

## JEFF NIU

Email: [jeffniu22@gmail.com](mailto:jeffniu22@gmail.com)  
Mobile: +1 650 864 3841

Github: [github.com/mogball](https://github.com/mogball)  
LinkedIn: [in/jeffniu22](https://in.linkedin.com/in/jeffniu22)

work experience	<b>MODULAR</b> - Lead architect for Mojo, a Pythonic systems programming language for high performance heterogeneous computing, with first class Python and C++ FFI - Gave <a href="#">a bunch of talks</a> about this work  - Built Mojo's MLIR-based implementation from scratch: LLVM code generation, mid-level optimizer (inlining, SROA, mem2reg, arg promotion, AA, memory SSA, etc.), metaprogramming features, compile time interpreter, type system, compilation model - Implemented tons of major language features across the stack: async/await, traits and generics, closures, type families/constraints, parameter inference, metatypes, exception/error handling, split GPU compilation, lifetimes and memory safety  - Prototyped key Modular technologies in Mojo that were handed to other teams: NDBuffer with partially static properties, thread pool and async runtime, structural codegen prototype, "engine extensibility API"  - Lead integration of Mojo and Modular's graph compiler, tuned the mid-level optimizer for autofusion, and drove engineering excellence across both teams - Drove >1000x reductions in compile time over 2 years (Mojo compiles llama3 with hundreds of specialized kernels for SoTA performance in less than 40 sec)  - Technical Lead of the Mojo Compiler Team, driving language design and evolution with 6 other engineers across the frontend and backend, primary design and technical PoC with all other teams at Modular, DRI for many projects (debugger, Python interop, GPU bringup) - Performed various managerial/TPM functions: "therapy" and growth mentorship for engineers on my team, quarterly work planning, prioritization, and allocation, bug triage, etc. - Hosted an intern to build first-class library optimizations, paper targeting PLDI 2025  <b>GOOGLE (CoreML)</b> - Built high-level control-flow and schedule optimizer for TensorFlow, delivering ~2% improvements to Google AdBrain models (\$millions of compute/year) - Implemented optimizations to the TensorFlow compiler, saving ~\$700k/yr - Built <a href="#">MLIR dataflow analysis framework</a> , collaborated with JSIR team to implement security analyses over JavaScript code running in Chrome and other Google environments - Contributed to upstream MLIR core infrastructure, dialects and tooling	2022-now
internships	<b>ETH ZURICH (Scalable Parallel Computing Lab)</b> - Built <a href="#">IRDL</a> , a DSL for defining compilers, published in <a href="#">PLDI 2022</a>  <b>CITADEL</b> - Bringup and evaluation of Cerebras hardware for internal HPC applications  <b>GOOGLE (Brain)</b> - Built <a href="#">PDL</a> , an optimized pattern rewriter, now used by various hardware companies  <b>APPLE (Silicon validation)</b> - GPU assembly programming to stress SoCs as much as possible  <b>COREAVI (Embedded graphics)</b> - Hacking on embedded GPU driver for various Khronos standards, VxWorks 6/7, etc.  <b>YAHOO! (Big data)</b> - Contributed to <a href="#">Apache Superset</a> , built <a href="#">anomaly detector</a> , created <a href="#">ember-localforage</a>	2020  Summer 2020  Fall 2019  Winter 2019  Summer 2018  Fall 2017
projects	<b>gg-mlir</b> : An MLIR-based distributed workload optimizer and executor ( <a href="#">slides</a> ) <b>WATERLOOP</b> : C++ <a href="#">STL</a> and <a href="#">package manager</a> , embedded networking and architecture <b>UWNRG</b> : Computer vision, pathfinding, <a href="#">robotics stuff</a>	
education	UNIVERSITY OF WATERLOO: B.A.Sc. in Mechatronics Engineering (2021)	