

Mirko Mariotti ^{1,2} Giulio Bianchini ¹ Loriano Storchi ^{3,2} Giacomo Surace ² Daniele Spiga ²

¹Dipartimento di Fisica e Geologia, Universitá degli Studi di Perugia

²INFN sezione di Perugia

³Dipartimento di Farmacia, Universitá degli Studi G. D'Annunzio

Third ML-INFN Hackathon: Advanced Level



Introduction

FPGA HDL workflow HLS Workflow Concepts Cloud

2 The BondMachine project

Architectures handling Architectures molding Bondgo Basm API 3 Misc Project timeline

4 Machine Learning

BondMachine creation Simulation Accelerator Benchmark

5 Optimizations

6 Conclusions and Future directions Conclusions Ongoing Future

Hand-on sessions

Some topic will have a hands-on sessions, you can either:

Install the tools in you laptop following quickstart at http://bondmachine.fisica.unipg.it/docs

Use the container of the hackaton

After that just clone the examples repository:

git clone https://github.com/BondMachineHQ/bmexamples.git



Current challenges in computing

Von Neumann Bottleneck:

New computational problems show that current architectural models has to be improved or changed to address future payloads.

Energy Efficient computation:

Not wasting "resources" (silicon, time, energy, instructions). Using the right resource for the specific case

Edge/Fog/Cloud Computing: Making the computation where it make sense Avoiding the transfer of unnecessary data Creating consistent interfaces for distributed systems

Third ML-INFN Hackathon: Advanced Level

Current challenges in computing

Von Neumann Bottleneck:

New computational problems show that current architectural models has to be improved or changed to address future payloads.

Energy Efficient computation: Not wasting "resources" (silicon, time, energy, instructions). Using the right resource for the specific case

Edge/Fog/Cloud Computing: Making the computation where it make sense Avoiding the transfer of unnecessary data Creating consistent interfaces for distributed systems

Third ML-INFN Hackathon: Advanced Level

Current challenges in computing

Von Neumann Bottleneck:

New computational problems show that current architectural models has to be improved or changed to address future payloads.

Energy Efficient computation: Not wasting "resources" (silicon, time, energy, instructions). Using the right resource for the specific case

Edge/Fog/Cloud Computing: Making the computation where it make sense Avoiding the transfer of unnecessary data Creating consistent interfaces for distributed systems

Third ML-INFN Hackathon: Advanced Level

A field programmable gate array (FPGA) is an integrated circuit whose logic is re-programmable.

- Parallel computing Highly specialized
- Energy efficient





- Array of programmable logic blocks
 - Logic blocks configurable to perform complex functions
- The configuration is specified with the hardware description language



General reference Cell model



A Cell usually contains:

- A **look-up table**: allow to map any combinatorial function of 4 inputs and 1 output.
- a FF of type D: to store persistent data.
- A **mux 2 -> 1**: to eventually bypass the FF for purely combinatorial cells.



CLOCK A RES

Programmable element

Third ML-INFN Hackathon: Advanced Level



Third ML-INFN Hackathon: Advanced Level



The use of FPGA in computing is growing due several reasons:

can potentially deliver great performance via massive parallelism

can address payloads which are not performing well on uniprocessors (Neural Networks, Deep Learning)

can handle efficiently non-standard data types

Third ML-INFN Hackathon: Advanced Level



The use of FPGA in computing is growing due several reasons:

can potentially deliver great performance via massive parallelism

can address payloads which are not performing well on uniprocessors (Neural Networks, Deep Learning)

can handle efficiently non-standard data types

Third ML-INFN Hackathon: Advanced Level



The use of FPGA in computing is growing due several reasons:

can potentially deliver great performance via massive parallelism

can address payloads which are not performing well on uniprocessors (Neural Networks, Deep Learning)

can handle efficiently non-standard data types

Third ML-INFN Hackathon: Advanced Level

Integration of neural networks on FPGA

FPGAs are playing an increasingly important role in the industry sampling and data processing.





Deep Learning

In the industrial field

- Intelligent vision;
- Financial services;
- Scientific simulations;
- Life science and medical data analysis;

In the scientific field

- Real time deep learning in particle physics;
- Hardware trigger of LHC experiments;
- And many others ...

CPU vs GPU vs FPGA: Instructions, memory and parallelism

CPU	CPU GPU		
Fixed architecture	Fixed architecture	Adaptable Architecture	
Predefined instruction set	Predefined instruction set	No fixed instruction set	
Fixed memory hierarchy	Fixed memory hierarchy	Customizable memory hierarchy	
Thread-level parallelism	SIMD parallelism	Excels at all types of parallelism	

Third ML-INFN Hackathon: Advanced Level



FPGA Performance is predictable

There is no context switch, garbage collector or any background process

The bitstream will be executed the same number of clock cycles every time

The number of clock cycles needed can be computed easily

Slide credits: M.Barbone

Third ML-INFN Hackathon: Advanced Level



Slide credits: M.Barbone

Third ML-INFN Hackathon: Advanced Level



On the other hand the adoption on FPGA has several drawbacks:

Porting of legacy code is usually hard.

Interoperability with standard applications is problematic.

Third ML-INFN Hackathon: Advanced Level



On the other hand the adoption on FPGA has several drawbacks:

Porting of legacy code is usually hard.

Interoperability with standard applications is problematic.

FPGA Programming workflow	HDL code	







CPU vs	s GPU vs FPGA	(HDL): Coding	g difficulty	
		CPU	GPU	FPGA
	Engineering time	short	medium	very long
	Compilation time	short	short	very long
	Debugging	easy	medium	hard
	Maintainability	easy	medium	hard

Slide credits: M.Barbone



Slide credits: M.Barbone

Third ML-INFN Hackathon: Advanced Level

CPU vs GPU vs FPGA (HLS): Coding difficulty

		CPU	GPU	FPGA Sim	FPGA Hw
	Engineering time	short	medium	medium	medium
	Compilation time	short	short	short	Very long
	Debugging	easy	medium	medium	medium
	Maintainability	easy	medium	medium	medium

Slide credits: M.Barbone

Third ML-INFN Hackathon: Advanced Level

Firmware generation

Many projects have the goal of abstracting the firmware generation and use process.



Third ML-INFN Hackathon: Advanced Level



Third ML-INFN Hackathon: Advanced Level



Throughput: quantity of outputs in the unit of time

Third ML-INFN Hackathon: Advanced Level

Latency in FPGA

FPGAs achieve lower latency than software: latency below the microsecond latency constant and predictable

CPUs are not suitable as thread overhead is order of 10 microseconds GPUs are even worse since PCIe bus sys-call could require milliseconds

Examples:

- Ultra low latency trading
- CERN trigger system

Slide credits: M.Barbone

Third ML-INFN Hackathon: Advanced Level

Is the amount of the various resources used on and FPGA (LUTs, DSP, etc) Problems:

- The algorithms has to create a circuit that can be contained in the available resources
- Tools to shrink or grow the occupancy
- The need for a runtime, platforms



Occupancy



Microsoft's Project Brainwave is a deep learning platform for real-time AI inference in the cloud and on the edge. A soft Neural Processing Unit (NPU), based on a high-performance field-programmable gate array (FPGA).

Cloud resources





Third ML-INFN Hackathon: Advanced Level



Third ML-INFN Hackathon: Advanced Level



Building a new kind of computer architecture (multi-core and heterogeneous both in cores types and interconnections) which dynamically adapt to the specific computational problem rather than be static.



Third ML-INFN Hackathon: Advanced Level




Third ML-INFN Hackathon: Advanced Level





Third ML-INFN Hackathon: Advanced Level





Third ML-INFN Hackathon: Advanced Level





Third ML-INFN Hackathon: Advanced Level



Basm API 6 Conclusions and Future directions Conclusions Ongoing Future

The BondMachine is a software ecosystem for the dynamic generation of computer architectures that:

Are composed by many, possibly hundreds, computing cores.

- Have very small cores and not necessarily of the same type (different ISA and ABI).
- Have a not fixed way of interconnecting cores.
- May have some elements shared among cores (for example channels and shared memories).

The BondMachine is a software ecosystem for the dynamic generation of computer architectures that:

Are composed by many, possibly hundreds, computing cores.

- Have very small cores and not necessarily of the same type (different ISA and ABI).
- Have a not fixed way of interconnecting cores.
- May have some elements shared among cores (for example channels and shared memories).

The BondMachine is a software ecosystem for the dynamic generation of computer architectures that:

Are composed by many, possibly hundreds, computing cores.

Have very small cores and not necessarily of the same type (different ISA and ABI).

Have a not fixed way of interconnecting cores.

May have some elements shared among cores (for example channels and shared memories).

The BondMachine is a software ecosystem for the dynamic generation of computer architectures that:

Are composed by many, possibly hundreds, computing cores.

Have very small cores and not necessarily of the same type (different ISA and ABI).Have a not fixed way of interconnecting cores.

May have some elements shared among cores (for example channels and shared memories).

The BondMachine is a software ecosystem for the dynamic generation of computer architectures that:

Are composed by many, possibly hundreds, computing cores.

- Have very small cores and not necessarily of the same type (different ISA and ABI).
 Have a not fixed way of interconnecting cores.
- May have some elements shared among cores (for example channels and shared memories).



The computational unit of the BM

The atomic computational unit of a BM is the "connecting processor" (CP) and has:

Some general purpose registers of size RSize. Some I/O dedicated registers of size Rsize. A set of implemented opcodes chosen among many availabl Dedicated ROM and RAM.

Three possible operating modes.

The computational unit of the BM

The atomic computational unit of a BM is the "connecting processor" (CP) and has:

Some general purpose registers of size Rsize. Some I/O dedicated registers of size Rsize. A set of implemented opcodes chosen among many available. Dedicated ROM and RAM.

Three possible operating modes.

General purpose registers

 2^R registers: r0,r1,r2,r3 ... r 2^R

Third ML-INFN Hackathon: Advanced Level

The computational unit of the BM

The atomic computational unit of a BM is the "connecting processor" (CP) and has:

Some general purpose registers of size Rsize

Some I/O dedicated registers of size Rsize.

A set of implemented opcodes chosen among many available.

Dedicated ROM and RAM.

Three possible operating modes.

I/O specialized registers

N input registers: i0,i1 ... iN M output registers: o0,o1 ... oM

Third ML-INFN Hackathon: Advanced Level

The computational unit of the BM

The atomic computational unit of a BM is the "connecting processor" (CP) and has:

Some general purpose registers of size Rsize.

A set of implemented opcodes chosen among many available.

Dedicated ROM and RAM

Three possible operating modes.

Full set of possible opcodes

adc,add,addf,addi,and,chc,chw,cil,cilc,cir,cirn,clc,clr,cpy,cset,dec,div,divf,dpc,expf,hit hlt,i2r,i2rw,incc,inc,j,jc,je,jgt0f,jlt,jlte,jr,jz,lfsr82,lfsr162r,m2r,mod,mulc,mult,multf nand,nop,nor,not,or,r2m,r2o,r2owa,r2owaa,r2s,r2v,r2vri,ro2r,ro2rri,rsc,rset,sic,s2r,saj,sbc sub,wrd,wwr,xnor,xor

Third ML-INFN Hackathon: Advanced Level

The computational unit of the BM

The atomic computational unit of a BM is the "connecting processor" (CP) and has:

Some general purpose registers of size Rsize. Some I/O dedicated registers of size Rsize. A set of implemented opcodes chosen among many availab

Dedicated ROM and RAM.

Three possible operating modes.

RAM and ROM

2^L RAM memory cells.

2⁰ ROM memory cells.

The computational unit of the BM

The atomic computational unit of a BM is the "connecting processor" (CP) and has:

Some general purpose registers of size RSize. Some I/O dedicated registers of size Rsize. A set of implemented opcodes chosen among many available. Dedicated ROM and RAM.

Three possible operating modes.

Operating modes

Full Harvard mode.

- Full Von Neuman mode.
- Hybrid mode.

Third ML-INFN Hackathon: Advanced Level

Shared Objects (SO)

The non-computational element of the BM

Alongside CPs, BondMachines include non-computing units called "Shared Objects" (SO). Examples of their purposes are:

- Data storage (Memories)
- Message passing.
- CP synchronization.

A single SO can be shared among different CPs. To use it CPs have special instructions (opcodes) oriented to the specific SO.

Four kind of SO have been developed so far: the Channel, the Shared Memory, the Barrier and a Pseudo Random Numbers Generator.

Third ML-INFN Hackathon: Advanced Level

Alongside CPs, BondMachines include non-computing units called "Shared Objects" (SO).

Examples of their purposes are:

Data storage (Memories).

Message passing

CP synchronization.

A single SO can be shared among different CPs. To use it CPs have special instructions (opcodes) oriented to the specific SO.

Four kind of SO have been developed so far: the Channel, the Shared Memory, the Barrier and a Pseudo Random Numbers Generator.

Third ML-INFN Hackathon: Advanced Level

Alongside CPs, BondMachines include non-computing units called "Shared Objects" (SO).

Examples of their purposes are:

Data storage (Memories).

Message passing.

CP synchronization.

A single SO can be shared among different CPs. To use it CPs have special instructions (opcodes) oriented to the specific SO.

Four kind of SO have been developed so far: the Channel, the Shared Memory, the Barrier and a Pseudo Random Numbers Generator.

Third ML-INFN Hackathon: Advanced Level

Alongside CPs, BondMachines include non-computing units called "Shared Objects" (SO).

Examples of their purposes are:

Data storage (Memories).

Message passing.

CP synchronization.

A single SO can be shared among different CPs. To use it CPs have special instructions (opcodes) oriented to the specific SO.

Four kind of SO have been developed so far: the Channel, the Shared Memory, the Barrier and a Pseudo Random Numbers Generator.

Third ML-INFN Hackathon: Advanced Level

Alongside CPs, BondMachines include non-computing units called "Shared Objects" (SO).

Examples of their purposes are:

Data storage (Memories).

Message passing.

CP synchronization.

A single SO can be shared among different CPs. To use it CPs have special instructions (opcodes) oriented to the specific SO.

Four kind of SO have been developed so far: the Channel, the Shared Memory, the Barrier and a Pseudo Random Numbers Generator.

Third ML-INFN Hackathon: Advanced Level

Alongside CPs, BondMachines include non-computing units called "Shared Objects" (SO).

Examples of their purposes are:

Data storage (Memories).

Message passing.

CP synchronization.

A single SO can be shared among different CPs. To use it CPs have special instructions (opcodes) oriented to the specific SO.

Four kind of SO have been developed so far: the Channel, the Shared Memory, the Barrier and a Pseudo Random Numbers Generator.

Third ML-INFN Hackathon: Advanced Level

The BM computer architecture is managed by a set of tools to:

build a specify architecture

modify a pre-existing architecture

simulate or emulate the behavior

generate the Hardware Description Language Code (HDL)

Processor Builder

Selects the single processor, assembles and disassembles, saves on disk as JSON, creates the HDL code of a CP BondMachine Builder

Connects CPs and SOs together in custom topologies, loads and saves on disk as JSON, create BM's HDL code Simulates the behaviour, emulates a BM on a standard Linux workstation

Third ML-INFN Hackathon: Advanced Level

The BM computer architecture is managed by a set of tools to:

build a specify architecture

modify a pre-existing architecture

simulate or emulate the behavior

generate the Hardware Description Language Code (HDL)

Processor Builder

Selects the single processor, assembles and disassembles, saves on disk as JSON, creates the HDL code of a CP BondMachine Builder

Connects CPs and SOs together in custom topologies, loads and saves on disk as JSON, create BM's HDL code Simulates the behaviour, emulates a BM on a standard Linux workstation

The BM computer architecture is managed by a set of tools to:

build a specify architecture

modify a pre-existing architecture

simulate or emulate the behavior

generate the Hardware Description Language Code (HDL)

Processor Builder

Selects the single processor, assembles and disassembles, saves on disk as JSON, creates the HDL code of a CP BondMachine Builder

Connects CPs and SOs together in custom topologies, loads and saves on disk as JSON, create BM's HDL code Simulation Framework Simulates the behaviour, emulates a BM on a standard Linux workstation

The BM computer architecture is managed by a set of tools to:

build a specify architecture

modify a pre-existing architecture

simulate or emulate the behavior

generate the Hardware Description Language Code (HDL)

Processor Builder

Selects the single processor, assembles and disassembles, saves on disk as JSON, creates the HDL code of a CP BondMachine Builder

Connects CPs and SOs together in custom topologies, loads and saves on disk as JSON, create BM's HDL code Simulation Framework

Simulates the behaviour, emulates a BM on a standard Linux workstation

A set of toolchains allow the build and the direct deploy to a target device of BondMachines

Bondgo Toolchain main targets

A file local.mk contains references to the source code as well all the build necessities make bondmachine creates the JSON representation of the BM and assemble its code make hdl creates the HDL files of the BM make show displays a graphical representation of the BM make simulate [simbatch] start a simulation [batch simulation] make bitstream [design_bitstream] create the firwware [accelerator firmware] make program flash the device into the destination target

Third ML-INFN Hackathon: Advanced Level

Toolchains

Simulation

An important feature of the tools is the possibility of simulating BondMachine behavior.

An event input file describes how BondMachines elements has to change during the simulation timespan and which one has to be be reported.

The simulator can produce results in the form of:

- Activity log of the BM internal.
- Graphical representation of the simulation.
- Report file with quantitative data. Useful to construct metrics

Graphical simulation in action

Third ML-INFN Hackathon: Advanced Level

Simulation

An important feature of the tools is the possibility of simulating BondMachine behavior.

An event input file describes how BondMachines elements has to change during the simulation timespan and which one has to be be reported.

The simulator can produce results in the form of:

- Activity log of the BM internal.
- Graphical representation of the simulation.
- Report file with quantitative data. Useful to construct metrics

Graphical simulation in action

Third ML-INFN Hackathon: Advanced Level



An important feature of the tools is the possibility of simulating BondMachine behavior.

An event input file describes how BondMachines elements has to change during the simulation timespan and which one has to be be reported.

The simulator can produce results in the form of:

- Activity log of the BM internal.
- Graphical representation of the simulation.
- Report file with quantitative data. Useful to construct metrics

Graphical simulation in action

Third ML-INFN Hackathon: Advanced Level

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning

Third ML-INFN Hackathon: Advanced Level

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning

Third ML-INFN Hackathon: Advanced Level

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning

Third ML-INFN Hackathon: Advanced Level

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning

Third ML-INFN Hackathon: Advanced Level

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning

Third ML-INFN Hackathon: Advanced Level
Molding the BondMachine

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning

Third ML-INFN Hackathon: Advanced Level

Molding the BondMachine

As stated before BondMachines are not general purpose architectures, and to be effective have to be shaped according the specific problem.

Several methods (apart from writing in assembly and building a BondMachine from scratch) have been developed to do that:

bondgo: A new type of compiler that create not only the CPs assembly but also the architecture itself.

basm: The BondMachine Assembler.

A set of API to create BondMachine to fit a specific computational problems.

An Evolutionary Computation framework to "grow" BondMachines according some fitness function via simulation.

A set of tools to use BondMachine in Machine Learning.

Third ML-INFN Hackathon: Advanced Level



Mapping specific computational problems to BMs



more about these tools

Third ML-INFN Hackathon: Advanced Level



more about these tools

Third ML-INFN Hackathon: Advanced Level



more about these tools

Third ML-INFN Hackathon: Advanced Level



more about these tools

Third ML-INFN Hackathon: Advanced Level



more about these tools

Third ML-INFN Hackathon: Advanced Level



more about these tools

Third ML-INFN Hackathon: Advanced Level



more about these tools

Third ML-INFN Hackathon: Advanced Level



Bondgo is the name chosen for the compiler developed for the BondMachine.

The compiler source language is Go as the name suggest.

Bondgo

This is the standard flow when building computer programs

Bondgo

This is the standard flow when building computer programs

high level language source

Third ML-INFN Hackathon: Advanced Level





Bondgo

Bondgo does something different from standard compilers

Bondgo

Bondgo does something different from standard compilers ...

high level GO source

Third ML-INFN Hackathon: Advanced Level

































Bondgo

... it can do even much more interesting things when compiling concurrent programs.

Bondgo

... it can do even much more interesting things when compiling concurrent programs.

high level GO source














Compiling Architectures	

One of the most important result

The architecture creation is a part of the compilation process.

Third ML-INFN Hackathon: Advanced Level



The BondMachine assembler *Basm* is the compiler complementary tools. The BondMachine "fluid" nature gives the assembler some unique features: Support for template based assembly code

Combining and rewriting fragments of assembly code

Building hardware from assembly

Software/Hardware rearrange capabilities

Third ML-INFN Hackathon: Advanced Level



Support for template based assembly code

Combining and rewriting fragments of assembly code

Building hardware from assembly

Software/Hardware rearrange capabilities

Third ML-INFN Hackathon: Advanced Level



Support for template based assembly code

Combining and rewriting fragments of assembly code

Building hardware from assembly

Software/Hardware rearrange capabilities

Third ML-INFN Hackathon: Advanced Level



Support for template based assembly code

Combining and rewriting fragments of assembly code

Building hardware from assembly

Software/Hardware rearrange capabilities

Third ML-INFN Hackathon: Advanced Level



Support for template based assembly code

Combining and rewriting fragments of assembly code

Building hardware from assembly

Software/Hardware rearrange capabilities

Third ML-INFN Hackathon: Advanced Level

Abstract Assembly

The Assembly language for the BM has been kept as independent as possible from the particular CP.

Given a specific piece of assembly code Bondgo has the ability to compute the "minimum CP" that can execute that code.

·			
i2r	r0	i0	
i2r	r1	i1	
add	r0	r1	
r2o	r0	o0	
j 0			



These are Building Blocks for complex BondMachines.

Third ML-INFN Hackathon: Advanced Level

With these Building Blocks Several libraries have been developed to map specific problems on BondMachines: Symbond, to handle mathematical expression.

Boolbond, to map boolean expression.

Matrixwork, to perform matrices operations.

more about these tools

Third ML-INFN Hackathon: Advanced Level

Builders API

With these Building Blocks

Several libraries have been developed to map specific problems on BondMachines:

Symbond, to handle mathematical expression.

Boolbond, to map boolean expression.

Matrixwork, to perform matrices operations.

more about these tools

Third ML-INFN Hackathon: Advanced Level

Builders API

With these Building Blocks Several libraries have been developed to map specific problems on BondMachines:

Symbond, to handle mathematical expression.

Boolbond, to map boolean expression.

Matrixwork, to perform matrices operations.

more about these tools

Third ML-INFN Hackathon: Advanced Level

Builders API

With these Building Blocks Several libraries have been developed to map specific problems on BondMachines:

- Symbond, to handle mathematical expression.
- Boolbond, to map boolean expression.
- Matrixwork, to perform matrices operations.

more about these tools

Builders API



Talk with details about how the accelerator is build

Third ML-INFN Hackathon: Advanced Level





Bondge Basm API 6 Conclusions and Future directions Conclusions Ongoing Future

The BondMachine is a software ecosystem for the dynamical generation (from several HL types of origin) of computer architectures that can be synthesized of FPGA and

used as standalone devices,

as clustered devices,

and as firmware for computing accelerators.

Third ML-INFN Hackathon: Advanced Level



The BondMachine is a software ecosystem for the dynamical generation (from several HL types of origin) of computer architectures that can be synthesized of FPGA and

used as standalone devices,

as clustered devices,

and as firmware for computing accelerators.



Third ML-INFN Hackathon: Advanced Level

The BondMachine is a software ecosystem for the dynamical generation (from several HL types of origin) of computer architectures that can be synthesized of FPGA and

used as standalone devices,

as clustered devices,

and as firmware for computing accelerators.



Third ML-INFN Hackathon: Advanced Level

The BondMachine is a software ecosystem for the dynamical generation (from several HL types of origin) of computer architectures that can be synthesized of FPGA and

used as standalone devices,

as clustered devices,

and as firmware for computing accelerators.

Third ML-INFN Hackathon: Advanced Level

CCR 2015 First ideas, 2016 Poster, 2017 Talk

InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA Invited lectures at: "Advanced Workshop on Modern FPGA Based Technology for Scientific Computing", ICTP 2019

Invited lectures at: "NiPS Summer School 2019

- Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"
- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022



PON PHD program

Third ML-INFN Hackathon: Advanced Level

CCR 2015 First ideas, 2016 Poster, 2017 Talk InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA

Modern FPGA Based Technology for Scientific Computing", ICTP 2019

> Invited lectures at: "NiPS Summer School 2019 – Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"

- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022



PON PHD program

Third ML-INFN Hackathon: Advanced Level

CCR 2015 First ideas, 2016 Poster, 2017 Talk InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA

Invited lectures at: "Advanced Workshop on Modern FPGA Based Technology for Scientific Computing", ICTP 2019

Invited lectures at: "NiPS Summer School 2019 – Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"

- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022



PON PHD program

Third ML-INFN Hackathon: Advanced Level

- CCR 2015 First ideas, 2016 Poster, 2017 Talk InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA
- Invited lectures at: "Advanced Workshop on Modern FPGA Based Technology for Scientific Computing", ICTP 2019
- Invited lectures at: "NiPS Summer School 2019 – Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"
- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022

The BondMachine Toolkit

Mirko Mariotti

Department of Physics and Geology - University of Perugia INFN Perugia

NiPS Summer School 2019 Architectures and Algorithms for Energy-Efficient IoT and HPC Applications 3-6 September 2019 - Perugia





PON PHD program

Third ML-INFN Hackathon: Advanced Level

- CCR 2015 First ideas, 2016 Poster, 2017 Talk InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA
- Invited lectures at: "Advanced Workshop on Modern FPGA Based Technology for Scientific Computing", ICTP 2019
- Invited lectures at: "NiPS Summer School 2019 – Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"
- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022



PON PHD program

Third ML-INFN Hackathon: Advanced Level

- CCR 2015 First ideas, 2016 Poster, 2017 Talk InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA
- Invited lectures at: "Advanced Workshop on Modern FPGA Based Technology for Scientific Computing", ICTP 2019
- Invited lectures at: "NiPS Summer School 2019
- Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"
- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022



PON PHD program

Third ML-INFN Hackathon: Advanced Level

- CCR 2015 First ideas, 2016 Poster, 2017 Talk InnovateFPGA 2018 Iron Award, Grand Final at Intel Campus (CA) USA
- Invited lectures at: "Advanced Workshop on Modern FPGA Based Technology for Scientific Computing", ICTP 2019
- Invited lectures at: "NiPS Summer School 2019
- Architectures and Algorithms for Energy-Efficient IoT and HPC Applications"
- Golab 2018 talk and ISGC 2019 PoS
- Article published on Parallel Computing, Elsevier 2022
- PON PHD program

Third ML-INFN Hackathon: Advanced Level



Machine Learning with BondMachine

Architectures with multiple interconnected processors like the ones produced by the BondMachine Toolkit are a perfect fit for Neural Networks and Computational Graphs.

Several ways to map this structures to BondMachine has been developed:

- A native Neural Network library
- A Tensorflow to BondMachine translator
- An NNEF based BondMachine composer

Machine Learning with BondMachine

Architectures with multiple interconnected processors like the ones produced by the BondMachine Toolkit are a perfect fit for Neural Networks and Computational Graphs.

Several ways to map this structures to BondMachine has been developed:

- A native Neural Network library
- A Tensorflow to BondMachine translator
- An NNEF based BondMachine composer

Machine Learning with BondMachine Native Neural Network library

The tool *neuralbond* allow the creation of BM-based neural chips from an API go interface.

Neurons are converted to BondMachine connecting processors.

Tensors are mapped to CP connections.



Third ML-INFN Hackathon: Advanced Level





Machine Learning with BondMachine NNEF Composer

Neural Network Exchange Format (NNEF) is a standard from Khronos Group to enable the easy transfer of trained networks among frameworks, inference engines and devices

The NNEF BM tool approach is to descent NNEF models and build BondMachine multi-core accordingly

This approch has several advandages over the previous:

- It is not limited to a single framework
- NNEF is a textual file, so no complex operations are needed to read models


FPGA

- Digilent Zedboard
- Soc: Zynq XC7Z020-CLG484-1
- 512 MB DDR3
- Vivado 2020.2
- 100MHz
- PYNQ 2.6 (custom build)



BM inference: A first tentative idea

A neuron of a neural network can be seen as Connecting Processor of BM

H1

X1

X2

Х3

Χ4



 e^{z_j}

%section softmax .romtext iomode:sync
entry _start ; Entry point
_start:
mov r8, 0f0.0
<pre>{{range \$y := intRange "0" .Params.inputs}}</pre>
{{printf "i2r r1,i%d\n" \$y}}
mov r0, 0f1.0
mov r2, 0f1.0
mov r3, 0f1.0
mov r4, 0f1.0
mov r5, 0f1.0
<pre>mov r7, {{\$.Params.expprec}}</pre>
<pre>loop{{printf "%d" \$y}}:</pre>
multf r2, r1
multf r3, r4
addf r4, r5
mov r6, r2
divf r6,r3
addf r0, r6
dec r/
jz i/,ext((printi %d \$y))
j (cop((print) %d \$y))
exter(print) -xa- syy):
(152 := alot \$.Params.pos))
((() eq by \$2})
hondraction
wendset t ton

inputs hidden layer output layer outputs

S1

S2

Y1

Y2

Third ML-INFN Hackathon: Advanced Level

Machine Learning on FPGA

-6 -4 -2 0 2 4

 $\sigma(\vec{z})_i$

From idea to implementation

Starting from High Level Code, a NN model trained with **TensorFlow** and exported in a standard interpreted by **neuralbond** that converts nodes and weights of the network into a set of heterogeneous processors.



Third ML-INFN Hackathon: Advanced Level



A first test Dataset info:

- **Dataset name**: Banknote Authentication
- **Description**: Dataset on the distinction between genuine and counterfeit banknotes. The data was extracted from images taken from genuine and fake banknote-like samples.
- N. features: 4
- Classification: binary
- **Samples**: 1097

Neural network info: Class: Multilayer perceptron fully connected

Layers:

 An hidden layer with 1 linear neuron
One output layer with 2 softmax neurons

Graphic representation:



Third ML-INFN Hackathon: Advanced Level



To train the simple model and prepare it from the BM

Reference notebook: banknote-train.ipynb

Third ML-INFN Hackathon: Advanced Level



sw.csv is the software predictions over that dataset and will be used to check the BM inference probabilities and predictions

modelBM.json is the trained network that will use as BM source in the next demo

Third ML-INFN Hackathon: Advanced Level



Use modelBM created on the previous step to create a BondMachine

Reference notebook: proj zedboard ml creation /notebook.ipynb

Third ML-INFN Hackathon: Advanced Level



The outcome of this second part of the hands-on are:

bondmachine.json, a representation of the generated abstract machine

Al the HDL files needed to build the firmware for the given board



Simulate the test dataset with the BondMachine simulator

Compare the results with the keras prediction

Reference notebook: proj_zedboard_ml_simbatch /notebook.ipynb

Third ML-INFN Hackathon: Advanced Level



The outcome of this third part of the hands-on is:

simbatchoutput.csv, a simulated CSV files containing the output probabilities and the prediction

Third ML-INFN Hackathon: Advanced Level



Create the accelerator firmware

Reference notebook: proj_zedboard_ml_accelerator /notebook.ipynb

Third ML-INFN Hackathon: Advanced Level



Notebook on the board - predictions and correctness



Thanks to PYNQ we can easily load the bitstream and program the FPGA in real time.

With their APIs we interact with the memory addresses of the BM IP to send data into the inputs and read the outputs (not using BM kernel module)

Dump output results for future analysis

Open the notebook

Third ML-INFN Hackathon: Advanced Level

Benchcore

Fortunately we have a custom design and an FPGA.

We can put the benchmarks tool inside the accelerator.



Benchcore

Fortunately we have a custom design and an FPGA.

We can put the benchmarks tool inside the accelerator.



Benchcore

Fortunately we have a custom design and an FPGA.

We can put the benchmarks tool inside the accelerator.



Machine Learning on FPGA

Third ML-INFN Hackathon: Advanced Level

Benchcore

Fortunately we have a custom design and an FPGA.

We can put the benchmarks tool inside the accelerator.



Benchcore

Fortunately we have a custom design and an FPGA.

We can put the benchmarks tool inside the accelerator.



Inference evaluation

Evaluation metrics used:

Inference speed: time taken to predict a sample i.e. time between the arrival of the input and the change of the output measured with the **benchcore**; **Resource usage**: luts and registers in use;

Accuracy: as the average percentage of error on probabilities.



Third ML-INFN Hackathon: Advanced Level

 σ : 2875.94

Mean: 10268.45

Latency: 102.68 µs

resource	value	occupancy
regs	15122	28.42%
luts	11192	10.51%

Resource usage

Analysis r	notebook						
	Another no	tebook is i	used to	compare r	runs from	different ac	celerators
		Software				BondMachine	2
	prob0	prob1	class]	prob0	prob1	class
	0.6895	0.3104	0		0.6895	0.3104	0
	0.5748	0.4251	0		0.5748	0.4251	0

1

The output of the bm corresponds to the software output

0.4009

0.5990

1

Reference notebook: analysis /zedboard_banknote /analysis.ipynb

Third ML-INFN Hackathon: Advanced Level

0.4009

0.5990



2 The BondMachine project

Architectures handling Architectures molding Bondgo Basm API

5 Optimizations

6 Conclusions and Future directions Conclusions Ongoing Future

A first example of optimization

Remember the softmax function?

$$\underbrace{\sigma(z_i) = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}}_{$$

$$e^{x} = \sum_{l=0}^{K} \frac{x^{l}}{l!}$$

%section softmax .romtext iomode:sync entry _start ; Entry point _start: mov r8, 0f0.0 {{range \$y := intRange "0" .Params.inputs}} {{print f":2r r1,i%d(n" \$y)} mov r0, 0f1.0 mov r3, 0f1.0

r4, 0f1.0 mov mov r5, 0f1.0 r7, {{\$.Params.expprec}} mov loop{{printf "%d" \$y}} multf r2. r1 r3. r4 multf addf r4, r5 mov r6. r2 divf r6. r3 addf r0, r6 dec r7.exit{{printf "%d" \$v}} loop{{printf "%d" \$v}} exit{{printf "%d" \$v}}: {{\$z := atoi \$.Params.pos}} {{if eq \$v \$z}} mov r9, r0

%endsection



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...



Changing number of K of the exponential factors in the softmax function...





Reduced inference times by a factor of 10 ... only by decreasing the number of iterations.



Third ML-INFN Hackathon: Advanced Level

The tools (neuralbond+basm) create a graph of relations among fragments of assembly

Not necessarily a fragment has to be mapped to a single CP

- They can arbitrarily be rearranged into CPs
- The resulting firmwares are identical in term of the computing outcome, but differs in occupancy and latency.



The tools (neuralbond+basm) create a graph of relations among fragments of assembly Not necessarily a fragment has to be mapped to a single CP

They can arbitrarily be rearranged into CPs The resulting firmwares are identical in term of the computing outcome, but differs in occupancy and latency.



The tools (neuralbond+basm) create a graph of relations among fragments of assembly Not necessarily a fragment has to be mapped to a single CP

They can arbitrarily be rearranged into CPs

The resulting firmwares are identical in term of the computing outcome, but differs in occupancy and latency.



The tools (neuralbond+basm) create a graph of relations among fragments of assembly Not necessarily a fragment has to be mapped to a single CP

They can arbitrarily be rearranged into CPs

The resulting firmwares are identical in term of the computing outcome, but differs in occupancy and latency.




Reference notebook: proj_zedboard_ml_cp_pruning /notebook.ipynb

Third ML-INFN Hackathon: Advanced Level



Collapse processors and find out the consequences

Reference notebook: proj_zedboard_ml_cp_collapsing /notebook.ipynb

Third ML-INFN Hackathon: Advanced Level



Copy a project directory and try pruning, collapsing, simulating and the assembly of the neurons



Third ML-INFN Hackathon: Advanced Level



6 Conclusions and Future directions Conclusions Ongoing Future The BondMachine is a new kind of computing device made possible in practice only by the emerging of new re-programmable hardware technologies such as FPGA.

The result of this process is the construction of a computer architecture that is not anymore a static constraint where computing occurs but its creation becomes a part of the computing process, gaining computing power and flexibility.

Over this abstraction is it possible to create a full computing Ecosystem, ranging from small interconnected IoT devices to Machine Learning accelerators.

Conclusions



Documentation

- First DAQ use case
- Complete the inclusion of Intel and Lattice FPGAs
- ML inference in a cloud workflow



Different data types and operations, especially low and trans-precision

Different boards support, especially data center accelerator

Compare with GPUs

Include some real power consumption measures



Quantization

- More datasets: test on other datasets with more features and multiclass classification
- Neurons: increase the library of neurons to support other activation functions
- **Evaluate results**: compare the results obtained with other technologies (CPU and GPU) in terms of inference speed and energy efficiency



Extend the compiler to include more data structures

Assembler improvements, fragments optimization and others

Improve the networking including new kind of interconnection firmware

What would an OS for BondMachines look like ?

Third ML-INFN Hackathon: Advanced Level

Include new processor shared objects and currently unsupported opcodes Extend the compiler to include more data structures

Assembler improvements, fragments optimization and others

Improve the networking including new kind of interconnection firmware

What would an OS for BondMachines look like ?

Third ML-INFN Hackathon: Advanced Level

Future work

Include new processor shared objects and currently unsupported opcodes

Extend the compiler to include more data structures

Assembler improvements, fragments optimization and others

I Improve the networking including new kind of interconnection firmware

What would an OS for BondMachines look like ?

Third ML-INFN Hackathon: Advanced Level

Future work

Include new processor shared objects and currently unsupported opcodes Extend the compiler to include more data structures

Assembler improvements, fragments optimization and others

Improve the networking including new kind of interconnection firmware

What would an OS for BondMachines look like ?

Third ML-INFN Hackathon: Advanced Level

Future work

Include new processor shared objects and currently unsupported opcodes Extend the compiler to include more data structures

Assembler improvements, fragments optimization and others

Improve the networking including new kind of interconnection firmware

What would an OS for BondMachines look like ?

Future work



website: http://bondmachine.fisica.unipg.it code: https://github.com/BondMachineHQ parallel computing paper: link contact email: mirko.mariotti@unipg.it

Third ML-INFN Hackathon: Advanced Level