

Optimisation of embedded neural networks for the energy reconstruction of the Liquid Argon Calorimeter cells of ATLAS

-

CPPM Seminar 3rd year PhD
27.10.2025

Raphaël Bertrand

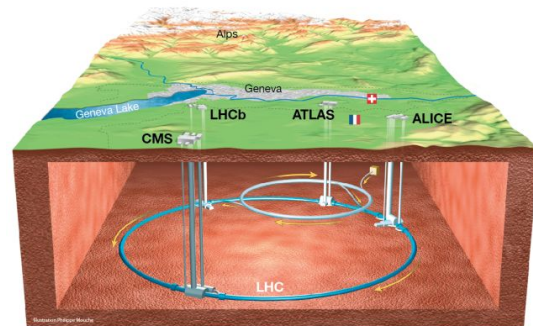


Introduction

Experimental Context

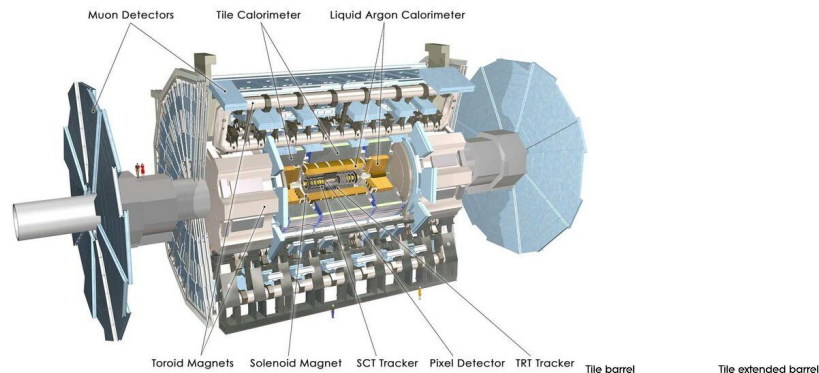
- Large Hadron Collider (LHC)

- Proton-proton collider at 13.6 TeV
- Protons accelerated via superconducting magnets
- Collisions at 40 MHz



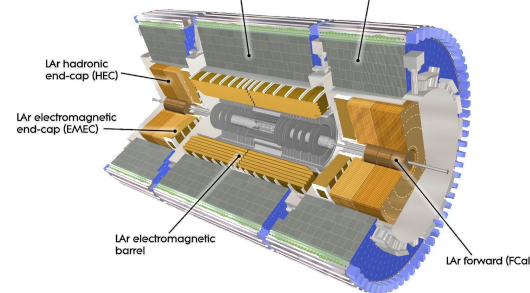
- ATLAS detector

- General-purpose experiment
- Very high data rate
 - On-the-fly event selection required



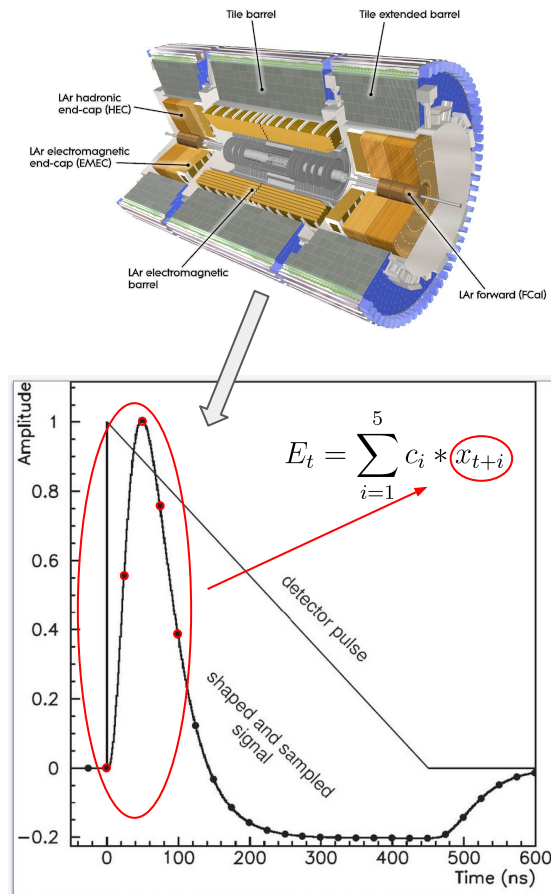
- Liquid Argon (LAr) Calorimeter

- ATLAS sub-detector for energy measurement (e^{\pm}, γ)
- Sampling in active LAr alternating with inactive metal (Cu, Pb, W)
 - Accordion shaper absorbers for EMB and EMEC
 - Ionization signal from particle interactions

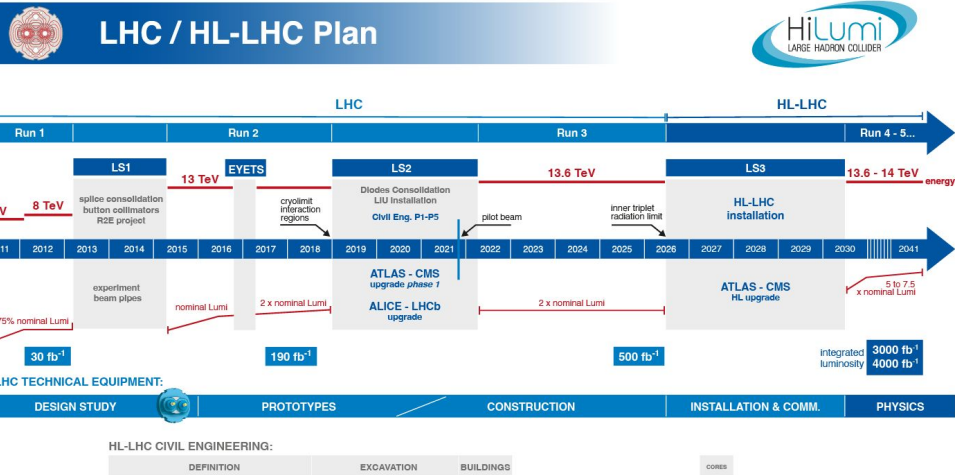


Signal processing and energy reconstruction

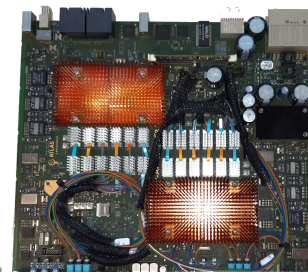
- **Electronic signal produced**
 - **Amplitude** \propto **true deposited energy** (E^{true})
 - **Spans** **~625 ns** (25 proton-proton **Bunch Crossings**)
 - **Shaped, sampled and digitised at 40 MHz**
- **Energy reconstruction with optimal filtering (OF) algorithm**
 - **Weighted sum** of samples around the pulse peak
 - **Max finder/Timing cut** to select the correct BC
- **Reconstruction algorithm requirements :**
 - **Online** computation (per BC)
 - **Max latency** : **~125 ns** (used in trigger system)
 - **Fit in FPGAs** : **O(500)** Multiply-Accumulate operations (**MAC units**)
 - 5 MAC units required to implement OF
 - **384 channels per FPGA** (many algorithm instances needed)



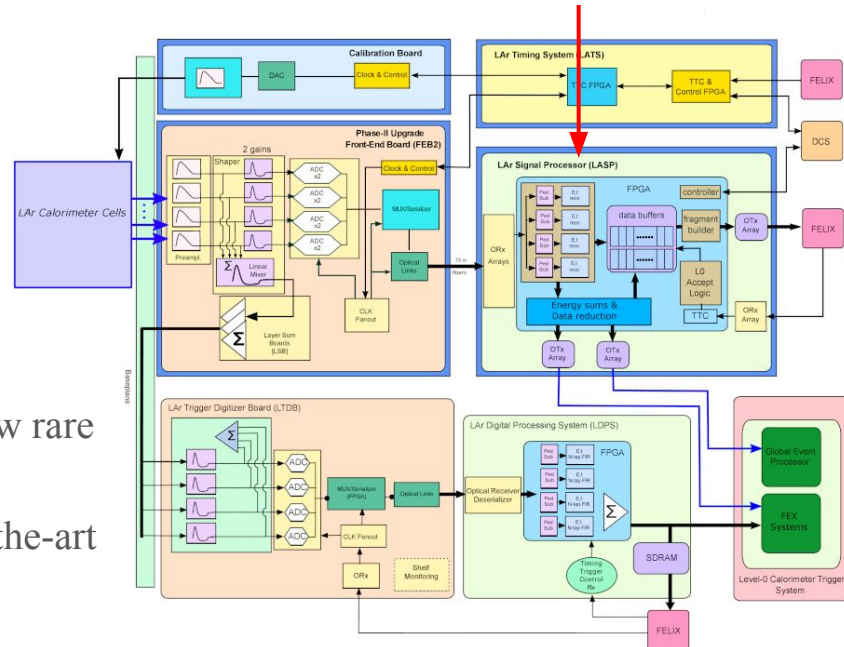
HL-LHC schedule and ATLAS Phase-II Upgrade



- HL-LHC \Rightarrow Luminosity \nearrow
- HL-LHC is needed to study Higgs properties and detect new rare processes
- Off-detector readout board (LASP) will carry two state-of-the-art FPGAs for energy computation
 - Opportunity to embark more complex algorithms

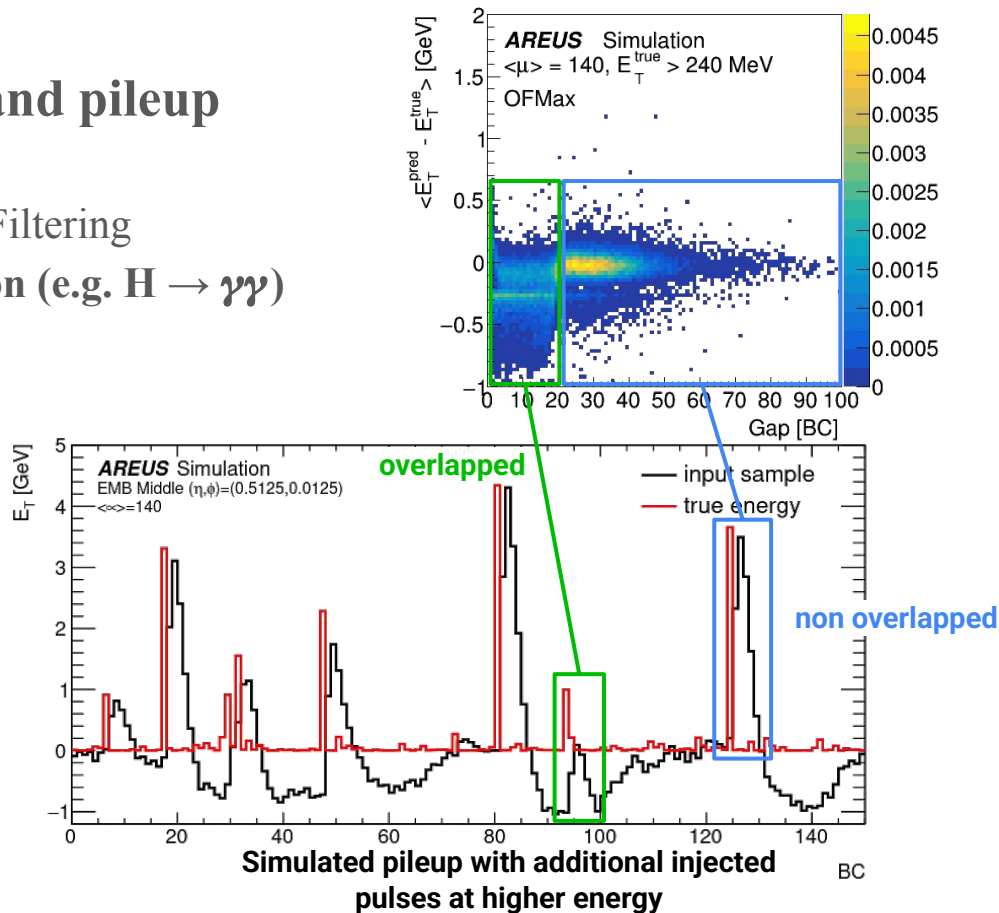
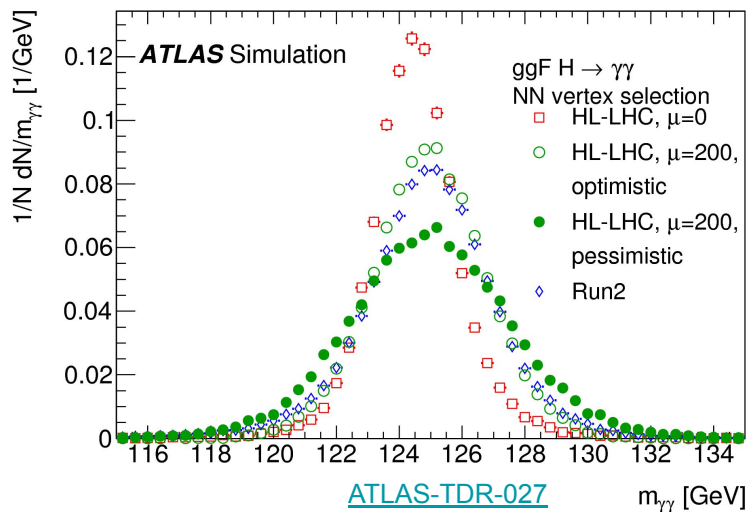


LASP board demonstrator



Impact of high luminosity

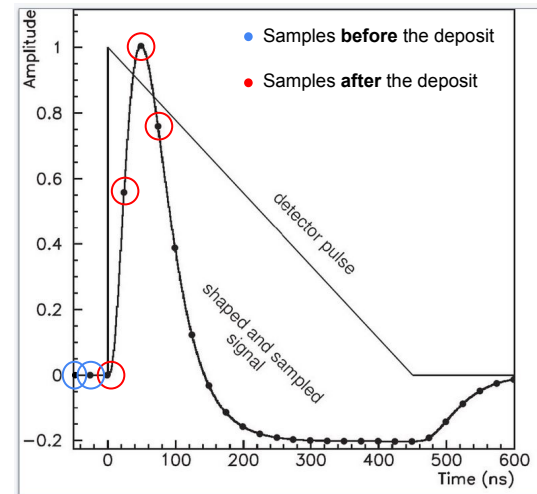
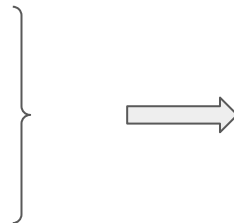
- HL-LHC \Rightarrow Increased luminosity and pileup
 - Increased rates of overlapping pulses
 - \hookrightarrow Degraded performance of Optimal Filtering
 - Significant impact on energy resolution (e.g. $H \rightarrow \gamma\gamma$)



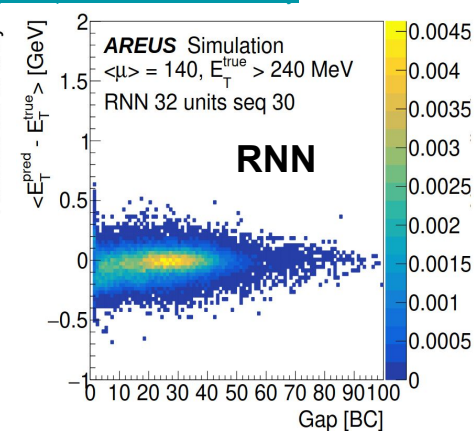
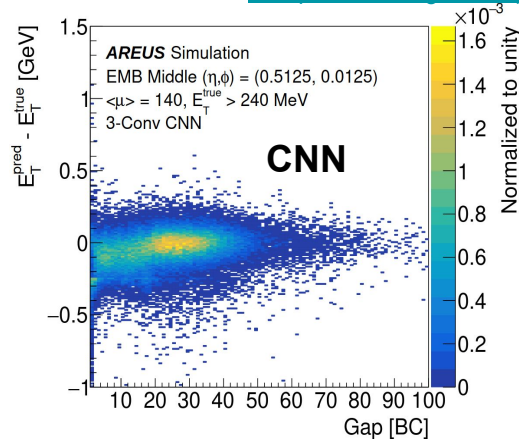
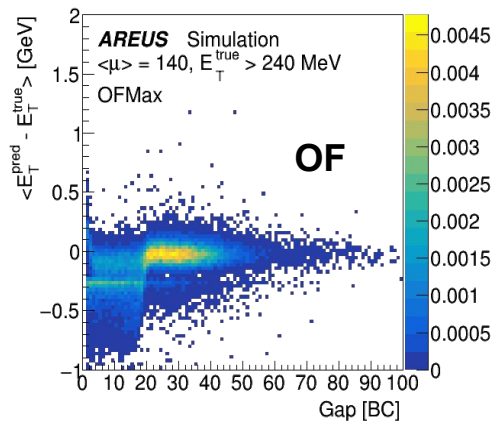
Neural network approaches as energy reconstruction algorithms

Neural networks

- Exploit samples before the energy deposit to **correct overlapping pulses**
- Several architectures tested : RNN, RNN+Dense, CNN, Dense
- Samples from **before** and **after** the energy deposit are used :
 - **After** the energy deposit (similar to OF inputs)
 - Capture the pulse amplitude
 - **Before** the energy deposit (additional inputs)
 - Correct for pulse distortions from previous deposits
- Preliminary studies done with high rate of pulse overlap
 - **Neural networks can correct for overlapping pulses**
 - The correction is **dependent on the size** of network

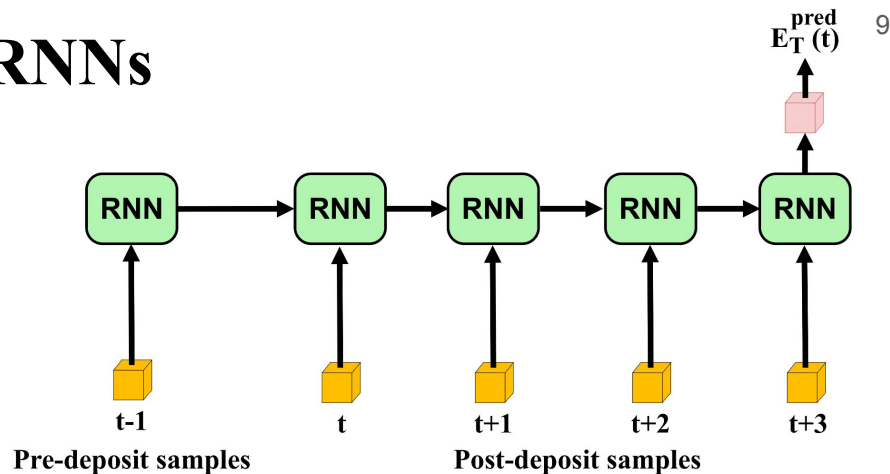


[Comput.Softw.Big Sci. 5 \(2021\). s41781-021-00066-y](#)

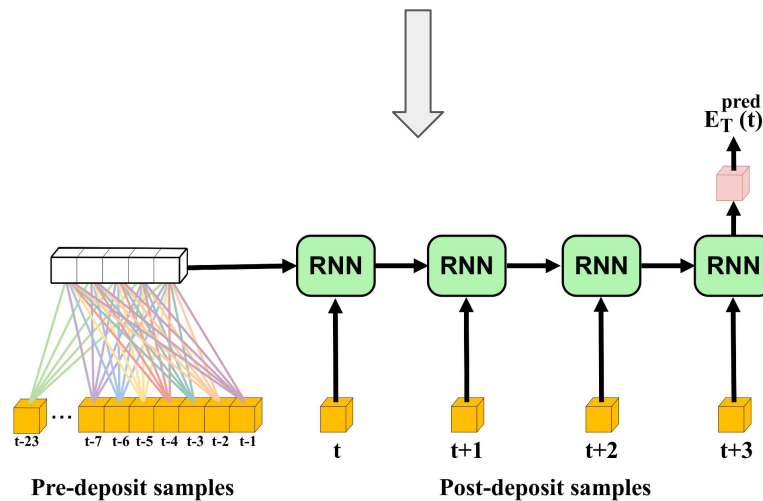


Neural network architecture - RNNs

- Architecture RNN
 - o Multiple RNN cells sharing the same parameters
 - One cell per sample
 - o One dense layer to concatenate output from last cell
 - Return predicted energy
- The RNN architecture assigns **equal importance to all samples**
- Start computations at first sample
 - o **Good latency**
- FPGA implementability :
 - o Number of **Multiply-Accumulate Operations (MAC units)**



- Architecture **Dense+RNN** → **additional dense layer** before RNN cells
 - o Computation for **pre-deposit samples**
 - o **Less MAC units** for a dense layer

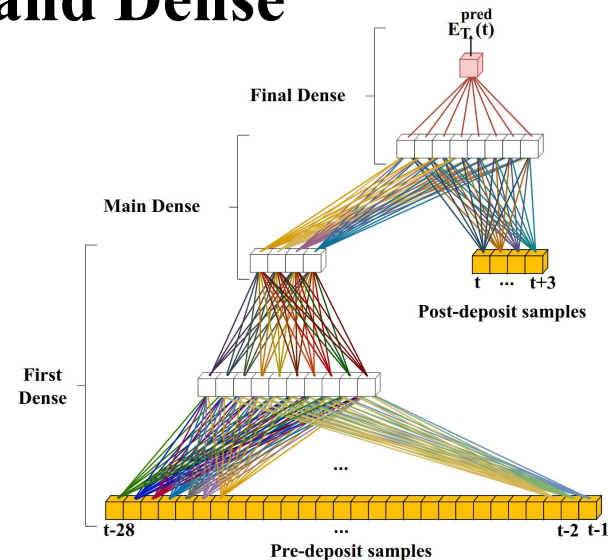
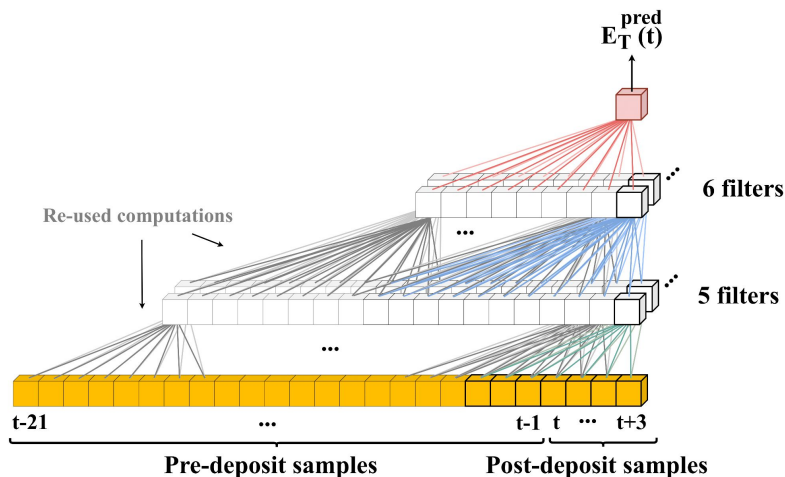


$$\Rightarrow \text{MAC units}_{\text{RNN layer}} \propto \text{units}^2 \times \text{nb of samples}$$

$$\Rightarrow \text{MAC units}_{\text{Dense layer}} \propto \text{units} \times \text{nb of samples}$$

Neural network architecture - CNN and Dense

- Architecture CNN \Rightarrow Three CNN layers
- Uncertain latency compliance
- FPGA implementability :
 - Re-used computations



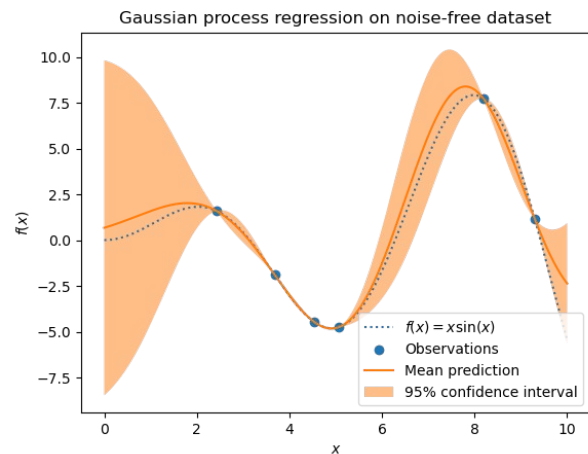
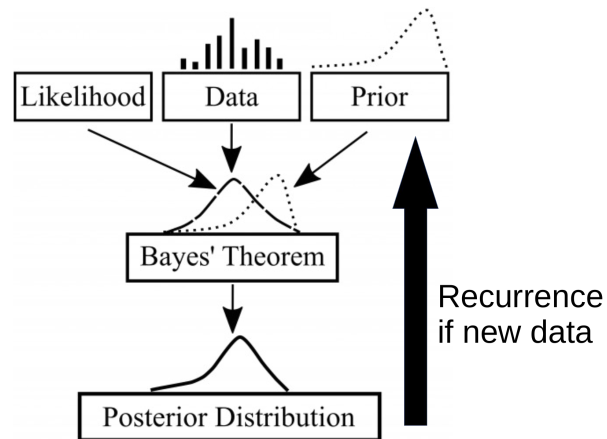
- Architecture **Dense**
 - One block of dense for **aggregation of pre-deposit samples**
 - One dense layer to add a layer of computations before concatenation
 - One dense layer to concatenate output from last cell
 - **Return predicted energy**
- **Good latency**
- **Smaller number of MAC units** with dense layers

Neural networks hyperparameters tuning

bayesian optimization

Bayesian optimisation

- Goal : **Find the best parameters** to maximize/minimize a **performance function** while **evaluating the function as few times as possible**
- **Initialization** with several random points
- **Iterations** to find the best parameters space
 - **Interpolation** between points
 - Based on a gaussian kernel with associated uncertainty
 - **Acquisition function** to determine where to evaluate next
 - Balance between **exploration** and **exploitation**
 - **Evaluation** of the performance function **at the chosen point**



Bayesian optimisation applied on energy reconstruction

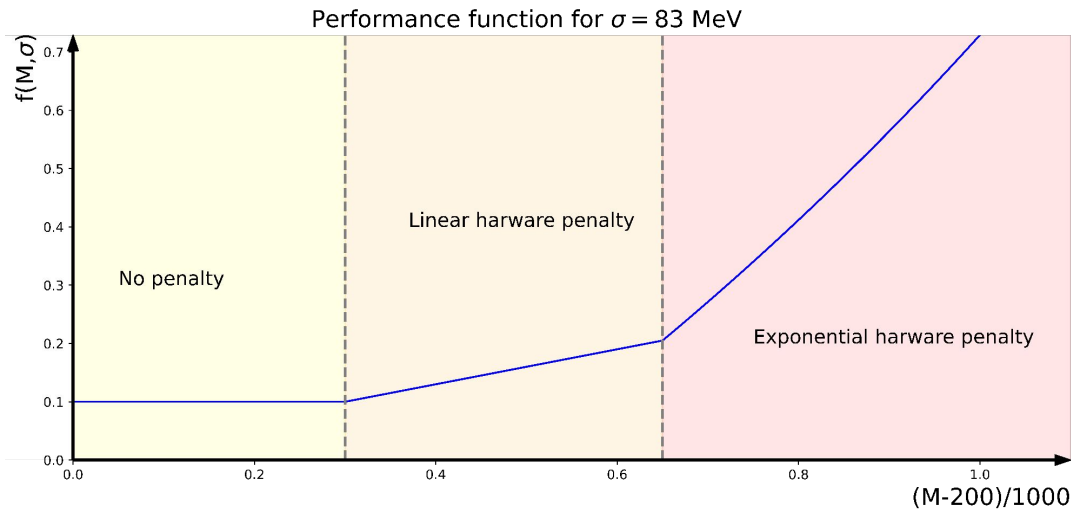
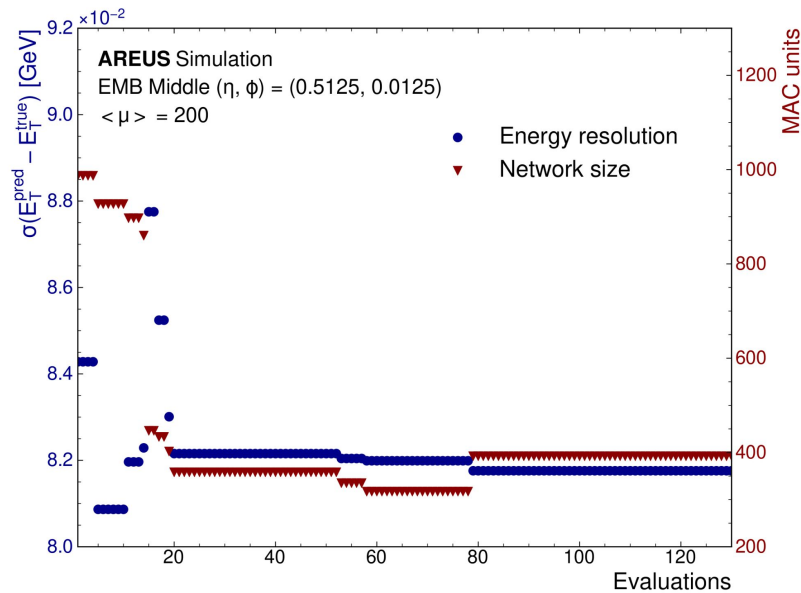
Performance function used for the bayesian optimization :

$$f(M, \sigma) = \frac{\sigma - 70}{130} \text{ for } M \leq 500$$

- Optimisation on both performance and hardware to fit in FPGAs
 - o **Energy resolution** (σ [MeV])
 - o **Number of MAC units** (M)
- Hyperparameters to be tuned (e.g. for the Dense architecture) :
 - o **Number of samples** (before the energy deposit)
 - o **Number of units** for the intermediate layers

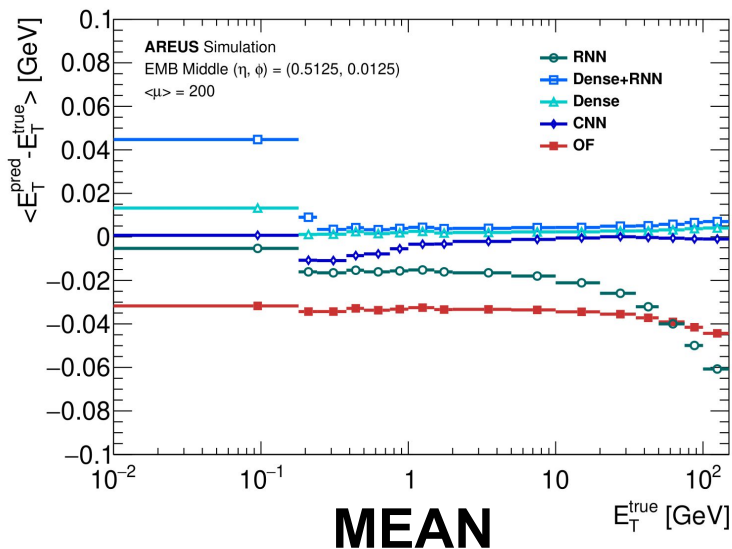
$$f(M, \sigma) = f(500, \sigma) + a * \frac{M - 500}{1000} \text{ for } M \in] 500 ; 850]$$

$$f(M, \sigma) = f(850, \sigma) + b * e^{\frac{M - 850}{1000}} - 1 \text{ for } M > 850$$



Energy scale and resolution as function of true energy

- **Better energy scale of $E_T^{\text{pred}} - E_T^{\text{true}}$ for Dense+RNN, Dense and CNN architectures compared to OF**
 - RNN energy scale falling with higher E_T^{true}
 - OF energy scale not centered around 0
- **Better energy resolution of $E_T^{\text{pred}} - E_T^{\text{true}}$ for Dense+RNN, Dense and CNN architectures compared to OF**
 - Visible for the whole energy range



MAC units

368

240

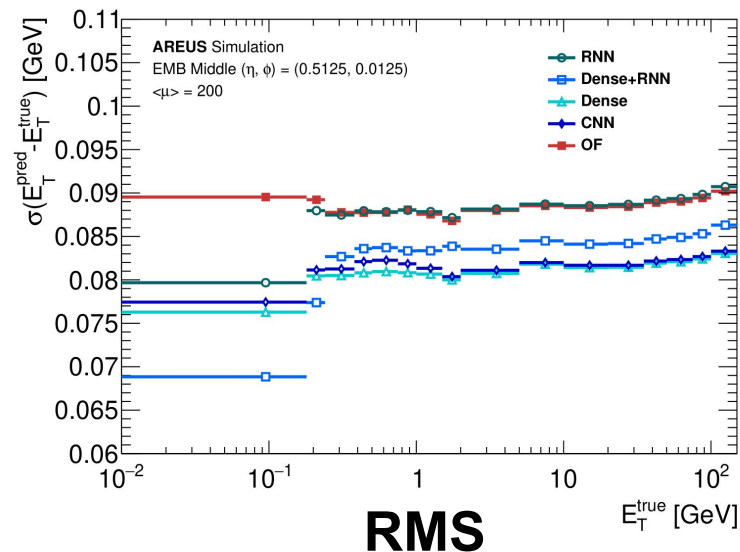
392

419

5

[Pre-print](#)
(submitted to
EPJC)

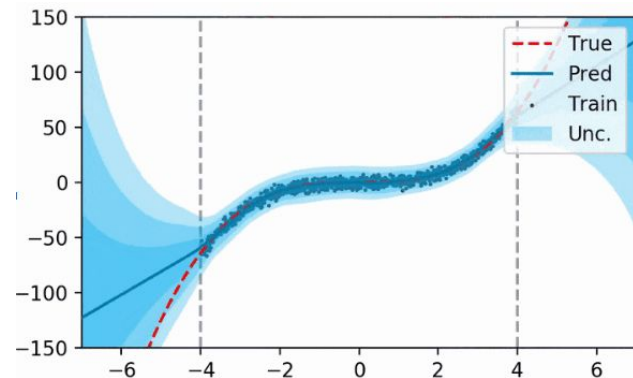
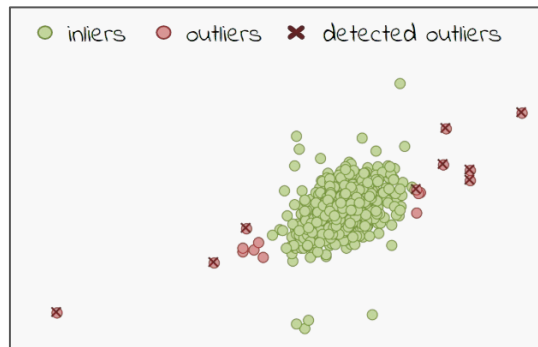
[Public code](#)



Uncertainty prediction using neural network with deep evidential regression

Deep evidential regression (DER)

- NNs are trained to minimize their prediction errors
 - o Unknown accuracy of the model for individual prediction
 - o It would be interesting to **know when the model is more likely to fail (or the opposite)**
- **Model the energy prediction as a distribution**
 - o Mean of the distribution → **energy prediction**
 - o Standard deviation of the distribution → **uncertainty**
- Differentiate uncertainties :
 - o **Epistemic**
 - Lack of knowledge, model uncertainty
 - Can be reduced
 - o **Aleatoric**
 - Inherent to data
 - Cannot be reduced



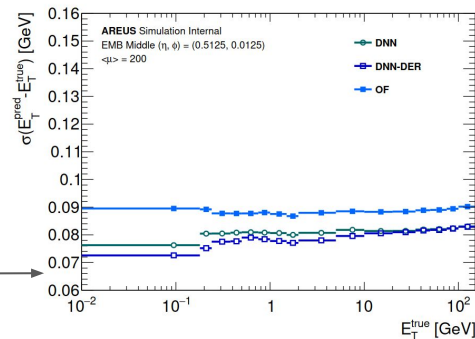
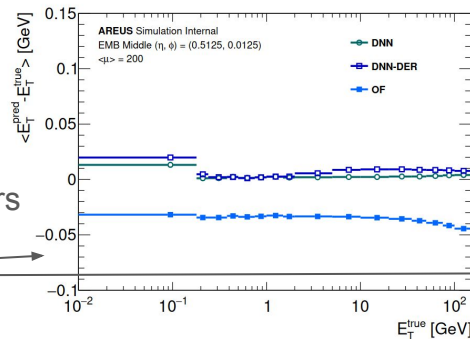
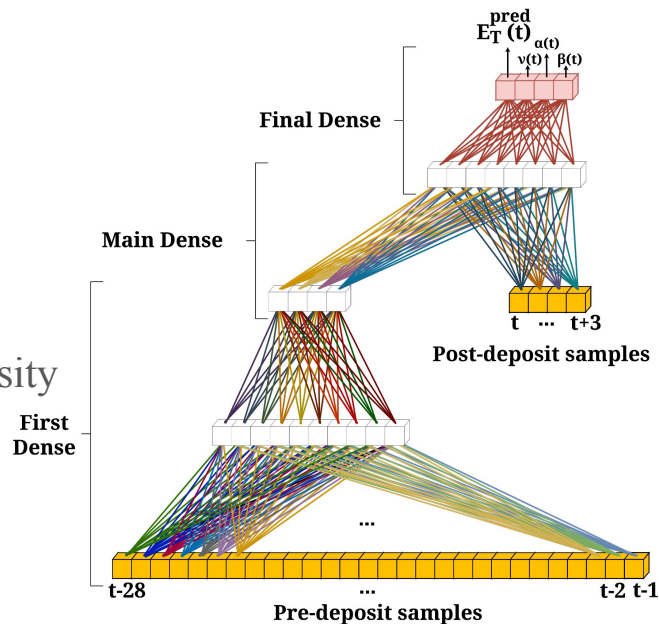
Deep evidential regression (DER)

- Normale-Inverse Gamma distribution to describe mean and uncertainty
 - **4 parameters** (γ, v, α, β) rather than one
 - Uncertainty computation
- **DER applied to LAr cells energy reconstruction**
 - Would allow to take into account instantaneous luminosity changes or bunch train structure
- Adapted to the Dense architecture
 - **Use of bias** compared to the other architectures
 - **Still possible to implement in FPGA**
 - **416 MAC units**
- Training loss function : ~~MSE~~
 - Likelihood + Regularisation

fit the distribution

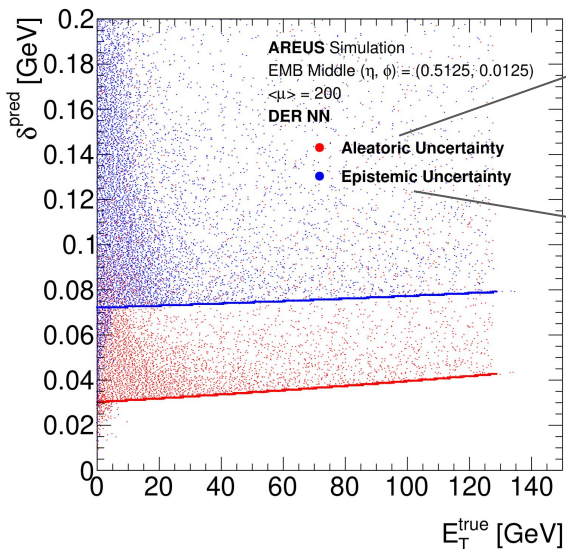
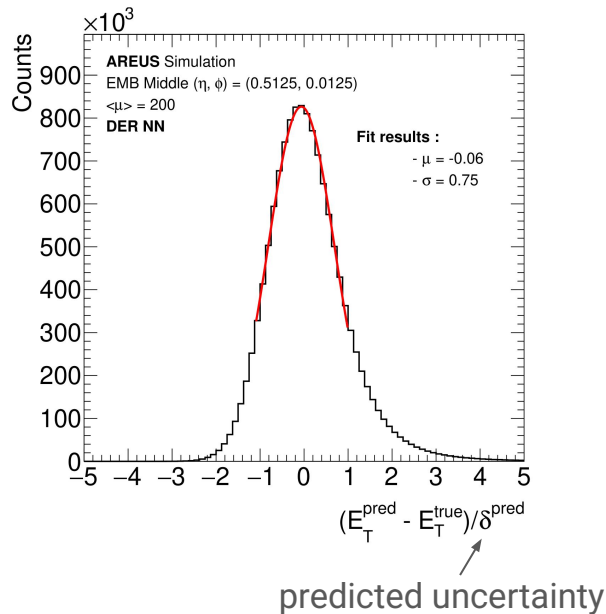
increase uncertainty on large errors

- **Similar performance**



Uncertainty prediction

- Overall good pull distribution
 - Estimated uncertainty comparable to $E_T^{\text{pred}} - E_T^{\text{true}}$
 - **Slightly biased**
 - Right tails
 - Uncertainty overestimated by 25%



Data uncertainty

$$\mathbb{E}[\sigma^2] = \frac{\beta}{\alpha - 1}$$

Model uncertainty

$$\text{Var}[E_T^{\text{pred}}] = \frac{\beta}{\nu(\alpha - 1)}$$

- Epistemic and aleatoric uncertainties are constant
- **Epistemic uncertainty is dominant** in the total uncertainty

Conclusion

- **Online energy reconstruction for LAr cells performed using neural networks**
 - Pre-print on arxiv : [Optimised neural networks for online processing of ATLAS calorimeter data on FPGAs](#)
- **Four neural network architectures were tested and optimized**
 - **CNN, RNN, RNN+Dense and Dense**
- **Hyperparameter tuning performed using bayesian optimization**
 - **Balance between performance and size of the network** to fit in FPGAs
 - **NNs outperform OF**
- **Uncertainty on energy prediction using deep evidential regression**
 - **Accurate uncertainty prediction**
 - **Possible to implement in FPGAs**

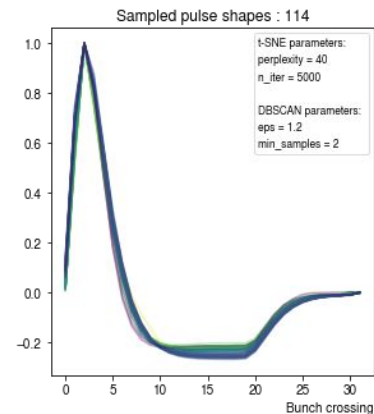
BACKUP

Bayesian optimisation code

- **Parallel iterations**
 - o **Multiple models (with different parameters) trained simultaneously**
 - **Min distance** defined to avoid to train models on the same peak/minimum
 - Fixed
 - Decreasing
 - At lowest
 - **Exploration parameter** (to favour high-uncertainty regions)
 - Fixed
 - Fixed
 - Decreasing
- **Noise considerations**
 - o Choice 1 : **multiple models trained at the same parameters, only the best one is considered**
 - o Choice 2 : multiple models trained at the same parameters, all used to train the Gaussian Process Regressor
- **Integers considerations :**
 - o **Only integer parameters are proposed**
 - o Two iterations can't be on the same set of parameters
- **Kernel choice : Matern 5/2**
 - o Twice differentiable
 - o **More realistic** than usual RBF/Gaussian kernel (infinitely differentiable)

Cells grouping

- Training one network per cell is not feasible
- LAr calorimeter cells grouped
 - One neural network \Leftrightarrow several cells
 - 182,468 cells \rightarrow **few hundreds groups**
- Grouping performed using techniques known as t-sne and dbscan
 - Groups are trained layer per layer
 - The grouping recovers the detector symmetry



Single cell training performance
fully recovered

