

Self-learning Monte Carlo method with equivariant transformer

The Information Technology Center,
The University of Tokyo



Yuki Nagai

Yuki Nagai and Akio Tomiya, "Self-Learning Monte Carlo with Equivariant Transformer", J. Phys. Soc. Jpn. 93, 114007 (2024) **Editors' choice**

Introduction

Speedup with machine learning

In field of machine learning

Image recognition, AI chat etc.

We do not have a theory of these. But the machine can imitate these



How to use machine learning in simulations?

Known heavy task is replaced



We replace the heavy tasks by neural networks

Self learning Monte Carlo

Exact and effective models

If you want to do simulations with very heavy computational cost,
the effective model to imitate the original model might be useful



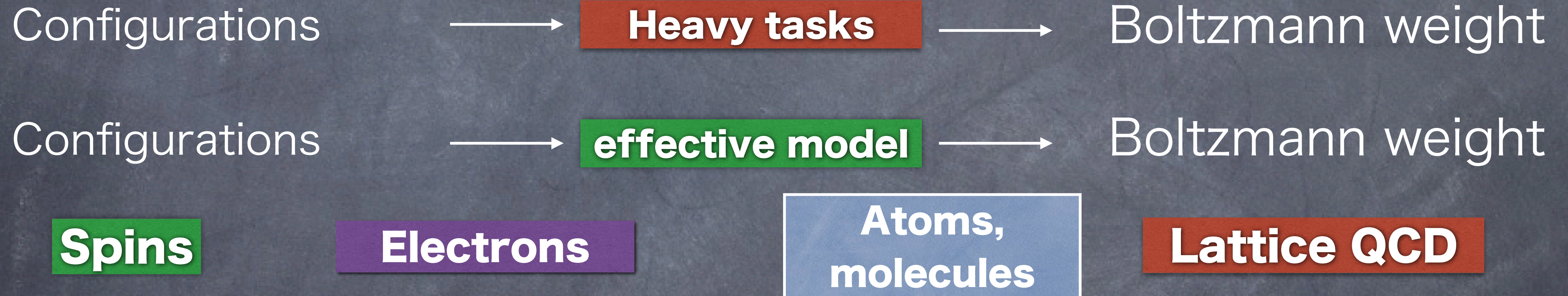
1. If the effective model is not good, B' is not good
2. How long do you have to train the model?

By using the self-learning Monte Carlo method,
the output with an effective model becomes **exact**

Self-learning Monte Carlo

We focus on Markov-chain Monte Carlo method

We usually calculate physical observables with a partition function $Z = \int \exp(-S)$ or $\sum \exp(-\beta H)$



In MCMC, there are two parts

1. proposing new configuration
2. evaluating Boltzmann weight

To propose a new configuration, we use the effective model

Self-learning Monte Carlo

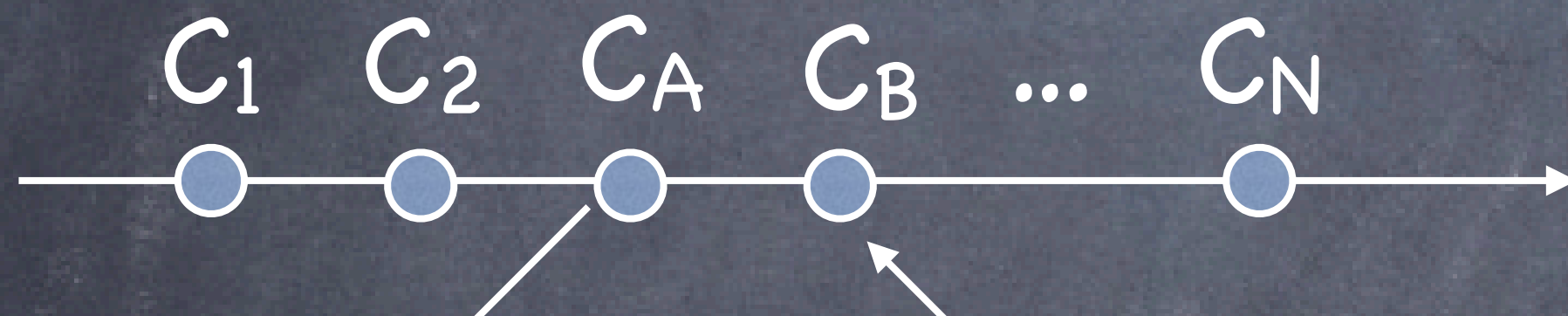
Self-learning Monte Carlo method (SLMC)

Self-learning Hybrid Monte Carlo method (SLHMC)

To speed up the Markov Chain Monte Carlo (MCMC) simulations

SLMC

Markov chain with the probability $W(C)$



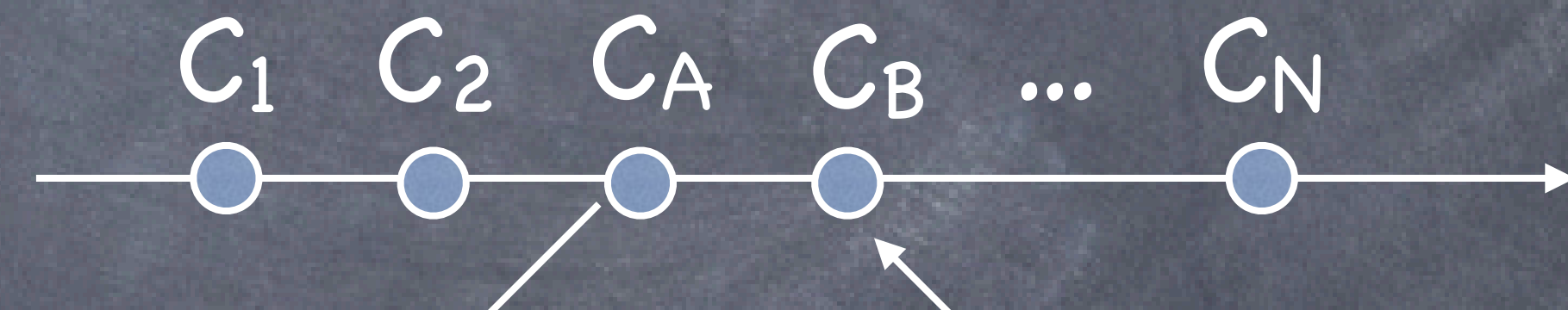
To propose C_B from C_A



Another Markov chain with the probability $W'(C)$

SLHMC

Markov chain with the probability $W(C)$



To propose C_B from C_A



Machine learning molecular dynamics

Machine learning techniques are used for proposing new configuration!

What is the self-learning Monte Carlo?

common simulation with machine learning:



Self-learning Monte Carlo method



We do three steps in same simulations

Num. of training data is drastically reduced (1/10)
because of efficient sampling

Self-learning Monte Carlo

Spin systems

J. Liu, Y. Qi, Z. Y. Meng, and L. Fu, Phys. Rev. B 95, 041101(R) (2017)

H. Kohshiro and YN,,
“Effective Ruderman–Kittel–Kasuya–Yosida-like Interaction in Diluted Double-exchange Model: Self-learning Monte Carlo Approach”,
J. Phys. Soc. Jpn. 90, 034711 (2021)

YN and A. Tomiya, “Self-learning Monte Carlo with equivariant Transformer”, J. Phys. Soc. Jpn. 93, 114007 (2024)

Fermion+classical spins

Electrons

YN, H. Shen, Y. Qi, J. Liu, and L. Fu
“Self-learning Monte Carlo method: Continuous-time algorithm”,
Physical Review B 96, 161102(R) (2017) *Editors’ Suggestion*

YN, M. Okumura, A. Tanaka
“Self-learning Monte Carlo method with Behler-Parrinello neural networks”,
Phys. Rev. B 101, 115111 (2020)

Continuous time Quantum Monte Carlo

Atoms/molecules

Machine-learning MD

YN, M. Okumura, K. Kobayashi, and M. Shiga,
“Self-learning Hybrid Monte Carlo: A First-principles Approach”,
Phys. Rev. B **102**, 041124(R) (2020)

K. Kobayashi, YN, M. Itakura, and M. Shiga,
“Self-learning hybrid Monte Carlo method for isothermal–isobaric ensemble: Application to liquid silica”,
J. Chem. Phys. **155**, 034106 (2021)

YN, Yutaka Iwasaki, Koichi Kitahara, Yoshiki Takagiwa, Kaoru Kimura, Motoyuki Shiga, “High-Temperature Atomic Diffusion and Specific Heat in Quasicrystals”, Phys. Rev. Lett. 132, 196301 (2024)

Bo Thomsen, YN, Keita kobayashi, Ikutaro Hamada, and Motoyuki Shiga, “Self-learning path integral hybrid Monte Carlo with mixed ab initio and machine learning potentials for modeling nuclear quantum effects in water”, J. Chem. Phys. 161, 204109 (2024)

Lattice QCD

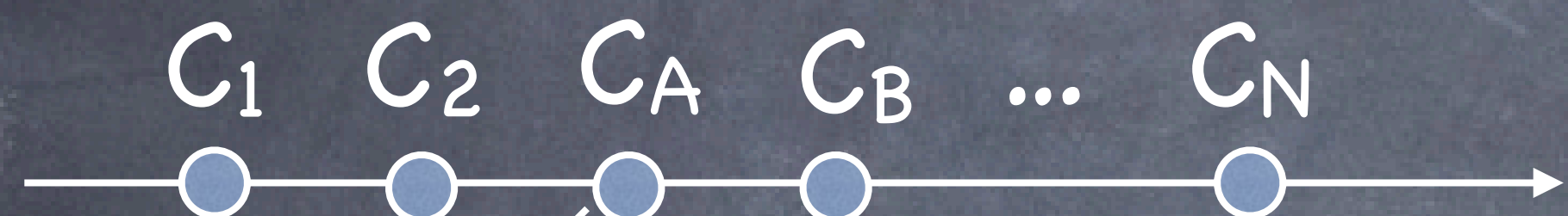
SU(N) Gauge theory on the lattice

YN, Akinori Tanaka, Akio Tomiya,
“Self-learning Monte-Carlo for non-abelian gauge theory with dynamical fermions”,
Phys. Rev. D 107, 054501 (2023)

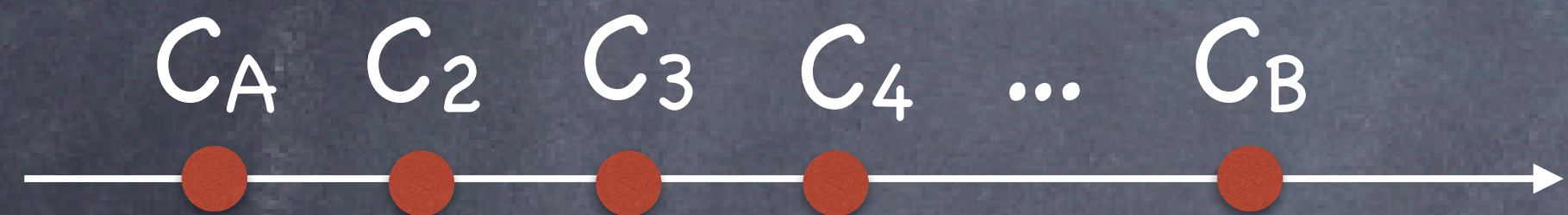
YN and Akio Tomiya,
“Gauge covariant neural network for 4 dimensional non-abelian gauge theory”,
Phys. Rev. D 111, 07450 (2025)

Concept of SLMC

Markov chain with the probability $W(C)$



To propose C_B from C_A

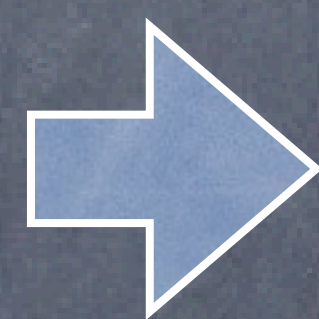


Another Markov chain with the probability $W'(C)$

Acceptance ratio for the Metropolis-Hastings algorithm

$$A(C_B, C_A) = \min \left(1, \frac{W(C_B) g(C_A | C_B)}{W(C_A) g(C_B | C_A)} \right)$$

$g(C_B | C_A)$: Proposal probability



$$A(C_B, C_A) = \min \left(1, \frac{W(C_B) W'(C_A)}{W(C_A) W'(C_B)} \right)$$

If $W'(C) = W(C)$,
the acceptance ratio is **one!**

If the computational cost of the proposal Markov chain is small,
 we can speed up the simulation

How to construct the Markov chain with $W'(C)$?

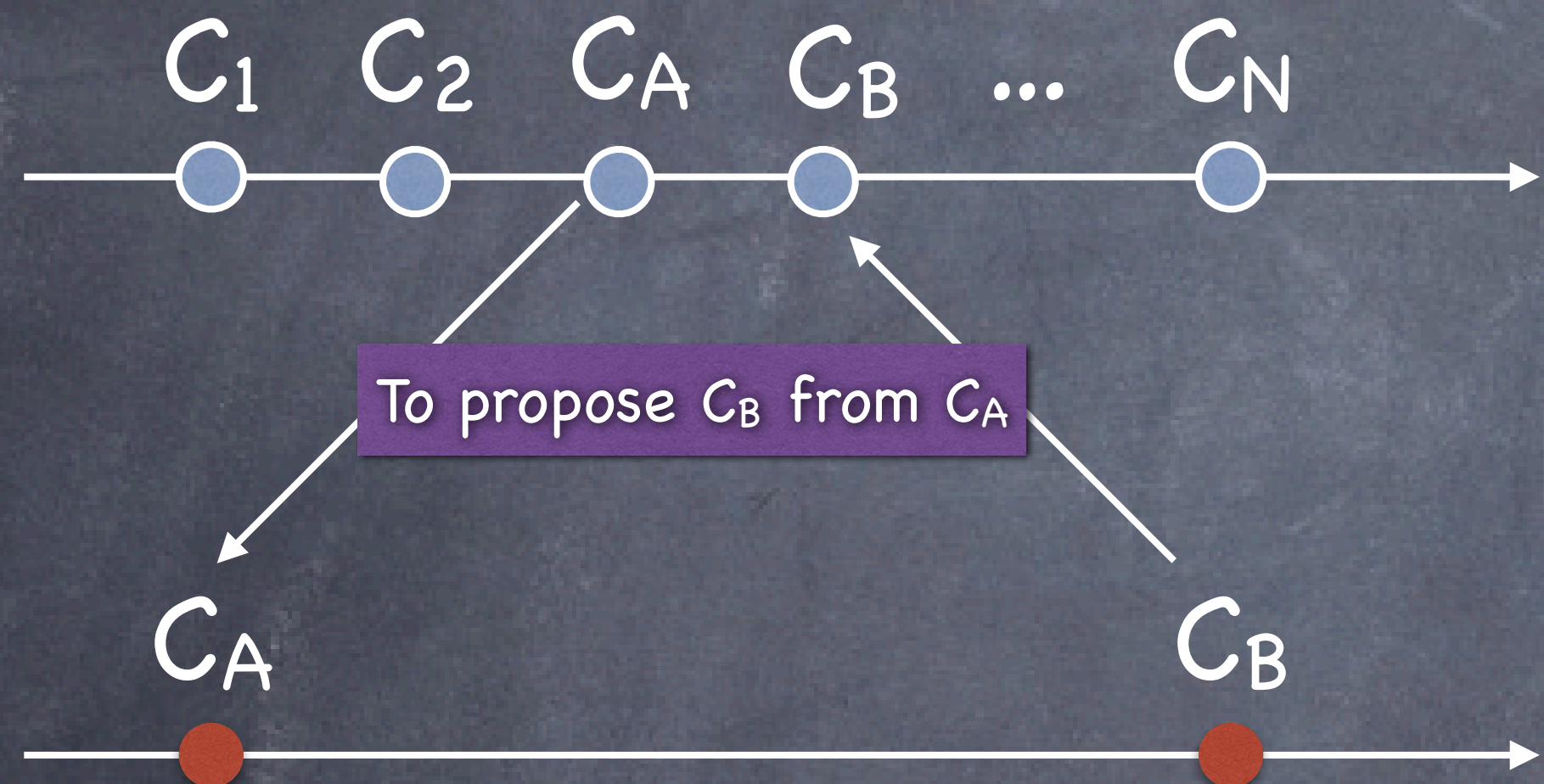
→ Machine learning technique!

$$W(C) = \exp(-\beta H(C)) \rightarrow W'(C) = \exp(-\beta H_{\text{eff}}(C))$$

We construct the effective Hamiltonian

Concept of SLHMC

Markov chain with the probability $W(C)$



Machine learning molecular dynamics (MLMD)

Acceptance ratio for the Metropolis-Hastings algorithm
in Hybrid Monte Carlo

$$A(C_B, C_A) = \min \left(1, \frac{W(C_B)}{W(C_A)} \right) \quad \text{if the MD is time-reversal symmetric}$$

MLMD conserves the energy of the effective model

MLMD DOES NOT conserve the energy of the original model

If the MD conserves the energy of the original model
the acceptance ratio is **one!**

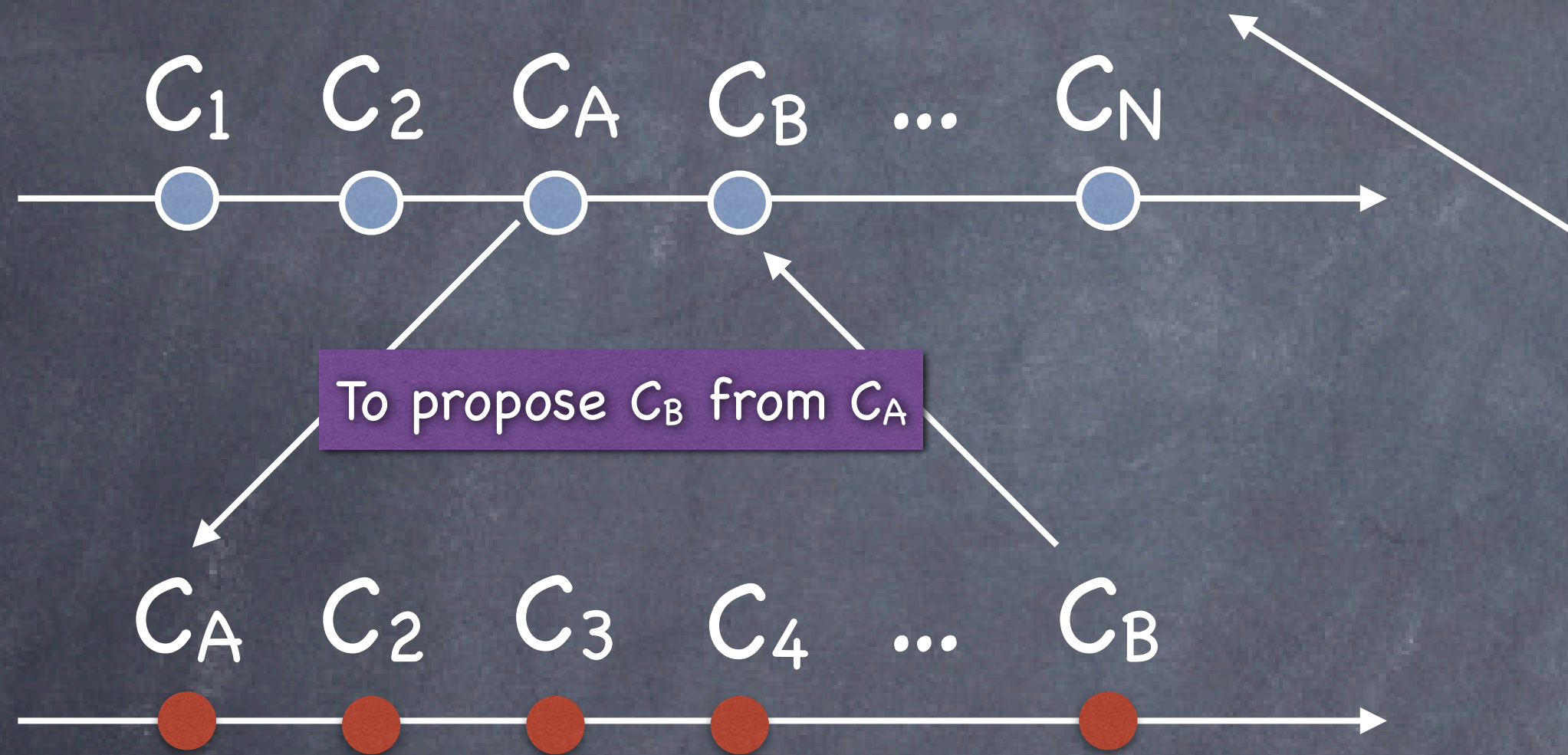
**If the computational cost of the MLMD is small,
we can speed up the simulation**

In the field of atom and molecular systems, machine learning molecular dynamics was proposed in 2007

Self-learning Monte Carlo

for lattice QCD

Markov chain with the probability $W(C)$



Another Markov chain with the probability $W'(C)$

YN, Akinori Tanaka, Akio Tomiya,
 “Self-learning Monte-Carlo for non-abelian gauge theory with dynamical fermions”,
 Phys. Rev. D 107, 054501 (2023)

$$S[U] = S_g[U] + S_f[U],$$

$$S_f[U] = -\log \det M^\dagger M,$$

integrated fermion action

effective model without fermion actions

$$\begin{aligned}
 S_{\text{eff}}^\theta[U] = & \sum_n \left[\beta_{\text{plaq}} \sum_{\mu=1}^4 \sum_{\nu>\mu} \left(1 - \frac{1}{2} \text{tr} U_{\mu\nu}(n) \right) + \beta_{\text{rect}} \sum_{\mu=1}^4 \sum_{\nu \neq \mu} \left(1 - \frac{1}{2} \text{tr} R_{\mu\nu}(n) \right) \right] \\
 & + \beta_{\text{Pol}}^{\mu=1} \sum_{n_2, n_3, n_4} \text{tr} \left[\prod_{n_1=0}^{N_1-1} U_1(\vec{n}, n_4) \right] + \beta_{\text{Pol}}^{\mu=2} \sum_{n_1, n_3, n_4} \text{tr} \left[\prod_{n_2=0}^{N_2-1} U_2(\vec{n}, n_4) \right] \\
 & + \beta_{\text{Pol}}^{\mu=3} \sum_{n_1, n_2, n_4} \text{tr} \left[\prod_{n_3=0}^{N_3-1} U_3(\vec{n}, n_4) \right] + \beta_{\text{Pol}}^{\mu=4} \sum_{n_1, n_2, n_3} \text{tr} \left[\prod_{n_4=0}^{N_4-1} U_4(\vec{n}, n_4) \right] + \beta_{\text{const}},
 \end{aligned}$$

Self-learning Monte Carlo

for lattice QCD

YN, Akinori Tanaka, Akio Tomiya,
 “Self-learning Monte-Carlo for non-abelian gauge theory with dynamical fermions”,
 Phys. Rev. D 107, 054501 (2023)

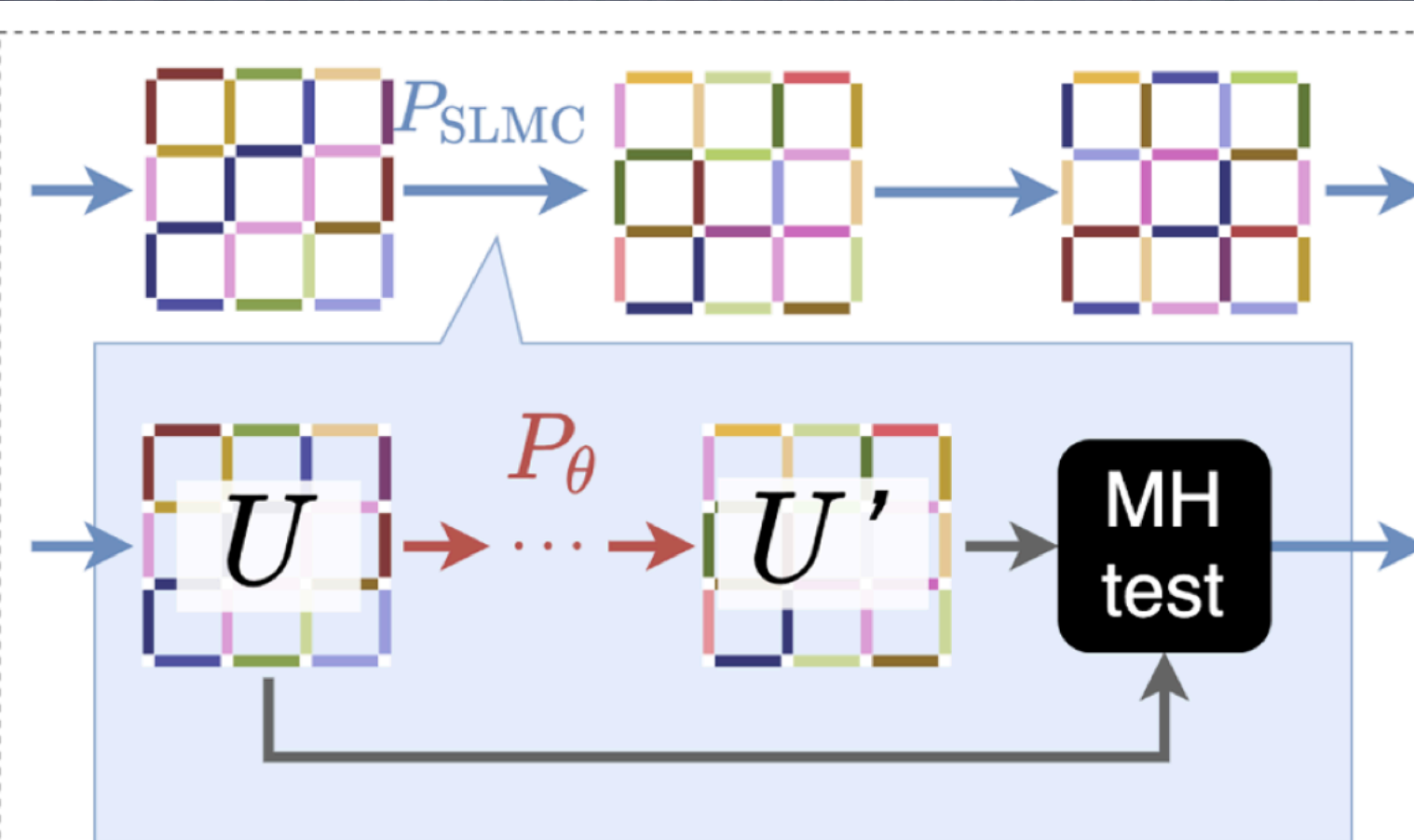
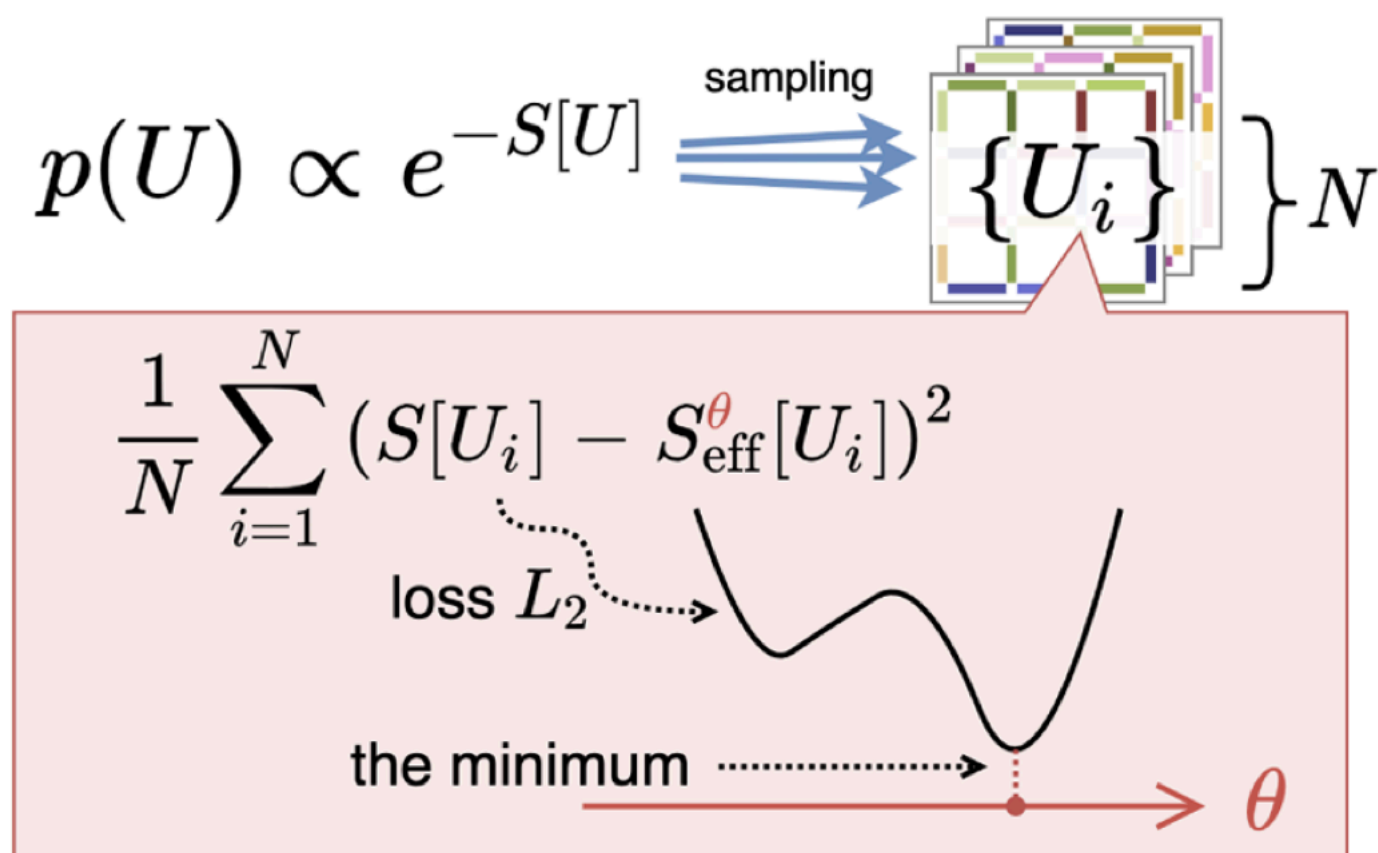
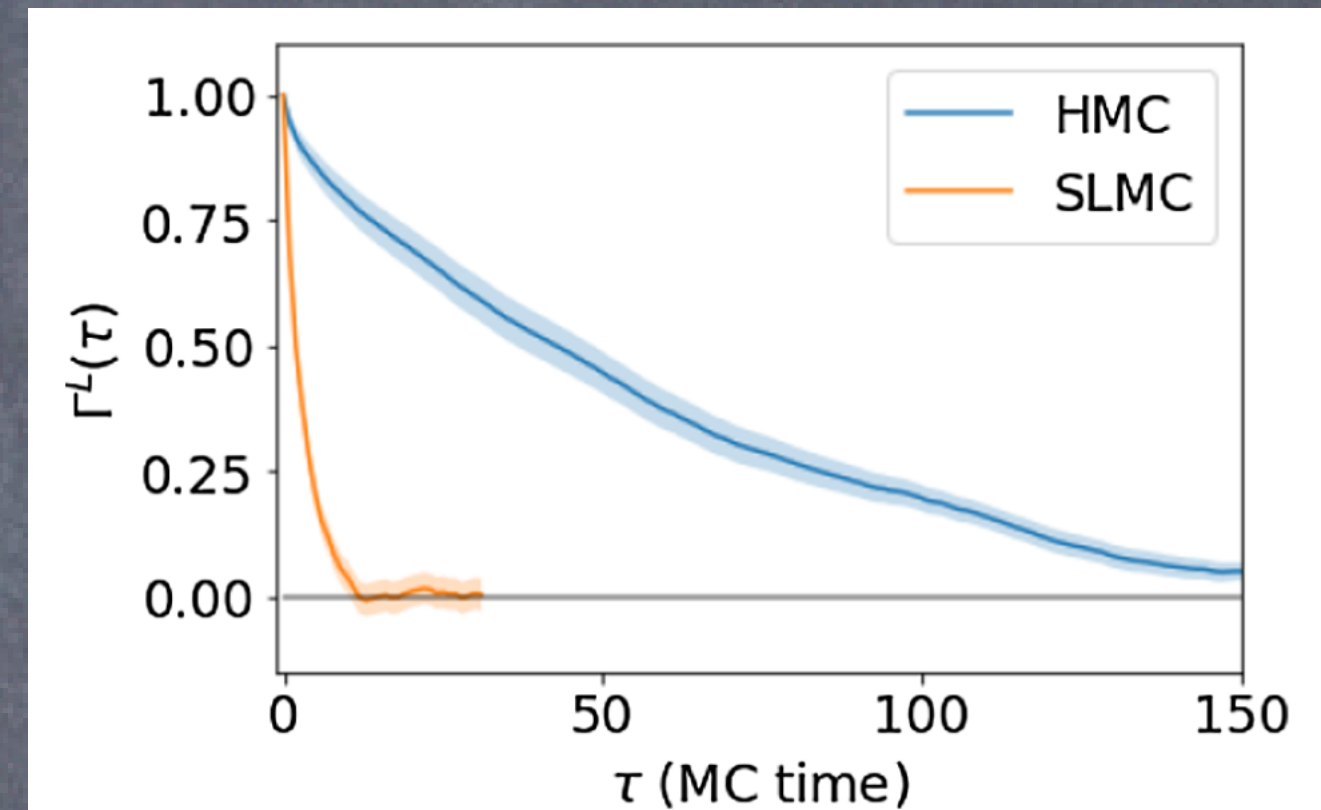
effective model without fermion actions

$$S[U] = S_g[U] + S_f[U],$$

$$S_f[U] = -\log \det M^\dagger M,$$

integrated fermion action

$$\begin{aligned}
 S_{\text{eff}}^\theta[U] = & \sum_n \left[\beta_{\text{plaq}} \sum_{\mu=1}^4 \sum_{\nu>\mu} \left(1 - \frac{1}{2} \text{tr} U_{\mu\nu}(n) \right) + \beta_{\text{rect}} \sum_{\mu=1}^4 \sum_{\nu \neq \mu} \left(1 - \frac{1}{2} \text{tr} R_{\mu\nu}(n) \right) \right] \\
 & + \beta_{\text{Pol}}^{\mu=1} \sum_{n_2, n_3, n_4} \text{tr} \left[\prod_{n_1=0}^{N_1-1} U_1(\vec{n}, n_4) \right] + \beta_{\text{Pol}}^{\mu=2} \sum_{n_1, n_3, n_4} \text{tr} \left[\prod_{n_2=0}^{N_2-1} U_2(\vec{n}, n_4) \right] \\
 & + \beta_{\text{Pol}}^{\mu=3} \sum_{n_1, n_2, n_4} \text{tr} \left[\prod_{n_3=0}^{N_3-1} U_3(\vec{n}, n_4) \right] + \beta_{\text{Pol}}^{\mu=4} \sum_{n_1, n_2, n_3} \text{tr} \left[\prod_{n_4=0}^{N_4-1} U_4(\vec{n}, n_4) \right] + \beta_{\text{const}},
 \end{aligned}$$

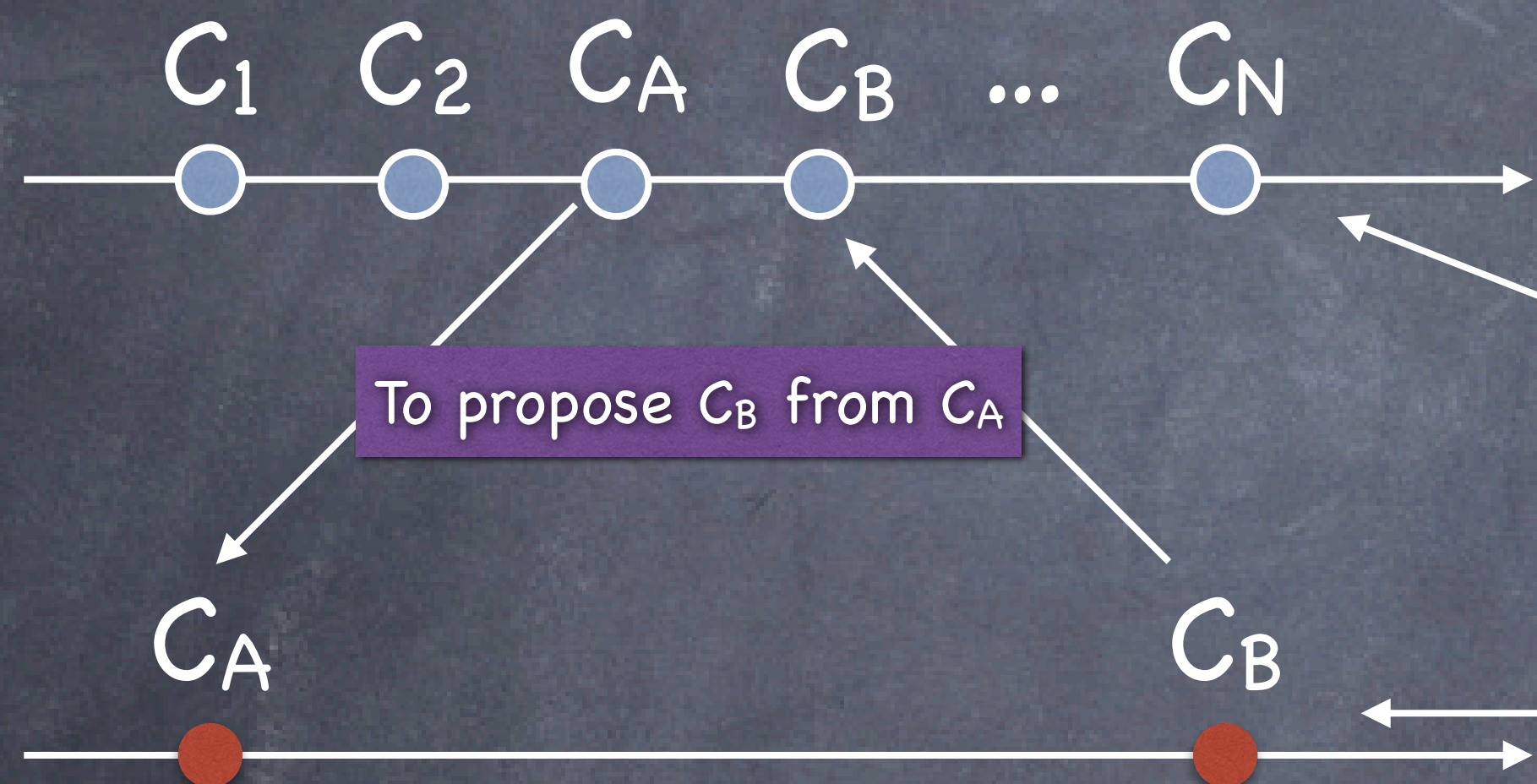


We use a linear interpolation
 how to improve effective action?

Self-learning Hybrid Monte Carlo

for lattice QCD

Markov chain with the probability $W(C)$



YN and Akio Tomiya,

“Gauge covariant neural network for 4 dimensional non-abelian gauge theory”,
arXiv:2103.11965

target action

$$S[U] = S_g[U] + S_f[\phi, U; m_l],$$

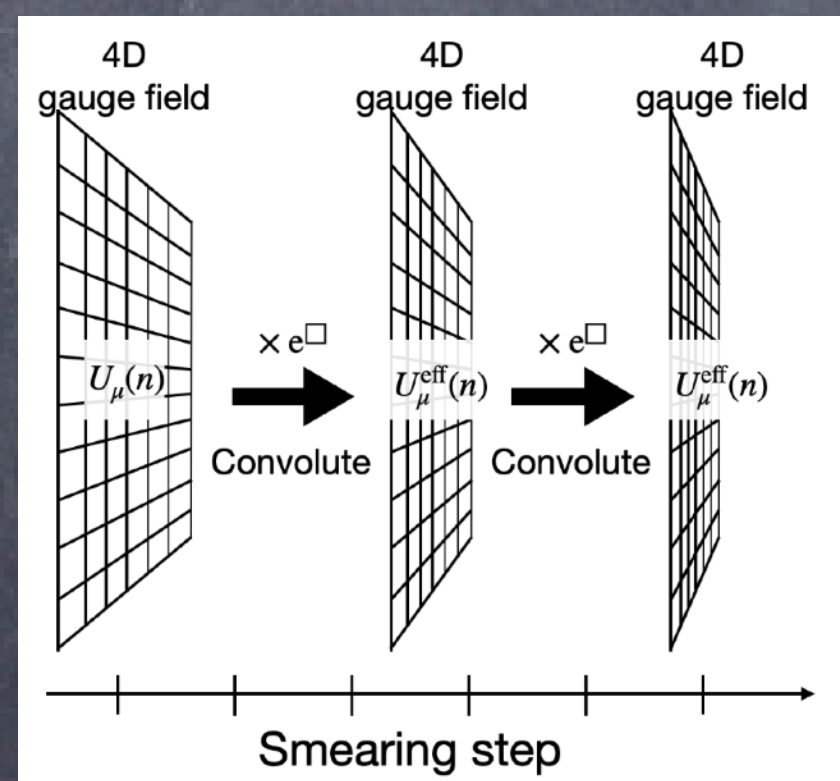
effective action

$$S_\theta[U] = S_g[U] + S_f[\phi, U_\theta^{\text{NN}}[U]; m_h],$$

Machine learning molecular dynamics (MLMD)

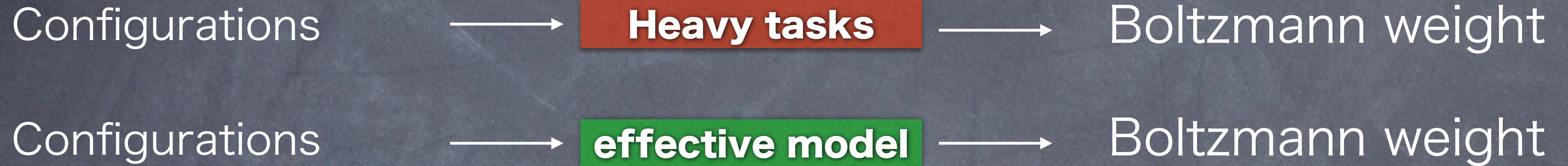
if $m_l < m_h$, the computational cost reduces

U^{NN} : trainable stout smearing



-> Dr. Akio Tomiya's talk

Problems of SLMC



How to construct effective models?

Quality of the effective model is very important

In previous studies,

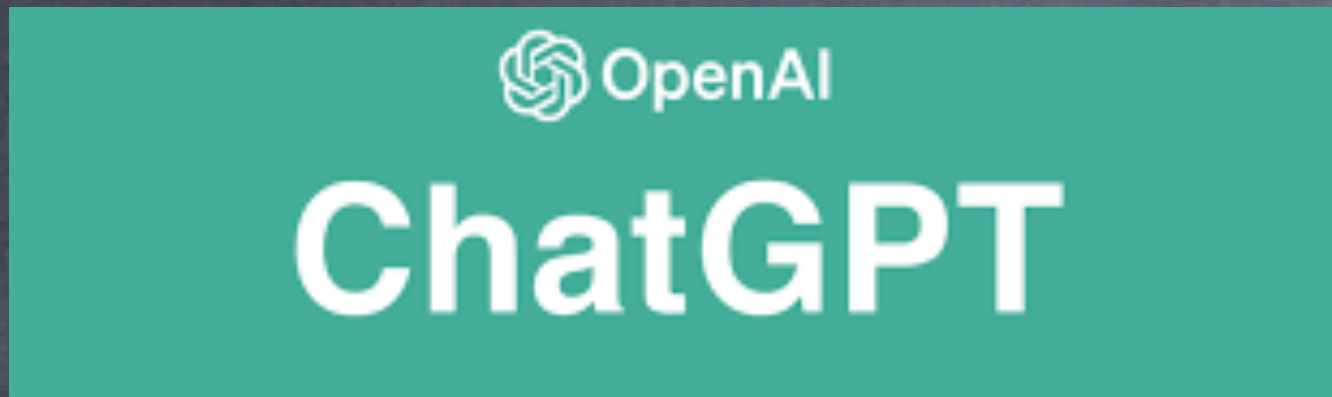
for example, a linear regression is used to construct the effective model inspired by the physical insight

Use Transformer!!

Transformer and Attention mechanism

Generative AIs

<https://arxiv.org/abs/1706.03762>



These AI have same architecture called Transformer

Transformer

AI Chat, Visualization, language translation

protein foldings etc.

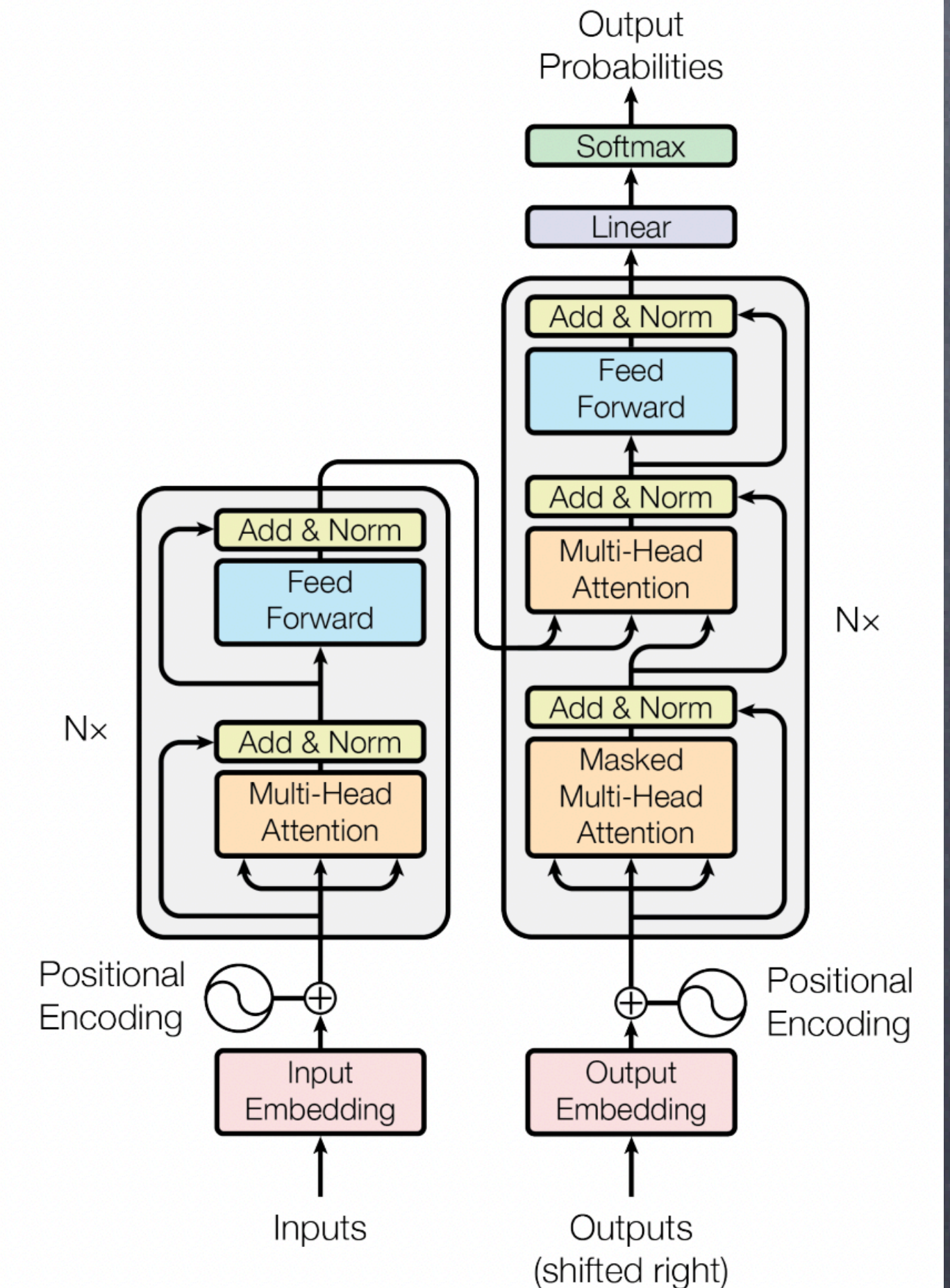


Figure 1: The Transformer - model architecture.

Scaling laws of Transformer

<https://arxiv.org/abs/2001.08361>

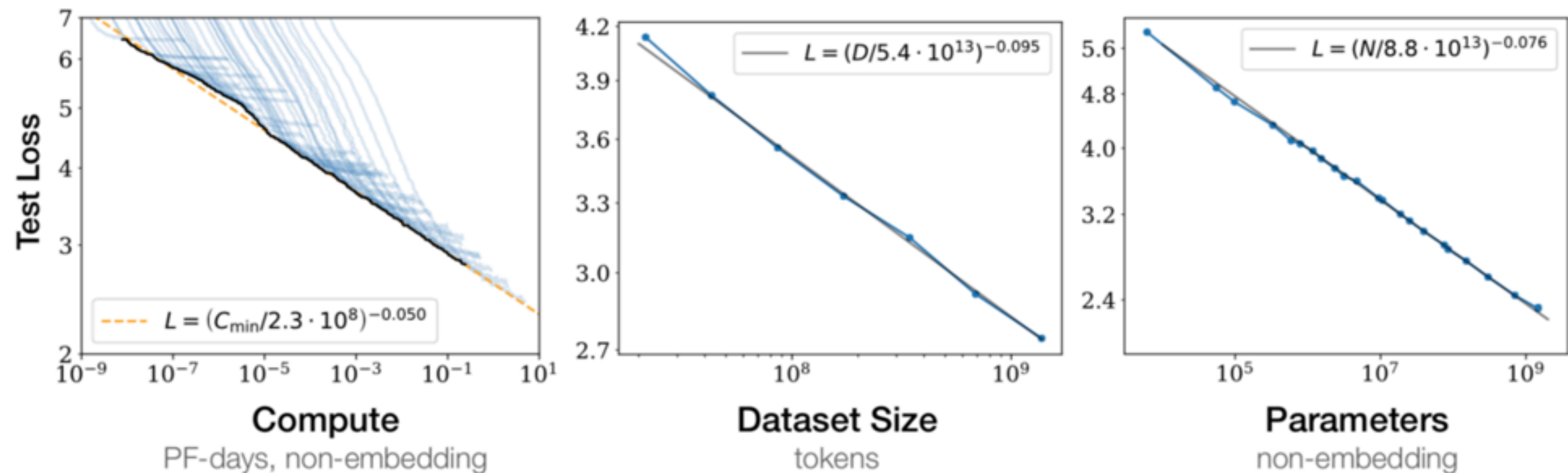


Figure 1 Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute² used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

It requires huge data (e.g. GPT uses all electric books in the world)
= weak inductive bias, large data makes prediction better

Attention learns Non-local correlations

“Attention is all you need” (arXiv:1706.03762)

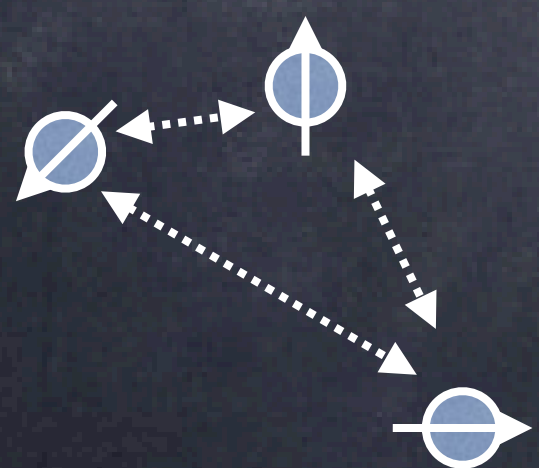
Attention is an important building block in the transformer

Language model (chat, translation etc..)

“Attention” layer can capture non-local relations between words

Physics

Can we use Attention layer to learn **non local correlation**?



If so, we can design a new model in physics

What is the attention mechanism?

There are many websites to explain the transformer and attention mechanism, in terms of language processing...

I try to explain the attention in terms of simple mathematics

This came from discussions with Dr. Tomiya

1. We consider a vector/matrix/tensor A A_i or A_{ij} or A_{ijk}

2. We make three variables K, Q, V from A

$$K = W^K A, Q = W^Q A, V = W^V A \quad W^K, W^Q, W^V: \text{trainable parameters}$$

3. We generate new vector/matrix/tensor B

$$B_l = A_l + \sum_i P_i^l V_i \quad P = \sigma(QK^T)$$

correlation between Q and K

$l=i$ or ij or ijk

What is the attention mechanism?

$$K = W^K A, Q = W^Q A, V = W^V A \quad W^K, W^Q, W^V: \text{trainable parameters}$$

3. We generate new vector/matrix/tensor B

$$B_l = A_l + \sum_i P_i^l V_i \quad P = \sigma(QK^T) \quad \sigma: \text{nonlinear function}$$

weighted sum
correlation between Q and K

self-attention mechanism

This is most simplest architecture

In generative AIs, they use the multi-head attention

Simple mechanism but very effective!

How can we use this in physics?

Equivariant transformer

Difference between AI and Physics

Generative AIs: huge number of the parameters (Billions, Trillions)

We have to gather huge amount of data

We need large computational resources and money

We do not want to gather and train data by ourselves

But, we want to use transformers because of good properties

Physics: small number of parameters

When there are too many parameters in a machine learning model,
I feel like I don't really understand what's going on.

If the computation is too heavy, it's impossible to speed up the simulation

We want to reduce the number of parameters in transformers

Incorporating natural laws can simplify the model -> Symmetries!

Fermion and spin model

We want to focus on a simple lattice model

fermions and classical spins

$$H = -t \sum_{\alpha, \langle i, j \rangle} (\hat{c}_{i\alpha}^\dagger \hat{c}_{j\alpha} + \text{h.c.}) + \frac{J}{2} \sum_i \mathbf{S}_i \cdot \hat{\boldsymbol{\sigma}}_i - \mu \sum_{\alpha, i} \hat{c}_{i\alpha}^\dagger \hat{c}_{i\alpha},$$

called double exchange model
in condensed matter physics

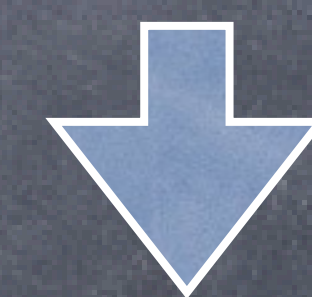
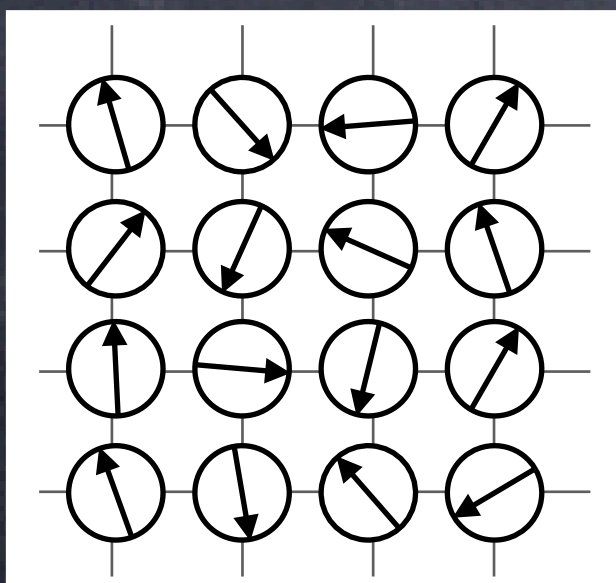
Partition function:

$$Z = \sum_{\{\mathbf{S}\}} \prod_n (1 + e^{-\beta(\mu - E_n(\{\mathbf{S}\}))})$$

Input: spin configurations $\{\mathbf{S}\}$

Configurations: classical spins $\{\mathbf{S}_i\}$

\mathbf{S}_i : i -th three dimensional vector in spin space



diagonalization

Output: Boltzmann weight

We want to replace the diagonalization

Fermion and spin model

We want to focus on a simple lattice model

fermions and classical spins

$$H = -t \sum_{\alpha, \langle i, j \rangle} (\hat{c}_{i\alpha}^\dagger \hat{c}_{j\alpha} + \text{h.c.}) + \frac{J}{2} \sum_i \mathbf{S}_i \cdot \hat{\sigma}_i - \mu \sum_{\alpha, i} \hat{c}_{i\alpha}^\dagger \hat{c}_{i\alpha},$$

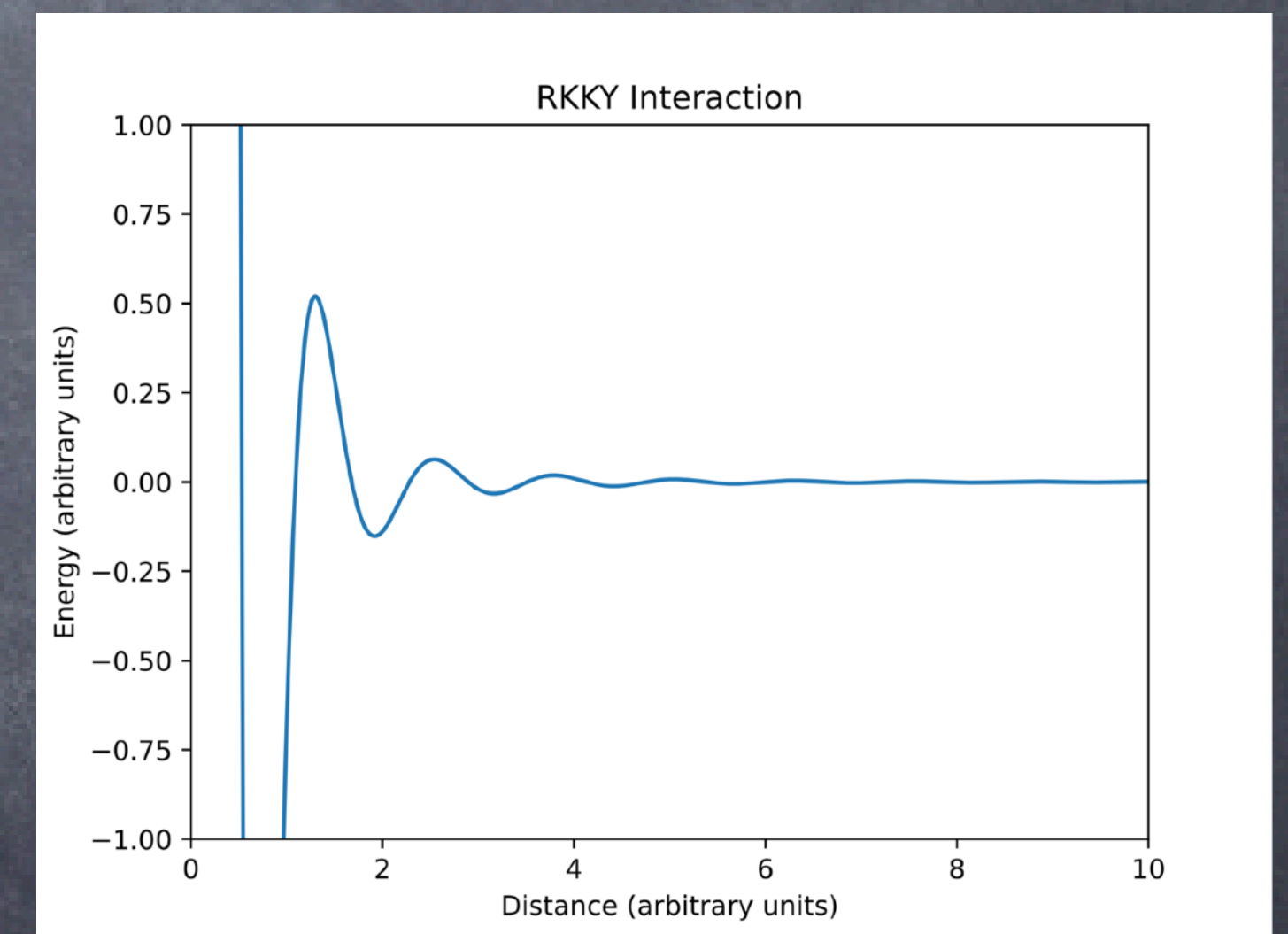
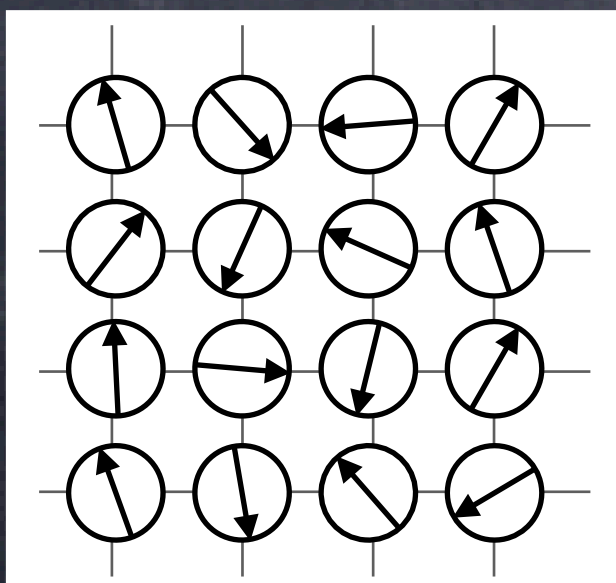
called double exchange model
in condensed matter physics

If J is small, we can use the perturbation theory

the **Ruderman–Kittel–Kasuya–Yosida (RKKY) interaction** models

$$H_{\text{RKKY}} = - \sum_{\langle i, j \rangle_n} J_n \mathbf{S}_i \cdot \mathbf{S}_j$$

We can integrate out fermion degrees of freedom
fermion + spin \rightarrow spin



Fermion and spin model

We want to focus on a simple lattice model

fermions and classical spins

$$H = -t \sum_{\alpha, \langle i, j \rangle} (\hat{c}_{i\alpha}^\dagger \hat{c}_{j\alpha} + \text{h.c.}) + \frac{J}{2} \sum_i \mathbf{S}_i \cdot \hat{\sigma}_i - \mu \sum_{\alpha, i} \hat{c}_{i\alpha}^\dagger \hat{c}_{i\alpha},$$

called double exchange model
in condensed matter physics

We want to consider large J region

Simple effective model

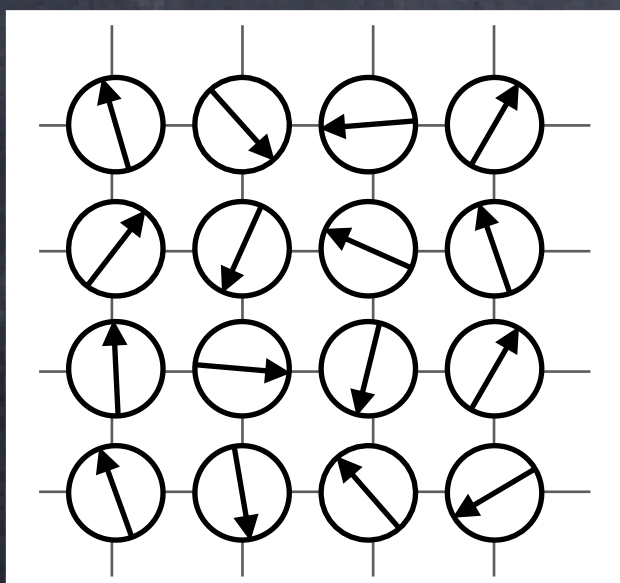
J. Liu, H. Shen, Y. Qi, Z. Y. Meng, and L. Fu, Phys. Rev. B 95, 241104(R)(2017)

J_n^{eff} : n-th nearest neighbor interaction

$$H_{\text{eff}}^{\text{Linear}} = - \sum_{\langle i, j \rangle_n} J_n^{\text{eff}} \mathbf{S}_i \cdot \mathbf{S}_j + E_0$$

This is a linear model

by integrating out fermion degrees of freedom



similar to RKKY model
derived by physicist

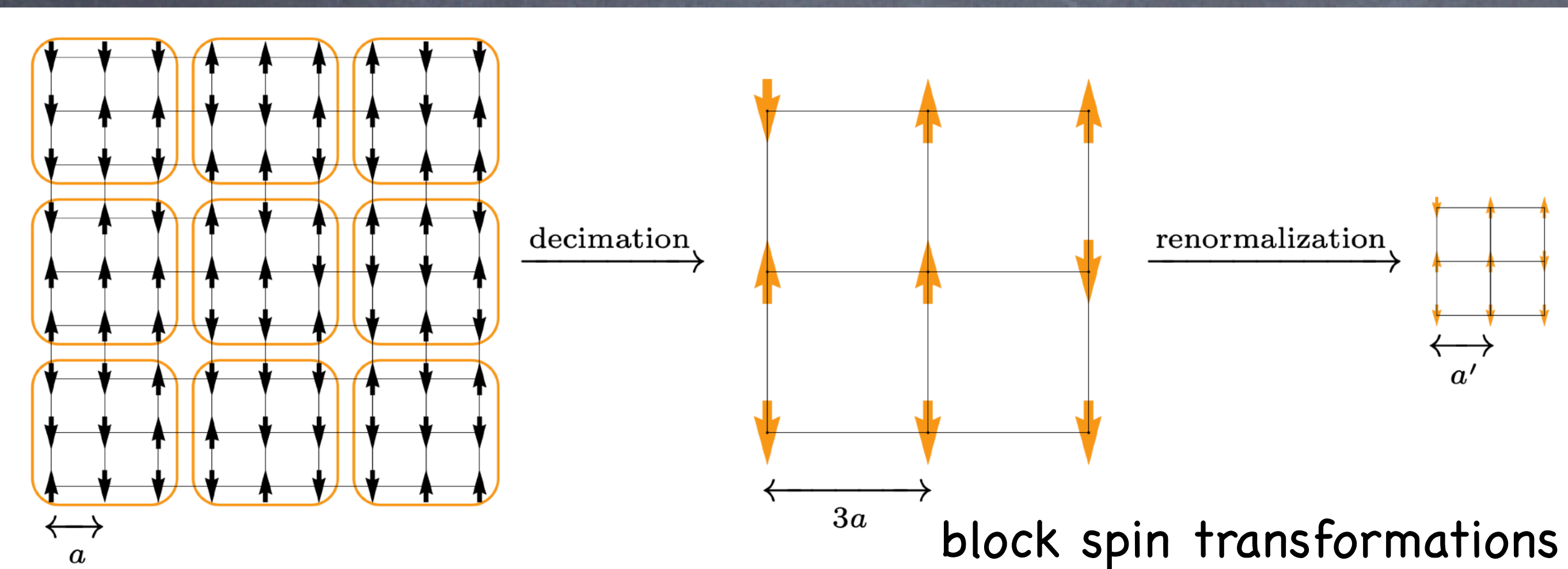
There are only few parameters J_n^{eff}

Num. of parameters is too small! How to improve this model?

Fermion and spin model

How to construct model?

In physics, we know the renormalization group analysis



Spins become “effective” spins

$$H_{\text{eff}} = - \sum_{\langle i,j \rangle_n} J_n^{\text{eff}} \mathbf{S}_i^{\text{NN}} \cdot \mathbf{S}_j^{\text{NN}} + E_0$$

Heisenberg model for effective spins

Spins are renormalized

Renormalized spin should have same symmetries

If we can construct effective spins, we can construct effective model!

We need an equivariant model

Charlie Duclut., “Nonequilibrium critical phenomena :exact Langevin equations, erosion of tilted landscapes” Université Pierre et Marie Curie – Paris VI, 2017.

What is equivariance?



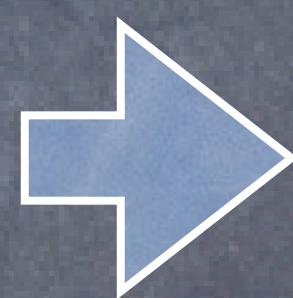
A. Cat or Not?

Rotate input image

→ Still recognized as a cat

Output unchanged → **Invariant**

B. Where are the cat's ears?



Rotate input image

→ Output (ear position) also rotates

Output changes accordingly

→ **Equivariant**

Equivariance =

Output transforms with input

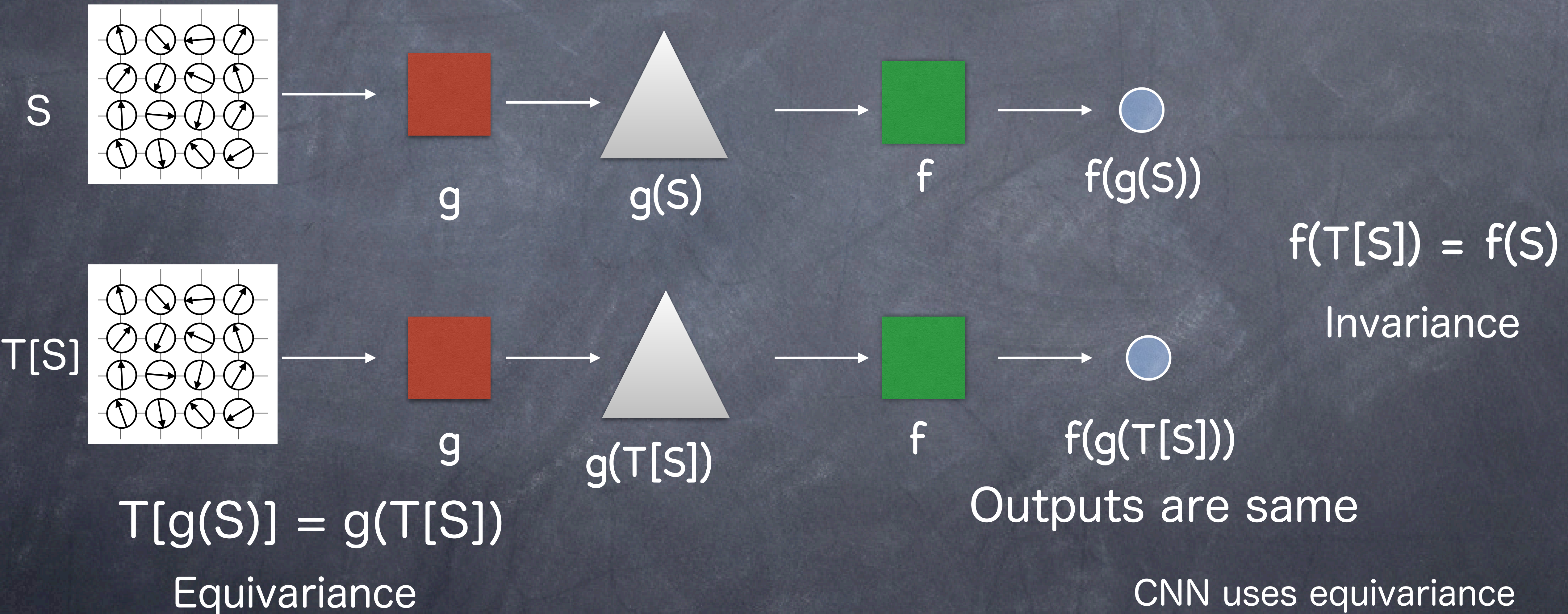
Hamiltonian: Rotational **invariant**

Effective spin:

must be rotational **equivariant**

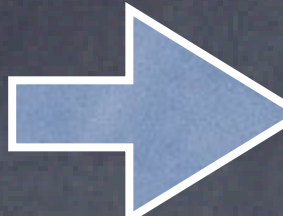
Invariance and equivariance

We make equivariant networks and make the output invariant



How to construct the attention layer

1. We consider a vector/matrix/tensor \mathbf{A} A_i or A_{ij} or A_{ijk}

 We consider spin "matrix" $[\hat{S}]_{i\mu} = [\vec{S}_i]_\mu$
 \vec{S}_i is a classical spin on i-th site (vector)

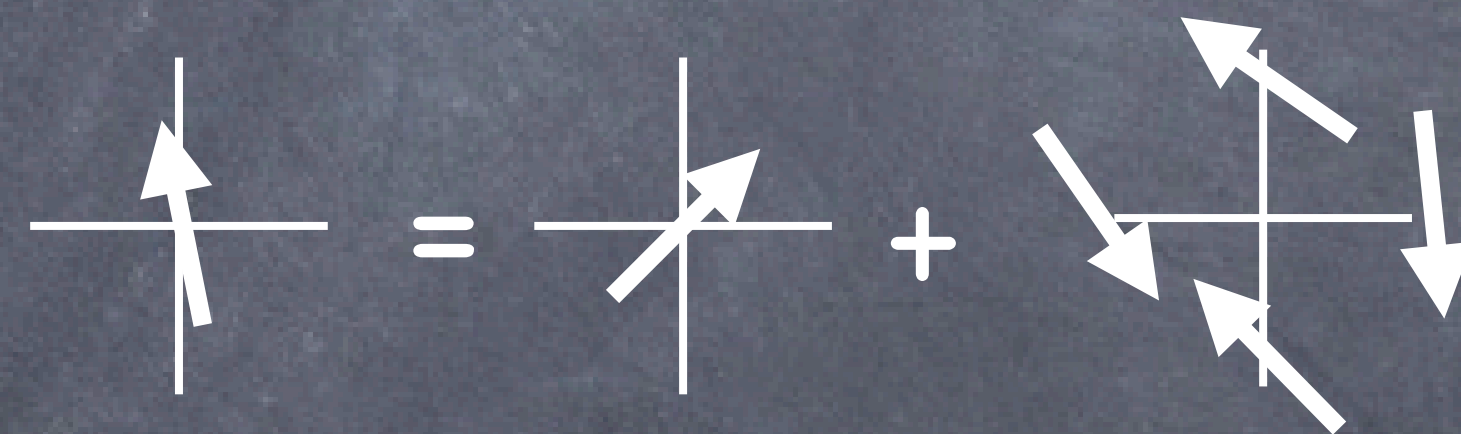
2. We make three variables $\mathbf{K}, \mathbf{Q}, \mathbf{V}$ from \mathbf{A}

$$\mathbf{K} = \mathbf{W}^{\mathbf{K}} \mathbf{A}, \quad \mathbf{Q} = \mathbf{W}^{\mathbf{Q}} \mathbf{A}, \quad \mathbf{V} = \mathbf{W}^{\mathbf{V}} \mathbf{A}$$

 We introduce "operators" $\hat{S}^{\mathbf{Q}} = \bar{\mathbf{W}}^{\mathbf{Q}} \hat{S}$ $\hat{S}^{\mathbf{K}} = \bar{\mathbf{W}}^{\mathbf{K}} \hat{S}$ $\hat{S}^{\mathbf{V}} = \bar{\mathbf{W}}^{\mathbf{V}} \hat{S}$

$$[\bar{\mathbf{W}}^\alpha \hat{S}]_{i\mu} \equiv \sum_{\langle i,j \rangle_n} W_n^\alpha \hat{S}_{j\mu}$$

n -th nearest neighbors



$\mathbf{W}^{\mathbf{K}}, \mathbf{W}^{\mathbf{Q}}, \mathbf{W}^{\mathbf{V}}$: trainable parameters

$\mathbf{W}^{\mathbf{K}}, \mathbf{W}^{\mathbf{Q}}, \mathbf{W}^{\mathbf{V}}$ do not depend on the site i (translational symmetry)

num. of parameters becomes a few

How to construct the attention layer

3. We generate new vector/matrix/tensor \mathbf{B}

$$B_l = A_l + \sum P_i^l V_i$$

➡ $\hat{S}^{(l)} \equiv \mathcal{N}(\hat{S}^{(l-1)} + \check{M} \hat{S}^V)$ $[\mathcal{N}(\hat{S})]_{i\mu} = [\vec{S}_i]_{\mu} / ||\vec{S}_i||$

$P = \sigma(QK^T)$ correlation between Q and K

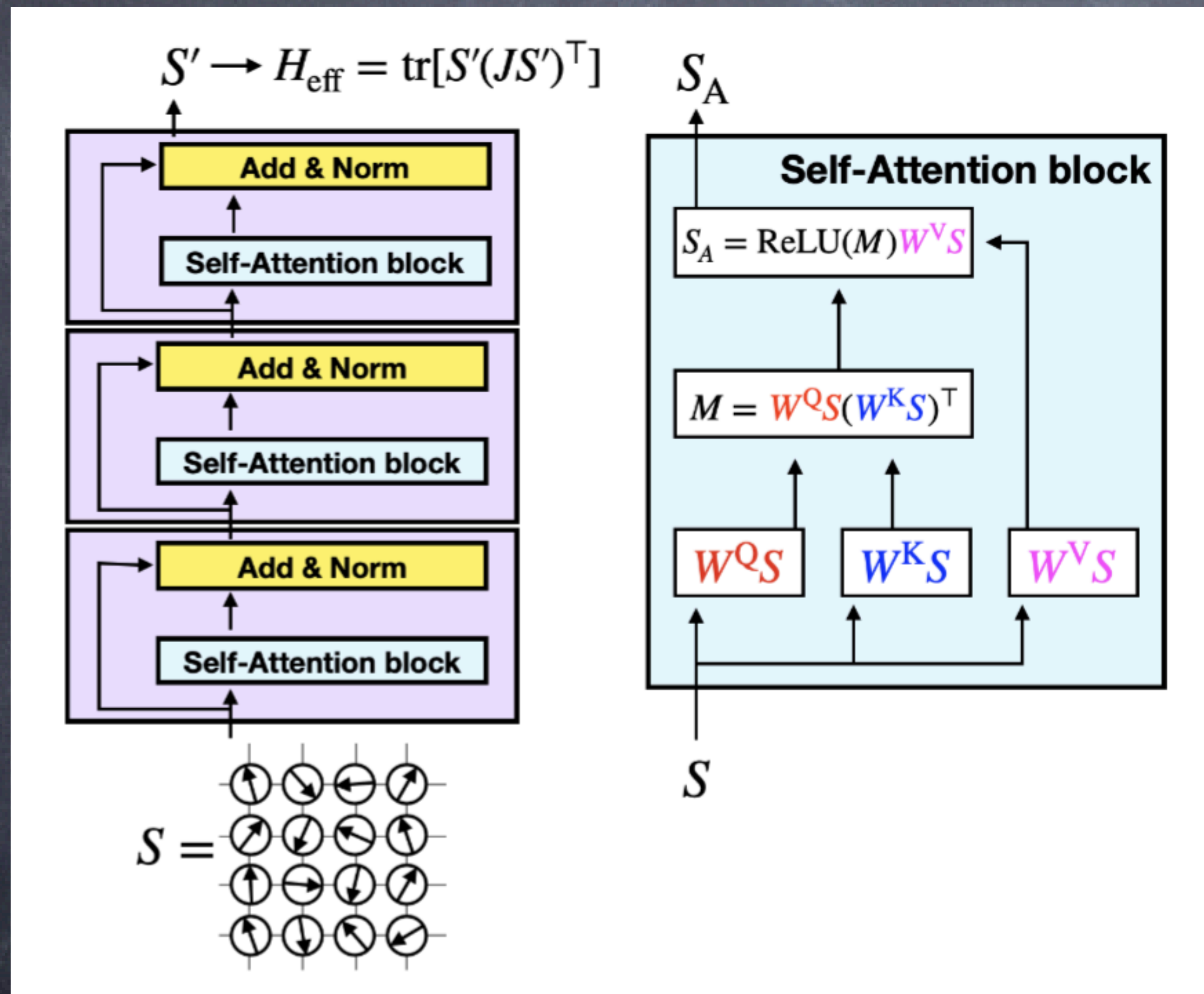
➡ $[\check{M}]_{ij} = \text{ReLU} \left(\frac{1}{\sqrt{3}} \sum_{\mu=1}^3 \hat{S}_{i\mu}^Q \hat{S}_{j\mu}^K \right)$

The “effective” spin $S^{(L)}$ can be regarded as a physical spin

$\hat{S}^{(l)}$ has spin-rotational equivariance
 $R[g(S)] = g(R[S])$

→ renormalized spin
We can build a model!

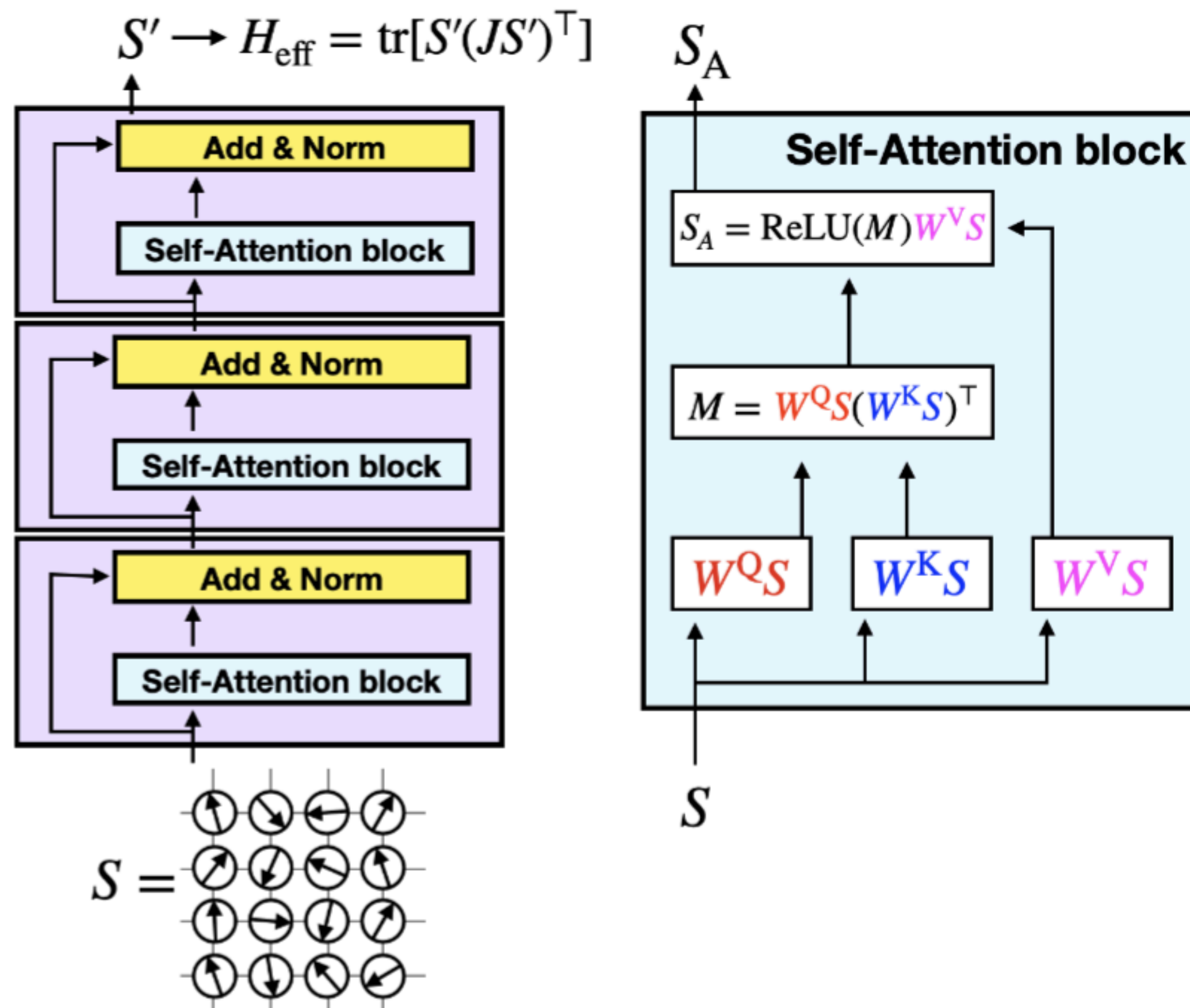
Equivariant Transformer for spin systems



$$\begin{aligned}
 S &\rightarrow \begin{cases} K = W^K S \\ Q = W^Q S \\ V = W^V S \end{cases} \\
 K_i &= \sum_l W_l S_{i+l} \\
 S &= \begin{array}{c} \uparrow \\ \downarrow \end{array} = \begin{array}{c} \uparrow \\ \downarrow \end{array} + \begin{array}{c} \uparrow \\ \downarrow \end{array} \\
 &\text{W only mixes neighbor spins (short range interaction)} \\
 &\text{like block spin transformations} \\
 M &= W^Q S S + W^K \\
 &\text{Rotational and translational invariant} \\
 S' &= S + \text{ReLU}(M) W^V S \\
 &\text{Long range correlation is included}
 \end{aligned}$$

Equivariant Transformer for spin systems

$$\mathcal{N}(\mathbf{S}_i) = \mathbf{S}_i / \|\mathbf{S}_i\|$$



Layer 1

$$S_1 = \mathcal{N}(S + \text{ReLU}(M^1(S)) W^{V1} S)$$

Layer 2

$$S_2 = \mathcal{N}(S_1 + \text{ReLU}(M^2(S_1)) W^{V2} S_1)$$

Layer 3

$$S_3 = \mathcal{N}(S_2 + \text{ReLU}(M^3(S_2)) W^{V3} S_2)$$

Last Heisenberg model with effective spins

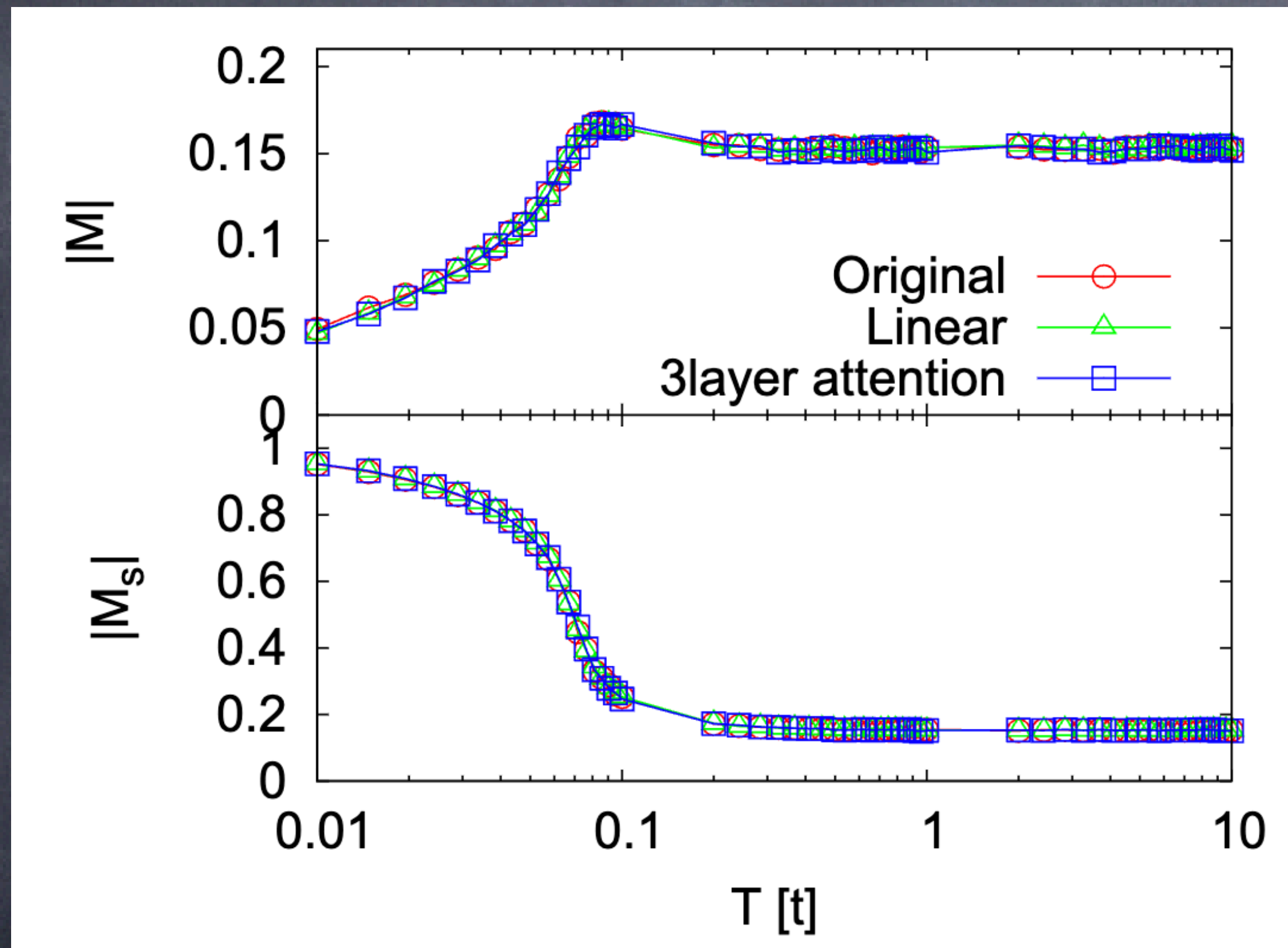
$$E = \sum_i \sum_l J_l \vec{S}_{3i} \cdot \vec{S}_{3i+l} + E_0$$

If the second term is zero

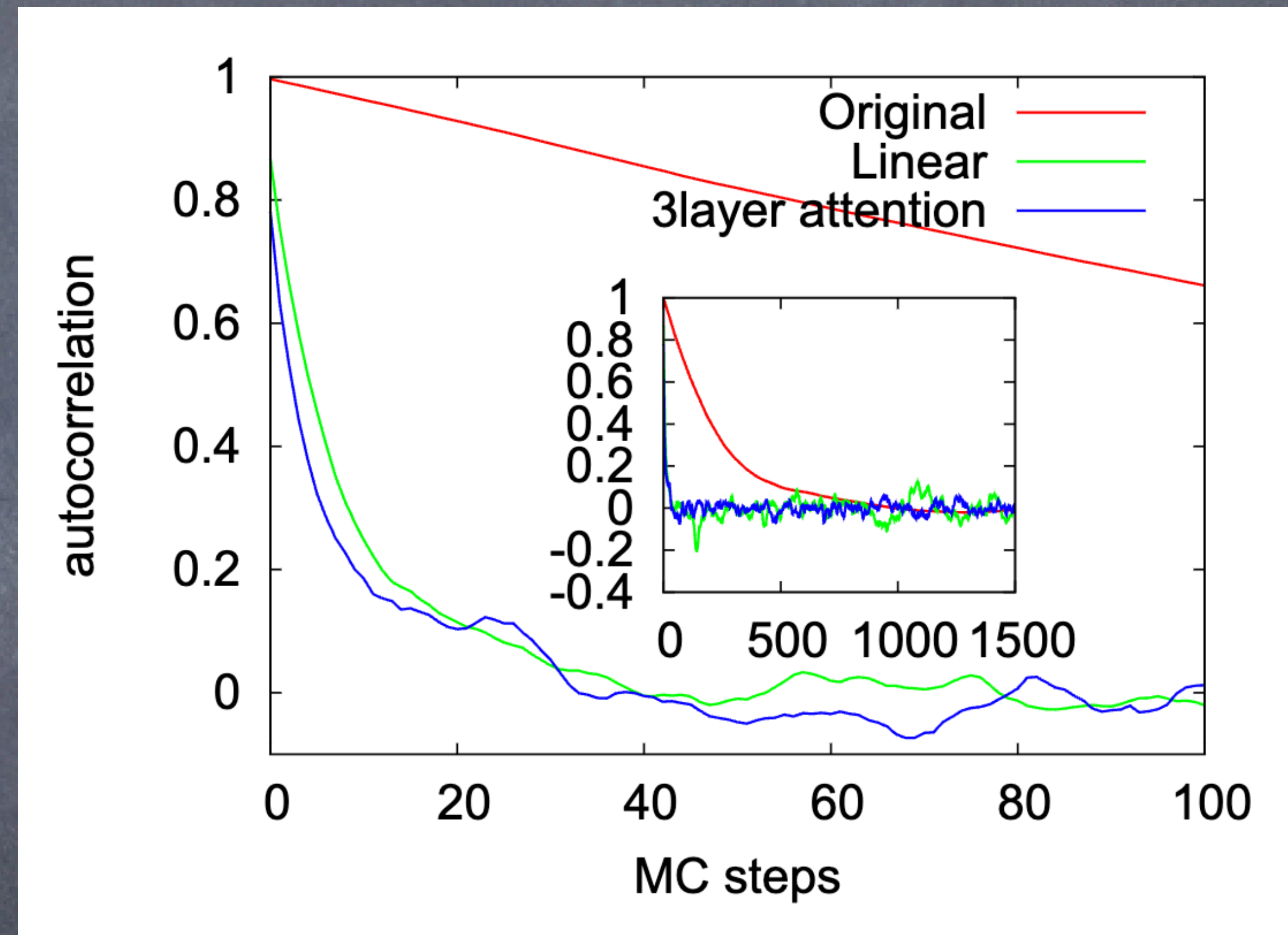
$$E = \sum_i \sum_l J_l \vec{S}_i \cdot \vec{S}_{i+l} + E_0 \quad \text{we get linearized model}$$

Results

2D double exchange model(fermion + classical spin)



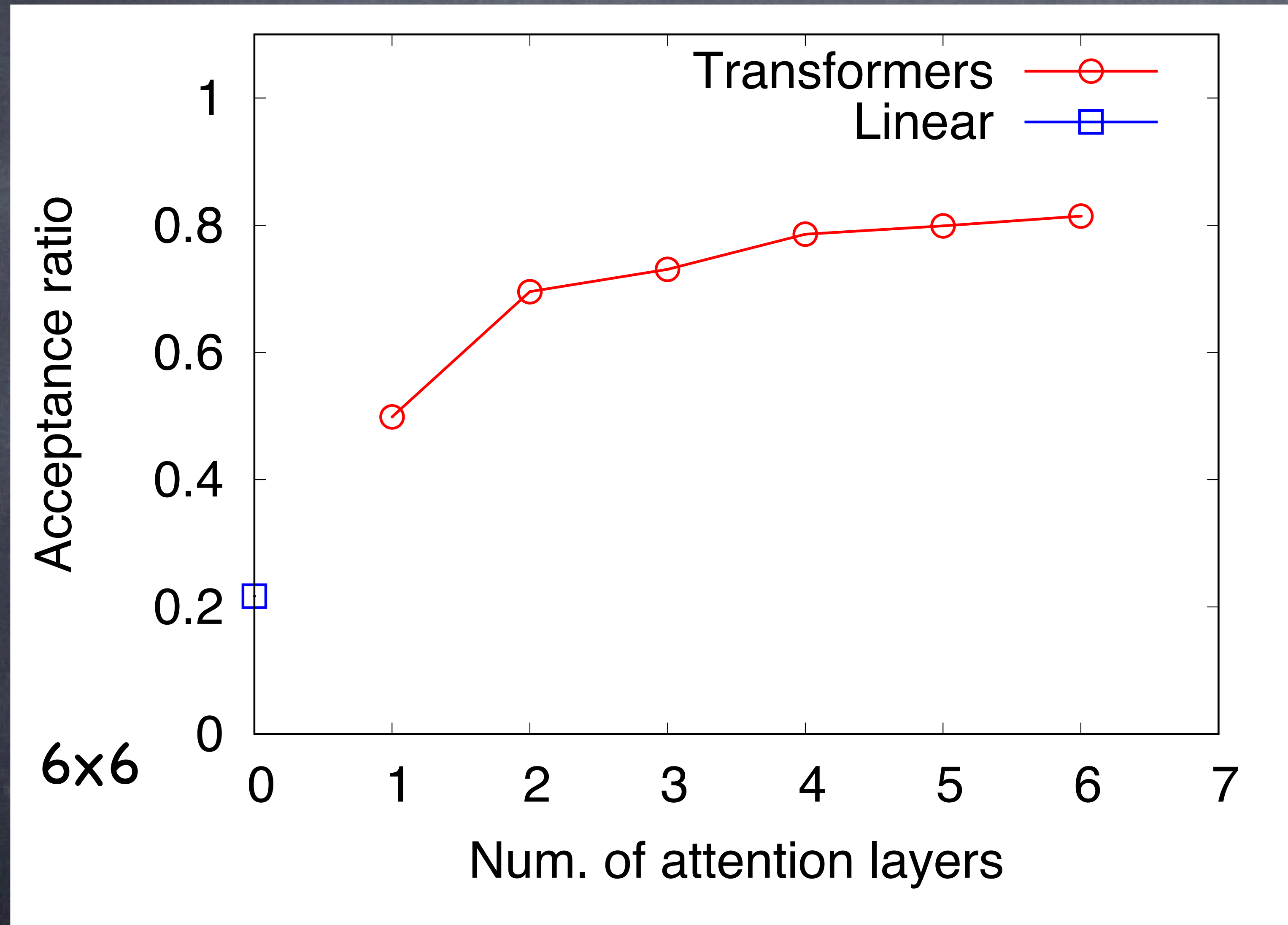
magnetization
and staggered magnetization



Autocorrelation time is reduced

Results

N=6



6-th nearest neighbors

$$K_i = \sum_l W_l S_{i+l}$$

Num. of parameters per layer

$$7+7+7 = 21$$

Last layer: nearest neighbors

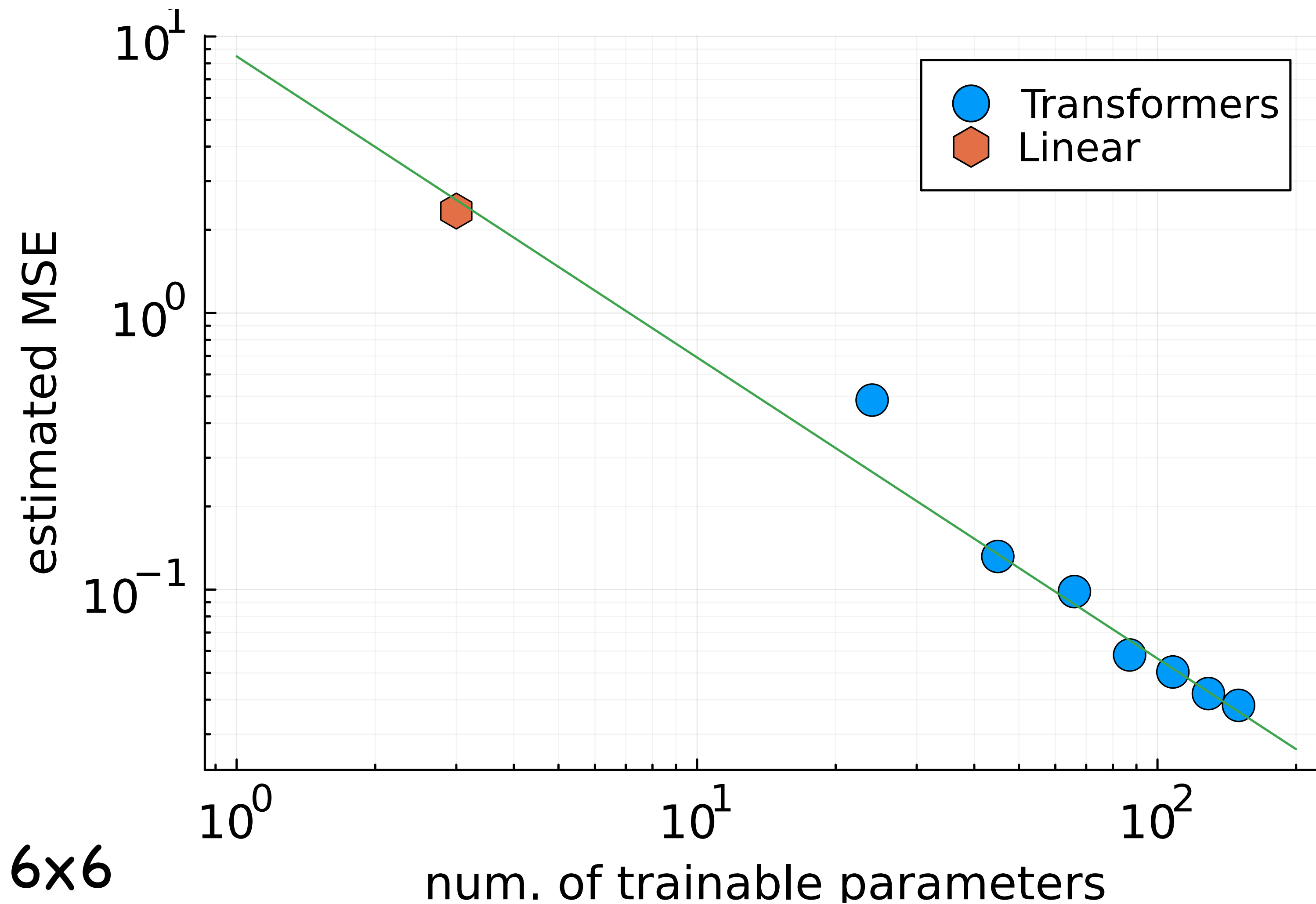
$$E = \sum_i \sum_l J_l \vec{S}_{3i} \cdot \vec{S}_{3i+l}$$

Num. of parameters is small

High acceptance ratio!

Results

arXiv: 2306.11527



6-th nearest neighbors

$$K_i = \sum_l W_l S_{i+l}$$

Num. of parameters per layer

$$7+7+7 = 21$$

Scaling law?

This is like the scaling laws in
Large Language Models

This is MC simulation
We generate data as we want

Application to LatticeQCD

2. We make three variables K, Q, V from A

$$K = W^K A, \quad Q = W^Q A, \quad V = W^V A$$

We introduced "operators" $\hat{S}^Q = \bar{W}^Q \hat{S}$  Effective gauge field U^Q is needed

3. We generate new vector/matrix/tensor B

$$B_l = A_l + \sum_i P_i^l V_i \quad P = \sigma(QK^T) \quad \text{correlation between } Q \text{ and } K$$

We introduced inner product of spins  What is "inner product" in gauge field?

We can use Wilson loop as an inner product

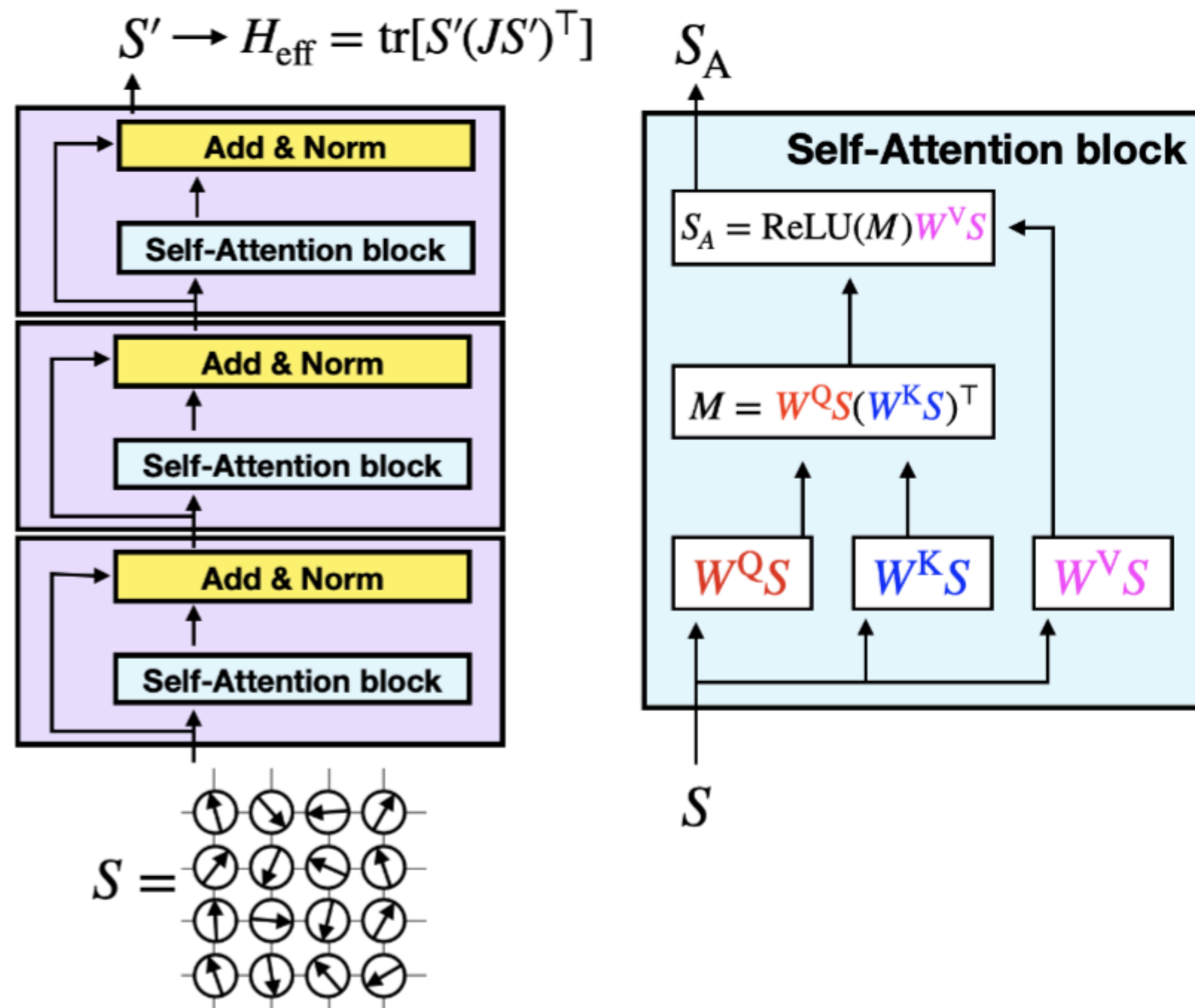
Yuki Nagai, Hiroshi Ohno, Akio Tomiya, "CASK: A Gauge Covariant Transformer for Lattice Gauge Theory", PoS, DOI:10.22323/1.466.0030, arXiv:2501.16955

Summary

Summary

Yuki Nagai and Akio Tomiya, “Self-Learning Monte Carlo with Equivariant Transformer”, J. Phys. Soc. Jpn. 93, 114007 (2024)

Equivariant Transformer in spin systems



Equivariant with respect to spin-rotational and translational symmetries

We found the scaling law!

We can improve models with increasing num. of layers

“Transformer and Attention” is very useful!