# Optimised Deep Learning - Jean Zay

## Practical work (of your choice !)

# Like you want !

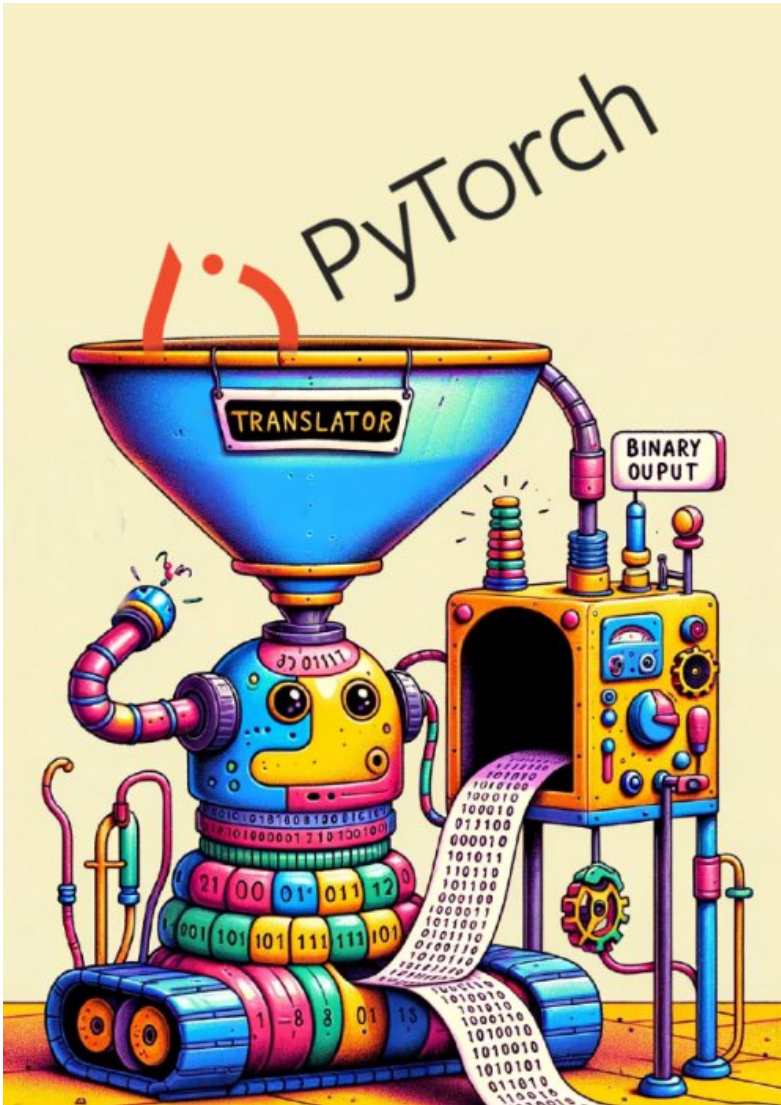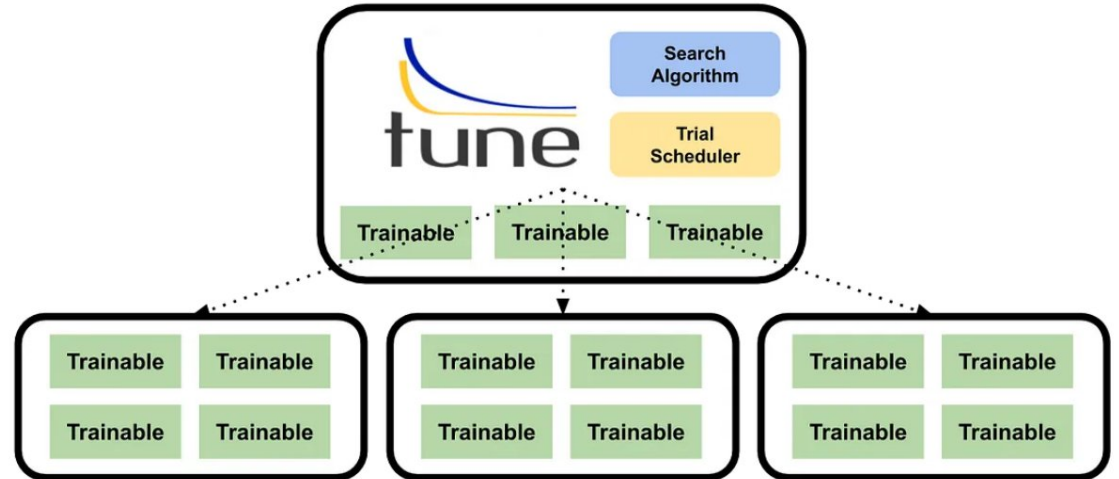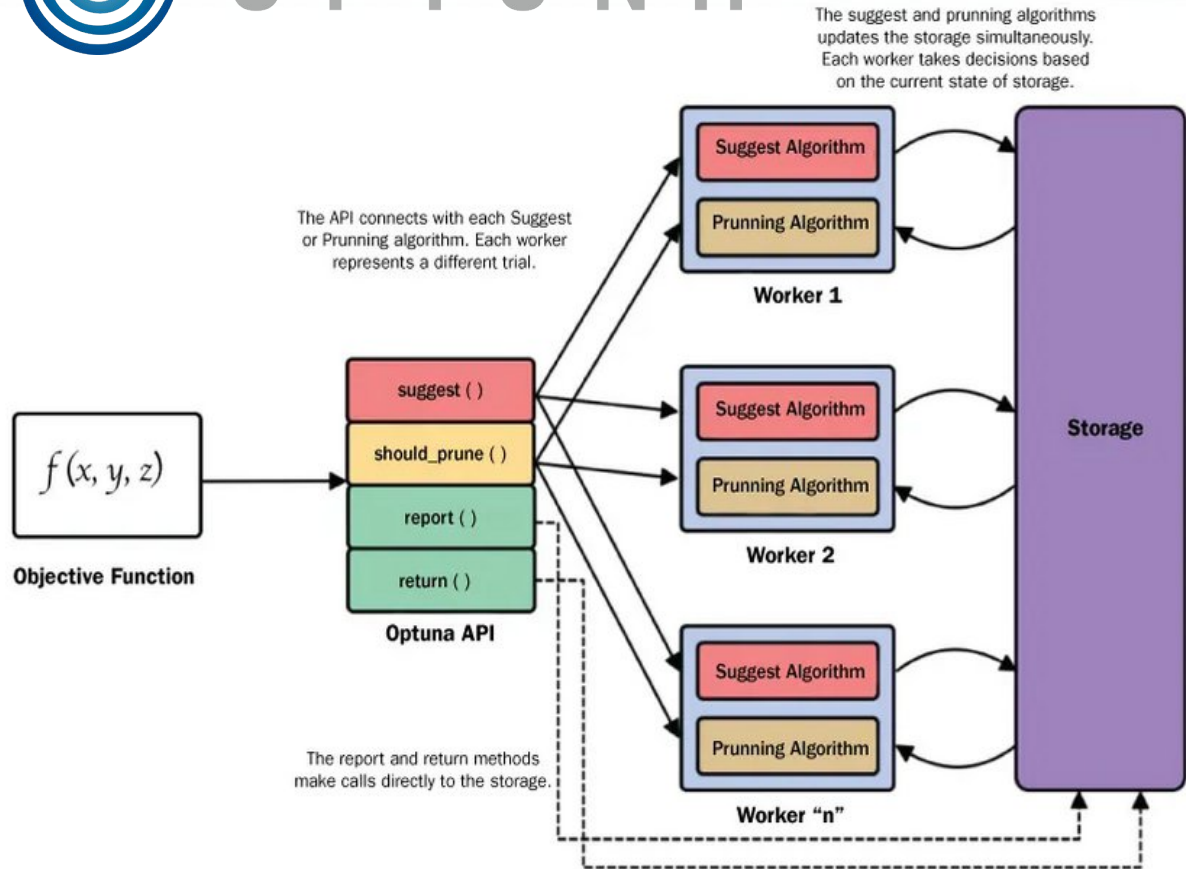Just-In Time Compilation

How to accelerate deep learning codes using *old school* compilation ?

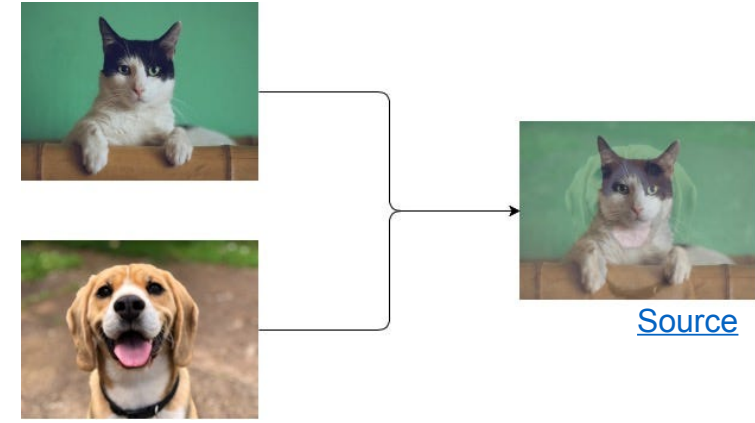**How to scale your HPO on supercomputer ?**
**How to use Optuna & Ray at scale ?**

**Is CPU enough?**
**How to delegate the Data Augmentation to the GPU?**
**How to optimize the Data Augmentation on the GPU?**

- **RandAugment** (torchvision): combinations of multiple transforms, either geometric or photometric, or both



Source

- **MixUp** (*mixup.py*): mixing up the features and their corresponding labels

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j, \qquad \text{where } x_i, x_j \text{ are raw input vectors}$$
$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j, \qquad \text{where } y_i, y_j \text{ are one-hot label encodings}$$



Source

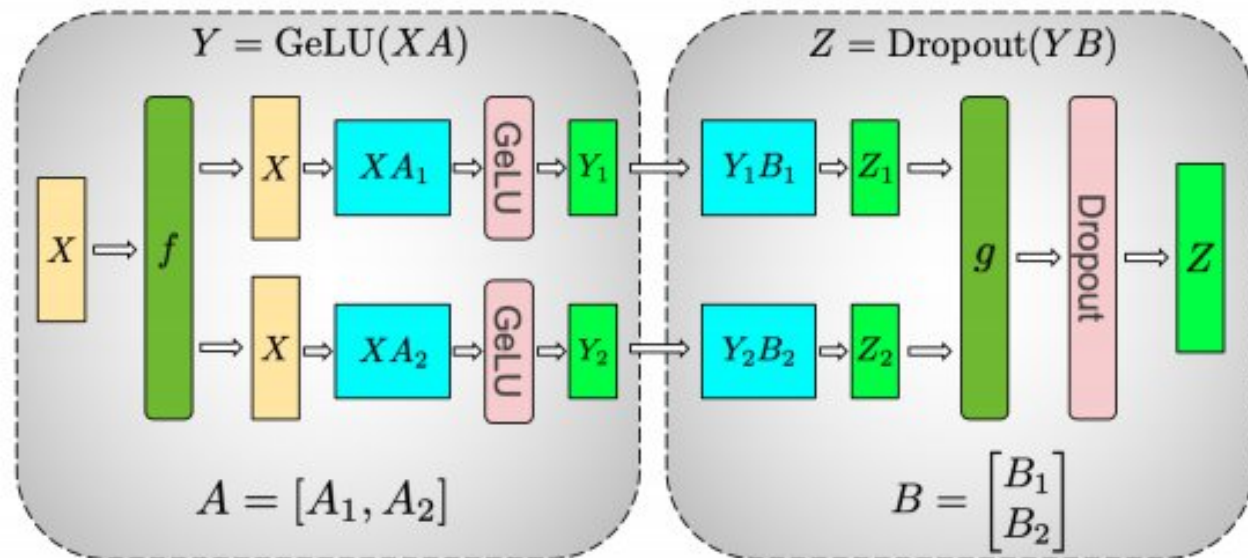- **CutMix** (*cutmix.py*): replacing an image region with a patch from another training image



Source

Pipeline Parallelism



Coatnet-7
2440M parameters

deepspeed

Optimized Data Parallelism
ZeRO - FSDP

# Conclusions

JIT ◄

Data Parallelism under the hood ◄

HPO ◄

Data Augmentation ◄

Large model ◄

Just-In Time Compilation
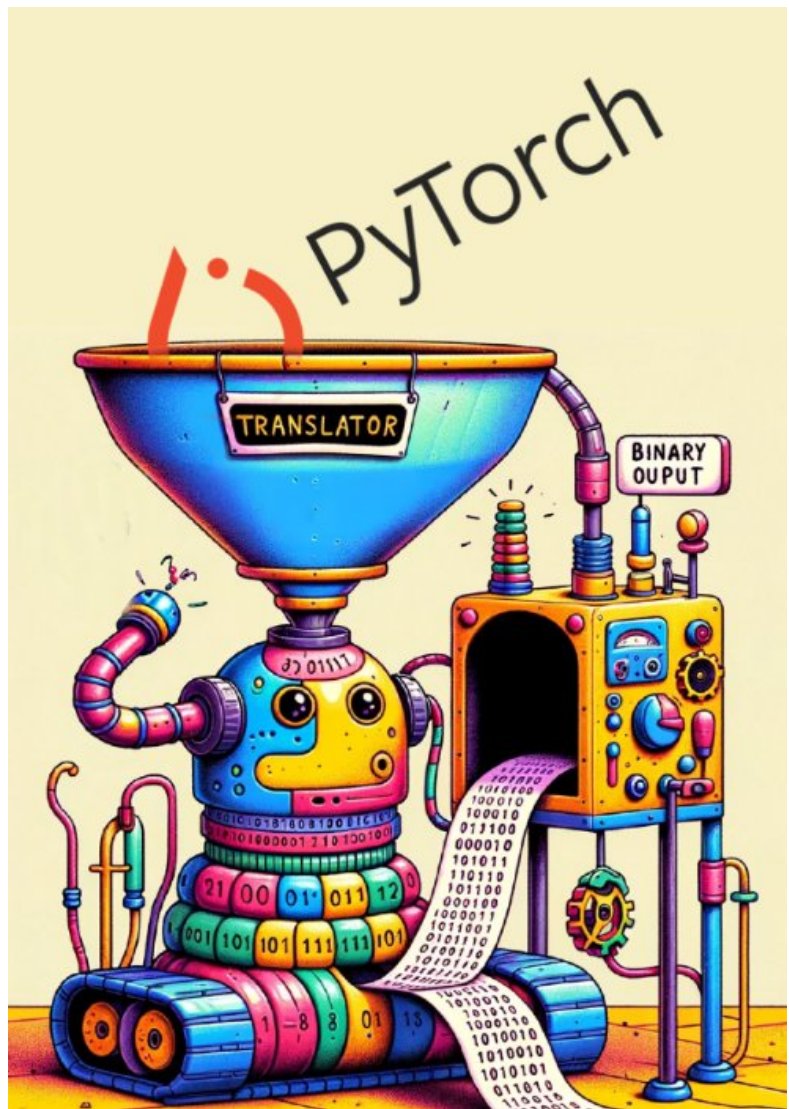
How to accelerate deep learning codes using *old school* compilation ?

```
model_opt = torch.compile(model, mode="reduce-overhead")
```

**How to scale your HPO on supercomputer ?**
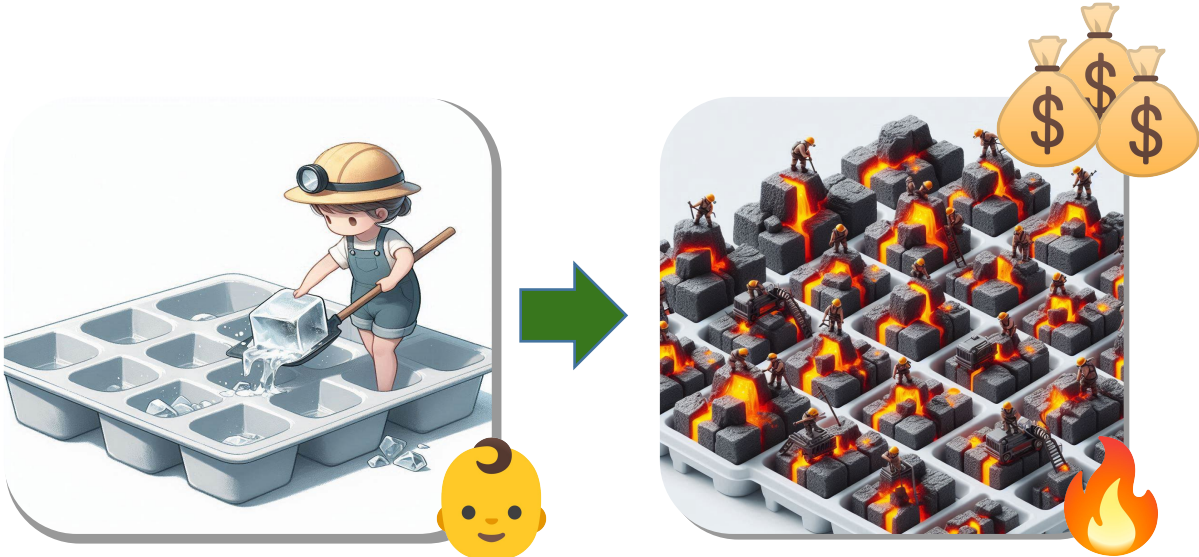**How to use Optuna & Ray at scale ?**



OPTUNA

- Database approach
- Not Multiprocess Native

RAY tune

- Multi-worker native
- Many parallel algorithms
- Difficult to use (especially with slurm)

**Is CPU enough?**
**How to delegate the Data Augmentation to the GPU?**
**How to optimize the Data Augmentation on the GPU?**
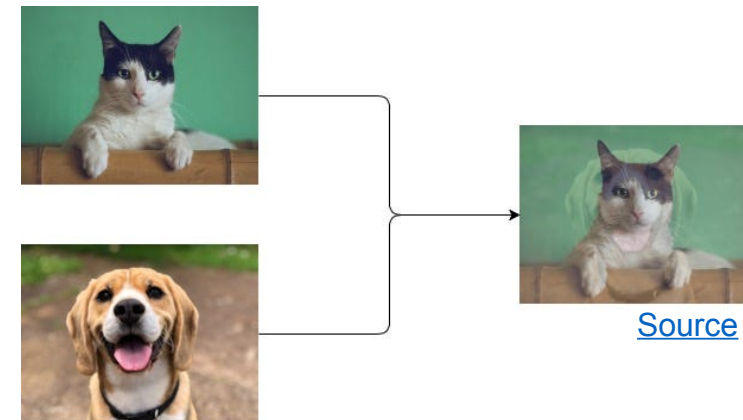
- **RandAugment** (torchvision): **CPU is enough**



Source

- **MixUp** (*mixup.py*): **much better on GPU**

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j, \qquad \text{where } x_i, x_j \text{ are raw input vectors}$$
$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j, \qquad \text{where } y_i, y_j \text{ are one-hot label encodings}$$



Source

- **CutMix** (*cutmix.py*): **tranform loop over images into loop over batches**
  **to improve parallelism and benifit from the GPU acceleration**
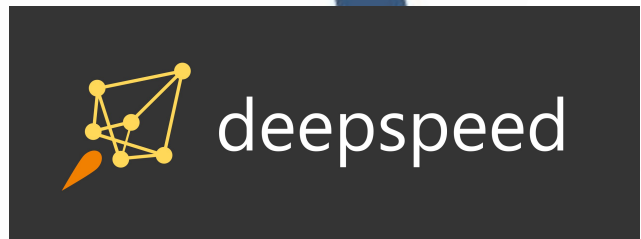


Source

11

Large model

Pipeline Parallelism from scratch
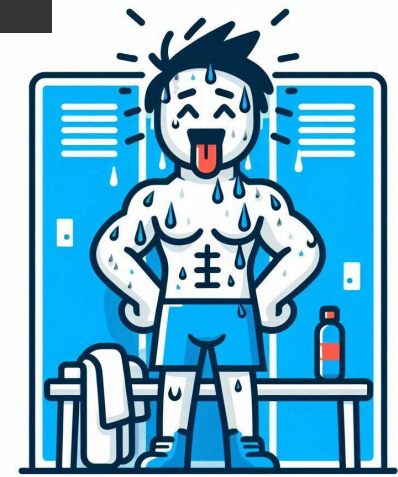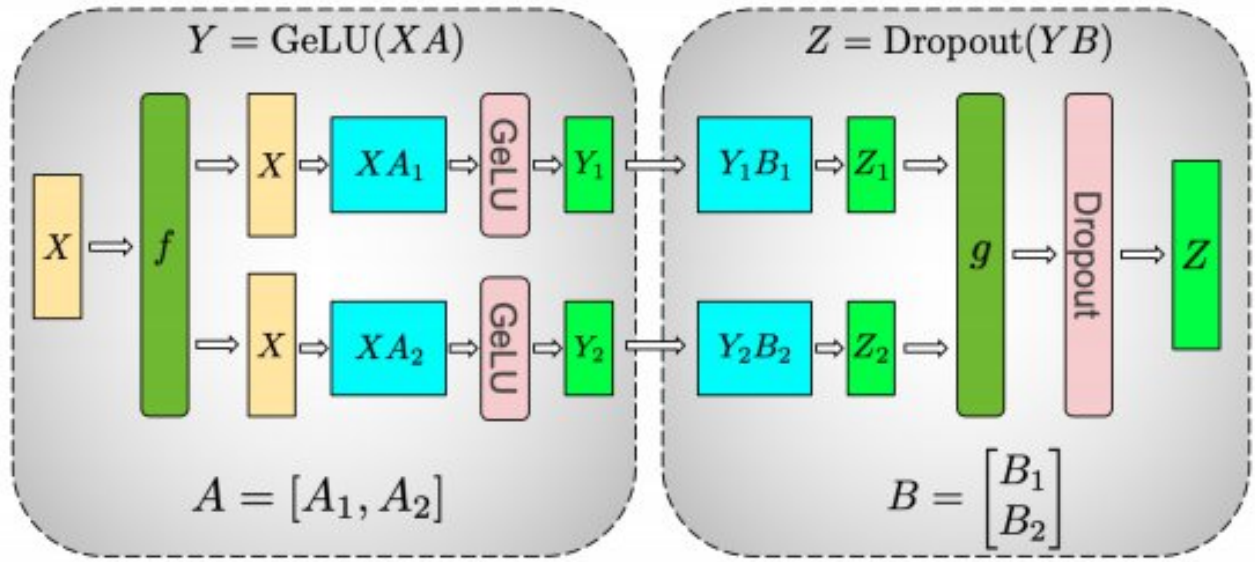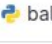
Pipeline Parallelism Deepspeed Implementation

PyTorch

deepspeed

FSDP

ZeRO

# Tensor Parallelism under the hood



$$Y = \text{GeLU}(XA)$$

$$A = [A_1, A_2]$$

$$Z = \text{Dropout}(YB)$$

$$B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

**tp tensor slicing**
Nathan CASSEREAU authored 3 months ago

| Name | Last commit |
|------|-------------|
| .. | |
| balise1.py | tp tensor slicing |
| balise10.py | tp tensor slicing |
| balise11.py | tp tensor slicing |
| balise2.py | tp tensor slicing |
| balise3.py | tp tensor slicing |
| balise4.py | tp tensor slicing |
| balise5.py | tp tensor slicing |
| balise6.py | tp tensor slicing |
| balise7.py | tp tensor slicing |
| balise8.py | tp tensor slicing |
| balise9.py | tp tensor slicing |