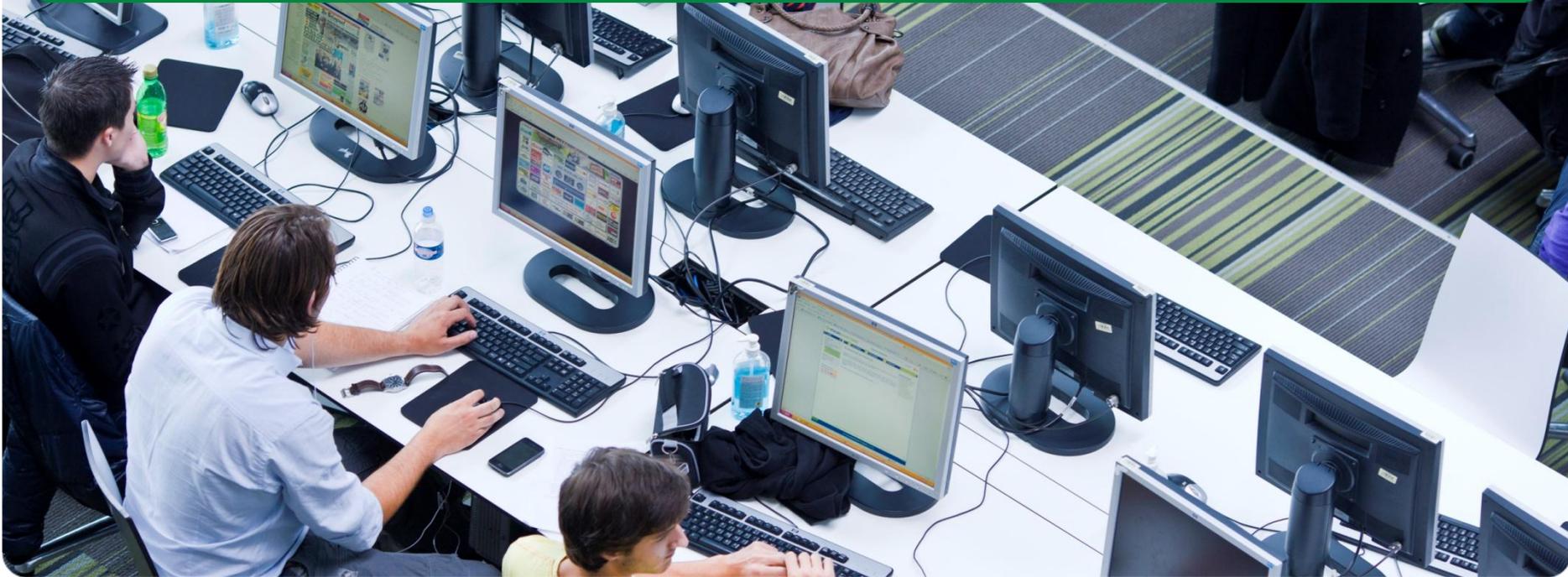


# Hadoop @ SURFsara

## USING THE CLUSTER



Jeroen Schot <[jeroen.schot@surfsara.nl](mailto:jeroen.schot@surfsara.nl)>



# Overview

- SURFsara in a nutshell
- The SURFsara Hadoop cluster
- **How to use the cluster**

# About SURF



**SURF SARA**

High-performance computing, data and visualisation for science



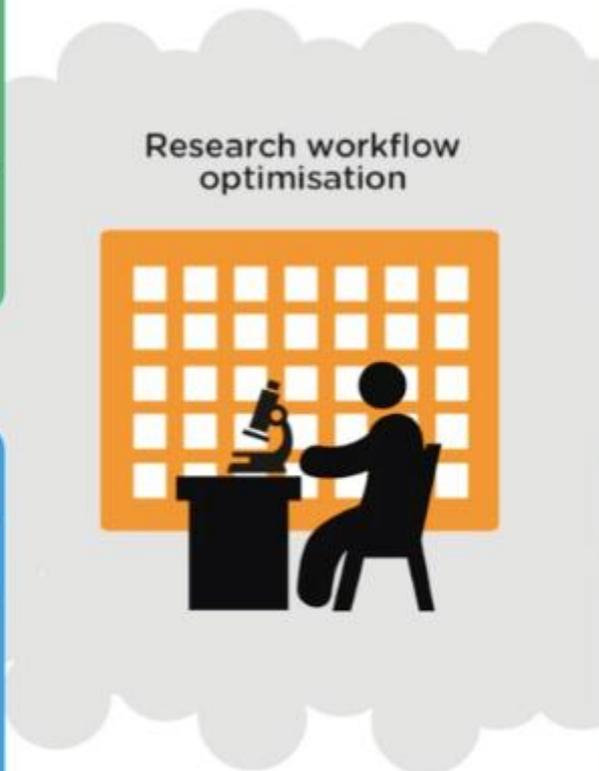
**SURF NET**

Connects users and ICT services and creates new functional possibilities



netherlands  
**eScience center**

Reinforces and accelerates multi-disciplinary and data-intensive research

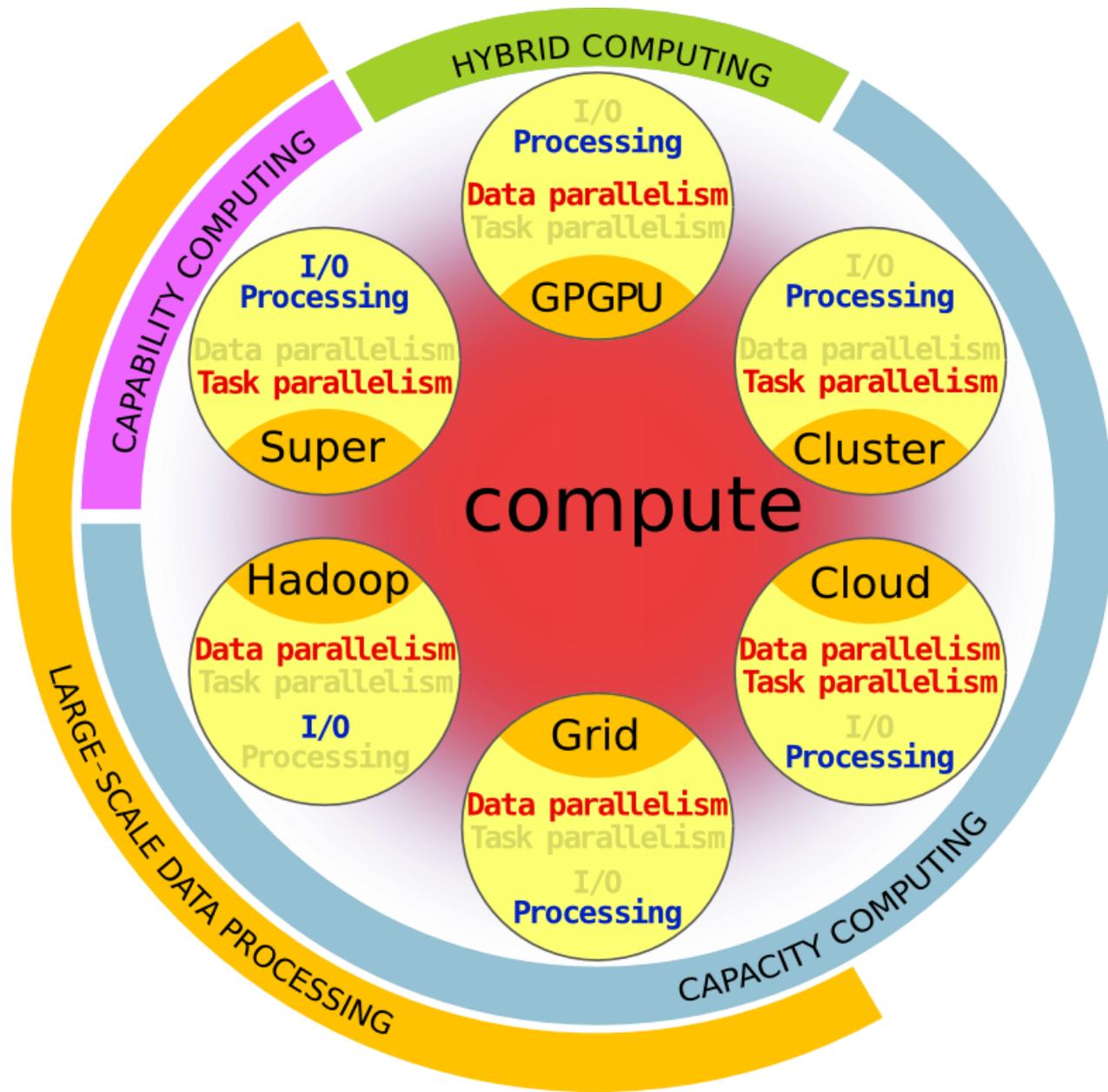


**SURF MARKET**

Favourable conditions for ICT services, software, content

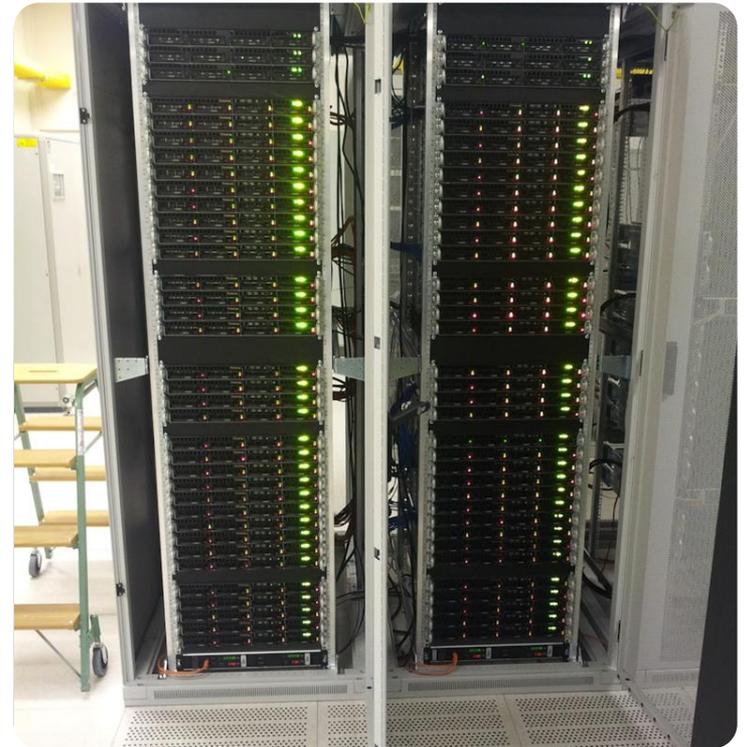
# SARA → SURFsara

- Founded in 1971 as SARA by UvA, VU and CWI as shared data-processing facility.
- Starting 1984 took the role as national supercomputing center.
- Became an independent foundation in 1995.
- Joined the SURF foundation as SURFsara in 2013.



# Hadoop cluster

- Started a test cluster in 2011 on six old machines
- Real cluster in 2012: 60 machines: Hadoop 0.20, MapReduce, Pig
- Now ~ 100 machines: Hadoop 2.6, MapReduce, Pig, Spark, Giraph, Tez, Cascading



# Hadoop 1.0

## Basic cluster components

- One of each:
  - Namenode (NN): master node for HDFS
  - Jobtracker (JT): master node for job submission
- Set of each per slave machine:
  - Tasktracker (TT): contains multiple task slots
  - Datanode (DN): serves HDFS data blocks

# Hadoop 2.0: YARN

No longer just MapReduce:

One ResourceManager (~ JobTracker)

Many NodeManagers (~ TaskTracker)

Job coordination is done by an ApplicationMaster (one per job)  
(Used to be the JobTracker)

# Using the cluster

- Command-line: from your own computer or our login node
- Resource manager web-interface
  
- Develop in your favorite IDE (Eclipse, IntelliJ)
- Package your jobs as jar files
- Submit the jar file using 'hadoop jar' or 'yarn jar'

# Dependency management

Your code probably depends on libraries.

These libraries need to be available on the cluster machines.

Multiple options:

1. Specify on command line:  
- `yarn jar myjar.jar -libjars foo.jar,bar.jar`
2. Bundle the jars inside the lib folder or your jar.
3. Extract all dependency class files (maven shade plugin)

Build tools like maven, ivy and ant can help you with this.

Example Maven POM-file (using method 2):

<http://beehub.nl/surfsara-hadoop/public/lsde-pom.xml>

You don't need to include the Hadoop/MapReduce dependencies.

# Step 1 – Login node

Access via SSH:

```
ssh lsdeXX@login.hathi.surfsara.nl (replace lsdeXX with your username)
```

Optionally enable X-Forwarding for graphical applications:

```
ssh -X lsdeXX@login.hathi.surfsara.nl
```

## Step 2 – Interacting with the HDFS

Use the 'hdfs dfs' command to access the distributed filesystem. Some common commands include:

<code>hdfs dfs -ls dir</code>	list contents of 'dir'
<code>hdfs dfs -rm file</code>	remove file
<code>hdfs dfs -cat file</code>	print file
<code>hdfs dfs -copyFromLocal src dest</code>	copy src on login node to dest on HDFS
<code>hdfs dfs -copyFromLocal src dest</code>	copy src on HDFS to dest on login node

The full list can be found at <http://hadoop.apache.org/docs/r2.6.0/hadoop-project-dist/hadoop-common/FileSystemShell.html>

# Step 3 – Submitting jobs

(MapReduce) jobs can be submitted using the 'yarn jar' command.

This runs one of the standard jobs bundled with the Hadoop framework:

```
yarn jar /usr/hdp/2.2.0.0-2041/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 10 10
```

Generally, you build your jar file on your desktop, use scp to copy it to the login node and use:

```
yarn jar JARFILE MAINCLASS ARGUMENTS
```

# Step 4 – ResourceManager web interface

You can look at the progress of your job and the log files of individual process on the web interface of the ResourceManager.

This can be accessed via 'firefox' started on the login node (X-Forwarding needed)

You will need to change one setting in Firefox, see <https://surfsara.nl/systems/hadoop/usage>

# Need help?

Problems using the SURFsara Hadoop cluster?

Contact either your course instructors or  
[hadoop.support@surfsara.nl](mailto:hadoop.support@surfsara.nl)