# REPORT

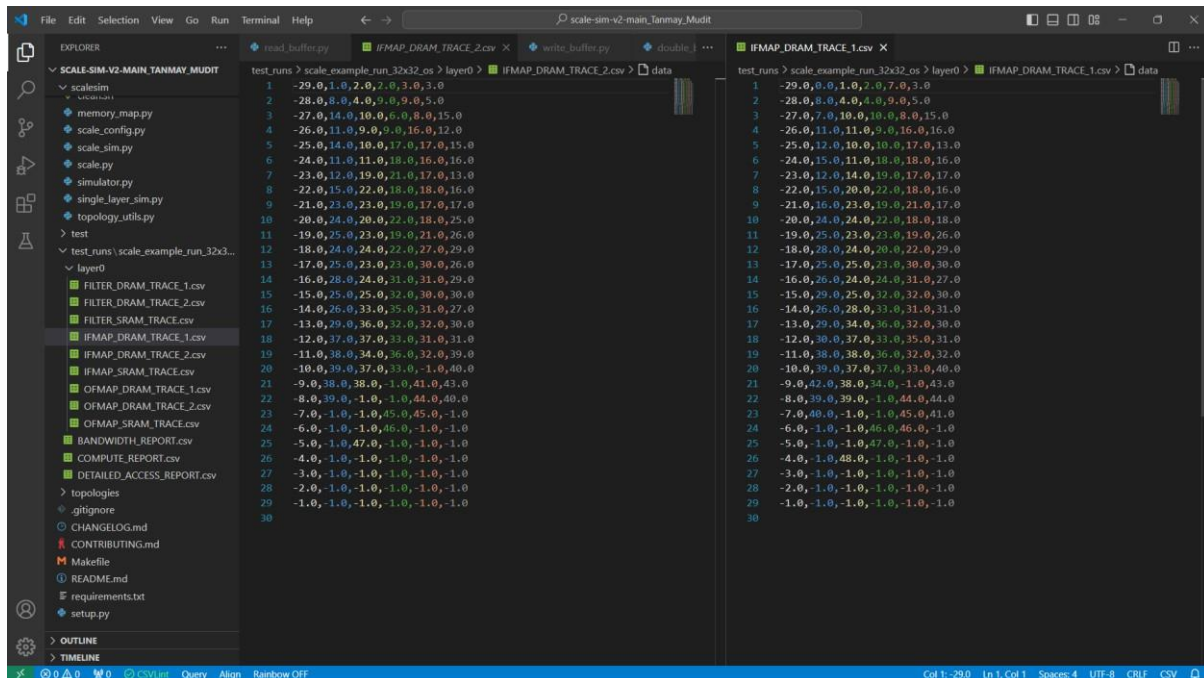## - Tanmay Aron (2023JVL2240)

## Primitive Mapping Implementation:

Exploring multiple memory bank modes in Scale- sim-v2. Data from local memory is distributed across different off-chip DRAMs. As a primitive mapping technique, data is cyclically distributed into 2 banks.

**Read Buffer: -**

- In this implementation, I've sent a row of requests into 2 DRAMs.

For each row, I am dividing addresses alternatively into both DRAMs.

Both the DRAMs are accessing elements parallelly.



- In the below function prefetch_active_buffer, separate prefetch requests and response cycle arrays for both DRAMS are created.

They are used to fill addresses and cycles in the respective trace matrices

```python
class read_buffer:
    def prefetch_active_buffer(self, start_cycle):
        response_cycles_arr2 = self.backing_buffer2.service_reads(incoming_cycles_arr=cycles_arr2,
                                                                  incoming_requests_arr_np=prefetch_requests2)
        # 4. Update the variables
        self.last_prefect_cycle = int(response_cycles_arr[-1][0])

        self.last_prefect_cycle1 = int(response_cycles_arr1[-1][0])
        self.trace_matrix1 = np.concatenate((response_cycles_arr1, prefetch_requests1), axis=1)

        self.last_prefect_cycle2 = int(response_cycles_arr2[-1][0])
        self.trace_matrix2 = np.concatenate((response_cycles_arr2, prefetch_requests2), axis=1)

        # Update the trace matrix
        self.trace_matrix = np.concatenate((response_cycles_arr, prefetch_requests), axis=1)
        self.trace_valid = True

        # Set active buffer contents
        active_buf_start_line_id = 0
        active_buf_end_line_id = self.num_active_buf_lines
        self.active_buffer_set_limits = [active_buf_start_line_id, active_buf_end_line_id]

        prefetch_buf_start_line_id = active_buf_end_line_id
        prefetch_buf_end_line_id = prefetch_buf_start_line_id + self.num_prefetch_buf_lines
        self.prefetch_buffer_set_limits = [prefetch_buf_start_line_id, prefetch_buf_end_line_id]

        self.active_buf_full_flag = True

        # Set the line to be prefetched next
        # The module operator is to ensure that the indices wrap around
        if requested_data_size > self.active_buf_size:  # Some elements in the current idx is left out in this case
            self.next_line_prefetch_idx = num_lines % self.fetch_matrix.shape[0]
        else:
            self.next_line_prefetch_idx = (num_lines + 1) % self.fetch_matrix.shape[0]
```

- Below are given functions created in read_buffer.py to read both DRAM's trace files, These functions are called from double_buffered_scratchpad_mem.py:

```python
 10    class read_buffer:
498        def get_external_access_start_stop_cycles(self):
502
503            return start_cycle, end_cycle
504    |
505
506        def print_trace_1(self, filename):
507            if not self.trace_valid:
508                print('No trace has been generated yet')
509                return
510
511
512            np.savetxt(filename, self.trace_matrix1, fmt='%s', delimiter=",")
513
514        def print_trace_2(self, filename):
515            if not self.trace_valid:
516                print('No trace has been generated yet')
517                return
518
519            np.savetxt(filename, self.trace_matrix2, fmt='%s', delimiter=",")
```

So, if I've set bandwidth=10 (in config file), then 5 addresses are read from 1st DRAM and other 5 addresses are read from 2nd DRAM.

This helps decrease bandwidth requirements per DRAM ({Bandwidth/n} for n DRAMs).

```
configs >  scale.cfg
  1    [general]
  2    run_name = scale_example_run_32x32_os
  3
  4    [architecture_presets]
  5    ArrayHeight:    32
  6    ArrayWidth:     32
  7    IfmapSramSzkB:   64
  8    FilterSramSzkB:  64
  9    OfmapSramSzkB:   64
 10    IfmapOffset:    0
 11    FilterOffset:   0
 12    OfmapOffset:    0
 13    Bandwidth : 10
 14    Dataflow : ws
 15    MemoryBanks:    1
 16
 17    [run_presets]
 18    InterfaceBandwidth: USER
 19
```

PROBLEMS    OUTPUT    DEBUG CONSOLE    **TERMINAL**    PORTS

```
Dataflow:         Weight Stationary
CSV file path:    ../scale-sim-v2-main_Tanmay_Mudit/topologies/conv_nets/test.csv
Bandwidth:        10
Working in USE USER BANDWIDTH mode.
========================================================

Running Layer 0
100%|
Compute cycles: 1235
Stall cycles: 0
Overall utilization: 6.83%
Mapping efficiency: 78.12%
Average IFMAP DRAM BW: 5.000 words/cycle
Average Filter DRAM BW: 5.000 words/cycle
Average OFMAP DRAM BW: 9.988 words/cycle
Saving traces: Done!
*********** SCALE SIM Run Complete ****************
```

As both the DRAMs are fetching data parallelly, the number of accesses is halved.

```
response_cycles_arr1 = self.backing_buffer1.service_reads(incoming_cycles_arr=cycles_arr1,
                                                          incoming_requests_arr_np=prefetch_requests1)
response_cycles_arr2 = self.backing_buffer2.service_reads(incoming_cycles_arr=cycles_arr2,
                                                          incoming_requests_arr_np=prefetch_requests2)
# 4. Update the variables

self.last_prefect_cycle1 = int(response_cycles_arr1[-1][0])
self.trace_matrix1 = np.concatenate((response_cycles_arr1, prefetch_requests1), axis=1)

self.last_prefect_cycle2 = int(response_cycles_arr2[-1][0])
self.trace_matrix2 = np.concatenate((response_cycles_arr2, prefetch_requests2), axis=1)

# Update the trace matrix
self.trace_matrix = np.concatenate((response_cycles_arr, prefetch_requests), axis=1)
self.trace_valid = True
```

- To make them parallel, I've created 2 read ports connected by 2 backing buffers. Now, both of the DRAMs are accessed independently in the same cycle.



- I've made changes in multiple files to print both the trace matrices of IFMAP, FILTER and OFMAP for every layer.

- In double_buffered_mem_scratchpad.py, we can reshape the matrices to double the number of columns and reduce the number of rows to half before calling service_read function.

```
def service_memory_requests(self, ifmap_demand_mat, filter_demand_mat, ofmap_demand_mat):
    ifmap_demand_mat = np.reshape(ifmap_demand_mat, (ifmap_demand_mat.shape[0] // 2, ifmap_demand_mat.shape[1] * 2))
    filter_demand_mat = np.reshape(filter_demand_mat, (filter_demand_mat.shape[0] // 2, filter_demand_mat.shape[1] * 2))
    ofmap_demand_mat = np.reshape(ofmap_demand_mat, (ofmap_demand_mat.shape[0] // 2, ofmap_demand_mat.shape[1] * 2))


    assert self.params_valid_flag, 'Memories not initialized yet'
    ifmap_lines = ifmap_demand_mat.shape[0]
    print("IFMAP LINES", ifmap_lines)
    ofmap_lines = ofmap_demand_mat.shape[0]

    self.total_cycles = 0
    self.stall_cycles = 0

    ifmap_hit_latency = self.ifmap_buf.get_hit_latency()
    filter_hit_latency = self.filter_buf.get_hit_latency()

    ifmap_serviced_cycles = []
    filter_serviced_cycles = []
    ofmap_serviced_cycles = []
    ifmap=int(ifmap_lines/2)
    pbar_disable = not self.verbose
    for i in tqdm(range(ifmap_lines), disable=pbar_disable):

        cycle_arr = np.zeros((1,1)) + i + self.stall_cycles

        ifmap_demand_line = ifmap_demand_mat[i, :].reshape((1,ifmap_demand_mat.shape[1]))
        ifmap_cycle_out = self.ifmap_buf.service_reads(incoming_requests_arr_np=ifmap_demand_line,
                                                        incoming_cycles_arr=cycle_arr)
        ifmap_serviced_cycles += [ifmap_cycle_out[0]]
        ifmap_stalls = ifmap_cycle_out[0] - cycle_arr[0] - ifmap_hit_latency
```

This will reduce the compute cycles to half of the original as more data is sent in a single cycle.

```
Running Layer 1
IFMAP LINES 1236
100%|                                                                        | 1236/1236 [00:00<00:00, 10641.88i
100%|                                                                        | 1236/1236 [00:00<00:00, 22636.57i
Compute cycles: 1235    ←
Stall cycles: 0
Overall utilization: 6.83%
Mapping efficiency: 78.12%
Average IFMAP DRAM BW: 5.000 words/cycle
Average Filter DRAM BW: 5.000 words/cycle
Average OFMAP DRAM BW: 9.988 words/cycle
Saving traces: Done!
*********** SCALE SIM Run Complete ***************
PS C:\Users\Tanmay\Documents\IITD\SEM2\Synthesis of Digital systems'\scale-sim-v2-main_Tanmay_Mudit>
```

*Figure 1 Compute cycles before making the changes*

```
Running Layer 1
IFMAP LINES 618
100%|                                                                        | 618/618 [00:00<00:00, 10098.37
100%|                                                                        | 618/618 [00:00<00:00, 16710.05
Compute cycles: 617
Stall cycles: 0
Overall utilization: 13.68%
Mapping efficiency: 78.12%
Average IFMAP DRAM BW: 5.000 words/cycle
Average Filter DRAM BW: 5.000 words/cycle
Average OFMAP DRAM BW: 9.988 words/cycle
Saving traces: Done!
*********** SCALE SIM Run Complete ***************
PS C:\Users\Tanmay\Documents\IITD\SEM2\Synthesis of Digital systems'\scale-sim-v2-main_Tanmay_Mudit>
```

*Figure 2 Compute cycles after making the changes*

## Write Buffer:

- Creating 2 trace matrices for the write buffer:

```python
class write_buffer:
    def __init__(self):
        self.req_gen_bandwidth = 100

        # Status of the buffer
        self.free_space = self.total_size_elems
        self.drain_buf_start_line_id = 0
        self.drain_buf_end_line_id = 0

        # Helper data structures for faster execution
        self.line_idx = 0
        self.current_line1 = np.ones((1, 1)) * -1
        self.current_line2 = np.ones((1, 1)) * -1
        self.max_cache_lines = 2 ** 10
        self.trace_matrix_cache1 = np.zeros((1, 1))
        self.trace_matrix_cache2 = np.zeros((1, 1))

        # Access counts
        self.num_access = 0

        # Trace matrices
        self.trace_matrix1 = np.zeros((1, 1))
        self.trace_matrix2 = np.zeros((1, 1))
        self.cycles_vec = np.zeros((1, 1))

        # Flags
        self.state = 0
        self.drain_end_cycle = 0

        self.trace_valid = False
        self.trace_valid_2 = False
        self.trace_matrix_cache1_empty = True
        self.trace_matrix_cache2_empty = True
        self.trace_matrix_empty = True

    def set_params(self, backing_buf_obj,
                   total_size_bytes=128, word_size=1, active_buf_frac=0.9,
```

*Figure 3. Creating 2 trace matrices for write buffer*

- Both trace matrices are filled in the functions store_to_trace_mat_cache and append_to_trace_mat

```
  7   class write_buffer:
 94       def store_to_trace_mat_cache(self, elem):
116                   self.trace_matrix_cache1 = self.current_line1
117                   self.trace_matrix_cache1_empty = False
118               else:
119                   self.trace_matrix_cache1 = np.concatenate((self.trace_matrix_cache1, self.current_line1), axis=0)
120
121               if self.trace_matrix_cache2_empty:
122                   self.trace_matrix_cache2 = self.current_line2
123                   self.trace_matrix_cache2_empty = False
124               else:
125                   self.trace_matrix_cache2 = np.concatenate((self.trace_matrix_cache2, self.current_line2), axis=0)
126
127               self.current_line1 = np.ones((1, 1)) * -1
128               self.current_line2 = np.ones((1, 1)) * -1
129               self.line_idx = 0
130
131               if not self.trace_matrix_cache1.shape[0] < self.max_cache_lines:
132                   self.append_to_trace_mat()
133
134               if not self.trace_matrix_cache2.shape[0] < self.max_cache_lines:
135                   self.append_to_trace_mat(force=True)
136
137       def append_to_trace_mat(self, force=False):
138           if force:
139               if not self.line_idx == 0:
140                   if self.trace_matrix_cache1_empty:
141                       self.trace_matrix_cache1 = self.current_line1
142                       self.trace_matrix_cache2 = self.current_line2
143                       self.trace_matrix_cache_empty = False
144                   else:
145                       self.trace_matrix_cache1 = np.concatenate((self.trace_matrix_cache1, self.current_line1), axis=0)
146                       self.trace_matrix_cache2 = np.concatenate((self.trace_matrix_cache2, self.current_line2), axis=0)
147
148                   self.current_line1 = np.ones((1, 1)) * -1
149                   self.current_line2 = np.ones((1, 1)) * -1
150                   self.line_idx = 0
```

*Figure 4. Both trace matrices are filled in the functions store_to_trace_mat_cache and append_to_trace_mat*

- These trace matrices are filled and printed using these functions:
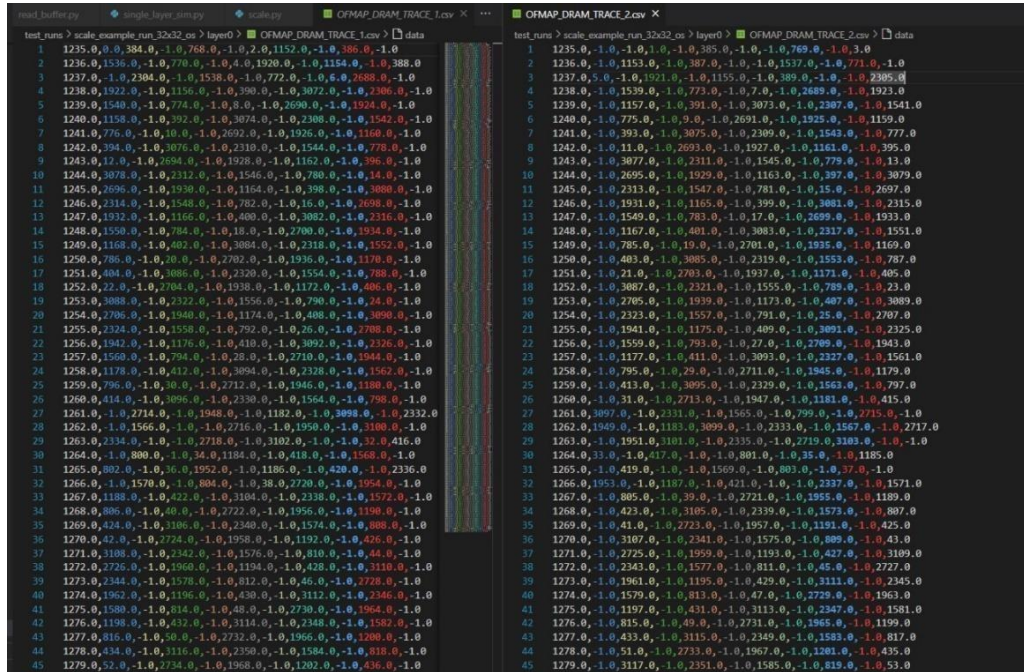
```
  7   class write_buffer:
277
278       def get_trace_matrix1(self):
279           if not self.trace_valid:
280               print('No trace has been generated yet')
281               return
282
283           trace_matrix1 = np.concatenate((self.cycles_vec, self.trace_matrix1), axis=1)
284           return trace_matrix1
285
286       def get_trace_matrix2(self):
287           if not self.trace_valid:
288               print('No trace has been generated yet')
289               return
290
291           trace_matrix2 = np.concatenate((self.cycles_vec, self.trace_matrix2), axis=1)
292           return trace_matrix2
293
294       def get_free_space(self):
295           return self.free_space
296
297       def get_num_accesses(self):
298           assert self.trace_valid, 'Traces not ready yet'
299           return self.num_access
300
301       def get_external_access_start_stop_cycles(self):
302           assert self.trace_valid, 'Traces not ready yet'
303           start_cycle = self.cycles_vec[0][0]
304           end_cycle = self.cycles_vec[-1][0]
305           return start_cycle, end_cycle
306
307       def print_trace1(self, filename):
308           if not self.trace_valid:
309               print('No trace has been generated yet')
310               return
311           trace_matrix1 = self.get_trace_matrix1()
312           np.savetxt(filename, trace_matrix1, fmt='%s', delimiter=",")
313
314       def print_trace2(self, filename):
315           if not self.trace_valid:
316               print('No trace has been generated yet')
317               return
318           trace_matrix2 = self.get_trace_matrix2()
319           np.savetxt(filename, trace_matrix2, fmt='%s', delimiter=",")
320
```

*Figure 5. The trace matrices are filled and printed using these functions*

These trace matrices are generated in write_buffer.py



*Figure c. OFMAP tracematrix*

➔ Please follow discussions regarding parameters (Avg. Bandwidth, utilization, etc.) in the Novel Implementation report

**Results (for read_buffer):-**

1. Total number of cycles remains the same.
2. The bandwidth per DRAM is halved.

**Results (for write_buffer):-**

1. Total number of cycles is halved.
2. The bandwidth per DRAM is the same.

# Novel Mapping Implementation: -

- In this implementation, I've divided the data row-wise alternatively in 2 DRAMs.
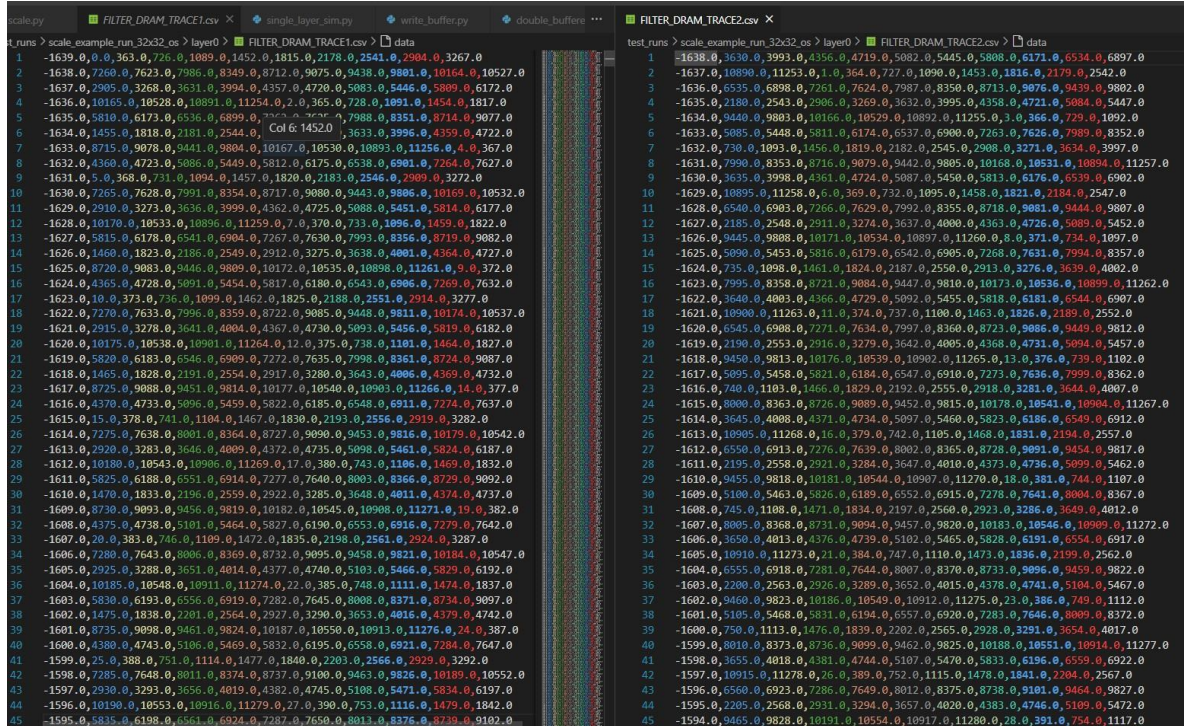


*Figure 6 Trace Matrices*

- During memory mapping, if a data block is previously fetched from DRAM i, and if it is now requested from DRAM j, then the trace matrix of DRAM j will show '-1'. And, data block is accessed from DRAM i.

```
class read_buffer:
    def prefetch_active_buffer(self, start_cycle):
        prefetch_requests = self.fetch_matrix[start_idx:end_idx, :]
        first_occurrences = {}

        for i, row in enumerate(prefetch_requests):
            row_type = "even" if i % 2 == 0 else "odd"
            for j, element in enumerate(row):
                if element in first_occurrences:
                    # If already seen and the current row type doesn't match the first occurrence row type
                    if first_occurrences[element] != row_type:
                        row[j] = -1  # Modify the element
                else:
                    # First occurrence, store the type of row it was first seen
                    first_occurrences[element] = row_type

        prefetch_requests1=[]
        prefetch_requests2=[]
        num1=0
        num2=0
        alt=0
        for i, row in enumerate(prefetch_requests):
            if alt == 0:
                alt=1
                prefetch_requests1.append(row)
                num1=num1+1
            else:
                alt=0
                prefetch_requests2.append(row)
                num2 += 1

        prefetch_requests1 = np.array(prefetch_requests1)
        prefetch_requests2 = np.array(prefetch_requests2)
```

*Figure 7. Memory mapping and division of rows*



*Figure 8. '-1's are put if the blocks are previously fetched from different DRAM*

**Results and Observations: -**

1. Mapping efficiency- Depends on number of MACs, array height and width of systolic arrays. This can't be improved by modifying memory mapping.
   Although, by altering row and column size we can increase mapping efficiency.

```python
class systolic_compute_is:
    def create_ifmap_demand_mat(self):
        this_fold_demand = np.concatenate((this_fold_demand, null_req_mat), axis=0)

        # The IFMAP elems are needed to be filled in reverse order to ensure that
        # top element is pushed in last to maintain alignment with the input elements
        this_fold_demand = np.flip(this_fold_demand, 0)

        # Account for the cycles for partial sum generation and accumulation
        this_fold_demand = np.concatenate((this_fold_demand, inter_fold_gap_suffix_mat), axis=0)

        # Calculate the mapping efficiency
        row_used = min(self.arr_row, row_end_idx - row_start_id)
        col_used = min(self.arr_col, col_end_idx - col_start_id)
        mac_used = row_used * col_used
        mapping_eff_this_fold = mac_used / (self.arr_row * self.arr_col)

        cycles_this_fold = this_fold_demand.shape[0] + this_fold_demand.shape[1] - 1
        compute_cycles_this_fold = mac_used * self.T
        compute_util_this_fold = compute_cycles_this_fold / (self.arr_row * self.arr_col * cycles_this_fold)

        self.mapping_efficiency_per_fold.append(mapping_eff_this_fold)
        self.compute_utility_per_fold.append(compute_util_this_fold)

        if fr == 0 and fc == 0:
            self.ifmap_demand_matrix = this_fold_demand
        else:
            self.ifmap_demand_matrix = np.concatenate((self.ifmap_demand_matrix, this_fold_demand), axis=0)
```

2. Compute cycles- Compute cycles are total cycles after writing the output feature map values.

```python
    def get_total_compute_cycles(self):
        assert self.traces_valid, 'Traces not generated yet'
        return self.total_cycles
```

Ways compute cycles are decreased: -
   a) Increasing row size sent per cycle.
   b) Decreasing cycles during Read and Write:
      In our implementation, total cycles by half by fetching 2 rows in a single cycle through 2 DRAMs.

*Figure 9. Original Scale-sim results*



*Figure 10. New Scale-sim results*

## 3. Overall Utilization-

```
self.overall_util = (self.num_compute * 100) / (self.total_cycles * self.num_mac_unit)
```

This formula computes what portion of the maximum possible computational capacity is actually being utilized. This can provide insights into how efficiently the computational resources are being used.

➔ As shown in the above figures, there's about twice the increase in overall utilization.

## 4. Average DRAM BW:-

In the config file, we have input the bandwidth of DRAMs.

So, the average DRAM BW should not be more than this value. Whereas, the total DRAM bandwidth will be:-

Total DRAM BW= Avg. DRAM BW * No. of DRAMs

```
1    [general]
2    run_name = scale_example_run_32x32_os
3
4    [architecture_presets]
5    ArrayHeight:     32
6    ArrayWidth:      32
7    IfmapSramSzkB:    64
8    FilterSramSzkB:   64
9    OfmapSramSzkB:    64
10   IfmapOffset:     0
11   FilterOffset:    0
12   OfmapOffset:     0
13   Bandwidth  :  10
14   Dataflow : ws
15   MemoryBanks:     1
16
17   [run_presets]
18   InterfaceBandwidth: USER
19
```

Bandwidth is dependent on reads/writes and start, stop cycles.

```
class single_layer_sim:
    def calc_report_data(self):
        self.ofmap_sram_writes = self.compute_system.get_ofmap_requests()
        self.avg_ifmap_sram_bw = self.ifmap_sram_reads / self.total_cycles
        self.avg_filter_sram_bw = self.filter_sram_reads / self.total_cycles
        self.avg_ofmap_sram_bw = self.ofmap_sram_writes / self.total_cycles

        # Detail report
        self.ifmap_sram_start_cycle, self.ifmap_sram_stop_cycle \
            = self.memory_system.get_ifmap_sram_start_stop_cycles()

        self.filter_sram_start_cycle, self.filter_sram_stop_cycle \
            = self.memory_system.get_filter_sram_start_stop_cycles()

        self.ofmap_sram_start_cycle, self.ofmap_sram_stop_cycle \
            = self.memory_system.get_ofmap_sram_start_stop_cycles()

        self.ifmap_dram_start_cycle, self.ifmap_dram_stop_cycle, self.ifmap_dram_reads \
            = self.memory_system.get_ifmap_dram_details()

        self.filter_dram_start_cycle, self.filter_dram_stop_cycle, self.filter_dram_reads \
            = self.memory_system.get_filter_dram_details()

        self.ofmap_dram_start_cycle, self.ofmap_dram_stop_cycle, self.ofmap_dram_writes \
            = self.memory_system.get_ofmap_dram_details()

        # BW calc for DRAM access
        self.avg_ifmap_dram_bw = self.ifmap_dram_reads / (self.ifmap_dram_stop_cycle - self.ifmap_dram_start_cycle + 1)
        self.avg_filter_dram_bw = self.filter_dram_reads / (self.filter_dram_stop_cycle - self.filter_dram_start_cycle + 1)
        self.avg_ofmap_dram_bw = self.ofmap_dram_writes / (self.ofmap_dram_stop_cycle - self.ofmap_dram_start_cycle + 1)

        self.report_items_ready = True
```

Reads/Writes are computed using number of accesses either for fetching or writing on addresses.

```
#
def get_ifmap_dram_details(self):
    assert self.traces_valid, 'Traces not generated yet'

    self.ifmap_dram_reads = self.ifmap_buf.get_num_accesses()
    self.ifmap_dram_start_cycle, self.ifmap_dram_stop_cycle \
        = self.ifmap_buf.get_external_access_start_stop_cycles()

    return self.ifmap_dram_start_cycle, self.ifmap_dram_stop_cycle, self.ifmap_dram_reads
```

num_accesses are calculated using addition of number of data blocks.

- This is a simulation of CNN accelerator. We are performing DRAM accesses parallelly (by creating different read ports for different DRAMs).

So, I've modified num_access to the total number of data blocks _accessed by 1 DRAM_ as other DRAMs are accessed in parallel.

→ For demonstration purposes, let's change num_access to total number of data blocks fetched for writing in OFMAP.

```python
def empty_drain_buf(self, empty_start_cycle=0):
    lines_to_fill_dbuf = int(math.ceil(self.drain_buf_size / self.req_gen_bandwidth))
    self.drain_buf_end_line_id1 = self.drain_buf_start_line_id1 + lines_to_fill_dbuf
    self.drain_buf_end_line_id1 = min(self.drain_buf_end_line_id1, self.trace_matrix1.shape[0])
    self.drain_buf_end_line_id2 = self.drain_buf_start_line_id2 + lines_to_fill_dbuf
    self.drain_buf_end_line_id2 = min(self.drain_buf_end_line_id2, self.trace_matrix2.shape[0])

    requests_arr_np1 = self.trace_matrix1[self.drain_buf_start_line_id1: self.drain_buf_end_line_id1, :]
    requests_arr_np2 = self.trace_matrix2[self.drain_buf_start_line_id2: self.drain_buf_end_line_id2, :]
    num_lines1 = requests_arr_np1.shape[0]
    num_lines2 = requests_arr_np2.shape[0]

    data_sz_to_drain1 = num_lines1 * requests_arr_np1.shape[1]
    data_sz_to_drain2 = num_lines2 * requests_arr_np2.shape[1]
    for elem in requests_arr_np1[-1,:]:
        if elem == -1:
            data_sz_to_drain1 -= 1
    for elem in requests_arr_np2[-1,:]:
        if elem == -1:
            data_sz_to_drain2 -= 1
    self.num_access += data_sz_to_drain1 + data_sz_to_drain2
    #self.num_access += data_sz_to_drain1
```

Avg DRAM BW when accessed parallelly: -

```
Running Layer 0
IFMAP LINES 588
100%|
100%|
Compute cycles: 587
Stall cycles: 0
Overall utilization: 6.39%
Mapping efficiency: 78.12%
Average IFMAP DRAM BW: 10.000 words/cycle
Average Filter DRAM BW: 10.000 words/cycle
Average OFMAP DRAM BW: 9.974 words/cycle
Saving traces: Done!
*********** SCALE SIM Run Complete ****************
```

_Figure 11. Correct Avg DRAM BW_

Avg DRAM BW when data blocks are accessed sequentially (Incorrect) : -

```
Running Layer 0
IFMAP LINES 588
100%|
100%|
Compute cycles: 587
Stall cycles: 0
Overall utilization: 6.39%
Mapping efficiency: 78.12%
Average IFMAP DRAM BW: 10.000 words/cycle
Average Filter DRAM BW: 10.000 words/cycle
Average OFMAP DRAM BW: 19.948 words/cycle
Saving traces: Done!
```