

PIPEORGAN: Efficient Inter-operation Pipelining with Flexible Spatial Organization and Interconnects

Raveesh Garg¹, Hyoukjun Kwon², Eric Qin³, Yu-Hsin Chen³, Tushar Krishna¹, and Liangzhen Lai⁴

¹Georgia Tech, ²UC Irvine, ³Meta, ⁴Witmem

Abstract—Because of the recent trends in Deep Neural Networks (DNN) models being memory-bound (e.g., large language models and convolutional neural networks with heavy skip connections), inter-operator pipelining for DNN accelerators is emerging as a promising optimization. Inter-operator pipelining reduces costly on-chip global memory and off-chip memory accesses by forwarding the output of a layer as the input of the next layer within the compute array, which is proven to be an effective optimization by previous works.

However, the design space of inter-operator pipelining is huge, and the space is not yet fully explored. In particular, identifying the right depth and granularity of pipelining (or no pipelining at all) is significantly dependent on the layer shapes and data volumes of weights and activations, and these are different even within a domain. For instance, AR/VR applications have 6 orders of magnitude swing in the activation to weight ratios. Another factor affecting the right depth and granularity is the dependencies in the form of skip connections, which increase the activation accesses and vary in reuse distance and density across applications.

Moreover, works divide the substrate into large chunks and map one layer onto each chunk, which requires communicating halfway through or through the global buffer. However, for fine-grained inter-operation pipelining, placing the corresponding consumer of the next layer tile close to the producer tile of the current layer is a better way to exploit fine-grained spatial reuse.

In order to support variable number of layers (ie the right depth) and support multiple spatial organizations of layers (in accordance with the pipelining granularity) on the substrate, we propose PIPEORGAN, a new class of spatial data organization strategy for energy efficient and congestion-free communication between the PEs for various pipeline depth and granularity. PIPEORGAN takes advantage of flexible spatial organization and can allocate layers to PEs based on the granularity of pipelining. We also propose changes to the conventional mesh topology to improve the performance of coarse-grained allocation. PIPEORGAN achieves 1.95x performance improvement over the state-of-the-art pipelined dataflow on XR-bench workloads.

I. INTRODUCTION

Deep Neural Networks (DNN) are gaining popularity due to their use in applications including natural language processing (NLP) [1], [7], computer vision [11], [12], [28] and personalized recommendations [30]. The most compute-intensive Deep Neural Networks operators are the Einstein summation (einsum)-based operators, which refers to dot product-based operations generalized to arbitrary dimensions (e.g., general matrix multiplications, or GEMM). Example DNN operators based on Einstein summation include linear layer, convolution, and batched matrix multiplication. Because the einsum-based operators typically account for 70% of total latency on GPUs [32], many accelerators focusing on matrix

multiplication [5], [13], [16], [24], [32] and design-space exploration tools [15], [17], [21], [22], [31], [41] emerged.

DNN accelerators employ various dataflows and scheduling strategies to map the matrix multiplication spatially over the processing elements (PEs) and temporally. However, optimizing individual matrix multiplication operators does not always translate to optimal execution of the whole application. This is because it misses out on the opportunity to reuse the output feature map in the next operator that certain layers have. Therefore, current works are actively investigating inter-layer pipelining or inter-operator (operator means tensor-operator) pipelining in order to reuse the portions of the intermediate feature maps [2], [3], [8], [35], [37], [40], [44]. Moreover, recent prior works also explore design-space exploration with inter-operator pipelining [3], [37], [43], [44]. Prior work FLAT [18] and TANGRAM [8], respectively claim 1.5x and 2x performance gains from pipelining, over their respective state-of-the-arts. However, the benefits of pipelining heavily depend on two aspects: (1) pipelining depth (i.e., how many layers do we pipeline) and (2) pipelining granularity (i.e., the tile size of the pipelining).

Pipelining Depth Fig. 1 shows the impact of depth on two opposite kinds of layers. The right pipeline depth depends on two factors - A/W ratio (*Activation volume/Weight volume*) and skip connections.

A/W ratio: Activation-heavy layers prefer higher pipelining depth, since, pipelining essentially reuses activations, which are shared between two consecutive layers, while shallow pipelining results in a large overhead of re-fetching big activations. On the other hand, weight-heavy layers do not favor pipelining, and prefer intra-operator reuse, since weights are not shared between the layers, and deeper pipelines, involves those unshared weights being together at a time. In case of models [4], [6], [10], [11], [25], [27], [27], [29], [33], [34], [36], [38], [39], [42] inside XR-bench [23], *Activation volume/Weight volume* ratio roughly spans from 10^{-3} to 10^3 as Fig. 5 shows in Sec. III.

Skip connections: Skip connections make the otherwise activation-heavy layer, even more activation-heavy, since these combine activations from multiple preceding layers. The last layer of the DenseNet [14] (used in RITNet [4] model for eye segmentation) block combines four activations, and op-by-op operator in the middle means, re-fetching all these activations. Multiple XR-bench models have skip connections with different densities and reuse distances as shown in Fig. 6 in Sec. III. Skip connections also diminish the benefits if the pipeline is

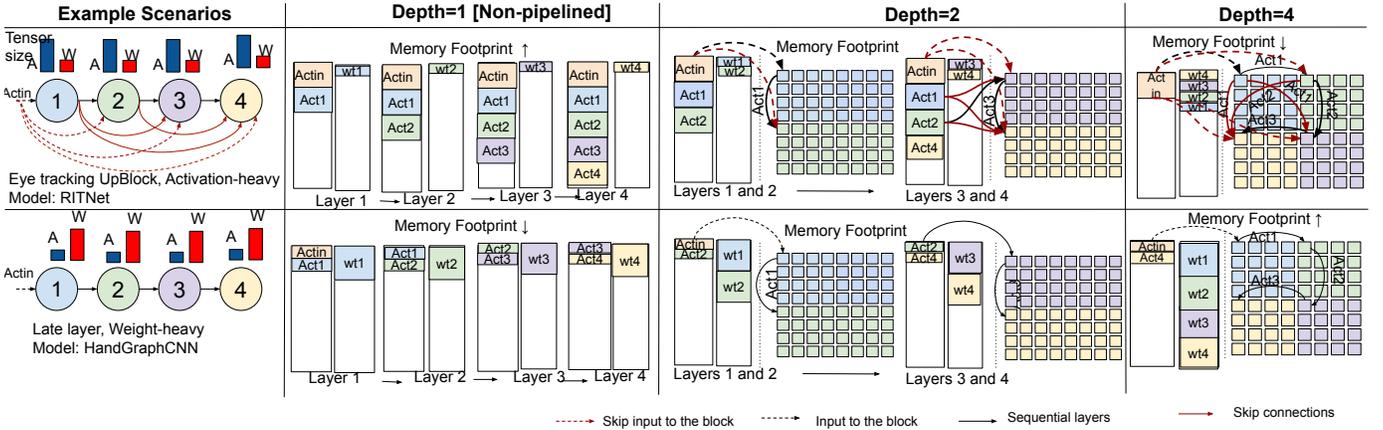


Fig. 1. Impact of pipeline depth for different sets of layers. We only show data movement with one PE from each layer in the figure. The boxes represent memory footprint, and not the space allocated in the buffer. Tensors with larger volumes get reused within PE array with deeper pipelining.

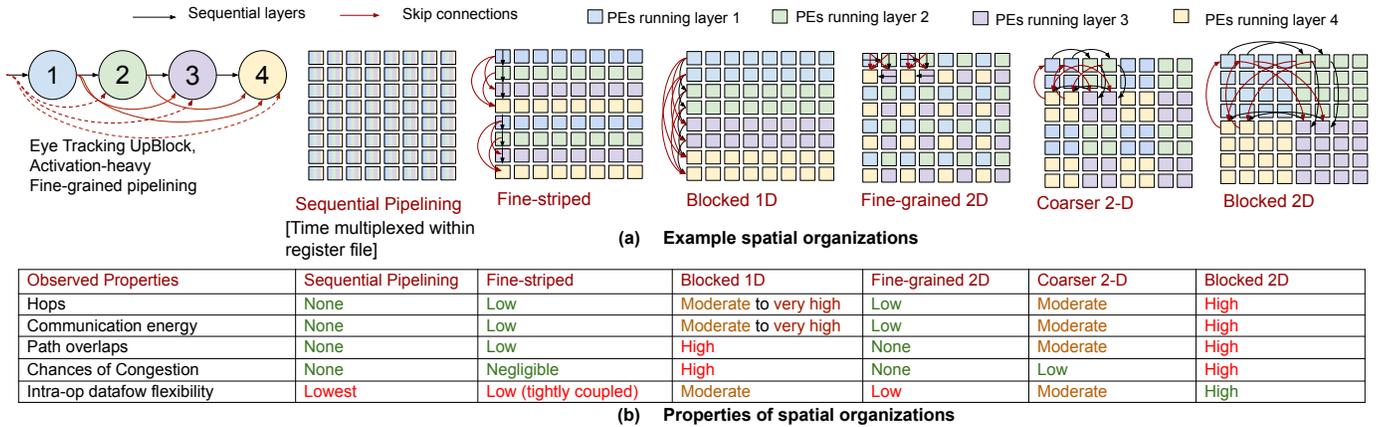


Fig. 2. High-level example showing the impact of spatial organizations on energy and latency. We consider depth = 4, along with the traffic that results after running an activation-heavy UpBlock in RITNet [4], used extensively in eye-tracking. We only show data movement with 2 PEs from each layer. Cycle-level analysis on traffic is shown in Fig. 8-11

not deep enough. As Fig. 1 shows, with pipeline depth of 2, the output of layer 1 eventually needs to be written back because of a downstream dependency, and we only observe the benefit on output of layer 3.

Pipelining Granularity: Pipeline granularity refers to the size/portion of the intermediate tensor consumed by the consumer. In prior works on inter-operator pipelining, one layer is mapped on a block of PEs with each partition running an individual layer. This is inefficient for highly fine-grained pipelining between layers, and this spatial organization can lead to higher hop energy and possibly NoC congestion. Fig. 2 Blocked-2D and Fine-grained-2D show examples of aforementioned spatial arrangements. Interleaving creates a trade-off between locality and flexibility since finest-grained interleaving makes the layers tightly coupled to each other and poses constraints on the tile sizes. On the other hand, coarse grained interleaving, may not exploit enough locality but it provides, flexibility in tiling layers. We show that arrangement of the layers spatially depends on the preferred dataflow of the layers which also depends on the layer shape.

In this work, we propose PIPEORGAN¹, a new class of spatial organization strategies for inter-operator pipelining, and a systematic optimization methodology for that. Fig. 2 shows some examples of different organization strategies for depth=4.

We also characterize the possible traffic patterns that arise from plethora of factors, like depth, granularity, skip connections, load balancing layers with unequal MACs etc. Mesh topology does not achieve the best energy efficiency different kinds of traffic patterns, specially ones which involve skip connections. Moreover it requires us to trade-off on-chip energy for intra-operator flexibility, in cases where fine-grained pipelining is not possible, since coarse-grained allocation requires large hops. Also, spatial organizations on mesh topology trade-off inter-operator and intra-operator reuse specially in cases of large pipelining depths. Flattened butterfly [19] topology allows direct communication between distant PEs, however is an overkill, and increases the link complexity to $O(N \log N)$. Thus we propose changes to mesh topology that reduce congestion and hop count for mesh without adding too many or too long links.

¹Pipelining Organization

Our contributions are as follows:

- We show the importance of considering variable depth and pipelining, given variation across applications and variation within an application itself.
- We quantitatively show that activation/weight ratio is the key metric that affects the pipeline depth and granularity.
- We propose PIPEORGAN, a new class of spatial organization strategies, with different granularities at which multiple layers are arranged, ranging from fine-grained checkerboard and striped to traditional blocked 1-D or 2-D arrangement as Fig. 2 shows. This is the first work (to the best of our knowledge) that proposes finer-grained spatial organization strategies between different layers.
- We characterize the traffic patterns resulting from different spatial organization strategies and different factors like depth, granularity, skip connections, unequal PE allocation between layers etc. We identify the bottlenecks in mesh and propose AMP², a modified mesh for supporting coarse-grained patterns as well.

II. BACKGROUND AND RELATED WORK

A. Individual Operation Dataflows and Mappings

Majority of the expensive operations in DNNs involve Convolutions and matrix multiplication operations. These operations can be represented using einsums as follows:

$$O_{m,n} = \sum_k A_{m,k} \times B_{k,n} \quad (1)$$

$$O_{n,h,w,k} = \sum_{c,r,s} I_{n,h+r,w+s,c} \times W_{r,s,c,k} \quad (2)$$

Equation 1 represents matrix multiplication with M,K and K,N as the dimensions of the input matrices and M,N as the dimensions of the output matrix. Equation 2 represents Convolution layer, with H and W as the feature map height and width respectively, R and S as the filter height and width respectively, N as the batch size and C and K as the number of input and output filters respectively.

These operations can be represented as loop nests as shown below.

```

1. for n in range(N):
2.   for h in range(H):
3.     for w in range(W):
4.       for k in range(K):
5.         for c in range(C):
6.           for r in range(R):
7.             for s in range(S):
8.               O(n,h,w,k) += I(n,h+r,w+s,c) * W(r,s,c,k)

```

Dataflow is used to describe the loop transformations for staging the operations in space and time on a spatial accelerator. Here, the order of the temporal loops in the above example, is NHWKCRS where N is the outermost loop. The dataflow determines the compute utilization and the locality of the tensors in the memory hierarchy, thus choice of the right dataflow is crucial for performance and energy efficiency. In

this work, term "dataflow" is to refer to hardware agnostic loop transformations.

One example of a loop transformation is shown below-

```

1. for n in range(N):
2.   for h in range(H):
3.     for w in range(W):
4.       for k in range(K):
5.         for c1 in range(C/C0):
6.           for r in range(R):
7.             for s in range(S):
8.               parfor c0 in range(C0):
9.                 c=c1*C0+c0
10.                O(n,h,w,k) += I(n,h+r,w+s,c) * W(r,s,c,k)

```

Moreover, C is spatially parallelized across multiple processing elements or PEs. C0 is based on the PE array size. Hardware dependent loop transformations determine the complete *mapping*. We discuss the space for inter-operation dataflows and mappings in Sec. III

B. Inter-Operation Pipelining

Layer-by-layer computation leads to the whole feature map being written into the memory which leads to high occupancy of the data in the memory hierarchy and also causes excess roundtrip memory accesses to fetch the data for the next layer. To overcome this, prior works have proposed inter-operation pipelining, also known as inter-layer pipelining or layer fusion or inter-operation pipelining where a portion of the feature map is produced and used by the next layer, decreasing the occupancy of the data inside the memory hierarchy and consequently reducing the overall roundtrip memory accesses. Since the next layer consumes only a portion that the previous layer produces, the dataflows of both the layers must be staged to ensure that, thus leading to an interdependence of dataflows [9].

Fig. 3 shows the producer-consumer relationship in inter-operation pipelining. The producer produces a portion of the intermediate data in one timestep which is consumed within the next timestep with producer producing the next piece of data in parallel. Once the data has been consumed, that data is no longer needed. The first stage where the consumer has not started consuming is referred to as "init" while the rest of the stages are collectively referred to as "steady-state". We refer to the duration of the timestep as "interval". The portion of the intermediate data produced/consumed in a timestep is referred to as "granularity" here.

As we can see in Fig. 3 pipelining can be expressed as a waterfall diagram with vertical axis showing different operations and horizontal axis showing time. With variable granularities and load imbalance, we compute producer side delay as the previous interval delay normalized by the ratio of the number of operations in the current interval and the previous interval. The interval delay is the maximum of the producer side and the consumer side delays. The overall latency is the summation of all the interval delays once (which accounts for all the init delay) and the steady state delay of the last operation. These subtleties are captured in the equations in Fig. 3.

²Augmented Mesh for Pipelining

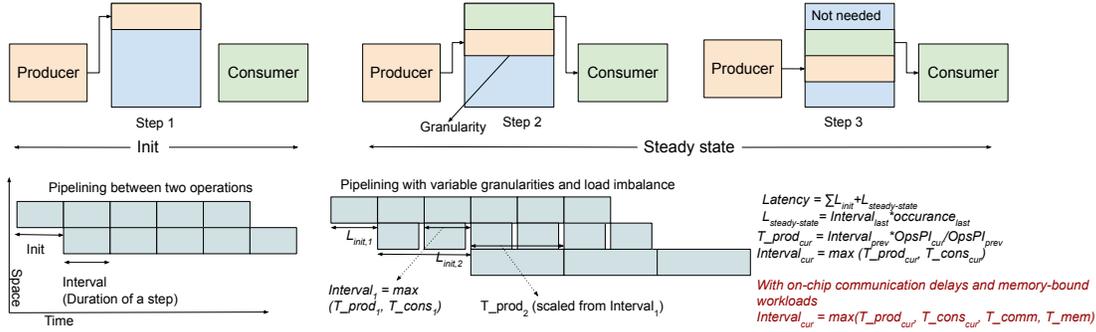


Fig. 3. Inter-operation pipelining between producer and consumer and latency equation for pipelines with arbitrary depth. For fine-grained pipelining, where tiles are small enough to fit inside the local memory of PEs, the data can also be moved through the NoC as opposed to accessing the global buffer.

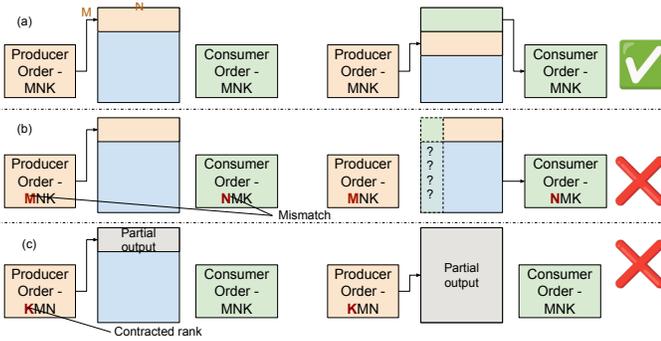


Fig. 4. Conditions for inter-operation pipelining. (a) Conditions met (b) Violation of the same outermost loop (c) Contracted rank of the producer in the outermost loop

Conditions (shown in Fig. 4) for making inter-operation pipelining between producer and consumer possible are-

- For a tensor shared between the producer and the consumer, atleast the outermost loop should be the same. This is needed to divide the producer and consumer into stages.
- The contracted rank should not be the outermost rank for the producer, since complete sums are needed earlier for consumption. Similarly, the unshared rank of the consumer must not be in the outermost loop since it would nullify the benefits of pipelining by having to use the complete intermediate tensor in the inner loops multiple times.

C. Related Work

Various prior works like Atomic Dataflow [43], Stream [37], HyGCN [40], OMEGA [9], TANGRAM [8] and FusedCNNs [2] work on inter-operation pipelining between two tensors. Some mapping frameworks like CoSA [15] and GAMMA [17] focus on layer-by-layer computation but explore multiple dataflows on a flexible accelerator.

Most of the prior works on pipelining do not consider both variable pipeline depth and variable granularity together. We show the prior works in terms of their support for fine grained pipelining and depth awareness in Table I. Please note, that each work presents a different set of contributions, and our work primarily focuses on interconnect and spatial organization. This is the first work, within DNN accelerators that also explores finer-grained spatial organization of layers on PEs (Fig. 2).

TABLE I

RELATED WORK TABLE SHOWING PRIOR WORKS ON INTER-LAYER PIPELINING. REGARDING ORANGE CHECKMARK, - 1) TILEFLOW CANONICALLY SUPPORTS VARIABLE DEPTH BUT FIXES IT WITHIN AN APPLICATION, AND 2) HYGCN HAS TWO GRANULARITIES.

| Prior Work | Variable Depth | Variable Granularity | Contribution |
|----------------------|----------------|----------------------|----------------------------------------|
| Atomic Dataflow [43] | ✓ | ✗ | Inter-operation mapper |
| Stream [37] | ✗ | ✓ | Heterogeneous Dataflow Framework |
| HyGCN [40] | ✗ | ✓ | Accelerator |
| EnGN [26] | ✗ | ✓ | Accelerator |
| OMEGA [9] | ✗ | ✓ | Cost model |
| TANGRAM [8] | ✓ | ✗ | New dataflow |
| FusedCNN [2] | ✗ | ✗ | Accelerator |
| SET [3] | ✓ | ✗ | Mapper |
| TileFlow [44] | ✓ | ✓ | Mapper |
| SIMBA [35] | ✓ | ✗ | Mapper |
| PIPEORGAN | ✓ | ✓ | Spatial organization, NoC and Dataflow |

D. Motivation toward Flexible Pipelining: Heterogeneity in AR/VR models

As XR-bench [23] shows, DNN models have high heterogeneity in terms of kinds of layers (convolution, depthwise convolution, GEMM, RPN, ROIAlign), in terms of DAG (skip connections with various reuse distances and densities in a block), layer dimensions etc. Fig. 5 shows the activation/weight ratios of layers within the XR-bench models. These ratios range six orders of magnitude from activation-dominant to weight-dominant layers. Moreover, unlike traditional models like ResNet [12], the location of activation heavy and weight heavy layers inside the model isn't predictable. Fig. 6 shows the skip connections within five X-R bench models. These skip connections vary in reuse distance and density. For example, RITNet has dense skip connections of multiple reuse distances and midas one skip connection per block with varying reuse distance. Therefore, high heterogeneity in layers requires more attention to finding the right depth and granularity of pipelining, and a flexible NoC that can efficiently execute each scenario.

III. PIPELINING DESIGN-SPACE

The pipelining dataflow space consists of four aspects - depth, intra-operator dataflow, granularity and spatial organization.

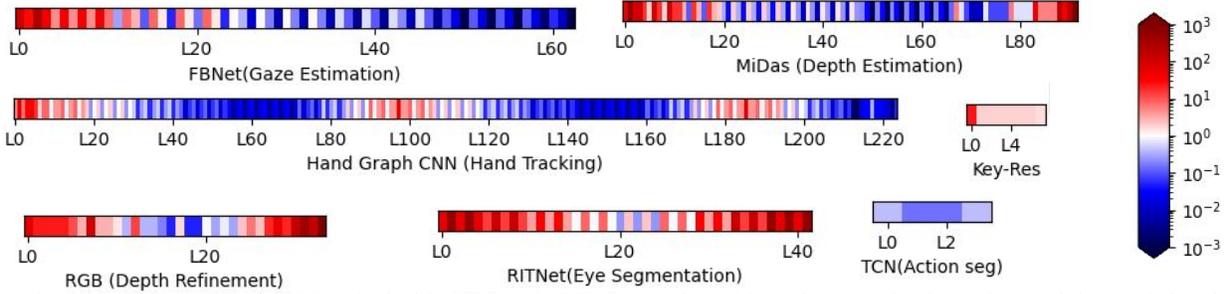


Fig. 5. Activation/weight ratios of CNN tasks inside XRbench [23]. Layers in red have larger activations, than weights, and the blue ones have larger weights. This excludes the skip connection activation traffic.

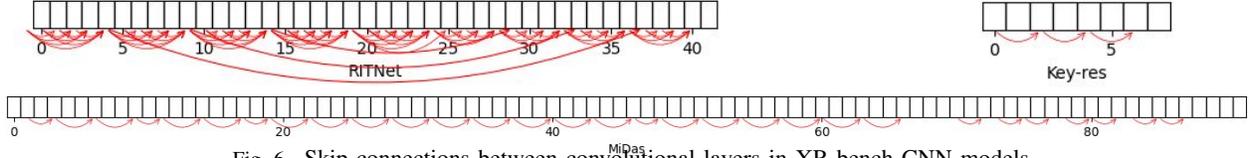


Fig. 6. Skip connections between convolutional layers in XR-bench CNN models

Depth. Refers to the number of layers being pipelined at a time. The model is divided into "pipeline segments" of various depths.

Intra-operator Dataflow. Refers to the dataflow of the individual operator.

Granularity. Refers to the size/portion of the intermediate tensor consumed by the consumer. For example, granularity of pipelining in Fig. 3 is one row.

Spatial Organization. Patterns in which different layers can be arranged on the PEs. Fig. 2 shows examples of different spatial organization strategies.

DNN models can be represented as a DAG of layers. We divide the model into segments of variable depths, which is dependent on shapes of multiple layers and DAG dependencies. We choose the intra-operator dataflow, which depends on the shape of that layer. Based on the intra-operator dataflow, we determine the granularity at which we can pipeline a pair of producer and consumer, within the confines of the depth. Then, we determine the right spatial organization strategy, which depends on the depth and the granularity of pipelining.

The first three namely, depth, granularity and intra-operator dataflow are agnostic of specific hardware details like NoC topology etc. Spatial organization on the other hand, comprises of mapping as it depends upon specific hardware details like NoC topology, PE array dimensions etc.

A. Factors Affecting Pipelining Depth

The choice of pipelining depth impacts inter-operation vs intra-operation reuse. Right pipeline depth depends on two main factors - A/W ratios and skip connections

A/W ratios: Deeper pipeline implies reusing intermediate activations, but at the same time increasing the memory footprint of weights from multiple layers as Fig. 1 also shows. Therefore, for pipeline depth of D , weight memory footprint is $\sum_{i=l}^{l+D} W_i$ where l is the first layer of that segment. Activation memory footprint on the other hand is reduced significantly

in comparison, its $A_l + A_{l+D} + \sum_{i=l+1}^{l+D-1} Granularity_i$, where granularity is the portion of the intermediate matrix between the producer and the adjacent consumer. If certain pairs of layers are fine-grained pipelined, the granularity component can be taken care of by PE-to-PE communication leaving the footprint $A_l + A_{l+D}$. Thus larger depth implies, more activations in the middle can simply be skipped from calculation. Although, in case of weights, larger depth implies incurring the footprint of D layers throughout the execution, at all times, reducing the available tile size for weights. Thus, layers with large activation/weight ratios benefit from deep pipelining while the ones with small ratios benefit from shallow or no pipelining.

Skip connections: Skip connections were specific to ResNet [12] but have become a norm, and are used in various DNN models [4], [6], [33], [38], [39]. Skip connections can vary in density as well as reuse distance. The activation footprint in presence of the skip connections changes to $A_l + A_{l+D} + \sum A_i$ such that $i \notin (l, l+D)$ and there exists skip connection between i and $j \in (l, l+D)$. This includes both outgoing and incoming skip connections. Thus, for deeper pipelining, the instances of skip connections from outside $(l, l+D)$ are likely to be lower. Therefore, presence of skip connections skews towards deeper pipelining, to absorb those connections within the depth of pipelining.

B. Factors Affecting Intra-operator Dataflow

The performance of intra-operator dataflow depends primarily on the tensor dimensions [5], [21], since the purpose is to get reuse. For the off-chip dataflow for operators with extreme A/W ratios, larger tensors should be stationary and the smaller tensors should be streaming, as they can stream from on-chip.

C. Factors affecting Pipelining Granularity

Granularity is determined by intra-operator dataflow. For a finer-grained pipelining, the tensor should be consumed in the same order in which it is being produced. For example,

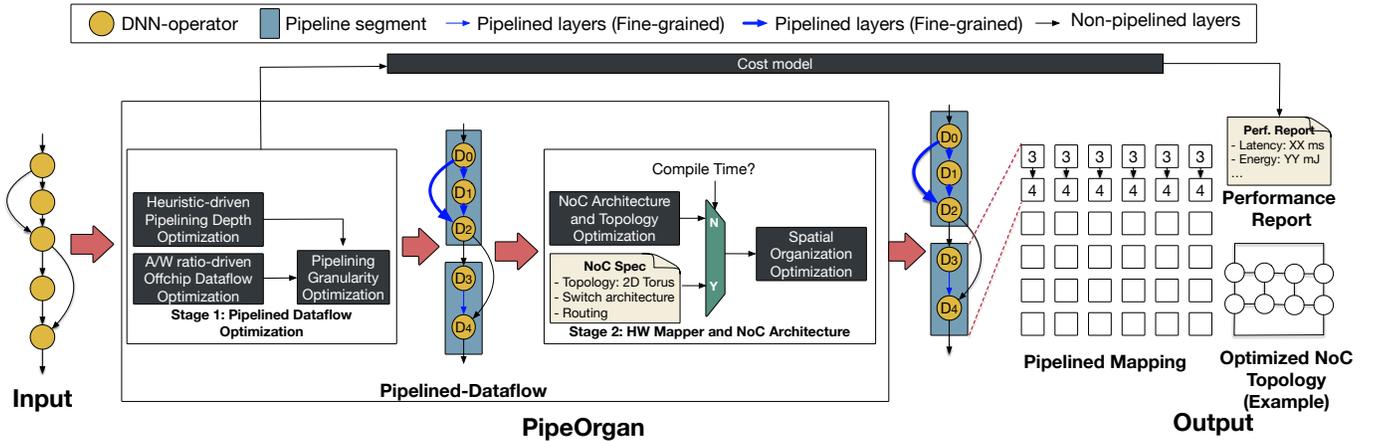


Fig. 7. Entire flow, including mapping heuristics, PIPEORGAN and AMP.

in convolutions, the finest grained pipelining is between the pair NHWKCRS-NHWCKRS, given that the data is being consumed exactly as produced. However, the pair NHWKCRS and NHKWCRS has a coarser granularity since layers can only be staged by NH. Similarly for a GEMM, for example, MNK-MKN is the finest grained pipelining possible while MNK-MNK is a coarser grained pipelining. Depending on the loop order, we can calculate the portion of the intermediate tensor and hence the granularity.

Tile sizes can have an impact of granularity particularly if they are unequal. Lets consider the pair NHWKCRS-NHWCKRS. Just based on the loop order, the pipeline granularity is that of a filter. However, if, for example tile size of H is different, the producer and consumer will only be synchronize only when $LCM(Tile_H_{producer}, Tile_H_{consumer})$ rows have been computed. Hence regardless of the loop order, difference in tile sizes can affect the granularity of pipelining.

D. Factors affecting Spatial Organization

The right spatial organization depends on the pipeline depth, which determines how many different layers are allocated at a time, and on pipeline granularity, which determines whether the inter-layer communication should fine-grained or coarse-grained. Prior works divide the PE array into blocked chunks based on the depth and assigned a layer to the chunk, however this is inefficient for fine-granularity pipelining. Therefore we propose a new class of mappings where the layers can be organized flexibly on the PE array.

In Sec. IV, we focus on spatial organization in more detail.

IV. PIPEORGAN: FLEXIBLE SPATIAL ORGANIZATION

Fig. 7 shows the whole flow of PIPEORGAN. The first stage involves pipelined dataflow optimization, and the second stage involves hardware-aware mapping and NoC architecture.

Stage 1: Pipelined Dataflow Optimization This stage involves using heuristics to choose the right depth and intra-operator dataflow, which then helps determine granularity of pipelining.

Stage 2: HW mapping and NoC architecture: In the second stage on the depth and granularity, we determine

the spatial organization strategy to determine the complete mapping. At design-time we study the traffic patterns of various pipelined-dataflows on various spatial organization strategies, identify bottlenecks in the mesh and propose AMP. At compile-time, we determine the spatial organization strategy based for AMP based on pipelined-dataflow. We discuss the details of AMP in Sec. IV-D.

A. Pipelined Dataflow

First stage of PIPEORGAN takes in an input DAG and uses heuristics to (a) partition the whole model into segments of flexible depth, (b) determine intra-operation dataflows and (c)

Determining Depth: In this work, we determine depth of a segment (starting at layer l) by comparing the memory footprints $A_l + A_{l+D}$ with $\sum_{i=l}^{l+D} W_i$ (see Sec. III-A), increasing the value of D . We stop adding more depth, the moment $\sum_{i=l}^{l+D} W_i$ is greater. In case of skip connections, we also add additional activations due to skip connections, thus skip connections skew the decision towards deeper pipeline. We also cut the depth if we encounter a complex layer like ROIAlign. The depth is also limited by the size of the substrate. The maximum depth we consider is $\sqrt{\text{numPEs}}$.

Determining Intra-operation Dataflows (Loop order): Intra-operation dataflows greatly influence pipelining and even the ability to pipeline as Sec. III-C shows. Ideal intra-operation dataflows can depend on layer shapes as prior works [5], [21] have shown. For the scope of the paper, we simply choose a few dataflows depending on the ratio of the weight and activation volumes. In case of larger weights, we use weight stationary dataflow, where ranks from weights form the outermost loop, to get more reuse on weights. This dataflow is not friendly to pipelining. While for the activation-heavy layers, we choose the activation stationary dataflow. Depending on how large activation is compared to the weight, we decide whether to make the dataflow completely activation stationary (for example, NHWKCRS) or we allow some reuse on weights (for example, NHKWCRS). We validate our heuristic on XR-bench usage scenarios. We are able to achieve the best possible arithmetic

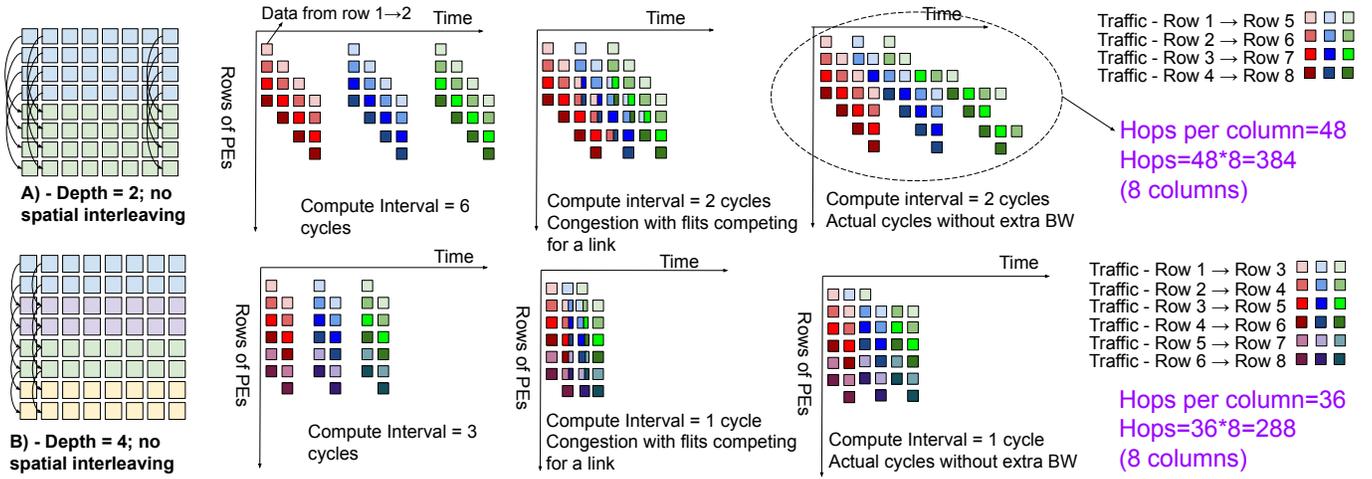


Fig. 8. Fine-grained inter-operation pipelining with coarse-grained (blocked) spatial allocation for depth=2 and depth=4.

intensity³ in case of 99.94% of the layers with on-chip buffer size of 512KB and 97.2% of the layers with on-chip buffer size of 256KB. Note that this only determines the order of dimensions in the outer loops, not including spatial parallelism and spatial tiling.

Determining Finest Possible Granularity Based on Loop Order We determine granularity from intra-operation dataflows using Alg. 1 model in this section. The algorithm compares

Algorithm 1: Determination of granularity from intra-operation dataflows (within a pipeline segment)

```

1 Granularity = Output size
2 pipnest=0
3 %Within each pipeline segment
4 for layer ∈ CUR to CUR+DEPTH-1 do
5   for loop ∈ 0 to NUM_UNCONTRACTED_LOOPS do
6     if (loop==0 || tilesz[layer][loop-1]==tilesz[layer+1][loop-1])
7       && (looppair(layer,layer+1)==(N,N)||(H,H)||(W,W)||(K,C))
8         then
9           Granularity = Granularity/Dimension[layer][loop] *
              lcm(tilesz[layer][loop],tilesz[layer+1][loop])
              pipnest += 1

```

the loop pairs to determine if they can be fused and does the same till the loop nests are fusible. However it stops if there is a mismatch in the tile size as discussed in Sec. III-C. This step determines the finest possible granularity, but the granularity can also change depending on the parallelization strategy determined in the spatial organization.

B. HW Mapping: Determining Spatial Organization Strategy

To determine spatial organization, we use the depth to ascertain the number of layers to organize spatially. We allocate number of PEs for each layer based on the ratio of MACs.

We determine the arrangement of those layers based on granularity. We compare the total register file (RF) size with the granularity of pipelining. If $RF_{total} < Granularity$, the data is moved through the Global Buffer. In mappings with

coarse granularity of pipelining, the intermediate data is moved through the Global Buffer (GB). This is always done in a blocked organization.

On the contrary if the total register file size of the producer is larger than the granularity, the spatial organization is decided based on how fine the granularity is, relative to the PE register file. For example, the finest possible granularity can be executed on a checkerboard organization or via sequential pipelining. If the granularity is almost at par with total producer register file size, then pipelining can be done in a blocked organization. The number of PEs involved on the producer side is determined by $Granularity/RF_{per_PE}$.

The key idea is that once the granularity is ascertained, we allow each pipeline interval (defined in Fig. 3) to be mapped flexibly to allow for optimal intra-operation reuse. Fine-grained spatial organization for coarse-grained pipelining constrains the parallelization and tiling strategies that the individual layers use. For example, in case of sequential pipeline fine-grained checkerboard allocation and sequential pipelining, only the part of the intermediate tensor that is produced, would be consumed. For example, if the activations parallelize K and W dimension, they also have to be consumed in the same way, with fine-grained spatial allocation.

Hence coarser-grained pipelining also should use coarser-grained spatial organization. Additionally, 1-D vs 2-D organization is decided based on the depth and on the reuse within the pipeline stage (whether its more along one dimension, or along multiple dimensions). Once the spatial organization strategy is decided, PEs could be allocated to the layers in ratios that ensure load balancing and maximum utilization. Parallelization strategy is also decided at this step, which could potentially increase the granularity from stage 1, but this change does impact the spatial organization decision.

C. HW Mapping: Design-time Traffic Analysis

Fig. 8 shows the traffic patterns generated by fine grained pipelining with blocked 1D spatial organization. Here, we compare the total hop time against the interval time. The

³Best-case arithmetic intensity is obtained by considering only cold misses.

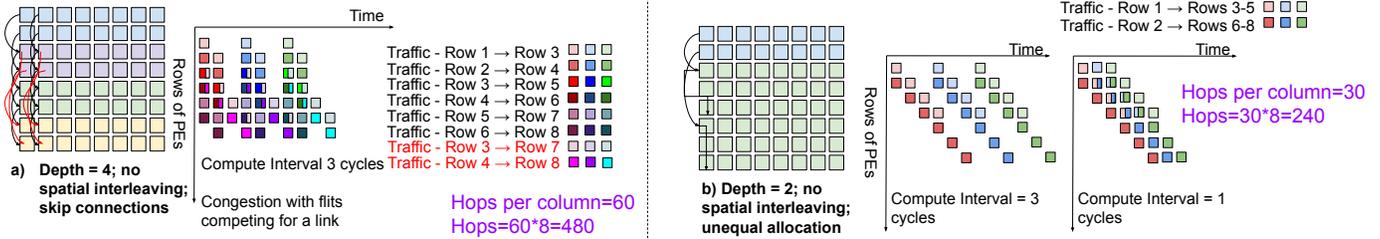


Fig. 9. Effect of (a) skip connection (b) unequal PE allocation due to load balancing on traffic pattern

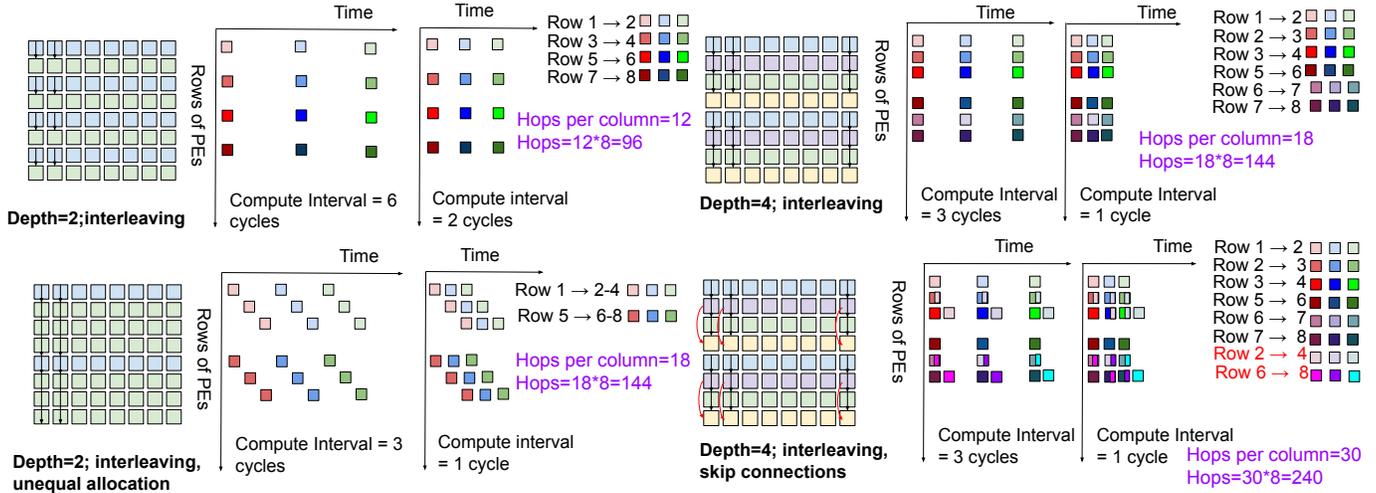
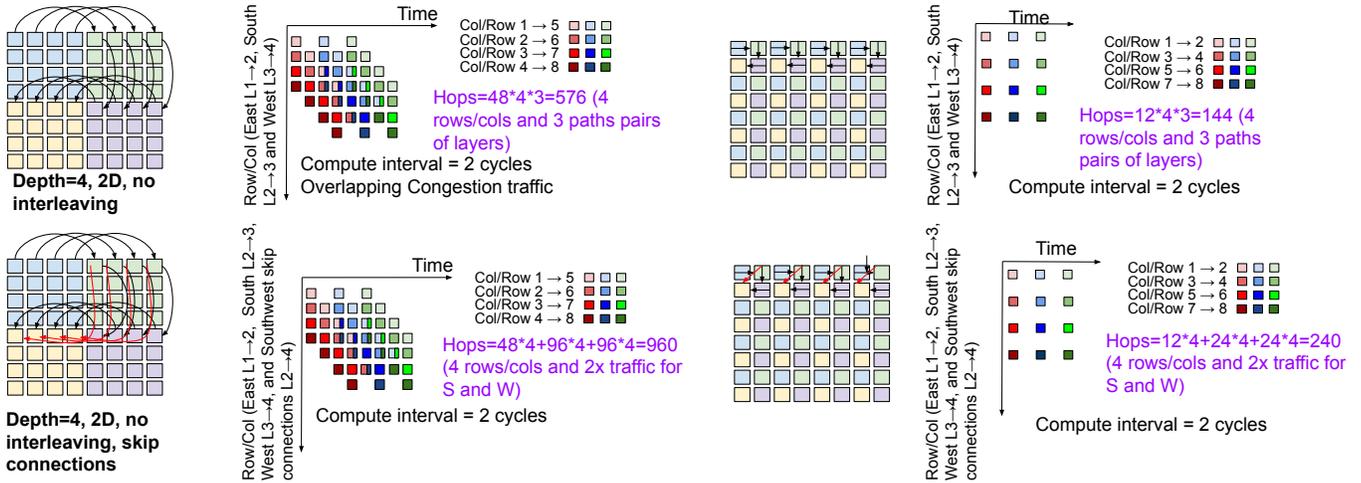


Fig. 10. Congestion-free traffic with 1-D interleaving (fine-striped configuration). Significant hop reduction is observed corresponding to counterparts from Fig. 8 to 9.



For south and west parts, the volume of traffic increases 2x (leading to more congestion) because of the skip connection from L2→4.

Fig. 11. 2-D spatial organization with and without interleaving and with and without skip connections. Pipelining with 2-D allocation can be broken down into multiple 1-D paths. In this example, eastern path spans communication from layer 1 to 2, southern path - layer 2 to layer 3 and western path - layer 3 to layer 4. Skip connection requires traversal along multiple directions.

compute interval time stems from the temporal reduction that needs to be done within the PEs to produce an output that can be consumed. If this time is greater than the hop count, congestion does not happen. However, if this time is less, it leads to congestion. This is due to the contention caused by the new traffic getting generated at a faster rate. On resolving this congestion, we find that the latency is limited by the hop count

rather than the compute interval. This traffic pattern also has a high overall hop count. Moreover, coarse-grained pipelined dataflows benefit from intra-operation reuse.

Fig. 9a) shows additional congestion that is caused by skip connections in residual blocks of ResNet. This would be even more detrimental, in case of RITNet where there is a skip connection to each later layer inside the block, with highly

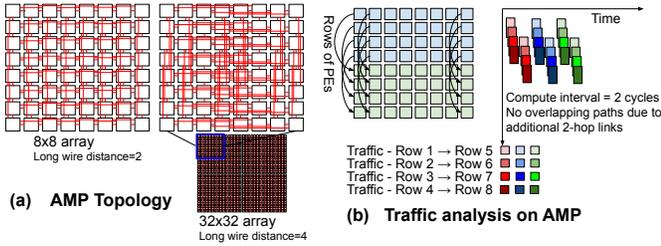


Fig. 12. (a)AMP topology (b) Traffic analysis on depth=2, no spatial interleaving with AMP topology.

activation-heavy layers.

Fig. 9b) shows, a case where PEs for the producer and consumer are unequally allocated. The traffic hotspot is at the boundary of the producer and the consumer, so it is hard to pinpoint the exact row of PEs where it happens. A common example of this is ResNet, with 1x1 and 3x3 filter sizes and varying number of output channels inside a residual block.

Finally, Fig. 10 shows, how 1-D interleaving can avoid congestion by co-locating the producer and the consumer tiles, which is beneficial for fine-grained exchange between the layers. It also shows the traffic benefits in case of skip connections and unequal allocation. Significant reduction in congestion and hop count is observed compared to the blocked organization counterparts in Fig. 8-9.

Fig. 11 compares blocked and fine-grained spatial organization for 2-D allocation, with depth=4. In 2-D allocation, the traffic can be broken down into multiple sets of paths involving 1-D communication. In this example, communication from layer 1 to 2 is along east in the top half, and four rows are involved. Communication from layer 2 to 3 is along the south in the right half and four columns are involved. Similarly communication from layer 3 to 4 is along west in the bottom half and four rows are involved. Skip connections involve traversal along both southern and western paths doubling the traffic for those.

Notice that coarse grained pipelining incurs less congestion even with blocked organization, however, it incurs high hop count. Reducing the hop count would require us to organize the coarse grained pipelined layers in a tightly coupled manner, which reduces intra-op mapping flexibility and efficiency. This motivates us to propose AMP to account for the traffic patterns, in order to not trade-off intra-operation flexibility for on-chip energy. We summarize the bottlenecks caused by various scenarios on mesh in Table II.

TABLE II
SUMMARY OF MESH BOTTLENECKS OBSERVED IN FIG. 8-11.

| Cause | Effect | Prevalent in |
|----------------------------------|-------------------|--------------------|
| Many long overlapping paths | High Congestion | Blocked 1D and 2D |
| Many long overlapping paths | High hop energy | Blocked 1D and 2D |
| Extra BW for skip connections | High congestion | All organizations |
| Extra hops with skip connections | High hop energy | All configurations |
| Routing in multiple directions | Higher hop energy | 2D organizations |

D. AMP Topology

Table II shows the bottlenecks with the mesh topology. These bottlenecks are caused by overlapping paths that could go all the way through a row/column in case of coarse-grained spatial organization (as Sec. IV discusses, fine-grained spatial organization can resolve this, but is not always possible). This bottleneck is made worse by skip connections, as they introduce extra congestion. Flattened butterfly [18] adds links to increase the bandwidth but is an overkill and the number of links increase in complexity from $O(N)$ to $O(N \log N)$. SMART [20] adds single-cycle multi-hop support to mesh, and can help in reducing energy, however, it does not resolve congestion by increasing the available bandwidth.

Fig. 12a shows the AMP topology for 8×8 array. We modify the conventional mesh topology by adding wires of length $\text{Round}(\sqrt{\frac{\text{Rows}}{2}})^4$, in each PE, from each direction. This limits the length of the wires and it scales by $O(\sqrt{N})$ (the wire length spans 4 PEs for a 32×32 PE array and 8 PEs for a 64×64 PE array). By the virtue of having long enough wires itself, the congestion decreases because of less overlapping paths, without links going all the way unlike torus. AMP increases the number of links compared to mesh by under 2x. This also reduces the hop count. Fig. 12b shows the resulting reduction in congestion and hop count.

V. EXPERIMENTAL METHODOLOGY

A. Evaluation Framework

We develop an in-house simulation framework for evaluating the spatial organizations.

PIPEORGAN: The framework models PIPEORGAN stage 1 using the heuristics described in detail in Sec. IV-A. Based on the pipeline depth ie, number of different layers, and the pipeline granularity, the framework determines the right spatial organization strategy at compile time as Sec. IV-B explains, and evaluate performance and energy. At design-time, the framework evaluates the performance and energy of different spatial organization strategies for different stage 1 outputs, for traffic analysis and NoC design.

Performance modeling: The framework models the communication energy and latency cost using in-house NoC simulator to model traffic patterns, topology and routing to compute the hops and estimate the congestion. The NoC simulation automates the NoC and traffic analysis visually shown in Fig. 8-11. We obtain the producer and consumer compute interval latency based on Fig. 3. We also factor in the additional stalls due to limited on-chip memory and limited memory bandwidth. Using compute, communication and memory components, we are able to obtain inter-layer pipelining, based on Fig. 3.

B. Workloads and Datasets

We primarily evaluate PIPEORGAN and AMP on XR-bench [23] CNN tasks, capturing different tasks related to eye tracking, hand tracking, keyword spotting, world locking,

⁴We choose this as its the geomean of single hop and $\frac{\text{Rows}}{2}$ hop case in Fig. 8.

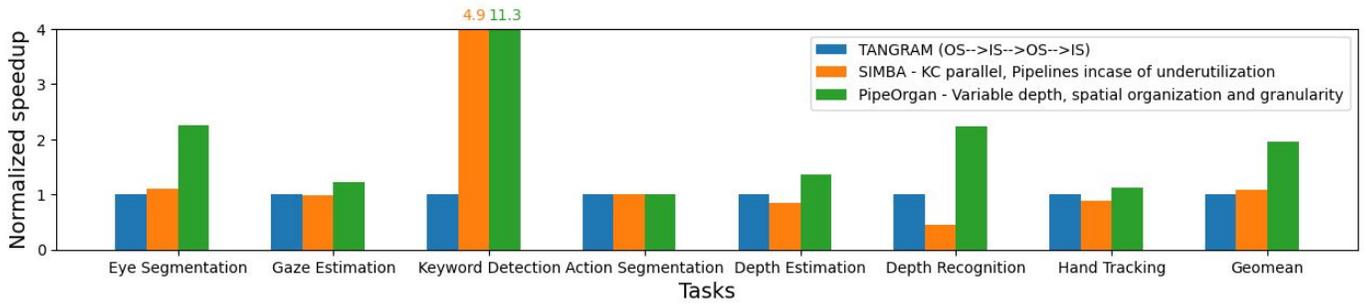


Fig. 13. End-to-end performance benefits for each task in XR-bench [23]. Results are normalized to TANGRAM-like dataflow. Higher is better.

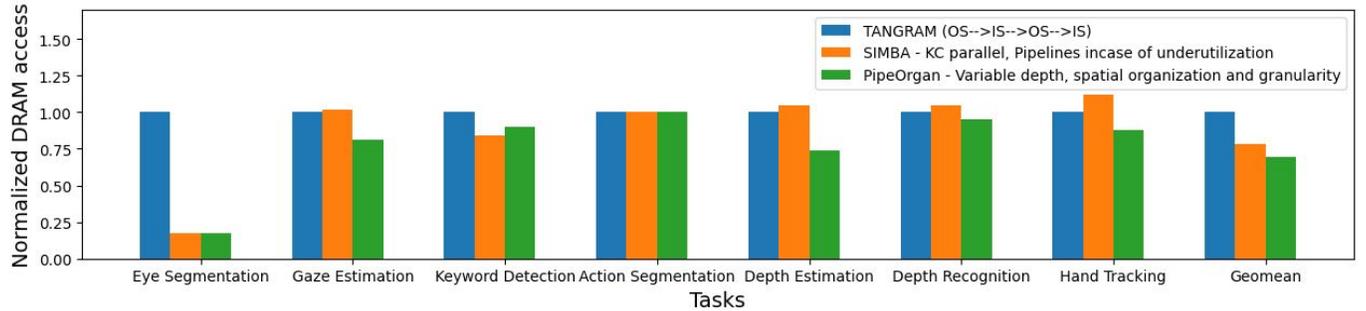


Fig. 14. End-to-end normalized DRAM accesses for each task in XR-bench [23]. Results are normalized to TANGRAM-like dataflow. Lower is better.

object detection, action segmentation etc. As Fig. 5-6 show, these workloads have wide variety of layer shapes and different DAG structures due to skip connections.

C. Baselines

We compare the end-to-end performance of PIPEORGAN on AMP against the dataflows used by TANGRAM [8] and SIMBA [35], for each XR-bench task [23]. TANGRAM-like dataflow alternates between output stationary and input stationary, thus uses fine-grained pipelining with depth=2. SIMBA like dataflow parallelizes input and output channels and does pipelining only when these two dimensions cannot utilize the substrate.

We also show the depth and granularity as an output of stage 1 for each of the XR-bench [23] task in Fig. 16 and 17.

D. Architecture Parameters

We consider the architecture parameters shown in Table III.

TABLE III
CONFIGURATIONS FOR WORKLOADS AND ARCHITECTURE

| Parameter | Value |
|----------------------------------|-------------|
| Bytes per word/element | 1B (8 bits) |
| PE array size | 32×32 |
| PE dot product size ⁵ | 8 |
| SRAM capacity | 1MB |
| Memory bandwidth | 256 GB/s |

VI. RESULTS

A. Performance

Fig. 13 shows the end-to-end performance benefits of PIPEORGAN over SIMBA-like [35] and TANGRAM-like [8]

dataflows for each task. SIMBA-like dataflow pipelines if one layer is unable to utilize the substrate. SIMBA experiences latency in cases where parallelizing input channels and output channels is not sufficient. Other major source of latency is load imbalance, specially in cases with different filter sizes. TANGRAM-like dataflow uses fine-grained pipelining alternating between output stationary and input stationary. However, blocked spatial allocation leads to congestion, hence deteriorating the performance, as it becomes on-chip NoC bound. KD-resnet in particular, is the most affected in case of TANGRAM-like, since the compute interval duration is 1-cycle, and blocked organization is too coarse. PIPEORGAN uses fine-grained spatial organization to avoid this issue. This is similar to the behavior observed in Fig. 8. Action segmentation and hand tracking are mostly weight heavy, with large channels, and therefore do not favor pipelining. Gaze estimation and depth estimation do better with deeper pipelining in the activation heavy regions, given that DWCONV layers are memory bound. PIPEORGAN exploits reuse through flexible depth and spatial organization. All in all, PIPEORGAN gains geomean 1.95x in performance (this includes all the layers and not just the layers that benefit from pipelining).

B. Normalized DRAM accesses

Fig. 14 shows end-to-end normalized DRAM accesses, which are normalized to TANGRAM-like. High DRAM access reduction was achieved on eye segmentation due to flexible depth which absorbs the dense skip connections. Similarly, gaze estimation and depth estimation have memory-bound layers, where PIPEORGAN is able to reduce memory accesses in earlier segments. All in all, DRAM accesses were reduced by 31%

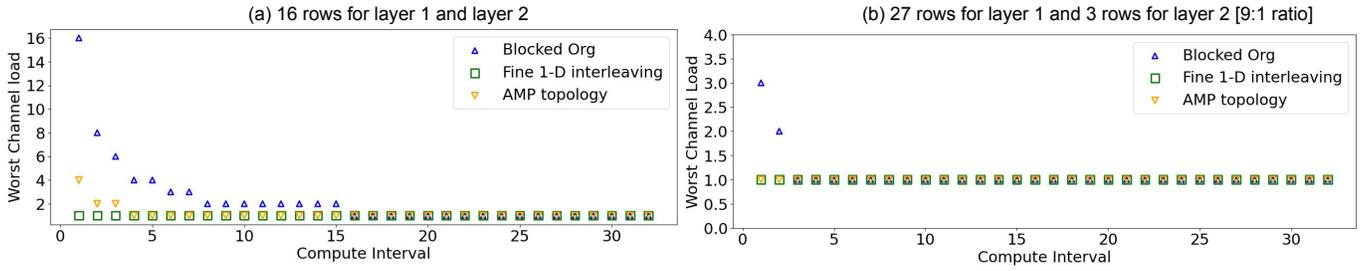


Fig. 15. Worst case channel load as a function of compute interval for 1-D spatial allocation with depth=2 for 32×32 mesh/AMP. We compare blocked organization, PIPEORGAN fine 1-D organization and AMP.

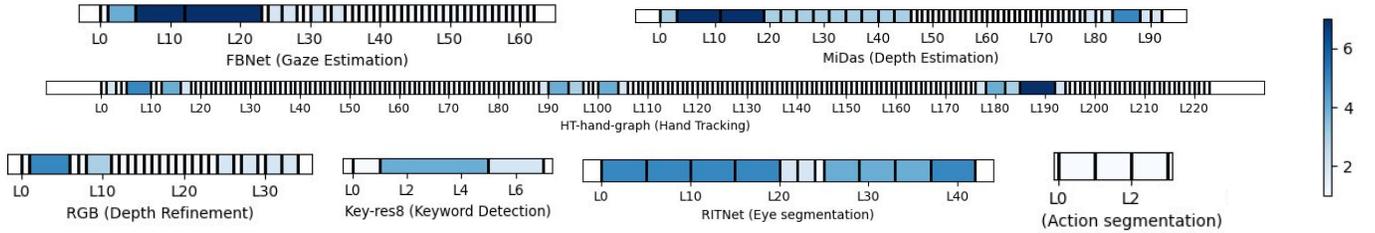


Fig. 16. Depths of XR-bench CNN tasks. L0 stands for Layer 0 and so on.

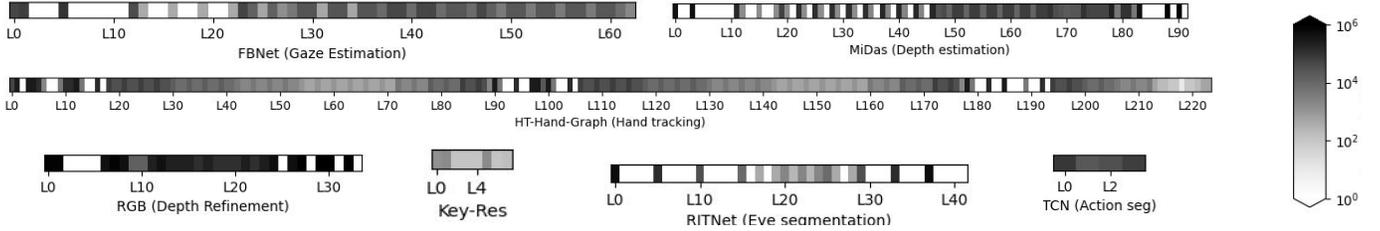


Fig. 17. Finest possible granularity based on depth and intra-operation dataflows across CNN tasks. L0 stands for Layer 0 and so on.

geomean. Note that this is end-to-end reduction, which also includes non-pipeline friendly layers (which in some cases, might also have larger tensors, and hence more accesses).

C. Congestion analysis: PIPEORGAN and AMP

Fig. 15 shows the variation of congestion with compute interval and spatial organization strategies. We show 1-D allocation with depth=2 since we observe most of the congestion in those cases, specially in TANGRAM-like dataflow. Two highly observed cases were equal allocation, and unequal allocation with 1×1 and 3×3 filters. Blocked organization has a large delay per packet for equal allocation case. The overall interval delay is *Worst case channel load* \times *Compute Interval*. For compute interval of 2 cycles, the overall communication delay increases by a factor of 8, as a result its 16. Fine-grained 1-D interleaving resolves this by avoiding congestion since it brings the consumer closer. Likewise, AMP topology, by virtue of long links, reduces the congestion delay and only incurs congestion when compute interval is below 4 cycles. The unequal allocation case incurs lower latency compared to the equal allocation counterpart, nevertheless, blocked organization is still more likely to cause congestion.

D. Pipeline Depth and Granularity

1) *Pipelining Depth across Tasks*: Fig. 16 shows the depths of CNN layers obtained after applying the depth heuristic, for

the entire models. Eye segmentation [4] task has the most regions with deep pipelining, primarily because of the high A/W ratios and skip connections. Keyword detection [38] in particular prefers pipelining despite nominal A/W ratios because of skip connections. Depth estimation [33] also has multiple deep pipelined regions. This is primarily because of DWCONV layers, which have a high A/W ratio, given that, weights are only along one channel. Moreover, these layers actually need pipelining since these layers are memory bound compared to a CONV with same number of multiplications.

2) *Finest Pipelining Granularity across Tasks*: Fig. 17 obtain the finest possible granularities assigned by stage 1, for the entire CNN models. However, this does not include the final granularities after the spatial organization is finalized. Granularity also depends on activation size, thus some cases appear fine-grained despite no pipelining. High-depth regions can have long regions of finer granularity.

VII. CONCLUSION

Modern DNN application domains consist of multiple models with vastly different layer shapes, layer dependencies and types of operations used. The right depth and granularity of inter-operation pipelining or lack thereof depends heavily on these model characteristics. Prior works on pipelining miss out on the opportunity to take advantage of fine-grained inter-operation pipelining opportunities and regions are allocated

to layers in a coarse-grained inter-layer manner, which has a significant on-chip communication overhead. In this work, we propose PIPEORGAN, a new class of spatial organization strategies for inter-operation pipelining, and a systematic methodology to choose the right depth, granularity and consequently the spatial organization. We also propose AMP, a topology that reduces the hops and congestion for coarse-grained spatial allocation.

ACKNOWLEDGEMENT

Part of this work was supported by ACE, one of the seven centers in JUMP 2.0, a Semiconductor Research Corporation (SRC) program sponsored by DARPA.

REFERENCES

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [2] M. Alwani, H. Chen, M. Ferdman, and P. Milder, “Fused-layer cnn accelerators,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2016, pp. 1–12.
- [3] J. Cai, Y. Wei, Z. Wu, S. Peng, and K. Ma, “Inter-layer scheduling space definition and exploration for tiled accelerators,” in *Proceedings of the 50th Annual International Symposium on Computer Architecture*, ser. ISCA '23. New York, NY, USA: Association for Computing Machinery, 2023. [Online]. Available: <https://doi.org/10.1145/3579371.3589048>
- [4] A. K. Chaudhary, R. Kothari, M. Acharya, S. Dangi, N. Nair, R. Bailey, C. Kanan, G. Diaz, and J. B. Pelz, “Ritnet: Real-time semantic segmentation of the eye for gaze tracking,” in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 2019, pp. 3698–3702.
- [5] Y.-H. Chen, J. Emer, and V. Sze, “Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks,” in *Proceedings of the 43rd International Symposium on Computer Architecture*, ser. ISCA '16. IEEE Press, 2016, p. 367–379. [Online]. Available: <https://doi.org/10.1109/ISCA.2016.40>
- [6] X. Dai, A. Wan, P. Zhang, B. Wu, Z. He, Z. Wei, K. Chen, Y. Tian, M. Yu, P. Vajda *et al.*, “Fbnetv3: Joint architecture-recipe search using predictor pretraining,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 276–16 285.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [8] M. Gao, X. Yang, J. Pu, M. Horowitz, and C. Kozyrakis, “Tangram: Optimized coarse-grained dataflow for scalable nn accelerators,” in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 807–820. [Online]. Available: <https://doi.org/10.1145/3297858.3304014>
- [9] R. Garg, E. Qin, F. Muñoz-Martínez, R. Guirado, A. Jain, S. Abadal, J. L. Abellán, M. E. Acacio, E. Alarcón, S. Rajamanickam, and T. Krishna, “Understanding the design-space of sparse/dense multiphase gnn dataflows on spatial accelerators,” in *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2022.
- [10] L. Ge, Z. Ren, Y. Li, Z. Xue, Y. Wang, J. Cai, and J. Yuan, “3d hand shape and pose estimation from a single rgb image,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 833–10 842.
- [11] J. Gu, H. Kwon, D. Wang, W. Ye, M. Li, Y.-H. Chen, L. Lai, V. Chandra, and D. Z. Pan, “Multi-scale high-resolution vision transformer for semantic segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 094–12 103.
- [12] K. He *et al.*, “Deep Residual Learning for Image Recognition,” in *CVPR*, 2016.
- [13] K. Hegde, H. Asghari-Moghaddam, M. Pellauer, N. Crago, A. Jaleel, E. Solomonik, J. Emer, and C. W. Fletcher, “Extensor: An accelerator for sparse tensor algebra,” in *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*, ser. MICRO '52. New York, NY, USA: Association for Computing Machinery, 2019, p. 319–333. [Online]. Available: <https://doi.org/10.1145/3352460.3358275>
- [14] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [15] Q. Huang, M. Kang, G. Dinh, T. Norell, A. Kalaiah, J. Demmel, J. Wawrzyniek, and Y. S. Shao, “Cosa: Scheduling by $\langle u \rangle$ -constrained $\langle o \rangle$ -optimization for $\langle s \rangle$ -spatial $\langle a \rangle$ -accelerators,” in *Proceedings of the 48th Annual International Symposium on Computer Architecture*, ser. ISCA '21. IEEE Press, 2021, p. 554–566. [Online]. Available: <https://doi.org/10.1109/ISCA52012.2021.00050>
- [16] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, R. Boyle, P. I. Cantin, C. Chao, C. Clark, J. Coriell, M. Daley, M. Dau, J. Dean, B. Gelb, T. V. Ghemmaghami, R. Gottipati, W. Gulland, R. Hagmann, C. R. Ho, D. Hogberg, J. Hu, R. Hundt, D. Hurt, J. Ibarz, A. Jaffey, A. Jaworski, A. Kaplan, H. Khaitan, D. Killebrew, A. Koch, N. Kumar, S. Lacy, J. Laudon, J. Law, D. Le, C. Leary, Z. Liu, K. Lucke, A. Lundin, G. MacKean, A. Maggiore, M. Mahony, K. Miller, R. Nagarajan, R. Narayanaswami, R. Ni, K. Nix, T. Norrie, M. Omernick, N. Penukonda, A. Phelps, J. Ross, M. Ross, A. Salek, E. Samadiani, C. Severn, G. Sizikov, M. Snellman, J. Souter, D. Steinberg, A. Swing, M. Tan, G. Thorson, B. Tian, H. Toma, E. Tuttle, V. Vasudevan, R. Walter, W. Wang, E. Wilcox, and D. H. Yoon, “In-datacenter performance analysis of a tensor processing unit,” in *Proceedings of the 44th Annual International Symposium on Computer Architecture (ISCA)*, 2017.
- [17] S.-C. Kao and T. Krishna, “Gamma: automating the hw mapping of dnn models on accelerators via genetic algorithm,” in *2020 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*. IEEE, 2020, pp. 1–9.
- [18] S.-C. Kao, S. Subramanian, G. Agrawal, and T. Krishna, “An optimized dataflow for mitigating attention performance bottlenecks,” *arXiv preprint arXiv:2107.06419*, 2021.
- [19] J. Kim, J. Balfour, and W. Dally, “Flattened butterfly topology for on-chip networks,” in *40th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO 2007)*, 2007, pp. 172–182.
- [20] T. Krishna, C.-H. O. Chen, W. C. Kwon, and L.-S. Peh, “Breaking the on-chip latency barrier using smart,” in *2013 IEEE 19th International Symposium on High Performance Computer Architecture (HPCA)*, 2013, pp. 378–389.
- [21] H. Kwon, P. Chatarasi, M. Pellauer, A. Parashar, V. Sarkar, and T. Krishna, “Understanding reuse, performance, and hardware cost of dnn dataflow: A data-centric approach,” in *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 2019, pp. 754–768.
- [22] H. Kwon, L. Lai, M. Pellauer, T. Krishna, Y.-H. Chen, and V. Chandra, “Heterogeneous dataflow accelerators for multi-dnn workloads,” in *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2021, pp. 71–83.
- [23] H. Kwon, K. Nair, J. Seo, J. Yik, D. Mohapatra, D. Zhan, J. Song, P. Capak, P. Zhang, P. Vajda *et al.*, “Xrbench: An extended reality (xr) machine learning benchmark suite for the metaverse,” *Proceedings of Machine Learning and Systems*, vol. 5, 2023.
- [24] H. Kwon, A. Samajdar, and T. Krishna, “MAERI: enabling flexible dataflow mapping over dnn accelerators via programmable interconnects,” in *Proceedings of the Twenty-Third International Conference on Architectural Support for Programming Languages and Operating Systems*. ACM, 2018, p. 461–475.
- [25] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager, “Temporal convolutional networks for action segmentation and detection,” in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 156–165.
- [26] S. Liang, Y. Wang, C. Liu, L. He, L. Huawei, D. Xu, and X. Li, “Engn: A high-throughput and energy-efficient accelerator for large graph neural networks,” *IEEE Transactions on Computers*, 2020.
- [27] C. Liu, K. Kim, J. Gu, Y. Furukawa, and J. Kautz, “Planercnn: 3d plane detection and reconstruction from a single image,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4450–4459.
- [28] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 10 012–10 022.
- [29] F. Ma and S. Karaman, “Sparse-to-dense: Depth prediction from sparse depth samples and a single image,” in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 4796–4803.
- [30] M. Naumov, D. Mudigere, H.-J. M. Shi, J. Huang, N. Sundaraman, J. Park, X. Wang, U. Gupta, C.-J. Wu, A. G. Azzolini, D. Dzhulgakov, A. Malleevich, I. Cherniavskii, Y. Lu, R. Krishnamoorthi, A. Yu, V. Kondratenko, S. Pereira, X. Chen, W. Chen, V. Rao, B. Jia, L. Xiong, and M. Smelyanskiy, “Deep learning recommendation model for personalization and recommendation systems,” *arXiv preprint arXiv:1906.00091*, 2019.
- [31] A. Parashar, P. Raina, Y. S. Shao, Y.-H. Chen, V. A. Ying, A. Mukkara, R. Venkatesan, B. Khailany, S. W. Keckler, and J. Emer, “Timeloo: A systematic approach to dnn accelerator evaluation,” in *2019 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2019, pp. 304–315.
- [32] E. Qin, A. Samajdar, H. Kwon, V. Nadella, S. Srinivasan, D. Das, B. Kaul, and T. Krishna, “Sigma: A sparse and irregular gemm accelerator with flexible interconnects for dnn training,” in *2020 IEEE International*

- Symposium on High Performance Computer Architecture (HPCA)*, 2020, pp. 58–70.
- [33] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, “Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 3, pp. 1623–1637, 2022.
- [34] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *Advances in neural information processing systems*, vol. 28, 2015.
- [35] Y. S. Shao, J. Clemons, R. Venkatesan, B. Zimmer, M. Fojtik, N. Jiang, B. Keller, A. Klinefelter, N. Pinckney, P. Raina, S. G. Tell, Y. Zhang, W. J. Dally, J. Emer, C. T. Gray, B. Khailany, and S. W. Keckler, “Simba: Scaling deep-learning inference with multi-chip-module-based architecture,” in *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*, ser. MICRO ’52. New York, NY, USA: Association for Computing Machinery, 2019, p. 14–27. [Online]. Available: <https://doi.org/10.1145/3352460.3358302>
- [36] Y. Shi, Y. Wang, C. Wu, C.-F. Yeh, J. Chan, F. Zhang, D. Le, and M. Seltzer, “Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 6783–6787.
- [37] A. Symons, L. Mei, S. Coleman, P. Houshmand, S. Karl, and M. Verhelst, “Towards heterogeneous multi-core accelerators exploiting fine-grained scheduling of layer-fused deep neural networks,” *arXiv preprint arXiv:2212.10612*, 2022.
- [38] R. Tang and J. Lin, “Deep residual learning for small-footprint keyword spotting,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5484–5488.
- [39] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, “Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 10734–10742.
- [40] M. Yan, L. Deng, X. Hu, L. Liang, Y. Feng, X. Ye, Z. Zhang, D. Fan, and Y. Xie, “Hygcn: A gcn accelerator with hybrid architecture,” in *2020 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2020, pp. 15–29.
- [41] X. Yang, M. Gao, Q. Liu, J. Setter, J. Pu, A. Nayak, S. Bell, K. Cao, H. Ha, P. Raina, C. Kozyrakis, and M. Horowitz, “Interstellar: Using halide’s scheduling language to analyze dnn accelerators,” in *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 369–383. [Online]. Available: <https://doi.org/10.1145/3373376.3378514>
- [42] H. You, C. Wan, Y. Zhao, Z. Yu, Y. Fu, J. Yuan, S. Wu, S. Zhang, Y. Zhang, C. Li *et al.*, “Eyecod: eye tracking system acceleration via flatcam-based algorithm & accelerator co-design,” in *Proceedings of the 49th Annual International Symposium on Computer Architecture*, 2022, pp. 610–622.
- [43] S. Zheng, X. Zhang, L. Liu, S. Wei, and S. Yin, “Atomic dataflow based graph-level workload orchestration for scalable dnn accelerators,” in *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2022, pp. 475–489.
- [44] S. Zheng, S. Chen, S. Gao, L. Jia, G. Sun, R. Wang, and Y. Liang, “Tileflow: A framework for modeling fusion dataflow via tree-based analysis,” in *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, 2023.