ECE 6775 High-Level Digital Design Automation Fall 2024

More Scheduling



Cornell University



Announcements

- Lab 3 due Monday
- HW 2 will be released soon
 - Two problems related to pipelining need next week's lecture content

Agenda

- ILP for time-constrained scheduling
- Heuristic algorithms for constrained scheduling
 - List scheduling
 - SDC-based scheduling

Recap: ILP Formulation for Dependence Constraints

 Using ASAP and ALAP, the non-trivial inequalities are: (assuming no chaining and single-cycle ops)

$$\begin{aligned} 2x_{7,2} + 3x_{7,3} - x_{6,1} - 2x_{6,2} &\geq 1 \\ 4x_{5,4} - 2x_{7,2} - 3x_{7,3} &\geq 1 \\ 2x_{9,2} + 3x_{9,3} + 4x_{9,4} - x_{8,1} - 2x_{8,2} - 3x_{8,3} &\geq 1 \\ 2x_{11,2} + 3x_{11,3} + 4x_{11,4} - x_{10,1} - 2x_{10,2} - 3x_{10,3} &\geq 1 \end{aligned}$$



assume L=4 and no chaining

Recap: ILP Formulation for Resource Constraints

Resource constraints (assuming 2 multipliers and 1 ALU)

$$\begin{aligned} x_{1,1} + x_{2,1} + x_{6,1} + x_{8,1} &\leq 2 \\ x_{3,2} + x_{6,2} + x_{7,2} + x_{8,2} &\leq 2 \\ x_{7,3} + x_{8,3} &\leq 2 \end{aligned}$$

$$\begin{aligned} x_{10,1} &\leq 1 \\ x_{9,2} + x_{10,2} + x_{11,2} &\leq 1 \\ x_{4,3} + x_{9,3} + x_{10,3} + x_{11,3} &\leq 1 \\ x_{5,4} + x_{9,4} + x_{11,4} &\leq 1 \end{aligned}$$



assume L=4 and no chaining

Another Exercise: Formulating ILP

- Minimize the number of classrooms that the school must allocate for the following courses
- Steps to formulate the ILP
 - 1 Create variables
 - 2 Each course to be scheduled to exactly one of the preferred slots
 - ③ Determine the number of rooms required (by creating derived variables)

4 Set up the objective function

Course	Preferred Slots
А	(1) (2)
В	(1) (3)
С	(2) (3)
D	(2)

(1) 8:00 – 10:00am
(2) 10:00am - 12:00pm
(3) 12:00 – 2:00pm

1 x_{i,s} : course i uses slot s

(2)
$$x_{A,1} + x_{A,2} = 1$$

 $x_{B,1} + x_{B,3} = 1$
 $x_{C,2} + x_{C,3} = 1$
 $x_{D,2} = 1$
(3) $r_1 = x_{A,1} + x_{B,1}$
 $r_2 = x_{A,2} + x_{C,2} + x_{D,2}$
 $r_3 = x_{B,3} + x_{C,3}$

$$\begin{array}{l} \text{min } \mathsf{R} \\ \mathsf{R} \geq \mathsf{r}_1, \, \mathsf{R} \geq \mathsf{r}_2, \, \mathsf{R} \geq \mathsf{r}_3 \end{array}$$

5

Time-Constrained Scheduling (TCS)

- Dual problem of resource-constrained scheduling
 - Overall latency is given as a constraint (deadline)
 - Minimize the total cost in terms of area (or resource usage), power, etc.
- NP-hard problem
 - ILP formulation is exact but is not a polynomial-time solution
 - Force-directed scheduling is a well-known heuristic for TCS (see De Micheli chapter 5.4.4)

Example: ILP Formulation for TCS

ILP for time-constrained scheduling minimize c^Ty

$$\begin{split} & x_{1,1} + x_{2,1} + x_{6,1} + x_{8,1} \leq y_1 \\ & x_{3,2} + x_{6,2} + x_{7,2} + x_{8,2} \leq y_1 \\ & x_{7,3} + x_{8,3} \leq y_1 \\ & x_{5,4} + x_{9,4} + x_{11,4} \leq y_2 \end{split}$$

. . .

What does the y vector represent?



Constrained Scheduling in HLS

- Constrained scheduling
 - General case NP-hard
 - Resource-constrained scheduling (RCS)
 - Minimize latency given constraints on area or resources
 - Time-constrained scheduling (TCS)
 - Minimize resources subject to bound on latency

Exact methods

- Integer linear programming (ILP)
- Symbolic scheduling using BDDs
- Hu's algorithm for a very restricted problem
- Heuristics
 - List scheduling
 - Force-directed list scheduling
 - SDC-based scheduling

Symbolic Scheduling: Representing Resource Constraints with BDD (an example)

- Assume 2 multipliers are available and there are 4 potential multiply operations at control step k
- The following Boolean expression captures the resource constraint at step k

$$x_{1,k}' \bullet x_{2,k}' + x_{1,k}' \bullet x_{3,k}' + x_{1,k}' \bullet x_{4,k}' + x_{2,k}' \bullet x_{3,k}' + x_{2,k}' \bullet x_{4,k}' + x_{3,k}' \bullet x_{4,k}'$$

This expression indicates that at least (4 - 2) multiplication operations (among 4 potential operations in step k) cannot be scheduled to the same step



List Scheduling

- A widely-used heuristic algorithm for RCS
 - Schedule one control step (cycle) at a time
 - Maintain a list of "ready" operations considering dependence
 - Assign priorities to operations; most "critical" operations (with the highest priorities) go first
- Often refers to a family of algorithms
 - Typically classified by the way priority function is calculated
 - Static priority: Priorities are calculated once before scheduling
 - Dynamic priority calculation: Priorities are updated during scheduling

Static Priority Example: Node Height



Nodes are labelled with distance to sink (height)

Ready operations are colored in green

- Assumptions:
 - All operations have unit delay
 - 2 MULTs, 1 AddSub, and 1 CMP available

Ready Nodes with Highest Priorities Picked First



- Assumptions:
 - All operations have unit delay
 - 2 MULTs, 1 AddSub, and 1 CMP available

Update Ready Nodes and Repeat for Each Step



- Assumptions:
 - All operations have unit delay
 - 2 MULTs, 1 AddSub, and 1 CMP available

Update Ready Nodes and Repeat for Each Step



- Assumptions:
 - All operations have unit delay
 - 2 MULTs, 1 AddSub, and 1 CMP available

Repeat Until All Nodes Scheduled





- Assumptions:
 - All operations have unit delay
 - 2 MULTs, 1 AddSub, and 1 CMP available

A Special Case

- With the following (very) restrictive conditions:
 - All operations have unit delay (i.e., single cycle)
 - All operations (and resources) are of the same type
 - Graph is a forest
- List scheduling with static height-based priorities guarantees optimality
- This is known as Hu's algorithm
 - T. C. Hu, Parallel sequencing and assembly line problems. Operations Research, 9(6), 841-848, 1961
 - Guarantees

HLS Scheduling: Tension between Scalability and Quality



More Realistic Scheduling Problems

- Operation chaining
 - More compact schedule
- Multi-cycle operations
 - Nonpipelined or pipelined
 - Higher frequency
- Mutually exclusive operations
 - Scheduled in the same step, but with mutually exclusive execution conditions
 - Higher resource utilization
- Other timing constraints
 - Frequency constraints, latency constraints, relative time constraints





A Simple Operation Chaining Problem

Given: A chain of n operations. Without any registers, the cycle time equals the total combinational delay, which is $D = sum(d_i)$.



Question: How to place TWO registers on the chain to achieve the minimum cycle time?

Example:
$$\rightarrow 5 \rightarrow 1$$
 $\rightarrow 6 \rightarrow 2$ $\rightarrow 7 \rightarrow$

SDC-Based Scheduling

SDC = System of difference constraints



- \mathbf{s}_i : schedule variable for operation *i*
 - Dependence constraints $\vec{}$ $\Rightarrow < v_0, v_4 > : s_0 - s_4 \le 0$

Cycle time constraints

) constraints

Timing

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

 $v_1 \rightarrow v_5 : s_1 - s_5 \le -1$ $\downarrow v_2 \rightarrow v_5 : s_2 - s_5 \le -1$

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

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

Exercise: Latency Constraint in SDC



How to enforce that operations v_3 and v_4 are not chained and at most two cycles apart?

Difference Constraints

- A difference constraint is a formula in the form of x − y ≤ b or x − y < b for numeric variables x and y, and constant b
- With scheduling variables, we use integer difference constraints to model a variety of scheduling constraints
 - x and y must have integral values
 - Thus *b* only needs to be an integer => form x y < b is redundant

SDC Constraint Matrix

- The constraint matrix of SDC(X, C) is a totally unimodular matrix (TUM):
 - Every nonsingular square submatrix has a determinant of -1/+1.

$$\begin{bmatrix} 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 \\ 0 & 0 & 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & 0 & 0 & -1 \end{bmatrix} \begin{pmatrix} s_0 \\ s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{bmatrix} \leq \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ -1 \\ -1 \\ -1 \end{bmatrix}$$

- Theorem (Hoffman & Kruskal, 1956): If A is totally unimodular and b is a vector of integers, every extreme point of polyhedron {x : Ax ≤ b} is integral.
 - Solving linear programming (LP) relaxation leads to integral solutions

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)



Handling Resource Constraints (NP-Hard in General)

Resource constraints cannot be represented exactly in integer difference form*



Resource constraint
 Two read ports

- Resource constraints
 → Heuristic partial orderings
 v₀ → v₂ : s₀ s₂ ≤ -1 3 cycle latency
 - $\begin{array}{l} v_1 \rightarrow v_0: s_1 s_0 \leq -1 \\ v_2 \rightarrow v_0: s_2 s_0 \leq -1 \end{array} \hspace{0.1 cm} \text{2 cycle latency} \end{array}$

OR

*A more recent SDC scheduling paper combined SAT and SDC to solve RCS exactly [Dai et al. FPGA'2018]

Linear Objectives

- ASAP: min $\sum_{i \in V} s_i$
- ALAP: max $\sum_{i \in V} s_i$
- Minimum latency: min $max_{i \in V} \{s_i\}$ ►
- Minimum average case latency (control-intensive design)

v₅(st

Many other ...



Control Flow Graphs

- Control dependencies can also be honored
 - If bb₂ is control dependent on bb₁, the operation nodes of bb₂ are not allowed to be scheduled before those of bb₁
 - Polarize each basic block bb_i with two scheduling variables (head and tail)
 - $\forall v \in bb_i$, $s_h(bb_i) s_h(v) \le 0$
 - $\forall v \in bb_i$, $s_t(v) s_t(bb_i) \leq 0$
 - If $e_c(bb_i, bb_j) \in E_c$ and e_c is not a back edge
 - $s_t(bb_i) s_h(bb_j) \leq 0$



 $s_t(B_1) - s_h(B_2) \le 0$

Example: Greatest Common Divisor (GCD)



GCD in SSA form



Interpreting the LP Solution of SDC Scheduling



Operations and Predicates

 $x_0 = in1$ $y_0 = in2$ $cond1 = (x_0 != y_0)$

 $\begin{array}{l} x_1 = \Phi \; (x_0, \, x_1, \, x_2) \\ y_1 = \Phi \; (y_0, \, y_1, \, y_2) \\ cond2 = (x_1 > y_1) \\ x_2 = x_1 - y_1 \\ cond3 = (x_2 \; != \; y_1) \\ y_2 = y_1 - x_1 \\ cond4 = (x_1 \; != \; y_2) \\ x_3 = \Phi \; (x_0, \, x_1, \, x_2) \\ ^* out = x_3 \end{array}$





!cond1 && (cond3 || cond4)

Scheduling Summary

► ILP

- Exact, but exponential worst-case runtime
- Hu's algorithm
 - Optimal and polynomial
 - Only works in very restricted cases
- List scheduling
 - Extension to Hu's for general cases
 - Greedy (fast) but suboptimal
- SDC-based scheduling
 - A versatile heuristic based on LP formulation with different constraints
 - Amenable to global optimization

Next Lecture

Pipelining

Acknowledgements

- These slides contain/adapt materials developed by
 - Ryan Kastner (UCSD)