

---

 Research Interest Bioinformatics • Machine Learning • Artificial Intelligence
 

---

Skills	<ul style="list-style-type: none"> <li>• Intern experience on <b>recommender systems</b>, C++ backend, <b>Hive/Presto</b>, and <b>Python</b> script</li> <li>• Familiar with machine learning tools (<b>scikit-learn</b>, <b>lightGBM</b>, <b>PyTorch</b>, <b>TensorFlow</b>, <b>Numpy</b>) and achieved top rankings in Kaggle competitions with <b>natural language processing</b>, and <b>computer vision</b> topic.</li> <li>• Good class performance on <b>algorithm design</b>, <b>machine learning</b>, <b>data mining</b>, and <b>software engineering</b></li> <li>• Design algorithms to analyze big biological data (<b>Hidden Markov Model</b>, <b>sequence comparison</b>)</li> </ul>
--------	--

Education	<b>Ph. D., Computer Science</b> <span style="float: right;"><b>May 2019 (Expected)</b></span> Michigan State University, Michigan, USA <ul style="list-style-type: none"> <li>• GPA: 4.0/4.0</li> </ul>
	<b>M. Sc., Condensed Matter Physics</b> <span style="float: right;"><b>May 2015</b></span> Michigan State University, Michigan, USA <ul style="list-style-type: none"> <li>• GPA: 3.6/4.0</li> </ul>
	<b>B. Sc., Physics</b> <span style="float: right;"><b>June 2011</b></span> Fudan University, Shanghai, China

Experience	<b>Software Engineer Ph.D. Intern, Machine Learning</b> <span style="float: right;"><b>June 2018 - present</b></span> Facebook, California, USA <ul style="list-style-type: none"> <li>• Improving Facebook Marketplace ranking by introducing new NLP features</li> <li>• Implementing data pipeline for offline test in Hive and Presto</li> <li>• Implementing feature fetcher in PHP and C++</li> <li>• Conducting offline experiment and online AB test</li> </ul>
	<b>Research Assistant, Bioinformatics Lab</b> <span style="float: right;"><b>Jan 2015 – June 2018</b></span> Michigan State University, Michigan, USA <ul style="list-style-type: none"> <li>• Working on developing the new algorithm for the challenge of the long erroneous sequence.</li> <li>• Designed an algorithm use Hidden Markov Model combined with one popular error correction method based on the directed acyclic graph to infer the possible errors in the reads and correct it. After correction, the alignment length of protein profile can be improved 40% on low coverage dataset.</li> <li>• Designed an algorithm that uses statistic method to determine the threshold to filter the raw data and chaining the data points to make the prediction of the similarity of two long error reads. We expect this method can boost 5% of the sensitivity without losing the accuracy.</li> </ul>

Projects	<b>Silver Medal, Avito Demand Prediction Challenge</b> <span style="float: right;"><b>April 2018 – June 2018</b></span> Kaggle Competition <ul style="list-style-type: none"> <li>• Final ranking <b>68 of 1917</b> with another teammate</li> <li>• Based on user purchase history and the product information (image and post) to predict the CTR</li> <li>• Extract features from customer activities, image embedding and text features and embeddings</li> <li>• Combine gradient boosting method using lightGBM and neural network model implemented by Keras</li> </ul>
	<b>Silver Medal, Toxic Comment Classification Challenge</b> <span style="float: right;"><b>Jan 2018 – March 2018</b></span> Kaggle Competition <ul style="list-style-type: none"> <li>• Final ranking <b>69 of 4551</b> with two team mates</li> <li>• Classified online comment to 7 toxic comment class</li> <li>• Use various neural network model based on LSTM, GRU, and Text CNN</li> </ul>
	<b>Silver Medal, Google Cloud &amp; YouTube-8M Video Understanding Challenge</b> <span style="float: right;"><b>March 2017 – June 2017</b></span> Kaggle Competition <ul style="list-style-type: none"> <li>• Final ranking <b>35 of 650</b> with team RandomForest</li> <li>• Predict tag of video (sequence of processed feature matrix) from YouTube (video classification)</li> <li>• Data processing on large YouTube video dataset (1.7TB)</li> <li>• Ensemble based on variations of LSTM and MoE models</li> </ul>

Publication	<ul style="list-style-type: none"> <li>• Du N., Chen J., and Sun Y., Improving the sensitivity of detecting long read overlaps using grouped short k-mer matches, submitted to <i>APBC 2019</i></li> <li>• Pak, D., Du, N., Kim, Y., Sun, Y., &amp; Burton, Z. F. (2018). Rooted tRNAomes and evolution of the genetic code. <i>Transcription</i>, 9(3), 137-151.</li> <li>• Du, Nan, and Sun, Yanni. "Improve homology search sensitivity of PacBio data by correcting frameshifts." <i>Bioinformatics</i> 32.17 (2016): i529-i537.</li> <li>• Improve homology search sensitivity of Pacbio data by correcting frameshifts (Proceeding Talk), <i>15th European Conference on Computational Biology (ECCB 2016)</i>, The Hague, Netherlands</li> </ul>
-------------	---

---