

- *Discuss VPU/TPU/APU/GPU/FPGA/ASIC memory architecture and how it handles matrix sparsity*

- *Show ineffectivity of pruning on hardware without optimisations for sparse matrices*

The explosion of Deep Neural Network applications in recent years has prompted the production of a wave of specialised hardware architectures to improve the efficiency and compute of these kinds of workloads. The mainstay of this form of processing has been until recently been dominated by GPUs.