



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ

ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ  
ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΥΠΟΛΟΓΙΣΤΩΝ

## «1η Εργαστηριακή Άσκηση»

Εργαστηριακή Αναφορά στο μάθημα

### Σχεδιασμός Ενσωματωμένων Συστημάτων

των φοιτητών

Αβραμίδης Σταύρος Α.Μ. 17811

Λάζου Μαρία-Αργυρώ Α.Μ. 20192

Διδάσκοντες:

Σωτήριος Ξύδης, Δημήτριος Σούντρης, Σωτήριος Κοκόσης

Νοέμβριος 2024

# Μέρος 1ο

## Εισαγωγή

Στο 1ο μέρος της άσκησης ζητήθηκε η βελτιστοποίηση ενός αλγορίθμου Parallel Hierarchical One Dimensional Search (PHODS). Ο αλγόριθμος PHODS είναι ένας αλγόριθμος εκτίμησης κίνησης (Motion Estimation), ο οποίος έχει στόχο να ανιχνεύσει τη κίνηση των αντικειμένων μεταξύ δύο διαδοχικών εικόνων (frame) του βίντεο. Ο αλγόριθμος βελτιστοποιήθηκε για το S7G2, έναν μικροελεγκτή της Renesas.

## Ζητούμενα

Μέσω του Datasheet και SpecSheet του dev-board εξάγουμε τις παρακάτω πληροφορίες:<sup>1</sup>

- 640KB SRAM
- 4MB Flash
- 240MHz CPU clock frequency
- Cortex-M4 CPU
- Flash Cache
  - ▶ Flash cache 1 - Instruction cache: 256 bytes
  - ▶ Flash cache 2 - CPU operands: 16 bytes
  - ▶ Prefetch Buffer: 32 bytes

Από την σελίδα 1763 του «*Groups User manual*» για τον εν λόγω μικροελεγκτή, μπορούμε να δούμε ότι η FCACHE είναι απενεργοποιημένη. Πάραυτα, το BSP ενεργοποιεί την cache, όταν καλείται η ρουτίνα του Reset Handle, κατά την αρχικοποίηση των ρολογιών (bps\_clocks.c).

## Υλοποίηση

Για τον σκοπό της άσκησης, η σημαία βελτιστοποιήσεων του μεταγλωττιστή τέθηκε σε «-O0». Για τον λόγο αυτό, πέρα των loop transformation, που είναι ο κύριος σκοπός της άσκησης, έγιναν και μερικές περαιτέρω βελτιστοποιήσεις, που ανώτερα επίπεδα βελτιστοποίησης θα μπορούσαν να εξαλείψουν. Την ίδια στιγμή, όμως, αποφεύχθηκαν βελτιστοποιήσεις, που θα άλλαζαν την βασική δομή του αλγορίθμου (λ.χ. μετατροπή int σε uint8\_t), θεωρώντας τες εκτός του πλαισίου της άσκησης.

Αναλυτικότερα στον δοσμένο κώδικα, έγιναν οι εξής βελτιστοποιήσεις:

1. Loop Merging (x3) στις των loop i, k, l. (Ακόλουθο αυτού, sub-expression elimination του p1)
2. Loop Merging των 2 εξωτερικών loop (i,j) και (x,y) του κύριου αλγορίθμου και τις αρχικοποίησης των vectors\_x και vectors\_y.
3. Απαλοιφή των πολλαπλών memory accesses των vectors\_x[x][y] και vectors\_y[x][y].
4. Loop Unrolling του loop i από -S έως S+1. \*
  - Αυτό έχει ως αποτέλεσμα την απλοποίηση της περίπτωσης για i=0, αφού οι συνθήκες των if/else είναι κοινές και κατα συνέπεια  $p2 = q2$  και κατ'επέκταση  $p1-q2 = p1-p2$ .

---

<sup>1</sup>«<https://www.renesas.com/en/products/microcontrollers-microprocessors/renesas-synergy-platform-mcus/s7g2-240-mhz-arm-cortex-m4-cpu#documents>»

5. SubExpression Elimination για όλες τις πράξεις που γίνονται πάνω από μια φορά εντός των εσωτερικών των loop (i, k, l).
6. Υλοποίηση των πολλαπλασιασμών στα x,y loop, ως running sum, αντί για πολλαπλασιασμό κάθε φορά.
7. Τέλος, παρατηρούμε ότι η εντολές των if/else που προσανξάνουν τα distx και disty, πραγματοποιούν πρακτικά μια πράξη abs. Για τον λόγο αυτό, η πράξη της αφαίρεσης και αφαίρεσης του προσήμου υλοποιήθηκε με βέλτιστο τρόπο σε assembly. (Συνοπτική επεξήγηση, ακολουθεί παρακάτω).

\* Δεν εξετάστηκε, αν αντί για unroll, ήταν ευνοϊκότερη αντ' αυτού η χρήση if-else για την περίπτωση του i=0, με την προϋπόθεση ότι το loop body θα χωράει στην instructions cache.

**Η ορθότητα του αλγορίθμου επαληθεύτηκε με την χρήση του randomized unit testing, στον προσωπικό μας υπολογιστή.**

Ο τελικός κώδικας, μετά τις βελτιστοποιήσεις, είναι ο εξής:

```
phods.h

1  #include "inttypes.h"
2
3  #define N (10) /* Frame dimension for QCIF format */
4  #define M (10) /* Frame dimension for QCIF format */
5  #define B (5) /* Block size*/
6  #define p
7  (7) /* Search space. Restricted in a [-p,p] region around the original \
8      location of the block. */ \
9
10
11 __attribute__((always_inline))
12 inline int32_t sub_and_abs(int32_t a, int32_t b) {
13     int32_t result;
14     __asm__ ("subs %0, %1, %2\n" // result = a - b, sets flags
15             "it mi\n" // If result is negative (MI), execute next instruction
16             "negmi %0, %0\n" // If negative, result = -result
17             : "=r"(result) // output
18             : "r"(a), "r"(b) // inputs
19             : "cc" // clobbers the condition codes
20     );
21     return result;
22 }
23
24 void phods_motion_estimation(const int current[N][M], const int previous[N][M],
25                             int vectors_x[N / B][M / B],
26                             int vectors_y[N / B][M / B]) {
27     int x, y, i, k, l, p1, p2, q2, distx, disty, S, min1, min2, bestx, besty;
28
29     distx = 0;
30     disty = 0;
31
32     /*For all blocks in the current frame*/
33     int B_times_x = 0;
34
35     for (x = 0; x < N / B; x++) {
36
37         int B_times_y = 0;
38
39         for (y = 0; y < M / B; y++) {
40
41             /*Initialize the vector motion matrices*/
42             vectors_x[x][y] = 0;
43             vectors_y[x][y] = 0;
44
45             S = 4;
46
47             while (S > 0) {
48                 min1 = 255 * B * B;
49                 min2 = 255 * B * B;
50
51             }
```

```

51      /*For all candidate blocks in X and Y dimension*/
52
53      const int vec_x_xy = vectors_x[x][y];
54      const int vec_y_xy = vectors_y[x][y];
55      const int B_times_x_plus_vec_x_xy = B_times_x + vec_x_xy;
56      const int B_times_y_plus_vec_y_xy = B_times_y + vec_y_xy;
57
58      {
59          const int i = -S;
60          distx = 0;
61          disty = 0;
62
63          const int B_times_x_plus_vec_x_xy_plus_i =
64              B_times_x_plus_vec_x_xy + i;
65          const int B_times_y_plus_vec_y_xy_plus_i =
66              B_times_y_plus_vec_y_xy + i;
67
68          /*For all pixels in the block*/
69          for (k = 0; k < B; k++) {
70              const int B_times_x_plus_vec_x_xy_plus_i_plus_k =
71                  B_times_x_plus_vec_x_xy_plus_i + k;
72
73              const int B_times_x_plus_vec_x_xy_plus_k =
74                  B_times_x_plus_vec_x_xy + k;
75
76              const int B_times_x_plus_k = B_times_x + k;
77
78              for (l = 0; l < B; l++) {
79                  const int B_times_y_plus_vec_y_xy_plus_i_plus_l =
80                      B_times_y_plus_vec_y_xy_plus_i + l;
81
82                  const int B_times_y_plus_vec_y_xy_plus_l =
83                      B_times_y_plus_vec_y_xy + l;
84
85                  p1 = current[B_times_x_plus_k][B_times_y + l];
86
87                  if ((B_times_x_plus_vec_x_xy_plus_i_plus_k) < 0 ||
88                      (B_times_x_plus_vec_x_xy_plus_i_plus_k) > (N - 1) ||
89                      (B_times_y_plus_vec_y_xy_plus_l) < 0 ||
90                      (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
91                      p2 = 0;
92                  } else {
93                      p2 = previous[B_times_x_plus_vec_x_xy_plus_i + k]
94                          [B_times_y_plus_vec_y_xy_plus_l];
95                  }
96
97                  if ((B_times_x_plus_vec_x_xy_plus_k) < 0 ||
98                      (B_times_x_plus_vec_x_xy_plus_k) > (N - 1) ||
99                      (B_times_y_plus_vec_y_xy_plus_i_plus_l) < 0 ||
100                      (B_times_y_plus_vec_y_xy_plus_i_plus_l) > (M - 1)) {
101                      q2 = 0;
102                  } else {
103                      q2 = previous[B_times_x_plus_vec_x_xy_plus_k]
104                          [B_times_y_plus_vec_y_xy_plus_i_plus_l];
105                  }
106
107                  distx += sub_and_abs(p1, p2);
108                  disty += sub_and_abs(p1, q2);
109              }
110          }
111
112          if (distx < min1) {
113              min1 = distx;
114              bestx = i;
115          }
116
117          if (disty < min2) {
118              min2 = disty;
119              besty = i;
120          }
121      }
122
123      {
124          const int i = 0;
125
126          distx = 0;
127          disty = 0;
128
129          /*For all pixels in the block*/
130          for (k = 0; k < B; k++) {
131              const int B_times_x_plus_vec_x_xy_plus_k =

```

```

132     B_times_x_plus_vec_x_xy + k;
133
134     const int B_times_x_plus_k = B_times_x + k;
135
136     for (l = 0; l < B; l++) {
137         const int B_times_y_plus_vec_y_xy_plus_l =
138             B_times_y_plus_vec_y_xy + l;
139
140         p1 = current[B_times_x_plus_k][B_times_y + l];
141
142         if ((B_times_x_plus_vec_x_xy_plus_k) < 0 ||
143             (B_times_x_plus_vec_x_xy_plus_k) > (N - 1) ||
144             (B_times_y_plus_vec_y_xy_plus_l) < 0 ||
145             (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
146             p2 = 0;
147         } else {
148             p2 = previous[B_times_x_plus_vec_x_xy_plus_k]
149                 [B_times_y_plus_vec_y_xy_plus_l];
150         }
151
152         int res = sub_and_abs(p1, p2);
153         distx += res;
154         disty += res;
155     }
156 }
157
158 if (distx < min1) {
159     min1 = distx;
160     bestx = i;
161 }
162
163 if (disty < min2) {
164     min2 = disty;
165     besty = i;
166 }
167 }
168
169 {
170     const int i = S;
171     distx = 0;
172     disty = 0;
173
174     const int B_times_x_plus_vec_x_xy_plus_i =
175         B_times_x_plus_vec_x_xy + i;
176     const int B_times_y_plus_vec_y_xy_plus_i =
177         B_times_y_plus_vec_y_xy + i;
178
179     /*For all pixels in the block*/
180     for (k = 0; k < B; k++) {
181         const int B_times_x_plus_vec_x_xy_plus_i_plus_k =
182             B_times_x_plus_vec_x_xy_plus_i + k;
183
184         const int B_times_x_plus_vec_x_xy_plus_k =
185             B_times_x_plus_vec_x_xy + k;
186
187         const int B_times_x_plus_k = B_times_x + k;
188
189         for (l = 0; l < B; l++) {
190             const int B_times_y_plus_vec_y_xy_plus_i_plus_l =
191                 B_times_y_plus_vec_y_xy_plus_i + l;
192
193             const int B_times_y_plus_vec_y_xy_plus_l =
194                 B_times_y_plus_vec_y_xy + l;
195
196             p1 = current[B_times_x_plus_k][B_times_y + l];
197
198             if ((B_times_x_plus_vec_x_xy_plus_i_plus_k) < 0 ||
199                 (B_times_x_plus_vec_x_xy_plus_i_plus_k) > (N - 1) ||
200                 (B_times_y_plus_vec_y_xy_plus_l) < 0 ||
201                 (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
202                 p2 = 0;
203             } else {
204                 p2 = previous[B_times_x_plus_vec_x_xy_plus_i + k]
205                     [B_times_y_plus_vec_y_xy_plus_l];
206             }
207
208             if (B_times_x_plus_vec_x_xy_plus_k < 0 ||
209                 B_times_x_plus_vec_x_xy_plus_k > (N - 1) ||
210                 (B_times_y_plus_vec_y_xy_plus_i_plus_l) < 0 ||
211                 (B_times_y_plus_vec_y_xy_plus_i_plus_l) > (M - 1)) {
212                 q2 = 0;

```

```

213         } else {
214             q2 = previous[B_times_x_plus_vec_x_xy_plus_k]
215                 [B_times_y_plus_vec_y_xy_plus_i_plus_l];
216         }
217
218         distx += sub_and_abs(p1, p2);
219         disty += sub_and_abs(p1, q2);
220     }
221 }
222
223 if (distx < min1) {
224     min1 = distx;
225     bestx = i;
226 }
227
228 if (disty < min2) {
229     min2 = disty;
230     besty = i;
231 }
232 }
233
234 S = S / 2;
235 vectors_x[x][y] += bestx;
236 vectors_y[x][y] += besty;
237 }
238
239 // Reduce B * y to addition
240 B_times_y += B;
241 }
242
243 // Reduce B * x
244 B_times_x += B;
245 }
246 }

```

## hal\_entry.c

```

1  #include "stdio.h"
2  #include "string.h"
3
4  #include "common_utils.h"
5  #include "hal_data.h"
6
7
8  #include "image.h"
9  #include "phods.h"
10
11 void hal_entry(void) {
12
13     // Initialize your variables here
14     int motion_vectors_x[N / B][M / B], motion_vectors_y[N / B][M / B], i, j;
15
16     uint32_t timer;
17
18     // Code to initialize the DWT->CYCCNT register
19     CoreDebug->DEMCR |= 0x01000000;
20     ITM->LAR = 0xC5ACCE55;
21     DWT->CYCCNT = 0;
22     DWT->CTRL |= 1;
23
24     // Initial
25
26     APP_PRINT("Embedded Systems - Lab1\n\n");
27     APP_PRINT("Omada 24 - Maria Lazou, Avramidis Stavros\n\n");
28
29     APP_PRINT("System Running @%u MHz\n\n", SystemCoreClock / 1000000);
30
31     APP_PRINT("Running for B: %d\n", Bx, By);
32
33     for (int i = 0; i < 10; i++) {
34         timer = DWT->CYCCNT;
35         phods_motion_estimation_naive(current, previous, motion_vectors_x,
36                                     motion_vectors_y);
37         timer = DWT->CYCCNT - timer;
38
39         APP_PRINT("Time for naive loop[%d]: %u\n", i, timer);
40     }
41 }

```

```

42  APP_PRINT("\n");
43
44  for (int i = 0; i < 10; i++) {
45      timer = DWT->CYCCNT;
46      phods_motion_estimation_opt(current, previous, motion_vectors_x,
47                                  motion_vectors_y);
48      timer = DWT->CYCCNT - timer;
49
50      APP_PRINT("Time for opt loop[%d]: %u\n", i, timer);
51  }
52
53  while (1)
54      ;
55  }

```

Όσον αφορά την υλοποίηση της πράξης της αφαίρεσης με βέλτιστο τρόπο σε assembly, έγινε η ακόλουθη υλοποίηση:

```

SUBS %0, %1, %2
TI
NEGMI %0, %0

```

Αρχικά, υλοποιούμε την αφαίρεση, η οποία ενημερώνει το αντίστοιχο flag, αν το αποτέλεσμα της είναι αρνητικό. Κατόπιν εισάγουμε ένα IT (If/Then) block<sup>2</sup>, που θα εκτελεστεί μόνο αν η συνθήκη του Negative flag, έχει τεθεί. Έτσι, μπορούμε να κάνουμε τις πράξεις της αφαίρεσης και abs(), με μόνο 2 εντολές κι είναι και branchless, μη επιβαρύνοντας έτσι τον branch predictor.

## Χρόνοι εκτέλεσης

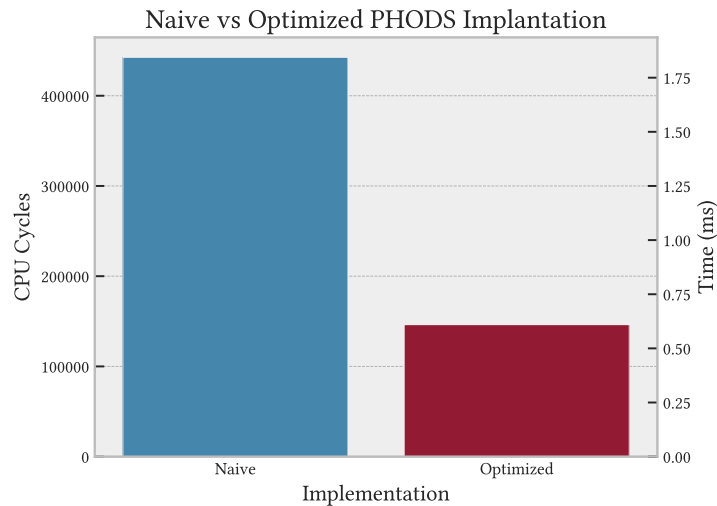
Ο αρχικός χρόνος εκτέλεσης του αλγορίθμου ήταν 442,530 κύκλοι ή 1.84ms. Μετά τις βελτιστοποιήσεις, ο χρόνος εκτέλεσης μειώθηκε στους 146,283 κύκλους ή 0.61ms.

Υλοποίηση	Κύκλοι	Χρόνος Εκτέλεσης (ms)
Naive	442530	1.84
Optimized	146283	0.61

Πίνακας 1: Χρόνοι εκτέλεσης πριν και μετά τις βελτιστοποιήσεις

<sup>2</sup>To IT block (If-Then block) στους επεξεργαστές ARM είναι ένα χαρακτηριστικό του instruction set της αρχιτεκτονικής ARM Thumb-2, που επιτρέπει την εκτέλεση υπό συνθήκη μιας ομάδας εντολών

<https://developer.arm.com/documentation/den0042/a/Unified-Assembly-Language-Instructions/Instruction-set-basics/Conditional-execution>



Σχήμα 1: Σύγκριση χρόνων εκτέλεσης

Καταλήγουμε λοιπόν να έχουν Speedup = 3.025.

Ενδιαφέρον επίσης παρουσιάζει η απενεργοποίηση της FCACHE, κι η σύγκριση του χρόνου του αλγορίθμου. Με την επίδοση να πέφτει στους 184285 κύκλους (26% Αύξηση!).

## Μεταβολή του Block Size

Στην συνέχεια ζητήθηκε, να συγκριθεί και να αναλυθεί η επίδραση του Block Size (B) στον χρόνο εκτέλεσης του αλγορίθμου.

B	Κύκλοι	Χρόνος Εκτέλεσης (ms)
1	273050	1.14
2	179422	0.75
5	146283	0.61
10	139770	0.58

Πίνακας 2: Συγκριτικά αποτελέσματα ανά Block Size

Παρατηρούμε ότι όσο μεγαλύτερο είναι το Block Size, τόσο μειώνεται ο χρόνος εκτέλεσης του αλγορίθμου. Αυτό είναι λογικό, αφού στον αλγόριθμο έχουμε 2 ομάδες loop, τα εξωτερικά(x,y) που είναι αντιστρόφως ανάλογα του B και τα εσωτερικά (k,l), που είναι ανάλογα του B. Ο συνολικός αριθμός επαναλήψεων συνεπώς μένει ίδιος, όμως όσο το φόρτο των υπολογισμών μετατοπίζεται προς το εσωτερικό, τόσο τα memory accesses, αλλά και διάφοροι υπολογισμοί μειώνονται. Την ίδια στιγμή δίνεται η ευκαιρία να «χτυπήσουμε» περισσότερο την instruction cache, αφού ο αλγόριθμός περνάει μεγαλύτερο χρονικό διάστημα στον ίδιο κώδικα.



## Ορθογώνιο Παράθυρο

Στην συνέχεια, ζητήθηκε η υλοποίηση του αλγορίθμου με ορθογώνιο παράθυρο, διαστάσεων  $B_x, B_y$ . Η υλοποίηση τροποποιείται όπως επισημαίνεται παρακάτω:

### phods\_rect.h

```
1  #include "inttypes.h"
2
3  #define N (10) /* Frame dimension for QCIF format */
4  #define M (10) /* Frame dimension for QCIF format */
5  #define B (5) /* Block size*/
6  #define p
7  (7) /* Search space. Restricted in a [-p,p] region around the original \
8      location of the block. */ \
9
10 #define Bx (5)
11 #define By (5)
12
13 __attribute__((always_inline))
14 inline int32_t sub_and_abs(int32_t a, int32_t b) {
15     int32_t result;
16     __asm__ ("subs %0, %1, %2\n" // result = a - b, sets flags
17             "it mi\n" // If result is negative (MI), execute next instruction
18             "negmi %0, %0\n" // If negative, result = -result
19             : "=r"(result) // output
20             : "r"(a), "r"(b) // inputs
21             : "cc" // clobbers the condition codes
22     );
23     return result;
24 }
25
26 void phods_motion_estimation_opt(const int current[N][M],
27                                 const int previous[N][M],
28                                 int vectors_x[N / B][M / B],
29                                 int vectors_y[N / B][M / B]) {
30     int x, y, i, k, l, p1, p2, q2, distx, disty, S, min1, min2, bestx, besty;
31
32     distx = 0;
33     disty = 0;
34
35     /*For all blocks in the current frame*/
36     int B_times_x = 0;
37
38     for (x = 0; x < N / Bx; x++) {
39
40         int B_times_y = 0;
41
42         for (y = 0; y < M / By; y++) {
43
44             /*Initialize the vector motion matrices*/
45             vectors_x[x][y] = 0;
46             vectors_y[x][y] = 0;
47
48             S = 4;
49
50             while (S > 0) {
51                 min1 = 255 * B * B;
52                 min2 = 255 * B * B;
53
54                 /*For all candidate blocks in X and Y dimension*/
55
56                 const int vec_x_xy = vectors_x[x][y];
57                 const int vec_y_xy = vectors_y[x][y];
58                 const int B_times_x_plus_vec_x_xy = B_times_x + vec_x_xy;
59                 const int B_times_y_plus_vec_y_xy = B_times_y + vec_y_xy;
60
61                 {
62                     const int i = -S;
63                     distx = 0;
64                     disty = 0;
65
66                     const int B_times_x_plus_vec_x_xy_plus_i =
67                         B_times_x_plus_vec_x_xy + i;
68                     const int B_times_y_plus_vec_y_xy_plus_i =
69                         B_times_y_plus_vec_y_xy + i;
70
71                     /*For all pixels in the block*/
```

```

72     for (k = 0; k < Bx; k++) {
73         const int B_times_x_plus_vec_x_xy_plus_i_plus_k =
74             B_times_x_plus_vec_x_xy_plus_i + k;
75
76         const int B_times_x_plus_vec_x_xy_plus_k =
77             B_times_x_plus_vec_x_xy + k;
78
79         const int B_times_x_plus_k = B_times_x + k;
80
81     for (l = 0; l < By; l++) {
82         const int B_times_y_plus_vec_y_xy_plus_i_plus_l =
83             B_times_y_plus_vec_y_xy_plus_i + l;
84
85         const int B_times_y_plus_vec_y_xy_plus_l =
86             B_times_y_plus_vec_y_xy + l;
87
88         p1 = current[B_times_x_plus_k][B_times_y + l];
89
90         if ((B_times_x_plus_vec_x_xy_plus_i_plus_k) < 0 ||
91             (B_times_x_plus_vec_x_xy_plus_k) > (N - 1) ||
92             (B_times_y_plus_vec_y_xy_plus_l) < 0 ||
93             (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
94             p2 = 0;
95         } else {
96             p2 = previous[B_times_x_plus_vec_x_xy_plus_i + k]
97                 [B_times_y_plus_vec_y_xy_plus_l];
98         }
99
100        if ((B_times_x_plus_vec_x_xy_plus_k) < 0 ||
101            (B_times_x_plus_vec_x_xy_plus_k) > (N - 1) ||
102            (B_times_y_plus_vec_y_xy_plus_i_plus_l) < 0 ||
103            (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
104            q2 = 0;
105        } else {
106            q2 = previous[B_times_x_plus_vec_x_xy_plus_k]
107                [B_times_y_plus_vec_y_xy_plus_i_plus_l];
108        }
109
110        distx += sub_and_abs(p1, p2);
111        disty += sub_and_abs(p1, q2);
112    }
113 }
114
115 if (distx < min1) {
116     min1 = distx;
117     bestx = i;
118 }
119
120 if (disty < min2) {
121     min2 = disty;
122     besty = i;
123 }
124 }
125
126 {
127     const int i = 0;
128
129     distx = 0;
130     disty = 0;
131
132     /*For all pixels in the block*/
133     for (k = 0; k < Bx; k++) {
134         const int B_times_x_plus_vec_x_xy_plus_k =
135             B_times_x_plus_vec_x_xy + k;
136
137         const int B_times_x_plus_k = B_times_x + k;
138
139     for (l = 0; l < By; l++) {
140         const int B_times_y_plus_vec_y_xy_plus_l =
141             B_times_y_plus_vec_y_xy + l;
142
143         p1 = current[B_times_x_plus_k][B_times_y + l];
144
145         if ((B_times_x_plus_vec_x_xy_plus_k) < 0 ||
146             (B_times_x_plus_vec_x_xy_plus_k) > (N - 1) ||
147             (B_times_y_plus_vec_y_xy_plus_l) < 0 ||
148             (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
149             p2 = 0;
150         } else {
151             p2 = previous[B_times_x_plus_vec_x_xy_plus_k]
152                 [B_times_y_plus_vec_y_xy_plus_l];

```

```

153     }
154
155     int res = sub_and_abs(p1, p2);
156     distx += res;
157     disty += res;
158 }
159 }
160
161 if (distx < min1) {
162     min1 = distx;
163     bestx = i;
164 }
165
166 if (disty < min2) {
167     min2 = disty;
168     besty = i;
169 }
170 }
171
172 {
173     const int i = S;
174     distx = 0;
175     disty = 0;
176
177     const int B_times_x_plus_vec_x_xy_plus_i =
178         B_times_x_plus_vec_x_xy + i;
179     const int B_times_y_plus_vec_y_xy_plus_i =
180         B_times_y_plus_vec_y_xy + i;
181
182     /*For all pixels in the block*/
183     for (k = 0; k < Bx; k++) {
184         const int B_times_x_plus_vec_x_xy_plus_i_plus_k =
185             B_times_x_plus_vec_x_xy_plus_i + k;
186
187         const int B_times_x_plus_vec_x_xy_plus_k =
188             B_times_x_plus_vec_x_xy + k;
189
190         const int B_times_x_plus_k = B_times_x + k;
191
192         for (l = 0; l < By; l++) {
193             const int B_times_y_plus_vec_y_xy_plus_i_plus_l =
194                 B_times_y_plus_vec_y_xy_plus_i + l;
195
196             const int B_times_y_plus_vec_y_xy_plus_l =
197                 B_times_y_plus_vec_y_xy + l;
198
199             p1 = current[B_times_x_plus_k][B_times_y + l];
200
201             if ((B_times_x_plus_vec_x_xy_plus_i_plus_k) < 0 ||
202                 (B_times_x_plus_vec_x_xy_plus_i_plus_k) > (N - 1) ||
203                 (B_times_y_plus_vec_y_xy_plus_l) < 0 ||
204                 (B_times_y_plus_vec_y_xy_plus_l) > (M - 1)) {
205                 p2 = 0;
206             } else {
207                 p2 = previous[B_times_x_plus_vec_x_xy_plus_i + k]
208                     [B_times_y_plus_vec_y_xy_plus_l];
209             }
210
211             if (B_times_x_plus_vec_x_xy_plus_k < 0 ||
212                 B_times_x_plus_vec_x_xy_plus_k > (N - 1) ||
213                 (B_times_y_plus_vec_y_xy_plus_i_plus_l) < 0 ||
214                 (B_times_y_plus_vec_y_xy_plus_i_plus_l) > (M - 1)) {
215                 q2 = 0;
216             } else {
217                 q2 = previous[B_times_x_plus_vec_x_xy_plus_k]
218                     [B_times_y_plus_vec_y_xy_plus_i_plus_l];
219             }
220
221             distx += sub_and_abs(p1, p2);
222             disty += sub_and_abs(p1, q2);
223         }
224     }
225
226     if (distx < min1) {
227         min1 = distx;
228         bestx = i;
229     }
230
231     if (disty < min2) {
232         min2 = disty;

```

```

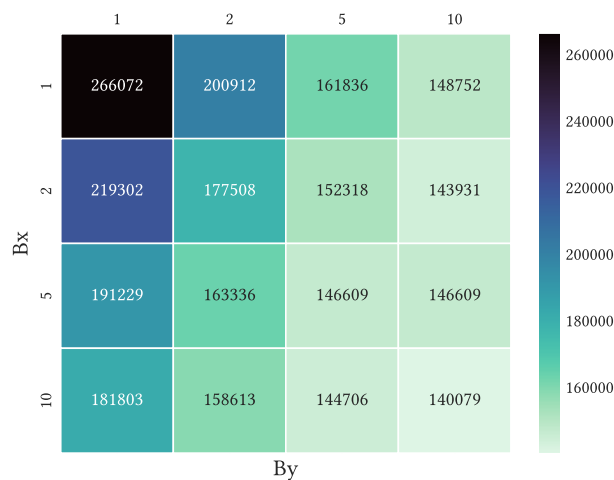
233         besty = i;
234     }
235 }
236
237 S = S / 2;
238 vectors_x[x][y] += bestx;
239 vectors_y[x][y] += besty;
240 }
241
242 // Reduce B * y to addition
243 B_times_y += By;
244 }
245
246 // Reduce B * x
247 B_times_x += Bx;
248 }
249 }

```

Οι χρόνοι εκτέλεσης του αλγορίθμου με ορθογώνιο παράθυρο, για διάφορες τιμές των ( $B_x$ ,  $B_y$ ), παρουσιάζονται στον παρακάτω πίνακα:

$B_x \backslash B_y$	1	2	5	10
1	266072	200912	161836	148752
2	219302	177508	152318	143931
5	191229	163336	146609	146609
10	181803	158613	144706	140079

Πίνακας 3: Συγκεντρωτικά αποτελέσματα ανά  $(B_x, B_y)$



Σχήμα 2: Χρόνοι εκτέλεσης ανά  $(B_x, B_y)$

Για τους ίδιους λόγους που αναφέρθηκαν και παραπάνω, καλύτερο χρόνο επιτυγχάνει η υλοποίηση για  $B_x = 10$ ,  $B_y = 10$ .

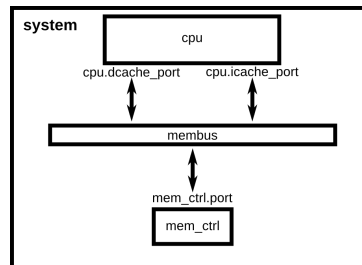
### Ευριστικός Αλγόριθμος.

Μία ευριστική μέθοδος, για την ανίχνευση κινήσεων στα frames, θα ήταν η τυχαία αναζήτηση block, κι όχι η αναζήτηση όλων των block. Αν το block, εμφανίζει αλλαγές μπορούμε να κινηθούμε ανάλογα με αυτές. Πρακτικά, ο αριθμός ο αριθμός των τυχαίων block που επιλέγουμε σε κάθε βήμα, αποτελεί και το quality factor του αλγορίθμου. Τέλος, η δημιουργία sub-blocks, θα μπορούσε να βελτιώσει την ποιότητα της συμπίεσης (όπως λ.χ. το H.265/HEVC).

## Μέρος 2ο

1.

Προσομοιώνουμε πρώτα στο gem5 την κάτωθι αρχιτεκτονική (CPU + main memory):



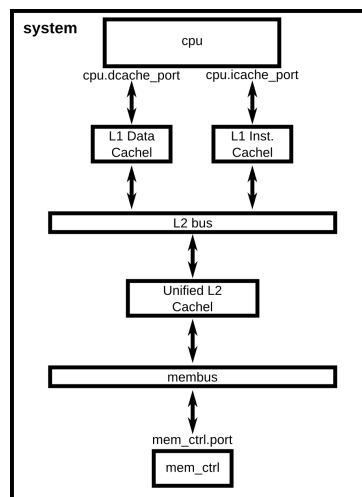
Σχήμα 3: Αρχιτεκτονική 1

και καταγράφουμε τους κύκλους που απαιτήθηκαν:

```
root@6262902d3b2f:/gem5# cat m5out/stats.txt | grep system.cpu.numCycles
system.cpu.numCycles      858895958      # Number of cpu cycles simulated (Cycle)
```

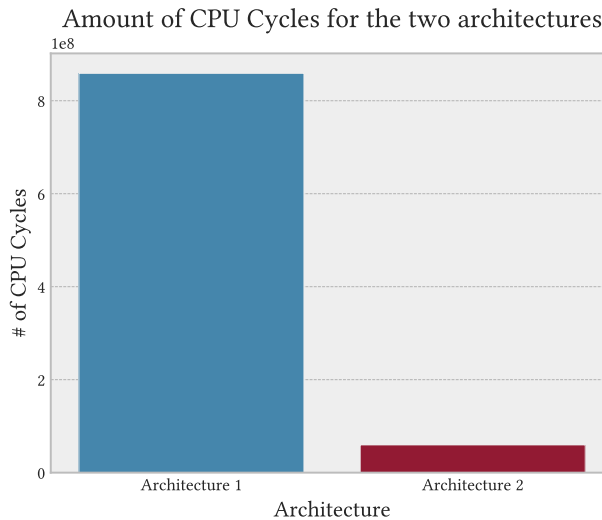
2.

Ομοίως για την ακόλουθη αρχιτεκτονική:



Σχήμα 4: Αρχιτεκτονική 2

```
root@6262902d3b2f:/gem5# cat m5out/stats.txt | grep system.cpu.numCycles
system.cpu.numCycles      59723234      # Number of cpu cycles simulated (Cycle)
```



Σχήμα 5: Σύγκριση συνολικών κύκλων εκτέλεσης αρχιτεκτονικών

Είναι προφανές, ότι με την προσθήκη κρυφής μνήμης μειώνεται σημαντικά το latency. Τα δεδομένα βρίσκονται πιο κοντά στον επεξεργαστή και δεν χρειάζεται να πληρώνουμε το overhead για την μεταφορά τους από την κύρια μνήμη σε κάθε single operation παρά μόνο όταν έχουμε misses. Επιπρόσθετα, ο επεξεργαστής που προσημειώνουμε (X86TimingSimpleCPU<sup>3</sup>), δεν έχει pipeline και χρειάζεται να περιμένει την ολοκλήρωση των προηγούμενων εντολών πριν εκτελέσει τις επόμενες. Συνεπώς, η προσθήκη κρυφής μνήμης και κατ'επέκταση μείωση του χρόνου αναμονής των δεδομένων, μειώνει το stall time και τον συνολικό αριθμό κύκλων εκτέλεσης. Τέλος, η προσθήκη της Instruction L1 cache, είναι σίγουρα οφέλη, αφού το πρόγραμμα υπό εξέταση, απαρτίζεται κατα κύριο λόγο δύο for-loop, οπότε εκτελούνται συνεχώς τα ίδια κομμάτια κώδικα.

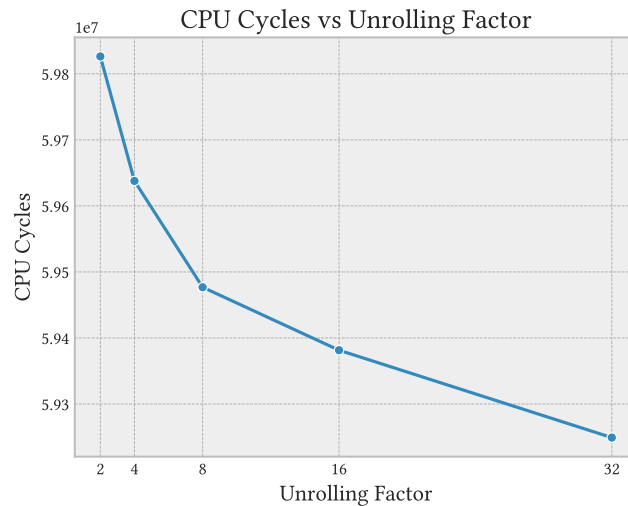
### 3.

Με το ακόλουθο script τρέχουμε τις προσομοιώσεις με τις διαφορετικές τιμές unrolling factor (2, 4, 8, 16, 32) και απεικονίζουμε τα αποτελέσματα:

#### unrolling.sh

```
1 #!/bin/bash
2
3 outfile="unrolling.txt"
4
5 rm -f $outfile
6 touch $outfile
7 echo -e "Unrolling Factor system.cpu.numCycles " >>$outfile
8 for i in 2 4 8 16 32; do
9     build/X86/gem5.opt configs/learning_gem5/part1/two_level.py /gem5/tables_UF/tables_uf$i.exe --
10     lli_size=8kB --lld_size=8kB --l2_size=128kB
11     echo -ne "${i}:\t\t\t" >>$outfile
12     cat m5out/stats.txt | grep system.cpu.numCycles | sed 's/[^0-9]*//g' >>$outfile
13 done
```

<sup>3</sup>[https://www.gem5.org/documentation/general\\_docs/cpu\\_models/SimpleCPU#timingsimplecpu](https://www.gem5.org/documentation/general_docs/cpu_models/SimpleCPU#timingsimplecpu)



Σχήμα 6: Συγκριτικό διάγραμμα για το Unrolling Factor

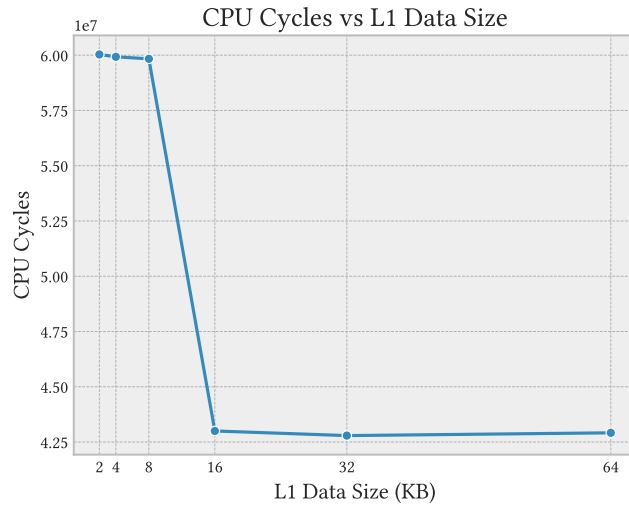
Ο αριθμός των κύκλων μειώνεται σημαντικά όσο το unrolling factor αυξάνεται. Δεδομένου ότι, έχουμε single threaded εκτέλεση σε επεξεργαστή χωρίς pipelining, η βελτίωση οφείλεται αποκλειστικά στην μείωση του overhead των μηχανισμών των loops (κόστος αύξησης loop counter, σύγκριση κλπ.)

#### 4.

Με το ακόλουθο script τρέχουμε τις προσομοιώσεις με τα διαφορετικά μεγέθη L1 Data Cache (2kB, 4kB, 8kB, 16kB, 32kB, 64kB) και απεικονίζουμε τα αποτελέσματα:

##### cache.sh

```
1 #!/bin/bash
2
3 outfile="cache.txt"
4
5 rm -f $outfile
6 touch $outfile
7 echo -e "L1D Size(KB) system.cpu.numCycles " >>$outfile
8 for i in 2 4 8 16 32 64; do
9     build/X86/gem5.opt configs/learning_gem5/part1/two_level.py /gem5/tables_UF/tables.exe --
10     lli_size=8kB --l1d_size=${i}kB --l2_size=128kB
11     echo -ne "${i}:\t\t\t" >>$outfile
12     cat m5out/stats.txt | grep system.cpu.numCycles | sed 's/^[^0-9]*//g' >>$outfile
13 done
```



Σχήμα 7: Συγκριτικό διάγραμμα για το L1 Data Cache Size

Αυξάνοντας το μέγεθος της L1 Data Cache παρατηρούμε μικρή βελτίωση από τα 2KB έως 8KB και ραγδαία επιτάχυνση από τα 8KB στα 16KB, ενώ μετά δεν υπάρχει κλιμάκωση. Ανατρέχοντας στην υλοποίηση της εν λόγω L1 cache<sup>4</sup>, βρίσκουμε ότι πρόκειται για μία 2-way associative cache. Εικάζουμε, λοιπόν, πώς η αύξηση του μεγέθους της cache, επιτρέπει την αποθήκευση περισσότερων δεδομένων και την μείωση των conflict misses. Στα 16KB ολόκληρο το Data Set κατάφερε τελικά να χωρέσει στην κρυφή μνήμη πληρώνοντας πλέον μόνο τα compulsory misses, οπότε περαιτέρω αύξηση της L1 Data Cache επιφέρει μικρή βελτίωση.

#### Εξαντλητική Αναζήτηση.

Με την μέθοδο της εξαντλητικής αναζήτησης, τρέχουμε τις προσομοιώσεις για όλους τους πιθανούς συνδυασμούς των μνημών του συστήματος και του unrolling factor που προκύπτουν από τα σύνολα που δίνονται (συνολικά  $6 \times 6 \times 4 \times 5 = 720$ ):

Parameters	Values
L1D Cache Size	2kB, 4kB, 8kB, 16kB, 32kB, 64kB
L1I Cache Size	2kB, 4kB, 8kB, 16kB, 32kB, 64kB
L2 Cache Size	128kB, 256kB, 512kB, 1024kB
Unrolling Factor	2, 4, 8, 16, 32

Πίνακας 4: Συνδυασμοί παραμέτρων

#### exhaustive.py

```

1 import sys
2 import subprocess
3 import csv
4
5 # search space values
6 L1D = [2, 4, 8, 16, 32, 64]
7 L1I = [2, 4, 8, 16, 32, 64]
8 L2 = [128, 256, 512, 1024]
9 unrolling = [2, 4, 8, 16, 32]

```

<sup>4</sup>[https://github.com/gem5/gem5/blob/stable/configs/learning\\_gem5/part1/caches.py](https://github.com/gem5/gem5/blob/stable/configs/learning_gem5/part1/caches.py)



```

10
11 outfile = "exhaustive.csv"
12
13 grep_cmd = f"cat m5out/stats.txt | grep system.cpu.numCycles | sed 's/[^0-9]*/g'"
14
15
16 with open(outfile, "w", newline="") as csvfile:
17     csv_writer = csv.writer(csvfile)
18     csv_writer.writerow(["L1D", "L1I", "L2", "Unrolling", "Cycles"])
19
20     for i in L1D:
21         for j in L1I:
22             for k in L2:
23                 for u in unrolling:
24                     cmd = f"build/X86/gem5.opt configs/learning_gem5/part1/two_level.py /gem5/
tables_UF/tables_uf{u}.exe --l1i_size={j}kB --l1d_size={i}kB --l2_size={k}kB"
25
26                     print("Running: ", cmd)
27                     try:
28                         subprocess.run(cmd, shell=True, check=True)
29                         result = subprocess.run(
30                             grep_cmd,
31                             shell=True,
32                             check=True,
33                             capture_output=True,
34                             text=True,
35                         )
36                         cc = result.stdout.strip()
37                         print(cc)
38                         if cc is not None:
39                             csv_writer.writerow(
40                                 [f"{i}", f"{j}", f"{k}", f"{u}", f"{cc}"]
41                             )
42                         else:
43                             print(f"No cycles found for command: {cmd}")
44
45                     except subprocess.CalledProcessError as e:
46                         print(f"Error executing command: {cmd}\n{e}")

```

Για την εξαγωγή των pareto βέλτιστων σημείων του χώρου εξερεύνησης όσον αφορά στο συνολικό μέγεθος μνήμης (L1D + L1I + L2) και τον αριθμό κύκλων ταυτόχρονα εφαρμόζουμε την built-in μέθοδο `paretoset` με συνάρτηση ελαχιστοποίησης και στα 2 μεγέθη και καταγράφουμε τα αποτελέσματα:

#### find\_pareto.py

```

1 """
2 @Authors: Maria Lazou, Stavros Avramidis
3 @Date: 11/11/2024
4 """
5
6 from paretoset import paretoset
7 import pandas as pd
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10
11 df = pd.read_csv("exhaustive.csv")
12 df["totalMemoryKB"] = df["L1D"] + df["L1I"] + df["L2"]
13 result_df = df[["totalMemoryKB", "Cycles"]].rename(columns={"Cycles": "Latency"})
14 print(result_df.head())
15
16 mask = paretoset(result_df, sense=["min", "min"])
17 # apply it on initial df
18 optimal = df[mask]
19 print(optimal)
20 optimal.to_csv("optimal_results.csv", index=False)
21
22 df["Pareto_optimal"] = mask
23
24 # Setup seaborn
25 sns.set_theme(style="whitegrid")
26 plt.style.use("bmh")
27 plt.rcParams["font.family"] = "Libertinus Serif"
28 plt.rcParams["font.size"] = 12
29 bmh_colors = plt.rcParams["axes.prop_cycle"].by_key()["color"]
30

```

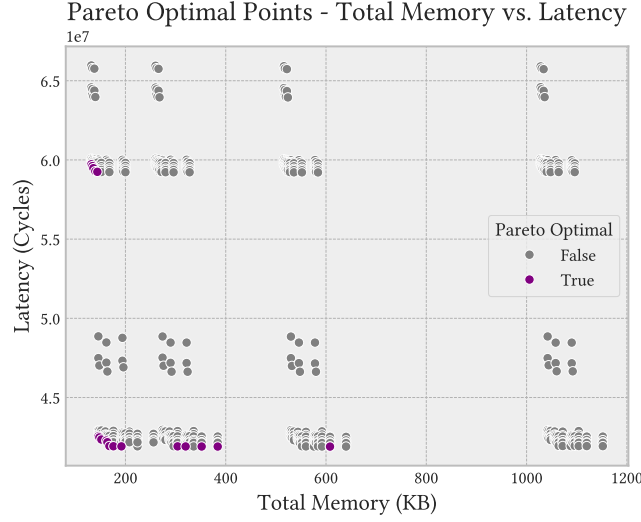
```

31 # Trick seaborn to bring pareto optimal points to the front
32 g = sns.scatterplot(
33     data=df[df["Pareto_optimal"] == False],
34     x="totalMemoryKB",
35     y="Cycles",
36     hue="Pareto_optimal",
37     palette={True: bmh_colors[1], False: "darkgray"},
38 )
39 sns.scatterplot(
40     data=df[df["Pareto_optimal"] == True],
41     x="totalMemoryKB",
42     y="Cycles",
43     hue="Pareto_optimal",
44     palette={True: bmh_colors[1], False: "darkgray"},
45 )
46 g.set_title("Pareto Optimal Points - Total Memory vs. Latency")
47 g.set_xlabel("Total Memory (KB)")
48 g.set_ylabel("Latency (Cycles)")
49 g.legend(title="Pareto Optimal")
50 plt.savefig("pareto_optimal.svg")

```

L1D	L1I	L2	Unrolling Factor	Total Memory (KB)	Latency(CC)
2	2	128	8	132	59725456
2	4	128	16	134	59608257
2	8	128	32	138	59447087
4	4	128	16	136	59477547
4	8	128	32	140	59287222
8	8	128	32	144	59249179
16	2	128	8	146	42547387
16	4	128	16	148	42480342
16	8	128	32	152	42348777
32	2	128	8	162	42278590
32	4	128	16	164	42176201
32	8	128	32	168	41936302
32	16	128	32	176	41922370
32	16	256	32	304	41915840
32	32	128	32	192	41915945
32	32	256	32	320	41910851
32	64	512	32	608	41901786
64	32	256	32	352	41906326
64	64	256	32	384	41904339

Πίνακας 5: Συγκεντρωτικός πίνακας των *pareto-optimal* σημείων



Σχήμα 8: Σύγκριση pareto-optimal σημείων

## 5.

Η ευριστική εξερεύνηση του ίδιου χώρου αναζήτησης ολοκληρώθηκε σε εκθετικά μικρότερο χρόνο και για τον ακριβή αριθμό των evaluations που έγιναν προσθέσαμε μια εκτύπωση της μεταβλητής `result.algorithm.evaluator.n_eval` στο `genOptimizer.py` η οποία επέστρεψε αποτέλεσμα **30**.

Η συνάρτηση επέστρεψε πολύ λιγότερα pareto-optimal που είναι όμως υποσύνολο της εξαντλητικής αναζήτησης όπως φαίνεται στην συνέχεια:

L1I	L1D	L2	Unrolling Factor	Total Memory(KB)	Latency(CC)
4	4	128	16	136	59477547
16	32	128	32	176	41922370
4	32	128	8	164	42227791
2	4	128	4	134	59780093

Πίνακας 6: Συγκεντρωτικός πίνακας των ευριστικών σημείων

Παρατηρούμε ότι δεν βρήκε το βέλτιστο Latency που είναι 41901786CC ούτε την ελάχιστη μνήμη 132kB. Άλλοι συνδυασμοί που δεν εντοπίστηκαν διαφέρουν κυρίως στα unrolling factors που είναι ανεξάρτητο της αρχιτεκτονικής και την συνολικής μνήμης που καλείται ο αλγόριθμος να ελαχιστοποιήσει.