

Saket Choudhary

CONTACT INFORMATION	2677 Ellendale Place Apartment #218 Los Angeles California 90007	<i>Email:</i> saketkc@gmail.com <i>Homepage:</i> http://saket-choudhary.me [Google Scholar] [CrossValidated] [Github]
EDUCATION	University Of Southern California [USC] , Los Angeles, USA <i>PhD Student, Computational Biology and Bioinformatics</i>	[2014 – Ongoing]
	University Of Southern California [USC] , Los Angeles, USA <i>Masters in Computer Science (Data Science Track)</i> Department of Computer Science	[2016 – 2019 (Expected)]
	University Of Southern California [USC] , Los Angeles, USA <i>Masters in Statistics, Department of Mathematics</i>	[2016 – 2018]
	Graduate Diploma in Statistics , Royal Statistical Society, London, England <i>Cleared four out of five modules</i>	[2016 – 2017]
	Indian Institute of Technology Bombay [IITB] , Mumbai India <i>Bachelor of Technology, Master of Technology, Chemical Engineering</i> Masters Thesis: Pattern Recognition in Clinical Data	[2009 – 2014]
HONORS AND AWARDS	<ul style="list-style-type: none">- Provost Fellowship, awarded to outstanding incoming PhD students at USC- Gandhian Young Technological Innovation Award by Indian Institute of Management Ahmedabad, for designing a low cost water impurity detection device- Institute Technical Special Mention for three consecutive years at IITB- Undergraduate Research Award, for developing ‘Scilab on Cloud’- Kishor Vaignyanik Protsahan Yojana (KVPY) Fellowship by Department of Science and Technology (DST), Government of India- Homi Bhabha Young Scientists’ Gold Medal	[2014] [2013] [2010-12] [2012] [2007] [2005]
OLYMPIADS	<ul style="list-style-type: none">- Top 6 to be selected for Indian National Mathematics Olympiad(INMO), selection level exam for International Mathematical Olympiad(IMO)- Top 30 in Regional Mathematics Olympiad(RMO)- Top 250 in Indian National Physics Olympiad (INPhO)- Top 300 in Indian National Astronomy Olympiad (INAO)	[2008] [2009] [2009] [2009]
PUBLICATIONS /PREPRINTS	<ol style="list-style-type: none">1. Rahman, Syed Asad, Gilliean Torrance, Lorenzo Baldacci, Sergio Martínez Cuesta, Franz Fenninger, Nimish Gopal, Saket Choudhary, John May, Gemma L. Holliday, Christoph Steinbeck and Janet M Thornton. <i>Reaction Decoder Tool (RDT): extracting features from chemical reactions</i>. <i>Bioinformatics</i> 32, no. 13 (2016): 2065-2066. [Online]2. Choudhary, Saket, Leyla Garcia, Andrew Nightingale, and Maria-Jesus Martin. <i>BioJS-HGV Viewer: Genetic Variation Visualizer</i>. <i>bioRxiv</i> (2015): 032573. [Preprint]	

3. Syed, Parvez, Shabarni Gupta, **Saket Choudhary**, Narendra Goud Pandala, Apurva Atak, Annie Richharia, Heng Zhu et al. *Autoantibody Profiling of Glioma Serum Samples to Identify Biomarkers Using Human Proteome Arrays* Scientific reports 5 (2015). [Online]
4. Yachdav, Guy, Tatyana Goldberg, Sebastian Wilzbach, David Dao, Iris Shih, **Saket Choudhary**, Steve Crouch et al. *Anatomy of BioJS, an open source community for the life sciences*. eLife 4 (2015): e07009. [Preprint]
5. **Choudhary, Saket**, and Santosh B. Noronha. *GalDrive: Pipeline for comparative identification of driver mutations using the Galaxy framework*. bioRxiv (2014): 010538. [Preprint]
6. **Choudhary, Saket**, Vishnu Raj, K. Sanmugasundaram, Gyan Singh Patel, and Kannan Moudgalya. *Scilab on Cloud and Textbook Companion Project: A Web 2.0 Service for Open Source Education*. In 2013 International Conference on Cloud Computing and Big Data. [Online]
7. Gatkine, Pradip, Swati Gatkine, Sushanth Poojary, **Saket Choudhary**, and Santosh Noronha. *Development of piezo-electric sensor based noninvasive low cost Arterial Pulse Analyzer*. In Biomedical Engineering International Conference (BMEiCON), 2013 6th, pp. 1-4. IEEE, 2013. [Online]
8. Dilip Save, Yogesh, R. Rakhi, N. D. Shambhulingayya, Amit Srivastava, Manas Ranjan Das, **Saket Choudhary**, and Kannan M. Moudgalya. *Oscad: An open source EDA tool for circuit design, simulation, analysis and PCB design*. In Electronics, Circuits, and Systems (ICECS), 2013 IEEE 20th International Conference on, pp. 851-854. IEEE, 2013. [Online]

RESEARCH
EXPERIENCE

Detecting cancer in Histopathology Images, Research Intern May 2018 - July 2018
Guide: Dr. Radhakrishna Bettadapura Strand Life Sciences, India

I developed a core library for applied machine learning on histopathology images. Using unsupervised segmentation and random forests, our method performs at par with other deep learning approaches that have been applied to this problem in literature.

Evolution of post-transcriptional regulation, PhD Project Oct 2016 - Ongoing
Guide: Prof. Andrew Smith Computational Biology and Bioinformatics, USC

Post transcriptional gene regulation is a key mechanism that determines the final protein abundance levels. My work has focused on using statistical models to identify sites of ribosomal pausing. Another area of research has involved characterizing the evolutionary signature of RNA binding protein HuR. We have characterized how conserved HuR binding sites lead to induced mRNA stability.

Source: <https://github.com/saketkc/rna-seq-snakemake>;
<http://saketkc.github.io/riboraptor>

Tools for Motif Conservation Analysis, PhD Project May 2015 - Feb 2016
Guide: Prof. Anton Valouev Dept. of Preventive Medicine, Keck School of Medicine, USC

Motifs predicted by motif discovery tools can often not be the ‘true motifs’ and can have significant p-value(or E-values) for even ‘false motifs’. We hypothesized that a ‘true motif’ should exhibit high evolutionary conservation scores. MoCA makes use of PhyloP and Gerp scores to assess the conservation profile of motif bases.

We used MoCA to analyze ENCODE Chip-Seq datasets and found that the ‘true motifs’(ones which have been validated experimentally) do exhibit high conservation scores and that these are statistically significant when compared to the scores of flanking regions or randomly sampled regions.

Source: <https://github.com/saketkc/moca>
Poster: <https://doi.org/10.6084/m9.figshare.1565626.v5>

Predicting protein coding boundaries using Deep Learning
Course Project

Nov 2017 - Ongoing
USC

I explored how recurrent neural networks (RNNs) can be used to predict protein coding domains in a gene. A word embedding approach along with bi-directional LSTMs gave promising results using the entire pool of protein coding genes in human achieving an overall accuracy of 0.67. This model when used to predict protein coding domains in a different species, mouse, achieved an overall accuracy of 0.70 when tested on non-orthologous genes (where orthogonality implies a gene in mouse shares significant sequence from a human gene owing to descent from a common ancestor). Preprint: <https://doi.org/10.6084/m9.figshare.5902726.v1>

Pattern Recognition in Clinical Data, Masters Thesis

Apr 2013 - Jul 2013

Guide: Prof. Santosh Noronha

Dept. of Chemical Engineering, IIT Bombay

Awarded Outstanding Thesis Award

Multiple methods exist for determining oncogenic 'driver' mutations. These tools often have non overlapping predictions and input format is tool specific.

We developed a Galaxy based toolbox to run such prediction tools in parallel with a standard input format. The end results are presented as an intuitive heatmap indicating mutations which are predicted to be drivers by a majority of the tools. Preprint: "[GalDrive: Pipeline for comparative identification of driver mutations using the Galaxy framework](#)"

In a separate project, we analyzed proteomics data from Glioblastoma patients and predicted a smaller set of marker genes. Paper: "[Autoantibody Profiling of Glioma Serum Samples to Identify Biomarkers Using Human Proteome Arrays](#)".

Automated Mining of Reaction Patterns

May 2012 - Jul 2012, Jan 2014 - Mar 2014

Guide: Dr. Syed Asad Rahman

Dr. Dame Janet Thornton Lab, EMBL-EBI, Cambridge(UK)

EC-BLAST is a novel tool to compare enzymes and map reactions. We used clustering based approaches to highlight misclassified enzymes in the established enzyme classification system(EC).

We developed a web-service that facilitated automated job submissions to back end clusters at EBI that led to significant reduction in job runtime.

Next Generation Sequencing, Supervised Learning Project

Jul 2012 - Dec 2012

Guide: Prof. Santosh Noronha

Dept. of Chemical Engineering, IIT Bombay & ACTREC

We developed automated pipelines using Python to analyze whole genome sequence data of cancer tumors. As part of the project, I also contributed BWA and samtools wrappers to Biopython, a Python based open source library for bioinformatics.

Scilab On Cloud

May 2012 - Jul 2012

Guide: Prof. Kannan Moudagalaya

Dept. of Chemical Engineering, IIT Bombay

Scilab is an open source software for numerical computation and is primarily command line/GUI based. We developed a back-end that allowed running Scilab through browser much like the modern

day IPython notebooks. This enabled accessing Scilab remotely, even on low configuration devices.

Presented at: [IEEE Conference Cloud Computing and Big Data \(CloudCom-Asia\), 2013](#)

PROFESSIONAL
EXPERIENCE

Google Summer of Code 2015 | Mixed Effect Models for *statsmodels* May 2015 - Jul 2015
Student Contract Developer

- *statsmodels* is a Python based library for statistical modeling
- Implemented IPython based notebooks illustrating varied applications of Mixed Effects Models
- Implemented likelihood ratio tests
- Progress Report: <http://statsmodels-mlm-gsoc2015.blogspot.com>

Google Summer of Code | BioJavascript Jul 2014 - Sep 2014
Student Contract Developer

- BioJavascript is an open source library to facilitate biological data visualization
- Developed 'Human Genetic Variation Viewer', a d3.js based component to visualize genetic variations from SNP databases
- Preprint: [BioJS-HGV Viewer: Genetic Variation Visualizer](#)

Google Summer of Code | Galaxy Project Jul 2013 - Sep 2013
Student Contract Developer

- Galaxy Project is an open source web-based platform used for reproducible bioinformatics analysis
- Implemented 'nested workflows' that allows users to run a workflow inside a workflow, obviating the need to replicate steps
- Added 'edit on the go' functionality to edit default parameters before runtime
- Progress Report: <http://galaxy-gsoc2013.blogspot.com>

Google Summer of Code | Connexions Project Jul 2012 - Sep 2012
Student Contract Developer

- Developed a Python module to allow embedding slide-shows in online notebooks
- Created functionality to add user defined quiz as an additional achievement
- Progress Report: <http://oerpub.github.io/oerpub.rhaptoslabs.slideimporter/>

OTHER PROJECTS

Solutions to various examinations in Statistics June 2015 - Ongoing
Personal Project

- Royal Statistical Society Examinations Solutions:
<http://www.saket-choudhary.me/rss-graduate-diploma-solutions/>
- *Piddling Pertinent* - Solutions to several trivial problems in statistics:
<http://www.saket-choudhary.me/pertinent-blog/>
- *Screening Exam Solutions* - Solutions to screening examinations held at USC:
<http://www.saket-choudhary.me/usc-math-505A-screening-solutions/>;
<http://www.saket-choudhary.me/usc-math-541A-screening-solutions/>;
<http://www.saket-choudhary.me/usc-math-541B-screening-solutions/>

sklearn-hogsvd Jan 2019 - Ongoing
Personal Project

- Scikit-learn compatible python implementation of higher order generalized singular value decomposition. <https://github.com/saketkc/sklearn-hogsvd>

pysradb Nov 2018 - Ongoing
Personal Project

- Python package for interacting with SRADB and downloading datasets from NCBI Sequence Read Archive (SRA). <https://github.com/saketkc/pysradb>

pyseqlogo Nov 2017 - Ongoing
Personal Project

- Python package to plot sequence logos <https://github.com/saketkc/pyseqlogo>

Image Analysis of Tuberculosis Samples Jan 2013 - Apr 2013
Supervised Learning Project, Collaborator: Hinduja Hospital, Mumbai

- Used image processing algorithms to detect probable cases of TB from sputum images
- Developed a user friendly GUI to aid histologists thus reducing the overall delay in analysis

Pratham, Student Satellite Program May 2010 - Oct 2010
India's First Students' Satellite Team, IIT Bombay

- Executed hardware testing of the On-board Computer system
- Implemented signal processing pipeline for communications subsystem

TEACHING
EXPERIENCE

- Teaching Assistant, Computer Programming and Utilization Fall 2011
- Teaching Assistant, Artificial Intelligence in Process Engineering Fall 2013
- Teaching Assistant, How the Body Works Fall 2017
- Teaching Assistant, How the Body Works Spring 2018

POSITIONS OF
RESPONSIBILITY

Web Manager, UG Academic Council Jul 2012 - Apr 2013

- Initiated a number of web portals, thus improving online accessibility of academic resources
- Awarded **Institute Organizational Award**

TechniC, Core Group Member Jul 2010 - Apr 2011

- Organized institute wide technical events; mentored students

STANDARDIZED
TEST SCORES

- GRE: Quantitative: 170/170 Verbal: 153/170 Analytical Writing: 3.5/6
- TOEFL: Reading: 29/30 Listening: 28/30 Speaking: 24/30 Writing: 28/30 Total: 109/120

RELEVANT
COURSEWORK AT
USC

- Machine Learning
- Deep Learning
- Wavelets
- Time Series Analysis
- Mathematical Statistics
- Applied Probability
- Numerical Analysis
- Analysis of Algorithms
- Introduction to Computational Biology
- Biostatistics
- Molecular Biology
- Seminar in Statistical Consulting
- Methods of Statistical Inference

RELEVANT
COURSEWORK THROUGH
COURSERA (VERIFIED)

- Machine Learning
- Mathematical Biostatistics Boot Camp 1
- Applied Logistic Regression
- Reproducible Research
- Case-Based Introduction to Bio-
- statistics
- Linear Algebra
- Exploratory Data Analysis