# Flexible Filters for High Performance Embedded Computing

Rebecca Collins and Luca Carloni
Department of Computer Science
Columbia University

## **Motivation**

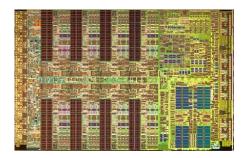


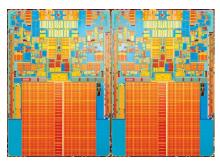




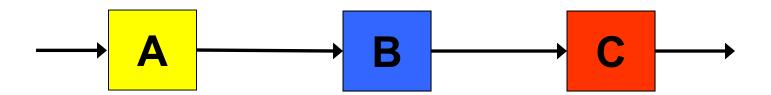




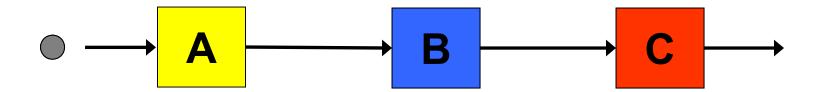




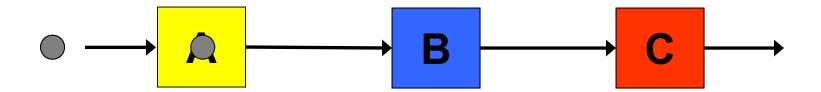




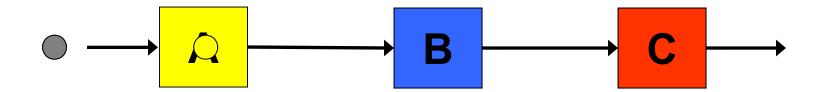
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



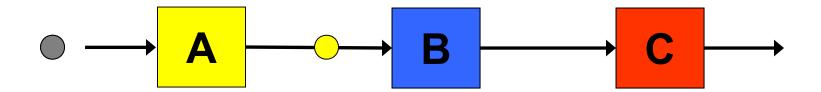
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



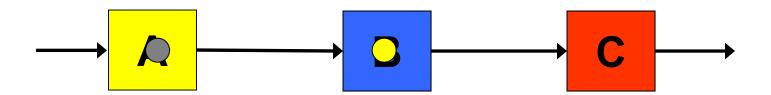
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



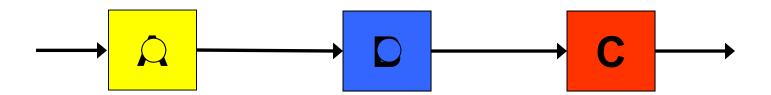
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



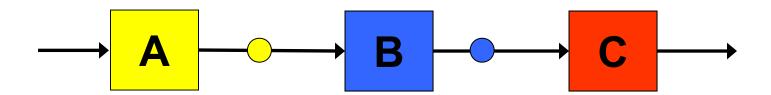
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



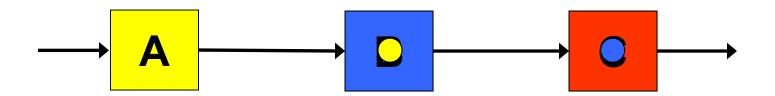
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



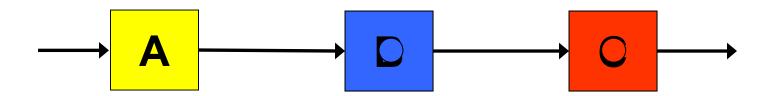
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



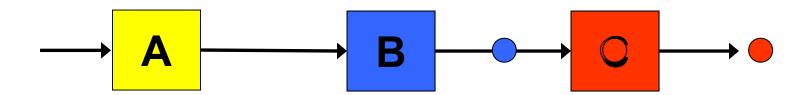
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



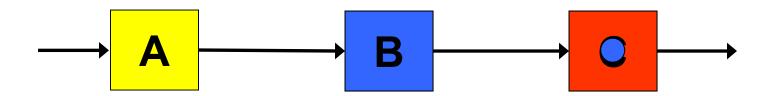
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



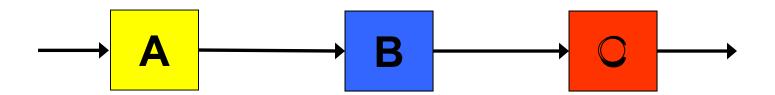
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



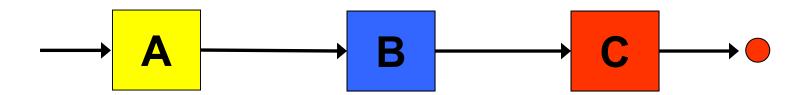
- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications

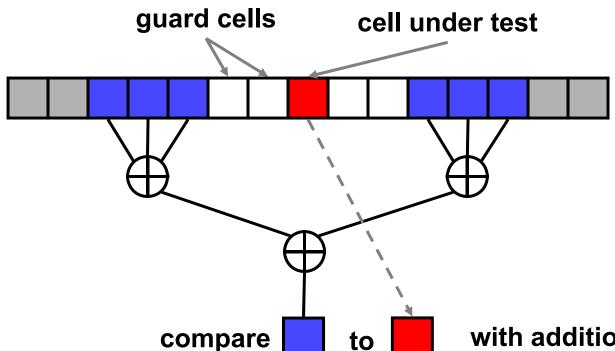


- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications



- Stream Programming model
  - filter: a piece of sequential programming
  - channels: how filters communicate
  - token: an indivisible unit of data for a filter
- Examples: Signal processing, image processing, embedded applications

## **Example: Constant False-Alarm** Rate (CFAR) Detection

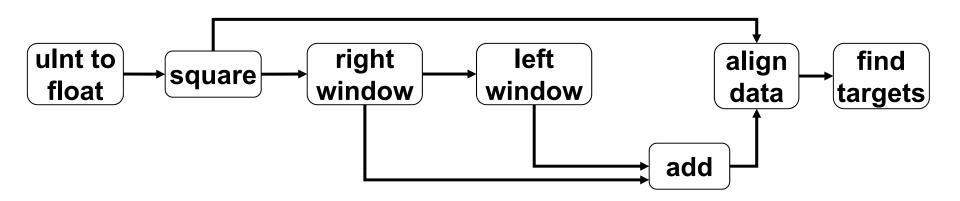


#### **HPEC Challenge**

http://www.ll.mit.edu/HPECchallenge/

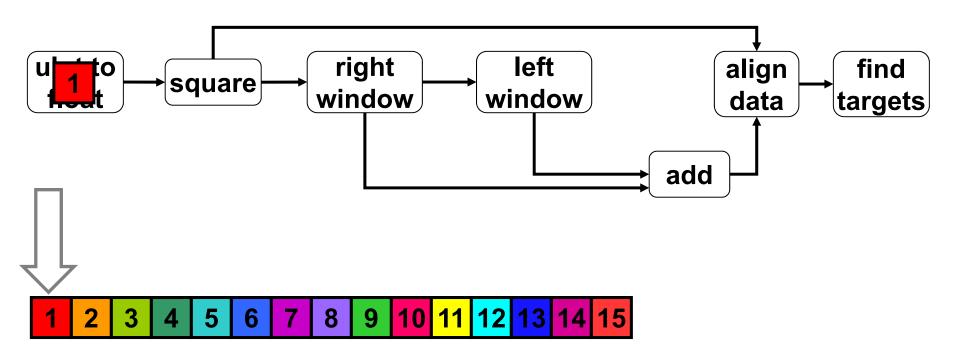
with additional factors

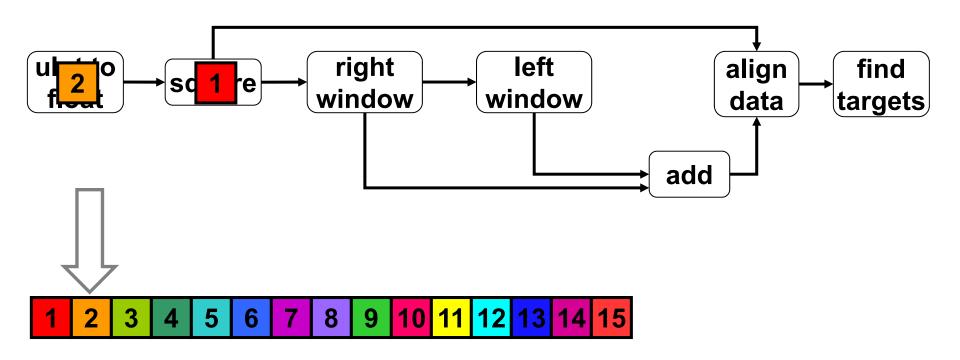
- number of gates
- threshold (µ)
- number of guard cells
- rows and other dimensional data

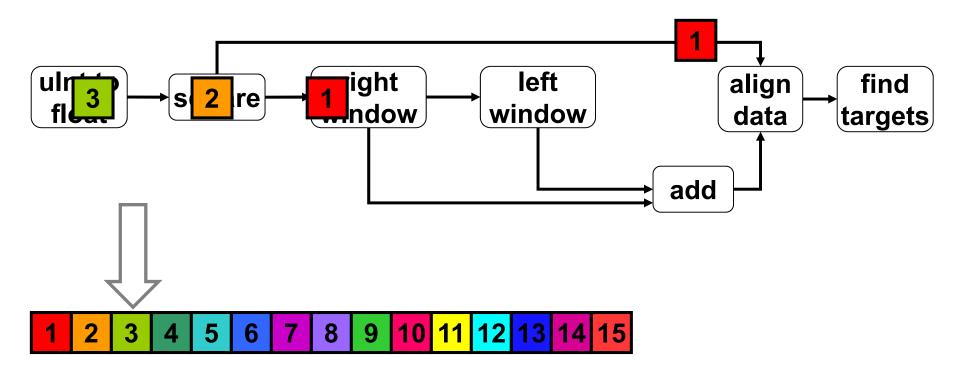


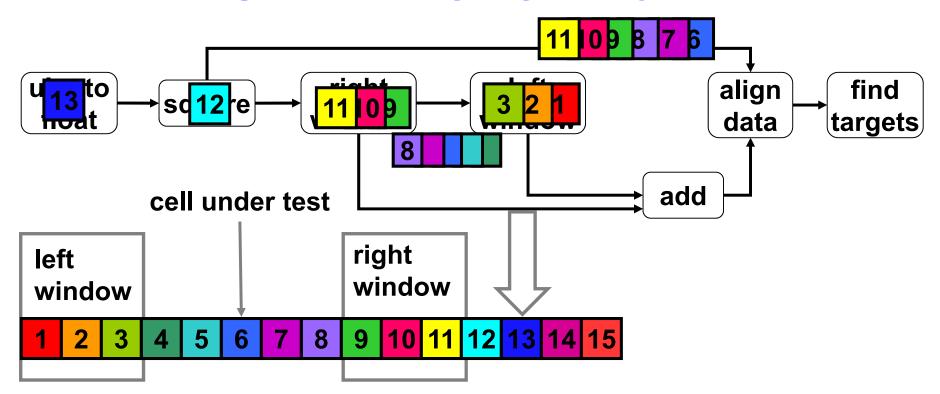
#### **Data Stream:**

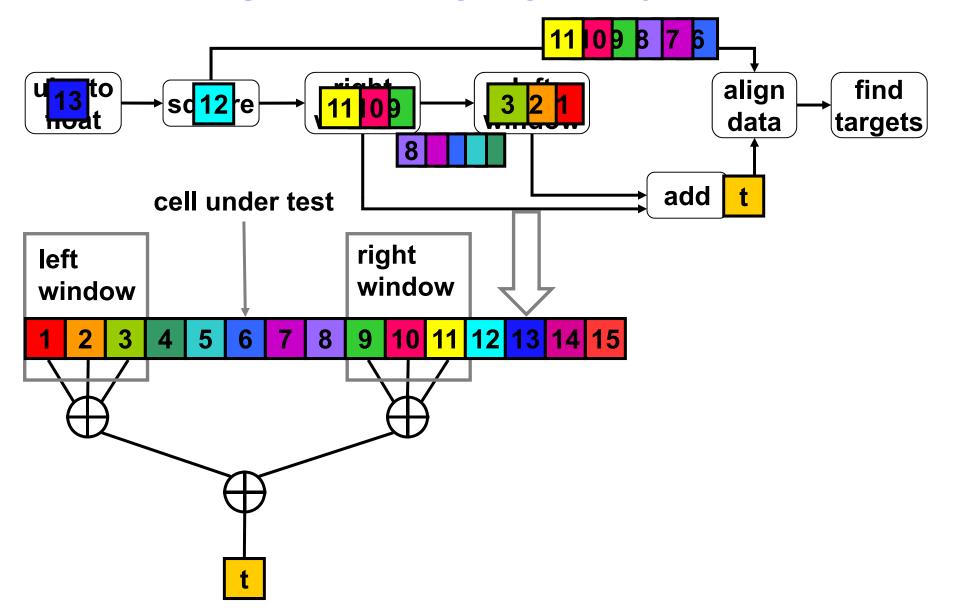


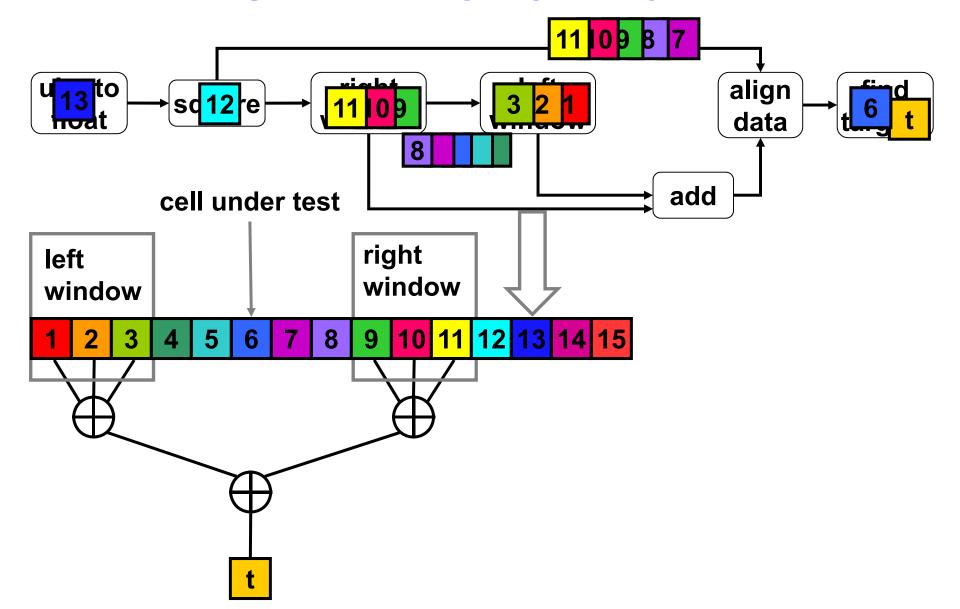




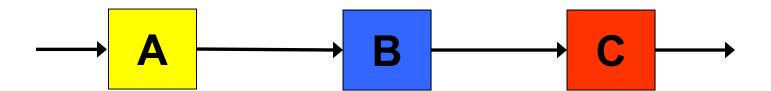




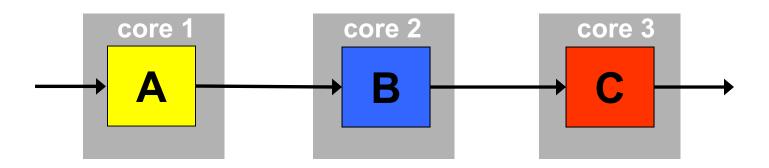


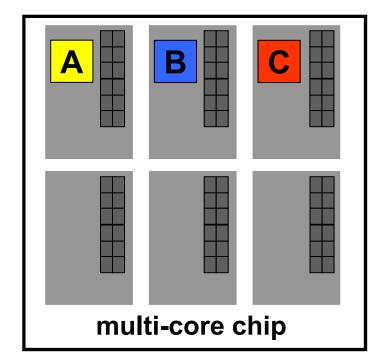


# **Mapping**



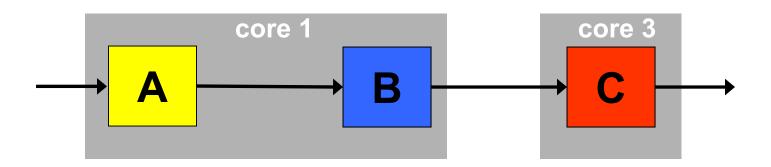
## **Mapping**

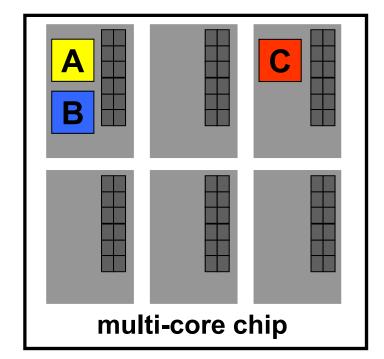




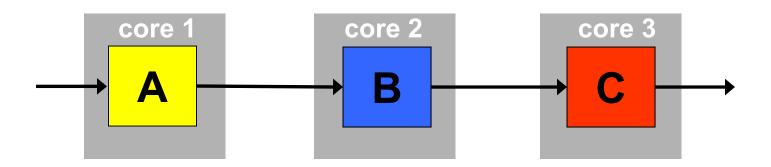
Throughput: rate of processing data tokens

## **Mapping**

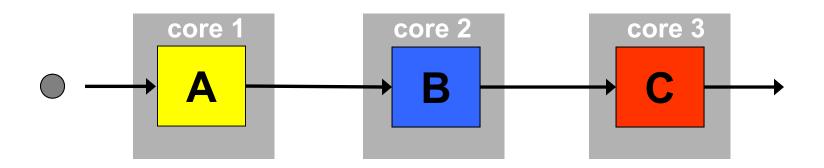




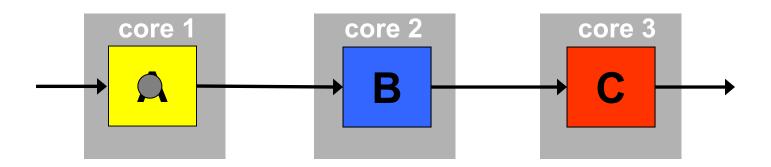
Throughput: rate of processing data tokens



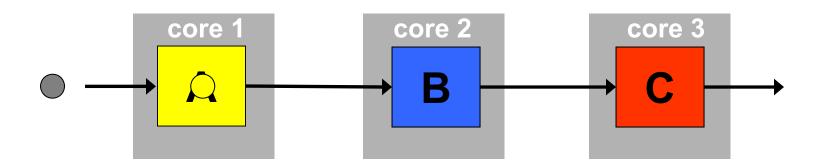
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



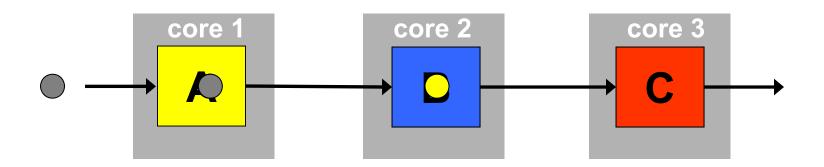
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



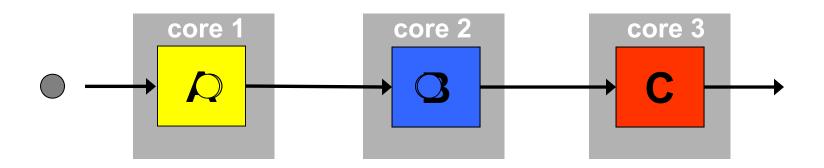
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



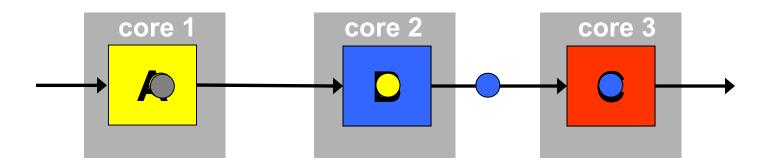
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



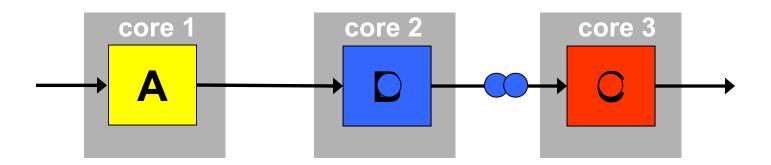
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



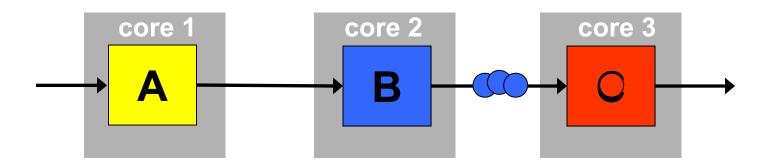
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes

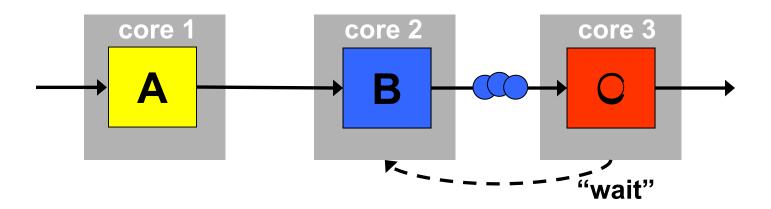


- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes



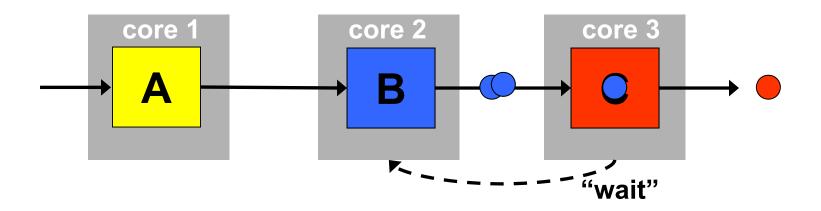
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes

#### **Unbalanced Flow**



- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes

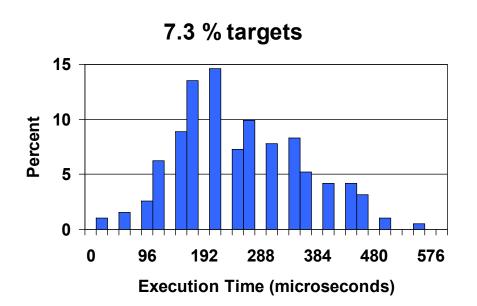
#### **Unbalanced Flow**

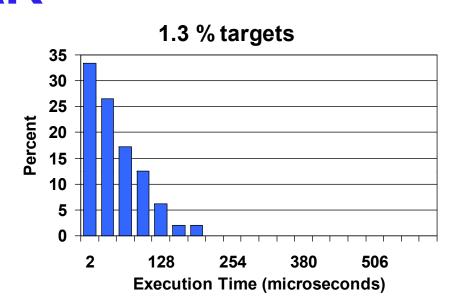


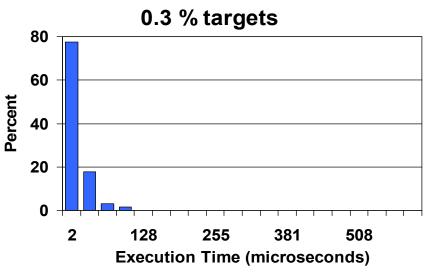
- Bottlenecks reduce throughput
  - Caused by backpressure
    - inherent algorithmic imbalances
    - data-dependent computational spikes

# Data Dependent Execution Time: CFAR

- Using Set 1 from the HPEC CFAR Kernel Benchmark
- Over a block of about 100 cells
- Extra workload of 32 microseconds per target



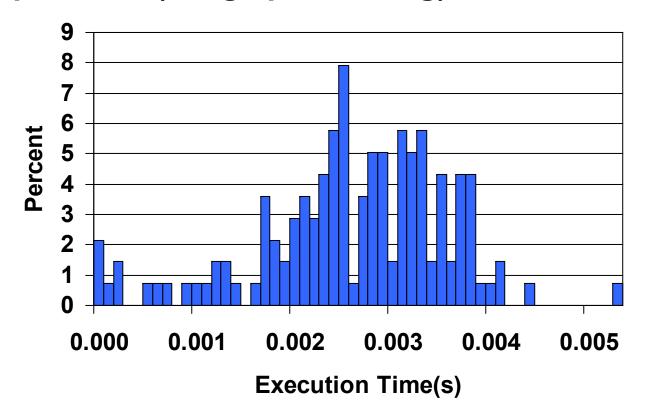


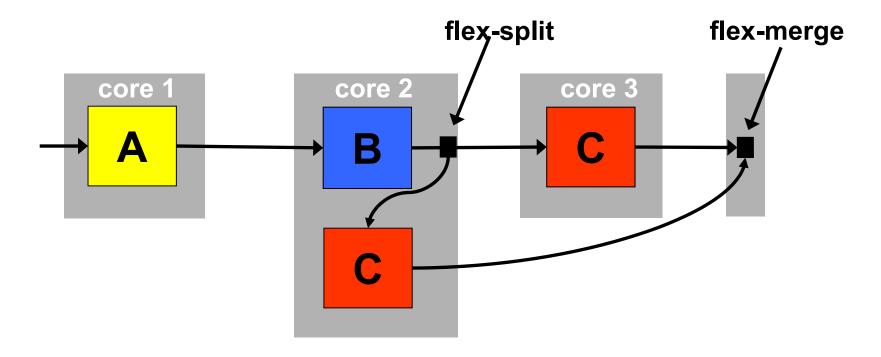


### **Data Dependent Execution Time**

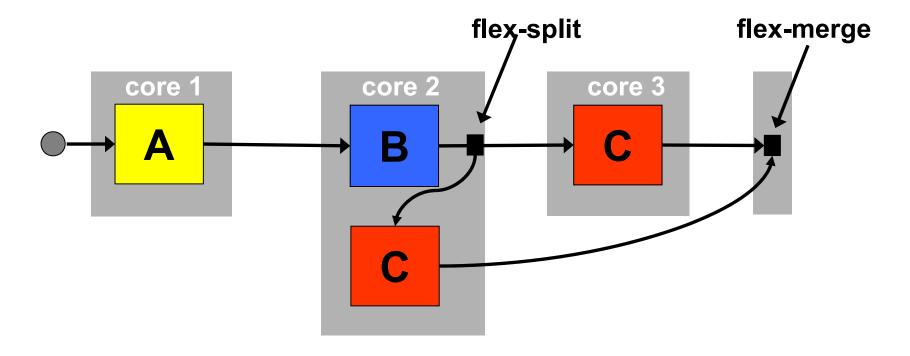
#### Some other examples:

- Bloom Filters (Financial, Spam detection)
- Compression (Image processing)

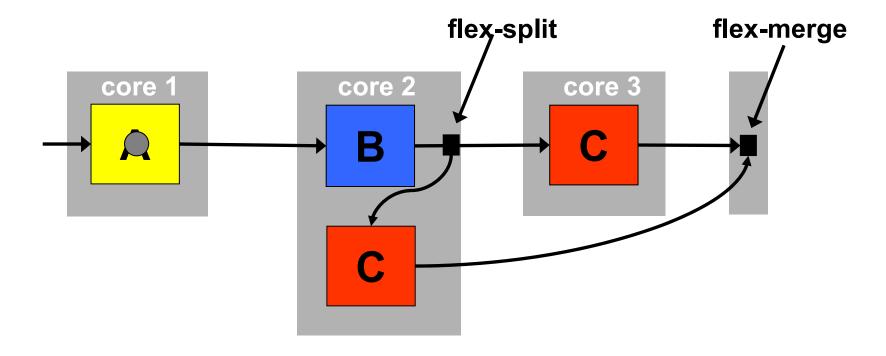




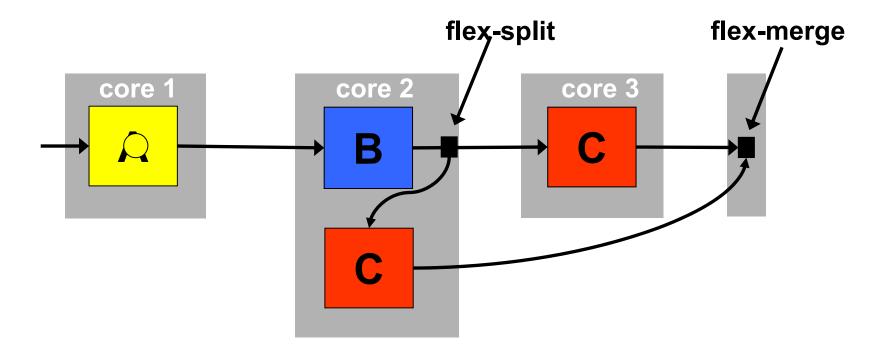
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



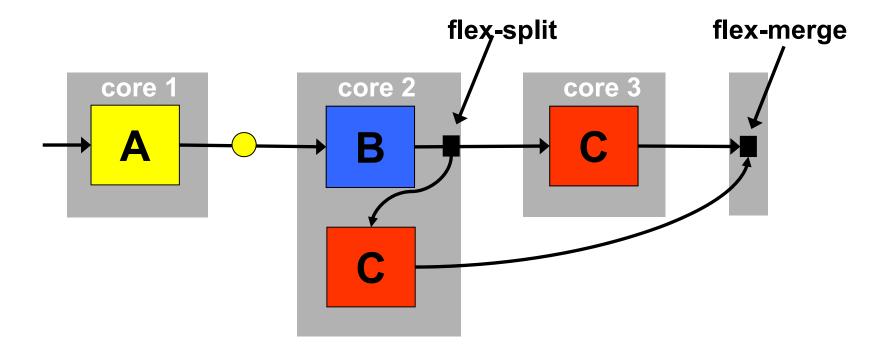
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



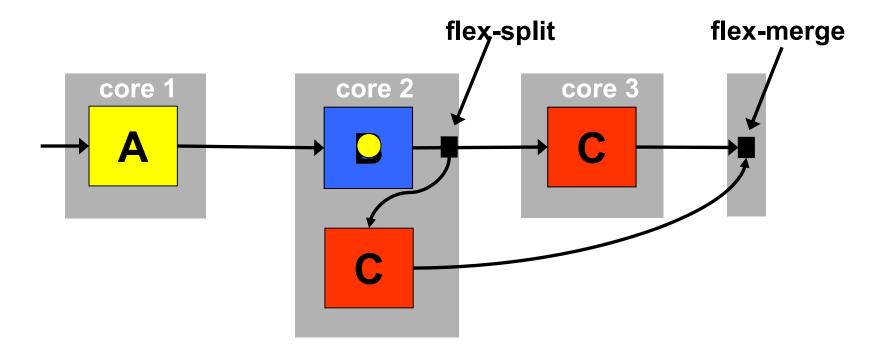
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



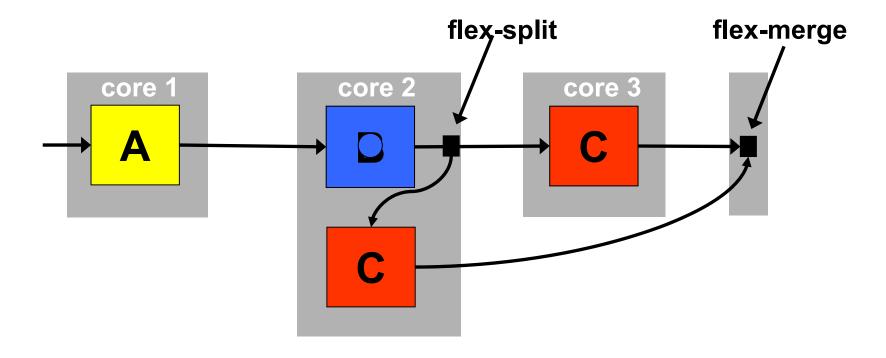
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



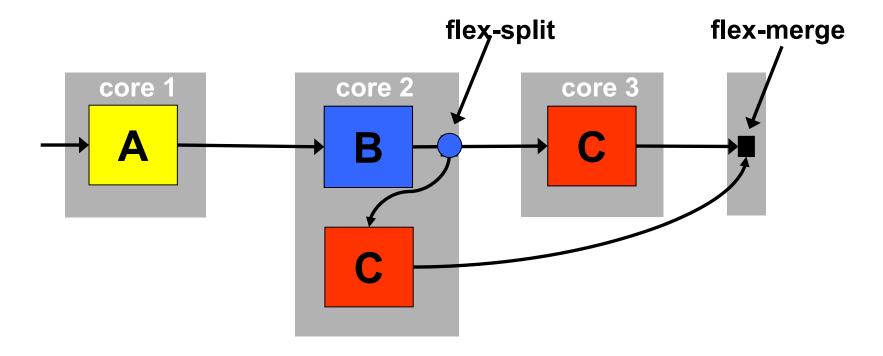
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



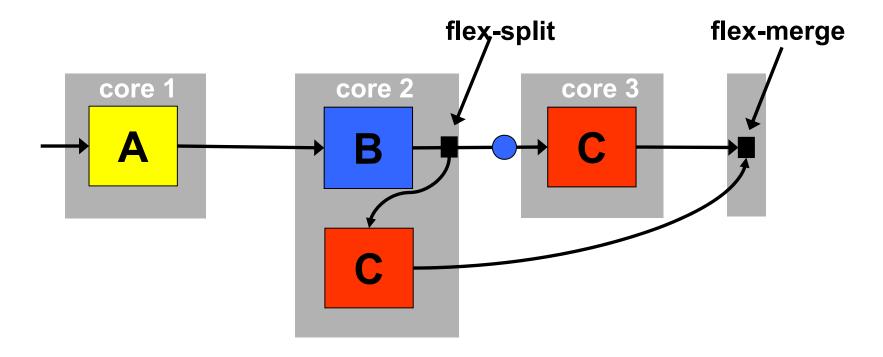
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



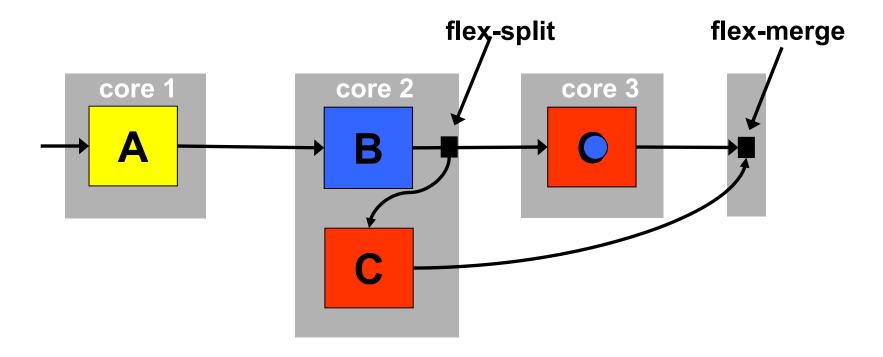
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



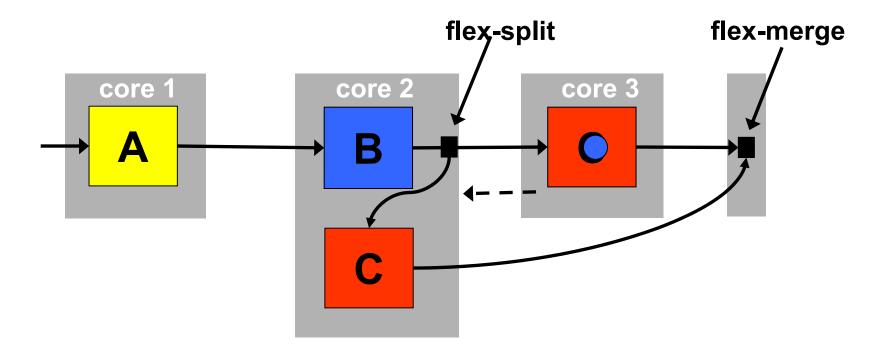
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



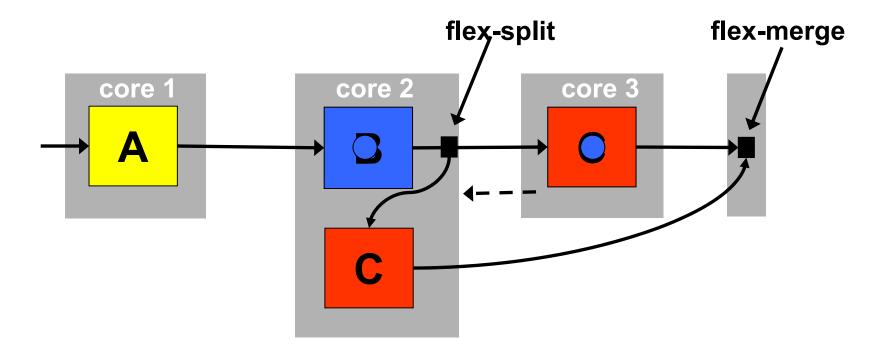
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



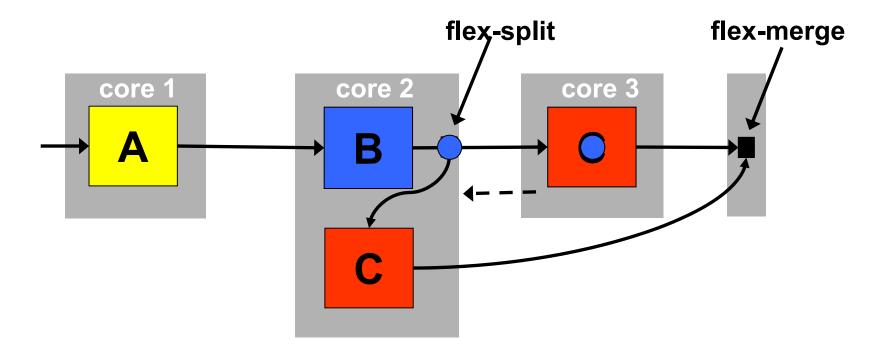
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



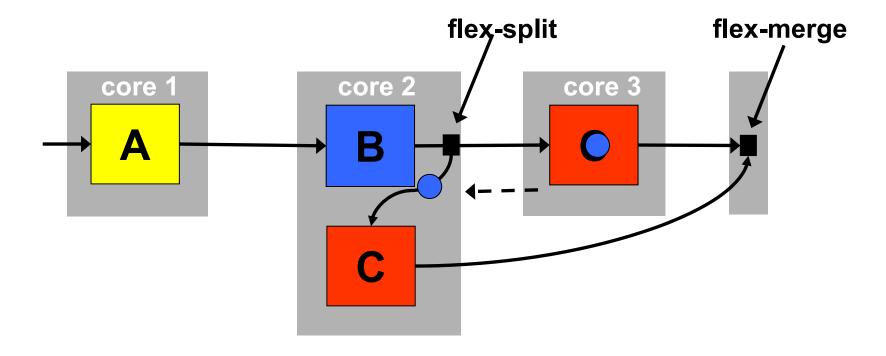
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



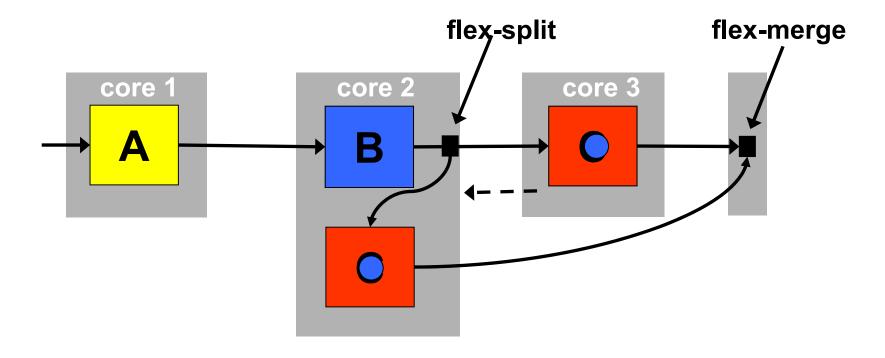
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



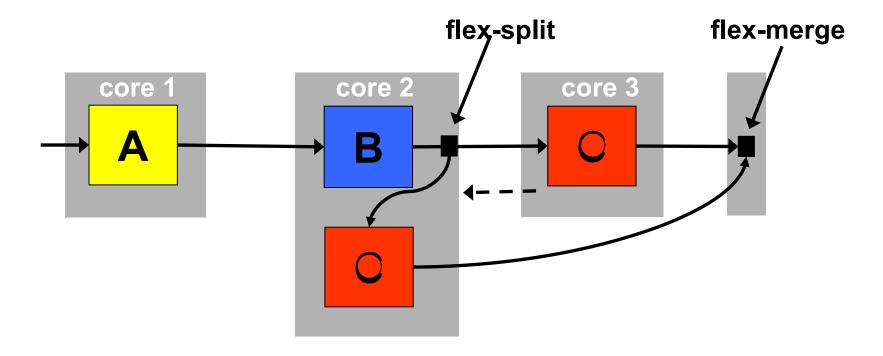
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



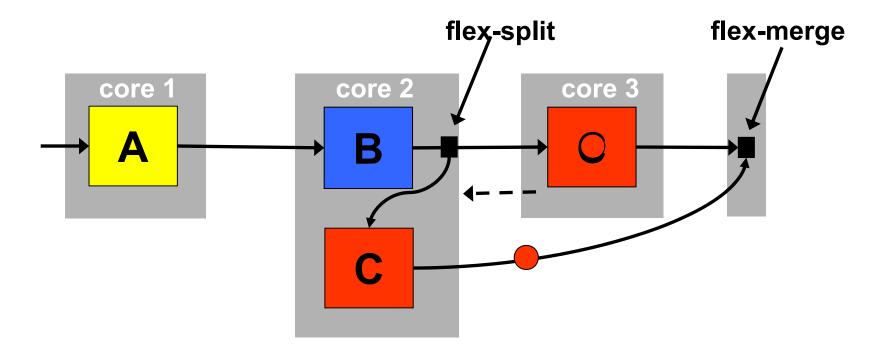
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



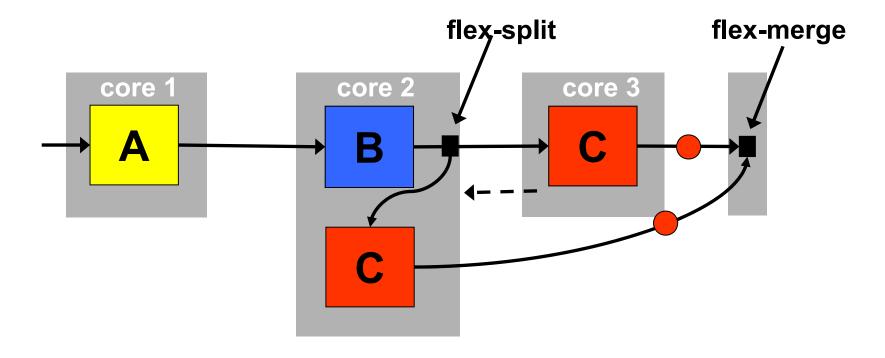
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



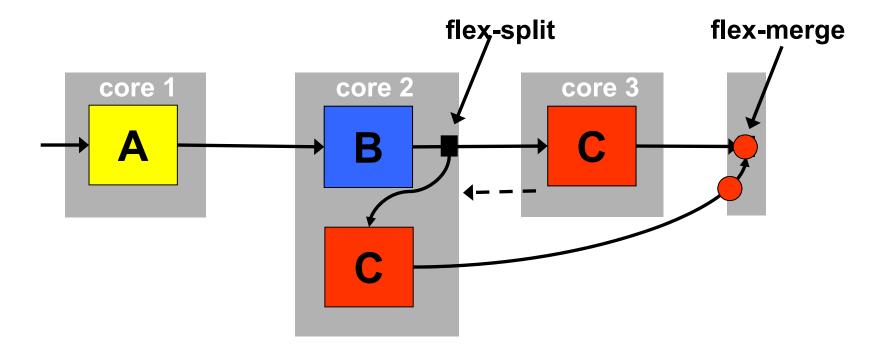
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



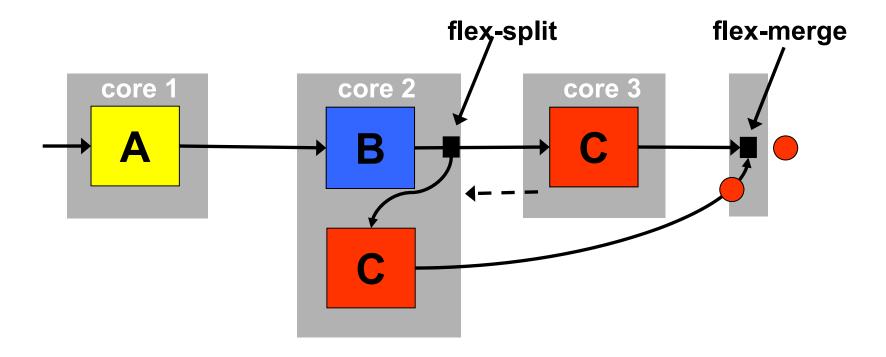
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



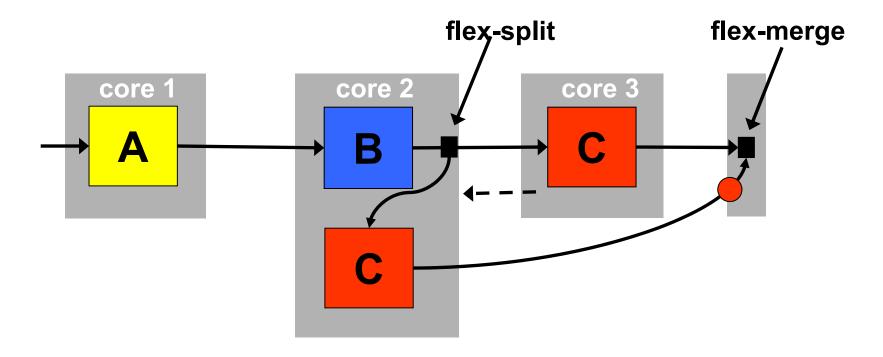
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



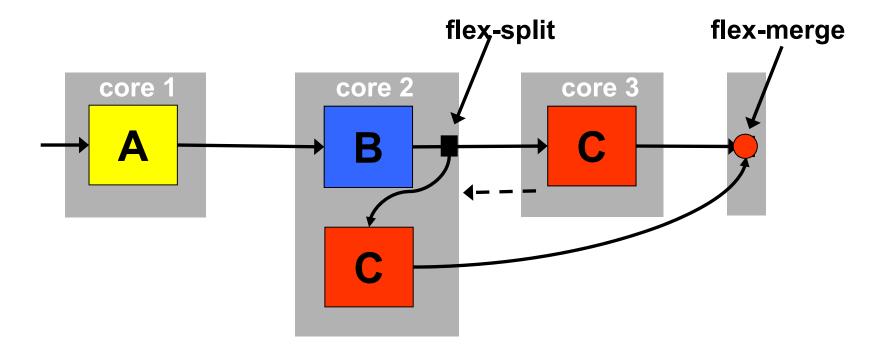
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



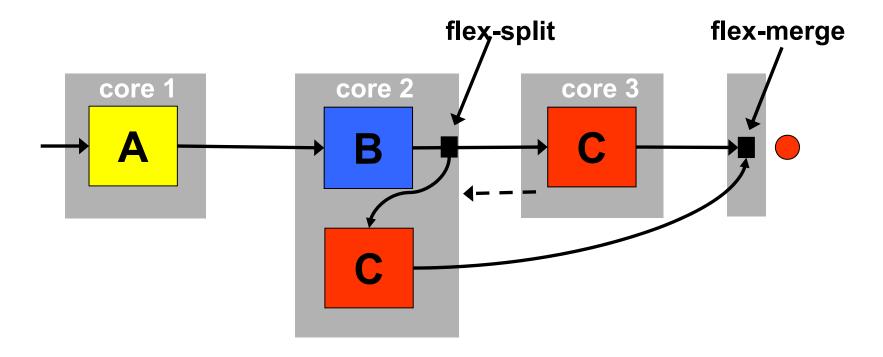
- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation



- Unused cycles on B are filled by working ahead on filter C with the data already present on B
- Push stream flow upstream of a bottleneck
- Semantic Preservation

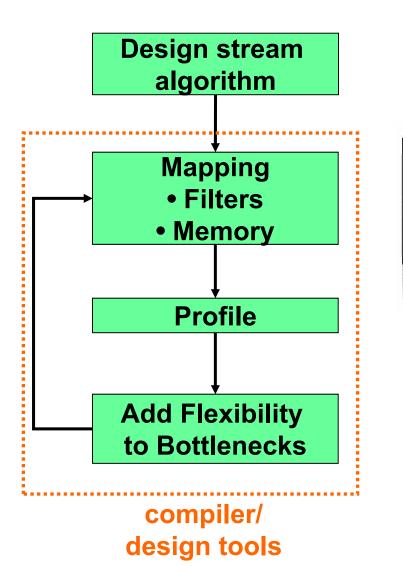
#### **Related Works**

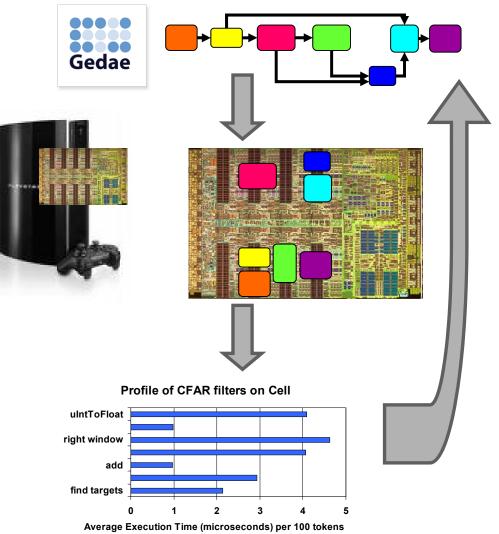
- Static Compiler Optimizations
  - StreamIt [Gordon et al, 2006]
- Dynamic Runtime Load Balancing
  - Work Stealing/Filter Migration
    - [Kakulavarapu et al., 2001]
    - Cilk [Frigo et al., 1998]
    - Flux [Shah et al., 2003]
    - Borealis [Xing et al., 2005]
  - Queue based load balancing
    - Diamond [Huston et al., 2005] distributed search, queue based load balancing, filter re-ordering
- Combination Static+Dynamic
  - FlexStream [Hormati et al., 2009] multiple competing programs

#### **Outline**

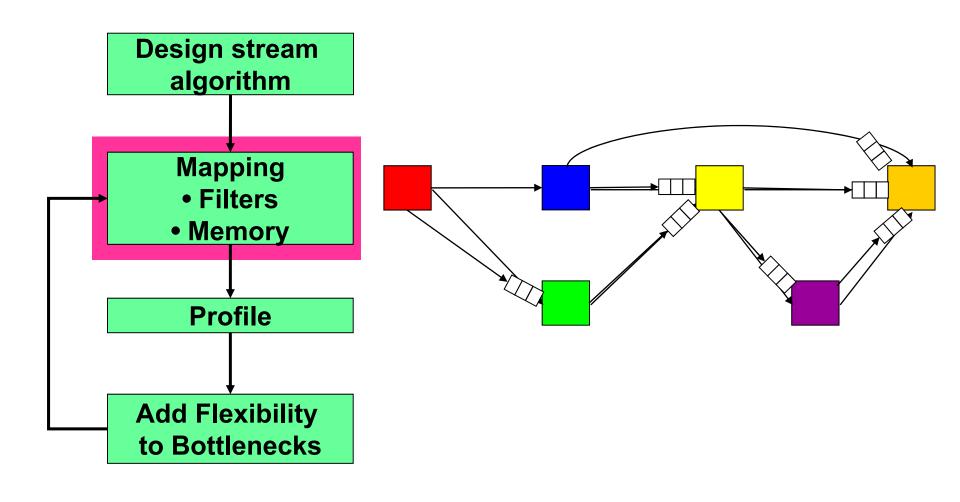
- Introduction
- Design Flow of a Stream Program with Flexibility
- Performance
- Implementation of Flexible Filters
- Experiments
  - CFAR Case Study

# Design Flow of a Stream Program with Flexible Filters

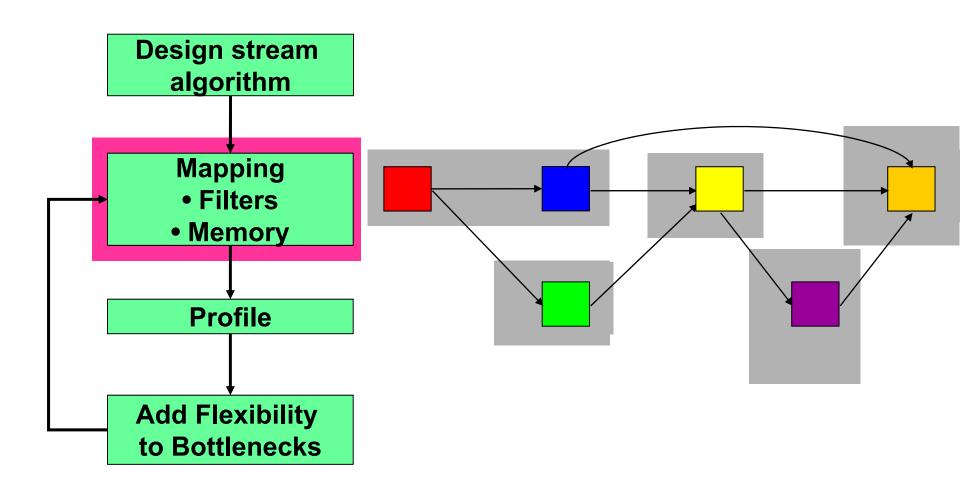




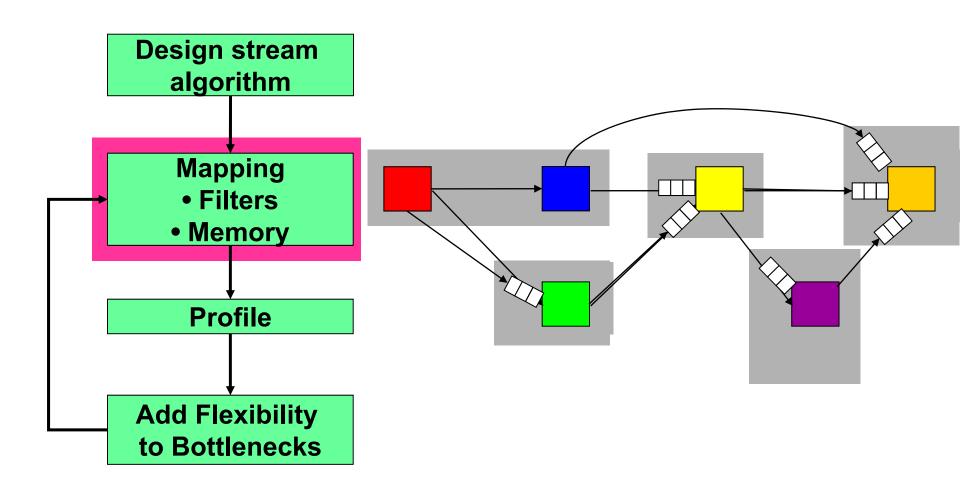
### **Design Considerations**



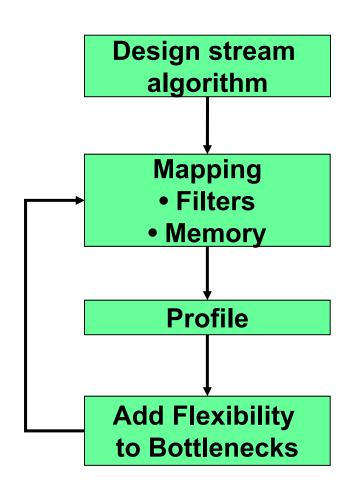
# **Design Considerations**

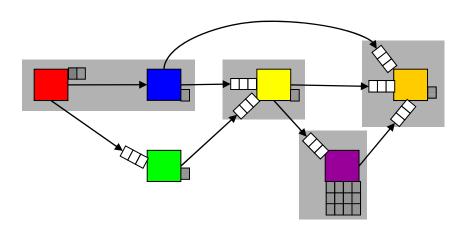


# **Design Considerations**

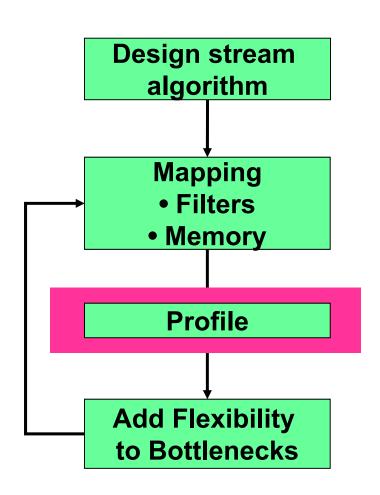


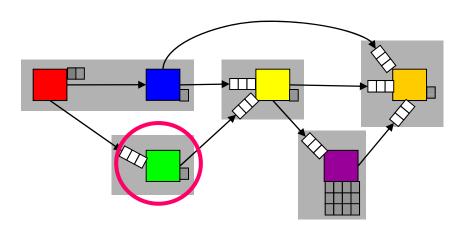
#### **Adding Flexibility**



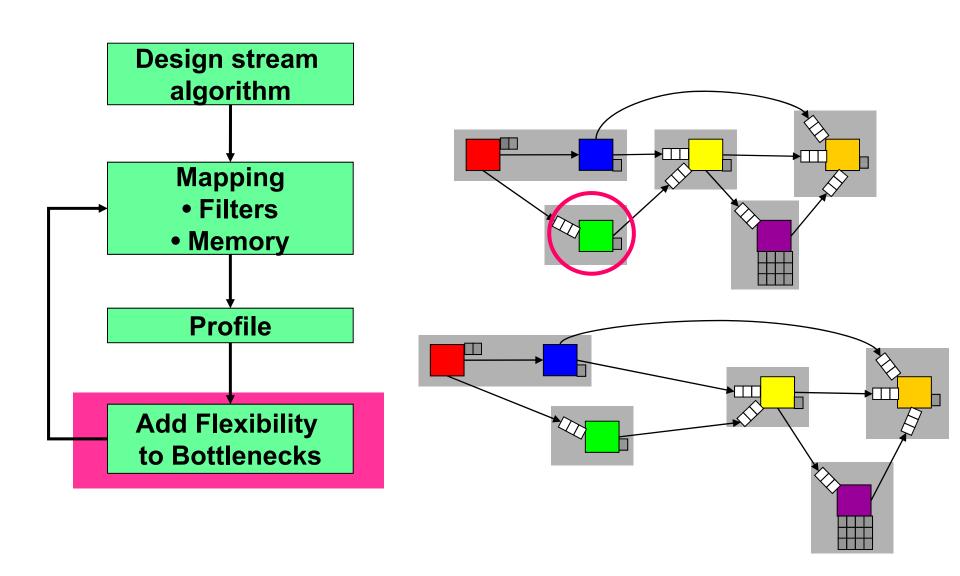


#### **Adding Flexibility**

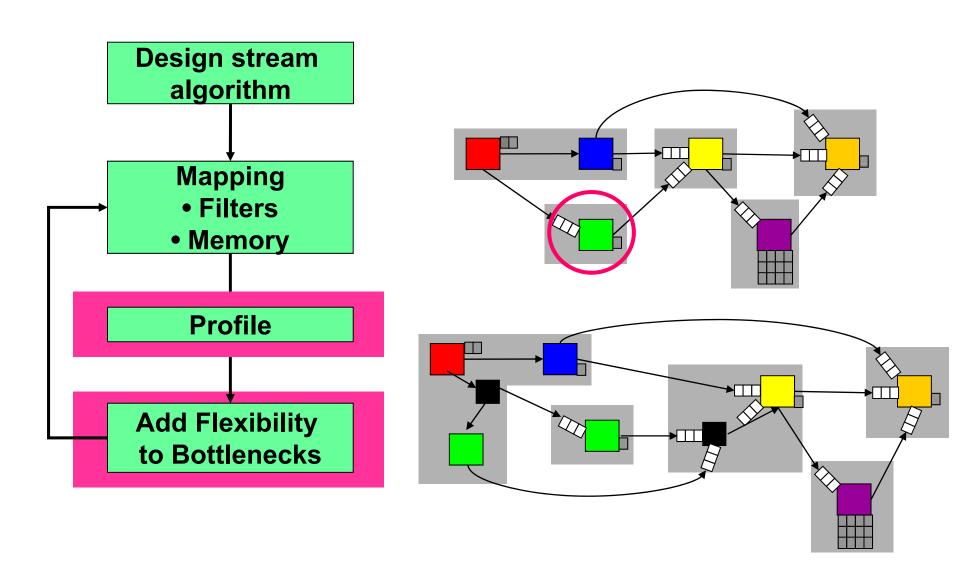




#### **Adding Flexibility**



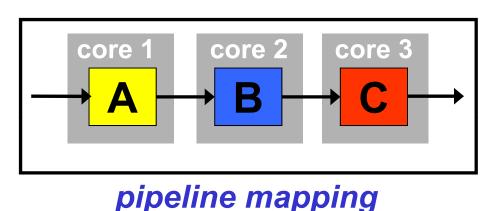
### **Adding Flexibility**

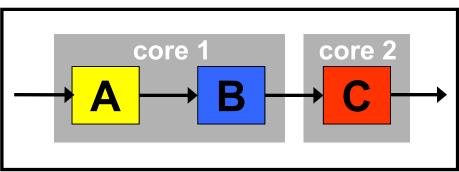


#### **Outline**

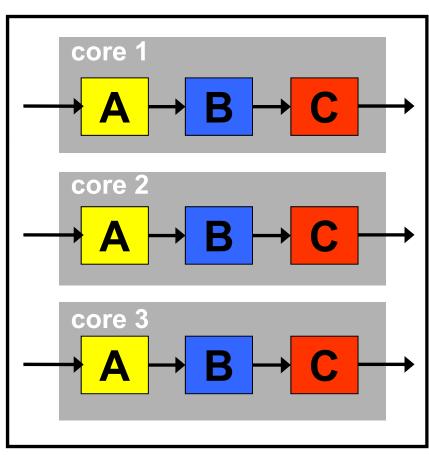
- Introduction
- Design Flow of a Stream Program with Flexibility
- Performance
- Implementation of Flexible Filters
- Experiments
  - CFAR Case Study

## Mapping Stream Programs to Multi-Core Platforms

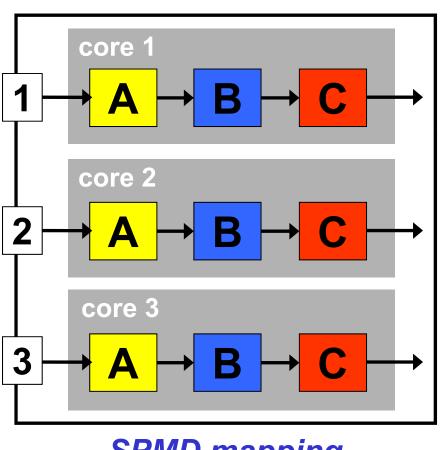






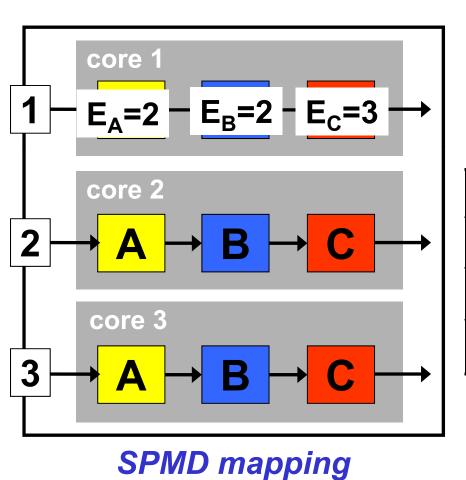


SPMD mapping



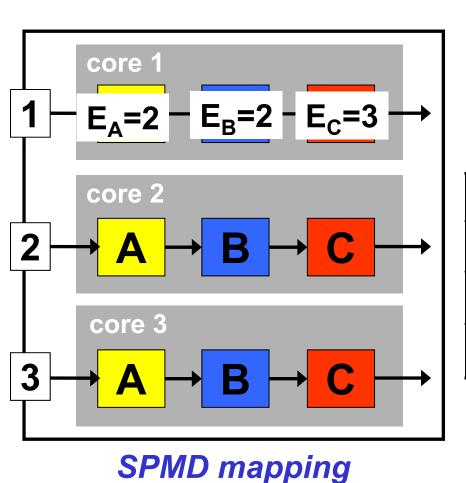
SPMD	map	ping
------	-----	------

	t0	t1	t2	t3	t4	t5	t6
core 1							
core 2							
core 3							



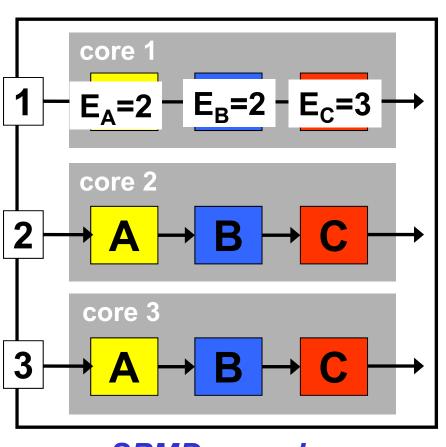
Suppose  $E_A = 2$ ,  $E_B = 2$ ,  $E_C = 3$ 

	t0	t1	<b>t2</b>	t3	t4	t5	t6
core 1							
core 2							
core 3							



Suppose  $E_A = 2$ ,  $E_B = 2$ ,  $E_C = 3$ 

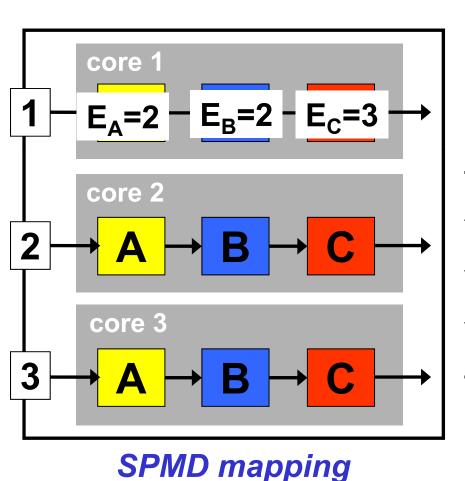
	t0	t1	t2	t3	t4	t5	t6
core 1	A	1					
core 2	A	2					
core 3	Α	3					



Suppose  $E_A = 2$ ,  $E_B = 2$ ,  $E_C = 3$ 

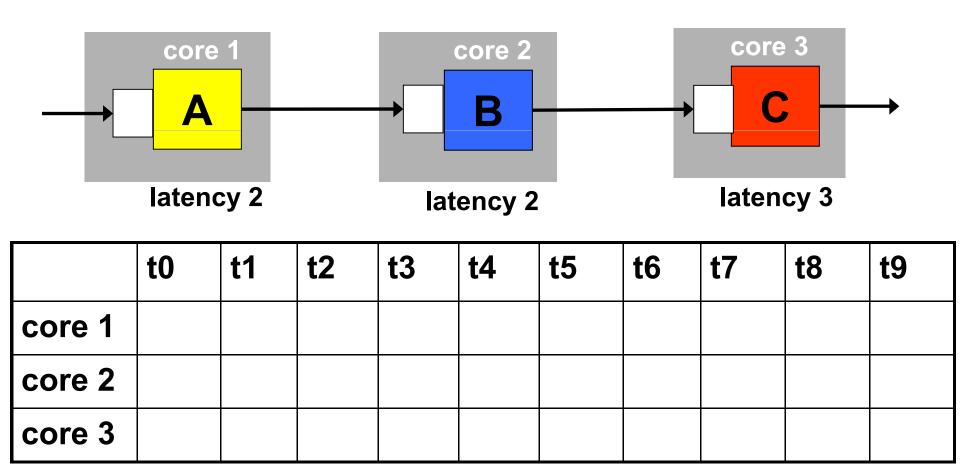
	t0	t1	t2	t3	t4	t5	t6
core 1	A	1	В	1			
core 2	A	2	В	2			
core 3	A	3	В	3			

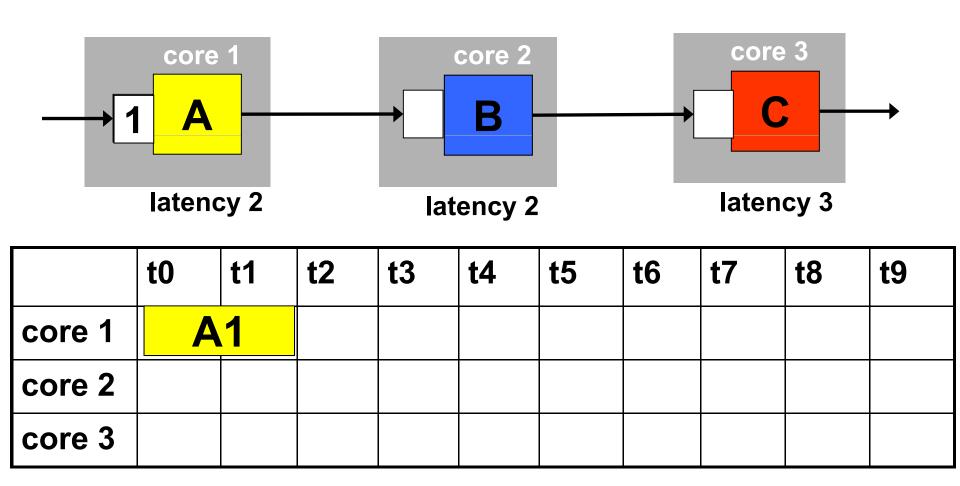
**SPMD** mapping

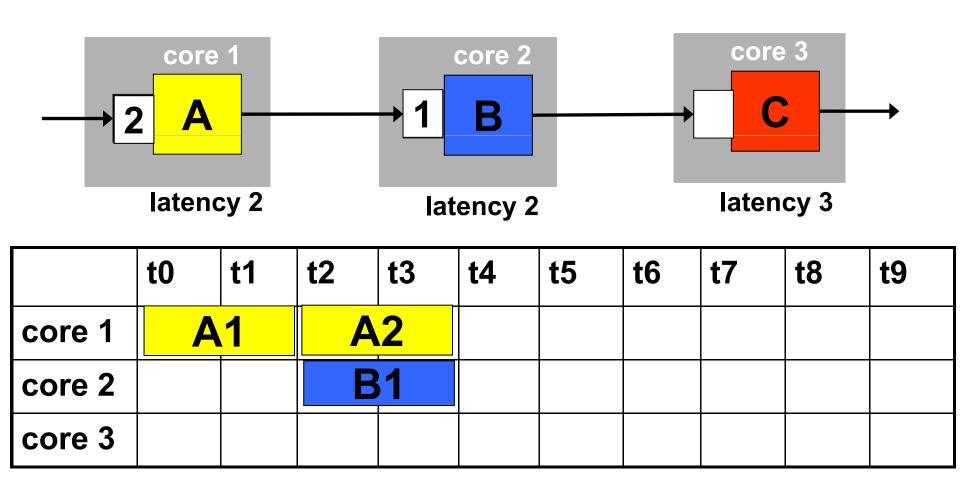


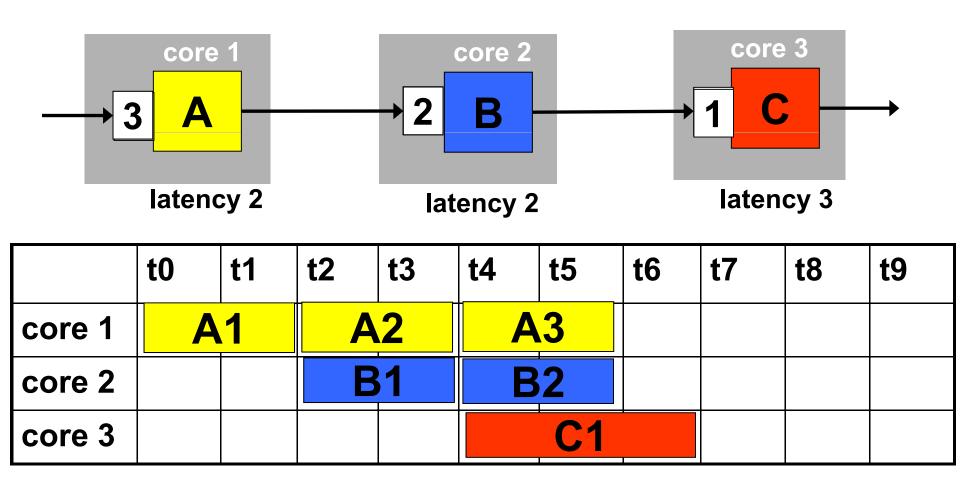
Suppose  $E_A = 2$ ,  $E_B = 2$ ,  $E_C = 3$ 

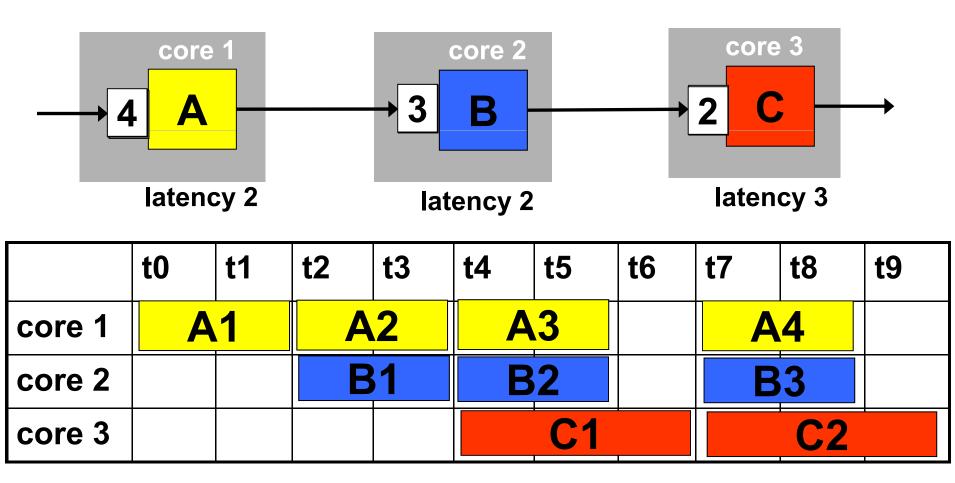
	t0	t1	t2	t3	t4	t5	t6
core 1	A	1	В	1		C1	
core 2	A	2	В	2		C2	
core 3	A	3	В	3		<b>C</b> 3	



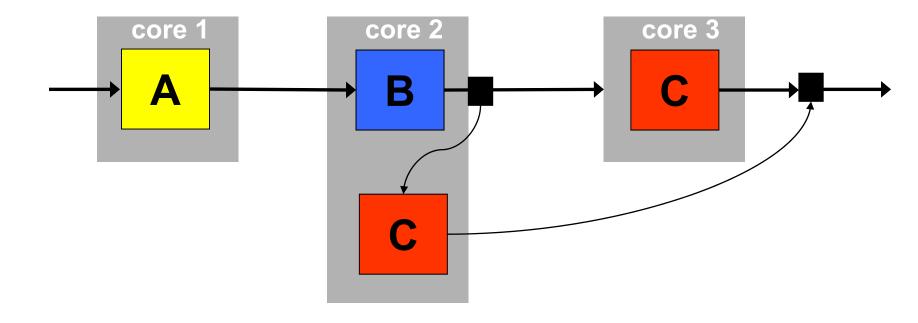


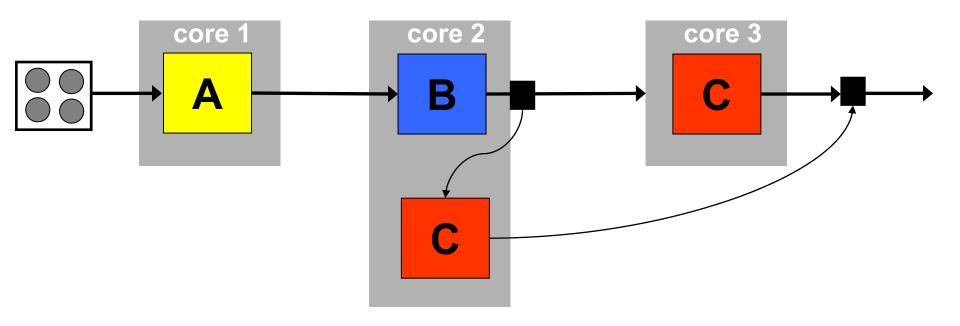


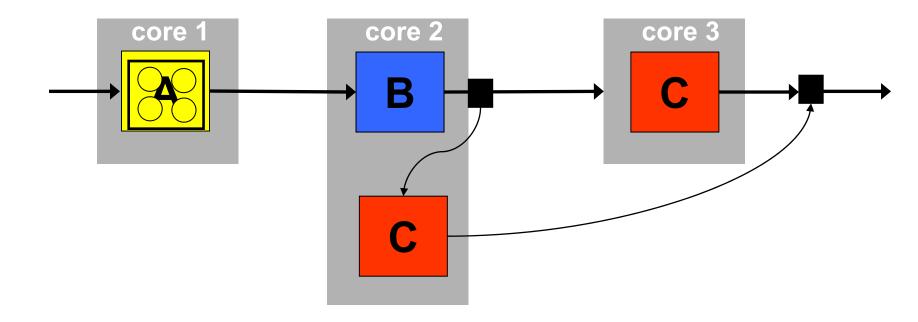


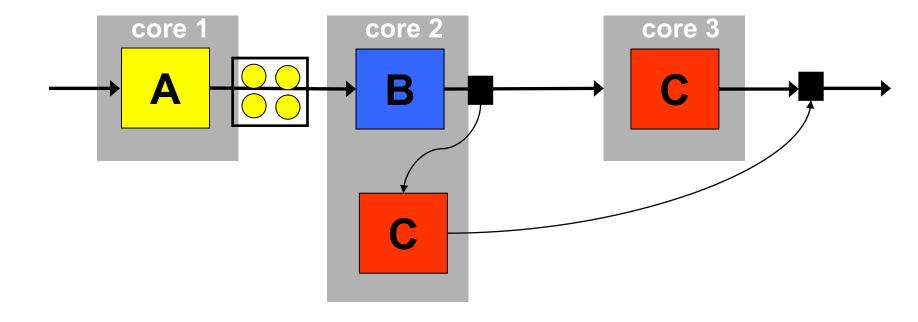


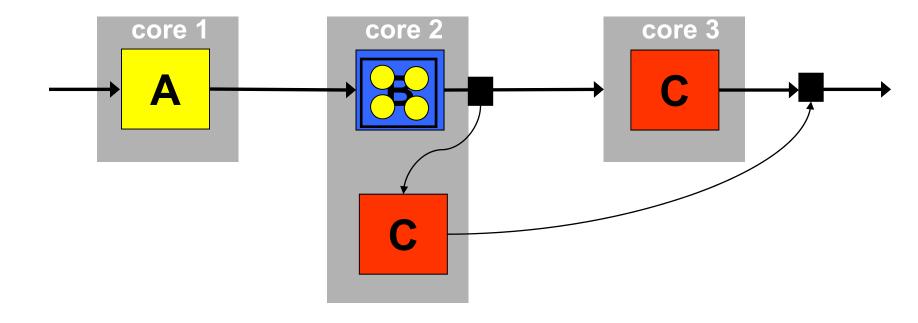
throughput = 1/3 = 0.333 < 0.429

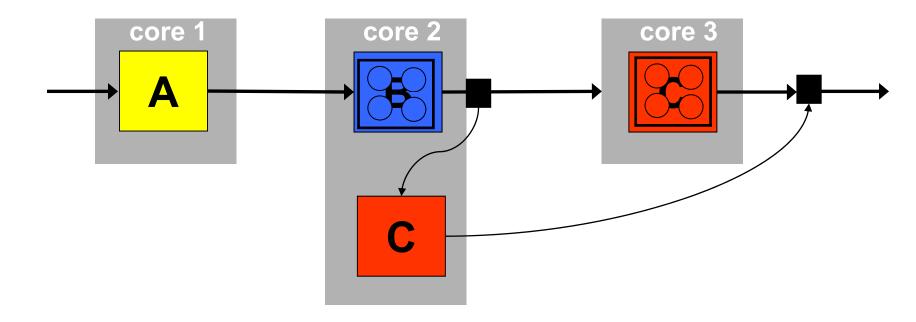


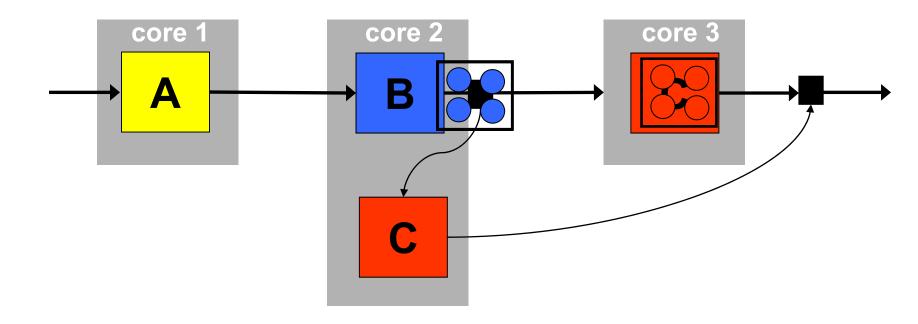


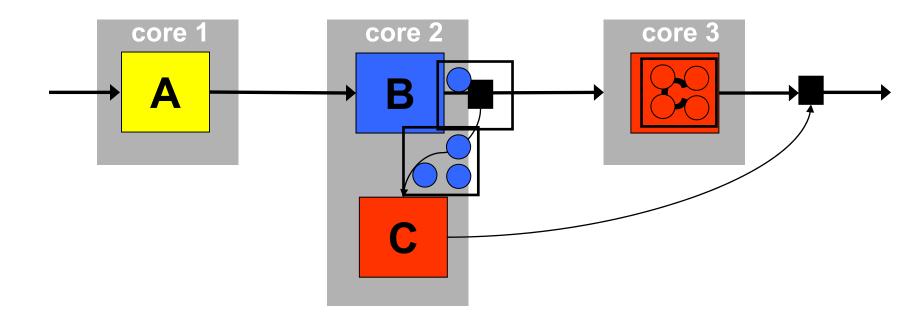


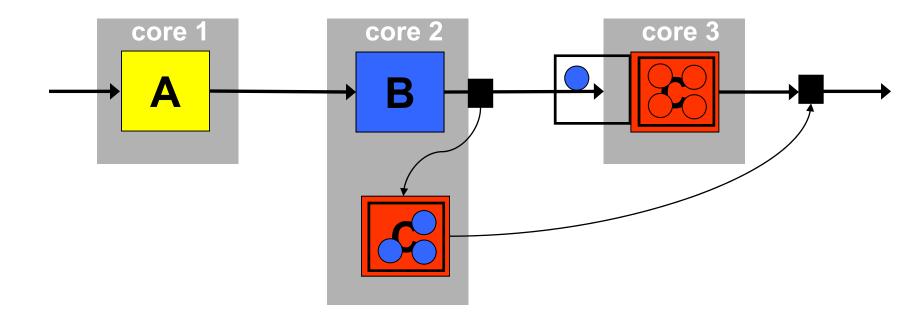


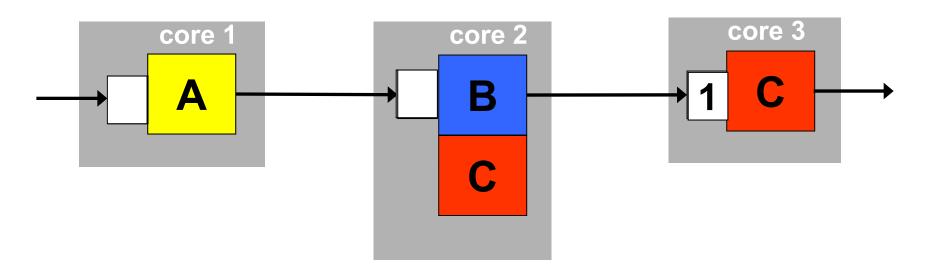




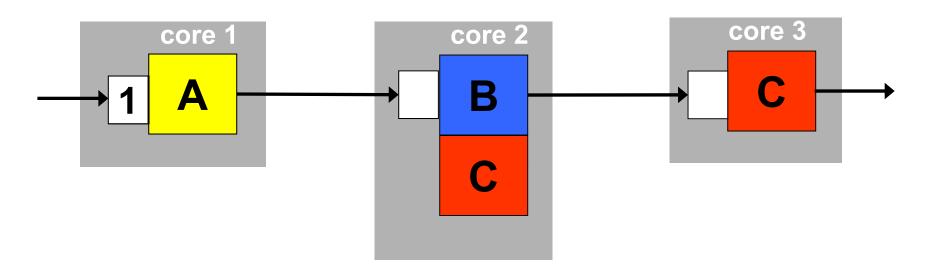




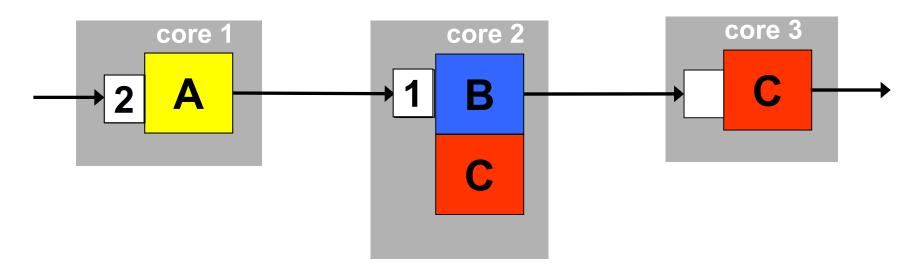




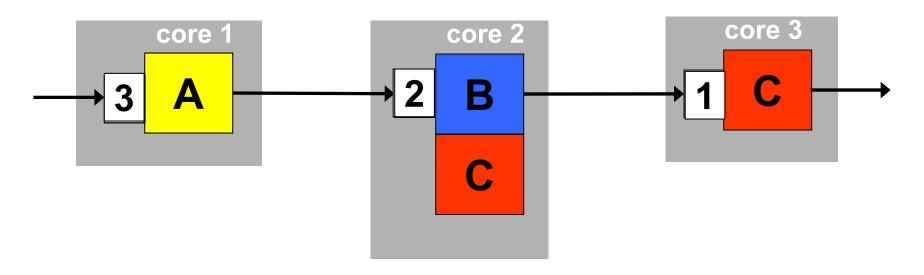
	t0	t1	t2	t3	t4	t5	t6	t7	t8	t9
core 1										
core 2										
core 3										



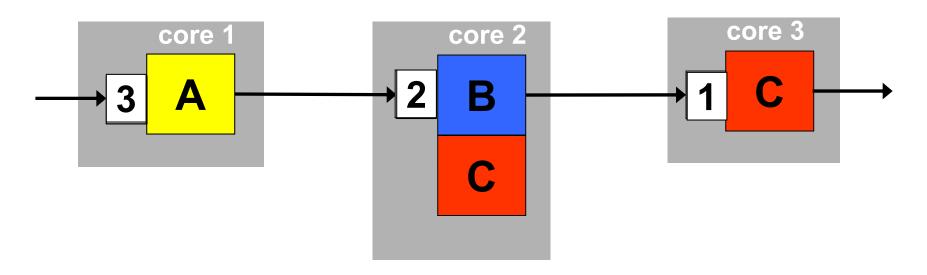
	t0	t1	t2	t3	t4	t5	t6	t7	t8	t9
core 1	A	1								
core 2										
core 3										



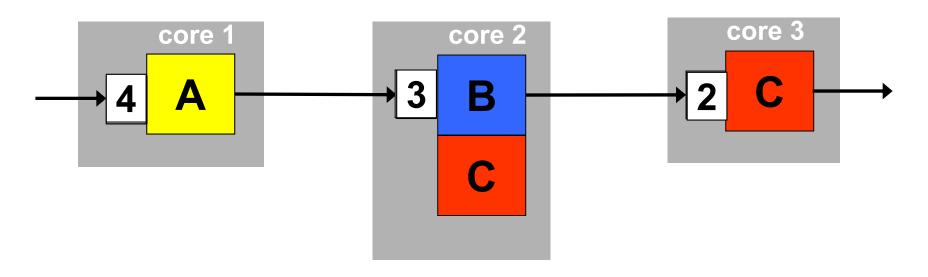
	t0	t1	t2	t3	t4	t5	t6	t7	t8	t9
core 1	A	1	A	2						
core 2			В	1						
core 3										



	t0	t1	t2	t3	t4	t5	t6	t7	t8	t9
core 1	Д	1	-	<b>\2</b>	A	\3				
core 2			•	31	E	<b>32</b>				
core 3						<b>C1</b>				



	t0	t1	t2	t3	t4	t5	t6	t7	t8	t9
core 1	Д	1	-	<b>\2</b>	A	\3				
core 2				31	E	<b>32</b>	<b>C2</b>			
core 3						<b>C1</b>				



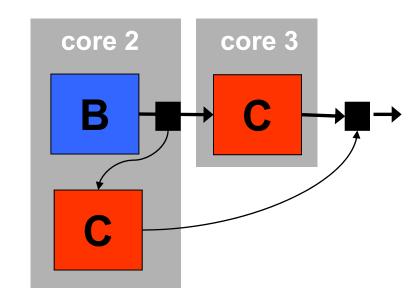
	t0	t1	t2	t3	t4	t5	t6	t7	t8	t9
core 1	A	1		42	F	\3		A	4	
core 2				<b>31</b>	E	32	<b>C2</b>	E	3	
core 3						<b>C1</b>			2	

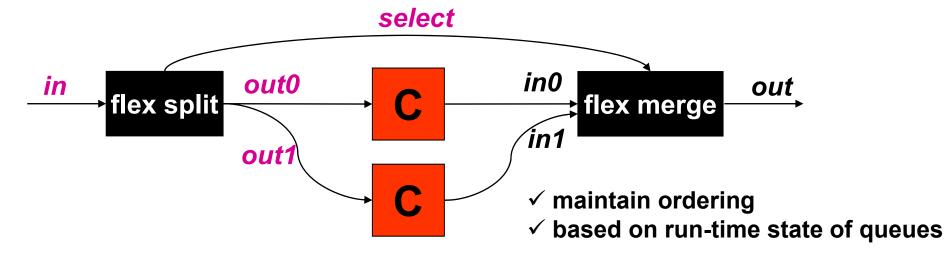
#### **Outline**

- Introduction
- Design Flow of a Stream Program with Flexibility
- Performance
- Implementation of Flexible Filters
- Experiments
  - CFAR Case Study

### Flex-Split

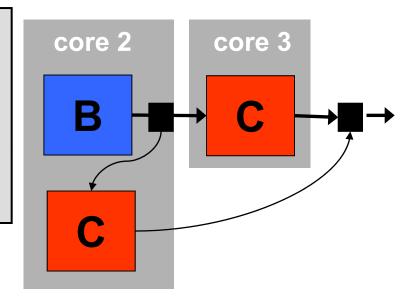
```
Flex-split
  pop data block b from in
  n0 = available space on out0
  n1 = |b| - n0
  send n0 to out0, n1 to out1
  send n0 0's, then n1 1's to
    select
```

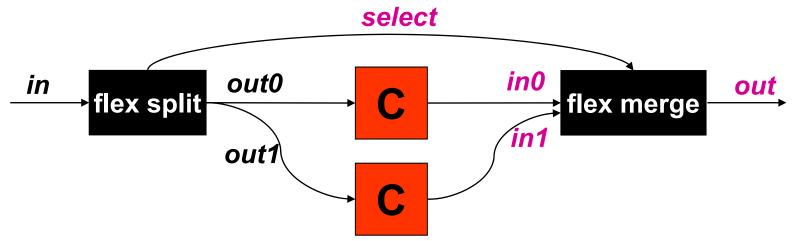




### Flex-Merge

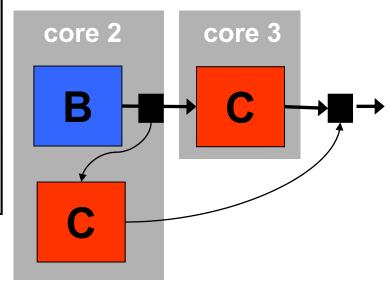
```
Flex-merge
  pop i from select
  if i is 0, pop token from in0
  if i is 1, pop token from in1
  push token to out
```



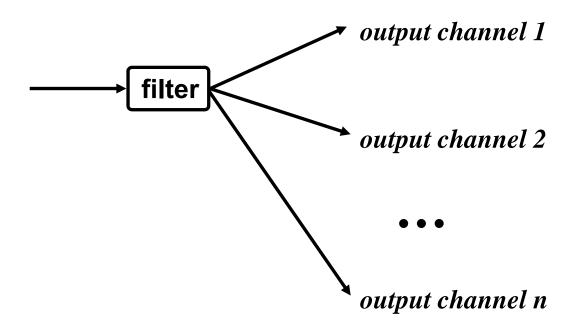


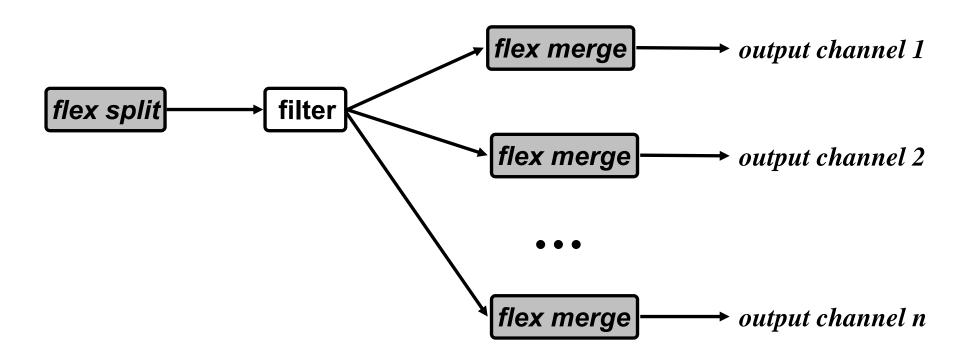
### Flex-Merge

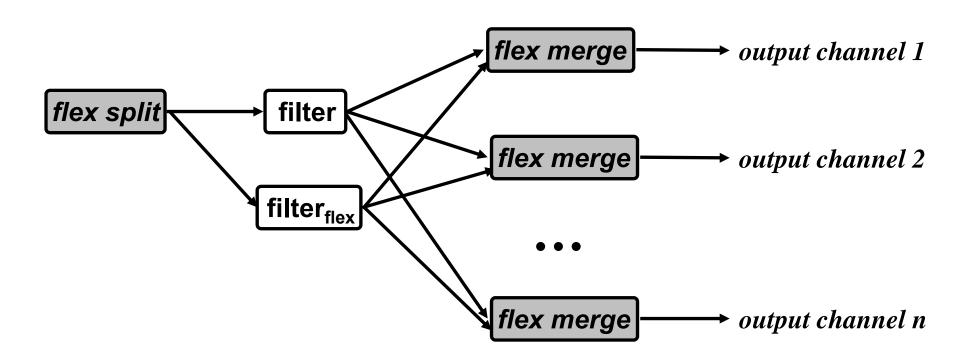
```
Flex-merge
  pop i from select
  if i is 0, pop token from in0
  if i is 1, pop token from in1
  push token to out
```

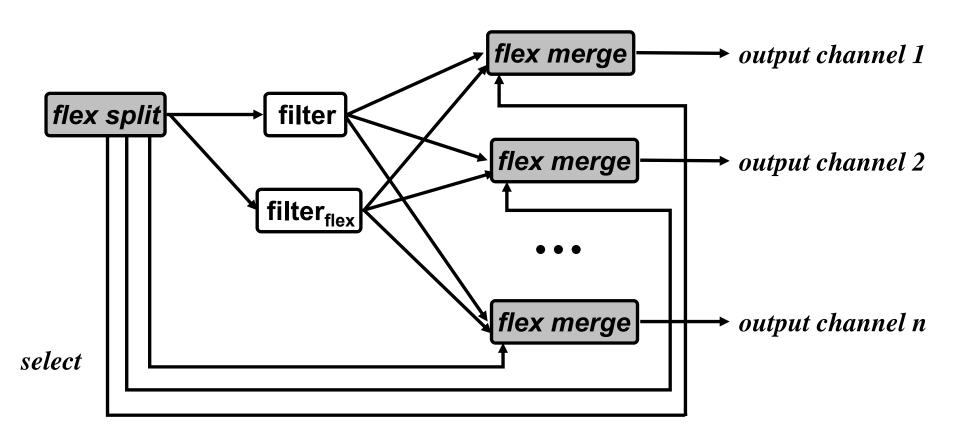


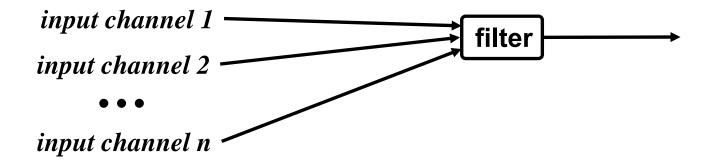
in Overhead of Flexibility?

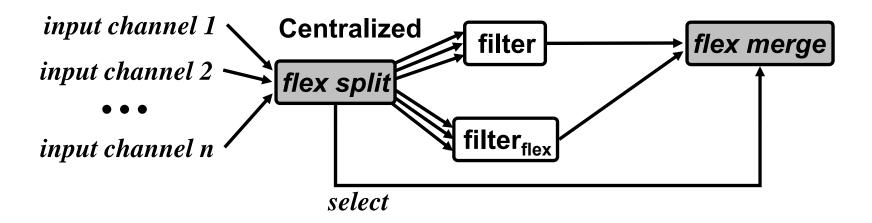


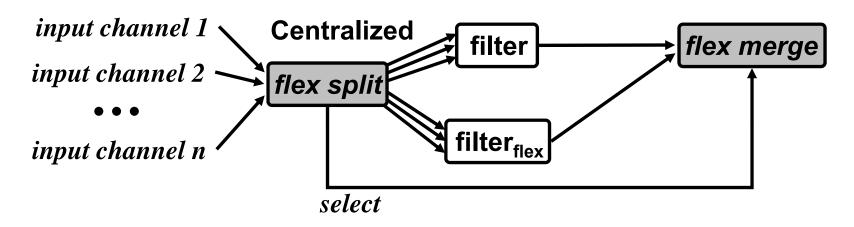


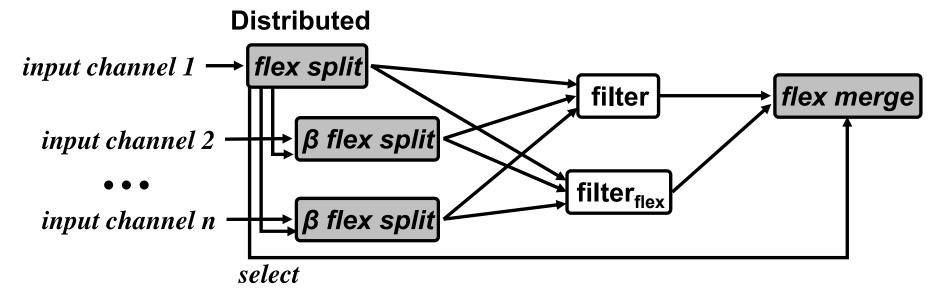










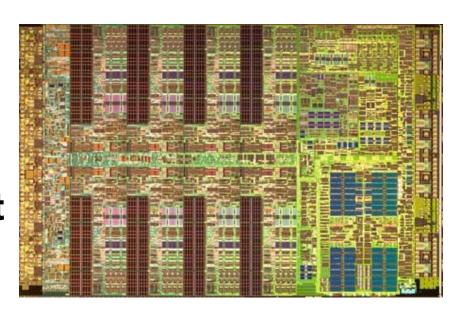


#### **Outline**

- Introduction
- Design Flow of a Stream Program with Flexibility
- Performance
- Implementation of Flexible Filters
- Experiments
  - CFAR Case Study

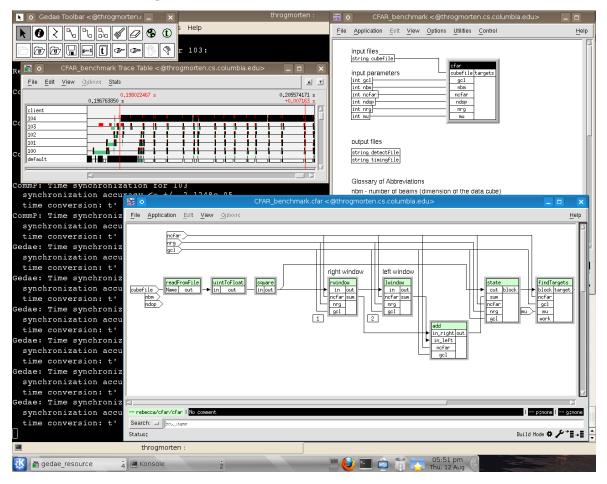
### **Cell BE Processor**

- Distributed Memory
- Heterogeneous
  - 8 SIMD (SPU) cores
  - 1 PowerPC (PPU)
- Element Interconnect
   Bus
  - 4 rings
  - 205 Gb/s
- Gedae Programming Language Communication Layer

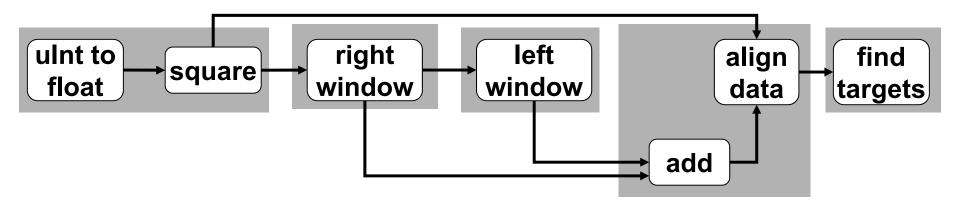


### Gedae

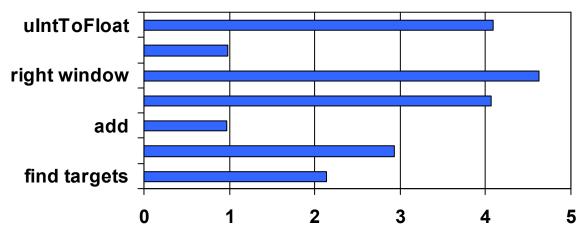
- Commercial data-flow language and programming GUI
- Performance analysis tools



#### **CFAR Benchmark**



#### **Profile of CFAR filters on Cell**



Average Execution Time (microseconds) per 100 tokens

### **Data Dependency**

- By changing threshold, change % targets
  - **1.3 %**
  - **7.3 %**
- Additional workload per target
  - 16 µs
  - 32 µs
  - $-64 \mu s$

#### % Targets

Additional Workload

	1.3	7.3
16 µs	0.82	1.45
32 µs	1.06	1.39
64 µs	1.27	1.47

### **More Benchmarks**

Benchmark	Field	Results	
Dedup Information Theory	Rabin block/ max chunk size	Speedup	
	Theory	4096/512	2.00
JPEG Image Processing	Image width x height		
	128x128	1.31	
	Processing	256x256	1.16
		512x512	1.25
Value-at- Risk	Finance	stocks/walks/timesteps	
		16/1024/1024	0.98
		64/1024/1024	1.56
		128/1024/1024	1.55

### Conclusions

#### Flexible filters

- adapt to data dependent bottlenecks
- distributed load balancing
- provide speedup without modification to original filters
- can be implemented on top of general stream languages