

# HW1 Report

תכנות מקבילי וմבוזר לעיבוד נתונים  
236370

מיכל עזרי 325719052  
גיא סודאי 214300550

## 1. גלגול

```
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 2 --pty python3 main.py
Epoch 1, accuracy 91.81 %.
Epoch 2, accuracy 95.24 %.
Epoch 3, accuracy 96.24 %.
Epoch 4, accuracy 96.62 %.
Epoch 5, accuracy 97.0 %.
Time to train using np.matmul: 23.219115257263184 seconds
Test Accuracy: 96.72%
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$
```

## חלק 2

הרצה עם ליבת אחת:

```
[michal.ozeri@132.68.39.159:22 - Bitvise xterm - michal.ozeri@lambda: ~/hw1_cdp]
(tf23-gpu) michal.ozeri@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: 12.886113958898932
[*] Numba: 0.03010775800794363
[*] CUDA: 0.08389257686212659
```

הסביר על שימוש סדרת max\_kernel ו-`max_gpu`:  
בפונקציה זו מוחלים את 2 מערכיו הקלט- A ו-B לסקום. יוצרים מערך פלט C ישירות על הסקום כדי לחסוך זמן. שכן יצרתו על הסקום הייתה גוררת את העברתו לסקום ללא צורך- אין בו מידע קודם שיש להעביר מהסקום.

לפי הדרישה, נקבע 1000 בלוקים כך שכל בלוק מכיל 1000 חוטים ונורץ את `max_kernel` עם `dev`, `dev_B`, `dev_C`. בחרורה מ-`max_kernel` נחזיר לנמען את מערך הפלט C נחזיר אותו ונוציאים. בפונקציה `max_kernel`, מכיוון שMOVEDת לנו שהמטריצות A ו-B הן מממדים  $1000 \times 1000$  ומכיוון שיש לנו 1000 בלוקים ובכל בלוק 1000 חוטים, נגדיר שכל חוט בכל בלוק אחראי לעדכן את התא  $C[bx, tx]$  כך שהוא הבלוק בו החוט נמצא ו- $ax$  זה מזהה החוט בבלוק. בכך נמנע גישות משותפת לתאים ובettaח נוכנות וניתוח מקסימלי של משאבים.

הרצה עם כמה ליבות:

```
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 2 --pty python3 max_functions.py
[+] max_cpu passed
[+] max_numba passed
[+] max_gpu passed
[+] All tests passed

[*] CPU: 12.814635635819286
[*] Numba: 0.029517433140426874
[*] CUDA: 0.08398842392489314
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 4 --pty python3 max_functions.py
[+] max_cpu passed
[+] max_numba passed
[+] max_gpu passed
[+] All tests passed

[*] CPU: 12.8831029930152
[*] Numba: 0.022172821685671806
[*] CUDA: 0.08307844772934914
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 8 --pty python3 max_functions.py
[+] max_cpu passed
[+] max_numba passed
[+] max_gpu passed
[+] All tests passed

[*] CPU: 12.748099931050092
[*] Numba: 0.012952613178640604
[*] CUDA: 0.08278299774974585
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ █
```

чисוב *speedup*:

<i>cores / speedup of _</i>	$\frac{max\_gpu}{max\_numba}$	$\frac{max\_gpu}{max\_cpu}$
1	$\frac{0.08389}{0.03010} \approx 2.787$	$\frac{0.08389}{12.886} \approx 6.51E - 3 \approx \frac{1}{153}$
2	$\frac{0.08398}{0.02951} \approx 2.846$	$\frac{0.08398}{12.814} \approx 6.55E - 3 \approx \frac{1}{152}$
4	$\frac{0.08307}{0.02217} \approx 3.747$	$\frac{0.08307}{12.883} = 6.45E - 3 \approx \frac{1}{155}$
8	$\frac{0.08278}{0.01295} \approx 6.392$	$\frac{0.08278}{12.748} \approx 6.49E - 3 \approx \frac{1}{154}$

נבחין כי כשהעלנו את מספר ה *cores*, ריצת *numba* *max* לא השתנהה. הדבר צפוי כיוון שהפעולה ממומשת באופן סדרתי לחלוטין, אך לא ניתן למקבל אותה עם מספר *cores*. לעומת זאת ניכול של ה *cores* הנוספים, והריצה על *core* יחיד, זמני הריצה זהים. גם ריצת *gpu* *max* לא השתנהה, זאת כיוון שהגדלת ה- *-cores* ב \**numa*\* לא תגדיל את כוח החישוב ב *gpu*. לכן, חישוב המקסימום המתבצע בסעופ יפועל באופן זהה. (יש פעולות הקורות גם על המעבד כמו שליחת הנתונים, ותקבלם, אך גם הן לא יכולות להפיק תועלות משמעותית מריבוי *cores*, מהסבירים דומים ל *numa*).

עד כה, הסבכנו שזמן הריצה של *gpu* *max* לא משתנים בהגדלת מספר ה *cores* וכן השיפוע נשאר קבוע. הסיבה לשיפוע זה היא יש המוני כניסה בלתי תלויות שיש לחשב, וזה מושתמש ב  $1000 \times 1000$  רכיבים מקבילים לחישוב (אחד לכל כניסה), ככה שההאצה בחישוב המקביל משמעותית בהרבה מהתקורתה של התקשרות עם ה *gpu*, והפעולה מהירה יותר ב *gpu*.

כאשר העלנו את מספר ה *cores* עבור *numba* *max*, זמן החישוב התקצר. הדבר הגיוני כיוון שבשימוש ב *numba* עם *prange*, הינו יוכל לחלק ביניהם את העבודה ה- *range*, וכך יותר עבודה נעשתה במקביל.

כיוון שגם הריצה של *gpu* *max* לא משתנה בהגדלת מספר ה *cores*, ואילו של *numba* כן, קורה השינוי ב *speedup* כשמגדילים את מספר ה *cores*. הסיבה שבאופן כללי *numba* יותר מהיר מה *gpu*, היא שתקורת העברת הנתונים בסעופ ותקבלם אחרת, לא משתמשת עבור מספר רכיבים מקבילים גדול יותר, ואילו מקבל על ה *gpu* אשר הזכרן כבר מוקן, אין את תקורת העברת זו, וה *cores* מהירים יותר (וגם מעטים יותר) עדיפה.

## חלק 3

הסביר על שימוש `matmul` ו-`matmul_transpose_gpu` בפונקציה `matmul_transpose_gpu` שולחים את מערך הקלט- $X$  לנוק. יוצרים מערך פلت C במימדים  $rows \times rows$  (נובע מכך שמיידי  $X$  הם  $cols \times rows$  ומיידי  $X^T$  הם  $rows \times cols$ ) ישרות על הנק  $\text{cp}$  היהgorת את העברתו לנוק ללא צורך. אין בו מידע קודם שיש להעביר מהסק. לפיה הדרישה, נקבע בלוק אחד המכיל 1024 חוטים ונರיץ את `matmul_kernel` עם `dev_X`, `dev_C`. בחרזה מ-`matmul_kernel` נחזיר לנק `cp` את מערך הפלט C נחזיר אותו ונסימם. בפונקציה `matmul_kernel`, מכיוון שיש לנו 1024 חוטים, נגדיר שכל חוט אחראי לעדכן את התאים  $C[i, j]$  כך ש-

$$(i * \text{num\_columns\_of\_C} + j) \bmod 1024 = tx$$

$$\text{iff}, \quad (i * \text{rows} + j) \bmod 1024 = tx$$

כך ש- $tx$  זה מזזה החוט בבלוק. כלומר, כל חוט יתחילה בעדכון התא  $C[tx // \text{rows}, tx \% \text{rows}]$  ויעדכו את התאים הבאים בקייזות של 1024. בכך נמנע גישות משותפת לתאים ונבטיח נכונות וኒzuל מקסימלי של משאבם (כלומר, חלוקה שווה בקריבוב של העומסים).

```
michal.ozeri@132.68.39.159:22 - Bitvise xterm - michal.ozeri@lambda: ~/hw1_cdp
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 1 --pty python3 matmul_functions.py
[+] matmul_transpose_trivial passed
[+] matmul_transpose_numba passed
[+] matmul_transpose_gpu passed
[+] All tests passed

Numpy: 0.4209525021724403
Numba: 5.80656470824033
CUDA: 6.987547372933477
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 2 --pty python3 matmul_functions.py
[+] matmul_transpose_trivial passed
[+] matmul_transpose_numba passed
[+] matmul_transpose_gpu passed
[+] All tests passed

Numpy: 0.4210456367582083
Numba: 5.805208188015968
CUDA: 6.983507165219635
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 4 --pty python3 matmul_functions.py
[+] matmul_transpose_trivial passed
[+] matmul_transpose_numba passed
[+] matmul_transpose_gpu passed
[+] All tests passed

Numpy: 0.2318826108239591
Numba: 3.120263021904975
CUDA: 6.987735772039741
(tf23-gpu) michal.ozeri@lambda:~/hw1_cdp$ srun --gres=gpu:1 -c 8 --pty python3 matmul_functions.py
[+] matmul_transpose_trivial passed
[+] matmul_transpose_numba passed
[+] matmul_transpose_gpu passed
[+] All tests passed

Numpy: 0.1492176498286426
Numba: 1.8676453148946166
CUDA: 6.988251808099449
```

נסביר את התוצאות:

:numpty

בядוע מומש אלגוריתם יעיל מאוד לחישוב מכפלת מטריצות (יותר יעיל מהימוש הטריוויאלי הלוקח  $O(n^3)$ )

מקומפל מראש, הידוע לנצל מספר cores (מקבל את החישוב).  
לכן ריצתו מהירה מאוד, ומשתפרת עם העלתה cores.

:numba

בnumba מומש האלגוריתם הטריוויאלי לחישוב מכפלת מטריצות (הלוקח  $O(n^3)$ ), הידוע לנצל מספר cores (מקבל את החישוב עם prange).

אין תקורת העברת נתונים כי המידע הדרוש כבר נמצא בזיכרון המעבד.  
לכן ריצתו משתפרת עם העלתה cores, אך איטית יותר מnumpty בגלל הבדלי מימוש האלגוריתם.

:gpus

ב gpus מומש האלגוריתם הטריוויאלי לחישוב מכפלת מטריצות (הלוקח  $O(n^3)$ ), ובנוסף לא מסוגל לנצל ריבוי cores ב gpu (רכיב חומרה אחר, והפעולות הרכוכות ב gpu כמו בהעברת המידע לא "נהנות" מריבוי cores).

בנוסף, יש תקורת העברת נתונים כי המידע הדרוש נמצא בזיכרון המעבד, ויש להעיבו לזכרון gpu לצורך החישוב, וכן להציג את התוצאה למעבד. פועלות החישוב עצמה מתבצעת באופן מקבילי על מספר גדול של threads, אך לא בצורה המשتمלת עבור תקורת gpu, בהשוואה להרצה מקבילתית על gpu כמו numba.  
לסיכום, ריצתו אינה משתפרת עם העלתה cores, וכן איטית יותר מ numba בגלל שהאלגוריתם, אך תקורת העברת הנתונים משמעויות יותר מאשר המיקובל.