Fast Machine Learning for Science Imperial College London Real-time and accelerated ML 25-28 September 2023

Real-time and accelerated ML for fundamental sciences

Contribution ID: 63

Type: Standard Talk

Real-time Fitting and Materials Characterization in Band- Excitation Piezoresponse Force Microscopy

Tuesday 26 September 2023 16:15 (15 minutes)

Increased development and utilization of multimodal scanning probe microscopy (SPM) and spectroscopy techniques have led to an orders-of-magnitude increase in the volume, velocity, and variety of collected data. While larger datasets have certain advantages, practical challenges arise from their increased complexity including the extraction and analysis of actionable scientific information. In recent years, there has been an increase in the application of machine and deep learning techniques that use batching and stochastic methods to regularize statistical models to execute functions or aid in scientific discovery and interpretation. While this powerful method has been applied in a variety of imaging systems (e.g., SPM, electron microscopy, etc.), analysis alone takes on the order of weeks to months due to scheduling and IO overhead imposed by GPU and CPU based systems which limits streaming inference rates to speeds above 50ms. This latency precludes the possibility of real-time analysis in SPM techniques such as band-excitation piezoresponse force spectroscopy (BE PFM), where typical measurements of cantilever resonance occur at 64Hz.

One method to accelerate machine learning inference is to bring computational resources as close to the data acquisition source as possible to minimize latencies associated with I/O and scheduling. Therefore, we leverage the National Instruments PXI platform to establish a direct, peer-to-peer channel over PCIe between an analog-to-digital converter and a Xilinx field programmable gate array (FPGA). Through the LabVIEW FPGA design suite, we develop this FPGA-based pipeline using cantilever resonances acquired in BE PFM to conduct real-time prediction of the simple harmonic oscillator (SHO) fit. To accomplish this, we use hls4ml to compile a high-level synthesis (HLS) representation of the neural network. Once this HLS model is synthesized to a register transfer level description (RTL), we implement the design on the FPGAs programmable logic. The parallelizable nature of FPGAs allows for heavily pipelined neural network implementations to achieve latencies on the order of microseconds. We currently benchmark our implementation at 36 us per inference with a fourier transformation accounting for an additional 330 us. At the expense of FPGA resources, we overlap data acquisition with computation to enable continuous acquisition and processing of response data. This work provides a foundation for deploying on-sensor neural networks using specialty hardware for real- time analysis and control of materials imaging systems. To further enhance the performance and capabilities we discuss our progress implementing this system on an RFSoC 4x2 (Radio Frequency System-on-Chip) which integrates both the FPGA and RF data converters on a single chip, effectively combining analog and digital processing capabilities and reducing latencies associated with I/O and scheduling.

Author: OBUTE, Veronica (Drexel)

Presenter: OBUTE, Veronica (Drexel)

Session Classification: Contributed Talks

Track Classification: Contributed Talks