ECE 5997 Hardware Accelerator Design & Automation Fall 2021

Resource Sharing Pipelining



Cornell University



Announcements

- Yichi Zhang (PhD TA) will give a tutorial on neural networks next Tuesday in Rhodes 310
- Lab 3 is released (due 12/1)
 - Go through the CORDIC tutorial asap

Agenda

- Resource sharing overview
 - Sub-problems: functional unit, register, and connectivity binding problems
 - Key concepts: compatibility and conflict graphs
- Introduction to pipelining
 - Common forms in hardware accelerators
 - Throughput restrictions
 - Dependence types

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

Operation Binding



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 7		

Register Binding



Variable Lifetime Analysis



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



Note: A compatibility/conflict graphs for variables/registers can be constructed in a similar way

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



Example: Meeting Assignment Problem

Meeting	Schedule (am)
А	9:00~11:00
В	9:30~10:00
С	10:00~11:00
D	11:00~11:30

Conflict graph (chromatic number?)



Compatibility graph (max clique size?)

Perfect Graphs

- Clique partitioning and graph coloring problems are NP-hard on general graphs, with the exception of perfect graphs
- Definition of perfect graphs
 - For every induced subgraph, the size of the maximum (largest) clique equals the chromatic number of the subgraph
 - Examples: bipartite graphs, chordal graphs, etc.
 - <u>Chordal graphs</u>: every cycle of four or more vertices has a chord, i.e., an edge between two vertices that are not consecutive in the cycle.

Interval Graph

- Intersection graphs of a (multi)set of intervals on a line
 - Vertices correspond to intervals
 - Edges correspond to interval intersection
 - A special class of chordal graphs



[Figure source: en.wikipedia.org/wiki/Interval_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(nlogn) time

Left Edge Demonstration





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 17	

Binding Summary

- 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

Parallelization Techniques

- Parallel processing
 - Emphasizes concurrency by <u>replicating</u> a hardware structure several times (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 <u>decomposing</u> the function to be performed into smaller stages and allocating separate hardware to each stage (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

- Loop pipelining is one of the most important optimizations for high-level synthesis
 - Key metric: Initiation Interval (II) in # cycles
 - Allows a new iteration to begin processing every II cycles, before the previous iteration is complete



Pipeline Performance

- Given a 100-iteration loop with the loop body taking 50 cycles to execute
 - If we pipeline the loop with II = 1, how many cycles do we need to complete execution of the entire loop ?
 - What about II = 2?

Function Pipelining

 Function pipelining: Entire function is becomes a pipelined datapath



Pipeline the entire function of the FIR filter (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



Id – Load st – Store Assume chaining is not possible on memory reads (i.e., Id) and writes (i.e., st) due to cycle time constraint

Next Lecture

More Pipelining

Acknowledgements

- These slides contain/adapt materials developed by
 - Prof. Deming Chen (UIUC)
 - Prof. Jason Cong (UCLA)