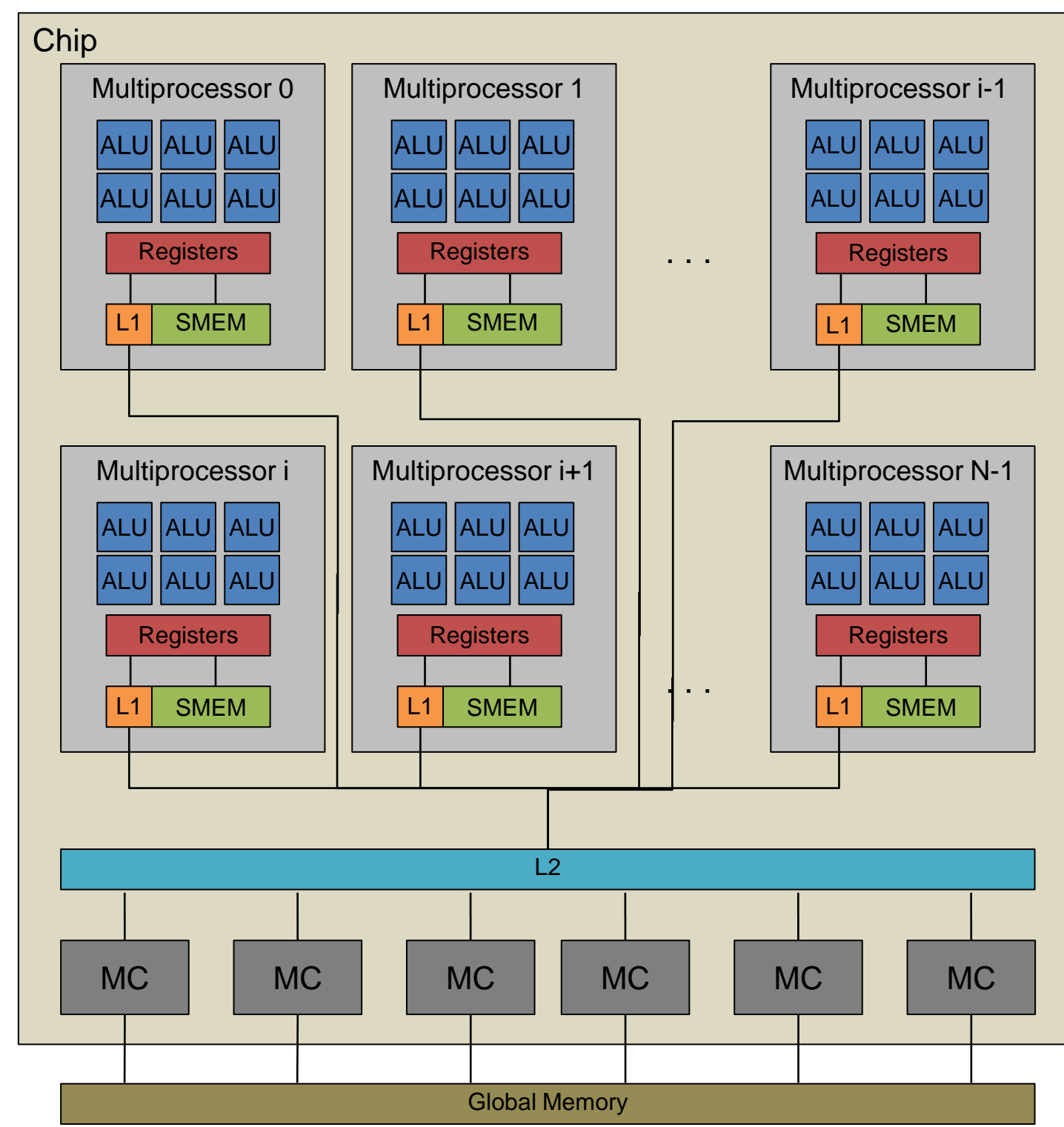
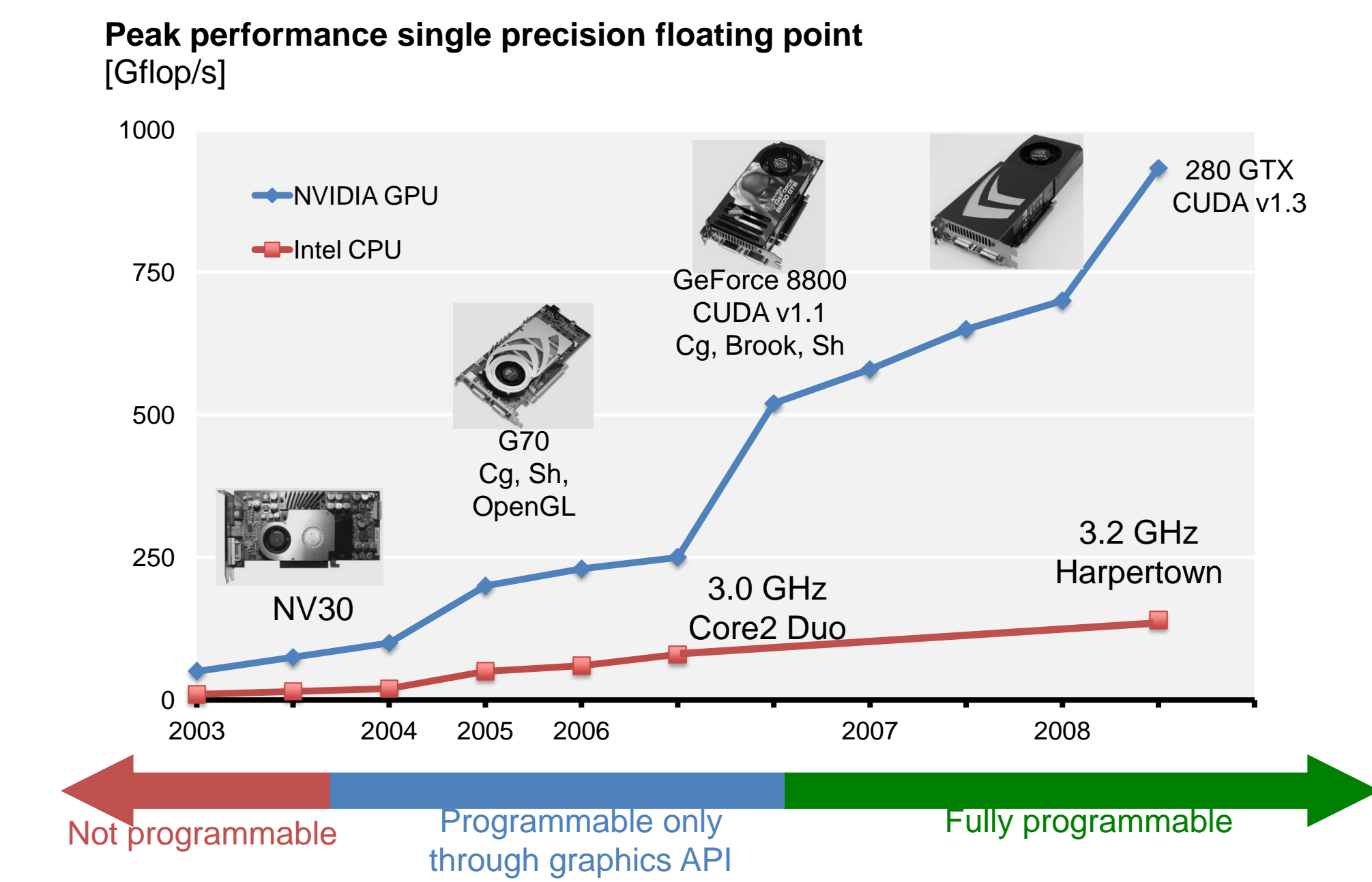


Automatic Generation of FFT Libraries for GPUs

Christos Angelopoulos, Franz Franchetti and Markus Püschel



GPUs and Programmability GPU Architecture Model



Shared Memory Optimized GPU DFT Algorithm

Original Stockham :

$$DFT_{r,n} \rightarrow \prod_{i=0}^{n-1} (DFT_r \otimes I_{r^{n-1-i}}) T_i^{r^n} (L_r^{n-1} \otimes I_r)$$

GPU Stockham :

$$DFT_{r,n} \rightarrow \prod_{i=0}^{n-1} ((D_i^r \cdot DFT_r \cdot Z_4^i \otimes I_p \otimes I_{\mu} \otimes I_{r^{n-1-i}}) L_r^{n-1})$$

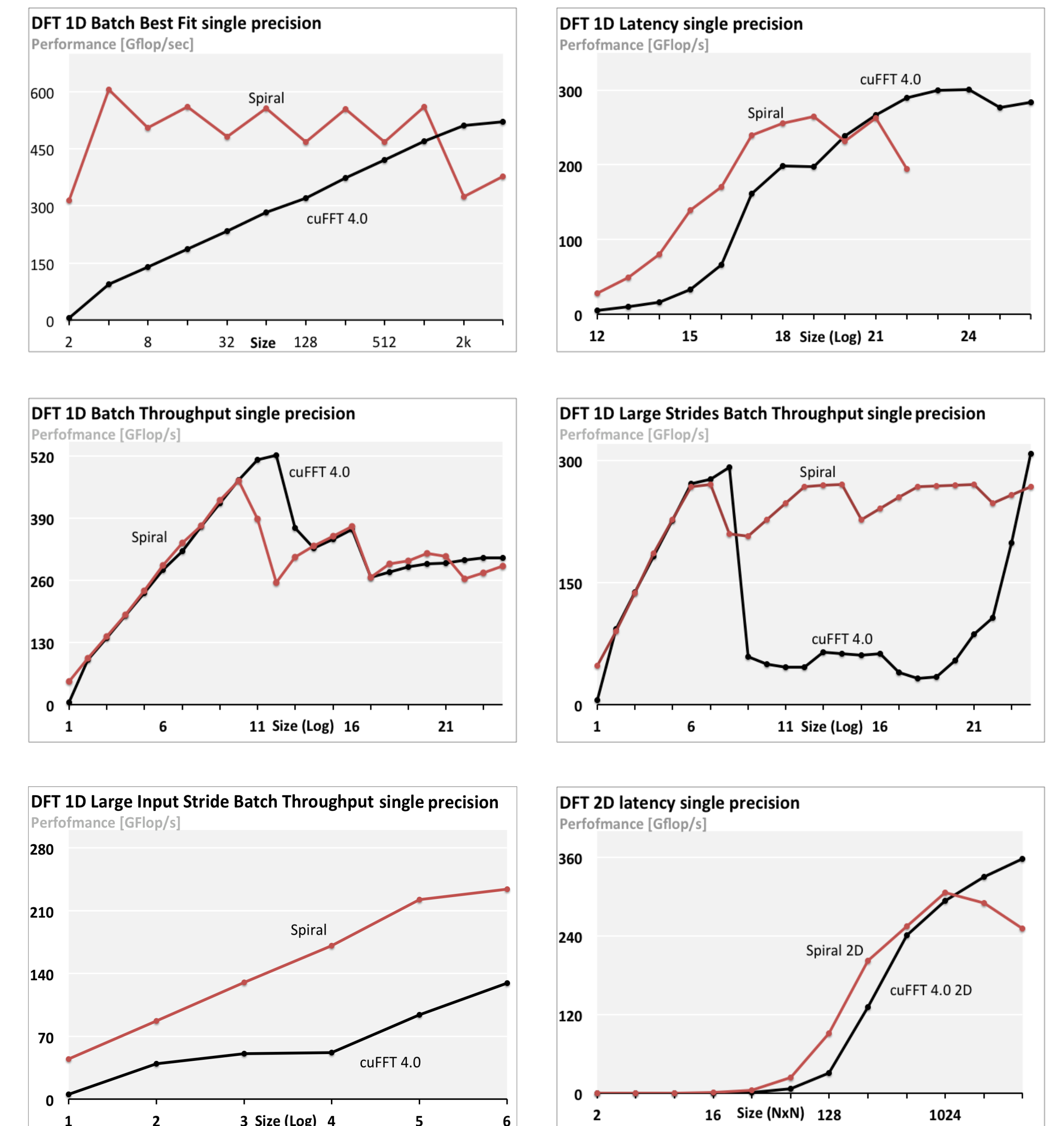
Rules:

- Loop Splitting
- Loop Interchange
- Loop Tiling for unit stride outputs
- Cyclic Shift Property to avoid bank conflicts

Advantages:

- Inplace algorithm
- Minimize global memory transfers
- Use Shared memory
- Register to register computations

Results on the GTX 480



Philosophy

Vendor Libraries

- Expertise required
- Platform specific
- Manual re-tuning

Auto-tuned Libraries

- Extensions to: New architectures, New algorithms

Computer Generated Libraries: SPIRAL for GPUs

Specification → Adaptive GPU library → Auto-tuning at runtime

Architecture

- 15 Multiprocessors
- 32 cores per multiprocessor
- 32 K registers per multiprocessor
- 48 KB of shared memory
- 16 KB of L1 cache
- 768 KB of L2 cache
- 1.5 GB of GPU Memory

Restrictions

- Banked Shared Memory
 - 32 banks
 - Within one warp resolve bank conflicts
 - Every thread in the warp Reads/Writes at different bank
 - 32 threads in a warp to 32 banks
- Register pressure
 - Max registers per MP = 32K/# of threads per MP
- Uncommon Architectural Model
 - Size of registers > Size of caches
 - Global Memory
 - Only block transfers, using caches
 - Double buffering

Key problems: Parallelism, Vectorization, Memory Hierarchy

Domain Specific Language: $A \otimes B = [\alpha_{k,l}]B$, for $A = [\alpha_{k,l}]$

Address Hardware Characteristics: $DFT_{mn} = (DFT_k \otimes I_m) T_m^n L_k^n (DFT_m \otimes I_k)$

Advantages: Efficient handling of complexity, Efficient porting to new platforms

Identify hardware parameters: $vector(v)$

Identify formulas: $A \otimes I_v, D_n, L_v^2$

Derive algorithm: $DFT_{mn} = ((DFT \otimes I_{n/v} \otimes I_v) T_m^{mn} (I_{m/v} \otimes (I_{n/v} \otimes L_v^2)) (L_{n/v}^{mn/v} \otimes I_v) (DFT_n \otimes I_v) (L_{m/v}^{mn/v} \otimes I_v))$

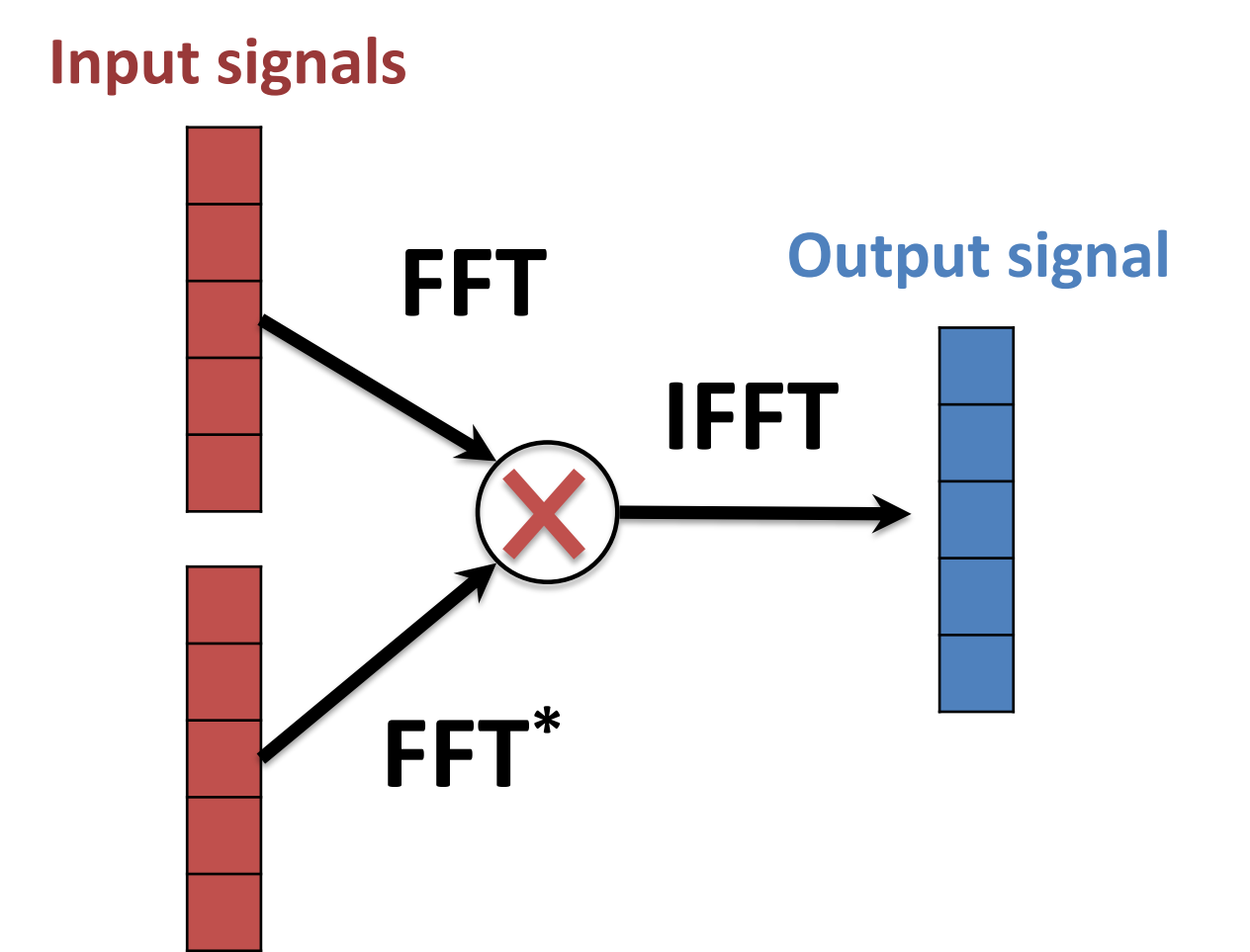
Generate platform oriented algorithm: $y = (I_1 \otimes DFT_2)x, y = (DFT_2 \otimes I_1)x$

Generate Platform Tuned Code

Next Step

Correlation (Frequency Domain)

- Code generator**
 - Only one data transfer from CPU DRAM to GPU
 - Minimize GPU DRAM memory roundtrips
- Application Scenario**
 - PDE solvers
 - Huge correlations



Future Work

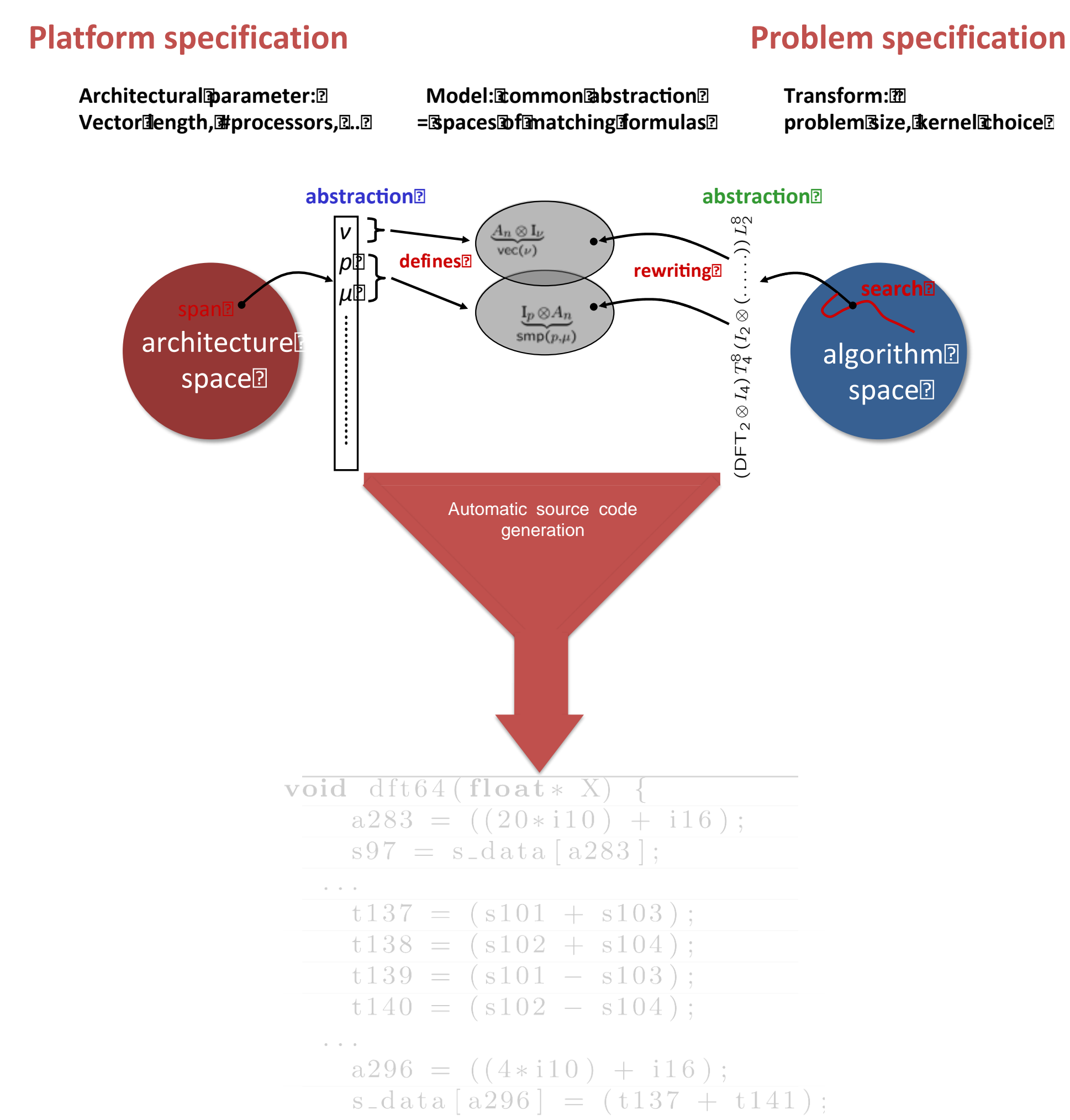
- Fast PDE solvers on GPUs

This work was supported by DARPA DESA Program and Nvidia

References:

- F. Franchetti, M. Püschel, Y. Voronenko, Sr. Chellappa and J. M. F. Moura "Discrete Fourier Transform on Multicore" IEEE Signal Processing Magazine, special issue on "Signal Processing on Platforms with Multiple Cores", Vol. 26, No. 6, pp. 90-102, 2009
- M. Püschel, J. M. F. Moura, J. Johnson, D. Padua, M. Veloso, B. Singer, J. Xiong, F. Franchetti, A. Gacic, Y. Voronenko, K. Chen, R. W. Johnson and N. Rizzolo "SPIRAL: Code Generation for DSP Transforms Proceedings of the IEEE, special issue on "Program Generation, Optimization, and Adaptation", Vol. 93, No. 2, pp. 232-275, 2005
- F. Franchetti, Y. Voronenko and M. Püschel "FFT Program Generation for Shared Memory: SMP and Multicore Proc. Supercomputing (SC), 2006

Forward Problem: Match Algorithm to Architecture



Customized, high-performance code for GPU platform
Mapped to target architecture by construction

Automatic Library Generation With Spiral

GPU Architectural Constrains in Formulas

Tile for local memory: Number of Memory Banks, Length of SIMD unit

Parallelize: For large sizes re-distribute data, Memory Hierarchy optimizations

Orchestrate permutation shuffles for nicer data accesses

Vectorize packet exchanges

Algorithm & Program Generation

GPU Code Through Formula Rewriting

Transform user specified: DFT_{64}

Fast algorithm in SPL choices: $(DFT_4 \otimes I_4) (I_4 \otimes L_4^4) D_4^4 (D_2 DFT_4 Z_4^2 \otimes I_4 \otimes I_4) (L_4^4 \otimes I_4) (D_4^4 \otimes I_4) (L_4^4 \otimes I_4)$

Σ-SPL: $\sum_j S_{j,i} DFT_4 \text{diag}(t_j^2 \odot d_j) G_{i,j}$

Optimization at all abstraction levels: parallelization, vectorization, memory architecture optimizations, loop optimizations, rewriting rules, software pipelining, warp scheduling, constant folding, scheduling

CUDA Code:

```
void dft64(float* X) {
    a283 = ((20+i10) + i16);
    s97 = s_data[a283];
    ...
    t137 = (s101 + s103);
    t138 = (s102 + s104);
    t139 = (s101 - s103);
    t140 = (s102 - s104);
    ...
    a296 = ((4+i10) + i16);
    s_data[a296] = (t137 + t141);
}
```

Iteration of this process to search for the fastest