YoseUe: "trimming, Random Forest's training towards resource-constrained inference on FPGA

A novel approach for inference on large decision tree models, tailored for resource-constrained FPGAs

Alessandro Verosimile

alessandro.verosimile@polimi.it

Andrea Damiani

andrea.damiani@polimi.it



POLITECNICO MILANO 1863 Alessandro Tierno alessandro.tierno@mail.polimi.it

alessand o.tlerno@mail.pointi.tt

Marco Domenico Santambrogio

<u>marco.santambrogio@polimi.it</u>

ASP-DAC Conference 2024 Incheon, 22-Jan-2024





Problem definition

Decision Tree Ensemble models

Inference optimization Embedded systems

Resource-constrained scenario

YoseUe contributions

- A domain-specific architecture for the inference of Random Forests on Embedded systems
- A new Decision tree based ML model, designed to optimize its inference on the developed architecture

Results

YoseUe is able to support DT Ensemble models with 2 orders of magnitude more trees

Agenda



Agenda



Artificial Intelligence of Things

Combinations of two worlds:





Artificial Intelligence of Things

Combinations of two worlds:





It finds application in multiple fields, such as:

- Video surveillance
- Automotive

Weaknesses & Risks of Alot

- Instability of latency/throughput
- Privacy and security risks





Cloud Infrastructure

Ideal solution

Perform all the computations (the inference) on the device



Ideal solution

Perform all the computations (the inference) on the device



Machine Learning Models

Artificial Neural Networks

Decision Tree - based models



Random Forest

"Tree ensembles are arguably among the most accurate ML models in use nowadays" [1]



[1] A. B. Arrieta, N. D' Iaz-Rodr' Iguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garc' Ia, S. Gil-Lo' pez, D. Molina, R. Benjamins *et al.*, "Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai," *Information fusion*, vol. 58, pp. 82–115, 2020.

Random Forest



Random Forest



Alternative architectures



Alternative architectures



Logic-Centric

Memory-Centric





Agenda



State of the art about Logic-Centric approach



- Conifer [2]: the entire model is hardcoded in the hardware leading to high throughput but high resource consumption
- Entrée [3]: employs partial dynamic reconfiguration to fit larger models

[2] S. Summers, G. Di Guglielmo, J. Duarte, P. Harris, D. Hoang, S. Jin- dariani, E. Kreinar, V. Loncar, J. Ngadiuba, M. Pierini et al., "Fast inference of boosted decision trees in fpgas for particle physics," Journal of Instrumentation, vol. 15, no. 05, p. P05026, 2020.

[3] A. Damiani, E. Del Sozzo, and M. D. Santambrogio, "Large forests and where to "partially" fit them," in 2022 27th Asia and South Pacific Design Automation Conference (ASP-DAC). IEEE, 2022, pp. 550–555.



• RF-RISA [4]: implements a specific co-processor and instruction set for RF acceleration. However, it still as significant room for improvement regarding resource-consumption

[4] S. Zhao, S. Chen, H. Yang, F. Wang, and Z. Wei, "Rf-risa: A novel flexible random forest accelerator based on fpga," Journal of Parallel and Distributed Computing, vol. 157, pp. 220–232, 2021







Algorithm









Formally, the ratio of unused memory of RF-RISA is:

$$1 - \frac{|NoInst| \cdot (2^{\hat{d}} - 1)}{|B| \cdot \hat{d}}$$

Legend

- |*B*| is the size of a memory block
- \hat{d} is the maximum depth of the DTs
- |NoInst| is the size of a Node instruction

Agenda











With our improvement, the ratio of unused memory becomes:

$$1 - T \frac{|NoInst| \cdot (2^{\hat{d}} - 1)}{|B| \cdot \hat{d}}$$

Where *T* is the number of trees that can fit into a set of PEs, that can be calculated as:

$$\frac{|B|}{|NoInst| \cdot 2^{\hat{d}-1}}$$



Agenda



NI 0-0	NI 1-0	NI 2-0	NI 3-0
NI 0-0	NI 1-1	NI 2-1	NI 3-1
	NI 1-0	NI 2-2	NI 3-2
	NI 1-1	NI 2-3	NI 3-3
		NI 2-0	NI 3-4
		NI 2-1	NI 3-5
		NI 2-2	NI 3-6
		NI 2-3	NI 3-7
			NI 3-0
			NI 3-1
			NI 3-2
			NI 3-3
			NI 3-4
			NI 3-5
			NI 3-6
			NI 3-7
BRAM 0	BRAM 1	BRAM 2	BRAM 3

Waste of resources



NI 0-0	NI 1-0	NI 2-0
NI 0-0	NI 1-1	NI 2-1
NI 0-0	NI 1-0	NI 2-2
NI 0-0	NI 1-1	NI 2-3
NI 0-0	NI 1-0	NI 2-0
NI 0-0	NI 1-1	NI 2-1
NI 0-0	NI 1-0	NI 2-2
NI 0-0	NI 1-1	NI 2-3
	NI 1-0	NI 2-0
	NI 1-1	NI 2-1
	NI 1-0	NI 2-2
	NI 1-1	NI 2-3
	NI 1-0	NI 2-0
	NI 1-1	NI 2-1
	NI 1-0	NI 2-2
	NI 1-1	NI 2-3
BRAM 0	BRAM 1	BRAM 2

NI 3-0				
NI 3-1				
NI 3-2				
NI 3-3				
NI 3-4				
NI 3-5				
NI 3-6				
NI 3-7				
NI 3-0				
NI 3-1				
NI 3-2				
NI 3-3				
NI 3-4				
NI 3-5				
NI 3-6				
NI 3-7				
BRAM 3				

NI 0-0	NI 1-0	NI 2-0	NI 3-0
NI 0-0	NI 1-1	NI 2-1	NI 3-1
NI 0-0	NI 1-0	NI 2-2	NI 3-2
NI 0-0	NI 1-1	NI 2-3	NI 3-3
NI 0-0	NI 1-0	NI 2-0	NI 3-4
NI 0-0	NI 1-1	NI 2-1	NI 3-5
NI 0-0	NI 1-0	NI 2-2	NI 3-6
NI 0-0	NI 1-1	NI 2-3	NI 3-7
	NI 1-0	NI 2-0	NI 3-0
	NI 1-1	NI 2-1	NI 3-1
	NI 1-0	NI 2-2	NI 3-2
	NI 1-1	NI 2-3	NI 3-3
	NI 1-0	NI 2-0	NI 3-4
	NI 1-1	NI 2-1	NI 3-5
	NI 1-0	NI 2-2	NI 3-6
	NI 1-1	NI 2-3	NI 3-7
BRAM 0	BRAM 1	BRAM 2	BRAM 3





Multi-depth RF



The tests are performed over 5 datasets, Satellite, Shuttle, Accelerometer, Vehicles, Vowel from the UCI ML Repository [5]



[5] K. N. Markelle Kelly, Rachel Longjohn, "The uci machine learning repository." [Online]. Available: https://archive.ics.uci.edu















Agenda





Axi-stream

 ∇

BRAM port





Axi-stream

 \setminus /

BRAM port





Experimental setup

AVNET Ultra96-v2



141,160 70,560 FFs LUTs

36K BRAMs 216

Zynq Ultrascale+ MPSoC XCZU3EG

Experimental results

Throughput



RF-RISA [4] re-implementation on Ultra96-v2

YoseUe on Ultra96-v2

 $M = x10^9$

Tables refer to an implementation with samples composed of 5 attributes (ap_fixed<32,16>)

[4] S. Zhao, S. Chen, H. Yang, F. Wang, and Z. Wei, "Rf-risa: A novel flexible random forest accelerator based on fpga," Journal of Parallel and Distributed Computing, vol. 157, pp. 220–232, 2021

Experimental results

Synthesizable configurations



RF-RISA [4] re-implementation on Ultra96-v2

YoseUe on Ultra96-v2

Tables refer to an implementation with samples composed of 5 attributes (ap_fixed<32,16>)

[4] S. Zhao, S. Chen, H. Yang, F. Wang, and Z. Wei, "Rf-risa: A novel flexible random forest accelerator based on fpga," Journal of Parallel and Distributed Computing, vol. 157, pp. 220–232, 2021

Conclusions

• MDRFC: new Ensemble DT-based model for efficient inference

 A DSA to infer the MDRFC, allowing the inference of two orders of magnitudes more trees w.r.t. SoA

Future Work

1.

Allowing the execution of an unbounded number of trees by adding pre-fetching modules for dynamically loading new instructions into memories.

2.

Study different topologies of the architecture to carry out inference in a more efficient way.

Thanks for your attention!

Alessandro Verosimile alessandro.verosimile@polimi.it Alessandro Tierno alessandro.tierno@mail.polimi.it

Andrea Damiani andrea.damiani@polimi.it Marco Domenico Santambrogio <u>marco.santambrogio@polimi.it</u>

CREDITS: This presentation template was created by Slidesgo, including icons by Flaticon, infographics & images by Freepik and illustrations by Stories

Thanks for your attention!

Alessandro Verosimile alessandro.verosimile@polimi.it Alessandro Tierno alessandro.tierno@mail.polimi.it

Andrea Damiani andrea.damiani@polimi.it Marco Domenico Santambrogio <u>marco.santambrogio@polimi.it</u>

CREDITS: This presentation template was created by Slidesgo, including icons by Flaticon, infographics & images by Freepik and illustrations by Stories



FURTHER DETAILS



LUTs utilization



RF-RISA¹ re-implementation on Ultra96-v2



The total number of LUTs on the Ultra96-v2 is 70560

In each configuration, an additional 10% of LUTs is consumed due to the experimental setup

YoseUe on Ultra96-v2

Tables refer to an implementation with samples composed of 5 attributes (ap_fixed<32,16>)

¹ S. Zhao, S. Chen, H. Yang, F. Wang, and Z. Wei, "Rf-risa: A novel flexible random forest accelerator based on fpga," Journal of Parallel and Distributed Computing, vol. 157, pp. 220–232, 2021



FFs utilization



RF-RISA¹ re-implementation on Ultra96-v2



The total number of FFs on the Ultra96-v2 is 141120

In each configuration, an additional 7.7% of FFs is consumed due to the experimental setup

YoseUe on Ultra96-v2

Tables refer to an implementation with samples composed of 5 attributes (ap_fixed<32,16>)

¹ S. Zhao, S. Chen, H. Yang, F. Wang, and Z. Wei, "Rf-risa: A novel flexible random forest accelerator based on fpga," Journal of Parallel and Distributed Computing, vol. 157, pp. 220–232, 2021



BRAMs utilization



YoseUe on Ultra96-v2

RF-RISA¹ re-implementation on Ultra96-v2



The total number of BRAMs on the Ultra96-v2 is 216

Tables refer to an implementation with samples composed of 5 attributes (ap_fixed<32,16>)

¹ S. Zhao, S. Chen, H. Yang, F. Wang, and Z. Wei, "Rf-risa: A novel flexible random forest accelerator based on fpga," Journal of Parallel and Distributed Computing, vol. 157, pp. 220–232, 2021



PE's resource utilization

How does the number of attributes of the samples and their dimension influences LUTs and FFs utilization?

LUTs



FFs





The weights are obtained through the training of a small neural network















