Embedded Intelligent System and Novel Computer Architecture

Lecture 05 – Systolic Array

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

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

Very Different Architectures Exist



Consider a von Neumann program

- what is the significance of the instruction order?
- what is the significance of the storage locations?

Temporal v.s Spatial Architecture



Systolic Array

$$A[i] = aA^2[i] + b$$

What is the performance (MOPS) of these two Architectures ?



Systolic Computers are a pipelined array architecture, by replace single processor with an array of regular processing elements(PEs), and orchesteate data flow for high throughput with less memory access.

- Synchrony
- Modularity
- Regularity

- Spatial Locality
- Temporal Locality
- Pipelinability
- Parallel Computing

Applications of Systolic Arrays

The various applications of systolic arrays are:

- **1.** Matrix Inversion and Decomposition. UTU 2023
- 2. Polynomial Evaluation.
- 3. Convolution.
- 4. Matrix multiplication.
- 5. Image Processing.
- 6. Systolic lattice filters used for speech and seismic signal processing.
- 7. Artificial neural network
 - **Case Studies:**
 - DFT
 - **1D CONV**
 - Pattern Matching
 - **Top-K elements of an input data-stream**
 - VM and MM Mul

DFT (Discrete Fourier Transform)



DFT (Discrete Fourier Transform)

$$y(x) = \sum_{k=0}^{n-1} a_k x^k |_{x=w^j}$$

A five elements DFT is:

$$y(x = w^j) = a_4(w^j)^4 + a_3(w^j)^3 + a_2(w^j)^2 + a_1(w^j)^1 + a_0$$

$$y(x = w^0) = a_4(w^0)^4 + a_3(w^0)^3 + a_2(w^0)^2 + a_1(w^0)^1 + a_0$$

$$y(x = w^1) = a_4(w^1)^4 + a_3(w^1)^3 + a_2(w^1)^2 + a_1(w^1)^1 + a_0$$

$$y(x = w^2) = a_4(w^2)^4 + a_3(w^2)^3 + a_2(w^2)^2 + a_1(w^2)^1 + a_0$$

$$y(x = w^3) = a_4(w^3)^4 + a_3(w^3)^3 + a_2(w^3)^2 + a_1(w^3)^1 + a_0$$

$$y(x = w^4) = a_4(w^4)^4 + a_3(w^4)^3 + a_2(w^4)^2 + a_1(w^4)^1 + a_0$$

DFT (Discrete Fourier Transform)

$$y(x = w^{0}) = a_{4}(w^{0})^{4} + a_{3}(w^{0})^{3} + a_{2}(w^{0})^{2} + a_{1}(w^{0})^{1} + a_{0}$$

$$y(x = w^{1}) = a_{4}(w^{1})^{4} + a_{3}(w^{1})^{3} + a_{2}(w^{1})^{2} + a_{1}(w^{1})^{1} + a_{0}$$

$$y(x = w^{2}) = a_{4}(w^{2})^{4} + a_{3}(w^{2})^{3} + a_{2}(w^{2})^{2} + a_{1}(w^{2})^{1} + a_{0}$$

$$y(x = w^{3}) = a_{4}(w^{3})^{4} + a_{3}(w^{3})^{3} + a_{2}(w^{3})^{2} + a_{1}(w^{3})^{1} + a_{0}$$

$$y(x = w^{4}) = a_{4}(w^{4})^{4} + a_{3}(w^{4})^{3} + a_{2}(w^{4})^{2} + a_{1}(w^{4})^{1} + a_{0}$$

$$y(w^{0}) = y_{0} = \left(\left((a_{4}(w^{0}) + a_{3})w^{0} + a_{2}\right)w^{0} + a_{1}\right)w^{0} + a_{0}$$

$$y(w^{1}) = y_{1} = \left(\left((a_{4}(w^{1}) + a_{3})w^{1} + a_{2}\right)w^{1} + a_{1}\right)w^{1} + a_{0}$$

$$y(w^{2}) = y_{2} = \left(\left((a_{4}(w^{3}) + a_{3})w^{2} + a_{2}\right)w^{2} + a_{1}\right)w^{2} + a_{0}$$

$$y(w^{3}) = y_{3} = \left(\left((a_{4}(w^{4}) + a_{3})w^{4} + a_{2}\right)w^{3} + a_{1}\right)w^{3} + a_{0}$$

$$y(w^{4}) = y_{4} = \left(\left((a_{4}(w^{4}) + a_{3})w^{4} + a_{2}\right)w^{4} + a_{1}\right)w^{4} + a_{0}$$

DFT Hardware implementation



CONV (Convolution Operation)

Given a Weights Array: $w = (w_1, w_2, ..., w_n)$ and the Inputs Array: $x = (x_1, x_2, ..., x_m)$,

The Output Array: $y = (y_1, y_2, ..., y_{m-n+1})$, Where y_j is as following:

$$y_j = \sum_{i=1}^n w_i x_{i+j-1}, 1 \le j \le m-n-1$$

If n=4, m=6 then: $y = (y_1, y_2, y_3)$, as following:

$$y_{1} = x_{1}w_{1} + x_{2}w_{2} + x_{3}w_{3} + x_{4}w_{4}$$

$$y_{2} = x_{2}w_{1} + x_{3}w_{2} + x_{4}w_{3} + x_{5}w_{4}$$

$$y_{3} = x_{3}w_{1} + x_{4}w_{2} + x_{5}w_{3} + x_{6}w_{4}$$

CONV Hardware Implementation



CONV Hardware Implementation



Combine CONV & DFT Hardware Implementation



Pattern Matching can also using same Arch

Given a Weights Array:
$$w = (w_1, w_2, ..., w_n)$$
CONVand the Inputs Array: $x = (x_1, x_2, ..., x_m)$,...The Output Array: $y = (y_1, y_2, ..., y_{m-n+1})$, Where y_j is as following: $y_j = \sum_{i=1}^n w_i x_{i+j-1}$, $1 \le j \le m - n - 1$ Given a Pattern : $w = (w_1, w_2, ..., w_n)$ Pattern Matchingand the Inputs String: $x = (x_1, x_2, ..., x_m)$,The 1-bit Output Matching Result: $y = (y_1, y_2, ..., y_{m-n+1})$, Where y_j is as following $y_j = \prod_{i=1}^n Comp(w_i, x_{i+j-1})$, $1 \le j \le m - n - 1$

Pattern Matching can also using same Arch



Top K elements of a input queue ?



Top K elements of a input queue ?



Other Examples : Vector Matrix Mul



Other Examples : Matrix Matrix Mul



Systolic Array Design Methodology (1)

- 1. Represent the Algorithm as a Dependence Graph (DG)
- 2. Applying Processor, Scheduling and Projection Vectors (Space-Time Representation)
- 3. Edge Mapping
- 4. Construct the Final Systolic Architecture

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$



Mapping from (i j) axis to (j' t), where j' is location(Processor) and t is time

Systolic Array Design Methodology (2)

- 1. Represent the Algorithm as a Dependence Graph (DG)
- Applying Processor, Scheduling and Projection Vectors (Space-Time Representation) 2.
- 3. Edge Mapping
- 4. Construct the Final Systolic Architecture

ITU 20' Processor Space Vector : $\vec{p}^T = (p_1 p_2)$ Any node with index $I^T = (i, j)$ would be <u>executed by processor</u>: $\vec{p}^T I$

- Scheduling Vector: $\vec{s}^T = (s_1 \ s_2)$ Any node with index $I^T = (i, j)$ would be <u>executed at time</u> : $\vec{s}^T I$
- Projection Vector ; $\vec{d}^T = (d_1 d_2)$

Two nodes that are displaced by \vec{d} or multiples of \vec{d} are <u>executed by the same processor</u> e.g, $\vec{d}^T = (1 \ 0) \rightarrow$ node (1 0) and (2 0) will map to the same PE

Hardware Utilization Efficiency: HUE = $1/|\vec{s}^T d|$.

This is because two tasks executed by the same processor are spaced $|\vec{s}^T d|$ time units apart

Systolic Array Design Methodology (3)

- 1. Represent the Algorithm as a Dependence Graph (DG)
- 2. Applying Processor, Scheduling and Projection Vectors (Space-Time Representation)
- 3. Edge Mapping
- 4. Construct the Final Systolic Architecture

Processor Space Vector and Projection Vector must be orthogonal to each other and $\vec{p}^T \vec{d} = 0$

Proof: I_A and I_B are two computational nodes in the DG

Therefor processor for $I_A(\vec{p}^T I_A)$ and $I_B(\vec{p}^T I_B)$ are same, if $I_A - I_B = \vec{d} \rightarrow \vec{p}^T \vec{d} = 0$

- If A(I_A) and B(I_B) are mapped to the same processor, then they cannot be executed at the same time, i.e., $\vec{s}^T I_A \neq \vec{s}^T I_B$, that is, $\vec{s}^T d \neq 0$.
- Edge mapping: If an edge \vec{e} exists in the space representation or DG, then an edge $\vec{p}^T \vec{e}$ is introduced in the systolic array with $\vec{s}^T \vec{e}$ delays.
- Hardware utilization: $\vec{s}^T d$ = time different between successive nodes (computation) on same processor

e.g,. If $\vec{s}^T d=1$ -> new computation on array nodes with 100% HUE $\vec{s}^T d=2$ -> .. 50% HUE

FIR Filter Design B1 (Broadcast Inputs, Move Results, Weights Stay)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

 $\vec{d}^T = (1 \ \mathbf{0}), \vec{p}^T = (\mathbf{0} \ \mathbf{1}), \vec{s}^T = (\mathbf{1} \ \mathbf{0})$

Any node in the DG with index $I^T = (i, j)$:

- Is mapped to processor $\vec{p}^T I = j$
- Is executed at time $\vec{s}^T I = i$
- Since $\vec{s}^T d = 1$ we have HUE = $1/|\vec{s}^T d| = 100\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

112023

$$\vec{w}: (1 \ 0)$$

$$\vec{y}: (1 \ -1)$$

$$\vec{x}: (0 \ 1)$$

\vec{e}	$\overrightarrow{p}^T \overrightarrow{e}$	$\vec{s}^T \vec{e}$
weight w=(1 0)	0	1
Inputs x̄=(0 1)	1	0
Result	-1	1

FIR Filter Design B1 (Broadcast Inputs, Move Results, Weights Stay)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$



All the horizontal notes execute at the same processor

All the vertical notes execute at the same time $(\vec{s}^T I = i)$

FIR Filter Design B1 (Broadcast Inputs, Move Results, Weights Stay)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$



Alternative Designs

- B1 (Broadcast inputs, Move results, Weight Stay)
- B2 (Broadcast inputs, Move Weight, Results stay)
- F (Fan-in results, Move inputs, Weight stay)
- R1 (Results stay, Inputs and Weight move in opposite directions)
- R2 and Dual R2 (Results stay, Inputs and Weights move in the same direction but at different speeds)
- W1 (Weights stay, Inputs and Results move in opposite directions)
- W2 and Dual W2 (Weights stay, Inputs and Results move in same direction but at different speeds)

FIR Filter Design B2 (Broadcast Inputs, Move Weights, Results Stay)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

 $\vec{d}^T = (1 - 1), \vec{p}^T = (1 \ 1), \vec{s}^T = (1 \ 0)$

11202:3

Any node in the DG with index $I^T = (i, j)$:

- **I**s mapped to processor $\vec{p}^T I_i = i + j$
- Is executed at time $\vec{s}^T I_i = i$
- Since $\vec{s}^T d=1$ we have HUE = $1/|\vec{s}^T d| = 100\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

$$\vec{x} : (0 \ 1) \qquad \vec{y} : (1 \ -1) \qquad \vec{x} : (0 \ 1) \qquad \vec{x} = (0 \ 1) \qquad \vec{x$$

\vec{e}	$\overrightarrow{p}^T \overrightarrow{e}$	$\vec{s}^T \vec{e}$
weight w =(1 0)	1	1
Inputs x=(0 1)	1	0
Result	0	1

FIR Filter Design B2 (Broadcast Inputs, Move Weights, Results Stay)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$



FIR Filter Design F(Move Inputs, Weights Stay, Fan-In Results)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

 $\vec{d}^T = (1 \ 0), \ \vec{p}^T = (0 \ 1), \vec{s}^T = (1 \ 1)$

Any node in the DG with index $I^T = (i, j)$:

- **I**s mapped to processor $\vec{p}^T I_i = j$
- **I**s executed at time $\vec{s}^T I_i = i+j$
- Since $\vec{s}^T d = 1$ we have HUE = $1/|\vec{s}^T d| = 100\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

TU 202:



\vec{e}	$ec{p}^T ec{e}$	$\vec{s}^T \vec{e}$
weight w=(1 0)	0	1
Inputs x =(0 1)	1	1
Result	-1	0

FIR Filter Design F(Move Inputs, Weights Stay, Fan-In Results)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$



FIR Filter Design R1(Results Stay, Inputs and Weights Move in Opposite Direction)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

$$\vec{d}^T = (1 - 1), \ \vec{p}^T = (1 \ 1), \ \vec{s}^T = (1 \ -1)$$

TU 202:

Any node in the DG with index $I^T = (i, j)$:

- **I**s mapped to processor $\vec{p}^T I_i = i + j$
- **I**s executed at time $\vec{s}^T I_i = i j$
- Since $\vec{s}^T d=1$ we have HUE = $1/|\vec{s}^T d| = \frac{1}{2} = 50\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

$$\vec{x} : (0 \ 1) \qquad \vec{y} : (1 \ -1) \qquad \vec{x} : (0 \ 1) \qquad \vec{x} = (0 \ 1) \qquad \vec{x$$

\vec{e}	$ec{p}^T ec{e}$	$\vec{s}^T \vec{e}$
weight w=(1 0)	1	1
Inputs x̄=(0 1)	-1	1
Result	0	2

FIR Filter Design R1(Results Stay, Inputs and Weights Move in Opposite Direction)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$



FIR Filter Design R2 and Dual R2 (Results Stay, Inputs

and Weights Move in Same Direction but at Different Speeds)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

 $\vec{d}^T = (1 - 1), \ \vec{p}^T = (1 \ 1), R2: \ \vec{s}^T = (2 \ 1); \ Dual \ R2: \ \vec{s}^T = (1 \ 2)$

Any node in the DG with index $I^T = (i, j)$:

Is mapped to processor $\vec{p}^T I_i = i + j$

Is executed at time R2: $\vec{s}^T I_i = 2i + j$; Dual R2: $\vec{s}^T I_i = i + 2j$

Since $\vec{s}^T d = 1$ we have HUE = $1/|\vec{s}^T d| = 100\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

$$\vec{w} : (1 - 0)$$

 $\vec{y} : (1 - 1)$
 $\vec{x} : (0 1)$

	R2		Dual R2		
\vec{e}	$\vec{p}^T \vec{e}$	$\vec{s}^T \vec{e}$	$\vec{p}^T \vec{e}$	$\vec{s}^T \vec{e}$	
weight \vec{w} =(1 0)	1	2	1	1	
Inputs \vec{x} =(0 1)	1	1	1	2	
Result \vec{y} =(1 -1)	0	1	0	1	

33

FIR Filter Design W1(Weights Stay, Inputs and Results Move in Opposite Direction)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

 $\vec{d}^T = (1 \ 0), \ \vec{p}^T = (0 \ 1), \vec{s}^T = (2 \ 1)$

Any node in the DG with index $I^T = (i, j)$:

- **I**s mapped to processor $\vec{p}^T I_i = j$
- **I**s executed at time $\vec{s}^T I_i = 2i + j$
- Since $\vec{s}^T d=1$ we have HUE = $1/|\vec{s}^T d| = \frac{1}{2} = 50\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

TU 202:



\vec{e}	$\overrightarrow{p}^T \overrightarrow{e}$	$\vec{s}^T \vec{e}$
weight w =(1 0)	0	2
Inputs x =(0 1)	1	1
Result	-1	1

FIR Filter Design W2 and Dual W2 (Weights Stay, Inputs

and Results Move in Same Direction but at Different Speeds)

FIR filter: $y(n) = w_0 x(n) + w_1 x(n-1) + w_2 x(n-2)$

 $\vec{d}^T = (1 \ 0), \ \vec{p}^T = (0 \ 1), W2: \ \vec{s}^T = (1 \ 2); \ Dual W2: \ \vec{s}^T = (1 \ -1)$

Any node in the DG with index $I^T = (i, j)$:

Is mapped to processor $\vec{p}^T I_i = i$

Is executed at time R2: $\vec{s}^T I_i = i + 2j$; Dual R2: $\vec{s}^T I_i = i - j$

Since $\vec{s}^T d = 1$ we have HUE = $1/|\vec{s}^T d| = 100\%$.

Edge mapping: The 3 fundamental edges corresponding to *weight, input* and *result* can be mapped to corresponding edges in the systolic array as per the following table:

$$\vec{x} : (0 \ 1)$$

	W2		Dual W2	
e	$\vec{p}^T \vec{e}$	$\vec{s}^T \vec{e}$	$\vec{p}^T \vec{e}$	$\vec{s}^T \vec{e}$
weight \vec{w} =(1 0)	0	1	0	1
Inputs $\vec{x} = (0 \ 1)$	1	2	-1	1
Result \vec{y} =(1 -1)	1	1	-1	2

Relationship to Different Transformations

Systolic array architectures with same project vector and processor space vector, but different scheduling vectors can be derived from other transformations



Selection of Schedule Vector

Choose the schedule vector \vec{s} first, then \vec{d} and \vec{p} can be selected according to:

Procedure to get \vec{s} :

Capture all the fundamental edges in *reduced dependence* graph (RDG) constructed by regular iteration algorithm (RIA) description

Construct the scheduling inequalities

Solve feasible \vec{s}^T ($\vec{s}^T \vec{d} = 1$ or $\vec{s}^T \vec{d} = iteration bound$ is optimal)

Scheduling Inequalities

For an edge U to V:

$$U: I_u = \begin{pmatrix} i_u \\ j_u \end{pmatrix} \quad \textcircled{} \qquad \qquad \textcircled{} V: I_v = \begin{pmatrix} i_v \\ j_v \end{pmatrix}$$

The scheduling inequality is:

Linear scheduling:

Affine scheduling

inequality is: $S_{v} \geq S_{u} + T_{u}$ $S_{v} \geq S_{u} + T_{u} = (s_{1} \quad s_{2}) \begin{pmatrix} i_{u} \\ j_{v} \end{pmatrix}$ $S_{u} \equiv S^{T} \quad I_{u} + \gamma_{u} = (s_{1} \quad s_{2}) \begin{pmatrix} i_{u} \\ j_{u} \end{pmatrix} + \gamma_{u}$ $S_{\nu} = S^T I_{\nu} + \gamma_{\nu} = (s_1 \quad s_2) \begin{pmatrix} i_{\nu} \\ i_{\nu} \end{pmatrix} + \gamma_{\nu}$

So the scheduling inequalities:

 $S^T I_n + \gamma_n \geq S^T I_n + \gamma_n + T_n \Rightarrow S^T e_{u \to v} + \gamma_v - \gamma_u \geq T_u$

RDG constructed by RIA

For the FIR filter

w(i + 1, j) = w(i, j)x(i, j + 1) = x(i, j)y(i+1, j-1) = y(i, j) + w(i+1, j-1)x(i+1, j-1)Standard output Regular Iteration Algorithm (RIA) form w(i,j) = w(i-1,j)x(i,j) = x(i,j-1)y(i,j) = y(i-1,j+1) + w(i,j)x(i,j)The Reduced Dependence Graph (RDG): (0,1)(1,0)(0,0) (0,0)

Scheduling Vector Selection

Assume:
$$T_{mult} = 5$$
, $T_{add} = 2$, $T_{trans} = 1$, $s = \begin{pmatrix} s_1 \\ s_2 \end{pmatrix}$



Set
$$\gamma_x = \gamma_w = \gamma_y = 0$$
: $s_1 \ge 1$, $s_2 \ge 1$, $s_1 - s_2 \ge 8$

One solution:
$$s = \begin{pmatrix} 9 \\ 1 \end{pmatrix}$$
, then choose $d = (1, -1), p = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

Systolic Architecture Mapping

$$s = \binom{9}{1}, d = (1, -1), p = \binom{1}{1}$$



Matrix-Matrix multiplication and 2D Systolic Array Design (1)

 $\mathbf{C} = \mathbf{A} \times \mathbf{B} \qquad \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$



$$a(i, j, k) = a(i, j - 1, k)$$

$$b(i, j, k) = b(i - 1, j, k)$$

$$c(i, j, k) = c(i, j, k - 1) + a(i, j, k)b(i, j, k)$$

Matrix-Matrix multiplication and 2D Systolic Array Design (2)



One solution:
$$s = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$
, then choose $d = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$, $p = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix}$

Matrix-Matrix multiplication and 2D Systolic Array Design (S.Y.Kung)

$$\vec{d}^T = (0 \ 0 \ 1), \vec{p}^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \vec{s}^T = (1 \ 1 \ 1)$$



Matrix-Matrix multiplication and 2D Systolic Array Design (S.Y.Kung)

Γ4

4 7

$$\vec{d}^{T} = (1 \ 1 - 1), \vec{p}^{T} = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}, \vec{s}^{T} = (1 \ 1 \ 1)$$

$$\vec{p}^{T} \vec{e} \qquad \vec{p}^{T} \vec{e} \qquad \vec{s}^{T} \vec{e}$$

$$\vec{a} : (0 \ 1 \ 0)^{T} \qquad (0 \ 1)^{T} \qquad 1$$

$$\vec{b} : (1 \ 0 \ 0)^{T} \qquad (1 \ 0)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 0)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

$$\vec{c} : (0 \ 0 \ 1)^{T} \qquad (1 \ 1)^{T} \qquad 1$$

Matrix-Matrix multiplication and 2D Systolic Array Design (S.Y.Kung)

$$\vec{d}^T = (1 \ 1 \ -1), \vec{p}^T = \begin{bmatrix} 1 & -1 & -1 \\ 0 & 1 & -1 \end{bmatrix}, \vec{s}^T = (1 \ 1 \ 1)$$



Multi-projection

Systolic Design for Space Representations with Delays Can be used for multilevel systolic mapping \Box Define N': number of nodes mapped to a processor \Box Iteration period: $N'|S^Td|$

l-th iteration

(l + w) - th iteration

$$U: I_u = \begin{pmatrix} i_u \\ j_u \\ k_u \end{pmatrix} \quad \textcircled{V} \quad w \square \quad \swarrow \quad V: I_v = \begin{pmatrix} i_v \\ j_v \\ k_v \end{pmatrix}$$

Scheduling inequality $S^{T}I_{v} + (l+w)N'|S^{T}d| \ge S^{T}I_{u} + lN'|S^{T}d| + T_{u} \qquad []{} S^{T}e_{u \to v} + w N'|S^{T}d| \ge T_{x}$

All the mapping equations are the same except for the edge delay mapping $\vec{s}^T \vec{e} \Rightarrow \vec{s}^T \vec{e} + w N' |S^T d|$

Example of DG with Delays (1)



Example of DG with Delays (2)



Google TPU-V1 (Systolic Array)



1. H. T. Kung, "Why systolic architectures?", IEEE Computer, 1982.

2. Norman P. Jouppi, et al."In-Datacenter Performance Analysis of a Tensor Processing Unit" ISCA'17 50

Google TPU-V1 (Systolic Array)



- Simplicity reduce the manufacturing Cost & Energy
- When Matrix is mismatching with the size of the underlying MXU, will waste memory bandwidth.

Hardware Utilization

No Direct Support for Sparsity

Summary: Pros and Cons of Systolic Array

Pros:

- Regularity and modular design(Perfect for VLSI implementation)
- Local interconnections
- High degree of pipelining
- Simple I/O subsystem
- High speed and low cost

Cons:

- Global synchronization limits due to signal delays (Improved: Wavefront)
- High bandwidth requirements both for RAM and between PEs
- Poor run-time fault tolerance due to lack of interconnection protocol

Is there a hardware suit for Systolic Array?

YES! Reconfigurable Architecture:

- Spatial parallel computing like 2D-PE array
- Implement hardware structures dynamically
 Connection is programmable

Next Lecture : FPGA & CGRA (Field Programmable Gate Array Coarse Grained Reconfigurable Arch)