

# Linhua WANG

## Computational Biologist | Multi-omics data scientist

+1 412 956 7722

linhuaw@bcm.edu

linkedin.com/in/linhuaw

Personal website

github.com/linhuawang

8181 Fannin St, Houston, TX 77054 - USA

I am passionate about using computational methods to advance biomedical discoveries from genomics and transcriptomics data. My Ph.D. research focused on developing techniques for analyzing spatial transcriptomics and multi-omics data from single cells, which has taught me the importance of collaborating with biologists and clinicians to create practical, usable tools. With a background in both biology and computer science, I have the skills and expertise to apply computational approaches such as Machine Learning, Deep Learning, and Convex Optimization to address biological questions. I have developed several tools and packages that have proven beneficial in a range of areas, including cancer research, genetic disease diagnosis, neurodegenerative diseases, and protein function predictions.

## EDUCATION

- |         |  |
|---------|--|
| Current | <b>Baylor College of Medicine</b><br>PhD. Quantitative & Computational Biosciences<br>Research : Spatial Transcriptomics, single-cell Multi-omics, AI-aided rare-disease diagnosis |
| 2017    | <b>Carnegie Mellon University</b><br>MSc. Computational Biology<br>Concentration : Machine Learning, Programming Language, Algorithm   |
| 2013    | <b>Sun Yat-sen University</b><br>BSc. Biotechnology<br>Concentration : Molecular Biology, Cell Biology, Cancer Biology   |

## RESEARCH EXPERIENCE

- |                          |  |
|--------------------------|--|
| Current<br>December 2019 | <b>Texas Children's Hospital - ZHANDONG LIU'S LAB   DATA SCIENCE - Graduate Student</b><br><i>Department of Pediatrics</i> – MENTOR : DR. ZHANDONG LIU <ul style="list-style-type: none"><li>Developed a computational tool, MIST, to detect molecular regions and impute gene expression values for Spatial Transcriptomics (ST) data (Published at <b>Nature Communications</b>)</li><li>Designed a strategy, ReSort, to generate an internal reference for ST's cell type deconvolution.</li><li>Simulated ST datasets and demonstrated that ReSort increased the accuracies of state-of-the-art reference-based ST deconvolution methods.</li><li>Identified macrophage polarization in Breast Cancer with Epithelial-Mesenchymal Transitions using ReSort.</li><li>Validated ReSort's discoveries using external The Cancer Genome Atlas (TCGA) datasets and immunohistochemical staining of cell type markers (Under review at <b>Genome Biology</b>)</li><li>Constructed a machine learning pipeline, MARRVEL AI, to aid rare genetic-disease diagnosis by collaborating with domain experts. Resulted in &gt;10% top-5 accuracies than other academic and commercial tools. (Oral presentation by coauthor at <b>ASHG, 2022</b>)</li><li>Published a bioinformatics Python package, SEAGAL, for Spatial Enrichment Analysis of Gene Associations using L-index, after understanding the needs by interviewing biologists.</li><li>SEAGAL allows identifying and visualizing spatial co-localization or exclusion of immune cell types at specific spatial niches (Submitted to <b>Bioinformatics</b>).</li></ul> <div>Computational Biology Spatial Transcriptomics Single-cell Multi-omics Genetic Diagnosis Data Visualization</div> |
| Sep 2022<br>May 2022     | <b>Ancestry - DNA SCIENCES - Genomics Data Scientist Intern</b><br><i>Applied Machine Learning</i> – MENTOR : DR. MILOŠ PAVLOVIĆ <ul style="list-style-type: none"><li>Processed feature vectors with 30 million rows and 3000 features using chunked data processing.</li><li>Reduced the number of features from &gt;2000 to &lt;100 using feature selection techniques.</li><li>Constructed scalable XGBoost-based machine learning pipelines that reduced the number of classifiers (&gt;2000) to a log-scale (&lt;20).</li><li>Delivered accurate geographical community assignments to 30 million customers based on genetic features using XGBoost with improved precision and recall scores.</li></ul> <div>DNA Science Pandas AWS Scikit-learn Sparse Matrix Processing</div>   |

April 2019	Icahn School of Medicine at Mount Sinai - THE PANDEY LAB   MACHINE LEARNING - Full-Time Bioinformatician
July 2017	<i>Department of Genetics and Genomic Sciences – SUPERVISOR : DR. GAURAV PANDEY</i> <ul style="list-style-type: none"> <li>➤ Built and maintained ensemble models combining 11 base WEKA classifiers, including Naive Bayes, SVM, Logistic Regression, and so on.</li> <li>➤ Scaled up the tool leveraging parallel and distributed computing on high-performance computing (HPC) systems.</li> <li>➤ Improved predictions' performance (F-score) on 277 gene ontology terms for 63,449 amino acid sequences from 19 clinically relevant bacterial pathogens.</li> <li>➤ Packaged and maintained the tool on GitHub and wrote Docker files for version control.</li> </ul> Machine Learning Version Control Protein Function Prediction High-Performance Computing (HPC)
Dec 2016	University of Pittsburgh - LU LAB   CANCER BIOLOGY - MS research
May 2016	<i>Department of Biomedical Informatics – MENTOR : DR. XINGHUA LU</i> <ul style="list-style-type: none"> <li>➤ Performed cancer marker discovery analysis for Brain Tumors, including Lower Grade Glioma and Glioblastoma samples from The Cancer Genome Atlas (TCGA) using R programming.</li> </ul> R programming Cancer Research Statistical Analysis

## PUBLICATIONS

### Unraveling Spatial Gene Associations with SEAGAL : a Python Package for Spatial Transcriptomics Data Analysis and Visualization UNDER REVIEW

[Linhua Wang](#), Chaozhong Liu, Zhandong Liu

Bioinformatics

Spatial Transcriptomics Immune Co-localization Bi-variate Spatial Correlation

### scGREAT : Graph-based regulatory element analysis tool for single-cell multi-omics data UNDER REVIEW

Chaozhong Liu, [Linhua Wang](#), Zhandong Liu

Bioinformatics

Single-cell Multi-omics Data Analysis Cis-regulatory Elements

### Accurate cell type deconvolution in spatial transcriptomics using a batch effect-free strategy UNVER REVIEW

[Linhua Wang](#), Ling Wu, Chaozhong Liu, Wanli Wang, Xiang H.-F. Zhang, Zhandong Liu

Genome Biology

Spatial Transcriptomics Cell Type Deconvolution Tumor Microenvironment

### Region-specific denoising identifies spatial co-expression patterns and intra-tissue heterogeneity in spatially resolved transcriptomics data 2022

[Linhua Wang](#), Mirjana Maletic-Savatic, Zhandong Liu

 [Nature Communications](#)

Spatial Transcriptomics Clustering Modularity Detection Imputation Low-rank approximation

### Single-cell multi-omics integration for unpaired data by a siamese network with graph-based contrastive loss 2023

Chaozhong Liu, [Linhua Wang](#), Zhandong Liu

 [BMC Bioinformatics](#)

Deep Learning Single-cell Multi-omics Integration kNN-graph Imputation

### Integrating multimodal data through interpretable heterogeneous ensembles 2022

Yan Chak Li, [Linhua Wang](#), Jeffrey N Law, T M Murali, Gaurav Pandey

 [Bioinformatics Advances](#)

Ensemble Learning Data Integration COVID-19 Mortality Prediction

### Predicting youth diabetes risk using NHANES data and machine learning 2022

Nita Vangeepuram, Bian Liu, Po-hsiang Chiu, [Linhua Wang](#), Gaurav Pandey

 [Scientific Reports](#)

Diabetes Machine Learning

### Large-scale protein function prediction using heterogeneous ensembles 2018

[Linhua Wang](#), Jeffrey Law, Shiv D Kale, TM Murali, Gaurav Pandey


 [F1000Research](#)

Ensemble Learning Machine Learning Protein Function Predictions

## SKILLS

Programming	Python, OOP, R, MATLAB, Bash, GitHub, Docker, LaTeX
Data Science	Pandas, ScanPy, AnnData, Scikit-Learn, AWS, TensorFlow, PyTorch
Data Visualization	Seaborn, Matplotlib, ggplot2
Research	Spatial Transcriptomics Analysis, Machine Learning, Computational Biology, Single-cell Multi-omics Analysis, Genetics
High Performance Computing	LSF, Distributed Computing, Parallel Computing
Soft Skills	Communication, Creativity, Critical Thinking

## INVITED TALKS

- Large-scale assessment of protein function prediction using heterogeneous ensembles 2018  
 [The 26th Intelligent Systems for Molecular Biology \(ISMB\)](#)  
[Linhua Wang](#), Jeffrey Law, Shiv D Kale, TM Murali, Gaurav Pandey

## POSITION OF RESPONSIBILITY

2022-present	<b>MDPI - Multidisciplinary Digital Publishing Institute</b> <i>Invited Reviewer</i> <ul style="list-style-type: none"><li>➢ Review at least two-paper per month for various journals at MDPI, including Cells, Cancers, Diagnostics, and so on.</li></ul>
2020	<b>MLCB - Machine Learning in Computational Biology</b> <i>Scientific Program Committee</i> <ul style="list-style-type: none"><li>➢ Invited for review of conference papers.</li></ul>

## ACHIEVEMENTS & RECOGNITIONS

- 2015-2017 **Academic Achievement Fellowship** by Department of Computational Biology, CMU  
2022 **Honorable Mentions** in 2022 Empower22 HACKATHON, Ancestry, Lehi, UT.  
2021 **Second place** (over 250 participants, subtask) in 2021 **Multimodal Single-Cell Data Integration** (NeurIPS Competition 2021)  
2018 **F1000 Outstanding Presentation Prize** (one in six) at **Intelligent Systems for Molecular Biology (ISMB)**, Chicago.  
2018 **Travel Fellowship** for **Intelligent Systems for Molecular Biology (ISMB)**, Chicago

## LANGUAGES

English ● ● ● ● ●  
Chinese ● ● ● ● ●