

# USC Viterbi School of Engineering

## **CSCI 648: Advanced Information Integration**

**Units: 4**

**Term—Day—Time:**

**Spring 2015 – MW – 5:30-6:50pm**

**Location:** THH 114

**Instructor: Jose Luis Ambite**

**Office:** Outside classroom

**Office Hours:** Immediately after class

**Contact Info:** [ambite@isi.edu](mailto:ambite@isi.edu), 310-448-8472.

**Instructor: Craig Knoblock**

**Office:** AFH B55a

**Office Hours:** Wednesdays 10-11am

**Contact Info:** [knoblock@usc.edu](mailto:knoblock@usc.edu), 310-448-8786.

**Teaching Assistant: Bo Wu**

**Office:** plaza between RTH cafe and EEB

**Office Hours: Tuesday 10am-12pm**

**Contact Info:** [wubo@usc.edu](mailto:wubo@usc.edu)

### **Catalogue Course Description**

Foundations, techniques, and algorithms for information integration. Topics include Semantic Web, linked data, data integration, entity linkage, source modeling, and information extraction.

### **Expanded Course Description**

This course focuses on foundations, techniques, and algorithms for information extraction, modeling and integration. Topics covered include semantic web (RDF, OWL, SPARQL), linked data and services, mash-ups, theory of data integration, schema mappings, record/entity linkage, data cleaning, source modeling, and information extraction. The class will be run as a lecture course with lots of student participation and significant hands-on experience. As an integral part of the course each student will do a project using the research and tools covered in the class.

### **Learning Objectives**

The learning objectives for this course are:

- Understand the foundations and techniques of the Semantic Web, including RDF, OWL, SPARKL, linked data, and linked services
- Understand the theory and techniques of traditional data integration, including view integration, schema mapping, record linkage
- Understand the algorithms and techniques for data cleaning, source modeling, building mashups, semi-structured extraction, and information extraction
- Understand the theory and application of the state-of-the-art software and tools for information extraction
- For any given integration problem, be able to select and apply the most relevant information integration techniques to solve that problem

**Prerequisite(s):** CSCI 561

**Co-Requisite (s):** none

**Concurrent Enrollment:** none

**Recommended Preparation:** CSCI 585 and some programming experience

## **Course Notes**

The course will be run as a lecture class with student participation strongly encouraged. There are weekly readings and students are encouraged to do the readings prior to the discussion in class. All of the course materials, including the readings, lecture slides, homeworks will be posted online. The class project is a significant aspect of this course and at the end of the semester, students will present their projects in class.

## **Technological Proficiency and Hardware/Software Required**

Students are expected to know how to program in a language such as Java, C++, or Python. Students are also expected to have their own laptop or desktop computer where they can install and run software to do the weekly homework assignments.

## **Required Readings and Supplementary Materials**

Required Textbook: Principles of Data Integration by Doan, Halevy, & Ives, Morgan Kaufmann, 2012

The book is available online from the USC library and is available for purchase.

All of the required readings are listed in the course schedule.

## **Description and Assessment of Assignments**

### **Homework Assignments**

There will be weekly homework assignments for the first 12 weeks of class. The assignments must be done individually. The homework assignments are expected to take 6-8 hours per week. Each assignment is graded on a scale of 0-100 and the specific rubric for each assignment is given in the assignment. The homework topics are listed in the Course Schedule.

### **Course Project**

An integral part of this course is the course project, which builds on the topics and techniques covered in the class. Students can work in teams of up to two people on this project. They will write a written proposal for the project, conduct the project, and then write a paper about the project, create a video demonstration of the work, and present the project in class.

#### *Project Timeline:*

- Week 7: Project proposals due (team members, topics, presentation date)
- Week 8: Students receive feedback on proposed projects.
- Week 12: Mid-term report due (progress to date)
- Week 15: Project presentation in class (short talk and video demonstration)
- Week 15: Project papers due at end of week

*Sample project: "Geotagging Ansel Adams' Photographs"* Ansel Adams was one of the most famous American photographers. However, there is no single coherent source, which provides a structured collection of all of Ansel's photographs. Also most of his photos are not geotagged and thus there is no easy way to visualize Ansel's journey across the globe. This project extracted Ansel's photos along with their metadata from various sources into a single coherent consolidated schema and geotagged them by extracting and identifying location entities from each photo's metadata. The end result is a web application, which used the Google Maps API and Timeline JS to visualize his geotagged and time-stamped photos in a very impressive manner.

*Grading breakdown of the course project:*

- Proposal: 0%
- Mid-term report: 0%
- Project paper: 25%
- Project video: 25%
- Presentation: 25%
- Overall project innovation: 25%

## Grading Breakdown

**Quizzes:** There will be weekly quizzes based on the material from the week before. There is no mid-term for this class.

**Homework:** There will be weekly homework based on the topics of the class each week.

**Final Exam:** There is a final exam at the end of the semester covering all of the material covered in the class.

**Class Project:** Each student will do an independent class project based on the topics covered in the class. Students will propose their own project, do the research and build a proof-of-concept, write a paper about the work, present the work in class, and create a video demonstration of the work.

Grading Schema:

Quizzes	20%
Homework	20%
Final:	20%
Class Project	40%
<hr/>	
Total	100%

Grades will range from A through F. The following is the breakdown for grading:

94 - 100 = A	74 - 76 = C
90 - 93 = A-	70 - 73 = C-
87 - 89 = B+	67 - 69 = D+
84 - 86 = B	64 - 66 = D
83 - 83 = B-	60 - 63 = D-
77 - 79 = C+	Below 60 is an F

## Assignment Submission Policy

Homework assignments are due at 11:59pm on the due date and should be submitted in Blackboard. You can submit homework up to one week late, but you will lose 20% of the possible points for the assignment. After one week, the assignment cannot be submitted.

## Course Schedule: A Weekly Breakdown

	Topics/Daily Activities	Readings	Quizzes & Homeworks	Instructor
<b>Week 1</b> Jan 12	<b>Course Introduction</b>	Frank Manola and Eric Miller. Rdf primer. Technical report, W3C, February 2004. <a href="http://www.w3.org/TR/2004/REC-rdf-primer-20040210/">http://www.w3.org/TR/2004/REC-rdf-primer-20040210/</a> .	Homework 0: Academic Integrity (due on Friday)	Professor Ambite
Jan 14	<b>RDF, Graph data model</b>	Tim Berners-Lee. Why rdf model is different from the xml model. Technical report, W3C, 1998.		Professor Ambite

		<a href="http://www.w3.org/DesignIssues/RDF-XML.html">http://www.w3.org/DesignIssues/RDF-XML.html</a> .		
<b>Week 2</b> Jan 21	<b>RDF Schema</b>	<p>Rdf vocabulary description language 1.0: Rdf schema. Technical report, W3C, February 2004.  <a href="http://www.w3.org/TR/2004/REC-rdf-schema-20040210/">http://www.w3.org/TR/2004/REC-rdf-schema-20040210/</a>.</p> <p>Ben Adida, Ivan Herman, Manu Sporny, and Mark Birbeck. Rdfa 1.1 primer rich structured data markup for web documents. Technical report, W3C, June 2012.  <a href="http://www.w3.org/TR/2012/NOTE-rdfa-primer-20120607/">http://www.w3.org/TR/2012/NOTE-rdfa-primer-20120607/</a>.</p>	Quiz 1 (on Wednesday) Homework 1: Creating a Wrapper (due on Friday)	Professor Knoblock
<b>Week 3</b> Jan 26	<b>SPARQL query language</b>	<p>Steve Harris and Andy Seaborne. Sparql 1.1 query language. Technical report, W3C, January 2012. <a href="http://www.w3.org/TR/2012/PR-sparql11-query-20121108/">http://www.w3.org/TR/2012/PR-sparql11-query-20121108/</a>.</p>	Quiz 2 (on Monday)  Homework 2: RDF (due on Friday)	Professor Ambite
Jan 28	<b>OWL 2</b>	<p>Krtzsch Markus, Simancik Frantisek, and Horrocks Ian. A description logic primer. 2012.  <a href="http://arxiv.org/pdf/1201.4089.pdf">http://arxiv.org/pdf/1201.4089.pdf</a>.</p>		Professor Ambite
<b>Week 4</b> Feb 2	<b>Linked Data</b>	<p>Aduna B.V. Http communication protocol for sesame 2. In System documentation for Sesame 2.x, chapter 8. October 2013.  <a href="http://www.csee.umbc.edu/courses/graduate/691/spring14/01/examples/sesame/openrdf-sesame-2.6.10/docs/system/ch08.html">http://www.csee.umbc.edu/courses/graduate/691/spring14/01/examples/sesame/openrdf-sesame-2.6.10/docs/system/ch08.html</a>.</p> <p>Chimezie Ogbuji. Sparql 1.1 graph store http protocol. Technical report, W3C, May 2012.  <a href="http://www.w3.org/TR/sparql11-http-rdf-update/">http://www.w3.org/TR/sparql11-http-rdf-update/</a>.</p>	Quiz 3 (on Monday)  Homework 3 SPARQL (due on Friday)	Professor Knoblock
Feb 4	<b>Data Cleaning</b>	<p>Wrangler: Interactive visual specification of data transformation scripts. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, 2011. <a href="http://vis.stanford.edu/papers/wrangler">http://vis.stanford.edu/papers/wrangler</a>.</p> <p>Bo Wu, Pedro Szekely, and Craig A. Knoblock. Minimizing user effort in transforming data by example. In Proceedings of the International Conference on Intelligent User Interface, 2014.  <a href="http://www.isi.edu/integration/papers/wu14-iui.pdf">http://www.isi.edu/integration/papers/wu14-iui.pdf</a>.</p> <p>Open Refine, Explore data.  <a href="http://youtu.be/B70J_H_zAWM">http://youtu.be/B70J_H_zAWM</a>.</p> <p>Open Refine, Clean and transform data.  <a href="http://youtu.be/cO8NVCs_Ba0">http://youtu.be/cO8NVCs_Ba0</a>.</p> <p>Open Refine, Reconcile and match data.  <a href="http://youtu.be/5tsyz3ibYzk">http://youtu.be/5tsyz3ibYzk</a>.</p>		Mr. Wu

<b>Week 5</b> Feb 9	<b>Database theory basics</b>	AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapter 2.1, 2.2, 2.3 and 2.4. Morgan Kaufmann, 2012.	Quiz 4 (on Monday)	Professor Ambite
	<b>Logical Data Integration</b>	<a href="http://www.sciencedirect.com/science/book/9780124160446">http://www.sciencedirect.com/science/book/9780124160446</a>	Homework 4: OWL (due on Friday)	Professor Ambite
<b>Week 6</b> Feb 18	<b>Scalable data integration</b>	Alon Halevy and Rachel Pottinger. A scalable algorithm for answering queries using views. The VLDB Journal The International Journal on Very Large Data Bases, 2001. <a href="http://www.vldb.org/conf/2000/P484.pdf">http://www.vldb.org/conf/2000/P484.pdf</a> .  Scalable query rewriting: a graph-based approach, 2001. <a href="http://www.isi.edu/~ambite/konstantinidis2011-sigmod.pdf">http://www.isi.edu/~ambite/konstantinidis2011-sigmod.pdf</a> .	Quiz 5 (on Monday)  Homework 5: Data Cleaning (due on Friday)	Mr. Konstantinidis
<b>Week 7</b> Feb 23	<b>Linked Services</b>	Barry Norton and Reto Krummenacher. Consuming dynamic linked data. In Proceedings of the 1st International Workshop on Consuming Linked Data, 2010. <a href="http://ceur-ws.org/Vol-665/NortonEtAlCOLD2010.pdf">http://ceur-ws.org/Vol-665/NortonEtAlCOLD2010.pdf</a> .  Mohsen Taheriyani, Craig A. Knoblock, Pedro Szekely, and Jose Luis Ambite. Rapidly integrating services into the linked data cloud. In Proceedings of the 11th International Semantic Web Conference (ISWC 2012), 2012. <a href="http://www.isi.edu/integration/papers/taheriyani12-iswc.pdf">http://www.isi.edu/integration/papers/taheriyani12-iswc.pdf</a> .	Quiz 6 (on Monday)  Homework 6: Logical Data Integration (due on Friday)	Mr. Taheriyani
	<b>Semi-Automatic Source Modeling</b>	Craig A. Knoblock, Pedro Szekely, Jose Luis Ambite, , Aman Goel, Shubham Gupta, Kristina Lerman, Maria Muslea, Mohsen Taheriyani, and Parag Mallick. Semi-automatically mapping structured sources into the semantic web. In Proceedings of the Extended Semantic Web Conference, Crete, Greece, 2012. <a href="http://www.isi.edu/integration/papers/knoblock12-eswc.pdf">http://www.isi.edu/integration/papers/knoblock12-eswc.pdf</a> .		Professor Knoblock
<b>Week 8</b> Mar 2	<b>RDF Mapping Tools</b>	[2] R2rml: Rdb to rdf mapping language. <a href="http://www.w3.org/TR/r2rml/">http://www.w3.org/TR/r2rml/</a> .	Quiz 7 (on Monday)	Professor Ambite
	<b>Schema Mapping</b>	Jayant Madhavan, Shawn R. Jeffery, Shirley Cohen, Xin (Luna) Dong, David Ko, Cong Yu, and Alon Halevy. Web-scale data integration: You can only afford to pay as you go, 2007. <a href="http://www.docin.com/p-47000224.html">http://www.docin.com/p-47000224.html</a> .  AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapter 5. Morgan Kaufmann, 2012. <a href="http://www.sciencedirect.com/science/book/9780124160446">http://www.sciencedirect.com/science/book/9780124160446</a>	Homework 7: Triple Stores (due on Friday)	Professor Ambite
<b>Week 9</b> Mar 9	<b>Source Modeling</b>	Mark James Carman and Craig A. Knoblock. Learning semantic descriptions of web information sources. In Proceedings of the Twentieth International Joint	Quiz 8 (on Monday)	Professor Ambite

Mar 11	<b>String Matching</b>	<p>Conference on Artificial Intelligence (IJCAI), January 2007.  <a href="http://www.isi.edu/integration/papers/carman07-ijcai.pdf">http://www.isi.edu/integration/papers/carman07-ijcai.pdf</a>.</p> <p>José Luis Ambite, Sirish Darbha, Aman Goel, Craig A. Knoblock, Kristina Lerman, Rahul Parundekar, and Thomas Russ. Automatically constructing semantic web services from online sources. In Proceedings of the 8th International Semantic Web Conference (ISWC 2009), 2009.  <a href="http://www.isi.edu/integration/papers/ambite09-iswc.pdf">http://www.isi.edu/integration/papers/ambite09-iswc.pdf</a>.</p> <p>[1] AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapters 4. Morgan Kaufmann, 2012.  <a href="http://www.sciencedirect.com/science/book/9780124160446">http://www.sciencedirect.com/science/book/9780124160446</a></p>	Homework 8: Karma (due on Friday)	Professor Knoblock
<b>Week 10</b> Mar 23  Mar 25	<b>Record Linkage</b>  <b>Data Warehousing</b>	<p>[1] AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapters 7. Morgan Kaufmann, 2012.  <a href="http://www.sciencedirect.com/science/book/9780124160446">http://www.sciencedirect.com/science/book/9780124160446</a></p> <p>AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapter 10. Morgan Kaufmann, 2012.  <a href="http://www.sciencedirect.com/science/book/9780124160446">http://www.sciencedirect.com/science/book/9780124160446</a></p>	Quiz 9 (on Monday)  Homework 9: String Similarity (due on Friday)	Professor Knoblock  Professor Ambite
<b>Week 11</b> Mar 30  Apr 1	<b>Mashup principles</b>  <b>Mashup tools</b>	<p>Shubham Gupta and Craig A. Knoblock. Building geospatial mashups to visualize information for crisis management. In Proceedings of the 7th International Conference on Information Systems for Crisis Response and Management, 2010.  <a href="http://www.isi.edu/integration/papers/gupta10-iscram.pdf">http://www.isi.edu/integration/papers/gupta10-iscram.pdf</a>.</p> <p>Jeffrey Wong and Jason I. Hong. Making mashups with marmite: towards end-user programming for the web. In ACM SIGMOD Record, 2007.  <a href="http://repository.cmu.edu/cgi/viewcontent.cgi?article=1063&amp;context=hcii">http://repository.cmu.edu/cgi/viewcontent.cgi?article=1063&amp;context=hcii</a>.</p> <p>Rob Ennals, Eric Brewer, Minos Garofalakis, Michael Shadle, and Prashant Gandhi. Intel mash maker: join the web. 2007. <a href="http://23.30.224.201/publications/intel-mash-maker-join-web">http://23.30.224.201/publications/intel-mash-maker-join-web</a>.</p> <p>Huynh David, Mazzocchi Stefano, and Karger David. Piggy bank: Experience the semantic web inside your web browser. 2007. <a href="http://simile.mit.edu/papers/iswc05.pdf">http://simile.mit.edu/papers/iswc05.pdf</a>.</p>	Quiz 10 (on Monday)  Homework 10: Record Linkage (due on Friday)	Professor Knoblock  Professor Knoblock

		Leo Sauermann and Richard Cyganiak. Cool uris for the semantic web. Technical report, 2008. <a href="http://www.w3.org/TR/cooluris/">http://www.w3.org/TR/cooluris/</a> .		
<b>Week 12</b> Apr 6	<b>Information Extraction</b>	Matthew Michelson and Craig A. Knoblock. Semantic Annotation of Unstructured and Ungrammatical Text. In Proceedings of the 19th International Joint Conference on Artificial Intelligence (IJCAI-2005), Edinburgh, Scotland, 2005. <a href="http://www.isi.edu/integration/papers/michelson05-ijcai.pdf">http://www.isi.edu/integration/papers/michelson05-ijcai.pdf</a>  Andrew McCallum. Information Extraction: Distilling Structured Data from Unstructured Text . ACM Queue, volume 3, Number 9, November 2005. <a href="http://people.cs.umass.edu/~mccallum/papers/acm-queue-ie.pdf">http://people.cs.umass.edu/~mccallum/papers/acm-queue-ie.pdf</a>	Quiz 11 (on Monday)  Homework 11: Mashups (due on Friday)	Professor Knoblock
Apr 8		Charles Elkan, Tutorial on Log-linear Models and Conditional Random Fields. <a href="http://videlectures.net/cikm08_elkan_llmacrf/">http://videlectures.net/cikm08_elkan_llmacrf/</a>		Professor Knoblock
<b>Week 13</b> Apr 13	<b>OWL Profiles</b>	Diego Calvanese, Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini, and Riccardo Rosati. DI-lite: tractable description logics for ontologies. In Proc. of the 20th National Conference on Artificial Intelligence, 2005. <a href="http://www.aaai.org/Papers/AAAI/2005/AAAI05-094.pdf">http://www.aaai.org/Papers/AAAI/2005/AAAI05-094.pdf</a> .	Quiz 12 (on Monday)  Homework 12: Information Extraction (due on Friday)	Professor Ambite
Apr 15	<b>Wrapper Learning</b>	Ion Muslea, Steve Minton, and Craig A. Knoblock. A hierarchical approach to wrapper induction. In Proceedings of the 3rd International Conference on Autonomous Agents, Seattle, WA, 1999. <a href="http://www.isi.edu/integration/papers/muslea99-agents.pdf">http://www.isi.edu/integration/papers/muslea99-agents.pdf</a> .  AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapter 9. Morgan Kaufmann, 2012. <a href="http://www.sciencedirect.com/science/book/9780124160446">http://www.sciencedirect.com/science/book/9780124160446</a>		Professor Knoblock
<b>Week 14</b> Apr 20	<b>Ontology-based Data Integration</b>	Hector Prez-Urbina, Ian Horrocks, and Boris Motik. Efficient query answering for owl 2. In International Semantic Web Conference, 2009. Efficient Query Answering for OWL 2. <a href="https://www.cs.ox.ac.uk/boris.motik/pubs/puhm09query-OWL2.pdf">https://www.cs.ox.ac.uk/boris.motik/pubs/puhm09query-OWL2.pdf</a> .	Quiz 13 (on Monday)	Professor Ambite
Apr 22	<b>Wrapper Generation</b>	W. Crescenzi, G. Mecca, and P. Merialdo. RoadRunner. Towards automatic data extraction from large web sites. 2001. <a href="http://www.vldb.org/conf/2001/P109.pdf">http://www.vldb.org/conf/2001/P109.pdf</a> .  B. Cenk Gazen and Steven Minton. Overview of autofeed: An unsupervised learning system for generating webfeeds. In Proceedings of AAAI, 2006.		Professor Knoblock

		<a href="http://www.isi.edu/integration/courses/csci548/Papers/gazen06-aaai.pdf">http://www.isi.edu/integration/courses/csci548/Papers/gazen06-aaai.pdf</a> .		
<b>Week 15</b> Apr 27	<b>Intellectual Property</b>	Thomas P. Vartanian and Robert H. Ledig. Scrape it, scrub it and show it: The battle over data aggregation. <a href="http://web.archive.org/web/20070818130311/http://www.ffhsj.com/bancmail/bmarts/aba%20art.html">http://web.archive.org/web/20070818130311/http://www.ffhsj.com/bancmail/bmarts/aba art.html</a> .	Quiz 14	Prof. Knoblock
Apr 29	<b>Student Presentations</b>	Kembrew McLeod. Intellectual property law, freedom of expression, and the web, 2003. <a href="http://www.electronicbookreview.com/thread/technocapitalism/proprietary">http://www.electronicbookreview.com/thread/technocapitalism/proprietary</a> .  Electronic frontier foundation. <a href="http://www.eff.org/issues/intellectual-property">http://www.eff.org/issues/intellectual-property</a> .		Students
<b>FINAL</b> May 6 4:30-6:30pm	<b>Final Exam</b>		During assigned time in the <i>Schedule of Classes</i> at <a href="http://www.usc.edu">www.usc.edu</a> u/soc.	

## Statement on Academic Conduct and Support Systems

### Academic Conduct

Plagiarism – presenting someone else’s ideas as your own, either verbatim or recast in your own words – is a serious academic offense with serious consequences. Please familiarize yourself with the discussion of plagiarism in *SCampus* in Section 11, *Behavior Violating University Standards* <https://scampus.usc.edu/1100-behavior-violating-university-standards-and-appropriate-sanctions>. Other forms of academic dishonesty are equally unacceptable. See additional information in *SCampus* and university policies on scientific misconduct, <http://policy.usc.edu/scientific-misconduct>.

Discrimination, sexual assault, and harassment are not tolerated by the university. You are encouraged to report any incidents to the *Office of Equity and Diversity* <http://equity.usc.edu> or to the *Department of Public Safety* <http://capsnet.usc.edu/departments/departments-public-safety/online-forms/contact-us>. This is important for the safety of the whole USC community. Another member of the university community – such as a friend, classmate, advisor, or faculty member – can help initiate the report, or can initiate the report on behalf of another person. *The Center for Women and Men* <http://www.usc.edu/student-affairs/cwm/> provides 24/7 confidential support, and the sexual assault resource center webpage <http://sarc.usc.edu> describes reporting options and other resources.



## Support Systems

A number of USC's schools provide support for students who need help with scholarly writing. Check with your advisor or program staff to find out more. Students whose primary language is not English should check with the *American Language Institute* <http://dornsife.usc.edu/ali>, which sponsors courses and workshops specifically for international graduate students. *The Office of Disability Services and Programs* [http://sait.usc.edu/academicsupport/centerprograms/dsp/home\\_index.html](http://sait.usc.edu/academicsupport/centerprograms/dsp/home_index.html) provides certification for students with disabilities and helps arrange the relevant accommodations. If an officially declared emergency makes travel to campus infeasible, *USC Emergency Information* <http://emergency.usc.edu> will provide safety and other updates, including ways in which instruction will be continued by means of blackboard, teleconferencing, and other technology.