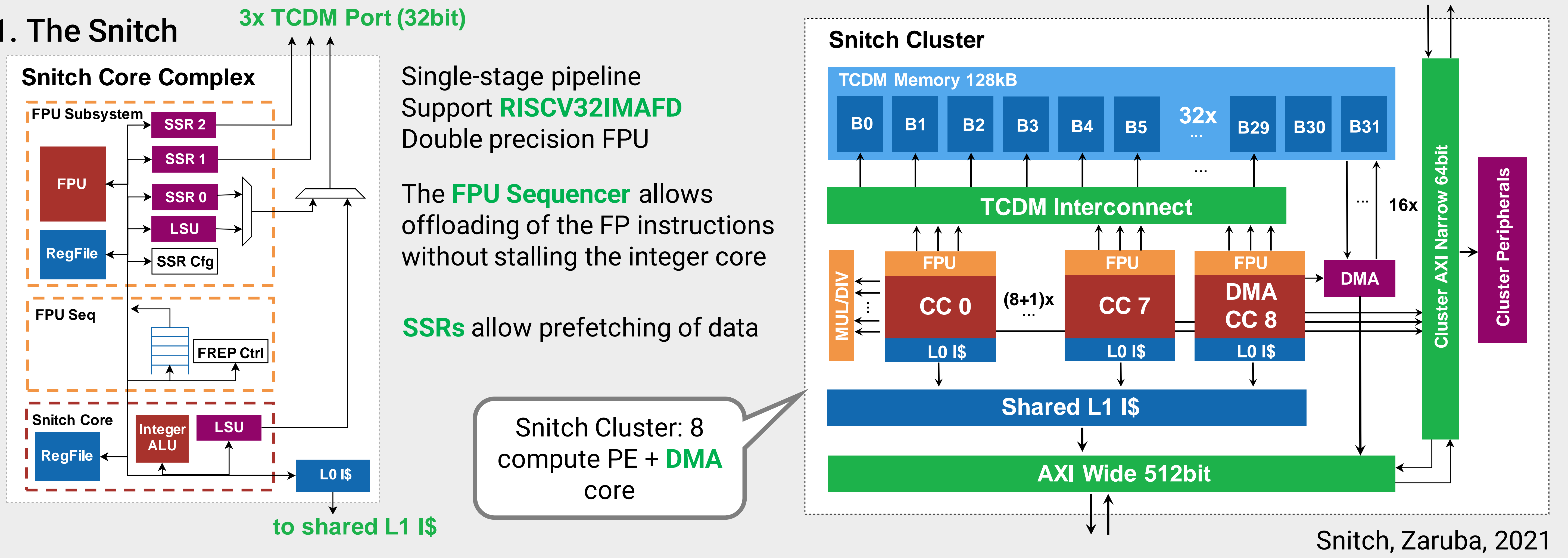


Parallel Sparse Deep Learning Operators on Lightweight RISC-V Processors

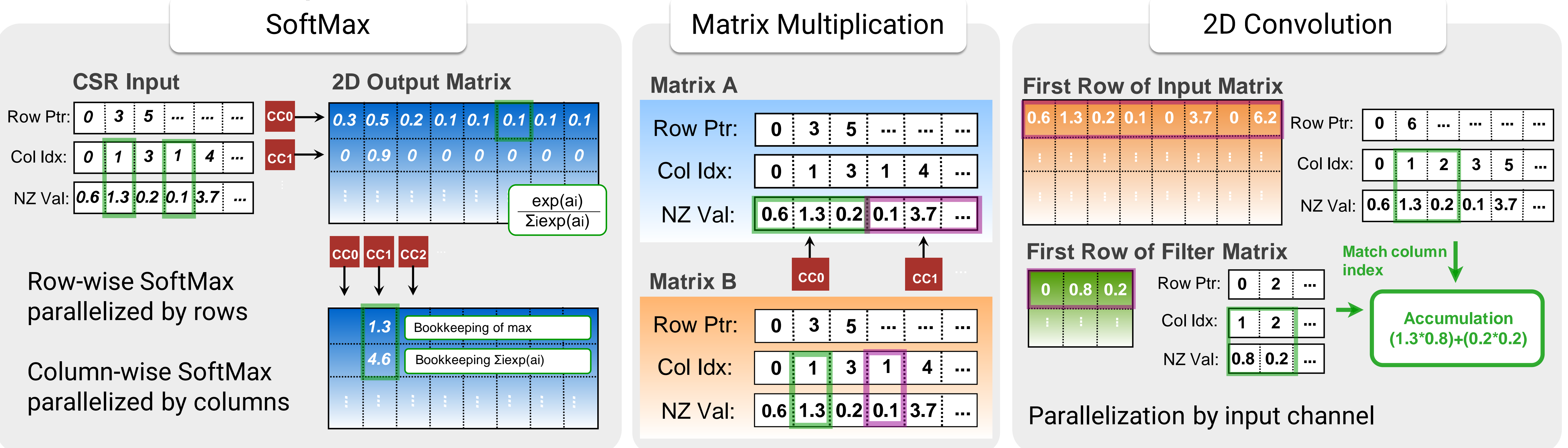
Marco Bertuletti¹, Tim Fischer¹, Yichao Zhang¹, Luca Benini^{1,2}
¹Integrated Systems Laboratory, ETH Zurich; ²Università di Bologna, Italy;



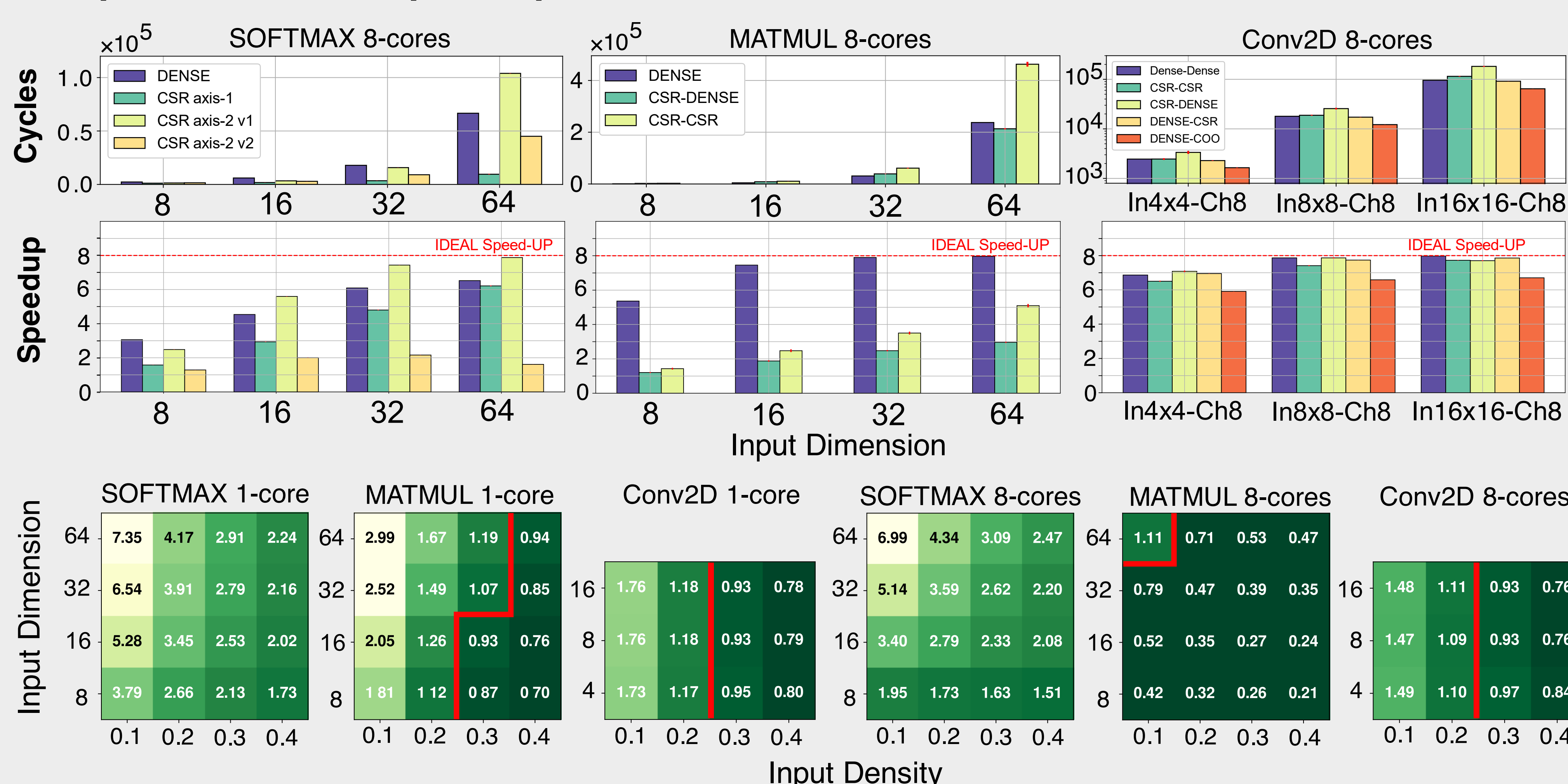
1. The Snitch



2. Parallelization of Sparse kernels



3. Sparse kernels speedup



SW Implementation of Sparse kernels:

Quasi-ideal parallel speedup

CSR vs DENSE SPEEDUP

x7.3 - x3.0 - x1.7 1 core
x7.0 - x1.1 - x1.5 8 cores

When the density of the inputs is high we still benefit from **reduced memory footprint**

Further improvement via prefetching