

FFT Compiler: From Math to Efficient Hardware

James C. Hoe
Department of ECE
Carnegie Mellon University

joint work with Peter A. Milder, Franz Franchetti, and
Markus Pueschel in the SPIRAL project

with support from NSF ACR-0234293, ITR/NGS-0325687, DARPA NBCH-105000, Intel

The SPIRAL Project

- ◆ High performance implementations of linear DSP transforms (DFT, DCT, DWT, filters, etc) are an important class of design problems
- ◆ Hand design and tuning is tricky and expensive
 - needs both math and implementation knowledge
 - time-consuming and tedious
 - needs to repeat effort for every new context
- ◆ **SPIRAL research goal:** A flexible push-button design generation framework that produces SW and HW implementations equal or better than expert hand design
Today I focus on FFT for HW

Why we can do better than hand design

- ◆ SPIRAL is only focused on linear DSP transforms
- ◆ These transforms are highly structured, highly regular and very well understood mathematically
- ◆ Algorithmic implementations of a transform can be enumerated following a known set of rules
- ◆ For a given objective function and mapping target, a computer generates a solution at least as good as the best human effort—by trying enough implementations

Outline

- ◆ SPIRAL Formula Framework
- ◆ SPIRAL for HW FFT cores
- ◆ Conclusions

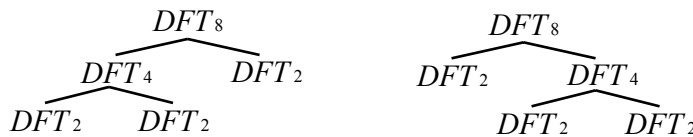
“Fast” Fourier Transform Algorithms

- ◆ Recursively factorize by the Cooley-Tukey rule until only leaf cases remain (e.g. DFT_r for radix-r)

$$DFT_8 = (DFT_2 \otimes I_4) D_2^8 (I_2 \otimes DFT_4) L_2^8$$

$$= (DFT_2 \otimes I_4) D_2^8 (I_2 \otimes ((DFT_2 \otimes I_2) D_2^4 (I_2 \otimes DFT_2) L_2^4)) L_2^8$$

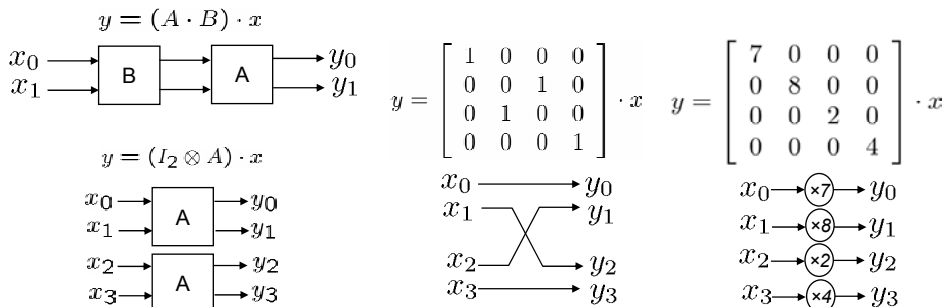
- ◆ Exponential number of alternatives



- ◆ Each ruletree corresponds a different algorithm
- ◆ All cost $O(N \log(N))$

Formula to HW (Combinational)

- ◆ Given $y = M \cdot x$ where M is:
 - $M = A \cdot B$ apply B , then A
 - M is a permutation permute x
 - $M = I_n \otimes A$ apply A , n times in parallel
 - M is a diagonal scale x



How about good HW?

- ◆ Matrix formulas have a natural mapping to dataflow and hence combinational datapath
- ◆ However, real hardware designs must fit a given resource constraint
 - ⇒ *sequential datapath that reuse available HW*
 - identify repeated kernels
 - instantiate kernels under resource constraints
 - schedule computation to reuse instantiated kernels

We want to do the analysis, mapping and scheduling at formula level, with high-level algorithm knowledge

Regular Structure for HW

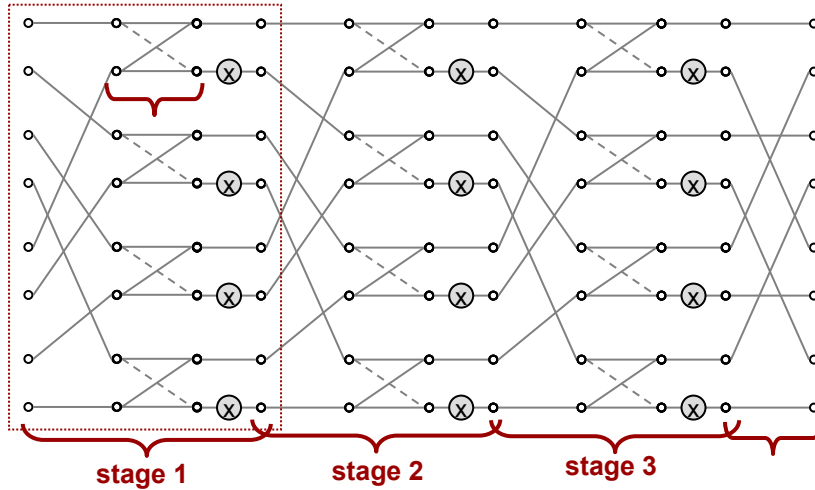
- ◆ Simple regular structure embodied in Pease FFT

$$\text{DFT}_{2^k} = R_{2^k} \left(\prod_{i=0}^{k-1} \underbrace{T_i(I_{2^{k-1}} \otimes F_2)}_{\text{kernel}} L_{2^{k-1}}^{2^k} \right)$$

- ◆ Example:

$$\text{DFT}_8 = R_8 \left(\underbrace{T_0(I_4 \otimes F_2)}_{\text{kernel}} L_4^8 \right) \left(\underbrace{T_1(I_4 \otimes F_2)}_{\text{kernel}} L_4^8 \right) \left(\underbrace{T_2(I_4 \otimes F_2)}_{\text{kernel}} L_4^8 \right)$$

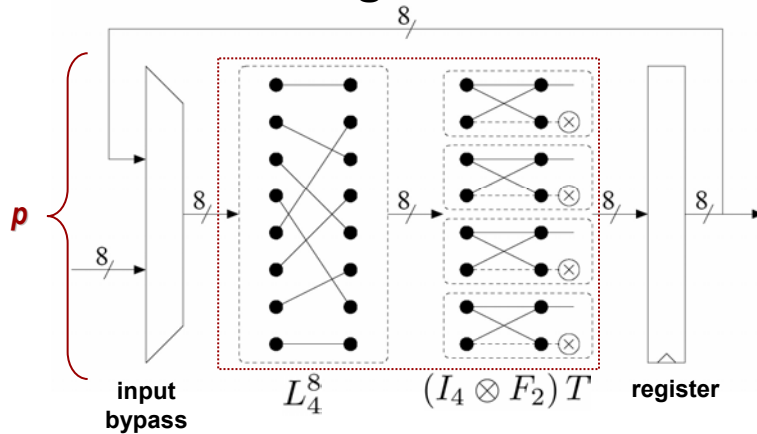
Pease DFT Example: DFT₈



$$DFT_8 = R_8 \cdot T_2(I_4 \otimes F_2)L_4^8 \cdot T_1(I_4 \otimes F_2)L_4^8 \cdot T_0(I_4 \otimes F_2)L_4^8$$

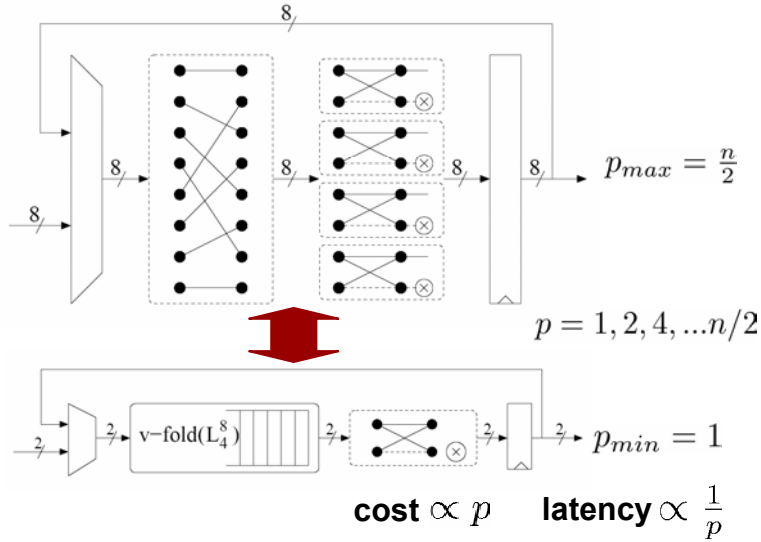
(formula is applied from right to left)

Horizontal folding



- ◆ a baseline datapath for DFT (*sans bit-reverse*)
- ◆ degree of freedom: vertical parallelism
 - parameter p

Vertical (V-)folding according to p



Fine-grained control over cost/latency tradeoff

DFTgen [Milder, et al., DAC'05, FPGA'06]

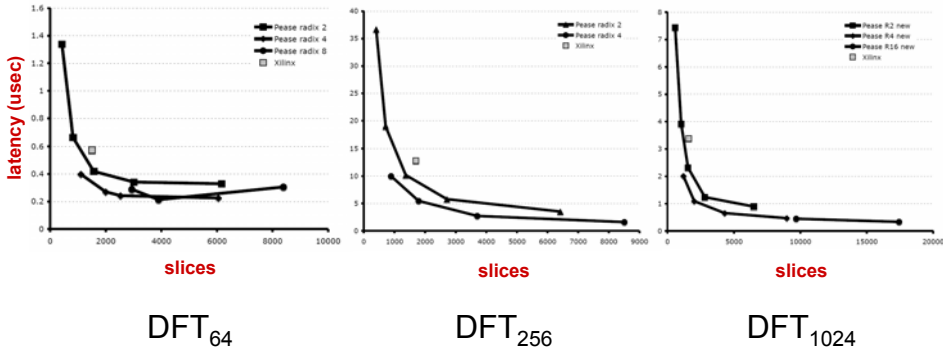
Input form

	parameter	value	range	usage
functionality	n	2048	4-16384	DFT size (?)
	data width	16	6-32	datapath width, fixed point precision (?)
	twiddle width	16	2-16	twiddle factor bitwidth (?)
	data ordering	bit-reversed out		bit-reversed input or output (?)
control area/latency tradeoff	scaling mode	unscaled		unscaled or scaled arithmetic (?)
	p	1	1-1024	degree of parallelism (?)
	twiddle storage	block RAM		store twiddles in dist. RAM or BRAM (?)
	FIFO threshold	Select	2-1024	FIFO of \geq depth will be built as BRAM. (?)

Core characteristics and resource usage (dynamically updated)

	char./res.	value	how obtained
resource model	cycles	11297	exact formula
	slices	876	estimated with linear fit formula
	BRAM	44	exact formula
	multipliers	4	exact formula
	num. error		to be implemented

Cost Performance Tradeoff

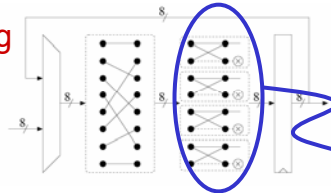


Structures of Interest for Parameterized Cores

- ◆ Simple regular structure embodied in formula

$$DFT_{2^k} = R_{2^k} \left(\prod_{i=0}^{s-1} T_i \left(I_{2^{k-1}} \otimes F_2 \right) L_{2^{k-1}} \right)$$

- ◆ Product: Horizontal folding



- ◆ Tensor Product: parallelism or streaming vertical folding



- ◆ Stride permutation: rewiring or buffer/delay

Applicability to other transforms?

◆ DFT radix 2
$$R_{2^k} \prod_{i=0}^{k-1} \left[T_i \left(I_{2^{k-1}} \otimes DFT_2 \right) L_{2^{k-1}}^{2^k} \right]$$

◆ DFT radix 2^r
$$R_{2^k} \prod_{i=0}^{k/r-1} \left[T_i \left(I_{2^{k-r}} \otimes DFT_{2^r} \right) L_{2^{k-r}}^{2^k} \right]$$

◆ 2-D DFT_{nxn}
$$\prod_{i=0}^1 \left[L_n^{n^2} \left(I_n \otimes DFT_n \right) \right]$$

◆ WHT
$$\prod_{i=0}^{k/r-1} \left[\left(I_{2^{k-r}} \otimes WHT_{2^r} \right) L_{2^{k-r}}^{2^k} \right]$$

◆ DCT (type II)
$$DP \prod_{i=0}^{k-1} \left[A_{k-i} L_2^{2^k} \right] L_{2^{k-1}}^{2^k} L_{2^{k-1}}^{2^k} P_{2^k}^H$$

Conclusions

- ◆ Mathematical approach to DSP transform implementation
- ◆ Encapsulating domain knowledge in a domain specific tool for **100% automatic** design generation and optimization
- ◆ Generalizable to other linear DSP transforms
- ◆ Thank you (Arvind)