# Bottleneck Elimination from Stream Graphs
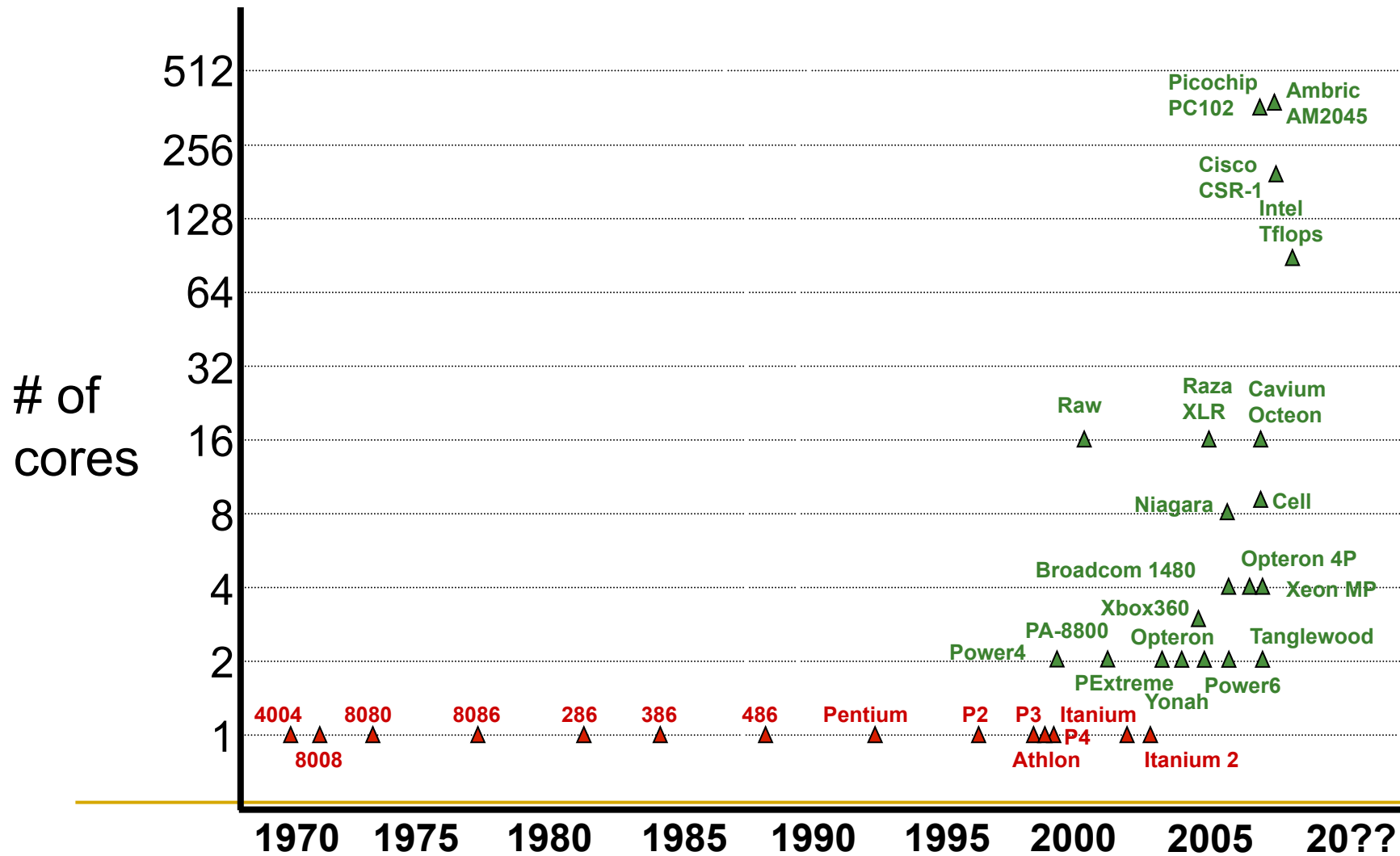
S. M. Farhad

The University of Sydney

Joint work with
Yousun Ko
Bernd Burgstaller
Bernhard Scholz
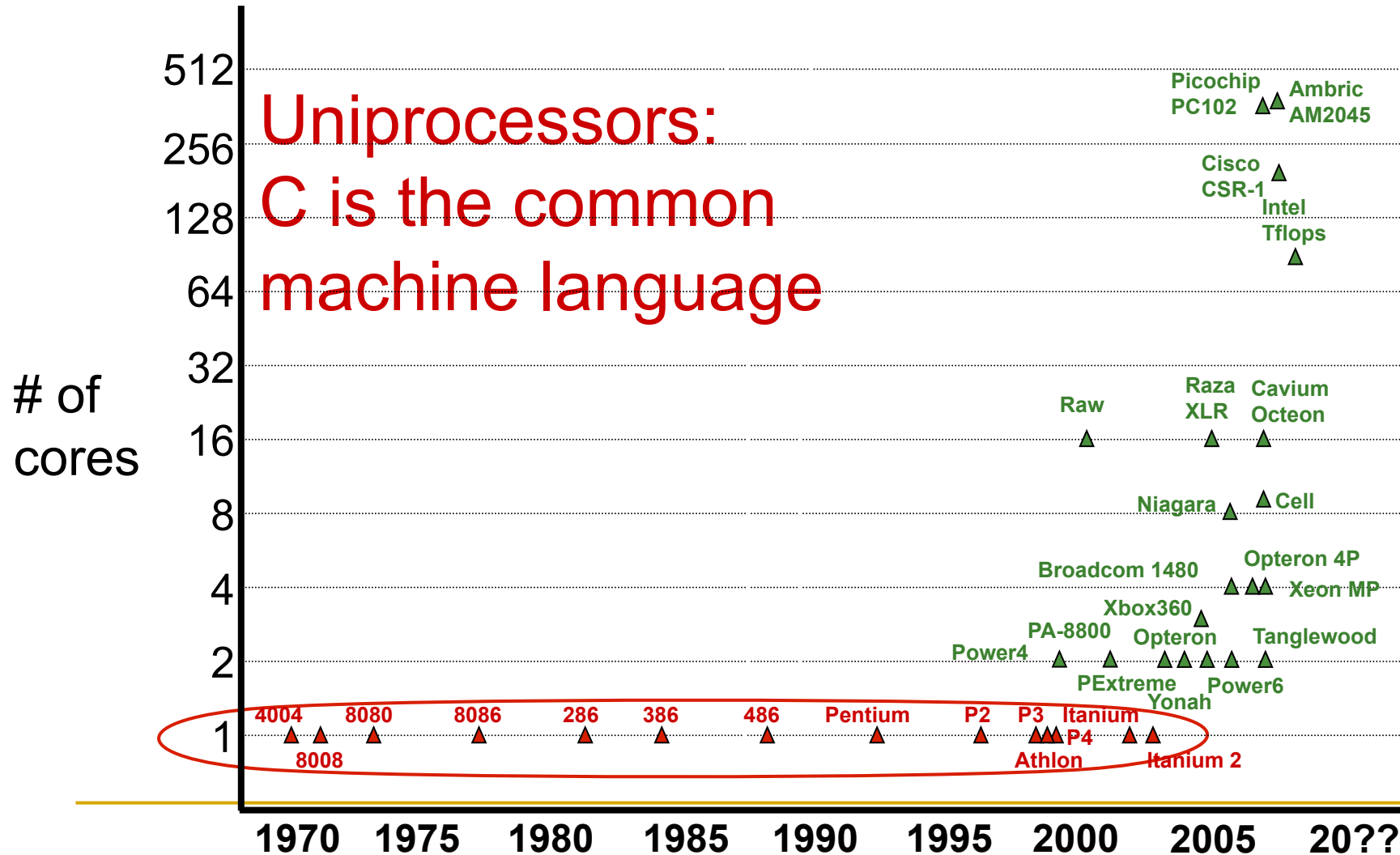
# Outline

- **Motivation**
  - Multicore
  - Stream programming
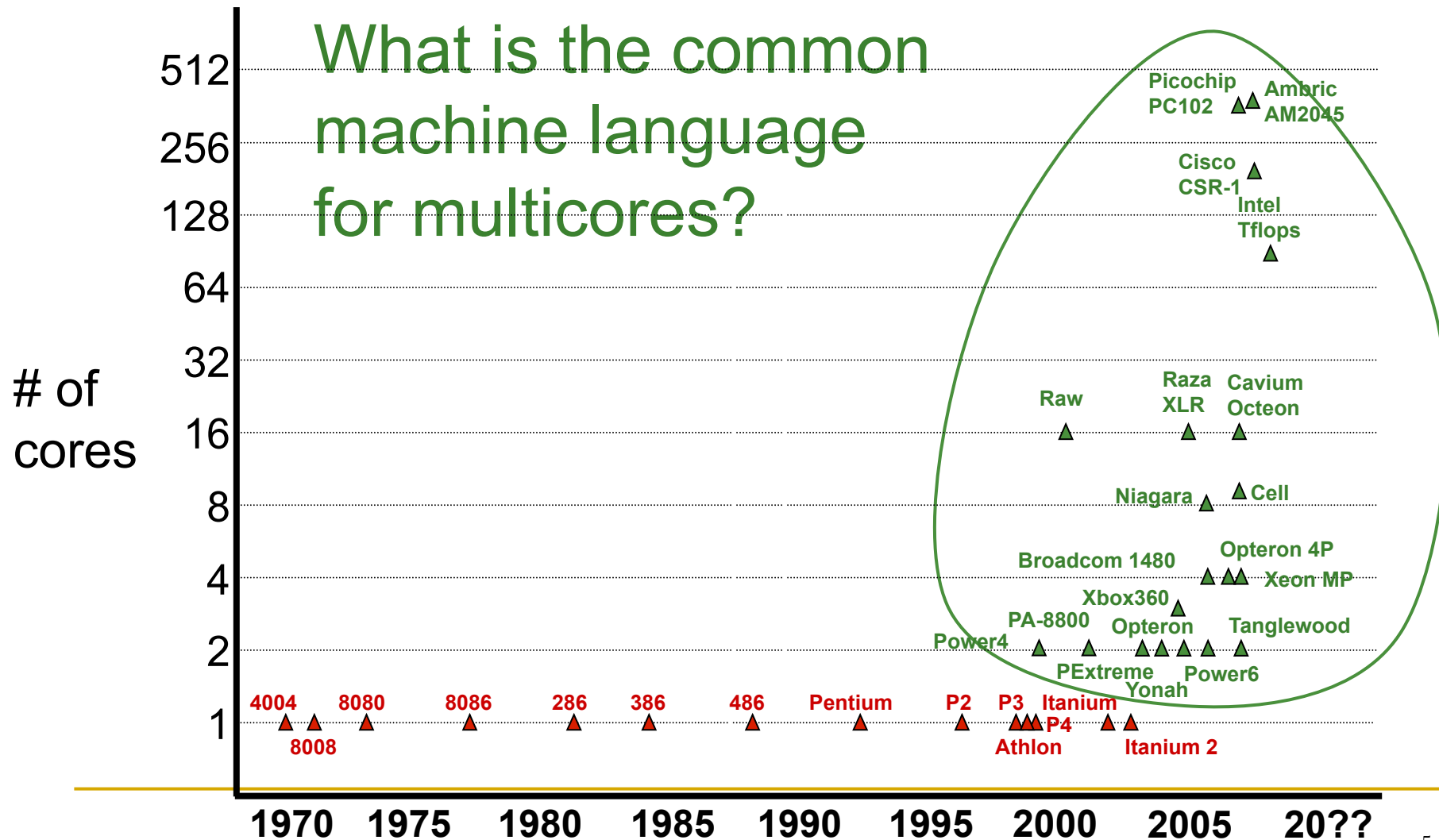- **Research question**
- **Our work**
- **Summary**

# Multicores Are Here!



# of cores

512 — Picochip PC102 / Ambric AM2045

256 — Cisco CSR-1

128 — Intel Tflops

64

32 — Raza XLR / Cavium Octeon

16 — Raw

8 — Niagara / Cell

4 — Broadcom 1480 / Opteron 4P / Xeon MP

2 — Power4, PA-8800, Opteron, Tanglewood, Xbox360, PExtreme, Power6

1 — 4004, 8008, 8080, 8086, 286, 386, 486, Pentium, P2, P3, P4, Athlon, Itanium, Itanium 2, Yonah

1970  1975  1980  1985  1990  1995  2000  2005  20??

3

# Multicores Are Here!



Uniprocessors:
C is the common
machine language

**# of cores**

512, 256, 128, 64, 32, 16, 8, 4, 2, 1

Picochip PC102, Ambric AM2045, Cisco CSR-1, Intel Tflops, Raza XLR, Cavium Octeon, Raw, Niagara, Cell, Broadcom 1480, Opteron 4P, Xeon MP, Xbox360, PA-8800, Opteron, Tanglewood, Power4, PExtreme, Power6, Yonah

4004, 8080, 8086, 286, 386, 486, Pentium, P2, P3, Itanium, P4, 8008, Athlon, Itanium 2

1970  1975  1980  1985  1990  1995  2000  2005  20??

# Multicores Are Here!



What is the common machine language for multicores?

**# of cores** (y-axis): 512, 256, 128, 64, 32, 16, 8, 4, 2, 1

x-axis: 1970, 1975, 1980, 1985, 1990, 1995, 2000, 2005, 20??

Labels (green):
Picochip PC102, Ambric AM2045, Cisco CSR-1, Intel Tflops, Raza XLR, Cavium Octeon, Raw, Niagara, Cell, Opteron 4P, Broadcom 1480, Xeon MP, Xbox360, PA-8800, Opteron, Tanglewood, Power4, PExtreme, Yonah, Power6

Labels (red):
4004, 8080, 8086, 286, 386, 486, Pentium, P2, P3, Itanium, P4, 8008, Athlon, Itanium 2

# Stream Programming Paradigm

- **Research topic in parallel programming**
- **Various forms of parallelism**
  - Pipeline, task, and data
- **Applications**
  - Signal Processing
  - Multi-media
  - High-Performance Computing
- **Programs expressed as stream graphs**
  - Streams
    - Infinite sequence of data elements (aka. Tokens)
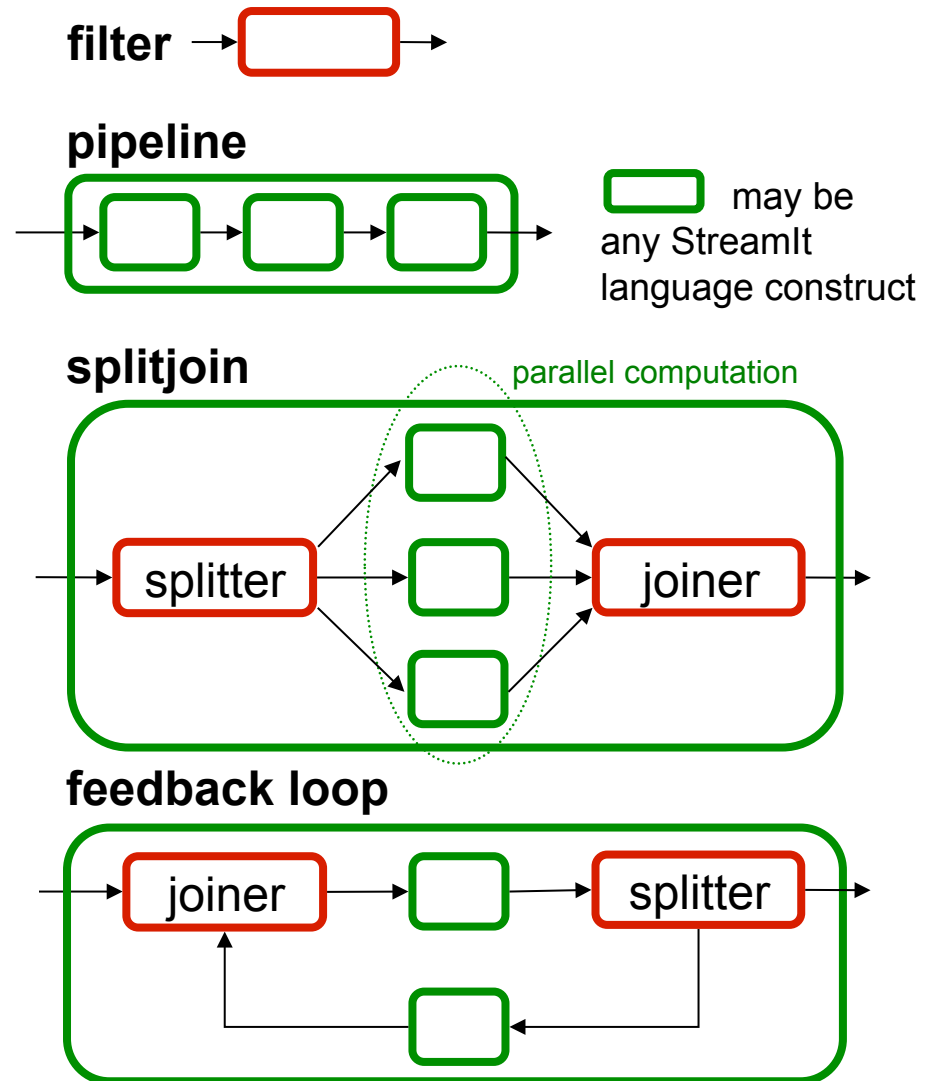  - Actor
    - Functions applied to streams

Stream

Actor

Stream

# Properties of Stream Program

- **Regular and repeating computation**

- **Independent actors with explicit communication**
  - Producer / Consumer dependencies

AtoD

FMDemod

Splitter

LPF$_1$   LPF$_2$   LPF$_3$
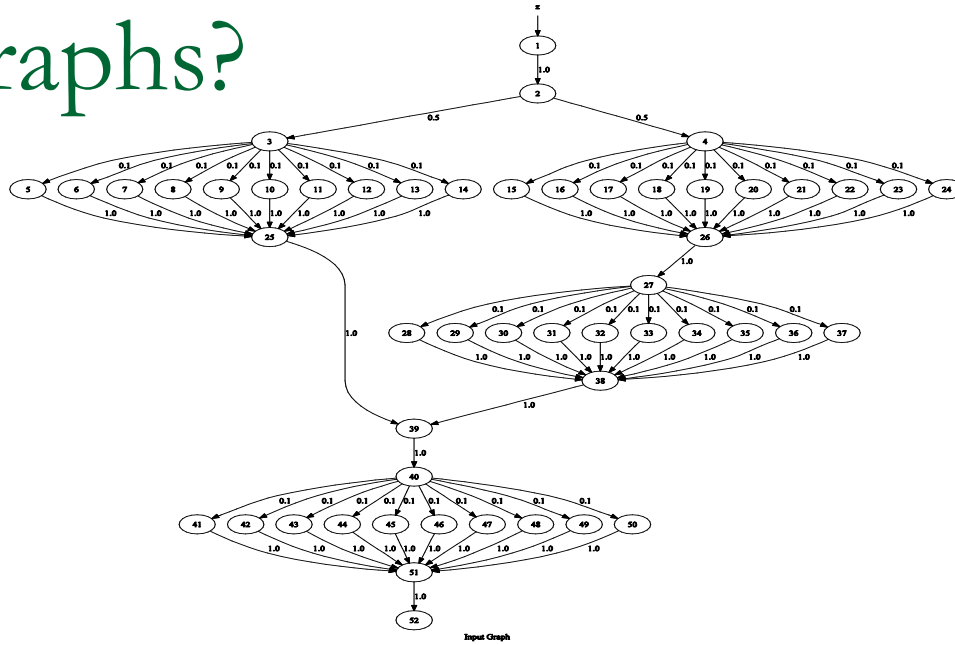
HPF$_1$   HPF$_2$   HPF$_3$

Joiner
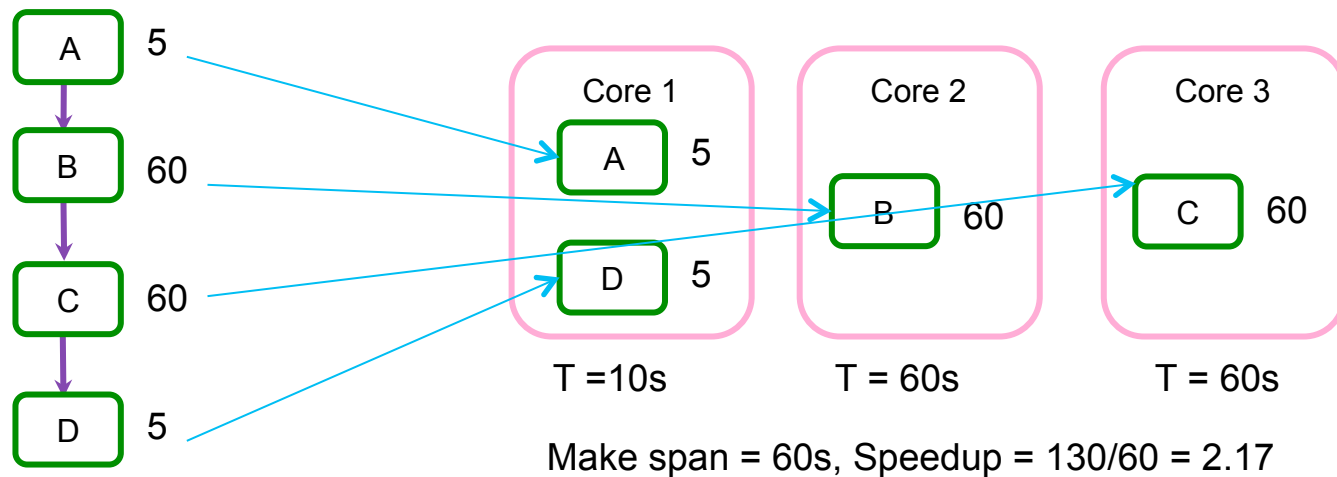
Adder

Speaker

# StreamIt Language [ASPLOS'2&6, PLDI'3]

- An implementation of stream prog.
- Each construct has single input/output stream
- Hierarchical structure
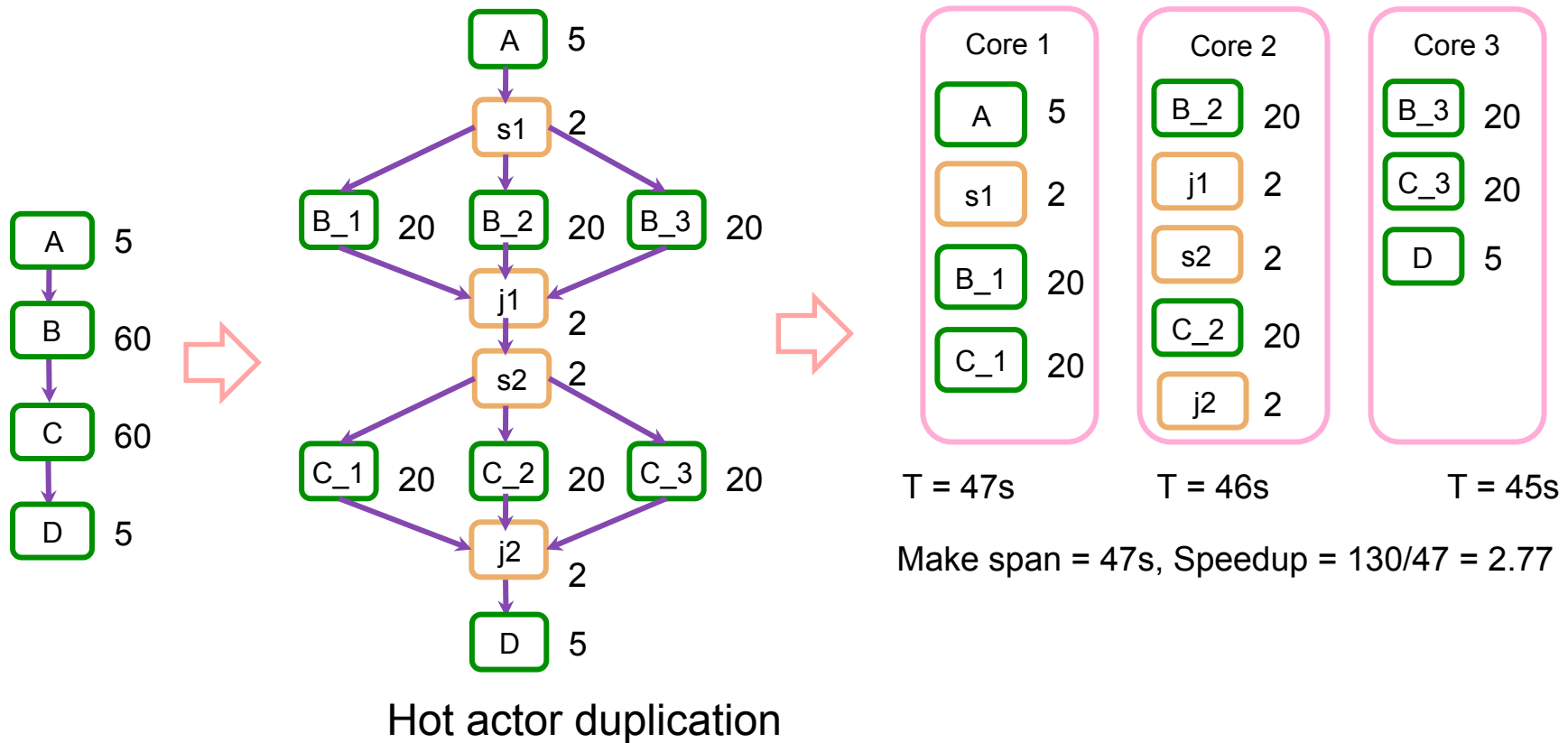- Filters can be stateful/stateless

**filter**

**pipeline**

may be any StreamIt language construct

**splitjoin**

parallel computation

splitter

joiner

**feedback loop**

joiner

splitter

# Research Question: How to Eliminate Bottlenecks (Hot Actors) from Stream Graphs?



Input Graph

# Mapping Actors



A   5

B   60

C   60

D   5

Core 1
A   5
D   5
T =10s

Core 2
B   60
T = 60s

Core 3
C   60
T = 60s

Make span = 60s, Speedup = 130/60 = 2.17

# Bottleneck Actors Limit the Performance



Hot actor duplication

Make span = 47s, Speedup = 130/47 = 2.77

# Bottleneck Resolving of Stream Program Contd.

- **Current state of the art**
  - Integer Linear Programming
    - Intractable

- **How to find a fast and good solution?**
  - Heuristics
  - Optimal

# Our Work

- A data rate transfer model to detect and eliminate bottlenecks

- We separate the bottleneck elimination from the actor allocation

- Heuristics to solve bottleneck problem efficiently

# Our Data Transfer Model

- Throughput depends on the data rate of the actors (maximize)

- Data transfer model forms a system of sim. functional linear equation

- Compute a closed form of the output data rate

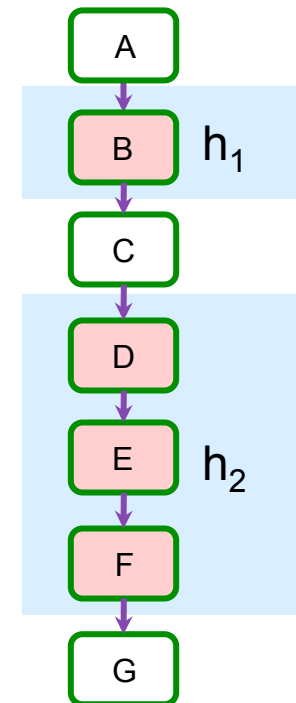- We also consider a processor utilization function for each actor

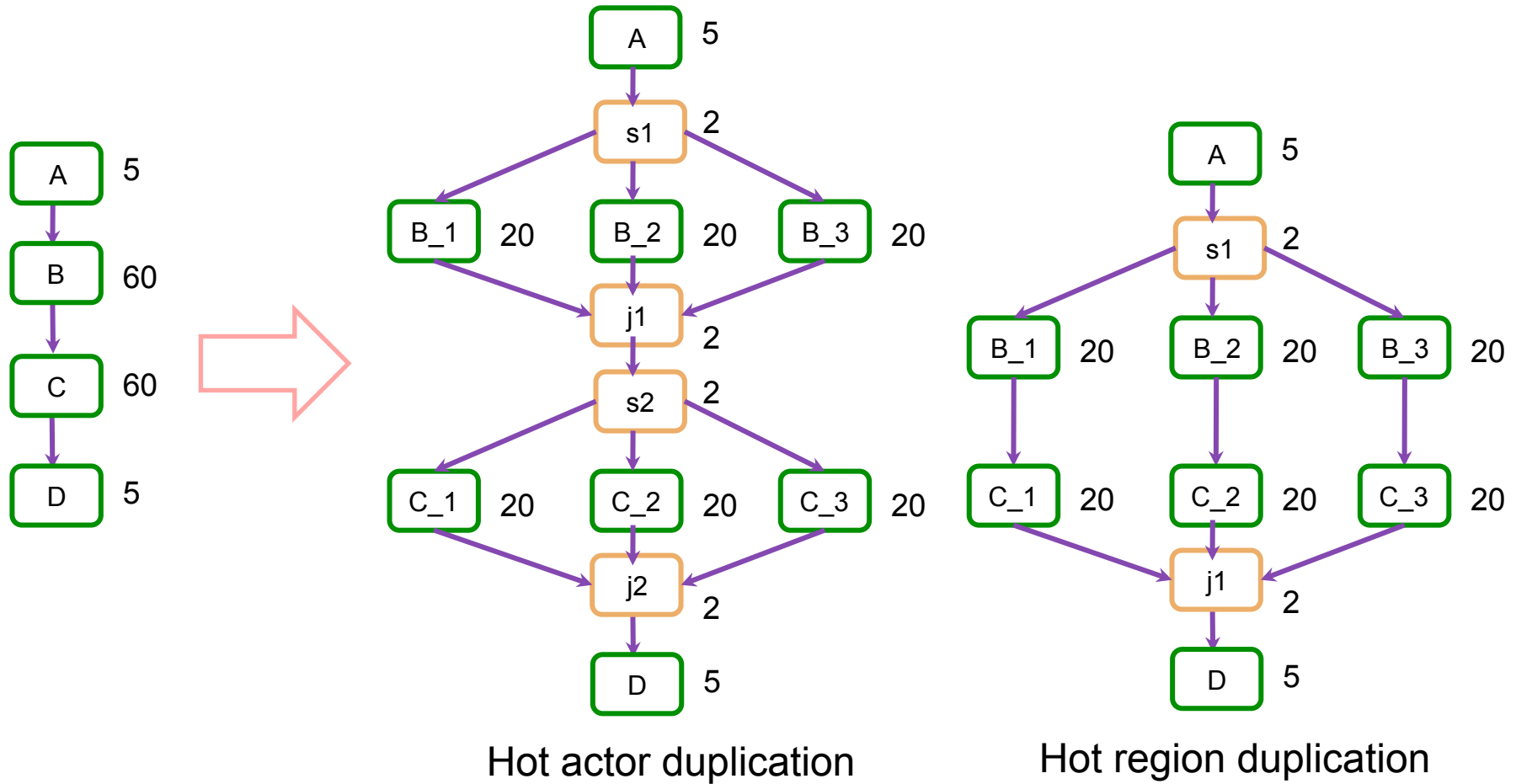$$z \longrightarrow \boxed{A} \xrightarrow[\phantom{x}]{1 \quad 5} \boxed{B} \xrightarrow[\phantom{x}]{1 \quad 1} \boxed{C}$$

$$x_A = z \qquad x_B = 0.2x_A \qquad x_C = x_B$$

$$x'_A = z \qquad x'_B = 0.2z \qquad x'_C = 0.2z$$

# Bottleneck Analysis

- ## The throughput is limited by
  - Processor capacity of the cores
  - Memory bandwidth
- ## A quantitative analysis determines
  - An upper bound of the throughput imposed by an actor
  - An upper bound of the throughput imposed by the parallel system
- ## Hot actor
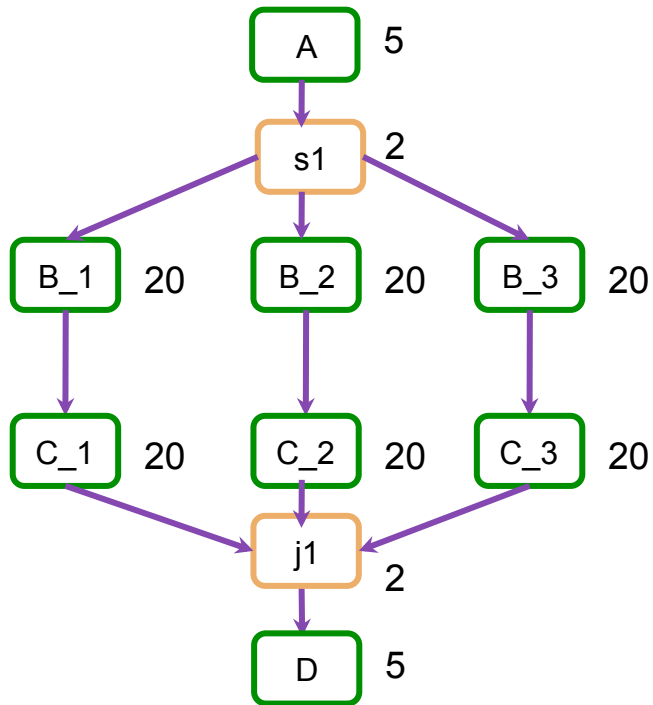  - Upper bound (actor) < upper bound (system)

# Hot Region

- Maximal connected subgraph
  $h = (V', E')$ where $V' \subseteq V, E' \subseteq E$
  and each $i \in V'$ is hot and
  stateless

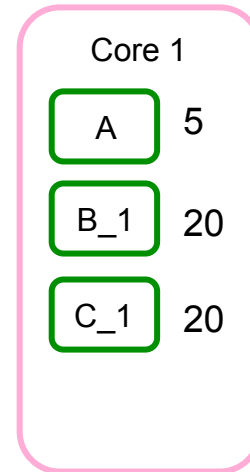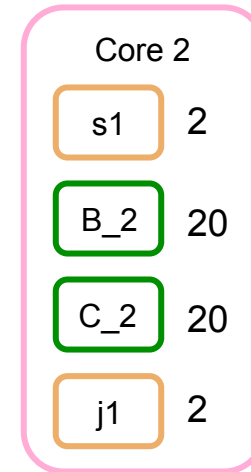# Resolving Bottleneck Options



Hot actor duplication

Hot region duplication

# Region Duplication further Increases Performance



A — 5
s1 — 2
B_1 — 20  B_2 — 20  B_3 — 20
C_1 — 20  C_2 — 20  C_3 — 20
j1 — 2
D — 5

Mapping

Core 1
A — 5
B_1 — 20
C_1 — 20
T = 45s
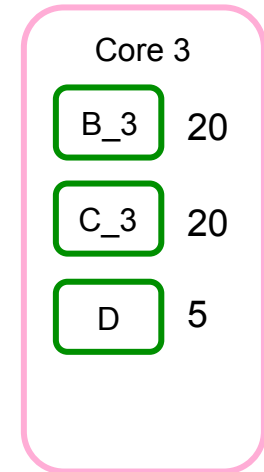
Core 2
s1 — 2
B_2 — 20
C_2 — 20
j1 — 2
T = 44s
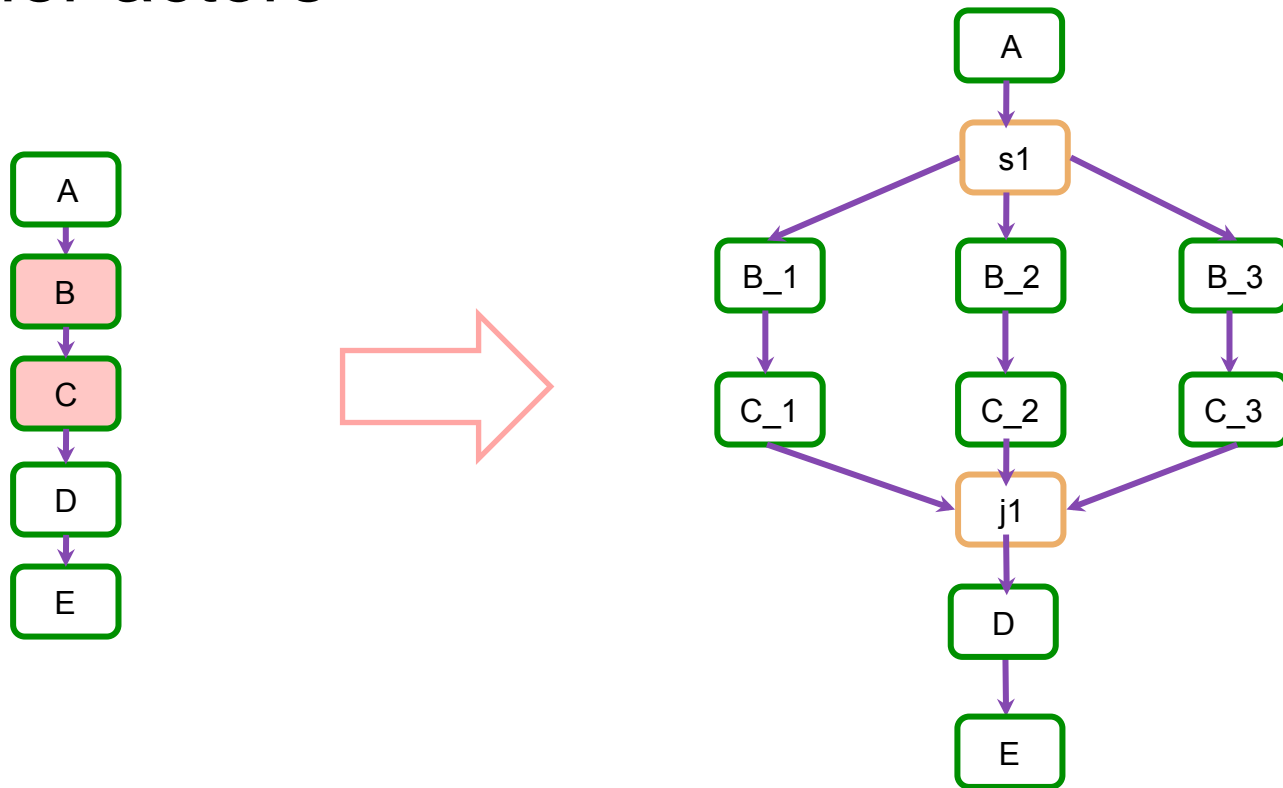
Core 3
B_3 — 20
C_3 — 20
D — 5
T = 45s

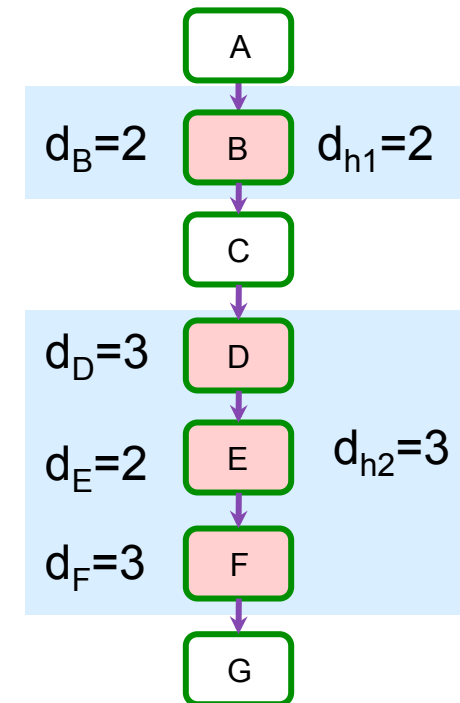Make span = 45s, Speedup = 130/45 = 2.89

# Cascading Effect of Duplication

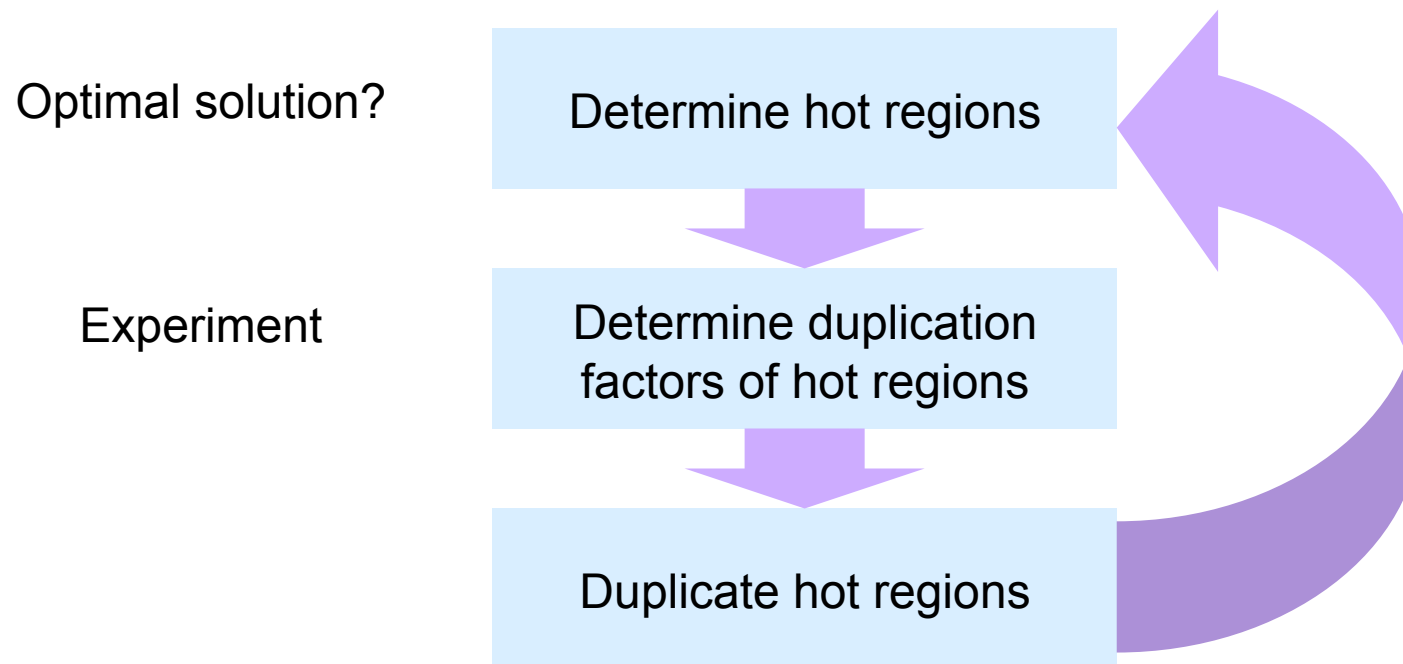- Actors may become hot due to duplication of other actors

# Duplication Factor of an Actor and a Hot Region

- The # of times the actor needs to be duplicated $d_i > 1$

- Maximum duplication factor of the actors of the hot region



$d_B=2$  B  $d_{h1}=2$

A

C

$d_D=3$  D

$d_E=2$  E  $d_{h2}=3$

$d_F=3$  F

G

# Heuristics to Resolve Bottlenecks

Optimal solution?

Determine hot regions

Experiment

Determine duplication
factors of hot regions

Duplicate hot regions

# Summary

- A simple quantitative analysis to detect and eliminate bottlenecks

- We separate the bottleneck elimination from the actor allocation

- Heuristics to eliminate bottlenecks

# Related Works

[1] Static Scheduling of SDF Programs for DSP [Lee '87]

[2] StreamIt: A language for streaming applications [Thies '02]

[3] Phased Scheduling of Stream Programs [Thies '03]

[4] Exploiting Coarse Grained Task, Data, and Pipeline Parallelism in Stream Programs [Thies '06]

[5] Orchestrating the Execution of Stream Programs on Cell [Scott '08]

[6] Software Pipelined Execution of Stream Programs on GPUs [Udupa '09]

[7] Synergistic Execution of Stream Programs on Multicores with Accelerators [Udupa '09]

# Questions?