# Deep Learning Optimisé - Jean Zay

## Conclusion

**Synchronous : num_worker = 0**

DataLoader          Forward/Backward
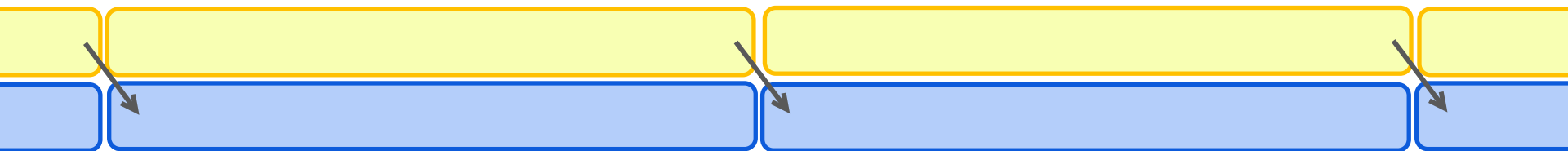
Time

**Asynchronous : num_worker > 0**
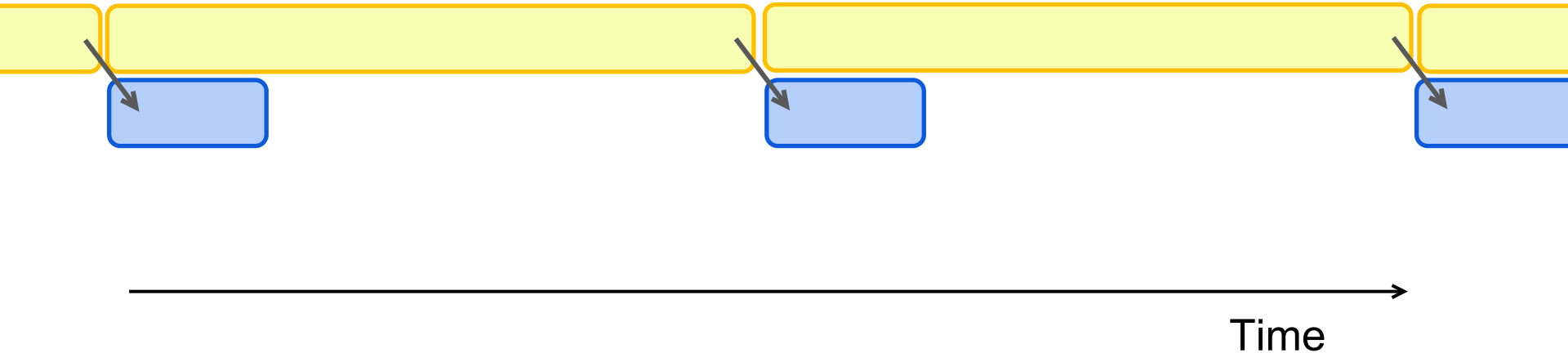
DataLoader          Forward/Backward



Time

# Conclusion

**GPU Computing, Mixed Precision, torch.compile, ...**
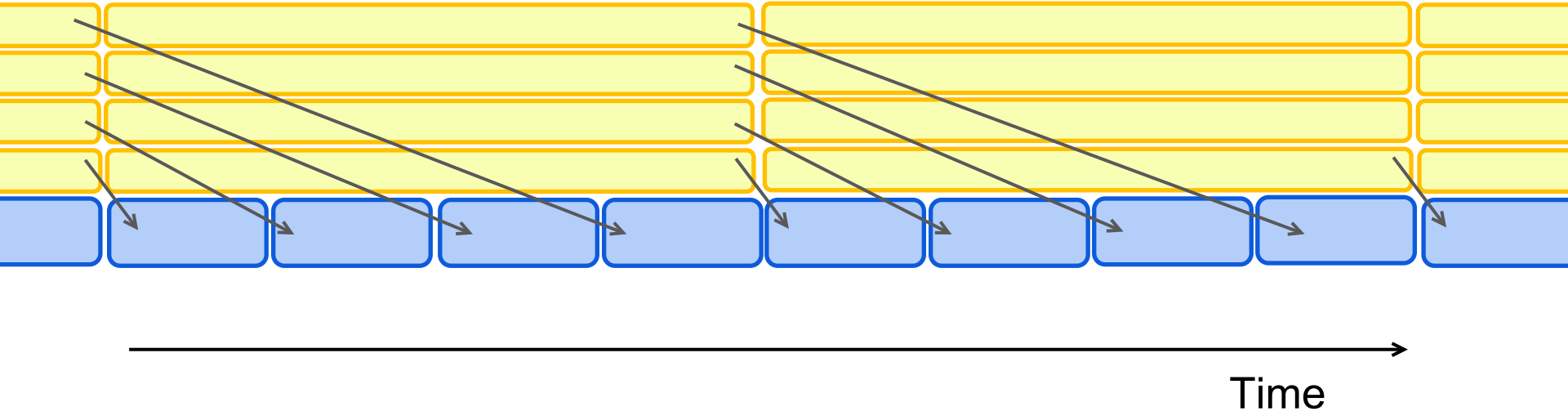
DataLoader      Forward/Backward

Time

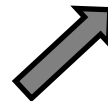**DataLoader Optim. : num_worker > 1, ...**

DataLoader Forward/Backward



Time

# 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

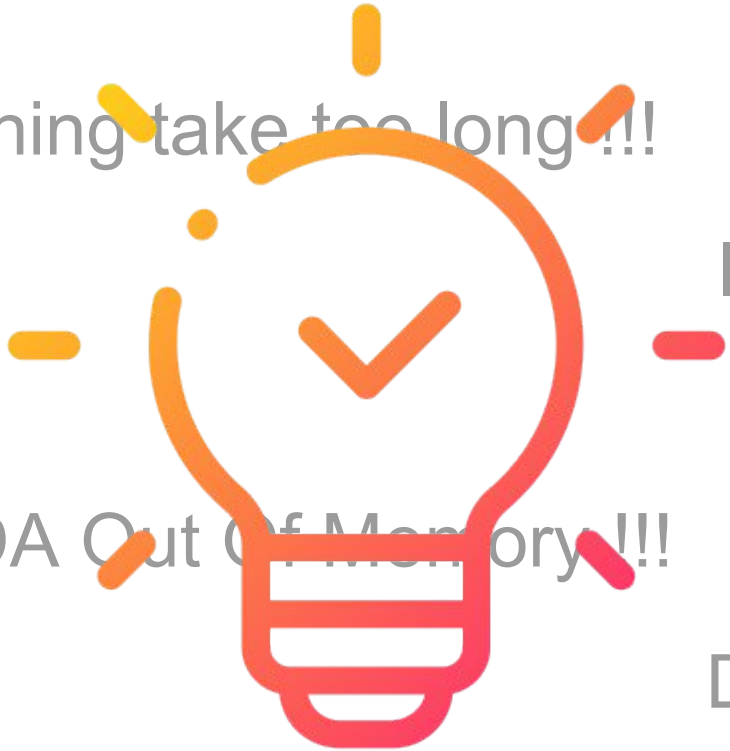# Conclusion

Training take too long !!!  →  Increase your batch size

CUDA Out Of Memory !!!  →  Decrease your batch size

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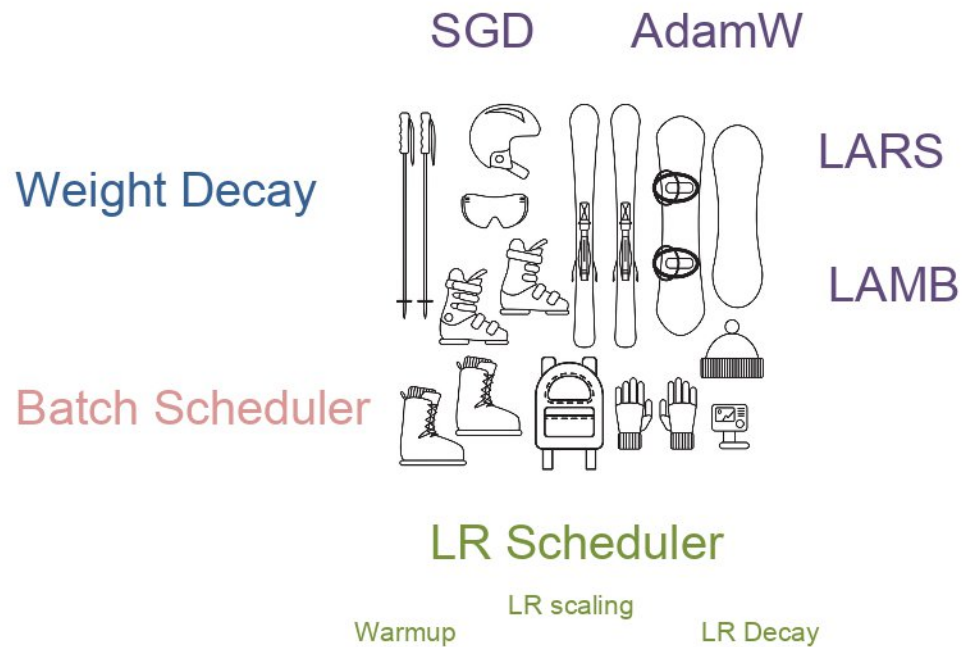10s or 100s M Params

# **Distributed Data Parallelism**

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

**Distributed Data Parallelism**

⟹ Large Batch !!

# Conclusion



SGD    AdamW

LARS

Weight Decay

LAMB

Batch Scheduler

LR Scheduler

LR scaling

Warmup            LR Decay

Sharp
Minima

# Conclusion

For Large Model !!!
> 1G Params

ZeRO

FSDP

**Model Parallelisms**
Pipeline Parallelism
Tensor Parallelism