

# Generative AI as the new frontier of sampling algorithms in Lattice QFT

---

Kim A. Nicoli

University of Bonn, TRA Matter, HISKP (Helmholtz-Instituts für Strahlen- und Kernphysik)

Talk based on [PRE 101, 023304 \(2020\)](#), [PRL 126, 032001 \(2021\)](#), [PRD 108, 114501 \(2023\)](#)

In collaboration with: P. [Kessel](#), S. [Nakajima](#), C. [Anders](#), P. [Stornati](#), L. [Funcke](#), T. [Hartung](#), K. [Jansen](#)

KI

write a poem which explains what is the standard model of particle physics



The universe is filled with particles, so small and so precise,  
Each with a role to play, each with a unique size.  
To understand them all, we need a model that can describe,  
The way they interact, the way they survive.

The Standard Model is that model, a framework to behold,  
It tells us what's inside, what's hot, what's cold.  
It's built upon particles, with quarks and leptons too,  
And bosons that carry force, like glue.

Six types of quarks, and six types of leptons too,  
All with different masses, all with different hues.  
The bosons that they use, are force carriers all,  
The photon, the W and Z, the gluon, and the Higgs, standing tall.

The photon is the light, that illuminates our world,  
The W and Z, the weak force, that makes some particles swirled.  
The gluon, the strong force, that binds quarks together tight,  
And the Higgs, that gives mass, without which particles take flight.



The universe is filled with particles, so small and so precise,  
Each with a role to play, each with a unique size.  
To understand them all, we need a model that can describe,  
The way they interact, the way they survive.

The Standard Model is that model, a framework to behold,  
It tells us what's inside, what's hot, what's cold.  
It's built upon particles, with quarks and leptons too,  
And bosons that carry force, like glue.

Six types of quarks, and six types of leptons too,  
All with different masses, all with different hues.  
The bosons that they use, are force carriers all,  
The photon, the W and Z, the gluon, and the Higgs, standing tall.

The photon is the light, that illuminates our world,  
The W and Z, the weak force, that makes some particles swirled.  
The gluon, the strong force, that binds quarks together tight,  
And the Higgs, that gives mass, without which particles take flight.

The Standard Model explains, how particles combine,

To form all we see, to make the universe shine.  
It's not complete, with mysteries to unfold,  
But it's a mighty tool, in the hands of the bold.

So when you think of particles, and all they can achieve,  
Remember the Standard Model, and what it can conceive.  
For in its tiny world, lies the secrets of the universe,  
And with it, we can unlock, a future that's diverse.



The universe is filled with particles, so small and so precise,  
Each with a role to play, each with a unique size.  
To understand them all, we need a model that can describe,  
The way they interact, the way they survive.

The Standard Model is that model, a framework to behold,  
It tells us what's inside, what's hot, what's cold.  
It's built upon particles, with quarks and leptons too,  
And bosons that carry force, like glue.

Six types of quarks, and six types of leptons too,  
All with different masses, all with different hues.  
The bosons that they use, are force carriers all,  
The photon, the W and Z, the gluon, and the Higgs, standing tall.

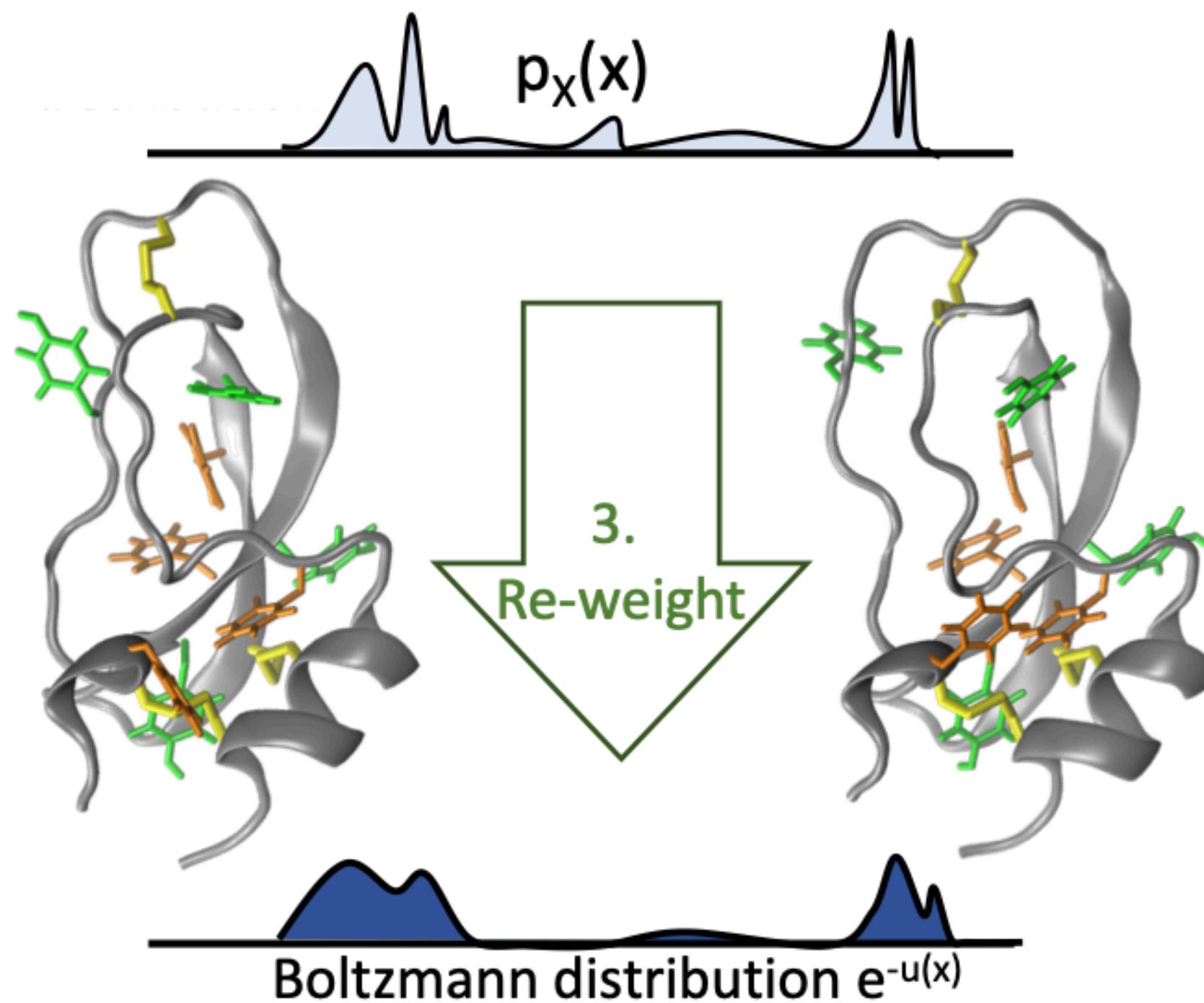
The photon is the light, that illuminates our world,  
The W and Z, the weak force, that makes some particles swirled.  
The gluon, the strong force, that binds quarks together tight,  
And the Higgs, that gives mass, without which particles take flight.

The Standard Model explains, how particles combine,  
To form all we see, to make the universe shine.  
It's not complete, with mysteries to unfold,  
But it's a mighty tool, in the hands of the bold.

So when you think of particles, and all they can achieve,  
Remember the Standard Model, and what it can conceive.  
For in its tiny world, lies the secrets of the universe,  
And with it, we can unlock, a future that's diverse.

# Generative AI for Physics : Past, Present and Future

**F. Noé, et al., Science, eaaw1147 (2019)**



(Image credits: F. Noé, et al., Science, eaaw1147 (2019))

# Generative AI for Physics : Past, Present and Future

## Scalar Field Theories ( $\phi^4$ )

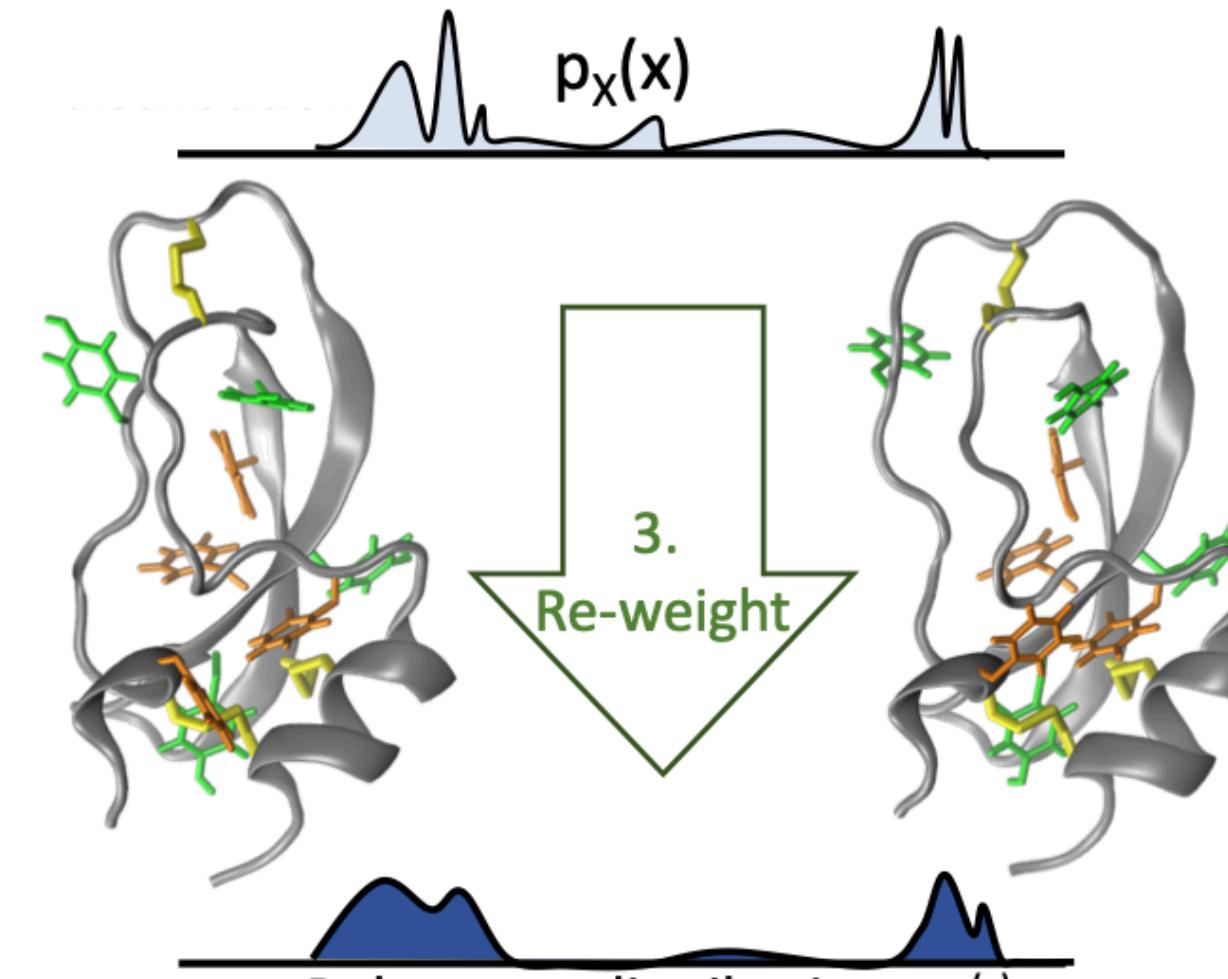
- [M. S. Albergo, et al., Phys. Rev. D 100, 034515 \(2019\)](#)
- [K. A. Nicoli, et al., Phys. Rev. Lett. 126, 032001 \(2021\)](#)
- [P. deHaan, et al., arXiv: 2110.02673 @ ML4Pys workshop \(2021\)](#)
- [L. Vaitl et al., arXiv: 2206.09016 @ ICML \(2022\)](#)
- [A. Matthews et al., arXiv:2201.13117 @ ICML \(2022\)](#)
- [M. Caselle, et al., J. High Energ. Phys. 2022, 15 \(2022\)](#)
- [M. Gerdes, et al., SciPost Phys. 15, 238 \(2023\)](#)
- [A. Singha, et al., Phys. Rev. D 107, 014512 \(2023\)](#)

## Lattice Gauge Theories (U(1), SU( $N$ ))

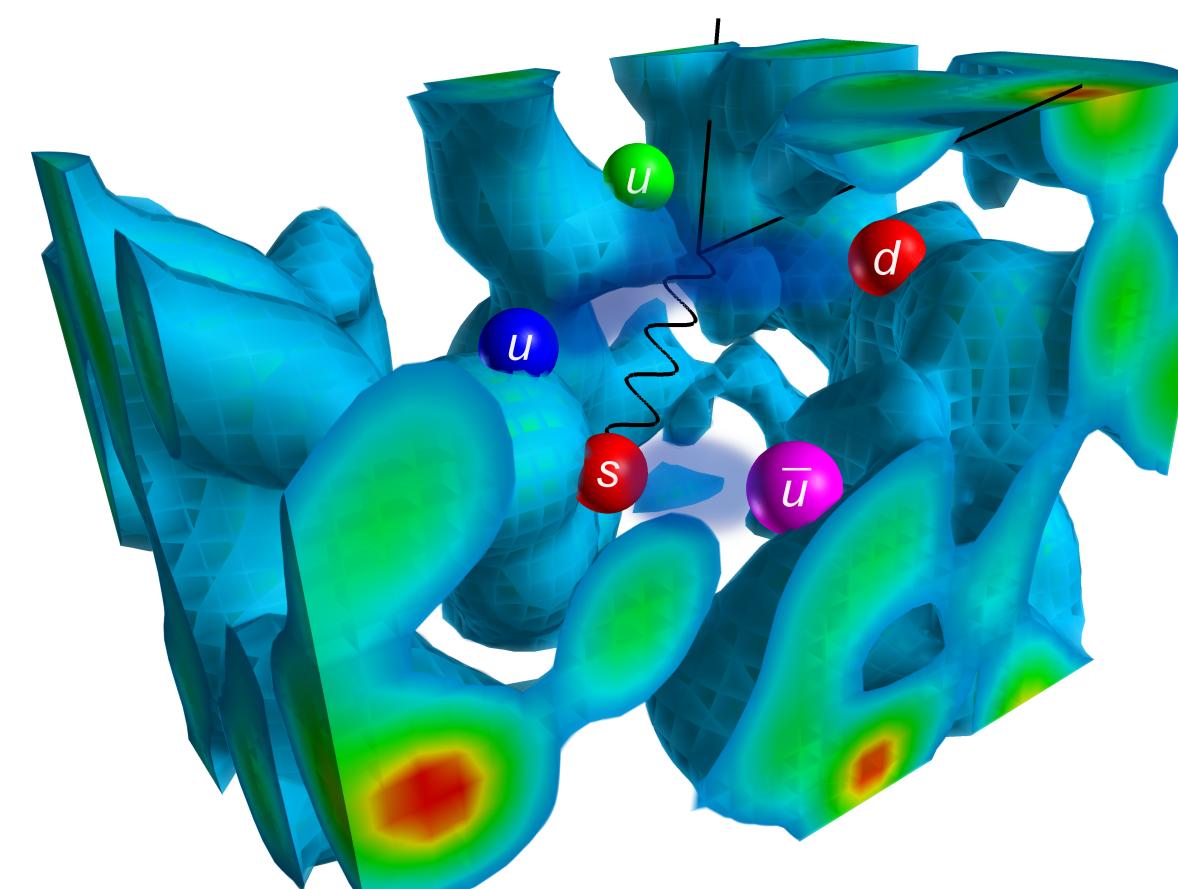
- [G. Kanwar, et al., Phys. Rev. Lett. 125, 121601 \(2020\)](#)
- [S. Bacchio, et al., Phys. Rev. D 107, L051504 \(2023\)](#)
- [R. Abbott et al., Phys. Rev. D 106, 074506 \(2022\)](#)
- [M. S. Albergo, et al., Phys. Rev. D 106, 014514 \(2022\)](#)
- [R. Abbott, et al., arXiv:2305.02402 \(2023\)](#)
- [J. Finkenrath, arXiv: 2201.02216 \(2022\)](#)

## Sampling Multimodal Densities in QFT

- [D.C. Hackett et al., arXiv:2107.00734 \(2021\)](#)
- [K. A. Nicoli, et al., arXiv: 2111.11303 @ LATTICE21 \(2021\)](#)
- [K. A. Nicoli, et al., Phys. Rev. D 108, 114501 \(2023\)](#)
- [B. Maté et al., TMLR 2835-8856 \(2023\)](#)
- [V. Kanaujia et al., arXiv:2401.15948 \(2024\)](#)



(Image credits: F. Noé, et al., Science, eaaw1147 (2019))



(Image credits: Lattice QCD © Derek Leinweber/CSSM/University of Adelaide)

## Scaling to Larger Lattices

- [L. Del Debbio, et al., Phys. Rev. D 104, 094507 \(2021\)](#)
- [R. Abbott et al., Eur. Phys. J. A 59, 257 \(2023\)](#)
- [A. Faraz et al., arXiv:2308.08615 \(2022\)](#)
- [B. Maté et al., arXiv: 2401.00828 \(2024\)](#)
- [R. Abbott et al., arXiv: 2401.10874v1 \(2024\)](#)
- [J. Finkenrath, arXiv: 2402.12176 \(2024\)](#)

## Autoregressive Models in Stat. Mech.

- [D. Wu et al., Phys. Rev. Lett. 122 \(8\), 080602 \(2019\)](#)
- [K. A. Nicoli, et al., Phys. Rev. E 101 \(2\), 023304 \(2019\)](#)
- [P. Bialas et al., Computer Physics Communications 281 \(2022\)](#)
- [P. Bialas et al., Phys. Rev. E 107 \(1\), 015303 \(2023\)](#)

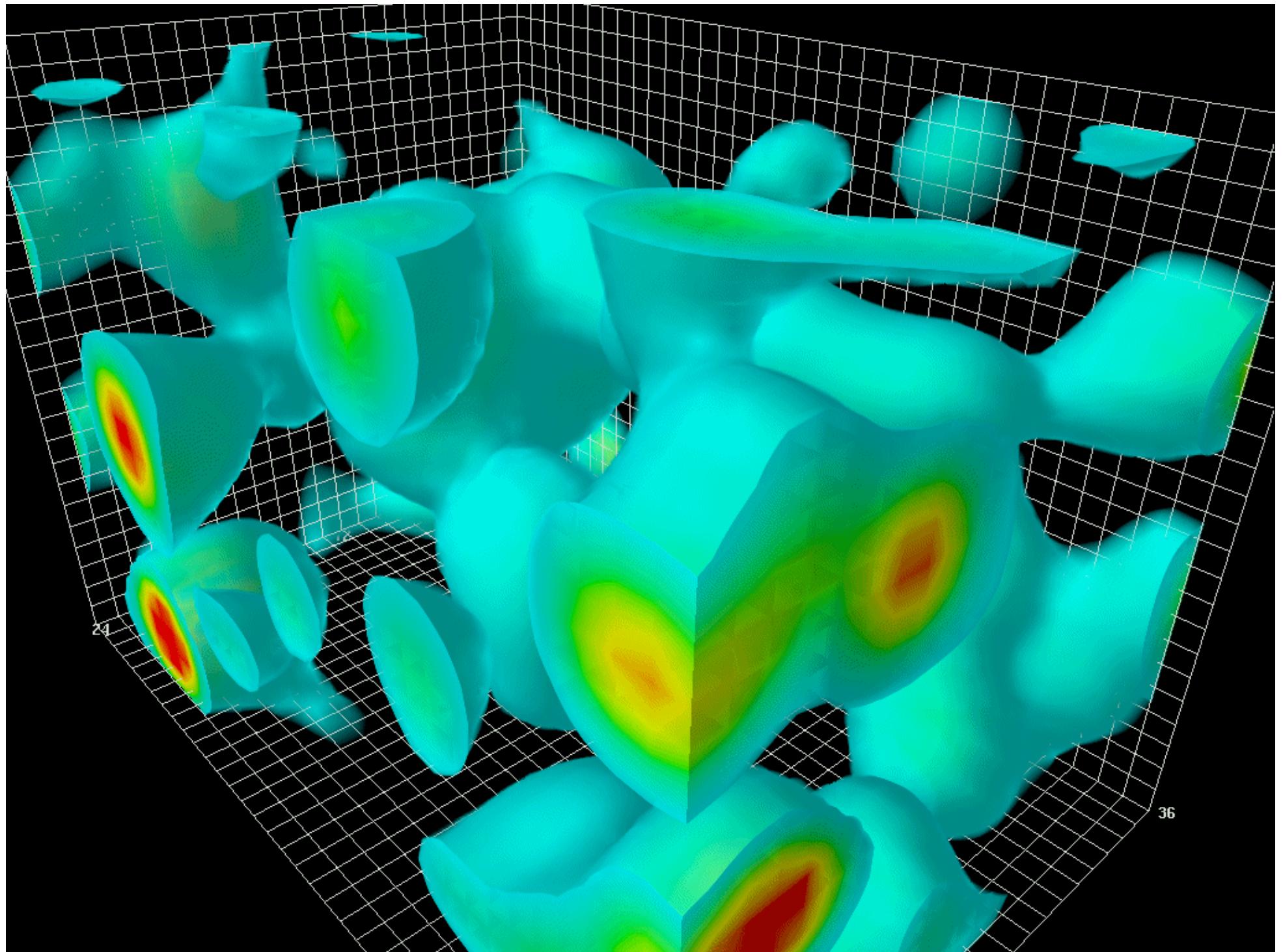
## Interesting Readings

- [F. Noé, et al., Science, eaaw1147 \(2019\)](#)
- [S. Chen, et al., Phys. Rev. D 107, 056001 \(2022\)](#)
- [B. Maté et al., arXiv: 2210.13772 \(2022\)](#)
- [L. Vaitl et al., MLST 3 \(4\), 045006 \(2022\)](#)
- [Caselle, M., et al., J. High Energ. Phys. 2024, 48 \(2024\)](#)
- [Cranmer K. et al., Nature Reviews Physics 5 \(9\), 526-535](#)

And much more...

# Quantum Field Theory

---



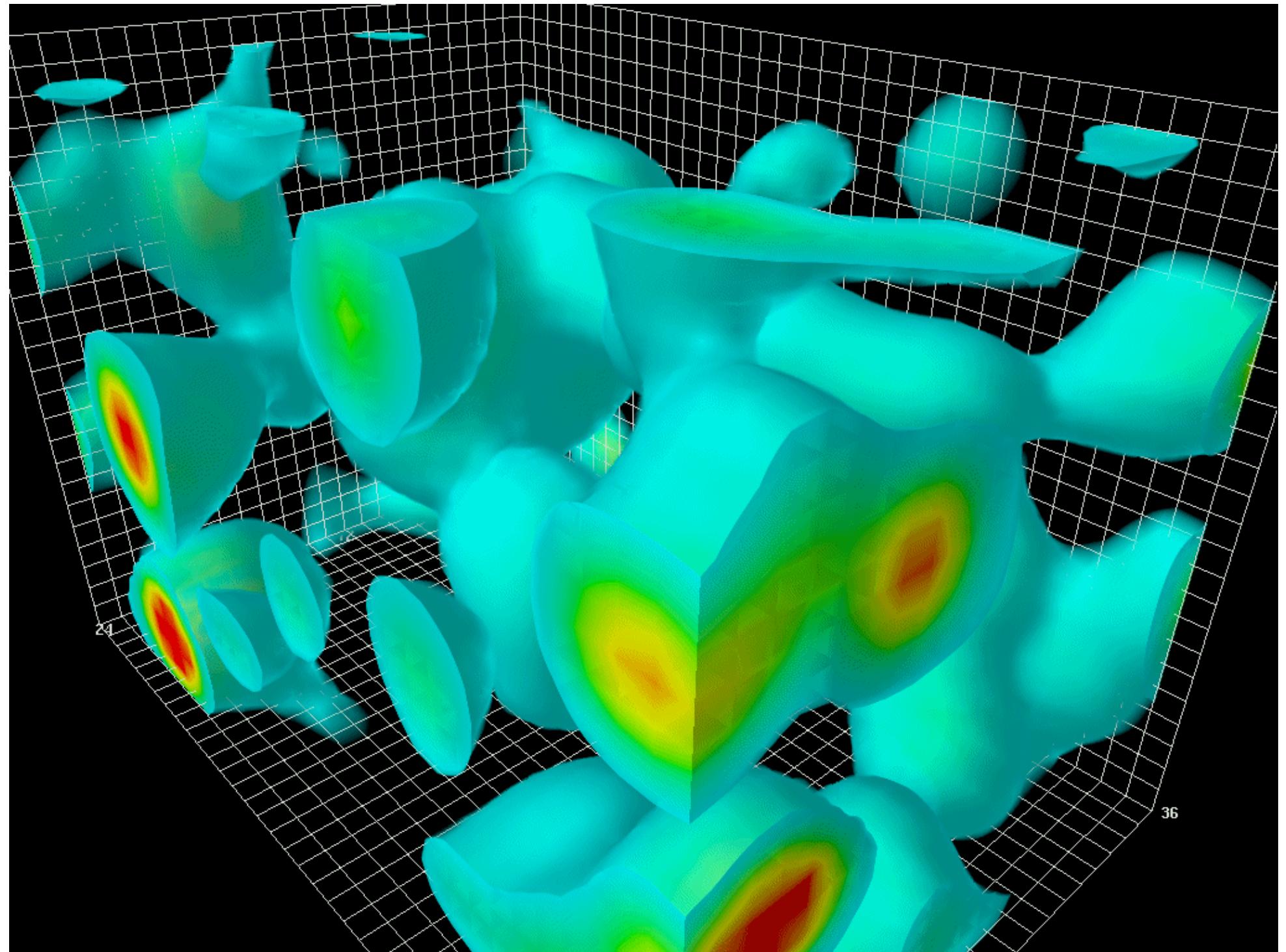
(credits: [Derek Leinweber](#))

“ [...] [A quantum field] [...] a bubbling soup of particles and anti-particles, constantly created and destroyed ”

D. Tong

# Quantum Field Theory

---

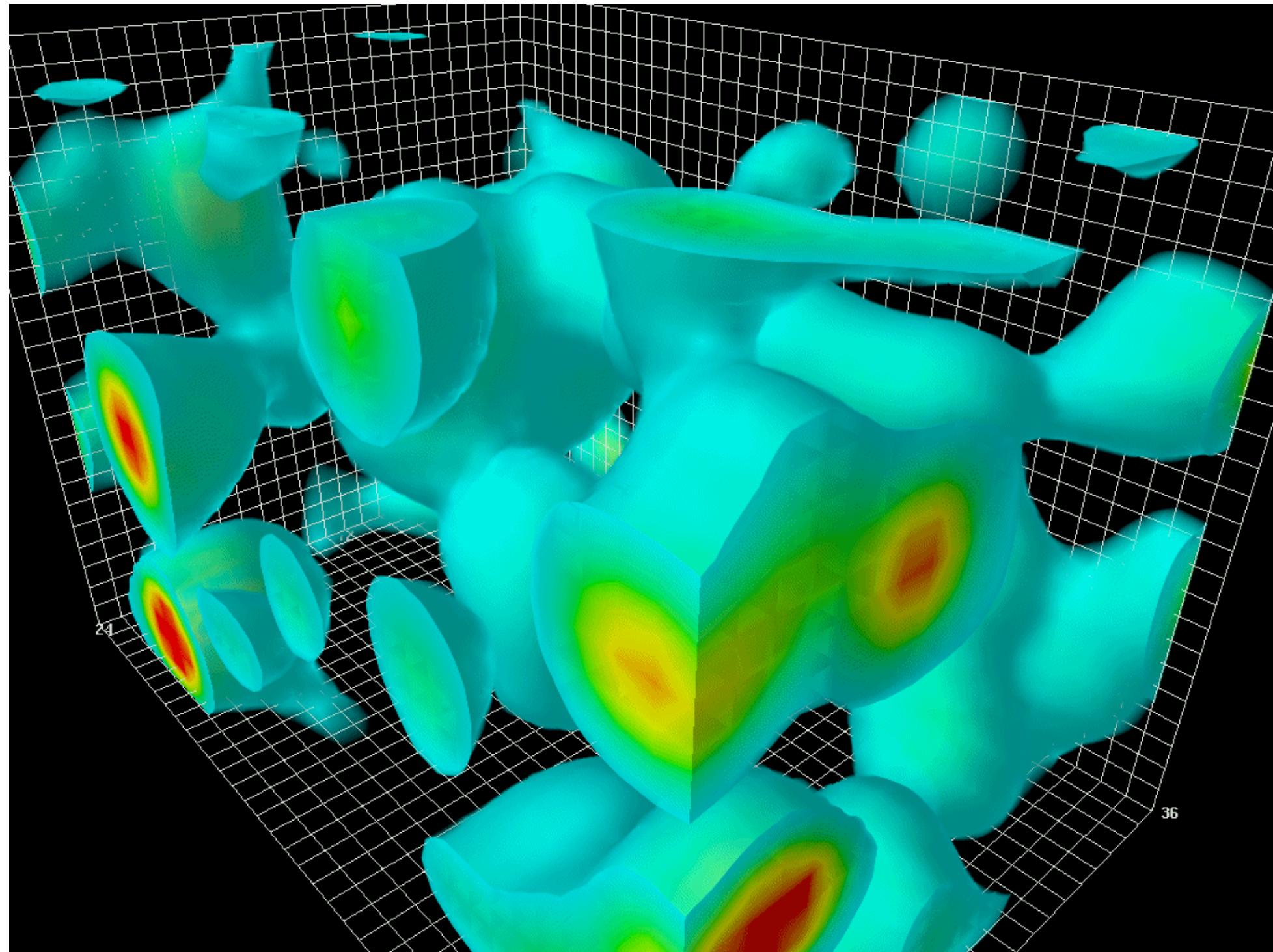


(credits: [Derek Leinweber](#))

“ [...] [A quantum field] [...] a bubbling soup of particles and anti-particles, constantly created and destroyed ”

D. Tong

# Quantum Field Theory



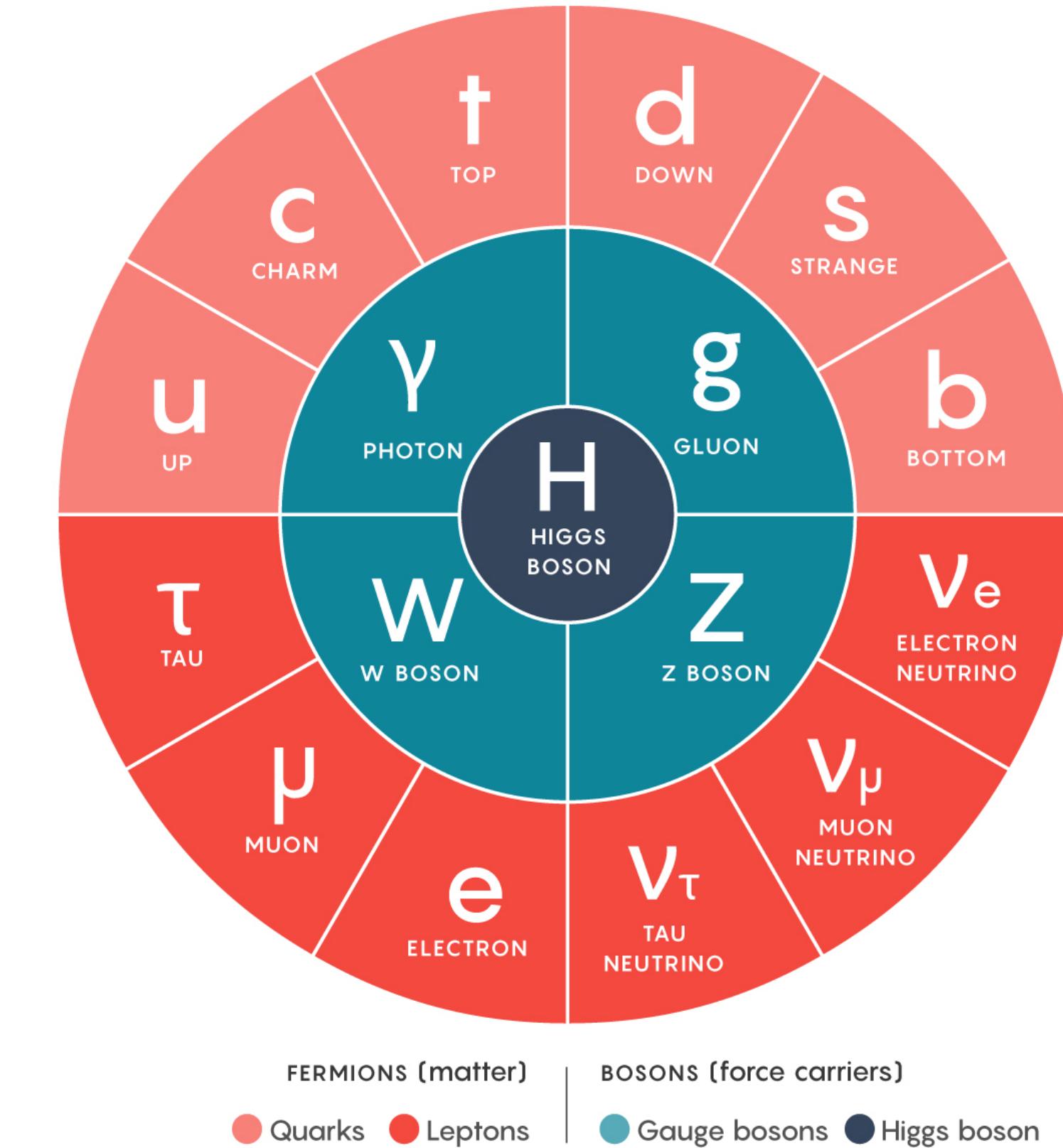
(credits: Derek Leinweber)

“ [...] [A quantum field] [...] a bubbling soup of particles and anti-particles, constantly created and destroyed ”

D. Tong

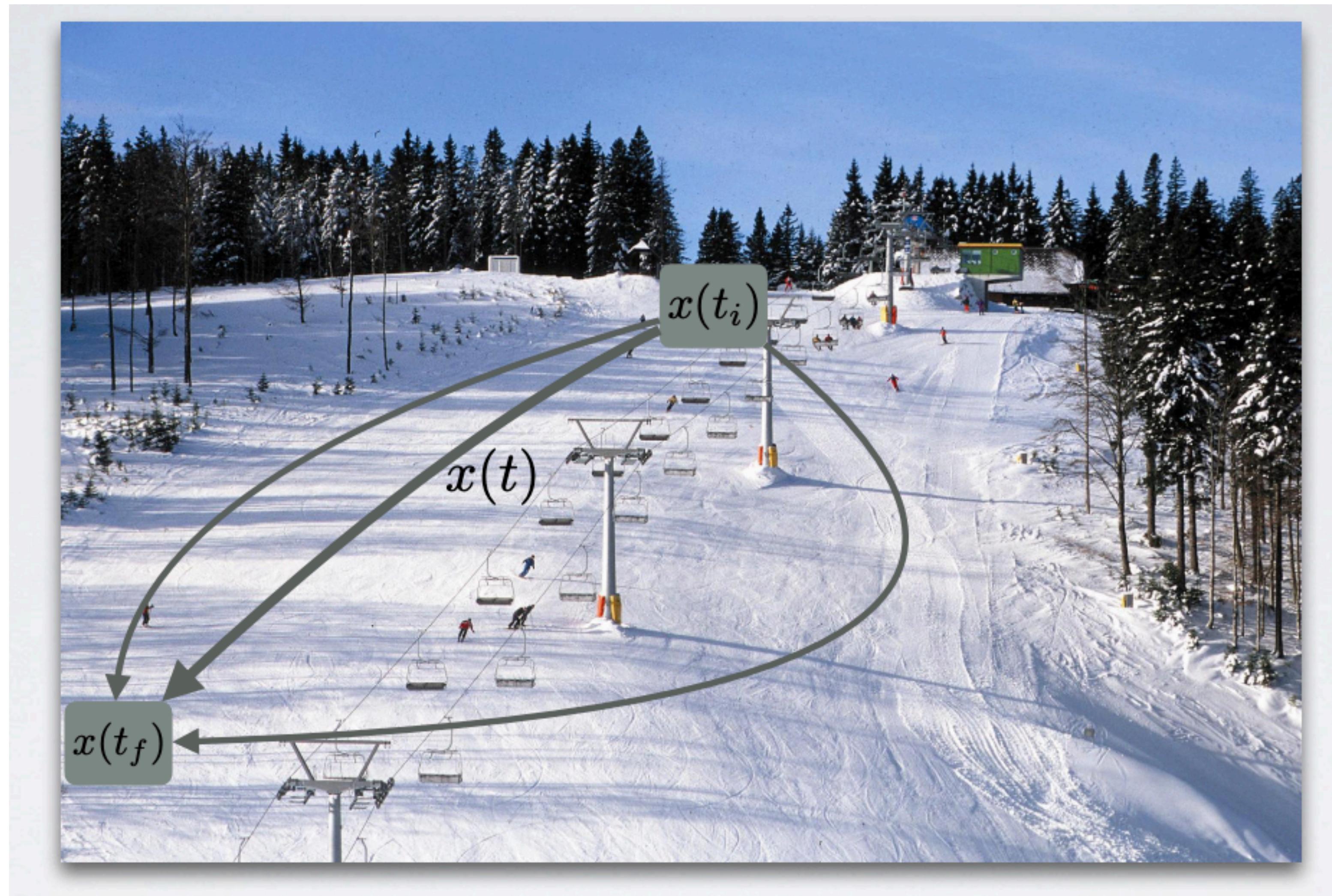


The Standard Model



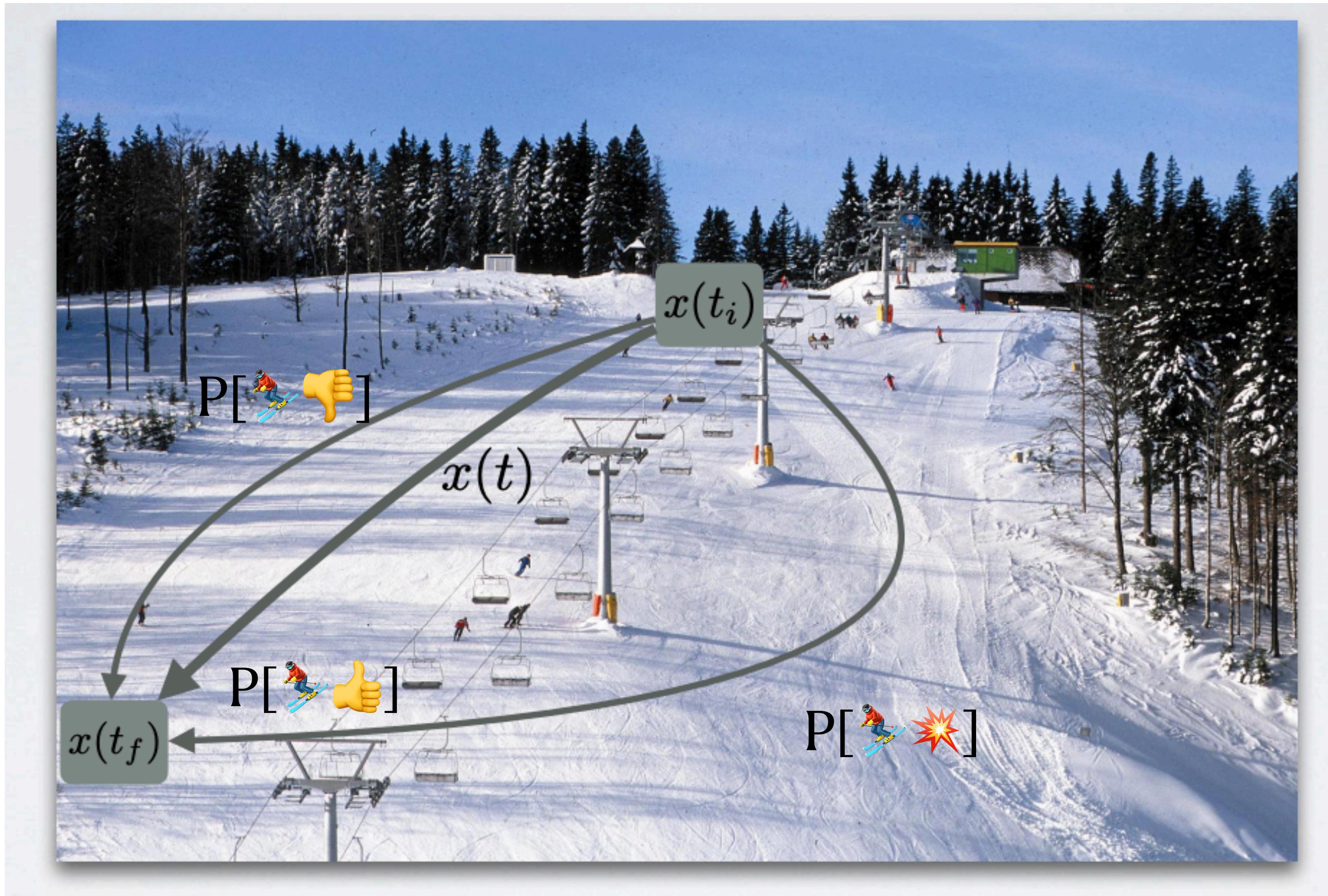
(Image credits: [Quanta Magazine](#))

# Path Integral Formulation



(Image credits: adapted from [Tom Luu](#))

# Path Integral Formulation



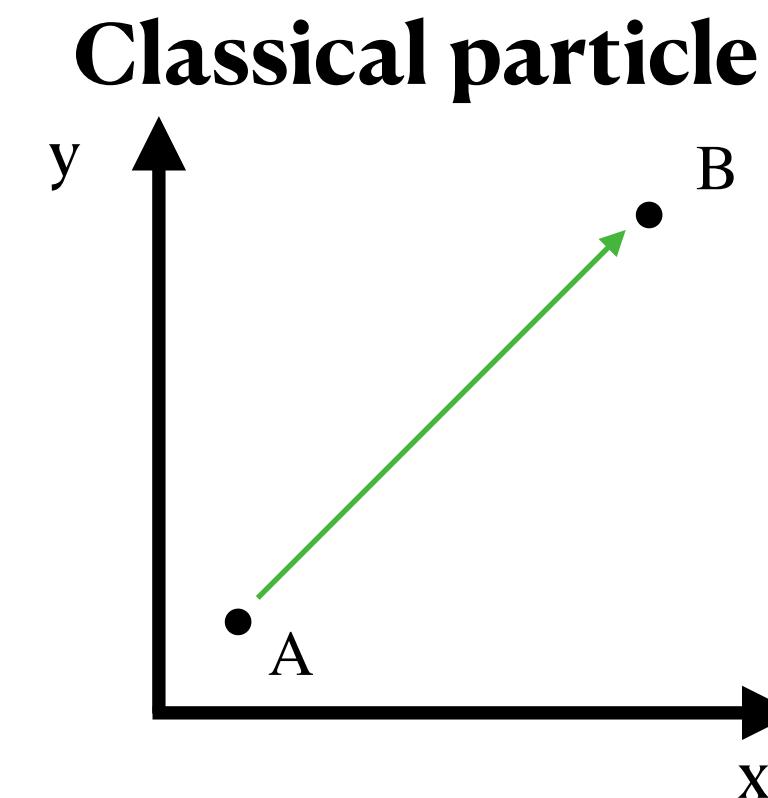
(Image credits: adapted from [Tom Luu](#))

# Path Integral Formulation

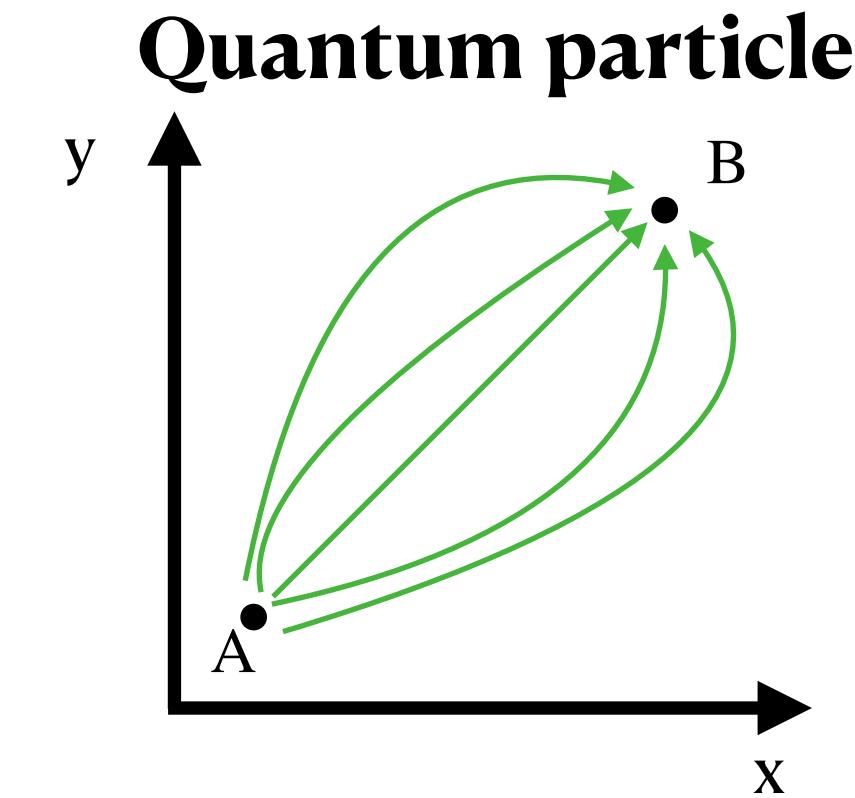
---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$



Least action  $\equiv$  Most probable



# Path Integral Formulation

---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$

# Path Integral Formulation

---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$

# Path Integral Formulation

---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$

# Path Integral Formulation

---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$

# Path Integral Formulation

---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$

# Path Integral Formulation

---

**Path-integral** is the basic tool for quantising fields and compute expectation values of physical observables  $\mathcal{O}$ :

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int D[\phi] \mathcal{O}(\phi) \exp\{-S(\phi)\}$$

computed over a Boltzmann-like density:

$$p(\phi) = \frac{e^{-S(\phi)}}{Z} \longrightarrow$$

known in **closed form** up to a **numerically intractable** normalisation

$$Z = \int D[\phi] e^{-S(\phi)}$$

# Lattice Quantum Field Theory

---

- The **path-integral** reduces from a **functional integral** to a **high-dimensional** ordinary integral

$$\langle \mathcal{O} \rangle = \int \prod_{x \in \Lambda} d[\phi(x)] \mathcal{O}(\phi) p(\phi)$$

- The field configuration  $\phi(x)$  is now a **random variable** of size  $\Lambda$  (lattice volume).
- With Markov-Chain Monte-Carlo algorithm we can **sample** and estimate observables.

$$\langle \mathcal{O} \rangle_p = \int D[\phi] \mathcal{O}(\phi) p(\phi) \approx \frac{1}{N} \sum_{i=1}^N \mathcal{O}(\phi_i)$$

where  $\phi_i \sim p$  (Boltzmann-like density)

# Markov Chain Monte-Carlo

---

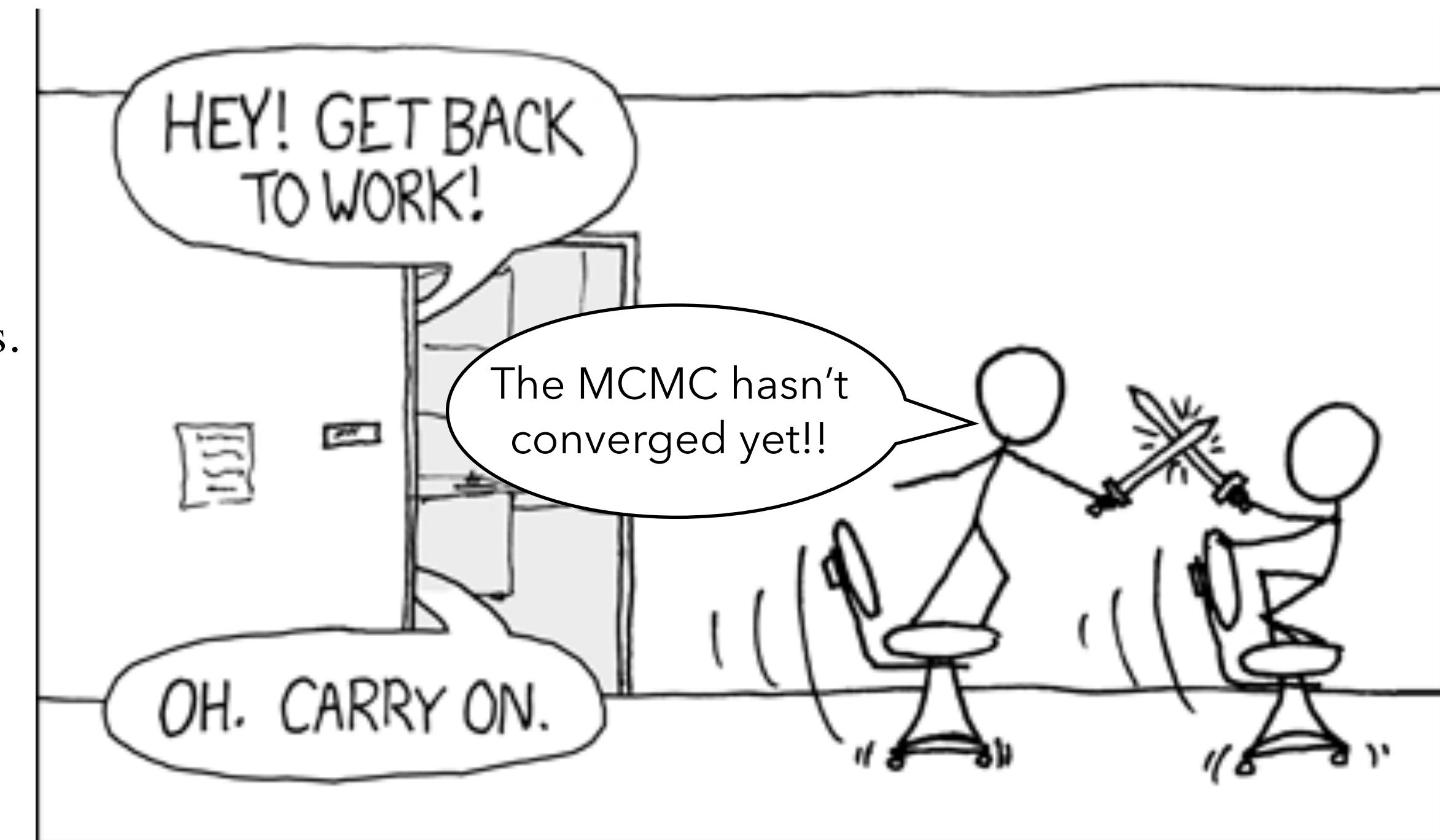
**MCMC**: sequentially proposeS the next sample and guarantees to eventually converge to a target density.

# Markov Chain Monte-Carlo

**MCMC**: sequentially proposeS the next sample and guarantees to eventually converge to a target density.

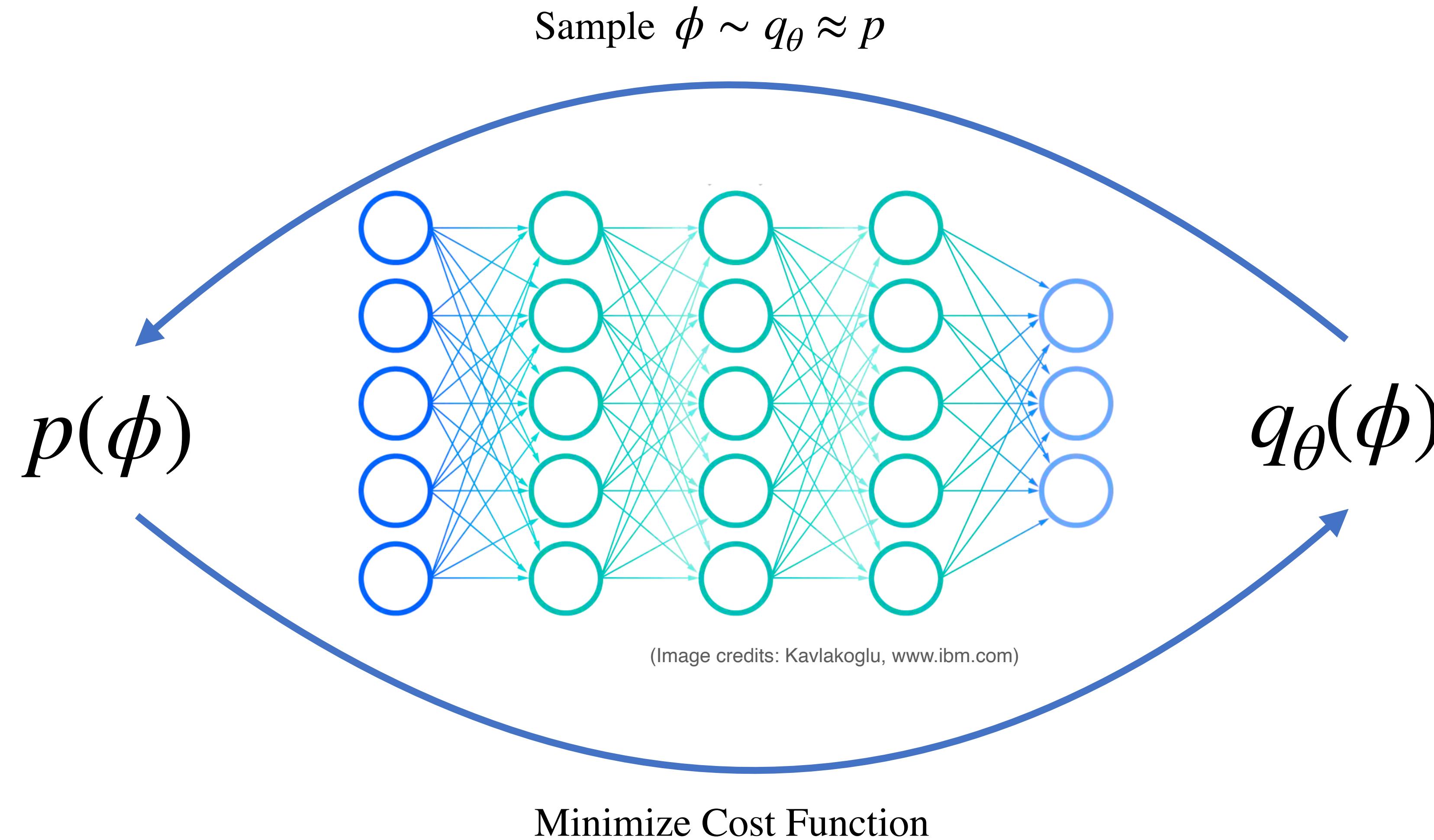
However MCMC algorithms come at a **cost**:

- 👎 **Sequential**  $\Rightarrow$  MCMC chains can't be parallelized.
- 👎 **Critical slowing down**  $\Rightarrow$  Struggles around phase transitions.
- 👎 Long range autocorrelations  $\Rightarrow$  **large statistical errors**.
- 👎 The partition function  $Z$  is **unknown**.
- 👎 No direct estimation of **thermodynamic observables**.

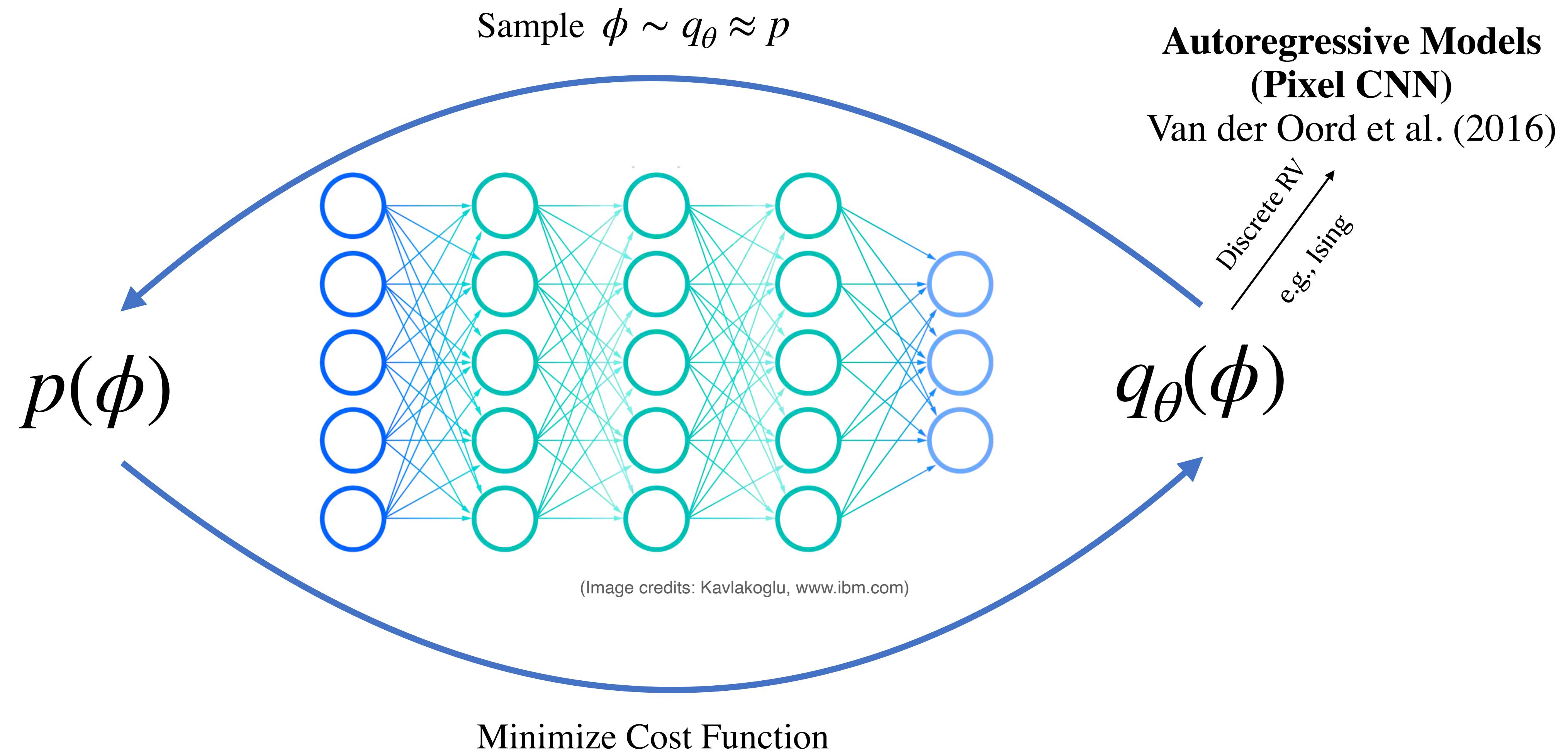


(Adapted from: [xkcd/303](#))

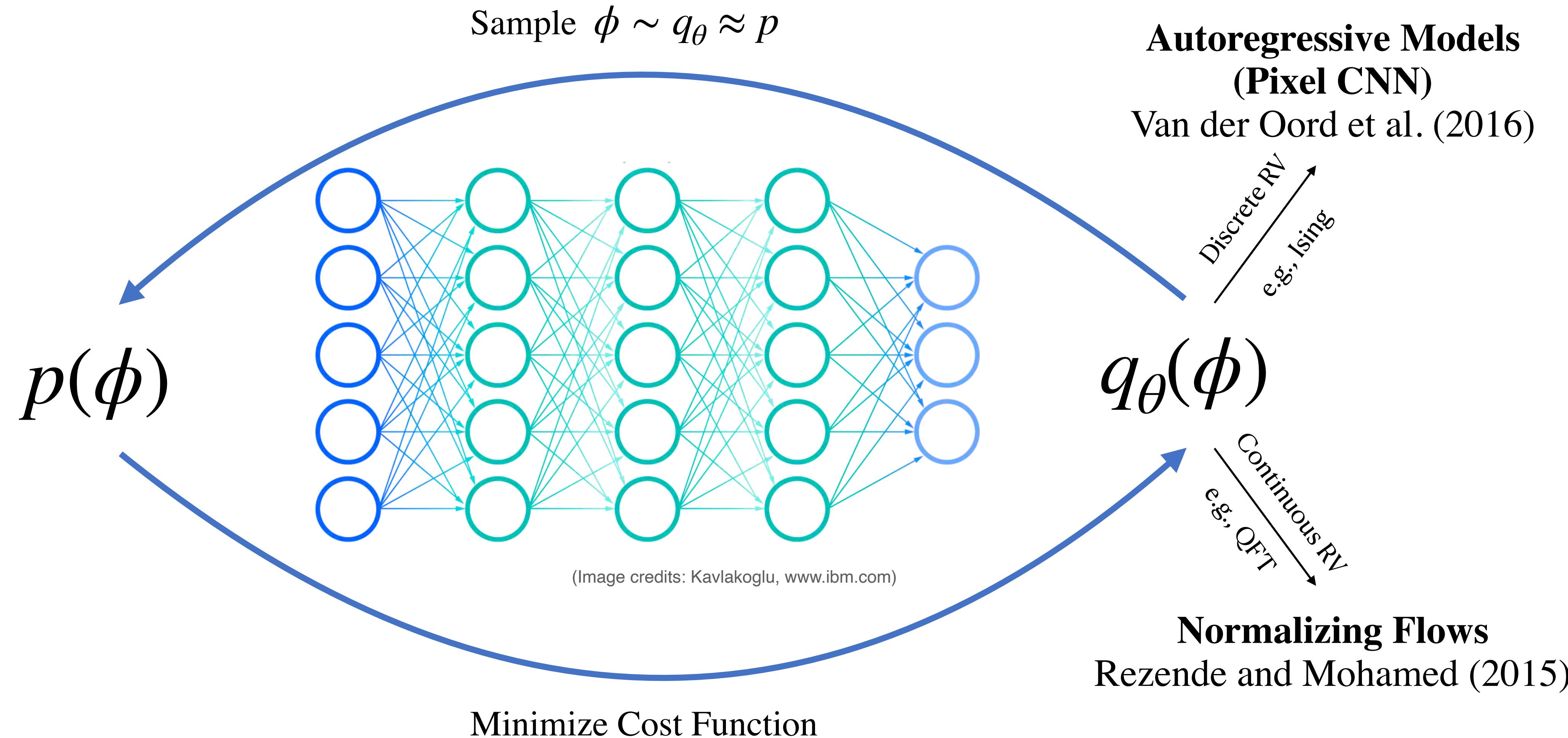
# Density Estimation with Deep Generative Models



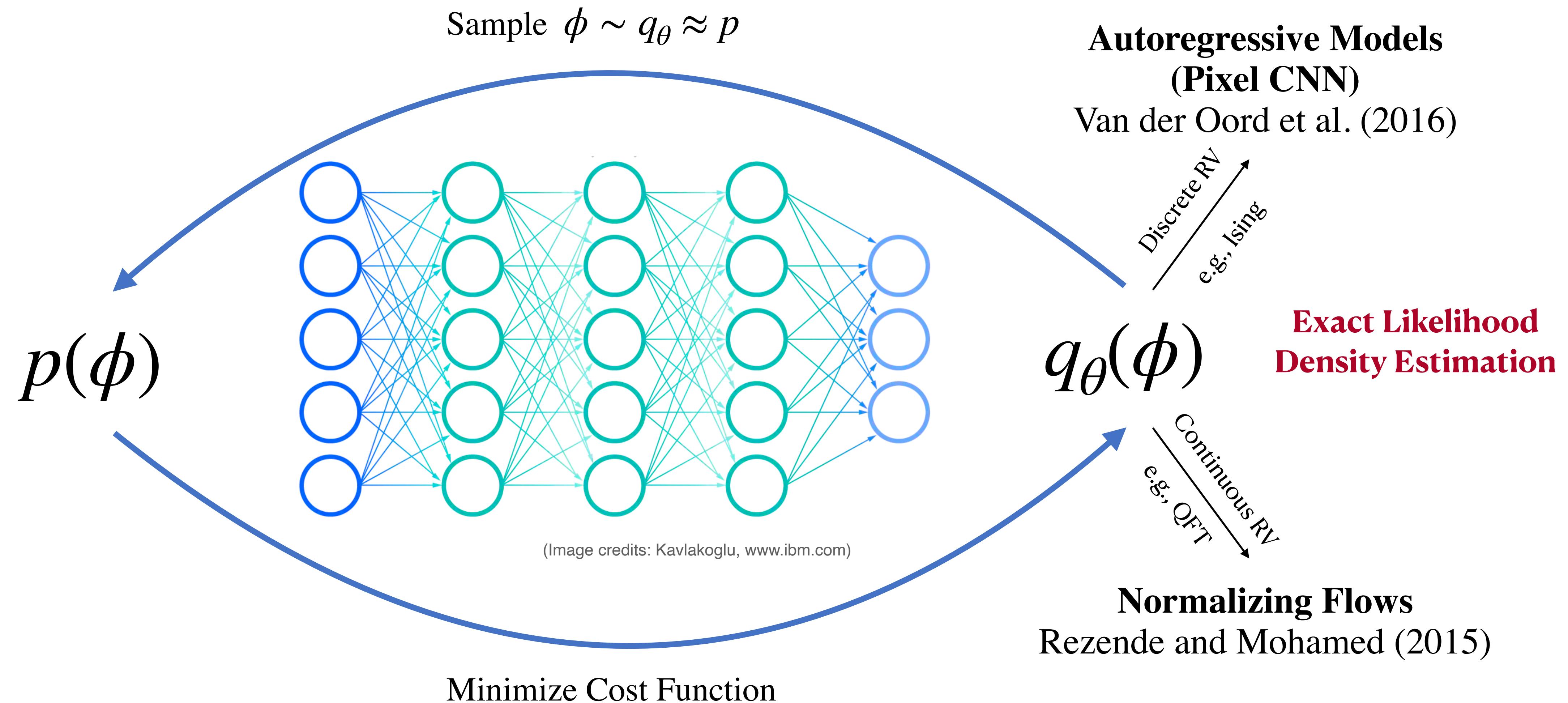
# Density Estimation with Deep Generative Models



# Density Estimation with Deep Generative Models



# Density Estimation with Deep Generative Models



# Taxonomy of Generative Models

---

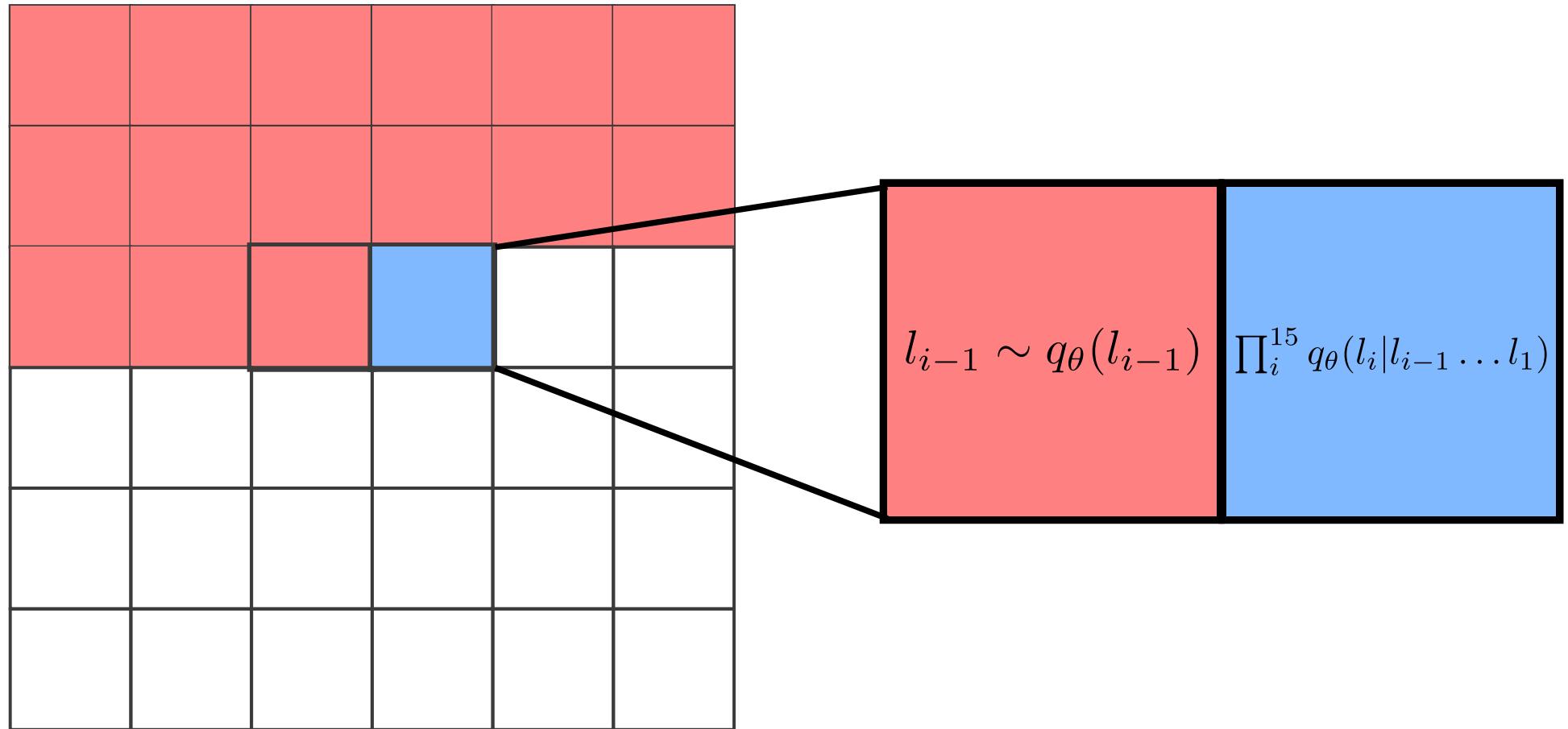
Table 2.1. **Exact likelihood estimation for different types of generative models:** overview on different types of deep generative models (**DGMs**) with few examples of their application in the physics domain. Generative adversarial networks (generative adversarial network (**GAN**)) usually do not have access to sampling probabilities while variational auto-encoders (**VAEs**) only allows for approximate ones. The cases of normalizing flows (**NFs**) and autoregressive neural networks (**ARNNs**) are instead substantially different as they allow access to normalized tractable densities.

DGM	sampling probability	Normalization
GAN	none	✗
VAE	approximate	✓
ARNN	exact	✓
NF	exact	✓
DDPM	approximate*	✓

\* can be exact under some circumstances, e.g., Sec 4.3 from [arXiv:2209.00796](https://arxiv.org/abs/2209.00796).

# Autoregressive Models

- ▶ Probability density of sampling a configuration  $L$  (**ancestral sampling**)



Nicoli, PhD diss., Technische Universität Berlin, 2023.

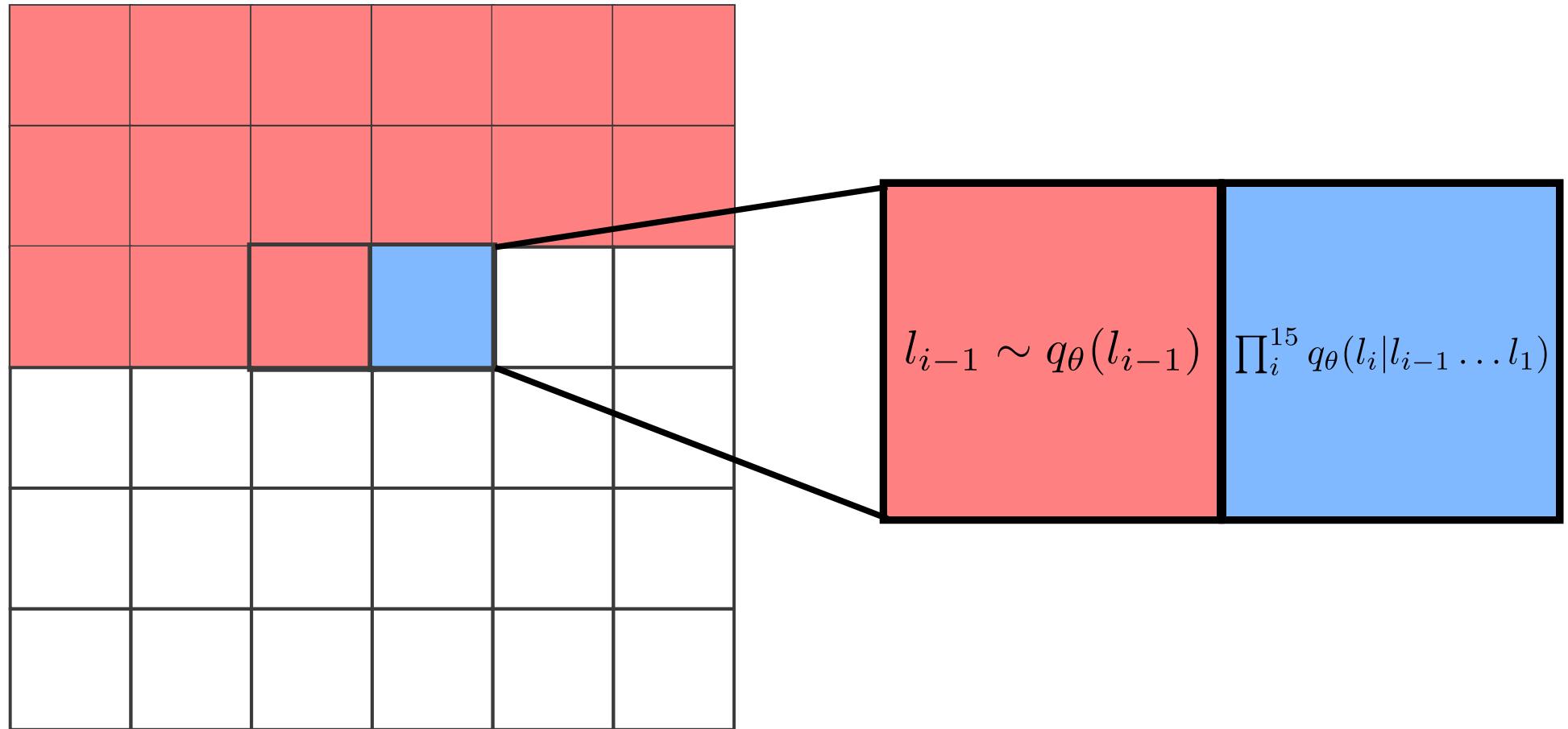
$$q_\theta(L) = \prod_{i=1}^{N^2} q_\theta(l_i | l_{i-1} \dots l_1)$$

## Images (PixelCNN)

A. Van der Oord et al., arXiv:1601.06759 (2016)

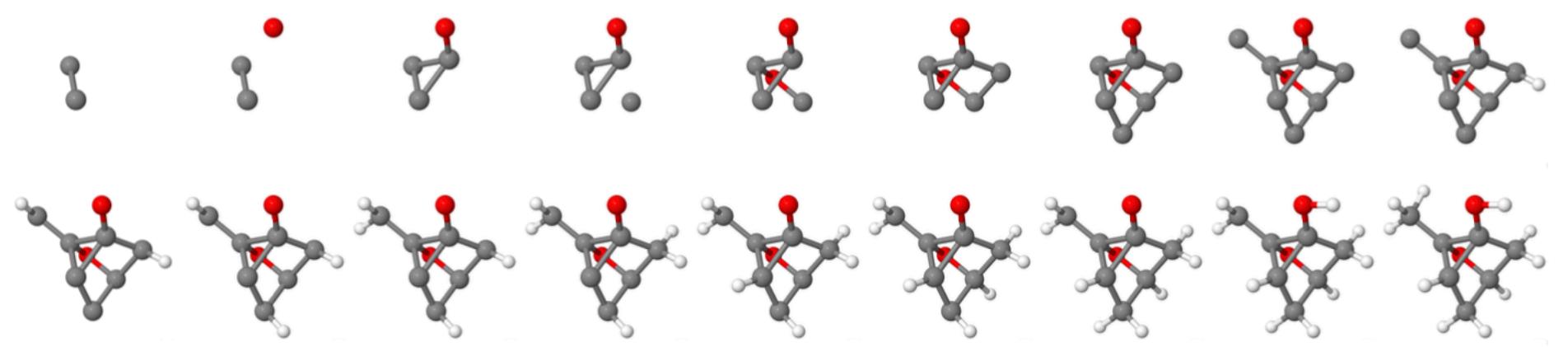
# Autoregressive Models

- Probability density of sampling a configuration  $L$  (**ancestral sampling**)



Nicoli, PhD diss., Technische Universität Berlin, 2023.

## Molecular Graphs



N. Gebauer et al., arXiv:1810.11347 (2018)

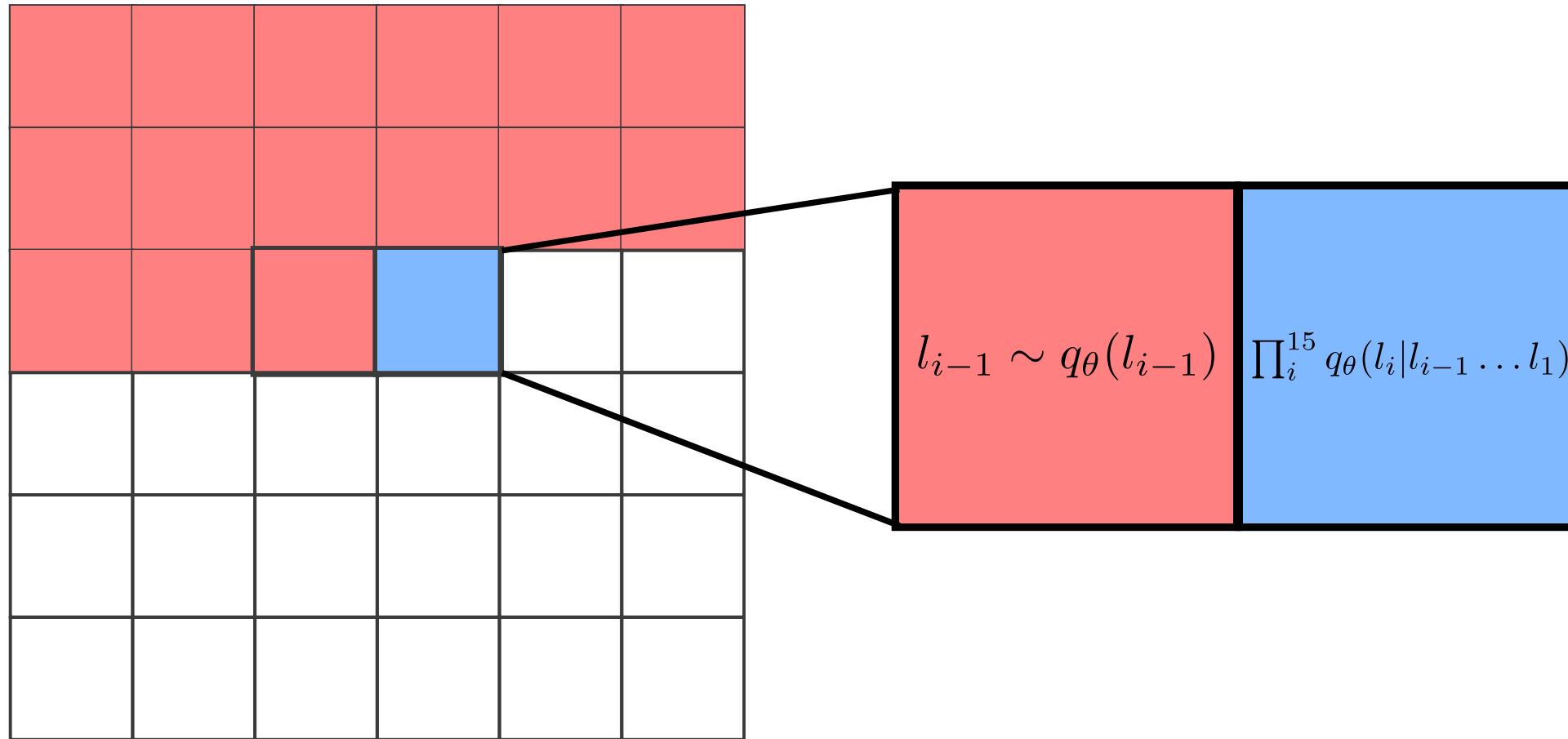
$$q_\theta(L) = \prod_{i=1}^{N^2} q_\theta(l_i | l_{i-1} \dots l_1)$$

## Images (PixelCNN)

A. Van der Oord et al., arXiv:1601.06759 (2016)

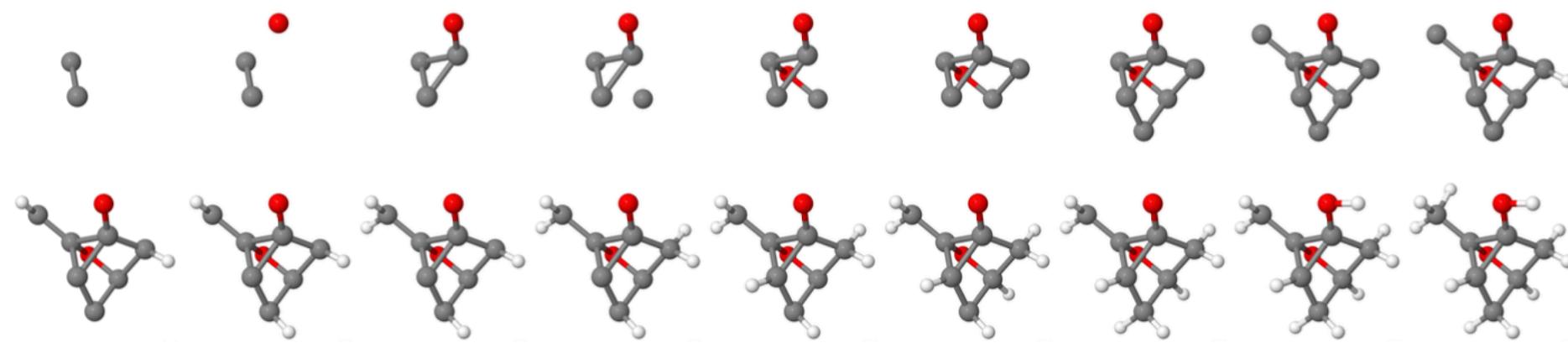
# Autoregressive Models

- Probability density of sampling a configuration  $L$  (**ancestral sampling**)



Nicoli, PhD diss., Technische Universität Berlin, 2023.

## Molecular Graphs



N. Gebauer et al., arXiv:1810.11347 (2018)

## Language Models (GPT)

“... quick brown fox **jumps** ...”  
 $p(\text{jumps} | \dots)$

T. Brown et al., arXiv: 2005.14165 (2020)

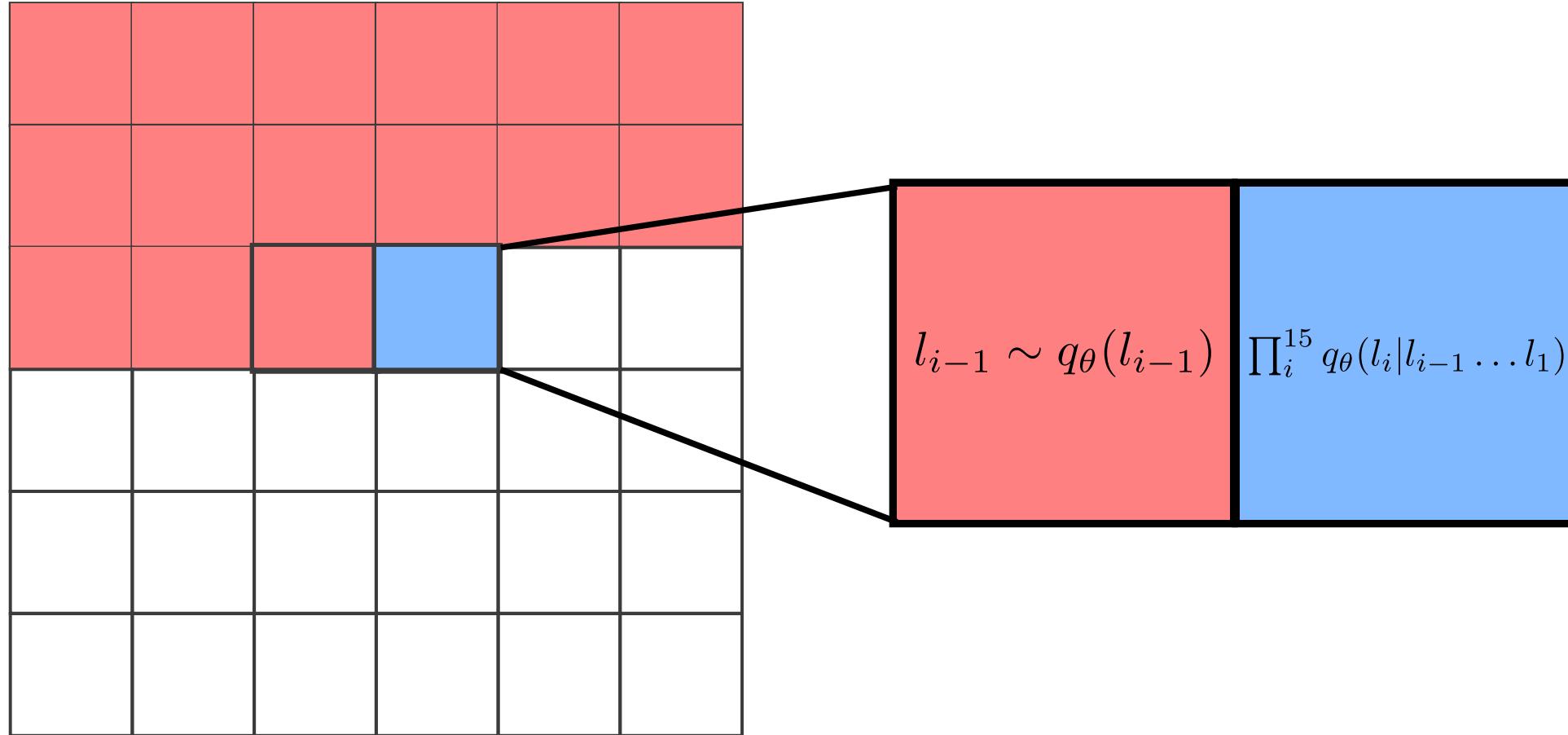
$$q_\theta(L) = \prod_{i=1}^{N^2} q_\theta(l_i | l_{i-1} \dots l_1)$$

## Images (PixelCNN)

A. Van der Oord et al., arXiv:1601.06759 (2016)

# Autoregressive Models

- Probability density of sampling a configuration  $L$  (**ancestral sampling**)



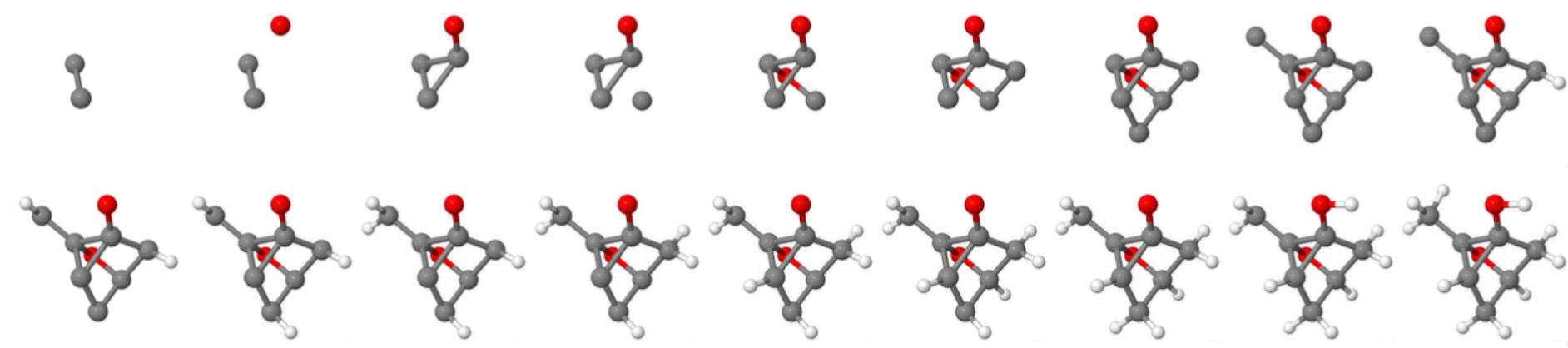
Nicoli, PhD diss., Technische Universität Berlin, 2023.

$$q_\theta(L) = \prod_{i=1}^{N^2} q_\theta(l_i | l_{i-1} \dots l_1)$$

## Images (PixelCNN)

A. Van der Oord et al., arXiv:1601.06759 (2016)

## Molecular Graphs



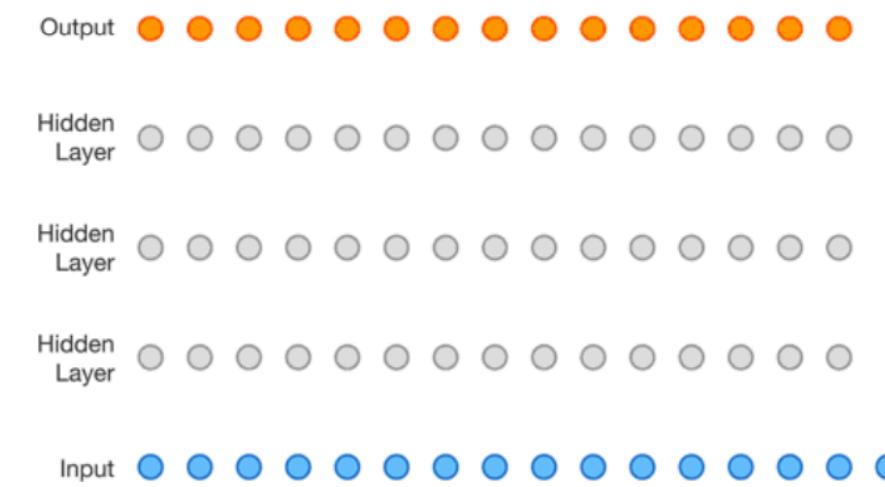
N. Gebauer et al., arXiv:1810.11347 (2018)

## Language Models (GPT)

“... quick brown fox **jumps** ...”  
 $p(\text{jumps} | \dots)$

T. Brown et al., arXiv: 2005.14165 (2020)

## Speech (WaveNet)



A. Van der Oord et al., arXiv:1609.03499 (2016)

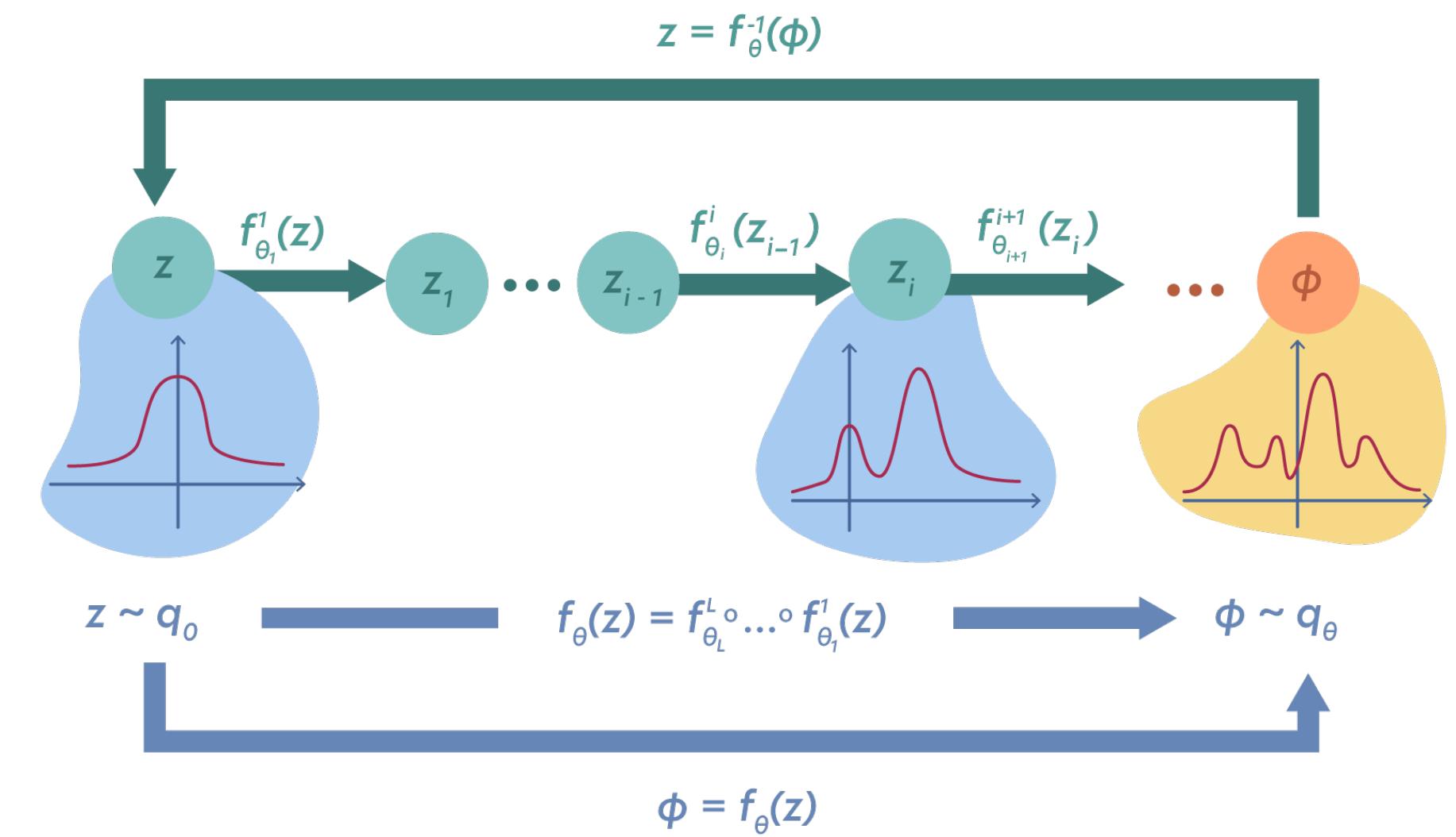
# Sampling with Normalizing Flows

We use a parametric function  $f_\theta$  (a **diffeomorphism**) to transform Gaussian samples  $z \sim q_0$  into physical configurations  $\phi \sim q_\theta$

$$f_\theta : z \in \mathcal{Z} \sim q_z \rightarrow x = f_\theta(z) \in \mathcal{X} \sim q_\theta .$$

The parametric function needs to fulfill certain criteria:

- **Bijective** transformation  $\phi = f_\theta(z)$
- **Invertible** and **differentiable**
- **Tractable** Jacobian



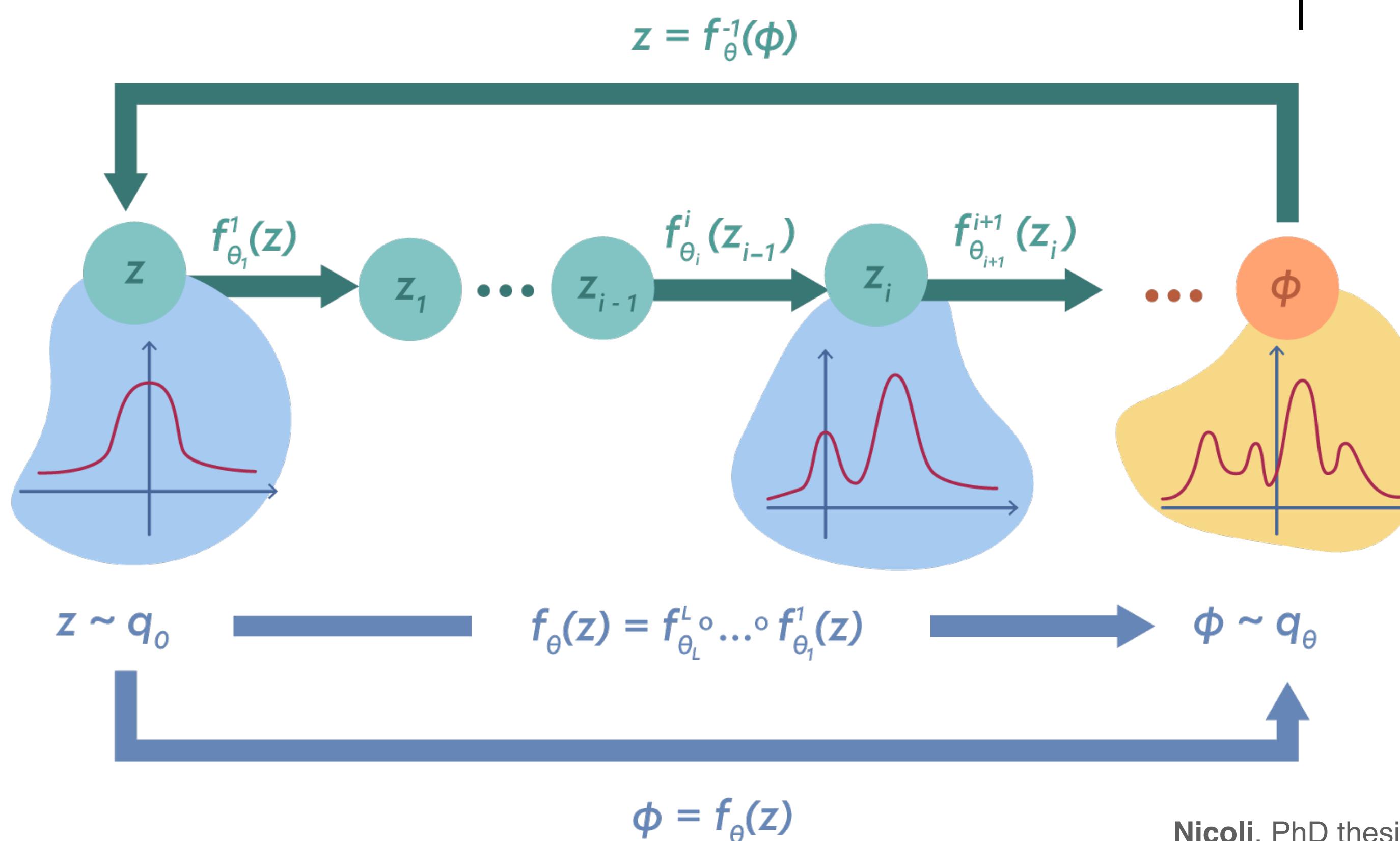
N.B. The function  $f_\theta$  is parametrized by a Neural Network.

Nicoli, PhD thesis (2023)

# Sampling with Normalizing Flows

The likelihood of  $q_\theta$  can be computed exactly:

$$q_\theta(\phi) = q_0(f_\theta^{-1}(\phi)) \left| \det \left( \frac{\partial f_\theta}{\partial z} \right) \right|^{-1}$$



# How do we train a generative model?

---

Often the variational density  $q_\theta$  is trained by minimizing the **Reverse-KL** divergence:

$$KL(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)} \equiv \mathbb{E}_{q_\theta} \left[ \ln \frac{q_\theta(\phi)}{p(\phi)} \right]$$

# How do we train a generative model?

---

Often the variational density  $q_\theta$  is trained by minimizing the **Reverse-KL** divergence:

$$KL(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)} \equiv \mathbb{E}_{q_\theta} \left[ \ln \frac{q_\theta(\phi)}{p(\phi)} \right]$$

since we know the target  $p(\phi)$  is a Boltzmann distribution  $p(\phi) = Z^{-1} \exp\{-S(\phi)\}$

$$KL(q_\theta || p) = \mathbb{E}_{q_\theta} \left[ \ln \frac{q_\theta(\phi)}{p(\phi)} \right] = \mathbb{E}_{q_\theta} [\ln q_\theta(\phi) + S(\phi) + \ln Z]$$

# How do we train a generative model?

---

Often the variational density  $q_\theta$  is trained by minimizing the **Reverse-KL** divergence:

$$KL(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)} \equiv \mathbb{E}_{q_\theta} \left[ \ln \frac{q_\theta(\phi)}{p(\phi)} \right]$$

since we know the target  $p(\phi)$  is a Boltzmann distribution  $p(\phi) = Z^{-1} \exp\{-S(\phi)\}$

$$\nabla_\theta KL(q_\theta || p) = \mathbb{E}_{q_\theta} \left[ \nabla_\theta \ln q_\theta(\phi) + \nabla_\theta S(\phi) + \cancel{\ln Z} \right]$$

# How do we train a generative model?

---

Often the variational density  $q_\theta$  is trained by minimizing the **Reverse-KL** divergence:

$$KL(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)} \equiv \mathbb{E}_{q_\theta} \left[ \ln \frac{q_\theta(\phi)}{p(\phi)} \right]$$

since we know the target  $p(\phi)$  is a Boltzmann distribution  $p(\phi) = Z^{-1} \exp\{-S(\phi)\}$

$$\nabla_\theta KL(q_\theta || p) = \mathbb{E}_{q_\theta} \left[ \nabla_\theta \ln q_\theta(\phi) + \nabla_\theta S(\phi) + \cancel{\ln Z} \right]$$

Training can be performed by **self-sampling** from the model we are training!

# Shall we Trust These Samplers?

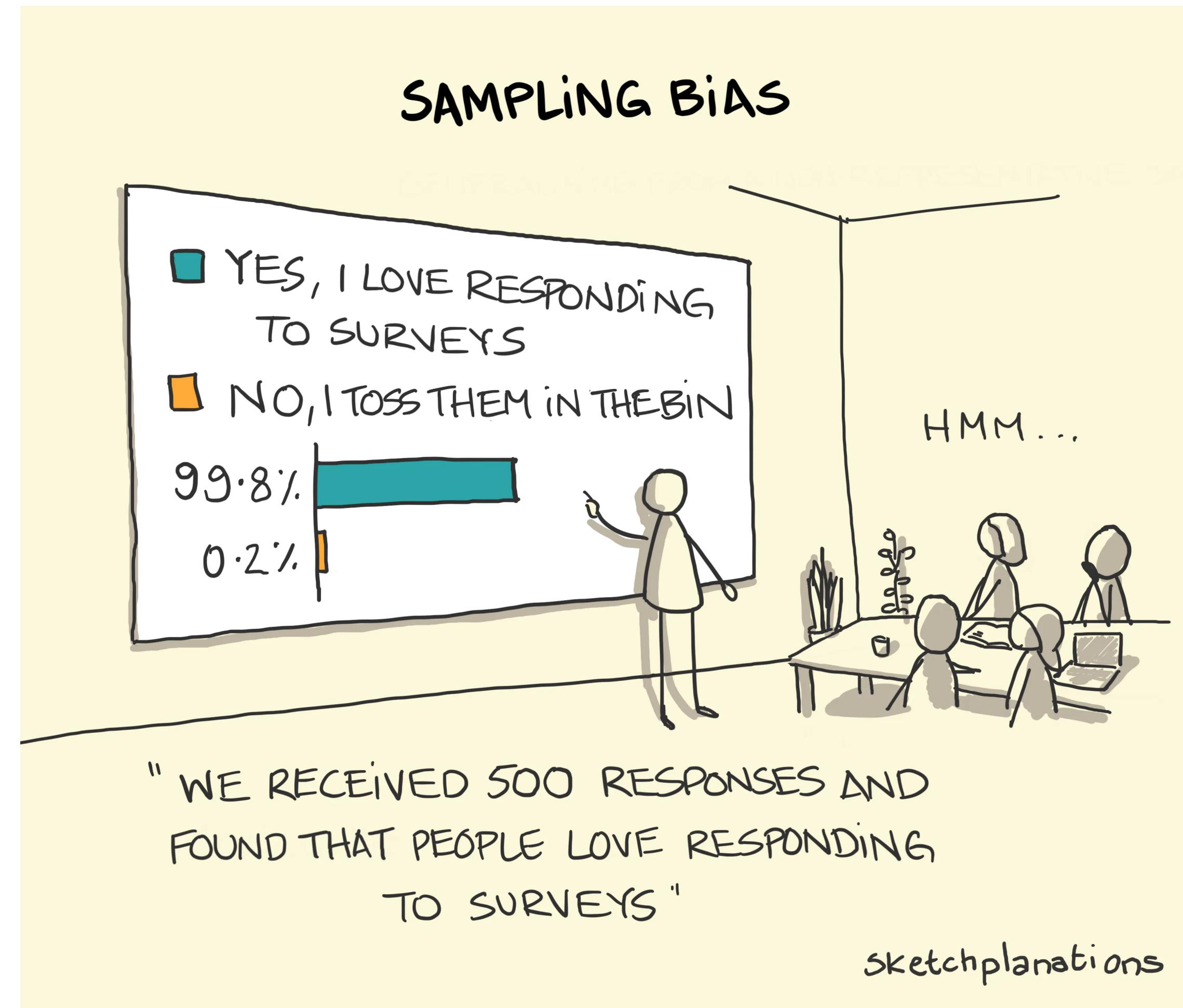
---



(Image credits: [Gymclimber](#))

# YES! But...

---



(Image credits: [sketchplanations](#))

# Neural Importance Sampling (NIS)

---

$$p(\phi) \approx q_\theta \sim \phi_i$$

where

$$p(\phi) = \frac{\exp\{-S(\phi)\}}{Z}$$

$$Z = \int D[\phi] \exp\{-S(\phi)\} = \int D[\phi] q_\theta(\phi) \tilde{w}(\phi) \quad \text{where} \quad \tilde{w}(\phi) = \frac{\exp\{-S(\phi)\}}{q_\theta(\phi)}$$

# Neural Importance Sampling (NIS)

---

$$p(\phi) \approx q_\theta \sim \phi_i$$

where

$$p(\phi) = \frac{\exp\{-S(\phi)\}}{Z}$$

$$Z = \int D[\phi] \exp\{-S(\phi)\} = \int D[\phi] q_\theta(\phi) \tilde{w}(\phi) \quad \text{where} \quad \tilde{w}(\phi) = \frac{\exp\{-S(\phi)\}}{q_\theta(\phi)}$$

# Neural Importance Sampling (NIS)

---

$$p(\phi) \approx q_\theta \sim \phi_i$$

where

$$p(\phi) = \frac{\exp\{-S(\phi)\}}{Z}$$

$$Z = \int D[\phi] \exp\{-S(\phi)\} = \int D[\phi] q_\theta(\phi) \tilde{w}(\phi) \quad \text{where} \quad \tilde{w}(\phi) = \frac{\exp\{-S(\phi)\}}{q_\theta(\phi)}$$

$$Z \stackrel{\text{MC}}{\approx} \hat{Z} = \frac{1}{N} \sum_{i=1}^N \tilde{w}(\phi_i) \quad \phi_i \sim q_\theta$$

# Asymptotically Unbiased Estimation of Physical Observables

---

$$\langle \mathcal{O} \rangle_p = \langle {}^w \mathcal{O} \rangle_{q_\theta} \stackrel{\text{MC}}{\approx} \frac{1}{N} \sum_{i=1}^N {}^w \mathcal{O}(\phi_i) \quad \phi_i \sim q_\theta$$

# Asymptotically Unbiased Estimation of Physical Observables

---

$$w(\phi) = \frac{p(\phi)}{q_\theta(\phi)}$$
$$\uparrow$$
$$\langle \mathcal{O} \rangle_p = \langle {}^w \mathcal{O} \rangle_{q_\theta} \stackrel{\text{MC}}{\approx} \frac{1}{N} \sum_{i=1}^N w(\phi_i) \mathcal{O}(\phi_i) \quad \phi_i \sim q_\theta$$

# Asymptotically Unbiased Estimation of Physical Observables

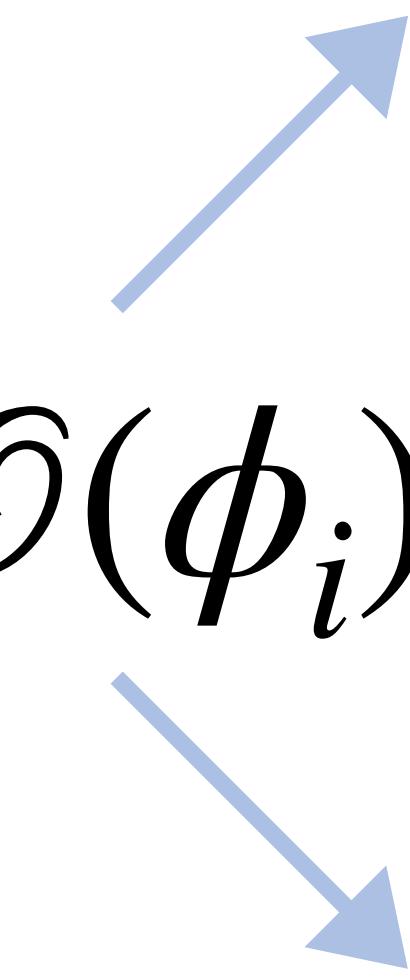
---

$$\langle \mathcal{O} \rangle_p = \langle {}^w \mathcal{O} \rangle_{q_\theta} \stackrel{\text{MC}}{\approx} \frac{1}{N} \sum_{i=1}^N {}^w \mathcal{O}(\phi_i) \quad \phi_i \sim q_\theta$$

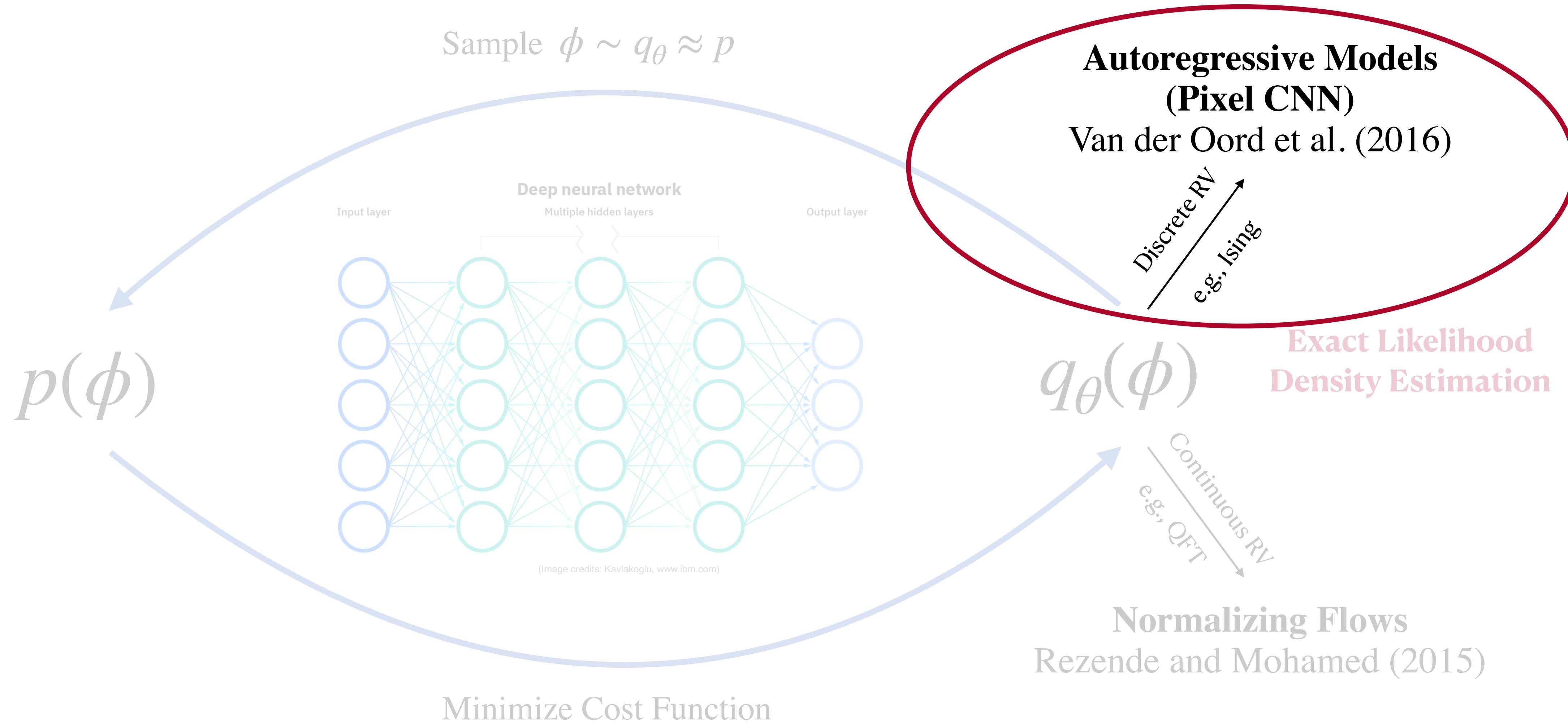
# Direct Estimation of Thermodynamic Quantities

Nicoli, Anders, Funcke, et al., Phys. Rev. Lett. (2021)

$$\hat{F} = -T \ln \hat{Z} \quad !$$

$$\langle \mathcal{O} \rangle_p = \langle {}^w \mathcal{O} \rangle_{q_\theta} \stackrel{\text{MC}}{\approx} \frac{1}{N} \sum_{i=1}^N w(\phi_i) \mathcal{O}(\phi_i) \quad \phi_i \sim q_\theta$$


# The 2D Ising Model

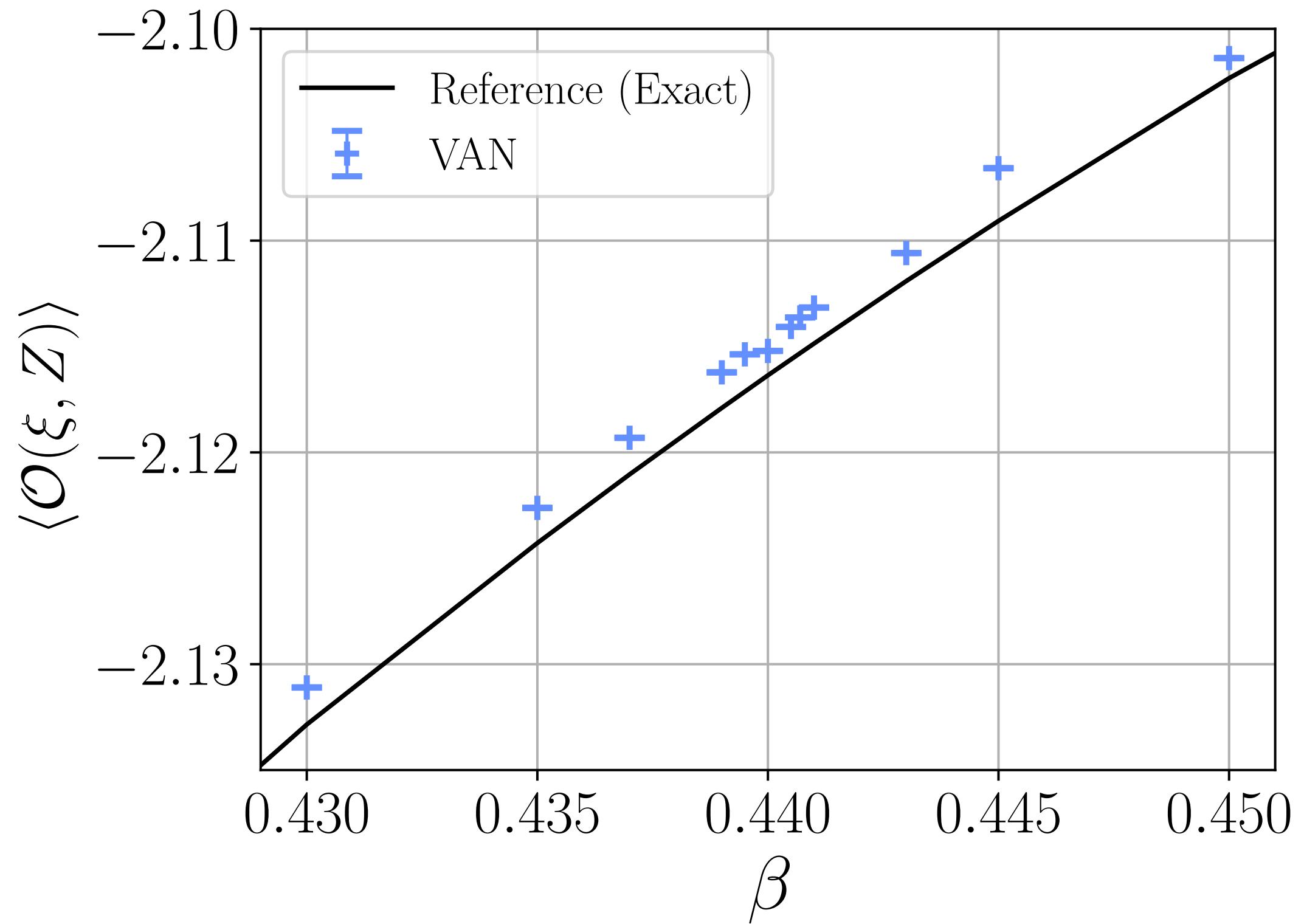


# The 2D Ising Model

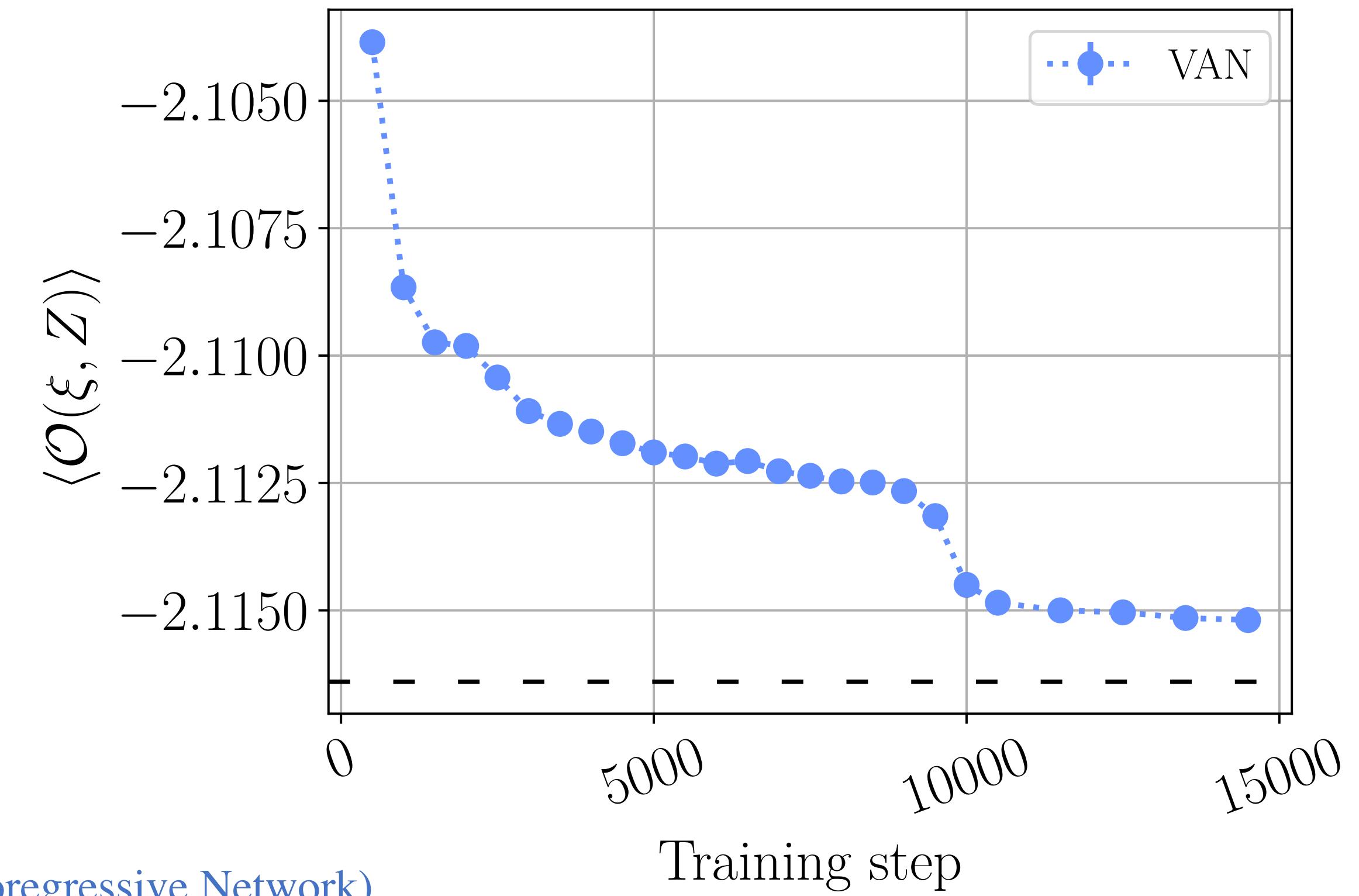
$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j$$

$$p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}, \quad \beta \propto \frac{1}{T}$$

$$\langle \mathcal{O} \rangle_{q_\theta} \neq \langle \mathcal{O} \rangle_p$$

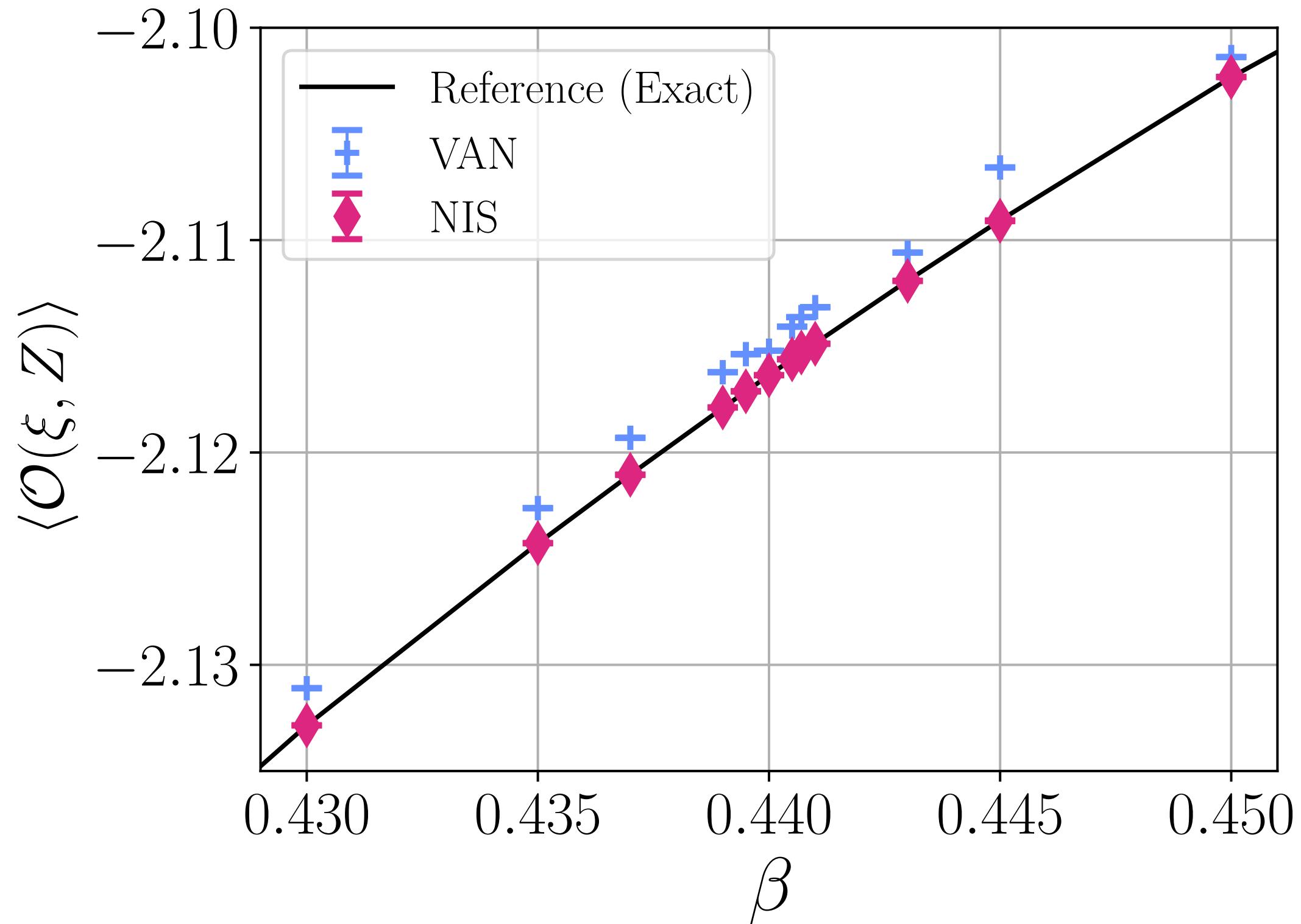


Wu, Wang, Zhang, Phys. Rev. Lett. (2019)  $\implies$  VAN (Variational Autoregressive Network)



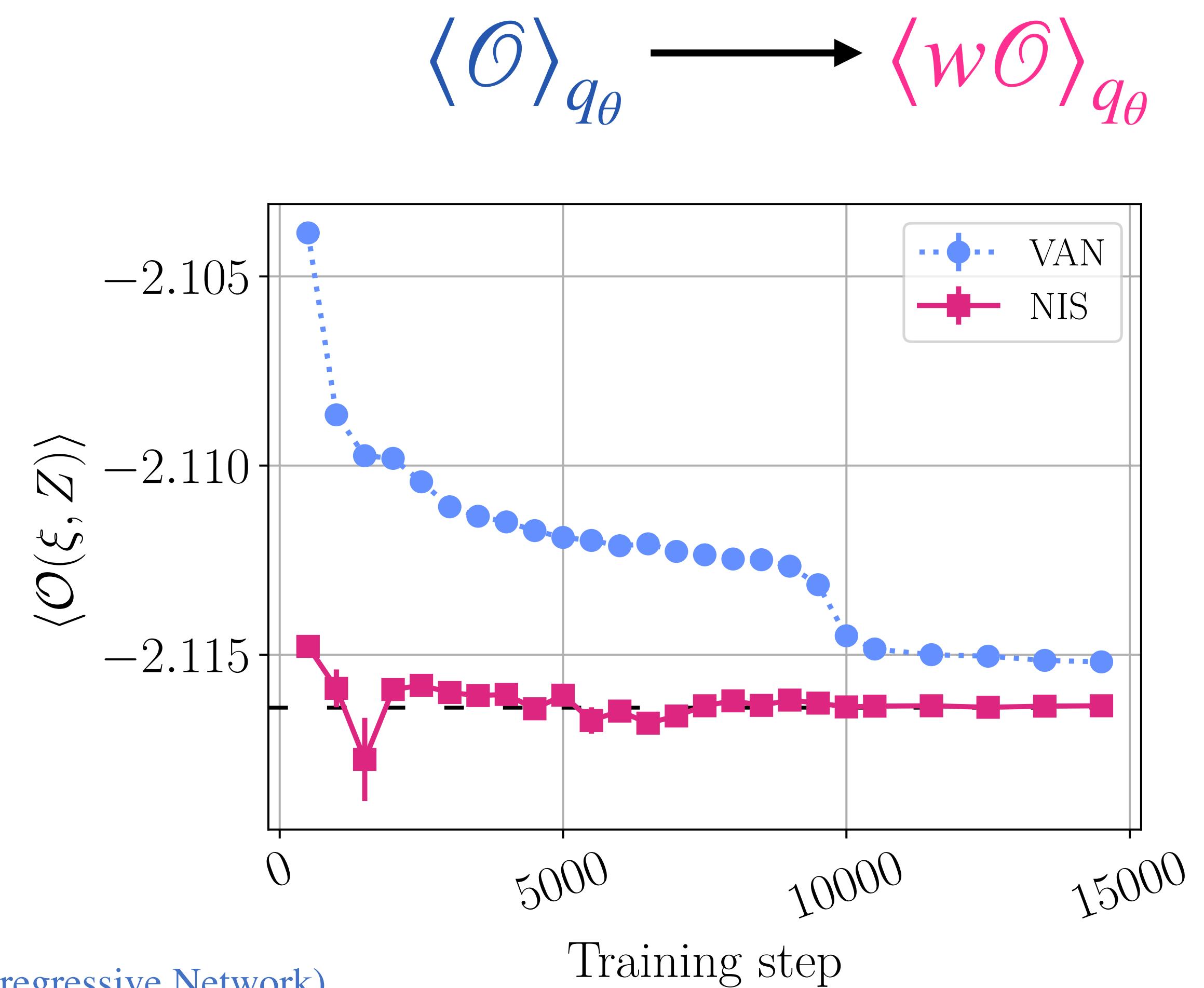
# The 2D Ising Model

$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j \quad p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}$$

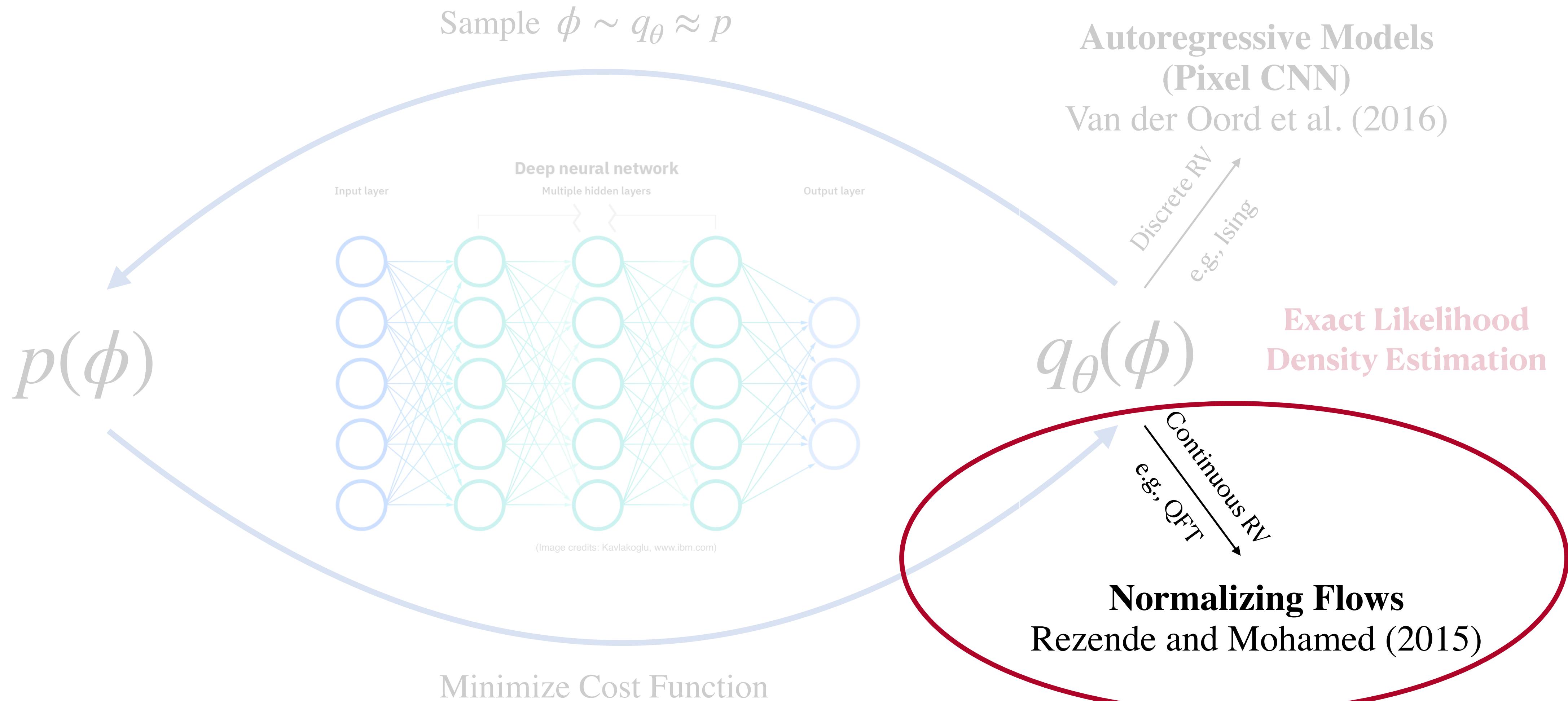


Wu, Wang, Zhang, Phys. Rev. Lett. (2019)  $\implies$  VAN (Variational Autoregressive Network)

Nicoli, Nakajima, Strodthoff, Samek, Müller, and Kessel, Phys. Rev. E (2020)  $\implies$  NIS (Neural Importance Sampling)



# Application to lattice quantum field theory



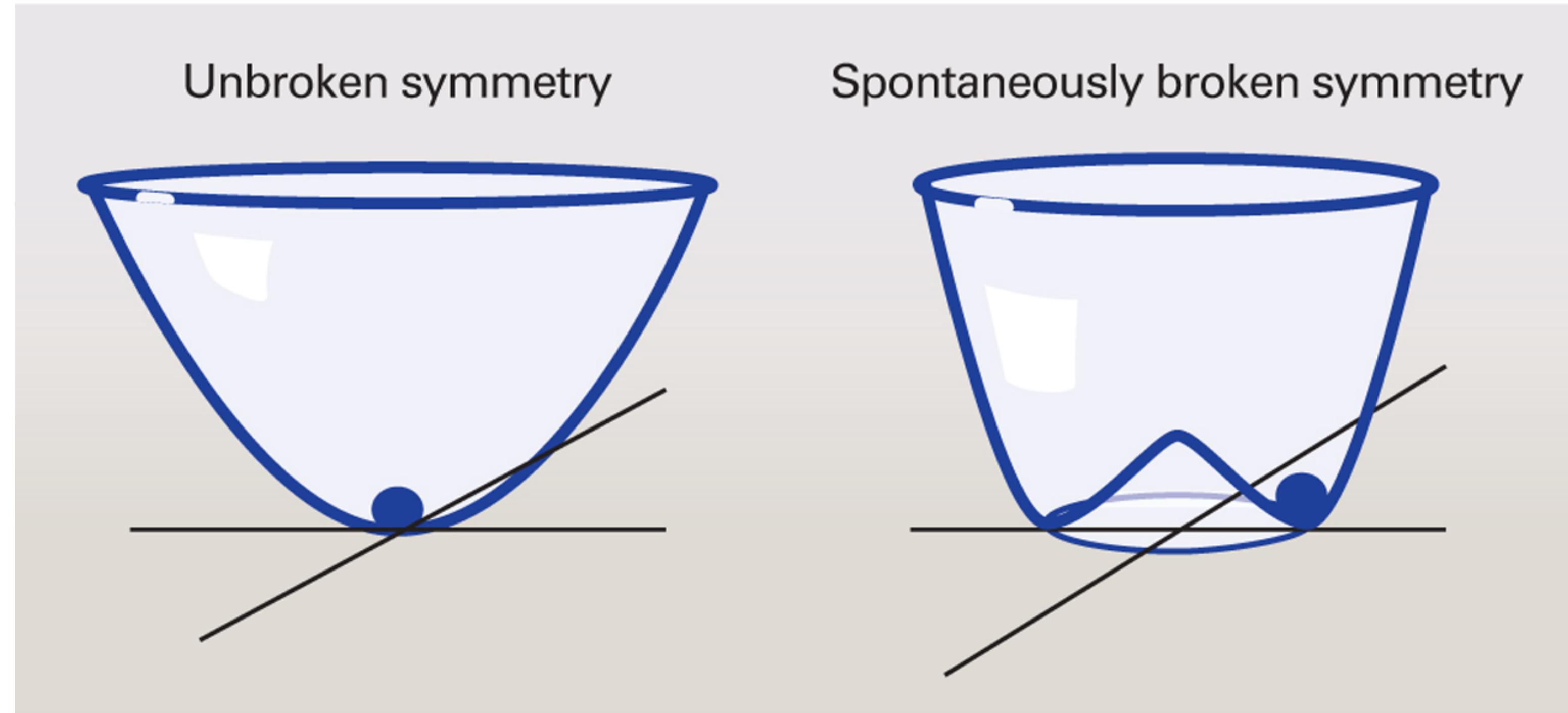
# Real Scalar $\phi^4$ -Theory in (1+1) D

---

$$S(\phi) = \sum_{x \in \Lambda} \left\{ -2\kappa \sum_{\mu=1}^d \phi(x)\phi(x + \hat{\mu}) + (1 - 2\lambda)\phi(x)^2 + \lambda\phi(x)^4 \right\}$$

# Real Scalar $\phi^4$ -Theory in (1+1) D

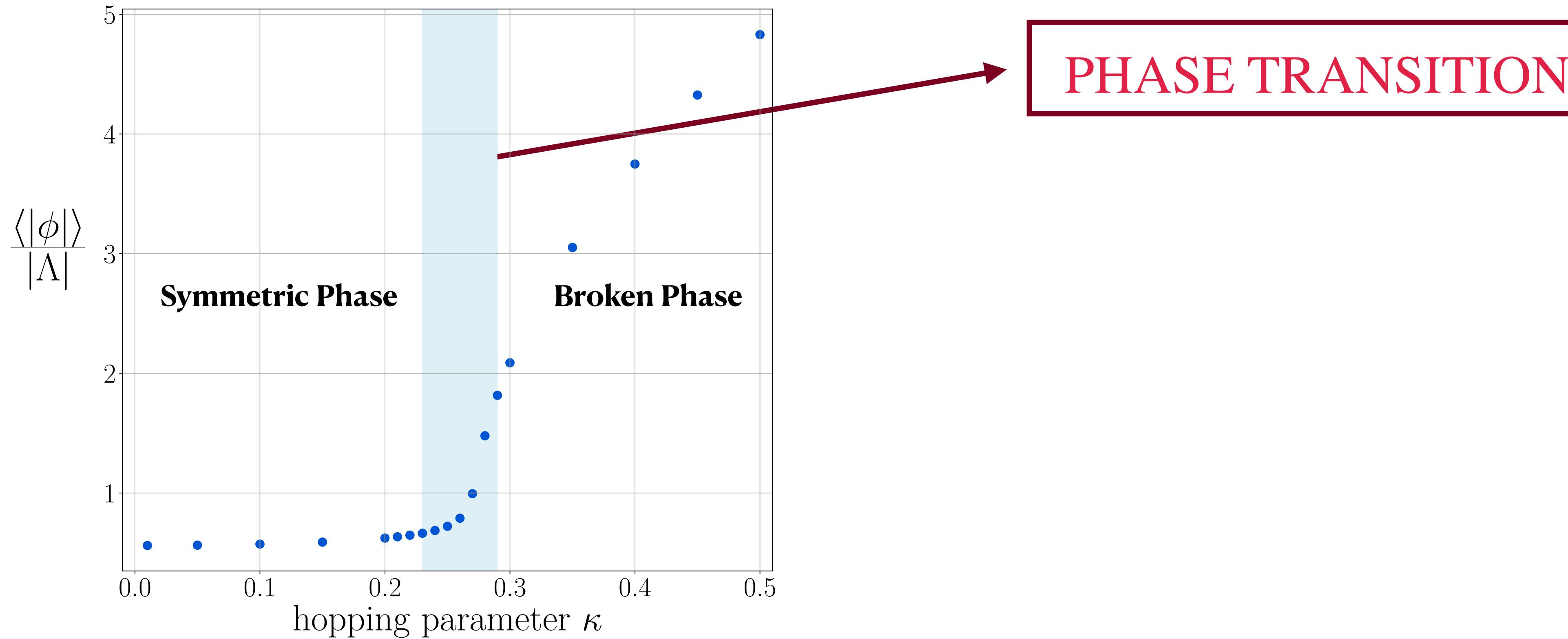
$$S(\phi) = \sum_{x \in \Lambda} \left\{ -2\kappa \sum_{\mu=1}^d \phi(x)\phi(x + \hat{\mu}) + (1 - 2\lambda)\phi(x)^2 + \lambda\phi(x)^4 \right\}$$



(Image credits: [Science](#))

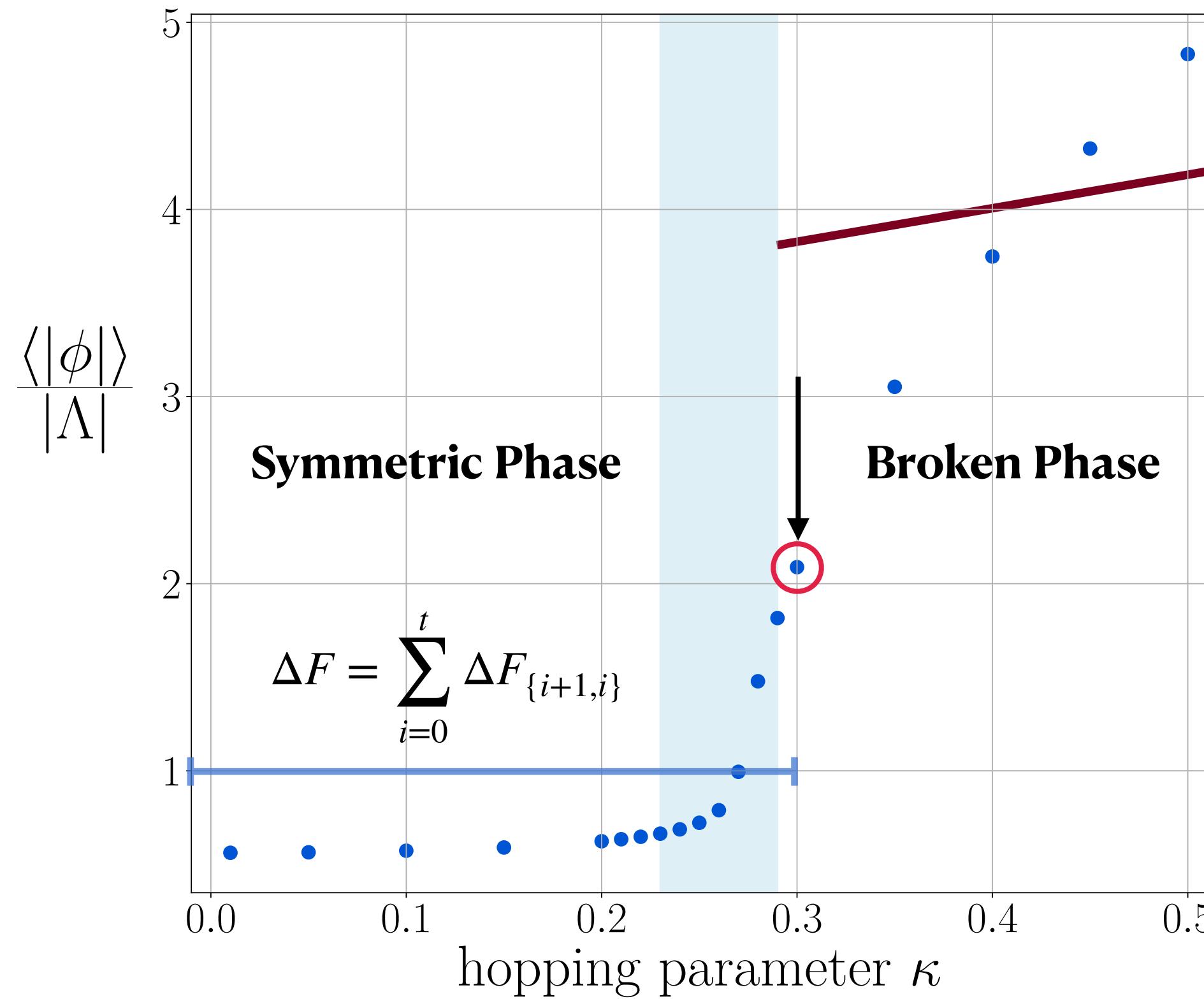
# Spontaneous Breaking of $\mathbb{Z}_2$ Symmetry

$$S(\phi) = \sum_{x \in \Lambda} \left\{ -2\kappa \sum_{\mu=1}^d \phi(x)\phi(x + \hat{\mu}) + (1 - 2\lambda)\phi(x)^2 + \lambda\phi(x)^4 \right\}$$



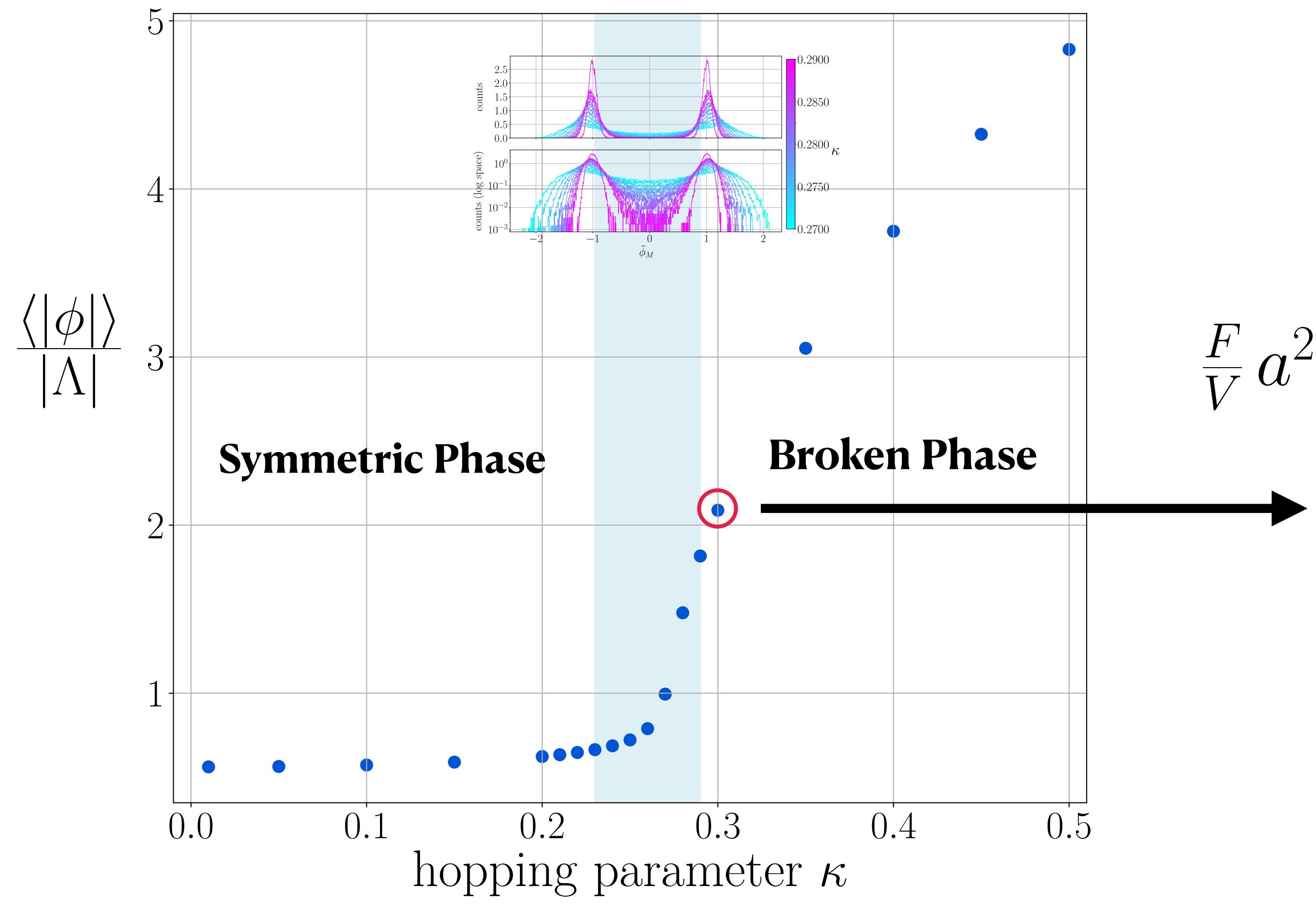
# Spontaneous Breaking of $\mathbb{Z}_2$ Symmetry

$$S(\phi) = \sum_{x \in \Lambda} \left\{ -2\kappa \sum_{\mu=1}^d \phi(x)\phi(x + \hat{\mu}) + (1 - 2\lambda)\phi(x)^2 + \lambda\phi(x)^4 \right\}$$

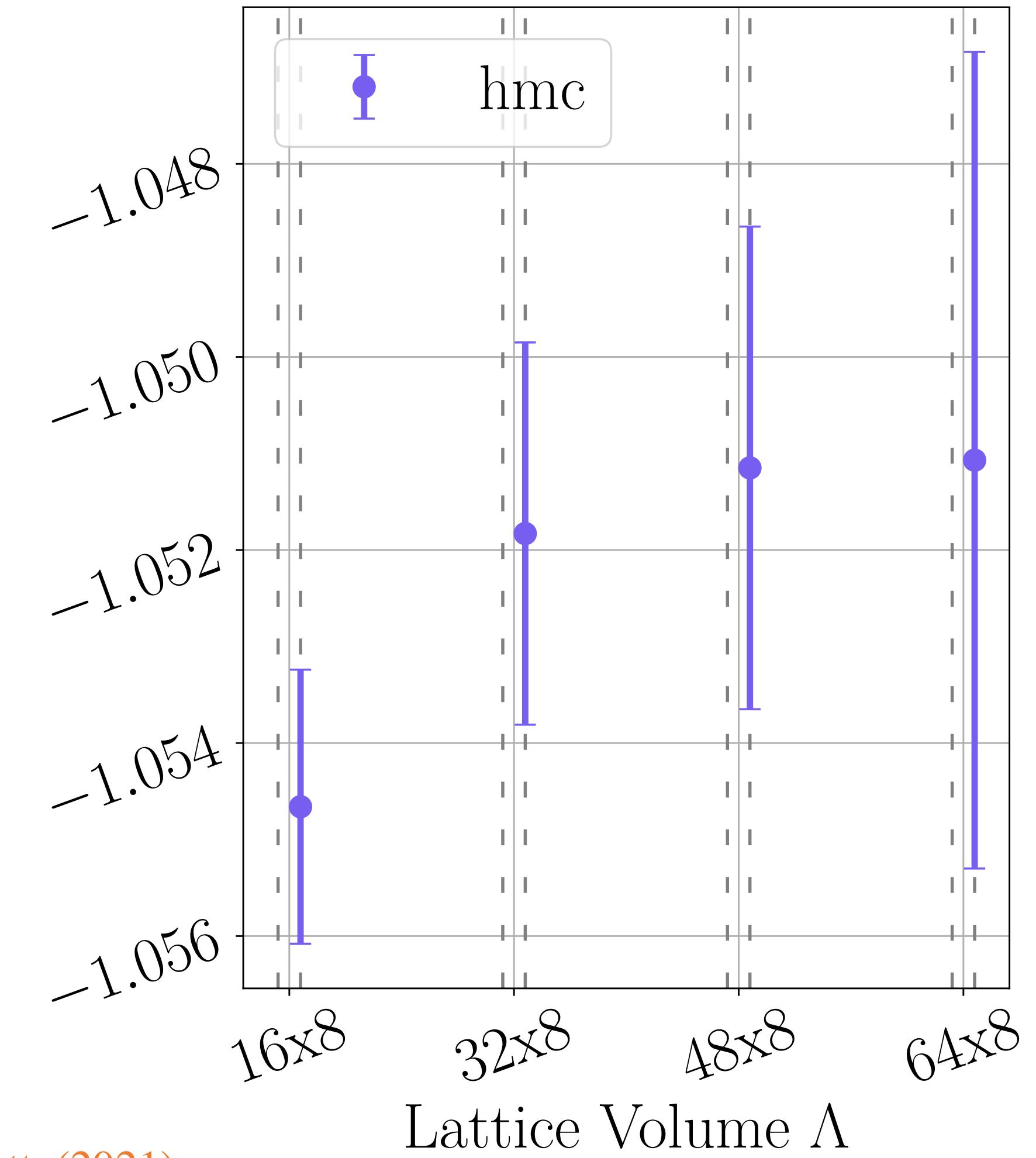


- $\Delta F = F_t - F_0 = -T \ln Z_t/Z_0$
- $F_0$  is needed to compute  $F_{target}$
- Integration through phase space (MCMC)
- Intermediate steps  $\Delta F_{\{i+1,i\}} = -T \ln \frac{Z_{i+1}}{Z_i}$

# Estimation of Thermodynamic Observables

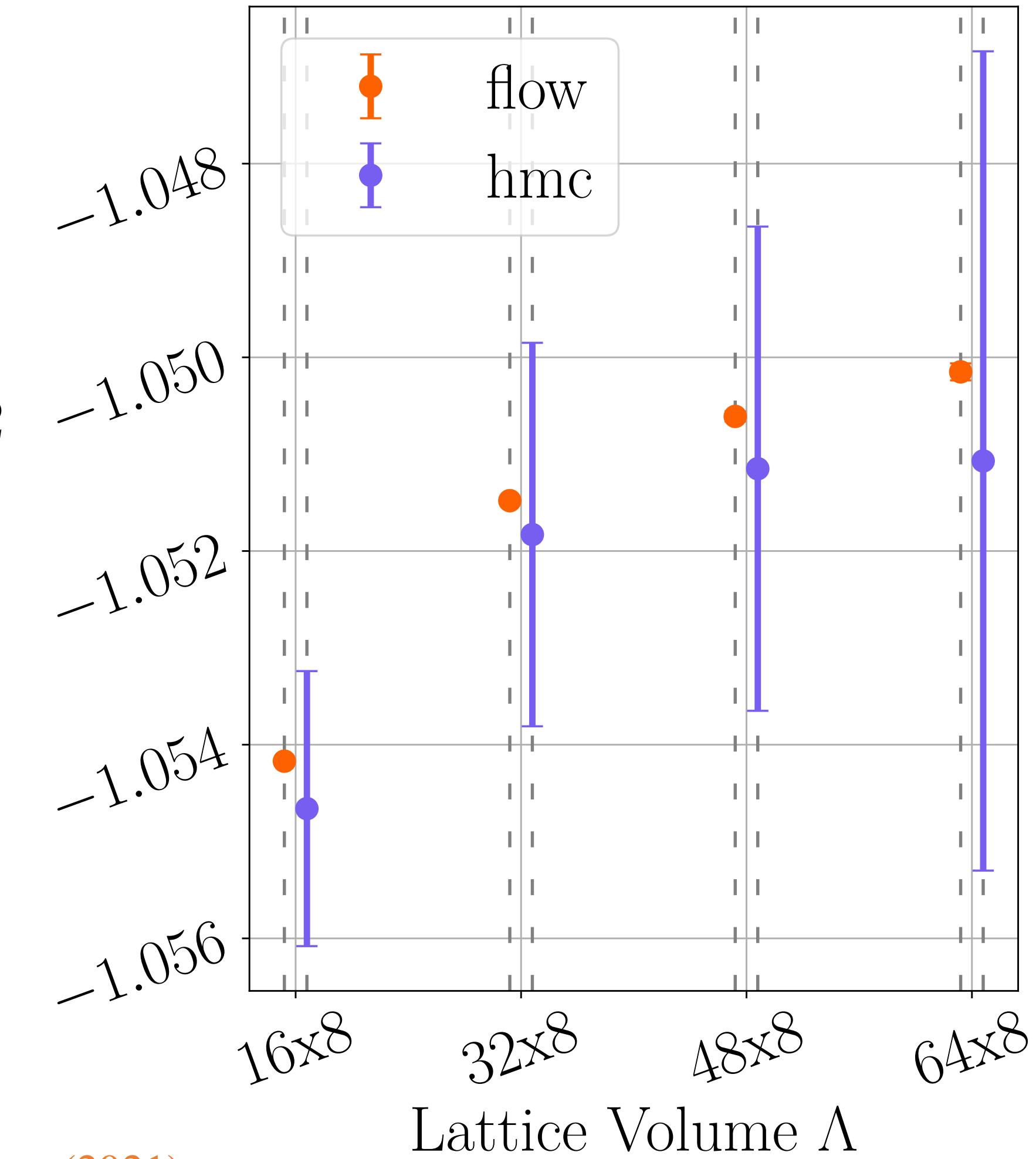
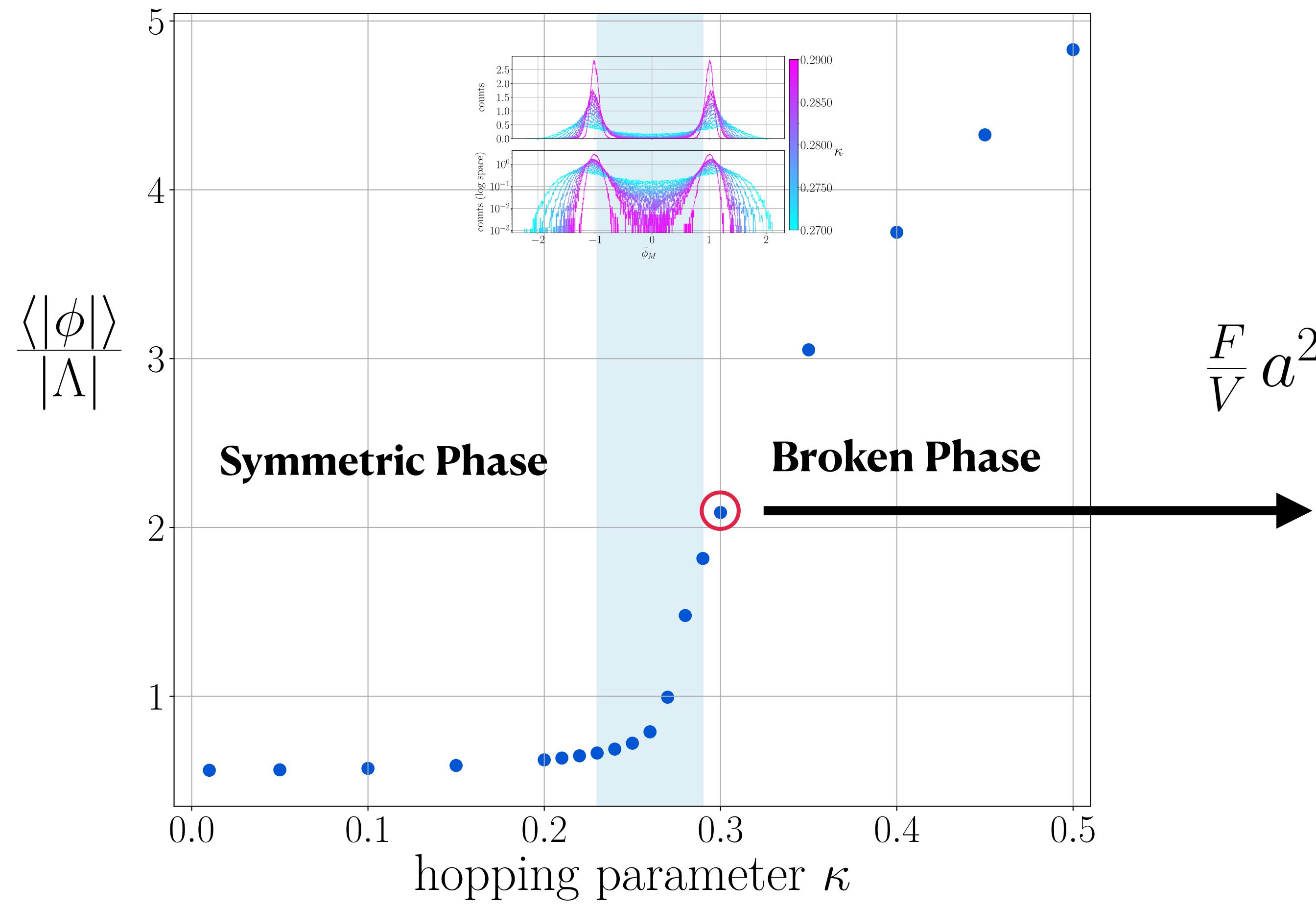


$$\frac{F}{V} a^2$$



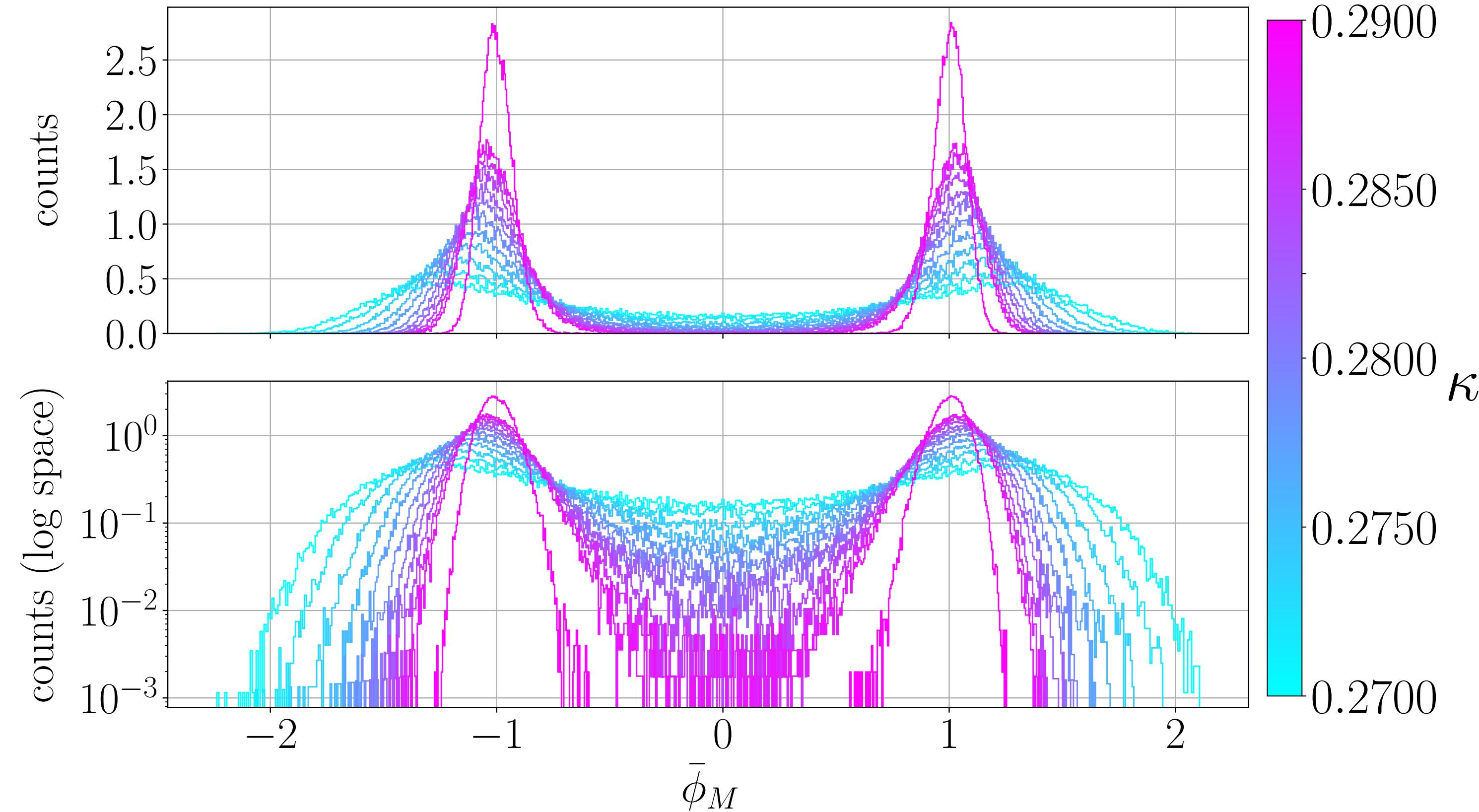
Nicoli, Anders, Funcke, Hartung, Jansen, Kessel, Nakajima, Stornati, Phys. Rev. Lett. (2021)

# Estimation of Thermodynamic Observables



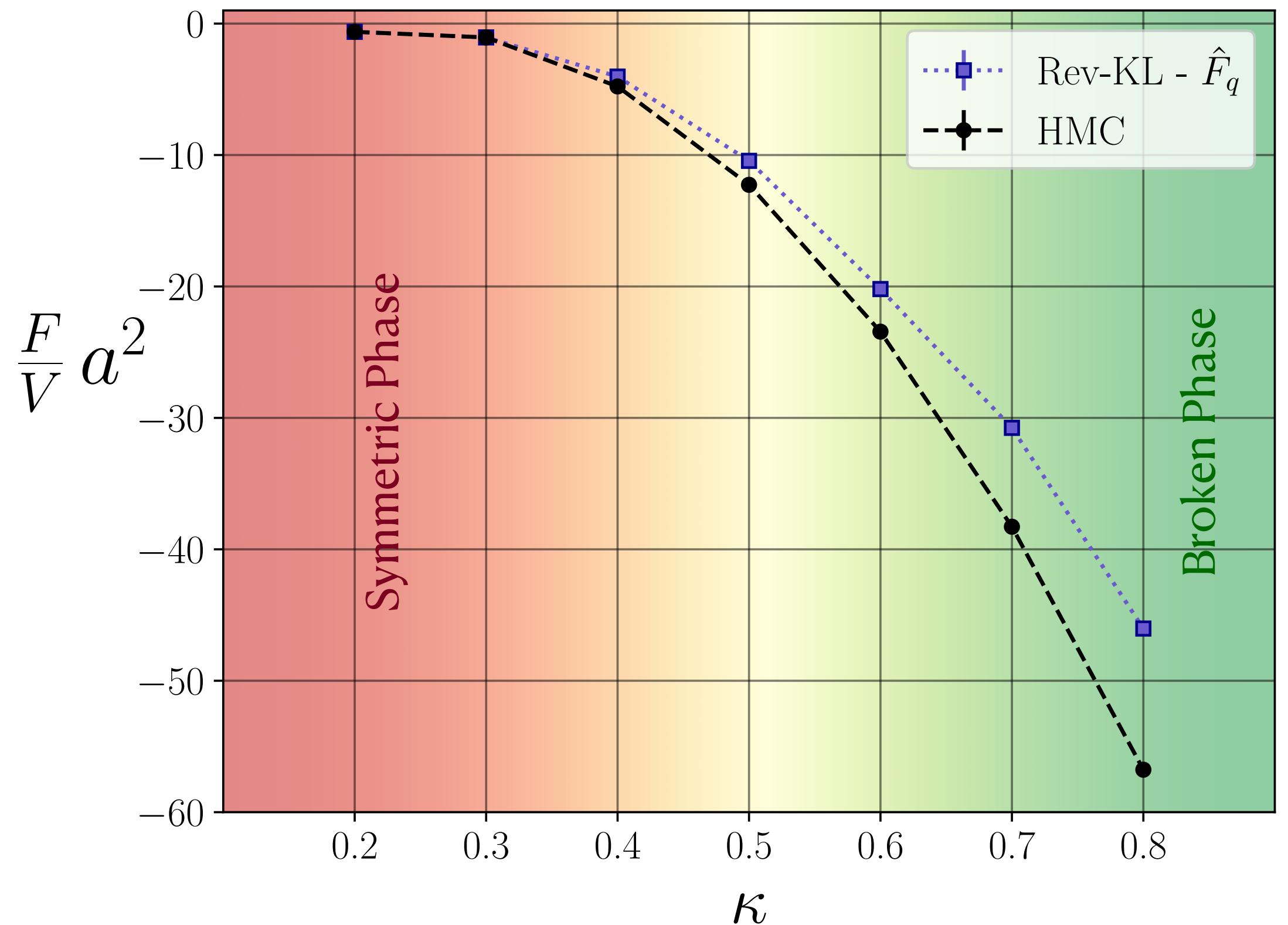
Nicoli, Anders, Funcke, Hartung, Jansen, Kessel, Nakajima, Stornati, Phys. Rev. Lett. (2021)

# Real Scalar $\phi^4$ -Theory in (1+1) D



Nicoli, Anders, Hartung, et al., Phys. Rev. D (2023)

# Real Scalar $\phi^4$ -Theory in (1+1) D



$q$ - estimator for the free energy

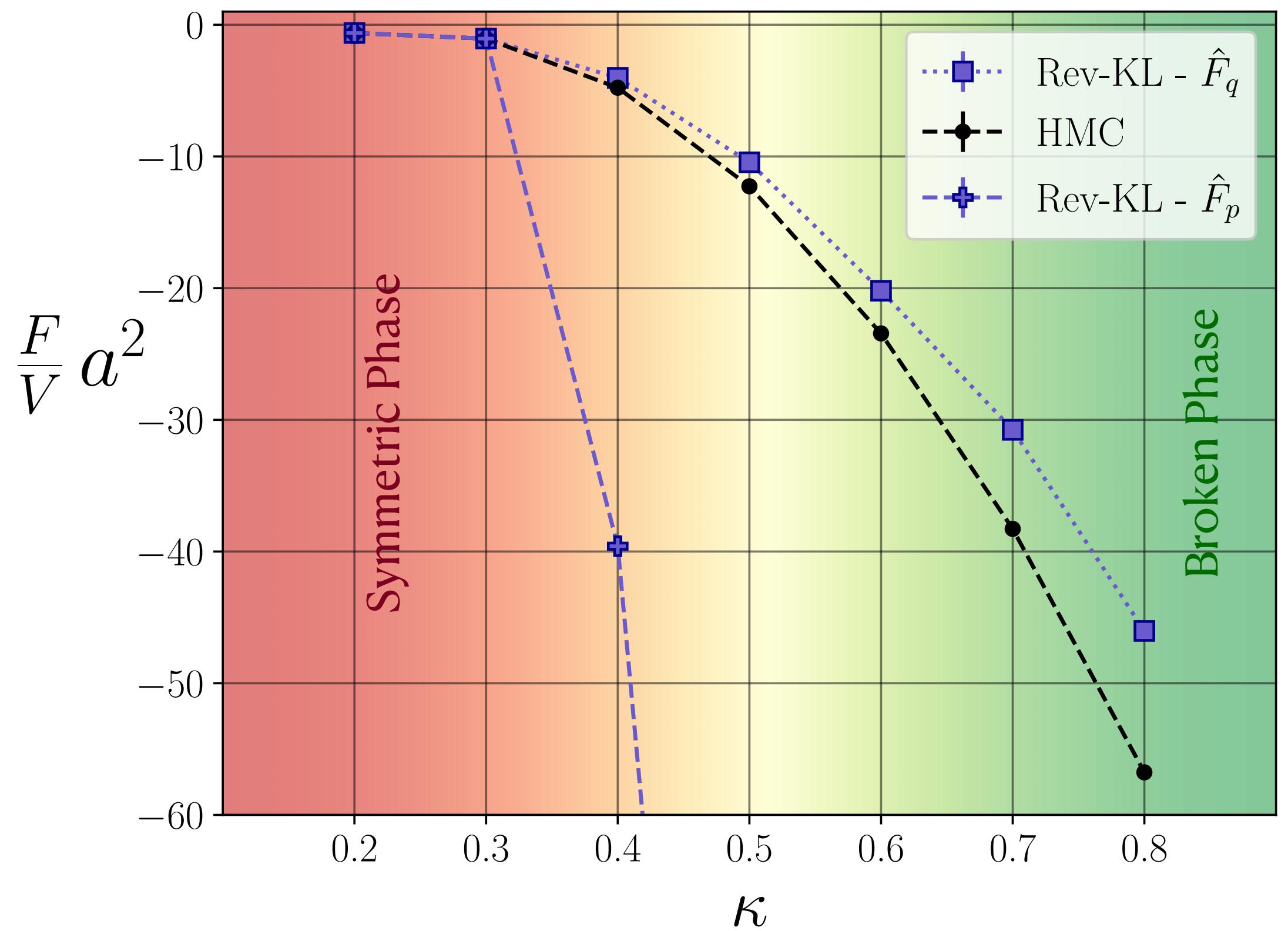
$$Z = \mathbb{E}_{\phi \sim q_\theta} \left[ \frac{e^{-S(\phi)}}{q_\theta(\phi)} \right] \approx \frac{1}{N} \sum_{i=1}^N \frac{e^{-S(\phi_i)}}{q_\theta(\phi_i)} \equiv \hat{Z}_{q_\theta}$$

$$\phi_i \sim q_\theta$$

$$\hat{F}_q = -T \log(\hat{Z}_{q_\theta})$$

Nicoli, Anders, Hartung, et al., Phys. Rev. D (2023)

# Real Scalar $\phi^4$ -Theory in (1+1) D



$q$ - estimator for the free energy

$$Z = \mathbb{E}_{\phi \sim q_\theta} \left[ \frac{e^{-S(\phi)}}{q_\theta(\phi)} \right] \approx \frac{1}{N} \sum_{i=1}^N \frac{e^{-S(\phi_i)}}{q_\theta(\phi_i)} \equiv \hat{Z}_{q_\theta}$$

$$\phi_i \sim q_\theta$$

$$\hat{F}_q = -T \log(\hat{Z}_{q_\theta})$$

$p$ - estimator for the free energy

$$Z^{-1} = \mathbb{E}_{\phi \sim p} \left[ \frac{q_\theta(\phi)}{e^{-S(\phi)}} \right] \approx \frac{1}{N} \sum_{j=1}^N \frac{q_\theta(\phi_j)}{e^{-S(\phi_j)}} \equiv \hat{Z}_p^{-1}$$

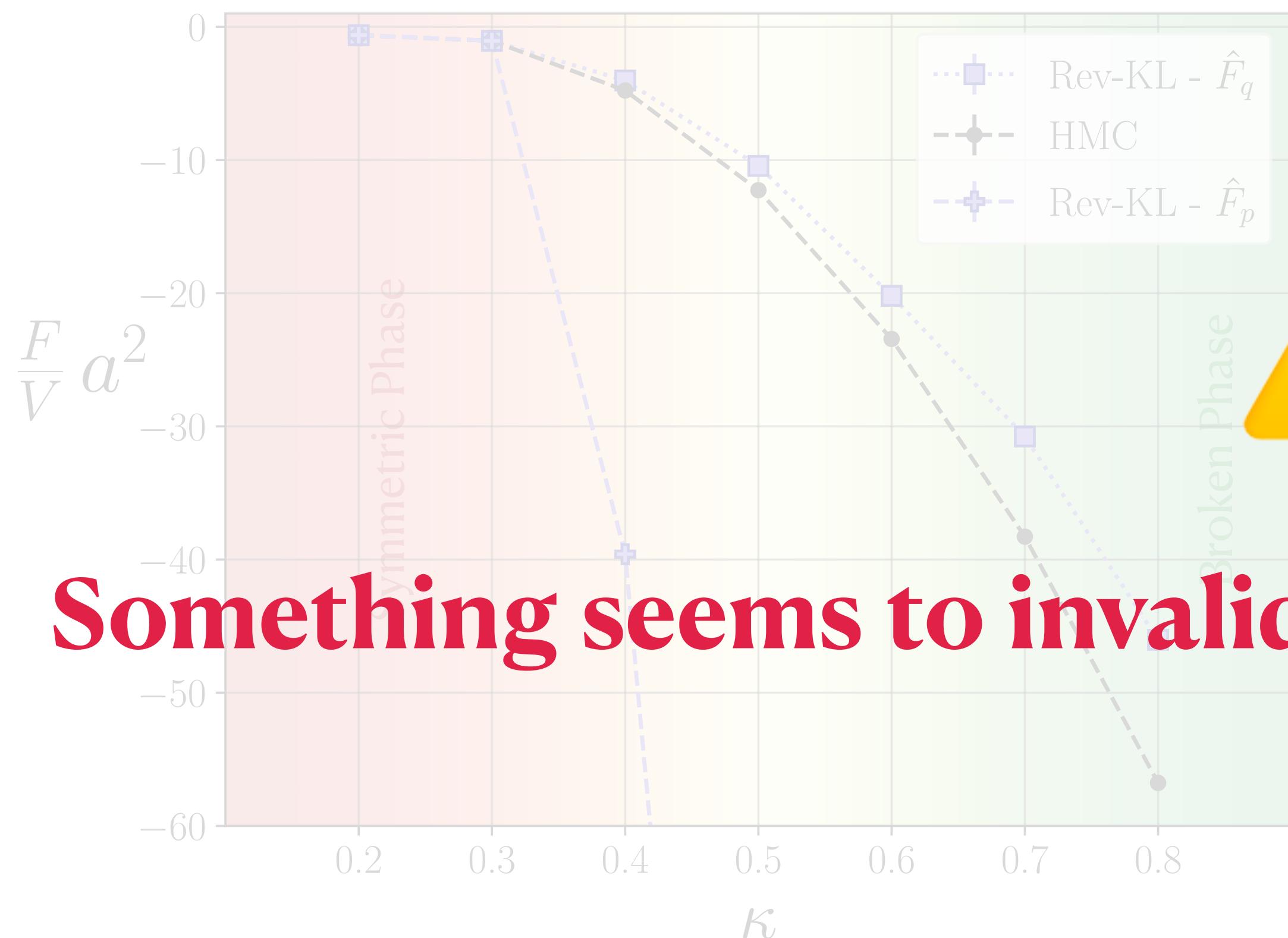
$$\phi_j \sim p$$

$$\hat{F}_p = T \log(\hat{Z}_p^{-1})$$

\*Samples are obtained from HMC away from criticality

Nicoli, Anders, Hartung, et al., Phys. Rev. D (2023)

# Real Scalar $\phi^4$ -Theory in (1+1) D



Nicoli, Anders, Hartung, et al., Phys. Rev. D (2023)

$$Z = \mathbb{E}_{\phi \sim q_\theta} \left[ \frac{e^{-S(\phi)}}{q_\theta(\phi)} \right] \approx \frac{1}{N} \sum_{i=1}^N \frac{e^{-S(\phi_i)}}{q_\theta(\phi_i)} \equiv \hat{Z}_{q_\theta}$$

$$\phi_i \sim q_\theta$$

$$\hat{F}_q = -T \log(\hat{Z}_{q_\theta})$$

$$Z = \mathbb{E}_{\phi \sim q_\theta} \left[ \frac{q_\theta(\phi)}{e^{-S(\phi)}} \right] = \frac{1}{N} \sum_{j=1}^N \frac{q_\theta(\phi_j)}{e^{-S(\phi_j)}} \approx \hat{Z}_p$$

$$\phi_j \sim p$$

$$\hat{F}_p = T \log(\hat{Z}_p^{-1})$$

# What's going wrong then?

---

## Reverse-KL Div.

$$KL_R(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)}$$

# What's going wrong then?

---

## Reverse-KL Div.

$$KL_R(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)}$$

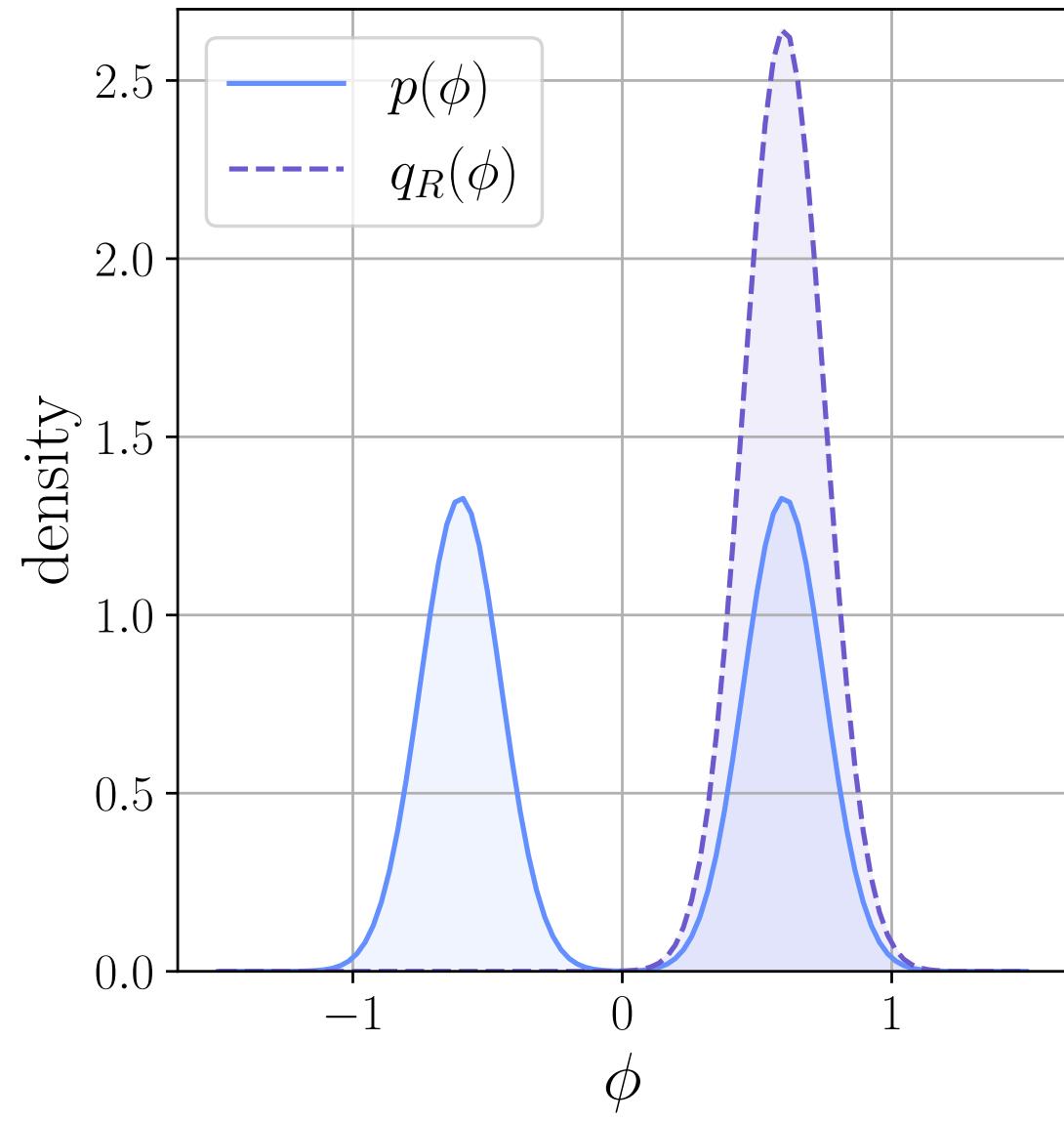
## Forward-KL Div.

$$KL_F(p || q_\theta) = \int D[\phi] p(\phi) \ln \frac{p(\phi)}{q_\theta(\phi)}$$

# What's going wrong then?

## Reverse-KL Div.

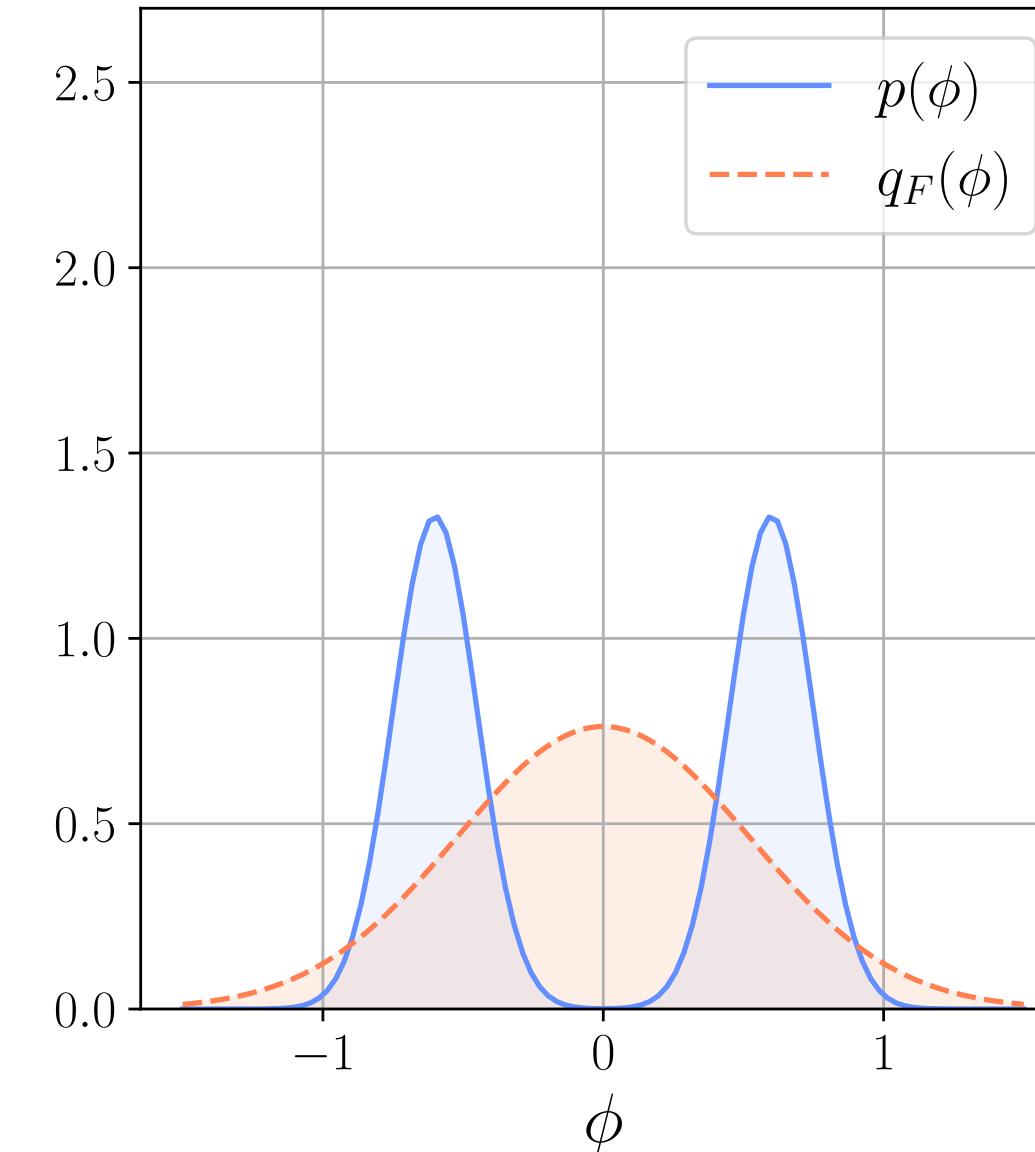
$$KL_R(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)}$$



- Self-Sampling (efficient).
- No need for training data.
- Mode-dropping.
- $q_R$  flow trained with Rev. KL

## Forward-KL Div.

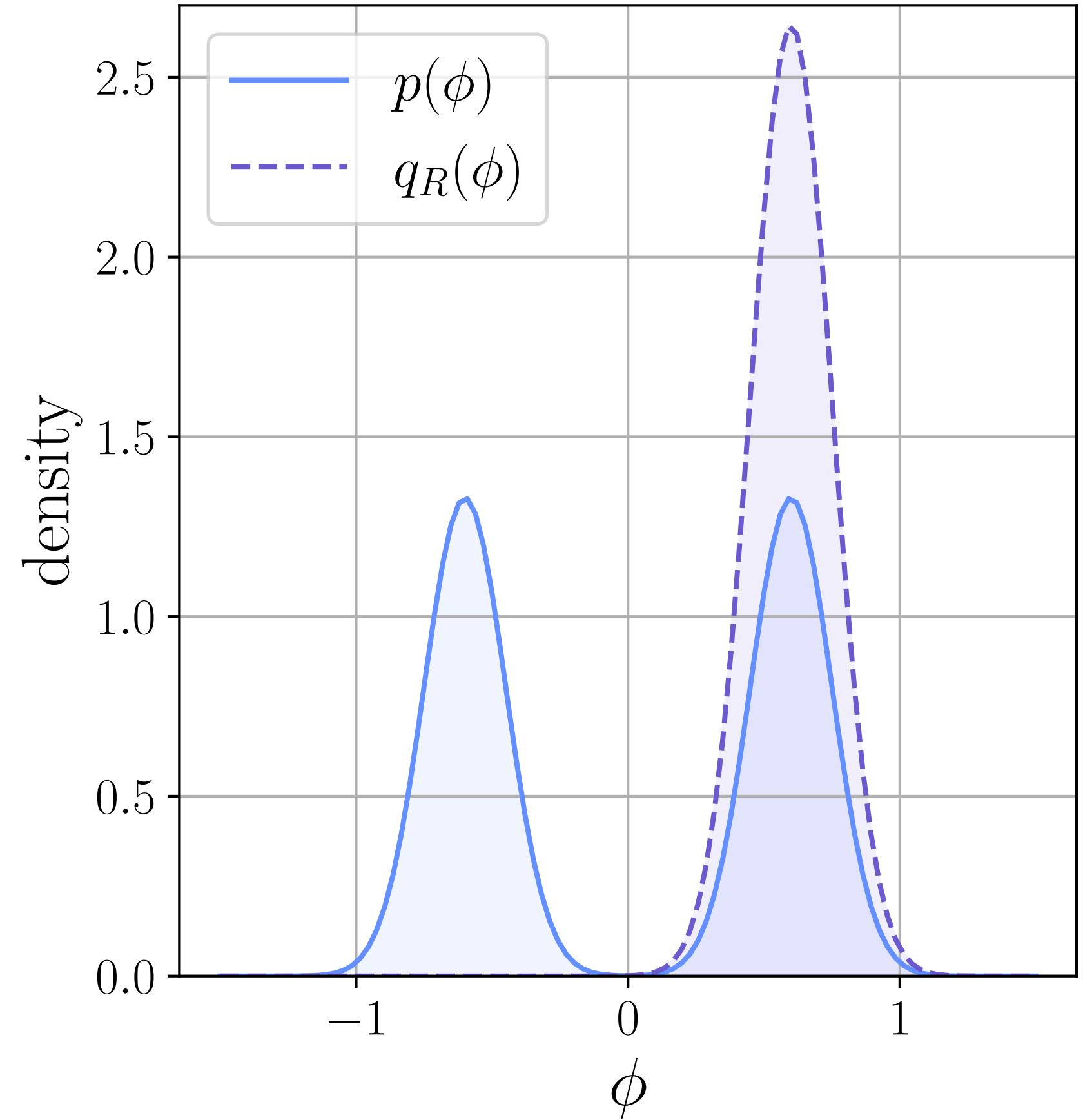
$$KL_F(p || q_\theta) = \int D[\phi] p(\phi) \ln \frac{p(\phi)}{q_\theta(\phi)}$$



- Maximum Likelihood.
- Requires training data.
- Fake modes.
- $q_F$  flow trained with Fwd. KL

# What's going wrong then?

---

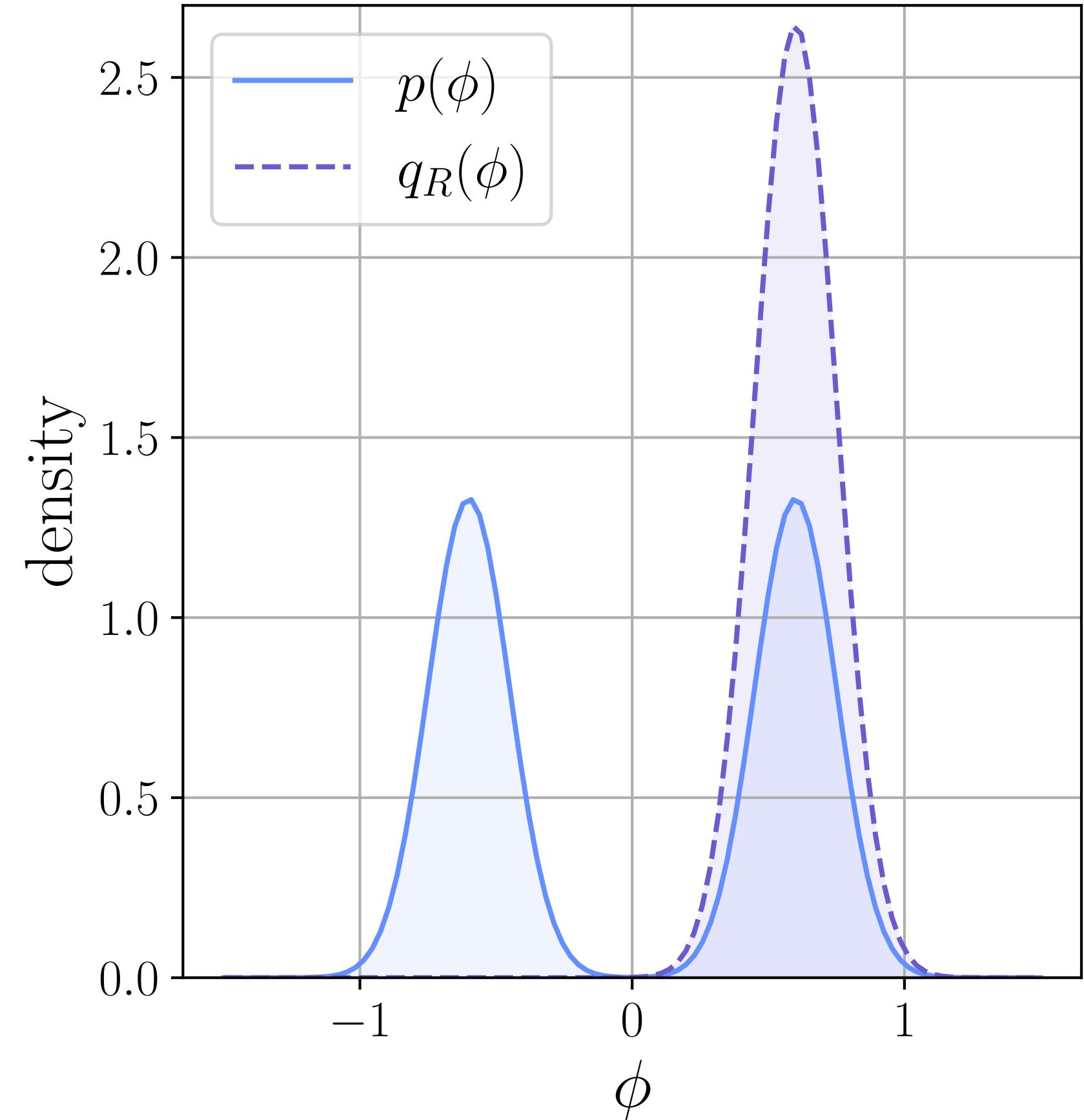


$$\text{ESR} = \frac{\text{ESS}}{N} = \frac{1}{\mathbb{E}_{q_\theta}[w^2]}$$

Where

$$w = \frac{p(\phi)}{q_\theta(\phi)}$$

# What's going wrong then?



The Effective sampling ratio (ESR) is not a good metric!

$$\text{ESR} = \frac{\text{ESS}}{N} = \frac{1}{\mathbb{E}_{q_\theta}[w^2]}$$

Where

$$w = \frac{p(\phi)}{q_\theta(\phi)}$$

The model is **blind** with respect to one (or more) of the modes of the target density.

# The mode-dropping estimator

---

When  $q_\theta$  has **full effective support** on the domain of  $p$

$$w^* = \mathbb{E}_{q_\theta} \left[ \frac{p(\phi)}{q_\theta(\phi)} \right] = \int_{\text{supp}(q_\theta)} q_\theta(\phi) \frac{p(\phi)}{q_\theta(\phi)} \mathcal{D}[\phi] = \int_{\text{supp}(q_\theta)} p(\phi) \mathcal{D}[\phi] = 1.$$

# The mode-dropping estimator

---

When  $q_\theta$  has **full effective support** on the domain of  $p$

$$w^* = \mathbb{E}_{q_\theta} \left[ \frac{p(\phi)}{q_\theta(\phi)} \right] = \int_{\text{supp}(q_\theta)} q_\theta(\phi) \frac{p(\phi)}{q_\theta(\phi)} \mathcal{D}[\phi] = \int_{\text{supp}(q_\theta)} p(\phi) \mathcal{D}[\phi] = 1.$$

however if  $q_\theta$  is **effectively mode-dropping** this expectation value becomes

$$\bar{w} \equiv \frac{1}{Z} \mathbb{E}_{\phi \sim \tilde{q}_\theta} \left[ \frac{e^{-S(\phi)}}{q_\theta(\phi)} \right] \in (0, 1]$$

# The mode-dropping estimator

---

When  $q_\theta$  has **full effective support** on the domain of  $p$

$$w^* = \mathbb{E}_{q_\theta} \left[ \frac{p(\phi)}{q_\theta(\phi)} \right] = \int_{\text{supp}(q_\theta)} q_\theta(\phi) \frac{p(\phi)}{q_\theta(\phi)} \mathcal{D}[\phi] = \int_{\text{supp}(q_\theta)} p(\phi) \mathcal{D}[\phi] = 1.$$

however if  $q_\theta$  is **effectively mode-dropping** this expectation value becomes

$$\bar{w} \equiv \frac{1}{Z} \mathbb{E}_{\phi \sim \tilde{q}_\theta} \left[ \frac{e^{-S(\phi)}}{q_\theta(\phi)} \right] \in (0, 1]$$

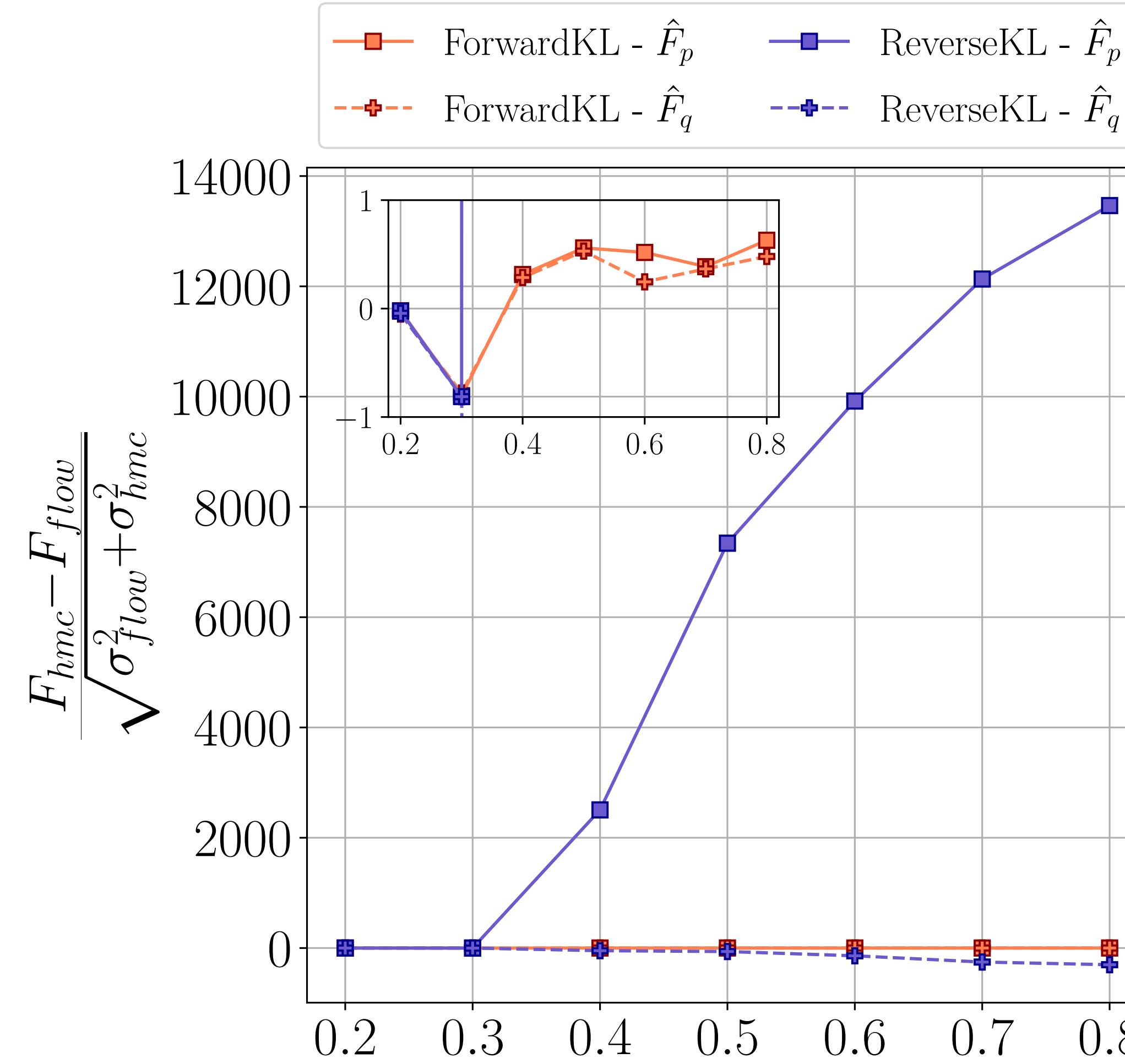
for which we can get the corresponding Monte Carlo estimator, i.e., the **mode-dropping estimator**

$$\bar{w} \approx \frac{1}{\hat{Z}_p} \left( \frac{1}{N} \sum_{i=1}^N \frac{e^{-S(\phi_i)}}{q_\theta(\phi_i)} \right) = \left( \frac{1}{N} \sum_{j=1}^N \frac{q_\theta(\phi_j)}{e^{-S(\phi_j)}} \right) \left( \frac{1}{N} \sum_{i=1}^N \frac{e^{-S(\phi_i)}}{q_\theta(\phi_i)} \right) \equiv \hat{w}$$

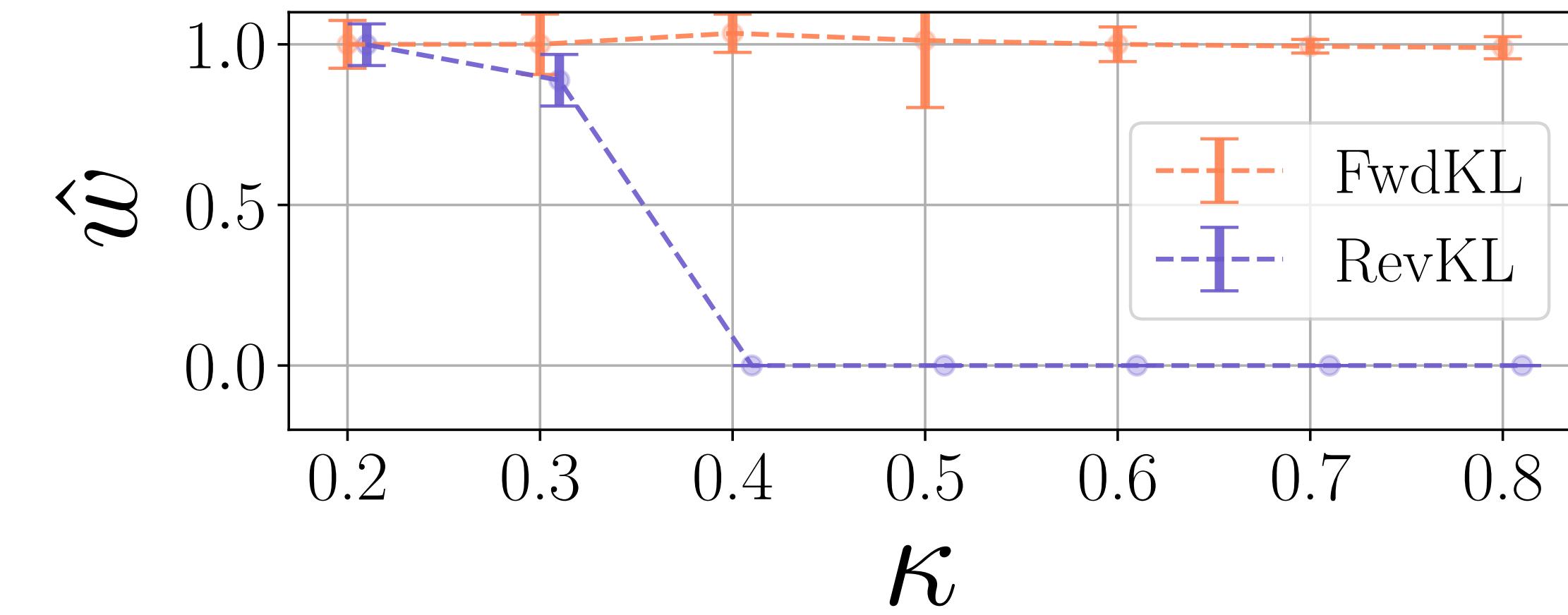
$$\phi_j \sim p$$

$$\phi_i \sim q_\theta$$

# Estimation of Mode Dropping Across Criticality



$\Lambda = 64 \times 8, \lambda = 0.022$



Nicoli, Anders, Hartung, et al., Phys. Rev. D (2023)

# Summary and Conclusions

---

- ▶ **Asymptotically unbiased samplers** can be constructed from trained DGMs (NIS or NMCMC).
- ▶ Direct estimation of the **partition function** and **thermodynamic observables**.
- ▶ Improved DGMs can use **inductive biases**, e.g. symmetries, bootstrapping, annealing, etc.
- ▶ Sampling from DGMs is **embarrassingly parallelizable (i.i.d)**  $\neq$  MCMC (**sequential**).
- ▶ Training with **forward-KL** leads to better models though requires training samples.
- ▶ Derivation of **mode-dropping estimator** to reliably assess the goodness of the model.
- ▶ **Mitigation** of mode-dropping using different **objectives** (FWD-KL) or **stochastic** approaches (SNFs).

# Summary and Conclusions

---

- ▶ Asymptotically unbiased samplers can be constructed from trained DGMs (NIS or NMCMC).

**TL;DR**

- ▶ Direct estimation of the partition function and thermodynamic observables.

- ▶ Improved DGMs can use **inductive biases**, e.g. symmetries, bootstrapping, annealing, etc.

**Deep Generative Models (DGMs) are promising candidates**

- ▶ Sampling from DGMs is embarrassingly parallelizable (i.i.d)  $\neq$  MCMC (sequential).

**for the next generation of sampling algorithms**

- ▶ Training with **forward-KL** leads to better models though requires training samples.

- ▶ Derivation of **mode-dropping estimator** to reliably assess the goodness of the model.

- ▶ **Mitigation** of mode-dropping using different objectives (FWD-KL) or stochastic approaches (SNFs).

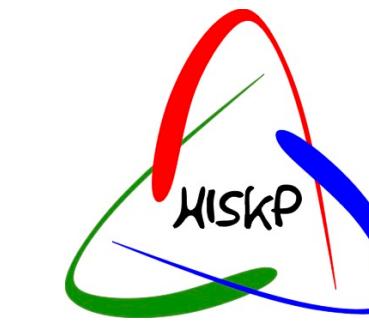
# Thank You!!

---

“It is nice to know that the computer understands the problem. But I would like to understand it too.”  
- E. Wigner



Bonn:



Berlin:

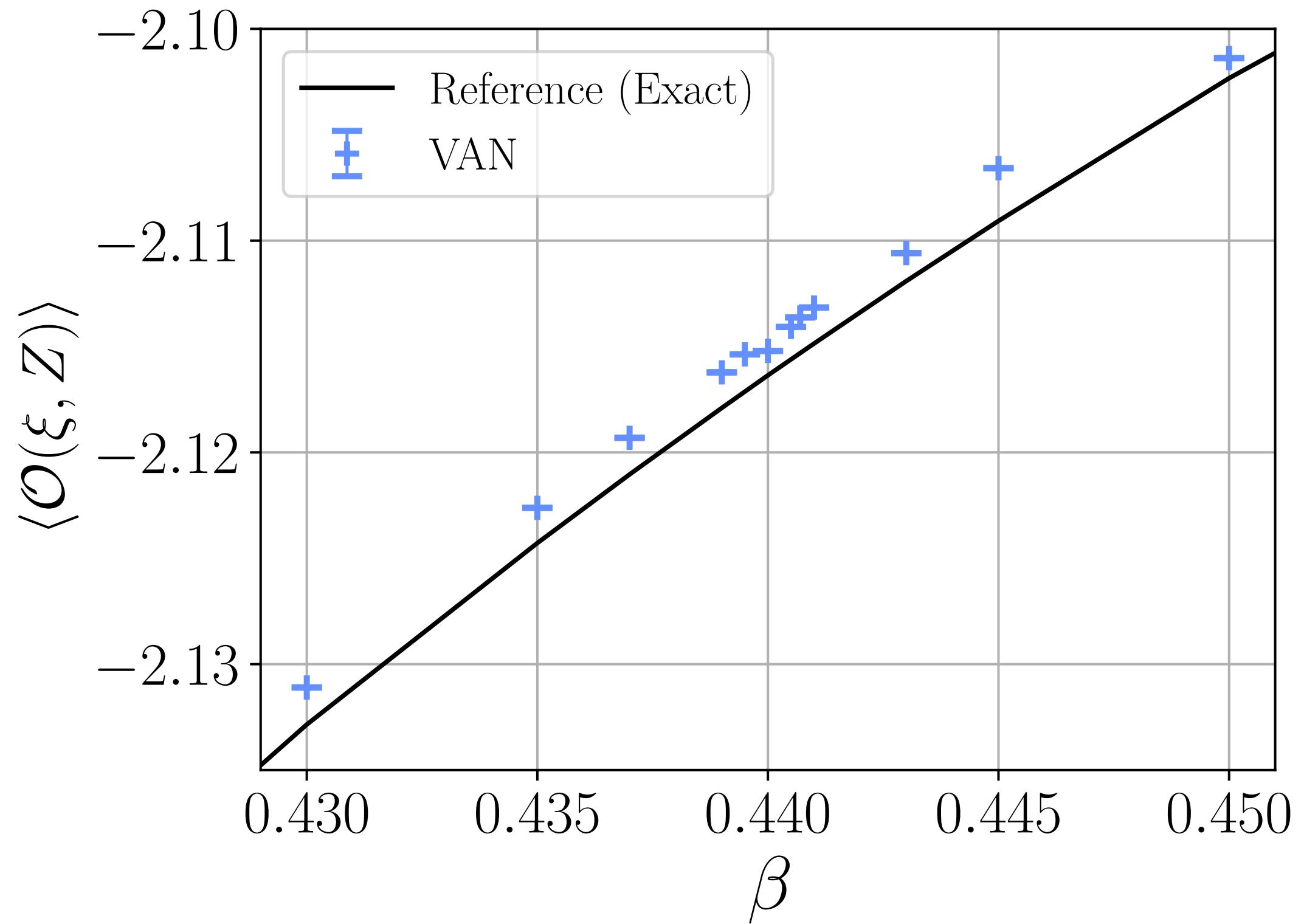


# The 2D Ising Model

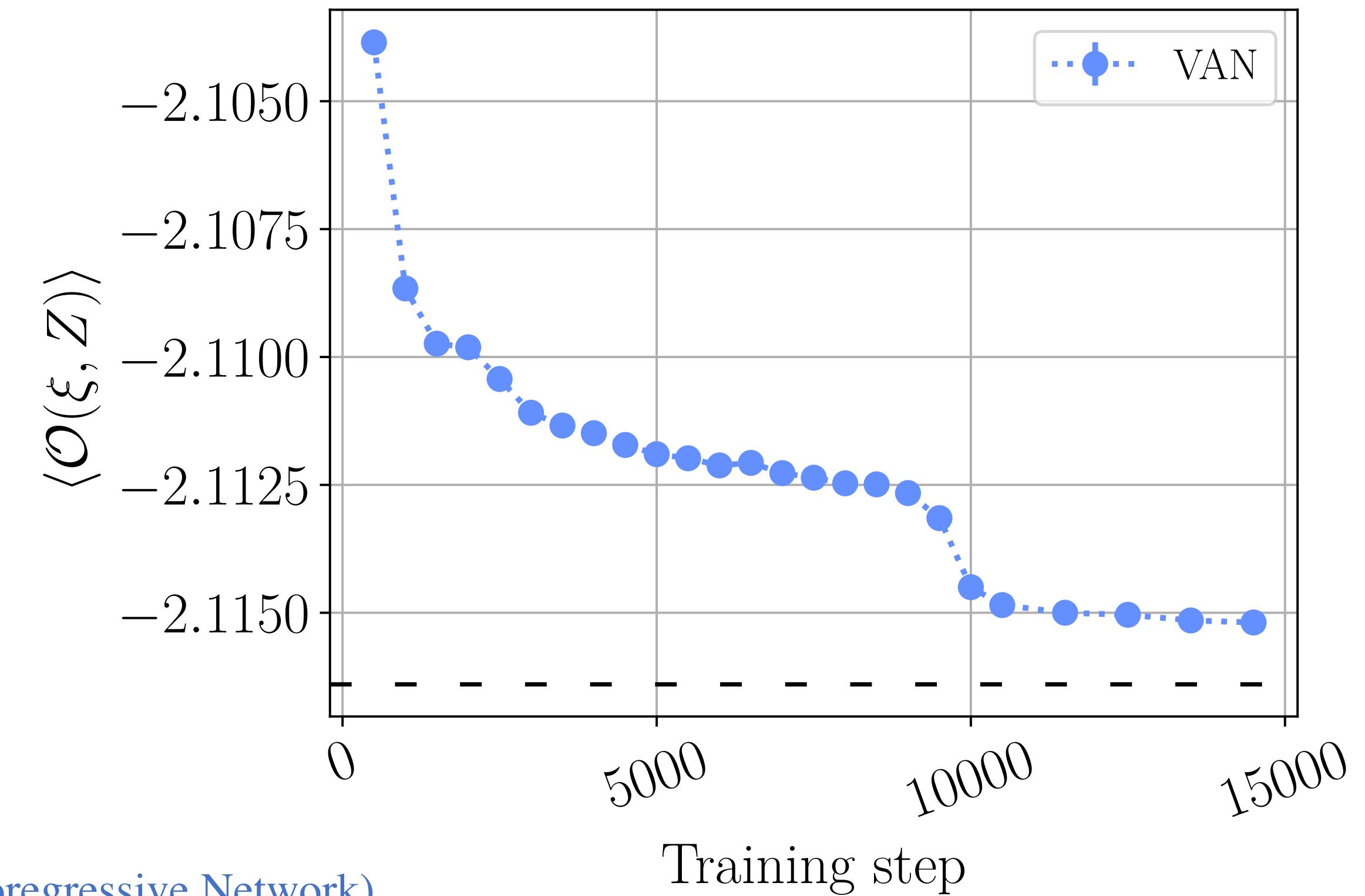
$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j$$

$$p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}, \quad \beta \propto \frac{1}{T}$$

$$\langle \mathcal{O} \rangle_{q_\theta} \neq \langle \mathcal{O} \rangle_p$$



Wu, Wang, Zhang, Phys. Rev. Lett. (2019)  $\implies$  VAN (Variational Autoregressive Network)



# Neural Markov chain Monte-Carlo (NMCMC)

---

$$p(\phi) \approx q_\theta \sim \phi_i$$

# Neural Markov chain Monte-Carlo (NMCMC)

---

$$p(\phi) \approx q_\theta \sim \phi_i$$

**Standard Metropolis**

$$A(\phi', \phi_i) = \min\left(1, \frac{p(\phi') s(\phi_i | \phi')}{p(\phi_i) s(\phi' | \phi_i)}\right) \xrightarrow{s(\phi' | \phi_i) \rightarrow q_\theta(\phi')} A(\phi', \phi_i) = \min\left(1, \frac{q_\theta(\phi_i) e^{-S(\phi')}}{q_\theta(\phi') e^{-S(\phi_i)}}\right)$$

**Neural MCMC**

# Neural Markov chain Monte-Carlo (NMCMC)

---

## Standard Metropolis

$$A(\phi', \phi_i) = \min\left(1, \frac{p(\phi') s(\phi_i | \phi')}{p(\phi_i) s(\phi' | \phi_i)}\right) \xrightarrow{s(\phi' | \phi_i) \rightarrow q_\theta(\phi')} A(\phi', \phi_i) = \min\left(1, \frac{q_\theta(\phi_i) e^{-S(\phi')}}{q_\theta(\phi') e^{-S(\phi_i)}}\right)$$

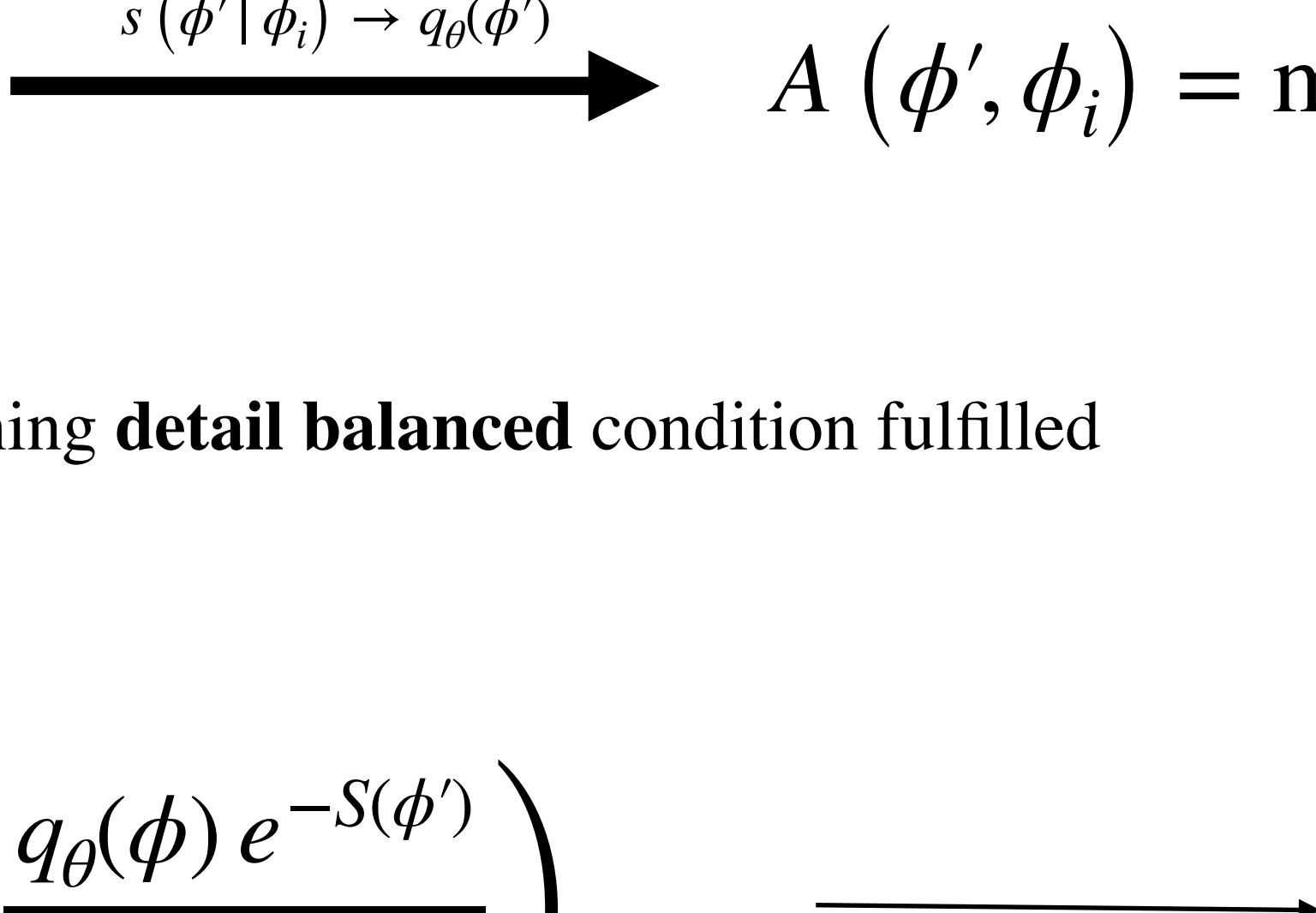
## Neural MCMC

NMCMC allows for **independent sampling** while retaining **detail balanced** condition fulfilled

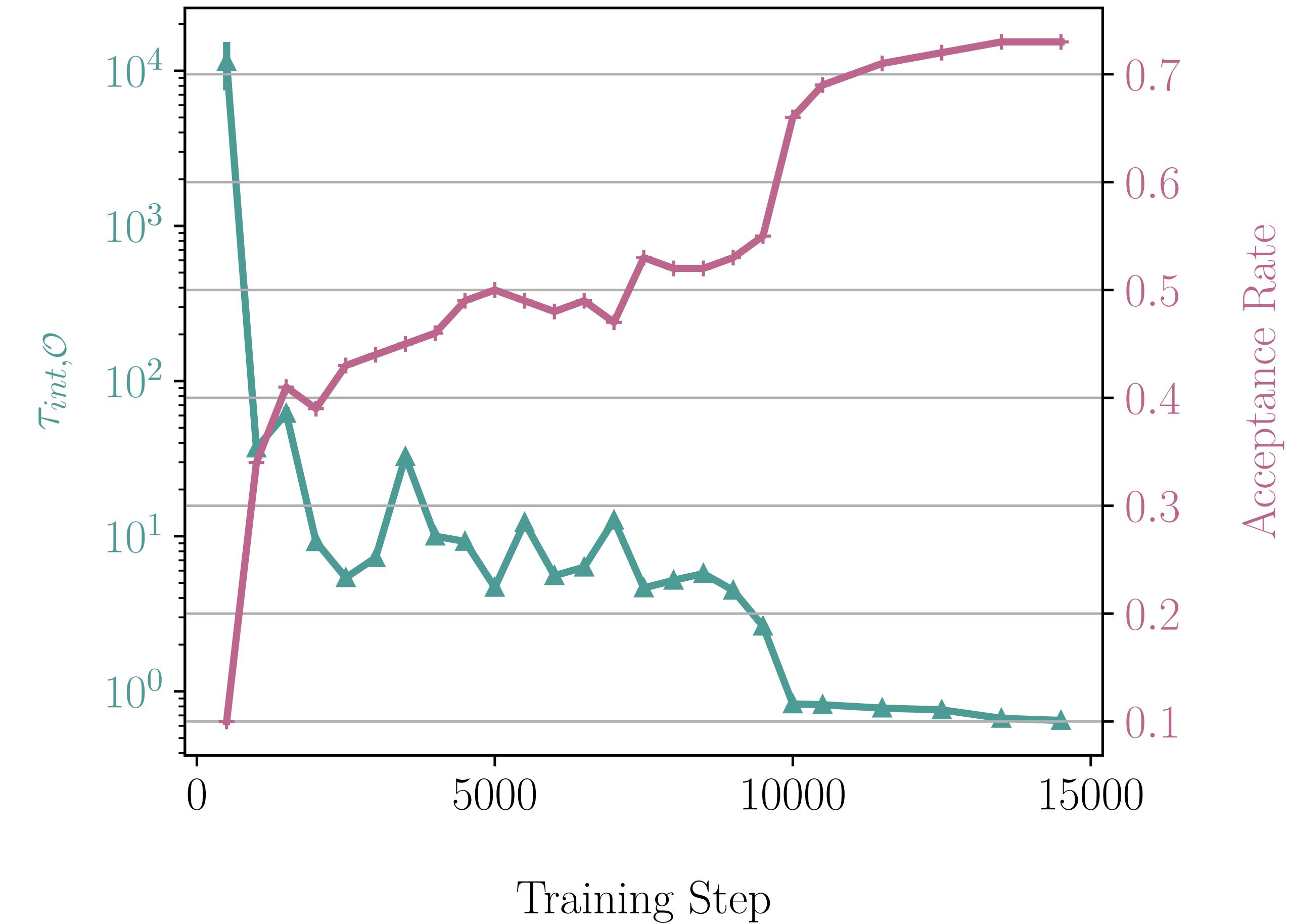
$$p(\phi' | \phi) e^{-S(\phi)} = q_\theta(\phi') A(\phi', \phi) e^{-S(\phi)}$$

$$= q_\theta(\phi') e^{-S(\phi)} \min\left(1, \frac{q_\theta(\phi) e^{-S(\phi')}}{q_\theta(\phi') e^{-S(\phi)}}\right)$$

$$= q_\theta(\phi) e^{-S(\phi')} \min\left(\frac{q_\theta(\phi') e^{-S(\phi)}}{q_\theta(\phi) e^{-S(\phi')}}, 1\right) = p(\phi | \phi') e^{-S(\phi')}$$

- 
- min operator is symmetric.
  - All factor are strictly positive.

# Neural Markov chain Monte-Carlo (NMCMC)



Nicoli, Nakajima, Strodthoff, Samek, Müller, and Kessel, Phys. Rev. E (2020)

# The effect of mode-dropping

---

**Definition.** *The effective support of the variational density  $q_\theta$  relative to  $p$*

$$\widetilde{\text{supp}}_{p,\epsilon}(q_\theta) = \{\phi \in \text{supp}(q_\theta); q_\theta(\phi) > \epsilon p(\phi)\}$$

*for a given numerical threshold  $\epsilon$ . The mode dropping set is then given by*

$$\mathcal{S} := \text{supp}(p) \setminus \widetilde{\text{supp}}_{p,\epsilon}(q_\theta)$$

# The effect of mode-dropping

---

**Definition.** *The effective support of the variational density  $q_\theta$  relative to  $p$*

$$\widetilde{\text{supp}}_{p,\epsilon}(q_\theta) = \{\phi \in \text{supp}(q_\theta); q_\theta(\phi) > \epsilon p(\phi)\}$$

for a given numerical threshold  $\epsilon$ . The mode dropping set is then given by

$$\mathcal{S} := \text{supp}(p) \setminus \widetilde{\text{supp}}_{p,\epsilon}(q_\theta)$$

if the flow is **effectively mode-dropping**, the importance-weighted estimator, with a finite number of samples  $N$ , will miss a contribution from the mass  $\int_{\mathcal{S}} p(\phi)d\phi$  with approximately the probability  $1 - \epsilon N \int_{\mathcal{S}} p(\phi)d\phi$ .

# The effect of mode-dropping

---

**Definition.** We define the effective sampler distribution

$$\tilde{q}_\theta(\phi) = \begin{cases} q_\theta(\phi)/\zeta & \text{if } \phi \in \widetilde{\text{supp}}_{p,\epsilon}(q_\theta) \\ 0 & \text{otherwise,} \end{cases} \quad \text{where} \quad \zeta = \int_{\widetilde{\text{supp}}_{p,\epsilon}} \mathcal{D}[\phi] q_\theta(\phi) \leq 1$$

is the multiplicative renormalization factor necessary to guarantee the normalization of  $\tilde{q}_\theta(\phi)$ .

# The effect of mode-dropping

---

**Definition.** We define the effective sampler distribution

$$\tilde{q}_\theta(\phi) = \begin{cases} q_\theta(\phi)/\zeta & \text{if } \phi \in \widetilde{\text{supp}}_{p,\epsilon}(q_\theta) \\ 0 & \text{otherwise,} \end{cases} \quad \text{where} \quad \zeta = \int_{\widetilde{\text{supp}}_{p,\epsilon}} \mathcal{D}[\phi] q_\theta(\phi) \leq 1$$

is the multiplicative renormalization factor necessary to guarantee the normalization of  $\tilde{q}_\theta(\phi)$ .

It follows that the importance-weighted estimator misses the contribution from the **mode-dropping set**  $\mathcal{S}$

$$\hat{\mathcal{O}} \equiv \frac{1}{N} \sum_{i=1}^N \frac{p(\phi_i)}{q_\theta(\phi_i)} \mathcal{O}(\phi_i) \approx \mathbb{E}_{\phi \sim \tilde{q}_\theta} \left[ \frac{p(\phi)}{q_\theta(\phi)} \mathcal{O}(\phi) \right] \equiv \bar{\mathcal{O}}$$

# The effect of mode-dropping

**Definition.** We define the effective sampler distribution

$$\tilde{q}_\theta(\phi) = \begin{cases} q_\theta(\phi)/\zeta & \text{if } \phi \in \widetilde{\text{supp}}_{p,\epsilon}(q_\theta) \\ 0 & \text{otherwise,} \end{cases} \quad \text{where} \quad \zeta = \int_{\widetilde{\text{supp}}_{p,\epsilon}} \mathcal{D}[\phi] q_\theta(\phi) \leq 1$$

is the multiplicative renormalization factor necessary to guarantee the normalization of  $\tilde{q}_\theta(\phi)$ .

It follows that the importance-weighted estimator misses the contribution from the **mode-dropping set**  $\mathcal{S}$

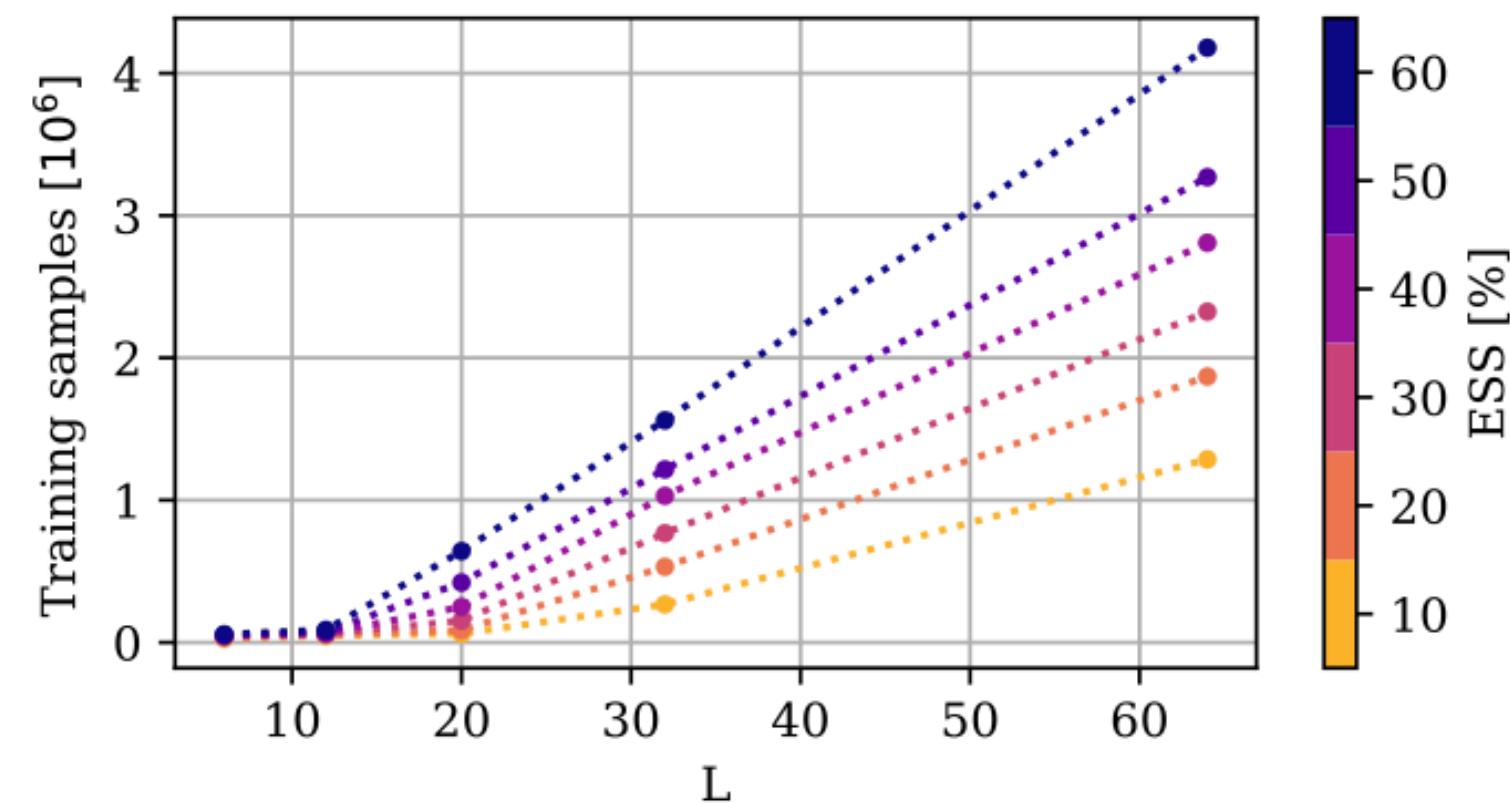
$$\hat{\mathcal{O}} \equiv \frac{1}{N} \sum_{i=1}^N \frac{p(\phi_i)}{q_\theta(\phi_i)} \mathcal{O}(\phi_i) \approx \mathbb{E}_{\phi \sim \tilde{q}_\theta} \left[ \frac{p(\phi)}{q_\theta(\phi)} \mathcal{O}(\phi) \right] \equiv \bar{\mathcal{O}}$$

the typical values of the estimator  $\hat{\mathcal{O}} \approx \bar{\mathcal{O}}$  can be **significantly different** from the true expectation value!



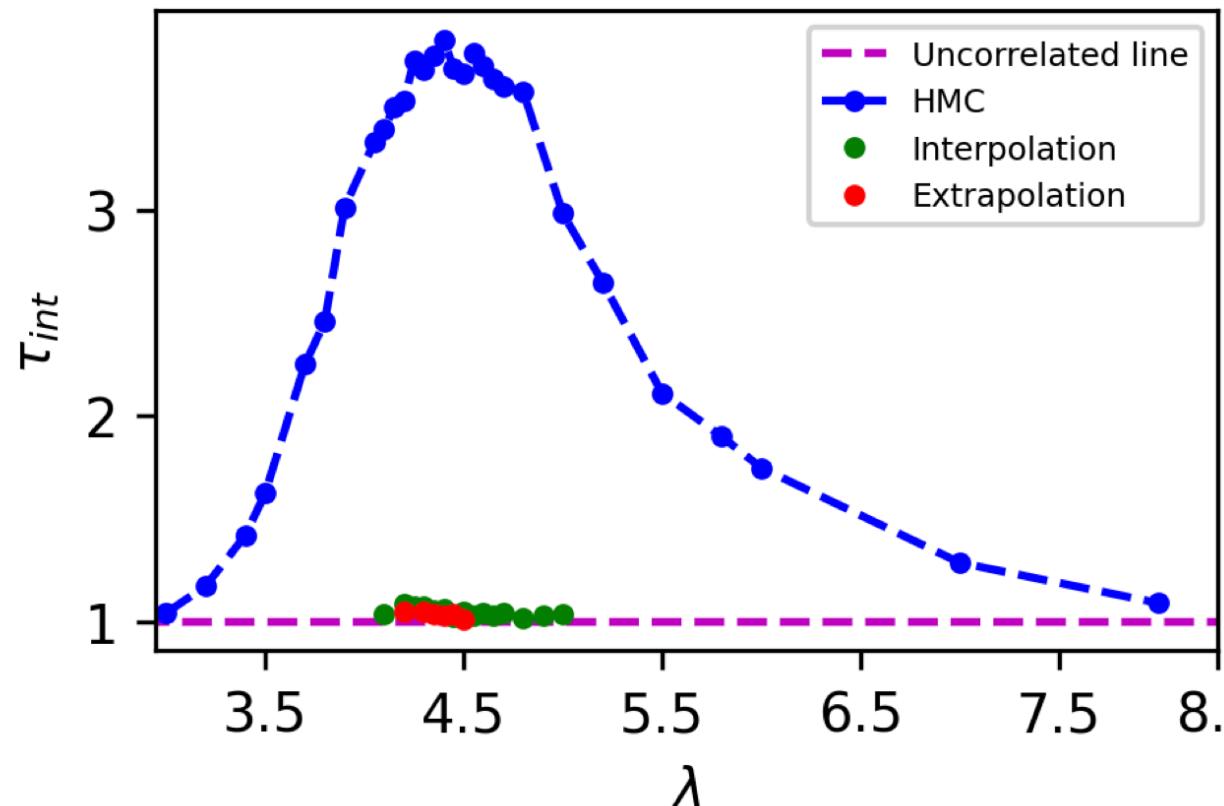
# Can We Do Better?

## Bootstrapping



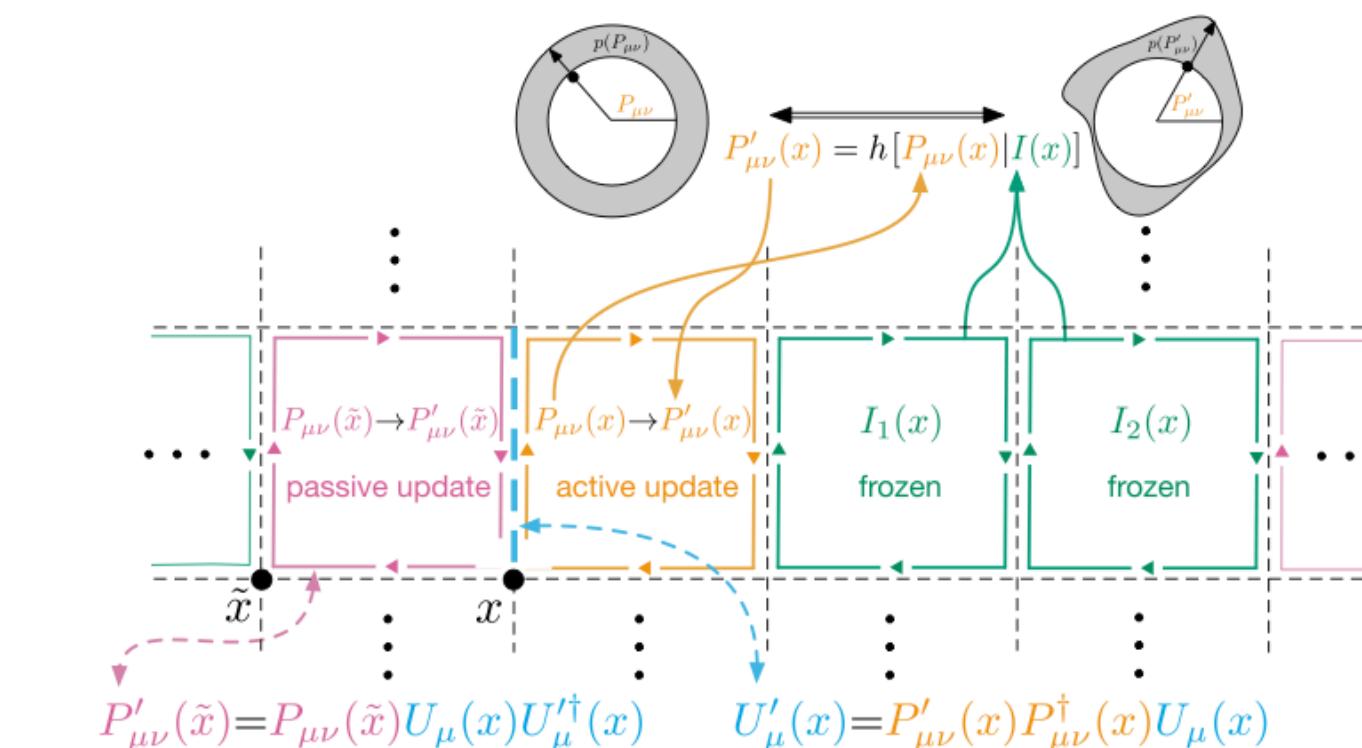
Kanwar et al., Phys. Rev. Lett. (2020)  
Gerdes et al., arXiv (2022)

## Conditioning



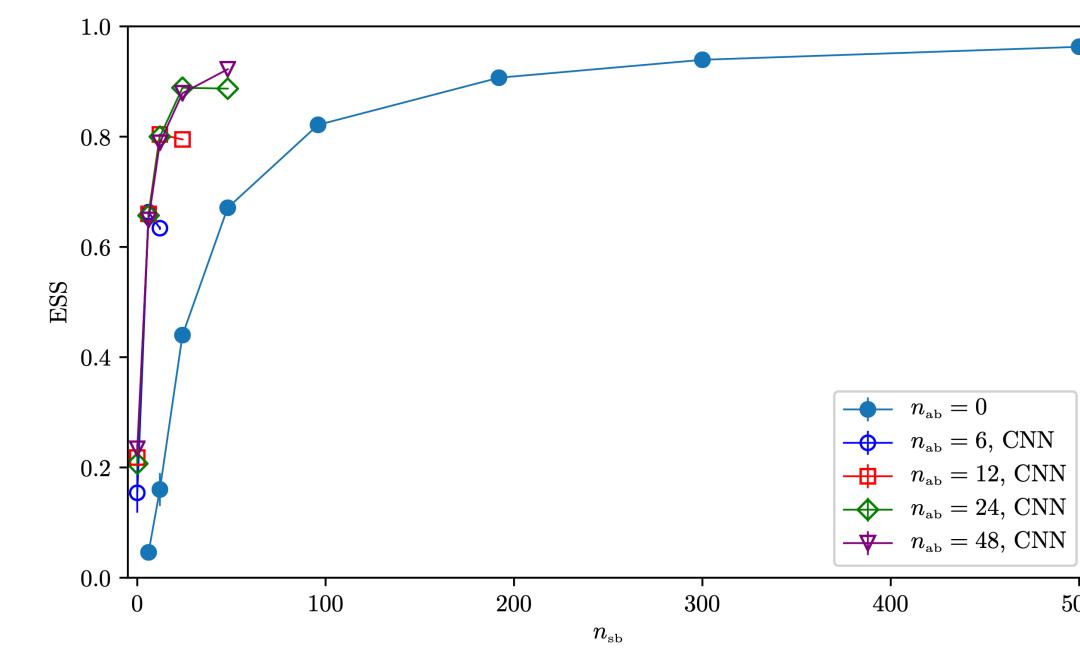
Singha et al., Phys. Rev. D (2023)  
Gerdes et al., arXiv (2022)

## Symmetries



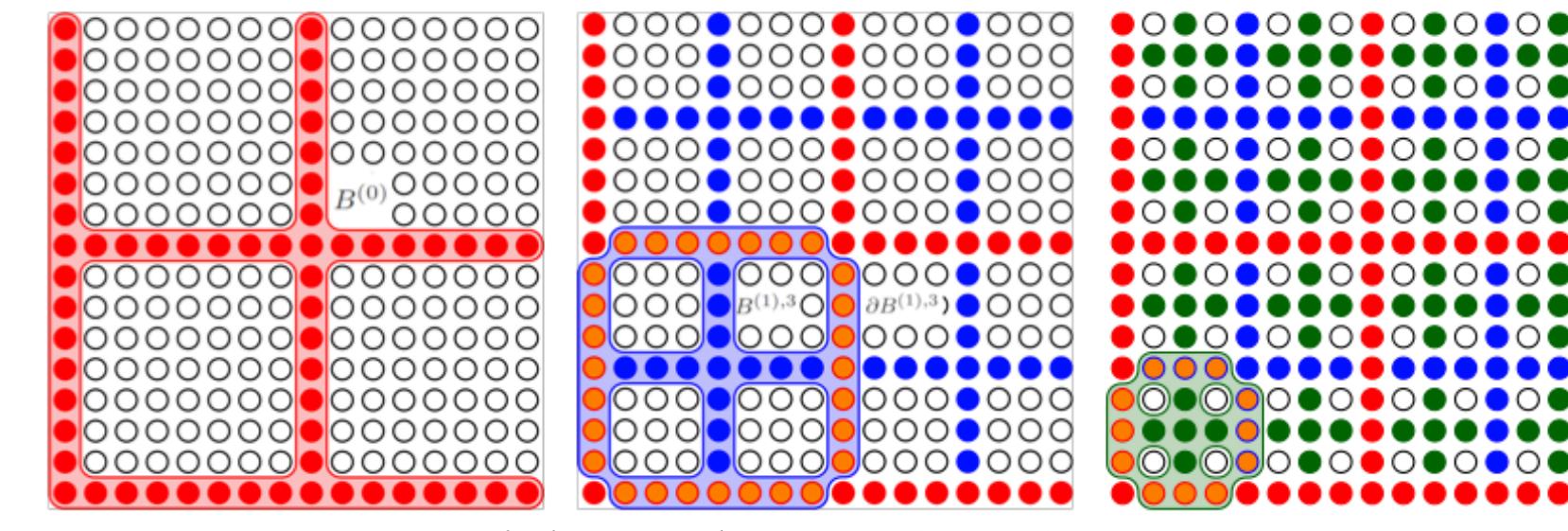
Kanwar et al., Phys. Rev. Lett. (2020)  
De Haan et al., arXiv (2021)

## Annealing Schemes and SNF



Caselle et al., JHEP. (2022)  
Caselle et al., LATTICE22 (2022)  
Matthews et al., ICML (2022) ~ CRAFT  
Wu et al., (2022) ~ SNF  
Midgley et al., ICLR (2023) ~ FAB

## Hierarchical scaling



Bialas et al., CPC (2022)

and much more...

# Estimation of Thermodynamic Observables

	<b>Training</b>	<b>Sampling*</b>
<b>FLOW</b>	20+ hrs	<1 min
<b>HMC</b>	None	$\sim 1.8 \text{ hrs} \times 14 = 25 \text{ hrs}$

\*sampling refers to draw the same number of 5.6 million configurations

⚠ Variance of the NIS estimator scales with  $O(N^{-1})$  → The flow is a much more efficient sampler!

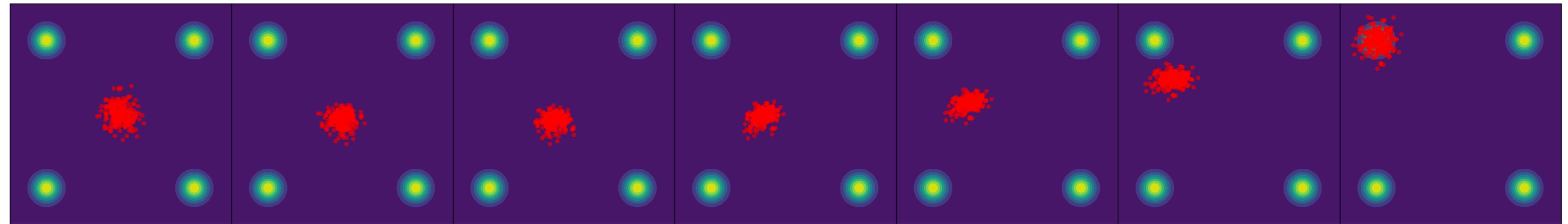
# Outlook Overview

---

- Scaling to larger volumes (Abbott et al., 2022).
- Leveraging better training schemes (i.e. different divergences and modified objectives).
- Develop physics informed coupling layers to leverage known symmetries.
- Use approximated density to study the change of topology in the presence of a phase transition.

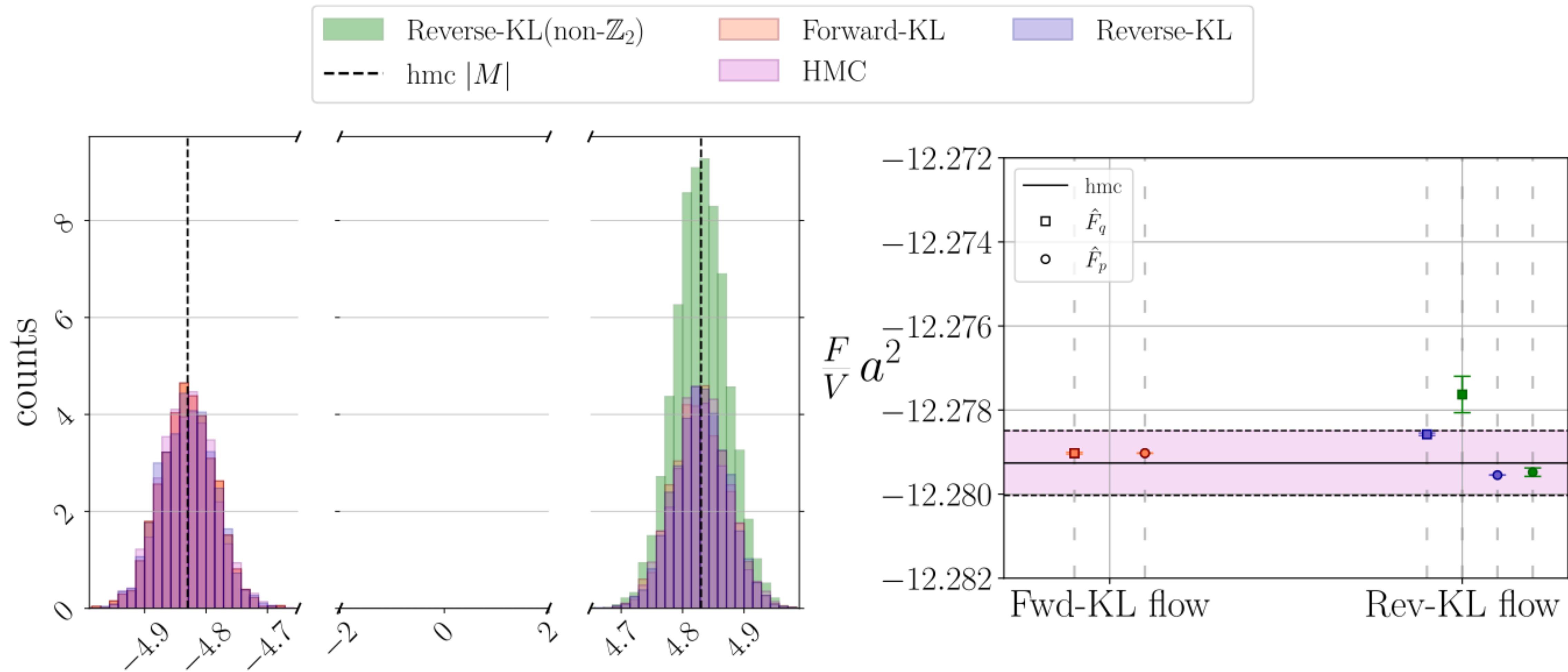
# Multimodality Makes Things Harder

---

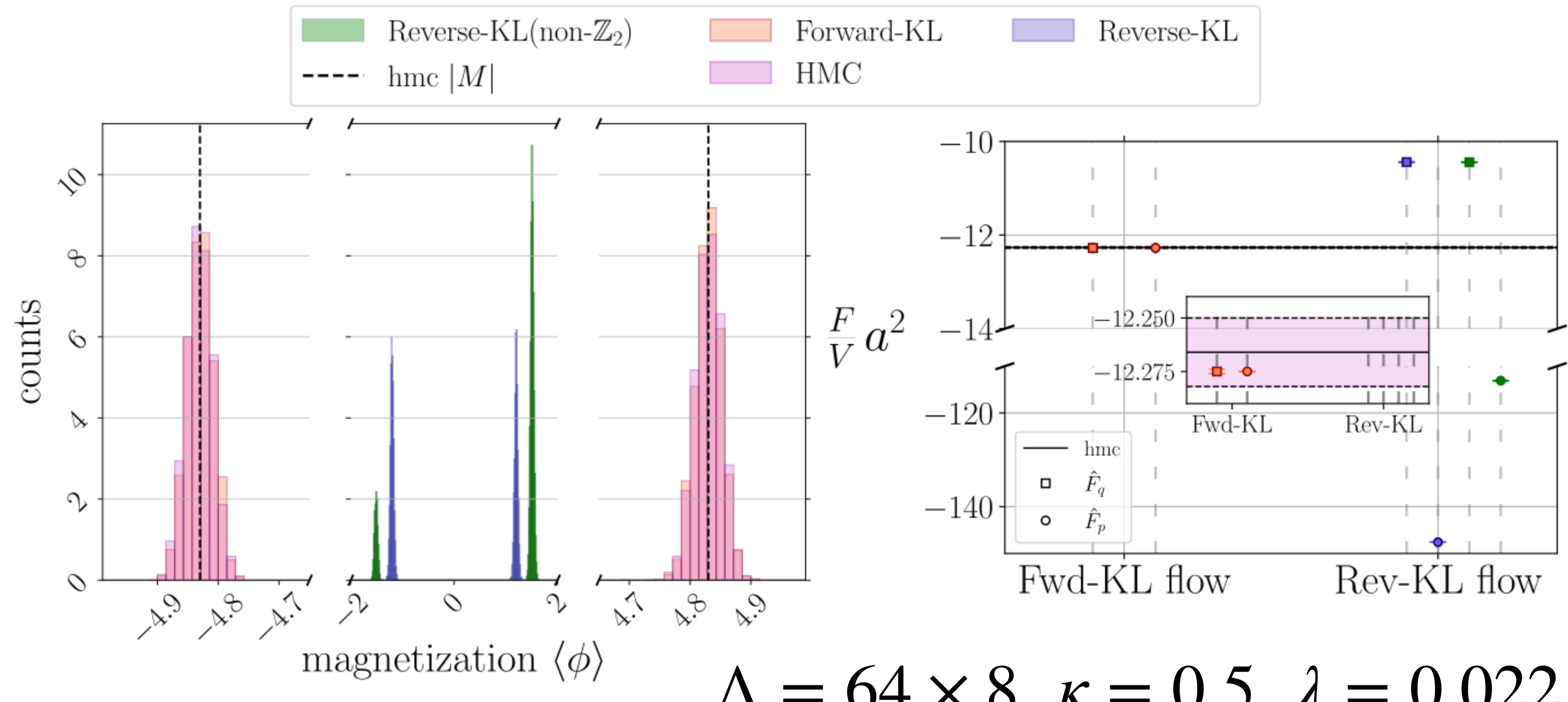


Maté B and Fleuret F., arXiv:2301.07388 (2023)

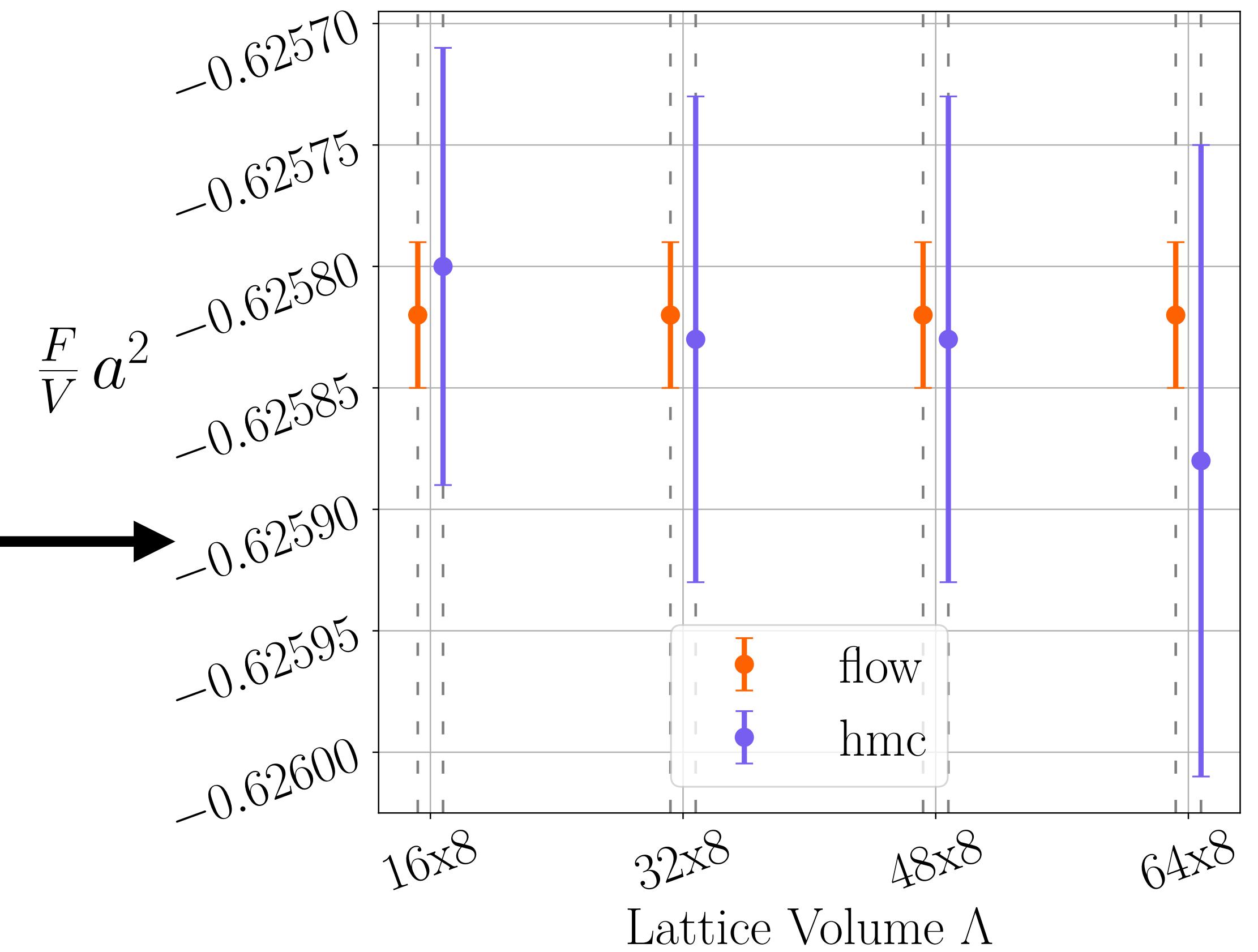
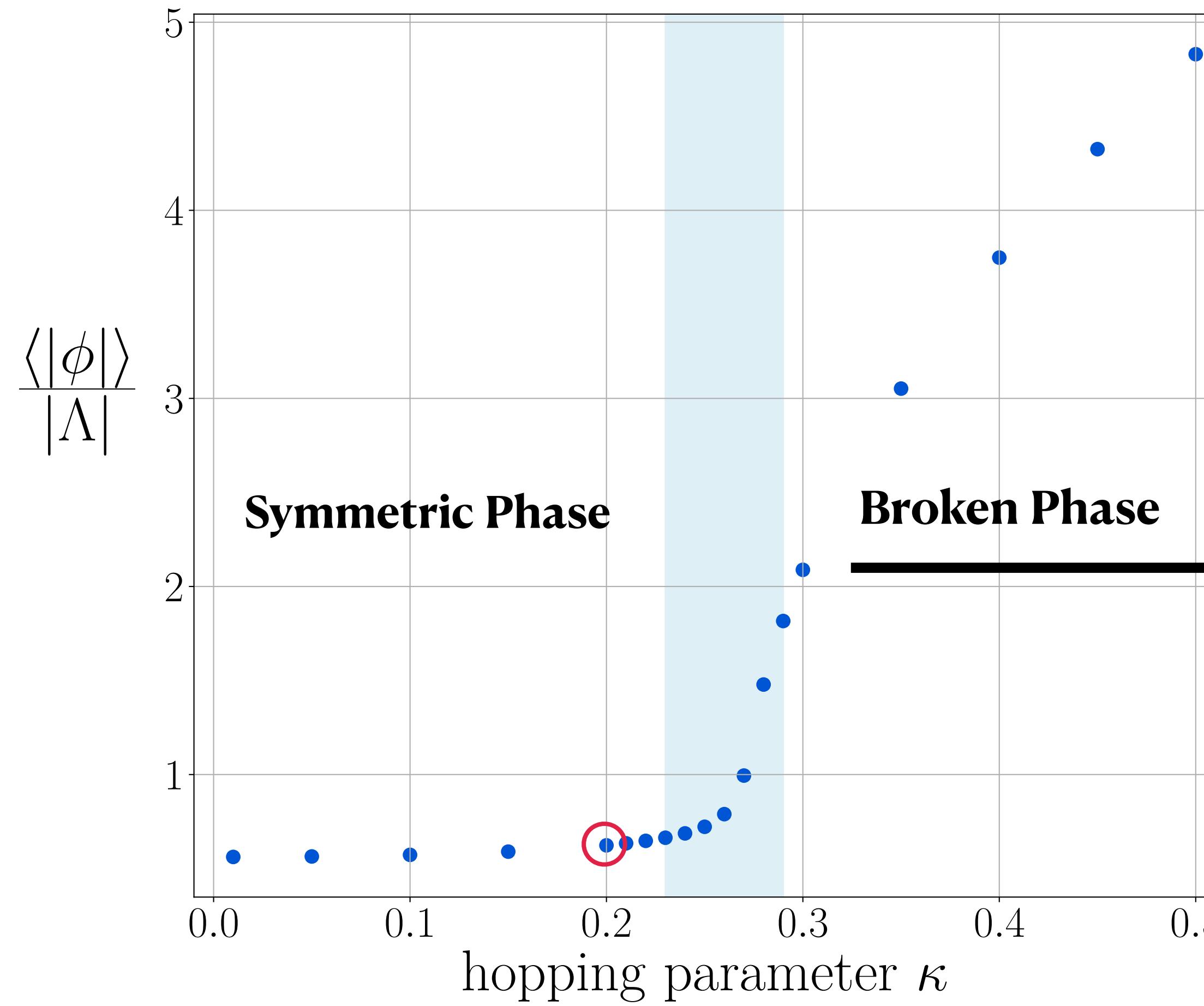
# Intentionally Mode Dropping Flow



# Bounds on the Free Energy Estimator

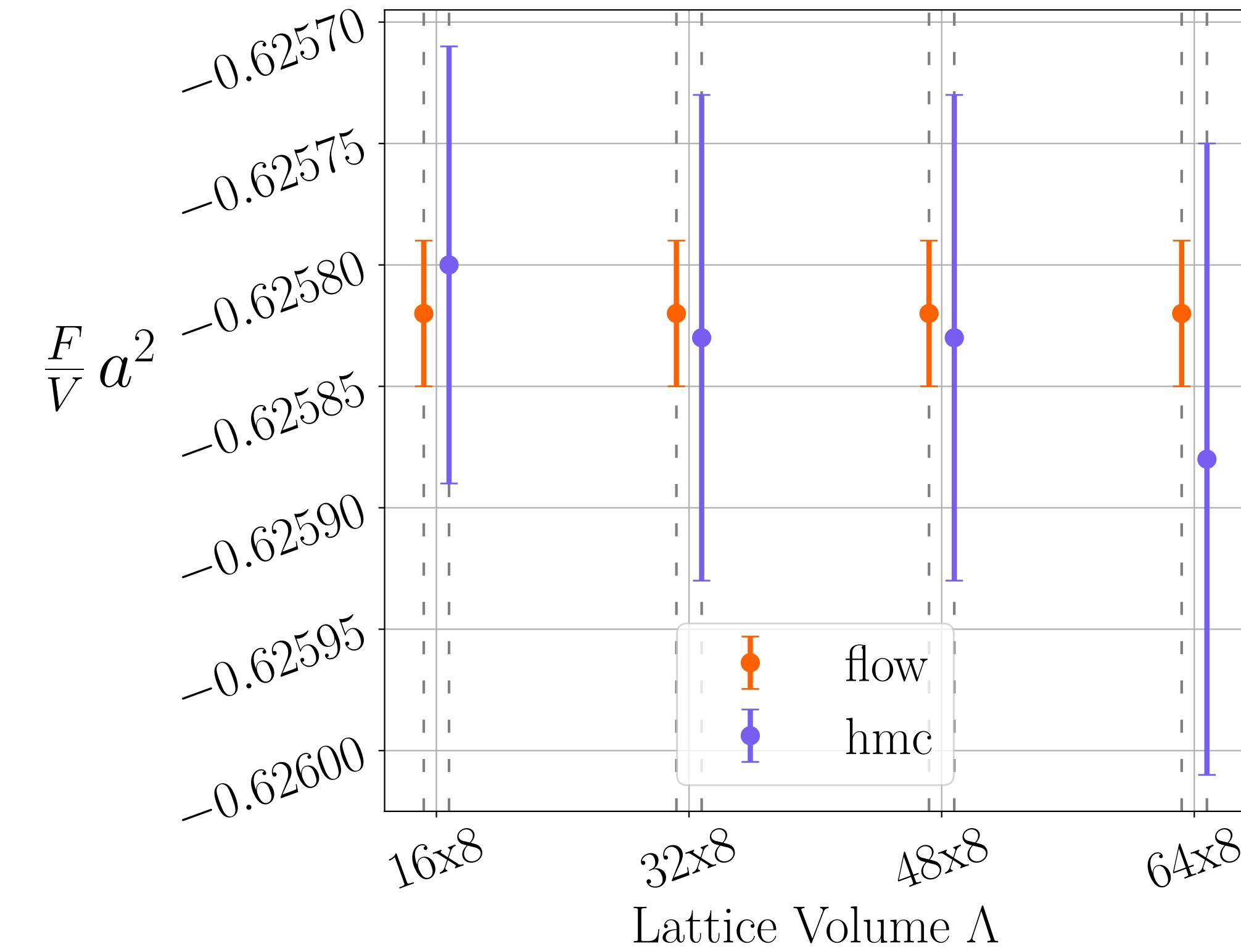
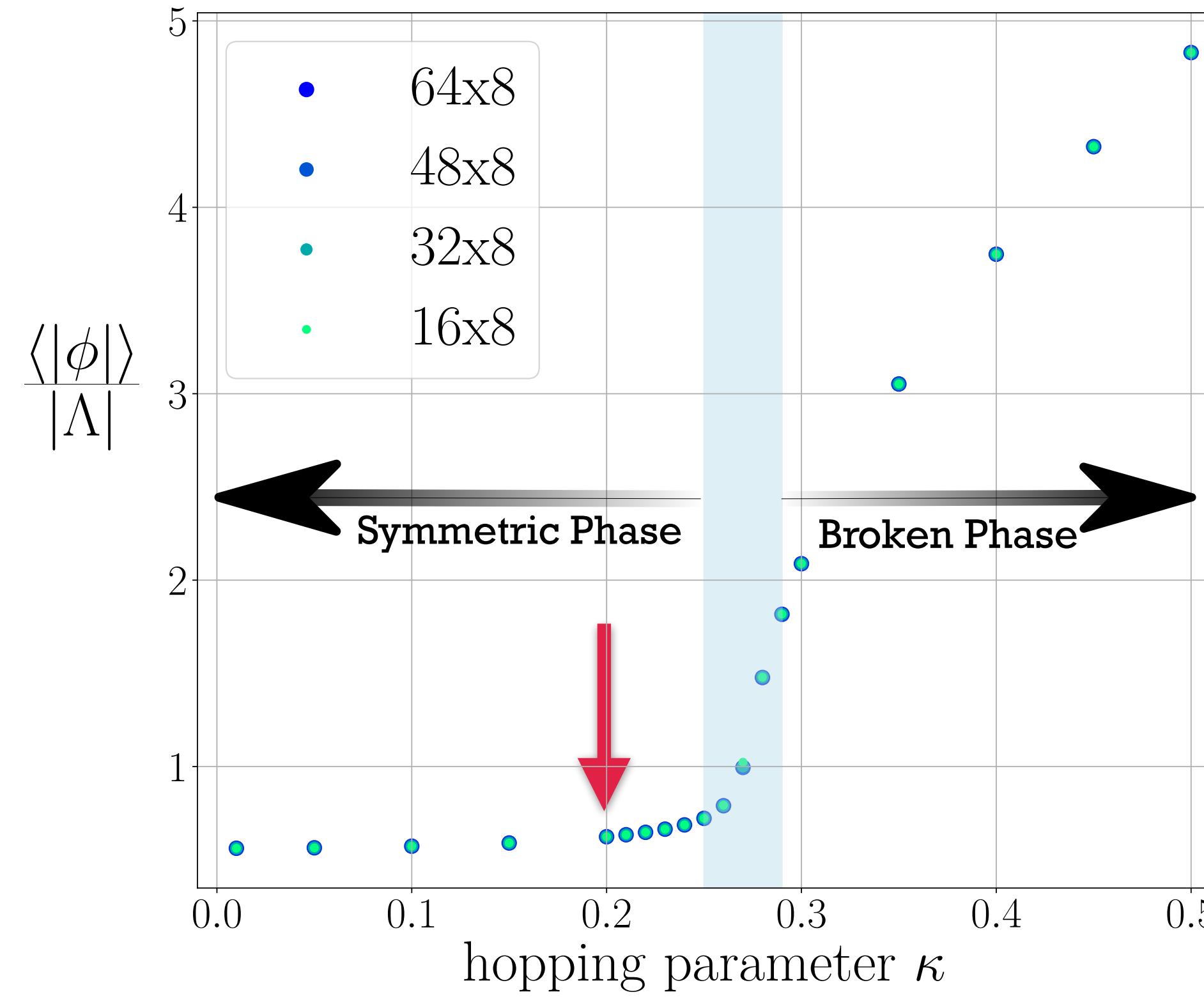


# Estimation of Thermodynamic Observables



NKA, Anders, Funcke, Hartung, Jansen, Kessel, Nakajima, Stornati, Phys. Rev. Lett. (2021)

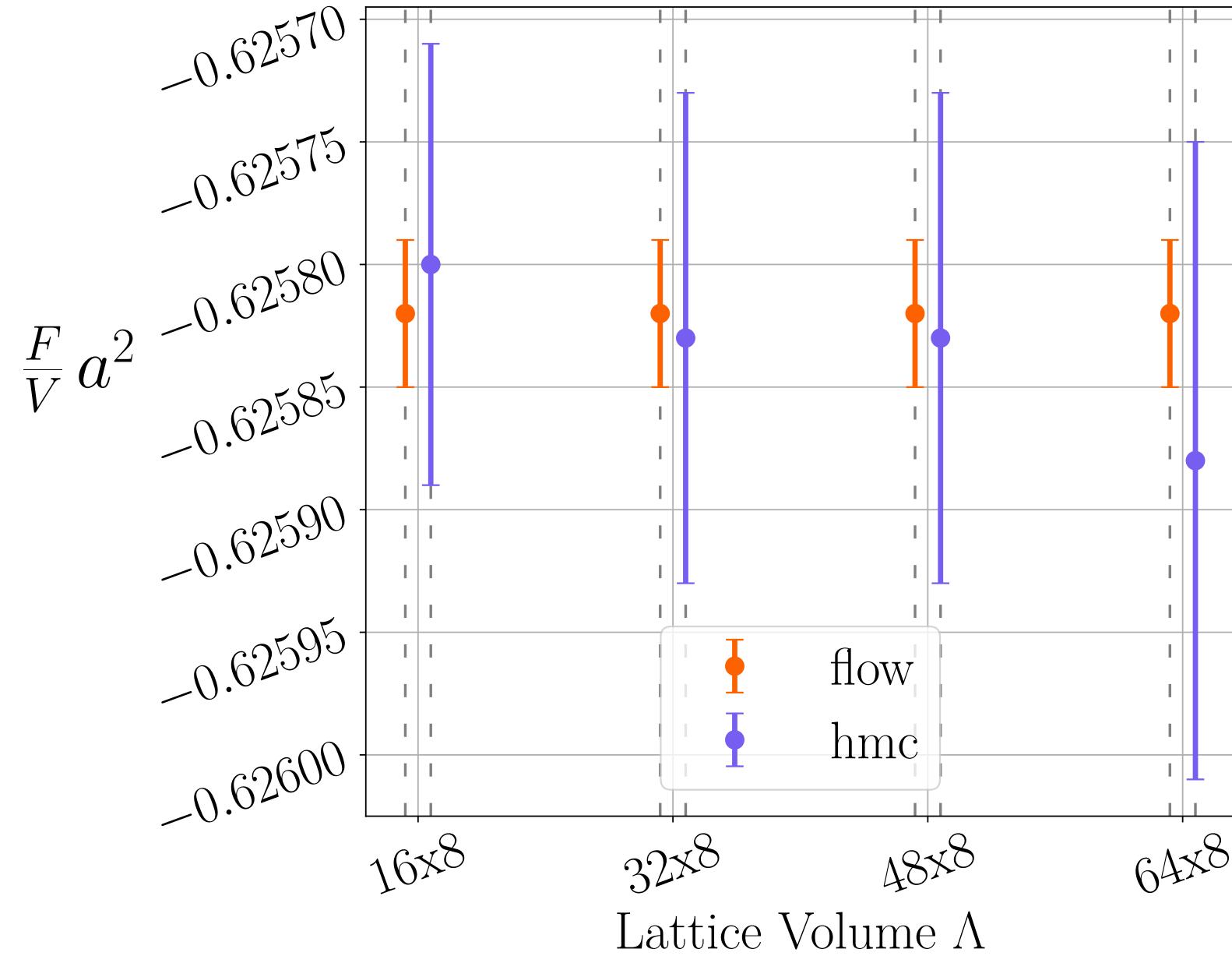
# Estimation of Thermodynamic Observables



KAN, Anders, Funcke, Hartung, Jansen, Kessel, Nakajima, Stornati, Phys. Rev. Lett. (2021)

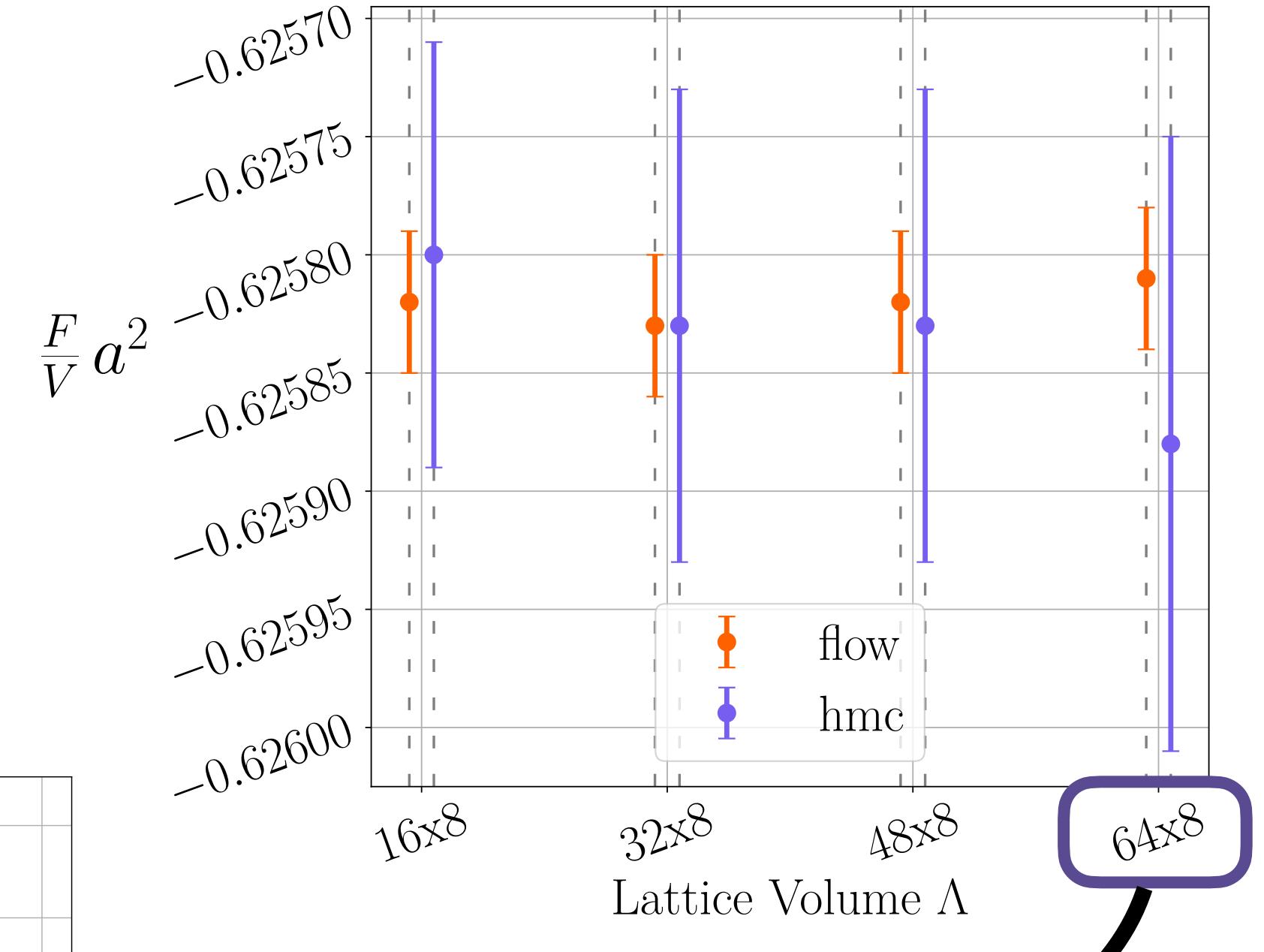
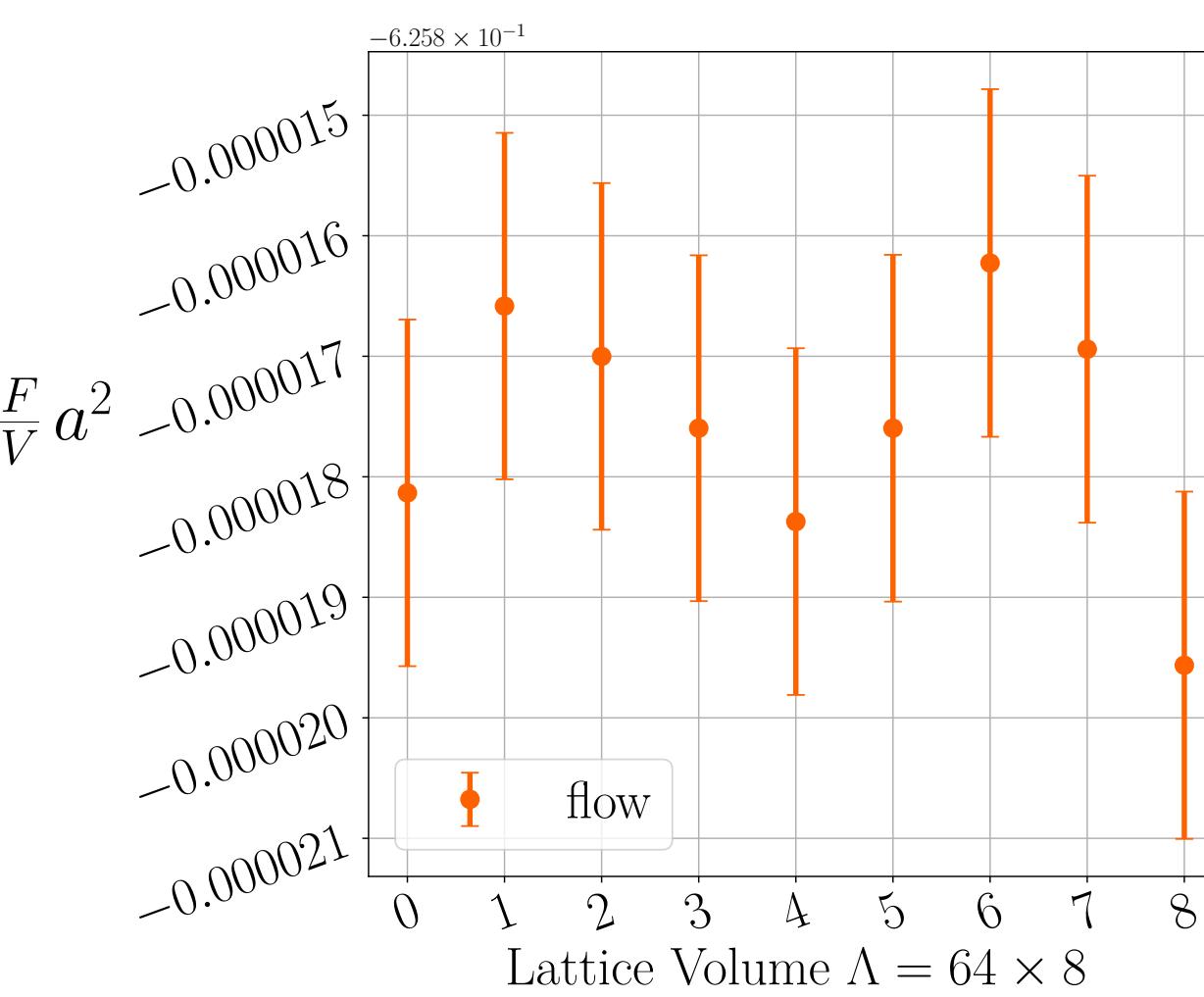
- $\lambda = 0.022, \kappa = 0.2$
- HMC:  $4 \times 10^5$  samples per chain (4 chains)
- NIS: same number of samples

# Fluctuations of the Estimator for $\kappa = 0.2$



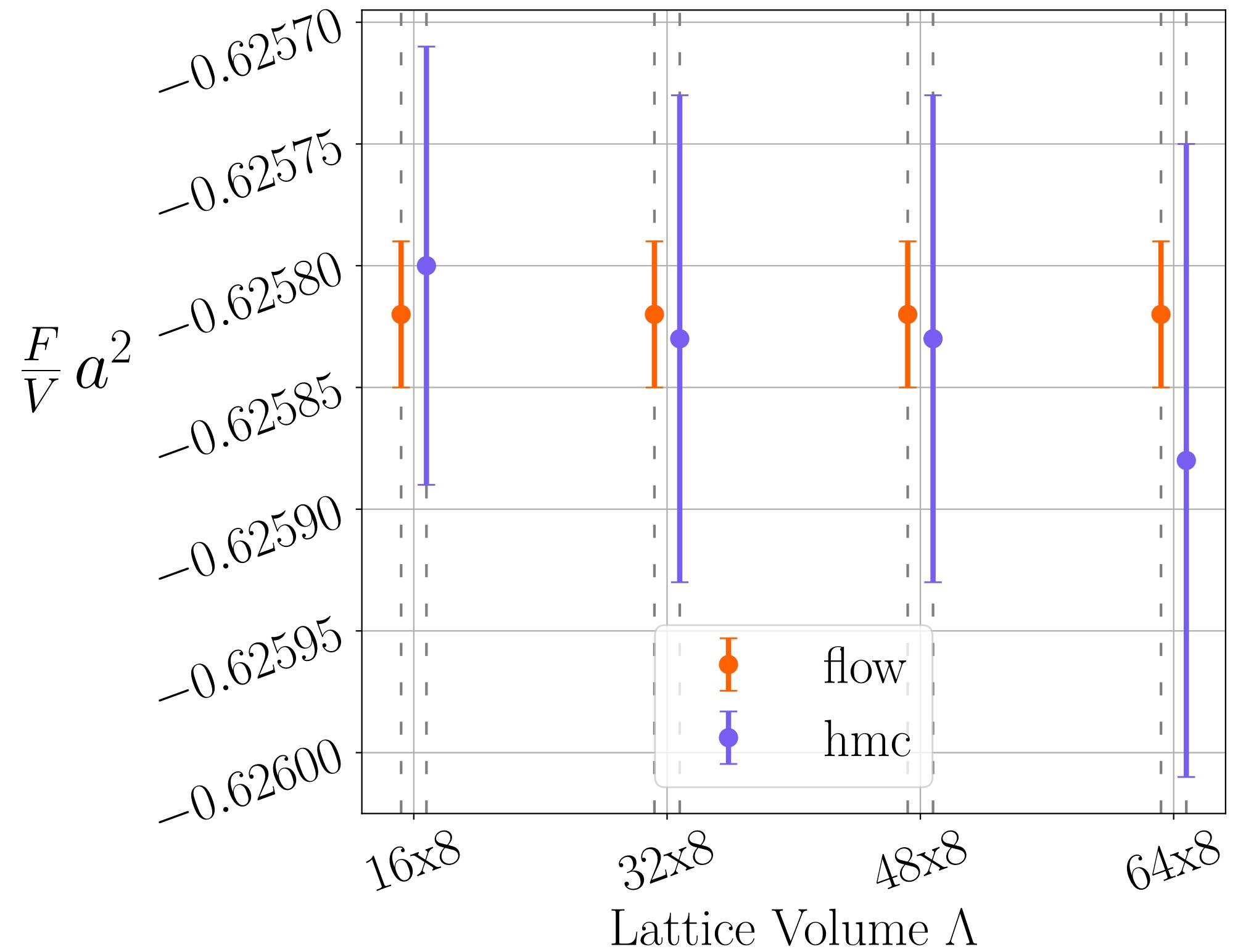
- $\lambda = 0.022, \kappa = 0.2$
- HMC:  $4 \times 10^5$  samples per chain (4 chains)
- NIS: same number of samples (1.6 M)

Repeating the experiment slightly changes the points (is a random process).



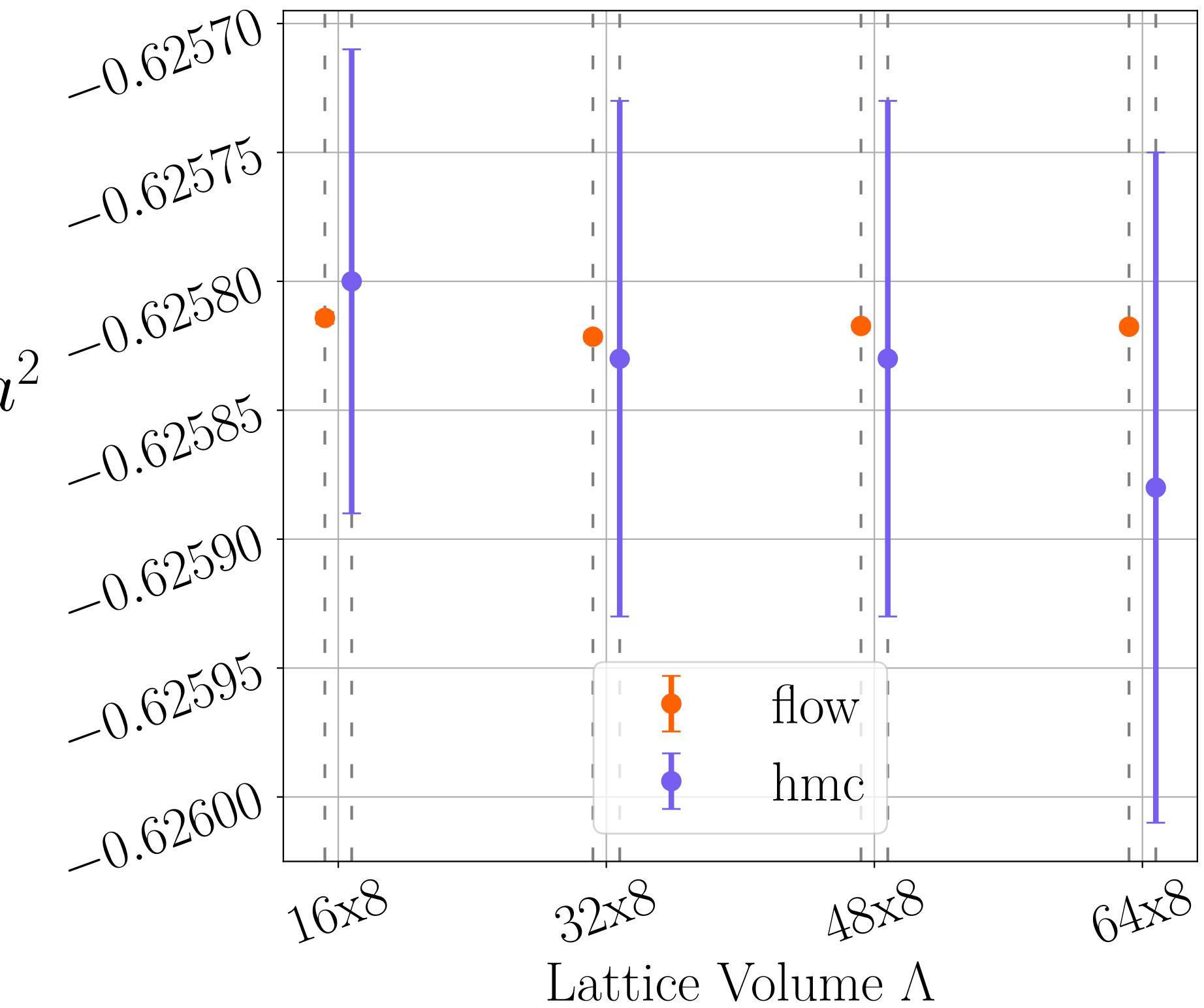
Confirmed by fixing the volume and repeating the experiment 9 times. The mean fluctuates leading to statistically consistent results.

# Better Model, Smaller Uncertainties



Model trained for 300k steps

**Same number of  
samples  $N = 1.6 \text{ M}$**



Model trained for 700k steps

# The $\phi^4$ -theory explained

---

$$S(\varphi) = \sum_{x \in \Lambda} a^2 \frac{1}{2} \sum_{\hat{\mu}=1}^2 \frac{(\varphi(x + a\hat{\mu}) - \varphi(x))^2}{a^2} + \frac{m_0^2}{2} \varphi^2(x) + \frac{g_0}{4!} \varphi^4(x)$$



$$\varphi = (2\kappa)^{\frac{1}{2}}\phi$$

$$(am_0)^2 = \frac{1 - 2\lambda}{\kappa} - 4$$

$$a^2 g_0 = \frac{6\lambda}{\kappa^2}$$

$$S(\phi) = \sum_{x \in \Lambda} -2\kappa \sum_{\hat{\mu}=1}^2 \phi(x)\phi(x + \hat{\mu}) + (1 - 2\lambda)\phi(x)^2 + \lambda\phi(x)^4$$

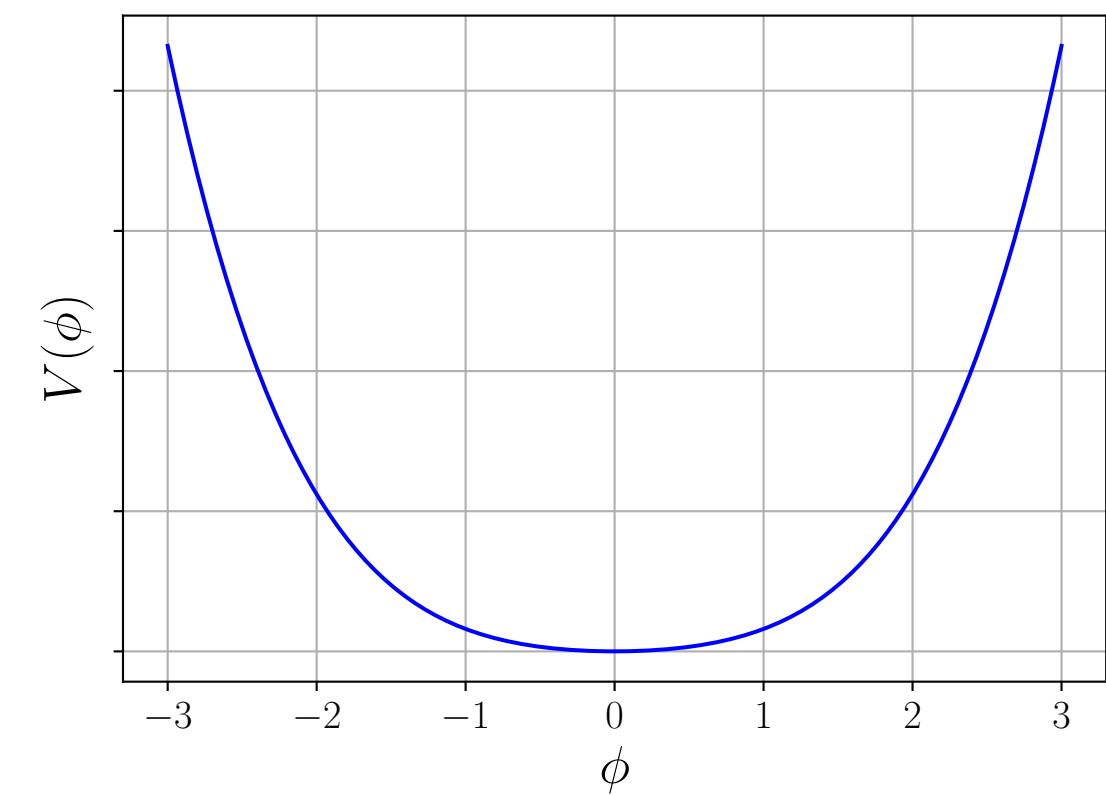
# The $\phi^4$ -theory explained

$m_0^2 = -m^2 + T^2$  is the effective mass.

$$m^2 \ll T^2$$

Early Universe

The temperature term dominates and the effective mass ends up positive

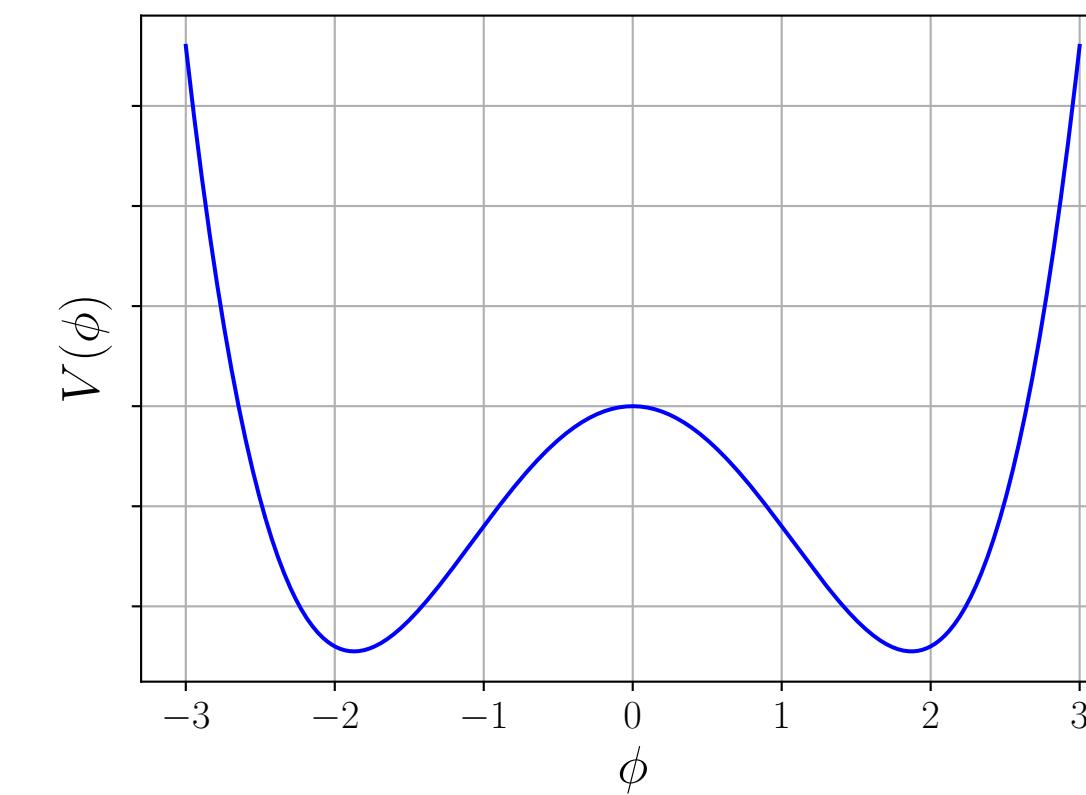


$$m^2 \gg T^2$$

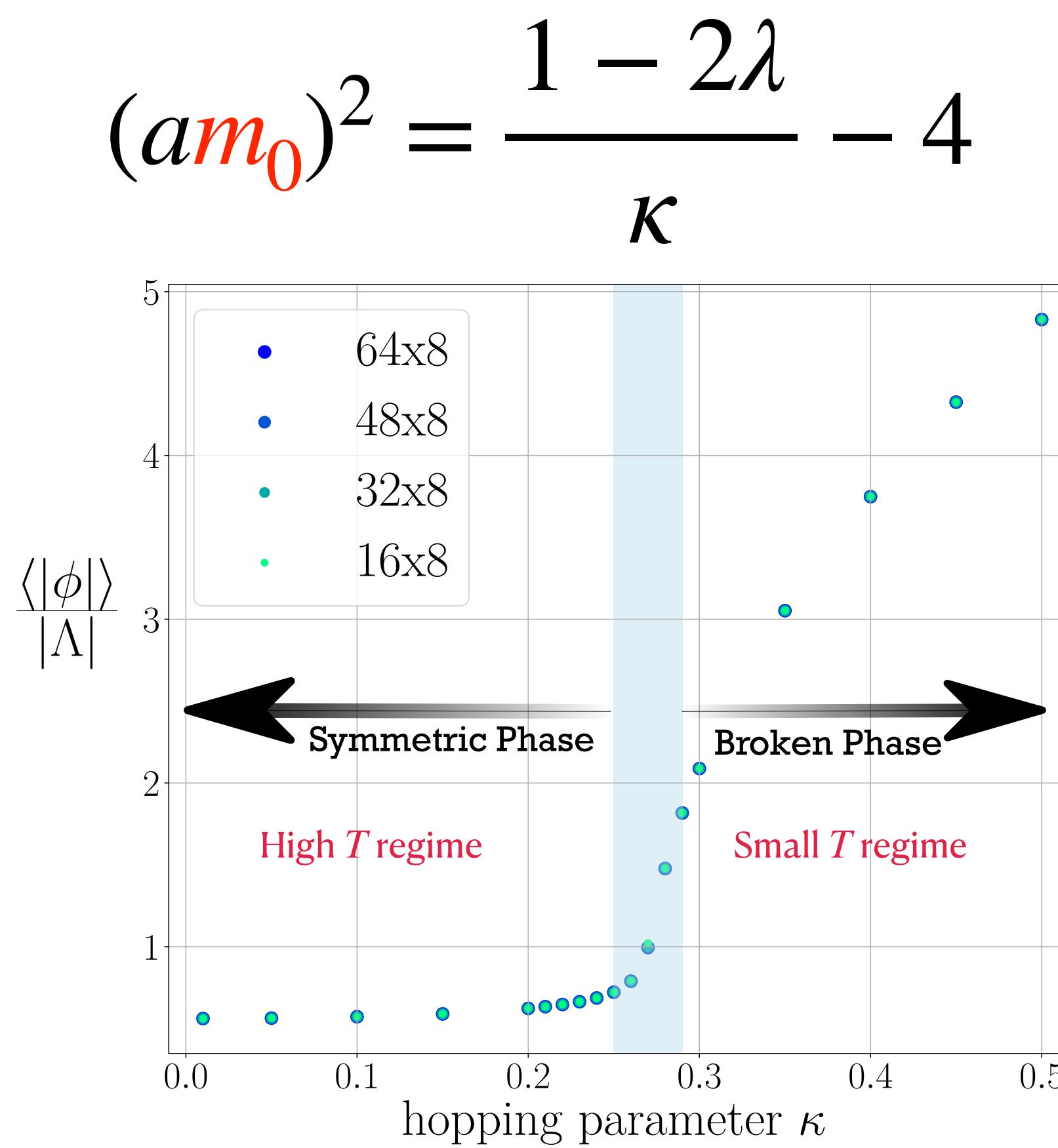
Our Universe

The mass term dominates and the effective mass ends up negative

Spontaneous Symmetry Breaking  
Temperature Decreases



# The $\phi^4$ -theory explained



$$\kappa > \kappa_c$$

The **second term** dominates and the reparametrized effective mass is **negative**

$$\kappa < \kappa_c$$

The **first term** dominates and the reparametrized effective mass is **positive**

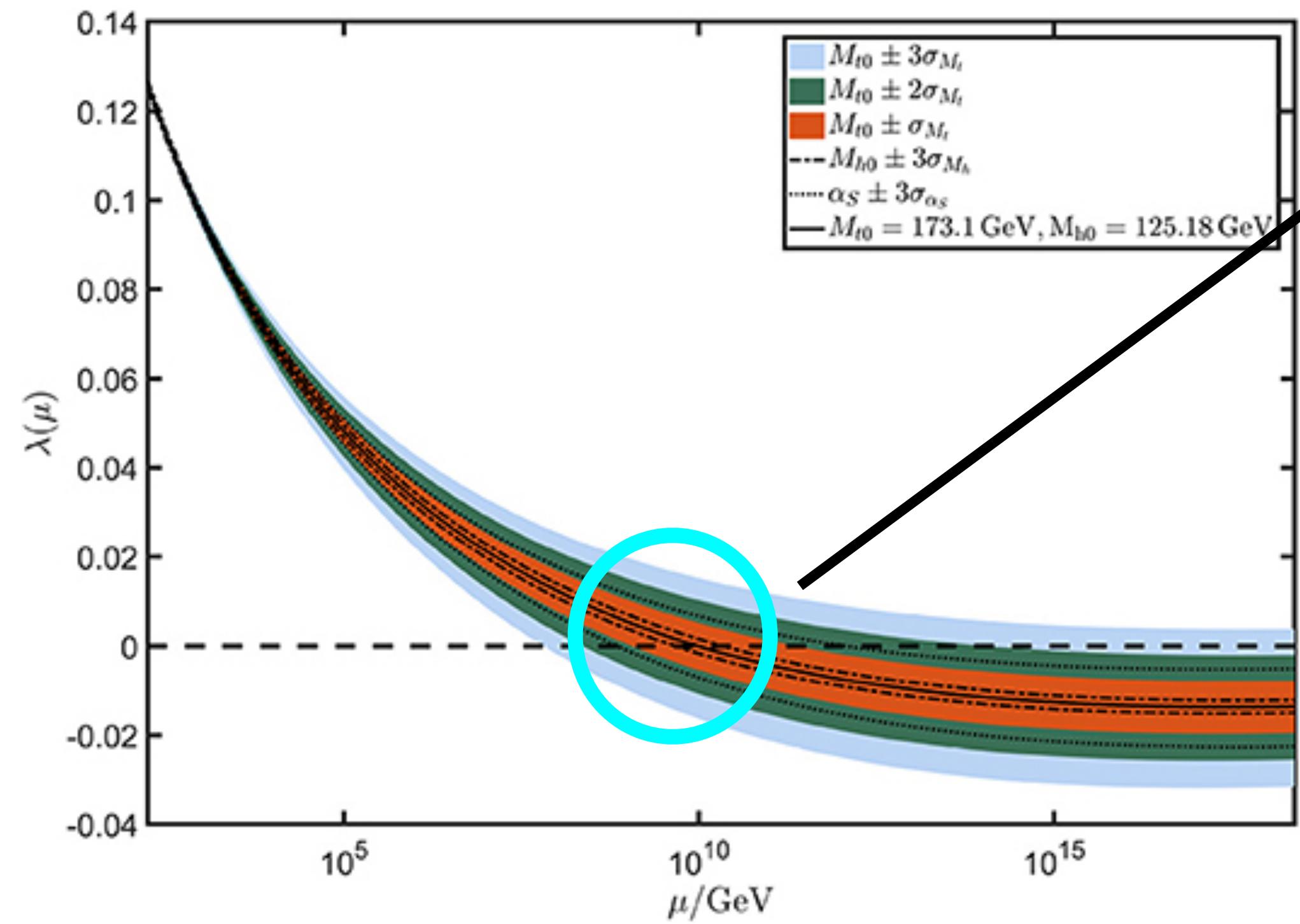
Lower  $\kappa \implies$  higher  $T$



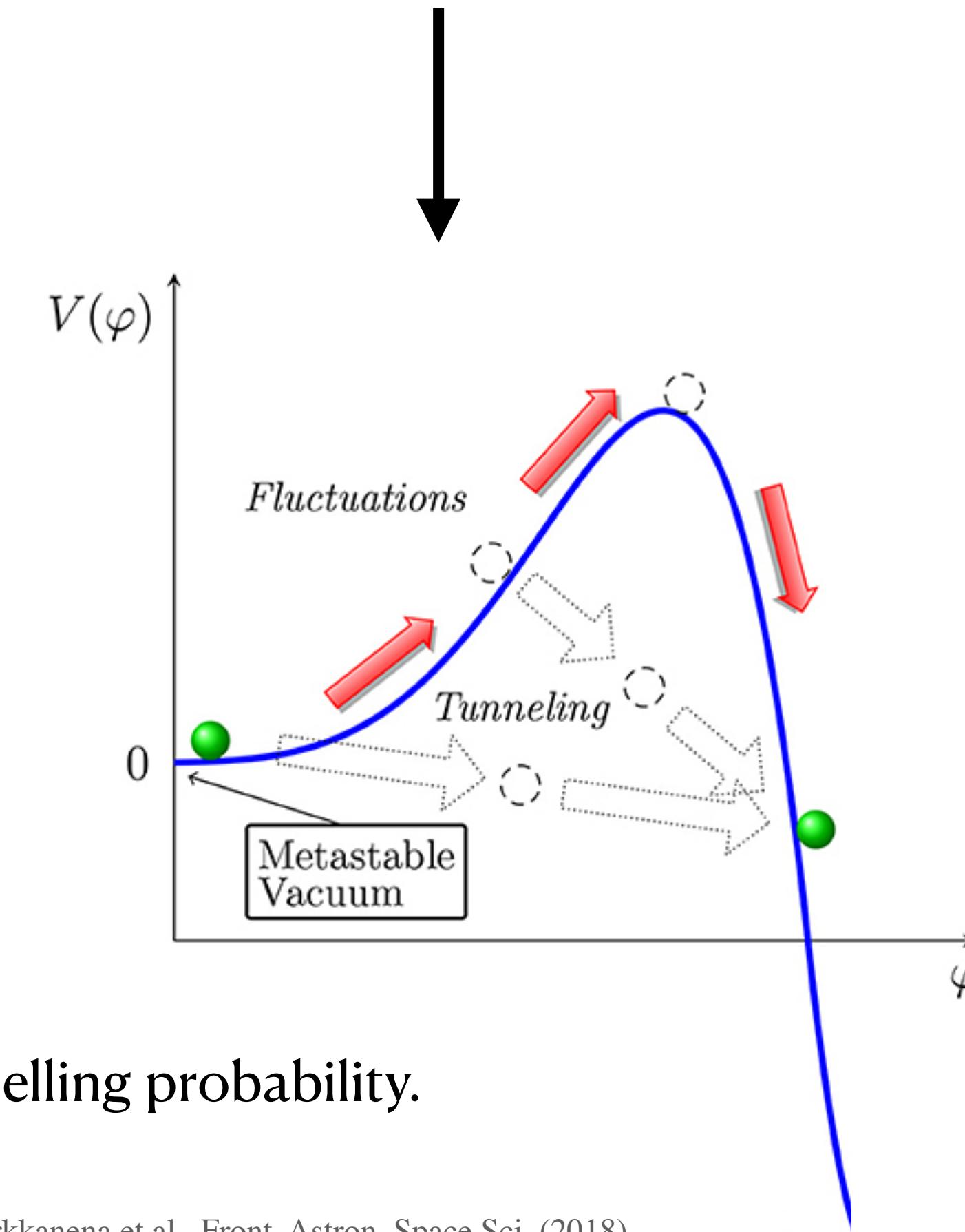
Cool down of the universe leads to **spontaneous symmetry breaking in the Higgs**

Higher  $\kappa \implies$  lower  $T$

# Triviality of $\phi^4$ -theory



Turning point of  $\lambda(\mu)$  becoming negative as the energy increases.



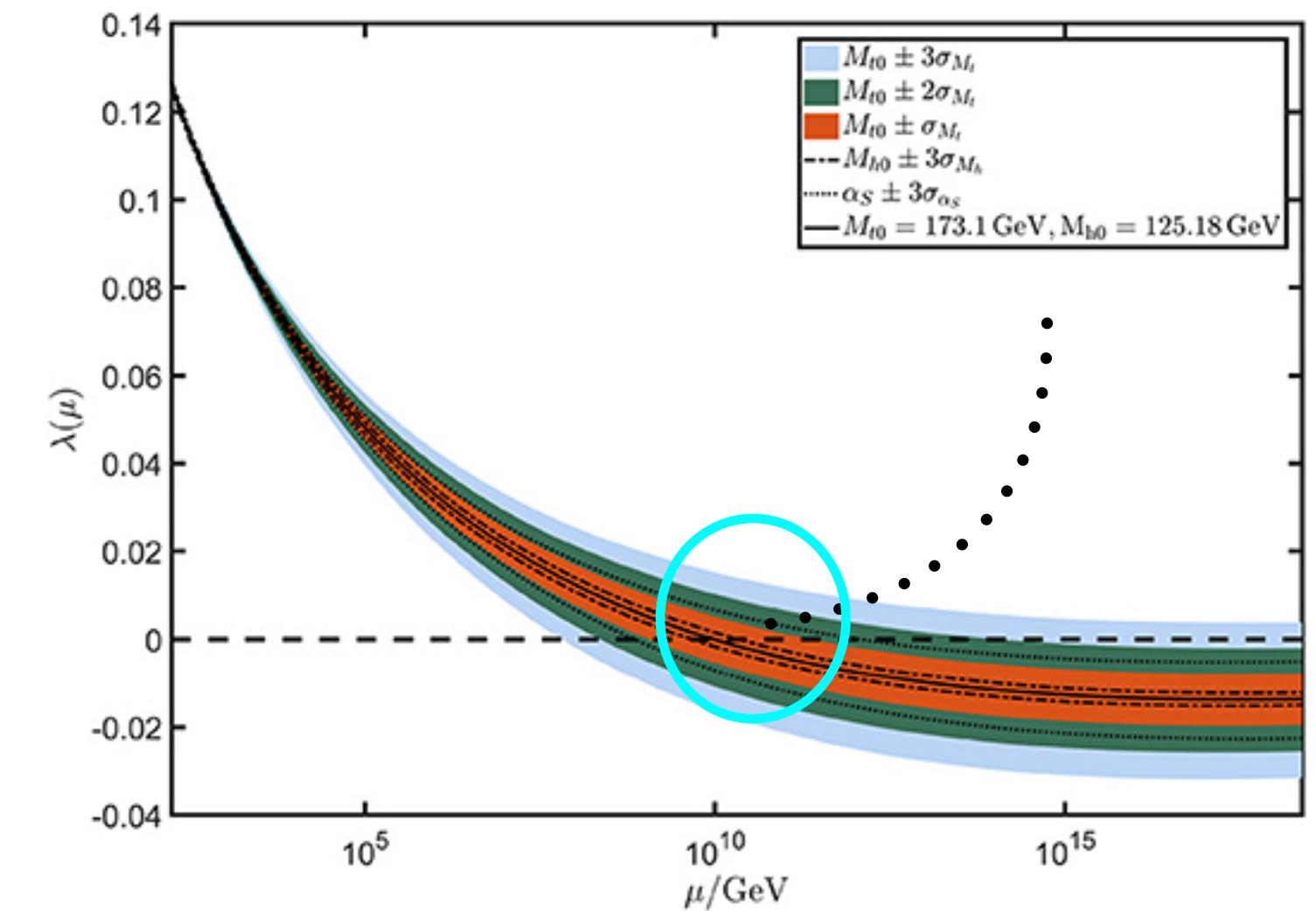
- In the Lagrangian this **flips the potential** thus making a non vanishing tunnelling probability.
- We live in a **metastable state of the vacuum**.

Markkanena et al., Front. Astron. Space Sci. (2018)

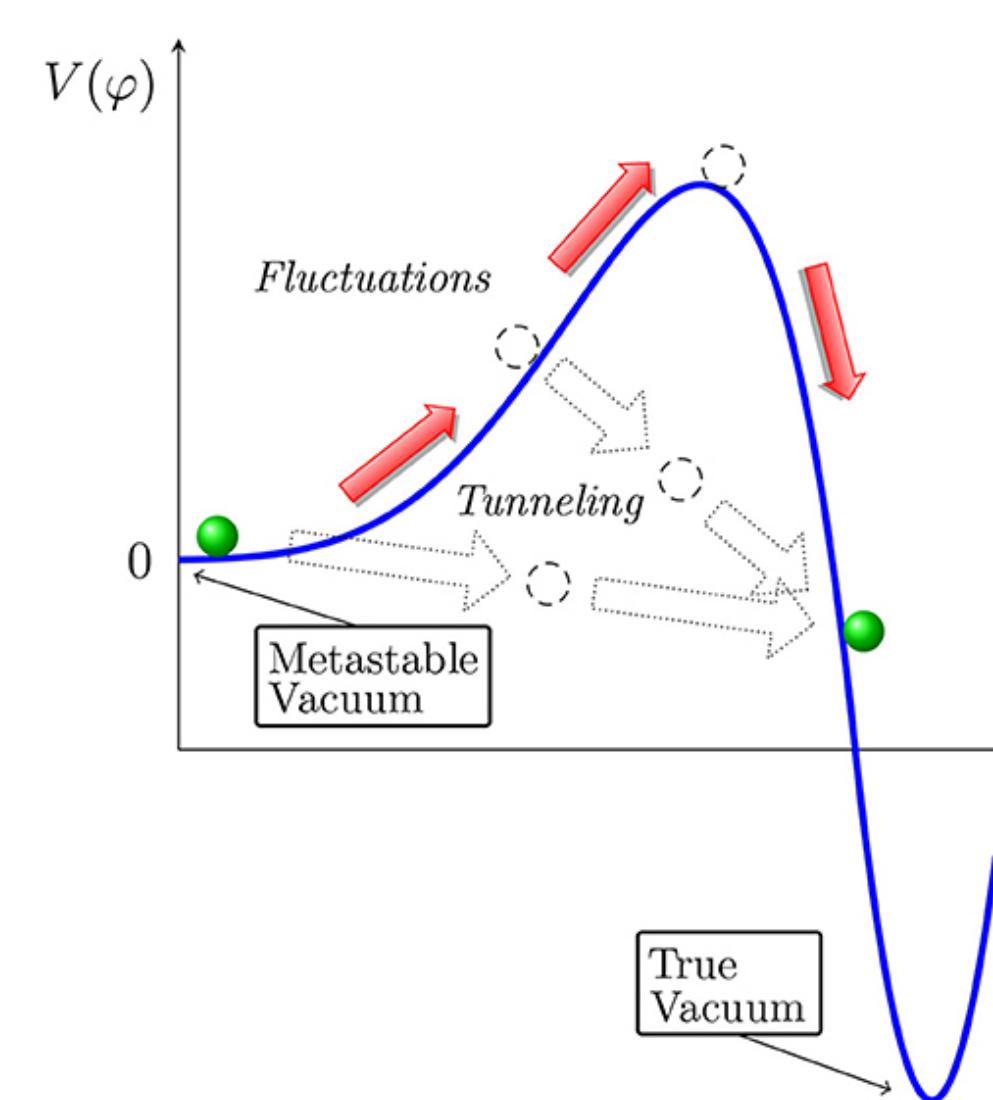
# Triviality of $\phi^4$ -theory

The theory is **well defined** as long as it is bounded from below. When  $\lambda$  becomes negative this is not true anymore.

Two directions to argue against the possibility of facing the end of the universe:



- As soon as the  $\lambda$  reaches its zero value the theory becomes **ill-defined**.
- At this point is reached, the theory becomes **non interacting** (e.g. a free theory) and therefore is **trivialised**.
- Theories beyond the standard should prevent  $\lambda < 0$  while  $\lambda = 0$  is a measure zero point



Markkanena et al., Front. Astron. Space Sci. (2018)

- There should be **physics beyond the standard model** that prevents the potential to go to -infinity (e.g. unknown particles.)
- Even though we might tunnel through the metastable state there's a **true vacuum** at a **lower energy** introduced by this new physics.

# Annealed Importance Sampling

---

Similar lines of work leveraging Annealed Importance Sampling

- SNF (Wu et al.)
  - CRAFT (Matthews et al.)
  - **FAB (Midgley et al.)**
- } Flows + stochastic sampling steps

$$D_\alpha(p\|q) = -\frac{\int_x p(\mathbf{x})^\alpha q(\mathbf{x})^{1-\alpha} d\mathbf{x}}{\alpha(1-\alpha)}$$

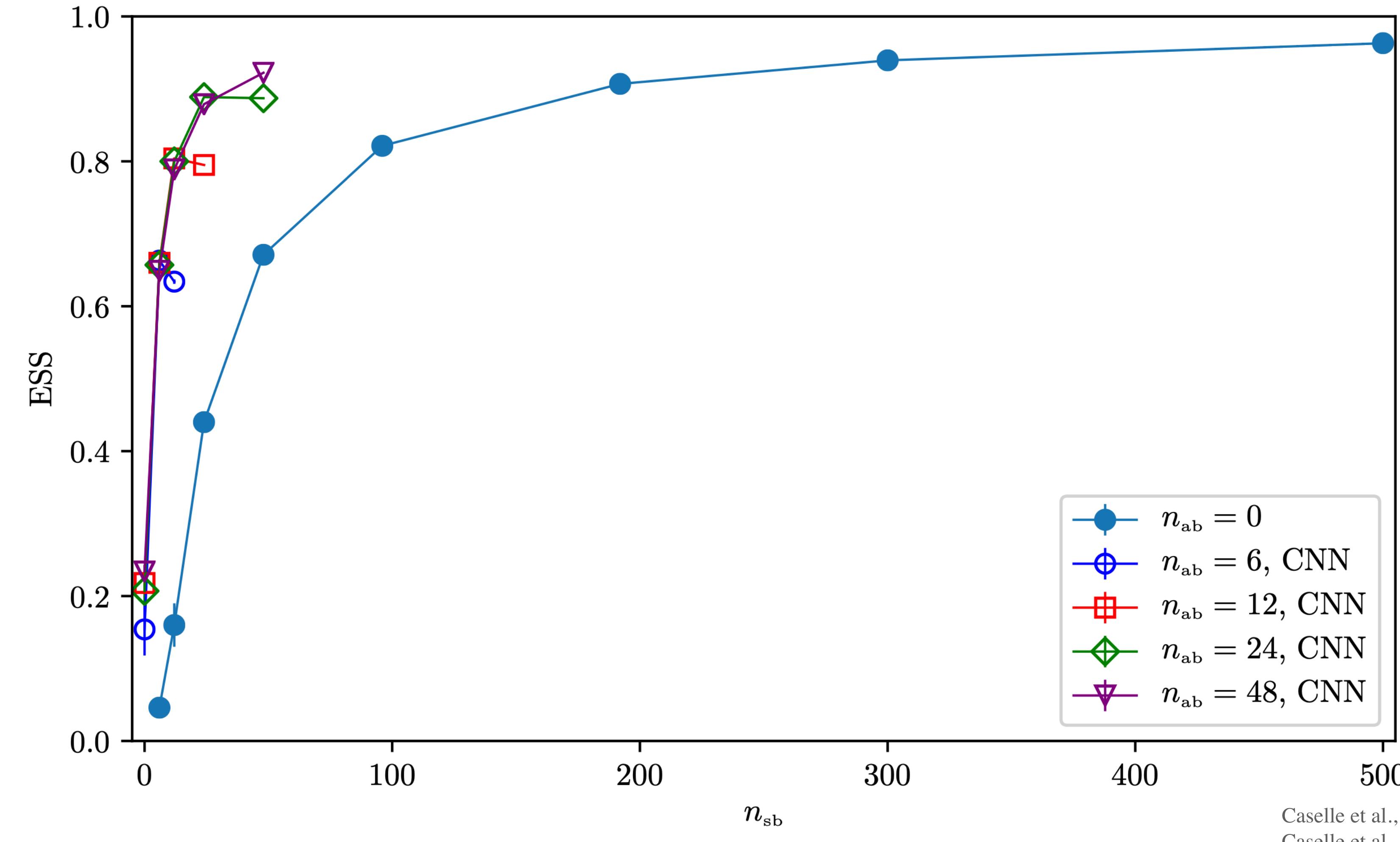
Minimising  $\alpha$ -divergence for  $\alpha = 2$  corresponds to  
minimise the variance of importance weights



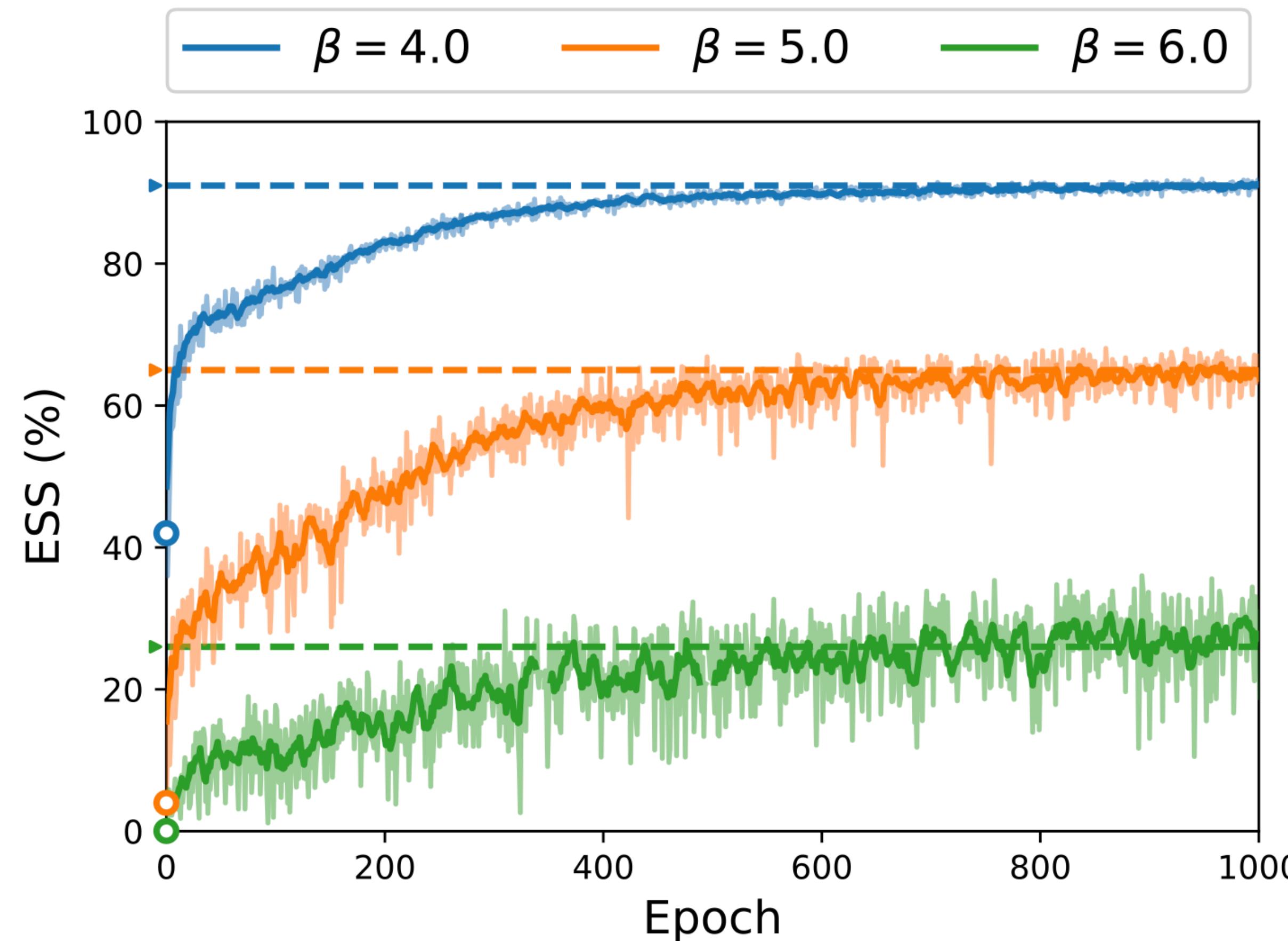
$$D_{\alpha=2}(p\|q_\theta) \propto \mathcal{L}(\theta) = \int \frac{p(\mathbf{x})^2}{q_\theta(\mathbf{x})} d\mathbf{x} = \mathbb{E}_{g(\mathbf{x})} \left[ \frac{p(\mathbf{x})^2}{q_\theta(\mathbf{x})g(\mathbf{x})} \right]$$

↓  
Target  $\propto p^2/q_\theta$

# SNFs for Lattice Field Theory



# Learning Trivialising Maps



Ref.	$N_{\text{params}}$	ESS at $\beta$		
		4.0	5.0	6.0
Lüscher, NL [3]	8 non-zero values	42%	4%	<1%
<b>This work</b>	A	$14 \equiv 2_t \times 7_W$	91%	65%
	B	$420 \equiv 10_t \times 42_W$	98%	88%
Boyda et al. [8]	$\mathcal{O}(10^6)$ estimated	88%	75%	48%

[A, B] Bacchio et al., arXiv. (2022)  
[8] Boyda et al., Phys. Rev. D (2021)

# Incorporating Symmetries

---

- Permutations

$$p(\dots, x_i, \dots, x_j, \dots) = p(\dots, x_j, \dots, x_i, \dots) \quad \text{Symmetric permutation}$$

$$p(\dots, x_i, \dots, x_j, \dots) = -p(\dots, x_j, \dots, x_i, \dots) \quad \text{Antisymmetric permutation}$$

- Translations, Rotations

Special Euclidean Group SE(3)

Octahedral Symmetries

- Gauge Symmetries

Special Unitary Group SU( $N$ )

Unitary Group U( $N$ )

$$(\Omega \cdot U)_\mu(x) = \Omega(x)U_\mu(x)\Omega^\dagger(x + \hat{\mu})$$

# Incorporating Symmetries

---

Knowing the symmetries of the system one can build a flow which is **invariant or equivariant** under some group action.

## Gauge Equivariant Flow

$$Y_\mu(x) = f(U_\mu(x); \theta)$$

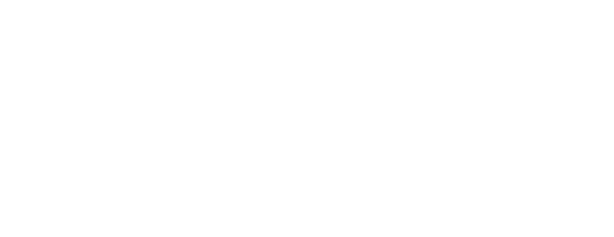
$$U_\mu(x) \rightarrow \Omega(x)U_\mu(x)\Omega(x + \hat{\mu})^\dagger$$

$$Y_\mu(x) \rightarrow \Omega(x)Y_\mu(x)\Omega(x + \hat{\mu})^\dagger$$

Kanwar et al., Phys. Rev. Lett. (2020)  
Boyda et al., Phys. Rev. D (2021)

## $\mathbb{Z}_2$ Invariant Flow

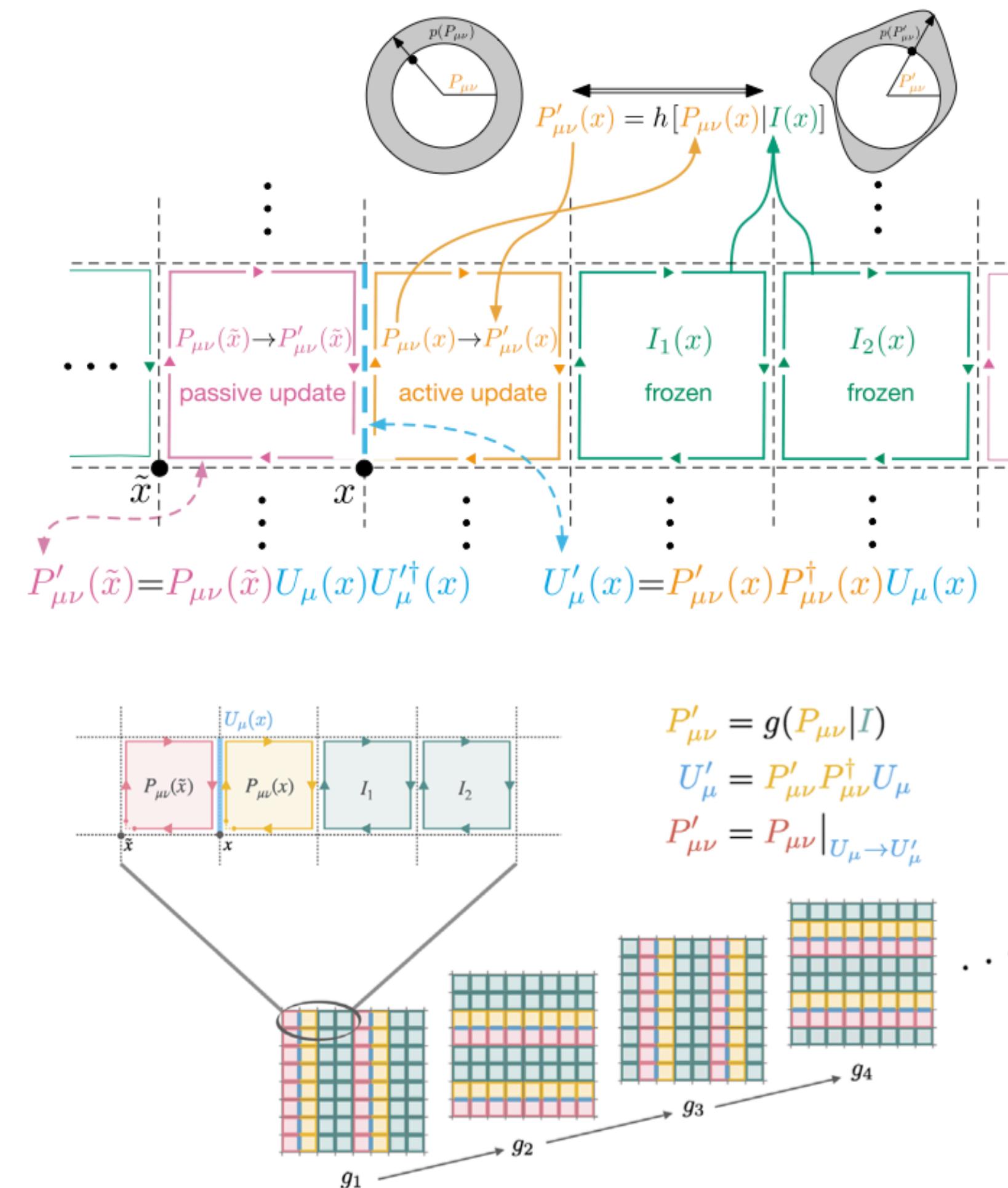
$$S(\phi) = S(-\phi)$$



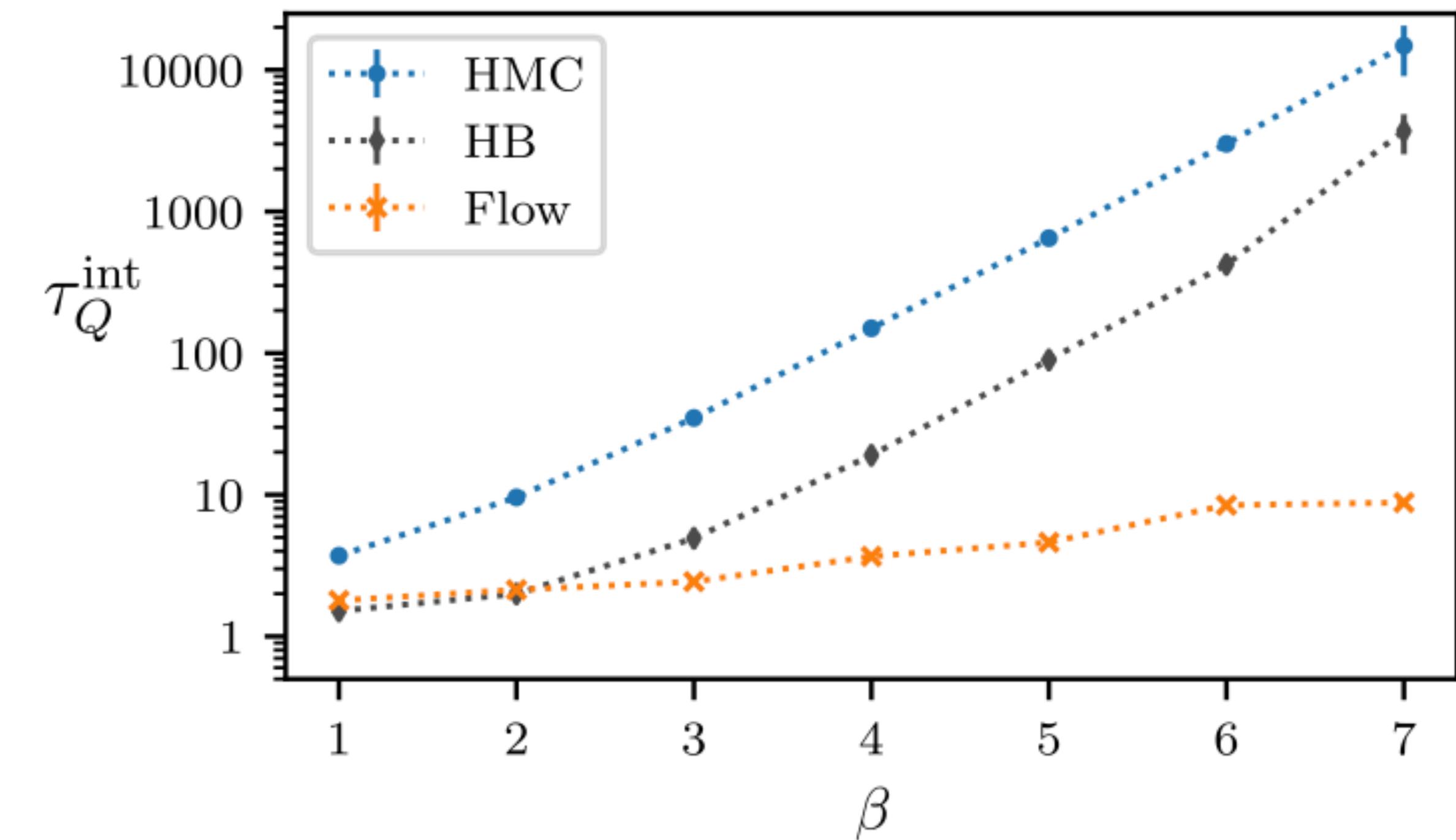
$$S(f(z; \theta)) = S(f(-z; \theta))$$

Nicoli et al., Phys. Rev. Lett. (2021)

# Gauge-Equivariant Flows for LGT

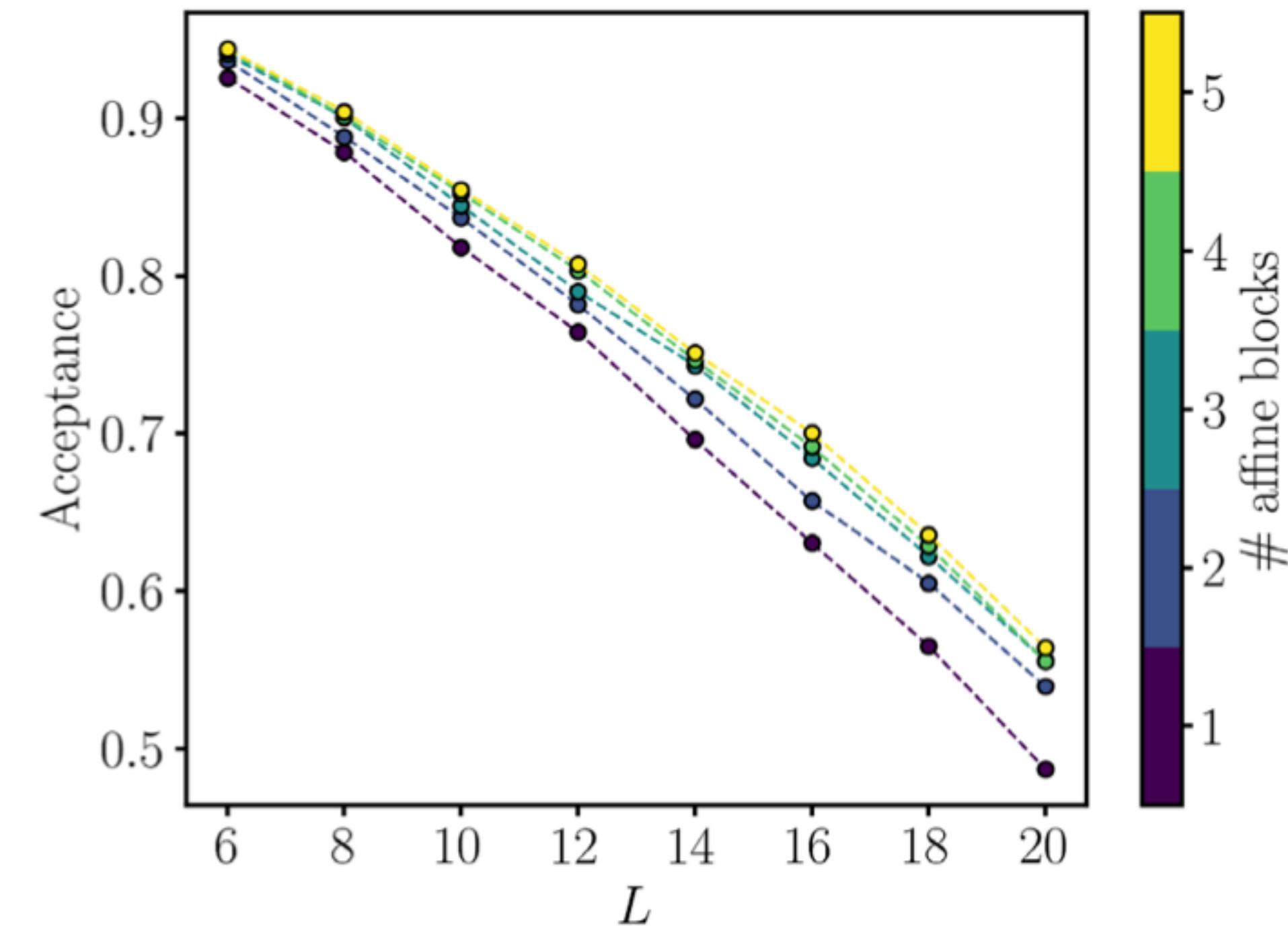


$U(1)$  Gauge Theory on  $16 \times 16$  lattice

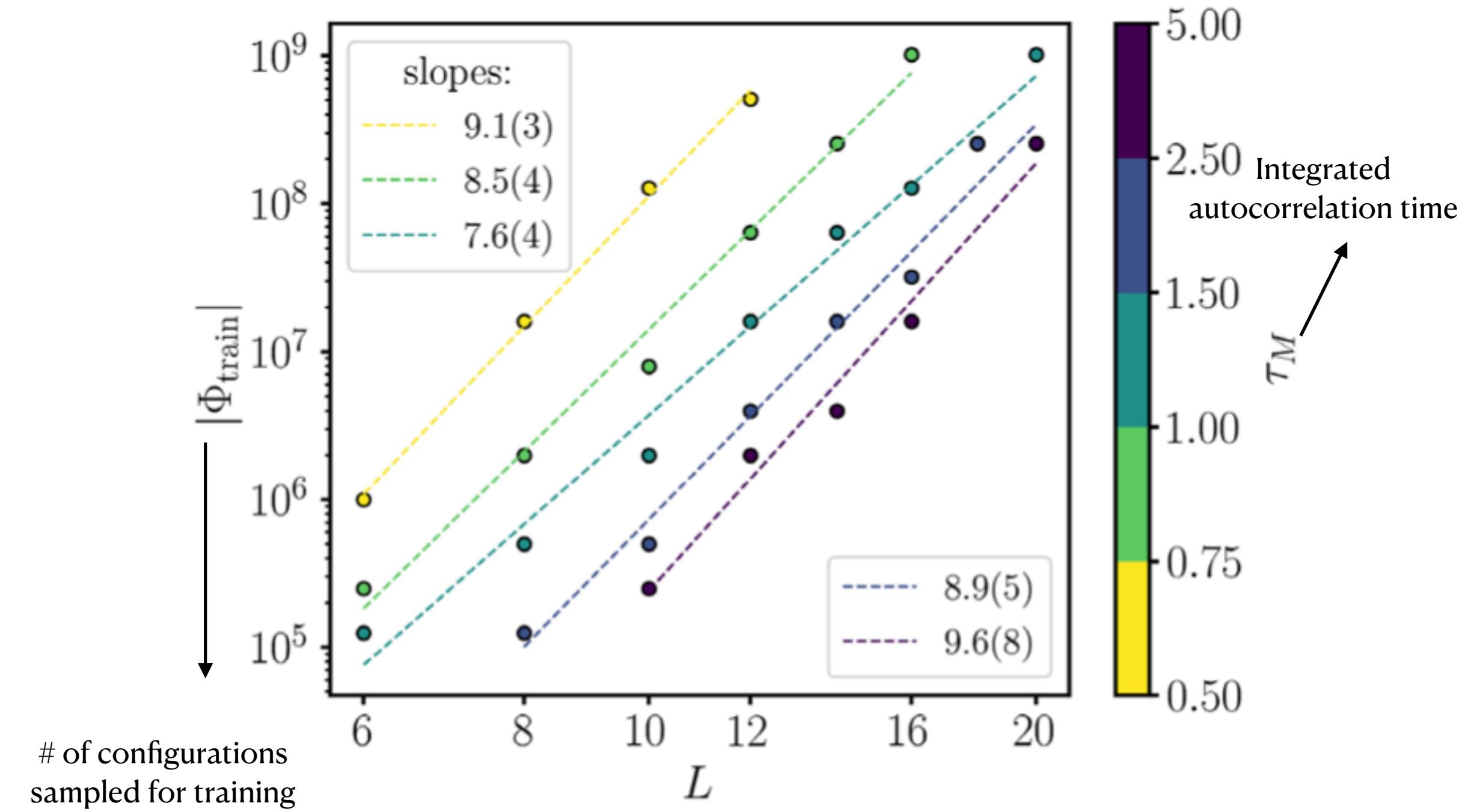


Kanwar et al., Phys. Rev. Lett. (2020)  
Boyda et al., Phys. Rev. D (2021)

# Scaling to Larger Lattices



Acceptance rate drastically decays as the lattice shape increases



Exponential increase in the training cost as the lattice shape increases

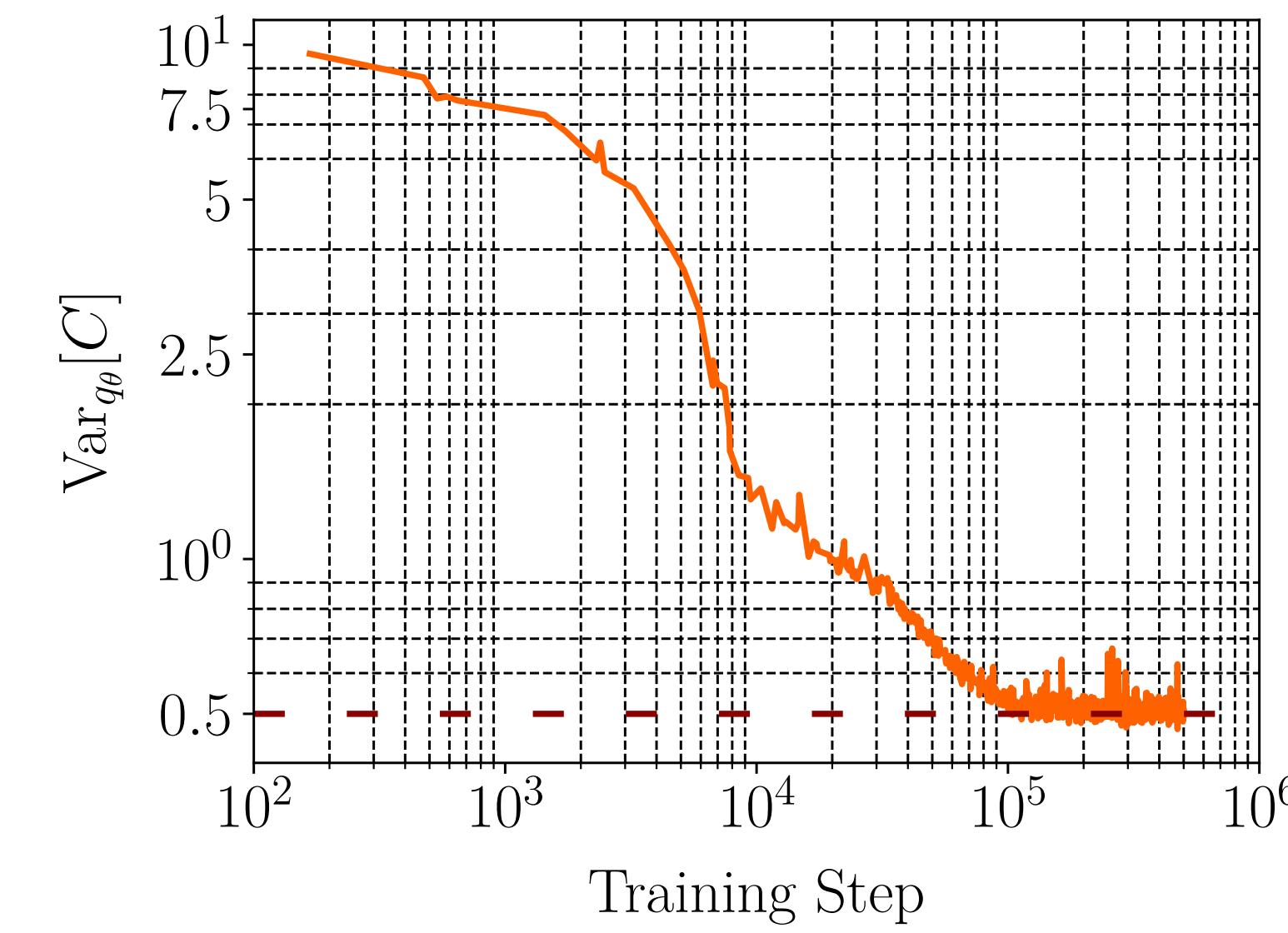
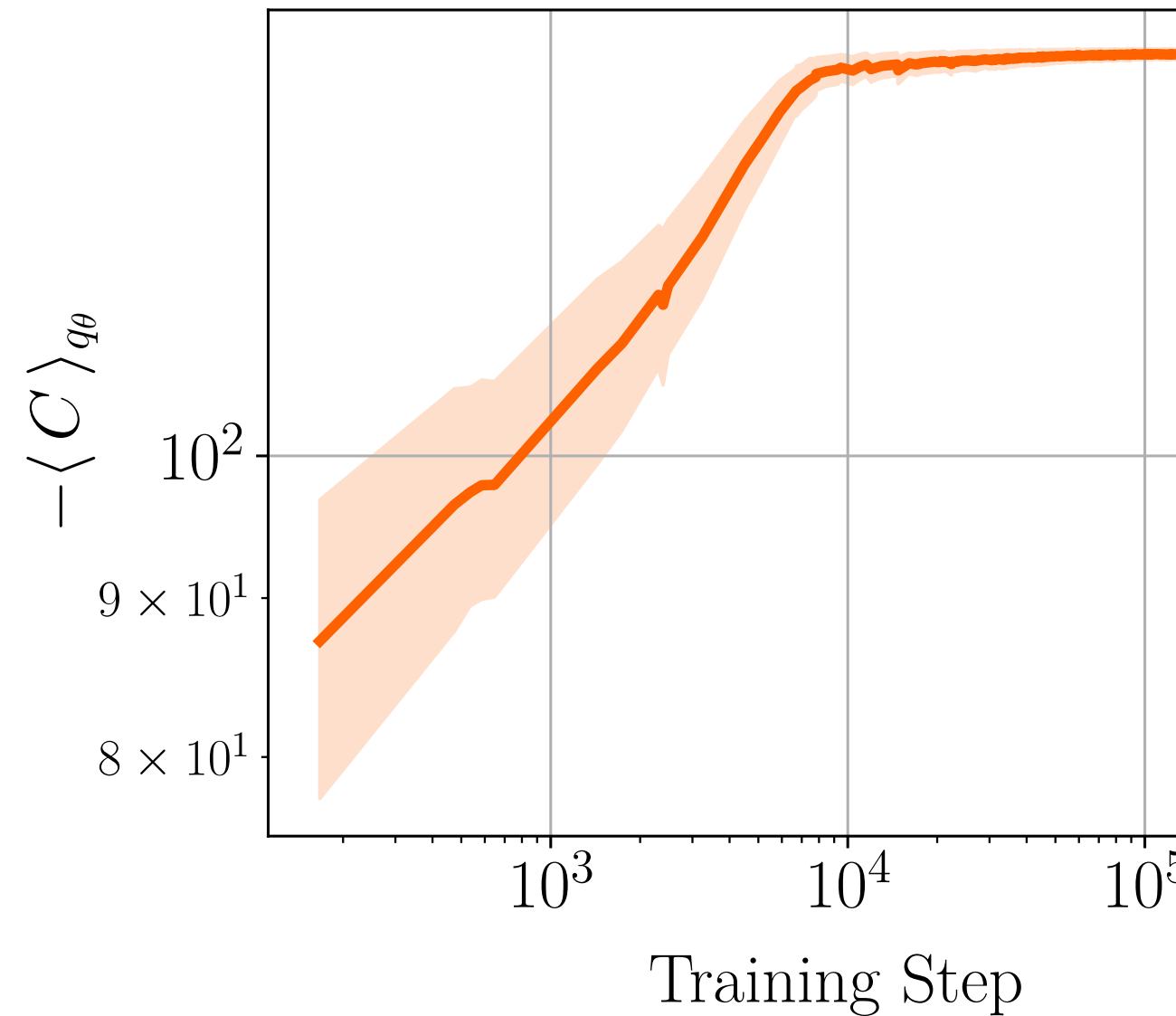
Del Debbio et al., Phys. Rev. D (2021)

# KL is Related to the Variance of Importance Weights

**Theorem** Let  $C(\phi) = S(\phi) + \ln q_\theta(\phi)$ . The following relation between the *KL* divergence and the variance of the random variable  $C$  holds:

$$KL(q_\theta || p) = \frac{1}{2} \text{Var}_{q_\theta} [C] + \mathcal{O} (\mathbb{E}_{q_\theta} [|w - 1|^3]) ,$$

where  $w(\phi) = \frac{p(\phi)}{q_\theta(\phi)}$  is the normalized importance weight.



# Relation to Gibbs Inequality

---

The **Gibbs Inequality** states that given a (discrete) probability distribution  $P = \{p_1, \dots, p_N\}$ , then for any other probability distribution  $Q = \{q_1, \dots, q_N\}$  the following inequality holds

$$-\sum_{i=1}^n p_i \ln p_i \leq -\sum_{i=1}^n p_i \ln q_i$$

Which implies that **the information entropy of a distribution  $P$  is always smaller or equal to its cross-entropy with respect to another distribution  $Q$ .**

Rearranging this inequality

$$+\sum_{i=1}^n p_i \ln p_i - \sum_{i=1}^n p_i \ln q_i \leq 0 \implies \sum_{i=1}^n p_i \ln \frac{p_i}{q_i} = \text{KL}(P || Q) \geq 0$$

Which tells us that the Kullback-Leibler divergence is always **positive!**

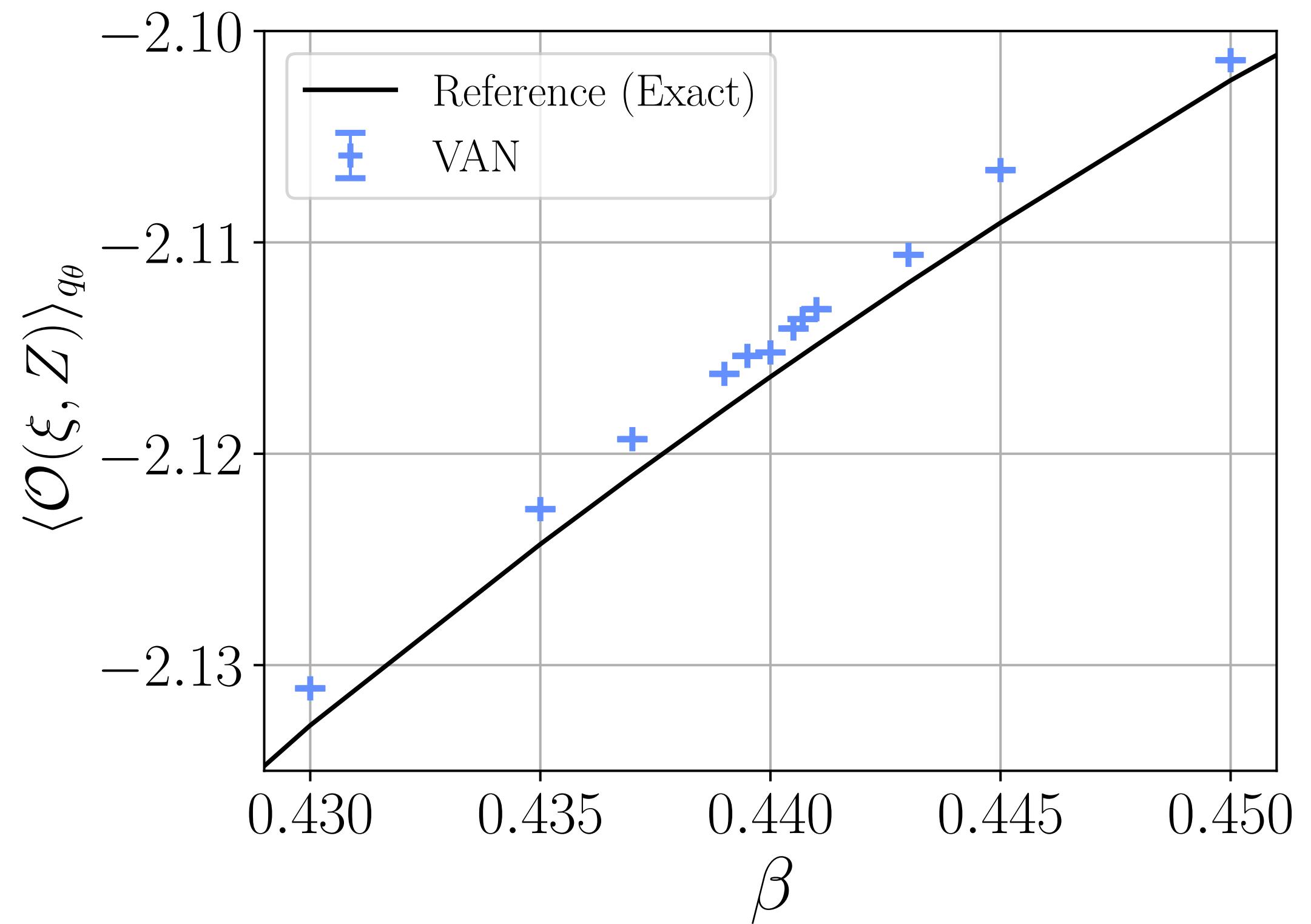
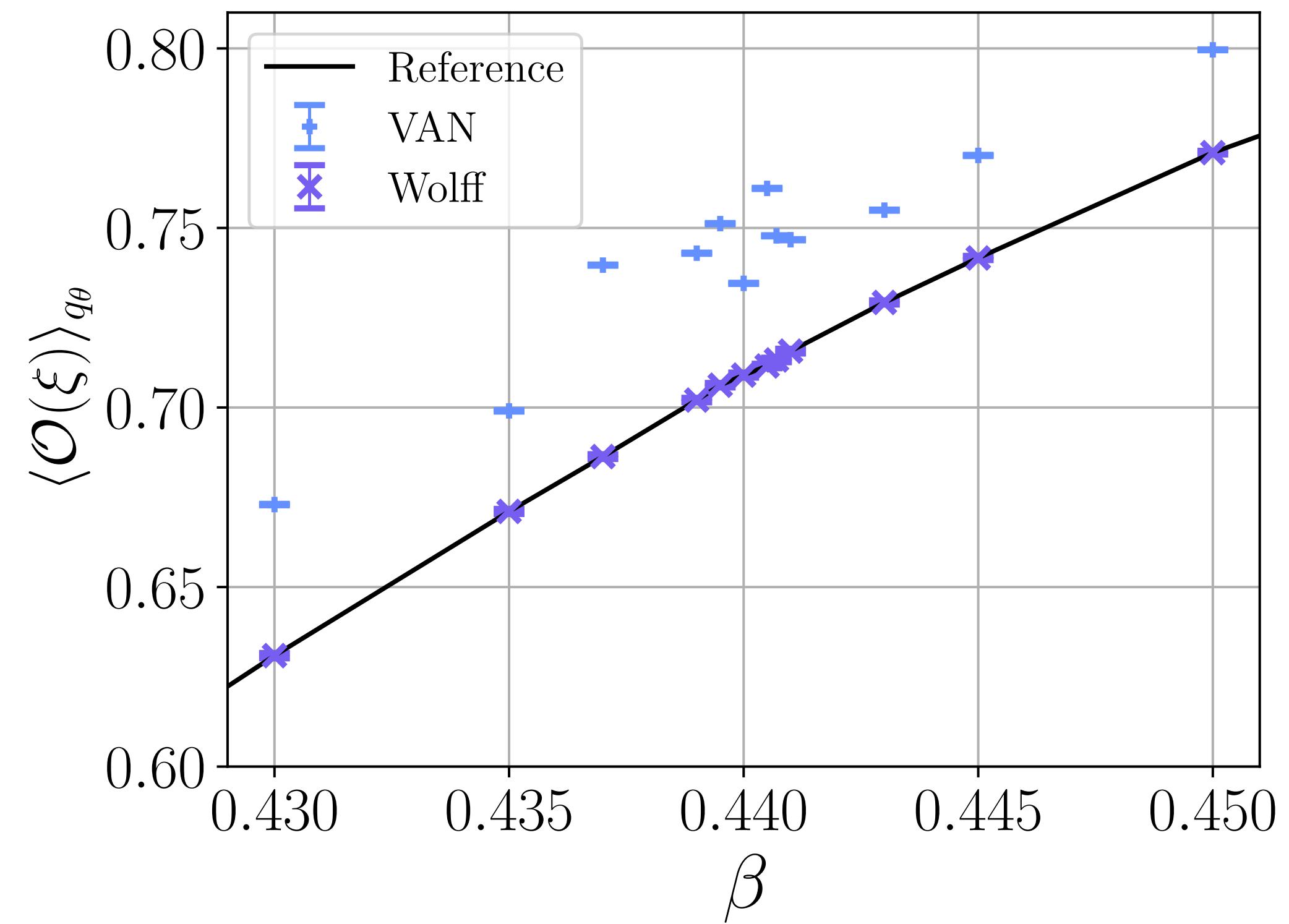
$$\text{KL}(q_\theta || p) = \beta(F_q - F) \geq 0 \implies \text{One can minimise } F_q \text{ instead of } F_q - F$$

# Application: the 2D Ising Model

$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j$$

$$p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}$$

$$\langle \mathcal{O} \rangle_{q_\theta} \neq \langle \mathcal{O} \rangle_p$$



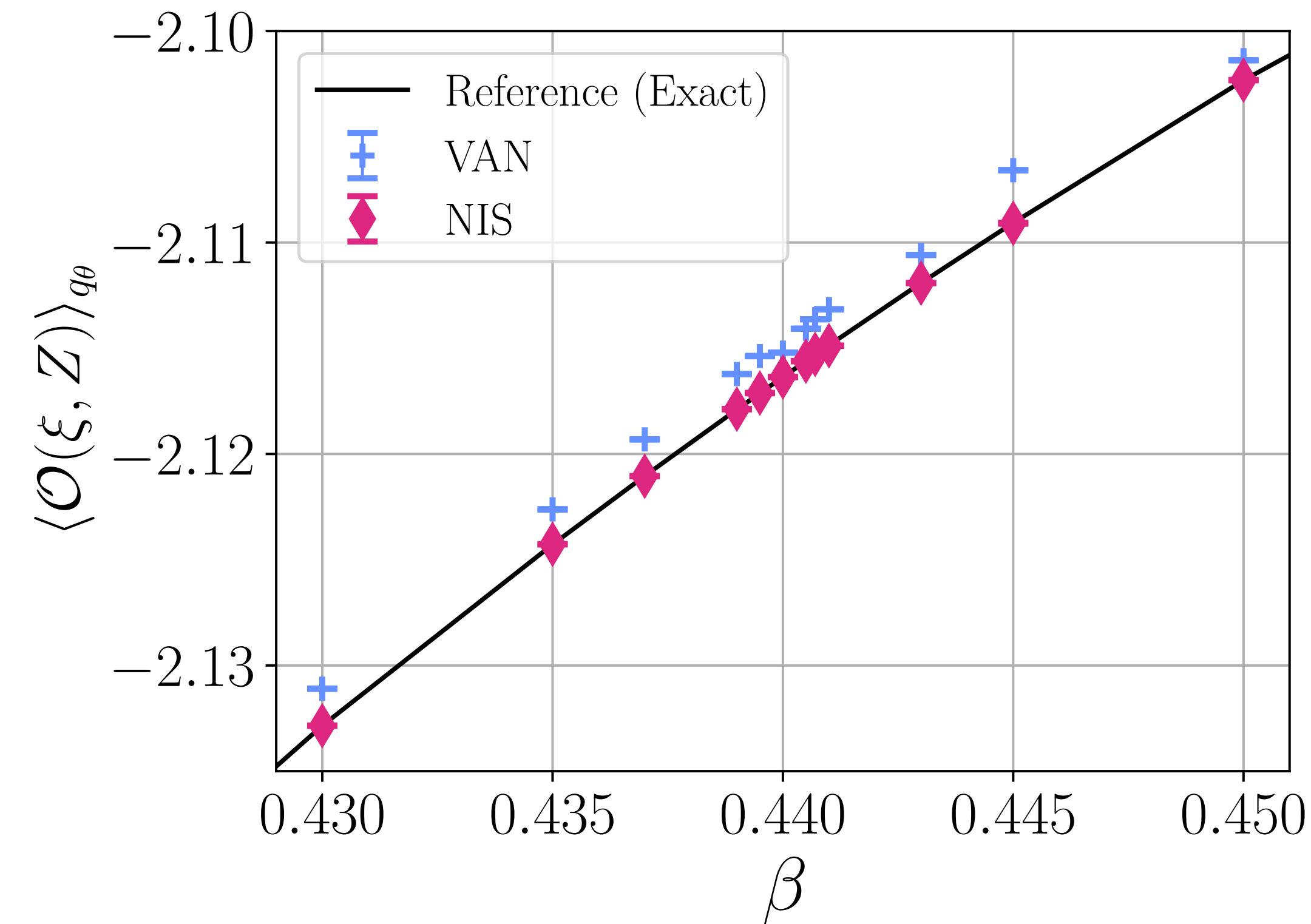
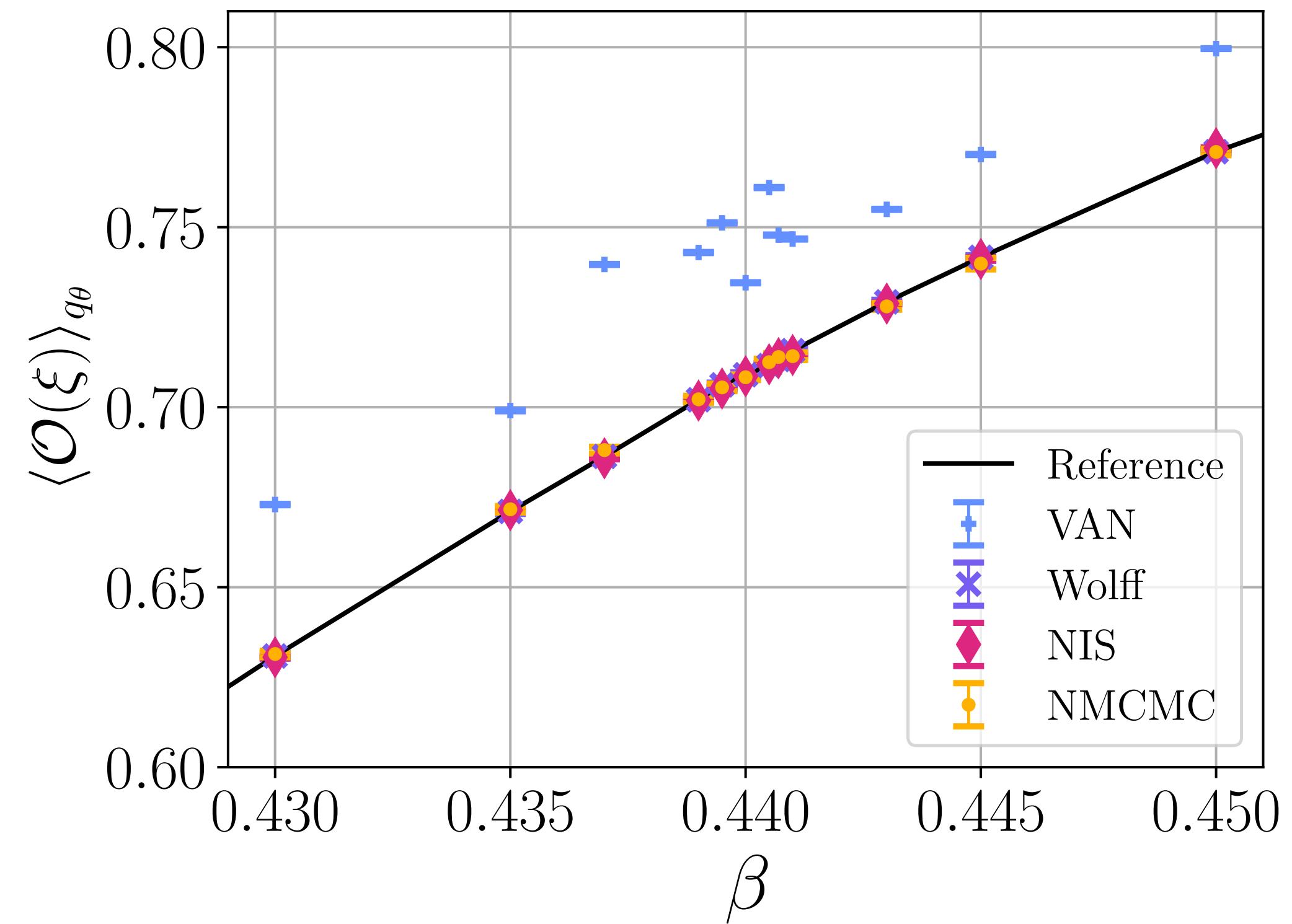
Wu, Wang, Zhang, Phys. Rev. Lett. (2019).  $\implies$  VAN (Variational Autoregressive Network)

# Application: the 2D Ising Model

$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j$$

$$p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}$$

$$\langle \mathcal{O} \rangle_{q_\theta} \longrightarrow \langle w\mathcal{O} \rangle_{q_\theta}$$



KAN, Nakajima, Strodthoff, Samek, Müller, and Kessel, Phys. Rev. E (2020).

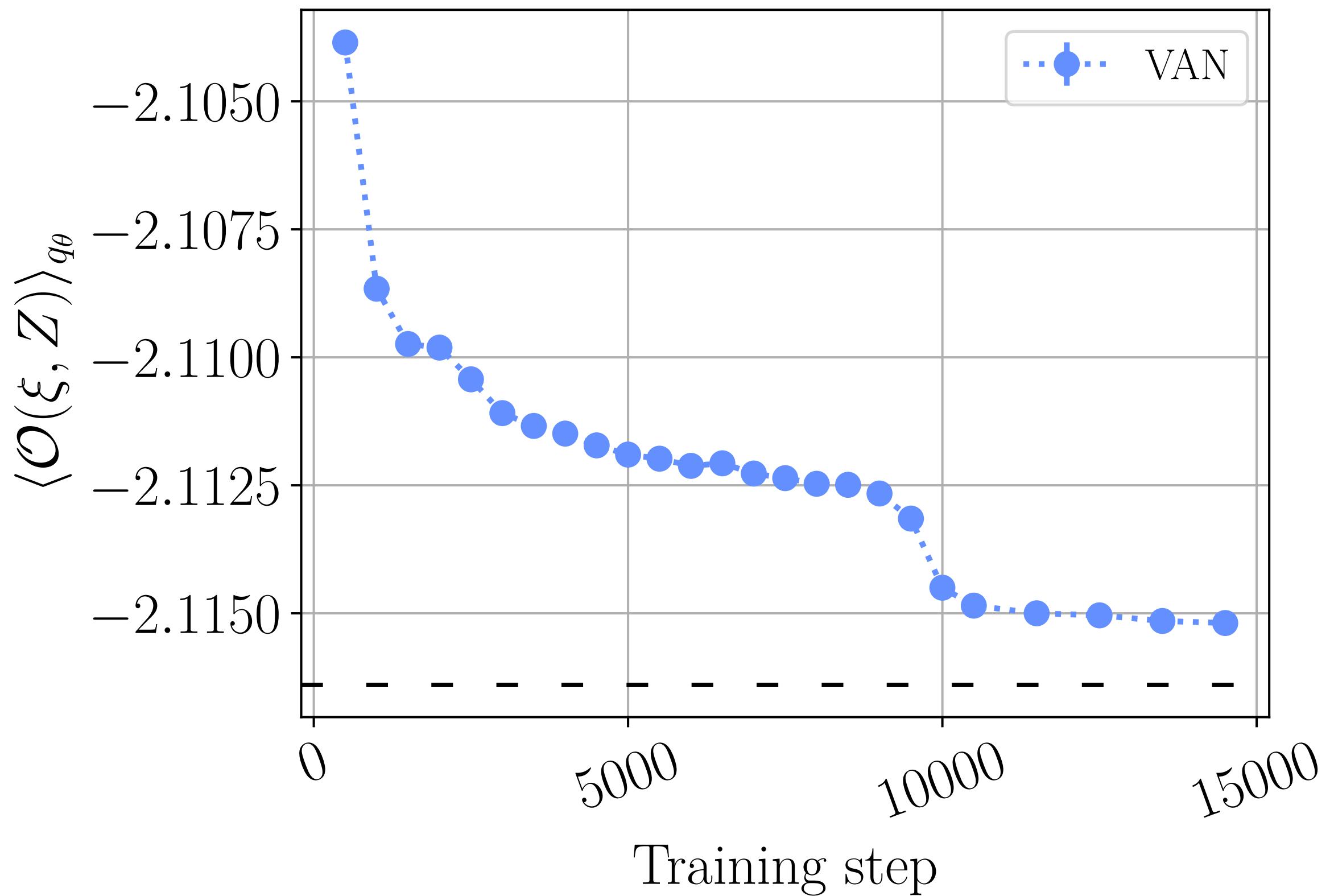
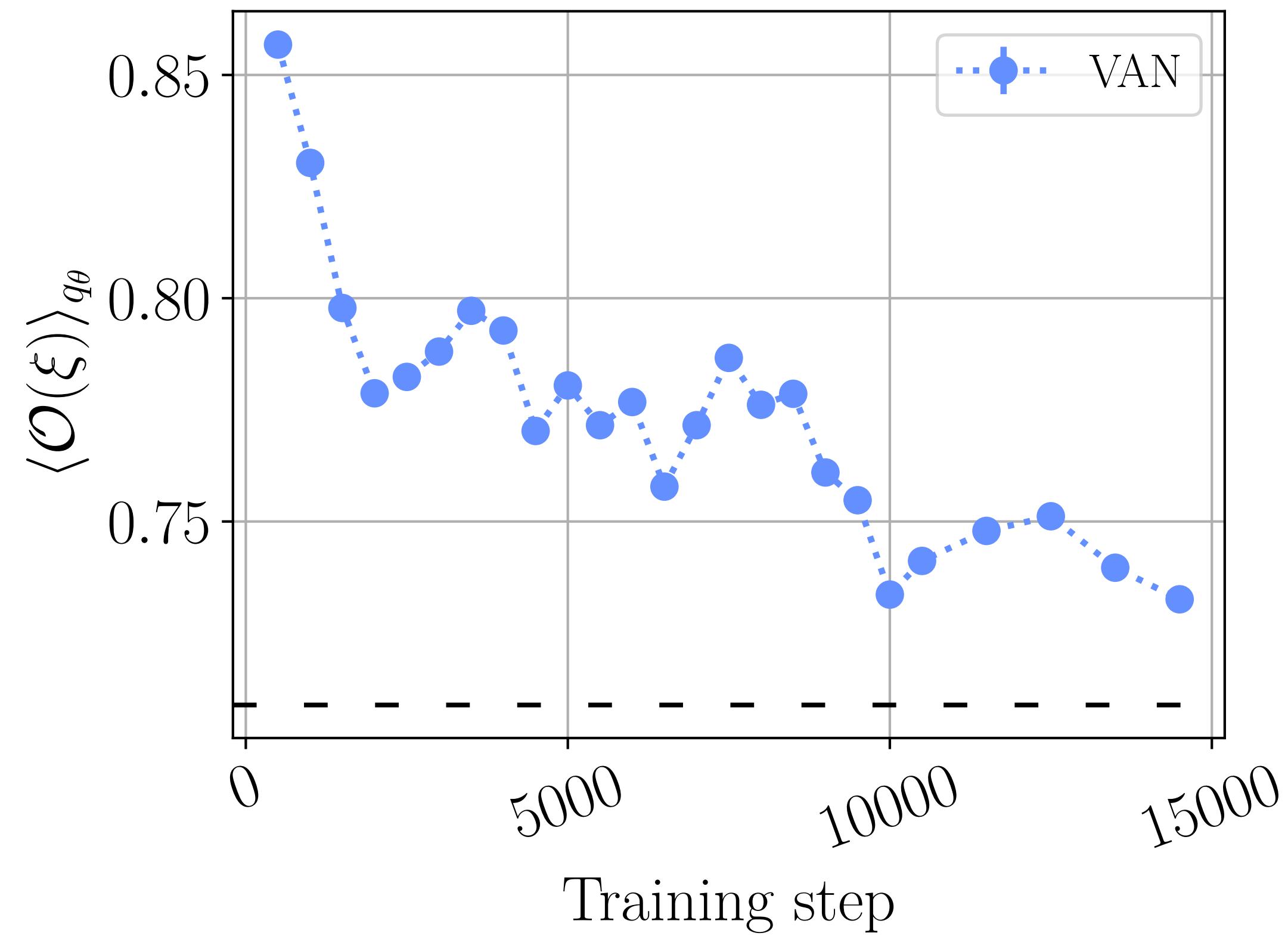
Wu, Wang, Zhang, Phys. Rev. Lett. (2019).

# Application: the 2D Ising Model

$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j$$

$$p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}$$

$$\langle \mathcal{O} \rangle_{q_\theta} \neq \langle \mathcal{O} \rangle_p$$



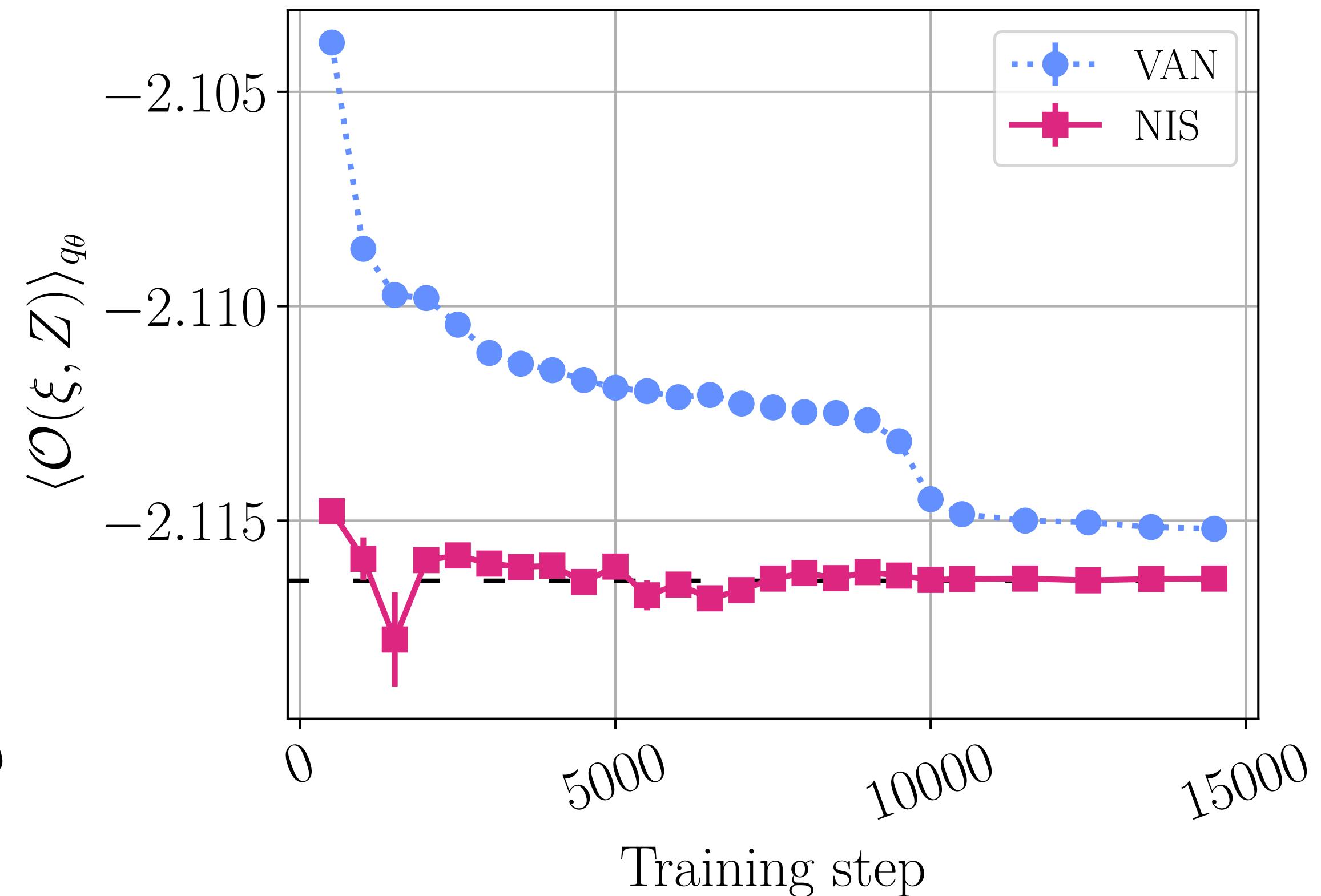
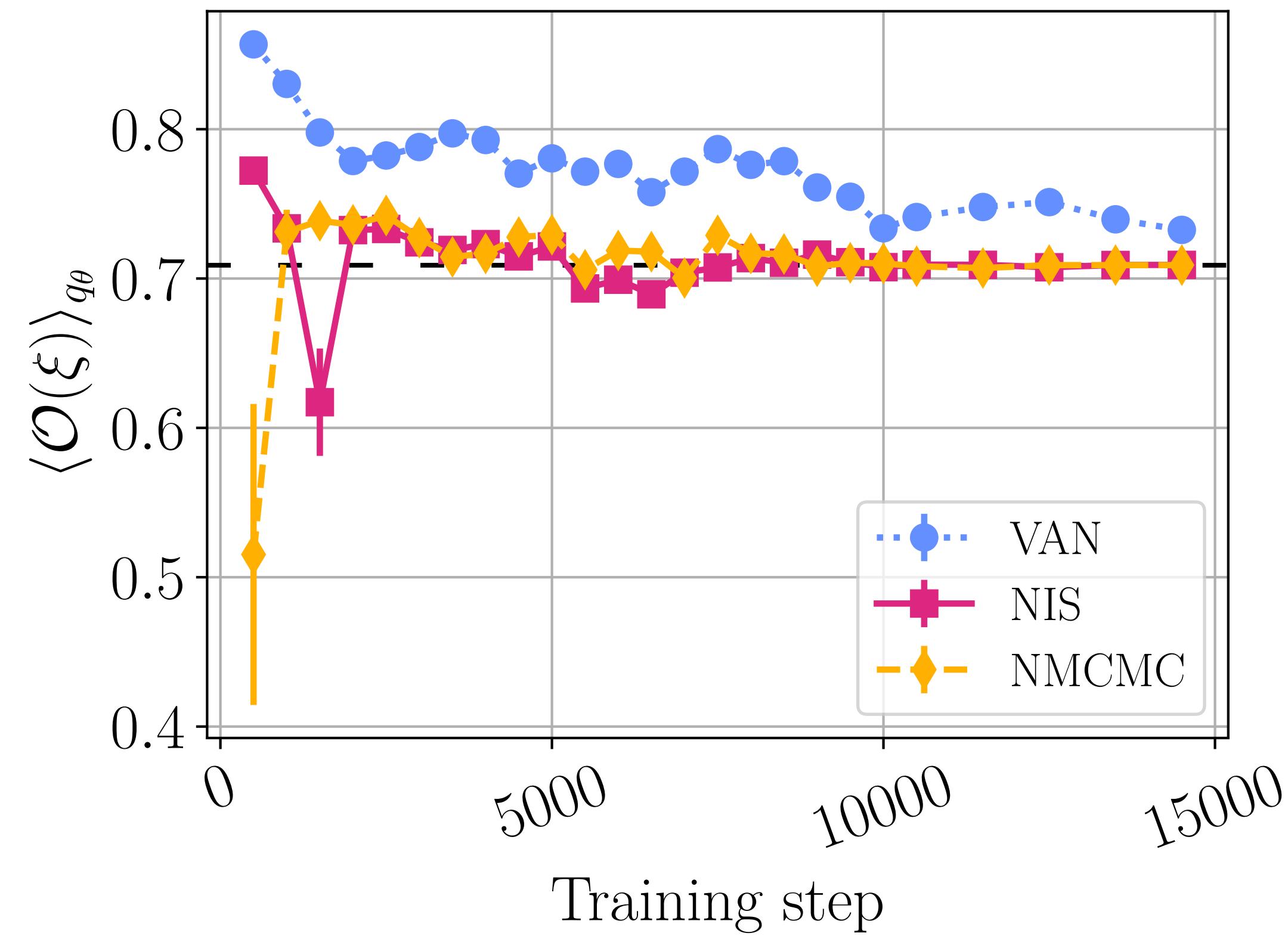
Wu, Wang, Zhang, Phys. Rev. Lett. (2019).

# Application: the 2D Ising Model

$$H(\xi) = -J \sum_{\langle i,j \rangle} \xi_i \xi_j$$

$$p(\xi) = \frac{e^{-\beta H(\xi)}}{Z}$$

$$\langle \mathcal{O} \rangle_{q_\theta} \longrightarrow \langle w\mathcal{O} \rangle_{q_\theta}$$



KAN, Nakajima, Strodthoff, Samek, Müller, and Kessel, Phys. Rev. E (2020).

Wu, Wang, Zhang, Phys. Rev. Lett. (2019).

# Delta Method of Moments

---

**Theorem 2.** Let  $\hat{X}_N = \frac{1}{N} \sum_{i=1}^N X_i$  be the sample mean of independent and identically distributed random variables  $X_i$  with  $\mathbb{E}[X_i^{2k+2}] < \infty$  for  $k \in \{0, 1\}$ . Let  $h$  be a real-valued function with uniformly bounded derivatives. It then holds that

$$\mathbb{E}[h(\hat{X}_N)] = c_0 + \frac{c_1}{N} + O\left(\frac{1}{N^2}\right), \quad (4.10)$$

where

$$c_0 = h(\mu), \quad c_1 = h''(\mu) \frac{\sigma^2}{2}, \quad (4.11)$$

with  $\sigma^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$  and  $\mu = \mathbb{E}[X]$ .

There's a subtlety here. The delta method requires the derivatives of the function  $h$ . For a generic LQFT, this will not be the case for  $h(x) = \ln(x)$  since its derivatives diverge for  $x \rightarrow 0^+$ . To overcome this problem one may require that the action of the considered lattice quantum field theory is bounded. For example, this can be ensured by putting the field theory in a box potential. Since only very high energy configurations (very unlikely to be sampled) are affected by this modification, there's no practical effect on the numerical experiments.

# Bias and Variance of the Free Energy Estimator

---

**Theorem 3.** *The bias of the estimator  $\hat{F}$  is given by*

$$\text{Bias}[-\beta\hat{F}] = -\frac{1}{2N} \frac{\mathbb{E}_{q_\theta}[(\tilde{w} - \mathbb{E}_{q_\theta}[\tilde{w}])^2]}{(\mathbb{E}_{q_\theta}[\tilde{w}])^2} + O(N^{-2}),$$

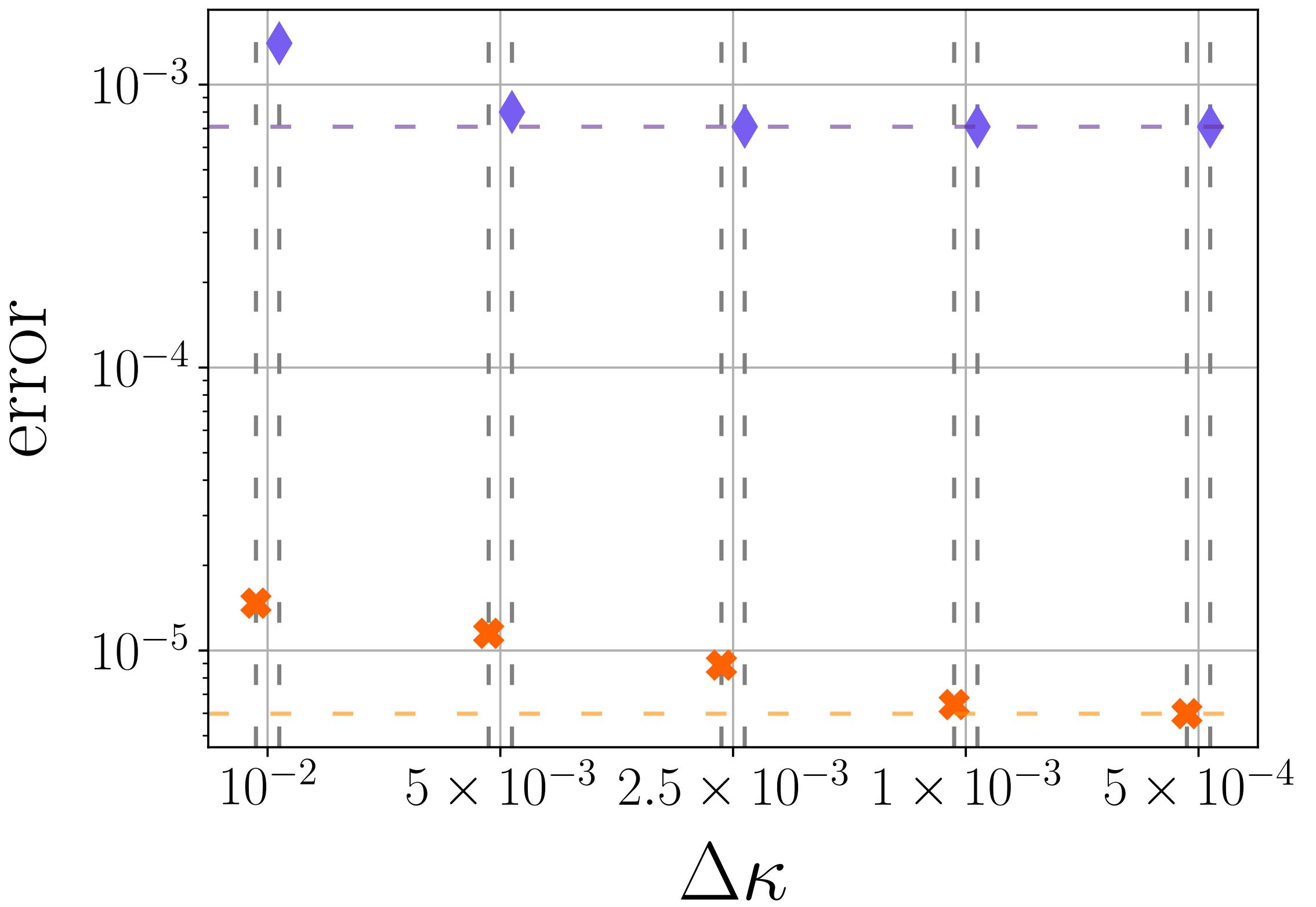
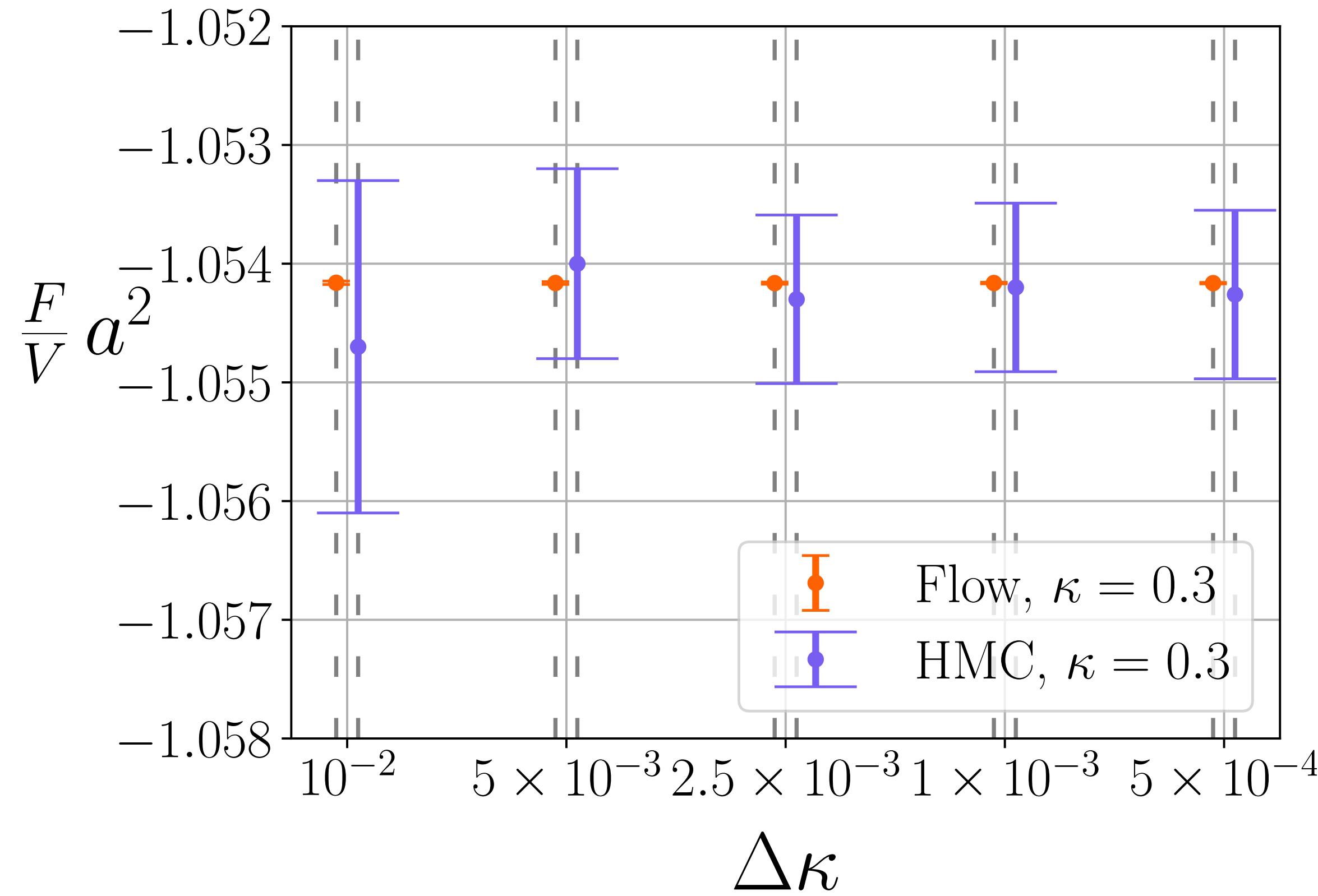
*assuming that  $\mathbb{E}_q[\tilde{w}^{2k+2}] < \infty$  for  $k \in \{0, 1\}$ .*

**Theorem 4.** *The variance of the estimator  $\hat{F}$  is given by*

$$\text{Var}[-\beta\hat{F}] = \frac{1}{N} \frac{\mathbb{E}_{q_\theta}[(\tilde{w} - \mathbb{E}_{q_\theta}[\tilde{w}])^2]}{(\mathbb{E}_{q_\theta}[\tilde{w}])^2} + O(N^{-2}),$$

*assuming that  $\mathbb{E}_q[\tilde{w}^{2k+2}] < \infty$  for  $k \in \{0, 1\}$ .*

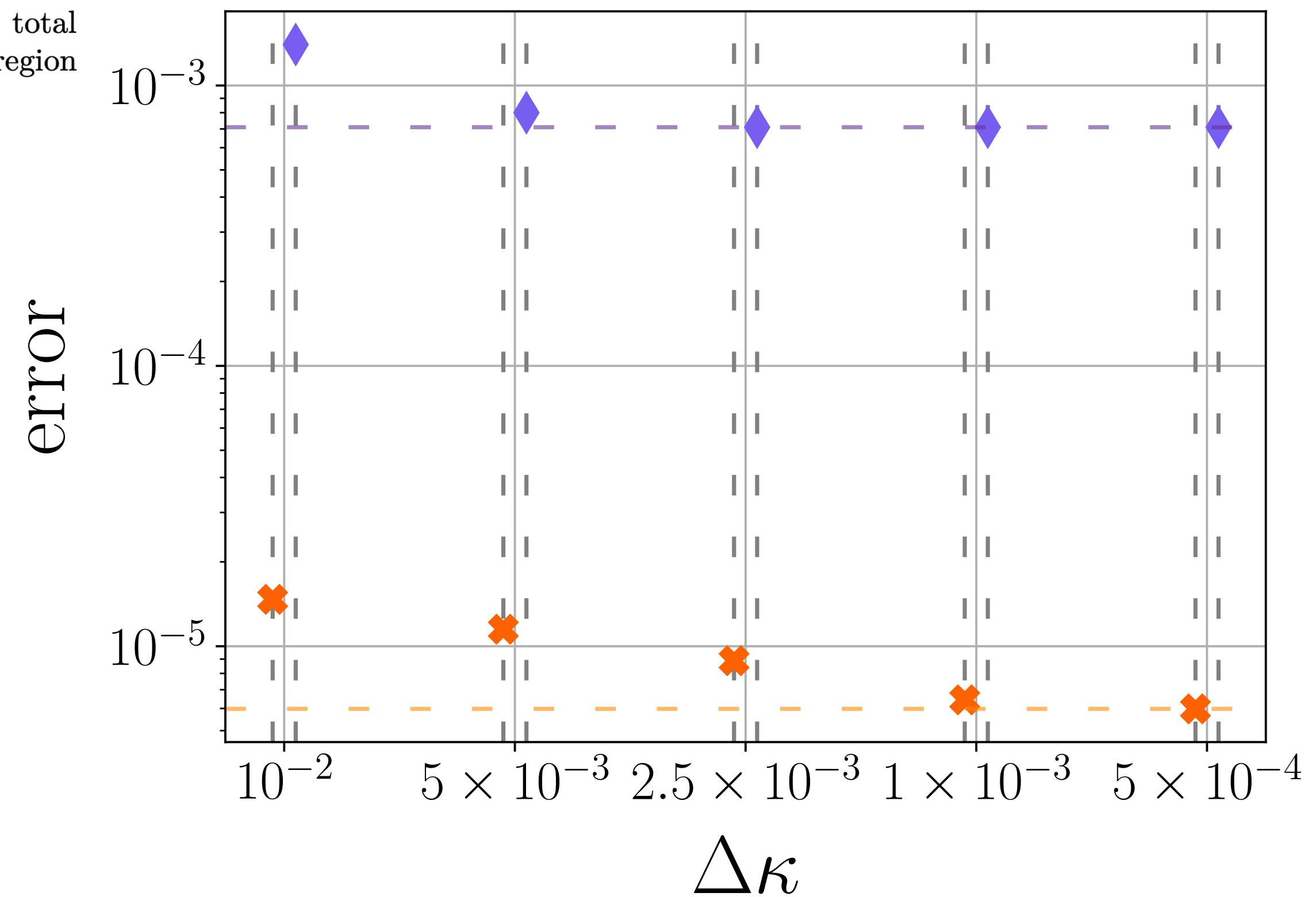
# Error Scaling of Free Energy



# Error Scaling of Free Energy

Table 4.1. **Details on the refinement analysis for smaller  $\Delta\kappa$ :** in each refinement stage, we take steps  $\Delta\kappa$  in a certain sub-region of the hopping parameter  $\kappa$  trajectory (see the last column). The focus is around the critical regime where the biggest error contribution for HMC, see fig. 4.5, is caused by critical slowing down phenomena. The step size taken in this region is shown in the first column. Outside this critical region, the same step sizes for all refinements are taken. Therefore, the total number of samples (third column) grows inversely proportional to the step size of the critical region (first column). Adapted from Nicoli *et al.* (2021).

$\Delta\kappa$	# chains	# samples	$\kappa$ refined region
0.01	14	5.6 M	0.20-0.30
0.005	24	9.6 M	0.20-0.30
0.0025	40	16 M	0.22-0.30
0.001	76	30.4 M	0.24-0.30
0.0005	88	35.2 M	0.267-0.279



# The Effective Sampling Size

---

$$N_{\text{eff}} = \frac{N}{1 + \text{Var}_{q_\theta} [w]}$$

Since

$$\text{Var}_{q_\theta} [w] = \mathbb{E}_{q_\theta} [w^2] - \mathbb{E}_{q_\theta} [w]^2 = \mathbb{E}_{q_\theta} [w^2] - 1$$

one recovers the definition of effective sample ratio (**ESR**)

$$N_{\text{eff}} = \frac{N}{\mathbb{E}_{q_\theta} [w^2]} \implies \text{ESR} = \frac{N_{\text{eff}}}{N} = \frac{1}{\mathbb{E}_{q_\theta} [w^2]}.$$

We refer to [Kong \(1992\)](#) for a mathematical derivation of the **ESS** from the variance of the estimator.

# Is Reweighting All You Need?

---

- Our estimator is **consistent** which implies the bias scales with  $O(N^{-1})$
- This might have **potentially large** (or infinite) **variance** depending on the **goodness** of our variational approximation  $q$ .
- We can lower the bias by taking  $N$  very large **BUT** our estimator might still have prohibitively high variance if  $q$  is bad.

$$\text{Bias}[\hat{Z}] \approx O(N^{-1})$$

$$\text{Std}[\hat{Z}] \approx O(1/\sqrt{N}) \longrightarrow \text{Dominant term.}$$

- The variance of the importance weights, also relates to the effective sampling ratio (ESR)

$$\text{ESR} = \frac{N_{eff}}{N} = \frac{1}{\mathbb{E}_q[w^2]} \text{ relates to } \text{ESR} = \frac{N_{eff}}{N} = \frac{1}{2\tau_{int,\mathcal{O}}}$$

**the goodness of the sampler relates to the degree of iid sampling from the generative model**

---

# Example: the Delta Function

$Z = \int dx \tilde{p}(x) = 1 \rightarrow$  where  $\tilde{p}(x)$  is the Delta function centred in  $x = 0$

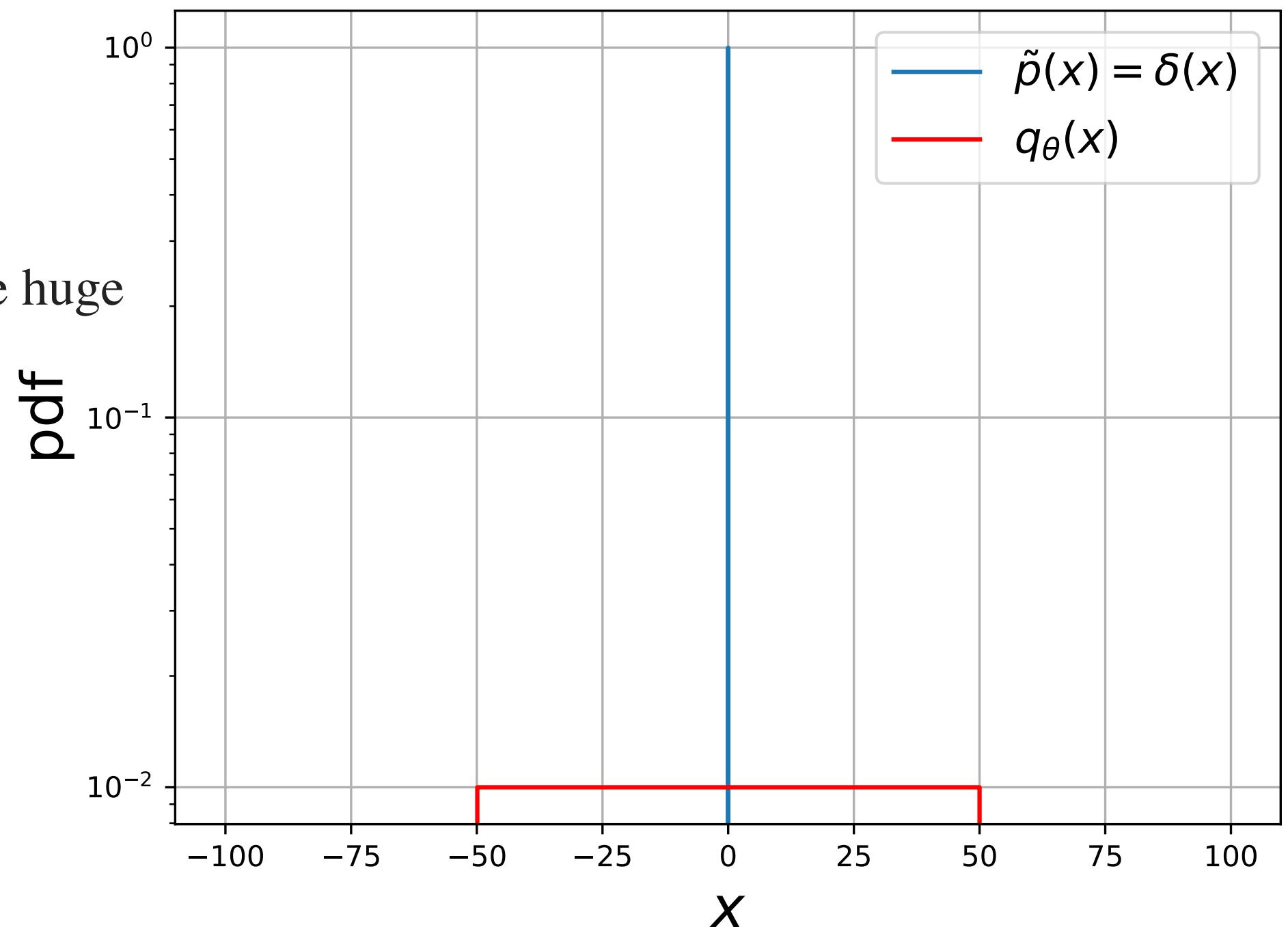
$$Z \approx N^{-1} \sum_n \tilde{p}(x_n) \quad x_n \sim q_\theta = \text{Uniform}(-a, a)$$

$x_n$  rarely hits the low energy valley of  $\tilde{p}(x) \implies$  for finite  $N$ , the bias and variance are huge

$$Z \approx \hat{Z} = N^{-1} \sum_n \tilde{p}(x_n) \approx 0 \quad p(x_n) = 0 \quad \forall x_n \sim q$$

$$\text{Bias}[Z, \hat{Z}] = Z - \mathbb{E}_q[Z] = Z - \hat{Z} \approx Z = 1$$

$$\text{Var}[\hat{Z}] = (\mathbb{E}_q[Z]^2 - \mathbb{E}_q[Z^2]) = 0$$



- sampling **finite** number of samples: we hit the **high probability region** of  $p(x)$  with very **low probability**. Hence  $\tilde{p}(x_n) \approx 0 \rightarrow \hat{Z} \approx 0$
- sampling **Infinite** number of samples: we have a scaling  $O(N^{-1})$ . When  $p(x) = \delta(x)$  **additional care is needed** for mathematical rigorousness around  $x = 0$ .

# Flow-Based and HMC-Based sampling

---

	PRO	CONS
HMC	<ul style="list-style-type: none"><li>• No training required</li><li>• Asymptotically unbiased</li></ul>	<ul style="list-style-type: none"><li>• Sequential Sampling</li><li>• Struggles at criticality</li><li>• Critical Slowing down</li><li>• No direct access to Z</li></ul>
FLOW	<ul style="list-style-type: none"><li>• Parallel sampling</li><li>• Asymptotically unbiased</li><li>• Direct access to Free Energy</li><li>• Point-wise estimates</li><li>• Smaller statistical errors</li><li>• Can leverage inductive biases</li></ul>	<ul style="list-style-type: none"><li>• Up-front expensive training</li><li>• Bias from imperfect training</li></ul>

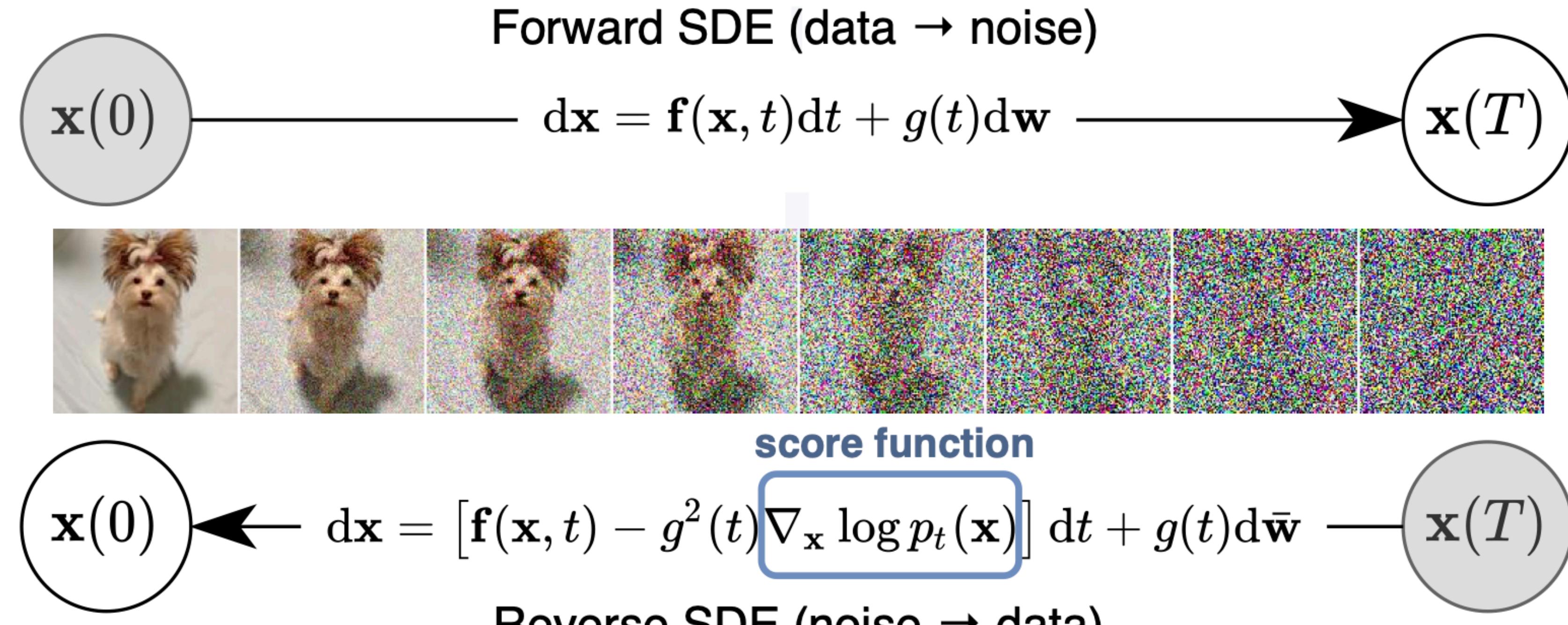
# Generalisation to Other Deep Generative Models

---

Table 3.1. **Applicability of NMCMC and NIS to different families of DGMs:**  $h$  refers to the explicit dependence of the operator  $\mathcal{O}(\phi, Z)$  on the partition function  $Z$ , see eq. (3.25). The observations from the main text are summarized using the following notation about our methods: ✓: applicable and the estimator is *asymptotically unbiased*; (✓): applicable but the estimator is still *biased*; X: not applicable. Adapted from [Nicoli \*et al.\* \(2020\)](#).

DGM	Sampling probability	NMCMC, NIS( $h \equiv 0$ )	NIS( $h \neq 0$ )
GAN	none	X	X
RBM	approximate, unnormalized	(✓)	X
VAE	approximate, normalized	(✓)	(✓)
ARNN, NF	exact, normalized	✓	✓

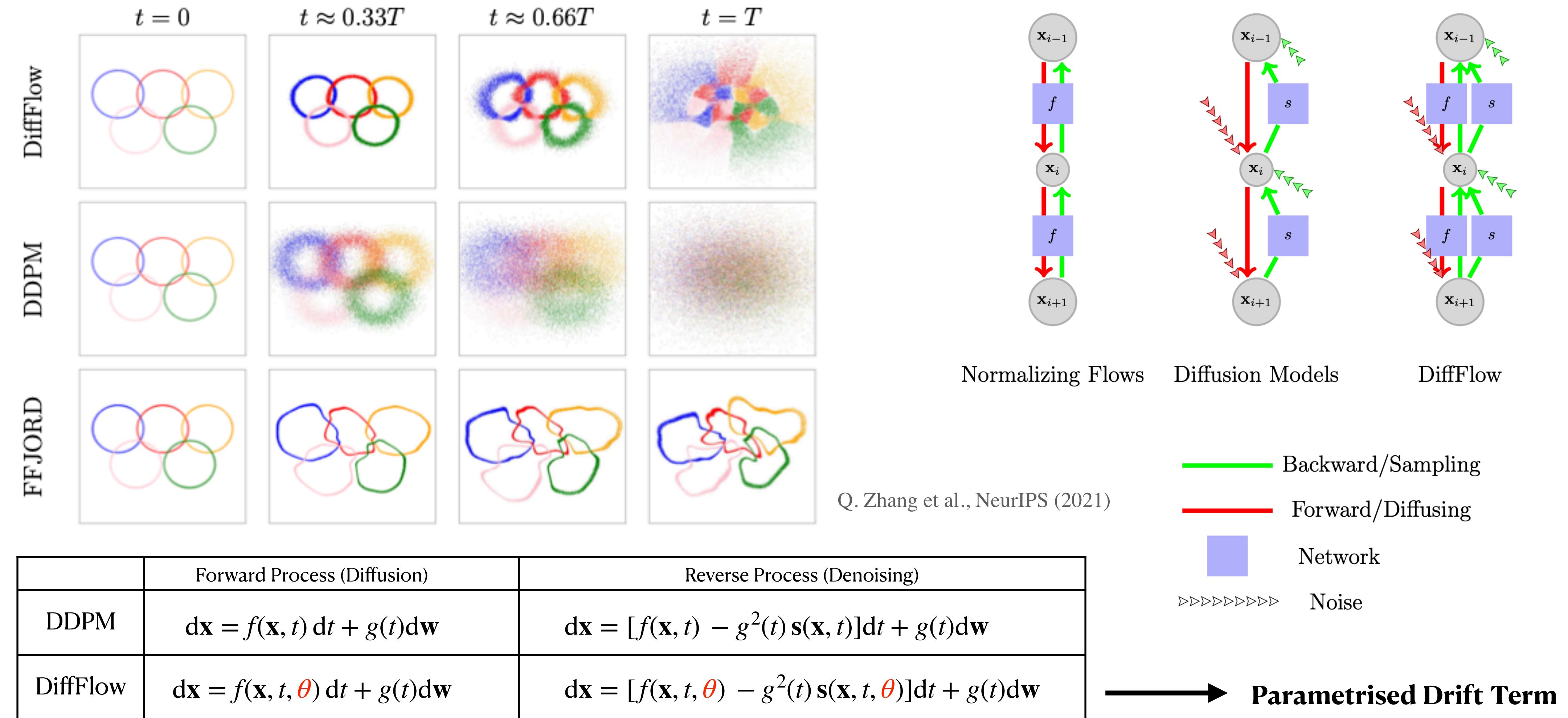
# From Normalizing Flows to Diffusion Models



In principle one can compute the **exact likelihood** on any input data leveraging instantaneous change of variable transformation from **Neural ODEs**.

However, Diffusion Models: **I**) are slow to train and sample from, **II**) require data, **III**) are memory intensive **IV**) need a lot of intermediate steps.

# DiffFlow: Diffusion Normalizing Flows



# How Big is Enough?

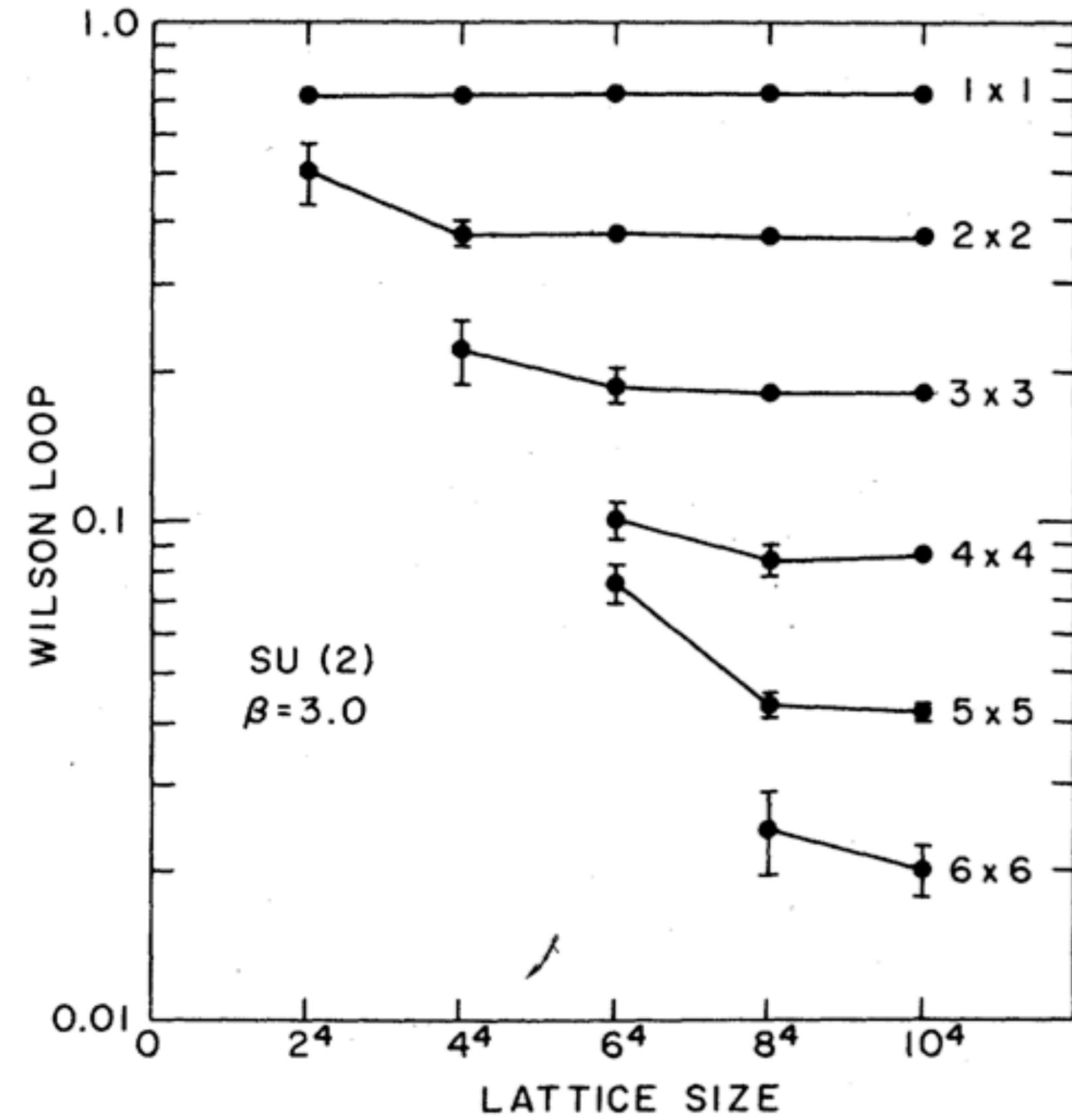


FIG. 3. Wilson loops at  $\beta = 3$  as a function of lattice size.  
M. Creutz, Phys. Rev. D (1980)

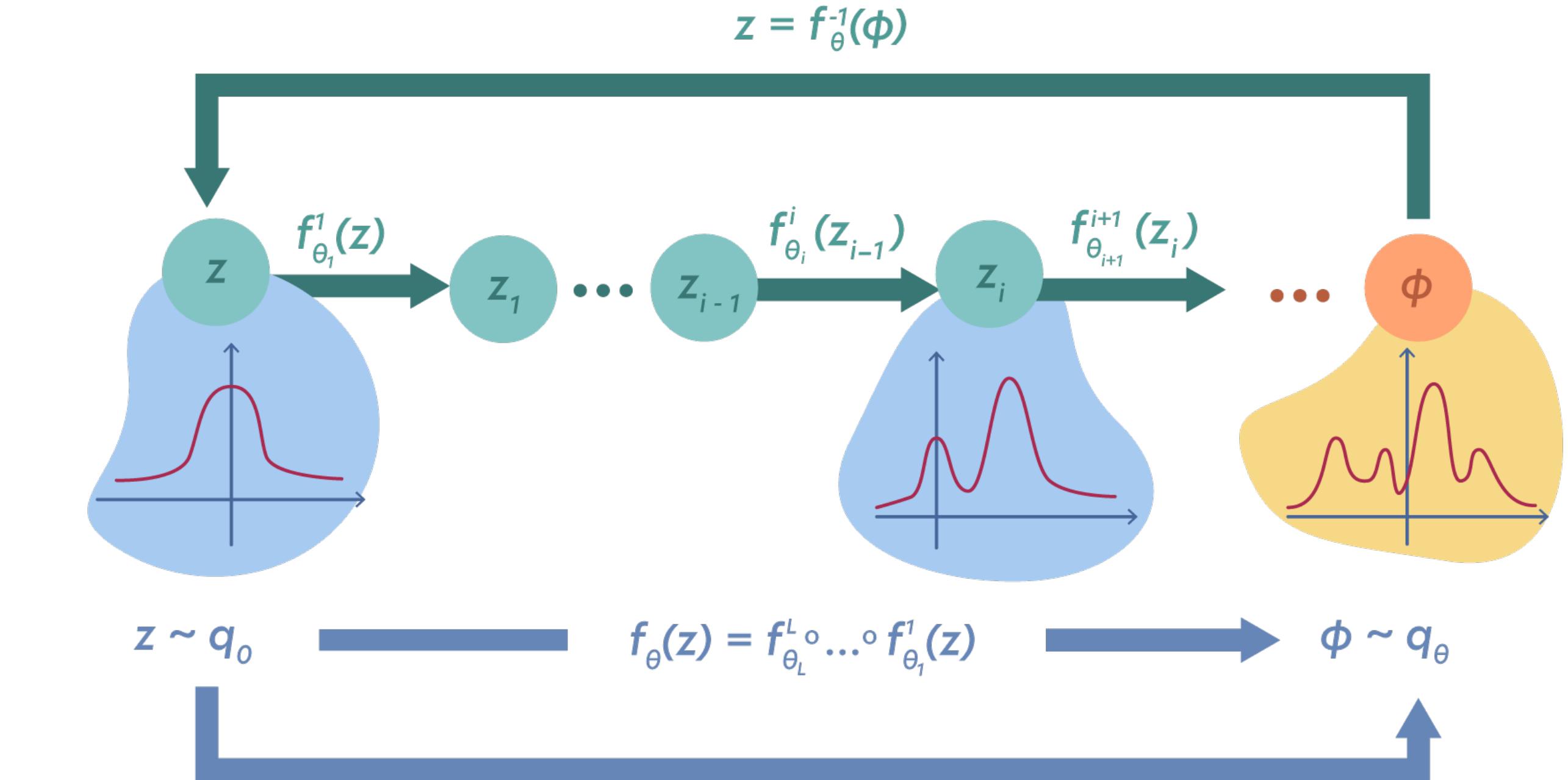
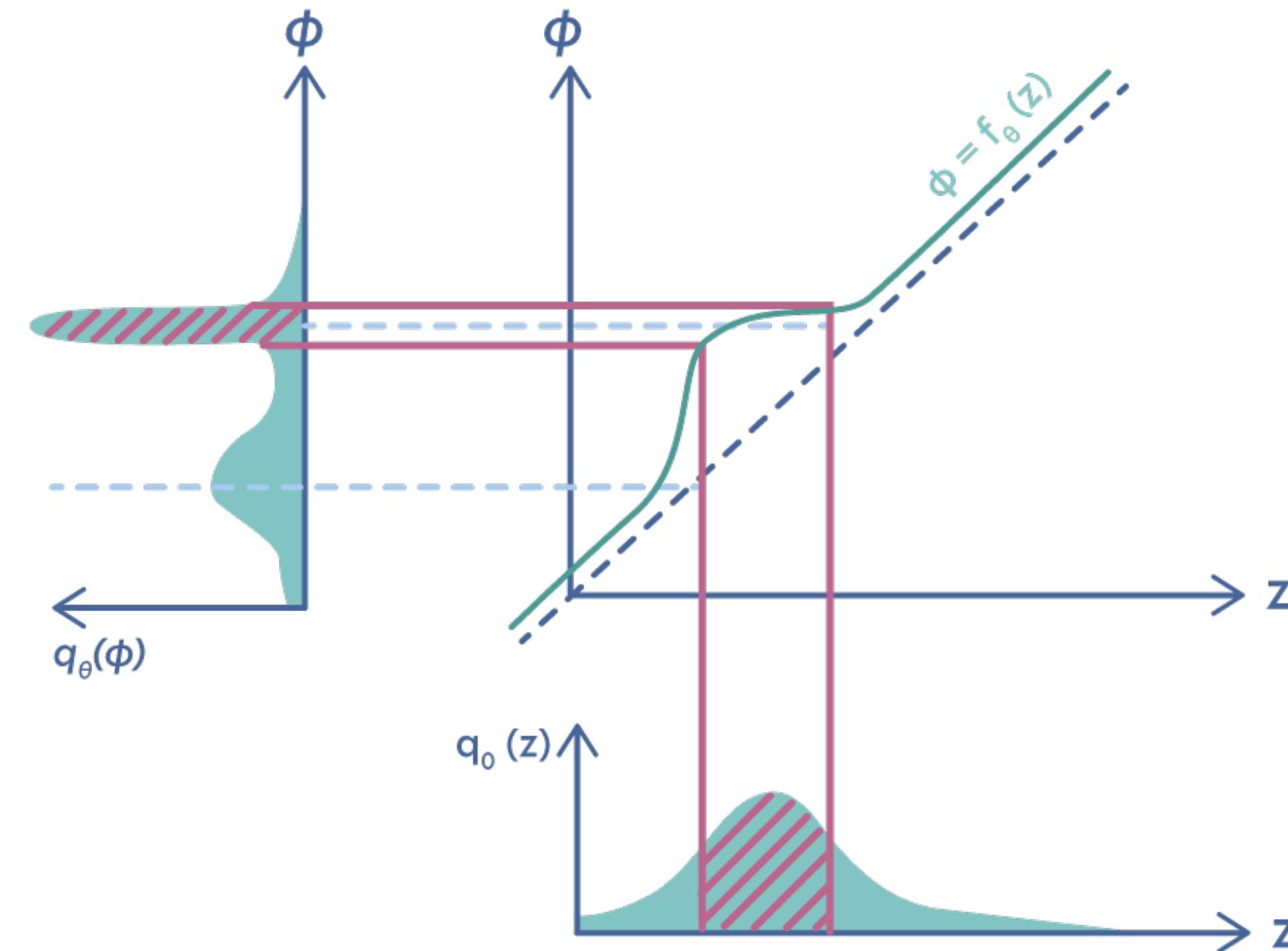
**$64 \times 8$  is not small !!!**

- $\sim 40$  years were needed to scale up LQCD.
- To start small is totally justified.
- Large room for improvement.
- Initial experiments as proof of concepts.

**Full QCD goal:**  $64^3 \times 128 \times 4 \times 8 \times 2 \approx 10^9$

Temporal extent  
Spatial extent  
Links from vertex  
SU(3) dimensionality  
Complex valued

# Normalizing Flows



- Bijective transformation  $\phi = f_{\theta}(z)$
- Invertible and differentiable\*
- Tractable Jacobian



$$q_{\theta}(\phi) = q_0(f_{\theta}^{-1}(\phi)) \left| \det \left( \frac{\partial f_{\theta}}{\partial z} \right) \right|^{-1}$$

\*Enough to be a **diffeomorphism** and **differentiable almost everywhere** wrt the Lebesgue measure of the base space

# Normalizing Flows: a different intuition

---

Let be  $z \in \mathbb{R}^D$  be a **random variable** sampled from a **base** distribution  $p_z : \mathbb{R}^D \rightarrow \mathbb{R}$

Let be  $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be a **diffeomorphism** on  $\mathbb{R}^D$  such that the base r.v. transforms like  $\phi = f(z)$ .

Using the change of variable formula one can compute the probability density of the transformed variable  $\phi$

$$q_\theta(\phi) = q_0(f_\theta^{-1}(\phi)) \left| \det \left( \frac{\partial f_\theta}{\partial z} \right) \right|^{-1}$$

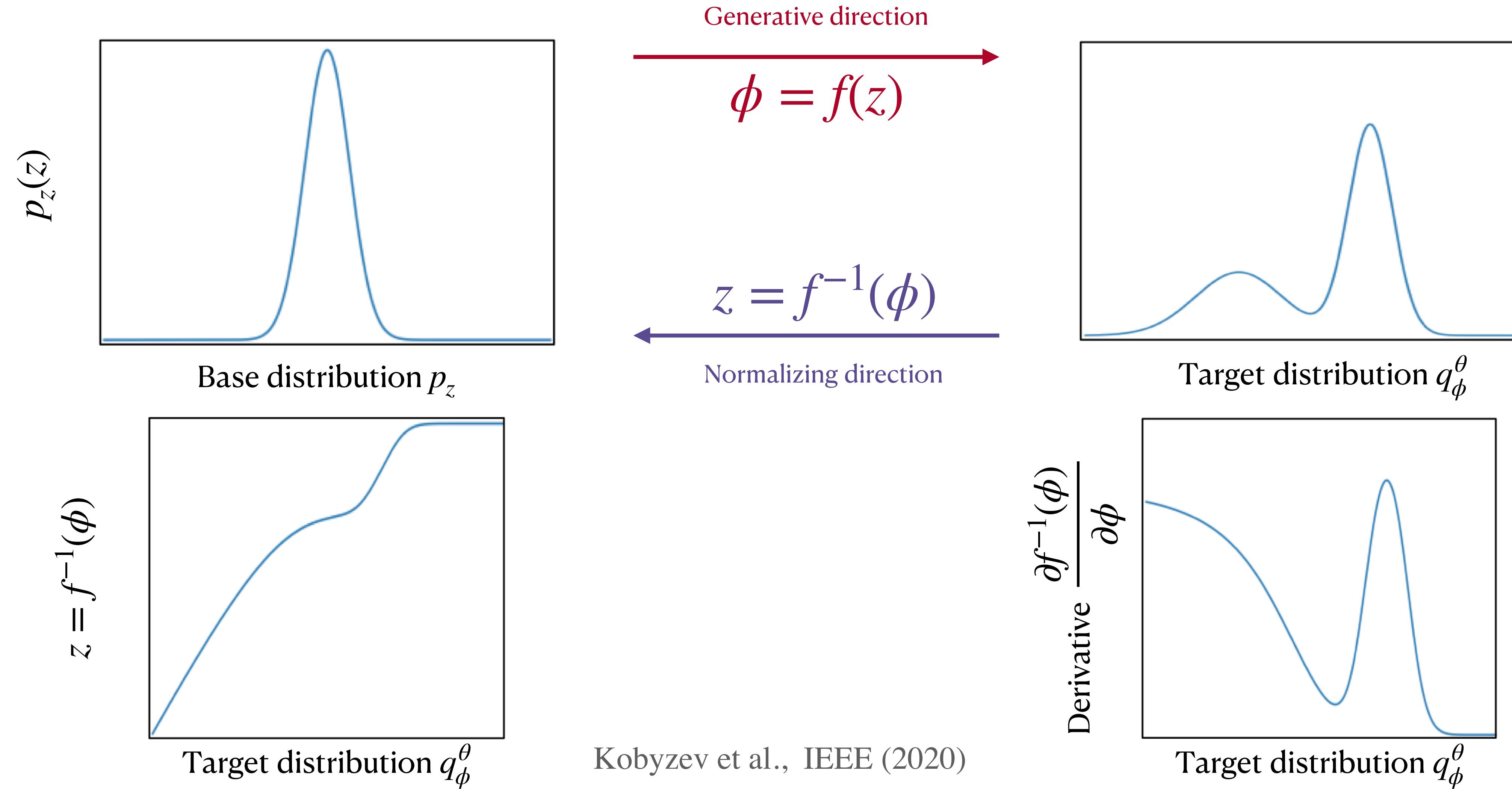
This new density is called the **pushforward** of the base density  $p_z$  by the function  $f$  and can be denoted by  $f_* p_z$

**Generative Direction:** the function  $f$  pushes forward the base density  $p_z$  (noise) to a more complicated one  $q_\phi^\theta$  parametrised by  $\theta$ .

One can also define a **pullback**  $f_*^{-1} q_\phi^\theta$  which maps the non-trivial density back into the base one.

**Normalizing Direction:** the inverse function  $f^{-1}$  flows backwards pulling back the complicated density towards the more regular form of the base measure  $p_z$ .

# Normalizing Flows: a different intuition



# Normalizing Flows: a different intuition

---

This view is what gives rise to the name “**normalizing flows**” as  $f$  is “normalizing” the data distribution.

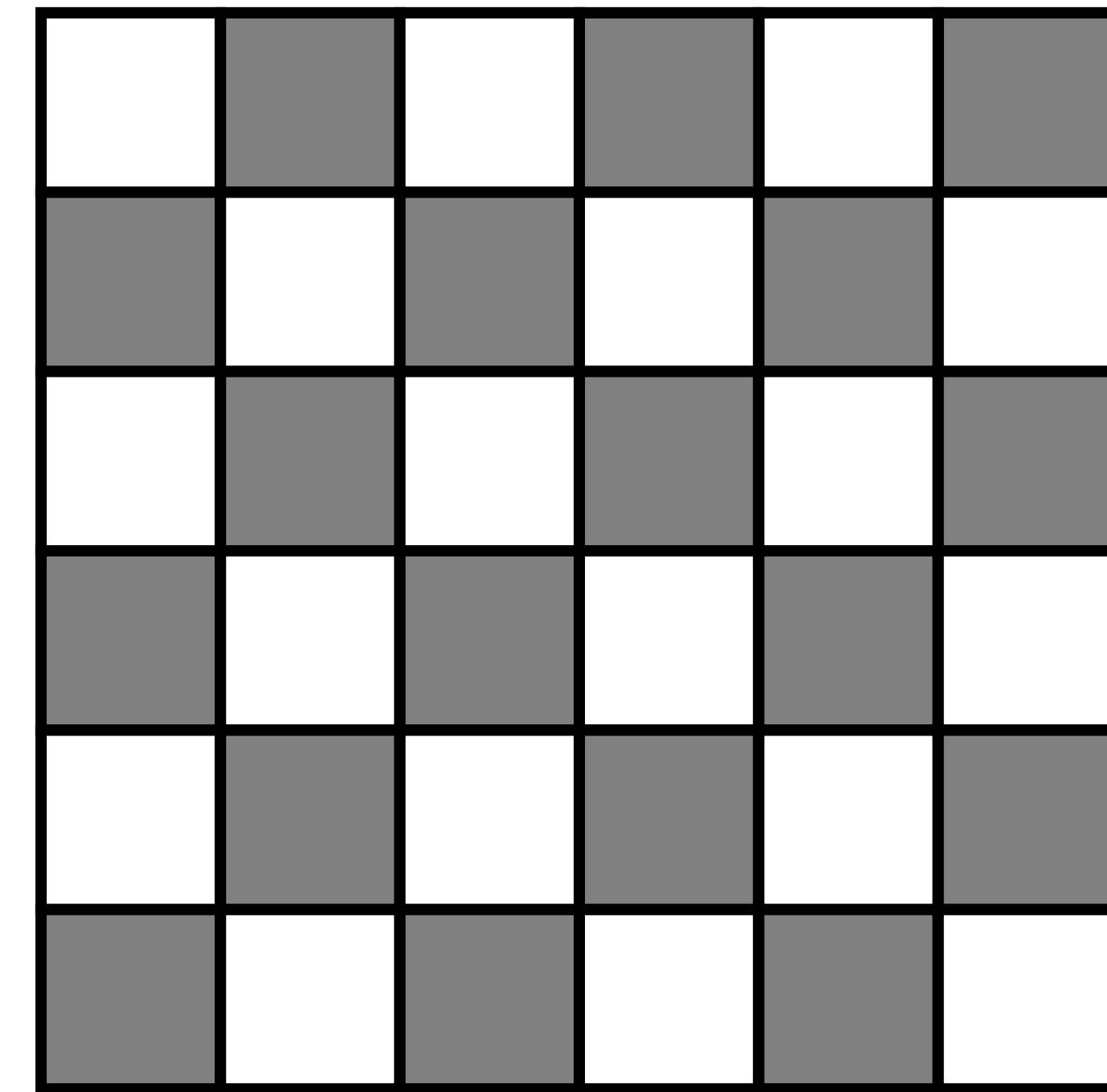
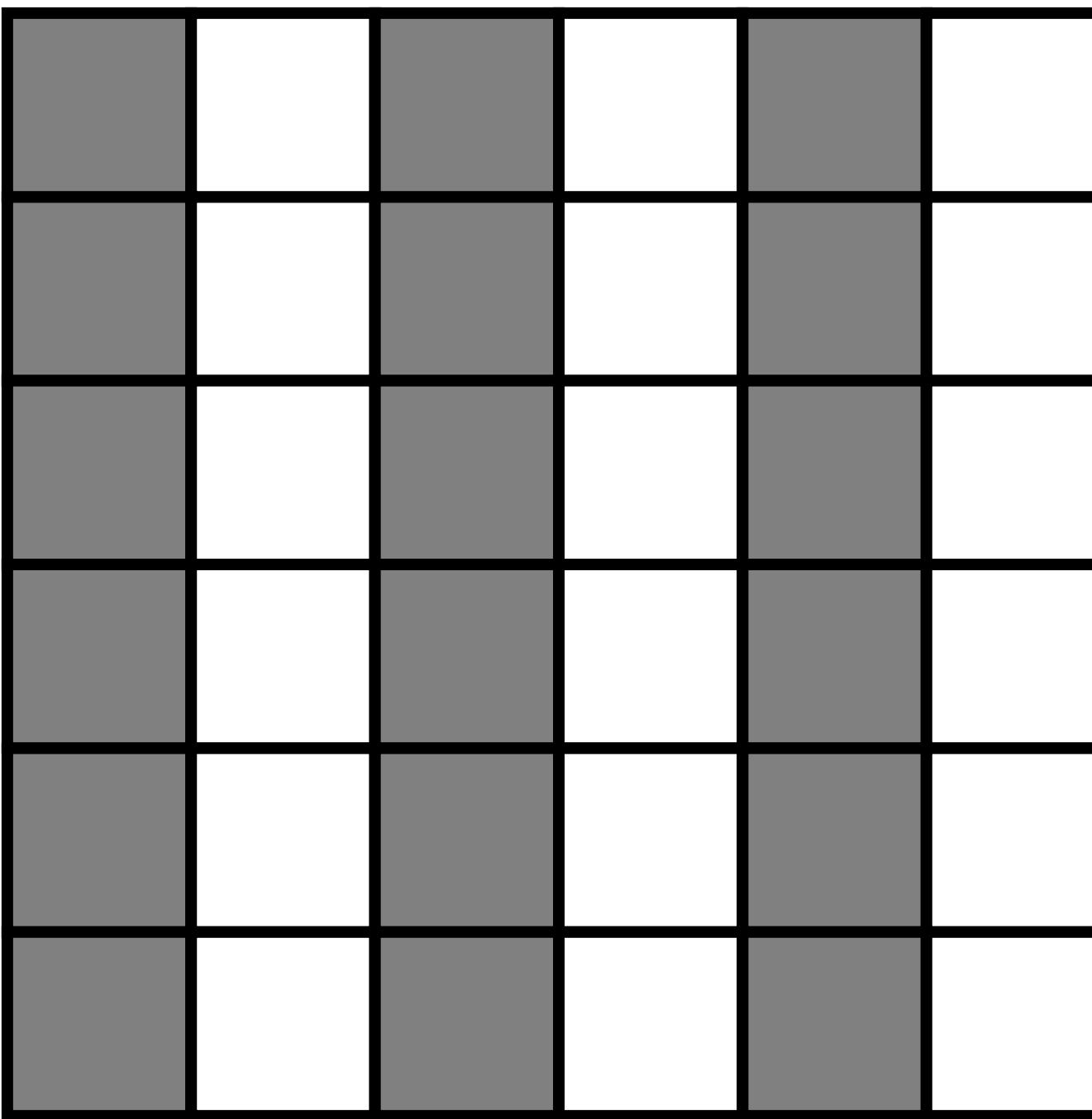
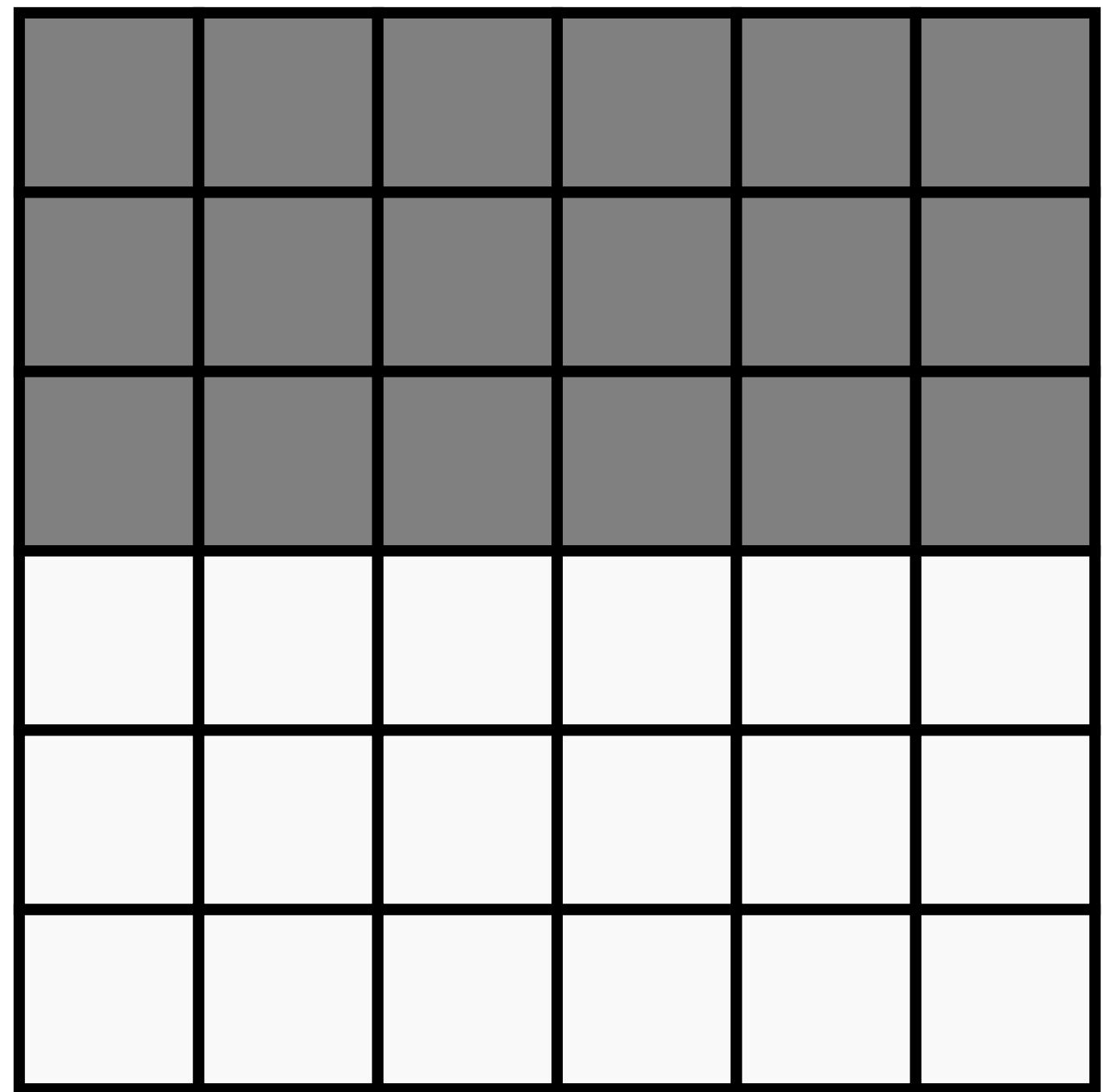
Normalizing Flows  $\implies$  ~~bijections~~  Formally incorrect

Normalizing Flows  $\implies$  **diffeomorphism** differentiable **almost everywhere** with respect to the Lebesgue measure on  $\mathbb{R}^D$ .

This allows for piecewise differentiable functions to be used to construct the transformation  $f$ .

# NICE Splitting Strategies

---



# Training a Flow with Reverse-KL

---

- **Minimize** the Reverse KL divergence

$$KL(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)} = \mathbb{E}_{q_\theta} [\ln q_\theta(\phi) - \ln p(\phi)] \geq 0$$

- **Maximize** the variational lower bound on the free energy  $F_q$

$$KL(q_\theta || p) = \underbrace{\mathbb{E}_{q_\theta} [\ln q_\theta(\phi) + S(\phi)]}_{\beta F_q} + \underbrace{\ln Z}_{F = -\beta \ln Z} = \beta(F - F_q)$$

# Training a Flow with Forward-KL

---

---

**Input:**

- prior density, e.g.,  $q_z \sim \mathcal{N}(0, \mathbb{1})$
- parametric model with parameters  $\theta$
- parametric action  $S(\phi, \kappa, \lambda)$  with fixed coupling parameters  $\{\lambda, \kappa\}$
- empty tensor for storing a batch of  $B$  configurations  $\Phi \in \mathbb{R}^{B \times N_S \times N_T}$

**Result:**

- learned bijective transformation  $f_\theta$  s.t.  $\phi_i = f_\theta(z_i)$
- exact likelihood function  $q_\theta$

```
1: begin
2:   /* Generate training samples from HMC */  

3:   for  $c$  in  $1, \dots, C_{\max}$  do
4:      $\phi$  = sampleHMC ( $\lambda, \kappa$ ) ; // sample configurations  $\phi \in \mathbb{R}^{n \times N_S \times N_T}$ 
5:      $\Phi$  = concatenate( $\Phi, \phi$ ) ; // concatenate configurations
6:   end
7:   /* Flow training using the generated dataset  $\Phi$  */  

8:   for  $t$  in  $1, \dots, T_{\max}$  do
9:     /* Iterate over dataset to collect batches of configurations */  

10:    draw samples  $\phi = \{\phi_i\}_{i=1}^m$  from  $\Phi$  where  $\phi_i \in \mathbb{R}^{N_S \times N_T} \forall i \in \{1, \dots, m\}$ 
11:    use  $\phi$  to evaluate  $\frac{1}{m} \sum_{i=1}^m \ln q_\theta(\phi_i)$ ;
12:     $\Delta\theta = \operatorname{argmin}_\theta \mathbb{E}_p [\ln q_\theta(\phi)]$ ;
13:    update  $f_\theta$  with  $\theta \leftarrow \theta + \eta \Delta\theta$  ;
14:   end
15: end
16: return  $f_\theta, q_\theta$ 
```

---

# Two Sides of the Same Coin

## Learning from data

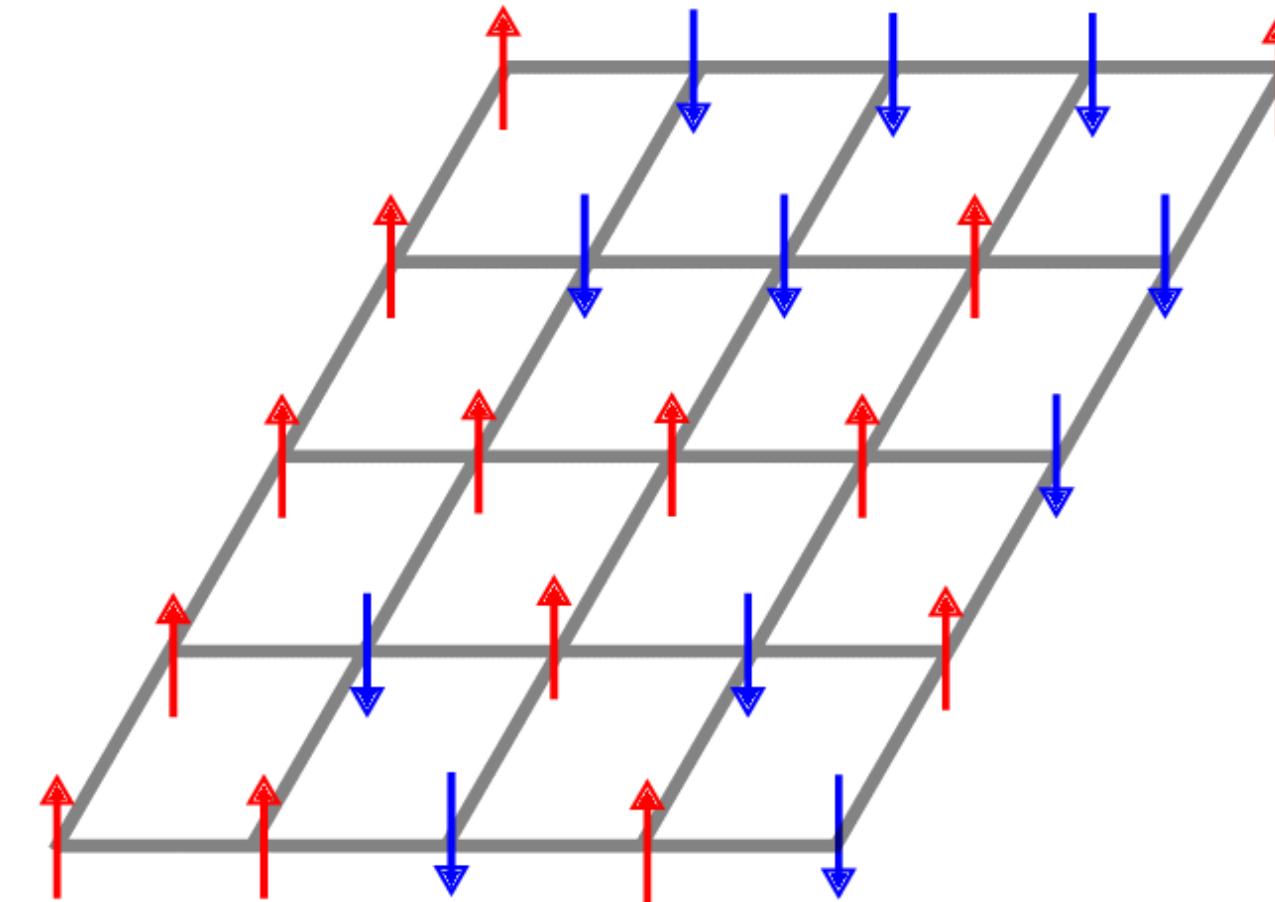


(Image credits: [medium.com](https://medium.com))

- Maximum likelihood (needs data)
- Minimising **Forward KL**
- Target: generating distribution

$$\mathcal{L} = \mathbb{E}_{\phi \sim p} [\ln q_{\theta}(\phi)]$$

## Learning from physics



(Image credits: Wald 2017)

- Self-sampling (needs action/hamilt.)
- Minimising **Reverse KL**
- Target: samples, normalisation  $Z$

$$\mathcal{L} = \mathbb{E}_{\phi \sim q_{\theta}} [\ln q_{\theta}(\phi) + S(\phi) + \ln Z]$$

(Slide inspired by Lei Wang lecture notes)

# The Free Energy as a Natural Cost Function

---

$$F = \underset{\text{Energy}}{E} - \underset{\text{Entropy}}{TS}$$

Free Energy is a **Cost Function** given by nature!

$$KL(q_\theta || p) = \int D[\phi] q_\theta(\phi) \ln \frac{q_\theta(\phi)}{p(\phi)} \approx \mathbb{E}_{q_\theta} [S(\phi)] - \mathbb{E}_{q_\theta} [\ln q_\theta(\phi)] \simeq F_q \geq 0$$

The optimisation problem translates into minimising a variational form of the free energy  $F_q$

(Slide inspired by Lei Wang lecture notes)

# HMC Estimate of Free Energies

---

$$\Delta F_{eb} = F_e - F_b \implies F_e = \Delta F_{eb} + F_b$$

In the limit of vanishing hopping parameter:

$$F_0(\lambda) = -|\Lambda| T \ln z_0(\lambda) = -T \ln Z_0(\lambda) \quad \text{with} \quad z_0(\lambda) = \sqrt{\frac{1-2\lambda}{4\lambda}} \exp\left(\frac{(1-2\lambda)^2}{8\lambda}\right) K_{\frac{1}{4}}\left(\frac{(1-2\lambda)^2}{8\lambda}\right)$$

Difference of free energies can also be understood in terms of ratios of partition functions:

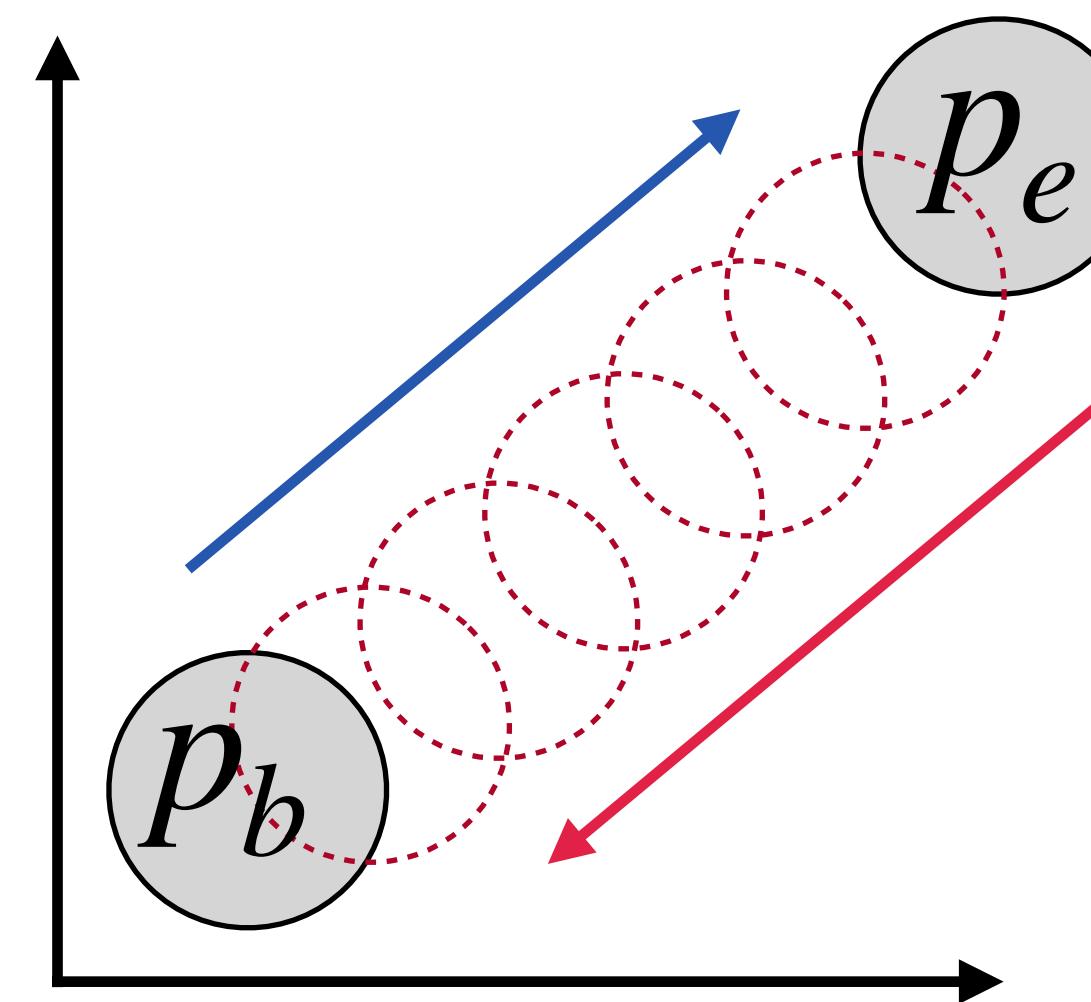
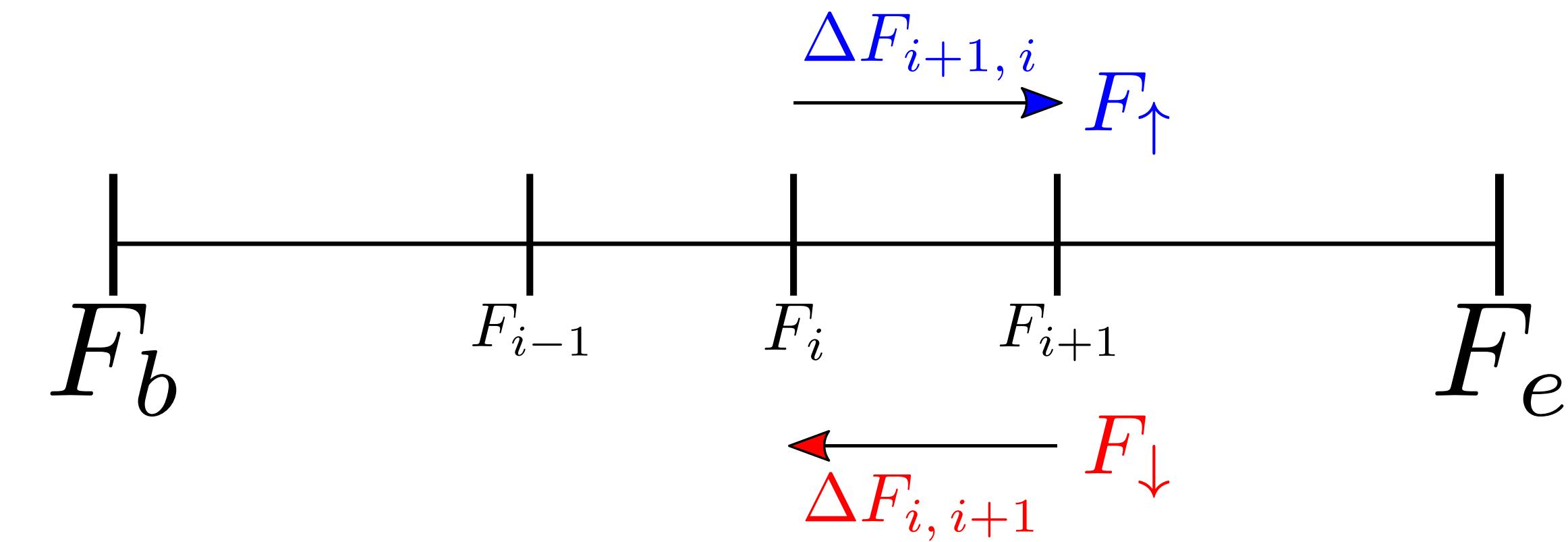
$$\Delta F_{i+1,i} = F_{i+1} - F_i = -T \ln \frac{Z_{i+1}}{Z_i}$$

⚠ Actions should be similar to ensure small variance!

Which can be estimated by MCMC sampling:

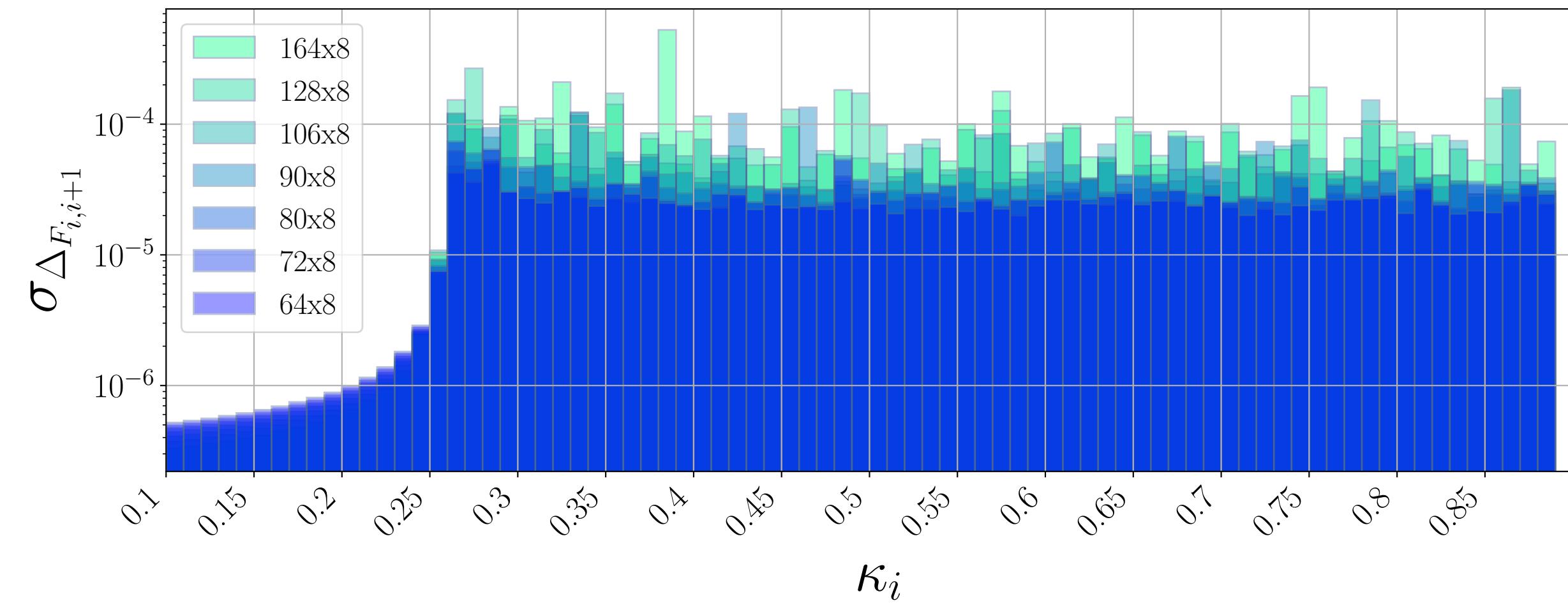
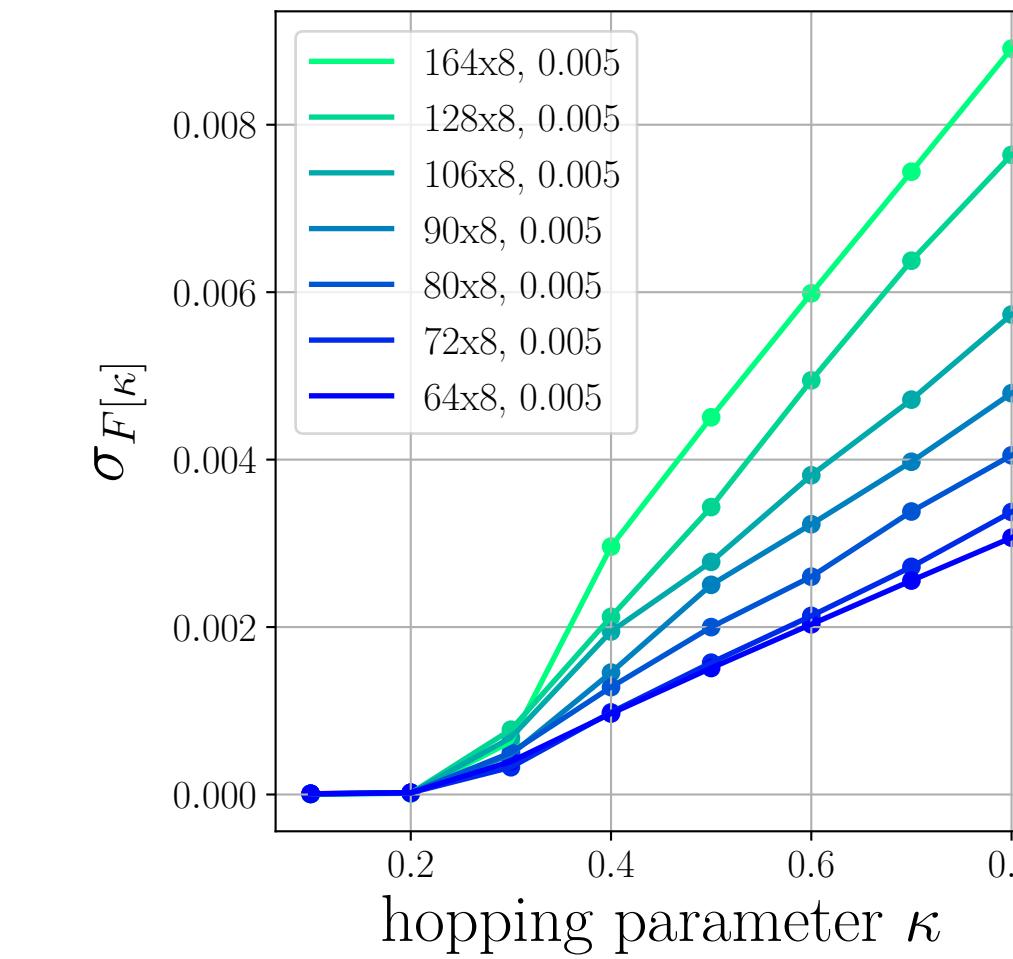
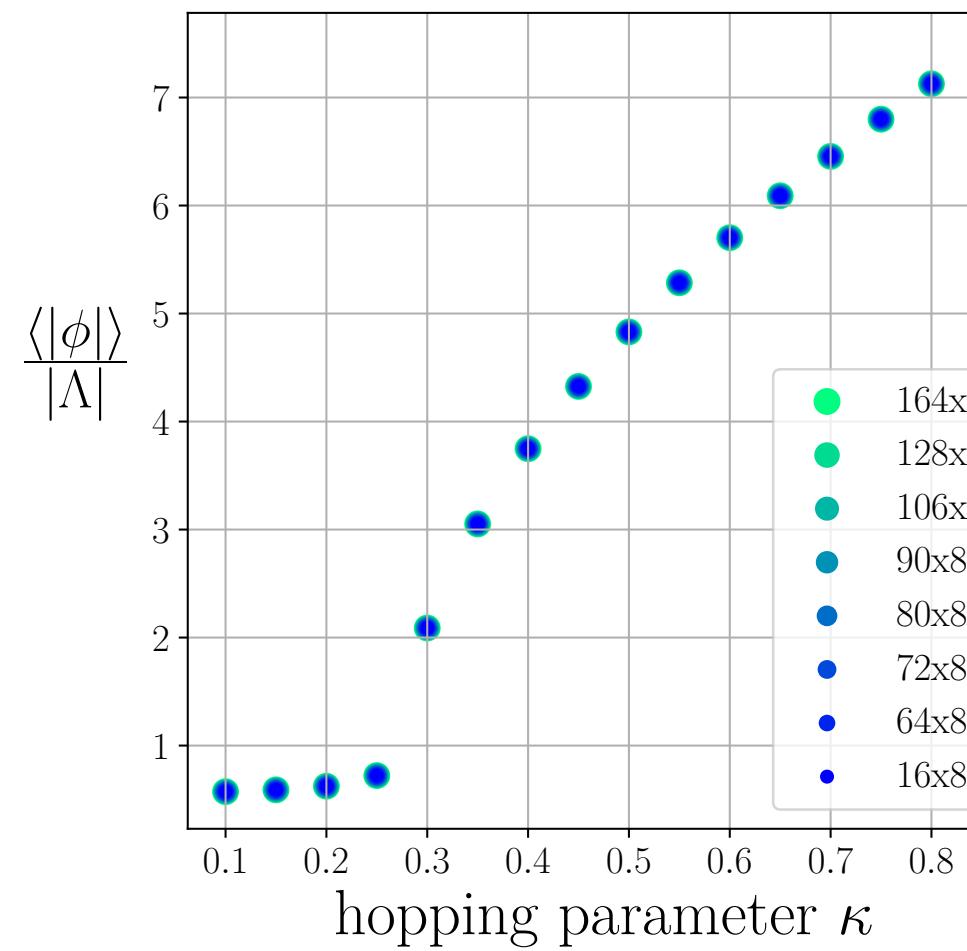
$$\mathbb{E}_{p_i} \left[ \frac{\exp(-S_{i+1})}{\exp(-S_i)} \right] = \frac{1}{Z_i} \int \mathcal{D}[\phi] e^{-S_i(\phi)} \overbrace{\frac{e^{-S_{i+1}(\phi)}}{e^{-S_i(\phi)}}}^{\sim} = \frac{Z_{i+1}}{Z_i}$$

# HMC Estimate of Free Energies

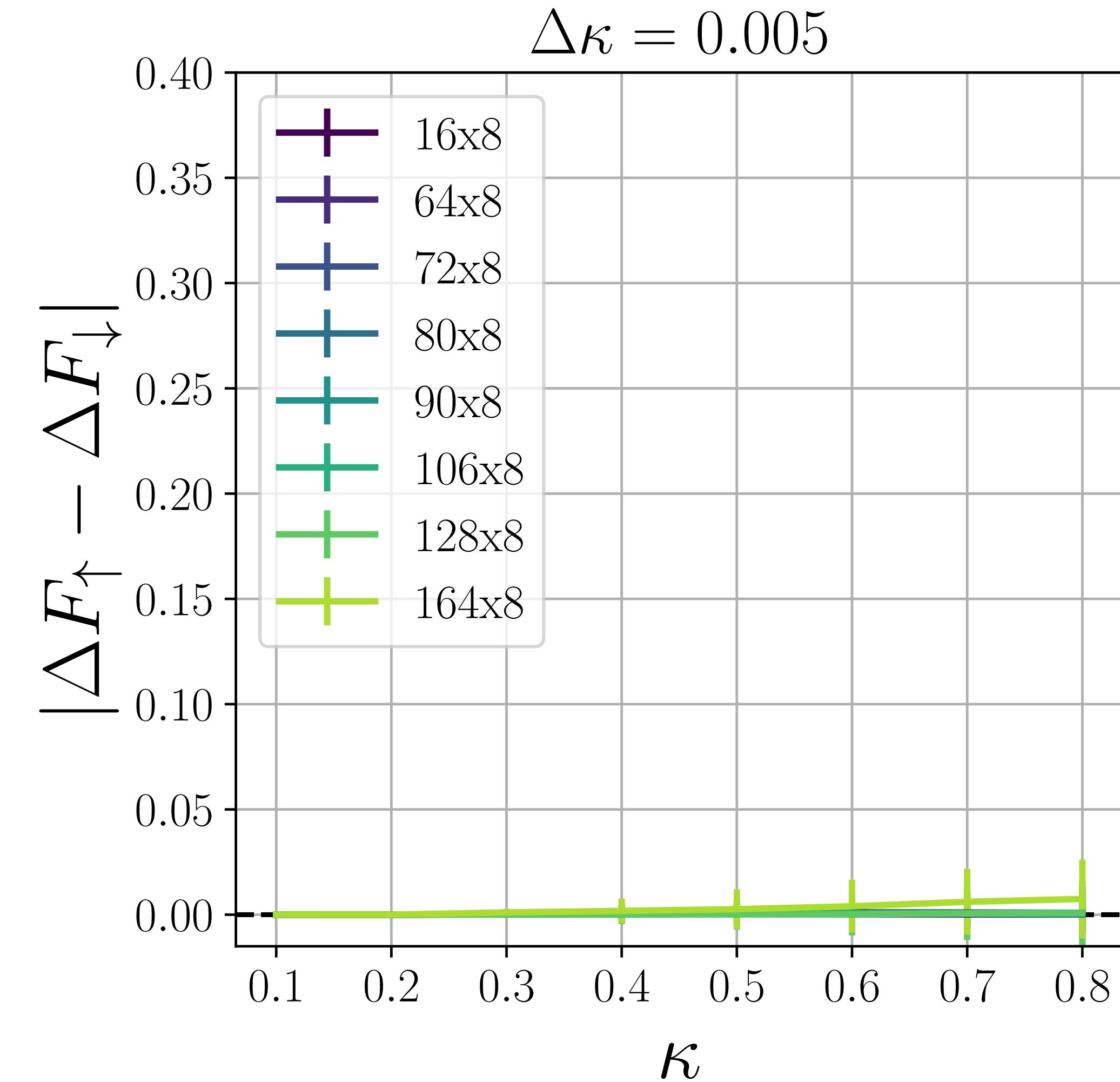
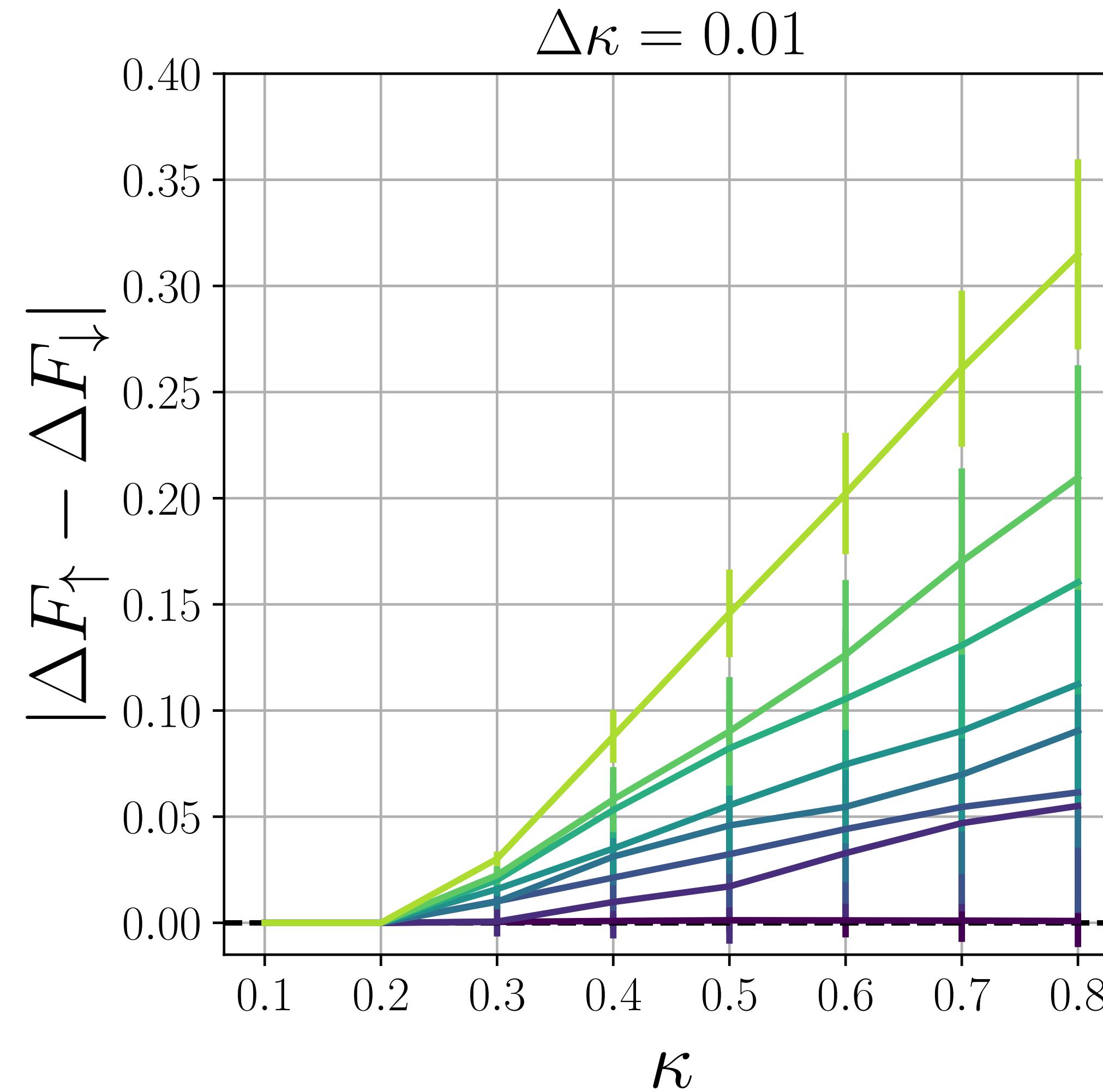


(Image inspired by Weinsberger et al. (2020))

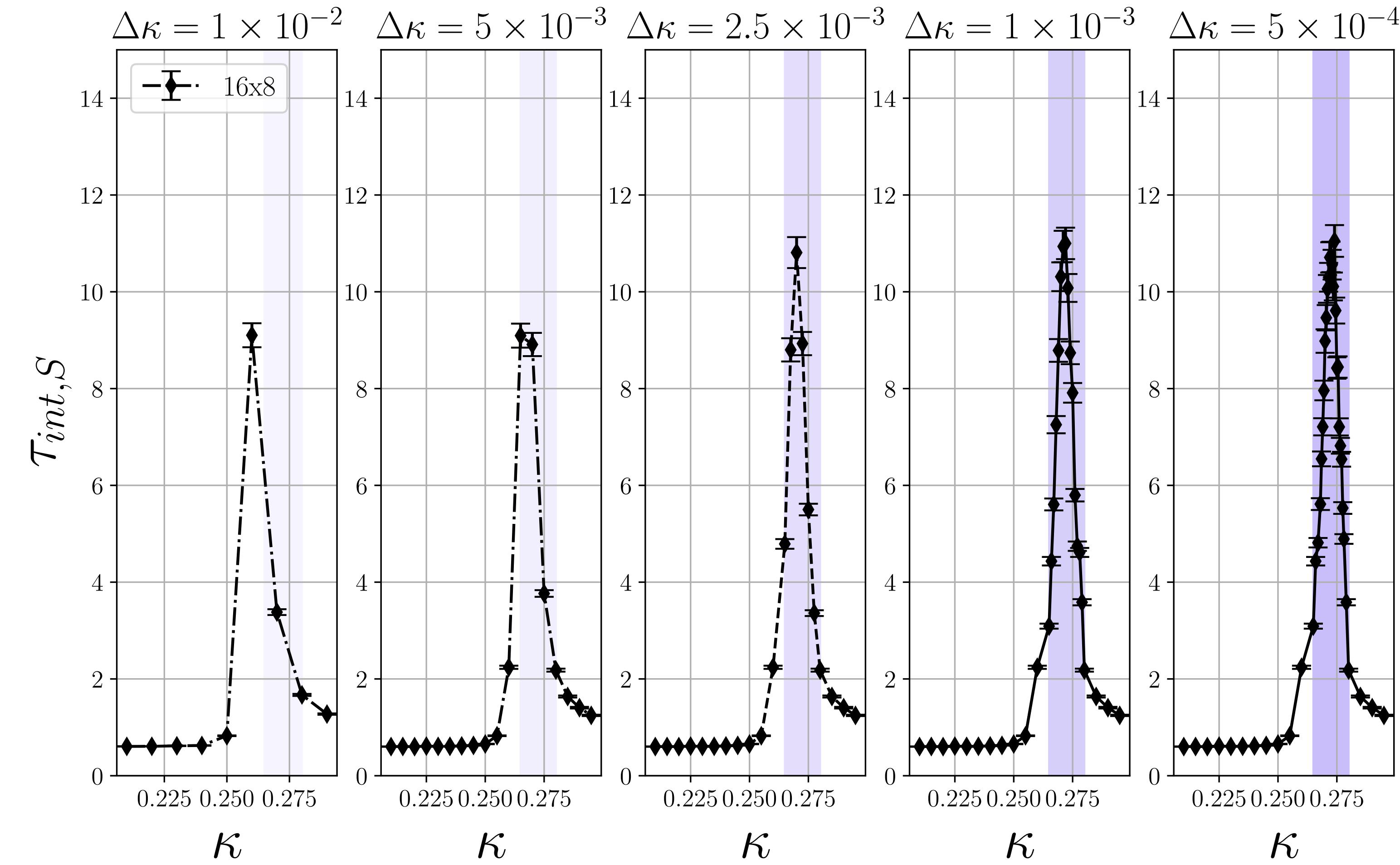
# HMC Estimate of Free Energies



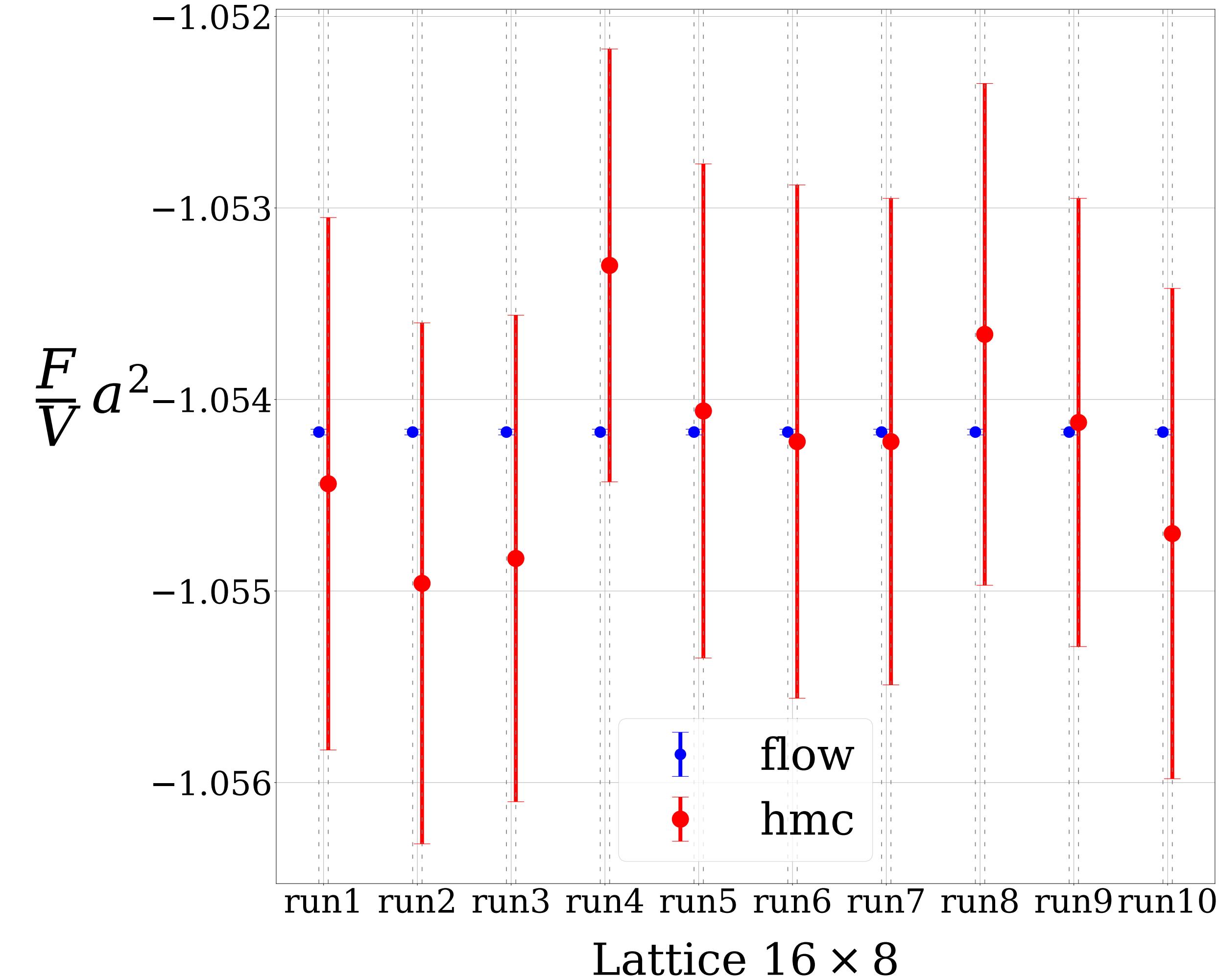
# HMC Estimate of Free Energies



# Integrated Autocorrelation



# Stability of the Estimator



KAN, et al., Phys. Rev. Lett. (2021)

# Reparametrization Trick

---

$$\text{KL}(q_\theta || p) = \mathbb{E}_{x \sim q_\theta} [S(x) + \log q_\theta(x)] + \text{const.}$$

Reparametrization Trick  $\implies x \sim q_\theta \rightarrow z \sim q_z$

$$\mathbb{E}_{x \sim q_\theta} [f(x)] = \mathbb{E}_{z \sim q_z} [f(g_\theta(z))]$$

Compute Gradient  $\frac{d}{d\theta}$

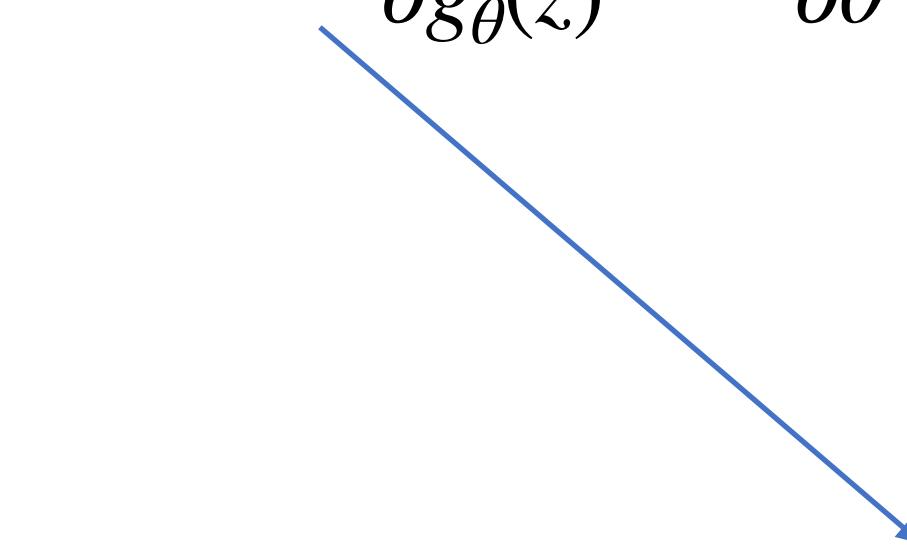
$$\frac{d}{d\theta} \text{KL}(q_\theta || p) = \mathbb{E}_{z \sim q_z} \left[ \frac{d}{d\theta} S(g_\theta(z)) + \frac{d}{d\theta} \log q_\theta(g_\theta(z)) \right] \underset{MC}{\approx} G_{Rep}$$

# Total Gradient Estimator

$$G_{RepQP} = \frac{1}{N} \sum_{i=1}^N \left( \frac{d}{d\theta} S(g_\theta(z_i)) + \frac{d}{d\theta} \log q_\theta(g_\theta(z_i)) \right), \quad z_i \sim q_Z$$

Path derivative

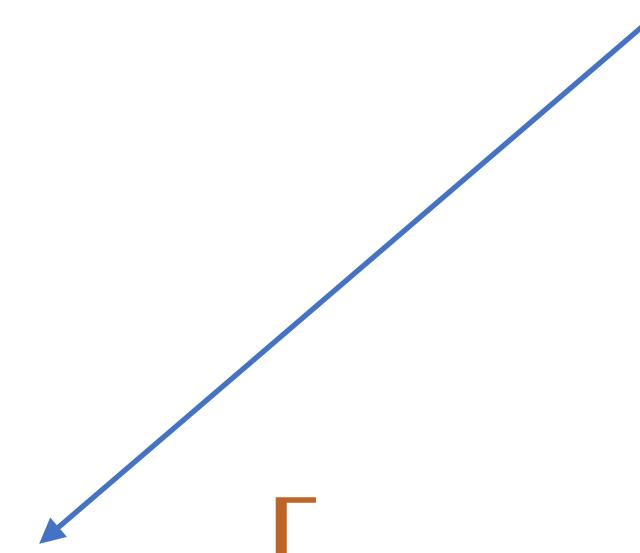
$$\nabla_\theta f(g_\theta(z), \theta) = \frac{\partial f(g_\theta(z), \theta)}{\partial g_\theta(z)} \frac{\partial g_\theta(z)}{\partial \theta}$$



$$\frac{d}{d\theta} \text{KL}(q_\theta || p) = \mathbb{E}_{z \sim q_Z} [ \nabla_\theta S(g_\theta(z)) + \nabla_\theta \log q_\theta(g_\theta(z)) ]$$

Total derivative

$$\frac{d}{d\theta} f(g_\theta(z), \theta) = \nabla_\theta f(g_\theta(z), \theta) + \frac{\partial}{\partial \theta} f(x, \theta) \Big|_{x=g_\theta(z)}$$



$$\left[ \frac{\partial}{\partial \theta} \log q_\theta(x) \Big|_{x=g_\theta(z)} \right]$$

Partial derivative

# Path and Score Terms

---

$$G_{RepQP} = \frac{1}{N} \sum_{i=1}^N \left( \frac{d}{d\theta} S(g_\theta(z_i)) + \frac{d}{d\theta} \log q_\theta(g_\theta(z_i)) \right), \quad z_i \sim q_Z$$

Path derivative

$$\nabla_\theta f(g_\theta(z), \theta) = \frac{\partial f(g_\theta(z), \theta)}{\partial g_\theta(z)} \frac{\partial g_\theta(z)}{\partial \theta}$$

$$\frac{d}{d\theta} \text{KL}(q_\theta || p) = G_{RepQP} = G_{PathQP} + G_{Score}$$

Total derivative

$$\frac{d}{d\theta} f(g_\theta(z), \theta) = \nabla_\theta f(g_\theta(z), \theta) + \frac{\partial}{\partial \theta} f(x, \theta) \Big|_{x=g_\theta(z)}$$

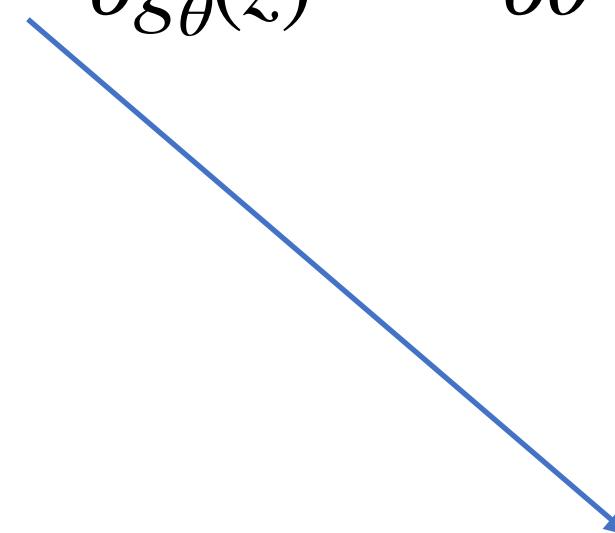
Partial derivative

# Path and Score Terms

$$G_{RepQP} = \frac{1}{N} \sum_{i=1}^N \left( \frac{d}{d\theta} S(g_\theta(z_i)) + \frac{d}{d\theta} \log q_\theta(g_\theta(z_i)) \right), \quad z_i \sim q_Z$$

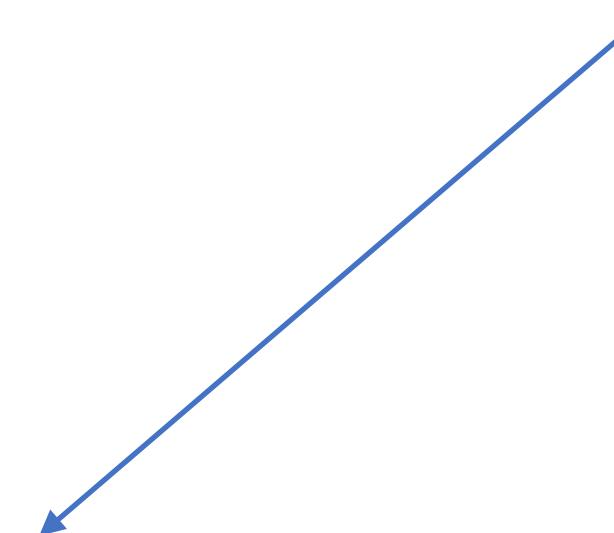
Path derivative

$$\nabla_\theta f(g_\theta(z), \theta) = \frac{\partial f(g_\theta(z), \theta)}{\partial g_\theta(z)} \frac{\partial g_\theta(z)}{\partial \theta}$$



Total derivative

$$\frac{d}{d\theta} f(g_\theta(z), \theta) = \nabla_\theta f(g_\theta(z), \theta) + \frac{\partial}{\partial \theta} f(x, \theta) \Big|_{x=g_\theta(z)}$$

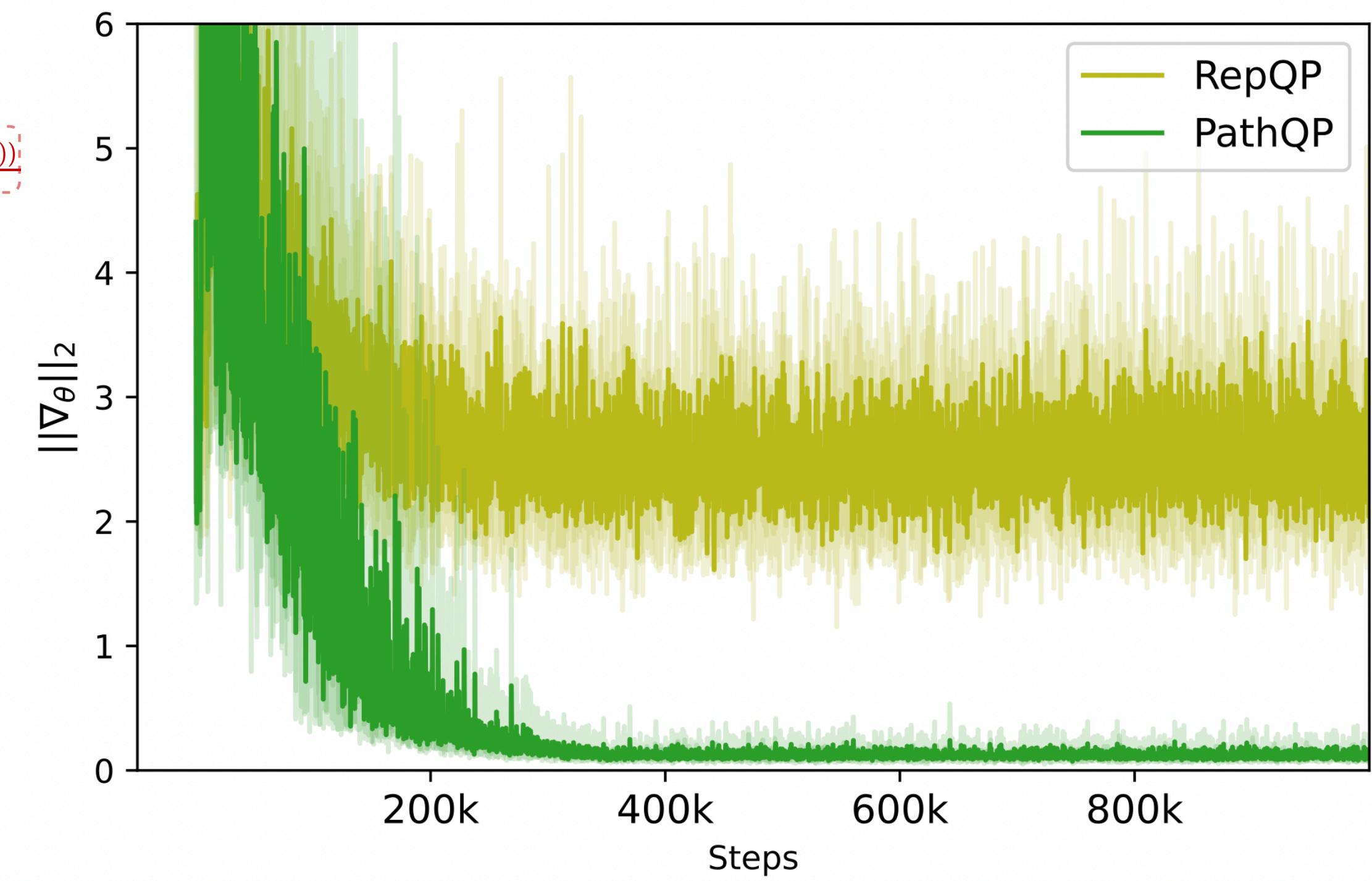
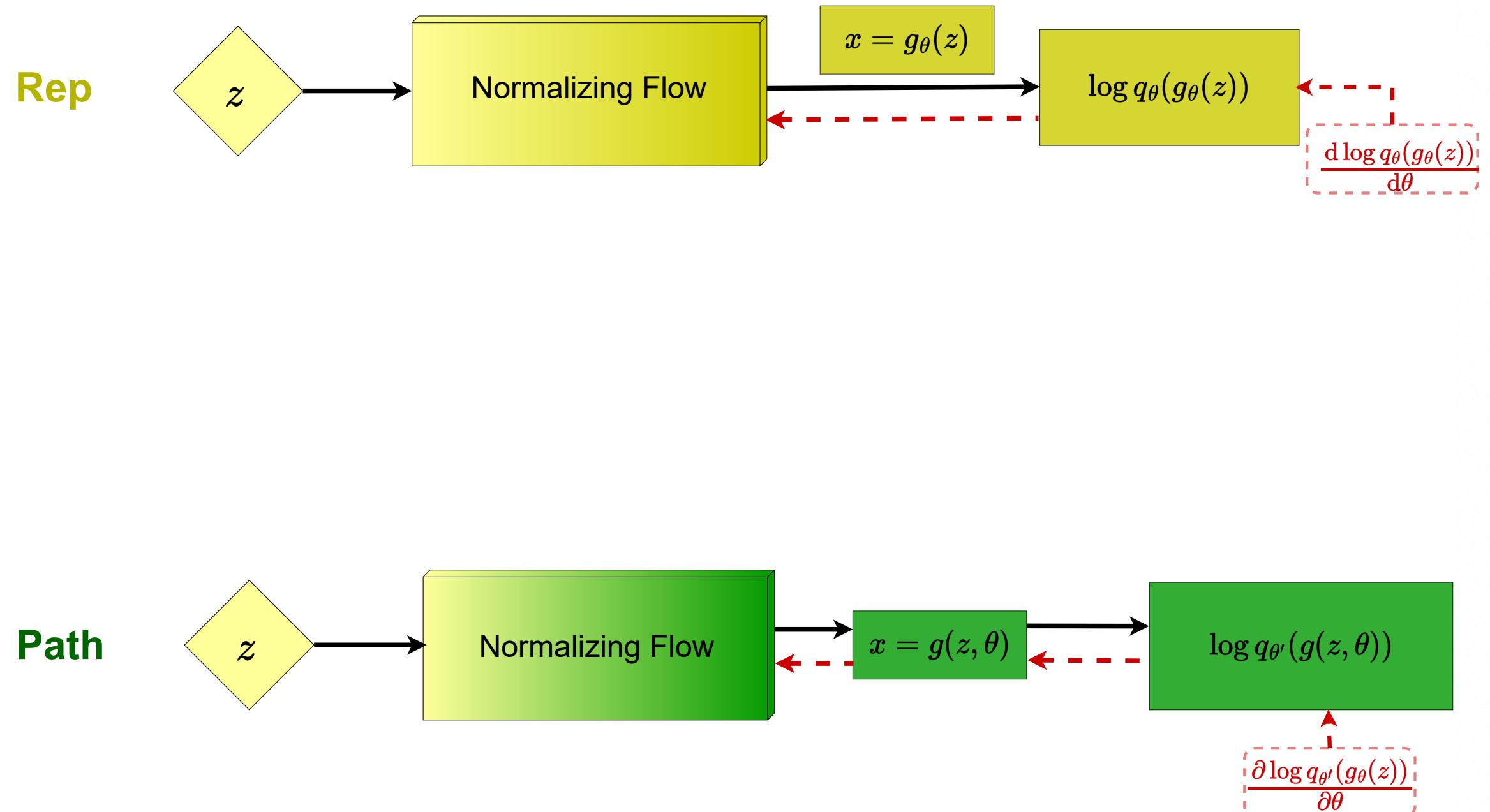


Partial derivative

$$G_{RepQP} \approx G_{PathQP} + \cancel{G_{Score}}$$

Roeder, G. et al., NeurIPS (2017)

# Stable Gradients Mean Better Convergence



# Stable Gradients Mean Better Convergence

---

---

**Algorithm 1:** Path gradient  $\nabla_{\theta} \log q_{\theta}(g_{\theta}(z))$ .

---

**Input:** base sample  $z \sim q_Z$

```
x' ← stop_gradient(gθ(z))          # forward pass of z through the flow without gradients
qθ(x') ← qZ(gθ-1(x'))  $\left| \frac{\partial g_{\theta}^{-1}(x')}{\partial x'} \right|$       # reverse pass to calculate density
G ←  $\frac{\partial \log(q_{\theta}(x'))}{\partial x'}$           # compute gradient with respect to x'
x ← gθ(z)                            # standard forward pass
return  $\frac{d}{d\theta} (\text{stop\_gradient}(G)^T x)$       # for path-gradient, contract  $\frac{\partial x}{\partial \theta}$  with  $\partial \log q_{\theta}(x)/\partial x$ 
```

---

# Quantum Mechanical Particle in Double Well

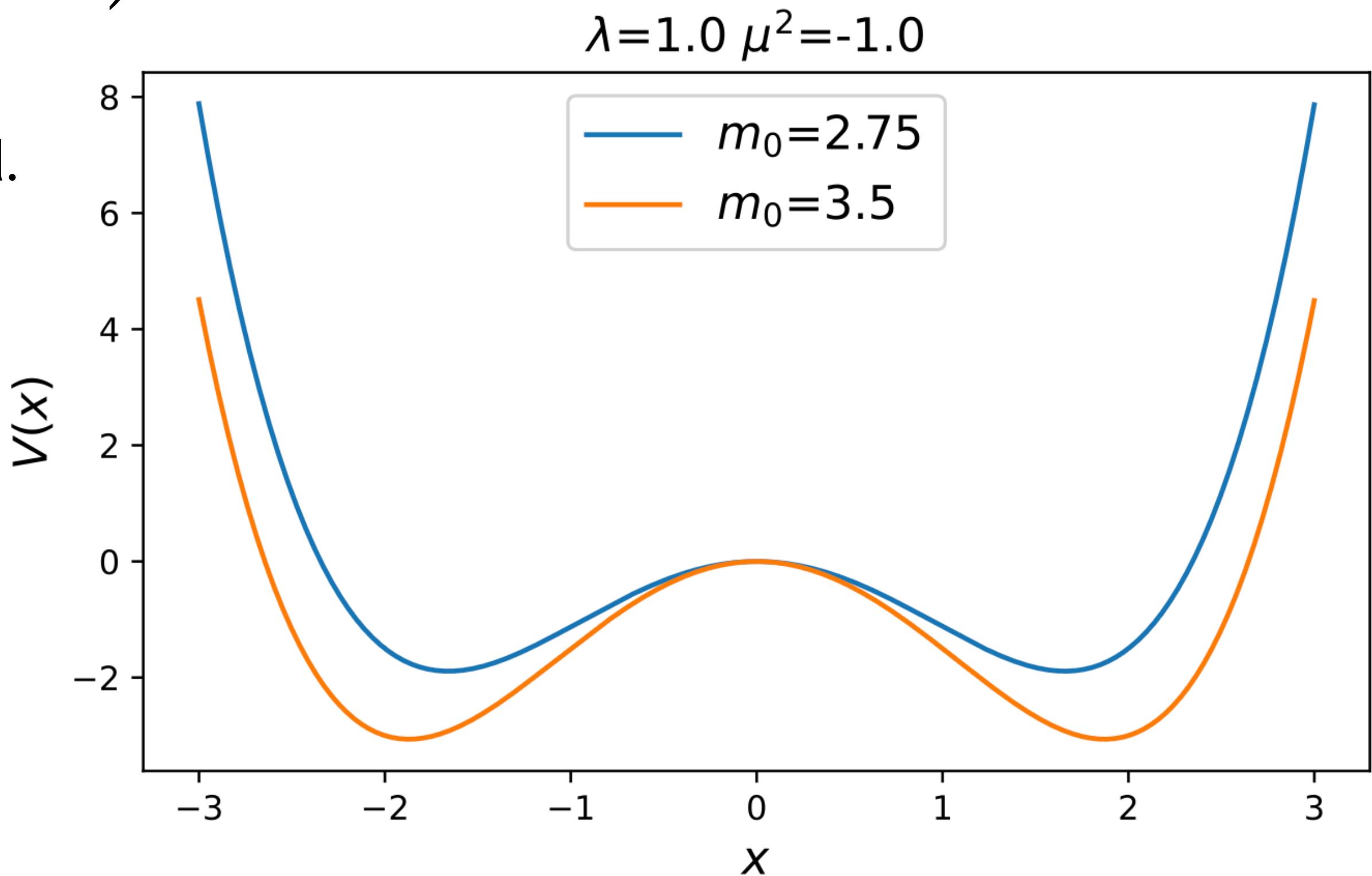
$$S(x) = a \sum_{t=0}^{T-1} \left( \frac{m_0}{2} (x_{t+1} - x_t)^2 + \frac{m_0 \mu^2}{2} x^2 + \frac{\lambda}{4} x^4 \right)$$

Where  $\{m_0, \mu, \lambda\}$  control the shape of the potential.

The probability density of the path is given by

$$p(\mathbf{x}) = \frac{1}{Z} \exp(-S(\mathbf{x}))$$

for a discretised trajectory  $\mathbf{x} = (x_0, x_1, \dots, x_{T-1})$



# A Few Words on the Effective Sampling Size

---

MCMC

$$\text{ESS} = N_{eff} = \frac{N}{2\tau_{int,\mathcal{O}}}$$

Reverse

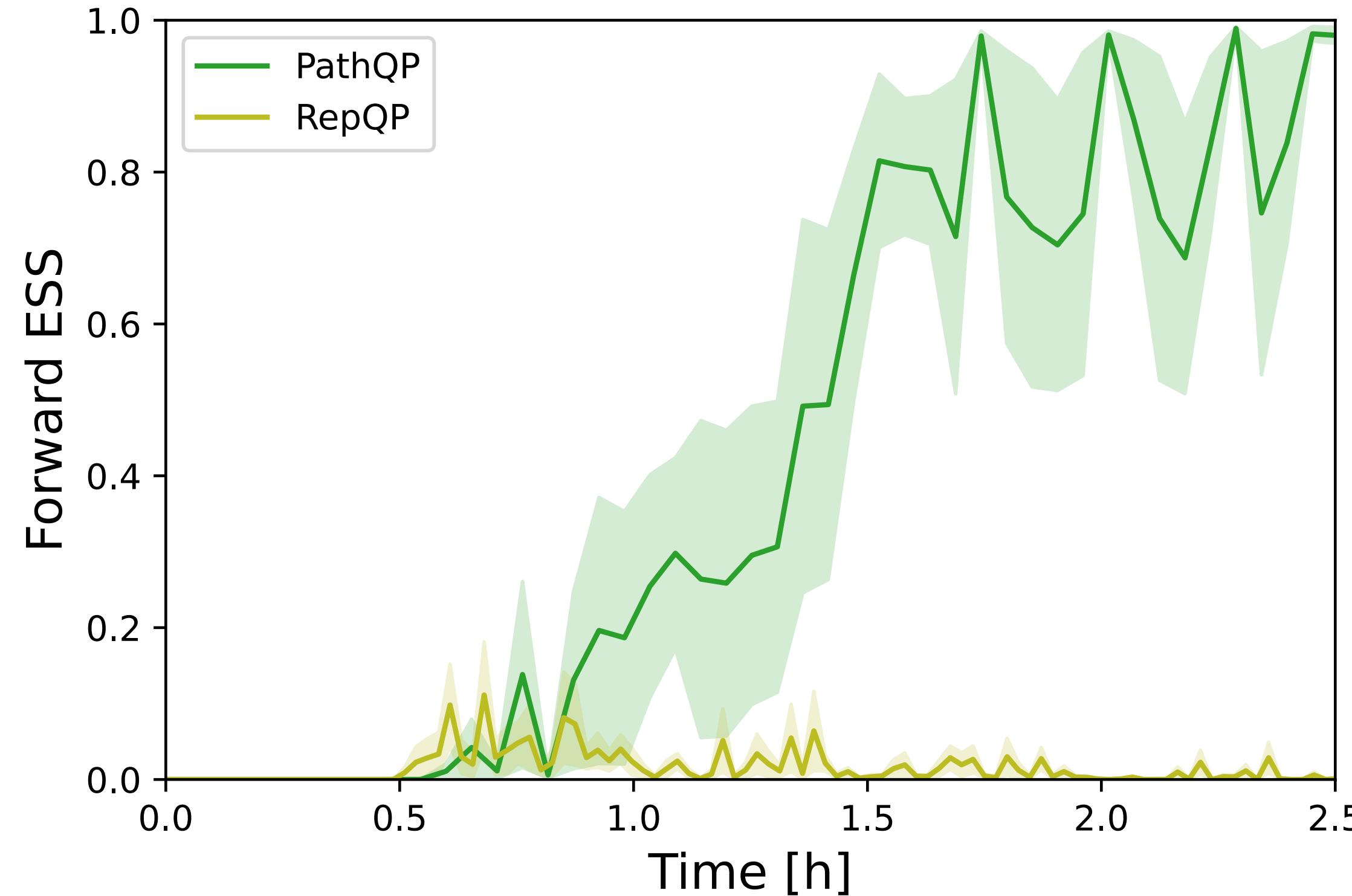
$$\text{ESS} = N_{eff} = \frac{N}{\mathbb{E}_{q_\theta}[w^2]}$$

Forward

$$\text{ESS} = N_{eff} = \frac{N}{\mathbb{E}_p[w]}$$

# Forward ESS for a Double-Well

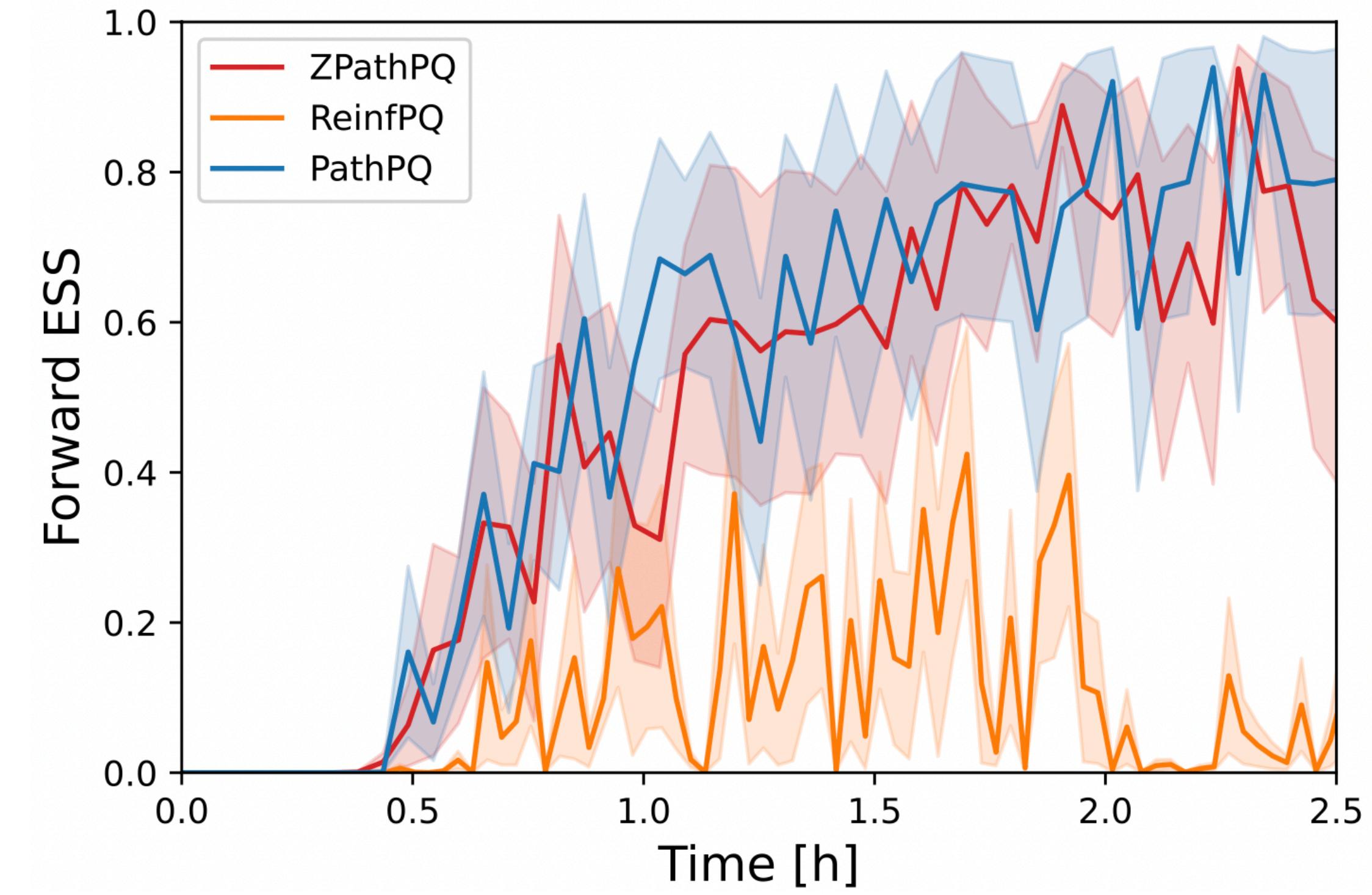
$$\{m_0 = 2.75, \mu^2 = -1, \lambda = 1\}$$



Unimodal target density

Reverse KL training approach

$$\{m_0 = 3.0, \mu^2 = -1, \lambda = 1\}$$

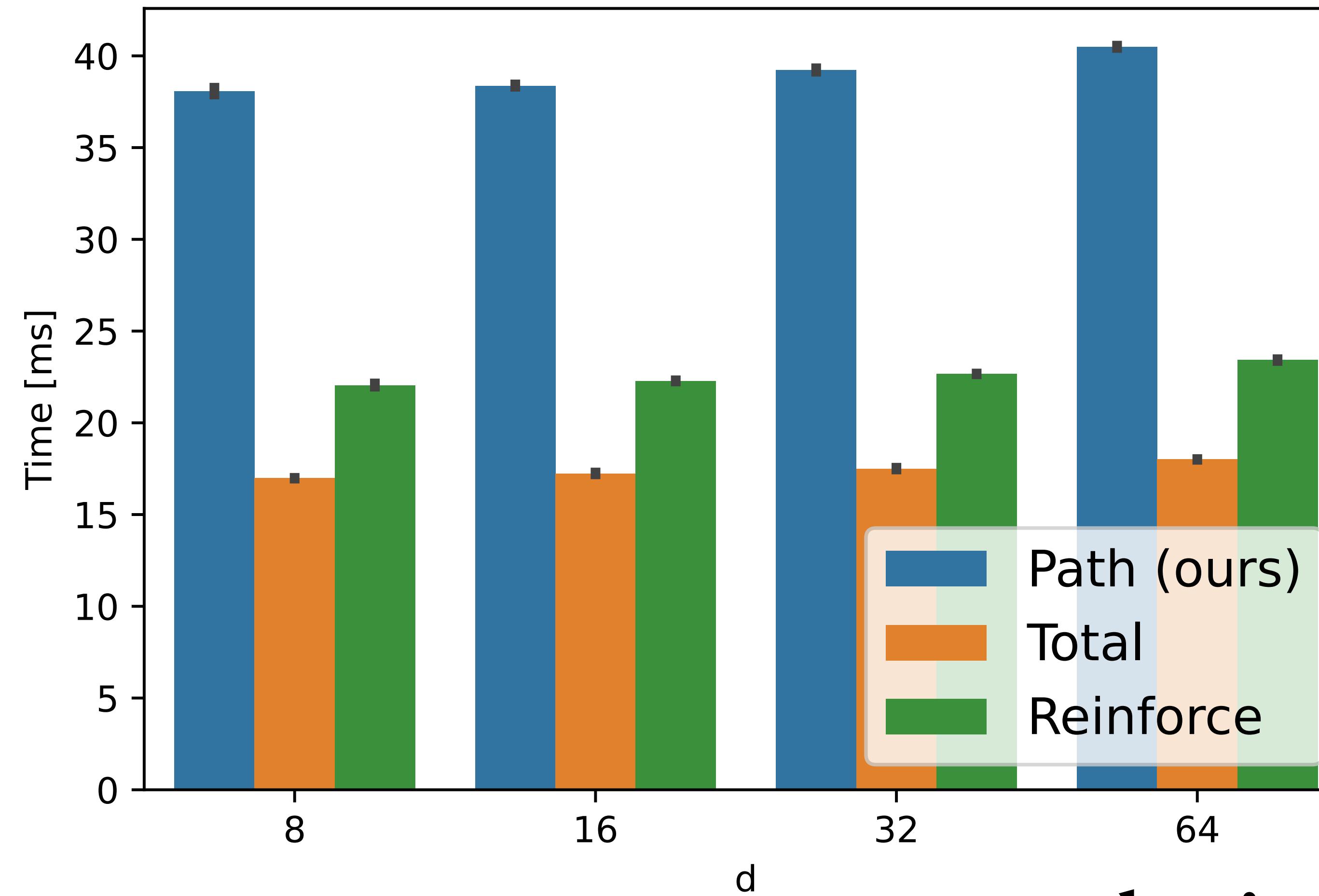


Multimodal target density

Forward KL training approach

# There's no free lunch...

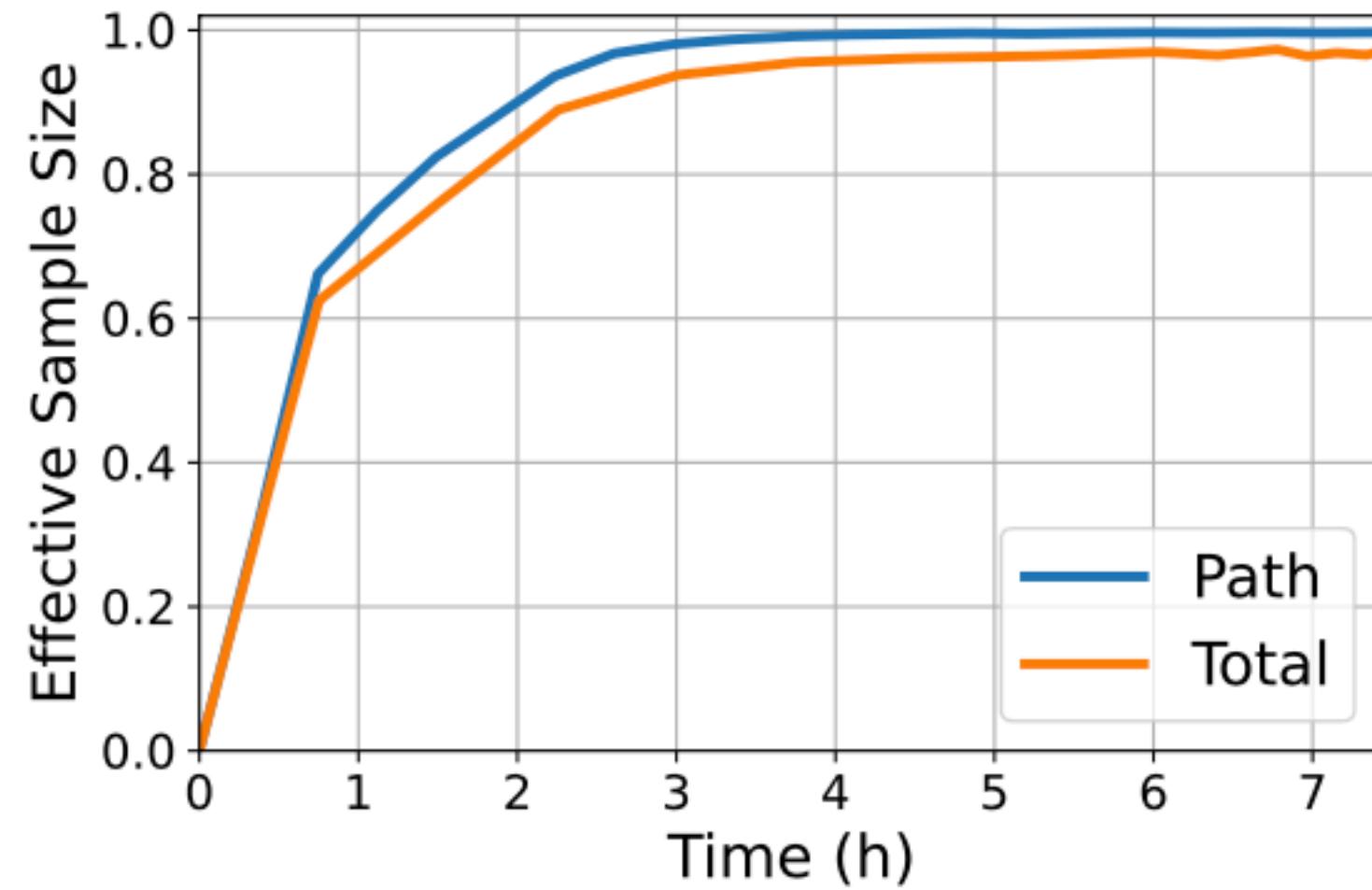
---



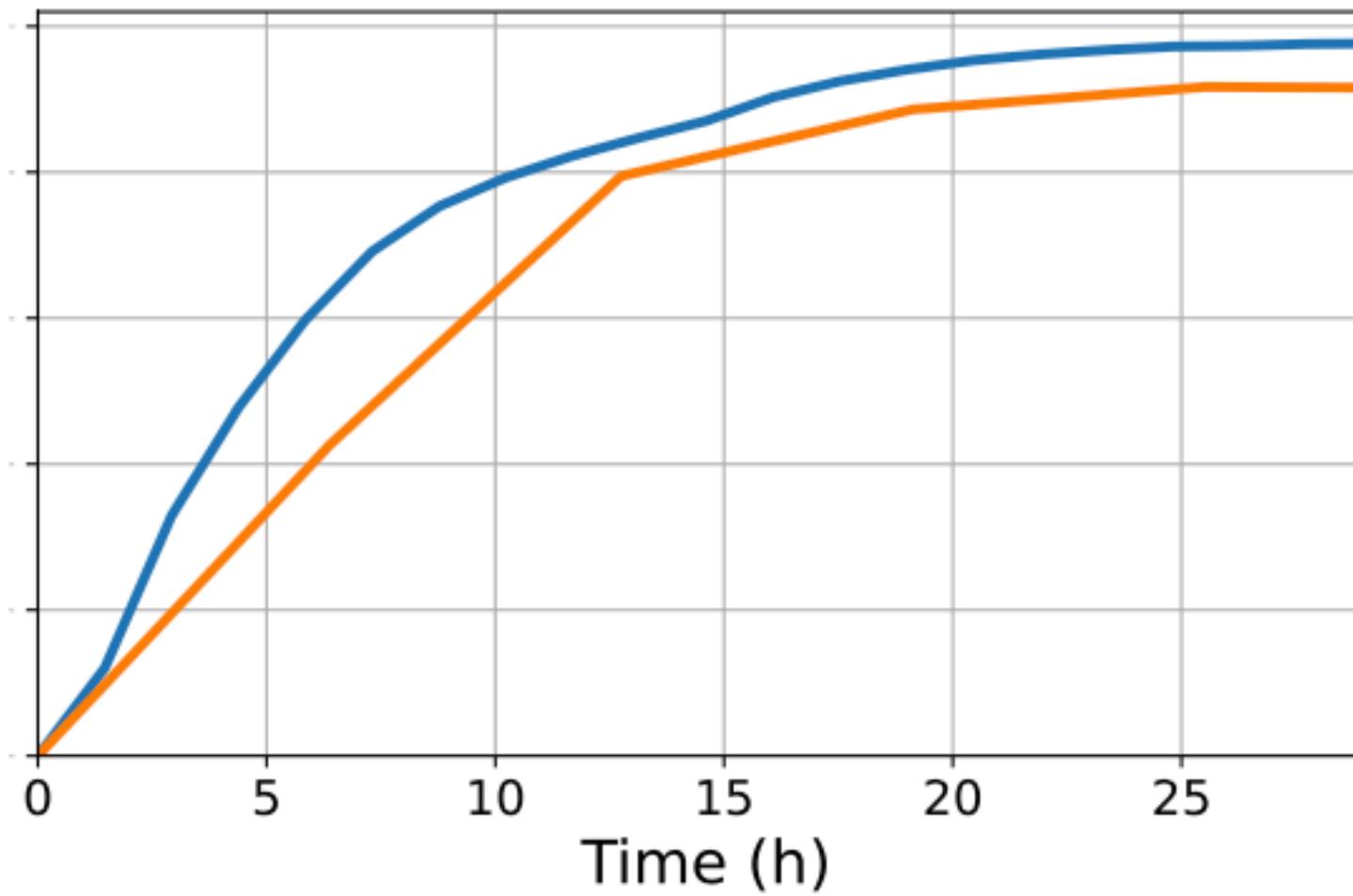
...but it can be cheap!

# Continuous Normalizing Flows for $\phi^4$

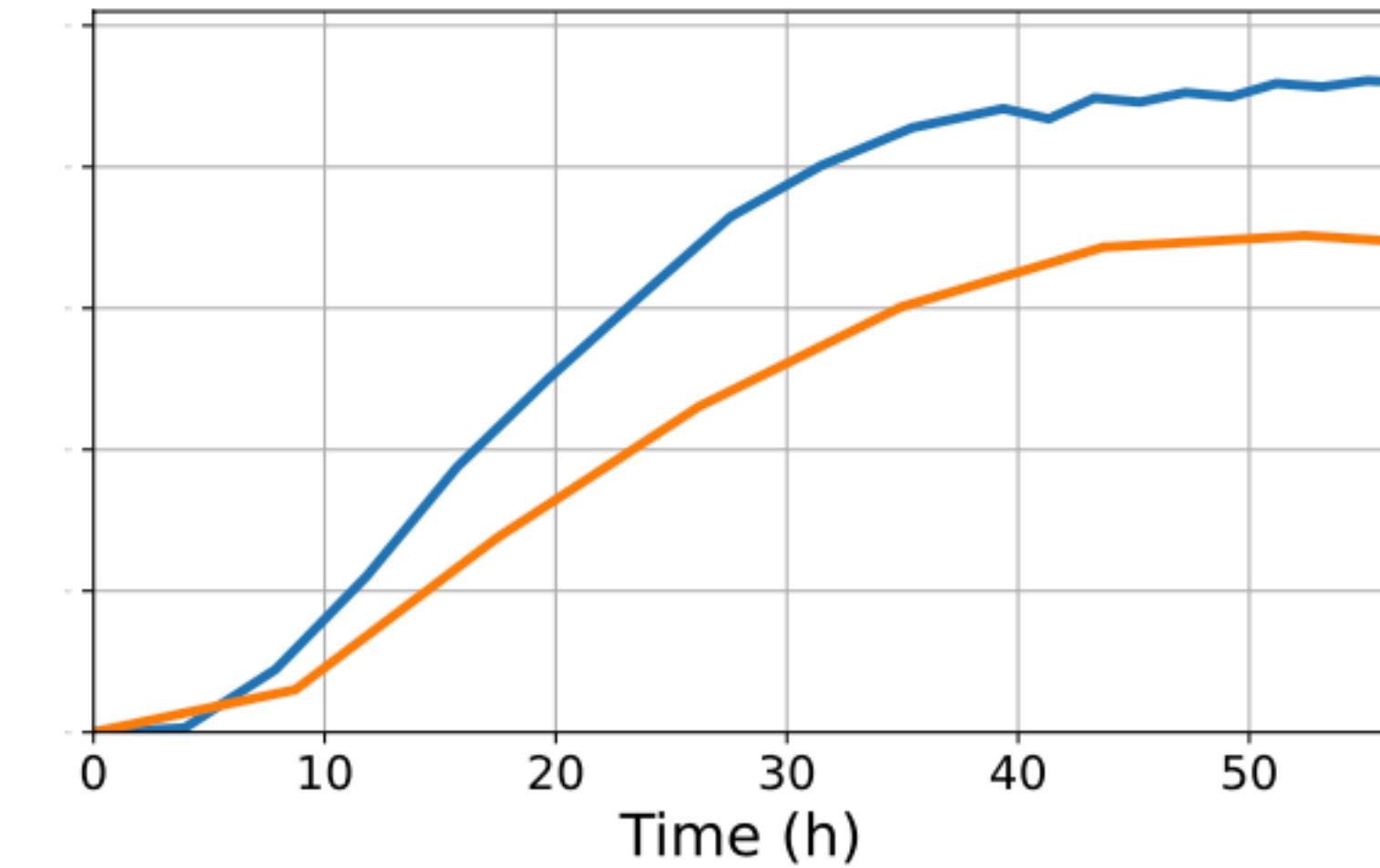
$12 \times 12$



$20 \times 20$



$32 \times 32$



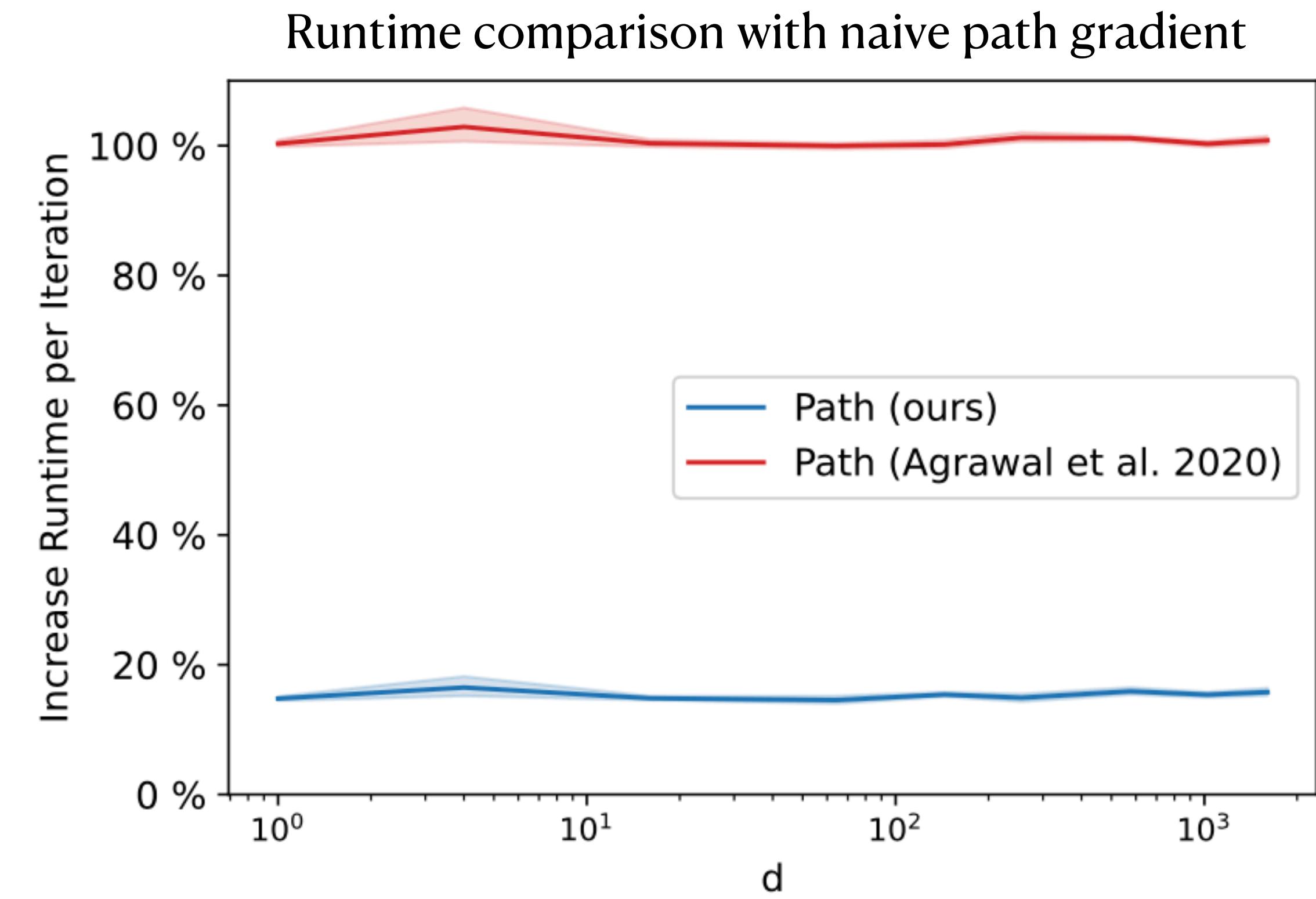
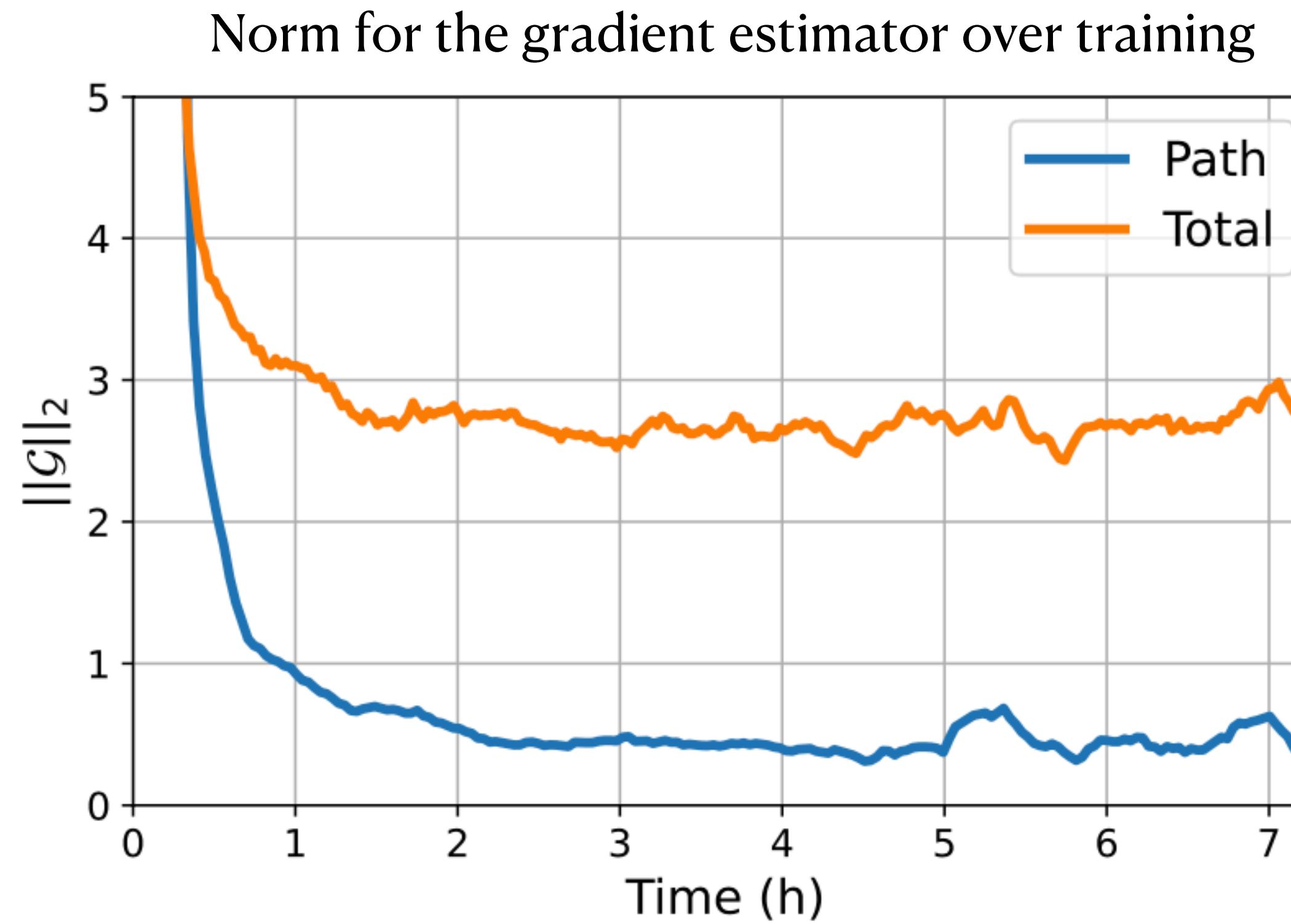
$$\text{ESS} = \frac{N_{\text{eff}}}{N} \approx \frac{\left( \frac{1}{N} \sum_{i=1}^N \tilde{w}_i \right)^2}{\frac{1}{N} \sum_{i=1}^N \tilde{w}_i^2} \in [0,1]$$

LATTICE SIZE	PATH [2]	TOTAL [1]
$12 \times 12$	<b>99.66</b> $\pm$ 0.07	$98.01 \pm 0.44$
$20 \times 20$	<b>97.65</b> $\pm$ 0.14	$91.56 \pm 1.13$
$32 \times 32^\dagger$	<b>91.81</b> $\pm$ 1.32	$69.53 \pm 5.59$

[1] de Haan et al., [NeurIPS workshop \(2021\)](#)

[2] Vaitl L., NKA, Nakajima S., Kessel P., [ICML \(2022\)](#)

# Path Gradient Estimators



LATTICE SIZE	PATH	TOTAL
12x12	<b>99.66</b> $\pm$ 0.07	98.01 $\pm$ 0.44
20x20	<b>97.65</b> $\pm$ 0.14	91.56 $\pm$ 1.13
32x32 <sup>†</sup>	<b>91.81</b> $\pm$ 1.32	69.53 $\pm$ 5.59

LATTICE SIZE	TIME PER ITERATION
12x12	+14% $\pm$ 0
20x20	+12% $\pm$ 1
32x32	+5% $\pm$ 1

# Path Gradient for CNFs: Technical Details

---

Continuous normalizing flows (Chen T.Q 2018) transform a sample  $z_0 \sim q_Z$  from a base-density  $q_Z$

$$\phi \equiv z_T = g_\theta(z_0) = z_0 + \int_0^T dt f_\theta(z_t, t)$$

Under the mild assumptions of the Picard–Lindelöf theorem, the map  $g_\theta$  is bijective.

This map is implemented using a Neural Ordinary Differential Equation (NODE) of the form

$$\frac{dz_t}{dt} = f_\theta(z_t, t)$$

where:

- $f_\theta$  is a (**not necessarily** invertible) neural network with parameters  $\theta$ .
- $z_t$  is the intermediate state at time  $t$ .

# Path Gradient for CNFs: Technical Details

---

The path gradient previously defined now becomes

$$\ln \left| \det \frac{dg_\theta}{dz_0} \right| = \int_0^T \text{tr} \left( \frac{df_\theta(z_t, t)}{dz_t} \right) dt$$

And the log variational density is thus given by

$$\ln q_\theta(\phi) = \ln q_Z(g_\theta^{-1}(\phi)) - \ln \left| \det \frac{dg_\theta}{dz_0} \right|_{z_0=g_\theta^{-1}(\phi)} = \ln q_Z(z_0) - \int_0^T \text{tr} \left( \frac{\partial f_\theta(z_t, t)}{\partial z_t} \right) dt$$

Using the adjoint state method gradients can be computed with

- **Constant** memory requirements
- At the cost of 2 **additional** forward passes.

# Path Gradient for CNFs: Technical Details

The path gradient previously introduced

$$\frac{d}{d\theta} \text{KL}(q_\theta, p) = \mathbb{E}_{z \sim q_Z} [\nabla_\theta S(g_\theta(z)) + \nabla_\theta \log q_\theta(g_\theta(z))] = \mathbb{E}_{z \sim q_Z} \left[ \frac{\partial}{\partial g_\theta(z)} \left( \ln \frac{q_\theta(g_\theta(z))}{p(g_\theta(z))} \right) \frac{\partial g_\theta(z)}{\partial \theta} \right]$$

Now becomes

$$\frac{d}{d\theta} \text{KL}(q_\theta, p) = \left[ \left( \frac{\partial}{\partial g_\theta(z)} \ln q_\theta(g_\theta(z)) - \frac{\partial}{\partial g_\theta(z)} \ln p(g_\theta(z)) \right) \frac{\partial g_\theta(z)}{\partial \theta} \right]$$



$$\frac{\partial}{\partial z_T} \ln q_\theta(z_T)$$


This derivative is non-trivial to compute for  
NFs and requires additional care!

# Path Gradient for CNFs: Technical Details

---

## Theorem

*The derivative  $\frac{\partial \ln q_\theta(z_T)}{\partial z_T}$  can be obtained by solving the initial value problem*

$$\frac{d}{dt} \frac{\partial \ln q_\theta(z_t)}{\partial z_t} = - \frac{\partial \ln q_\theta(z_t)^\top}{\partial z_t} \frac{\partial f_\theta(z_t, t)}{\partial z_t} - \partial_{z_t} \text{tr} \left( \frac{\partial f_\theta(z_t, t)}{\partial z_t} \right),$$

*with initial condition*

$$\frac{\partial \ln q_\theta(z_0)}{\partial z_0} = \frac{\partial \ln q_Z(z_0)}{\partial z_0}.$$

# Path Gradient Estimators

The only **non trivial term** to be computed for the path gradient is

$$\nabla_{\theta} \ln q_{\theta}(g_{\theta}(z)) = \frac{\partial \ln q_{\theta}(z_T)^T}{\partial z_T} \frac{\partial z_T}{\partial \theta}$$

where  $\nabla_{\theta} = \frac{d}{d\theta}$   
and  $z_T \equiv x = g_{\theta}(z_0)$

The **derivative of the log density** wrt the final state is computed using an ODE solved forward in time

- Compute the derivative  $\frac{\partial \ln q_{\theta}(z_T)}{\partial z_T}$  solving **an initial value problem**
- with **initial condition**  $\frac{\partial \ln q_{\theta}(z_0)}{\partial z_0} = \frac{\partial \ln q_Z(z_0)}{\partial z_0}$
- **Adjoint state method** (From CNF, Chen et al. 2018)

- 4 Forward passes**
- 1 Initial value problem
  - 1 Forward integration
  - 2 Adj. state method

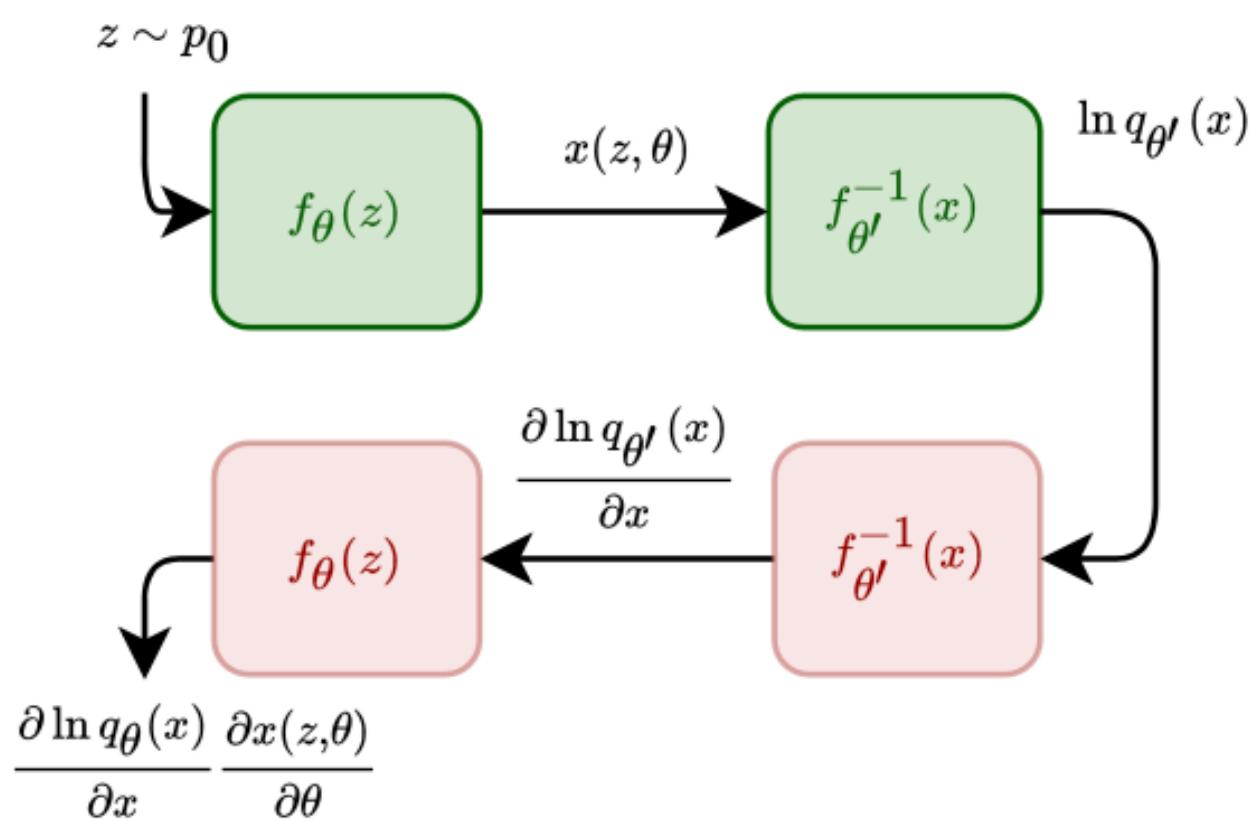
# Path Gradient Estimators

## 4 Forward passes

- 2 for forward mode CNF (adjoint state method for NODEs)
- 2 for Vector Jacobian product for Path gradient

### Naive Path Gradient

(Agrawal et al. 2020)

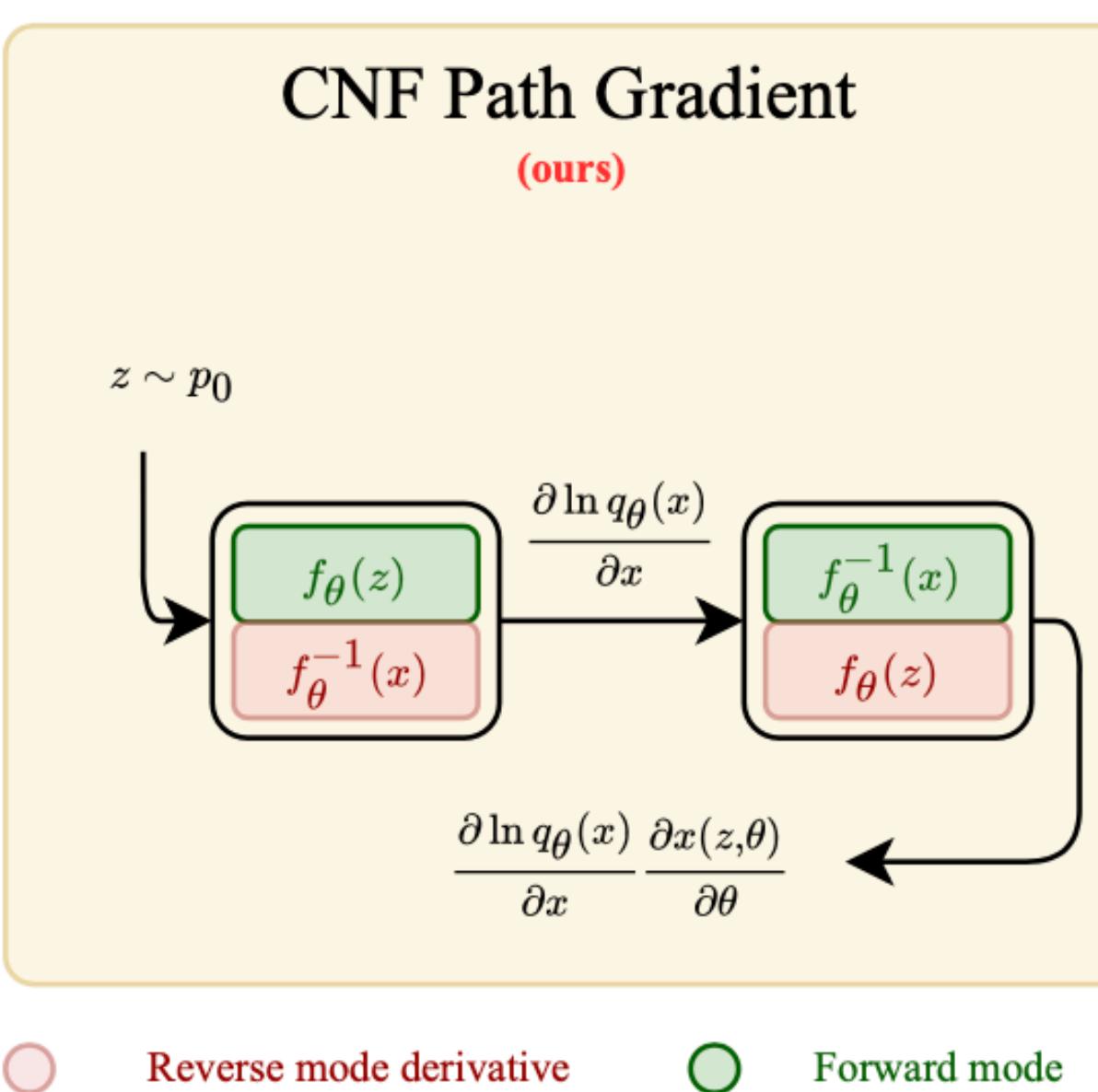


## 6 Forward passes

- Twice the memory footprint
- Considers 2 copy of the model for density estimation and sampling
- 2 forward passes
- 2 backward with the cost of 2 forward passes each

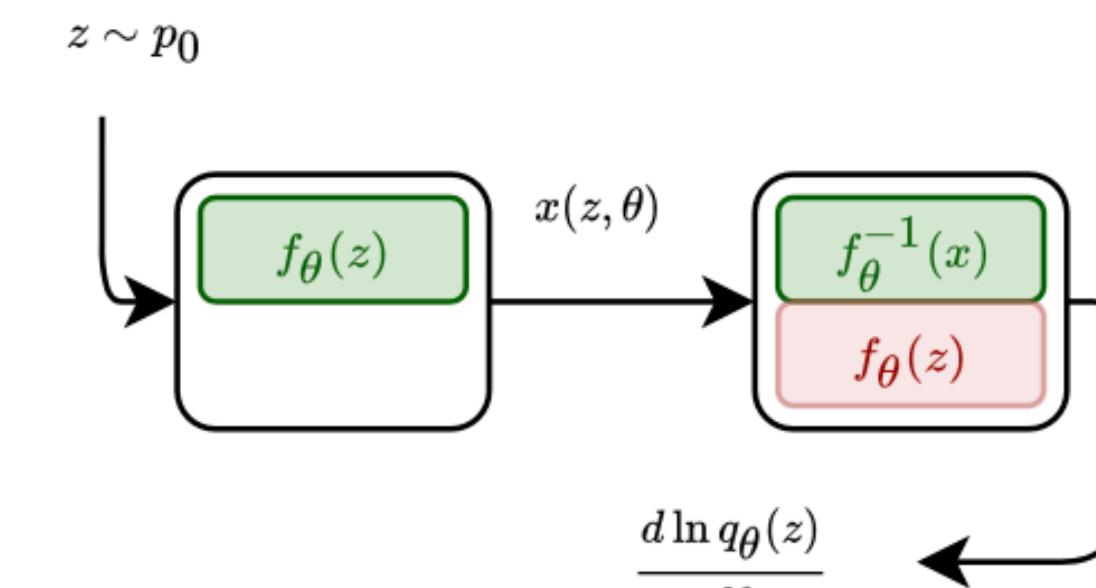
### CNF Path Gradient

(ours)



### CNF Total Gradient

(Chen et al. 2018)



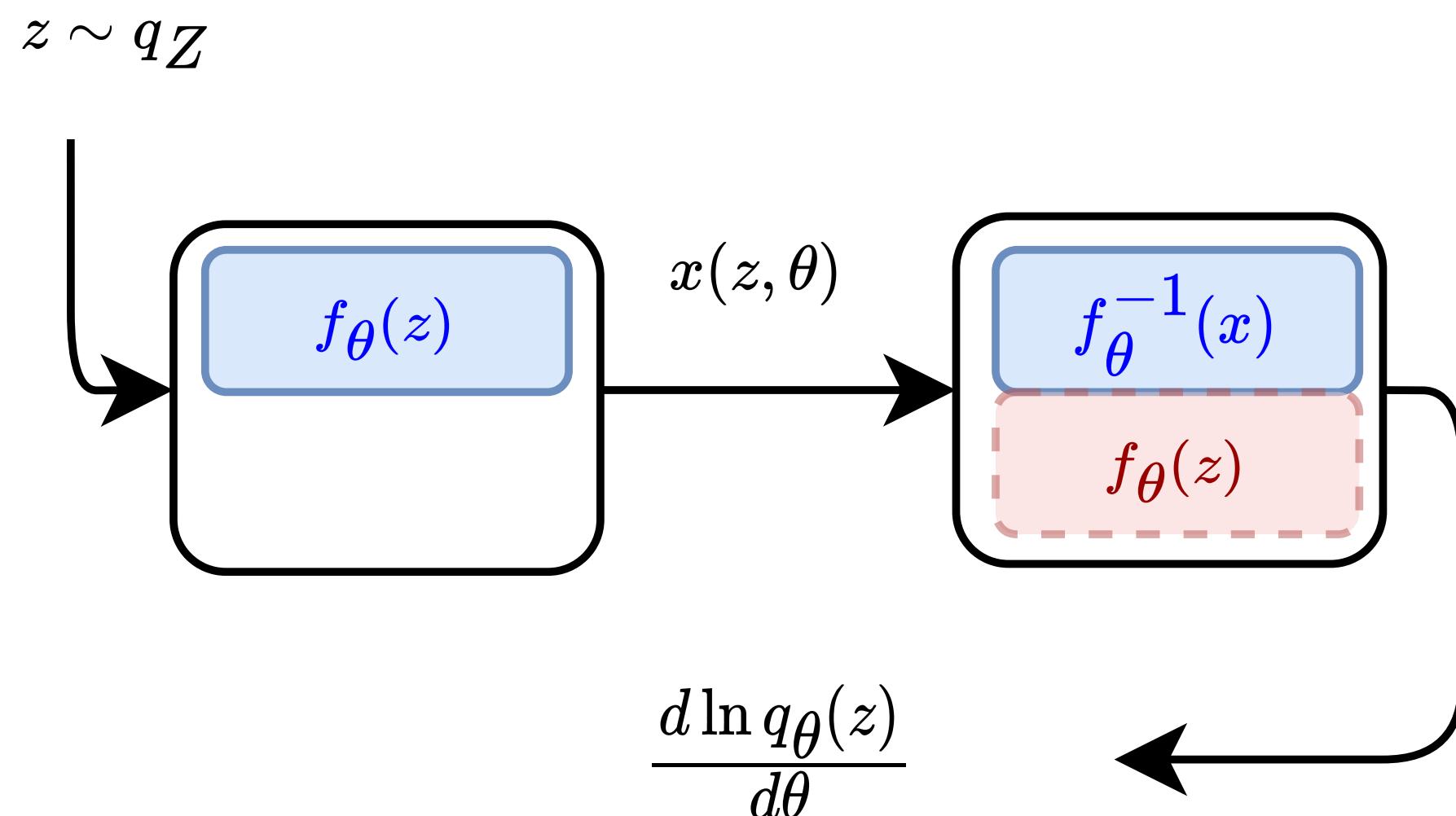
## 3 Forward passes

- 2 for forward mode (adjoint state method for NODEs)
- 1 for Total gradient

# Path Gradient for CNFs: Technical Details

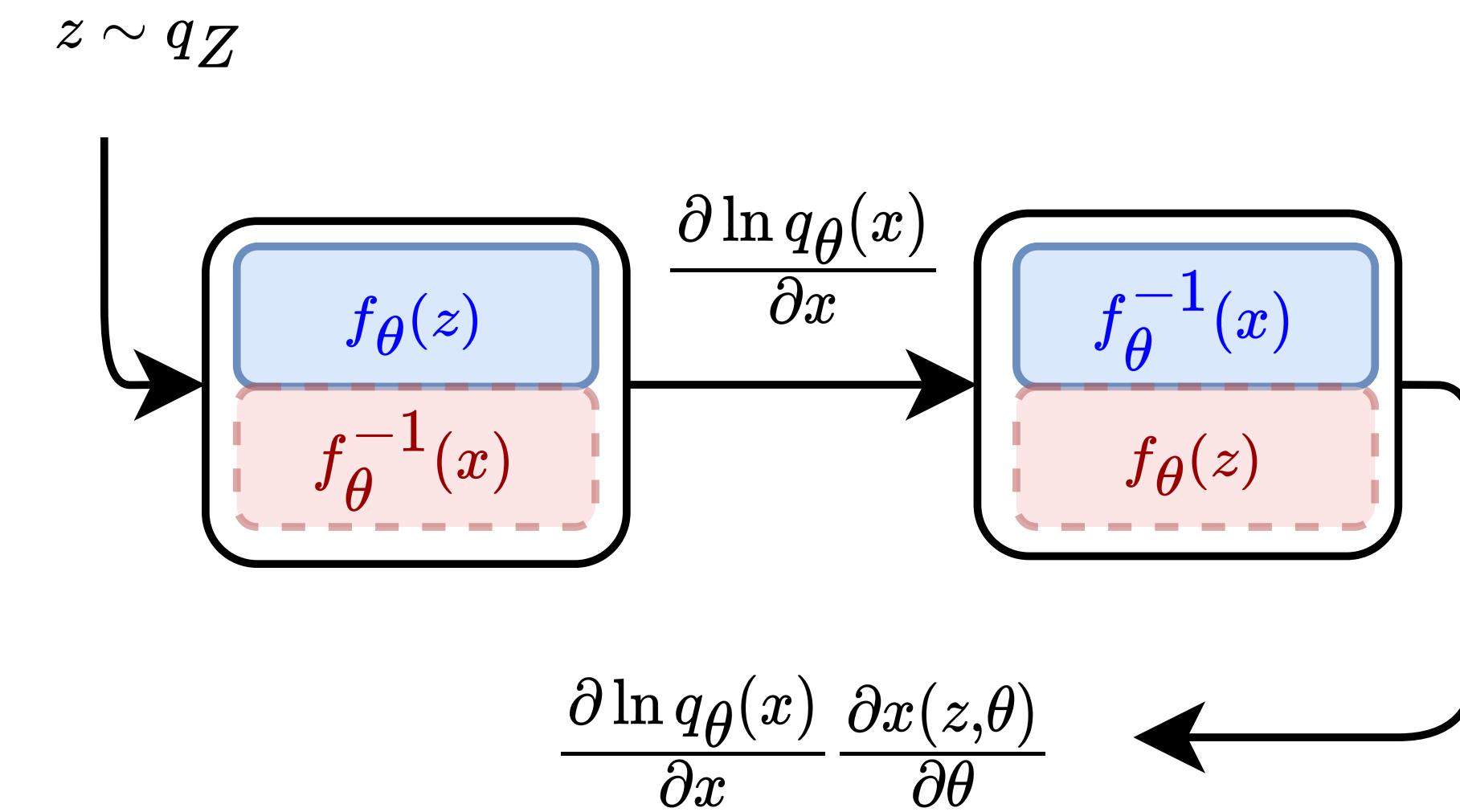
## CNF Total Gradient

(Chen et al. 2018)



## CNF Path Gradient

(ours)



● Reverse mode derivative

● Forward mode

● Reverse mode derivative

● Forward mode

# Path Gradient Estimators

---

- We proposed a path-gradient estimator which can be efficiently implemented for continuous normalizing flows.
- We expect the path-gradient estimator to have slightly higher costs per iteration.
- We find that this effect is more than compensated by better convergence properties of the path-gradient.
- The difference in iteration costs decreases as the size of the lattice increases (quadratic scaling of model params with lattice size).
- A limitation of our work is that so far there is only a theoretical explanation for the lower variance of the path-gradient estimator in the final training phase.
- However, our experiments suggest that the entire training process benefits from the path-gradient estimator. Further investigation is needed.