

Pengcheng Xu

+1(217)550-1337 | pengcheng.xu@uci.edu | <https://github.com/explcre> | <https://linkedin.com/in/pengcheng-xu-ryan>
Google Scholar: <https://scholar.google.com/citations?user=LXzbdZEAAAAJ&hl=en>
Personal Website: <https://explcre.github.io/> | ChatGPT version of me website: <https://explcre.github.io/mychat/>

EDUCATION

University of California Irvine Sep 2024 - 2029(Expected)
Ph.D. in Computer Science. GPA: 4.0/4.0
Research Interest: Large Language Model, Reinforcement Learning, Multi-modal Language Model, Bioinformatics, AI for Science, Large Language Model for Science, DNA Language Model, Single Cell, Spatial Transcriptomics Foundation Model.

University of Illinois Urbana-Champaign Aug 2022 - Dec 2023
Master's Degree in Computer Engineering, GPA: 3.8/4.0. Research Scientist at **AI@UIUC**.
Core Courses: Distributed Systems, Communication Network(A+), Computer Vision(A), Advanced Distributed Systems, Transfer Learning(A), Engineering Entrepreneurship(A+), Machine Learning(A), Applied Parallel Programming.

University of Michigan - Shanghai Jiao Tong University Joint Institute Sep 2018 - Aug 2022
Bachelor's Degree in Electrical and Computer Engineering. Minor in Data Science. Outstanding Graduate (School Level).
Core Courses: Data structures and Algorithms, Intro to Computer Organization (A, Teaching Assistant), Computational Methods for Statistics and Data Science (A+), Intro to Operating Systems, Intro to Data Science (A), Undergraduate Research (A+).

PUBLICATIONS

(Submitted to NeurIPS 2026) Pengcheng Xu, Junhao Liu, Peter Tan Dat Le, Haoyi Cai, Bruhati Aarushi Kuchi, Arnav Abhijit Dhariya, Chaoyang Wang, Hokin Deng, Jing Zhang. "DNAThinker: A Reasoning-Grounded Multimodal Foundation Model for Cell-Type-Conditioned DNA Generation, Editing, and Pair Prediction with Mask-Diffusion Reinforcement Learning".

(Under Review for Science Advances Journal) Pengcheng Xu, Junhao Liu, Yi Dai, Kainoa Andrew Nishida, Yutong Lei, Yaqi Hu, Chaoyang Wang, Jing Zhang. OmegaGenome: Toward Better Sub-Million Scale Expert Models from Large Genomic Language Models via Knowledge Distillation.

(Accepted to KDD 2026) Dongbo Sun, Junhao Liu, Pengcheng Xu, Yaqi Hu, Martin Jinye Zhang, Jing Zhang. [MUGO: Differentiable Combinatorial Optimization for Causal Variant Discovery in the Non-coding Genome](#).

(Accepted to ICML 2026) Maijunxian Wang, ..., [Pengcheng Xu](#), Hokin Deng. [A Very Big Video Reasoning Suite](#). <https://arxiv.org/abs/2602.20159>

(Under Review for some image signal processing conference) Yutong Lei, Yaqi Hu, Chaoyang Wang, Pengcheng Xu, Junhao Liu, Henry Cheng, Jing Zhang. SPATREG: A Hierarchical Framework for Histology Image Registration in Spatial Transcriptomics.

(To be Submitted to CV Conference) Xi Li, Ziheng Duan, Diptanshu Sikdar, Pengcheng Xu, Jing Zhang. BRIDGE: Spatially Aware Bimodal Representation Integration with Hierarchical Graphs for Spatial Transcriptomics Resolution Enhancement.

(Accepted to ACL-BioNLP Workshop 2025, Oral) Pengcheng Xu, Sicheng Shen, Jieli Zhou, Hongyi Xin. Team XSZ at BioLay-Summ2025: Section-Wise Summarization, Retrieval-Augmented LLM, and Reinforcement Learning Fine-Tuning for Lay Summaries. <https://aclanthology.org/2025.bionlp-share.33v2.pdf>

(Accepted to ACM BCB, ACM SIGBio Best Paper Award (1/204 accepted)) Ziqi Rong, Jinpu Cai, Jiahao Qiu, Pengcheng Xu, Lana Garmire, Qiuyu Lian, Hongyi Xin. "L2 Normalization and Geodesic Distance for Enhanced Information Preservation in Visualizing High-dimensional Single-cell Sequencing Data", 2023. <https://dl.acm.org/doi/abs/10.1145/3698587.3701361>

(Accepted to KDD-AIDSH workshop 2024, Oral) Pengcheng Xu, Tianfan Fu, Wenhao Gao, Jimeng Sun. "REINVENT-Transformer: Molecular de-novo design through Transformer-based Reinforcement Learning", 2024. <http://arxiv.org/abs/2310.05365>

(Accepted to ACL-BioNLP workshop 2024) Jieli Zhou*, Cheng Ye*, Pengcheng Xu, Hongyi Xin. "Adapting Large Language Models for Biomedical Lay Summarization". <https://aclanthology.org/2024.bionlp-1.76.pdf>

(Under Review) Pengcheng Xu*, Jinpu Cai*, Yulin Gao, Ziqi Rong, Hongyi Xin. "MIRACLE: Multi-task learning based Interpretable Regulation of Autoimmune diseases through Common Latent Epigenetics", 2024. <https://arxiv.org/abs/2306.13866>

(To be submitted) Jinpu Cai, Ziqi Rong, Luting Zhou, Xinzhu Jiang, Pengcheng Xu, Rui Gao, Yu Zhao, Bing He, Jianhua Yao, Qiuyu Lian*, Hongyi Xin* "Cell-ontology-aware marker identification and stratification with TAMER", 2023.

(To be submitted) Pengcheng Xu, Kaiyang Chen, Yuanrui Zhang, Indranil Gupta. "Pipe-Dejavu: Hardware-aware Latency Predictable, Differentiable Search for Faster Config and Convergence of Distributed ML Pipeline Parallelism", 2023. https://explcre.github.io/files/Pipe_Dejavu.pdf

Zhikai Yang, Pengcheng Xu, Dekun Yang, Yufeng Chen, Yancong Ma. "Vascular Intervention Training System Based on Electromagnetic Tracking Technology", ICVRV, 2020. <https://ieeexplore.ieee.org/document/9479727>

RESEARCH EXPERIENCE

DNAThinker: A Reasoning-Grounded Multimodal Foundation Model for Cell-Type-Conditioned DNA Generation, Editing, and Pair Prediction

Feb 2026 - Present

Lead researcher & primary author (I wrote most of the code + paper). NeurIPS 2026 submission.

- ▷ Built DNAThinker, a unified regulatory-DNA foundation model coupling a Qwen3.5-0.8B + LoRA reasoning LLM with an NTV3-100M-post + MDLM DNA generator through an SD3/FLUX-style cross-attention bridge and delta-residual conditioning path over a 1.5M-row fine-tuned regulatory-DNA corpus.
- ▷ Developed tool-grounded reasoning prompts with FIMO motif scans and per-cell expression context, so rationales cite verifiable transcription factors and JASPAR motifs across T1 enhancer generation, T2 pair prediction, and T3 enhancer editing.
- ▷ Applied strict-holdout grounded-RL post-training with TF-grounding reward and hallucination penalty; under the refined v2 detector on a scaled 7-cell holdout (n=700), improved TF-grounding from 0.702 to 0.919 on T1, 0.714 to 0.956 on T2, and 0.799 to 0.920 on T3, while reducing hallucinated TFs by 74%/67%/67%.
- ▷ Introduced vector-level TF-knockdown intervention on the 17,415-d cell-expression slot plus knockdown-aware policy-gradient training, achieving 63% conditional motif depletion on the Mic T3 test slice and 56% depletion / 71% depleting rows across the 7-cell stratified aggregate.
- ▷ Developed MDLM-RL, a policy-gradient framework for masked-diffusion DNA generators with GRPO, DAPO, GSPO, and SVGSPO variants; audited reward, algorithm, capacity, and learning-rate regimes, reporting T1 single-shot FID 8.78 (4.4x better than LEONINE), a 14x edit-tight T3 trunk-only motif lift, and clear saturated-trunk/no-RL regimes.

A Very Big Video Reasoning Suite: Public Task Generators and Dataset Construction

2025 - 2026

Contributing researcher. Accepted to ICML 2026; paper PDF: <https://arxiv.org/abs/2602.20159>.

- ▷ Implemented public Python data-generator repositories for [gear-chain direction](#), [arrow maze](#), [scale-balance tilt](#), [dot quantity comparison](#), [layer-depth ordering](#), [water-level reasoning](#), [binary counter](#), and [odd-one-out](#) task families.
- ▷ Built generation and validation logic that maps symbolic task states into controllable video prompts/assets and answer labels, supporting scalable data-suite construction for fine-grained temporal, spatial, counting, and physical reasoning evaluation.
- ▷ Helped with Wan video generation model fine-tuning experiments for dataset generation.

OmegaGenome: Toward Better Sub-Million Scale Expert Models from Large Genomic Language Models via Knowledge Distillation.

Mar 2025 - Feb 2026

Lead researcher & primary author (I wrote most of the code + paper).

- ▷ Designed and implemented OmegaGenome, a two-stage knowledge distillation framework that compresses billion-parameter genomic LMs – Nucleotide Transformer (NT-2.5B), Enformer (CNN-Transformer), DNABERT-2 (BPE tokenization), and Caduceus (BiMamba) – into sub-million-parameter expert models, achieving up to 20,000x compression while retaining 98% teacher performance across 18 genomic benchmarks.
- ▷ Fine-tuned all four teacher architectures using full-parameter training for Enformer, DNABERT-2, and Caduceus, and LoRA (rank=48) for NT-2.5B, reducing trainable parameters to <1% of the 2.5B backbone; distilled teachers into dilated CNN students (BPNetClassifier, 2K-3.6M params) via a composite objective combining KL divergence, MSE feature alignment, and cross-entropy.
- ▷ Benchmarked distilled students against teacher models and size-matched baselines, demonstrating 31-1,671x inference speedup, 7.5-513x GPU memory reduction, and consistent 6.8-10.8% MCC improvement over scratch-trained BPNet across epigenetic mark prediction, enhancer/promoter classification, and splice site detection.
- ▷ Implemented and compared four KD algorithms (OmegaGenome, DKD, DIST, Logit Standardization) with systematic grid search over 45+ configurations; quantified representation transfer via cosine similarity, t-SNE/UMAP, and logit-level R2 (distilled student 0.779 vs. scratch baseline 0.677), confirming faithful feature inheritance from teacher to student.

Adapting Large Language Models for Biomedical Lay Summarization

Apr 2024 - June 2024

- ▷ Fine-tuned Llama3 using Low-Rank Adaptation (LoRA) and optimized model performance, achieving a 68.8% improvement in the LENS readability score, leading to the first place in readability at the 2024 BioLaySumm workshop.
- ▷ Implemented K-shot prompting based on semantic similarity, enhancing factuality scores by 15.0% (AlignScore) and 7.9% (SummaC) through the integration of contextually relevant examples, leading to more accurate and relevant summaries.
- ▷ Developed techniques to resolve repeated word issues post-fine-tuning, resulting in a 12% increase in coherence and conciseness of the generated summaries.

Molecular de-novo design through Decision Transformer and Oracle-feedback reinforcement learning

May 2023 - May 2024

Advisor: Tianfan Fu (Assistant Professor at Rensselaer Polytechnic Institute), Jimeng Sun (Professor at CS, UIUC).

- ▷ Implemented a decision transformer architecture to improve the AUC for over fifteen molecular optimization tasks for 5% each on average.
- ▷ Applied Oracle-feedback reinforcement learning on the downstream tasks to reach higher performance than pretrained model.
- ▷ Carried out ablation study and investigation into loss curve and conditional probability over the next token as a function of previously chosen ones according to the model.

Hardware-aware Latency Predictable, Differentiable Search for Faster Config and Convergence of Distributed ML Pipeline Parallelism

Feb 2023 - May 2023

Advisor: Indranil Gupta, Professor of CS UIUC | Advanced Distributed Systems | Researcher.

- ▷ Implemented a predictive model that considers communication cost, model computational cost, and hardware information to predict latency and resources of parallel configurations, saving time on pre-profiling before searching the parallel configuration.
- ▷ Proposed a differentiable parallel configuration search space inspired by DARTS, can potentially reach optimal configuration faster than the original dynamic programming.
- ▷ Employed parallel random initialization using sampling algorithms like Bayesian Optimization for faster train loss convergence.

Multi-modal target detection with zero-shot depth estimation, Multi-modal NAS

May 2022 - Sep 2022

Shanghai AI Lab | Intelligent Photoelectric Department | Multi-modal Cognitive Computing Algorithm Intern.

- ▷ Literature review on nature/science papers on multi-modal learning, and conference papers on Neural Architecture Search, multi-modal learning, image and sound feature fusion.
- ▷ The model inputs image and 2 audio channels, and outputs the behavior category and distance. Improve the estimation of distance using zero-shot monocular depth estimation like MiDaS added to the model structure.

Methylation Multi-task Learning for Autoimmune Disease Diagnosis Using Biological Graph Informed Networks | Research Assistant

Sep 2021 - Dec 2023

Advisor: Hongyi Xin, Associate Professor of UM-SJTU Joint Institute, Shanghai Jiao Tong University.

- ▷ Explored adaptable and interpretable neural network to find common genotype given 480k dimension sites, hundreds of sample.
- ▷ Designed an explainable site-gene-pathway ontology constraint to NN to discover new biomarkers by checking weights.
- ▷ Implemented a Variational Auto-Encoder to support gene-level embedding shared among datasets to obtain multi-task learning.
- ▷ Optimized a pretrain-finetune training scheme to increase accuracy by over 10%, wrote the paper under review.

Augmented reality simulation of cardiovascular interventional surgery | Research Assistant

Mar 2020 - May 2021

Advisor: Lixu Gu, Professor of Biomedical Engineering, Shanghai Jiao Tong University.

- ▷ Developed the framework of an augmented reality surgery training assistant system for medical student and surgery.
- ▷ Predicted the operation trajectory using LSTM and used KD-Tree to calculate the distance for operation safety warning.
- ▷ Displayed vascular model in AR with OpenGL and designed the UI interface to support translation.
- ▷ Used the aruco library in OpenCV to coordinate positioning of the QR code.
- ▷ Published Vascular Intervention Training System Based on Electromagnetic Tracking Technology on ICVRV as second author.

Normal fundus image generation based on Generative Adversarial Network | Research Assistant

Apr 2019 - Oct 2019

Advisor: Lisheng Wang, Professor at SEIEE, Shanghai Jiao Tong University.

- ▷ Utilized GAN to design an image generation system to generate normal fundus image for comparisons to identify lesion areas.
- ▷ Used datasets including MESSIDOR-2 for data preprocessing like normalization, contrast enhancement, de-noising, clipping.
- ▷ Implemented a U-Net-based segmentation network to identify and remove fundus blood vessels as part of de-noising process.
- ▷ Built a CycleGAN framework based on segmentation to generate the normal fundus image to identify the lesion area.

PROFESSIONAL EXPERIENCE

XtalPi Inc (QuantumPharm Inc) | XAB-Antibody IDD | Research Scientist Intern

Aug 2024 - July 2025

- ▷ Designed and implemented AI models for antibody binding affinity prediction. Reproduced training pipelines for binding-ddg-predictor and GearBind, developed data visualizations for the SKEMPI dataset, and established robust dataset splits to prevent leakage by ensuring complexes remain within the same fold. Conducted comparative analyses of focal loss (mutation-based), weighted MSE/BCE, balanced sampling, and differentiable ranking strategies.
- ▷ Developed a balanced sampling strategy for binding-ddg-predictor to ensure equal representation of mutations with and without alanine. Achieved performance improvements of approximately +10% in Pearson correlation and +5% in Spearman correlation.
- ▷ Integrated Carbondesign encoder into the binding-ddg-predictor framework to enhance encoder feature representations. Augmented the SKEMPIv2 dataset by generating multi-mutation variants, increasing sample size by ~10x. Successfully reproduced GearBind and conducted benchmarking experiments on the augmented SKEMPIv2 dataset, demonstrating improved model generalization.

Amazon Web Services | VMware Cloud on AWS (Brio) | Software Engineer Intern | Seattle, Washington

May 2023 - Aug 2023

- ▷ Designed UI to automate the workflow to update Rate Card for Brio resource console using React, TypeScript and JavaScript.
- ▷ Integrated UI with backend API using Java, API Gateway and Amazon Lambda. Supported validation and creating and updating rate cards through csv. Preview visualizes modifications and identifies pricing errors. (5k+ lines of code in total).
- ▷ Designed anomaly detection algorithms for disbursement generator to identify usage spikes or subscription anomalies using Java, SQL, Amazon Athena, S3 Buckets, DynamoDB Table. Utilized CloudWatch to create tickets and alarms to engineers.

4Paradigm Co., Ltd | OpenMLDB | GitLink Code Camp 2022 | Open-Source Developer | Shanghai

July 2022 - Oct 2022

- ▷ Developed an automated feature engineering pipeline, including feature generation and selection, based on AutoX and OpenMLDB sql with Python (600+ lines of code). See pull request: <https://github.com/4paradigm/OpenMLDB/pull/2381>.
- ▷ Saving data scientists in dozens of companies' weeks of time by automatic feature selection pipeline.
- ▷ The data is transformed by OpenMLDB sql to get time series and statistics features, and then we select the most important K features based on algorithms like Adversarial Validation, GRN feature selection or Reinforcement Learning.

- ▷ Implemented new features to [Intel® Neural Compressor](#), Cross Layer Equalization, a data free quantization to rescale different layers' weight range to reduce over 5% in the drop in accuracy after quantization like FP32 to INT8, and edited documents.
- ▷ Studied open-source software like Nvidia Triton Inference Server (around 70k lines of code) and AI Model Efficiency Tool and gave presentations to around 100 colleagues in whole department to introduce the design and technique detail of them.
- ▷ Cooperated to design and implement class of AI inference server software by C++, which enables the team to deploy trained AI models from multiple frameworks and deploy more models on GPU or CPU based infrastructure to simplify AI inferencing.

ShuKun Technology Co., Ltd | R & D Department | Algorithm Intern | Shanghai

Dec 2020 - Apr 2021

- ▷ Implemented multi-node and multi-GPU training with horovod framework, NVIDIA clara train sdk, OpenMPI, and NCCL2; 4GPU (2 nodes each 2GPU) can reach 2.5x speed up compared to 1GPU (Discovered that communication is bottleneck).
- ▷ Compared the efficiencies of models like 3D-UNet when multi-node training with different GPU configurations in Python and wrote the training process documents, saving algorithm scientists days in each training job (~0.65n times faster with n nodes).
- ▷ Used Java to add multi-machine and multi-GPU training functions to the company's back-end web page for algorithm scientists.

TEACHING EXPERIENCE

VE370 Intro to Computer Organizations | Professor: Gang Zheng | Teaching Assistant

Sep 2021 - Dec 2021

COURSE PROJECTS

CS 425 Distributed Systems: Distributed ML Inference System (C++ and Python) (100/100)

Aug 2022 - Dec 2022

- ▷ Implemented scheduling algorithm to obtain fair-time inference, making each ML job query rate within 20% difference.
- ▷ Design a distributed file system to maintain datasets with SWIM-like membership protocol (nearly 1e-2 false positive rate).
- ▷ Supported fault-tolerance and can recover when N members and coordinator leave or fail. Leader re-election is within 10s.

CS 543 Computer Vision (A): (Python)

Aug 2022 - Dec 2022

- ▷ Channel Alignment: Designed multi-scale algorithm to align 3 channels of images, evaluating with NCC and Fourier Transform.
- ▷ Laplacian Blob Detection: Implemented a scale-space blob detection with a Laplacian scale space using Laplacian of Gaussian filter and performed non-maximum suppression in scale space.
- ▷ Stitching: Stitching pairs of images using OpenCV, SIFT descriptors and RANSAC to estimate homography mapping of images.
- ▷ 3D Shape from shading: Estimate the albedo and surface normal and compute the surface height map by integration to reconstruct 3D surface from image of shading.

CS 438 Communication Network (A+): TCP Protocol (C++)

Aug 2022 - Dec 2022

- ▷ Implemented transport protocol with properties equivalent to TCP based on unreliable UDP, which can tolerate packet drops, allow other concurrent connections fair chance (0.5x to 2x of TCP), doesn't give up entire bandwidth to other connections.

Sonar detection water pipe state system based on machine vision and neural network compression

Sep 2021 - Dec 2021

ECE4710 Advanced Embedded Systems (MDE, Thesis). | Instructor: Zou An, Assistant Professor of UM-SJTU Joint Institute.

- ▷ Processed the source data of sonar as a matrix and use Python to generate polar images in real time.
- ▷ Used OpenCV to perform subtraction and corrosion operations to complete hole or bulge detection.
- ▷ Designed a compressed and binarized YOLO architecture to realize defect area detection with accuracy 99%.
- ▷ Used HLS to deploy to FPGA, matrix parallel operation, real-time embedded system is implemented.

HONORS & AWARDS

2022 Shanghai Jiao Tong University Outstanding Graduate (School Level) – Aug 2022

2021 Microsoft Imagine Cup Global Competition - Third Prize in China – Jan 2021

2020 Mathematical Contest in Modeling - Meritorious Winner (Top 6%) [[see pdf report](#)] – Apr 2020

2020 "Jidong Cup" CCVR China Virtual Reality Competition, Product Creative Group - Second Prize – Nov 2020

2018-2019 and 2020-2021 Academic Year Undergraduate Excellence Scholarship – Nov 2019/2021

SKILLS

- ▷ Programming: C/C++ (Proficient), Python (Proficient), MATLAB, R, Verilog, Java, RISC-V Assembly, SQL, Shell, TypeScript, CUDA Programming.
- ▷ Frameworks & Libraries: PyTorch, TensorFlow, horovod, Keras, Sk-learn, Pandas, NumPy, OpenCV, Matplotlib, Selenium.
- ▷ Developer Tools: Docker, Git, LaTeX, PSpice, Linux.