

Flow Annealed Importance Sampling Bootstrap meets Differentiable Particle Physics

NeurIPS Workshop ML4PS



Annalena Kofler^{1,2,3}



Vincent Stimper^{1,4,5}



Mikhail Mikhasenko^{6,7}



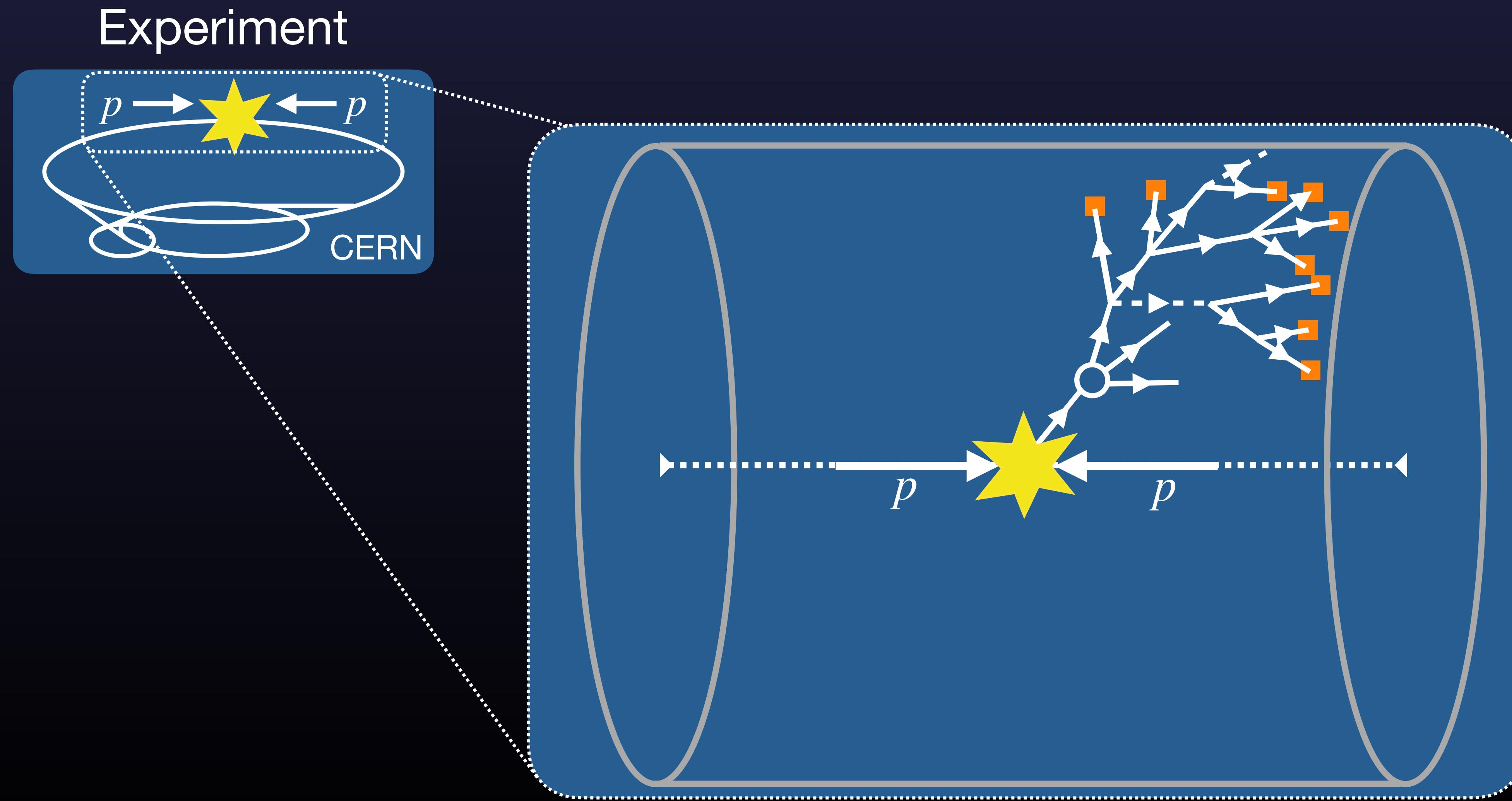
Michael Kagan⁸



Lukas Heinrich³

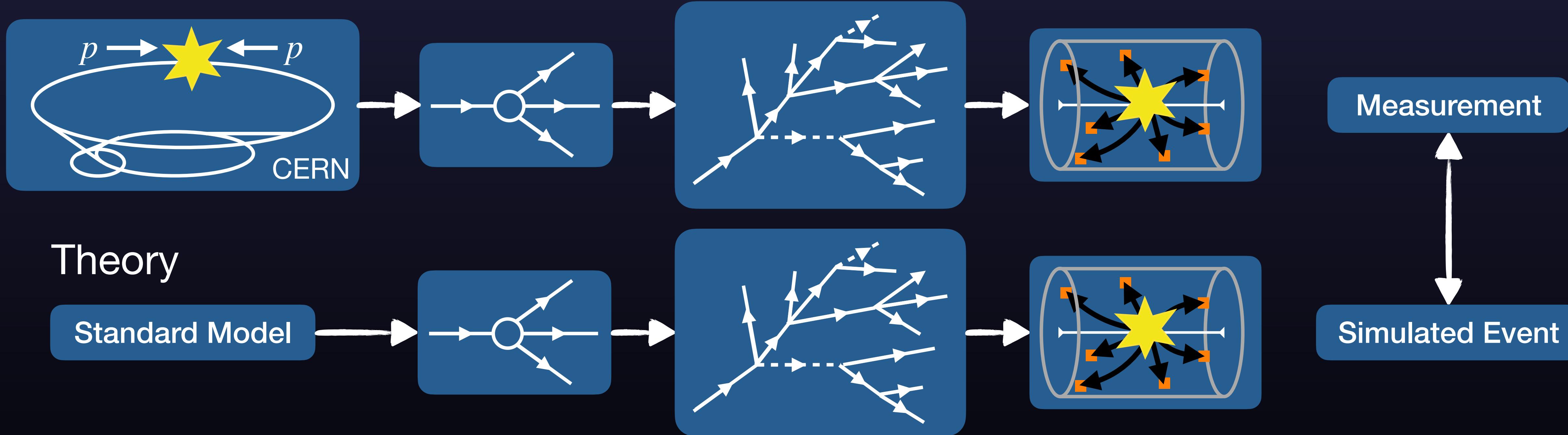
¹Max Planck Institute for Intelligent Systems Tübingen, ²Max Planck Institute for Gravitational Physics Potsdam, ³Technical University of Munich,
⁴Isomorphic Labs, ⁵University of Cambridge, ⁶University of Bochum, ⁷ORIGINS Excellence Cluster Munich, ⁸SLAC National Accelerator Laboratory

How can we learn more about particles?



How can we analyze data?

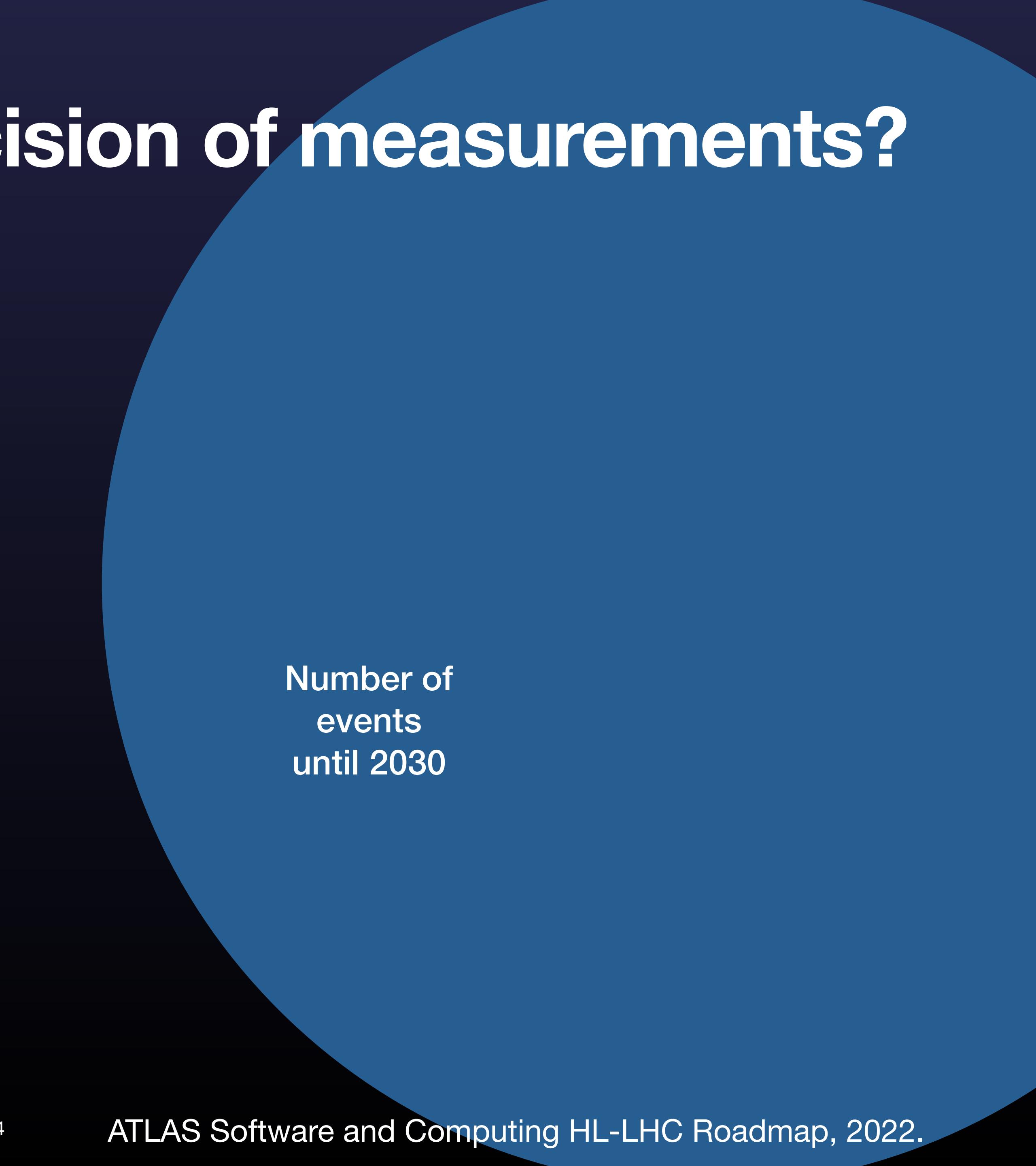
Experiment



→ Find deviations from theory

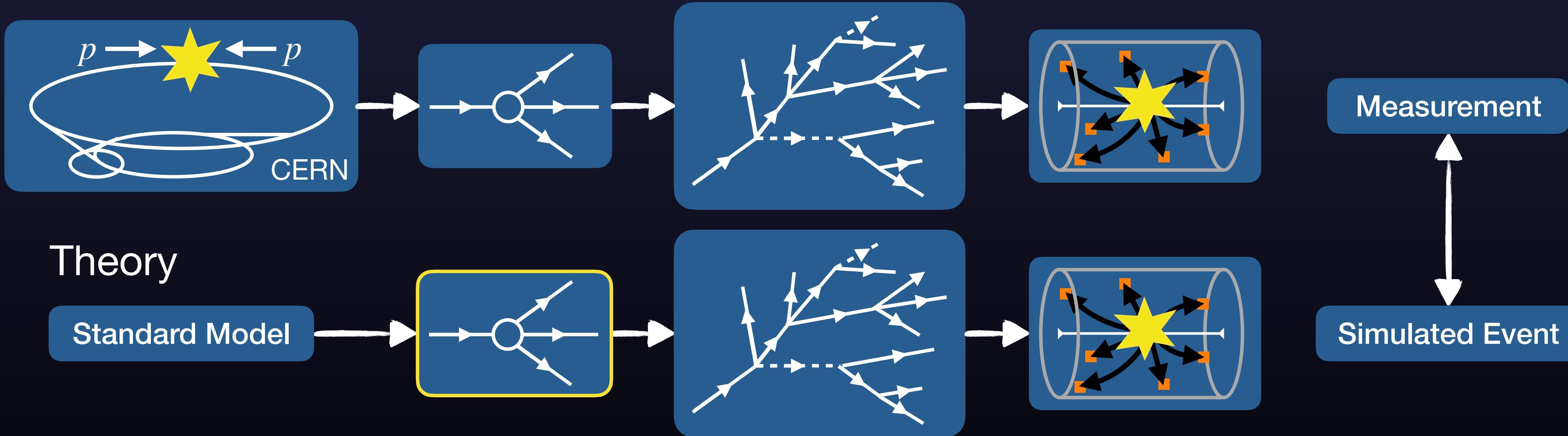
How to improve the precision of measurements?

→ Take more data!



What does this mean for the simulations?

Experiment



→ We need to speed up the simulation!

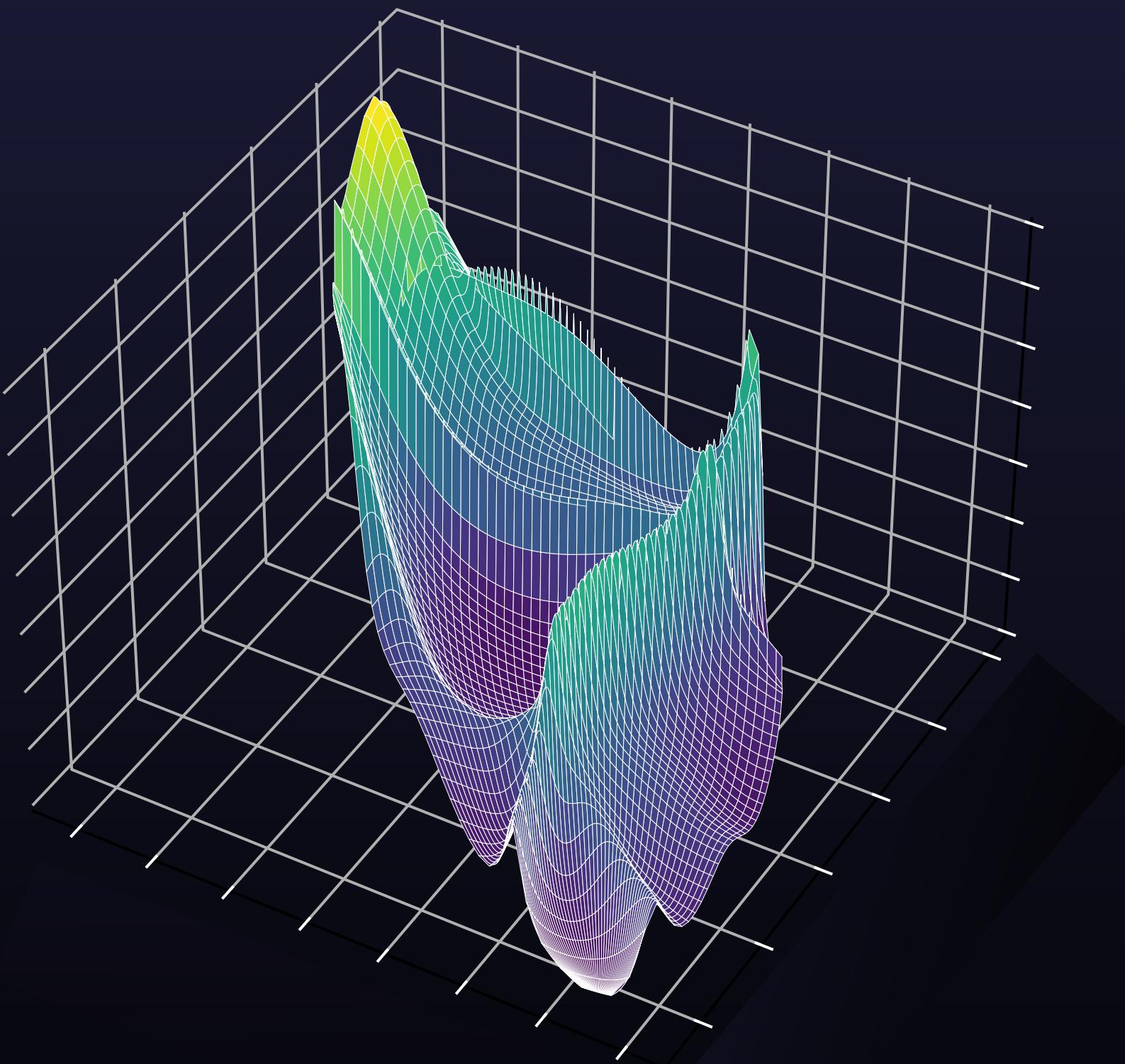
→ In this work: Event generation

What is event generation from a ML perspective?

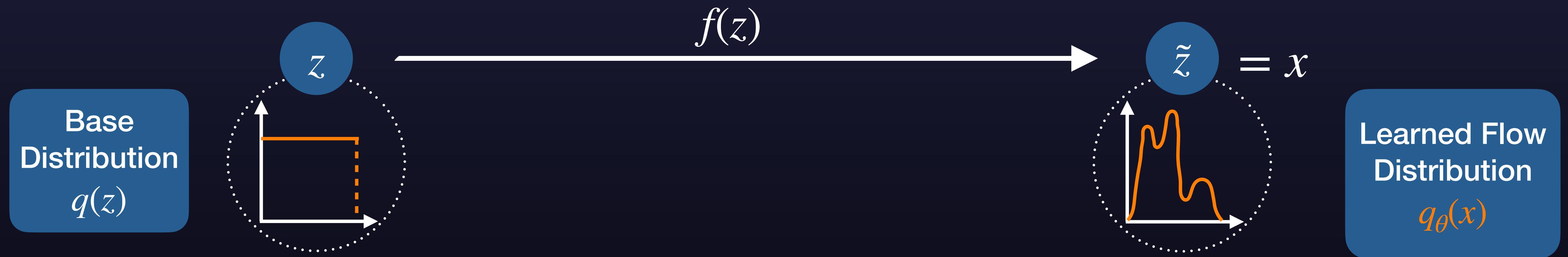
Theory



- Analytic formula for unnormalized distribution $p(x)$ (“matrix element”)
- $p(x)$ describes the outgoing particles
- Sample from this distribution
- **Normalizing flow instead of standard methods**

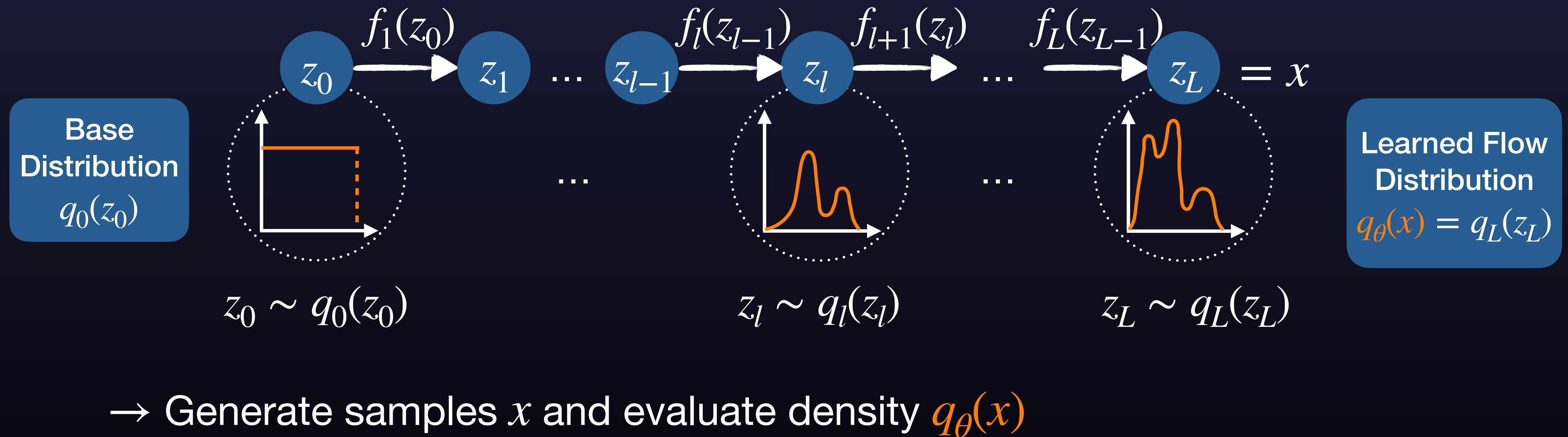


What is a normalizing flow?



Rezende and Mohamed, “Variational Inference with Normalizing Flows.” ICML’15.

What is a normalizing flow?



Rezende and Mohamed, “Variational Inference with Normalizing Flows.” ICML’15.

How to train a normalizing flow?

3 Methods

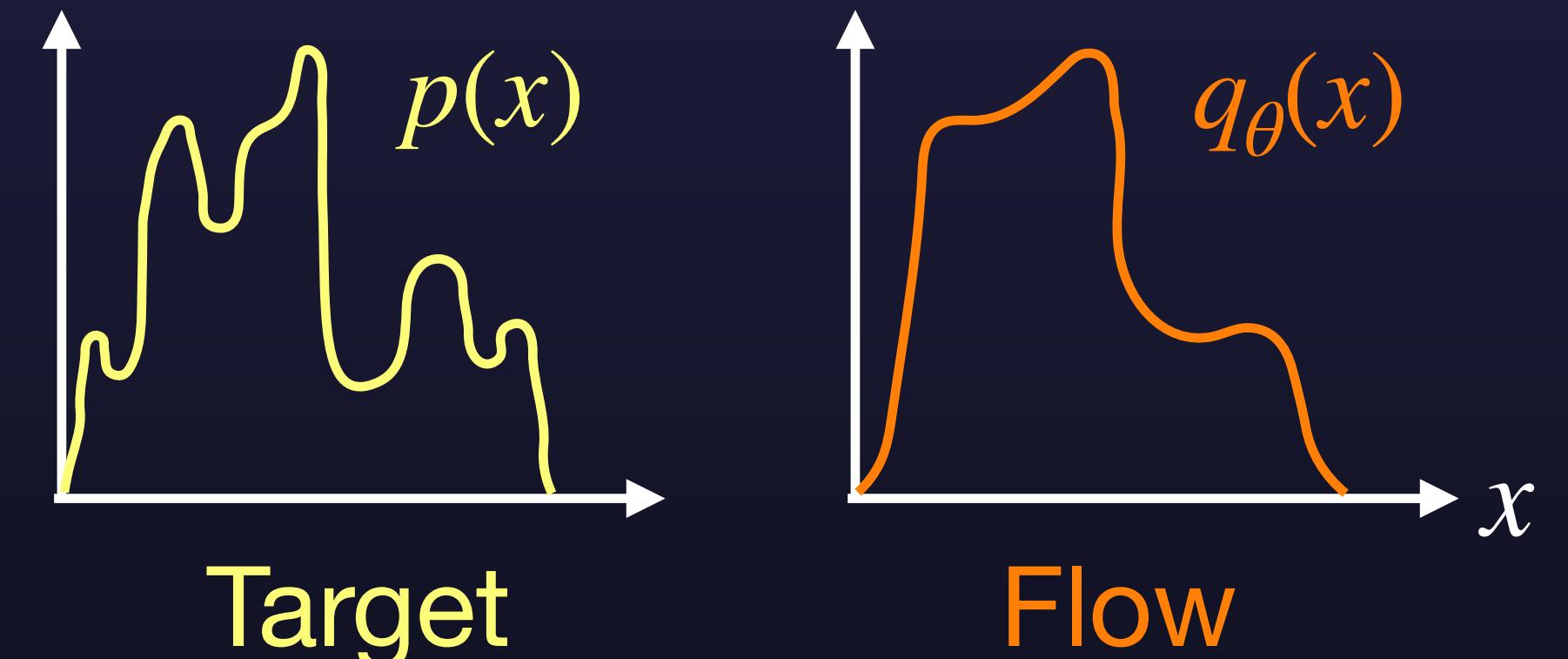
Samples from target $x \sim p(x)$

(1) Forward KL Divergence (fKLD)

Samples from flow $x \sim q_\theta(x)$ and density evaluation of $p(x)$

(2) Reverse KL Divergence (rKLD)

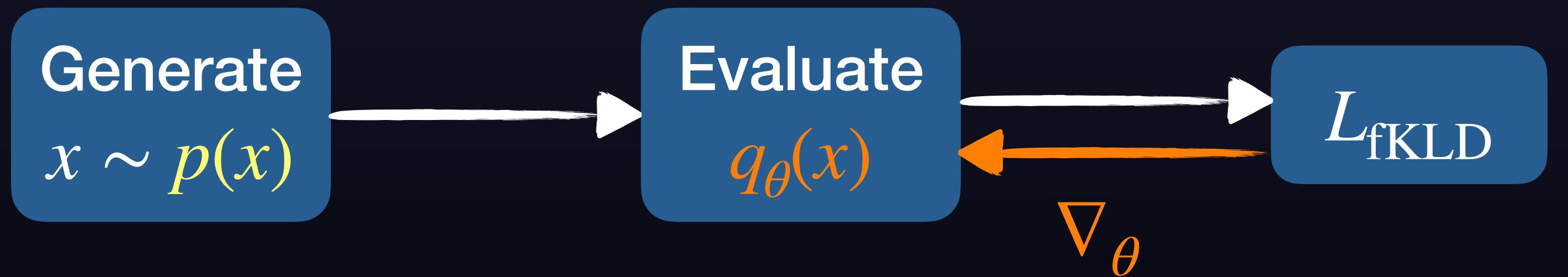
(3) Flow Annealed Importance Sampling Bootstrap (FAB)



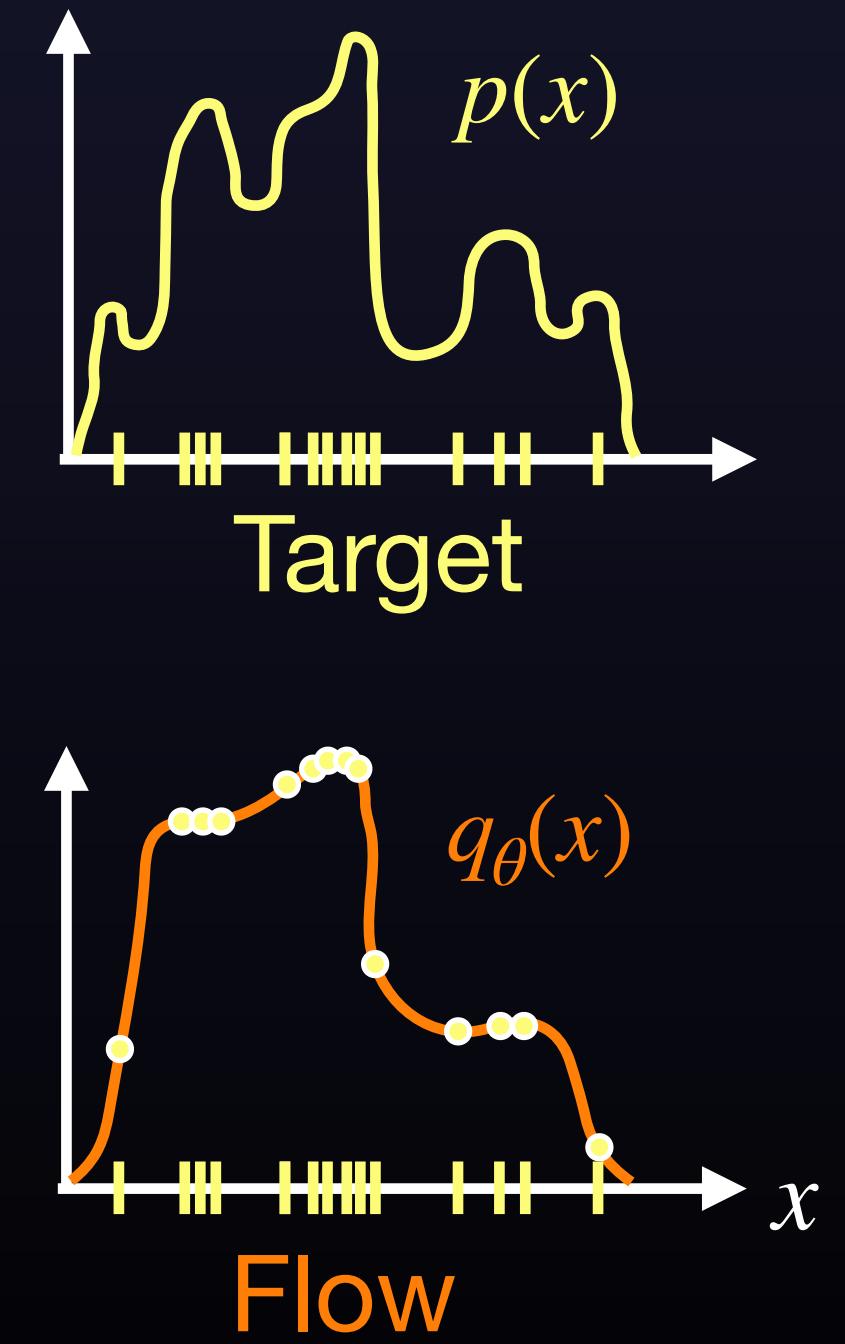
Training with Samples

Forward KL Divergence (fKLD)

$$L_{\text{fKLD}} = D_{\text{KL}}(p \parallel q_{\theta}) = \mathbb{E}_{x \sim p(x)} \left[\log \frac{p(x)}{q_{\theta}(x)} \right] = - \sum_{i=1}^N \log q_{\theta}(x_i) + \text{const.}$$



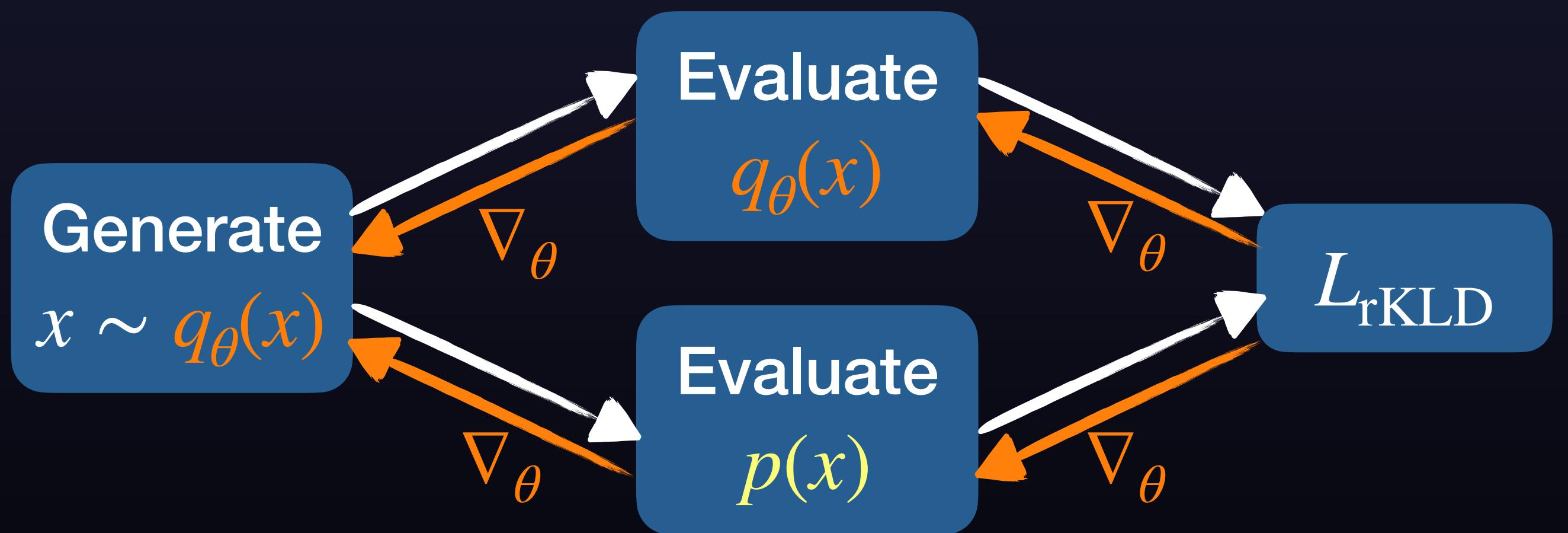
→ Expensive to generate training data



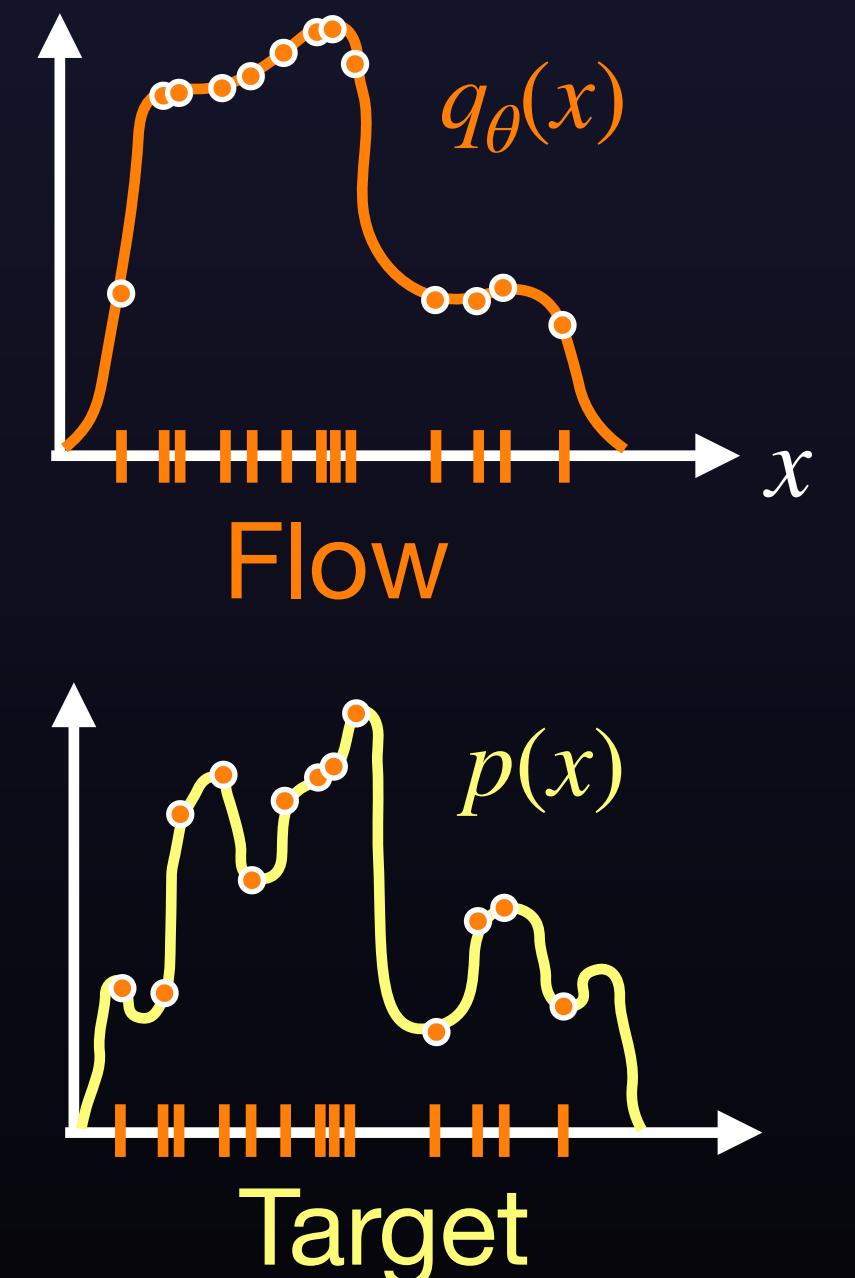
Training with Density Evaluation

Reverse KL Divergence (rKLD)

$$L_{\text{rKLD}} = D_{\text{KL}}(q_{\theta} \parallel p) = \mathbb{E}_{x \sim q_{\theta}(x)} \left[\log \frac{q_{\theta}(x)}{p(x)} \right]$$



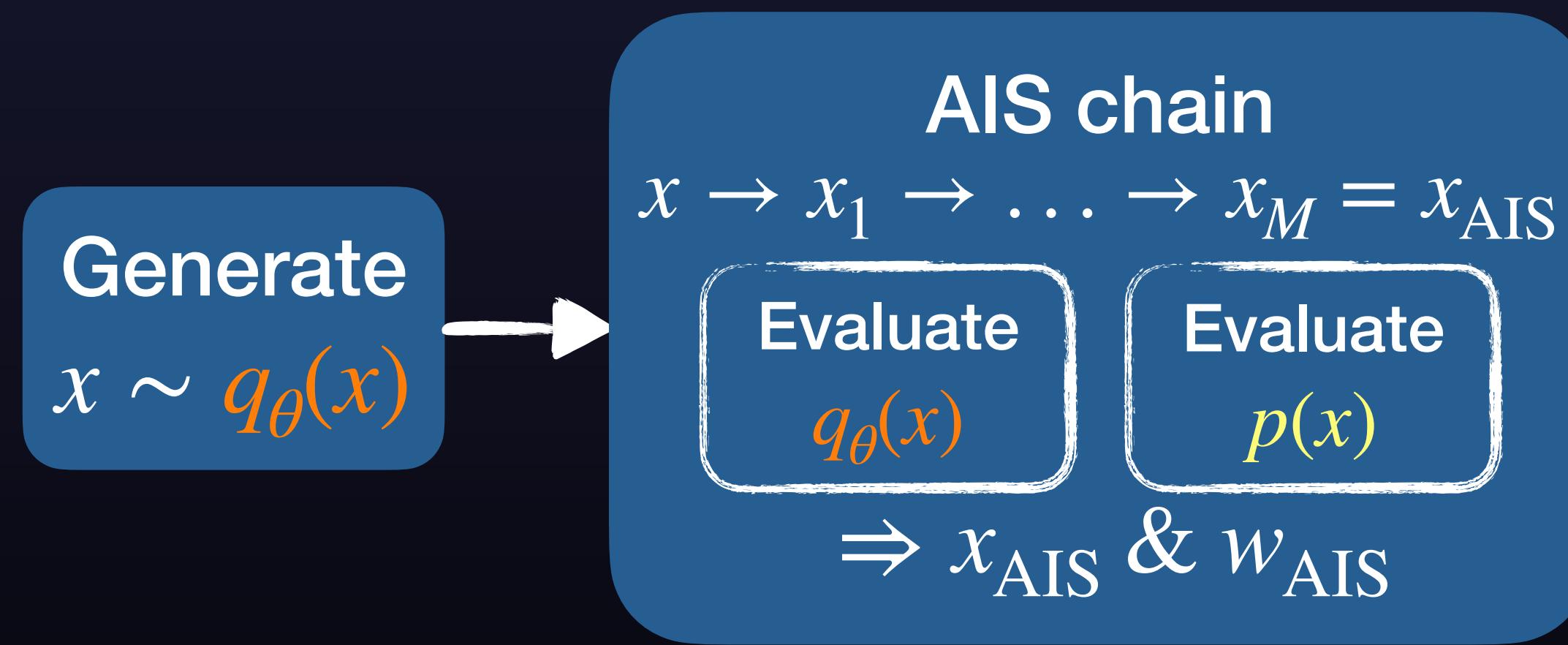
→ Requires differentiable target distribution $p(x)$



Training with Density Evaluation

Flow Annealed Importance Sampling Bootstrap (FAB)

Improve flow samples $x \sim q_\theta(x)$ with Annealed Importance Sampling (AIS)

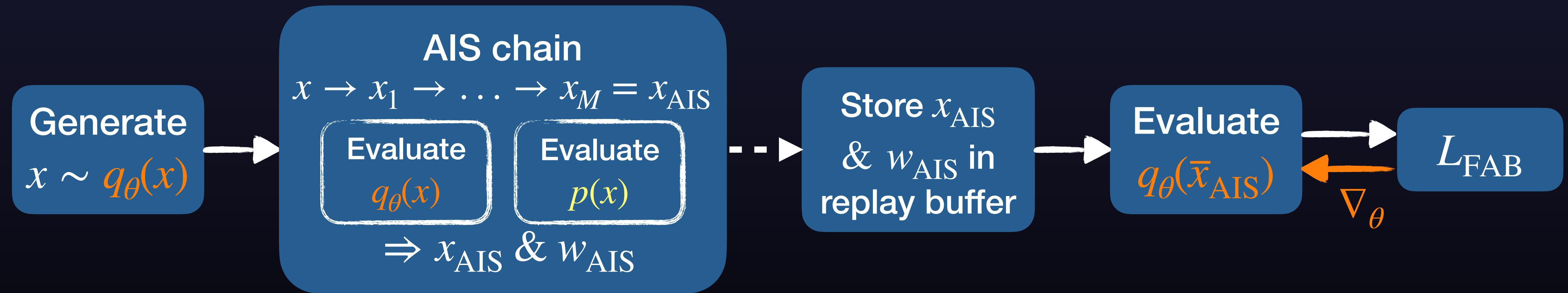


AIS chain implemented via Hamiltonian Monte Carlo (HMC)
→ Requires differentiable target $p(x)$

Training with Density Evaluation

Flow Annealed Importance Sampling Bootstrap (FAB)

Improve flow samples $x \sim q_\theta(x)$ with Annealed Importance Sampling (AIS)



$$L_{\text{FAB}} = - \sum_{i=1}^N \frac{\bar{w}_{\text{AIS}}^{(i)}}{\sum_j \bar{w}_{\text{AIS}}^{(j)}} q_\theta \left(\bar{x}_{\text{AIS}}^{(i)} \right) \rightarrow \text{minimizes variance of importance weights}$$

Differentiable target distributions

Examples

Recent work in particle physics: differentiable implementations of $p(x)$

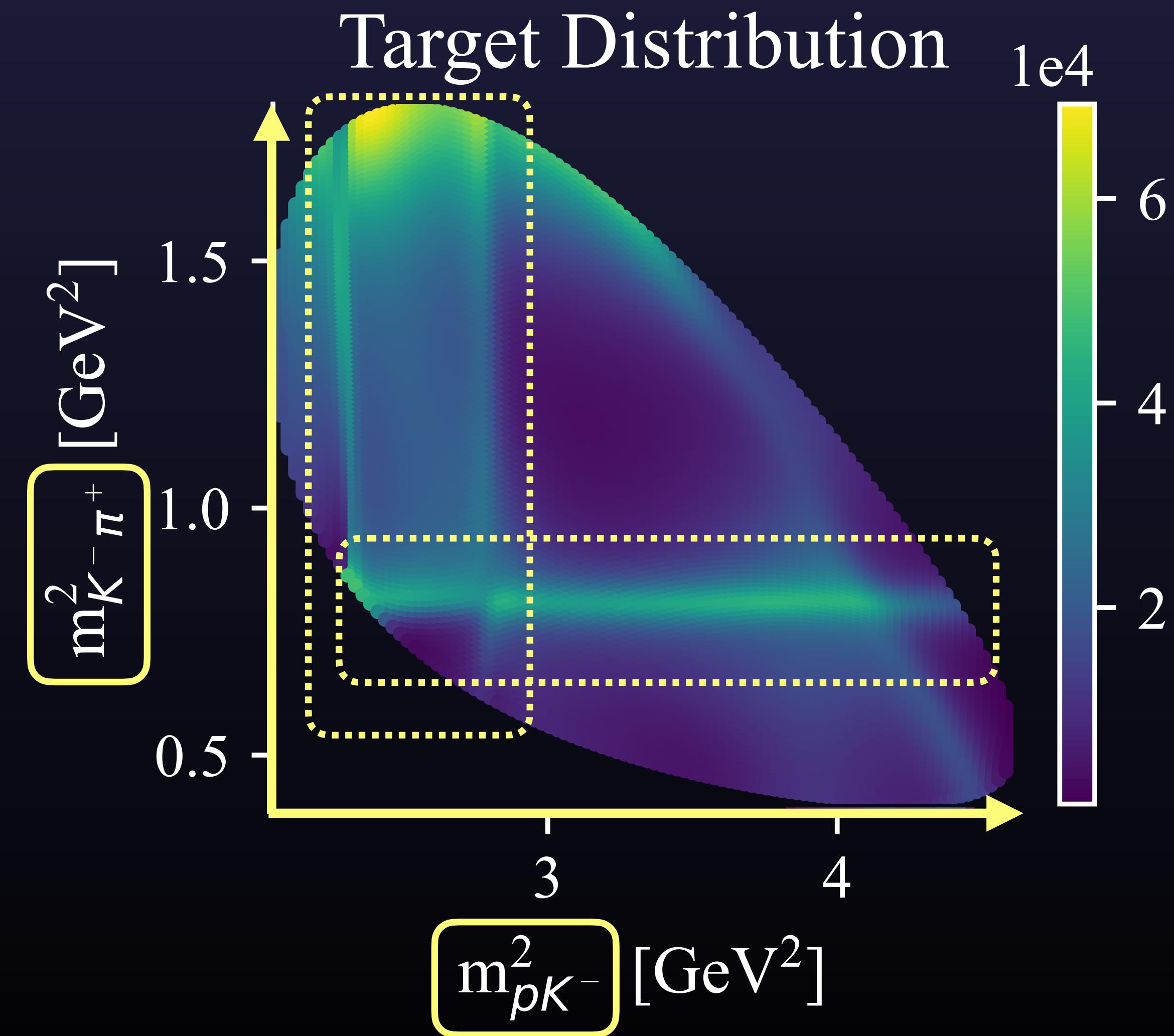
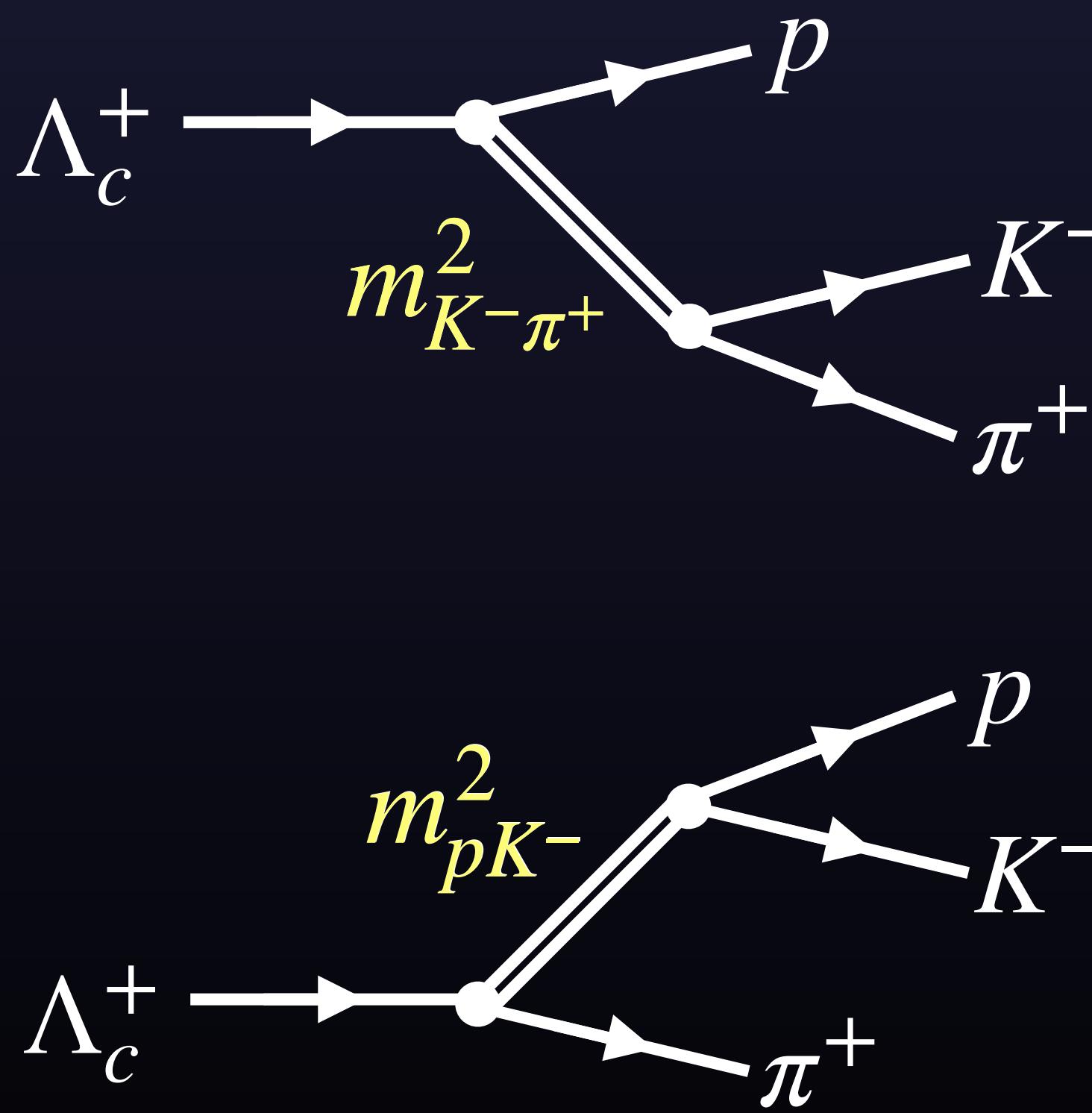
- **ComPWA:** $\Lambda_c^+ \rightarrow p K^- \pi^+$ (2D)
- **MadJAX:** $e^+ e^- \rightarrow t\bar{t}, t\bar{t} \rightarrow W^+ b, \bar{t} \rightarrow W^- \bar{b}$ (8D)

Common Partial Wave Analysis: A collaboration-independent organisation for amplitude analysis software, doi 10.5281/zenodo.6908150

Heinrich and Kagan, “Differentiable Matrix Elements with Madjax.” J. Phys. Conf., 2023.

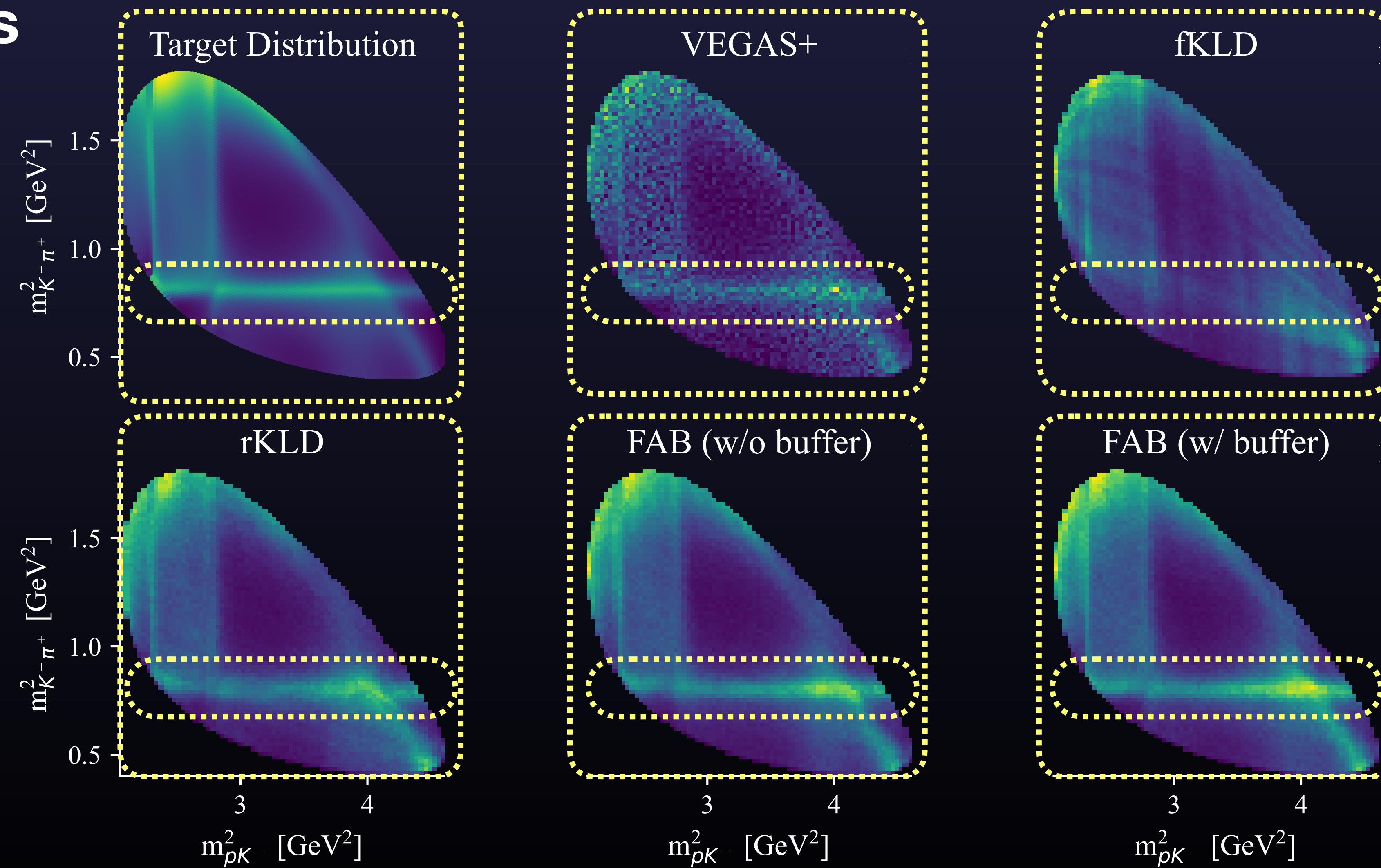
2D: $\Lambda_c^+ \rightarrow p K^- \pi^+$

What do we see?



2D: $\Lambda_c^+ \rightarrow p K^- \pi^+$

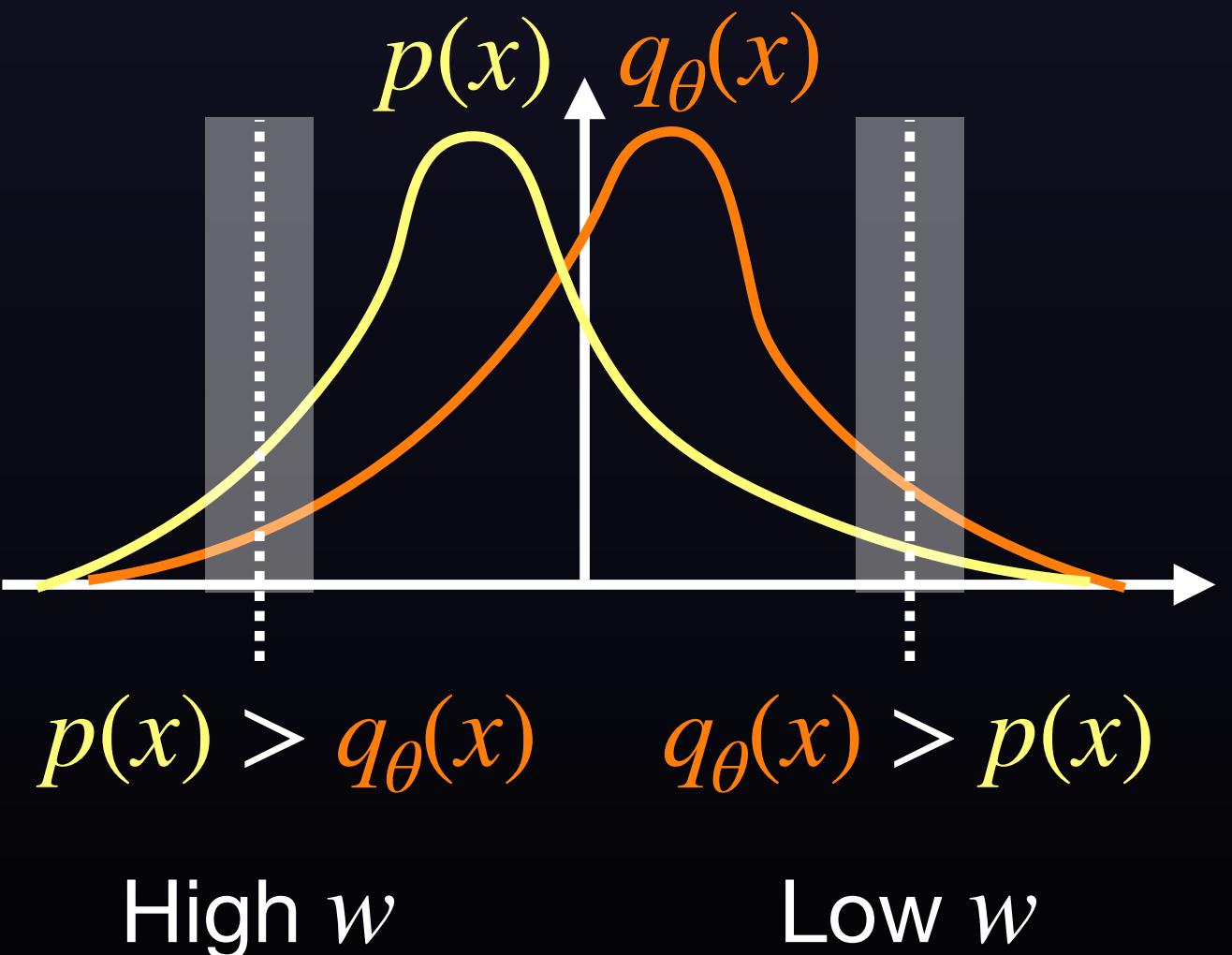
Histograms

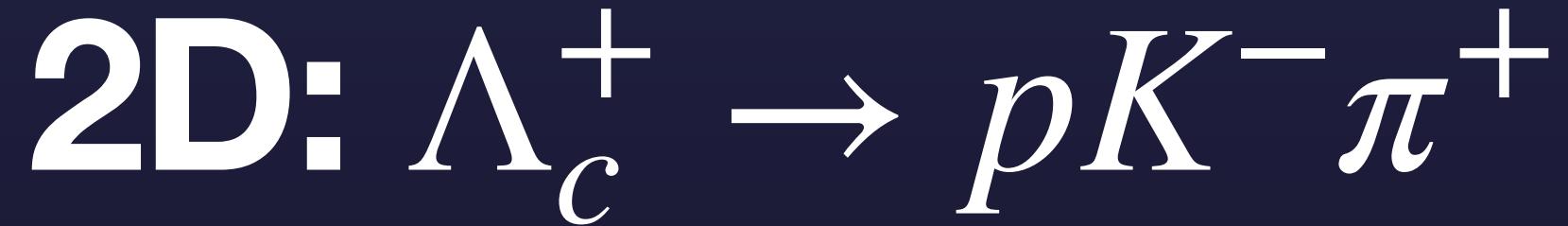


How to compare?

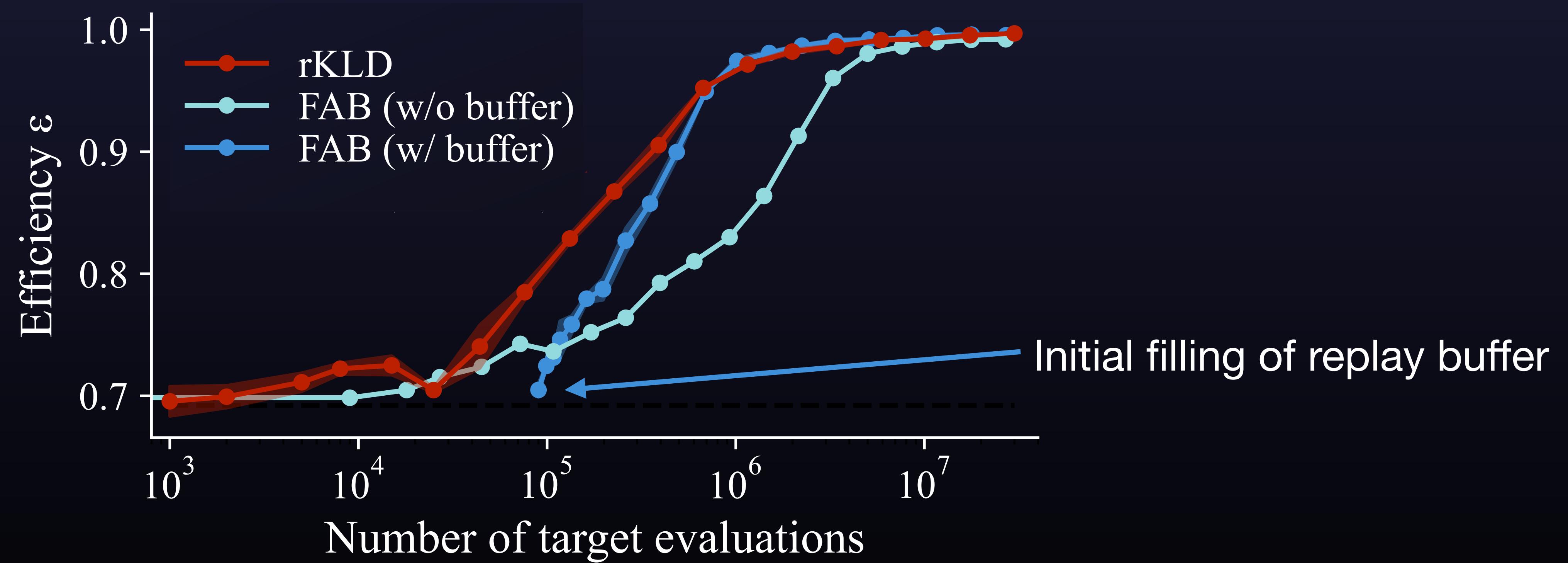
- Efficiency ϵ

$$\epsilon = \frac{1}{N} \frac{\left(\sum_i w_i \right)^2}{\sum_i w_i^2} \in [0,1] \quad \text{with importance weight} \quad w_i = \frac{p(x_i)}{q_\theta(x_i)}$$



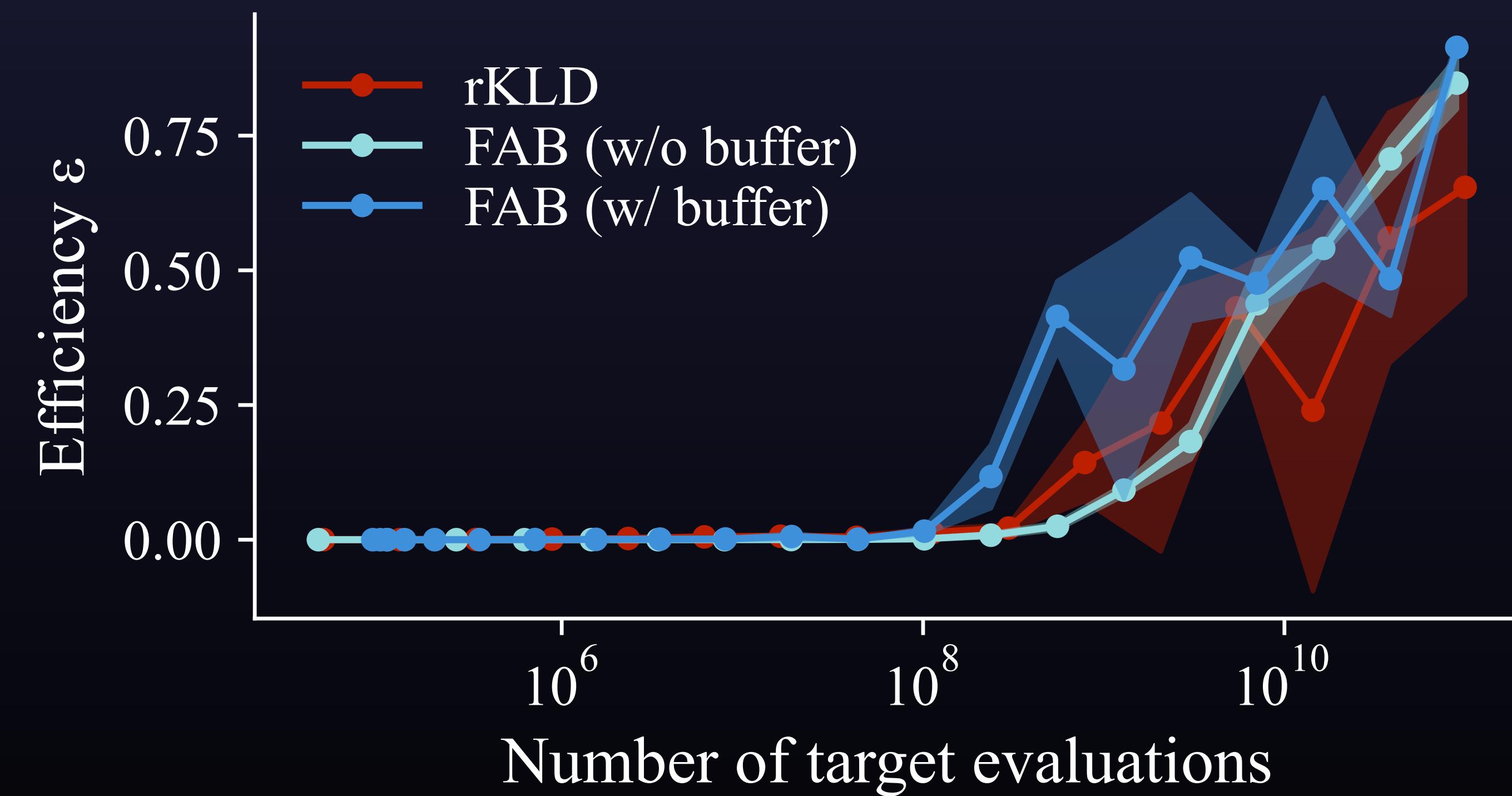


Efficiency vs. Number of target evaluations



8D: $e^+e^- \rightarrow t\bar{t}, t\bar{t} \rightarrow W^+b, \bar{t} \rightarrow W^-\bar{b}$

Efficiency vs. Number of target evaluations



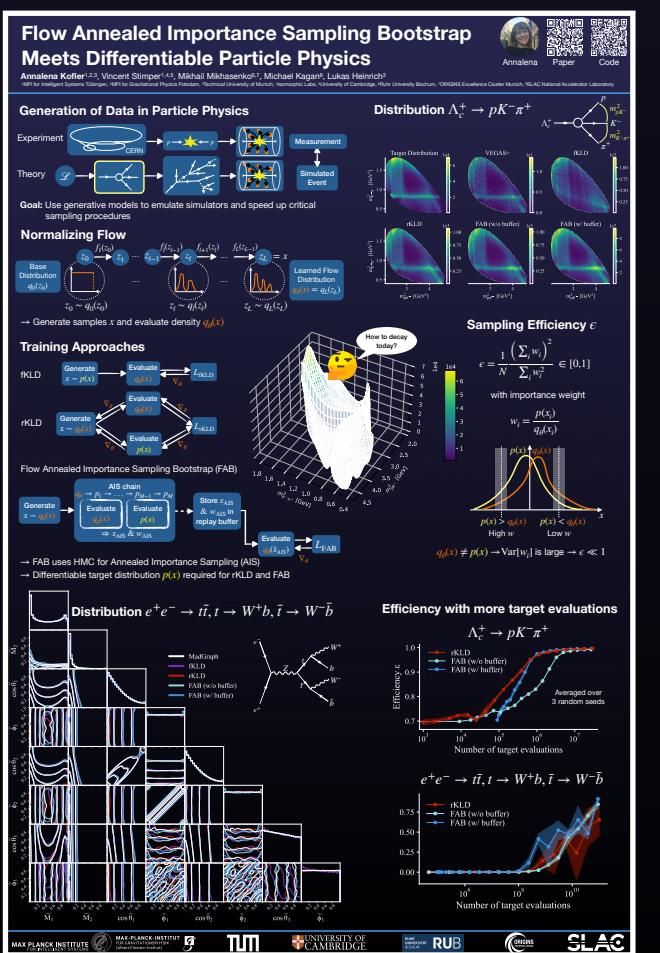
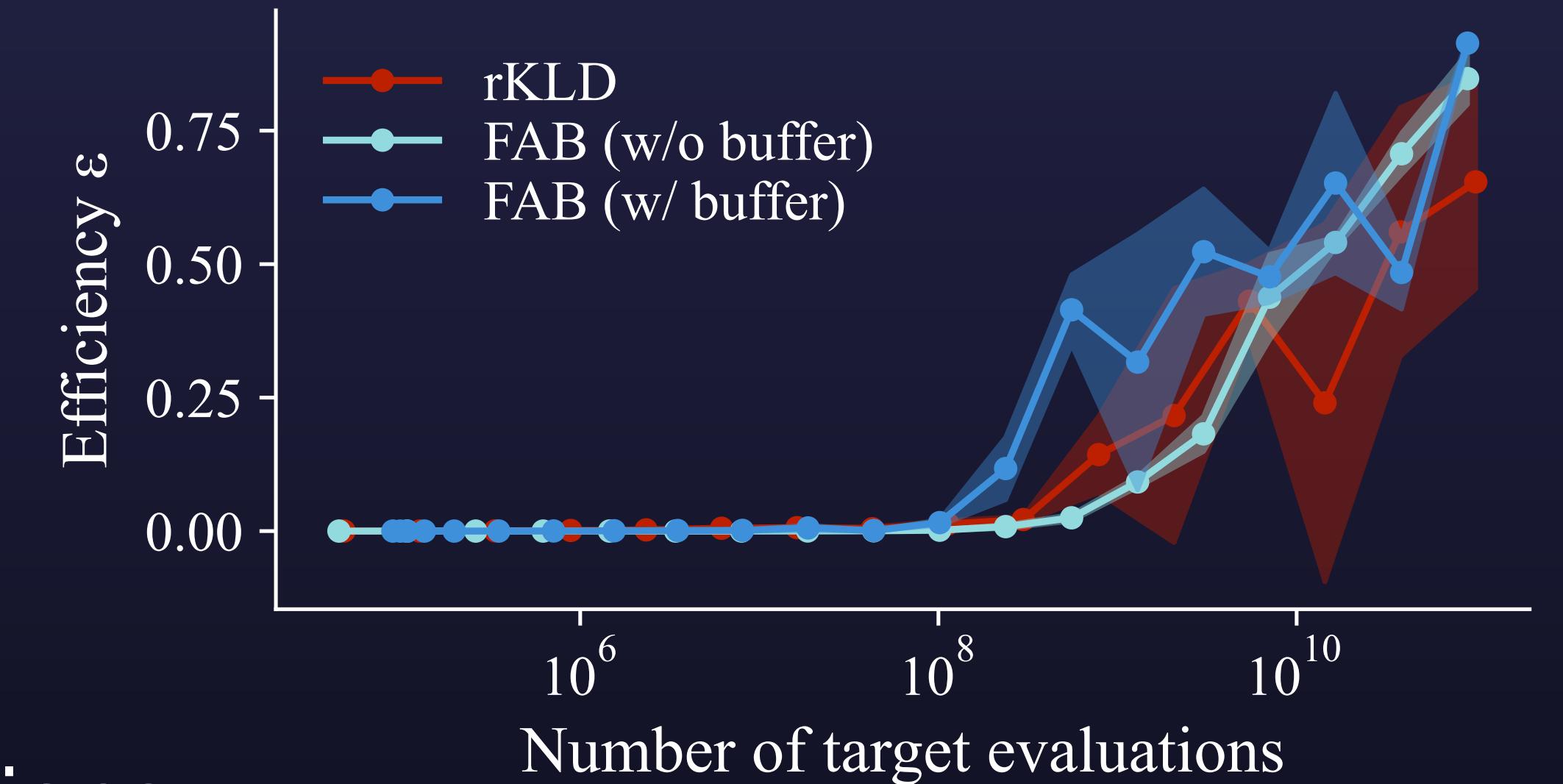
Take-Aways

FAB (w/ buffer) ...

... successfully adapted to particle physics

... outperforms other methods in high dimensions

... achieves higher sampling efficiency with fewer target evaluations



Visit my poster!
(16:15 - 17:25)



Paper



Code

Do you have any questions?

You want to have an in-depth discussion?
 → Visit my poster (16:15 - 17:25)



Paper



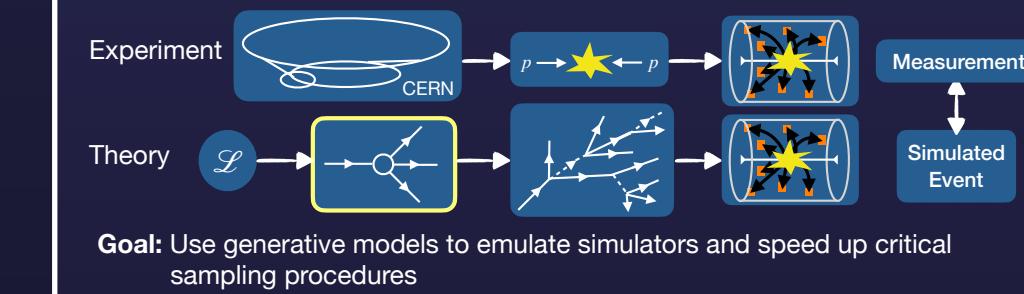
Code

Flow Annealed Importance Sampling Bootstrap Meets Differentiable Particle Physics

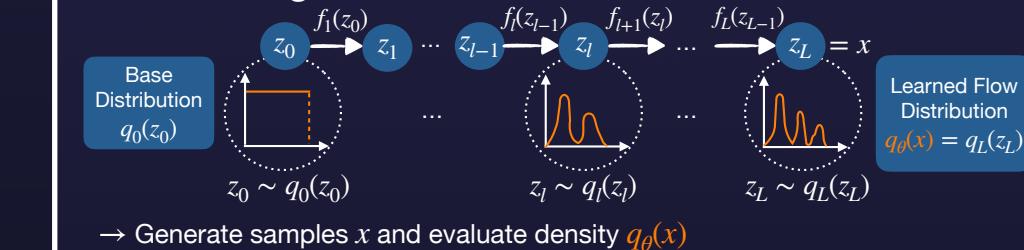
Annalena Kofler^{1,2,3}, Vincent Stumper^{1,4,5}, Mikhail Mikhaseko^{6,7}, Michael Kagan⁸, Lukas Heinrich³
¹MPI for Intelligent Systems Tübingen, ²MPG for Gravitational Physics Potsdam, ³Technical University of Munich, ⁴Isomorphic Labs, ⁵University of Cambridge, ⁶Fuhr University Bochum, ⁷ORIGINS Excellence Cluster Munich, ⁸SLAC National Accelerator Laboratory

Annalena
 Paper
 Code

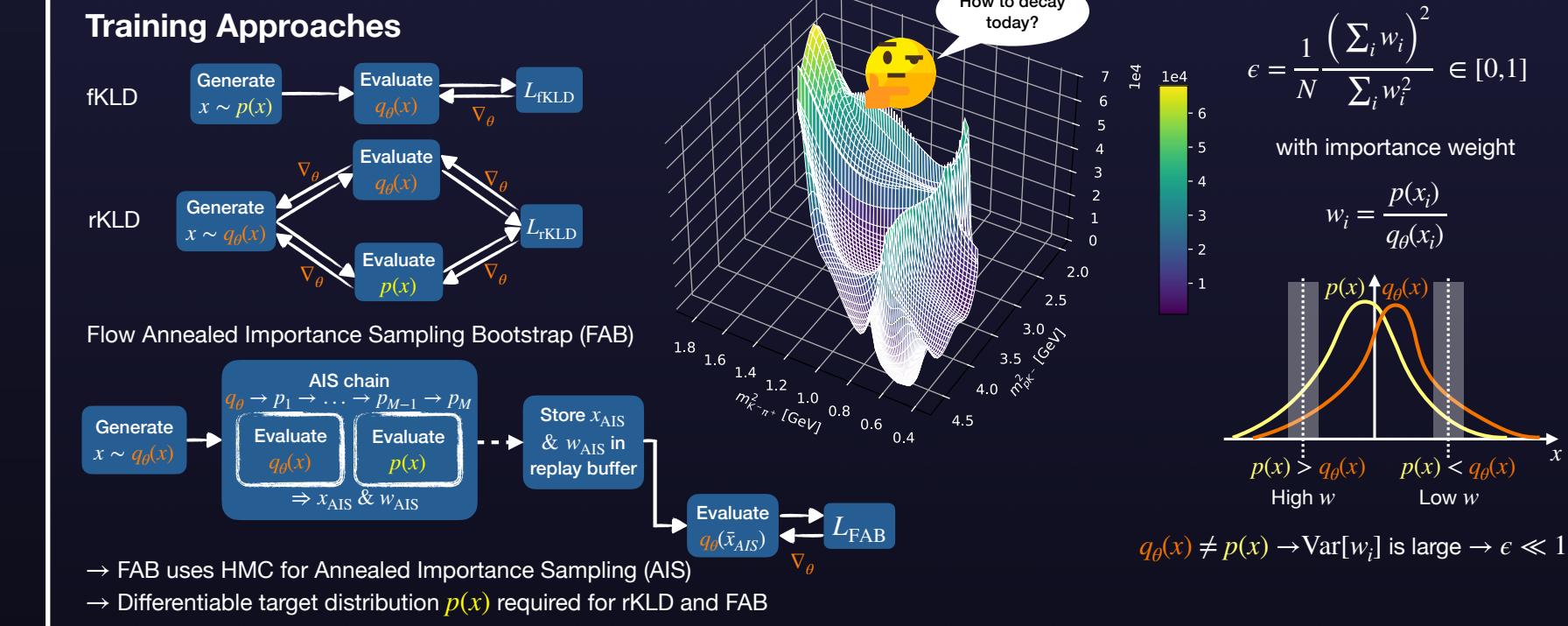
Generation of Data in Particle Physics



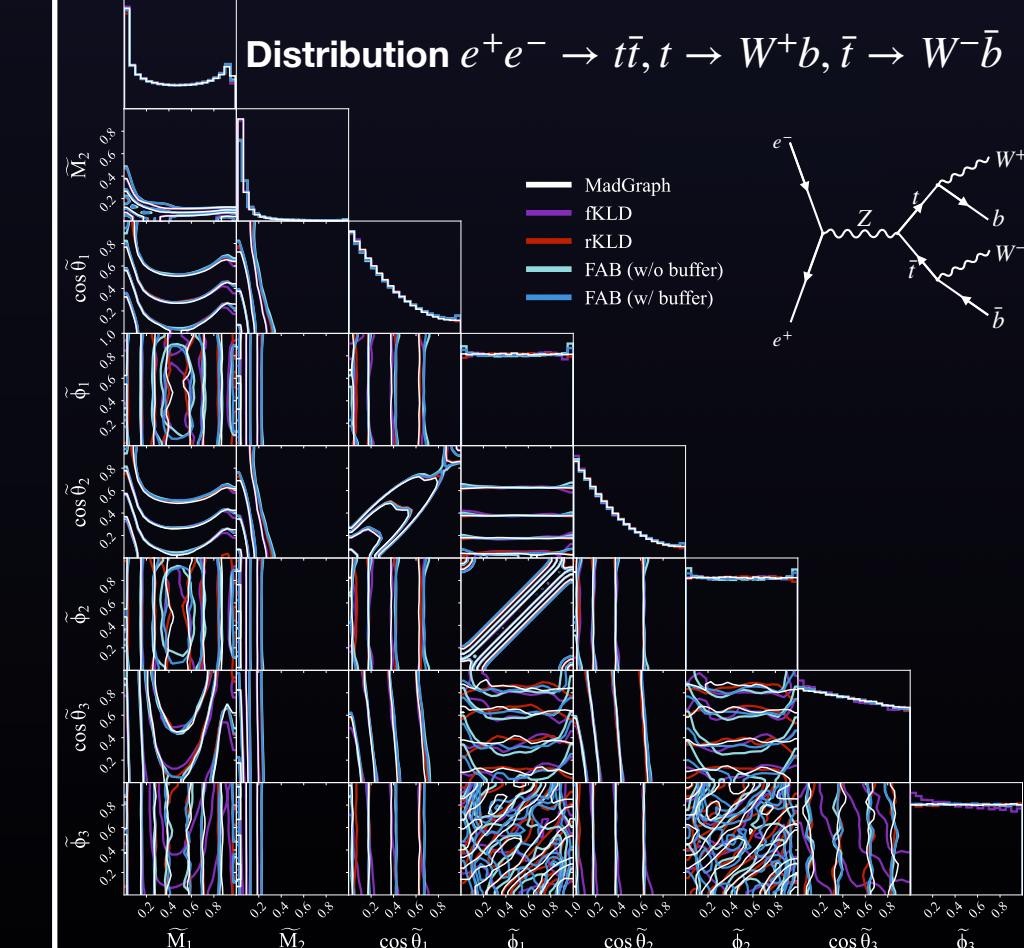
Normalizing Flow



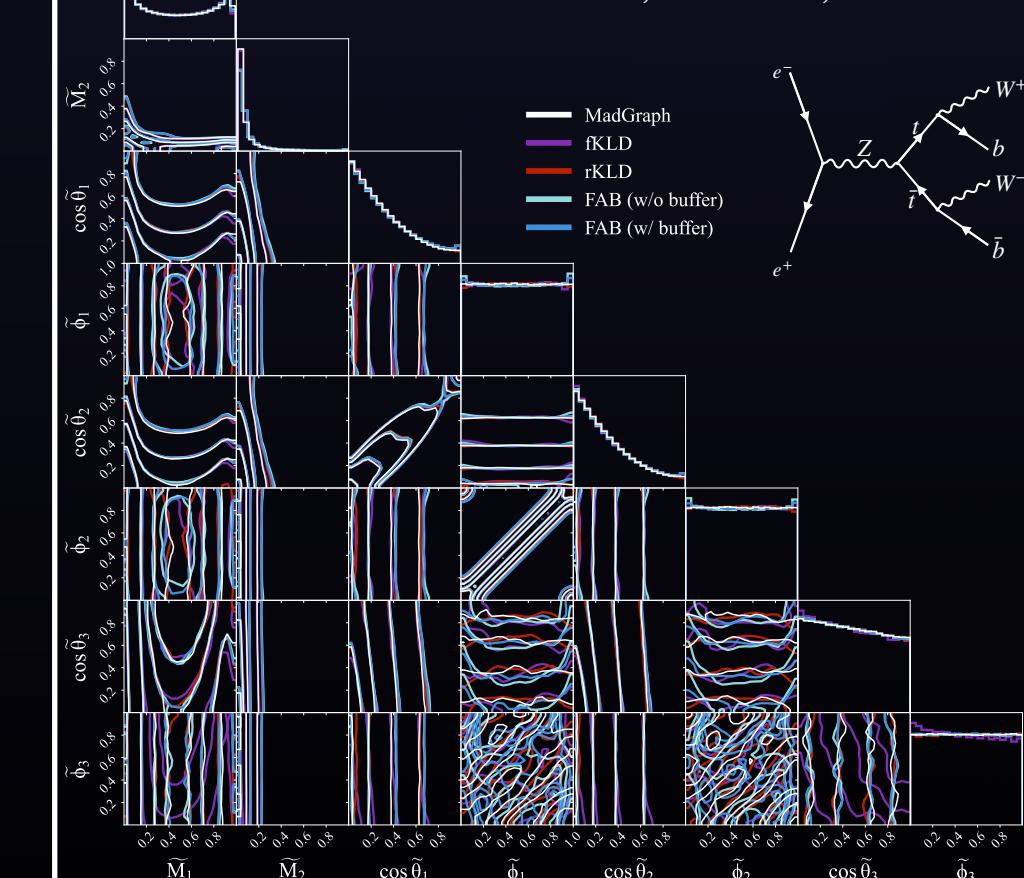
→ Generate samples x and evaluate density $q_\theta(x)$



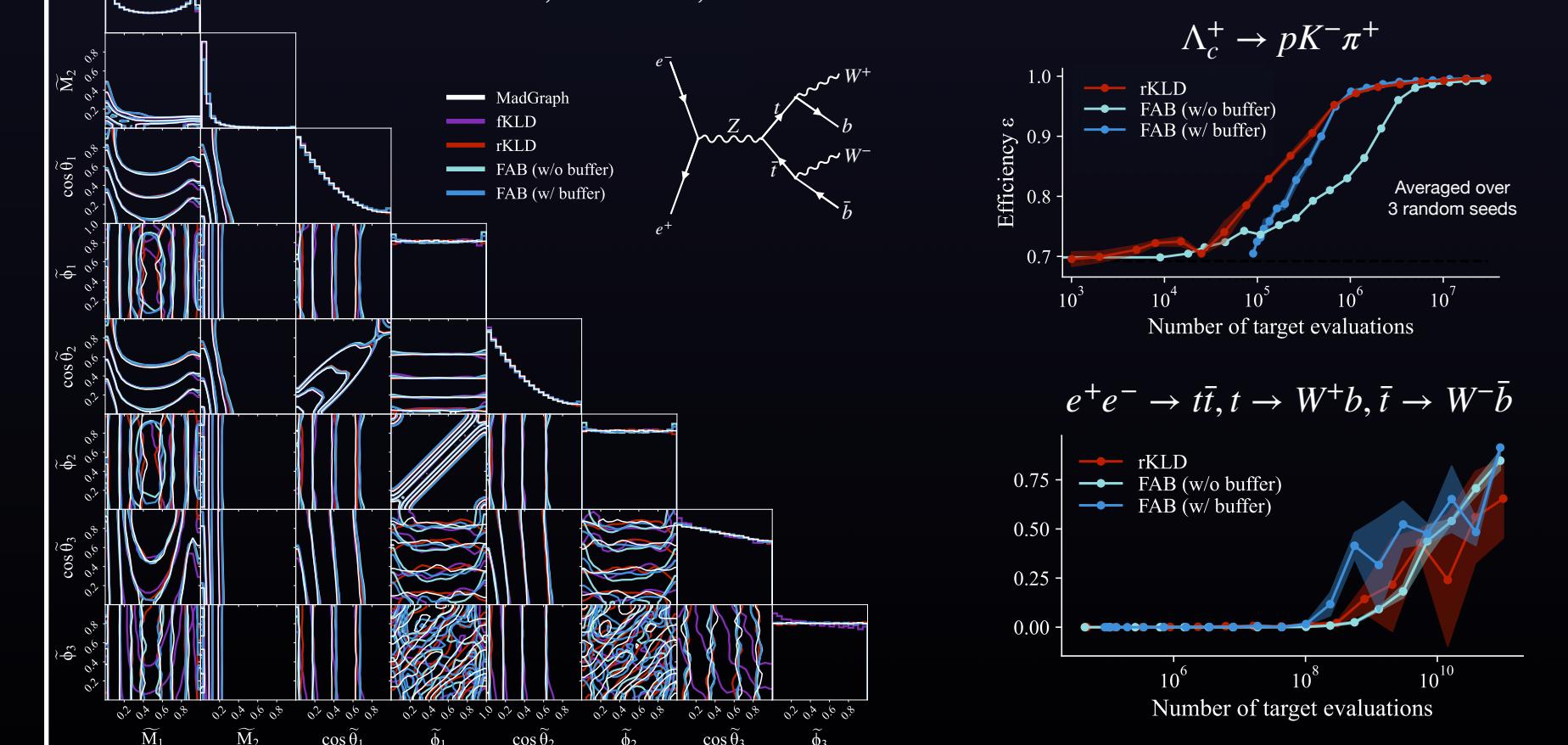
\rightarrow FAB uses HMC for Annealed Importance Sampling (AIS)
 \rightarrow Differentiable target distribution $p(x)$ required for rKLD and FAB



Distribution $e^+e^- \rightarrow t\bar{t}, t \rightarrow W^+b, \bar{t} \rightarrow W^-\bar{b}$



Efficiency with more target evaluations



References

- ATLAS Software and Computing HL-LHC Roadmap, 2022.
- Rezende and Mohamed, “Variational Inference with Normalizing Flows.” ICML’15.
- Durkan et al., “Neural Spline Flows.”, NeurIPS’19.
- LHCb, “Amplitude analysis of the $\Lambda_c^+ \rightarrow p K^- \pi^+$ decay and Λ_c^+ baryon polarization measurement in semileptonic beauty hadron decays.”, Phys. Rev. D 108, 2023.
- Heinrich and Kagan, “Differentiable Matrix Elements with MadJax.” J. Phys. Conf., 2023.
- Midgley, Stimper, et al., “Flow Annealed Importance Sampling Bootstrap”, ICLR, 2023.

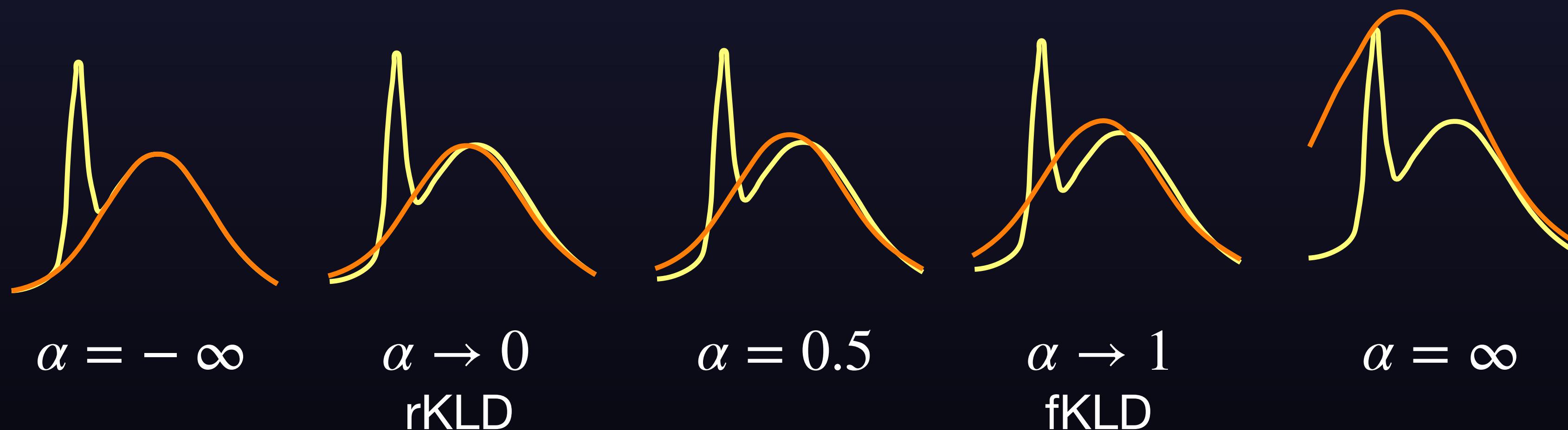
Appendix

Flow Annealed Importance Sampling Bootstrap (FAB)

What is special about the loss function?

- α -divergence

$$D_\alpha(p\|q_\theta) = -\frac{1}{\alpha(1-\alpha)} \int p(x)^\alpha q_\theta(x)^{1-\alpha} dx$$



$$D_{\alpha=2}(p\|q_\theta) \propto \int \frac{p(x)^2}{q_\theta(x)} dx$$

Approximate

$$L_{\text{FAB}} = - \sum_{i=1}^N \frac{\bar{w}_{\text{AIS}}^{(i)}}{\sum_j \bar{w}_{\text{AIS}}^{(j)}} q_\theta(\bar{x}_{\text{AIS}}^{(i)})$$

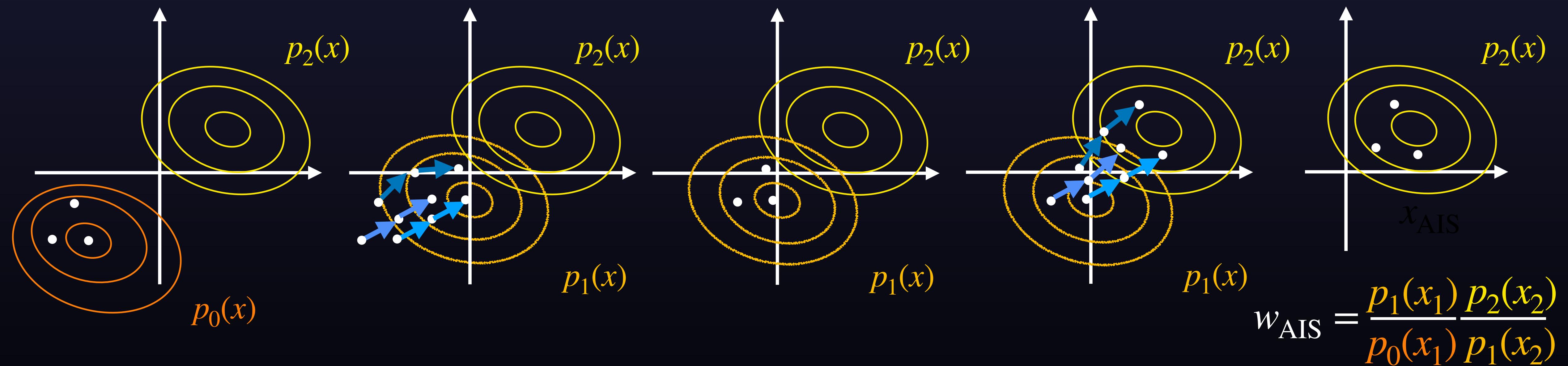
Flow Annealed Importance Sampling Bootstrap (FAB)

How does Annealed Importance Sampling (AIS) work?

- Interpolate between two distributions p_0 and p_M
 - Intermediate distribution $\log p_m(x) = \beta_m \log p_0(x) + (1 - \beta_m) \log p_M(x)$ $\beta_0 = 1 > \dots > \beta_m > \dots > \beta_M = 0$
- Move samples from previous p_m to next p_{m+1} via HMC
 - Track importance weight $w_{\text{AIS}} = \frac{p_1(x_1)}{p_0(x_1)} \frac{p_2(x_1)}{p_1(x_1)} \dots \frac{p_M(x_M)}{p_{M-1}(x_M)}$

Flow Annealed Importance Sampling Bootstrap (FAB)

How does Annealed Importance Sampling (AIS) work?



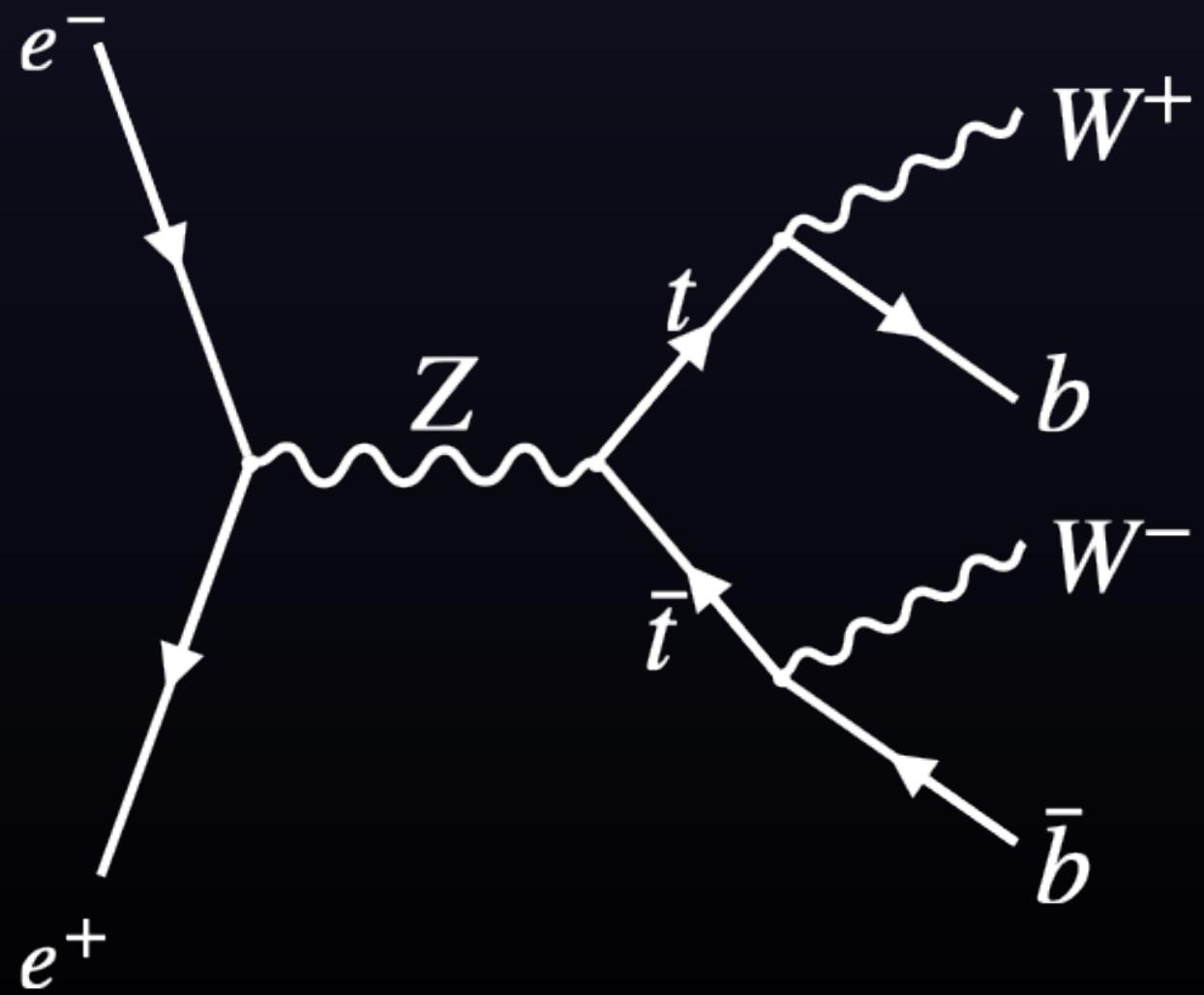
Flow Annealed Importance Sampling Bootstrap (FAB)

Replay buffer

Algorithm 7: Prioritized Replay Buffer for FAB

```
1 Initialize flow  $q_\theta(x)$ 
2 Initialize replay buffer to a fixed maximum size
3 for  $i = 0$  to  $K$  do
4   # Generate batch  $B$  of AIS samples and add to buffer:
5   Sample batch  $x^{(1:B)}$  from flow  $q_\theta(x)$  and evaluate  $\log q_\theta(x^{(1:B)})$ 
6   Run AIS with target  $d^2/q_\theta$  starting from  $x^{(1:B)}$  and  $\log q_\theta(x^{(1:B)})$  to get  $x_{\text{AIS}}^{(1:B)}$  and
     $\log w_{\text{AIS}}^{(1:B)}$ 
7   Add  $x_{\text{AIS}}^{(1:B)}$ ,  $\log w_{\text{AIS}}^{(1:B)}$ , and  $\log q_\theta(x_{\text{AIS}}^{(1:B)})$  to replay buffer
8   for  $j = 1$  to  $J$  do
9     # Sample batch  $b$  from buffer and update flow:
10    Sample  $x_{\text{AIS}}^{(1:b)}$  and  $\log q_\theta(x_{\text{AIS}}^{(1:b)})$  from buffer with probability  $p \propto w_{\text{AIS}}^{(1:b)}$ 
11    Calculate  $\log w_{\text{correction}}^{(1:b)} = \log q_{\theta_{\text{old}}}(x_{\text{AIS}}^{(1:b)}) + \text{stopgrad}(\log q_\theta(x_{\text{AIS}}^{(1:b)}))$ 
12    Update  $\log w_{\text{AIS}}^{(1:b)} \leftarrow \log w_{\text{AIS}}^{(1:b)} + \log w_{\text{correction}}^{(1:b)}$  in buffer
13    Update  $\log q_{\theta_{\text{old}}}(x_{\text{AIS}}^{(1:b)}) \leftarrow \log q_\theta(x_{\text{AIS}}^{(1:b)})$ 
14    Calculate loss  $S'(\theta) = -\frac{1}{N} \sum_i w_{\text{correction}}^{(i)} \log q_\theta(x_{\text{AIS}}^{(1:b)})$ 
15    Perform gradient descent on  $S'(\theta)$  to update  $\theta$ 
```

Results: $e^+e^- \rightarrow t\bar{t}$,
 $t \rightarrow W^+b, \bar{t} \rightarrow W^-b$
(8D)



8D: $e^+e^- \rightarrow t\bar{t}$,
 $t\bar{t} \rightarrow W^+b, \bar{t} \rightarrow W^- \bar{b}$

Corner Plot

