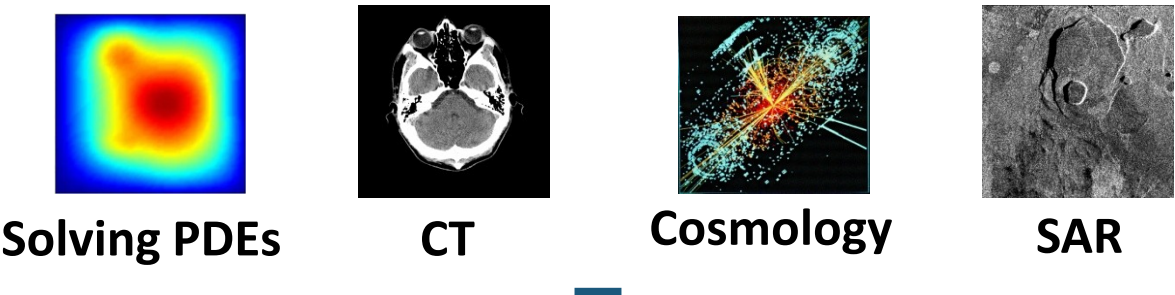


Overview

Large size 2D Fast Fourier Transform

- Used in image processing, scientific computing



- Typical datasets are **large** and **high precision!**

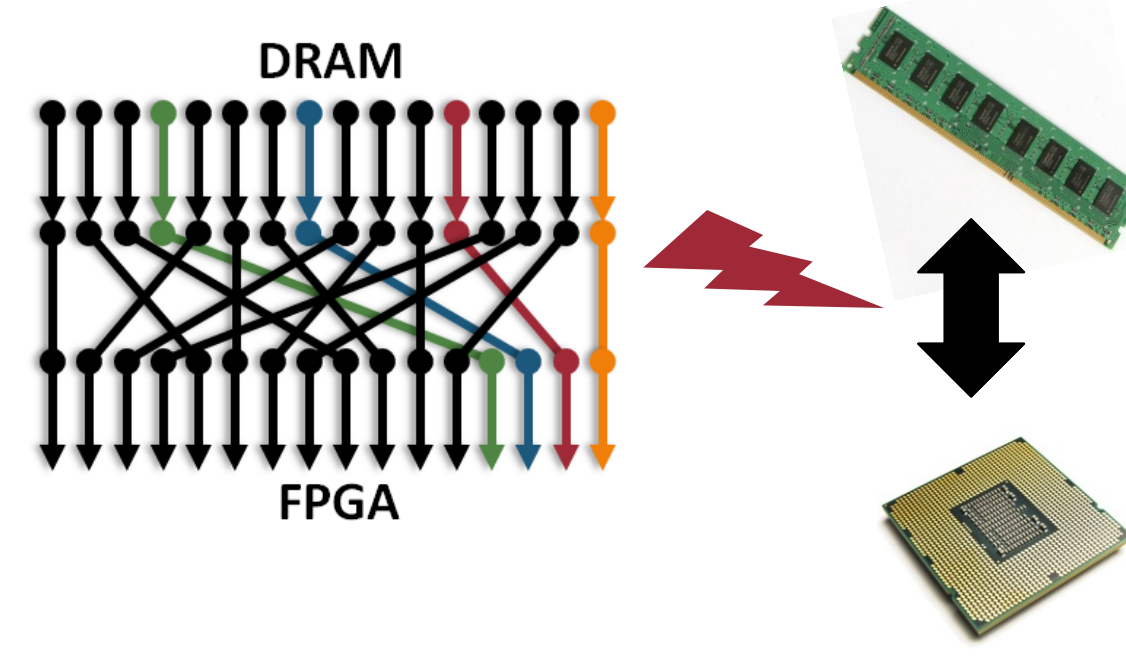
- e.g. 2K-by-2K double precision 2D-FFT:
 - Input dataset: **64 MB**
 - # of operations: **~461.4 Mflop**

- Does not fit on-chip
- Stored **off-chip**



Memory access pattern and achieved bandwidth

- Have large **strided** DRAM access pattern
- Does not exploit DRAM **row-buffer locality**



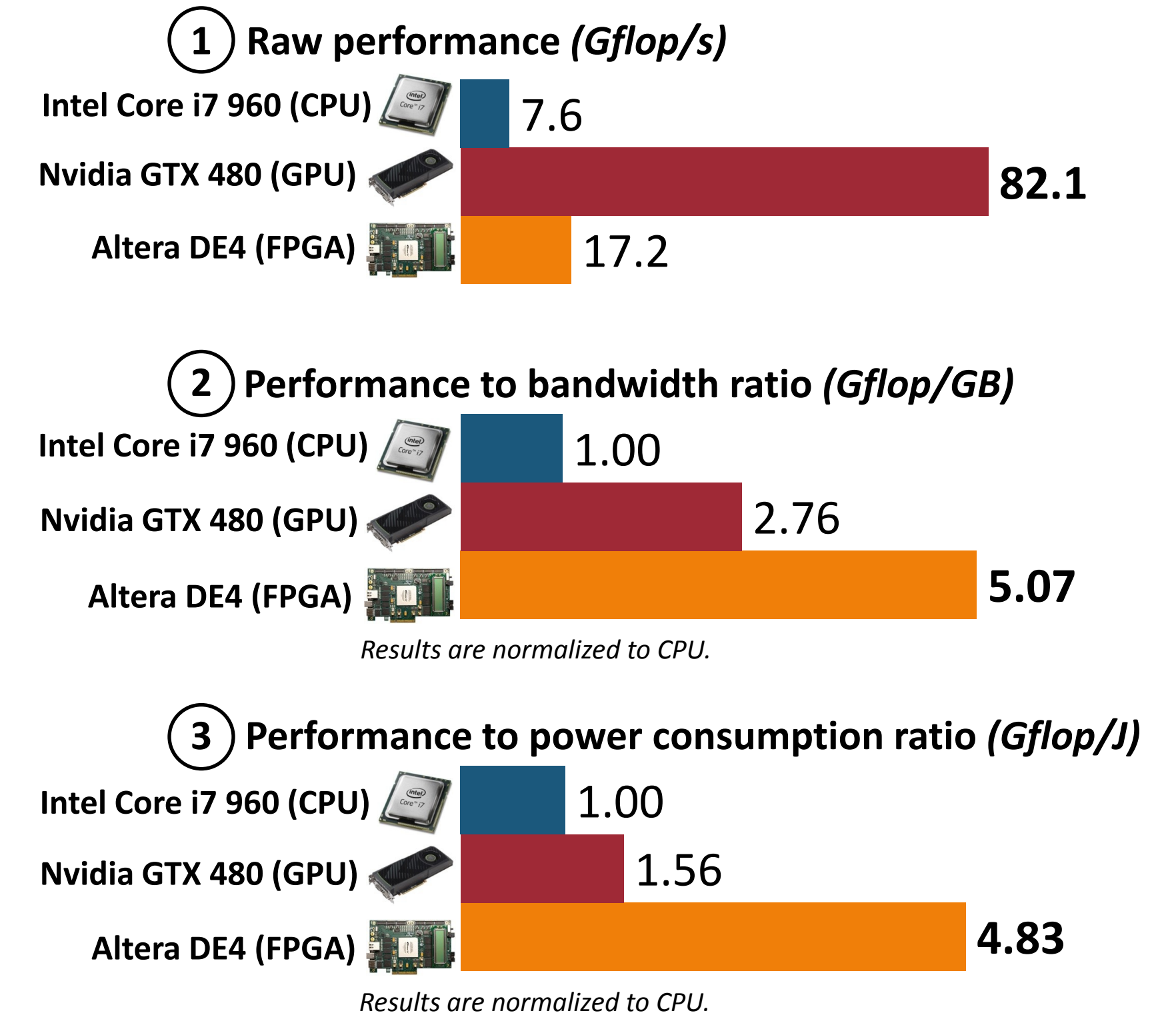
- Results in **low** memory bandwidth utilization!
- Memory bandwidth becomes **bottleneck** for achieving high performance



- Effective **bandwidth orchestration** is required for:

- Performance
- Bandwidth Efficiency
- Power Efficiency

1024-by-1024 double precision 2D-FFT



Background

DFT is matrix-vector multiplication

$$y = \text{DFT}_n \cdot x, \quad \text{DFT}_n = [e^{-2\pi i k l / n}]_{0 \leq k, l < n}$$

FFT algorithm is factorization of DFT matrix

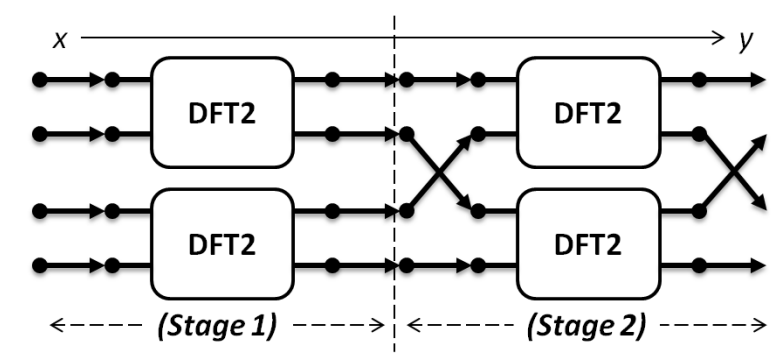
$$\text{DFT}_4 = (\text{DFT}_2 \otimes \text{I}_2) \text{T}_2^4 (\text{I}_2 \otimes \text{DFT}_2) \text{L}_2^4$$

tensor twiddle factors permutation

2D-FFT operates on 2D data, e.g. images

$$\text{DFT}_{2 \times 2} = (\text{DFT}_2 \otimes \text{I}_2) (\text{I}_2 \otimes \text{DFT}_2)$$

Stage 2 Stage 1

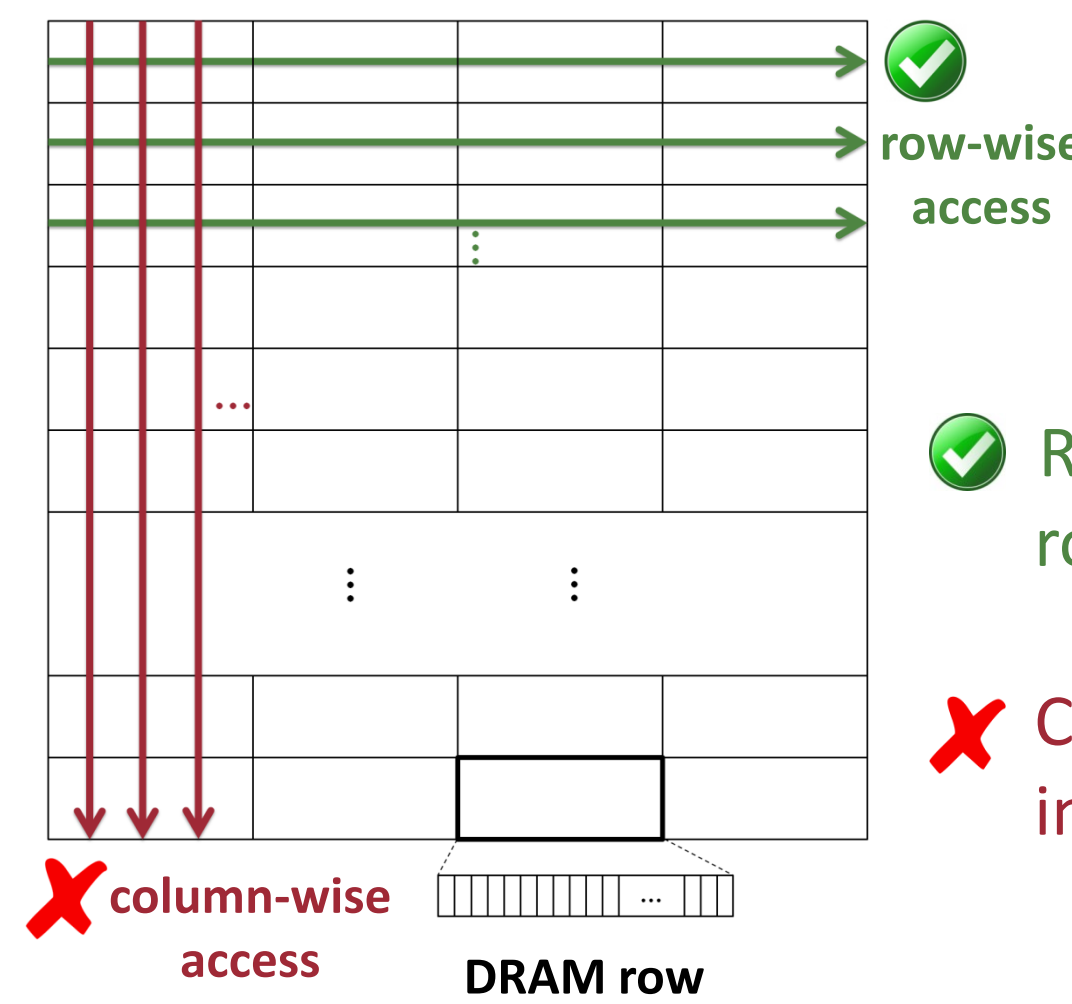


2D-FFT algorithms

- Row column algorithm:

$$\text{DFT}_{n \times n} = (\text{DFT}_n \otimes \text{I}_n) (\text{I}_n \otimes \text{DFT}_n)$$

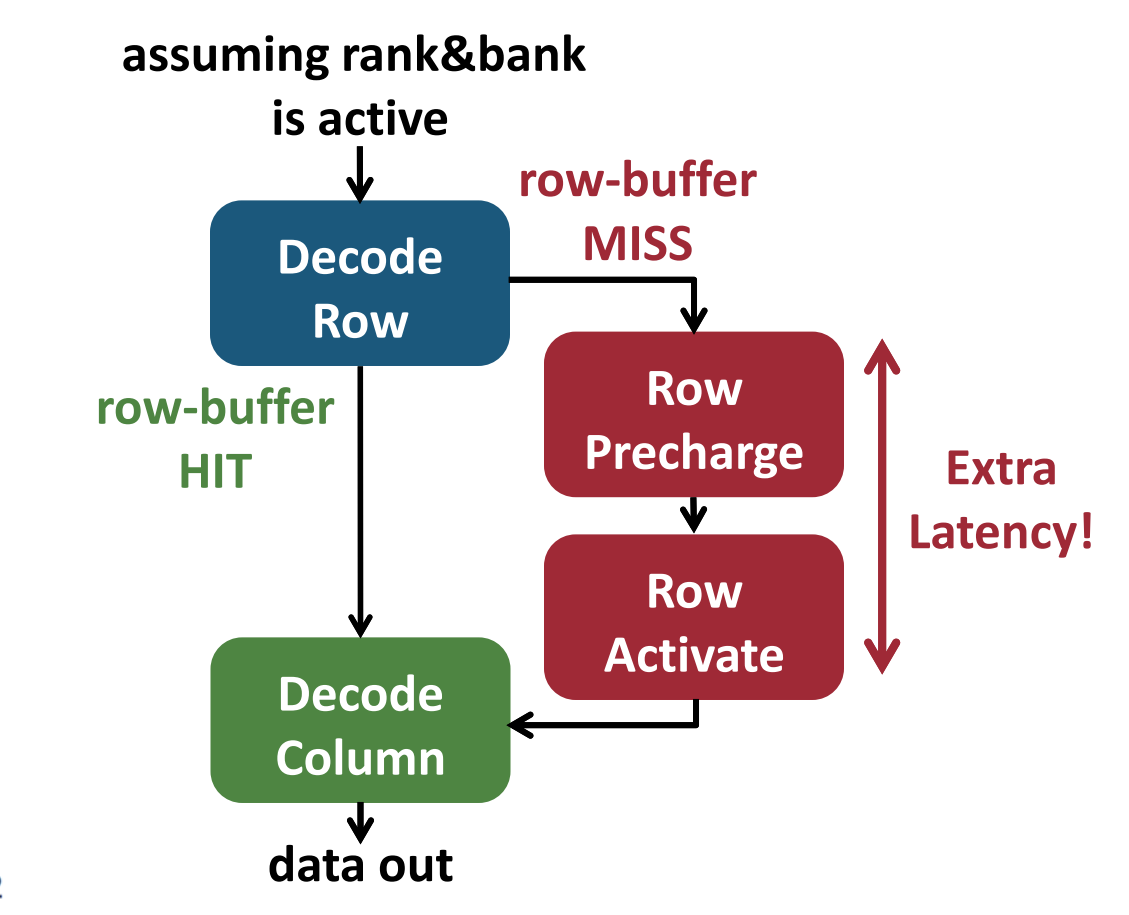
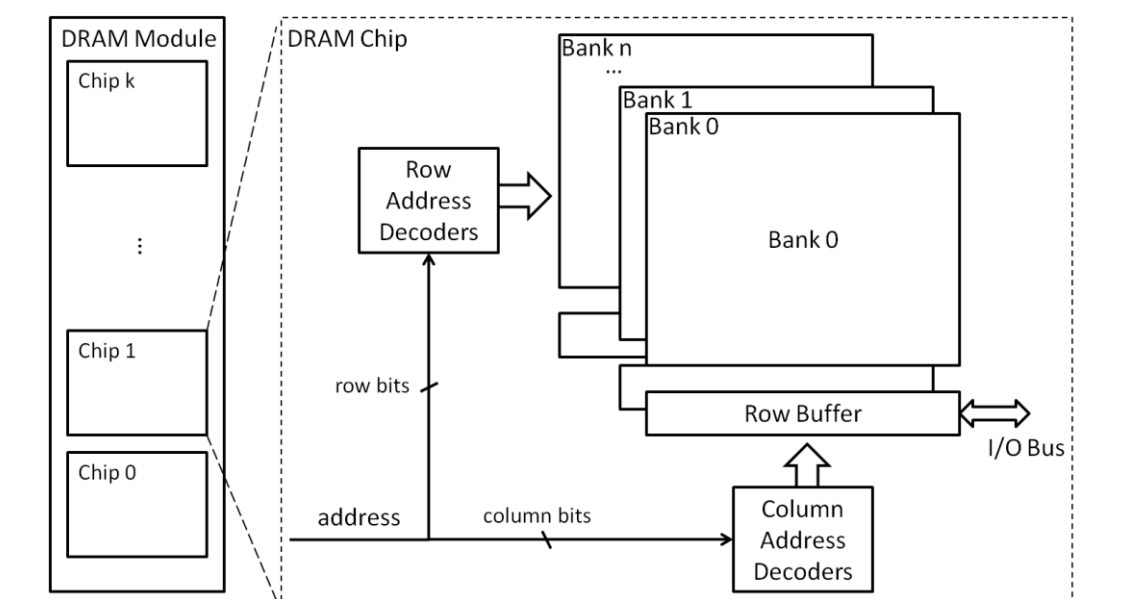
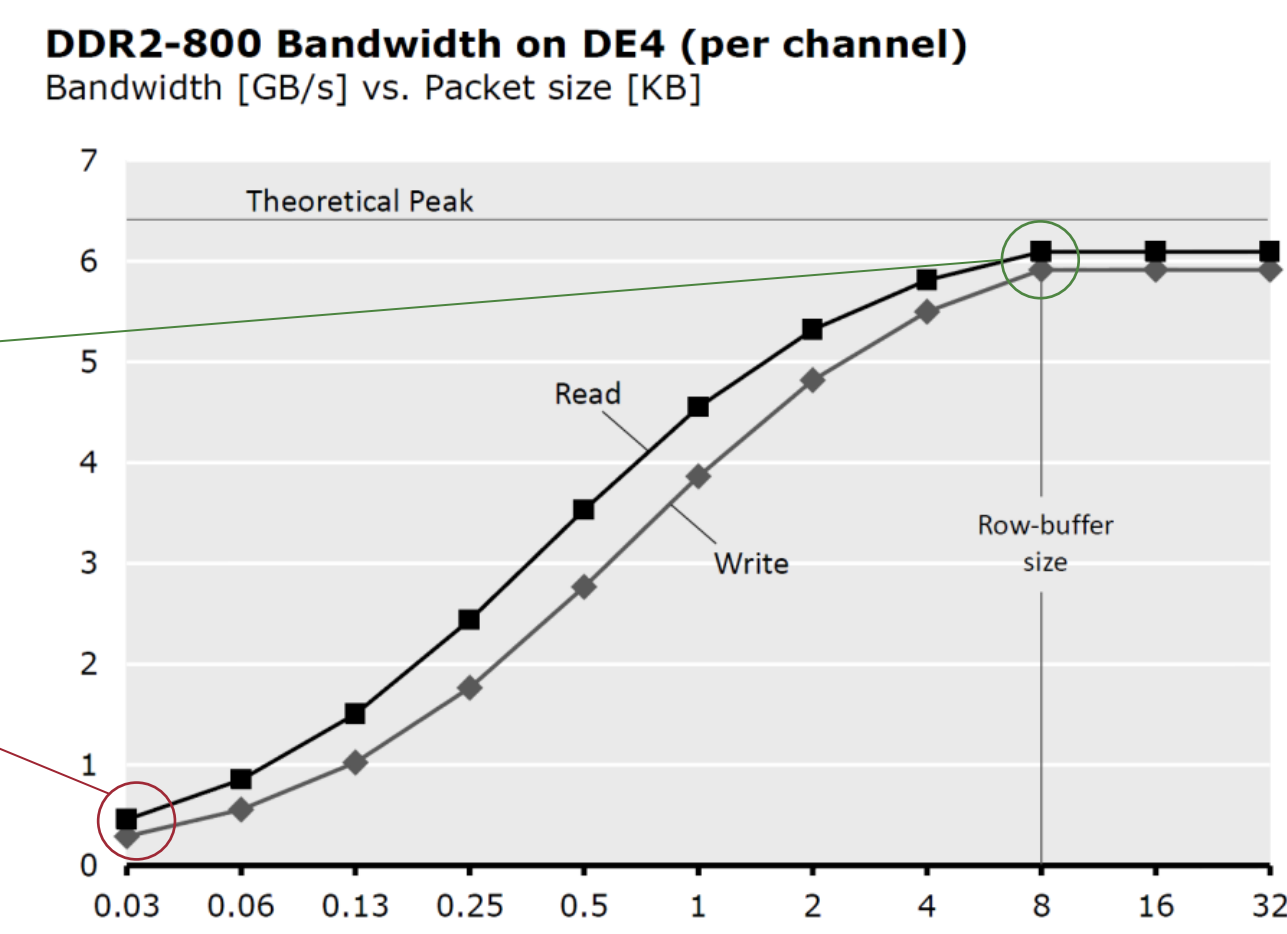
Row-wise and column-wise accesses!



- Row-wise access exploits row-buffer locality
- Column-wise access results in row-buffer misses

DRAM operation

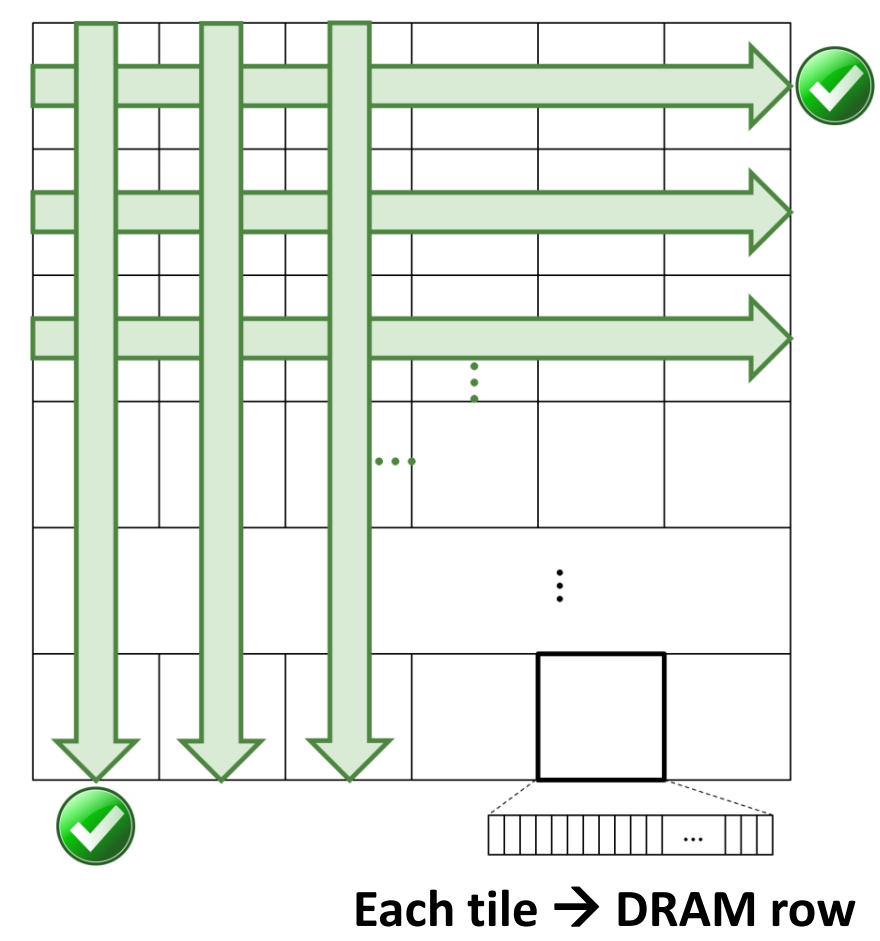
- Need to make use of every row touched to **maximize bandwidth**
- Large strides result in **small packets** of transferred data



Solution: Algorithm and Architecture

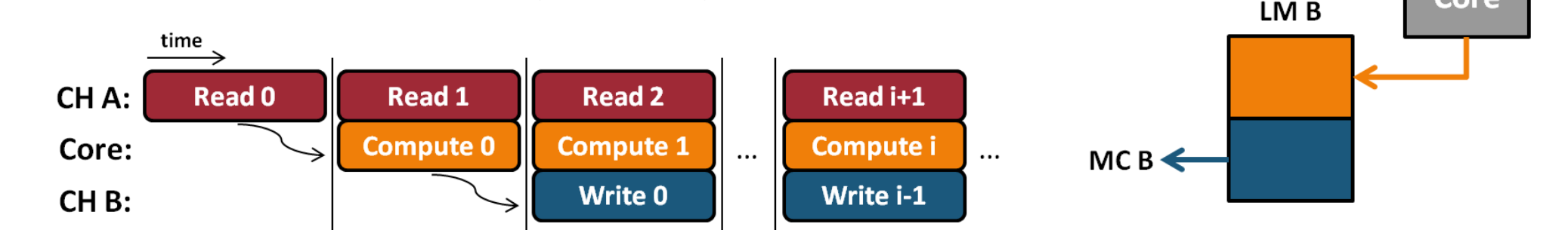
Restructured algorithm

- Linear data mapping** in DRAM causes row and column-wise accesses
- Use **2D tiled** data mapping where each tile is mapped to a DRAM row
- Restructure** the algorithm given 2D data mapping



- Data is accessed as **tiles**, not row and column-wise
- Row-buffer misses are minimized!

- Double-buffering:
 - Overlapped **computation** and **I/O**
 - All modules kept busy



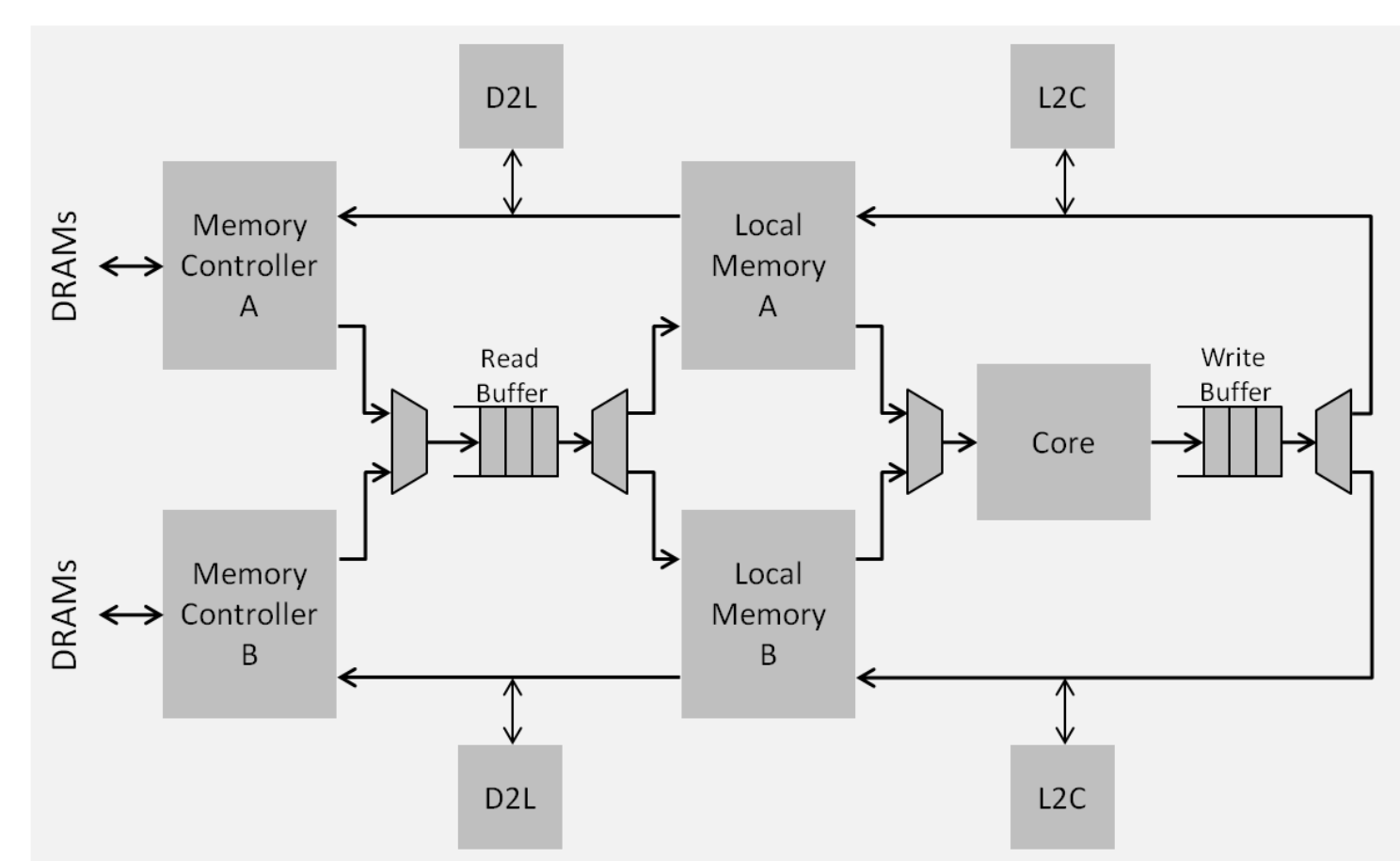
$$\text{DFT}_{n \times n} = \prod_{i=0}^{n-1} (\text{L}_{n^2/k}^i (\text{I}_n \otimes \text{DFT}_n) \text{I}_{n^2})$$

hardware aware manipulation

$$\text{DFT}_{n \times n} = \prod_{i=0}^{n-1} ((\text{L}_{n^2/k}^i \otimes \text{I}_k) (\text{I}_{n/k} \otimes \text{L}_{n^2/k}^i) (\text{I}_{n/k} \otimes \text{I}_k \otimes \text{DFT}_n) (\text{I}_{n/k} \otimes \text{L}_{n^2/k}^i) (\text{I}_{n/k} \otimes (\text{L}_{n^2/k}^i \otimes \text{I}_k)) (\text{I}_{n/k} \otimes (\text{L}_{n^2/k}^i \otimes \text{I}_k)))$$

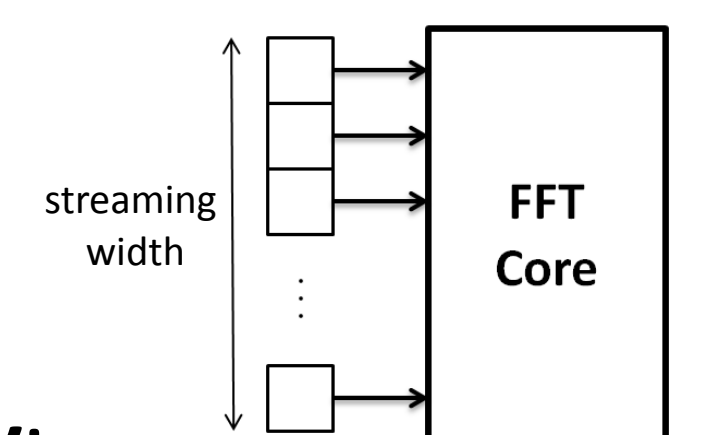
- write tiles into DRAM → DRAM write address generation
- on-chip permutation → Local Memory address generation
- DFT computation → Streaming FFT Core
- on-chip permutation → Local Memory address generation
- read tiles from DRAM → DRAM read address generation

From algorithm to hardware



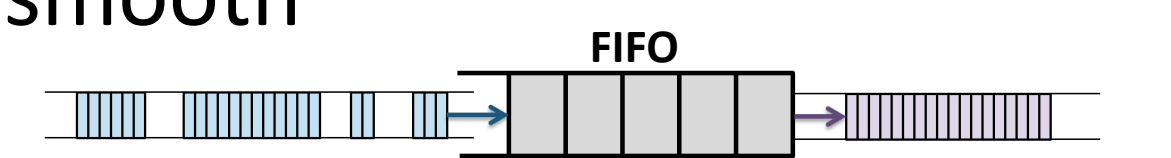
- Matching **throughput** to memory **bandwidth**:

- Achieved via fine-grain control over datapath **parallelism**
- Results in **balanced** design



- Ensuring **continuous dataflow**:

- Buffers are used to smooth the flow of data.



Evaluation

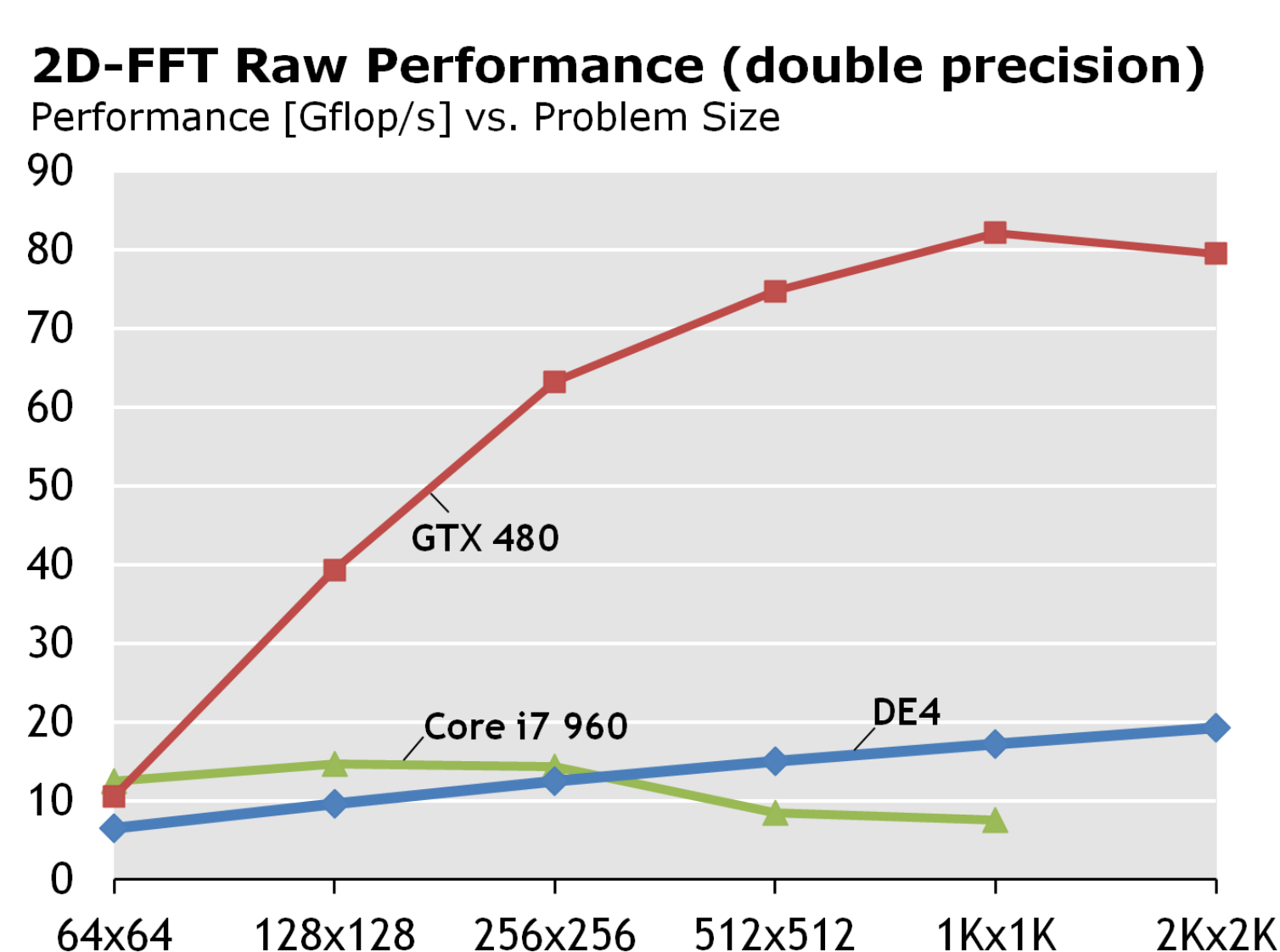
Target application:

- Double precision complex 2D-FFT
- Data sizes up to 2,048-by-2,048

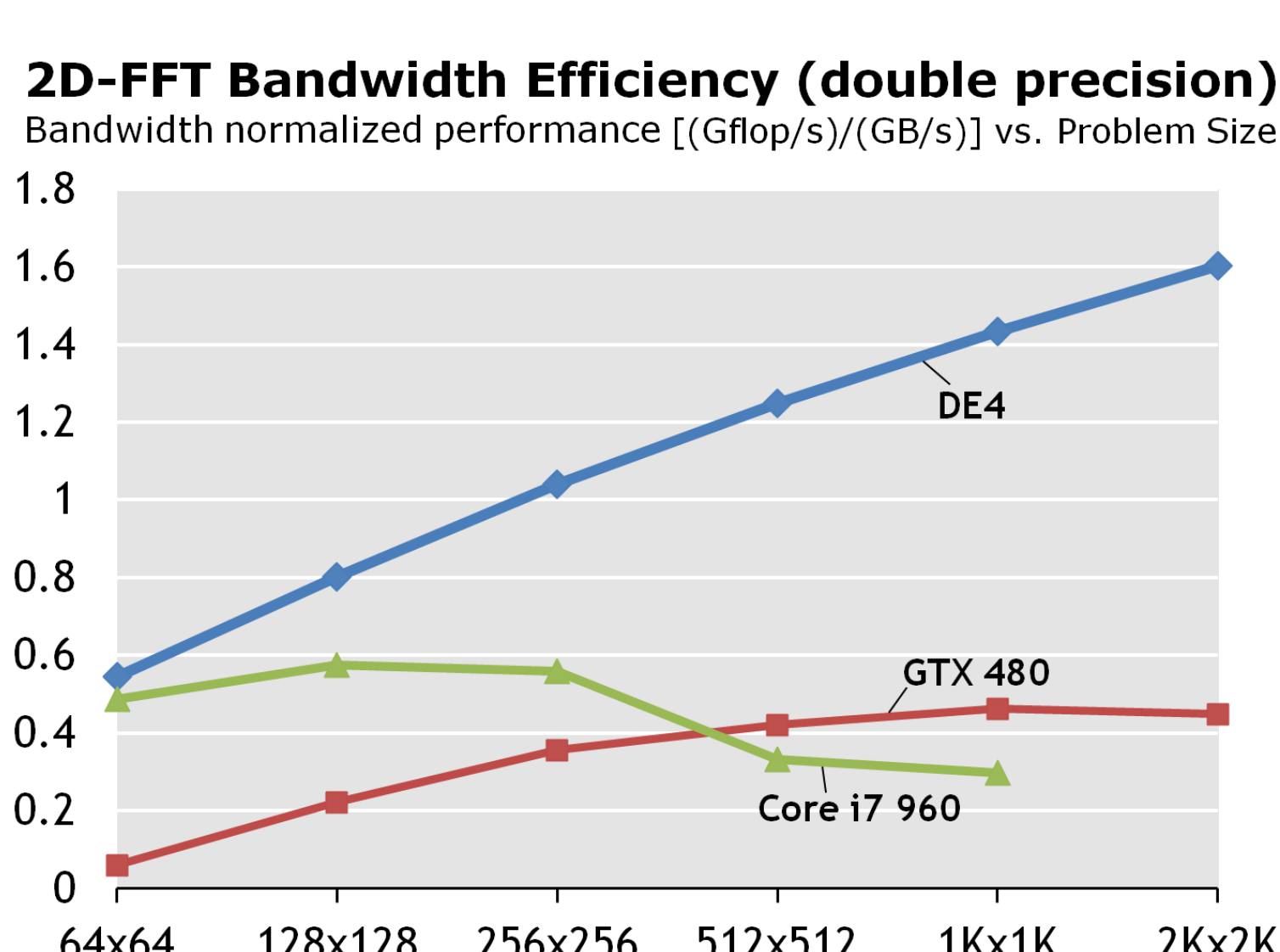
Target platforms:

	Core i7 960	GTX 480	Stratix IV (DE4) EP4SGX530
DRAM Type	DDR3	GDDR5	DDR2
# of Memory Channels	3	6	2
Memory BW (GB/s)	25.6	177.4	12
On-chip Memory (MB)	8	1.69	2.53
Proc. Freq. (MHZ)	3,200	1,401	200
# of Cores	4	480	N/A
Technology Node (nm)	45	40	40
Application Infrastructure	Spiral	CUDA 4.0	Spiral/Verilog

Raw performance:



Bandwidth Efficiency:



Power* Efficiency:

