

DRMap: A Generic DRAM Data Mapping Policy for Energy-Efficient Processing of Convolutional Neural Networks

Rachmad Vidya Wicaksana Putra, Muhammad Abdullah Hanif, Muhammad Shafique

Technische Universität Wien, Vienna, Austria

Email: {rachmad.putra, muhammad.hanif, muhammad.shafique}@tuwien.ac.at

Abstract—Many convolutional neural network (CNN) accelerators face performance- and energy-efficiency challenges which are crucial for embedded implementations, due to high DRAM access latency and energy. Recently, some DRAM architectures have been proposed to exploit subarray-level parallelism for decreasing the access latency. Towards this, we present a design space exploration methodology to study the latency and energy of different mapping policies on different DRAM architectures, and identify the pareto-optimal design choices. The results show that the energy-efficient DRAM accesses can be achieved by a mapping policy that orderly prioritizes to maximize the row buffer hits, bank- and subarray-level parallelism.

Index Terms—DRAM mapping, DRAM architectures, subarray-level parallelism, convolutional neural networks, CNNs, CNN accelerators.

I. INTRODUCTION

The widespread use of machine learning (ML) algorithms for organizing, analyzing, and inferring information from digital data is growing fast. Among many ML algorithms, convolutional neural network (CNN) algorithms have demonstrated state-of-the-art performance in data analytic tasks, such as image classification, object recognition, smart environment, health care, and automotive [1]. Since the CNN algorithms require data-intensive processing, CNN hardware accelerators are typically required to expedite the inference process. Over the past few years, several CNN accelerators have been proposed [2]–[12]. These accelerators offer higher performance- and energy-efficiency as compared to general-purpose CPUs. However, many CNN accelerators still face performance- and energy-efficiency challenges due to the high off-chip memory (i.e., DRAM) access latency and energy, which are higher than the latency and energy for the other compute operations [13]. *Therefore, reducing the DRAM access latency and energy are required for improving the performance- and energy-efficiency of CNN accelerators.*

A. The State-of-the-Art and Limitations

Previous works have proposed different techniques to reduce the DRAM access energy, by minimizing the number of DRAM accesses [3] [14] [15]. Their main ideas are similar, i.e., (i) defining layer partitioning¹ and then transferring each partition from DRAM to on-chip memory/buffer in a defined schedule, and (ii) maximally reusing the data that are already in the on-chip buffer. The state-of-the-art [14] considers adaptive layer partitioning and scheduling, to minimize the number of DRAM accesses, by adaptively switching the reuse priority between different data types: input activations/feature maps (*ifms*), output activations/feature maps (*ofms*), and weights (*wghs*), across the layers of a network. *Although all these works result in a reduced number of DRAM accesses (which also means reduced DRAM access energy), they do not consider improving (i) the DRAM latency-per-access, and (ii) DRAM energy-per-access. Therefore, performance- and energy-efficiency improvements achieved by the state-of-the-art are sub-optimal, thereby limiting the CNN accelerators to achieve*

¹Layer partitioning determines portions of data in the form of block/tile to be accessed from DRAM to on-chip memory at one time. A detailed explanation is provided in Section II-A.

further performance- and energy-efficiency improvements. We will illustrate this with the help of the following motivational case study.

B. Motivational Case Study and Associated Research Challenges

Motivational Case Study: Although there are various types of commodity DRAM (e.g., DDR3, DDR4, etc.), they have similar internal organization and operations [16] (detailed DRAM organization and operations are provided in Section II-B). Therefore, different types of commodity DRAM have similar behavior regarding latency-per-access and energy-per-access. The DRAM latency-per-access and energy-per-access vary depending upon whether a single DRAM access faces a *row buffer hit*, a *row buffer miss*, or a *row buffer conflict*. A row buffer hit means that the requested row is already available in the row buffer, hence the data access can be performed directly without additional operations. In case of a row buffer miss or conflict, the requested row has to be opened first before a data access can be performed. In this manner, a row buffer miss and conflict require higher latency-per-access and energy-per-access than a row buffer hit. To illustrate this, we performed an experimental analysis to observe the DRAM latency-per-access and energy-per-access for different conditions (i.e., a row buffer hit, miss, and conflict), and the experimental results are presented in Fig. 1. Furthermore, in commodity DRAM, each request that goes to a DRAM bank can only access a single DRAM subarray at a time, although each bank is composed of multiple subarrays. This limits the DRAM capability to offer lower DRAM access latency and energy. Recently, several DRAM architectures that offer subarray-level parallelism (SALP) in a DRAM bank, have been proposed in the literature. In [17], three variants of SALP architectures are presented, i.e., SALP-1, SALP-2, and SALP-MASA (detailed SALP architectures are provided in Section II-C). Our observation results in Fig. 1 show that SALP architectures have the potential to further reduce the DRAM latency-per-access and energy-per-access as compared to commodity DRAM.

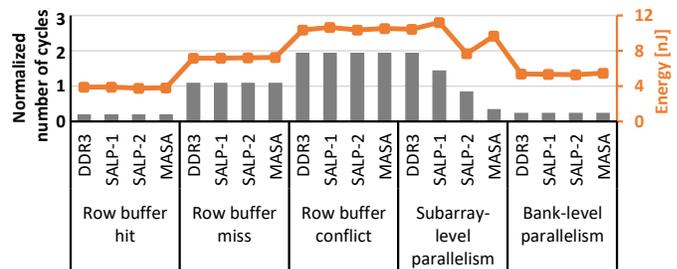


Fig. 1. DRAM latency-per-access and DRAM energy-per-access for different conditions (a row buffer hit, a row buffer miss, a row buffer conflict, subarray- and bank-level parallelism) in different DRAM architectures (DDR3, SALP-1, SALP-2, and SALP-MASA). Data are obtained from our experiments using a state-of-the-art cycle-accurate DRAM simulators [18] [19] for DDR3-1600 2Gb x8 and SALP 2Gb x8 with 8 subarrays-per-bank.

Associated Research Challenges: From above observations, the energy efficiency of DRAM accesses for CNN accelerators can be

improved by minimizing the DRAM latency-per-access and energy-per-access. Therefore, there is a need of a generic DRAM mapping policy that can achieve maximum row buffer hits while exploiting subarray- and bank-level parallelism. Furthermore, to justify that the proposed DRAM mapping policy is applicable to different design choices, a design space exploration (DSE) is required. This DSE explores different DRAM mapping policies in different DRAM architectures with different layer partitioning and scheduling schemes, to find the minimum energy-delay-product (EDP) of DRAM accesses. This EDP is used as a measure of the energy-efficiency of a CNN accelerator. Therefore, an analytical model for estimating the EDP of different DRAM mapping policies in the DSE, is also required.

C. Our Novel Contributions

In this paper, we make the following novel contributions (the overview is illustrated in Fig. 2) to overcome the associated challenges.

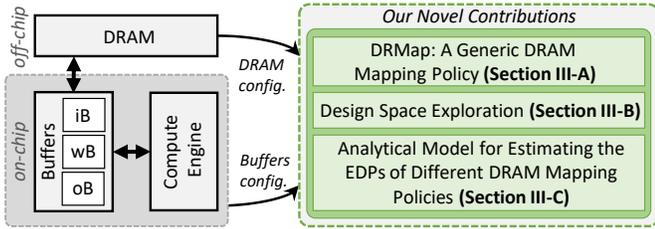


Fig. 2. The overview of our novel contributions, highlighted in the green boxes. We use separate on-chip buffers for different data types: input buffer (iB) for *ifms*, weight buffer (wB) for *wghs*, and output buffer (oB) for *ofms*.

- 1) We propose **DRMap: a generic DRAM data Mapping policy** that offers minimum energy-delay-product (EDP) of DRAM accesses, for a given DRAM architecture, layer partitioning and scheduling scheme. DRMap orderly prioritizes to maximize row buffer hits, bank- and subarray-level parallelism.
- 2) We propose a **design space exploration (DSE) algorithm** to find a DRAM mapping that offers minimum EDP, while considering different DRAM architectures, different layer partitioning and scheduling schemes.
- 3) We propose an **analytical model for estimating EDPs of different DRAM mapping policies**, which will be used in the DSE. The EDP for each DRAM mapping is estimated by multiplying the number of DRAM accesses with the respective number of cycles and energy values.

Key results: DRMap orderly prioritizes to maximize the row buffer hits, bank- and subarray-level parallelism. It improves the EDP compared to the other mapping policies, up to 96% for DDR3, 94% for SALP-1, 91% for SALP-2, and 80% for SALP-MASA on AlexNet [20].

II. PRELIMINARIES

A. Layer Partitioning and Scheduling in CNNs

The full CNN processing usually cannot be mapped at once on the accelerator fabric due to the limited on-chip buffer capacity (i.e., 100KB-500KB [13]), hence layer partitioning and scheduling are required. To illustrate this, a pseudo-code of a convolutional layer processing in a CNN accelerator is shown in Fig. 3. It has two parts, i.e., inner loops and outer loops. The inner loops represent the on-chip processing. The outer loops represent the scheduling of processing different portions of data (from all data types: *ifms*, *wghs*, and *ofms*), whose sizes have to be less than or equal to the sizes of respective buffers (iB, wB, and oB). These data are partitioned in the form of blocks/tiles which are represented with the step sizes. Furthermore,

the sequence of the outer loops represents the order in which the tiles are accessed from DRAM to the on-chip buffer. It thereby reflects the number of DRAM accesses required to process a layer of a network.

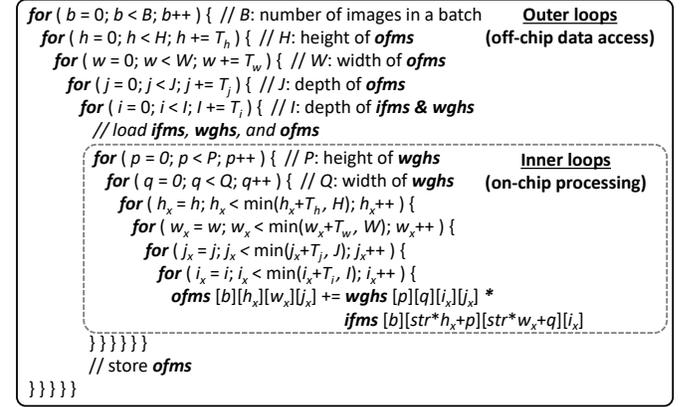


Fig. 3. Pseudo-code of the tiled convolutional neural network processing.

B. DRAM Overview

DRAM Organization: From top to bottom perspective, the organization of a commodity DRAM is composed of *channel*, *rank*, *chip*, *bank*, *row*, and *column* [16] [21], as shown in Fig. 4(a). In commodity DRAM, banks are the lowest hierarchy in DRAM, which can be accessed in parallel, and referred to as *bank-level parallelism* [22]. Actually, a DRAM bank is *not* implemented in a monolithic design (a large array of cells with a single row buffer). Instead, it is implemented in multiple *subarrays*, each of which has its local row buffer, as shown in Fig 4(b). Multiple subarrays in a bank share (i) global bitlines, which connect local row buffers to a global row buffer, and (ii) a global row address decoder [17].

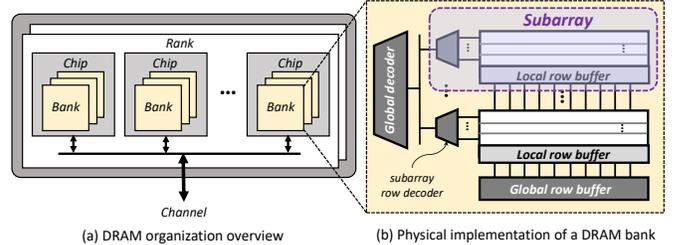


Fig. 4. (a) DRAM organization overview, and (b) physical implementation of a DRAM bank, showing multiple subarrays in a bank.

DRAM Operations: If there is a single DRAM request, a rank will respond and thereby multiple chips within this rank can be accessed in parallel, contributing to a DRAM word. For each chip, the request is routed to a specific bank and decoded into addresses of a row and a column. The *Activation (ACT)* command triggers a row activation, and data from the requested row are copied to the row buffer. Afterward, a *read (RD)* or *write (WR)* command can be issued to a specific column in the activated row buffer. If the requested row is already activated, then the data in this row is already in the row buffer (*a row buffer hit*). Therefore, it does not need a new row activation. If the requested row is not activated yet, then it is either *a row buffer miss* or *a row buffer conflict*. In a *row buffer miss*, there is no activated row yet in the row buffer and it requires to activate the requested row. Meanwhile, in a *row buffer conflict*, there is an activated row in the row buffer, but it is not the one that the request is expecting. Therefore, this condition requires to close the activated row first using the *precharging (PRE)* command, and then activate the requested row using the *activation (ACT)* command.

DRAM Data Mapping: The default data mapping prioritizes to map the subsequently accessed data in the different columns of the same row of a bank (for increasing row buffer hits) and the different banks of the same rank (for exploiting bank-level parallelism). However, it does not exploit subarray-level parallelism and does not consider different possible layer partitioning and scheduling. Therefore, the default data mapping solution is suboptimal.

C. DRAM Architectures that Exploit Subarray-level Parallelism

In a commodity DRAM, each request that goes to a DRAM bank, can only access a single subarray at a time. This limits the potential to reduce the DRAM access latency and energy. To address this limitation, [17] has proposed three DRAM architectures and mechanisms that exploit subarray-level parallelism (SALP) in the same bank, called SALP-1, SALP-2, and SALP-MASA. Following are the key ideas of these SALP architectures.

- **SALP-1** reduces the DRAM service time by overlapping the *precharging* of one subarray with the *activation* of another subarray, since mostly the *precharging* and *activation* are local to a subarray. To enable this mechanism, re-interpretation of the existing timing constraint for *precharging*, is required.
- **SALP-2** reduces the DRAM service time even more than SALP-1, by overlapping the *write-recovery* latency of an active subarray, with the *activation* of another subarray. To enable this, additional circuitry to activate two subarrays at the same time is required.
- **Multitude of Activated Subarrays (MASA)** reduces the DRAM service time even more than SALP-2, by activating multiple subarrays at the same time (the *activations* of different subarrays are overlapped). To enable this, additional circuitry (more than SALP-2) to activate multiple subarrays at the same time is required.

III. OUR DESIGN METHODOLOGY FOR DRAM MAPPING IN CNN ACCELERATORS

A. DRMap: A Generic DRAM Data Mapping Policy

Our observations from the results in Fig. 1 show that different DRAM architectures have similar behavior in terms of latency-per-access and energy-per-access. Therefore, we propose DRMap, a generic DRAM mapping policy for energy-efficient DRAM accesses in CNN accelerators. Its main idea is to orderly prioritize the data mapping that maximizes DRAM row buffer hit, bank- and subarray-level parallelism. The flowchart of DRMap mechanism in a DRAM chip is presented in Fig. 5, while its pseudo-code and physical representation of mapping policy are illustrated in Fig. 6. *DRMap* considers *tile-based partitioning* in its mechanism, thereby *DRMap* can be performed for each data tile using the following steps:

- 1) If we consider accessing a DRAM bank, then DRMap prioritizes to map a data partition to different columns in the same row to achieve maximum row buffer hits. If multiple chips are available within a rank, then this step can be performed in different chips for exploiting the chip-level parallelism.
- 2) If all columns in the same row of a bank are fully filled, then the remaining data are mapped to different banks in the same chip, to exploit bank-level parallelism. If multiple chips are available, then this step can be performed in different chips.
- 3) If all columns in the same row of all banks are fully filled, then the remaining data are mapped to a different subarray in the same bank, to exploit subarray-level parallelism. If multiple chips are available, then this step can be performed in different chips.
- 4) If there are remaining data left, then step 1) to 3) can be performed again for different subarray, until all data are mapped within the same rank. In this manner, DRMap can achieve maximum row

buffer hits, while maximally exploiting bank- and subarray-level parallelism within a DRAM rank.

- 5) If there are remaining data left, they can be mapped in different rank (channel) if available, using the same steps as 1) to 4). In this manner, our DRMap can achieve maximum row buffer hits, while maximally exploiting bank- and subarray-level parallelism in another DRAM rank (channel) as well.

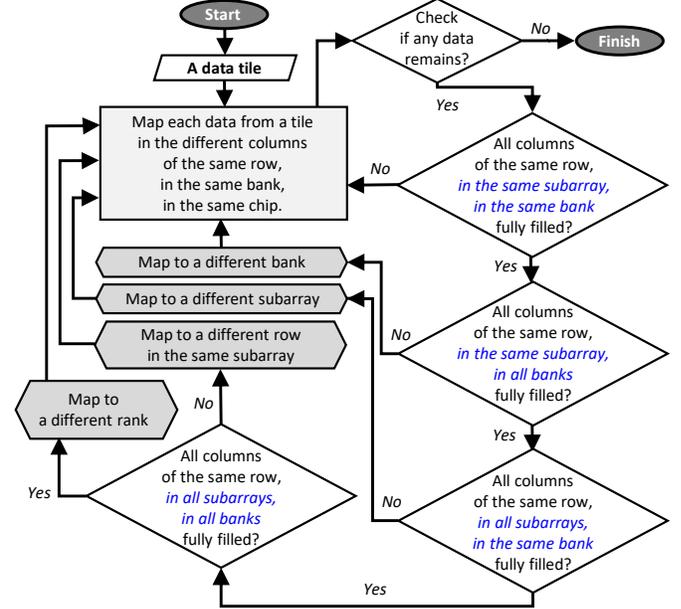


Fig. 5. Flowchart of the DRMap that illustrates how the mapping policy is performed in a DRAM chip.

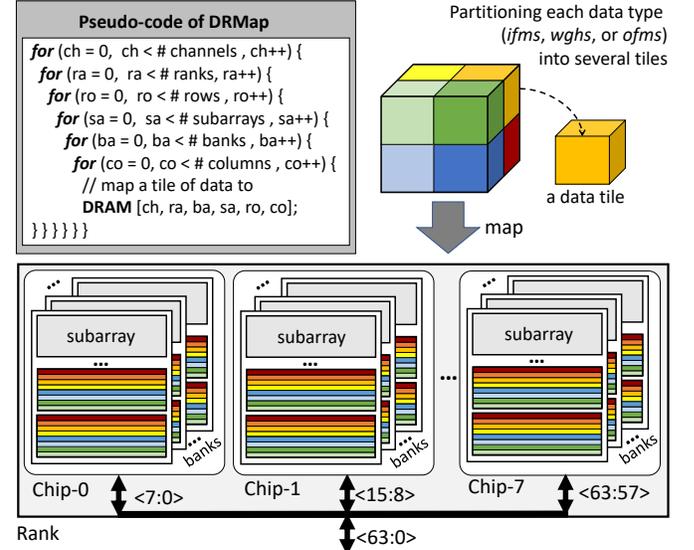


Fig. 6. Pseudo-code of DRMap and its conceptual implementation in DRAM.

To illustrate that our DRMap always achieves the minimum energy-delay-product (EDP) of DRAM accesses in different possible conditions, we perform an extensive design space exploration (DSE). The DSE investigates different DRAM mapping policies, different DRAM architectures, as well as different layer partitioning and scheduling schemes on CNN, and estimates EDP for these different combinations. *This DSE is important to corroborate that the best solution that provides the minimum EDP in each given combination is always the same as provided by our DRMap technique.*

B. Design Space Exploration for Evaluating Different DRAM Mapping Policies

To evaluate the impact of different DRAM mapping policies and see the performance of DRMap as compared to others, we performed an extensive design space exploration (DSE). An overview of the DSE is shown in Fig. 7 and its algorithm is presented in Algorithm 1.

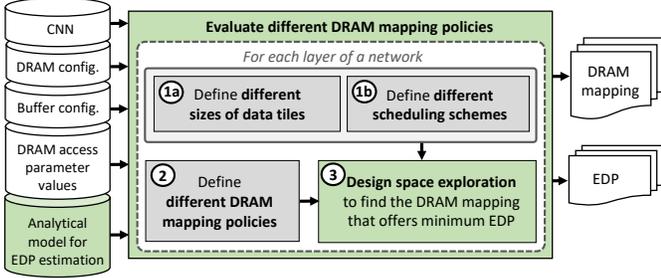


Fig. 7. The operational flow of the our DSE methodology. Our contributions are highlighted in the green boxes.

Algorithm 1 Pseudo-code of the proposed DSE algorithm

```

INPUT: (1) CNN configuration: number of layers ( $L$ );
          (2) Buffer size:  $ifms$  ( $iB$ ),  $wghs$  ( $wB$ ),  $ofms$  ( $oB$ );
          (3) Analytical models of EDP ( $EDP$ );
          (4) Layer partitioning for  $ifms$ ,  $wghs$ , and  $ofms$  ( $Partitioning$ );
          (5) DRAM access scheduling ( $Scheduling$ );
          (6) DRAM mapping policies ( $DRAMmaps$ );
OUTPUT: (1) Efficient DRAM mapping ( $map$ );
           (2) Minimum EDP ( $minEDP$ );
BEGIN
  Initialization:
  1:  $T_p = P$ ;
  2:  $T_q = Q$ ;
  3:  $EDP[] = 0$ ;
  4:  $minEDP[] = 0$ ;
  Process:
  5: for ( $l = 1$  to  $L$ ) do
  6:   for (each  $Partitioning$ ) do
  7:     for (each  $Scheduling$ ) do
  8:       for (each  $DRAMmaps$ ) do
  9:         if ( $ifms$  tile size  $\leq iB$ ) and ( $wghs$  tile size  $\leq wB$ ) and
           ( $ofms$  tile size  $\leq oB$ ) then
 10:          Calculate  $EDP[l]$ ;
 11:          if (first loop) then
 12:             $minEDP[l] = EDP[l]$ ;
 13:          else if ( $EDP[l] \leq minEDP[l]$ ) then
 14:             $minEDP[l] = EDP[l]$ ;
 15:            Save  $map$ ,  $minEDP$ ;
 16:          end if
 17:        end if
 18:      end for
 19:    end for
 20:  end for
 21: end for
 22: return (1)  $map$ ; (2)  $minEDP$ ;
END

```

For each layer of a network, the DSE performs three key steps: (1) defining different sizes of data tiles and scheduling schemes, (2) defining different DRAM mapping policies, (3) performing the DSE to find a DRAM mapping policy that offers minimum EDP. The operational flow of the DSE is explained in the following points:

Step-1. Define (1a) different sizes of data tiles for all data types ($ifms$, $wghs$, and $ofms$), and (1b) different scheduling schemes. The tile sizes are determined by the step sizes in the outer loops of the Fig. 3. The tile sizes of $ifms$, $wghs$, and $ofms$ have to fit in the corresponding buffers (iB , wB , and oB). Each combination of the tile sizes for all data types defines one possible partitioning, which will be considered in the DSE. The scheduling schemes are determined by the sequence of the outer loops of the Fig. 3. In this work, we consider four scheduling schemes, based on the reuse priority of the data type: $ifms$ -reuse, $wghs$ -reuse, $ofms$ -reuse, and $adaptive$ -reuse scheduling schemes. The $ifms$ -reuse scheduling means that $ifms$ data type will be maximally reused when the data are available in the on-chip buffer. Similar definition is also applied for $wghs$ -reuse and $ofms$ -reuse. Meanwhile, the $adaptive$ -reuse scheduling means that the reuse priority changes across different layers of a network, according to which one among $ifms$ - $wghs$ - $ofms$ -reuse scheduling that offers minimum number of DRAM accesses.

Step-2. Define different DRAM mapping policies, by determining the different orders of mapping loops to different columns, rows, subarrays, and banks in the same DRAM chip. For DDR3, orders of mapping loops are permutation of banks, rows, and columns, in the same DRAM chip; meanwhile for SALP, orders of mapping loops are permutation of banks, subarrays, rows, and columns, in the same DRAM chip. Here, we narrow down the design space by selecting the DRAM mapping policies that have the least frequent subsequent accesses to different rows, since it is the most expensive access in the same DRAM chip, for both latency and energy (as validated by Fig. 1). Therefore, there are six mapping policies to be explored in the DSE, as presented in Table I.

Step-3. Perform the DSE to find a DRAM mapping policy that offers minimum EDP, across different DRAM architectures, different layer partitioning and scheduling schemes. The minimum EDP and the corresponding DRAM mapping are the outputs of the DSE, for a given DRAM architecture, layer partitioning and scheduling.

Note that the time and energy are already included in the DSE, for determining the EDP in the final results. For each layer of a network, EDP is obtained by multiplying the DRAM access energy and latency consumed by each combination of different DRAM mapping policies, different DRAM architectures, as well as different sizes of layer partitioning and scheduling schemes. Therefore, DSE will be able to find the combination that offers minimum EDP for each layer of a network and minimum total EDP for a whole network.

TABLE I
DIFFERENT DRAM MAPPING POLICIES FOR THE DSE.

Mapping	Inner-most- to outer-most-loops
1	column, subarray, bank, row
2	subarray, column, bank, row
3	column, bank, subarray, row
4	bank, column, subarray, row
5	subarray, bank, column, row
6	bank, subarray, column, row

C. Analytical Model of Energy-Delay-Product (EDP) Estimation for Different DRAM Mapping Policies

Based on the proposed DSE, the optimization problem is formulated to minimize the EDP of DRAM accesses for each layer of a network and can be stated as

$$\text{Objective : minimize } (EDP_{layer}) \quad (1)$$

The EDP-per-layer (EDP_{layer}) is obtained by multiplying the energy-per-layer and latency-per-layer. The energy-per-layer is ob-

tained by accumulating all access energy values incurred from the DRAM accesses for all data tiles. The latency-per-layer is obtained by accumulating all access latency values incurred from the DRAM accesses for all data tiles. The access latency and energy are calculated on the basis of DRAM accesses for each data tile since we consider layer partitioning approach. Therefore, for each tile, the number of cycles required for DRAM accesses can be formulated as Eq. 2 and the DRAM access energy can be formulated as Eq. 3.

$$N_{cycle_{tile}} = N_{access_{dif_column}} \cdot N_{cycle_{dif_column}} + N_{access_{dif_rows}} \cdot N_{cycle_{dif_rows}} + N_{access_{dif_subarrays}} \cdot N_{cycle_{dif_subarrays}} + N_{access_{dif_banks}} \cdot N_{cycle_{dif_banks}} \quad (2)$$

$$E_{tile} = N_{access_{dif_column}} \cdot E_{dif_column} + N_{access_{dif_rows}} \cdot E_{dif_rows} + N_{access_{dif_subarrays}} \cdot E_{dif_subarrays} + N_{access_{dif_banks}} \cdot E_{dif_banks} \quad (3)$$

Term $N_{access_{dif_x}}$ denotes the number of accesses to different DRAM- x . $N_{cycle_{dif_x}}$ denotes the number of cycles incurred when accessing different DRAM- x . E_{dif_x} denotes the access energy incurred when accessing different DRAM- x . For all terms, $x \in \{\text{columns, rows, subarrays, banks}\}$.

IV. EVALUATION METHODOLOGY

To evaluate our proposed methodology, we built the experimental setup, as presented in Fig. 8.

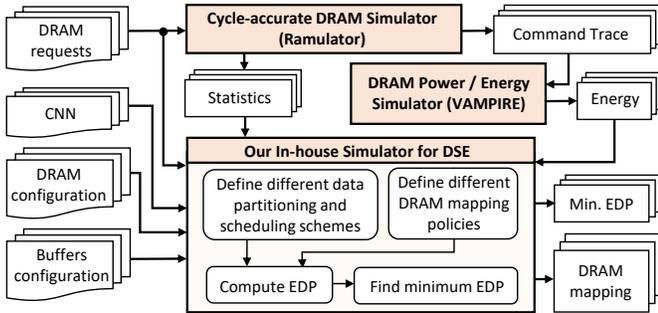


Fig. 8. Experimental setup and tool flow.

TABLE II
CONFIGURATION OF THE CNN ACCELERATOR.

Module	Description
CNN Processing Array	Size = 8×8 MACs
On-chip Buffers	iB : 64KB, wB : 64KB, oB : 64KB
Memory Controller	Policy = open row, scheduler = FCFS
DDR3-1600	Configuration: 2Gb x8 1 channel, 1 rank-per-channel, 1 chip-per-rank, 8 banks-per-chip
SALP	Configuration: 2Gb x8 1 channel, 1 rank-per-channel, 1 chip-per-rank, 8 banks-per-chip, 8 subarrays-per-bank

Tool flow: We used a cycle-accurate DRAM simulator, *Ramulator* [18], to obtain the statistics (i.e, number of cycles) of different DRAM access conditions: a row buffer hit, row buffer miss, row buffer conflict, subarray- and bank-level parallelism. To profile the energy, we used a real experiments-based DRAM energy simulator, *VAMPIRE* [19]. Information of energy and number of cycles are used for the DSE, which considers different DRAM mapping policies, different DRAM architectures, different layer partitioning and scheduling schemes to find the DRAM mapping policy that offers

minimum EDP. For DSE, we considered a state-of-the-art Tensor Processing Unit (TPU) [7]-like CNN accelerator with a reduced size of on-chip buffers and MAC array engine, as specified in Table II. To represent different DRAM architectures, we used DDR3 and SALP architectures (SALP-1, SALP-2, and SALP-MASA). For scheduling, we considered *ifms-reuse*, *wghs-reuse*, *ofms-reuse*, and *adaptive-reuse* scheduling schemes. For mapping, we considered the six mapping policies presented in Table I. For the input, we used AlexNet [20] with ImageNet dataset.

V. RESULTS AND DISCUSSIONS

A. Comparisons of Different DRAM Mapping Policies

We evaluated the impact of different DRAM mapping policies and the results are presented in Fig. 9.

Key Observation-①: Our DRMap (Mapping-3) achieves the lowest EDP across different layers of the network, across different DRAM architectures, and across different scheduling schemes. It indicates that the DRMap is the most effective DRAM mapping policy for different possible conditions. According to Table I, DRMap (Mapping-3) orderly prioritizes to map the data to different columns in the same row (leading to row buffer hits in SALP and DDR3), to different banks in the same chip (exploiting bank-level parallelism in SALP and DDR3), to different subarrays in the same bank (exploiting subarray-level parallelism in SALP, but leading to row buffer conflicts in DDR3), and to different rows in the same subarray (leading to row buffer conflicts in SALP and DDR3). *Therefore, the DRMap is proven as a generic DRAM mapping policy that offers the lowest EDP.* Moreover, different DRAM access scheduling schemes can make use of the DRMap, so that the CNN accelerators with different scheduling schemes can optimize their DRAM access latency and energy. DRMap improves the EDP up to 96% in DDR3, 94% in SALP-1, 91% in SALP-2, and 80% in SALP-MASA, as compared to other mapping policies.

Key Observation-②: Mapping-2 and Mapping-5 obtain worse EDPs (across different layers of the network, across different DRAM architectures, and across different scheduling schemes) than rest of the mapping policies. The reason is that, Mapping-2 and Mapping-5 prioritize to map data across different subarrays in the same bank (exploiting subarray-level parallelism in SALP, but leading to row buffer conflicts in DDR3), that incurs higher latency and energy, as compared to row buffer hits and exploiting bank-level parallelism.

Key Observation-③: Mapping-1 and Mapping-3 obtain comparable EDPs. The reason is that, Mapping-1 and Mapping-3 prioritize to map data across different columns in the same row (leading to row buffer hits in SALP and DDR3). The difference comes when Mapping-1 prioritizes to exploit subarray-level parallelism over bank-level parallelism, while Mapping-3 is the opposite. From Fig. 1, it is apparent that exploiting subarray-level parallelism incurs higher latency and energy than exploiting bank-level parallelism.

B. Comparisons of Employing Different DRAM architectures

In general, employing SALP architectures provides EDP improvements as compared to employing DDR3. It is mainly due to latency and energy saving that are offered when exploiting subarray-level parallelism. **Key Observation-④:** For instance, if we consider *adaptive-reuse* scheduling, EDP improvements achieved by employing SALP architectures as compared to DDR3 are:

- **For Mapping-1:** 0.59% (SALP-1), 3.89% (SALP-2), and 1.05% (SALP-MASA).
- **For Mapping-2:** 29.18% (SALP-1), 19.91% (SALP-2), and 81.04% (SALP-MASA).

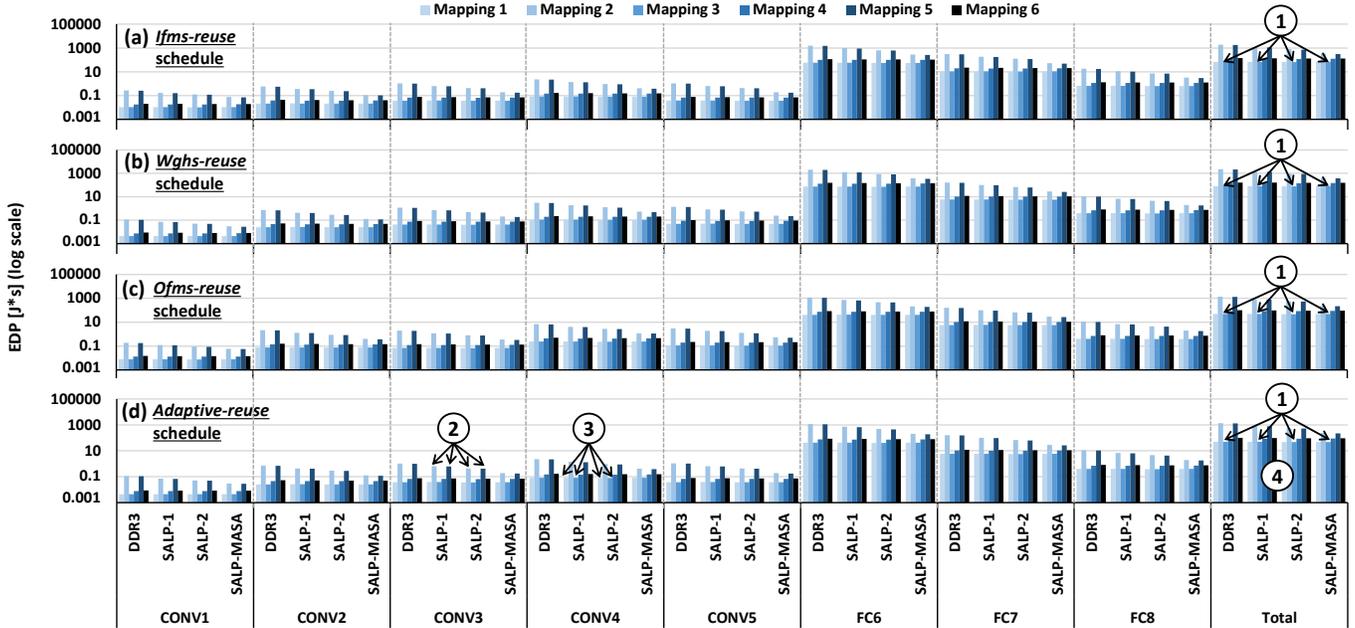


Fig. 9. The EDP in AlexNet for different DRAM mapping policies across different DRAM architectures (DDR3, SALP-1, SALP-2, and SALP-MASA), while considering different scheduling schemes: (a) *ifms-reuse* scheduling, (b) *wghs-reuse* scheduling, (c) *ofms-reuse* scheduling, and (d) adaptive reuse scheduling.

- **For Mapping-3 (DRMap):** 0.6% (SALP-1), 3.87% (SALP-2), and 1.01% (SALP-MASA).
- **For Mapping-4:** 0.71% (SALP-1), 0.54% (SALP-2), and 1.41% (SALP-MASA).
- **For Mapping-5:** 29.67% (SALP-1), 19.79% (SALP-2), and 81.76% (SALP-MASA).
- **For Mapping-6:** 3.15% (SALP-1), 3.39% (SALP-2), and 7.62% (SALP-MASA).

The results show that employing SALP architectures is beneficial for improving energy-efficiency of DRAM accesses, as long as an effective mapping policy like DRMap is employed. The EDP of employing different DRAM architecture would be different, due to the different DRAM access energy and latency. However, since the internal organization of all DRAM architectures is similar (i.e., it is composed of channel, rank, chip, bank, subarray, row, and column as seen from top to bottom perspective), our DRMap can also be employed for all DRAM architectures to achieve the energy-efficient processing of convolutional neural networks in CNN accelerators.

VI. CONCLUSION

In this paper, we present DRMap, a generic DRAM mapping policy that offers the lowest EDP of DRAM accesses for CNN accelerators, as compared to other mapping policies. It is proven through an extensive design space exploration that study the latency and energy of different mapping policies, in different DRAM architectures as well as different layer partitioning and scheduling schemes. We expect that this work could enable further studies on energy-efficient CNN accelerators and help the existing CNN accelerators to optimize their DRAM access latency and energy.

VII. ACKNOWLEDGMENT

Authors acknowledge the scholarship granted by Indonesia Endowment Fund for Education (IEFE/LPDP), Ministry of Finance, Republic of Indonesia.

REFERENCES

[1] Y. LeCun *et al.*, “Deep learning,” *Nature*, vol. 521, no. 7553, 2015.

[2] T. Chen *et al.*, “Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning,” in *ASPLOS*, 2014.

[3] C. Zhang *et al.*, “Optimizing fpga-based accelerator design for deep convolutional neural networks,” in *FPGA*, 2015.

[4] Y. Chen *et al.*, “Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks,” *IEEE JSSCC*, vol. 52, no. 1, pp. 127–138, Jan 2017.

[5] J. Albericio *et al.*, “Cnvlutin: Ineffectual-neuron-free deep neural network computing,” in *ISCA*, 2016.

[6] T. Luo *et al.*, “Dadiannao: A neural network supercomputer,” *IEEE TC*, vol. 66, no. 1, pp. 73–88, Jan 2017.

[7] N. P. Jouppi *et al.*, “In-datacenter performance analysis of a tensor processing unit,” in *ISCA*, 2017.

[8] A. Parashar *et al.*, “Senn: An accelerator for compressed-sparse convolutional neural networks,” in *ISCA*, 2017.

[9] W. Lu *et al.*, “Flexflow: A flexible dataflow accelerator architecture for convolutional neural networks,” in *HPCA*, 2017.

[10] H. Kwon *et al.*, “Maeri: Enabling flexible dataflow mapping over dnn accelerators via reconfigurable interconnects,” in *ASPLOS*, 2018.

[11] M. A. Hanif *et al.*, “MPNA: A massively-parallel neural array accelerator with dataflow optimization for convolutional neural networks,” *CoRR*, vol. abs/1810.12910, 2018.

[12] H. Sharma *et al.*, “Bit fusion: Bit-level dynamically composable architecture for accelerating deep neural network,” in *ISCA*, 2018.

[13] V. Sze *et al.*, “Efficient processing of deep neural networks: A tutorial and survey,” *Proc. of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.

[14] J. Li *et al.*, “Smartshuttle: Optimizing off-chip memory accesses for deep learning accelerators,” in *DATE*, 2018.

[15] A. Stoutchinin *et al.*, “Optimally scheduling CNN convolutions for efficient memory access,” *CoRR*, vol. abs/1902.01492, 2019.

[16] S. Ghose *et al.*, “Demystifying complex workload-dram interactions: An experimental study,” in *SIGMETRICS*, 2019.

[17] Y. Kim *et al.*, “A case for exploiting subarray-level parallelism (salp) in dram,” in *ISCA*, 2012.

[18] —, “Ramulator: A fast and extensible dram simulator,” *IEEE CAL*, vol. 15, no. 1, pp. 45–49, Jan 2016.

[19] S. Ghose *et al.*, “What your dram power models are not telling you: Lessons from a detailed experimental study,” *Proc. ACM MACS*, vol. 2, no. 3, pp. 38:1–38:41, Dec. 2018.

[20] A. Krizhevsky *et al.*, “Imagenet classification with deep convolutional neural networks,” in *NIPS*, 2012.

[21] J. Liu *et al.*, “Raidr: Retention-aware intelligent dram refresh,” in *ISCA*, 2012.

[22] Y. Kim *et al.*, “Thread cluster memory scheduling,” *IEEE Micro*, vol. 31, no. 1, pp. 78–89, 2011.