

---

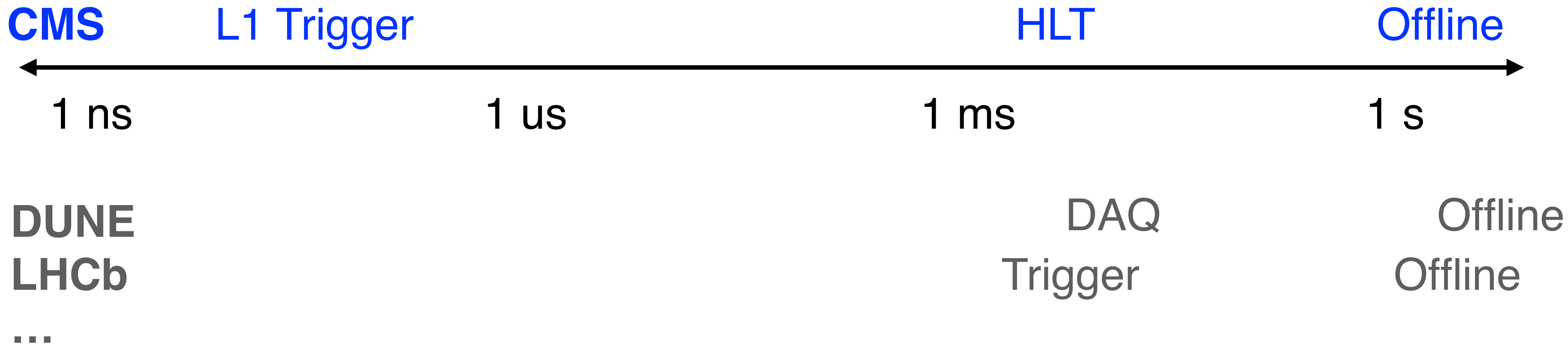
# Machine learning in FPGAs

Javier Duarte, Burt Holzmann, Sergo Jindariani, Benjamin Kreis,  
Kevin Pedro, Ryan Rivera, Nhan Tran

**Fermilab**

+ collaborators from UIC, MIT, CERN and industry

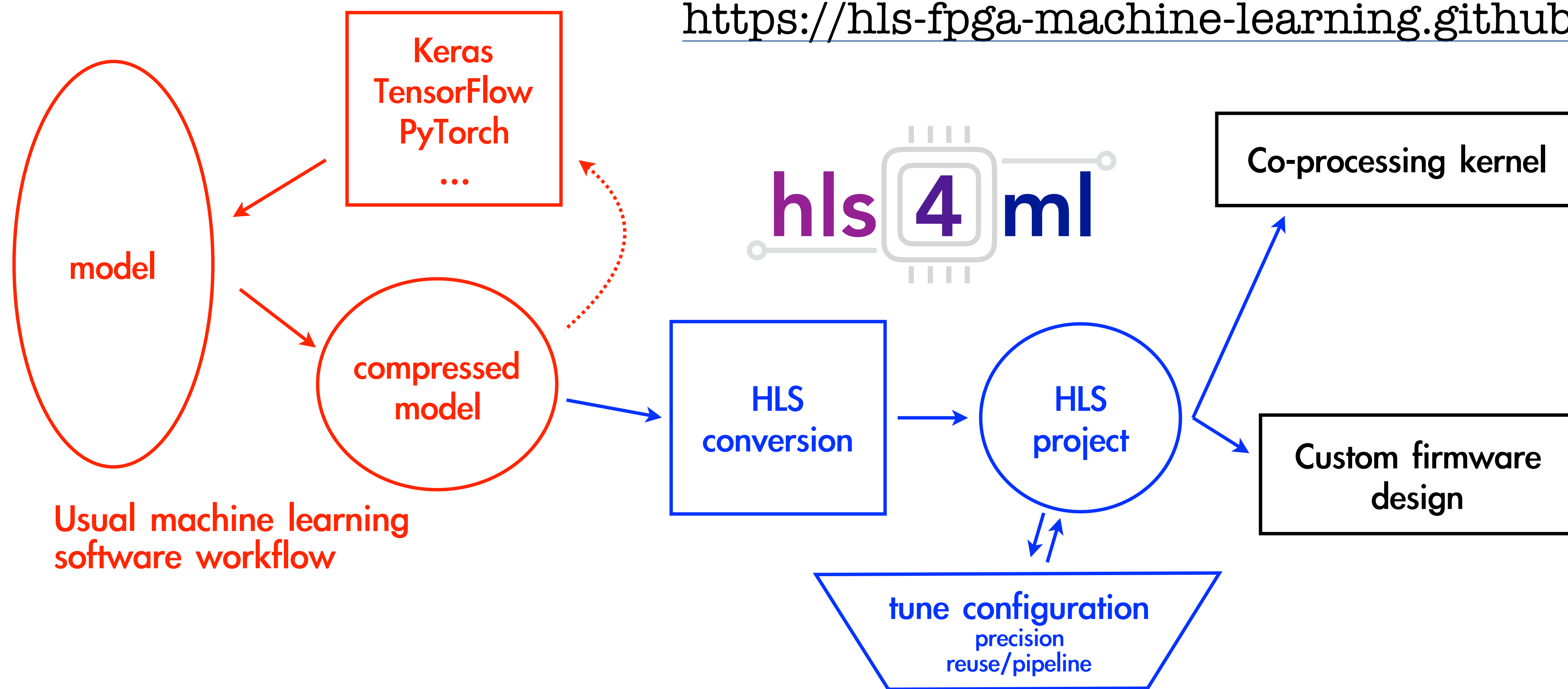
# LATENCY LANDSCAPE



Machine learning is being used to solve a wide array of problems across a large range of latency constraints

## *Assumptions*

- most algorithms can be formulated as machine learning problems**
- new specialized hardware will be optimized for machine learning**



## Fast inference of deep neural networks in FPGAs for particle physics

Javier Duarte<sup>a</sup>, Song Han<sup>b</sup>, Philip Harris<sup>b</sup>, Sergo Jindariani<sup>a</sup>, Edward Kreinar<sup>c</sup>, Benjamin Kreis<sup>a</sup>, Jennifer Ngadiuba<sup>d</sup>, Maurizio Pierini<sup>d</sup>, Ryan Rivera<sup>a</sup>, Nhan Tran<sup>a</sup>, Zhenbin Wu<sup>e</sup>

<sup>a</sup>Fermi National Accelerator Laboratory, Batavia, IL 60510, USA

<sup>b</sup>Massachusetts Institute of Technology, Cambridge, MA 02139, USA

<sup>c</sup>HawkEye360, Herndon, VA 20170, USA

<sup>d</sup>CERN, CH-1211 Geneva 23, Switzerland

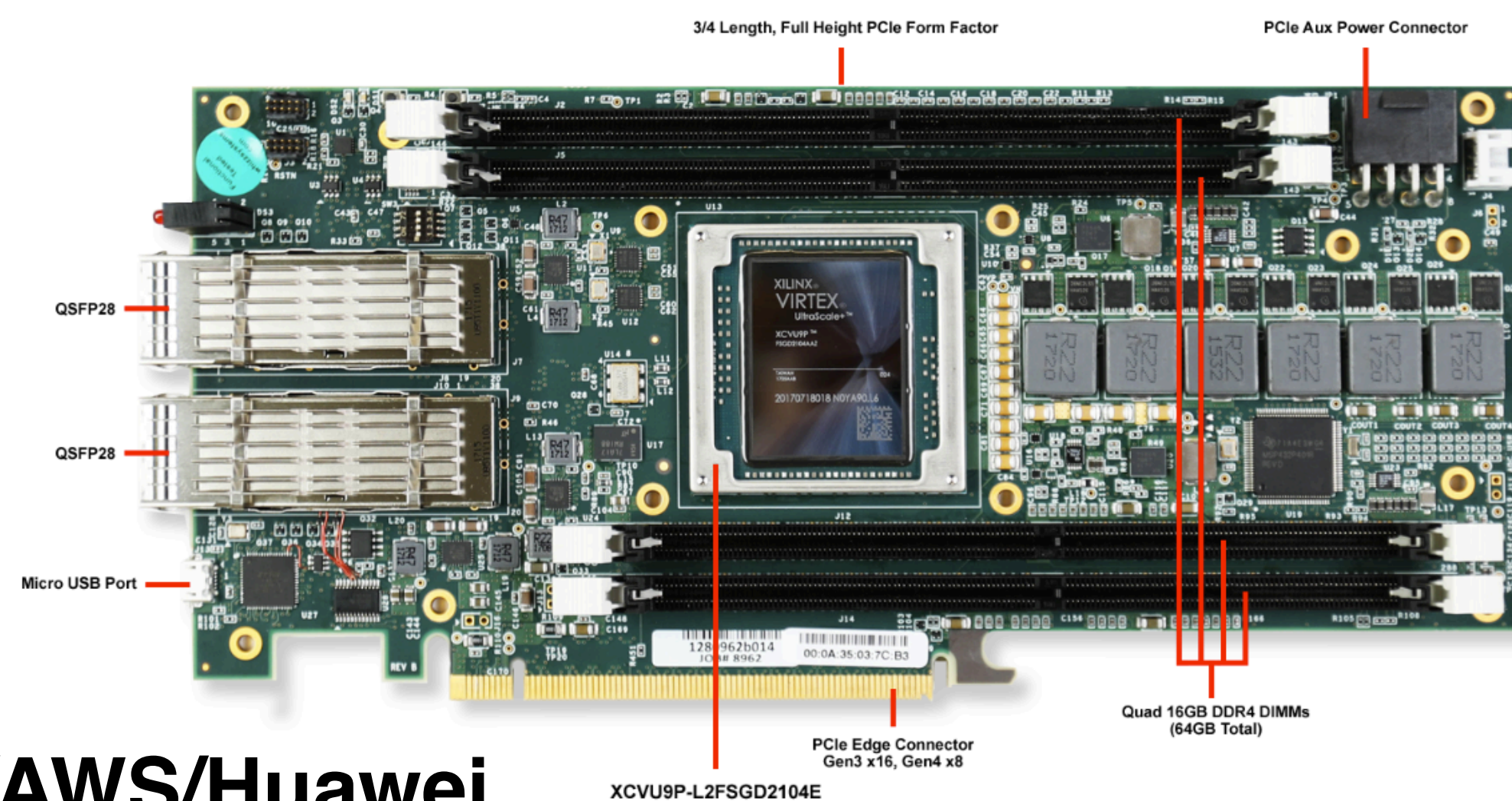
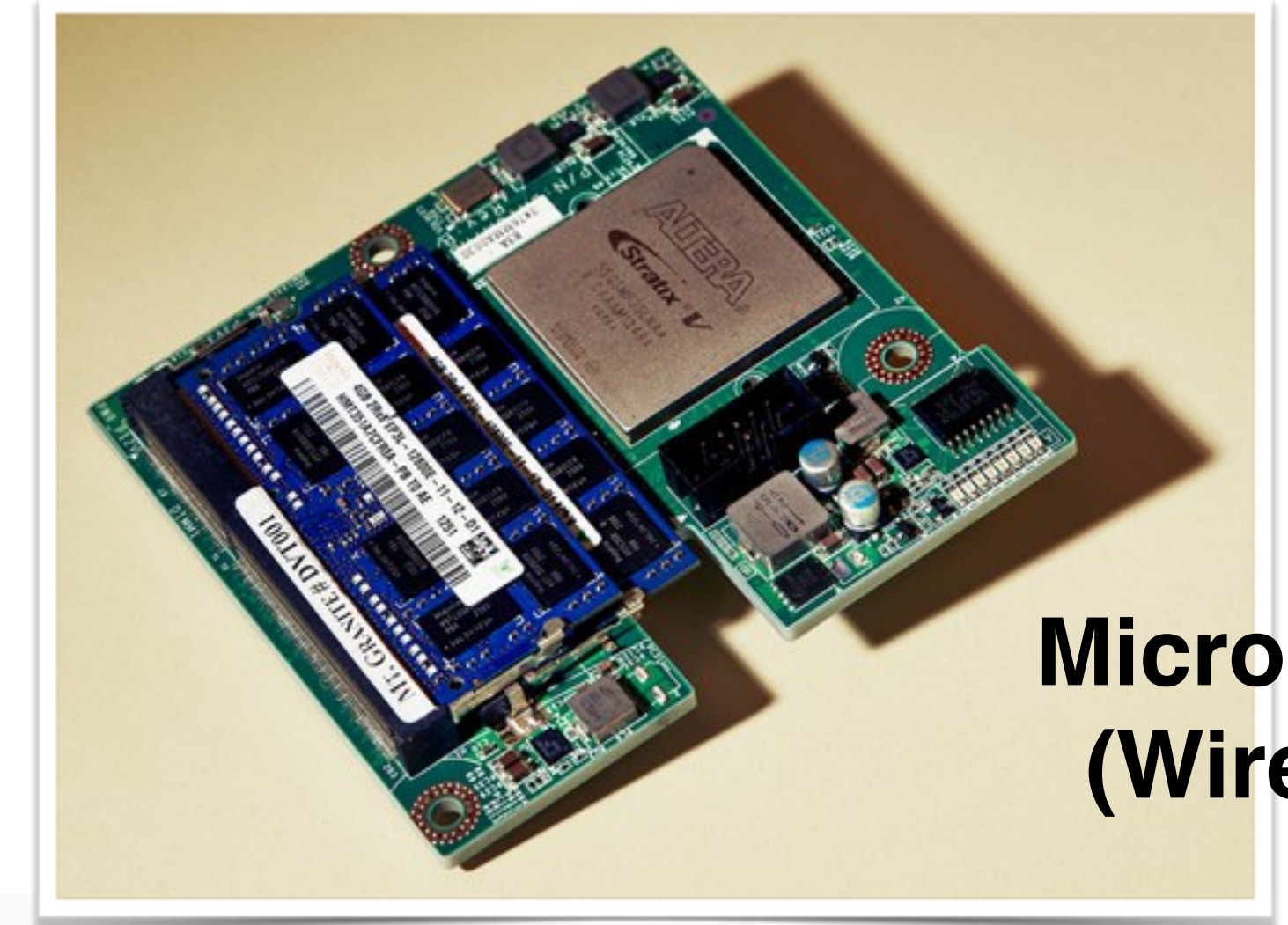
<sup>e</sup>University of Illinois at Chicago, Chicago, IL 60607, USA

<https://arxiv.org/abs/1804.06913>

SEE RESEARCH TECHNIQUES SEMINAR  
BY JAVIER DUARTE,  
APRIL 24 @ 1030 AM



**LARGE SPEED/ENERGY GAINS  
OVER CPU/GPU!**



**Xilinx/AWS/Huawei**



Efficiently interface between CMSSW and accelerator hardware?

Some pilot projects in progress with: Microsoft Azure + MSR , Amazon Web Services

Exploring connections through CERN OpenLab with: Intel

Academia?

***These systems will only improve (higher throughput, more DSPs/memory)***

Goals, benchmark performance of co-processor hardware against other computing architectures (CPU. GPU. etc.)

**SEE MACHINE LEARNING FORUM TALK, ANDREW PUTNAM (MICROSOFT RESEARCH), MAY 14 @ 1PM**