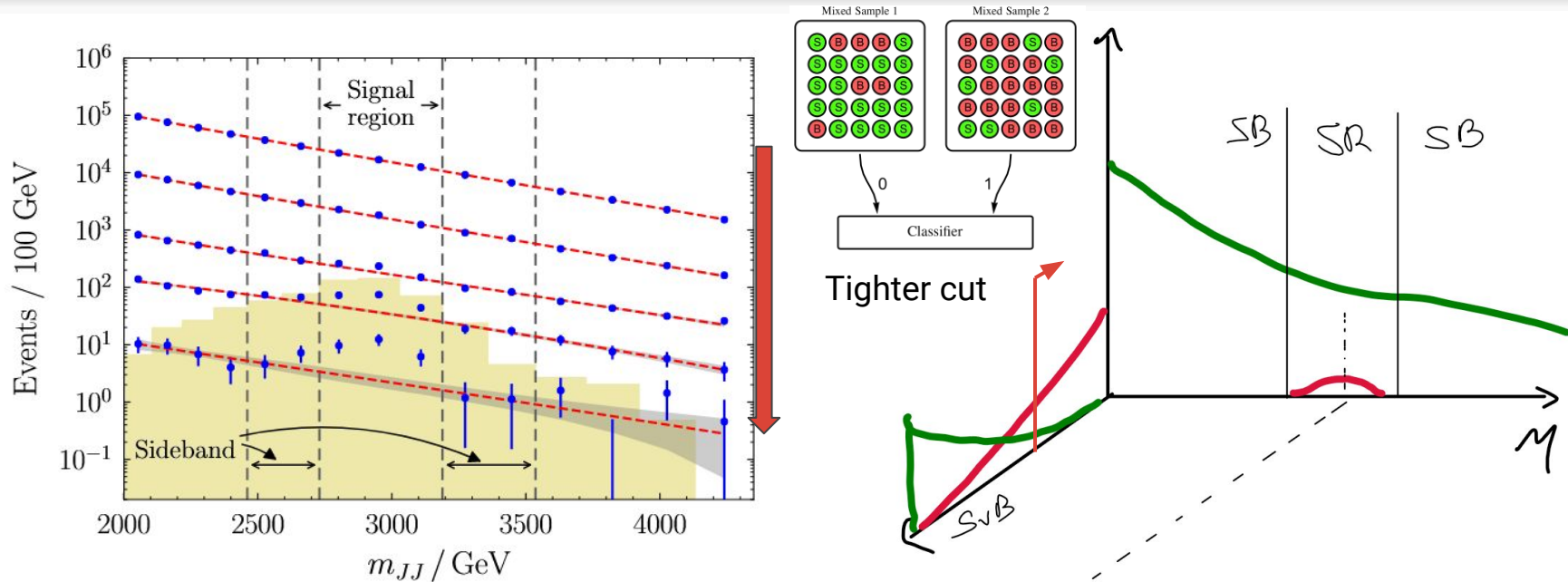
What do we train on?



Control region data  
vs  
Signal region

Noisy labels - weakly supervised!

# Extending Bump hunts with CWoLa



Collins et al Metodiev et al

# Extending Bump hunts

Works really well unless observables in classifier are ***correlated with Mass!***

How to take into account?

- Optimise choice of observables
- Bring in additional ML approaches for producing the background

# Extending Bump hunts

Approaches showing great improvements over standard CWoLa bump hunt

## ANODE (Nachman & Shih)

- Direct density estimation with normalizing flows, using base density for anomaly detection

## CATHODE (Hallin et al)

- Normalizing flows trained outside of signal window, generate background data in signal window

## SALAD (Andreassen et al) - Not using normalizing flows but also very good performance

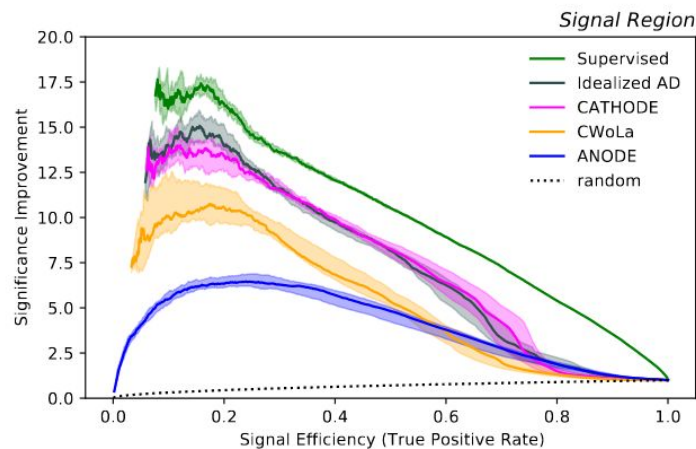
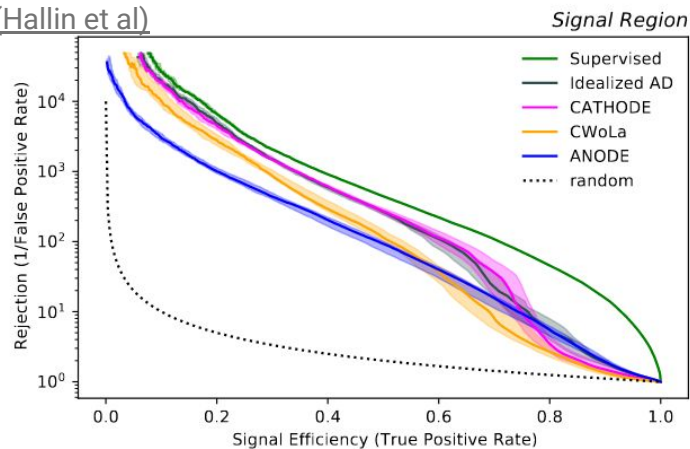
- Use simulation to transfer classifier to data with density ratio estimation

# Extending Bump hunts

Other approaches build on the idea and show great improvements over standard CWoLa bump hunt

ANODE (Hallin et al)

CATHODE



indow

# Introducing CURTAINS

*“What would this datapoint look like if it had a different value of mass?”*

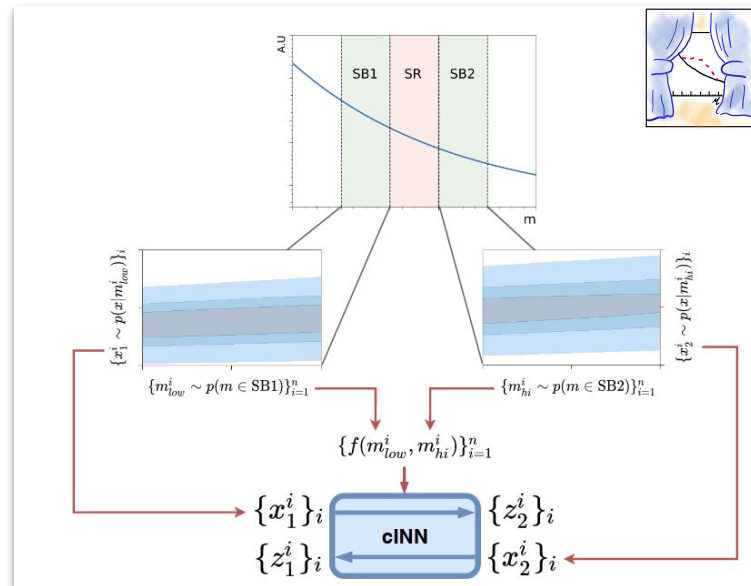
Train an INN to map data between sidebands

- Condition on the input and target mass
- Learn to account for changing mass

Once trained by choosing target mass values

**Transport sideband data into signal region!**

- No need to sample from base distribution
- Only estimate mass distribution in signal region





# Training CURTAINS

Draw data  $x$  from SB1 and SB2

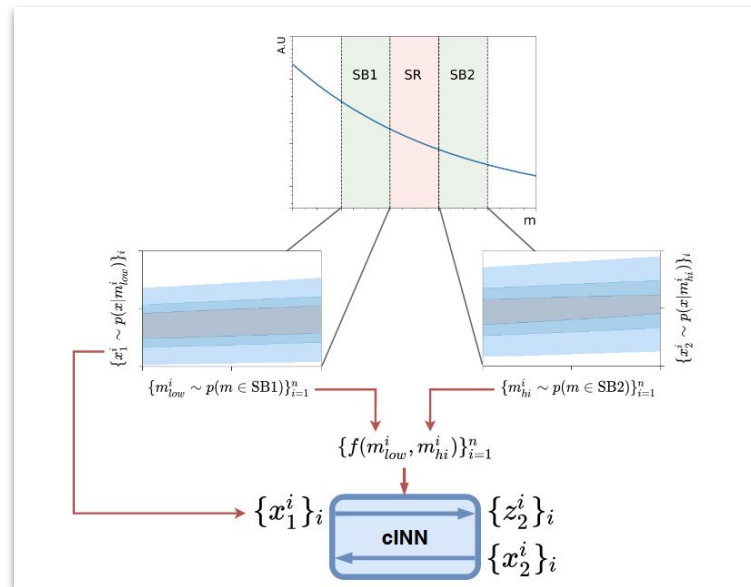
Assign target masses based on batch

Transport data from SB1 $\leftrightarrow$ SB2

- Can **train bidirectionally!**

Use an Optimal Transport loss to measure difference between  $z_2$  and  $x_2$

**Compares distributions not datapoints**



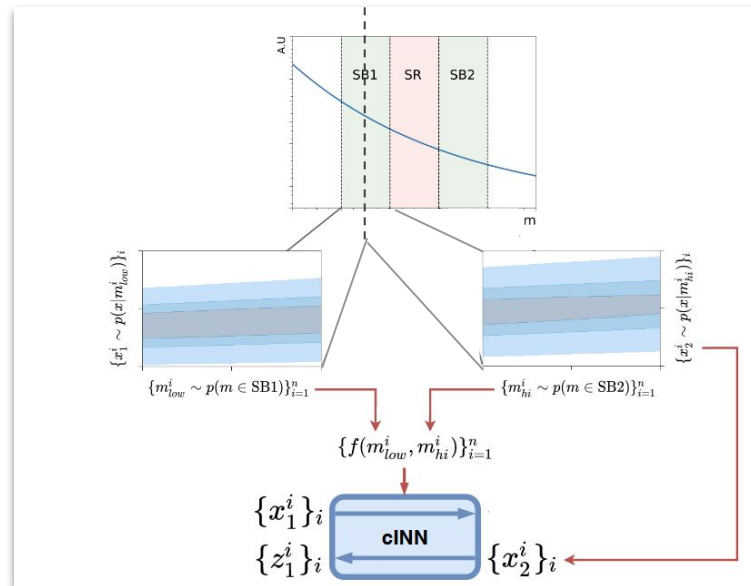
# Training CURTAINS - technicality

Using  $f(m, m') = m' - m$

- During training **min value is width of SR**
- To transport data from SB->SR **min value is 0**
- Outside of training domain...

Solution: Split sidebands into two

- Train between lower and upper half of SB1/2
- Now min value is also 0

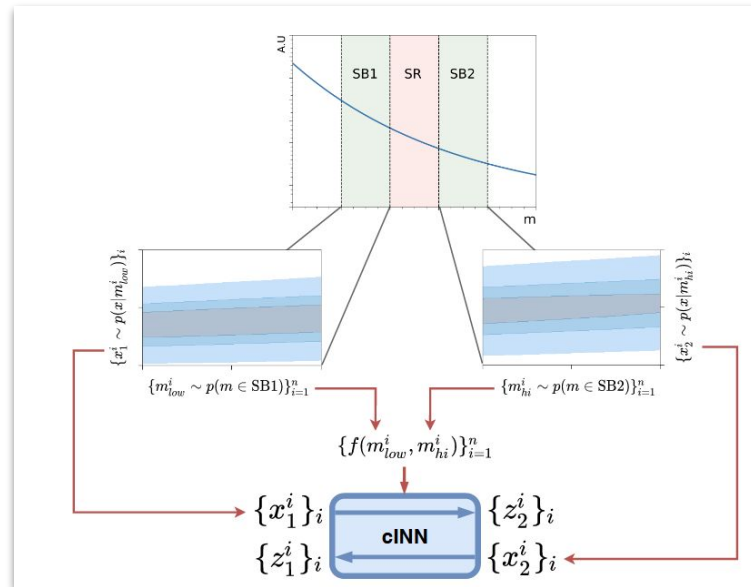


# Training CURTAINS - technicality

Using Sinkhorn loss to compare  $x$  to  $z$

- Difference in distribution over events
- Minimised if distributions and correlations correct after transport
- **Does not** strictly enforce correct mass conditioning for each data point
- Slow convergence due to stochastic sampling of target batch

**Not ideal** but **empirically works**



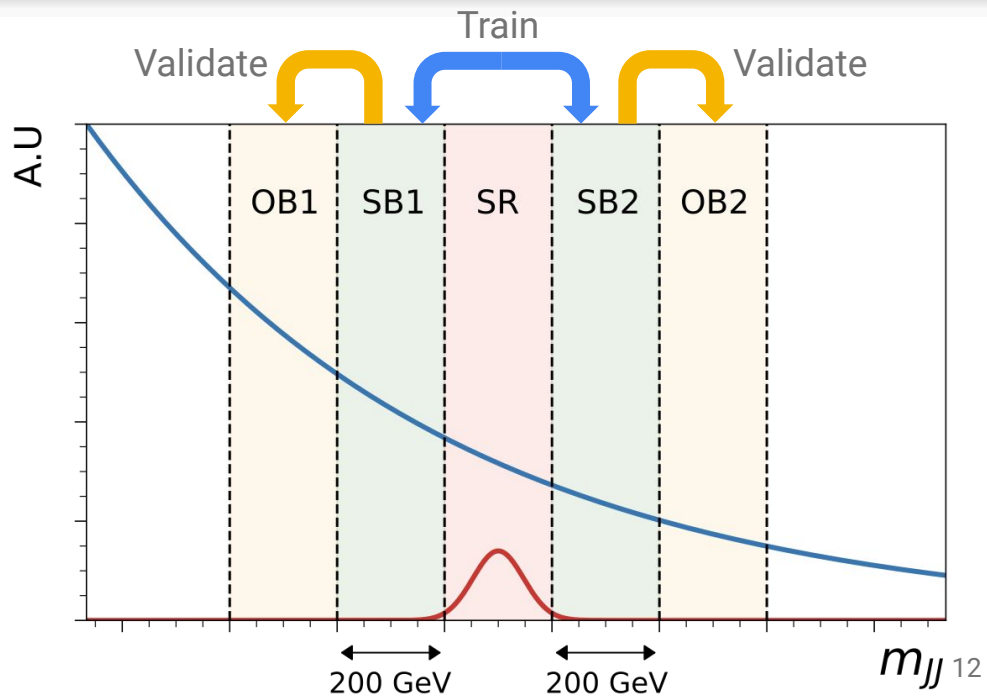
# CURTAINS Validation

Fix sidebands

Define Outer-Band (OB) validation regions

Train CURTAINS transformer

Validate on OBs



# CURTAINS - Training regions

Training on the LHC0 R&D anomaly detection dataset

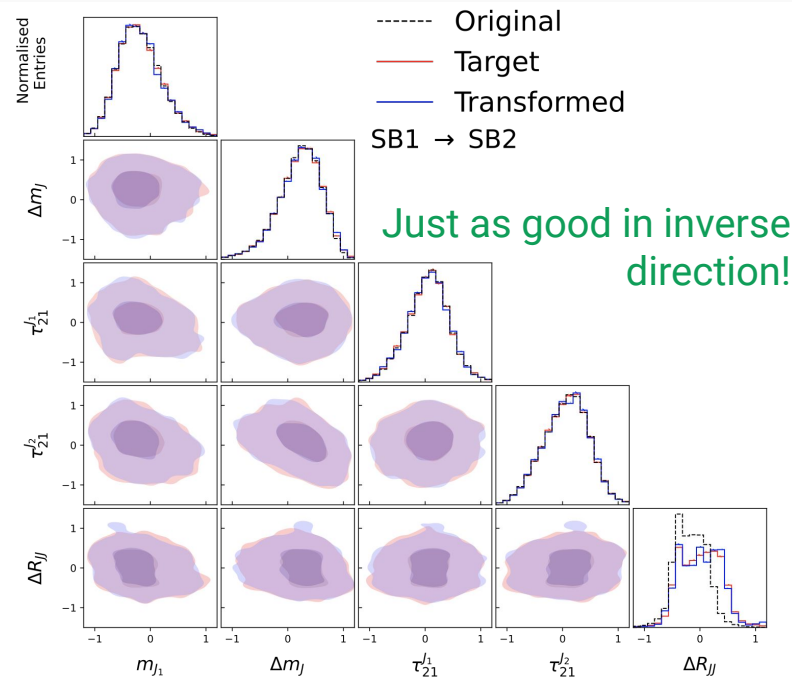
Sideband 1: [3200, 3400]

Sideband 2: [3600, 3800]

Five observables

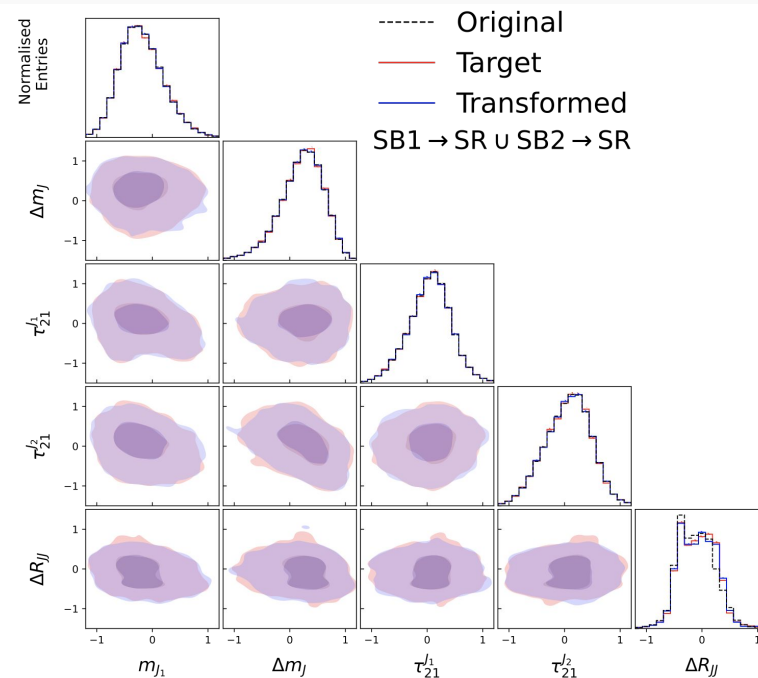
$M_{J1}, M_{J1} - M_{J2}, \tau_{J1}^{21}, \tau_{J2}^{21}$

Plus  $\Delta R_{JJ}$  due to correlation to  $M_{jj}$



# CURTAINS - Signal region

Nearly perfect matching implies near perfect background template!



*\*Can't look at this in the real analysis or application!*

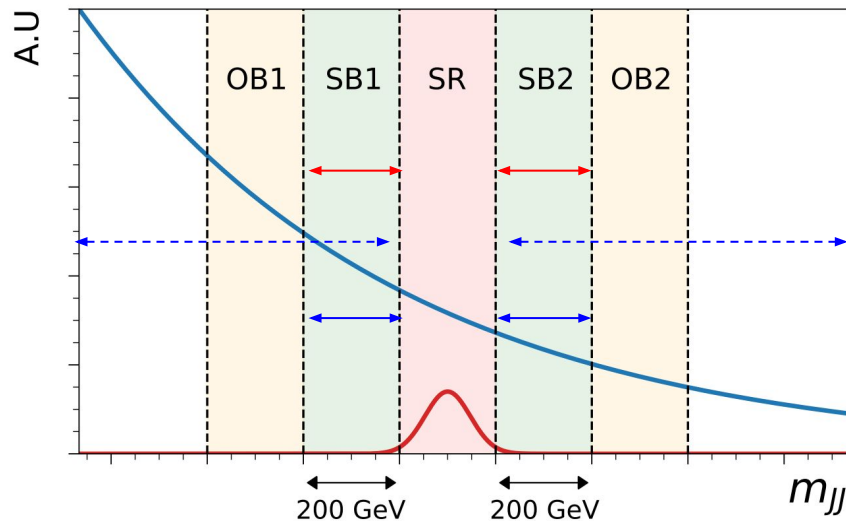
# CURTAINS - CWoLa Performance

Compare to CATHODE method

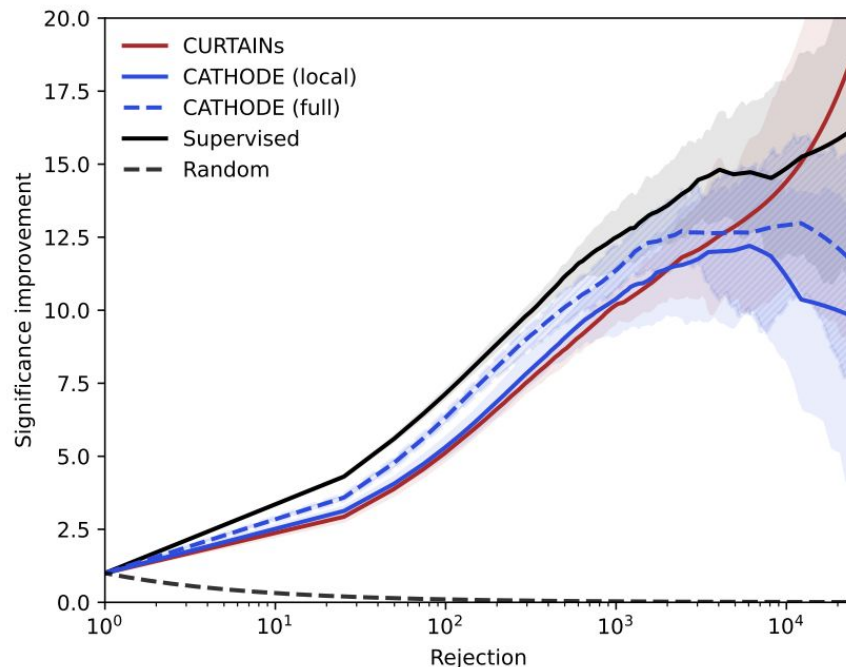
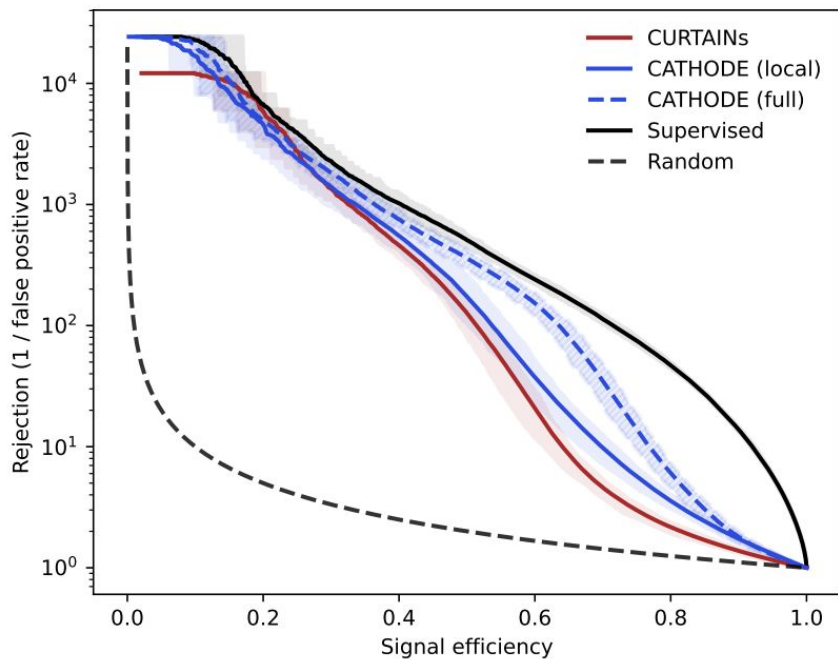
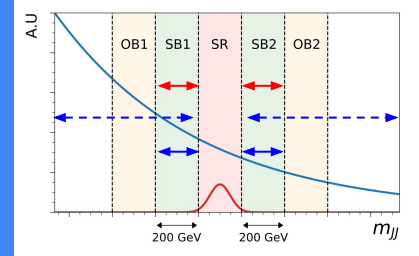
- Equivalent training window (local)
- All available data outside of SR (full)

Same number of generated bkg for both methods

- “Oversample” CURTAINS by transporting same data to multiple values of  $m$

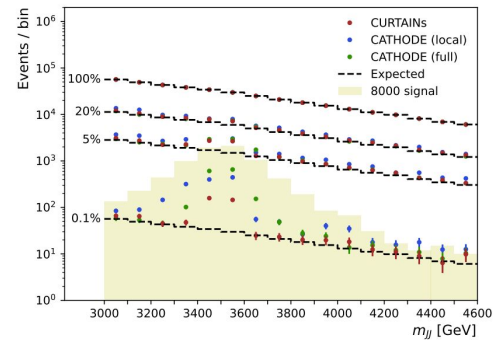
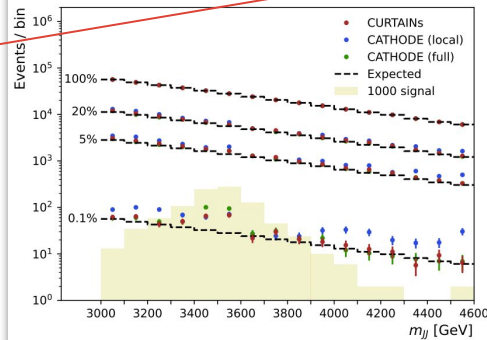
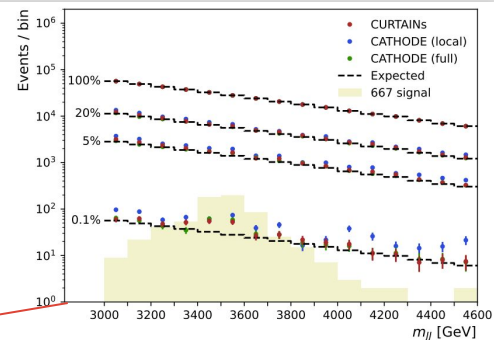
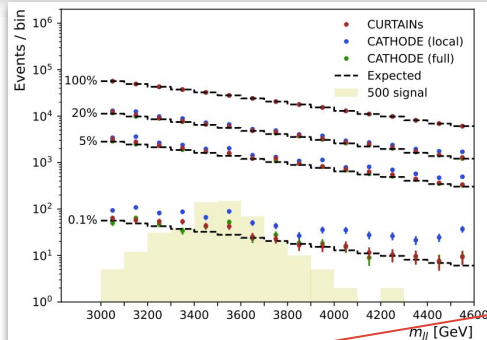
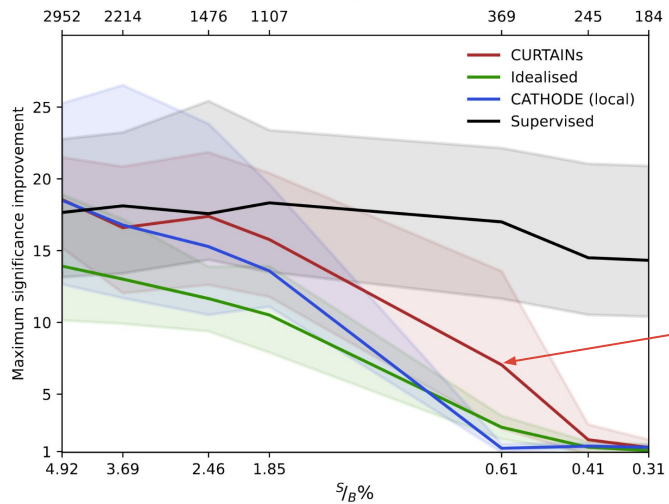


# CURTAINS - CWoLa Performance





# CURTAINS - Bump hunt



# CURTAINS

# Flows4Flows

# Flows4Flows

Use normalizing flows to parametrise base distribution! No more approximate OT loss!

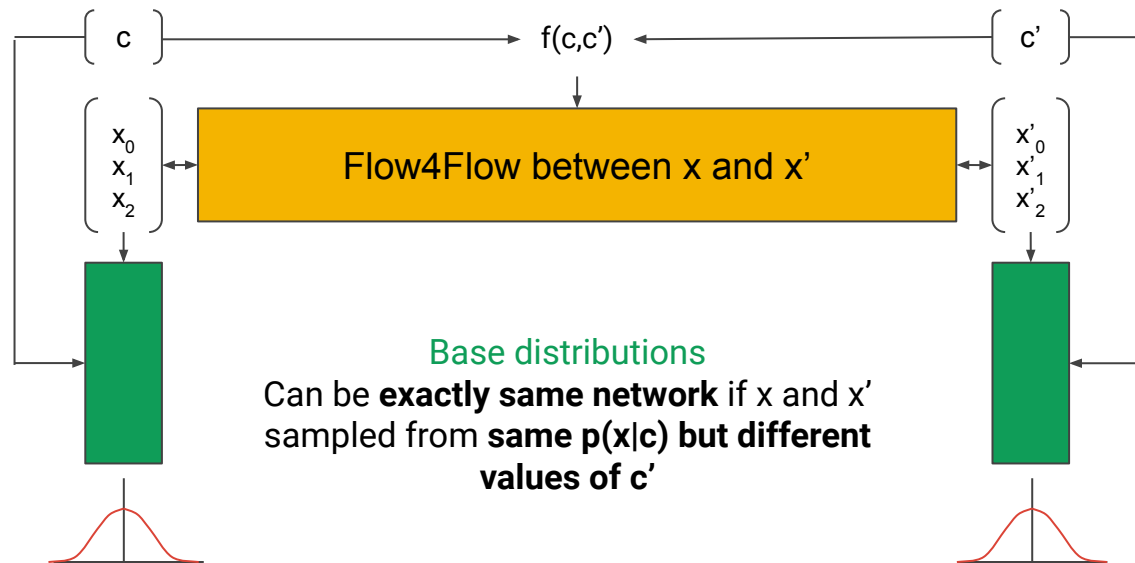
## Flows4flows

Train a flow between arbitrary distributions

Simply another change of variables for  $p(z)$  in normalizing flows!

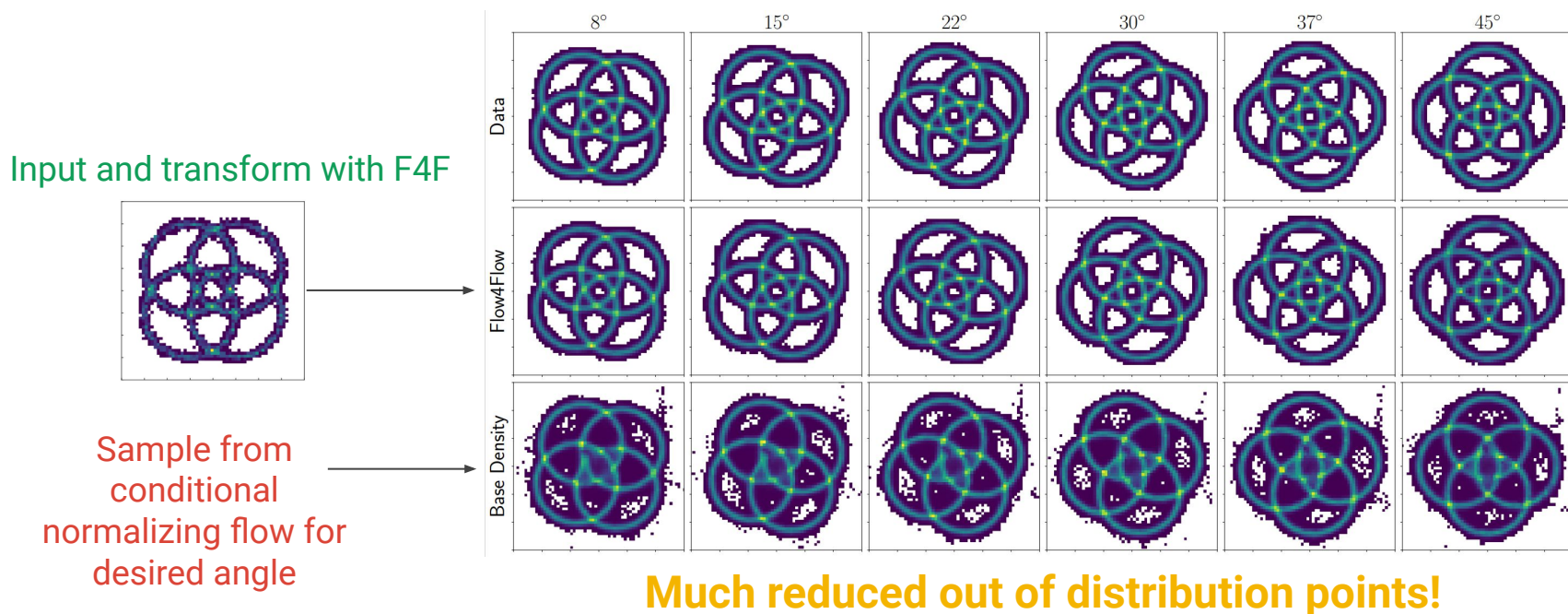
Pretrain base distribution(s)

Use base distribution for loss in exact maximum likelihood

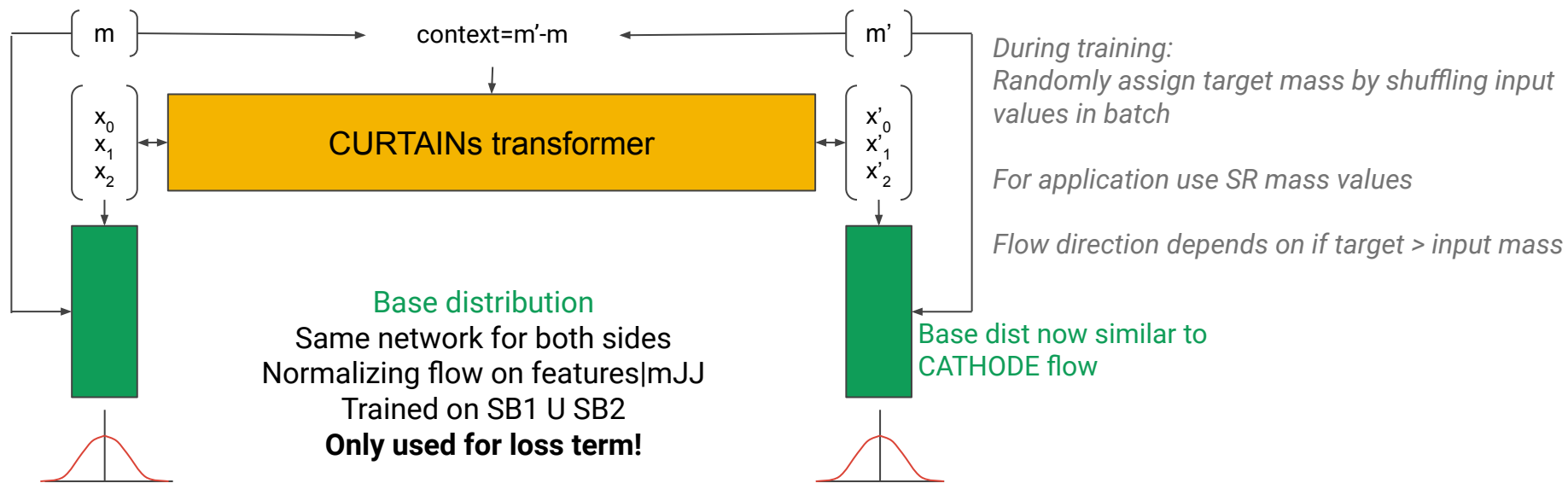


**Can now train CURTAINS with exact maximum likelihood!**

# Why Flows4Flows

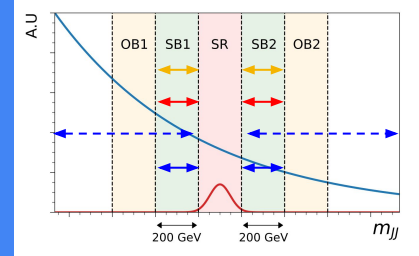


# CURTAINS - Flows4Flows



$$\text{Loss: } \log(p(x)) = \log \det |J(f(x|m, m'))| + \log \det |J(g(f(x|m, m')|m'))| + p(g(f(x|m, m')|m'))$$

# CURTAINs - Flows4Flows



**Significant improvement** with new loss!

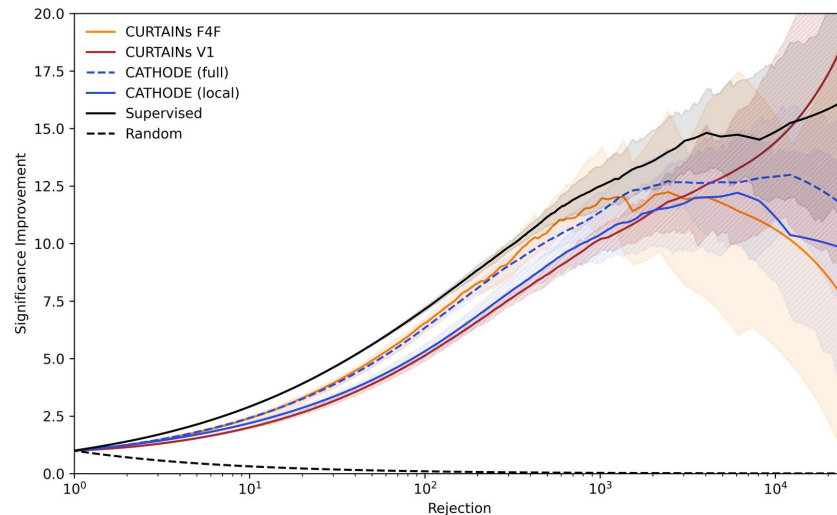
Much faster to train, *including* base density

Still trained on a very local window

- Only 200GeV either side of SR
- Matches CATHODE (full) now over most of the range

Compared to CURTAINs v1

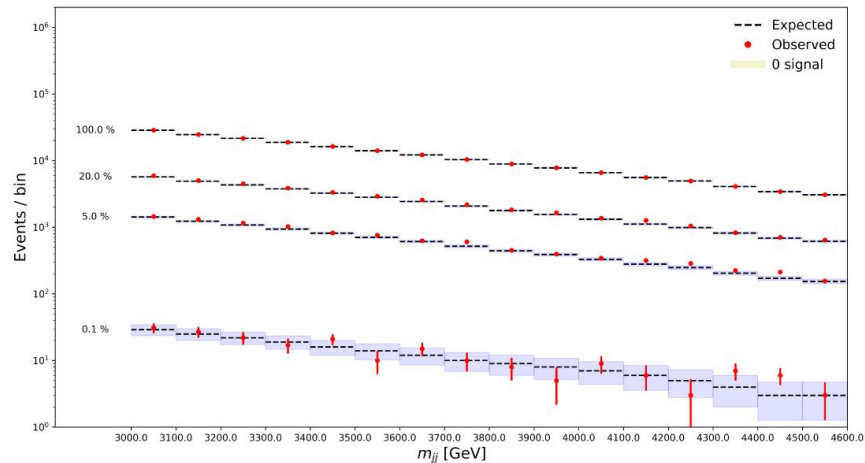
- **Simpler** to set up and train
- Features can be even more strongly correlated to resonant feature



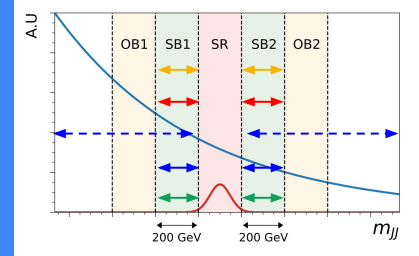
# CURTAINs - Flows4Flows

Compared to CURTAINs v1

- **Simpler** to set up and train
- Features can be even more strongly correlated to resonant feature
- Still **robust** to case where there is no signal

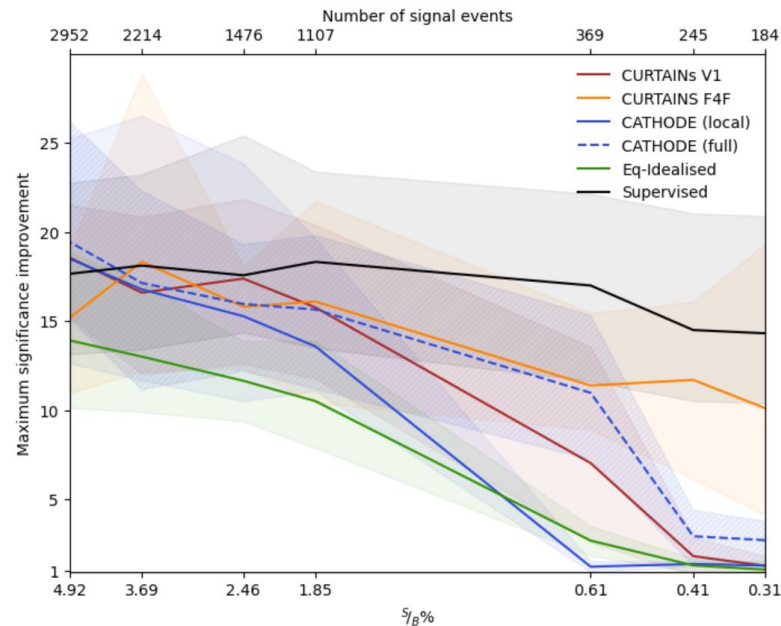


# CURTAINS - Flows4Flows



Compared to CURTAINS v1

- **Simpler** to set up and train
- Features can be even more strongly correlated to resonant feature
- Still **robust** to case where there is no signal
- **Much more sensitive** to even small amounts of signal





# Summary

**CURTAINs** is a new method for enhancing the Bump Hunt with CWoLa style classifiers

- Transforms data from sidebands into signal region
- Bypass need of going via an intermediate distribution
- Produces background data in SR and is complementary to other anomaly techniques

**CURTAINs** matches the performance of leading approaches without needing to train on the full mJJ spectrum

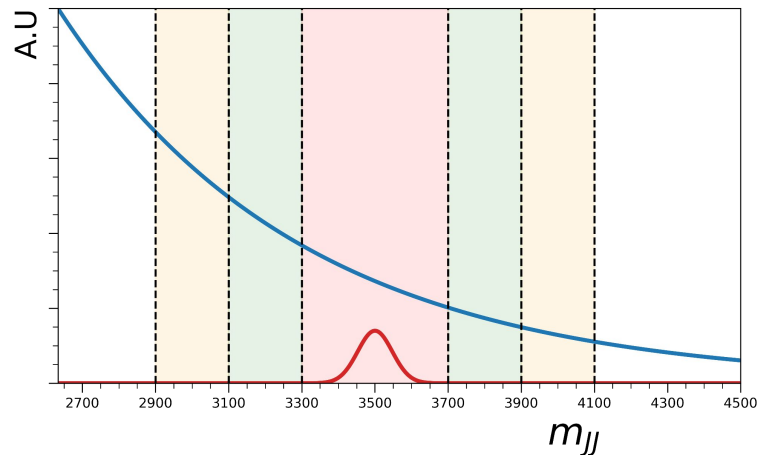
- Leading performance in a local setup

**CURTAINs+Flows4Flows** can reach even higher levels of performance - preprint soon!

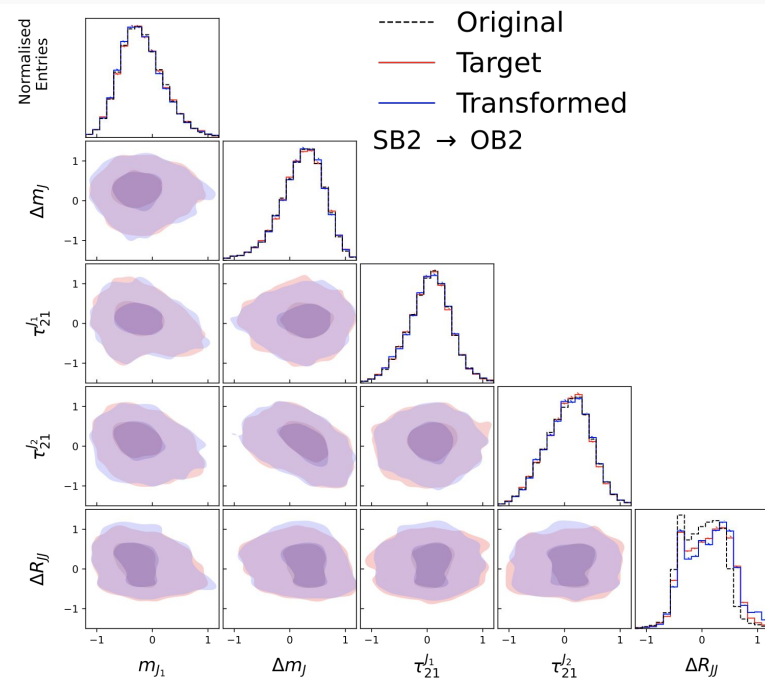
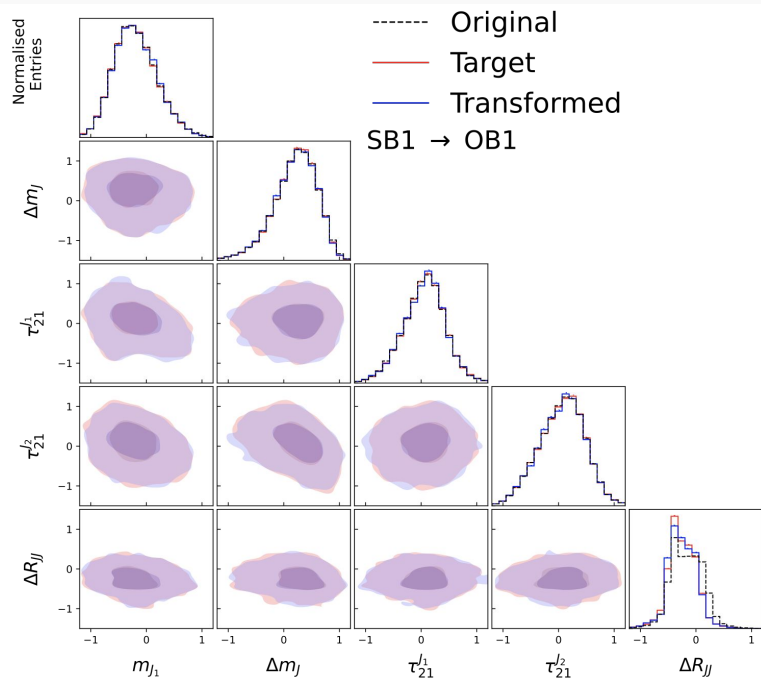
# Backup

# CURTAINs - CWoLa Performance

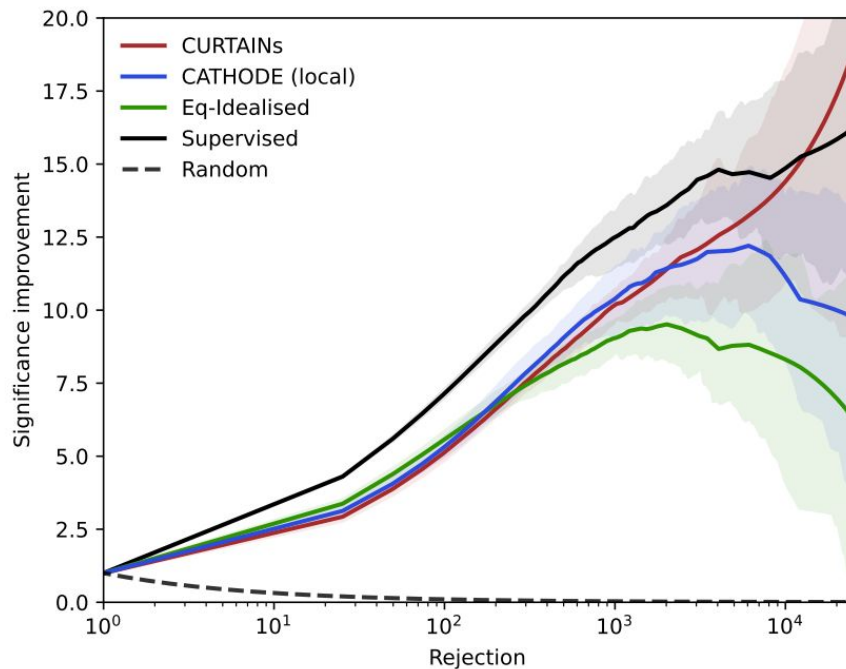
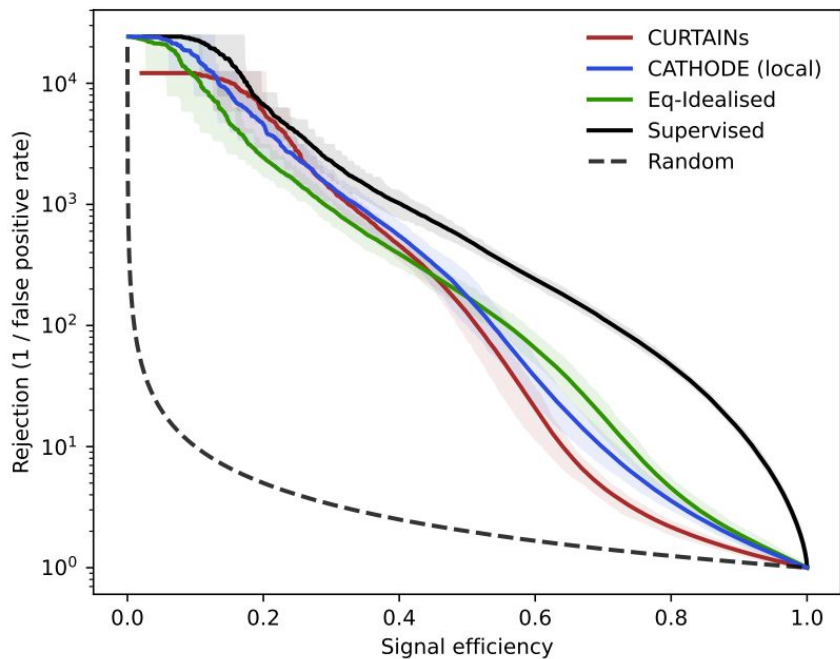
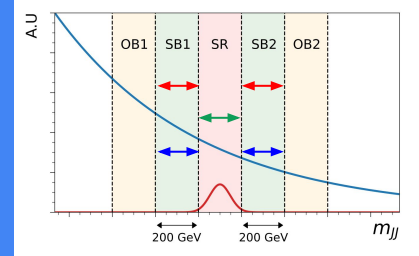
- Fix the signal region such that it contains almost all of the signal
- Train a CURTAINs transformer
- Train a classifier SR data vs CURTAINs



# CURTAINS - Validation regions



# CURTAINS - CWoLa Performance

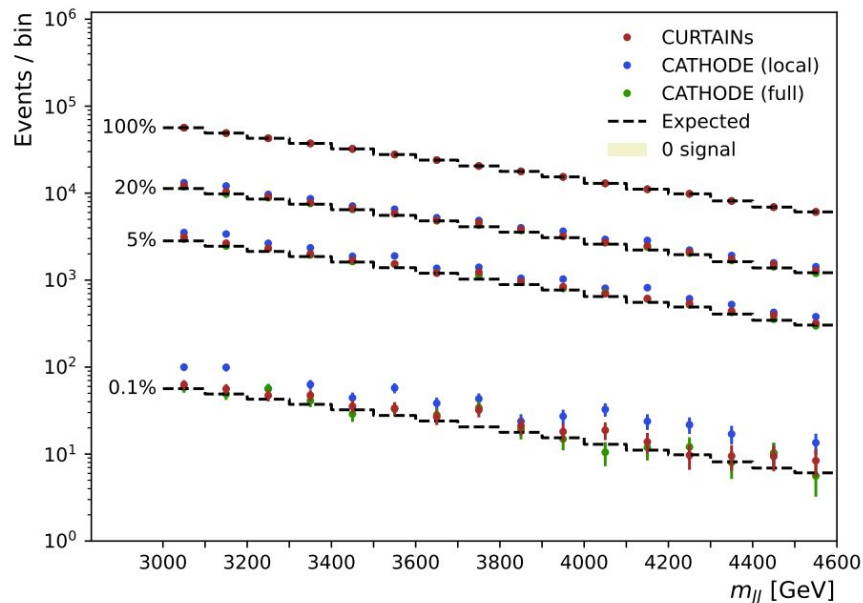


Idealised = take true Background data, and train S+B vs B with equal statistics

# CURTAINS - Bump hunt

Repeat CWoLa setup with  
non-overlapping 200 GeV steps

Apply cuts on classifier trained with e.g.  
CURTAINS and look for a bump



# CURTAINS - Flows4Flows

Train normalizing flow on  $p(x, m) = \text{SB1} \cup \text{SB2}$  for **base distribution** - like CATHODE

- But unlike CATHODE, use this for training another flow, not generating samples

Construct a **flow4flow from  $x$  to  $x'$**  conditioned on current and target masses ( $m, m'$ )

- Transform  $x \sim p(x|m)$  to  $p(x|m' \sim p(m))$  in flow  $f(x|m, m')$  - like CURTAINS

But now loss given by maximum likelihood:

$$\log(p(x)) = \log \det |J(f(x|m, m'))| + \log \det |J(g(f(x|m, m'))|m')| + p(g(f(x|m, m'))|m')$$

Flow4flow transform

Base density transform

Base density probability

# CURTAINs - Flows4Flows

If a datapoint has  $m' \geq m$ , transform from “left to right”

If a datapoint has  $m' < m$ , transform from “right to left”

No longer need to split sideband, train on combination of all data

- Guaranteed widest support of conditioning variable!
- Leads to much faster training, no longer iterating
  - Forward and inverse pass based on input/target, both passes done per batch



# CURTAINS - Flows4Flows

