

# Application of Bayesian Neural Networks to clustering for LUXE experiment

Roman Urmanov

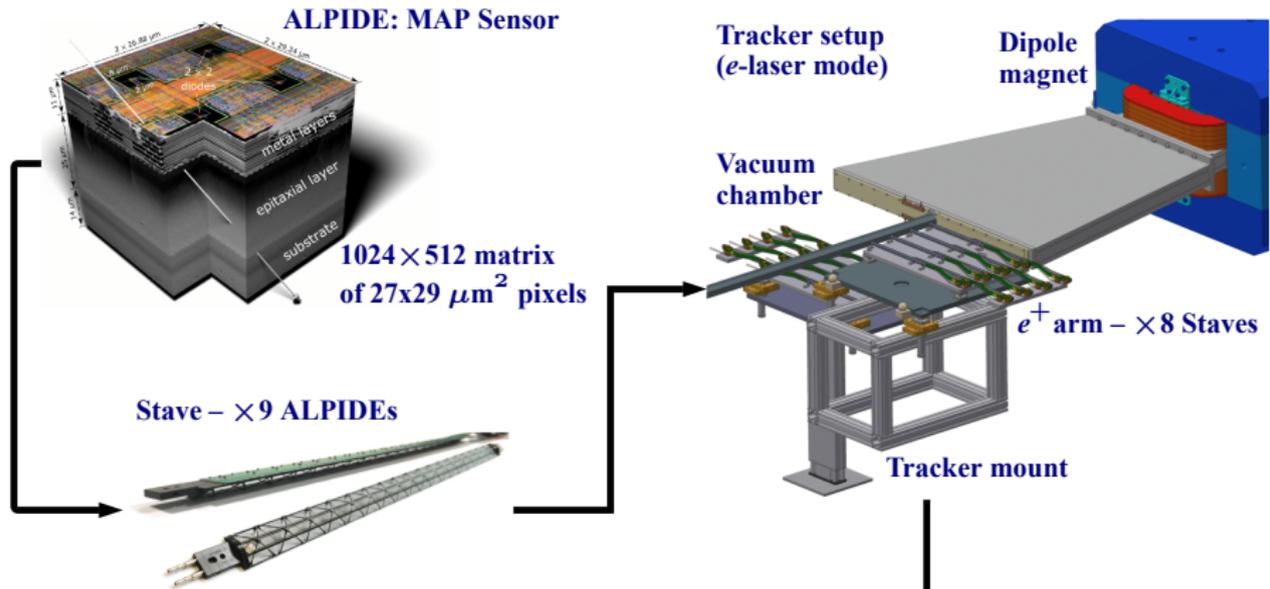
05.06.2023

WEIZMANN  
INSTITUTE  
OF SCIENCE

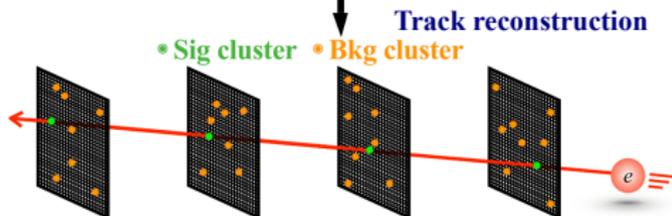
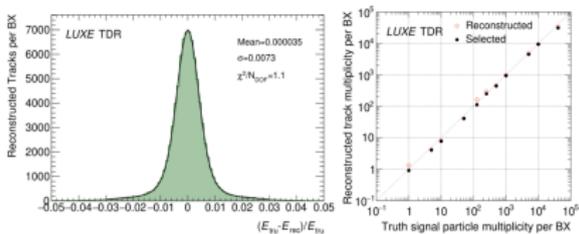


LUXE

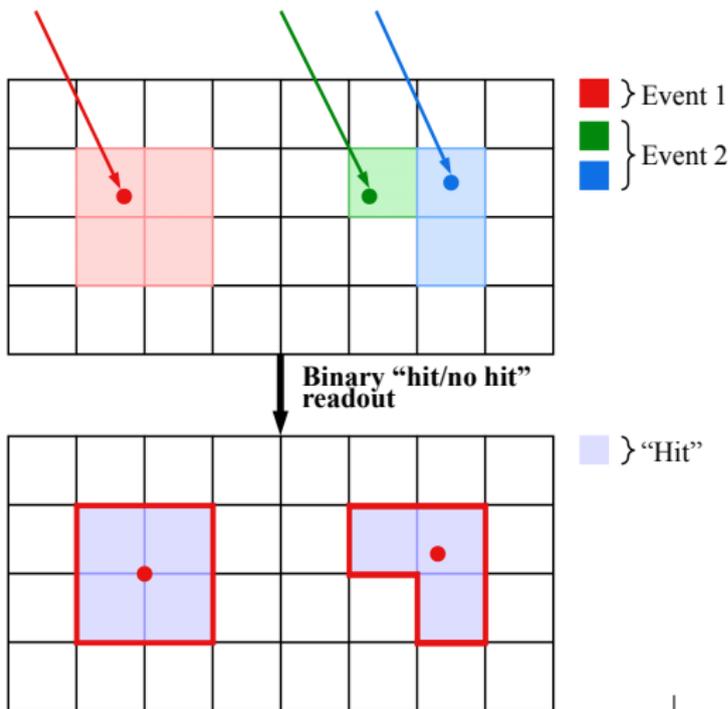
# Introduction: Tracker



## Physics analysis



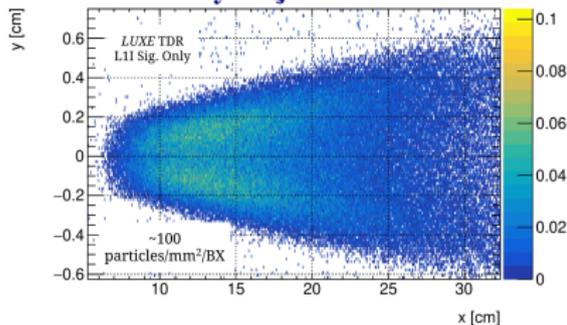
# Clustering: ML motivation



**"Pac-Man": adjacent hits combined into clusters. Cluster's "center of mass" is taken for track inference**

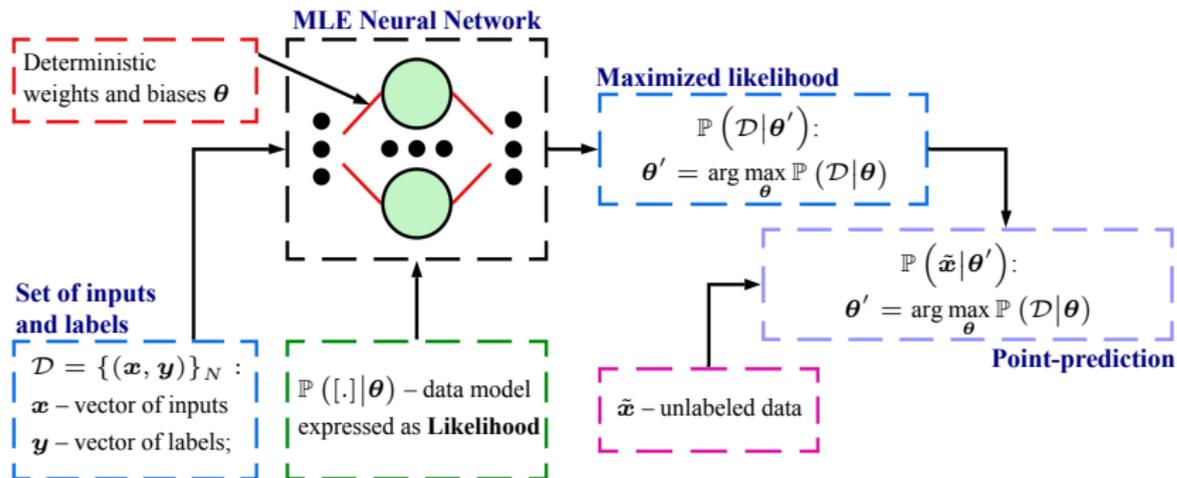
- Binary output – no information on number of particles
- LUXE will work in a wide range of pixel occupancies
- "Pac-Man" – up to  $\lesssim 10^4$   $e$  per BX
- Up to  $3 \times 10^5$   $e$  per BX in later stages

## Hit density at $\xi=7$

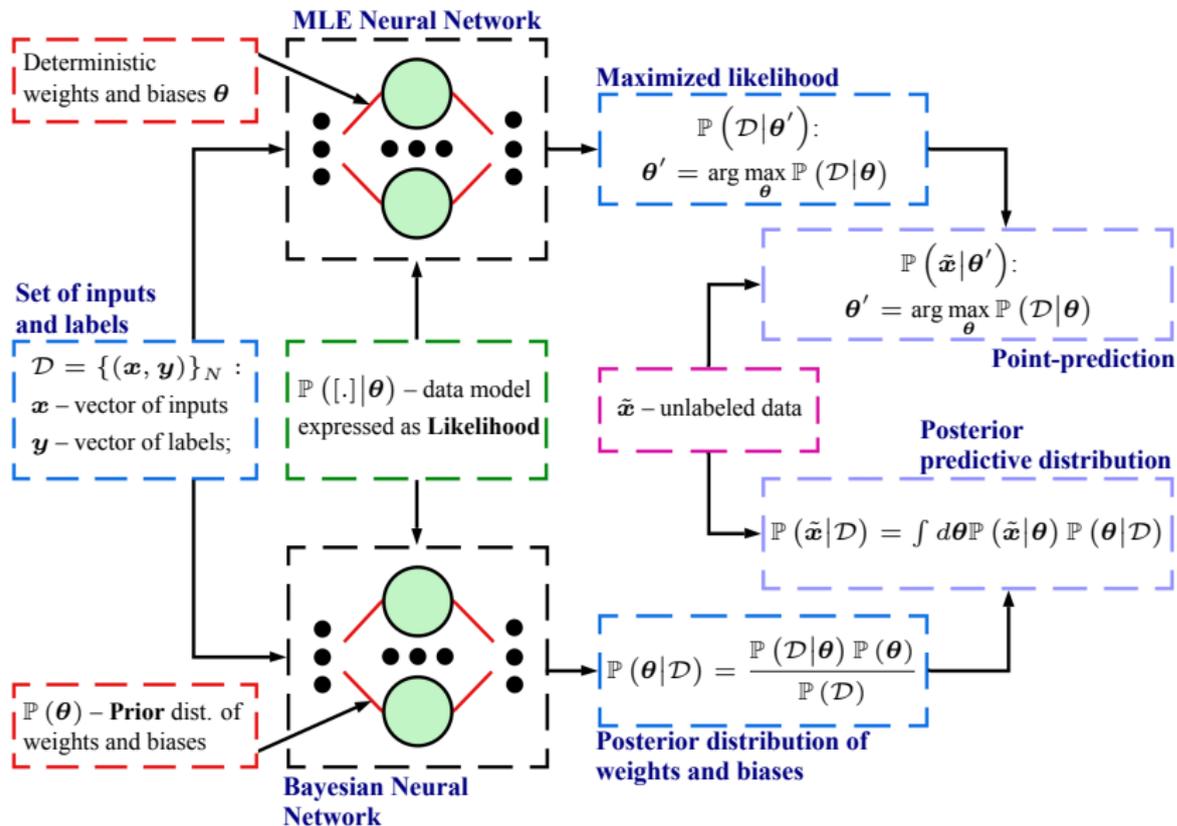


Detector reference	Hit density [mm <sup>-2</sup> ]		
	MCD	ATLAS ITk	ALICE ITS3
Pixel Layer 0	3.68	0.643	0.85
Pixel Layer 1	0.51	0.022	0.51

# Bayesian Neural Network: General structure

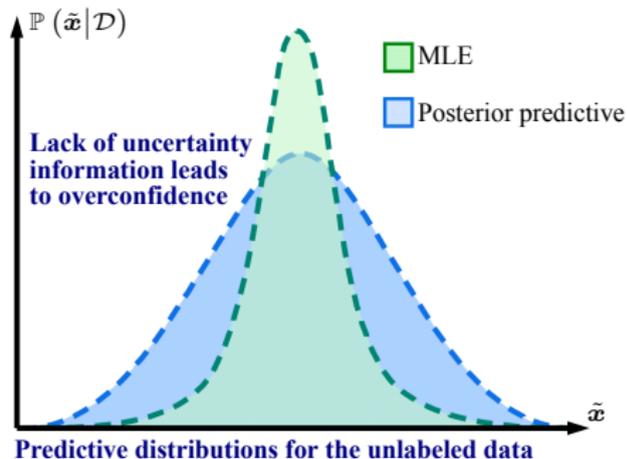
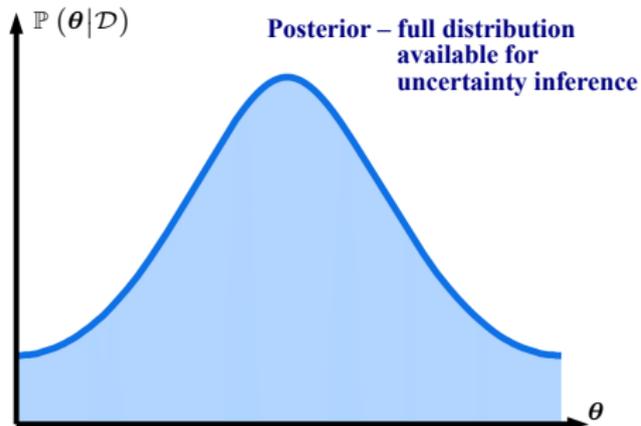
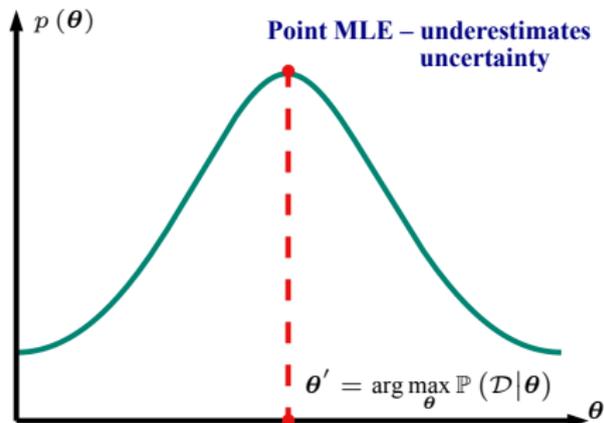


# Bayesian Neural Network: General structure

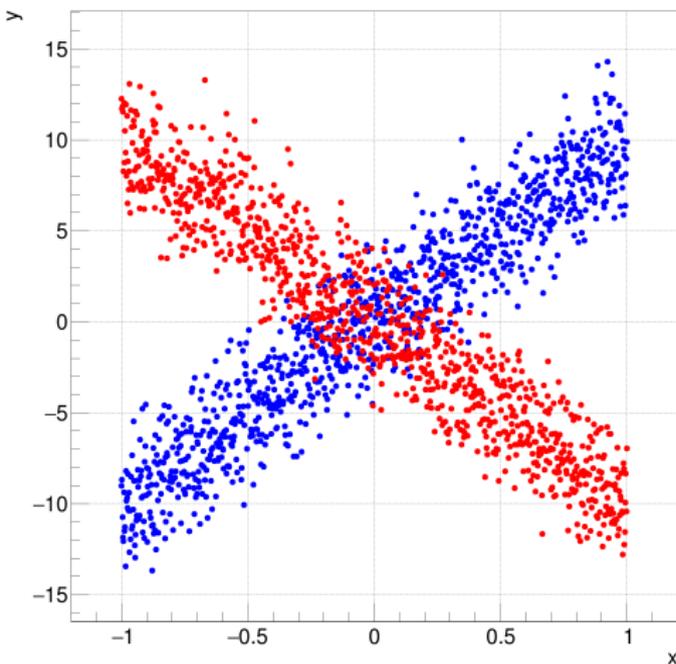


- Bayesian Neural Network – Point-Prediction network with stochastic weights and biases

# Bayesian Neural Networks: Motivation



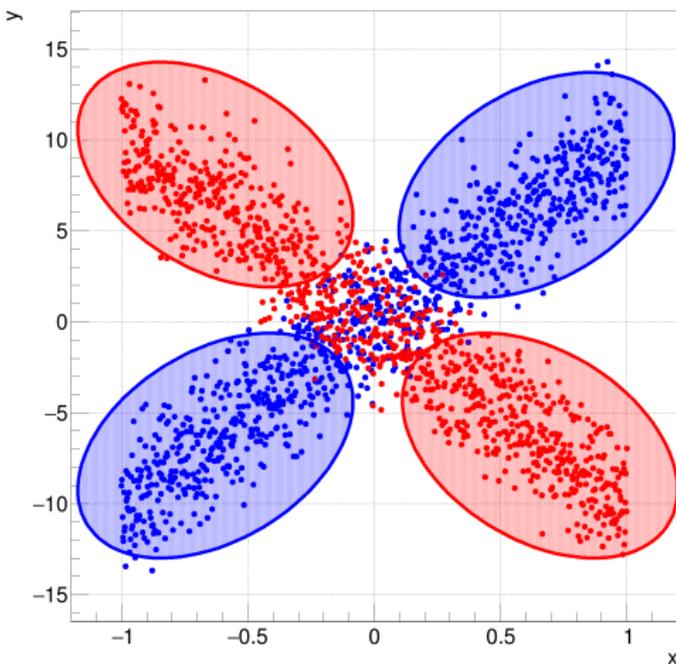
# Bayesian Neural Networks: Motivation



- Classification task: based on the  $(x, y)$  position assign one of the two classes to an event

**Example of a classification problem with a significant class overlap**

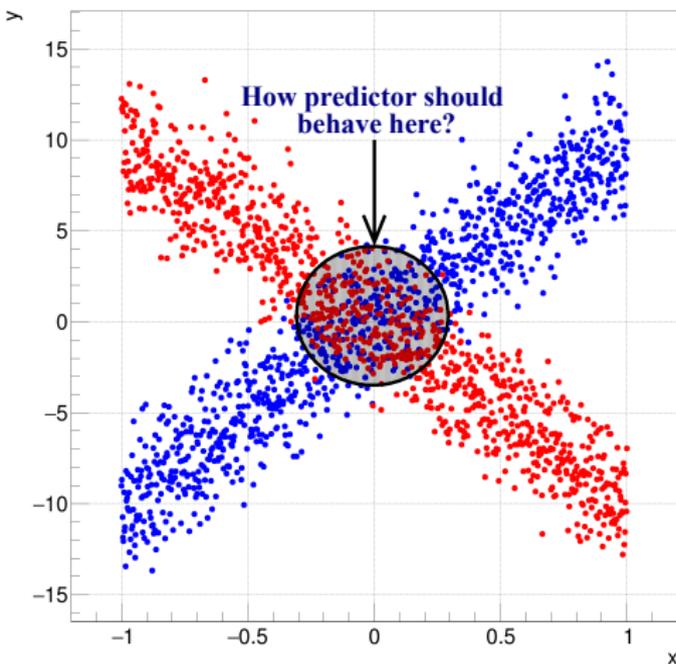
# Bayesian Neural Networks: Motivation



**Example of a classification problem with a significant class overlap**

- Classification task: based on the  $(x, y)$  position assign one of the two classes to an event
- It is clear what we want from a predictor in regions with no overlap

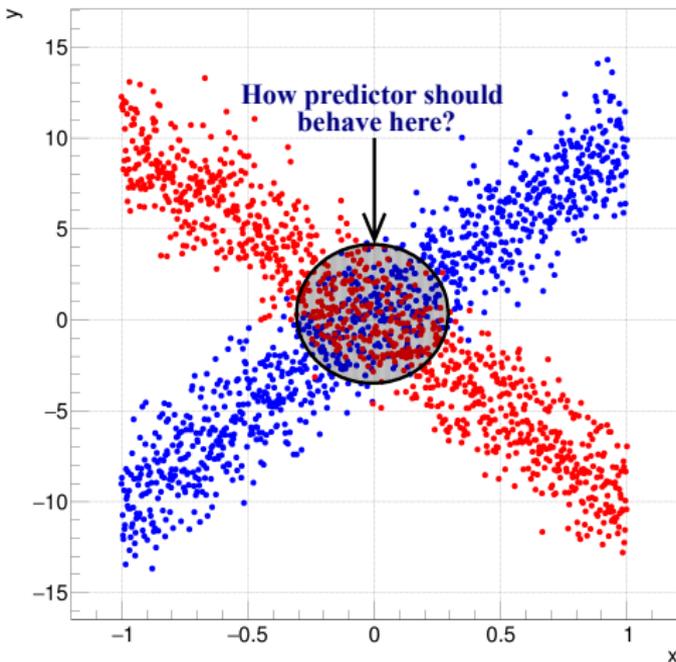
# Bayesian Neural Networks: Motivation



**Example of a classification problem with a significant class overlap**

- Classification task: based on the  $(x, y)$  position assign one of the two classes to an event
- It is clear what we want from a predictor in regions with no overlap
- What should the predictor do in the overlap region?
- If no more information is available about the data, one can as well believe whatever predictor infers

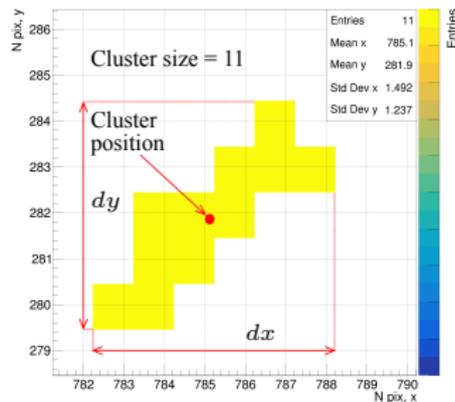
# Bayesian Neural Networks: Motivation



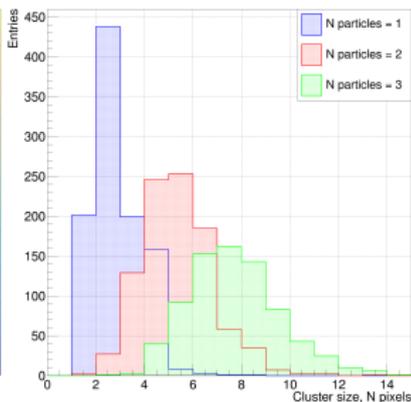
**Example of a classification problem with a significant class overlap**

- Classification task: based on the  $(x, y)$  position assign one of the two classes to an event
- It is clear what we want from a predictor in regions with no overlap
- What should the predictor do in the overlap region?
- If no more information is available about the data, one can as well believe whatever predictor infers
- If know about the data more than we're able to supply to the predictor uncertainty estimation becomes meaningful

# Bayesian Neural Networks: Motivation

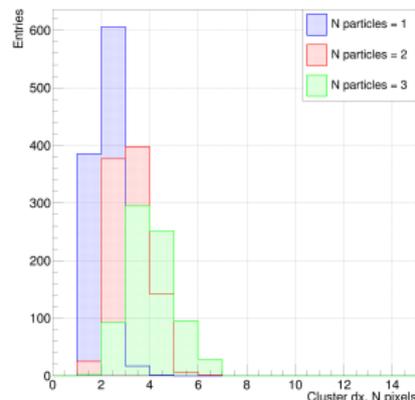


Generated cluster

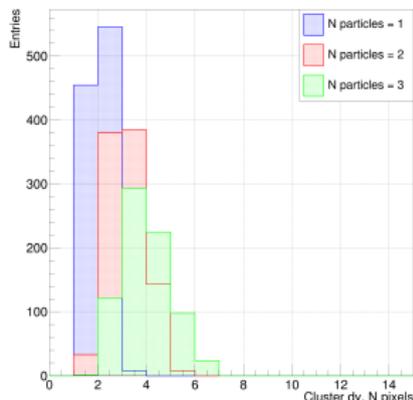


Cluster size distribution

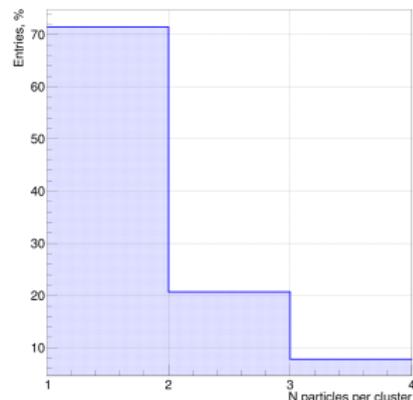
- Classification task: infer the number of particles in a cluster (out of 3)
- Significant parametric overlap between classes
- In LUXE we expect  $\sim 3$  times more 2-particle clusters than 3-particle clusters and  $\sim 4$  times more 1-particle clusters than 2-particle clusters



Cluster dx distribution



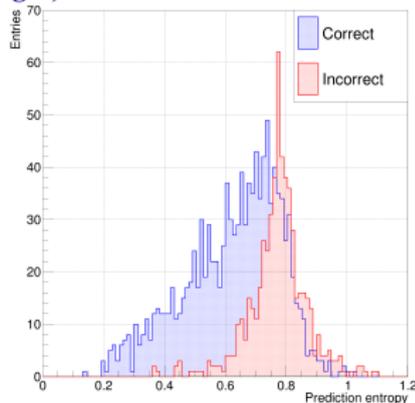
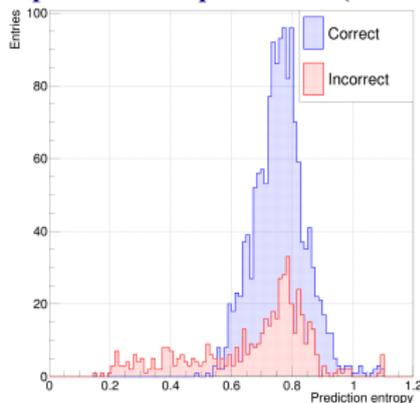
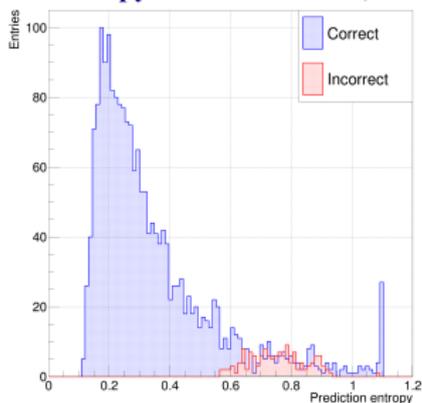
Cluster dy distribution



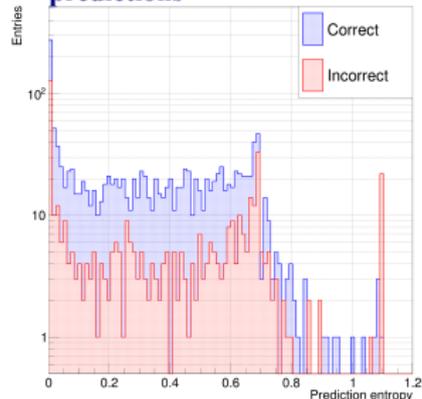
N particles per cluster distribution

# Performance analysis: accuracy and confidence

## Entropy distribution for 1,2 and 3 particle BNN predictions (left-to-right)



## Entropy distribution for 2-particle MLE NN predictions



• Prediction:  $C' = \arg \max_C \mathbb{P}(\tilde{\mathbf{x}}|\mathcal{D}), C \in \{1, 2, 3\}$

• Distribution's entropy:  $H(\mathbb{P}) = -\sum_n P_n \log(P_n)$

•  $H(\mathbb{P}) = 0 \iff F = \{1, 0, 0\}$

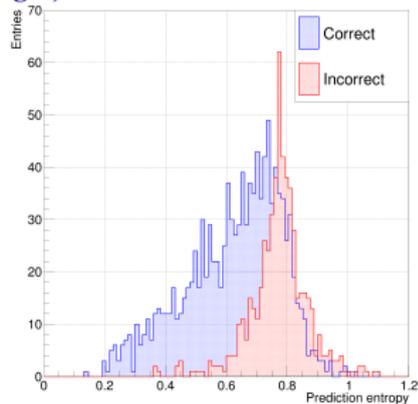
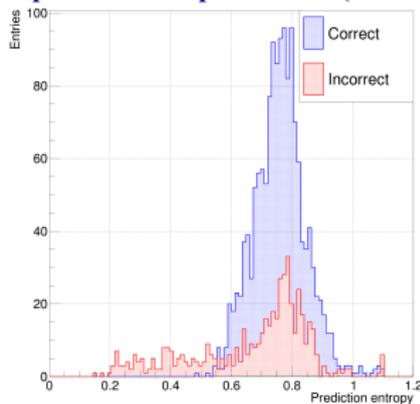
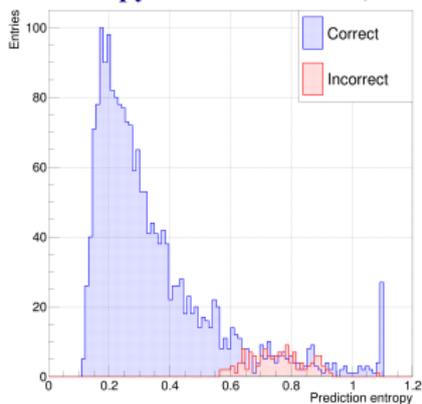
$H(\mathbb{P}) \rightarrow \max \iff \mathbb{P} = \{\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\}$

**Performance:**

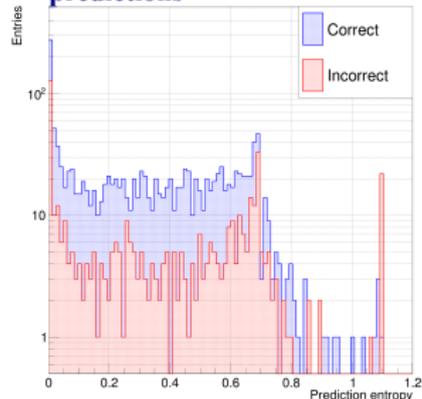
	Efficiency	BNN,%	MLE NN,%
$\varepsilon_1$		93.41	92.69
$\varepsilon_2$		73.89	75.13
$\varepsilon_3$		67.38	66.63

# Performance analysis: accuracy and confidence

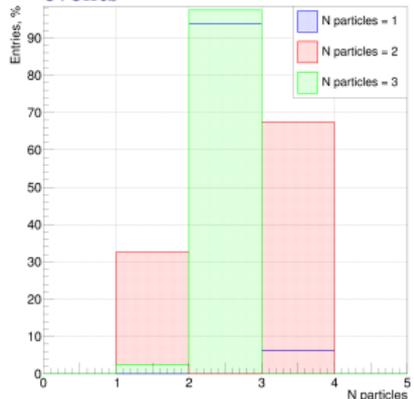
## Entropy distribution for 1,2 and 3 particle BNN predictions (left-to-right)



## Entropy distribution for 2-particle MLE NN predictions

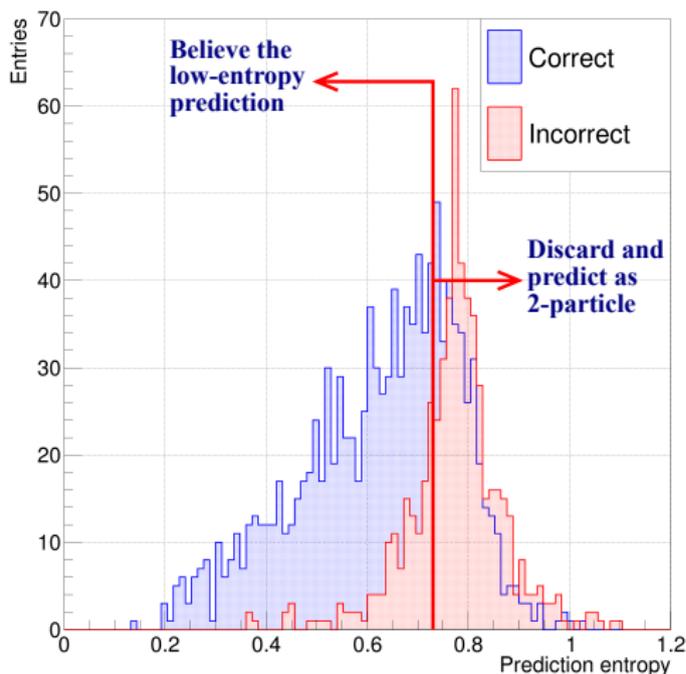


## Contamination of misclassified events



- There's no clear separation in MLE NN entropy distribution
- BNN's entropy distribution allows additional analysis to be employed
- Expected distribution between classes and contamination levels can be leveraged

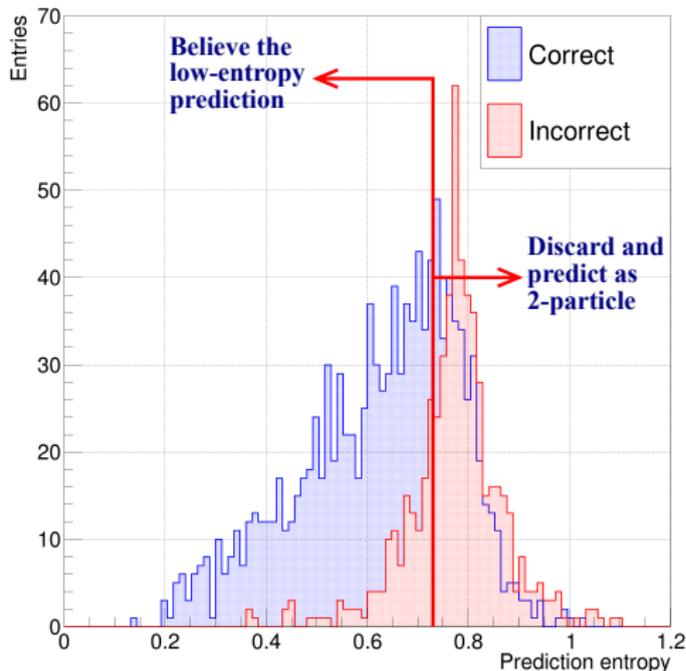
# Efficiency improvement



**Example of the entropy-thresholded decision making process**

- Decision making procedure:
  1. Set the decision threshold on the prediction's entropy  $H_0$
  2. If for a given 2 or 3 particle prediction  $H(\mathbb{P}) < H_0$  proceed with the network's inference
  3. If for a given 2 or 3 particle prediction  $H(\mathbb{P}) \geq H_0$  enforce classification as a 2-particle cluster
  4. Find the optimal  $H_0$  by observing the efficiency  $\varepsilon$  on the validation set

# Efficiency improvement



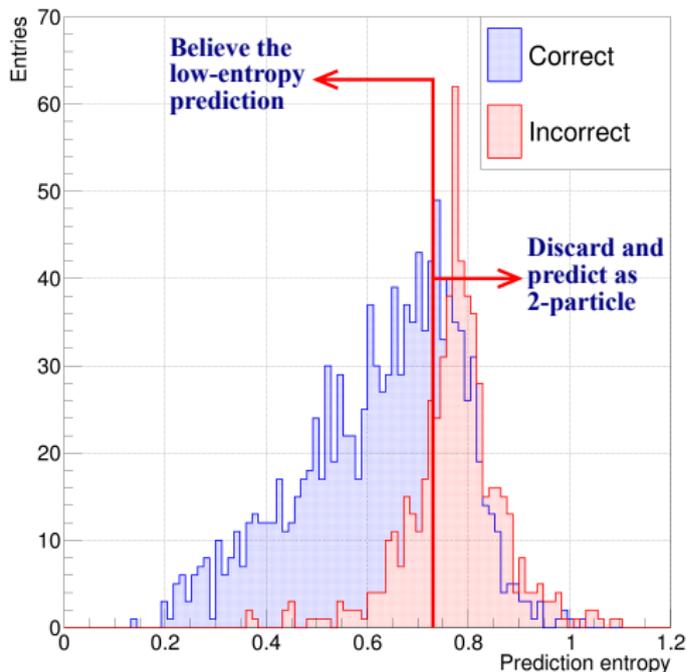
Example of the entropy-thresholded decision making process

- A separate dataset with the distribution of events between cluster classes has been prepared
- Overall efficiency  $\epsilon_{\text{tot}}$  of 87.24% and combined efficiency for 2 and 3 particle clusters  $\epsilon_{2\cup 3}$  of 74.55% have been achieved on the MLE network
- Thresholded BNN improves required computational time for tracking by  $\sim 5\text{-}10\%$

## BNN's efficiency for different entropy thresholds

$H_0$	0.65	0.68	0.71	0.74
$\epsilon_{\text{tot}}, \%$	90.22	<b>90.49</b>	90.33	90.07
$\epsilon_{2\cup 3}, \%$	83.27	<b>83.76</b>	83.12	82.52
$\epsilon_2, \%$	<b>97.14</b>	95.52	93.44	90.64
$\epsilon_3, \%$	40.63	45.21	49.58	<b>54.58</b>

# Efficiency improvement



Example of the entropy-thresholded decision making process

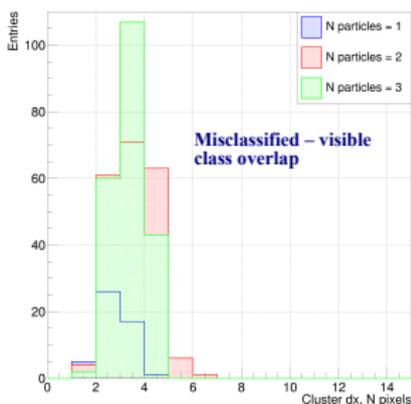
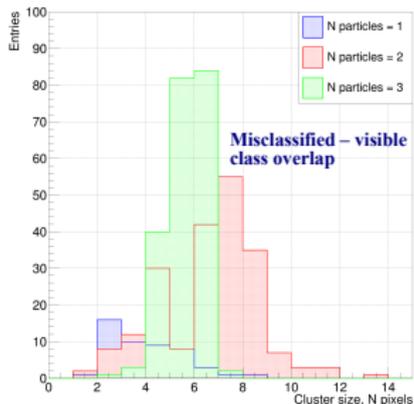
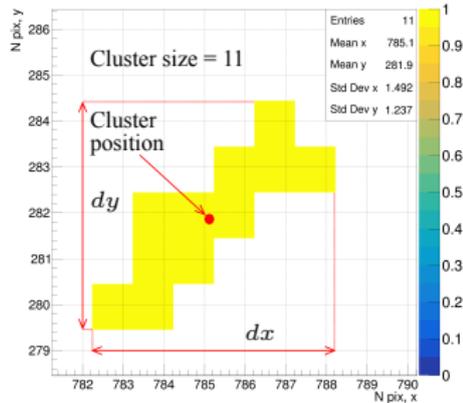
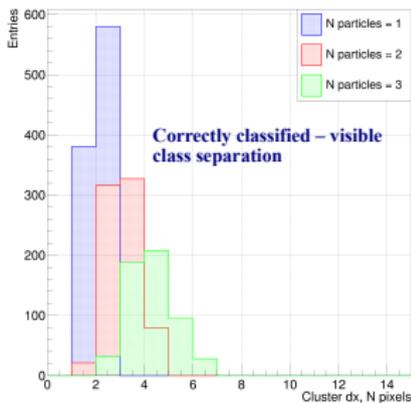
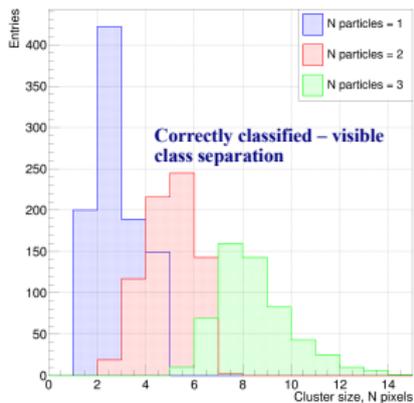
- Procedure is very crude – a lot of room for improvement
- Bernoulli-like decision tree can be implemented
- Distribution of above-threshold-entropy predictions can be accounted for
- More sensitive to the  $\varepsilon_3$  criterion of optimal  $H_0$  can be derived
- Architectures with better entropy separation can be searched for

## BNN's efficiency for different entropy thresholds

$H_0$	0.65	0.68	0.71	0.74
$\varepsilon_{\text{tot}}, \%$	90.22	<b>90.49</b>	90.33	90.07
$\varepsilon_{2\cup 3}, \%$	83.27	<b>83.76</b>	83.12	82.52
$\varepsilon_2, \%$	<b>97.14</b>	95.52	93.44	90.64
$\varepsilon_3, \%$	40.63	45.21	49.58	<b>54.58</b>

- Cluster merging causes track reconstruction issues. Clustering algorithm is required to perform over high range of occupancies
- MLE and Bayesian Neural Networks are applied to clustering
- Accuracies of  $\sim 93\%$ ,  $\sim 70\%$  and  $\sim 75\%$  are achieved for 1, 2 and 3-particle events respectively for both networks
- BNN's predictions allow for meaningful uncertainty estimation
- Ad-hoc procedure of alternative decision making has shown promising performance

# Backup: detailed observables investigation (1)

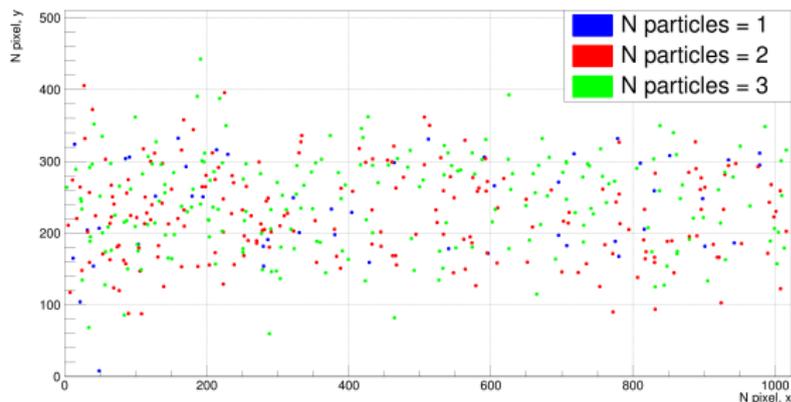
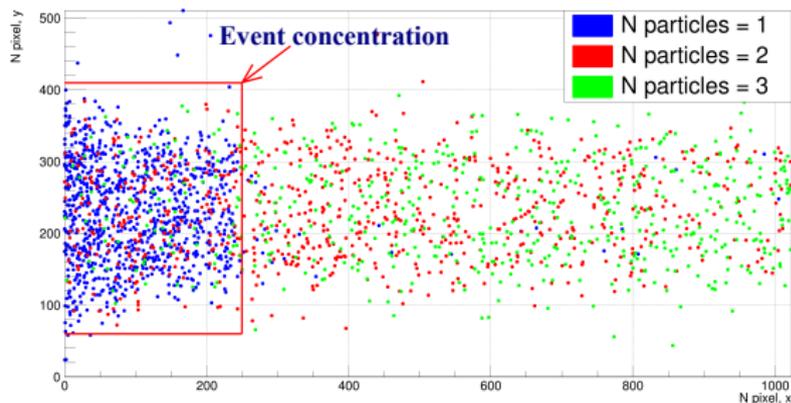


Cluster size for correct and misclassified predictions

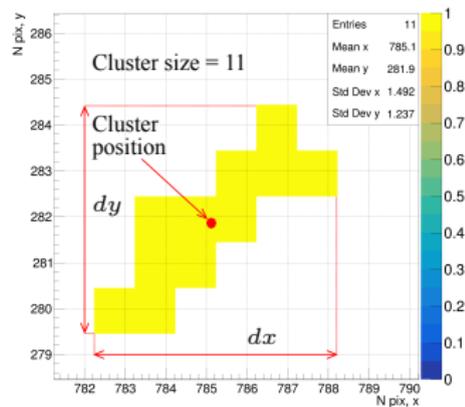
Cluster dx for correct and misclassified predictions

- BNN detects geometrical features
- Errors happen in regions of parametric overlap
- Better entropy separation requires architecture modifications

# Backup: detailed observables investigation (2)

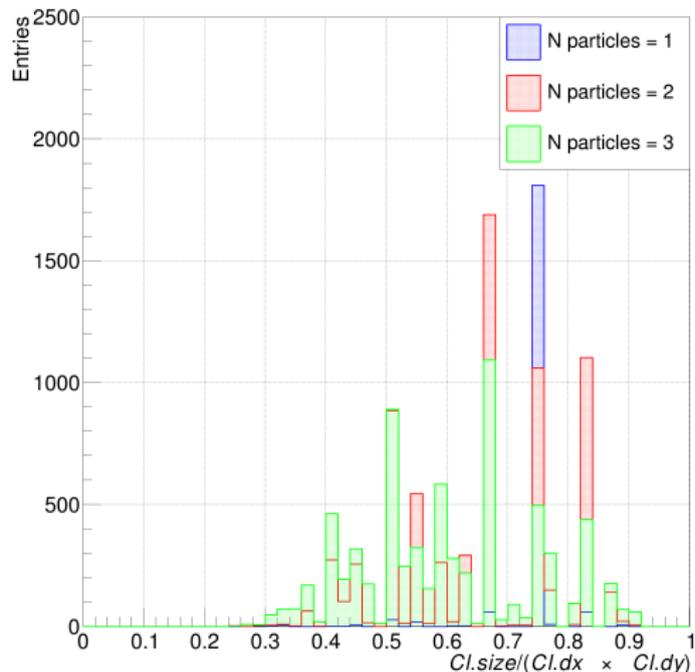


Cluster position distribution for correct (top) and misclassified (bottom) predictions



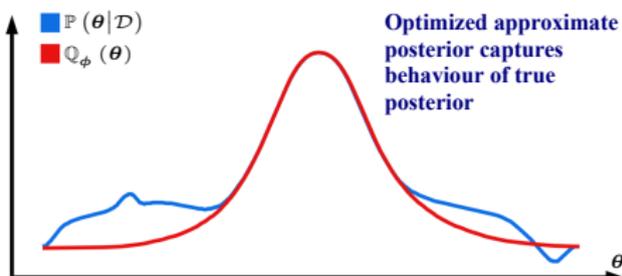
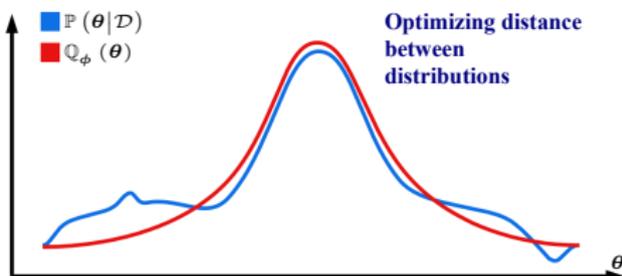
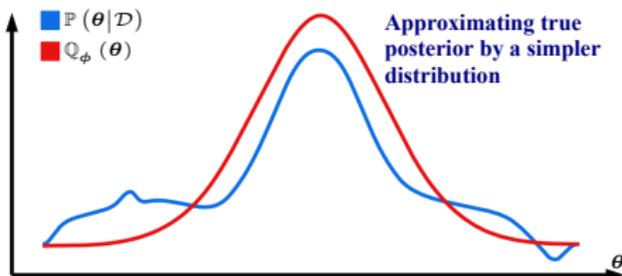
- BNN detects spatial features
- Higher multiplicity events are badly separated – additional parametrs, more complex network

# Backup: detailed observables investigation (3)



Percentage of area filled with active pixels over cluster's dimension rectangle

# Backup: SVI



- Predictions are Posterior-based.  
Have to infer  $\mathbb{P}(\theta|\mathcal{D})$
- MC methods are usually inefficient due to high dimensionality
- SVI – approximate + high variance gradients

Simple approximate posterior  $\mathbb{Q}_\phi(\theta)$



KL divergence:

$$D_{\text{KL}}(\mathbb{Q}_\phi(\theta) \parallel \mathbb{P}(\theta|\mathcal{D})) = \mathbb{E}_{\theta \sim \mathbb{Q}_\phi(\theta)} \left[ \log \left( \frac{\mathbb{Q}_\phi(\theta)}{\mathbb{P}(\theta|\mathcal{D})} \right) \right]$$

$$D_{\text{KL}}(\mathbb{Q}_\phi(\theta) \parallel \mathbb{P}(\theta|\mathcal{D})) \rightarrow \min_{\phi}$$

↑ unknown



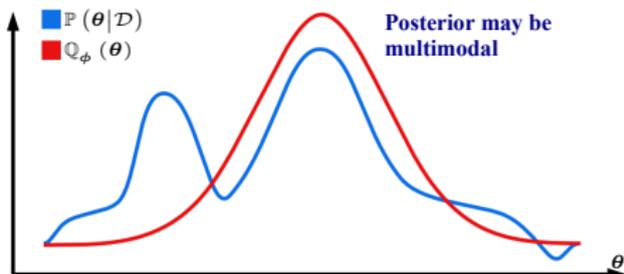
Evidence Lower Bound (ELBO):

$$\mathcal{L}(\mathbb{Q}) = \mathbb{E}_{\theta \sim \mathbb{Q}_\phi(\theta)} \left[ \log \left( \frac{\mathbb{P}(\mathcal{D}|\theta) \mathbb{P}(\theta)}{\mathbb{Q}_\phi(\theta)} \right) \right]$$

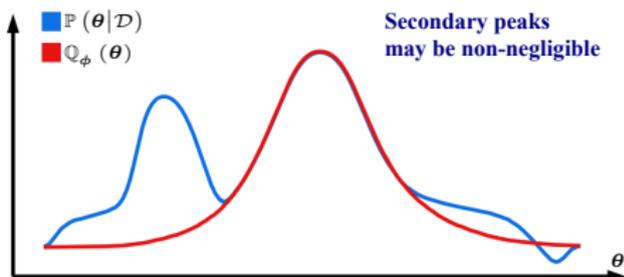
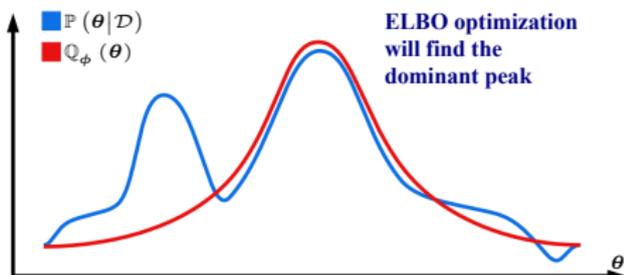
BNN output  
↓

$$\mathcal{L}(\mathbb{Q}) = \log \mathbb{P}(\mathcal{D}) - D_{\text{KL}}(\mathbb{Q}_\phi(\theta) \parallel \mathbb{P}(\theta|\mathcal{D})) \rightarrow \max_{\phi}$$

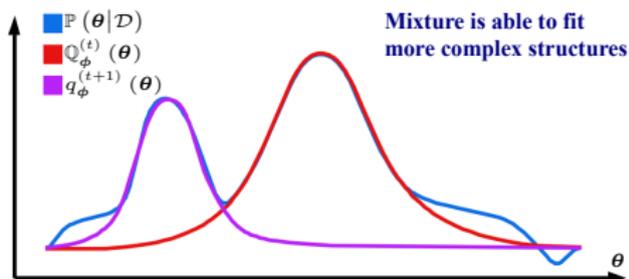
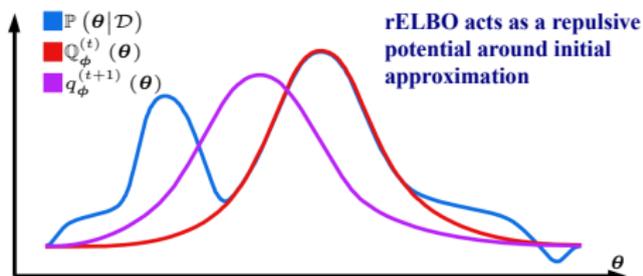
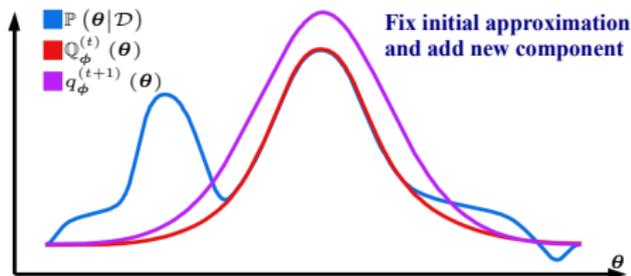
# Backup: Boosted BBVI (1)



- SVI depends on our assumption about the true posterior shape
- Due to high-variance gradients, it is hard to approximate multimodality by fitting a mixture
- Boosted BBVI algorithm:
  1. Find initial approximation by maximizing ELBO
  2. Fix initial approximation parameters



# Backup: Boosted BBVI (2)



- SVI depends on our assumption about the true posterior shape
- Due to high-variance gradients, it is hard to approximate multimodality by fitting a mixture
- Boosted BBVI algorithm:
  1. Find initial approximation by maximizing ELBO
  2. Fix initial approximation parameters
  3. Take next approximation component and construct a mixture
  4. Optimize added component parameters by maximizing rELBO

ELBO: approximation step has to be close to the posterior

$$\text{rELBO: } \mathcal{L}'(Q) = \mathbb{E}_{\theta \sim q_{\phi}^{(t+1)}} \left[ \log \left( \frac{\mathbb{P}(\mathcal{D}, \theta)}{q_{\phi}^{(t+1)}(\theta)} \right) \right] - \underbrace{\mathbb{E}_{\theta \sim q_{\phi}^{(t+1)}} \left[ \log Q_{\phi}^{(t)}(\theta) \right]}_{\text{residual: but not too close to the previous approximation step}}$$

residual: but not too close to the previous approximation step

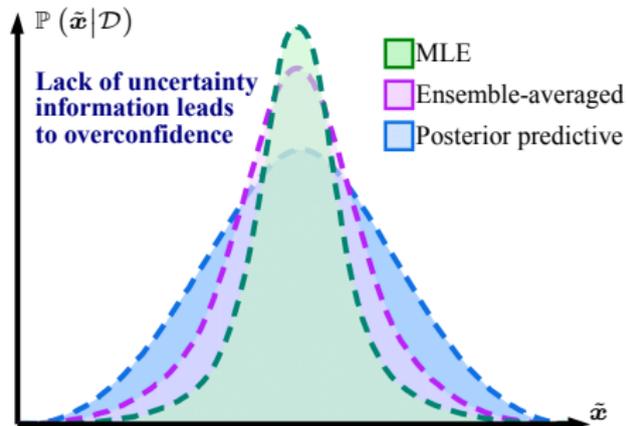
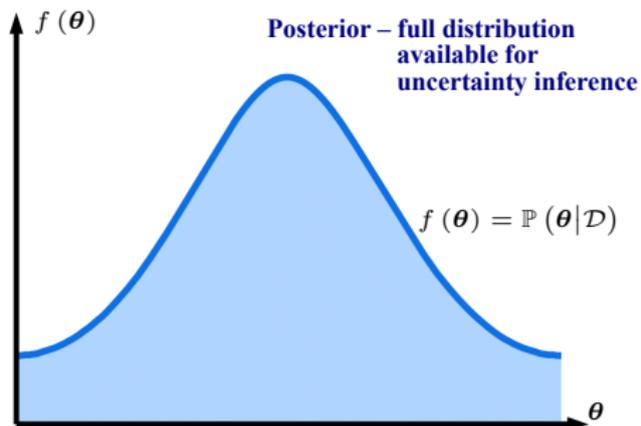
# Backup: additional BNN motivation

- BNNs are hard to overtrain:  $\mathbb{P}(\theta|\mathcal{D})$ , while  $\mathbb{P}(\mathcal{D}|\theta)$

↑  
**Emphasis on parameters**

↓  
**Emphasis on data**

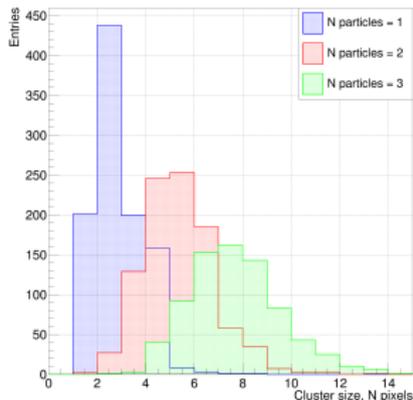
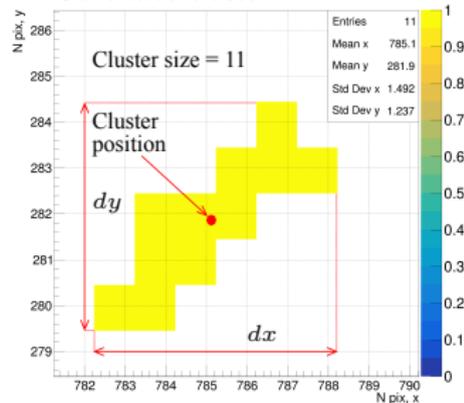
- BNNs perform better on OOD data
- BNNs are found to be generally more calibrated
- Bayesian formalism scales intuitively to multi-stage problems



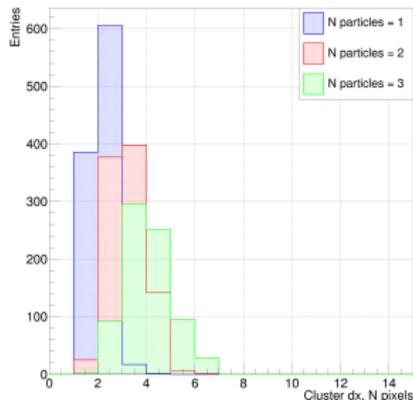
# Backup: BNN realization

- Data: Geant4 full physical simulation. Allpix2 digitization
- Network: three fully connected ReLU layers
- Inputs:  $7 \times 7$  matrix with active pixels, cluster position within tracker
- Weights: Laplace prior; Biases: Normal prior;
- Approx. posterior: Laplace distribution

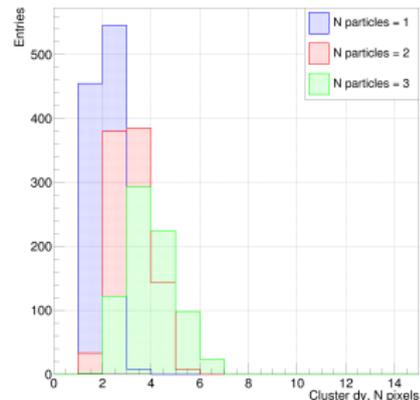
## Generated cluster



Cluster size distribution



Cluster dx distribution



Cluster dy distribution