

Sampling $SU(3)$ pure gauge theory with Stochastic Normalizing Flows

Alessandro Nada

Università degli Studi di Torino

41st International Symposium on Lattice Field Theory

Liverpool, 28th July-3rd August 2024

in collaboration with

Andrea Bulgarelli and Elia Cellini



Long autocorrelation times characterize several observables when $a \rightarrow 0$

Typical example are **topological observables**: for $a \rightarrow 0$ sectors characterized by different values of the topological charge Q emerge

Using standard MCMC algorithms the transition between these sectors is strongly suppressed

Long autocorrelation times characterize several observables when $a \rightarrow 0$

Typical example are **topological observables**: for $a \rightarrow 0$ sectors characterized by different values of the topological charge Q emerge

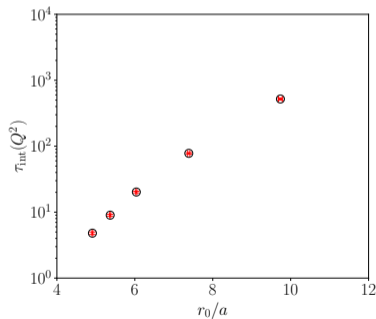
Using standard MCMC algorithms the transition between these sectors is strongly suppressed

This talk: focus on $SU(3)$ in 4 dimensions

Update algorithm of choice: 1 heat-bath step + 4 over-relaxation steps

Objective: mitigate freezing at $\beta = 6.5$ ($r_0/a \sim 11$)

$$\tau_{\text{int}}(Q^2) \sim 10^3$$



[Plot courtesy of C. Bonanno]

Flow-based approach

mapping between the target $p(\phi)$ and some tractable distribution $q_0(z)$

→ novel approach to fight critical slowing down

Lots of progress in Normalizing Flows in the last 5 years!

→ see **R. Abbott's** and **K. Javad's** talks in this session
+ **Tej's plenary from Lattice23**

However: NFs do not appear to scale well with the volume (i.e. with the degrees of freedom)

But: same approach is possible stochastically! → better scaling?

Out-of-equilibrium Monte Carlo evolutions

Out-of-equilibrium evolutions

sampling each consecutive step from a sequence of distributions

$$q_0 \simeq e^{-S_c(0)} \rightarrow e^{-S_c(1)} \rightarrow \dots \rightarrow p \simeq e^{-S_c(n_{\text{step}})}$$

Out-of-equilibrium evolutions

sampling each consecutive step from a sequence of distributions

$$q_0 \simeq e^{-S_{c(0)}} \rightarrow e^{-S_{c(1)}} \rightarrow \dots \rightarrow p \simeq e^{-S_{c(n_{\text{step}})}}$$

- ▶ $c(n)$ is a parameter of the action $S_{c(n)}$ of the model
- ▶ start **at equilibrium** from a distribution $q_0 = e^{-S_{c(0)}}/Z_0$, the **prior**
- ▶ n_{step} intermediate steps
- ▶ at each step: MC update with transition probability $P_{c(n)}(U_n \rightarrow U_{n+1})$
- ▶ $P_{c(n)}$ changes along the evolution according to the **protocol** $c(n)$
- ▶ end at the **target** probability distribution $p = e^{-S_{c(n_{\text{step}})}}/Z_{n_{\text{step}}} \equiv e^{-S}/Z$

Out-of-equilibrium evolutions

sampling each consecutive step from a sequence of distributions

$$q_0 \simeq e^{-S_{c(0)}} \rightarrow e^{-S_{c(1)}} \rightarrow \dots \rightarrow p \simeq e^{-S_{c(n_{\text{step}})}}$$

- ▶ $c(n)$ is a parameter of the action $S_{c(n)}$ of the model
- ▶ start **at equilibrium** from a distribution $q_0 = e^{-S_{c(0)}}/Z_0$, the **prior**
- ▶ n_{step} intermediate steps
- ▶ at each step: MC update with transition probability $P_{c(n)}(U_n \rightarrow U_{n+1})$
- ▶ $P_{c(n)}$ changes along the evolution according to the **protocol** $c(n)$
- ▶ end at the **target** probability distribution $p = e^{-S_{c(n_{\text{step}})}}/Z_{n_{\text{step}}} \equiv e^{-S}/Z$

"forward" transition probability

$$\mathcal{P}_f[U_0, \dots, U] = \prod_{n=1}^{n_{\text{step}}} P_{c(n)}(U_{n-1} \rightarrow U_n)$$

Crooks' theorem for MCMC [Crooks; 1999]: if the update algorithm satisfies detailed balance

$$\frac{q_0(U_0) \mathcal{P}_f[U_0, \dots, U_{n_{\text{step}}}]}{p(U) \mathcal{P}_r[U_{n_{\text{step}}}, \dots, U_0]} = \frac{q_0(U_0) \prod_{n=1}^{n_{\text{step}}} P_{c(n)}(U_{n-1} \rightarrow U_n)}{p(U_{n_{\text{step}}}) \prod_{n=1}^{n_{\text{step}}} P_{c(n)}(U_n \rightarrow U_{n-1})} = \exp(W - \Delta F)$$

Crooks' theorem for MCMC [Crooks; 1999]: if the update algorithm satisfies detailed balance

$$\frac{q_0(U_0)\mathcal{P}_f[U_0, \dots, U_{n_{\text{step}}}]}{p(U)\mathcal{P}_r[U_{n_{\text{step}}}, \dots, U_0]} = \frac{q_0(U_0) \prod_{n=1}^{n_{\text{step}}} P_{c(n)}(U_{n-1} \rightarrow U_n)}{p(U_{n_{\text{step}}}) \prod_{n=1}^{n_{\text{step}}} P_{c(n)}(U_n \rightarrow U_{n-1})} = \exp(W - \Delta F)$$

with the generalized **work**

$$W = \sum_{n=0}^{n_{\text{step}}-1} \{S_{c(n+1)}[U_n] - S_{c(n)}[U_n]\}$$

and the **free energy** difference

$$\exp(-\Delta F) = \frac{Z_{c(n_{\text{step}})}}{Z_{c(0)}}$$

Jarzynski's equality for MCMC

Integrating over all paths gives

$$\int [dU_0 \dots dU_{n_{\text{step}}}] q_0(U_0) \mathcal{P}_f[U_0, \dots, U_{n_{\text{step}}}] \exp(-(W - \Delta F)) = 1$$

Formal derivation of **Jarzynski's equality** [Jarzynski; 1997] for MCMC

$$\langle \exp(-W) \rangle_f = \exp(-\Delta F) = \frac{Z}{Z_0}$$

Ratio of partition functions computed directly with an **average** over "forward" non-equilibrium evolutions

→ see talk by **A. Bulgarelli** (Tue 14:35)

Integrating over all paths gives

$$\int [dU_0 \dots dU_{n_{\text{step}}}] q_0(U_0) \mathcal{P}_f[U_0, \dots, U_{n_{\text{step}}}] \exp(-(W - \Delta F)) = 1$$

Formal derivation of Jarzynski's equality [Jarzynski; 1997] for MCMC

$$\langle \exp(-W) \rangle_f = \exp(-\Delta F) = \frac{Z}{Z_0}$$

Ratio of partition functions computed directly with an **average** over "forward" non-equilibrium evolutions

→ see talk by A. Bulgarelli (Tue 14:35)

Using Jensen's inequality $\langle \exp x \rangle \geq \exp \langle x \rangle$

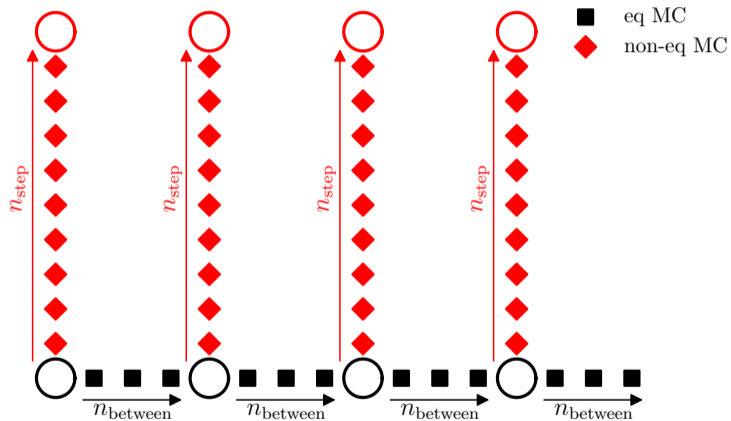
$$\langle W \rangle_f \geq \Delta F$$

→ **Second Law of Thermodynamics**

The same derivation holds if you want to compute v.e.v. of an observable for the target distribution p

$$\langle \mathcal{O} \rangle = \frac{\langle \mathcal{O} \exp(-W) \rangle_f}{\langle \exp(-W) \rangle_f} = \langle \mathcal{O} \exp(-W_d) \rangle_f$$

A non-equilibrium paradigm to perform MCMC



How far are we from equilibrium?

However we can measure the **similarity of forward and reverse processes**

$$\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \int [dU_0 \dots dU] q_0(U_0) \mathcal{P}_f[U_0, \dots, U] \log \frac{q_0(U_0) \mathcal{P}_f[U_0, \dots, U]}{p(U) \mathcal{P}_r[U, U_{n_{\text{step}}-1}, \dots, U_0]}$$

How far are we from equilibrium?

However we can measure the **similarity of forward and reverse processes**

$$\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \int [dU_0 \dots dU] q_0(U_0) \mathcal{P}_f[U_0, \dots, U] \log \frac{q_0(U_0) \mathcal{P}_f[U_0, \dots, U]}{p(U) \mathcal{P}_r[U, U_{n_{\text{step}}-1}, \dots, U_0]}$$

Clear "thermodynamic" interpretation

$$\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \langle W \rangle_f + \log \frac{Z}{Z_0} = \underbrace{\langle W \rangle_f - \Delta F}_{\text{Second Law of Thermodynamics!}} \geq 0$$

→ measure of how reversible the process is!

Upper bound for the divergence used for NFs

$$\tilde{D}_{\text{KL}}(q \| p) \leq \tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r)$$

How far are we from equilibrium?

However we can measure the **similarity of forward and reverse processes**

$$\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \int [dU_0 \dots dU] q_0(U_0) \mathcal{P}_f[U_0, \dots, U] \log \frac{q_0(U_0) \mathcal{P}_f[U_0, \dots, U]}{p(U) \mathcal{P}_r[U, U_{n_{\text{step}}-1}, \dots, U_0]}$$

Clear "thermodynamic" interpretation

$$\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \langle W \rangle_f + \log \frac{Z}{Z_0} = \underbrace{\langle W \rangle_f - \Delta F}_{\text{Second Law of Thermodynamics!}} \geq 0$$

→ measure of how reversible the process is!

Upper bound for the divergence used for NFs

$$\tilde{D}_{\text{KL}}(q \| p) \leq \tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r)$$

Another figure of merit is the Effective Sample Size

$$\text{ESS} = \frac{\langle \exp(-W) \rangle_f^2}{\langle \exp(-2W) \rangle_f}$$

Out-of-equilibrium evolutions in β for SU(3) in 4 dimensions

Evolution from thermalized MC at β_0 to a target β

$$q_0 \simeq e^{-S_{\beta(0)}} \rightarrow e^{-S_{\beta(1)}} \rightarrow \dots \rightarrow p \simeq e^{-S_{\beta(n_{\text{step}})}}$$

Objectives

- ▶ Analyze scaling with volume $(L/a)^4$
- ▶ Set MCMC standard for flow-based approach
- ▶ No topology yet (charge not frozen yet)

Setup

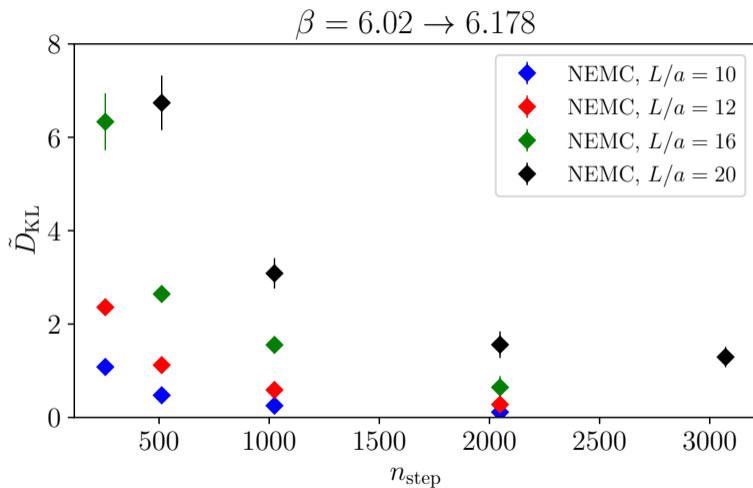
- ▶ Increasingly large lattices, from $L/a = 10$ to $L/a = 20$
- ▶ "Jump" in β :

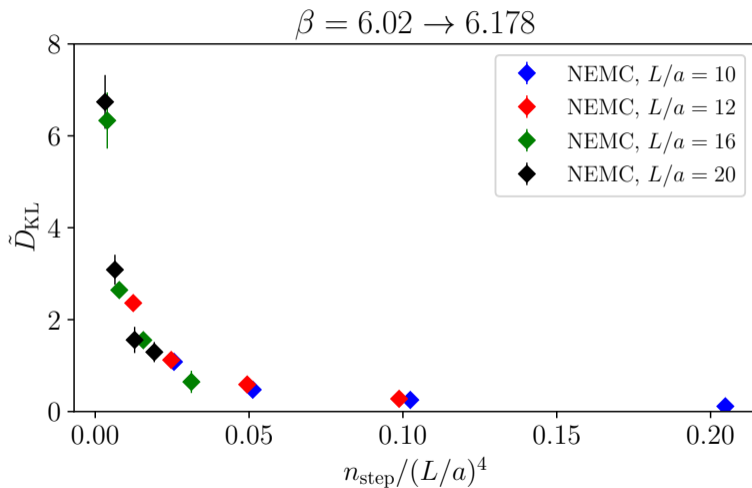
$$6.02 \rightarrow 6.178$$

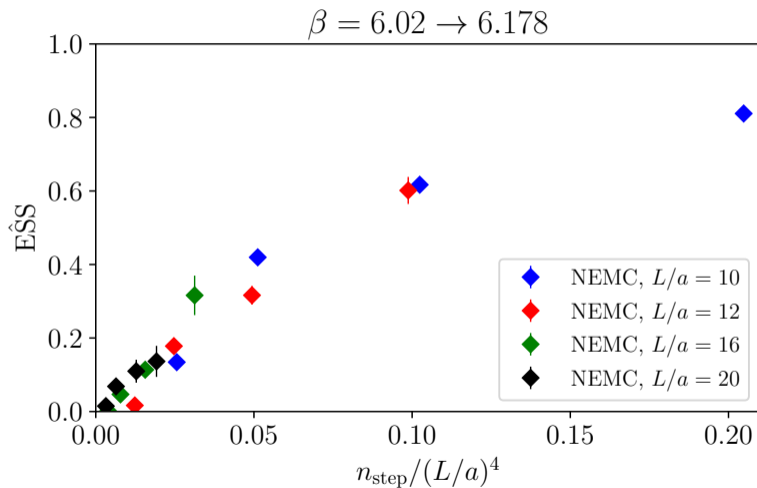
corresponding to $(1.8\text{fm})^4 \rightarrow (1.4\text{fm})^4$ for $L/a = 20$

This work: inverse coupling increased linearly

$$\beta(n) = \beta_0 + (\beta - \beta_0) \frac{n}{n_{\text{step}}}$$



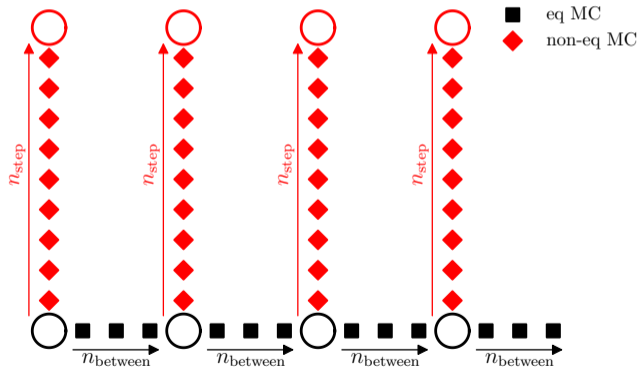




Stochastic Normalizing Flows

SNFs as systematic improvement of non-equilibrium evolutions

What if you introduce the same transformations used in NFs **between** the non-equilibrium Monte Carlo updates?

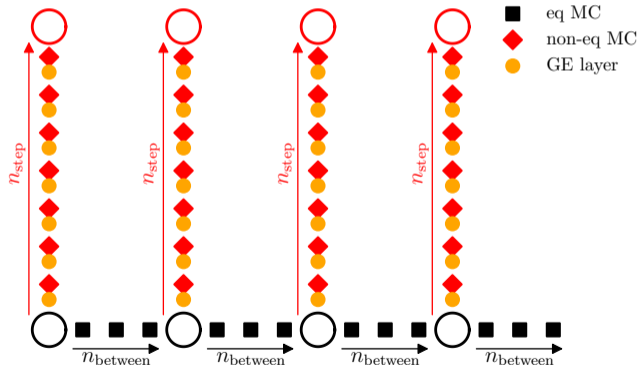


SNFs as systematic improvement of non-equilibrium evolutions

What if you introduce the same transformations used in NFs **between** the non-equilibrium Monte Carlo updates?

Stochastic Normalizing Flows (introduced in [Wu et al.; 2020])

$$U_0 \xrightarrow{g_1} g_1(U_0) \xrightarrow{P_{c(1)}} U_1 \xrightarrow{g_2} g_2(U_1) \xrightarrow{P_{c(2)}} U_2 \xrightarrow{g_3} \dots \xrightarrow{P_{c(n_{\text{step}})}} U_{n_{\text{step}}}$$



What if you introduce the same transformations used in NFs **between** the non-equilibrium Monte Carlo updates?

Stochastic Normalizing Flows (introduced in [Wu et al.; 2020])

$$U_0 \xrightarrow{g_1} g_1(U_0) \xrightarrow{P_{c(1)}} U_1 \xrightarrow{g_2} g_2(U_1) \xrightarrow{P_{c(2)}} U_2 \xrightarrow{g_3} \dots \xrightarrow{P_{c(n_{\text{step}})}} U_{n_{\text{step}}}$$

The (generalized) work now is

$$W = \sum_{n=0}^{n_{\text{step}}-1} \underbrace{S_{c(n+1)}(g_n(U_n)) - S_{c(n)}(g_n(U_n))}_{\text{stochastic}} - \underbrace{\log |\det J_n(U_n)|}_{\text{deterministic}}$$

- ▶ use gauge-equivariant layers to effectively decrease n_{step}
- ▶ how to do training? advantages from the architecture
- ▶ same scaling with the volume?

Implementation of the coupling layers introduced in [Nagai and Tomiya; 2021] and the link-level flow used in [Abbott et al.; 2023]

Essentially a stout-smearing transformation [Morningstar and Peardon; 2003] with masks to make it invertible (and compute $\log J$)

$$U'_\mu(x) = g_l(U_\mu(x)) = \exp\left(Q_\mu^{(l)}(x)\right) U_\mu(x)$$

with the algebra-valued

$$Q_\mu^{(l)}(x) = 2 \left[\Omega_\mu^{(l)}(x) \right]_{\text{TA}}$$

$$\Omega_\mu(x) = \underbrace{C_\mu(x)}_{\text{frozen}} \underbrace{U_\mu^\dagger(x)}_{\text{active}}$$

Sum of frozen staples

$$C_\mu(x) = \sum_{\nu \neq \mu} \rho \underbrace{S_{\mu\nu}(x)}_{\text{staple}}$$

Architecture: (1 gauge-equivariant + 1 full MC update) $\times n_{\text{step}}$

Training: minimizing $\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \langle W \rangle_f + \text{const}$

Architecture: (1 gauge-equivariant + 1 full MC update) $\times n_{\text{step}}$

Training: minimizing $\tilde{D}_{\text{KL}}(q_0 \mathcal{P}_f \| p \mathcal{P}_r) = \langle W \rangle_f + \text{const}$

Short trainings: 200-1000 epochs

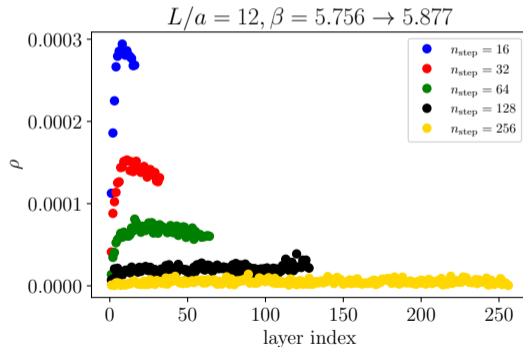
Memory issues for large n_{step} and large volumes

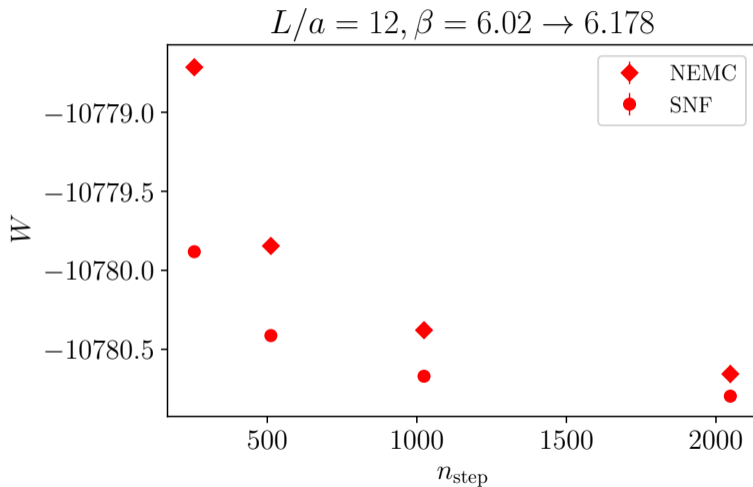
Practical solution: train each layer separately during the non-equilibrium evolution \rightarrow reminiscent of CRAFT [Matthews et al.; 2022]

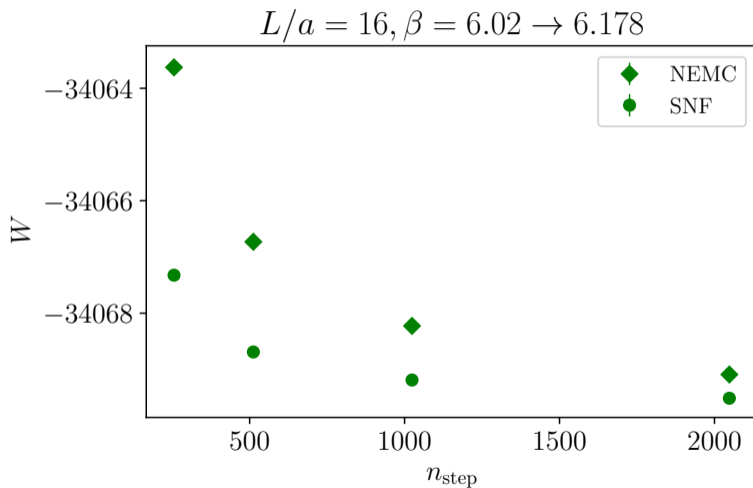
Heavy use of **transfer learning** for each $\beta_0 \rightarrow \beta$ evolution:

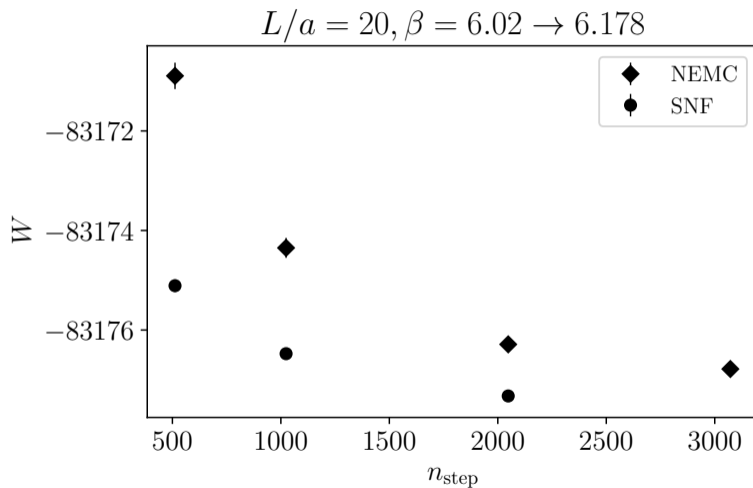
- ▶ training only at small volumes
- ▶ training only with small n_{step} : global interpolation of ρ

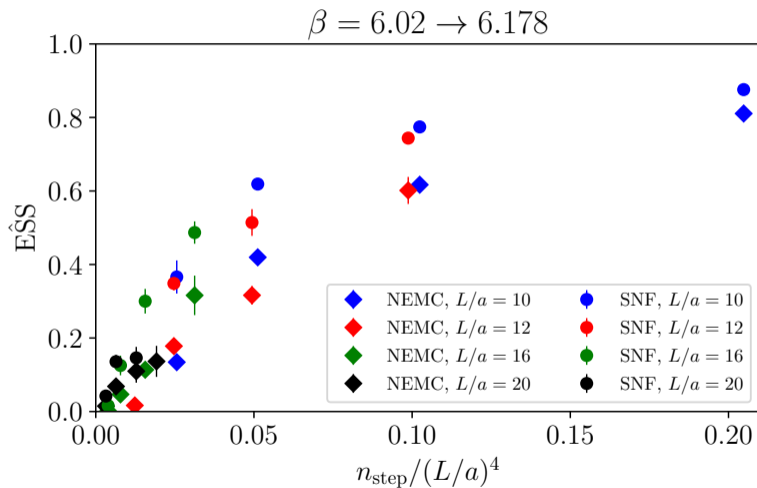
No retraining!











Stochastic approach guarantees a **clear scaling** with the degrees of freedom

$$n_{\text{step}} \sim \text{d.o.f.} \rightarrow \text{fixed } \tilde{D}_{\text{KL}} \text{ or ESS}$$

while providing a thermodynamic understanding of the flow

Stochastic approach guarantees a **clear scaling** with the degrees of freedom

$$n_{\text{step}} \sim \text{d.o.f.} \rightarrow \text{fixed } \tilde{D}_{\text{KL}} \text{ or ESS}$$

while providing a thermodynamic understanding of the flow

Overall strategy

systematically improve on stochastic approach by machine-learning deterministic transformations between MC steps

Future improvements

Better **protocols** (huge literature from non-eq SM):
only linear protocols were used in this work!

Better and deeper **layers**: include larger loops beyond
the plaquette + ρ as a neural network [Abbott et al.;
2023]

Future implementations

Implement **SNF for evolutions in the BC**

→ see poster by D. VDACCHINO

Push **SNFs/evolutions in β** at finer lattice spacings

Thank you for your attention!

Several applications in the last 8 years!

- ▶ calculation of the interface free-energy in the Z_2 gauge theory [Caselle et al.; 2016]
- ▶ $SU(3)$ pure gauge equation of state in 4d from the pressure [Caselle et al.; 2018]
- ▶ renormalized coupling for $SU(N)$ YM theories [Francesconi et al.; 2020]
- ▶ entanglement entropy [Bulgarelli and Panero; 2023]
- ▶ connection with Stochastic Normalizing Flows: ϕ^4 scalar field theory [Caselle et al.; 2022] and Nambu-Goto effective string model [Caselle et al.; 2023]
- ▶ Topological unfreezing for $CP(N - 1)$ model [Bonanno et al.; 2023]

The effective sample size

Effective Sample Size: defined in general as the ratio between the "theoretical" variance and the actual variance of the NE observable

$$\frac{\text{Var}(\mathcal{O})_{\text{NE}}}{n} = \frac{\text{Var}(\mathcal{O})_p}{n \text{ESS}}$$

but difficult to compute

We use the (customary) approximate estimator

$$\text{E}\hat{\text{SS}} = \frac{\langle \exp(-W) \rangle_f^2}{\langle \exp(-2W) \rangle_f} = \frac{1}{\langle \exp(-2W_d) \rangle_f}$$

Easy to understand in terms of the variance of $\exp(-W)$:

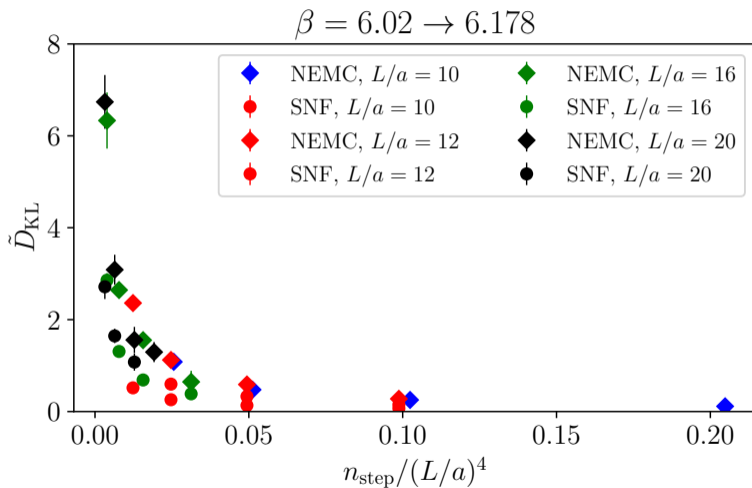
$$\text{Var}(\exp(-W)) = \left(\frac{1}{\text{E}\hat{\text{SS}}} - 1 \right) \exp(-2\Delta F) \geq 0$$

which leads to

$$0 < \text{E}\hat{\text{SS}} \leq 1$$

How to sample frozen topological observables at β_{target} on a L^4 lattice?

	Evolution in the boundary conditions	Evolution in β
Prior	thermalized Markov Chain at β_{target} with OBC on a L_d^3 defect	thermalized Markov Chain at $\beta_0 < \beta_{\text{target}}$ ($a_0 > a_{\text{target}}$)
Protocol	Gradually switch on PBC	Gradually increase β (compress the volume)
d.o.f.	$\sim (L_d/a)^3$	$\sim (L/a)^4$
Intermediate sampling	—	possible at any intermediate β



The Second Law of Thermodynamics

Clausius inequality for an (isothermal) transformation from state A to state B

$$\frac{Q}{T} \leq \Delta S$$

If we use

$$\begin{cases} Q = \Delta E - W & \text{(First Law)} \\ F \stackrel{\text{def}}{=} E - ST \end{cases}$$

the Second Law becomes

$$W \geq \Delta F$$

where the equality holds for reversible processes.

Moving from thermodynamics to statistical mechanics we know that actually

$$\langle W \rangle_f \geq \Delta F = F_B - F_A$$

for a given "forward" process f from A to B

A connection to traditional reweighting

A typical reweighting procedure is meant to sample a distribution p using a (close enough) distribution q_0 . It can be written as

$$\langle \mathcal{O} \rangle_{\text{RW}} = \frac{\langle \mathcal{O}(\phi) \exp(-\Delta S) \rangle_{q_0}}{\langle \exp(-\Delta S) \rangle_{q_0}}$$

It is just Jarzynski's equality for $n_{\text{step}} = 1$, see the work

$$W = \sum_{n=0}^{n_{\text{step}}-1} \{S_{c(n+1)}[\phi_n] - S_{c(n)}[\phi_n]\} = \Delta S(\phi_0)$$

with ϕ_0 sampled from q_0

- ▶ It's important to note that there is no issue with the fact that ΔS itself can be large
- ▶ The real issue is that the *distribution* of ΔS (and in general of W) can lead to an extremely poor estimate of $\Delta F \rightarrow$ highly inefficient sampling
- ▶ The exponential average can be tricky when very far from equilibrium!

A common framework: Stochastic Normalizing Flows

Jarzynski's equality is the same formula used to extract Z in NFs

$$\frac{Z}{Z_0} = \langle \tilde{w}(\phi) \rangle_{\phi \sim q_N} = \langle \exp(-W) \rangle_f$$

The exponent of the weight is always of the form

(note that for NFs $\langle \dots \rangle_{\phi \sim q_N} = \langle \dots \rangle_f$)

$$W(\phi_0, \dots, \phi_N) = S(\phi_N) - S_0(\phi_0) - Q(\phi_1, \dots, \phi_N)$$

Normalizing Flows

$$\phi_0 \rightarrow \phi_1 = g_1(\phi_0) \rightarrow \dots \rightarrow \phi_N$$

$$\text{"Q"} = \log J = \sum_{n=0}^{N-1} \log |\det J_n(\phi_n)|$$

stochastic non-equilibrium evolutions

$$\phi_0 \xrightarrow{P_{c(1)}} \phi_1 \xrightarrow{P_{c(2)}} \dots \xrightarrow{P_{c(N)}} \phi_N$$

$$Q = \sum_{n=0}^{N-1} S_{c(n+1)}(\phi_{n+1}) - S_{c(n+1)}(\phi_n)$$

Stochastic Normalizing Flows (introduced in [Wu et al.; 2020])

$$\phi_0 \rightarrow g_1(\phi_0) \xrightarrow{P_{c(1)}} \phi_1 \rightarrow g_2(\phi_1) \xrightarrow{P_{c(2)}} \dots \xrightarrow{P_{c(N)}} \phi_N$$

$$Q = \sum_{n=0}^{N-1} S_{c(n+1)}(\phi_{n+1}) - S_{c(n+1)}(g_n(\phi_n)) + \log |\det J_n(\phi_n)|$$

Large-scale application: computation of the SU(3) equation of state [Caselle et al.; 2018]

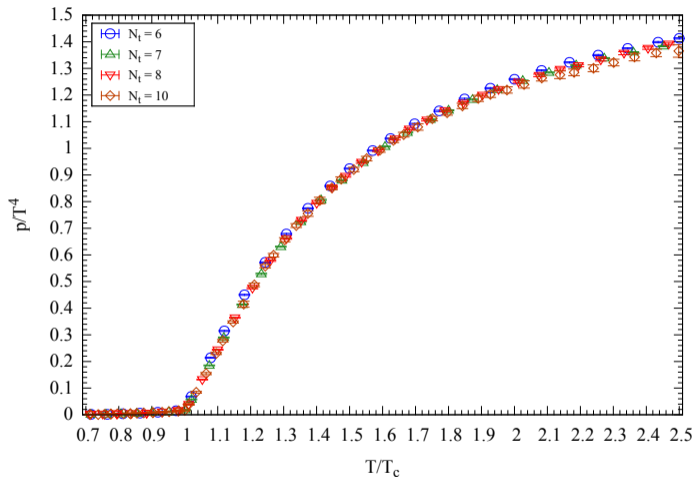
Goal: extract the pressure with Jarzynski's equality

$$\frac{p(T)}{T^4} - \frac{p(T_0)}{T_0^4} = \left(\frac{N_t}{N_s}\right)^3 \log \langle e^{-W_{\text{SU}(N_c)}} \rangle_f$$

evolution in β_g (inverse coupling) \rightarrow changes lattice spacing $a \rightarrow$ changes temperature $T = 1/(aN_t)$

Prior: thermalized Markov chain at a certain $\beta_g^{(0)}$

For systems with many d.o.f. (i.e. large volumes), JE works when N is large, i.e. evolution is slow (and expensive)



Large volumes (up to $160^3 \times 10$) and very fine lattice spacings $\beta \simeq 7$