

## PROGRAMA DE CURSO

Código	Nombre							
CC5212	Procesamiento Masivo de Datos							
<b>Nombre en Inglés</b>								
<b>Massive Data Processing</b>								
SCT	Unidades Docentes	Horas de Cátedra	Horas Docencia Auxiliar	Horas de Trabajo Personal				
6	10	3	1,5	5,5				
<b>Requisitos</b>			<b>Carácter del Curso</b>					
CC3201 Bases de Datos			Electivo					
<b>Resultados de Aprendizaje</b>								
<p>Al finalizar el curso el alumno entenderá los fundamentos del procesamiento masivo de datos, enfocado a la forma en que esto es realizado por grandes empresas como Google, Facebook, Twitter, Amazon, o grandes laboratorios de investigación como CERN.</p> <p>The goal of the course will be to provide a solid foundation in the traditional design aspects relating to Distributed Computing and Distributed Databases, showing how they have influenced modern developments in cloud computing, including distributed data storage (e.g., NoSQL storage techniques) and data processing abstractions (e.g., MapReduce/Hadoop, Pregel/Giraph).</p> <p>Practical assignments will reinforce course material and provide introductory hands-on programming experience developing distributed applications for handling massive data.</p>								

Metodología Docente	Evaluación General
Clases expositivas de 90 minutos cada una y sesiones prácticas de 90 minutos.	La evaluación contemplará Controles, Examen y Actividades Complementarias.

### Unidades Temáticas

Número	Nombre de la Unidad	Duración en Semanas	
1	Distribución y Paralelismo	3	
	Contenidos	Resultados de Aprendizajes de la Unidad	Referencias a la Bibliografía
<ul style="list-style-type: none"> <li>Course Introduction/Motivation for handling massive data.</li> <li>Introduction to Distributed Systems.</li> <li>Goals of a Distributed System</li> <li>Fundamentals of parallel processing and distributed computing.</li> <li>Traditional distributed computing architectures.</li> <li>Modern distributed computing instantiations.</li> </ul>	<p>Understand the fundamental design principles of Distributed Systems (fault tolerance, performance, horizontal scale, economy). See why parallel processing and distributed storage are key to handling massive data.</p> <p>Learn about the different types of Distributed Systems: learn about traditional abstractions (P2P, client-server, grid, OSI layers), and see how they influence modern instantiations (Bittorrent, Internet, Web, cloud computing, Google, etc.).</p>	<p>[TS06, cap. 1,2]  [KDF11, cap. 1,2]</p>	

Número	Nombre de la Unidad	Duración en Semanas	
2	Modelamiento de Procesamiento Distribuido	4	
	Contenidos	Resultados de Aprendizajes de la Unidad	Referencias a la Bibliografía
<ul style="list-style-type: none"> <li>Traditional distributed communication methods (RPC/MPI)</li> <li>Modern distributed processing abstractions, focusing on MapReduce with introduction to others (including Graph)</li> <li>Course assignments using Java RMI for introductory distributed communication methods</li> <li>Course assignments using Hadoop (MapReduce) to parallelize and solve a basic distributed computing task</li> </ul>	<p>Learn basics of distributed communication (RPC/MPI): gain hands-on experience of communicating between machines through introductory Java RMI assignment (intermediate Java skills a prerequisite).</p> <p>Learn modern distributed (cloud) computation abstractions, including MapReduce and Pregel (as used by Google et al.): gain hands-on experience with introductory programming assignments using Hadoop and Giraph (intermediate Java skills a prerequisite).</p>	<p>[TS06, cap. 2]  [W12]  [MABDHLC10]  [OV11, cap. 18]</p>	

Número	Nombre de la Unidad	Duración en Semanas
3	Conceptos de Manejo de Datos Distribuidos	4
Contenidos	Resultados de Aprendizajes de la Unidad	Referencias a la Bibliografía
<ul style="list-style-type: none"> <li>• Introduction to Distributed Databases</li> <li>• Goals of distributed databases</li> <li>• Data-placement strategies: sharding (partitioning), replication, duplication</li> <li>• Fundamentals of distributed query processing and optimization</li> <li>• Programming assignment to build a basic distributed database using Java RMI and on-disk files</li> </ul>	<p>Learn the fundamentals of Distributed Databases, including the trade-offs between fault-tolerance, scalability, performance and economy.</p> <p>Learn the different ways in which data can be spread out and managed over multiple machines, so as to enable efficient and reliable querying over massive data.</p> <p>Understand the different types of guarantees a distributed database can make, and their formal limitations.</p> <p>Build a basic distributed database with hash-based sharding and replication: assignment will use Java RMI (intermediate Java skills a prerequisite).</p>	[TS06, cap. 6,7] [OV11]

Número	Nombre de la Unidad	Duración en Semanas
4	Modelos de Almacenamiento Escalable	4
Contenidos	Resultados de Aprendizajes de la Unidad	Referencias a la Bibliografía
<ul style="list-style-type: none"> <li>• Introduction to modern NoSQL stores</li> <li>• Taxonomy of different NoSQL data models: key-value stores, column stores, document stores, graph stores, relational/tabular stores, triple stores, inverted indexes.</li> <li>• Brief introduction to modern NoSQL query and imperative languages.</li> </ul>	<p>Cover the taxonomy of current NoSQL stores commonly used for large-scale data management in cluster/cloud computing environments.</p> <p>Compare and contrast the strengths and weaknesses of different data models employed by stores.</p> <p>Learn about the different query languages employed by different stores.</p>	[SF12]

<ul style="list-style-type: none"> <li>• Programming assignment combining MapReduce (Hadoop) and a NoSQL store (TBD) for a distributed processing task.</li> </ul>	<p>Hands-on assignment will combine a NoSQL store and MapReduce to solve a distributed processing task, bringing together aspects of the different modules (intermediate Java skills a prerequisite).</p>	
--	---	--

Bibliografía	
<ul style="list-style-type: none"> <li>• [TS06] A. S. Tanenbaum, M. Van Steen. <i>Distributed Systems: Principles and Paradigms</i> (2nd Edition). Prentice Hall, 2006.</li> <li>• [MABDHLC10] G. Malewicz, M. H. Austern, A. J. C. Bik, J. C. Dehnert, I. Horn, N. Leiser, G. Czajkowski. Pregel: a system for large-scale graph processing. SIGMOD Conference 2010: 135-146.</li> <li>• [KDF11] K. Hwang, J. Dongarra, G. C. Fox. <i>Distributed and Cloud Computing: From Parallel Processing to the Internet of Things</i> (1st Edition). Morgan Kaufmann, 2011.</li> <li>• [OV11] M. T. Özsü, P. Valduriez. <i>Principles of Distributed Database Systems</i>. Springer, 2011.</li> <li>• [W12] T. White. <i>Hadoop: The Definitive Guide</i>. O'Reilly, 2012.</li> <li>• [SF12] P. J. Sadalage, M. Fowler. <i>NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence</i>. Addison-Wesley Professional, 2012.</li> </ul>	

Vigencia desde:	Otoño 2014
Elaborado por:	Pablo Barcelo / Aidan Hogan