

# Elementy GPGPU

Przemysław Kiciak

1

Istnieje kilka „wysokopoziomowych” języków programowania GPU; najbardziej znane to GLSL i HLSL („graficzne”) oraz OpenGL i CUDA („ogólnego stosowania”).

Język GLSL (*OpenGL shading language*) powstał jako uzupełnienie standardu OpenGL — pisze się w nim tzw. **szadery**, czyli programy, które są wmontowywane w tzw. **potok przetwarzania grafiki** (*rendering pipeline*). Zadaniem jego kolejnych etapów jest przetworzenie danych geometrycznych i innych na piksele obrazu, zgodnie z opisem podanym w specyfikacjach standardów OpenGL lub Vulkan.

Standardy te umożliwiają wykorzystywanie także tzw. **szaderek obliczeniowych** (*compute shaders*), działających poza potokiem przetwarzania grafiki. Szadery te mogą wykonywać dowolne obliczenia, których wyniki mogą być użyte później do wykonywania obrazów lub przesłane z RAM GPU do RAM CPU.

3

## Wprowadzenie

GPU (*graphics processing unit*) to „karta graficzna” wyposażona w pamięć RAM i zestaw procesorów specjalnie zaprojektowanych do wytwarzania obrazów.

Natura wykonywanych przy tym obliczeń umożliwia maszyną równoległość.

Procesory (rdzenie obliczeniowe) GPU są prostsze niż CPU, ale ich liczba w jednym urządzeniu jest o 2–3 rzędy wielkości większa niż liczba rdzeni CPU. Dlatego moc obliczeniowa GPU też jest o ok. 2 rzędy wielkości większa.

GPGPU (*general purpose GPU programming*) jest to zastosowanie GPU do obliczeń niekoniecznie związanych z grafiką. Celem jest pełne wykorzystanie mocy obliczeniowej GPU.

2

Aplikacja OpenGL-a musi utworzyć tzw. **kontekst OpenGL-a**, a następnie tzw. **programy szaderek** i **bufory**, w których będą umieszczone dane, wyniki pośrednie i wyniki końcowe obliczeń. Program szaderek powstaje z tekstów źródłowych w GLSL-u, które należy (przy użyciu procedur obecnych w bibliotece OpenGL) skompilować i złączyć.

Treść szadera obliczeniowego opisuje działanie **jednego wątku**; wywołanie szadera za pomocą procedury `glDispatchCompute` rozpoczyna równoległe wykonywanie **grupy roboczej** wątków.

**Lokalna grupa robocza** jest jedno-, dwu- lub trójwymiarową tablicą wątków, które wymiary są zadeklarowane w treści szadera (i nie można ich zmienić bez ponownej kompilacji szadera).

**Globalna grupa robocza** jest jedno-, dwu- lub trójwymiarową tablicą lokalnych grup roboczych. Jej wymiary są określone przez parametry procedury `glDispatchCompute` (i za każdym razem mogą być inne).

4

Każdy wątek otrzymuje informacje o wymiarach grupy globalnej i o swoim położeniu w grupie roboczej (lokalnej i globalnej) w **zmiennych wbudowanych**, które są trójkami indeksów do wspomnianych tablic.

Podstawowym problemem w programowaniu GPU jest wymaganie tzw. **jednolitości obliczeń**. Poszczególne rdzenie GPU wykonują **wątki obliczeniowe** zorganizowane w **grupy robocze**. Wszystkie procesory działające w grupie albo wykonują jednocześnie tę samą instrukcję albo czekają — co umożliwia realizację instrukcji warunkowych, ale im więcej rdzeni czeka, tym mniej efektywnie jest wykorzystywana GPU.

Zadaniem sterownika GPU jest przydzielenie wątków poszczególnym rdzeniom GPU, a także zapewnianie synchronizacji obliczeń.

5

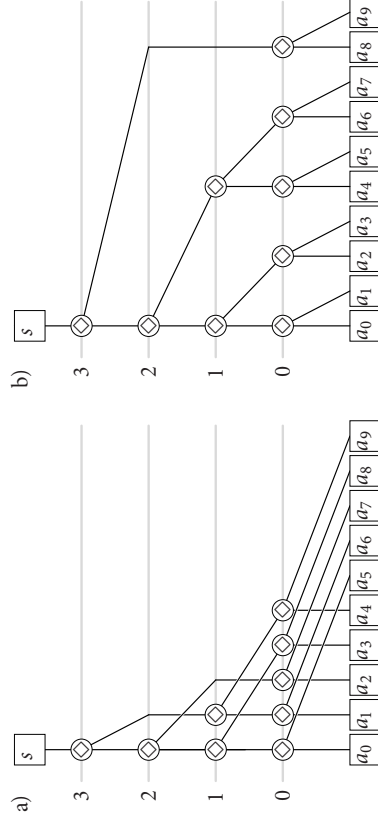
Najważniejsze zmienne, do których shader obliczeniowy ma dostęp, to

- **zmienne jednolite** (*uniform variables*) — wartości nadaje im aplikacja działająca na CPU, wszystkie wątki obliczeniowe „widzą” w danej chwili tę samą wartość każdej takiej zmiennej,
- **zmienne w blokach magazynowych** (*shader storage blocks*), które mogą być zarówno czytane, jak i pisane przez shader. Należy dbać o to, aby poszczególne wątki przypisywały wartości różnym zmiennym (np. różnym elementom tablicy),
- **zmienne wbudowane** (*built-in variables*), które opisują wielkości grupy roboczych (tzw. lokalnej i globalnej) i numer lokalnej grupy roboczej oraz numer wątku w grupie,
- **zmienne współdzielone** (*shared variables*), za pomocą których wątki shadera mogą przekazywać sobie informacje w obrębie lokalnej grupy roboczej.

6

## Działania parami

Działania dwuargumentowe często są iterowane, gdy trzeba obliczyć np. sumę wielu składników lub iloczyn wielu czynników. Aby takie obliczenia mogły być zrównoleżone, realizowane działanie „ $\diamond$ ” musi być co najmniej **łącznie**; jeśli jest także **przemienne**, to możliwe są różne algorytmy.



7

Obliczenie sumy lub iloczynu elementów ciągu o długości  $n$  można wykonać w  $\lceil \log_2 n \rceil$  krokach. Pierwszy algorytm jest realizowany przez shader

```

1: #version 450 core
2:
3: #define GROUP_SIZE 64
4: layout (local_size_x = GROUP_SIZE) in;
5:
6: uniform uint n;
7:
8: void AddTwoTerms ( uint i, uint j ); /* działanie łączne i przemienne */
9:
10: void main ( void )
11: {
12:     uint i, j;
13:
14:     i = uint ( gl_GlobalInvocationID.x );
15:     if ( (j = i+(n+1)/2) < n )
16:         AddTwoTerms ( i, j );
17: } /*main*/

```

8

Procedura `AddTwoTerms` realizuje działanie; jeśli ma to być dodawanie liczb całkowitych umieszczonych w tablicy, to może ona wyglądać tak:

```
1: layout (std430, binding=0) buffer Data { int a[]; } data;
2:
3: uniform int n0;
4: void AddTwoTerms ( int i, int j )
5: {
6:   data.a[n0+i] += data.a[n0+j];
7: } /*AddTwoTerms*/
```

Jak widać, „psujemy” zawartość tablicy. Końcowy wynik, tj. suma liczb danych początkowo na pozycjach od `n0` do `n0+n-1` ma się znaleźć na pozycji `n0`.

9

Szader ten jest wywoływany przez procedurę

```
1: static GLuint program_id, GLuint uloc[2];
2: static GLint lgsz[3];
3:
4: void GPUsumUp ( GLuint n, GLuint n0, GLuint databuf )
5: {
6:   glUseProgram ( program_id );
7:   glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, databuf );
8:   glUniform1ui ( uloc[1], n0 );
9:   while ( n > 1 ) {
10:    glUniform1ui ( uloc[0], n );
11:    glDispatchCompute ( (n/2+lgsz[0]-1)/lgsz[0], 1, 1 );
12:    n = (n+1)/2;
13:    glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
14:   }
15:   ExitIfGLError ( "GPUSumUp" );
16: } /*GPUSumUp*/
```

10

Drugi algorytm sumowania parami (który zakłada tylko łączność działania, w związku z czym nie zmienia kolejności argumentów) jest realizowany przez szader

```
1: #version 450 core
2: #define GROUP_SIZE 64
3: layout (local_size_x = GROUP_SIZE) in;
4: uniform uint n, q;
5:
6: void AddTwoTerms ( uint i, uint j ); /* dowolne działanie łączne */
7:
8: void main ( void )
9: {
10:   uint i, j;
11:
12:   i = uint ( (q+q)+lg_GlobalInvocationID.x );
13:   if ( (j = i+q) < n )
14:     AddTwoTerms ( i, j );
15: } /*main*/
```

W zmiennej `q` jest odstęp w tablicy elementów, które mają być dodane — w kolejnych krokach to są kolejne całkowite potęgi liczby 2.

11

Działająca na CPU procedura sumowania realizująca drugi algorytm jest taka:

```
1: static GLuint program_id, GLuint uloc[3];
2: static GLint lgsz[3];
3:
4: void GPUAltSumUp ( GLuint n, GLuint n0, GLuint databuf )
5: {
6:   GLuint q;
7:
8:   glUseProgram ( program_id );
9:   glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, databuf );
10:  glUniform1ui ( uloc[0], n ); /* uniform n0 = n0; */
11:  glUniform1ui ( uloc[1], n0 ); /* uniform n = n; */
12:  for ( q = 1; n > 1; q += q, n = (n+1)/2 ) {
13:    glUniform1ui ( uloc[2], q ); /* uniform q = q; */
14:    glDispatchCompute ( (n/2+lgsz[0]-1)/lgsz[0], 1, 1 );
15:    glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
16:   }
17: } /*GPUAltSumUp*/
```

12

**Uwagi:** Z powodu „psucia” początkowej zawartości tablicy to są algorytmy o sporej złożoności pamięciowej, którą można zamienić na czas — jeśli oryginalne dane są później potrzebne, to należy albo je skopiować i zepsuć kopię, albo zepsuć oryginał, odczytać wynik obliczeń i ponownie wygenerować dane.

Dodawanie i mnożenie liczb zmiennopozycyjnych (lub macierzy o współczynnikach zmiennopozycyjnych) jest łączne z dokładnością do błędów zaokrążeń — czyli *nie jest* łączne, ale się stara. Obliczony wynik będzie sumą lub iloczynem danych zaburzonych na poziomie błędów zaokrążeń, co oznacza numeryczną poprawność.

Algorytm numerycznie poprawny jest tym lepszy im mniejsze są stałe kumulacji. W przypadku sumowania po kolei  $n$  liczb one są rzędu  $n$ , a algorytmy sumowania parami mają stałe rzędu  $\log n$  — czyli są lepsze.

13

Wybór mniejszego lub większego elementu jest też działaniem łącznym i przemianym, zatem można w czasie rzędu  $\log n$  znaleźć najmniejszy lub największy element ciągu. Zobaczymy, jak oba skrajne elementy znaleźć jednocześnie.

W pierwszym kroku algorytmu użyjemy procedury zwanej **komparatorem**, która porówna dwa elementy na wskazanych miejscach i jeśli drugi jest mniejszy niż pierwszy, to je przestawi. W ten sposób w przetworzonej przez komparator parze pierwszy element jest mniejszy lub równy drugiemu.

Porządkujemy w ten sposób pary elementów sąsiednich. Dodatkowo, jeśli liczba  $n$  jest nieparzysta, to ostatni element ciągu (ten bez pary) trzeba porównać z elementami w pierwszej parze i ewentualnie zastąpić nim jeden z nich.

W kolejnych krokach działamy na parach par, znajdując w nich elementy najmniejsze i największe. W każdym kroku liczba par zawierających istotne dane maleje dwukrotnie.

14

Procedury pomocnicze:

```
1: void CompSwap ( uint i, uint j ) /* komparator */
2: {
3:   int x;
4:
5:   if ( data.a[i] += n0 ] > data.a[j] += n0 ] )
6:     { x = data.a[i]; data.a[i] = data.a[j]; data.a[j] = x; }
7: } /*CompSwap*/
8:
9: void ChooseMin ( uint i, uint j )
10: {
11:   if ( data.a[i] += n0 ] > data.a[j] += n0 ] )
12:     data.a[i] = data.a[j];
13: } /*ChooseMin*/
14:
15: void ChooseMax ( uint i, uint j )
16: {
17:   if ( data.a[i] += n0 ] < data.a[j] += n0 ] )
18:     data.a[i] = data.a[j];
19: } /*ChooseMax*/
```

15

Teraz procedura main szadera:

```
1: uniform uint n;
2: uniform bool s;
3:
4: void main ( void )
5: {
6:   uint i, j;
7:
8:   i = uint ( gl_GlobalInvocationID.x ); i += i;
9:   if ( s ) {
10:    if ( ( j = i+1 ) < n )
11:      CompSwap ( i, j );
12:    if ( i == 0 && ( n & 0x01 ) != 0 ) { /* n nieparzyste */
13:      ChooseMin ( 0, n-1 );
14:      ChooseMax ( 1, n-1 );
15:    }
16:  }
17:  else if ( ( j = i + 2*((n+3)/4) ) < n ) {
18:    ChooseMin ( i, j );
19:    ChooseMax ( i+1, j+1 );
20:  }
21: } /*main*/
```

16

Procedura na CPU wywołująca szader w kolejnych krokach:

```

1: static GLuint program_id, uloc[3];
2: static GLint lgsz[3];
3:
4: void GPUFindMinMax ( GLuint n, GLuint n0, GLuint databuf )
5: {
6:     glUseProgram ( program_id );
7:     glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, databuf );
8:     glUniform1ui ( uloc[0], n );
9:     glUniform1ui ( uloc[1], n0 );
10:    glUniform1ui ( uloc[2], true ); /* uniform s = true; */
11:    glDispatchCompute ( (n/2+lgsz[0]-1)/lgsz[0], 1, 1 );
12:    glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
13:    glUniform1ui ( uloc[2], false ); /* uniform s = false; */
14:    for ( n &= ~0x01; n > 2; n = 2*((n+3)/4) ) {
15:        glUniform1ui ( uloc[0], n );
16:        glDispatchCompute ( (n/4+lgsz[0]-1)/lgsz[0], 1, 1 );
17:        glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
18:    }
19:    ExitIfGLError ( "GPUFindMinMax" );
20: } /*GPUFindMinMax*/

```

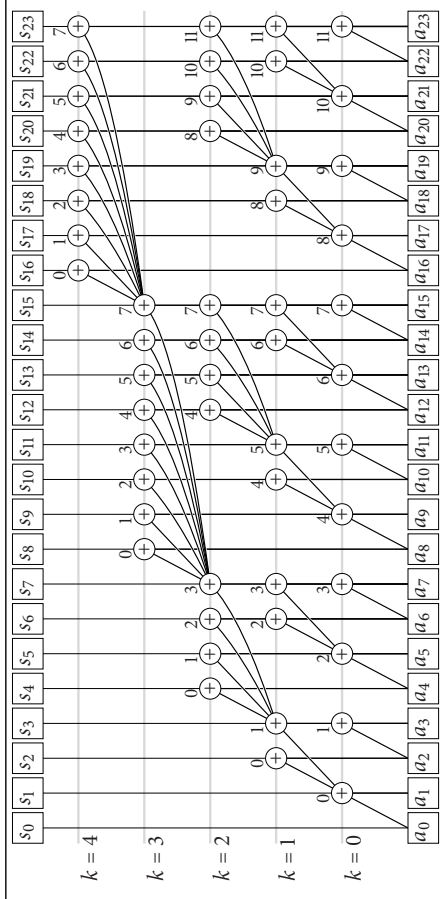
## Obliczanie sum prefiksowych

Dany jest ciąg  $a_0, \dots, a_{n-1}$  i działanie łączne w zbiorze, którego to są elementy — dla ustalenia uwagi dodawanie. Należy obliczyć ciąg sum prefiksowych  $s_0, \dots, s_{n-1}$  którego elementy są zdefiniowane wzorem

$$s_i = \sum_{j=0}^i a_j.$$

To zadanie również można wykonać w  $\lceil \log_2 n \rceil$  krokach, mając do dyspozycji  $\lceil n/2 \rceil$  procesorów (albo raczej wątków) działających równolegle.

Ciąg dany w tablicy zostanie zastąpiony przez wynik.



W  $k$ -tym kroku kolejne wątki obliczają kolejne sumy podciągów. Wątek  $i$ -ty ma obliczyć sumę elementów o indeksach

$$i_a = 2^{k+1} \lfloor i/2^k \rfloor + 2^k - 1, \quad i_b = i_a + (i \bmod 2^k) + 1$$

i zapamiętać ją na miejscu  $i_b$ .

Procedura na GPU realizująca krok obliczania sum prefiksowych nie ma nazwy `main`, bo zazwyczaj jest ona częścią dużo większego szadera, w którym pełni rolę pomocniczą (rozwiązuje podzadanie).

```

1: layout (std430, binding=0) buffer Data { int a[]; } data;
2: uniform uint prN0, prN, prStep;
3:
4: void iPrefixSum ( uint i )
5: {
6:     uint ii, m0, m1, ia, ib;
7:
8:     ii = i+i; m0 = 0x01 << prStep; m1 = m0-1;
9:     ia = (ii & ~m0) | m1;
10:    if ( (ib = ia + (i & m1) + 1) < prN )
11:        data.a[prN0 + ib] += data.a[prN0 + ia];
12: } /*iPrefixSum*/

```

Z tego samego powodu poniższy podprogram na CPU obliczający sumy prefiksowe nie wywołuje procedury `glUseProgram`, natomiast nadaje zmiennej jednolitej `stage` wartość powodującą wywołanie procedury wykonującej krok sumowania, a potem wywołuje szader w pętli.

```

1: static void iPrefixSum ( GLuint NO, GLuint N )
2: {
3:     unsigned int k, m, d;
4:     glUniform1i ( uvloc[0], 0 ); /* uniform stage = 0; */
5:     glUniform1ui ( uvloc[1], NO ); /* uniform prNO = NO; */
6:     glUniform1ui ( uvloc[2], N ); /* uniform prN = N; */
7:     d = (N+2*lgsize[0]-1)/(2*lgsize[0]);
8:     for ( k = 0, m = N-1; m > 0; k++, m >= 1 ) {
9:         glUniform1ui ( uvloc[3], k ); /* uniform prStep = k; */
10:         glDispatchCompute ( d, 1, 1 );
11:         glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
12:     }
13: }
14: ExitIfGLError ( "iPrefixSum" );
15: } /*iPrefixSum*/

```

21

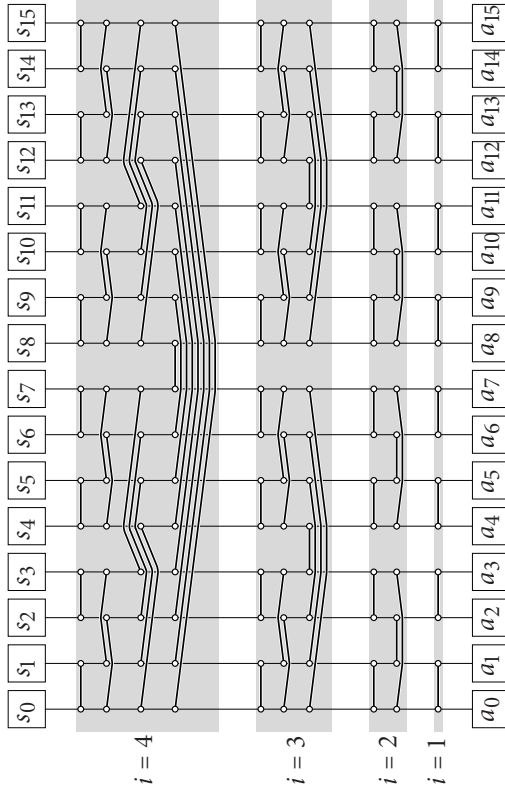
## Sortowanie

Kolejne zadanie, które chciałoby się rozwiązać algorytmem równoległym, to **sortowanie**: dany jest ciąg  $n$  elementów zbioru, w którym jest określony porządek liniowy. Elementy należy poprzestawiać tak, aby powstał ciąg niemalejący.

Mamy do dyspozycji komparator; trzeba go wywołać dla odpowiednich par elementów tak, aby na końcu otrzymać poprawny wynik sortowania niezależnie od początkowej permutacji danych elementów. Komparatory zorganizujemy w **sieć sortującą** — w kolejnych krokach  $\lfloor n/2 \rfloor$  komparatorów będzie jednocześnie porządkować  $\lfloor n/2 \rfloor$  par elementów.

Zastosowana strategia to **sortowanie przez scalanie**. W  $k$ -tym etapie (licząc od  $k = 1$ ) mamy posortowane kolejne podciągi o długości  $2^{k-1}$ . Kolejne pary takich podciągów scalimy w podciągi posortowane o długości  $2^k$ . Jeśli liczba  $n$  nie jest całkowitą potęgą dwójki, to ostatni podciąg może być krótszy. Wynik otrzymamy po  $\lceil \log_2 n \rceil$  etapach scalania.

22



Sortowanie ciągu o długości  $n$  składa się z  $\lceil \log_2 n \rceil$  etapów;  $k$ -ty etap składa się z  $k$  kroków, a więc całość wymaga wykonania  $\frac{1}{2} \lceil \log_2 n \rceil (\lceil \log_2 n \rceil + 1)$  kroków.

24

**Ciąg bitoniczny** jest połączeniem dwóch ciągów monotonicznych — nierosnącego i niemalejącego, o dowolnych długościach (któryś z nich może być nawet pusty).

Okazuje się (dowód w książce Cormena i in.), że jeśli podzielimy ciąg bitoniczny o parzystej długości na połowy, a potem zastosujemy komparator do par elementów ciągu bitonicznego —  $i$ -ty w pierwszej połowie z  $i$ -tym w drugiej dla każdego  $i$ , to połowy otrzymanego ciągu będą ciągami bitonicznymi, przy czym wszystkie elementy pierwszej połowy będą mniejsze lub równe elementom drugiej połowy.

Zatem, mając ciąg bitoniczny o długości  $2^k$ , dla  $k > 0$ , możemy go zamienić na dwa ciągi bitoniczne o długości  $2^{k-1}$  i dalej, indukcyjnie, na  $2^k$  ciągów bitonicznych o długości 1 — to razem da ciąg posortowany. Takie sortowanie ciągu bitonicznego, wykonane w  $k$  krokach, nazywa się **czyszczeniem**.

Jeśli teraz mamy scalić dwa sąsiednie ciągi posortowane o długości  $2^{k-1}$ , to zauważamy, że po odwróceniu kolejności elementów drugiego ciągu powstaje ciąg bitoniczny o długości  $2^k$ . Wystarczy go tylko wyczyścić.

23

Procedura main szadera sieci sortującej:

```
GLSL
1: #version 450 core
2:
3: #define GROUP_SIZE 64
4: layout (local_size_x = GROUP_SIZE) in;
5:
6: uniform uint n, h;
7: uniform bool reverse;
8:
9: void CompSwap ( uint i, uint j );
10:
11: void main ( void )
12: {
13:     uint iid, i, j, l, h2;
14:
15:     iid = uint ( gl_GlobalInvocationID.x );
16:     h2 = h >> 1; l = iid % h2; iid /= h2;
17:     i = iid*h2+1; j = reverse ? (iid+1)*h-1-1 : i+h2;
18:     if ( j < n )
19:         CompSwap ( i, j );
20: } /*main*/
```

25

Procedura sortowania w C niestety na jednym slajdzie się nie mieści.

```
C
1: void GPUWetSort ( GLuint n, GLuint n0, GLuint dbuf )
2: {
3:     GLuint steps, nn, h, h2, h4, k, kk, gsize, i;
4:
5:     if ( n < 2 )
6:         return;
7:     glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, dbuf );
8:     glUseProgram ( program_id );
9:     for ( nn = n-1, steps = 0; nn; nn >>= 1, steps ++ )
10: ;
11:     nn = 1 << steps;
12:     glUniform1ui ( uloc[0], n );
13:     glUniform1ui ( uloc[1], n0 );
14:     glUniform1ui ( uloc[2], false ); /* uniform reverse = false; */
15:     glUniform1ui ( uloc[3], 2 ); /* uniform h = 2; */
16:     glDispatchCompute ( gsize = (nn/2+1)lgsize[0]-1, 1, 1 );
17:     glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
```

26

## Mnożenie macierzy rzadkiej przez wektor

Macierz rzadka  $m \times n$  jest to (duża) macierz, której większość współczynników jest równa 0. Liczba niezerowych współczynników może być rzędu  $O(m+n)$  lub może to być kilka lub kilkanaście procent liczby  $mn$ .

Tu zajmują się rzadkimi macierzami nieregularnymi, których niezerowe współczynniki są rozmieszczone w dowolnych miejscach. Użyję oszczędnej reprezentacji takiej macierzy w postaci wykazu miejsc, w których są niezerowe współczynniki.

Mając daną macierz rzadką  $A$  i wektor  $x \in \mathbb{R}^n$ , chcę obliczyć wektor  $y = Ax$ . Wiersze i kolumny numeruję od 0. Prowadząc obliczenia według wzoru

$$y_i = \sum_{j=0}^{n-1} a_{ij}x_j$$

GPU obliczy i zsumuje tylko te składniki, w których  $a_{ij} \neq 0$ .

28

```

18: for ( i = 1, h2 = 2, h = 4, k = nn >> 2;
19:     i < steps;
20:     i++, h2 = h, h <= 1, k >>= 1 ) {
21:     glUniform1ui ( uloc[2], true ); /* uniform reverse = true; */
22:     glUniform1ui ( uloc[3], h );
23:     glDispatchCompute ( gsize, 1, 1 );
24:     glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
25:     glUniform1ui ( uloc[2], false ); /* uniform reverse = false; */
26:     for ( h2 = h >> 1, h4 = h2 >> 1, kk = k << 1;
27:         h2 > 1;
28:         h2 = h4, kk <= 1, h4 >>= 1 ) {
29:         glUniform1ui ( uloc[3], h2 ); /* uniform h = h2; */
30:         glDispatchCompute ( gsize, 1, 1 );
31:         glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
32:     }
33: }
34: ExitIfGLError ( "GPUWetSort" );
35: } /*GPUWetSort*/
```

27

Reprezentacja macierzy jest umieszczona w trzech tablicach: tablica  $r$  o długości  $m + 1$  zawiera liczby całkowite  $r_0, \dots, r_m$ , przy czym  $r_0 = 0$ , a dla  $i = 0, \dots, m - 1$  różnica  $r_{i+1} - r_i$  jest liczbą niezerowych współczynników w wierszu  $i$ -tym.

Liczba  $N = r_m$  jest liczbą wszystkich niezerowych współczynników macierzy  $A$ .

W tablicy  $c$  jest  $N$  liczb całkowitych,  $c_0, \dots, c_{N-1}$ . Liczby  $c_{r_i}, c_{r_i+1}, \dots, c_{r_{i+1}-1}$  są numerami tych kolumn, w których  $i$ -ty wiersz zawiera niezerowe współczynniki.

Same współczynniki (liczby rzeczywiste) są podane w tablicy  $a$  o długości  $N$ .

Na pozycjach  $r_i, r_i + 1, \dots, r_{i+1} - 1$  są przechowywane niezerowe współczynniki

z  $i$ -tego wiersza:  $a_{i,c_{r_i}}, a_{i,c_{r_i+1}}, \dots, a_{i,c_{r_{i+1}-1}}$ .

W pamięci GPU tablice  $r$  i  $c$  umieszczą jedną za drugą w jednym bloku magazynowym, a tablicę  $a$  w drugim bloku.

29

Szader obliczeniowy ma też dostęp do bloku z tablicą  $x$ , w której są podane współrzędne wektora  $x$ , i do bloku z tablicą  $y$ , do której ma być wpisany wynik.

Oprócz tego będą potrzebne dwie tablice robocze. W pierwszej z nich, o długości  $m$ , będą liczby niezerowych współczynników w kolejnych wierszach, tj. liczby składników do zsumowania.

W drugiej tablicy zmieści się  $N$  liczb rzeczywistych. Wątki szadera *jednocześnie* obliczą i zapamiętają w niej *wszystkie* iloczyny  $a_{i,j}x_j$ , a potem będą obliczały sumy — wszystkie równoległe i stosując algorytm sumowania parami.

30

Początek szadera obliczeniowego jest taki:

```

1: #version 450 core
2:
3: #define GROUP_SIZE 64
4: layout (local_size_x = GROUP_SIZE) in;
5:
6: layout (std430, binding = 0) buffer RowsCols { uint rc[]; } rc;
7: layout (std430, binding = 1) buffer Coeff { float a[]; } a;
8: layout (std430, binding = 2) buffer Xvec { float x[]; } x;
9: layout (std430, binding = 3) buffer Yvec { float y[]; } y;
10: layout (std430, binding = 4) buffer RowL { uint l[]; } lgt;
11: layout (std430, binding = 5) buffer Prod { float b[]; } b;
12:
13: uniform uint m, nnz, t;
14: uniform int stage;
15:
16: #define r(I) rc.rc[I]
17: #define c(I) rc.rc[m+1+(I)]

```

31

Macro  $r$  i  $c$  zapewniają dostęp do tablic  $r$  i  $c$ , parametr  $I$  jest indeksem.

Zmienne  $m$  i  $nnz$  mają wartości  $m$  i  $N$ .

Wartość zmiennej  $t$  określa liczbę składników w etapie sumowania iloczynów  $a_{i,j}x_j$ .

Zmienna jednolita  $stage$  służy do wyboru etapu obliczeń. Szader obliczeniowy będzie wywoływany za pośrednictwem makra EXECSTEP lub EXECSTAGE — to ostatnie nadaje wartość zmiennej  $stage$ ).

```

1: static GLuint program_id, uvloc[4]
2: static GLint lgsz[3];
3:
4: #define EXECSTEP(SIZE) \
5:   glDispatchCompute ( SIZE, 1, 1 ); \
6:   glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
7: #define EXECSTAGE(STAGE,SIZE) \
8:   { glUniform1i ( uvloc[0], STAGE ); EXECSTEP ( SIZE ) }

```

32



Procedura `main` szadera jest jedną wielką instrukcją przełącznika — wybierającego obliczenia dla kolejnych etapów.

```
1: void main ( void )
2: {
3:     uint iid, j, k, p;
4:
5:     iid = gl_GlobalInvocationID.x;
6:     switch ( stage ) {
```

Początkowe instrukcje procedury w C to przygotowanie — nadanie wartości zmiennym jednolitym `m` i `nnz`, przywiązanie buforów do punktów dowiązania bloków magazynowych i utworzenie buforów roboczych.

33

Etap 0 — znalezienie liczb niezerowych współczynników w poszczególnych wierszach:

```
21: EXECSTAGE ( 0, gsize = (m+lgsz[0]-1)/lgsz[0] )
22:
23: case 0:
24:     if ( iid < m )
25:         lgt.l[iid] = r(iid+1) - r(iid);
26:     return;
```

35

```
1: void GPU_MULTISPARSE_MATRIX_VECTORF ( GLuint m, GLuint n, GLuint nnz,
2:     GLuint *lmax, GLuint rowscols, GLuint coeff,
3:     GLuint x, GLuint y )
4: {
5:     GLuint auxb[2], gsize, t;
6:
7:     glUseProgram ( program_id );
8:     glUniform1i ( uvloc[1], m );
9:     glUniform1i ( uvloc[2], nnz );
10:    glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, rowscols );
11:    glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 1, coeff );
12:    glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 2, x );
13:    glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 3, y );
14:    glGenBuffers ( 2, auxb );
15:    glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 4, auxb[0] );
16:    glBindBufferData ( GL_SHADER_STORAGE_BUFFER, m*sizeof(GLuint),
17:        NULL, GL_DYNAMIC_DRAW );
18:    glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 5, auxb[1] );
19:    glBindBufferData ( GL_SHADER_STORAGE_BUFFER, nnz*sizeof(GLfloat),
20:        NULL, GL_DYNAMIC_DRAW );
```

34

Etap 1 — znalezienie największej liczby niezerowych współczynników w wierszu, wykonywany tylko jeśli zmienna `*lmax` ma wartość 0, w przeciwnym razie zakłada się, że ta liczba jest wartością zmiennej `*lmax`. Ponieważ to obliczenie psuje zawartość tablicy, etap 0 jest powtarzany.

```
22: if ( !*lmax ) {
23:     glUniform1i ( uvloc[0], 1 ); /* stage = 1 */
24:     for ( t = m; t > 1; t = (t+1)/2 ) {
25:         glUniform1i ( uvloc[3], t );
26:         EXECSTEP ( (t/2+lgsz[0]-1)/lgsz[0] )
27:     }
28:     glGetBufferSubData ( GL_SHADER_STORAGE_BUFFER, 0, sizeof(GLuint),
29:         lmax );
30:     EXECSTAGE ( 0, gsize )
31: }
32:
33: case 1:
34:     if ( ( j = iid + (t+1)/2 ) < t )
35:         if ( lgt.l[j] > lgt.l[iid] )
36:             lgt.l[iid] = lgt.l[j];
37:     return;
```

36

Etap 2 — obliczanie iloczynów  $a_{ij}x_j$ :

```

32: EXECSTAGE ( 2, (nnz+lgsize0[0]-1)/lgsize0[0] )
33: for ( t = *lmax; t > 1; t = (t+1)/2 ) {
34:     glUniform1ui ( uvloc[3], t );
35:     EXECSTEP ( (m*(t/2)+lgsize0[0]-1)/lgsize0[0] )
36: }
37:
20: case 3:
21:     p = t/2; j = iid % p; iid /= p;
22:     if ( iid < m ) {
23:         k = j + t - p;
24:         if ( k < lgt.l[iid] && k < t )
25:             b.b[r(iid)+j] += b.b[r(iid)+k];
26:     }
27:     return;

```

37

Etap 3 — sumowanie iloczynów parami. Wartość zmiennej  $t$  jest maksymalną liczbą składników, ale poszczególne sumy mogą mieć ich mniej. Dlatego dodawanie wykonuje się tylko dla takich par, w których drugi składnik ma numer mniejszy niż  $t$  i mniejszy niż  $r_{i+1} - r_i$  (po wykonaniu instrukcji w linii 21 zmienna  $iid$  ma wartość  $i$ ).

```

33: glUniform1i ( uvloc[0], 3 ); /* stage = 3 */
34: for ( t = *lmax; t > 1; t = (t+1)/2 ) {
35:     glUniform1ui ( uvloc[3], t );
36:     EXECSTEP ( (m*(t/2)+lgsize0[0]-1)/lgsize0[0] )
37: }
20: case 3:
21:     p = t/2; j = iid % p; iid /= p;
22:     if ( iid < m ) {
23:         k = j + t - p;
24:         if ( k < lgt.l[iid] && k < t )
25:             b.b[r(iid)+j] += b.b[r(iid)+k];
26:     }
27:     return;

```

38

Etap 4 — kopiowanie obliczonych współrzędnych wektora  $y$  z tablicy roboczej do tablicy  $y$ . Jeśli  $i$ -ty wiersz macierzy  $A$  jest zerowy, to następuje przypisanie  $y_i = 0$ .

```

38: EXECSTAGE ( 4, gsize )
39: glDeleteBuffers ( 2, auxb );
40: ExitIfGLError ( "GPUmultSparseMatrixVector" );
41: } /*GPUmultSparseMatrixVector*/
28: case 4:
29:     if ( iid < m )
30:         y.y[iid] = lgt.l[iid] > 0 ? b.b[r(iid)] : 0.0;
31:     return;
32: default:
33:     return;
34: }
35: } /*main*/

```

Na końcu procedura w C sprząta, tj. likwiduje buforzy robocze.

39

**Uwaga:** Jeśli macierz jest symetryczna, to można ją reprezentować w jeszcze mniejszej ilości pamięci, przechowując tylko dolny (albo górny) trójkąt. Niestety, taka oszczędna reprezentacja znacznie skomplikowałaby i spowolniła algorytm mnożenia — iloczynny współczynników  $a_{ij}$  z górnego trójkąta z liczbami  $x_j$  trzeba by sortować względem  $i$  przed sumowaniem.

40

## Mnożenie macierzy rzadkich

Spodziewając się, że iloczyn macierzy rzadkich,  $A$  i  $B$  o wymiarach  $m \times n$  i  $n \times l$ , też jest macierzą rzadką, chcemy obliczyć ten iloczyn możliwie małym kosztem, przy czym dla wszystkich macierzy użyjemy reprezentacji opisanej wcześniej. Trzeba znaleźć miejsca, w których macierz  $C = AB$  ma niezerowe współczynniki.

$i$ -ty wiersz macierzy  $C$  jest sumą wierszy macierzy  $B$  pomnożonych przez kolejne współczynniki z  $i$ -tego wiersza macierzy  $A$ ; można pomnożyć i zsumować tylko te wiersze, które są mnożone przez niezerowe współczynniki. Dalej, niezerowe współczynniki w  $i$ -tym wierszu macierzy  $C$  mogą być tylko w tych kolumnach, w których wybrane wiersze macierzy  $B$  mają niezerowe współczynniki.

41

Implementacja algorytmu mnożenia macierzy rzadkich składa się z procedury działającej na CPU i wywołanego przez nią programu szaderek, zawierającego szader obliczeniowy. W tej implementacji przydadzą się *wszystkie* algorytmy przedstawione wcześniej, do rozwiązania podzadania na kolejnych etapach.

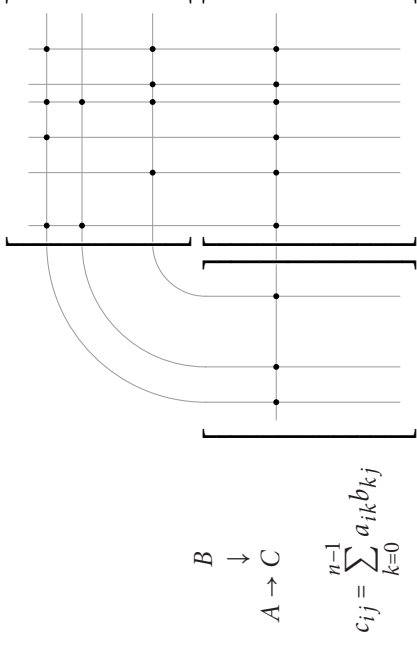
Początek szadera — punkty dowiązania buforów:

```

1: #version 450 core
2:
3: #define GROUP_SIZE 64
4:
5: layout (local_size_x = GROUP_SIZE) in;
6:
7: layout (std430, binding = 0) buffer RCA { uint rc[]; } rca;
8: layout (std430, binding = 1) buffer CoeffA { float a[]; } aa;
9: layout (std430, binding = 2) buffer RCB { uint rc[]; } rcb;
10: layout (std430, binding = 3) buffer CoeffB { float a[]; } ab;
11: layout (std430, binding = 4) buffer Auxb0 { uint a[]; } aux0;
12: layout (std430, binding = 5) buffer Auxb1 { float a[]; } aux1;
13: layout (std430, binding = 6) buffer Auxb2 { uint a[]; } aux2;
14:

```

43



Oczywiście, współczynnik  $c_{ij}$ , który jest sumą niezerowych iloczynów, może być zerem, ale to zignoruję. W razie potrzeby można zrobić dodatkowy etap obliczeń, który „wyczyści” otrzymaną macierz  $C$  ze współczynników zerowych lub mających wartość bezwzględną na poziomie błędów zaokrążeń.

42

Każda z macierzy  $A$  i  $B$  jest podana w dwóch blokach, przy czym dostęp do tablic  $r$  i  $c$  jest realizowany za pomocą odpowiednich makr. Są też trzy bloki robocze o nazwach `aux0`, `aux1` i `aux2` zawierające tablice, do których dostęp też jest realizowany przez makra objaśnione dalej.

GLSL

```

15: uniform int stage;
16: uniform uint prN0, prN, prStep;
17: uniform uint ma, mnza, mb, nprod, nnzc, h, tablgt;
18: uniform bool reverse;
19:
20: #define ra(I) rca.rc[I]
21: #define ca(I) rca.rc[ma+1+(I)]
22: #define rb(I) rcb.rc[I]
23: #define cb(I) rcb.rc[mb+1+(I)]
24: #define pairi(I) aux2.a[2*(I)]
25: #define pairj(I) aux2.a[2*(I)+1]
26: #define tab1(I) aux0.a[I]
27: #define tab2(I) aux0.a[tablgt+(I)]
28: #define tab3(I) aux0.a[tablgt+tablgt+(I)]
29:

```

44

Procedura `main` szadera jest przełącznikiem wybierającym instrukcje do wykonania na etapie określonym przez wartość zmiennej jednolitej `stage`. Jeśli jest to 0, to należy obliczyć ciąg sum prefiksowych w tablicy `aux0.a`, od miejsca `n0` do `n0+n-1`.

```

30: void iPrefixSum ( uint i ) { ... } /* procedura opisana wcześniej */
31:
32: void main ( void )
33: {
34:     uint i, j, k, l, m, p, q;
35:
36:     i = gl_GlobalInvocationID.x;
37:     switch ( stage ) {
38:     case 0:
39:         iPrefixSum ( i );
40:         return;

```

45

Procedura w C posługuje się takimi samymi makrami `EXECSTEP` i `EXECSTAGE` jak procedura mnożenia macierzy przez wektor. Parametry wejściowe opisują macierze  $A$  i  $B$ , a parametry wyjściowe są wskaźnikami zmiennych, którym będą przypisane liczba niezerowych współczynników oraz identyfikatory buforów (zarezerwowanych przez tę procedurę) z reprezentacją macierzy  $C$ .

```

1: static GLuint shader_id, program_id;
2: static GLuint uvloc[12];
3: static GLint lgsze[3];
4:
5: #define EXECSTEP(SIZE) \
6:     glDispatchCompute ( SIZE, 1, 1 ); \
7:     glMemoryBarrier ( GL_SHADER_STORAGE_BARRIER_BIT );
8: #define EXECSTAGE(STAGE,SIZE) \
9:     { glUniform1i ( uvloc[0], STAGE ); EXECSTEP ( SIZE ) }

46: char GPU_MULTISPARSE_MATRICESF ( GLuint m, GLuint n, GLuint l,
47:     GLuint nnza, GLuint rca, GLuint ca,
48:     GLuint nnzb, GLuint rcb, GLuint cb,
49:     GLuint *nnzc, GLuint *rcc, GLuint *cc )
50: {

```

46

W etapie 1 jest obliczana liczba  $P$  iloczynów niezerowych współczynników — dla każdego współczynnika  $a_{ij}$  bierzemy liczbę niezerowych współczynników w  $j$ -tym wierszu macierzy  $B$  i obliczamy sumę tych liczb, następnie odczytujemy ją z bufora. Jeśli suma jest zerem, to koniec obliczeń — macierz  $C$  jest zerowa.

```

65: EXECSTAGE ( 1, (nnza+lgsze[0]-1)/lgsze[0] )
66: iPrefixSum ( 1, nnza );
67: glGetBufferSubData ( GL_SHADER_STORAGE_BUFFER,
68:     nnza*sizeof(GLuint), sizeof(GLuint), &nprod );
69: if ( nprod ) {
70:     glDeleteBuffers ( 5, auxb );
71:     *nnzc = *rcc = *cc = 0;
72:     return false;
73: }

56: case 1:
57:     if ( i == 0 ) aux0.a[i] = 0;
58:     if ( i < nnza ) {
59:         j = ca(i);
60:         aux0.a[i+1] = rb(j+1)-rb(j);
61:     }
62:     return;

```

48

Bufory z danymi macierzami są przywoływane, zmienne jednolite `m`, `nnza` i `n` mają przypisywane wartości, a następnie procedura rezerwuje 5 buforów — dwa z nich pomniejszą wynik, a pozostałe trzy będą robocze. W pierwszym z nich ma być przechowywane  $N_A + 1$  liczb całkowitych.

```

51: GLuint auxb[5], nprod, _nnzc, maxnt, tablgt, i;
52:
53: glUseProgram ( program_id );
54: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, rca );
55: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 1, ca );
56: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 2, rcb );
57: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 3, cb );
58: glUniform1ui ( uloc[4], m );
59: glUniform1ui ( uloc[5], nnza );
60: glUniform1ui ( uloc[6], n );
61: glGenBuffers ( 5, auxb );
62: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 4, auxb[0] );
63: glBindBufferData ( GL_SHADER_STORAGE_BUFFER,
64:     (nnza+1)*sizeof(GLuint), NULL, GL_DYNAMIC_DRAW );

```

47

Rezerwujemy pozostałe dwa bufor robocze — pierwszy o pojemności  $2P$  liczb całkowitych, a drugi przeznaczony na  $P$  liczb zmiennopozycyjnych. Dla każdego iloczynu  $a_{ik}b_{kj}$  trzeba zapamiętać trójkę  $(i, j, a_{ik}b_{kj})$  — dostęp do elementów  $i, j$  trójki zapewniamą makra `pairi` i `pairj`.

```

74: glUniform1ui ( uLoc[7], nprod );
75: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 6, auxb[1] );
76: glBindBufferData ( GL_SHADER_STORAGE_BUFFER,
77:   2*nprod*sizeof(GLuint), NULL, GL_DYNAMIC_DRAW );
78: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 5, auxb[2] );
79: glBindBufferData ( GL_SHADER_STORAGE_BUFFER,
80:   nprod*sizeof(GLfloat), NULL, GL_DYNAMIC_DRAW );

```

W etapie 2 obliczane są iloczynny i zapamiętywane trójki, które je opisują, ale wcześniej shader musi znaleźć, dla każdego iloczynu, liczby  $i, j$  oraz wybrać z tablic współczynników do pomnożenia.

49

Robi się to metodą binarnego wyszukiwania — zmienna  $p$  otrzymuje indeks współczynnika  $a_{ik}$ , a zmienna  $q$  indeks współczynnika  $b_{kj}$  w odpowiedniej tablicy.

```

81: EXECSTAGE ( 2, (nprod+lgsize[0]-1)/lgsize[0] )
      GLSL
82: case 2:
83:   if ( i < nprod ) {
84:     for ( p = 0, k = mnza; k-p > 1; ) {
85:       l = p + (k-p)/2;
86:       if ( i >= aux0.a[l] ) p = l; else k = l;
87:     }
88:     for ( j = 0, k = ma; k-j > 1; ) {
89:       l = j + (k-j)/2;
90:       if ( p >= ra(l) ) j = l; else k = l;
91:     }
92:     pairi(i) = j;
93:     q = rb(ca(p))+i-aux0.a[p];
94:     pairj(j) = cb(q);
95:     aux1.a[i] = aa.a[p]*ab.a[q];
96:   }
97:   return;

```

50

Kolejne trzy etapy mają na celu posortowanie trójek tak, aby iloczynny do zsumowania znalazły się obok siebie. Trójki są już posortowane w kolejności niemalejących indeksów  $i$  — trzeba posortować niemalejąco względem  $j$  podciągów trójek o tych samych indeksach  $i$ .

W etapie 3, za pomocą binarnego wyszukiwania znajdowane są początki podciągów do posortowania.

```

88: EXECSTAGE ( 3, (m+lgsize[0]-1)/lgsize[0] )
      GLSL
89: case 3:
90:   if ( i == 0 ) { tab1(0) = 0; tab1(ma) = nprod; }
91:   else if ( i < ma ) {
92:     for ( j = 0, k = nprod; k-j > 1; ) {
93:       l = j + (k-j)/2;
94:       if ( pairi(l) < i ) j = l; else k = l;
95:     }
96:     tab1(i) = k;
97:   }
98:   return;

```

52

Następuje realokacja bufora na blok `aux0`; musi on teraz pomieścić trzy tablice o długościach  $T, T$  i  $N_C$ , gdzie  $T = \max(P, m) + 1$ . Odtąd dozwolone jest użycie makr `tab1`, `tab2` i `tab3` — używana w nich zmienna `tabLgt` otrzymuje wartość  $T$ . Liczba  $N_C$  nierzerowych współczynników macierzy  $C$  nie jest jeszcze znana, ale jest  $N_C \leq P$ .

```

82: glUniform1ui ( uLoc[11], tabLgt = nprod > m ? nprod+1 : m+1 );
83: glDeleteBuffers ( 1, &auxb[0] );
84: glGenBuffers ( 1, &auxb[0] );
85: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 4, auxb[0] );
86: glBindBufferData ( GL_SHADER_STORAGE_BUFFER,
87:   (3*tabLgt-1)*sizeof(GLuint), NULL, GL_DYNAMIC_DRAW );

```

51

Etap 4 oblicza długości tych podciągów, a w etapie 5 jest znajdowana największa długość (odczytywana z początku tablicy tab2 do zmiennej maxnt).

```

C
89: EXECSTAGE ( 4, (m+lgsize[0]-1)/lgsize[0] );
90: glUniform1i ( uloc[0], 5 ); /* uniform stage = 5; */
91: for ( i = m; i > 1; i = (i+1)/2 ) {
92:   glUniform1ui ( uloc[2], i );
93:   EXECSTEP ( (i/2+lgsize[0]-1)/lgsize[0] );
94: }
95: glGetBufferSubData ( GL_SHADER_STORAGE_BUFFER,
96:   tab1gt*sizeof(GLuint), sizeof(GLuint), &maxnt );
GLSL
89: case 4:
90:   if ( i < ma )
91:     tab2(i) = tab1(i+1)-tab1(i);
92:   return;
93: case 5:
94:   if ( (j = i+(prN+1)/2) < prN ) {
95:     if ( tab2(i) < tab2(j) )
96:       tab2(i) = tab2(j);
97:   }
98:   return;

```

53

W etapie 7 do tablicy tab1 są wpisywane jedynki — dla każdego pierwszego iloczynu w podciągu, który należy zsumować i zera dla pozostałych. Liczba tych jedynek jest liczbą  $N_C$  niezerowych współczynników macierzy C.

```

C
97: NetSort ( m, maxnt );
98: EXECSTAGE ( 7, (nprod+lgsize[0]-1)/lgsize[0] )
99: iPrefixSum ( 1, nprod );
100: glGetBufferSubData ( GL_SHADER_STORAGE_BUFFER,
101:   nprod*sizeof(GLuint), sizeof(GLuint), &nnzc );
GLSL
99: case 6:
100:   if ( (j = i / prN0) < ma )
101:     SortIt ( j, i % prN0 );
102:   return;
103: case 7:
104:   if ( i == 0 )
105:     { tab1(0) = 0; tab1(1) = 1; }
106:   else if ( i < nprod )
107:     tab1(i+1) = pairi(i-1) != pairi(i) || pairj(i-1) != pairj(i) ?
108:       1 : 0;
109:   return;

```

55

Kolejny etap, to sortowanie trójek — nie zamieszczam procedury NetSort, która realizuje algorytm opisany wcześniej. Ta implementacja sortuje wszystkie ciągi jednocześnie. Na podstawie numeru wątku w globalnej grupie roboczej szader ustala numer ns podciągu, który sortuje i numer np pary w tym podciągu, którą ma uporządkować. Procedura SortIt oblicza indeksy elementów pary i wykonuje zadanie komparatora:

```

GLSL
32: void SortIt ( uint ns, uint np )
33: {
34:   uint i, j, l, h2;
35:   float x;
36:
37:   h2 = h >> 1; l = np % h2; np /= h2;
38:   i = tab1(ns)+np*h+1; j = reverse ? tab1(ns)+(np+1)*h-1-1 : i + h2;
39:   if ( j < tab1(ns+1) ) {
40:     if ( pairj(i) > pairj(j) ) {
41:       l = pairj(i); pairj(i) = pairj(j); pairj(j) = l;
42:       x = aux1.a[i]; aux1.a[i] = aux1.a[j]; aux1.a[j] = x;
43:     }
44:   }
45: } /*SortIt*/

```

54

Znając liczbę  $N_C$ , można zarezerwować bufor, w których zostanie umieszczona macierz C. Bufory są przywiązywane do punktów wcześniej używanych do dostępu do macierzy A.

```

C
102: glUniform1ui ( uloc[8], nnzc );
103: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 1, auxb[4] );
104: glBindBufferData ( GL_SHADER_STORAGE_BUFFER,
105:   nnzc*sizeof(GLfloat), NULL, GL_DYNAMIC_DRAW );
106: glBindBufferBase ( GL_SHADER_STORAGE_BUFFER, 0, auxb[3] );
107: glBindBufferData ( GL_SHADER_STORAGE_BUFFER,
108:   (m+1+nnzc)*sizeof(GLuint), NULL, GL_DYNAMIC_DRAW );

```

56

W etapie 8 do tablicy `c` są wpisywane numery kolumn, w których są niezerowe współczynniki macierzy `C`. Do tablicy `tab2` jest wpisywany ciąg liczb — indeksów początków podciągów iloczynów do zsumowania. Różnice obliczane w etapie 9 są długościami podciągów.

```

109: EXECSTAGE ( 8, (nprod+lgsz[0])/lgsz[0] )
110: EXECSTAGE ( 9, (_nmzc+lgsz[0]-1)/lgsz[0] )
111: case 8:
112:   if ( i <= nprod ) {
113:     if ( i == nprod )
114:       tab2(nmzc) = nprod;
115:     else if ( tab1(i+1) > tab1(i) ) {
116:       ca(tab1(i)) = pairj(i);
117:       tab2(tab1(i)) = i;
118:     }
119:   }
120:   return;
121: case 9:
122:   if ( i < nmzc )
123:     tab3(i) = tab2(i+1) - tab2(i);
124:   return;

```

57

Etap 11 to sumowanie parami składników.

```

120: glUniform1i ( uloc[0], 11 ); /* uniform stage = 11; */
121: for ( i = maxnt; i > 1; i = (i+1)/2 ) {
122:   glUniform1ui ( uloc[2], i ); /* uniform prN = i; */
123:   EXECSTEP ( (_nmzc*(i/2)+lgsz[0]-1)/lgsz[0] )
124: }
125: case 11:
126:   p = prN/2; j = i % p; i /= p;
127:   if ( i < nmzc ) {
128:     k = j + prN - p;
129:     if ( k < tab3(i) && k < prN )
130:       aux1.a[tab2(i)+j] += aux1.a[tab2(i)+k];
131:   }
132:   return;

```

59

W etapie 10 znajdowana jest największa liczba składników do zsumowania. Potem jest odczytywana i etap 9 jest powtarzany, aby znów w tablicy `tab3` były te liczby składników.

```

111: glUniform1i ( uloc[0], 10 ); /* uniform stage = 10; */
112: for ( i = _nmzc; i > 1; i = (i+1)/2 ) {
113:   glUniform1ui ( uloc[2], i ); /* uniform prN = i; */
114:   EXECSTEP ( (i/2+lgsz[0]-1)/lgsz[0] )
115: }
116: glBindBuffer ( GL_SHADER_STORAGE_BUFFER, auxb[0] );
117: glGetBufferSubData ( GL_SHADER_STORAGE_BUFFER,
118:   (tablgt+tablgt)*sizeof(GLuint), sizeof(GLuint), &maxnt );
119: EXECSTAGE ( 9, (_nmzc+lgsz[0]-1)/lgsz[0] )
120: case 10:
121:   if ( ( j = i+(prN+1)/2 ) < prN ) {
122:     if ( tab3(i) < tab3(j) )
123:       tab3(i) = tab3(j);
124:   }
125:   return;

```

58

Etap 12 ma na celu przepisanie obliczonych sum, tj. współczynników macierzy `C`, na docelowe miejsca w buforze. Do tablicy pomocniczej `tab1` są wpisywane numery wierszy, w których są te współczynniki.

```

125: EXECSTAGE ( 12, (_nmzc+lgsz[0]-1)/lgsz[0] )
126: case 12:
127:   if ( i < nmzc ) {
128:     aa.a[i] = aux1.a[tab2(i)];
129:     tab1(i) = pairi(tab2(i));
130:   }
131:   return;

```

60

W etapie 13 wypełniana jest tablica  $r$  reprezentacji macierzy  $C$ , numerami miejsc, od których w tablicach  $c$  i  $a$  zaczynają się opisy kolejnych wierszy tej macierzy.

Tu znów konieczne okazało się wyszukiwanie binarne odpowiednich miejsc w tablicy  $tab1$ .

```

126: EXECSTAGE ( 13, (m+lgsize[0]-1)/lgsize[0] )
      |_____ C
      |_____ GlSl
144: case 13:
145:   if ( i == 0 ) { ra(0) = 0; ra(ma) = nnzc; }
146:   else if ( i < ma ) {
147:     for ( p = 0, k = nnzc; k-p > 1; ) {
148:       l = p + (k-p)/2;
149:       if ( i > tab1(1-1) ) p = 1; else k = 1;
150:     }
151:     ra(i) = p;
152:   }
153:   return;
154: default:
155:   return;
156: }
157: } /*main*/

```

61

Na końcu procedura w C musi już tylko przypisać odpowiednie wartości (liczbę  $N_C$  i identyfikatory buforów) zmiennym wskazywanym przez parametry i posprzątać (tj. zlikwidować bufor y z tablicami roboczymi).

```

127: *nnzc = _nnzc;
128: *rcc = auxb[3];
129: *cc = auxb[4];
130: glUseProgram ( 0 );
131: glDeleteBuffers ( 3, auxb );
132: ExitIfGLError ( "GPUmultSparseMatricesf" );
133: return true;
134: } /*GPUmultSparseMatricesf*/

```

62

*Dziękuję*

63