

# Suppression of Continuum Background with Neural Networks for Belle II

Bela Urbschat

Max Planck Institute for Physics, Technical University of Munich

December 19, 2023

2023-12-19

CS with NNs for Belle II

Suppression of Continuum Background with Neural Networks  
for Belle II

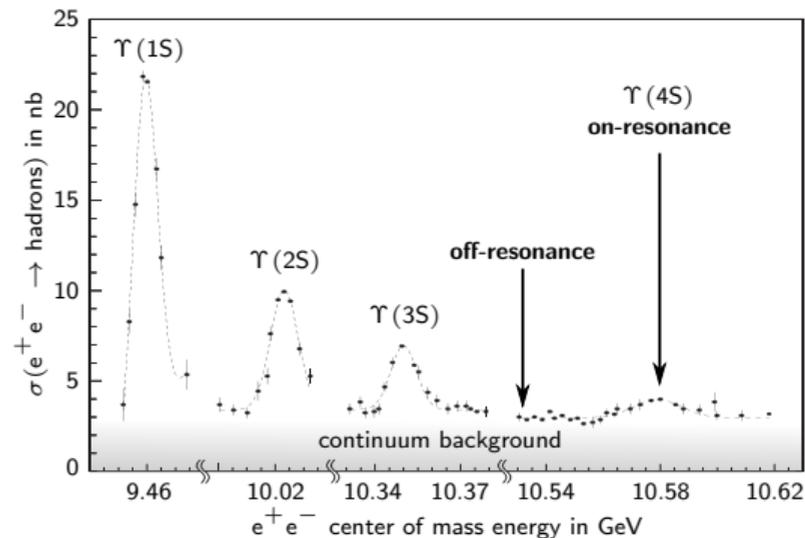
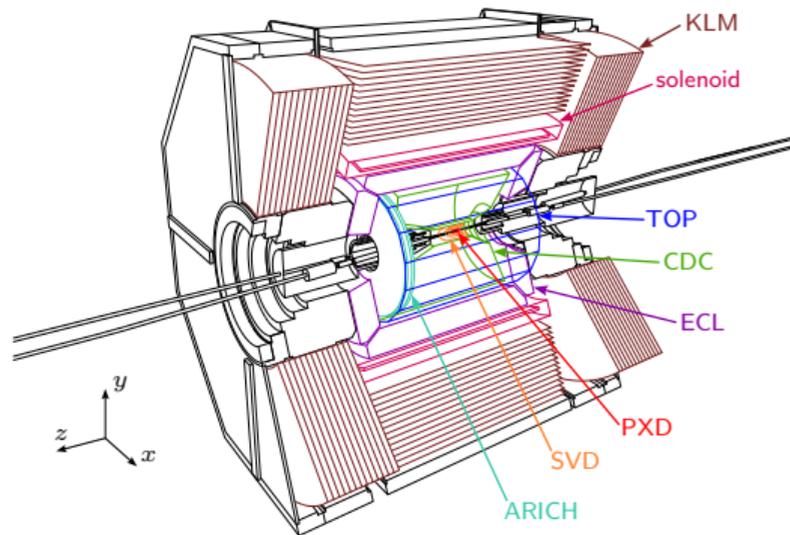
Bela Urbschat

Max Planck Institute for Physics, Technical University of Munich

December 19, 2023

# Belle II/SuperKEKB Overview

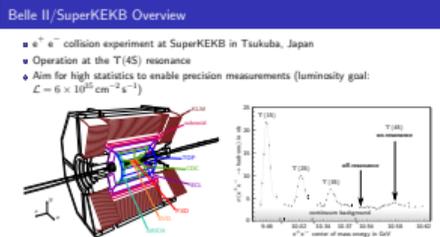
- $e^+ e^-$  collision experiment at SuperKEKB in Tsukuba, Japan
- Operation at the  $\Upsilon(4S)$  resonance
- Aim for high statistics to enable precision measurements (luminosity goal:  $\mathcal{L} = 6 \times 10^{35} \text{ cm}^{-2} \text{ s}^{-1}$ )



2023-12-19

## CS with NNs for Belle II

Belle II/SuperKEKB Overview



## Theoretical Motivation

Theoretical Motivation

SM Null Test ("Isospin Sum Rule")

$$2\mathcal{A}_{CP}(\pi^0 K^+) \frac{\mathcal{B}(\pi^0 K^+) \tau_{B^0}}{\mathcal{B}(\pi^- K^+) \tau_{B^+}} - \mathcal{A}_{CP}(\pi^+ K^0) \frac{\mathcal{B}(\pi^+ K^0) \tau_{B^0}}{\mathcal{B}(\pi^- K^+) \tau_{B^+}} - \mathcal{A}_{CP}(\pi^- K^+) + 2\mathcal{A}_{CP}(\pi^0 K^0) \frac{\mathcal{B}(\pi^0 K^0)}{\mathcal{B}(\pi^- K^+)} = \mathcal{O}(1\%)$$

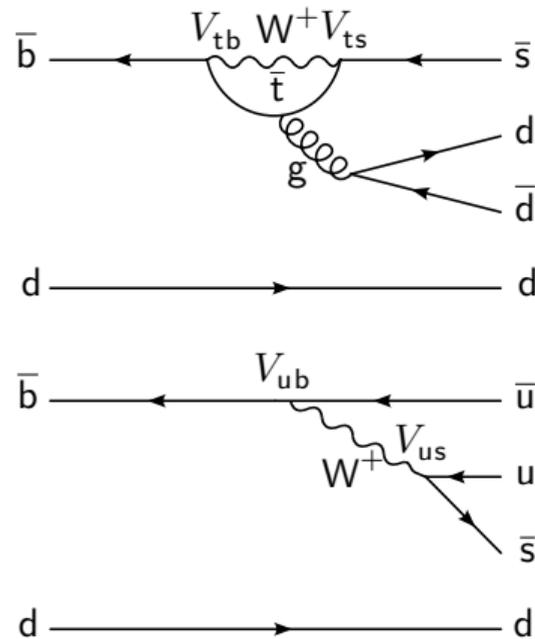
decay	$\mathcal{B} [10^{-6}]$	$\mathcal{A}_{CP}$
$B^0 \rightarrow K^+ \pi^-$	$20.67 \pm 0.37 \pm 0.62$	$-0.072 \pm 0.019 \pm 0.007$
$B^+ \rightarrow K^0 \pi^+$	$24.37 \pm 0.71 \pm 0.86$	$0.046 \pm 0.029 \pm 0.007$
$B^+ \rightarrow K^+ \pi^0$	$13.93 \pm 0.38 \pm 0.71$	$0.013 \pm 0.027 \pm 0.005$
$B^0 \rightarrow K^0 \pi^0$	$10.40 \pm 0.66 \pm 0.60$	$-0.06 \pm 0.15 \pm 0.04$

### SM Null Test ("Isospin Sum Rule")

$$2\mathcal{A}_{CP}(\pi^0 K^+) \frac{\mathcal{B}(\pi^0 K^+) \tau_{B^0}}{\mathcal{B}(\pi^- K^+) \tau_{B^+}}$$

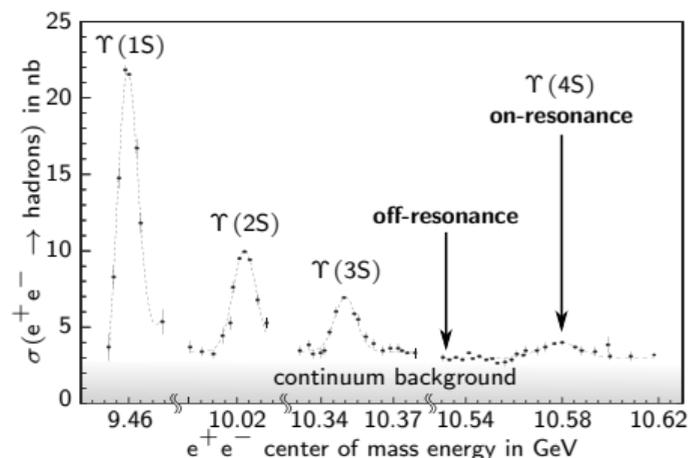
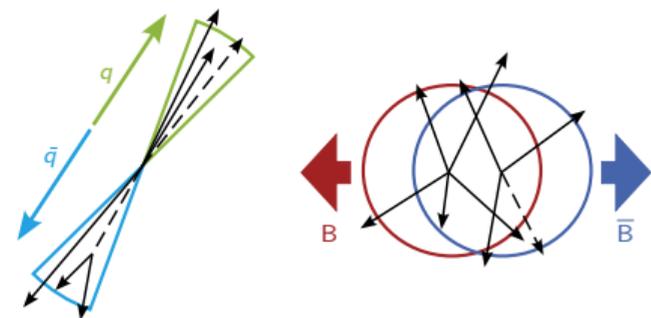
$$- \mathcal{A}_{CP}(\pi^+ K^0) \frac{\mathcal{B}(\pi^+ K^0) \tau_{B^0}}{\mathcal{B}(\pi^- K^+) \tau_{B^+}}$$

$$- \mathcal{A}_{CP}(\pi^- K^+) + 2\mathcal{A}_{CP}(\pi^0 K^0) \frac{\mathcal{B}(\pi^0 K^0)}{\mathcal{B}(\pi^- K^+)} = \mathcal{O}(1\%)$$



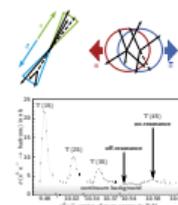
decay	$\mathcal{B} [10^{-6}]$	$\mathcal{A}_{CP}$
$B^0 \rightarrow K^+ \pi^-$	$20.67 \pm 0.37 \pm 0.62$	$-0.072 \pm 0.019 \pm 0.007$
$B^+ \rightarrow K^0 \pi^+$	$24.37 \pm 0.71 \pm 0.86$	$0.046 \pm 0.029 \pm 0.007$
$B^+ \rightarrow K^+ \pi^0$	$13.93 \pm 0.38 \pm 0.71$	$0.013 \pm 0.027 \pm 0.005$
$B^0 \rightarrow K^0 \pi^0$	$10.40 \pm 0.66 \pm 0.60$	$-0.06 \pm 0.15 \pm 0.04$

1. Sum rule as null-test to the SM.
2. Holds in isospin symmetry limit (equal quark masses) (right?)
3. Not exactly = 0, but expected deviation from zero is still much smaller than experimental uncertainties.
4. Highlight the  $B \rightarrow K\pi$  decay modes appearing in sum rule.
5. Highlight that  $B^0 \rightarrow K^0 \pi^0$  is measured worst (also as not self tagging)
6. NP (particles) could contribute to loops.



- $e^+e^- \rightarrow q\bar{q}$  where  $q = u, d, c, s$
- dominating background for B decay measurements (other backgrounds easily rejected)
- excess energy results in hadronic jets
- topology distinct from signal decays

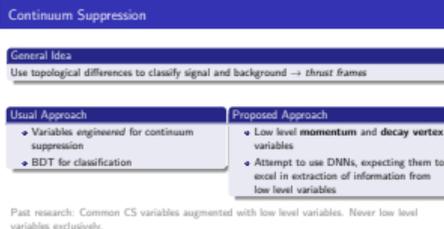
## Continuum Background



- $e^+e^- \rightarrow q\bar{q}$  where  $q = u, d, c, s$
- dominating background for B decay measurements (other backgrounds easily rejected)
- excess energy results in hadronic jets
- topology distinct from signal decays

1. Point to the event shape figure.
2. Explain uniform  $q\bar{q}$  background in resonances figure.

## Continuum Suppression



### General Idea

Use topological differences to classify signal and background → *thrust frames*

### Usual Approach

- Variables *engineered* for continuum suppression
- BDT for classification

### Proposed Approach

- Low level **momentum** and **decay vertex** variables
- Attempt to use DNNs, expecting them to excel in extraction of information from low level variables

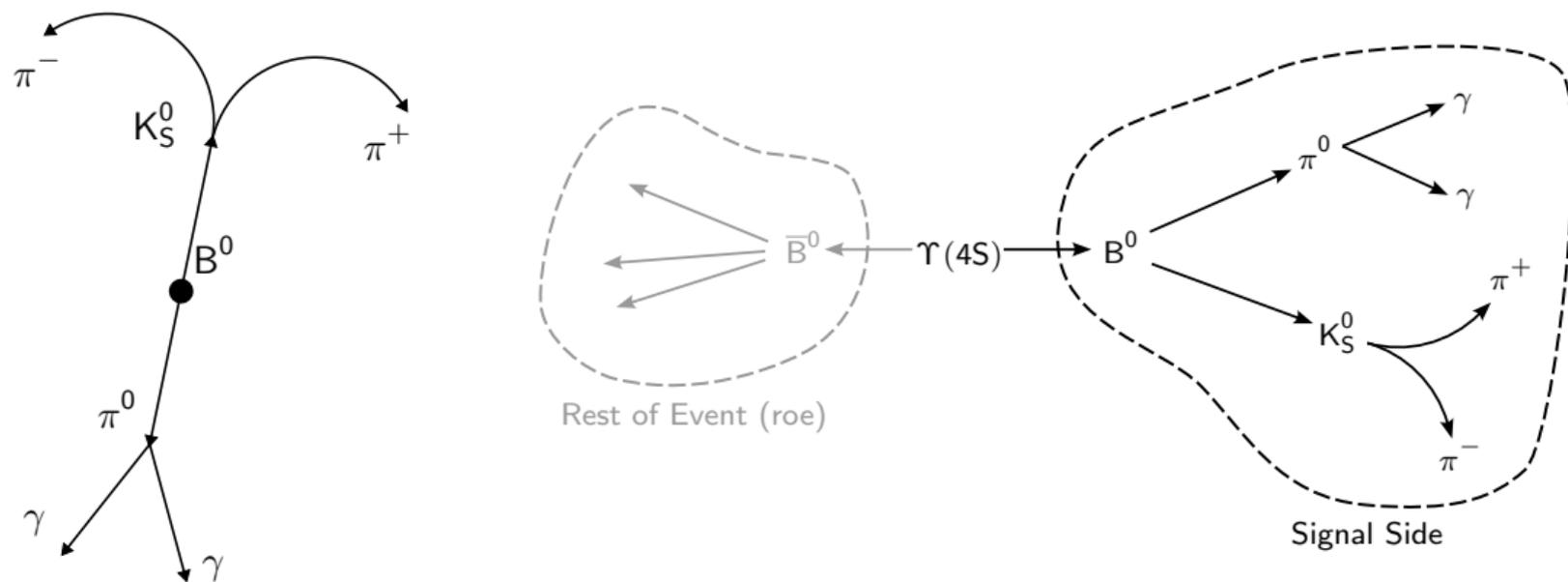
Past research: Common CS variables augmented with low level variables. Never low level variables exclusively.

1. Make sure to explain thrust frames!
2. Momentum/vertex variables in theory should contain all the information of event shape.

# Reconstruction and Data

Chose  $B^0 \rightarrow K_S^0(\pi^+\pi^-)\pi^0(\gamma\gamma)$  as an example

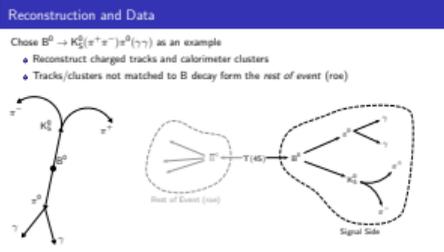
- Reconstruct charged tracks and calorimeter clusters
- Tracks/clusters not matched to B decay form the *rest of event* (roe)



2023-12-19

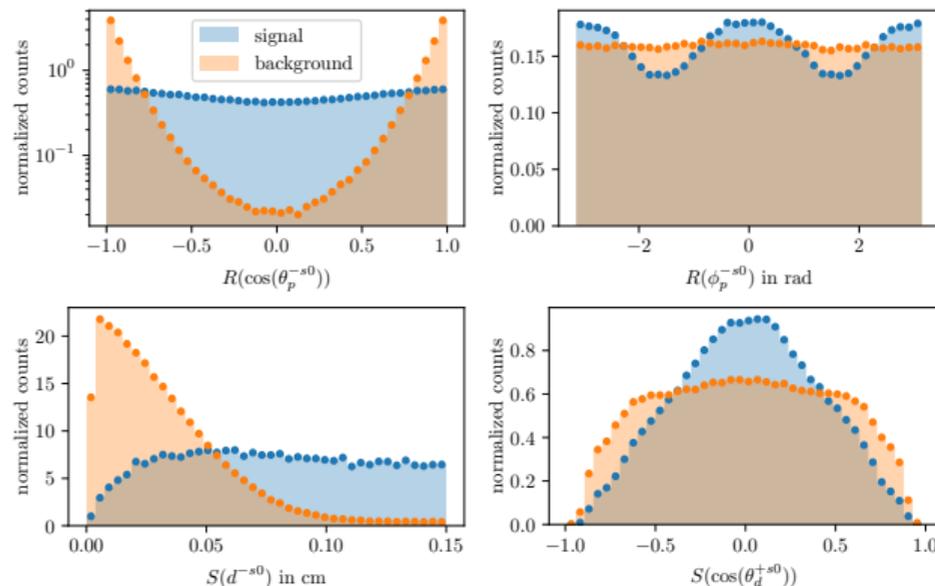
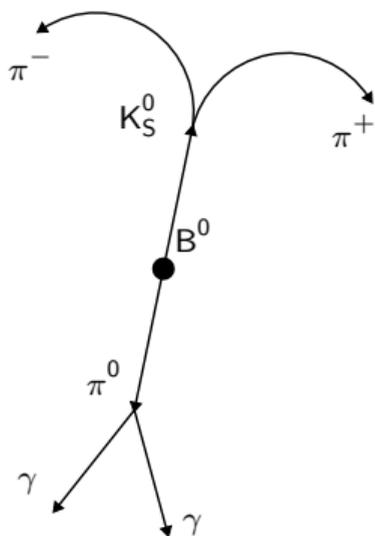
CS with NNs for Belle II

└ Reconstruction and Data

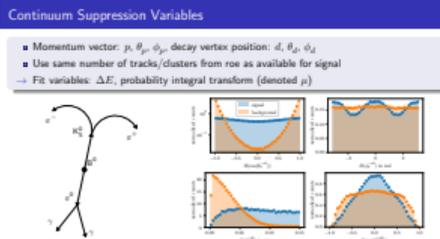


1. Explain signal thrust/roe thrust using figure on the right

- Momentum vector:  $p, \theta_p, \phi_p$ , decay vertex position:  $d, \theta_d, \phi_d$
  - Use same number of tracks/clusters from roe as available for signal
- Fit variables:  $\Delta E$ , probability integral transform (denoted  $\mu$ )



## Continuum Suppression Variables



1. Note that we attempted to use more variables from roe which did not result in a significant performance gain
2. Explain chosen orders tracks/clusters for variables
3. Explain notation (briefly)
4. Explain variables that do not fall under the naming scheme
5. Explain intuition for polar angle distribution based on antiparallel/random alignment of thrust axes.

## Classifiers Used

### Boosted Decision Trees (BTDs)

- Robust classifiers
- Give good baseline for expected performance
- Here no in-depth hyperparameter tuning

### Deep Neural Networks (DNNs)

- Initial motivation: Possibly better at utilizing information from low level variables → better performance?
- Turn out to be much more delicate/difficult to handle
- Main subject of studies for this thesis

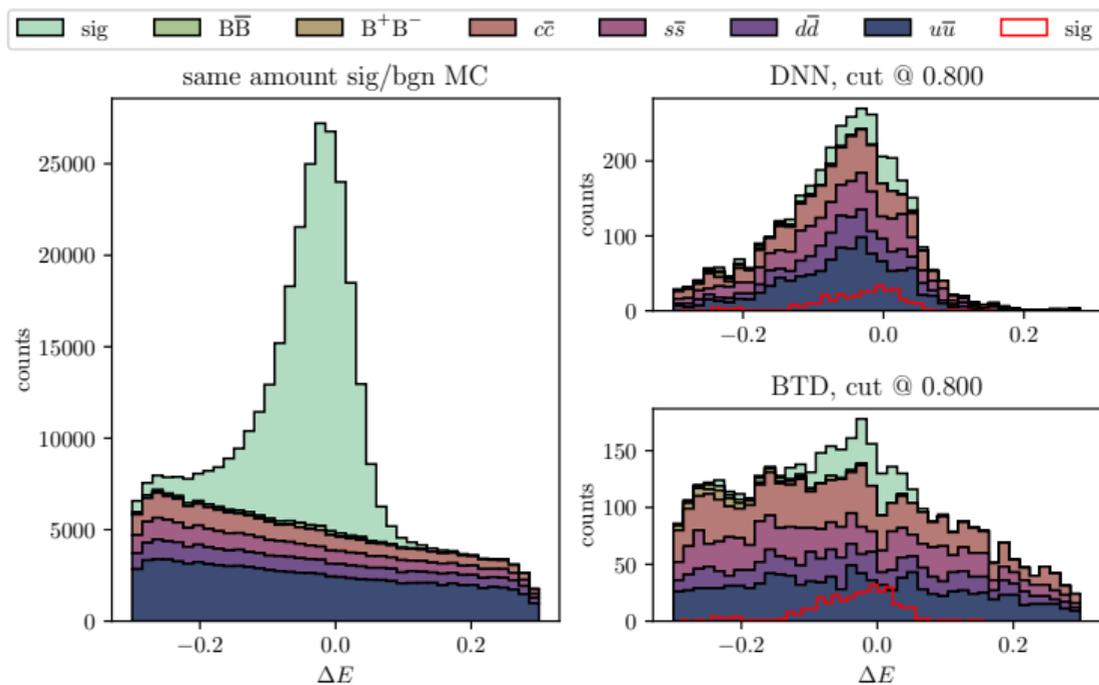
### Boosted Decision Trees (BTDs)

- Robust classifiers
- Give good baseline for expected performance
- Here no in-depth hyperparameter tuning

### Deep Neural Networks (DNNs)

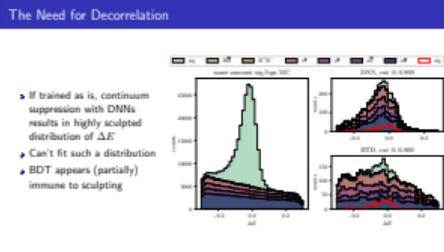
- Initial motivation: Possibly better at utilizing information from low level variables → better performance?
- Turn out to be much more delicate/difficult to handle
- Main subject of studies for this thesis

- If trained as is, continuum suppression with DNNs results in highly sculpted distribution of  $\Delta E$
- Can't fit such a distribution
- BDT appears (partially) immune to sculpting



## The Need for Decorrelation

1. Explain expected shape using left plot.
2. Highlight that fit with observed level of sculpting is clearly impossible.



- Efficiently estimable correlation metric, capturing also non-linear correlations
- Only one further hyperparameter introduced

Total loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{classifier}}(\vec{y}, \vec{y}_{\text{true}}) + \lambda \cdot \text{dCorr}(\vec{z}, \vec{y})$$

However tuning still difficult:

- Too large  $\lambda$  degrades performance
- Effectiveness of decorrelation also influenced by other hyperparameters (batch size, network architecture)
- Systematic tuning extremely difficult due to conflicting objectives

→ Studies with preliminary hyperparameters to better understand behavior

### └ Tools(s) for Decorrelation

1. Also mention that adversary networks have been implemented, but could not be sufficiently tuned for this thesis.
2. Explain symbols in the equation!
3. Mention that classifier loss is binary cross-entropy.
4. Explain the conflicting objectives of best performance and effective decorrelation (problem: performance always better for correlated classifier).

Tools(s) for Decorrelation

Distance Correlation

- Efficiently estimable correlation metric, capturing also non-linear correlations
- Only one further hyperparameter introduced

Total loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{classifier}}(\vec{y}, \vec{y}_{\text{true}}) + \lambda \cdot \text{dCorr}(\vec{z}, \vec{y})$$

However tuning still difficult:

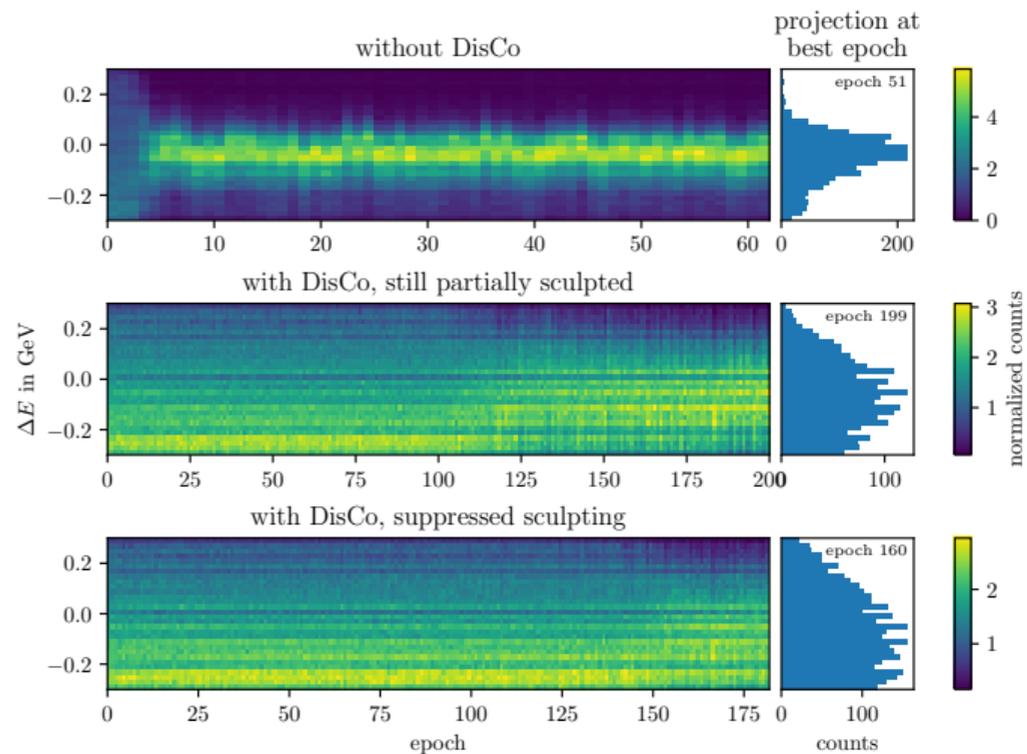
- Too large  $\lambda$  degrades performance
- Effectiveness of decorrelation also influenced by other hyperparameters (batch size, network architecture)
- Systematic tuning extremely difficult due to conflicting objectives

→ Studies with preliminary hyperparameters to better understand behavior

# Monitoring DNN Training

## Evolution of $\Delta E$ (Background) Distribution

- Preliminary hyperparameters with different values for  $\lambda$  (0, 1, 1.8)
- Achieved decorrelation still not satisfactory
- Sculpting (partially suppressed) suddenly starts after sufficient number of epochs

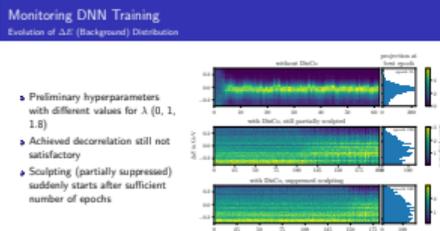


2023-12-19

CS with NNs for Belle II

## Monitoring DNN Training

1. Highlight that after sufficient training (or epochs), correlation (more or less suddenly) starts  $\rightarrow$  decorrelation is unstable.
2. Mention that here the goal was to reach lower sculpting than BDT in hope of this improving fit quality (i.e. lowering the statistical uncertainties). Thus the best decorrelation is still not satisfactory.
3. Distributions are normalized at each epoch!



	prelim. value	final value	description
$n_{\text{layers}}$	5	5	number of layers
$n_{\text{neurons},0}$	100	100	1st dense layer neurons
$n_{\text{neurons},1}$	100	100	2nd dense layer neurons
$n_{\text{neurons},2}$	4	6	3rd dense layer neurons
$n_{\text{neurons},3}$	100	100	4th dense layer neurons
$n_{\text{neurons},4}$	100	100	5th dense layer neurons
weight decay	0.000142	0.000142	Weight decay for AdamW
learning rate	0.002	0.015	learning rate
dCorr on bgn	True	True	choice to compute dCorr on only background events
$\lambda$	1.8	2	scale of dCorr in total loss
$s_{\lambda}$	7.5	7.5	scale factor for $\lambda$ when dCorr computed on bgn only
batch size	2048	16384	number of events in a minibatch

→ In the following DNN with applied decorrelation and final hyperparameters is referred to as *DisCoDNN*

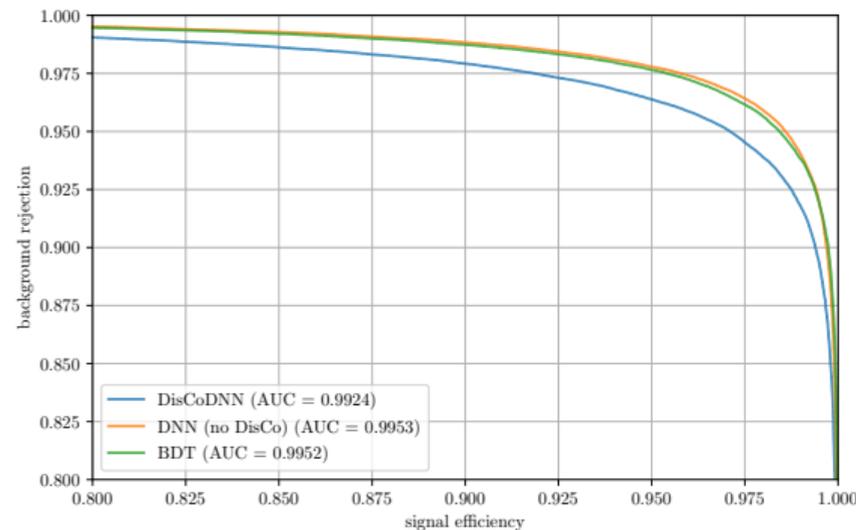
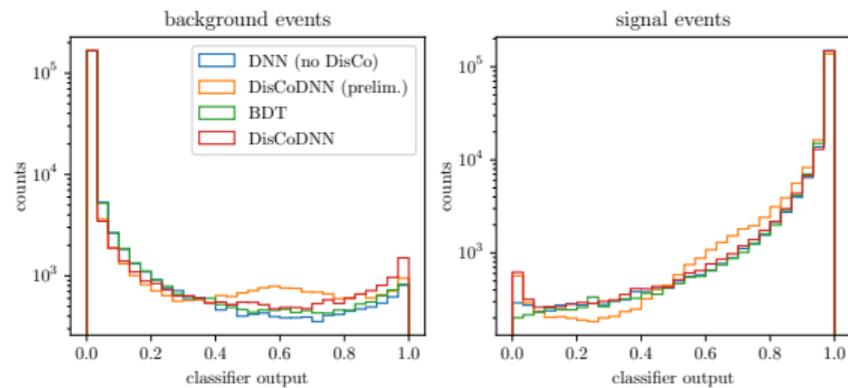
## Choice of Hyperparameters

	prelim. value	final value	description
$n_{\text{layers}}$	5	5	number of layers
$n_{\text{neurons},0}$	100	100	1st dense layer neurons
$n_{\text{neurons},1}$	100	100	2nd dense layer neurons
$n_{\text{neurons},2}$	4	6	3rd dense layer neurons
$n_{\text{neurons},3}$	100	100	4th dense layer neurons
$n_{\text{neurons},4}$	100	100	5th dense layer neurons
weight decay	0.000142	0.000142	Weight decay for AdamW
learning rate	0.002	0.015	learning rate
dCorr on bgn	True	True	choice to compute dCorr on only background events
$\lambda$	1.8	2	scale of dCorr in total loss
$s_{\lambda}$	7.5	7.5	scale factor for $\lambda$ when dCorr computed on bgn only
batch size	2048	16384	number of events in a minibatch

→ In the following DNN with applied decorrelation and final hyperparameters is referred to as *DisCoDNN*

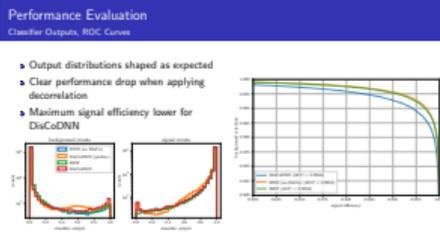
1. Highlight the "unusual" hyperparameters: Large batch size, bottleneck architecture

- Output distributions shaped as expected
- Clear performance drop when applying decorrelation
- Maximum signal efficiency lower for DisCoDNN



2023-12-19

Performance Evaluation



1. Note that prelim. DisCoDNN only shown as reference for *not good* output distribution.

# $\Delta E$ and $\mu$ after Continuum Suppression

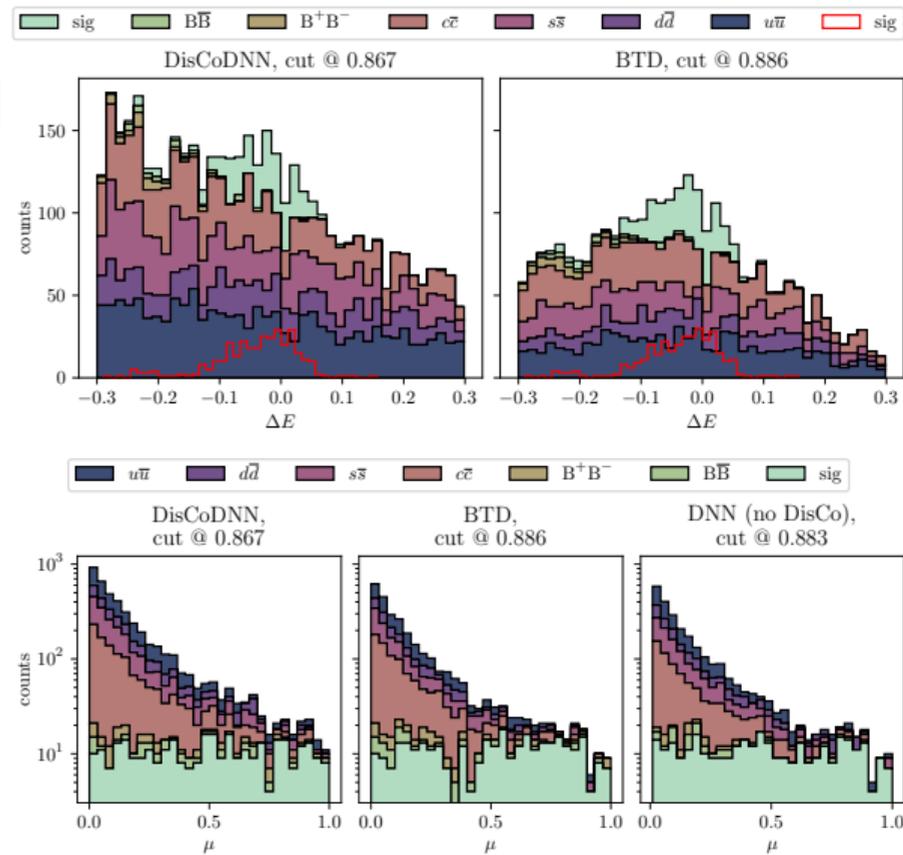
Cuts always chosen for 90% signal efficiency

$\Delta E$ :

- Effective decorrelation with DisCoDNN
- Remaining (but acceptable) sculpting for BDT
  - Could further investigate decorrelation for BDTs
- Overall better background suppression with BDT at same signal efficiency

$\mu$ :

- Shapes unaffected by decorrelation
- Reasonably flat signal contributions

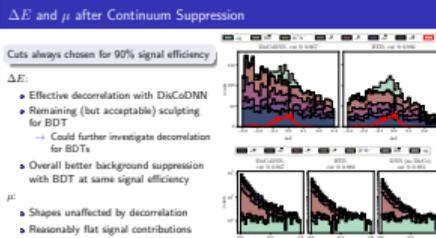


2023-12-19

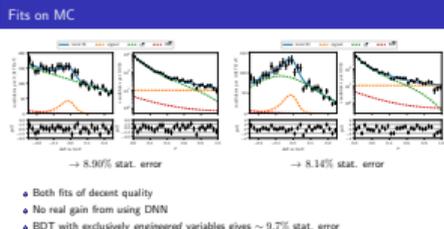
## CS with NNs for Belle II

$\Delta E$  and  $\mu$  after Continuum Suppression

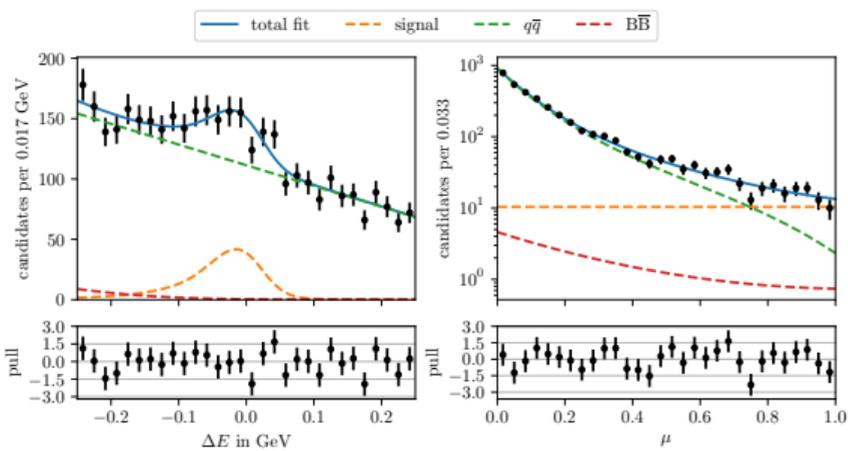
1. Maybe mention how cut positions were determined/that they were determined using an appropriate procedure.



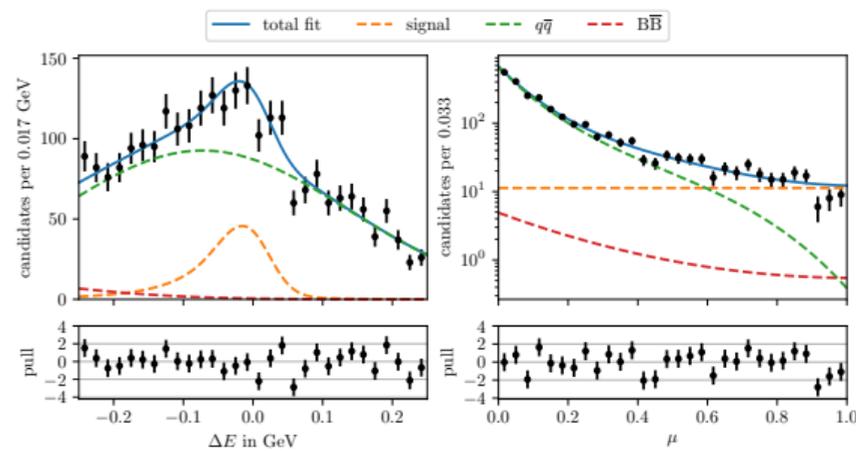
## Fits on MC



- Both fits of decent quality
- No real gain from using DNN
- BDT with exclusively *engineered* variables gives  $\sim 9.7\%$  stat. error



→ 8.90% stat. error



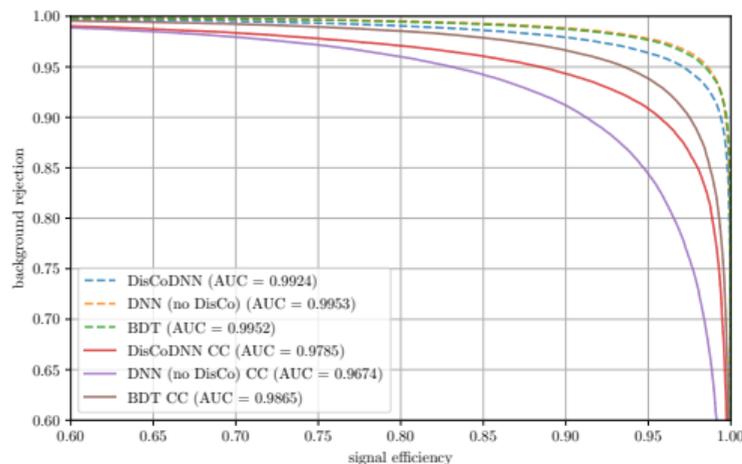
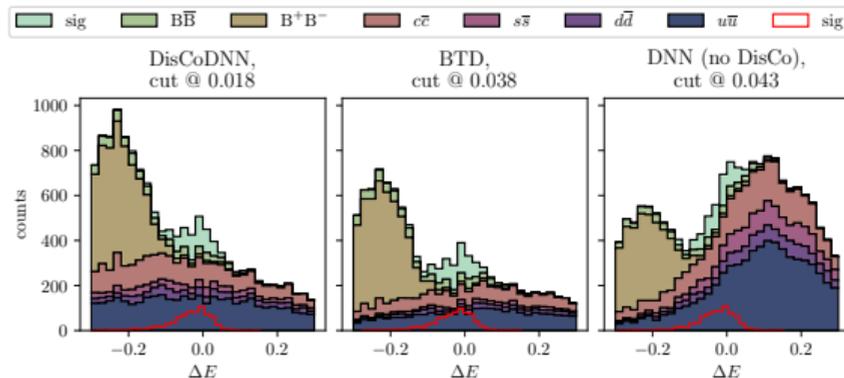
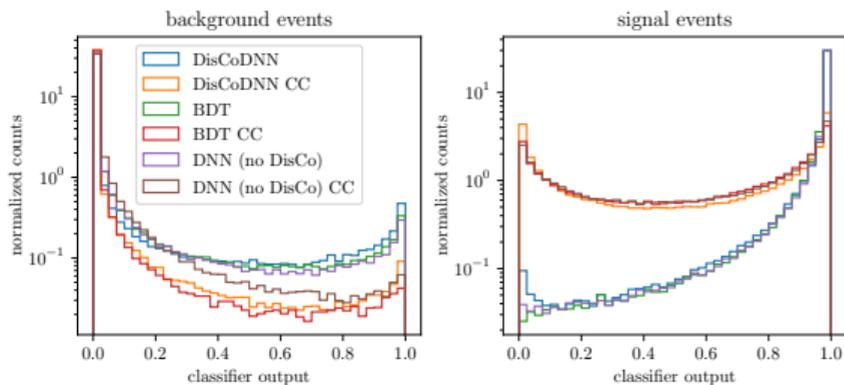
→ 8.14% stat. error

- Both fits of decent quality
- No real gain from using DNN
- BDT with exclusively *engineered* variables gives  $\sim 9.7\%$  stat. error

# Classifier Generalizability

Apply to topologically similar control channel  
 $B^0 \rightarrow \bar{D}^0(K^+\pi^-\pi^0(\gamma\gamma))$

- All classifiers fail to identify signal
- Surprisingly good continuum suppression possible with very loose cuts
- DNN without decorrelation fails spectacularly

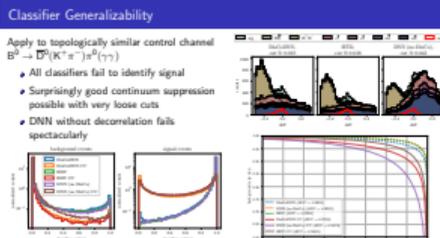


# CS with NNs for Belle II

2023-12-19

## Classifier Generalizability

1. Mention that this demonstrates the problem of generalizability!
2. Note that DNN (no DisCo) seems *not* to just "compute" or estimate  $\Delta E$  and then more or less cut on that, as  $B\bar{B}$  background remains!
3. Possibly the correlations are then what allows the DNN to sculpt  $\Delta E$ . This would make sense as DisCoDNN does not really rely on correlations.



- Introduced set of low level continuum suppression variables
  - Prepared BDT and DNNs using introduced variables, expecting DNN to profit from those
  - DNNs require decorrelation, which most likely limits their performance
  - Fits on MC show similar accuracies for BDT/DNN but slightly better than BDT with common CS variables
- Low level CS variables could reduce statistical errors but further investigation (e.g. systematics etc.) needed for final judgement

## For the Future

- Study influence of single variables on sculpting (to possibly exclude them)
- Impact on performance with alternative decorrelation method (e.g. adversarial networks)
- Application of similar decorrelation to BDT
- Application within a fully fledged analysis (including systematics etc.)

## Conclusion & Outlook

1. In fact the sculpting also happens with only *engineered* variables. It's just that so far everyone always used BDTs which are not subject to that issue.

- Introduced set of low level continuum suppression variables
  - Prepared BDT and DNNs using introduced variables, expecting DNN to profit from those
  - DNNs require decorrelation, which most likely limits their performance
  - Fits on MC show similar accuracies for BDT/DNN but slightly better than BDT with common CS variables
- Low level CS variables could reduce statistical errors but further investigation (e.g. systematics etc.) needed for final judgement

### For the Future

- Study influence of single variables on sculpting (to possibly exclude them)
- Impact on performance with alternative decorrelation method (e.g. adversarial networks)
- Application of similar decorrelation to BDT
- Application within a fully fledged analysis (including systematics etc.)

Backup

- Generic (run independent) MC ( $q\bar{q}$  where  $q = u, d, s, c$  &  $B\bar{B}$ ):  $1 \text{ ab}^{-1}$
- Pure signal MC for signal channel and control channel:  $4 \times 10^6$  and  $2 \times 10^6$  events produced resulting in 1 019 638 and 523 183 reconstructed events respectively
- Physics data:  $361.65 \text{ fb}^{-1}$
- Off-resonance generic MC ( $q\bar{q}$  where  $q = u, d, s, c$ ):  $169.328 \text{ fb}^{-1}$
- Off-resonance data:  $42.28 \text{ fb}^{-1}$

## MC Modeling

- Problems with the available samples ( $\tau^- \tau^+$ , momentum corrections) remain
  - MC modeling overall not bad, considering the above
- Further investigation needed for final judgment

## Data Samples

- Generic (run independent) MC ( $q\bar{q}$  where  $q = u, d, s, c$  &  $B\bar{B}$ ):  $1 \text{ ab}^{-1}$
- Pure signal MC for signal channel and control channel:  $4 \times 10^6$  and  $2 \times 10^6$  events produced resulting in 1 019 638 and 523 183 reconstructed events respectively
- Physics data:  $361.65 \text{ fb}^{-1}$
- Off-resonance generic MC ( $q\bar{q}$  where  $q = u, d, s, c$ ):  $169.328 \text{ fb}^{-1}$
- Off-resonance data:  $42.28 \text{ fb}^{-1}$

### MC Modeling

- Problems with the available samples ( $\tau^- \tau^+$ , momentum corrections) remain
  - MC modeling overall not bad, considering the above
- Further investigation needed for final judgment

# Continuum Suppression Variables

All Input Variables

$\Delta z$	$R(\cos(\theta_p^{+s0}))$	$R(\phi_p^{0s0})$	$S(\cos(\theta_p^{-s0}))$	$S(\phi_d^{-r0})$	$S(p^{+s0})$
$\cos(\theta_{SR})$	$R(\cos(\theta_d^{-r0}))$	$R(\phi_p^{0s1})$	$S(\cos(\theta_p^{+r0}))$	$S(\phi_d^{-s0})$	$S(\phi_p^{0r0})$
$\cos(\theta_{Sz})$	$R(\cos(\theta_d^{-s0}))$	$R(\phi_p^{-r0})$	$S(\cos(\theta_p^{+s0}))$	$S(\phi_d^{+r0})$	$S(\phi_p^{0r1})$
$M'_{bc}$	$R(\cos(\theta_d^{+r0}))$	$R(\phi_p^{-s0})$	$S(\cos(\theta_d^{-r0}))$	$S(\phi_d^{+s0})$	$S(\phi_p^{0s0})$
$R(\cos(\theta_p^{0r0}))$	$R(\cos(\theta_d^{+s0}))$	$R(\phi_p^{+r0})$	$S(\cos(\theta_d^{-s0}))$	$S(p^{0r0})$	$S(\phi_p^{0s1})$
$R(\cos(\theta_p^{0r1}))$	$R(\phi_d^{-r0})$	$R(\phi_p^{+s0})$	$S(\cos(\theta_d^{+r0}))$	$S(p^{0r1})$	$S(\phi_p^{-r0})$
$R(\cos(\theta_p^{0s0}))$	$R(\phi_d^{-s0})$	$S(\cos(\theta_p^{0r0}))$	$S(\cos(\theta_d^{+s0}))$	$S(p^{0s0})$	$S(\phi_p^{-s0})$
$R(\cos(\theta_p^{0s1}))$	$R(\phi_d^{+r0})$	$S(\cos(\theta_p^{0r1}))$	$S(d^{-r0})$	$S(p^{0s1})$	$S(\phi_p^{+r0})$
$R(\cos(\theta_p^{-r0}))$	$R(\phi_d^{+s0})$	$S(\cos(\theta_p^{0s0}))$	$S(d^{-s0})$	$S(p^{-r0})$	$S(\phi_p^{+s0})$
$R(\cos(\theta_p^{-s0}))$	$R(\phi_p^{0r0})$	$S(\cos(\theta_p^{0s1}))$	$S(d^{+r0})$	$S(p^{-s0})$	
$R(\cos(\theta_p^{+r0}))$	$R(\phi_p^{0r1})$	$S(\cos(\theta_p^{-r0}))$	$S(d^{+s0})$	$S(p^{+r0})$	

2023-12-19

CS with NNs for Belle II

Continuum Suppression Variables

Continuum Suppression Variables  
All Input Variables

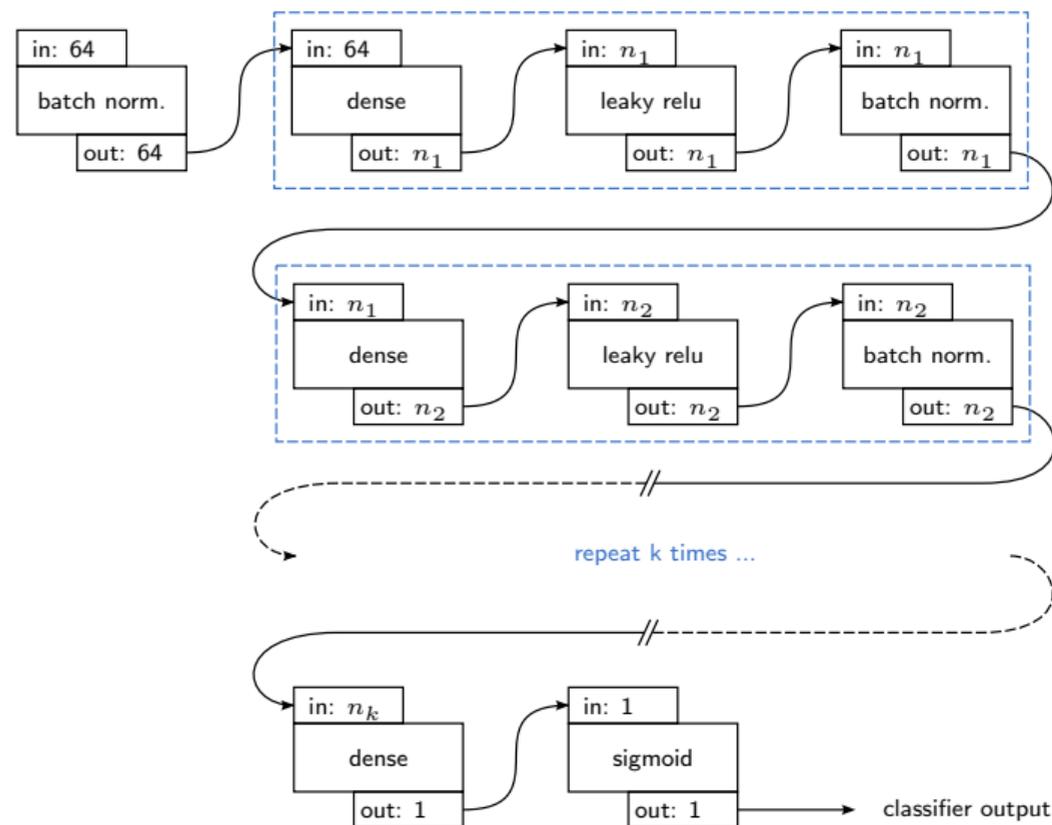
$\Delta z$	$R(\cos(\theta_p^{+s0}))$	$R(\phi_p^{0s0})$	$S(\cos(\theta_p^{-s0}))$	$S(\phi_d^{-r0})$	$S(p^{+s0})$
$\cos(\theta_{SR})$	$R(\cos(\theta_d^{-r0}))$	$R(\phi_p^{0s1})$	$S(\cos(\theta_p^{+r0}))$	$S(\phi_d^{-s0})$	$S(\phi_p^{0r0})$
$\cos(\theta_{Sz})$	$R(\cos(\theta_d^{-s0}))$	$R(\phi_p^{-r0})$	$S(\cos(\theta_p^{+s0}))$	$S(\phi_d^{+r0})$	$S(\phi_p^{0r1})$
$M'_{bc}$	$R(\cos(\theta_d^{+r0}))$	$R(\phi_p^{-s0})$	$S(\cos(\theta_d^{-r0}))$	$S(\phi_d^{+s0})$	$S(\phi_p^{0s0})$
$R(\cos(\theta_p^{0r0}))$	$R(\cos(\theta_d^{+s0}))$	$R(\phi_p^{+r0})$	$S(\cos(\theta_d^{-s0}))$	$S(p^{0r0})$	$S(\phi_p^{0s1})$
$R(\cos(\theta_p^{0r1}))$	$R(\phi_d^{-r0})$	$R(\phi_p^{+s0})$	$S(\cos(\theta_d^{+r0}))$	$S(p^{0r1})$	$S(\phi_p^{-r0})$
$R(\cos(\theta_p^{0s0}))$	$R(\phi_d^{-s0})$	$S(\cos(\theta_p^{0r0}))$	$S(\cos(\theta_d^{+s0}))$	$S(p^{0s0})$	$S(\phi_p^{-s0})$
$R(\cos(\theta_p^{0s1}))$	$R(\phi_d^{+r0})$	$S(\cos(\theta_p^{0r1}))$	$S(d^{-r0})$	$S(p^{0s1})$	$S(\phi_p^{+r0})$
$R(\cos(\theta_p^{-r0}))$	$R(\phi_d^{+s0})$	$S(\cos(\theta_p^{0s0}))$	$S(d^{-s0})$	$S(p^{-r0})$	$S(\phi_p^{+s0})$
$R(\cos(\theta_p^{-s0}))$	$R(\phi_p^{0r0})$	$S(\cos(\theta_p^{0s1}))$	$S(d^{+r0})$	$S(p^{-s0})$	
$R(\cos(\theta_p^{+r0}))$	$R(\phi_p^{0r1})$	$S(\cos(\theta_p^{-r0}))$	$S(d^{+s0})$	$S(p^{+r0})$	

## Network Architecture:

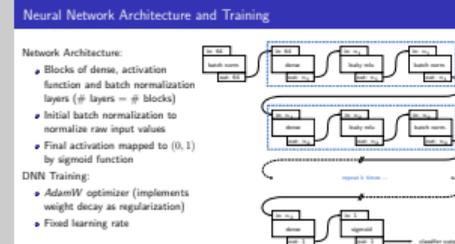
- Blocks of dense, activation function and batch normalization layers (# layers = # blocks)
- Initial batch normalization to normalize raw input values
- Final activation mapped to (0, 1) by sigmoid function

## DNN Training:

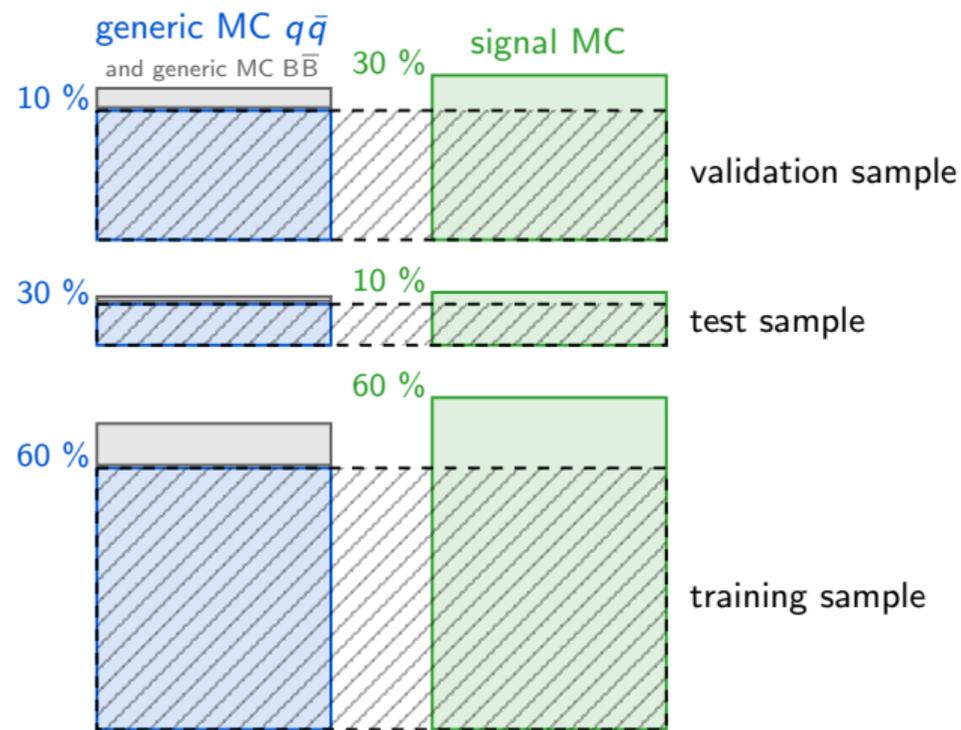
- *AdamW* optimizer (implements weight decay as regularization)
- Fixed learning rate



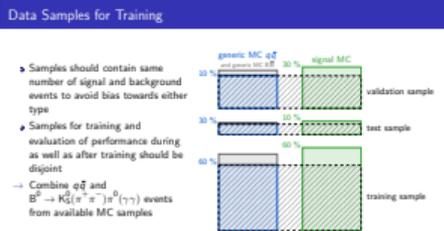
## Neural Network Architecture and Training



- Samples should contain same number of signal and background events to avoid bias towards either type
  - Samples for training and evaluation of performance during as well as after training should be disjoint
- Combine  $q\bar{q}$  and  $B^0 \rightarrow K_S^0(\pi^+\pi^-\pi^0)(\gamma\gamma)$  events from available MC samples



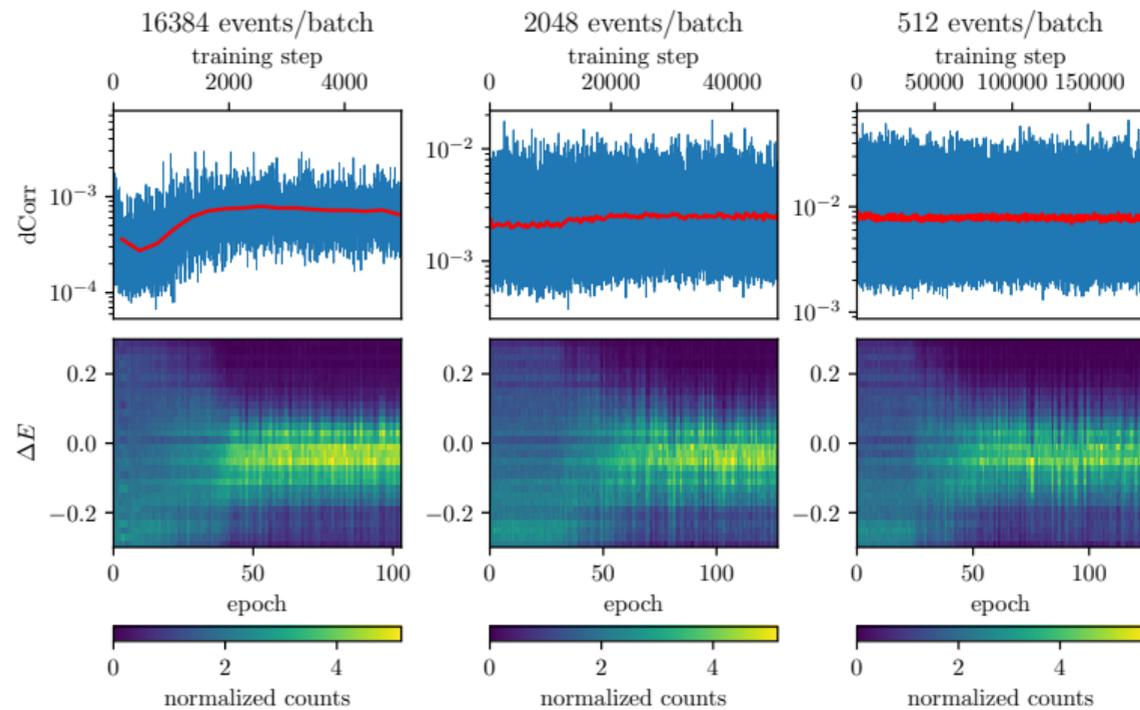
## Data Samples for Training



# Monitoring DNN Training

## Coincidence of dCorr Increase and Sculpting

- Very large batch sizes required for numerical stability
- Clear coincidence of start of sculpting and dCorr increase (if observable)



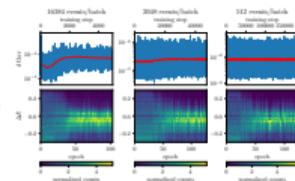
2023-12-19

CS with NNs for Belle II

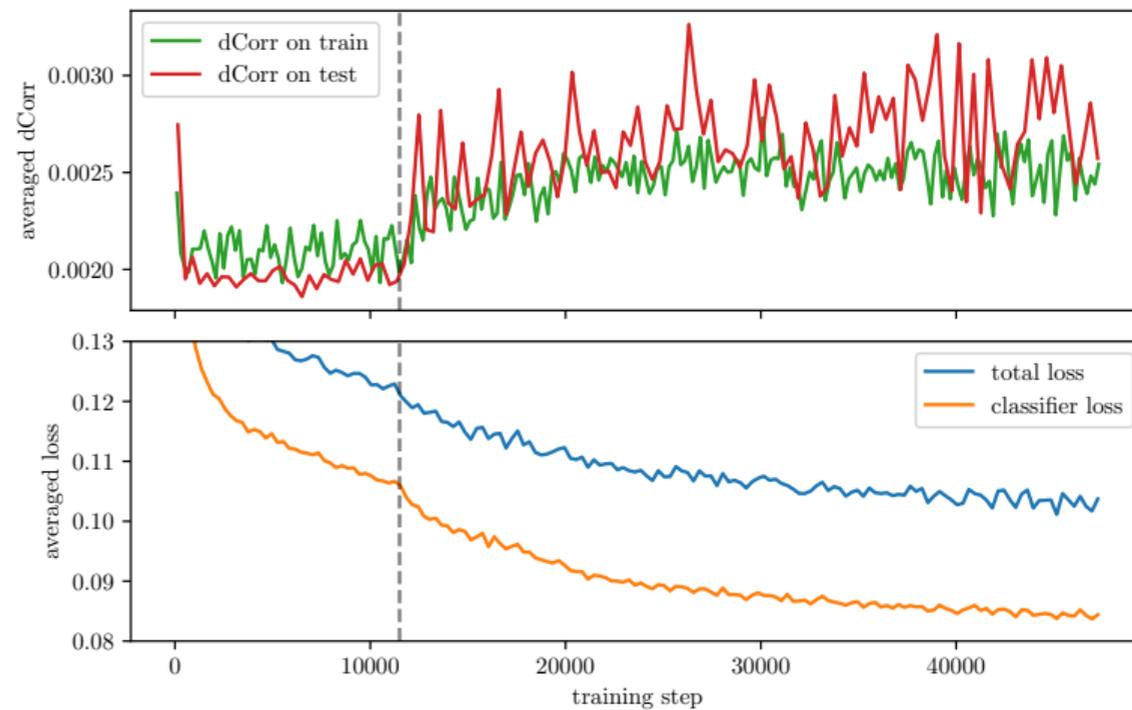
Monitoring DNN Training

Monitoring DNN Training  
Coincidence of dCorr Increase and Sculpting

- Very large batch sizes required for numerical stability
- Clear coincidence of start of sculpting and dCorr increase (if observable)

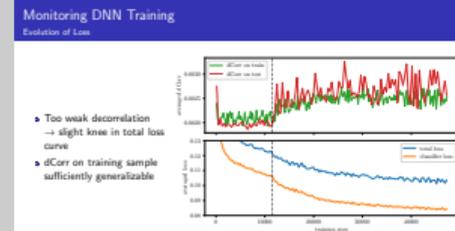


- Too weak decorrelation  
→ slight knee in total loss curve
- dCorr on training sample sufficiently generalizable

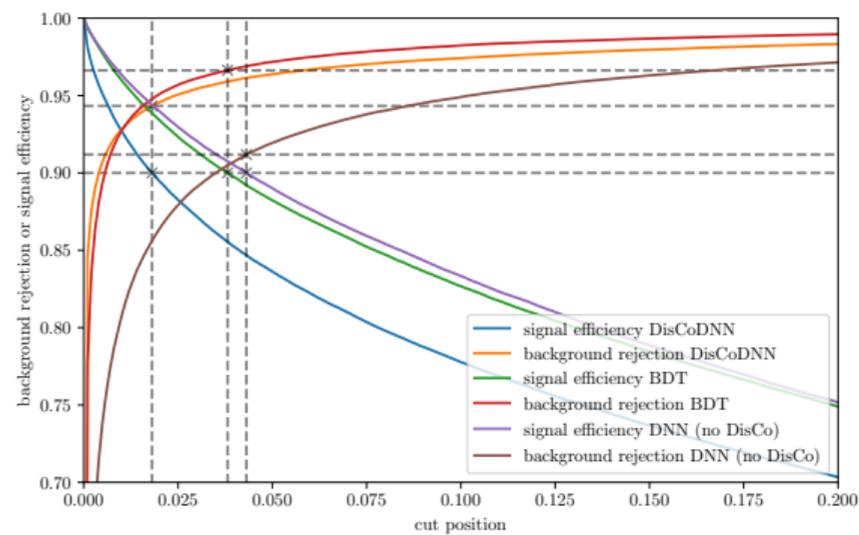
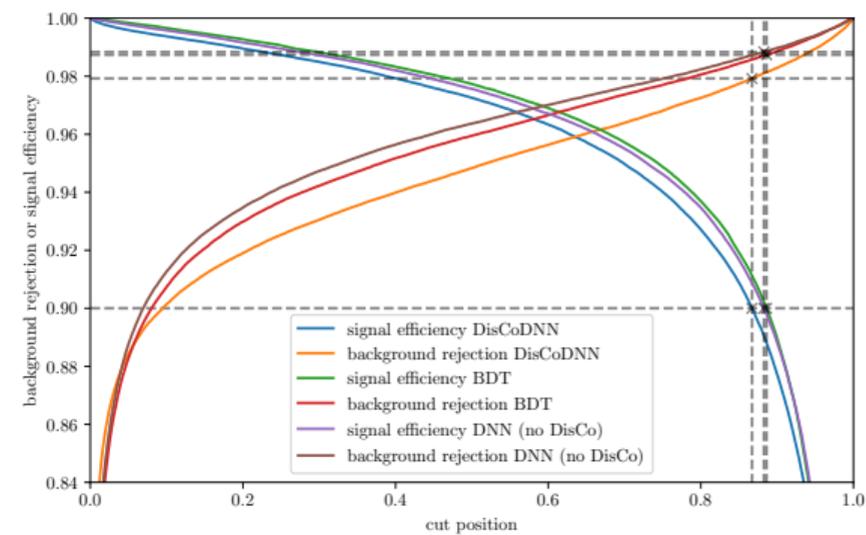


### Monitoring DNN Training

1. Talk about intuition of *barrier* in parameter space. DisCo appear to introduce barrier but never really plane the global (correlated) minimum.



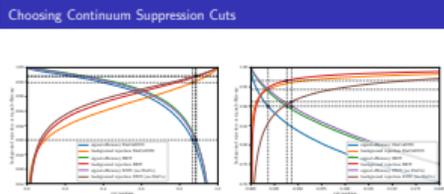
# Choosing Continuum Suppression Cuts



2023-12-19

CS with NNs for Belle II

Choosing Continuum Suppression Cuts



	signal	$q\bar{q}$	$B\bar{B}$
true yield DisCoDNN	318	3313	71
true yield BDT	321	2134	75
yield DisCoDNN	$310.6 \pm 28.3$	$3343 \pm 39$	$49.30 \pm 31.28$
yield BDT	$337.5 \pm 26.1$	$2149 \pm 35$	$43.52 \pm 27.83$
rel. fit error DisCoDNN in %	8.902	1.178	44.06
rel. fit error BDT in %	8.144	1.626	37.1
rel. true error DisCoDNN in %	$2.335 \pm 8.902$	$0.897 \pm 1.178$	$30.57 \pm 44.06$
rel. true error BDT in %	$5.133 \pm 8.144$	$0.710 \pm 1.626$	$41.97 \pm 37.10$
pull DisCoDNN in $\sigma$	-0.2623	0.7619	-0.6937
pull BDT in $\sigma$	0.6302	0.4367	-1.131

└ Fits on MC

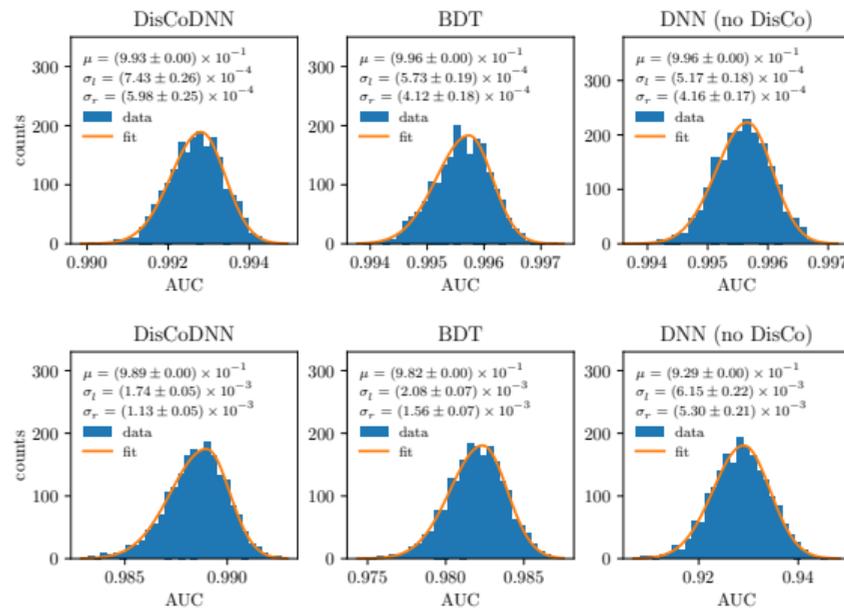
	signal	$q\bar{q}$	$B\bar{B}$
true yield DisCoDNN	318	3313	71
true yield BDT	321	2134	75
yield DisCoDNN	$310.6 \pm 28.3$	$3343 \pm 39$	$49.30 \pm 31.28$
yield BDT	$337.5 \pm 26.1$	$2149 \pm 35$	$43.52 \pm 27.83$
rel. fit error DisCoDNN in %	8.902	1.178	44.06
rel. fit error BDT in %	8.144	1.626	37.1
rel. true error DisCoDNN in %	$2.335 \pm 8.902$	$0.897 \pm 1.178$	$30.57 \pm 44.06$
rel. true error BDT in %	$5.133 \pm 8.144$	$0.710 \pm 1.626$	$41.97 \pm 37.10$
pull DisCoDNN in $\sigma$	-0.2623	0.7619	-0.6937
pull BDT in $\sigma$	0.6302	0.4367	-1.131

## Bootstrapping

- Models fluctuations of occurrences of event types, not numerical fluctuations
- All classifiers remain reasonably stable

## Uncorrelated Toys

- Do not model correlations, as nearly impossible
- Classifiers that do not significantly sculpt  $\Delta E$  barely utilize correlations between input variables



## Classifier Performance Stability & Input Variable Correlations

