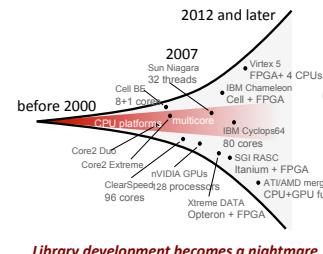


# DFT Transform on the Fermi (GTX480): Automatic Program Generation

Christos Angelopoulos, Franz Franchetti and Markus Pueschel



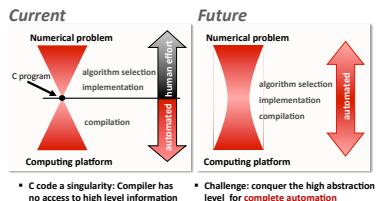
## Problem and Goal



## The Complete Picture

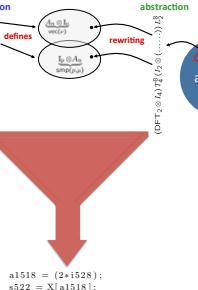
**Spiral:**  
Automating high performance  
library development

## Philosophy

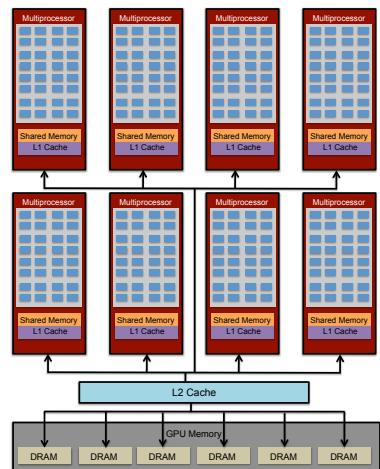


## Forward Problem: Match Algorithm to Architecture

Architectural parameter: Vector length, #processors, ...  
Model: common abstraction = spaces of matching formulas  
Transform: problem size, kernel choice



## GPU Abstraction Model



## Algorithm & Program Generation

### Program Generation in Spiral (Sketched)

Transform  
user specified

$$\text{DFT}_{64} \\ (\text{DFT}_4 \otimes \text{I}_4) (\text{I}_4 \otimes \text{L}_4^{\text{DFT}}) \text{DFT}_{64} \\ (\text{D}_4 \text{DFT}_4 \text{Z}_4^T \otimes \text{I}_4 \otimes \text{I}_4) (\text{I}_4 \otimes \text{L}_4^{\text{DFT}}) \\ (\text{D}_4^T \text{DFT}_4 \text{Z}_4^T \otimes \text{I}_4 \otimes \text{I}_4) (\text{I}_4 \otimes \text{L}_4^{\text{DFT}}) \text{L}_4^{\text{DFT}}$$

Fast algorithm  
in SPL  
many choices

$$\Sigma\text{-SPL}: \\ \sum \sum S_{i,j} \text{DFT}_4 \text{diag}(\text{I}_4^T \odot \text{d}_4) G_{i,j} \\ \sum \sum S_{i,j} \text{DFT}_4 \text{diag}(\text{I}_4^T \odot \text{d}_4) G_{i,j, i,j} \\ \sum \sum S_{i,j} \text{DFT}_4 G_{i,j, i,j}$$

CUDA Code:

$$... \\ \text{int } \text{I17} = (\text{texIdIdx.x / 4}) \% 4; \\ \text{int } \text{I19} = (\text{texIdIdx.y / 4}) \% 4; \\ \text{a89} = ((\text{I17} * \text{I19}) + 193); \\ \text{a41} = +\text{aData}[\text{a89}]; \\ \text{a80} = +\text{aData}[\text{a41}]; \\ \text{a42} = +\text{aData}[\text{a80}]; \\ \text{a43} = (\text{a80} + 32); \\ \text{a44} = (\text{a80} + 15); \\ \text{a92} = (\text{a80} + 96); \\ \text{a44} = +\text{aData}[\text{a92}]; \\ ...$$

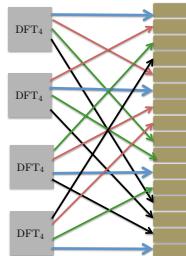
Iteration of this process to search for the fastest

### GPU Stockham Algorithm

Original Stockham :  $\text{DFT}_{In} \rightarrow \prod_{i=0}^{n-1} (\text{DFT}_i \otimes \text{I}_{In-i}) \text{I}_i^{\text{DFT}} (\text{I}_i^{\text{In}-i} \otimes \text{I}_i)$

GPU Stockham:  $\text{DFT}_{In} \rightarrow \prod_{i=0}^{n-1} (((\text{D}_i^T \cdot \text{DFT}_i \cdot \text{Z}_i^T \otimes \text{I}_p \otimes \text{I}_p \otimes \text{I}_{\frac{In}{Mp}}) \text{I}_i^{\text{DFT}}) \\ \otimes \text{I}_p \otimes \text{I}_p \otimes \text{I}_{\frac{In}{Mp}}) \text{I}_i^{\text{DFT}} (\text{I}_i^{\text{In}-i} \otimes \text{I}_i)$

Idea :  $\text{I}_n \otimes \text{DFT}_r = \text{I}_n \otimes_i \text{D}_r^i \cdot \text{DFT}_r \cdot \text{Z}_r^i$



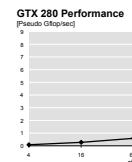
- uses DFT cyclic shift property
- avoids bank conflicts by adding extra operations
- avoids extra operations by merging twiddle diagonals from different stages

## Code Example

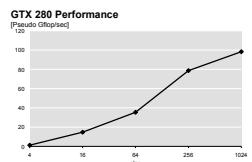
```
global<__shared__> void sub(float* X)
{
    ...
    a1518 = (2 * a1528);
    s522 = X[a1518];
    a1519 = (a1518 + 1);
    s523 = X[a1519];
    a1520 = (a1518 + 256);
    s524 = X[a1520];
    a1521 = (a1518 + 257);
    s525 = X[a1521];
    t1122 = (s522 + s525);
    t1123 = (s523 + s525);
    t1124 = (s522 - s524);
    t1125 = (s523 - s525);
    a1522 = (a1518 + 128);
    ...
}
```

## Results

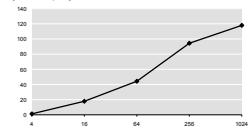
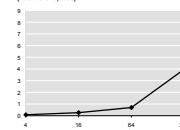
### SM to SM Latency



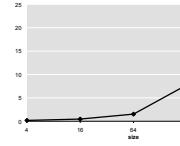
### Throughput



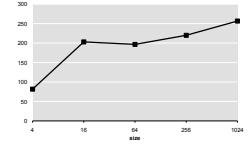
### Quadro FX 5800 Performance



### GTX 480 Performance



### GTX 480 Performance



## Future Work

- N-Dimensional DFTs
- Big (Global Memory) DFTs from GPU Memory to GPU Memory
- Overlapping DFT algorithms from Host to Device
- Work still in progress...

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
  - F. Franchetti, M. Püschel, Y. Voronenko, K. Chen, R. W. Johnson and N. Rizzoli. *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