

High-Productivity Stream Programming for High-Performance Systems

Rodric Rabbah, Bill Thies, Michael Gordon, Janis Sermulins,
and Saman Amarasinghe

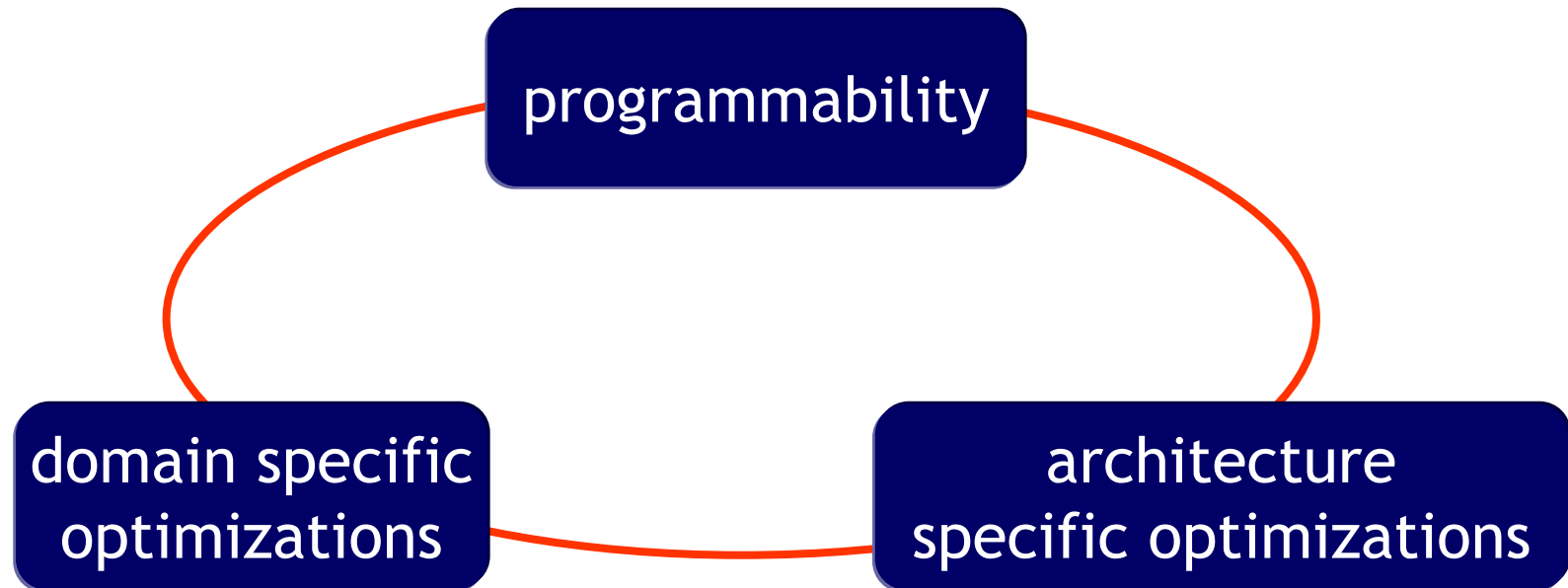
Massachusetts Institute of Technology

The logo for StreamIt, featuring the word "StreamIt" in a blue, stylized font. A red horizontal line with an arrowhead pointing to the right is positioned above the "m" and "i" characters.

<http://cag.lcs.mit.edu/streamit>

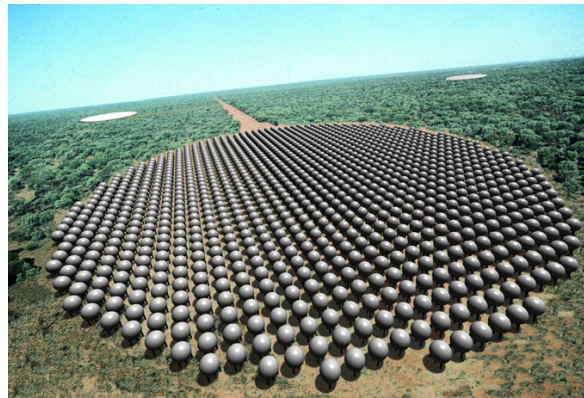
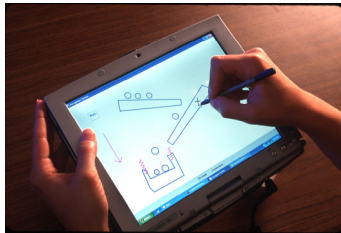
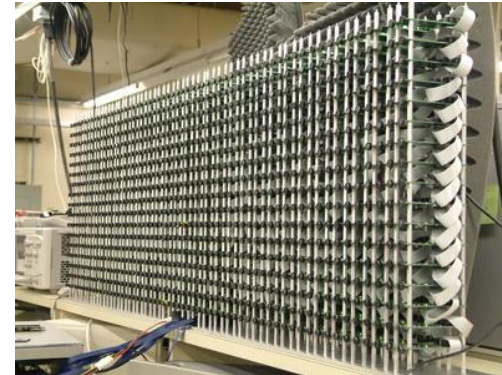
The StreamIt Vision

- Boost productivity, enable faster development and rapid prototyping



- Simple and effective optimizations for streams
- Targeting tiled architectures, clusters of workstations, DSPs, and traditional uniprocessors

Why an Emphasis on Streaming?

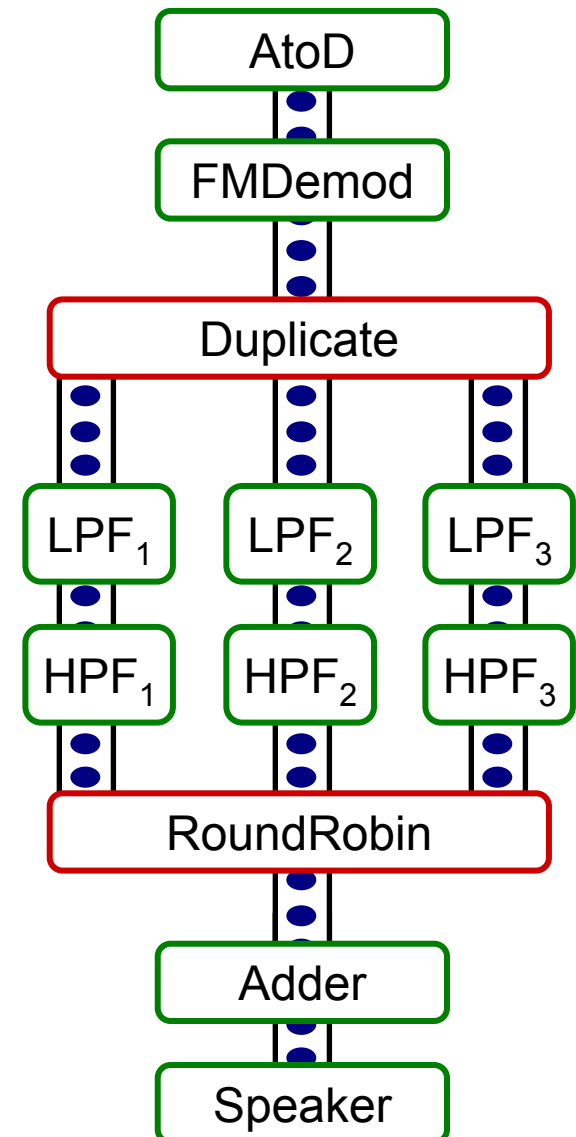


Streaming in other Domains as well

- Cryptography
 - Databases
 - Face recognition
 - Network processing and security
 - Scientific codes
 - ...
-
- Attractive programming model because of a simple mapping from specification to implementation

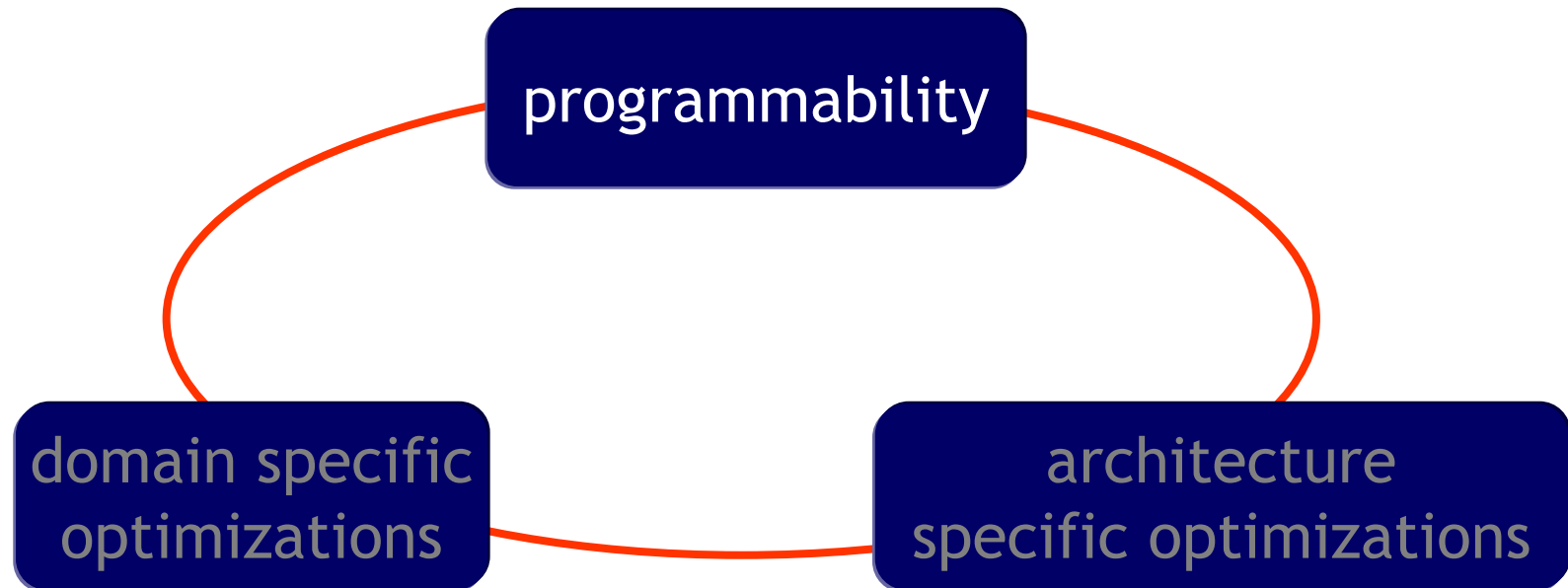
Properties of Stream Programs

- Mostly regular and repeating computation
- Parallel, independent computation with explicit communication
- Amenable to aggressive compiler optimizations
[ASPLOS '02, PLDI '03, LCTES'03, LCTES '05]



The StreamIt Vision

- Boost productivity, enable faster development and rapid prototyping



- Simple and effective optimizations for streams
- Targeting tiled architectures, clusters of workstations, DSPs, and traditional uniprocessors

Programming in StreamIt

```
void->void pipeline FMRadio(int N, float freq1, float freq2) {
```

```
  add AtoD();
```

- Natural correspondence
between text and
application graph

AtoD

```
  add FMDemod();
```

FMDemod

```
  add splitjoin {
    split duplicate;
    for (int i=0; i<N; i++) {
      add pipeline {
```

```
        add LowPassFilter(freq1 + i*(freq2-freq1)/N);
```

LPF₁

LPF₂

LPF₃

```
        add HighPassFilter(freq2 + i*(freq2-freq1)/N);
```

HPF₁

HPF₂

HPF₃

```
      }
    }
  }
```

```
  join roundrobin();
```

RoundRobin

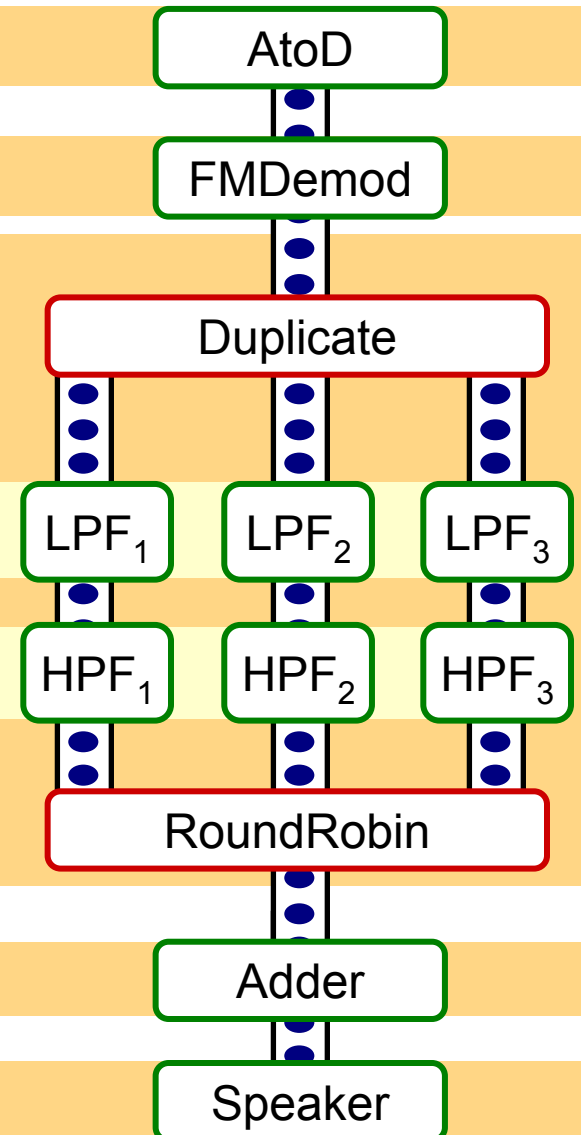
```
  add Adder();
```

Adder

```
  add Speaker();
```

Speaker

```
}
```



Programming in StreamIt

```
void->void pipeline FMRadio(int N, float freq1, float freq2) {
```

```
  add AtoD();
```

- Streams are easily composed

AtoD

```
  add FMDemod();
```

FMDemod

```
  add splitjoin {
    split duplicate;
    for (int i=0; i<N; i++) {
      add pipeline {
```

```
        add LowPassFilter(freq1 + i*(freq2-freq1)/N);
```

LPF₁

LPF₂

LPF₃

```
        add HighPassFilter(freq2 + i*(freq2-freq1)/N);
```

HPF₁

HPF₂

HPF₃

```
      }
```

```
    }
```

```
    join roundrobin();
```

RoundRobin

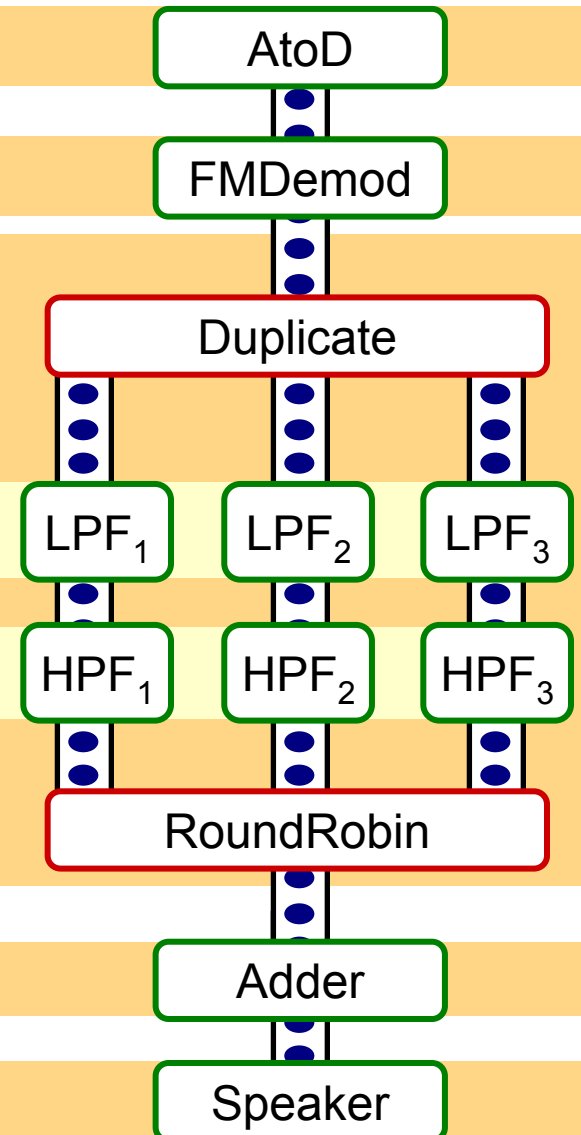
```
  add Adder();
```

Adder

```
  add Speaker();
```

Speaker

```
}
```



Programming in StreamIt

```
void->void pipeline FMRadio(int N, float freq1, float freq2) {
```

```
  add AtoD();
```

- Streams are
parameterized, and
malleable

AtoD

```
  add FMDemod();
```

FMDemod

```
  add splitjoin {
    split duplicate;
    for (int i=0; i<N; i++) {
      add pipeline {
```

```
        add LowPassFilter(freq1 + i*(freq2-freq1)/N);
```

LPF₁

LPF₂

LPF₃

```
        add HighPassFilter(freq2 + i*(freq2-freq1)/N);
```

HPF₁

HPF₂

HPF₃

```
      }
    }
  }
  join roundrobin();
```

RoundRobin

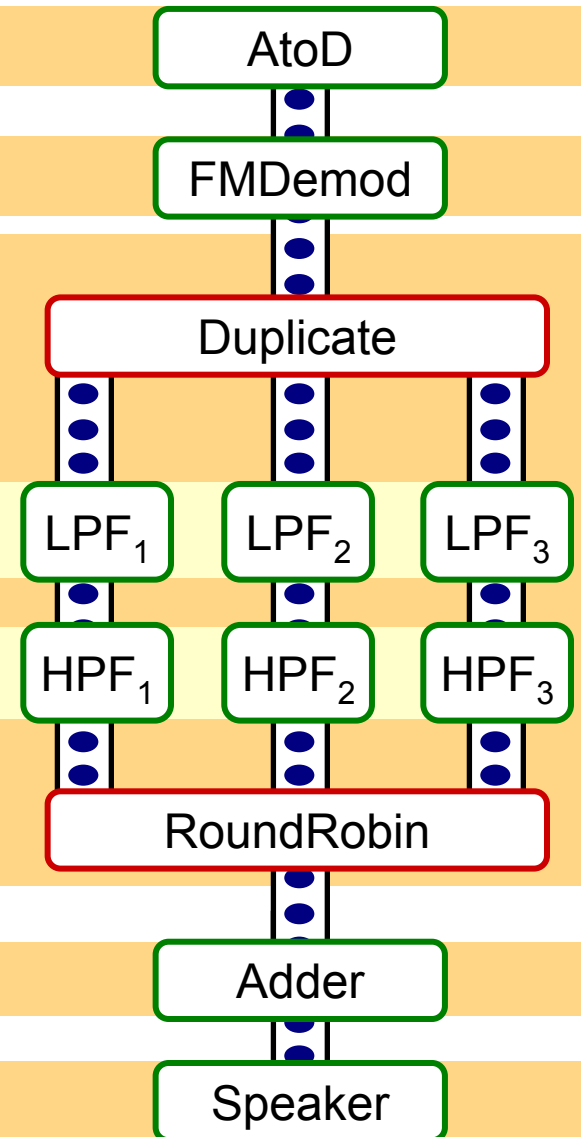
```
  add Adder();
```

Adder

```
  add Speaker();
```

Speaker

```
}
```



Programming in StreamIt

```
void->void pipeline FMRadio(int N, float freq1, float freq2) {
```

```
  add AtoD();
```

- Application is
architecture independent
(i.e., portable)

AtoD

```
  add FMDemod();
```

FMDemod

```
  add splitjoin {
    split duplicate;
    for (int i=0; i<N; i++) {
      add pipeline {
```

```
        add LowPassFilter(freq1 + i*(freq2-freq1)/N);
```

LPF₁

LPF₂

LPF₃

```
        add HighPassFilter(freq2 + i*(freq2-freq1)/N);
```

HPF₁

HPF₂

HPF₃

```
      }
    }
  }
```

```
  join roundrobin();
```

RoundRobin

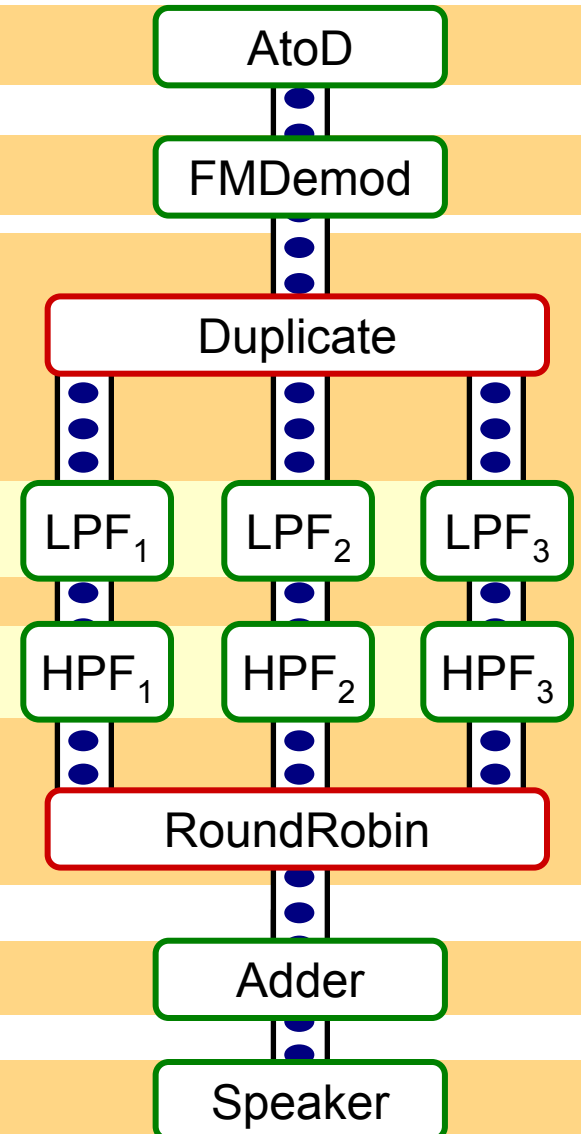
```
  add Adder();
```

Adder

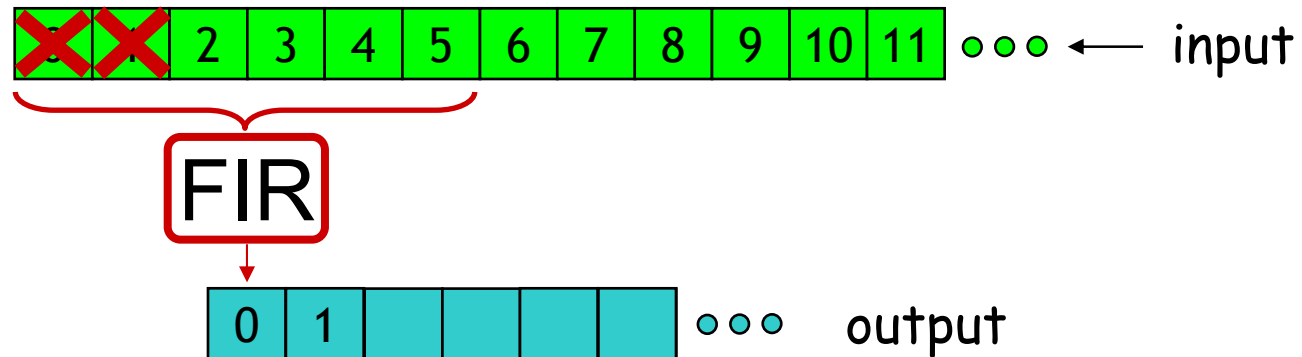
```
  add Speaker();
```

Speaker

```
}
```



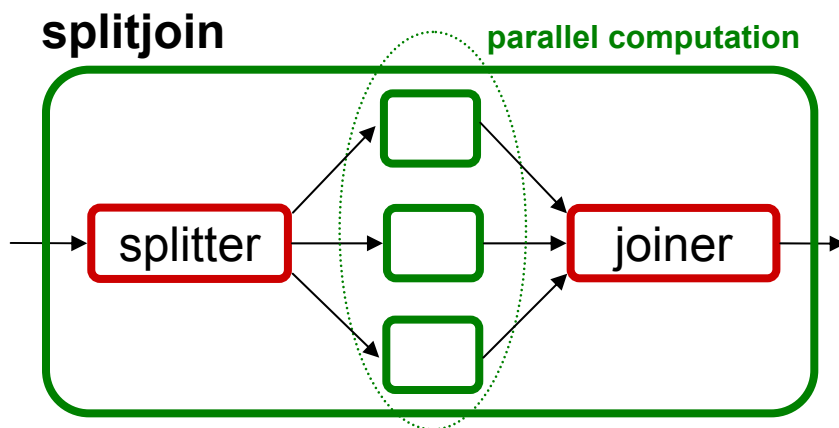
Filters as Computational Elements



```
float→float filter FIR (int N) {
    work push 1 pop 1 peek N {
        float result = 0;
        for (int i = 0; i < N; i++) {
            result += weights[i] * peek(i);
        }
        push(result);
        pop();
    }
}
```

Benefits of StreamIt

- Communication is exposed and pipeline parallelism is more readily discovered
- Flow of data provides a frame of reference for reasoning about “time” [PPoPP '05]
 - Powerful advantage when debugging parallel programs



versus

- Multiple threads with independent program counters
- Non-deterministic execution

StreamIt Development Environment ¹³

The screenshot displays the Eclipse IDE with the StreamIt development environment. The main window is titled "Debug - HelloWorld6.str - Eclipse Platform". It features several panels:

- Debug Console:** Shows the execution of the application, including threads like "System Thread [Finalizer] (Running)" and "Thread [main] (Suspended (breakpoint at line 20 in IntPrinter))".
- Breakpoints:** A list of breakpoints set at various lines of code, such as "[line: 7] - Complex" and "[line: 8] - real".
- Stream Graph:** A visual representation of the data flow graph. It shows a sequence of components: a "Pass (d=51)" component, followed by a "Pass (d=53)" component, and finally an "IntPrinter (d=19)" component. Each component has associated statistics like "Work Executions", "Pop Count", and "Push Count".
- Code Editor:** Displays the source code for "HelloWorld6.java". The code defines three filters: "IntSource", "IntPrinter", and "Pass", and a pipeline "Passer" that connects them.
- Overview of Stream Graph:** A small thumbnail view of the stream graph.
- Console:** Shows the output of the application, including the numbers 5, 6, 7, 8, 9, 10, and 11.

Annotations in red text highlight key features:

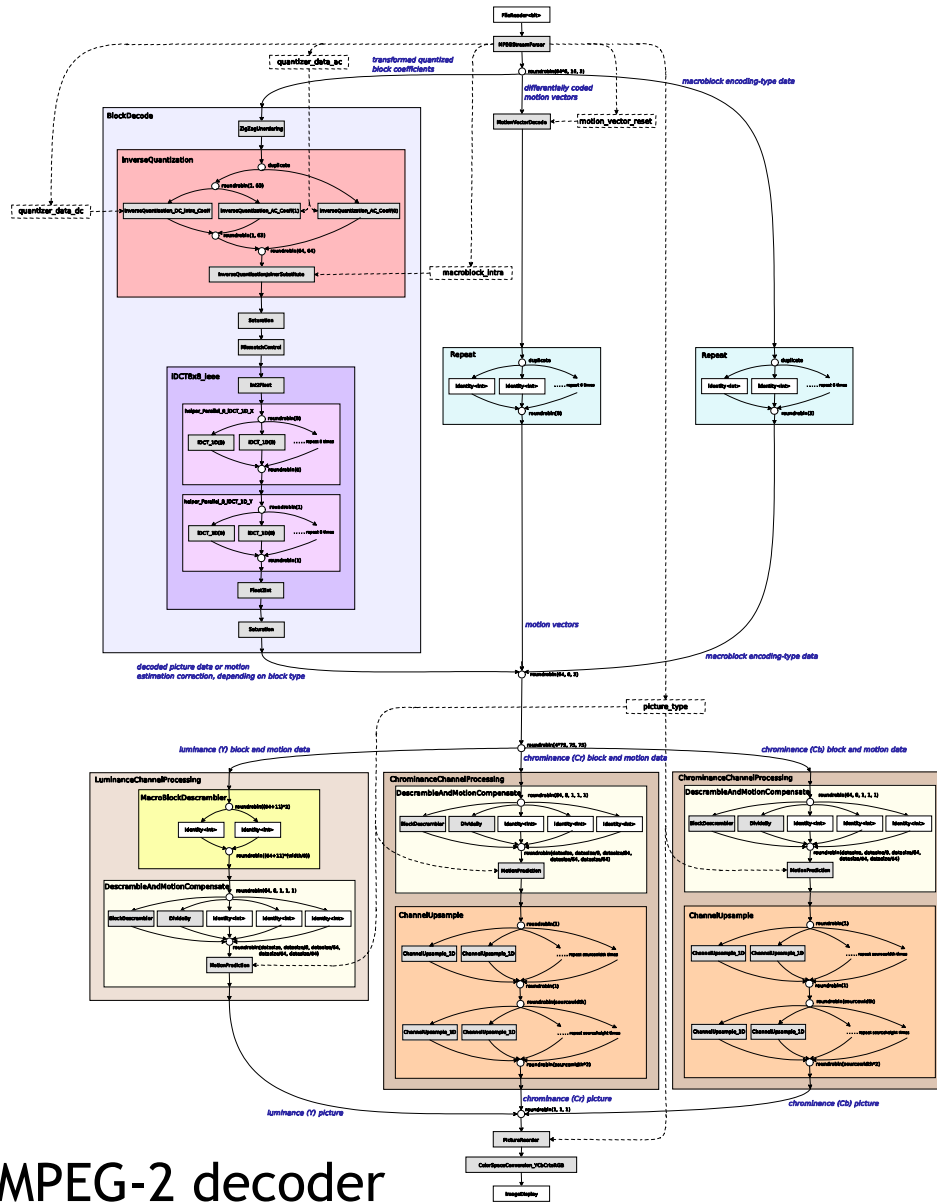
- General Debugging Information:** Points to the Breakpoints panel.
- StreamIt Graph Components:** Points to the Stream Graph panel.
- expanded and collapsed views of basic programmable unit:** Points to the expanded view of a "Pass" component in the Stream Graph.
- communication buffer with live data:** Points to the data flow between components in the Stream Graph.
- StreamIt Text Editor:** Points to the code editor.
- StreamIt Graph Zoom Panel:** Points to the Overview of Stream Graph panel.
- not shown: the StreamIt On-Line Help Manual:** Points to the area where the help manual would be displayed.
- Compiler and Output Consoles:** Points to the Console panel.

[PHEC '05]

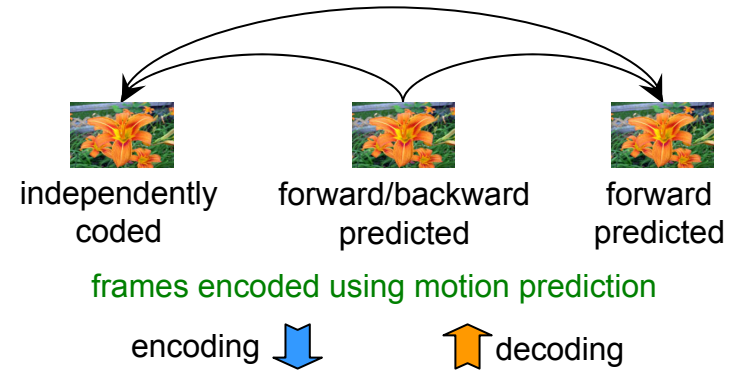
StreamIt Applications

- Software radio
- Frequency hopping radio
- Acoustic beam former
- Vocoder
- GMTI (ground moving target indicator)
- DES and Serpent blocked ciphers
- Sorting
- FFTs and DCTs
- JPEG
- ...

MPEG: Motion Video Codec



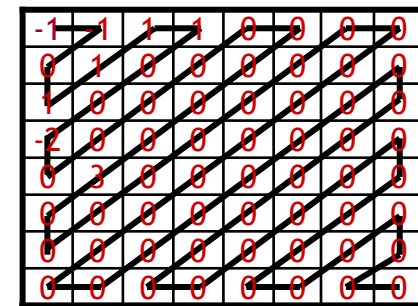
MPEG-2 decoder



frames encoded using motion prediction



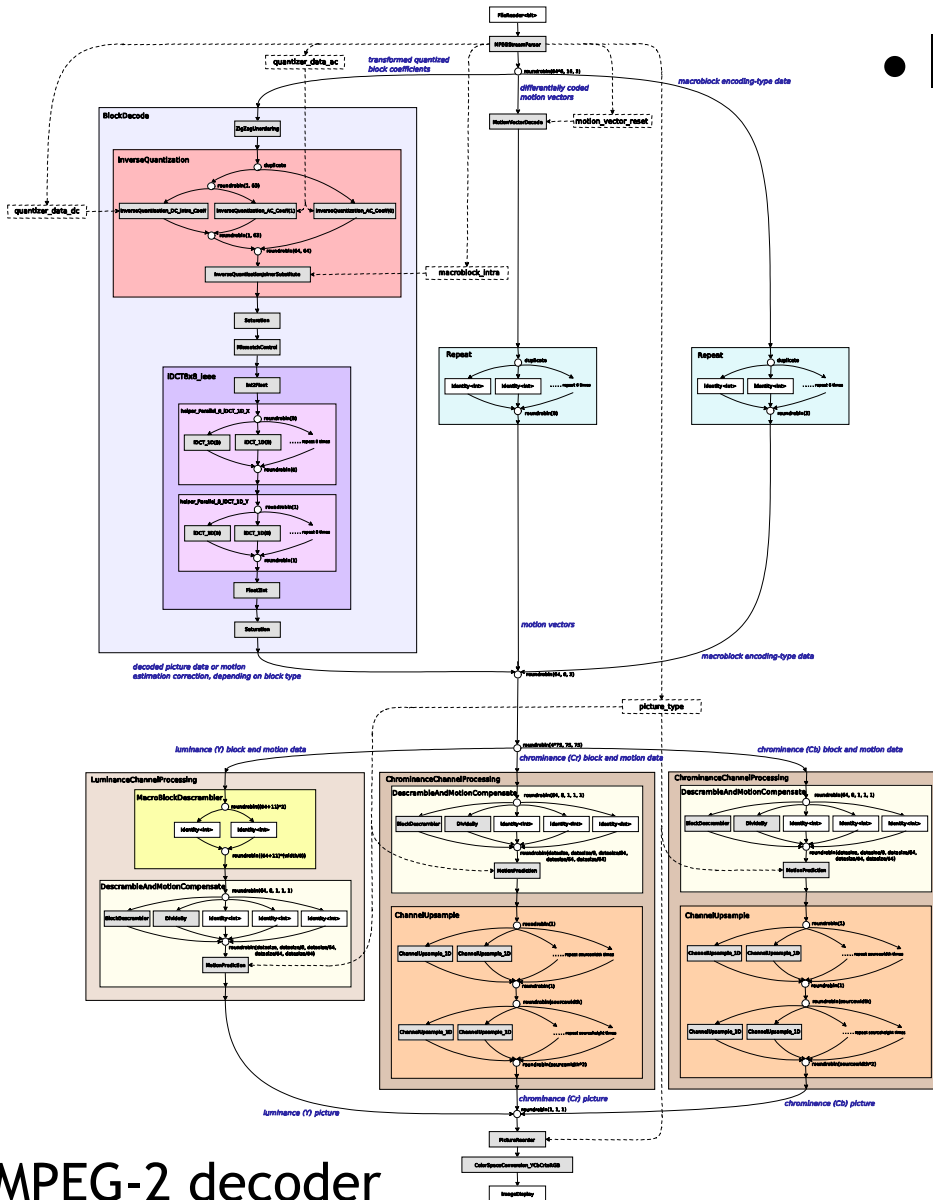
luminance and chrominance color data are separated



DCT and quantization of 8x8 image block

MPEG: Motion Video Codec

- Implementation statistics
 - 4921 lines of code
 - 48 static streams
 - Compile to ~2150 filters
 - 352x240 resolution
 - Reference C implementation has 9832 lines of code
 - Supports interlacing and multi-layer streams
 - 8 weeks of development
 - 1 programmer with no prior MPEG-2 experience



MPEG-2 decoder

Excerpt from StreamIt Implementation

Specification in Section 7.4.1: $F''[0][0] = \text{intra_dc_mult} \times \text{QF}[0][0]$

Table 7-4 - Relation between `intra_dc_precision` and `intra_dc_mult`

<code>intra_dc_precision</code>	<code>bits_of_precision</code>	<code>intra_dc_mult</code>
0	8	8
1	9	4
2	10	2
3	11	1

```
int->int filter InverseQuantization() {
    int[4] intra_dc_mult = {8, 4, 2, 1};
    int intra_dc_precision;

    work pop 1 push 1 {
        push(intra_dc_mult[intra_dc_precision] * pop());
    }
}
```

Excerpt from Reference Implementation

Specification in Section 7.4.1: $F''[0][0] = \text{intra_dc_mult} \times \text{QF}[0][0]$

Table 7-4 - Relation between `intra_dc_precision` and `intra_dc_mult`

<code>intra_dc_precision</code>	<code>bits_of_precision</code>	<code>intra_dc_mult</code>
0	8	8
1	9	4
2	10	2
3	11	1

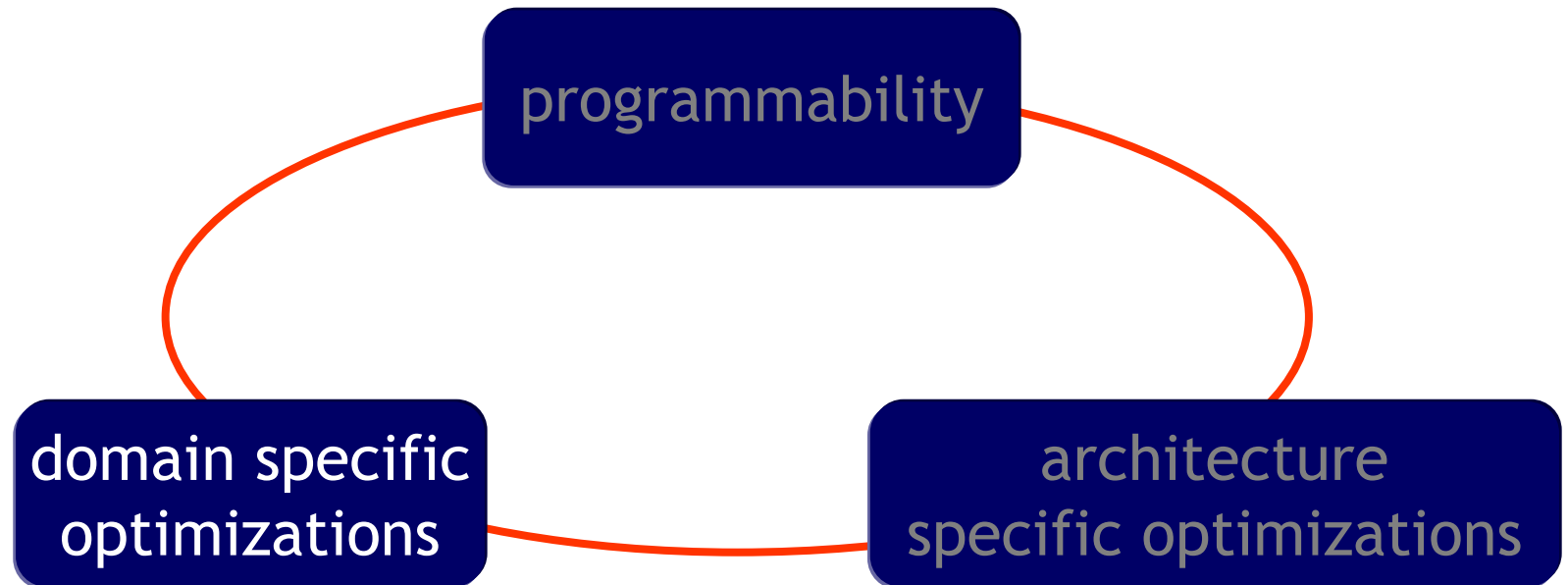
```
int[4] intra_dc_mult = {8, 4, 2, 1};

for (int m = 0; m < W*H/(16*16); m++)
    // six components for chrominance and luminance
    for (int comp = 0; comp < 6; comp++)
        if (macroblock[m].intra)
            macroblock[m].block[comp][0] *= intra_dc_mult[intra_dc_precision];

    // and many lines later
    if (cc == 0)
        val = (dc_dct_pred[0] += Get_Luma_DC_dct_diff());
    else if (cc == 1)
        val = (dc_dct_pred[1] += Get_Chroma_DC_dct_diff());
    else
        val = (dc_dct_pred[2] += Get_Chroma_DC_dct_diff());
    if (Fault_Flag) return;
    bp[0] = val << (3-intra_dc_precision);
```

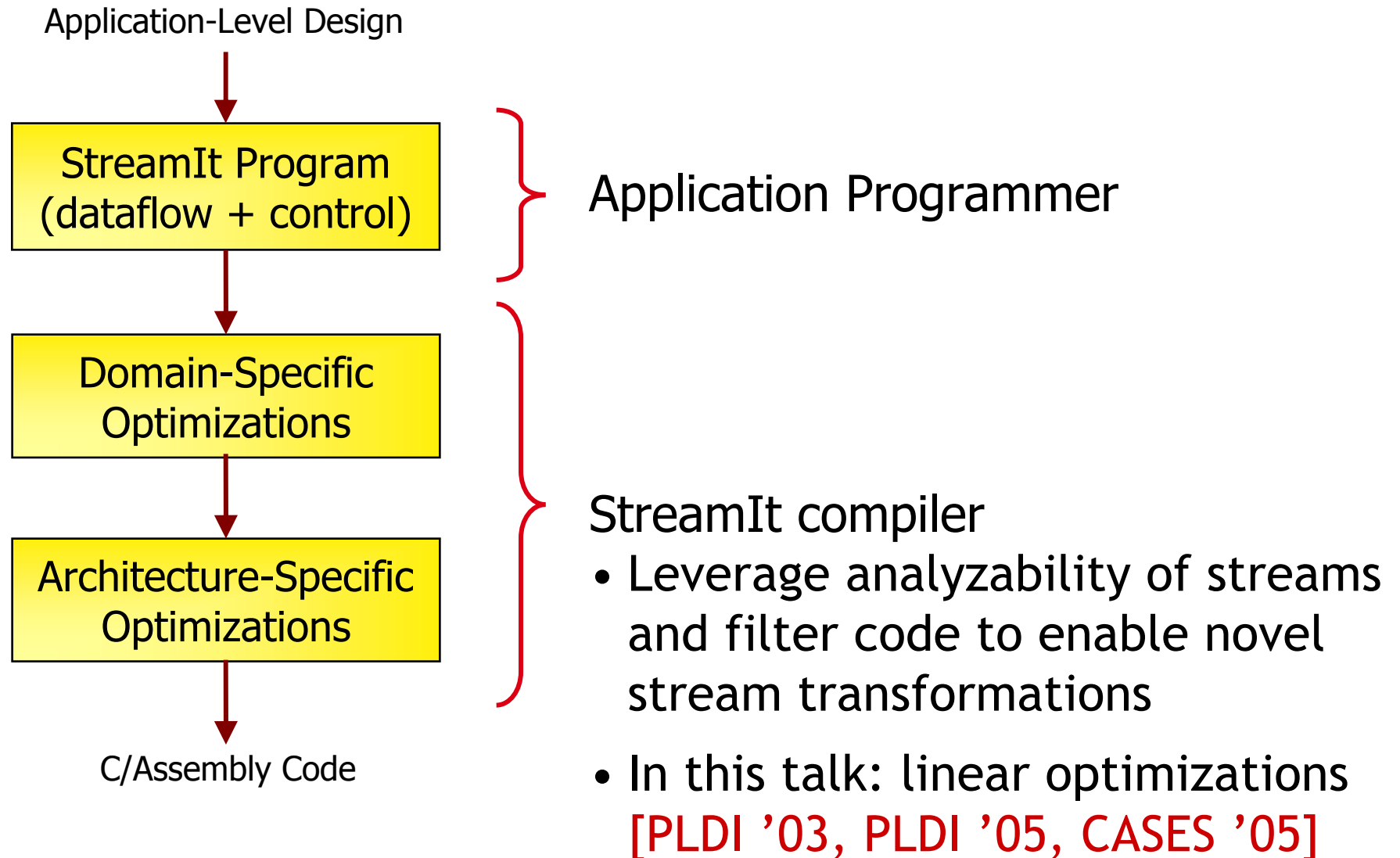
The StreamIt Vision

- Boost productivity, enable faster development and rapid prototyping



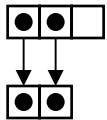
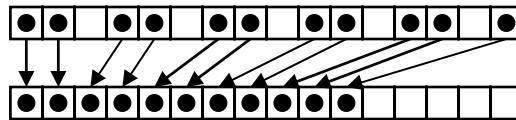
- Simple and effective optimizations for streams
- Targeting tiled architectures, clusters of workstations, DSPs, and traditional uniprocessors

Design Flow with StreamIt



Linear Filter Example

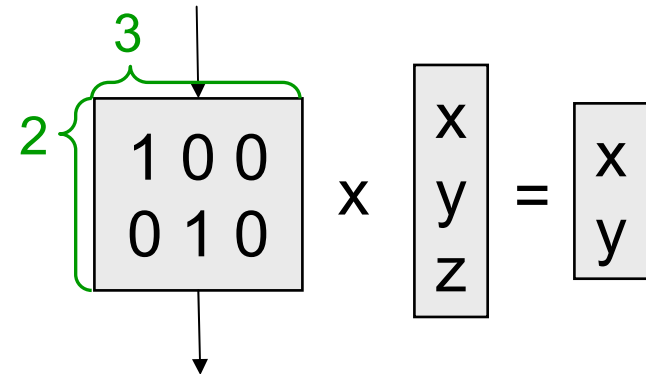
- “Drop every third bit in the bit stream”



```

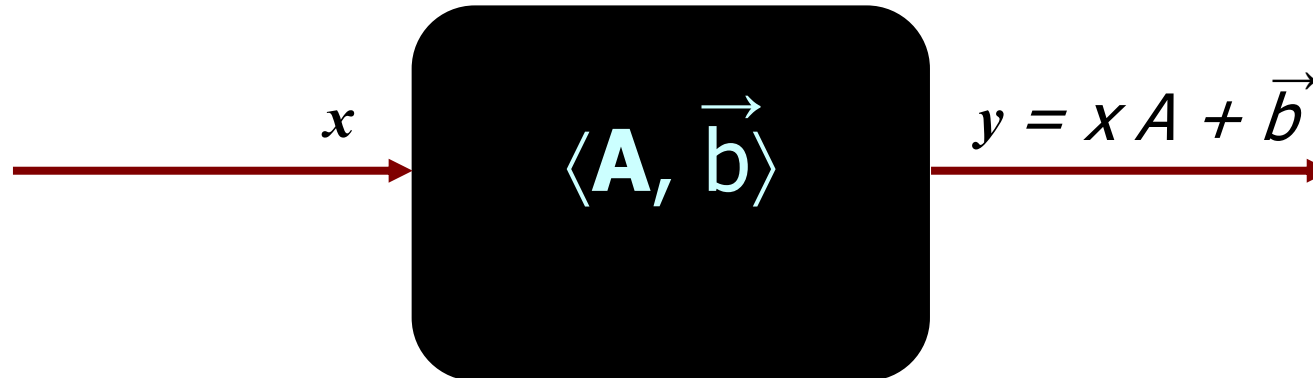
bit → bit filter DropThirdBit {
  work push 2 pop 3 {
    push(pop());
    push(pop());
    pop();
  }
}

```



In General

- A linear filter is a tuple $\langle \mathbf{A}, \vec{\mathbf{b}} \rangle$
 - \mathbf{A} : matrix of coefficients
 - $\vec{\mathbf{b}}$: vector of constants
- Example

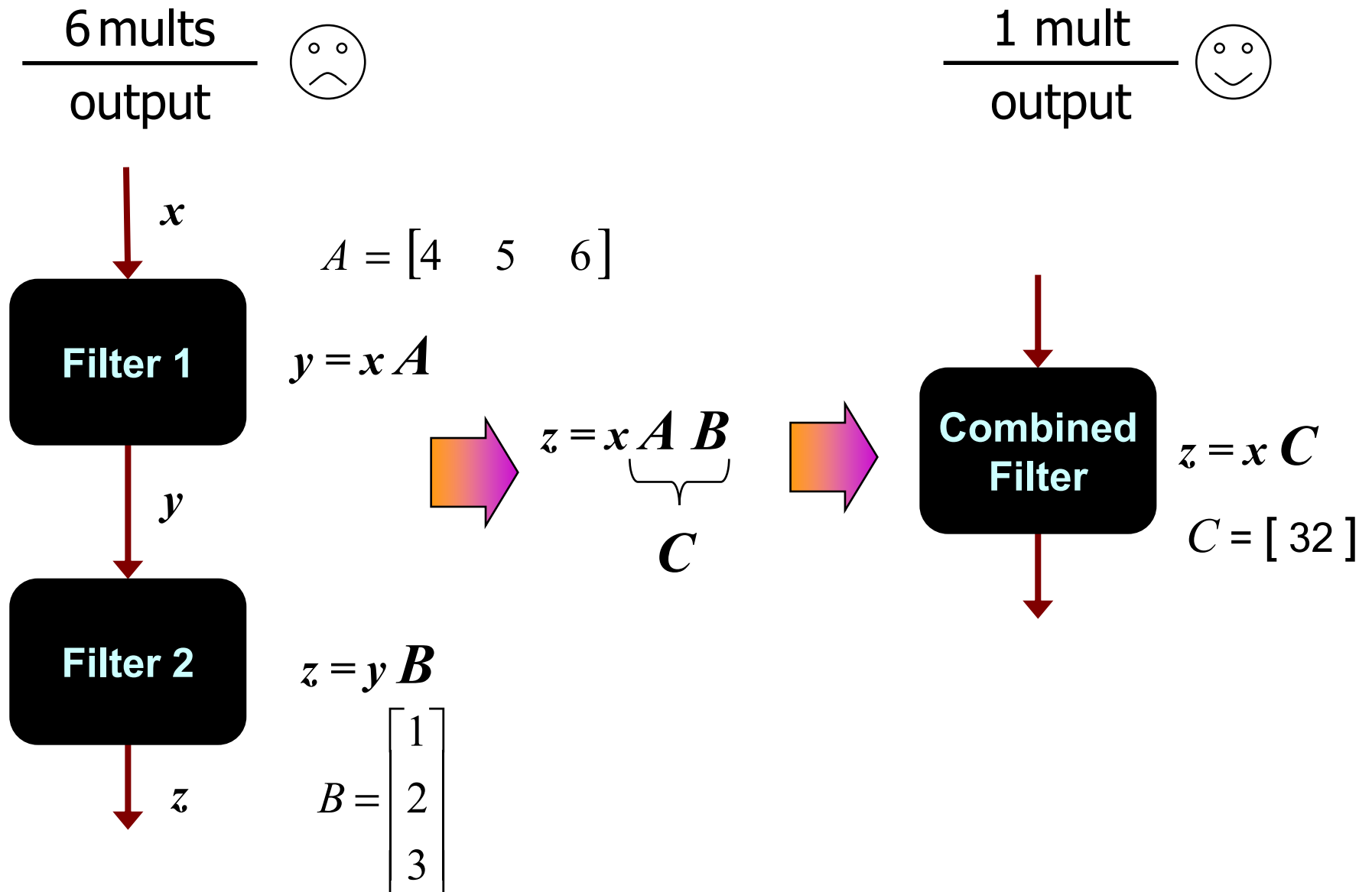


- Linear dataflow analysis resembles constant propagation

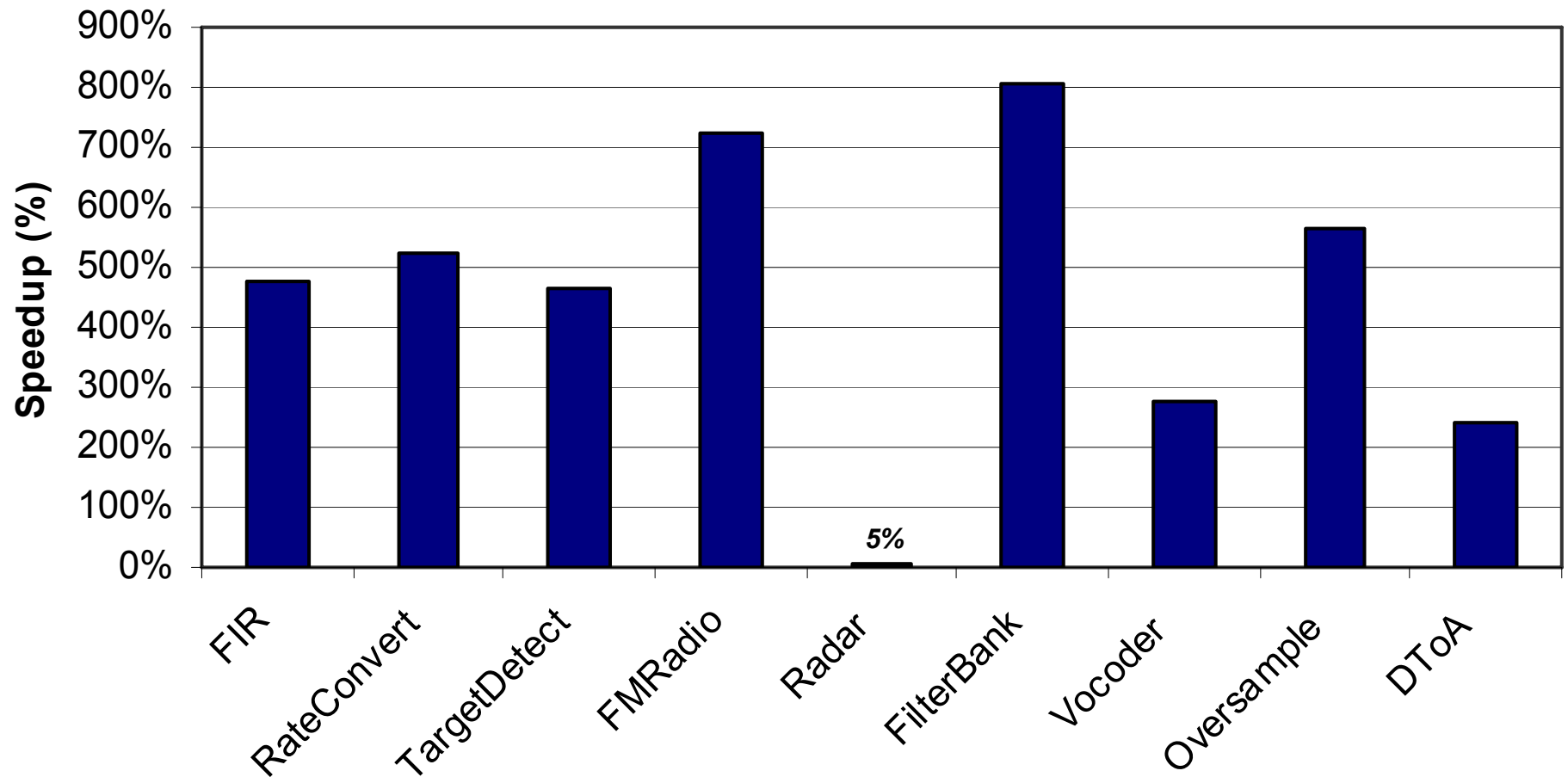
Opportunities for Linear Optimizations

- Occur frequently in streaming codes
 - FIR filters
 - Compressors
 - Expanders
 - DFT/DCT
 - Bit permutations in encryption algorithms
 - JPEG and MPEG codecs
 - ...
- Example optimizations
 - Combining adjacent nodes
 - Also, translating to frequency domain when profitable

Combining Linear Filters



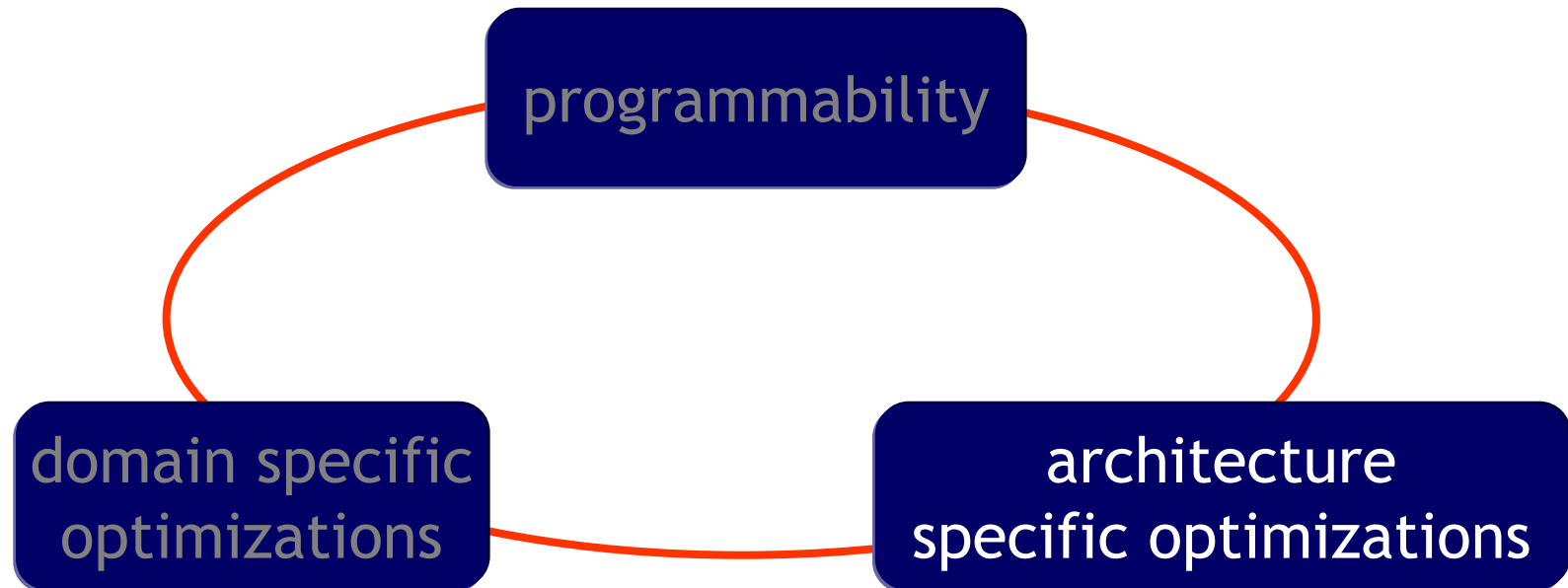
Results from Linear Optimizations



Pentium 4 results compared to baseline StreamIt

The StreamIt Vision

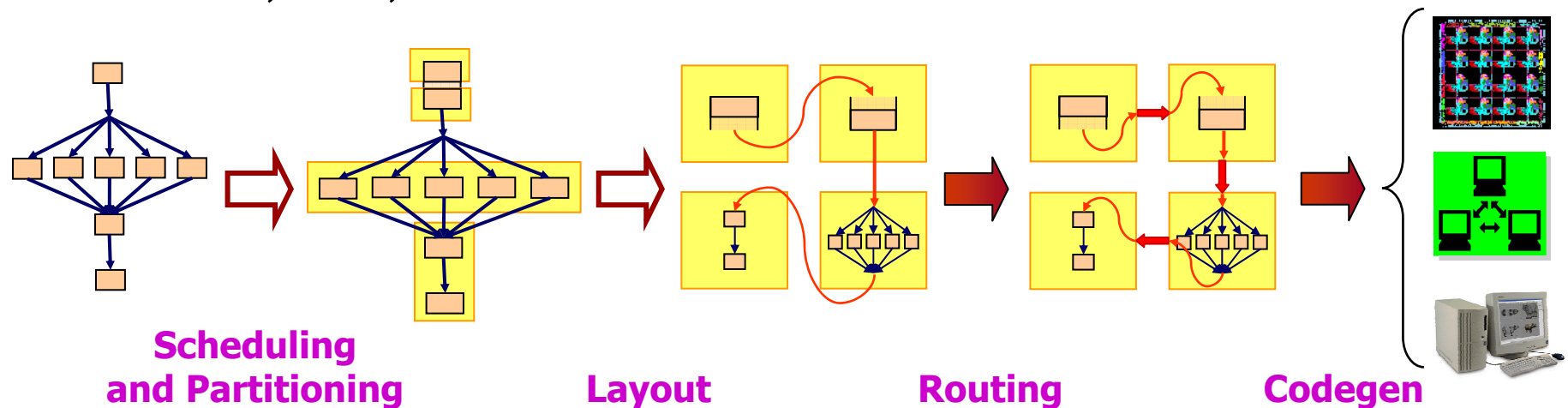
- Boost productivity, enable faster development and rapid prototyping



- Simple and effective optimizations for streams
- Targeting tiled architectures, clusters of workstations, DSPs, and traditional uniprocessors

Core Compilation Technology

- Focused on a common challenges in modern and future architectures
 - MIT Raw fabric architecture
 - Clusters of workstations
 - ARM, x86, and IA-64

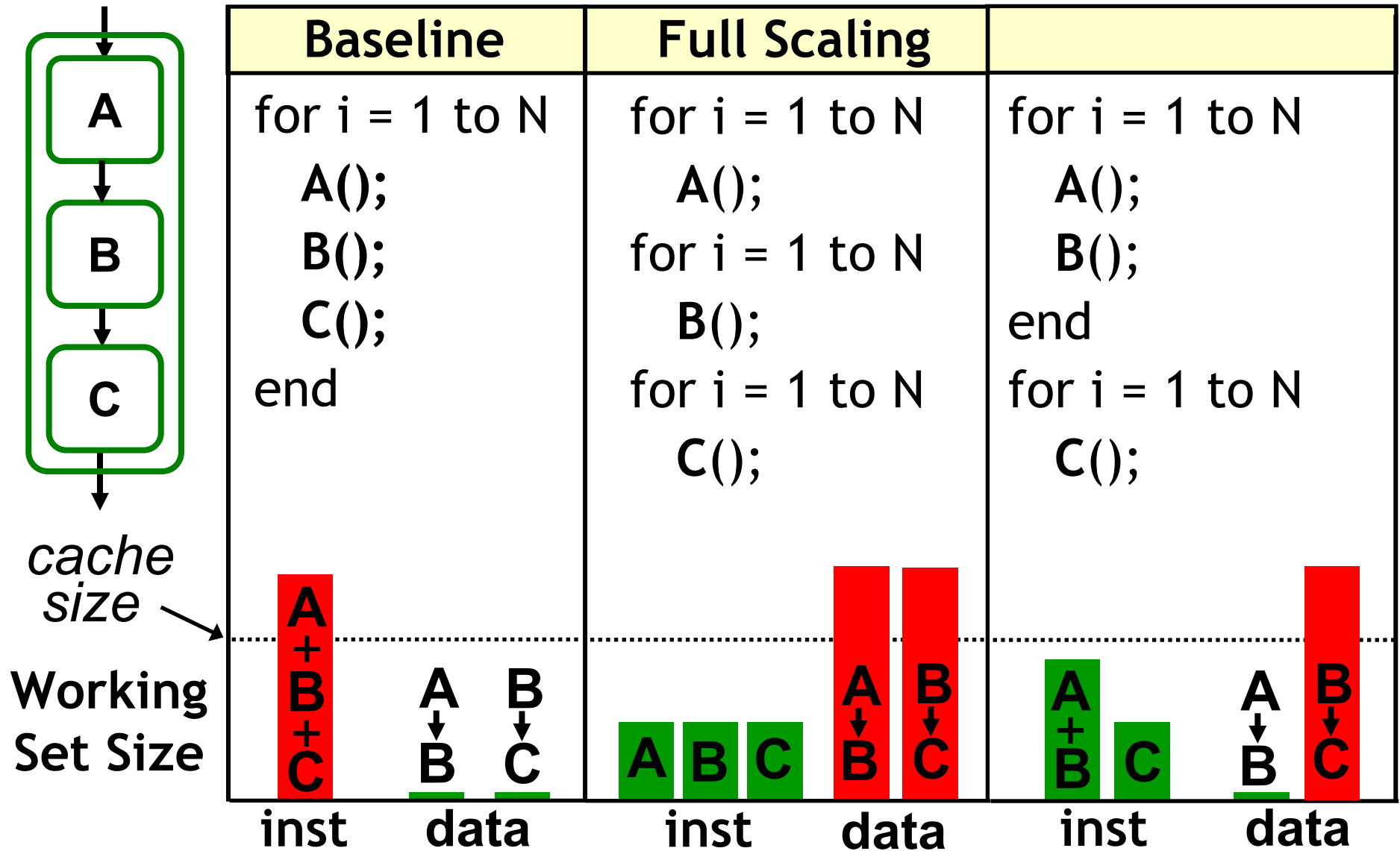


- Compiler's role: map the computation and communication pattern to processors, memories, and communication substrates

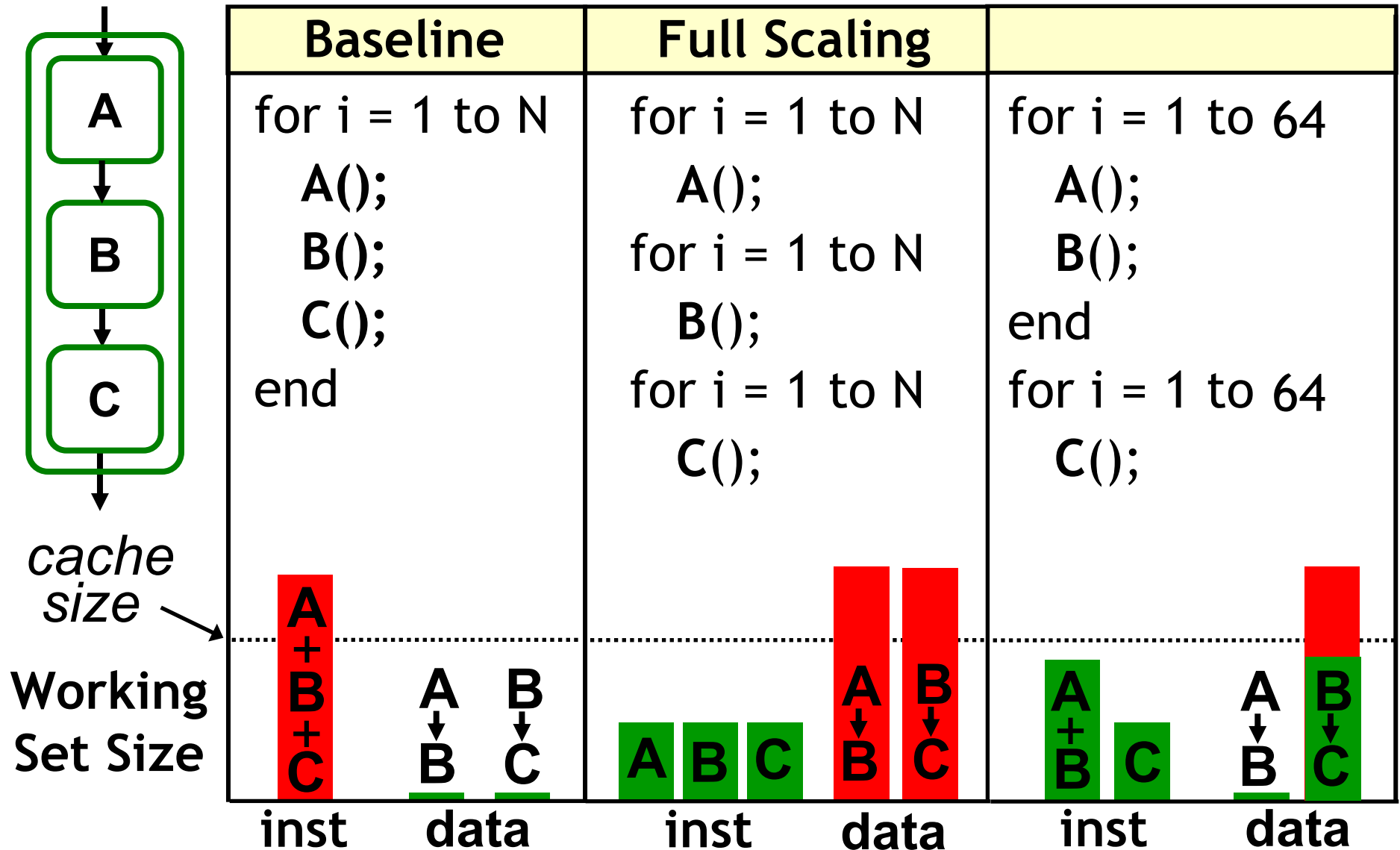
Compiler Issues

- Load balancing
 - Resource utilization
 - Fault tolerance
 - Dynamic reconfiguration
 - ...
-
- In this talk: cache aware scheduling and partitioning [LCTES '05]

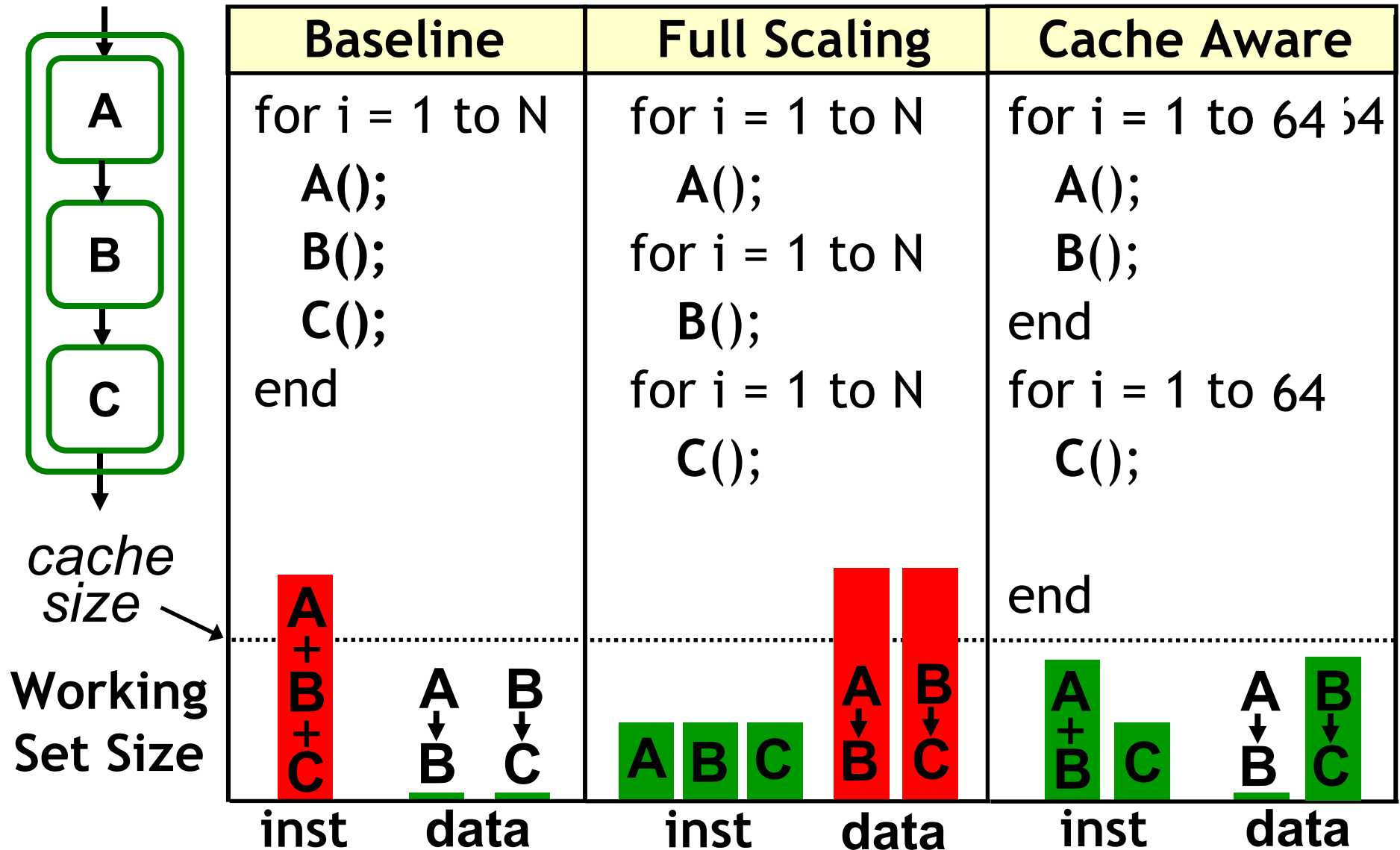
Example Cache Optimization



Example Cache Optimization



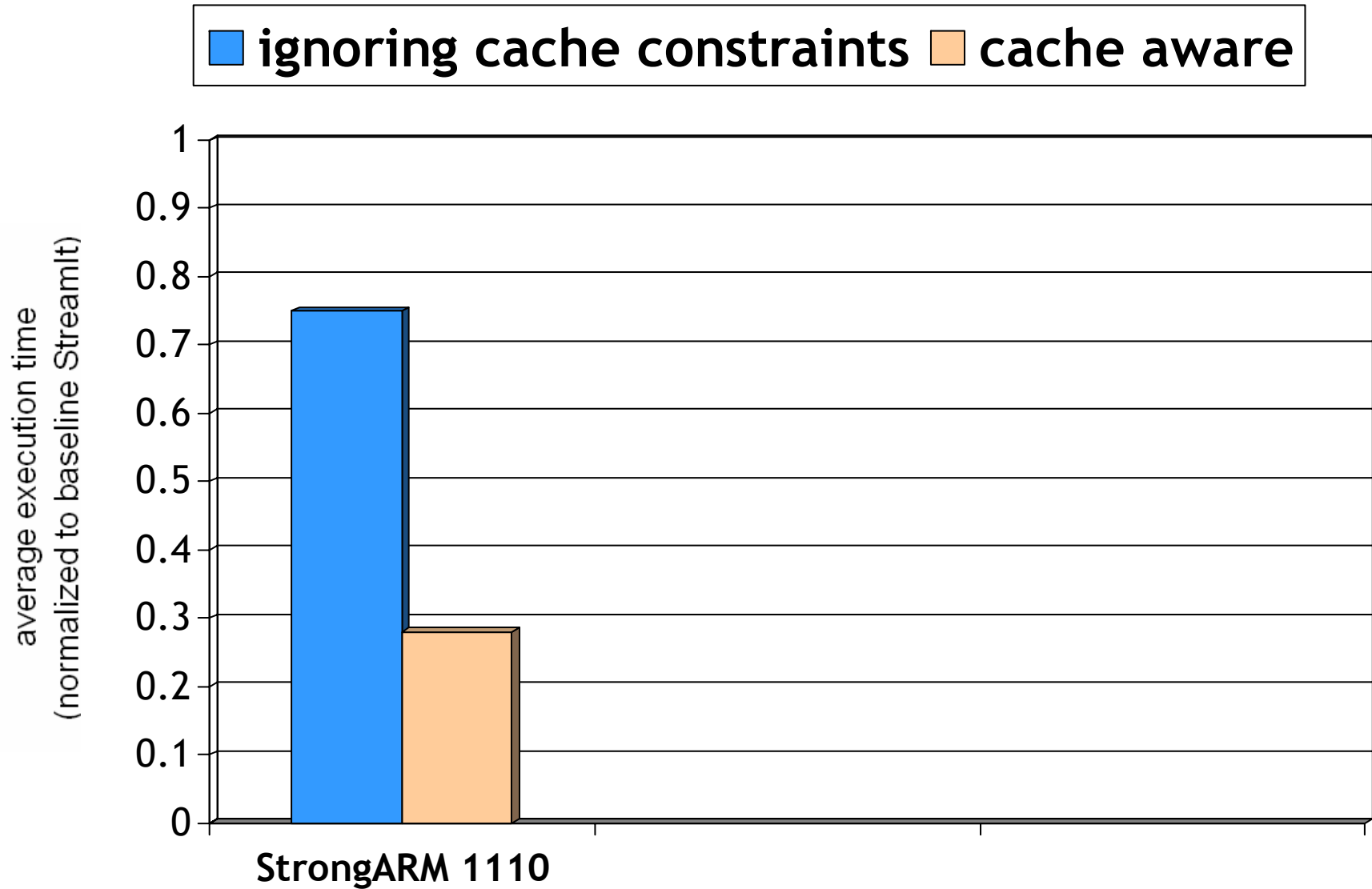
Example Cache Optimization



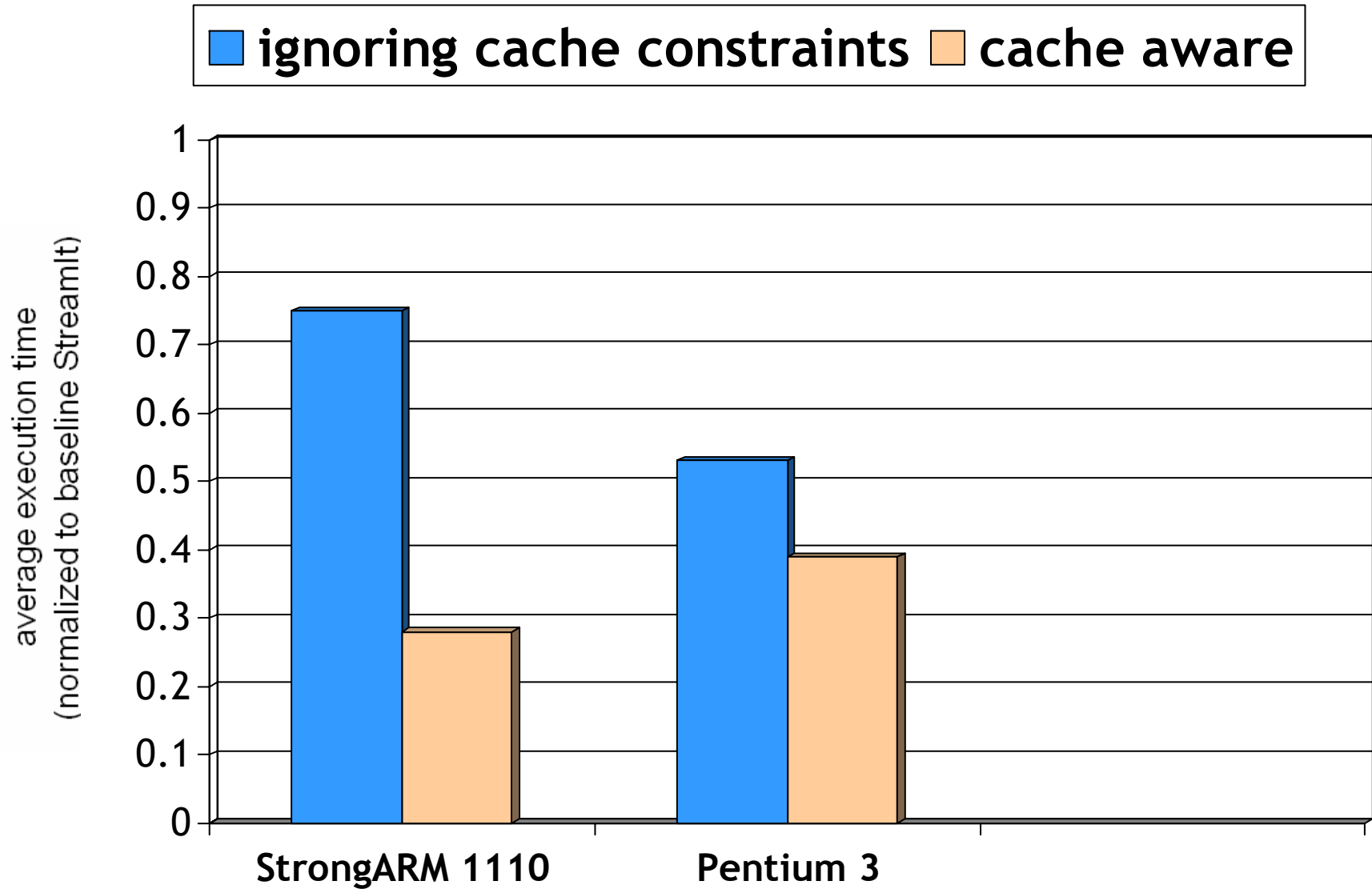
Evaluation Methodology

- StreamIt compiler generates C code
 - Baseline StreamIt optimizations
 - Unrolling, constant propagation
 - Compile C code with gcc-v3.4 with -O3 optimizations
- StrongARM 1110 (XScale) embedded processor
 - 370MHz, 16Kb I-Cache, 8Kb D-Cache
 - No L2 Cache (memory 100× slower than cache)
 - Median user time
- Also Pentium 3 and Itanium 2 processors
- Suite of 11 StreamIt Benchmarks

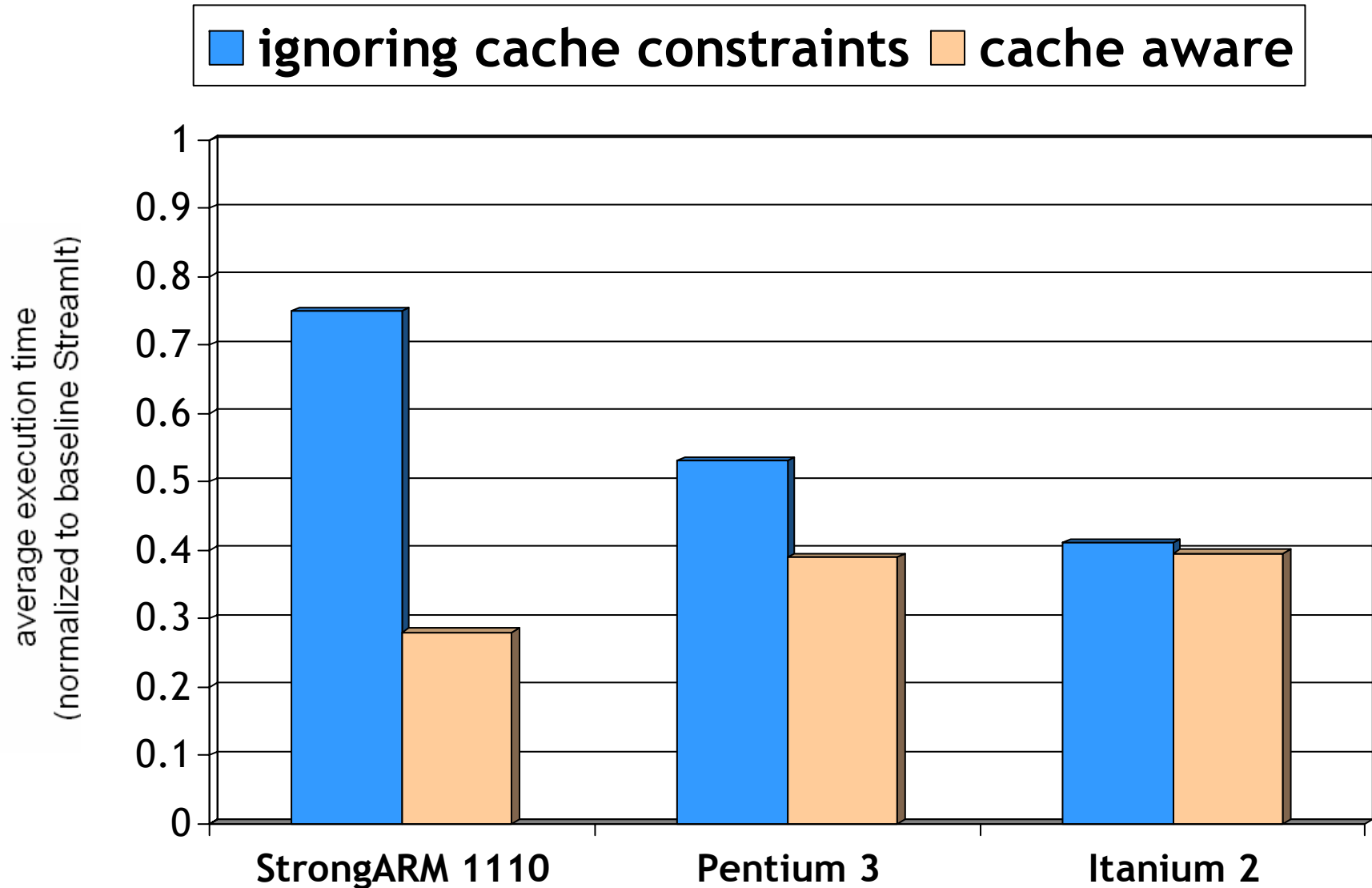
Cache Optimizations Results



Cache Optimizations Results

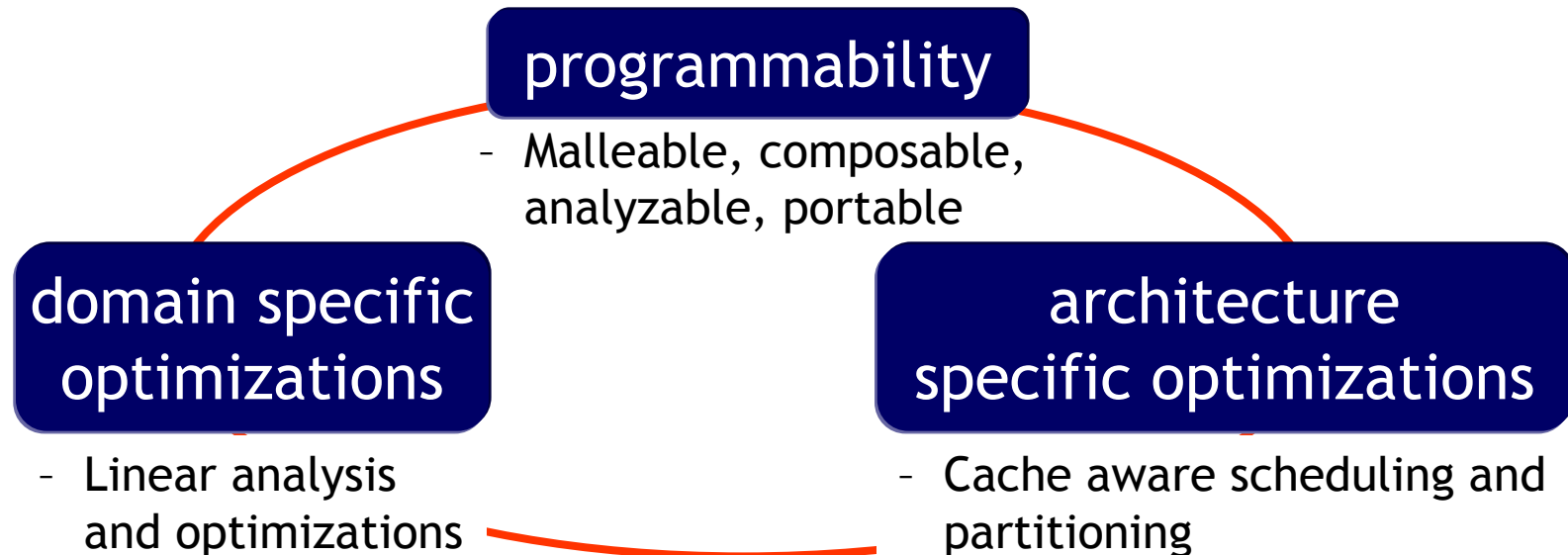


Cache Optimizations Results



Concluding Remarks

- StreamIt improves programmer productivity without compromising performance
 - Easily identify pipeline and data parallelism
 - Expose information for domain specific and architecture specific optimizations



Broader Impact

- Integration into future HPCS languages
 - IBM: X10
- StreamIt for graphics applications
 - Programmable graphics pipeline [Graphics Hardware '05]
- StreamIt for emerging architectures
- Looking for users with interesting applications

High-Productivity Stream Programming for High-Performance Systems

Rodric Rabbah, Bill Thies, Michael Gordon, Janis Sermulins,
and Saman Amarasinghe

Massachusetts Institute of Technology

The logo for StreamIt features the word "StreamIt" in a blue, sans-serif font. A red horizontal line with an arrowhead pointing to the right is positioned above the "m" and "i".

StreamIt

<http://cag.lcs.mit.edu/streamit>