



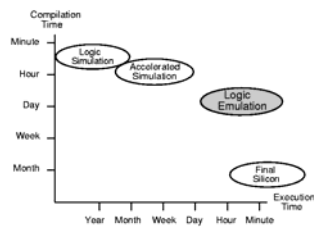
CprE / ComS 583 Reconfigurable Computing

Prof. Joseph Zambreno
Department of Electrical and Computer Engineering
Iowa State University

Lecture #12 – Systolic Computing



Recap – Logic Emulation



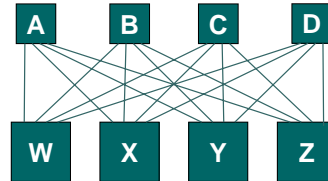
- Emulation takes a sizable amount of resources
- Compilation time can be large due to FPGA compiles

September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.3



Recap – Multi-FPGA Systems

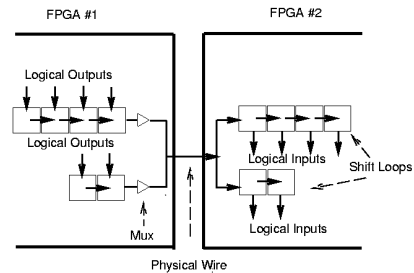
- Crossbar topology:
 - Devices A-D are routing only
 - Gives predictable performance
 - Potential waste of resources for near-neighbor connections



September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.2



Recap – Virtual Wires



- Overcome pin limitations by multiplexing pins and signals
- Schedule when communication will take place

September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.4



Outline

- Recap
- Introduction and Motivation
- Common Systolic Structures
- Algorithmic Mapping
- Mapping Examples
 - Finite impulse response
 - Matrix-vector product
 - Banded matrix-vector product
 - Banded matrix multiplication

September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.5



Systolic Computing

sys-to-le (sīs'tə-lē) *n.* – the rhythmic contraction of the heart, especially of the ventricles, by which blood is driven through the aorta and pulmonary artery after each dilation or diastole

[Greek *systolē*, from *systellein* to contract, from *syn-* + *stellein* to send]

– **sys-tol-ic** (sīs-tōl'ik) *adj.*

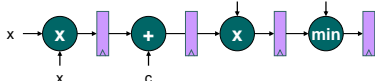
Data flows from memory in a rhythmic fashion, passing through many processing elements before it returns to memory.

[Kung, 1982]

September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.6

••• Systolic Architectures

- Goal – general methodology for mapping computations into hardware (spatial computing) structures
- Composition:
 - Simple compute cells (e.g. add, sub, max, min)
 - Regular interconnect pattern
 - Pipelined communication between cells
 - I/O at boundaries



September 27, 2007

CprE 583 – Reconfigurable Computing

Lect-12.7

••• Motivation

- Effectively utilize VLSI
- Reduce “Von Neumann Bottleneck”
- Target compute-intensive applications
- Reduce design cost
 - Simplicity
 - Regularity
- Exploit concurrency
- Local communication
 - Short wires (small delay, less area)
 - Scalable

September 27, 2007

CprE 583 – Reconfigurable Computing

Lect-12.8

••• Why Study?

- Original motivation – specialized accelerator for an application
- Model/goals is a close match to reconfigurable computing
- Target algorithms match
- Well-developed theory, techniques, and solutions
- One big difference – Kung’s approach targeted custom silicon (not a reconfigurable fabric)
 - Compute elements needed to be more general

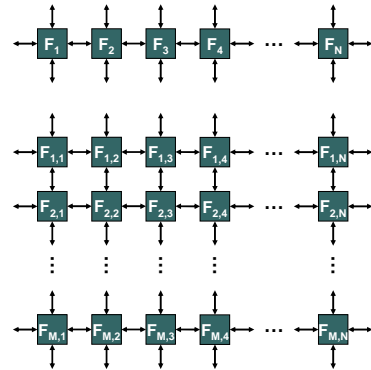
September 27, 2007

CprE 583 – Reconfigurable Computing

Lect-12.9

••• Common Systolic Structures

- One-dimensional linear array
- Two-dimensional mesh

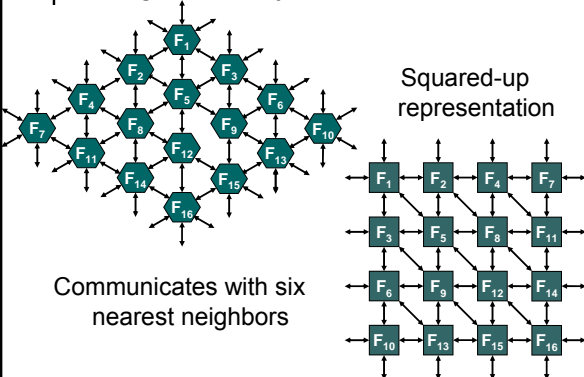


September 27, 2007

CprE 583 – Reconfigurable Computing

Lect-12.10

••• Hexagonal Array

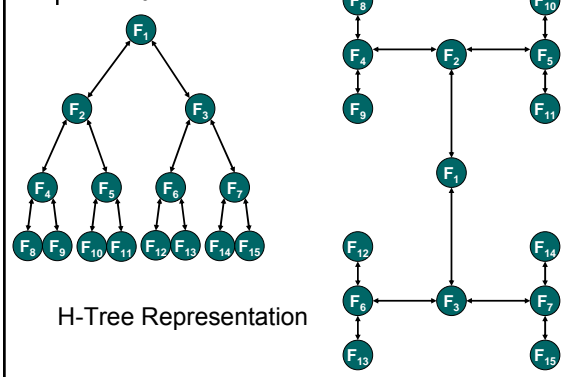


September 27, 2007

CprE 583 – Reconfigurable Computing

Lect-12.11

••• Binary Tree



September 27, 2007

CprE 583 – Reconfigurable Computing

Lect-12.12

Mapping Approach

- Allocate PEs
- Schedule computation
 - Schedule PEs
 - Schedule data flow
- Optimize
- Available Transformations:
 - Preload repeated values
 - Replace feedback loops with registers
 - Internalize data flow
 - Broadcast common input

September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.13

Example – Finite Impulse Response

- A Finite Impulse Response (FIR) filter is a type of digital filter
 - Finite – response to an impulse eventually settles to zero
 - Requires no feedback

$$y_i = w_1 \cdot x_i + w_2 \cdot x_{i+1} + \dots + w_k \cdot x_{i+k-1}$$

$$= \sum_{j=1}^k w_j \cdot x_{i+j-1}$$

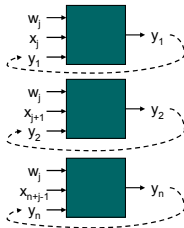
```
for (i=1; i<=n; i++)
  for (j=1; j <=k; j++)
    y[i] += w[j] * x[i+j-1];
```

September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.14

FIR Attempt #1

- Parallelize the outer loop


```
for (i=1; i<=n; i++)      in parallel
  for (j=1; j <=k; j++)  sequential
    y[i] += w[j] * x[i+j-1];
```

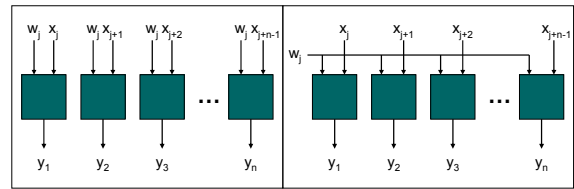


September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.15

FIR Attempt #1 (cont.)

- Broadcast common inputs


```
for (i=1; i<=n; i++)      in parallel
  for (j=1; j <=k; j++)  sequential
    y[i] += w[j] * x[i+j-1];
```

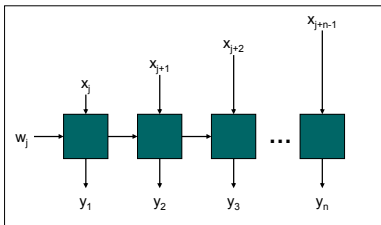


September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.16

FIR Attempt #1 (cont.)

- Retime to eliminate broadcast


```
for (i=1; i<=n; i++)      in parallel
  for (j=1; j <=k; j++)  sequential
    y[i] += w[j] * x[i+j-1];
```

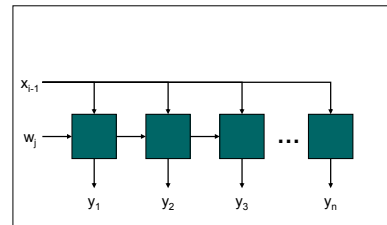


September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.17

FIR Attempt #1 (cont.)

- Broadcast common values


```
for (i=1; i<=n; i++)      in parallel
  for (j=1; j <=k; j++)  sequential
    y[i] += w[j] * x[i+j-1];
```

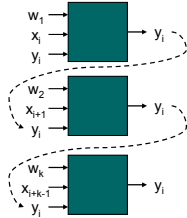


September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.18

FIR Attempt #2

- Parallelize the inner loop

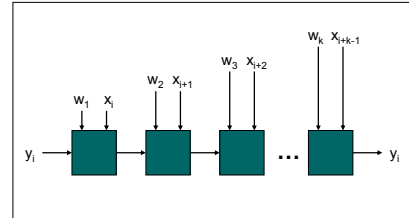

```
for (i=1; i<=n; i++) sequential
  for (j=1; j <=k; j++) in parallel
    y[i] += w[j] * x[i+j-1];
```



FIR Attempt #2 (cont.)

- Internalize data flow

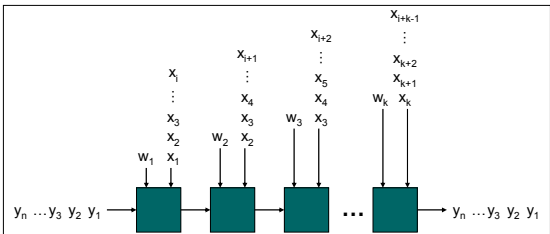

```
for (i=1; i<=n; i++) sequential
  for (j=1; j <=k; j++) in parallel
    y[i] += w[j] * x[i+j-1];
```



FIR Attempt #2 (cont.)

- Allocation schedule

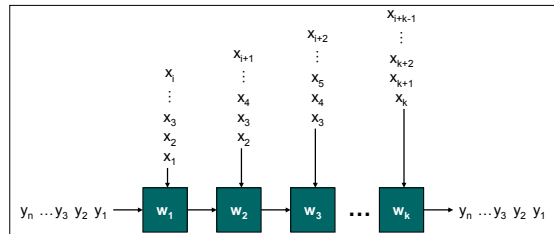

```
for (i=1; i<=n; i++) sequential
  for (j=1; j <=k; j++) in parallel
    y[i] += w[j] * x[i+j-1];
```



FIR Attempt #2 (cont.)

- Preload repeated values

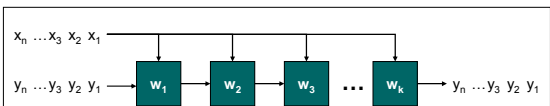

```
for (i=1; i<=n; i++) sequential
  for (j=1; j <=k; j++) in parallel
    y[i] += w[j] * x[i+j-1];
```



FIR Attempt #2 (cont.)

- Broadcast common values

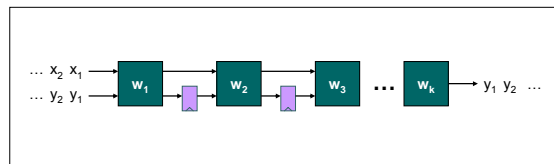

```
for (i=1; i<=n; i++) sequential
  for (j=1; j <=k; j++) in parallel
    y[i] += w[j] * x[i+j-1];
```

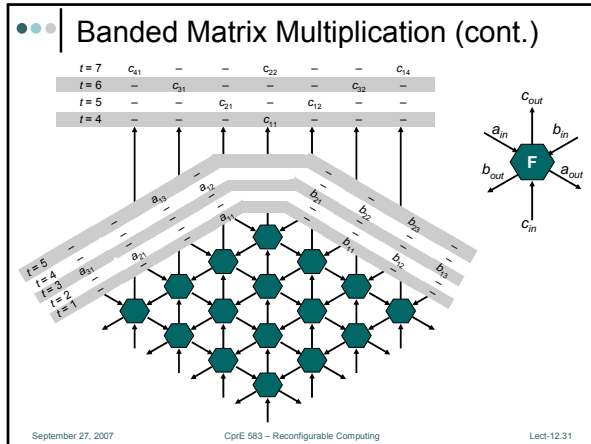


FIR Attempt #2 (cont.)

- Retime to eliminate broadcast


```
for (i=1; i<=n; i++) sequential
  for (j=1; j <=k; j++) in parallel
    y[i] += w[j] * x[i+j-1];
```





- ### Summary
- Systolic structures are good for computation-bound problems
 - Models costs in VLSI systems
 - Minimize number of memory accesses
 - Emphasize local interconnections (long wires are bad)
 - Candidate algorithms
 - Makes multiple use of input data (ex: n inputs, $O(n^3)$ computations)
 - Concurrency
 - Simple control flow, simple processing elements
- September 27, 2007 CprE 583 – Reconfigurable Computing Lect-12.32