



HARDWARE-EFFICIENT MACHINE LEARNING

DESIGNING ACROSS THE CIRCUIT-ARCHITECTURE-ALGORITHM STACK



3 - 5 June 2019

ABOUT THE COURSE

Machine Learning is rapidly gaining importance. Due to the advance in computational power more and more applications for machine learning and deep learning systems are becoming reality and are increasingly deployed into embedded, resource-constrained devices. This puts stringent requirements on electronic and integrated system design. In this course we will focus on the hardware-efficient implementation of machine learning, more specifically deep neural networks. It will become clear that a truly efficient design, optimized across the complete circuit / architecture / algorithmic level design space. The course will go in depth on all these aspects.

WHO SHOULD ATTEND

Engineers, IC designers and engineering manager who are interested on the hardware implementation of machine learning and deep learning systems should follow this course. Furthermore, this course is of great interest to people who work on the software implementation of machine learning and artificial intelligence algorithms, and want to understand the implications of algorithmic choices within a complete embedded system.

COURSE OBJECTIVES

In this 3 day program the participant will learn about:

- Deep learning concepts and algorithm-driven efficiency enhancement techniques
- Processor and datapath architectures for neural network execution
- Exploiting mixed-signal processing for machine learning
- Exploiting in-memory computations and emerging memory devices for machine learning
- Cross-layer dataflow optimizations across algorithms - architecture - circuits

REQUIRED BACKGROUND KNOWLEDGE

Some knowledge in CMOS technology, analog and digital IC design and signal processing is required to follow this course..

LECTURING TEAM

Professor Marian Verhelst, KU Leuven, Belgium
Professor Boris Murmann, Stanford, USA





MON 3.06.19

DEEP LEARNING ALGORITHMS AND PROCESSOR BASICS

1. General overview on the history of machine learning, neural networks and deep learning

- from neural networks to deep neural networks
- from training to inference
- applications

2. Trends in neural network topologies

- Basic concepts of deep neural networks
- Classes and evolution of deep neural networks
- Algorithmic efficiency enhancement techniques

3. Processor architectures for deep neural networks

- Basic processor components
- Efficiency enhancement techniques exploiting sparsity
- Efficiency enhancement techniques exploiting reduced and variable precision processing
- Illustrations with various recent chips from the international state-of-the-art



PROFESSOR MARIAN VERHELST



TUE 4.06.19

EXPLOITING MIXED-SIGNAL AND WITHIN-MEMORY COMPUTATION

1. Exploiting mixed-signal processing for machine learning

- Benchmarking of analog versus digital computing
- Mixed-signal circuits for hand-crafted classifiers
- Mixed-signal circuits for deep neural networks

2. Exploiting in-memory processing for machine learning

- In-memory computing using SRAM
- In-memory computing using RRAM



PROFESSOR BORIS MURMANN



WED 5.06.19

OPTIMIZING ACROSS THE COMPLETE ALGORITHMS - ARCHITECTURE – CIRCUITS STACK

1. The impact of data flow on system efficiency

- Parameters impacting processor efficiency
- Data flow impacting data reuse
- Data-flow driven processing architectures and memory architectures

2. Algorithm-hardware co-optimization for energy efficient inference

- Cross-layer optimization strategy
- Bringing hardware into the loop

3. Discussion and wrap up, trends and outlook



PROFESSOR MARIAN VERHELST

PROGRAM OVERVIEW

HARDWARE-EFFICIENT MACHINE LEARNING: DESIGNING ACROSS THE CIRCUIT-ARCHITECTURE-ALGORITHM STACK



MON 3.06.19

9.00 - 10.30
10.30 - 11.00
11.00 - 12.30
12.30 - 13.30
13.30 - 15.00
15.00 - 15.30
15.30 - 17.00

General overview on the history of machine learning, neural networks and deep learning
Coffee break
Trends in neural network topologies
Lunch break
Processor architectures for deep neural networks
Coffee break
Processor architectures for deep neural networks (part 2)



TUE 4.06.19

9.00 - 10.30
10.30 - 11.00
11.00 - 12.30
12.30 - 13.30
13.30 - 15.00
15.00 - 15.30
15.30 - 17.00

Exploiting mixed-signal processing for machine learning
Coffee break
Exploiting mixed-signal processing for machine learning (Part 2)
Lunch break
Exploiting in-memory processing for machine learning
Coffee break
Exploiting in-memory processing for machine learning (Part 2)



WED 5.06.19

9.00 - 10.30
10.30 - 11.00
11.00 - 12.30
12.30 - 13.30
13.30 - 15.00
15.00 - 15.30
15.30 - 17.00

The impact of data flow on system efficiency
Coffee break
The impact of data flow on system efficiency (Part 2)
Lunch break
Algorithm-hardware co-optimization for energy efficient inference
Coffee break
Discussion and wrap-up, trends and outlook