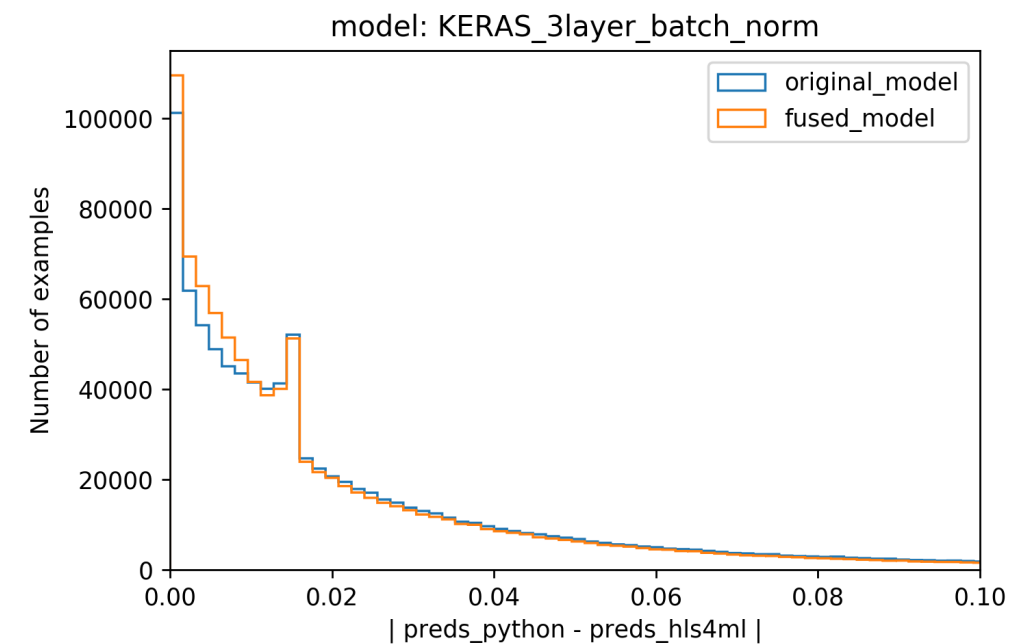


HLS4ML

Vladimir, Sioni, Jennifer
Petr, Adrian, Hasib, Dejan

Completed

- Fixed Convolutional layers issues
 - Implemented support for *channels first* and *channels last* options
- Fusion of dense and batch normalization layers
 - Include BN parameters in the preceding dense layer
 - Remove BN layer
 - $\text{fused_W} = \text{gamma} * \text{W} / \text{sqrt}(\text{variance} + \text{eps})$
 - $\text{fused_b} = \text{gamma} * (\text{b} - \text{mean}) / \text{sqrt}(\text{variance} + \text{eps}) + \text{beta}$
 - Fewer layers on the board = reduced latency
 - (Slightly) increased precision
- Started developing support for multiple backends
 - Vivado, Mentor Catapult, Intel



Latency [clock cycles]	
Original	27
Fused	19

Next Steps

- Fuse convolutional layers with batch normalization
- Develop additional features for HLS4ML
- Compare HLS4ML performance for 2D CNN with Micron