



**E-Infrastructures
H2020-EINFRA-2014-2015**

**EINFRA-4-2014: Pan-European High Performance Computing
Infrastructure and Services**

PRACE-4IP

PRACE Fourth Implementation Phase Project

Grant Agreement Number: EINFRA-653838

**D4.4
MOOC Pilot for HPC**

Final

Version: 1.0
Author(s): Eli Shmueli, IUCC
Date: 17.02.2017

Project and Deliverable Information Sheet

PRACE Project	Project Ref. №: EINFRA-653838	
	Project Title: PRACE Fourth Implementation Phase Project	
	Project Web Site: http://www.prace-project.eu	
	Deliverable ID: < D4.4 >	
	Deliverable Nature: <DOC_TYPE: Report >	
	Dissemination Level: PU*	Contractual Date of Delivery: 28 / February / 2017
		Actual Date of Delivery: 28 / February / 2017
EC Project Officer: Leonardo Flores Añover		

* - The dissemination level are indicated as follows: PU – Public, CO – Confidential, only for members of the consortium (including the Commission Services) CL – Classified, as referred to in Commission Decision 2991/844/EC.

Document Control Sheet

Document	Title: MOOC Pilot for HPC	
	ID: D4.4	
	Version: <1.0>	Status: <i>Final</i>
	Available at: http://www.prace-project.eu	
	Software Tool: Microsoft Word 2010	
	File(s): D4.4.docx	
Authorship	Written by:	Eli Shmueli, IUCC
	Contributors:	Jussi Enkovaara, CSC David Henty, EPCC Leon Kos, UL Janez Povh, UL
	Reviewed by:	Nico Sanna, CINECA Veronica Teodor, FZJ
	Approved by:	MB/TB

Document Status Sheet

Version	Date	Status	Comments
0.1	28/01/2017	Draft	First draft
0.2	17/02/2017	Draft	Draft for review
1.0	28/02/2017	Final version	

Document Keywords

Keywords:	PRACE, HPC, Research Infrastructure, MOOC , FutureLearn, Online learning
------------------	--

Disclaimer

This deliverable has been prepared by the responsible Work Package of the Project in accordance with the Consortium Agreement and the Grant Agreement n° EINFRA-653838. It solely reflects the opinion of the parties to such agreements on a collective basis in the context of the Project and to the extent foreseen in such agreements. Please note that even though all participants to the Project are members of PRACE AISBL, this deliverable has not been approved by the Council of PRACE AISBL and therefore does not emanate from it nor should it be considered to reflect PRACE AISBL's individual opinion.

Copyright notices

© 2017 PRACE Consortium Partners. All rights reserved. This document is a project document of the PRACE project. All contents are reserved by default and may not be disclosed to third parties without the written consent of the PRACE partners, except as mandated by the European Commission contract EINFRA-653838 for reviewing and dissemination purposes.

All trademarks and other rights on third party products mentioned in this document are acknowledged as own by the respective holders.

Table of Contents

Project and Deliverable Information Sheet	i
Document Control Sheet.....	i
Document Status Sheet	i
Document Keywords	ii
Table of Contents	iii
List of Figures.....	iii
List of Tables.....	iii
References and Applicable Documents	iv
List of Acronyms and Abbreviations.....	v
List of Project Partner Acronyms.....	vi
Executive Summary	1
1 Introduction	1
2 MOOC platform selection	2
3 Supercomputing	4
3.1 Course objectives.....	4
3.2 Course syllabus	5
3.3 Implementation on FutureLearn	5
3.4 Current status.....	6
4 Managing Big Data with R and Hadoop	6
4.1 Course objectives.....	6
4.1 Course syllabus	7
4.2 Current status	8
5 Summary	8

List of Figures

Figure 1: Age distribution of Supercomputing participants	6
---	---

List of Tables

Table 1: xMOOC vs. cMOOC (Source: [3]).....	3
---	---

References and Applicable Documents

- [1] Shah, D. "By the Numbers: MOOCs in 2016" [Online]. Available: <https://www.class-central.com/report/mooc-stats-2016/>
- [2] "FutureLearn," [Online]. Available: <https://www.futurelearn.com/>
- [3] Aparicio, F. Bacao and T. Oliveira, "MOOC's business models: turning black swans into gray swans," in *Proceedings of the International Conference on Information Systems and Design of Communication*, 2014.
- [4] "Coursera," [Online]. Available: <https://www.coursera.org/>
- [5] "edX," [Online]. Available: <https://www.edx.org/>

List of Acronyms and Abbreviations

aisbl	Association International Sans But Lucratif (legal form of the PRACE-RI)
BCO	Benchmark Code Owner
CoE	Center of Excellence
CPU	Central Processing Unit
CUDA	Compute Unified Device Architecture (NVIDIA)
DARPA	Defense Advanced Research Projects Agency
DEISA	Distributed European Infrastructure for Supercomputing Applications EU project by leading national HPC centres
DoA	Description of Action (formerly known as DoW)
EC	European Commission
EESI	European Exascale Software Initiative
EoI	Expression of Interest
ESFRI	European Strategy Forum on Research Infrastructures
GB	Giga (= $2^{30} \sim 10^9$) Bytes (= 8 bits), also GByte
Gb/s	Giga (= 10^9) bits per second, also Gbit/s
GB/s	Giga (= 10^9) Bytes (= 8 bits) per second, also GByte/s
GÉANT	Collaboration between National Research and Education Networks to build a multi-gigabit pan-European network. The current EC-funded project as of 2015 is GN4.
GFlop/s	Giga (= 10^9) Floating point operations (usually in 64-bit, i.e. DP) per second, also GF/s
GHz	Giga (= 10^9) Hertz, frequency = 10^9 periods or clock cycles per second
GPU	Graphic Processing Unit
HET	High Performance Computing in Europe Taskforce. Taskforce by representatives from European HPC community to shape the European HPC Research Infrastructure. Produced the scientific case and valuable groundwork for the PRACE project.
HMM	Hidden Markov Model
HPC	High Performance Computing; Computing at a high performance level at any given time; often used synonym with Supercomputing
HPL	High Performance LINPACK
ISC	International Supercomputing Conference; European equivalent to the US based SCxx conference. Held annually in Germany.
KB	Kilo (= $2^{10} \sim 10^3$) Bytes (= 8 bits), also KByte
LINPACK	Software library for Linear Algebra
MB	Management Board (highest decision making body of the project)
MB	Mega (= $2^{20} \sim 10^6$) Bytes (= 8 bits), also MByte
MB/s	Mega (= 10^6) Bytes (= 8 bits) per second, also MByte/s
MFlop/s	Mega (= 10^6) Floating point operations (usually in 64-bit, i.e. DP) per second, also MF/s
MooC	Massively open online Course
MoU	Memorandum of Understanding.
MPI	Message Passing Interface
NDA	Non-Disclosure Agreement. Typically signed between vendors and customers working together on products prior to their general availability or announcement.
PA	Preparatory Access (to PRACE resources)

PATC	PRACE Advanced Training Centres
PRACE	Partnership for Advanced Computing in Europe; Project Acronym
PRACE 2	The upcoming next phase of the PRACE Research Infrastructure following the initial five year period.
PRIDE	Project Information and Dissemination Event
RI	Research Infrastructure
TB	Technical Board (group of Work Package leaders)
TB	Tera (= 240 ~ 1012) Bytes (= 8 bits), also TByte
TCO	Total Cost of Ownership. Includes recurring costs (e.g. personnel, power, cooling, maintenance) in addition to the purchase cost.
TDP	Thermal Design Power
TFlop/s	Tera (= 1012) Floating-point operations (usually in 64-bit, i.e. DP) per second, also TF/s
Tier-0	Denotes the apex of a conceptual pyramid of HPC systems. In this context the Supercomputing Research Infrastructure would host the Tier-0 systems; national or topical HPC centres would constitute Tier-1
UNICORE	Uniform Interface to Computing Resources. Grid software for seamless access to distributed resources.

List of Project Partner Acronyms

BADW-LRZ	Leibniz-Rechenzentrum der Bayerischen Akademie der Wissenschaften, Germany (3 rd Party to GCS)
BILKENT	Bilkent University, Turkey (3 rd Party to UYBHM)
BSC	Barcelona Supercomputing Center - Centro Nacional de Supercomputacion, Spain
CaSToRC	Computation-based Science and Technology Research Center, Cyprus
CCSAS	Computing Centre of the Slovak Academy of Sciences, Slovakia
CEA	Commissariat à l'Energie Atomique et aux Energies Alternatives, France (3 rd Party to GENCI)
CESGA	Fundacion Publica Gallega Centro Tecnológico de Supercomputación de Galicia, Spain, (3 rd Party to BSC)
CINECA	CINECA Consorzio Interuniversitario, Italy
CINES	Centre Informatique National de l'Enseignement Supérieur, France (3 rd Party to GENCI)
CNRS	Centre National de la Recherche Scientifique, France (3 rd Party to GENCI)
CSC	CSC Scientific Computing Ltd., Finland
CSIC	Spanish Council for Scientific Research (3 rd Party to BSC)
CYFRONET	Academic Computing Centre CYFRONET AGH, Poland (3 rd party to PNSC)
EPCC	EPCC at The University of Edinburgh, UK
ETHZurich (CSCS)	Eidgenössische Technische Hochschule Zürich – CSCS, Switzerland
FIS	FACULTY OF INFORMATION STUDIES, Slovenia (3 rd Party to ULFME)
GCS	Gauss Centre for Supercomputing e.V.
GENCI	Grand Equipement National de Calcul Intensif, France
GRNET	Greek Research and Technology Network, Greece

INRIA	Institut National de Recherche en Informatique et Automatique, France (3 rd Party to GENCI)
IST	Instituto Superior Técnico, Portugal (3 rd Party to UC-LCA)
IT4I	IT4Innovations National Supercomputing Center, Czech Republic
IUCC	INTER UNIVERSITY COMPUTATION CENTRE, Israel
JKU	Institut fuer Graphische und Parallele Datenverarbeitung der Johannes Kepler Universitaet Linz, Austria
JUELICH	Forschungszentrum Juelich GmbH, Germany
KIFU	Governmental Information Technology Development Agency, Hungary
KTH	Royal Institute of Technology, Sweden (3 rd Party to SNIC)
LiU	Linkoping University, Sweden (3 rd Party to SNIC)
NCSA	NATIONAL CENTRE FOR SUPERCOMPUTING APPLICATIONS, Bulgaria
NTNU	The Norwegian University of Science and Technology, Norway (3 rd Party to SIGMA)
NUI-Galway	National University of Ireland Galway, Ireland
PRACE	Partnership for Advanced Computing in Europe aisbl, Belgium
PSNC	Poznan Supercomputing and Networking Center, Poland
RISCSW	RISC Software GmbH
RZG	Max Planck Gesellschaft zur Förderung der Wissenschaften e.V., Germany (3 rd Party to GCS)
SIGMA2	UNINETT Sigma2 AS, Norway
SNIC	Swedish National Infrastructure for Computing (within the Swedish Science Council), Sweden
STFC	Science and Technology Facilities Council, UK (3 rd Party to EPSRC)
SURFsara	Dutch national high-performance computing and e-Science support center, part of the SURF cooperative, Netherlands
UC-LCA	Universidade de Coimbra, Laboratório de Computação Avançada, Portugal
UCPH	Københavns Universitet, Denmark
UHEM	Istanbul Technical University, Ayazaga Campus, Turkey
UiO	University of Oslo, Norway (3 rd Party to SIGMA)
ULFME	UNIVERZA V LJUBLJANI, Slovenia
UmU	Umea University, Sweden (3 rd Party to SNIC)
UnivEvora	Universidade de Évora, Portugal (3 rd Party to UC-LCA)
UPC	Universitat Politècnica de Catalunya, Spain (3 rd Party to BSC)
UPM/CeSViMa	Madrid Supercomputing and Visualization Center, Spain (3 rd Party to BSC)
USTUTT-HLRS	Universitaet Stuttgart – HLRS, Germany (3 rd Party to GCS)
WCNS	Politechnika Wroclawska, Poland (3 rd party to PNSC)

Executive Summary

PRACE, at its current stage, provides a rich variety of online learning resources via PRACE Training portal, including lecture slides, videos and exercise material. Although we are yet to reach the ultimate goal of a complete online training experience in the form of online courses, we have made progress and the training work package (WP4) of the PRACE fourth Implementation Phase (PRACE-4IP) is in the process of creating two Massive Open Online Courses (MOOCs). These MOOCs will be open to the general public, free of charge, and will introduce the audience to the subjects “Supercomputing” and “Managing Big Data with R and Hadoop”. After due consideration of pedagogy, design and other aspects, the FutureLearn platform has been selected to deliver the PRACE MOOCs. Such open courses, delivered and advertised by a well-known platform, increase not only the public’s interest in high performance computing but also the visibility of the PRACE brand. We expect the courses to be launched in March 2017.

1 Introduction

PRACE has been providing online learning resources already for a long time via the PRACE Training portal. Lecture slides and exercise material for PRACE Advanced Training Centre courses and seasonal schools are normally made available via the Training portal, and there are also video tutorials on selected topics. However, so far no courses designed specifically for online learning have been developed within PRACE. The purpose of this work is to extend the PRACE training by offering courses where the syllabus and all the content are directed to online learning.

Online learning requires new approaches when compared to traditional face-to-face training in class room. Simply video recording lectures and publishing lecture slides do not normally offer very good online training experience, but syllabus and material have to be often designed and developed from beginning. Typically, online courses consist of videos (short, 5-10 minute, are preferred), textual material, quizzes and interaction for example via discussion forums. There can be also exercises which can be submitted and checked automatically for correctness. There is lots of ongoing research related to the cognitive scientific aspects of online training.

Massively Open Online Courses (MOOCs) are online courses which aim for unlimited and open participation via web. The main characteristics of MOOCs are their intended scale, at best tens of thousands of students can participate in a MOOC. The openness means that there are no restrictions in attending a course, and the courses are normally also free of charge. MOOCs resemble traditional courses in the sense they are often run during a specified time, and guidance is provided during that time in some online form. Duration of online courses varies typically between two and eight weeks, and the course is often divided into weekly modules.

Many universities and non-profit educational organisations, as well as commercial companies are currently providing MOOCs in several different platforms. The business model in most commercial MOOCs is that the courses themselves are free, but there is a fee for obtaining a certificate of attendance or certificate for passing the course. In 2016, 6850 courses have been offered by over 700 institutes. [1]

This deliverable describes the work that has been done in PRACE 4IP for developing MOOCs. The main objective for PRACE is to enlarge the number of people that can benefit

from expertise of PRACE trainers and thus enhance European research. High-quality online training strengthens also the PRACE brand.

This document is structured as follows: section 2 describes the FutureLearn [2] MOOC platform and the process for selecting it. Section 3 describes the “Supercomputing” MOOC and Section 4 the “Management of massive data” MOOC. Finally, a summary is presented in Section 5.

2 MOOC platform selection

In this section, we discuss the factors which have led us to select the FutureLearn MOOC platform as the appropriate platform for the project. The importance of the platform cannot be underestimated, and the selection demanded thorough characterisation and examination of the other options, all to make certain that the educational material conversion and transfer to an online platform will be done effectively. We would like to note that a delay due to contract negotiation has postponed this deliverable from M21 to M25.

During the process, we have examined the parameters of each system, taking a close look at its main features and unique characteristics, including its technological infrastructure and annual costs.

Several platforms have been considered, among them world-leading platforms and niche platforms. In the process, we have examined:

1. The efficiency of the learning process. Materials must be presented in a clear, understandable manner.
2. The platform’s ability to allow learning interactions between students.
3. The implementation of automatic tools which monitor the learning process and the learners and analyse their performances.
4. The probability of strengthening the PRACE brand by delivering high quality educational material. The platform acts as a showcase for PRACE, enabling exposure to new audiences which might have not been aware of its existence.

Two MOOC models

There are two major models of MOOC delivery: cMOOC (connectivist MOOC) and xMOOC (extended MOOC), with xMOOC being the most popular MOOC model. xMOOCs are delivered on special platforms and are designed to handle a large number of participants. The tutoring model is that of one-to-many. They use video lectures, computer-graded quizzes and exams, peer assessment and often award recognition (e.g. certificate) upon successful completion of the course.

cMOOC emphasizes networking and participants’ contributions. Many times, cMOOCs are not supported by institutes, and use general social media platforms (e.g. Twitter) for community interaction. The content is mostly participants-driven. The differences between the xMOOC and the cMOOC are summarized in Table 1.

Since the planned courses are beginner courses which will be delivered to a large audience, the model chosen for the PRACE MOOCs was the xMOOC.

Areas	Dimensions	cMOOCs	xMOOCs
Product Innovation	Value Proposition	Knowledge creation, autonomy, social network, social recognition, informal learning	Knowledge acquisition, certification, tutoring, collaborative groups, access certified experts
	Customer	Students, practitioners, peers	Students, practitioners, enterprises
Customer/ user Relationship	Distribution Chanel	Internet, Web, social learning platforms, MOOCs platforms	Internet, Web, MOOCs platforms
	Relationship	Communities of Practice	Global learning environment
	Partnership	Universities, schools, enterprises	Universities, schools, enterprises, other MOOCs platforms
Infrastructure Management	Value Configuration	Promotion, sharing, collaboration	Platform, universities, organizations
	Capability	Promotion of motivational processes to continuous usage	Promoting process of usage, cloud hosted MOOC services
Financial Aspects	Revenue	Sponsoring, platform analytics	Sponsoring, platform analytics, certification, potential on-campus students, tailored training courses for enterprises, SaaS services, tuition fees
	Cost	Platform infrastructure, maintenance	Platform infrastructure, maintenance, tutoring, courses development

Table 1: xMOOC vs. cMOOC (Source: [3])

xMOOC platforms for PRACE

We have examined the existing options and platforms available for PRACE, including the presentation of material, management of learning processes, study interaction, design and costs. We decided to focus on three central, cutting-edge platforms which offer international infrastructure for delivering the PRACE MOOCs.

FutureLearn [2]

FutureLearn has been founded by the Open University of UK at the end of 2012. It is the fourth biggest MOOC platform, with 5.2 million registered users, 2.3 million of them joined in 2016. It offers about 500 courses provided by over 100 partners. FutureLearn has, in addition to universities, other prominent partners such as the British Museum and the European Space Agency. Students can learn by watching videos, listening to audio resources and reading articles. Each article, video or piece of audio has a dedicated space to allow learners to comment and ask questions.

Coursera [4]

Founded by Stanford professors in 2012, Coursera is currently the biggest MOOC provider, with 23 million registered users and over 1700 active courses. It has added six million new users in 2016. Coursera has partners such as Yale University and Stanford University. Courses are interactive (during a lecture students get questions to answer on site), there are deadlines, quizzes, and sometimes a final exam (depends on the course).

edX [5]

edX was founded by MIT and Harvard University in 2012 and is the second largest MOOC provider, after Coursera. EdX offers about 1300 courses. Four million new users joined edX in 2016 and it currently has 10 million users. Among the edX partners are the University of Oxford and the University of California, Berkeley. The structure of courses is similar to Coursera, with deadlines, quizzes, and exams.

Why FutureLearn?

FutureLearn is the largest European MOOC platform. It is an English-speaking platform, which allows it to address a large audience worldwide. FutureLearn emphasizes social

constructivist learning in which the learner community plays a key role. FutureLearn courses are shorter, and use less video and more written material than Coursera and edX. It not only provides more flexibility to many learners, but also decreases production costs per course, which enables us to provide more in the long run. In comparison with Coursera and edX, FutureLearn has a higher proportion of discussions in its courses, and they are considered part of the learning activities. While in Coursera and edX there is a designated part of the course for discussions, in FutureLearn comments appear next to the materials, which allows for a more active discussion. The platform follows social media principles, so that learners can create profiles, reply directly to other learners, “Like” comments and filter information. The use of such social elements and features gives room to live conversations as part of the classes, in a way that not only enriches the learning experience, but also contributes to building bridges between learners around the globe. Another aspect is that FutureLearn is the key European solution to xMOOC, and this initiative may strengthen its role in the MOOC world with regards to the United States counterparts.

3 Supercomputing

3.1 Course objectives

This course introduces what supercomputers are, how they are used and how one can exploit their full computational potential to make scientific breakthroughs. This course is designed for anyone interested in leading-edge computing technology, supercomputers or the role that computer simulation takes in modern science and engineering.

All of the technical aspects will be covered at a conceptual level and there is no requirement to be able to write computer programs. However, anyone with existing programming experience will learn how programming modern supercomputers differs from programming a home PC.

The learning outcomes are that, by the end of the course, students will be able to:

- Understand how the performance of modern supercomputers is measured and achieved;
- Explain why they are built from thousands of simple processors;
- Understand the differences between of shared-memory and distributed-memory computers;
- Compare the architecture of a typical modern supercomputer with a desktop PC;
- Explain why computer simulation is a fundamental component of modern scientific discovery;
- Work with simple cellular automaton models;
- Analyse simple problems and look for opportunities for parallel processing;
- Explain the limitations of parallel computing;
- Give examples of scientific areas where computer simulation is used.

3.2 Course syllabus

The duration of the course is five weeks and the estimated work load per week is three hours. The high-level breakdown over weeks is:

Week 1: Supercomputers

This will cover what modern supercomputers look like and why they are designed that way. There will be a short history of supercomputers over the years to illustrate how performance has increased over time

Week 2: Parallel Computers

The concepts of shared and distributed memory architectures will be explained. We will cover the similarities and differences between specialist supercomputers and more general purpose computers such as PCs, laptops and games machines.

Week 3: Parallel Computing

The basic parallel programming models (shared and distributed memory) will be explained at a conceptual level using a traffic modelling example. This will be used to motivate the basics of parallel performance, including Amdahl's and Gustafson's laws, but without using any equations.

Week 4: Computer Simulation

This will cover what supercomputers are used for, e.g. the kinds of simulations that are done in the areas of nanotechnology, engineering and climate research. The traffic model serves as a very basic example of computer simulation, illustrating how computers effectively run virtual experiments where scientists choose the input parameters to model real situations.

Week 5: Case Studies

The final week will comprise a number of case studies of real scientific problems being tackled using supercomputers. Ideally, these will have associated material that makes the talks interesting, e.g. visualisations of real simulations. We will explain the approaches taken in terms of the concepts introduced in the previous weeks' lectures.

3.3 Implementation on FutureLearn

The choice of FutureLearn as a platform did not alter the learning outcomes or overall week-by-week breakdown of the MOOC. However, we had initially envisaged that most material would be delivered as recordings of short (10-15 minute) lectures, re-using a lot of existing lecture material in the form of PowerPoint slides from introductory training courses. The FutureLearn style is very different, favouring a series of very short steps each comprising a short article, quiz or discussion topic. Although use of video is encouraged, it is normally used to highlight specific topics and not as the major method of content delivery. Any videos are usually only a few minutes long and do not resemble traditional lectures.

Although this format fits very well with the introductory nature of the Supercomputing MOOC, it did mean that more new material was required. Although the overall volume of material is perhaps slightly less than initially planned, it has been more work to develop. The two sites involved in developing the MOOC have been EPCC and SURFsara.

3.4 Current status

The course was announced and opened for registration in early December 2016. Currently (as of 24 January 2017) there are around 1,350 people registered for the upcoming March run. FutureLearn provide various statistics on those registered, for example their age profile (Figure 1):

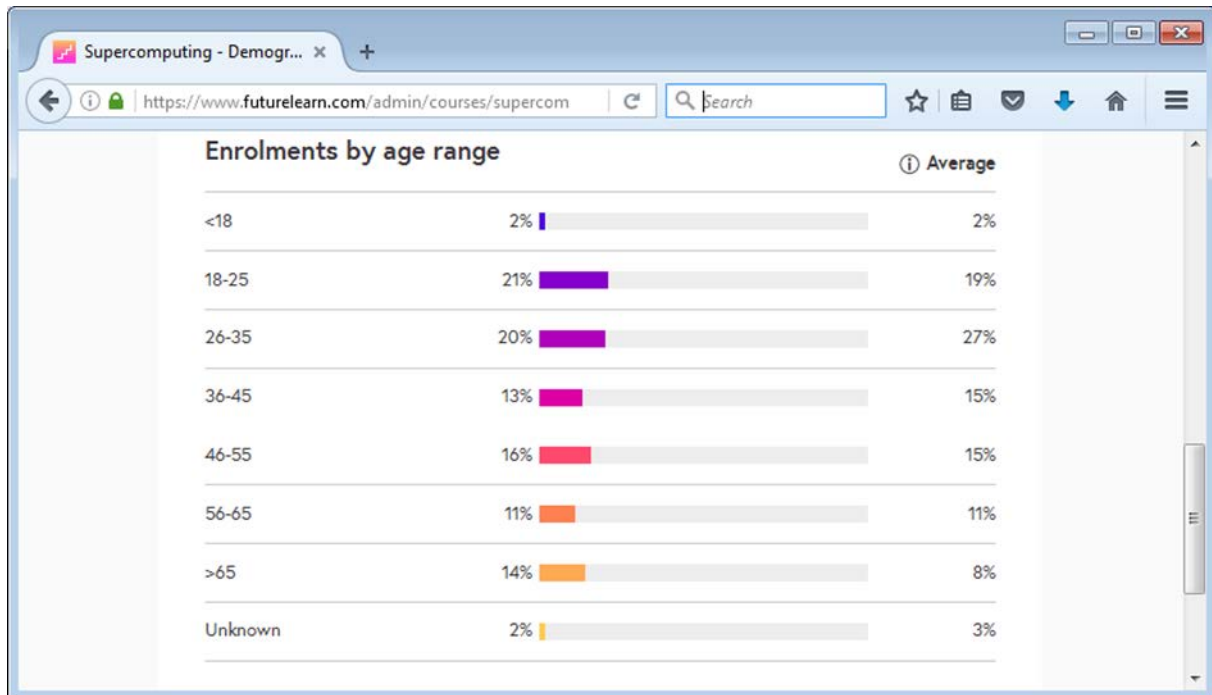


Figure 1: Age distribution of Supercomputing participants

The course is live at <https://www.futurelearn.com/courses/supercomputing>. The next step is a training session with FutureLearn where they will provide advice on supporting learners once the course has started.

4 Managing Big Data with R and Hadoop

4.1 Course objectives

The objective of this course is to introduce individuals with limited programming knowledge to various HPC facilities for big data analysis. At the end of the course, they will be able to use them, avoiding common pitfalls and thus saving them money and time. The course is especially suitable for participants interested in data science, computational statistics and machine learning. This course may be also useful for advanced undergraduate students and 1st year PhD students in data analysis, statistics, or bioinformatics who wish to gain a basic understanding in data management and HPC computing.

The main advantage and comparative variety to other similar MOOCs is that we use R environment for statistical computing and graphics. R is a programming language renowned for its simplicity, elegance, and the support of an outstanding community. R augmented with Hadoop allows data scientists to quickly utilize the enterprise-grade capabilities of Hadoop

with the analytic capabilities of R. The learning outcomes are that, by the end of the course, students will be able to:

- Understand how the performance of modern supercomputing is achieved;
- Understand the basic functionality of Bash terminal window;
- Understand the basic functionality of Apache Hadoop for scalable, distributed computing;
- Understand the basic functionality of RHadoop;
- Understand the basic problems of supervised and unsupervised learning;
- Perform basic clustering, regression and classification with RHadoop.

4.1 Course syllabus

The course will run for five weeks with estimation of four hours of working load per week. Weekly breakdown is:

Week 1: Welcome to BIG DATA MOOC

This will cover introduction to big data and distributed systems. Parallel computing systems and parallel databases will be explained in great details. Students will be introduced to **terminal window**, **sed**, **AWK**, and **grep** programs.

Week 2: First steps in R and RStudio

In this lecture we will cover an introduction to R and Rstudio. R is leading programming language for computational statistics and graphics and could be augmented for big data processing. Rstudio is a powerful and productive graphical user interface for R, make R easier and more efficient.

Week 3: Working with Apache Hadoop I (Fundamentals)

This will cover massively parallel computing with Hadoop. We will present Hadoop as a system for distributed computing and data storage, show how to use “ravro” library for reading and writing files in avro format and present powerful “plyrmr” library for data processing. We will also explain how to set up local virtual environment and prepare an image of virtual machine that the students will download and locally install.

Week 4: Working with Apache Hadoop II (RHadoop)

This lecture will cover introduction to functions providing management of the Hadoop Distributed File System (HDFS) using the “rhdfs” library. We will teach how R programmers can browse, read, write, and modify files stored in HDFS. We will also present how to perform statistical analysis via Hadoop MapReduce functionality on a Hadoop cluster. Finally, we will introduce HBASE distributed database and present how to browse, read, write, and modify tables stored in it.

Week 5: Statistical learning

This will cover introduction to machine learning and to classification framework in particular. We will present principles of supervised and unsupervised learning, the most frequent algorithms in linear regression, classification (discriminant analysis) and clustering (hierarchical clustering and k-means). For all these algorithms we will demonstrate how to implement them using R and RHadoop to handle big data.

Materials:

Students enrolled to this MOOC will have available:

- Virtual machine with Cloudera Hadoop installation file to make a local copy of such machine;
- Video material with short (up to 10 minutes) courses;
- Additional slides and other materials for deeper reading;
- Several assignments, including quizzes for self – assessment of the understanding.

Teachers

A group of three teachers will prepare the materials and run the MOOC:

- Prof. Leon Kos, PhD, University of Ljubljana, Faculty of mechanical engineering, Slovenia;
- Prof. Janez Povh, PhD, University of Ljubljana, Faculty of mechanical engineering, Slovenia;
- Prof. Biljana Mileva Boshkoska, Faculty of informations studies in Novo mesto, Slovenia.

4.2 Current status

The course is planned to run first time starting from 20 March 2017. About 560 learners have registered so far. Additional runs are planned at the moment for September-October 2017 and March – April 2018.

5 Summary

In the PRACE-4IP MOOC project, we have worked to define the MOOC goals and their alignment with PRACE goals. We have chosen the most appropriate platform out of three world-leading and niche-oriented platforms and defined the required characteristics for the PRACE courses. We have mapped the process of MOOC development, and decided on the developing institutes. Currently we are nearing the end of the development process, and intend to finalize the courses around March. The development process in the institutes is supported by the FutureLearn professionals. Once the courses are complete, we and the FutureLearn team will be able to advertise the courses among the general audience of potential MOOC students and the PRACE community.

The PRACE MOOC activity will become part of the overall online training offered by PRACE, with a wide range of online learning materials for the HPC community in general and the PRACE community in particular. This variety of online activities will allow us to offer the community a complete suite of online training.

Future activity as part of PRACE-5IP will allow development of additional courses and support of the existing courses, ensuring their sustainability. This will create a complete online training portal, which will include a variety of quality study materials from PRACE projects, establishing PRACE as a leading professional authority in the field of HPC.