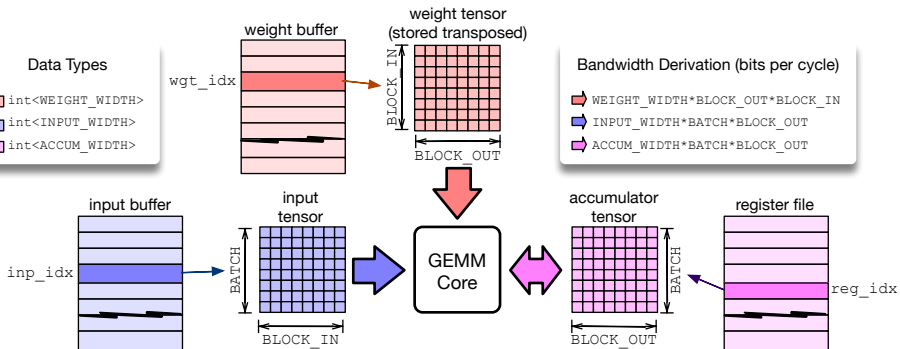


## Data Types

- $\text{int}\langle\text{WEIGHT\_WIDTH}\rangle$
- $\text{int}\langle\text{INPUT\_WIDTH}\rangle$
- $\text{int}\langle\text{ACCUM\_WIDTH}\rangle$

## Bandwidth Derivation (bits per cycle)

- ➡  $\text{WEIGHT\_WIDTH} \times \text{BLOCK\_OUT} \times \text{BLOCK\_IN}$
- ➡  $\text{INPUT\_WIDTH} \times \text{BATCH} \times \text{BLOCK\_OUT}$
- ➡  $\text{ACCUM\_WIDTH} \times \text{BATCH} \times \text{BLOCK\_OUT}$



## GEMM Instruction Pseudo-Code:

```
for i0 in range(0, end0):
    for i1 in range(0, end1):
        for uop_idx in range (uop_bgn, uop_end):
            x, y, z = decode_gemm_indices(uop_buffer[upc])
            reg_idx = i0 * x0 + i1 * x1 + x
            inp_idx = i0 * y0 + i1 * y1 + y
            wgt_idx = i0 * z0 + i1 * z1 + z
            reg_file[reg_idx] += GEMM(inp_buff[inp_idx], wgt_buff[wgt_idx])
```

LOG\_ACC\_BUFF\_DEPTH LOG\_INP\_BUFF\_DEPTH LOG\_WGT\_BUFF\_DEPTH

