

Les *CUDA Core* sont spécialisés pour le calcul vectoriel.

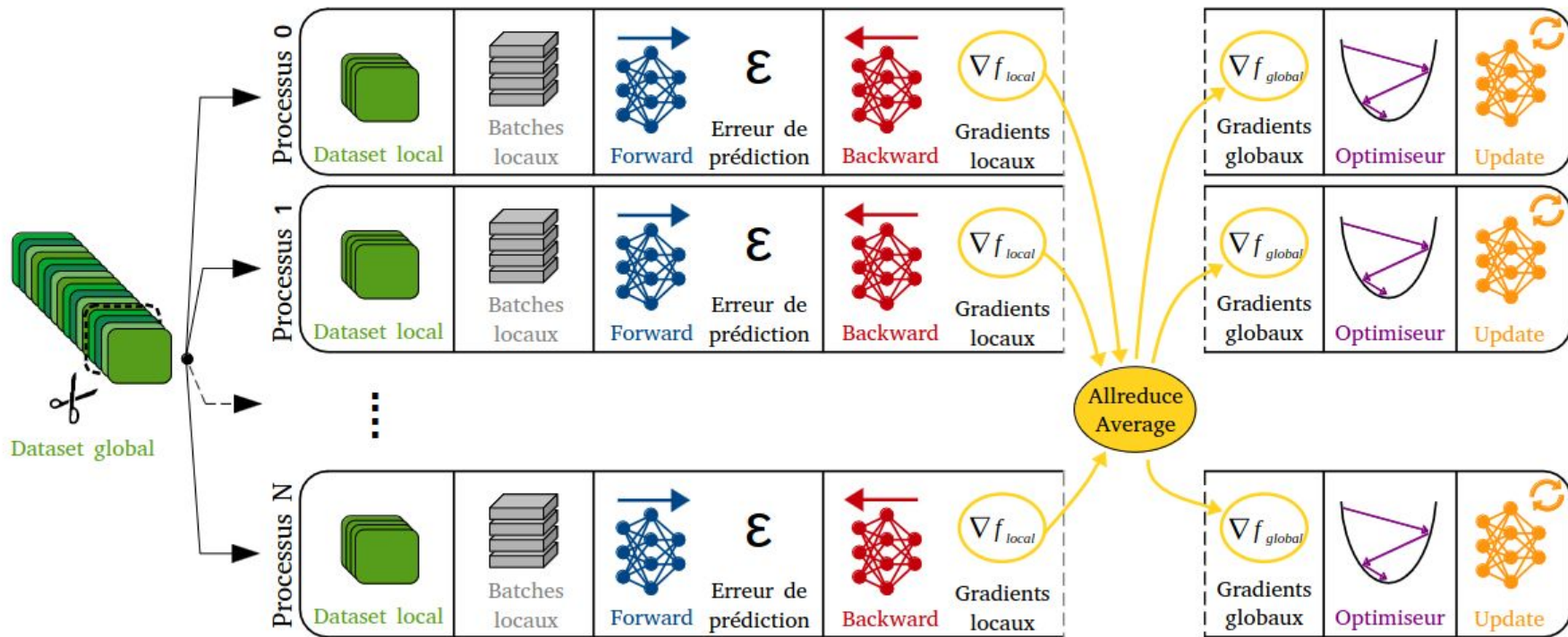
Les *Tensor Core* sont spécialisés pour le **calcul matriciel**.

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32                      FP16                      FP16 or FP32

Chaque *Tensor Core* est capable de traiter 64 opérations en 1 temps d'horloge.

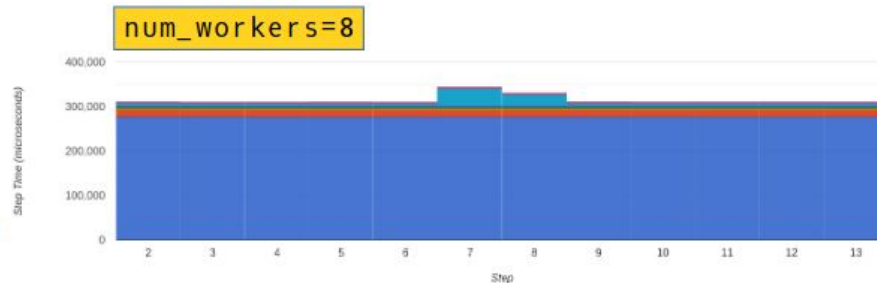
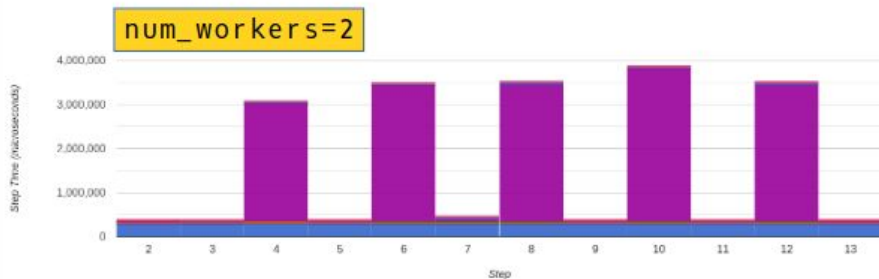
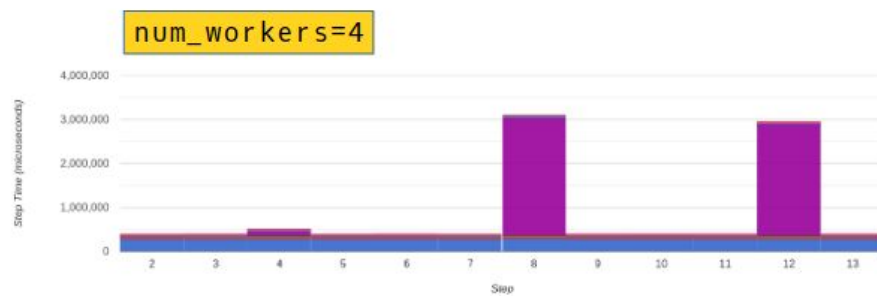
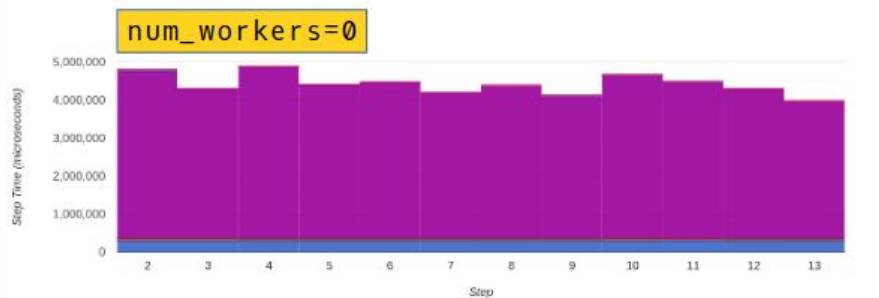
# Distribution : Parallélisme de données



# TP2\_1 : Optimisation du DataLoader



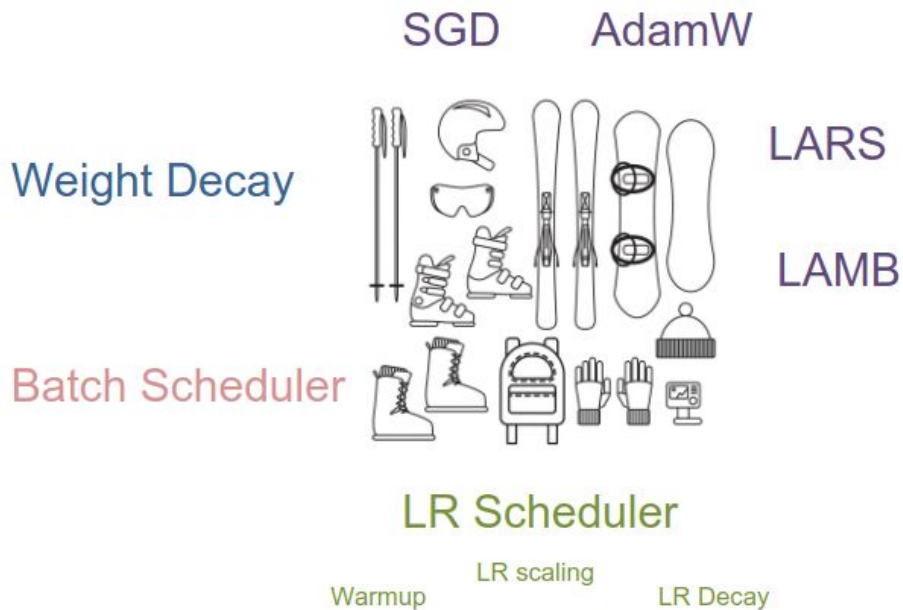
- L'optimisation la plus efficace est l'augmentation du nombre de *workers*



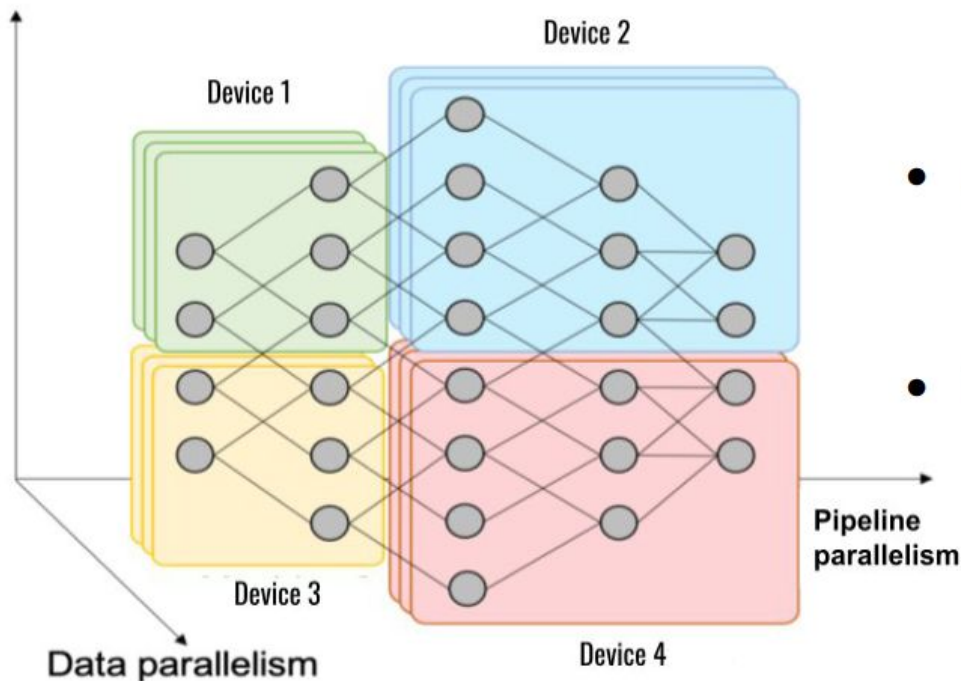
Kernel Memcpy Memset Runtime DataLoader CPU Exec Other



# Large Batches Rider



Tensor parallelism



- **Data Parallelism**

- Simple à implémenter
- Meilleure performance
- Augmente la taille du batch (problème de convergence)

- **Pipeline Parallelism**

- Effort d'implémentation.
- Équilibre entre mémoire, performance et convergence.

- **Tensor Parallelism**

- Effort important d'implémentation
- Bonne accélération des calculs
- **Bande passante très sollicitée** (implémentation Intra-nœud)

# Conclusion

- Pay attention to **Dataloader**. It is often the bottleneck. Profile it, Optimize it, Boost I/O !
- **For usual size model**, use Data Parallelism, Tensor Core, Mixed Precision, Large Batch optimizations, on 1 or few nodes.
- **For huge model**, use dedicated library like Megatron-LM, Deepspeed, Colossal-AI on several nodes.

