

# Combining ID & AD: simultaneous particle identification and anomaly detection at collider experiments using Bayesian neural networks

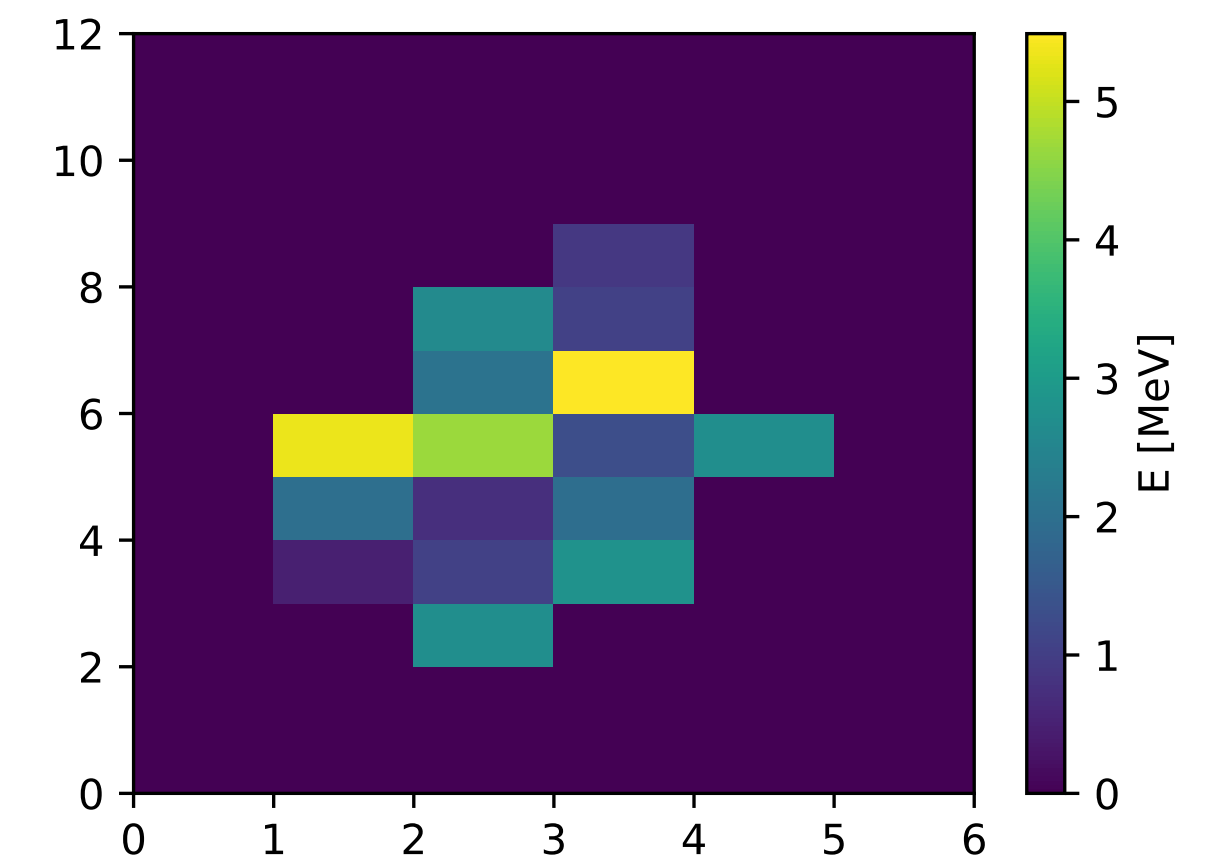
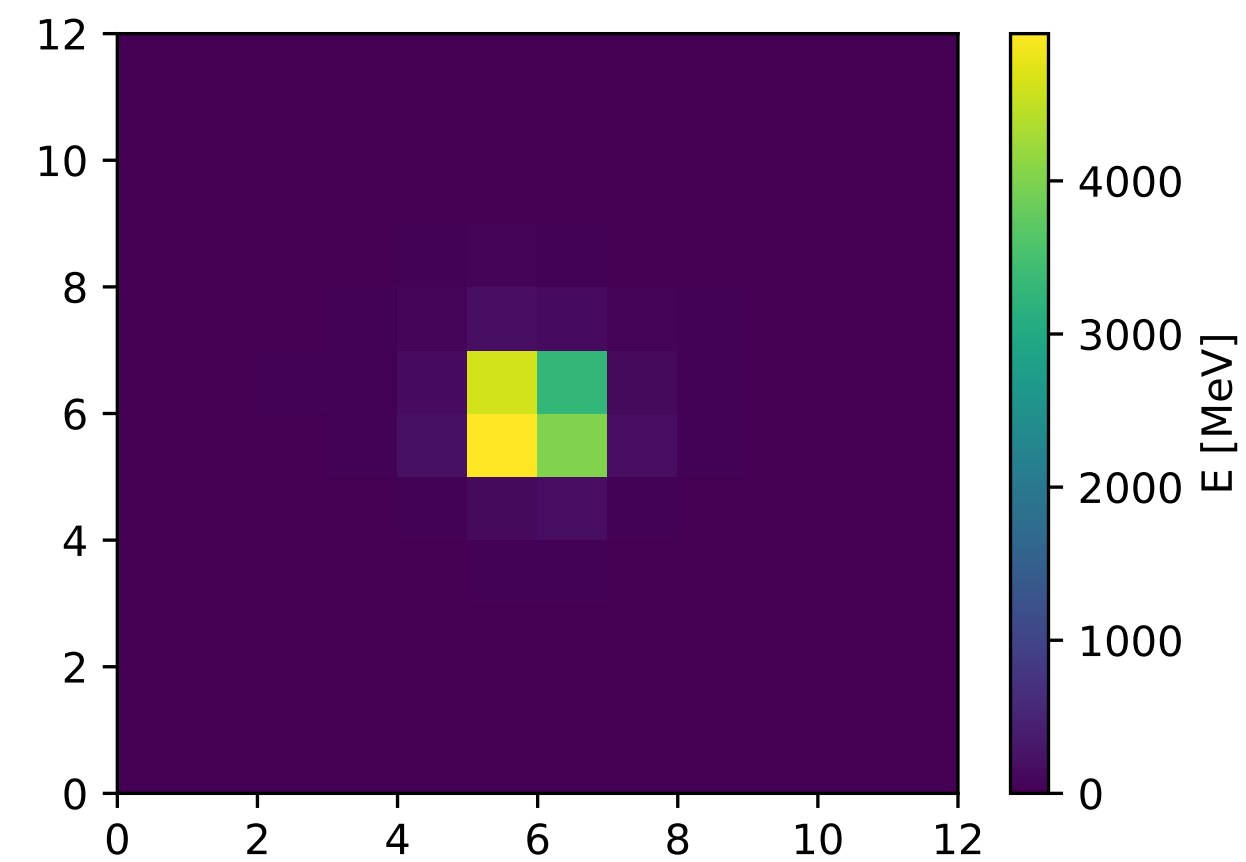
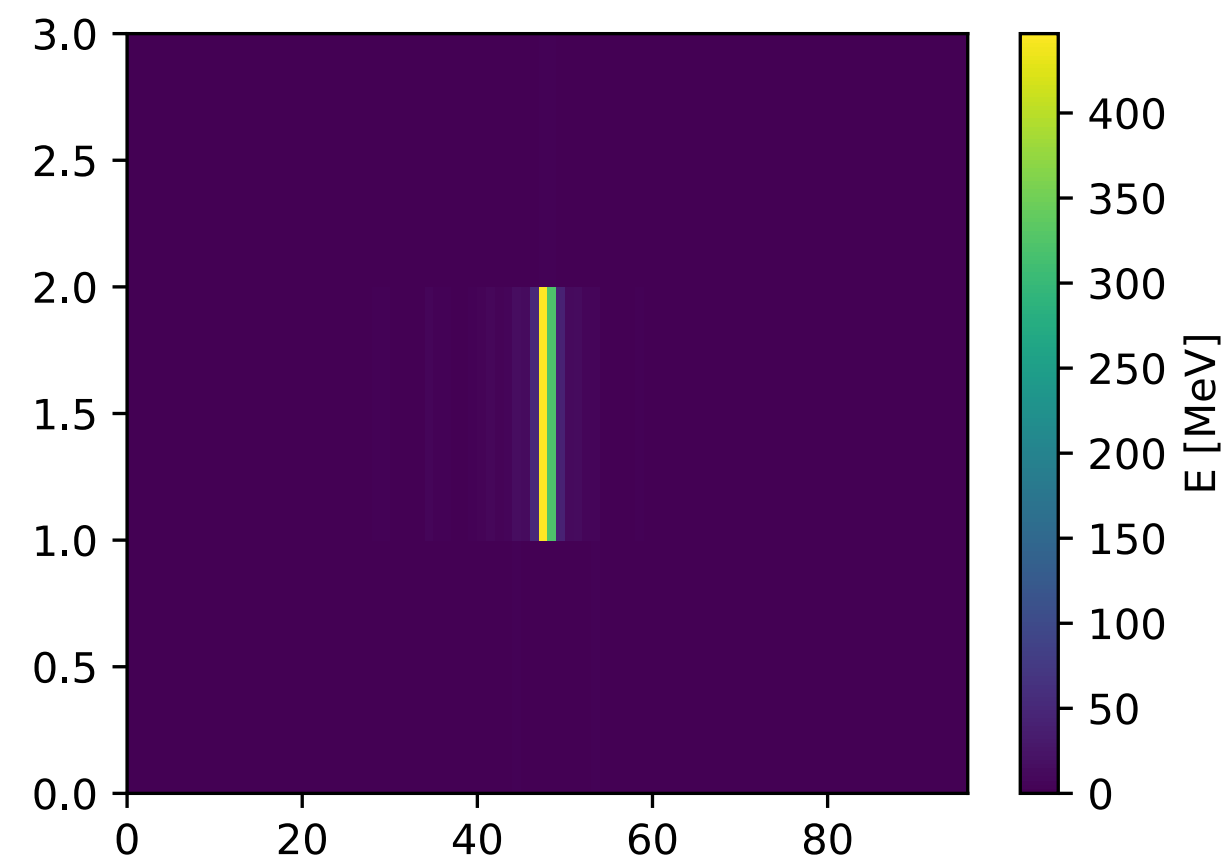
Wiehl 2022

Johannes Erdmann, Burim Ramosaj, **Daniel Wall**

**preliminary**

# Outline

- Bayesian vs. Deterministic NNs
- Our example use case:
  - Multiclass classification of EM calorimeter images
- Classification Performance
- Anomaly Detection
- Conclusions



# Bayesian Neural Networks

- Basic idea: learn the whole distribution over the NN weights  $w$  given training data  $\{X, Y\}$
- Which parameters  $w$  of a function  $f$  are likely to have generated  $Y$  from  $X$  ?

$$p(w|X, Y) = \frac{p(X, Y|w)p(w)}{p(X, Y)} = \frac{p(Y|X, w)\cancel{p(X|w)}p(w)}{p(Y|X)\cancel{p(X)}} = \frac{p(Y|X, w)p(w)}{\int p(Y|X, w)p(w)dw}$$

with for example  $p(y = d|\mathbf{x}, \omega) = \frac{\exp(f_d^\omega(\mathbf{x}))}{\sum_{d'} \exp(f_{d'}^\omega(\mathbf{x}))}$

- Often intractable  $\rightarrow$  approximate with simpler function  $q_\theta(\omega)$  minimizing KL divergence

$$\text{KL}(q_\theta(\omega) || p(\omega|\mathbf{X}, \mathbf{Y})) = \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega = - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \text{KL}(q_\theta(\omega) || p(\omega)) + \text{const}$$

- Sampling from optimal  $q_\theta^*(\omega)$   $\rightarrow$  distribution of predictions instead of a point estimate

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) \approx \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) q_\theta^*(\omega) d\omega$$

# Bayesian Neural Networks

- Basic idea: learn the whole distribution over the NN weights  $w$  given training data  $\{X, Y\}$
- Which parameters  $w$  of a function  $f$  are likely to have generated  $Y$  from  $X$  ?

$$p(w|X, Y) = \frac{p(X, Y|w)p(w)}{p(X, Y)} = \frac{p(Y|X, w)\cancel{p(X|w)}p(w)}{p(Y|X)\cancel{p(X)}} = \frac{p(Y|X, w)p(w)}{\int p(Y|X, w)p(w)dw}$$

with for example  $p(y = d|\mathbf{x}, \omega) = \frac{\exp(f_d^\omega(\mathbf{x}))}{\sum_{d'} \exp(f_{d'}^\omega(\mathbf{x}))}$

- Often intractable  $\rightarrow$  approximate with simpler function  $q_\theta(\omega)$  minimizing KL divergence

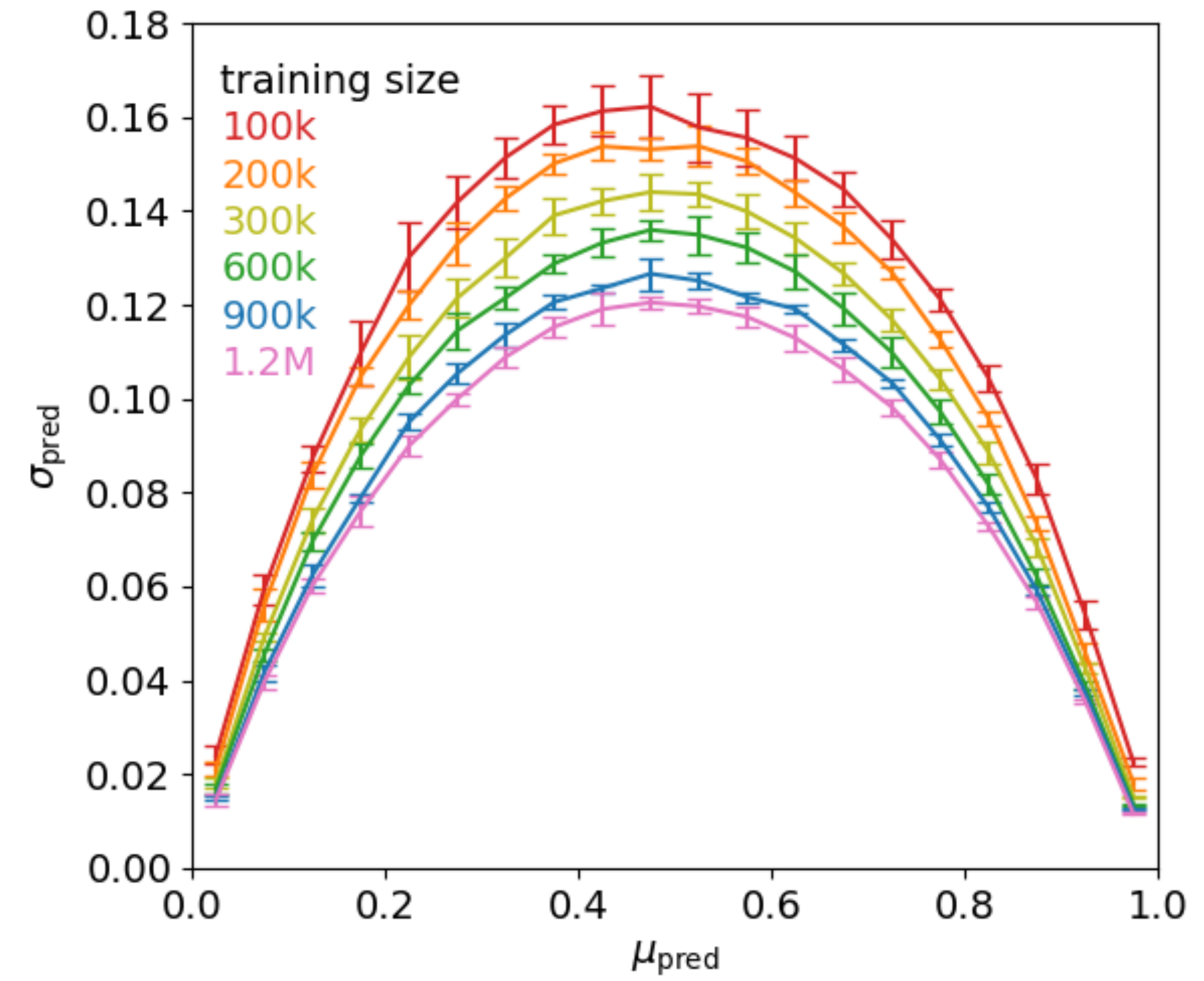
$$\text{KL}(q_\theta(\omega) || p(\omega|\mathbf{X}, \mathbf{Y})) = \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega = - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \text{KL}(q_\theta(\omega) || p(\omega)) + \text{const}$$

- Sampling from optimal  $q_{\theta^*}(\omega)$   $\rightarrow$  distribution of predictions instead of a point estimate

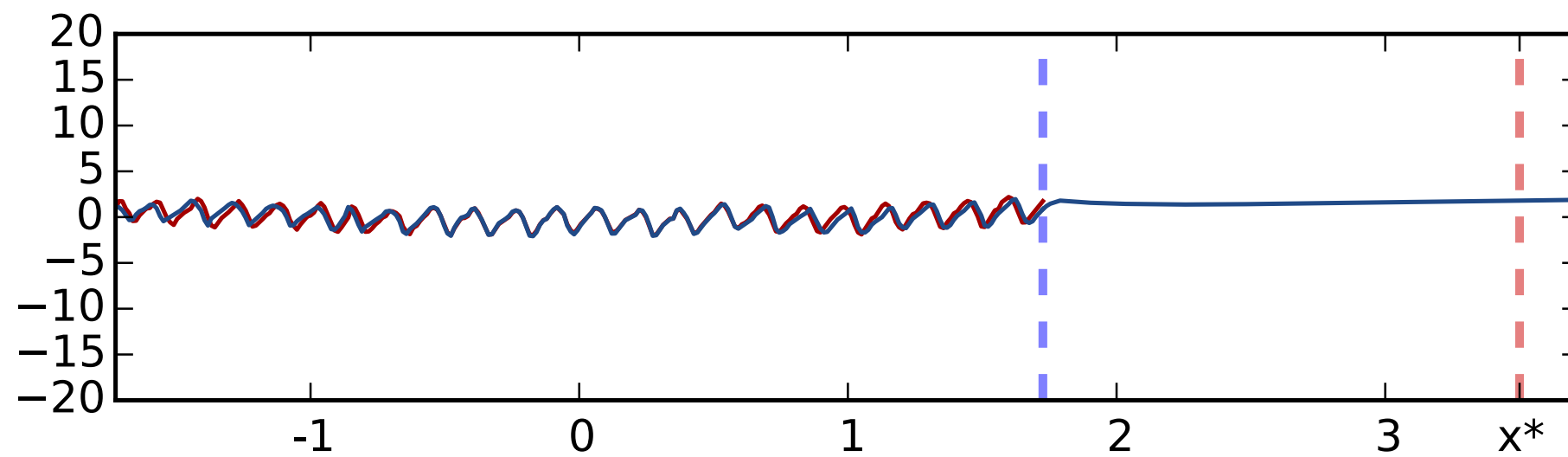
$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) \approx \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) q_{\theta^*}(\omega) d\omega$$

# Bayesian Neural Networks

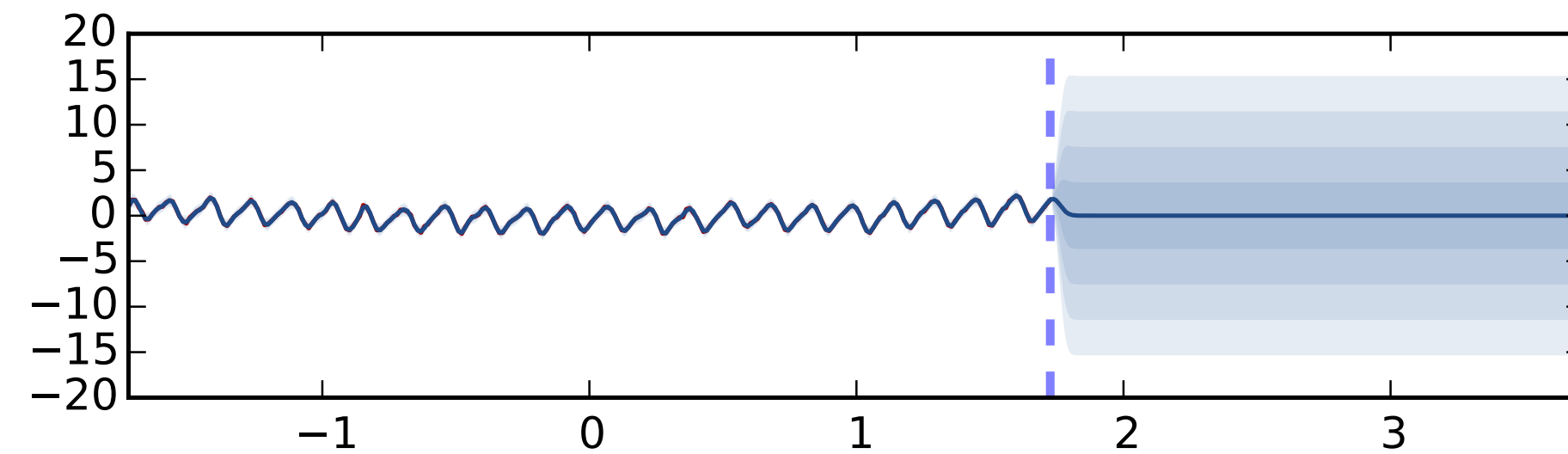
- Small training datasets lead to larger uncertainties in the BNN predictions (example from top-tagging in 1904.10004)
- BNNs predictions for out-of-distribution test samples can have large uncertainties



1904.10004



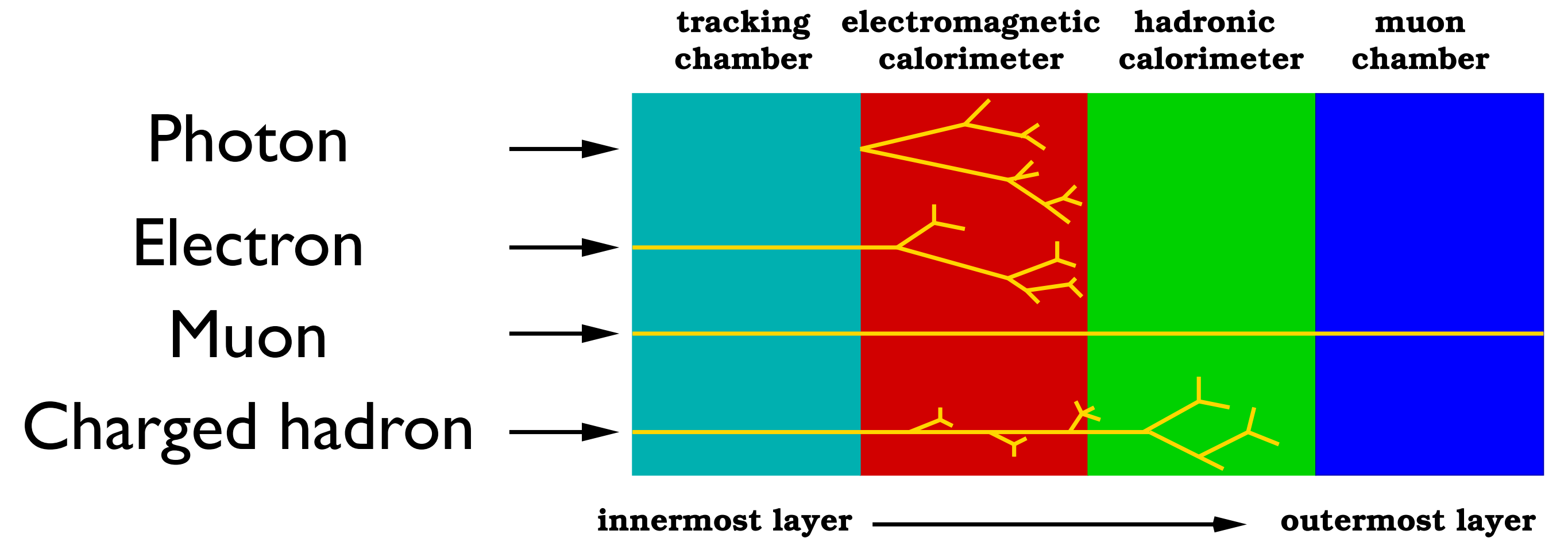
(a) Standard deep learning model



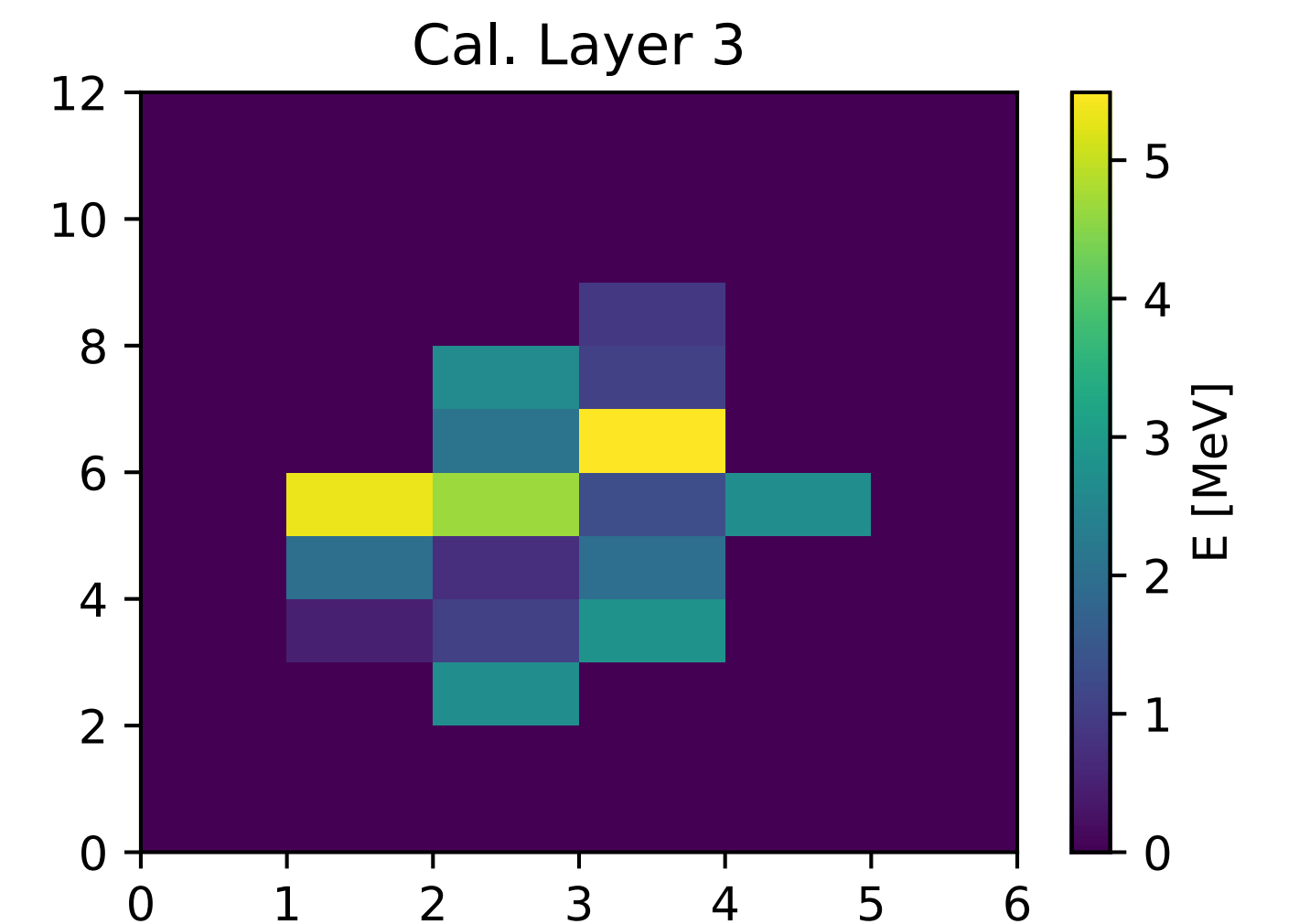
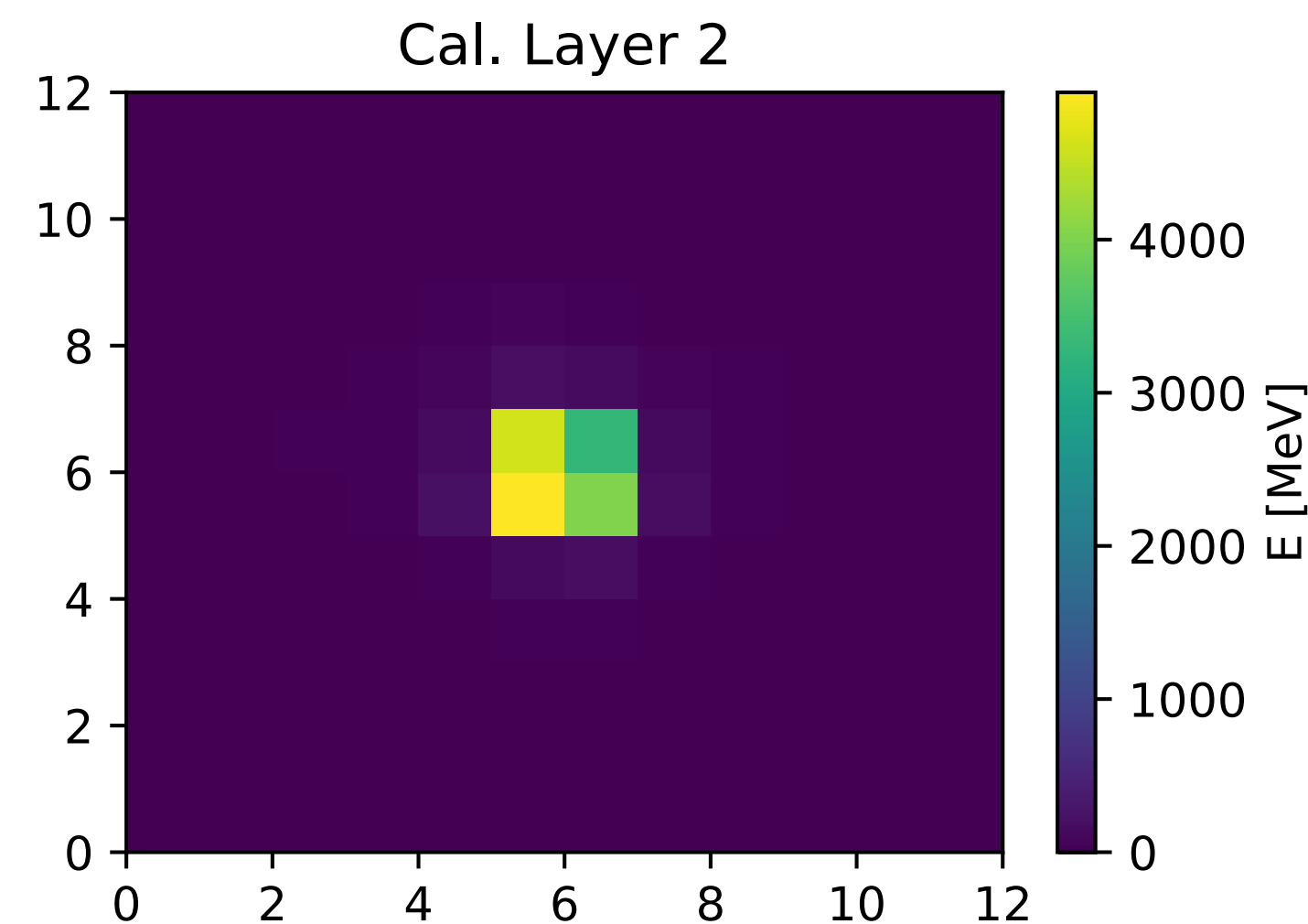
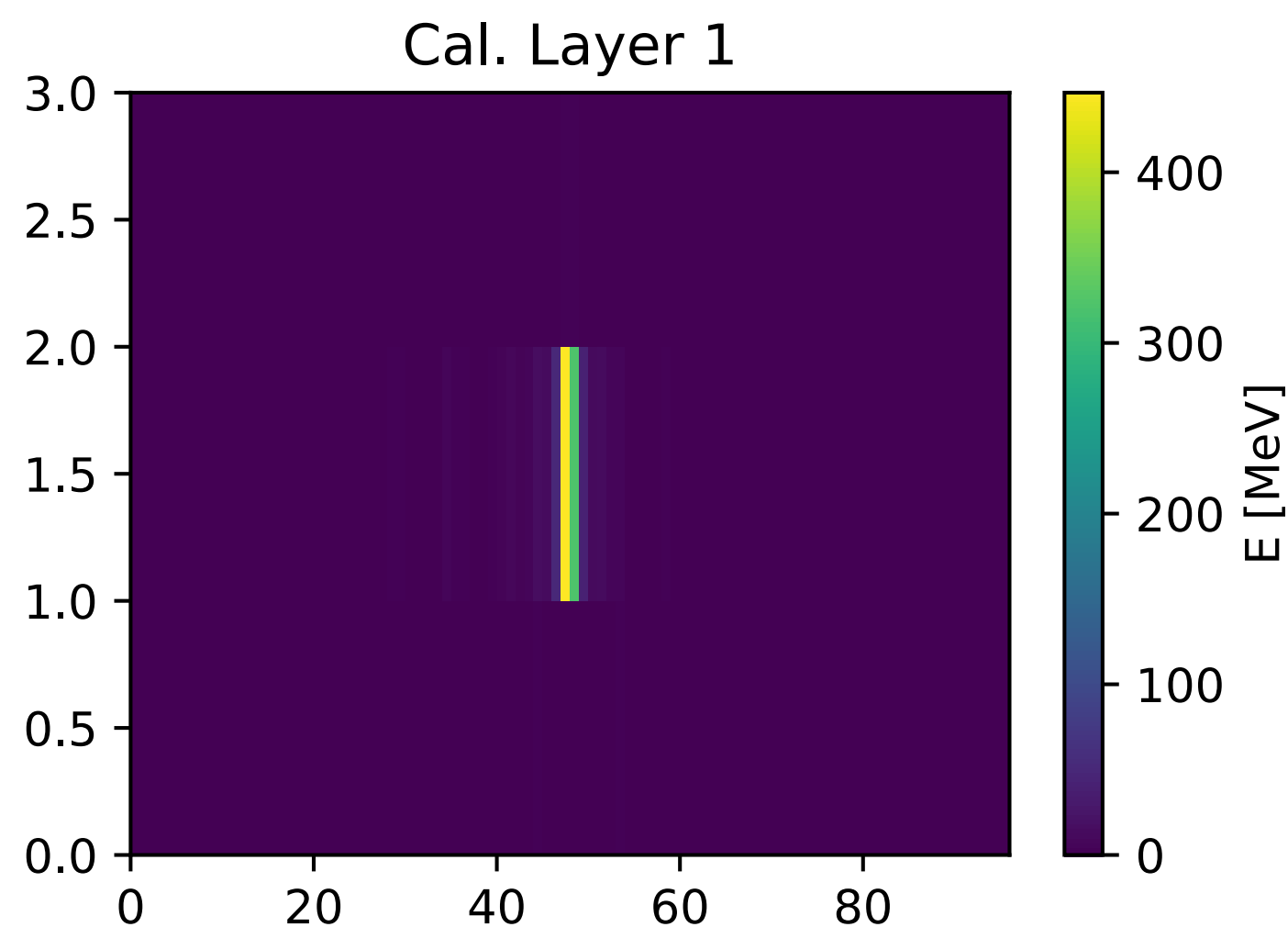
(b) Probabilistic model

# Our Example Use Case

- Classification of images in EM calorimeters = photon identification
- Main background:
  - High-energy  $\pi^0 \rightarrow \gamma\gamma$
- Toy EM calorimeter à la I712.10321
  - ATLAS-like: 3 layers of LAr+Pb
  - 1 m from Geant4 particle gun

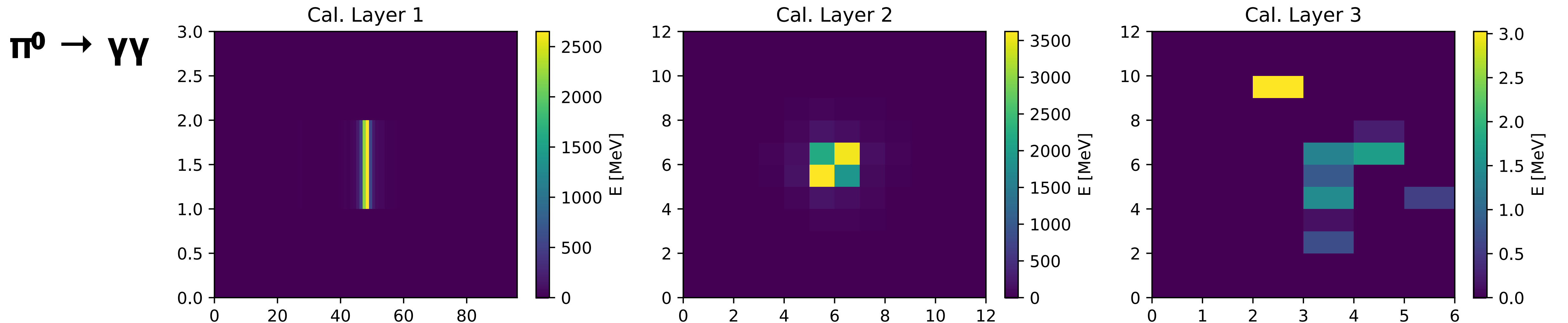
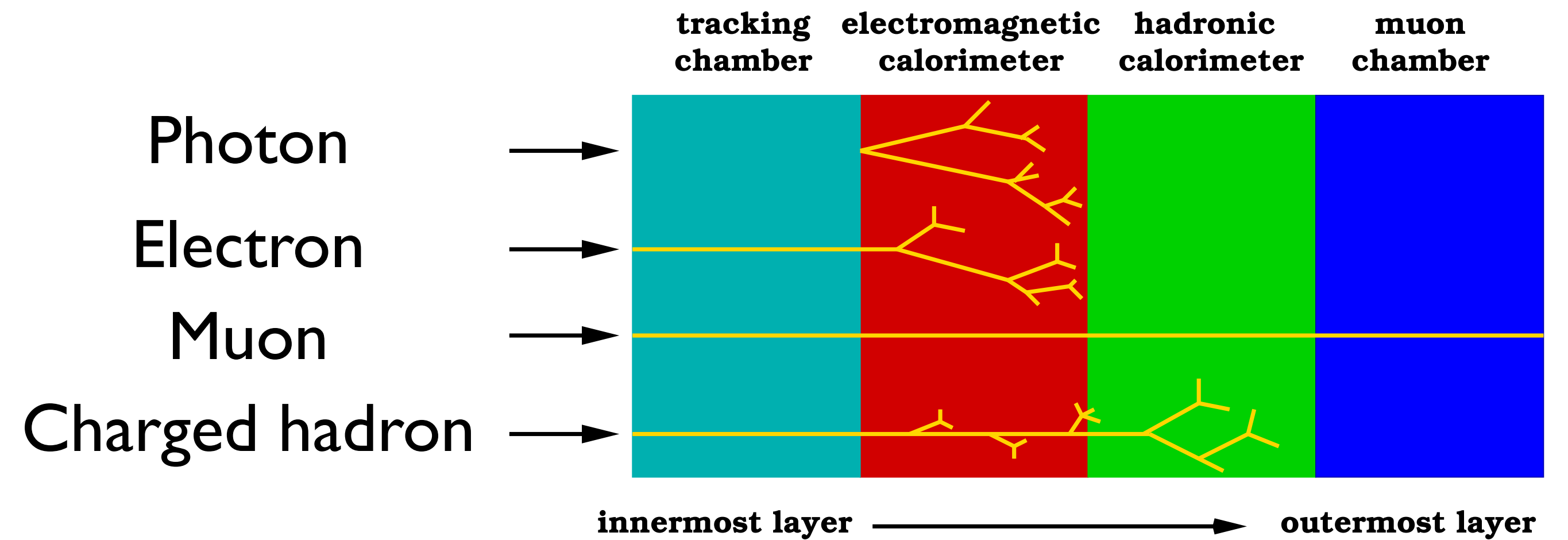


## Photon



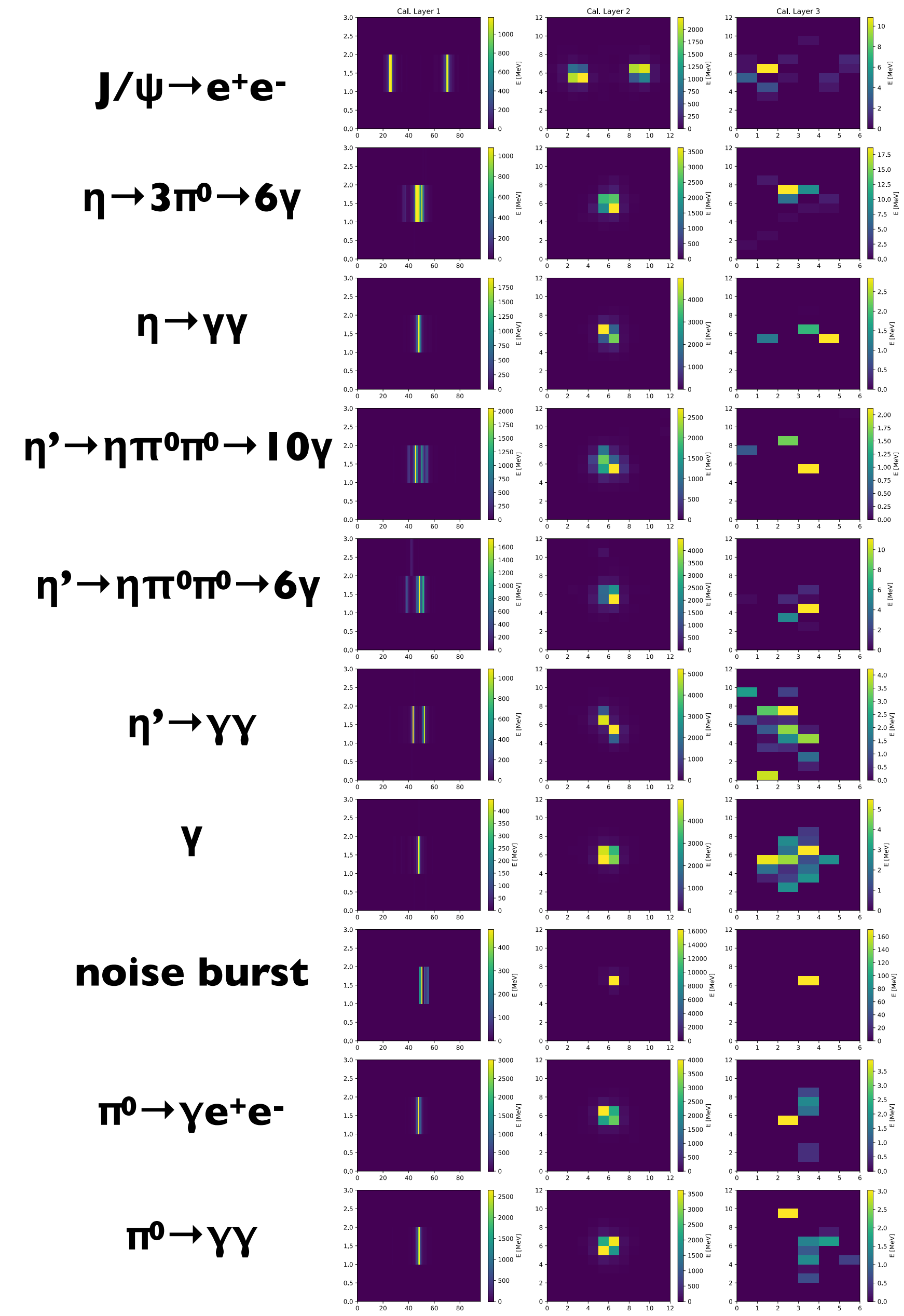
# Our Example Use Case

- Classification of images in EM calorimeters = photon identification
- Main background:
  - High-energy  $\pi^0 \rightarrow \gamma\gamma$
- Toy EM calorimeter à la I712.10321
  - ATLAS-like: 3 layers of LAr+Pb
  - 1 m from Geant4 particle gun



# Our Example Use Case

- Particle gun kinetic energy: 20 GeV
- Single photon (signal)
- 8 physical background classes:
  - Different purely EM decays of mesons with different masses  $\rightarrow$  different opening angles
- + noise background class (noise burst in 2nd layer with 1% cross talk to neighbouring cells)
- Photon ID algorithm at LHC would be trained on effective mixture via parton shower programs





# Classification Performance

- Setup: 2D CNN with 8 filters and 3x3 kernel size for each calorimeter layer  
+ 10 output nodes all using Flipout (1803.04386)
- Assume weights to be Gaussian distributed, uncorrelated and with Gaussian priors
- Very similar performance to deterministic NN

### Deterministic NN

$J/\psi \rightarrow e^+ e^-$	1	0	0	0	0	0	0	0	0	
$\eta \rightarrow 3\pi^0 \rightarrow 6\gamma$	0	0.99	0	0	0.01	0	0	0	0	
$\eta \rightarrow 2\gamma$	0	0	1	0	0	0	0	0	0	
$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 10\gamma$	0	0	0	0.99	0.01	0	0	0	0	
$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 6\gamma$	0	0.01	0	0.01	0.98	0	0	0	0	
$\eta' \rightarrow 2\gamma$	0	0	0	0	0	1	0	0	0	
$\gamma$	0	0	0	0	0	0	0.98	0.01	0.01	
Noiseburst	0	0	0	0	0	0	0	1	0	
$\pi^0 \rightarrow e^+ e^- \gamma$	0	0	0	0	0	0	0	0	0.75	0.24
$\pi^0 \rightarrow 2\gamma$	0	0	0	0	0	0	0	0	0.25	0.74
	$J/\psi \rightarrow e^+ e^-$	$\eta \rightarrow 3\pi^0 \rightarrow 6\gamma$	$\eta \rightarrow 2\gamma$	$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 10\gamma$	$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 6\gamma$	$\eta' \rightarrow 2\gamma$	$\gamma$	Noiseburst	$\pi^0 \rightarrow e^+ e^- \gamma$	$\pi^0 \rightarrow 2\gamma$

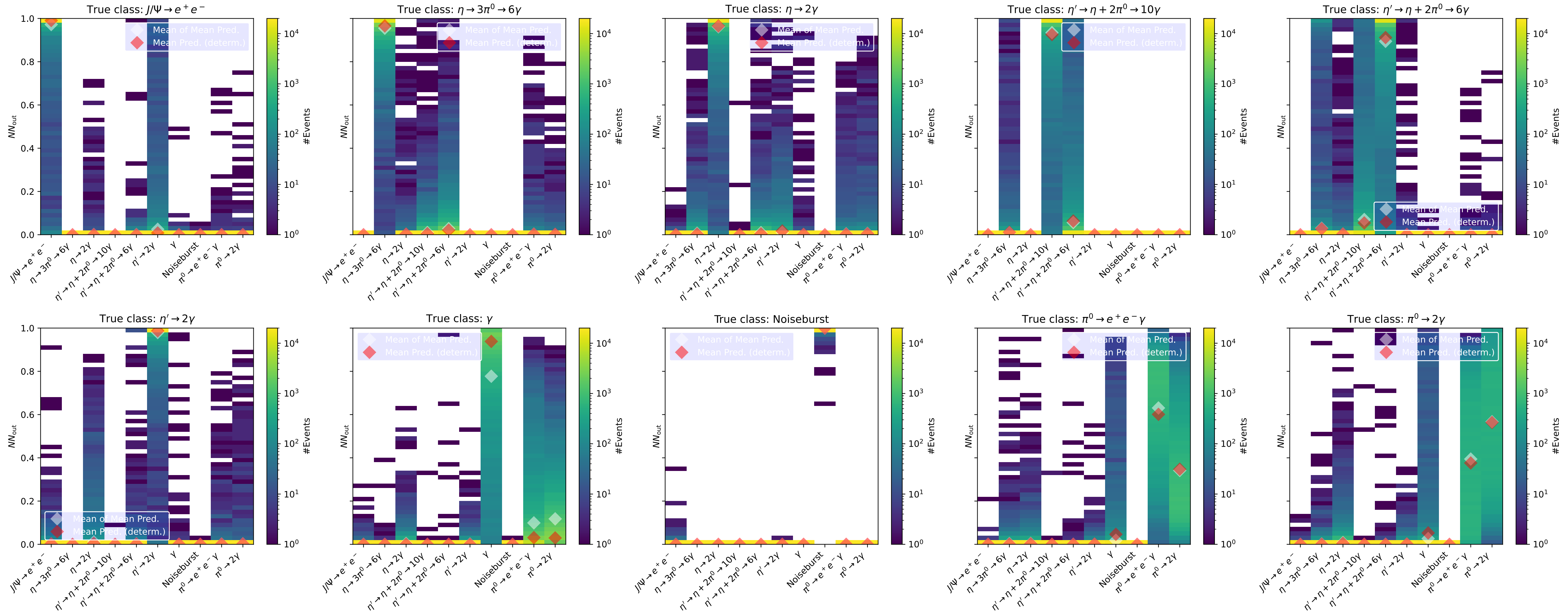
### Bayesian NN

$J/\psi \rightarrow e^+ e^-$	1	0	0	0	0	0	0	0	0	
$\eta \rightarrow 3\pi^0 \rightarrow 6\gamma$	0	0.99	0	0	0.01	0	0	0	0	
$\eta \rightarrow 2\gamma$	0	0	1	0	0	0	0	0	0	
$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 10\gamma$	0	0	0	0.99	0.01	0	0	0	0	
$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 6\gamma$	0	0.01	0	0.01	0.98	0	0	0	0	
$\eta' \rightarrow 2\gamma$	0	0	0	0	0	1	0	0	0	
$\gamma$	0	0	0	0	0	0	0.98	0.01	0.01	
Noiseburst	0	0	0	0	0	0	0	1	0	
$\pi^0 \rightarrow e^+ e^- \gamma$	0	0	0	0	0	0	0.02	0	0.75	0.23
$\pi^0 \rightarrow 2\gamma$	0	0	0	0	0	0	0.03	0	0.23	0.74
	$J/\psi \rightarrow e^+ e^-$	$\eta \rightarrow 3\pi^0 \rightarrow 6\gamma$	$\eta \rightarrow 2\gamma$	$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 10\gamma$	$\eta' \rightarrow \eta + 2\pi^0 \rightarrow 6\gamma$	$\eta' \rightarrow 2\gamma$	$\gamma$	Noiseburst	$\pi^0 \rightarrow e^+ e^- \gamma$	$\pi^0 \rightarrow 2\gamma$

# Results

- For each image, sample from the weight distributions  $\rightarrow$  mean & variance per image

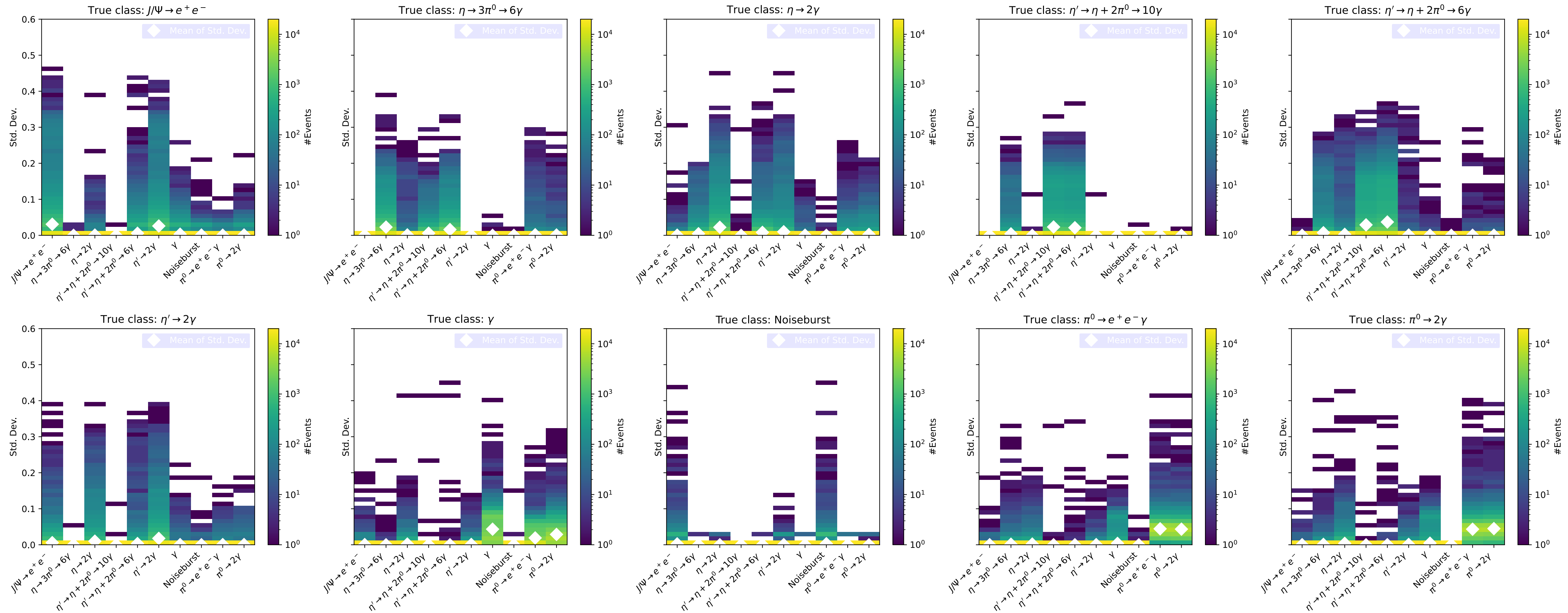
## Mean



# Results

- For each image, sample from the weight distributions  $\rightarrow$  mean & variance per image

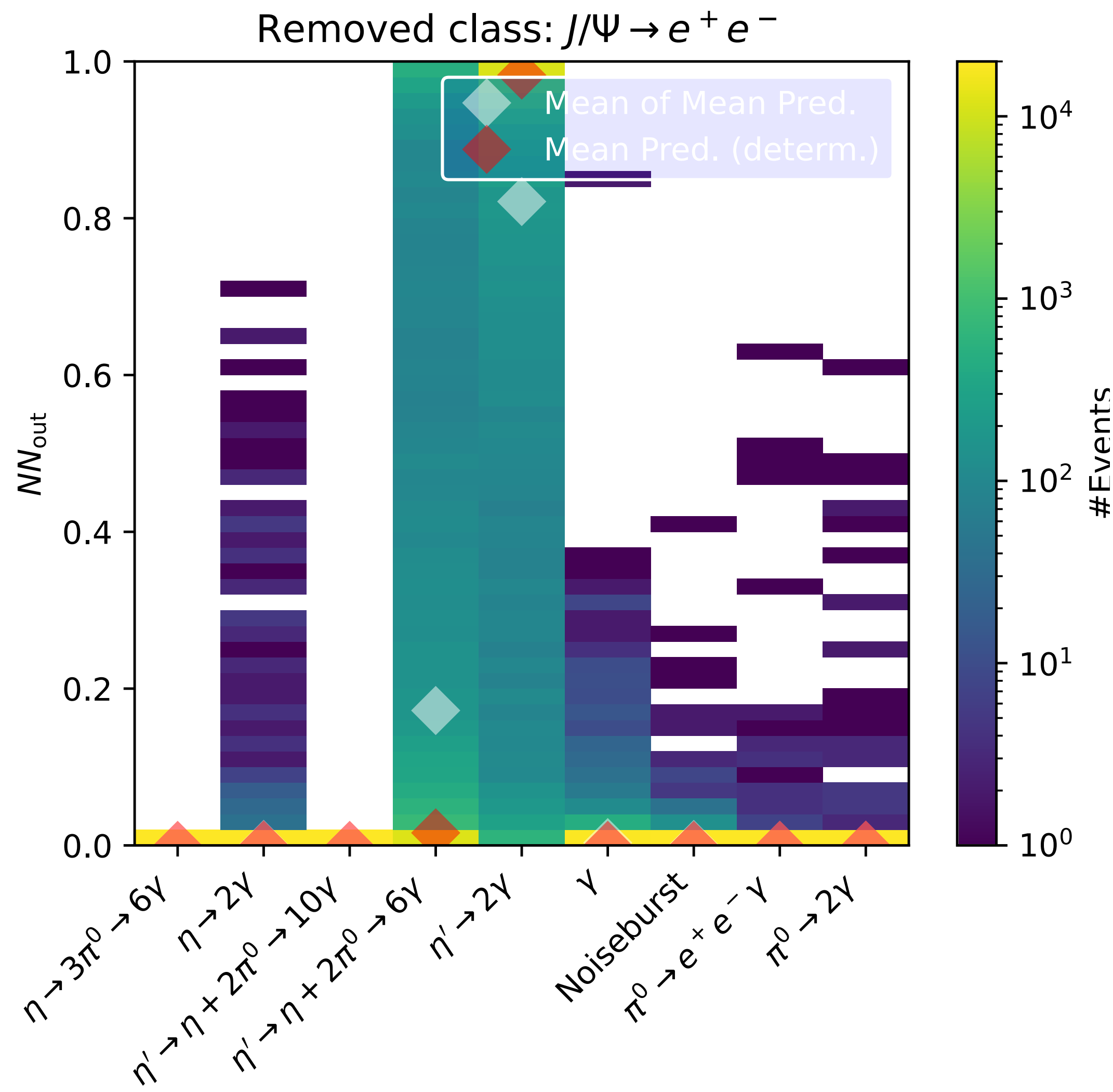
## Std. Dev.



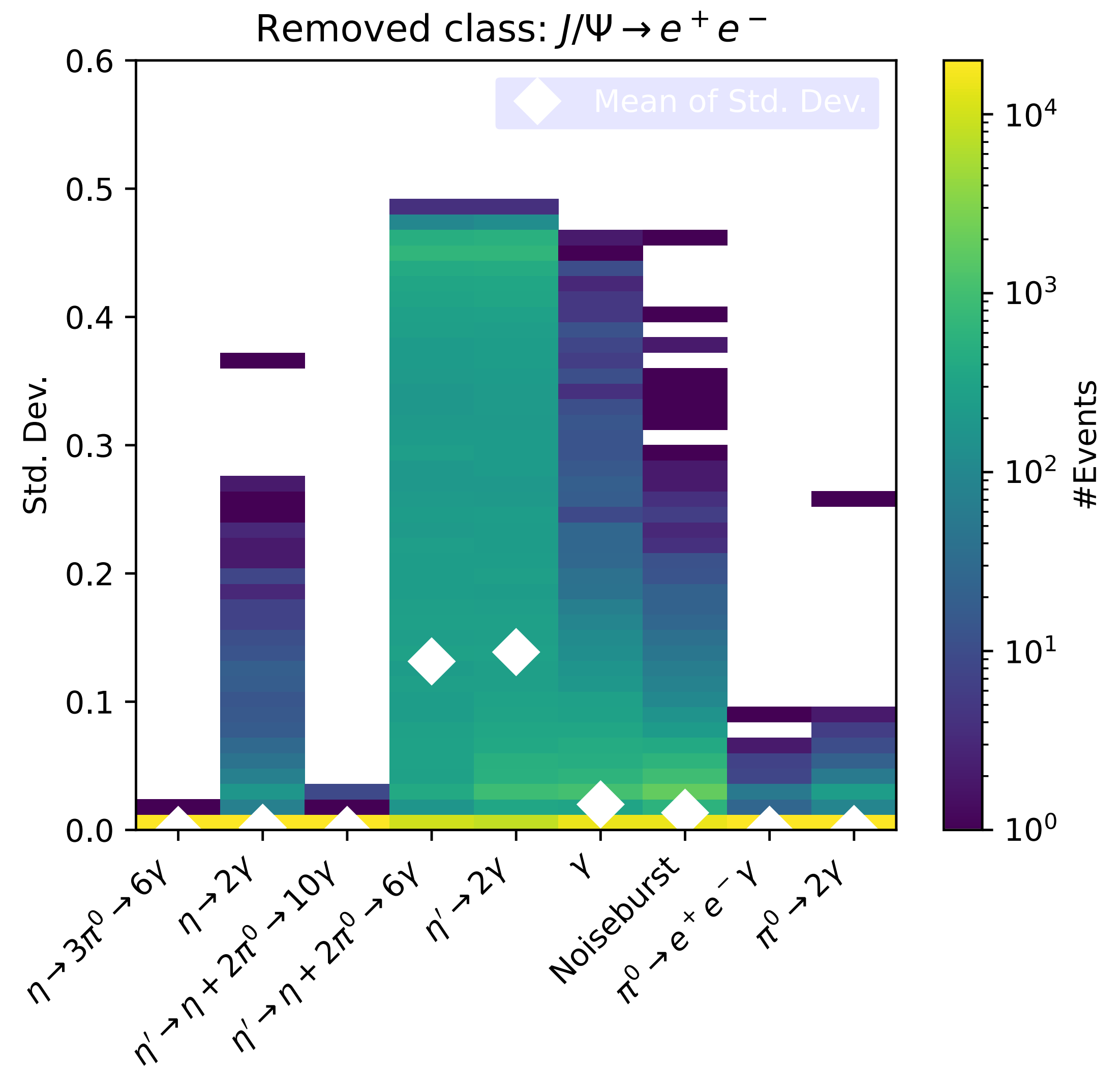
# Results

- Now: remove one class at a time during training = anomaly (here:  $J/\psi \rightarrow e^+e^-$ )

## Mean



## Std. Dev.



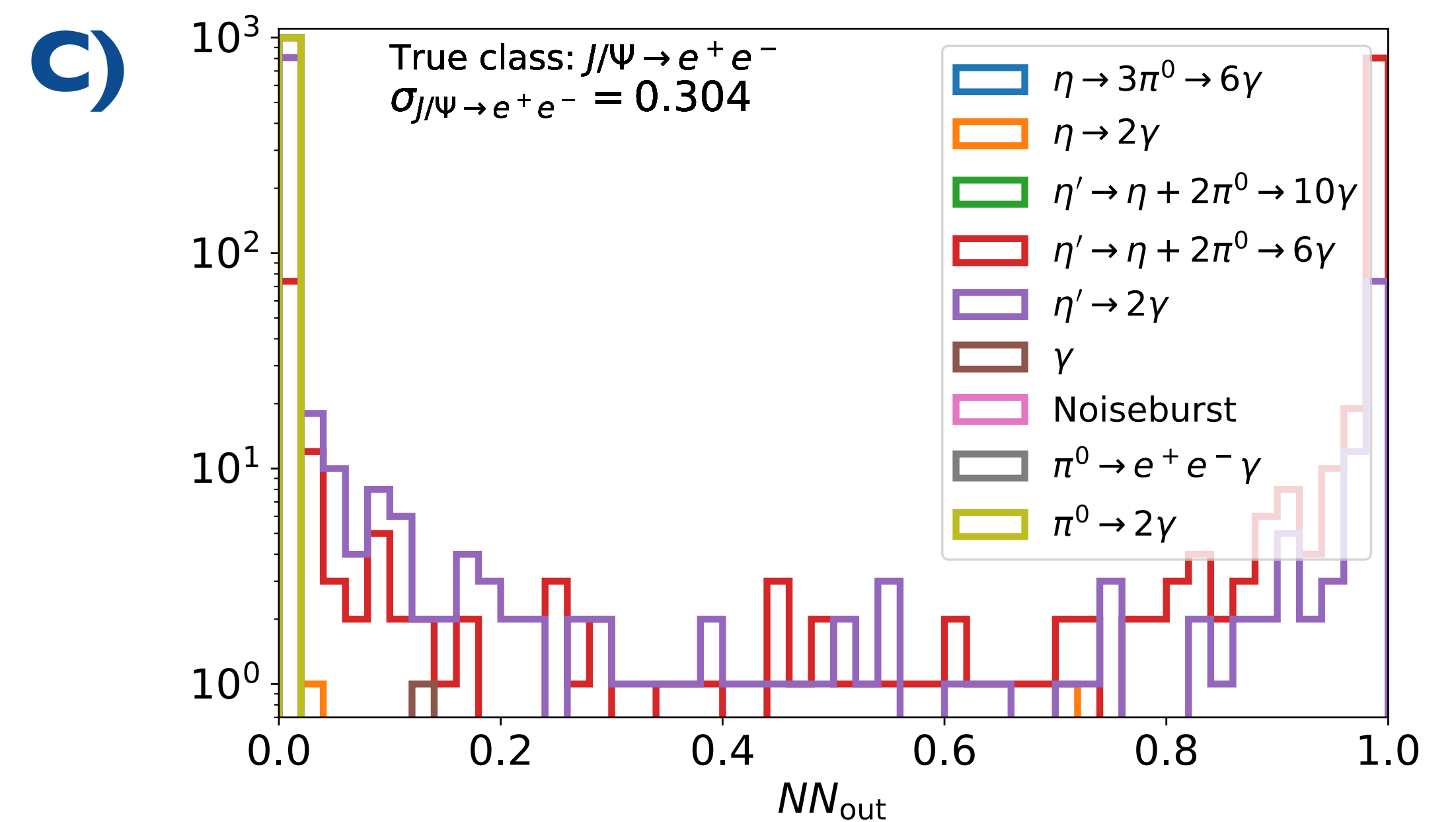
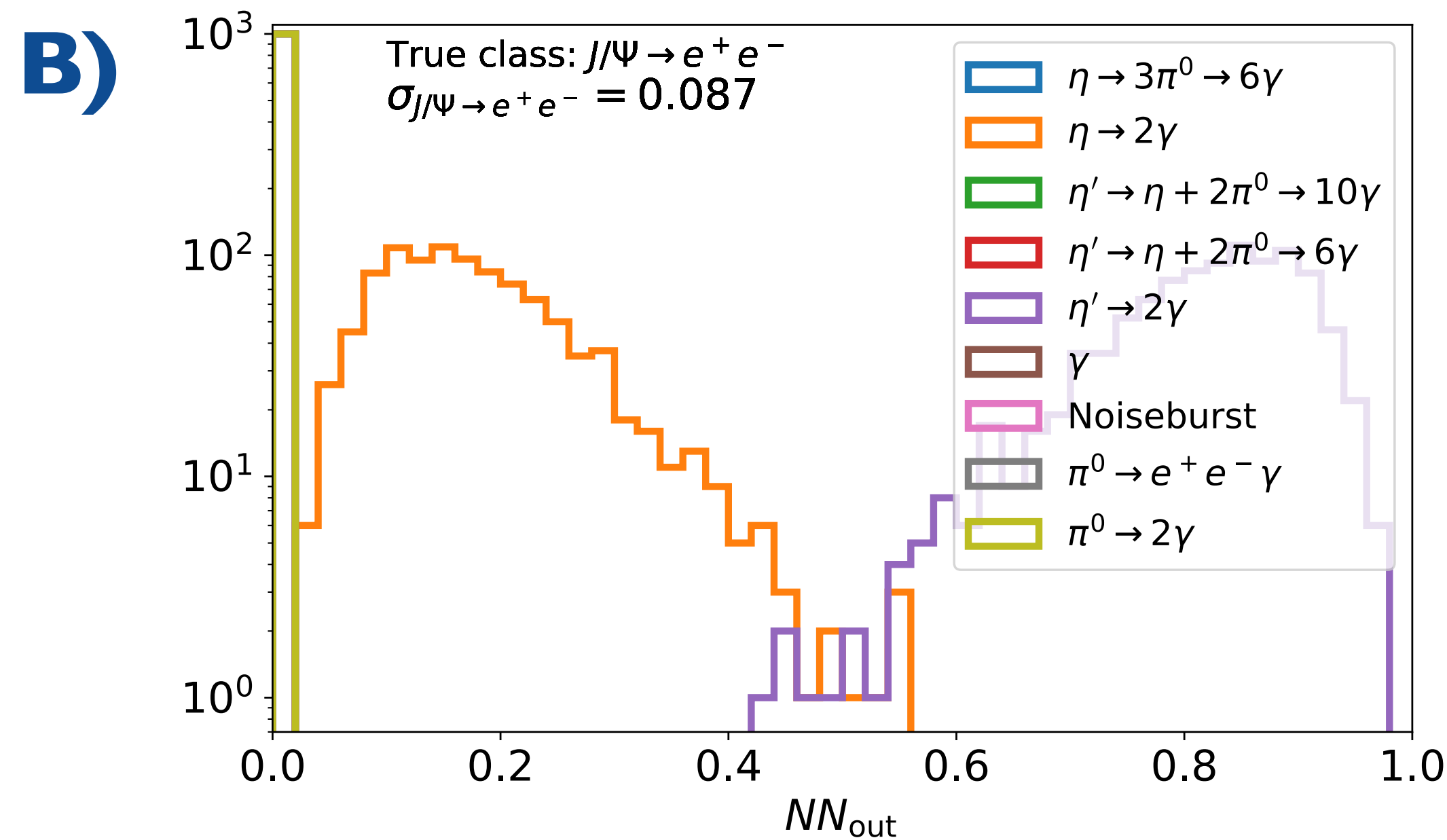
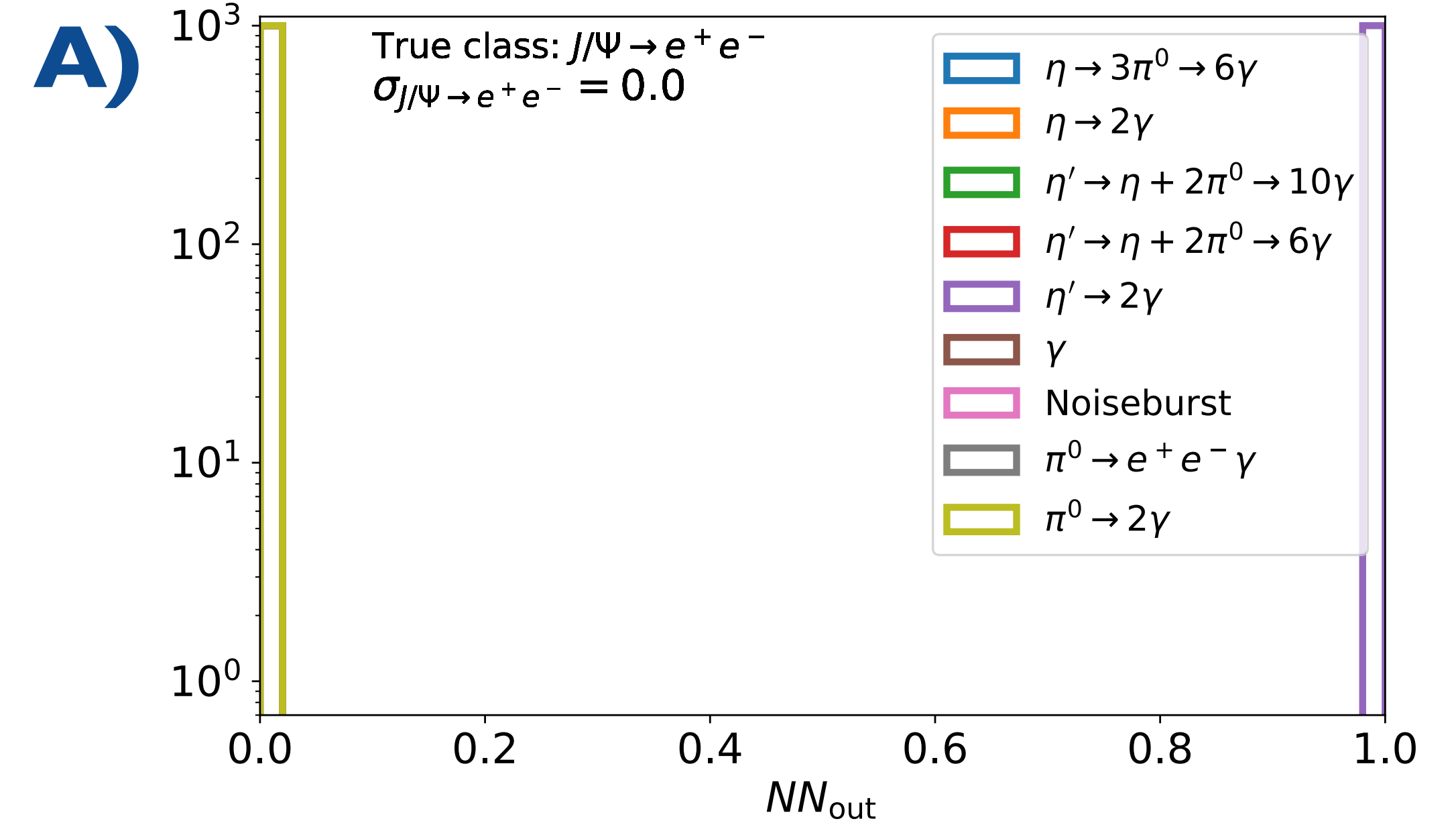
# Results

- In general, larger uncertainty than before
- Stems from three different cases:

A) Examples with  $\sim 0$  variance

B) Examples with  $> 1$  active output node

C) Examples with “jumpy decisions”



# Results

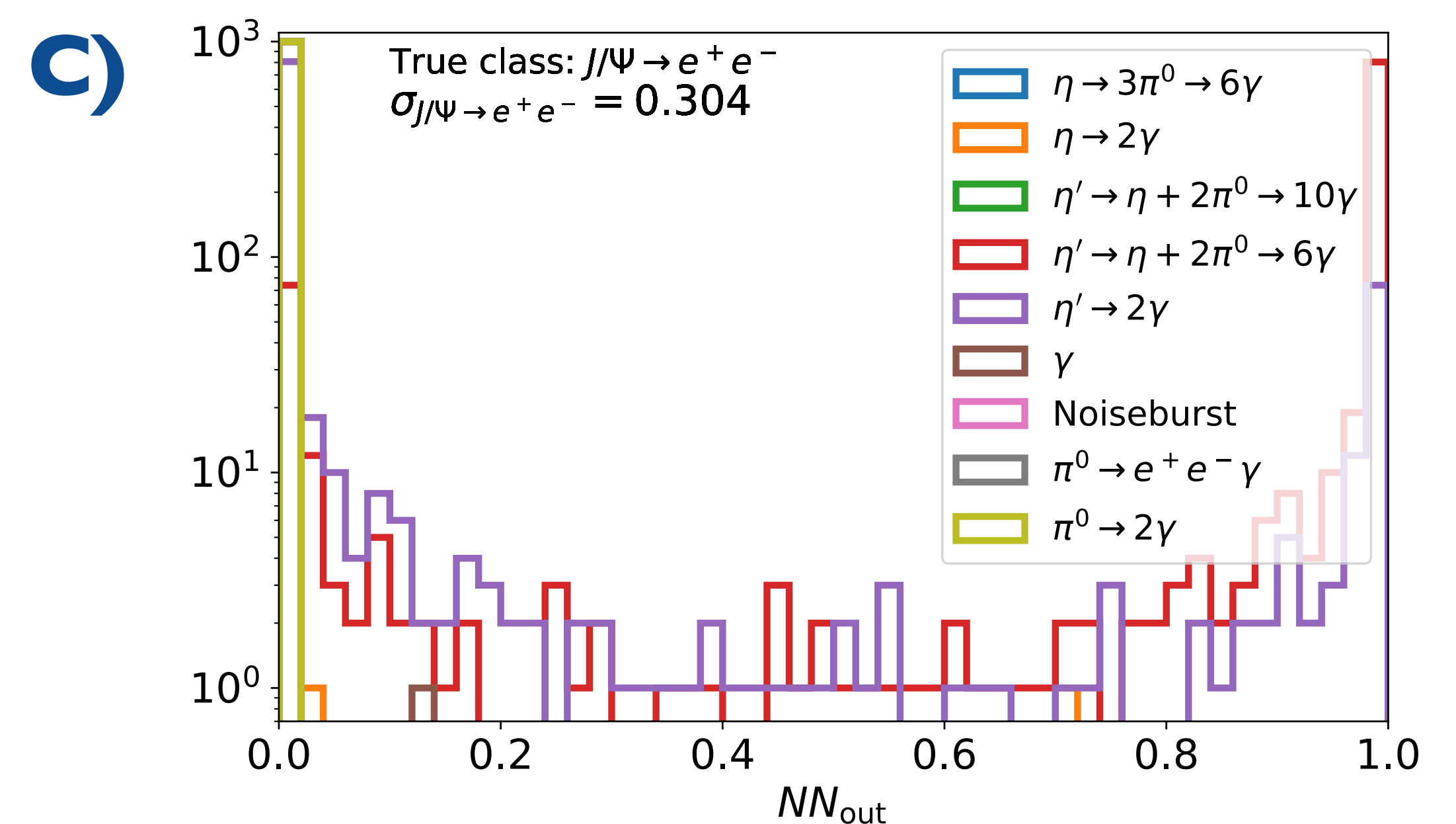
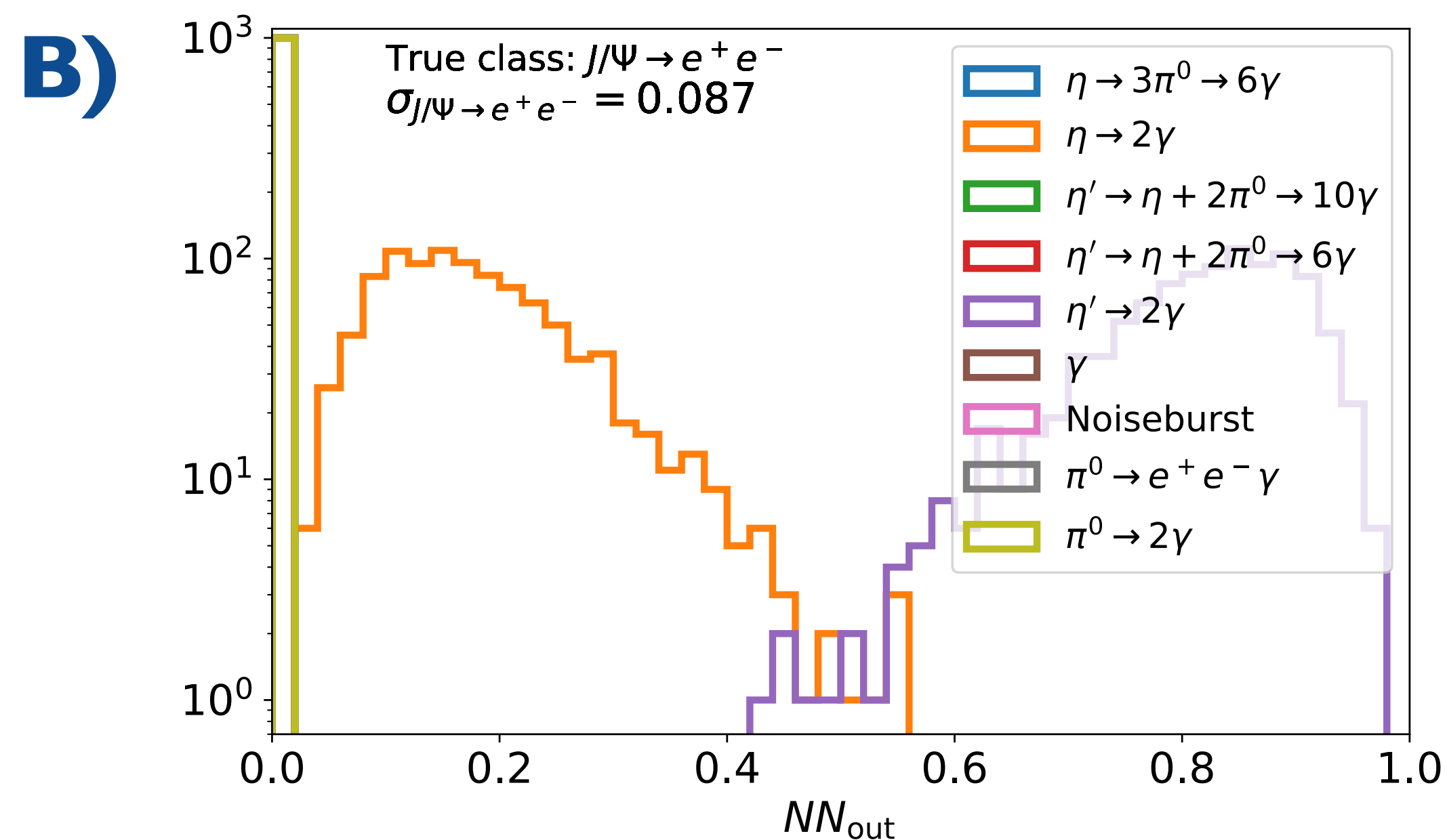
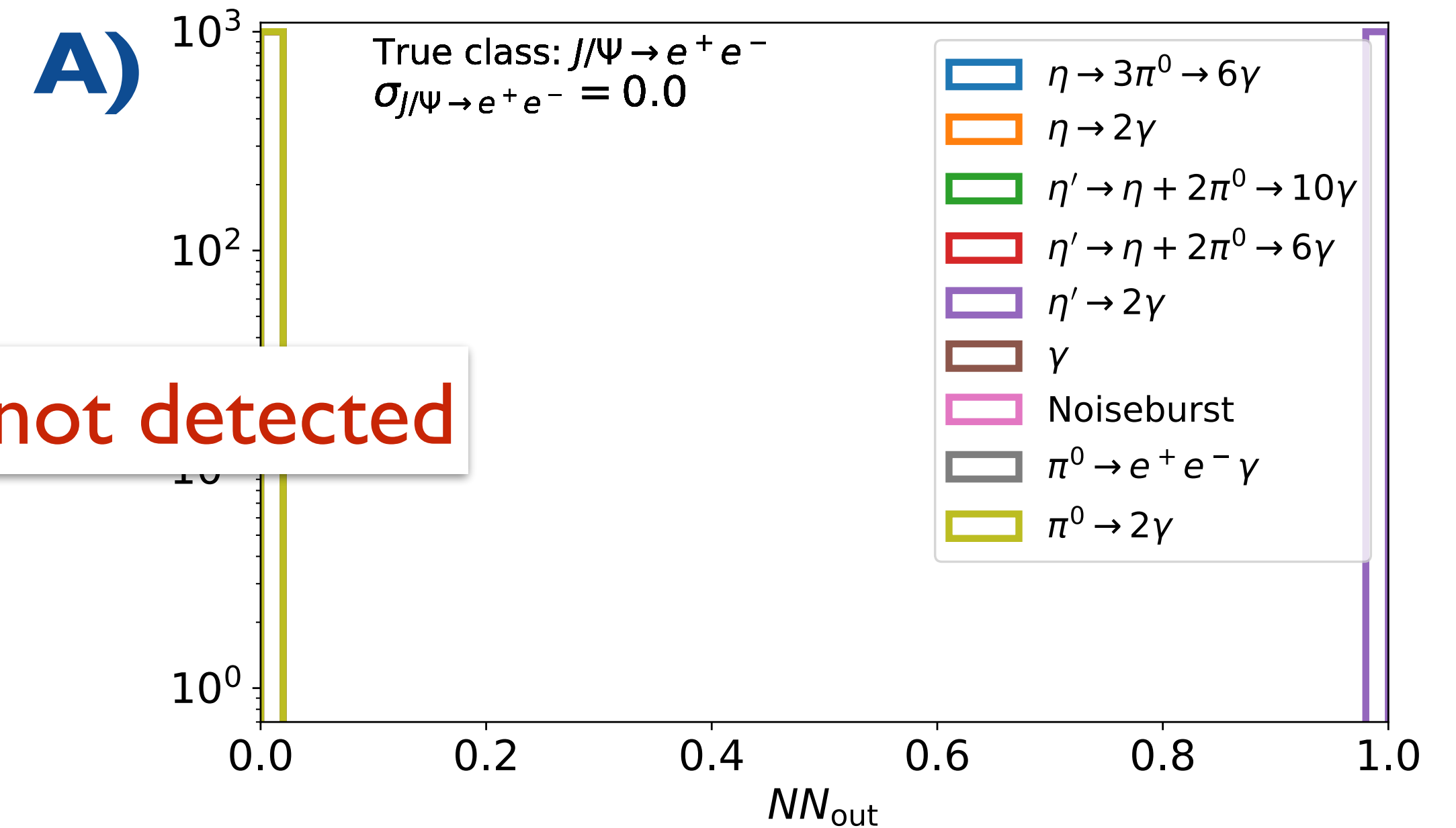
- In general, larger uncertainty than before
- Stems from three different cases:

A) Examples with  $\sim 0$  variance

← anomaly not detected

B) Examples with  $> 1$  active output node

C) Examples with “jumpy decisions”



# Results

- In general, larger uncertainty than before
- Stems from three different cases:

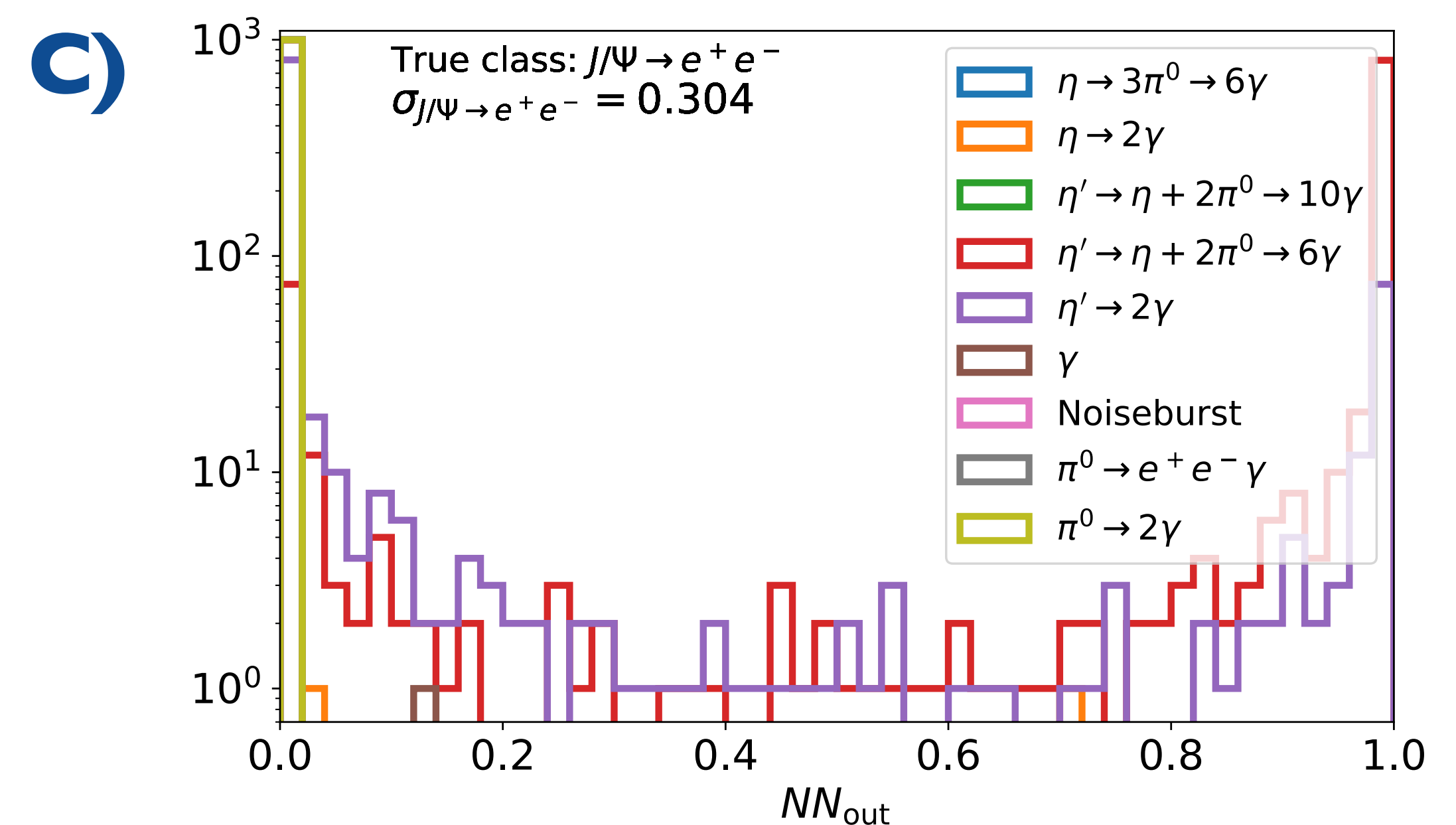
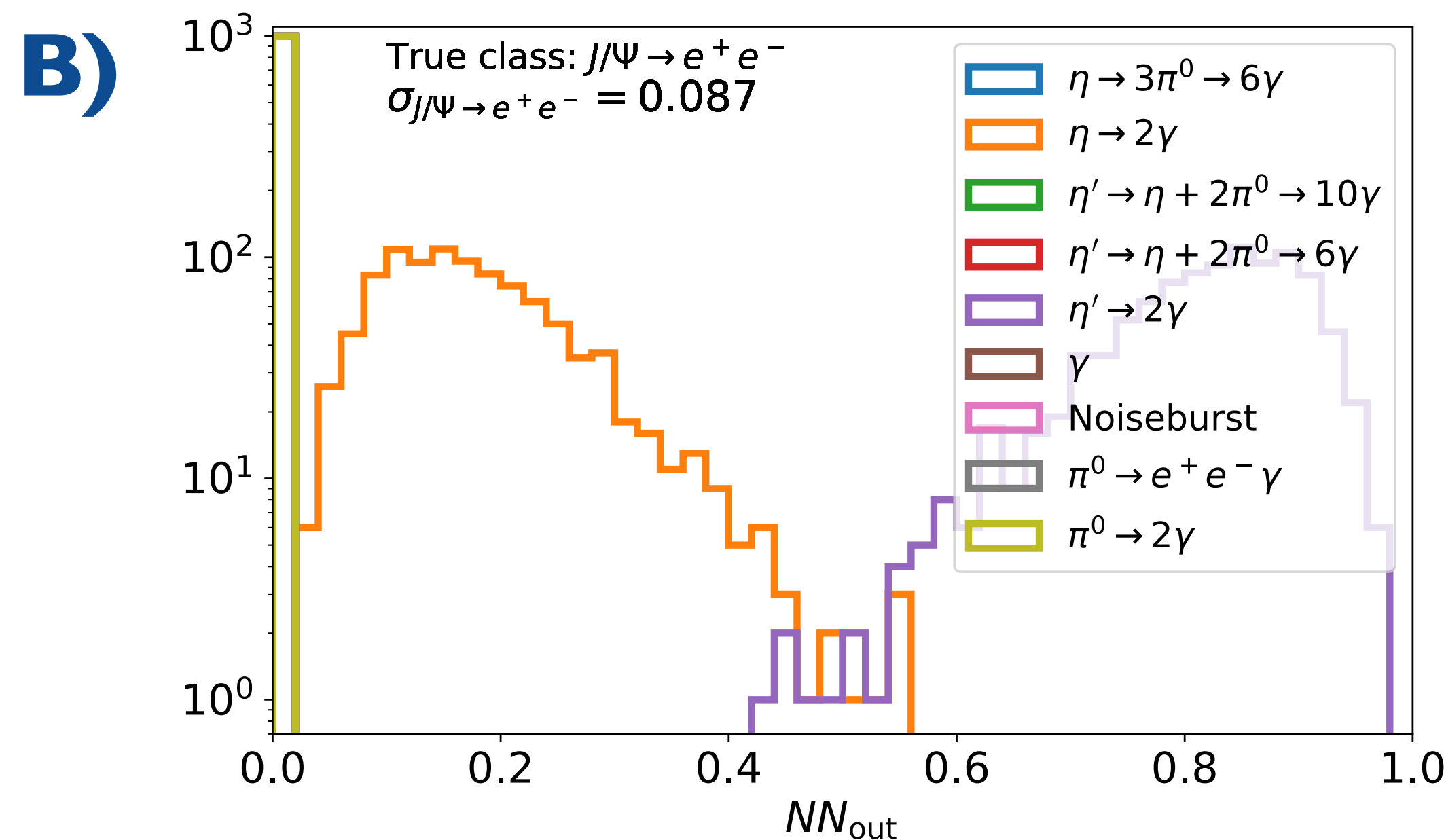
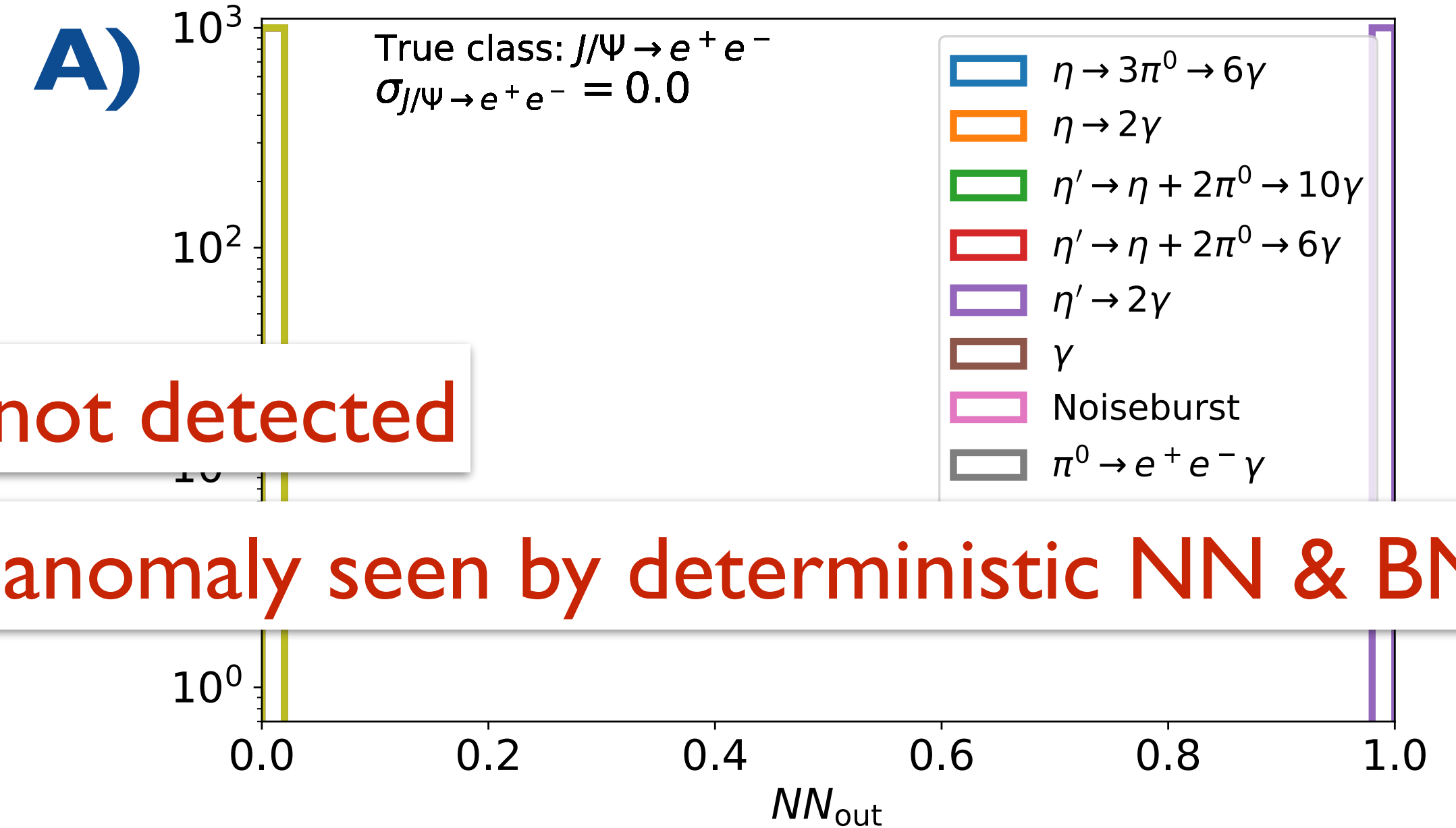
A) Examples with  $\sim 0$  variance

← anomaly not detected

B) Examples with  $> 1$  active output node

← anomaly seen by deterministic NN & BNN

C) Examples with “jumpy decisions”



# Results

- In general, larger uncertainty than before
- Stems from three different cases:

A) Examples with  $\sim 0$  variance

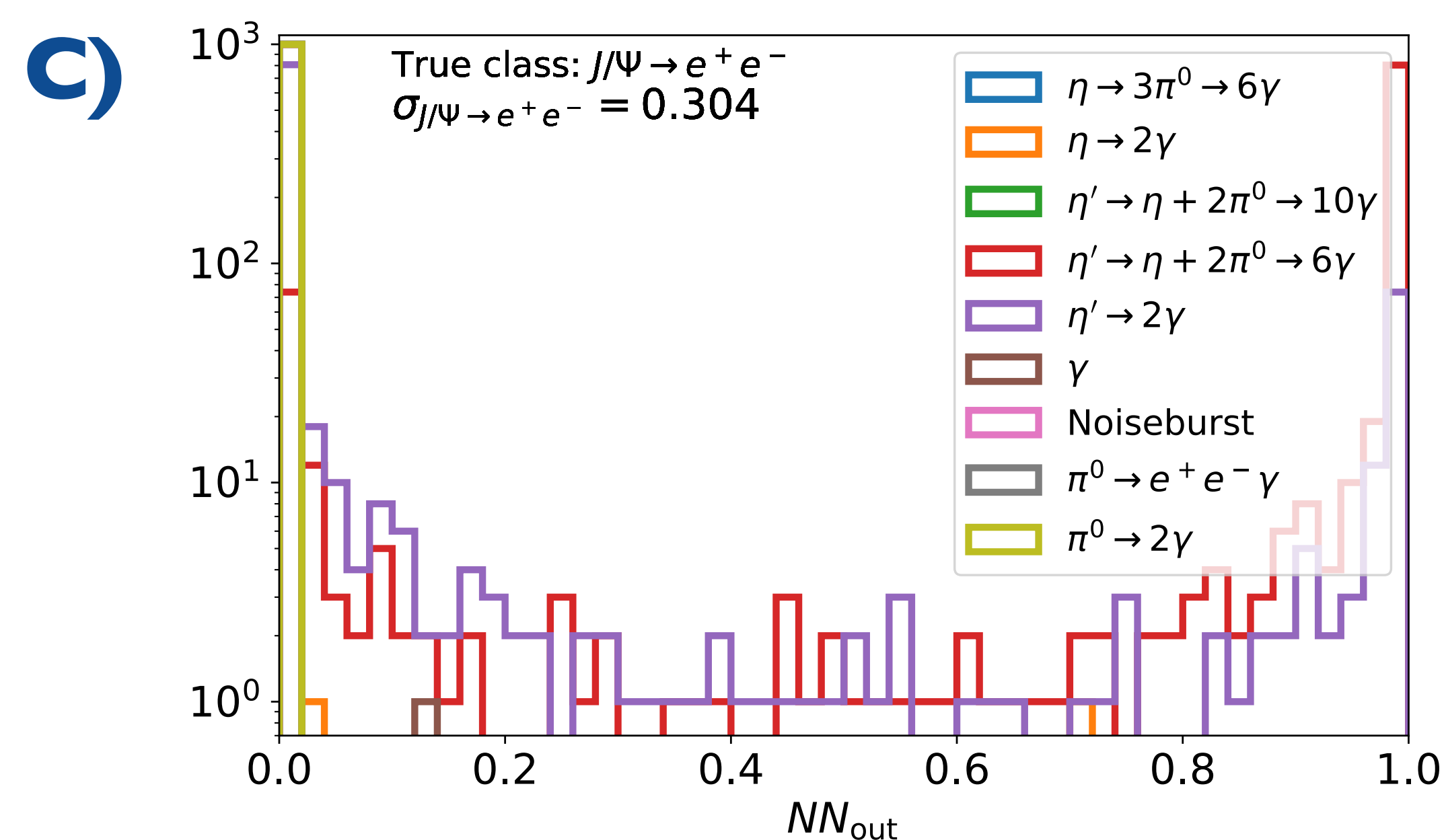
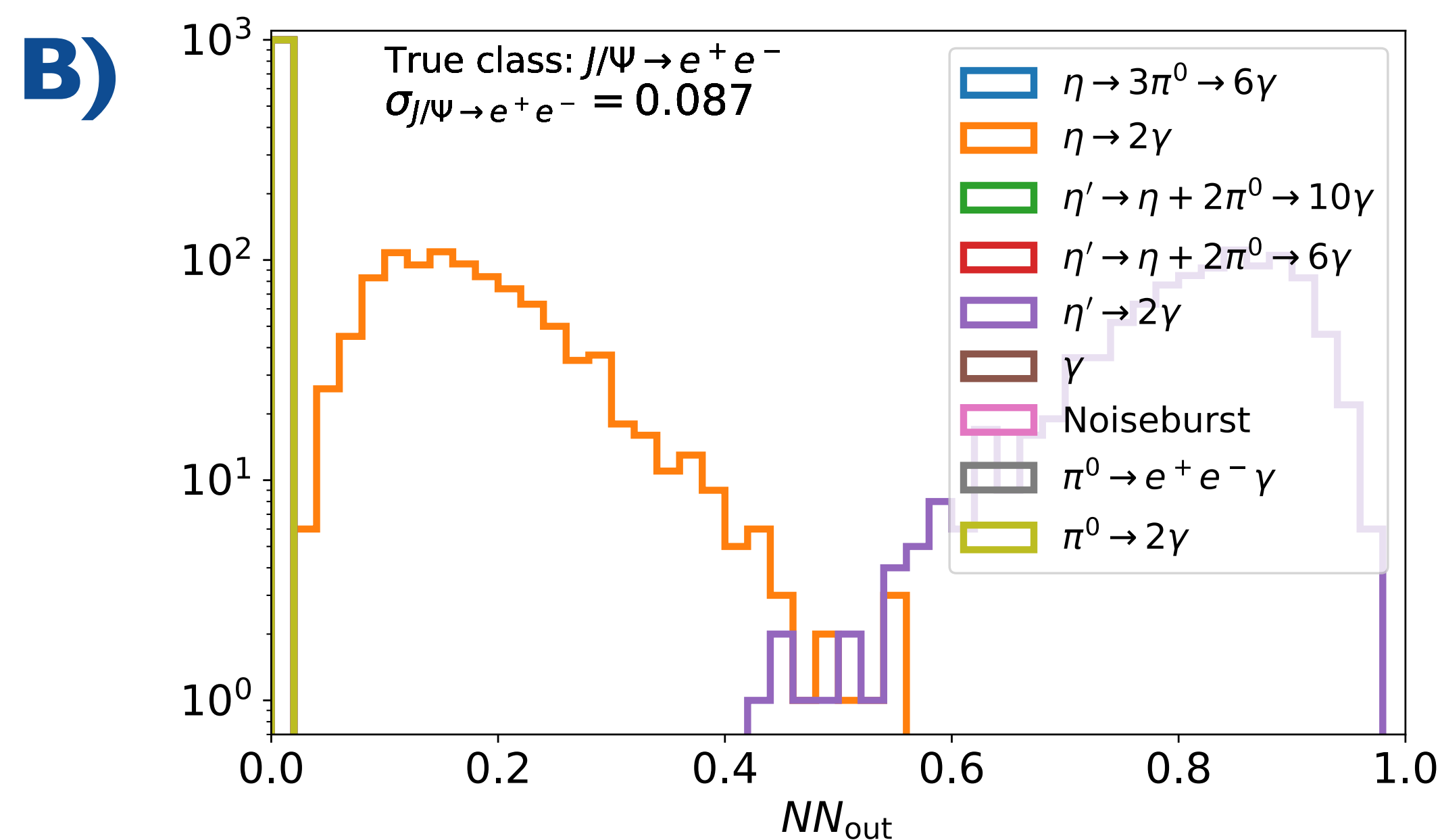
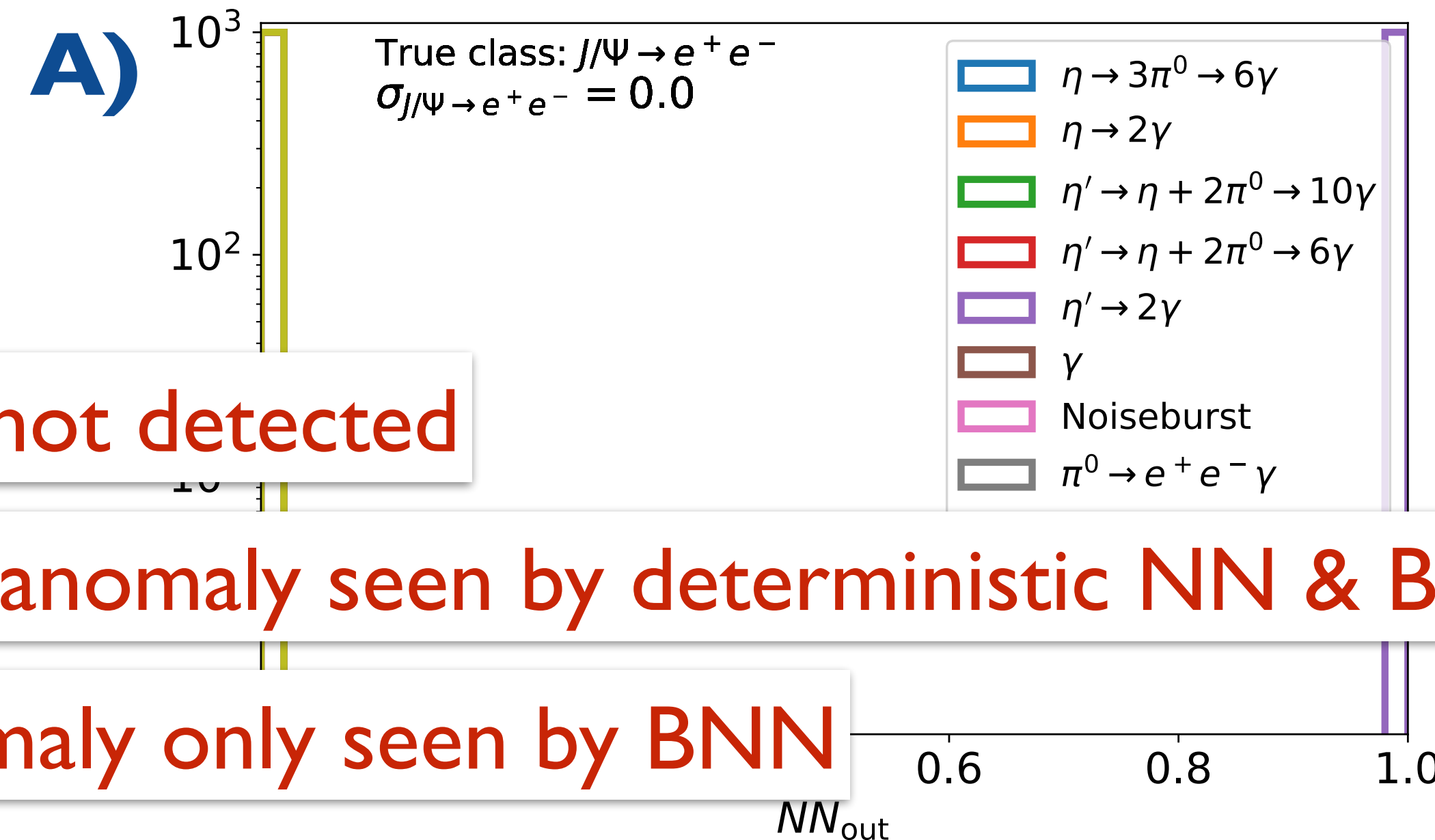
← anomaly not detected

B) Examples with  $> 1$  active output node

← anomaly seen by deterministic NN & BNN

C) Examples with “jumpy decisions”

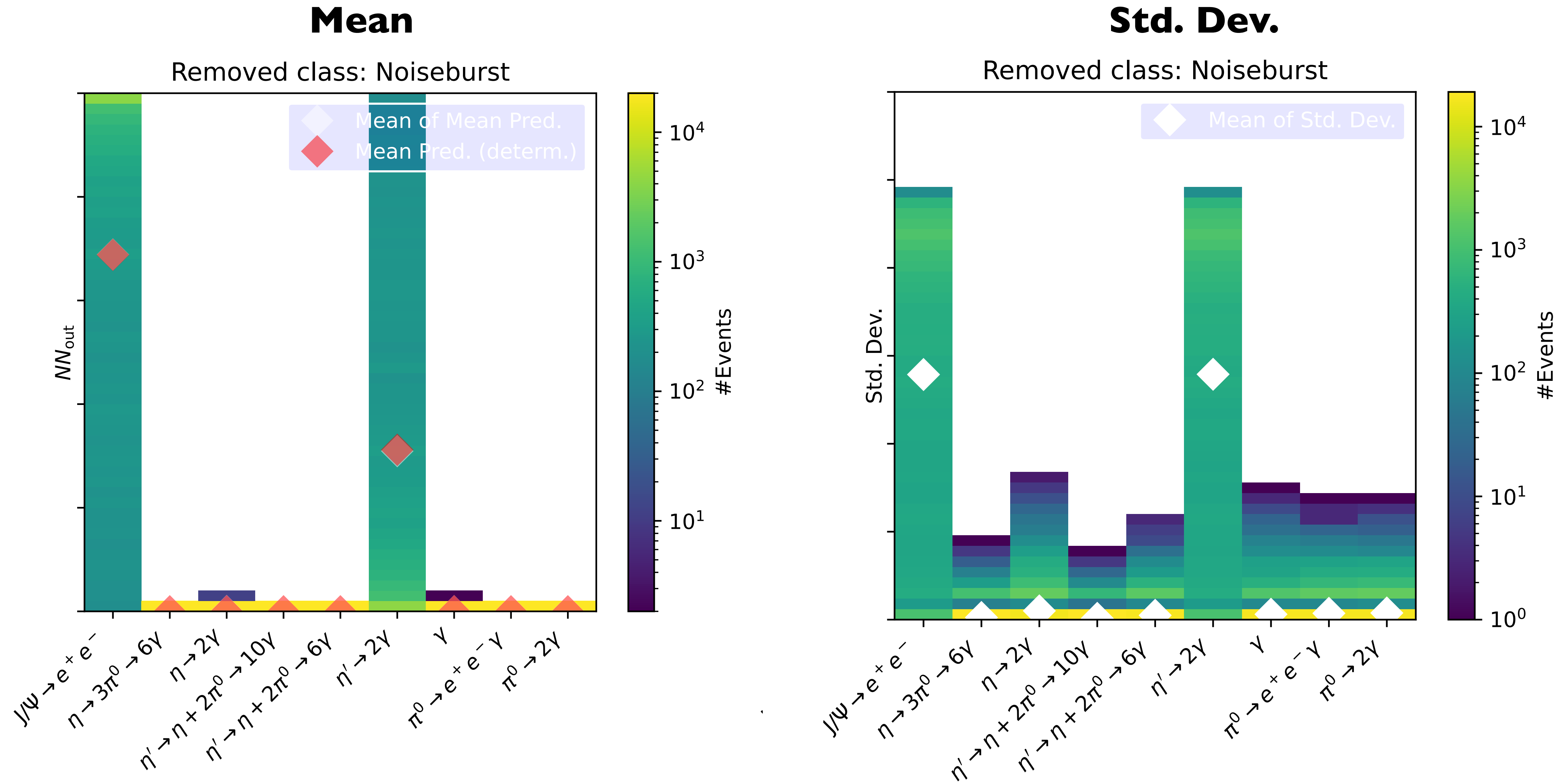
← anomaly only seen by BNN





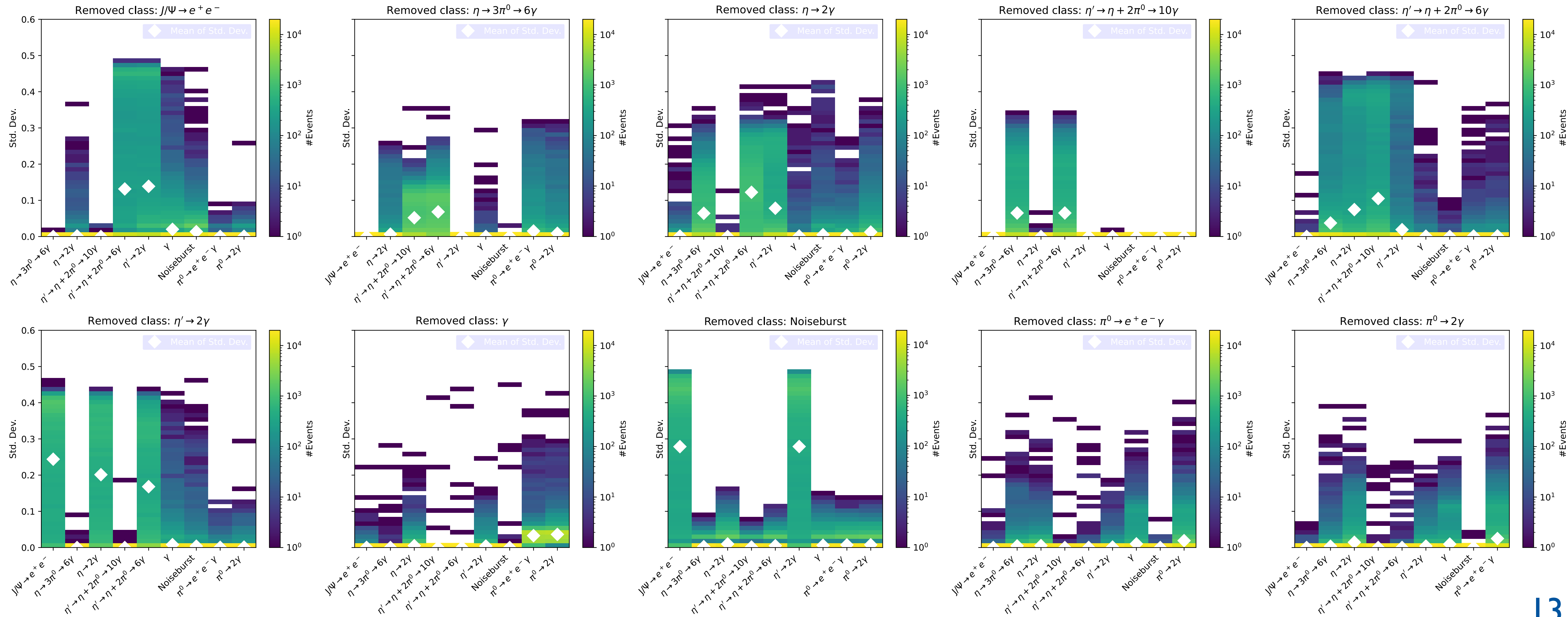
# Results

- Another example: noise burst as anomaly



# Results

- Standard deviations when each class is removed one-by-one
- Some anomalies are easier to identify than others



# Conclusions

- Bayesian network's uncertainty estimate may help to identify anomalies
- Semisupervised approach (LHC Olympics 2020, 2101.08320)
- We don't know, yet, if it can compete with specialized AD algorithms
- One advantage:  
Provides AD capabilities as a sanity/DQ cross check for standard classification tasks (such as ID)
- At some additional training and prediction costs

Section	Short Name	Method Type
3.1	VRNN	Unsupervised
3.2	ANODE	Unsupervised
3.3	BuHuLaSpa	Unsupervised
3.4	GAN-AE	Unsupervised
3.5	GIS	Unsupervised
3.6	LDA	Unsupervised
3.7	PGA	Unsupervised
3.8	Reg. Likelihoods	Unsupervised
3.9	UCluster	Unsupervised
4.1	CWoLa	Weakly Supervised
4.2	CWoLa AE Compare	Weakly/Unsupervised
4.3	Tag N' Train	Weakly Supervised
4.4	SALAD	Weakly Supervised
4.5	SA-CWoLa	Weakly Supervised
5.1	Deep Ensemble	Semisupervised
5.2	Factorized Topics	Semisupervised
5.3	QUAK	Semisupervised
5.4	LSTM	Semisupervised