# Introduction of Processor Design for AI Applications

## L05 – Retiming and Pipelining

Pengju Ren Institute of Artificial Intelligence and Robotics Xi'an Jiaotong University

http://gr.xjtu.edu.cn/web/pengjuren

#### LTI Systems

#### Linear systems

 $\Box$ Assume x1 (n)->y1 (n) and x2 (n)->y2 (n), where "->" denotes "lead to". If ax1 (n)+bx2 (n) -> ay1 (n)+by2 (n), then the systems is referred to as a "Linear System."

□ Homogenous and additive properties

Time-invariant (TI) systems
 \$\Box\$x(n-n0)->y(n-n0)
 LTI systems
 \$\Downline{u}y(n)=h(n)\*x(n)\$

Causal systems

 $\Box$  y(n<sub>0</sub>) depends only on x(n), where n <= n<sub>0</sub>

#### **Pipelining of FIR Digital Filters (1)**

Leads to a reduction in the critical path



The **critical path** (or the minimum time required for processing a new sample) is limited by 1 multiply and 2 add times. Thus, the "sample period" (or the "sample frequency") is given by:

$$T_{sample} \ge T_M + 2T_A$$
$$f_{sample} \le \frac{1}{T_M + 2T_A}$$

#### **Pipelining of FIR Digital Filters (2)**

Pipelining: reduce the effective critical path by introducing latches along the critical data path

FIR: y(n) = 1 a  $\pm x_0(n)x(n) + x_0(n x(n) + 1) + x_0(n x(n) - 2)$  x(n) x(n) x(n-1) a a b b a b c c f c d f

Clock	Input	Node1	Node2	Node3	Output	
0	x(n)	ax(n) + bx(n-1)	ax(n-1) + bx(n-2)	cx(n-3)	y(n-2)	
1	x(n+1)	ax(n+1) + bx(n)	ax(n) + bx(n-1)	cx(n-2)	y(n-1)	
2	x(n+2)	ax(n+2) + bx(n+1)	ax(n+1) + bx(n)	cx(n-1)	<i>y</i> ( <i>n</i> )	
3	x(n + 3)	ax(n+3) + bx(n+2)	ax(n+2) + bx(n+1)	cx(n)	y(n + 1)	

### **Pipelining of FIR Digital Filters (3)**

Pipelining reduces the critical path, but <u>leads to a penalty in</u> terms of an increased <u>latency</u> and the number of <u>latches</u>.

This longest path or the "critical path" can be reduced by suitably placing the pipelining latches in the DSP architecture – The pipelining latches can only be placed across any *feed-forward cutset* of the graph

*Cutset*: a set of edges of a graph such that if these edges are removed from the graph, the graph becomes disjoint

*Feed-forward cutset:* data move in the forward direction on all the edge of the cutset

#### **SFG representation of Cutset**



#### **Fine-grain Pipelining of FIR Filter**



Let  $T_M$  =10 units and  $T_A$  =2 units. If the multiplier is broken into 2 smaller units ( $T_{M1}$  and  $T_{M2}$ ) with processing times of 5 units, respectively (by placing the latches on the horizontal cutset across the multiplier), then the desired clock period can be achieved as 7 units

#### **Retiming and Example (1)**

A transformation technique used to **change the location of delay elements** in the circuit without affecting the input/output characteristics, changing 1) **the clock period** and 2) **the number of registers** 



#### **Retiming and Example (2)**



**Critical Path** :  $T_M + T_A = 3$  u.t Number of registers: 4 Critical Path:  $2T_A = 2 \ u.t$ Number of registers: 5

#### **Retiming Formulation (1)**

Cutset retiming only affects the weights of edges in the cutset.

If the 2 disconnected subgraphs are labeled  $G_1$  and  $G_2$ , then cutset retiming consists of adding k delays to each edge from  $G_1$ to  $G_2$ , and removing k delays to each edge from  $G_2$  to  $G_1$ .



#### **Retiming Formulation (2)**

Map circuit  $G \rightarrow G_r$  using retiming

Retiming can be presented with r(X), X is one of the nodes in the circuit

For an edge  $U \xrightarrow{e} V$  (U is Source, V is the destination)

w'(e) = w(e) + r(V) - r(U)w(e): weight (delay) of the edge  $U \xrightarrow{e} V$  in the origin circuit w'(e): weight of the edge  $U \xrightarrow{e} V$  in the retimed circuit

A retiming solution is feasible if : $w'(e) \ge 0, \forall e \in G$ 

 $w(e) + r(V) - r(U) \ge 0 \Longrightarrow r(U) - r(V) \le w(e)$ 

One inequality per edge due to the causality of the system

#### **Retiming Formulation: cutset contains 1 node**

If cutset = {
$$x_i$$
}: we set  $r(x_i) = \begin{cases} k & if \ x = x_i \\ 0, & others \end{cases}$   $x \text{ are nodes in } G$   
Therefore, for an edge  $U \stackrel{e}{\rightarrow} V$  in  $G$  (U is Source, V is destination)  
 $P = V(V) - r(U) = \begin{cases} k, & if \ V = x_i \\ -k, & if \ U = x_i \\ 0, & others \end{cases}$   
 $w'(e) = w(e) + r(V) - r(U) = \begin{cases} w(e) + k, & if \ V = x_i \\ 0, & others \end{cases}$  In-degree edges of  $x_i$   
 $w(e), & others \end{cases}$  Out-degree edges of  $x_i$ 

#### **Retiming Formulation: cutset contains one node**



**Feasible solution of cutset retiming:**  $-\min(w(e), e \in G_1 \rightarrow G_2) \le k \le \min(w(e), e \in G_2 \rightarrow G_1)$ 

#### **Retiming Formulation: cutset contains many nodes**



#### **Properties of Retiming**

The weight of the retimed path:

 $p = V_0 \xrightarrow{e_0} V_1 \xrightarrow{e_1} \dots \xrightarrow{e_{m-1}} V_m \text{ is given by } w'(p) = w(p) + r(V_m) - r(V_0)$ Proof:  $w'(p) = \sum_{i=0}^{m-1} w'(e_i)$   $= \sum_{i=0}^{m-1} (w(e_i) + r(V_{i+1}) - r(V_i))$   $= \sum_{i=0}^{m-1} w(e_i) + (\sum_{i=0}^{m-1} r(V_{i+1}) - \sum_{i=0}^{m-1} r(V_i))$   $= w(p) + r(V_m) - r(V_0)$ 

- Retiming does not change the number of delays in a cycle
   Retiming does not alter the *iteration bound* in a DFG
- Adding the constant value *j* to the retiming value of each node does not change the mapping from  $G to G_r$

w'(e) = w(e) + (r(V) + j) - (r(U) + j)) = w(e) + r(V) - r(U)

#### **Solving Systems of Inequalities (1)**

 Given a set of M equalities in N variables, use *shortest path algorithm* to solve the results

Step 1: draw a constraint graph

- Draw the node  $x_i$  for each of the N variables  $x_i$ , i=1,2,...,N
- Draw the node  $x_{N+1}$
- For each inequality  $r(x_i) r(x_j) \le k$ , draw the edge  $x_j \to x_i$  from the node  $x_j$  to node  $x_i$  with length k
- For each node  $x_i$ , i=1,2,...n, draw the edge  $x_{N+1} \rightarrow x_i$  from the node  $x_{N+1}$  to the node  $x_i$  with length 0

#### **Solving Systems of Inequalities (2)**

Step 2: Solve using a shortest path algorithm 2024

- The system of inequalities has a solution if and only if the constraint graph contains no negative cycles
- If a solution exists, one solution is where  $r(x_i)$  is the minimum-length path from the node  $x_{N+1}$  to the node  $x_i$

#### Example



Bellman-Ford shortest path algorithm:

$$R^{(6)} = \begin{bmatrix} \infty & \infty & 5 & 4 & \infty \\ 0 & \infty & 2 & 1 & \infty \\ \infty & \infty & \infty & -1 & \infty \\ \infty & \infty & \infty & \infty & \infty \end{bmatrix}$$

$$r_1 = 0, r_2 = 0, r_3 = 0, r_4 = -1$$



#### **Cutset Retiming and Slow-down (Lattice Filter)**

#### **Systolic Transformation**



Source: L. D. Van, "A new 2-D systolic digital filter architecture **without global broadcast**," IEEE Trans. VLSI Systs., vol. 10, pp. 477-486, Aug. 2002

#### The Usage of Retiming Techniques

Cutset retiming and pipelining
Retiming for clock period minimization
Retiming for register minimization

#### **Retiming for clock period minimization (1)**

Minimum feasible clock period or critical path:

$$\Phi(G) = max\{t(\mathbf{p}): U \xrightarrow{\mathbf{p}} V, w(\mathbf{p}) = 0\}$$

Introduce two quantities W(U,V) and D(U,V) to determine if there is a retiming solution that can achieve a desired clock period :  $\square \text{Minimum number of registers of } U \xrightarrow{P} V : W(U,V) = \min\{w(p): U \xrightarrow{P} V\}$   $\square \text{Maximum computation time of } U \xrightarrow{P} V \text{ with minimum number of registers:}$   $D(U,V) = \max\{t(p): U \xrightarrow{P} V \text{ and } w(p) = W(U,V)\}$ 



#### **Retiming for clock period minimization (2)**

Algorithm to compute W(U, V) and D(U, V)

- 1. Let  $M = t_{max}n$ , where  $t_{max}$  is the maximum computation time of the nodes in *G* and *n* is the #of nodes in *G*.
- 2. Form a new graph G' which is the same as G except the edge weights are replaced by w'(e) = Mw(e) - t(U) for all edges  $U \to V$ .
- 3. Solve for all pair shortest path problem on G' by using Floyd*Warshall algorithm*. Let  $S'_{UV}$  be the *shortest path* form  $U \rightarrow V$ .
- 4. If  $U \neq V$ , then  $W(U, V) = [S'_{UV}/M]$  and  $D(U, V) = MW(U, V) - S'_{UV} + t(V)$ . If U = V, then W(U, V) = 0 and D(U, V) = t(U)

#### **Retiming for clock period minimization (3)**





**Step1:**  $M = t_{max}n = 2 \times 4 = 8$  **Step2:** New G' with w'(e) = Mw(e) - t(U)**Step3:** pair shortest path problem on G'

**Step4:** If  $U \neq V$ , then  $W(U,V) = [S'_{UV}/M]$  and  $D(U,V) = MW(U,V) - S'_{UV} + t(V)$ . If U = V, then W(U,V) = 0 and D(U,V) = t(U)

$S'_{UV}$	1	2	3	4	W(U,V)	1	2	3	4	D(U,V)	1	2	3	4
1	12	5	7	15	1	0	1	1	2	1	1	4	3	3
2	7	12	14	22	2	1	0	2	3	2	2	1	4	4
3	5	-2	12	20	3	1	0	0	3	3	4	3	2	6
4	5	-2	12	20	4	1	0	2	0	4	4	3	6	2

#### **Retiming for clock period minimization (4)**

#### Feasibility constraints



#### Retiming for clock period minimization (c=3)

Using W(U, V) and D(U, V) the feasibility and critical path constraints are formulated to give certain inequalities.

The inequalities are solved using *constraint graphs* and if a feasible solution is obtained then the circuit can be clocked with a period 'c'

Feasibility constraints



#### Retiming for clock period minimization (c=2)

Using W(U, V) and D(U, V) the feasibility and critical path constraints are formulated to give certain inequalities.

The inequalities are solved using *constraint graphs* and if a feasible solution is obtained then the circuit can be clocked with a period 'c'

Feasibility constraints



#### **General Approach** (Retiming for clock period minimization)

- Compute W(U, V) and D(U, V)1.
- 2. Sort the values of D(U, V)
- 120:22 Perform *binary search* on these values to find the retiming 3. solution with the minimum clock period. pengiu Rend

#### **Retiming for Register Minimization(1)**

If a node has several output edges carrying the same signal, the number of registers required to implement these edges is the maximum number of registers on any one of the edges.



 $R_U = \max_{\substack{e \\ U \to ?}} \{w_r(e)\} \qquad \qquad R_U = 7$ 

#### **Retiming for Register Minimization(2)**

Minimize  $COST = \sum R_X$  subject to

- (fanout constraint)  $R_X \ge w_r(e)$  for all X and all edges  $X \xrightarrow{e} ?$ 1.
- 2.
- (feasibility constraint)  $r(U) r(V) \le w(e)$  for every edge  $U \xrightarrow{e} V$ (clock period constraint)  $r(U) r(V) \le W(U,V) 1$  for all 3. vertices U, V such that D(U, V) > c.

#### **Retiming and Pipelining**

Pipelining is a special case of cutset retiming where there are no edges in the cutset from G2 to G1, i.e., pipelining applies to graphs without loops

Pipelining is Equivalent to Introducing Many delays at the Input (or output) followed by Retiming



# Next Lecture : Parallel Architecture & Resource Sharing