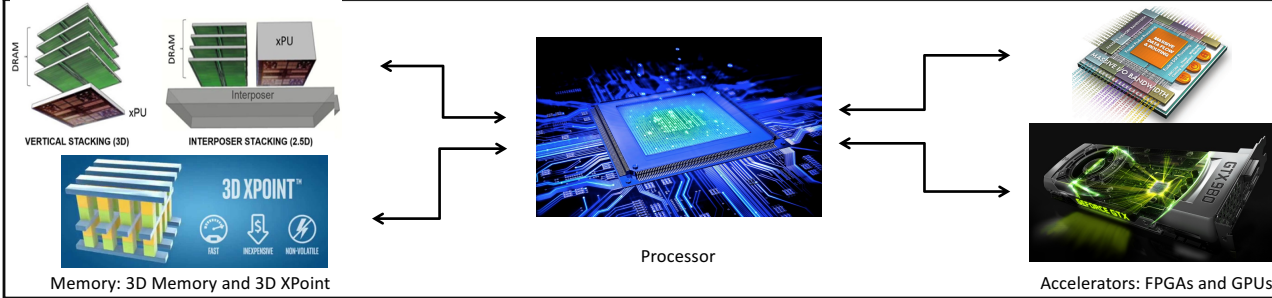


# Exploiting 3D Memory for Energy-Efficient Memory-Driven Computing

PI: Viktor K Prasanna

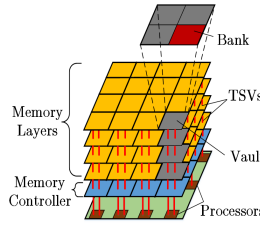
Co-PI: Charalampos Chelmis

## Motivation



## Challenges

- Traditional challenges still exist
- Row activation overhead
- Low page hit rate
- Optimizations specific to 3D Memory
- Structure and organization
- Access pattern of the application
- Data layout in the memory



## Modeling 3D Memory

- Timing Parameters of 3D memory are:
  - $t_{val}$ : accesses to different vaults
  - $t_{layer}$ : accesses to different layers in a vault
  - $t_{bank}$ : accesses to different banks in a layer in a vault
  - $t_{row}$ : accesses to different rows in a bank
  - $t_{col}$ : accesses to different columns in a row

## Signal Processing

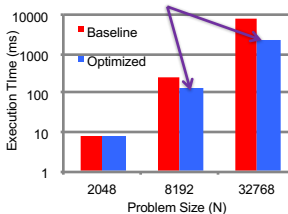
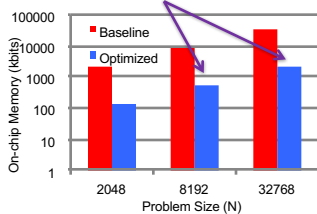
- 2D FFT: Row phase (sequential) and Column phase (strided)
- Trade row activations for on-chip memory
- Exploit parallelism at all levels: Inter-Layer Pipelining and Parallel Vault access
- Exploit large number of banks: Hide  $t_{col}$ ,  $t_{row}$ ,  $t_{bank}$
- Optimized data layout for 2D FFT on 3D memory
  - Streaming data sequential and strided accesses
  - $O(N)$  on-chip memory for  $N \times N$  problem size
  - $\sqrt{c} \times$  reduction in on-chip memory against state-of-the-art for  $c$  columns in a bank row
  - $2 \times$  to  $4 \times$  improvement in execution time on limited on-chip memory architectures

- On-chip memory reduced to  $O(N)$  from  $O(\sqrt{cN})$ <sup>[1]</sup>

- Execution time reduced by a factor of  $(t_{bank} / t_{layer})$  to  $(t_{col} / t_{layer})$ <sup>[1]</sup>

16x reduction in on-chip memory

2x - 4x reduction in execution time



<sup>[1]</sup> "On-chip Memory Efficient Data Layout for 2D FFT on 3D Memory Integrated FPGA", Shreyas G. Singapura, Rajgopal Kannan and Viktor K. Prasanna, HPEC 2016 (Best Paper Finalist)

## Graph Analytics

- 2 Level Main memory:
  - High bandwidth, small capacity scratchpad
  - Low bandwidth, large capacity DRAM
- Trade DRAM accesses for Scratchpad accesses
- Reuse data in scratchpad
- Optimized data placement for graph processing
  - $2 \times$  to  $3 \times$  reduction in algorithm iterations
  - $1.7 \times$  to  $2.7 \times$  reduction in DRAM accesses
  - $1.4 \times$  to  $2 \times$  improvement in total memory accesses

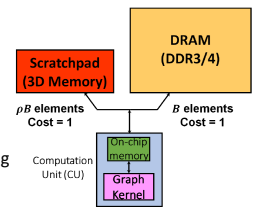
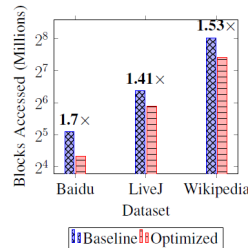
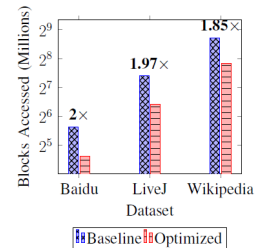


TABLE VI: Iterations for BFS

Dataset	DRAM		Scratchpad
	Baseline ( $d_1$ )	Optimized ( $d_2$ )	Optimized ( $d_{SA}^{min}$ )
Baidu	21	7	12
LiveJ	15	5	10
Wikipedia	48	11	41



(a) BFS



(b) SSSP