



Deep Learning Optimisé - Jean Zay

Conclusion



INSTITUT DU
DÉVELOPPEMENT ET DES
RESSOURCES EN
INFORMATIQUE
SCIENTIFIQUE



Conclusion

Synchronous : num_worker = 0

DataLoader

Forward/Backward



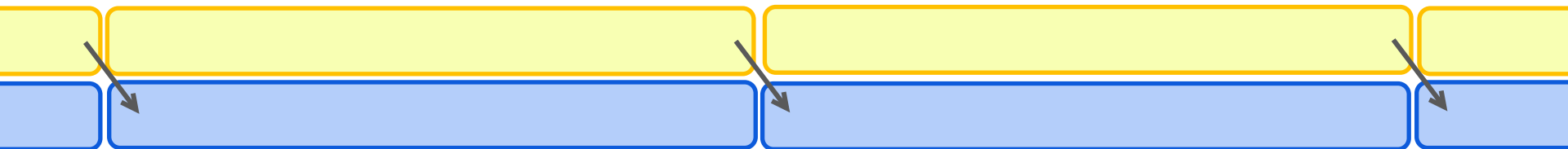
Time →

Conclusion

Asynchronous : $\text{num_worker} > 0$

DataLoader

Forward/Backward



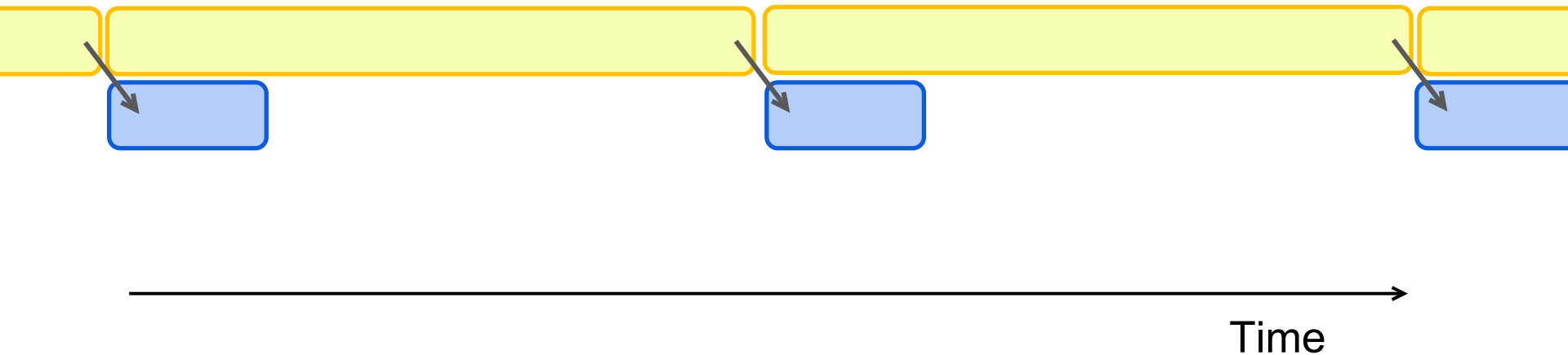
Time

Conclusion

GPU Computing, Mixed Precision,
torch.compile, ...

DataLoader

Forward/Backward

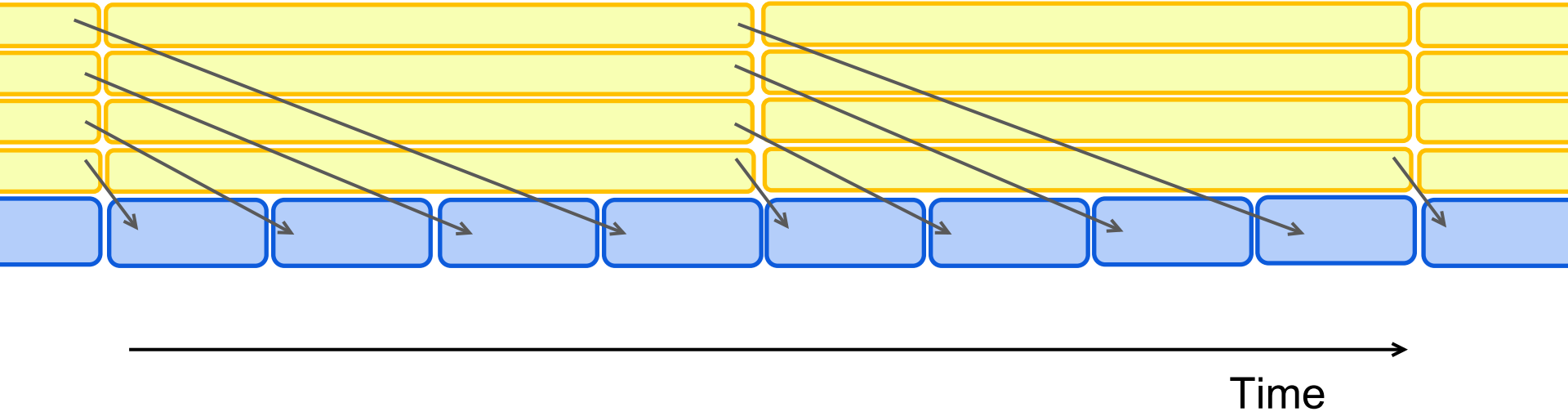


Conclusion

DataLoader Optim. : $\text{num_worker} > 1, \dots$

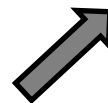
DataLoader

Forward/Backward



Conclusion

Training take too long !!!



Increase your batch size

Conclusion

Training take too long !!!



Increase _{your} batch size

CUDA Out Of Memory !!!



Decrease _{your} batch size

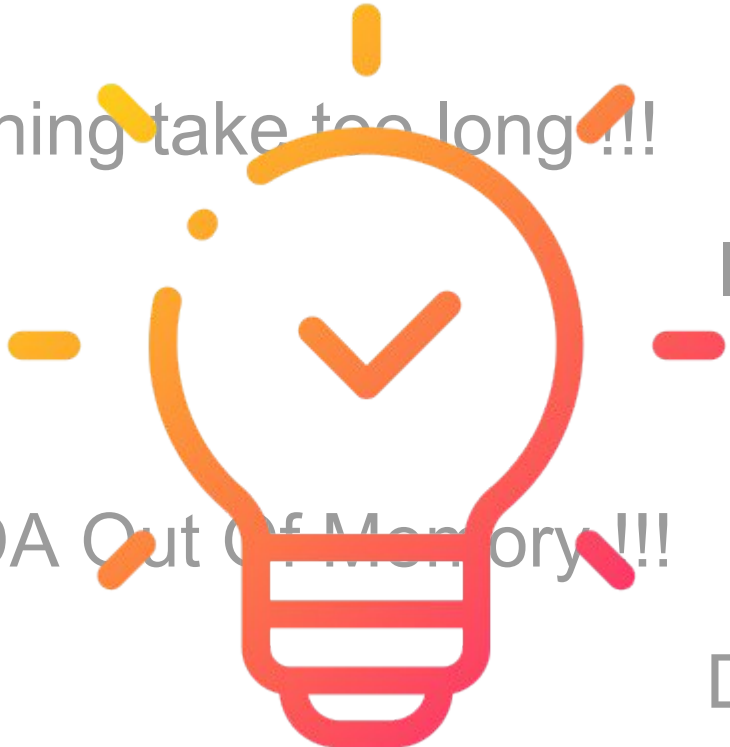
Conclusion

Training take too long !!!

↑
Increase your batch size

CUDA Out Of Memory !!!

↓
Decrease your batch size





For Small Model !!!
10s or 100s M Params

**Distributed Data
Parallelism**



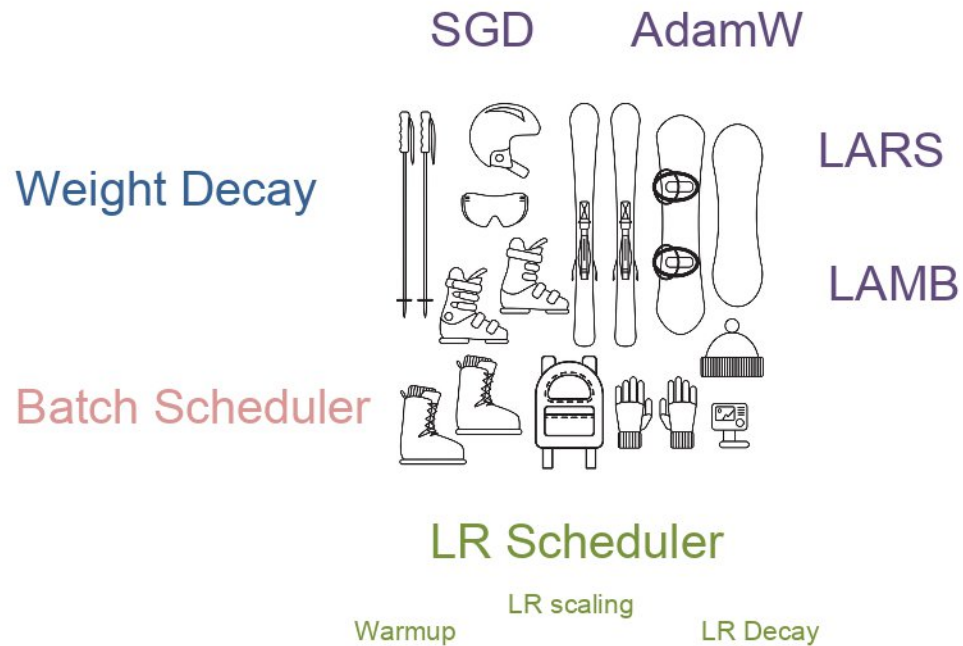
For Small Model !!!
10s or 100s M Params

Distributed Data Parallelism



→ Large Batch !!

Conclusion



Conclusion



For Large Model !!!
> 1G Params

ZeRO

FSDP

Model Parallelisms

Pipeline Parallelism

Tensor Parallelism