

Dynamic Resource Partitioning for Multi-Tenant Systolic Array Based DNN Accelerator

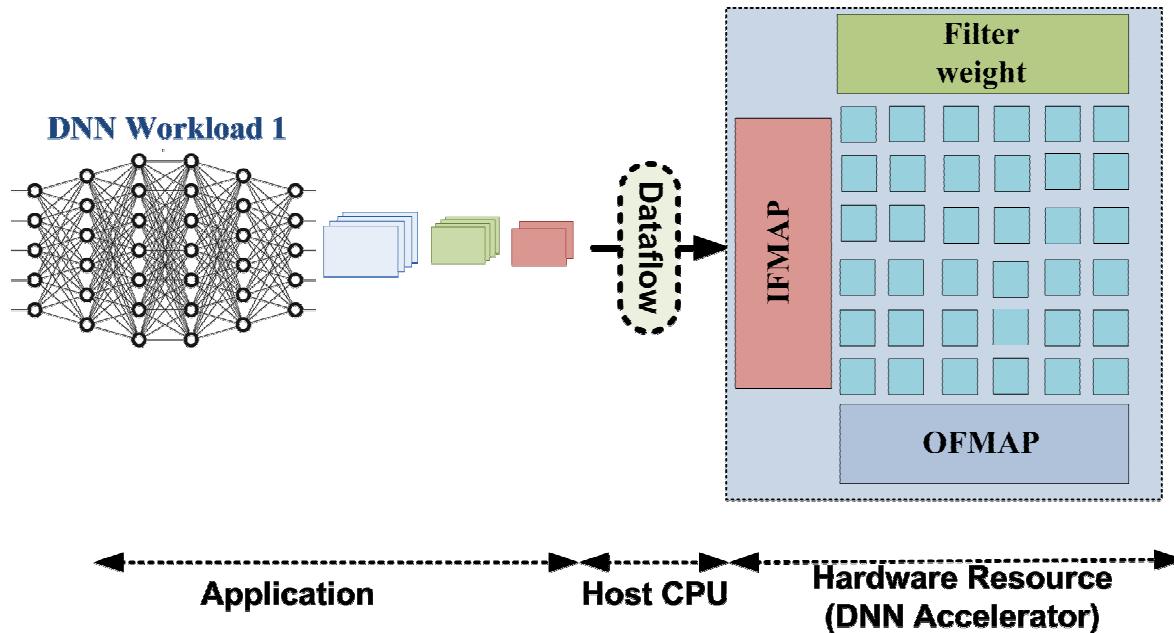
Midia Reshadi, David Gregg
School of Computer Science and Statistics,
Lero, Trinity College Dublin

Outline

- Introduction
- Preliminaries
- Dynamic Resource Partitioning Algorithm
- LOCAL mapping algorithm
- Simulation Results

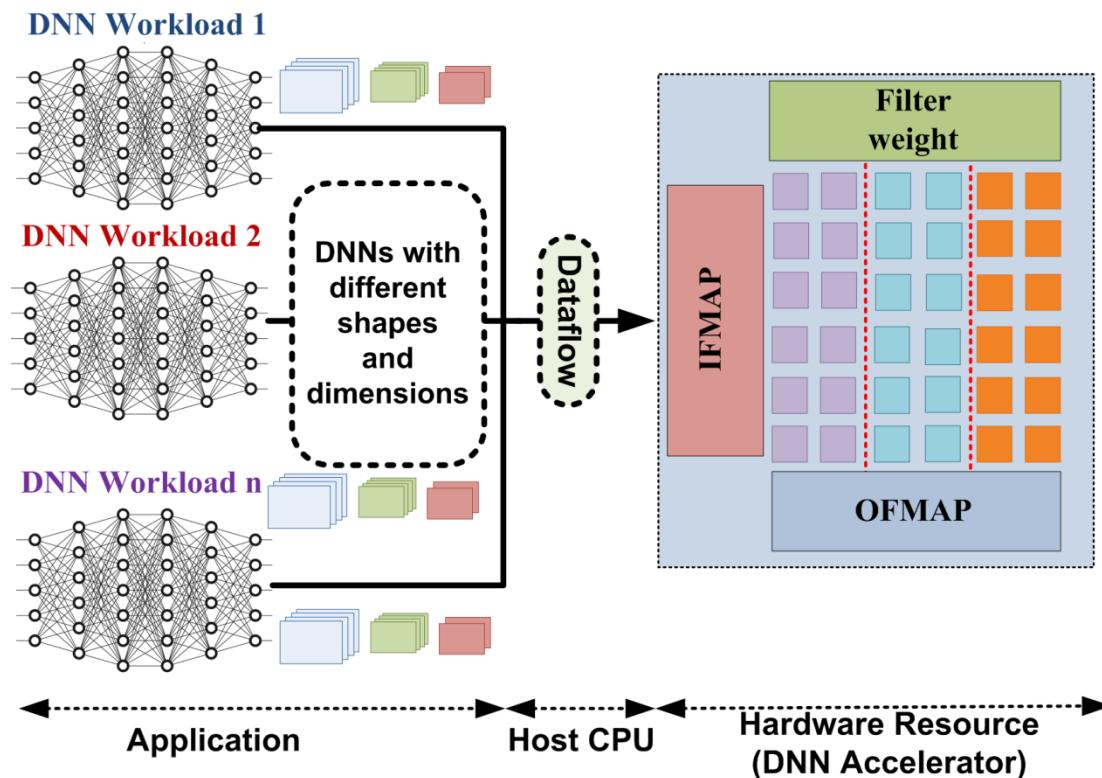
Introduction

- Single tenancy



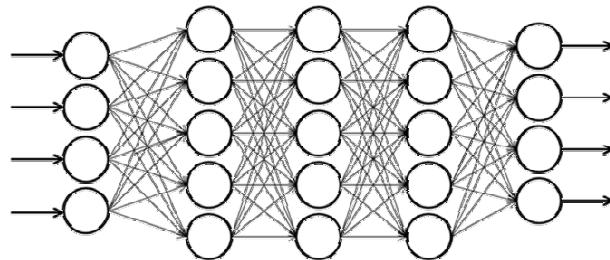
Introduction

- Multi tenancy



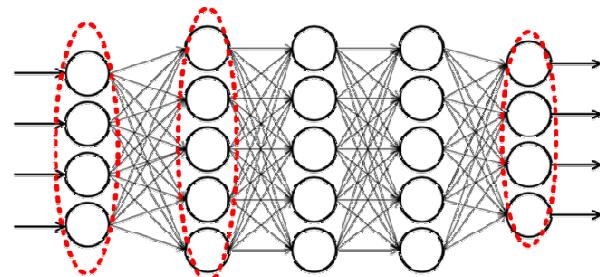
Preliminaries

- Deep Neural Networks Graph (DNNG)



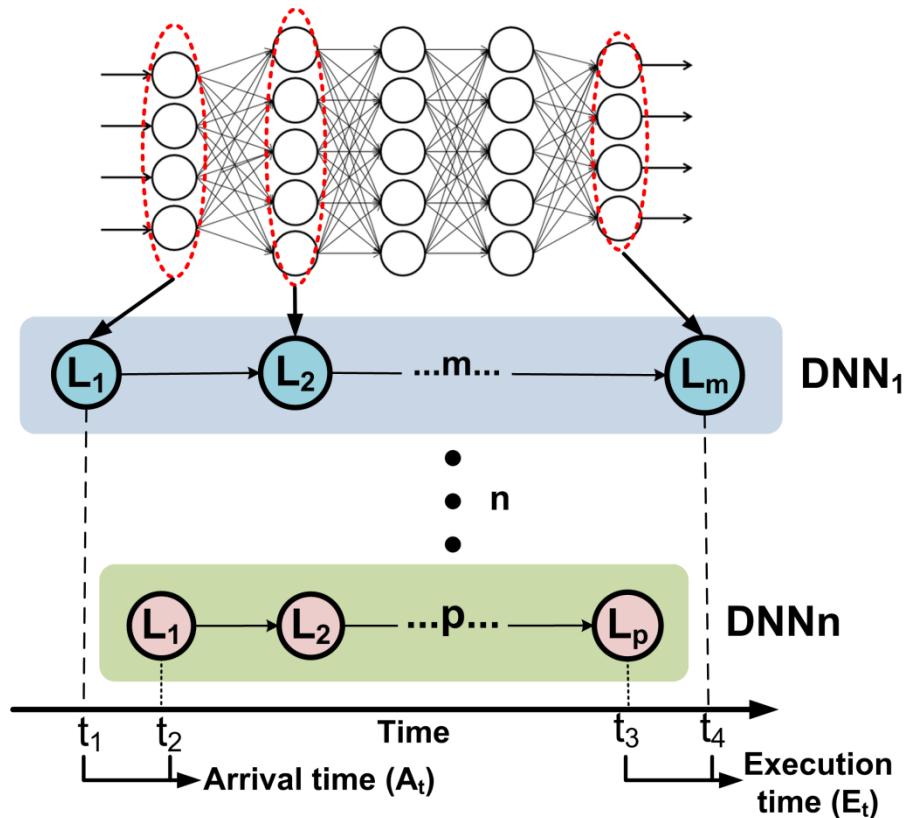
Preliminaries

- Deep Neural Networks Graph (DNNG)



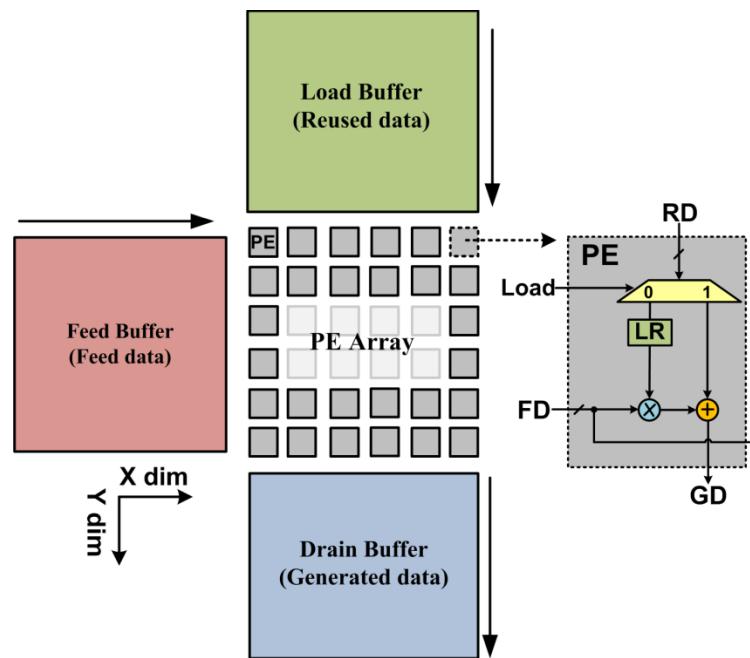
Preliminaries

- Deep Neural Networks Graph (DNNG)



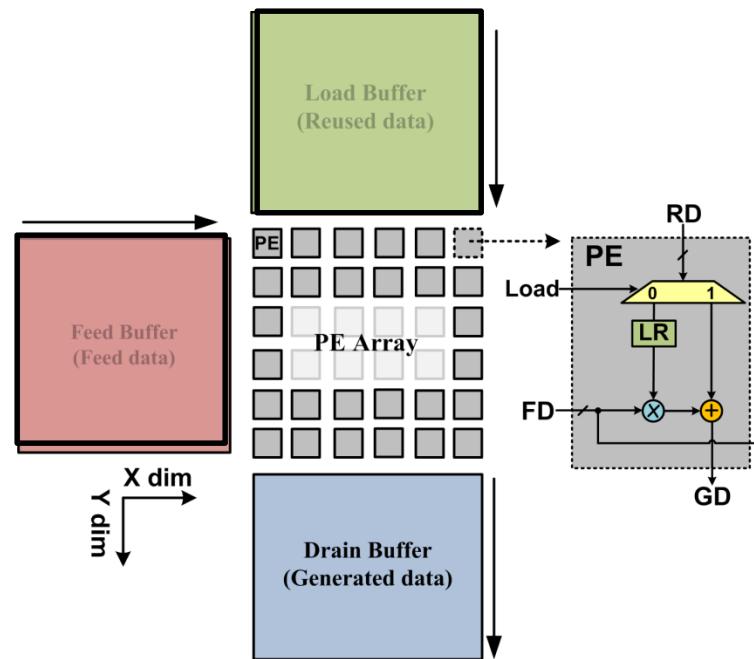
Preliminaries

- Systolic-array Based DNN Accelerator



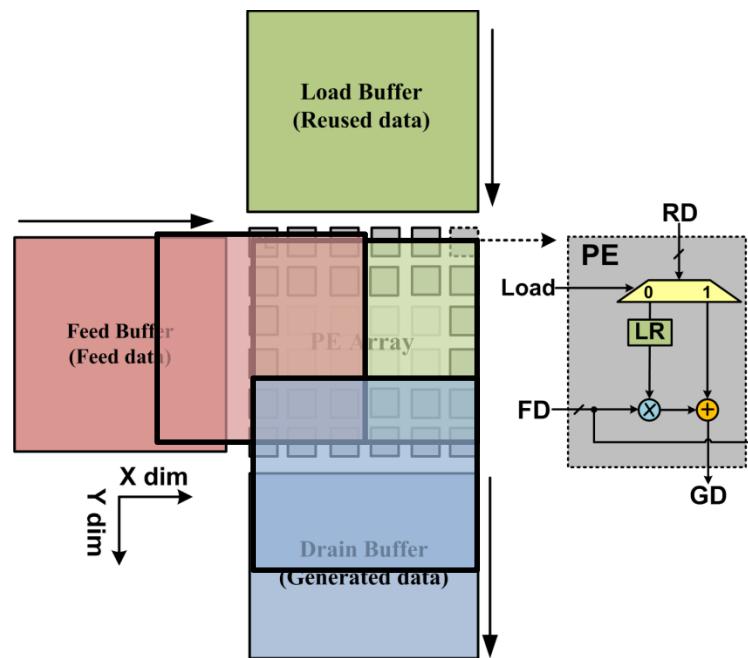
Preliminaries

- Systolic-array Based DNN Accelerator



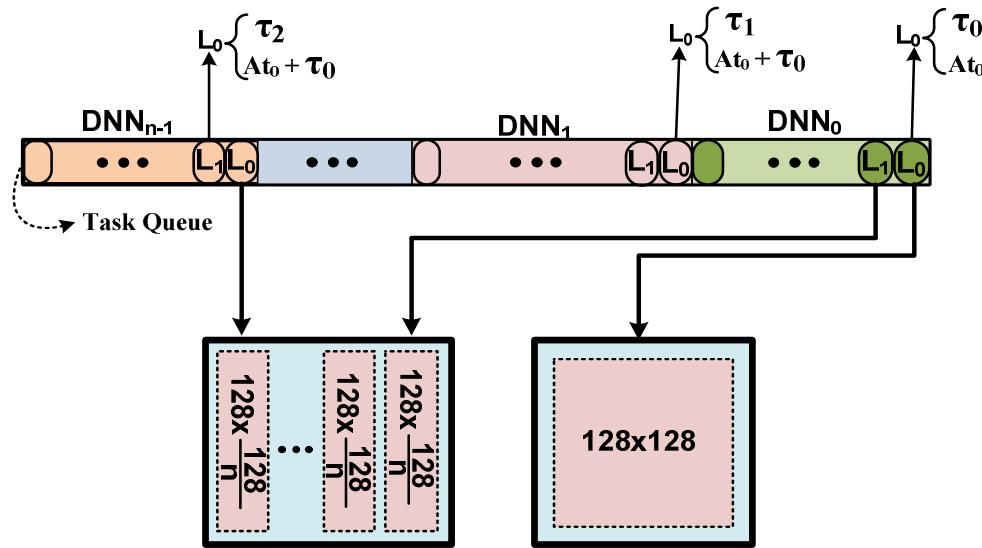
Preliminaries

- Systolic-array Based DNN Accelerator



Dynamic Resource Partitioning Algorithm

- Systolic Array Partitioning



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

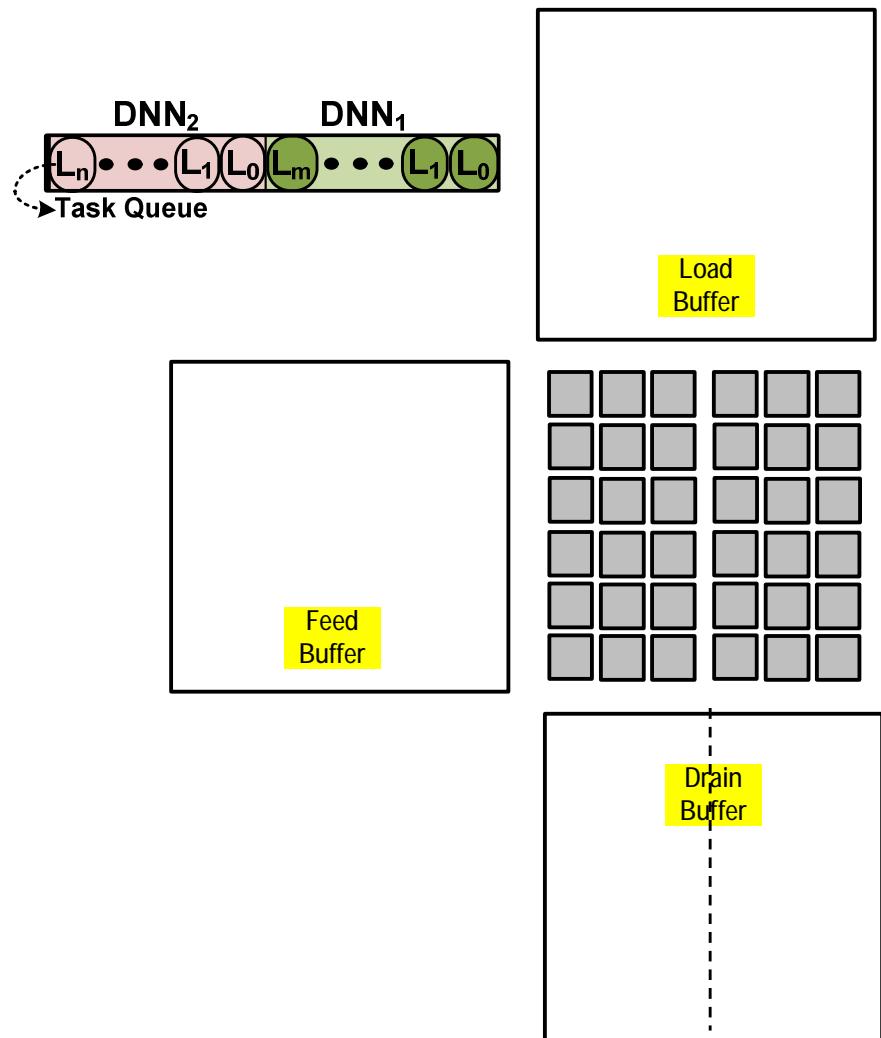
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

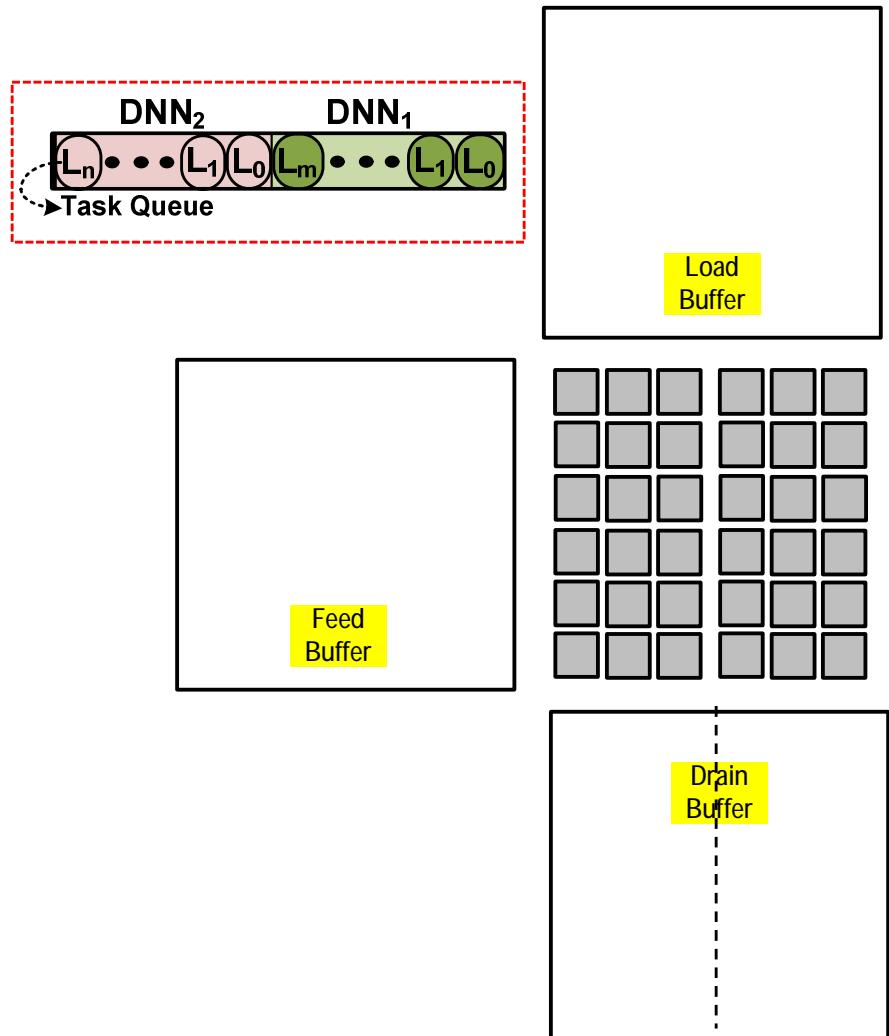
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition column] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition column]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ①
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition column]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

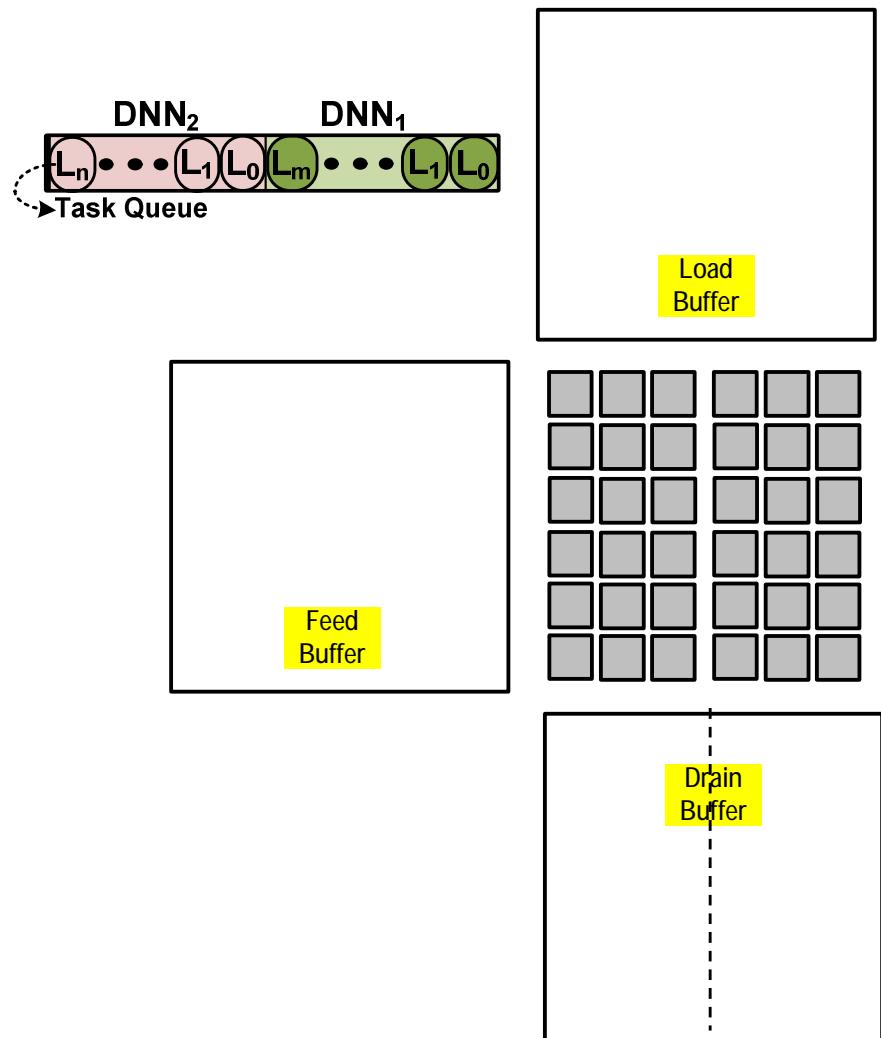
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

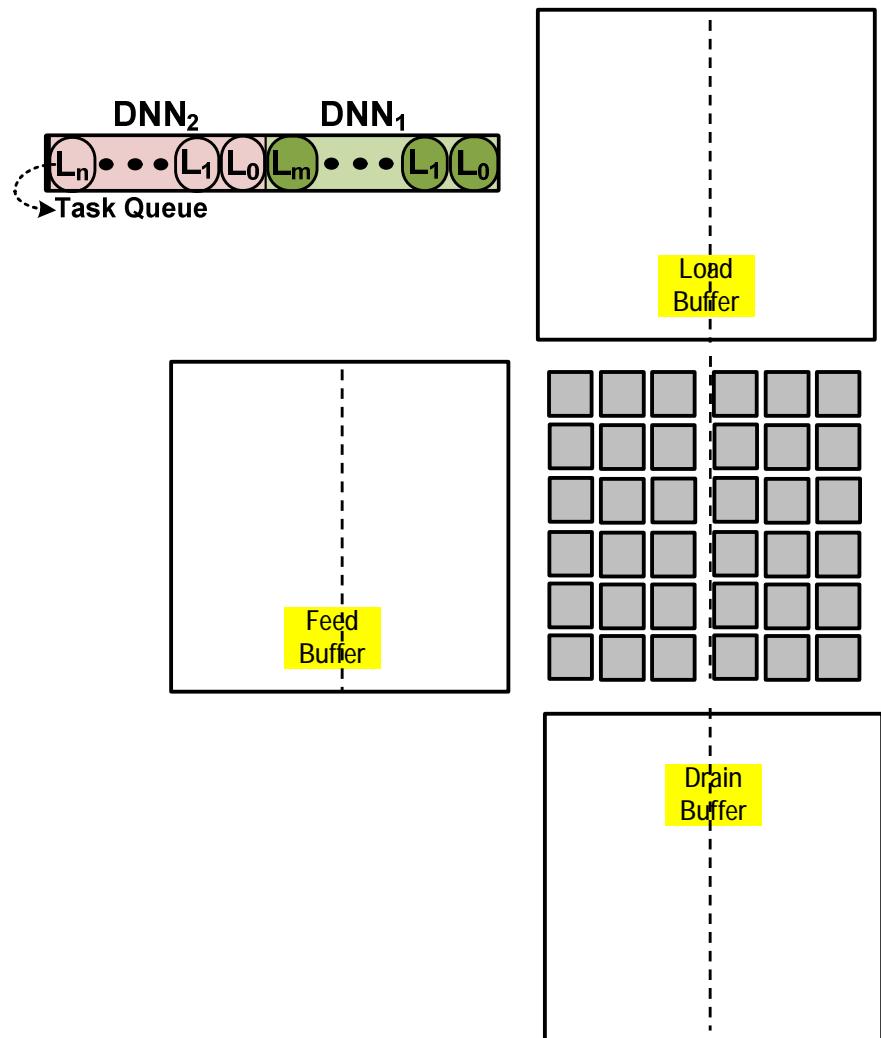
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

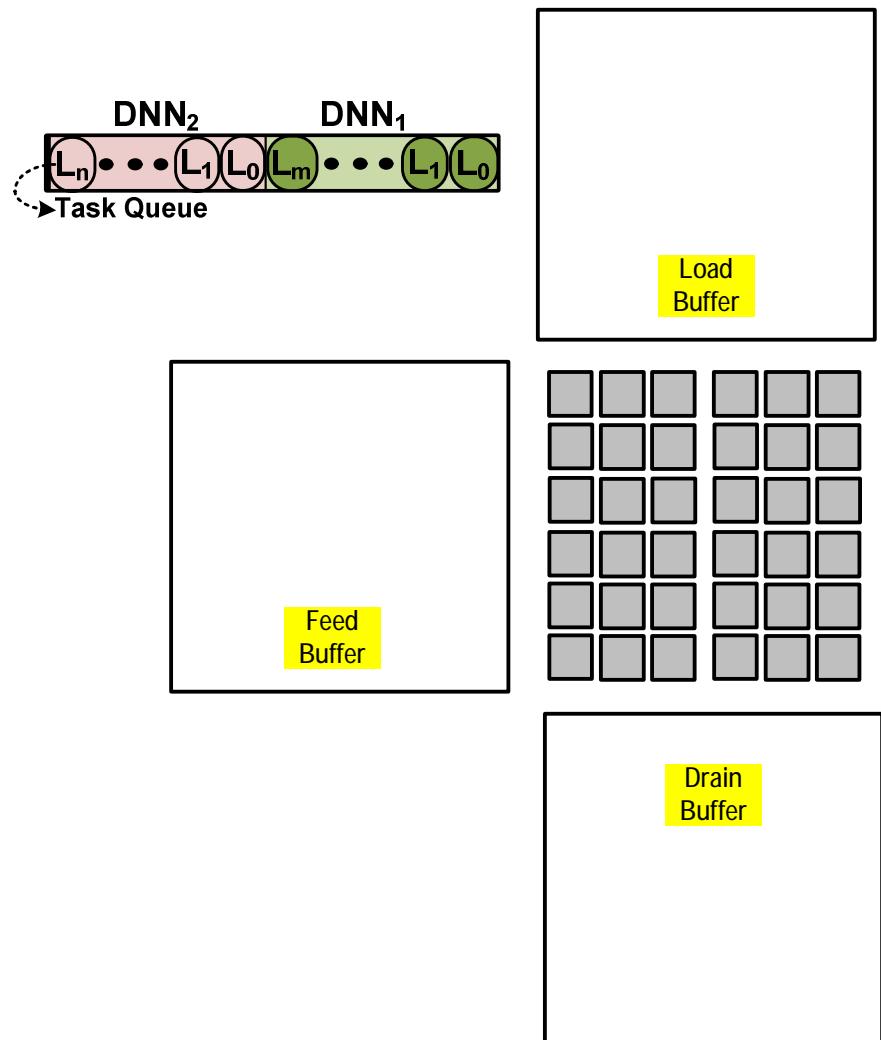
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

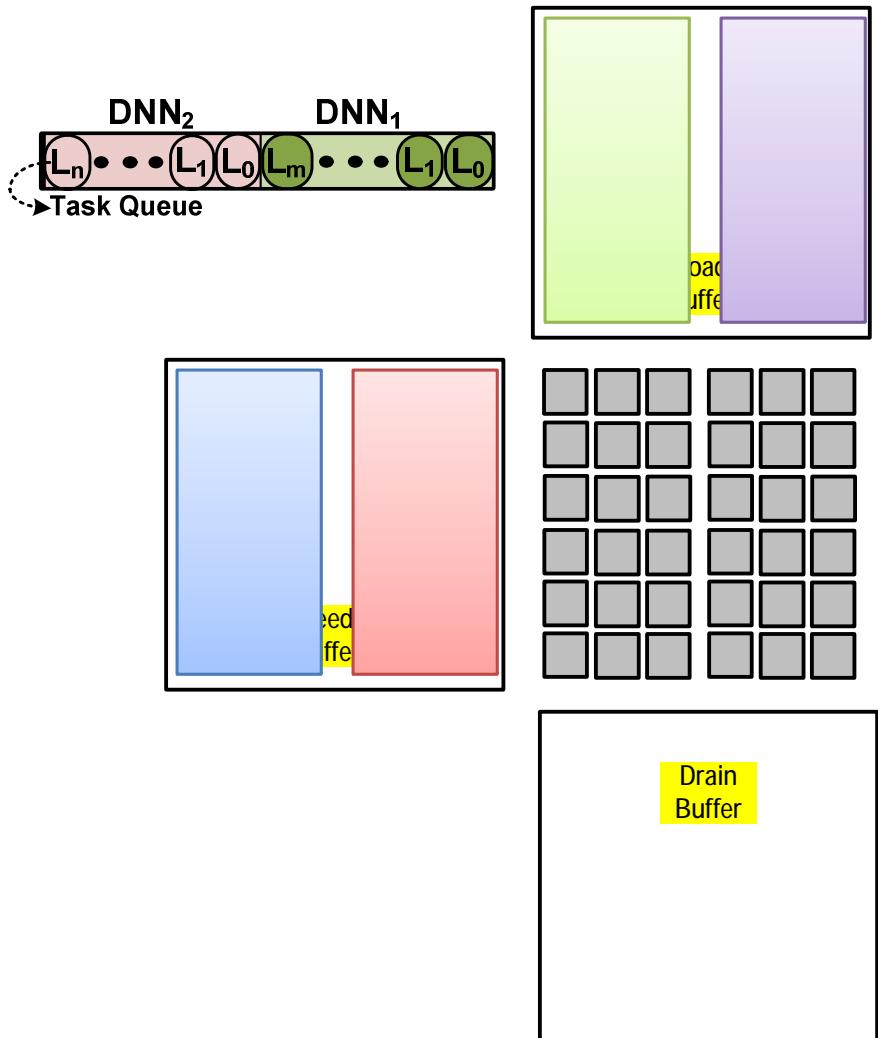
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ❶
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ❷
35:     Temporal_for PE [partition column] on PE [partition columny]
36:     Parallel_for PE [partition row] on PE [partition rowx]
37:     //step ❸
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

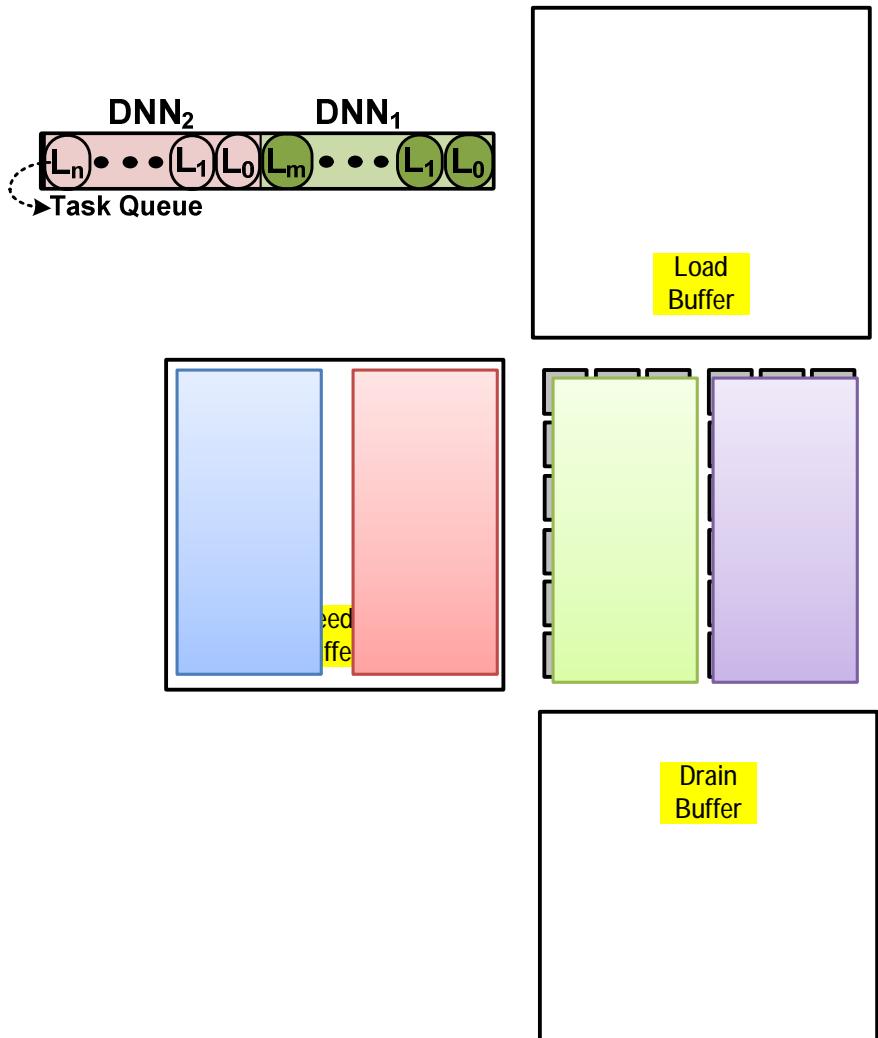
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ❶
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ❷
35:     Temporal_for Feed Buffer [partition column] on PE [partition column]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ❸
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

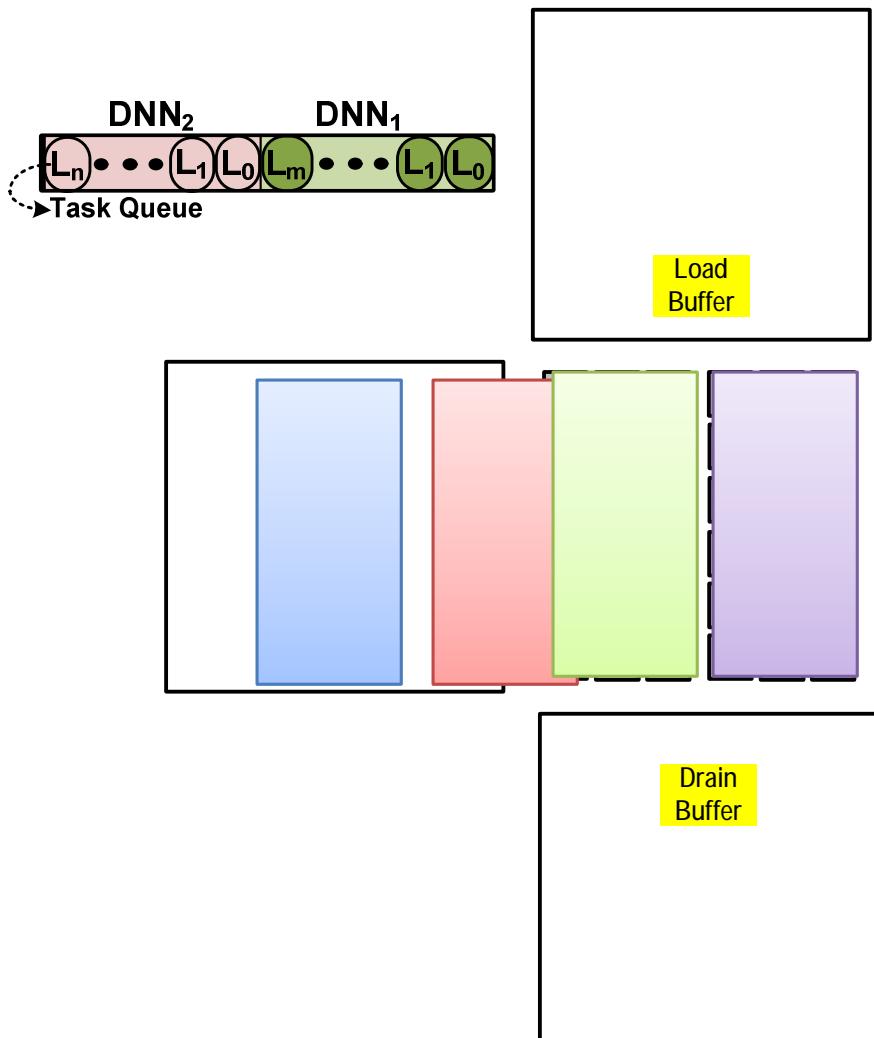
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition column]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition column]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

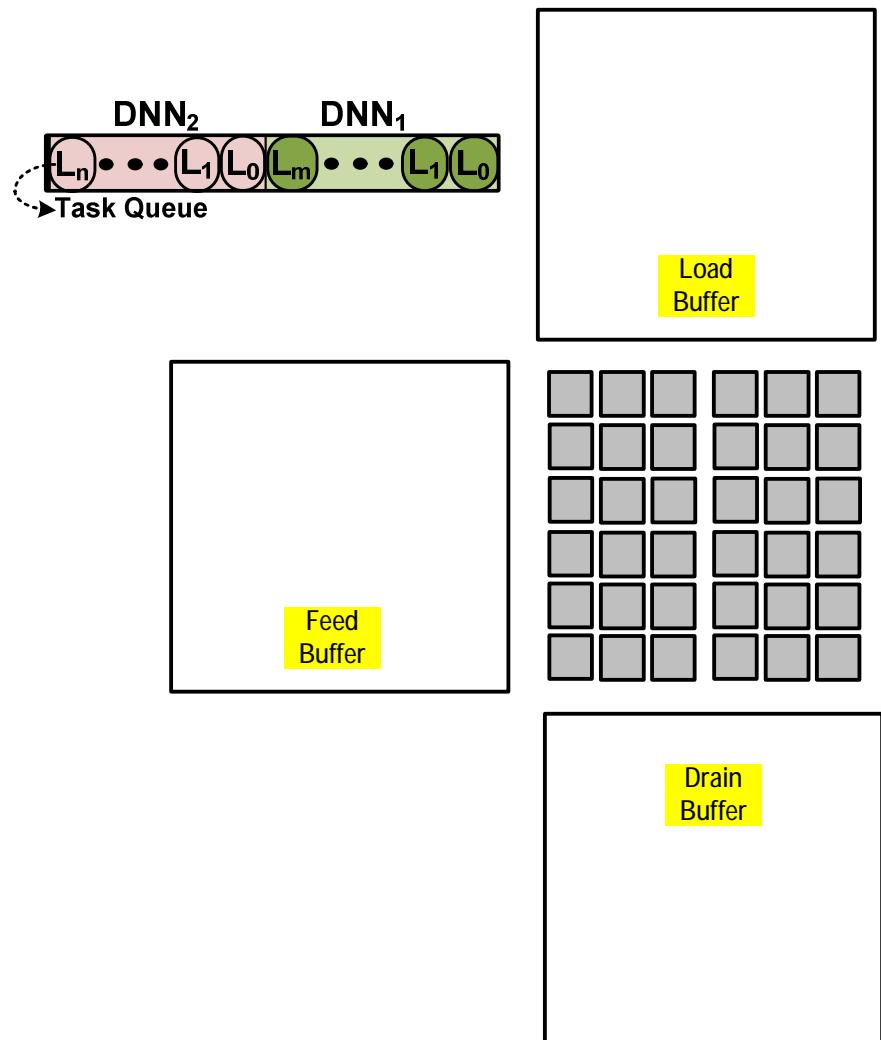
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

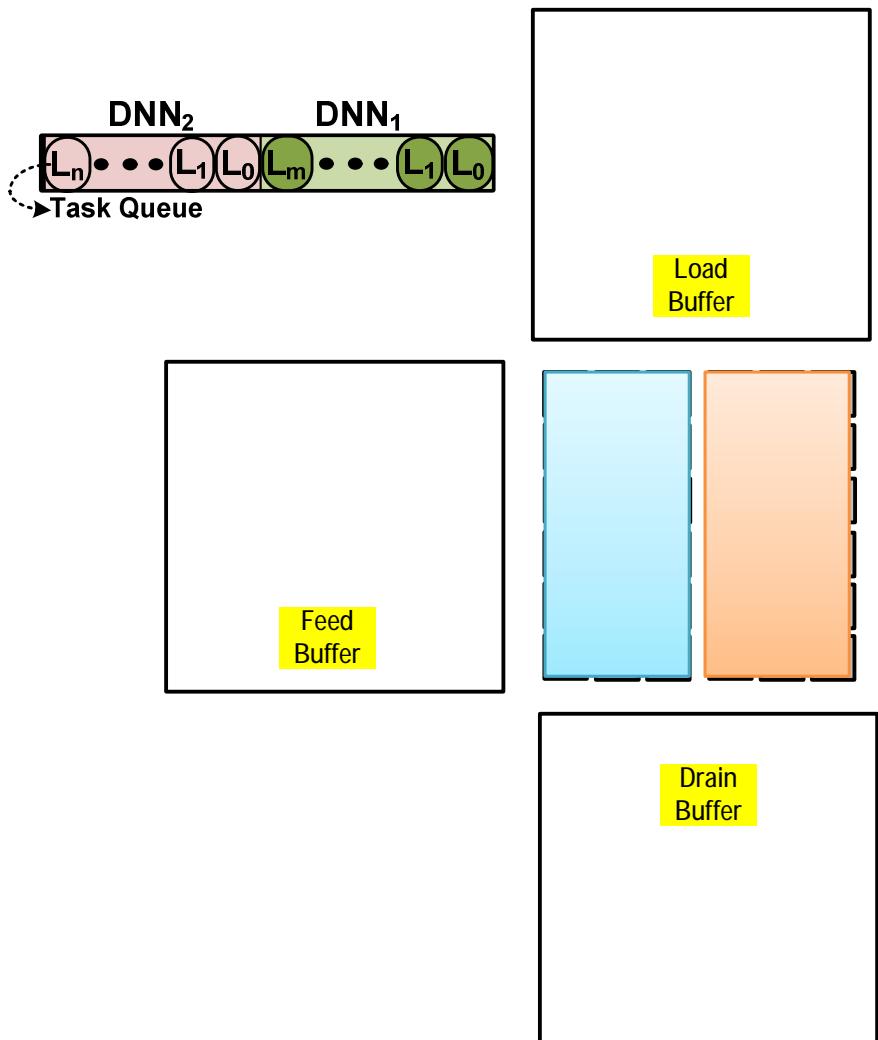
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



Dynamic Partitioning Algorithm

Algorithm1: Dynamic Resource Partitioning

```

1: Inputs: DNNG (m, n, At, Et)
   SA architecture (Storage(size), PE(x,y))
2: Outputs: Partition size estimation (PE(x',y'))
3:   Task assignment( $\tau_i \rightarrow PE(x',y')$ 

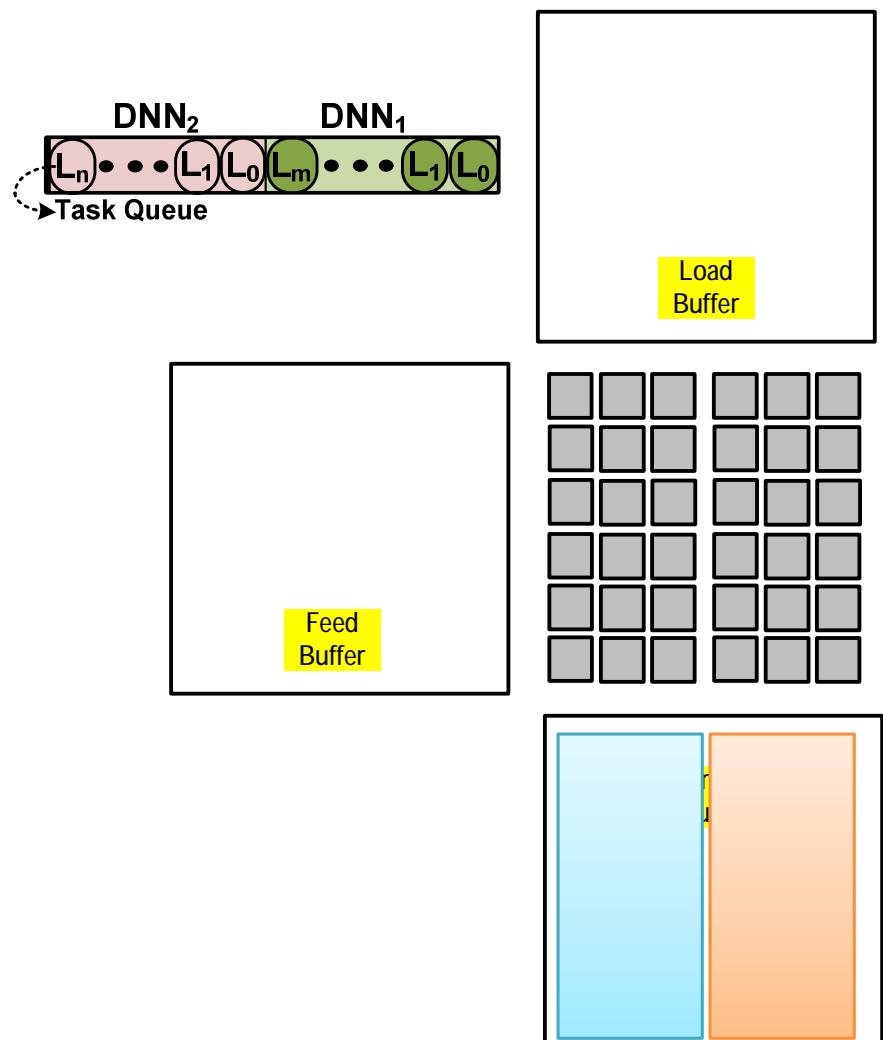
4: Queue [index]  $\leftarrow \emptyset$ 
5: if (index==1) //First DNNG inside queue
6:   Assign  $\tau_0$  to PE(x,y) //Assign to all PEs with no partitioning
7: else
8:   for All  $l_i$  with  $A_{ti} \leq E_{ti}$ 
9:     Number of partitions= Number of DNNGs inside Queue
10:    Call Partition_Calculation (Number of available layers)
11:    Call Task_Assignment (Available Layers)
12:    Call Partitioned_Weight_Stationery (Available Layers)
13:   end_for
14: end_if

15: Function Partition_Calculation (Number of available layers)
16:    $PE_{xi} = PE_x$  //Partition size in x dimension
17:    $PE_{yi} = \left\lfloor \frac{PE_y}{\text{Number of Available layers}} \right\rfloor$  //Partition size in y dimension
18:   return (PE(x',y'))
19: end_function

20: Function Task_Assignment (Available layers)
21:   for All available  $l_i$ 
22:     Number of operations of the layer  $l_i$ ,  $Opr(l_i) = \prod shapes$ 
23:     Sort  $Opr(l_i)$  from high to low
24:     Assign  $l_i$  with the highest Opr to available partition
25:   end_for
26:   return (Assigned  $l_i$ )
27: end_function

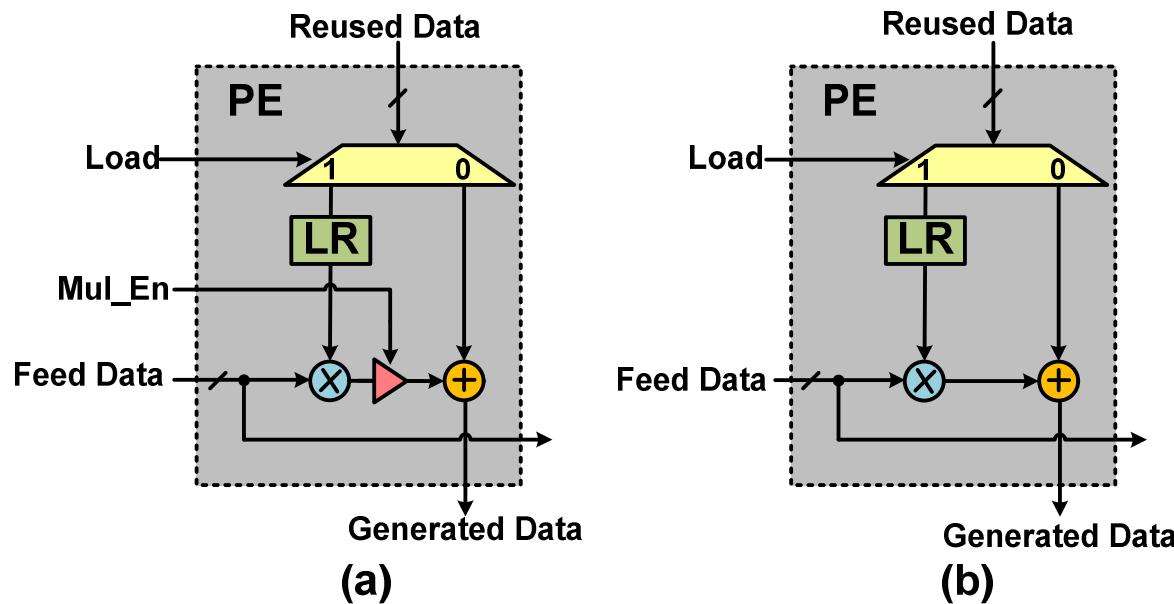
28: Function Partitioned_Weight_Stationery (Available layers)
29:   n = Number of available layers
30:   for All partitions n
31:     //step ①
32:     Temporal_for PE [partition columny] to Drain Buffer [partition column]
33:     Parallel_for PE [partition rowx] to Drain Buffer [partition row]
34:     //step ②
35:     Temporal_for Feed Buffer [partition column] on PE [partition columny]
36:     Parallel_for Feed Buffer [partition row] on PE [partition rowx]
37:     //step ③
38:     Parallel_for Load Buffer [partition row] to PE [partition rowx]
39:     Parallel_for Load Buffer [partition column] to PE [partition columny]
40:   end_parallel_temporal_for
41: end_for
42: end_function

```



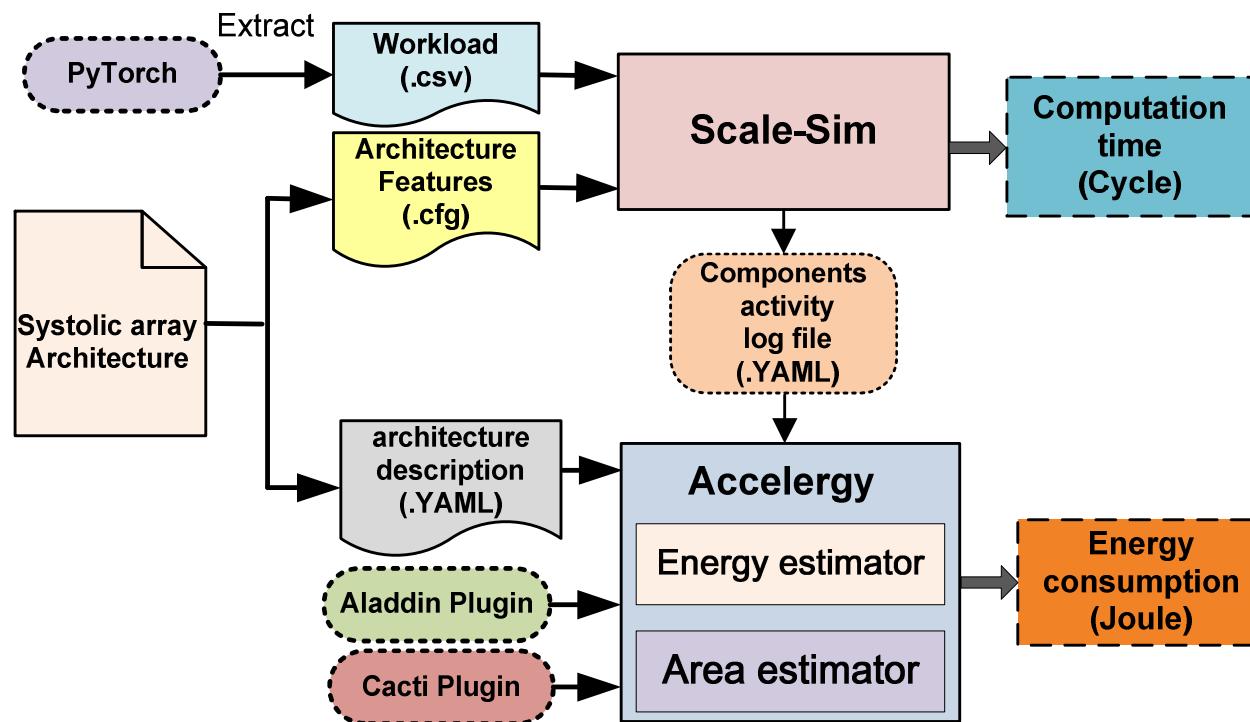
The proposed PE

- Systolic-array Based DNN Accelerator



The Simulation Toolchain

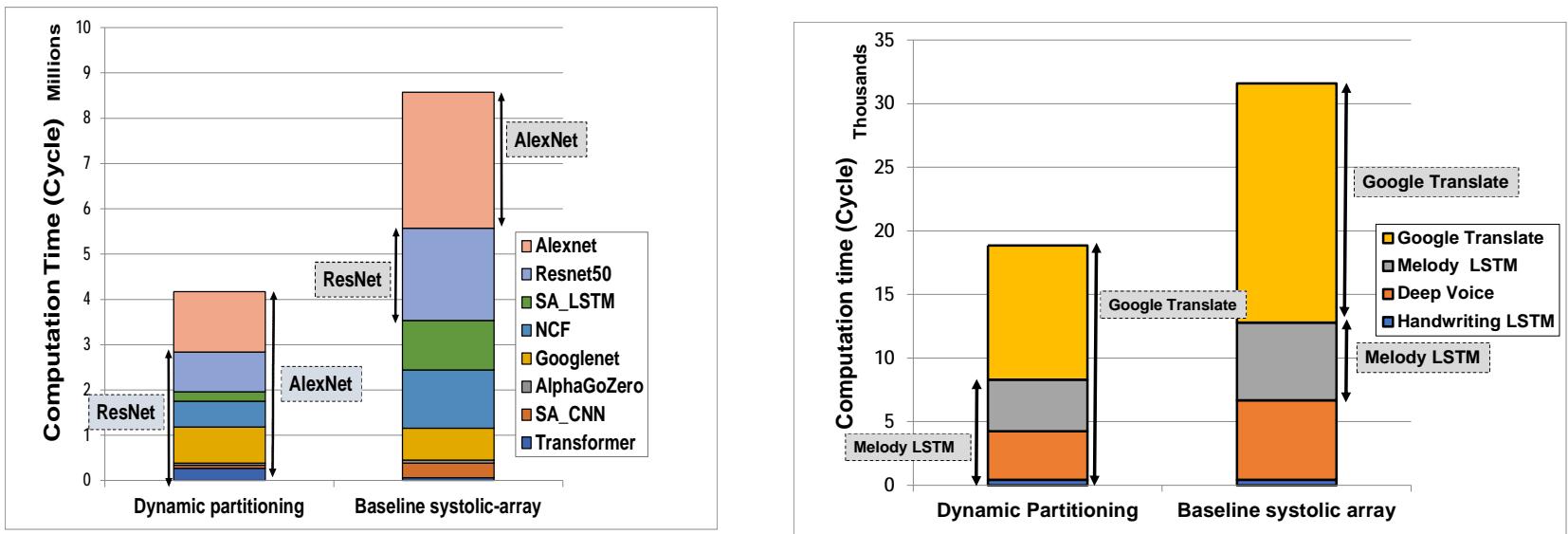
- Systolic-array Based DNN Accelerator



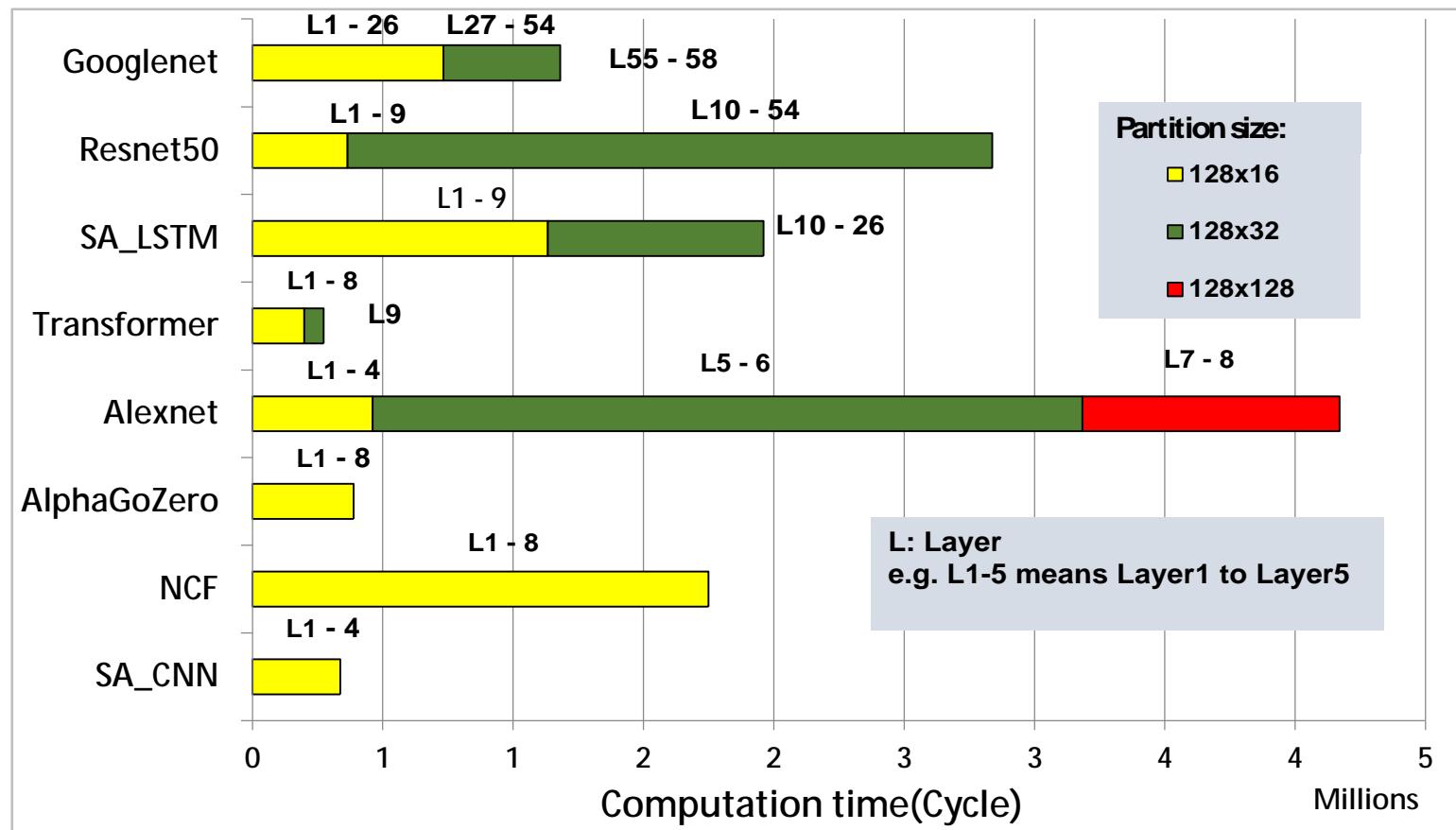
The Simulation Workloads

Type	Domain	Training model
Multi-domain / Heavy load workload	Image classification	AlexNet[16] ResNet50[17] GoogleNet[18]
	Sentiment analysis	Sentiment analysis CNN (SA_CNN)[19] Sentiment analysis LSTM (SA_LSTM)[20]
	Recommendation system	Neural collaborative filter (NCF)[21]
	Intelligent search	AlphaGoZero[22]
	Natural language processing	Transformer [23]
Recurrent Neural Network / Lightload workload	Melody extraction	Melody LSTM [24]
	Language translation	Google Translate [25]
	Text to speech	Deep voice[26]
	Handwriting recognition	Handwriting LSTM [27]

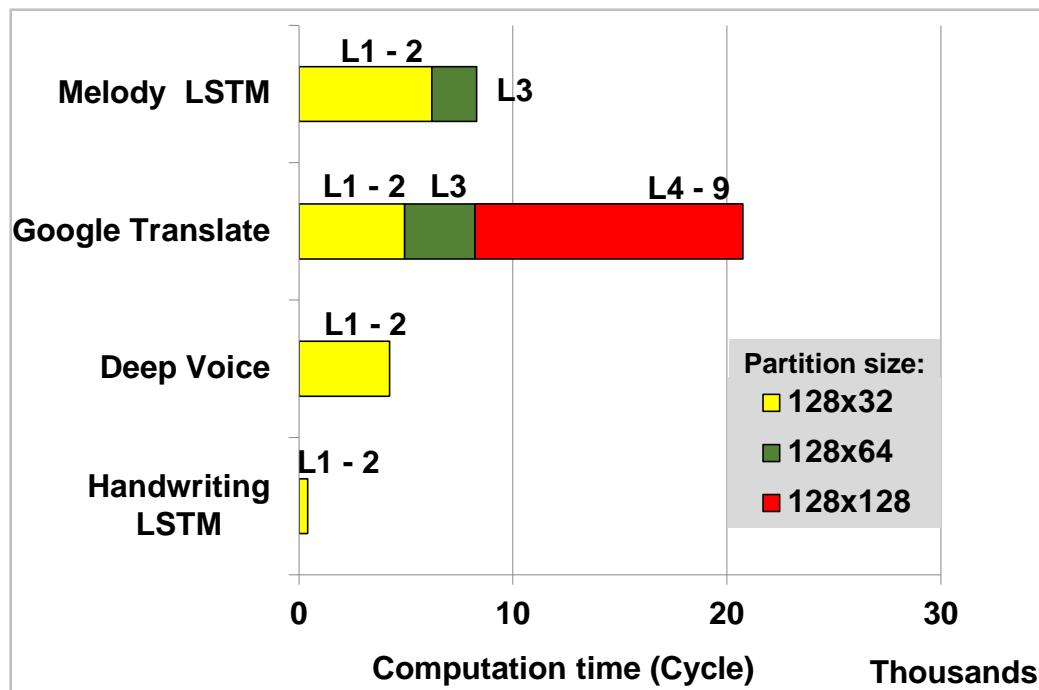
The Simulation Results



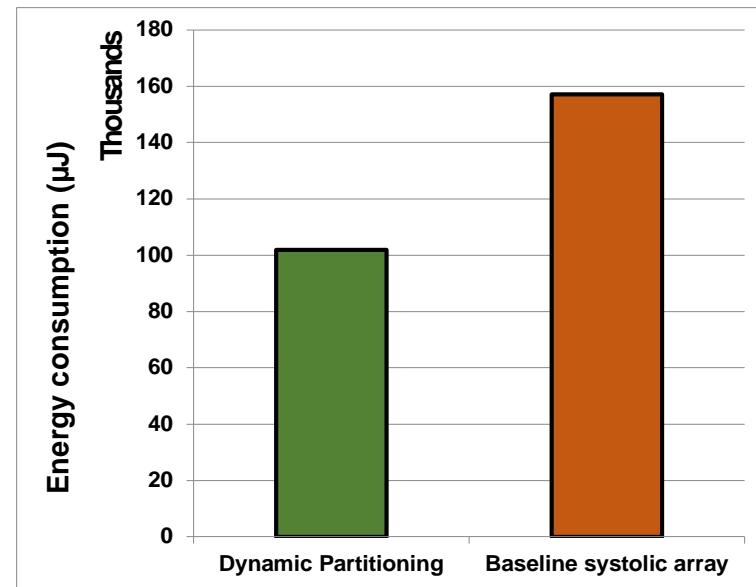
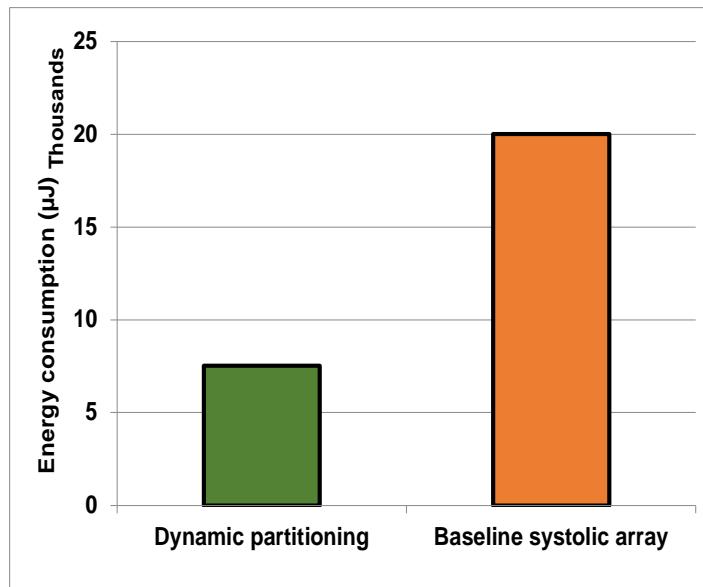
The Simulation workloads



The Simulation workloads



The Simulation workloads



Conclusion

- Concurrent running multiple DNNs
- Sharing accelerator's storage and processing resources
- Dynamic partitioning algorithm:
 - Increasing resource utilization
 - Reducing computation time
- Partitioned weight stationary dataflow
- Minor logic modification of systolic array

Thank you for your attention