# Near Memory Computing
# Spectral and Sparse Accelerators

**Franz Franchetti**

ECE, Carnegie Mellon University
`www.ece.cmu.edu/~franzf`
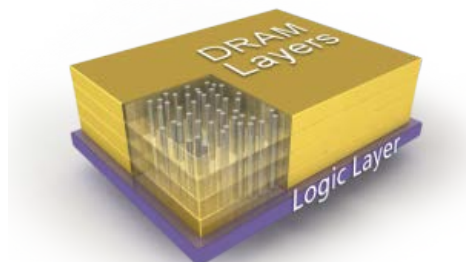
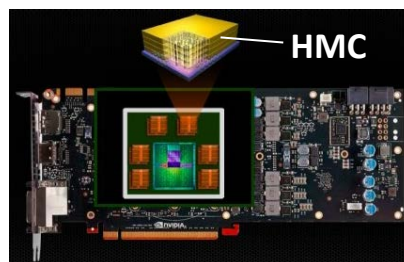Co-Founder, SpiralGen
`www.spiralgen.com`
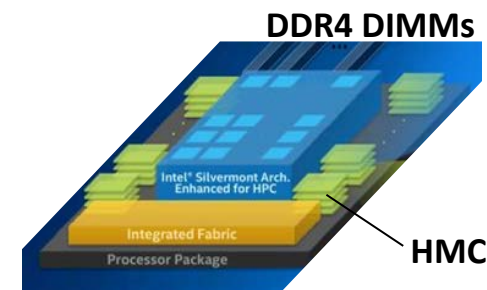
# DRAM-Optimized Near Memory Acceleration

- **Enabling technology: 3D stacked integration**
  - Logic and DRAM layers connected by TSVs
  - Better timing, lower area, advanced IO in the logic layer



**Micron HMC**



HMC

**Nvidia Volta**



DDR4 DIMMs

HMC

**Intel Knights Landing**

- **Accelerators in 3D DRAM, behind the conventional interface**
  - Bandwidth and latency concerns still exist to off-chip
  - Integration behind conventional interface opens up internal resources

- **DRAM operation and organization aware accelerators**
  - Conventional near memory computing only aim to reduce the communication distance between memory and processor

# Concept: Memory-Side Accelerators
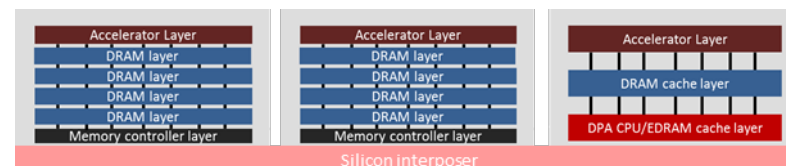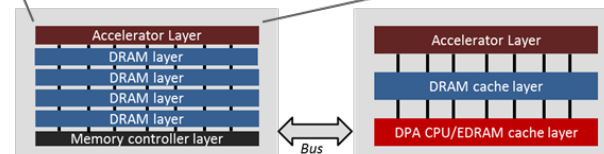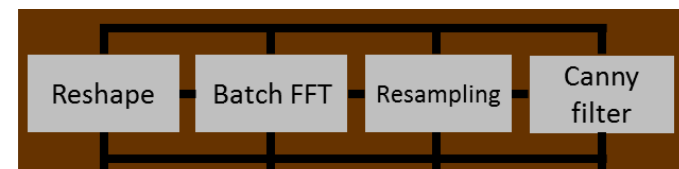
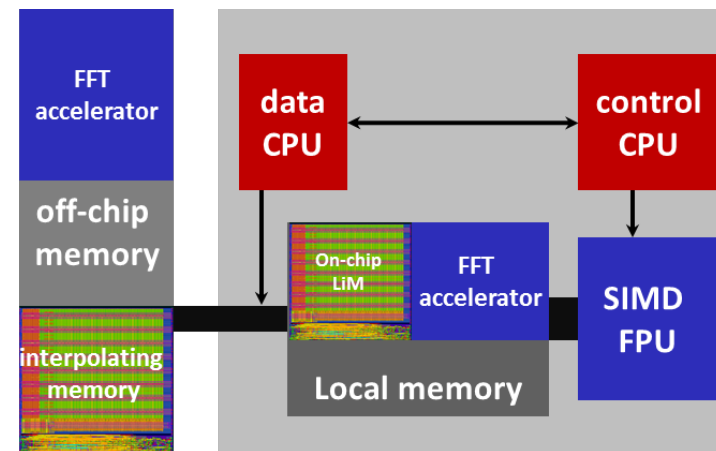- **Idea: Accelerator on DRAM side**
  - No off-DIMM data traffic
  - Huge problem sizes possible
  - 3D stacking is enabling technology
- **Configurable array of accelerators**
  - Domain specific, highly configurable
  - Cover DoD-relevant kernels
  - Configurable on-accelerator routing
- **System CPU**
  - Multicore/manycore CPU
  - CPU-side accelerators
  - Explicit memory management
  - SIMD and multicore parallelism

# Memory-Side Accelerator Architecture

- **DRAM-side Accelerator**

  - LiM layer in stacked DRAM

  - Array of acceleration engines

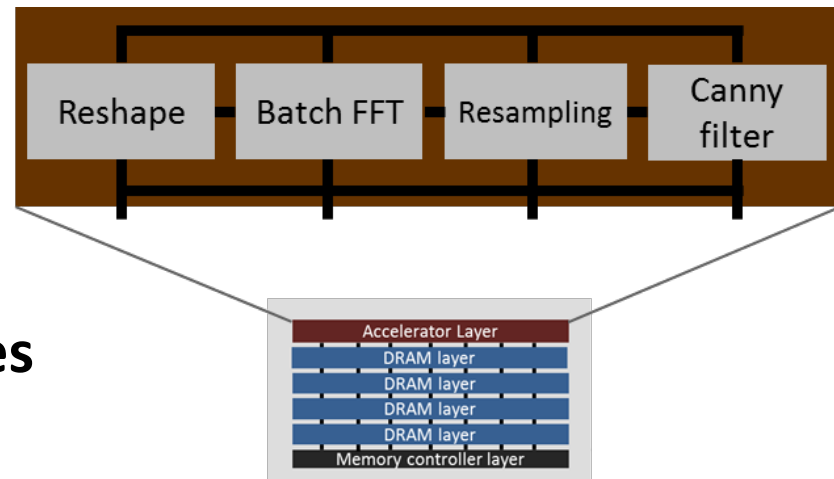  - Domain specific, highly configurable

- **Customized Acceleration Pipelines**

  - DRAM-to-DRAM streaming

  - Connect and configure multiple engines

  - Configurable on-LiM routing

- **Accelerator Cores**

  - **Reshape:** linear-to-block data layout changes

  - **Batch FFT:** primitive for large 1D and 2D FFTs

  - **Resampling:** interpolation, geometric transformations, image alignment

  - **Feature extraction:** edge and shape detection

- **Example:** PFA SAR requires Reshape, Batch FFT and Resampling

# Simulation, Emulation, and Software Stack

- **Accelerator Simulation**
  - **Timing:** Synopsis DesignWare
  - **Power:** DRAMSim2, Cacti, Cacti-3DD McPAT

- **Full System Evaluation**
  - Run code on real system (Haswell, Xeon Phi) or in simulator (SimpleScalar,…)
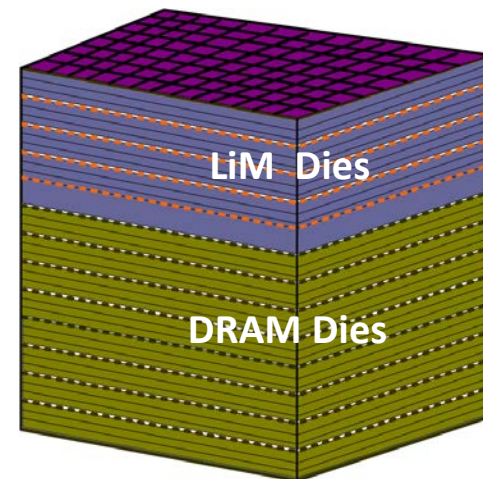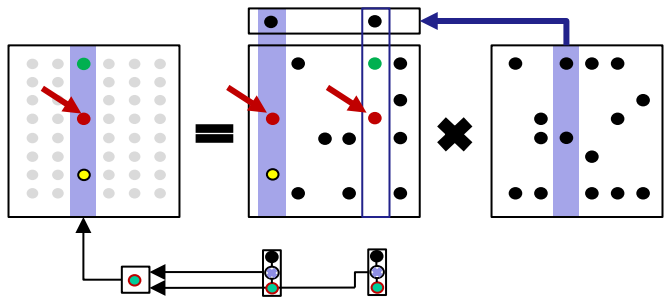  - Normal DRAM access for CPU, but trap accelerator command memory space, invoke simulator

```
#include <accelerator-fftw3.h>
...
{
  fftw_complex *in, *out;
  fftw_plan p;
  ...
  in = (fftw_complex*)
    fftw_malloc(sizeof(fftw_complex) * N);
  out = (fftw_complex*)
    fftw_malloc(sizeof(fftw_complex) * N);
  p = fftw_plan_dft_1d(N, in, out,
    FFTW_FORWARD, FFTW_ESTIMATE);
  ...
  fftw_execute(p); /* repeat as needed */
  ...
  fftw_destroy_plan(p);
  fftw_free(in); fftw_free(out);
}
```

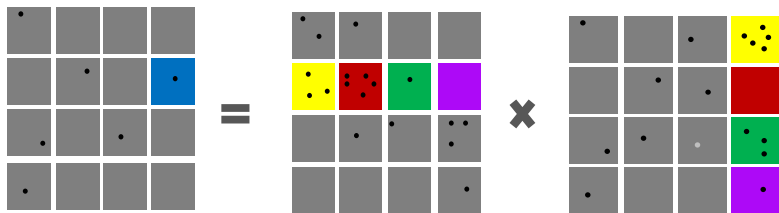- **API and Software stack**
  - **Accelerator:** memory mapped device with command and data address space
  - **User API:** C configuration library, standard API where applicable
  - **Virtual memory:** fixed non-standard logical-to-physical mapping
  - **Memory management:** Special `malloc/free`, Linux kernel support

Electrical & Computer
ENGINEERING

# Sparse Matrix Multiplication Accelerator

**On-chip:** regular SpGEMM kernel



**Off-chip:** SRUMMA (tiled shared memory algorithm)



**Scratchpad:** Hierarchical tiling

**SRAM**
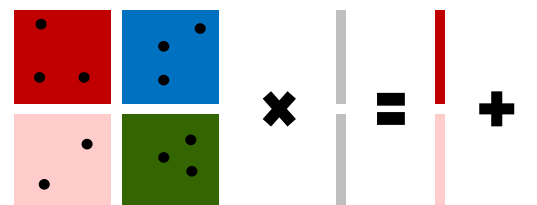logic-in-memory

**EDRAM**
scratchpad

**DRAM**

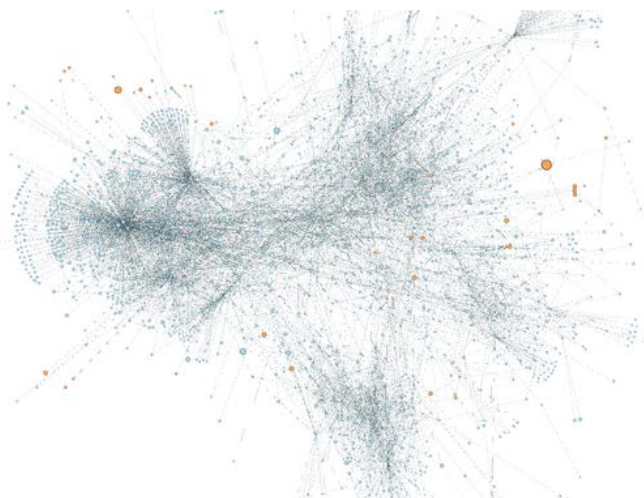# Sparse Matrix-Vector Product Accelerator

## SpMV Kernel



- Standard sparse matrix times dense vector
- Matrix is very sparse (a few non-zeroes per row)
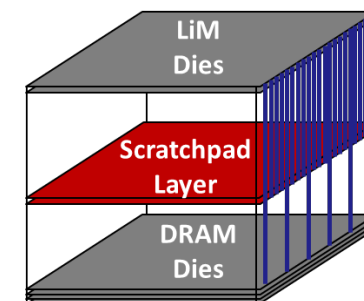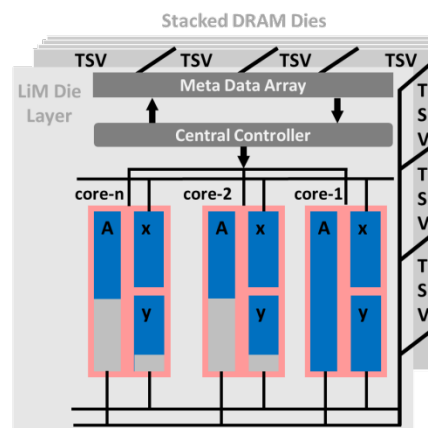- Matrix and vector size is large (>10M x 10M/GB)

## Algorithm



- Block-column multiply + merge step
- Partial results are streamed to DRAM
- Matrix streams from DRAM
- Vector segment is held in scratchpad/EDRAM
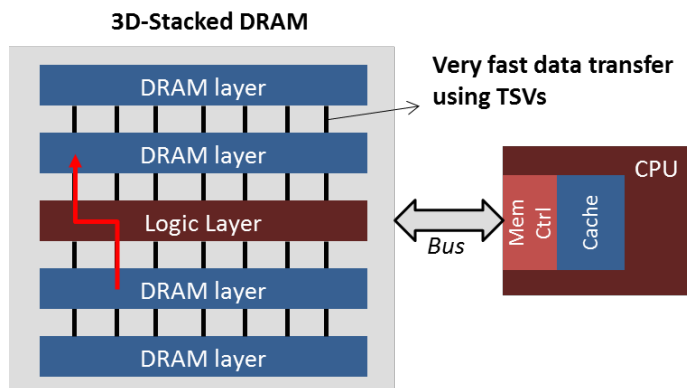
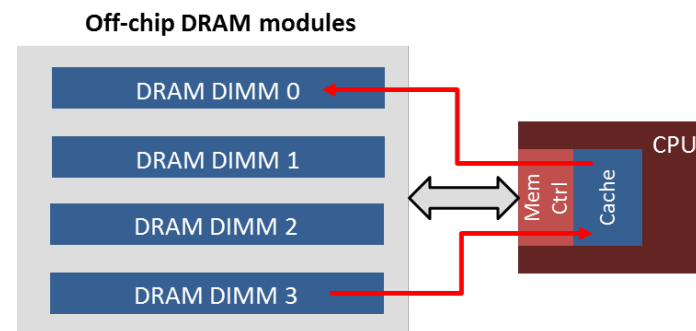## Target: Larges Sparse Graphs



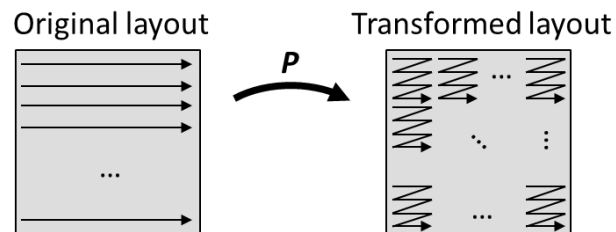## 3DIC Memory Side Accelerator
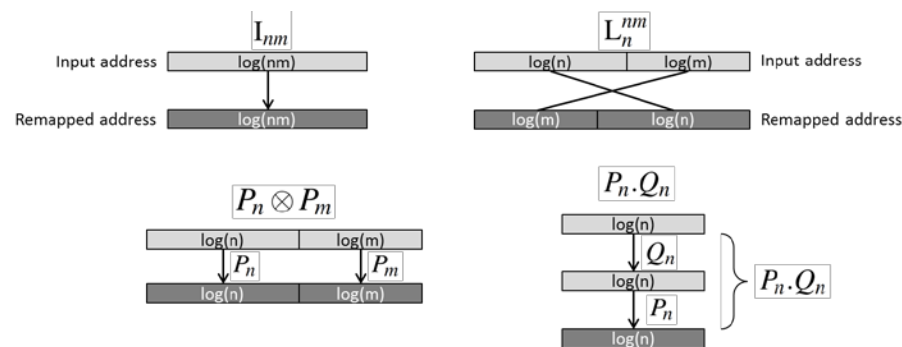
# In-DRAM Reshape Accelerator
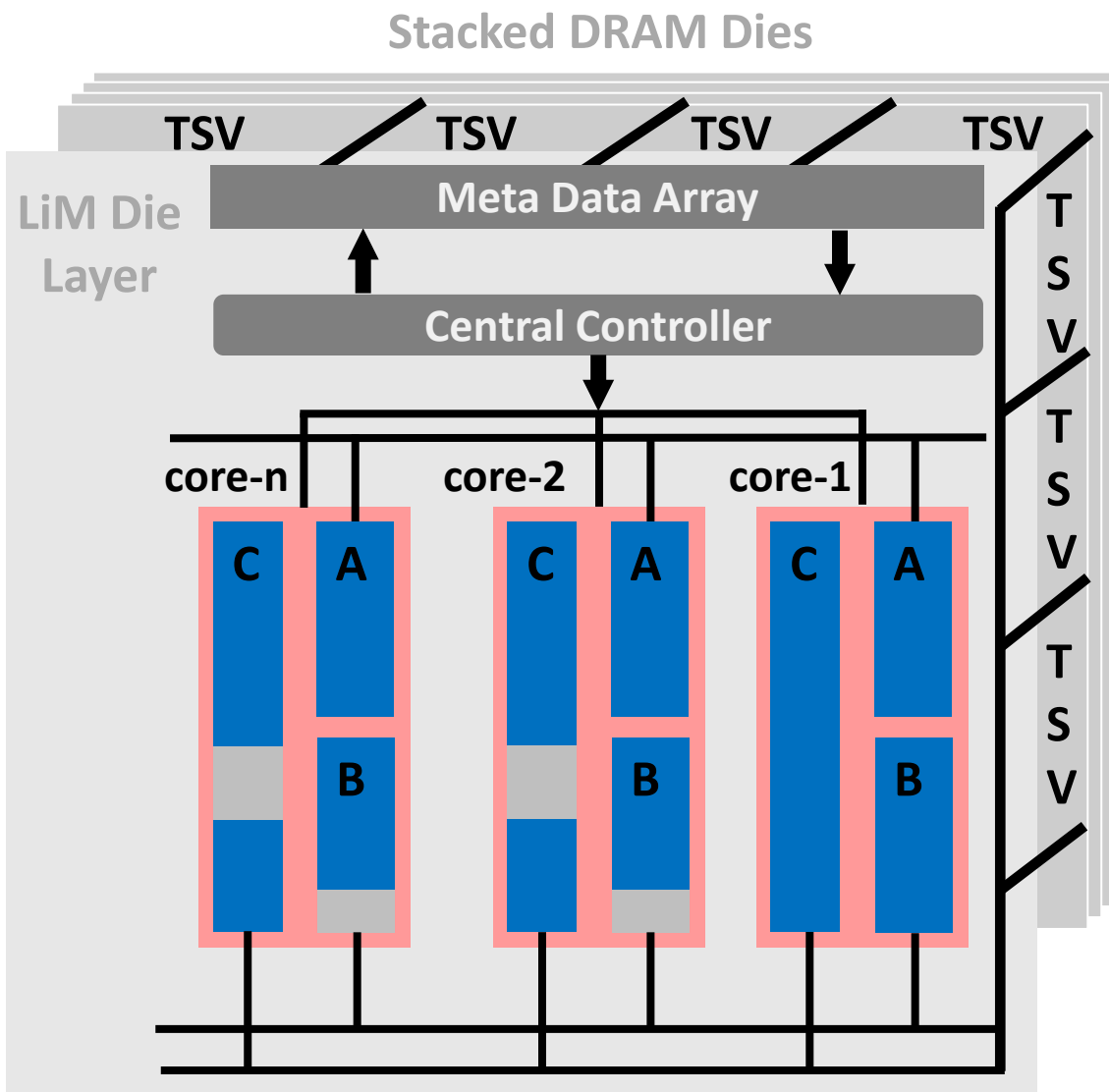
## 3DIC Reshape Stack
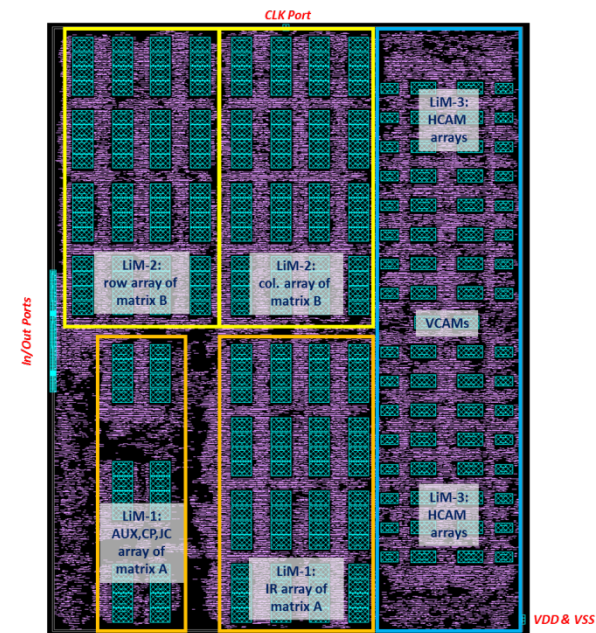


## Traditional System



## Reshape Operation



## Abstraction: Bit Permutations

# 3DIC DRAM + SpGEMM Multicore Design

**Stacked DRAM Dies**

**65nm test chip**



TSV   TSV   TSV   TSV

**LiM Die Layer**

**Meta Data Array**

**Central Controller**

core-n   core-2   core-1

C A   C A   C A

B   B   B

T S V T S V T S V T S V

- Taped out (65nm)
- Core Area: 1.003 x 1.292 mm$^2$
- Transistor count > 1M
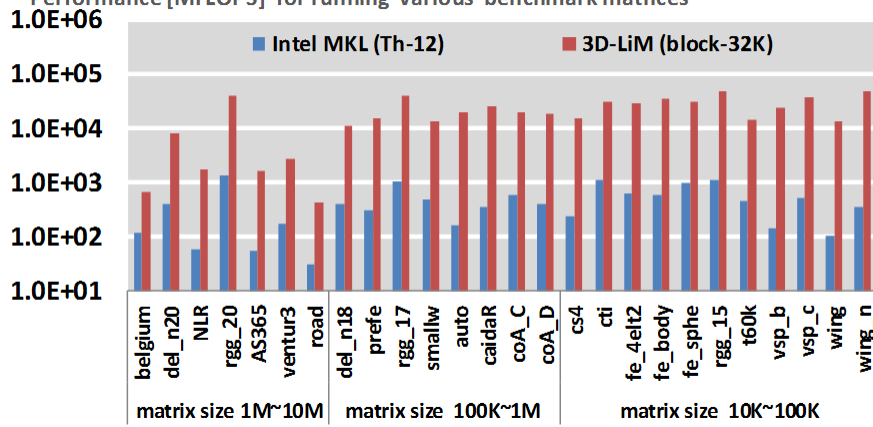- Frequency: 525MHz (SS @ICC)
- Power: 150mW (@DC)

Tape-out is funded under IARPA (PI Pileggi, Co-PIs Fedder, Franchetti, Piazza)

# SpGEMM: Experimental Results

## Intel Dual-CPU Intel Xeon E5-2430 system

- about 140W, 6 DRAM channels, 64 GB/s max, 6 DIMMs (48 chips), 210 GFLOPS peak

- Intel MKL 11.0 mkl_dcsrmultdcsr (unsorted):
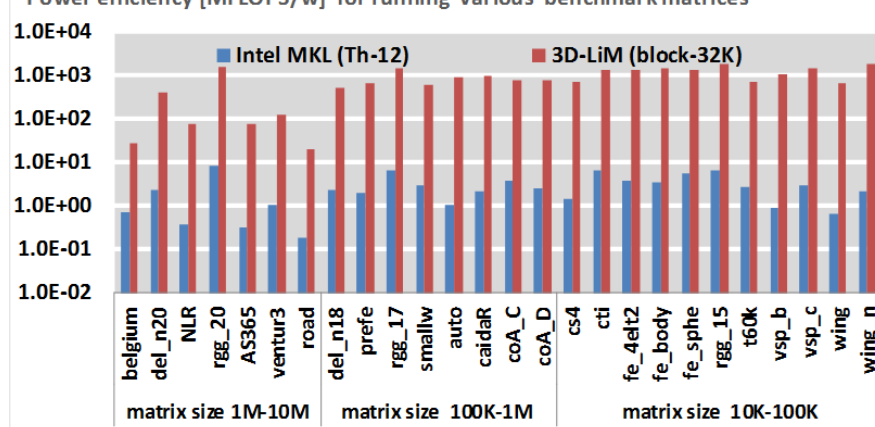  **100 MFLOPS – 1 GFLOPS, 1 – 10 MFLOPS/W**

## Our 3DIC System

- 4 DRAM layers@ 2GBit, 1 logic layer, 32nm, 30 – 50 cores, 1GHz, 30 – 50 GFLOPS peak

- **Performance: 10 – 50 GFLOPS** @ 668GB/s with 1024 TSV
  **Power efficiency: 1 GFLOPS/W** @ 350GB/s with 512TSV or 8GB/s with 1024 TSV



Matrices from University of Florida sparse matrix collection

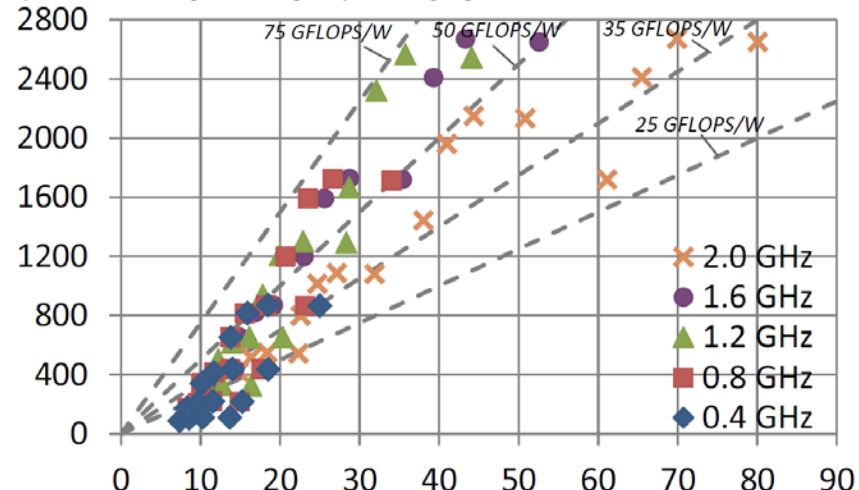# 3DIC for 2DFFT and PFA SAR



8192x8192 2D-FFT with 3D-stacked DRAM
performance [GFLOPS] vs power [W]

8192x8192 SAR with 3D-stacked DRAM
performance [GFLOPS] vs power [W]

- **3DIC:** 8Gbit, 4-layer DRAM + 1-layer logic,
  16 banks/layer, 512 TSV/bank, 1KB pages, 32nm (320GB/s max BW)

- **2DFFT:** tiled FFT, tile size matches DRAM row buffers

- **PFA SAR:** Polar to rectangular local interpolation with logic-in-memory
  More flops for the same # of accesses

Electrical & Computer
ENGINEERING

# Results: DRAM-aware FFT Algorithms

**DRAM-aware FFT algorithms optimized by SPIRAL**

$$\mathrm{DFT}_{n\times n} = \left(\left(R_{t2}^{\mathrm{T}}R_{t3}^{\mathrm{T}}(I_{n/k}\,\tilde{\otimes}\,I_k\otimes\mathrm{DFT}_n)R_{t3}R_{t2}\mid R_{t0}^{\mathrm{T}}R_{t1}^{\mathrm{T}}(I_{n/k}\,\tilde{\otimes}\,I_k\otimes\mathrm{DFT}_n)R_{t1}R_{t0}\right)^{\overleftarrow{Q}}\right)^{\overrightarrow{P}},\ \text{where } P=Q^{-1}=I_{n/k}\otimes L_{n/k}^{n}\otimes I_k.$$

$$\mathrm{DFT}_{n\times n\times n} = \left(\left((I_{n^2/k^2}\,\tilde{\otimes}\,I_{k^2}\otimes\mathrm{DFT}_n)^{R_{c2}}\mid(I_{n^2/k^2}\,\tilde{\otimes}\,I_{k^2}\otimes\mathrm{DFT}_n)^{R_{c1}}\mid(I_{n^2/k^2}\,\tilde{\otimes}\,I_{k^2}\otimes\mathrm{DFT}_n)^{R_{c0}}\right)^{\overleftarrow{Q}}\right)^{\overrightarrow{P}},\ P=Q^{-1}=(I_{n^2/k^2}\otimes L_{n/k}^n\otimes I_{k^2})(I_{n/k}\otimes L_{n/k}^n\otimes L_{n/k}^n\otimes I_k)$$

$$\mathrm{DFT}_{n^2} = \left(\left(R_{t2}^{\mathrm{T}}R_{t3}^{\mathrm{T}}(I_{n/k}\,\tilde{\otimes}\,I_k\otimes\mathrm{DFT}_n)R_{t3}R_{t2}\mid R_{t0}^{\mathrm{T}}R_{t1}^{\mathrm{T}}D_n^{n^2}(I_{n/k}\,\tilde{\otimes}\,I_k\otimes\mathrm{DFT}_n)R_{t3}R_{t2}\right)^{\overleftarrow{Q}}\right)^{\overrightarrow{P}},\ \text{where } P=Q^{-1}=I_{n/k}\otimes L_{n/k}^n\otimes I_k.$$

**Performance model and power/performance simulation results**



(a) 2D-FFT with Tiled Data Layout
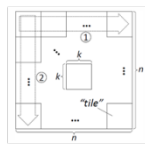
(b) 3D-FFT with Cubic Data Layout

(c) 1D-FFT with Tiled Data Layout



2D-FFT Power Consumption Breakdown (W)

2D tiles

3D bricks

# Near Memory Computing in DARPA PERFECT

## HW/SW Formalization

### *SPIRAL System*

### Hardware Synthesis

### Software Synthesis

## Algorithm+HW Design

### *Memory-Side Accelerators*

### 3DIC System

### Accelerator Architecture

**75 GFLOPS/W**

## Low Power Accelerator Cores

### *Logic-in-memory for sub-22nm CMOS*

### Design Tools and Simulation