

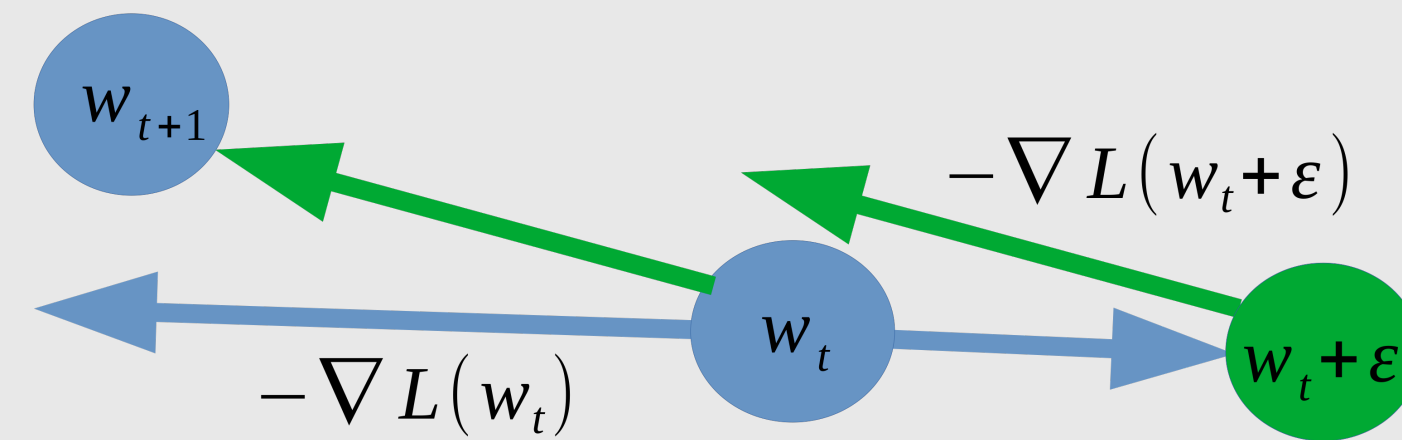
# Normalization Layers Are all that Sharpness-Aware Minimization Needs

## Sharpness-Aware Minimization

- Optimization method designed to implicitly maximize the flatness in weight space during gradient descent:

$$\min_{\mathbf{w}} \max_{\|\epsilon\|_p < \rho} \mathcal{L}(\mathbf{w} + \epsilon)$$

- In practice: approximate inner maximization with 1-step and use  $\nabla \mathcal{L}(\mathbf{w} + \epsilon)$  for batch-wise weight update instead of  $\nabla \mathcal{L}(\mathbf{w}) \rightarrow$  resulting SAM-algorithm requires additional forward-backward pass



- Adaptive variant (ASAM) has objective:

$$\min_{\mathbf{w}} \max_{\|T_w^{-1}\epsilon\|_p < \rho} \mathcal{L}(\mathbf{w} + \epsilon)$$

- $T_w$  is a normalization operator (diagonal matrix) making the perturbation (partly) invariant to rescaling of the parameters

- The 1-step solution of the inner maximization

$$\epsilon = \rho \frac{T_w^2 \nabla \mathcal{L}(\mathbf{w})}{\|T_w \nabla \mathcal{L}(\mathbf{w})\|_2} \text{ for } p = 2$$

$$\epsilon = \rho T_w \text{sign}(\nabla \mathcal{L}(\mathbf{w})) \text{ for } p = \infty$$

- Justifications for empirical success of SAM and its variants remain inconclusive

## SAM-ON

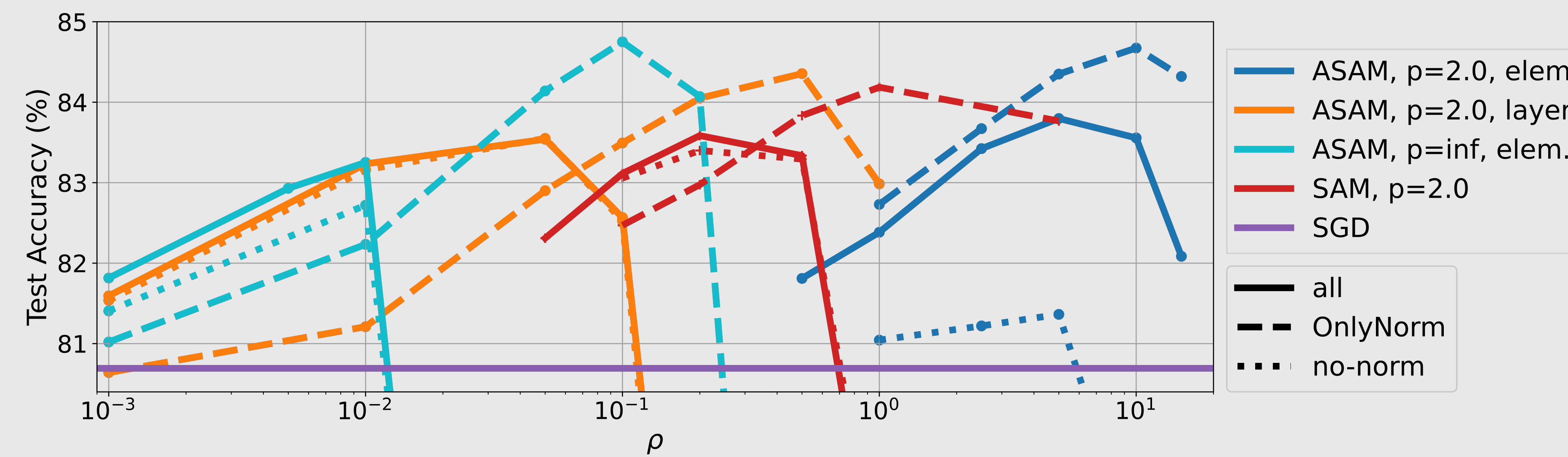
- Normalization layers normalize input  $x$  with mean  $\mu$  and variance  $\sigma^2$

$$\text{Norm}(x) = \gamma \times \frac{x - \mu}{\sigma} + \beta$$

- We propose SAM-ON, which perturbs **only**  $\gamma$  and  $\beta$  (typically 0.1% of all parameters) in the adversarial step of SAM

- For ResNets with BatchNorm and VisionTransformers with Layernorm we observe that SAM-ON either matches the performance of the conventional SAM algorithm (called SAM-all), or even outperforms it

## ResNets on CIFAR-100



- SAM-ON (dashed) improves generalization performance across a range of SAM-variants compared to SAM-all and vanilla SGD (shown is a WRN28-10, more models in paper!)
- omitting the normalization layers in the perturbation (no-norm, dotted) can harm training

## Central Message

Applying Sharpness-Aware Minimization **only to the normalization layers** of a network typically enhances its performance

## ViT-S/32 on ImageNet

- For AdamW as base optimizer, SAM-ON improves strongly over the vanilla optimizer and either improves over SAM-all or performs on par
- For Lion as base optimizer, SAM-ON always improves over SAM-all and the vanilla optimizer

ID	AdamW			Lion			
	vanilla	SAM-all	SAM-ON	vanilla	SAM-all	SAM-ON	
ImageNet	66.89±0.04	71.47±0.12	71.37±0.026	68.20±0.02	71.90±0.19	<b>72.64</b> ±0.14	
ImageNetV2	48.43±0.48	53.61±0.11	53.67±0.29	50.20±0.01	54.20±0.27	<b>55.38</b> ±0.09	
ImageNetR	25.04±0.04	31.56±0.48	<b>32.98</b> ±0.10	25.61±0.04	32.17±0.41	<b>33.87</b> ±0.47	
OOD ImageNetA	4.72±0.15	5.21±0.05	5.19±0.18	5.45±0.19	5.01±0.22	<b>5.77</b> ±0.21	
ImageNetSketch	13.68±0.24	18.50±0.44	<b>19.35</b> ±0.17	14.47±0.02	18.22±0.34	<b>20.48</b> ±0.12	
ObjectNet	11.32±0.39	13.75±0.12	13.55±0.25	12.06±0.02	13.93±0.40	<b>15.35</b> ±0.13	
adv. rob.	$l_2, \epsilon = 0.25$	19.67±0.47	37.53±0.69	<b>41.16</b> ±0.24	22.01±0.78	38.52±0.66	<b>43.12</b> ±0.97
	$l_2, \epsilon = 0.50$	5.47±0.18	17.71±0.61	<b>22.72</b> ±0.25	6.63±0.46	19.03±0.92	<b>24.27</b> ±1.34
	$l_\infty, \epsilon = 0.25/255$	33.45±0.80	48.08±0.14	<b>49.34</b> ±0.08	35.31±0.08	49.57±0.60	<b>51.37</b> ±0.99
	$l_\infty, \epsilon = 0.5/255$	14.98±0.18	29.68±0.09	<b>32.46</b> ±0.15	15.86±0.13	31.68±0.62	<b>34.23</b> ±1.73

## Sharpness

SAM-ON models are sharper, yet generalize better. Shown is logit-normalized  $l_\infty$   $m$ -sharpness, averaged over three models per method for a WRN-28 on CIFAR-100.

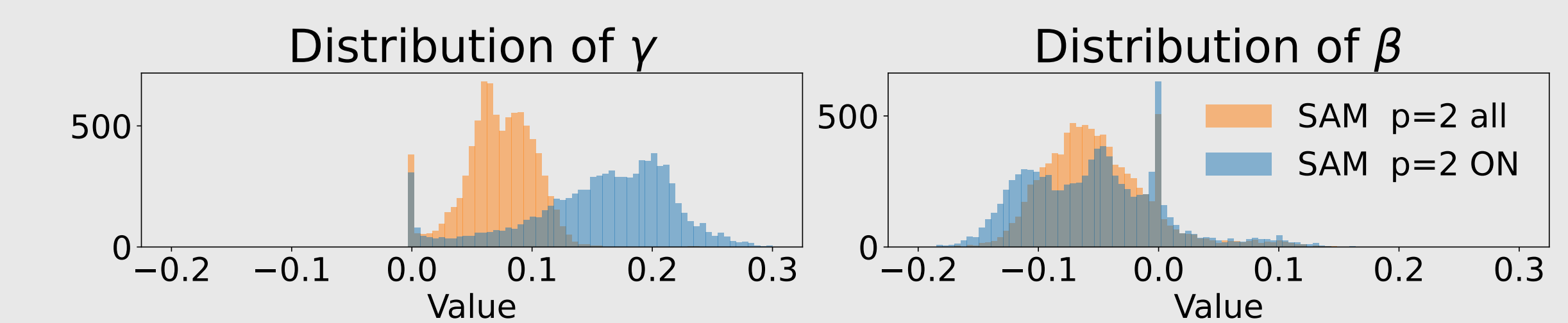
	SGD	SAM-all	SAM-ON	
Test Accuracy (%)	80.71±0.2	83.11±0.3	<b>84.19</b> ±0.2	
$l_\infty$ -sharp.	20 steps, $\rho = 0.003$	0.071±0.000	<b>0.048</b> ±0.001	0.090±0.005
	20 steps, $\rho = 0.007$	0.433±0.002	<b>0.309</b> ±0.011	0.585±0.018
	1 step, $\rho = 0.01$	0.204±0.005	<b>0.183</b> ±0.002	0.315±0.010
	1 step, $\rho = 0.03$	0.809±0.003	<b>0.769</b> ±0.017	0.843±0.007

## Other sparse perturbation approaches

Other sparse perturbation approaches are less effective than SAM-ON, especially when probed at very high sparsity levels.

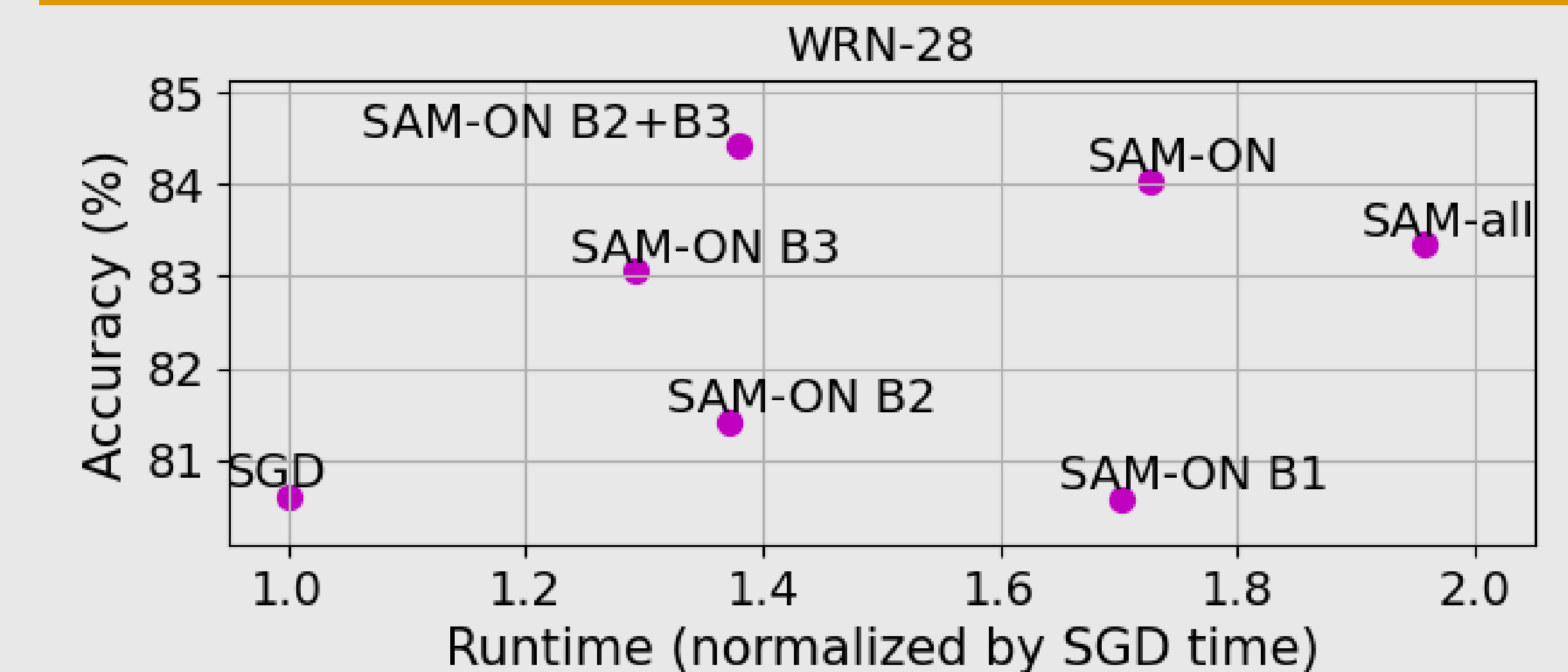
Sparsity	SAM	SAM-ON	Random Mask	SSAM-F	
0%		99.95%	99.95%	50%	99.95%
Test Accuracy (%)	83.11±0.3	<b>84.19</b> ±0.2	80.97±0.2	83.94±0.1	83.14±0.1

## SAM-ON induces a parameter shift



For SAM-ON the distribution of  $\gamma$  shifts towards larger values

## Computational savings



SAM-ON reduces the computational load in practice, and further gains might be achieved by perturbing the normalization layers of selected blocks (B1-B3)

## Paper & code

Find paper and code at [github.com/mueller-mp/SAM-ON](https://github.com/mueller-mp/SAM-ON).