# GPU acceleration on MetaCentrum nodes

**Jiří Vorel**

CESNET, MetaCentrum, 12 December 2022

vorel@cesnet.cz

https://metacentrum.cz

https://metavo.metacentrum.cz

- ... is the National Grid Infrastructure (NGI) operated by CESNET (part of the e-INFRA CZ)

- ... is a provider of computational resources, application tools (commercial and free/open source) and data storage (for data in active use)

- ... is free of charge

  - Users 'pay' by Acknowledgement in their research publications

    https://wiki.metacentrum.cz/wiki/Usage_rules/Acknowledgement

- ... can be used only for non-commercial (academic) research

- ... is primarily dedicated to students and employees from Czech universities, the Czech Academy of Science, non-commercial research facilities etc., but we can grant access to foreign researchers and partners

- Computational resources are available to users immediately after the registration

- Individual jobs are scheduled and managed via the PBS Pro batch system

- MetaCentrum offers...

  - cca 45,000 CPU cores (x86_64)

  - SMP servers with up to 3 TB RAM, special servers with 6 and 10 TB RAM, small servers with up to 32 CPU, etc...

  - cca 400 various GPU cards (NVIDIA A10, A40, A100, RTX A4000, Tesla T4 etc.)

- Preferably CLI (Debian 11 and CentOS 7), also a GUI environment

- Before you start, read the documentation

  https://wiki.metacentrum.cz

  https://wiki.metacentrum.cz/wiki/Beginners_guide

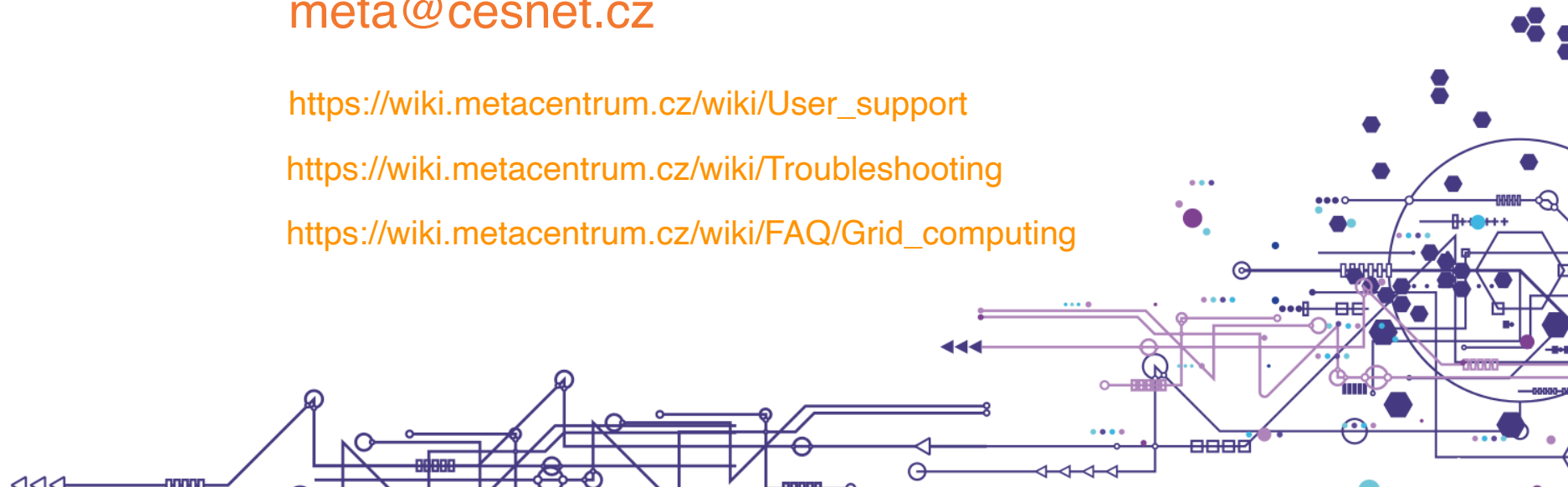  https://wiki.metacentrum.cz/wiki/Usage_rules

- If something goes wrong, do not hesitate to contact user support

  meta@cesnet.cz

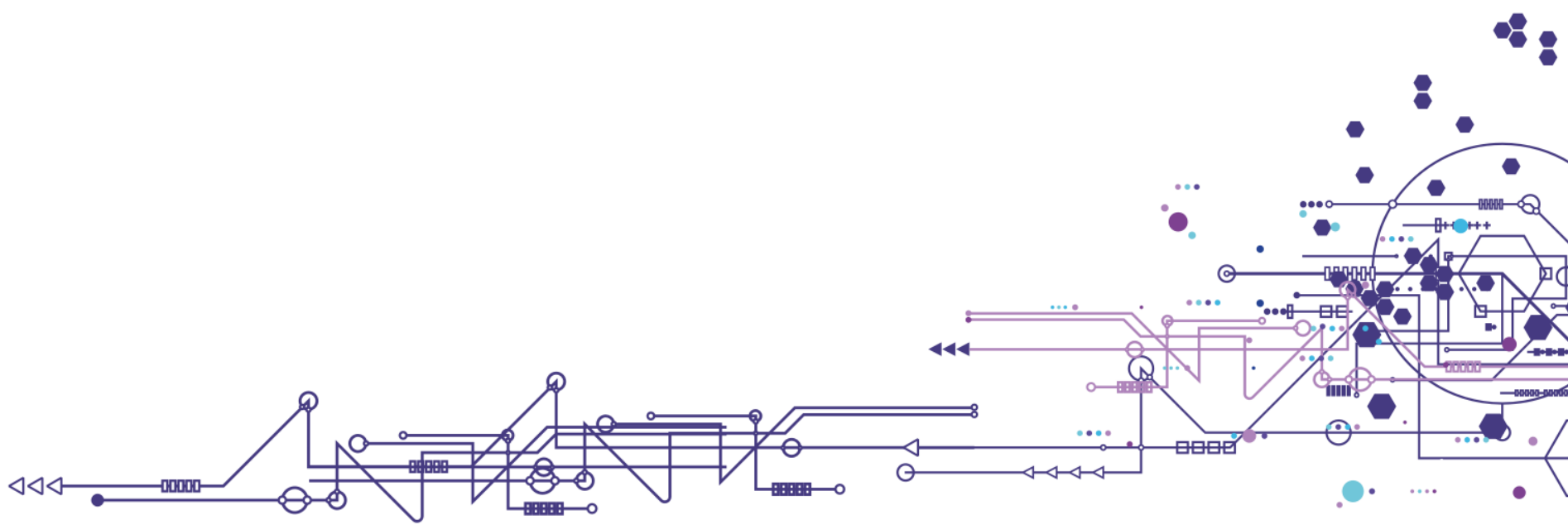  https://wiki.metacentrum.cz/wiki/User_support

  https://wiki.metacentrum.cz/wiki/Troubleshooting

  https://wiki.metacentrum.cz/wiki/FAQ/Grid_computing

# A brief introduction to how to use MetaCentrum with a preference for GPU nodes

# Frontend servers

- Gateway to the entire grid infrastructure (accessible via ssh with a password, no ssh tickets)

- Frontends submit jobs to PBS servers

- Frontends are small virtual machines mainly for purposes like writing scripts for batch jobs, checking applications and user data etc.

- **Do not run long and/or demanding calculations directly on frontends!**

- Frontend servers usually have different home directories

- All user home directories are available from all frontends

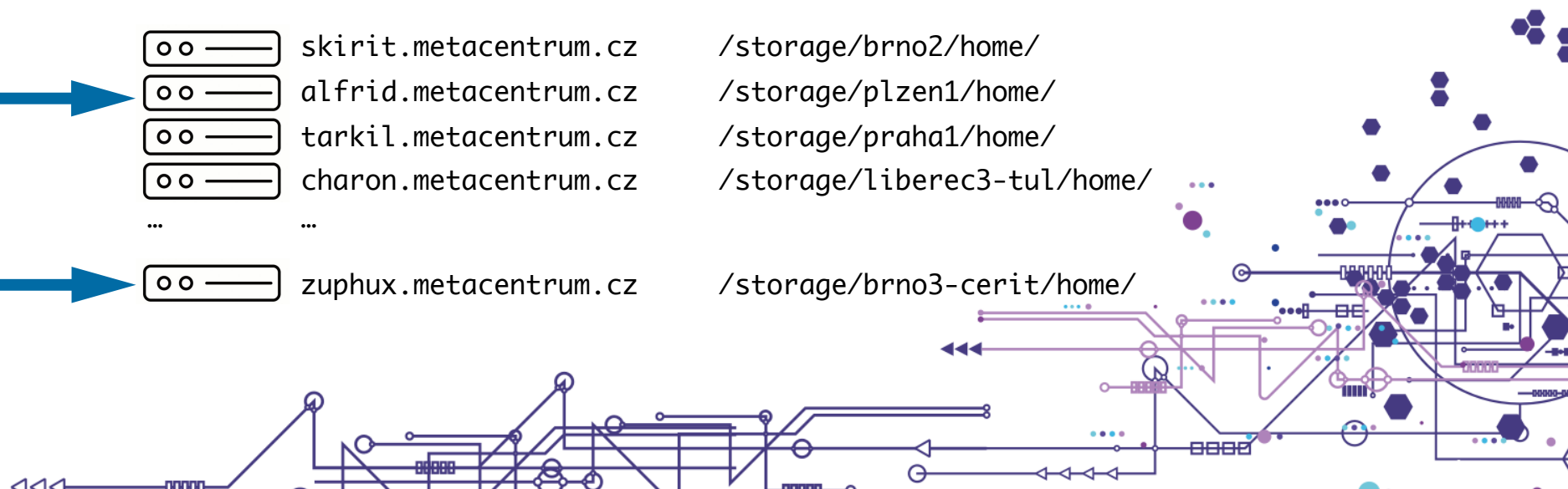| | | |
|---|---|---|
| meta-pbs.metacentrum.cz → | skirit.metacentrum.cz | /storage/brno2/home/ |
| | alfrid.metacentrum.cz | /storage/plzen1/home/ |
| | tarkil.metacentrum.cz | /storage/praha1/home/ |
| | charon.metacentrum.cz | /storage/liberec3-tul/home/ |
| | … … | |
| cerit-pbs.cerit-sc.cz → | zuphux.metacentrum.cz | /storage/brno3-cerit/home/ |

# Allocation of resources and qsub assembler

https://wiki.metacentrum.cz/wiki/About_scheduling_system

- Hardware resources (CPUs, GPUs, RAM, scratch, walltime,...) are reserved by PBS

- *qsub* command is used to submit jobs to the queue

- Users can use an interactive tool which assembles qsub command based on the selected criteria (requirements)



Go to metavo.metacentrum.cz - Current state - Personal view - **qsub assembler**

https://metavo.metacentrum.cz/pbsmon2/person

# GPU clusters

e-INFRA
CZ

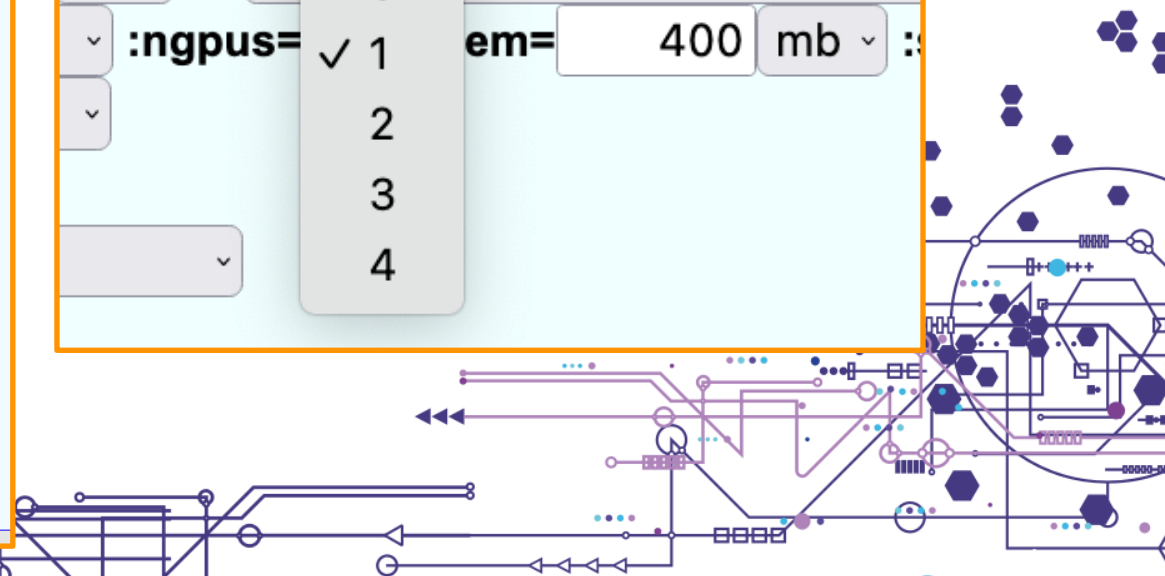| GPU clusters in MetaCentrum | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Cluster** | **Nodes** | **GPUs per node** | **Memory MiB** | **Compute Capability** | **CuDNN** | *gpu_cap=* | *cuda_version=* |
| galdor.metacentrum.cz | *galdor1.metacentrum.cz - galdor20.metacentrum.cz* | **4x A40** | 45 634 | 8.6 | YES | *cuda35,cuda61,cuda75,cuda80,cuda86* | *11.4* |
| luna2022.fzu.cz | *luna201.fzu.cz - luna206.fzu.cz* | **1x A40** | 45 634 | 8.6 | YES | *cuda35,cuda61,cuda75,cuda80,cuda86* | *11.4* |
| fer.natur.cuni.cz | *fer1.natur.cuni.cz - fer3.natur.cuni.cz* | **8x RTX A4000** | 16 117 | 8.6 | YES | *cuda35,cuda61,cuda75,cuda80,cuda86* | *11.2* |
| zefron.cerit-sc.cz | *zefron6.cerit-sc.cz* | **1x A10** | 22 731 | 8.6 | YES | *cuda35,cuda61,cuda75,cuda80,cuda86* | *11.2* |
| zia.cerit-sc.cz | *zia1.cerit-sc.cz - zia5.cerit-sc.cz* | **4x A100** | 40 536 | 8.0 | YES | *cuda35,cuda61,cuda75,cuda80* | *11.2* |
| fau.natur.cuni.cz | *fau1.natur.cuni.cz - fau3.natur.cuni.cz* | **8x Quadro RTX 5000** | 16 125 | 7.5 | YES | *cuda35,cuda61,cuda75* | *11.2* |
| cha.natur.cuni.cz | *cha.natur.cuni.cz* | **8x GeForce RTX 2080 Ti** | 11 019 | 7.5 | YES | *cuda35,cuda61,cuda75* | *11.2* |
| gita.cerit-sc.cz | *gita1.cerit-sc.cz - gita7.cerit-sc.cz* | **2x GeForce RTX 2080 Ti** | 11 019 | 7.5 | YES | *cuda35,cuda61,cuda75* | *11.2* |
| adan.grid.cesnet.cz | *adan1.grid.cesnet.cz - adan61.grid.cesnet.cz* | **2x Tesla T4** | 15 109 | 7.5 | YES | *cuda35,cuda61,cuda75* | *11.2* |
| glados.cerit-sc.cz | *glados2.cerit-sc.cz - glados7.cerit-sc.cz* | **2x GeForce RTX 2080** | 7 982 | 7.5 | YES | *cuda35,cuda61,cuda75* | *11.2* |
| glados.cerit-sc.cz | *glados1.cerit-sc.cz* | **1x TITAN V GPU** | 12 066 | 7.0 | YES | *cuda35,cuda61,cuda70* | *11.2* |
| konos.fav.zcu.cz | *konos1.fav.zcu.cz - konos8.fav.zcu.cz* | **4x GeForce GTX 1080 Ti** | 11 178 | 6.1 | YES | *cuda35,cuda61* | *11.2* |
| glados.cerit-sc.cz | *glados10.cerit-sc.cz - glados13.cerit-sc.cz* | **2x 1080Ti GPU** | 11 178 | 6.1 | YES | *cuda35,cuda61* | *11.2* |
| zefron.cerit-sc.cz | *zefron7.cerit-sc.cz* | **1x GeForce GTX 1070** | 8 119 | 3.5 | YES | *cuda35, cuda61* | *11.2* |
| black1.cerit-sc.cz | *black1.cerit-sc.cz* | **4x Tesla P100** | 16 280 | 6.0 | YES | *cuda35, cuda60* | *11.2* |
| grimbold.metacentrum.cz | *grimbold.metacentrum.cz* | **2x Tesla P100** | 12 198 | 6.0 | YES | *cuda35, cuda60* | *11.2* |
| zefron.cerit-sc.cz | *zefron8.cerit-sc.cz* | **1x Tesla K40c** | 11 441 | 3.5 | YES | *cuda35* | *11.2* |

# Qsub assembler for PBSPro

This page assist in assembling correct parameters for the qsub command that is used for submitting jobs in PBSPro planners.

Only computing resources avaliable to the user **vorel** are offered.

qsub -l walltime= `1` : `0` : `0` -q `default@meta-pbs.metacentrum.cz` \
   -l select= `1` :ncpus= `1` :ngpus= `0` :mem= `400` `mb` :scratch_ `local` = `400` `mb`

   cluster ... `[            ]`

   city ... `[            ]`

   SPECfp2017 per core ... `[            ]`

   other resources ...

     :arch= `[      ]`

     :biocev= `[      ]`

     :cgroups= `[      ]`

     :cluster= `[      ]`

     :cpu_flag= `[        ]`

     :cpu_vendor= `[      ]`

     :cuda_version= `[      ]`

     :gpu_cap= `[      ]`

     :host= `[        ]`

     :hyperthreading= `[      ]`

     :infiniband= `[       ]`

     :luna= `[     ]`

     :os= `[      ]`

     :osfamily= `[      ]`

     :pbs_server= `[         ]`

     :pruhonice= `[      ]`

     :scratch_shm= `[     ]`

     :vestec= `[     ]`

     :vnode= `[      ]`

`Find machines mathing the resource specification`

9

---

**Walltime dropdown (expanded):**

- 0
- 1
- 2
- 4
- ✓ 24
- 48
- 96
- 168
- 336
- 720
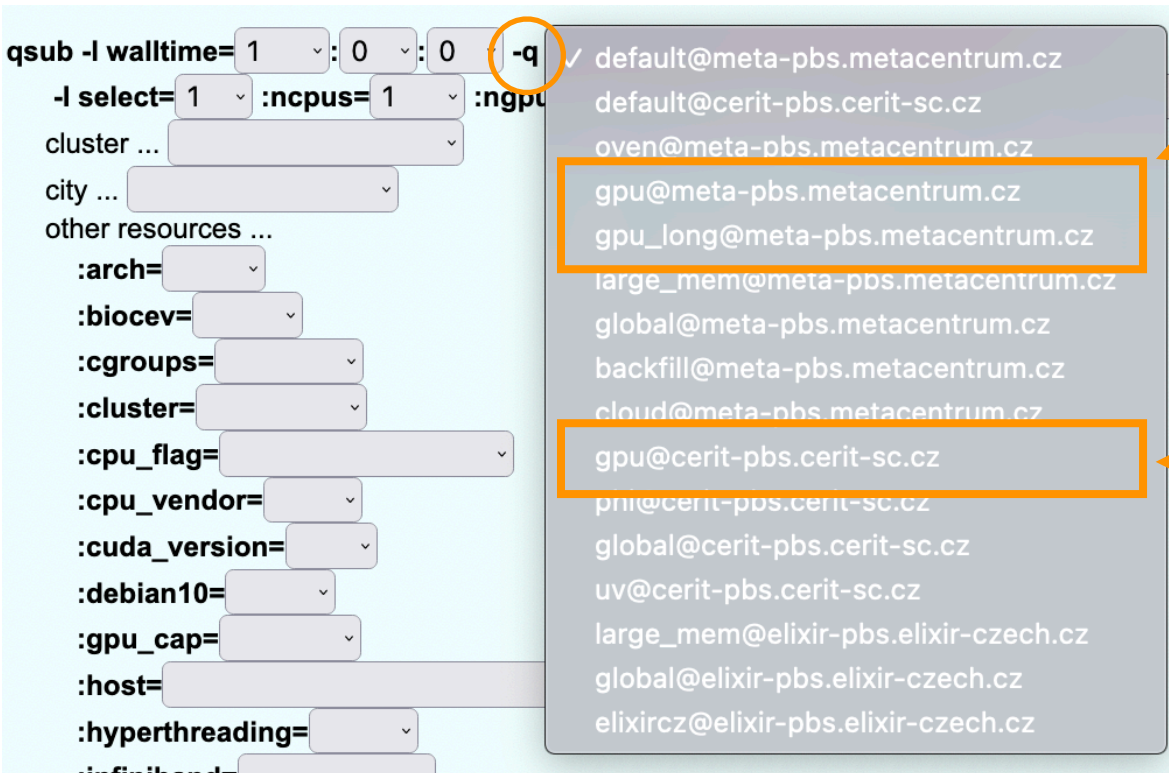
**ngpus dropdown (expanded):**

- 0
- ✓ 1
- 2
- 3
- 4

# Queues

- Not all visible queues are suitable for direct usage

- GPU calculations **must be** submitted to GPU queues
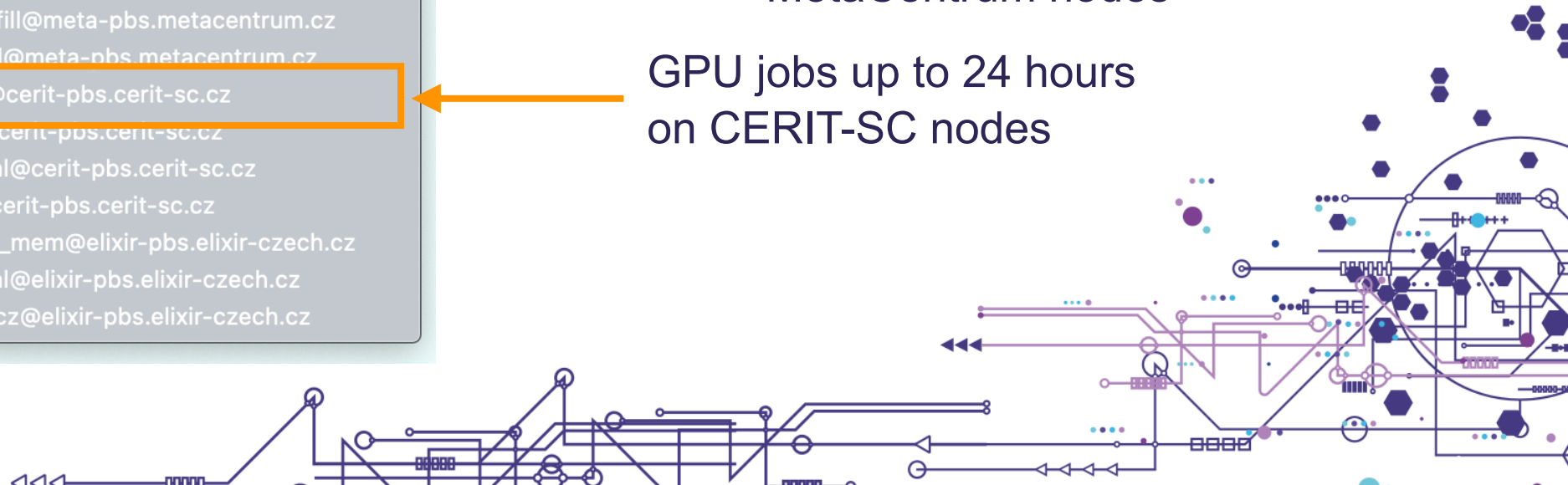
- Explore the *-q* option of the qsub assembler

GPU jobs with walltime up to 24 hours on MetaCentrum nodes

GPU jobs with walltime up to 336 hours on MetaCentrum nodes

GPU jobs up to 24 hours on CERIT-SC nodes

**qsub -l walltime=** `48` : `0` : `0` **-q** `gpu_long@meta-pbs.metacentrum.cz` \
**-l select=** `1` **:ncpus=** `1` **:ngpus=** `1` **:mem=** `20` `gb` **:scratch_** `local` **=** `200` `gb`

Find machines mathing the resource specification

**selection from command line**

qsub -l walltime=48:0:0 -q gpu_long@meta-pbs.metacentrum.cz -l select=1:ncpus=1:ngpus=1:mem=20gb:scratch_local=200gb

**selection in shell script**

```
#!/bin/bash
#PBS –q gpu_long@meta-pbs.metacentrum.cz
#PBS –l walltime=48:0:0
#PBS –l select=1:ncpus=1:ngpus=1:mem=20gb:scratch_local=200gb
#PBS –N my_awesome_job
```

**Result**

**OK**

The requirement is 1 machine, and 3 such machines are free, out of 19 machines matching the requirements. The job may be started immediately.
it.

**Machines available right now**

| adan4 (7 CPU, 5.9 SPEC, 164.6 GiB RAM, 710.6 GiB HDD) | galdor8 (7 CPU, 8.0 SPEC, 169.8 GiB RAM, 6.5 PiB HDD) | galdor10 (6 CPU, 8.0 SPEC, 455.7 GiB RAM, 6.5 PiB HDD) |
|---|---|---|

**All nodes matching the selection**

| adan1 (6 CPU, 5.9 SPEC, 174.6 GiB RAM, 568.5 GiB HDD) | adan2 | adan3 | adan4 (7 CPU, 5.9 SPEC, 164.6 GiB RAM, 710.6 GiB HDD) | adan5 |
|---|---|---|---|---|
| adan6 (7 CPU, 5.9 SPEC, 77.6 GiB RAM, 709.3 GiB HDD) | adan7 (7 CPU, 5.9 SPEC, 157.6 GiB RAM, 736.6 GiB HDD) | adan8 (5 CPU, 5.9 SPEC, 76.6 GiB RAM, 717.3 GiB HDD) | adan9 (6 CPU, 5.9 SPEC, 176.6 GiB RAM, 620.0 GiB HDD) | adan10 (7 CPU, 5.9 SPEC, 82.6 GiB RAM, 739.9 GiB HDD) |
| galdor1 (64 CPU, 8.0 SPEC, 503.8 GiB RAM, 6.5 PiB HDD) | galdor3 (6 CPU, 8.0 SPEC, 449.8 GiB RAM, 6.5 PiB HDD) | galdor4 (7 CPU, 8.0 SPEC, 136.8 GiB RAM, 6.5 PiB HDD) | galdor5 | galdor6 |
| galdor7 (4 CPU, 8.0 SPEC, 167.8 GiB RAM, 6.5 PiB HDD) | galdor8 (7 CPU, 8.0 SPEC, 169.8 GiB RAM, 6.5 PiB HDD) | galdor9 | galdor10 (6 CPU, 8.0 SPEC, 455.7 GiB RAM, 6.5 PiB HDD) | |

# Some hints for GPU reservation

- Each GPU calculation needs at least one CPU (ncpus=1)

- Remember that the newest GPU is NOT the best for all jobs

- GPU card can not be shared and is entirely dedicated to one calculation

- GPU calculations can be monitored on the same computation nodes by **nvidia-smi** command

- In most cases is not wise to target one specific cluster (e.g. **:cl_adan=True**), select a smaller set of machines using parameters:

  https://wiki.metacentrum.cz/wiki/GPU_clusters

  - **:gpu_mem=20gb**

  - **:gpu_cap=cuda80**

  - **:cuda_version=11.4**

# Example 1: Basecalling of ONT (Oxford Nanopore Technologies) reads in an interactive job

- Basecalling is a process how to determine individual nucleotides (DNA/RNA) from a characteristic electrical signal

- Requirements:

  - GPU card with at least 20 GB of memory

  - data processing toolkit Guppy with GPU support    https://wiki.metacentrum.cz/wiki/Guppy

  - input data in fast5 format (small data set for test)

- We will use an interactive job (calculation is waiting for individual commands typed by user)

- 1) Login to some frontend    https://wiki.metacentrum.cz/wiki/Beginners_guide#Run_interactive_job

```
jirivorel@MacBook ~$ ssh vorel@nympha.metacentrum.cz
vorel@nympha.metacentrum.cz's password:
Linux nympha.zcu.cz 5.10.0-13-amd64 #1 SMP Debian 5.10.106-1+zs1 (2022-03-28) x86_64
Last login: Thu Dec  8 19:03:00 2022 from dhcp17-232.ics.muni.cz

                                                                    cesnet
```

- 2) Check the availability of guppy software (via *module ava* command)

```
(BULLSEYE)vorel@nympha:~$ module ava guppy

------------------------------------------------------------ /packages/run/modules-2.0/modulefiles -------------------------------
guppy-3.0.3             guppy-3.6.0             guppy-4.5.4-gpu        guppy-6.0.1-cpu         guppy-6.0.6-cpu        guppy-6.3.8-gpu
guppy-3.4.5             guppy-4.4.1             guppy-5.0.15-cpu       guppy-6.0.1-gpu         guppy-6.0.6-gpu
guppy-3.5.1             guppy-4.5.4-cpu         guppy-5.0.15-gpu       guppy-6.0.1-gpu-singularity guppy-6.3.8-cpu
```

- 3) Start the interactive job with appropriate hardware resources and set the calculation

Start the interactive job instead of the regular batch job

```
(BULLSEYE)vorel@nympha:~$ qsub -I -l walltime=4:0:0 -q gpu@meta-pbs.metacentrum.cz -l select=1:ncpus=1:ngpus=1:mem=30gb:scratch_local=20gb:gpu_mem=20gb
qsub: waiting for job 13632463.meta-pbs.metacentrum.cz to start
qsub: job 13632463.meta-pbs.metacentrum.cz ready

(BULLSEYE)vorel@galdor4:~$ cd $SCRATCHDIR
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ cp -r /storage/praha5-elixir/home/vorel/ONT_input .
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ module add guppy-6.3.8-gpu
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ guppy_basecaller --version
: Guppy Basecalling Software, (C) Oxford Nanopore Technologies plc. Version 6.3.8+d9e0f64, minimap2 version 2.22-r1101
```

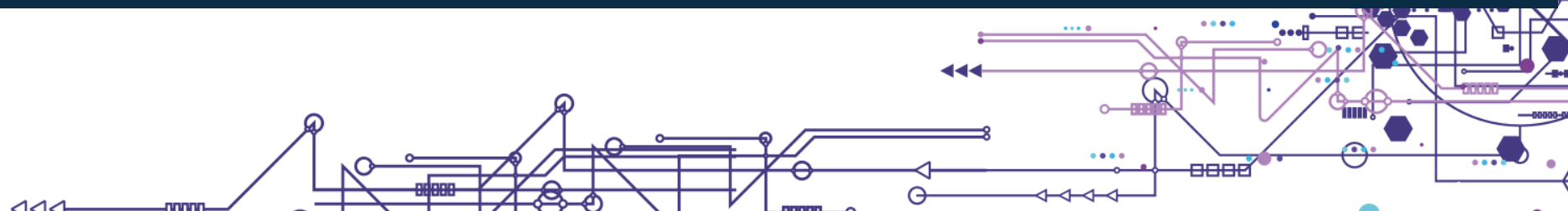Variable SCRATCHDIR is set automatically for each calculation

- ## 4) Run the calculation

```
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ guppy_basecaller -i ./ONT_input -r -s out_fastq_reads --flowcell FLO-MIN106 --kit SQK-LSK109 \
> -x auto --gpu_runners_per_device 16 --num_callers 16 --chunks_per_runner 2000 --trim_strategy none --disable_qscore_filtering
CRASHPAD MESSAGE:
ONT Guppy basecalling software version 6.3.8+d9e0f64, minimap2 version 2.22-r1101
config file:        /afs/ics.muni.cz/software/guppy/6.3.8-gpu/data/dna_r9.4.1_450bps_hac.cfg
model file:         /afs/ics.muni.cz/software/guppy/6.3.8-gpu/data/template_r9.4.1_450bps_hac.jsn
input path:         ./ONT_input
save path:          out_fastq_reads
chunk size:         2000
chunks per runner:  2000
records per file:   4000
num basecallers:    16
gpu device:         auto
kernel path:
runners per device: 16

Use of this software is permitted solely under the terms of the end user license agreement (EULA).By running, copying or accessing this software, you are demonstrating your acc
The EULA may be found in /afs/ics.muni.cz/software/guppy/6.3.8-gpu/bin
Found 2004 input read files to process.
Init time: 34337 ms

0%   10   20   30   40   50   60   70   80   90   100%
|----|----|----|----|----|----|----|----|----|----|
**************************************************
Caller time: 5003 ms, Samples called: 46685374, samples/s: 9.33148e+06
Finishing up any open output files.
Basecalling completed successfully.
```

- 5) In the meantime, when the calculation is running, you can open the second **e-INFRA CZ** terminal, login to the same node and check the GPU utilisation by *nvidia-smi* command

```
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 470.103.01   Driver Version: 470.103.01   CUDA Version: 11.4     |
|-------------------------------+----------------------+----------------------+
| GPU  Name        Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap|         Memory-Usage | GPU-Util  Compute M. |
|                               |                      |               MIG M. |
|===============================+======================+======================|
+-------------------------------+----------------------+----------------------+
|   3  NVIDIA A40          On   | 00000000:E1:00.0 Off |                    0 |
|  0%   48C    P0   267W / 300W |  25527MiB / 45634MiB |     98%      Default |
|                               |                      |                  N/A |
+-------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                                  |
|  GPU   GI   CI        PID   Type   Process name                  GPU Memory |
|        ID   ID                                                   Usage      |
|=============================================================================|
|    0   N/A  N/A      27676      C   ...x/DualSPHysics5.0_linux64     2001MiB |
|    1   N/A  N/A     728948      C   python                          17129MiB |
|    3   N/A  N/A     255837      C   guppy_basecaller                25523MiB |
+-----------------------------------------------------------------------------+
```

- 6) Check the result and clean everything



```
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ ls
ONT_input  guppy_basecaller-core-dump-db  out_fastq_reads
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ ls out_fastq_reads/
fastq_runid_78490aa79c827ee6f0554c0e8a22faedd299a6fb_0_0.fastq  guppy_basecaller_log-2022-12-12_00-35-55.log  guppy_basecaller_log-2022-12-12_00-38-28.log  sequencing_summary.txt
guppy_basecaller-core-dump-db                                   guppy_basecaller_log-2022-12-12_00-37-25.log  guppy_basecaller_log-2022-12-12_00-52-14.log  sequencing_telemetry.js
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ mv out_fastq_reads /storage/praha5-elixir/home/vorel/
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ rm -rf *
(BULLSEYE)vorel@galdor4:/scratch.ssd/vorel/job_13632463.meta-pbs.metacentrum.cz$ exit
logout

qsub: job 13632463.meta-pbs.metacentrum.cz completed
(BULLSEYE)vorel@nympha:~$
```

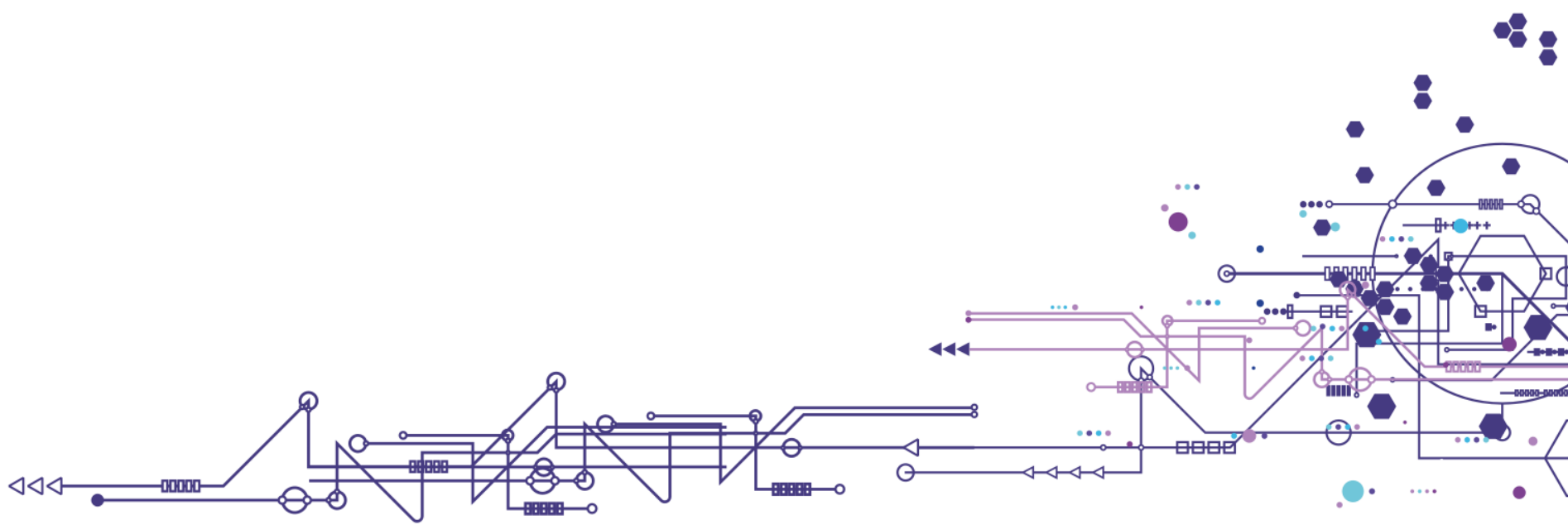Move only desired results back to the storage

Always remove everything unnecessary

Quit the interactive calculation

# Example 2: PyTorch MNIST training with Singularity container in batch job

- We will use the PyTorch Singularity image to train a MNIST model (Handwritten digit recognition)

- Requirements:

  - Basic test, no special HW requirements    https://wiki.metacentrum.cz/wiki/Singularity

  - Singularity

  - Torch

- We will use a batch job (all commands are in one shell script)

- 1) Login to some frontend    https://wiki.metacentrum.cz/wiki/Beginners_guide#Run_batch_jobs

- 2) Write a shell script for batch job

```bash
#!/bin/bash
#PBS -q gpu@meta-pbs.metacentrum.cz
#PBS -l walltime=1:0:0
#PBS -l select=1:ncpus=1:ngpus=1:mem=20gb:scratch_local=10gb:gpu_cap=cuda61
#PBS -N GPU_pytorch_test_job

# test if a scratch directory exists
test -n "$SCRATCHDIR" || { echo >&2 "Variable SCRATCHDIR is not set!"; exit 1; }

# move into the scratch directory
cd $SCRATCHDIR


# download test data
wget https://github.com/pytorch/examples/archive/refs/heads/master.zip
unzip master.zip
cd ./examples-main/word_language_model

# start the calculation
singularity exec --nv -B $SCRATCHDIR \
/cvmfs/singularity.metacentrum.cz/NGC/PyTorch\:22.10-py3.SIF \
python ./main.py --cuda --epochs 6

# clean the scratch automatically
clean_scratch
```
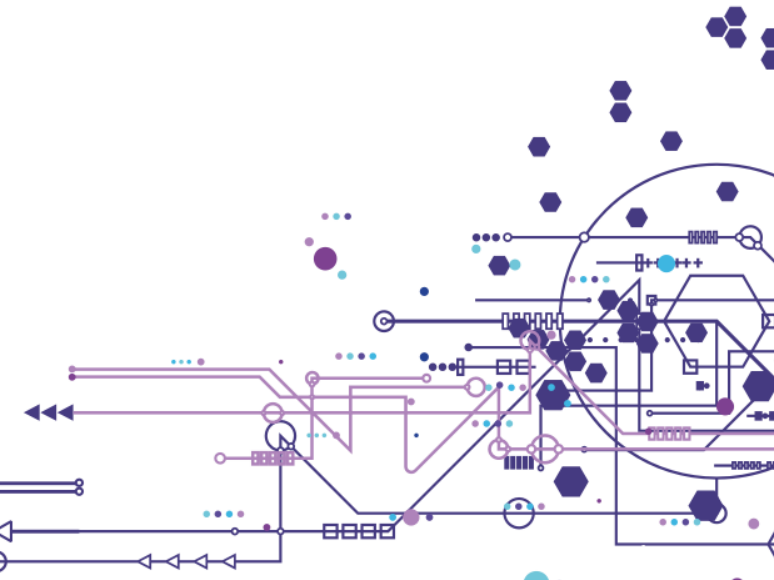
Let's specify a version of compute capability for demonstration purposes

Automatically remove data from the scratch directory

21

- 3) Submit the calculation and check logs

```
(BULLSEYE)vorel@nympha:/storage/praha1/home/vorel$ ls
GPU_pytorch.sh
(BULLSEYE)vorel@nympha:/storage/praha1/home/vorel$ qsub GPU_pytorch.sh
13633801.meta-pbs.metacentrum.cz
(BULLSEYE)vorel@nympha:/storage/praha1/home/vorel$ ls -l
celkem 33
-rwx------ 1 vorel meta   707 12. pro 09.26 GPU_pytorch.sh
-rw------- 1 vorel meta 11066 12. pro 09.28 GPU_pytorch_test_job.e13633801
-rw------- 1 vorel meta 13981 12. pro 09.28 GPU_pytorch_test_job.o13633801
(BULLSEYE)vorel@nympha:/storage/praha1/home/vorel$ less GPU_pytorch_test_job.e13633801
```

*qsub* command submits calculation to the PBS

Check standard outputs

```
| epoch   6 |   1200/ 2983 batches | lr 20.00 | ms/batch 15.58 | loss  4.68 | ppl   108.02
| epoch   6 |   1400/ 2983 batches | lr 20.00 | ms/batch 15.57 | loss  4.72 | ppl   112.40
| epoch   6 |   1600/ 2983 batches | lr 20.00 | ms/batch 15.56 | loss  4.80 | ppl   121.10
| epoch   6 |   1800/ 2983 batches | lr 20.00 | ms/batch 15.58 | loss  4.68 | ppl   107.95
| epoch   6 |   2000/ 2983 batches | lr 20.00 | ms/batch 15.56 | loss  4.72 | ppl   111.73
| epoch   6 |   2200/ 2983 batches | lr 20.00 | ms/batch 15.59 | loss  4.61 | ppl   100.78
| epoch   6 |   2400/ 2983 batches | lr 20.00 | ms/batch 15.58 | loss  4.65 | ppl   104.96
| epoch   6 |   2600/ 2983 batches | lr 20.00 | ms/batch 15.58 | loss  4.68 | ppl   107.83
| epoch   6 |   2800/ 2983 batches | lr 20.00 | ms/batch 15.59 | loss  4.62 | ppl   101.14
-----------------------------------------------------------------------------------------
| end of epoch   6 | time: 48.32s | valid loss  5.00 | valid ppl   148.28
-----------------------------------------------------------------------------------------
=========================================================================================
| End of training | test loss  4.94 | test ppl   139.93
=========================================================================================
```

Thank you for your attention

e-infra.cz