# HW4 Systolic Array Mapping & Mapping of Vector Quantization

---

VLSI DSP HW4

Shun-Liang Yeh, NCHU Lab612

4/29 2023

# INDEX

# I. Dependence graph mapping

## Linear schedule Checking



- From the selected Scheduling Vector and Projection vector, check two conditions.
- Check whether the inner product of scheduling vector s and projection vector d is not equal to zero.
- Check whether the inner product of scheduling vector and every edges on DG is equal or greater then zero.

## Projection Procedure



- Select a suitable processor vector P s.t. the inner product of processor vector P and projection vector d is 0.

# Space time transform



**Space-time transform**
- Map an N-dimensional DG index to another N-dimensional space time index
- (1D time index, (N-1)-D space index)
- Transform matrix $T = \begin{bmatrix} s^t \\ \overline{P} \end{bmatrix}$

**Mapping procedures**
- Node mapping, arc mapping, I/O mapping

**Node mapping**

schedule / Processor index $\rightarrow \begin{bmatrix} t(\underline{i}) \\ \underline{n} \end{bmatrix} = \begin{bmatrix} s^t \\ \overline{P} \end{bmatrix} \cdot [\underline{i}] \leftarrow$ DG node

**Arc mapping**
- Maps arcs of DG to the edges of SFG

edge delay / SFG edge $\rightarrow \begin{bmatrix} D(\underline{e}) \\ \underline{e} \end{bmatrix} = \begin{bmatrix} s^t \\ \overline{P} \end{bmatrix} \cdot [\underline{a}] \leftarrow$ DG arc

**I/O mapping**

schedule / Processor index $\rightarrow \begin{bmatrix} t(\underline{c}) \\ \underline{n} \end{bmatrix} = \begin{bmatrix} s^t \\ \overline{P} \end{bmatrix} \cdot [\underline{c}] \leftarrow$ I/O node

- From the derived s , d , P derive a transform matrix T for the N-th dimension DG graph.

- Start mapping every edges on the dependence graph according to the transform matrix T to get its space-time mapping results.

- From the mapping results of N-1 space-time dimension, sketch the edges on this new transformed space according to the mapping results.

- For more step by step algorithmic mapping examples, I suggest watching the lecture by Parhi, the video link is in the references.

# II. Problem & Solution

## Q1

### a)

Q1. For the convolution DG shown in Figure 1, assume each DG node performs a multiply-and-accumulate operation, where $b_i$'s stand for parameters, and $u(\cdot)$'s indicate input samples.

(a) Which of the following sets of scheduling and projection are permissible?

i. $s = [1\ 0]^T$, $d = [1\ 0]^T$
ii. $s = [0\ 1]^T$, $d = [1\ 0]^T$
iii. $s = [1\ 1]^T$, $d = [1\ 0]^T$
iv. $s = [1\ 1]^T$, $d = [0\ 1]^T$



### b)

(b) derive the mapping for each permissible set

Mapped Result

(i), (ii), (iv) — mapped DG diagrams with nodes $b_2$, $b_1$, $b_0$, $u(n)$, $y(n)$, etc.

c)

(c) reverse the direction of data accumulation in the DG, derive a systolic array mapping (all inter-PE data links should have at least one delay element) for it



c)

① Derive the DG.

② find suitable $\vec{s}$, $\vec{d}$ s.t. $\langle \vec{s}, \vec{d} \rangle \neq 0$

let $\vec{s} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$, $d = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

$\langle \vec{s}, d \rangle = 1 \neq 0 \checkmark$

③ Check $\forall e \in DG, \langle \vec{s}, e \rangle > 0$.
$\langle \vec{s}, e \rangle = -2 \times \leq 0$

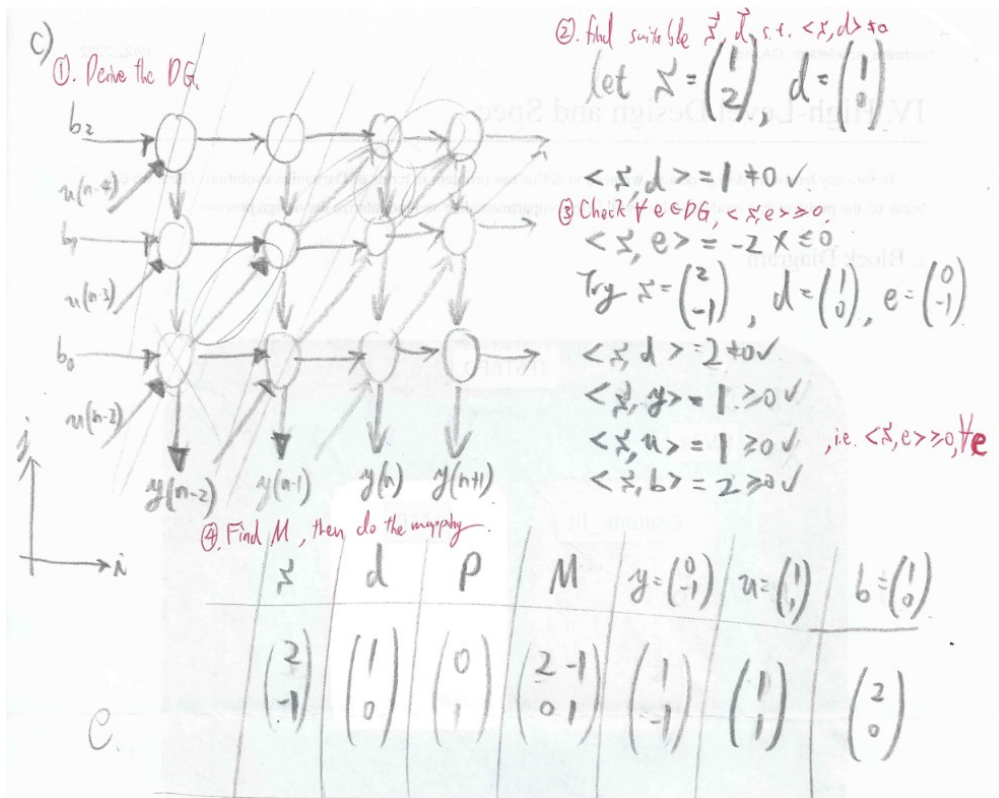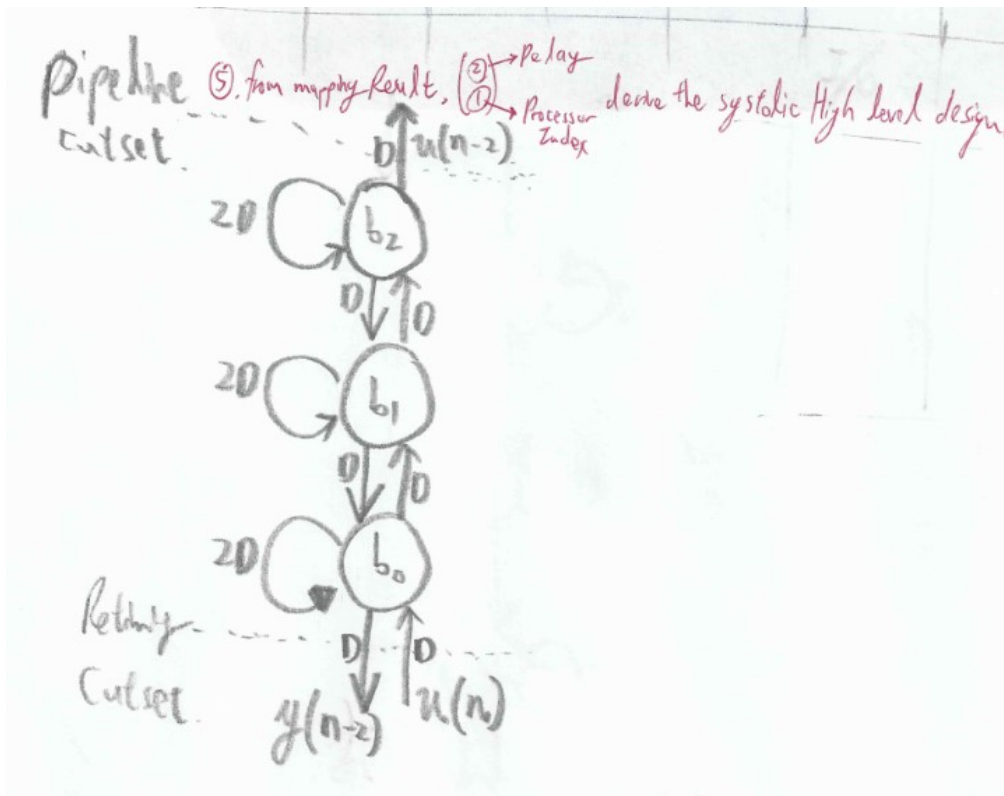Try $\vec{s} = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$, $d = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $e = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$

$\langle \vec{s}, d \rangle = 2 \neq 0 \checkmark$
$\langle \vec{s}, y \rangle = 1 \geq 0 \checkmark$
$\langle \vec{s}, u \rangle = 1 \geq 0 \checkmark$, i.e. $\langle \vec{s}, e \rangle > 0 \checkmark$e
$\langle \vec{s}, b \rangle = 2 \geq 0 \checkmark$

④ Find $M$, then do the mapping.

| $\vec{s}$ | $d$ | $P$ | $M$ | $y = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$ | $u = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ | $b = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ |
|---|---|---|---|---|---|---|
| $\begin{pmatrix} 2 \\ -1 \end{pmatrix}$ | $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ | $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ | $\begin{pmatrix} 2 & -1 \\ 0 & 1 \end{pmatrix}$ | $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$ | $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ | $\begin{pmatrix} 2 \\ 0 \end{pmatrix}$ |

C.

# Q2

**Q2.** Vector quantization design

Given an input $k$-dimensional column vector $\mathbf{r}_{k\times1}$, and a codebook $\mathcal{B}=\{\mathbf{b}_i \mid i=1\sim N\}$ consists of $N$ $k$-dimensional column vectors, vector quantization (VQ) is to find a vector in that codebook that has the shortest Euclidean distance from the input vector $\mathbf{r}_{k\times1}$. The Euclidean distance between two vectors is defined as

$$d'(\mathbf{x},\mathbf{y}) = ||\mathbf{x}-\mathbf{y}|| = \sqrt{\sum_{j=0}^{k-1}(x_j-y_j)^2} \tag{1}$$
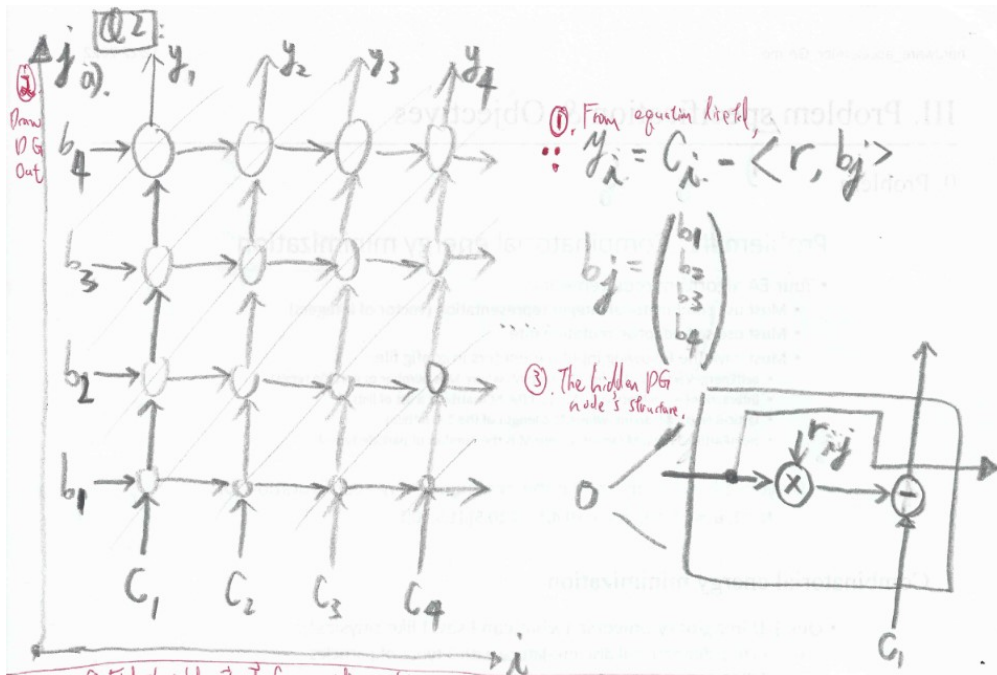
For simplicity, we may use the square distance instead.

$$d(\mathbf{x},\mathbf{y}) = ||\mathbf{x}-\mathbf{y}||^2 = \sum_{j=0}^{k-1}(x_j-y_j)^2 \tag{2}$$

Let $\mathbf{x}=\mathbf{r}$ and $\mathbf{y}=\mathbf{b}_i$, Eq(2) can be rewritten as $d(\mathbf{r},\mathbf{b}_i)=\|\mathbf{r}\|^2-2\mathbf{r}^t\mathbf{b}_i+\|\mathbf{b}_i\|^2$. Since $\|\mathbf{r}\|^2$ is a constant term in all distance calculations, and $\|\mathbf{b}_i\|^2$ can be precomputed, VQ calculation can be expressed as
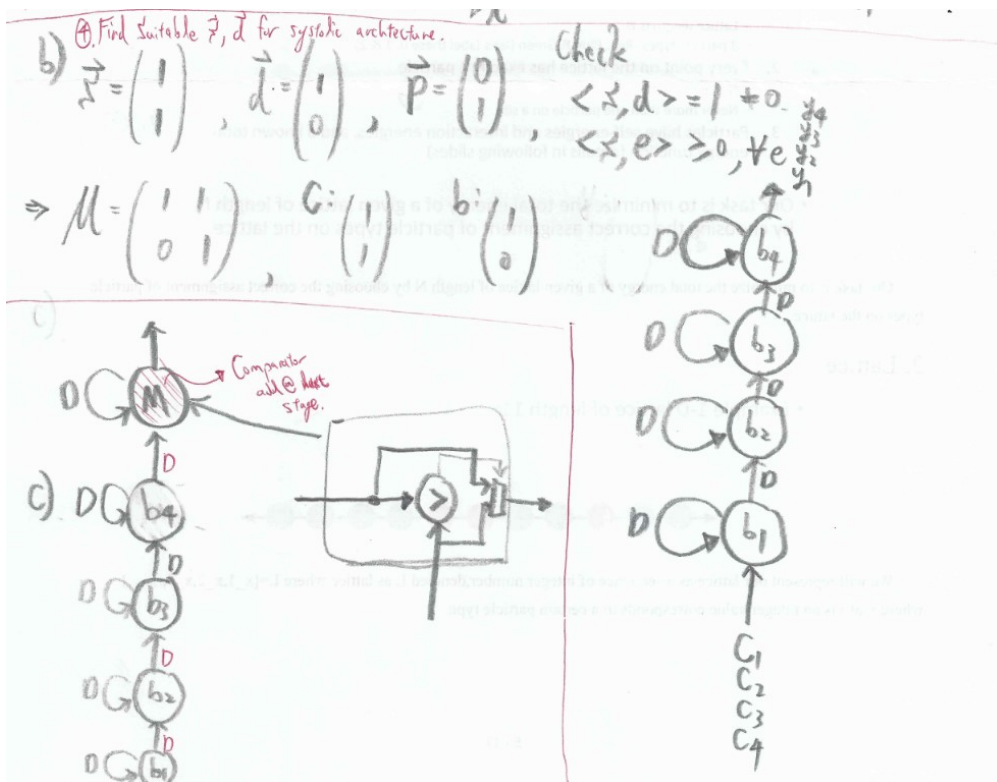
$$\arg\{\min\{c_i-\mathbf{r}^t\mathbf{b}_i \mid i=1,N\}\}, \tag{3}$$

where $c_i=\|\mathbf{b}_i\|^2/2$. In other words, VQ can be accomplished by calculating $c_i-\mathbf{r}^t\mathbf{b}_i$, for all $\mathbf{b}_i$'s in the codebook, and recording the one with the smallest distance. Note that $\mathbf{r}^t\mathbf{b}_i$ is an inner product operation, and $c_i$ can be input from a pre-computed table.

**(a)** Please draw the DG of the $y(i) = c_i - \mathbf{r}^t\mathbf{b}_i$ for $i = 1 \sim N$. For simplicity, assume the vector dimension $k$ is 4. In each iteration, $c_i$ and $\mathbf{b}_i$ are regarded as input and $y(i)$ is the output.



**(b)** Select a scheduling and projection scheme to obtain its systolic array design (at least one delay element in every inter-processor data link).

**(c)** Add a comparator module so that the vector index $i$ corresponding to the minimum Euclidean distance can be obtained at the end of $N$ iterations.

Note

1. Though when deriving the high level architecture for systolic array, only the edges needed to be used, however; the schedule of inputs and outputs are also important during implementation.
2. While transforming from algorithm to DG, beware of the indices, also try to check whether the indices you assigned make sense by going through some examples.

# III. References

[1] VLSI DIGITAL SIGNAL PROCESSING SYSTEMS DESIGN AND IMPLEMENTATION, CH7 Systolic Archtiecture design, Prof. Keshab Parhi, p189~p212

[2] VLSI DIGITAL SIGNAL PROCESSING SYSTEMS DESIGN AND IMPLEMENTATION, CH7 Systolic Archtiecture design, exercises 11, Prof. Keshab Parhi, p215

[3] UMN EE-5329 VLSI Signal Processing Lecture-15 (Spring 2019),Systolic Architecture Design , Prof. Keshab Parhi

[4] UMN EE-5329 VLSI Signal Processing Lecture-16 (Spring 2019),Systolic Architecture Design, Space-Time Mapping , Prof. Keshab Parhi

[5] Computer Architecture - Lecture 27: Systolic Arrays (ETH Zürich, Fall 2020), Prof. Onur Mutlu