



Research Training Group  
Physics of the Heaviest  
Particles at the LHC



Collaborative Research Center TRR 257



Particle Physics Phenomenology after the Higgs Discovery



Institute for  
Theoretical  
Particle Physics  
and Cosmology

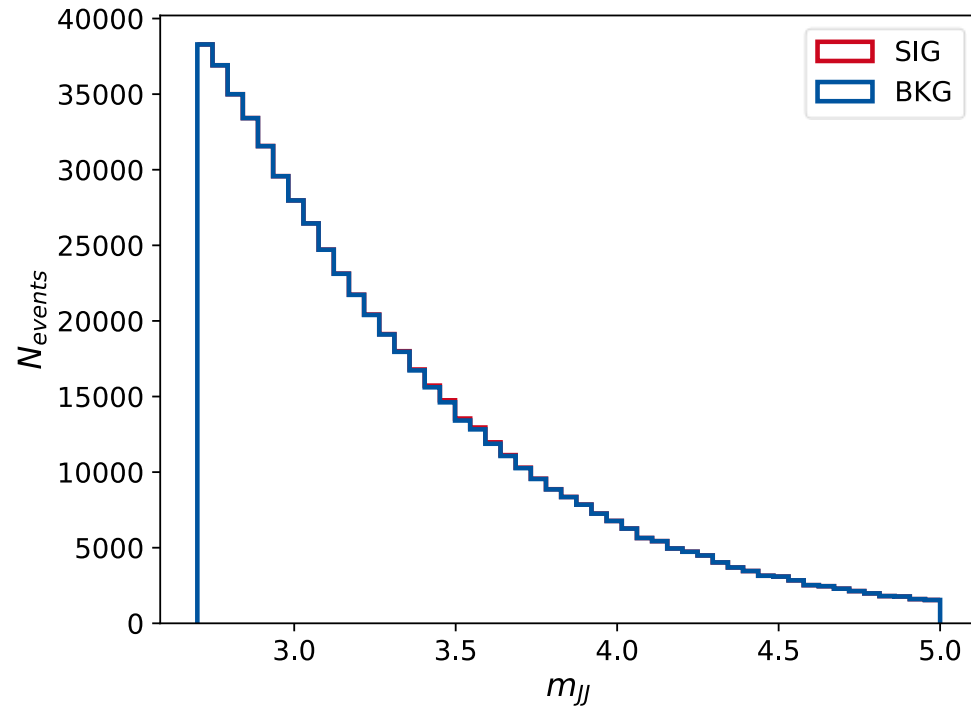


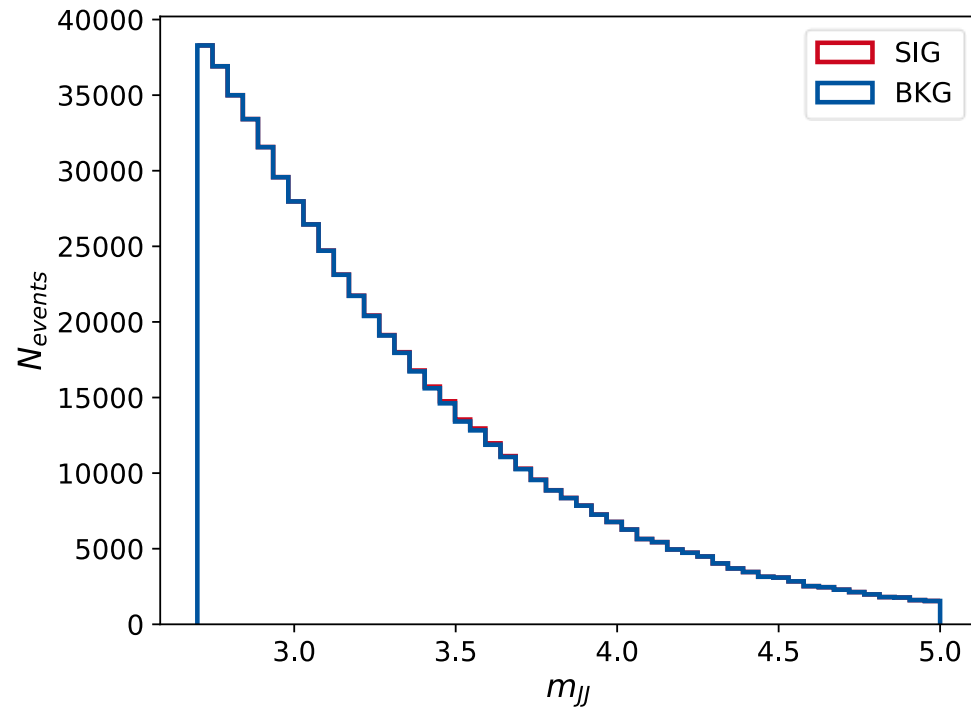
# Picking the right setup for anomaly detection

**Marie Hein**

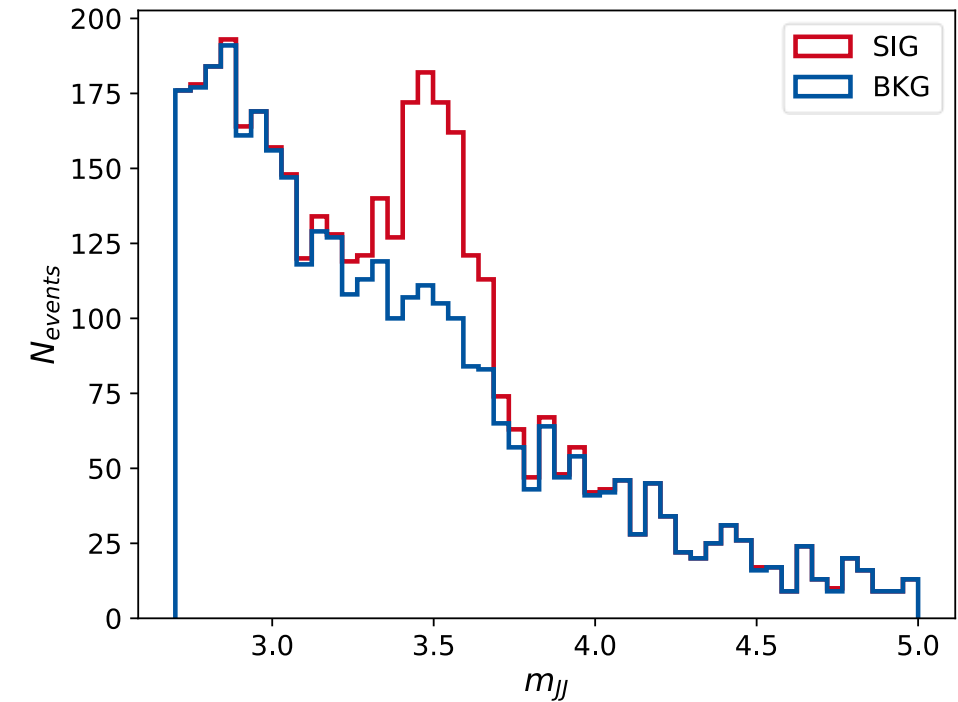
with Gregor Kasieczka, Michael Krämer, Louis Moureaux, Alexander Mück, Tobias Quadfasel and David Shih

CRC Young Scientists Meeting 2025



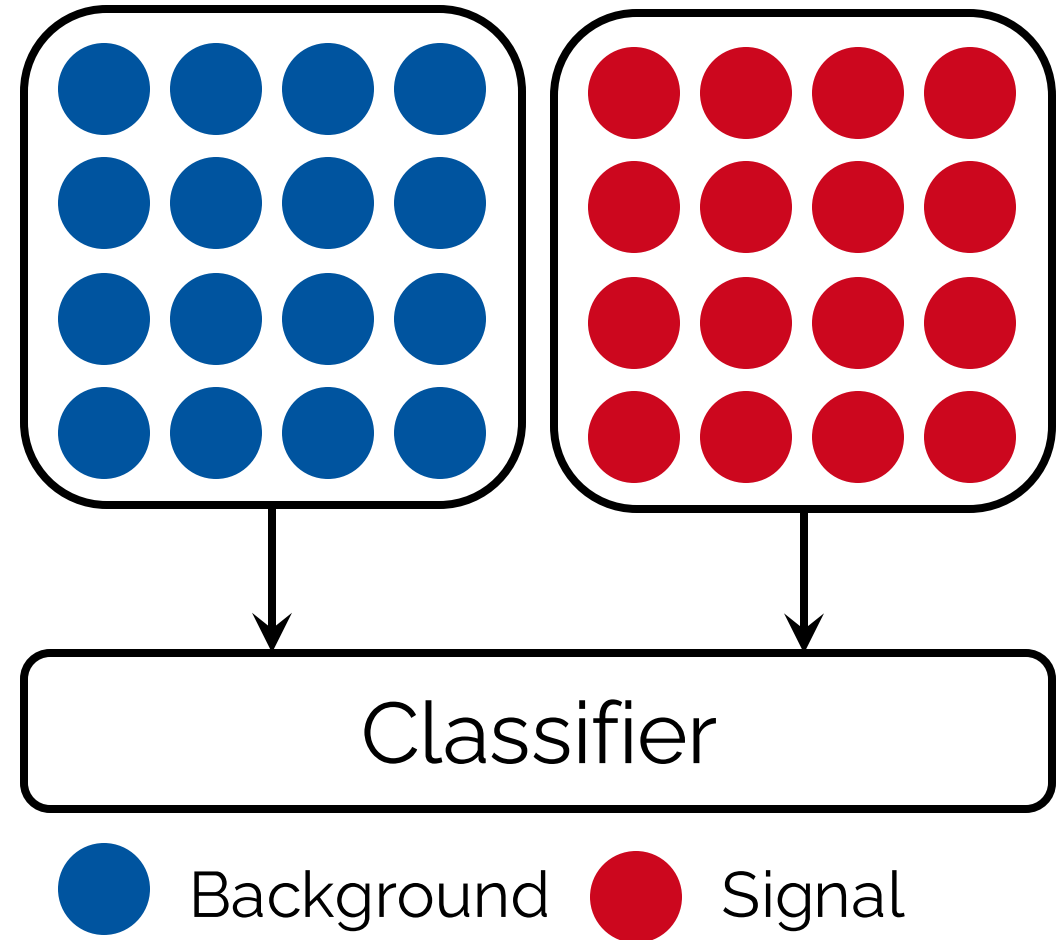


Anomaly  
Detector



- Optimal classifier

$$R_{\text{optimal}}(x) = \frac{p_S(x)}{p_B(x)}$$

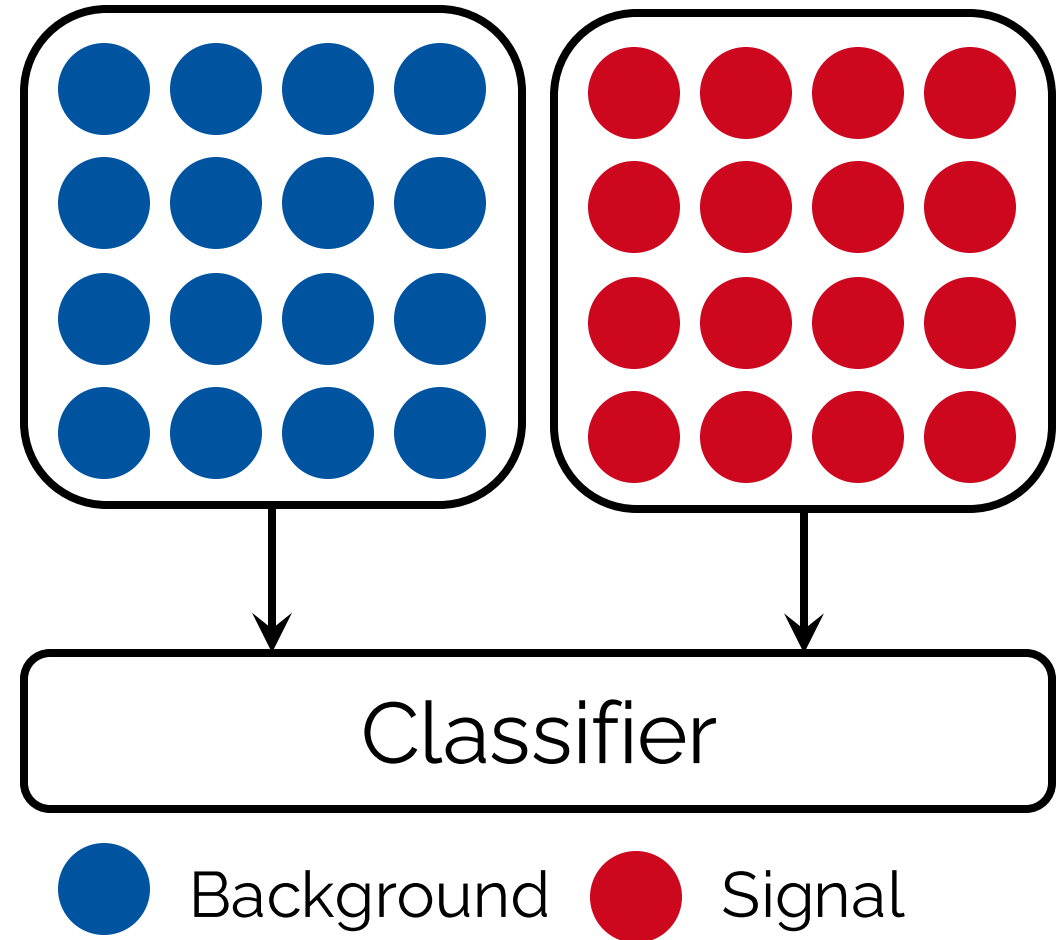


- Optimal classifier

$$R_{\text{optimal}}(x) = \frac{p_S(x)}{p_B(x)}$$

- For Machine Learning use binary cross entropy loss

$$BCE = -\log p_{\text{pred, true}}$$



- Optimal classifier

$$R_{\text{optimal}}(x) = \frac{p_S(x)}{p_B(x)}$$

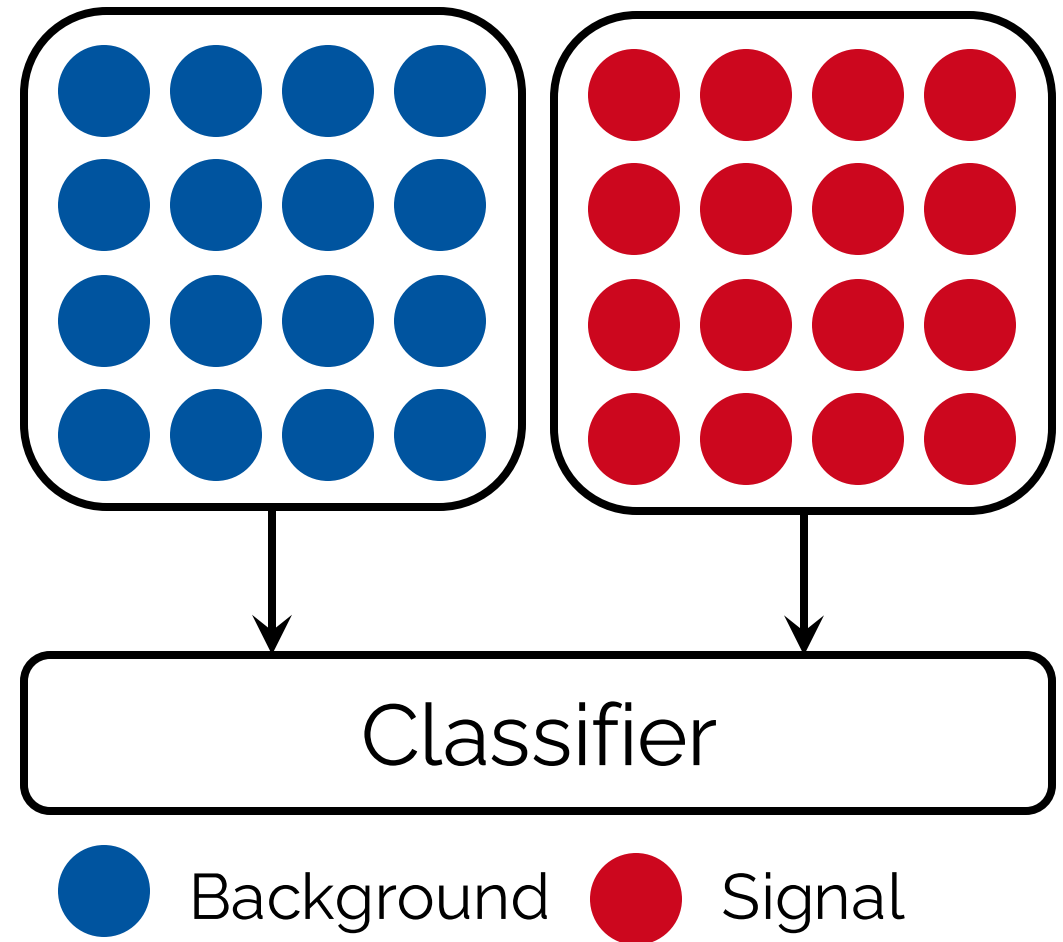
- For Machine Learning use binary cross entropy loss

$$BCE = -\log p_{\text{pred, true}}$$

→ Optimal solution function monotonically related to  $R_{\text{optimal}}$

$$f(x) = \frac{p_S(x)}{p_B(x) + p_S(x)}$$

→ Same decision boundaries



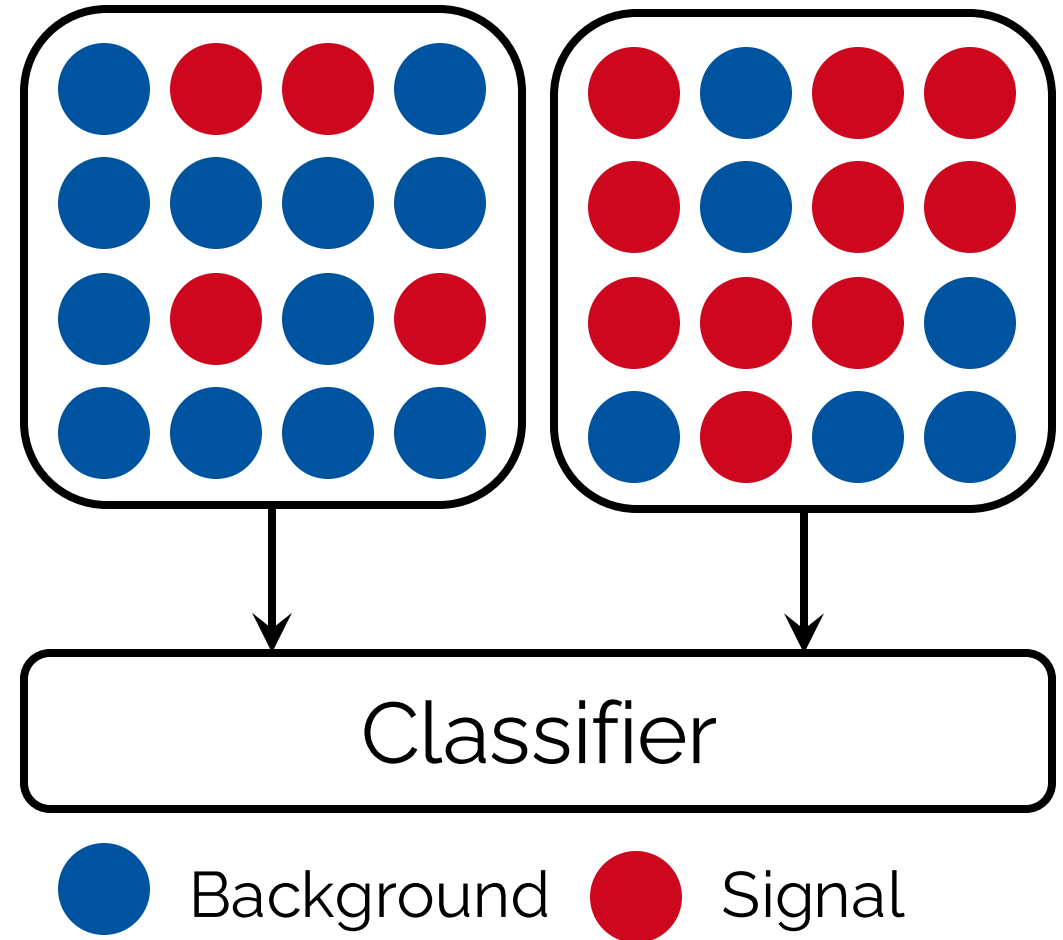
“Classification without labels: Learning from mixed samples in high energy physics” [\[1709.02949\]](#), E. Metodiev, B. Nachman, J. Thaler

- Optimal classifier

$$R_{\text{optimal}}(x) = \frac{p_S(x)}{p_B(x)}$$

- For mixed datasets with signal fractions  $f_i$

$$R_{\text{mixed}}(x) = \frac{f_1 R_{\text{optimal}}(x) + (1 - f_1)}{f_2 R_{\text{optimal}}(x) + (1 - f_2)}$$



“Classification without labels: Learning from mixed samples in high energy physics” [\[1709.02949\]](#), E. Metodiev, B. Nachman, J. Thaler

- Optimal classifier

$$R_{\text{optimal}}(x) = \frac{p_S(x)}{p_B(x)}$$

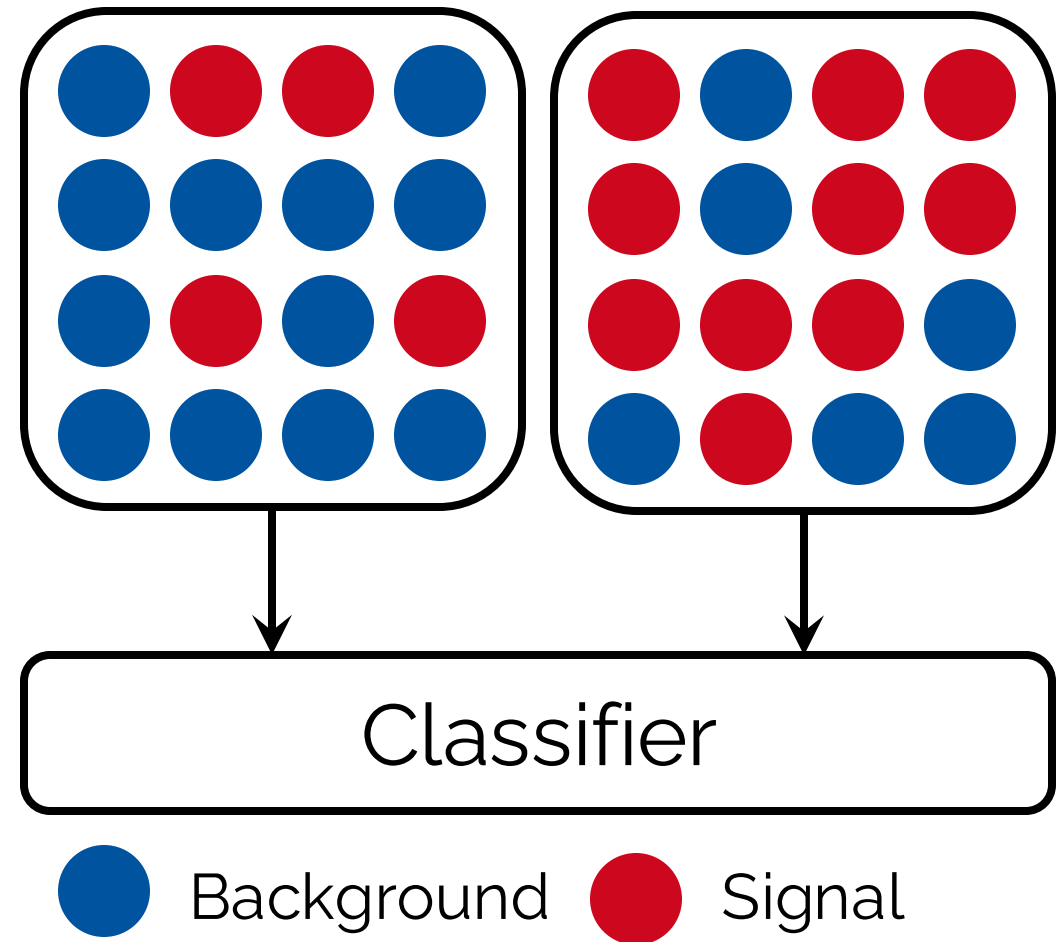
- For mixed datasets with signal fractions  $f_i$

$$R_{\text{mixed}}(x) = \frac{f_1 R_{\text{optimal}}(x) + (1 - f_1)}{f_2 R_{\text{optimal}}(x) + (1 - f_2)}$$

→ Monotonically increasing function of

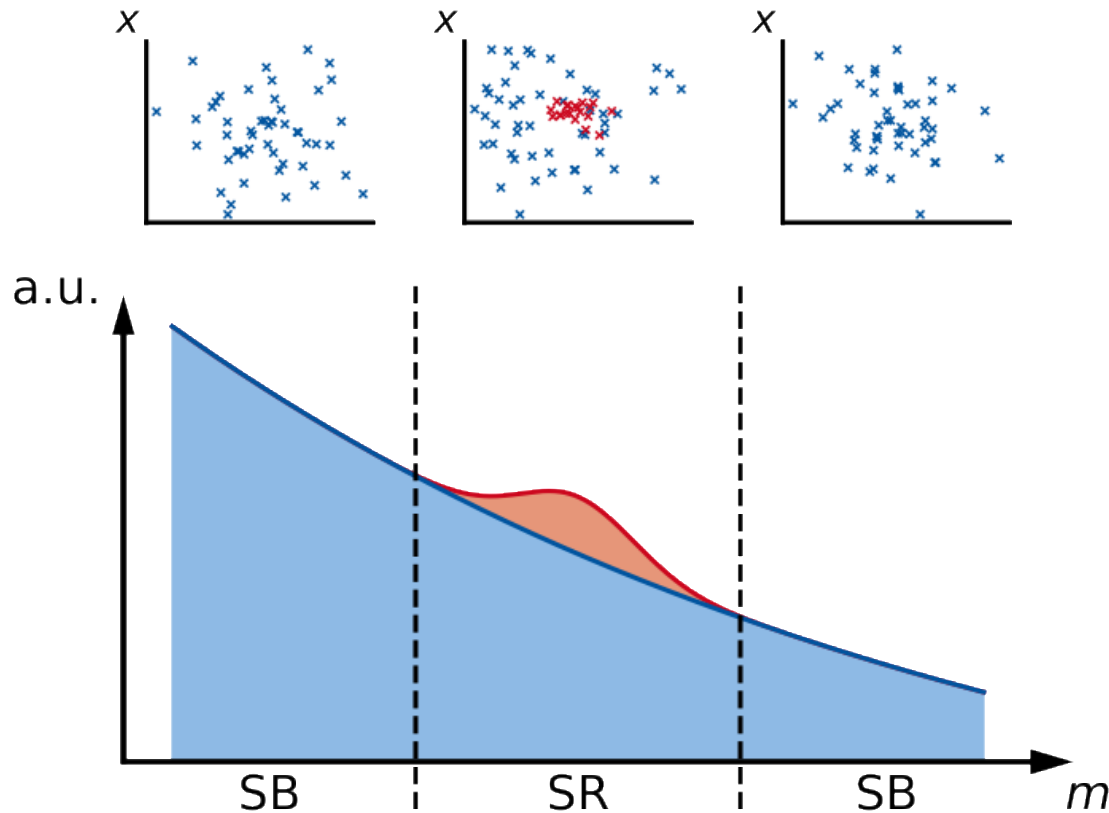
$R_{\text{optimal}}(x)$  as long as  $f_1 > f_2$

→ Same decision boundaries

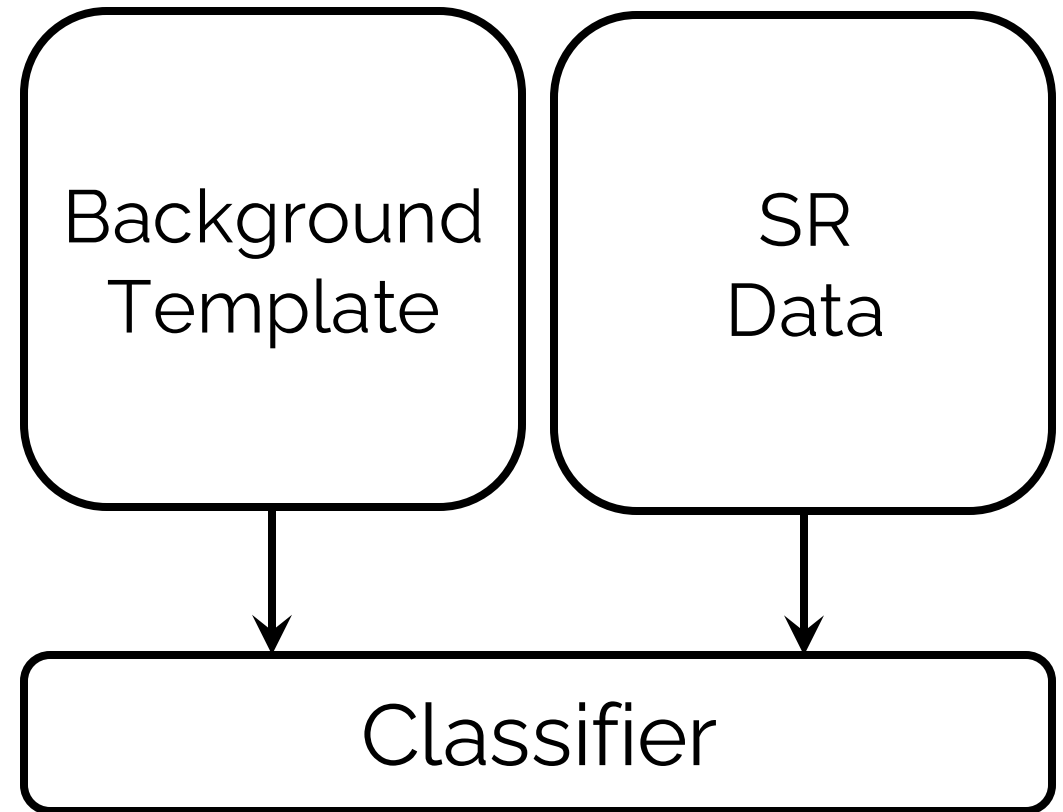




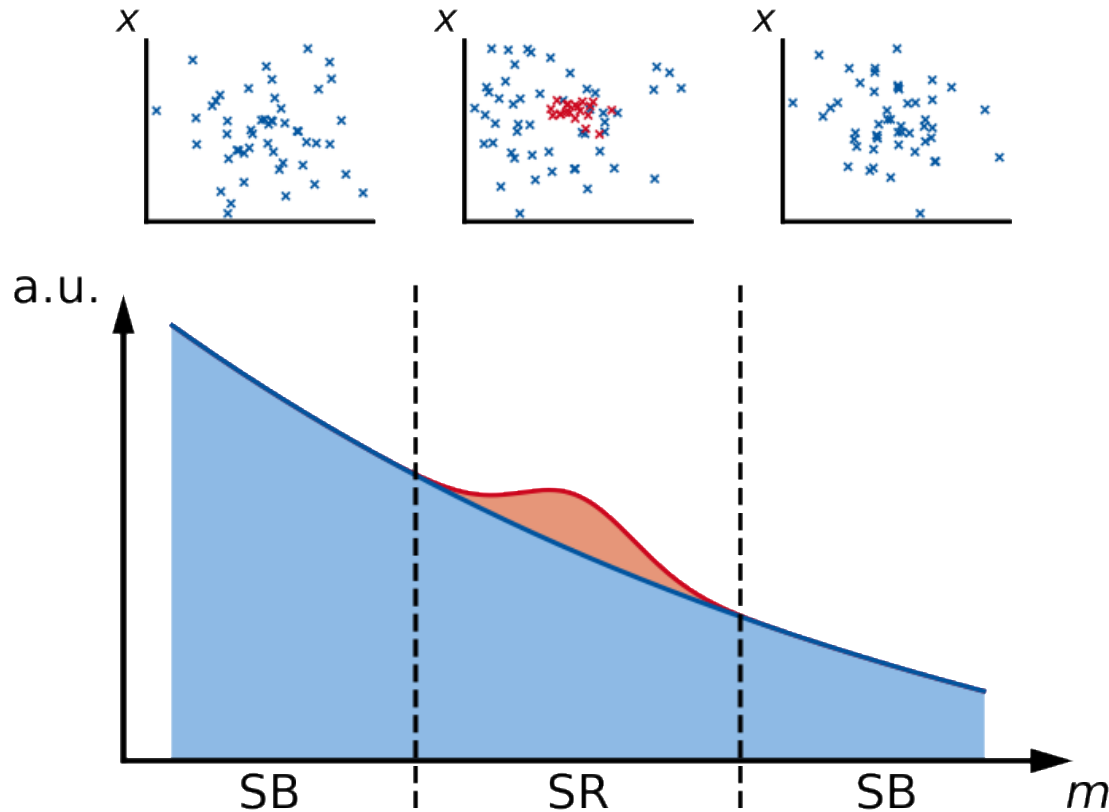
# Application to resonance searches



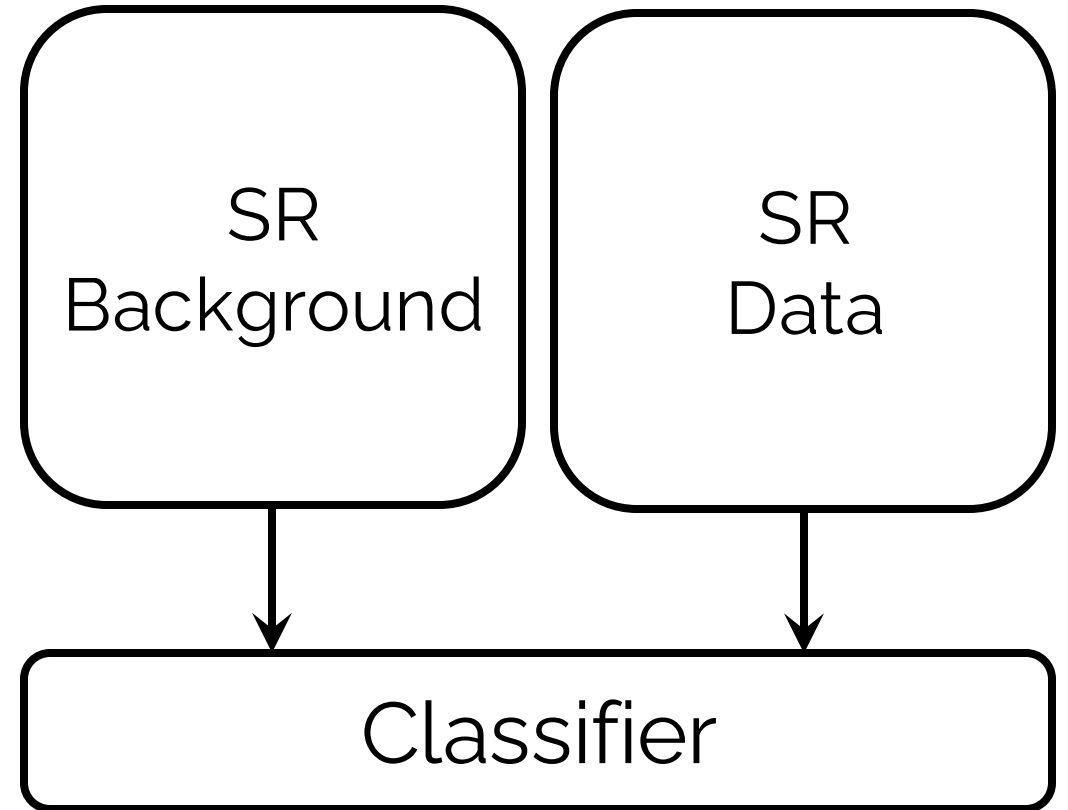
Recreated from [\[2109.00546\]](#)



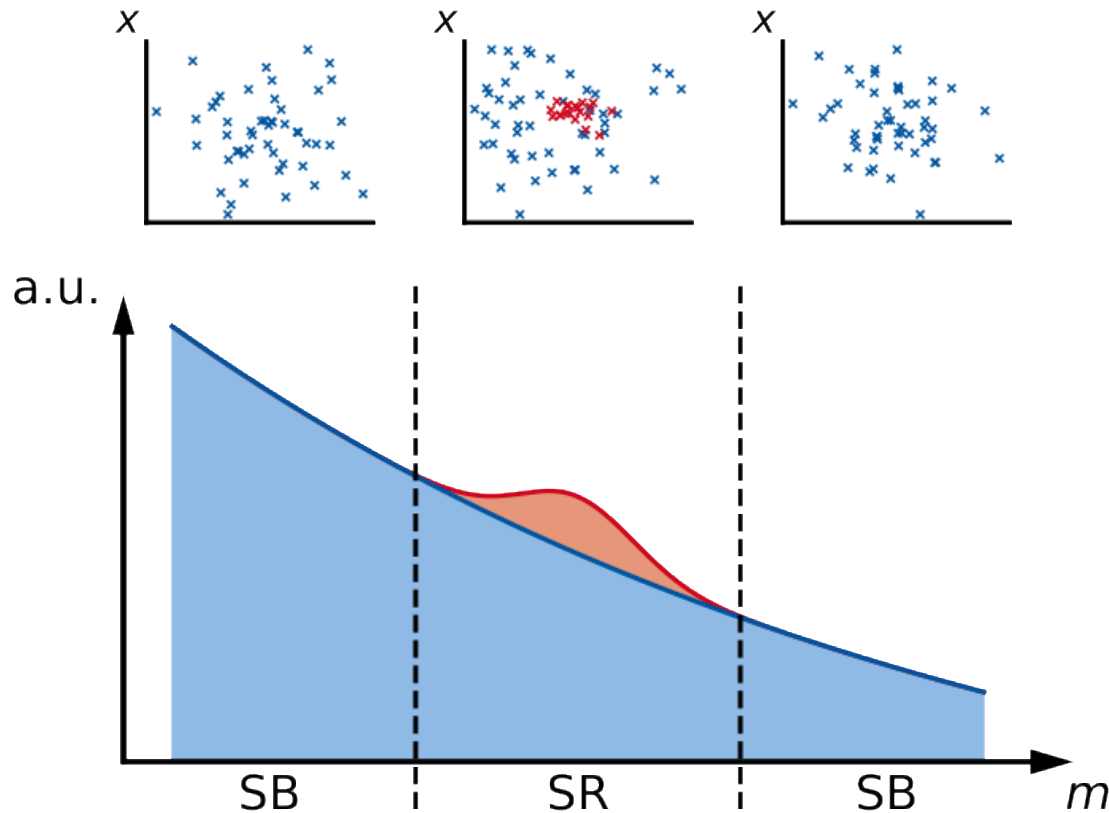
# Idealized Anomaly Detector



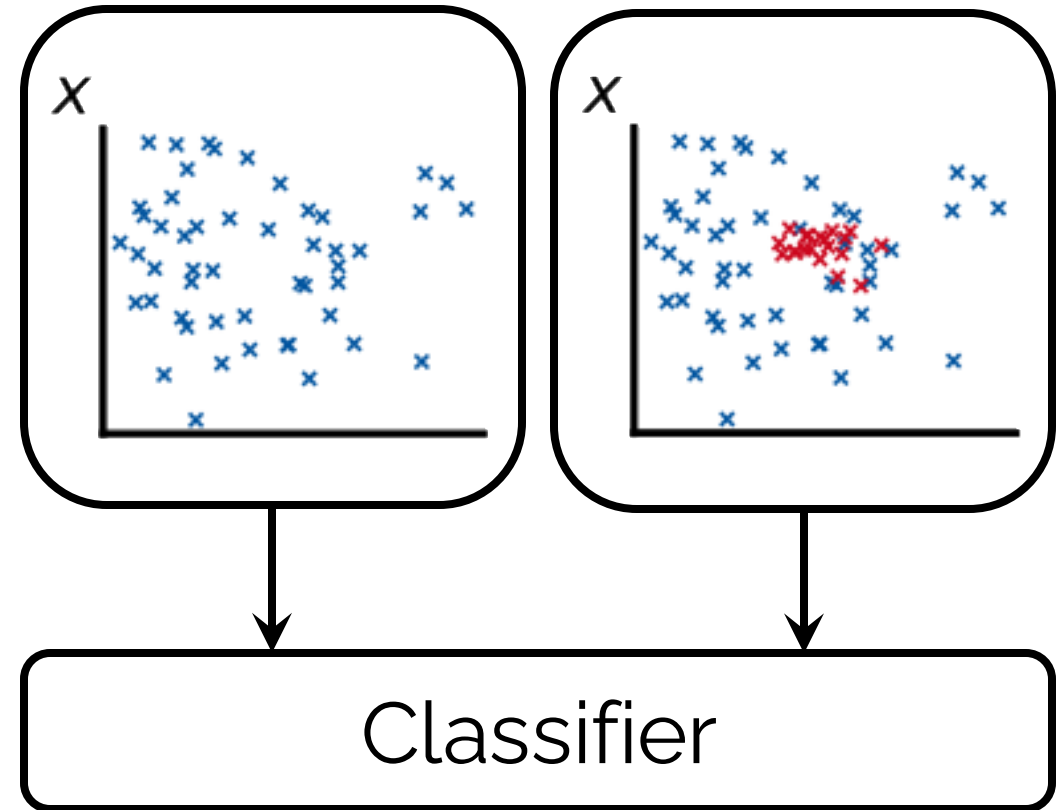
Recreated from [\[2109.00546\]](#)

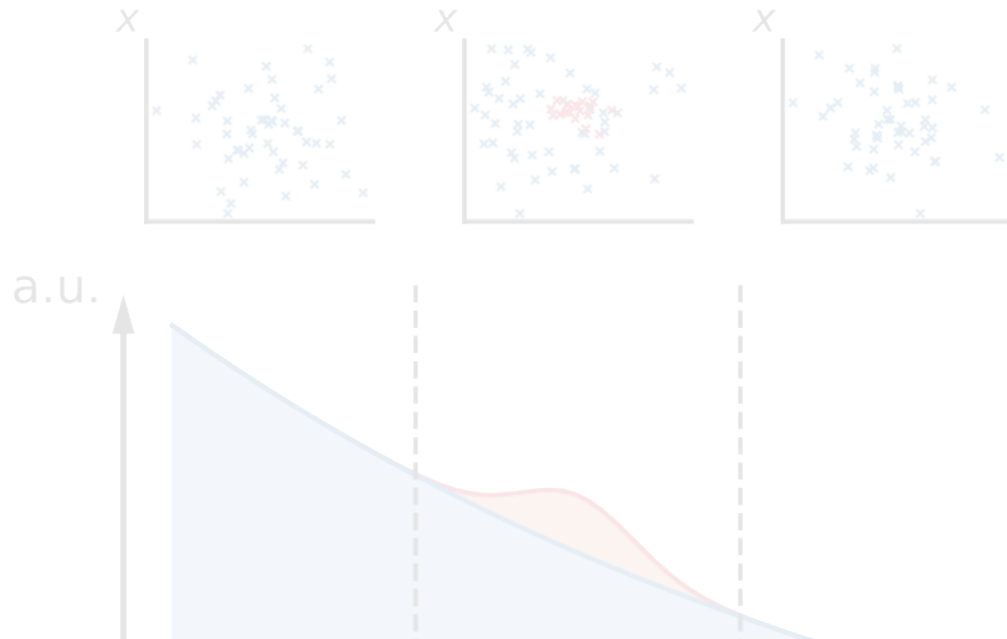


# Application to resonance searches



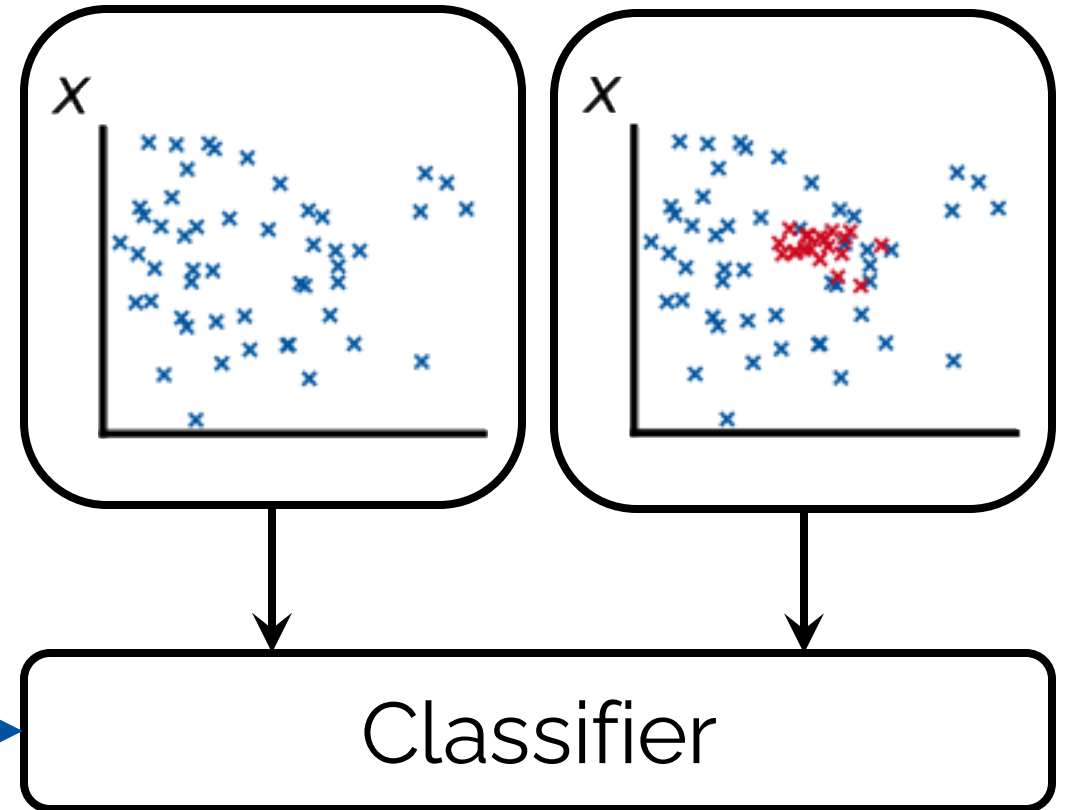
Recreated from [\[2109.00546\]](#)





How do we pick the best  
classifier architecture?

Recreated from [\[2109.00546\]](#)



# Picking the classifier architecture

**Pick on simulations**

**Pick on data**

## Pick on simulations

- **Advantages:** metrics that directly access background & signal labels

## Pick on data

## Pick on simulations

- **Advantages:** metrics that directly access background & signal labels
- **Disadvantages:** less model agnostic, dependent on specific simulation & chosen signals

## Pick on data

## Pick on simulations

- **Advantages:** metrics that directly access background & signal labels
- **Disadvantages:** less model agnostic, dependent on specific simulation & chosen signals

## Pick on data

- **Advantages:** more model agnostic

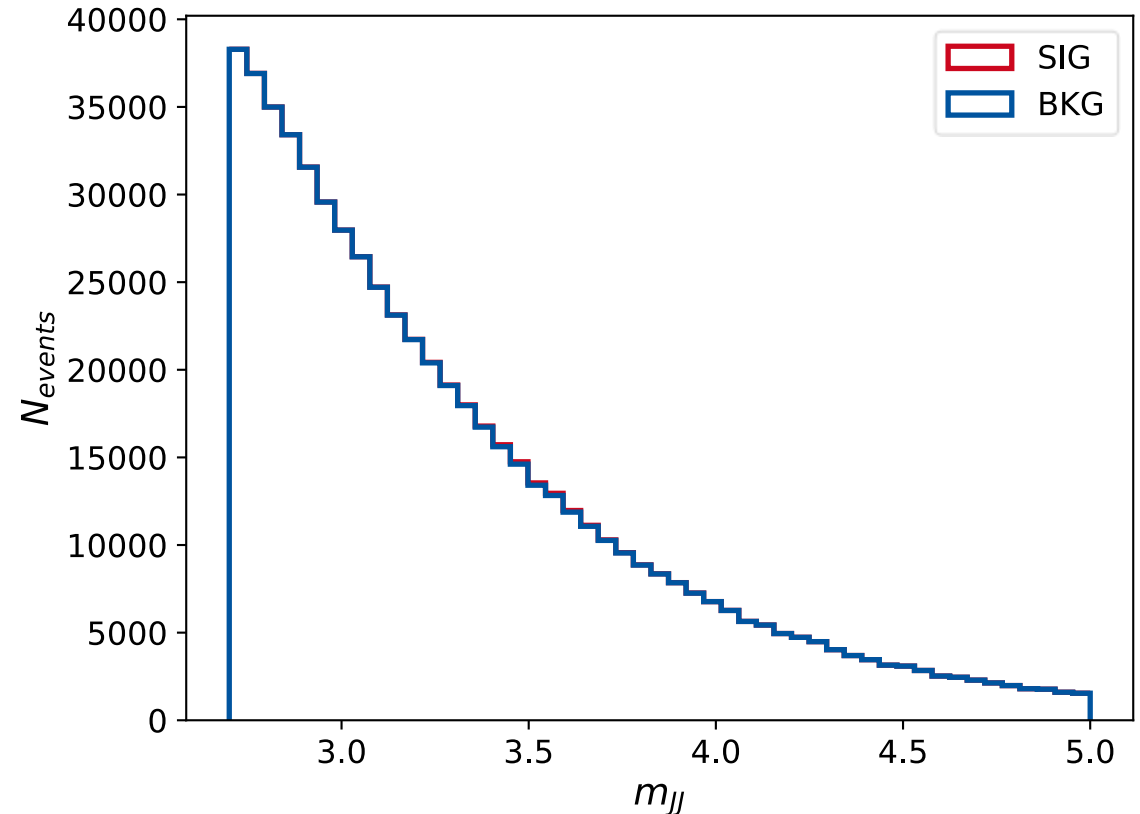


## Pick on simulations

- **Advantages:** metrics that directly access background & signal labels
- **Disadvantages:** less model agnostic, dependent on specific simulation & chosen signals

## Pick on data

- **Advantages:** more model agnostic
- **Disadvantages:** limited number of signal events results in noisy metrics

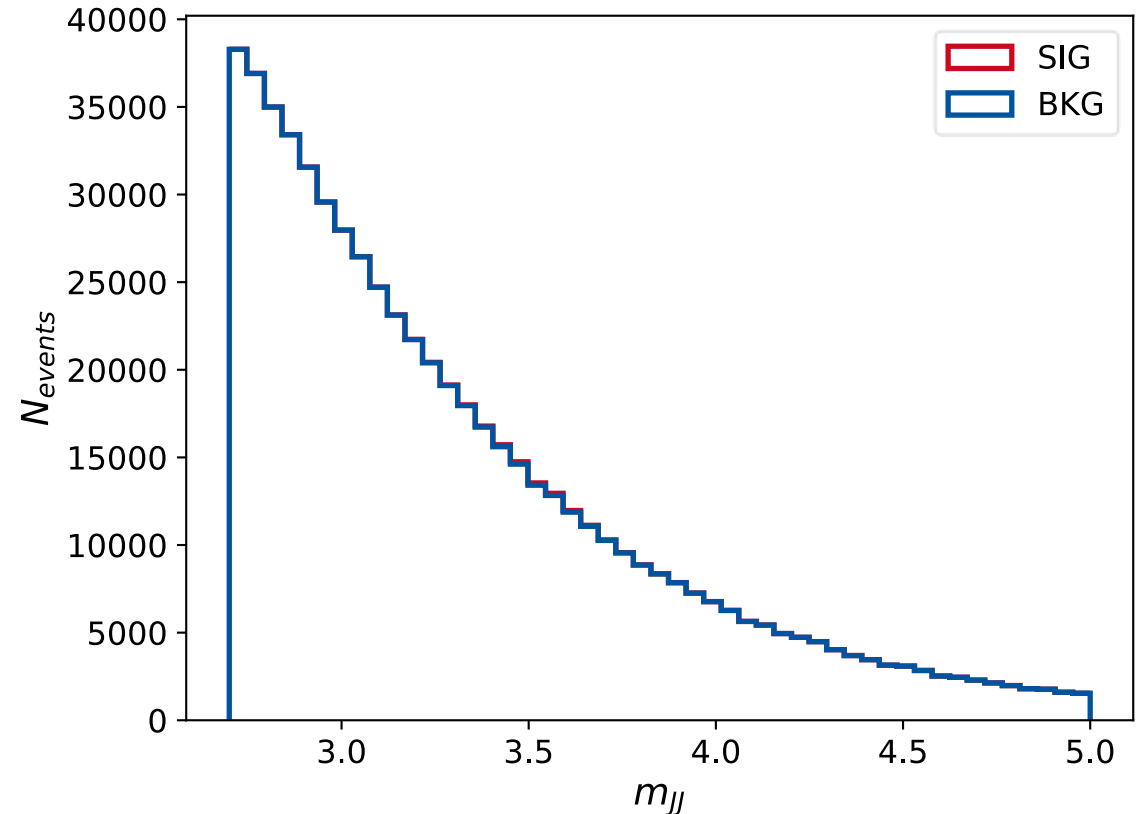


## Pick on simulations

- **Advantages:** metrics that directly access background & signal labels
- **Disadvantages:** less model agnostic, dependent on specific simulation & chosen signals

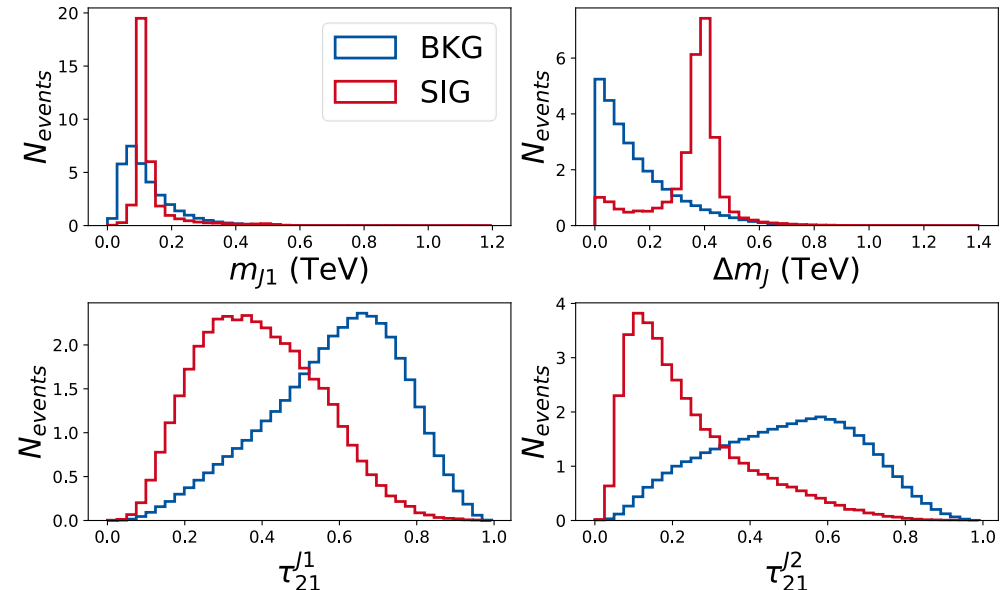
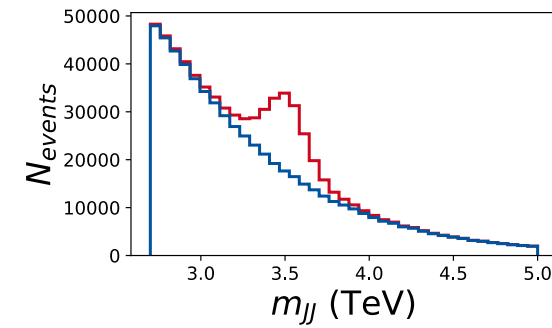
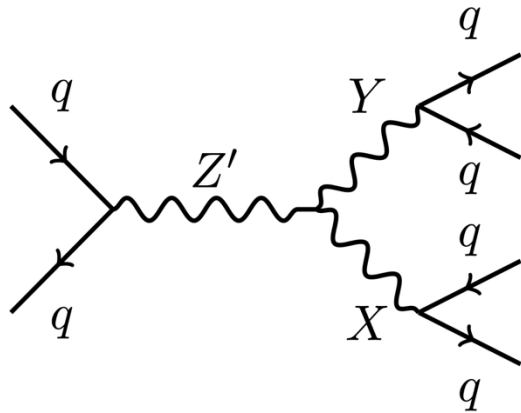
## Pick on data

- **Advantages:** more model agnostic
- **Disadvantages:** limited number of signal events results in noisy metrics
  - Investigate how problematic this noise is

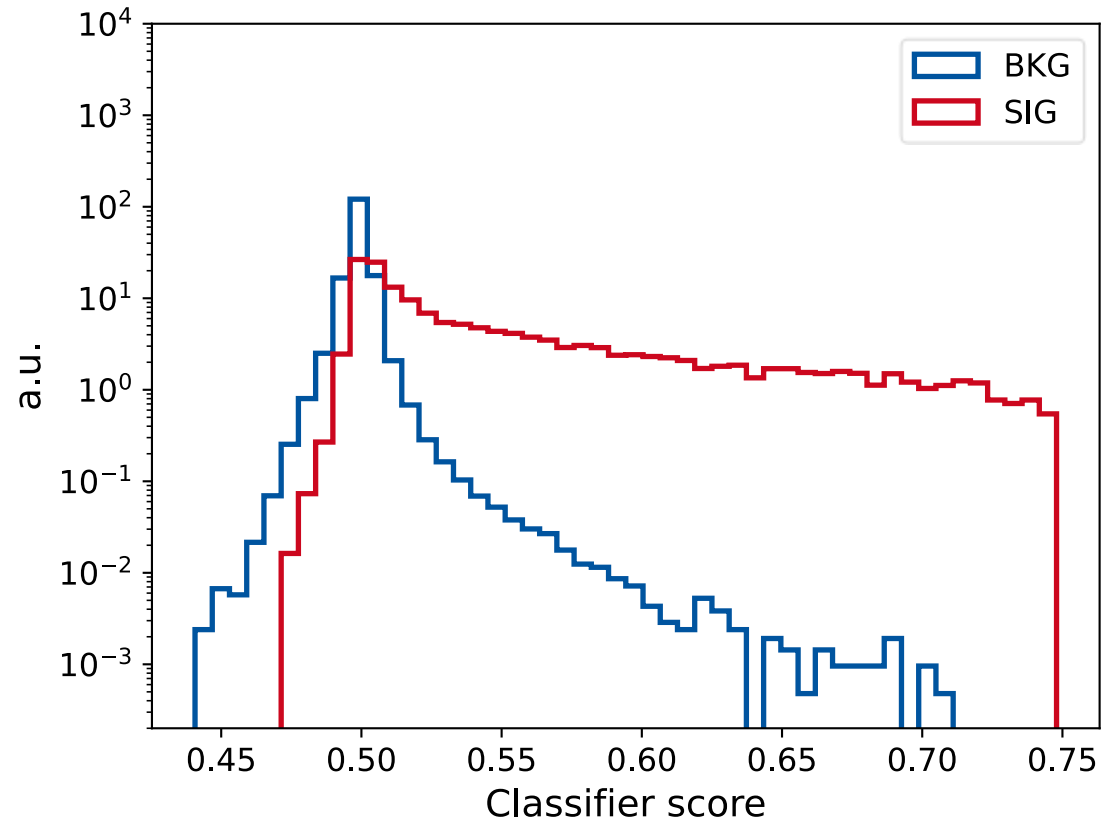


“The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics” [\[2101.08320\]](#), G. Kasieczka, B. Nachman, D. Shih et. al.

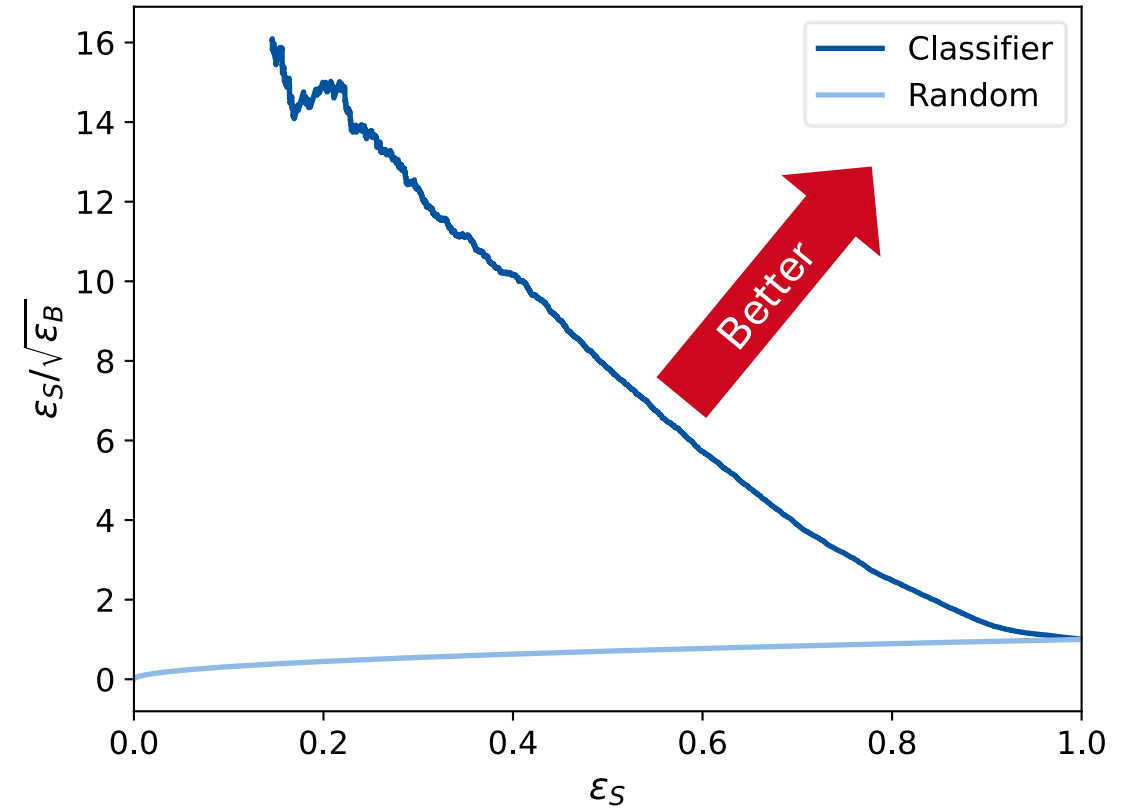
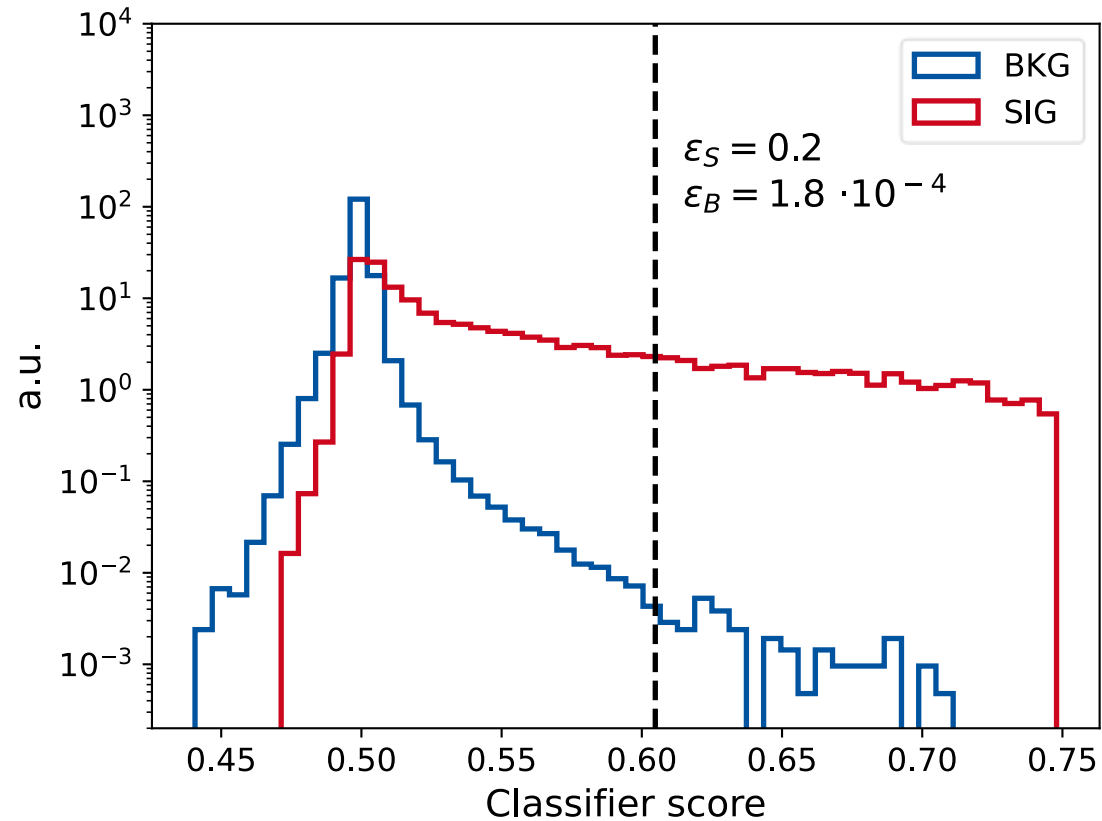
- Benchmark dataset for anomaly detection
- QCD dijet background (1M events)
- Signal ( $N_{sig}$  events)



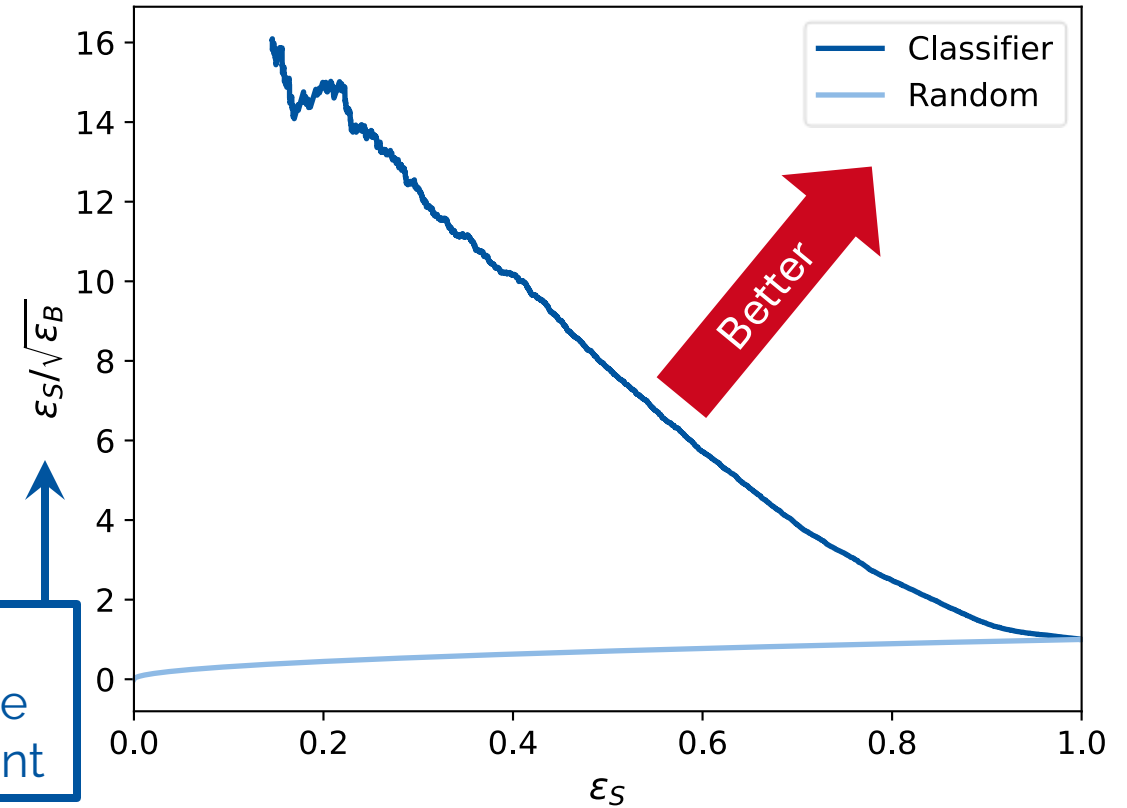
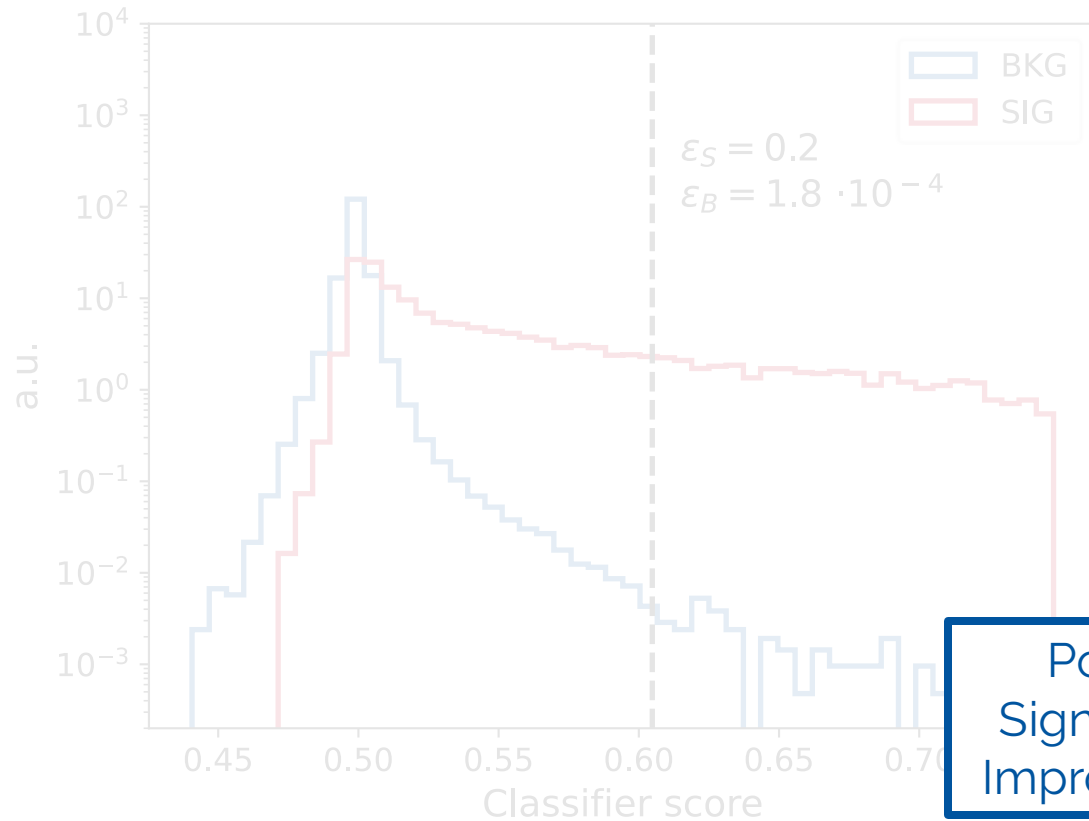
# Supervised metric: Max SIC



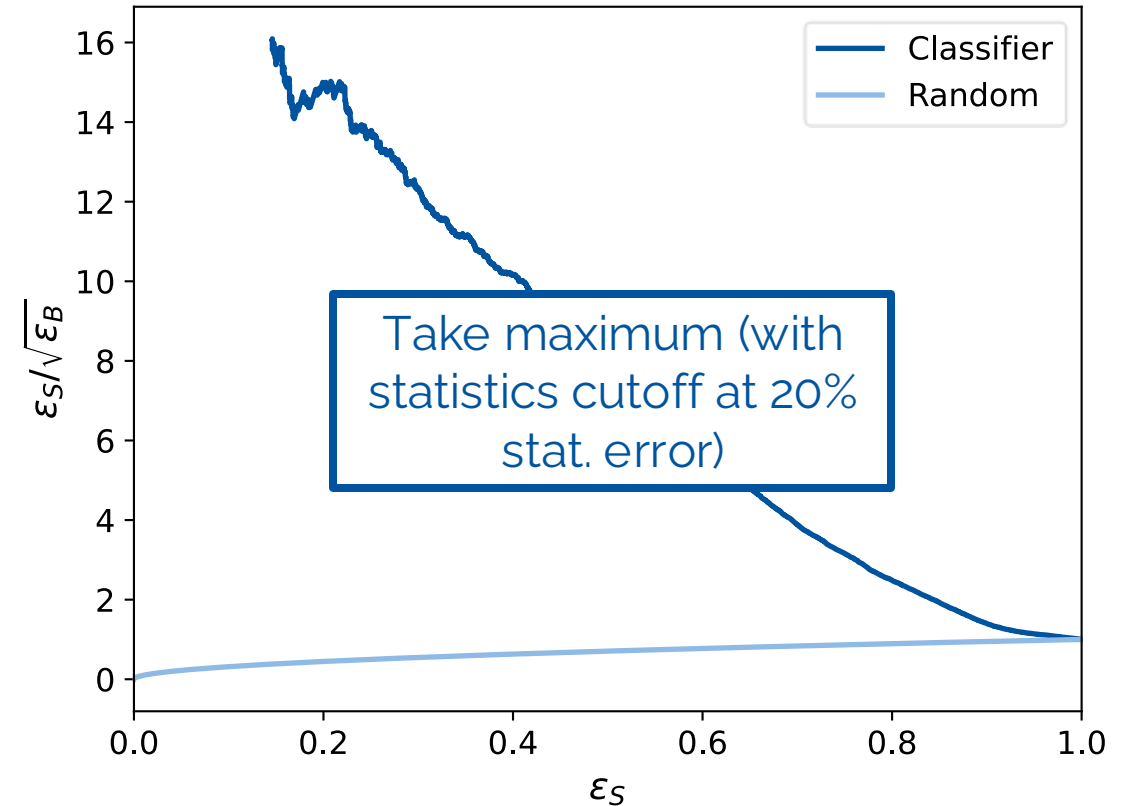
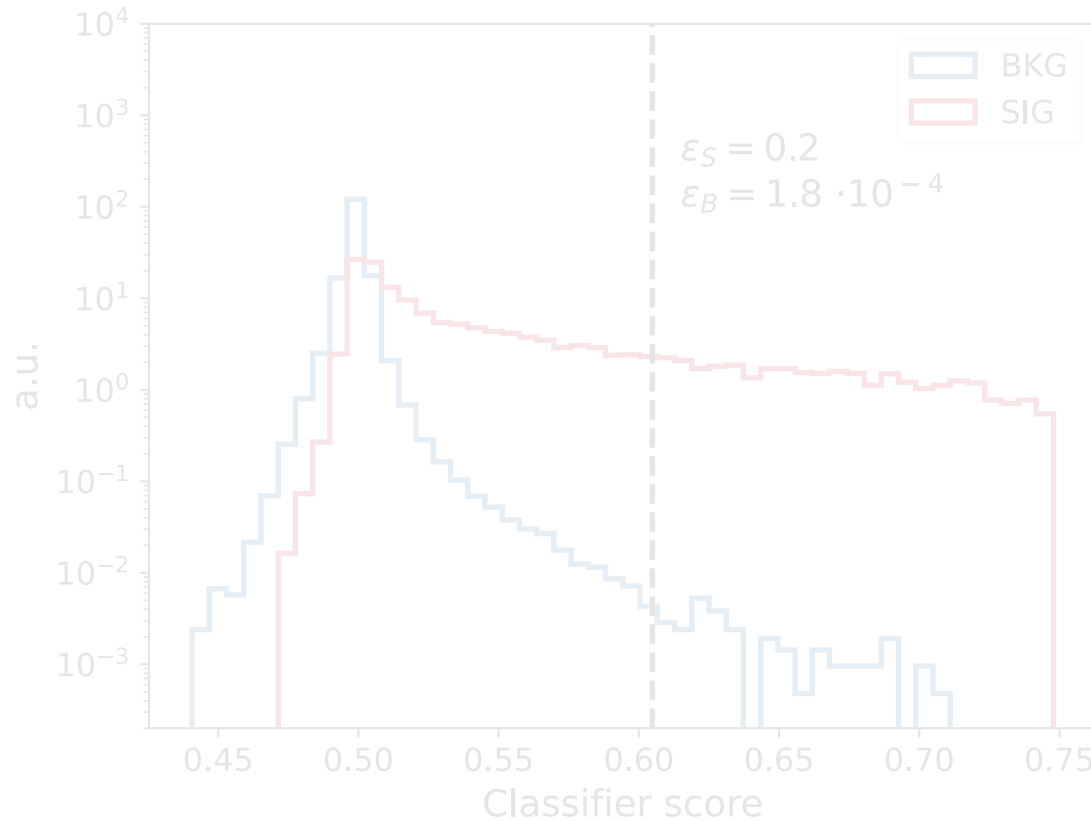
# Supervised metric: Max SIC



# Supervised metric: Max SIC



# Supervised metric: Max SIC



$$BCE = -\log p_{\text{pred, true}}$$

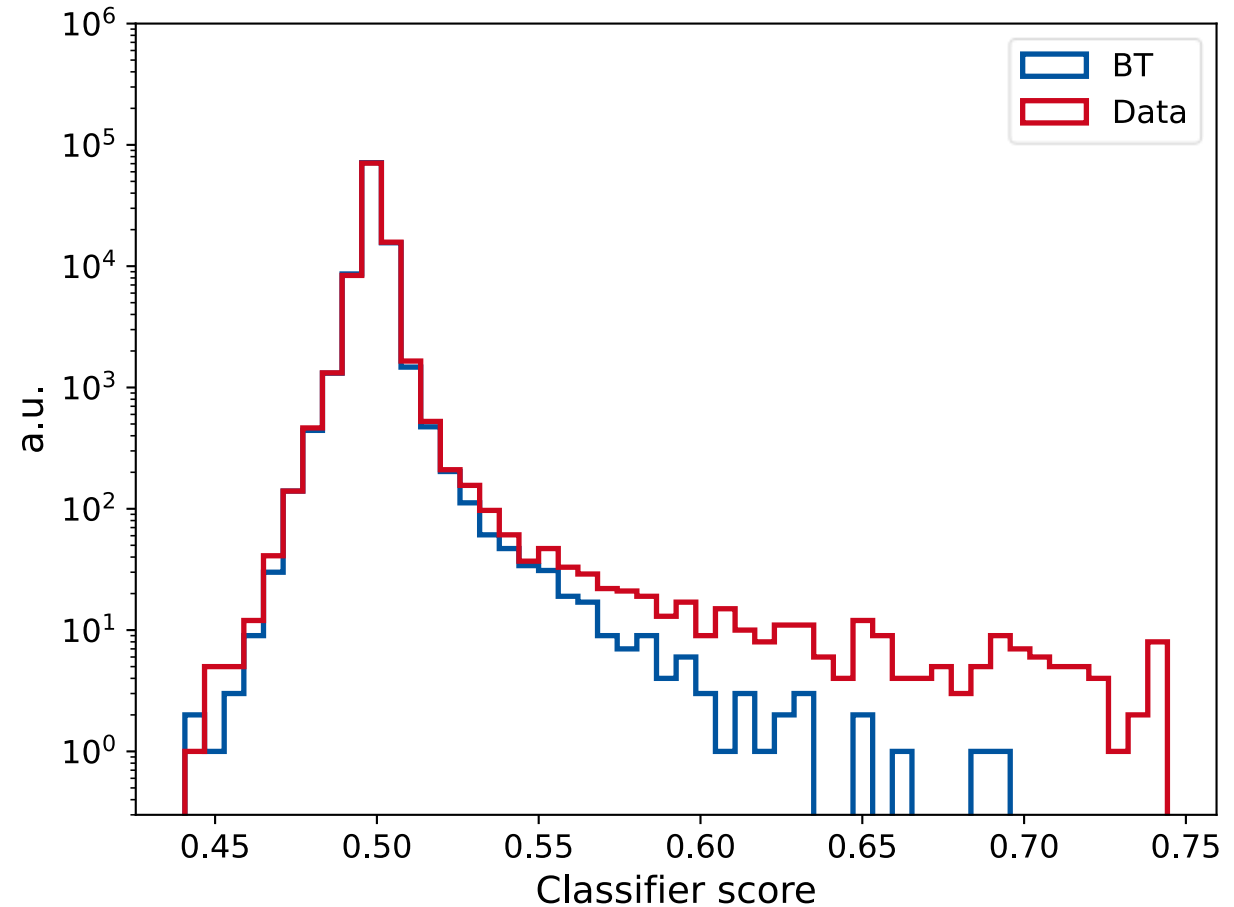
- Random scores: all events  $p_{\text{pred}} = 0.5$

$$BCE_{\text{random}} = \ln 2 \approx 0.6931$$

- Optimal scores:  $p_{\text{pred}, B} = 0.5$ ,  $p_{\text{pred}, S} = 1$ , e.g.

$$BCE_{\text{opt}} = \frac{10^5 - 1000}{10^5} \ln 2 + \frac{1000}{10^5} \ln 1 \approx 0.6924$$

→ Dominated by background (noisy)





$$BCE = -\log p_{\text{pred, true}}$$

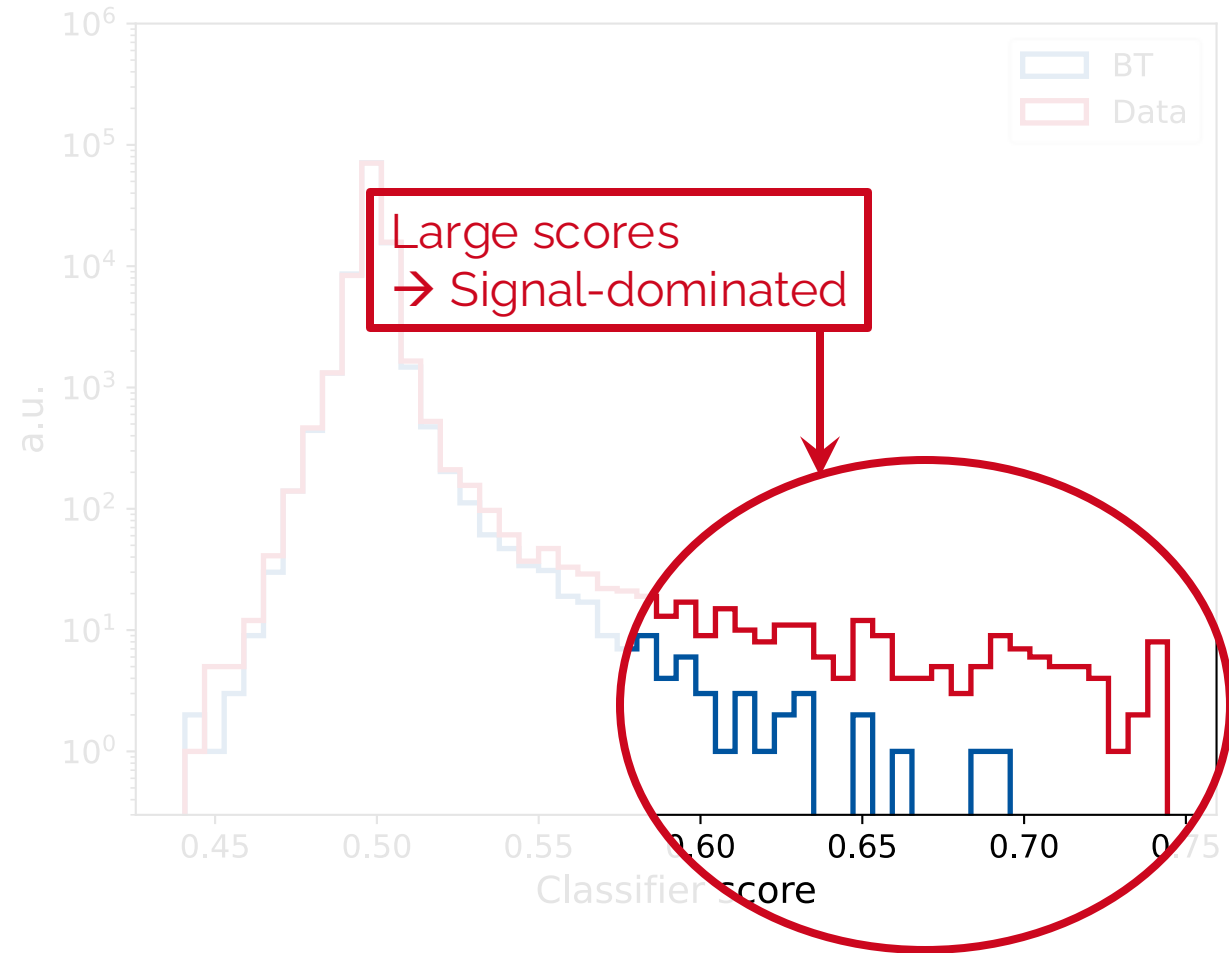
- Random scores: all events  $p_{\text{pred}} = 0.5$

$$BCE_{\text{random}} = \ln 2 \approx 0.6931$$

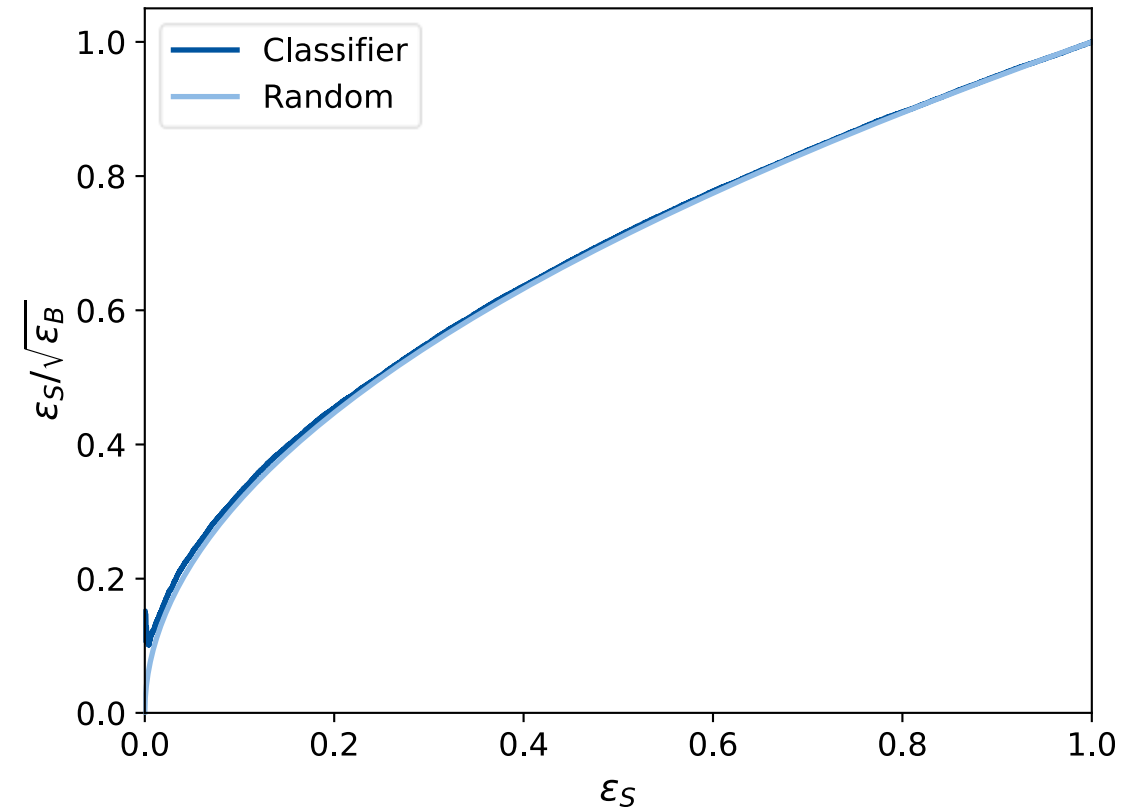
- Optimal scores:  $p_{\text{pred}, B} = 0.5$ ,  $p_{\text{pred}, S} = 1$ , e.g.

$$BCE_{\text{opt}} = \frac{10^5 - 1000}{10^5} \ln 2 + \frac{1000}{10^5} \ln 1 \approx 0.6924$$

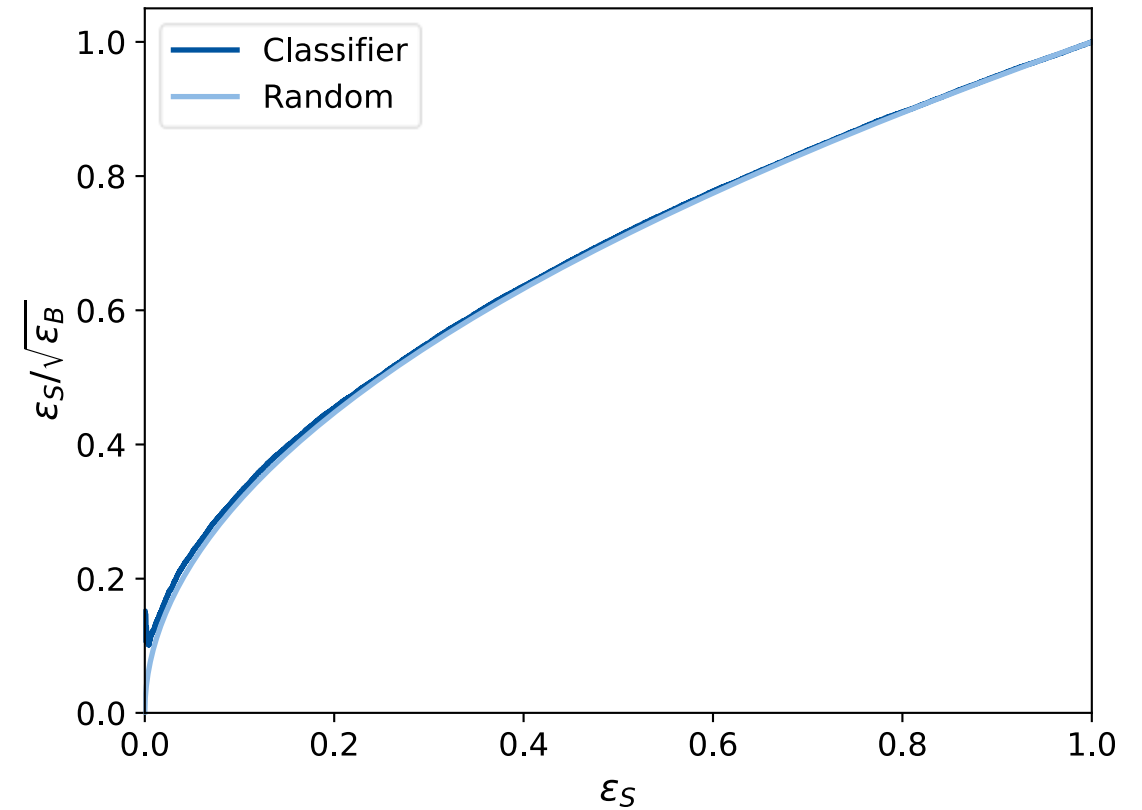
→ Dominated by background (noisy)



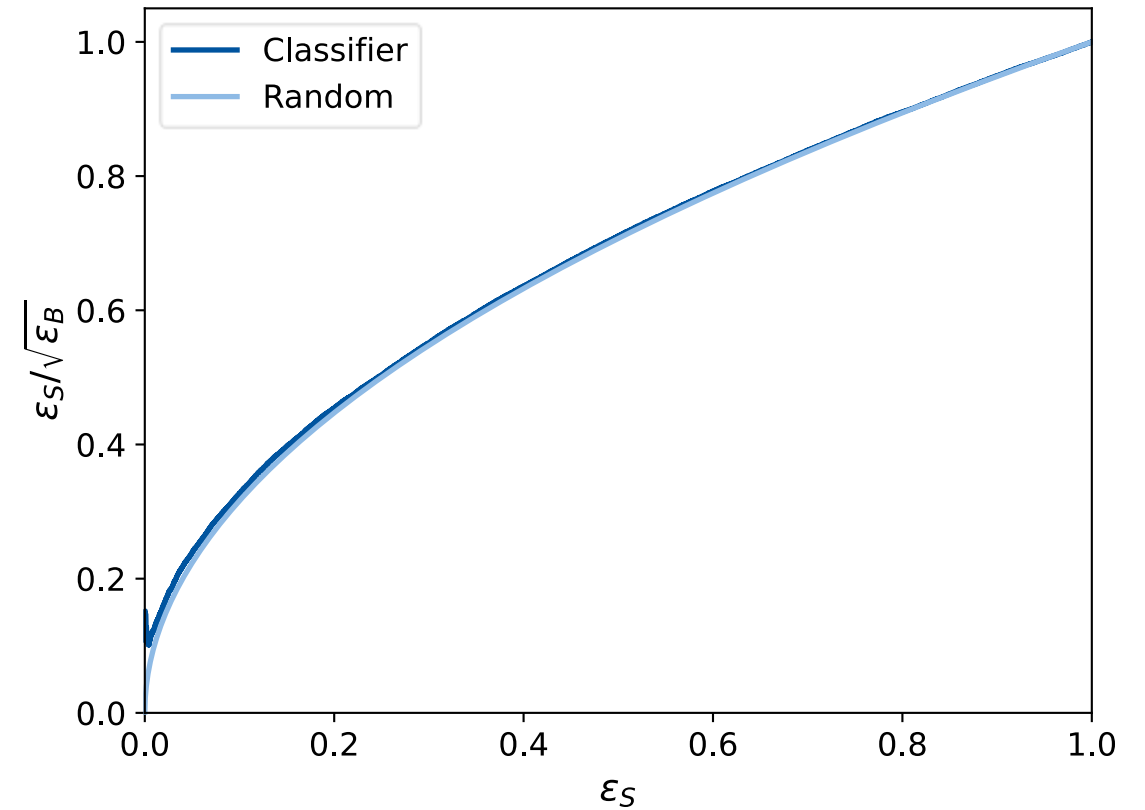
- SIC curve calculated at all classifier scores  
→ Can pick large threshold



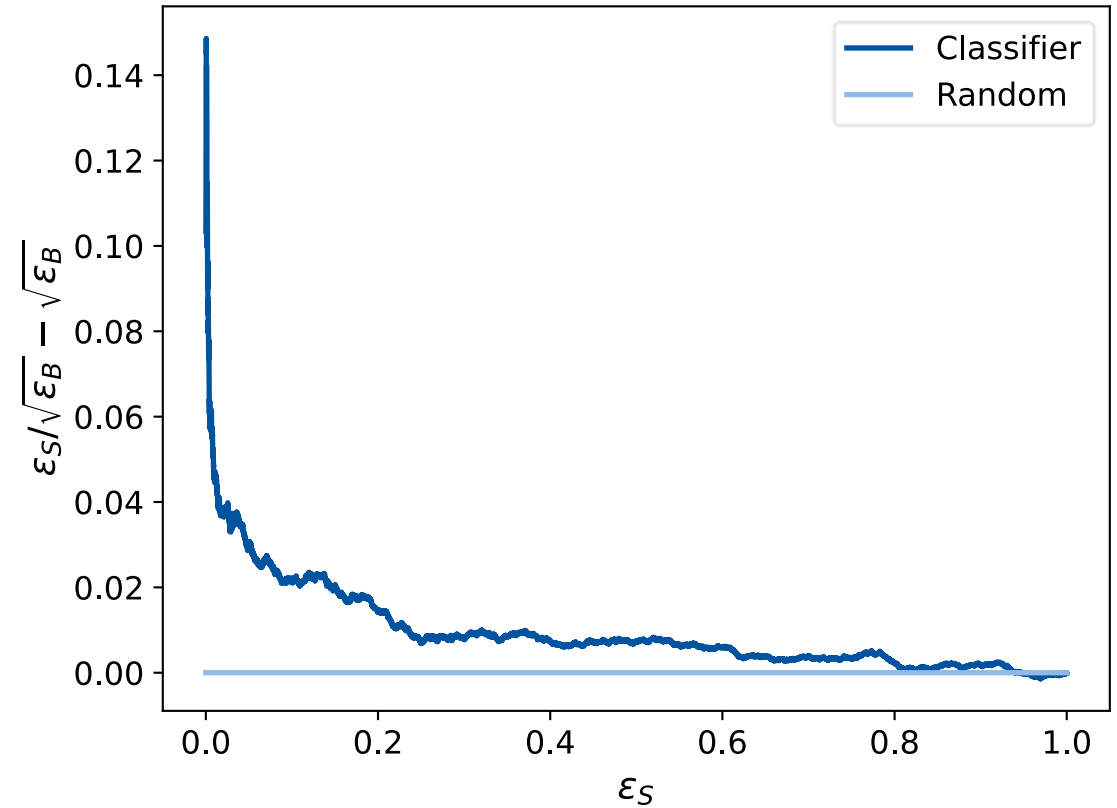
- SIC curve calculated at all classifier scores  
→ Can pick large threshold



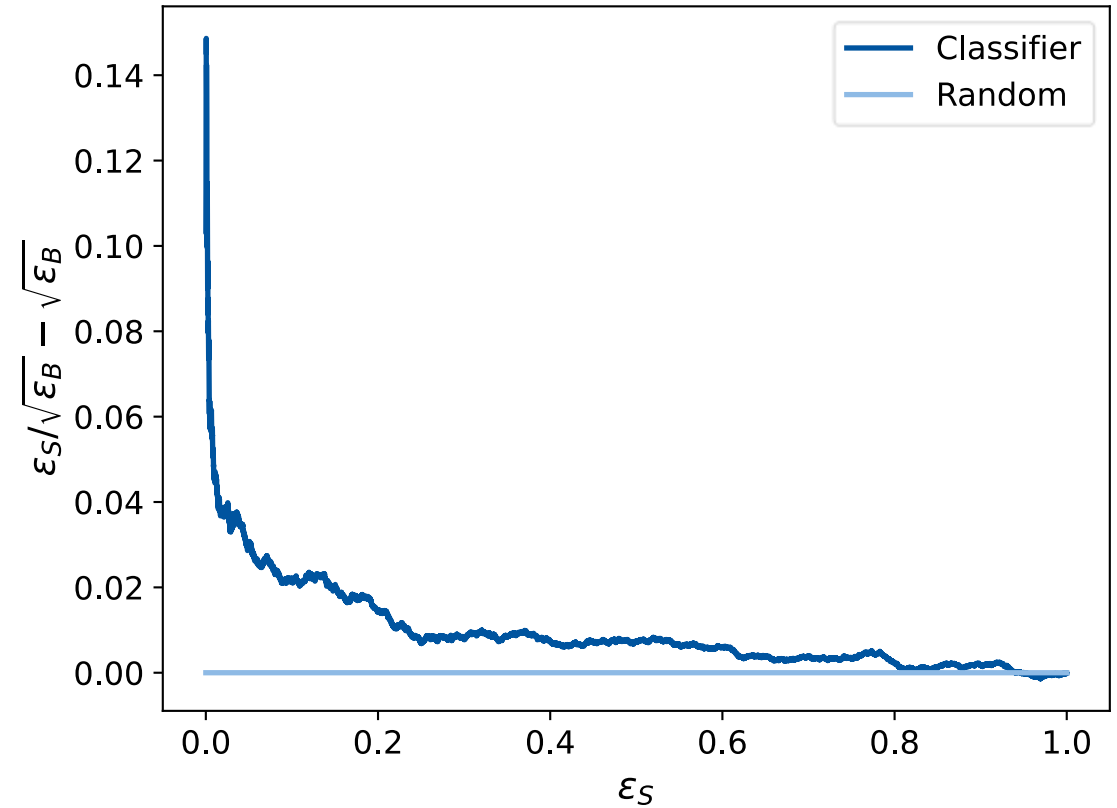
- SIC curve calculated at all classifier scores  
→ Can pick large threshold



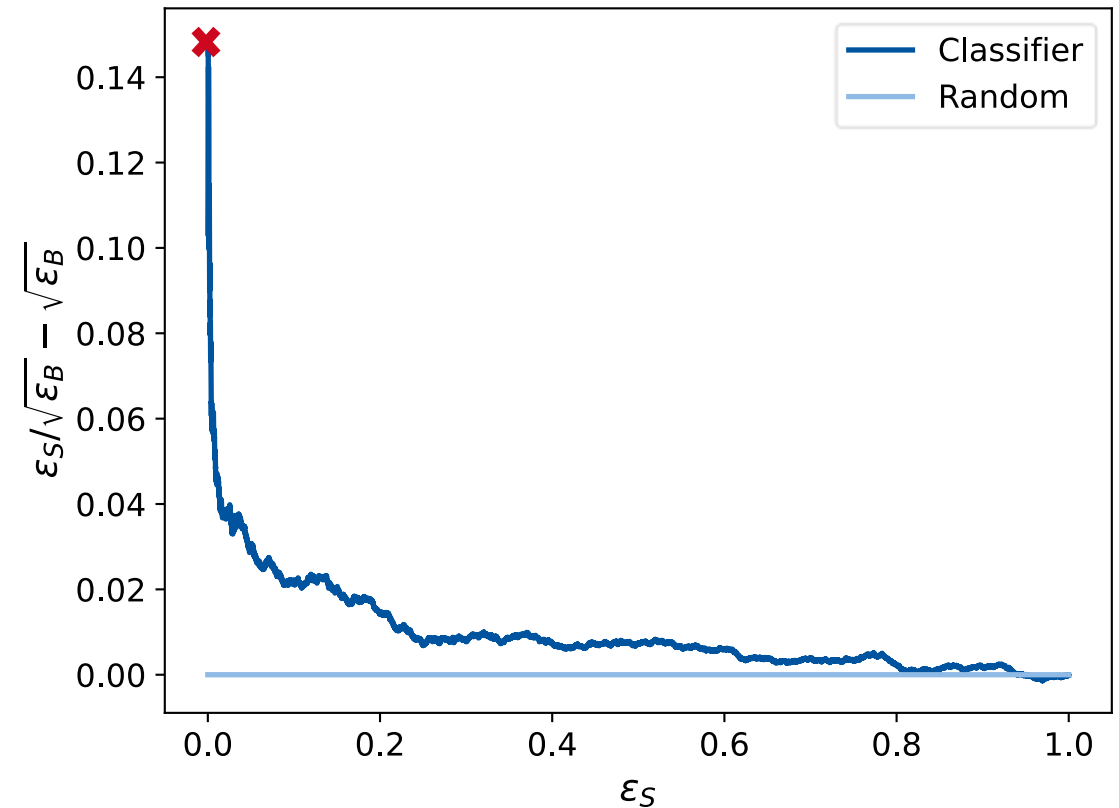
- SIC curve calculated at all classifier scores
  - Can pick large threshold
- Val SIC very close to random
  - Subtract random



- SIC curve calculated at all classifier scores
  - Can pick large threshold
- Val SIC very close to random
  - Subtract random
- For most signal-like events, pick maximum

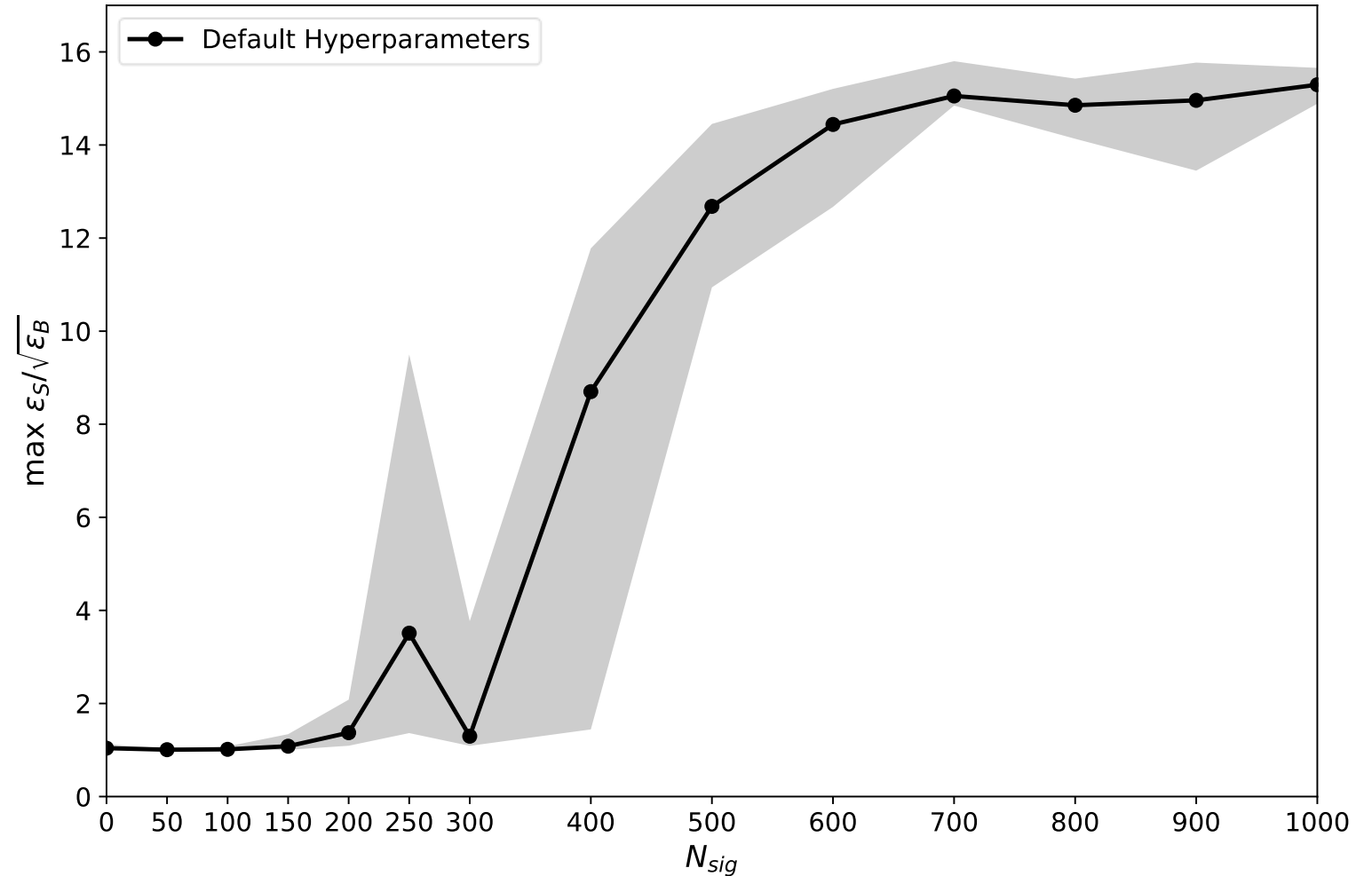


- SIC curve calculated at all classifier scores
  - Can pick large threshold
- Val SIC very close to random
  - Subtract random
- For most signal-like events, pick maximum
  - Will refer to this as **Val SIC**



## Idea:

- Test HP configurations and pick best based on each metric



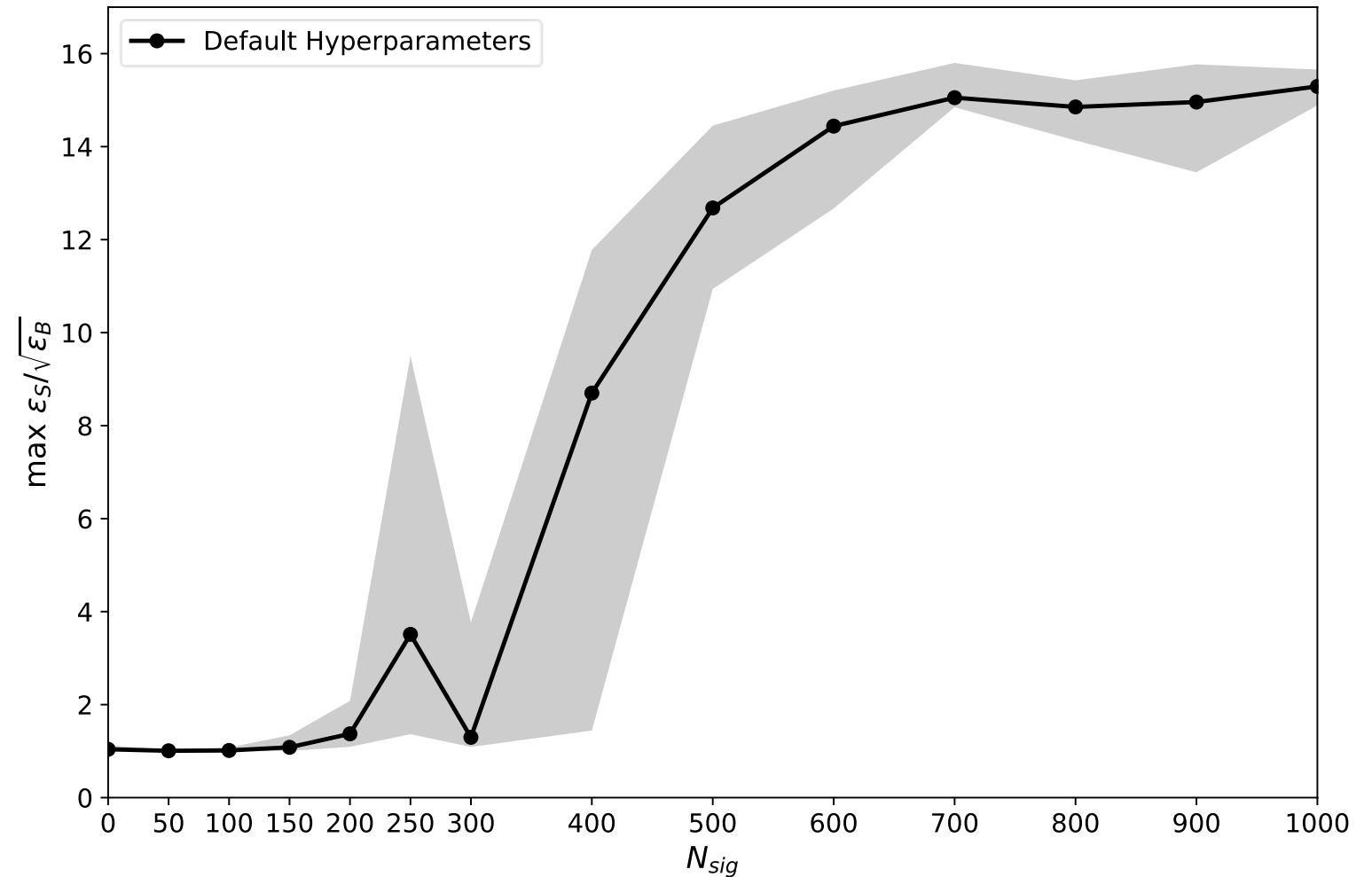


## Idea:

- Test HP configurations and pick best based on each metric

## Results:

- Benchmark: Default HP optimized for this setup

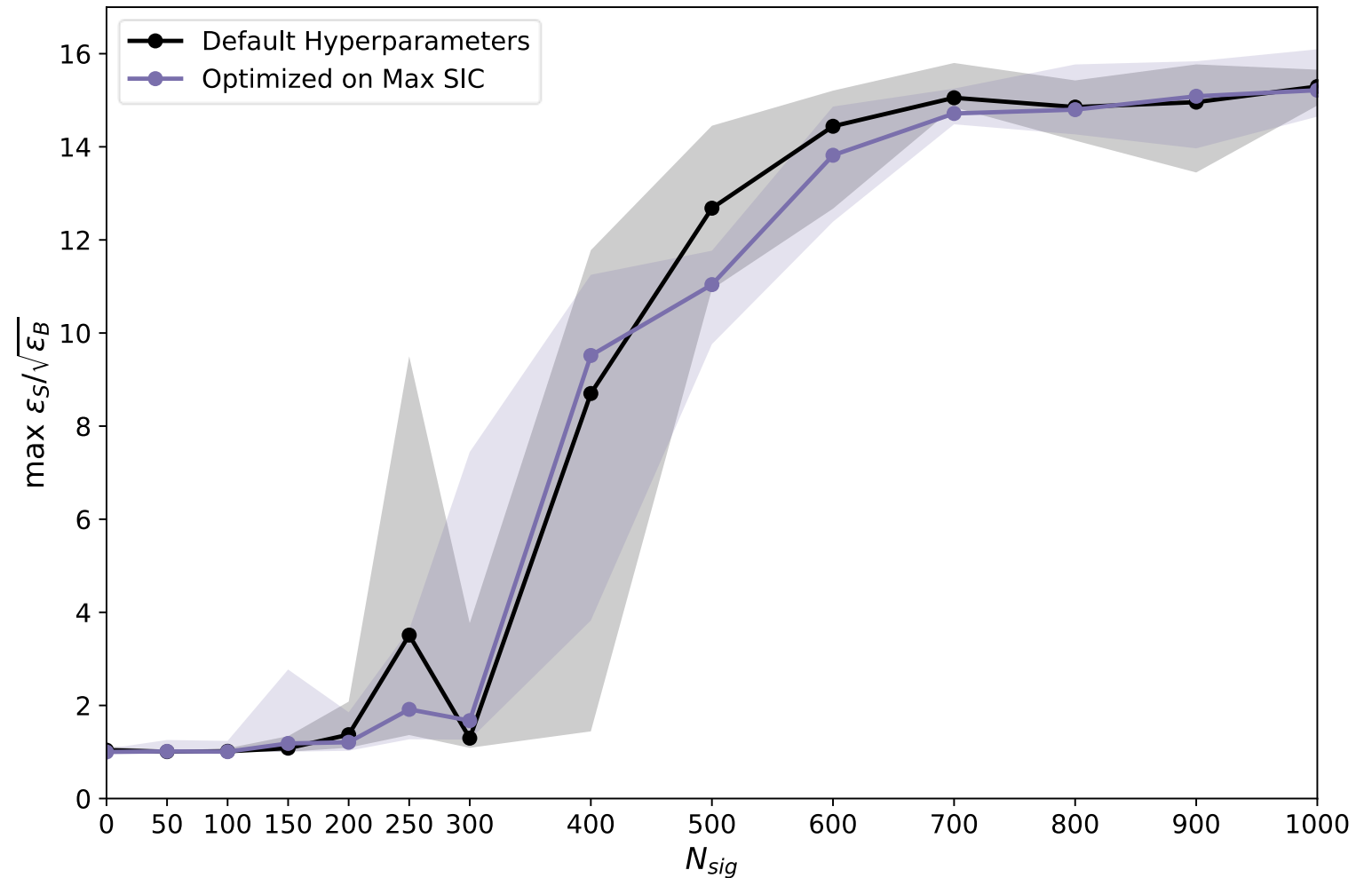


## Idea:

- Test HP configurations and pick best based on each metric

## Results:

- Benchmark: Default HP optimized for this setup
- Max SIC: Performance comparable to benchmark

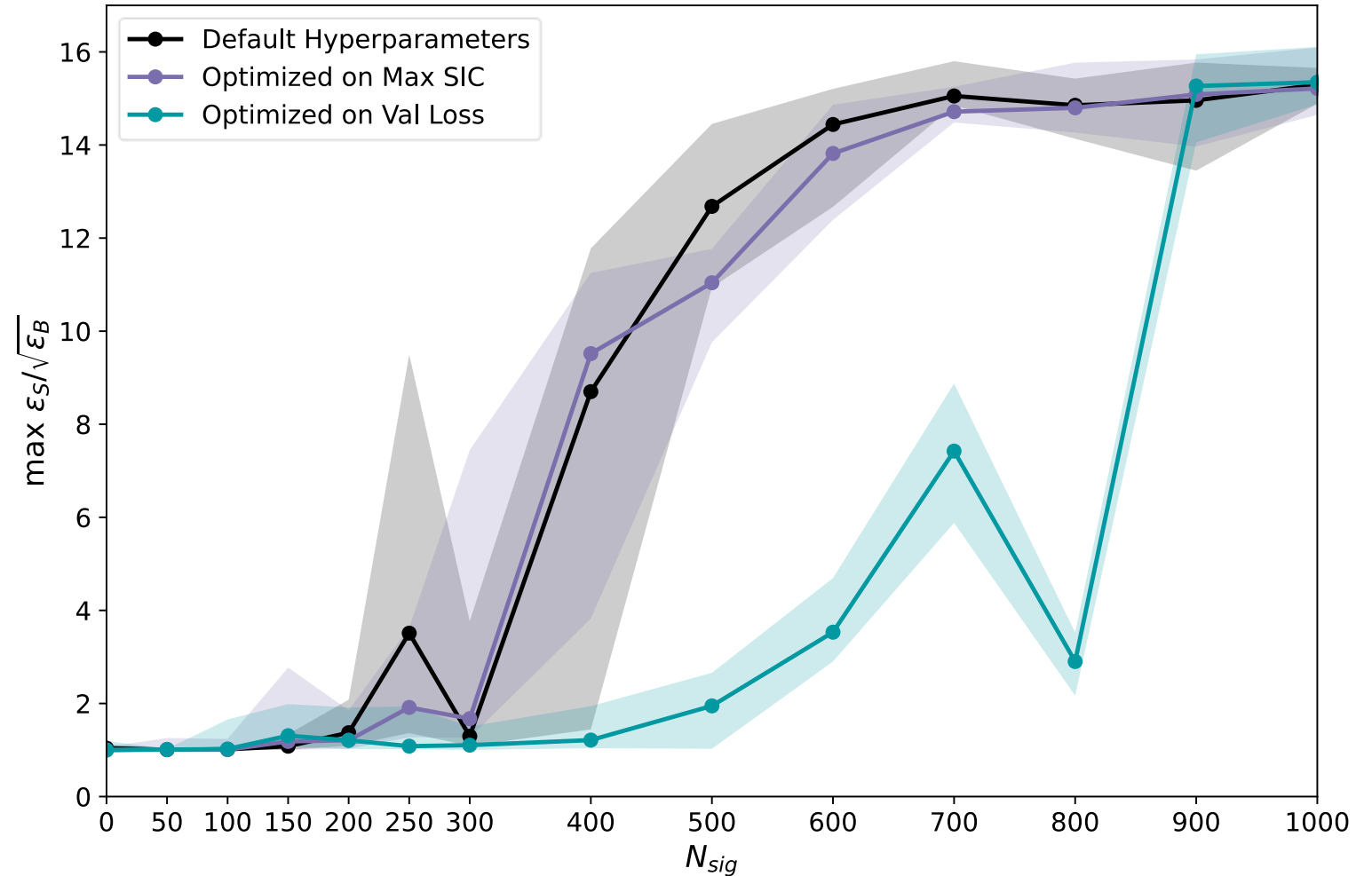


## Idea:

- Test HP configurations and pick best based on each metric

## Results:

- Benchmark: Default HP optimized for this setup
- Max SIC: Performance comparable to benchmark
- Val loss: fails at low  $N_{sig}$

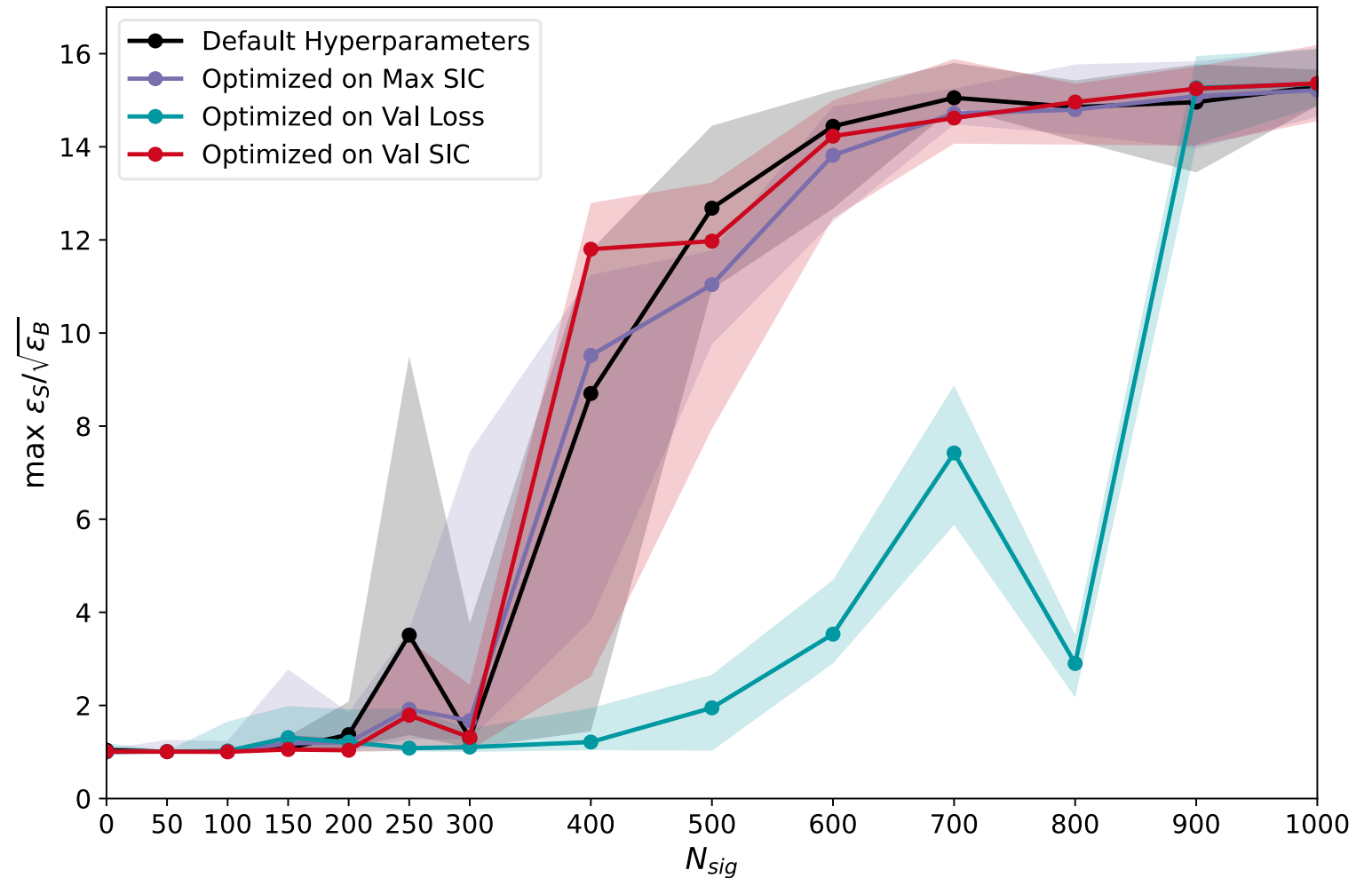


## Idea:

- Test HP configurations and pick best based on each metric

## Results:

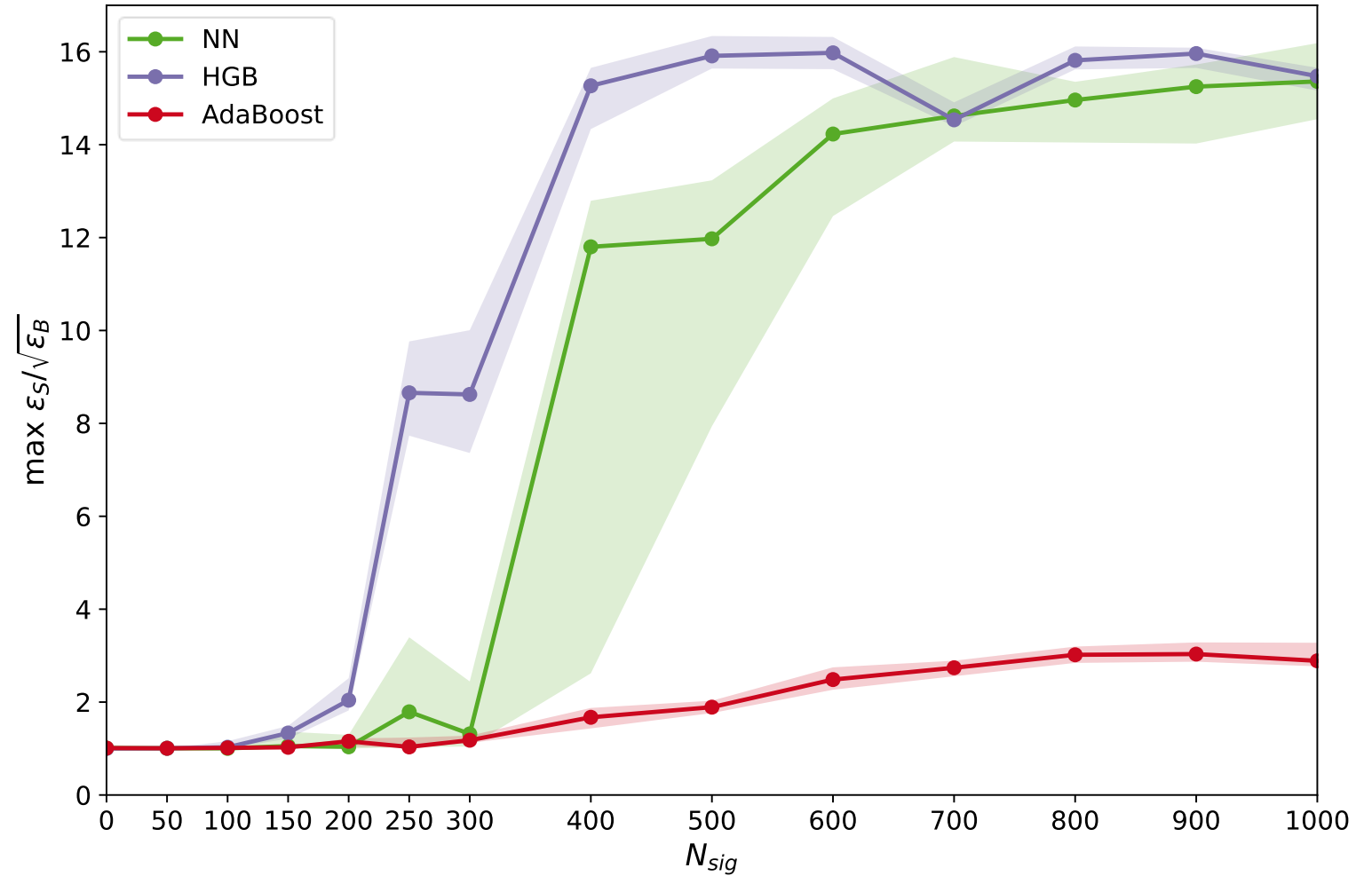
- Benchmark: Default HP optimized for this setup
- Max SIC: Performance comparable to benchmark
- Val loss: fails at low  $N_{sig}$
- Val SIC: Performance close to max SIC & benchmark



# Picking the setup

## Procedure:

- Pick several architectures
- Optimize HP for each
- Train each setup on  $\frac{1}{2}$  data, evaluate metrics on  $\frac{1}{2}$  data
- Pick best setup based on metric



# Picking the setup

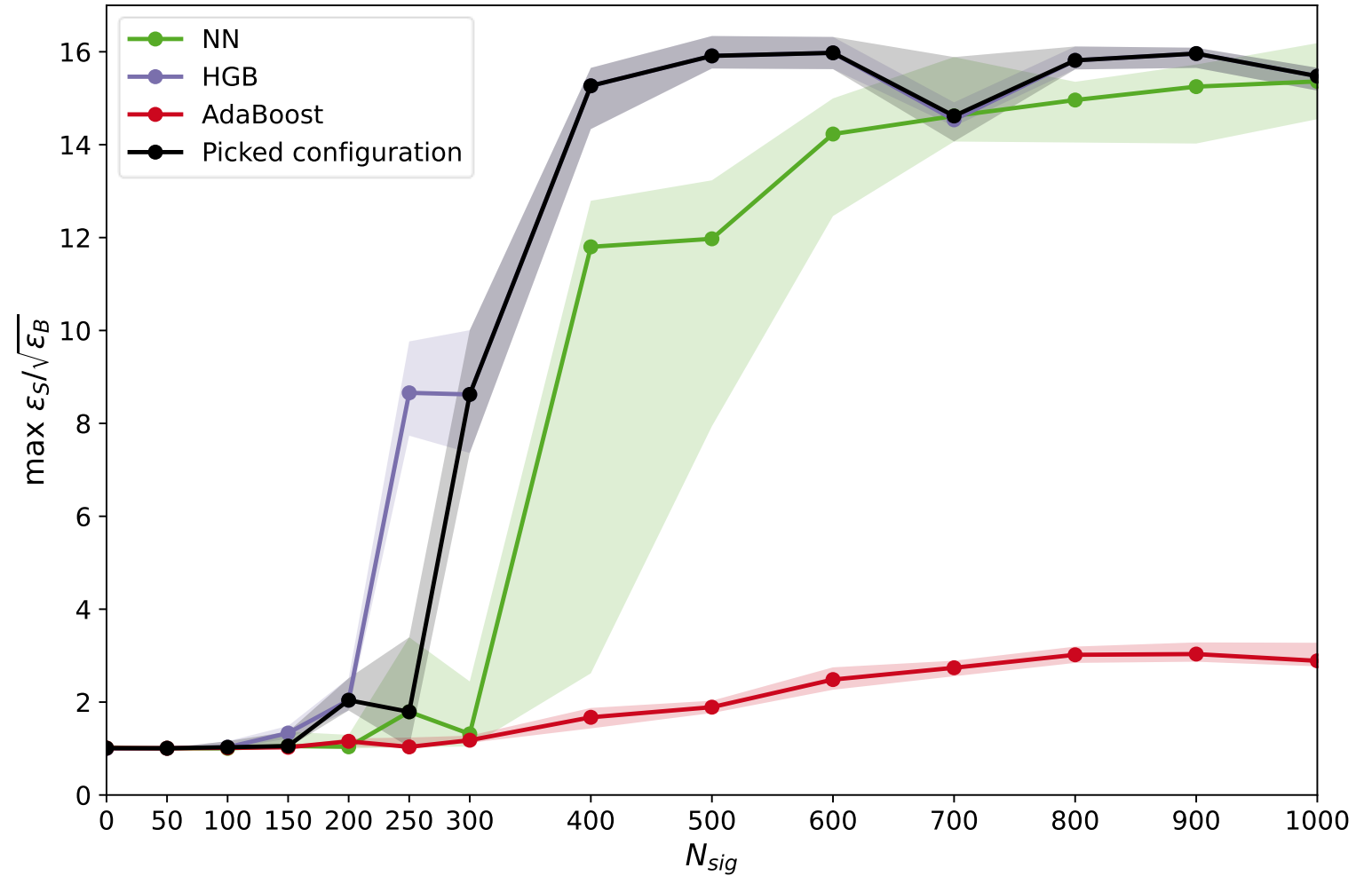
## Procedure:

- Pick several architectures
- Optimize HP for each
- Train each setup on  $\frac{1}{2}$  data, evaluate metrics on  $\frac{1}{2}$  data
- Pick best setup based on metric

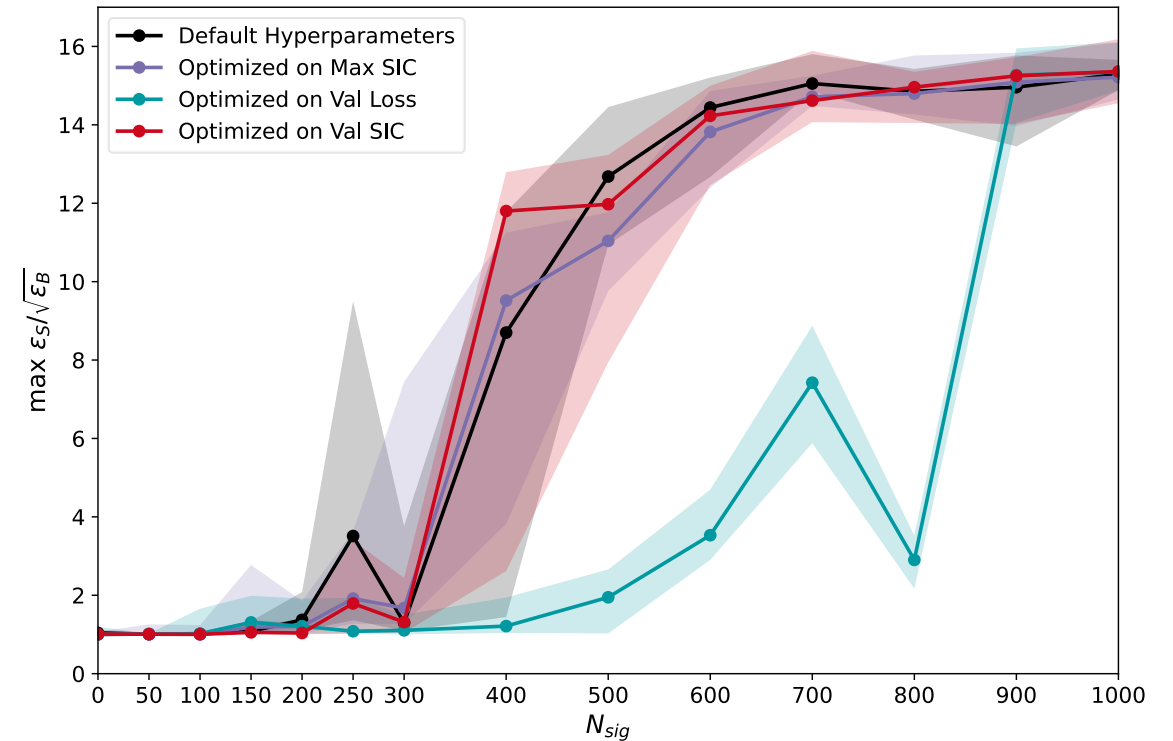
## Result:

Val SIC picks best model almost everywhere

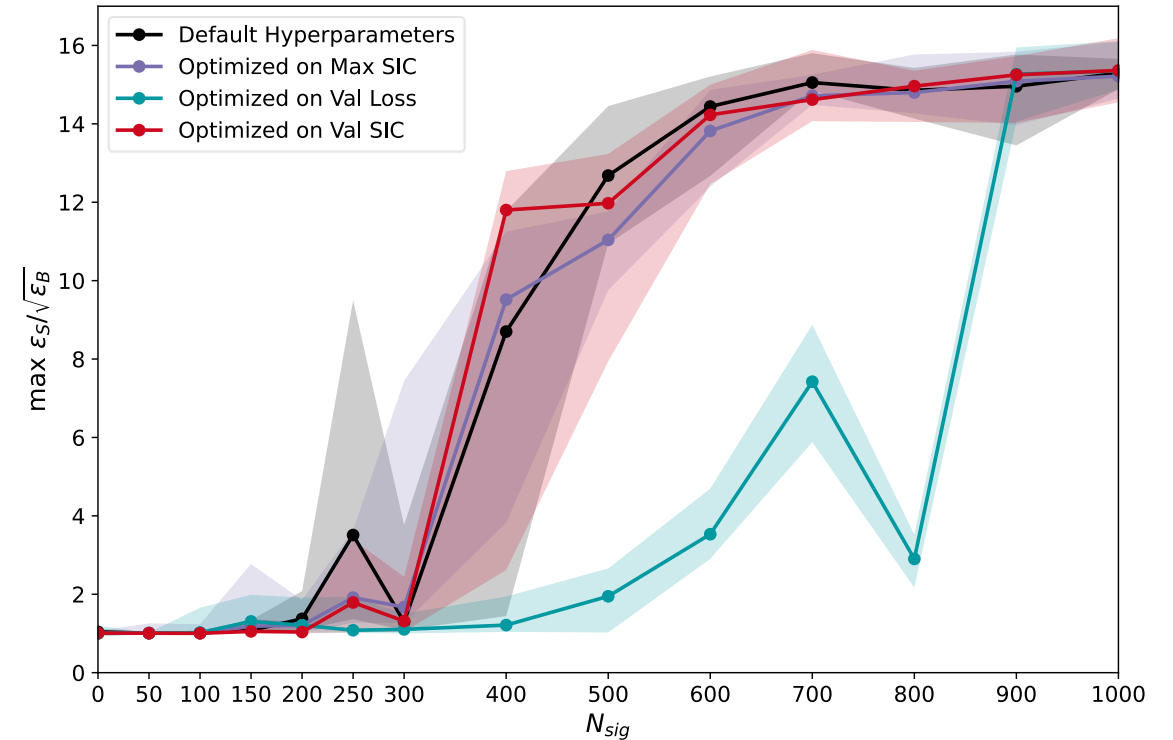
→ Fails at very low signal injections (expected)



- Investigated two data-driven metrics for model agnostic setup optimization

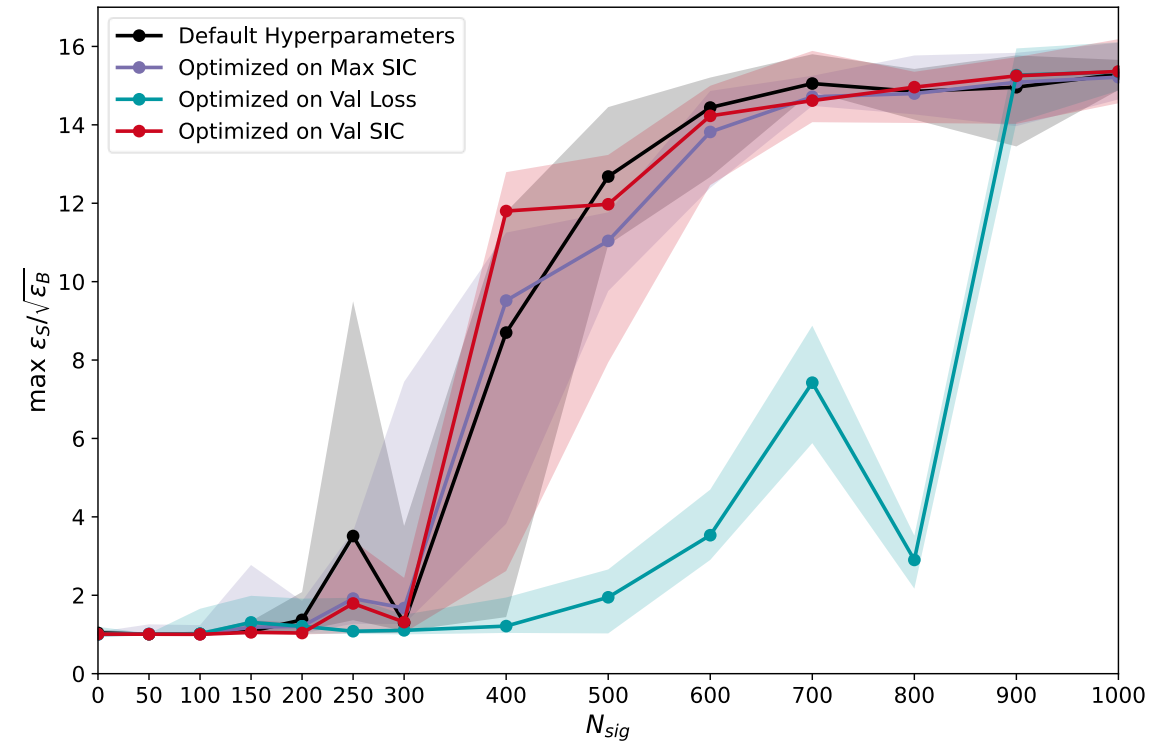


- Investigated two data-driven metrics for model agnostic setup optimization
  - Val loss: too sensitive to background distribution

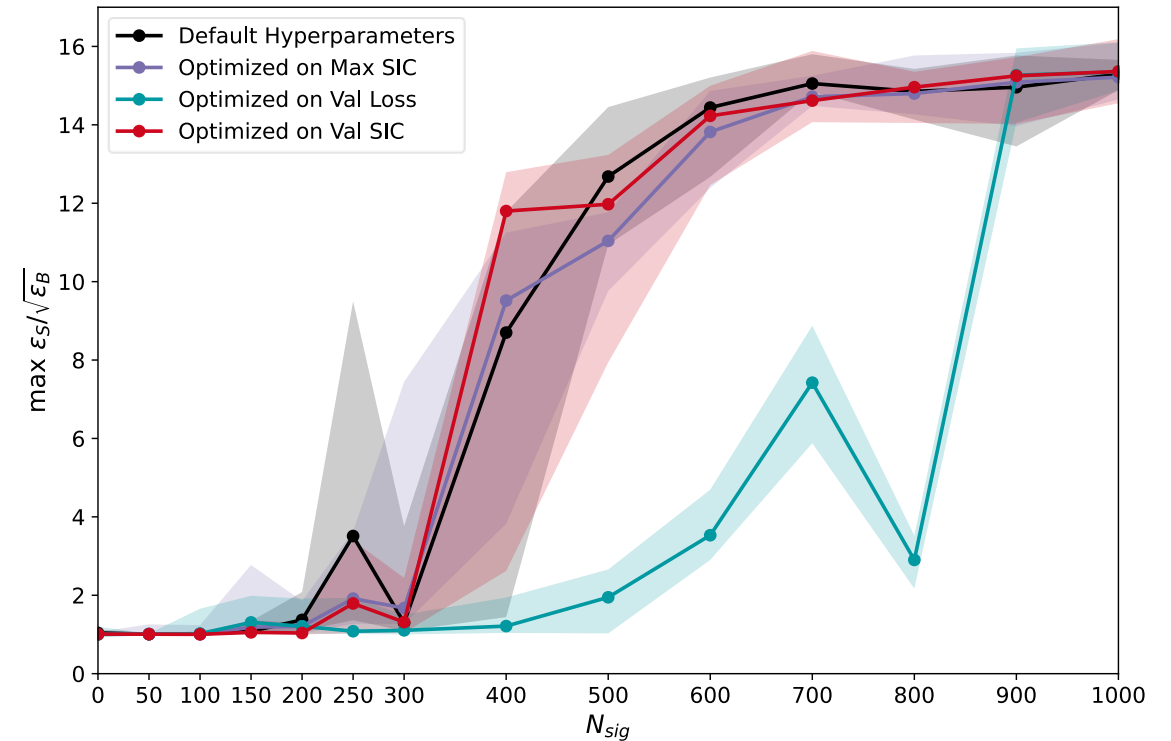




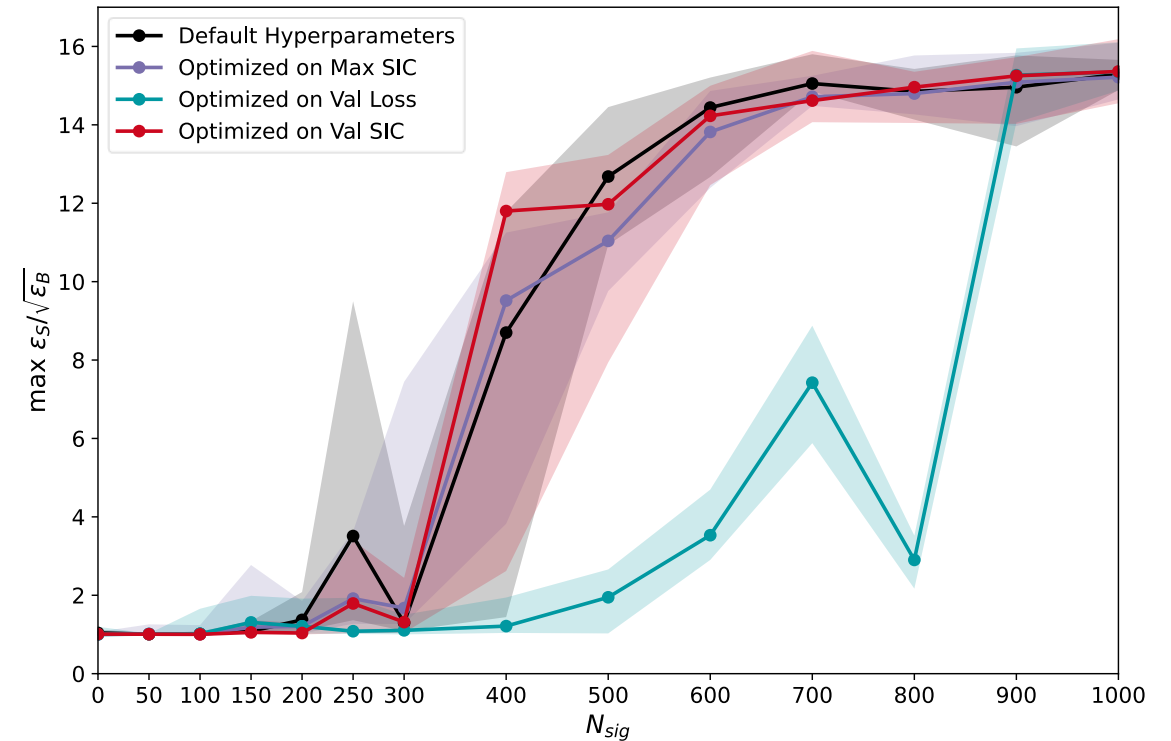
- Investigated two data-driven metrics for model agnostic setup optimization
  - Val loss: too sensitive to background distribution
  - Val SIC: highly correlated with max SIC by focusing on most signal-like events



- Investigated two data-driven metrics for model agnostic setup optimization
  - Val loss: too sensitive to background distribution
  - Val SIC: highly correlated with max SIC by focusing on most signal-like events
- Optimizing on val SIC leads to excellent anomaly detection performance

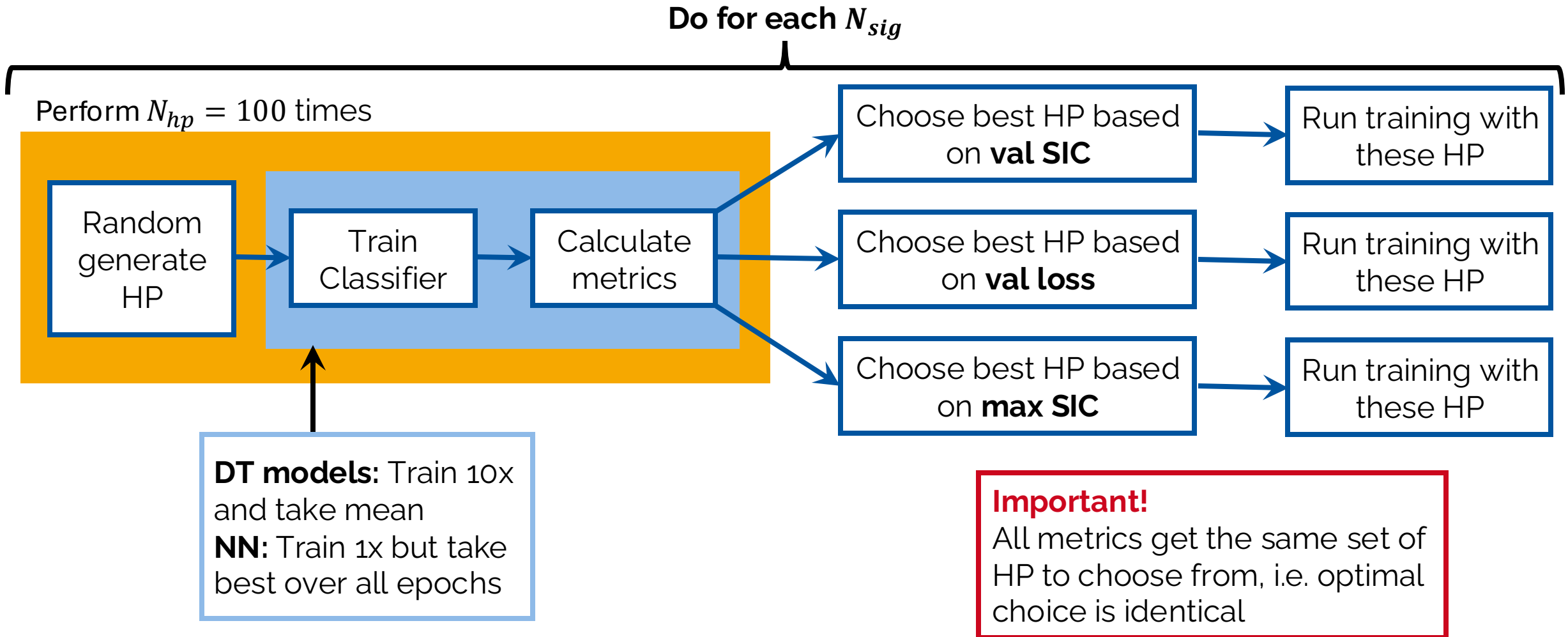


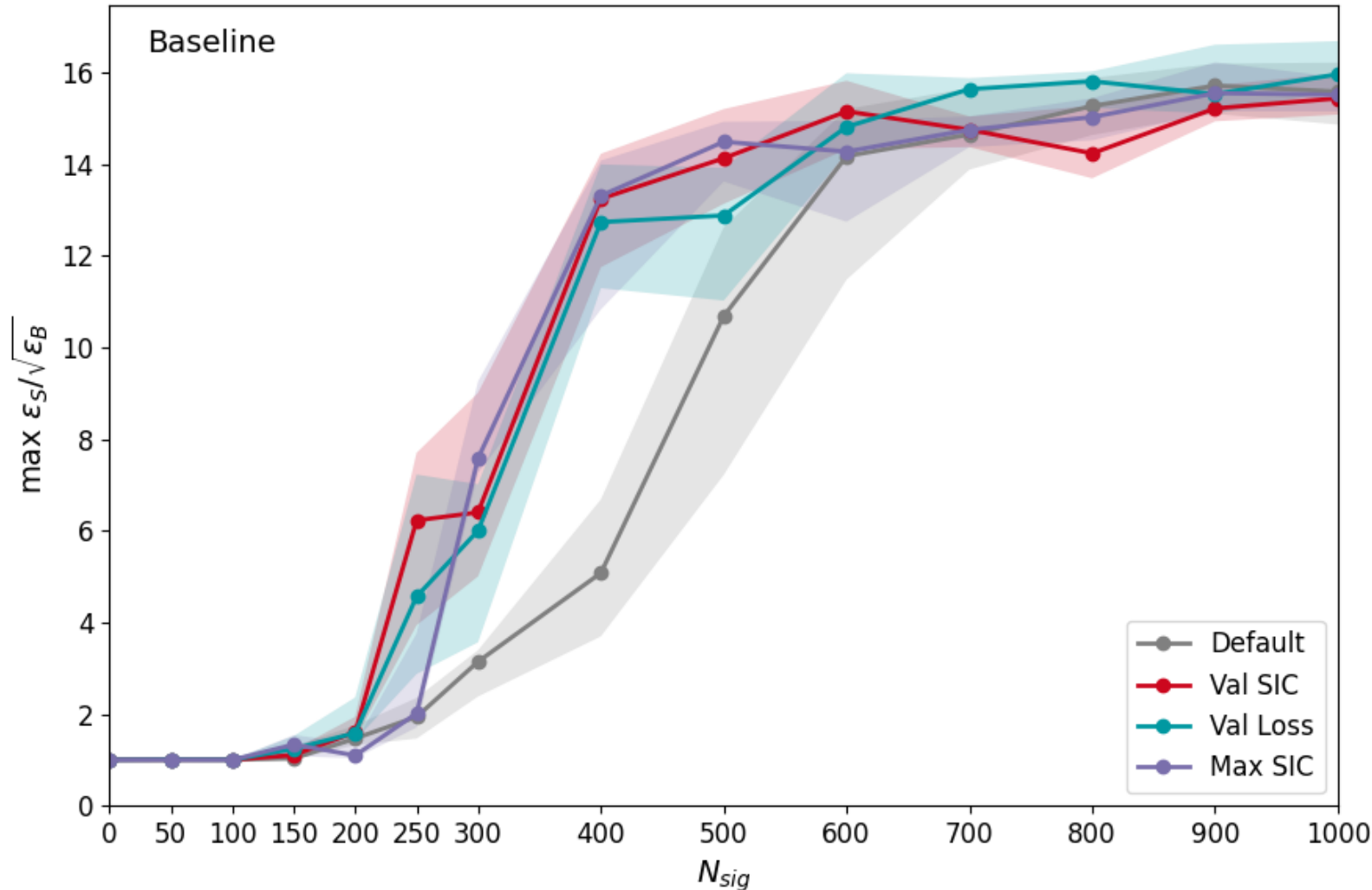
- Investigated two data-driven metrics for model agnostic setup optimization
  - Val loss: too sensitive to background distribution
  - Val SIC: highly correlated with max SIC by focusing on most signal-like events
- Optimizing on val SIC leads to excellent anomaly detection performance
- Behind the scenes: Seen comparable results for CWoLa Hunting and CATHODE



# Backup

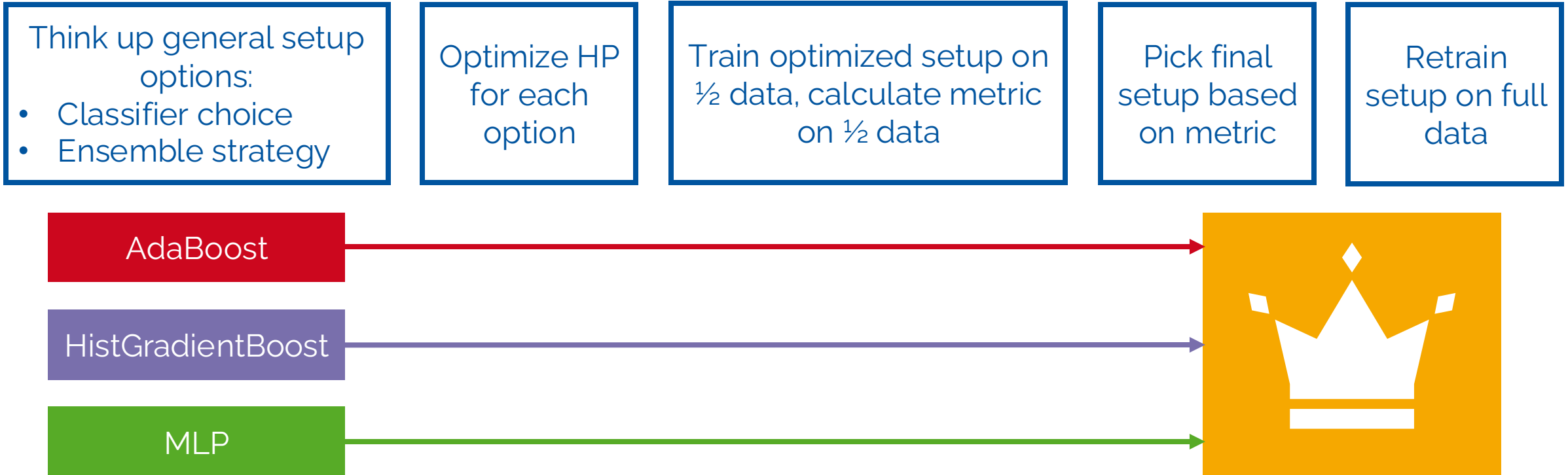
```
params = {}  
params["lr"] = float(np.random.choice([0.01, 0.005, 0.001, 0.0005, 0.0001]))  
params["batch_size"] = int(np.random.choice([64, 128, 256, 512, 1024, 2048, 5096]))  
params["layers"] = [64,64,64]#layers  
params["epochs"] = int(30)  
params["dropout"] = float(np.random.choice([0, 0.1, 0.2, 0.3, 0.4, 0.5]))  
params["weight_decay"] = float(np.random.choice([0,1e-4, 1e-3, 1e-2, 1e-5]))  
params["momentum"] = float(np.random.choice([0.9, 0.99, 0.8, 0.95]))
```





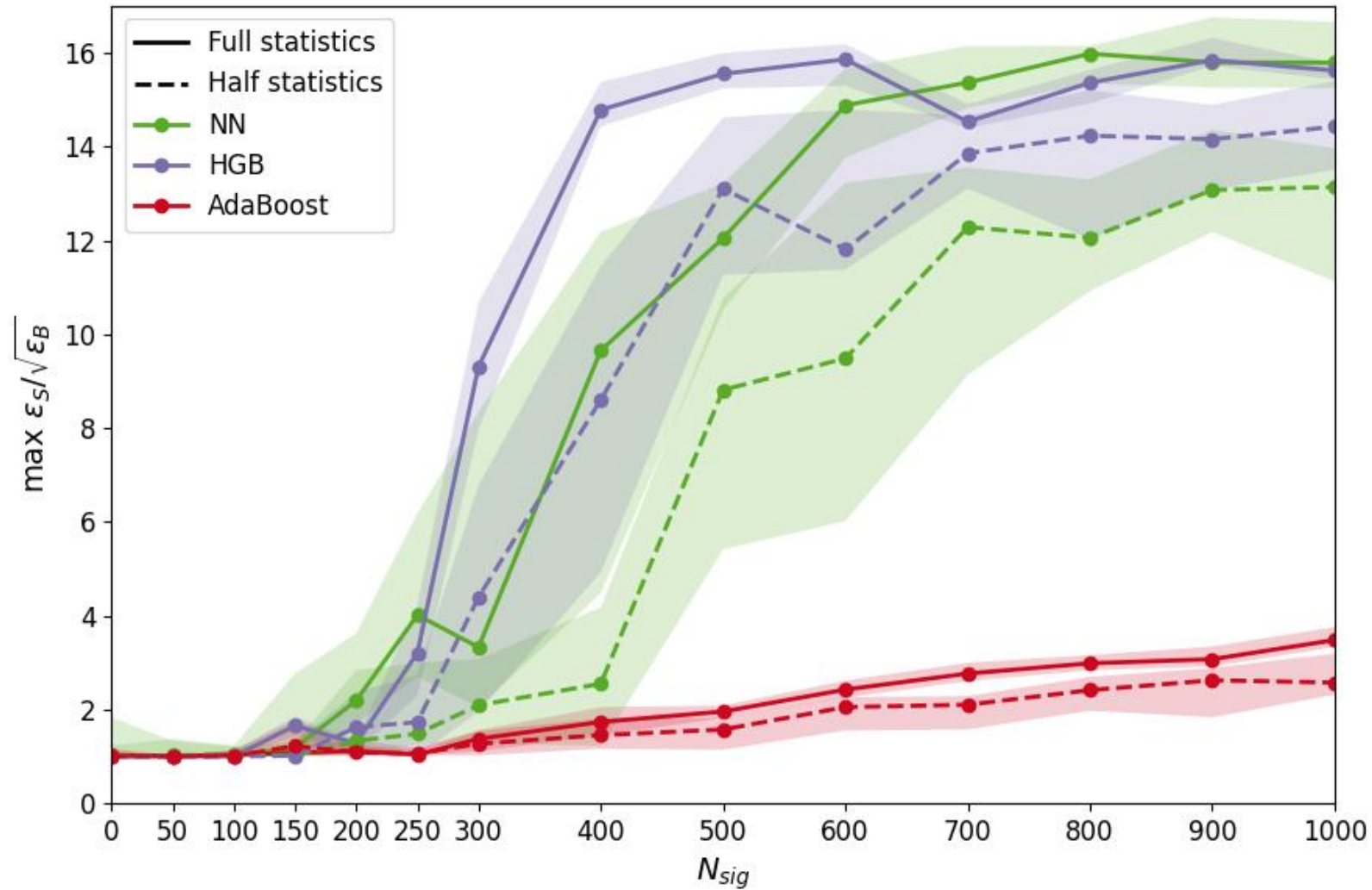
- Different metrics show very similar performance
- While performance at  $N_{sig} = 1000$  is optimal with default hyperparameters, optimizing for lower signal injections can result in slight performance gain
- Trend generally holds up for extended sets but performance gain decreases

# Optimization as a multi-step process



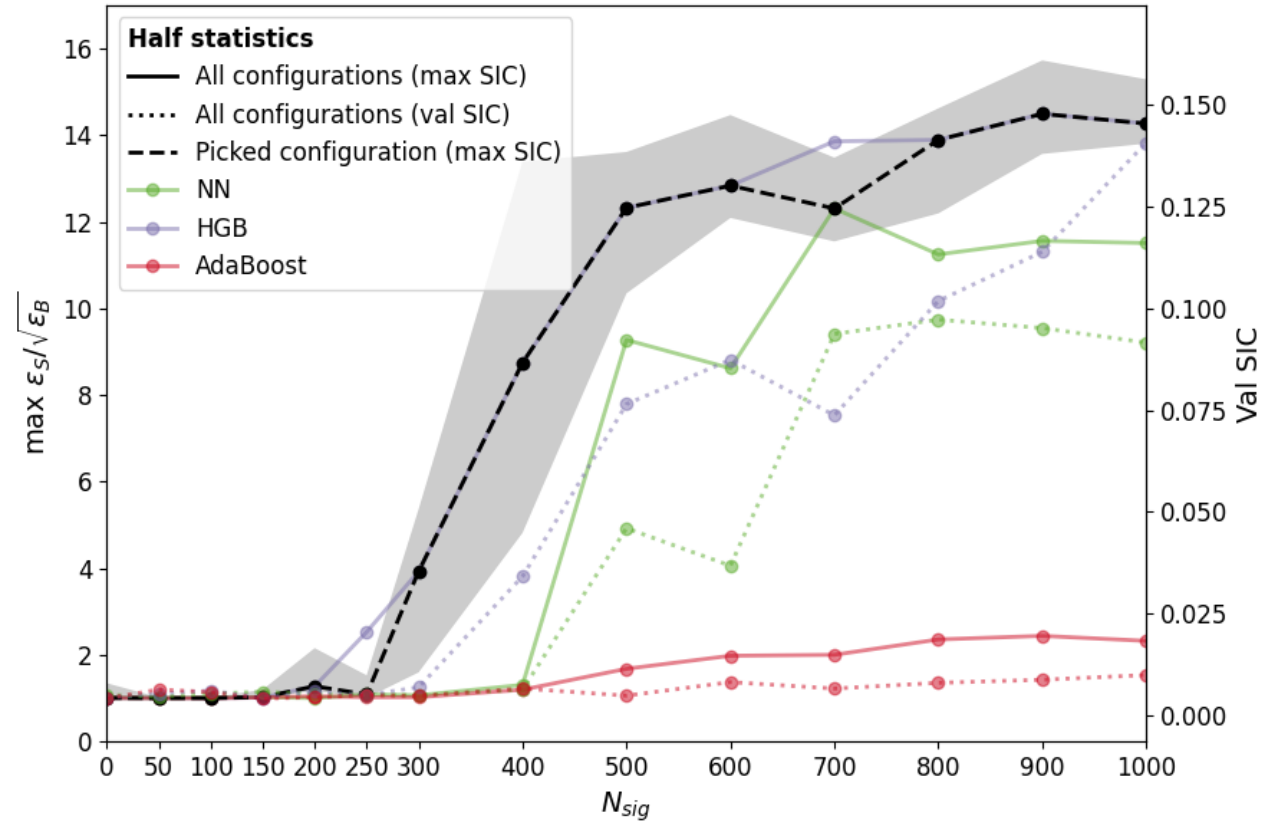


# Performance on half statistics



# Picking option based on val SIC

What we pick



Final performance

