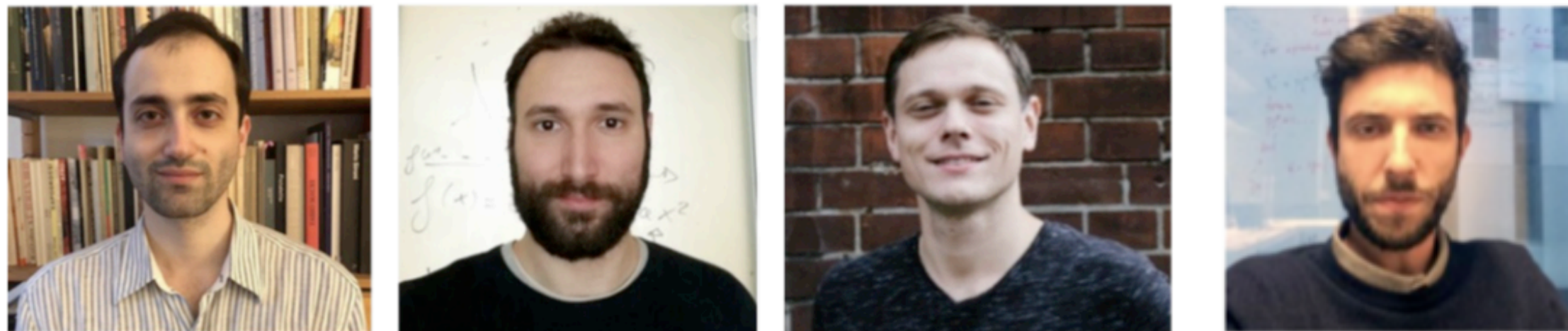


# Entropic Gradient Descent Algorithms and Wide Flat Minima

**Fabrizio Pittorino, Carlo Lucibello, Christoph Feinauer, Gabriele Perugini,  
Carlo Baldassi, Elizaveta Demyanenko, Riccardo Zecchina**



Artificial Intelligence Lab

Bocconi

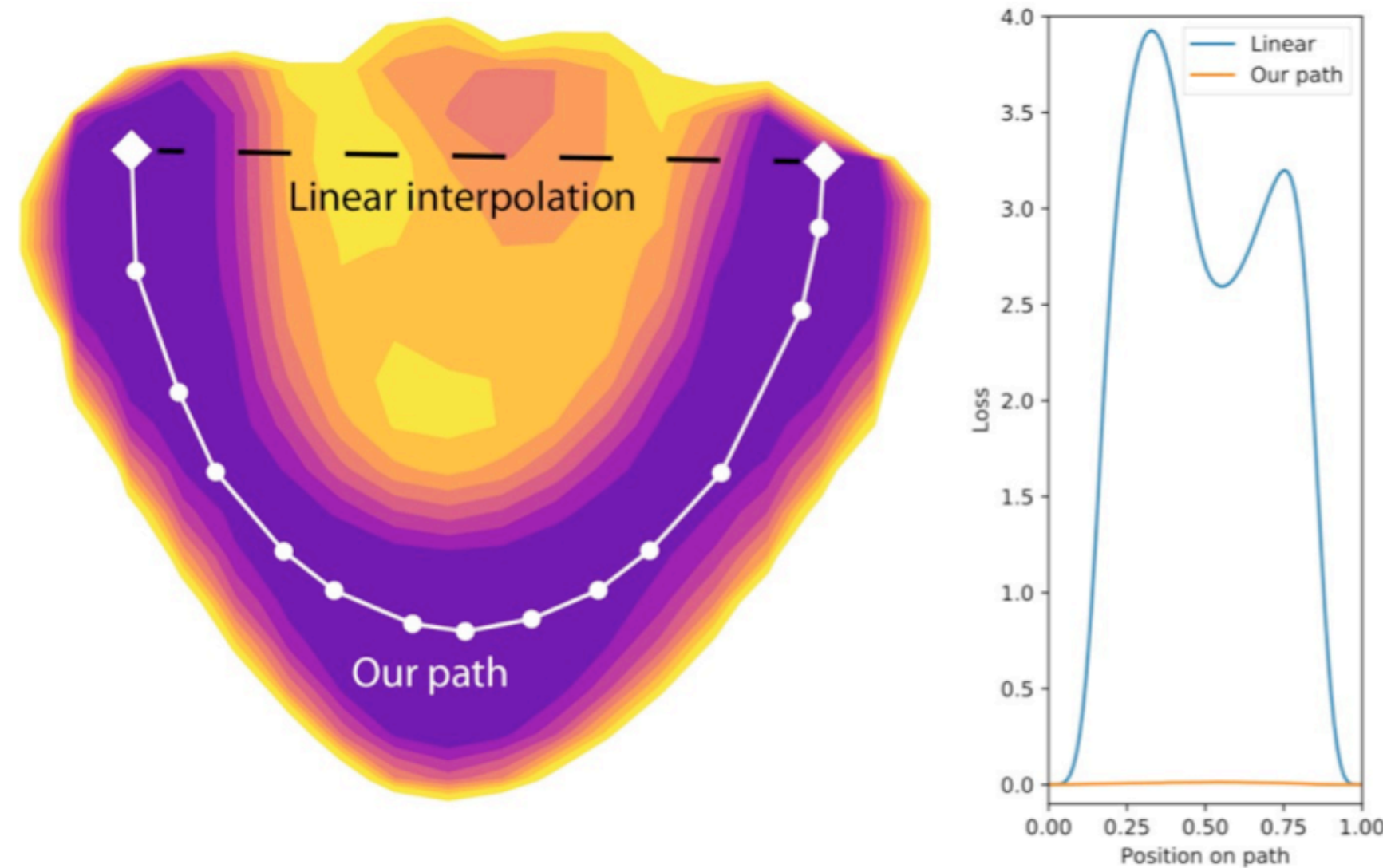




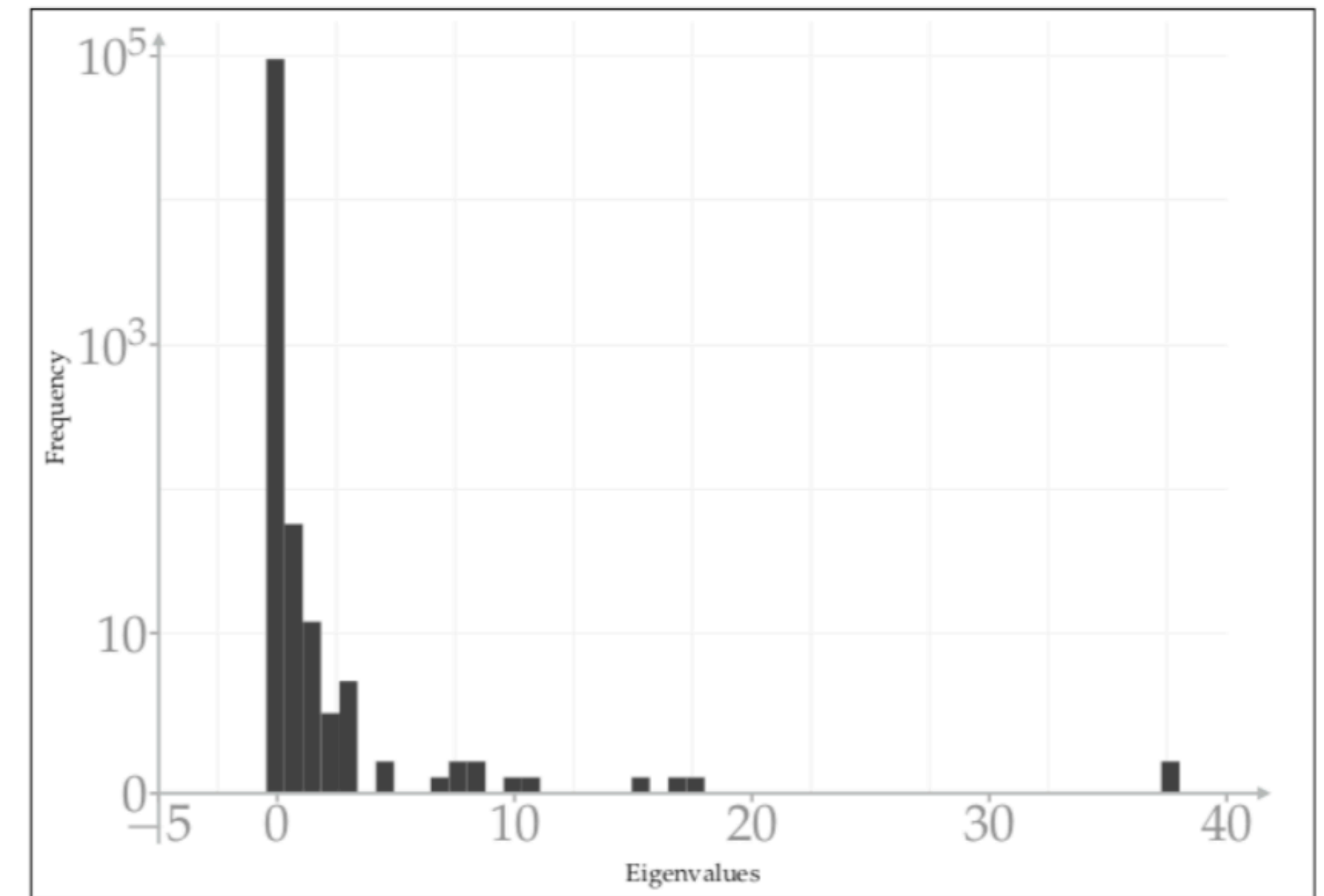
# Complex Loss Landscapes in Neural Networks

- There seems to be a flat non-convex “bottom” connecting “accessible” minimizers.

- The spectra of the Hessian in a quasi-minimum. Many flat directions.

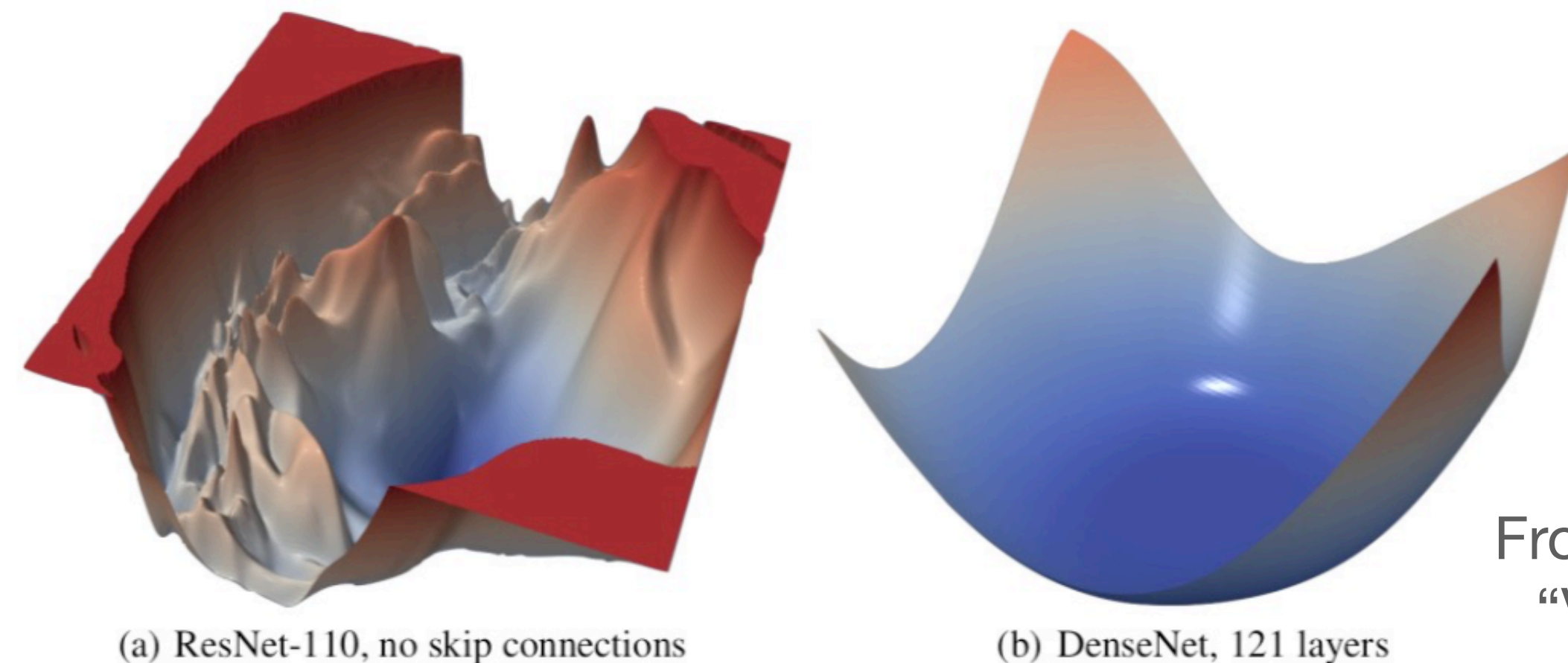


From Draxler et al. ‘18, “Essentially no barriers...”



From Chaudhary et al ‘17 “EntropySGD...”

- Architectural choices (e.g. loss, activations, batch-norm, skip-connections) influence the roughness and the large-scale structure of the landscape
- SGD batch size anti-correlates with minima width and with generalization



From Li et al. ‘17, “Visualizing...”

# Local Entropy and Local Energy

- How to tell apart good minima from bad minima?
- We conjecture that some geometrical properties of the training loss landscape, and in particular the **flatness** of minima, correlates well with generalization
- We define the **local entropy loss** as a way to characterize flatness, and as an auxiliary loss:

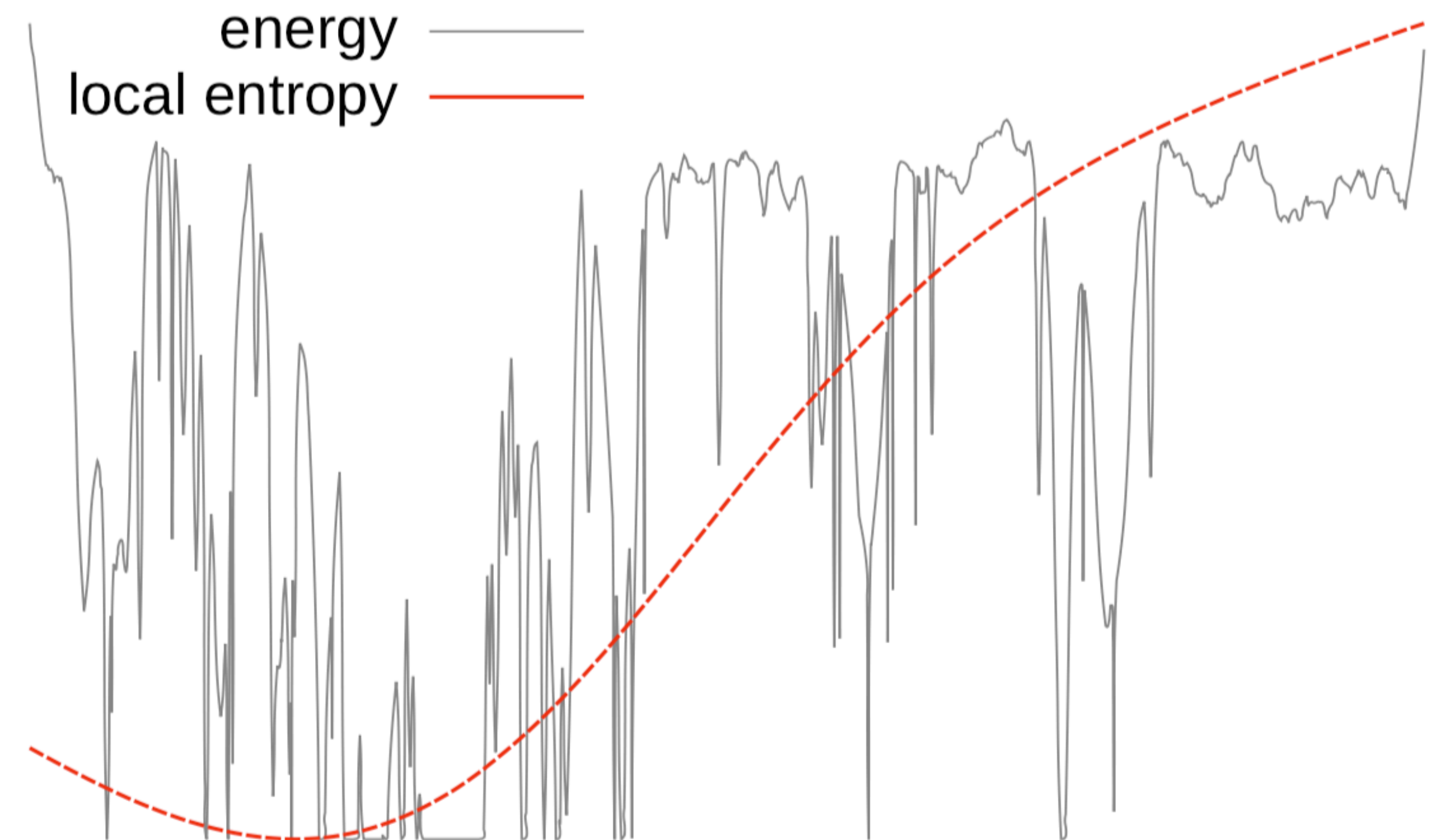
$$\mathcal{L}_{\text{LE}}(w) = -\frac{1}{\beta} \log \int dw' e^{-\beta \mathcal{L}(w') - \beta \gamma d(w', w)}$$

Where the squared distance is:  $d(w', w) = \frac{1}{2} \sum_{i=1}^N (w'_i - w_i)^2$

- The **local energy** is a simple flatness measure:

$$\delta E_{\text{train}}(w, \sigma) = \mathbb{E}_z E_{\text{train}}(w + \sigma z \odot w) - E_{\text{train}}(w)$$

Where the noise is:  $z \sim \mathcal{N}(0, I_N)$





# Entropic Algorithms

- Local Entropy hard to compute,
- but the gradient:  $\nabla \mathcal{L}_{\text{LE}}(w) = \gamma (w - \langle w' \rangle)$
- can be approximated by Stochastic Gradient Langevin Dynamics. The corresponding algorithm is called Entropy-SGD [1]

---

## Algorithm 1: Entropy-SGD (eSGD)

---

**Input**  $: w$   
**Hyper-parameters**  $: L, \eta, \gamma, \eta', \epsilon, \alpha$

```

1 for  $t = 1, 2, \dots$  do
2    $w', \mu \leftarrow w$ 
3   for  $l = 1, \dots, L$  do
4      $\Xi \leftarrow \text{sample minibatch}$ 
5      $dw' \leftarrow \nabla \mathcal{L}(w'; \Xi) + \gamma (w' - w)$ 
6      $w' \leftarrow w' - \eta' dw' + \sqrt{\eta'} \epsilon \mathcal{N}(0, I)$ 
7      $\mu \leftarrow \alpha \mu + (1 - \alpha) w'$ 
8    $w \leftarrow w - \eta (w - \mu)$ 

```

---



---

## Algorithm 2: Replicated-SGD (rSGD)

---

**Input**  $: \{w^a\}$   
**Hyper-parameters**  $: y, \eta, \gamma, K$

```

1 for  $t = 1, 2, \dots$  do
2    $\bar{w} \leftarrow \frac{1}{y} \sum_{a=1}^y w^a$ 
3   for  $a = 1, \dots, y$  do
4      $\Xi \leftarrow \text{sample minibatch}$ 
5      $dw^a \leftarrow \nabla \mathcal{L}(w^a; \Xi)$ 
6     if  $t = 0 \bmod K$  then
7        $dw^a \leftarrow dw^a + K \gamma (w^a - \bar{w})$ 
8      $w^a \leftarrow w^a - \eta dw^a$ 

```

---

- Another class of entropic algorithms can be derived starting from

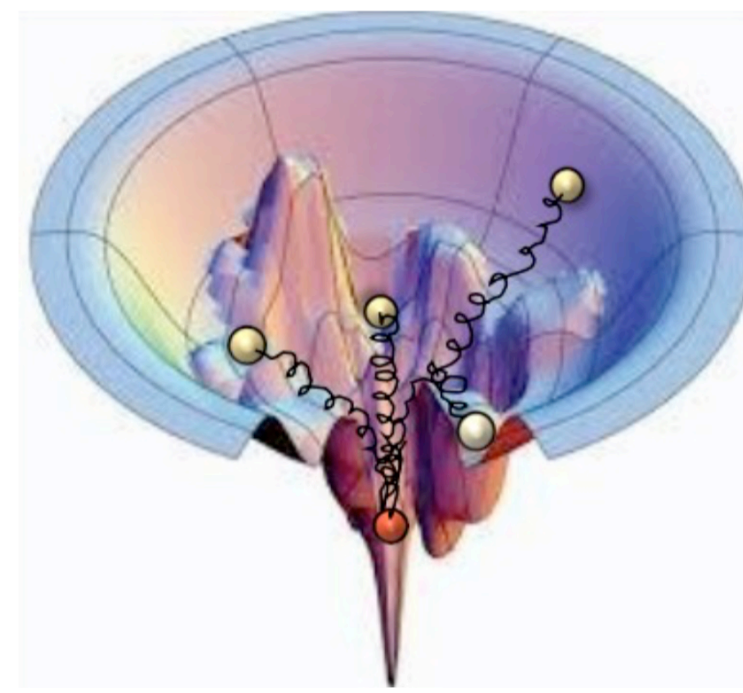
$$p(w) \propto e^{-\beta y \mathcal{L}_{\text{LE}}(w)}$$

- For  $y$  integer, one can use the local entropy definition to obtain the statistical measure of a system with  $y+1$  replicas, then integrate out the original one and obtain:

$$p(\{w^a\}_{a=1}^y) \propto e^{-\beta \mathcal{L}_R(\{w^a\})}$$

- Where  $\mathcal{L}_R(\{w^a\}_a) = \sum_{a=1}^y \mathcal{L}(w^a) + \gamma \sum_{a=1}^y d(w^a, \bar{w})$

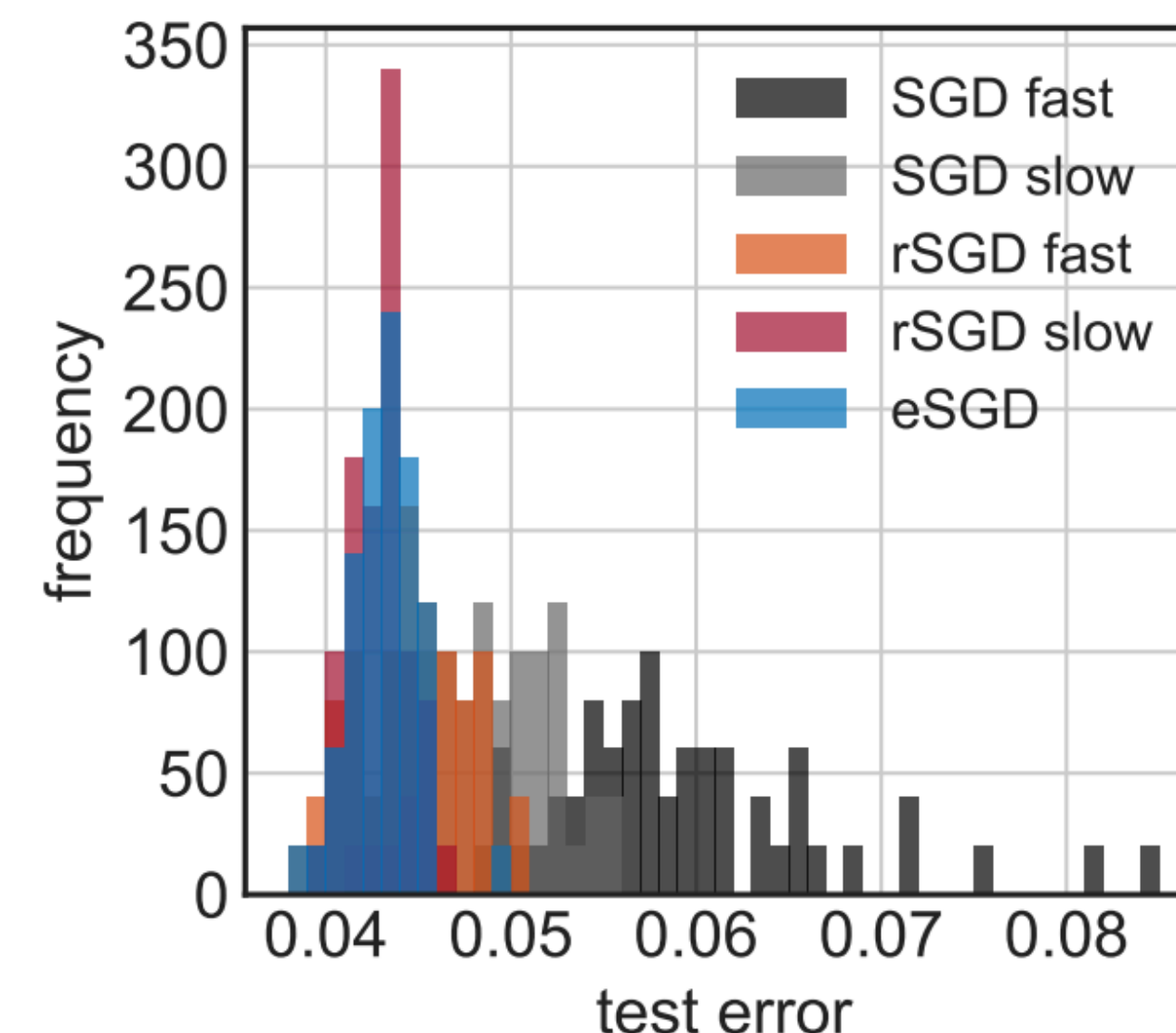
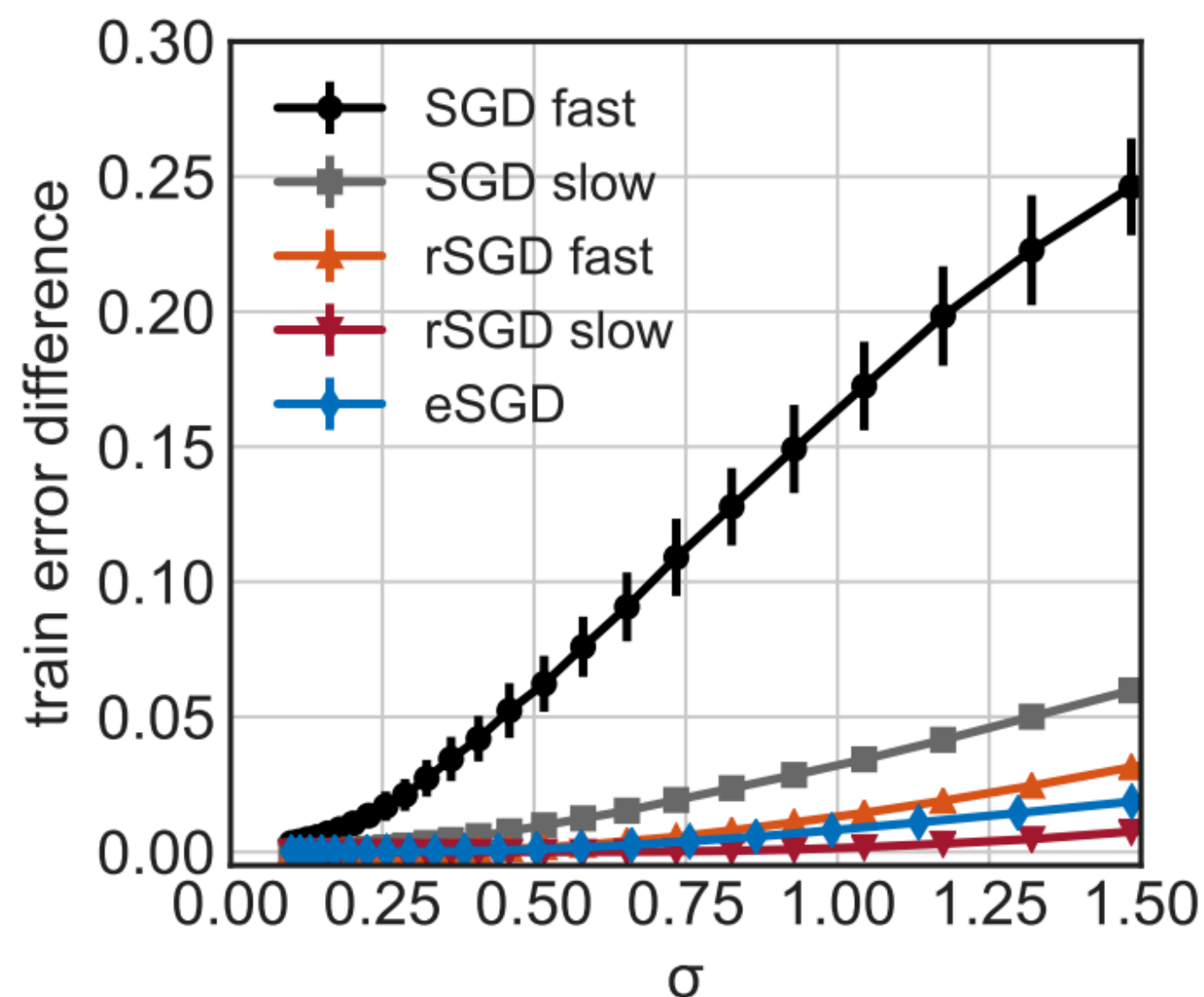
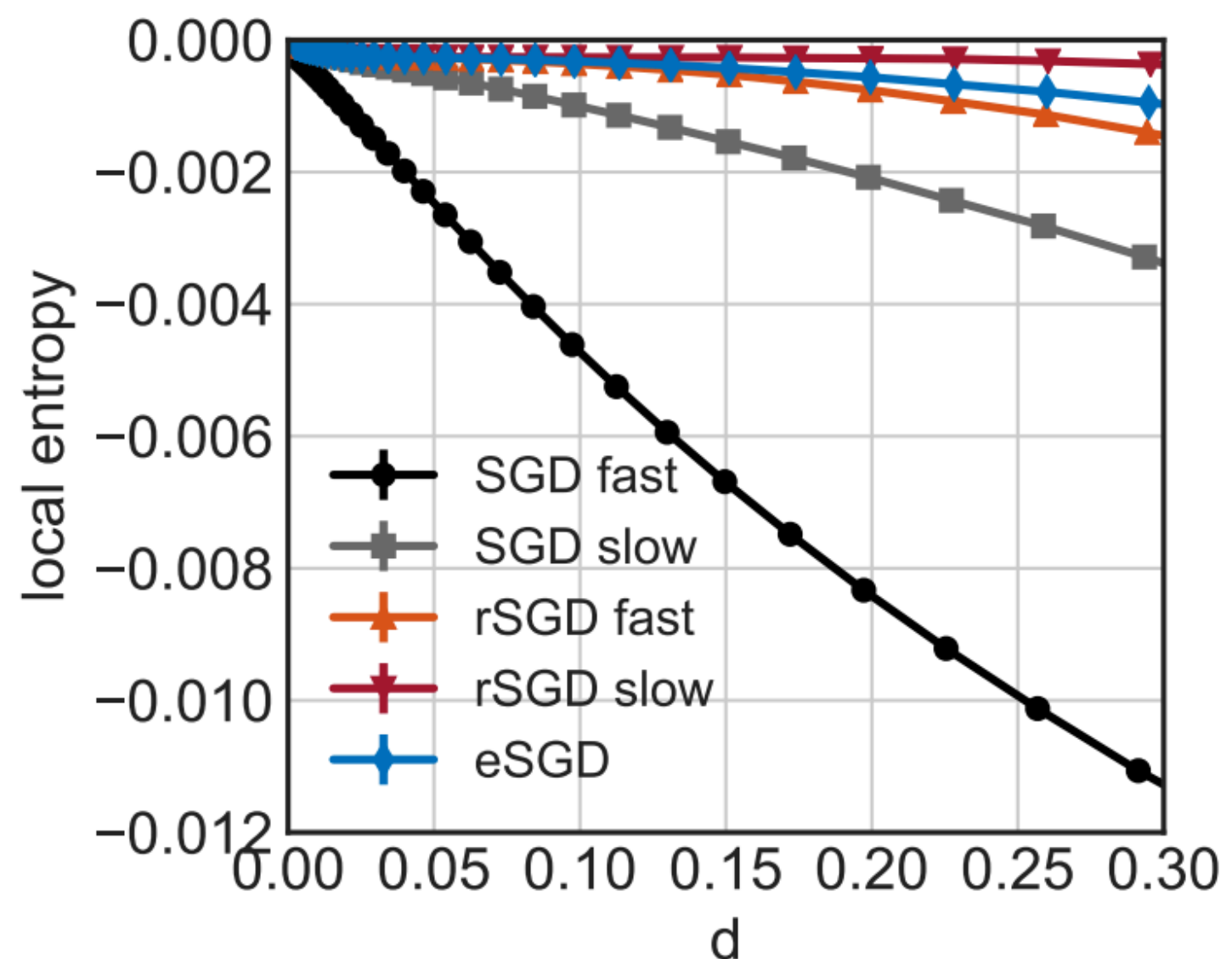
- with  $\bar{w} = \frac{1}{y} \sum_a w^a$ . Now one can perform SGD on the replicated loss.



# Shallow networks: estimation by Belief Propagation

- Shallow network performing a binary classification task on Fashion-MNIST
- We use Belief Propagation to estimate local entropy
- Local entropy and local energy correlate with each other and with generalisation

$$\hat{\sigma}(w, x) = \text{sign} \left[ \frac{1}{\sqrt{K}} \sum_{k=1}^K \text{sign} \left( \frac{1}{\sqrt{N}} \sum_{i=1}^N w_{ki} x_i \right) \right]$$





# Deep Networks: flatness and generalisation

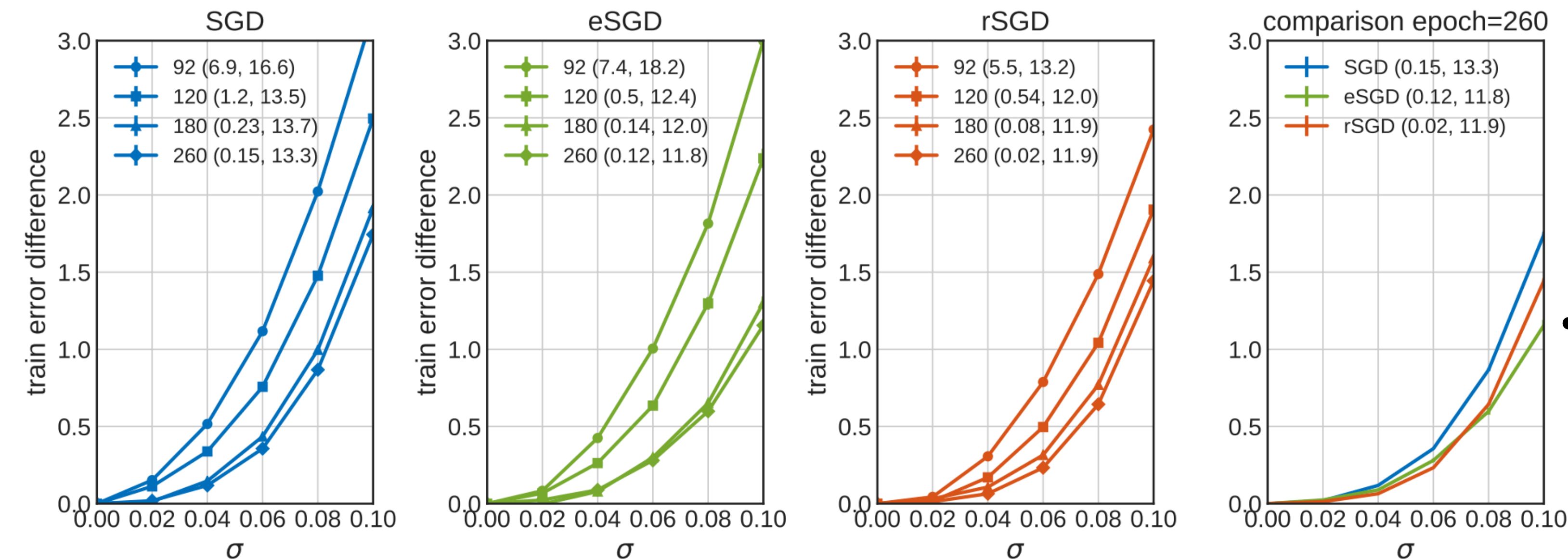
Dataset	Model	Baseline	rSGD	eSGD	rSGD $\times y$
<b>CIFAR-10</b>	SmallConvNet	$16.5 \pm 0.2$	$15.6 \pm 0.3$	$14.7 \pm 0.3$	$14.9 \pm 0.2$
	ResNet-18	$13.1 \pm 0.3$	$12.4 \pm 0.3$	$12.1 \pm 0.3$	$11.8 \pm 0.1$
	ResNet-110	$6.4 \pm 0.1$	$6.2 \pm 0.2$	$6.2 \pm 0.1$	$5.3 \pm 0.1$
	PyramidNet+ShakeDrop	$2.1 \pm 0.2$	$2.2 \pm 0.1$		1.8
<b>CIFAR-100</b>	PyramidNet+ShakeDrop	$13.8 \pm 0.1$	$13.5 \pm 0.1$		12.7
	EfficientNet-B0	20.5	20.6	$20.1 \pm 0.2$	19.5
<b>Tiny ImageNet</b>	ResNet-50	$45.2 \pm 1.2$	$41.5 \pm 0.3$	$41.7 \pm 1$	$39.2 \pm 0.3$
	DenseNet-121	$41.4 \pm 0.3$	$39.8 \pm 0.2$	$38.6 \pm 0.4$	$38.9 \pm 0.3$

- We want to verify that our entropic algorithm effectively find flatter minima.

- Local Entropy is expensive to compute, we compute the cheap Local Energy:

$$\delta\epsilon(w) = \mathbb{E}_z \epsilon(w(1 + \sigma z)) - \epsilon(w)$$

- Confirming that entropic algorithm find flatter minima and generalize better



# Conclusions

- Local entropy and local energy correlate with each other and with generalisation
- Detailed comparison on shallow networks (semi-analytical study)
- For deep networks, we showed entropic algorithms outperform standard ones, having enhanced generalisation and flatness