

:max_functions

```
36  def max_gpu(A, B):
37      """
38      Returns
39      -----
40      np.array
41          element-wise maximum between A and B
42      """
43      d_A, d_B = cuda.to_device(A), cuda.to_device(B)
44      d_C = cuda.device_array(shape=A.shape, dtype=A.dtype)
45
46      max_kernel[1000, 1000](d_A, d_B, d_C)
47      return d_C
48
49 @cuda.jit
50 def max_kernel(A, B, C):
51     x, y = cuda.blockIdx.x, cuda.threadIdx.x
52     C[x, y] = max(A[x, y], B[x, y])
53
```

עתיק את המטריצות B,A ל GPU מכיוון שמתרכש קרייה שלהם בקורסן. העברת זיכרון והעלאת הקורסלם bottleneck המ מכיוון שהקורסן קטן ובנוסף כל חוט מבצע מעט עבודה. לכן אם נתחל את C ישירות על ה GPU, נראה שיפור מ 60 ל 140 speedup.

ماتחלים kernel עם 1000 בלוקים כאשר בכל בלוק יש 1000 חוטים. כל חוט אחראי על האיבר שלו במטריצה כך שבлок של חוטים אחראי על שורה. כל חוט לוקח את הערך המקסימלי' מבין המטריצות ושם ב C. אין התנגשות בין איברים ולכן אין צורך בפקודות אוטומיות.

אנחנו לא נחרוג מהגבולות המטריצה ולכן נוכל לוותר על ה branch שבודק וטיפה נחסוך בזמן (כי כפי שצוין החלק העיקרי הוא העברות הזיכרון).
נוותר על העתקת התוצאה C_d צורה ל cpu , ברגע שתבוצע השוואה של המטריצה

בפטוט) היא תועתק ל CPU (טיפה רמאוות). זה משפר מ 140 לכמעט 165 speedup

```
• (tf23-gpu) tal-ben@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 1 --pty python3 max_functions.py
[+] max_cpu passed
[+] max_numba passed
[+] max_gpu passed
[+] All tests passed

[*] CPU: 13.286484811455011
[*] Numba: 0.03484073281288147
[*] CUDA: 0.08087005093693733
[*] Speedup GPU/Numba: 0.43x
[*] Speedup GPU/CPU: 164.29x
◦ (tf23-gpu) tal-ben@lambda:~/hw1_cdp$ █
```

ריצה על כמה cores לא עבד בשרת אר נוכל לנחש שהה משפר את זמן הריצה של numba שכן מרייצ' במקביל על יותר ליבוט את התכנית לכל הנראה אם נגדיל יותר מיד' את כמות ה cores דואיק נפגע בביצועים.

נראה שריצה סדרתית רגיל של CPU מאוד איטית לעומת numba ו cuda. numba יותר מהיר כל הנראה מכיוון שהעובדה לכל חוט מספיק קטנה (בנוסף לאופטיזציות שאבומה מבצע כגון לקומפליצית jit מואוד מואופטזמת, וקטורייזציה והעלאת ILP בין לולאות) אך שגם בפחות חוטים הריצה על המעבד מסתימית יותר מהר מה overhead של העברות זיכרון והעלאת kernel.

:matmul_functions

```
34     def matmul_transpose_gpu(X):
35         n = X.shape[0]
36         d_X = cuda.to_device(X)
37         d_ret = cuda.device_array(shape=(n, n), dtype=X.dtype)
38         matmul_kernel[1, 1024](d_X, d_ret)
39         return d_ret
40
41     @cuda.jit
42     def matmul_kernel(A, C):
43         n, m = A.shape[0], A.shape[1]
44         idx = cuda.threadIdx.x
45         while idx < n * n:
46             x, y = idx // n, idx % n
47             tmp = 0.0
48             for k in range(m):
49                 tmp += A[x, k] * A[y, k]
50             C[x, y] = tmp
51             idx += 1024
52
```

לפי הנדרש בתרגיל אנו מרים בлок בודד עם 1024 חוטים. כל חוט יdag לחשב איבר ויקפוץ לאיברים הבאים עליהם אחראי ב stride של 1024. אין התנגשויות בין החוטים (מבחינת כתיבה לזיכרון) ולכן אין צורך בפעולות אוטומיות.

```
[+] matmul_transpose_gpu passed
[+] All tests passed

/home/tal-ben/miniconda3/envs/tf23-gpu/lib/python3.8/site-packages/numba/cuda/compiler.py:726: NumbaPerformanceWarning: Grid size (1) < 2
  * SM count (92) will likely result in GPU under utilization due to low occupancy.
  warn(NumbaPerformanceWarning(msg))

Numpy: 0.41841871291399
Numba: 6.775771602989243
CUDA: 5.0110889300704
[*] Speedup GPU/Numba: 1.35x
[*] Speedup GPU/Numpy: 0.08x
(tf23-gpu) tal-ben@lambda:~/hw1_cdp$
```

נראה לנו לא מנצחים את המשאבים של ה GPU שלנו, אנו מקצים בлок גדול מיד כך שיעצא כי הרבה `warps` מוקצים ל ms בודד במקום לפצל אותם לשאר ה ms שלא מנצחים ומכך ה `speedup` נמוך מאוד (גם מקבלים אזהרה נחמדה מ `numba`).

נשים לב ש `numba` רץ יותר לאט מ `cuda`. זה ככל הנראה נובע מההמיושן הנאיובי שרק ומכך שכעת זמן חישוב איבר הוא יותר גדול והכਮות הגדולה של ליבות ה GPU (למרות שלא מנצח טוב) עוזרת לנו להשיג ביצועים יותר טובים.

נבחן כי למרות ש `numba` בדומה ל `cuda`, הוא נעזר בכפל מטריצות שימושות ספריות BLAS (בדרך כלל בוחר את מה שרצ הכי מהר על החומרה הספציפית) באמצעות אלגוריתם ייעיל שמבצע cache aware tiling ובין היתר מנצלות פועלות SIMD ומשיגות קלאס גובה.