

Checkpointing

Jonas Eppelt, Matthias Schnepf, Günter Quast | August 28, 2023

ETP



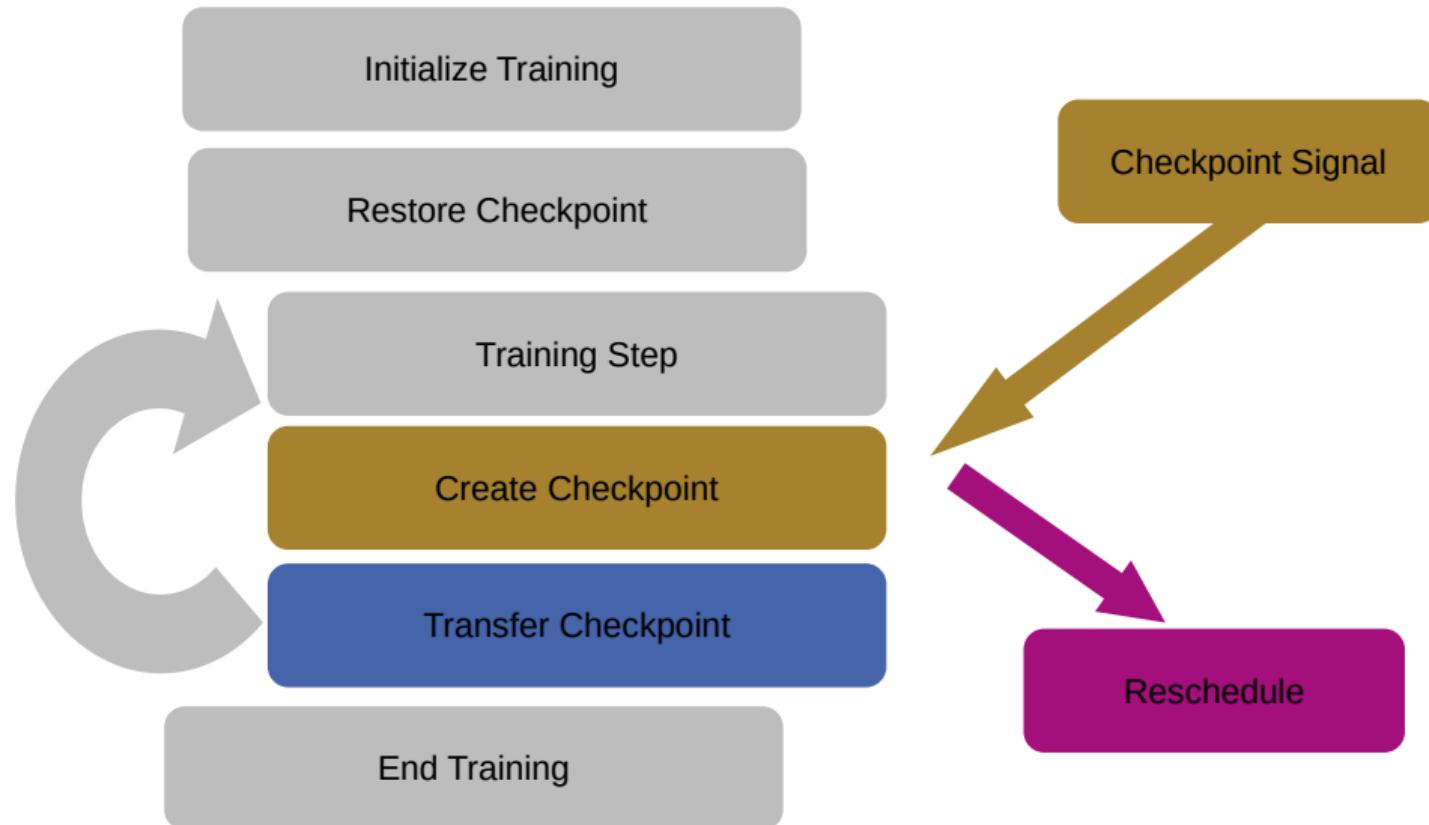
Motivation: Long-running ML Trainings

- not blocking resources for several days
- get intermediate output
- resistance to eviction/crashes of resources
- flexible usage of opportunistic resources → enables sustainability

Goals

- get it as approved as a service task (ST)
- develop self-checkpointing (training) jobs
- enable induced checkpointing by host (ST)
- provide examples for ETP, TOpAS, NAF (ST), KEK, Horeka
- provide convenience python class (ST)

Basic Training Program with checkpointing



How to Checkpoint: The ingredients

When to checkpoint:

- regularly (e.g. every 10 epochs)
- upon receiving a signal (SIGTERM, SIGINT, user-defined)

Where to store the checkpoint:

- shared filesystem (e.g. ETP)
- remote storage (e.g. xrootd)
- batch system (e.g. htcondor's transfer_checkpoint_files)

How to restart the job:

- manually (resubmit to the batch system, restart local process)
- batch system (htcondor's checkpoint_exit_code)
- Workflow Management (e.g. luigi)

Python class: almost 2 lines of code

```
1 from checkpointer.checkpointer import Checkpointer
2
3 checkpointer = Checkpointer(
4     # define checkpoint file
5     local_checkpoint_file=Path("checkpoint.pt"),
6     # define a function, that saves the checkpoint
7     checkpoint_function=lambda path, model: torch.save(model.state_dict(),
8             path),
9     # define a function, that restores the checkpoint
10    restore_function=lambda path: model.load_state_dict(torch.load(path)[
11        "model_state_dict"]),
12
13    # Reload checkpoint and give default value if there is none
14    model = checkpointer.restore(model)
15
16    # invoke the checkpoint_function
17    checkpointer.checkpoint(i)
```

When to checkpoint:

Regularly Checkpointing

```
1 checkpointer = Checkpointer(  
2     ...  
3     # define a frequency  
4     checkpoint_every = 100)  
5 ...  
6 for i in range(epochs):  
7     train(train_dataloader, model  
          , loss_fn, optimizer)  
8     test(test_dataloader, model,  
           loss_fn)  
9     checkpointer.step(i, model)
```

Induced Checkpointing

- python program sets traps on SIGTERM and SIGINT automatically
- To run on a batch system change the executable:

```
1 # push program to the  
2     background  
3 python3 counter.py &  
4 pid=$! # get pid  
5 # trap SIGINT & SIGTERM  
6 trap "kill -15 $pid" 15  
7 trap "kill -10 $pid" 10  
8 # wait for program to exit  
9 wait $pid
```

How to restart the job: HTCondor

- use the traps from induced checkpointing
- define a `checkpoint_exit_code` in the `jdl` file:

```
1 ...
2 # define an exit code on which condor copies the checkpoint
3 checkpoint_exit_code = 85
4 # only if also checkpoint transfer via htcondor
5 transfer_checkpoint_files = checkpoint.txt
6 when_to_transfer_output = ON_EXIT_OR_EVICT
```

- give the same code to the checkpointer:

```
1 checkpointer = Checkpointer(
2 ...
3     checkpoint_exit_code = 85 # default
4     # only if also checkpoint transfer via htcondor
5     checkpoint_transfer_mode="htcondor" # automatically infers
6         local_checkpoint_file from _CONDOR_JOB_AD
7 )
```

How to restart the job: Luigi

- Use checkpointer normally in the task.
- need to configure the luigi scheduler:

```
1 [scheduler]
2 retry_count=100
3 retry_delay=5
4 [worker]
5 keep_alive=true
6 max_reschedules=100
```

Where to store the checkpoints

Shared Filesystem

- set checkpoint_transfer_mode

```
1 checkpointer = Checkpointer(  
2     ...  
3     checkpoint_transfer_mode = "  
4         shared",  
5     checkpoint_transfer_target =  
6         Path("/work/jeppelt/  
7             checkpoint.txt")  
8 )
```

- First, the checkpoint is written locally.
- Then it is transferred to the target.

Remote storage

- set checkpoint_transfer_mode

```
1 checkpointer = Checkpointer(  
2     ...  
3     checkpoint_transfer_mode = "  
4         xrootd",  
5     checkpoint_transfer_target =  
6         "/pnfs/...",  
7     xrootd_server_name = "root  
8         ://..." )
```

- First, the checkpoint is written locally.
- Then it is transferred to the target.
- Don't forget your proxy!

Checkpointing in LSF

- has kernel-level checkpoints by default.
- works via additional scripts ‘echkpnt’ and ‘erestart’.
- is on application level only available in C and Fortran.

What is working so far

- transferring checkpoints to shared/remote filesystems → python programs can be checkpointed

Conclusion

Road to service Task

- see full examples in gitlab ✓
- provide examples with submitfiles for HTcondor systems (ETP, NAF) ✓
- demonstrate integration with workflow management system (b2luigi) ✓
- bring python prototype to real world application
 - used in my own project ✓
 - integrate in Flavourtagger pipeline

More Improvements

- provide example for KEK
- Bring to MonteCarlo Jobs
- enable integration to highlevel ML-frameworks (keras, lightning)
- think about benchmarking?