



# 并行计算大作业

## ——并行优化矩阵乘法

朱桐 10175102111  
周亦然 10175102207

2019.12.27



# Outline

Outline

Introduction

Photo

Conclusion

References



# code: buffered input

```
inline char next_char()
{
    static char buf[1000000], *p1 = buf,
               *p2 = buf;
    return p1 == p2 && (p2 = (p1 = buf)
        + fread(buf, 1, 1000000, fin), p1 == p2)
        ? EOF : *p1++;
}
```



## code: buffered output

```
inline void flush() {  
    fwrite(buffer, 1, s-buffer, stdout);  
    s = buffer;  
    fflush(stdout);  
}  
  
inline void print(const char ch) {  
    // putchar(ch); return;  
    if (s-buffer>OutputBufferSize-2) flush();  
    *s++ = ch;  
}
```



## code: share memory, block partition

```
__shared__ ld c_a[max_shared_size];  
int index = blockDim.x * blockIdx.x + threadIdx.x;  
if (index >= an * bm) return;  
int st = min(index, addi) * (workload+1) +  
max(0, index - addi) * workload, ed =  
st + workload + (index < addi ? 1 : 0);  
int shareda = min(am, max_shared_size);  
for (int p=st; p<ed; ++p) {  
    // ...  
}
```



# code: share memory, block partition

```
for (int p=st; p<ed; ++p) {  
    int i = p / bm, j = p % bm;  
    if (p % bm == 0) {  
  
        for (int j=0; j<shareda; ++j) {  
            c_a[j] = d_a[i * am + j];  
        }  
        __syncthreads();  
    }  
    ...  
}
```



# code: share memory, block partition

```
for (int p=st; p<ed; ++p) {  
    int i = p / bm, j = p % bm;  
    if (p % bm == 0) {  
  
        for (int j=0; j<shareda; ++j) {  
            c_a[j] = d_a[i * am + j];  
        }  
        __syncthreads();  
    }  
    ...  
}
```



# code: async IO, creating streams

```
int st = 0, ed = n * m;  
// printf("st=%d ed=%d, a=%p\n", st, ed, a);  
cudaStream_t stream[2];  
int mask = 0;  
cudaStreamCreate(&stream[0]);  
cudaStreamCreate(&stream[1]);  
int size;
```





## code: async IO, creating streams

```
cudaStream_t mainstream;  
cudaStreamCreate(&mainstream);  
// ...  
copyMatrixAsync(h_a, d_a, an, am, mainstream);  
// ... read h_b  
copyMatrixAsync(h_b, d_b, bn, bm, mainstream);  
// ...  
matrixMult<<<grids, block_size, 0, mainstream>>>  
(d_a, d_b, d_c, an, bm, am);  
// ...  
handleCudaError(cudaStreamSynchronize(mainstream));  
// ...  
outputMatrixAsync(h_c, d_c, n, m);  
// pipeline: memcpy from device & output each rows
```



## code: async IO, pipelines

```
for (; st<ed; st+=size, mask^=1) {  
    size = min(chunk_size, ed - st);  
    handleCudaError(cudaMemcpyAsync(a + st,  
        d_a + st, size * sizeof(ld),  
        cudaMemcpyDeviceToHost, stream[mask]));  
    // exit(0);  
    if (st - chunk_size >= 0) {  
        // printf("%d %d\n",st-chunk_size, st);  
        handleCudaError(cudaStreamSynchronize(stream));  
        outputinterval(a, st-chunk_size, st);  
    }  
}
```



# Backgrounds

- 11111111111111
- 22222222222222
- 33333333333333
- 44444444444444



# My Photo



图: hahahaha...



# Emmm...

## Sequence Tagging Loss

$$\mathcal{L}_p = - \sum_{i=1}^S \sum_{j=1}^N p_{i,j} \log(\hat{p}_{i,j})$$

## Language Classifier Loss

$$\mathcal{L}_a = - \sum_{i=1}^S l_i \log(\hat{l}_i)$$

## Bidirectional Language Model Loss

$$\mathcal{L}_l = - \sum_{i=1}^S \sum_{j=1}^N \log(P(w_{j+1}|f_j)) + \log(P(w_{j-1}|b_j))$$



# References



Xuezhe Ma and Eduard Hovy. (2016).

## **End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF.**

In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, pages 1064–1074, Berlin, Germany, August 7-12, 2016.



Marek Rei. (2017).

## **Semi-supervised Multitask Learning for Sequence Labeling.**

In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, pages 2121–2130, Vancouver, Canada, July 30 - August 4, 2017.