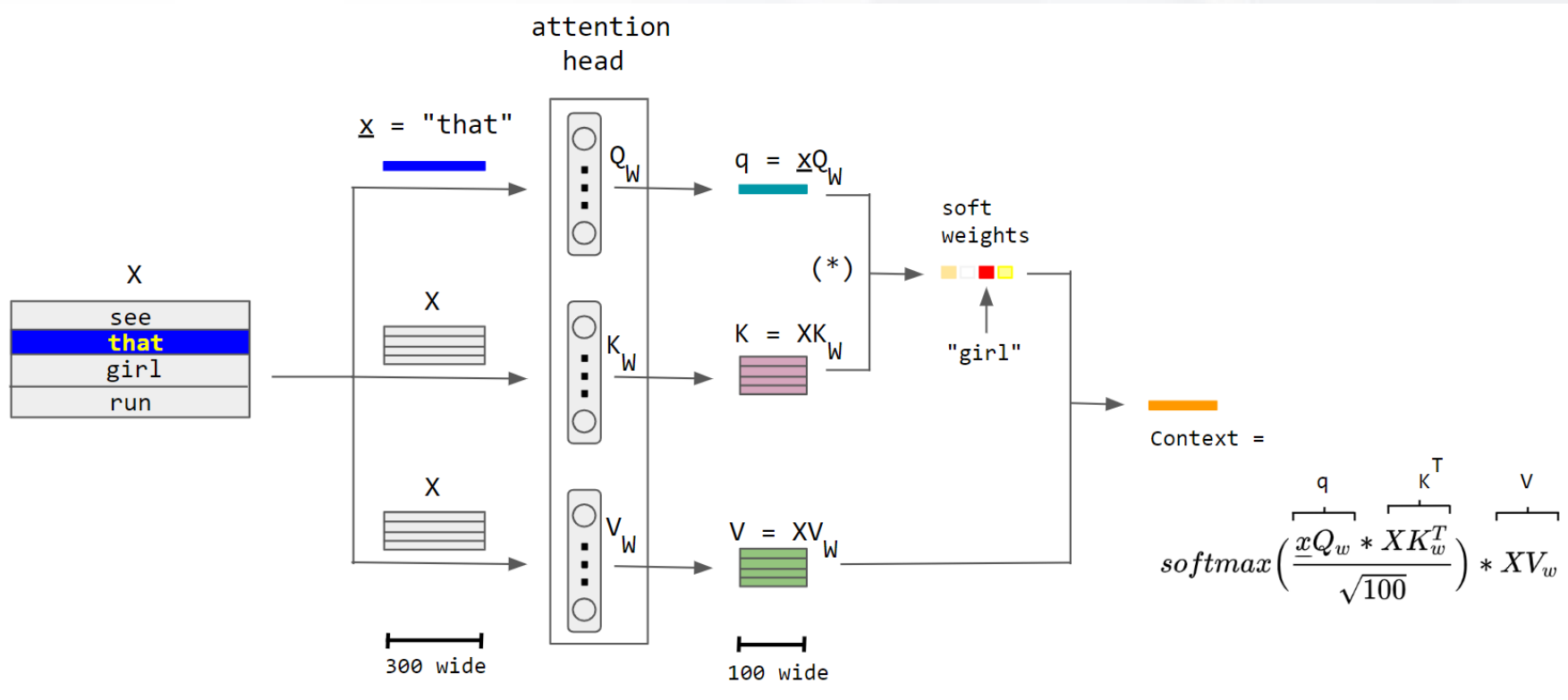


Performance Analysis of Parallelization on Vision Transformers

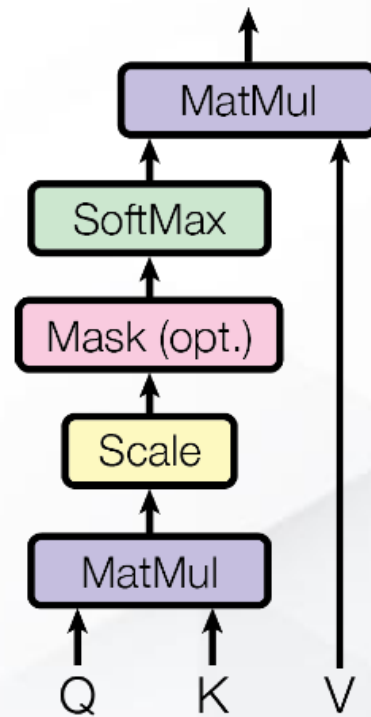
07.06.2023

Mustafa İzzet Muştu

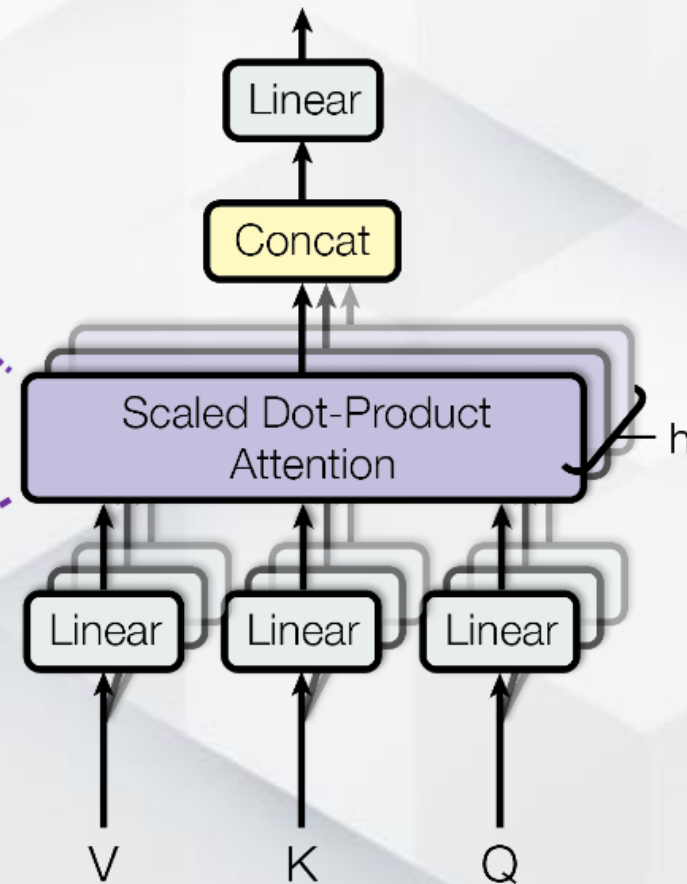
- Multi-Head Self-Attention
- Transformer
- Vision Transformer (ViT)
- Implementation Details
- Experiments & Results

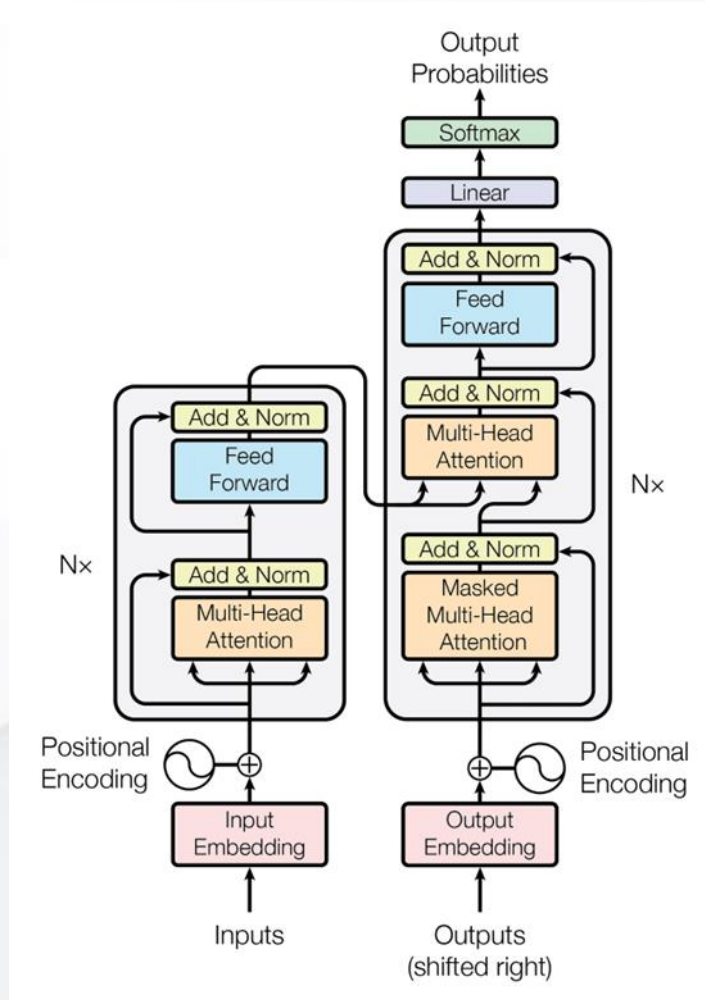


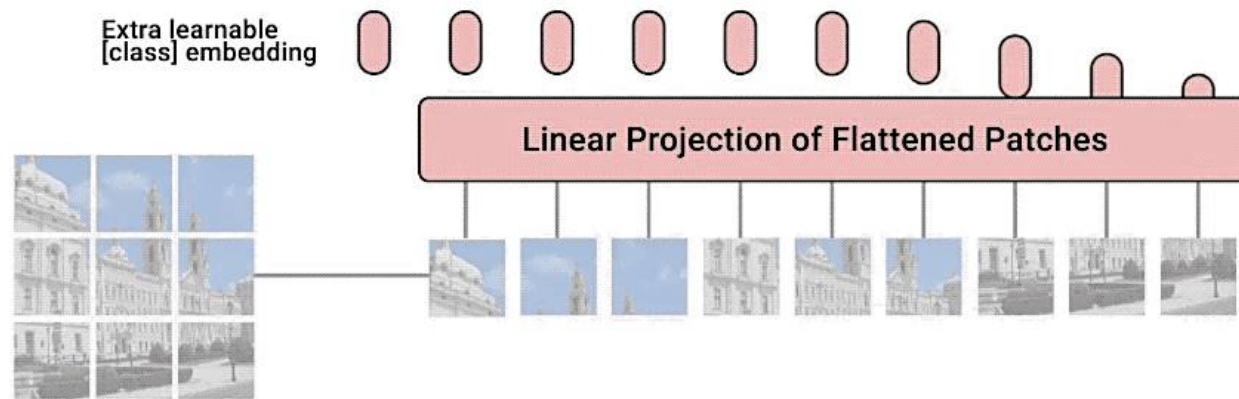
Scaled Dot-Product Attention

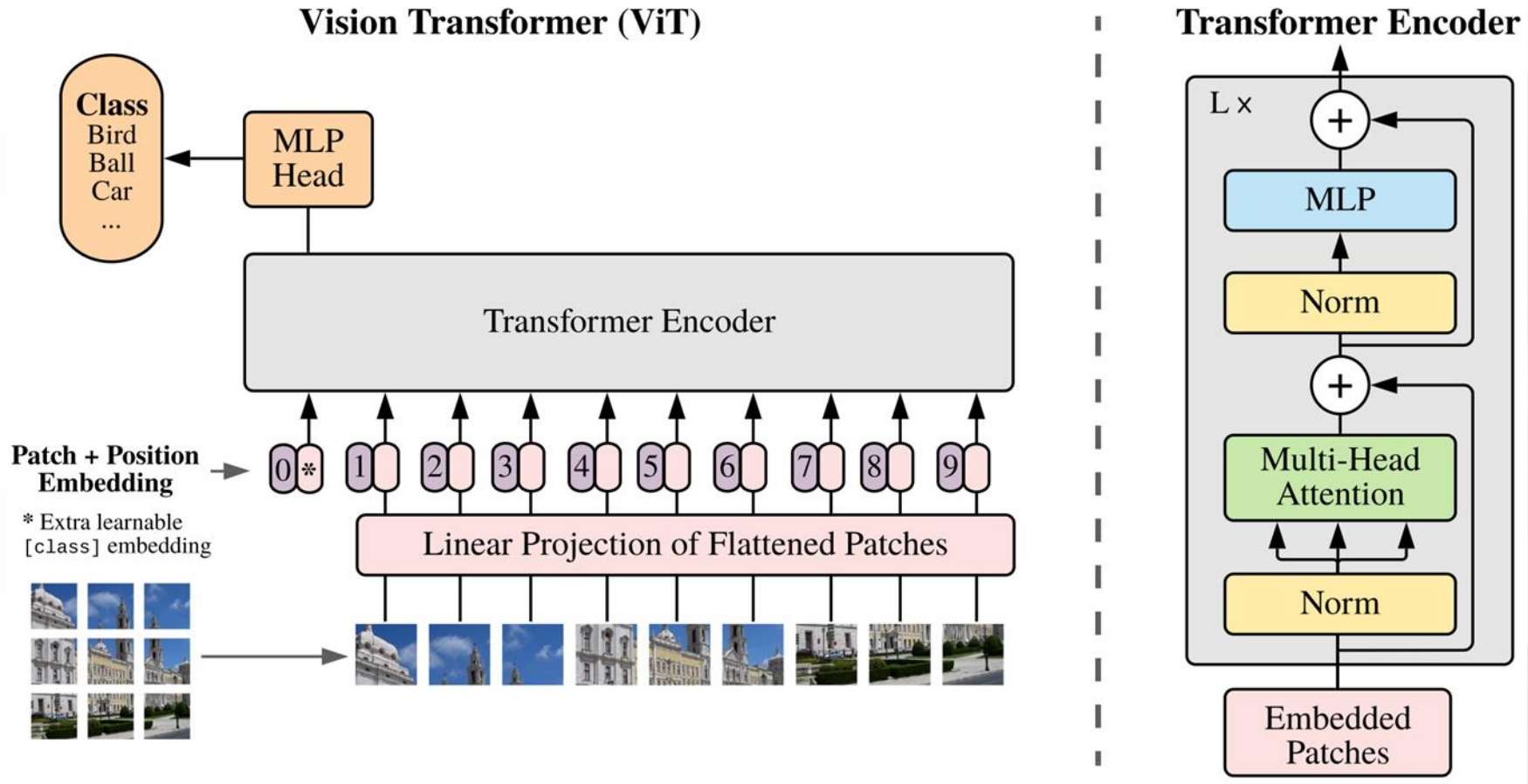


Multi-Head Attention









$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \quad (1)$$

$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \quad \ell = 1 \dots L \quad (2)$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, \quad \ell = 1 \dots L \quad (3)$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0) \quad (4)$$

$\mathcal{X} \in \mathbb{R}^{H \times W \times C} \rightarrow \mathcal{X}_p \in \mathbb{R}^{N \times (P^2 C)}$ where (H, W) is the resolution of the original image, C is the number of channels, (P, P) is the resolution of each image patch, and $N = HW/P^2$ is the resulting number of patches.

Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants.

- ViT-Base model without dropout and GELU in MLP
- Eigen Library for matrix representation
- Shared memory and a tiling approach for matrix multiplication
- 12 heads in MSA
- 12 layers in ViT encoder
- 12 CPU thread calls
- 12 CUDA streams
- $12(\text{layers}) * 12(\text{heads}) * 5 + 12$ multiplication kernel call for attention
- $12(\text{layers}) * 2$ multiplication kernel call for MLP
- 1 multiplication kernel call for embeddings

Algorithm 1 Tiled dense matrix-matrix multiplication kernel pseudocode

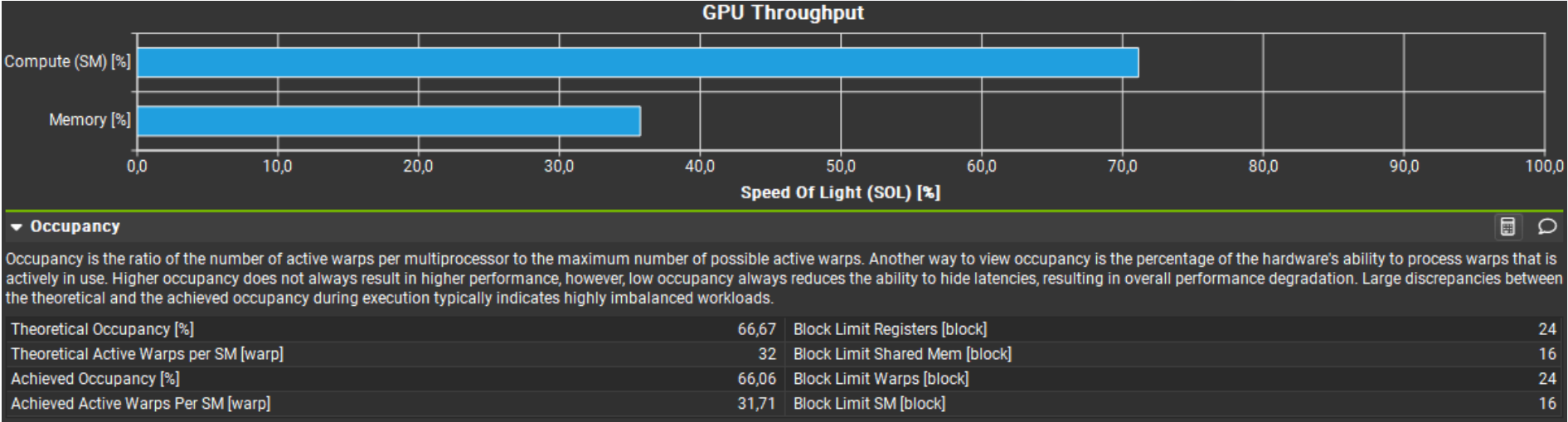
Require: $A, B, \text{TILE_WIDTH}, \text{Width}$

Ensure: $C = AB$

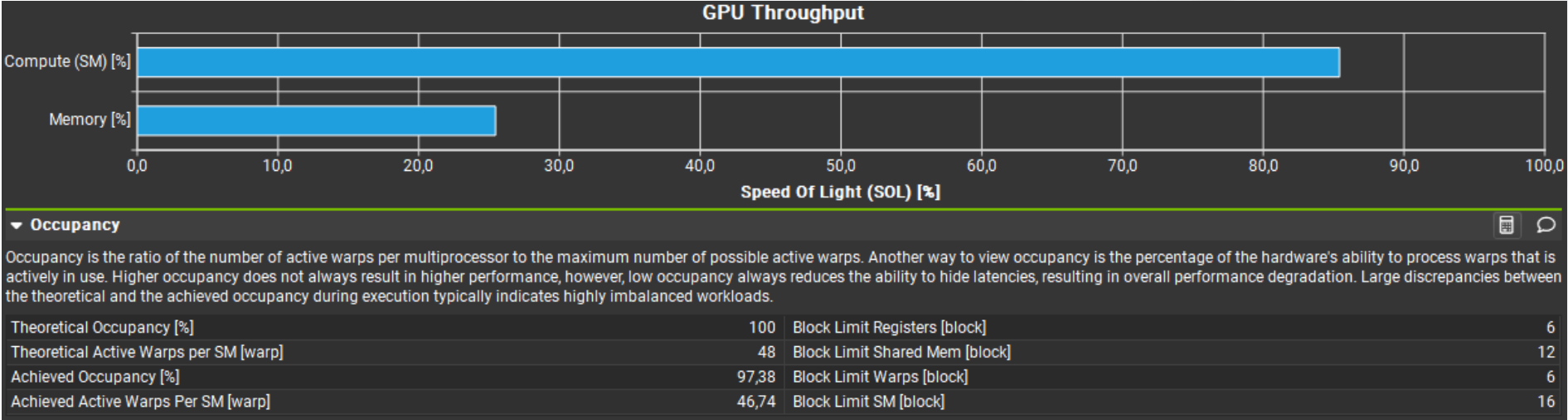
```

  sharedA[TILE_WIDTH][TILE_WIDTH]
  sharedB[TILE_WIDTH][TILE_WIDTH]
  bx ← blockIdx.x
  by ← blockIdx.y
  tx ← threadIdx.x
  ty ← threadIdx.y
  row ← ty + blockDim.y · by
  column ← tx + blockDim.x · bx
  result ← 0
  for i = 0, i < (Width-1)/TILE_WIDTH + 1, i += 1 do
    sharedA[ty][tx] ← A[row · Width + i · TILE_WIDTH + tx]
    sharedB[ty][tx] ← B[(i · TILE_WIDTH + ty) · Width + column]
    __syncthreads()
    for j = 0, j < TILE_WIDTH, j += 1 do
      result += sharedA[ty][j] · sharedB[j][tx]
    end for
    __syncthreads()
  end for
  C[row · Width + column] ← result

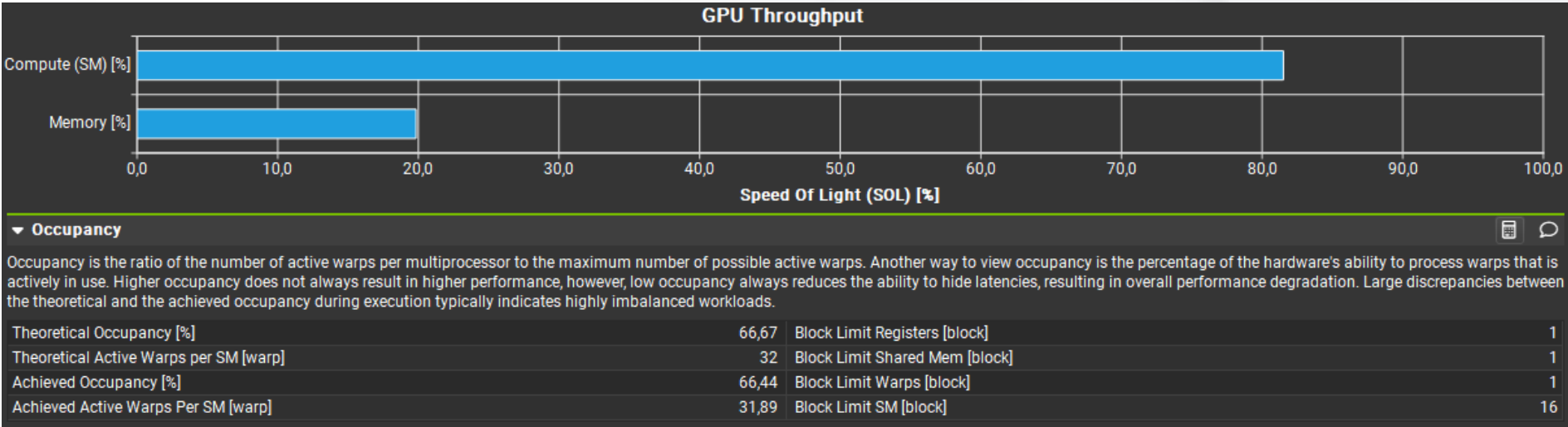
```



Tile width = 8



Tile width = 16



Tile width = 32

TABLE I

EXECUTION TIMES IN 1 LAYER OF THE ViT-BASE ENCODER (CPU).EXECUTION TIMES IN 1 LAYER OF THE ViT-BASE ENCODER (GPU-CPU).

Operation	Time (s)
Embeddings (Only before the first layer)	1.776
Layer Normalization	0.025
MSA (Multithreaded)	32.626
Residual	0.001
Layer Normalization	0.026
MLP	14.339
Residual	0.002
Layer Normalization (Only after the last layer)	0.025

TABLE II

Operation	Time (s)
Embeddings (Only before the first layer)	0.011
Layer Normalization	0.026
MSA (Multithreaded)	0.037
Residual	0.001
Layer Normalization	0.026
MLP	0.007
Residual	0.002
Layer Normalization (Only after the last layer)	0.026

Performance Analysis of Parallelization on Vision Transformers

Thank you

07.06.2023

Mustafa İzzet Muştu