# Rethinking Common Practices in Deep Learning

## Elad Hoffer



**TECHNION**
Israel Institute
of Technology

# Deep Learning practices

- Deep learning is evolving fast, while many fundamental questions remain unanswered
    - Many heuristics and intuition-guided decisions
    - It works!
    - But some may prove misguided

We'll cover 3 of these
1. The impact of batch-size on generalization
2. Early-stopping and determining "over-fitting"
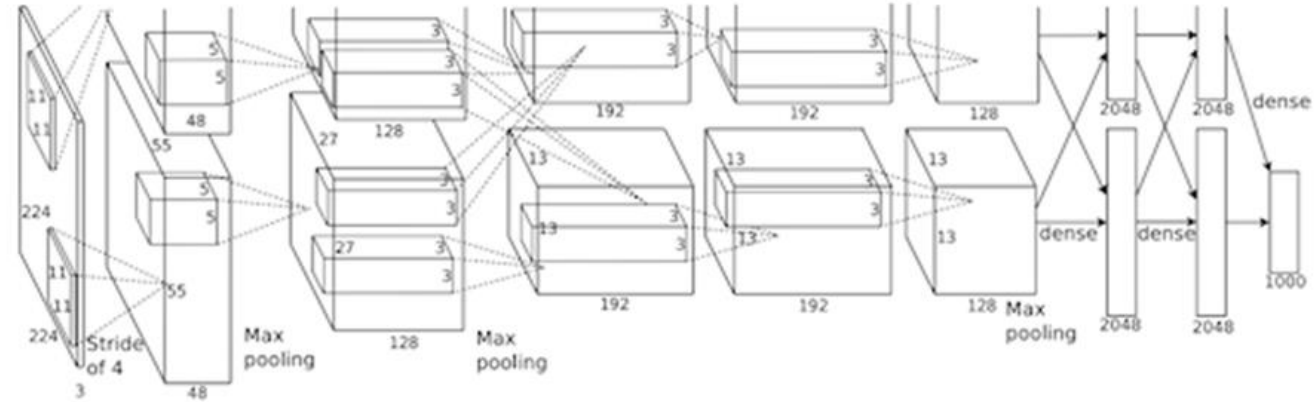3. The role of the last classification layer

# 1) Closing the "generalization gap"

**"Train longer, generalize better: closing the generalization gap in large batch training of neural networks" (NIPS 2017 oral)**

Elad Hoffer*, Itay Hubara*, Daniel Soudry

# Better models - parallelization is crucial



AlexNet [Krizhevsky et al. 2012]: model split on two GPUs

• Model parallelism:

Split model (same data)

• Data parallelism:

Split data (same model)

$$\Delta \mathbf{w} \propto -\frac{1}{b} \sum_{n=1}^{b} \nabla_{\mathbf{w}} L_n \left( \mathbf{w} \right)$$
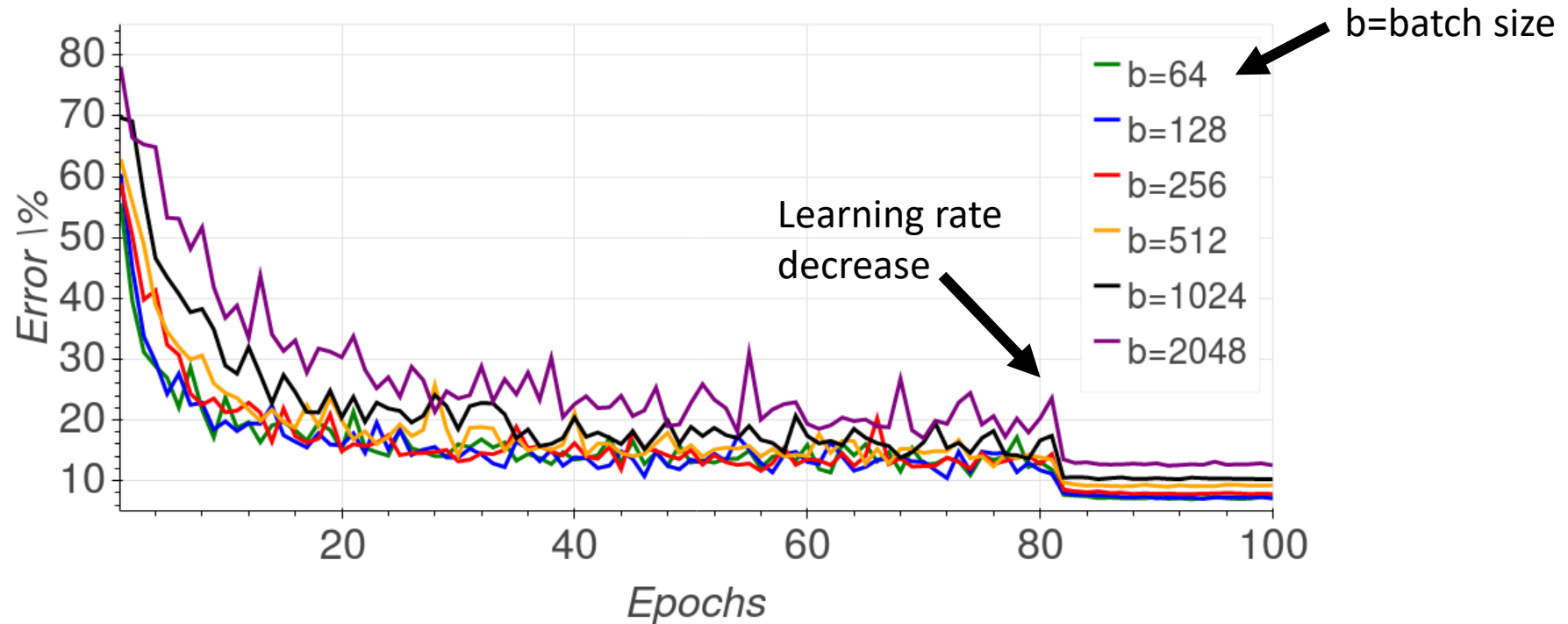
SGD: weight update proportional to gradients averaged over mini batch

**Can we increase batch size and improve parallelization?**

# Large batch size hurts generalization?

**Dataset:** CIFAR10, **Architecture:** Resnet44, **Training**: SGD + momentum (+ gradient clipping )
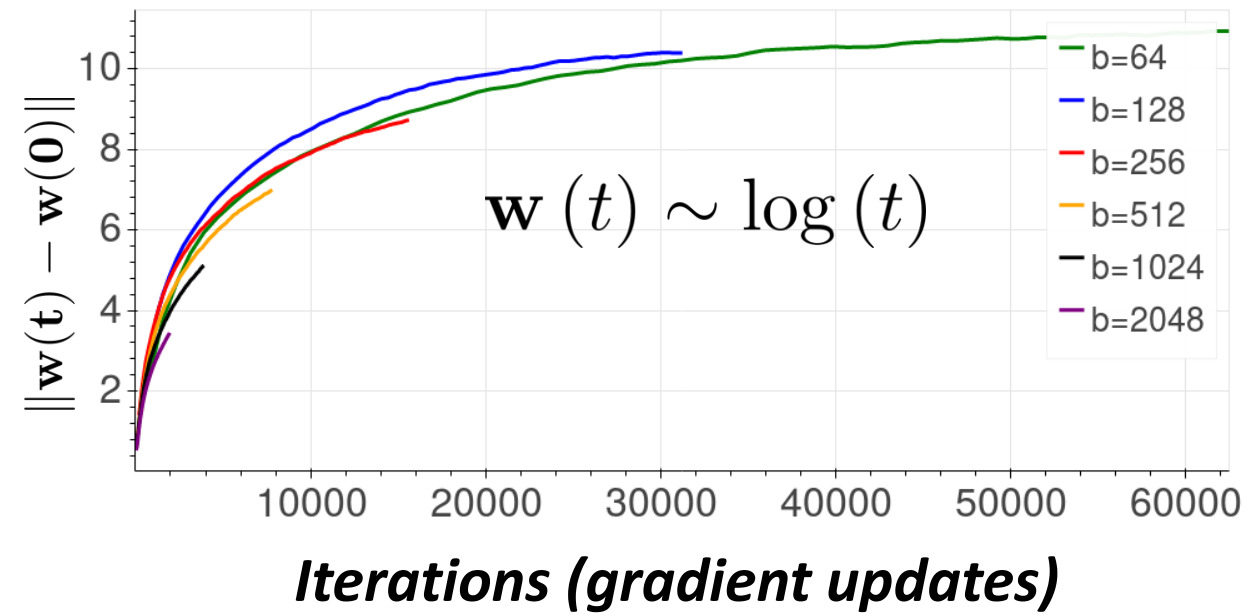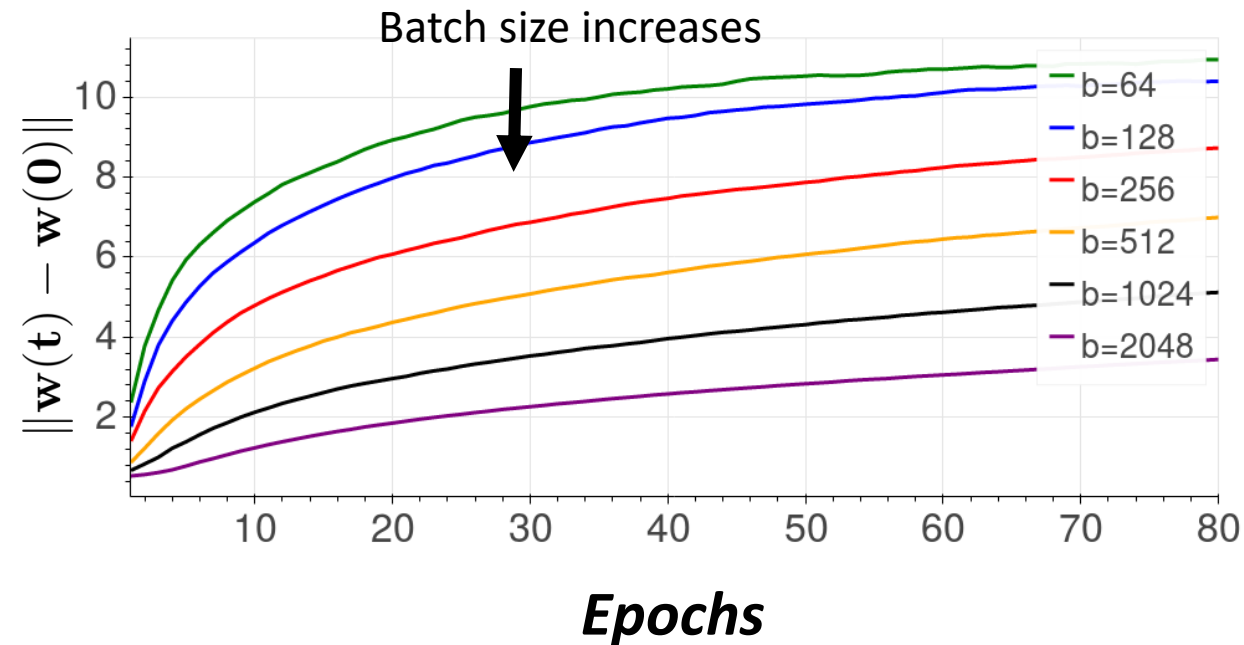
**Why?**



- Generalization gap persisted in models trained "without any budget or limits, until the loss function ceased to improve" [Keskar et al. 2017]

# Observation

Weight distances from initialization increase

logarithmically **with iterations**



Batch size increases

$\|\mathbf{w}(t) - \mathbf{w}(0)\|$

$\mathbf{w}(t) \sim \log(t)$

b=64
b=128
b=256
b=512
b=1024
b=2048

*Epochs*

*Iterations (gradient updates)*

Why logarithmic behavior? Theory later…

# Experimental details

- We experiment with various datasets and models
- Optimizing using SGD + momentum + gradient clipping
    - Usually generalize better than adaptive methods (e.g Adam)
    - Grad clipping effectively creates a "warm-up" phase
- Noticeable generalization gap between small and large batch

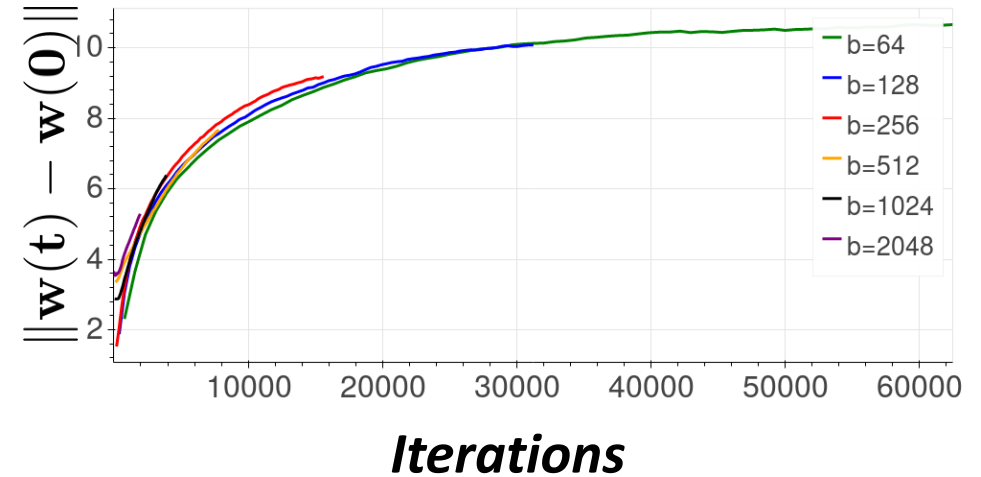| Network | Dataset | SB | LB |
|---|---|---|---|
| F1 (Keskar et al., 2017) | MNIST | 98.27% | 97.05% |
| C1 (Keskar et al., 2017) | Cifar10 | 87.80% | 83.95% |
| Resnet44 (He et al., 2016) | Cifar10 | 92.83% | 86.10% |
| VGG (Simonyan, 2014) | Cifar10 | 92.30% | 84.1% |
| C3 (Keskar et al., 2017) | Cifar100 | 61.25% | 51.50% |
| WResnet16-4 (Zagoruyko, 2016) | Cifar100 | 73.70% | 68.15% |

# Closing the generalization gap (2/4)

- Adapt learning rate. In CIFAR $\propto \sqrt{b}$
  - Idea: mimic small batch gradient statistics (dataset dependent)

- Noticeably improves generalization, the gap remains

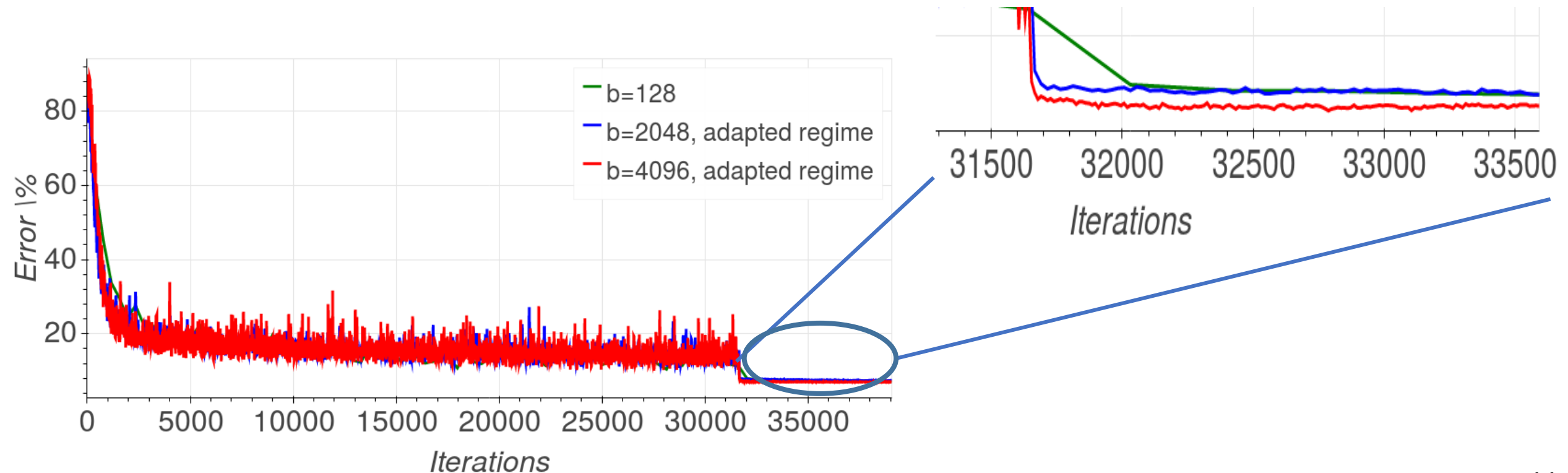| Network | Dataset | SB | LB | +LR |
| --- | --- | --- | --- | --- |
| F1 (Keskar et al., 2017) | MNIST | 98.27% | 97.05% | 97.55% |
| C1 (Keskar et al., 2017) | Cifar10 | 87.80% | 83.95% | 86.15% |
| Resnet44 (He et al., 2016) | Cifar10 | 92.83% | 86.10% | 89.30% |
| VGG (Simonyan, 2014) | Cifar10 | 92.30% | 84.1% | 88.6% |
| C3 (Keskar et al., 2017) | Cifar100 | 61.25% | 51.50% | 57.38% |
| WResnet16-4 (Zagoruyko, 2016) | Cifar100 | 73.70% | 68.15% | 69.05% |

# Graph indicates: not enough iterations?

- Using these modifications – distance from initialization now better matched

- However, graph indicates: insufficient iterations with large batch



*Iterations*

| Network | Dataset | SB | LB | +LR | +GBN |
|---|---|---|---|---|---|
| F1 (Keskar et al., 2017) | MNIST | 98.27% | 97.05% | 97.55% | 97.60% |
| C1 (Keskar et al., 2017) | Cifar10 | 87.80% | 83.95% | 86.15% | 86.4% |
| Resnet44 (He et al., 2016) | Cifar10 | 92.83% | 86.10% | 89.30% | 90.50% |
| VGG (Simonyan, 2014) | Cifar10 | 92.30% | 84.1% | 88.6% | 91.50% |
| C3 (Keskar et al., 2017) | Cifar100 | 61.25% | 51.50% | 57.38% | 57.5% |
| WResnet16-4 (Zagoruyko, 2016) | Cifar100 | 73.70% | 68.15% | 69.05% | 71.20% |

# Train longer, generalize better

- With sufficient iterations in "plateau" region,
  generalization gap vanish:

# Closing the generalization gap (4/4)

- Regime Adaptation – train so that the number of iterations is fixed for all batch sizes (train longer number of epochs)
  - Completely closes the generalization gap

| Network | Dataset | SB | LB | +LR | +GBN | +RA |
|---|---|---|---|---|---|---|
| F1 (Keskar et al., 2017) | MNIST | 98.27% | 97.05% | 97.55% | 97.60% | 98.53% |
| C1 (Keskar et al., 2017) | Cifar10 | 87.80% | 83.95% | 86.15% | 86.4% | 88.20% |
| Resnet44 (He et al., 2016) | Cifar10 | 92.83% | 86.10% | 89.30% | 90.50% | 93.07% |
| VGG (Simonyan, 2014) | Cifar10 | 92.30% | 84.1% | 88.6% | 91.50% | 93.03% |
| C3 (Keskar et al., 2017) | Cifar100 | 61.25% | 51.50% | 57.38% | 57.5% | 63.20% |
| WResnet16-4 (Zagoruyko, 2016) | Cifar100 | 73.70% | 68.15% | 69.05% | 71.20% | 73.57% |

ImageNet (AlexNet):

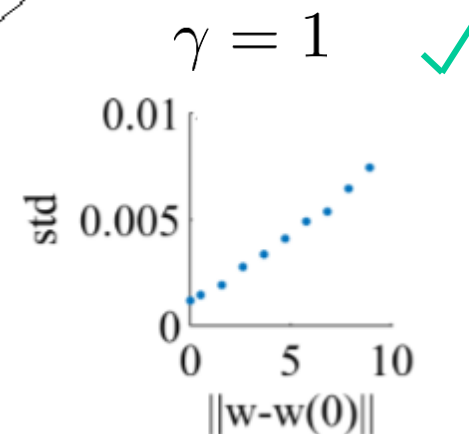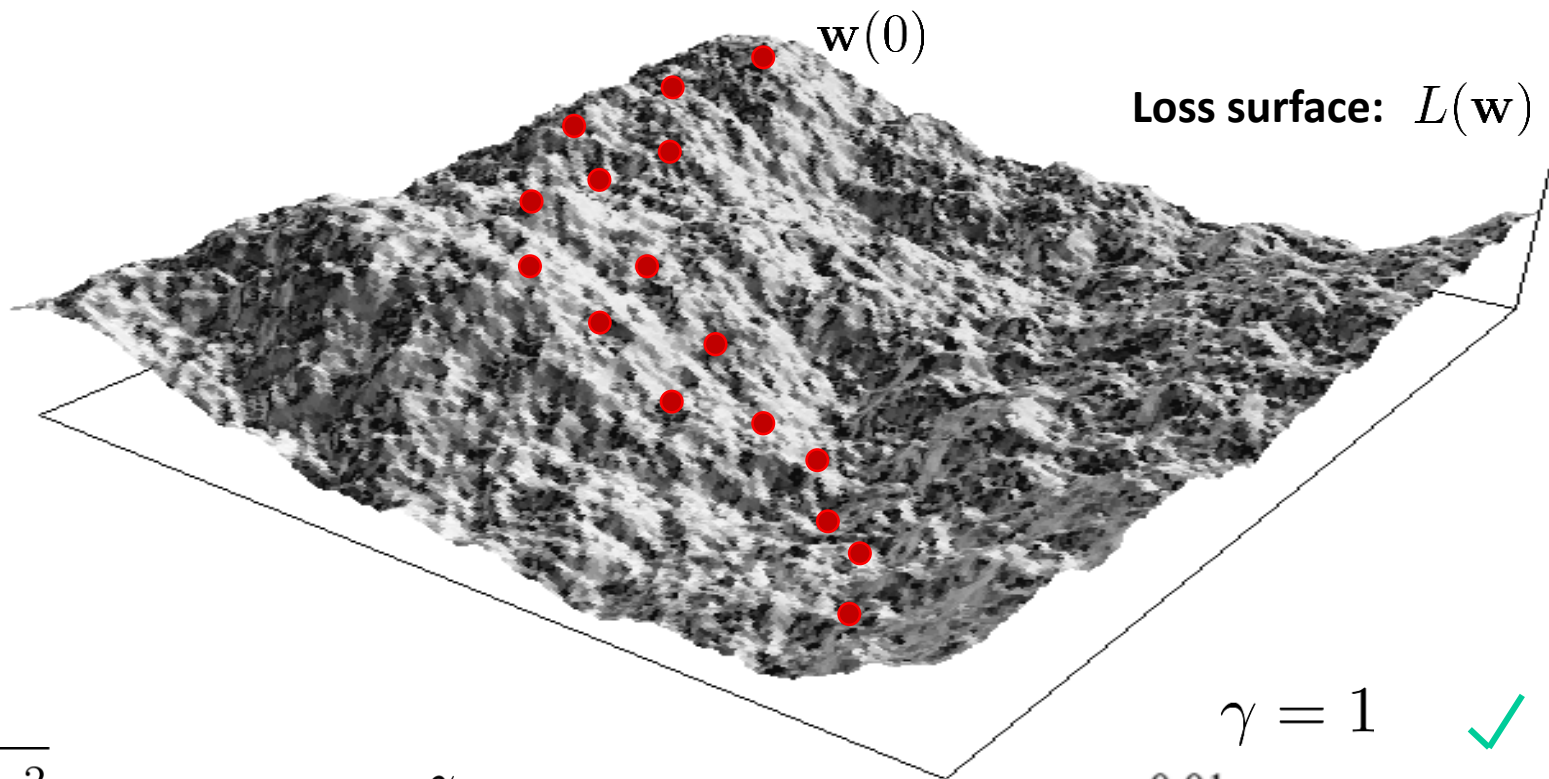| LB size | Dataset | SB | $LB^8$ | $+LR^8$ | +GBN | +RA |
|---|---|---|---|---|---|---|
| 4096 | ImageNet | 57.10% | 41.23% | 53.25% | 54.92% | 59.5% |
| 8192 | ImageNet | 57.10% | 41.23% | 53.25% | 53.93% | 59.5% |

# Why weight distances increase logarithmically?

**Hypothesis:**
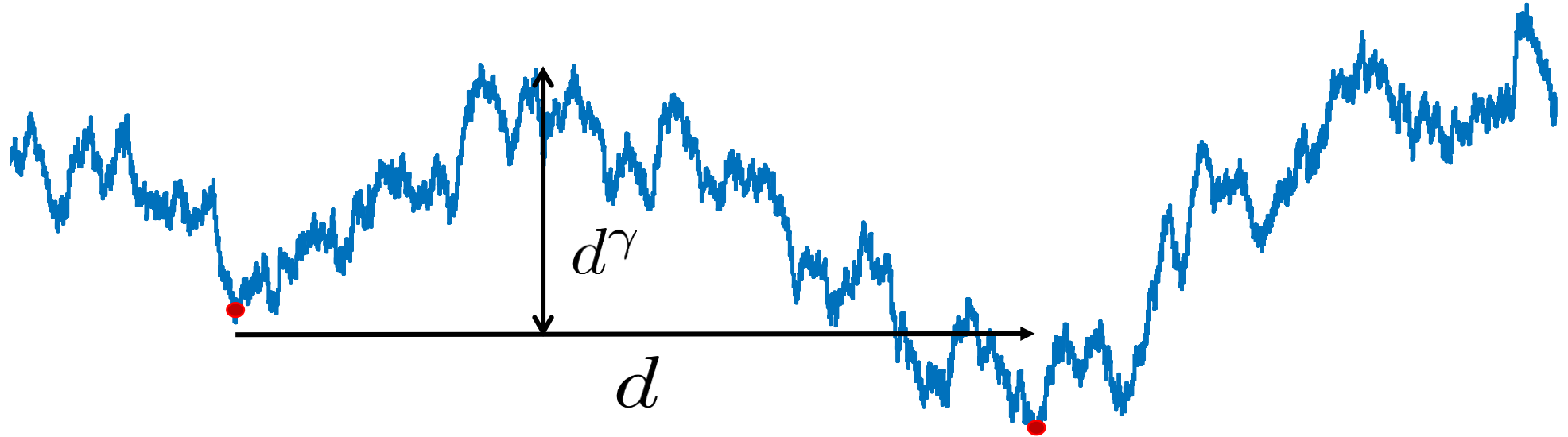
During initial high learning rate phase: "random walk on a random potential" where

$$\text{std} \triangleq \sqrt{\mathbb{E}\left(L\left(\mathbf{w}\right) - L\left(\mathbf{w}(0)\right)\right)^2} \sim \|\mathbf{w} - \mathbf{w}(0)\|^{\gamma}$$

**Marinari et al., 1983:** $\quad \mathbf{w}\left(t\right) \sim \log^{\frac{1}{\gamma}}\left(t\right) \quad$ **"ultra-slow diffusion"**



$\mathbf{w}(0)$

**Loss surface:** $L(\mathbf{w})$

$\gamma = 1 \quad \checkmark$

13

# Ultra-slow diffusion: Basic idea



Time to pass tallest barrier: $\quad t \propto \exp(d^{\gamma}) \qquad \Rightarrow d \propto \log^{\frac{1}{\gamma}}(t)$

# Summary so far

- **Q:** Is there inherent generalization problem with large batches?
  **A:** Observed: no, just adjust training regime.

- **Q:** What is the mechanism behind training dynamics?
  **A:** Hypothesis: "random walk on a random potential"

- **Q:** Can we reduce the total wall clock time?
  **A:** Yes, in some models
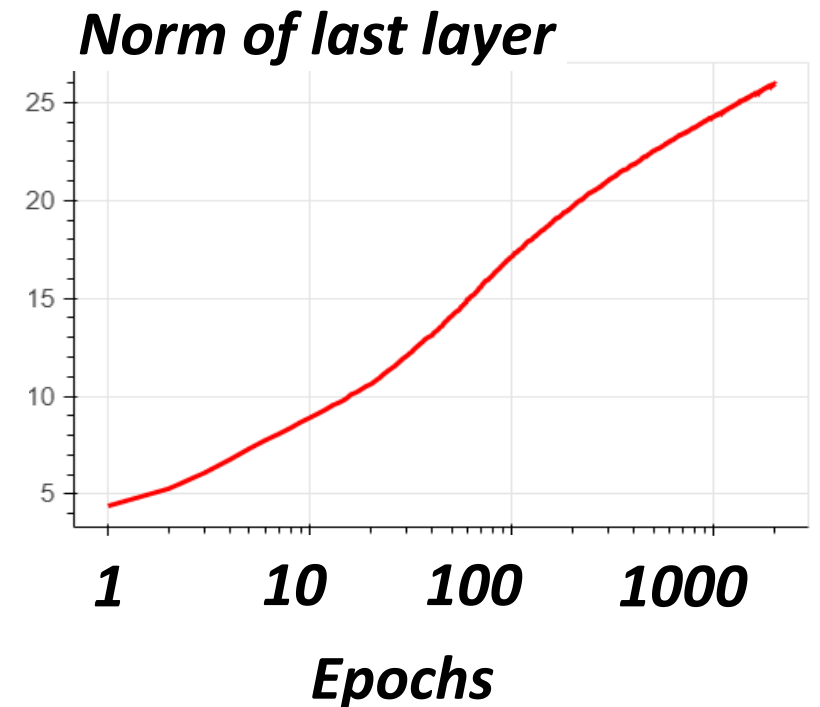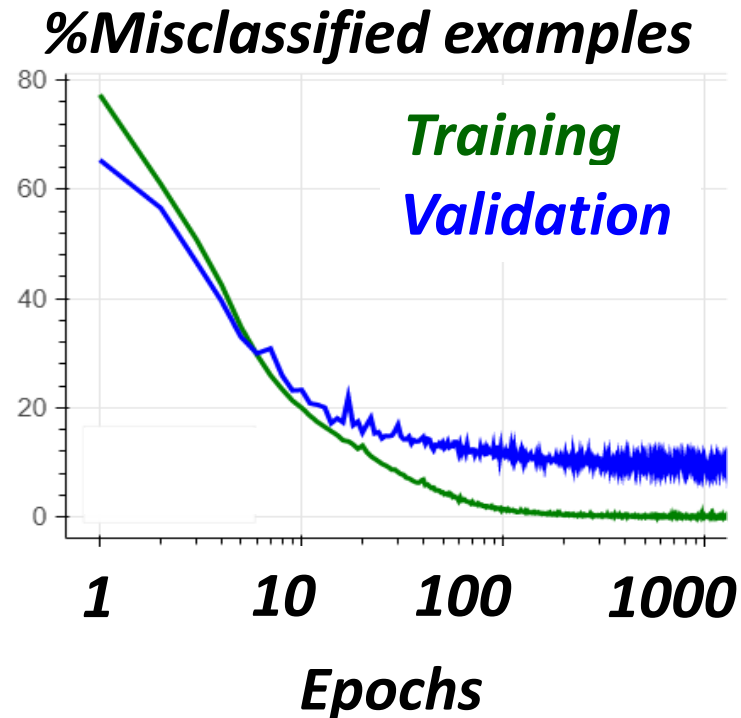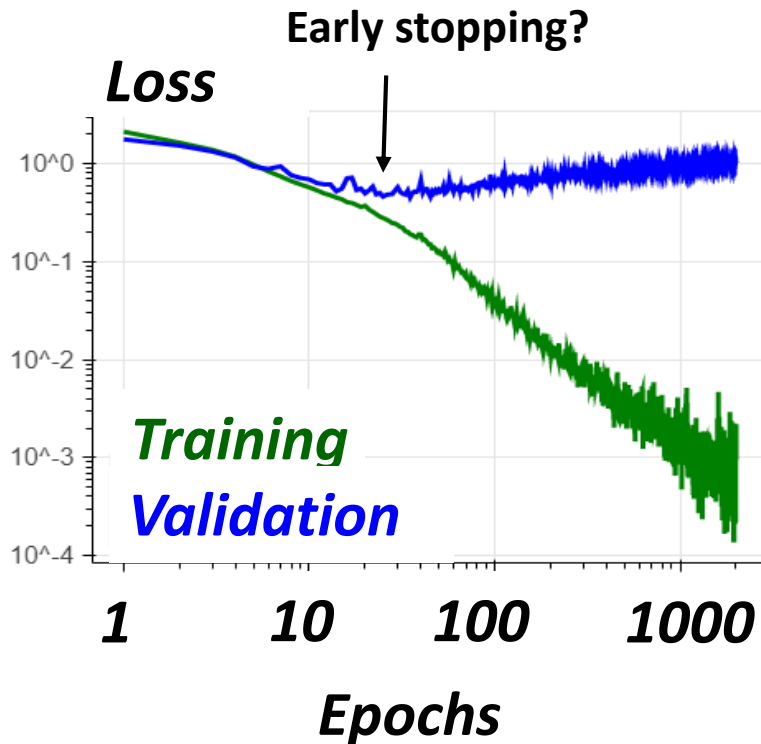
# Significant speed-ups possible

**Accurate, Large Minibatch SGD: Training ImageNet in <u>1 Hour</u>**

Goyal et al. (Facebook whitepaper, two weeks after us)

- Large scale experiments: ResNet over ImageNet, 256 GPUs
- Similar methods, except learning rate
- X29 times faster than a single worker


- More followed:
  - **Large Batch Training of Convolutional Networks** (You et al.)
  - **ImageNet Training in Minutes** (You et al.)
  - **Extremely Large Minibatch SGD: Training ResNet-50 on ImageNet in 15 Minutes** (Akiba et al.)

# 2) Why "Overfitting" is good for generalization?

- In contrast to common practice: good generalization results from many gradient updates in an "overfitting regime"



17

# Peculiar generalization dynamics - summary

- Validation Loss increases
- Training error + loss goes to zero
- Weight Norm diverges

Looks like we are overfitting... but

- Validation error (classification) seems to never stop decreasing (slowly)

**Conclusion:** No need for early stopping (!)

**How all of this makes sense?**

# Why "Overfitting" is good for generalization?

- Can be shown to happen for logistic regression on separable data!
- Slow convergence to max-margin solution

**The Implicit Bias of Gradient Descent on Separable Data** (ICLR 2018)
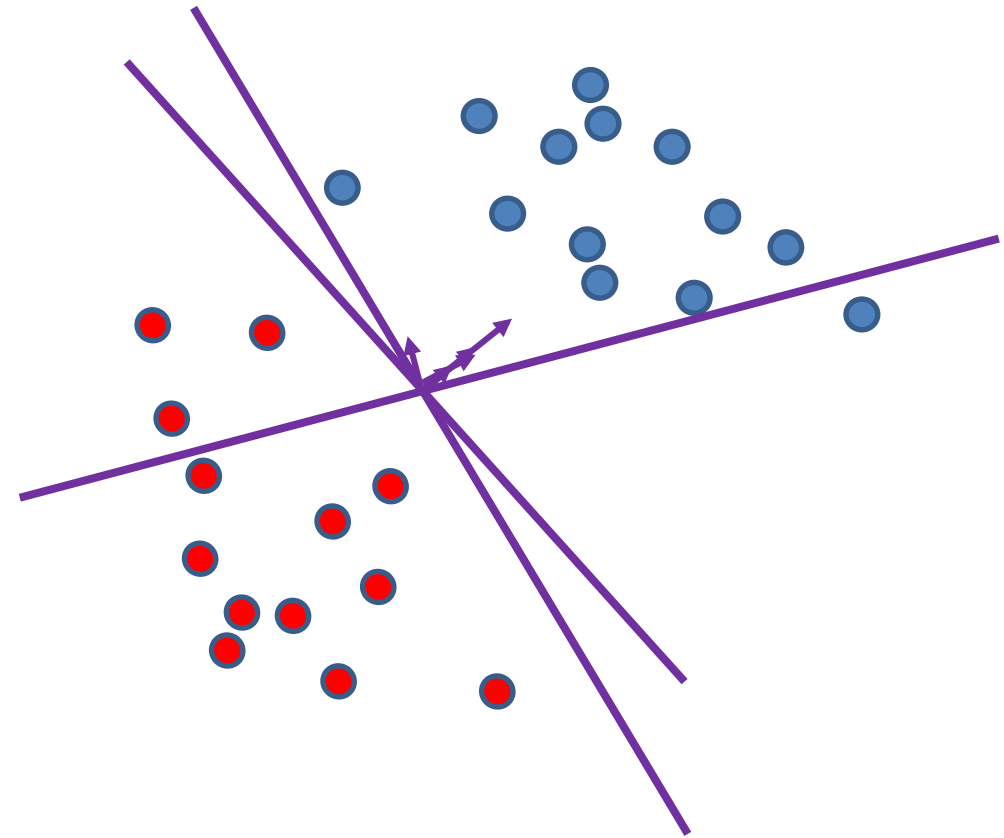- *Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Nati Srebro*

# Main Theorem

Gradient descent on logistic loss: $\Delta \mathbf{w} = -\eta \nabla \mathcal{L}(\mathbf{w})$

**Theorem 1:** $\mathbf{w}(t) = \hat{\mathbf{w}} \log t + \boldsymbol{\rho}(t)$ ,

$\hat{\mathbf{w}}$ is the (L2) max margin vector

$\boldsymbol{\rho}(t)$ is bounded, for almost every dataset.

Therefore: $\dfrac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} \longrightarrow \dfrac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}$
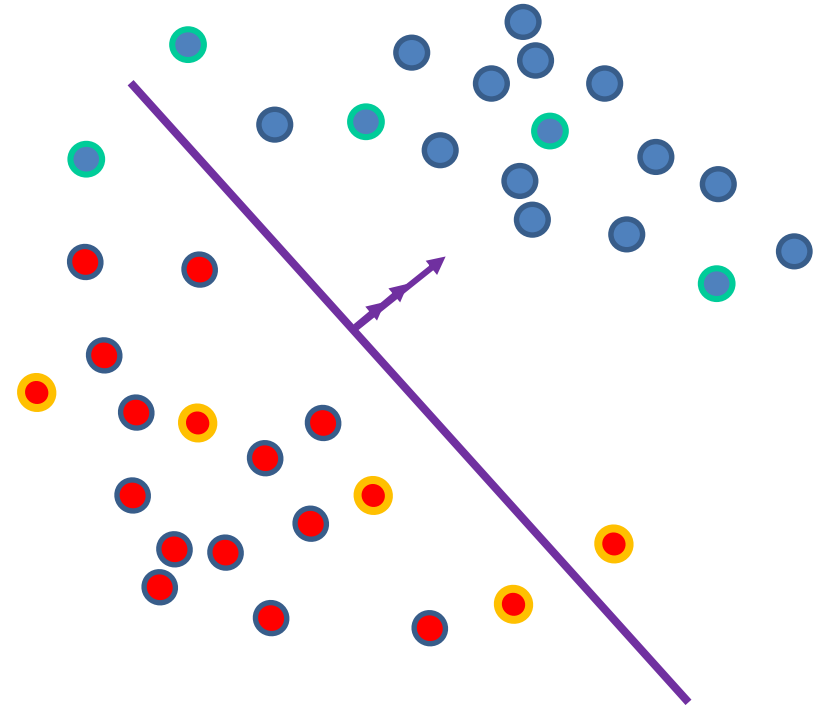
# … While expected loss (and test loss) increases

$$\mathbf{w}\,(t) = \hat{\mathbf{w}}\log t + \boldsymbol{\rho}\,(t)$$

**Expected loss:**

$$\mathbb{E}[\mathcal{L}\,(\mathbf{w})] = \Omega(\log t).$$

**Also true for test loss**

**Validation loss is expected to increase, although accuracy may still improve!**
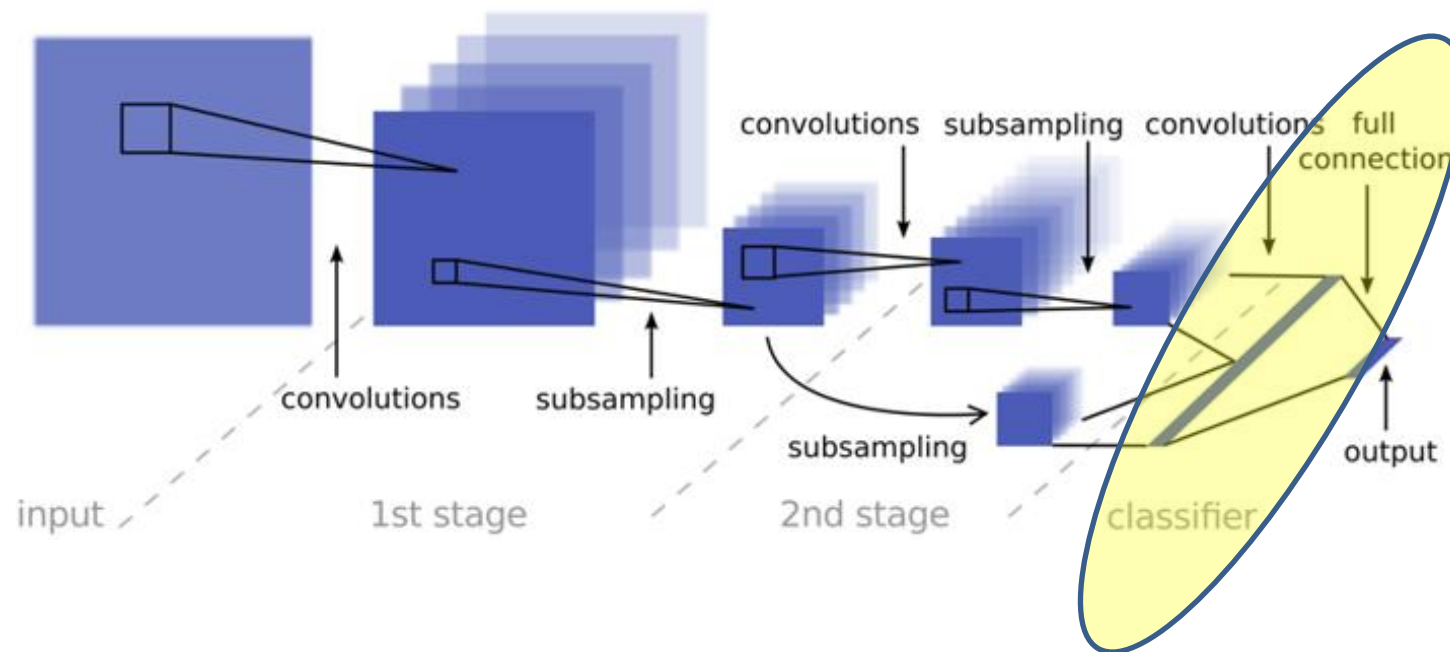
# 3) The role of the final classifier

**Fix your classifier: the marginal value of training the last weight layer**

Elad Hoffer, Itay Hubara, Daniel Soudry

# Fully connected classifier

We focus on the final representation obtained by the network $F$ before the classifier $x = F(z; \theta)$ (the last hidden layer) .

# Fully connected classifier

- In common NN models, this representation is followed by an additional affine transformation on $x \in \mathbb{R}^C$ to all possible classes $C$.

$$y = W^T x + b$$

- where the number of parameters is dependent on number of classes $W \in \mathbb{R}^{N \times C}$ (can grow to be very large)

# Fully connected classifier

Training is done using cross-entropy loss, by feeding the network outputs through a softmax activation

$$v_i = \frac{e^{y_i}}{\sum_j^C e^{y_j}}, \ \ i \in \{1, \dots, C\}$$

and reducing the expected negative log likelihood with respect to ground-truth target

$$\mathcal{L}(x, t) = -\log v_t = -w_t \cdot x - b_t + \log\left(\sum_j^C e^{w_j \cdot x + b_j}\right)$$

where $w_i$ is the $i$-th column of $W$.

# Fully connected classifier?

But the final fully-connected transform is a linear classifier:

- – The network learns features that are already separable at this point.
- – They fully-connected layers are also notoriously redundant -- easily compressed and discarded

Can they be removed completely?

# Fixed classifier

- To evaluate our conjecture, we replaced the trainable parameter matrix $W$ with a fixed orthonormal projection

$$Q \in \mathbb{R}^{N \times C} \text{ such that } QQ^T = I_n$$

- As the rows of classifier weight matrix are fixed with an equal $L_2$ norm, we also restrict the representation of $x$ to reside on the n-dimensional sphere

$$\hat{x} = \frac{x}{\|x\|_2}$$

# Fixed classifier

Since $-1 \leq q_i \cdot \hat{x} \leq 1$ and softmax function is scale-sensitive, we introduce another temperature scaling coefficient $\alpha$

$$v_i = \frac{e^{\alpha q_i \cdot \hat{x} + b_i}}{\sum_j^C e^{\alpha q_j \cdot \hat{x} + b_j}}, \ \ i \in \{1, \ldots, C\}$$

and we minimize the loss:

$$\mathcal{L}(x, t) = -\alpha q_t \cdot \frac{x}{\|x\|_2} + b_t + \log \left( \sum_{i=1}^{C} \exp \left( \alpha q_i \cdot \frac{x}{\|x\|_2} + b_i \right) \right)$$
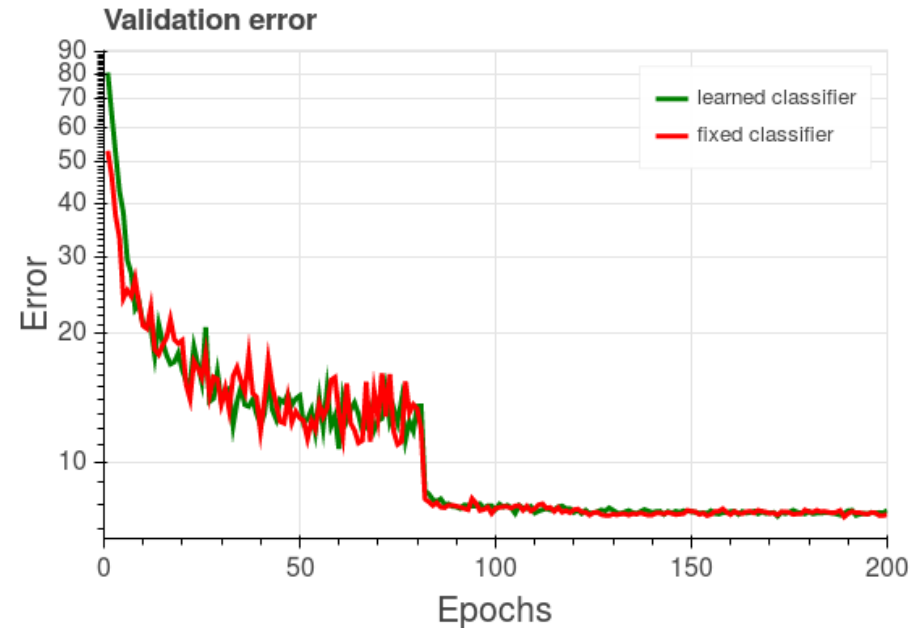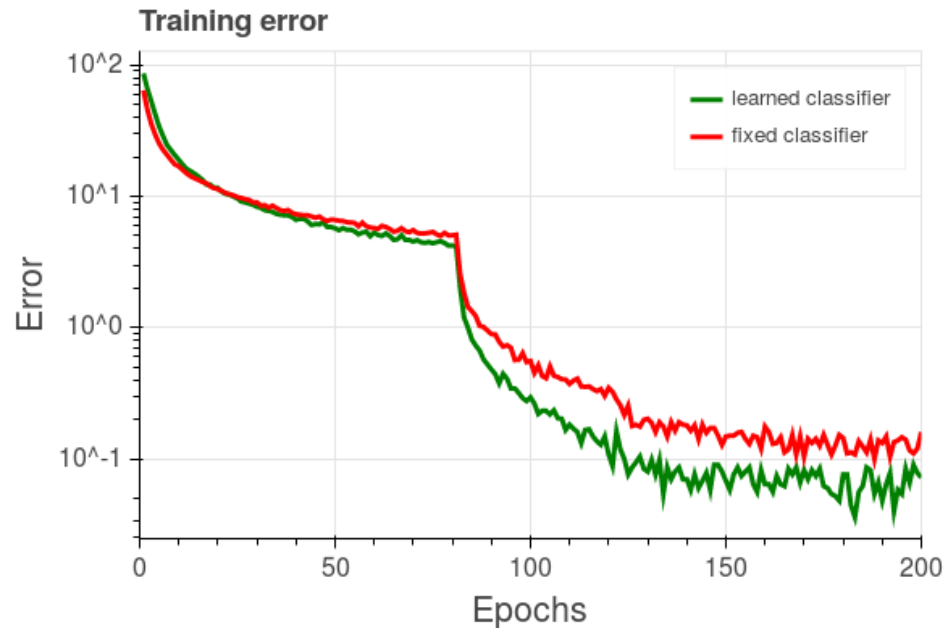
# Hadamard classifier

The fixed orthogonal weights can be chosen to be a Hadamard matrix:

$$H^T H = nI_n , \qquad H \in \{-1,1\}^n$$

- A deterministic, low-memory and easily generated matrix that can be used to classify.

- Removal of the need to perform a full matrix-matrix multiplication -- as multiplying by a Hadamard matrix can be done by simple sign manipulation and addition.

# Learned vs. Fixed classifier

- We compare the training of a **fully-learned classifier** with a **fixed classifier** (Cifar10, ResNet)
    - Training error is lower when using a learned classifier.
    - Both achieve the same accuracy on the validation set.

# Empirical results

We find that this behavior remains in other datasets and models

- Negligible decrease in accuracy when final layer is fixed
- Reduces number of weights -- e.g ShuffleNet, where most of the parameters are in the last layer

| Network | Dataset | Learned | Fixed | # Params | % Fixed params |
|---------|---------|---------|-------|----------|----------------|
| Resnet56 | Cifar10 | 93.03% | 93.14% | 855,770 | 0.07% |
| DenseNet(k=12) | Cifar100 | 77.73% | 77.67% | 800,032 | 4.2% |
| Resnet50 | ImageNet | 75.3% | 75.3% | 25,557,032 | 8.01% |
| DenseNet169 | ImageNet | 76.2% | 76% | 14,149,480 | 11.76% |
| ShuffleNet | ImageNet | 65.9% | 65.4% | 1,826,555 | 52.56% |

# Summary

- Large batch training $\not\Rightarrow$ generalization decrease
- Validation loss increases $\neq$ overfitting occurs
- Validation error monotonically decreases: no early stopping
- Linear classification layers have marginal effect on accuracy

# Thank you for your time! Questions?

For more information, visit my page at:
www.DeepLearning.co.il