ECE 6775 High-Level Digital Design Automation Fall 2024

Pipelining



Cornell University

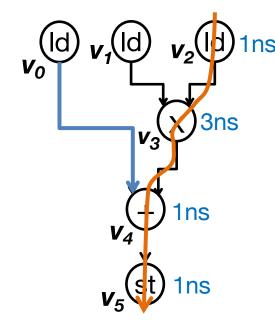


Announcements

- HW 2 posted, due Friday
 - NO penalty for late submissions, *up to 5 days late* (until Wed Oct 16); Solution will be released on Thursday Oct 17th.

Review: SDC-Based Scheduling

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



- \mathbf{s}_i : schedule variable for operation *i*
- $\begin{array}{c|c} \bullet & \text{Dependence constraints} \\ \hline & \checkmark & <v_0 \ , \ v_4 > \ : \ s_0 s_4 \leq 0 \\ \hline & <v_1 \ , \ v_3 > \ : \ s_1 s_3 \leq 0 \\ \hline & \text{chaining is} \ & <v_2 \ , \ v_3 > \ : \ s_2 s_3 \leq 0 \end{array}$
- - Cycle time constraints $v_1 \rightarrow v_5 : s_1 - s_5 \le -1$

 \downarrow $v_2 \rightarrow v_5 : s_2 - s_5 \le -1$

Timing constraints

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

have a minimum separation of one cycle

To meet the cycle time, v_2 and v_5 should

[J. Cong & Z. Zhang, DAC, 2006] [Z. Zhang & B. Liu, ICCAD, 2013]

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

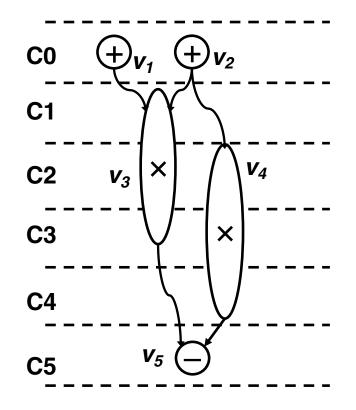
[source: Peter Kogge, The Architecture of Pipelined Computers]

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)

ld

 Allows a new iteration to begin processing, II cycles after the start of the previous iteration (II=1 means the loop is fully pipelined)

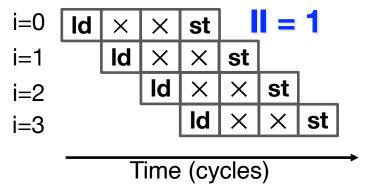
for (i = 0; i < N; ++i) p[i] = x[i] * y[i];

Х

st

ld

Dataflow of loop body Pipelined schedule



Id – Load (memory read)st – Store (memory write)

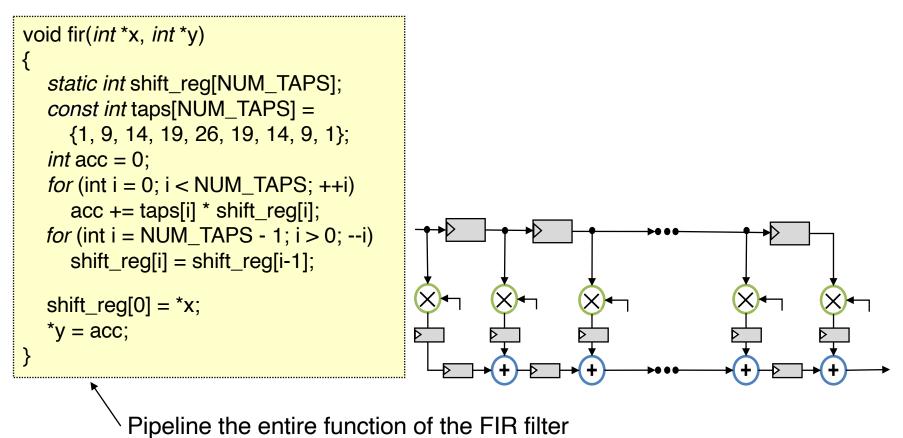
Here we assume multiplication (\times) takes two cycles

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?

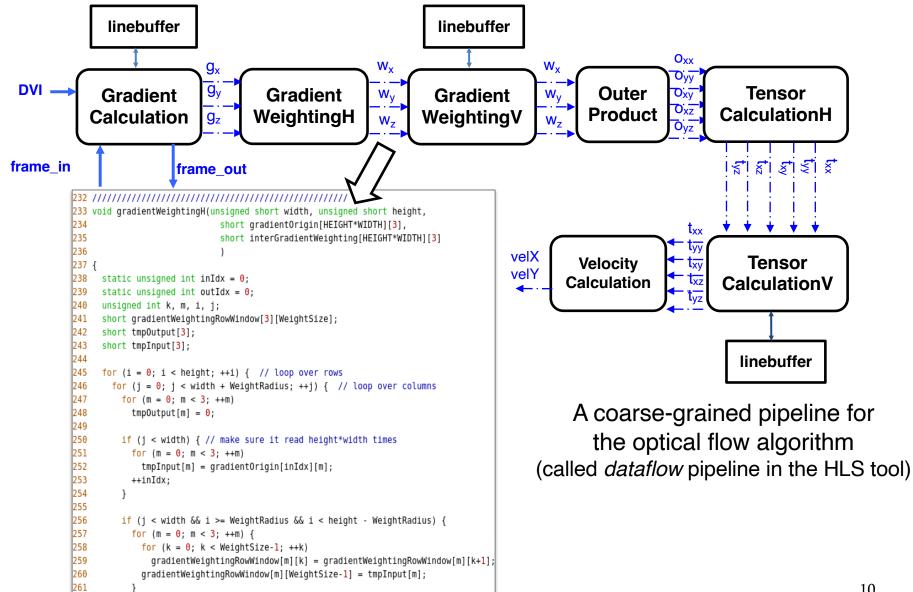
Function Pipelining

 Function pipelining: Entire function is becomes a pipelined datapath



(with all loops unrolled and arrays completely partitioned)

Task Pipelining



Restrictions of Pipeline Throughput

- Resource limitations
 - Limited compute resources
 - Limited memory resources (esp. memory port limitations)
 - Restricted I/O bandwidth
 - Low throughput of subcomponent

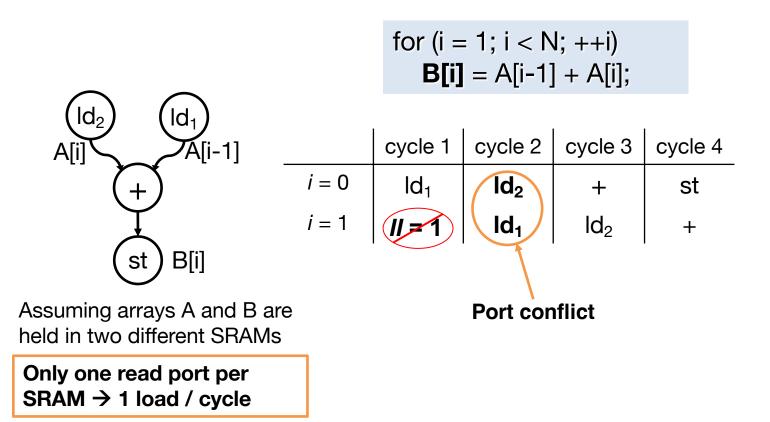
• • •

Recurrences

- Also known as feedbacks, carried dependences
- Fundamental limits of the throughput of a pipeline

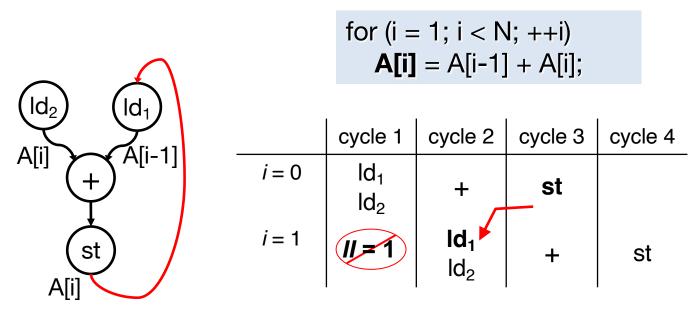
Resource Limitation

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



Recurrence Restriction

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



ld – Load st – Store

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 <u>Read After Write</u> (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 <u>Write After Read</u> (WAR) dependence

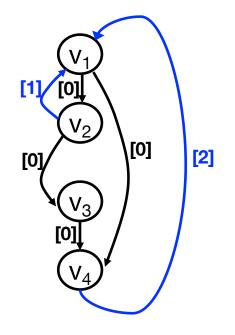
Example for (i = 1; i < N; i++) { A[i-1] = b - a; Inter-iteration anti-dependence on "A" B[i] = A[i] + 1 (distance = 1) }

 Output dependence: also known as <u>Write After Write</u> (WAW) dependence

Example Inter-iteration output dependence on "B" (distance = 2) $for (i = 0; i < N-2; i++) {$ B[i] = A[i-1] + 1 A[i] = B[i+1] + bB[i+2] = b - a

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

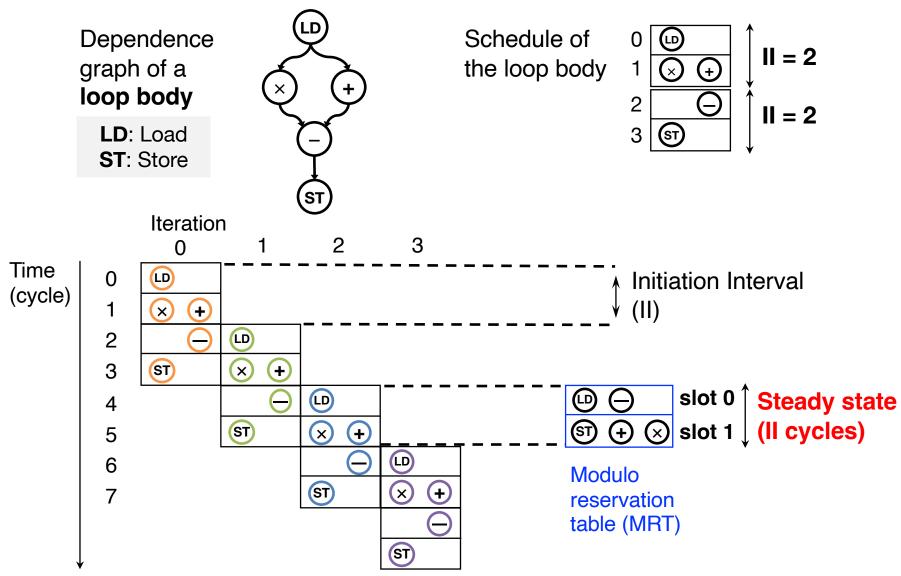


Edges annotated with distance values

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 <u>or</u> resource constraints

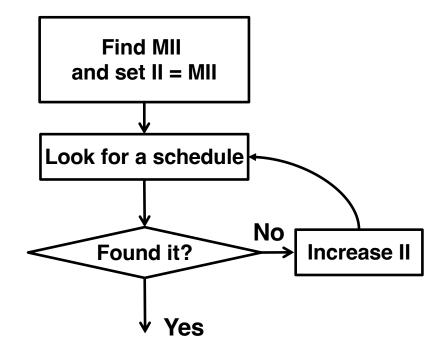
Modulo Scheduling Example



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 (ResMII, RecMII)
 - 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 (ResMII, RecMII)
 - A lower bound, not necessarily achievable
- Resource constrained MII (ResMII)
 - ResMII = max_i [OPs(r_i) / Limit(r_i)]
 OPs(r): number of operations that use resource of type r Limit(r): number of available resources of type r
- Recurrence constrained MII (RecMII)
 - RecMII = max_i [Latency(c_i) / Distance(c_i)]
 Latency(c_i): total latency in dependence cycle c_i
 Distance(c_i): total distance in dependence cycle c_i

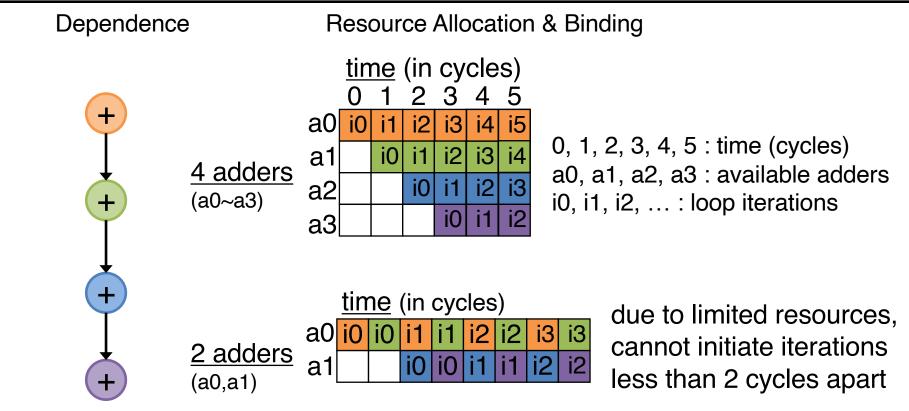
Minimum II due to Resource Limits (ResMII)

Compute ResMII: Max among all types of resources

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

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

Take the max ratio among all resource types

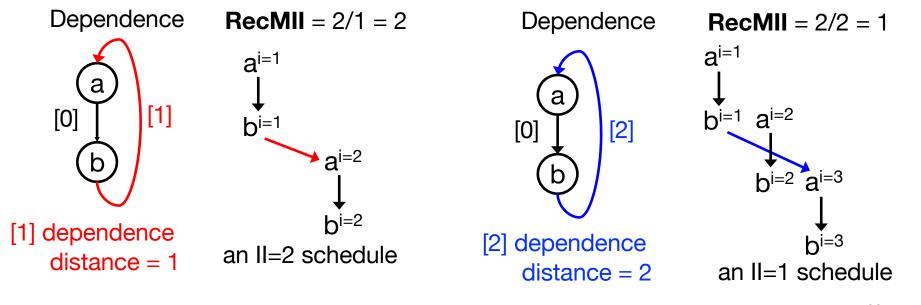


Minimum II due to Recurrences (RecMII)

Compute recurrence MII (RecMII)

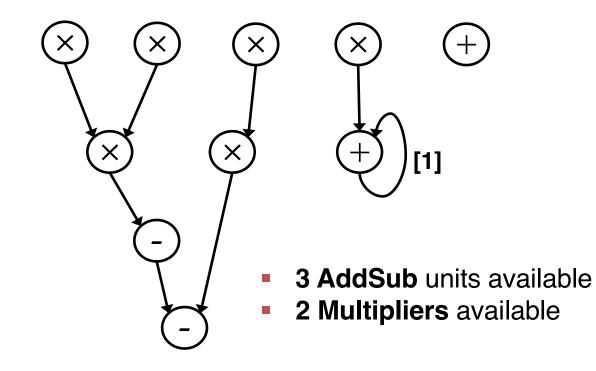
Take the max ratio among all dependence cycles **RecMII = max_i** $\begin{bmatrix} Latency(c_i) / Distance(c_i) \end{bmatrix}$

Latency(c): sum of operation latencies along cycle *c* Distance(c): sum of dependence distances along cycle *c*



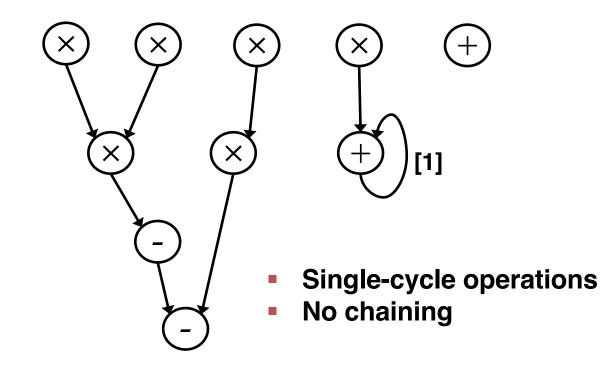
above examples assume single-cycle operations and no chaining

What's the ResMII



Analyze the ResMII for pipelining the above DFG

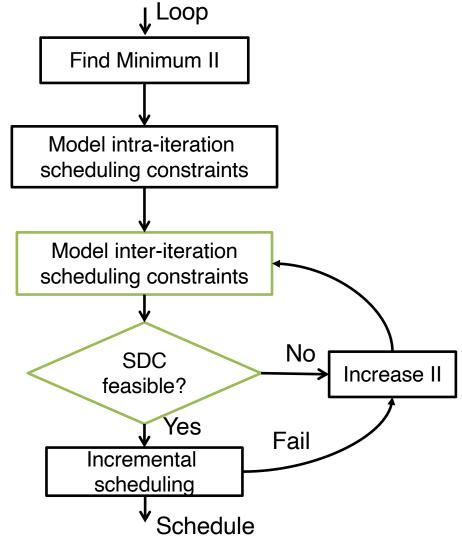
What's the RecMII



Analyze the RecMII for pipelining the above DFG

SDC-Based Modulo Scheduling

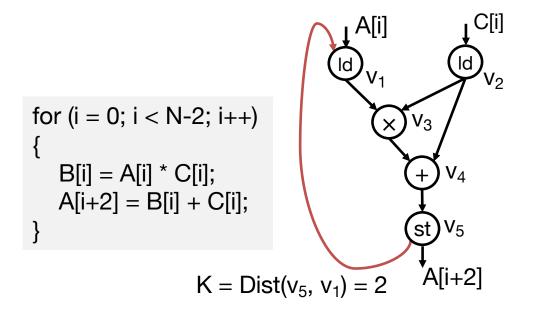
- The SDC formulation can be extended to support modulo scheduling
 - Unifies intra-iteration and interiteration scheduling constraints in a single SDC
 - Iterative algorithm with efficient incremental SDC update



[Z. Zhang & B. Liu, ICCAD 2013] ²⁶

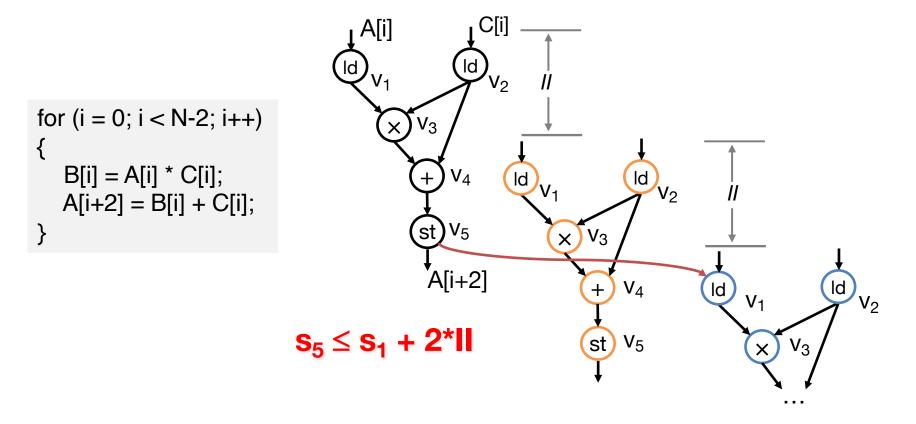
Modeling Loop-Carried Dependence with SDC

► Loop-carried dependence $u \rightarrow v$ with Distance(u, v) = K



Modeling Loop-Carried Dependence with SDC

Loop-carried dependence u → v with Distance(u, v) = K s_u + Latency_u ≤ s_v + K*II



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)