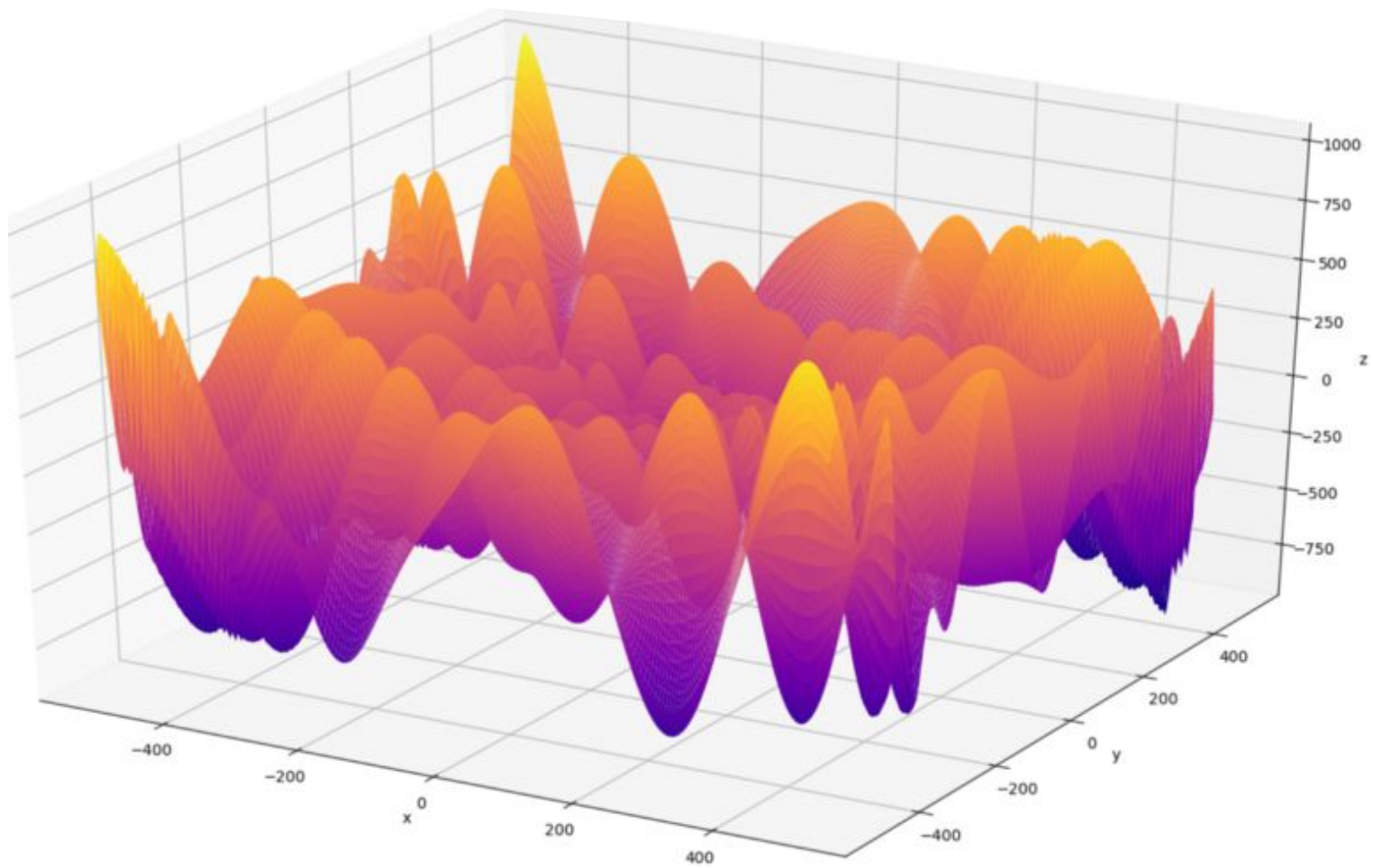
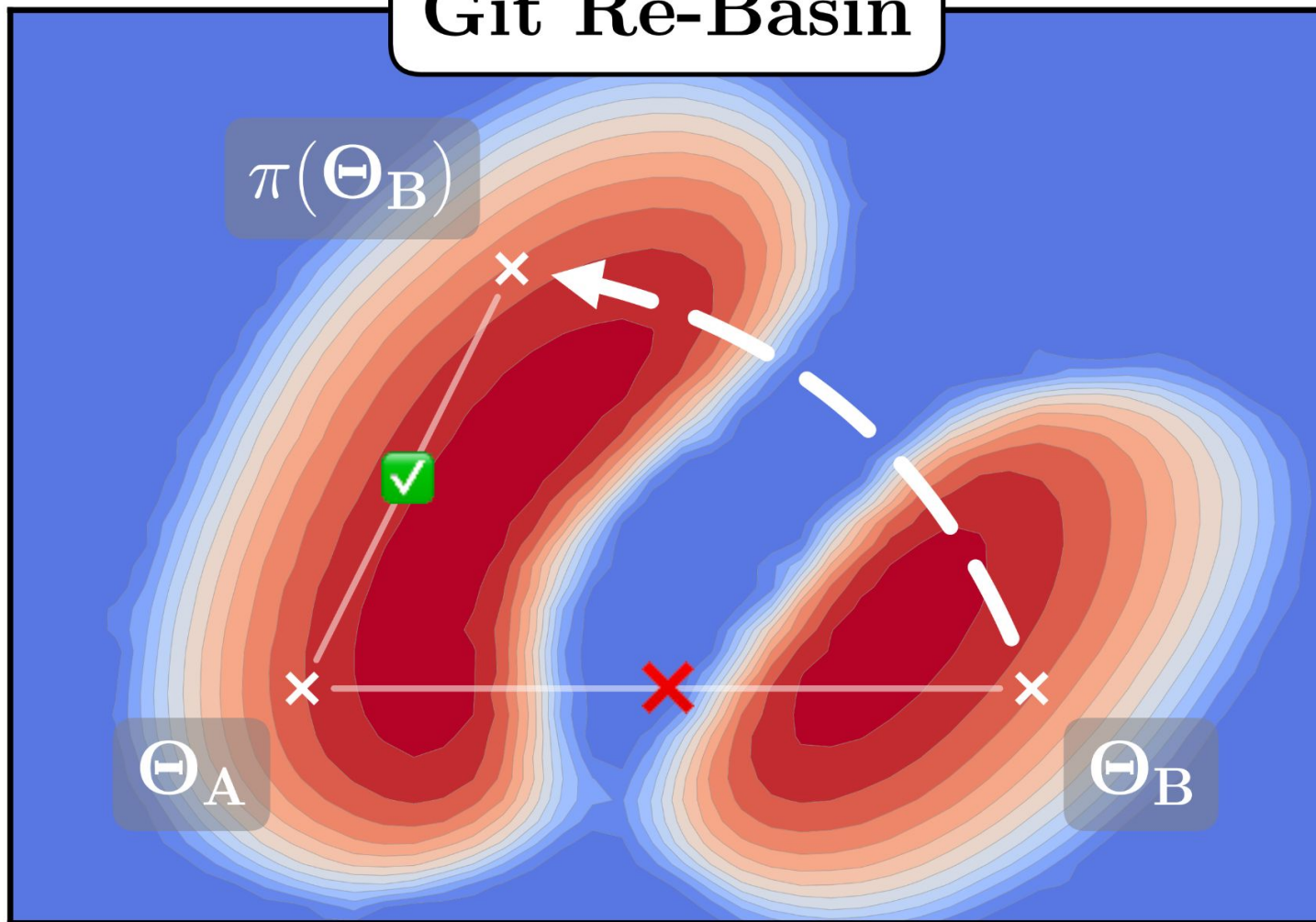


Git Re-Basin: Merging Models modulo Permutation Symmetries

Samuel K. Ainsworth, Jonathan Hayase, Siddhartha Srinivasa

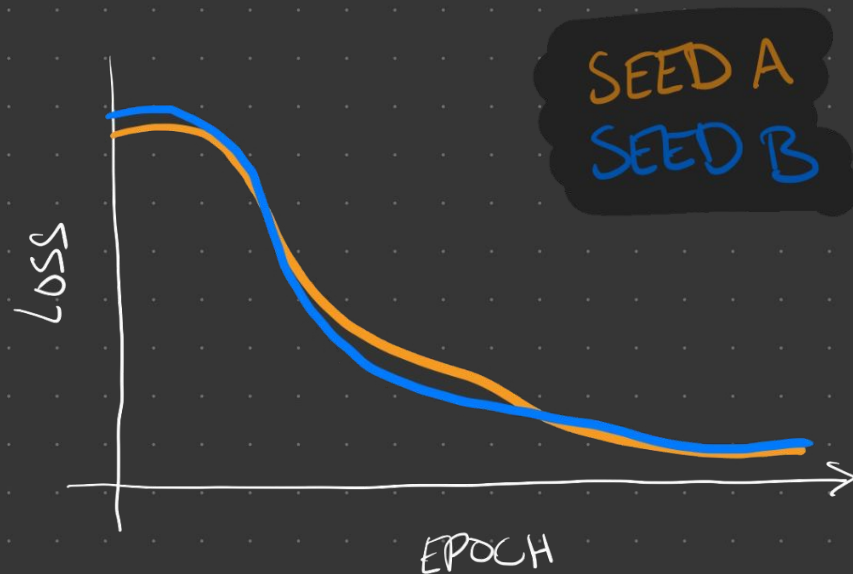


Git Re-Basin



Unreasonable effectiveness of SGD

1. Why does SGD work in deep learning but fail elsewhere (policy learning, traj. opt., recommender systems)?
2. Where are all the local minima?
3. Why do independently trained models have the same training dynamics?



- DIFFERENT BATCH ORDERS!
- DIFFERENT INITIALIZATIONS!

Idea: permutation symmetries to blame?



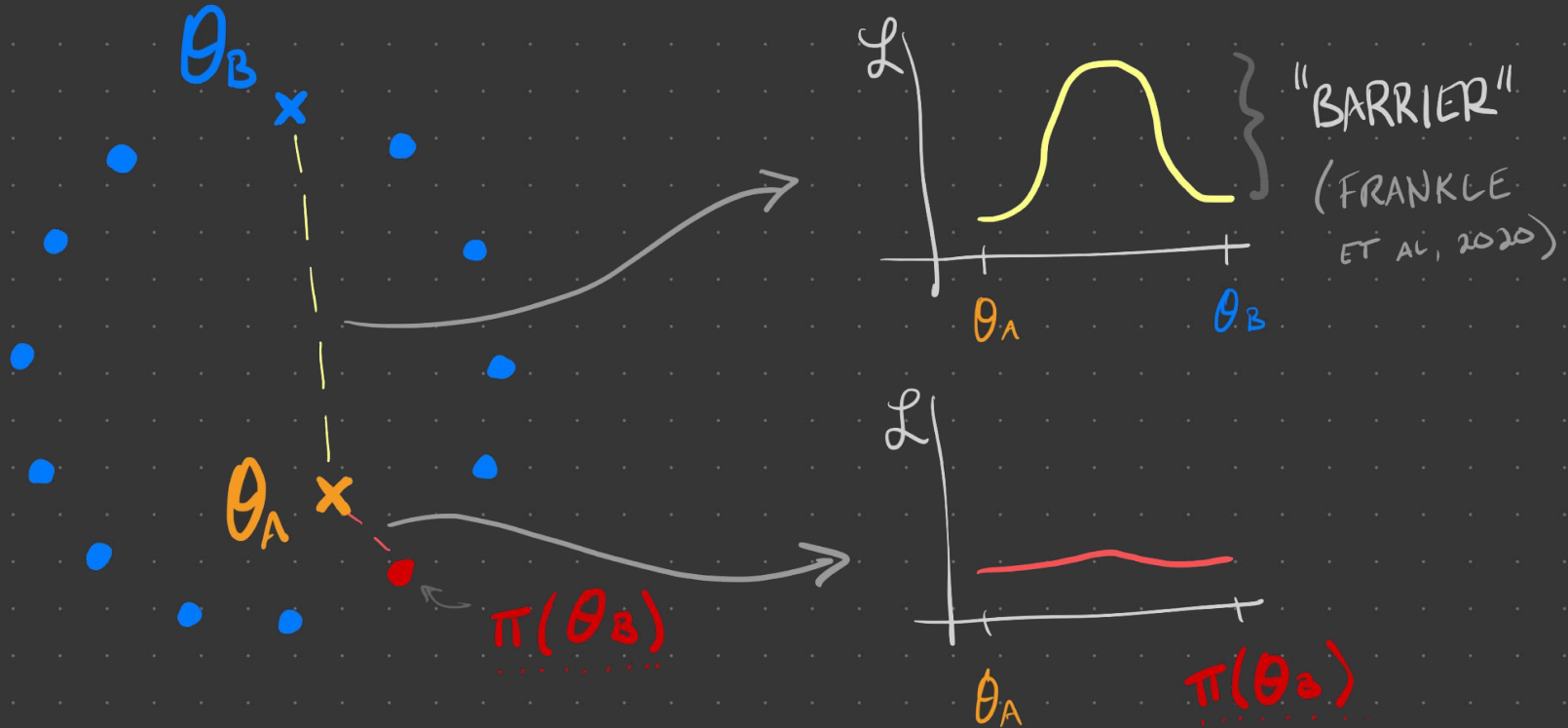
swap hidden units: functionally-equivalent models, different parameters.

ResNet50: 10^{55109}

atoms in universe: 10^{82}

See also *Entezari et al, 2021!*

Idea: permutation symmetries to blame?



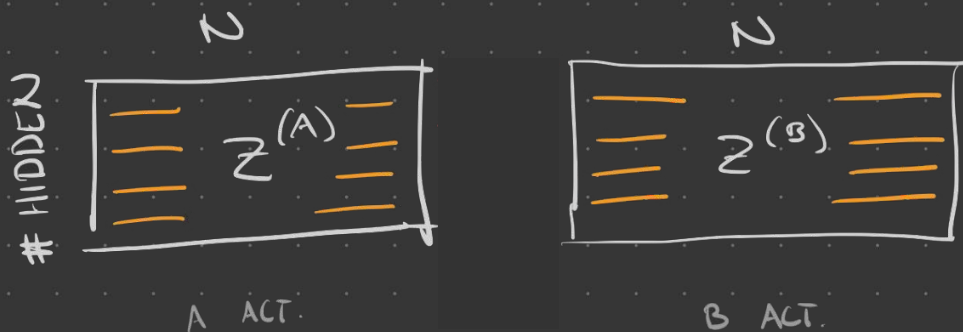
How to find π ?

methods

Alg. 1: Activation matching

IDEA: SIMILAR UNITS OUGHT TO HAVE CORRELATED ACTIVATIONS!

ALGORITHM: AT EACH LAYER, COLLECT ACTIVATIONS FROM BOTH



OLS CONSTRAINED TO S_d :

$$\operatorname{argmin}_{P \in S_d} \sum_{i=1}^N \| z_{:,i}^{(A)} - P z_{:,i}^{(B)} \|^2$$

$$= \operatorname{argmax}_{P \in S_d} \left\langle P, z^{(A)} z^{(B)T} \right\rangle_F$$

"LINEAR ASSIGNMENT PROBLEM"

\Rightarrow SOLVE w/ HUNGARIAN ALGO, ETC.

Alg. 2: Weight matching

IDEA: SIMILAR UNITS SHOULD TO HAVE SIMILAR WEIGHTS.

$$\operatorname{argmin}_{\pi} \underbrace{\| \text{vec}(\theta_A) - \text{vec}(\pi(\theta_B)) \|^2}_{L2} = \operatorname{argmax}_{\pi} \underbrace{\text{vec}(\theta_A) \cdot \text{vec}(\pi(\theta_B))}_{\text{INNER PROD.}}$$

$$= \operatorname{argmax}_{\pi = \{P_i\}} \langle W_1^{(A)}, P_1 W_1^{(B)} \rangle_F + \langle W_2^{(A)}, P_2 W_2^{(B)} P_1^T \rangle_F + \dots + \langle W_L^{(A)}, W_L^{(B)} P_{L-1}^T \rangle_F$$

LEMMA: THIS IS NP-HARD!

ALGORITHM:

- RANDOMLY PICK P_i ,
- SOLVE JUST THAT ONE (REDUCES TO LINEAR ASSIGNMENT PROBLEM),
- REPEAT UNTIL CONVERGENCE.

Alg. 3: Straight-through estimator

IDEA: LEARN $\pi(\theta_B)$ BASED ON THE DATA DISTRIBUTION.

$$\min_{\tilde{\theta}_B} \mathcal{L}\left(\frac{1}{2}(\theta_A + \text{proj}(\tilde{\theta}_B))\right) \quad \text{proj}(\theta) \triangleq \underset{\pi(\theta_B)}{\text{argmax}} \text{VEC}(\theta) \cdot \text{VEC}(\pi(\theta_B))$$

PROBLEM: proj IS NOT DIFFERENTIABLE...

⇒ USE STRAIGHT-THRU ESTIMATOR!

⇒ SOLVE proj W/ WEIGHT MATCHING ALGO,

FWD PASS: USE $\text{proj}(\tilde{\theta}_B)$

BWD PASS: USE $\tilde{\theta}_B$

Bonus: Merging more than 2 models?

Algorithm 3: MERGEMANY

Given: Model weights $\Theta_1, \dots, \Theta_N$

Result: A merged set of parameters $\tilde{\Theta}$.

repeat

for $i \in \text{RANDOMPERMUTATION}(1, \dots, N)$ **do**

$\Theta' \leftarrow \frac{1}{N-1} \sum_{j \in \{1, \dots, N\} \setminus \{i\}} \Theta_j$

$\pi \leftarrow \text{PERMUTATIONCOORDINATEDDESCENT}(\Theta', \Theta_i)$

$\Theta_i \leftarrow \pi(\Theta_i)$

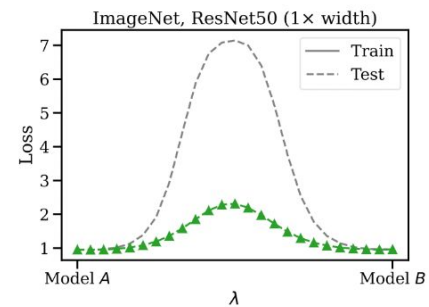
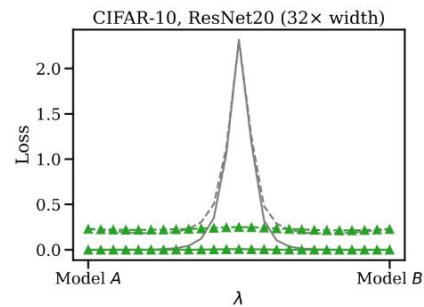
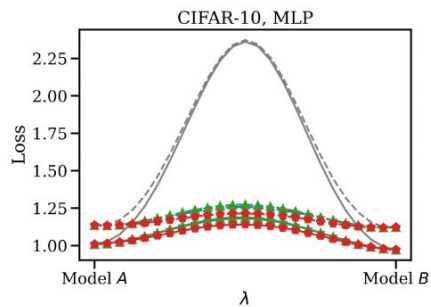
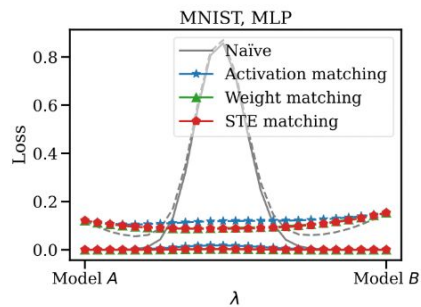
end

until *convergence*

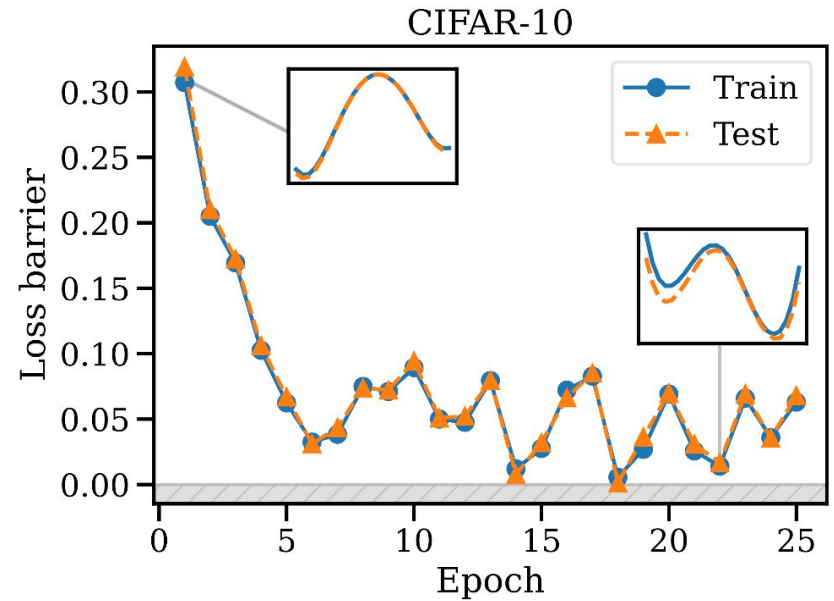
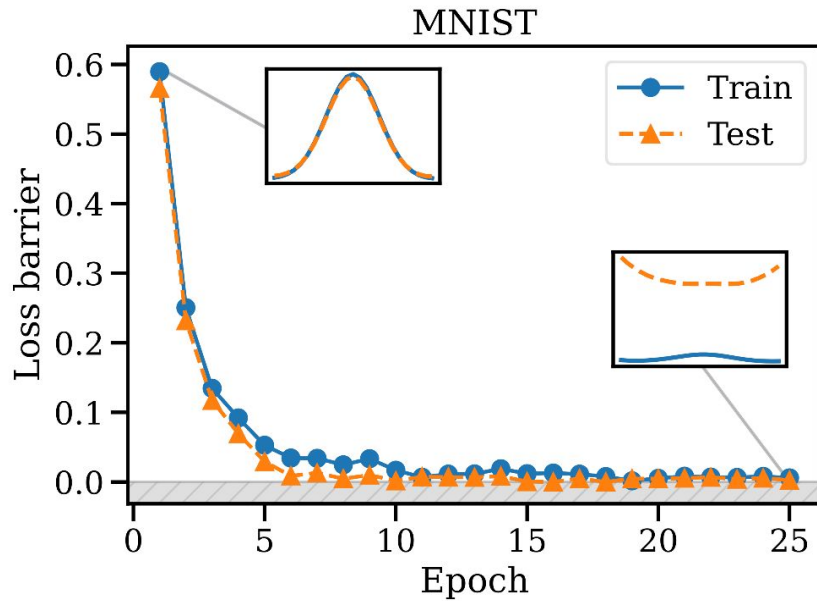
return $\frac{1}{N} \sum_{j=1}^N \Theta_j$

experiments!

LMC before/after matching

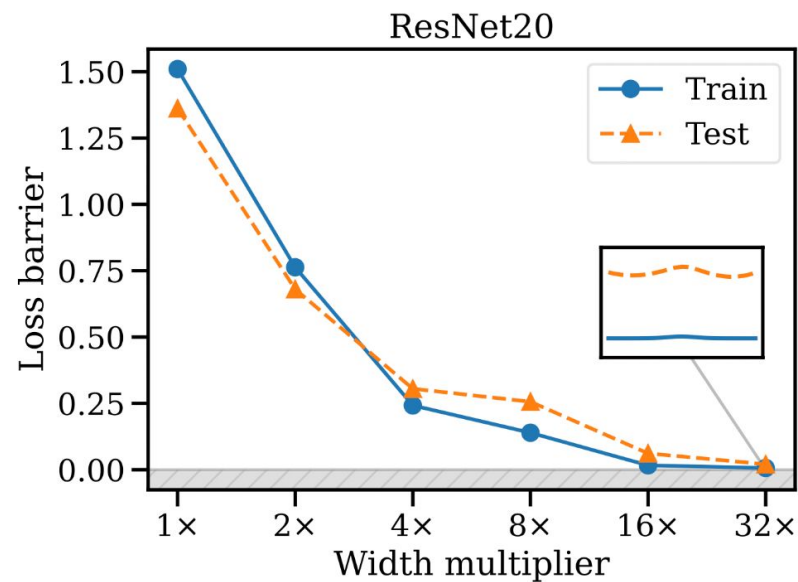
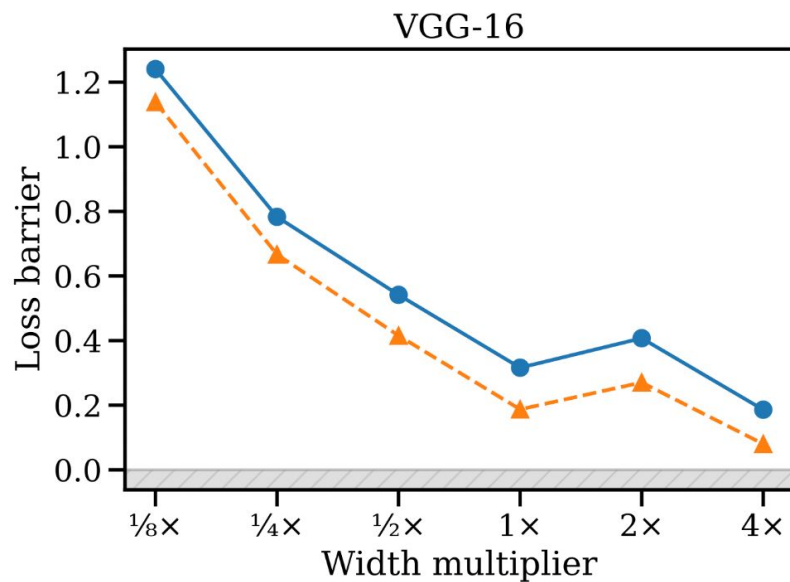


LMC is an emergent property of training

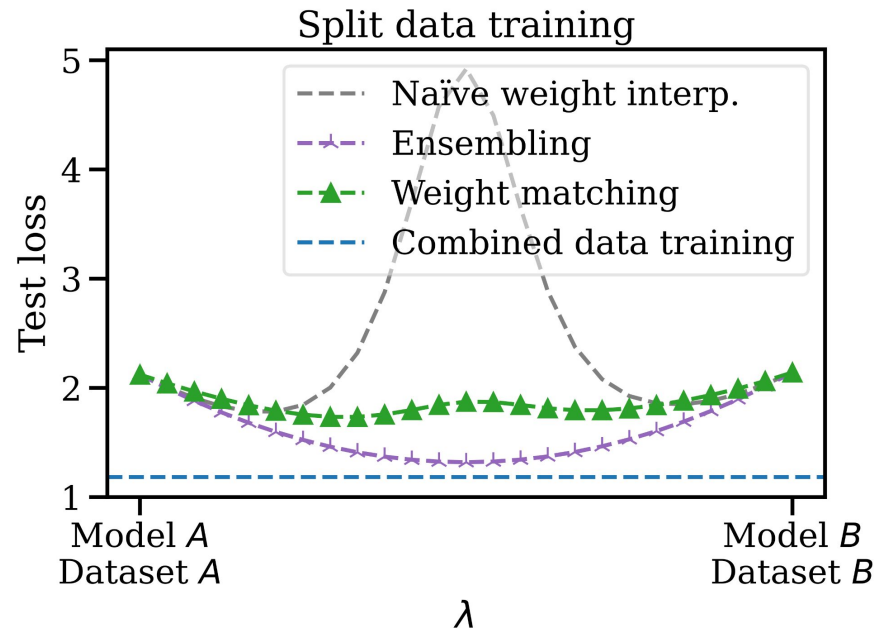


See also Benzing et al, 2022!

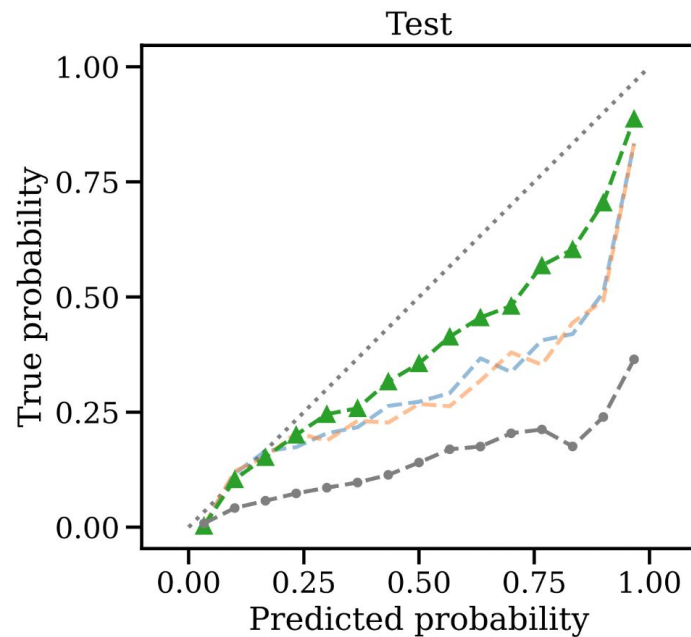
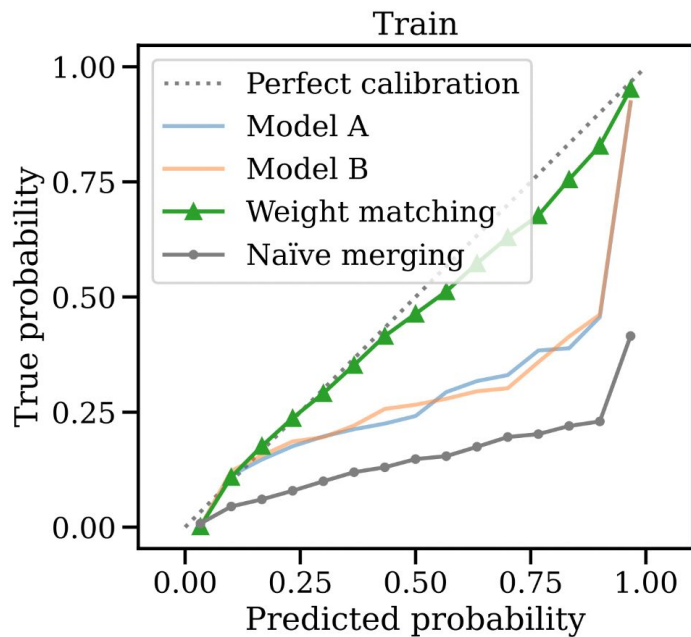
Wider models are better models



Model patching/disjoint datasets



Model patching/disjoint datasets, calibration

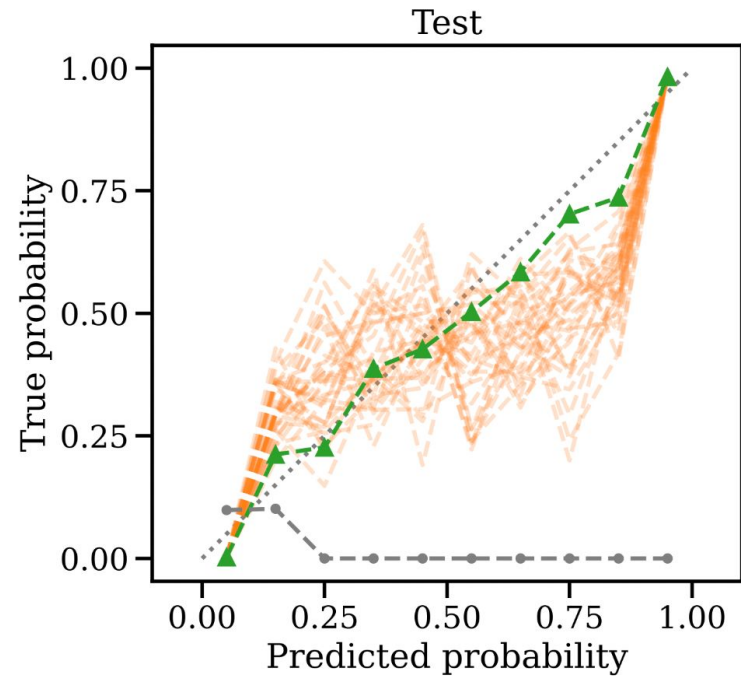
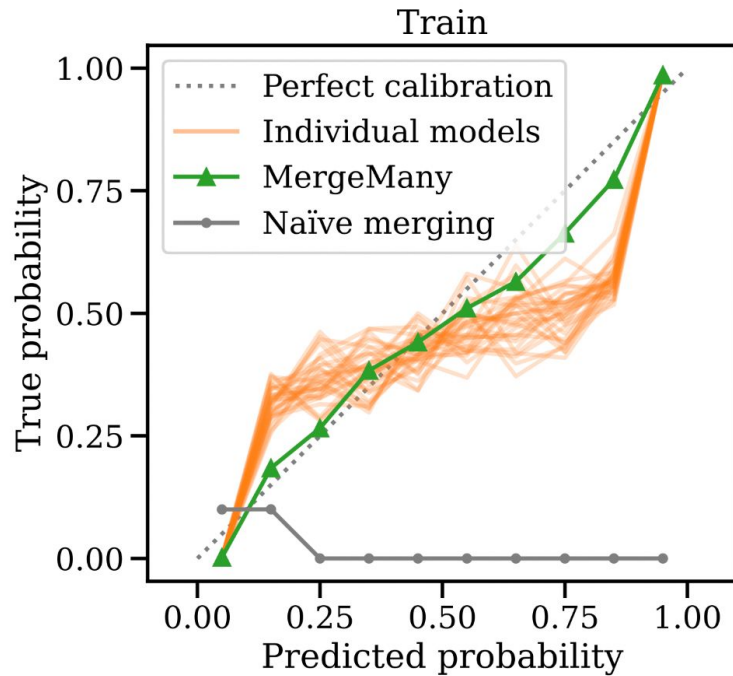


Bonus: MergeMany results

	Train Loss	Train Acc.	Test Loss	Test Acc.
Seed 1	0.0000	1.0000	0.1153	0.9856
Seed 2	0.0000	1.0000	0.1531	0.9854
Seed 3	0.0000	1.0000	0.1229	0.9855
Seed 4	0.0000	1.0000	0.1108	0.9865
Seed 5	0.0000	1.0000	0.1443	0.9871
MERGEMANY	0.0141	0.9952	0.0727	0.9831

43% decrease in test loss!

Bonus: MergeMany results (con't)



conclusion

1.

Loss landscapes seem to contain **only a single basin** mod. permutation symmetries in many settings. And wider is better.

2.

Independently trained models can be merged in **weight space** by teleporting into the same basin.

3.

Merged models are **better calibrated** and tend to outperform single models on test loss. Methodology for improving model performance?

Open questions and future work

- Federated learning? Distributed training?
- Cross validation for deep learning?
- What about thin models? They seem to exhibit similar behavior in training but don't work as well as wide ones...
 - Are thin models just [wide models in superposition](#)? Connection to optimal transport?
 - Hypothesis: Activations between different models can be linearly related. In the infinite width limit it just so happens that a sufficient permutation relationship exists.
- When does Git Re-Basin fail and why? Why is SGD implicitly biased towards solutions that admit LMC?
- Security implications for model merging? How safe is your data?

Thanks to my collaborators!



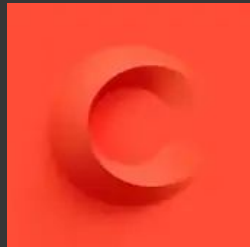
Samuel Ainsworth
(UW, Cruise AI Research)



Jonathan Hayase
(UW)



Sidd Srinivasa
(UW)



thank you!

questions?