



Faculty of Computer Science

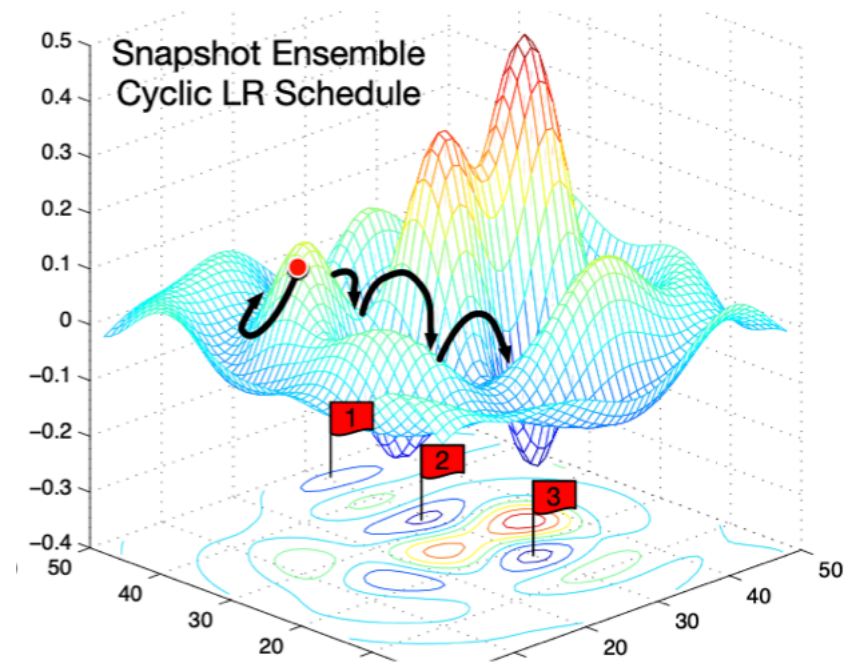
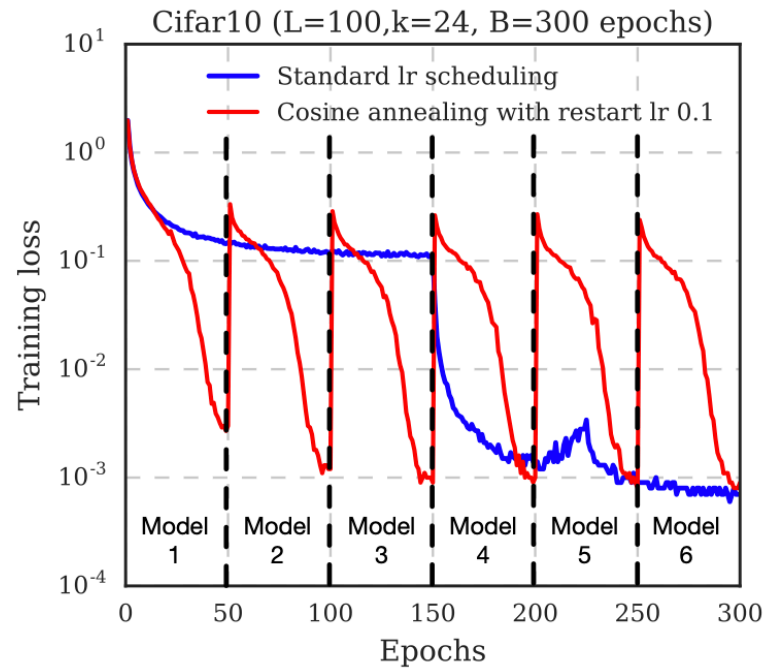
Applied Mathematics &  
Computer Science

Moscow 2024

# Analysis of Neural Networks Internal Representations During Transfer Learning

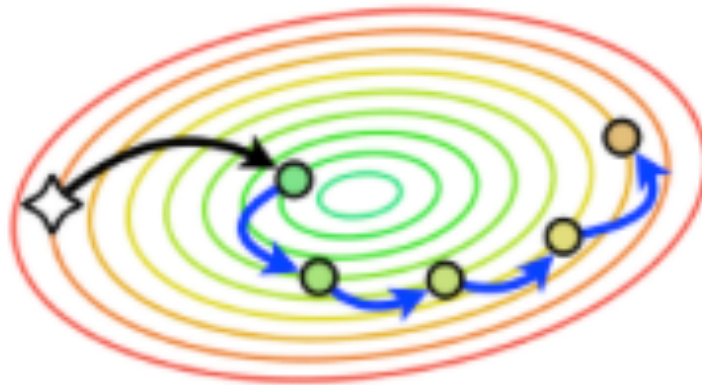
Gritsaev Timofei Grigorievich  
under Ildus Sadrtudinov supervision

# SNAPSHOT ENSEMBLES (SSE): TRAIN 1, GET M FOR FREE

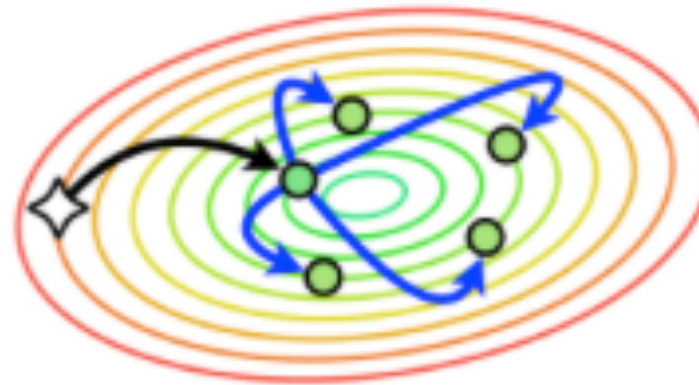


# “To Stay or Not to Stay in the Pre-train Basin: Insights on Ensembling in Transfer Learning”

SSE

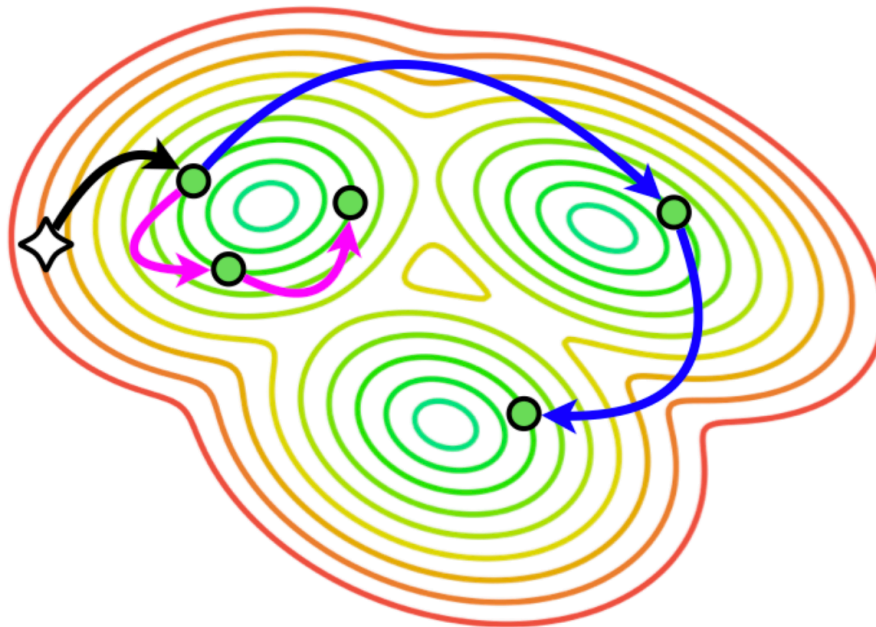


StarSSE



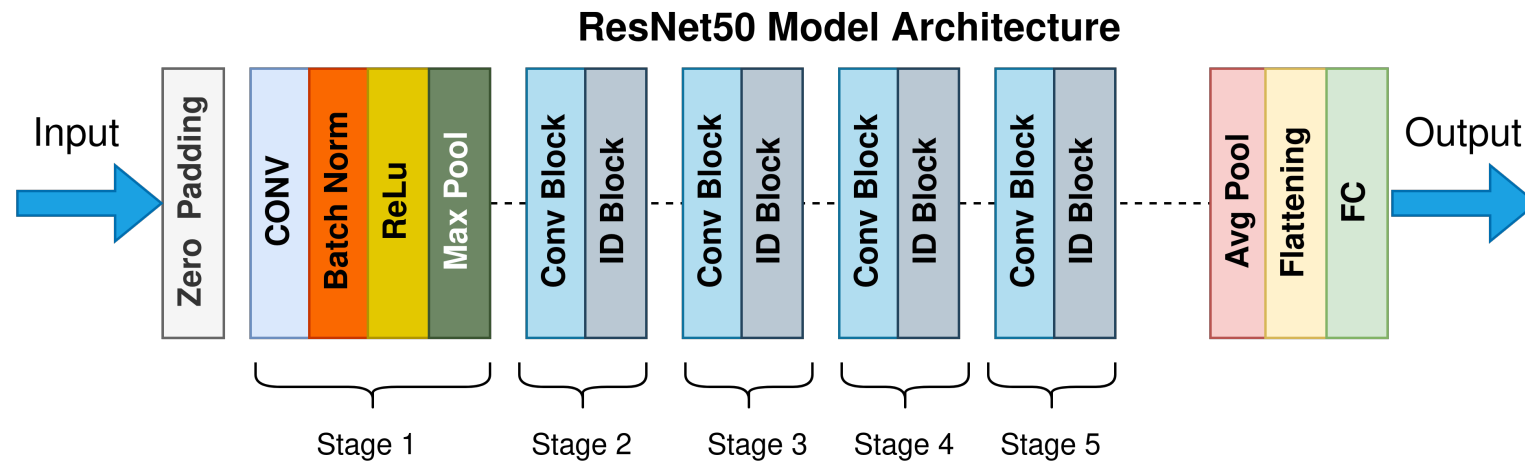


## More local vs Semi-local



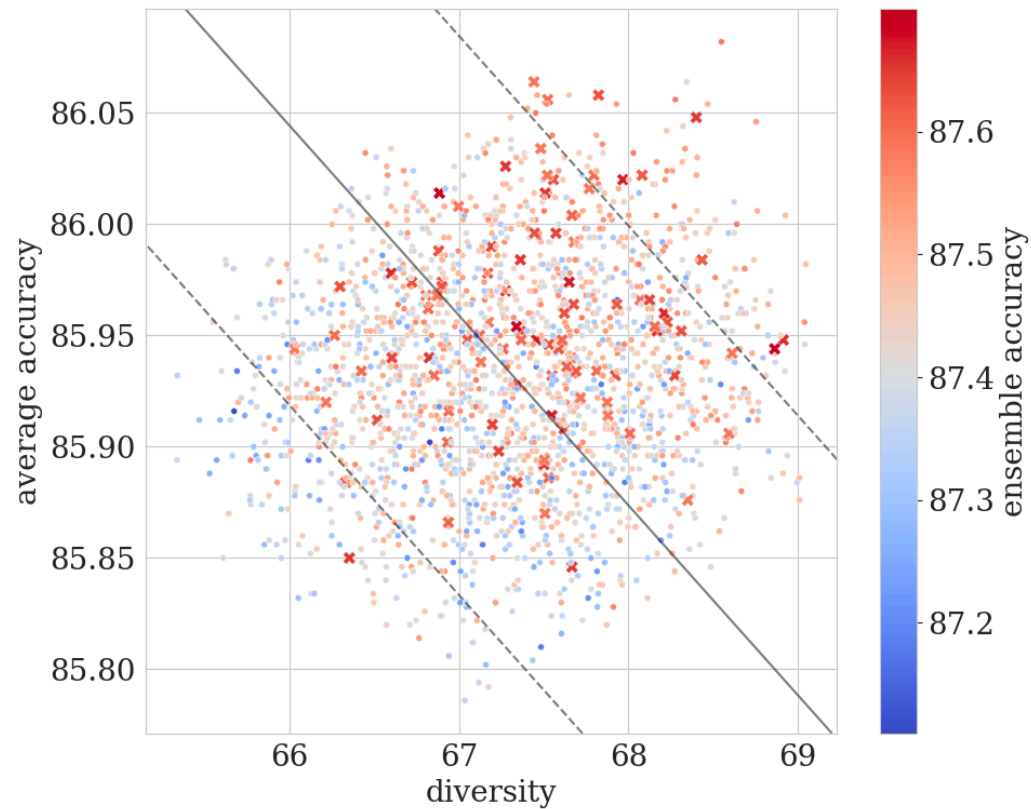
# Methodology

- ResNet50
- CIFAR-100
- Following the protocol from  
“To Stay or Not to Stay in the Pre-train Basin: Insights on Ensembling in Transfer Learning”



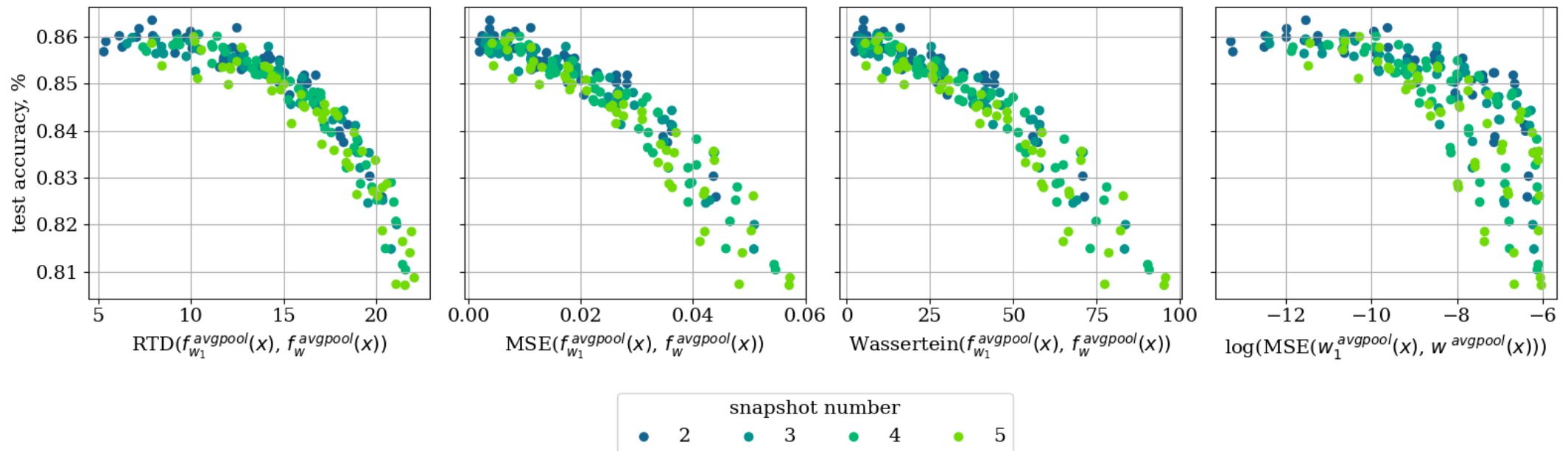
# Average accuracy, diversity, ensemble accuracy

StarSSE



Algorithm		worst 5%	best 5%
<b>StarSSE</b>	avg. acc.	85.91 $\pm$ 0.04	85.95 $\pm$ 0.05
	div.	67.11 $\pm$ 0.78	67.44 $\pm$ 0.6
	<b>ens. acc.</b>	87.24 $\pm$ 0.03	87.62 $\pm$ 0.03

# SSE regularization motivation

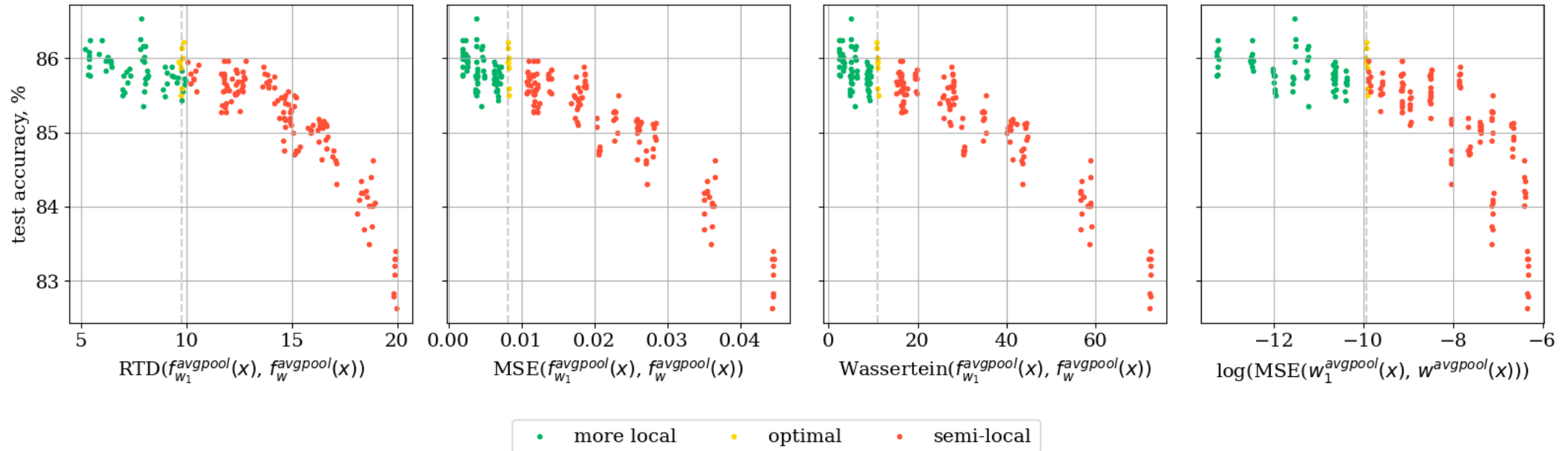


## Regularized SSE results

Algorithm	avg. acc.	div.	ens. acc.
<b>SSE</b>	85.52 $\pm$ 0.34	71.57 $\pm$ 5.94	87.11 $\pm$ 0.06
<b>StarSSE</b>	85.88 $\pm$ 0.25	68.12 $\pm$ 1.15	<b>87.41<math>\pm</math>0.12</b>
<b>SSE-RTD</b>	85.71 $\pm$ 0.07	68.04 $\pm$ 5.78	87.36 $\pm$ 0.09
<b>SSE-MSE</b>	85.81 $\pm$ 0.221	67.62 $\pm$ 5.95	87.37 $\pm$ 0.02



# Analysis of StarSSE internal representations



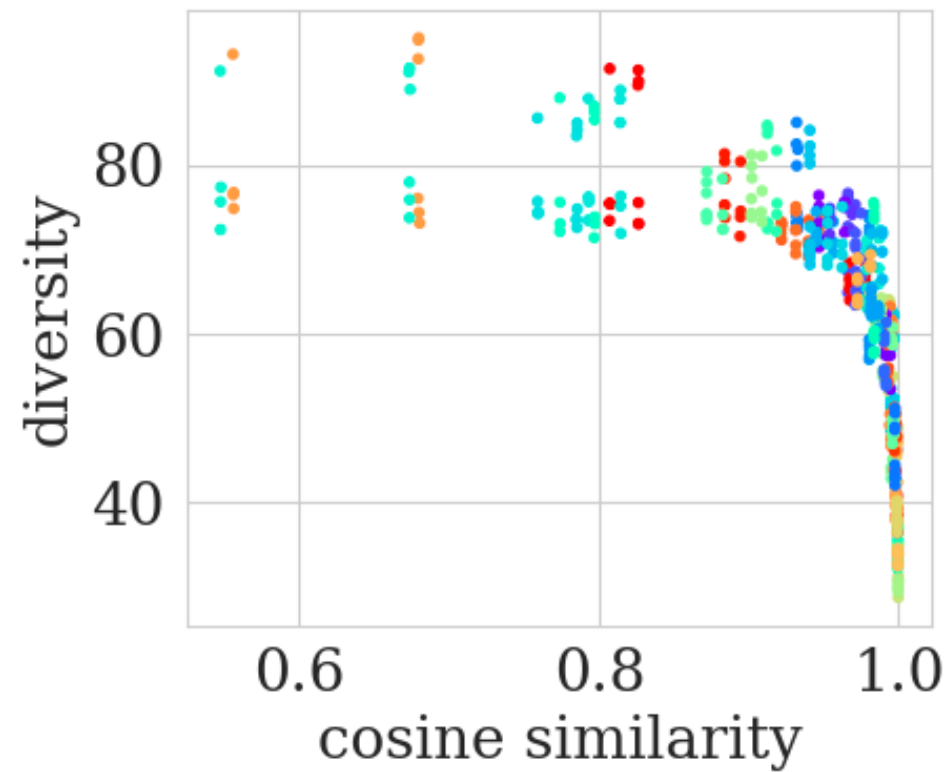


## Pairwise diversity

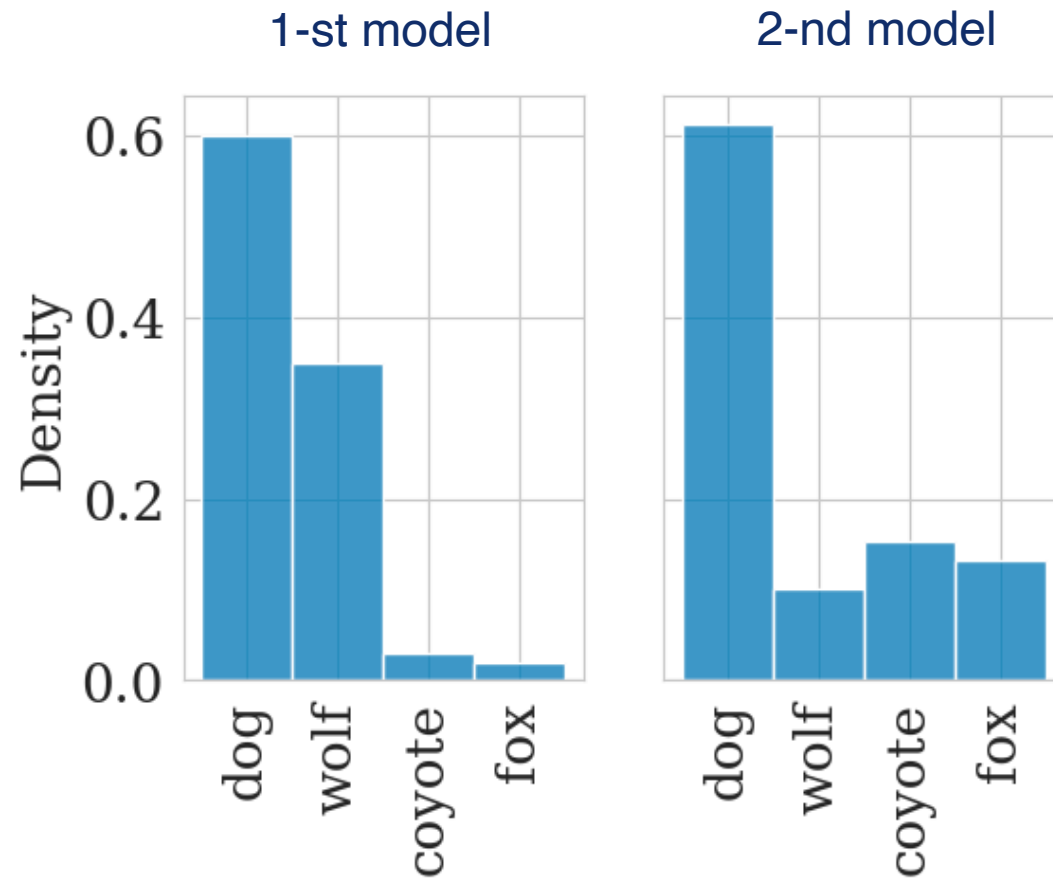
### Random ensemble

	I	II	III	IV	V
I	0.0	73.26	70.83	69.58	70.76
II	73.26	0.0	70.93	72.06	<b>74.49</b>
III	70.83	70.93	0.0	71.07	71.14
IV	69.58	72.06	71.07	0.0	71.57
V	70.76	<b>74.49</b>	71.14	71.57	0.0

## StarSSE—WO



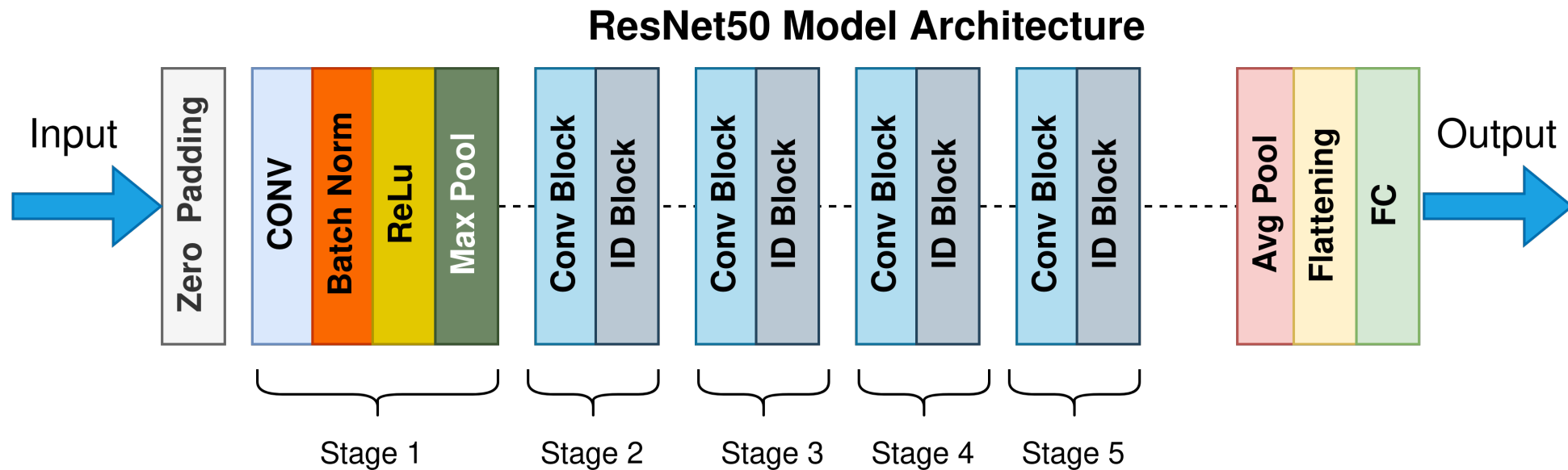
## StarSSE—CE



## The modifications of StarSSE results

Algorithm	avg. acc.	div.	ens. acc.
<b>SSE</b>	85.52 $\pm$ 0.34	71.57 $\pm$ 5.94	87.11 $\pm$ 0.06
<b>StarSSE</b>	85.88 $\pm$ 0.25	68.12 $\pm$ 1.15	<b>87.41<math>\pm</math>0.12</b>
<b>StarSSE–WO</b>	85.83 $\pm$ 0.24	67.31 $\pm$ 1.77	87.4 $\pm$ 0.16
<b>StarSSE–CE (2)</b>	85.4 $\pm$ 0.49	75.21 $\pm$ 3.79	<b>87.44<math>\pm</math>0.2</b>

## StarSSE with Parameter efficient Fine-tuning (PeFt)



## StarSSE with PeFt results

Algorithm	avg. acc.	div.	ens. acc.
<b>StarSSE</b>	<b>85.88<math>\pm</math>0.25</b>	<b>68.12<math>\pm</math>1.15</b>	<b>87.41<math>\pm</math>0.12</b>
StarSSE (1) x0.74 time	85.57 $\pm$ 0.356	71.37 $\pm$ 2.02	87.26 $\pm$ 0.06
StarSSE (2) x0.6 time	85.94 $\pm$ 0.225	60.94 $\pm$ 1.25	87.1 $\pm$ 0.14
StarSSE (3) x0.55 time	86.01 $\pm$ 0.195	17.38 $\pm$ 1.02	86.09 $\pm$ 0.22



## Conclusion

- Individual quality and diversity are the key
- Increasing StarSSE individual quality reduces ensemble performance
- Increasing StarSSE diversity corrupts individual models, but it is possible
- StarSSE is the best algorithm, gives the best diversity with insignificant individual quality decrease