

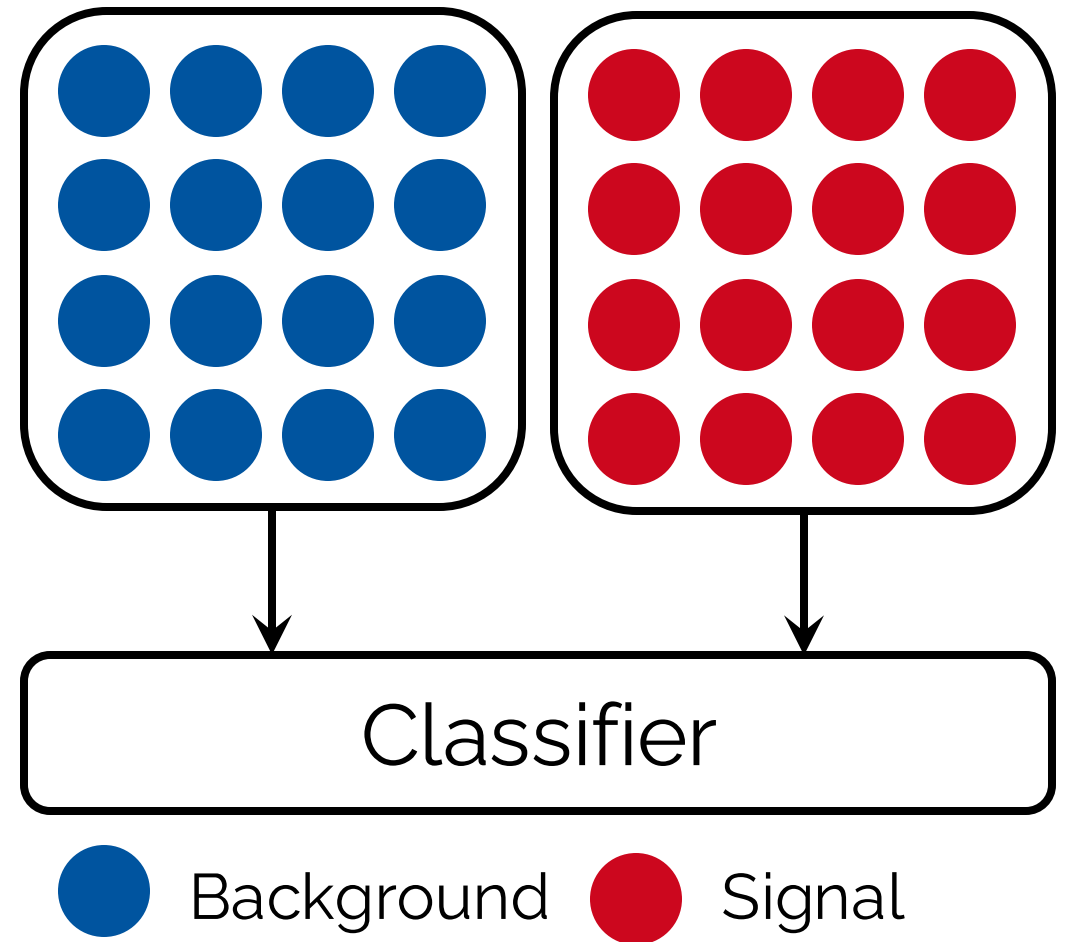
Choosing the right features for weak supervision

Marie Hein

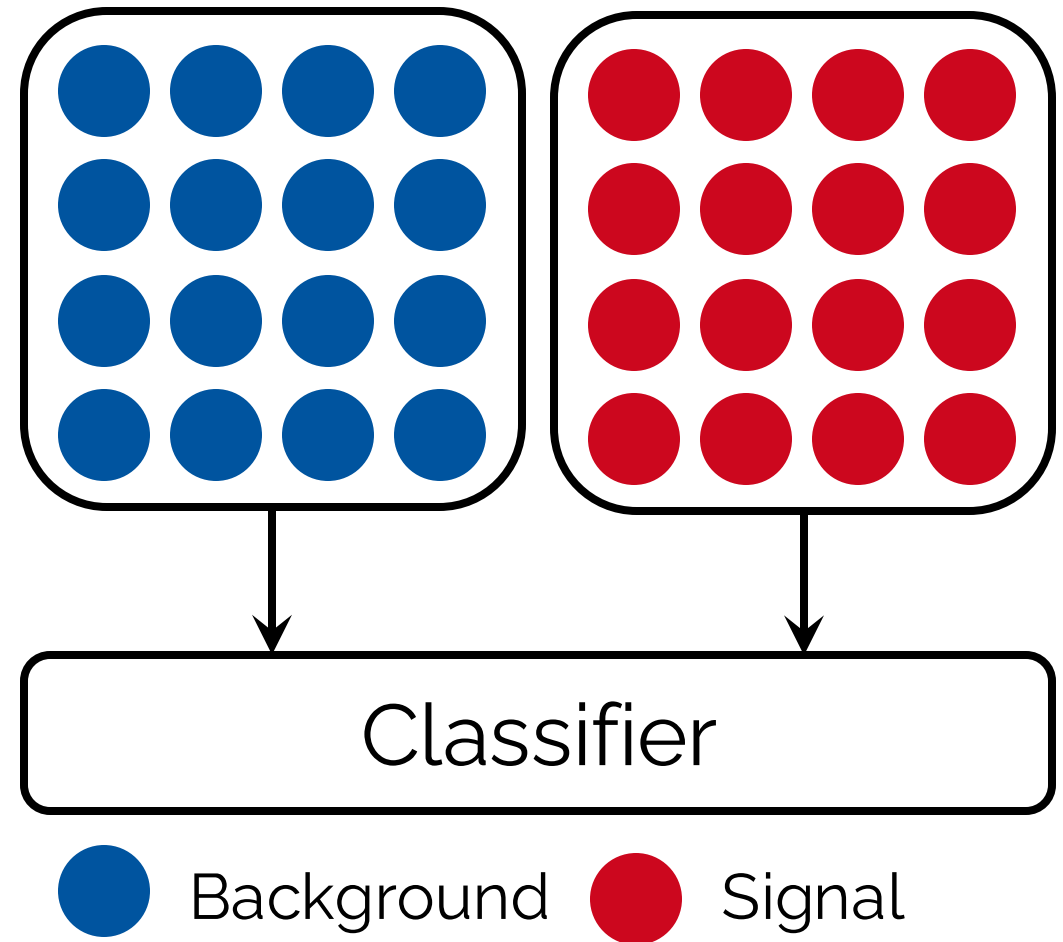
In collaboration with Joep Geuskens, Michael Krämer, Lukas Lang, Radha Mastandrea & Alexander Mück

CRC Young Scientists Meeting 2024

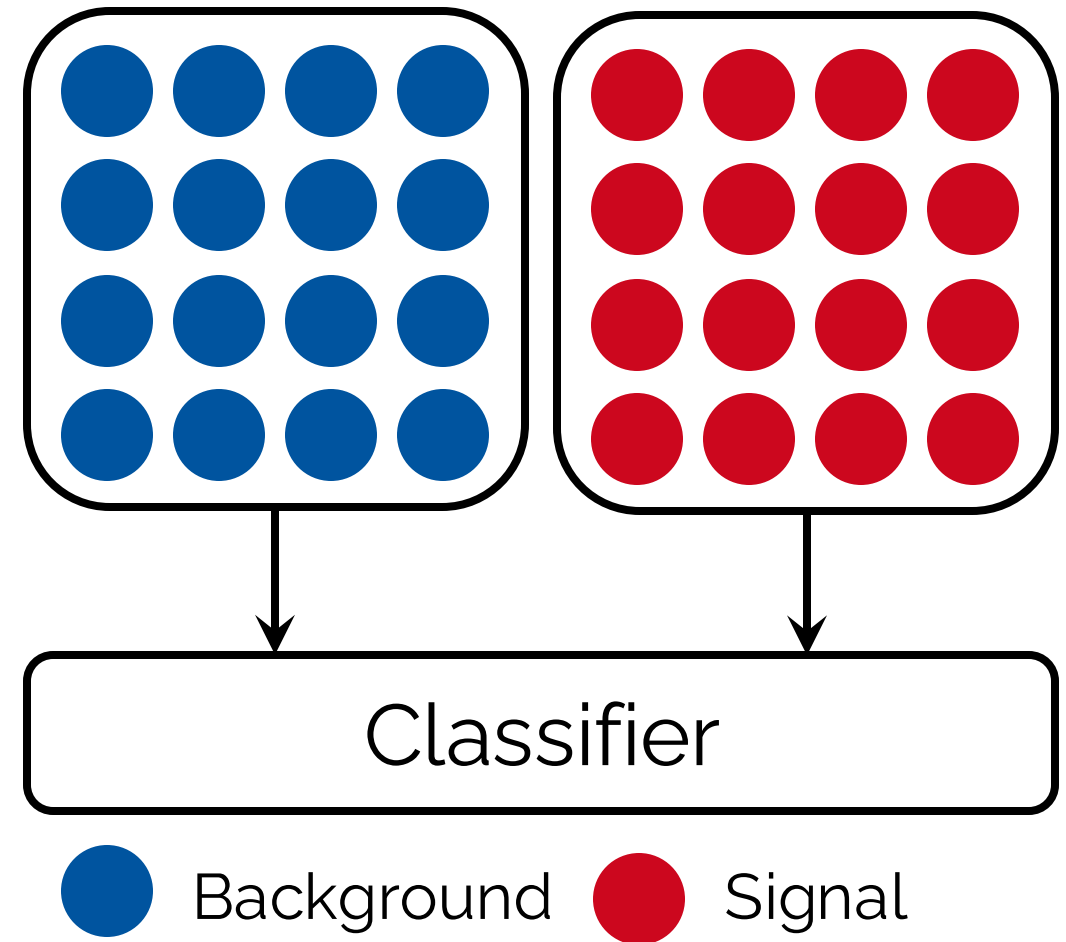
- Goal: To achieve a better signal to background ratio



- Goal: To achieve a better signal to background ratio
- Ansatz: Perform classification task

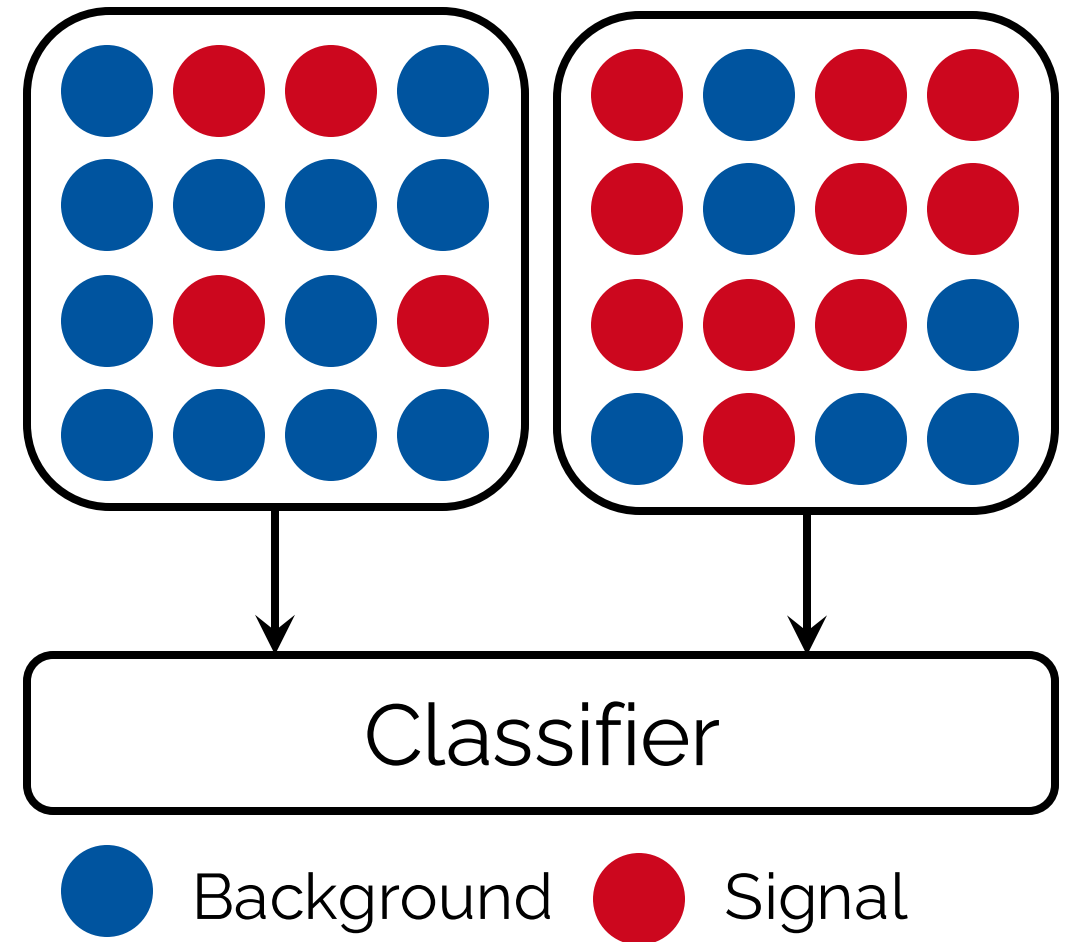


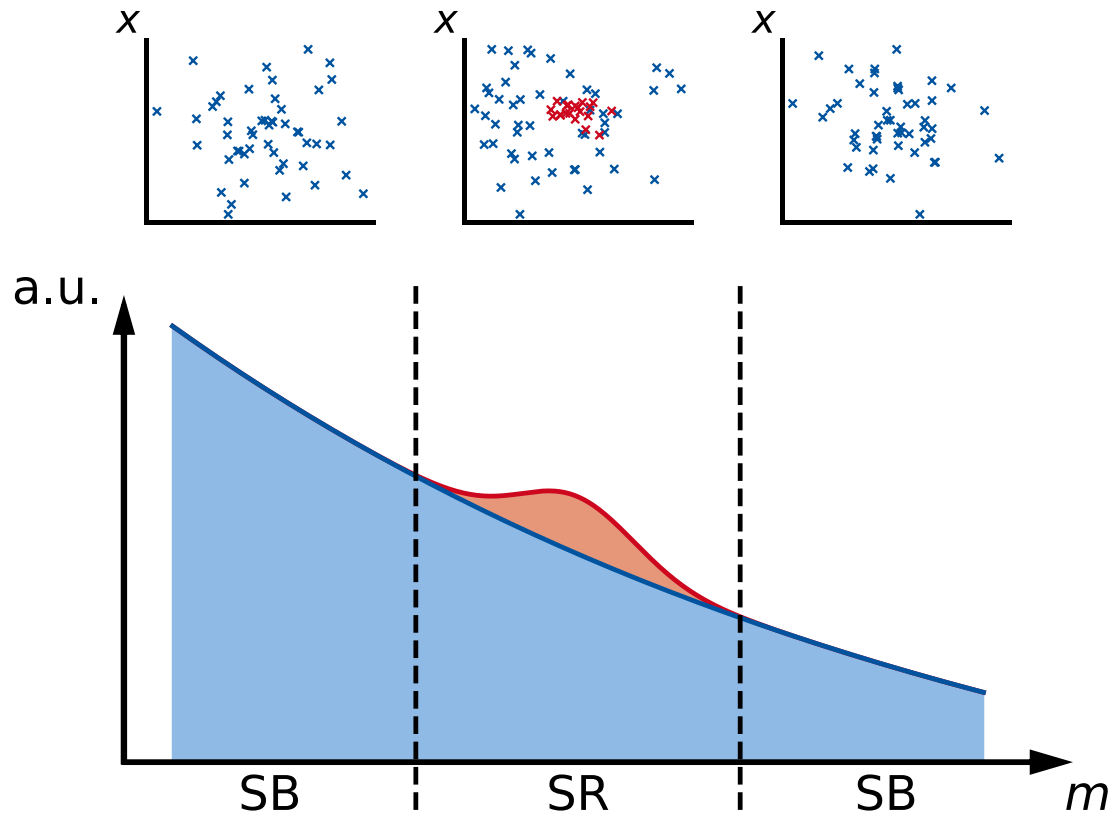
- **Goal:** To achieve a better signal to background ratio
- **Ansatz:** Perform classification task
- **Problem:** Labels are not available on real data



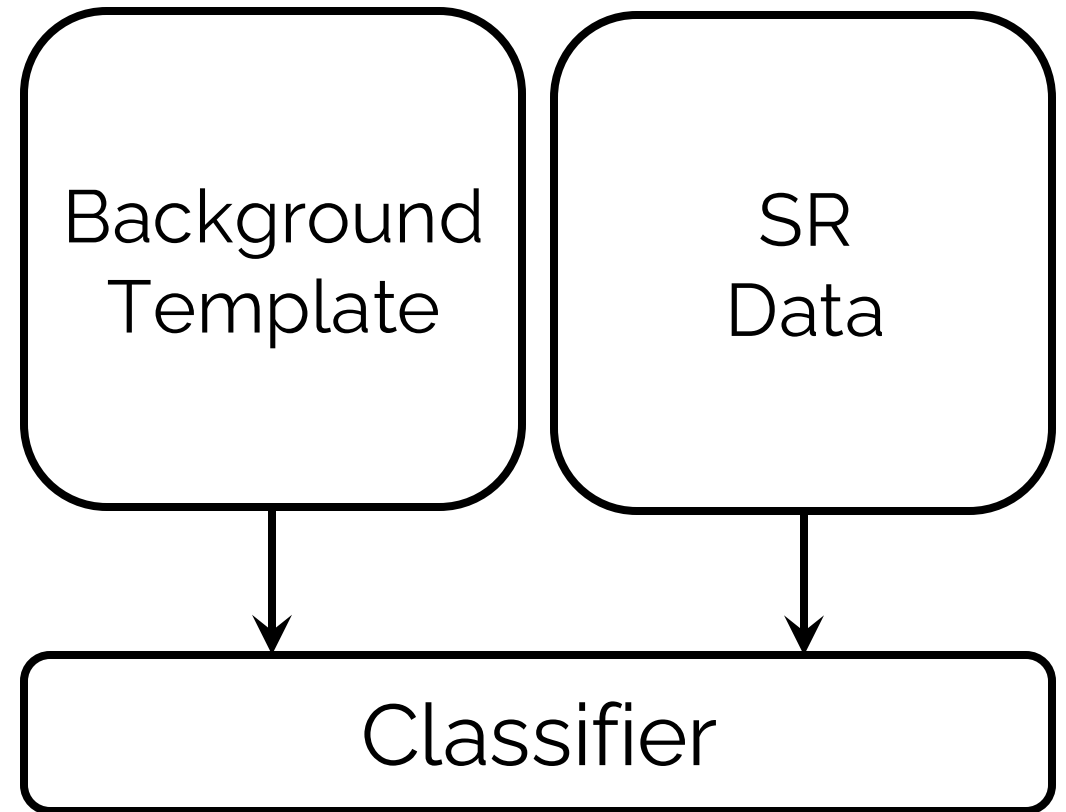
“Classification without labels: Learning from mixed samples in high energy physics” [1709.02949], E. Metodiev, B. Nachman, J. Thaler

- **Goal:** To achieve a better signal to background ratio
- **Ansatz:** Perform classification task
- **Problem:** Labels are not available on real data
- **Solution:** Classify between mixed classes
 - Fundamentally, both problems are equivalent

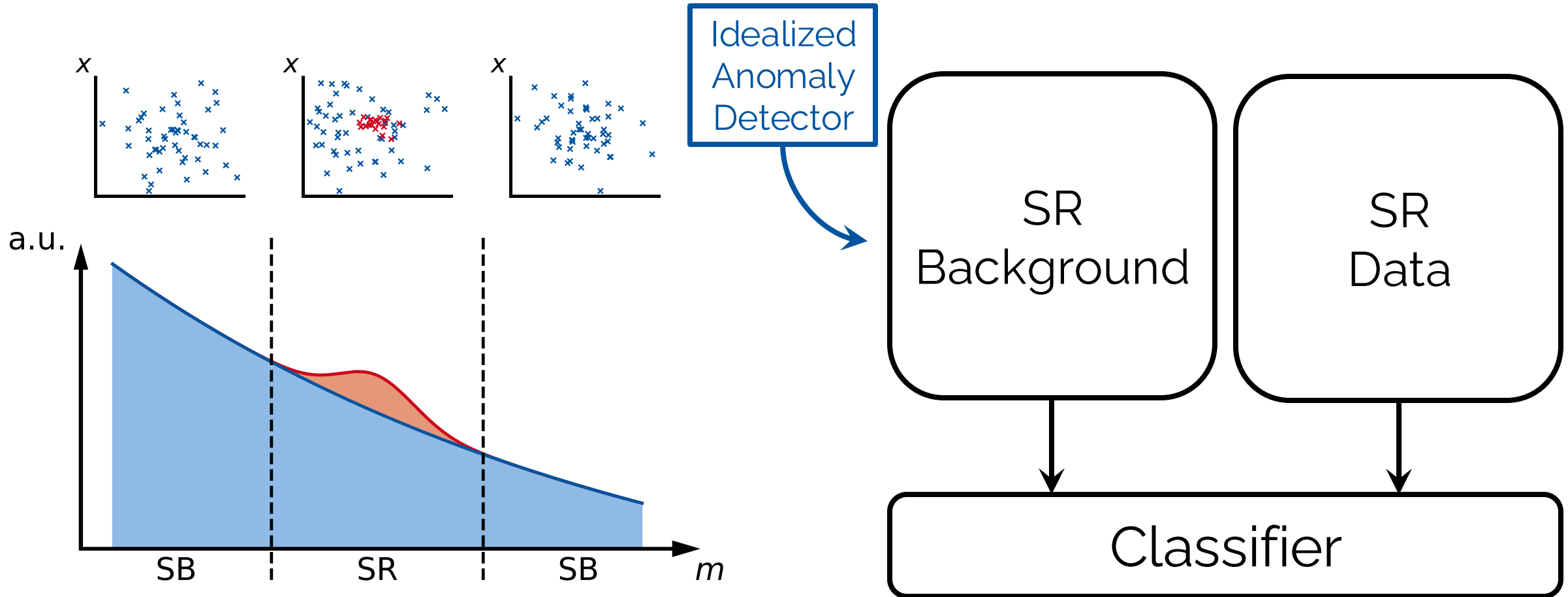




Recreated from [\[2109.00546\]](#)

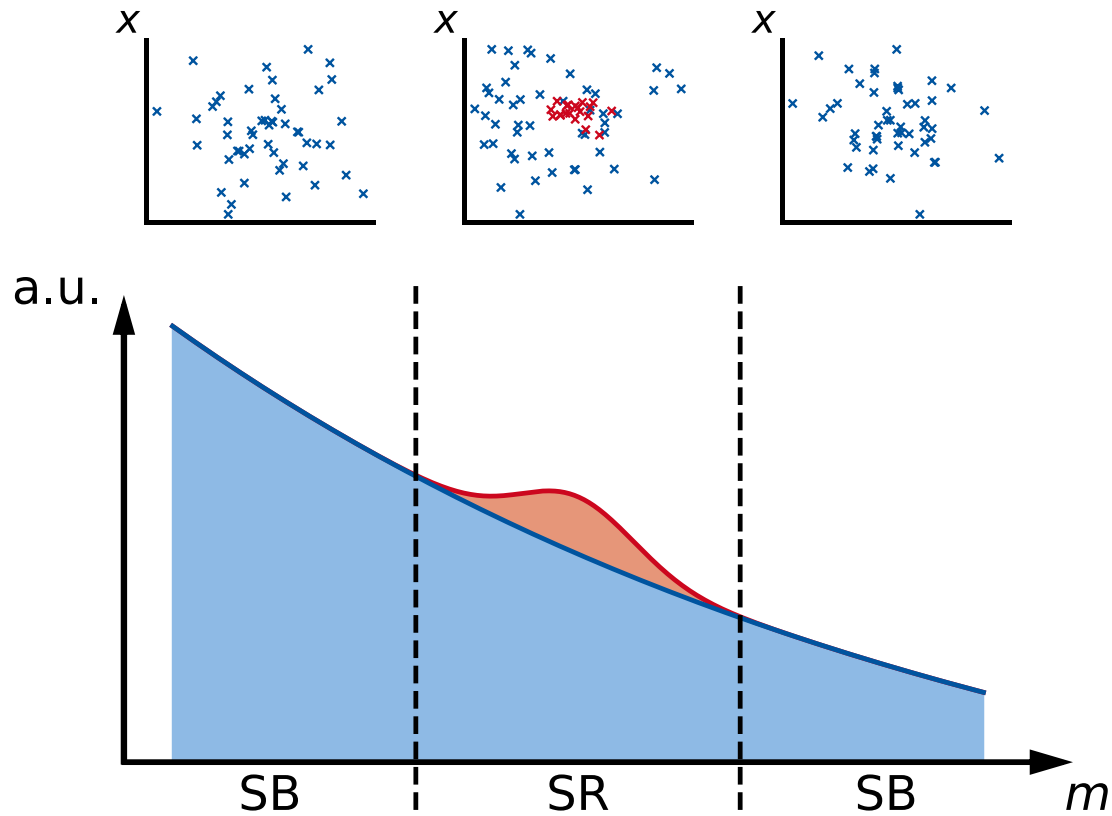


Application to resonance searches

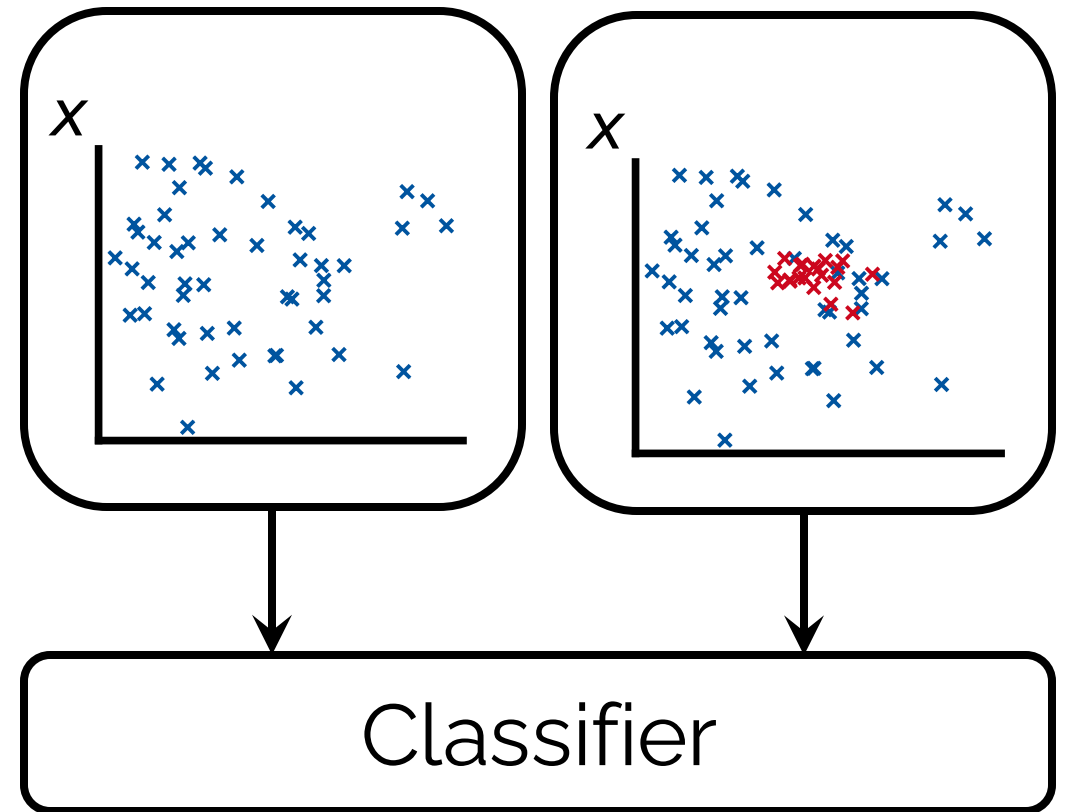


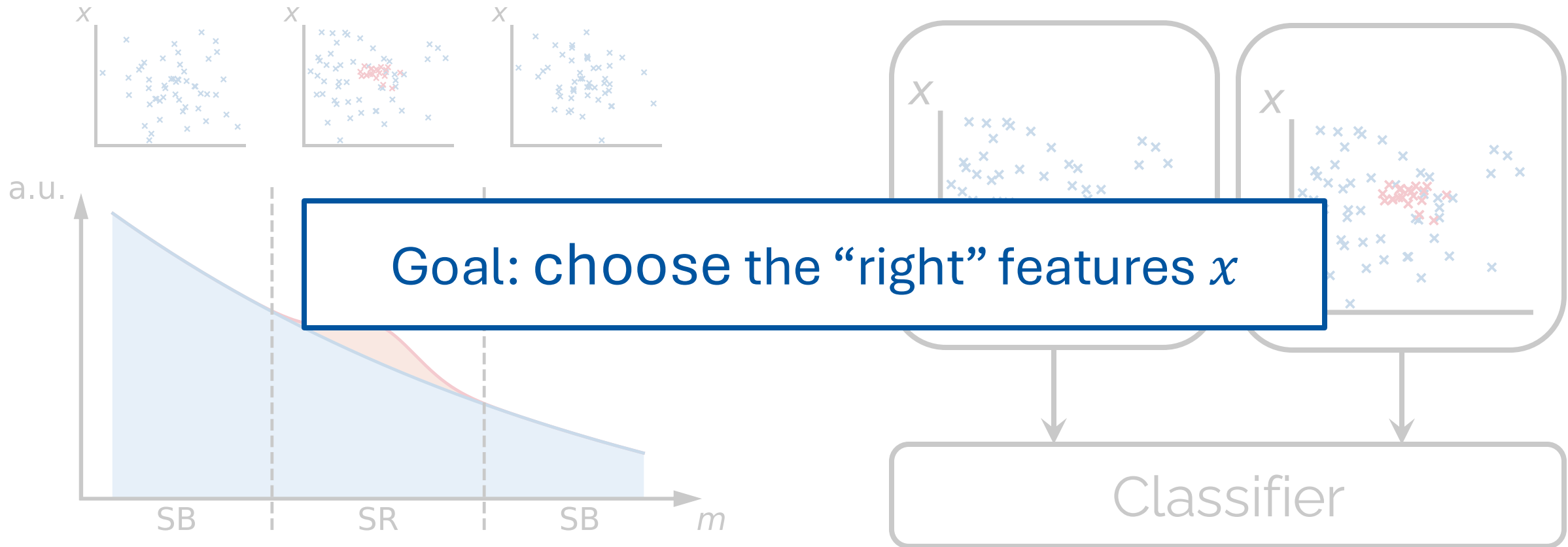
Recreated from [2109.00546]

Application to resonance searches



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



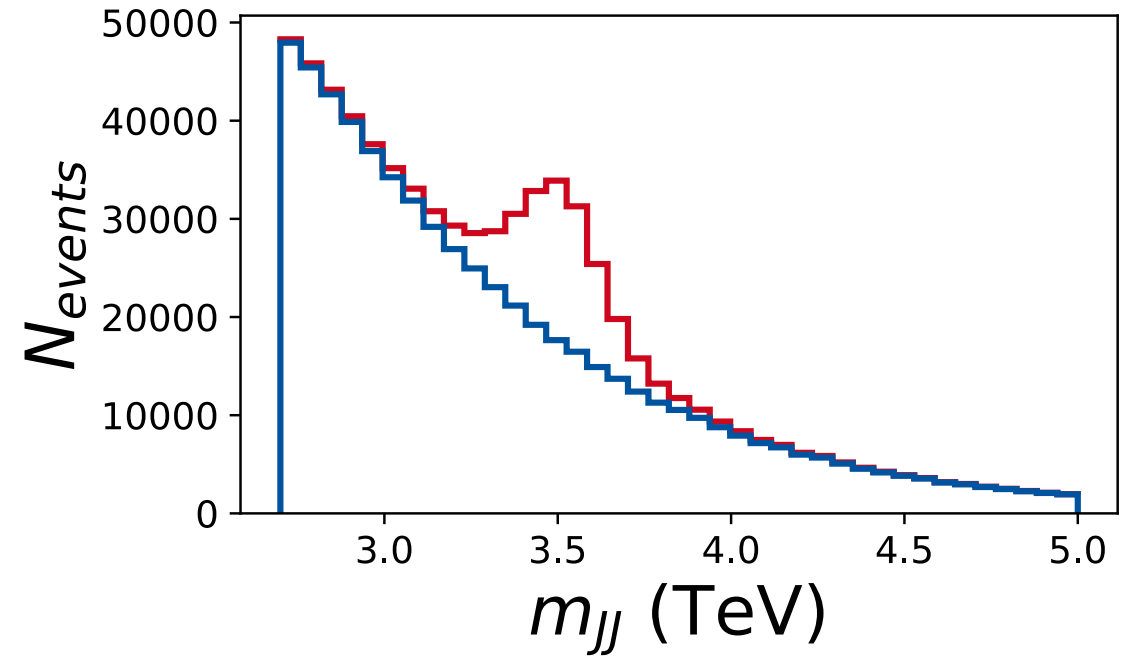
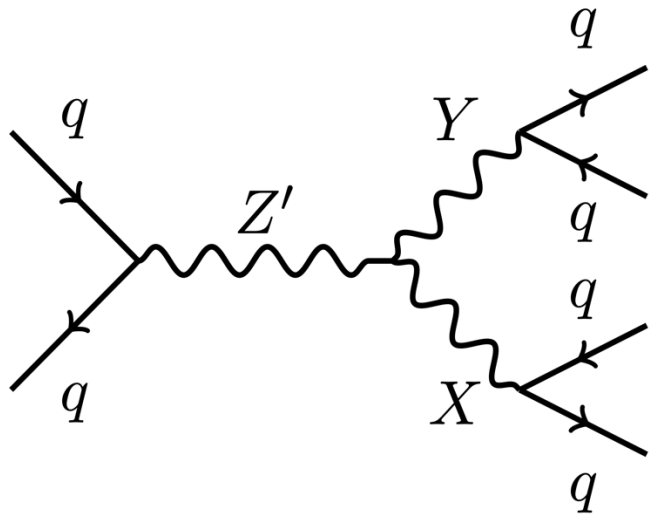


Recreated from [\[2109.00546\]](#)

The Dataset & Features

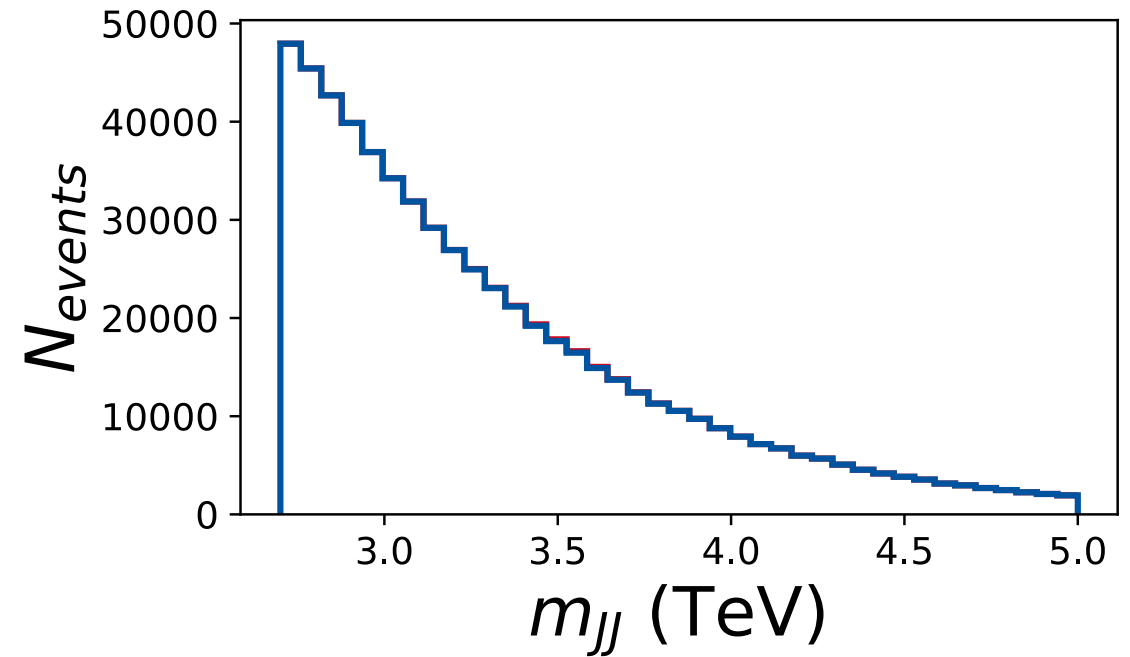
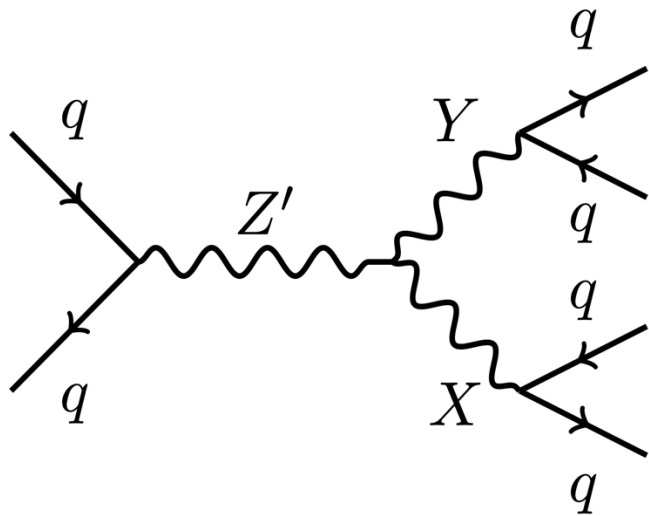
“The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics” [2101.08320], G. Kasieczka, B. Nachman, D. Shih et. al.

- Benchmark dataset for anomaly detection
- QCD dijet background
- Signal



“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
- Signal



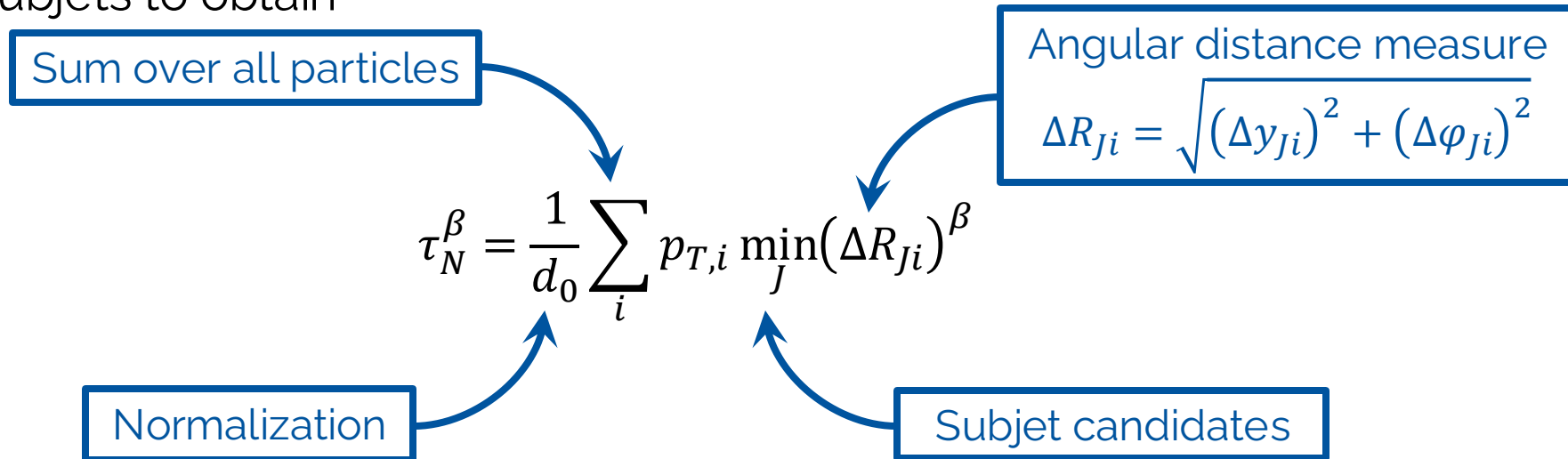
1. **Low level features:** use particle four-momenta in jet-separated LorentzNet
 - Very model agnostic

2. **High level Features:** derive observables from low-level features
 - Less model agnostic
 - Easier classification task (more closely related to problem to be solved)
 - a. **N-Subjettiness**
 - b. **Energy Flow Polynomials**

“Identifying Boosted Objects with N-subjettiness” [1011.2268], J. Thaler, K. Van Tilburg

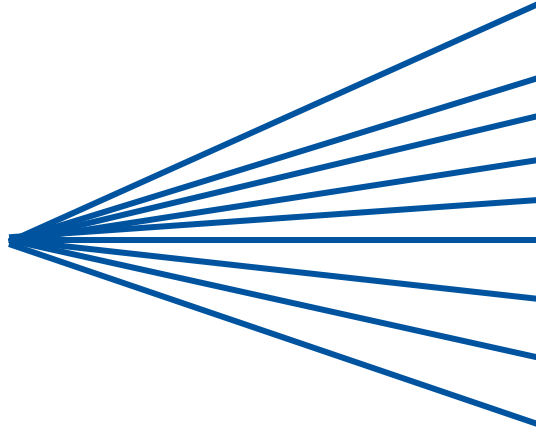
“Maximizing Boosted Top Identification by Minimizing N-subjettiness” [1108.2701], J. Thaler, K. Van Tilburg

- Cluster into N subjects to obtain



- “Momentum-weighted sum of angular distance of all particles to closest subjet”

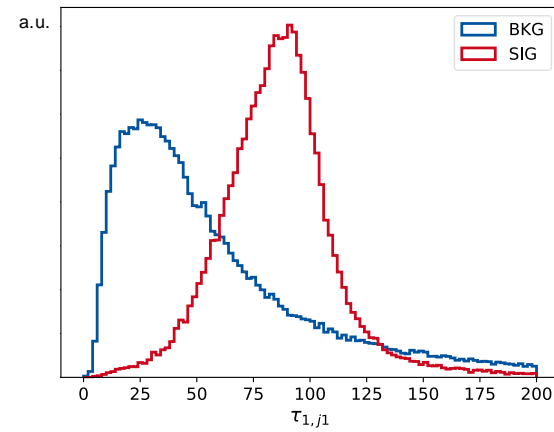
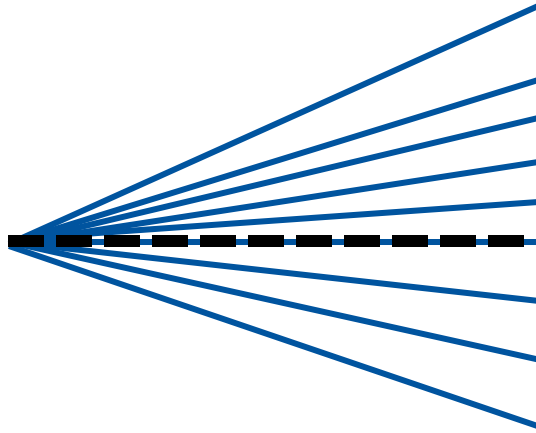
Background



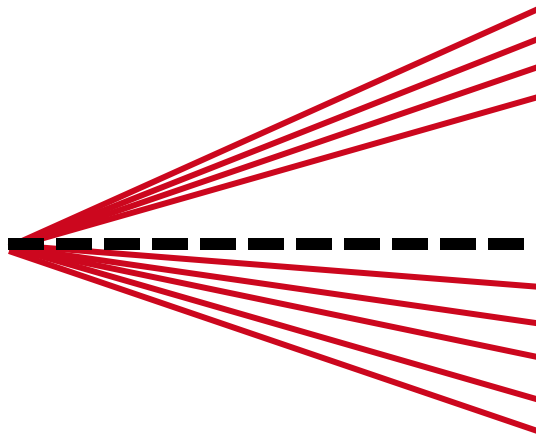
Signal



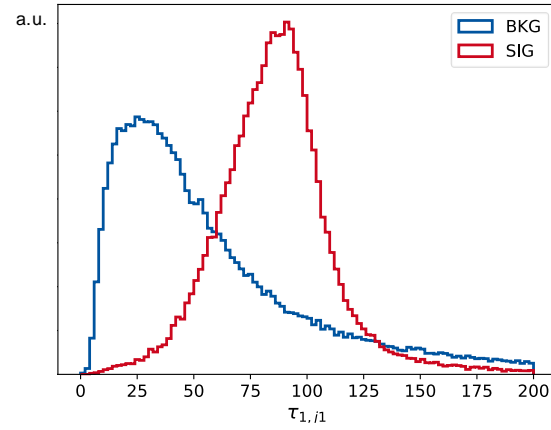
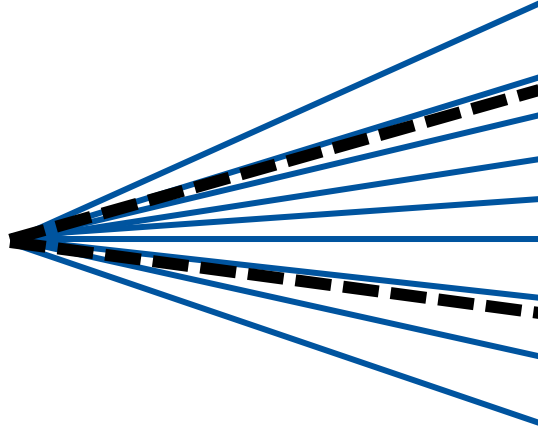
Background



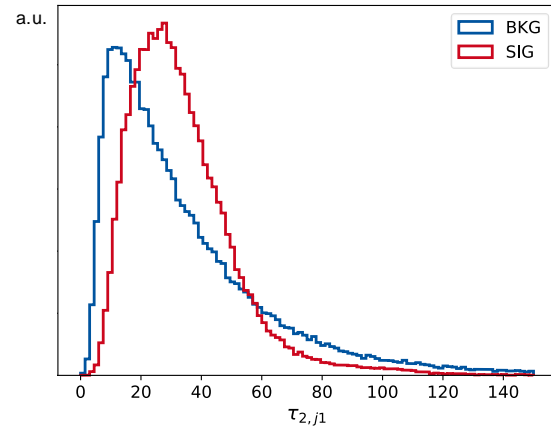
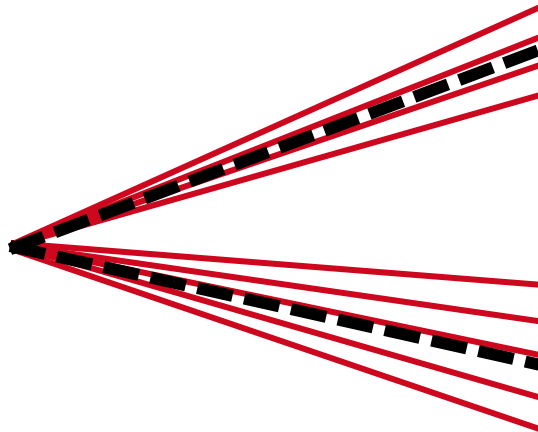
Signal



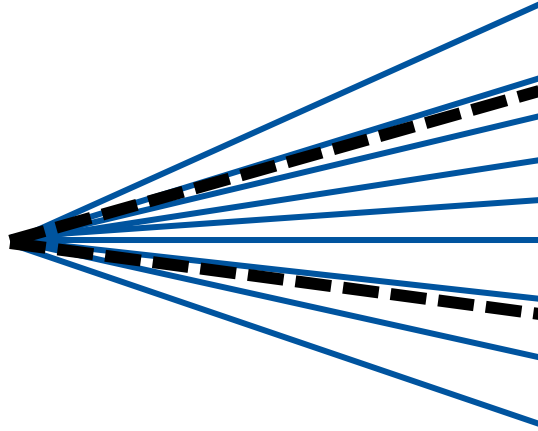
Background



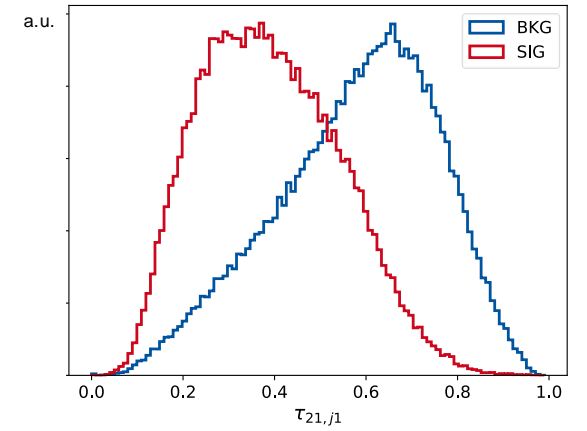
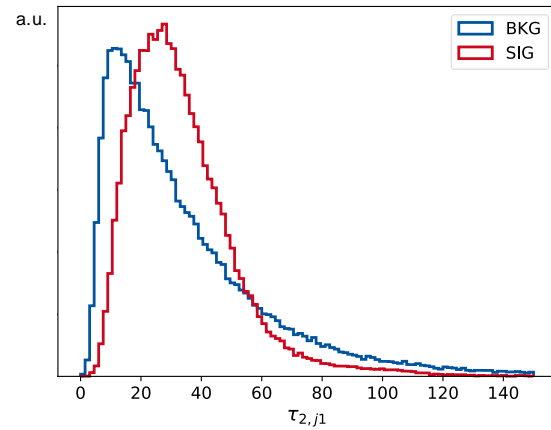
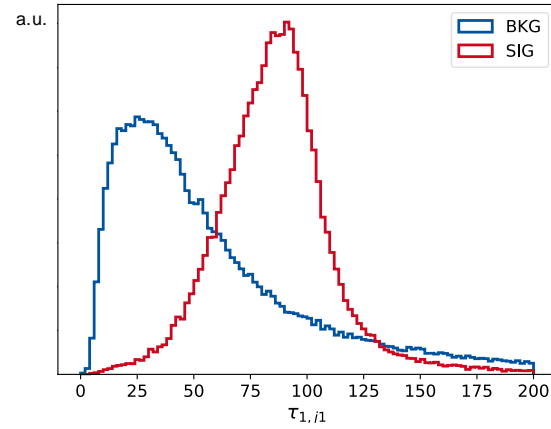
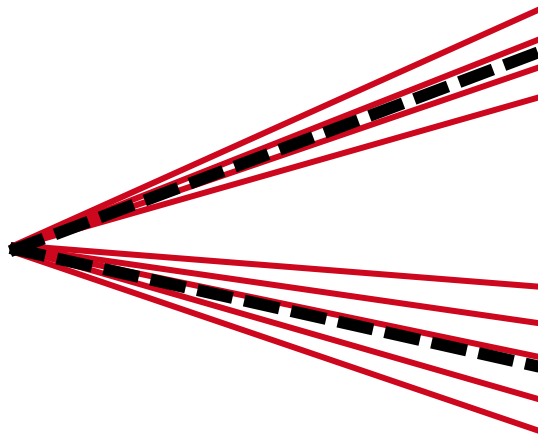
Signal



Background



Signal



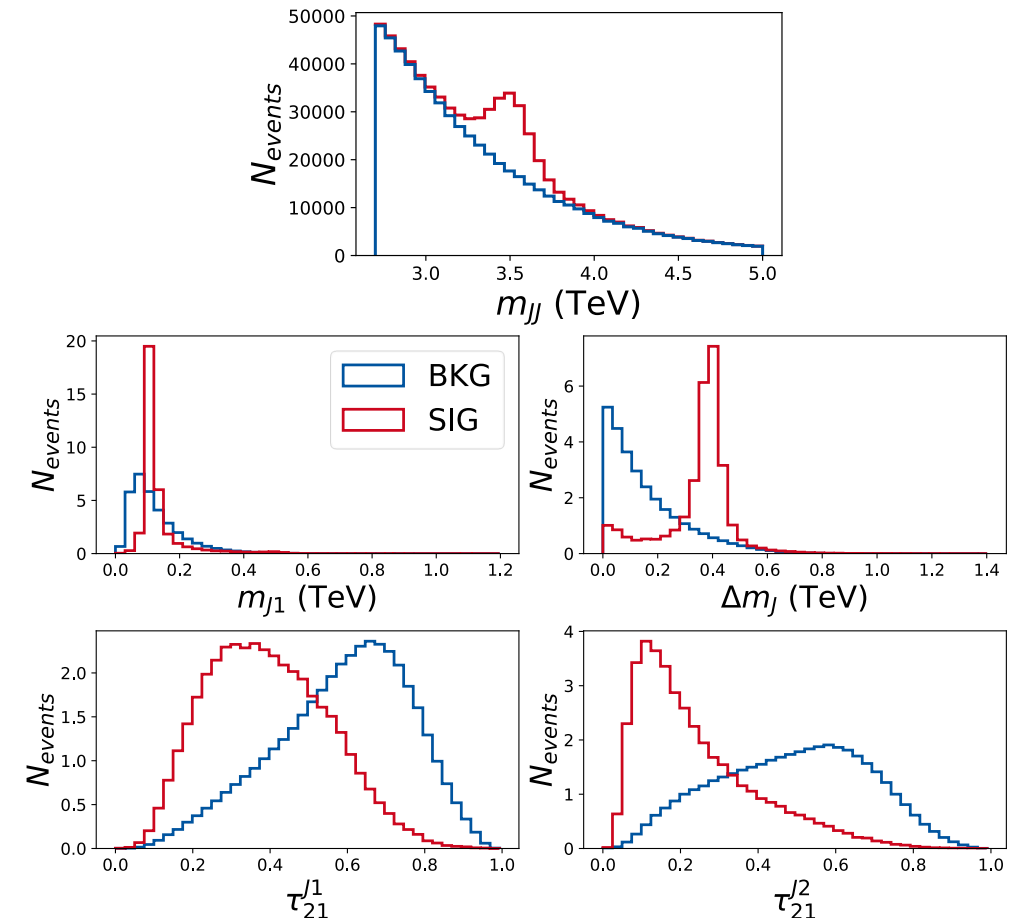
1. Baseline feature set

- Jet masses $m_{J1}, \Delta m_J$
- 21-Subjettiness ratio $\tau_{21,J1}, \tau_{21,J2}$

2. Extended feature set

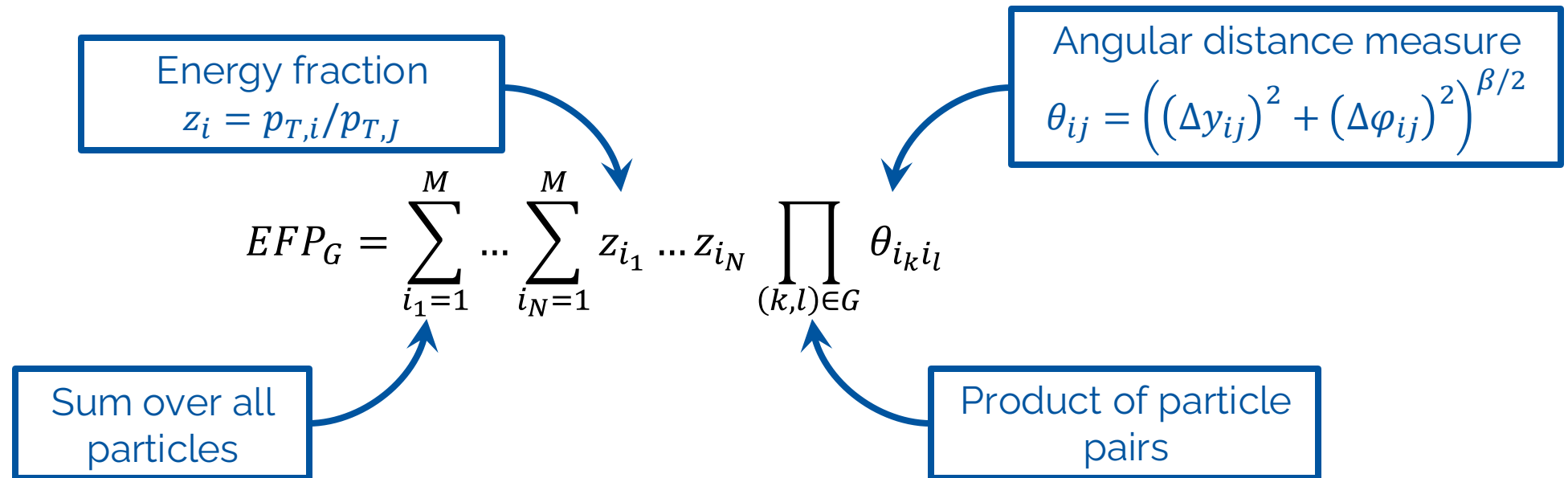
- Jet masses $m_{J1}, \Delta m_J$
- Use 54 different subjettiness features (varying N and β)

Baseline Features

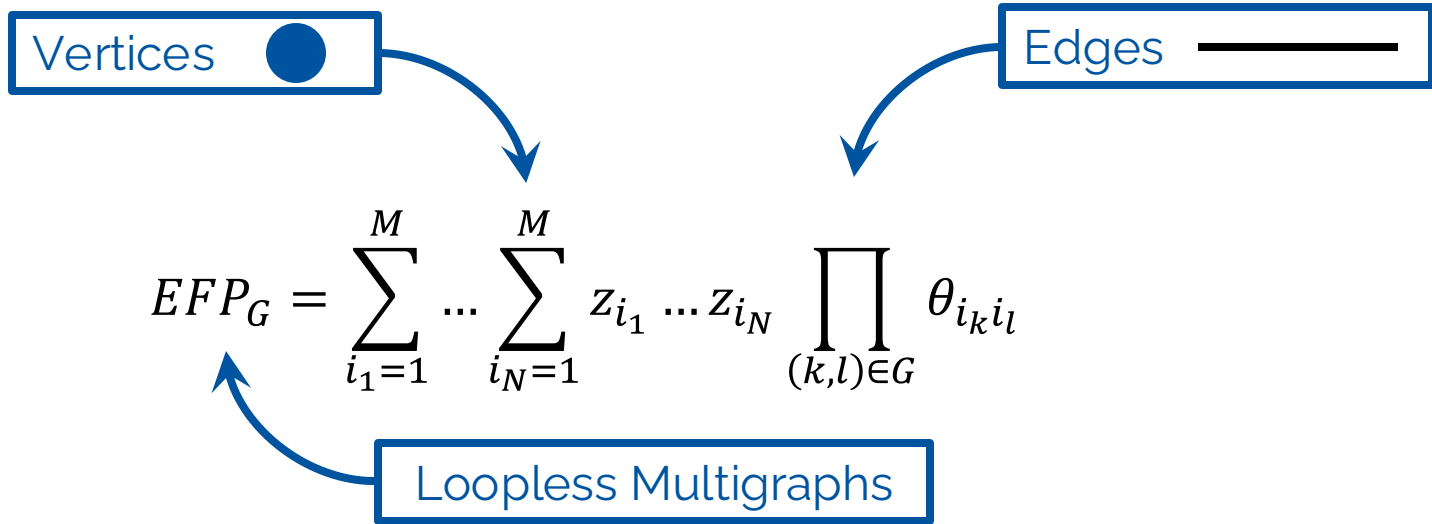


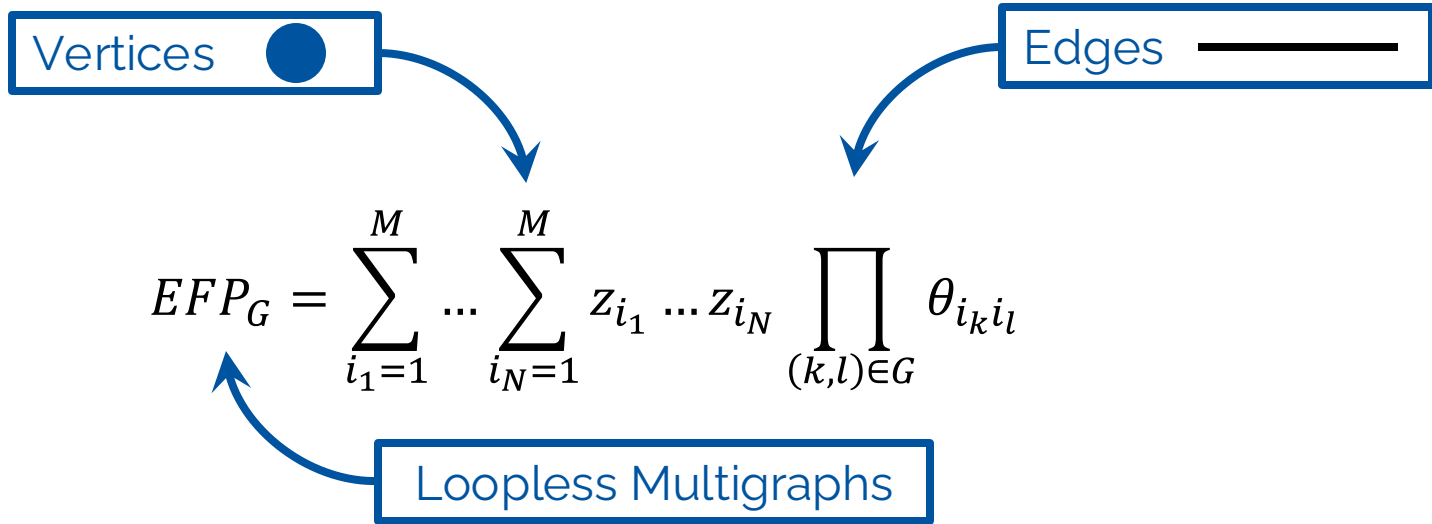
“Energy flow polynomials: A complete linear basis for jet substructure” [1712.07124], P. Komiske, E. Metodiev, J. Thaler
“Energy Flow Networks: Deep Sets for Particle Jets” [1810.05165], P. Komiske, E. Metodiev, J. Thaler


- Complete linear basis of jet substructure observables



EFP-Multigraph Correspondence

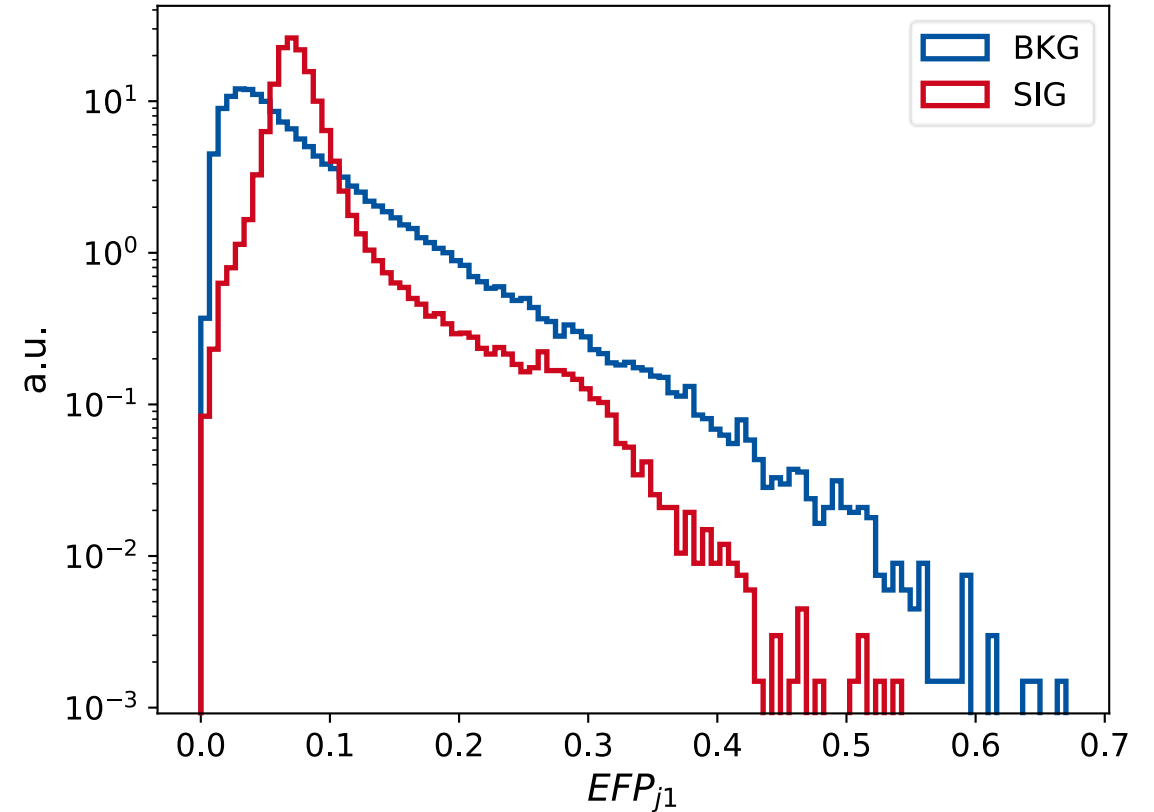
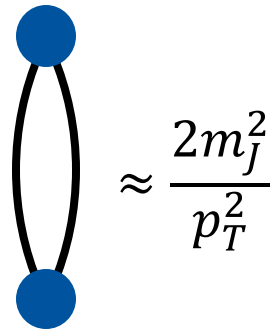






$$= \sum_{i=1}^M \sum_{j=1}^M z_i z_j \theta_{ij}^2 = \sum_{i=1}^M \sum_{j=1}^M \frac{p_{T,i} p_{T,j}}{p_{T,J}^2} \left((\Delta y_{ij})^2 + (\Delta \varphi_{ij})^2 \right) \approx \frac{2m_J^2}{p_T^2}$$

$$EFP_G = \sum_{i_1=1}^M \dots \sum_{i_N=1}^M z_{i_1} \dots z_{i_N} \prod_{(k,l) \in G} \theta_{i_k i_l}$$

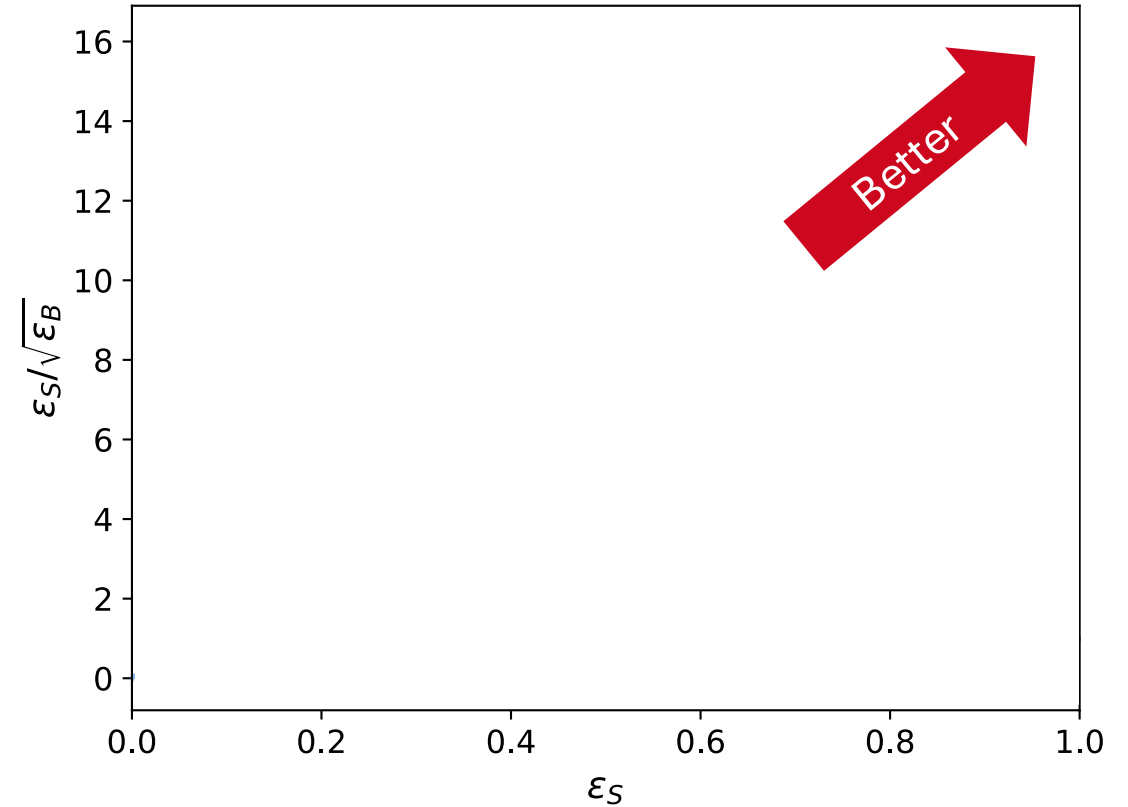
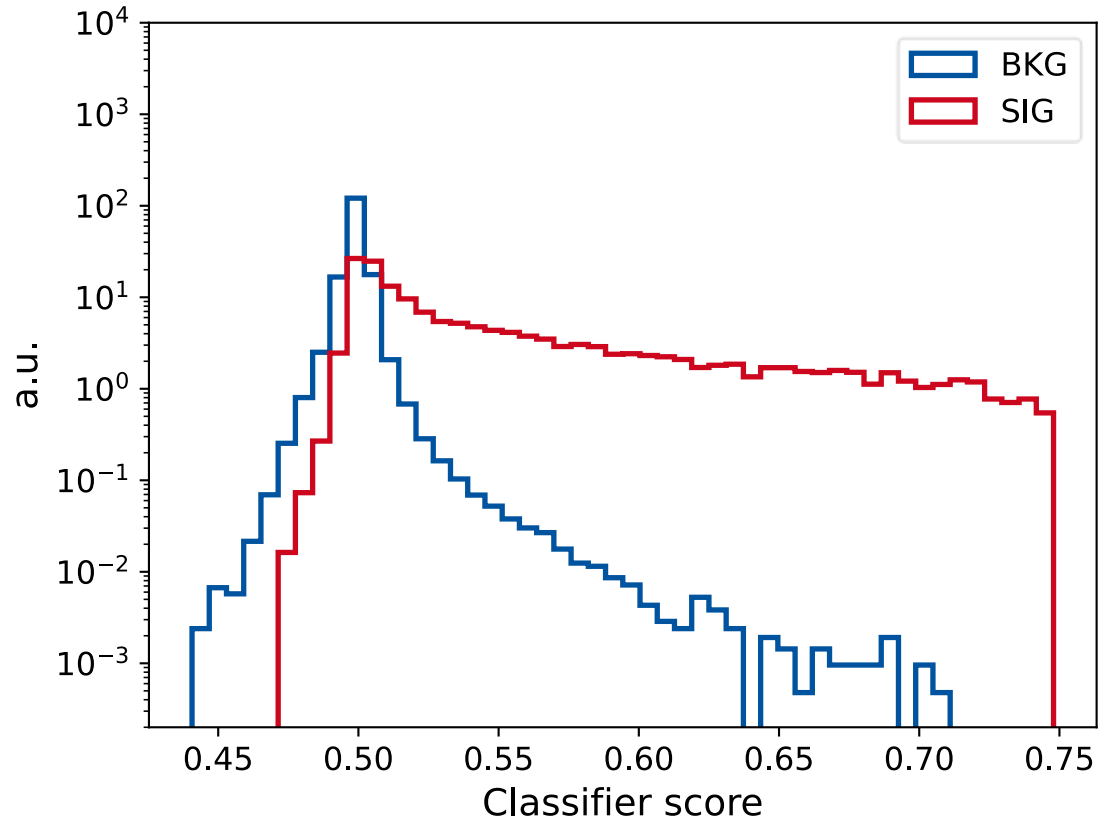


1. EFP feature set
 - Jet masses $m_{J1}, \Delta m_J$
 - 490 EFPs per jet (up to 7 edges)

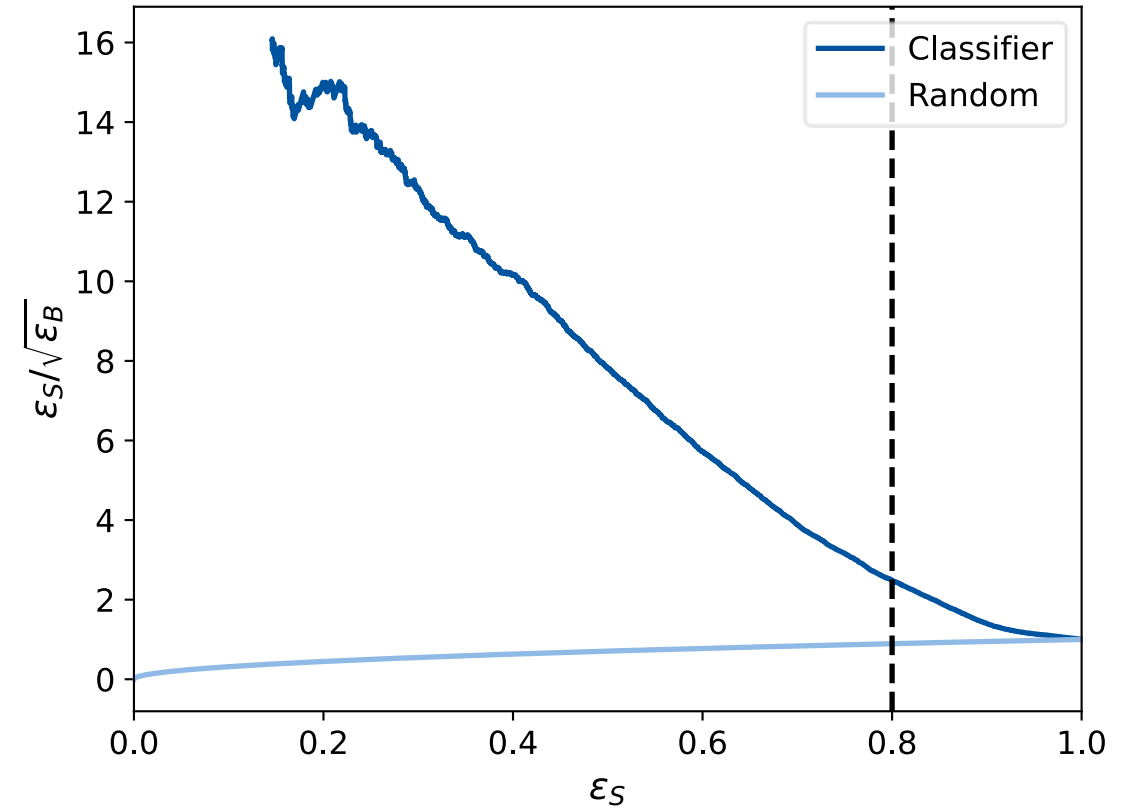
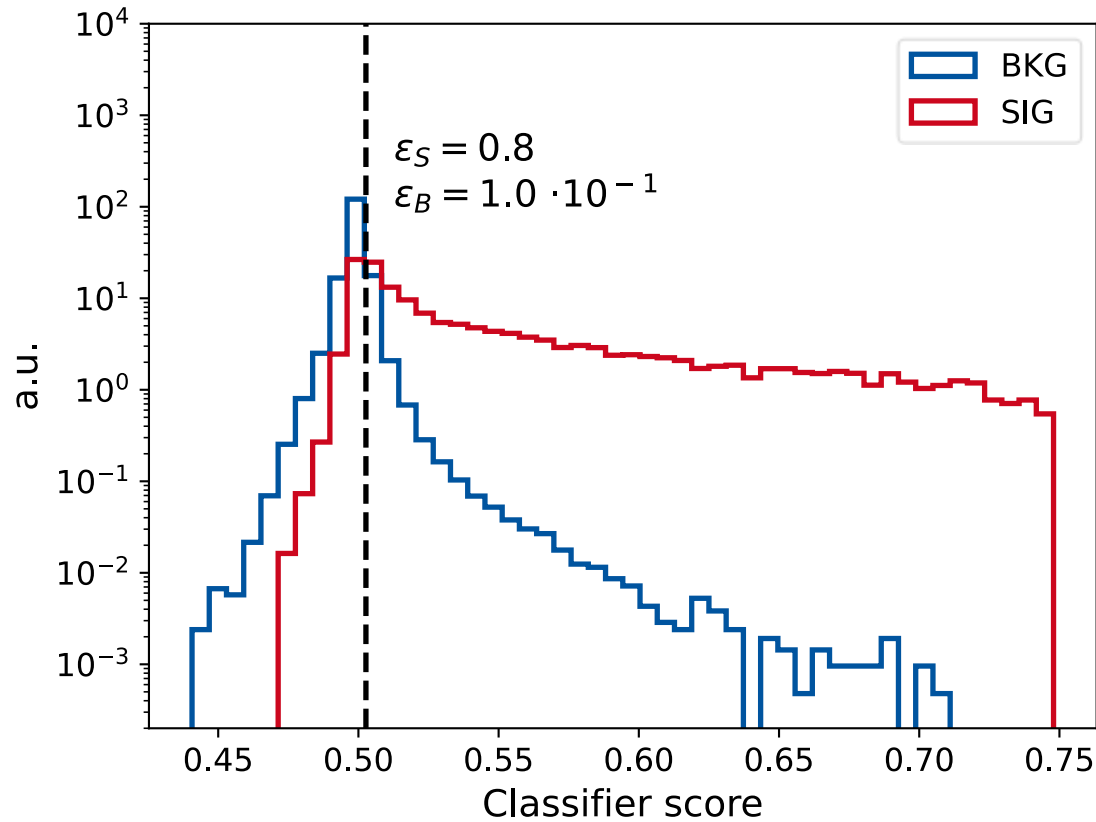
1. Low level features
2. High level features
 - a. Baseline feature set (Subjettiness) → 4 features
 - b. Extended feature set (Subjettiness) → 56 features
 - c. EFP feature set → 982 features

Performance Measure

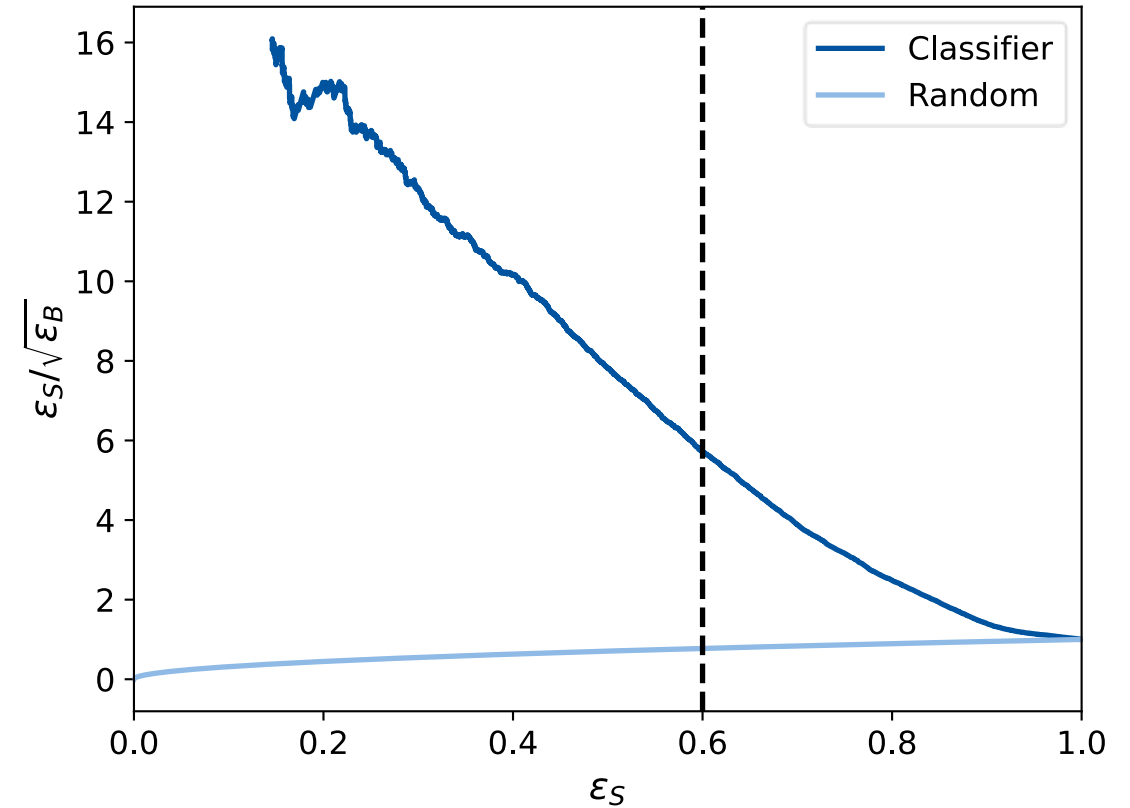
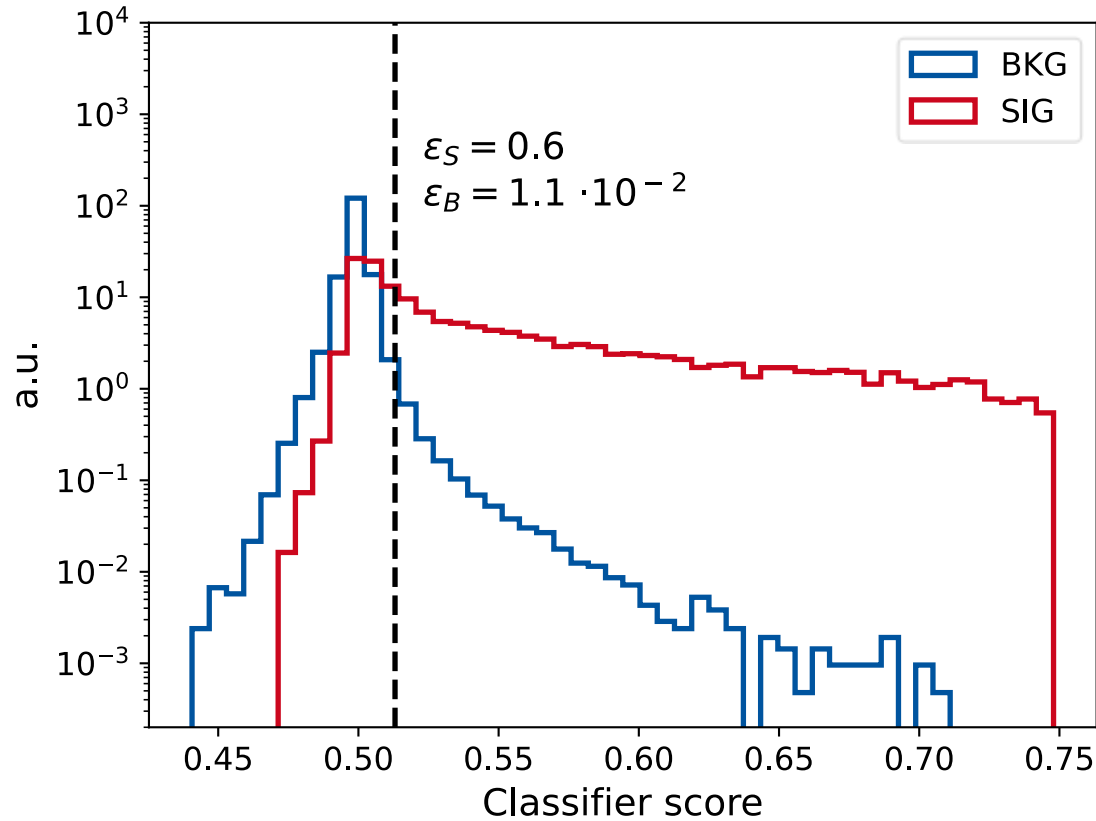
Classifier score to anomaly detection



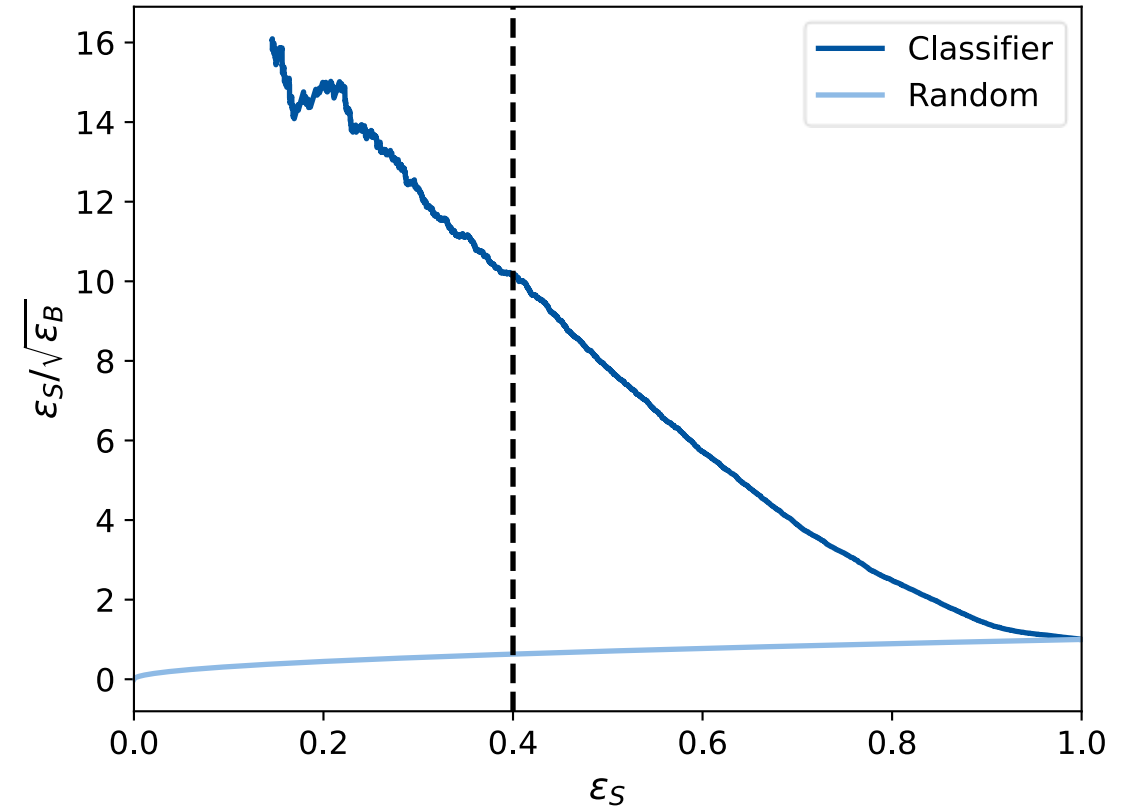
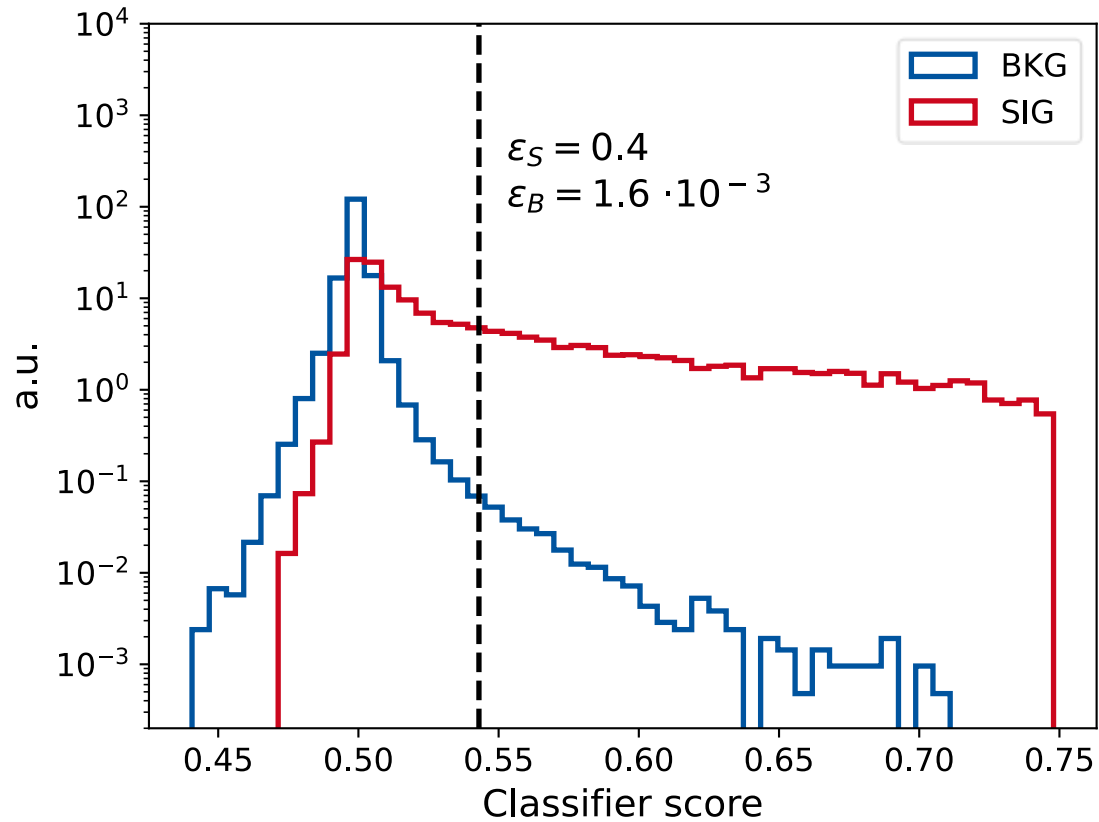
Classifier score to anomaly detection



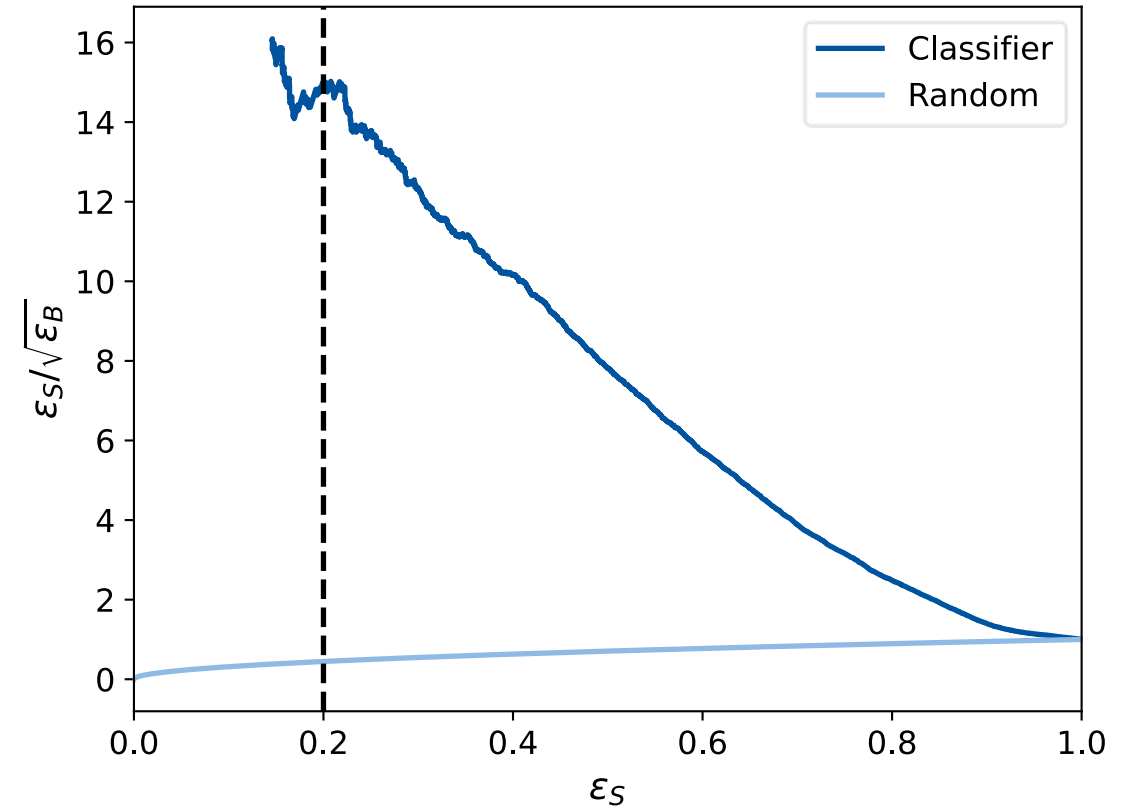
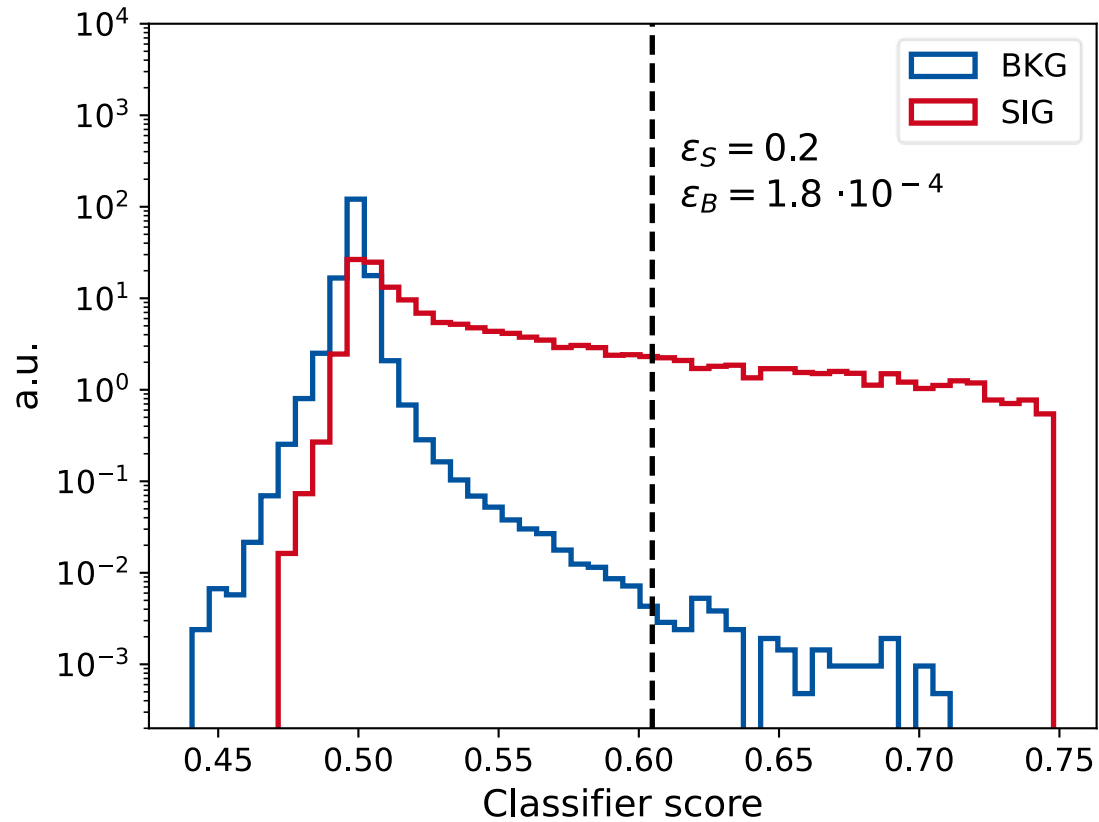
Classifier score to anomaly detection

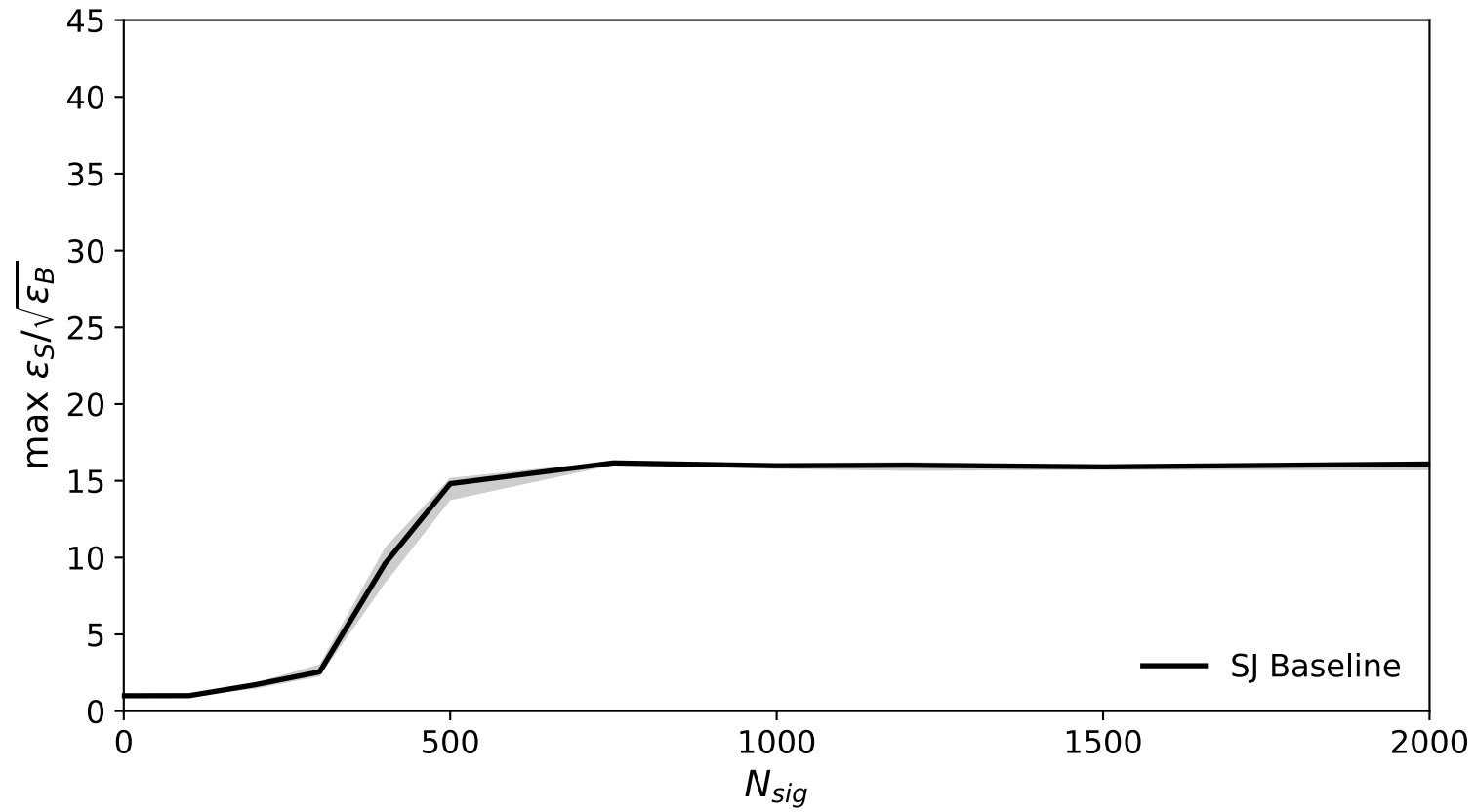


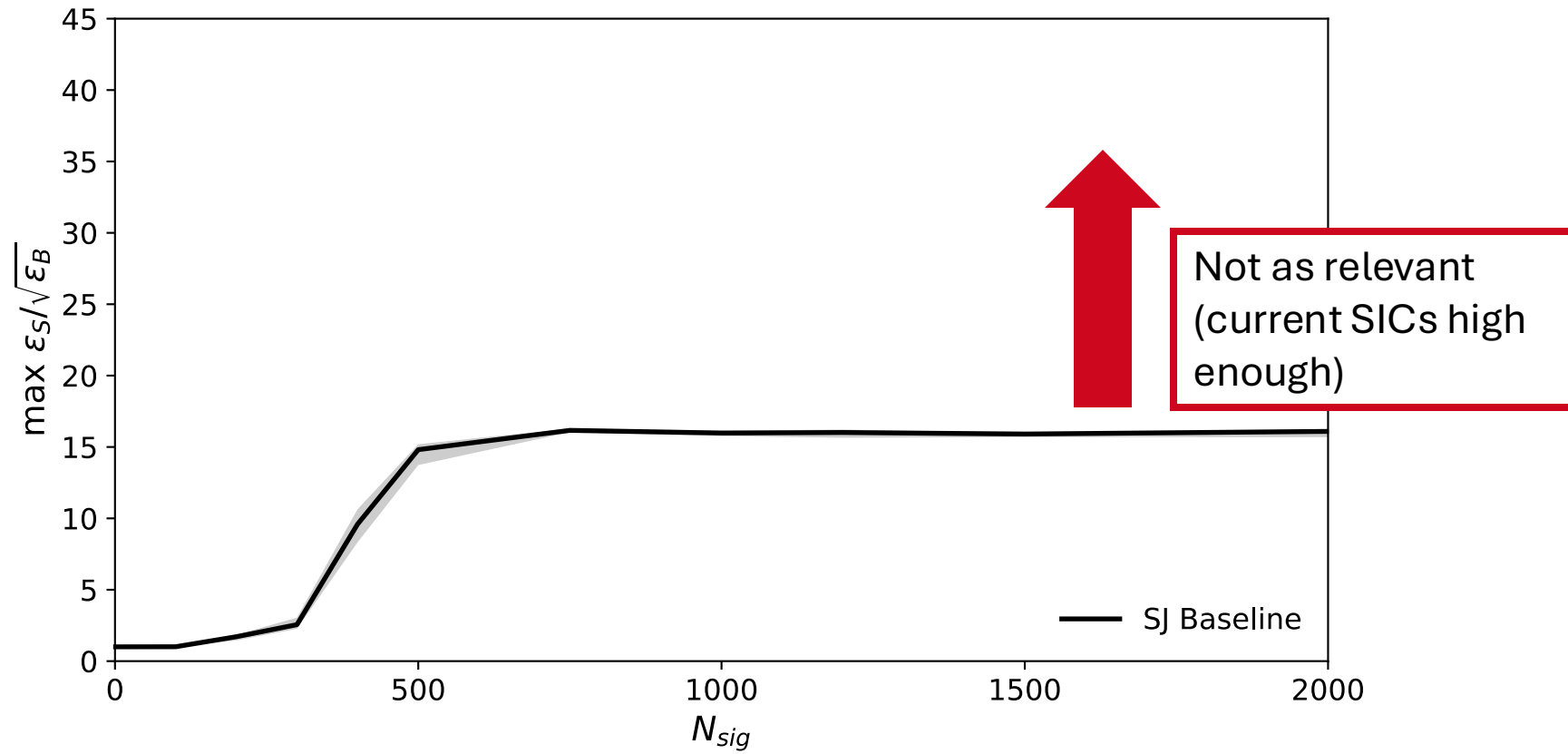
Classifier score to anomaly detection

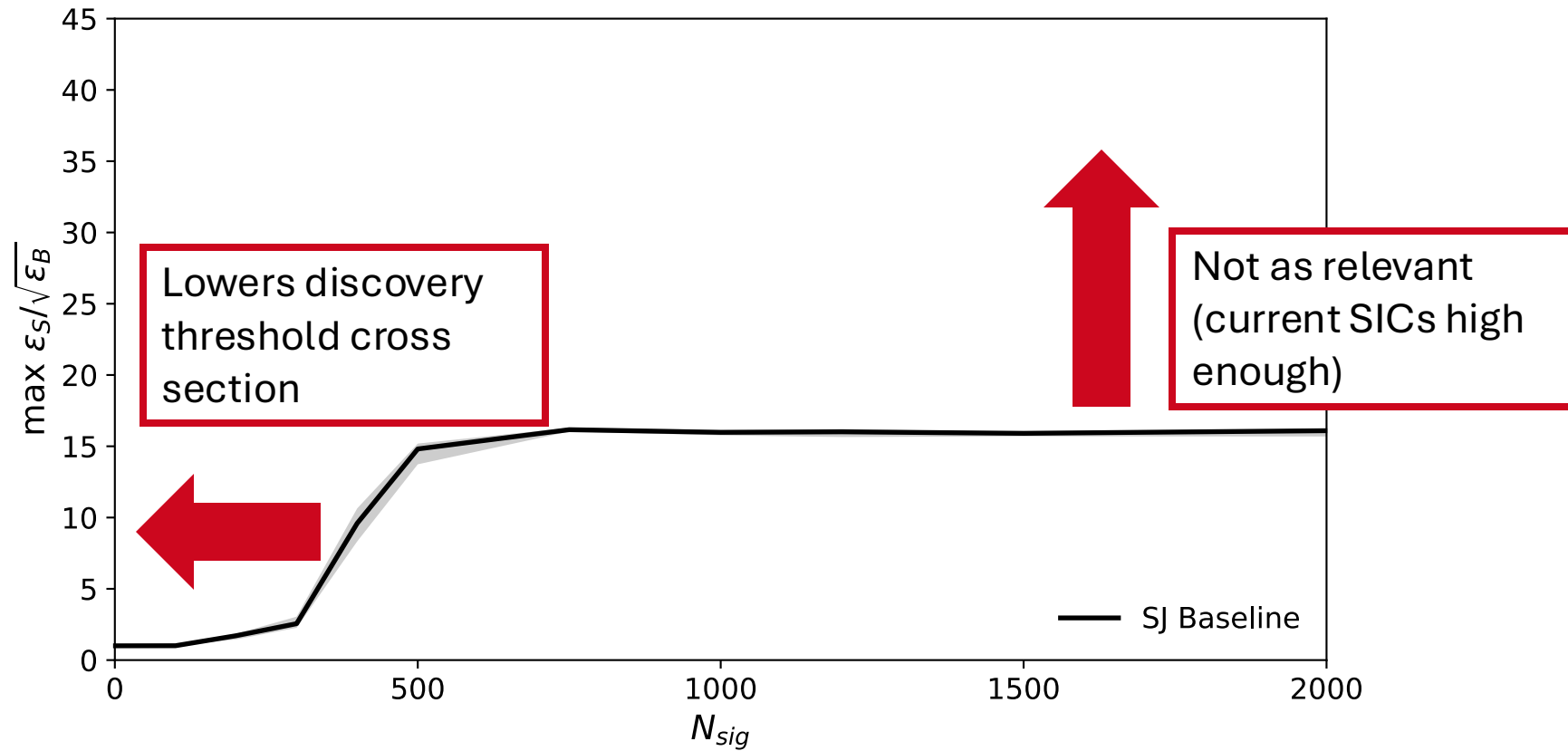


Classifier score to anomaly detection





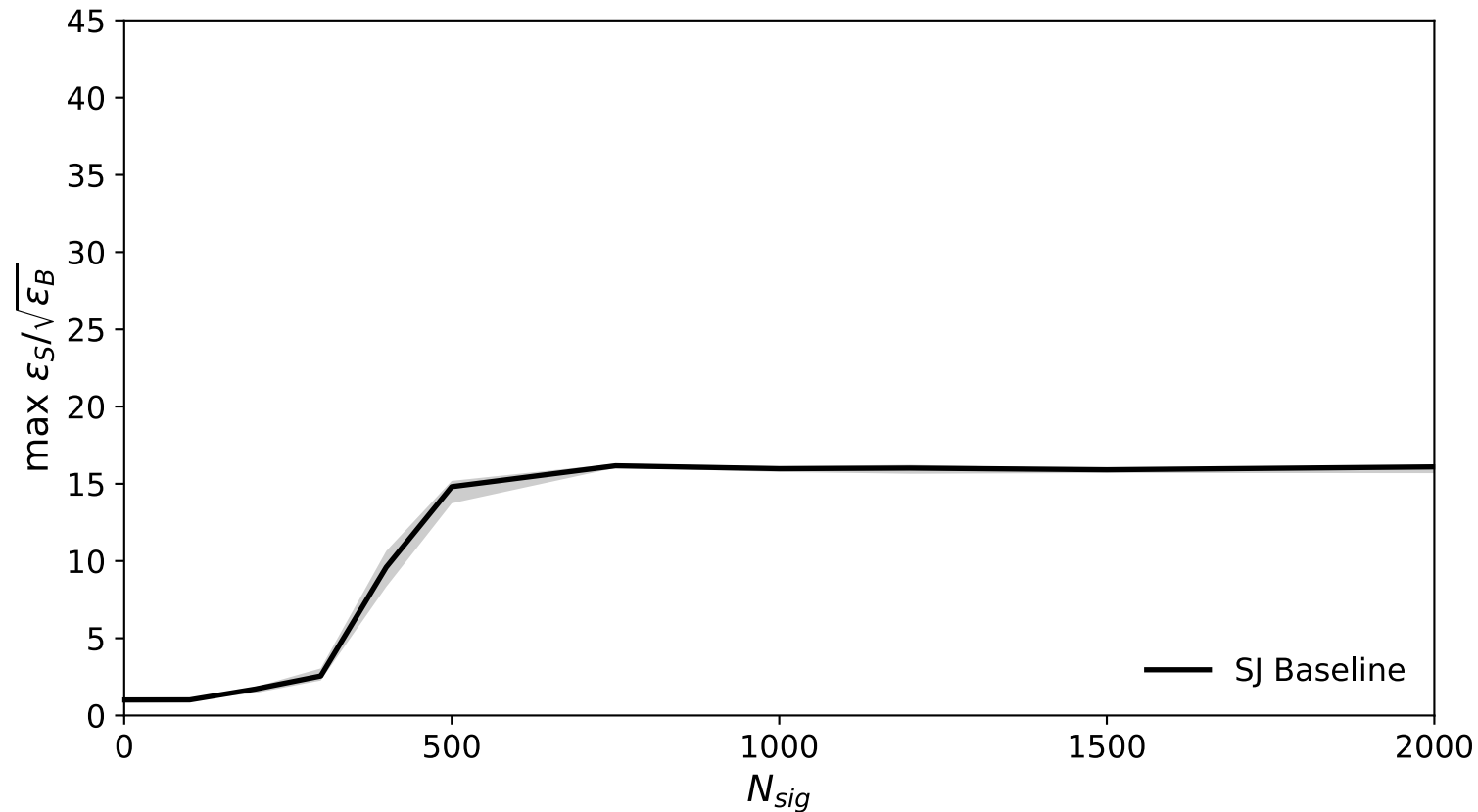




Results

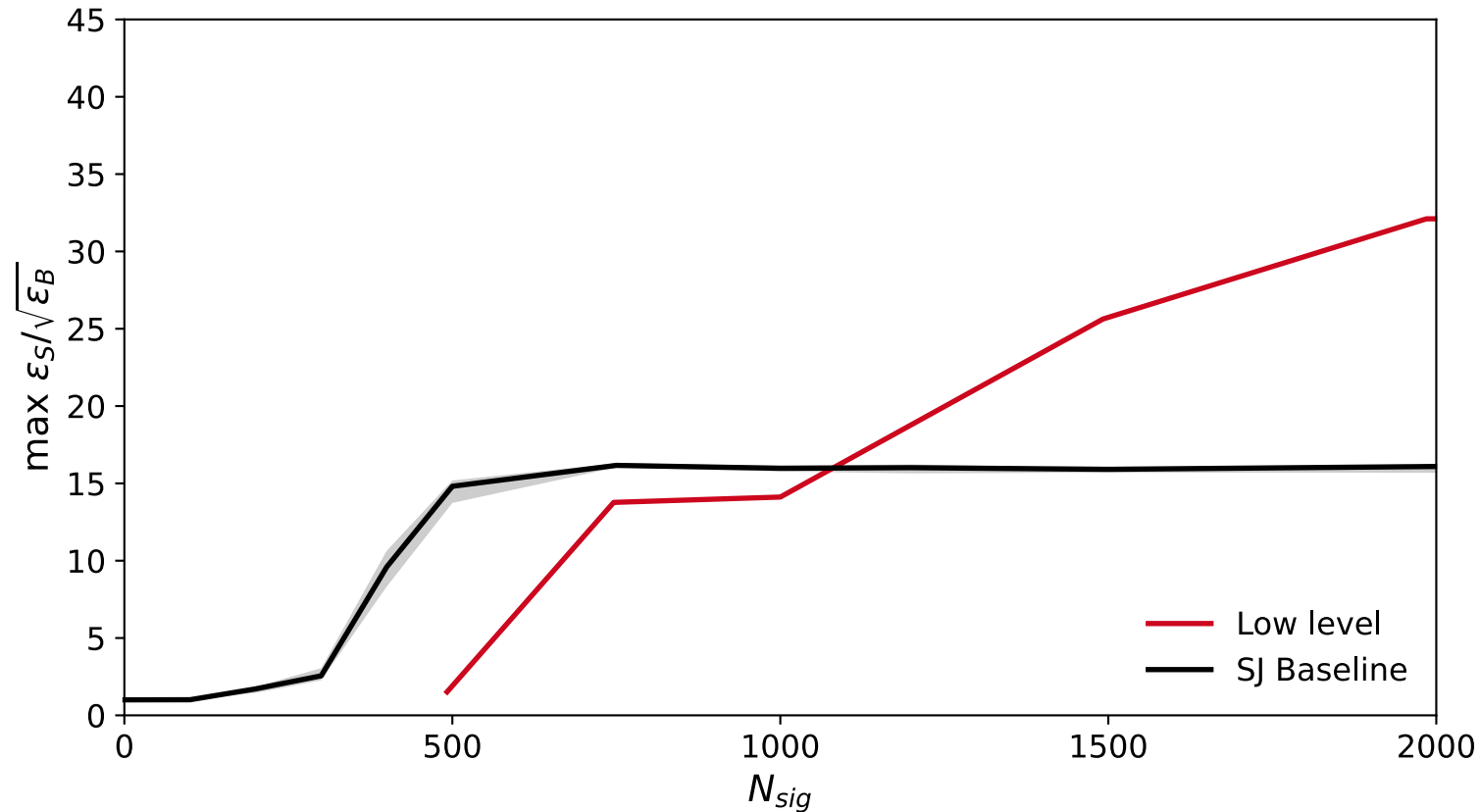
Comparison of Feature Sets

"Tree-Based Algorithms for Weakly Supervised Anomaly Detection" [[2309.13111](#)], T. Finke, **MH** et. al.
"Identifying Anomalous Events Using Low-Level LHC Data", Master Thesis of Joep Geuskens (2023)
Master Thesis of Lukas Lang (2024)



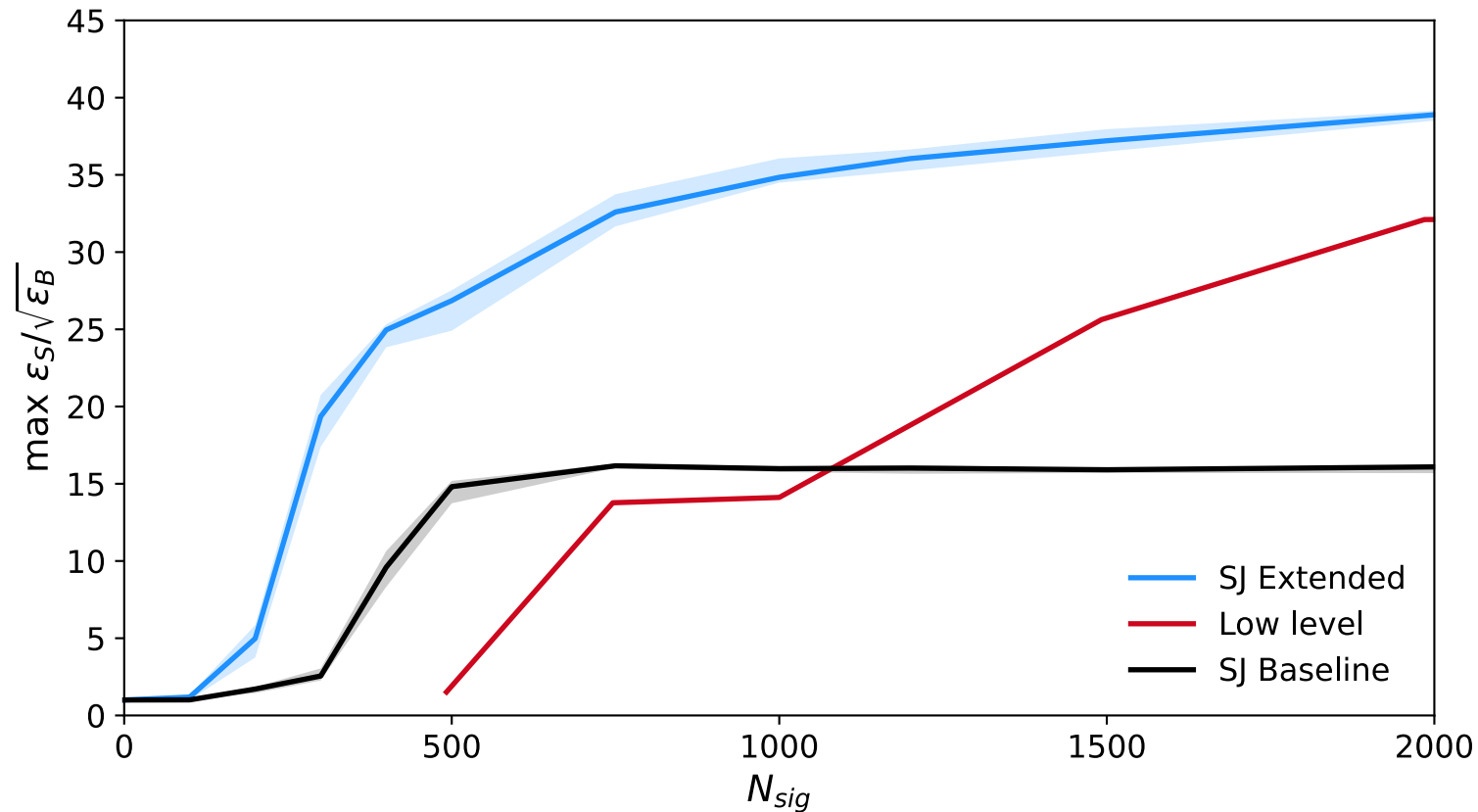
Comparison of Feature Sets

"Tree-Based Algorithms for Weakly Supervised Anomaly Detection" [2309.13111], T. Finke, **MH** et. al.
"Identifying Anomalous Events Using Low-Level LHC Data", Master Thesis of Joep Geuskens (2023)
Master Thesis of Lukas Lang (2024)



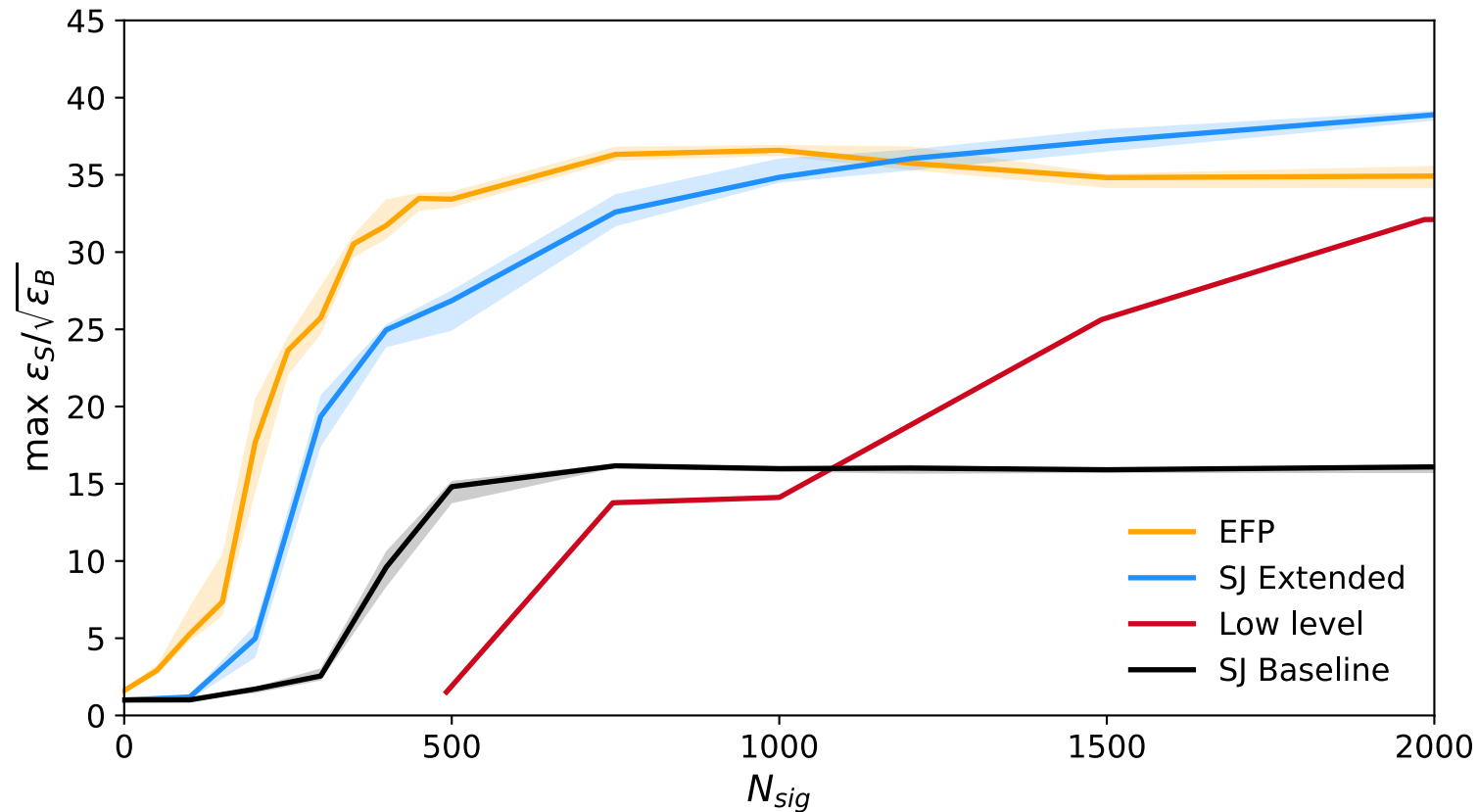
Comparison of Feature Sets

"Tree-Based Algorithms for Weakly Supervised Anomaly Detection" [2309.13111], T. Finke, **MH** et. al.
"Identifying Anomalous Events Using Low-Level LHC Data", Master Thesis of Joep Geuskens (2023)
Master Thesis of Lukas Lang (2024)



Comparison of Feature Sets

"Tree-Based Algorithms for Weakly Supervised Anomaly Detection" [2309.13111], T. Finke, **MH** et. al.
"Identifying Anomalous Events Using Low-Level LHC Data", Master Thesis of Joep Geuskens (2023)
Master Thesis of Lukas Lang (2024)



Summary

- EFPs are useful for anomaly detection
- By choosing the right observables, we can...
 - be more model agnostic
 - be sensitive to lower signal cross sections

Outlook

- Understand why EFPs work so well
 - Currently using interpretable ML methods
- Test EFPs for other signal types
 - Currently working on semi-visible jet
- Test EFPs in more realistic setup