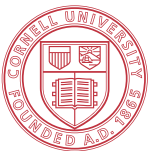




ECE 6775  
High-Level Digital Design Automation  
Fall 2025

**More Pipelining  
Resource Sharing**



Cornell University



# Announcements

- ▶ Lab 3 due tomorrow
- ▶ HW 2 will be posted today, due Friday Oct 17
  - Lower penalty (1%) for late submissions, ***up to 4 days late*** (until Tue Oct 21); Solution will be released on Wed Oct 22nd
- ▶ 2<sup>nd</sup> paper reading session on Tuesday Oct 21
  - C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, [“Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks”](#), FPGA 2015.

# Agenda

- ▶ Modulo scheduling case studies
- ▶ Systolic arrays: combining parallel processing and pipelining
  - Uniform recurrence equations
  - Case study on matrix multiplication
- ▶ Resource sharing basics
  - Sub-problems: functional unit, register, and connectivity binding problems
  - Key concepts: compatibility and conflict graphs

# Recap: Calculating Lower Bound of II

- ▶ Minimum possible II (MII)
  - $MII = \max(\text{ResMII}, \text{RecMII})$
  - A lower bound, not necessarily achievable
- ▶ Resource constrained MII (ResMII)
  - $\text{ResMII} = \max_i \lceil \text{OPs}(r_i) / \text{Limit}(r_i) \rceil$   
OPs(r): number of operations that use resource of type r  
Limit(r): number of available resources of type r
- ▶ Recurrence constrained MII (RecMII)
  - $\text{RecMII} = \max_i \lceil \text{Latency}(c_i) / \text{Distance}(c_i) \rceil$   
Latency( $c_i$ ): total latency in dependence cycle  $c_i$   
Distance( $c_i$ ): total distance in dependence cycle  $c_i$

# Recap: Minimum II due to Resource Limits

- ▶ Compute ResMII: Max among all types of resources

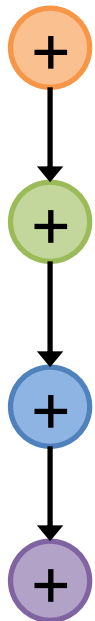
$$\text{ResMII} = \max_i \lceil \text{OPs}(r_i) / \text{Limit}(r_i) \rceil$$

OPs(r): # of operations that use resource r  
 Limit(r): # of available resources of type r

Take the max ratio among all resource types

Dependence

Resource Allocation & Binding



4 adders  
(a0~a3)

	time (in cycles)					
	0	1	2	3	4	5
a0	i0	i1	i2	i3	i4	i5
a1		i0	i1	i2	i3	i4
a2			i0	i1	i2	i3
a3				i0	i1	i2

0, 1, 2, 3, 4, 5 : time (cycles)  
 a0, a1, a2, a3 : available adders  
 i0, i1, i2, ... : iterations

2 adders  
(a0,a1)

	time (in cycles)							
a0	i0	i0	i1	i1	i2	i2	i3	i3
a1			i0	i0	i1	i1	i2	i2

due to limited resources,  
 cannot initiate iterations  
 less than 2 cycles apart

# Case Study: Prefix Sum

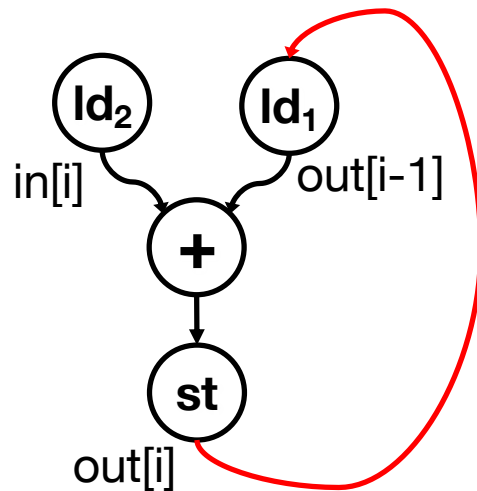
- ▶ Prefix sum computes a cumulative sum of a sequence of numbers
  - commonly used in many applications such as radix sort, histogram, etc.

```
void prefixsum ( int in[N], int out[N] )  
  out[0] = in[0];  
  for ( int i = 1; i < N; i++ ) {  
    #pragma HLS pipeline II=?  
    out[i] = out[i-1] + in[i];  
  }  
}
```

```
out[0] = in[0];  
out[1] = in[0] + in[1];  
out[2] = in[0] + in[1] + in[2];  
out[3] = in[0] + in[1] + in[2] + in[3];  
  
...
```

# Prefix Sum: RecMII

- ▶ Loop-carried dependence exists between reads on 'out'
  - Assume chaining is not possible on memory reads (ld) and writes (st) due to target cycle time
  - RecMII = 3



```

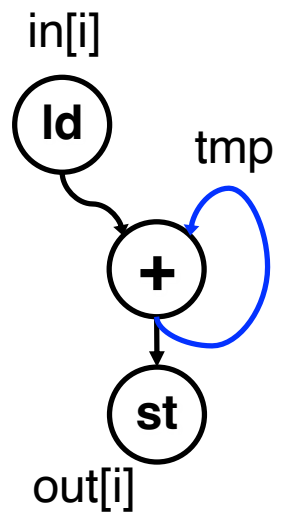
out[0] = in[0];
for ( int i = 1; i < N; i++ )
    out[i] = out[i-1] + in[i];
    
```

	cycle 1	cycle 2	cycle 3	cycle 4
$i = 0$	ld <sub>1</sub> ld <sub>2</sub>	+	st	
$i = 1$	<del>// = 1</del>	ld <sub>1</sub> ld <sub>2</sub>	+	st

ld – Load  
st – Store

# Prefix Sum: Code Optimization

- ▶ Introduce an intermediate variable 'tmp' to hold the running sum from the previous 'in' values
  - Shorter dependence cycle leads to RecMII = 1



**ld** – Load  
**st** – Store

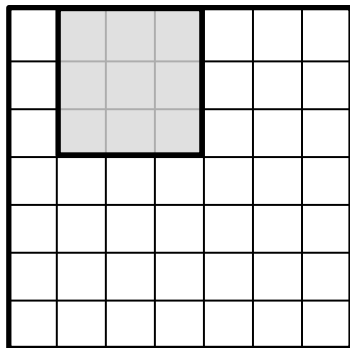
```

int tmp = in[0];
for ( int i = 1; i < N; i++ ) {
    tmp += in[i];
    out[i] = tmp;
}
    
```

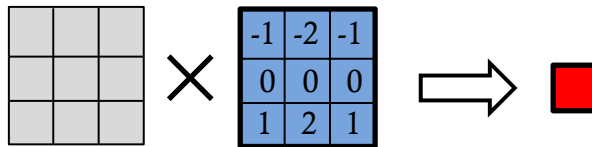
	cycle 1	cycle 2	cycle 3	cycle 4
$i = 0$	ld	+	st	
$i = 1$	<b>// = 1</b>	ld	+	st

# Case Study: 2D Convolution

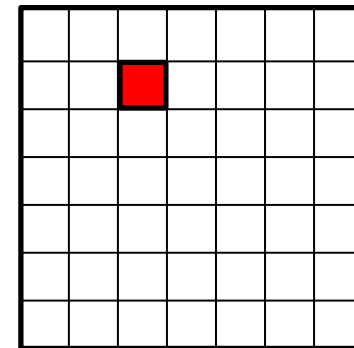
```
for (r = 1; r < R; r++) {  
  for (c = 1; c < C; c++) {  
    #pragma HLS pipeline II=?  
    for (i = 0; i < 3; i++)  
      for (j = 0; j < 3; j++)  
        sum += img[r+i-1][c+j-1] * f[i][j];  
    out[r][c] = sum;  
  }  
}
```



Input image  
frame



A K by K **dot product** is performed  
for each output pixel (K=3 here)



Output image  
frame

## 3x3 Convolution: MII

```
for (r = 1; r < R; r++) {  
  for (c = 1; c < C; c++) {  
    #pragma HLS pipeline II=?  
    for (i = 0; i < 3; i++)  
      for (j = 0; j < 3; j++)  
        sum += img[r+i-1][c+j-1] * f[i][j];  
    out[r][c] = sum;  
  }  
}
```

- ▶ Inner loops, “i” & “j”, are automatically unrolled
- ▶ The 3x3 filter array, “f”, is partitioned into 9 registers
- ▶ The entire input image, “img”, is stored in an on-chip SRAM with **two read ports**

ResMII = ?

What about RecMII?

# Revisiting Systolic Arrays

- ▶ An array of processing elements (PEs) that process data in a systolic manner using nearest-neighbor communication
  - Systolic means “data flows from memory in a rhythmic fashion, passing through many processing elements before it returns to memory” – H.T. Kung

**Systolic Arrays (for VLSI)**

**H. T. Kung† and Charles E. Leiserson†**

*And now I see with eye serene  
The very pulse of the machine.  
--William Wordsworth*

**Abstract**

A **systolic** system is a network of processors which rhythmically compute and pass data through the system. Physiologists use the word “systole” to refer to the rhythmically recurrent contraction of the heart and arteries which pulses blood through the body. In a **systolic** computing system, the function of a processor is analogous to that of the heart. Every processor regularly pumps data in and out, each time performing some short computation, so that a regular flow of data is kept up in the network.

Many basic matrix computations can be pipelined elegantly and efficiently on **systolic** networks having an array structure. As an example, hexagonally connected processors can optimally perform matrix multiplication. Surprisingly, a similar **systolic** array can compute the LU-decomposition of a matrix. These **systolic arrays** enjoy simple and regular communication paths, and almost all processors used in the networks are identical. As a result, special purpose hardware devices based on **systolic arrays** can be built inexpensively using the **VLSI** technology.

**1. Introduction**

Developments in microelectronics have revolutionized computer design. Integrated circuit technology has increased the number and complexity of components that can fit on a chip or a printed circuit board. Component density has been doubling every one-to-two years and already, a multiplier can fit on a very large scale integrated

In Sparse Matrix Proceedings, 1978



parallel processing + pipelining

- + Simple & regular design
- + Massive parallelism
- + Short interconnection
- + Balancing compute with I/O

# Uniform Recurrence Equations (UREs)

- ▶ Any **systolic algorithm** can be described by a set of UREs
  - i.e., an n-dimensional loop nest where the recurrences (inter-iteration dependences) must have constant distances

---

**c = a \* b**

```
c = 0;
for (int i = 0; i < N; i++)
  c += a[i] * b[i];
```

**Dot product in UREs**

$$Z[0] = 0$$
$$Z[i] = Z[i - 1] + a[i] \cdot b[i], \quad 0 < i < N$$
$$c = Z[N - 1]$$

---

**c = A \* b**

```
for (int i = 0; i < N; i++)
  c[i] = 0;
for (int j = 0; j < N; j++)
  c[i] += A[i, j] * b[j]
```

**Matrix Vector Multiplication (MV) in UREs**

$$Z[i, 0] = 0$$
$$Z[i, j] = Z[i, j - 1] + A[i, j] \cdot x[j], \quad 0 < j < N$$
$$c[i] = Z[i, N - 1]$$

---

**C = A \* B**

```
for (int i = 0; i < N; i++)
  for (int j = 0; j < N; j++)
    C[i, j] = 0;
for (int k = 0; k < N; k++)
  C[i, j] += A[i, k] * B[k, j]
```

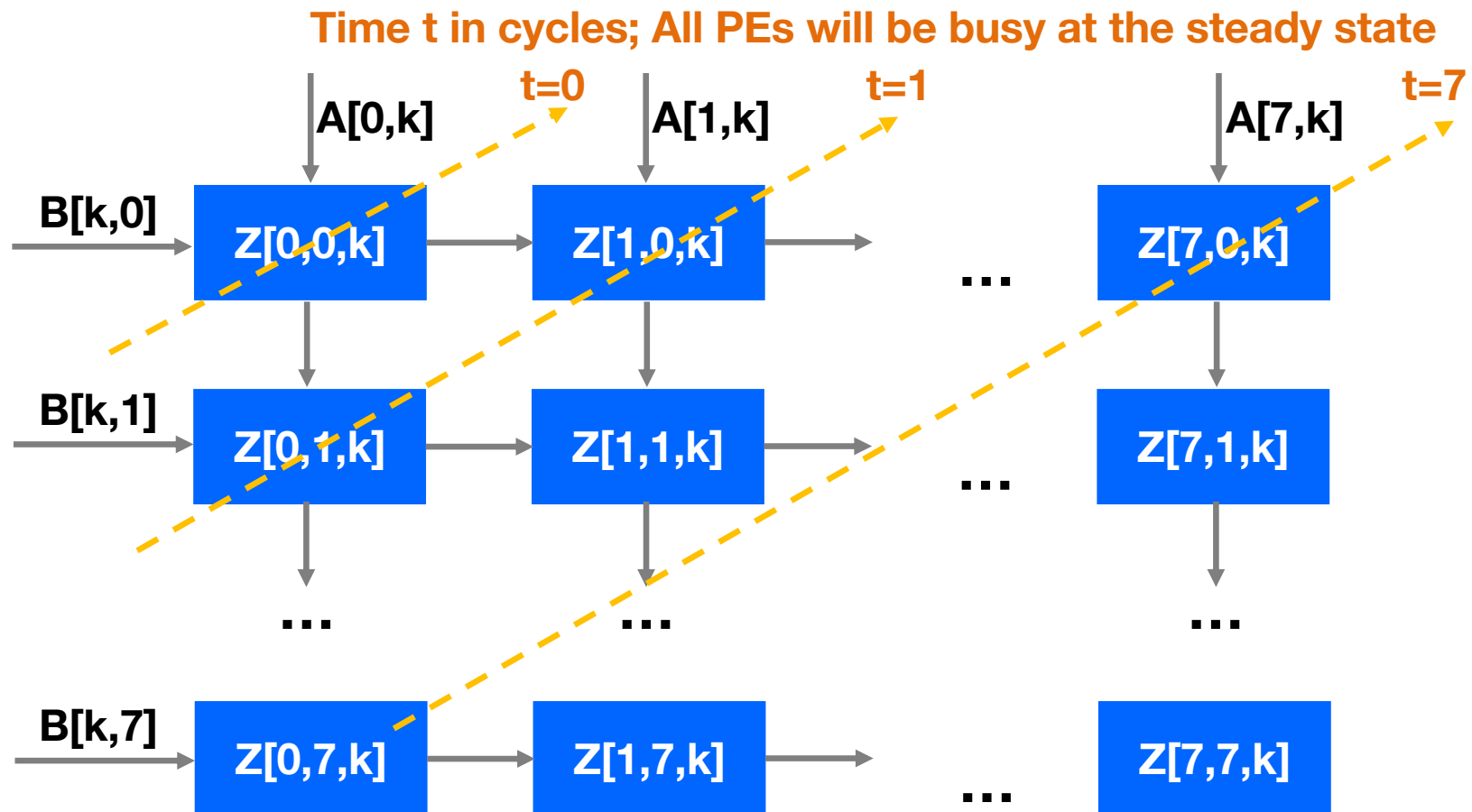
**Matrix Matrix Multiplication (MM) in UREs**

$$Z[i, j, 0] = 0$$
$$Z[i, j, k] = Z[i, j, k - 1] + A[i, k] \cdot B[k, j], \quad 0 < k < N$$
$$C[i, j] = Z[i, j, N - 1]$$

# Mapping MM to Systolic Array

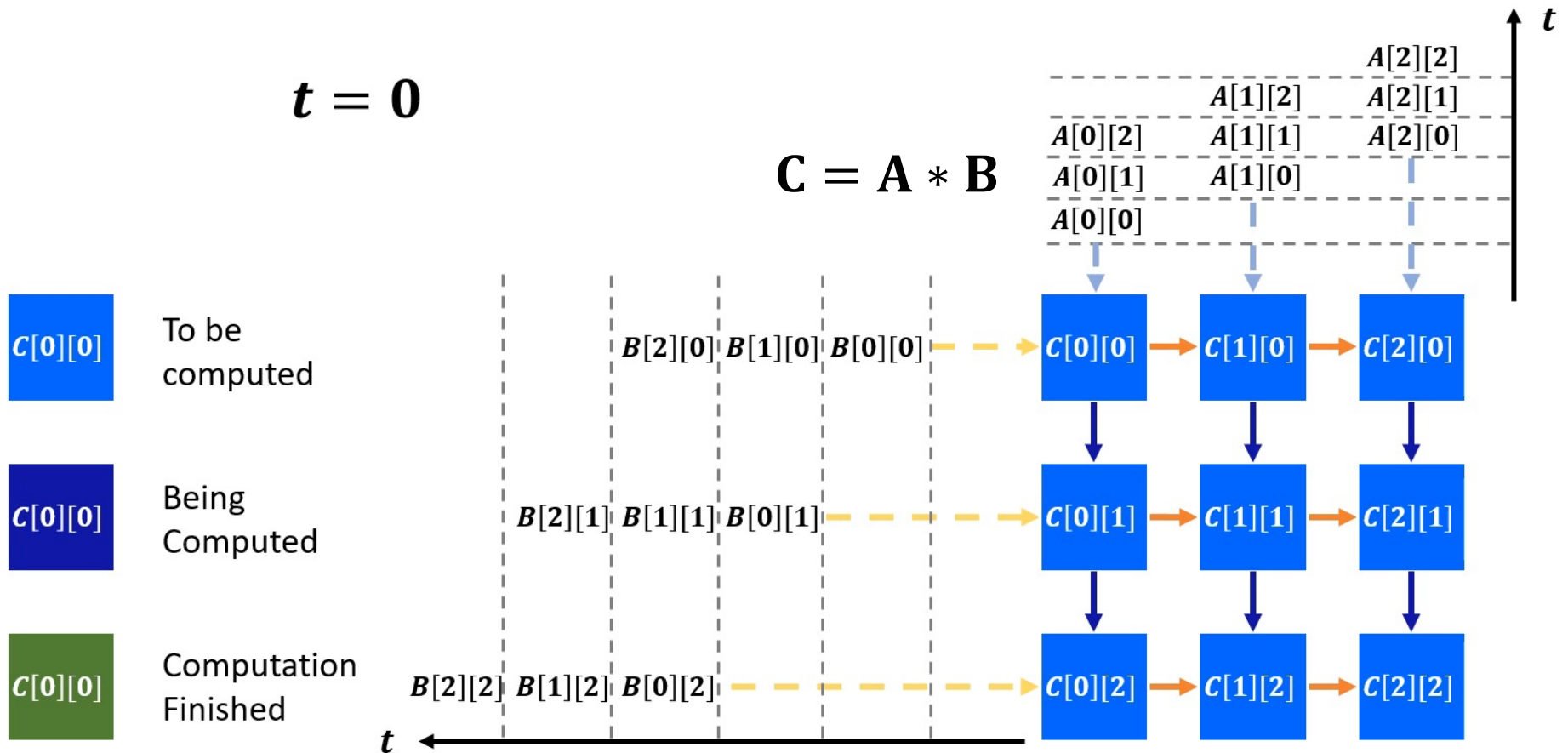
$$\begin{aligned}
 \mathbf{C} &= \mathbf{A} * \mathbf{B} \\
 Z[i, j, 0] &= 0 \\
 Z[i, j, k] &= Z[i, j, k - 1] + A[i, k] \cdot B[k, j], \quad 0 < k < N \\
 C[i, j] &= Z[i, j, N - 1]
 \end{aligned}$$

Map the n-dimensional iteration space into a physical array of PEs

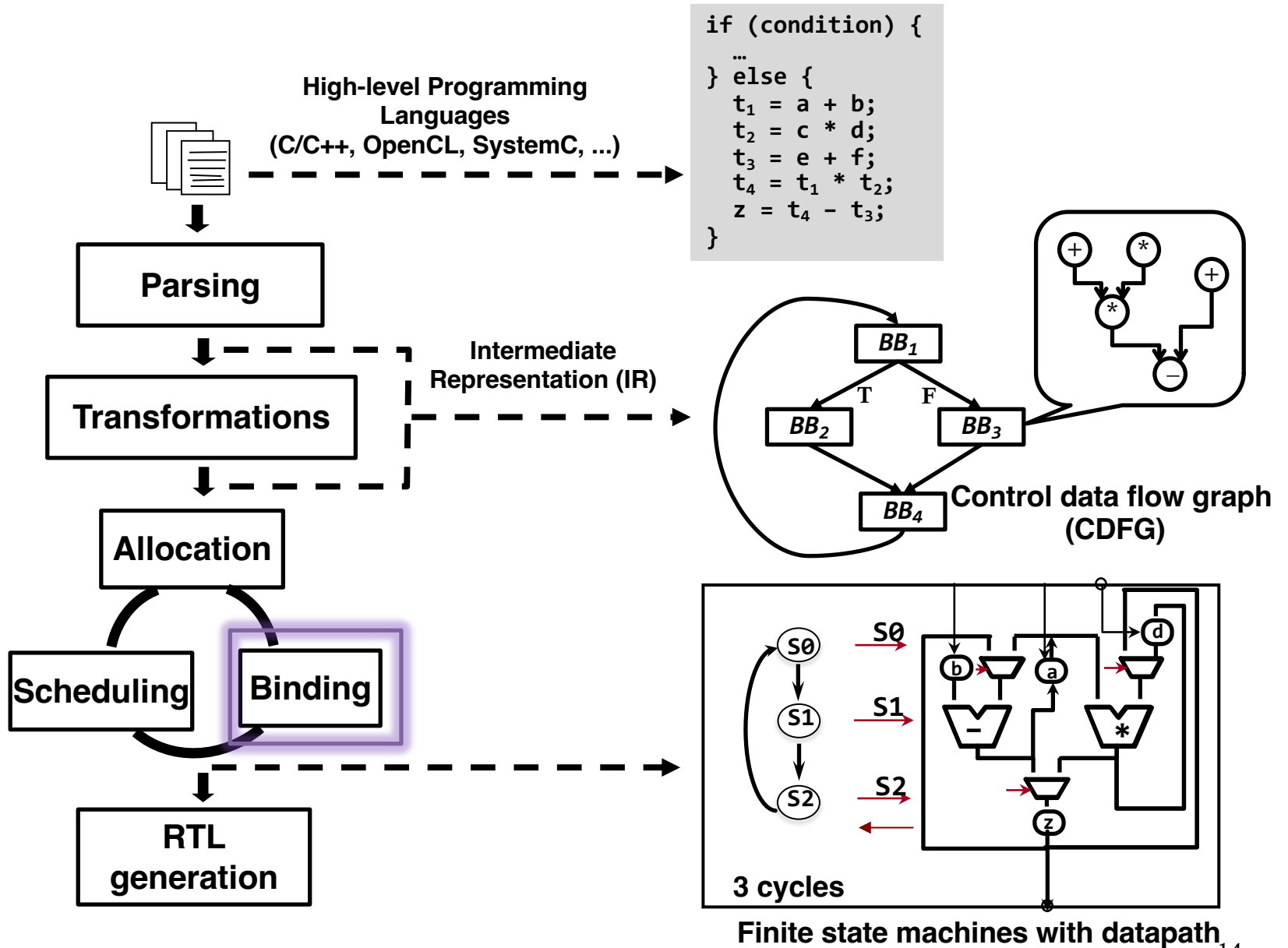


# Running MM on Systolic Array: Another Look

- ▶ An array of processing elements that process data in a systolic manner



# Recap: A Typical HLS Flow



# Resource Sharing and Binding

- ▶ **Resource sharing** enables reuse of hardware resources to minimize cost, in resource usage/area/power
  - Typically carried out by binding in HLS
  - Other subtasks such allocation and scheduling greatly impact the resource sharing opportunities
- ▶ **Binding** maps operations, variables, and/or data transfers to the available resources
  - After scheduling: decide resource usage and detailed architecture (**focus of this lecture**)
  - Before scheduling: affect both area and delay
  - Simultaneous scheduling and binding: better result but more expensive

# Binding Sub-problems

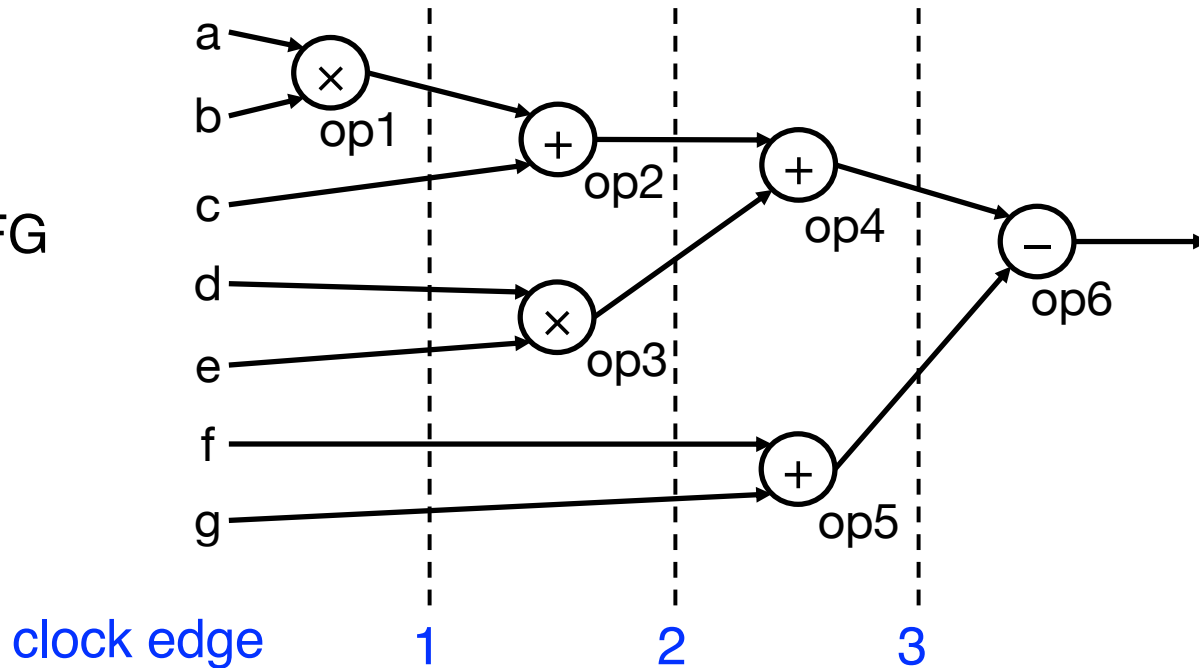
- ▶ Functional unit (FU) binding
  - Primary objective is to minimize the number of FUs
  - Considers connection cost
- ▶ Register binding
  - Primary objective is to minimize the number of registers
  - Considers connection cost
- ▶ Connectivity binding
  - Minimize connections by exploiting the commutative property of some operations / FUs
  - NP-hard

# Sharing Conditions

- ▶ Functional units (registers) are shared by operations (variables) of same type whose *lifetimes* do not overlap
- ▶ **Lifetime:** [birth-time, death-time)
  - Operation: The whole execution time (if unpipelined)
  - Variable: From the time this variable is defined to the time it is last used
- ▶ **Here we assume no pipelining to simplify discussion**
  - With pipelining (modulo scheduling), we use slots to determine overlaps rather than control steps

# Functional Unit Binding

a scheduled DFG  
(unpipelined)



Functional Unit	Operations
Mul1	op1, op3
AddSub1	op2, op4
AddSub2	op5, <u>op6</u>

Binding 1

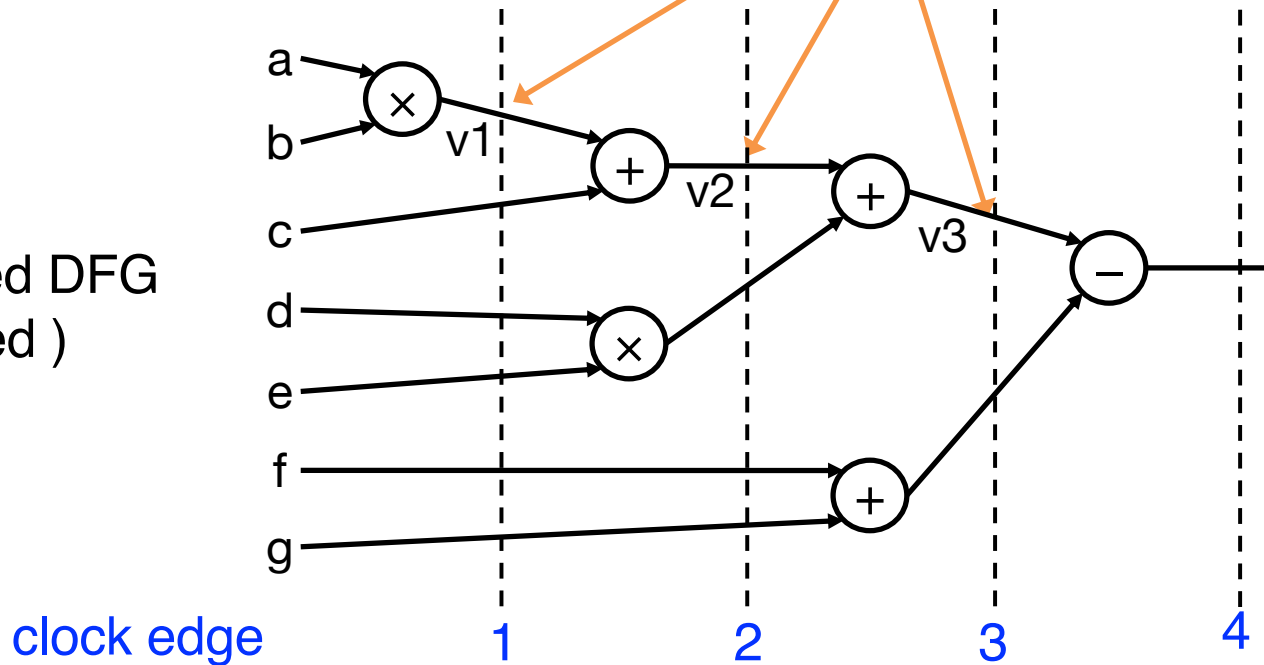
Functional Unit	Operations
Mul1	op1, op3
AddSub1	op2, op4, <u>op6</u>
AddSub2	op5

Binding 2

# Register Binding

Lifetimes crossing at least one clock edge  
→ register(s) required

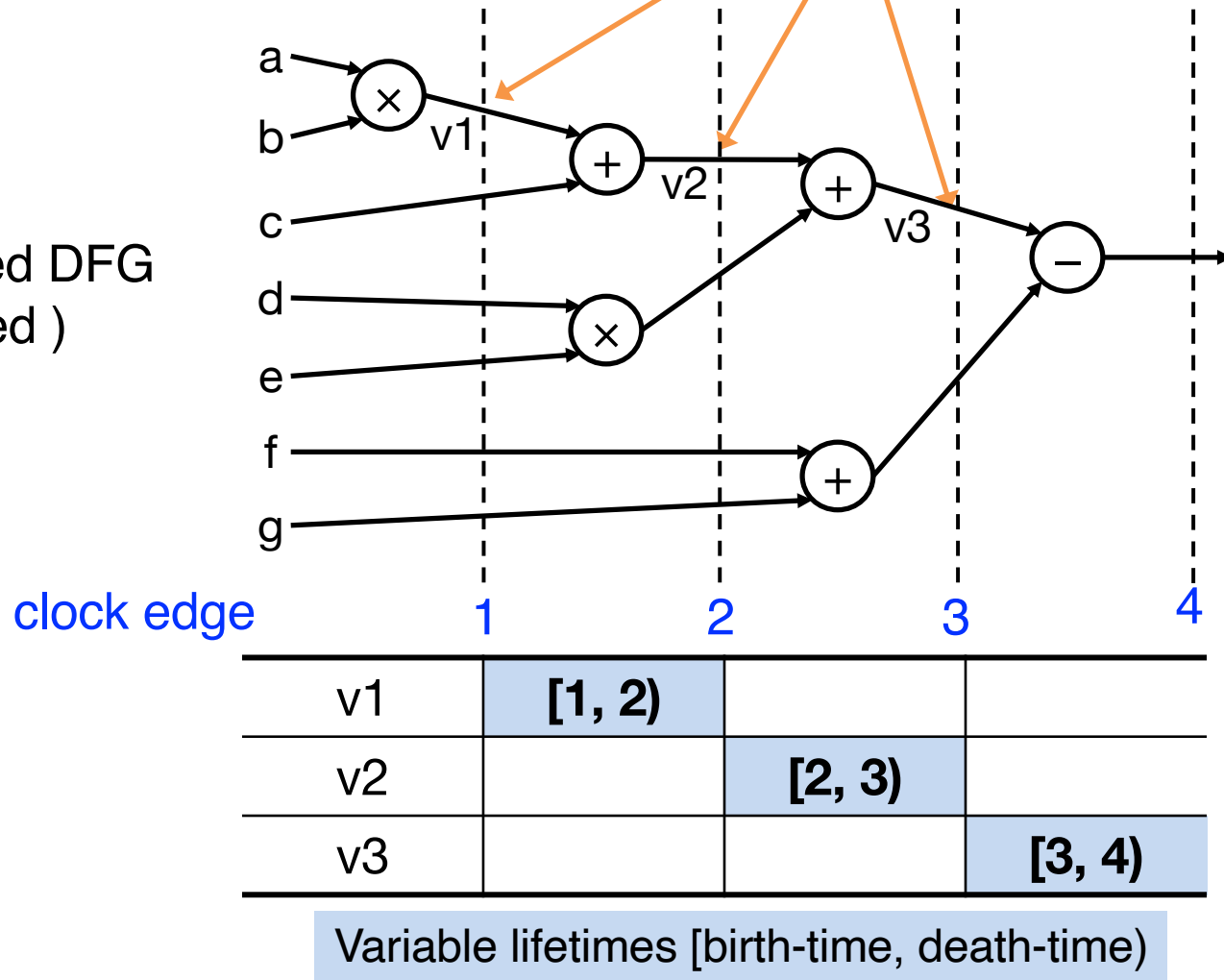
a scheduled DFG  
(unpipelined)



# Variable Lifetime Analysis

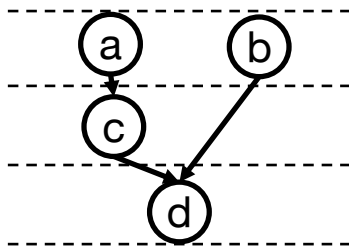
Variables v1, v2, and v3 can share the same register

a scheduled DFG  
(unpipelined)

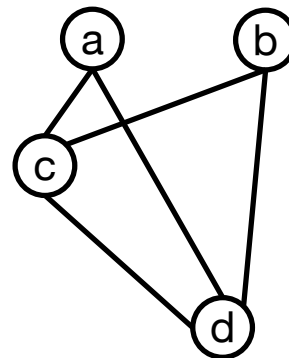


# Compatibility and Conflict Graphs

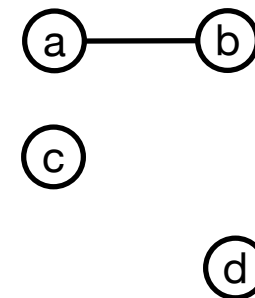
- ▶ Operation/variables **compatibility**
  - **Same type, non-overlapping lifetimes**
- ▶ Compatibility graph
  - Vertices: operations/variables
  - Edges: compatibility relation
- ▶ Conflict graph: Complement of compatibility graph



A scheduled DFG  
(unpipelined; operations  
have the same type)



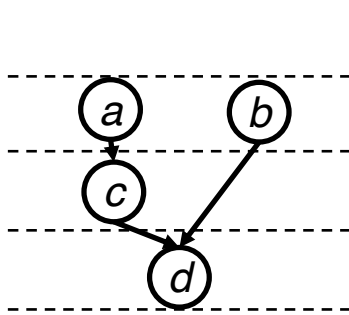
Compatibility graph



Conflict graph

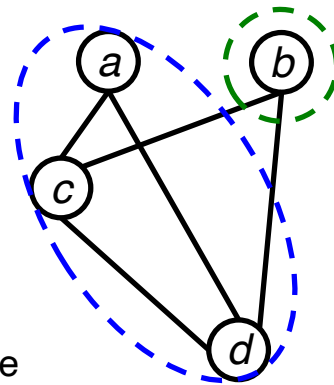
# Clique Cover Number and Chromatic Number

- ▶ Compatibility graph
  - Partition the graph into a **minimum number of cliques**
    - Clique in an undirected graph is a subset of its vertices such that every two vertices in the subset are connected by an edge
- ▶ Conflict graph
  - Color the vertices by a **minimum number of colors** (chromatic number), where adjacent vertices cannot use the same color

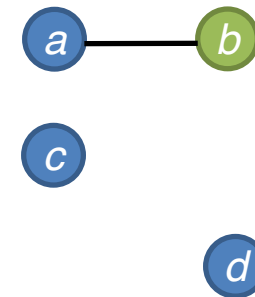


Operations have same type

A scheduled DFG



**Clique partitioning** on compatibility graph



**Coloring** on conflict graph

# Left Edge Algorithm

## ► Problem statement

- Given: Input is a group of intervals with starting and ending time
- Goal: Minimize the number of colors of the corresponding interval graph

### **Repeat**

create a new color group  $c$

### **Repeat**

assign leftmost feasible interval to  $c$

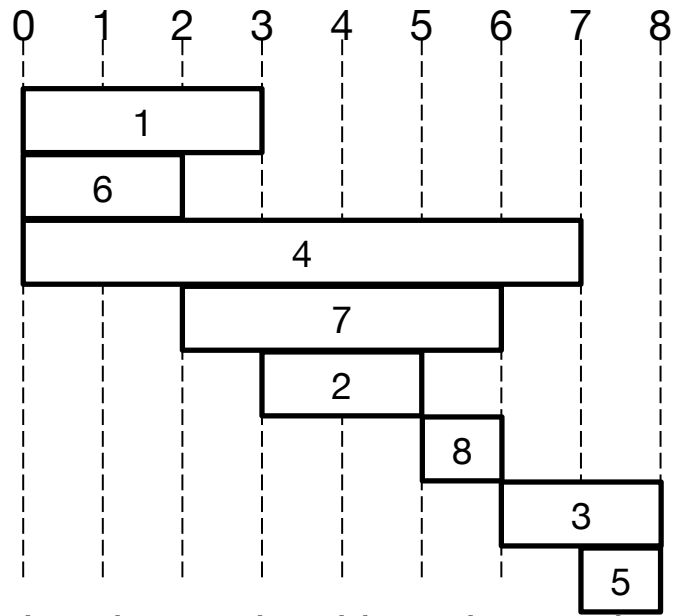
**until** no more feasible interval

**until** no more interval

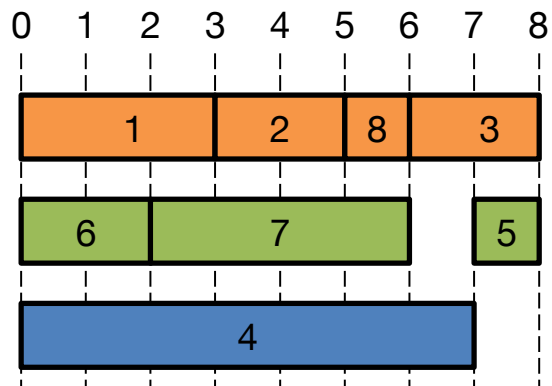
Interval are **sorted** according to their **left endpoints**

**Greedy algorithm,  $O(n \log n)$  time**

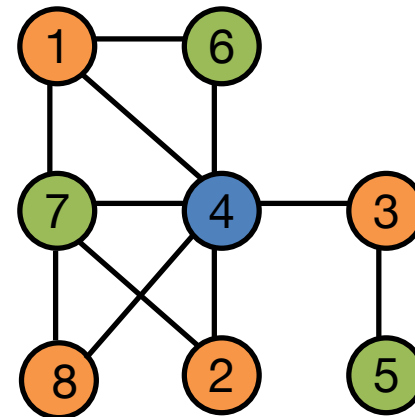
# Left Edge Demonstration



Lifetime intervals with a given schedule



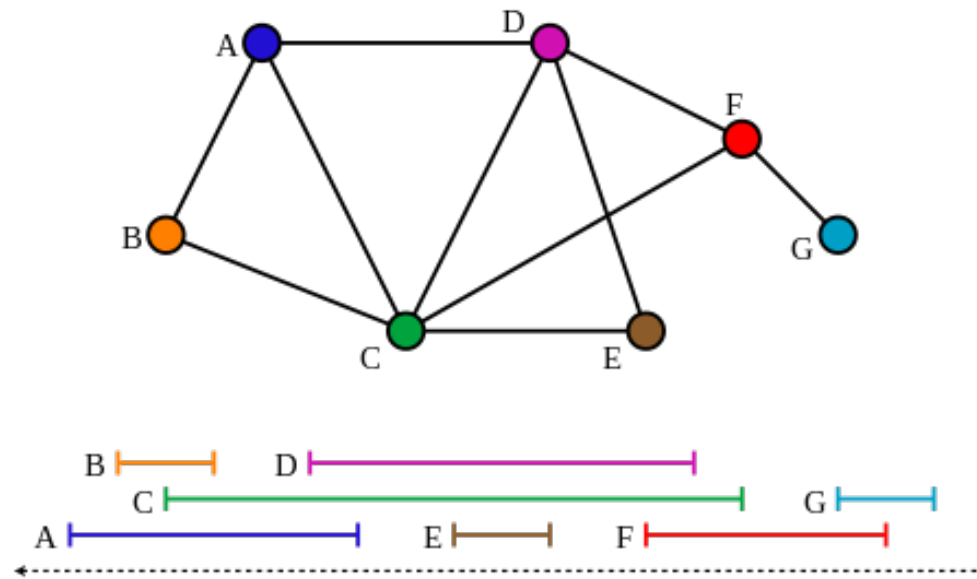
Assign colors (or tracks) using left edge algorithm



Corresponding colored conflict graph

# Interval Graph

- ▶ Conflict graphs of a set of intervals on a line (i.e., lifetimes in DFGs)
  - Vertices correspond to intervals
  - Edges correspond to interval overlap
- ▶ Interval graph belongs to the class of *perfect graphs*
  - **Clique partitioning and graph coloring are NP-hard for general graphs but can be solved efficiently for perfect graphs**

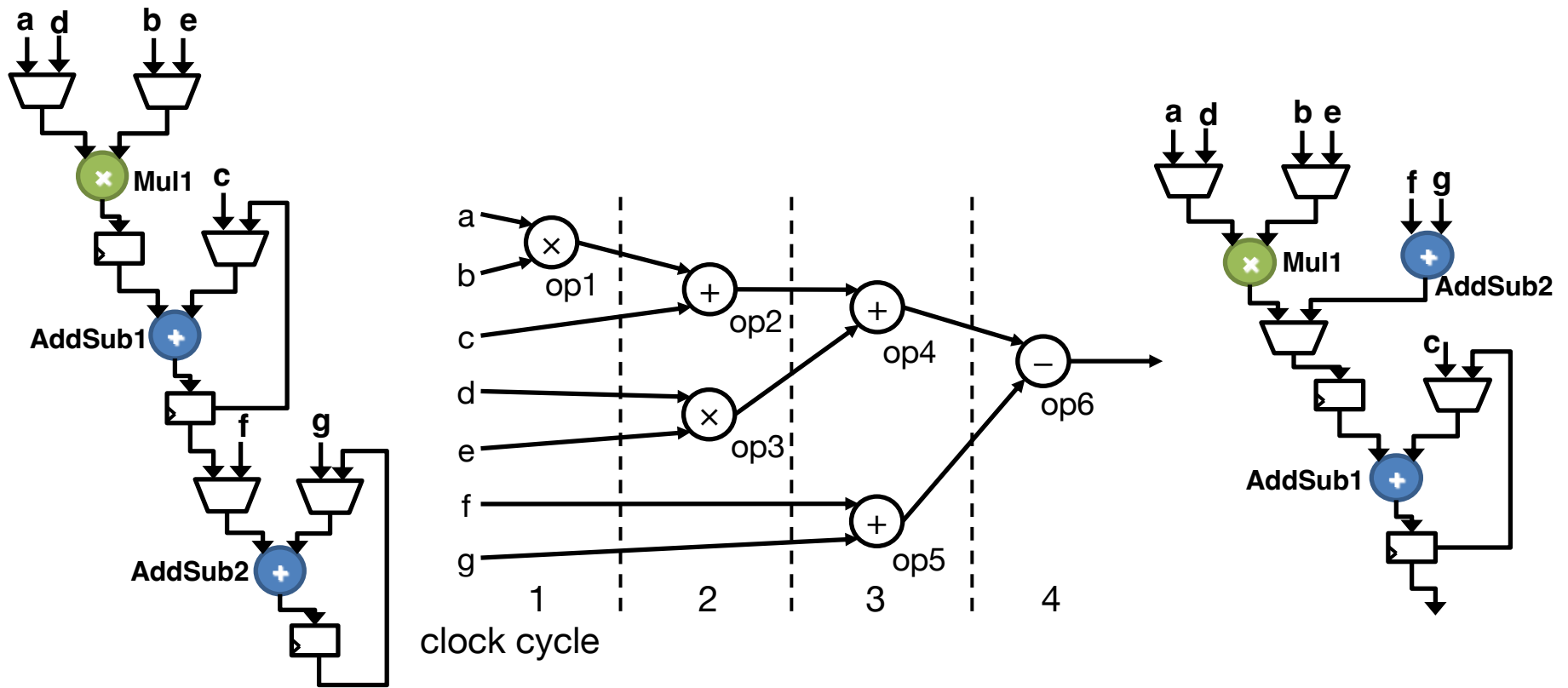


[Figure source: [en.wikipedia.org/wiki/Interval\\_graph](https://en.wikipedia.org/wiki/Interval_graph)]

# Binding Isn't Always Simple

- ▶ Resource sharing directly impacts the complexity of the resulting datapath
  - # of functional units and registers, multiplexer networks, etc.
- ▶ Binding for resource usage minimization
  - Left edge algorithm: greedy but optimal for DFGs
  - **NP-hard problem with the general form of CDFG**
    - Polynomial-time algorithm exists for SSA-based register binding, although more registers are required
- ▶ **Connectivity binding problem** (e.g., multiplexer minimization) **is NP-Hard**

# Example: Binding Impact on Multiplexer Network



Functional Unit	Operations
Mul1	op1, op3
AddSub1	op2, op4
AddSub2	op5, <u>op6</u>

Binding 1

Functional Unit	Operations
Mul1	op1, op3
AddSub1	op2, op4, <u>op6</u>
AddSub2	op5

Binding 2

# Acknowledgements

- ▶ These slides contain/adapt materials developed by
  - Prof. Jason Cong (UCLA)
  - Prof. Ryan Kastner (UCSD)
  - Dr. Stephen Neuendorffer (AMD Xilinx)