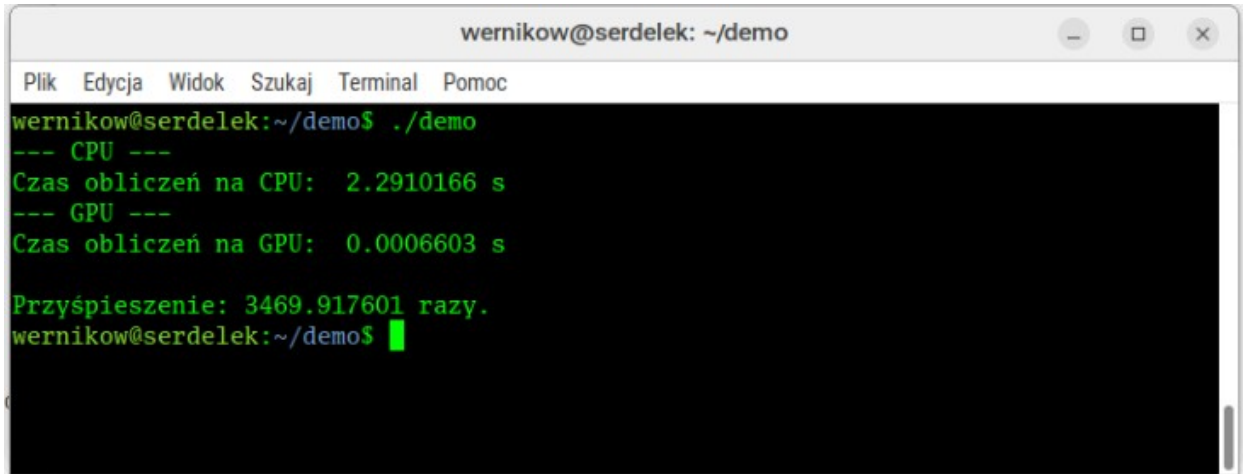


Profiler NVIDIA – szybki start

1. Dla celów niniejszego wprowadzenia użyto prostego kodu z jednym kernelem – treść pliku [demo.cu](#) znajduje się na końcu tego tekstu.
2. Plik skompilowano do pliku wykonalnego o nazwie `demo` i testowo uruchomiono.




```
wernikow@serdelek: ~/demo
Plik Edycja Widok Szukaj Terminal Pomoc
wernikow@serdelek:~/demo$ ./demo
--- CPU ---
Czas obliczeń na CPU:  2.2910166 s
--- GPU ---
Czas obliczeń na GPU:  0.0006603 s

Przyspieszenie: 3469.917601 razy.
wernikow@serdelek:~/demo$
```

W celu sprofilowania powyższego kodu należy kolejno:

1. uruchomić program *NVIDIA nSight Systems*
2. wykonać sekwencję *File* → *New Project* (tworzony jest projekt o nazwie *Project n*, gdzie *n* jest kolejnym numerem projektu)
3. w oknie po prawej stronie wybieramy *serdelek* jako *target for profiling*
4. pojawiają się ustawienia parametrów profilowania – dwa z nich podlegają obowiązkowemu wypełnieniu:
 1. *command line with arguments* – tutaj wprowadzamy pełną (kanoniczną!) ścieżkę do programu, który chcemy profilować; niedozwolone jest użycie znaku `~` (tylda) dla oznaczenia katalogu domowego;
 2. *working directory* – wypełnia się automatycznie na podstawie tekstu wprowadzonego do poprzedniego pola, ale może wymagać korekty po dokonaniu jakichkolwiek poprawek w pierwszym z pól



Specify process launch options below

Command line with arguments: [Edit arguments](#)
/home/wernikow/demo/demo

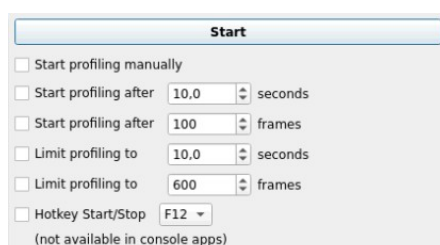
Working directory:
/home/wernikow/demo

5. następnie należy upewnić się, że wśród zbieranych statystyk znajduje się pozycja *Collect CUDA trace*



Collect CUDA trace

6. w tym momencie można już uruchomić profilowanie klikając w przycisk *Start*



Start

Start profiling manually

Start profiling after 10.0 seconds

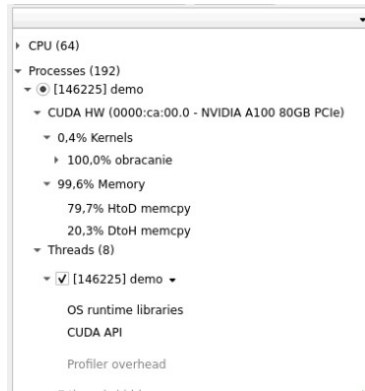
Start profiling after 100 frames

Limit profiling to 10.0 seconds

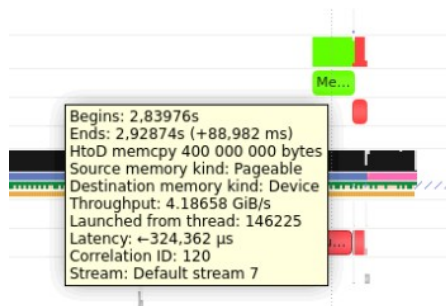
Limit profiling to 600 frames

Hotkey Start/Stop F12
(not available in console apps)

7. Po zakończeniu pracy programu i zignorowaniu ewentualnych ostrzeżeń pojawi się ekran wyników profilowania (może to zająć kilka sekund) – należy w nim rozwinąć gałęzie odpowiedzialne za prezentację aktywności urządzenia CUDA



8. W prawej części wyników ukaże się wykres czasowy, z którego można zgrubnie odczytać proporcje czasów zajętych przez transfery danych z/do urządzenia oraz wykonań kerneli. Dokładne wartości ukażą się po najechaniu kursorem myszy na wybrane odcinki na osi czasu.



Załącznik 1 – kod testowy użyty w przykładzie.

```
1  #include <stdio.h>
2  #include <math.h>
3  #include <assert.h>
4
5  #define N 5000000
6
7  struct WIERZCHOLEK {
8      float x,y;
9  } Figura[N];
10
11 struct timespec stoper_start(){
12     struct timespec start;
13     clock_gettime(CLOCK_MONOTONIC, &start);
14     return start;
15 }
16
17 long stoper_stop(struct timespec start){
18     struct timespec koniec;
19     clock_gettime(CLOCK_MONOTONIC, &koniec);
20     long nanosekundy = (koniec.tv_sec * 1E9 + koniec.tv_nsec) -
21     (start.tv_sec * 1E9 + start.tv_nsec);
22     return nanosekundy;
23 }
24
25 __global__ void obracanie(WIERZCHOLEK *W, float alfa) {
26     int i = blockDim.x * blockIdx.x + threadIdx.x;
27
28     if (i < N) {
29         float x = W[i].x * acos(alfa) - W[i].y * asin(alfa);
30         float y = W[i].x * asin(alfa) + W[i].y * acos(alfa);
31         W[i].x = x;
32         W[i].y = y;
33     }
34 }
35
36 int main(void) {
37     float          alfa = 3.1415;
38     struct         timespec stoper;
39     long           cpu, gpu;
40     int            ile_cudow;
41     WIERZCHOLEK   *d_W;
42
43     puts("--- CPU --- ");
44     stoper = stoper_start();
45     for(int i = 0; i < N; i++) {
46         float x = Figura[i].x * acos(alfa) - Figura[i].y * asin(alfa);
47         float y = Figura[i].x * asin(alfa) + Figura[i].y * acos(alfa);
48         Figura[i].x = x;
49         Figura[i].y = y;
50     }
51     cpu = stoper_stop(stoper);
52     printf("Czas obliczeń na CPU: %10.7f s\n", cpu / 1e9);
53     cudaGetDeviceCount(&ile_cudow);
54     if(ile_cudow == 0) {
55         perror("Nie ściemniaj - nie masz CUDY");
56         return 1;
57     }
58     puts("--- GPU ---");
59     cudaMalloc(&d_W, sizeof(Figura));
60     cudaMemcpy(d_W, Figura, sizeof(Figura), cudaMemcpyHostToDevice);
61     int watki_na_blok = 1024;
62     int bloki_na_siatke = (N + watki_na_blok - 1) / watki_na_blok;
63     stoper = stoper_start();
64     obracanie<<<bloki_na_siatke, watki_na_blok>>>(d_W, alfa);
65     cudaError_t err = cudaGetLastError();
66     assert(err == cudaSuccess);
67     cudaDeviceSynchronize();
68     gpu = stoper_stop(stoper);
69     cudaMemcpy(Figura, d_W, sizeof(Figura), cudaMemcpyDeviceToHost);
70     printf("Czas obliczeń na GPU: %10.7f s\n", gpu / 1e9);
71     cudaFree(d_W);
72     puts("");
73     printf("Przyspieszenie: %f razy.\n", (double)cpu/(double)gpu);
74     return 0;
75 }
```