# Flow-based networks and their benefits for us in high-energy physics

ERLANGEN CENTRE
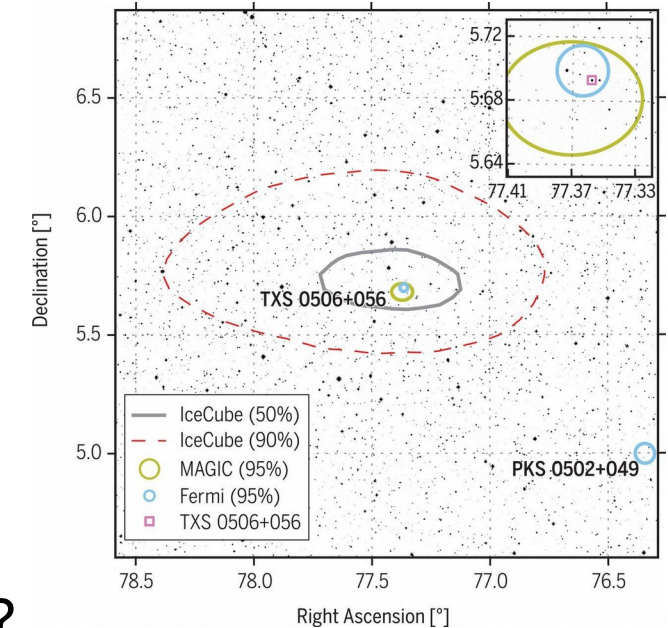FOR ASTROPARTICLE
PHYSICS

**Based on
arXiv:2008.05825**

Thorsten Glüsenkamp, ErUM-Data meeting, Sep. 22nd, 2020

# Overview

- Motivation
- Joint KL-divergence
- Flows generalize MSE
- Coverage
- Systematics
- Goodness of fit

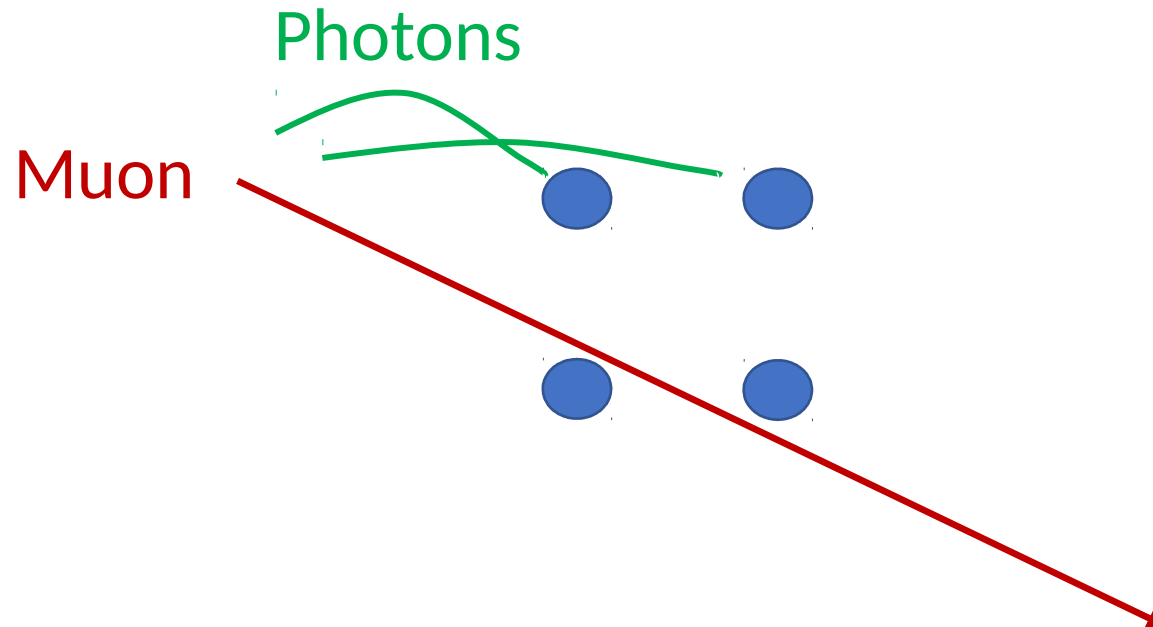# Motivation (here astronomical Posteriors)

- Standard recos: Coverage + correct systematics big issue (since years)
  (what about goodness-of-fit !? )



- Can we maybe solve all issues with neural networks?
  **Indeed, just using a simple upgrade of our existing networks,
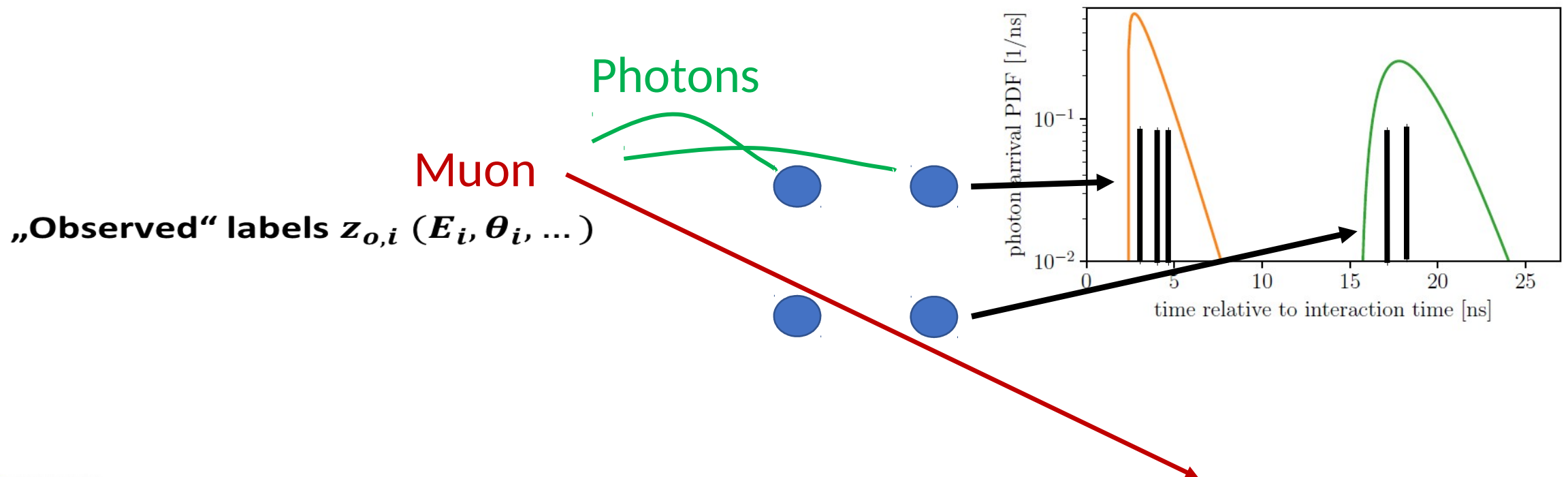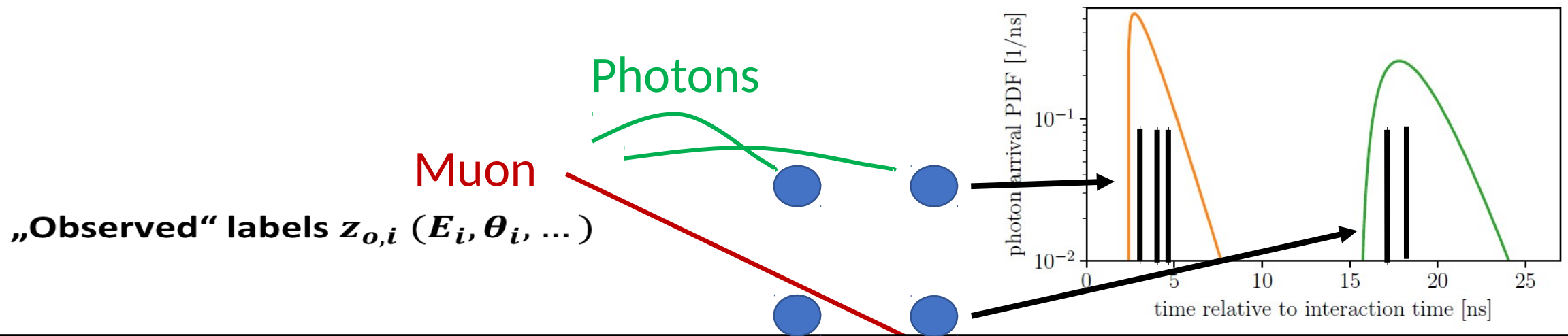  using so called „normalizing flows".**

# What is a Monte Carlo simulation?

An answer: **Samples $x_i, z_{o,i}$ from a „true" (intractable) joint distribution** $\mathcal{P}_t(x, z_o)$

Photons

Muon

# What is a Monte Carlo simulation?

- An answer: **Samples $x_i, z_{o,i}$ from a „true" (intractable) joint distribution $P_t(x, z_o)$**

Photons

Muon

„Observed" labels $z_{o,i}$ $(E_i, \theta_i, \dots)$



data $x_i$

# What is a Monte Carlo simulation?

- An answer: **Samples $x_i, z_{o,i}$ from a „true" (intractable) joint distribution $\mathcal{P}_t(x, z_o)$**

**Photons**

**Muon**

**data $x_i$**

„**Observed" labels $z_{o,i}$ $(E_i, \theta_i, \dots)$**



photon arrival PDF [1/ns]

time relative to interaction time [ns]

---

⟶ **Can evaluate expectation of arbitrary function $f(x, z_o)$ via samples!**

$$E_{x, z_o}[f(x, z_o)] = \int\limits_{x, z_o} \mathcal{P}_t(x, z_o) f(x, z_o) \, dx dz_o \approx \frac{1}{N} \sum_{x_i, z_{o,i}} f(x_i, z_{o,i})$$

# Supervised learning loss

$$f(x, z_o) = \ln \frac{P_t(x, z_o)}{q(x, z_o)}$$

$$\Rightarrow \quad E_{x, z_o}[f(x, z_o)] = D_{\text{KL,joint}(x, z_o)}(\mathcal{P}_t; q) = \int_x \int_{z_o} \mathcal{P}_t(z_o, x) \cdot \ln \frac{\mathcal{P}_t(z_o, x)}{q(z_o, x)} dz_o dx$$

# Supervised learning loss

- A particular choice of $f(x, z_o)$ yields the loss fuction in supervised learning!

$$f(x, z_o) = \ln \frac{P_t(x, z_o)}{q(x, z_o)}$$

➡ $$E_{x,z_o}[f(x, z_o)] = D_{\text{KL,joint}(x,z_o)}(\mathcal{P}_t; q) = \int_x \int_{z_o} \mathcal{P}_t(z_o, x) \cdot \ln \frac{\mathcal{P}_t(z_o, x)}{q(z_o, x)} dz_o dx$$

**Use samples, parametrize $q(z_o; x)$ with neural network as $q_\phi$, minimize result over $\phi$**

➡ $$\arg\min_\phi \hat{D}_{\text{KL,joint}(x,z_o)}(\mathcal{P}_t; q_\phi) = \ldots = \arg\min_\phi \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln\left(q_\phi(z_{o,i}; x_i)\right)$$

# Supervised learning loss

$$f(x, z_o) = \ln \frac{P_t(x, z_o)}{q(x, z_o)}$$

➡ $$E_{x,z_o}[f(x, z_o)] = D_{\text{KL,joint(x,z_o)}}(\mathcal{P}_t; q) = \int_x \int_{z_o} \mathcal{P}_t(z_o, x) \cdot \ln \frac{\mathcal{P}_t(z_o, x)}{q(z_o, x)} dz_o dx$$

**Use samples, parametrize $q(z_o; x)$ with neural network as $q_\phi$, minimize result over $\phi$**

➡ $$\arg\min_\phi \hat{D}_{\text{KL,joint(x,z_o)}}(\mathcal{P}_t; q_\phi) = \ldots = \arg\min_\phi \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln(q_\phi(z_{o,i}; x_i))$$

**MSE-Loss:** $\sum (\mu_\phi(x_i) - z_{o,i})^2$ ⟵

$q_\phi = N(\mu; 1)$ (standard Normal)

9

# Meaning of the KL-divergence viewpoint

$$\arg\min_{\phi} \hat{D}_{\text{KL,joint}(x,z_o)}(\mathcal{P}_t; q_\phi) = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} \ln\left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)}\right) + \ln\left(\frac{\mathcal{P}_t(x_i)}{q(x_i)}\right)$$

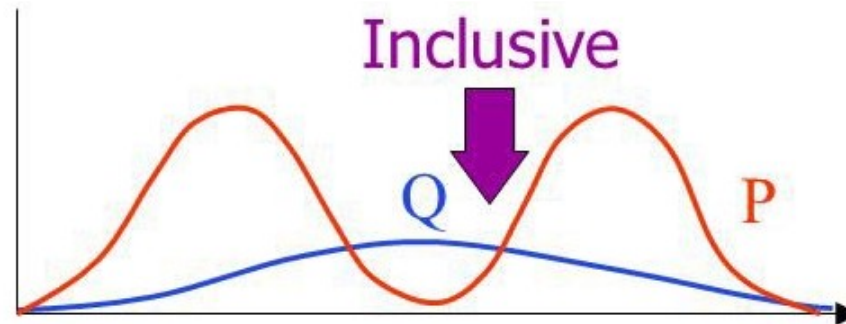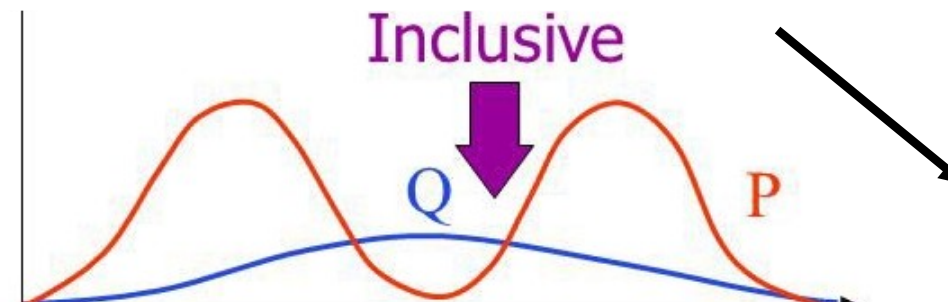$$\dots = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln\left(q_\phi(z_{o,i}; x_i)\right)$$

# Meaning of the KL-divergence viewpoint

$$\arg\min_{\phi} \hat{D}_{\text{KL,joint}(\mathbf{x},\mathbf{z}_o)}(\mathcal{P}_t; q_\phi) = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} \ln\left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)}\right) + \ln\left(\frac{\mathcal{P}_t(x_i)}{q(x_i)}\right)$$

$$\dots = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln\left(q_\phi(z_{o,i}; x_i)\right)$$

„Minimizing KL-divergence between **True Posterior P_t** and **approximate Posterior $q_\phi$ over $\phi$** = minimizing supervised learning loss"

Minimising

KL$(P||Q)$
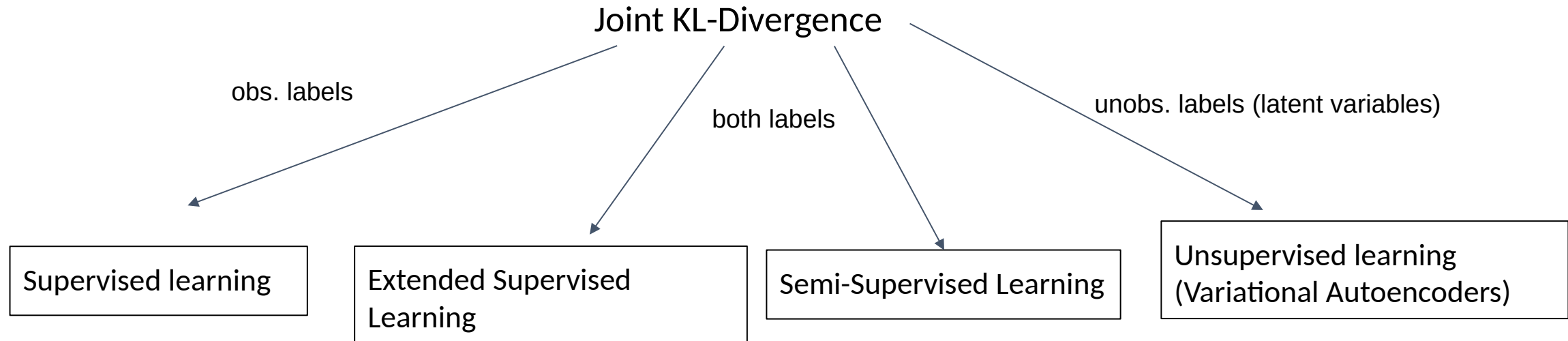
$$= \sum_H P(H|V) \ln \frac{P(H|V)}{Q(H)}$$



Inclusive

Q          P

# Meaning of the KL-divergence viewpoint

$$\arg\min_{\phi} \hat{D}_{\mathrm{KL,joint(x,z_o)}}(\mathcal{P}_t; q_\phi) = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} \ln\left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)}\right) + \ln\left(\frac{\mathcal{P}_t(x_i)}{q(x_i)}\right)$$

$$\ldots = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln\left(q_\phi(z_{o,i}; x_i)\right)$$

„Minimizing KL-divergence between **True Posterior P_t** and **approximate Posterior $q_\phi$ over $\phi$** = minimizing supervised learning loss"

Minimising

$\mathrm{KL}(P\|Q)$

$$= \sum_H P(H|V)\ln\frac{P(H|V)}{Q(H)}$$

Inclusive

Q        P

**! This is a VERY useful viewpoint for us in physics !**

**Standard supervised learning is performing approximate Likelihood-free inference**

**„Neural networks learn to approximate the true posterior"**

# This viewpoint unifies various approaches!

Joint KL-Divergence

obs. labels

both labels

unobs. labels (latent variables)

| Supervised learning | Extended Supervised Learning | Semi-Supervised Learning | Unsupervised learning (Variational Autoencoders) |

# This viewpoint unifies various approaches!

Joint KL-Divergence

obs. labels

both labels

unobs. labels (latent variables)

| Supervised learning | Extended Supervised Learning | Semi-Supervised Learning | Unsupervised learning (Variational Autoencoders) |

classification

regression

**Systematics**

**Goodness-of-Fit**

**Normalizing Flows**

**Coverage of posterior regions**

$q_\Phi\left(z_o; x\right)$

1) Data Encoding
- CNN
- LSTMs
- Graph NN
- ...

2) PDF description

# This viewpoint unifies various approaches!

Joint KL-Divergence

obs. labels

both labels

unobs. labels (latent variables)

Supervised learning

Extended Supervised Learning

Semi-Supervised Learning

Unsupervised learning (Variational Autoencoders)

**This talk**

classification

regression

**Systematics**

**Goodness-of-Fit**

**Normalizing Flows**

**Coverage of posterior regions**



1) Data Encoding
- CNN
- LSTMs
- Graph NN
- ...

$q_\Phi(z_o; x)$

2) PDF description

# What are good PDF Aproximators $q_\phi$ ?

Predicting parameters of any complex distribution? E.g. a sum of gaussians? Possible, but there is something better …

**Normalizing flows: (**1912.02762 for a recent review)

# What are good PDF Aproximators $q_\phi$ ?

Predicting parameters of any complex distribution? E.g. a sum of gaussians? Possible, but there is something better …

**Normalizing flows: (**1912.02762 for a recent review)

**Generic Flow:**

Invertible mapping (inverse required for density evaluation)



Base distribution
(typically standard normal)

Final distribution

Parameters of mapping
are output of Neural Network (here conditional on x)

# What are good PDF Aproximators $q_\phi$ ?

Predicting parameters of any complex distribution? E.g. a sum of gaussians? Possible, but there is something better ...

**Normalizing flows: (**1912.02762 for a recent review)

**Generic Flow:**

Invertible mapping (inverse required for density evaluation)

**A particular Flow:**

A Gaussian PDF is an affine normalizing flow:



Base distribution
(typically standard normal)

Final distribution

Parameters of mapping
are output of Neural Network (here conditional on x)



Final Gaussian distribution

Generalizes MSE loss ...
MSE loss corresponds
to an affine flow with no scaling!

18

# Example Posteriors



Complex flow

Affine flow
(MSE + sigma)

# Example Posteriors

# Example Posteriors



**Data encoding is bottleneck, not the flow!**

# Coverage



- • Can calculate coverage of arbitrary PDF at the base using standard $\chi^2$- test
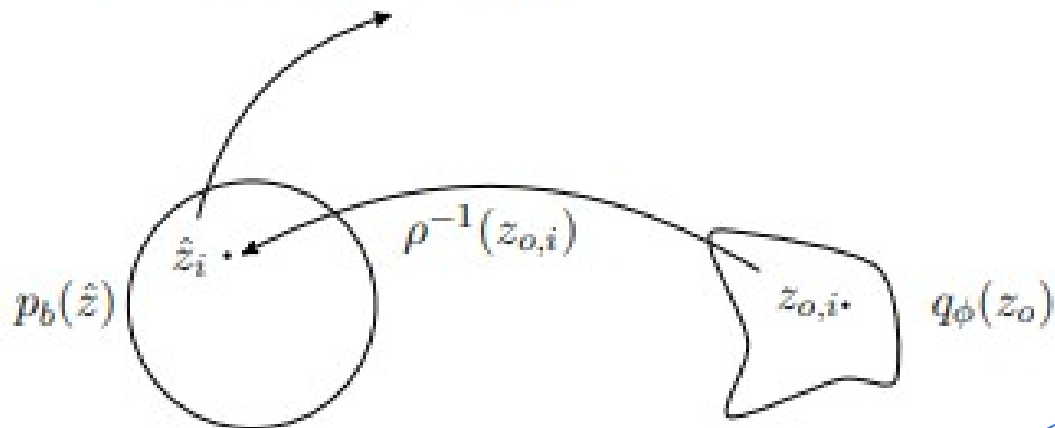
$$-2 \cdot \left(\ln p_b(\hat{z}_i) - \ln p_b(0)\right) \sim \chi^2$$



(a) Coverage of 3-d posteriors using dataset 3 for different stages of training.
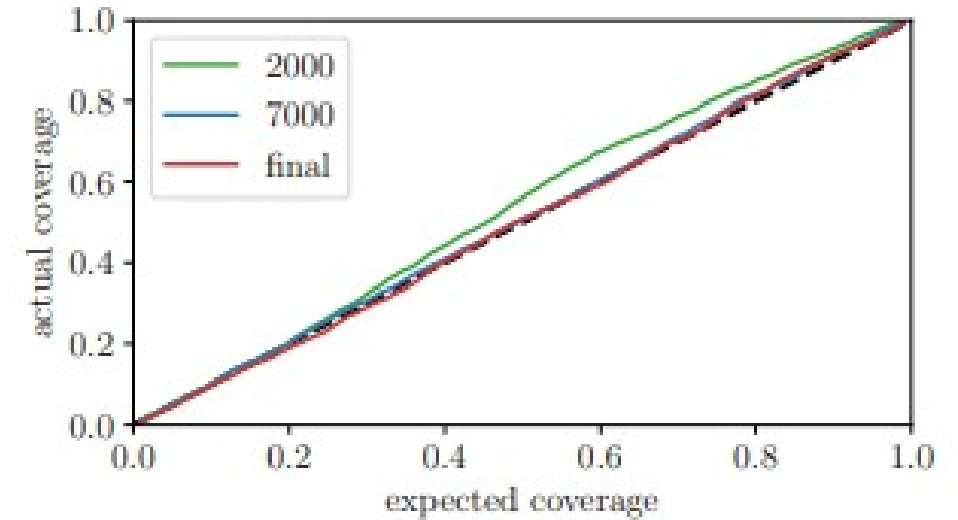
# Coverage

- Can calculate coverage of arbitrary PDF at the base using standard $\chi^2$- test

$$-2 \cdot \left(\ln p_b(\hat{z}_i) - \ln p_b(0)\right) \sim \chi^2$$



(a) Coverage of 3-d posteriors using dataset 3 for different stages of training.

**Rapid initial training phase to obtain coverage**

**Slow training phase (diffusive phase) shrinks posterior regions while maintaining coverage!**

# Coverage



**Also works for arb. posteriors of directions (on spheres)**



$$\rho_{\Psi,\text{intrinsic}}$$

$$\frac{K(\theta_1,\ldots,\theta_{d-2})}{S_d} \cdot \sin^{d-1}(\theta_{d-1})$$

$$q_\Psi(\theta_1,\ldots,\theta_{d-1},\phi_d)$$

$$\rho_{\text{tot}} = \rho_2 \circ \rho_1$$

$$\theta_{d-1}$$

$$\rho_2(r_f) = \theta_{d-1} = \arccos\left(\frac{r_f^2-1}{r_f^2+1}\right)$$

$$\frac{K(\theta_1,\ldots,\theta_{d-2})}{(2\pi)^{d/2}} \cdot r_g^{d-1} \cdot \exp\left(-\frac{r_g^2}{2}\right)$$

$$\frac{K(\theta_1,\ldots,\theta_{d-2})}{S_d} \cdot \left(\frac{2}{r_f^2+1}\right)^d \cdot r_f^{d-1}$$

$$\rho_1(r_g) = r_f = \text{CDF}_{r,f}^{-1}\left(\text{CDF}_{r,g}(r_g)\right)$$

(a) Coverage of 3-d posteriors using dataset 3 for different stages of training.

**Rapid initial training phase to obtain coverage**

**Slow training phase (diffusive phase) shrinks posterior regions while maintaining coverage!**

24

Red: True Label      White: True Posterior 68% region      Black: Predicted Posterior 68% region

**Red: True Label**   **White: True Posterior 68% region**   **Black: Predicted Posterior 68% region**

**White: True Posterior 68% region**   **Black: Predicted Posterior 68% region**
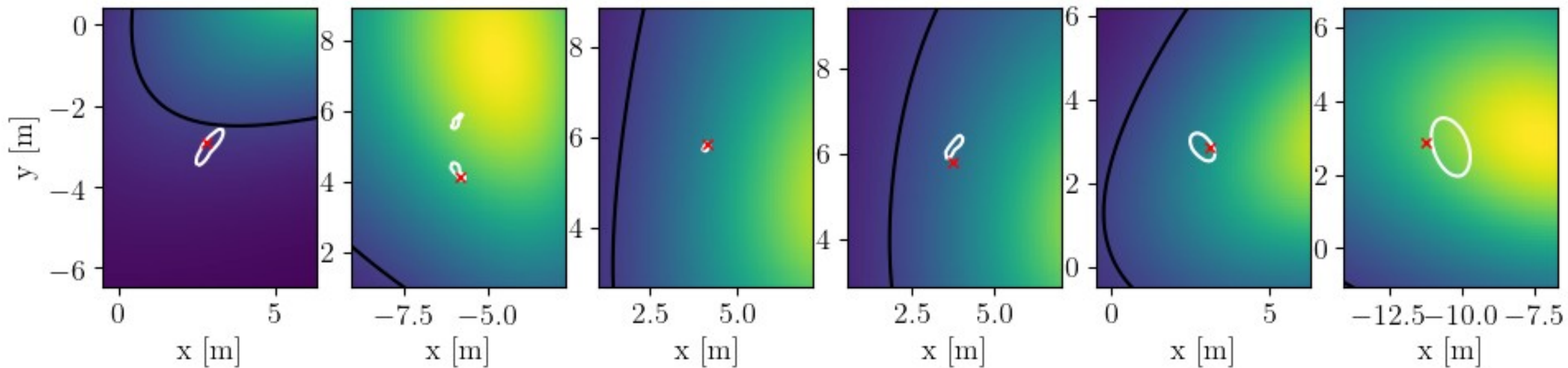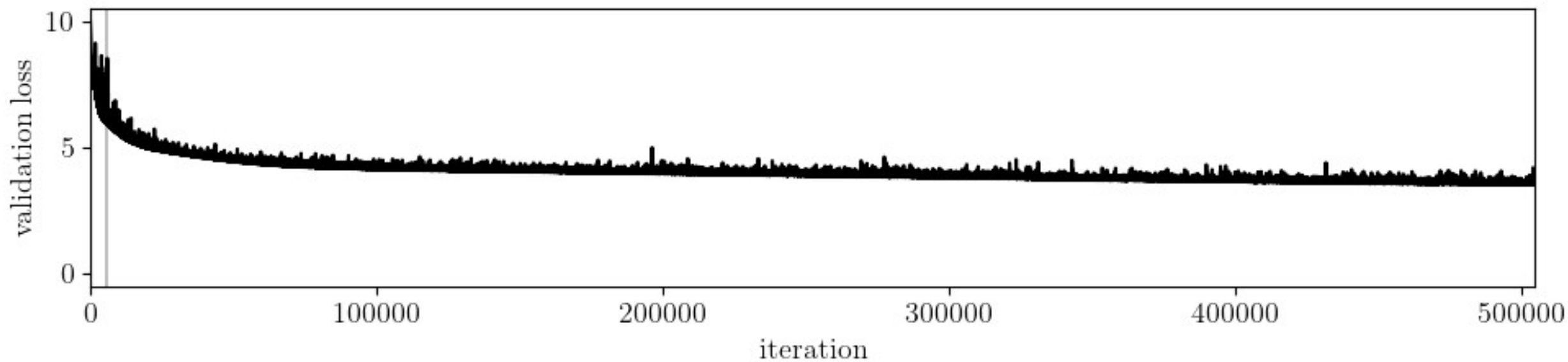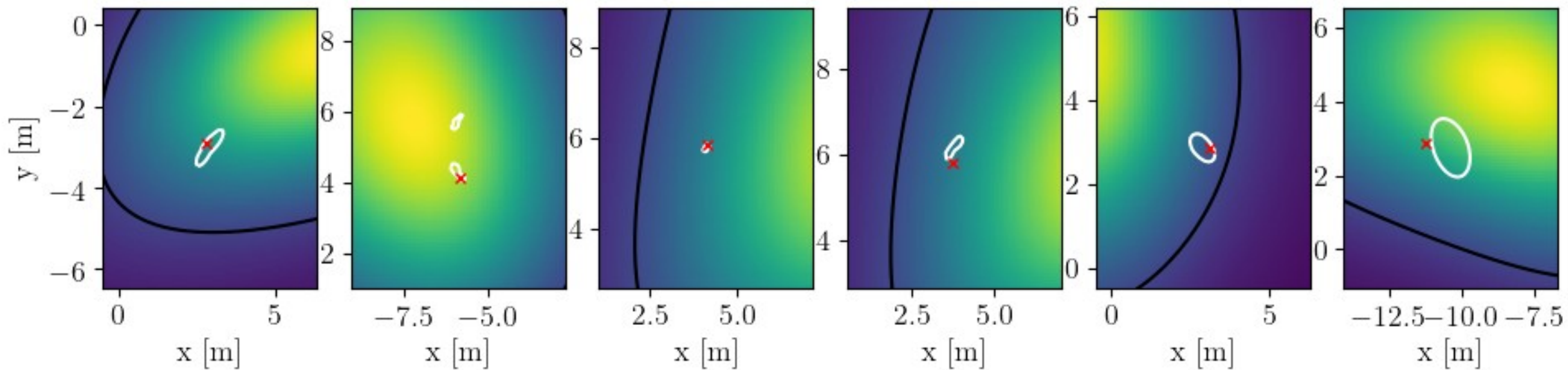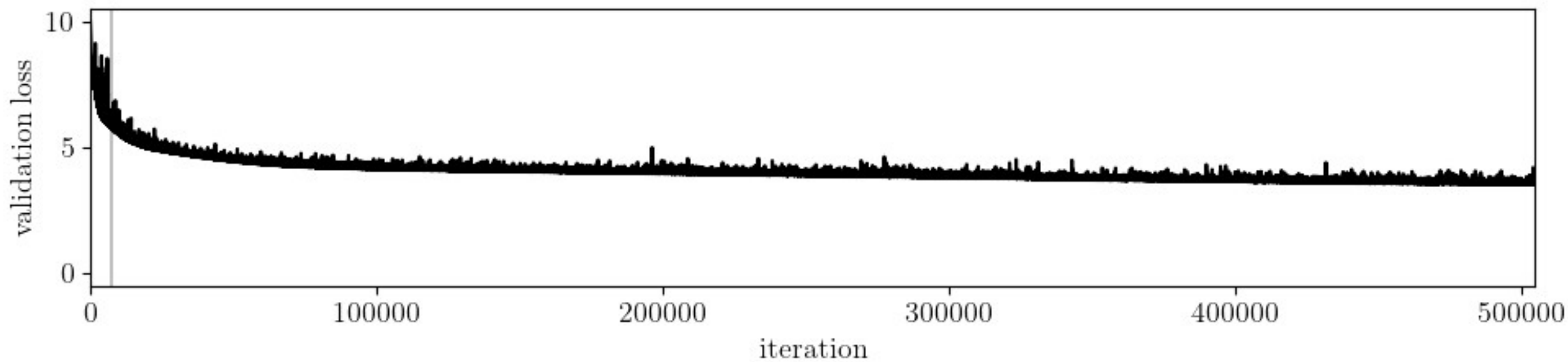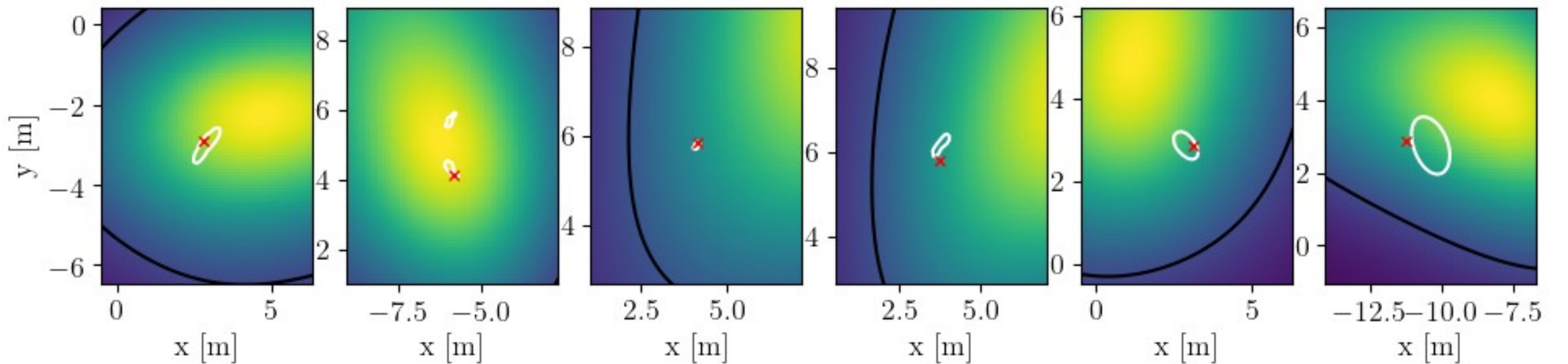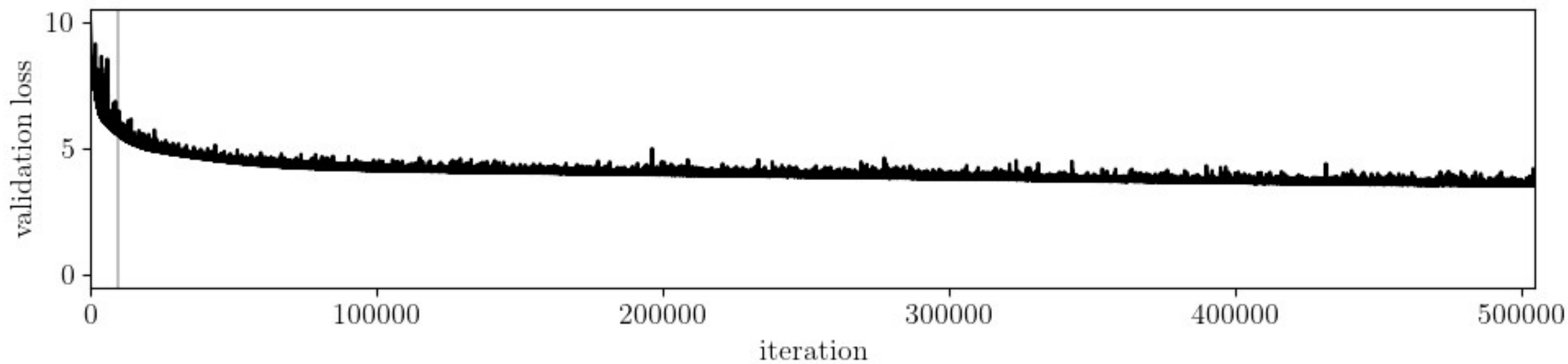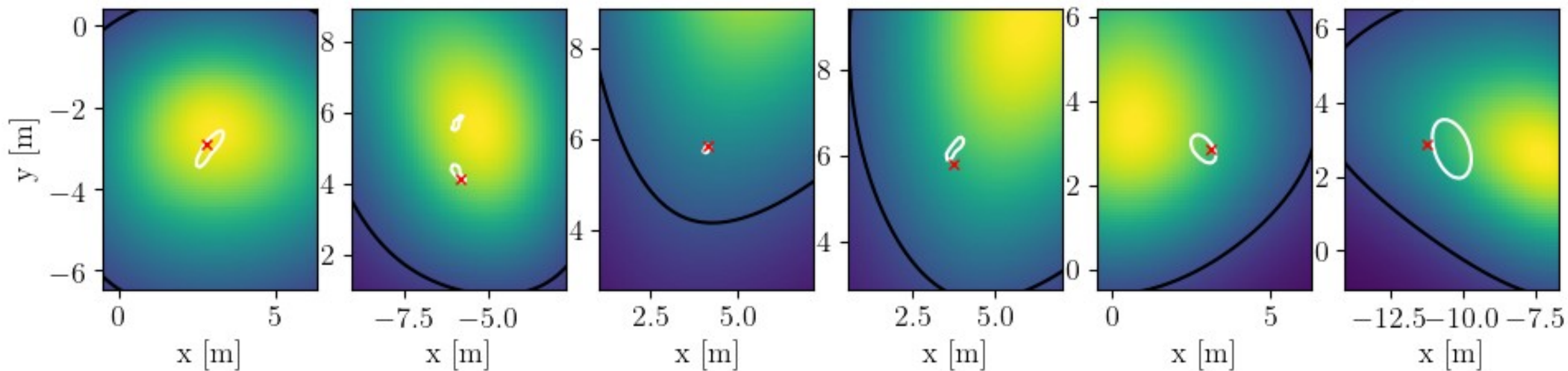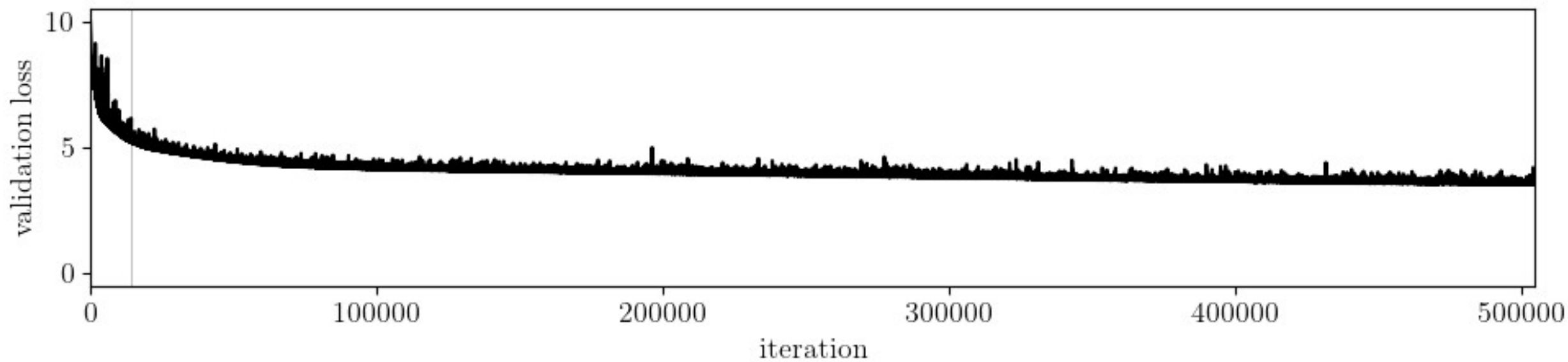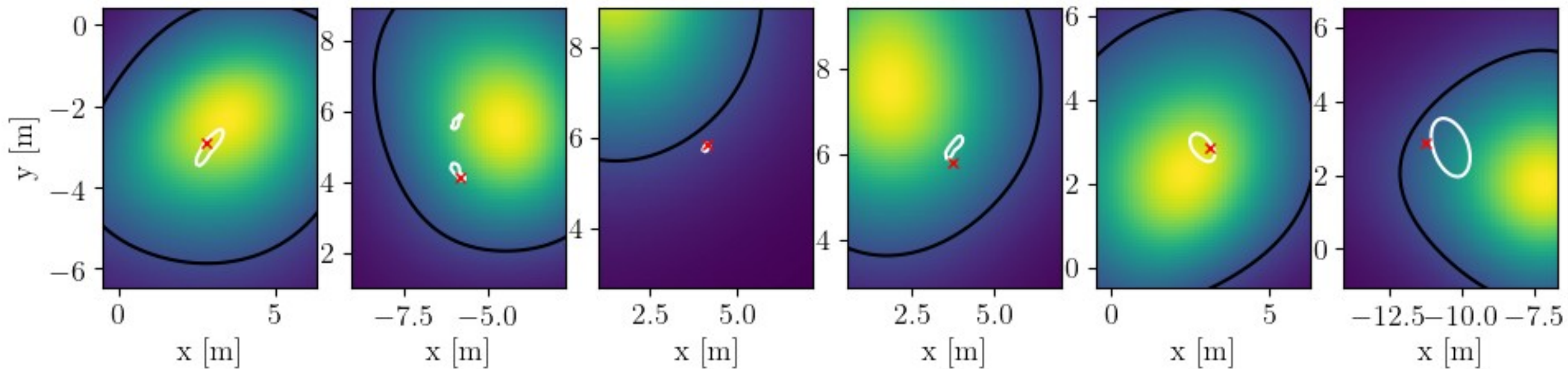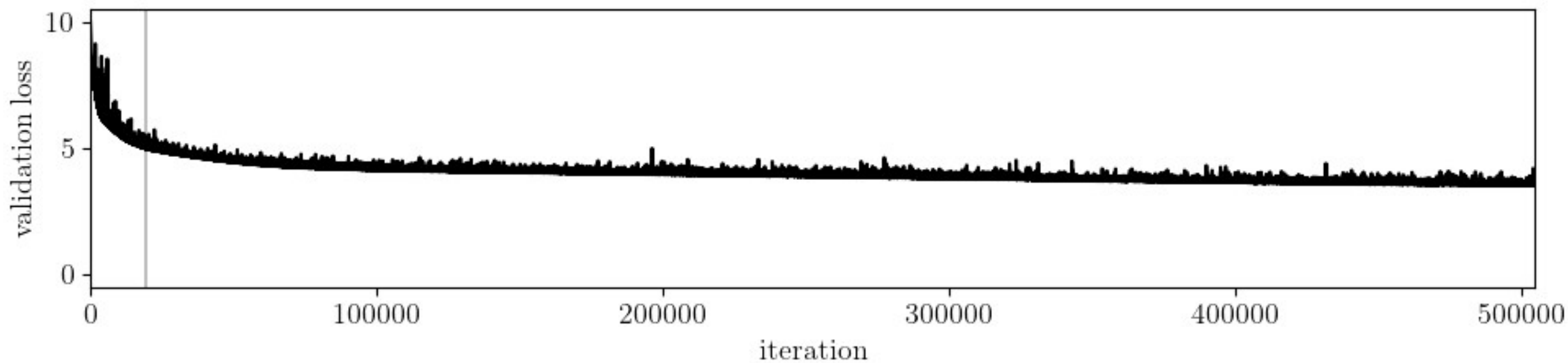
Red: True Label          White: True Posterior 68% region          Black: Predicted Posterior 68% region

**Red: True Label**  **White: True Posterior 68% region**  **Black: Predicted Posterior 68% region**
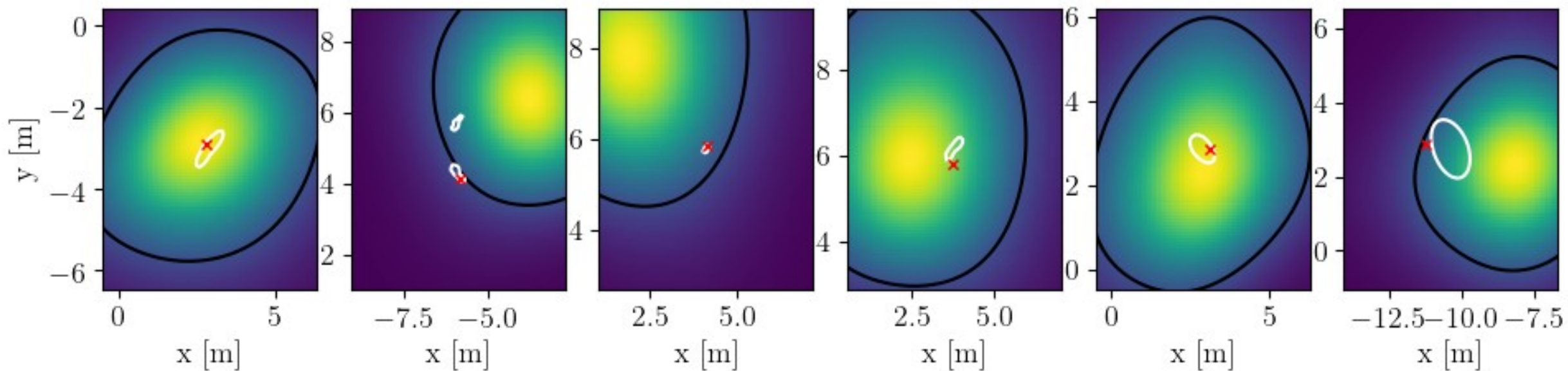
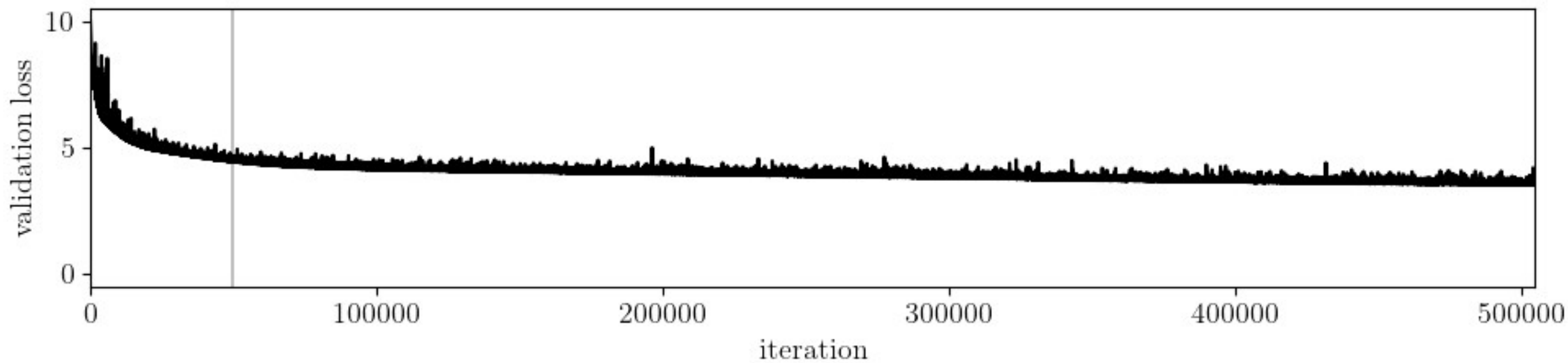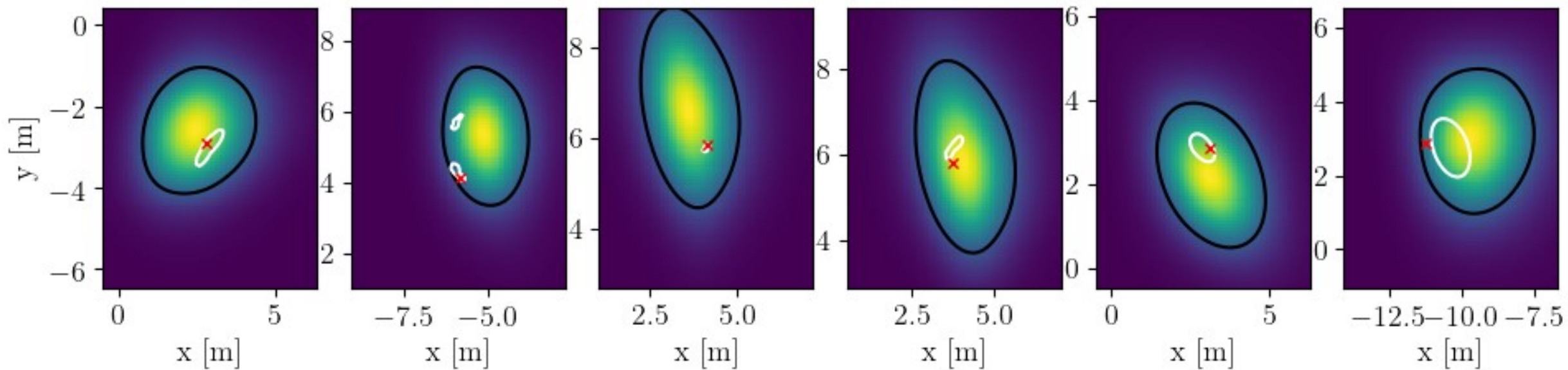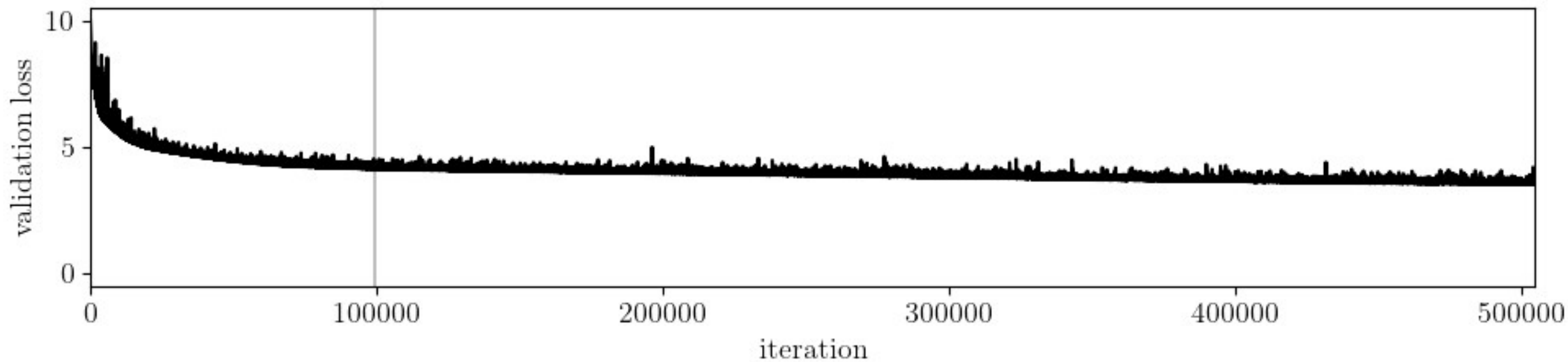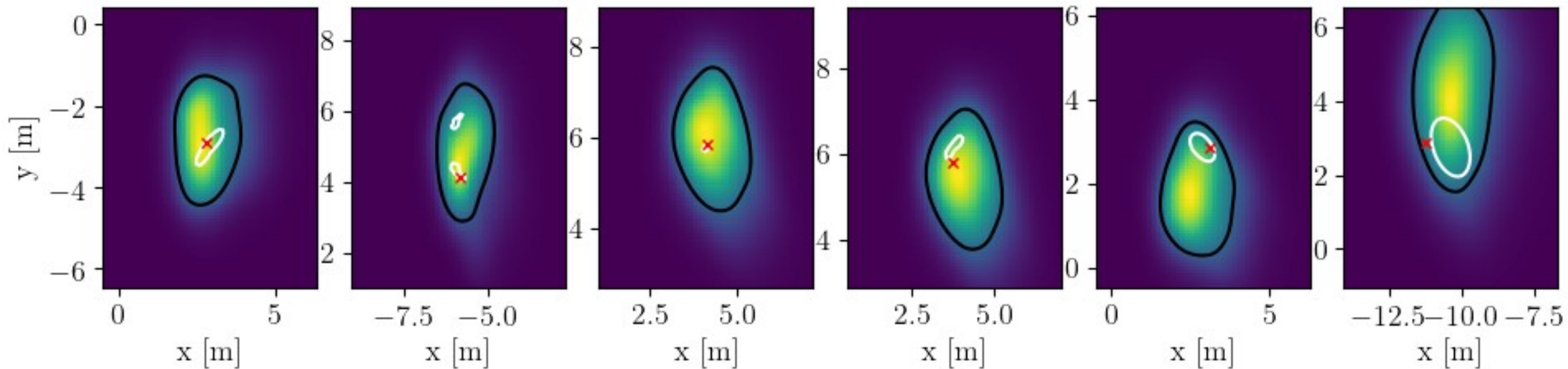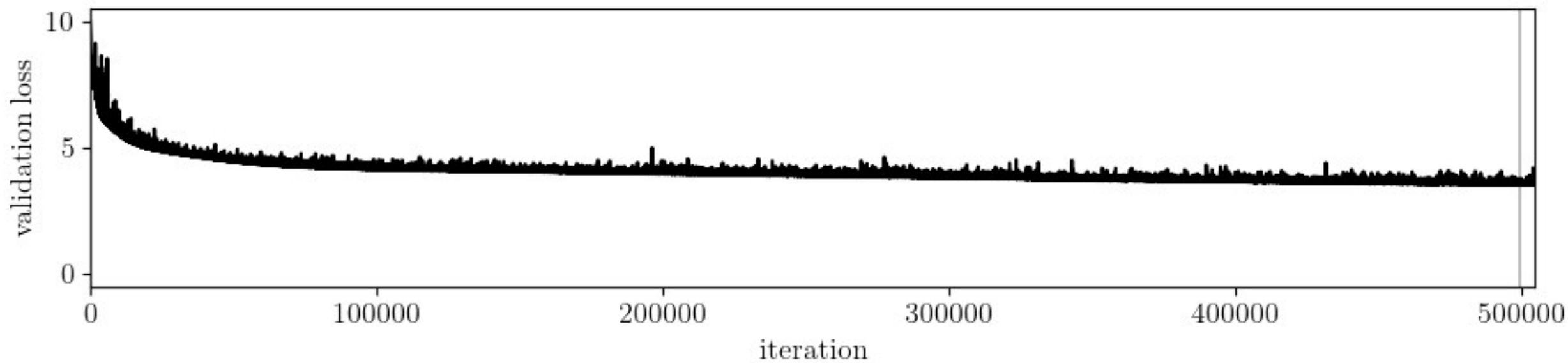Red: True Label          White: True Posterior 68% region          Black: Predicted Posterior 68% region
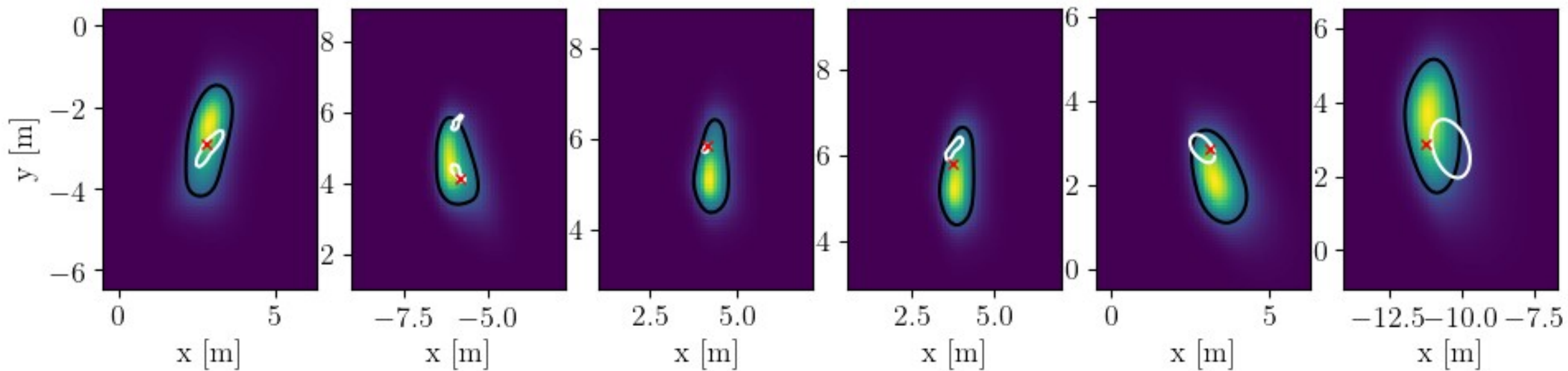
Red: True Label    White: True Posterior 68% region    Black: Predicted Posterior 68% region

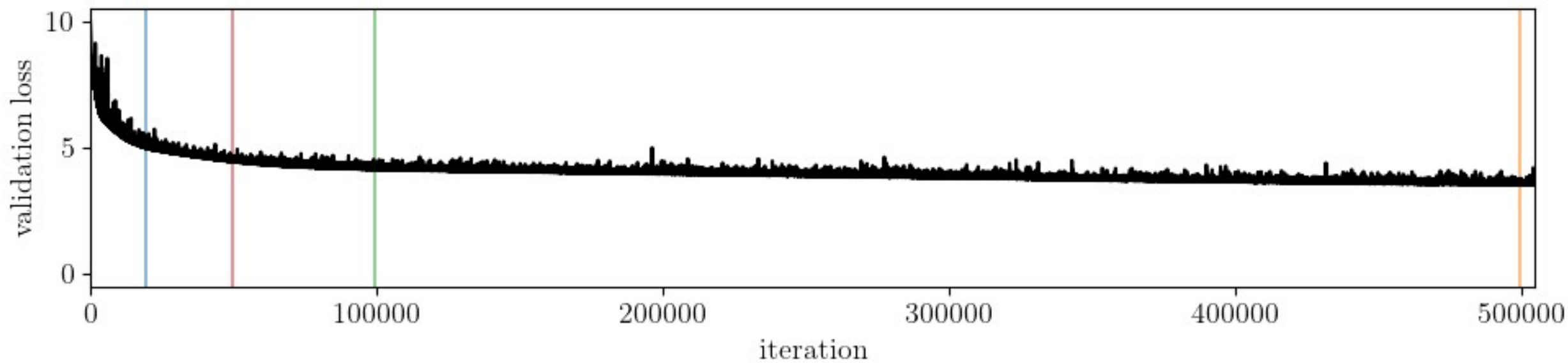**Red: True Label**  **White: True Posterior 68% region**  **Black: Predicted Posterior 68% region**

Red: True Label          True Posterior region: Black          Colors: Approximated Posterior at various stages of learning

# Systematics

Draw events from prior $p(\nu)$ before event generation

This will effectively produce samples from the marginalized True distribution

$$\mathcal{P}_{t,M}(\theta; x) = \int \mathcal{P}_t(\theta; x, \nu) \cdot p(\nu) d\nu$$

, which again is approximated by the neural network

Coverage for a distribution fitted with systematics (red) applied to a standard dataset with a fixed systematic value (green)

Overcoverage
is desired when including
systematics

# Goodness-of-fit

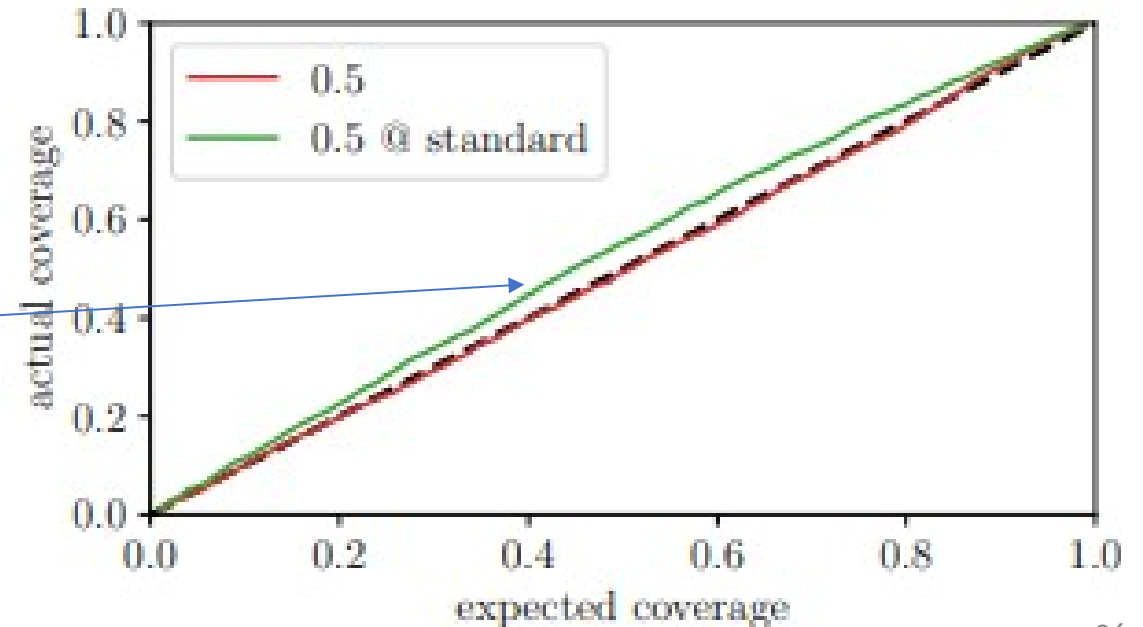- There is a well-motivated construction that upgrades the supervised loss to an extended loss , also fitting a likelihood $p_\theta$ (generative model) and some other terms (see 2008.05825 for details)

- Having both a generative model and the posterior allows to calculate a Bayesian p-value

$$p = \int_{x,z} \boldsymbol{I}_{T(x,z)>T(x_{\mathrm{obs}},z)} p_\theta(x;z) q_\phi(z;x_{\mathrm{obs}}) dx dz$$

$$T(x,z) = \ln p_\theta(x;z)/N_d$$

**Any event sufficiently dissimilar to training data has a low p-value …**

# Summary

- **Supervised Learning learns to approximate the true posterior with a conditional PDF**
  -> **Bayesian systematics**: approximates marginalized Posterior

- Normalizing flows (NFs) can be designed to make this conditional PDF as precise as possible
  -> **supervised learning** can be „upgraded" to behave as usable **likelihood-free inference** (standard MSE loss is not good enough for that)

# Summary

- **Supervised Learning learns to approximate the true posterior with a conditional PDF** -> **Bayesian systematics**: approximates marginalized Posterior

- Normalizing flows (NFs) can be designed to make this conditional PDF as precise as possible -> **supervised learning** can be „upgraded" to behave as usable **likelihood-free inference** (standard MSE loss is not good enough for that)

- Furthermore, normalizing flows allow to
  - **Calculate exact coverage** of the approximate posterior of ANY shape
    - Coverage is obtained very quickly in the training, long before it finishes!
  - **Calculate a goodness-of-fit** that can potentially be used in event selection

- We can upgrade our existing supervised learning models (CNN/LSTM) with NFs, **get improved performance (most for non-gaussian Posteriors**), and get coverage/g-o-f (incl. systematics)

**More info: arXiv:2008.05825**