



UNIVERSITÉ DE BORDEAUX

ind-47:
Rendu OpenMP/OpenCL

MENADJLIA SOFIANE

May 3, 2022

Contents

1	Version Hybride OpenMP/OpenCL	3
1.1	Initialisation du contexte	3
1.1.1	Code	3
1.2	Détection de terminaison	3
1.2.1	code	4
1.3	Appel du traitement GPU	4
1.3.1	Code	4
1.4	Traitement CPU	5
1.4.1	Code	5
1.5	Communication et échange de données	5
1.5.1	Code	6
1.6	Kernel OpenCL GPU	6
1.6.1	Code	6
2	Performances	7
2.1	Carte de chaleur	7
2.2	Trace	9
3	Rapport de bugs	9
3.1	Code	9
4	Conclusion	10

1 Version Hybride OpenMP/OpenCL

- La version hybride OpenMP/OpenCL consiste à diviser la charge du travail du noyau synchrone entre le CPU ainsi que le GPU afin de maximiser le rendement.
- Voici donc les différentes parties de l'implémentation:

1.1 Initialisation du contexte

- Afin d'implémenter cette version hybride il suffit de configurer la queue d'exécution du GPU au sein de la fonction `ssandPile_init_ocl_hybrid`.
- Après avoir appelé de la fonction `ssandPile_init`, la répartition de charges doit aussi être faite dans cette fonction.

1.1.1 Code

- `cpu_y_part` et `gpu_y_part` représentent les parties à faire par le CPU et le GPU respectivement.
- `size` représente le nombre de tuiles dans la partie à faire par le GPU.
- `TABLE_BUFF` et `new_buff` sont initialisés afin de permettre la détection de la terminaison du programme.

```
1 void ssandPile_init_ocl_hybrid (void)
2 {
3     ssandPile_init();
4
5     if (GPU_TILE_H != TILE_H)
6         exit_with_error ("CPU and GPU Tiles should have the same height (%d != %d)",
7                           GPU_TILE_H, TILE_H);
8
9     cpu_y_part = (NB_TILES_Y / 2) * GPU_TILE_H; // Start with fifty-fifty
10    gpu_y_part = DIM - cpu_y_part;
11
12    int size = (DIM / GPU_TILE_W) * (gpu_y_part / GPU_TILE_H);
13    TABLE_BUFF = (TYPE *) malloc(size * sizeof(TYPE));
14
15    if (!TABLE_BUFF)
16        exit_with_error ("unable to allocate new buffer table!!\n");
17
18    new_buff = clCreateBuffer (context, CL_MEM_READ_WRITE,
19                              sizeof (TYPE) * size, NULL, NULL); // to check size
20
21    if (!new_buff)
22        exit_with_error ("unable to allocate new buffer!!\n");
23 }
24
```

1.2 Détection de terminaison

- La version OpenCL du rendu précédent n'était pas fini car elle n'arrivait pas à détecter la stabilisation du tas de sable, il fallait donc préciser le bon nombre d'itérations à effectuer au lancement du programme.

- Étant donné que la version hybride est exécutée sur deux supports différents il faut qu'il communique l'information sur l'état de stagnation courant à travers la variable change.
- Il se trouve que en OpenCL la lecture et écriture dans un buffer est très coûteuse, il faut donc trouver un système pour éviter de récupérer cette information à chaque itération.
- Dans ma version j'ai choisi de vérifier cela chaque **DIM/8** car ça donne le meilleur rendement avec différentes tailles d'images.

1.2.1 code

- Je commence par initialiser TABLE_BUFF et new_buff dans la fonction `ssandPile_init_ocl_hybrid` avec le même `size = (DIM / GPU_TILE_W) * (gpu_y_part / GPU_TILE_H)` qui représente le nombre de tuiles dans la partie GPU.
- Chaque élément dans ce buffer contient le changement cumulé de toutes les cellules dans la tuile correspondante.
- Après avoir exécuté le code des deux parties on teste que la variable change est à 0 dans la partie CPU ainsi que le nombre d'itérations est bien divisible par (DIM/8) avant de lire les changements du côté GPU en faisant un READ de new_buff dans TABLE_BUFF.
- Enfin il suffit de parcourir TABLE_BUFF et récupérer le changement total dans la variable change qu'on pourra après tester pour savoir si le programme s'est terminé.

```

1  if (change == 0 && iteration % (DIM/8) == 0){
2      err = clEnqueueReadBuffer (queue, new_buff, CL_TRUE, 0, sizeof(unsigned
        ) * size, TABLE_BUFF, 0, NULL, NULL);
3
4      check (err, "Failed to read to buffer");
5
6      #pragma omp parallel for schedule(runtime) reduction(|: change)
7      for (int i = 0; i < size ; i++){
8          change |= TABLE_BUFF[i];
9      }
10
11     if (!change){
12         diff = it;
13         break;
14     }
15 }

```

1.3 Appel du traitement GPU

- Cette partie concerne l'implémentation du code du côté CPU qui permet d'invoquer le kernel OCL du côté GPU.
- Pour cela, après avoir reparté la charge de travail, il faut ajouter les arguments à passer au kernel OCL puis l'appeler.

1.3.1 Code

Les paramètres à passer sont les suivants:

- **cur_buffer** : buffer en cours d'utilisation
- **next_buffer** : buffer de la prochaine itération
- **new_buffer** : buffer permettant de détecter la terminaison du programme
- **cpu_y_part** : servant de offset pour le kernel afin de ne pas commencer au début de l'image.
- **itération** : permettant de vérifier les itérations avant decrire dans **new_buff**

```

1 // Set kernel arguments
2
3 err = 0;
4 err |= clSetKernelArg (compute_kernel, 0, sizeof (cl_mem), &cur_buffer);
5 err |= clSetKernelArg (compute_kernel, 1, sizeof (cl_mem), &next_buffer);
6 err |= clSetKernelArg (compute_kernel, 2, sizeof (cl_mem), &new_buff);
7 err |= clSetKernelArg (compute_kernel, 3, sizeof (unsigned), &cpu_y_part);
8 err |= clSetKernelArg (compute_kernel, 4, sizeof (unsigned), &iteration);
9
10 check (err, "Failed to set kernel arguments");
11
12 // Launch GPU kernel
13 err = clEnqueueNDRangeKernel (queue, compute_kernel, 2, NULL, global, local,
14                               0, NULL, &kernel_event);
15 check (err, "Failed to execute kernel");
16
17 clFlush (queue);

```

1.4 Traitement CPU

- Le traitement du cote CPU est fait grâce a OpenMP et est très similaire a ceux des rendus précédents, avec l'exception de parcourir les lignes jusqu'à **cpu_y_part** au lieu de **DIM** pour ce limiter a la partie CPU.
- les variable t1 et t2 permettent de connaître la durée de temps de calcul de la partie cpu.

1.4.1 Code

```

1
2 t1 = what_time_is_it ();
3 int change = 0;
4 #pragma omp parallel for collapse(2) schedule(runtime) reduction (|: change)
5 for (int y = 0; y < cpu_y_part; y += TILE_H)
6     for (int x = 0; x < DIM; x += TILE_W)
7         change |=
8             do_tile(x + (x == 0), y + (y == 0),
9                     TILE_W - ((x + TILE_W == DIM) + (x == 0)),
10                     TILE_H - ((y + TILE_H == DIM) + (y == 0)), omp_get_thread_num
11                     ());
12 t2 = what_time_is_it ();

```

1.5 Communication et échange de données

- Afin de permettre le bon fonctionnement de programme il faut que le kernel OCL communique les résultat des données traitées au traitement CPU et vice versa.

- Pour cela il suffit de communiquer les deux lignes correspondantes a l'intersection entre les deux parties, notamment les lignes `cpu_y_part` et `gpu_y_part`.
- Concrètement, j'utilise les fonctions `clEnqueueReadBuffer` et `clEnqueueWriteBuffer` qui servent a lire la première ligne du renvoyée par le gpu ainsi qu'écrire la dernière ligne traitée par le cpu.

1.5.1 Code

```

1 //READ WRITE
2
3 err = clEnqueueReadBuffer (queue, next_buffer, CL_TRUE, (cpu_y_part) * DIM *
4     sizeof(unsigned) ,
5     DIM * sizeof (unsigned), &table(out, (cpu_y_part),
6     0), 0,
7     NULL, NULL);
8
9 check (err, "Failed to read to buffer");
10
11 err = clEnqueueWriteBuffer (queue, next_buffer, CL_TRUE, (cpu_y_part - 1) *
12     DIM * sizeof(unsigned) ,
13     DIM * sizeof (unsigned), &table(out, (cpu_y_part -
14     1), 0), 0,
15     NULL, NULL);
16
17 check (err, "Failed to write to buffer");

```

1.6 Kernel OpenCL GPU

- Le traitement du cote GPU est fait grâce a OpenCL et est très similaire a celui du rendu OCL précédent.
- Comme pour la partie CPU, il faut commencer a partir de la bonne ligne par rapport a la partie désigné pour le traitement GPU.
- il faut aussi gérer l'écriture dans le buffer `new_buff` pour permettre la lecture du cote CPU.

1.6.1 Code

- Le code ne diffère pas trop sauf qu'il faut incrémenter y de `offset` qui vaut toujours `cpu_y_part`, avant de commencer afin de se limiter a la partie GPU.
- Il faut trouver la tuile correspondante de chaque cellule en normalisant le x et y avant de pouvoir mettre le cumul du changement total dans le buffer `new_buff`.

```

1
2 __kernel void ssandPile_ocl_hybrid (__global unsigned *in, __global unsigned *out
3     , __global unsigned *new, unsigned offset, unsigned iteration)
4 {
5     int y = get_global_id (1) + offset;
6     int x = get_global_id (0);
7
8     int index = y * DIM + x;
9
10    unsigned local_y = get_local_id (1);

```

```

10 unsigned local_x = get_local_id (0);
11
12 if(x > 0 && y > 0 && x < DIM-1 && y < DIM-1)
13 {
14
15     out[index] =
16     (in[y * DIM + x] % 4) +
17     (in[(y+1) * DIM + x] /4) +
18     (in[(y-1) * DIM + x] /4) +
19     (in[y * DIM + (x+1)] /4) +
20     (in[y * DIM + (x-1)] /4);
21
22     barrier (CLK_LOCAL_MEM_FENCE);
23
24     y = y / GPU_TILE_H;
25     x = x / GPU_TILE_W;
26     int tile_index = y * (DIM / GPU_TILE_W) + x;
27     if (iteration % (DIM/8) == 0)
28         new[tile_index] |= in[index] != out[index];
29 }
30 }

```

2 Performances

2.1 Carte de chaleur

- En ce qui concerne cette version hybride on remarque un speedup assez grand avec un threading de 44 a 48.
- on remarque que la version statique donne un rendement nettement supérieure a celui de la dynamique.
- Le tuilage optimale et de $32 * 16$ avec 44 threads, $16 * 32$ avec 46 et $32 * 16$ avec 48.



Figure 1: ssandPile version ocl_hybrid avec do_tile_opt size 2048 mode alea tile size 32*16 machine dali(48 threads)

2.2 Trace

Voici la trace générée grâce à l'option -t de EasyPAP ou on peut voir que les threads cpu ne sont très bien exploités et prennent beaucoup de temps par rapport au gpu.

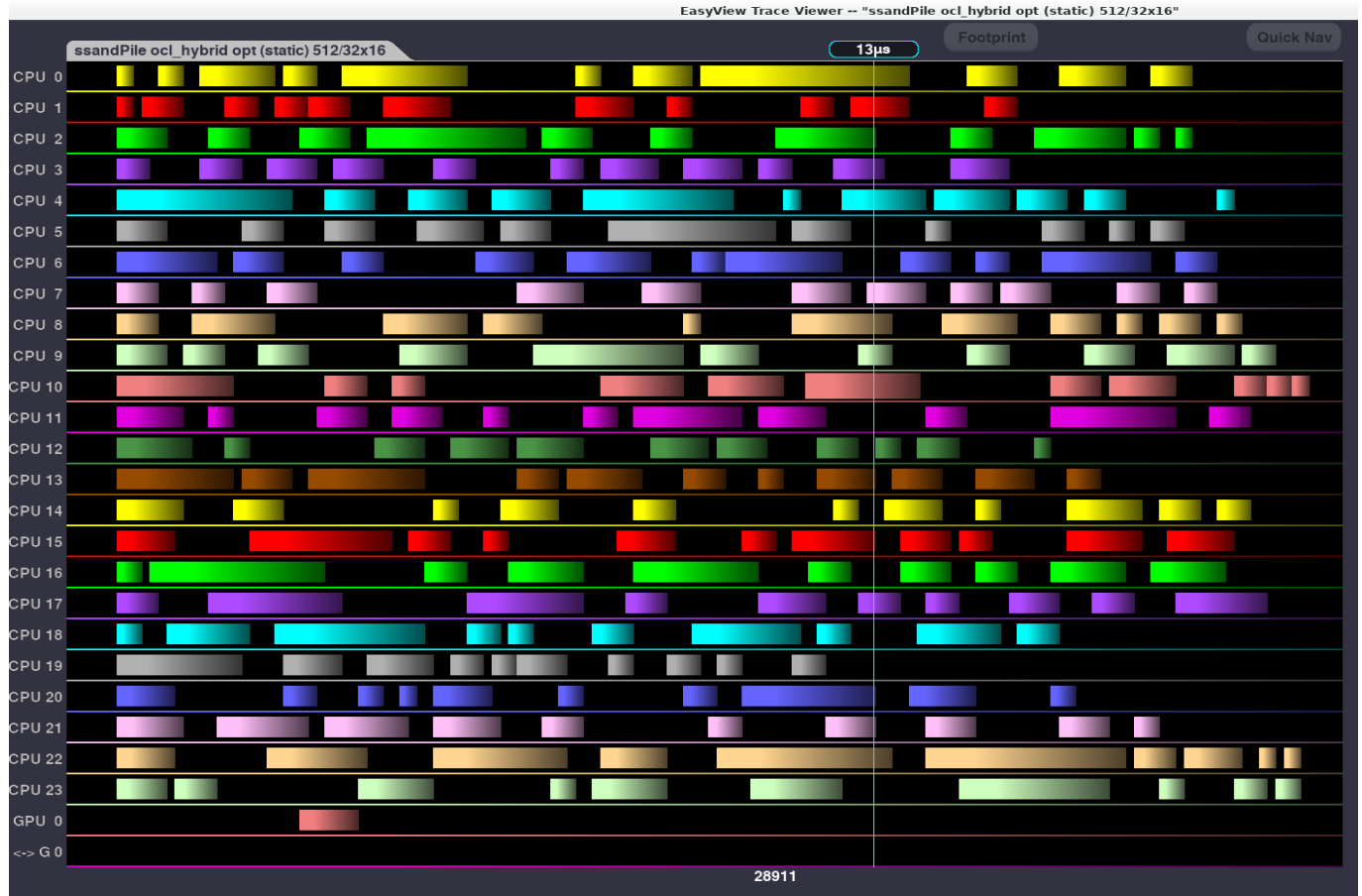


Figure 2: trace-ssand-omp-tiled-512/32x16

3 Rapport de bugs

- Malheureusement, je n'ai pas réussi à implémenter la partie concernant le rééquilibrage de charges, même après avoir changé les valeurs `cpu_y_part` et `gpu_y_part` ainsi que mis à jour le tableau global et avoir ré-alloué les buffers de terminaison avec la bonne taille.
- le programme dans ce cas semble de seulement marcher correctement quand appelé avec l'option "-i" et ceci est dû au fait que le programme est appelé plusieurs fois avec un `nb_iter` mis à 1 ce qui impose l'appel de la fonction `ssandPile_init_ocl_hybrid` avec la répartition de travail mise à 50% pour chacun.
- Même en initialisant la valeur `cpu_y_part` à $(NB_TILES_Y / 4) * GPU_TILE_H$ au lieu de $(NB_TILES_Y / 2) * GPU_TILE_H$.
- Le code correspondant aux bugs mentionnés ci-dessus est commenté sur le code.

3.1 Code

```

1 // Load balancing
2 // if (gpu_duration != 0) {
3 // if (much_greater_than (gpu_duration, cpu_duration) &&
4 // gpu_y_part > GPU_TILE_H) {
5 // gpu_y_part -= GPU_TILE_H;
6 // cpu_y_part += GPU_TILE_H;
7 // global[1] = gpu_y_part;
8 // }
9 // else
10 // if (much_greater_than (cpu_duration, gpu_duration) &&
11 // cpu_y_part > GPU_TILE_H) {
12 // gpu_y_part += GPU_TILE_H;
13 // cpu_y_part -= GPU_TILE_H;
14 // global[1] = gpu_y_part;
15 // }
16 // }
17 // }
18 size_t size = (DIM / GPU_TILE_W) * (gpu_y_part / GPU_TILE_H);
19 // TABLE_BUFF = realloc(TABLE_BUFF, size * sizeof(TYPE));
20 // cl_mem new = clCreateBuffer (context, CL_MEM_READ_WRITE, sizeof (TYPE) *
    size , NULL, NULL);

```

4 Conclusion

- Dans ce rendu j'ai réussi à implémenter la version hybride OpenMP et OpenCL indiqué dans le rendu 4 ainsi que la détection de terminaison du programme.
- Je tiens à préciser que j'ai appris beaucoup sur la programmation parallèle grâce à cette UE, ça m'a permis de mieux connaître mes préférences ainsi qu'approfondir mes compétences dans le domaine.