

Neural Network Preconditioners for 2D U(1) Wilson-type Dirac Operators

Brian Xiao, Phiala Shanahan, Daniel Hackett, Salvatore Cali, and Yin Lin

Massachusetts Institute of Technology



Introduction

A large part of the cost of a lattice QCD computation: calculation of propagators

- Requires solving a Dirac operator D onto sources
- Large systems \rightarrow direct solution is impractical; must use iterative solvers (e.g. CG) instead
- Can improve solver performance by preconditioning (e.g. AMG, IC, Jacobi)

Idea: train a neural network to produce sparse preconditioners based on input operators

- One-time high training cost
- Cheap preconditioner generation once network is trained

Preconditioning

$$\mathbf{Ax} = \mathbf{b} \quad \Rightarrow \quad \mathbf{M}^{-1}\mathbf{Ax} = \mathbf{M}^{-1}\mathbf{b}$$

$$\kappa(\mathbf{A}) \quad \Rightarrow \quad \kappa(\mathbf{M}^{-1}\mathbf{A}) \ll \kappa(\mathbf{A})$$

Solver convergence controlled by condition number κ

Neural network: maps $\mathbf{A} \rightarrow \mathbf{M}^{-1}$

Loss function: $\kappa(\mathbf{M}^{-1}\mathbf{A})$ or suitable substitute

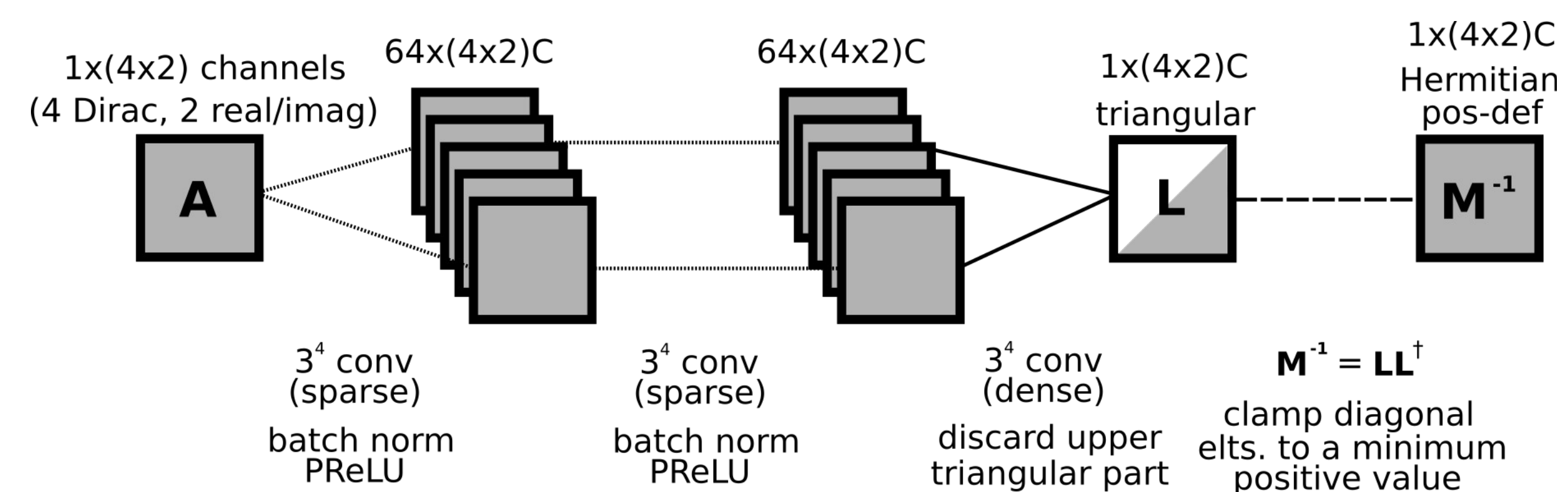
Desirable characteristics for \mathbf{M}^{-1} :

- approximates inverse of \mathbf{A}
- computationally cheap to apply (e.g. sparse)

Another simple preconditioning scheme: even-odd. Based on subspace decomposition; cuts problem difficulty by $\sim 50\%$

Architecture & Training

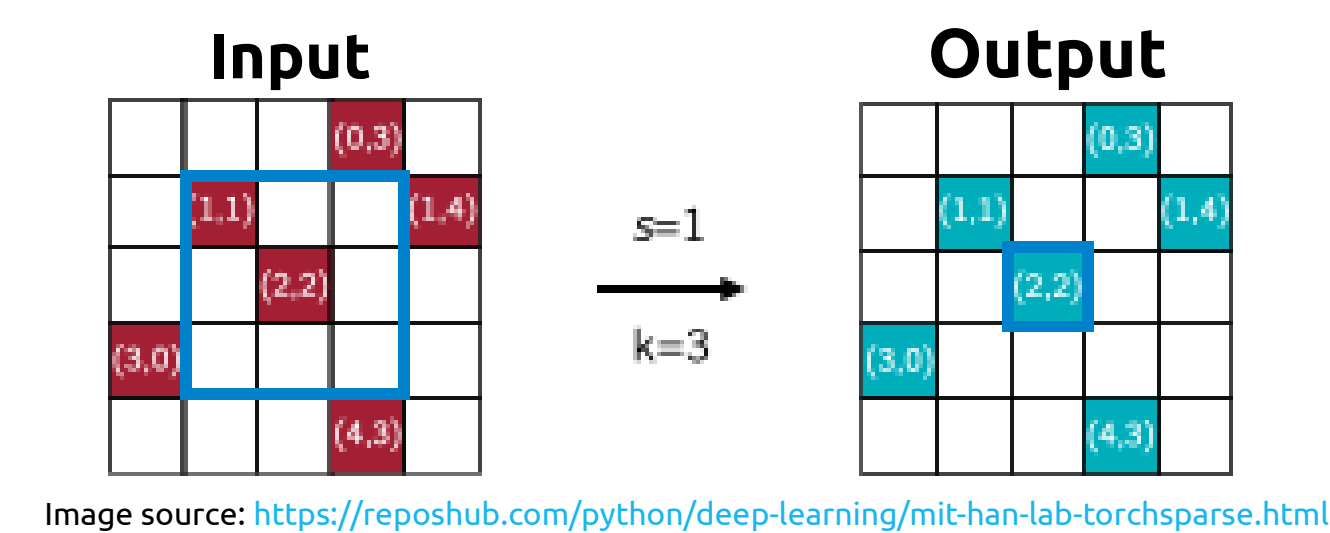
Network architecture:



- 4D convolutions over spatial indices (\mathbf{x}, \mathbf{y}) ; Dirac indices and real/imag parts in channels

- 2D: γ_μ is $2 \times 2 \rightarrow$ Dirac indices $\mu, \nu = 0, 1$

Sparse convolutions:



Each output entry: weighted sum of nearby input entries, parameterized by trainable convolution kernel

Sparse convolutions: only update *nonzero* input entries – preconditioner will be sparse if input Dirac operator is sparse

No fixed size: single network can produce preconditioners for any lattice volume

Loss function: K-condition number; cheaper to compute and better training performance than κ

$$K(\mathbf{A}) = \frac{\text{tr } \mathbf{A}}{n(\det \mathbf{A})^{1/n}} \quad \text{loss}(\mathbf{A}, \mathbf{M}^{-1}) = \frac{K[(\mathbf{M}^{-1})^\dagger \mathbf{A} \mathbf{M}^{-1}]}{K(\mathbf{A}^\dagger \mathbf{A})} = \frac{\text{tr}[(\mathbf{M}^{-1})^\dagger \mathbf{A} \mathbf{M}^{-1}]}{(\text{tr } \mathbf{A}^\dagger \mathbf{A}) |\det \mathbf{M}^{-1}|^{2/n}}$$

Training parameters:

With and without even-odd preprocessing

Dirac operator D from pure-gauge configurations; quenched approximation. Two datasets:

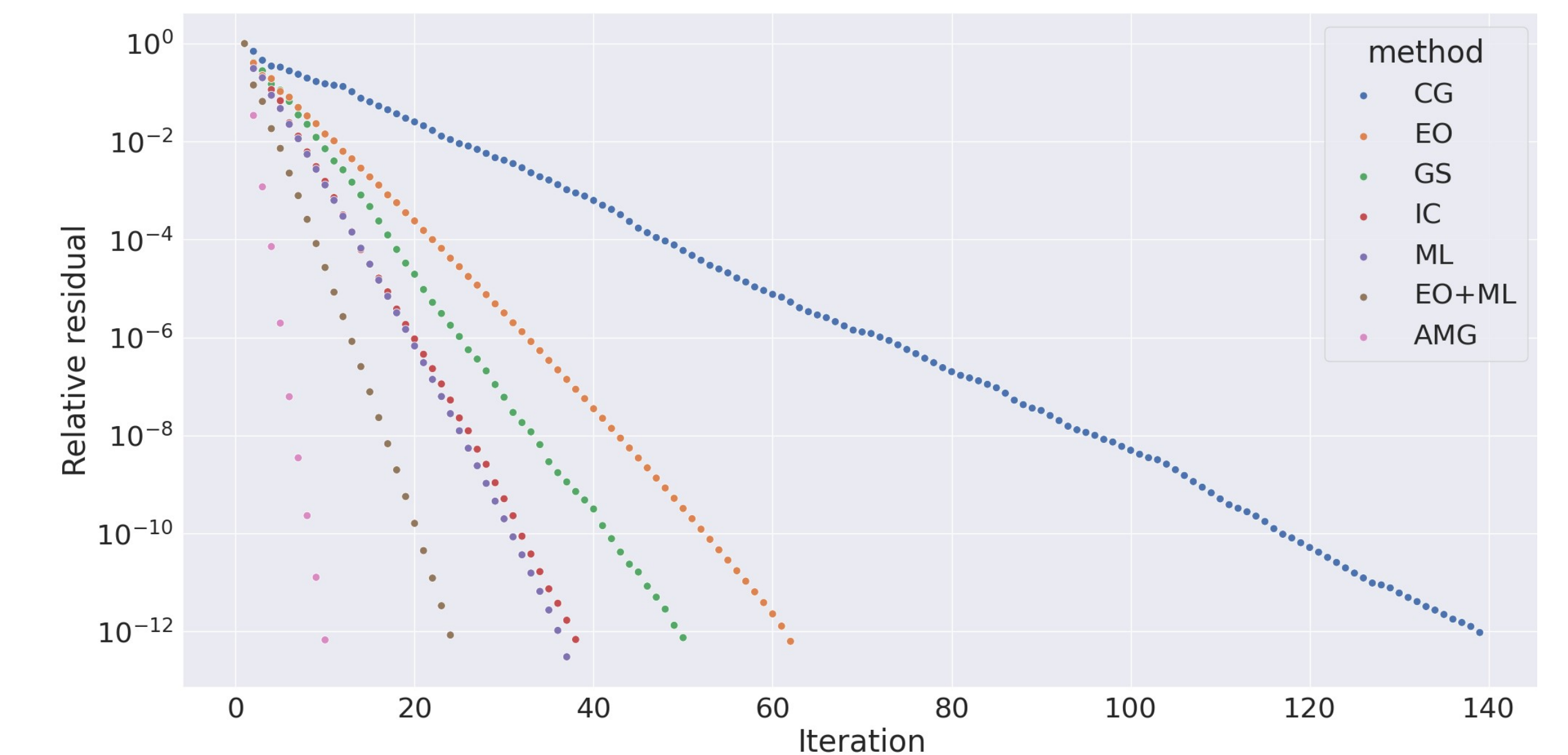
- 1) $\kappa=0.24$: $A = DD^\dagger$ (solved w/ CG)
- 2) $\kappa=0.32$: $A = D\gamma_5$ (solved w/ BiCGSTAB)

- 2800 train / 200 validate / 1000 test

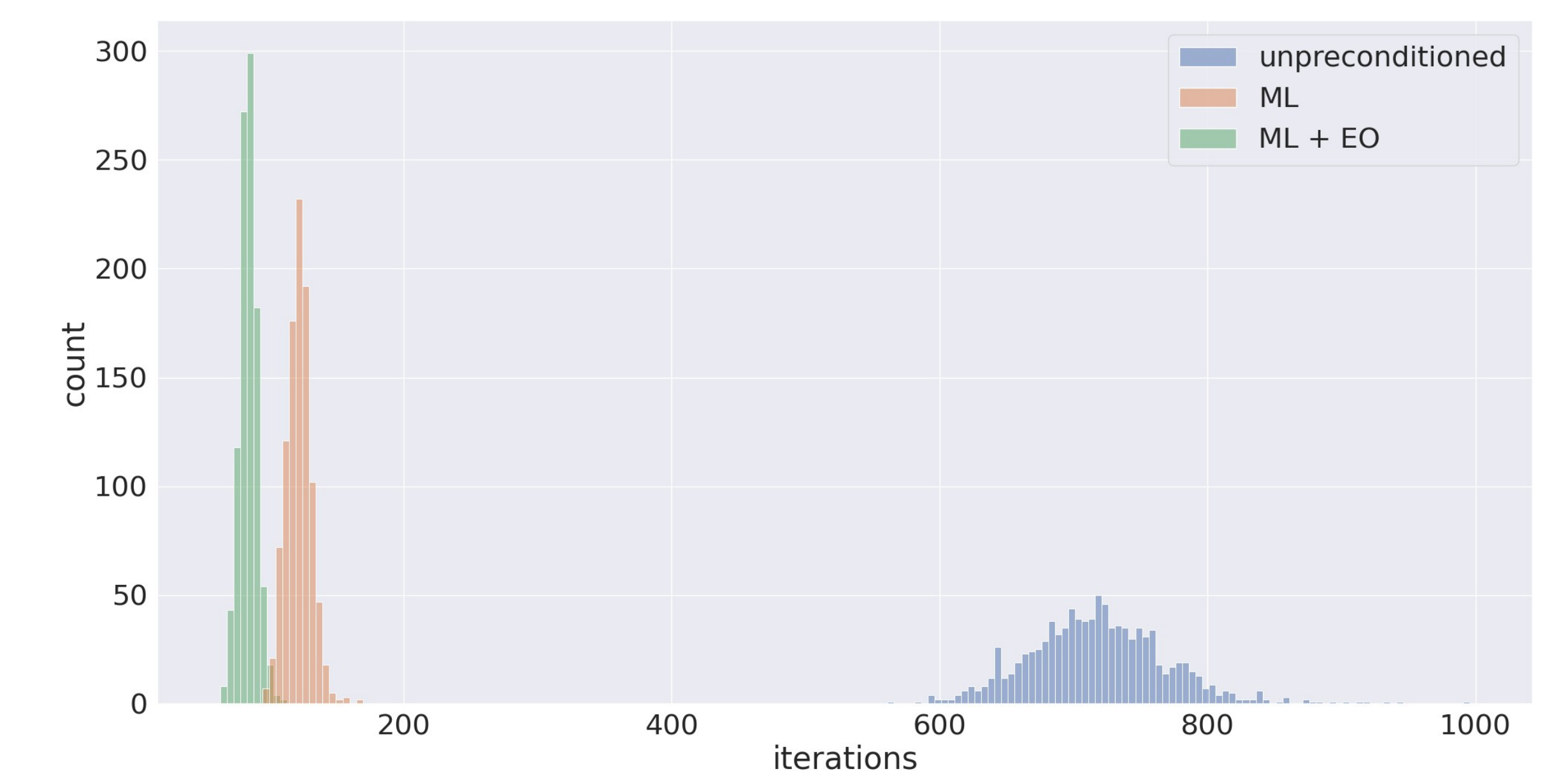
Adam optimizer @ lr 2×10^{-5} , weight decay 10^{-5}

50 epochs; batch size 8×4 processes

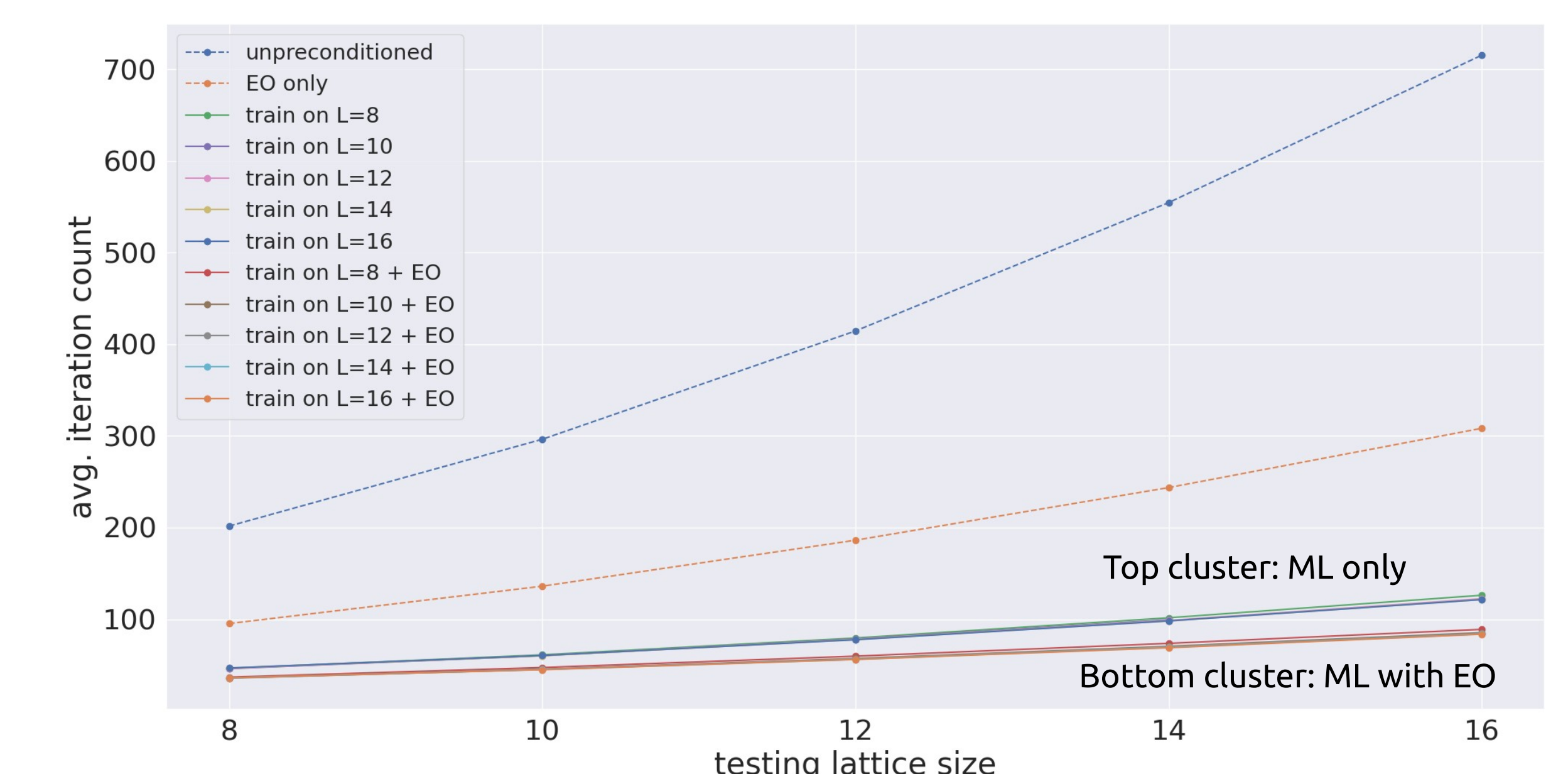
Results



CG convergence, one sample configuration
Dataset: #1; $L = 16$; solver tolerance 10^{-12}



Iteration count histogram
Dataset: #2; $L = 16$



Volume scaling: apply network across different lattice sizes without any retraining
Dataset: #2