



SITCON '21

Floyd-Warshall-GPU

- 觀察 Data dependency

$$f_k(i, j) = \min(f_{k-1}(i, j) , f_{k-1}(i, k) + f_{k-1}(k, j))$$



• 裡面兩層沒有關聯

●將裡面兩層迴圈平行化

```
1  for (int k = 0; k < n; k++) {
2      for (int i = 0; i < n; i++) {
3          for (int j = 0; j < n; j++) {
4              g[i][j] = min(g[i][j], g[i][k] + g[k][j]);
5          }
6      }
7  }
```

```
1  __global__ void gpu(int **g, int n, int k) {
2      int id = blockIdx.x * blockDim.x + threadIdx.x;
3      int i = id / n;
4      int j = id % n;
5      g[i][j] = min(g[i][j], g[i][k] + g[k][j]);
6  }
7
8  // bs * 1024 = n^2
9  for (int k = 0; k < n; k++) {
10     gpu<<<bs, 1024>>>(g, n, k);
11 }
```

Floyd-Warshall - GPU

- 觀察 Data dependency

$$f_k(i, j) = \min(f_{k-1}(i, j), f_{k-1}(i, k) + f_{k-1}(k, j))$$

- 裡面兩層 loop 沒有關聯
- 將裡面兩層迴圈平行化

```
1  for (int k = 0; k < n; k++) {
2      for (int i = 0; i < n; i++) {
3          for (int j = 0; j < n; j++) {
4              g[i][j] = min(g[i][j], g[i][k] + g[k][j]);
5          }
6      }
7  }
```

```
1  __global__ void gpu(int **g, int n, int k) {
2      int id = blockIdx.x * blockDim.x + threadIdx.x;
3      int i = id / n;
4      int j = id % n;
5      g[i][j] = min(g[i][j], g[i][k] + g[k][j]);
6  }
7
8  // bs * 1024 = n^2
9  for (int k = 0; k < n; k++) {
10     gpu<<<bs, 1024>>>(g, n, k);
11 }
```


Performance Issue

- 無連續讀取記憶體讀取

```
min(g[i][j], g[i][k] + g[k][j]);
```

- 無共用記憶體

- 效能瓶頸為記憶體頻寬

