

# Majid Daliri

370 Jay St, Brooklyn, NY  
11<sup>th</sup> Floor, Desk A1119

Cell: +1-646-750-4667  
[daliri.majid@nyu.edu](mailto:daliri.majid@nyu.edu)  
[majid-daliri.github.io](https://github.com/majid-daliri)

Education	New York University, New York, USA Ph.D. in Computer Science (in progress) Student in the <a href="#">Theoretical Computer Science</a> group Advised by <a href="#">Prof. Christopher Musco</a>	2022 - present
	New York University, New York, USA M.S. in Computer Science	2022 - 2024
	University of Tehran, Tehran, Iran B.S. in Computer Engineering (Cumulative GPA: <b>3.97/4.0</b> )	2017 - 2022
Publications	<div><div>(1) <a href="#">Coupling without Communication and Drafter-Invariant Speculative Decoding</a> (Arxiv) Majid Daliri, Christopher Musco, Ananda Theertha Suresh</div><div>(2) <a href="#">Unlocking the Theory Behind Scaling 1-Bit Neural Networks</a> (CPAL) 2025 Majid Daliri, Zhao Song, Chiwun Yang</div><div>(3) <a href="#">QJL: 1-Bit Quantized JL transform for KV Cache Quantization</a> (AAAI) 2025 Amir Zandieh, Majid Daliri, Insu Han</div><div>(4) <a href="#">Matrix Product Sketching via Coordinated Sampling</a> (ICLR) 2025 Majid Daliri, Juliana Freire, Danrong Li, Christopher Musco</div><div>(5) <a href="#">Sampling Methods for Inner Product Sketching</a> (VLDB) 2024 Majid Daliri, Juliana Freire, Christopher Musco, Aécio Santos, Haoxiang Zhang</div><div>(6) <a href="#">Simple Analysis of Priority Sampling</a> (SOSA) 2024 Majid Daliri, Juliana Freire, Christopher Musco, Aécio Santos, Haoxiang Zhang</div><div>(7) <a href="#">KDEformer: Accelerating Transformers via Kernel Density Estimation</a> (ICML) 2023 Amir Zandieh, Insu Han*, Majid Daliri*, Amin Karbasi (* equal contribution)</div><div>(8) <a href="#">Weighted Minwise Hashing Beats Linear Sketching for Inner Product Estimation</a> (PODS) 2023 Aline Bessa, Majid Daliri, Juliana Freire, Cameron Musco, Christopher Musco, Aécio Santos, Haoxiang Zhang</div><div>(9) <a href="#">Efficient Approximations for Cache-conscious Data Placement</a> (PLDI) 2022 Ali Ahmadi, Majid Daliri, Amir Kafshdar Goharshady, Andreas Pavlogiannis</div><div>(10) <a href="#">A 10-Approximation of the <math>\frac{\pi}{2}</math>-MST</a> (STACS) 2022 Ahmad Biniaz, Majid Daliri, AmirHossein Moradpour</div></div>	
Research Internship	Machine Learning Intern, Max Planck Institute, Advised by Dr. Amir Zandieh	2021-2022
	<ul style="list-style-type: none"><li>Implemented Fast Attention in Transformers, optimizing accuracy and efficiency in sequence modeling tasks.</li><li>Designed an efficient GPU-compatible LSH method, boosting performance in attention approximation.</li><li>Technical Stack: BigGAN, PyTorch, Transformer architectures, and advanced sequence modeling tools.</li></ul>	
	Research Intern, HKUST, Advised by <a href="#">Prof. Amir Goharshady</a>	2021-2022
	<ul style="list-style-type: none"><li>Made a pioneering theoretical contribution to Cache-conscious Data Placement (CDP), addressing a longstanding challenge in optimizing cache hits.</li><li>Implemented and tested various cache management policies and algorithms, outperforming previous heuristics.</li><li>Our method emerged as the most effective solution, setting a new standard for cache optimization.</li></ul>	
	Research Intern, University of Salzburg, Advised by <a href="#">Prof. Ana Sokolova</a>	Summer 2022
	<ul style="list-style-type: none"><li>Developed algorithms for Distribution Bisimilarity, focusing on finite bisimulation up to convex hull witness.</li><li>Extended research to probabilistic system verification and collaborated on using Quantum Annealers for program verification.</li></ul>	

Awards and Honors	<b>Travel Grant, ACM-SIAM Symposium on Discrete Algorithms</b>	2024
	Awarded a travel grant to present a paper.	
	<b>Research Grant, University of Salzburg</b>	Summer 2022
	Awarded a €5,000 grant for a research internship focusing on algorithms for distribution bisimilarity, probabilistic systems verification, and quantum annealing projects.	
	<b>Hong Kong PhD Fellowship Scheme (declined to attend NYU)</b>	2022
	Awarded \$185,000 Ph.D. Fellowship as one of the top 300 students selected across all fields for academic excellence and research potential.	
	<b>ACM ICPC - Regional (University of Tehran)</b>	2019
	Ranked 6 <sup>th</sup> among more than 100 team all around the Iran.	
	<b>Iranian National Olympiad in Informatics Finalist (IOI, Iran)</b>	2016
	Awarded to 50 students after a year long competition involving 10,000 students.	
Service	<ul style="list-style-type: none"> <li>• Reviewer for <a href="#">International Conference on Learning Representations (ICLR 2025)</a></li> <li>• Reviewer for <a href="#">Conference on Neural Information Processing Systems (NeurIPS 2024)</a></li> <li>• Reviewer for <a href="#">ACL Rolling Review (ACL, EMNLP, NAACL 2024-2025)</a></li> <li>• Reviewer for <a href="#">International Conference on Machine Learning (ICML 2024-2025)</a></li> <li>• Reviewer for <a href="#">Royal Society Open</a></li> <li>• External Reviewer for <a href="#">Canadian Conference on Computational Geometry (CCCG 2023)</a></li> </ul>	
Conference Presentations	<ul style="list-style-type: none"> <li>• Simple Analysis of Priority Sampling</li> <li>• Accelerating Transformers via Kernel Density Estimation</li> <li>• Weighted MinHash for Inner Product Estimation</li> <li>• Efficient Approximations for Cache-conscious Data Placement</li> </ul>	Presentation, (SOSA) 2024 Poster, (ICML) 2023 Poster, (PODS) 2023 Presentation, (PLDI) 2022
Teaching	<ul style="list-style-type: none"> <li>• Section Leader for <a href="#">NYU CSCI-UA 310 Basic Algorithms</a></li> <li>• Teaching Assistant <a href="#">NYU CS-GY 6763 Algorithmic Machine Learning</a></li> <li>• Teaching Assistant University of Tehran Design and Analysis of Algorithms</li> </ul>	Spring 2023 Fall 2022 Fall 2020-2021
Work Experience	<b>Site Reliability Engineer at <a href="#">Cafebazaar</a></b> <ul style="list-style-type: none"> <li>• Designed, implemented, and maintained both Redis-as-a-Service/PostgreSQL-as-a-Service on a Kubernetes-based cloud.</li> <li>• Achieved consistent performance benchmarks for both services with 100% uptime and a 99.9% response rate.</li> <li>• Technical Stack: Kubernetes, Docker, Sentry, S3, Prometheus</li> </ul>	2021 - 2022
Skills	<b>Theoretical Background:</b> Proficient in Machine Learning Theory, Neural Networks, Linear Algebra, and Probability.  <b>Technical Skills:</b> Highly skilled in C/C++, CUDA, Go, Python, Bash-Scripting, PHP, JavaScript. Experience with PyTorch, TensorFlow, Django, CSS3, HTML5, and git.	