# Hashtag Post Counting

01418343 Parallel Computing with CUDA
6610402141 Patiput Ukham

# About the project

Parallel Hashtag Counter using CUDA - เปรียบเทียบ
ประสิทธิภาพระหว่าง Sequential และ GPU Parallel Processing

# Dataset

# Idea

I don't feel good  #fb                                      1 :#fb
Sad end to the game  #Canucks                    1 :#Canucks
I can get back in time for #SNL                      1 :#SNL                                       count

                                                                                                     f  b

I   d o n ' t   f e e l   g o o d   # f b
1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9

# Sequential Solving

```c
void count_hashtag(char** h_str, int* h_len, int numstr, char hashtags[MAX_TAGS][MAX_TAGS_LEN], int tag_count[MAX_TAGS], int* tags_count)
{
    int cht = 0;
    for (int i = 0; i < numstr; i++)
    {
        char* s = h_str[i];
        for (int j = 0; s[j] != '\0'; j++)
        {
            if (s[j] == '#')
            {
                char hashtag[MAX_TAGS_LEN];
                int tag_len = 0;
                for (int k = j; s[k] != '\0'; k++)
                {
                    if (!isValidHashtagChar(s[k]) || (s[k] == '#' && k != j) || tag_len >= MAX_TAGS_LEN-1)
                    {
                        break;
                    }
                    hashtag[tag_len++] = s[k];
                }
                hashtag[tag_len] = '\0';

                // counting
                bool tag_exist = false;
                for(int l = 0; l < (*tags_count); l++)
                {
                    if (strcmp(hashtag, hashtags[l]) == 0)
                    {
                        tag_count[l]++;
                        tag_exist = true;
                        break;
                    }
                }

                // init new hashtag
                if (!tag_exist && (*tags_count) < MAX_TAGS)
                {
                    strcpy(hashtags[(*tags_count)], hashtag);
                    tag_count[(*tags_count)] = 1;
                    (*tags_count)++;
                }
            }
        }
    }
}
```

f e e l   g o o d   # f b
g a m e   # C a n u c k s

Time Complexity: O(n × m × k × t)
- n = จำนวนของข้อความ
- m = ความยาวเฉลี่ยของแต่ละบรรทัด
- k = จำนวน hashtag เฉลี่ยต่อบรรทัด
- t = จำนวน unique hashtags ที่เจอแล้ว

# Parallel optimize

```
__global__ void parallel_hashtag_count(char **d_str, int *d_len, int numstr, char (*d_hashtags)[MAX_TAGS_LEN], int *d_tag_count, int *d_tags_count)
{
    int idx = blockIdx.x * blockDim.x + threadIdx.x;

    if (idx < numstr)
    {
        char* s = d_str[idx];

        for (int i = 0; s[i] != '\0'; i++)
        {
            if (s[i] == '#')
            {
                char hashtag[MAX_TAGS_LEN];
                int tag_len = 0;
                for (int j = i+1; s[j] != '\0'; j++)
                {
                    if (!isValidHashtagChar(s[j]) || tag_len >= MAX_TAGS_LEN-1)
                    {
                        break;
                    }

                    hashtag[tag_len++] = s[j];
                }
                hashtag[tag_len] = '\0';

                // counting
                bool tag_exist = false;
                for (int k = 0; k < *d_tags_count; k++)
                {
                    bool match = true;
                    // strcmp
                    for (int l = 0; l < tag_len; l++)
                    {
                        if (hashtag[l] != d_hashtags[k][l])
                        {
                            match = false;
                            break;
                        }
                    }

                    // counting exists hashtag
                    if (match && d_hashtags[k][tag_len] == '\0')
                    {
                        atomicAdd(&d_tag_count[k], 1);
                        tag_exist = true;
                        break;
                    }
                }
```

t0 f e e l    g o o d    # f b

t1 g a m e    # C a n u c k s

tn

Time Complexity: $O(m \times k \times t / P)$
- m = ความยาวบรรทัด
- k = จำนวน hashtags ต่อบรรทัด
- t = จำนวน unique hashtags (linear search)
- P = จำนวน threads ที่ active พร้อมกัน

6

# Parallel optimize

```
70
71            // init new hashtag
72            if (!tag_exist)
73            {
74                int cur_len_tags = atomicAdd(d_tags_count, 1);
75                if (cur_len_tags < MAX_TAGS)
76                {
77                    for (int l = 0; l < tag_len; l++)
78                    {
79                        d_hashtags[cur_len_tags][l] = hashtag[l];
80                    }
81                    d_hashtags[cur_len_tags][tag_len] = '\0';
82                    d_tag_count[cur_len_tags] = 1;
83                }
84            }
85        }
86    }
87  }
88 }
89
```

t0   f e e l   g o o d     # f b

t1   g a m e     # C a n u c k s

tn

Time Complexity: O(m × k × t / P)
- m = ความยาวบรรทัด
- k = จำนวน hashtags ต่อบรรทัด
- t = จำนวน unique hashtags (linear search)
- P = จำนวน threads ที่ active พร้อมกัน

# Parallel optimize reduce

phase 1

Time Complexity: O(m x n / P)

```
1    __global__ void find_hashtag_positions(char* d_buffer, int buffer_size, int* d_hashtag_positions, int *d_hashtag_count)
2    {
3        int idx = blockIdx.x * blockDim.x + threadIdx.x;
4
5        if (idx < buffer_size && d_buffer[idx] == '#')
6        {
7            int cur_len_tag = atomicAdd(d_hashtag_count, 1);
8            d_hashtag_positions[cur_len_tag] = idx;
9        }
10   }
```

t6 -> checks '#' -> atomicAdd -> position[0] = 6
t20 -> checks '#' -> atomicAdd -> position[1] = 20
t40 -> checks '#' -> atomicAdd -> position[2] = 40

# Parallel optimize reduce

## phase 2

```
1   __global__ void parallel_hashtag_extracting(char* d_buffer, int buffer_size,
2                                                int *d_hashtag_positions,
3                                                int hashtag_count,
4                                                char (*d_hashtags)[MAX_TAGS_LEN])
5   {
6       int idx = blockIdx.x * blockDim.x + threadIdx.x;
7
8       if (idx < hashtag_count)
9       {
10          int start_position = d_hashtag_positions[idx];
11          int tag_len = 0;
12
13          for (int i = start_position+1; i < buffer_size && tag_len < MAX_TAGS_LEN-1; i++)
14          {
15              if (!isValidHashtagChar(d_buffer[i]) || d_buffer[i] == '\n')
16              {
17                  break;
18              }
19              d_hashtags[idx][tag_len++] = d_buffer[i];
20          }
21          d_hashtags[idx][tag_len] = '\0';
22      }
23  }
```

Time Complexity: O(L)
- L = avg hashtag length

Thread 0 → position[0] = 6
        -> Extract from index 7: "world"
        -> d_hashtags[0] = "world"

# Parallel optimize reduce

phase 3

Time Complexity: O(k × t × L / P)
- k = total hashtags
- t = unique hashtags
- L = avg hashtag length

Thread 0 -> "world"
    -> Not found -> Add to position 0
    -> unique_hashtags[0] = "world", count[0] = 1
Thread n -> "CUDA"
    -> Found at position m -> atomicAdd(&count[m])
    -> count[m]++

```cuda
__global__ void unique_count_hashtags(char (*d_hashtags)[MAX_TAGS_LEN], int hashtag_count,
                                      char (*d_unique_hashtags)[MAX_TAGS_LEN],
                                      int *d_unique_counts,
                                      int *d_unique_count)
{
    int idx = blockIdx.x * blockDim.x + threadIdx.x;

    if (idx < hashtag_count)
    {
        char *hashtag = d_hashtags[idx];
        int tag_len = 0;
        while (hashtag[tag_len] != '\0' && tag_len < MAX_TAGS_LEN)
        {
            tag_len++;
        }

        if (tag_len <= 1)
        {
            return;
        }

        // tag exist?
        bool tag_exist = false;
        for (int i = 0; i < *d_unique_count; i++)
        {
            bool match = true;
            for (int j = 0; j < tag_len; j++)
            {
                if (hashtag[j] != d_unique_hashtags[i][j])
                {
                    match = false;
                    break;
                }
            }

            if (match && d_unique_hashtags[i][tag_len] == '\0')
            {
                atomicAdd(&d_unique_counts[i], 1);
                tag_exist = true;
                break;
            }
        }

        // init new hashtag
        if (!tag_exist)
        {
            int cur_len_tag = atomicAdd(d_unique_count, 1);
            if (cur_len_tag < MAX_TAGS)
            {
                for (int j = 0; j < tag_len; j++)
                {
                    d_unique_hashtags[cur_len_tag][j] = hashtag[j];
                }
                d_unique_hashtags[cur_len_tag][tag_len] = '\0';
                d_unique_counts[cur_len_tag] = 1;
            }
        }
    }
}
```

# Result

process on 1.6m post with
cpu: i5 13500
GPU: Nvidia 3070Ti

| Algorithm | Time process | Cmp |
|---|---|---|
| Sequential | 476.734 ms | 1X |
| Parallel optimize | 181.837 ms | 2.6X |
| Parallel optimize reduce | 8.07261 ms | 59X |

# Thank you