

Final Project Proposal

Smith-Waterman for Circular RNA Aligner

賴志杰、麥智詠、王晨峰

2022/06/22

<https://github.com/CHIHCHIEH-LAI/HLS>

Outline

- Smith-Waterman
- a first implementation
- systolic array architecture
- input compression
- shift register
- dual physical ports
- find max value index
- scalable design
- Future Work
 - Fully scalable design

Smith-Waterman Algorithm

設要比對的兩序列為 $A = a_1 a_2 \dots a_n$ 和 $B = b_1 b_2 \dots b_m$ ，其中 n 和 m 分別為序列 A 和 B 的長度。

1. 確定置換矩陣和空位罰分
 - $s(a, b)$ - 組成序列的元素之間的相似性得分
 - W_k - 長度為 k 的空位罰分
2. 建立得分矩陣 H 並初始化其首行和首列。該矩陣的大小為 $n + 1$ 行 $m + 1$ 列（注意計數從 0 開始）。

$$H_{k0} = H_{0l} = 0, (0 \leq k \leq n, 0 \leq l \leq m)$$

3. 從左到右，從上到下進行打分，填充得分矩陣 H 剩餘部分

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ \max_{k \geq 1} \{H_{i-k,j} - W_k\}, \\ \max_{l \geq 1} \{H_{i,j-l} - W_l\}, \\ 0 \end{cases} \quad (1 \leq i \leq n, 1 \leq j \leq m)$$

其中，

$H_{i-1,j-1} + s(a_i, b_j)$ 表示將 a_i 和 b_j 比對的相似性得分，

$H_{i-k,j} - W_k$ 表示 a_i 位於一段長度為 k 的刪除的末端的得分，

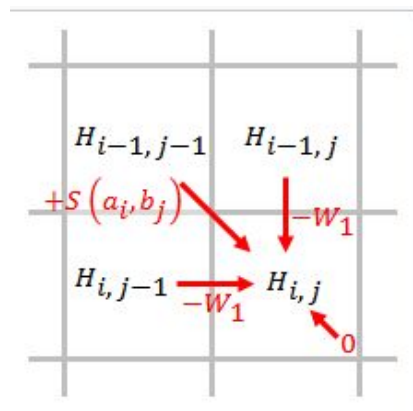
$H_{i,j-l} - W_l$ 表示 b_j 位於一段長度為 l 的刪除的末端的得分，

0 表示 a_i 和 b_j 到此為止無相似性。

4. 回溯。從矩陣 H 中得分最高的元素開始根據得分的來源回溯至上一位置，如此反覆直至遇到得分為 0 的元素。

Smith-Waterman Score Matrix

得分矩陣的作用是對兩序列中兩兩位置進行打分以逐步記錄最佳比對。



$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ H_{i-1,j} - W_1, \\ H_{i,j-1} - W_1, \\ 0 \end{cases} \quad s(a_i, b_j) = \begin{cases} 1, & a_i = b_j \\ -1, & a_i \neq b_j \end{cases}$$

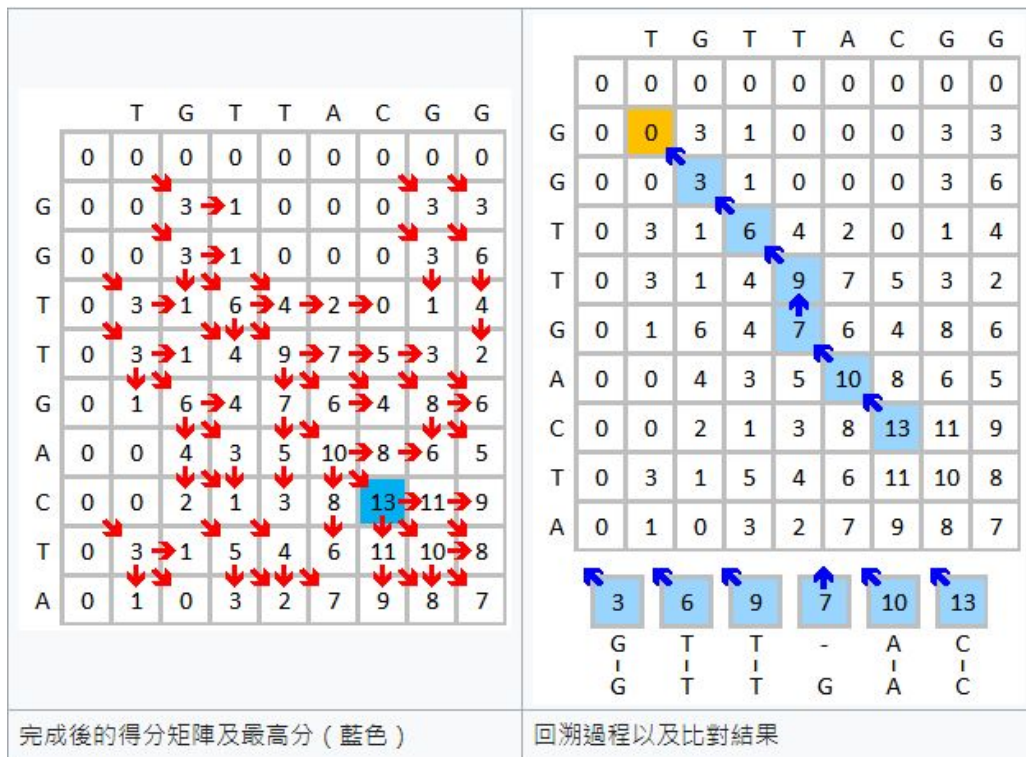
		T	G	T	...
G	0	0	0	0	
G	0				
⋮					

		T	G	T	...
G	0	0	0	0	
G	0	-3	0	-2	
⋮					

		T	G	T	...
G	0	0	0	0	
G	0	0	0	-2	
⋮					

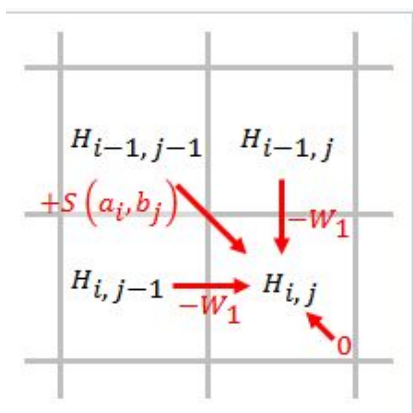
		T	G	T	...
G	0	0	0	0	
G	0	0	0	-2	
⋮					

Smith-Waterman Traceback



constant

```
// scores used for Smith Waterman similarity computation
static const short GAP_i = -1;
static const short GAP_d = -1;
static const short MATCH = 2;
static const short MISS_MATCH = -1;
```



$$H_{ij} = \max \begin{cases} H_{i-1, j-1} + s(a_i, b_j), \\ H_{i-1, j} - W_1, \\ H_{i, j-1} - W_1, \\ 0 \end{cases} \quad s(a_i, b_j) = \begin{cases} 1, & a_i = b_j \\ -1, & a_i \neq b_j \end{cases}$$

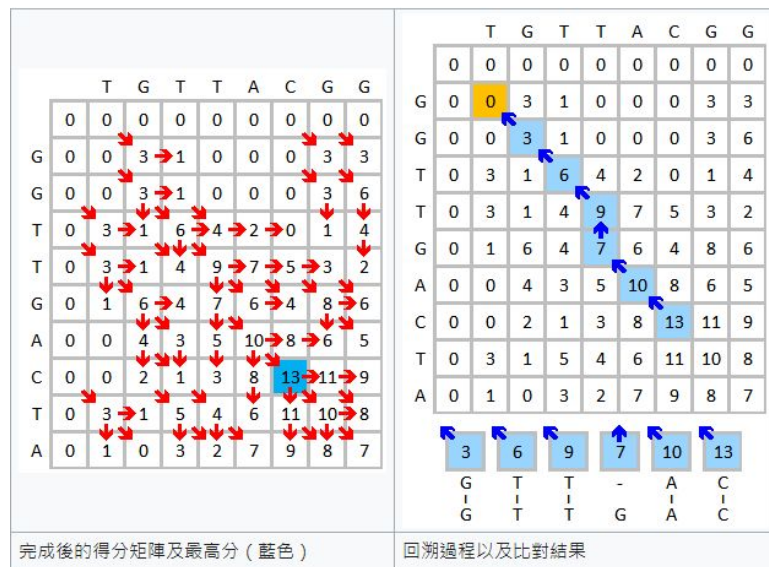
空位罰分: $\text{GAP}_i = \text{GAP}_d = -w_1 = -1$

元素之間的相似得分: $\text{Match} = +S(a_i, b_j) = 2$ or $\text{MISSMATCH} = -S(a_i, b_j) = -1$

constant

```
// directions codes
static const int CENTER = 0;
static const int NORTH = 1;
static const int NORTH_WEST = 2;
static const int WEST = 3;
```

用來記錄得分來源，並儲存結果於direction_matrix



A first implementation

```

20 void compute_matrices(
21     char *string1, char *string2,
22     int *max_index, int *similarity_matrix, short *direction_matrix)
23 {
24     #pragma HLS INTERFACE m_axi port=string1 offset=slave bundle=gmem
25     #pragma HLS INTERFACE m_axi port=string2 offset=slave bundle=gmem
26     #pragma HLS INTERFACE m_axi port=max_index offset=slave bundle=gmem
27     #pragma HLS INTERFACE m_axi port=similarity_matrix offset=slave bundle=gmem
28     #pragma HLS INTERFACE m_axi port=direction_matrix offset=slave bundle=gmem
29
30     #pragma HLS INTERFACE s_axilite port=string1 bundle=control
31     #pragma HLS INTERFACE s_axilite port=string2 bundle=control
32     #pragma HLS INTERFACE s_axilite port=max_index bundle=control
33     #pragma HLS INTERFACE s_axilite port=similarity_matrix bundle=control
34     #pragma HLS INTERFACE s_axilite port=direction_matrix bundle=control
35
36     #pragma HLS INTERFACE s_axilite port=return bundle=control
37
38     //here the real computation starts...
39     int index = 0;
40     int i = 0;
41     int j = 0;
42     short dir = CENTER;
43     short match = 0;
44     int val = 0;
45     int north = 0;
46     int west = 0;
47     int northwest = 0;
48     int max_value = 0;
49     int test_val = 0;
50
51     max_index[0] = 0;
52
53     for(index = N; index < MATRIX_SIZE; index++) {
54         dir = CENTER;
55         val = 0;
56
57         i = index % N; // column index
58         j = index / N; // row index

```

```

59
60     if(i == 0) {
61         // first column
62         west = 0;
63         northwest = 0;
64     } else {
65
66         // all columns but first
67         north = similarity_matrix[index - N];
68         match = ( string1[i] == string2[j] ) ? MATCH : MISS_MATCH;
69
70         test_val = northwest + match;
71         if(test_val > val){
72             val = test_val;
73             dir = NORTH_WEST;
74         }
75
76         test_val = north + GAP_d;
77         if(test_val > val){
78             val = test_val;
79             dir = NORTH;
80         }
81
82         test_val = west + GAP_i;
83         if(test_val > val){
84             val = test_val;
85             dir = WEST;
86         }
87
88         similarity_matrix[index] = val;
89         direction_matrix[index] = dir;
90         west = val;
91         northwest = north;
92         if(val > max_value) {
93             max_index[0] = index;
94             max_value = val;
95         }
96     }
97 }
98
99

```

	T	G	T	T	A	C	G	O	G
G	0	0	3	→ 1	0	0	0	3	3
G	0	0	3	→ 1	0	0	0	3	6
T	0	3	→ 1	6	→ 4	2	→ 0	1	4
T	0	3	→ 1	4	9	→ 7	5	3	2
G	0	1	6	→ 4	7	6	→ 4	8	6
A	0	0	4	3	5	10	→ 8	6	5
C	0	0	2	1	3	8	→ 13	11	9
T	0	3	→ 1	5	4	6	11	10	8
A	0	1	0	3	2	7	9	8	7

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ H_{i-1,j} - W_1, \\ H_{i,j-1} - W_1, \\ 0 \end{cases} \quad s(a_i, b_j) = \begin{cases} 1, & a_i = b_j \\ -1, & a_i \neq b_j \end{cases}$$

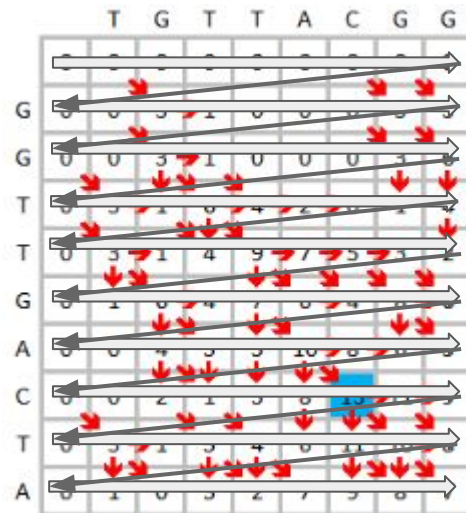
使用直觀演算法不做優化。

A first implementation

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices	II Violation	584786	1.949E6		584787		no	2	~0	0	0	2228	~0	3748	~0	0.00
VITIS_LOOP_53_1	II Violation	584714	1.949E6	148	145	4032	yes									

在最初版本中，每次只計算陣列中的一個值，共需要64x64個cycles，效率極低。

每個cycle中，讀取資料都要透過m-axi向memory讀取，因此需要每次cycle需要大量時間在資料傳輸。

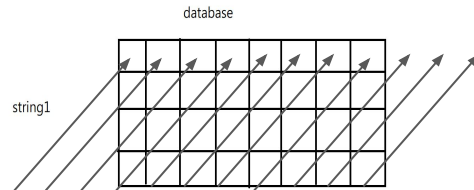


systolic array architecture

```
103     char string1[N];
104     #pragma HLS ARRAY_PARTITION variable=string1 complete dim=1
105     char string2[DATABASE_SIZE];
106     #pragma HLS ARRAY_PARTITION variable=string2 complete dim=1
107
108     // short direction_matrix[DIRECTION_MATRIX_SIZE];
109     // #pragma HLS ARRAY_PARTITION variable=direction_matrix complete dim=1
110
111     memcpy(string1, string1_g, N*sizeof(char));
112     memcpy(string2, string2_g, DATABASE_SIZE * sizeof(char));
113
114     int north[N+1];
115     #pragma HLS ARRAY_PARTITION variable=north complete dim=1
116     int west[N+1];
117     #pragma HLS ARRAY_PARTITION variable=west complete dim=1
118     int northwest[N+1];
119     #pragma HLS ARRAY_PARTITION variable=northwest complete dim=1
120
121     /*
122     short directionDiagonal[N];
123     #pragma HLS ARRAY_PARTITION variable=directionDiagonal complete dim=1
124     */
125
126     ap_uint<512> compressed_diag[1];
127
128     init_dep_for:for(int i = 0; i <= N; i++){
129         north[i] = 0;
130         west[i] = 0;
131         northwest[i] = 0;
132     }
133
134     int directions_index = 0;
135
136     num_diag_for: for(int num_diagonals = 0; num_diagonals < N + M - 1; num_diagonals++){
137         #pragma HLS inline region recursive
138         #pragma HLS PIPELINE
139
140         calculate_diagonal(num_diagonals, string1, string2, northwest, north, west, directions_index, compressed_diag);
141         store_diagonal(directions_index, direction_matrix_g, compressed_diag);
142         directions_index ++;
143     }
```

先將由AXI-master傳輸的query和database儲存於FPGA上的string1與string2。

方向只有3種可能值，因此只需2bit即可表達。



由於反對角線方向的運算彼此不相關，因此可以同時進行。

systolic array architecture

計算三種可能的value

```
33 void calculate_diagonal(int num_diagonals, char string1[N], char string2[DATABASE_SIZE],
34
35     int databaseLocalIndex = num_diagonals;
36     int from, to;
37     from = N * 2 - 2;
38     to = N * 2 - 1;
39
40     calculate_diagonal_for: for(int index = N - 1; index >= 0; index --){
41         int val = 0;
42
43         unsigned int q = string1[index];
44         unsigned int db = string2[databaseLocalIndex];
45
46         // if(num_diagonals < N - 1 && databaseLocalIndex < N - 1 - num_diagonals) db = 9;
47
48         const short match = (q == db) ? MATCH : MISS_MATCH;
49         const short val1 = northwest[index] + match;
50         const short val2 = north[index] + GAP_d;
51         const short val3 = west[index] + GAP_i;
```

kernel寫2 bit值可用range
當Left_variable

why loop starts from N-1?

取三種value最大者

```
53
54
55     if(val1 > val && val1 >= val2 && val1 >= val3){
56         //val1
57         northwest[index + 1] = north[index];
58         north[index] = val1;
59         west[index + 1] = val1;
60         compressed_diag[0].range(to,from) = NORTH_WEST;
61         directionDiagonal[index] = NORTH_WEST;
62     } else if (val2 > val && val2 >= val3) {
63         //val2
64         northwest[index + 1] = north[index];
65         north[index] = val2;
66         west[index + 1] = val2;
67         compressed_diag[0].range(to,from) = NORTH;
68         directionDiagonal[index] = NORTH;
69     } else if (val3 > val){
70         //val3
71         northwest[index + 1] = north[index];
72         north[index] = val3;
73         west[index + 1] = val3;
74         compressed_diag[0].range(to,from) = WEST;
75         directionDiagonal[index] = WEST;
76     } else{
77         //val
78         northwest[index + 1] = north[index];
79         north[index] = val;
80         west[index + 1] = val;
81         compressed_diag[0].range(to,from) = CENTER;
82         directionDiagonal[index] = CENTER;
83     }
84     databaseLocalIndex ++;
85     from -= 2;
86     to -= 2;
87 }
```

systolic array architecture

```
25 void store_diagonal(int directions_index, ap_uint<512> *direction_matrix_g, ap_uint<512> compressed_diag[1]) {  
26  
27     memcpy(direction_matrix_g + directions_index, compressed_diag, sizeof(ap_uint<512>));  
28  
29 }
```

```
30 short get(char data[], int key) {  
31     const int position = (key % 4) * 2;  
32     key /= 4;  
33     char mask = 0;  
34     mask &= 00000000;  
35  
36     mask |= 3 << position;  
37  
38     char fin_mask = 0;  
39     fin_mask &= 00000000;  
40  
41     fin_mask |= 3 << 0;  
42     return ((data[key] & mask) >> (position)) & fin_mask;  
43 }
```

Host 解讀2 bit 編碼的資料需透過 shift與mask運算。

systolic array architecture

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		504	1.680E3		505		no	60	2	0	0	16389	~0	200829	23	0.00
Loop 1		167	557.000	73	1	95	yes									
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		128	427.000	3	1	127	yes									

優化內容：

- 先將由AXI-master傳輸的query和database儲存於FPGA上的string1與string2。
- 反對角線方向的運算 unroll後可以同時執行，最多可以有效計算 N個值。
- 需要回傳的directions經過 2 bit 壓縮編碼，可減少數據傳輸量。

input compression

```
33 void set_char_main(unsigned int * array, int index, unsigned char val){
34     switch(val){
35         case 'A' : array[index / 16] |= (0 << ((index % 16) * 2));
36         break;
37         case 'C' : array[index / 16] |= (1 << ((index % 16) * 2));
38         break;
39         case 'G' : array[index / 16] |= (3 << ((index % 16) * 2));
40         break;
41         case 'T' : array[index / 16] |= (2 << ((index % 16) * 2));
42         break;
43     }
44 }
45 }
```

DNA有ATCG 4種鹼基對, RNA有AUCG 4種鹼基對。
無論分析DNA或RNA皆可只用 2 bit 表示一種鹼基對。

```
105     ap_uint<512> string1[N / NUM_ELEM + 1];
106     // #pragma HLS ARRAY_PARTITION variable=string1 complete dim=1
107     ap_uint<512> string2[(DATABASE_SIZE) / NUM_ELEM + 1];
108     // #pragma HLS ARRAY_PARTITION variable=string2 complete dim=1
109
110     // short direction_matrix[DIRECTION_MATRIX_SIZE];
111     // #pragma HLS ARRAY_PARTITION variable=direction_matrix complete dim=1
112
113     memcpy(string1, string1_g, (N/NUM_ELEM + 1) * 64);
114     memcpy(string2, string2_g, ((DATABASE_SIZE) / NUM_ELEM + 1) * 64);
```

可節省空間與傳輸量

```
44     const short q = string1[index/NUM_ELEM].range((index%NUM_ELEM) * 2 + 1, (index%NUM_ELEM) * 2);
45     short db = string2[databaseLocalIndex/NUM_ELEM].range((databaseLocalIndex%NUM_ELEM) * 2 + 1, (databaseLocalIndex%NUM_ELEM) * 2);
```

在硬體中可以直接獲得指定位置的 bit value,
因此不需要額外時間做解碼。

input compression

優化前 (systolic array):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
▼ ● compute_matrices		504	1.680E3		505		no	60	2	0	0	16389	~0	200829	23	0.00
● Loop 1		167	557.000	73	1	95	yes									
● init_dep_for		65	217.000	1	1	65	yes									
● num_diag_for		128	427.000	3	1	127	yes									

優化後 (systolic array + input compression):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
▼ ● compute_matrices		337	1.123E3		338		no	60	2	0	0	12307	~0	459216	52	0.00
● init_dep_for		65	217.000	1	1	65	yes									
● num_diag_for		129	430.000	4	1	127	yes									

優化內容:

將string1 與 string2 壓縮編碼, 節省空間與傳輸量。

string1與string2壓縮後, 長度各為1個ap_uint<512>, 因此Loop1(memcpy)縮短為2 cycle。

shift register

修改前:

```
139 num_diag_for: for(int num_diagonals = 0; num_diagonals < N + M - 1; num_diagonals++){
140     #pragma HLS inline region recursive
141     #pragma HLS PIPELINE
142
143     calculate_diagonal(num_diagonals, string1, string2, northwest, north, west, directions_index, compressed_diag);
144     store_diagonal(directions_index, direction_matrix_g, compressed_diag);
145     directions_index ++;
146 }

37 int databaseLocalIndex = num_diagonals;
38 int from, to;
39 from = N * 2 - 2;
40 to = N * 2 - 1;
41
42 calculate_diagonal_for: for(int index = N - 1; index >= 0; index --){
43     int val = 0;
44     const short q = string1[index/NUM_ELEM].range((index%NUM_ELEM) * 2 + 1, (index%NUM_ELEM) * 2);
45     short db = string2[databaseLocalIndex/NUM_ELEM].range((databaseLocalIndex%NUM_ELEM) * 2 + 1, (databaseLocalIndex%NUM_ELEM) * 2);
46
84     databaseLocalIndex ++;
85     from -= 2;
86     to -= 2;
87 }
```

原本string2的起始index隨著num_diagnals增加而增加, 因此 db存取的index會隨時間變化。

shift register

修改後:

```
162 num_diag_for: for(int num_diagonals = 0; num_diagonals < N + M - 1; num_diagonals++){
163     #pragma HLS inline region recursive
164     #pragma HLS PIPELINE
165
166     calculate_diagonal(num_diagonals, string1, string2, northwest, north, west, directions_index, compressed_diag, shift_db, direction_matrix_g);
167
168     update_database(string2, shift_db, num_diagonals);
169     // store_diagonal(directions_index, direction_matrix_g, compressed_diag);
170     // directions_index ++;
171 }

51 int databaseLocalIndex = 0;
52 int from, to;
53 from = N * 2 - 2;
54 to = N * 2 - 1;
55
56 calculate_diagonal_for: for(int index = N - 1; index >= 0; index --){
57     int val = 0;
58     const short q = string1[index/NUM_ELEM].range((index%NUM_ELEM) * 2 + 1, (index%NUM_ELEM) * 2);
59     unsigned int db = shift_db[databaseLocalIndex/NUM_ELEM].range((databaseLocalIndex % NUM_ELEM) * 2 + 1, (databaseLocalIndex % NUM_ELEM) * 2);

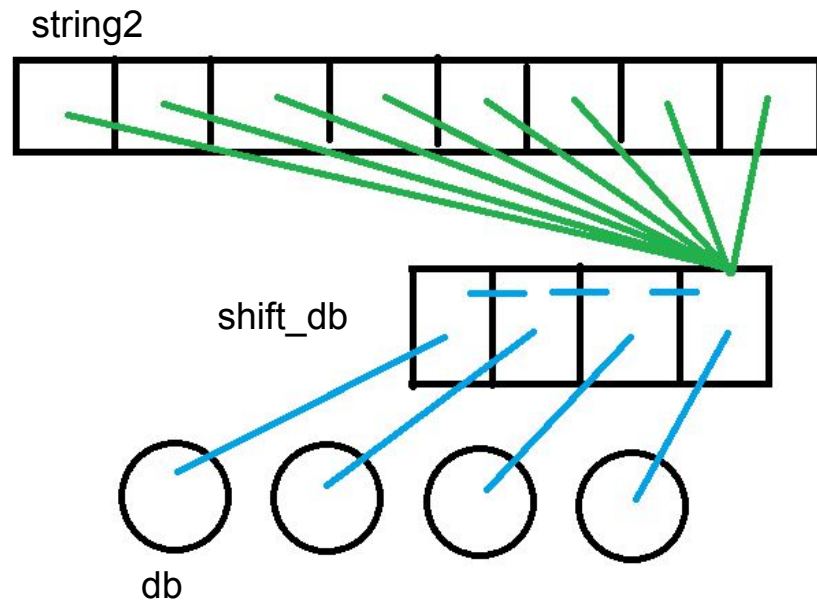
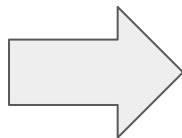
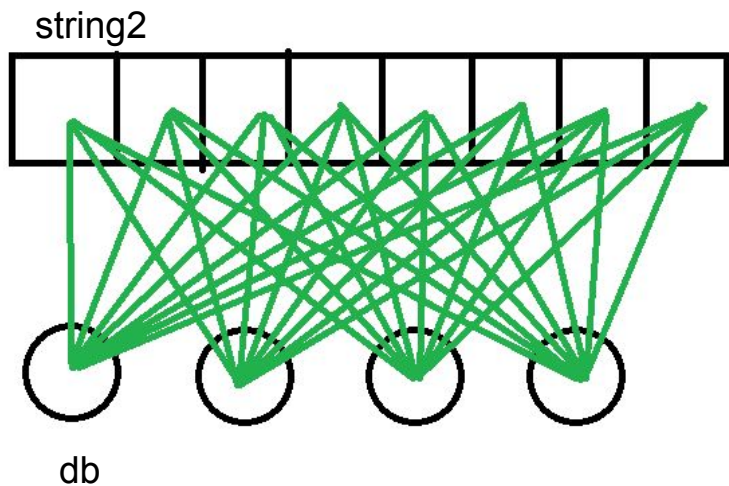
33 void update_database(ap_uint<512> *database, ap_uint<512> *shift_db, int num_diagonals){
34     int startingIndex = N + num_diagonals;
35     update_database:for(int i = 1; i < N; i++){
36         #pragma HLS PIPELINE
37         shift_db[(i-1)/NUM_ELEM].range(((i-1)%NUM_ELEM)*2+1, ((i-1)%NUM_ELEM)*2) = shift_db[i/NUM_ELEM].range((i%NUM_ELEM)*2+1, (i%NUM_ELEM)*2);
38         //set_char(shift_db, i-1, get_char(databaseLocal, i));
39         //databaseLocal[i-1] = databaseLocal[i];
40     }
41     shift_db[(N-1)/NUM_ELEM].range(((N-1)%NUM_ELEM) * 2 + 1, ((N-1)%NUM_ELEM) * 2) = database[startingIndex/NUM_ELEM].range((startingIndex%NUM_ELEM) * 2 + 1, (startingIndex%NUM_ELEM) * 2);
42     //set_2bit(shift_db, N-1, get_2bit(database, startingIndex));
43     //set_char(shift_db, N-1, get_char(database, startingIndex));
44     //databaseLocal[N-1] = database[startingIndex];
45 }
```

string2透過shift_db做媒介, 起始index固定在0。
存取的index固定, 不隨時間變化。

shift_db作為一個滑動的窗口。

shift register

data傳輸網路圖



可看到數據之間的連線數大幅減少！

shift register

優化前 (systolic array + input compression):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		337	1.123E3		338		no	60	2	0	0	12307	~0	459216	52	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		129	430.000	4	1	127	yes									

優化後 (systolic array + input compression + shift register):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices	II Violation	25354	8.451E4		25355		no	60	2	0	0	19376	1	211583	24	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		25146	8.381E4	198		127	no									
calculate_diagonal_for		66	220.000	3	1	64	yes									
update_database	II Violation	126	420.000	2	2	63	yes									

可看到總Latency大幅增加！

update_database出現II violation (原先認為此處應該 unroll)

shift register

num_diag_for	
calculate_diagonal_for	
update_database	II Violation

```
33 void update_database(ap_uint<512> *database, ap_uint<512> *shift_db, int num_diagonals){
34     int startingIndex = N + num_diagonals;
35     update_database:for(int i = 1; i < N; i++){
36         #pragma HLS PIPELINE
37         shift_db[(i-1)/NUM_ELEM].range(((i-1)%NUM_ELEM)*2+1, ((i-1)%NUM_ELEM)*2) = shift_db[i/NUM_ELEM].range((i%NUM_ELEM)*2+1, (i%NUM_ELEM)*2);
38         //set_char(shift_db, i-1, get_char(databaseLocal, i));
39         //databaseLocal[i-1] = databaseLocal[i];
40     }
41     shift_db[(N-1)/NUM_ELEM].range(((N-1)%NUM_ELEM) * 2 + 1, ((N-1)%NUM_ELEM) * 2) = database[startingIndex/NUM_ELEM].range((startingIndex%NUM_ELEM) * 2 +1, (startingIndex%
42     //set_2bit(shift_db, N-1, get_2bit(database, startingIndex));
43     //set_char(shift_db, N-1, get_char(database, startingIndex));
44     //databaseLocal[N-1] = database[startingIndex];
45 }
```

原因:

在update_database裡面出現 #pragma HLS PIPELINE，因此判斷是此 pragma 誤導 Vitis HLS。因為上層 num_diag_for 已經切 pipeline，因此會盡量讓下層的 update_database 平行或 unroll，但是 update_database 下面多了 #pragma HLS PIPELINE，使 update_database 的 for loop 無法被 unroll，造成 II violation。

解決方法:

我們將其註解調並重新合成。

shift register

優化前 (systolic array + input compression):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		337	1.123E3		338		no	60	2	0	0	12307	~0	459216	52	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		129	430.000	4	1	127	yes									

優化後 (systolic array + input compression + shift register):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		337	1.123E3		338		no	60	2	0	0	13706	~0	118187	13	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		129	430.000	4	1	127	yes									

可發現雖然Latency與先前完全相同，但因為 shift register很容易於硬體實現，也簡化了資料傳輸，因此資源消耗量大幅降低。

shift register

嘗試不同硬體規模：

原配置：M與N分別為128與256，此大小需要較多資源且合成極慢（需要數個小時）。

```
18 #define N 256
19 #define M 128
20 #define NUM_ELEM 256
21 #define DATABASE_SIZE (M + 2 * (N - 1))
22 #define DIRECTION_MATRIX_SIZE ((N + M - 1) * N)
23 #define MATRIX_SIZE (N * M)
```

若只修改M與N的定義，則其電路無法正確運作。

經分析後，發現code其他部分僅能在 $N = 256$ 與 $M = 128$ 的條件下成立。

因此我們將code改寫能夠適應各種N與M的配置。

dual physical ports

```
110 void compute_matrices( ap_uint<512> *string1_g, ap_uint<512> *string2_g, ap_uint<512> *direction_matrix_g)
111 {
112 #pragma HLS INTERFACE m_axi port=string1_g offset=slave bundle=gmem0
113 #pragma HLS INTERFACE m_axi port=string2_g offset=slave bundle=gmem0
114 #pragma HLS INTERFACE m_axi port=direction_matrix_g offset=slave bundle=gmem1
```



```
110 void compute_matrices( ap_uint<512> *string1_g, ap_uint<512> *string2_g, ap_uint<512> *direction_matrix_g)
111 {
112 #pragma HLS INTERFACE m_axi port=string1_g offset=slave bundle=gmem0
113 #pragma HLS INTERFACE m_axi port=string2_g offset=slave bundle=gmem1
114 #pragma HLS INTERFACE m_axi port=direction_matrix_g offset=slave bundle=gmem2
```

```
130 memcpy(string1, string1_g, (N/NUM_ELEM + 1) * 64);
131 memcpy(string2, string2_g, ((DATABASE_SIZE) / NUM_ELEM + 1) * 64);
```

原本string1_g與string2_g為同一個bundle，因此資料依序傳輸。
將string1_g跟string2_g設為不同的bundle以實現dual physical ports，同時傳輸兩者資料。

dual physical ports

優化前 (systolic array + input compression + shift register):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		337	1.123E3		338		no	60	2	0	0	13706	~0	118187	13	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		129	430.000	4	1	127	yes									

優化後 (systolic array + input compression + shift register + dual physical ports):

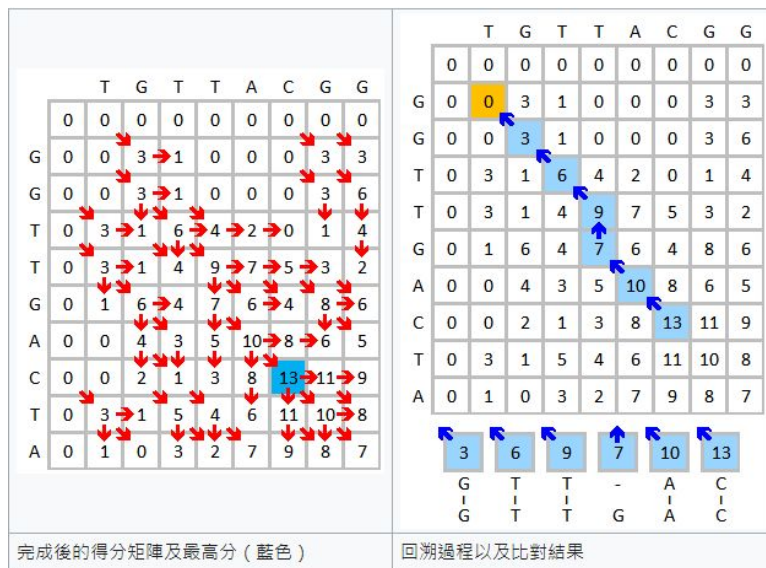
Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		336	1.120E3		337		no	90	3	0	0	15120	~0	119775	13	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		129	430.000	4	1	127	yes									

dual ports讓string1_g和string2_g可以同時傳，在本次實驗中 string1_g與string2_g長度皆為1個 ap_uint<512>，原先需要2個cycle，優化後僅需要1個cycle。

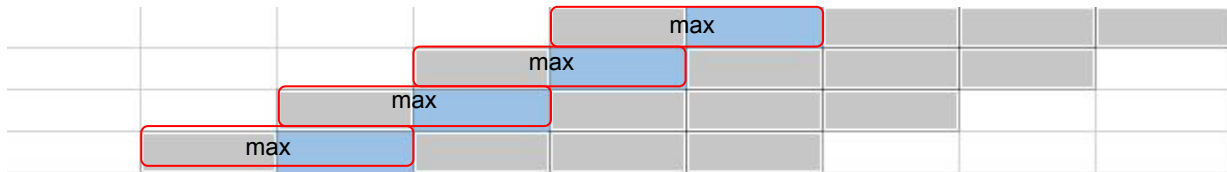
雖然看似只省去僅僅一個 cycle，但在做scalable design時，string1_g和string2_g如果很大，整體省下的 cycles還是很可觀的。

find max value index

在做smith-waterman時，不僅需要算direction matrix，還要算出哪裡有最大值，這樣才有辦法從最大值的那個點進行backtracking。



find max value index



```
249 num_diag_for: for(int num_diagonals = 0; num_diagonals < N + host_M - 1; num_diagonals++){
250 #pragma HLS inline region recursive
251 #pragma HLS PIPELINE
252
253     calculate_diagonal(num_diagonals, string1, string2, northwest, north, west, directions_index, compressed_diag, shift_db, direction_matrix_g);
254
255     compare(num_diagonals, north, max_index, max_value);
256
257     update_database(string2, shift_db, num_diagonals);
258     // store_diagonal(directions_index, direction_matrix_g, compressed_diag);
259     // directions_index++;
260     if(rest_size > 0){
261         if(num_diagonals - (N-1) - tail >= NUM_ELEM){
262             memcpy(string2 + replace_idx%3, string2_g + replace_idx, 64);
263             tail += NUM_ELEM;
264             replace_idx++;
265             rest_size--;
266         }
267     }
268
269     find_max(max_index, max_value);
270     memcpy(max_value_idx, max_index, sizeof(int));
271 }
```

```
void compare(int num_diagonals, int north[N + 1], int max_index[N + 1], int max_value[N + 1]){
    for(int i=0; i<N; i++){
        if(north[i] > max_value[i]){
            max_index[i] = num_diagonals;
            max_value[i] = north[i];
        }
    }
}
```

針對每條對角線，我們會對此對角線與前一條對角線上N個值各別去做同位置的比大小，使得max_index都會存所有目前對角線上N個最大值的index。

find max value index

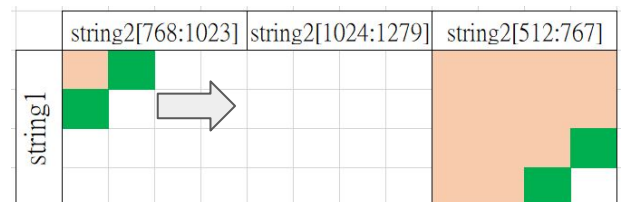
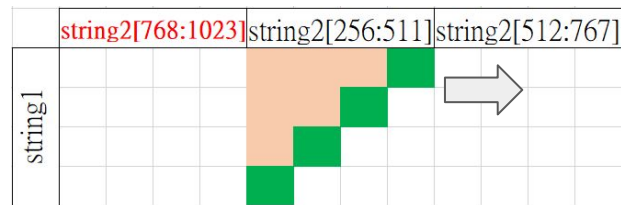
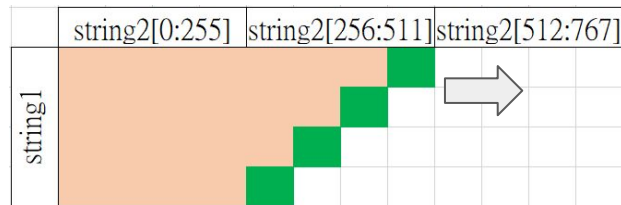
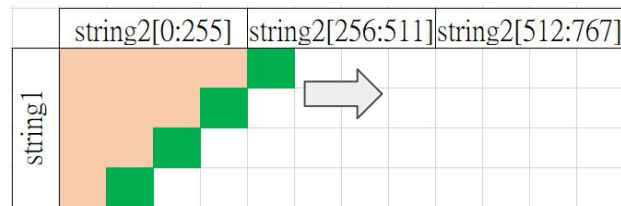
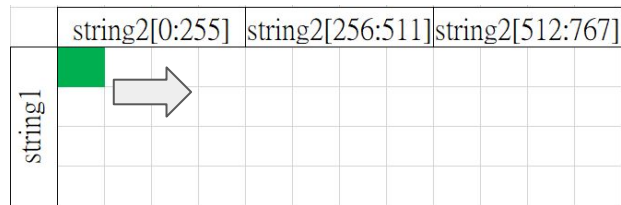
```
117 void find_max(int indexs[N + 1], int values[N + 1]){
118     #pragma HLS PIPELINE
119
120     int max_index_64[64], max_index_32[32], max_index_16[16], max_index_8[8], max_index_4[4], max_index_2[2];
121     int max_value_32[32], max_value_16[16], max_value_8[8], max_value_4[4], max_value_2[2];
122
123     for(int i=0;i<N;i++){
124         max_index_64[i] = indexs[i] * N + i;
125     }
126
127     for(int i = 0; i<32; i++){
128         if(values[i*2] >= values[i*2+1]){
129             max_index_32[i] = max_index_64[i*2];
130             max_value_32[i] = values[i*2];
131         }else{
132             max_index_32[i] = max_index_64[i*2+1];
133             max_value_32[i] = values[i*2+1];
134         }
135     }
136
137     for(int i = 0; i<16; i++){
138         if(max_value_32[i*2] >= max_value_32[i*2+1]){
139             max_index_16[i] = max_index_32[i*2];
140             max_value_16[i] = max_value_32[i*2];
141         }else{
142             max_index_16[i] = max_index_32[i*2+1];
143             max_value_16[i] = max_value_32[i*2+1];
144         }
145     }
146
147     for(int i = 0; i<8; i++){
148         if(max_value_16[i*2] >= max_value_16[i*2+1]){
149             max_index_8[i] = max_index_16[i*2];
150             max_value_8[i] = max_value_16[i*2];
151         }else{
152             max_index_8[i] = max_index_16[i*2+1];
153             max_value_8[i] = max_value_16[i*2+1];
154         }
155     }
156
157     for(int i = 0; i<4; i++){
158         if(max_value_8[i*2] >= max_value_8[i*2+1]){
159             max_index_4[i] = max_index_8[i*2];
160             max_value_4[i] = max_value_8[i*2];
161         }else{
162             max_index_4[i] = max_index_8[i*2+1];
163             max_value_4[i] = max_value_8[i*2+1];
```

計算完所有對角線後，我們會採用max tree，搭配 pipeline，分成6級運算，以比較N個值中的最大值，並取得max value index，這樣傳回host就可以做 backtracking。

scalable design

```
245 int rest_size = (host_DATABASE_SIZE+NUM_ELEM-1)/NUM_ELEM - 3;
246 int replace_idx = 3;
247 int tail = 0;
248
249 num_diag_for: for(int num_diagonals = 0; num_diagonals < N + host_M - 1; num_diagonals++){
250 #pragma HLS inline region recursive
251 #pragma HLS PIPELINE
252
253     calculate_diagonal(num_diagonals, string1, string2, northwest, north, west, directions_
254
255     compare(num_diagonals, north, max_index, max_value);
256
257     update_database(string2, shift_db, num_diagonals);
258     store_diagonal(directions_index, direction_matrix_g, compressed_diag);
259     directions_index++;
260
261     if(rest_size > 0){
262         if(num_diagonals - (N-1) - tail >= NUM_ELEM){
263             memcpy(string2 + replace_idx*3, string2_g + replace_idx, 64);
264             tail += NUM_ELEM;
265             replace_idx++;
266             rest_size--;
267         }
268     }
```

原本硬體的kernel_N與kernel_M需與string1與string2的長度一致。
而在真實的smithwaterman應用中，序列長度不固定。
我們透過實作circular buffer，將string2分段讀取，即可在固定的硬體資源運算任意長度的string2 (database)。



final version

優化前 (systolic array + input compression + shift register + dual physical ports):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices		336	1.120E3		337		no	90	3	0	0	15120	~0	119775	13	0.00
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for		129	430.000	4	1	127	yes									

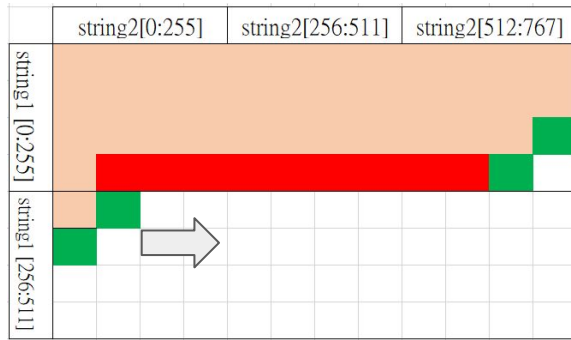
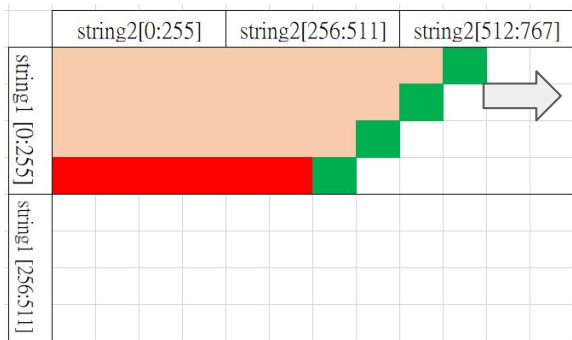
優化後 (systolic array + input compression + shift register + dual physical ports + find max + scalable design):

Name	Issue Type	Latency (cycles)	Latency (ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Slack
compute_matrices							no	62	2	0	0	39165	2	181938	20	0.00
find_max_1		2	6.666		1		yes	0	0	0	0	4920	~0	7282	~0	
Loop 1		4	13.332	2	1	3	yes									
init_dep_for		65	217.000	1	1	65	yes									
num_diag_for				76	1		yes									

雖然scalable design比之前的design多出一倍多的運算資源，但是scalable design可以在不更改硬體的情況下，應付不同的database長度，且可算出max index。

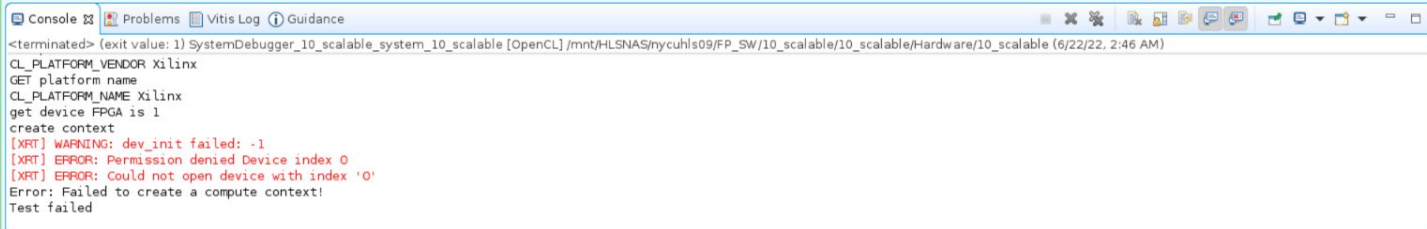
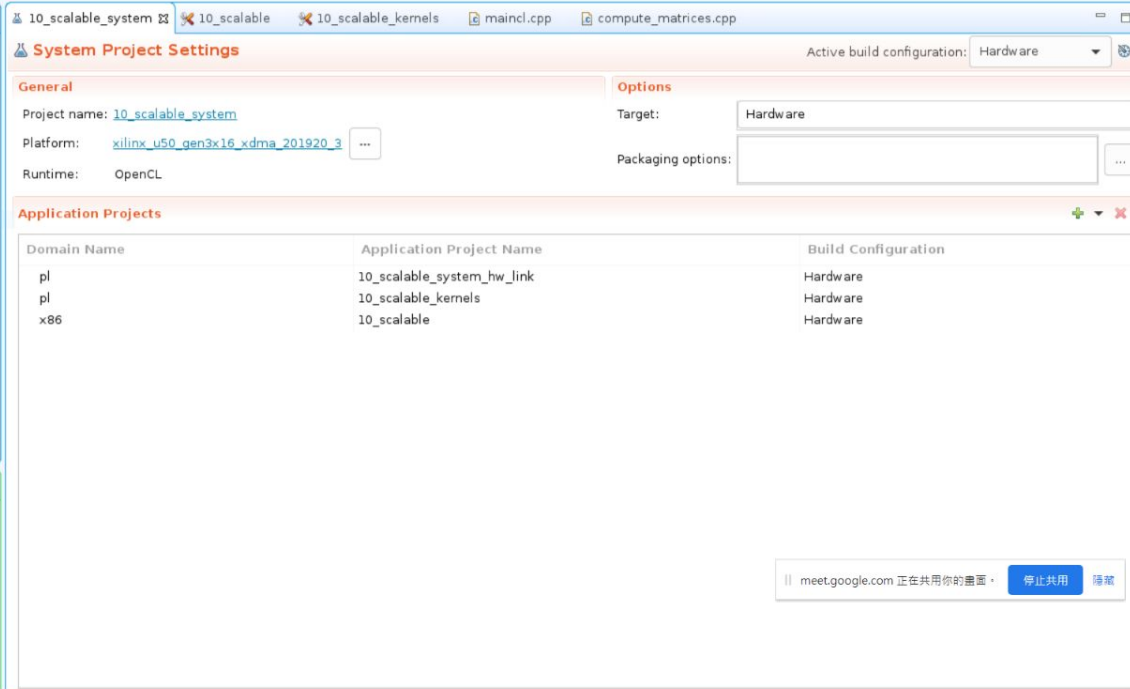
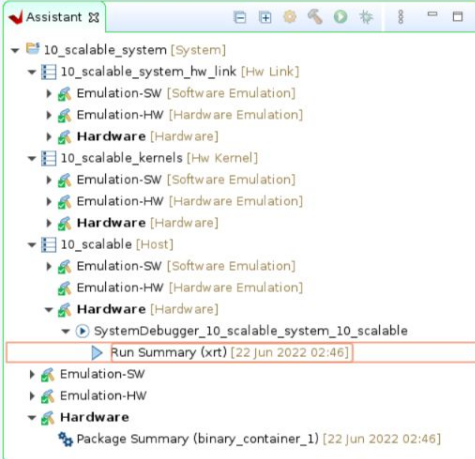
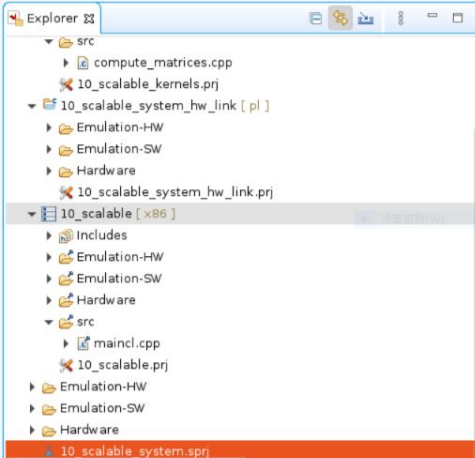
Future Work

Fully scalable design:



需要使用到長度 = $M - \text{NCU}$ 的buffer儲存分割處的value。

($M = \text{string2 length}$, $\text{NCU} = \text{number of compute unit}$)



Build Console [fp1_kernels, Emulation-HW]

INFO: [v++ 60-1306] Additional information associated with this v++ compile can be found at:

Reports: /home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/reports/compute_matrices

Log files: /home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/logs/compute_matrices

Running Dispatch Server on port: 40707

INFO: [v++ 60-1548] Creating build summary session with primary output /home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/compute_matrices.xo.compile_summary, at Tue Jun 21 17:45:13 2022

INFO: [v++ 60-1316] Initiating connection to rulecheck server, at Tue Jun 21 17:45:13 2022

INFO: [v++ 60-895] Target platform: /opt/xilinx/platforms/xilinx_u2_gen3x4_xdma_gc_2_202110_1/xilinx_u2_gen3x4_xdma_gc_2_202110_1.xpfm

INFO: [v++ 60-1578] This platform contains Xilinx Shell Archive '/opt/xilinx/platforms/xilinx_u2_gen3x4_xdma_gc_2_202110_1/hw/hw.xsa'

INFO: [v++ 74-78] Compiler Version string: 2021.2

INFO: [v++ 60-585] Compiling for hardware emulation target

INFO: [v++ 60-423] Target device: xilinx_u2_gen3x4_xdma_gc_2_202110_1

INFO: [v++ 60-242] Creating kernel: 'compute_matrices'

====>The following messages were generated while performing high-level synthesis for kernel: compute_matrices Log file: /home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/compute_matrices/compute_matrices/vitis_hls.log

INFO: [v++ 204-61] Pipelining loop 'Loop 1'.

INFO: [v++ 200-1470] Pipelining result : Target II = NA, Final II = 1, Depth = 2, loop 'Loop 1'

INFO: [v++ 204-61] Pipelining loop 'init_dep_for'.

INFO: [v++ 200-1470] Pipelining result : Target II = NA, Final II = 1, Depth = 1, loop 'init_dep_for'

INFO: [v++ 204-61] Pipelining loop 'num_diag_for'.

INFO: [v++ 200-1470] Pipelining result : Target II = NA, Final II = 1, Depth = 139, loop 'num_diag_for'

INFO: [v++ 204-61] Pipelining function 'find_max'.

INFO: [v++ 200-1470] Pipelining result : Target II = NA, Final II = 1, Depth = 6, function 'find_max'

INFO: [v++ 200-790] **** Loop Constraint Status: All loop constraints were satisfied.

INFO: [v++ 200-789] **** Estimated Fmax: 402.17 MHz

ERROR: [v++ 213-28] Failed to generate IP.

ERROR: [v++ 60-300] Failed to build kernel(ip) compute_matrices, see log for details: /home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/compute_matrices/compute_matrices/vitis_hls.log

ERROR: [v++ 60-773] In '/home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/compute_matrices/compute_matrices/vitis_hls.log', caught Tcl error: ERROR: '2206211747' is an invalid argument. Please specify an integer

ERROR: [v++ 60-773] In '/home/dcs41_2/fp/fp1_kernels/Emulation-HW/build/compute_matrices/compute_matrices/vitis_hls.log', caught Tcl error: ERROR: [IMPL 213-28] Failed to generate IP.

ERROR: [v++ 60-599] Kernel compilation failed to complete

ERROR: [v++ 60-592] Failed to finish compilation

INFO: [v++ 60-1653] Closing dispatch client.

makefile:56: recipe for target 'build/compute_matrices.xo' failed

make: *** [build/compute_matrices.xo] Error 1

17:47:29 Build Finished (took 2m:28s.724ms)