# The use of adversaries of optimal Neural Net training
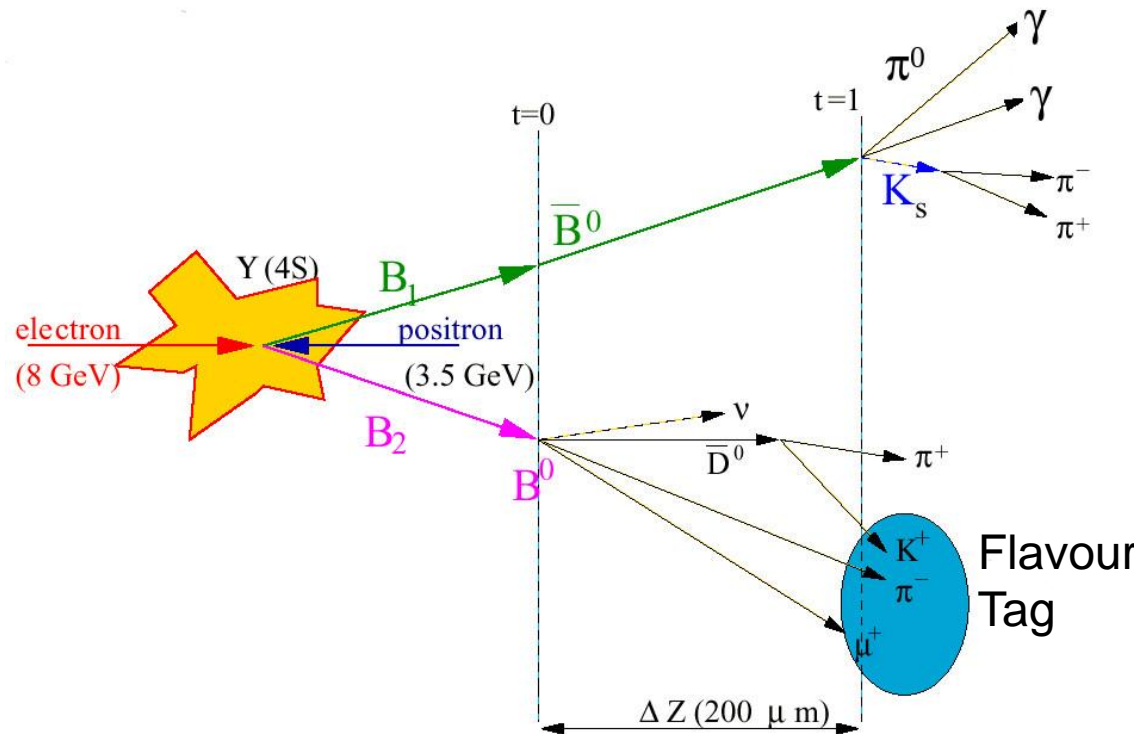


(https://travellingbuzz.com)

Anton Hawthorne and Martin Sevior
University of Melbourne and Belle(II) collaboration
Computing in High Energy Physics
Sophia, Bulgaria, July, 2018

# Introduction

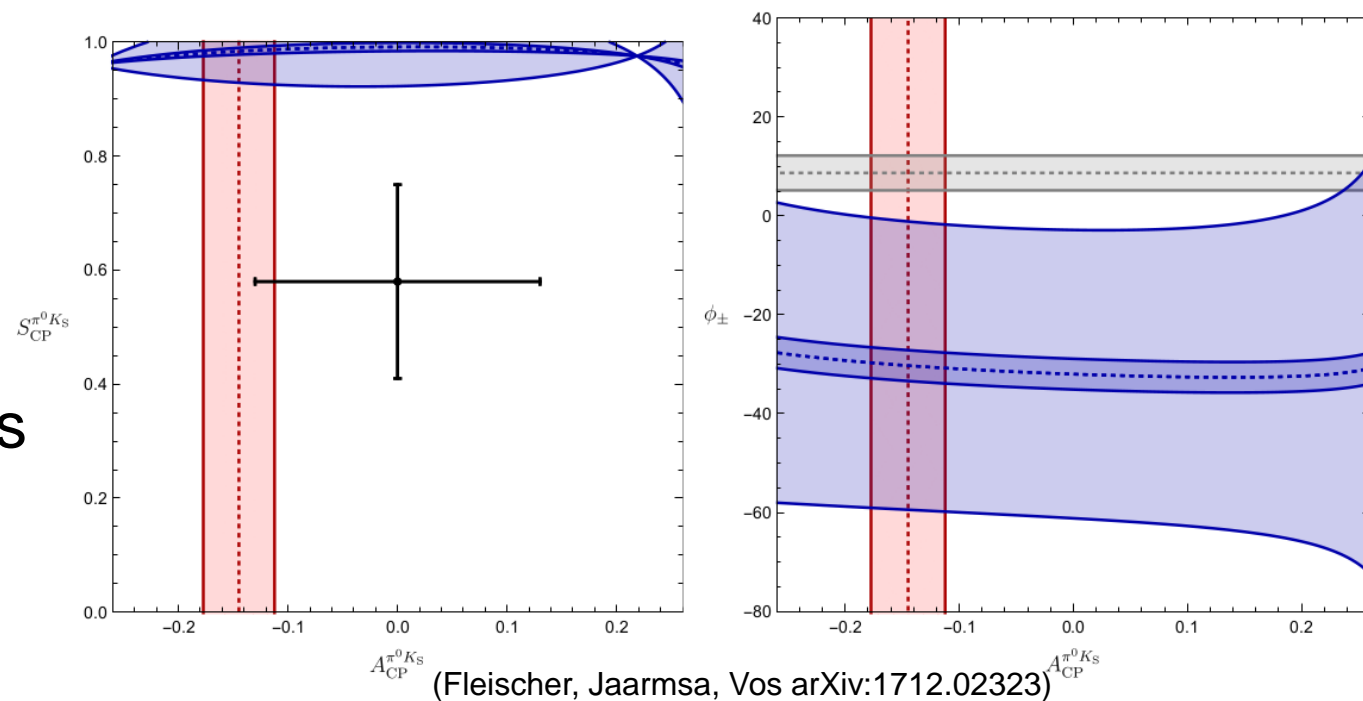This presentation summarizes the M.Phil.of Anton Hawthorne
Detailed write up in arXiv:1712.07790

- Belle
- CP-violation in $\overline{B^0} \to K_s\pi^0$ decays
- Analysis and Backgrounds to $\overline{B^0} \to K_s\pi^0$
- Deep Neural Network vs Shallow and BDT
- Data normalisation
- Performance of Deep Net
- Background sculpting
- Adversarial Neural Net
- Performance of Adversarial Net
- Conclusions

# CP-violation in $\overline{B^0} \to K_S \pi^0$ decays

Define: $A(t) = \dfrac{N\left(\overline{B^0} \to K_S \pi^0\right)(t) - N(B^0 \to K_S \pi^0)(t)}{N\left(\overline{B^0} \to K_S \pi^0\right)(t) + N(B^0 \to K_S \pi^0)(t)}$

$$A(t) = S_{CP}\sin(\Delta m t) + A_{CP}\cos(\Delta m t)$$



- SM prediction 2.2σ from measurements
- Measurements are statistically limited

(Fleischer, Jaarmsa, Vos arXiv:1712.02323)

# Kinematic Variables in B-Factory measurements

$$M_{\mathrm{bc}} = \sqrt{E_{beam}^{*2} - p_B^{*2}}$$

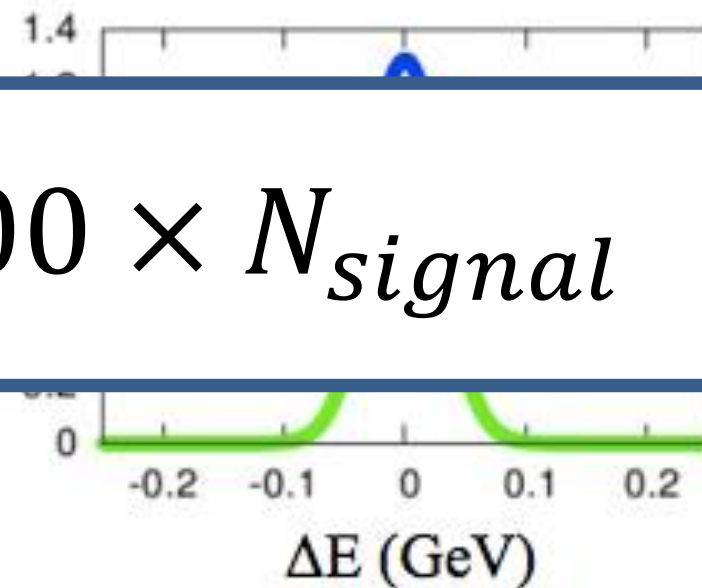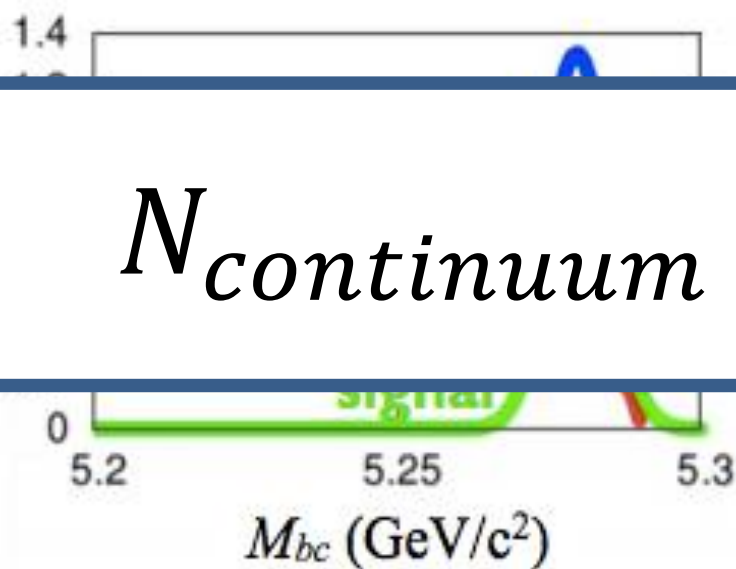$$\Delta E = E_B^* - E_{beam}^*$$



$M_{bc}$ peaks at B mass for fully reconstructed signal
$\Delta E$ peaks at zero for fully reconstructed signal

# Kinematic Variables in B-Factory measurements

$$M_{\mathrm{bc}} = \sqrt{E_{beam}^{*2} - p_B^{*2}} \qquad \Delta E = E_B^* - E_{beam}^*$$
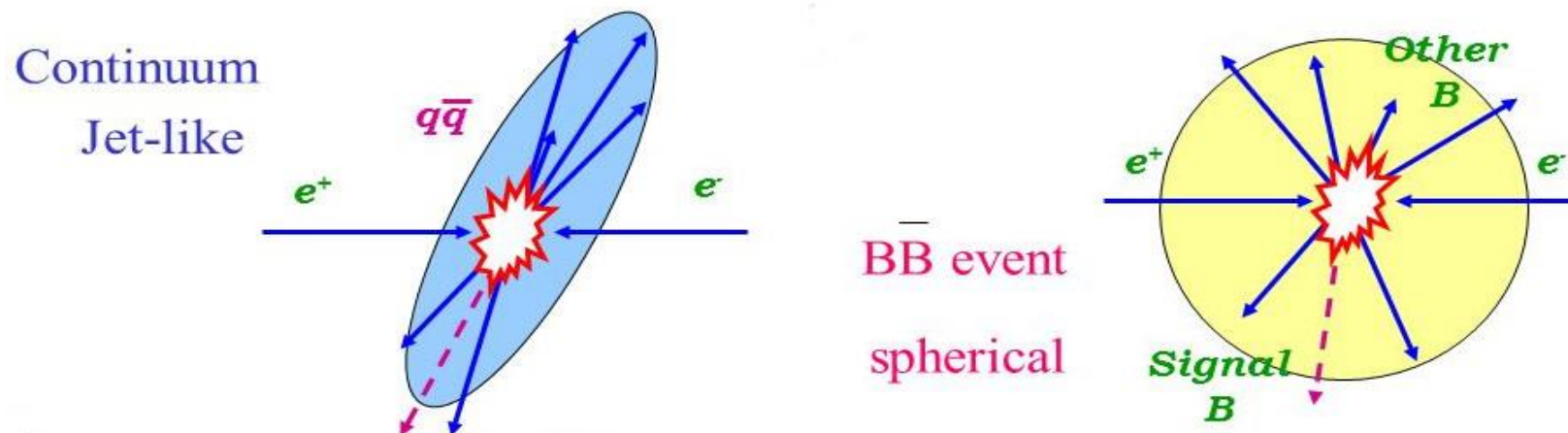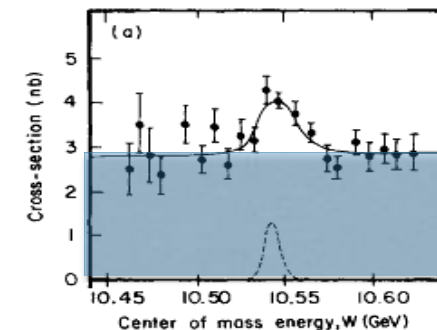
$$N_{continuum} = 400 \times N_{signal}$$

$M_{bc}$ peaks at B mass for fully reconstructed signal
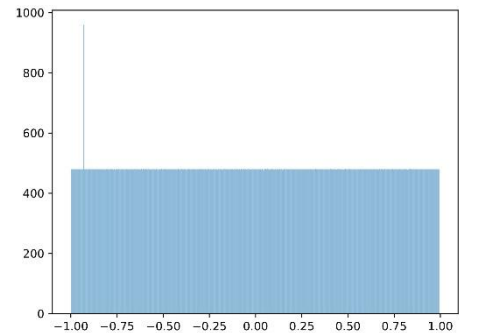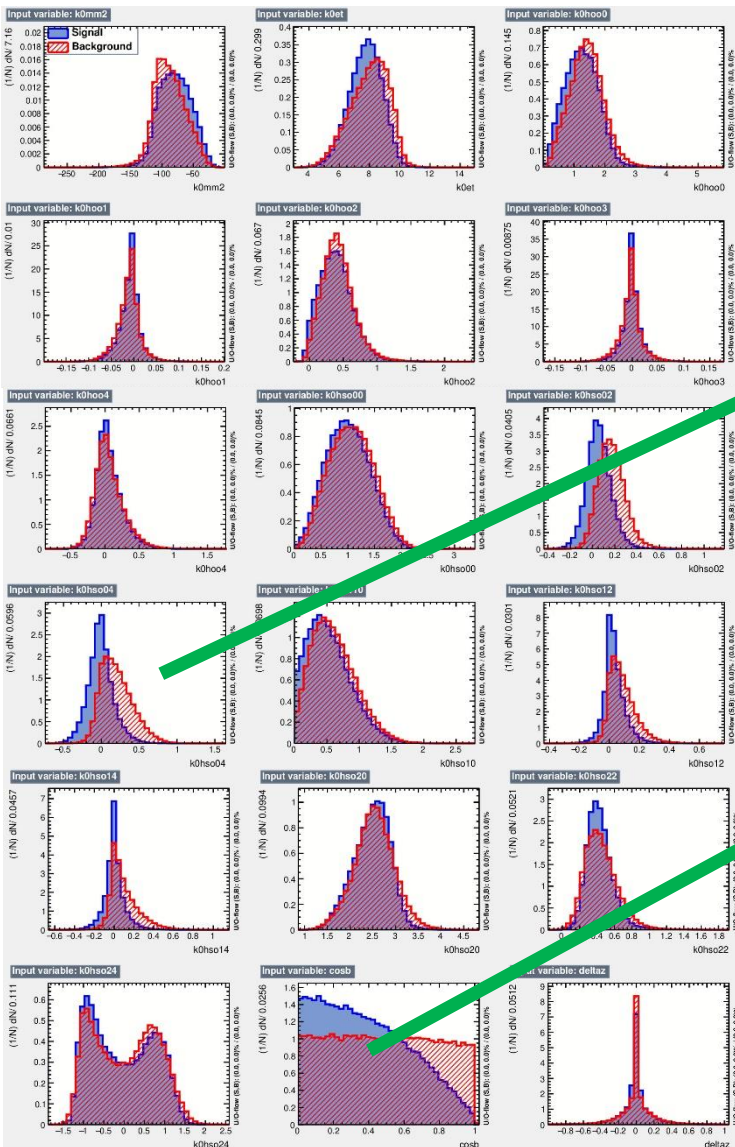$\Delta E$ peaks at zero for fully reconstructed signal

# Continuum Background

- Continuum background $e^+e^- \rightarrow q\bar{q}(u,d,s,c)$
  - Dominant background
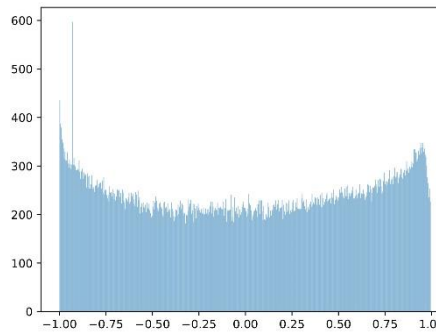  - Event topology differs from BB decays





  - Combine variables describing the event topology in a Multi-Variate analysis.
  - Investigate a Deep Neural Net for improved performance
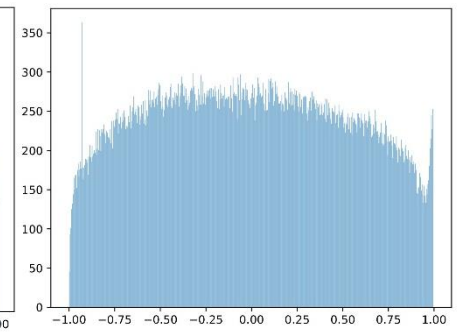
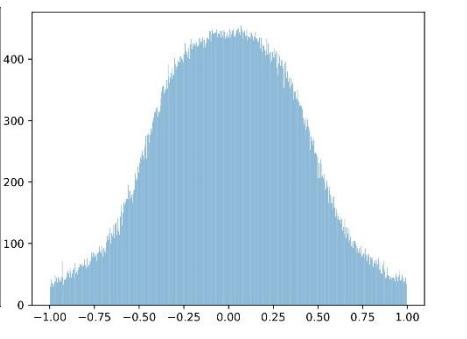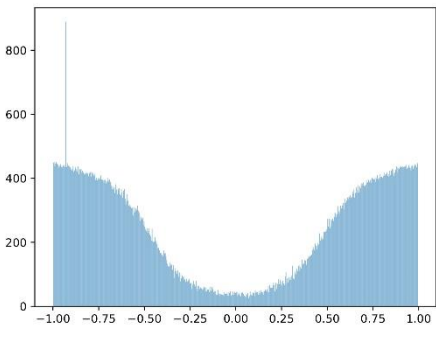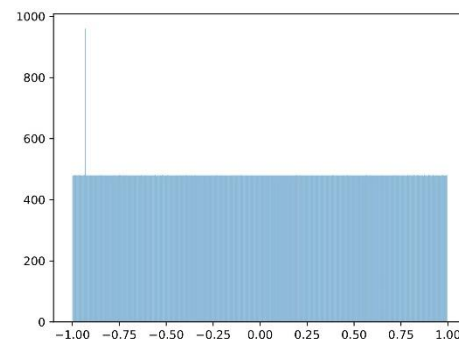# Normalization of continuum fighting variables
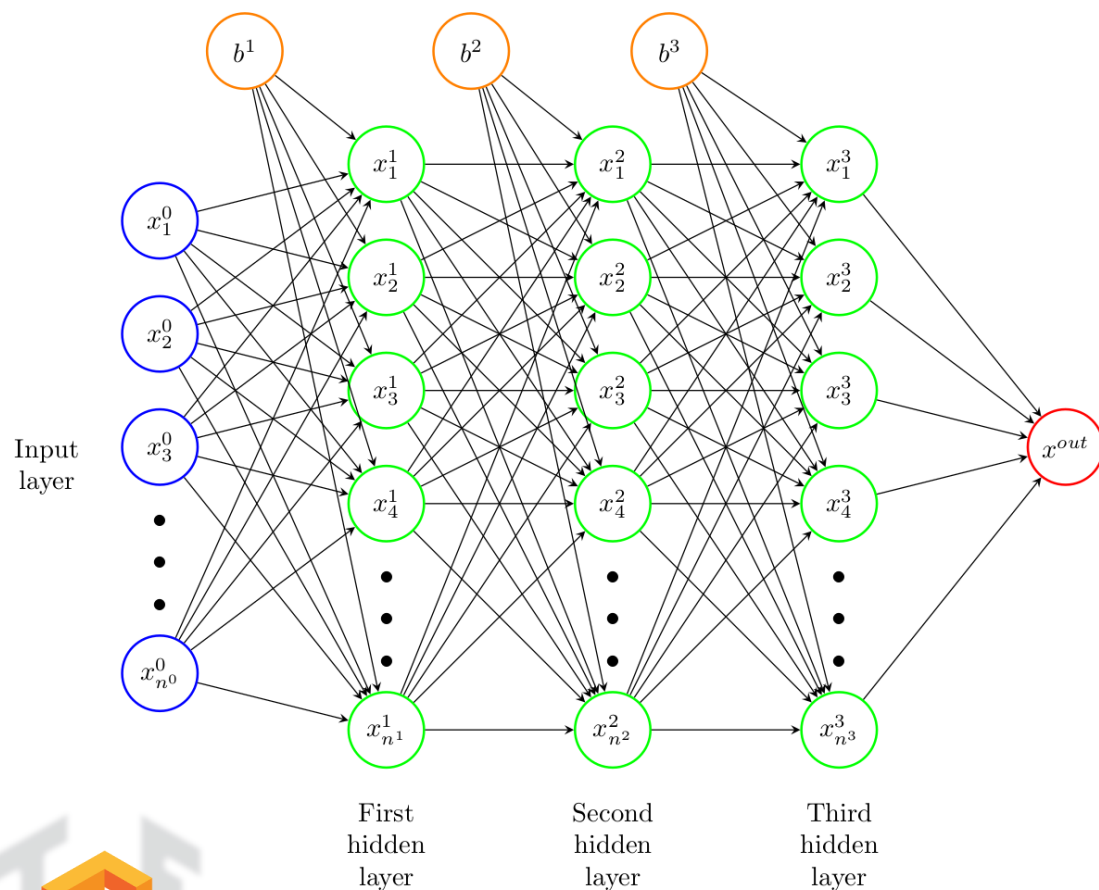


Signal and continuum          continuum          Signal

Signal

continuum

## Equal frequency binning used to map into range -1.0 to +1.0
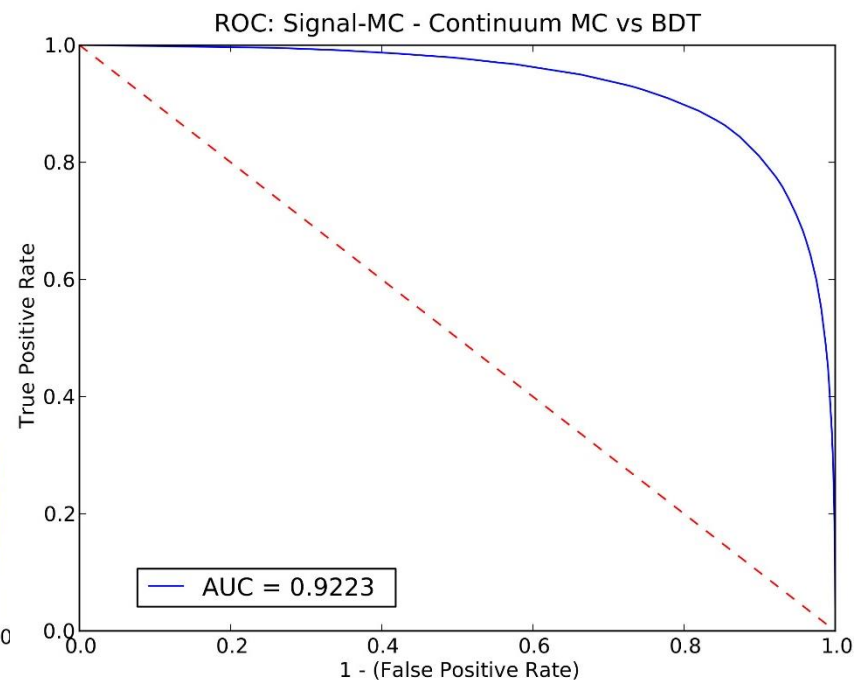
# Implementation of Deep Neural Network (TensorFlow)



- Deep Neural Net built from the ground up in TensorFlow
- Employed Hyperband to search for best hyper-parameters

- Trained with 125000 signal and continuum events
- Validated with 125000 signal and continuum events
- Tested with 62500 signal and continuum events
- Employed ADAM algorithm for training
- $L_{class}(\vec{x}, \hat{y}) = -\hat{y} \cdot \log(y(\vec{x})) - (1 - \hat{y}) \cdot \log(1 - y(\vec{x}))$

- A maximum number of epochs 600.
- 50 events per batch
- Learning rate of 0.0001.
- Six hidden layers.
- 47 nodes per hidden layer.
- Exponential linear unit activation function.

# Performance of Deep Neural Network



TensorFlow
Deep Neural Network
AUC = 0.9501
70% sig. eff. => 3.3% Continuum

Neurobayes
Shallow Neural Network
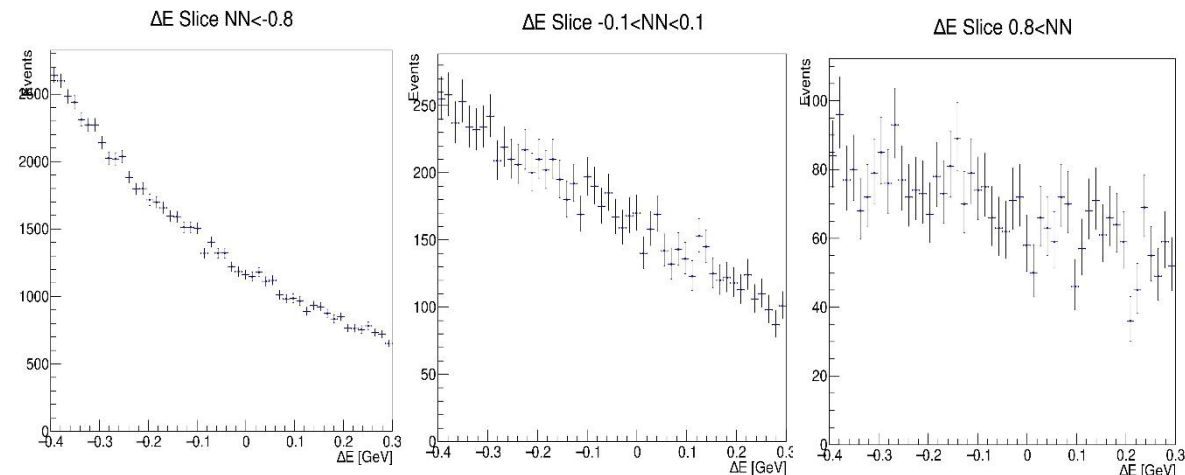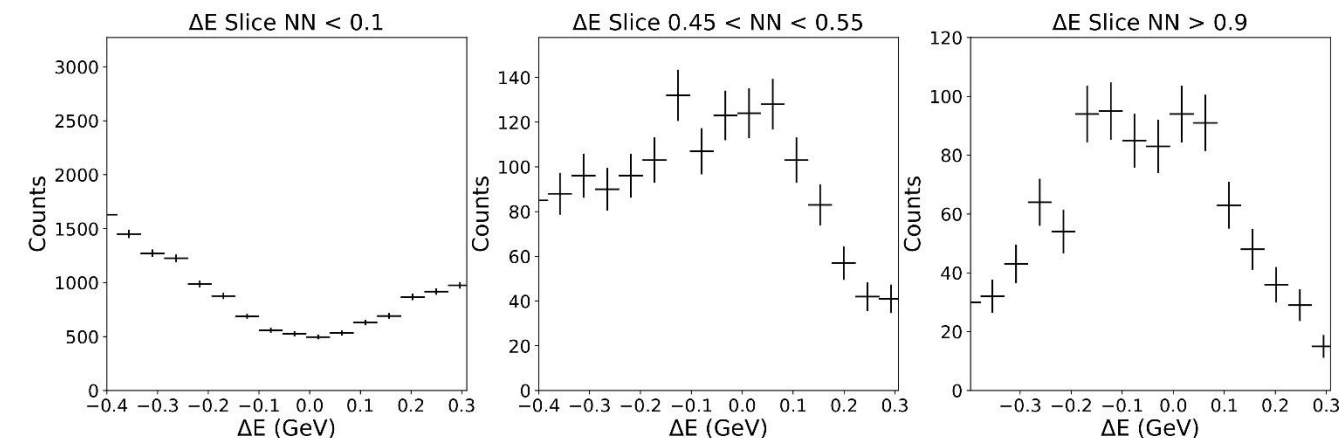AUC = 0.912
70% sig. eff => 7.6% continuum

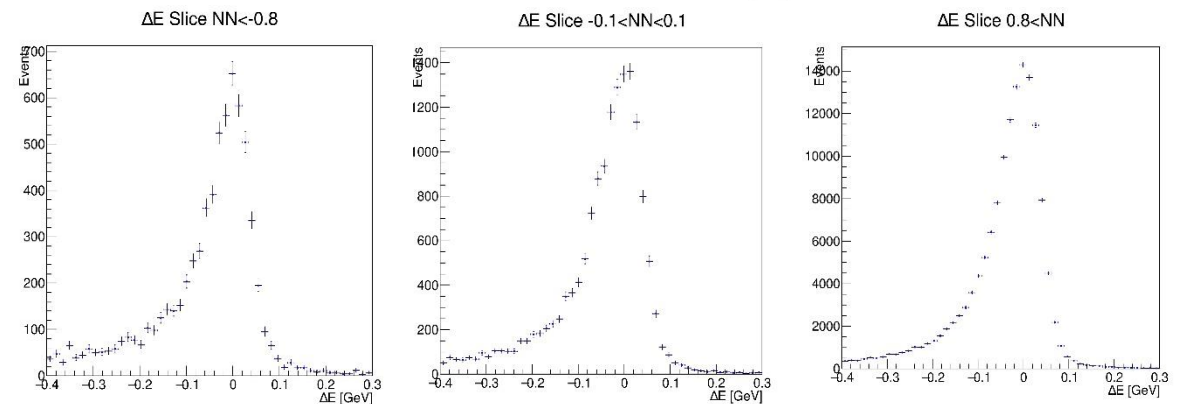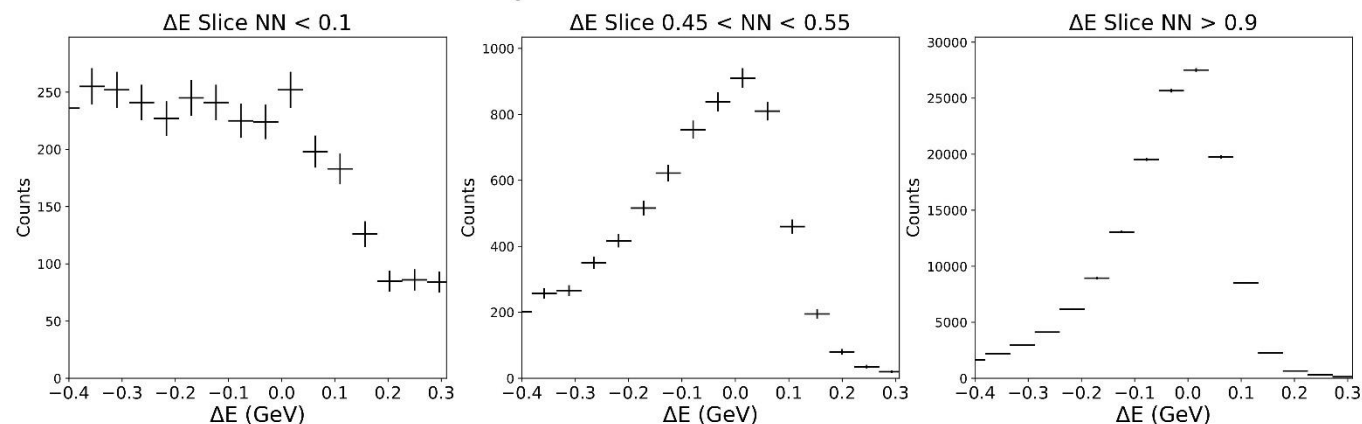TMVA
Boosted Decision Trees
AUC = 0.922
70% sig. eff => 7.0% continuum

# Sculpting in $\Delta E$ distribution



Deep Neural Net

NeuroBayes

# Sculpting in $\Delta E$ distribution



Off-resonance real-data
Deep Neural Net

Boosted Decision Trees

# Correlations between KSFW and $\Delta E$



Continuum - Scatter Plot of $\Delta E$ and $R_{20}^{so}$

Continuum - Scatter Plot of $\Delta E$ and $R_0^{oo}$

Continuum - Scatter Plot of $\Delta E$ and $R_2^{oo}$

Continuum - Scatter Plot of $\Delta E$ and $R_{22}^{so}$

| Classifier | AUC | Correl. |
|---|---|---|
| Deep Net All | 0.950 | 0.114 |
| Deep Net 1 removed | 0.938 | 0.073 |
| Deep Net 2 removed | 0.928 | 0.083 |
| Deep Net 3 removed | 0.923 | 0.056 |
| Deep Net 4 removed | 0.918 | 0.062 |
| NeuroBayes All | 0.912 | 0.058 |
| NeuroBayes (reduced) | 0.902 | -0.001 |
| BDT All | 0.922 | 0.262 |
| BDT (reduced) | 0.913 | 0.054 |

"reduced" means all four correlated variables removed

# Adversarial Neural Network

Build an Adversarial Neural Net to keep the correlated variables but remove the sculpting.



$$L_{adv}(NN, \Delta E) = -\log\left(\sum_{i=1}^{5} \frac{f_i'(NN)}{\sqrt{2\pi\sigma_i^2(NN)}} e^{\frac{-(\mu_i - \Delta E)^2}{2\pi\sigma_i^2(NN)}}\right)$$

$$L_{tot} = L_{class} - \lambda_{adv}L_{adv}$$

*Learning to Pivot with Adversarial Networks.* Gilles Louppe, Michael Kagan, and Kyle Cranmer. (arXiv:1611.01046)

# Performance of Adversarial Neural Network
## Adjust $\lambda_{adv}$



Anton Hawthorne and Martin Sevior, CHEP 2018, Sophia, Bulgaria

# Performance of Adversarial Neural Network

Background rejection for 92.5% signal acceptance in full and signal ($-0.1\ GeV < \Delta E < 0.1\ GeV$) $\Delta E$ regions

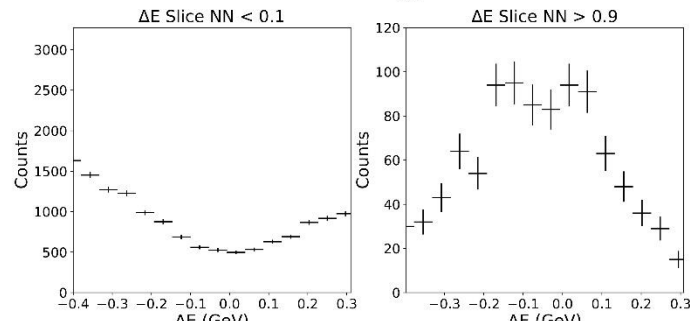| $\lambda_{adv}$ | Bck. Rej. Full-$\Delta E$ | Bck. Rej Signal-$\Delta E$ | Corr with $\Delta E$ | |
|---|---|---|---|---|
| 0.0 | 81.3% | 63.2% | 0.114 | |
| 0.25 | 80.6% | 65.1% | 0.080 | |
| 0.25 | 79.7% | 66.6% | 0.057 | |
| 0.75 | 80.0% | 67.6% | 0.025 | |
| 1.0 | 78.4% | 67.6% | 0.012 | |
| 1.5 | 74.6% | 67.4% | -0.025 | Working point |
| 2.0 | 70.2% | 66.6% | -0.057 | |
| 3.0 | 64.3% | 63.0% | -0.106 | |
| 4.0 | 57.1% | 57.6% | -0.145 | |
| 5.0 | 51.9% | 53.4% | -0.175 | |
| NeuroBayes | 66.7% | 64.1% | 0.058 | |
| NeuroBayes (Reduced) | 63.0% | 64.1% | -0.001 | |
| BDT | 73.7% | 67.2% | 0.262 | "reduced" means all four correlated variables removed |
| BDT (Reduced) | 66.9% | 66.5% | 0.054 | |

# Adversarial Neural Network $\lambda_{adv} = 1.5$



Signal $\Delta E$ Distributions for $\lambda_{adv} = 1.5$

Continuum $\Delta E$ Distributions for $\lambda_{adv} = 1.5$

Off resonance $\Delta E$ Distributions for $\lambda_{adv} = 1.5$

Signal, $\lambda_{adv} = 1.5$

Continuum MC, $\lambda_{adv} = 1.5$

Off-resonance, $\lambda_{adv} = 1.5$

# Conclusions

- Deep Neural Nets discovered and exploited a subtle correlation among the KSFW moments
- This sculpted the background $\Delta E$ distribution to resemble signal
- Removing the correlated discriminating variables reduces the effectiveness of the classification
- Built an adversarial neural net to counter-act this (negative feedback)
- Hyperparameter $\lambda_{adv}$ adjusts the strength of the negative feedback
- Increasing $\lambda_{adv}$ decreases the sculpting and correlation with $\Delta E$
- Too large $\lambda_{adv}$ causes negative correlation with $\Delta E$
- In optimal region the Deep Net with $\lambda_{adv}$ = 1.5 has the best background rejection with smallest correlation
- Still some sculpting $\lambda_{adv}$ = 1.5 but $\Delta E$ continuum distribution still significantly different from Signal
- Off-resonance data validates the MC distributions predicted with $\lambda_{adv}$ = 1.5

# Backup

Anton Hawthorne and Martin Sevior, CHEP 2018, Sophia, Bulgaria

# KEKB and Belle



KEKB maximum Luminosity $2.1\text{x}10^{34}\text{cm}^2\text{s}^{-1}$ => 21 B-pairs/sec
SuperKEKB $\rightarrow 8\text{x}10^{35}\text{cm}^2\text{s}^{-1}$ => 800 B-pairs/sec (Currently $2.1\text{x}10^{33}$ cm$^2$s$^{-1}$)

# The Belle experiment



## Integrated luminosity of B factories



> 1 ab$^{-1}$
**On resonance:**
$\Upsilon(5S)$: 121 fb$^{-1}$
$\Upsilon(4S)$: 711 fb$^{-1}$
$\Upsilon(3S)$: 3 fb$^{-1}$
$\Upsilon(2S)$: 25 fb$^{-1}$
$\Upsilon(1S)$: 6 fb$^{-1}$
**Off reson./scan:**
~ 100 fb$^{-1}$

~ 550 fb$^{-1}$
**On resonance:**
$\Upsilon(4S)$: 433 fb$^{-1}$
$\Upsilon(3S)$: 30 fb$^{-1}$
$\Upsilon(2S)$: 14 fb$^{-1}$
**Off resonance:**
~ 54 fb$^{-1}$

# Hyperparameters of Adversarial Neural Network

- 100 training steps.
- 125 events per batch.
-  A Learning rate of 0.01.
- Two hidden layers.
- 20 nodes per hidden layer.
- Exponential linear unit activation function for the nodes in the hidden layer.
- 15 output nodes (three output nodes corresponding to each Gaussian):

    – 5 output nodes corresponding to $\mu_i$ - no activation function (identity operator).
    – 5 output nodes corresponding to un-normalised fractions $f_i$ - no activation function (identity operator).
    – 5 output nodes corresponding to $\sigma_i$, where the 'activation' is the exponential function, to ensure that the widths of the Gaussians are positive.

# Training the Adversarial Neural Network

1. Train the NN to optimally separate signal and continuum. (TF1)
2. Create the ANN, and the classifying (the original) NN with the same architecture
3. For every 20,000 steps and a given choice of $\lambda_{adv}$ :
   (a) Train the ANN for the given number of adversary training steps (100 steps), where
       (i) For every event in the batch, get the NN output from the classifier.
       (ii) Using NN and $\Delta$E get the adversarial loss given by $L_{class}$ .
       (iii) Train the ANN given the adversarial loss and adversarial learning rate
   (b) Train the classifier for one training step, with the loss function given by $L_{tot}$
       (dependence on $\Delta$E, as well as NN)
4. This is the ANN-corrected Neural Net

# Adversarial Neural Network

## Correlation as training proceeds for $\lambda_{adv}$=0.5



NN DeltaE Correlation vs Step

- 4 epochs, of 5000 classifier-training steps
- adversarial network is trained for 125 steps per classifier training step.
- correlations are in the validation data sets, and calculated over the entire range 0 < NN < 1.