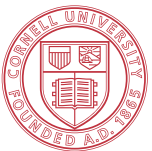




ECE 6775  
High-Level Digital Design Automation  
Fall 2025

# Pipelining



Cornell University

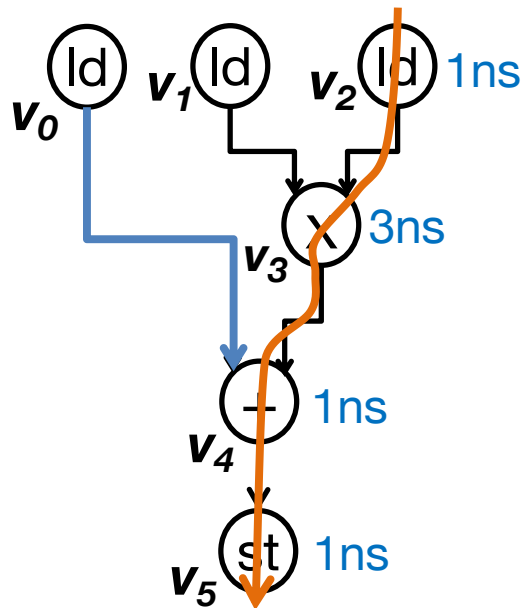


# Announcements

- ▶ Links to previous quizzes and solutions are posted on Ed (see pinned post)
  - Grades will be released later this semester

# Review: SDC-Based Scheduling

- ▶ A linear programming formulation based on system of **integer** difference constraints (SDC)



$s_i$  : schedule variable for operation  $i$

- Dependence constraints

➡  $\langle v_0, v_4 \rangle : s_0 - s_4 \leq 0$

$\langle v_1, v_3 \rangle : s_1 - s_3 \leq 0$

$\langle v_2, v_3 \rangle : s_2 - s_3 \leq 0$

$\langle v_3, v_4 \rangle : s_3 - s_4 \leq 0$

$\langle v_4, v_5 \rangle : s_4 - s_5 \leq 0$

- Cycle time constraints

$v_1 \rightarrow v_5 : s_1 - s_5 \leq -1$

➡  $v_2 \rightarrow v_5 : s_2 - s_5 \leq -1$

Timing constraints

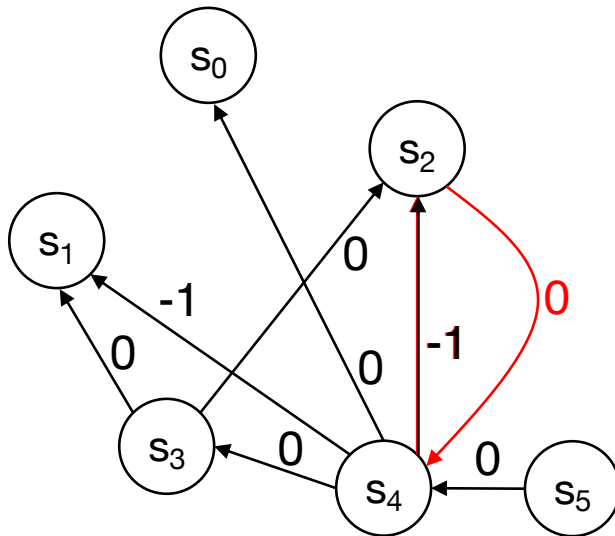
Operation chaining is naturally supported

- Target cycle time: 5ns
- Delay estimates
  - Mul (x): 3ns
  - Add (+): 1ns
  - Load/Store (ld/st): 1ns

To meet the cycle time,  $v_2$  and  $v_5$  should have a minimum separation of one cycle

# SDC Constraint Graph

- ▶ Difference constraints can be conveniently represented using **constraint graph**
  - Each vertex represents a variable, and each weighted edge corresponds to a different constraint
  - Detect infeasibility by the presence of negative cycle (by solving single-source shortest path)



$$s_2 - s_4 \leq -1$$

$$s_4 - s_2 \leq 0$$

---


$$0 \leq -1$$

$$s_0 - s_4 \leq 0$$

$$s_1 - s_3 \leq 0$$

$$s_2 - s_3 \leq 0$$

$$s_3 - s_4 \leq 0$$

$$s_4 - s_5 \leq 0$$

$$s_2 - s_4 \leq -1$$

$$s_1 - s_4 \leq -1$$

$$s_4 - s_2 \leq 0$$

# Agenda

- ▶ Introduction to pipelined scheduling
  - Parallel processing vs. Pipelining
  - Common forms in hardware accelerators
  - Throughput restrictions: resources and recurrences
- ▶ Modulo scheduling concepts
  - Recurrence and resource MII
  - Extending SDC formulation for pipelining

# Parallelization Techniques

- ▶ Parallel processing
  - Emphasizes concurrency by *replicating* a hardware structure several times (typically homogeneous)
    - High performance is attained by having all structures execute simultaneously on different parts of the problem to be solved
- ▶ Pipelining
  - Takes the approach of *decomposing* the function to be performed into smaller stages and allocating separate hardware to each stage (typically heterogeneous)
    - Data/instructions flow through the stage of a hardware pipeline at a rate (often) independent of the length of the pipeline

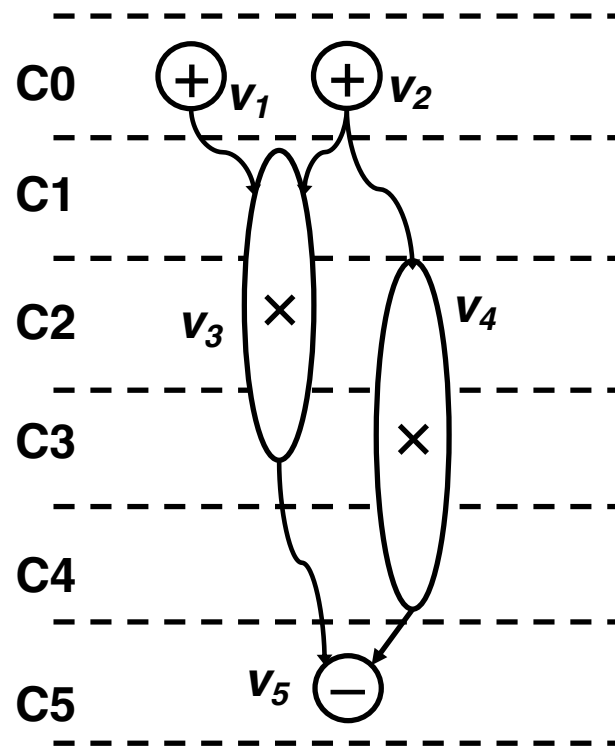
[source: Peter Kogge, The Architecture of Pipelined Computers]

# Common Forms of Pipelining

- ▶ Operator pipelining
  - Fine-grained pipeline (e.g., functional units, memories)
  - Execute a sequence of operations on a pipelined resource
- ▶ Loop/function pipelining (**focus of this class**)
  - Statically scheduled
  - Overlap successive loop iterations / function invocations at a fixed rate
- ▶ Task pipelining
  - Coarse-grained pipeline formed by multiple concurrent processes (often expressed in loops or functions)
  - Dynamically controlled
  - Start a new task before the prior one is completed

# Operator Pipelining

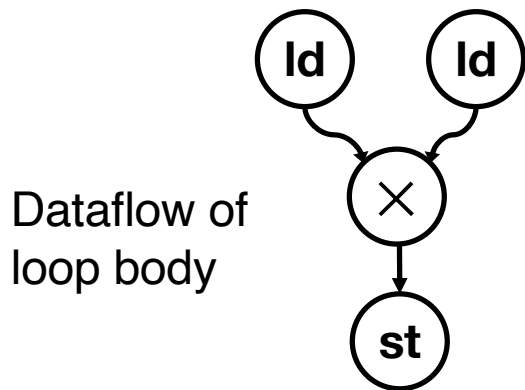
- ▶ Pipelined multi-cycle operations
  - $v_3$  and  $v_4$  can share the same pipelined multiplier (3 stages)



# Loop Pipelining

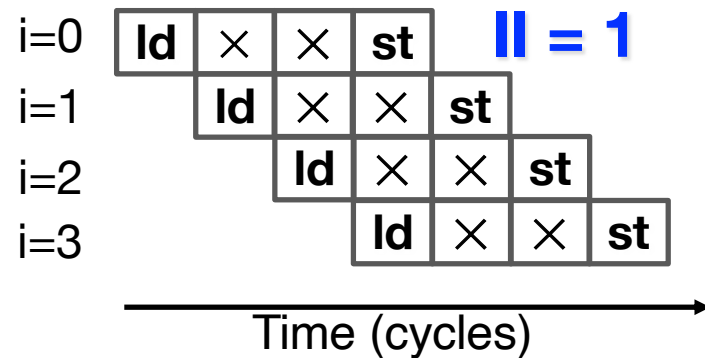
- ▶ Pipelining is one of the most important optimizations for HLS
  - Key factor: **Initiation Interval (II)**
  - Allows a new iteration to begin processing, II cycles after the start of the previous iteration (**II=1 means the loop is fully pipelined**)

```
for (i = 0; i < N; ++i)  
    p[i] = x[i] * y[i];
```



**ld** – Load (memory read)  
**st** – Store (memory write)

Pipelined schedule



Here we assume multiplication (×) takes two cycles

## Exercise: Loop Pipeline Performance

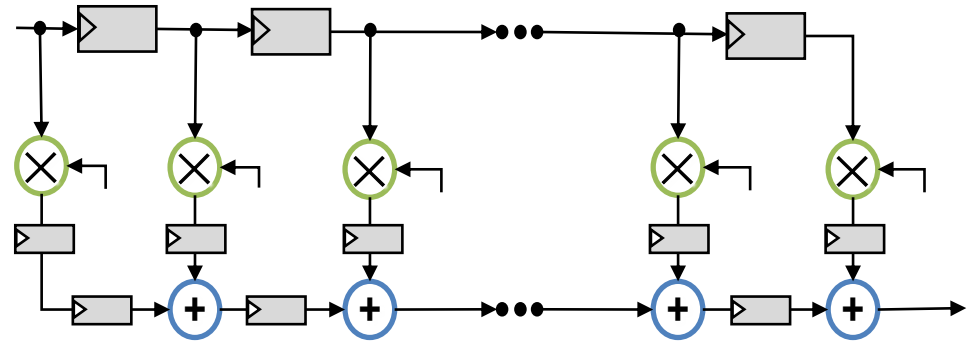
- ▶ Given a 100-iteration loop, where its loop body takes 50 cycles to execute
  - With  $II = 1$ , how many cycles is needed to complete execution of the entire loop?
  - What about  $II = 2$ ?

# Function Pipelining

- ▶ Function pipelining: Entire function is becomes a pipelined datapath

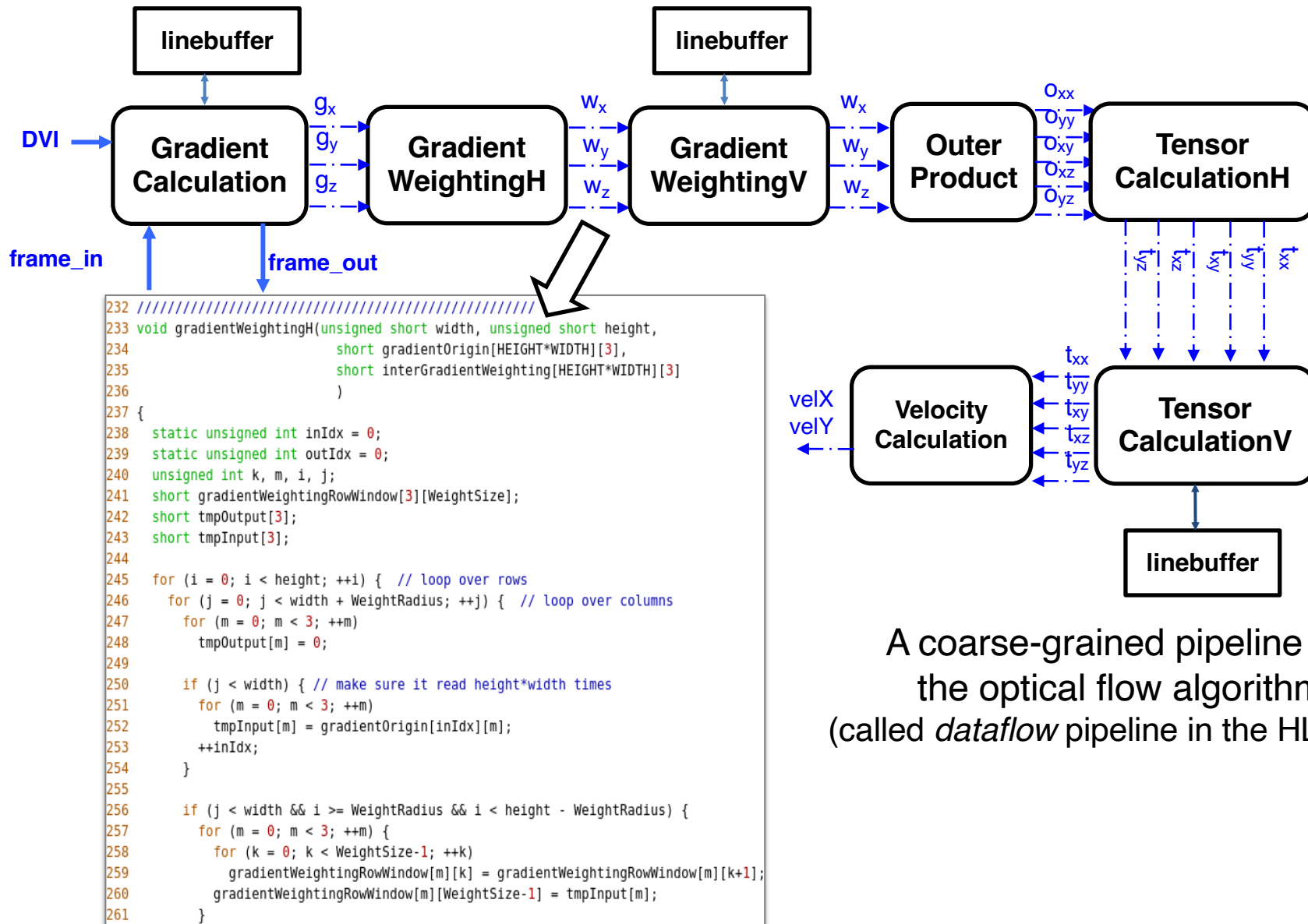
```
void fir(int *x, int *y)
{
    static int shift_reg[NUM_TAPS];
    const int taps[NUM_TAPS] =
        {1, 9, 14, 19, 26, 19, 14, 9, 1};
    int acc = 0;
    for (int i = 0; i < NUM_TAPS; ++i)
        acc += taps[i] * shift_reg[i];
    for (int i = NUM_TAPS - 1; i > 0; --i)
        shift_reg[i] = shift_reg[i-1];

    shift_reg[0] = *x;
    *y = acc;
}
```



Pipeline the entire function of the FIR filter  
(with all loops unrolled and arrays completely partitioned)

# Task Pipelining



A coarse-grained pipeline for the optical flow algorithm (called *dataflow* pipeline in the HLS tool)

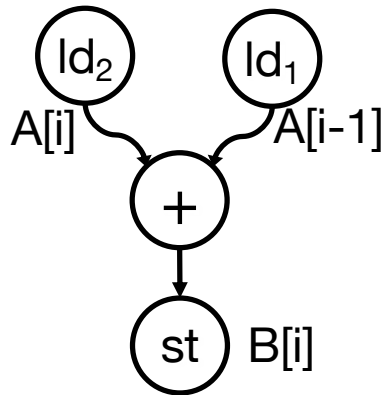
# Restrictions of Pipeline Throughput

- ▶ Resource limitations
  - Limited compute resources
  - Limited memory resources (esp. memory port limitations)
  - Restricted I/O bandwidth
  - Low throughput of subcomponent
  - ...
- ▶ Recurrences
  - Also known as feedbacks, carried dependences
  - **Fundamental limits of the throughput of a pipeline**

# Resource Limitation

- ▶ Memory is a common source of resource contention
  - e.g., memory port limitations

for ( $i = 1; i < N; ++i$ )  
 $B[i] = A[i-1] + A[i];$



	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>	+

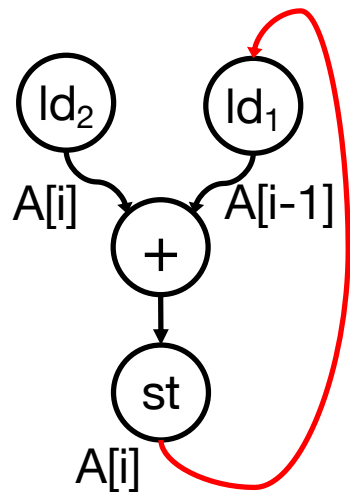
Port conflict

Assuming arrays A and B are held in two different SRAMs

Only one read port per SRAM → 1 load / cycle

# Recurrence Restriction

- ▶ Recurrences restrict pipeline throughput
  - Computation of a component depends on a previous result from the same component



**ld** – Load  
**st** – Store

```
for (i = 1; i < N; ++i)
    A[i] = A[i-1] + A[i];
```

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

Assume operation chaining is not allowed here due to cycle time constraint

# More on Recurrences

- ▶ **Recurrence** – if an operation from one iteration has *dependence* on the same operation in a previous iteration
  - Direct or indirect
  - Data or control dependence
- ▶ Types of **dependences**: true dependences, anti-dependences, output dependences
  - Intra-iteration, also known as loop-independent dependences (Lec 09)
  - **Inter-iteration**, also known as **loop-carried** dependences (focus of this lecture)
- ▶ Dependence **distance** – *number of iterations* separating the two dependent operations
  - Intra-iteration, distance = 0 (same iteration)
  - Inter-iteration, distance > 0

# True Dependences

- ▶ True dependence, also known as Read After Write (RAW) or flow dependence

**Example 1**    for (i = 0; i < N; i++)

          A[i] += A[i-1] - 1;



Inter-iteration true dependence on “A”  
(distance = 1)

**Example 2**    for (i = 0; i < N; i++)

          sum += A[i];



Inter-iteration true dependence on “sum”  
(distance = 1)

# Anti-Dependences and Output Dependences

- ▶ Anti-dependence, also known as Write After Read (WAR) dependence

**Example**

```
for (i = 1; i < N; i++) {  
    A[i-1] = b - a;  
    B[i] = A[i] + 1  
}
```

Inter-iteration anti-dependence on "A"  
(distance = 1)

- ▶ Output dependence: also known as Write After Write (WAW) dependence

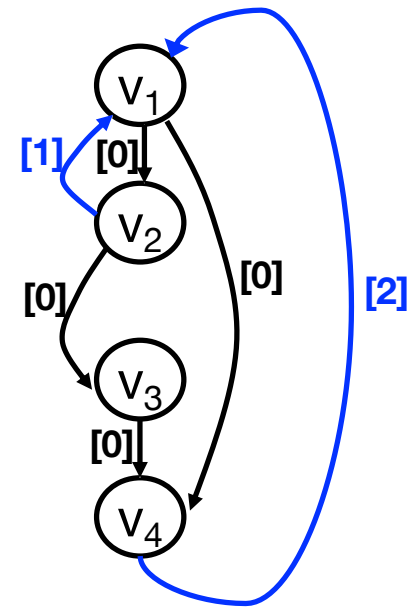
**Example**

```
for (i = 0; i < N-2; i++) {  
    B[i] = A[i-1] + 1  
    A[i] = B[i+1] + b  
    B[i+2] = b - a  
}
```

Inter-iteration output dependence on "B"  
(distance = 2)

# Dependence Graph

- ▶ Data dependences of a loop are often represented by a dependence graph
  - Forward edges: **Intra-iteration**
  - Back edges: **Inter-iteration**
  - Edges are annotated with **distance** values: number of iterations separating the two dependent operations involved
- ▶ Recurrence manifests itself as a **cycle** in the dependence graph



Edges annotated with distance values

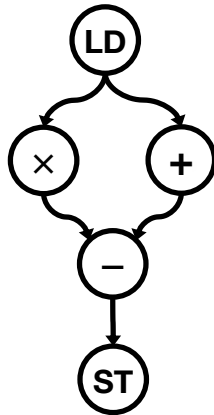
# Modulo Scheduling

- ▶ A regular form of loop (or function) pipelining technique
  - Also applies to software pipelining in compiler optimization
  - **Loop iterations use the same schedule, which are initiated at a constant rate**
  - Typical objective: Minimize initiation interval (II) under resource constraints
  
- ▶ Advantages of modulo scheduling
  - Cost efficient: No code or hardware replication
  - Easy to analyze: **Steady state determines II & resource**
  
- ▶ **NP-hard in general:** optimal polynomial time solution only exists without recurrences or resource constraints

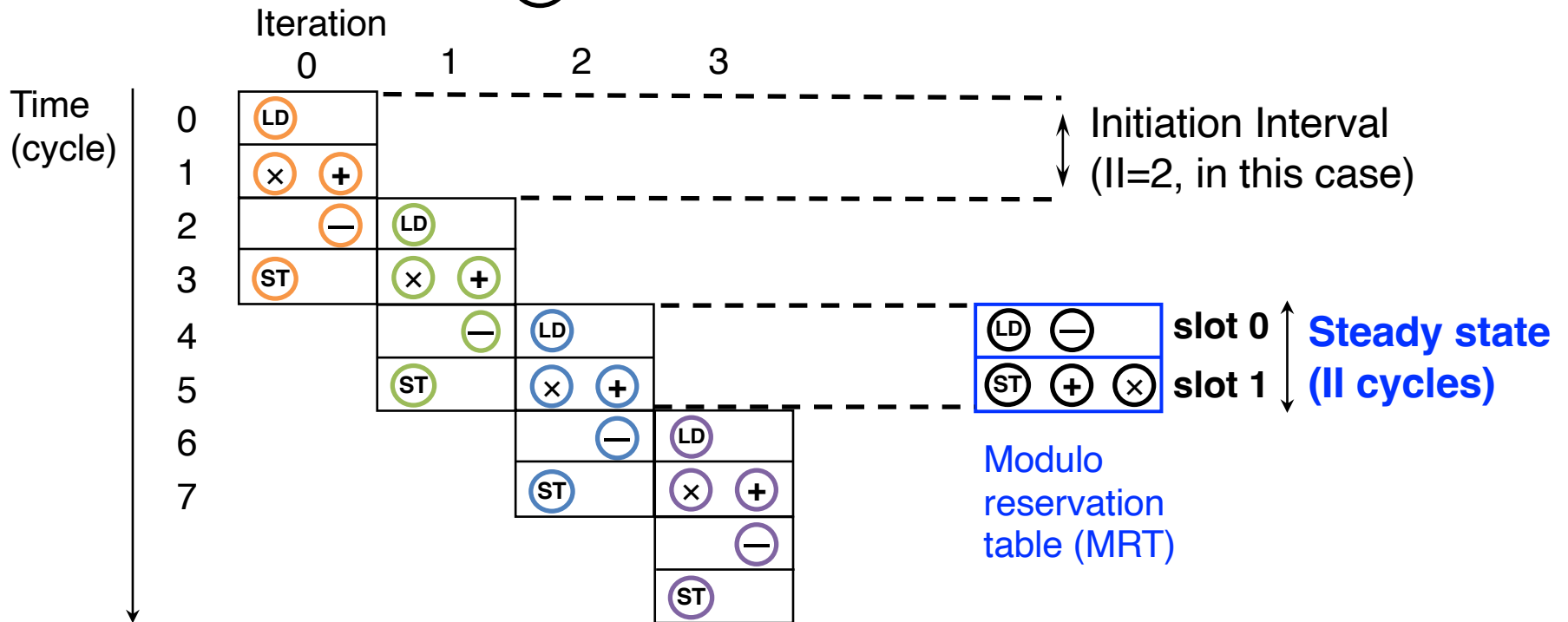
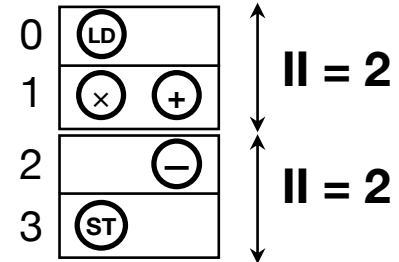
# Modulo Scheduling Example

Dependence graph of a loop body

LD: Load  
ST: Store



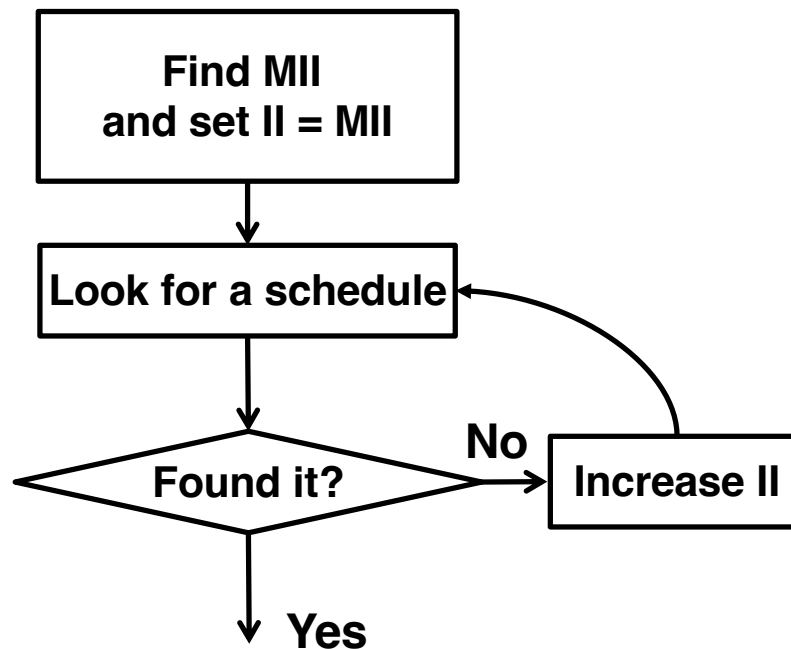
Schedule of the loop body (which is given in this example)



Steady state determines both performance and resource usage

# Heuristics for Modulo Scheduling

- ▶ A common, iterative scheme of heuristic algorithms
  - Find a lower bound on  $II$ :  $MII = \max(\text{Res}MII, \text{Rec}MII)$
  - Look for a schedule with the given  $II$
  - If a feasible schedule not found, increase  $II$  and try again



# Calculating Lower Bound of Initiation Interval

- ▶ 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  $c_i$

# Minimum II due to Resource Limits (ResMII)

- ▶ Compute ResMII: Max among all types of resources

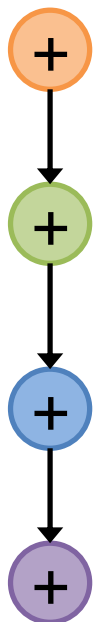
$$\text{ResMII} = \max_i \lceil \text{OPs}(r_i) / \text{Limit}(r_i) \rceil$$

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



4 adders  
(a0~a3)

Resource Allocation & Binding

	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

# Minimum II due to Recurrences (RecMII)

- Compute recurrence MII (RecMII)

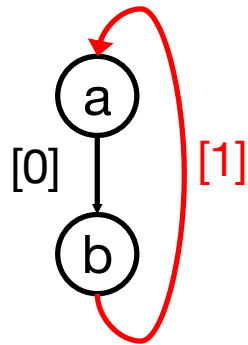
Take the max ratio among all dependence cycles

$$\text{RecMII} = \max_i \lceil \text{Latency}(c_i) / \text{Distance}(c_i) \rceil$$

Latency(c): sum of operation latencies along cycle  $c$

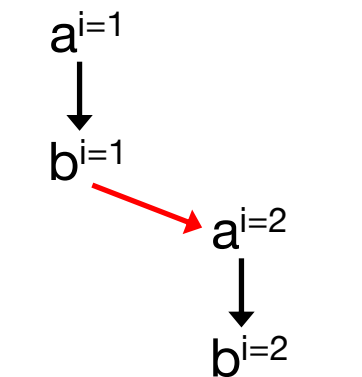
Distance(c): sum of dependence distances along cycle  $c$

Dependence



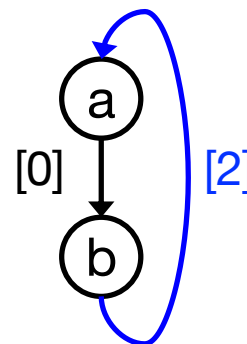
[1] dependence  
distance = 1

**RecMII** = 2/1 = 2



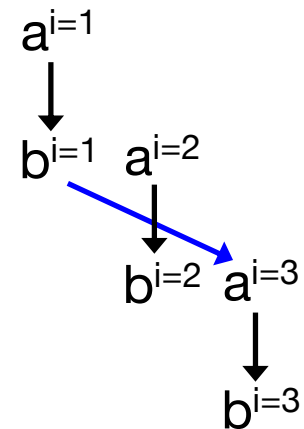
an II=2 schedule

Dependence



[2] dependence  
distance = 2

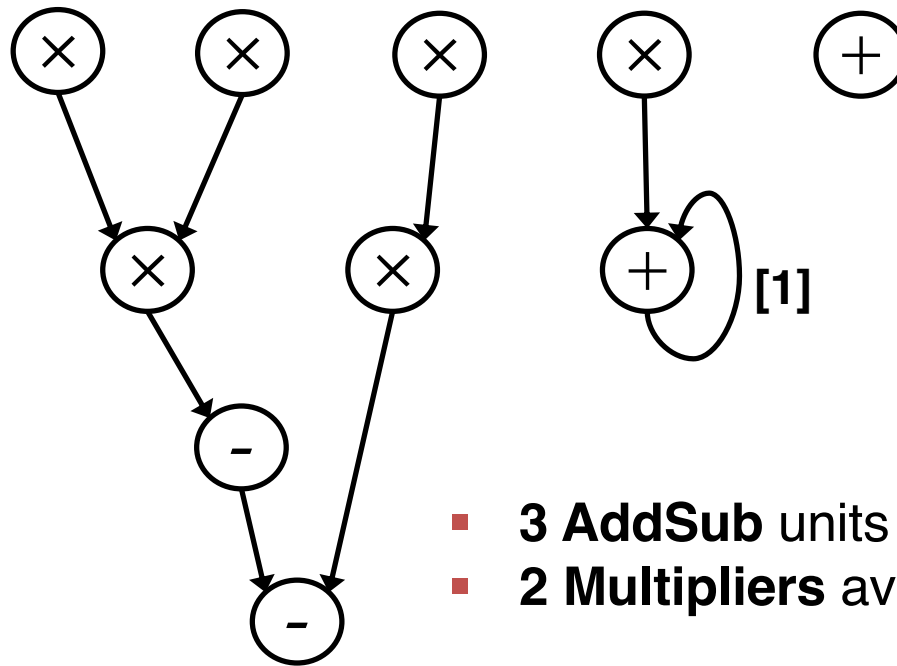
**RecMII** = 2/2 = 1



an II=1 schedule

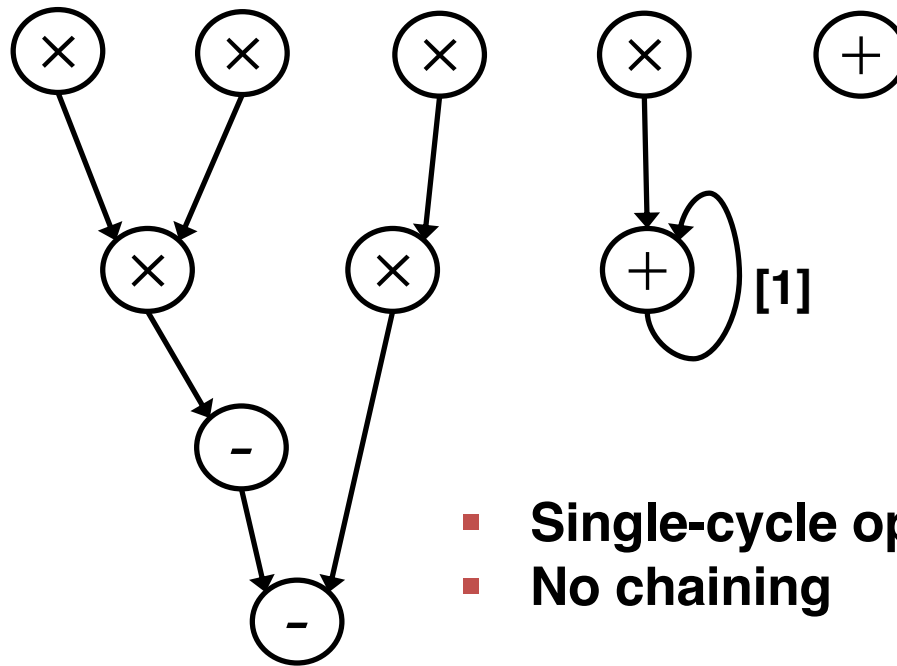
above examples assume single-cycle operations and no chaining

# What's the ResMII



Analyze the ResMII for pipelining the above DFG

# What's the RecMII

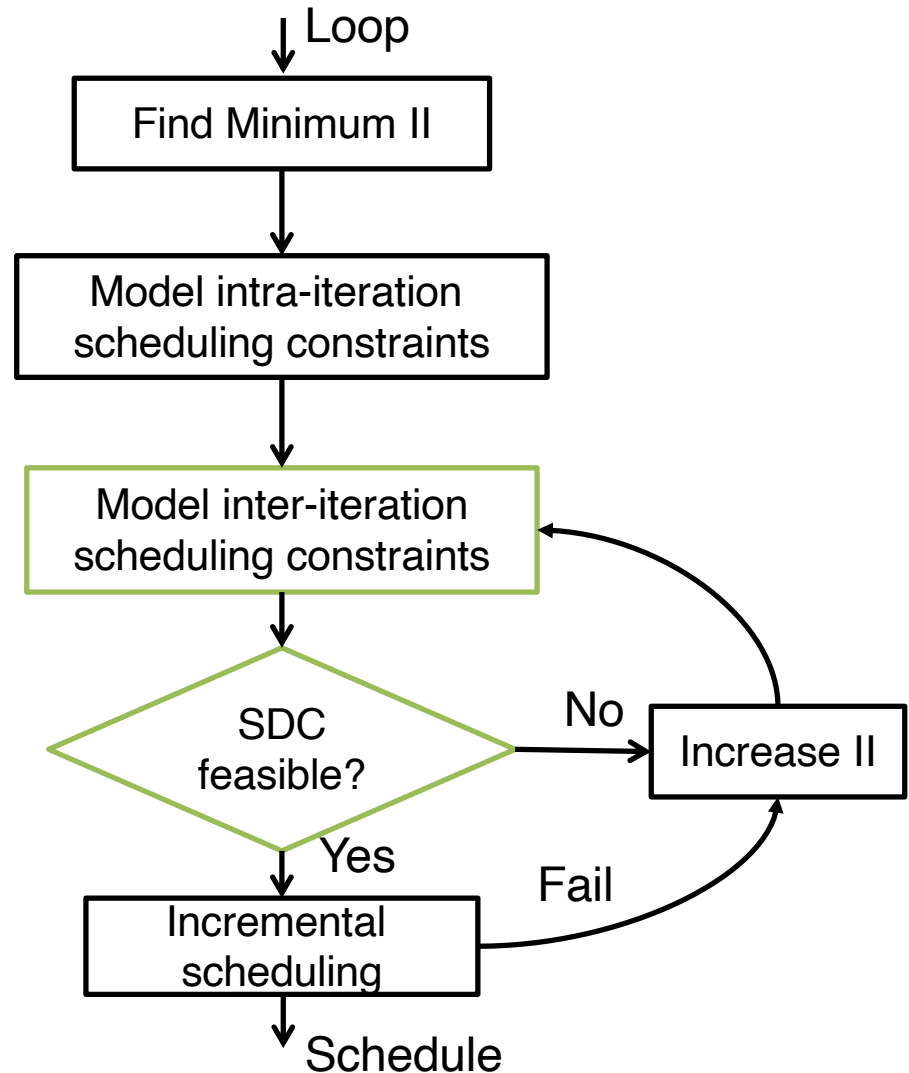


- **Single-cycle operations**
- **No chaining**

Analyze the RecMII for pipelining the above DFG

# SDC-Based Modulo Scheduling

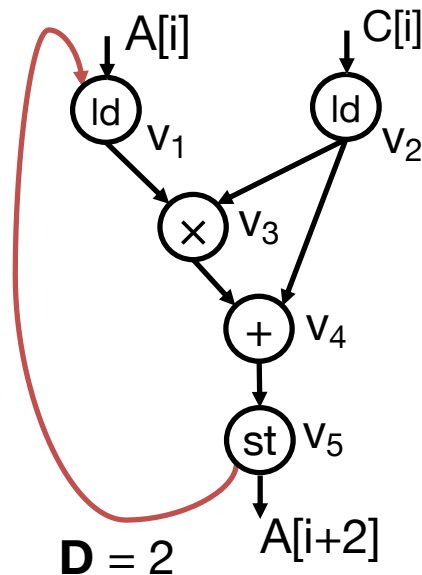
- ▶ The SDC formulation can be extended to support modulo scheduling
  - Unifies intra-iteration and inter-iteration scheduling constraints in a single SDC
  - Iterative algorithm with efficient incremental SDC update



# Modeling Loop-Carried Dependence with SDC

- ▶ Loop-carried dependence  $u \rightarrow v$  with **Distance(u, v)**, or **D** for short

```
for (i = 0; i < N-2; i++)  
{  
  B[i] = A[i] * C[i];  
  A[i+2] = B[i] + C[i];  
}
```

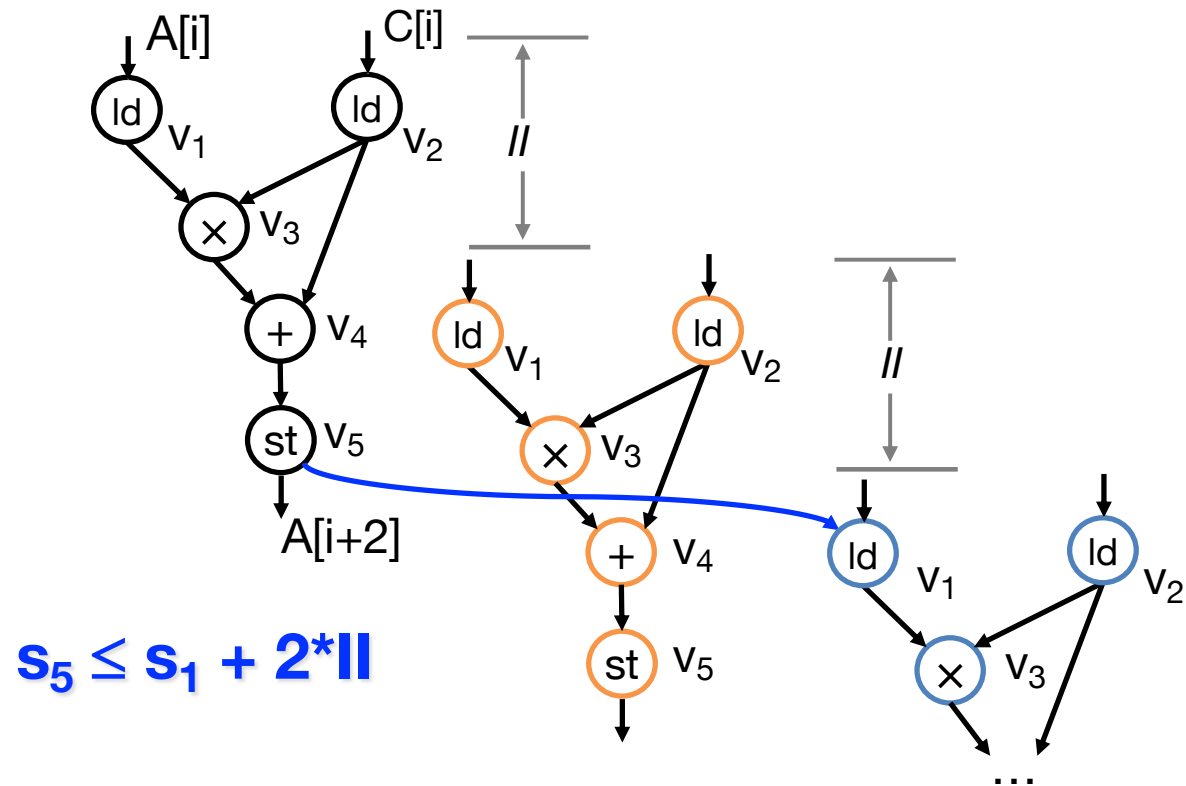


# Modeling Loop-Carried Dependence with SDC

- ▶ Loop-carried dependence  $u \rightarrow v$  with  $\text{Distance}(u, v) = D$   
 $s_u + \text{Latency}_u \leq s_v + D * II$

```

for (i = 0; i < N-2; i++)
{
    B[i] = A[i] * C[i];
    A[i+2] = B[i] + C[i];
}
    
```



# Summary

- ▶ Pipelining is one of the most commonly-used techniques in HLS to boost the performance
  - Recurrences and resource restrictions limit the pipeline throughput
- ▶ Modulo scheduling
  - A regular form of software pipeline technique
    - Also applies to loop pipelining for hardware synthesis
    - NP-hard problem in general
  - SDC-based approach provides an efficient heuristic which supports both nonpipelined and pipelined scheduling

## Next Lecture

- ▶ More Pipelining
- ▶ Resource Sharing

# Acknowledgements

- ▶ These slides contain/adapt materials developed by
  - Prof. Scott Mahlke (UMich)