



# Deep Learning Optimisé - Jean Zay

---

## Entraînement et large batches



INSTITUT DU  
DÉVELOPPEMENT ET DES  
RESSOURCES EN  
INFORMATIQUE  
SCIENTIFIQUE



# Loss Landscape

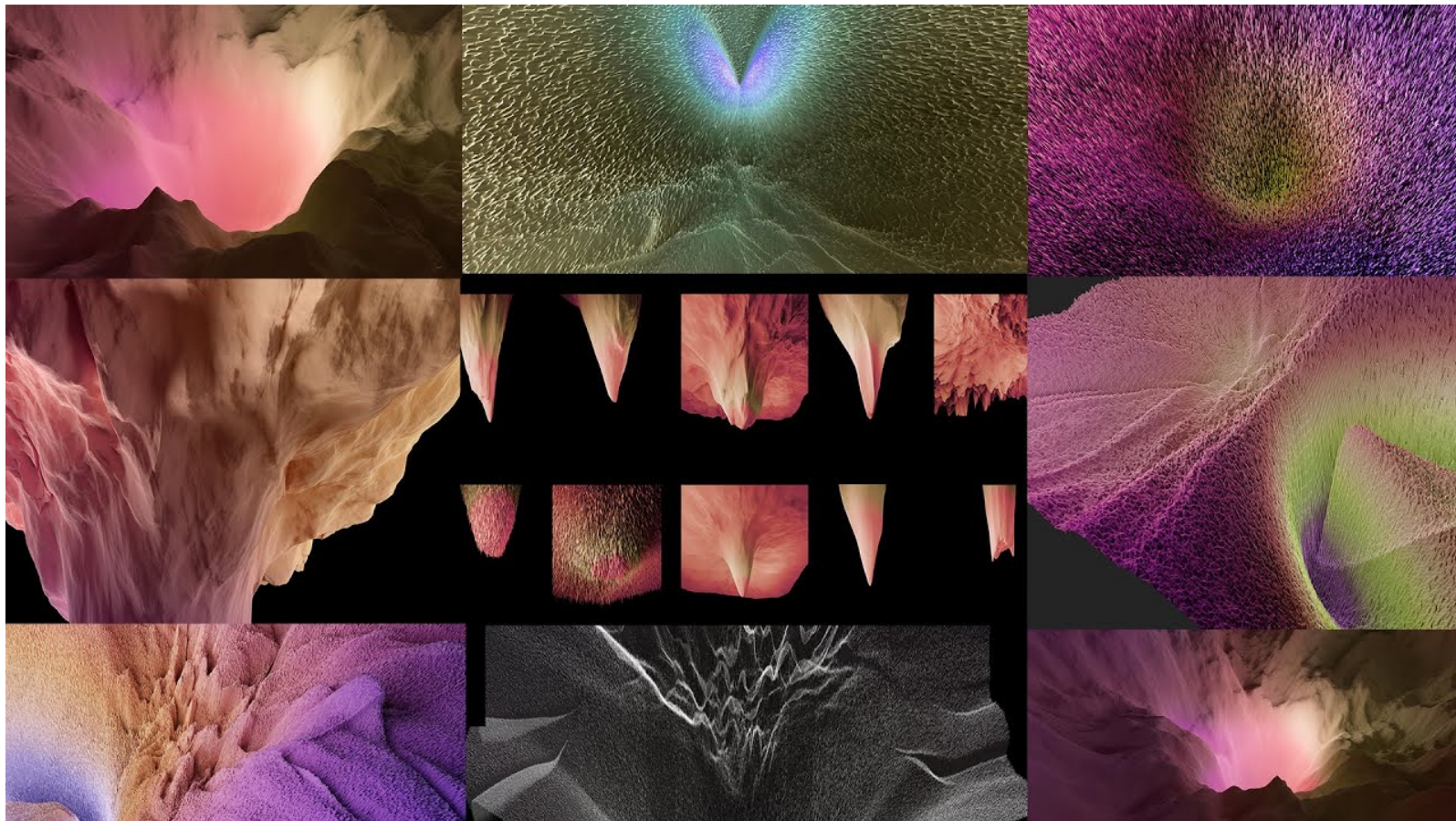
Loss Landscape ◀

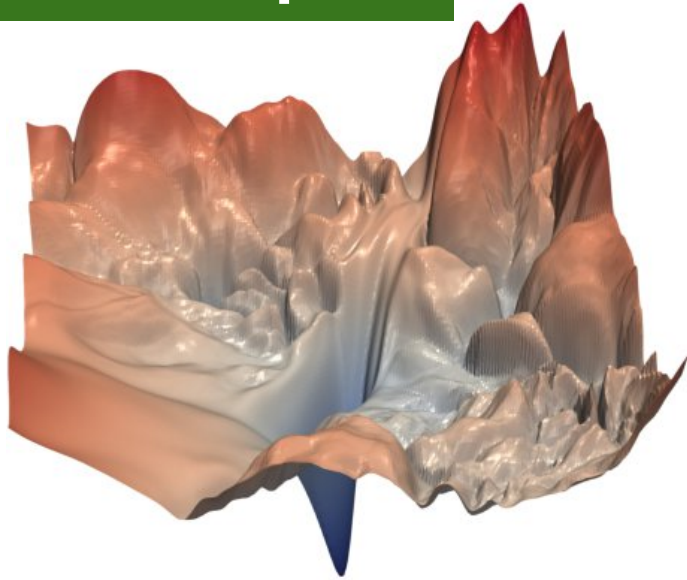
Residual Learning ◀

Initialization ◀

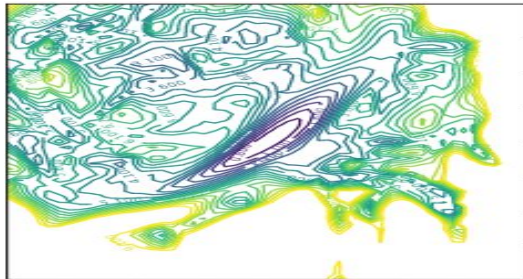
# Loss Landscape

<https://losslandscape.com/>

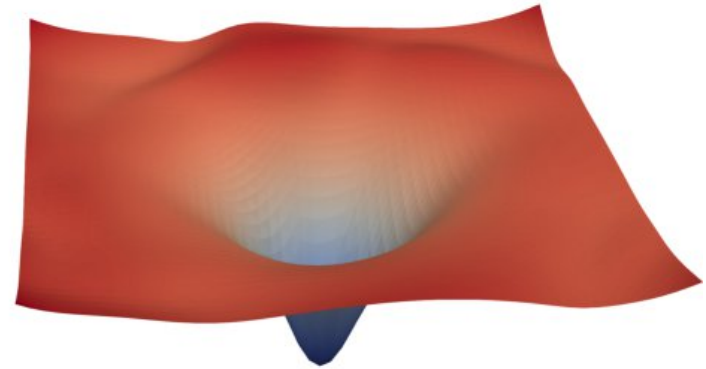




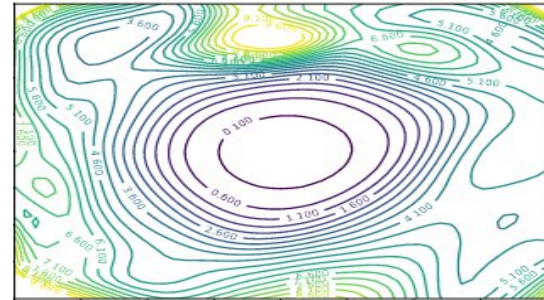
(a) without skip connections



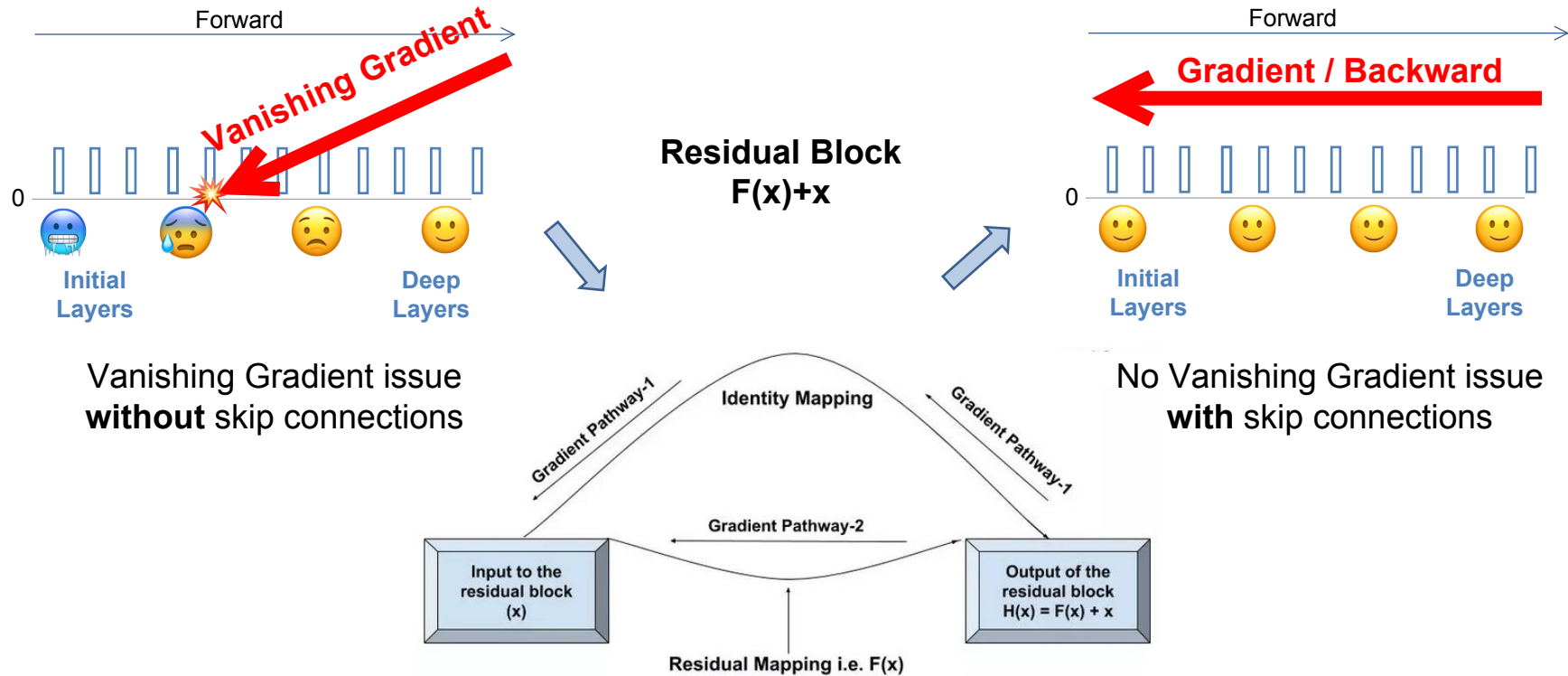
## Residual Learning Depuis les Resnets (2015) ...



(b) with skip connections



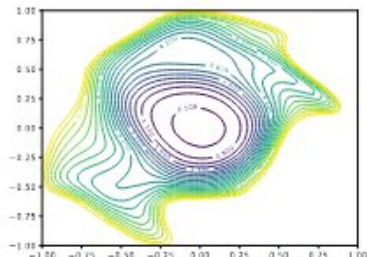
# Residual Learning



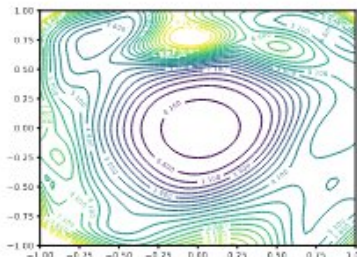
Gradient Pathways in ResNet

# Residual Learning – impact sur la profondeur

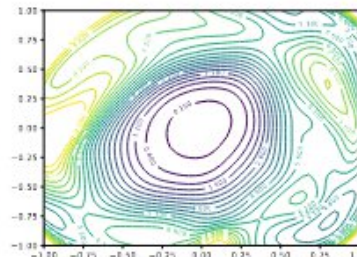
ResNet



(a) ResNet-20, 7.37%

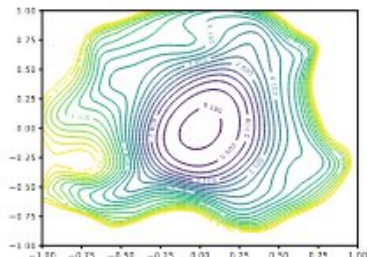


(b) ResNet-56, 5.89%

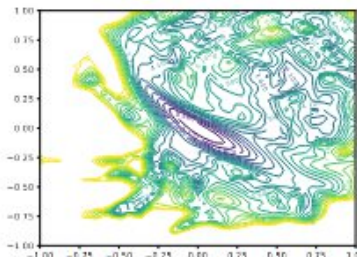


(c) ResNet-110, 5.79%

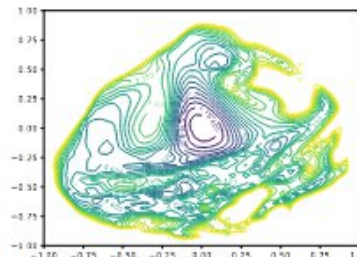
ResNet -  
No Short  
without skip connections



(d) ResNet-20-NS, 8.18%



(e) ResNet-56-NS, 13.31%



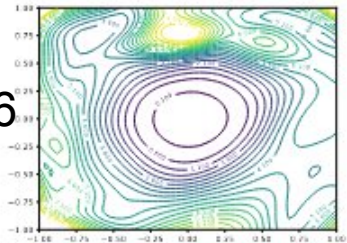
(f) ResNet-110-NS, 16.44%



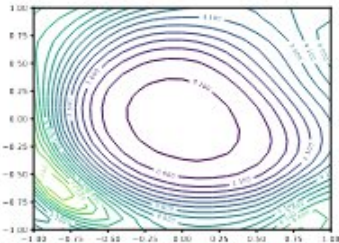
# Residual Learning – impact sur la largeur



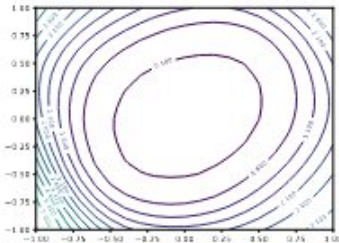
Wide-ResNet-56



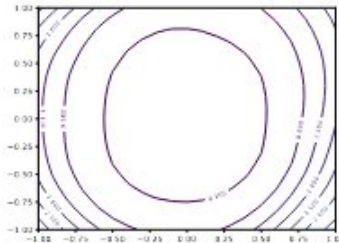
(a)  $k = 1, 5.89\%$



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

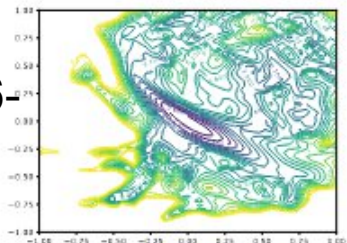


(c)  $k = 4, 4.34\%$

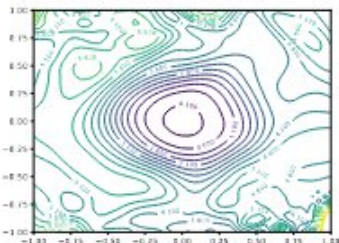


(d)  $k = 8, 3.93\%$

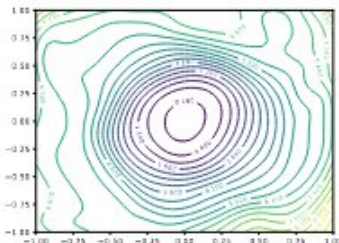
Wide-ResNet-56-  
No Short  
without skip connections



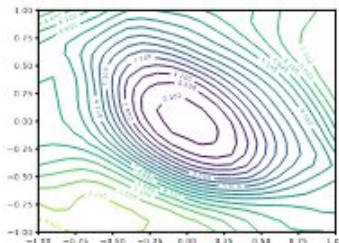
(e)  $k = 1, 13.31\%$



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



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



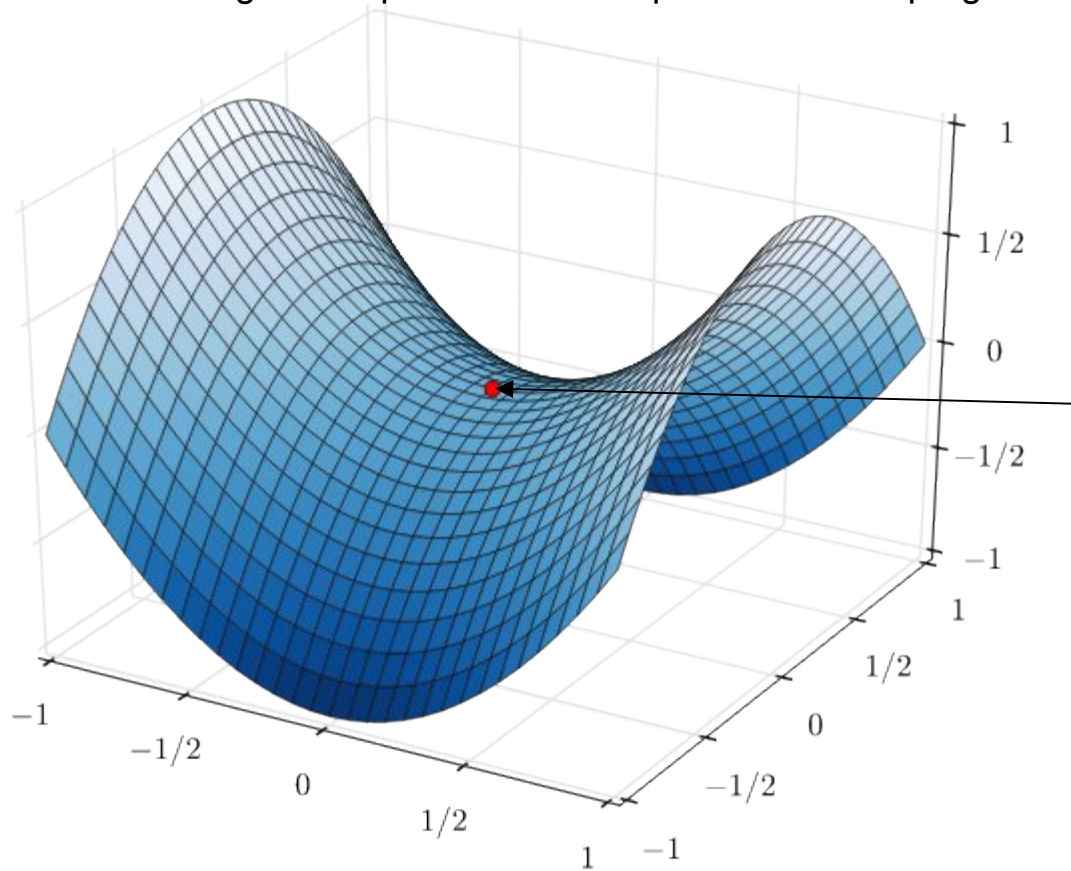
(h)  $k = 8, 8.70\%$



$k$  = coefficient largeur des *channels* par rapport à ResNet

# Problème du point-selle

Point-selle : Un gradient proche de zéro qui va rendre la progression du modèle très lente



Comment quitter le point-selle ?

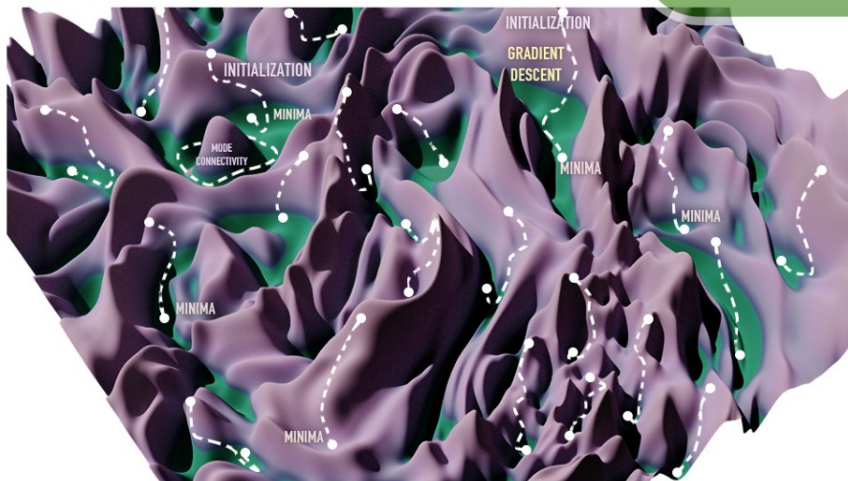


# Initialisation des paramètres du modèle

The Blessing of Dimensionality :

*Local*

NEARBY PATHS  
TO CONVERGENCE



FINDING A MINIMA BECOMES A "LOCAL" CHALLENGE



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

Par défaut dans *PyTorch* :

- Meilleur algorithme d'initialisation selon le type de couche (linéaire, convolutionnel, transformer, ...).
- Aujourd'hui, il n'est plus nécessaire de chercher à optimiser l'initialisation.

# Learning rate scheduler

Learning rate scheduler ◀

    Cyclic scheduler ◀

    One cycle scheduler ◀

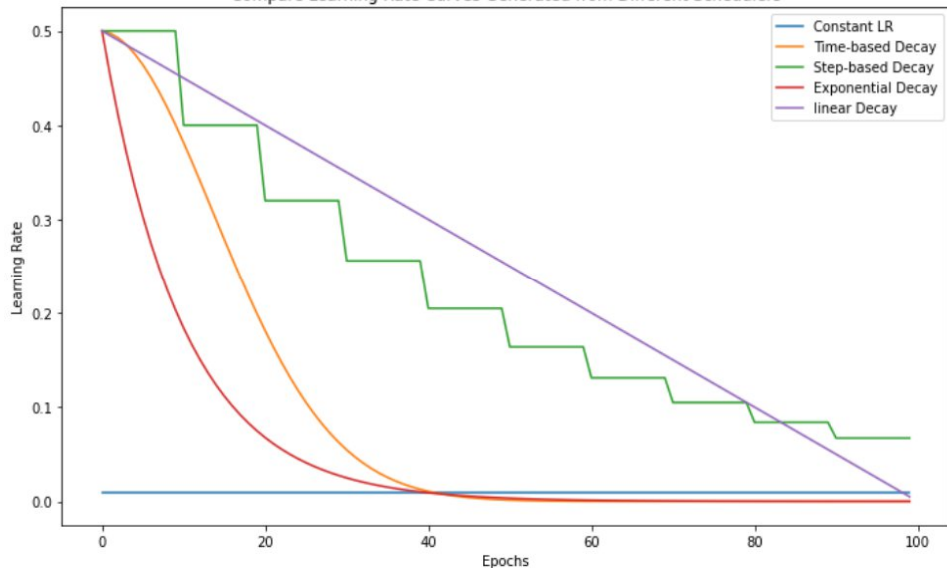
        LR finder ◀

        LR large batch ◀

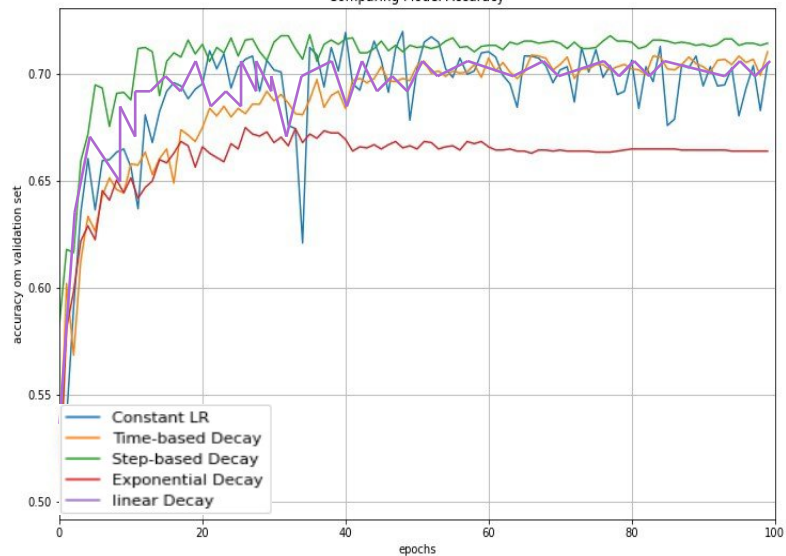
# Learning Rate Scheduler

## Learning rate decay

Compare Learning Rate Curves Generated from Different Schedulers



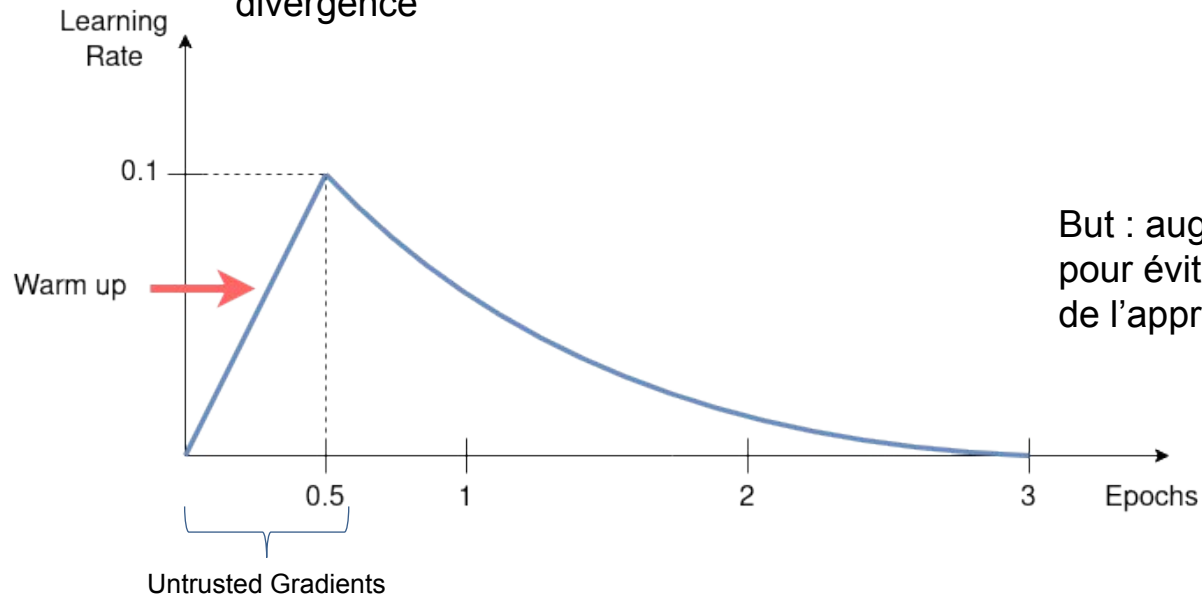
Comparing Model Accuracy



# Learning Rate Scheduler

## WARMUP pour *large batches*

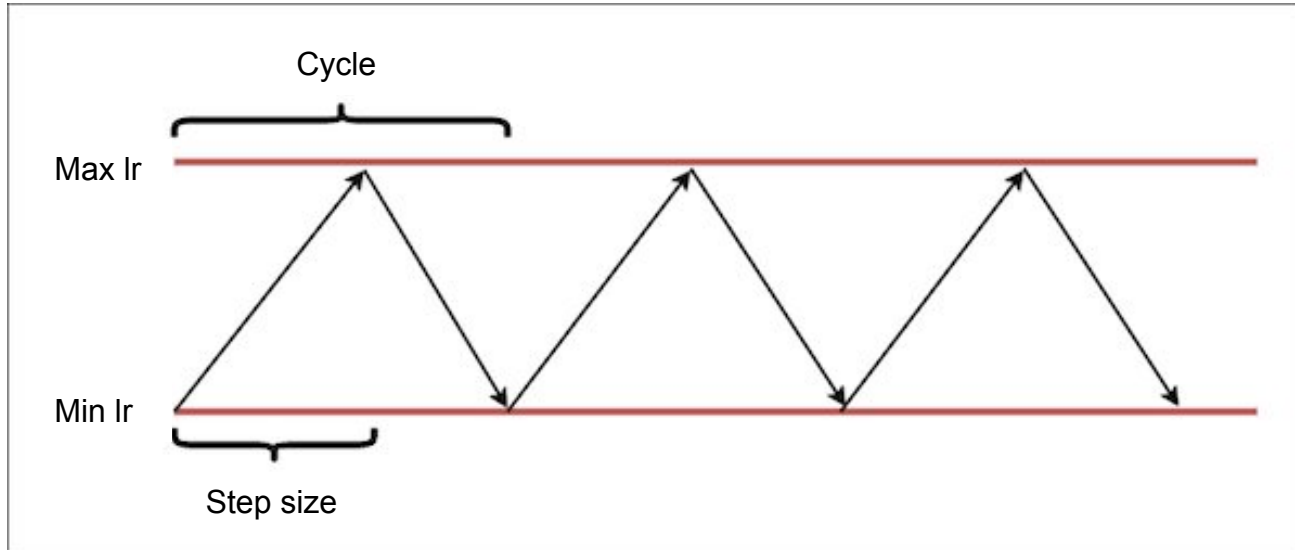
Problèmes : Les premières itérations ont trop d'effet sur le modèle (loss importantes, gradients élevés, biais, ...), un *learning rate* élevé peut provoquer une forte instabilité ou une divergence



But : augmenter petit à petit le *learning rate* pour éviter le risque de divergence au début de l'apprentissage

# Cyclic Learning Rate Scheduler

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

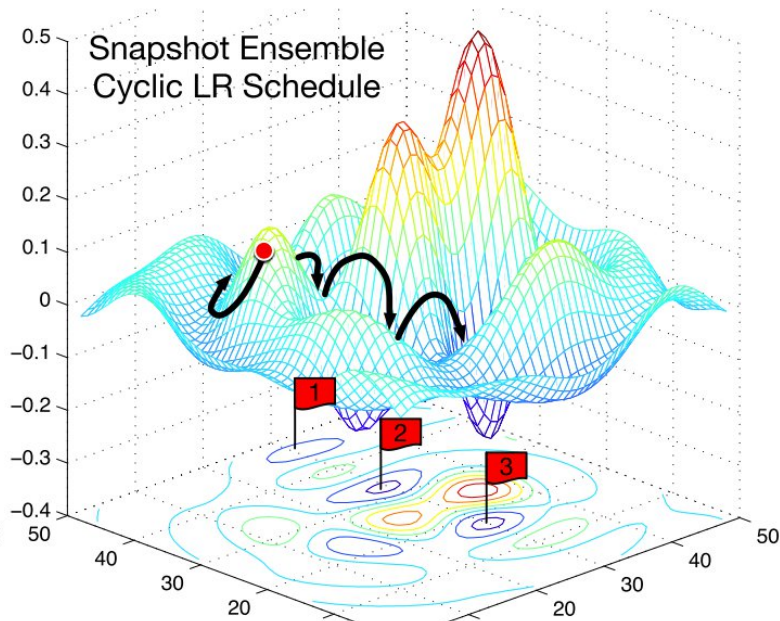
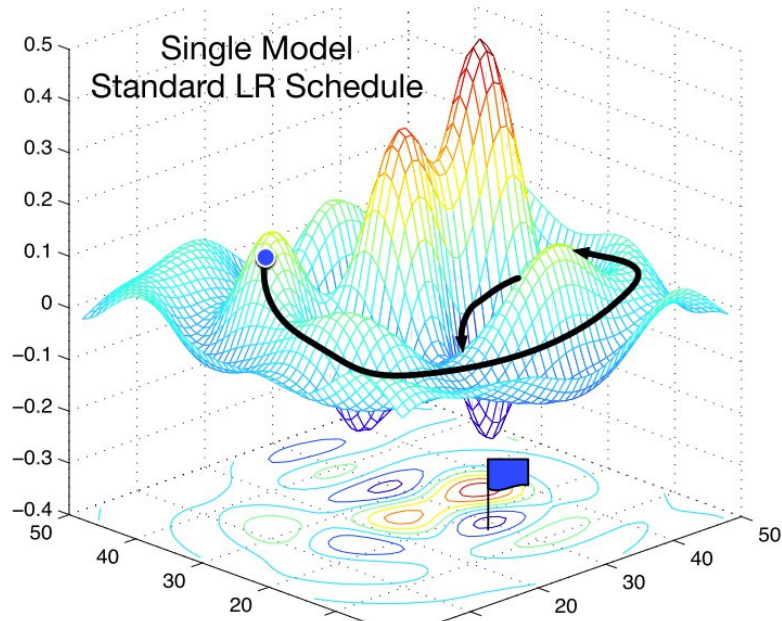


Paramètres :

- $\text{Step\_size} = x * \text{epoch}$  ( $2 \leq x \leq 10$ )
- $\text{Base\_lr}$  -> valeur minimum de convergence
- $\text{max\_lr}$  -> valeur maximale avant divergence

Succession de *warmups* et de *learning rate decays*

# Cyclic Learning Rate Scheduler

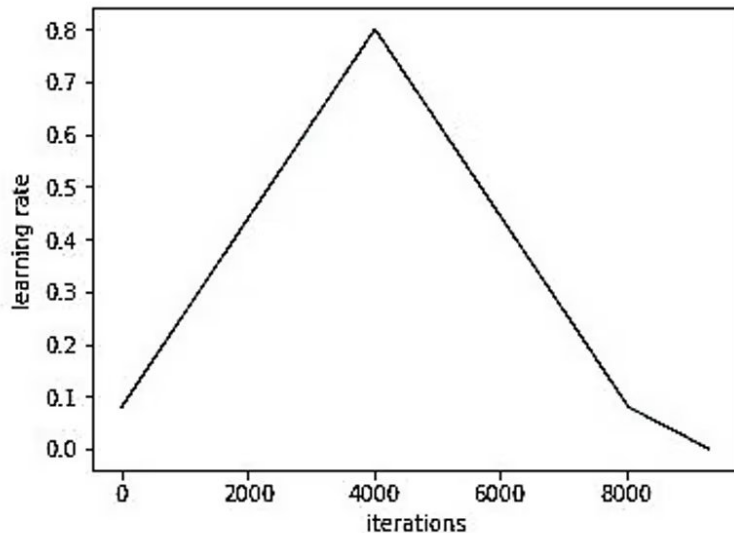


SNAPSHOT ENSEMBLES: TRAIN 1, GET M FOR FREE  
*Gao Huang, Yixuan Li, Geoff Pleiss*

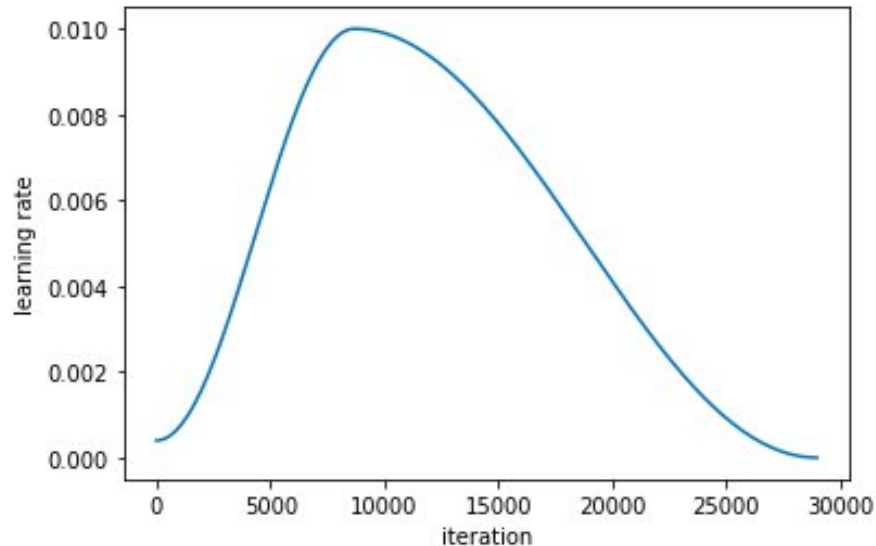
# One Cycle Learning Rate

Un seul cycle suffit ! A disciplined approach to neural network hyper-parameters -  
[Leslie N. Smith](#)

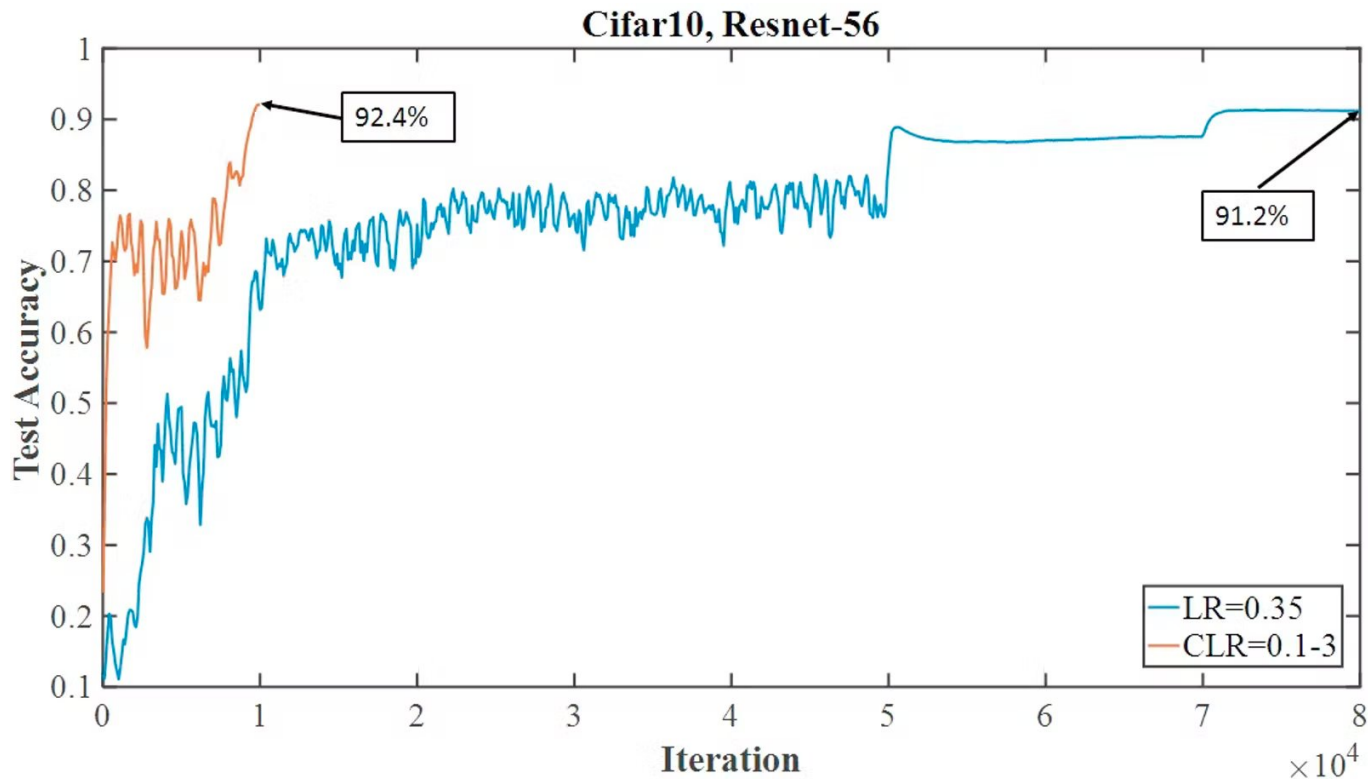
Proposition initiale



cosine annealing : Recommendation par FastAI



# One Cycle Learning Rate - Super convergence



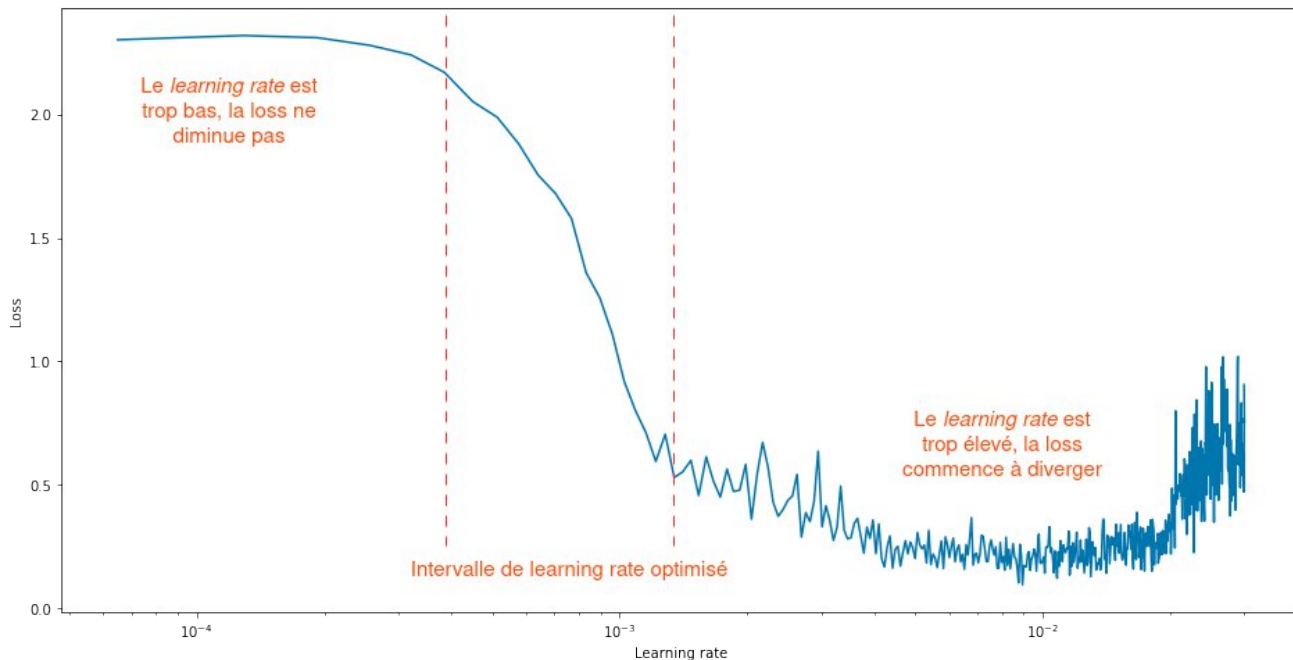
Une convergence plus rapide pour une précision finale équivalente



# Learning Rate Finder

But : Trouver les valeurs de *learning rate* optimales pour son modèle, particulièrement pour la valeur maximale d'un *cyclic scheduler*

- Faire tourner son modèle sur quelques *epochs* en faisant augmenter son *learning rate*
- Début de baisse de la *loss*  
→ *Learning rate minimal*
- Début de variation de la *loss*  
→ *Learning rate maximal*



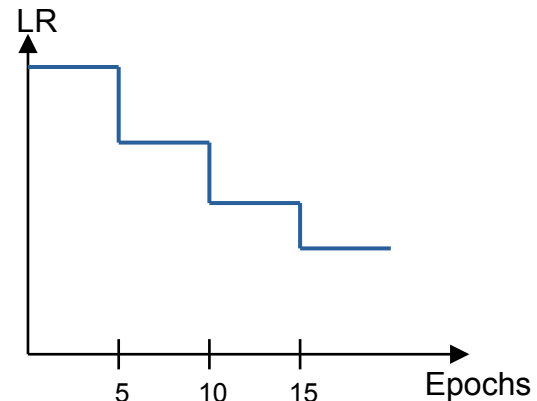
# Learning Rate Scheduler

Chaque *scheduler* a ses propres paramètres

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

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

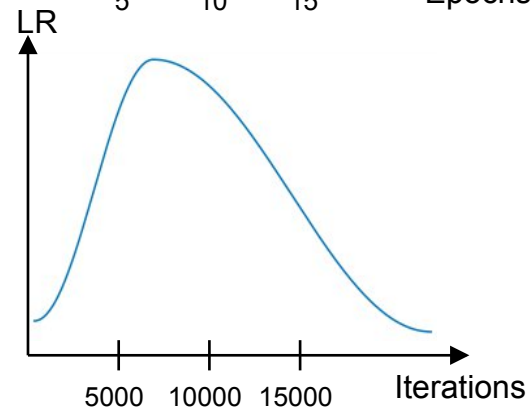
```
for epoch in range(100):  
    train(...)  
    validate(...)  
    scheduler.step()
```



```
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()
```



# Optimiseur de descente de gradient

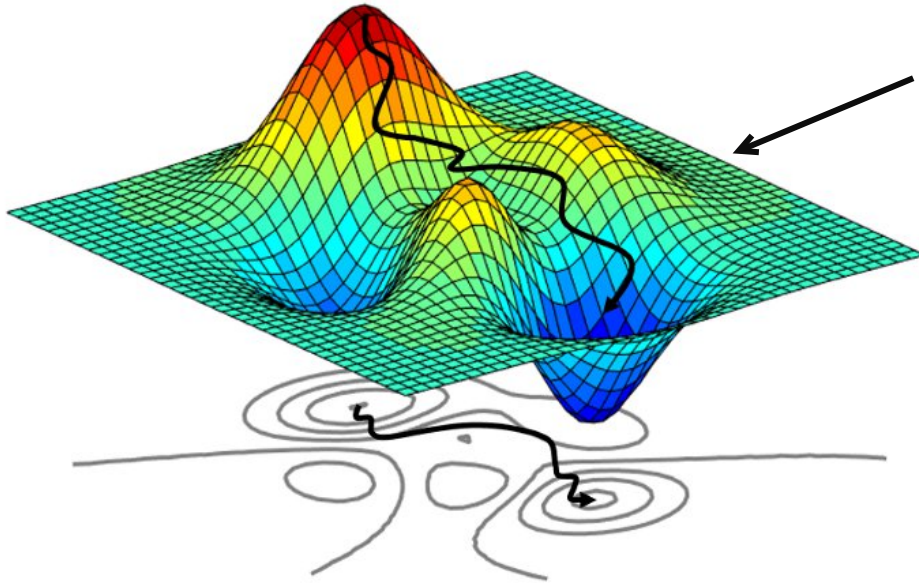
SGD ◀

ADAM ◀

ADAMW ◀

# Optimiseur - SGD

L'optimiseur est l'algorithme qui pilote la descente de gradient et la recherche de minimum avec pour but d'optimiser le temps d'apprentissage et la métrique finale.



SGD = *Stochastic Gradient Descent*  
Calcul du Gradient et mise à jour des poids  
à chaque *batch*

- + Taille de *batch* et *learning rate* adaptables selon les besoins contradictoires :
  - D'exploration pour trouver le meilleur minimum local
  - D'accélération de la descente de gradient

# SGD with Momentum

$$m_0 = 0$$

Coefficient de *momentum*

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

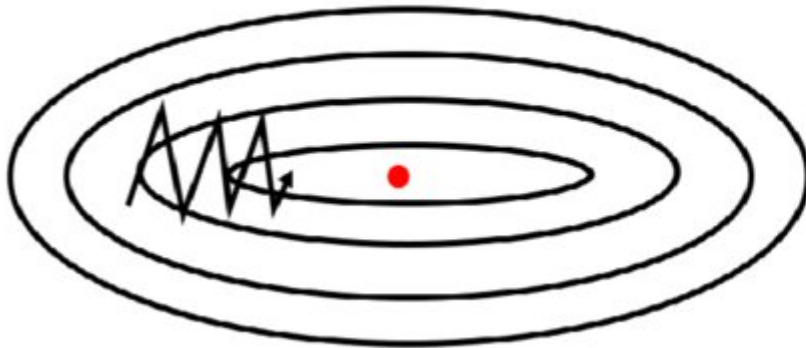
$$\theta_i = \theta_{i-1} - \alpha * m_i$$

Objectif : Prendre en considération les gradients précédents pour une descente de gradient plus rapide.

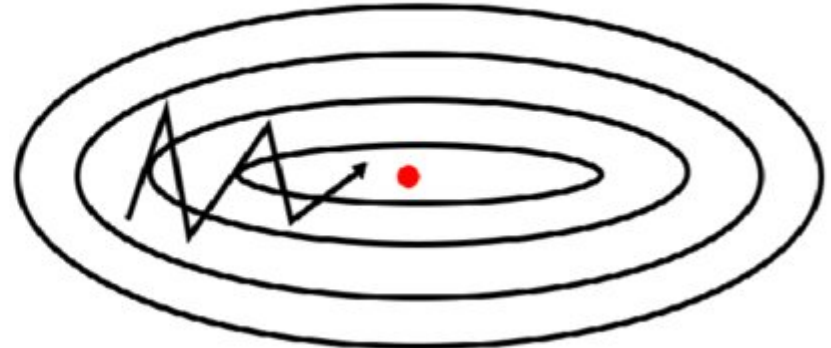
Valeur initiale conseillée : 0,9

$$0.85 < \beta < 0.95$$

SGD without momentum



SGD with momentum

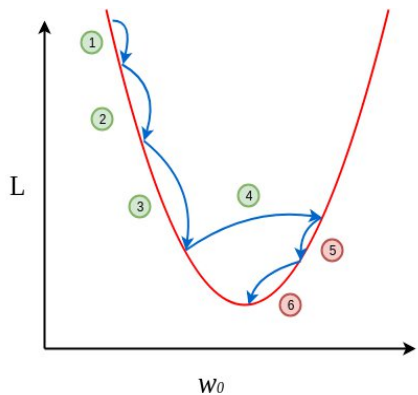


+ Permet de converger plus rapidement

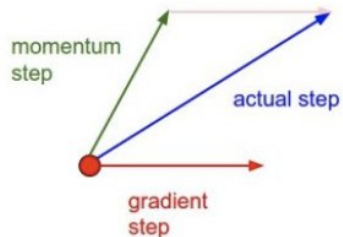
- Pas de garantie que le *momentum* nous amène dans la bonne direction

# Types de *Momentum*

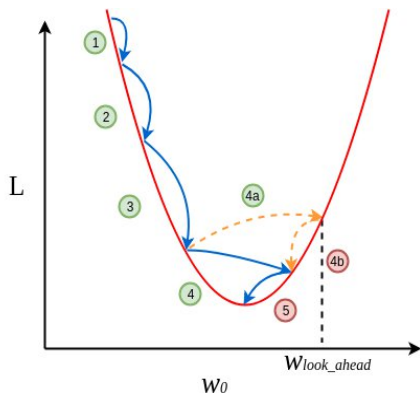
## Momentum



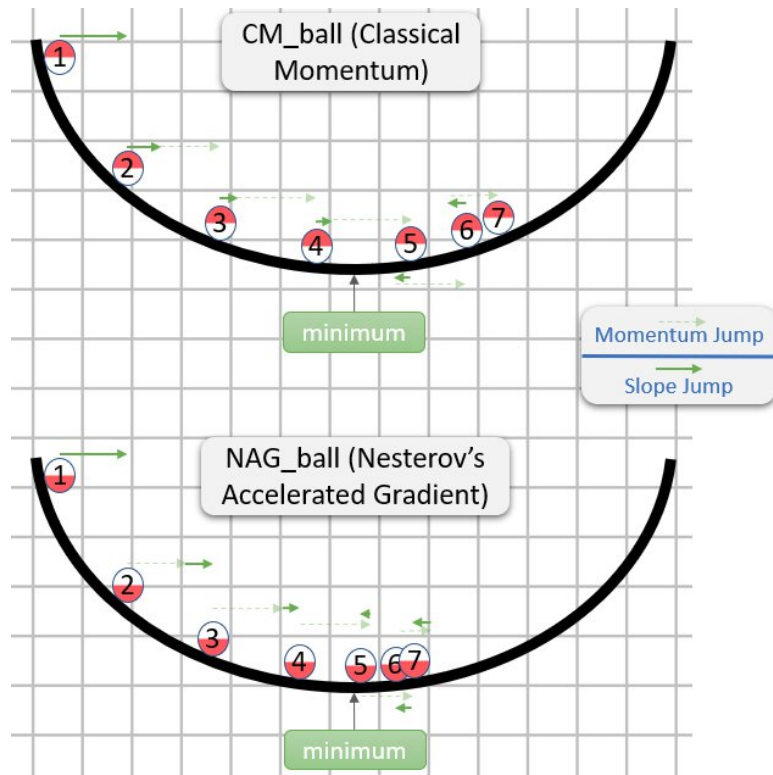
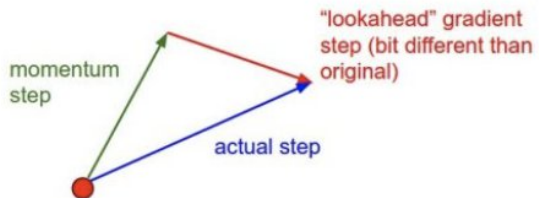
### Momentum update



## Nesterov momentum



### Nesterov momentum update



# Optimiseurs adaptatifs

Plutôt que de piloter la descente de gradient manuellement avec le *learning rate* ...

Nous pouvons adapter le *learning rate* **pour chaque poids du modèle** en fonction du gradient, du gradient<sup>2</sup>, ou de la norme des poids de la couche !!

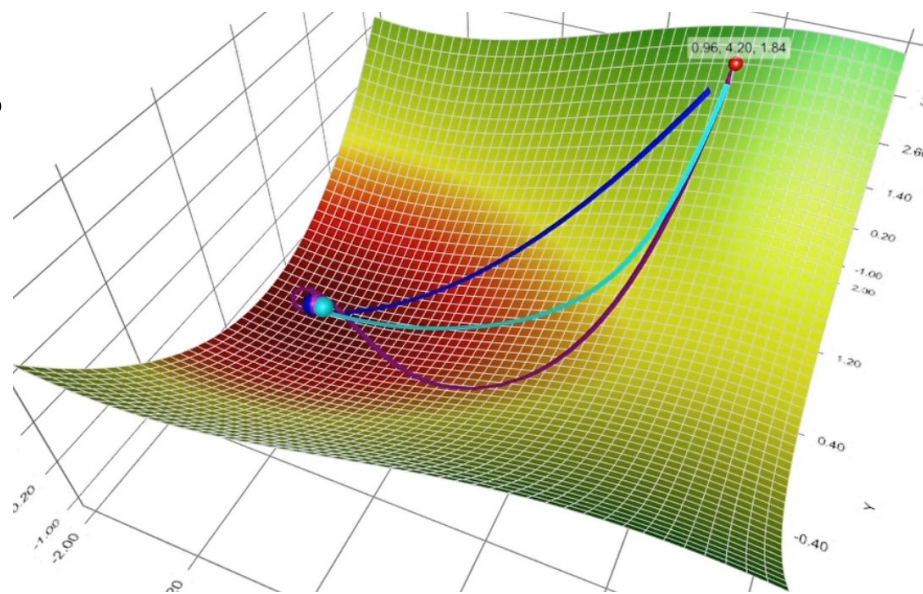
Exemples :

- AdaGrad,
- AdaDelta,
- RMSprop
- Adam

Spécialisés pour les batches larges :

- LARS
- LAMB

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



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

Premier moment : moyenne glissante

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

Second moment : variance non centrée glissante

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

Correction des biais des premières itérations

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

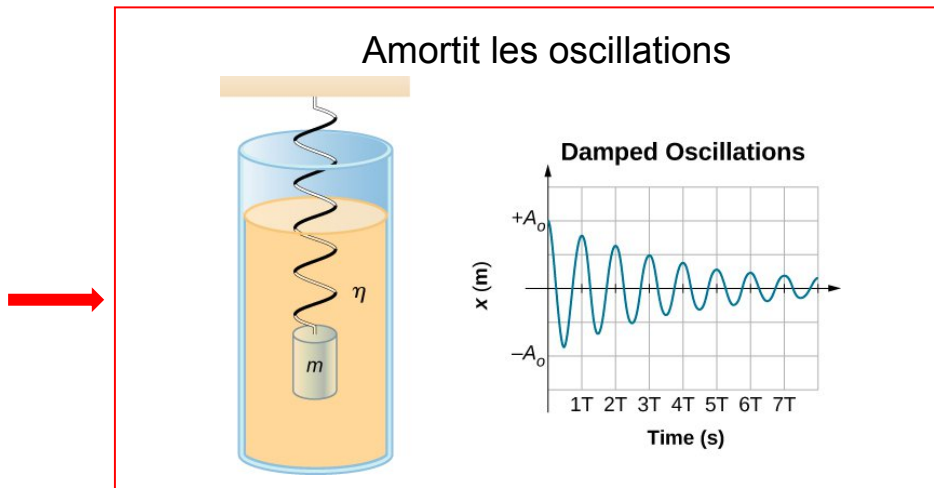
$$\theta_i = \theta_{i-1} - \frac{\alpha}{\sqrt{\hat{v}_i + \epsilon}} * \hat{m}_i$$

Paramètres :

$\beta_1$  &  $\beta_2$  = Taux de régression ( $\beta_1 = 0.9$  &  $\beta_2 = 0.999$ )

$\epsilon$  — Très petite valeur pour éviter une division par zéro

But : Adapter l'importance des mise à jours de poids en fonction des précédents gradients et de la variabilité du gradient.

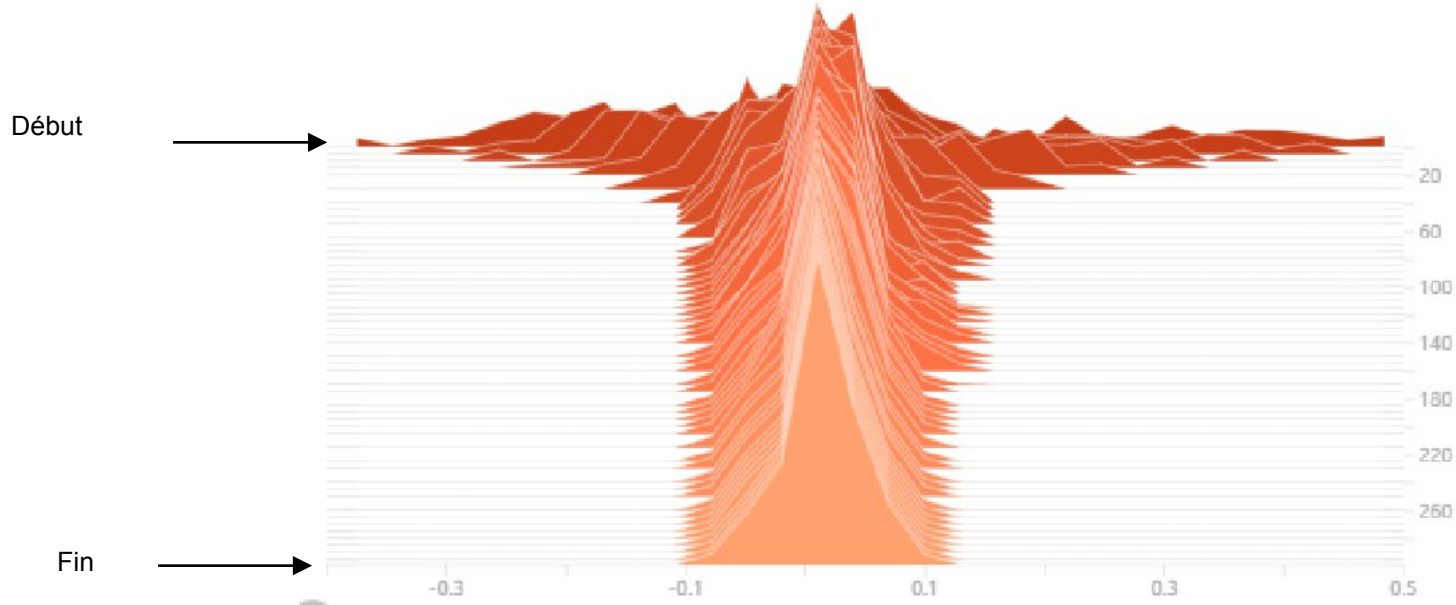




# Weight decay

Un réseau de neurones qui converge et généralise correctement (ni sous-apprentissage ni sur-apprentissage) a généralement **des poids qui tendent vers 0**.

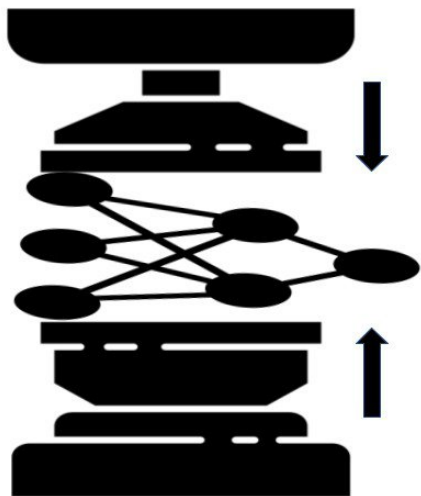
Distribution des poids durant l'apprentissage :



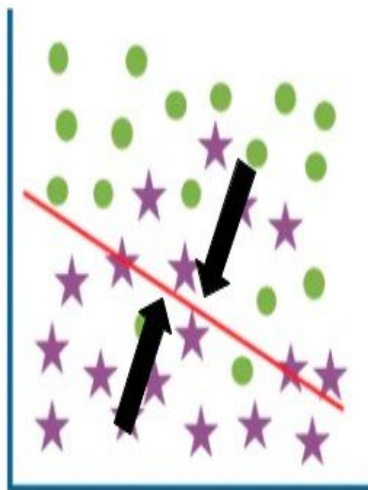
# Weight decay

Préférable à la régularisation L2 standard définie dans la fonction de perte

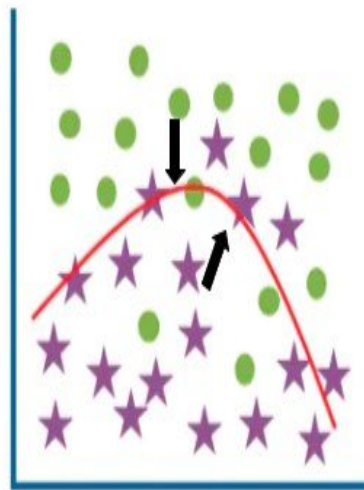
$\lambda$  : paramètre du *weight decay* (généralement entre 0 et 0.1)



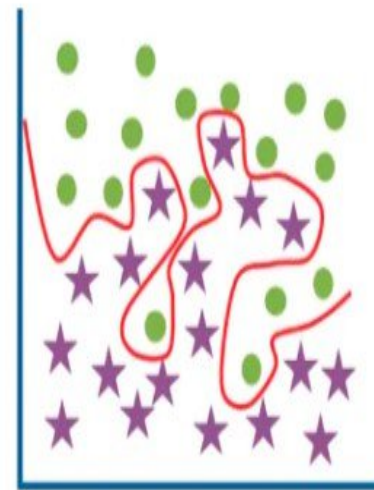
*Underfitting*  
( $\lambda$  trop élevé)



*Weight decay correct*  
(Valeur idéal de  $\lambda$ )



*Overfitting*  
( $\lambda$  trop faible ou nul)



La technique de *weight decay*, définie dans *l'optimizer* permet de forcer les poids à converger vers des valeurs proches de zéro.

# 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

Évolution du *weight decay*: *Decoupled weight decay* (découplé du momentum !!)

- SGD et Adam avec le *weight decay*
- SGD et AdamW avec le *decoupled weight decay*

SGD et SGD et AdamW sont à peu près équivalents en performance.

**Cependant AdamW est notablement meilleur que Adam !!**

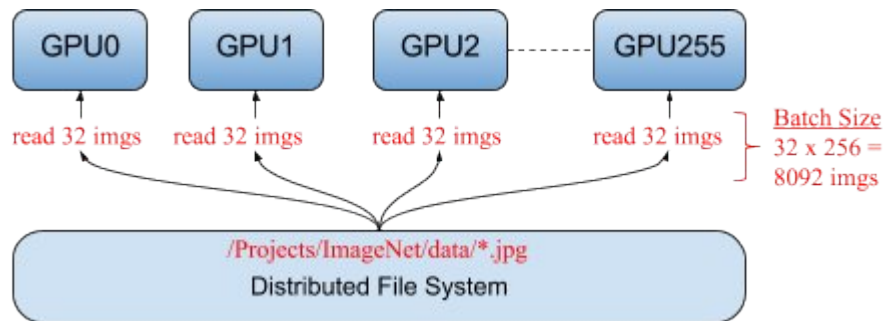
# Optimisation des larges batch

Problématiques larges batches ◀

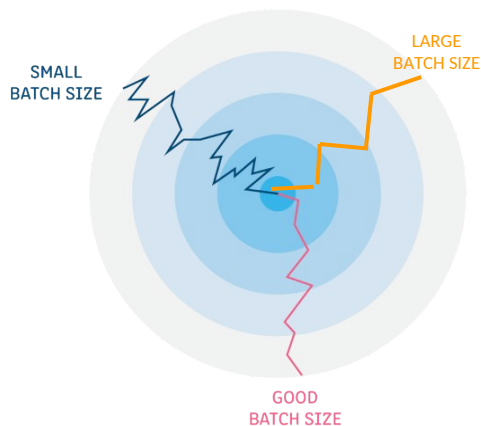
Learning Rate Scaling & Batch Schedulers ◀

Optimiseurs larges batches ◀

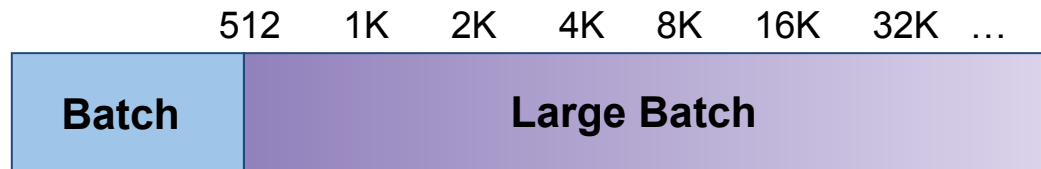
# Large Batches avec le parallélisme de données



Parallélisme de données : La parallélisation implique une grande taille de *batch*



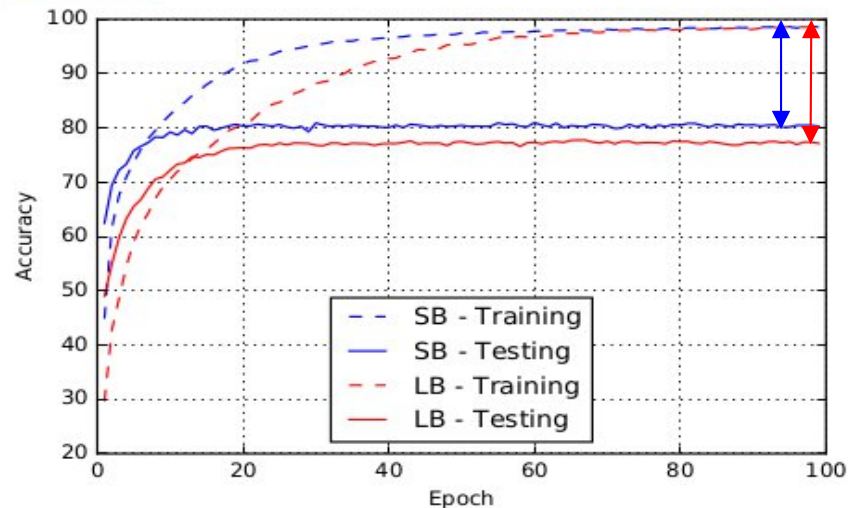
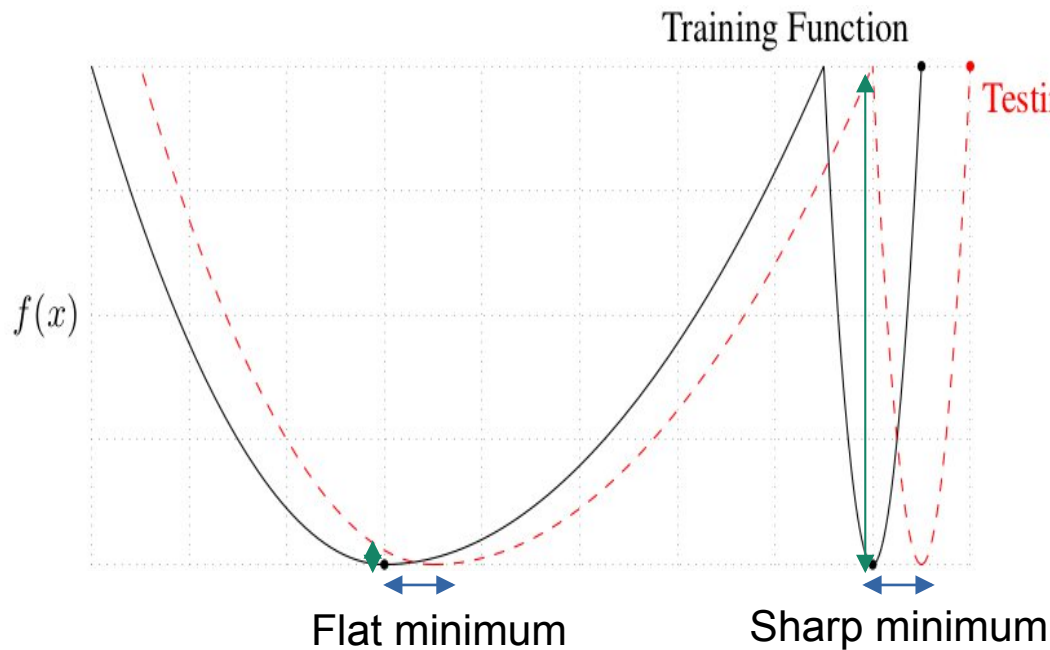
Problème : Les *batches* trop grands ( $> 512$ ) ont tendance à engendrer de moins bonnes performances



# 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



Comparaison d'entraînement d'un réseau convolutionnel avec des petits batch (SB) et large batch (LB) sur CIFAR 10

Plus le *batch* est grand, plus le modèle a tendance à converger vers des minimums pentus et étroits.

# Large Batches : Learning rate scaling

Lorsqu'on augmente considérablement la taille du *batch* global, il est souvent nécessaire de mettre le *learning rate* à l'échelle :

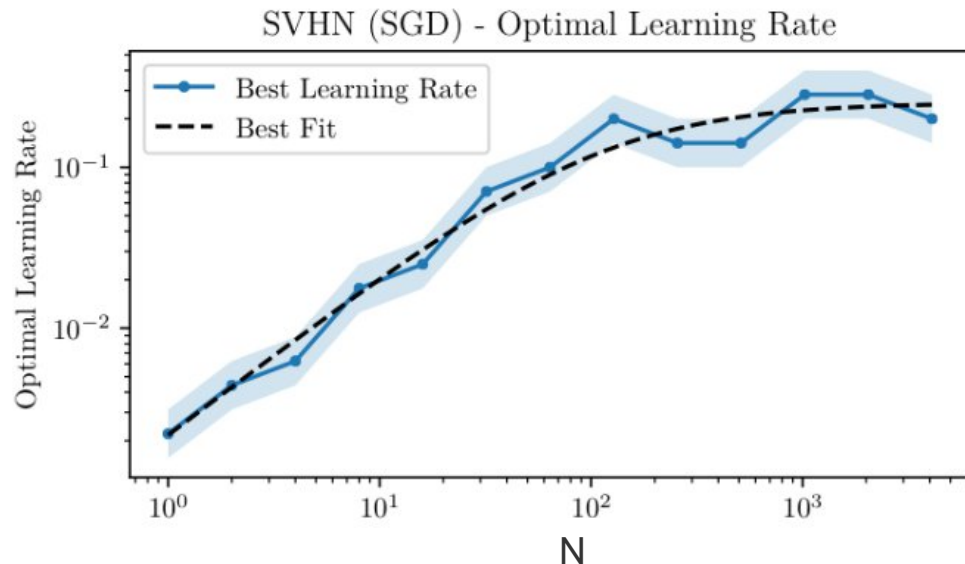
N = Nombre de processus parallèles

Croissance linéaire du *learning rate* :

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

Croissance en racine carrée du *learning rate* :

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



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

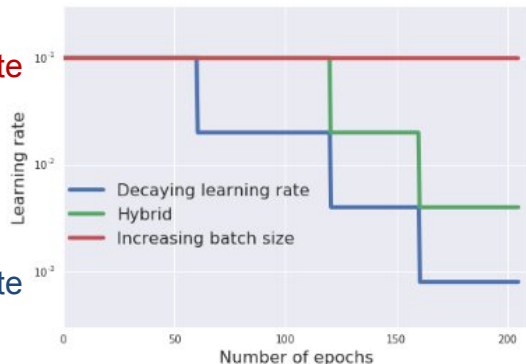
Optimal : croissance linéaire au début puis en racine carrée (recommandé par OpenAI)

# Batch Size Scheduler

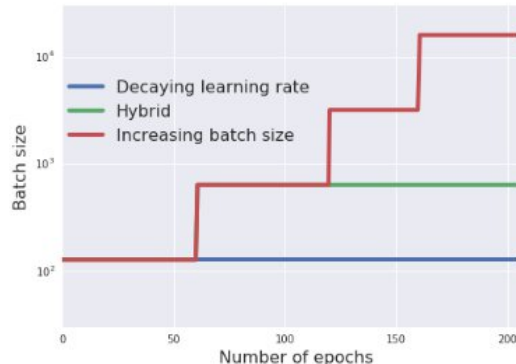
=> Alternative au *Learning Rate Scheduler*

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

High Learning Rate

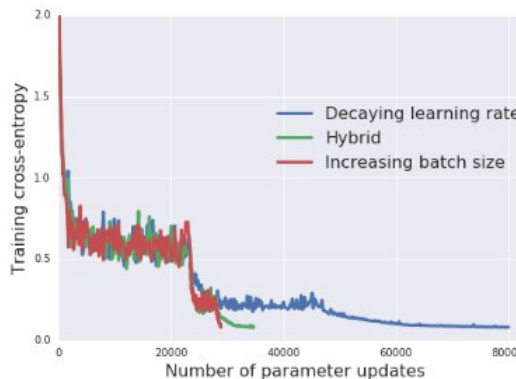
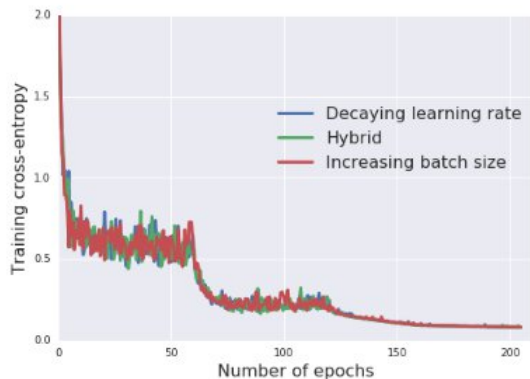


Low Learning Rate



Large Batch

Small Batch





# Large Batches

## SGD

## Adam

Weight  
Decay

Tendencies :

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

- Test Loss

+ Test Loss

Slow Descent

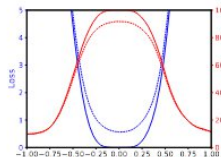
Fast Descent

← Small Batch | Large Batch →

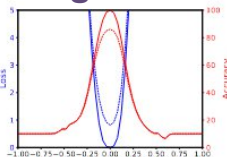
← SGD | ADAM →

← Weight Decay | Sans W Decay →

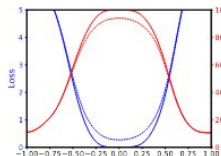
Small Batch | Large Batch



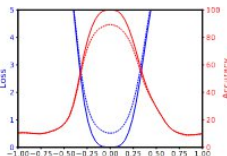
(a) SGD, 128, 8.26%



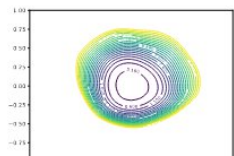
(b) SGD, 4096, 13.93%



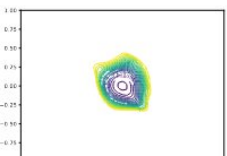
(e) SGD, 128, 5.89%



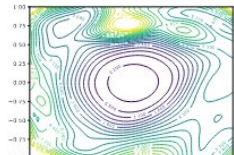
(f) SGD, 4096, 10.59%



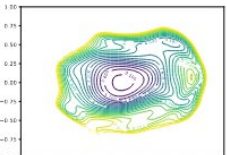
(i) SGD, 128, 8.26%



(j) SGD, 4096, 13.93%



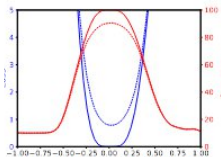
(m) SGD, 128, 5.89%



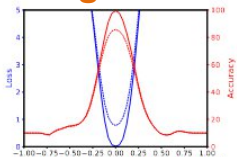
(n) SGD, 4096, 10.59%

= 0

Small Batch | Large Batch

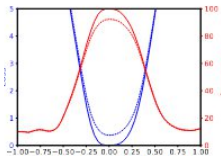


(c) Adam, 128, 9.55%

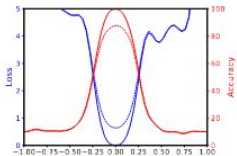


(d) Adam, 4096, 14.30%

=  $5e^{-4}$

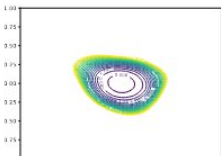


(g) Adam, 128, 7.67%

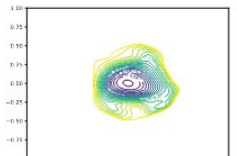


(h) Adam, 4096, 12.36%

= 0

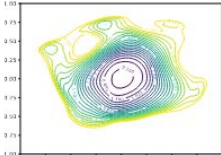


(k) Adam, 128, 9.55%

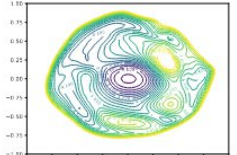


(l) Adam, 4096, 14.30%

=  $5e^{-4}$



(o) Adam, 128, 7.67%



(p) Adam, 4096, 12.36%

# Optimiseurs *Large Batches* - LARS

## LARS pour “Layer-wise Adaptive Rate Scaling.”

Adaptation de *SGD with momentum* avec l'ajout d'un trust ratio pour chaque couche qui dépend de l'évolution du gradient de la couche

$r$  = Trust ratio

$l$  = Numéro de couche

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

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

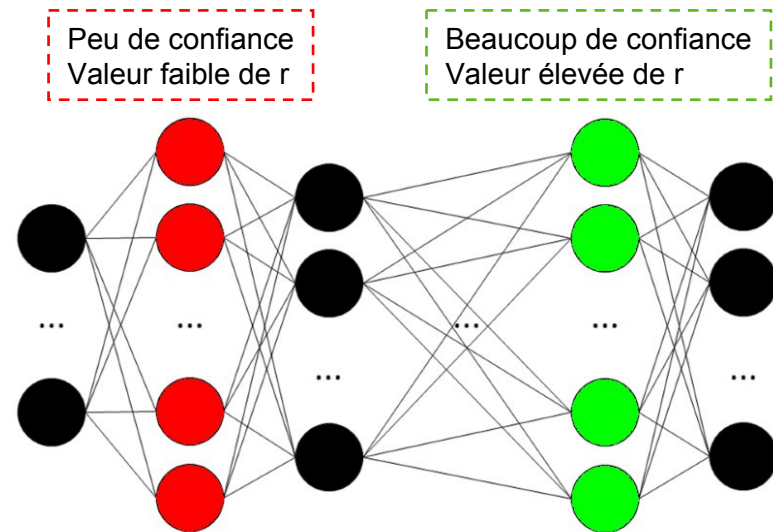
$$r_2 = \left\| m_i^l \right\|_2$$

$$r = r_1 / r_2$$

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

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

Weight decay



But : Adapter l'importance des mise à jours des poids en fonction d'un *trust ratio* calculé pour chaque couche du réseau.

# Optimiseurs *Large Batches* - LAMB

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

Adaptation de Adam avec l'ajout d'un trust ratio pour chaque couche qui dépend de l'évolution du gradient de la couche

$r$  = Trust ratio

$l$  = Numéro de couche

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

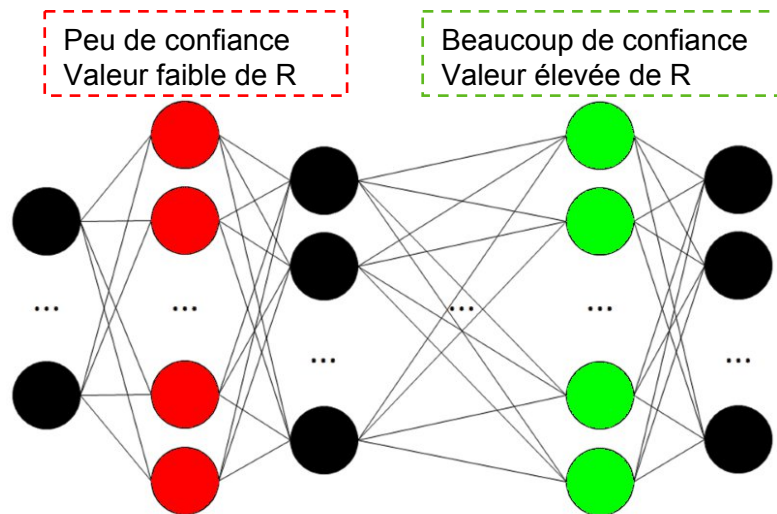
$$r_2 = \left\| \frac{\hat{m}_i^l}{\sqrt{\hat{v}_i^l + \epsilon}} + \lambda \theta_{i-1}^l \right\|_2$$

$$r = r_1 / r_2$$

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

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

Decoupled weight decay



But : Adapter l'importance des mise à jours des poids en fonction d'un *trust ratio* calculé pour chaque couche du réseau.

# Implémentation des optimiseurs

Chaque optimiseur a ses propres paramètres

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

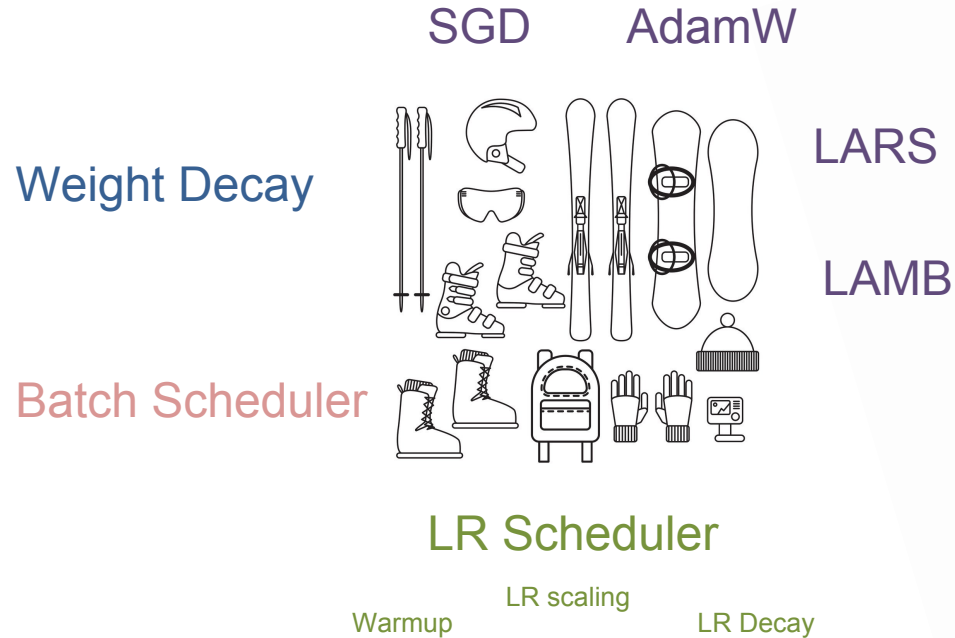
```
import torch.optim as opt
```

```
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)
```

Optimisation  
apex de LARS

# Large Batches Rider



# BLOOM example

**BLOOM**  
176B params 59 languages Open-access



95281 steps (116.8 days)

AdamW,

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

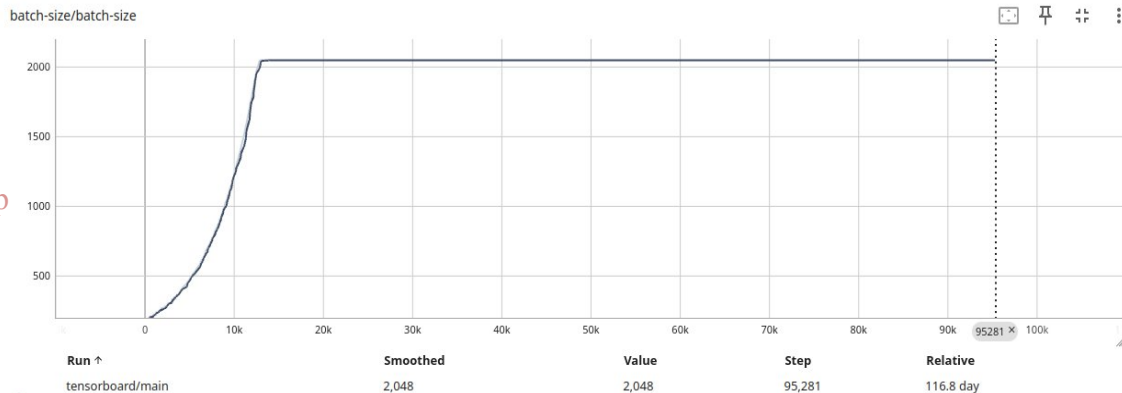
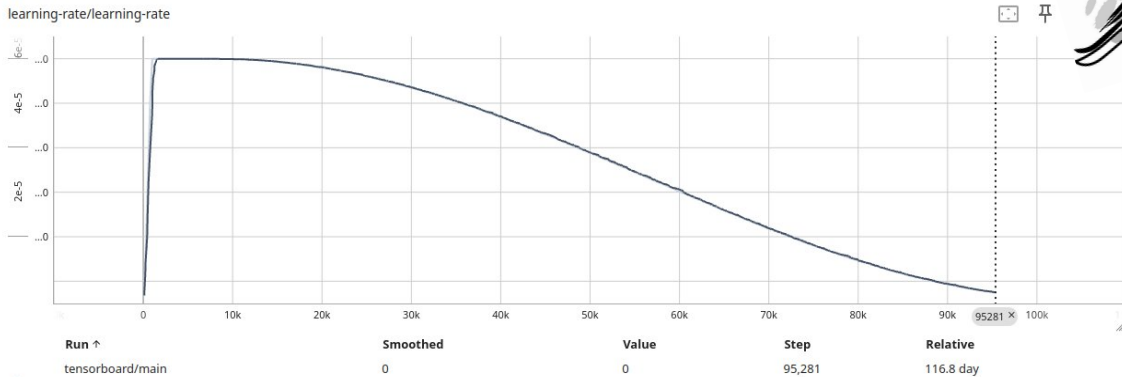
Weight Decay of 0.1

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  $\sim 20B$  tokens
- then continue at 4.2M tokens/step (GBS=2048) for 430B tokens



# Nouveaux optimiseurs

Nouvelle tendance : l'apprentissage d'optimiseur ◀

Gouffre des nouveaux optimiseurs ◀

LION : exemple d'une nouvelle approche ◀

# Nouvelle tendance : l'apprentissage d'optimiseur



**Learning to Optimize Neural Nets**  
Ke Li<sup>1</sup> Jitendra Malik<sup>1</sup>

**VeLO: Training Versatile Learned Optimizers**  
Luke Metz<sup>†</sup> James Harrison<sup>†</sup> C. Daniel Freeman, Amil Merchant,  
Lucas Beyer, James Bradbury, Naman Agarwal, Ben Poole,  
Igor Mordatch, Adam Roberts, Jascha Sohl-Dickstein<sup>‡</sup>  
Google Research, Brain Team

**Neural Optimizer Search with Reinforcement Learning**  
Irwan Bello<sup>∗1</sup> Barret Zoph<sup>∗1</sup> Vijay Vasudevan<sup>1</sup> Quoc V. Le<sup>1</sup>

**Learning to learn by gradient descent  
by gradient descent**  
Marcin Andrychowicz<sup>1</sup>, Misha Denil<sup>1</sup>, Sergio Gómez 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  
{andnil,sergomez,mwhoffman,pfau,schaul}@google.com  
brendan.shillingford@cs.ox.ac.uk, nandodefreitas@google.com





# LION : un nouvel optimiseur

---

## 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  $\theta_t$ 
```

---

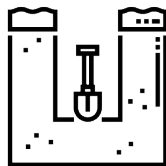
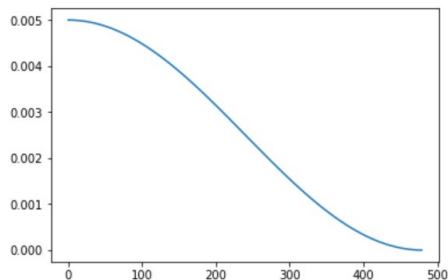
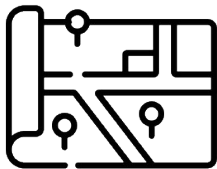
---

## Algorithm 2 Lion Optimizer (ours)

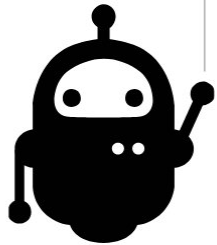
---

```
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  $\theta_t$ 
```

---



# TP : Learning rate + Optimiseurs



Objectifs :

- Modifier le learning rate scheduler
- Modifier l'optimiseur
- Faire des entraînements avec des large batches



Depuis JupyterHub :

- Lancer une instance SLURM CPU
- Allez dans le dossier `tp_optimiseurs`
- Ouvrez le notebook `DLO-JZ_Optimiseurs`