



**ICLR**  
**2022**

# Surrogate Gap Guided Minimization Improves Sharpness-Aware Training

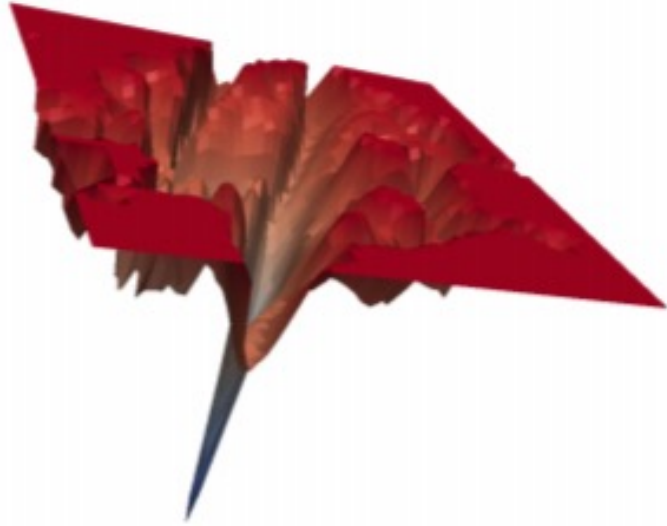
Juntang Zhuang<sup>1</sup>, Boqing Gong<sup>2</sup>, Liangzhe Yuan<sup>2</sup>, Yin Cui<sup>2</sup>, Hartwig Adam<sup>2</sup>,  
Nicha C. Dvornek<sup>1</sup>, Sekhar Tatikonda<sup>1</sup>, James S. Duncan<sup>1</sup>, Ting Liu<sup>2</sup>

1: Yale University

2: Google Research



# Generalization and curvature of loss surface



- Minimizing the training loss  $\rightarrow$  sharp local minima (left)  $\rightarrow$  poor generalization

[1] Chaudhari, Pratik, et al. "Entropy-sgd: Biasing gradient descent into wide valleys." *Journal of Statistical Mechanics: Theory and Experiment* 2019.12 (2019): 124018.

[2] Jiang, Yiding, et al. "Fantastic generalization measures and where to find them." *arXiv preprint arXiv:1912.02178* (2019).

[3] Niladri S Chatterji, et al. The intriguing role of module criticality in the generalization of deep networks. ICLR 2020



# Generalization and curvature of loss surface



- Minimizing the training loss  $\rightarrow$  sharp local minima (left)  $\rightarrow$  poor generalization
- Flatter local minima (right)  $\rightarrow$  better generalization [1,2,3]

[1] Chaudhari, Pratik, et al. "Entropy-sgd: Biasing gradient descent into wide valleys." *Journal of Statistical Mechanics: Theory and Experiment* 2019.12 (2019): 124018.

[2] Jiang, Yiding, et al. "Fantastic generalization measures and where to find them." *arXiv preprint arXiv:1912.02178* (2019).

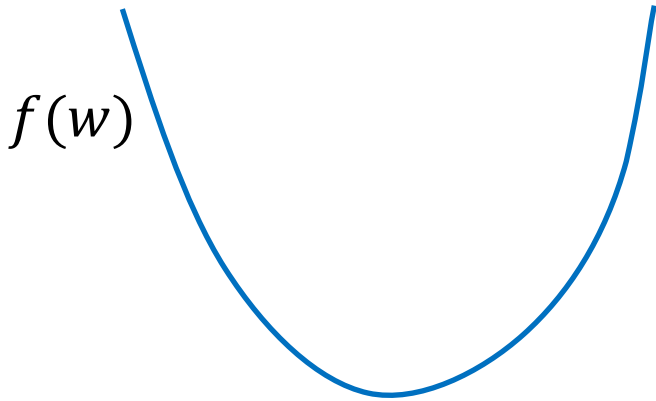
[3] Niladri S Chatterji, et al. The intriguing role of module criticality in the generalization of deep networks. ICLR 2020



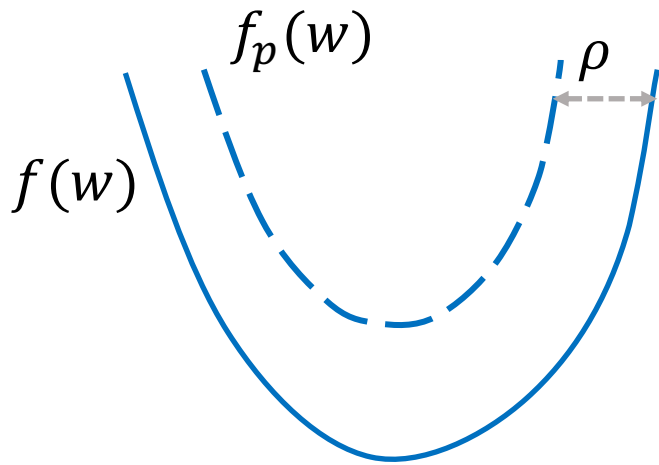
Surrogate Gap: an equivalent measure of curvature at local minima

# Surrogate Gap: an equivalent measure of curvature at local minima

$f(w)$ : loss function



# Surrogate Gap: an equivalent measure of curvature at local minima

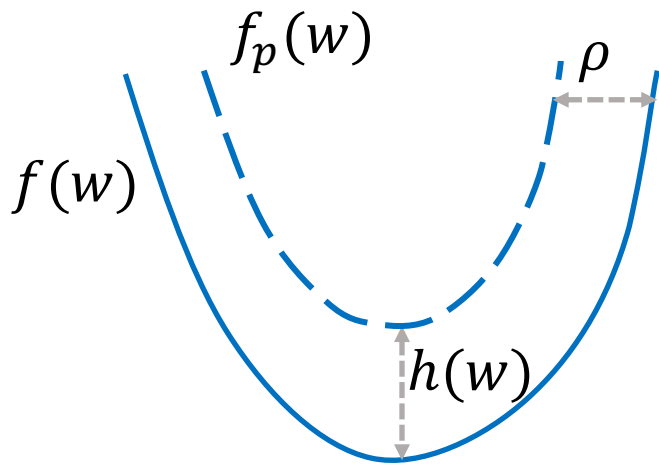


$f(w)$ : loss function

$f_p(w)$ : the maximum loss within the neighborhood (with a fixed radius  $\rho$ ) centered around  $w$

$$f_p(w) \approx f\left(w + \rho \frac{\nabla f(w)}{\|\nabla f(w)\|}\right)$$

# Surrogate Gap: an equivalent measure of curvature at local minima



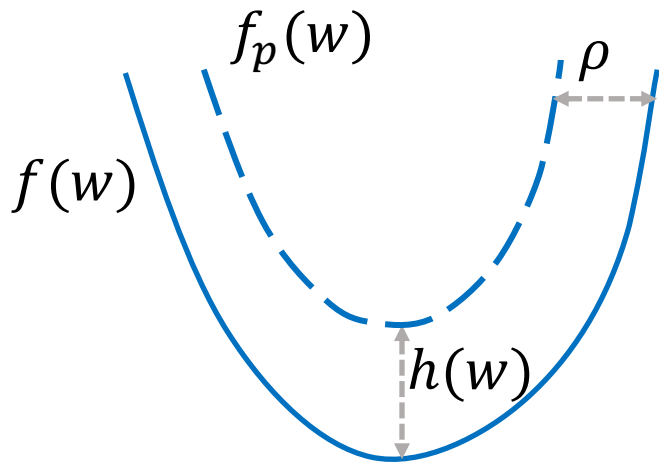
$f(w)$ : loss function

$f_p(w)$ : the maximum loss within the neighborhood (with a fixed radius  $\rho$ ) centered around  $w$

$$f_p(w) \approx f\left(w + \rho \frac{\nabla f(w)}{\|\nabla f(w)\|}\right)$$

$h(w)$ :  $f_p(w) - f(w)$  The “surrogate gap”

# Surrogate Gap: an equivalent measure of curvature at local minima



$f(w)$ : loss function

$f_p(w)$ : the maximum loss within the neighborhood (with a fixed radius  $\rho$ ) centered around  $w$

$$f_p(w) \approx f\left(w + \rho \frac{\nabla f(w)}{\|\nabla f(w)\|}\right)$$

$h(w)$ :  $f_p(w) - f(w)$  The “surrogate gap”

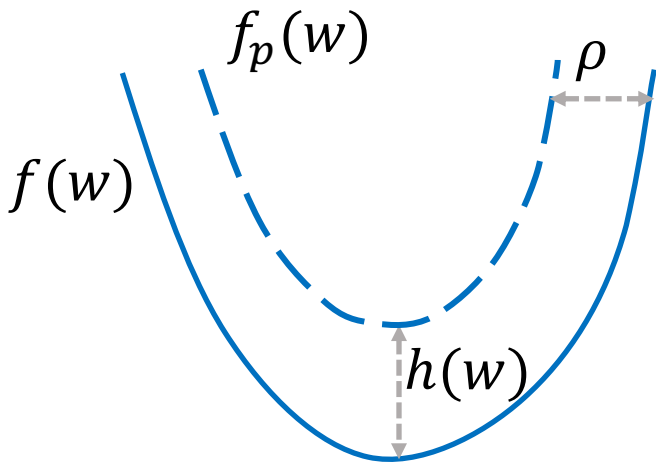
At a local minimum:

$$h(w) \approx \frac{1}{2} |\sigma_{max}| \rho^2 \propto |\sigma_{max}|$$

$\sigma_{max}$  : dominate eigenvalue of the Hessian



# Surrogate Gap: an equivalent measure of curvature at local minima



$f(w)$ : loss function

$f_p(w)$ : the maximum loss within the neighborhood (with a fixed radius  $\rho$ ) centered around  $w$

$$f_p(w) \approx f\left(w + \rho \frac{\nabla f(w)}{\|\nabla f(w)\|}\right)$$

$h(w)$ :  $f_p(w) - f(w)$  The “surrogate gap”

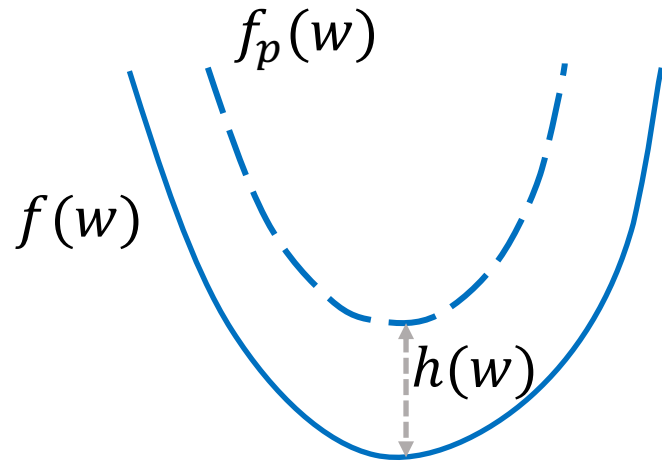
At a local minimum:

$$h(w) \approx \frac{1}{2} |\sigma_{max}| \rho^2 \propto |\sigma_{max}|$$

$\sigma_{max}$ : dominate eigenvalue of the Hessian

Larger  $h$ , sharper surface, worse generalization

# Sharpness-Aware Minimization (SAM)



Vanilla Training:

$$\min f(w)$$

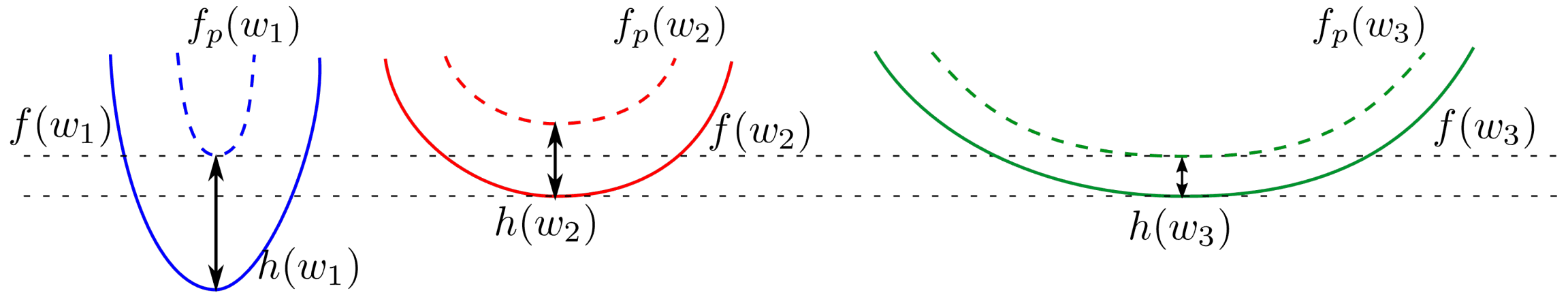
SAM:

$$\min f_p(w)$$

[1] Foret, Pierre, et al. "Sharpness-aware minimization for efficiently improving generalization." *arXiv preprint arXiv:2010.01412* (2020).



# Potential Caveat of Sharpness-Aware Minimization (SAM)



Left to right:

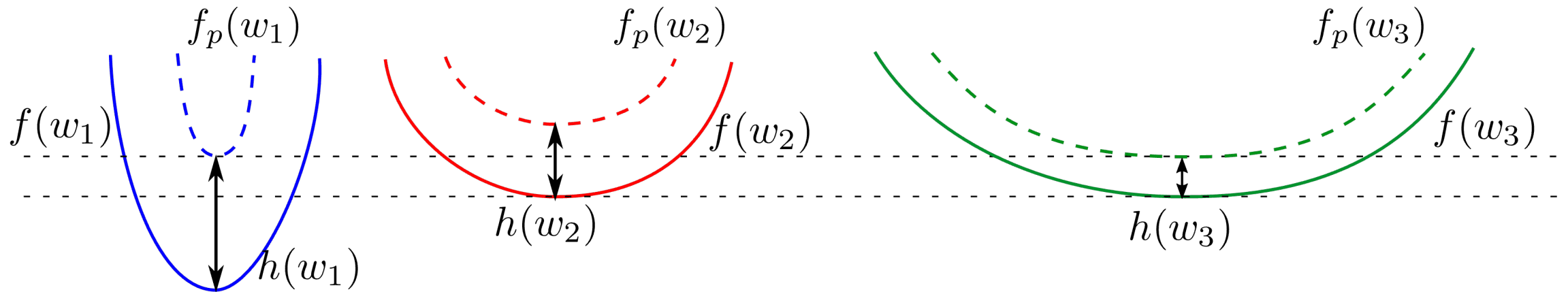
$$\text{Sharpness}(w_1) > \text{Sharpness}(w_2) > \text{Sharpness}(w_3)$$

$$h(w_1) > h(w_2) > h(w_3)$$

$$f_p(w_1) < f_p(w_2) > f_p(w_3)$$



# Potential Caveat of Sharpness-Aware Minimization (SAM)



Left to right:

$$\text{Sharpness}(w_1) > \text{Sharpness}(w_2) > \text{Sharpness}(w_3)$$

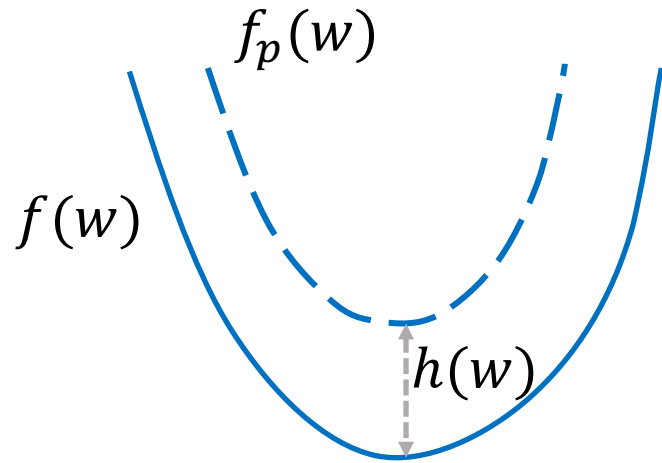
$$h(w_1) > h(w_2) > h(w_3)$$

$$f_p(w_1) < f_p(w_2) > f_p(w_3)$$

$f_p$  might disagree with sharpness, surrogate gap  $h$  agrees with sharpness  
(Lemma 3.1, 3.2, 3.3)



# Surrogate Gap Guided Sharpness-Aware Minimization (GSAM)



Vanilla Training:  $\min f(w)$

SAM:  $\min f_p(w)$

GSAM (ours):  $\min( f_p(w), h(w) )$

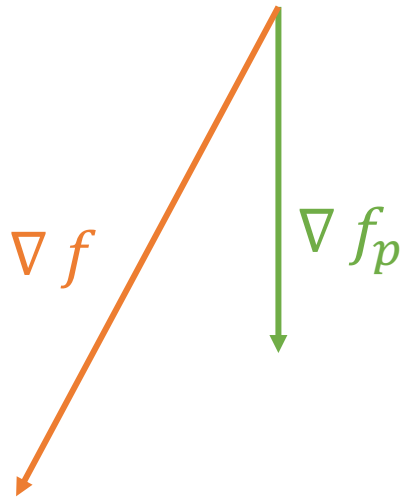
Minimize training  
loss

Minimize  
generalization gap



Simultaneously minimize two objectives that might **compete**

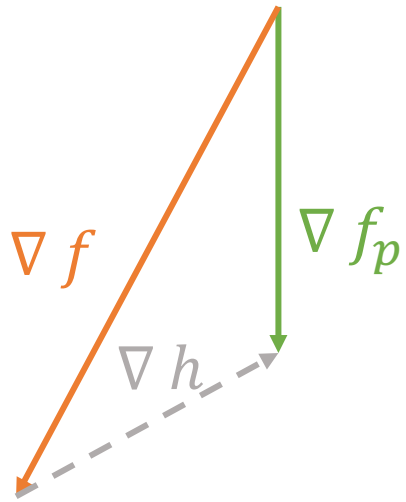
GSAM (ours):  $\min( f_p(w), h(w) )$



Simultaneously minimize two objectives that might **compete**

GSAM (ours):  $\min( f_p(w), h(w) )$

$$\nabla h = \nabla f_p - \nabla f$$

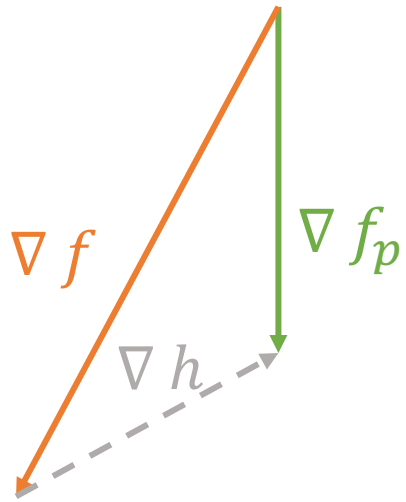


Simultaneously minimize two objectives that might **compete**

$$\text{GSAM (ours):} \quad \min( f_p(w), h(w) )$$

$$\nabla h = \nabla f_p - \nabla f$$

$\nabla h$  might have **negative** inner product with  $\nabla f_p$   
→ minimizing  $h$  would **increase**  $f_p$





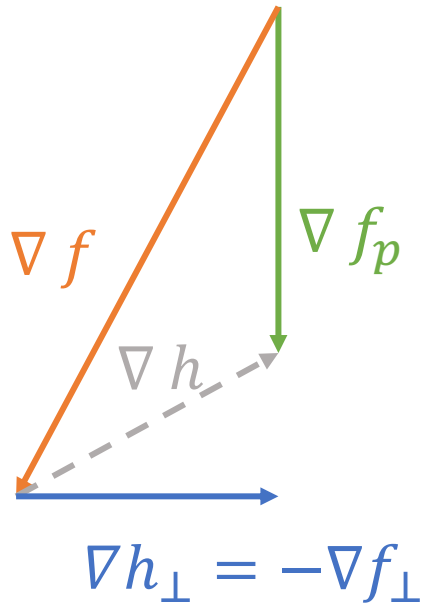
Simultaneously minimize two objectives that might **compete**

GSAM (ours):  $\min( f_p(w), h(w) )$

$$\nabla h = \nabla f_p - \nabla f$$

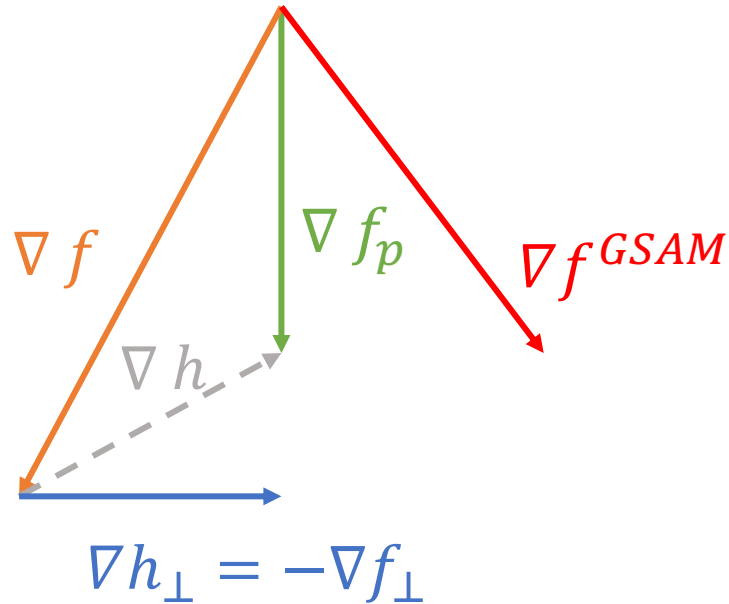
$\nabla h$  might have **negative** inner product with  $\nabla f_p$   
→ minimizing  $h$  would **increase**  $f_p$

Decompose  $\nabla h$  onto parallel and vertical to  $\nabla f_p$ ,  
update in  $\nabla h_{\perp}$  **not affects**  $f_p$ .



Simultaneously minimize two objectives that might **compete**

GSAM (ours):  $\min( f_p(w), h(w) )$



$$\nabla h = \nabla f_p - \nabla f$$

$\nabla h$  might have **negative** inner product with  $\nabla f_p$   
→ minimizing  $h$  would **increase**  $f_p$

Decompose  $\nabla h$  onto parallel and vertical to  $\nabla f_p$ ,  
update in  $\nabla h_{\perp}$  **not affects**  $f_p$ .

$$\nabla f^{GSAM} = \nabla f_p + \alpha \nabla h_{\perp} = \nabla f_p - \alpha \nabla f_{\perp}$$

min  $f_p$

min  $h$ , not  
affect  $f_p$



- SGD
- SAM
- GSAM

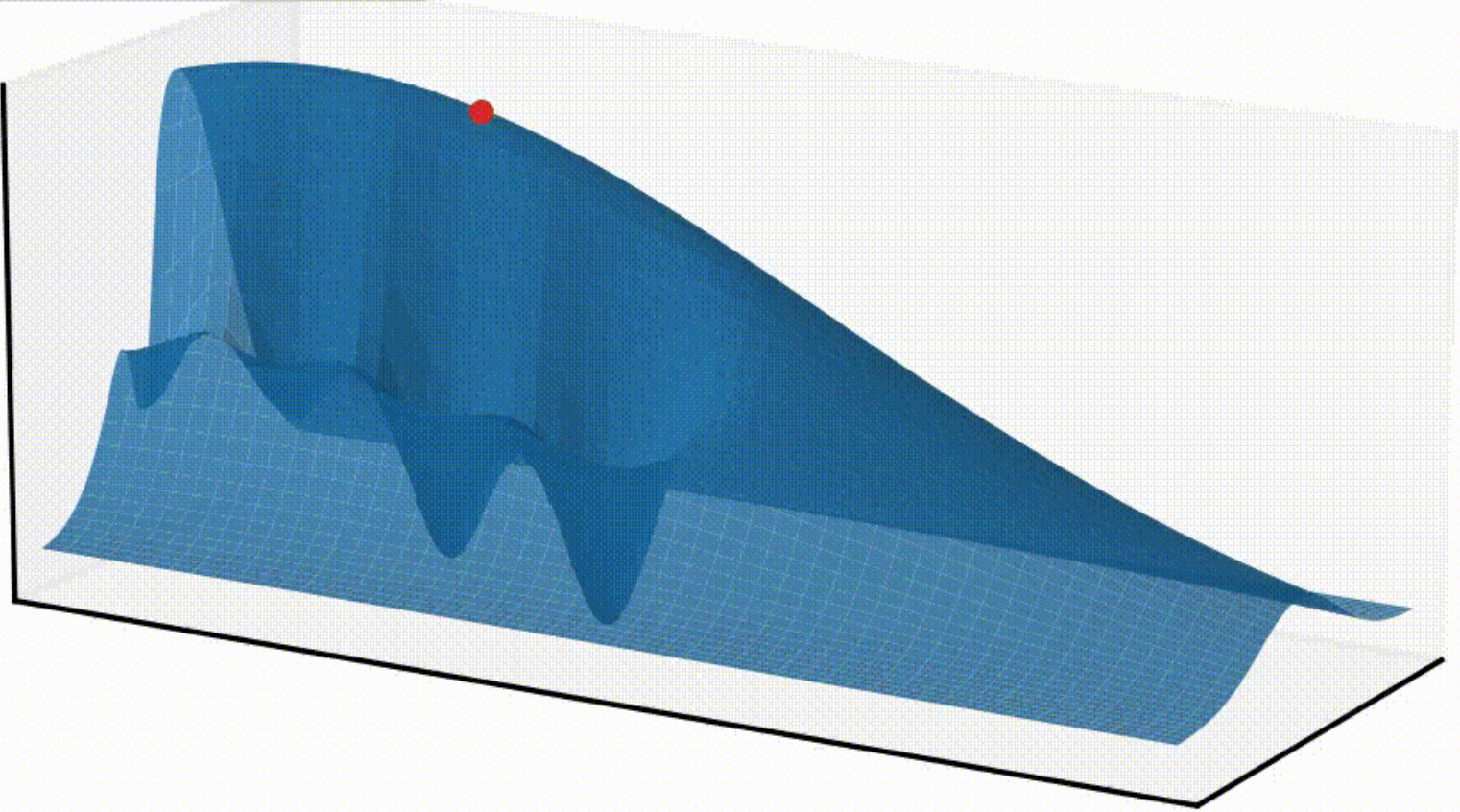


Table 1: Top-1 Accuracy (%) on ImageNet datasets for ResNets, ViTs and MLP-Mixers trained with Vanilla SGD or AdamW, SAM, and GSAM optimizers.

| Model     | Training      | ImageNet-v1 | ImageNet-Real | ImageNet-V2 | ImageNet-R  | ImageNet-C  |
|-----------|---------------|-------------|---------------|-------------|-------------|-------------|
| ResNet    |               |             |               |             |             |             |
| ResNet50  | Vanilla (SGD) | 76.0        | 82.4          | 63.6        | 22.2        | 44.6        |
|           | SAM           | 76.9        | 83.3          | 64.4        | <b>23.8</b> | 46.5        |
|           | <b>GSAM</b>   | <b>77.2</b> | <b>83.9</b>   | <b>64.6</b> | 23.6        | <b>47.6</b> |
| ResNet101 | Vanilla (SGD) | 77.8        | 83.9          | 65.3        | 24.4        | 48.5        |
|           | SAM           | 78.6        | 84.8          | 66.7        | 25.9        | 51.3        |
|           | <b>GSAM</b>   | <b>78.9</b> | <b>85.2</b>   | <b>67.3</b> | <b>26.3</b> | <b>51.8</b> |
| ResNet152 | Vanilla (SGD) | 78.5        | 84.2          | 66.3        | 25.3        | 50.0        |
|           | SAM           | 79.3        | 84.9          | 67.3        | 25.7        | 52.2        |
|           | <b>GSAM</b>   | <b>80.0</b> | <b>85.9</b>   | <b>68.6</b> | <b>27.3</b> | <b>54.1</b> |

| Vision Transformer |                 |             |             |             |             |             |
|--------------------|-----------------|-------------|-------------|-------------|-------------|-------------|
| ViT-S/32           | Vanilla (AdamW) | 68.4        | 75.2        | 54.3        | 19.0        | 43.3        |
|                    | SAM             | 70.5        | 77.5        | 56.9        | 21.4        | 46.2        |
|                    | <b>GSAM</b>     | <b>73.8</b> | <b>80.4</b> | <b>60.4</b> | <b>22.5</b> | <b>48.2</b> |
| ViT-S/16           | Vanilla (AdamW) | 74.4        | 80.4        | 61.7        | 20.0        | 46.5        |
|                    | SAM             | 78.1        | 84.1        | 65.6        | 24.7        | 53.0        |
|                    | <b>GSAM</b>     | <b>79.5</b> | <b>85.3</b> | <b>67.3</b> | <b>25.3</b> | <b>53.3</b> |
| ViT-B/32           | Vanilla (AdamW) | 71.4        | 77.5        | 57.5        | 23.4        | 44.0        |
|                    | SAM             | 73.6        | 80.3        | 60.0        | 24.0        | 50.7        |
|                    | <b>GSAM</b>     | <b>76.8</b> | <b>82.7</b> | <b>63.0</b> | <b>25.1</b> | <b>51.7</b> |
| ViT-B/16           | Vanilla (AdamW) | 74.6        | 79.8        | 61.3        | 20.1        | 46.6        |
|                    | SAM             | 79.9        | 85.2        | 67.5        | 26.4        | <b>56.5</b> |
|                    | <b>GSAM</b>     | <b>81.0</b> | <b>86.5</b> | <b>69.2</b> | <b>27.1</b> | 55.7        |



| Vision Transformer |                 |             |             |             |             |             |
|--------------------|-----------------|-------------|-------------|-------------|-------------|-------------|
| ViT-S/32           | Vanilla (AdamW) | 68.4        | 75.2        | 54.3        | 19.0        | 43.3        |
|                    | SAM             | 70.5        | 77.5        | 56.9        | 21.4        | 46.2        |
|                    | <b>GSAM</b>     | <b>73.8</b> | <b>80.4</b> | <b>60.4</b> | <b>22.5</b> | <b>48.2</b> |
| ViT-S/16           | Vanilla (AdamW) | 74.4        | 80.4        | 61.7        | 20.0        | 46.5        |
|                    | SAM             | 78.1        | 84.1        | 65.6        | 24.7        | 53.0        |
|                    | <b>GSAM</b>     | <b>79.5</b> | <b>85.3</b> | <b>67.3</b> | <b>25.3</b> | <b>53.3</b> |
| ViT-B/32           | Vanilla (AdamW) | 71.4        | 77.5        | 57.5        | 23.4        | 44.0        |
|                    | SAM             | 73.6        | 80.3        | 60.0        | 24.0        | 50.7        |
|                    | <b>GSAM</b>     | <b>76.8</b> | <b>82.7</b> | <b>63.0</b> | <b>25.1</b> | <b>51.7</b> |
| ViT-B/16           | Vanilla (AdamW) | 74.6        | 79.8        | 61.3        | 20.1        | 46.6        |
|                    | SAM             | 79.9        | 85.2        | 67.5        | 26.4        | <b>56.5</b> |
|                    | <b>GSAM</b>     | <b>81.0</b> | <b>86.5</b> | <b>69.2</b> | <b>27.1</b> | 55.7        |

| MLP-Mixer  |                 |             |             |             |             |             |
|------------|-----------------|-------------|-------------|-------------|-------------|-------------|
| Mixer-S/32 | Vanilla (AdamW) | 63.9        | 70.3        | 49.5        | 16.9        | 35.2        |
|            | SAM             | 66.7        | 73.8        | 52.4        | 18.6        | 39.3        |
|            | <b>GSAM</b>     | <b>68.6</b> | <b>75.8</b> | <b>55.0</b> | <b>22.6</b> | <b>44.6</b> |
| Mixer-S/16 | Vanilla (AdamW) | 68.8        | 75.1        | 54.8        | 15.9        | 35.6        |
|            | SAM             | 72.9        | 79.8        | 58.9        | 20.1        | 42.0        |
|            | <b>GSAM</b>     | <b>75.0</b> | <b>81.7</b> | <b>61.9</b> | <b>23.7</b> | <b>48.5</b> |
| Mixer-S/8  | Vanilla (AdamW) | 70.2        | 76.2        | 56.1        | 15.4        | 34.6        |
|            | SAM             | 75.9        | 82.5        | 62.3        | 20.5        | 42.4        |
|            | <b>GSAM</b>     | <b>76.8</b> | <b>83.4</b> | <b>64.0</b> | <b>24.6</b> | <b>47.8</b> |
| Mixer-B/32 | Vanilla (AdamW) | 62.5        | 68.1        | 47.6        | 14.6        | 33.8        |
|            | SAM             | 72.4        | 79.0        | 58.0        | 22.8        | 46.2        |
|            | <b>GSAM</b>     | <b>73.6</b> | <b>80.2</b> | <b>59.9</b> | <b>27.9</b> | <b>52.1</b> |
| Mixer-B/16 | Vanilla (AdamW) | 66.4        | 72.1        | 50.8        | 14.5        | 33.8        |
|            | SAM             | 77.4        | 83.5        | 63.9        | 24.7        | 48.8        |
|            | <b>GSAM</b>     | <b>77.8</b> | <b>84.0</b> | <b>64.9</b> | <b>28.3</b> | <b>54.4</b> |



## Conclusions:

- Surrogate Gap is an equivalent measure of sharpness
- Minimize both  $f_p$  (perturbed training loss) and  $h$  (sharpness)
- Gradient decomposition to avoid conflicts of multi-objective optimization
- Code available on project website (<https://sites.google.com/view/gsam-iclr22>)

