

# Machine Learning on FPGAs for Real-Time Processing for the ATLAS Liquid Argon Calorimeter

Lauri Laatu on behalf of the ATLAS Liquid Argon Calorimeter Group

12.01.2023

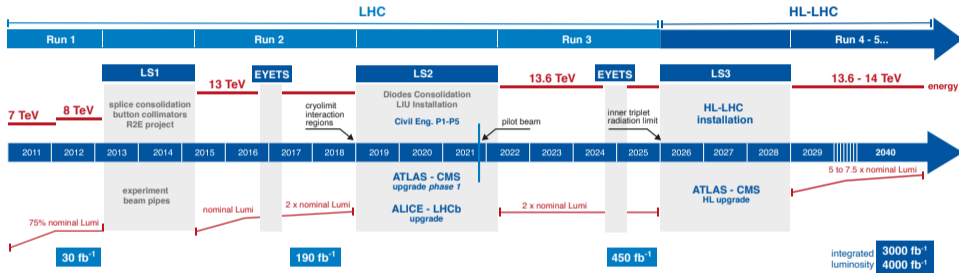


# Content

1. Background
2. Network Architectures
3. Network Performance
4. FPGA Implementation
5. Conclusion

# The Phase-II Upgrade of the LHC

## Upgrade of the ATLAS experiment

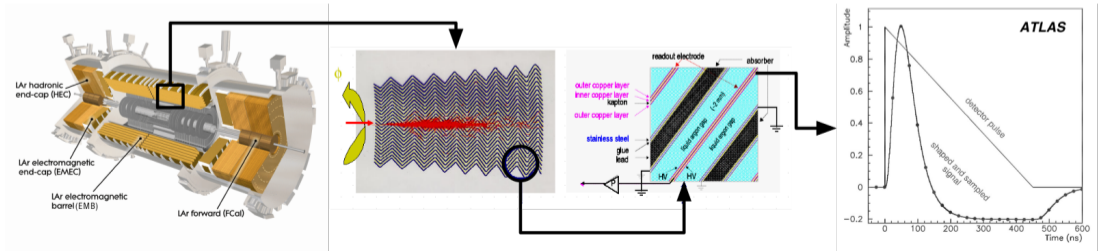


- The High Luminosity LHC (HL-LHC) is an important milestone for particle physics
  - To increase the luminosity to study rare processes
  - To increase the collision rate to up to 200 simultaneous p-p collisions (pileup) per bunch crossing
- The detectors will be upgraded to cope with the high collision rate at the HL-LHC
  - In particular the ATLAS calorimeter readout electronics will be completely replaced

# ATLAS Liquid Argon Calorimeter

## Energy reconstruction in the LAr calorimeter

- The Liquid Argon Calorimeter (LAr) mainly measures the energy deposited by electromagnetically interacting particles
  - Consisting of  $\approx 182\,000$  calorimeter cells
- Passing particles ionize the material
  - Bipolar pulse shape with total length of up to 750 ns (30 BCs)
  - Pulse is sampled and digitized at 40MHz
- Energy reconstruction is done in real-time and used in triggering decision
  - Using the digitized samples from the pulse



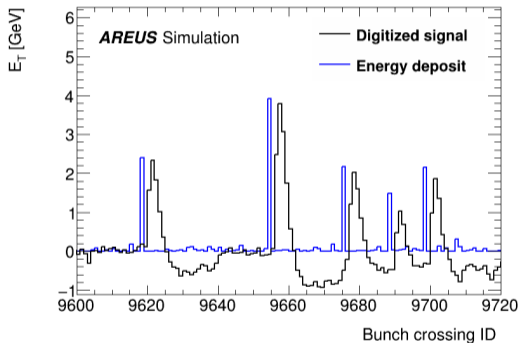
# Energy Reconstruction

Energy reconstruction in the LAr calorimeter

- Current energy reconstruction uses the Optimal Filtering Algorithm with maximum finder (OFMax)

$$E(t) = \sum_{i=1}^5 a_i \cdot s_i$$

- $a_i$  - Predefined coefficients to fit the pulse
- $s_i$  - Sampled signal
- Distorted pulses result in significantly decreased performance of OFMax





# Table of Contents

1. Background

2. Network Architectures

3. Network Performance

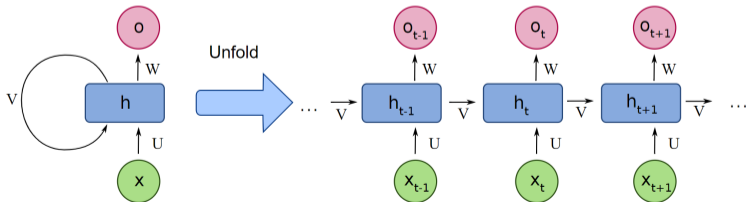
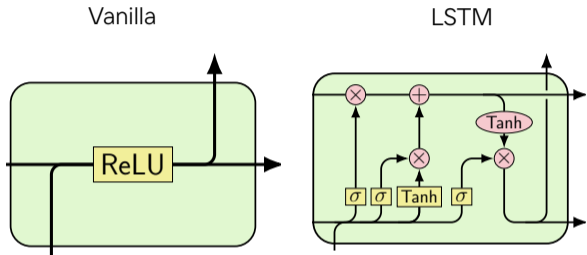
4. FPGA Implementation

5. Conclusion

# RNN Architecture

Time series processing with Recurrent Neural Networks (RNNs)

- Recurrent Neural Networks (RNNs) are designed to process time series data
- RNNs consist of neural network layers that process by combining new time input with past processed state
- Vanilla RNN is the smallest RNN structure
- Long Short-Term Memory (LSTM) network for efficiently handling past information

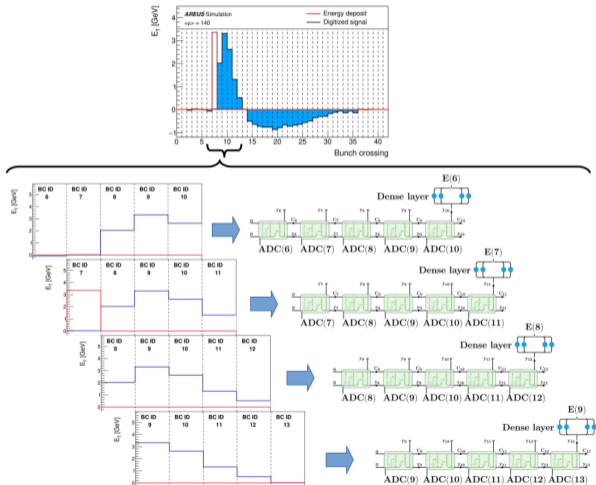




# RNNs for Energy Reconstruction

Using many-to-one and many-to-many networks for energy reconstruction

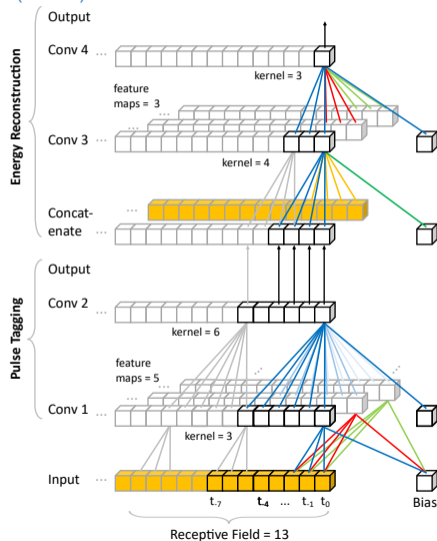
- Use digitized samples as inputs for the recurrent network
- Sliding window
  - Full sequence split into overlapping subsequences with a sliding window
  - One energy prediction per subsequence
  - Four samples in the peak, one in the past
  - Possible for Vanilla RNN and LSTM
- Single cell
  - Use the LSTM cell to process all digitized samples in one continuous chain instead of a sliding window
  - Full history of events available
  - Possible only for LSTM



# CNN Architecture

## Time series processing with Convolutional Neural Networks (CNNs)

- 1D convolutional network for time series regression
- Pulse tagging layers
  - Two layers to classify pulses above 240MeV
- Energy reconstruction
  - Add on top another layer for energy reconstruction
  - Conv3: 5 samples in the peak, 23 in the past with 3 total layers
  - Conv4: 5 samples in the peak, 8 in the past with 4 total layers



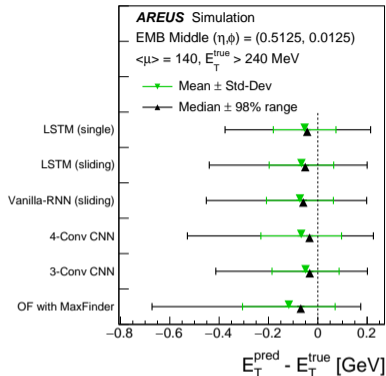
## Table of Contents

1. Background
2. Network Architectures
3. Network Performance
4. FPGA Implementation
5. Conclusion

# NN Performance

## Resolution and network size

- Overall better energy resolution than OFMax
  - Smaller tails and mean closer to zero
- Best performance with LSTM
  - Too large to fit on the FPGA
- CNNs and Vanilla RNN perform well with fewer parameters

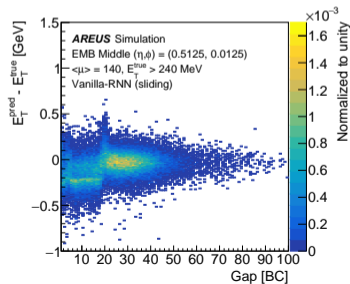
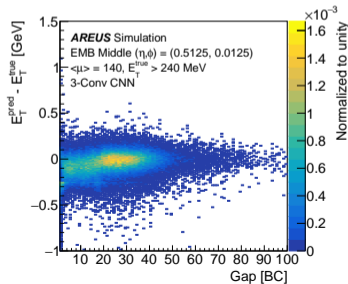
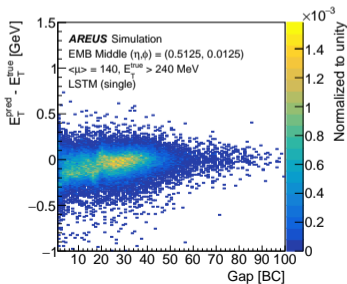
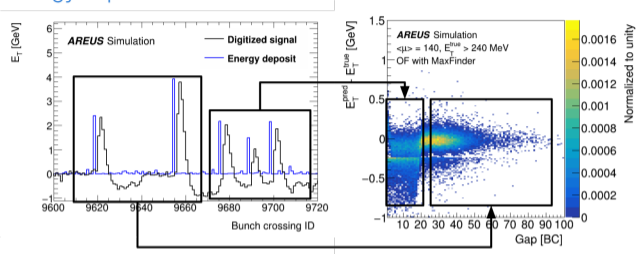


Algorithm	LSTM (single)	LSTM (sliding)	Vanilla (sliding)	CNN (3-conv)	CNN (4-conv)	Optimal filtering
Number of parameters	491	491	89	94	88	5
MAC units	480	2360	368	87	78	5

# NN Performance

Resolution as a function of gap to previous energy deposit in BCs

- Clear performance decrease with OFMax at low gap
- All NNs perform better with overlapping events
  - More past samples allows for better correction of overlapping events



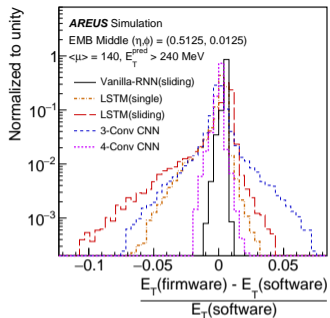
## Table of Contents

1. Background
2. Network Architectures
3. Network Performance
4. FPGA Implementation
5. Conclusion

# Firmware Implementation

## Running in the FPGA

- Implementation on Stratix 10 FPGA
- CNNs implemented in VHDL
- RNNs implemented in HLS
- O(1%) resolution
  - Fixed-point arithmetic
  - Look-up tables for activation functions
- Implementations are close to requirements in term of resource usage and latency, demonstrating feasibility but additional tuning is required



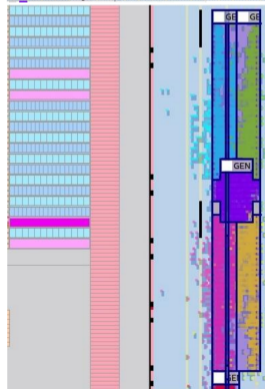
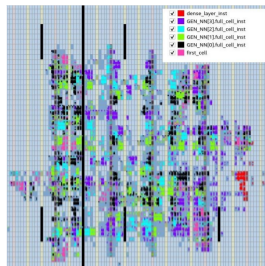
	3-Conv CNN	4-Conv CNN	Vanilla RNN
Multiplicity	6	6	15
Frequency			
$F_{\text{max}}$ [MHz]	344	334	640
Latency			
$\text{clk}_{\text{core}}$ cycles	81	62	120
Initiation interval			
$\text{clk}_{\text{core}}$ cycles	1	1	1
Max. Channels	390	352	576
Resource Usage			
#DSPs	46	42	152
	0.8%	0.7%	2.6%
#ALMs	14235	15627	5782
	1.5%	1.7%	0.6%

# RNN Implementation in VHDL

## Running in the FPGA

- Further optimisation of RNN implementation in VHDL for better tuning of the placement
- Incremental compilation with forced placement
  - Tackle timing violations
- Multiplexing to compute several networks simultaneously

	N networks x multiplexing	ALM	DSP	FMax	latency
<b>target</b>	<b>384 channels</b>	<b>30%*</b>	<b>70%*</b>	-	<b>125 ns</b>
HLS (no multiplexing)	384x1	226%	529%	-	322 ns
HLS optimized	37x10	23%	100%	414 MHz	302 ns
VHDL optimized	28x14	18%	66%	561 MHz	121 ns

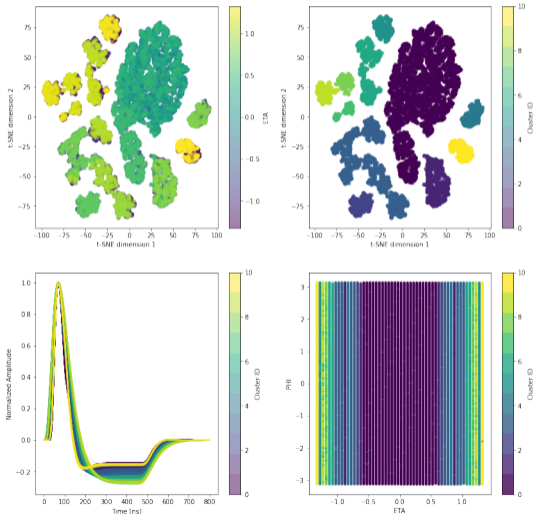




# Reconstruction for Full Detector

## Pulse Clustering

- Pulse shape differs in the detector
  - Reduced performance with differing pulse shapes
  - One NN training will not perform well for the full detector, nor is 182k NNs feasible
  - Need to reduce the number of NNs trained while maintaining accuracy
- Clustering method used to group detector regions
  - t-SNE from calibration pulses to acquire clustering
  - DBSCAN to automatically classify cluster
  - Separation correlates with  $\eta$  according to pulse shape differences



# Pulse Clustering

## Reconstruction in different regions

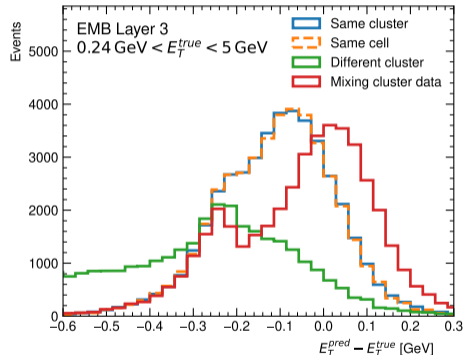
### Evaluate inside same cluster

- Train with one cell, test with another
- Same performance as with training and testing with the same cell

### Large performance drop when training with one cluster and testing with another

### Train with mixed data from all clusters, test with single cluster

- Mixing data across clusters slightly restores performance



## Table of Contents

1. Background
2. Network Architectures
3. Network Performance
4. FPGA Implementation
5. Conclusion

# Conclusion

## Energy reconstruction using recurrent neural networks

- Energy reconstruction with CNNs and RNNs overperforms legacy algorithms in Phase-II conditions
  - Better energy resolution overall
  - Better recovery of energy resolution with overlapping signals
- Implemented and validated in firmware and the implementations mostly fulfill the LAr real-time processing requirements
  - Testing on DevKits started and is showing good results
- Next step is to quantify the effect on object (electrons, photons) reconstruction and physics performance
- Paper published available [▶ Here](#)

