

GPU H/W and S/W

Hyesoon Kim

Arithmetic Intensity

$$\text{Arithmetic Intensity (AI)} = \frac{\text{Number of floating-point operations (FLOPs)}}{\text{Bytes of memory transferred}}$$

Code example: CUDA K-mean

```
__global__ void assignClusters(const float *points, const float *centroids, int *assignments) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < n_points) {
        int best_cluster = 0;
        float min_dist = FLT_MAX;

        for (int c = 0; c < k; c++) {
            float dist = 0.0f;
            for (int d = 0; d < dim; d++) {
                float diff = points[idx * dim + d] - centroids[c * dim + d];
                dist += diff * diff; // 1 FLOP multiply + 1 FLOP add per dimension
            }
            if (dist < min_dist) {
                min_dist = dist;
                best_cluster = c;
            }
        }
        assignments[idx] = best_cluster;
    }
}
```

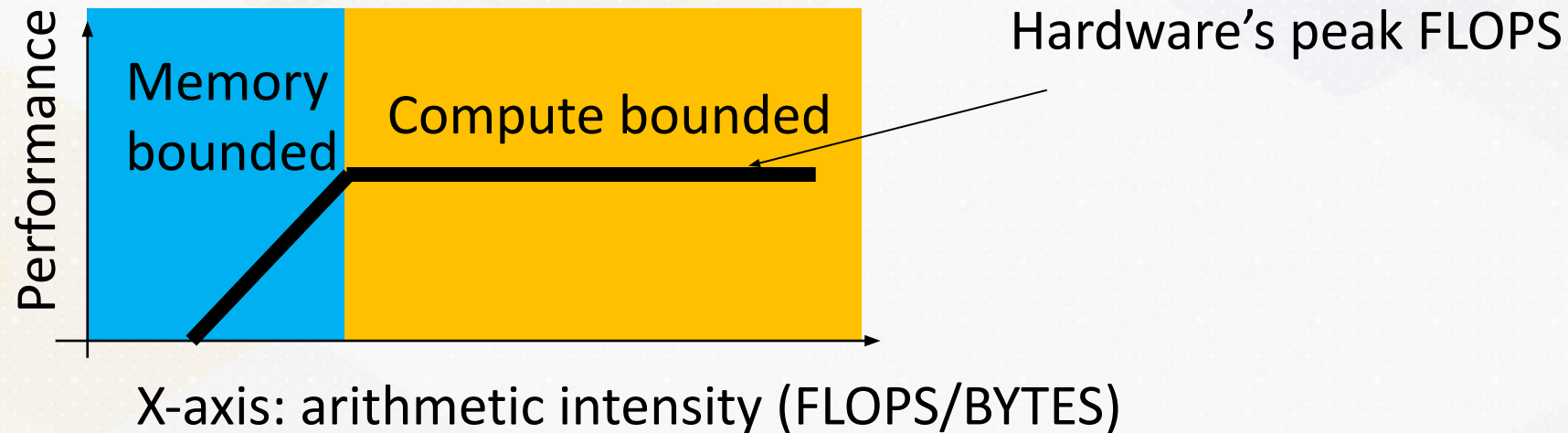
FLOPS = $\text{dim} \times (1 \text{ sub} + 1 \text{ mul} + 1 \text{ add}) = \text{dim} \times 3$
FLOPS per point = $k \text{ clusters} \times \text{dim} \times 3$
Total FLOPS. = $n * K * \text{dim} * 3$

Memory load:
Point data loads: $n * \text{dim}$ floats
Centroid data loads: $k * \text{dim}$ floats
Store: n
= $n * \text{dim} * 4 + n * 4$

- Can you find Arithmetic Intensity?

Code examples: (Perplexity AI, 2025)

Roofline Model



- A Visual performance model to determine whether an application (or a processor) is limited by the compute bandwidth or memory bandwidth

What to count for Memory Bytes

- Roofline model: use theoretical minimal memory bytes
- What if we count based on transactional size?
- Based on what to bring from memory or caches

```
#define N 4096
```

```
__global__ void matAdd(float *A, float *B, float *C) {  
    int idx = blockIdx.x * blockDim.x + threadIdx.x;  
    if (idx < N)  
        C[idx] = A[idx] + B[idx];  
}
```

```
#define N 4096
```

```
#define STRIDE 2
```

```
__global__ void matAdd(float *A, float *B, float *C) {  
    int idx = blockIdx.x * blockDim.x + threadIdx.x;  
    if (idx < N/STRIDE)  
        C[idx*STRIDE] = A[idx*STRIDE] + B[idx*STRIDE];  
}
```

Transactional size

- Between cache: cache block size
- From memory: minimum memory transaction size
- Significant performance impacts