



Leveraging Big Data Tools on HPC with HOD

Kenneth Hoste | DICT | HPC-UGent
20161024

hpc@ugent.be

http://users.ugent.be/~kehoste/hod_20161024.pdf

Useful links

- HPC-UGent website: <http://hpc.ugent.be>
- HPC-UGent userwiki: <http://hpc.ugent.be/userwiki>
- VSC website: <https://www.vscentrum.be>
- **HPC-UGent support team contact: hpc@ugent.be**
- HOD documentation: <http://hod.readthedocs.org>
- HPC-UGent site-specific details on HOD:
<http://hpc.ugent.be/userwiki/index.php/Tips:Software:hanythingondemand>
- HOD code repository & issue tracker:
<https://github.com/hpcugent/hanythingondemand>
- HOD mailing list: <https://lists.ugent.be/wws/info/hod>

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[11.00am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] Requirements & installation

[10.30am] 'hod' command line interface

[10.45am] Creating and using HOD clusters

(lunch)

[1.00pm] Submitting batch scripts to an HOD cluster

[1.30pm] Connecting to web services of an HOD cluster

[2.00pm] Creating own HOD cluster configs

[2.30pm] Troubleshooting

[3.00pm] Hands-on: BYOW (bring your own workloads)

HPC vs Hadoop & co

HPC-UGent infrastructure

- multiple *generic* clusters
- fast interconnect, shared filesystem(s)
- access via job submission system



<https://www.vscentrum.be/infrastructure/hardware/hardware-ugent>

Traditional Hadoop clusters



- specialised in nature: only Hadoop (& co)
- HDFS as 'shared filesystem' on local disks
- 'direct' access through Hadoop commands

Why not a dedicated Hadoop cluster?



There are *no* dedicated Hadoop clusters in the HPC-UGent infrastructure, for a number of reasons:

- specialised setup can only be used for Hadoop & co
- requires significant expertise to administer and support
- unclear whether there is enough (consistent) demand
- does not scale dynamically with demand
- what about support for other 'services', like Jupyter notebooks?

hanythingondemand



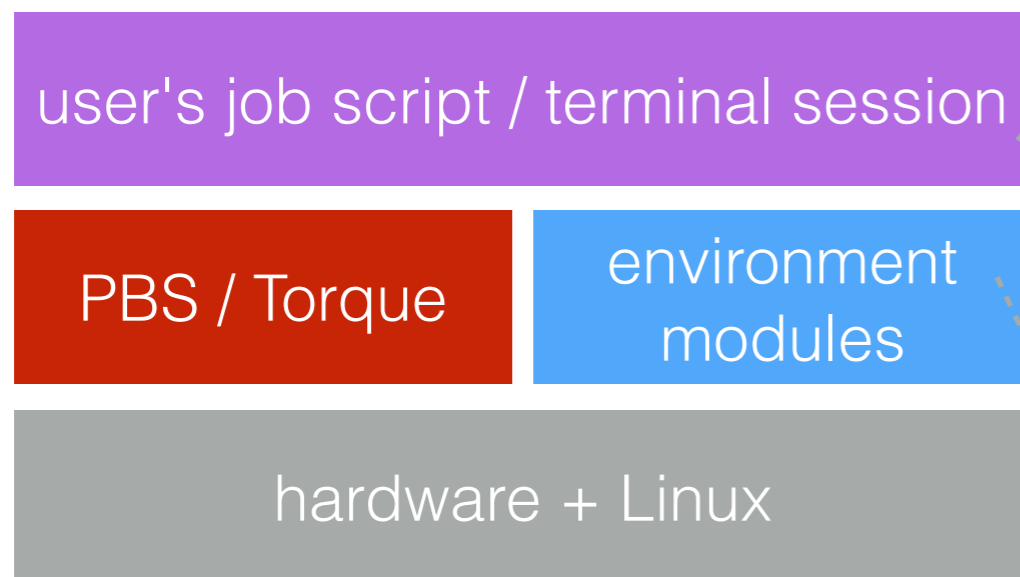
hanythingondemand (a.k.a. HOD) is a tool to set up and use an ad-hoc Hadoop cluster on HPC systems.

<http://hod.readthedocs.io/> - <https://github.com/hpcugent/hanythingondemand>

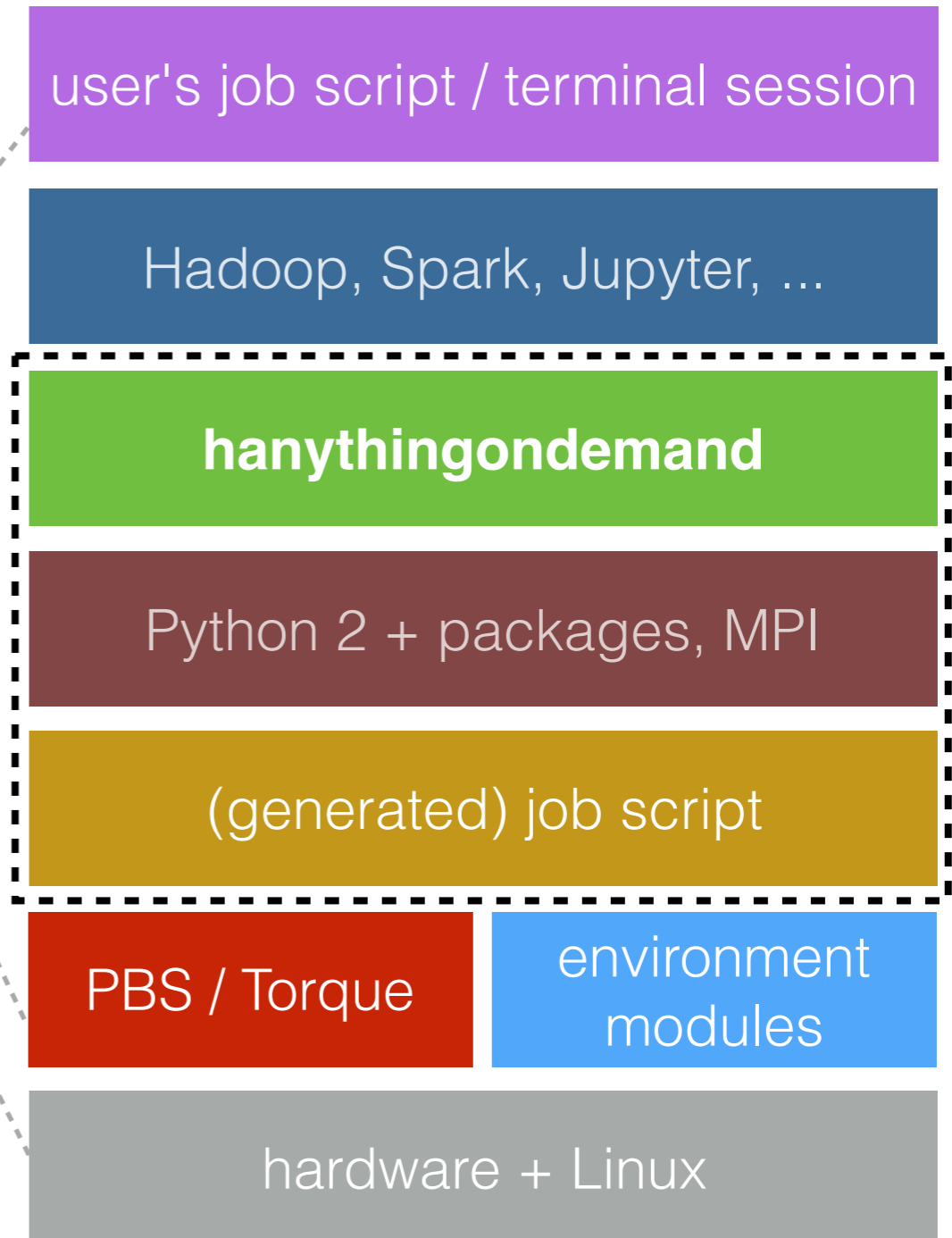
- focus on ease-of-use: hides complexity under the covers
- currently (only) compatible with PBS/Torque (HPC-UGent)
- can also be used for other services, e.g., Jupyter notebooks
- inspired by Hadoop On Demand; rewrite in Python 2
(used to be included with Hadoop: https://hadoop.apache.org/docs/r1.2.1/hod_scheduler.html)
- (kind of) similar to myHadoop (<https://github.com/glennklockwood/myhadoop>)

Bird's-eye view of HOD

traditional cluster:

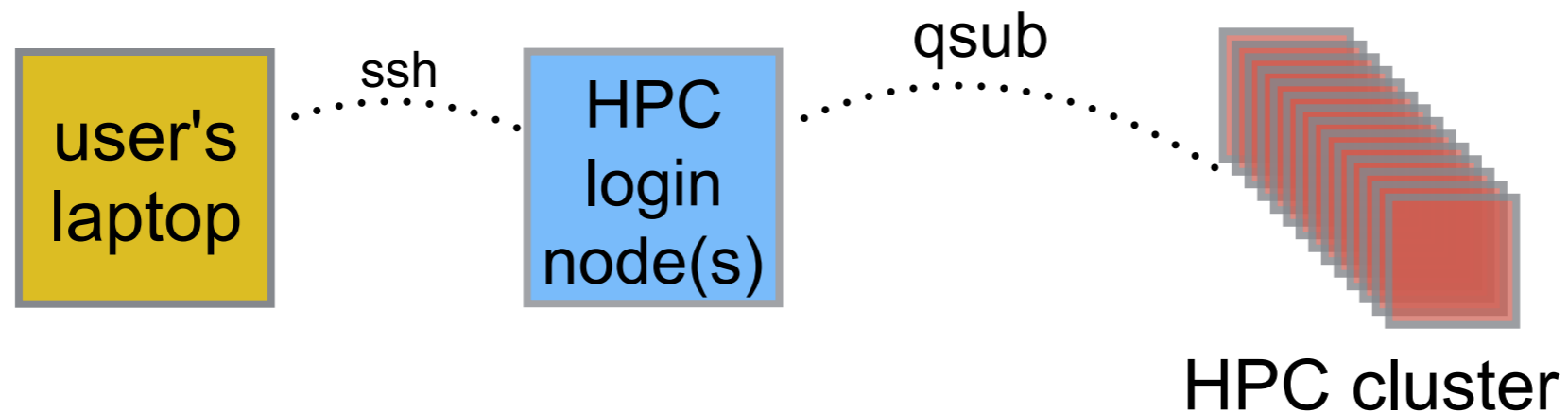


HOD on a cluster:

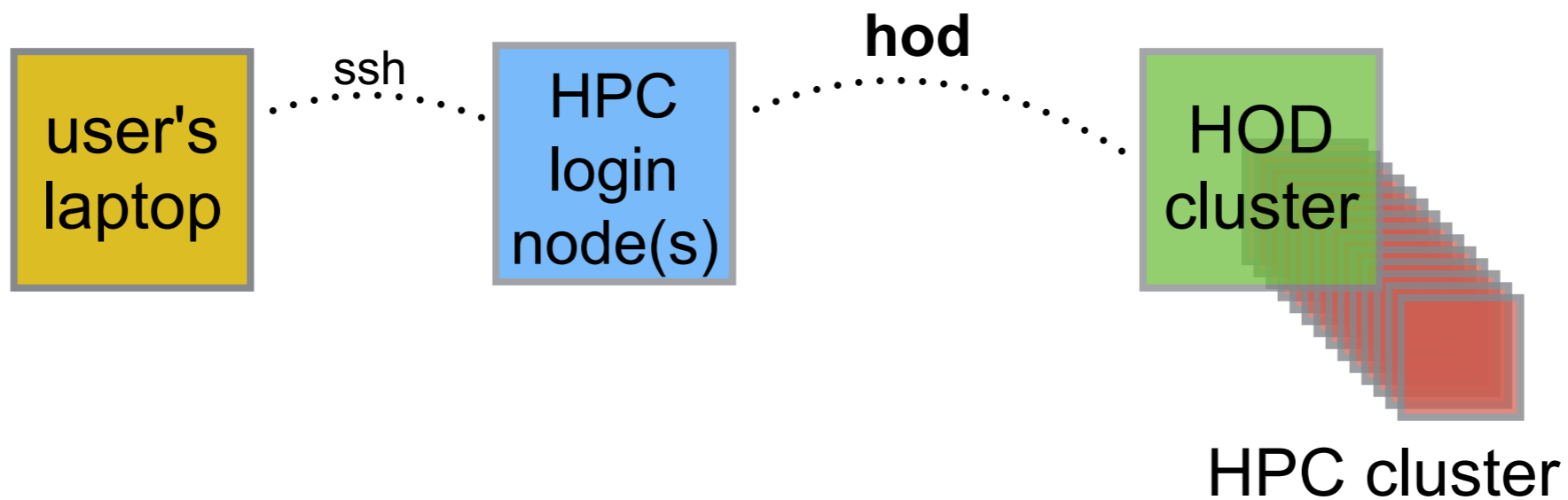


From a user's perspective

traditional cluster:



HOD on a cluster:



A little bit of HOD history

- original implementation by Stijn De Weirdt (*v2.0, May 2012*)
- major restructuring of codebase by Jens Timmerman (*v2.1, Jan 2014*)
- various enhancements by Ewan Higgs (*v2.2, Apr 2015*)
- major revision by Ewan Higgs & Kenneth Hoste (*v3.0, Oct 2015*)
 - redesigned/faster 'hod' CLI
 - support for 'hod batch', 'hod connect'
 - support for labelling HOD clusters
 - support for IPython notebooks
- support for more recent IPython/Jupyter versions (*v3.1, May 2016*)
- 'hod destroy', enhanced 'hod dists' (*v3.2, Oct 2016*)

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] 'hod' command line interface

[10.45am] Creating and using HOD clusters

(lunch)

[1.00pm] Submitting batch scripts to an HOD cluster

[1.30pm] Connecting to web services of an HOD cluster

[2.00pm] Creating own HOD cluster configs

[2.30pm] Troubleshooting

[3.00pm] Hands-on: BYOW (bring your own workloads)

HOD requirements

- HPC cluster running **Torque**
- **Python 2** (v2.6 or more recent)
- **MPI** library (e.g., OpenMPI, Intel MPI)
- Python packages:
netaddr, netifaces, mpi4py, vsc-base, vsc-mypirun, pbs_python
- environment modules tool (e.g., Lmod - <https://github.com/TACC/Lmod>)
- app modules installed using EasyBuild (<http://hpcugent.github.io/easybuild/>)
Hadoop, Spark, IPython/Jupyter, matplotlib, ...

HOD installation



- *highly* recommended to install HOD using EasyBuild
- same goes for modules that are required for Hadoop, IPython, etc.
- installation on HPC-UGent is split across *two* modules:
 - minimal installation provided by 'hod' module
 - full installation provided by 'hanythingondemand' module
 - <http://hpc.ugent.be/userwiki/index.php/Tips:Software:hanythingondemand>
- **long story short: just use 'module load hod' to get started**

```
$ module load hod
```

Hands-on

```
Last login: Thu Oct 20 10:43:35 2016 from example.ugent.be

STEVIN HPC-UGent infrastructure status on Thu, 20 Oct 2016 10:50:02

  cluster - full - free - part - total - running - queued
           nodes nodes free  nodes  jobs   jobs
-----
delcatty   0     0     0    160   N/A   N/A
golett     0     4     0    200   N/A   N/A
raichu     0     0     0     60   N/A   N/A
  muk     514     2     3    528   N/A   N/A
phanpy     2     4     0     16   N/A   N/A
swalot     0     1     0    128   N/A   N/A

For a full view of the current loads and queues see:
http://hpc.ugent.be/clusterstate/
Updates on maintenance and unscheduled downtime can be found on
https://www.vscentrum.be/en/user-portal/system-status

vsc40000@gligar01:~ $
```

- Log in to HPC-UGent login nodes.
- Which modules are available for HOD?
- Which modules are available for Hadoop, Spark, and IPython?
- Which other modules get loaded when loading the 'hod' module?

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

'hod' command line interface

http://hod.readthedocs.io/en/latest/Command_line_interface.html

- HOD command line interface is available via '**hod**' command
- several subcommands are supported

```
$ module load hod
$ hod
hanythingondemand version 3.2.0 - Run services within an HPC cluster
usage: hod <subcommand> [subcommand options]
Available subcommands (one of these must be specified!):
  batch          Submit a job to spawn a cluster on a PBS job controller,
                 run a job script, and tear down the cluster when it's done
  clean         Remove stale cluster info.
  clone         Write hod configs to a directory for editing purposes.
  connect       Connect to a hod cluster.
  create        Submit a job to spawn a cluster on a PBS job controller
  destroy       Destroy an HOD cluster.
  dists         List the available distributions
  genconfig     Write hod configs to a directory for diagnostic purposes
  help-template Print the values of the configuration templates based on
                 the current machine.
```


'hod' command line interface

http://hod.readthedocs.io/en/latest/Command_line_interface.html

- both global and subcommand-specific options are available
- some options are (semi-)mandatory, others are optional

```
$ hod create --help
Usage: hod create [options]

Options:
  --version          show program's version number and exit
  ...
  Create configuration:
  Configuration options for the 'create' subcommand
  (configfile section config)

  --dist=DIST        Prepackaged Hadoop distribution (e.g.
                    Hadoop/2.5.0-cdh5.3.1-native). This cannot be
                    set if --hodconf is set (type string)
  --hod-module=HOD-MODULE
                    Module to load for hanythingondemand in
                    submitted job (type string)
```

'hod' command line interface

http://hod.readthedocs.io/en/latest/Command_line_interface.html

- options for 'hod' subcommand can be specified either:
 - via configuration files
 - via environment variables (`$HOD_<SUBCMD>_<OPTION>`)
 - on the 'hod' command line
- CLI arguments override env. variables, which override conf. files

```
$ env | grep HOD  
HOD_CREATE_HOD_MODULE=hanythingondemand/3.2.0-intel-2016b-Python-2.7.12  
HOD_CREATE_WORKDIR=$VSC_SCRATCH/hod  
HOD_BATCH_WORKDIR=$VSC_SCRATCH/hod  
HOD_BATCH_HOD_MODULE=hanythingondemand/3.2.0-intel-2016b-Python-2.7.12  
  
$ hod create --workdir $VSC_SCRATCH_PHANPY/hod ...
```

Hands-on

```
Last login: Thu Oct 20 10:43:35 2016 from example.ugent.be

STEVIN HPC-UGent infrastructure status on Thu, 20 Oct 2016 10:50:02

  cluster - full - free - part - total - running - queued
           nodes nodes free  nodes  jobs   jobs
-----
delcatty   0     0     0    160   N/A   N/A
golett     0     4     0    200   N/A   N/A
raichu     0     0     0     60   N/A   N/A
  muk     514     2     3    528   N/A   N/A
phanpy     2     4     0     16   N/A   N/A
swalot     0     1     0    128   N/A   N/A

For a full view of the current loads and queues see:
http://hpc.ugent.be/clusterstate/
Updates on maintenance and unscheduled downtime can be found on
https://www.vscenrum.be/en/user-portal/system-status

vsc40000@gligar01:~ $
```

- Which global 'hod' options are mandatory?
- Will you need to specify these options *every* time you use 'hod'?
- Which global options are semi-mandatory for 'hod create' ?
- Which additional options are supported for 'hod batch' ?

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Creating an HOD cluster

http://hod.readthedocs.io/en/latest/Command_line_interface.html#hod-create

- to create an HOD cluster, use **'hod create'**
- you *must* specify either **'--dist'** or **'--hodconf'** (not both!)
- HOD clusters can be labeled for convenience (default label: job ID)

```
$ export HOD_CREATE_DIST=Hadoop-2.6.0-cdh5.4.5-native
```

```
$ env | grep HOD_CREATE
```

```
HOD_CREATE_DIST=Hadoop-2.6.0-cdh5.4.5-native
```

```
HOD_CREATE_HOD_MODULE=hanythingondemand/3.2.0-intel-2016b-Python-2.7.12
```

```
HOD_CREATE_WORKDIR=$VSC_SCRATCH/hod
```

```
$ hod create --label my_first_hod_cluster --job-walltime=1
```

```
Submitting HOD cluster with label 'my_first_hod_cluster'...
```

```
Job submitted: Jobid 12345.example state Q ehosts
```

Available HOD cluster configs

http://hod.readthedocs.io/en/latest/Example_use_cases.html#example-use-cases-common-available-dists

- HOD comes with a set of prepared cluster configs, a.k.a. 'dists'
- list of available cluster configs is printed by '**hod dists**'
- also mentions corresponding modules that will be loaded
- *note: used dist must be compatible with the value for --hod-module!*
(due to an unresolved bug in HOD v3.2, <https://github.com/hpcugent/hanythingondemand/issues/157>)

```
$ hod dists
* HBase-1.0.2
  modules: HBase/1.0.2, Hadoop/2.6.0-cdh5.4.5-native
...
* Hadoop-2.6.0-cdh5.4.5-native
  modules: Hadoop/2.6.0-cdh5.4.5-native
...
* Jupyter-notebook-5.1.0
  modules: Hadoop/2.6.0-cdh5.8.0-native, Spark/2.0.0,
           IPython/5.1.0-intel-2016b-Python-2.7.12,
           matplotlib/1.5.1-intel-2016b-Python-2.7.12
```

List of existing HOD clusters

http://hod.readthedocs.io/en/latest/Example_use_cases.html#example-hadoop-wordcount

- use '**hod list**' for an overview of existing HOD clusters
- also mentions state of cluster, and on which workernodes it is running
- HOD clusters that were terminated or submitted to a different system are also included

```
$ hod create --label hadoop_test --dist Hadoop-2.6.0-cdh5.4.5-native
Submitting HOD cluster with label 'hadoop_test'...
Job submitted: Jobid 12347.example state Q ehosts

$ hod list
Cluster label      Job ID           State           Hosts
2nd_test           12346.example    R               node1001.cluster
my_first_hod_cluster 12345.example    <job-not-found> <none>
hadoop_test        12347.example    Q
just_a_test        39592.other      <job-not-found> <none>
```

Connecting to an HOD cluster

http://hod.readthedocs.io/en/latest/Example_use_cases.html#example-hadoop-wordcount

- to connect to an HOD cluster, use '**hod connect <label>**'
- of course, the HOD cluster you specify must be *running*...
- once connected, you can work interactively using Hadoop, Spark, ...

```
$ hostname
login.hpc
$ hod connect example
Connecting to HOD cluster with label 'example'...
Job ID found: 123456.master.cluster
HOD cluster 'example' @ job ID 123456.master.cluster appears to be running...
Setting up SSH connection to node1001...
Welcome to your hanythingondemand cluster (label: example)

Relevant environment variables:
HADOOP_CONF_DIR=/scratch/me/hod/123456.master.cluster/me.node1001.5429/conf
HOD_LOCALWORKDIR=/scratch/me/hod/23456.master.cluster/me.node1001.5429
MODULEPATH=/scratch/modules/all:/etc/modulefiles

List of loaded modules:
Currently Loaded Modulefiles:
  1) Java/1.7.0_80                2) Hadoop/2.6.0-cdh5.4.5-native

$ hostname
node1001
```


Relabelling an HOD cluster

http://hod.readthedocs.io/en/latest/Command_line_interface.html#hod-relabel-old-label-new-label

- to relabel an existing cluster, use '**hod relabel <current> <new>**'
- useful to replace default label (job ID)
- works regardless of cluster state
- tip: always use a meaningful cluster label!

```
$ hod create --dist Hadoop-2.6.0-cdh5.4.5-native --job-walltime=1
Submitting HOD cluster with no label (job id will be used as a default label) ...
Job submitted: Jobid 54321.master.cluster state Q ehosts
```

```
$ hod list
```

Cluster label	Job ID	State	Hosts
54321.master.cluster	54321.master.cluster	R	node1001

```
$ hod relabel 54321.master.cluster pikachu
```

```
$ hod list
```

Cluster label	Job ID	State	Hosts
pikachu	54321.master.cluster	R	node1001

Destroying an HOD cluster

http://hod.readthedocs.io/en/latest/Command_line_interface.html#hod-destroy-cluster-label

- to get rid of an HOD cluster, use '**hod destroy <label>**'
- removes the job, cluster info directory and local working directory
- required confirmation for a *running* HOD cluster

```
$ hod destroy example
Destroying HOD cluster with label 'example'...
Job ID: 123456.master.cluster
Job status: R
Confirm destroying the *running* HOD cluster with label 'example'? [y/n]: y

Starting actual destruction of HOD cluster with label 'example'...

Job with ID 123456.master.cluster deleted.
Removed cluster localworkdir directory /scratch/me/hod/123456.master.cluster for
cluster labeled example
Removed cluster info directory /home/me/.config/hod.d/example for cluster labeled
example

HOD cluster with label 'example' (job ID: 123456.master.cluster) destroyed.
```

Cleaning up

http://hod.readthedocs.io/en/latest/Command_line_interface.html#hod-clean

- terminated HOD clusters will still be listed in the output of 'hod list'
- to clean up all *terminated* clusters at once, use '**hod clean**'
- only clusters submitted to the 'current' system will be cleaned up!

```
$ hod list
Cluster label      Job ID      State      Hosts
2nd_test           12346.example R          node1001.cluster
my_first_hod_cluster 12345.example <job-not-found> <none>
just_a_test        39592.other  <job-not-found> <none>

$ hod clean
...

$ hod list
Cluster label      Job ID      State      Hosts
2nd_test           12346.example R          node1001.cluster
just_a_test        39592.other  <job-not-found> <none>
```

Hands-on

```
Last login: Thu Oct 20 10:43:35 2016 from example.ugent.be

STEVIN HPC-UGent infrastructure status on Thu, 20 Oct 2016 10:50:02

  cluster - full - free - part - total - running - queued
           nodes nodes free  nodes  jobs   jobs
-----
delcatty   0     0     0    160   N/A   N/A
golett     0     4     0    200   N/A   N/A
raichu     0     0     0     60   N/A   N/A
  muk     514     2     3    528   N/A   N/A
phanpy     2     4     0     16   N/A   N/A
swalot     0     1     0    128   N/A   N/A

For a full view of the current loads and queues see:
http://hpc.ugent.be/clusterstate/
Updates on maintenance and unscheduled downtime can be found on
https://www.vscentrum.be/en/user-portal/system-status

vsc40000@gligar01:~ $
```

- Submit an HOD cluster for using HBase, *requesting 2h of walltime*.
- Make sure it starts.
- Relabel your (running) HOD cluster to 'hbase_test'.
- Connect to your HOD cluster.
- Run the Hadoop WordCount example.
http://hod.readthedocs.io/en/latest/Example_use_cases.html#interactively-using-a-hadoop-cluster
- Destroy the HOD cluster.

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Submitting batch scripts

http://hod.readthedocs.io/en/latest/Example_use_cases.html#running-a-batch-script-on-a-hadoop-cluster

- once you have defined your workflow, using a script is more efficient
- use '**hod batch**' to submit an HOD cluster that runs a script
- the HOD cluster will auto-destruct once the script is completed
- you *must* specify '**--script**', on top of the other mandatory options

```
$ export HOD_BATCH_DIST=Hadoop-2.6.0-cdh5.4.5-native

$ env | grep HOD_BATCH
HOD_BATCH_DIST=Hadoop-2.6.0-cdh5.4.5-native
HOD_BATCH_HOD_MODULE=hanythingondemand/3.2.0-intel-2016b-Python-2.7.12
HOD_BATCH_WORKDIR=$VSC_SCRATCH/hod

$ hod batch --script wordcount.sh --label wordcount --job-walltime=1
Submitting HOD cluster with label 'wordcount'...
Job submitted: Jobid 12345.example state Q ehosts
```

Hands-on

```
Last login: Thu Oct 20 10:43:35 2016 from example.ugent.be

STEVIN HPC-UGent infrastructure status on Thu, 20 Oct 2016 10:50:02

  cluster - full - free - part - total - running - queued
           nodes nodes free  nodes  jobs   jobs
-----
delcatty   0     0     0    160   N/A   N/A
golett     0     4     0    200   N/A   N/A
raichu     0     0     0     60   N/A   N/A
  muk     514     2     3    528   N/A   N/A
phanpy     2     4     0     16   N/A   N/A
swalot     0     1     0    128   N/A   N/A

For a full view of the current loads and queues see:
http://hpc.ugent.be/clusterstate/
Updates on maintenance and unscheduled downtime can be found on
https://www.vscentrum.be/en/user-portal/system-status

vsc40000@gligar01:~ $
```

- Create a script for running the Hadoop WordCount example.
- Submit an HOD cluster that runs that script (note: use short walltime!).
- Make sure you can access the result.

- **No cheating!**

http://hod.readthedocs.io/en/latest/Example_use_cases.html#running-a-batch-script-on-a-hadoop-cluster

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Connecting to web services

- Hadoop & co typically provide some web services
- likewise: IPython/Jupyter notebooks
- to connect to the web services of your HOD cluster, you need to jump through some hoops...
- SSH tunnel to direct traffic over HPC-UGent login node
- SOCKS proxy to make browser use SSH tunnel

http://hod.readthedocs.io/en/latest/Connecting_to_web_UIs.html



Spark Spark Worker at 10.168.193.41:41345

ID: worker-20140314184018-10.168.193.41-41345
Master URL: spark://10.160.137.165:7077
Cores: 2 (2 Used)
Memory: 6.3 GB (1024.0 MB Used)

[Back to Master](#)

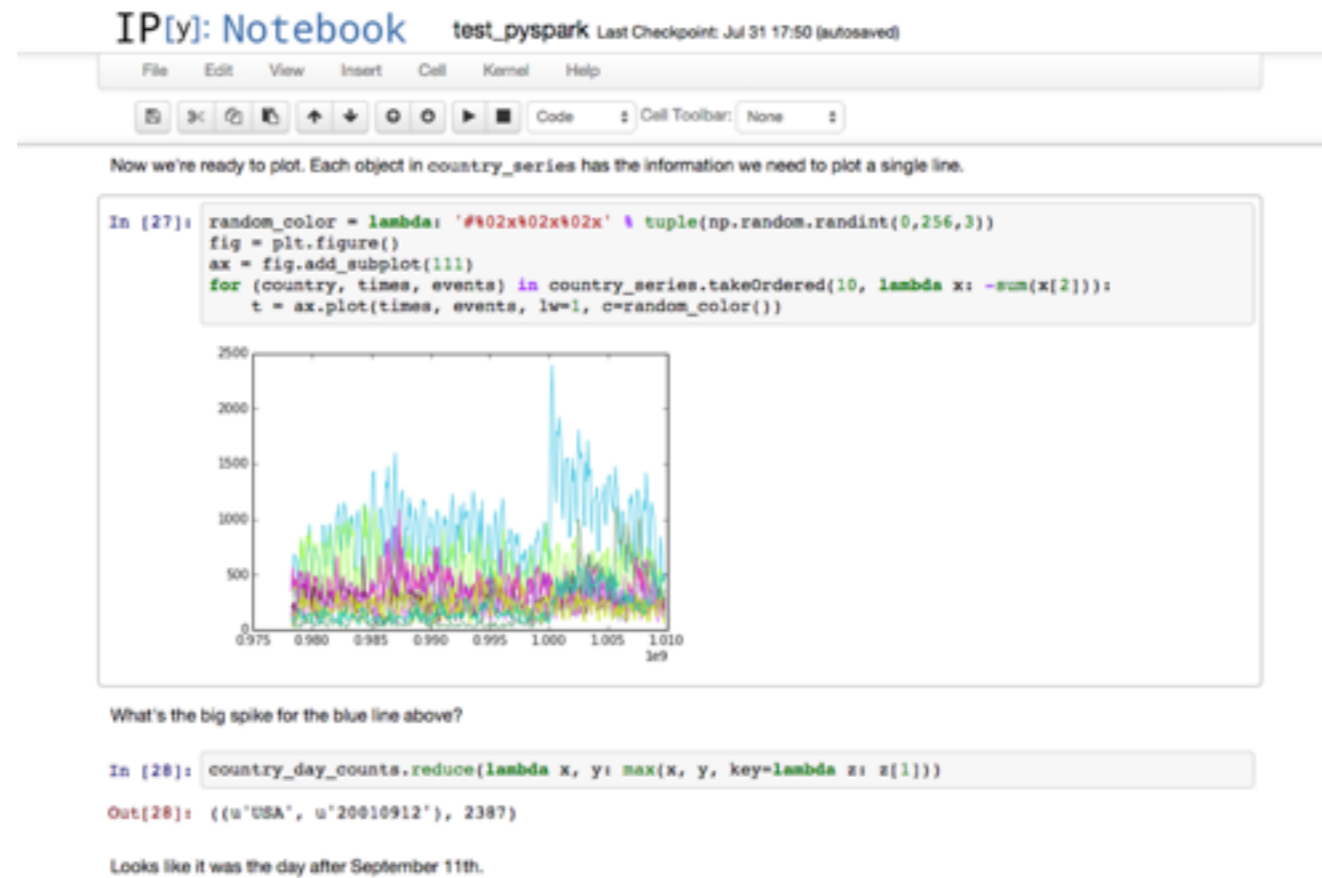
Running Executors 1

ExecutorID	Cores	Memory	Job Details	Logs
1	2	1024.0 MB	ID: app-20140314203709-0000 Name: @hark:ip-10-160-137-165 User: automaton	stdout stderr

Finished Executors

ExecutorID	Cores	Memory	Job Details	Logs
------------	-------	--------	-------------	------

Hands-on



- Submit an HOD cluster for running a Jupyter notebook.
- Set up an SSH tunnel to the hosting workernode.
- Set up a SOCKS proxy to direct your browser over the tunnel.
- Open the Jupyter notebook in your browser (<http://localhost:8888>).
- Destroy HOD cluster, undo putting SOCKS proxy in place.

http://hod.readthedocs.io/en/latest/Connecting_to_web_UIs.html

http://hod.readthedocs.io/en/latest/Example_use_cases.html#connecting-to-an-ipython-notebook-running-on-an-hod-cluster

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Creating own HOD cluster configs

<http://hod.readthedocs.io/en/latest/Configuration.html>

- the provided cluster configs included with HOD may not be sufficient
- you can compose your own cluster config and use that instead, via the '**--hodconf**' option for 'hod create' and 'hod batch'
- useful 'hod' subcommands:
 - '**hod clone**' to clone one of the provided 'dists' and start from that
 - '**hod genconfig**' for previewing the resulting cluster config files
 - '**hod help-template**' to show values for config templates

Hands-on

```
Last login: Thu Oct 20 10:43:35 2016 from example.ugent.be

STEVIN HPC-UGent infrastructure status on Thu, 20 Oct 2016 10:50:02

  cluster - full - free - part - total - running - queued
           nodes nodes free nodes jobs   jobs
-----
delcatty   0     0     0    160   N/A   N/A
golett     0     4     0    200   N/A   N/A
raichu     0     0     0     60   N/A   N/A
  muk     514     2     3    528   N/A   N/A
phanpy     2     4     0     16   N/A   N/A
swalot     0     1     0    128   N/A   N/A

For a full view of the current loads and queues see:
http://hpc.ugent.be/clusterstate/
Updates on maintenance and unscheduled downtime can be found on
https://www.vscentrum.be/en/user-portal/system-status

vsc40000@gligar01:~ $
```

- Compose a custom HOD cluster config for Hadoop 2.6 + Spark 2.
- Submit an HOD cluster using that `hod.conf`.
- Connect to the HOD cluster, check that Spark is available.
- Try to run a Spark example (<http://spark.apache.org/examples.html>).
- Destroy the HOD cluster.

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Troubleshooting

- every now and then, something may go wrong
- determining *what* went wrong (and how to fix it) may not be trivial...
- errors messages and logs are spread across multiple locations
- output in terminal, job script output, HOD logs, service logs, ...
- if you can't seem to figure it out: hpc@ugent.be, mention job ID

<http://hod.readthedocs.io/en/latest/Logging.html>

```
$ hod connect example
```

```
...
```

```
$ ls $HOD_LOCALWORKDIR/log
```

```
total 8192
```

```
drwxr-xr-x 2 vsc40000 vsc40000 512 Oct 21 17:41 userlogs
```

```
-rw-rw-r-- 1 vsc40023 vsc40023 43983 Oct 21 17:41 yarn-vsc40023-nodemanager
```

```
-rw-rw-r-- 1 vsc40023 vsc40023 51567 Oct 21 17:41 yarn-vsc40023-resourceman
```

Hands-on

```
Last login: Thu Oct 20 10:43:35 2016 from example.ugent.be

STEVIN HPC-UGent infrastructure status on Thu, 20 Oct 2016 10:50:02

  cluster - full - free - part - total - running - queued
           nodes nodes free  nodes  jobs   jobs
-----
delcatty   0     0     0    160   N/A   N/A
golett     0     4     0    200   N/A   N/A
raichu     0     0     0     60   N/A   N/A
  muk     514     2     3    528   N/A   N/A
phanpy     2     4     0     16   N/A   N/A
swalot     0     1     0    128   N/A   N/A

For a full view of the current loads and queues see:
http://hpc.ugent.be/clusterstate/
Updates on maintenance and unscheduled downtime can be found on
https://www.vscentrum.be/en/user-portal/system-status

vsc40000@gligar01:~ $
```

- Submit an HOD cluster using the following command:

```
hod create --dist Jupyter-notebook-5.1.0 --label pandas
           --modules pandas/0.18.1-intel-2016b-Python-3.5.2
           --job-walltime=1
```

- Figure out why this doesn't work as you may expect.

Outline

[10.05am] **HOD: what, why, how?**

[10.20am] **Requirements & installation**

[10.30am] **'hod' command line interface**

[10.45am] **Creating and using HOD clusters**

(lunch)

[1.00pm] **Submitting batch scripts to an HOD cluster**

[1.30pm] **Connecting to web services of an HOD cluster**

[2.00pm] **Creating own HOD cluster configs**

[2.30pm] **Troubleshooting**

[3.00pm] **Hands-on: BYOW (bring your own workloads)**

Hands-on: BYOW

- To finish up, try creating an HOD cluster for *your* workload...
- or try running the WordCount example with Spark
- or try using pyspark or scikit-learn in a Jupyter notebook
- or ...



Useful links

- HPC-UGent website: <http://hpc.ugent.be>
- HPC-UGent userwiki: <http://hpc.ugent.be/userwiki>
- VSC website: <https://www.vscentrum.be>
- **HPC-UGent support team contact: hpc@ugent.be**
- HOD documentation: <http://hod.readthedocs.org>
- HPC-UGent site-specific details on HOD:
<http://hpc.ugent.be/userwiki/index.php/Tips:Software:hanythingondemand>
- HOD code repository & issue tracker:
<https://github.com/hpcugent/hanythingondemand>
- HOD mailing list: <https://lists.ugent.be/wws/info/hod>

“Big Data sucks.”

–Robert McLay (TACC)