Domain Specific Optimizations

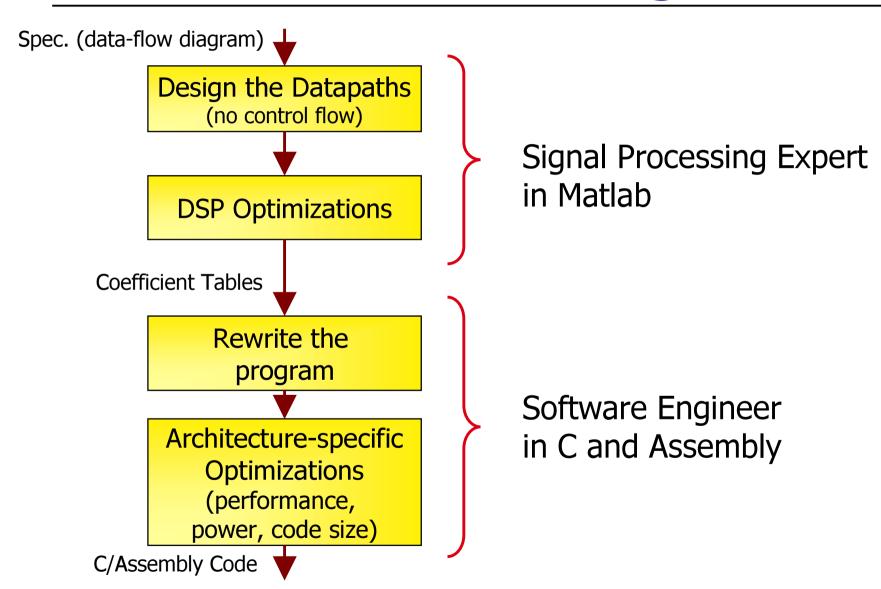
Saman Amarasinghe and William Thies Massachusetts Institute of Technology

> PACT 2003 September 27, 2003

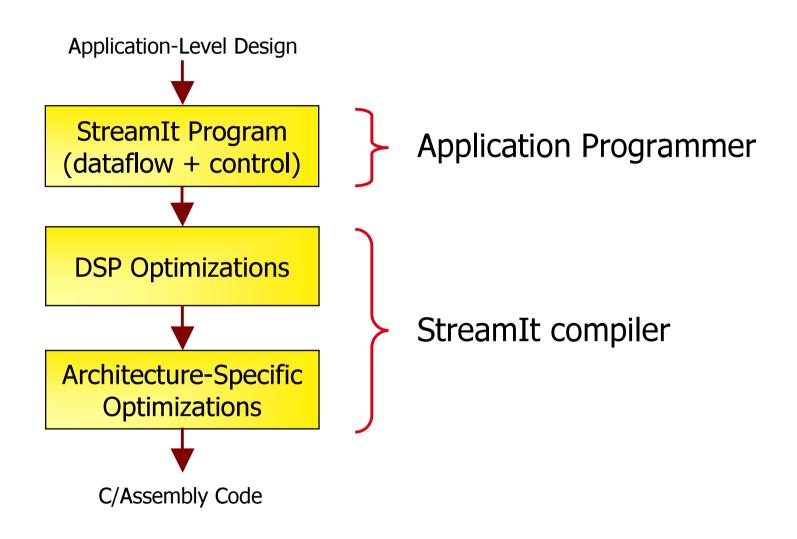
Schedule

1:30-1:40	Overview (Saman)
1:40-2:20	Stream Architectures (Saman)
2:20-3:00	Stream Languages (Bill)
3:00-3:30	Break
3:30-3:55	Stream Compilers (Saman)
3:55-4:20	Domain-specific
	Optimizations (Saman)
4:20-5:00	Scheduling Algorithms (Bill)

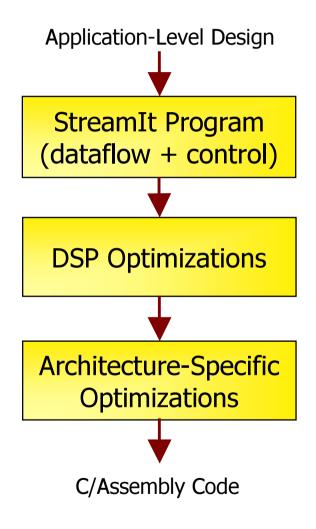
Conventional DSP Design Flow



Design Flow with StreamIt



Design Flow with StreamIt



- Benefits of programming in a single, high-level abstraction
 - Modular
 - Composable
 - Portable
 - Malleable
- The Challenge: Maintaining Performance
 - Replacing Expert DSP Engineer
 - Replacing Expert Assembly Hacker

Our Focus: Linear Filters

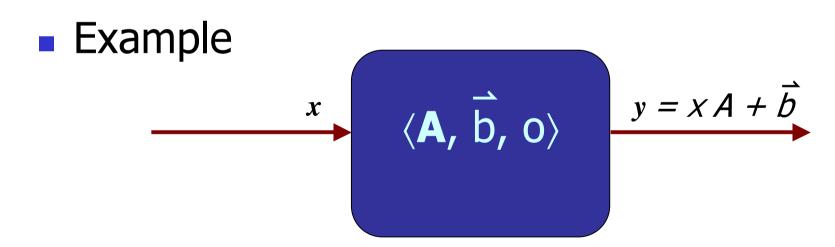
- Most common target of DSP optimizations
 - FIR filters
 - Compressors
 - Expanders
 - DFT/DCT

Output is weighted sum of inputs

- Example optimizations:
 - Combining Adjacent Nodes
 - Translating to Frequency Domain

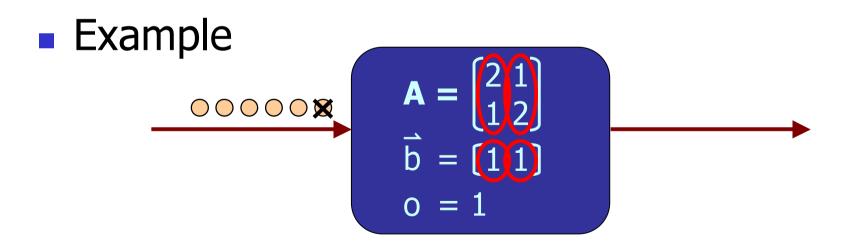
Representing Linear Filters

- A linear filter is a tuple $\langle \mathbf{A}, \overline{\mathbf{b}}, \mathbf{o} \rangle$
 - **A**: matrix of coefficients
 - b: vector of constants
 - o: number of items popped



Representing Linear Filters

- A linear filter is a tuple $\langle \mathbf{A}, \overline{\mathbf{b}}, \mathbf{o} \rangle$
 - **A**: matrix of coefficients
 - b: vector of constants
 - o: number of items popped



Extracting Linear Representation

```
work peek N pop 1 push 1 {
  float sum = 0;
  for (int i=0; i<N; i++) {
    sum += h[i]*peek(i);
  }
  push(sum);
  pop();
}</pre>
Linear

Dataflow

Analysis

    (A, b, o)
```

- Resembles constant propagation
- Maintains linear form $\langle \vec{v}, \vec{b} \rangle$ for each variable
 - Peek expression: generate fresh v̄
 - Push expression: copy v into A
 - Pop expression: increment o

Optimizations using Linear Analysis

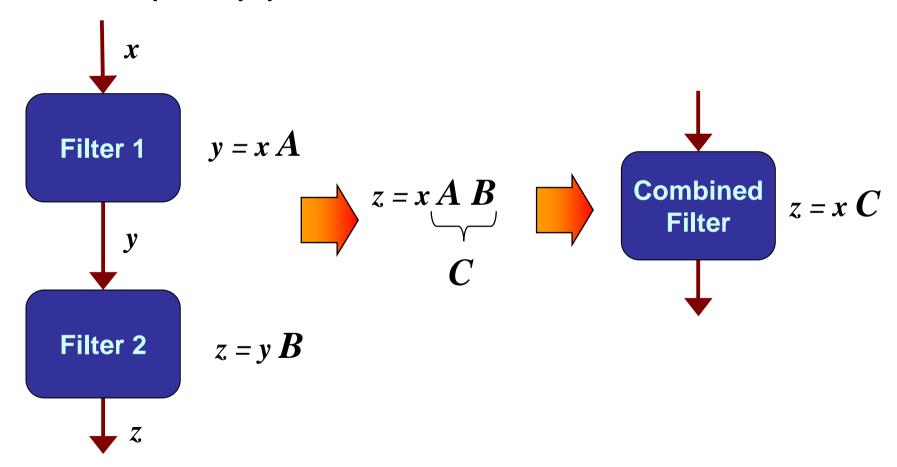
1) Combining adjacent linear structures

2) Shifting from time to the frequency domain

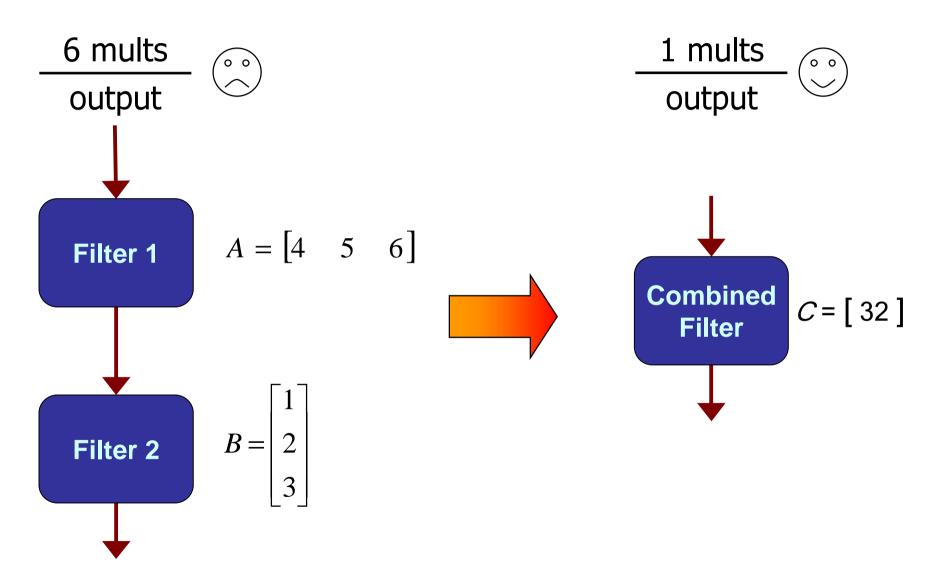
3) Selection of 'optimal' set of transformations

1) Combining Linear Filters

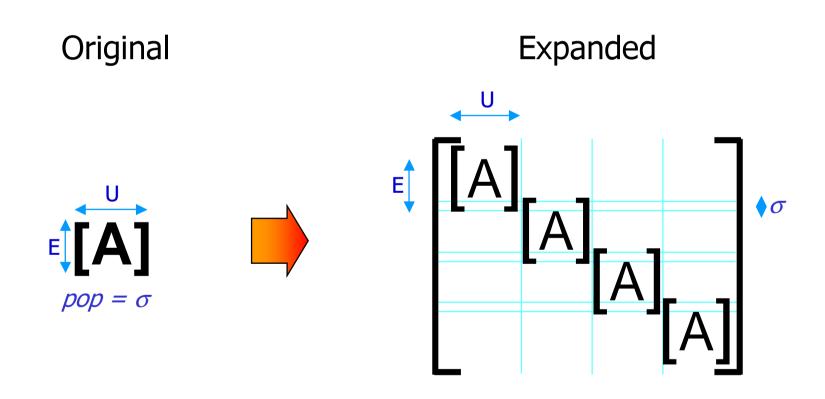
- Pipelines and splitjoins can be collapsed
- Example: pipeline

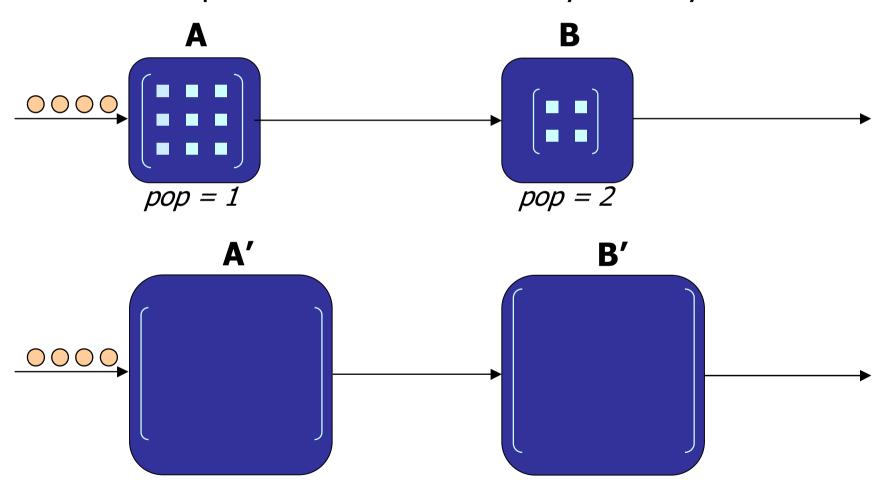


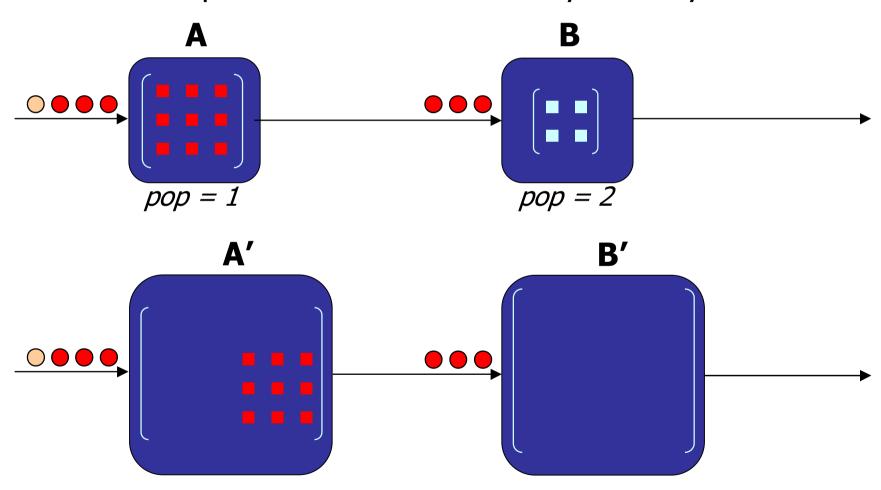
Combination Example

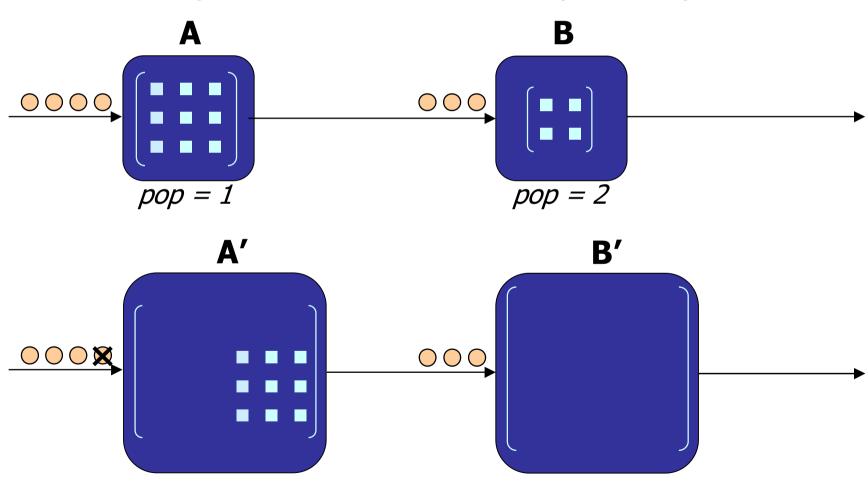


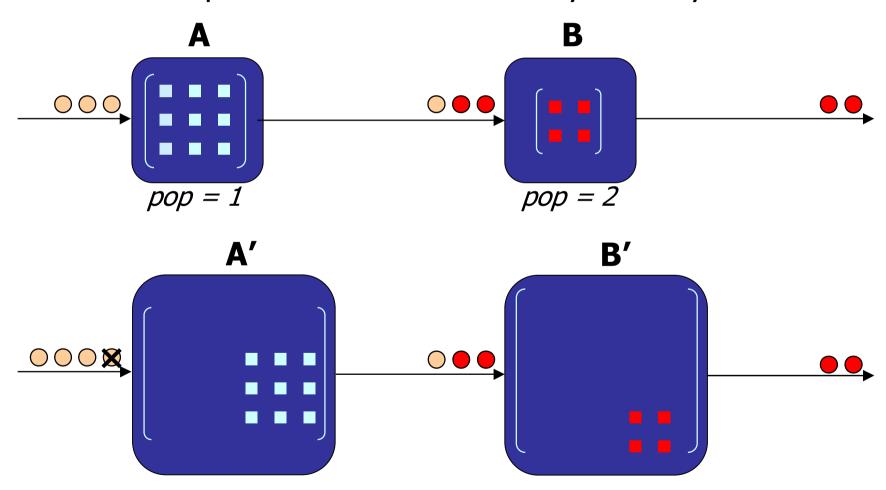
Linear Expansion

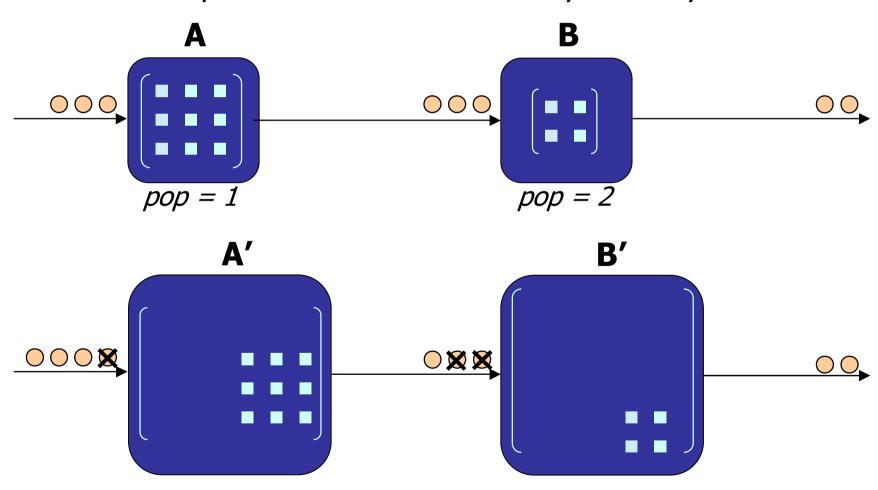


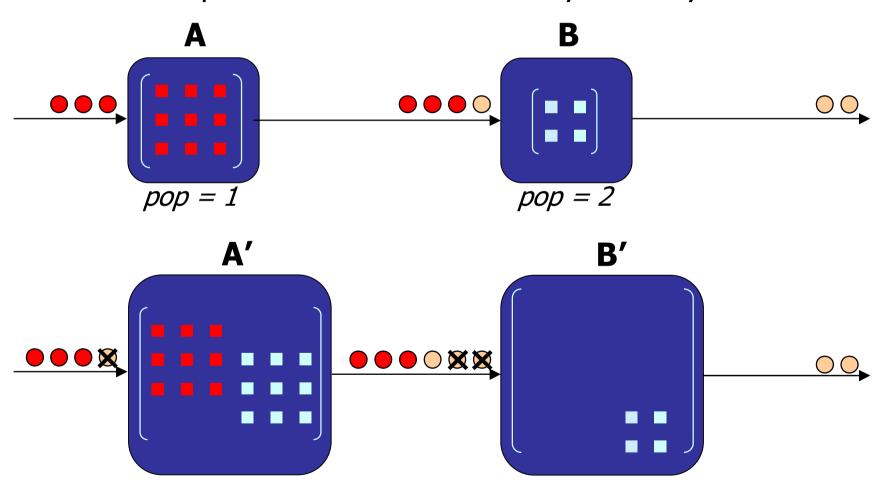


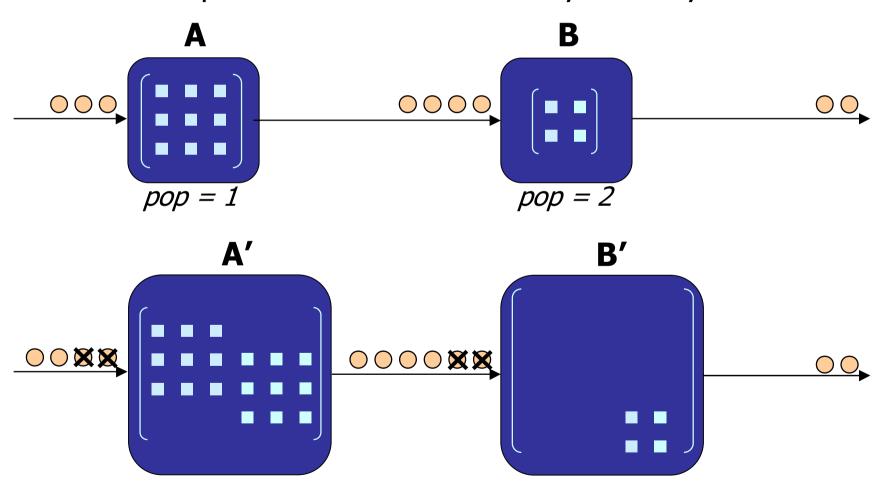


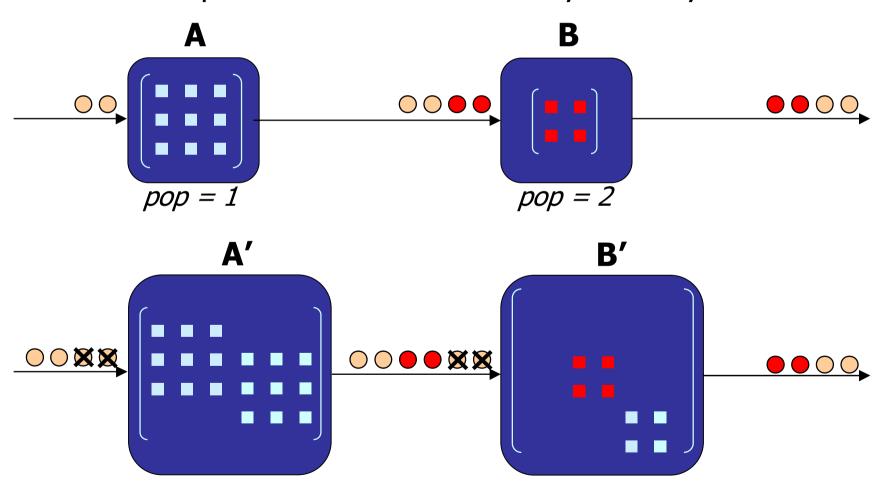


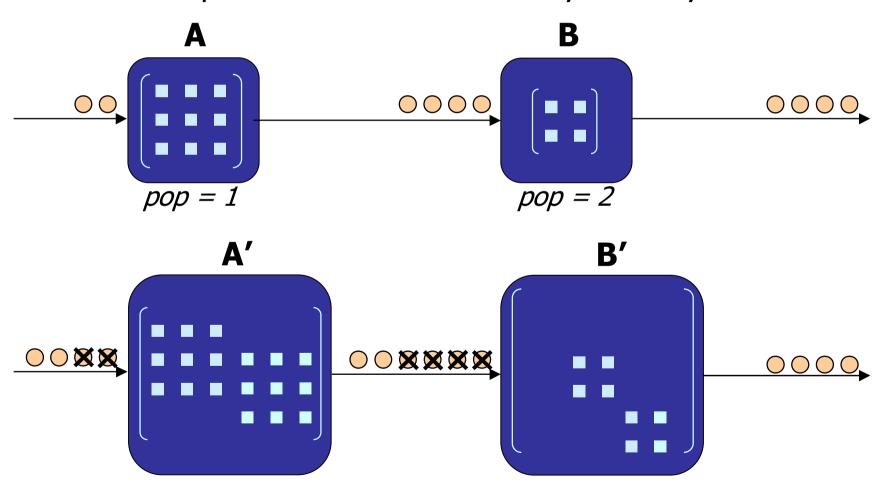


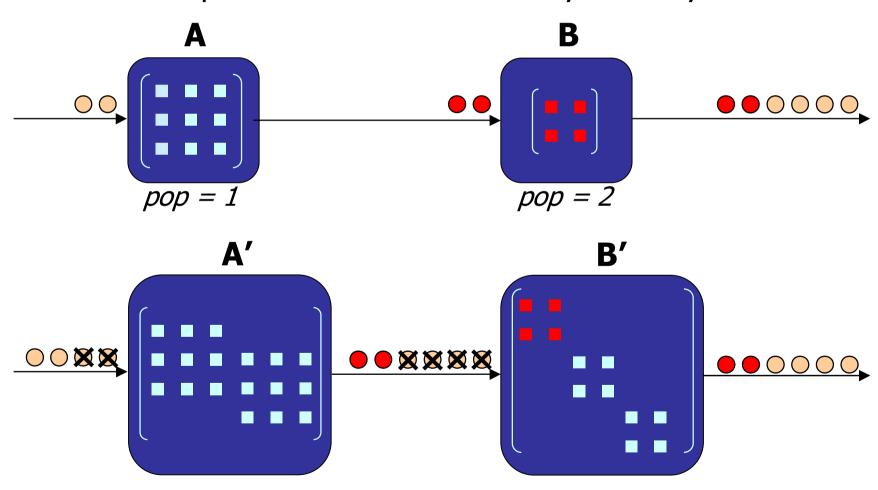


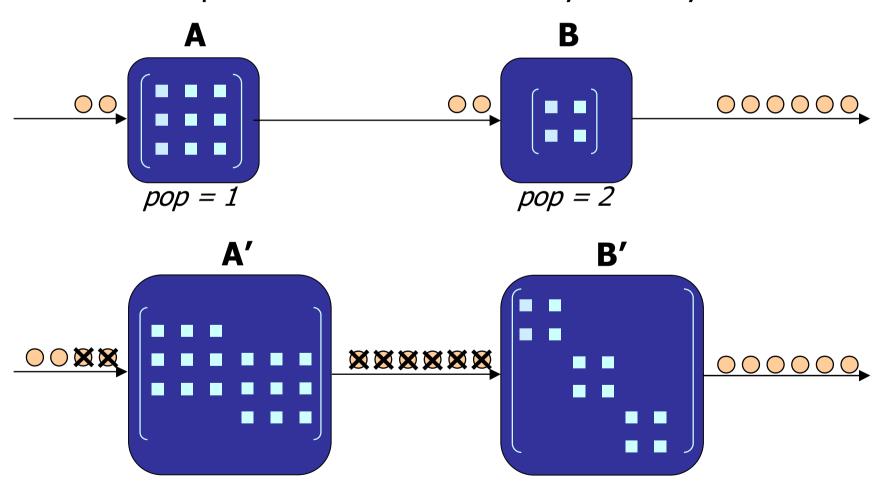




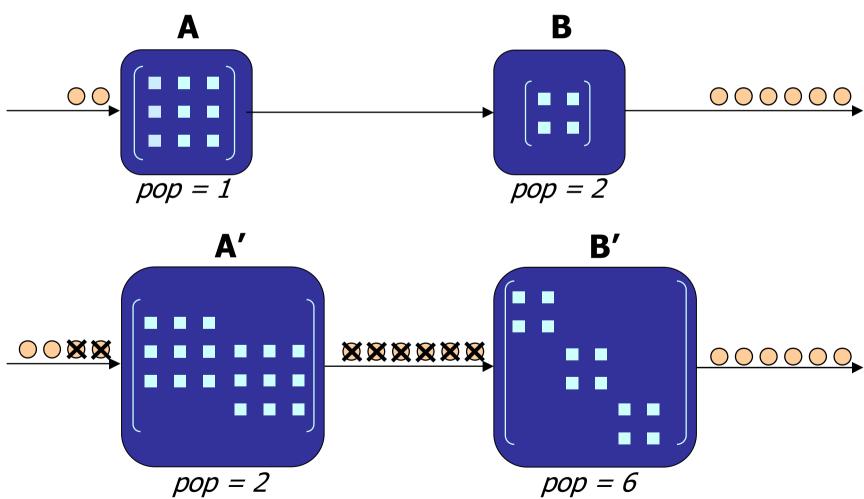






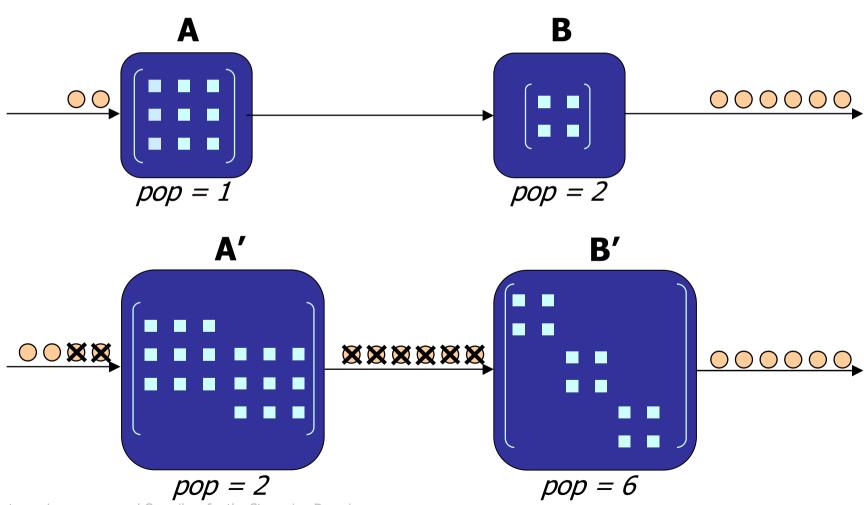


Need to "expand" matrices to a steady-state cycle:



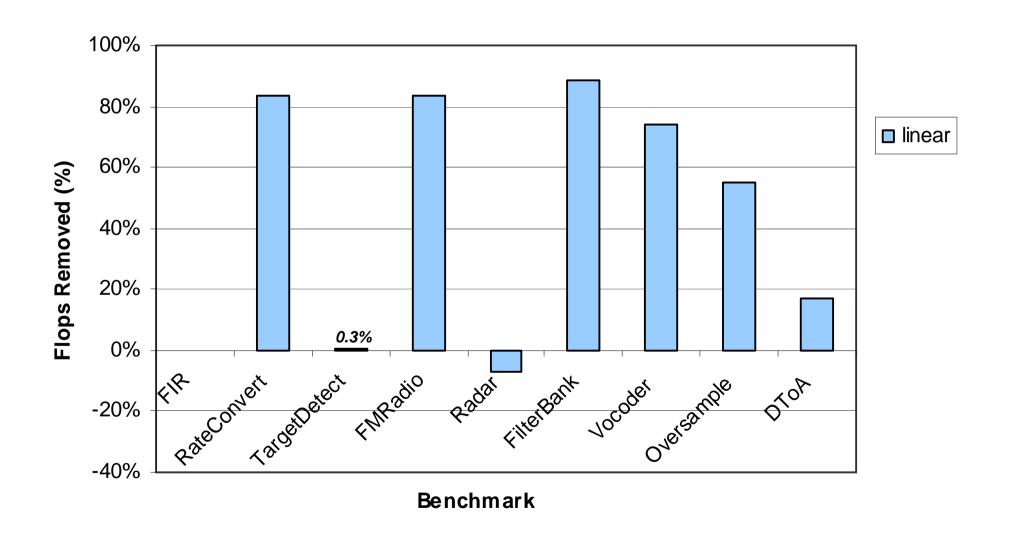
Architectures, Languages, and Compilers for the Streaming Domain PACT 2003 Tutorial - Saman Amarasinghe, William Thies - MIT CSAIL

Expanded dimensions match



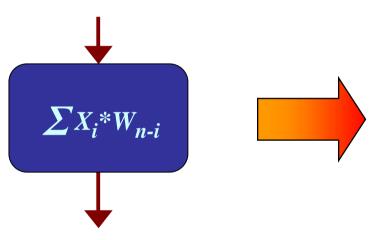
Architectures, Languages, and Compilers for the Streaming Domain PACT 2003 Tutorial - Saman Amarasinghe, William Thies - MIT CSAIL

Floating-Point Operations Reduction



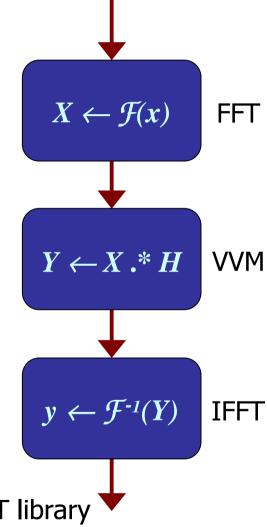
2) From Time to Frequency Domain

 Convolutions can be done cheaply in the Frequency Domain



- Painful to do by hand
 - Blocking

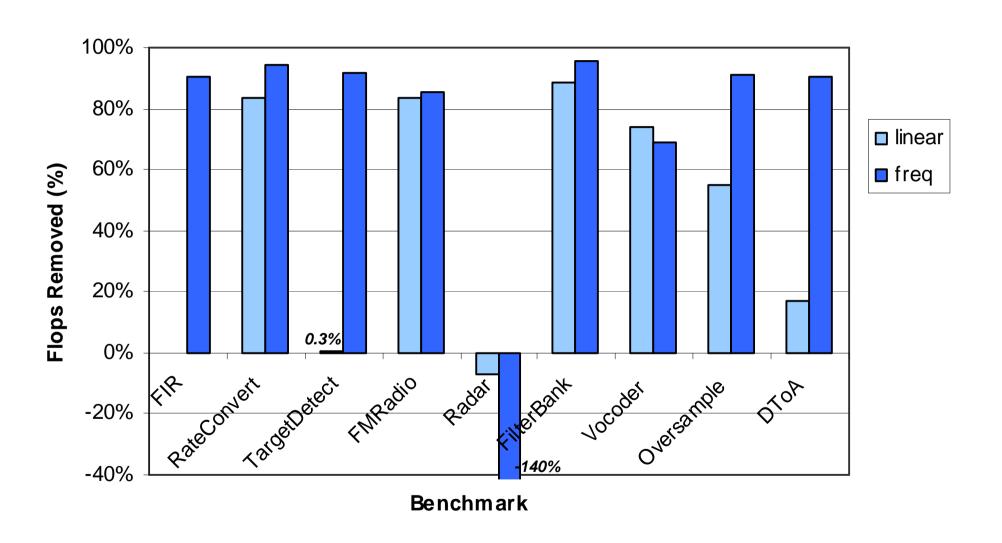
- Multiple outputs
- Coefficient calculations
 Interfacing with FFT library
- Startup



Generic Freq. Implementation

```
float \rightarrow float pipeline optimizedFreq (A, \vec{b}, e, o, u) {
  add float → float filter {
    N \leftarrow 2\lceil \lg(2e) \rceil
    m \leftarrow N - 2e + 1
    partials \leftarrow \text{new array}[0 \dots e-2, 0 \dots u-1]
                                                                             work peek r pop r push u * r {
    r \leftarrow m + e - 1
                                                                               \vec{x} \leftarrow pop(0 \dots m + e - 2)
                                                                               \vec{X} \leftarrow \mathbf{FFT}(N, \vec{x})
    init {
                                                                               for j = 0 to u - 1 {
      for i = 0 to u - 1
                                                                                 \vec{Y}[*,j] \leftarrow \vec{X}.*\vec{H}[*,j]
        \vec{H}[*,j] \leftarrow \mathbf{FFT}(N,A[*,u-1-j])
                                                                                 \vec{y}[*,j] \leftarrow \mathbf{IFFT}(N, \vec{Y}[*,j])
    prework peek r pop r push u*m {
                                                                               for i=0 to e-1
      \vec{x} \leftarrow pop(0 \dots m + e - 2)
                                                                                  for j = 0 to u - 1 {
      \vec{X} \leftarrow \mathbf{FFT}(N, \vec{x})
                                                                                    push(\vec{y}[i,j] + partials[i,j])
      for i = 0 to u - 1 {
                                                                                   partials[i,j] \leftarrow \vec{y}[m+e-1+i,j]
        \vec{Y}[*,j] \leftarrow \vec{X}.*\vec{H}[*,j]
        \vec{y}[*,j] \leftarrow \mathbf{IFFT}(N,\vec{Y}[*,j])
                                                                               for i=0 to m-1
        partials[*,j] \leftarrow \vec{y}[m+e-1...m+2e-3,j]
                                                                                  for i = 0 to u - 1
                                                                                   push(\vec{y}[i+e-1,j]+\vec{b}[j])
      for i=0 to m-1
        for i = 0 to u - 1
          push(\vec{y}[i+e-1, j] + \vec{b}[j])
                                                                           add Decimator (o, u)
```

Floating-Point Operations Reduction



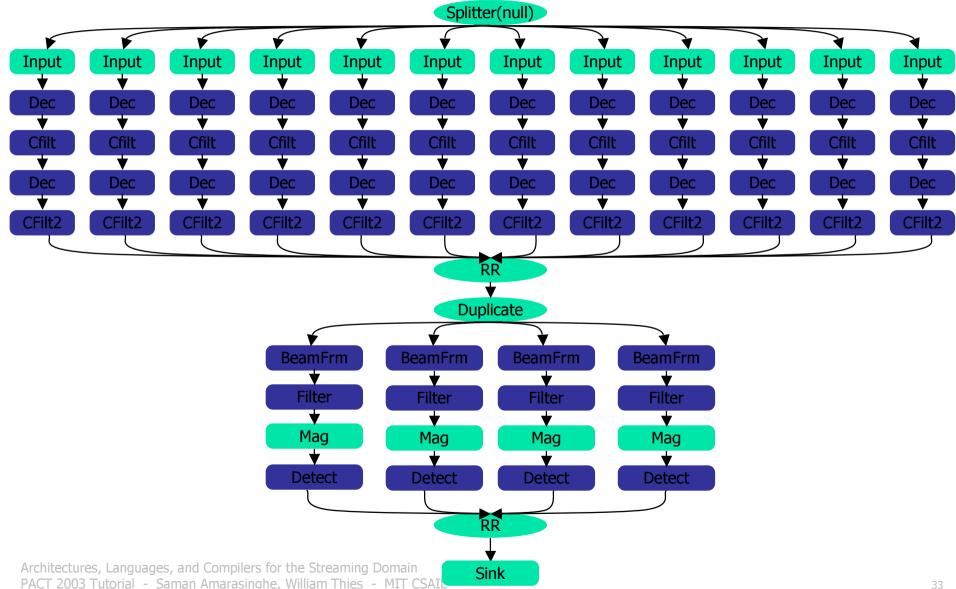
3) Transformation Selection

- When to apply what transformations?
 - Linear filter combination can increase the computation cost
 - Shifting to the Frequency domain is expensive for filters with pop > 1
 - Compute all outputs, then decimate by pop rate
 - Some expensive transformations may later enable other transformations, reducing the overall cost

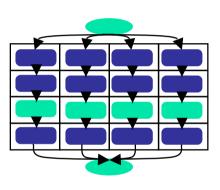
Selection Algorithm

- Estimate minimal cost for each structure:
 - Linear combination
 - Frequency translation
- Cost function based on profiler feedback

- No transformation
 - If hierarchical, consider all possible groupings of children
- Overlapping sub-problems allows efficient dynamic programming search

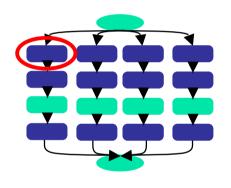


First compute cost of individual filters:

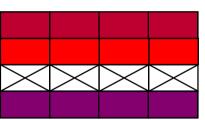


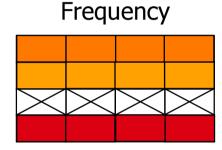
First compute cost of individual filters:







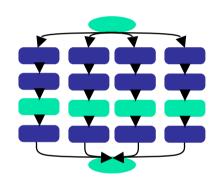


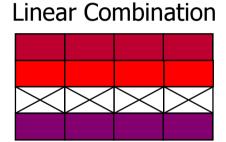


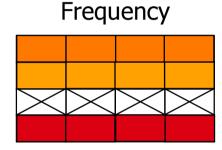


First compute cost of individual filters:







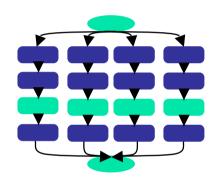




1x1

Then, compute cost of 1x2 nodes:



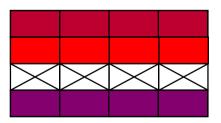


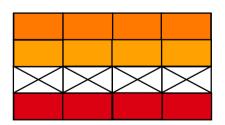
Linear Combination

Frequency

No Transform

1x1



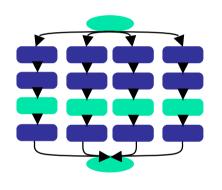


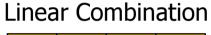


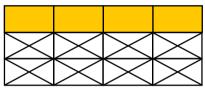
Then, compute cost of 1x2 nodes: low high No Transform **Linear Combination** Frequency min= 1x1

Then, compute cost of 1x2 nodes:

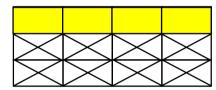








Frequency

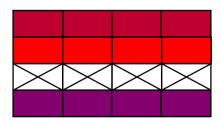


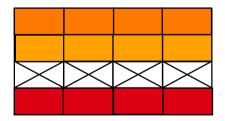
No Transform



1x2

1×1







Continue with 1x3 2x1 3x1 4x1

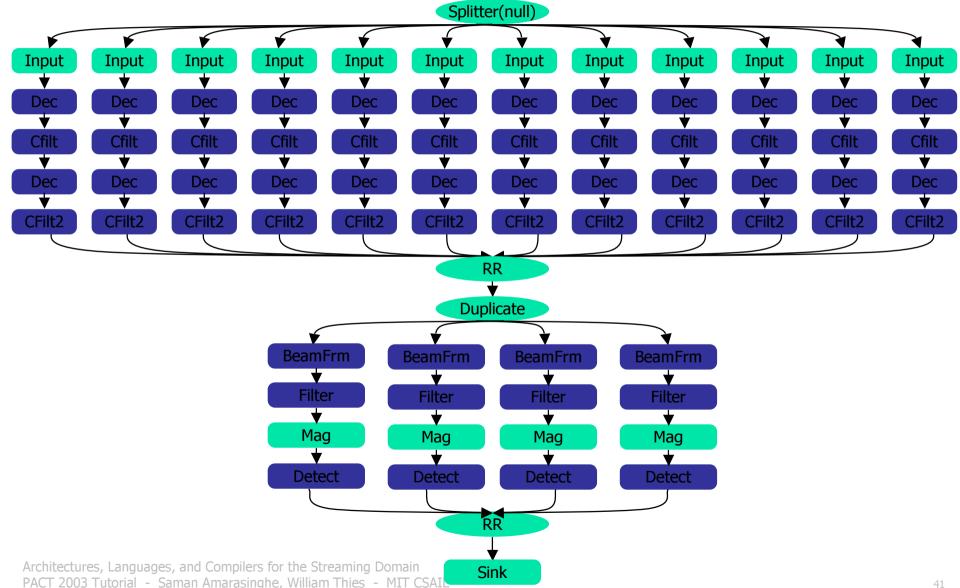
1x4 2x2 3x2 4x2

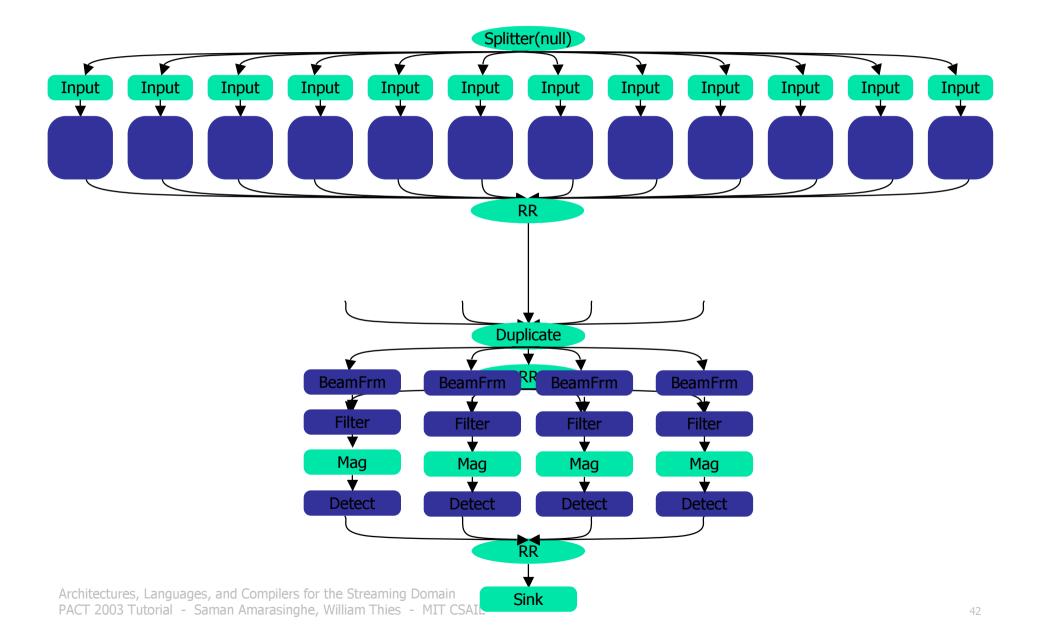
2x3 3x3 4x3

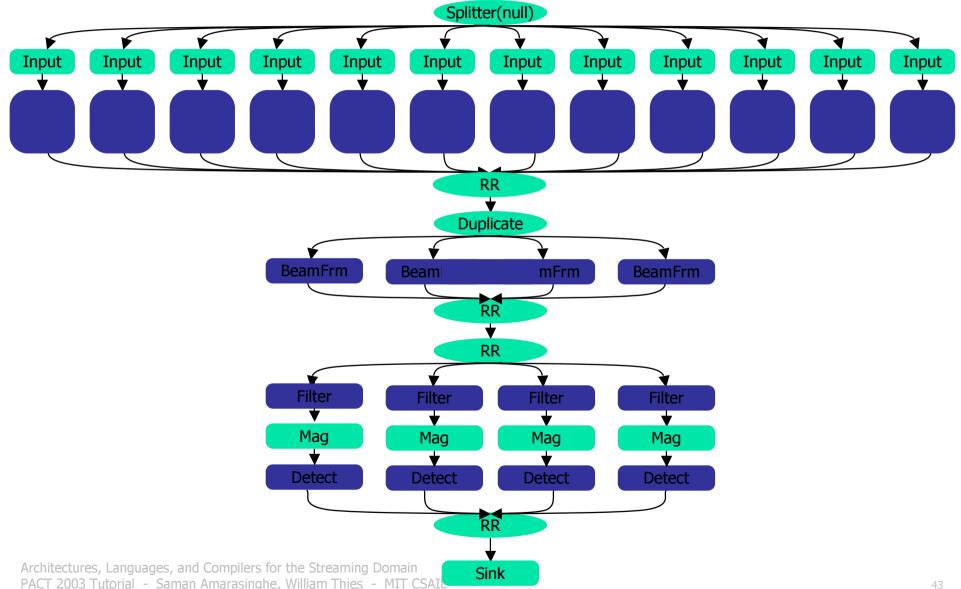
2x4 3x4(4x4)



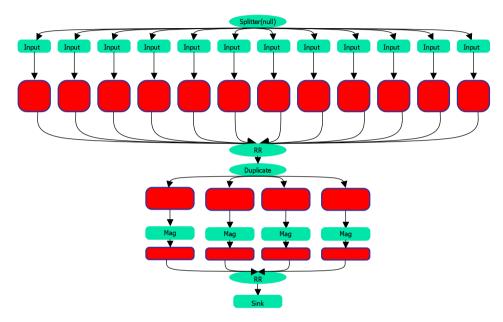
Overall solution







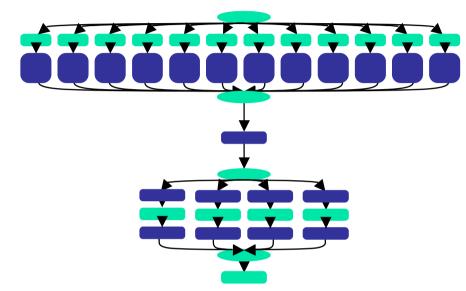
Radar







2.4 times as many FLOPS

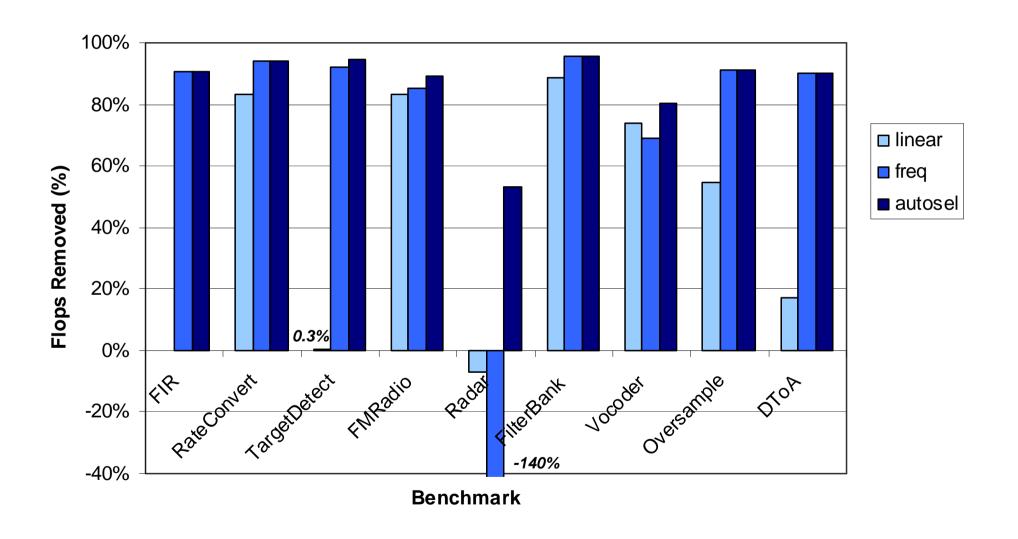


Using Transformation Selection



half as many FLOPS

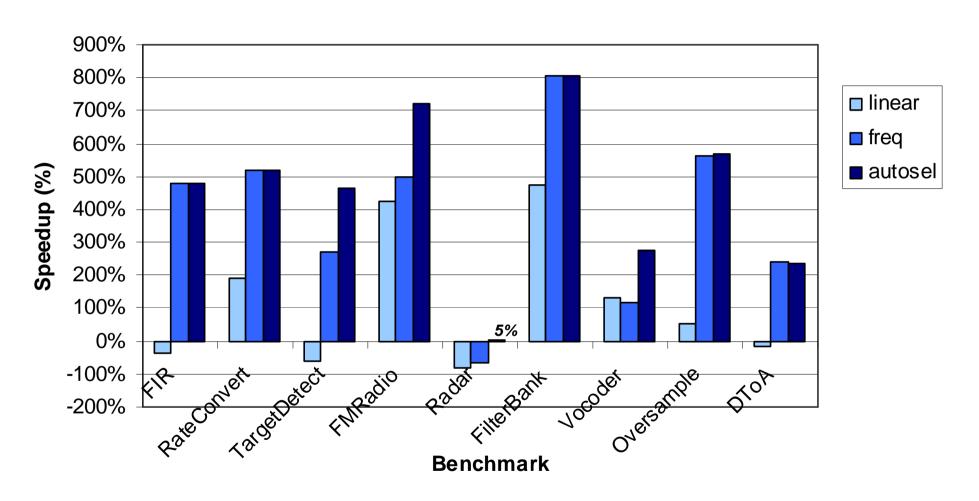
Floating-Point Operations Reduction



Experimental Results

- Fully automatic implementation
 - StreamIt compiler
- StreamIt to C compilation
 - FFTW for shifting to the frequency domain
- Benchmarks all written in StreamIt
- Measurements
 - Dynamic floating-point instruction counting
 - Speedups on a general purpose processor

Execution Speedup



On a Pentium IV

Related Work

- SPIRAL/SPL (Püschel et. al)
 - Automatic derivation of DSP transforms
- FFTW (Friego et. al)
 - Wicked fast FFT
- ADE (Covell, MIT PhD Thesis, 1989)
- Affine Analysis (Karr, Acta Informatica, 1976)
 - Affine relationships among variables of a program
- Linear Analysis (Cousot, Halbwatchs, POPL, 1978)
 - Automatic discovery of linear restraints among variables of a program

Conclusions

- A DSP Program Representation: Linear Filters
 - A dataflow analysis that recognizes linear filters
- Three Optimizations using Linear Information
 - Adjacent Linear Structure Combination
 - Time Domain to Frequency Domain Transformation
 - Automatic Transformation Selection
- First Step in Replacing the DSP Engineer from the Design Flow
 - On the average 90% of the FLOPs eliminated
 - Average performance speedup of 450%
- StreamIt: A Unified High-level Abstraction for DSP Programming
 - Increased abstraction does not have to sacrifice performance

http://cag.lcs.mit.edu/linear/