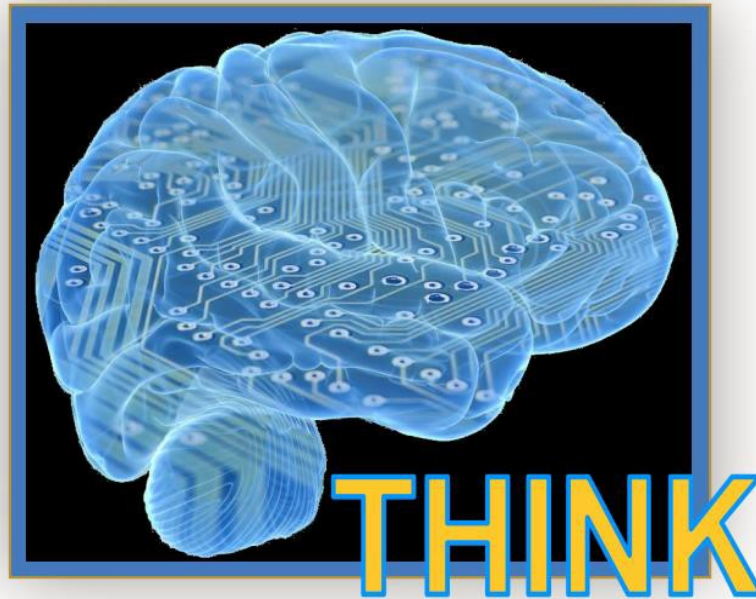


Projet transverse THINK

Testing Hardware Instantiations of Neural Kernels



J.-P. Cachemiche
CPPM

Motivation

Systèmes neuronaux et deep learning

- Montrent leur efficacité dans de nombreux domaines, principalement pour toute forme de reconnaissance
- Technique ancienne mais qui monte en puissance sous deux effets :
 - Disponibilité de larges bases de données
 - Augmentation de puissance des processeurs

Application dans le domaine de la physique

- Augmentation des luminosités et du bruit de fond rends les détections par algorithme classique plus difficiles
- Infléchissement de la loi de Moore
 - ➔ Besoins d'accélérateurs
- Premiers résultats très intéressants dans CMS pour l'extraction de Jets

Besoins de formation exprimés lors de l'ANR DAQ Emergents sur techniques sous-jacentes

- Calcul GPU
- Architectures MPPA type Kalray
- Langages de haut niveau pour FPGAs

But du projet

Augmenter le niveau de connaissance des ingénieurs et techniciens sur les techniques d'apprentissage et les techniques neuronales

- Organisation de webinaires avec les meilleurs spécialistes

Identifier quelques applications typiques

Définir une architecture neuronale appropriée

Effectuer une phase d'apprentissage

Effectuer un portage sur différents candidats matériels

- FPGA
- MPPA Kalray
- GPU
- Processeur neuromorphique (Movidius Intel, BrainChip d'Akira, ...)

Comparer les performances

Élaborer des méthodes, des blocs réutilisables

Diffuser les résultats à la communauté sous forme de workshops

Applications potentielles

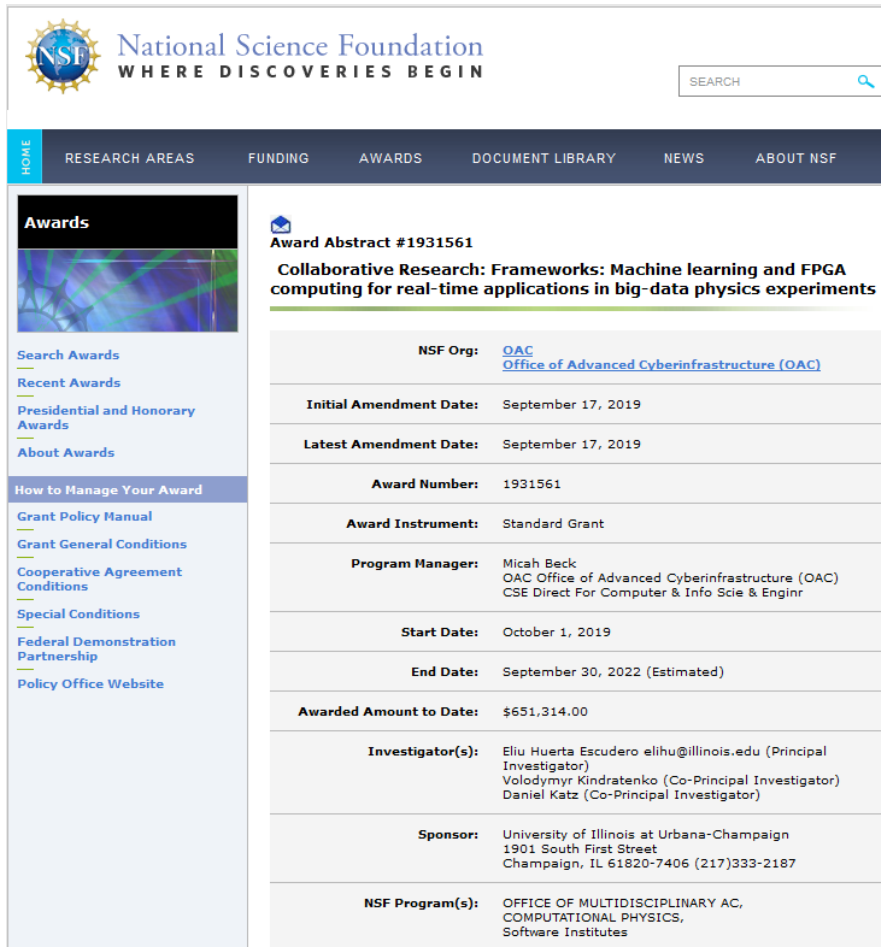
Le projet **Amidex OWEN** (Optimal Waveform recognition Electronic Node) qui consiste à développer un nouvel instrument pour traiter le signal venant d'un détecteur innovant, une TPC sphérique à haute pression. Son but est la recherche d'un phénomène rare tel que la détection directe de matière noire et l'observation de la décroissance double bêta sans neutrino. Dans ce contexte, il s'agit de développer un système d'acquisition intégrant un algorithme de problème inverse basé sur les réseaux de neurones pour l'**identification des formes d'ondes**

Le projet **RTA** (Real-Time Analysis) dans l'expérience LHCb qui consiste à traiter 40 Tb des données par seconde pour n'en garder que 80Gb/s pour une analyse plus profonde offline. Pour ce faire RTA doit à la fois utiliser efficacement les architectures modernes de calcul, et mettre en place des algorithmes avancés tels que les réseaux neurones.

Le projet **Amidex AIDAQ** qui consiste à implémenter des algorithmes de reconnaissance neuronale sur FPGA dans le calorimètre à argon liquide d'ATLAS pour réaliser les **fonctions de trigger de premier niveau** en environnement fortement bruité et avec des niveaux de pile-up variables.

Les projets d'imagerie médicale et en particulier ceux articulés autour des **tomographes** où les problèmes de reconnaissance sont cruciaux.

Projets similaires et premiers travaux



The screenshot shows the NSF website interface. At the top is the NSF logo and the tagline 'WHERE DISCOVERIES BEGIN'. Below is a navigation bar with links for HOME, RESEARCH AREAS, FUNDING, AWARDS, DOCUMENT LIBRARY, NEWS, and ABOUT NSF. A search box is located in the top right. The main content area displays an 'Award Abstract #1931561' for a collaborative research project. The project title is 'Collaborative Research: Frameworks: Machine learning and FPGA computing for real-time applications in big-data physics experiments'. The abstract details include the NSF Org (OAC Office of Advanced Cyberinfrastructure), initial and latest amendment dates (September 17, 2019), award number (1931561), award instrument (Standard Grant), program manager (Micha Beck), start and end dates (October 1, 2019 to September 30, 2022), awarded amount (\$651,314.00), investigator(s) (Eliu Huerta Escudero, Volodymyr Kindratenko, Daniel Katz), sponsor (University of Illinois at Urbana-Champaign), and NSF program(s) (OFFICE OF MULTIDISCIPLINARY AC, COMPUTATIONAL PHYSICS, Software Institutes).

NSF Org: [OAC Office of Advanced Cyberinfrastructure \(OAC\)](#)

Initial Amendment Date: September 17, 2019

Latest Amendment Date: September 17, 2019

Award Number: 1931561

Award Instrument: Standard Grant

Program Manager: Micha Beck
OAC Office of Advanced Cyberinfrastructure (OAC)
CSE Direct For Computer & Info Scie & Enginr

Start Date: October 1, 2019

End Date: September 30, 2022 (Estimated)

Awarded Amount to Date: \$651,314.00

Investigator(s): Eliu Huerta Escudero elihu@illinois.edu (Principal Investigator)
Volodymyr Kindratenko (Co-Principal Investigator)
Daniel Katz (Co-Principal Investigator)

Sponsor: University of Illinois at Urbana-Champaign
1901 South First Street
Champaign, IL 61820-7406 (217)333-2187

NSF Program(s): OFFICE OF MULTIDISCIPLINARY AC,
COMPUTATIONAL PHYSICS,
Software Institutes

Fast inference of deep neural networks in FPGAs for particle physics

Javier Duarte, Song Han, Philip Harris, Sergio Jindariani, Edward Kreinar, Benjamin Kreis, Jennifer Ngadiuba, Maurizio Pierini, Ryan Rivera, Nhan Tran, Zhenbin Wu

(Submitted on 16 Apr 2018 (v1), last revised 28 Jun 2018 (this version, v3))

Recent results at the Large Hadron Collider (LHC) have pointed to enhanced physics capabilities through the improvement of the real-time event processing techniques. Machine learning methods are ubiquitous and have proven to be very powerful in LHC physics, and particle physics as a whole. However, exploration of the use of such techniques in low-latency, low-power FPGA hardware has only just begun. FPGA-based trigger and data acquisition (DAQ) systems have extremely low, sub-microsecond latency requirements that are unique to particle physics. We present a case study for neural network inference in FPGAs focusing on a classifier for jet substructure which would enable, among many other physics scenarios, searches for new dark sector particles and novel measurements of the Higgs boson. While we focus on a specific example, the lessons are far-reaching. We develop a package based on High-Level Synthesis (HLS) called `hls4ml` to build machine learning models in FPGAs. The use of HLS increases accessibility across a broad user community and allows for a drastic decrease in firmware development time. We map out FPGA resource usage and latency versus neural network hyperparameters to identify the problems in particle physics that would benefit from performing neural network inference with FPGAs. For our example jet substructure model, we fit well within the available resources of modern FPGAs with a latency on the scale of 100 ns.

Planning

Projet sur 3 ans

	2020				2021				2022			
	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4
Trainings												
Choix d'applications pertinentes												
Apprentissage												
Implémentation matérielle												
Evaluation												
Mise à disposition outils et blocs réutilisables												

Participants

Physiciens

Nom	Prénom	Laboratoire	Statut	% ETP
Gligorov	Vladimir	LPNHE	DR	5%
Monnier	Emmanuel	CPPM	DR	5%
Aad	George	CPPM	CR	10%
Calvet	Thomas	CPPM	CR	10%
Boursier	Yannick	CPPM	MDC	5%

Ingénieurs

Nom	Prénom	Laboratoire	Statut	% ETP
Cachemiche	Jean-Pierre	CPPM	IR	10%
Le Dortz	Olivier	LPNHE	IR	10%
Druillole	Frédéric	CENBG	IR	15%
Bouet	Raphaël	CENBG	CDD IE	20%
Rebii	Abdel	CENBG	IR	15%
Etasse	David	LPC Caen	IR	10%
Hommet	Jean	LPC Caen	IR	10%
Bellachia	Fatih	LAPP	IE	10%
Lafrasse	Sylvain	LAPP	IE	10%
Frontera-Pons				5%

Responsabilités

- **LPC Caen** : portage sur MPPA, éventuellement sur carte développée par le laboratoire
- **LAPP** : portage sur processeur neuromorphique
- **LPNHE** : Portage sur FPGA et GPU
- **CENGB** : Portage sur FPGA Xilinx
- **IRFU/AIM** : Aspects théoriques et formation
- **CPPM** : coordination du projet, portage sur FPGA Intel et sur GPU

Conclusion

Projet en cours de soumission

→ Excellents feedbacks de la direction de l'IN2P3

Domaine évoluant très vite

→ Surveillance de l'état de l'art

Ouvert à de nouveaux participants tout au long du projet