



# Optimized Deep Learning - Jean Zay

---

Training and large batches



INSTITUT DU  
DÉVELOPPEMENT ET DES  
RESSOURCES EN  
INFORMATIQUE  
SCIENTIFIQUE



# Loss Landscape

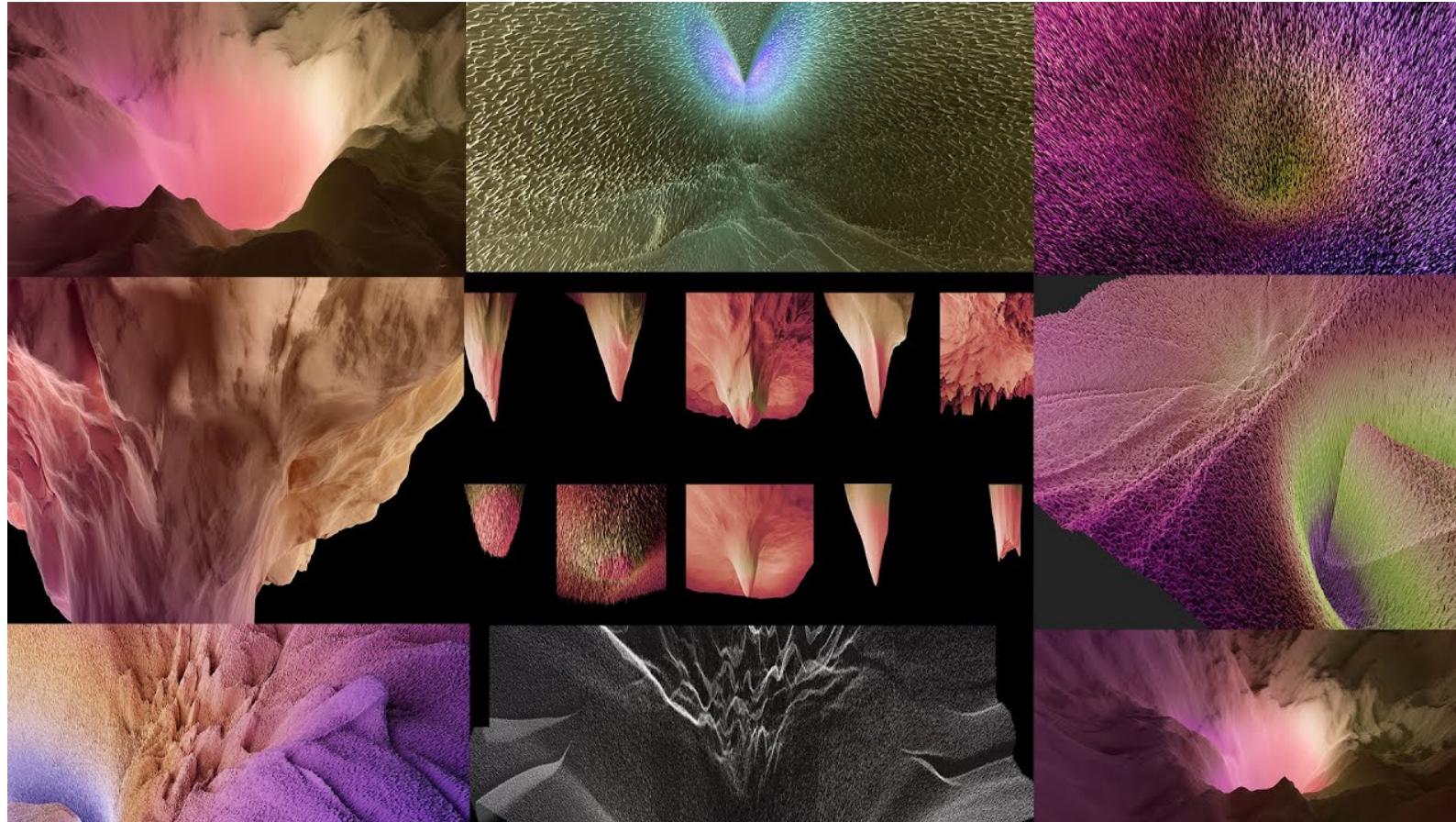
Loss Landscape ◀

Residual Learning ◀

Initialization ◀

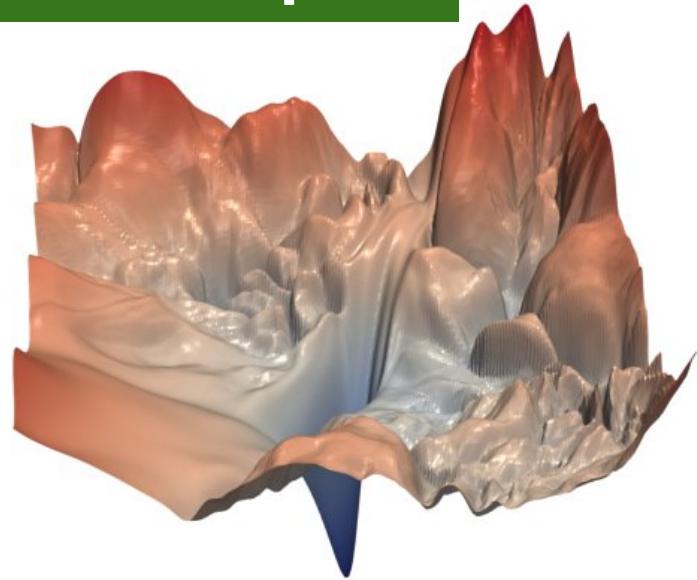
# Loss Landscape

<https://losslandscape.com/>

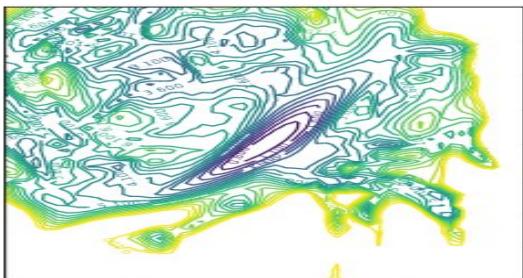


# Loss Landscape

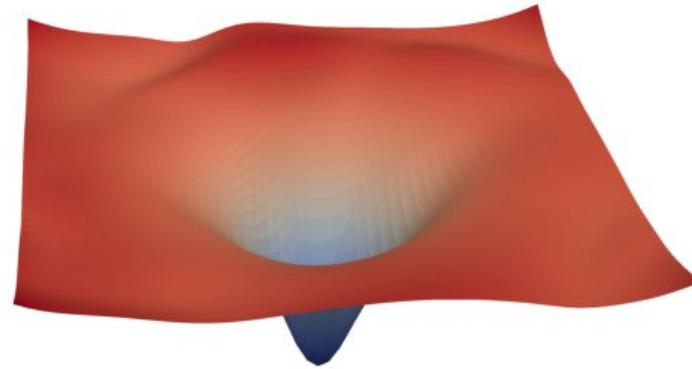
<https://arxiv.org/pdf/1712.09913.pdf>



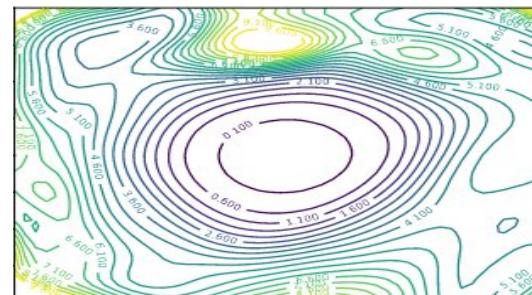
(a) without skip connections



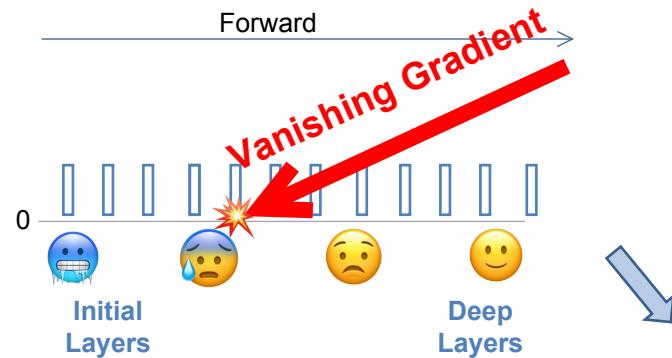
**Residual Learning**  
Since Resnets (2015) ...



(b) with skip connections

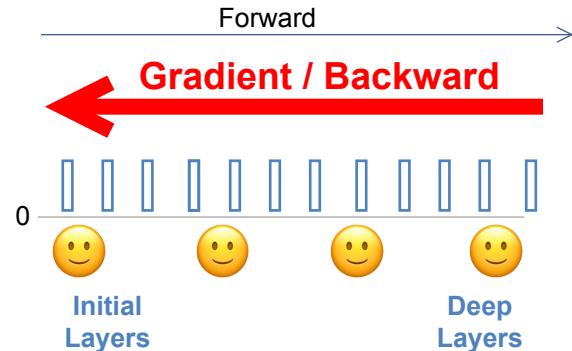


# Residual Learning

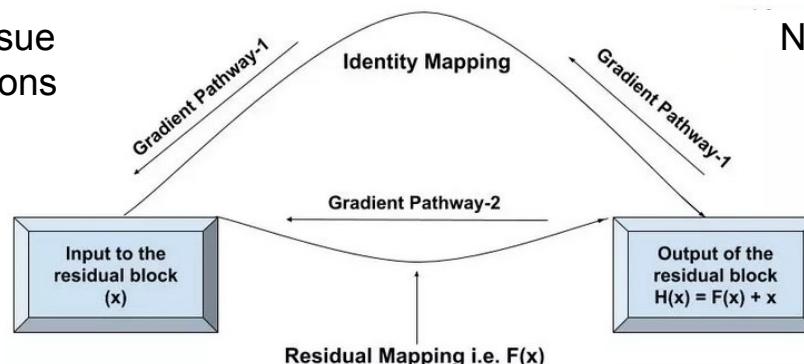


Vanishing Gradient issue  
without skip connections

**Residual Block**  
 $F(x) + x$



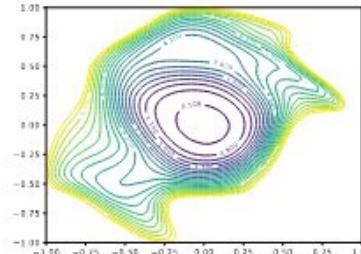
No Vanishing Gradient issue  
with skip connections



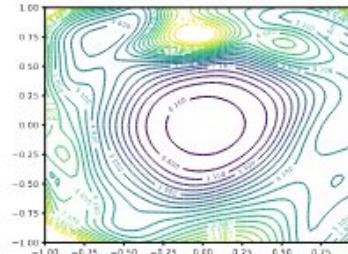
Gradient Pathways in ResNet

# Residual Learning – depth impact

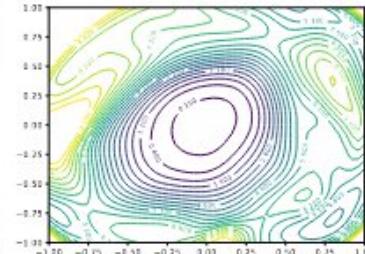
ResNet



(a) ResNet-20, 7.37%

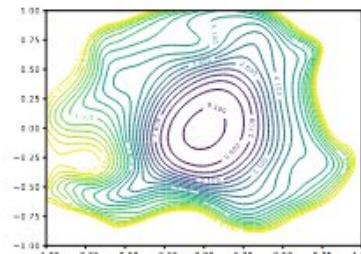


(b) ResNet-56, 5.89%

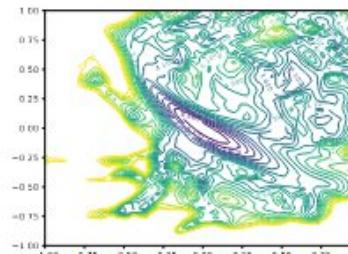


(c) ResNet-110, 5.79%

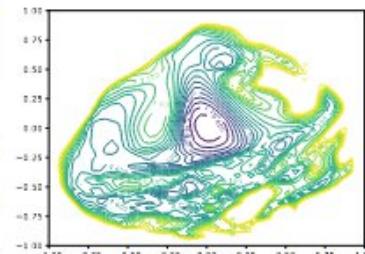
ResNet -  
No Short  
without skip connections



(d) ResNet-20-NS, 8.18%



(e) ResNet-56-NS, 13.31%



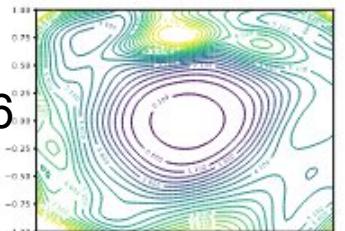
(f) ResNet-110-NS, 16.44%



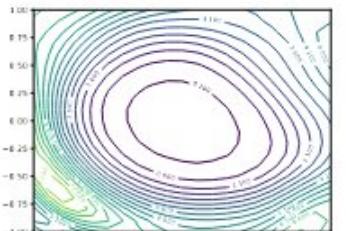
# Residual Learning – width impact



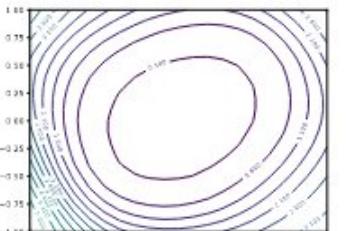
Wide-ResNet-56



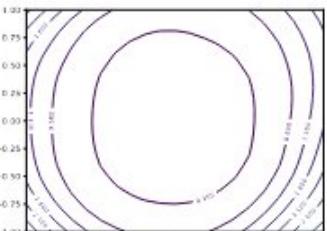
(a)  $k = 1, 5.89\%$



(b)  $k = 2, 5.07\%$

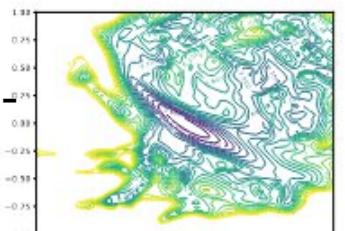


(c)  $k = 4, 4.34\%$

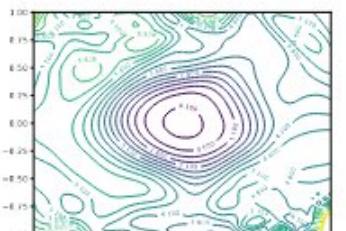


(d)  $k = 8, 3.93\%$

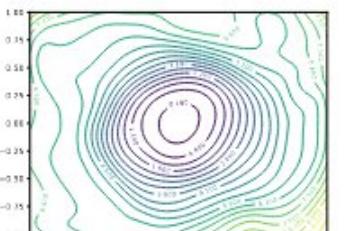
Wide-ResNet-56-  
No Short  
without skip connections



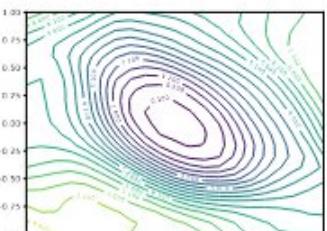
(e)  $k = 1, 13.31\%$



(f)  $k = 2, 10.26\%$



(g)  $k = 4, 9.69\%$



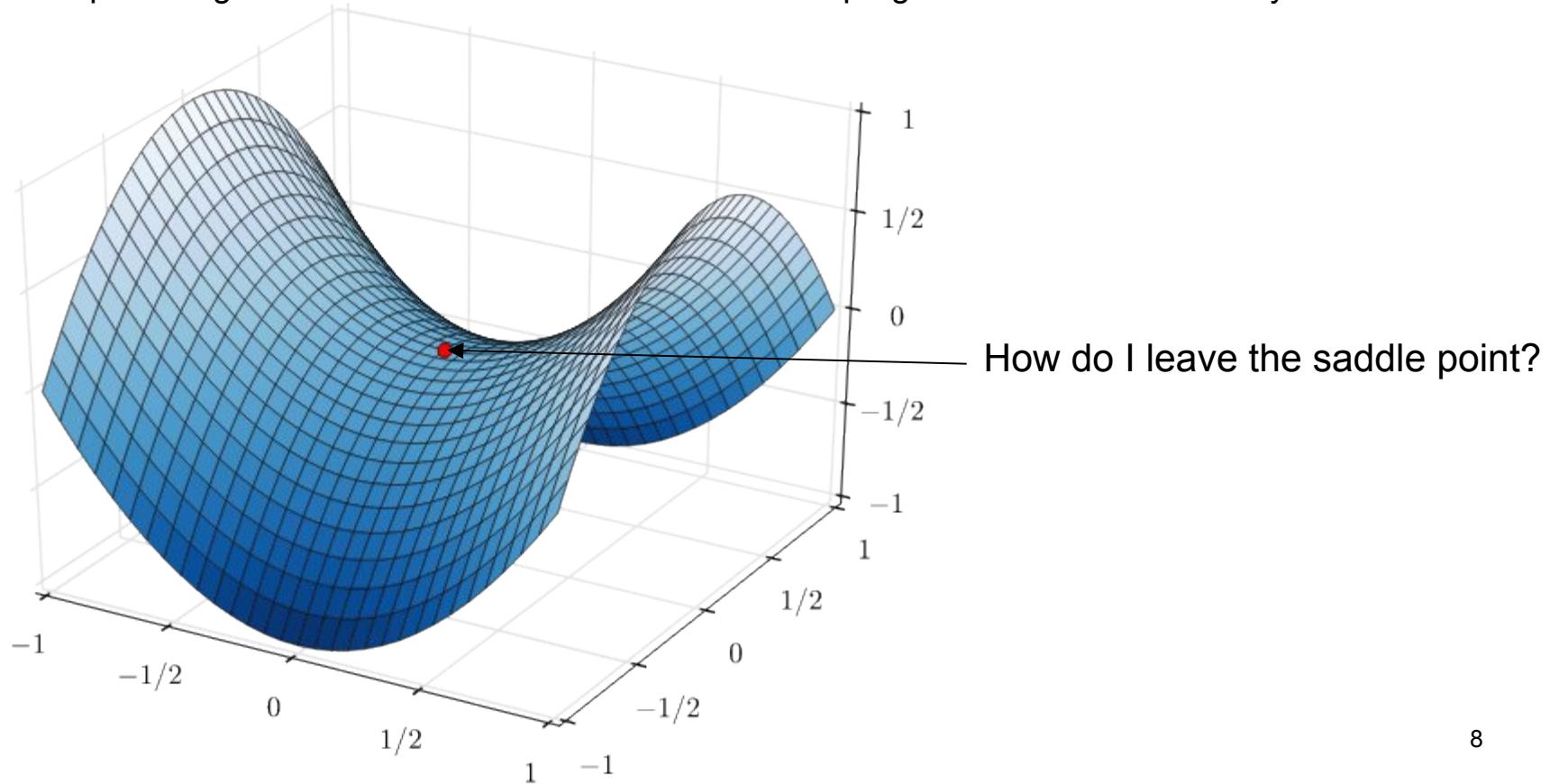
(h)  $k = 8, 8.70\%$



$k$  = channel width coefficient compared to ResNet

# Saddle point Problem

Saddle point: A gradient close to zero which will make the progression of the model very slow



# Model Parameters Initialization

The Blessing of Dimensionality :



- Xavier Initialization
  - uniform
  - normal
- Kaiming Initialization
  - uniform
  - normal

By default in PyTorch:

- Best initialization algorithm depending on the type of layer (linear, convolutional, transform, ...).
- Today, it is no longer necessary to try to optimize initialization.

# Learning Rate Scheduler

Learning rate scheduler ◀

Cyclic scheduler ◀

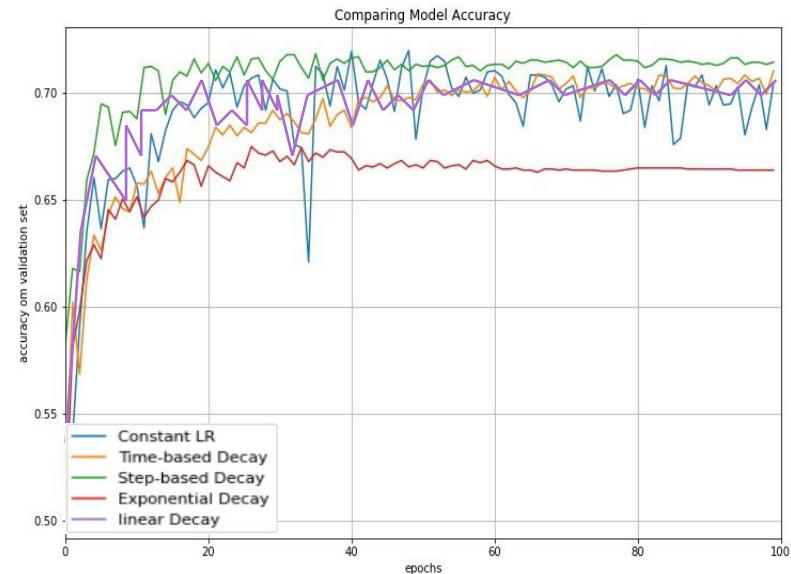
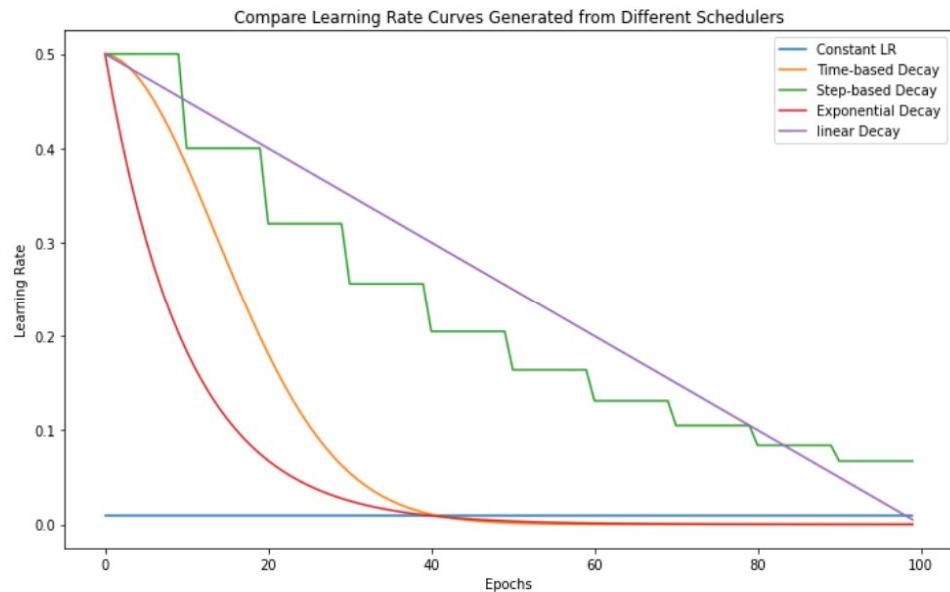
One cycle scheduler ◀

LR finder ◀

LR large batch ◀

# Learning Rate Scheduler

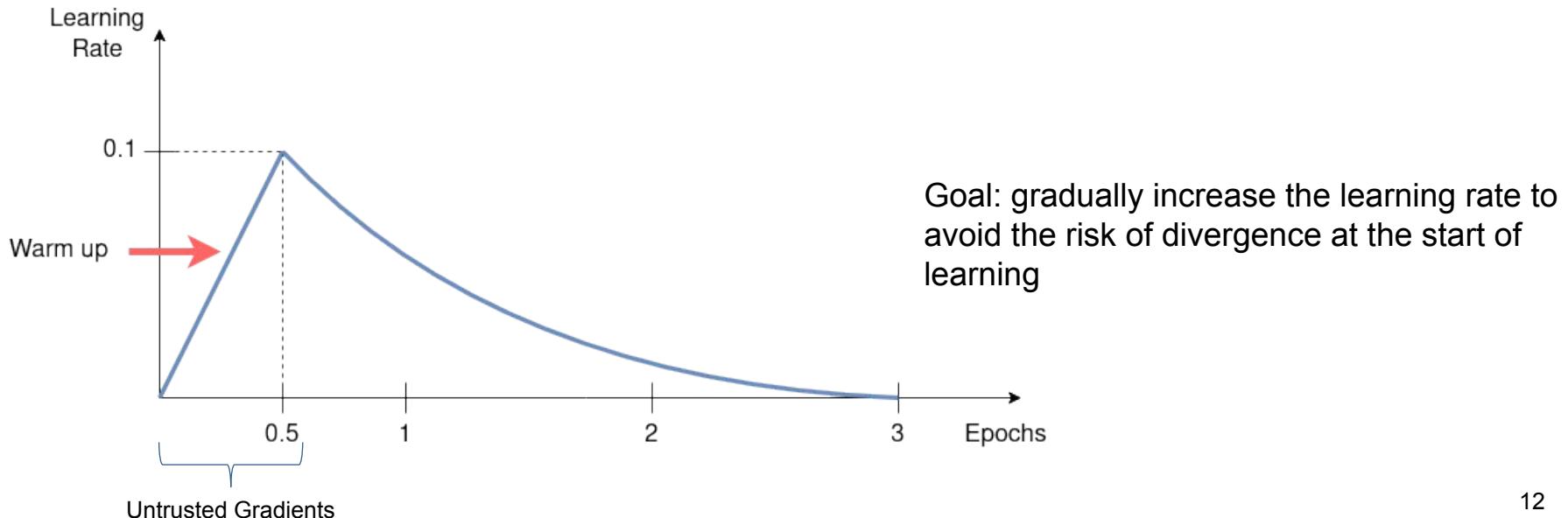
## Learning rate decay



# Learning Rate Scheduler

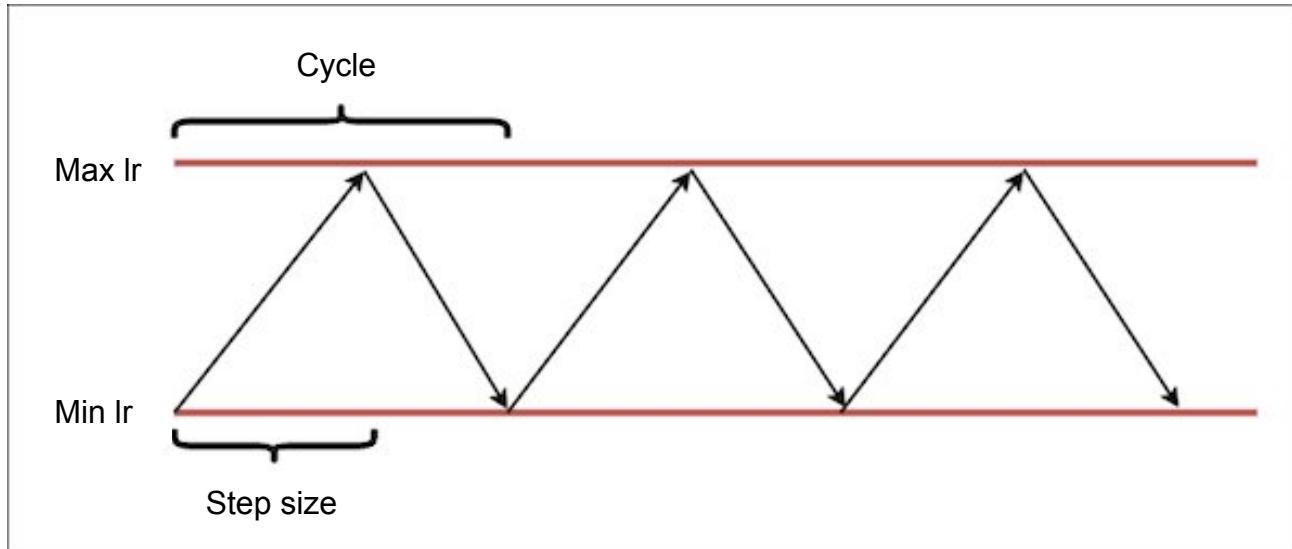
## WARMUP for *large batches*

Problems: The first iterations have too much effect on the model (significant losses, high gradients, bias, etc.), a high learning rate can cause strong instability or divergence



# Cyclic Learning Rate Scheduler

*Cyclical Learning Rates for Training Neural Networks - Leslie N. Smith 2017*

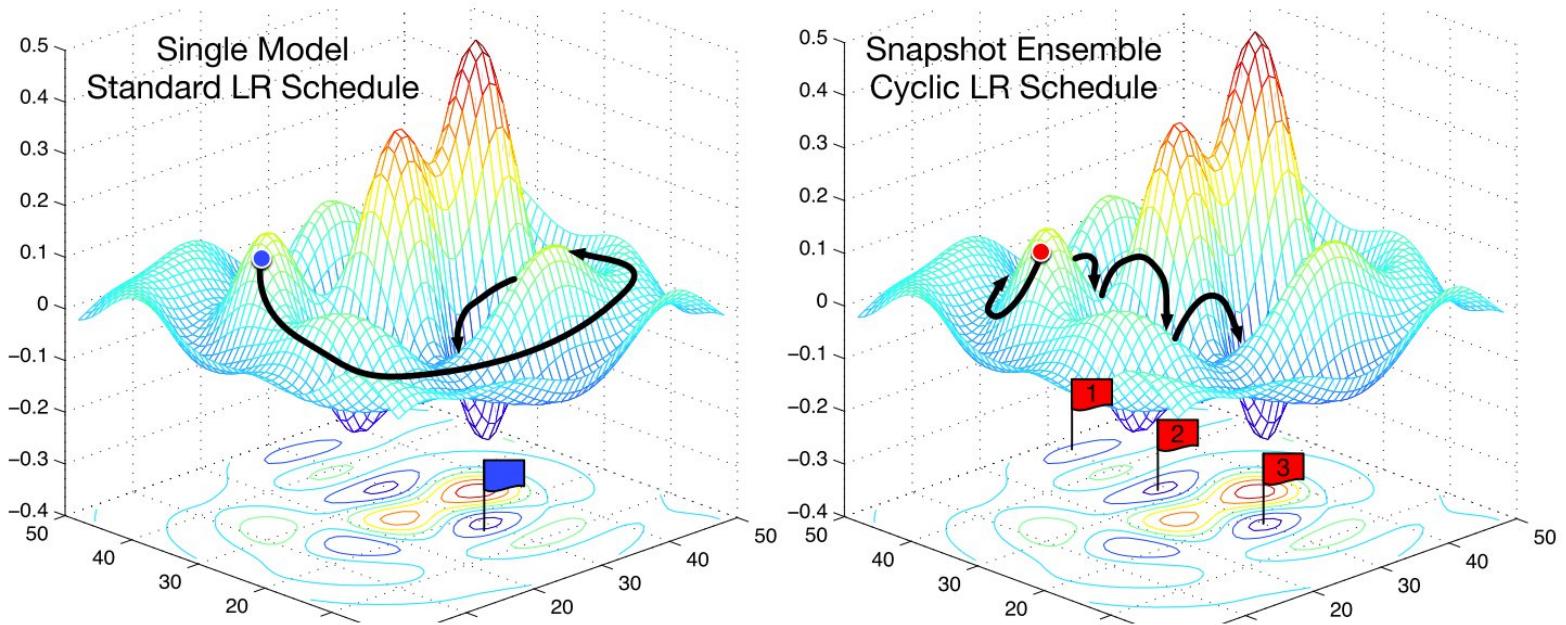


Paramètres :

- $\text{Step\_size} = x * \text{epoch}$  ( $2 \leq x \leq 10$ )
- $\text{Base\_lr} \rightarrow$  min convergence value
- $\text{max\_lr} \rightarrow$  max convergence value

Succession of warmups and learning rate decays

# Cyclic Learning Rate Scheduler



SNAPSHOT ENSEMBLES: TRAIN 1, GET M FOR FREE

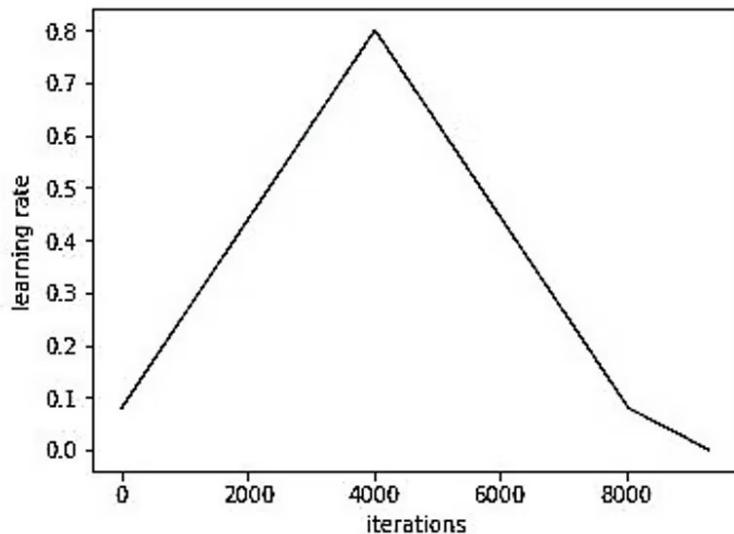
Gao Huang, Yixuan Li, Geoff Pleiss

# One Cycle Learning Rate

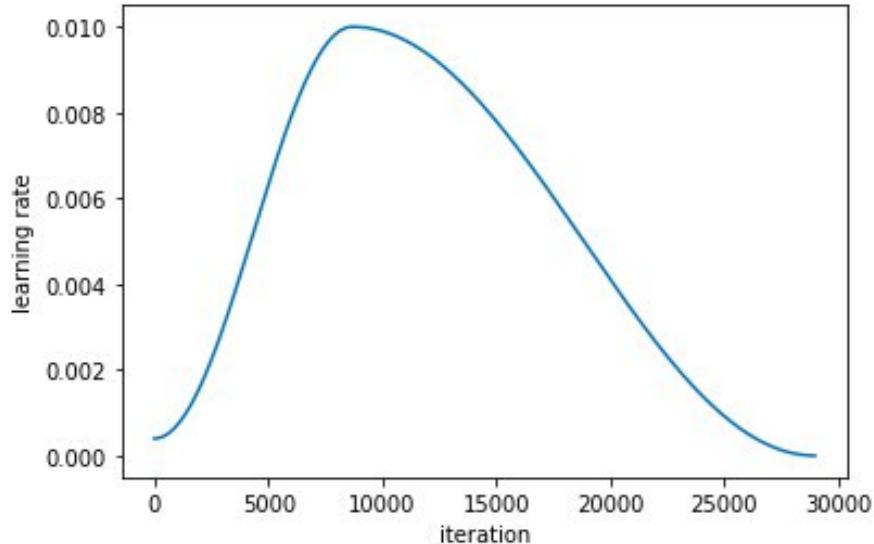
One cycle is enough!

A disciplined approach to neural network hyper-parameters -  
[Leslie N. Smith](#)

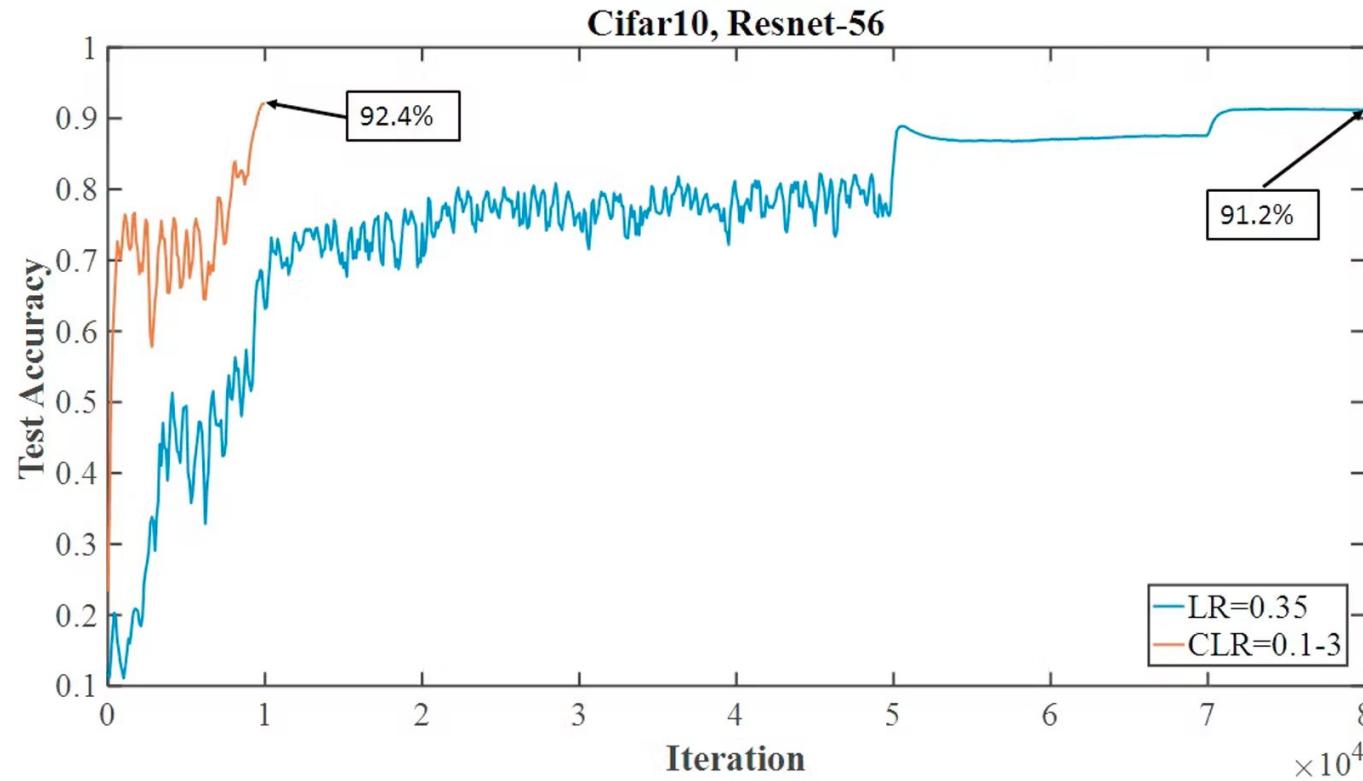
Proposition initiale



cosine annealing : Recommandation par FastAI



# One Cycle Learning Rate - *Super convergence*

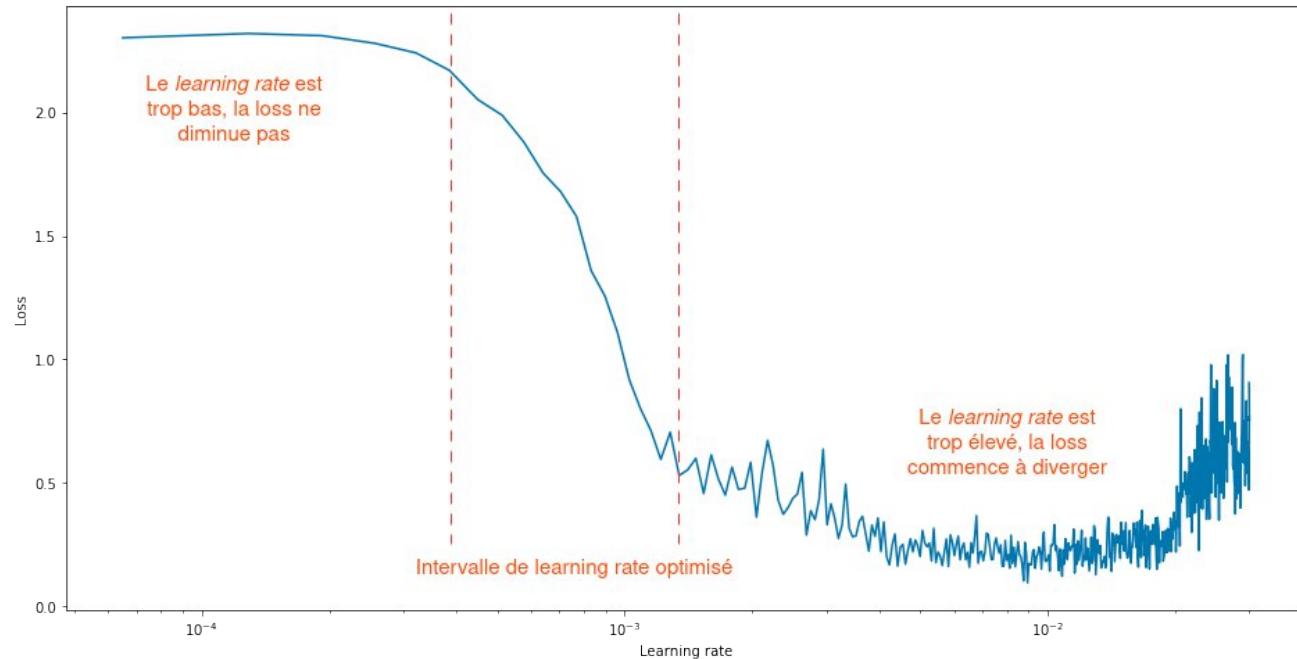


# Learning Rate Finder

Goal: Find the **optimal learning rate** values for your model, particularly for **the maximum value** of a *cyclic scheduler*

Run your model over a few epochs by increasing its learning rate

- Start of loss reduction → Minimal learning rate
- Start of loss variation → Maximum learning rate



# Learning Rate Scheduler

Each **scheduler** has *its own settings*

```
import torch.optim as opt
```

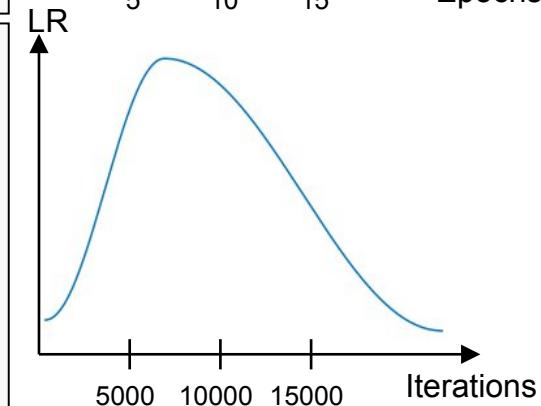
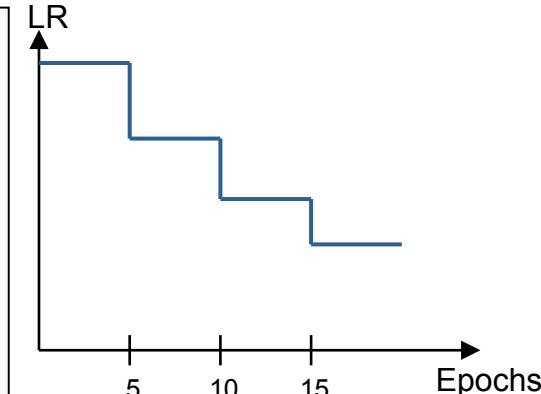
```
scheduler = opt.lr_scheduler.StepLR(optimizer, step_size=5, gamma=0.1)
```

```
for epoch in range(100):
    train(...)
    validate(...)
    scheduler.step()
```

```
import torch.optim as opt
```

```
scheduler = opt.lr_scheduler.CyclicLR(optimizer, base_lr=0.01, max_lr=0.1)
```

```
for epoch in range(10):
    for batch in data_loader:
        train_batch(...)
        scheduler.step()
```

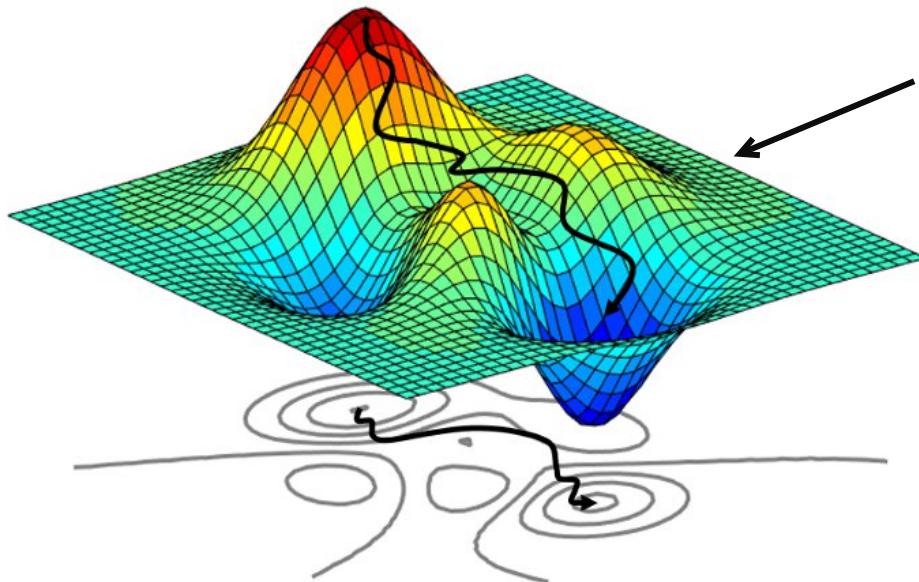


# Gradient Descent Optimizer

SGD ◀  
ADAM◀  
ADAMW◀

# Optimizer - SGD

The **optimizer** is the algorithm that **controls the gradient descent** and **the minimum search** with the aim of optimizing the learning time and the final metric.



**SGD = Stochastic Gradient Descent**

*Calculating the Gradient and updating the weights  
at each batch*

- + Batch size and learning rate adaptable according to conflicting needs:
  - Exploration to find the best local minimum
  - Acceleration of gradient descent

# SGD with Momentum

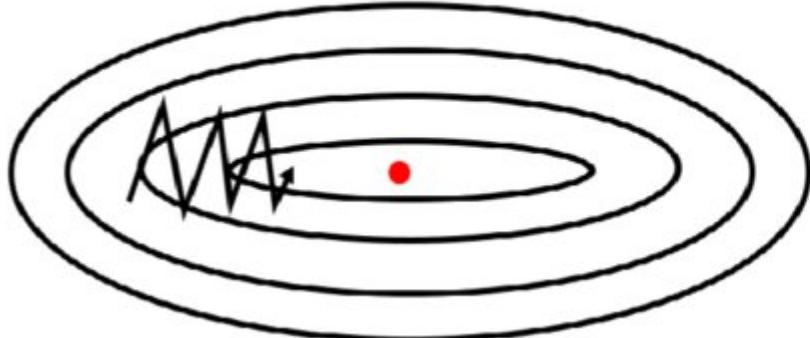
$$m_0 = 0$$

Momentum coefficient

$$m_i = \beta * m_{i-1} + (1 - \beta) * g_i$$

$$\theta_i = \theta_{i-1} - \alpha * m_i$$

SGD without momentum

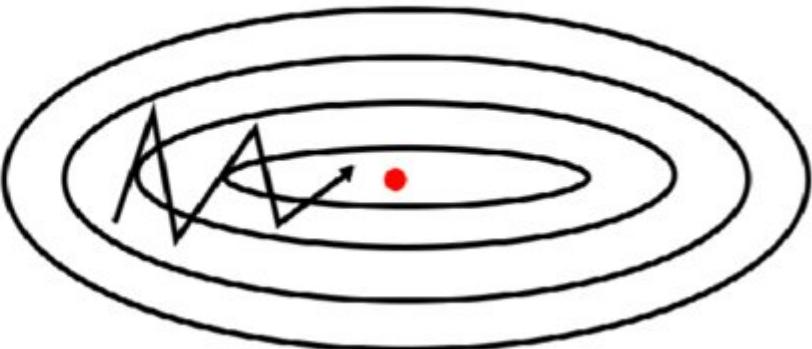


Goal: Take **previous gradients** into consideration for **faster** gradient descent.

Recommended initial value: 0.9

$$0.85 < \beta < 0.95$$

SGD with momentum

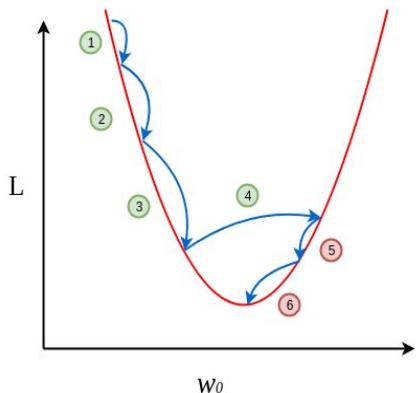


+ Allows you to **converge more quickly**

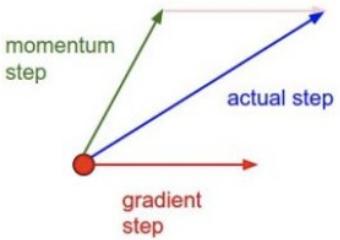
- **No guarantee** that momentum will take us in the right direction

# Momentum type

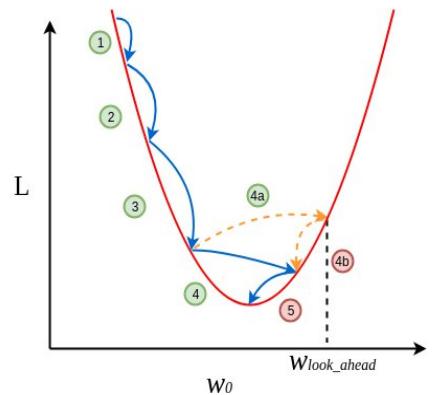
**Momentum**



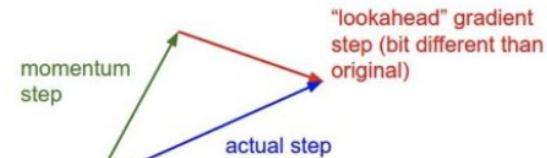
Momentum update



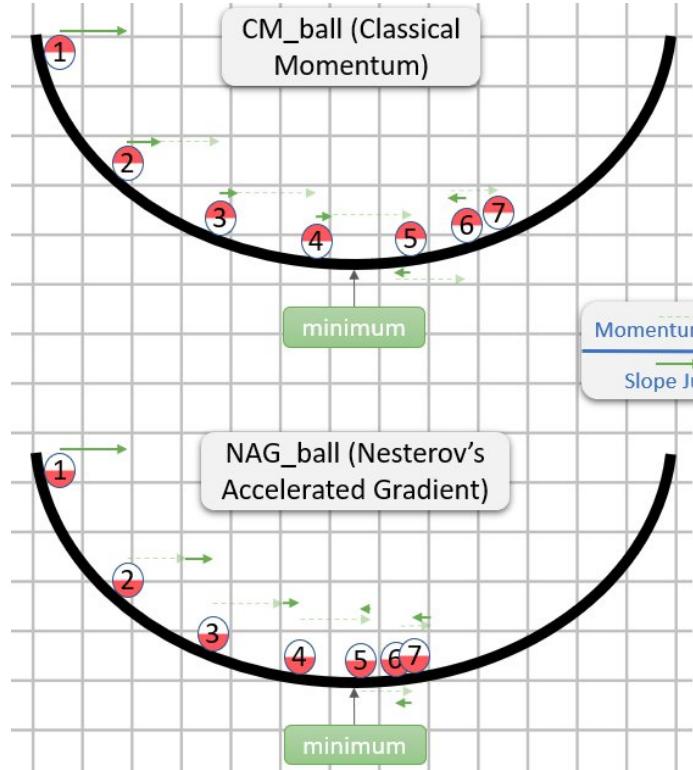
**Nesterov momentum**



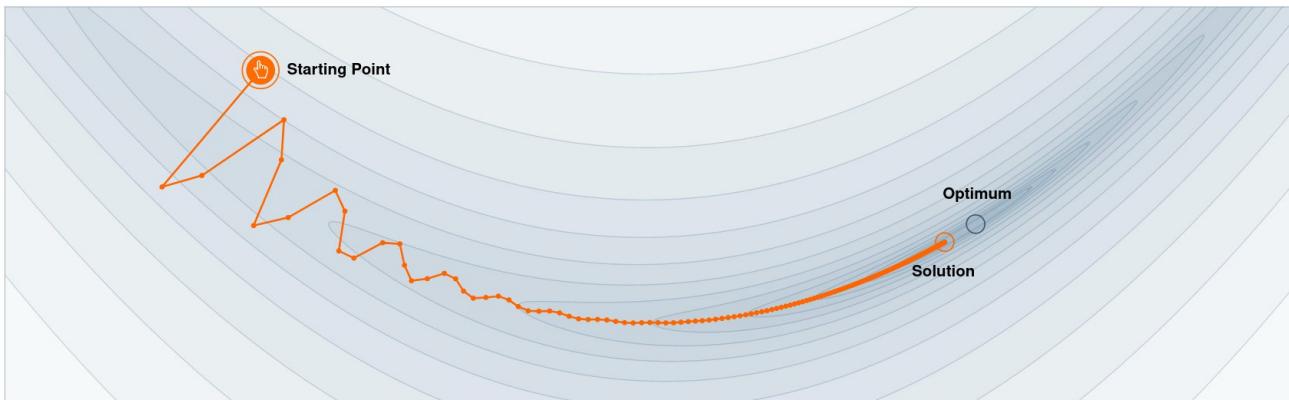
Nesterov momentum update



CM\_ball (Classical Momentum)



# Why Momentum Works ?



Step-size  $\alpha = 0.02$



Momentum  $\beta = 0.99$



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

# Adaptive Optimizers

Rather than controlling the gradient descent manually with the learning rate...

... We can adapt the *learning rate* for each weight of the model according to the **gradient**, the **gradient2**, or the **norm of the weights** of the layer!!!

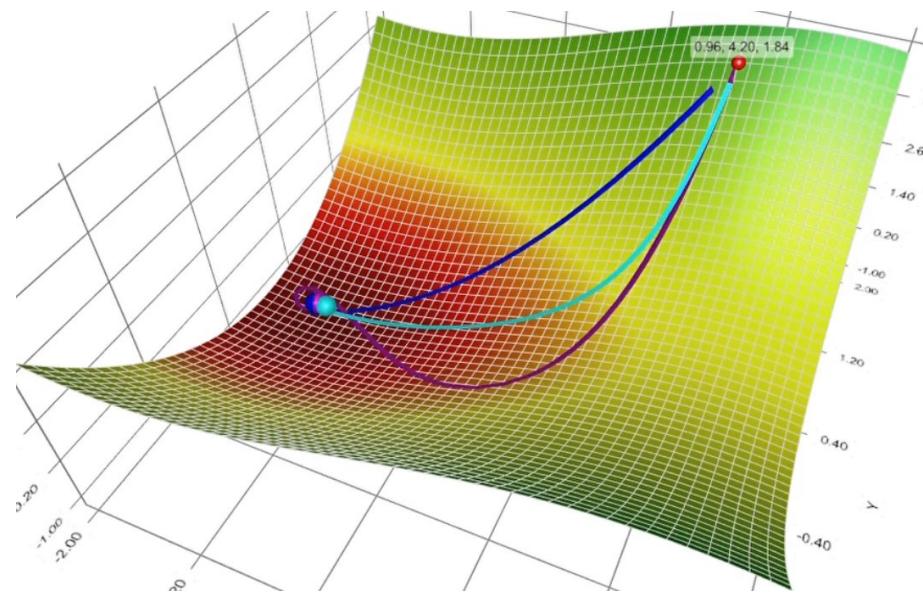
Examples :

- AdaGrad,
- AdaDelta,
- RMSprop
- Adam

Specialized for large batches :

- LARS
- LAMB

- SGD (no *momentum*)
- SGD (with *momentum*)
- Adam



# Adam

$$m_i = \beta_1 * m_{i-1} + (1 - \beta_1) * g_i$$

$$v_i = \beta_2 * v_{i-1} + (1 - \beta_2) * g_i^2$$

$$\hat{m}_i = \frac{m_i}{1 - \beta_1^i}$$

$$\hat{v}_i = \frac{v_i}{1 - \beta_2^i}$$

$$\theta_i = \theta_{i-1} - \frac{\alpha}{\sqrt{\hat{v}_i + \epsilon}} * \hat{m}_i$$

Parameters:

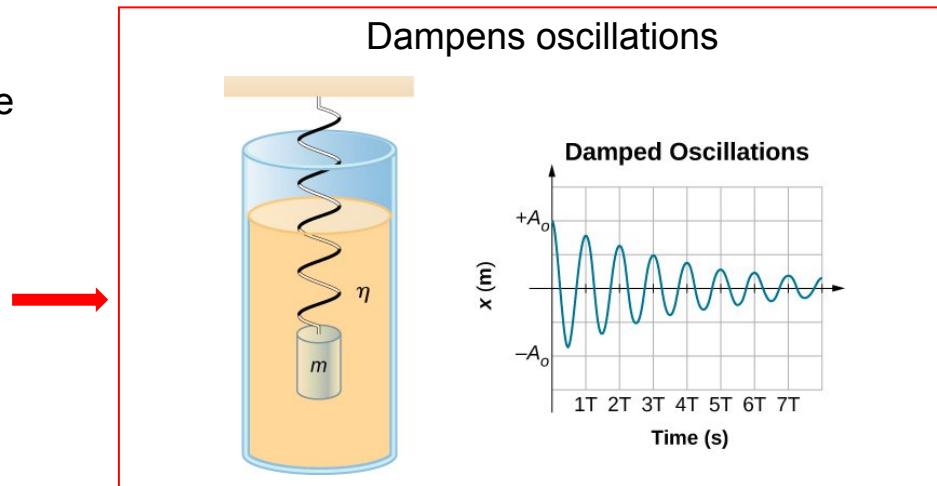
$\beta_1$  &  $\beta_2$  = Regression rate ( $\beta_1 = 0.9$  &  $\beta_2 = 0.999$ )

$\epsilon$  = Very small value to avoid division by zero

Adam : Adaptative moment estimation

First moment : sliding mean

Second moment : sliding non-centered variance

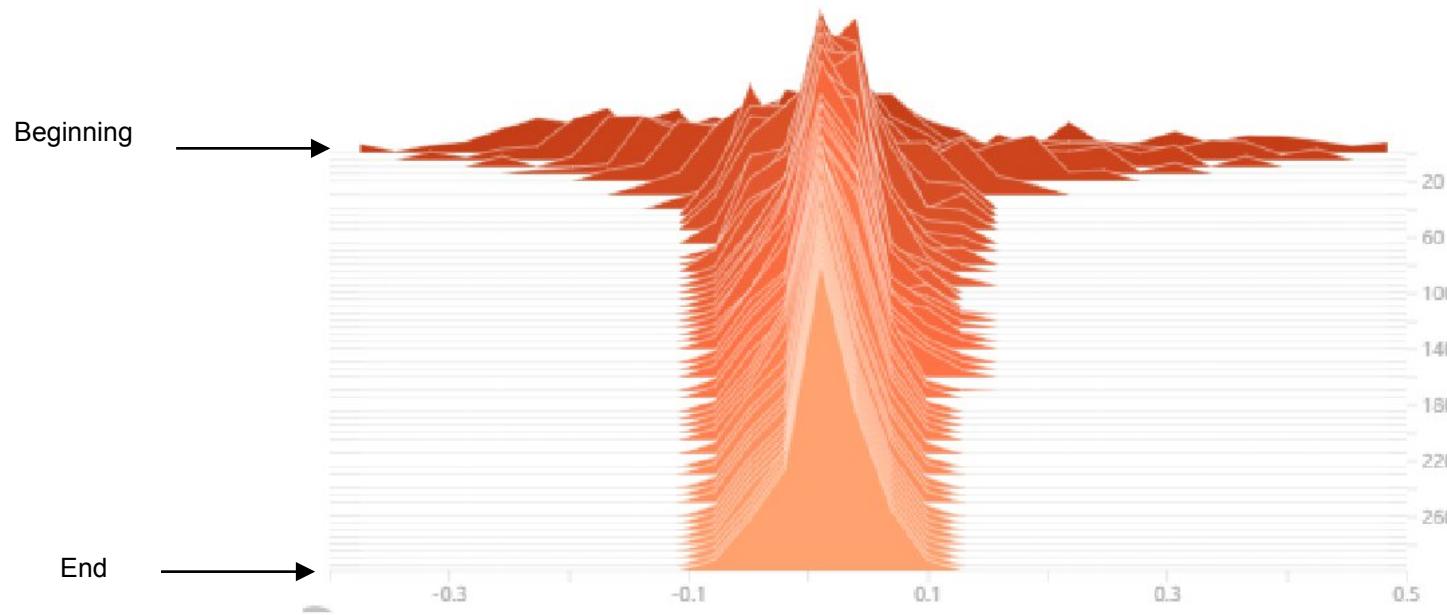


Goal: **Adapt the importance of weight** updates based on previous gradients and gradient variability.

# Weight decay

A neural network that **converges and generalizes correctly\*** generally has weights that tend to **0**. \*(neither underfitting nor overfitting)

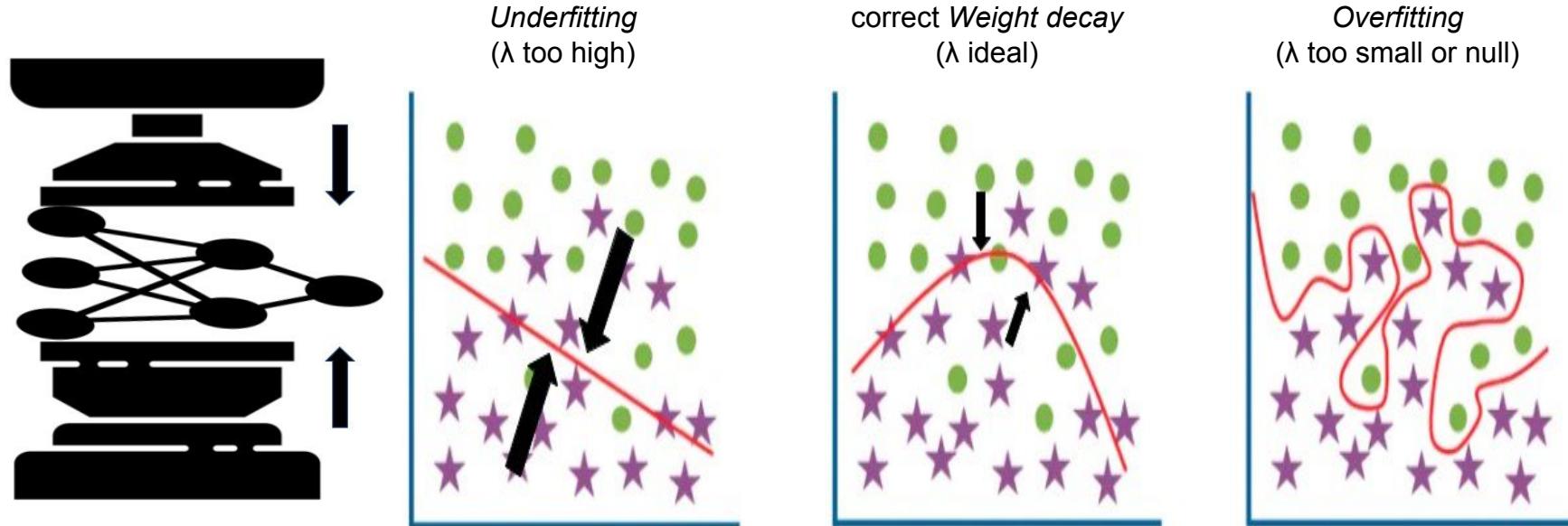
Distribution of weights during learning:



# Weight decay

Preferable to standard L2 regularization defined in loss function

$\lambda$  : weight decay parameter (between 0 and 0.1)



The weight decay technique, defined in the optimizer, makes it possible to force the weights to converge towards values close to zero.

# Weight decay and decoupled weight decay

ADAM

```
For i = 1 to ...
   $g_i = \nabla_{\theta} f_i(\theta_{i-1}) + \lambda \theta_{i-1}$ 
   $m_i = \beta_1 * m_{i-1} + (1 - \beta_1) * g_i$ 
   $v_i = \beta_2 * v_{i-1} + (1 - \beta_2) * g_i^2$ 
   $\hat{m}_i = \frac{m_i}{1 - \beta_1^i}$ 
   $\hat{v}_i = \frac{v_i}{1 - \beta_2^i}$ 
   $\theta_i = \theta_{i-1} - \frac{\alpha}{\sqrt{\hat{v}_i + \epsilon}} * \hat{m}_i$ 
Return  $\theta_i$ 
```

Weight decay

ADAMW

```
For i = 1 to ...
   $g_i = \nabla_{\theta} f_i(\theta_{i-1})$ 
   $m_i = \beta_1 * m_{i-1} + (1 - \beta_1) * g_i$ 
   $v_i = \beta_2 * v_{i-1} + (1 - \beta_2) * g_i^2$ 
   $\hat{m}_i = \frac{m_i}{1 - \beta_1^i}$ 
   $\hat{v}_i = \frac{v_i}{1 - \beta_2^i}$ 
   $\theta_i = \theta_{i-1} - \frac{\alpha}{\sqrt{\hat{v}_i + \epsilon}} * \hat{m}_i - \alpha \lambda \theta_{i-1}$ 
Return  $\theta_i$ 
```

Decoupled weight decay

Evolution of weight decay:

*Decoupled weight decay* (decoupled from momentum!!)

- SGD and Adam with weight decay
- SGDW and AdamW with decoupled weight decay

**SGD and SGDW are roughly equivalent in performance.**

However **AdamW is noticeably better than Adam!!**

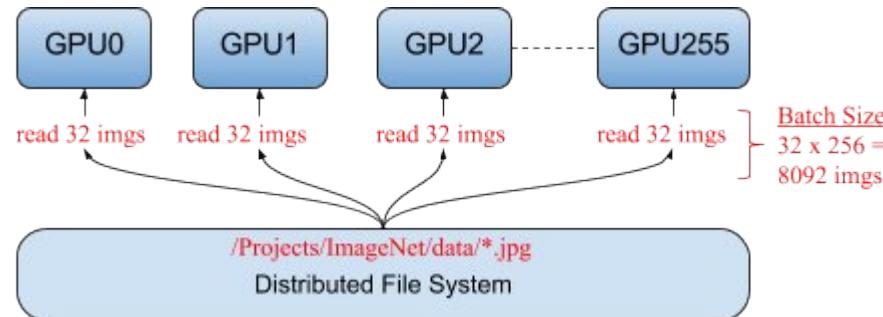
# Optimization of large batches

Large batches issues ◀

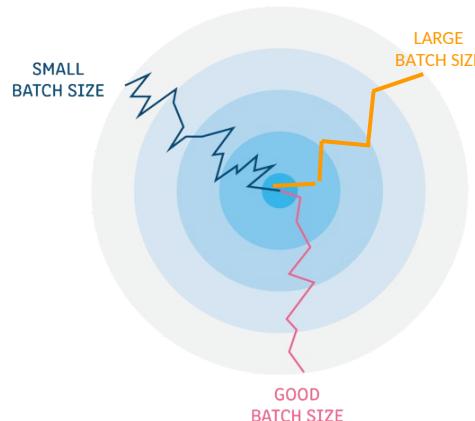
Learning Rate Scaling & Batch Schedulers ◀

Large batches optimizers ◀

# Large Batches with Data Parallelism



**Data Parallelism:** This parallelism generates very large batches



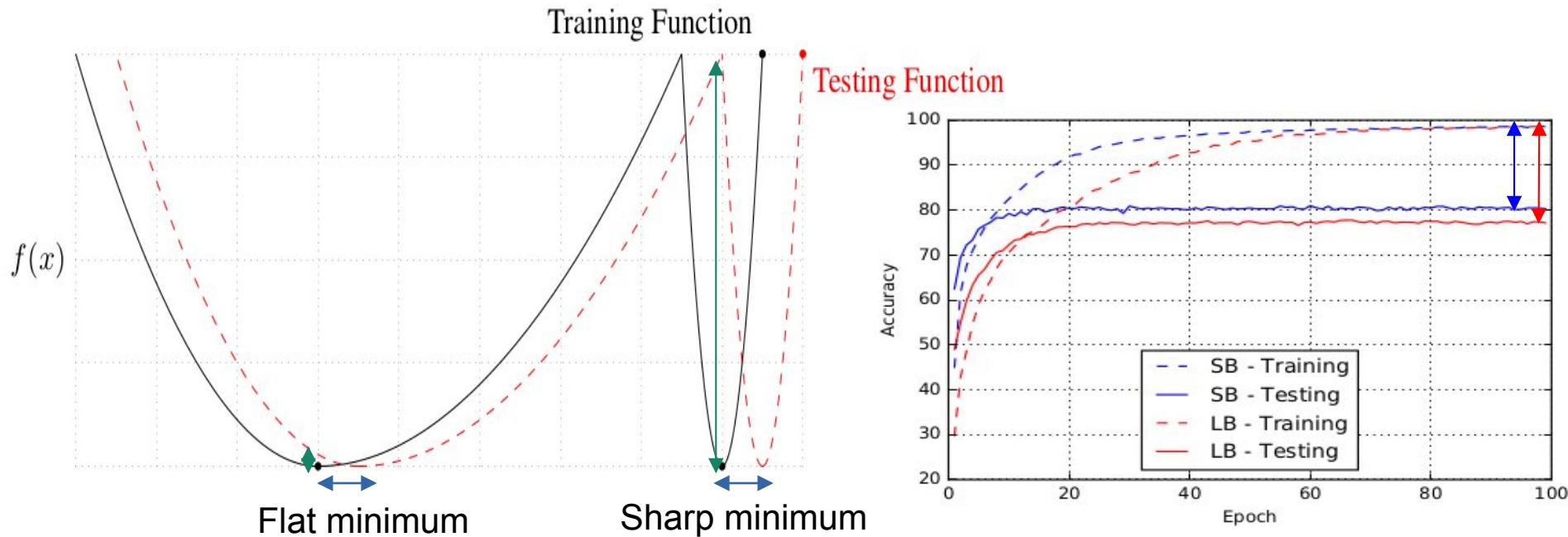
Problem: *Batch* that are too large (> 512) tend to result in **poorer performance**



# Large Batches

On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima

Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, Ping Tak Peter Tang



Comparison of training a convolutional network with small batch (SB) and large batch (LB) on CIFAR 10

The *larger the batch*, the more the model tends to converge towards **steep and narrow minima**.

# Large Batches : Learning rate scaling

When the *size of the overall batch is considerably increased*, it is often *necessary to scale the learning rate*:

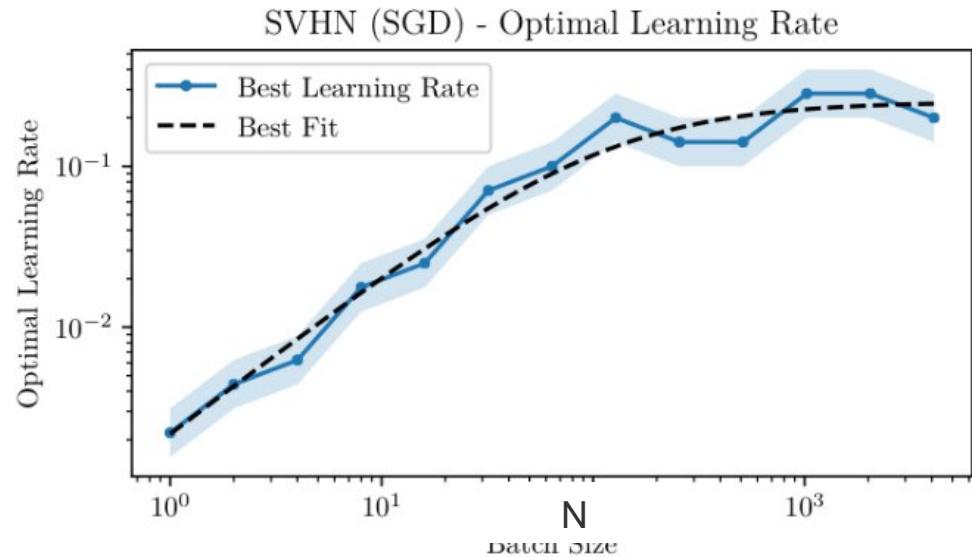
$N$  = Number of parallel processes

**Linear growth of learning rate:**

$$\alpha \rightarrow N * \alpha$$

**Square root growth of learning rate:**

$$\alpha \rightarrow \sqrt{N} * \alpha$$



*An Empirical Model of Large-Batch Training*  
Sam McCandlish, Jared Kaplan, Dario Amodei

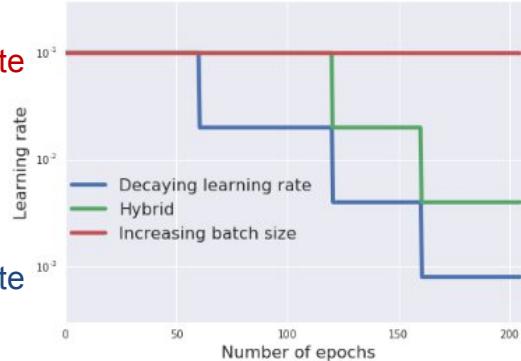
Optimal: **linear growth** at first then **square root**  
(recommended by OpenAI)

# Batch Size Scheduler

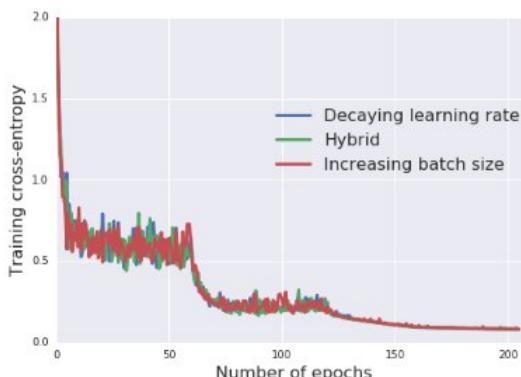
=> Alternative to Learning Rate Scheduler

DON'T DECAY THE LEARNING RATE, INCREASE THE BATCH SIZE

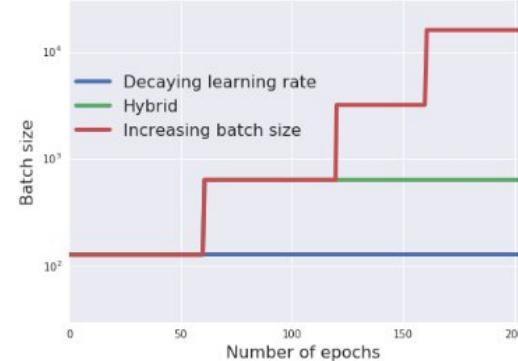
High Learning Rate



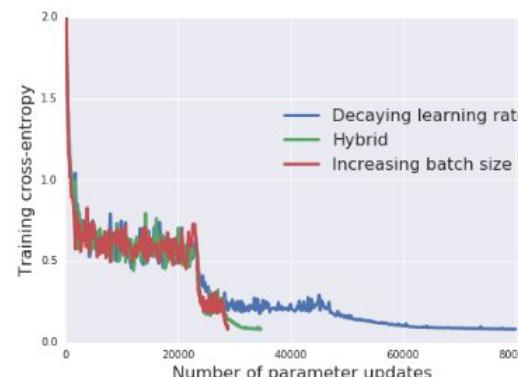
Low Learning Rate



Large Batch



Small Batch



# Large Batches

Trends :

**Flat Minimum** | **Sharp Minimum**

- Test Loss

+ Test Loss

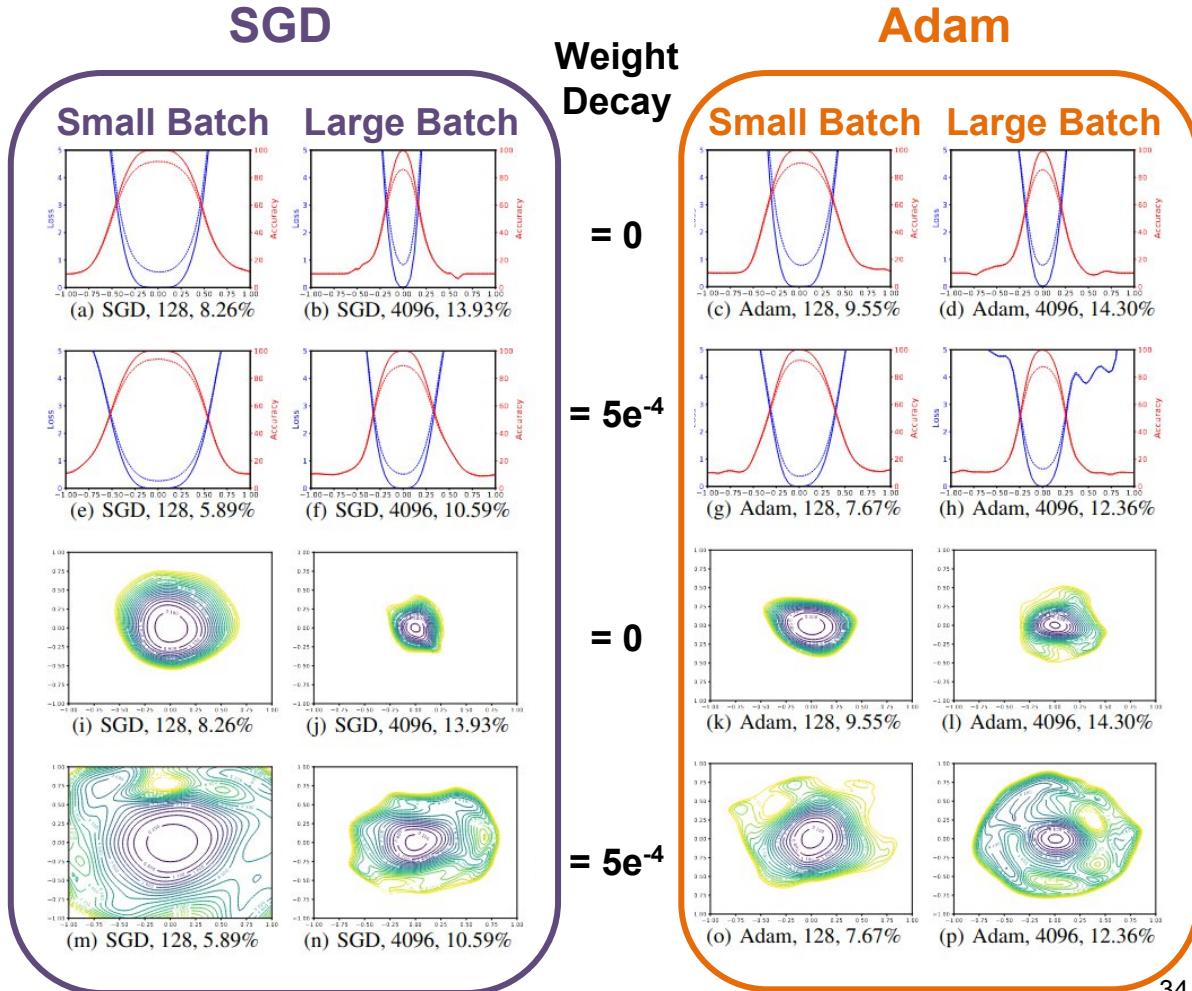
Slow Descent

Fast Descent

Small Batch Large Batch

SGD ADAM

Weight Decay w/o W.Decay



# Large Batches Optimisers - LARS

LARS = “Layer-wise Adaptive Rate Scaling.”

Adaptation of SGD with momentum with the addition of a **trust ratio** for each layer which depends on the evolution of the layer's gradient

$r$  = Trust ratio

$l$  = layer number

$$m_i = \beta * m_{i-1} + (1 - \beta) * (g_i + \lambda \theta_{i-1})$$

$$r_1 = \|\theta_{i-1}^l\|_2$$

$$r_2 = \|m_i^l\|_2$$

$$r = r_1 / r_2$$

$$\alpha^l = r * \alpha$$

$$\theta_i^l = \theta_{i-1}^l - \alpha^l * m_i^l$$

Weight decay

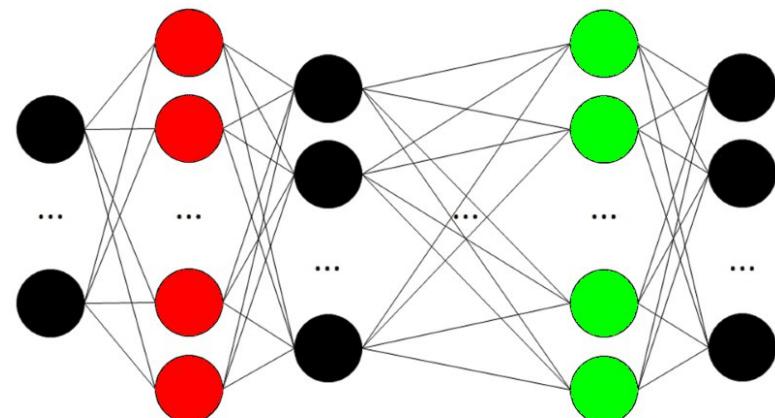


Low confidence

Low value of  $r$

High confidence

High value of  $r$



Goal: Adapt the importance of weight updates based on a **trust ratio** calculated for each layer of the network.

# Large Batches Optimisers - LAMB

LAMB pour “Layer-wise Adaptive Moments optimizer for Batch training.”

Adaptation of ADAM with momentum with the addition of a **trust ratio** for **each layer** which depends on the evolution of the layer's gradient

$r$  = Trust ratio

$l$  = layer number

$$r_1 = \|\theta_{i-1}^l\|_2$$

$$r_2 = \left\| \frac{\hat{m}_i^l}{\sqrt{\hat{v}_i^l + \epsilon}} + \underline{\lambda \theta_{i-1}^l} \right\|_2$$

$$r = r_1 / r_2$$

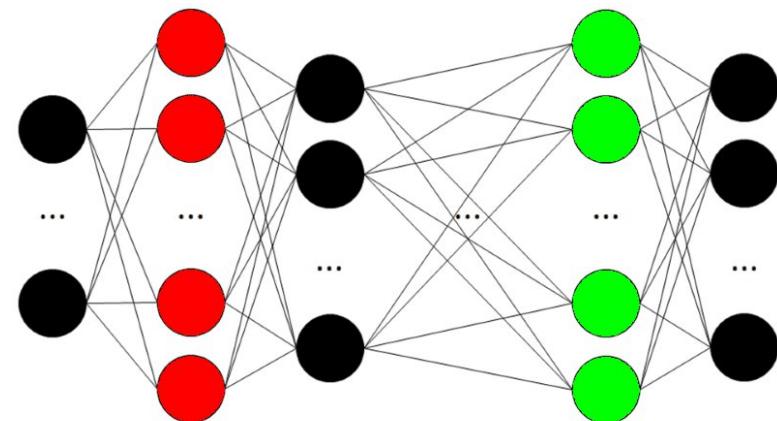
$$\alpha^l = r * \alpha$$

$$\theta_i^l = \theta_{i-1}^l - \alpha^l * \left( \frac{\hat{m}_i^l}{\sqrt{\hat{v}_i^l + \epsilon}} + \underline{\lambda \theta_{i-1}^l} \right)$$

Decoupled weight decay

Low confidence  
Low value of  $r$

High confidence  
High value of  $r$



Goal: Adapt the importance of weight updates based on a **trust ratio** calculated for each layer of the network.

# Optimizer implementation

SGD

```
import torch.optim as opt  
  
SGD_optimizer = opt.SGD(params, lr, momentum=0, weight_decay=0, nesterov=False, ...)
```

ADAMW

```
import torch.optim as opt  
  
ADAM_optimizer = opt.AdamW(params, lr=0.001, betas=(0.9, 0.999), weight_decay=0.05,...)
```

LAMB

```
from apex.optimizers import FusedLamb  
  
LAMB_optimizer = FusedLamb(params, lr=0.001, betas=(0.9, 0.999), weight_decay=0,  
adam_w_mode=True)
```

LARC

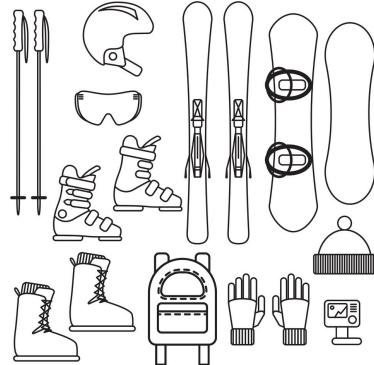
LARS  
optimisation  
from APEX

```
import torch.optim as opt  
from apex.parallel.LARC import LARC  
  
base_optimizer = opt.SGD(params, lr=0.001, momentum=0.9, weight_decay=0)  
optimizer = LARC(base_optimizer)  
scheduler = opt.lr_scheduler.CyclicLR(base_optimizer, base_lr=0.01, max_lr=0.1)
```

# *Large Batches Rider*

Weight Decay

SGD      AdamW



Batch Scheduler

LARS  
LAMB

LR Scheduler

Warmup      LR scaling      LR Decay



# BLOOM example

95281 steps (116.8 days)

AdamW,  
 $\beta_1=0.9$ ,  $\beta_2=0.95$ ,  $\text{eps}=1e-8$

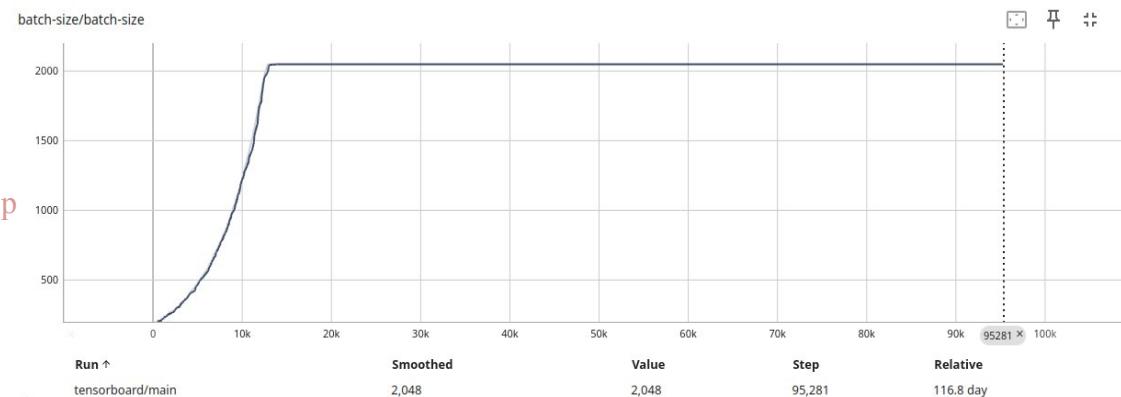
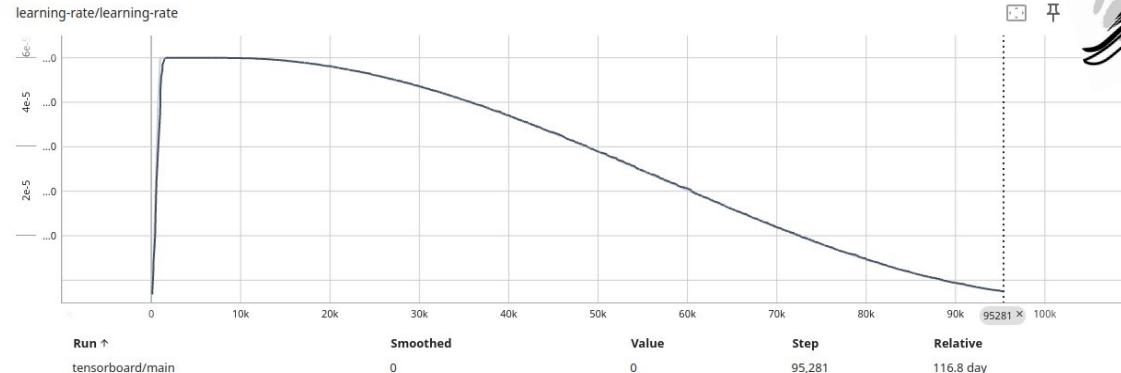
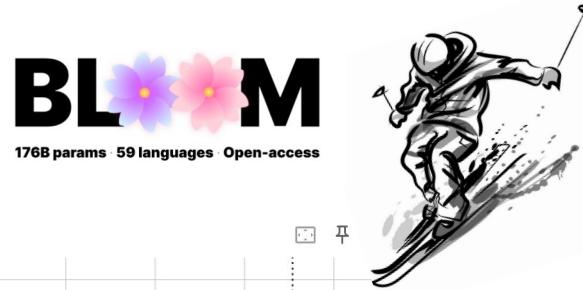
Weight Decay of 0.1

## LR Scheduler

- peak=6e-5
- warmup over 375M tokens
- cosine decay for learning rate down to 10% of its value, over 410B tokens

## Batch Scheduler

- start from 32k tokens (GBS=16)
- increase linearly to 4.2M tokens/step (GBS=2048) over ~20B tokens
- then continue at 4.2M tokens/step (GBS=2048) for 430B tokens



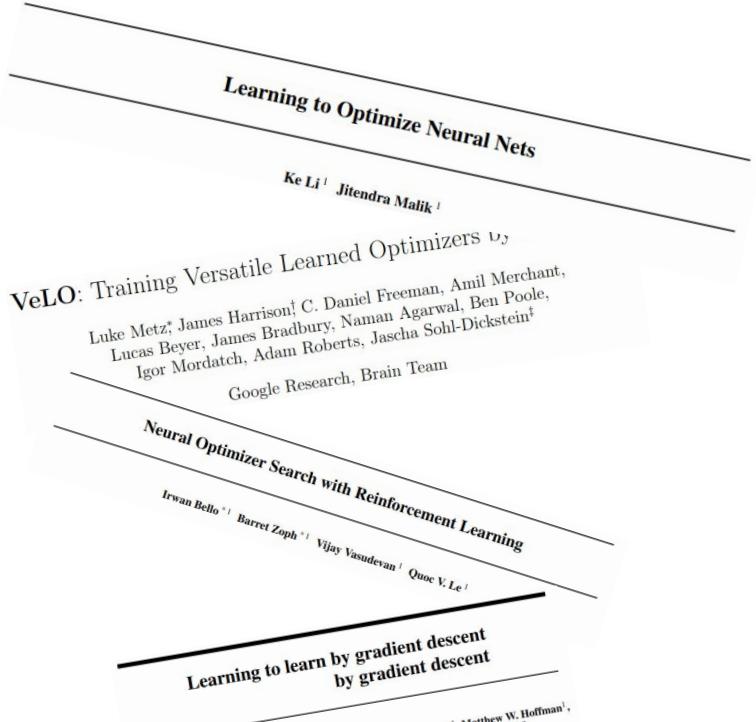
# New optimizers

New trend : optimizers learning ◀

New optimizers Abyss ◀

LION : example of a new approach ◀

# New trend : optimizers learning



Marcin Andrychowicz<sup>1</sup>, Misha Denil<sup>1</sup>, Sergio Gomez Colmenarejo<sup>1</sup>, Matthew W. Hoffman<sup>1</sup>,  
David Pfau<sup>1</sup>, Tom Schaul<sup>1</sup>, Brendan Shillingford<sup>1,2</sup>, Nando de Freitas<sup>1,2,3</sup>

<sup>1</sup>Google DeepMind    <sup>2</sup>University of Oxford    <sup>3</sup>Canadian Institute for Advanced Research

marcin.andrychowicz@gmail.com    {denil, sergomez, mwhoffman, pfau, schaul}@google.com

brendan.shillingford@cs.ox.ac.uk    nando.defreitas@google.com

# New optimizers Abyss

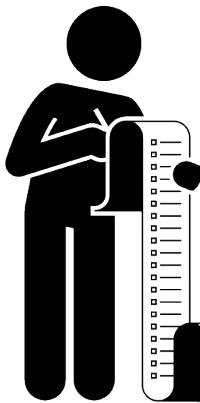
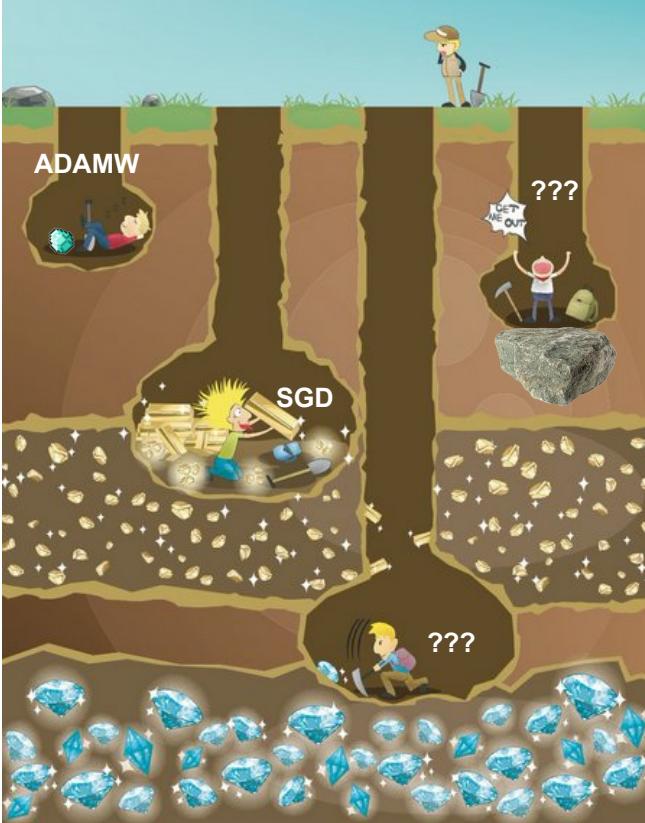


Table 2: List of optimizers considered for our benchmark. This is only a subset of all existing methods for deep learning.

Name	Ref.	Name	Ref.
AccGrad	(Levy et al., 2018)	HypGN	(Wang et al., 2018a)
ACCGD	(Zhang et al., 2020)	K-BFGS-BFGS(L)	(Gaidhane et al., 2020)
ADAM	(Kingma & Ba, 2014)	K-FQN-CNN	(Bridle et al., 2020)
Adadelta	(Diederik et al., 2017)	KFSG	(Molchan & Grosse, 2015)
Adadelta-Keller-Gill-SS	(Zhang et al., 2020)	L2Adam-PiNewton	(Belotti & Mazzanti, 2015)
AdadeltaK	(Vas et al., 2019)	LAAdam	(You et al., 2020)
Adadelta	(Gardner et al., 2018)	LAM	(Gardner et al., 2018)
Adadelta	(Chen et al., 2018)	LADPT	(You et al., 2017)
Adadelta	(Zhang et al., 2018)	LEAD	(You et al., 2016)
Adadelta	(Shazeer & Stern, 2018)	LockAdam	(Zhang et al., 2019)
Adadelta	(Zhang et al., 2019a)	MADGRAD	(Belotti & Jerez, 2021)
Adadelta	(Hochreiter & Schmidhuber, 2001)	MAOPT	(Dhariwal & Neelakantan, 2017)
Adadelta	(Duchi et al., 2011)	MEKA	(Chen et al., 2020b)
AdadeltaMAN	(Yuan et al., 2019)	MGR	(Hochreiter & Schmidhuber, 2001)
Adadelta	(Ou et al., 2020)	MVR-C(MVR-C-2)	(Chen & Zhou, 2020)
Adam	(Kingma & Ba, 2014)	NeA	(You et al., 2016)
Adam	(Tramer et al., 2017)	NAMSINAMSG	(Chen et al., 2019b)
Adam	(Kingma & Ba, 2015)	NDS-Adam	(Zhang et al., 2017)
Adam	(Lee et al., 2018)	NG-Adam	(You et al., 2017)
Adam	(Dong et al., 2017)	Nesterov	(Nesterov, 1983)
AdamHS	(Kingma & Ba, 2015)	Noisy AdamNoisy K-FAC	(Zhang et al., 2018)
AdamNC	(Kingma & Ba, 2014)	Novaladam	(Huang et al., 2019)
Adams	(Dong et al., 2018)	Nomadam	(Goh et al., 2019)
AdamPGP	(Ihler et al., 2021)	NT-SGD	(Zheng et al., 2020)
Adams	(Goh et al., 2020)	OASIS	(Huang et al., 2020)
AdamW	(Loshuk & Hutter, 2019)	PAGL	(Li et al., 2020b)
Adams	(Tolstikhin & Jegorov, 2020)	PAL	(Hochreiter & Schmidhuber, 2001)
ADAS	(Dziugaite et al., 2020)	Poly-Adam	(Orubisa et al., 2019)
Adasgd	(Hochreiter & Schmidhuber, 2001)	PowerDI/PowerSGD	(Vogelzang et al., 2019)
Adasgd	(Ioffe et al., 2015)	Probabilistic Poyal	(de Ros et al., 2021)
Adasgd	(Ogden & Wainwright, 2020)	Pegas	(Oehmichen et al., 2021)
Adasgd	(Ha et al., 2019)	Platonic	(Mu & Wang, 2020)
Adasgd	(Ho et al., 2019)	PLAdam/QLM	(Rabbani et al., 2019)
Adasgd	(Li et al., 2020a)	RAdam	(Li et al., 2020a)
ADGD	(Li et al., 2019)	Radam	(Olah et al., 2019)
ADMGd	(Bermúdez et al., 2020)	RangeLars	(Gaudie, 2020)
AMSGrad	(Rakhlin et al., 2017)	RankProp	(Tolstikhin et al., 2021)
AMSGrad	(Rakhlin et al., 2018)	RBMSnew	(Chen et al., 2019)
ANASG	(Wen et al., 2019)	REED	(Wang et al., 2020b)
ANASG	(Norouzi et al., 2019)	Safada	(Huang et al., 2020)
ANASG	(Chen et al., 2019)	Safada-MMSGd	(You et al., 2020)
ANASG	(Koen et al., 2021)	SALR	(Frenet et al., 2017)
ANASG	(Dai et al., 2021)	SC-Adagrad/SC-RMSProp	(Makarenko et al., 2017)
BFGS	(Schraudolph et al., 2019)	SDProp	(Ihsen et al., 2017)
BFGS	(Huang et al., 2017b)	SDF	(Robbins & Monro, 1951)
BFGSProp	(Huang et al., 2017b)	SGD-BB	(Chen et al., 2016)
CADAM	(Huang et al., 2017b)	SGLD	(Goh et al., 2020)
CADAM	(Huang et al., 2017b)	SGDGM	(Huang et al., 2020)
CADA	(Huang et al., 2017b)	SGBM	(Lin & Lee, 2020)
CGOptimizer	(Brooks et al., 2019)	SGBM	(Huang et al., 2020)
CGOptimizer	(Pereyra & Riquelme, 2019)	SGBMDgrad	(Huang et al., 2020)
CGOptimizer	(Huang et al., 2017b)	SGD	(Huang et al., 2020)
CGOptimizer	(Huang et al., 2017b)	SGD	(Huang et al., 2020)
Dense	(Ng et al., 2019)	Sgdm	(Huang et al., 2020)
Descent	(Ng et al., 2019)	Sgdm++	(Wang et al., 2019a)
DiCoProp	(Dai et al., 2021a)	SGQN-SQQN	(Huang et al., 2020)
DiCoProp	(Dai et al., 2021a)	SQ	(Wang et al., 2020)
DIFD	(Dai et al., 2020)	SQG	(Adu et al., 2019)
DIFD	(Dai et al., 2020)	SQM	(Huang et al., 2020)
DIFD	(Dai et al., 2020)	SQSM	(Zhou et al., 2020)
DKM	(Dai et al., 2018)	SqSmooth	(Huang et al., 2020)
EKFDC	(Huang et al., 2017a)	SGSgd	(Wang et al., 2020)
EKFDC	(Daly & Arora, 2019)	SGSGD	(Wang et al., 2020)
Exponential	(Dai et al., 2020)	SGSGD	(Wang et al., 2020)
Exponential	(Dai et al., 2020)	SGSGD SGD	(Wang et al., 2020)
Exponential	(Wang & Ye, 2020)	SGWAT	(Berak & Sclar, 2017)
Exponential	(Chukwuma et al., 2018)	SGXNS	(Ihsen et al., 2017)
Zhang & Gama, 2018)	(Zhang & Gama, 2018)	Tahar	(Ihsen et al., 2020)
Exponential	(Goh et al., 2020)	TMGR	(Ihsen et al., 2020)
Exponential	(Chen et al., 2021)	Vadara	(Xuan et al., 2019)
Exponential	(Rozsas et al., 2019)	VGR	(Shochet et al., 2015)
Exponential	(Rozsas et al., 2020)	ViGR	(Shochet et al., 2015)
Feverate	(Feverate et al., 2020)	VS-GRD	(Dai et al., 2019)
Feverate	(Feverate et al., 2020)	WGR	(Wang et al., 2019)
Fisher	(Bottou & Zhang, 2019)	YellowFin	(Zhang & Miličević, 2019)
Fisher	(Zhang & Miličević, 2019)	Yogi	(Zhang et al., 2019)

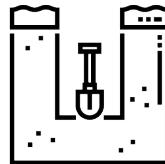
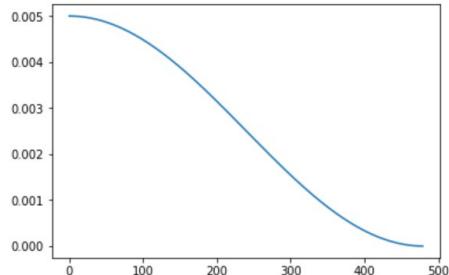
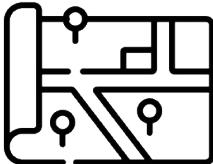


Schmidt, Robin M., Frank Schneider, and Philipp Hennig. "Descending through a crowded valley-benchmarking deep learning optimizers." *International Conference on Machine Learning*. PMLR, 2021.

# LION : example of a new approach

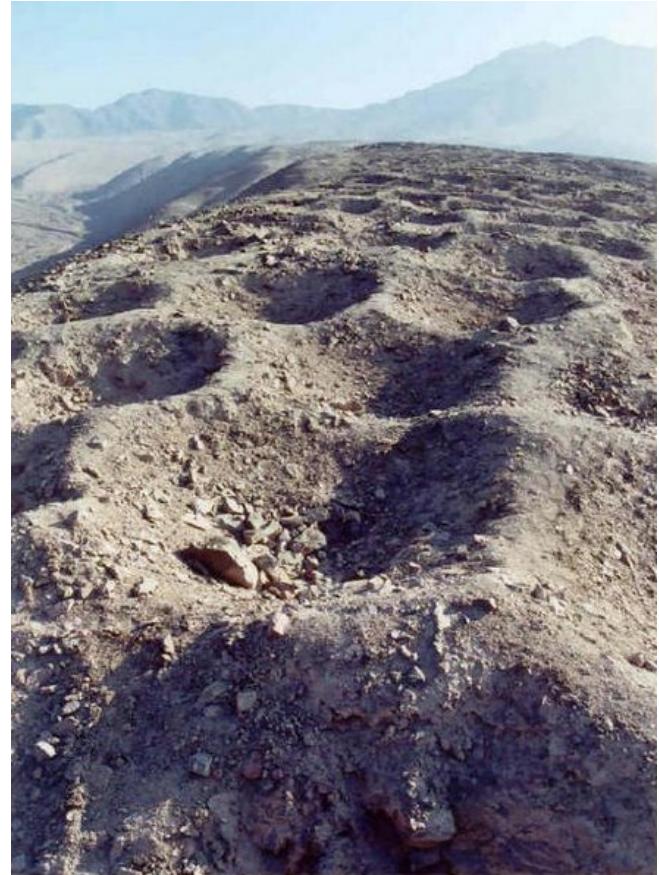
## Algorithm 1 AdamW Optimizer

```
given  $\beta_1, \beta_2, \epsilon, \lambda, \eta, f$ 
initialize  $\theta_0, m_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$ 
while  $\theta_t$  not converged do
     $t \leftarrow t + 1$ 
     $g_t \leftarrow \nabla_{\theta} f(\theta_{t-1})$ 
    update EMA of  $g_t$  and  $g_t^2$ 
     $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$ 
     $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ 
    bias correction
     $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ 
     $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ 
    update model parameters
     $\theta_t \leftarrow \theta_{t-1} - \eta_t (\hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon) + \lambda \theta_{t-1})$ 
end while
return  $\theta_t$ 
```

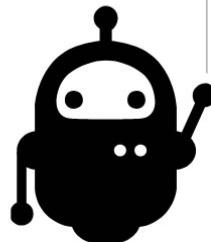


## Algorithm 2 Lion Optimizer (ours)

```
given  $\beta_1, \beta_2, \lambda, \eta, f$ 
initialize  $\theta_0, m_0 \leftarrow 0$ 
while  $\theta_t$  not converged do
     $g_t \leftarrow \nabla_{\theta} f(\theta_{t-1})$ 
    update model parameters
     $c_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$ 
     $\theta_t \leftarrow \theta_{t-1} - \eta_t (\text{sign}(c_t) + \lambda \theta_{t-1})$ 
    update EMA of  $g_t$ 
     $m_t \leftarrow \beta_2 m_{t-1} + (1 - \beta_2) g_t$ 
end while
return  $\theta_t$ 
```



# Pratice : Learning rate + Optimiseurs



Goals :

- Edit the **learning rate scheduler**
- Edit the **optimizer**
- Do training with **large batches**



From **JupyterHub**:

- Launch an interactive instance
- Go to the tp\_optimizers folder
- Open the DLO-JZ\_Optimizers notebook