

Accelerating Optimal Power Flow: Condensed-Space Interior-Point Methods and Automatic Differentiation on GPUs

Sungho Shin and Mihai Animescu
Mathematics and Computer Science Division
Argonne National Laboratory
Lemont, IL, USA
sshin@anl.gov, animescu@mcs.anl.gov

François Pacaud
Centre Automatique et Systèmes
Mines Paris - PSL
Paris, France
francois.pacaud@minesparis.psl.eu

While graphics processing units (GPUs) have showcased impressive capabilities in various computing domains, their utilization in large-scale constrained nonlinear optimization regimes, such as alternating current (AC) optimal power flow (OPF) problems, has been somewhat limited. This limitation stems from the challenges associated with parallel factorization of indefinite sparse matrices commonly encountered within constrained optimization algorithms (Animescu et. al. 2021). Although GPU computation can accelerate various other parts of the optimization process, including automatic differentiation (AD) and sparse matrix-vector multiplications, the slow data transfer between host and device memory hinders the ad-hoc implementation of GPU accelerations. To fully leverage the capabilities of modern GPU hardware, it is essential to implement a comprehensive computational framework on the GPU that incorporates AD, linear algebra, and optimization while minimizing data transfers to and from host memory.

This paper presents a comprehensive computational framework and the associated software implementations for solving AC OPF problems on GPUs. Our approach utilizes the following techniques: (i) condensed-space interior-point methods (IPMs) with an inequality relaxation strategy, (ii) sparse matrix factorization with a fixed pivot sequence, and (iii) a single-instruction, multiple-data (SIMD) abstraction of nonlinear programs. Specifically, our method relaxes power flow equality constraints by allowing small violations, which enables expressing the Karush-Kuhn-Tucker (KKT) system entirely in the primal space through the condensation procedure. Although this strategy is not new (see Nocedal and Wright, 2006), it has traditionally been considered less efficient than the standard full-space method due to increased fill-in in the sparse factorization. However, when implemented on GPUs, it offers the key advantage of guaranteeing positive definiteness in the condensed KKT system through the application of standard regularization techniques, which in turn, allows for the utilization of linear solvers with a fixed numerical pivot sequence (so-called refactorization). An efficient implemen-

tation of the sparse refactorization is available as part of the CUDA library, facilitating the implementation of efficient KKT system solutions on GPUs. Although this method is susceptible to ill-conditioning, our results demonstrate that the solver is robust enough to solve problems with a relative accuracy of 10^{-6} .

Furthermore, by leveraging the SIMD abstraction of nonlinear programs, which preserves the parallelizable structure in the model, the model functions and derivative evaluations can be parallelized, thereby facilitating evaluations on the GPU. We demonstrate that the AC power flow model is particularly well-suited for this abstraction as it involves repetitive expressions for each type of component in the model (e.g., buses, lines, generators), and the number of computational patterns does not increase with the network's size. This structure has been effectively utilized by Gravity (Hijazi et. al. 2018), demonstrating significant acceleration in model evaluations.

The paper will present comprehensive numerical benchmark results to demonstrate the efficiency of our proposed approach. Our method is implemented in our packages, a nonlinear optimization solver MadNLP.jl and an automatic differentiation tool SIMDDiff.jl, with the solution of the KKT system being carried out using the external CUSOLVER library. We compare our method against the standard CPU approach (Wächter and Beigler, 2006) and the recently developed reduced-space interior-point method on GPUs (Pacaud et. al., 2022) using the standard data available in pglib-opf. Preliminary results indicate that our proposed framework has the potential to accelerate the solution of AC OPF problems up to moderate tolerances (10^{-6}) by an order of magnitude compared to existing tools such as Ipopt (interfaced with MATPOWER or PowerModels.jl), especially for large-scale instances. We will conclude with a discussion on the future extensibility of our method for more complex optimization tasks, such as multi-period, security-constrained, and joint transmission-distribution optimization.