## Posit-based ML & DNN Acceleration for AI in EPI

The adoption of Machine Learning (ML) and Deep Neural Networks (DNN) to enable Artificial Intelligence (AI) techniques on edge devices in strategic EU markets (e.g. automotive, aerospace, industry4.0, biomedical, robots), is requiring HW and SW acceleration to achieve High Performance Computing with high energy efficiency. To this aim, in collaboration between University of Pisa and Kalray, the use of Posits in the EPI project aims at exploiting a new type of arithmetic that allows the reduction of power consumption and circuit complexity for data processing and storage. The fundamental application of the Posit research topic is the acceleration of Deep Neural Network and the acceleration of autonomous driving applications. Posits are an effective alternative to classic integer and floating-point (IEEE 754) arithmetic to preserve the accuracy of floats but with just half of the bits. This means that it is possible to use this novel format for data compression and for computation, i.e. doubling the bandwidth and/or halving the memory complexity, without losing accuracy.
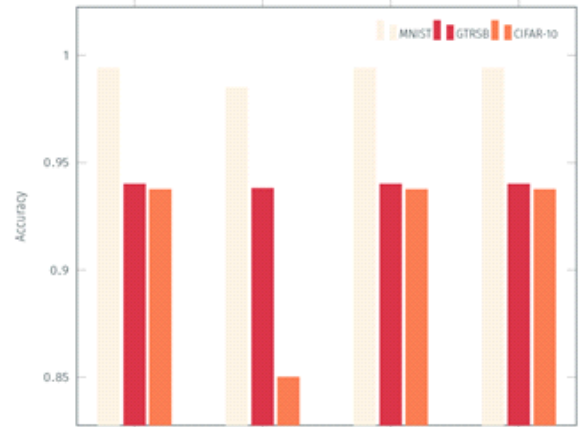


*Fig. 1. Accuracy for Posit 16, 12, 8 vs. FP32 for TinyDNN with MNIST, GTRSB and CIFAR data sets*

### Posit C++ library developed by UniPisa provides:

- Custom definition of any Posit type
- Different back-ends for accelerated emulation (floating point, fixed point, tabulation)
- Support for FPGA synthesis to provide HW acceleration of a Posit Processing Unit
- Support for different vector processors used in EPI: ARM v8.2 Scalable Vector Extension; RISC-V "V" 0.8 extension

### Posits in Kalray MPPA (Massively Parallel Processor Array):

- Posit8 numbers identified as an effective compressed representation for the Float32 parameters: the results of rounding can be restricted to Posit8,0 or Posit8,1 numbers, with the benefit of reducing by half the memory capacity and bandwidth
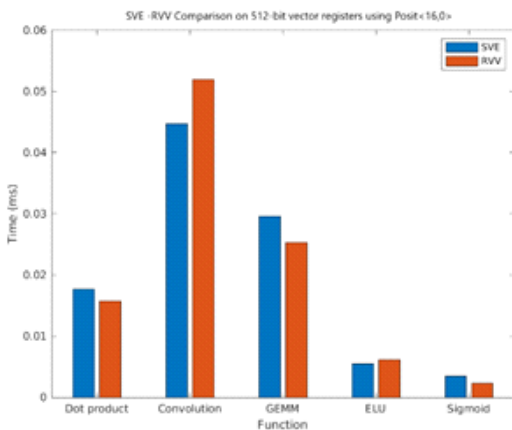


*Fig. 2. Execution time of Posit C++ library using ARM SVE and RISC-V V backends with 512-b vector registers for several AI functions: dot product, convolution, Exponential Linear Unit (ELU), sigmoid, General Matrix-Multiply (GEMM)*

Moreover, Posits can be further squeezed down to a quarter of the size of floats, losing little-to-none in terms of accuracy. In EPI, Posits have been implemented: i) through a SW library (Pisa CppPosit) and the SW framework is completed by the open source tinyDNN Deep Neural Network library extended to support the new cppPosit features; ii) in Kalray MPPA (Massively Parallel Processing Array).

Posit8 numbers have been identified by Kalray as an effective compressed representation for the Float32 network parameters: instead of rounding Float32 parameter values to Float16, the results of rounding can be restricted to Posit8,0 or Posit8,1 numbers, with the primary benefit of reducing by half the memory capacity and bandwidth required by the network parameters. Kalray focuses on the Posit8,0 and Posit8,1 numbers because they are exactly represented as Float16 numbers, and thus can benefit from the exact Float16.32 dot-product operator of the MPPA3 co-processors. This evaluation should lead to the inclusion of new arithmetic instructions to expand Posit8 to Float16 in the MPPA IP delivered to the H2020 European Processor Initiative.