# USC Viterbi School of Engineering

**INF 558/CSCI 563: Building Knowledge Graphs**
**Units: 4**
**Term—Day—Time:**
Fall 2018 – Friday – 2:00 - 5:20pm

**Location:** VKC 150

**Instructor:** Pedro Szekely
**Office:** ISI
**Office Hours:** After each class in VKC 150, or at ISI by appointment
**Contact Info**: szekely@usc.edu, 310-448-8641.

**Instructor:** Jay Pujara
**Office:** RTH 512
**Office Hours:** After each class in VKC 150, or by appointment
**Contact Info:** jpujara@usc.edu, 310-448-8482.

**Teaching Assistant:** Binh Vu
**Office:** SAL Lab
**Office Hours:** Tu, 4:00-5:30pm
**Contact Info:** binhvu@isi.edu

## Catalogue Course Description

Foundations, techniques, and algorithms for building knowledge graphs and doing so at scale. Topics include information extraction, data alignment, entity linking, and the Semantic Web.

## Expanded Course Description

This course focuses on foundations, techniques, and algorithms for building knowledge graphs. Students will learn the theory and applications of the techniques needed to build and query massive knowledge graphs. Topics include crawling web sites, wrapper learning, information extraction, source alignment, string matching, entity linking, graph databases, querying knowledge graphs, data cleaning, Semantic Web, linked data, graph analytics, and intellectual property. 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 algorithms and techniques for crawling web sites, structured data extraction, and information extraction from unstructured text.
* Understand the theory and techniques for cleaning, aligning, matching, and linking data.
* Understand the foundations and techniques of the Semantic Web, including RDF, ontologies, SPARQL, and linked data.
* Understand how to work with graph databases, including how to load massive datasets into such databases, how to organize the data for efficient access, and how to efficiently query the contents.
* Understand the entire process of how to design, construct, and query a knowledge graph to solve real-world problems.

- Understand how to apply the big data tools and infrastructure (e.g., Spark) to build and query knowledge graphs.

## Required Preparation:

Prerequisite(s):   INF 551 or CSCI 585
                          INF 552 or CSCI 567

Recommended Background: Experience programming in Python

## Course Notes

The course will be run as a lecture class with student participation strongly encouraged. The first 4-5 weeks of the course are structured as a quickstart to provide a shallow primer on the end-to-end process of knowledge graph construction, followed by deeper presentations and more technical material for the remainder of the course. 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, and homeworks will be posted online (https://bit.ly/2Llz3Cx).   The class project is a significant aspect of this course and at the end of the semester students will present their projects in class.

## Required Readings and Supplementary Materials

Required Textbook: none
We use a set of technical papers and book chapters that are all available online.  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 11 weeks of class.  The assignments must be done individually.  The homework assignments are expected to take 8-10 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 present their project proposals in class, conduct the project, and then create a video demonstration of the work and present the project in class.

*Project Timeline:*
- Week 4:  Project proposals presented in class (team members, topic)
- Week 8: Project status update due (1 page status report)
- Week 12: Project status update due (1 page status report)
- Week 15: Project presentation in class (short talk and video demonstration)

*Project description:*  Each project team will build a knowledge graph for a topic of their choice.  The knowledge graph must combine data from at least 3 different sources and at least 2 of those sites must be from online web sites.  The best projects build on many of the topics covered in the class.   The homeworks have been designed so that you can work on your projects in the process of doing your homework.

An example project would be to build a knowledge graph of used bicycles that could be purchased near the USC campus.  This project would combine data from used sources, such as Craig's List, new bike sources such as BikeNashbar, and bicycle review sites, such as bicycling.com.  The project would collect the data from each of these sources using wrapper techniques, extract the details of the used bicycle ads from Craig's List using information extraction techniques, align the data across these various sources to a domain ontology, link the entities across sources to combine the used data with the reviews from bicycling.com and

prices from BikeNashbar, store all of the data into a graph database such as elasticsearch, and then build a simple user interface to show the results by executing queries against the graph database.

*Grading breakdown of the course project:*
- Proposal: 10%
- Project video: 30%
- Presentation: 30%
- Overall project: 30%

## Grading Breakdown

**Quizzes:** There will be weekly quizzes based on the material from the week before. The lowest quiz grade will be dropped. Missed quizzes will receive a zero grade, and there will be no make-up quizzes for any reason.
**Midterm:** 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. The final exam will be on the date designated by USC in https://classes.usc.edu/term-20183/finals/
**Class Project:** Each student will do a group 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, create a video demonstration of the proof-of-concept, and present the project in class.

Grading Schema:

| | |
|---|---|
| Quizzes | 20% |
| Homework | 20% |
| Final: | 25% |
| Class Project | 35% |
| _____ | |
| Total | 100% |

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

| | |
|---|---|
| 94 - 100 = A | 74 – 76.9 = C |
| 90 – 93.9 = A- | 70 – 73.9 = C- |
| 87 – 89.9 = B+ | 67 – 69.9 = D+ |
| 84 – 86.9 = B | 64 – 66.9 = D |
| 80 – 83.9 = B- | 60 – 63.9 = D- |
| 77 – 79.9 =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 |
|---|---|---|---|---|
| **Week 1 Aug 24** | **Course Introduction & Use Case** | Pedro Szekely, et al. Building and using a knowledge graph to combat human trafficking. In Proceedings of the 14th | Homework 1: Crawling | Pujara Szekely |

| | | International Semantic Web Conference (ISWC 2015), 2015. http://iswc2015.semanticweb.org/sites/iswc2015.semanticweb.org/files/93670175.pdf<br><br>The Anatomy of a Large Scale Hypertextual Web Search Engine Sergey Brin and Lawrence Page, Seventh International World Wide Web Conference, 1998. http://ilpubs.stanford.edu:8090/361/1/1998-8.pdf<br><br>The Structure of the Web: http://science.sciencemag.org/content/sci/294/5548/1849.full.pdf<br><br>Optional: Web crawling and indexes: https://nlp.stanford.edu/IR-book/pdf/20crawl.pdf<br><br>Optional: Searching the Web Arvind Arasu et al., ACM Transactions on Internet Technology, 2001 http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.25.5183 | | |
| :-- | :-- | :-- | :-- | :-- |
| **Quickstart: Crawling** | | | | |
| **Week 2 Aug 31** | **Quickstart: Information Extraction** | D. C. Wimalasuriya and D. Dou. Ontology-based Information Extraction: An Introduction and a Survey of Current Approaches. J. Information Science, 36(3), 2010.<br><br>O. Etzioni, M. Banko, S. Soderland, and D. S. Weld. Open Information Extraction from the Web. Communications of the ACM, 51(12):68–74, 2008.<br><br><br>**Optional:**<br><br>J. Zhu, Z. Nie, X. Liu, B. Zhang, and J.-R. Wen. Statsnowball: A Statistical Approach to Extracting Entity Relationships. In Proceedings of the 18th International Conference on World Wide Web, pages 101–110. ACM, 2009.<br><br>S. Krause, H. Li, H. Uszkoreit, and F. Xu. Large-Scale Learning of Relation-Extraction Rules with Distant Supervision from the Web. | Quiz 1<br><br>Homework 2: Extracting information from web documents | Szekely |

| | | International Semantic Web Conference, pages 263–278, 2012.<br><br>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. http://www.isi.edu/integration/papers/muslea99-agents.pdf<br><br>AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapter 9. Morgan Kaufmann, 2012. http://www.sciencedirect.com/science/book/9780124160446<br><br>W. Crescenzi, G. Mecca, and P. Merialdo. RoadRunner. Towards automatic data extraction from large web sites. 2001. http://www.vldb.org/conf/2001/P109.pdf<br><br>B. Cenk Gazen and Steven Minton. Overview of autofeed: An unsupervised learning system for generating webfeeds. In Proceedings of AAAI, 2006. http://www.isi.edu/integration/courses/csci548/Papers/gazen06-aaai.pdf. | | |
| **Week 3 Sep 7** | **Quickstart: Knowledge Representation & Entity Linking** | A. Barr and J. Davidson. Representation of Knowledge, in Handbook of AI, volume 1, Chapter 3A-B, pages 141–160.<br><br>W. Cohen, P. Ravikumar, and S. Fienberg. A Comparison of String Distance Metrics for Name-matching Tasks. Conference on Information Integration on the Web, 2003.<br><br>J. Pujara and L. Getoor. Generic Statistical Relational Entity Resolution in Knowledge Graphs. StaRAI 2016.<br><br>Optional:<br><br>I. P. Fellegi and A. B. Sunter. A theory for record linkage. Journal of the American Statistical Association, 64(328):1183–1210, 1969.<br><br>W. E. Winkler. The state of record linkage and current research problems. In Statistical | Quiz 2<br><br>Homework 3: Representing knowledge and linking entities | Pujara |

| | | Research Division, US Census Bureau. Citeseer, 1999.

G. Navarro. A Guided Tour to Approximate String Matching. ACM Comput. Surv., 33 (1):31–88, Mar. 2001. ISSN 0360-0300. doi: 10.1145/375360.375365. URL http://doi.acm.org/10.1145/375360.375365

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. http://www.isi.edu/integration/papers/michelson05-ijcai.pdf

Andrew McCallum.  Information Extraction: Distilling Structured Data from Unstructured Text . ACM Queue, volume 3, Number 9, November 2005. http://people.cs.umass.edu/~mccallum/papers/acm-queue-ie.pdf

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. http://www.isi.edu/integration/papers/wu14-iui.pdf.

Open Refine, Explore data. http://youtu.be/B70J_H_zAWM.

Open Refine, Clean and transform data. http://youtu.be/cO8NVCs_Ba0.

Open Refine, Reconcile and match data. http://youtu.be/5tsyz3ibYzk.


G. Antoniou and F. Van Harmelen. A Semantic Web Primer. MIT press, 2004

F. van Harmelen, V. Lifschitz, and B. Porter. Handbook of Knowledge Representation. Elsevier Science, San Diego, USA, 2007. ISBN 0444522115, 9780444522115 | | |
| --- | --- | --- | --- | --- |

| | | P. Hitzler, M. Krotzsch, and S. Rudolph. Knowledge Representation for the Semantic Web. KI, 2009 | | |
|---|---|---|---|---|
| **Week 4 Sep 14** | **Quickstart: DIG, RDF, & Triplestores** | Wrangler: Interactive visual specification of data transformation scripts. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, 2011. http://vis.stanford.edu/papers/wrangler.<br><br>Frank Manola and Eric Miller. Rdf primer. Technical report, W3C, February 2004. http://www.w3.org/TR/2004/REC-rdf-primer-20040210/<br><br>Tim Berners-Lee. Why rdf model is different from the xml model. Technical report, W3C, 1998. http://www.w3.org/DesignIssues/RDF-XML.html.<br><br>Rdf vocabulary description language 1.0: Rdf schema. Technical report, W3C, February 2004. http://www.w3.org/TR/2004/REC-rdf-schema-20040210/<br><br>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. http://www.w3.org/TR/2012/NOTE-rdfa-primer-20120607/ | Quiz 3<br><br>Homework 4: Using DIG | Szekely |
| **Week 5 Sep 21** | **Quickstart: Querying & Analysis** | Oren, E., Delbru, R., Catasta, M., Cyganiak, R., Stenzhorn, H., & Tummarello, G. (2008). Sindice. com: a document-oriented lookup index for open linked data. *International Journal of Metadata, Semantics and Ontologies*, *3*(1), 37-52. http://wtlab.um.ac.ir/images/e-library/linked_data/other/Sindice.pdf<br><br>Freitas, A., Oliveira, J. G., Curry, E., O'Riain, S., & da Silva, J. C. P. (2011, June). Treo: combining entity-search, spreading activation and semantic relatedness for querying linked data. In *Proc. of 1st Workshop on Question Answering over Linked Data (QALD-1) at the 8th Extended Semantic Web Conference (ESWC 2011)*. | Quiz 4<br><br>Homework 5: KG Analysis | Pujara |

| | | https://www.deri.ie/sites/default/files/public ations/freitas_qald_2011_0.pdf<br><br>Steve Harris and Andy Seaborne. Sparql 1.1 query language. Technical report, W3C, January 2012. http://www.w3.org/TR/2012/PR-sparql11-query-20121108 | | |
|---|---|---|---|---|
| **Week 6 Sep 28** | **Large, Real-World KGs**<br><br><br>**Project Proposal** | S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. Dbpedia: A Nucleus for a Web of Open Data. The Semantic Web, 2007.<br><br>Bollacker K, Evans C, Paritosh P, Sturge T, Taylor J. Freebase: a collaboratively created graph database for structuring human knowledge. SIGMOD, 2008.<br><br>Optional:<br><br>X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmann, S. Sun, and W. Zhang. Knowledge Vault: A Web-Scale Approach to Probabilistic Knowledge Fusion. KDD, 2014.<br><br>F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: A Core of Semantic Knowledge. WWW, 2007.<br><br>F. Niu, C. Zhang, C. Re, and J. W. Shavlik. DeepDive: Web-Scale Knowledge-Base ́ Construction Using Statistical Learning and Inference. VLDS, 2012.<br><br>D. B. Lenat. Cyc: A Large-Scale Investment in Knowledge Infrastructure. Communications of the ACM, 38(11):33–38, 1995 | Quiz 5<br><br>Project Proposals Due | Pujara |
| **Week 7 Oct 5** | **Information Extraction** | D. C. Wimalasuriya and D. Dou. Ontology-based Information Extraction: An Introduction and a Survey of Current Approaches. J. Information Science, 36(3), 2010.<br><br>J. Zhu, Z. Nie, X. Liu, B. Zhang, and J.-R. Wen. Statsnowball: A Statistical Approach to Extracting Entity Relationships. In Proceedings of the 18th International Conference on World Wide Web, pages 101– | Quiz 6<br><br>Homework 7: Information Extraction II | Pujara |

| | | 110. ACM, 2009.<br><br>O. Etzioni, M. Banko, S. Soderland, and D. S. Weld. Open Information Extraction from the Web. Communications of the ACM, 51(12):68–74, 2008.<br><br>S. Krause, H. Li, H. Uszkoreit, and F. Xu. Large-Scale Learning of Relation-Extraction Rules with Distant Supervision from the Web. International Semantic Web Conference, pages 263–278, 2012.<br><br><br>AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapters 4. Morgan Kaufmann, 2012.<br>http://www.sciencedirect.com/science/book/9780124160446 | | |
| **Week 8**<br>**Oct 12** | **Structured Data** | AnHai Doan, Alon Y. Halevy, and Zachary G. Ives. Principles of Data Integration, chapters 7. Morgan Kaufmann, 2012.<br>http://www.sciencedirect.com/science/book/9780124160446<br><br><br>Limaye, G., Sarawagi, S., Chakrabarti, S.: Annotating and Searching Web Tables using Entities, Types and Relationships. Proc. VLDB Endow. 3(1-2), 1338-1347<br>http://www.vldb.org/pvldb/vldb2010/pvldb_vol3/R118.pdf<br><br>Cafarella, Michael J., Alon Halevy, Daisy Zhe Wang, Eugene Wu, and Yang Zhang. "Webtables: exploring the power of tables on the web." *Proceedings of the VLDB Endowment* 1, no. 1 (2008): 538-549.<br>http://tangra.si.umich.edu/~radev/767w10/papers/Week06/TextRepresentation/Cafarella.pdf<br><br>Crestan, Eric, and Patrick Pantel. "Web-scale table census and classification."<br>In *Proceedings of the fourth ACM international conference on Web search and data mining*, pp. 545-554. ACM, 2011.<br>http://www.patrickpantel.com/download/papers/2011/wsdm11.pdf | Quiz 7<br><br>Homework 8: Structured Data | Pujara |

| Week 9<br>Oct 19 | **Entity**<br>**Resolution** | I. P. Fellegi and A. B. Sunter. A theory for record linkage. Journal of the American Statistical Association, 64(328):1183–1210, 1969.<br><br>W. E. Winkler. The state of record linkage and current research problems. In Statistical Research Division, US Census Bureau. Citeseer, 1999.<br><br>G. Navarro. A Guided Tour to Approximate String Matching. ACM Comput. Surv., 33 (1):31–88, Mar. 2001. ISSN 0360-0300. doi: 10.1145/375360.375365. URL http: //doi.acm.org/10.1145/375360.375365<br><br>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. http://www.isi.edu/integration/papers/michelson05-ijcai.pdf<br><br>Andrew McCallum. Information Extraction: Distilling Structured Data from Unstructured Text . ACM Queue, volume 3, Number 9, November 2005. http://people.cs.umass.edu/~mccallum/papers/acm-queue-ie.pdf<br><br>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. http://www.isi.edu/integration/papers/wu14-iui.pdf.<br><br>Open Refine, Explore data. http://youtu.be/B70J_H_zAWM.<br><br>Open Refine, Clean and transform data. http://youtu.be/cO8NVCs_Ba0.<br><br>Open Refine, Reconcile and match data. http://youtu.be/5tsyz3ibYzk. | Quiz 8<br><br>Homework 9: Entity Resolution | Szekely |
| **Week 10** | **Ontologies &** | Frank Manola and Eric Miller. Rdf primer. | Quiz 9 | Szekely |

| Oct 26 | RDF | Technical report, W3C, February 2004. http://www.w3.org/TR/2004/REC-rdf-primer-20040210/ <br><br>Tim Berners-Lee. Why rdf model is different from the xml model. Technical report, W3C, 1998. http://www.w3.org/DesignIssues/RDF-XML.html. <br><br>Rdf vocabulary description language 1.0: Rdf schema. Technical report, W3C, February 2004. http://www.w3.org/TR/2004/REC-rdf-schema-20040210/ <br><br>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. http://www.w3.org/TR/2012/NOTE-rdfa-primer-20120607/ | Homework 10: Ontologies | |
| Week 11 Nov 2 | Semantic Modeling | Pham, M.; Alse, S.; Knoblock, C.; and Szekely, P, Semantic labeling: A domain-independent approach. In *Proceedings of the 15th International Semantic Web Conference*, 2016. http://usc-isi-i2.github.io/papers/pham16-iswc.pdf <br><br>Mark James Carman and Craig A. Knoblock. Learning semantic descriptions of web information sources. In Proceedings of the Twentieth International Joint Conference on Artificial Intelligence (IJCAI), January 2007. http://www.isi.edu/integration/papers/carman07-ijcai.pdf <br><br>Jos´e 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. http://www.isi.edu/integration/papers/ambite09-iswc.pdf <br><br>Craig A. Knoblock and Pedro Szekely. Exploiting semantics for big data integration. AI Magazine, 2015. http://usc-isi-i2.github.io/papers/knoblock15- | Quiz 10 <br><br>Homework 11: Semantic Modeling | |

| | | aimagazine.pdf<br><br>Krtzsch Markus, Simancik Frantisek, and Horrocks Ian. A description logic primer. 2012. http://arxiv.org/pdf/1201.4089.pdf.<br><br>R2rml: Rdb to rdf mapping language. http://www.w3.org/TR/r2rml/ | | |
|---|---|---|---|---|
| **Week 12 Nov 9** | **Graph Embeddings & Probabilistic Models** | J. Pujara, H. Miao, L. Getoor, and W. Cohen. Using Semantics & Statistics to Turn Data into Knowledge. AI Magazine, 36(1):65–74, 2015b<br><br>Antoine Bordes, Nicolas Usunier, Alberto GarciaDuran, Jason Weston, and Oksana Yakhnenko. 2013. Translating Embeddings for Modeling Multi-relational Data. In NIPS<br><br>Sebastian Riedel, Limin Yao, Andrew McCallum, and Benjamin M Marlin. 2013. Relation extraction with matrix factorization and universal schemas. In NAACL-HLT.<br><br>Richard Socher, Danqi Chen, Christopher D Manning, and Andrew Ng. 2013. Reasoning with neural tensor networks for knowledge base completion. In NIPS<br><br>Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning Entity and Relation Embeddings for Knowledge Graph Completion. In AAA<br><br>Random walk inference and learning in a large scale knowledge base<br>N Lao, T Mitchell, WW Cohen. EMNLP 2011<br>Incorporating vector space similarity in random walk inference over knowledge bases<br>M Gardner, P Talukdar, J Krishnamurthy, T Mitchell. EMNLP 2014.<br><br>M. Nickel, K. Murphy, V. Tresp, and E. Gabrilovich. A Review of Relational Machine Learning for Knowledge Graphs: From Multi-Relational Link Prediction to Automated Knowledge Graph Construction. | Quiz 11<br><br>Homework 12: Graph Embeddings | Pujara |

| | | | | |
|---|---|---|---|---|
| **Week 13 Nov 16** | **Search & Queries**<br><br>**Intellectual Property** | Rajaraman, J. Leskovec and J. D. Ullman, *Mining of Massive Datasets*, Cambridge University Press, 2012. http://infolab.stanford.edu/~ullman/mmds/ch10.pdf<br><br>Thomas P. Vartanian and Robert H. Ledig. Scrape it, scrub it and show it: The battle over data aggregation. http://web.archive.org/web/20070818130311/http:/www.ffhsj.com/bancmail/bmarts/aba_art.html.<br><br>Kembrew McLeod. Intellectual property law, freedom of expression, and the web, 2003. http://www.electronicbookreview.com/thread/technocapitalism/proprietary.<br><br>Electronic frontier foundation. http://www.eff.org/issues/intellectual-property. | Quiz 12 | Szekely |
| **Week 14 Nov 23** | **Thanksgiving break! No class.** | | | |
| **Week 15 Nov 30** | **Student Presentations** | | | Szekely Pujara Vu |
| **FINAL Dec 7 2-4pm** | **Final Exam** | | During assigned time in the *Schedule of Classes* at https://classes.usc.edu/term-20183/finals/ | |

# 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/department/department-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 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.