# **Preserving Properties of Neural Networks by Perturbative Updates**

Andreas Krämer, Jonas Köhler, Frank Noé

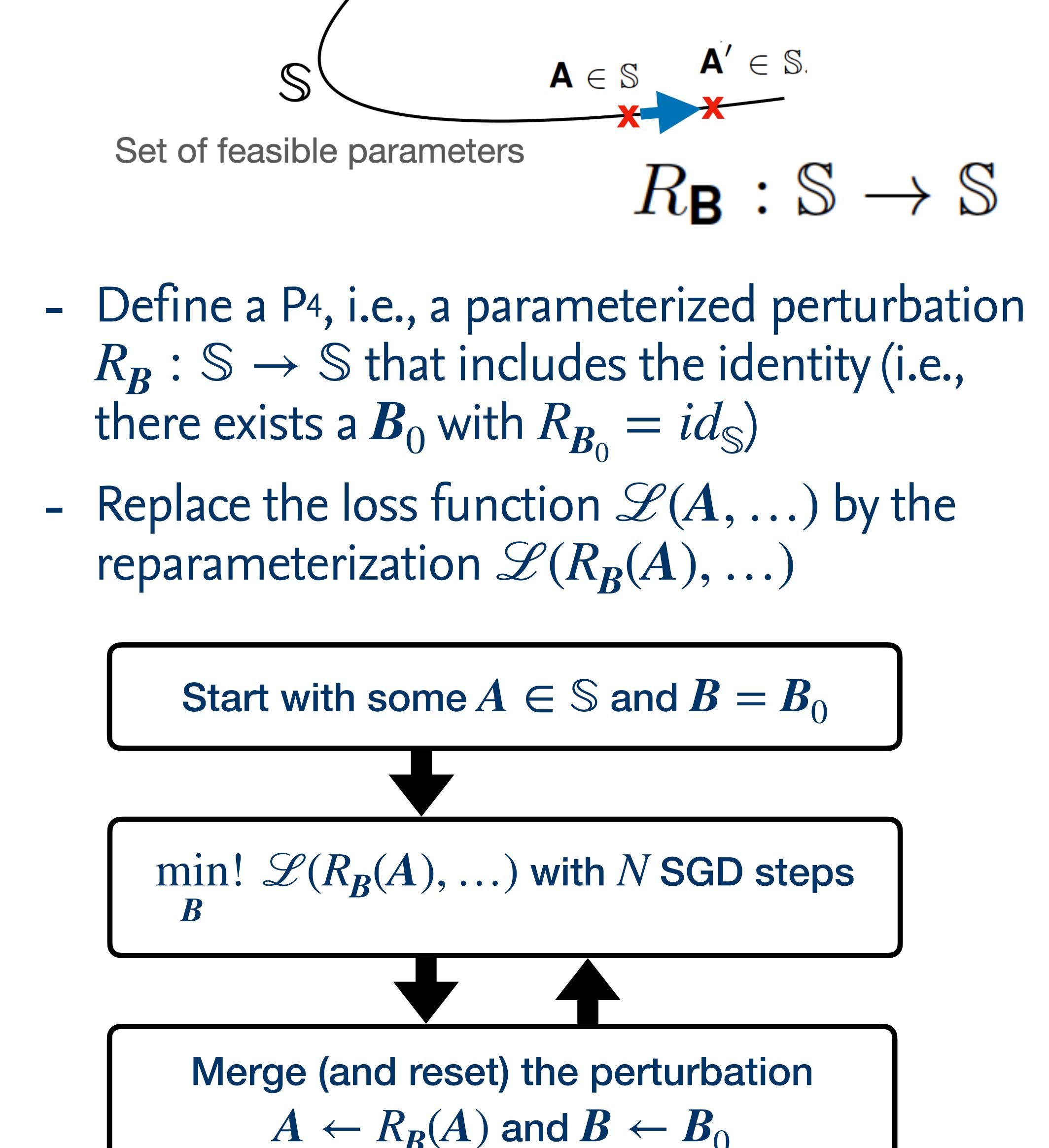
### Introduction

- Machine learning (especially in the physical sciences) often requires networks/layers that obey certain properties
  - Invariances, Equivariances
  - (Lipschitz) Regularity
  - Invertibility, Orthogonality, ...
- Retaining such properties during training can be challenging
- Existing approaches: enforce properties by network design; on-manifold optimization with Riemannian gradient descent; reparameterization; dynamic trivializations
- Idea: train neural networks through property-preserving parameter perturbations (P4)
  - 1. Freeze the network parameters
  - 2. Optimize a perturbation to these parameters that preserves the desired properties
  - 3. In regular intervals: merge the perturbation into the frozen parameters
- Here we train invertible linear layers through rank-one perturbations (P4Inv)

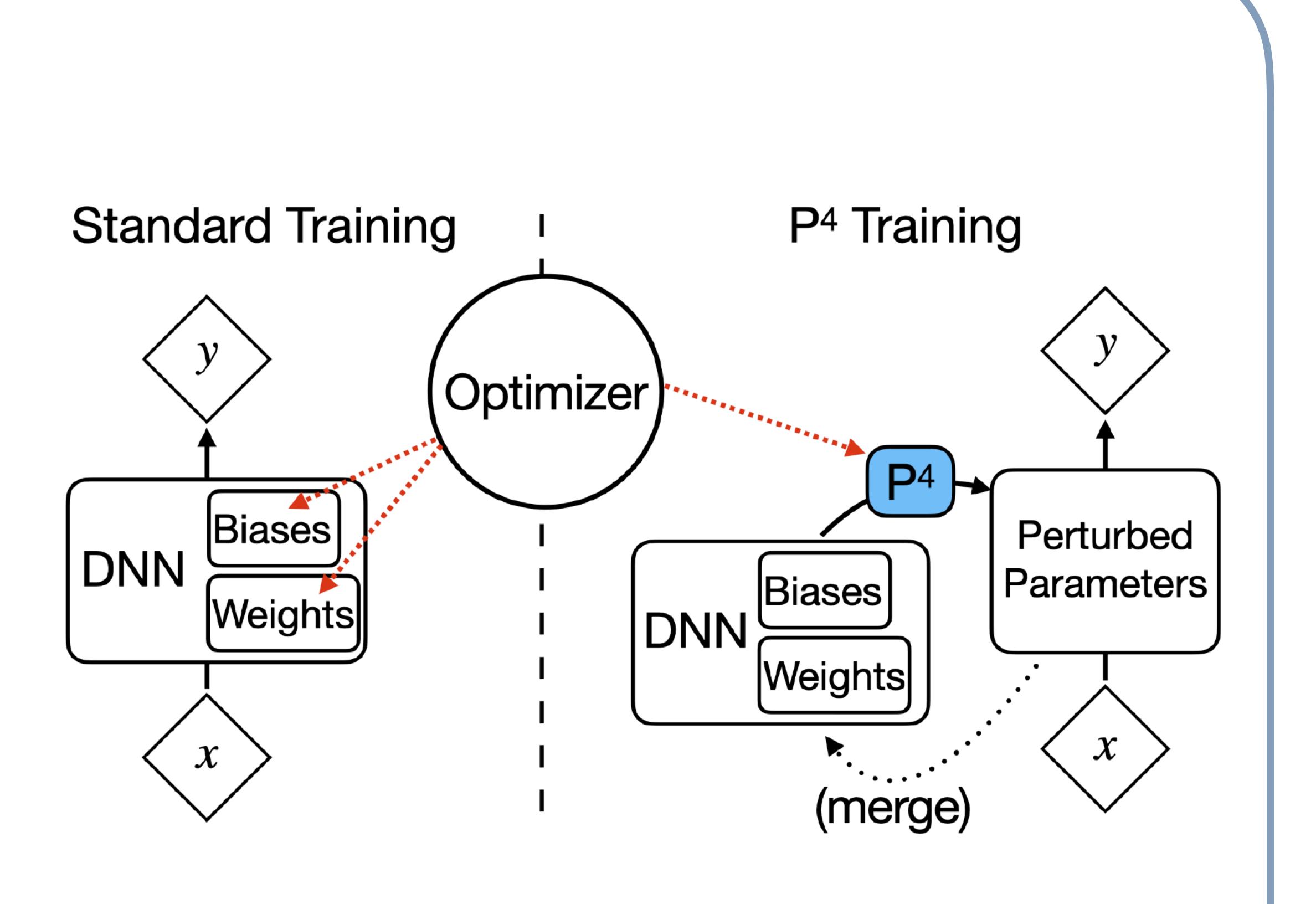
### Method



Parameter Space



arXiv preprint arXiv:2010.07033

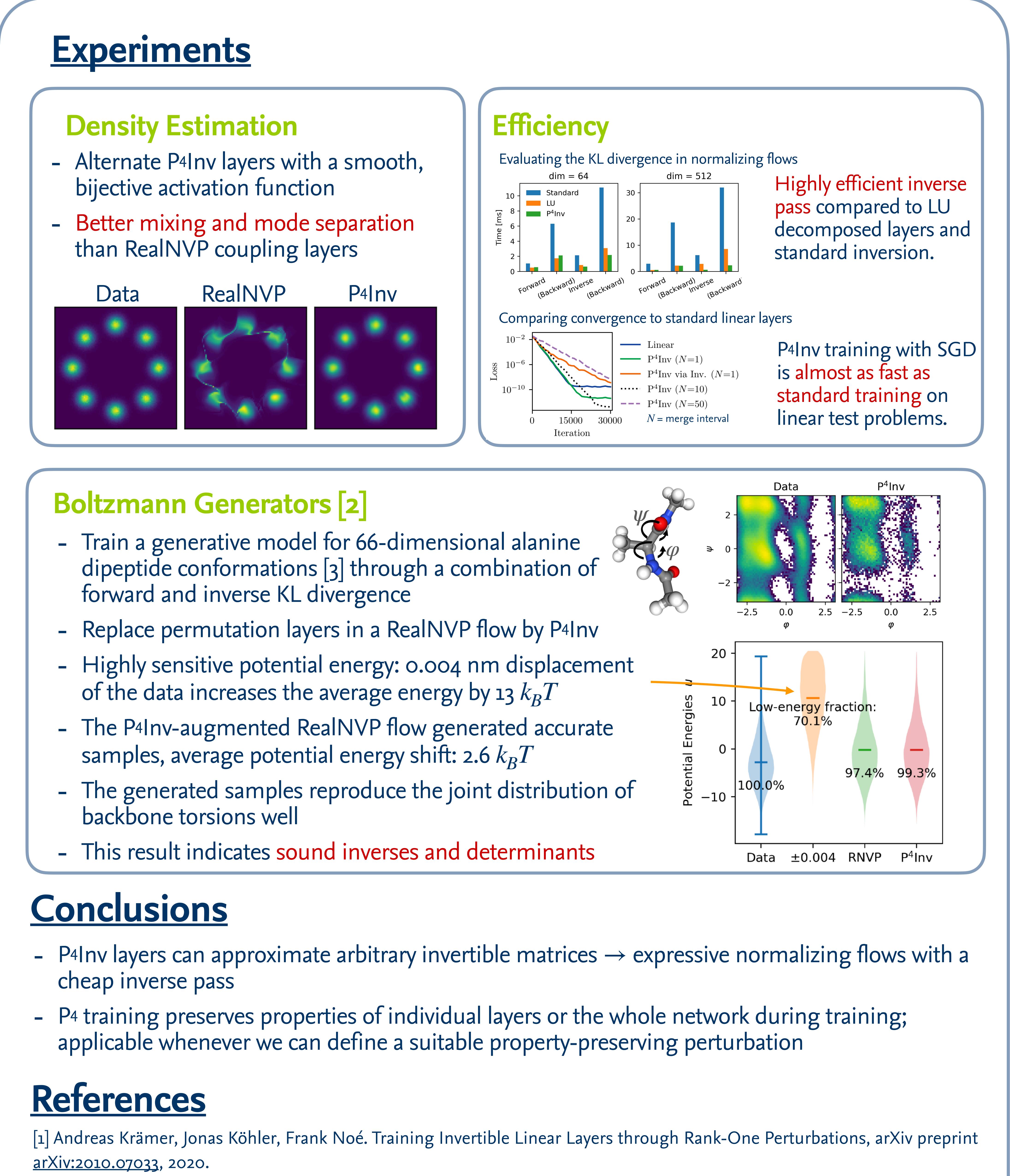


#### **P4Inv: Preserving Invertibility**

- Keep track of inverses and determinants
- Can be used in normalizing flows, reversible generative models that require tractable computation of inverse and Jacobian determinant
- S := set of invertible matrices with known inverses and determinants
- $R_R$  := rank-one perturbation with B = (u, v)

$$R_{u,v}(A) = A + uv^{T}$$
$$(A + uv^{T})^{-1} = A^{-1} - \frac{1}{1 + v^{T}A^{-1}u}A^{-1}uv^{T}A^{-1}u^{$$

- $B_0$  := reset to  $id_{\mathbb{S}}$  via  $v := 0, u_i \sim \mathcal{N}(0,1)$
- Merging step with numerical stabilization
- Skip merge if the perturbed determinant signals ill conditioning of the update
- Optionally augment the loss function by a penalty term that prevents vanishing determinants
- Infrequently run one iteration of an iterative matrix inversion to remove numerical errors in the inverse



[2] Frank Noé, Simon Olsson, Jonas Köhler, Hao Wu. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. Science, 365(6457):eaaw1147, 2019.

[3] Manuel Dibak, Leon Klein, and Frank Noé. Temperature-steerable flows. Machine Learning and the Physical Sciences, Workshop at the 34th Conference on Neural Information Processing Systems (NeurIPS), arXiv preprint arXiv:2012.00429, 2020.

## Freie Universität



Al<sub>4</sub>Science Group

Department of Mathematics and Computer Science