

HPC for Computational Genomics

The NCSA Genomics Group is a host for research into the use of high performance computing (HPC) for primary genomics analyses, such as alignment, variant calling, genome assembly, and RNASeq. By its nature, this research is highly collaborative. Every member of our team is affiliated with multiple departments or campus initiatives. The student participants in this group serve as a bond between the campus faculty using computational genomics analyses in their research, and the NCSA experts in HPC, storage, networking, databases, etc. Together we enable the use of advanced computing infrastructure in computational genomics. Explore this page to find out who is involved, how we are connected, and what projects are currently ongoing.

Staff members from the software directorate and other groups within research consulting are frequently collaborators on our projects.

NCSA Press:

[Crossing over, branching out: Meet the NCSA Genomics team](#)

[Engineering Open House Award](#)

[Collaborative efforts produce clinical workflows for fast, translational genetic analysis](#)

Table of Contents:

- [Active Projects](#)
- [Staff](#)
- [Graduate Students](#)
- [Undergraduate Students](#)
- [Alumni](#)
- [Other Collaborations](#)

Active Projects

Project name	Project description	Genomics staff	Collaborators
Biomedical Pipeline Development	Multiple projects involving workflow development and feature addition including determination on whether workflows can be deployed into the cloud.	Joshua Allen Mohith Manjunath Weihao Ge Raghd Alhazmy	The Mayo Clinic
CROPPS	GPU acceleration and refactoring of code for gRNA design for CRISPR assays.	David Bianchi	Cornell Computer Science Crop Science
NEAT	Software Development of sequence simulator with mutational models.	Josh Allen Raghd Alhazmy	The Mayo Clinic Ontario Cancer Research Center The Broad Institute University of Wyoming
Farm to Food Bank Mobile Application Development	Development of a mobile app for farmers to sell off-spec and extra produce to food banks.	Christina Fliege	Google

Christina Elizabeth Fliege

Technical Program Manager, National Center for Supercomputing Applications



NCSA Genomics, September 2017. Credit: Steve Deunsing





NCSA Genomics 'Best Original Undergraduate Research' Award

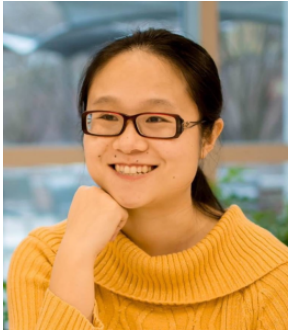




NCSA Genomics - Bluewaters tour, 17 June 2019





MINERVA	Developing an interface for combined genomic and diagnostic analysis for improved prediction, research and clinical interpretation of genomic variation.	David Bianchi Raghid Alhazmy	The Mayo Clinic
Investigator Portal	Developing a research compute portal for clinical diagnostics and genomics, where investigators can filter and query results against bioinformatic catalogs and applications, and generate re-useable datasets that can be viewed, analyzed and managed.	David Bianchi Misaël Lazaro	The Mayo Clinic
GWAS Study of Dairy Cattle	Analyze genotype data of over 11 cow farms, to find the cows' susceptibility to diseases and underlying genomic variants	Joshua Allen Weihao Ge	Prof. Sandra Rodríguez Zas
Metabolomics data analysis on microbiome	Understand gut microbiome products for personalized medicine.	Weihao Ge Misaël Lazaro David Bianchi	Prof. Issac Cann
genomic selection for maize and sorghum	Evaluate how much genomic variants will contribute to traits between maize and sorghum.	Weihao Ge	Prof. Alex Lipka


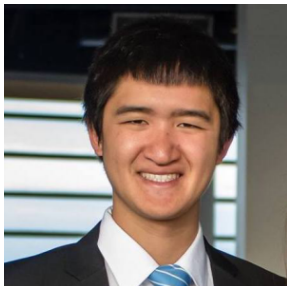

Staff




Research Interests		
	Joshua Allen Senior Research Programmer BA Mathematics and English (2001) MA English (2005) MS Bioinformatics (2019)	Sequence Simulation and Advanced Workflow Development. Modeling next-generation sequencing data, applying machine learning techniques to enhance models, writing production-ready code.
	David Bianchi Research Scientist Ph.D Physical Chemistry (2022)	Metabolic Engineering, Synthetic Biology, Gene Regulation, Genomics Analysis, Digital Agriculture, Spatial /Multi Omics, Personalized Medicine, Research Software Engineering, High-Performance Computing, GPU Computing
	Mohith Manjunath Research Programmer Ph.D Aersospace Engineering (2014)	Quantum Computing in Biology and Chemistry





	Weihao Ge B.S. Physics (2008) M.S. Physics (2011) Ph.D. Biophysics (2018)	Machine Learning in Epidemiology and Genomics. Biostatistics and Informatics.
	Raghid Alhazmy Research Programmer B.S. Biology (2021) M.S. Bioinformatics (2023)	Genomic data analysis, Machine Learning, and Computational Biology.
	Joao Paulo Gomes Viana	
	Misael Lazaro Academic Researcher B.S. Biochemistry (2019) M.S. Biochemistry (2023)	Genomics Data Analysis, Drug Development and Design, and Computational Biology
Graduate Students		
	Dhruvesh Shah School of Information	
	Yash Wasnik Information	Racial Health Disparities Yazhuo is involved in Racial Health Disparities project and researches with machine learning and data science skills. Her work is to do statistical analysis and write codes to build a pipeline on health datasets in collaboration with team members.
Undergraduate Students		



Alumni


	<p>Katherine Kendig</p> <p>Associate Project Manager</p> <p>B.A. Anthropology (2012)</p> <p>M.F.A. Creative Writing (2017)</p>	<h2>Project Management</h2> <p>Katherine is a project manager with the NCSA Industry Program, working primarily with biomedical partners.</p> <p>She benchmarked the Sentieon variant calling software for the Mayo Grand Challenge: https://www.biorxiv.org/content/10.1101/396325v1</p> <p>She has also contributed to NCSA's Public Affairs team, writing articles about NCSA and XSEDE research:</p> <p>After the storm; Bringing supercomputing to psychology; DISSCO Tech; ECSS: Profiles in Consulting; NCSA Genomics; History was here</p>
	<p>Ramshankar Venkatakrishnan</p> <p>Research Programmer</p> <p>B.S. Electronics & Communications (2012)</p> <p>M.S. Electrical & Computer Engineering (2015)</p>	<h2>Phillips 66 and Hardware support</h2> <p>Ram is developing code for the Phillips 66 project with the Data Analytics team. The idea of the code is to use Machine Learning to determine the best price to sell their petroleum products. The model considers a vast array of parameters to make the decision.</p> <p>Ram is working with the Innovative Systems Laboratory (ISL) at NCSA to create roofline model for a U250 Xilinx card using convolution as the code to plot the model.</p> <p>Ram also provides software and installation support for the HPC clusters at NCSA for a variety of clients.</p>
	<p>Dan Lanier, Research Programmer</p> <p>B.S. Applied Mathematics (2008)</p>	<h2>NCSA Industry</h2> <p>Dan supports biomedical partners in the NCSA Industry program.</p> <p>Dan provides a complementary mix of expertise in HPC and mathematical data analysis to enable pharmaceutical, agricultural and medical companies to utilize the high performance computing resources at NCSA.</p>
<p>blocked URL</p>	<p>Matthew Kendzior</p> <p>Research Programmer</p> <p>BS Crop Sciences (2016)</p> <p>MS Bioinformatics (2019)</p>	<h2>Mayo Grand Challenge</h2> <p>Mr. K is working as a researcher in the Mayo Grand Challenge, which aims to drastically speed up the time for detection of genomic variants, and to extract more information from whole genome sequencing data.</p> <h3>Genomic variant calling by assembly</h3> <p>Mr. K is focusing on a method to detect genomic variants by assembly.</p> <p>He is employing the software Cortex-var, which constructs de-novo genome assembly on multiple sequencing samples, and then compares the resultant de Bruijn graphs to detect where they diverge, indicating a potential variant. This could be a good method for detecting novel variants, especially repeats and complex rearrangements in complex genomes, such as polyploid plants and cancer. Mr. K is using his strong background in genomics to interpret, clean-up and validate the output.</p> <p>Mr. K is also working with Tiffany on the genomic analysis of HLHS for the Mayo Grand Challenge.</p> <p>Poster: Variant Calling by Assembly</p> <p>Poster: Reference-guided variant calling for non-repetitive sequences in Glycine Max</p>
	<p>Brian Bliss, Research Programmer</p>	<h2>Data compression</h2> <p>Brian will be working on data compression for the Mayo Grand Challenge project.</p>
	<p>Sushma Yellapragada</p> <p>Bachelor of Technology: Computer Science Engineering, Northcap University (2019)</p> <p>M.S. Computer Science, UIUC (2022)</p>	<h2>NEAT</h2> <p>Sushma is currently working on the NEAT project, contributing code and testing.</p>

	Angelo Santos	
	Yazhuo Zhang MS in Information Management	Racial Health Disparities <p>Yazhuo is involved in Racial Health Disparities project and researches with machine learning and data science skills. Her work is to do statistical analysis and write codes to build a pipeline on health datasets in collaboration with team members.</p>
	Sijia Huo B.S. Mathematics & Computer Science (2018) second major in Statistics third major in Economics	Parallelization of R <p>Sijia is working with NCSA Faculty Fellow Dr. Zeynep Madak-Erdogan to introduce parallel R code into her research.</p> <p>Dr. Madak-Erdogan is exploring racial disparities in breast cancer occurrence through the lens of diet and nutrition.</p>
	Ryan Chui B.S. Biochemistry (2016) M.S. Bioinformatics (2017) Department of Computer Science, UIUC	NCSA Industry <p>Ryan performed software installation, benchmarking, and development for a variety of industry partners. To investigate how the training time for deep neural networks (DNN's) can be affected, Ryan worked with TensorFlow, Google's deep learning library, to perform multi-label classification on a data set.</p> <p>He built an autoencoder – an unsupervised deep neural network - to extract salient features from the data.</p> <p>On Github:</p> <p>EpiQuant: Hadoop, C, Tensorflow - epistasis software prototypes</p> <p>MLCC - multi-label cancer classification</p> <p>q2b - binary representation of nucleotides</p> <p>ptgz - parallel tar gzip</p> <p>Usage Analyzer - log analyzer for HPC schedulers</p>
	Jennie Zermeno B.S. Integrative Biology (2017)	Benchmarking performance and accuracy of genomic variant calling software <p>Jennie collaborated to document our efforts in benchmarking variant calling on HPC systems. Jennie also participated in the debugging of the H3ABioNet GATK Germline Workflow.</p> Bioinformatics in the Cloud <p>Jennie is investigating the issues of portability, reproducibility and scaling of bioinformatics workflows in cloud infrastructure by instantiating containerized versions of workflows.</p> <p>Students Capitalize on Computational Genomics Research Using AWS</p>
blocked URL	Angela Chen M.S. Statistics (2017) Department of Statistics, UIUC CompGen fellow advised by Dr. Alexander Lipka	Accurate and scalable GWAS algorithms <p>Angela and Khory collaborated to improve the scalability and parallelization of the statistical software TASSEL5, widely used for conducting genome wide association studies (GWAS) in plants.</p> <p>Angela wrote a manuscript to demonstrate that her new stepwise epistatic model selection procedure has greater statistical power compared to other methods. However, the Java-based TASSEL5 cannot be easily parallelized across multiple nodes in a computational cluster, to run on modern, relevant datasets, which tend to be very large, such as the Alzheimer's SNP panel.</p>
	Khory Wagner advised by Dr. Vologymyr Kindratenko	<p>Khory provided the expertise in computer science to convert this Java code into C++ and parallelize it in HPC environment.</p>

blocked URL	<p>Nainika Roy</p> <p>B.S. Molecular and Cellular Biology (2017)</p> <p>minor in Informatics and Chemistry</p> <p>SPIN fellow</p>	<h2>Data formats and data structures in computational genomics</h2>
blocked URL	<p>Junyu Li</p> <p>B.S. Molecular and Cellular Biology (2017)</p> <p>minor in Computer Science</p> <p>SPIN fellow</p>	<h2>Genomic variant calling by assembly</h2> <p>Junyu worked with Mr. K in an interdisciplinary team, providing the expertise in math and computer science to automate the Cortex-var workflow and interpret the algorithm.</p> <p>Poster: Reference-guided variant calling for novel non-repetitive sequences in <i>Glycine max</i></p>
blocked URL	<p>Noah Flynn</p> <p>B.S. Bioengineering, Mathematics (2017)</p> <p>minor in computer science</p> <p>SPIN fellow</p>	<h2>Evolution of molecular networks and persistence of organisms</h2>
	<p>Jacob Heidenbrand</p> <p>Research Programmer</p> <p>B.S. Biochemistry (2014)</p> <p>M.S. Bioinformatics (2016)</p>	<h2>NCSA Industry</h2> <p>Jacob supports biomedical partners in the NCSA Industry program.</p> <p>Jacob provides a complementary mix of expertise in HPC and bioinformatics data analysis to enable pharmaceutical, agricultural and medical companies to utilize the high performance computing resources at NCSA.</p> <p>Jacob and Azza Ahmed (Ph. D. candidate, University of Khartoum) are exploring and evaluating the use of Swift T for variant calling.</p> <p>Github: Swift T Variant Calling</p> <p>Guide: Downloading large datasets with SRA Toolkit</p>
	<p>Matthew Weber</p> <p>B.S. Molecular and Cellular Biology (2016)</p> <p>M.S. Bioinformatics (2018)</p> <p>Department of Crop Sciences, UIUC</p> <p>CompGen fellow</p> <p>advised by Dr. Matthew Hudson</p>	<h2>Mutation profiles of cancer</h2> <p>Mr. Weber is developing machine learning methods to effectively stratify cancers based on the statistical properties of mutations found in afflicted individuals. Cancer stratification is predictive of disease outcomes, drug response and drug metabolism. Effective computational approaches based on total data acquired to-date can make this process cheaper in the clinic. Matt collaborates with the Ontario Institute for Cancer Research to make sure his models are realistic.</p> <p>Paper: Simulating Next-Generation Sequencing Datasets from Empirical Mutation and Sequencing Models</p> <p>Poster: Statistical models to capture mutational properties for NextGen Sequencing Data</p>
	<p>Aishwarya Raj</p> <p>B.S. Biochemistry (2019)</p> <p>minor in Bioinformatics</p> <p>Illinois Informatics Institute fellow</p>	<h2>Evolution of molecular networks and persistence of organisms</h2> <p>Construct and compare gene, metabolic and signaling networks from organisms across the tree of life.</p> <p>The goal of the project is to provide support for the general framework of persistence strategies.</p> <p>It postulates that persistence is achieved by biological systems via a tradeoff of traits that serve either economy, flexibility, or robustness. In this project we want to determine and quantify the molecular mechanisms that underlie these persistence strategies. Will analysis of the biomolecular networks allow us to differentiate between organisms of differing economy, flexibility, and robustness, and subsequently classify unknown, newly discovered, or modified organisms within such predefined classes?</p> <p>Poster: Persistence Strategies in Biomolecular Network Architecture</p> <p>NCUR Slides: Architecture and Dynamics of Biomolecular Networks Facilitate Evolution of Persistence Strategies in Living Organisms</p>


	<p>Cynthia Liu</p> <p>B.S. Bioengineering (2019)</p> <p>minor in Computer Science</p>	<h2>Workflow management comparisons</h2> <p>Cynthia worked to learn the Nextflow system for workflow management and to compare and contrast three competing workflow management options for bioinformatics in association with the work Ram is performing for the Mayo Grand Challenge.</p> <p>Poster: Comparative Analysis of Genomic Sequencing Workflow Management Systems</p>
	<p>Brian Rao</p> <p>B.S Integrative Biology (2018)</p> <p>Minor in Informatics</p>	<p>Brian wrote and tested the variant calling workflow code for the Mayo Grand Challenge. He focused on the accuracy and performance considerations of tumor variant detection in clinical settings.</p>
	<p>Angelynn Huang</p>	<p>Angelynn contributed to benchmarking the performance and accuracy of Minimap2 (Li, 2018) - a program used for analyzing sequencing read data in genomics.</p> <p>Minimap2 maps the sequencing reads against the reference genome for the species. Currently, BWA MEM (Li, 2013) is the most widely used tool for this purpose, with Novoalign (Hercus and Albertyn, 2012) coming as a close second. However, recent research (Li, 2018) suggests that Minimap2 is equally accurate yet also faster than BWA MEM. Are these claims true? Can we validate them independently using our own measurements? Sophia and Angelynn ran tests in AWS to answer these questions.</p> <p>Poster : Minimap2_BWA MEM</p> <p>Spotlight: http://www.ncsa.illinois.edu/news/story/ncsa_student_spotlight_angelynn_huang_and_sophia_torrellas</p>
<p>blocked URL</p>	<p>Sparsh Agarwal</p> <p>B.Tech + M.Tech in Biochemical Engineering and Biotechnology (2018)</p> <p>MS in Bioinformatics (2020)</p>	<p>Mayo Grand Challenge Project</p> <p>He is working on Mayo Grand Challenge project that aims to detect genomic variants in humans responsible for HLHS disease by using Cortex-var software as the de novo assembler and variant caller.</p>
	<p>Prakruthi Burra</p> <p>B. E. Computer Science (2018)</p> <p>M.S. Biological Sciences (2018)</p>	<h2>Human Heredity & Health in Africa</h2> <p>Prakruthi contributes to UIUC's work with the H3Africa Consortium. She is involved with projects on graph representations of genome assemblies and machine learning techniques applied to biological problems.</p> <h2>Workflow management for variant calling</h2> <p>Prakruthi is also implementing a variant calling workflow in Nextflow, an increasingly popular workflow manager. Prior to her workflow development work, she was briefly involved in testing the workflow developed for the Mayo Grand Challenge.</p>
	<p>Dave Istanto</p> <p>B.S. Crop Sciences (2018)</p>	<h2>Nextflow Cortex_Var Structural Variant Calling Workflow</h2> <p>Dave is responsible to develop a user-friendly and cluter-portable version of cortex_var workflow to detect large structural variants in given genomes using Nextflow workflow management language</p> <h2>Soybean Haplotype and Structural Variant Profiling and Analysis</h2> <p>Dave is responsible for both profiling of variants in 481 soybean lines, which later will be processed by correlating them to certain visible characteristics</p>



blocked URL	<p>Shubham Rawlani</p> <p>Bachelors in Electronics and Communication Engineering</p> <p>Masters in Information Management</p>	<h2>Space Search Reduction and EpiQuant</h2> <p>Shubham is involved in data analysis part where he writes code for data wrangling, extraction and cleaning to ease out the evaluation of statistical algorithms in the analysis of GWAS data for genomic variant epistasis</p> <p>Shubham is also involved in benchmarking the EpiQuant project and will collaborate to improve the scalability by testing on different datasets and nodes to achieve efficient results</p>
blocked URL	<p>Priya Balgi</p> <p>Bachelors in Information Technology Engineering</p> <p>Masters in Information Management</p>	<h2>Project Management</h2> <p>Priya is responsible for assisting in execution of Project Management tasks. Additionally, she performs genomics workflow testing using bash scripting in HPC environment and is developing a website using GitHub Pages/Jekyll for creation & auto-maintenance of project documentation.</p> <p>She also lead a student group of 8 for representing NCSA industry research during the Engineering Open House where the Genomics group won the Second Best Original Under Graduate Research Award and will also represent NCSA Industry research at the BioIT World Conference.</p> <p>Poster: NCSA Industry Research</p>
blocked URL	<p>Mingyu Yang</p> <p>B.E. Network Engineering</p> <p>M.S. Electrical and Computer Engineering</p>	<h2>Mayo Grand Challenge Project</h2> <p>Mingyu is working on optimize and test the performance of GABAC, which is a gene compression application.</p>
	<p>Yazhuo Zhang</p> <p>MS in Information Management</p>	<h2>Racial Health Disparities</h2> <p>Yazhuo is involved in Racial Health Disparities project and researches with machine learning and data science skills. Her work is to do statistical analysis and write codes to build a pipeline on health datasets in collaboration with team members.</p>
	<p>Dipro Ray</p> <p>B.S. Computer Science (2020)</p> <p>Minor in Mathematics</p>	<h2>Resolving Racial Disparities by Applying Statistics on Complex, Multidimensional Datasets</h2> <p>Dipro is working on turning a proof-of-concept prototype, of a statistical pipeline to analyze health data, into a well-structured open source package that is very portable, containerized and deployable through the cloud (like AWS), making such critical software available to researchers and collaborators with only a few commands.</p> <p>In pursuit of this goal, Dipro also works on refining the statistical pipeline in a modular manner and chalking out key design decisions for its implementation, and improving the package's computational efficiency (by making use of the host computer's architecture and resources)."</p>
blocked URL	<p>Tajesvi Bhat</p> <p>B.S. Computer Science (2020)</p> <p>Minor in Bioengineering</p>	<h2>Deployment of Variant Calling Workflows on Cloud Platform</h2> <p>Tajesvi is working on this that project aims to deploy variant calling workflows implemented using systems such as WDL and Nextflow in AWS and other cloud services.</p>

	<p>Tiffany Li B.S. Integrative Biology (2018) minor in Computer Science</p>	<p>Benchmarking performance and accuracy of genomic variant calling software</p> <p>Tiffany collaborates to document our efforts in benchmarking variant calling on HPC systems. We have run variant calling experiments on 500 genomes in parallel, on Blue Waters, to identify performance bottlenecks when using the GATK best practices workflow.</p> <p>We have also tested a number of alternative software, such as Isaac, Genallice, and Sentieon, as well as Dragen - a hardware solution. Tiffany is documenting the pros and cons of each of these excellent approaches in a separate manuscript.</p> <p>Validation and benchmarking on ParFu - a parallel file packaging utility</p> <p>Tiffany is also involved in testing and benchmarking of ParFu, an MPI tool for creating or extracting directory tree archives written by Dr. Craig Steffen, who works in the Blue Waters team.</p> <p>Github: Parfu Archive Tool</p>
---	--	--



Other Collaborations

blocked URL	<p>Dr. Matthew Hudson Bioinformatics Crop Science</p>	<p>HPCBio, Carver Biotechnology Center http://hpcbio.illinois.edu/</p>
blocked URL	<p>Dan Wickland Ph.D. Informatics (2019)</p>	
blocked URL	<p>Dr. Daniel Katz Computer Science</p>	<p>NCSA Scientific Software and Applications Portable variant calling workflow in Swift Github: Swift Variant Calling</p>
blocked URL	<p>Azza Ahmed Computer Science University of Khartoum advised by Dr. Faisal Fadlelmola</p>	
	<p>Dr. Zeynep Madak-Erdogan Food Science & Human Nutrition</p>	<p>Madak-Erdogan Lab</p> <p>Systems Biology of Estrogen Signaling</p> <ul style="list-style-type: none"> NCSA Faculty Fellow 2017-2018 Understanding Breast Cancer Disparities in African-American Women

	<p>Brandi Smith Ph.D. Food Science and Human Nutrition (2021)</p>	<p>H3Africa Consortium</p> <ul style="list-style-type: none"> • bioinformatics workflows in the cloud • custom genotyping chip for African populations • H3Africa bioinformatics node accreditation
	<p>Morgan Taschuk Bioinformatics blocked URL</p>	<p>OICR</p> <ul style="list-style-type: none"> • production infrastructure for primary genomics analyses • reproducibility of research in cancer genomics
<p>blocked URL</p>	<p>Paul Hatton HPC / Visualisation blocked URL</p>	<p>University of Birmingham</p>
<p>blocked URL</p>	<p>Nahil Sobh Machine Learning, AI</p>	<p>UIUC Beckman Institute Curriculum Vitae</p>
<p>blocked URL</p>	<p>Umberto Ravaoli Cyberinfrastructure, ECE</p>	<p>UIUC ECE, Beckman Institute Biosketch</p>
	<p>Lynn Hassan Jones Radiology</p>	<p>UIUC Resume</p>