

Explainable-DSE: An Agile and Explainable Exploration of Efficient Hardware/Software Codesigns of Deep Learning Accelerators Using Bottleneck Analysis

Shail Dave
Shail.Dave@asu.edu
Arizona State University
Tempe, AZ, USA

Tony Nowatzki
tjn@cs.ucla.edu
University of California, Los Angeles
CA, USA

Aviral Shrivastava
Arizona State University
Tempe, AZ, USA
Aviral.Shrivastava@asu.edu

ABSTRACT

Effective design space exploration (DSE) is paramount for hardware/software codesigns of deep learning accelerators that must meet strict execution constraints. For their vast search space, existing DSE techniques can require excessive number of trials to obtain valid and efficient solution because they rely on black-box explorations that do not reason about design inefficiencies. In this paper, we propose Explainable-DSE – a framework for DSE of DNN accelerator codesigns using bottleneck analysis. By leveraging information about execution costs from bottleneck models, our DSE is able to identify the bottlenecks and therefore the reasons for design inefficiency, and can therefore make mitigating acquisitions in further explorations. We describe the construction of such bottleneck models for DNN accelerator domain. We also propose an API for expressing such domain-specific models and integrating them into the DSE framework. Acquisitions of our DSE framework caters to multiple bottlenecks in executions of workloads like DNNs that contain different functions with diverse execution characteristics. Evaluations for recent computer vision and language models show that Explainable-DSE mostly explores effectual candidates, achieving codesigns of 6× lower latency in 47× fewer iterations vs. non-explainable techniques using evolutionary or ML-based optimizations. By taking minutes or tens of iterations, it enables opportunities for runtime DSEs.

ACM Reference Format:

Shail Dave, Tony Nowatzki, and Aviral Shrivastava. 2023. Explainable-DSE: An Agile and Explainable Exploration of Efficient Hardware/Software Codesigns of Deep Learning Accelerators Using Bottleneck Analysis. In *Proceedings of ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*. ACM, New York, NY, USA, 21 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Domain-specific accelerators, e.g., for deep learning models, are deployed from datacenters to edge. In order to meet strict constraints on execution costs (e.g., power and area) while minimizing

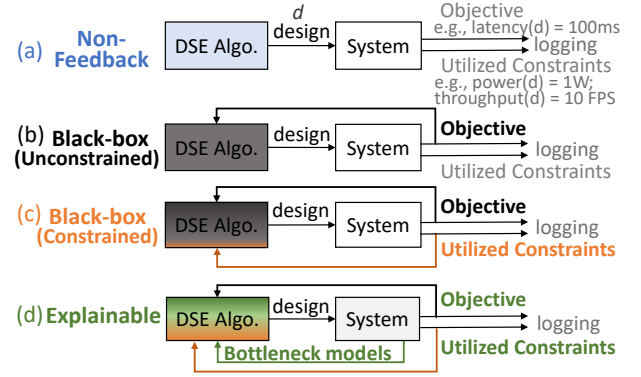


Figure 1: DSE with (a) Non-feedback, (b) Unconstrained black-box, (c) Constrained black-box, and (d) Explainable optimization, which leverages domain-specific bottleneck models.

an objective (e.g., latency), their hardware/software codesigns must be effectively explored using an effective *design space exploration* (DSE). However, the search space is vast, and it can contain $O(10^{29})$ solutions, with each evaluation taking milliseconds–minutes. For instance, [74] showed that a TPU-like architecture has 10^{14} hardware solutions with modest options for design parameters. For every hardware configuration, software space can also be huge. For example, DNN layers can be mapped on a spatial architecture in $O(10^{15})$ ways aka dataflows (as shown by [17] and Table 7). Clearly, an effective exploration is needed to achieve feasible and efficient solutions quickly.¹

Recent DSE techniques for deep learning accelerators use either non-feedback or black-box optimizations. *Non-feedback* optimizations include grid search (in [26, 49]) and random search (in [35, 45]). They evaluate different solutions for a pre-set number of iterations and terminate (Fig. 1a). *Black-box* optimizations, on the other hand, consider value of objective before acquiring² the next candidates (Fig. 1b-c). Thus, they can be more effective than non-feedback approaches. They include simulated annealing [63], genetic algorithm [59, 65, 73], Bayesian optimization [28, 40, 43, 46, 50, 69, 74], and reinforcement learning [10, 30, 67, 70, 77]. These optimizations can be unconstrained or constrained.

¹A feasible solution meets all constraints, and its hardware and software configurations are compatible; An efficient solution minimizes objective; Agility refers to DSE’s ability to find desired solutions quickly, which becomes crucial for exploring vast space in practical DSE budgets and runtime DSEs.

²Acquisition refers to a step in a DSE algorithm that selects next set of candidate designs to evaluate. §A.1 discusses terminology for the DSE techniques.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

ASPLOS ’24, April 27–May 01, 2024, San Diego, CA

© 2023 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

For vast accelerator hardware/software codesign space, existing techniques require excessive trials for convergence or even finding a feasible solution. We argue that this is because of lack of explainability during the exploration. By *explainability* of a DSE technique, we imply its ability to reason about, at each acquisition attempt, why a certain design corresponds to specific costs, and what are underlying inefficiencies, and how they can be ameliorated. Existing exploration techniques are non-explainable as in they lack information and reasoning about the quality of designs acquired during DSE. They may figure out which of the previous trials reduced the objective but they cannot determine *why*. In contrast, an explainable DSE would identify the inefficiencies of the acquired design that incur high costs and also estimate mitigation requirements that would improve designs and execution further. For instance, in reducing the latency of a DNN accelerator, an explainable DSE could reason that the latency is dominated by memory access time that cannot be hidden behind the time for computation or communicating data on-chip. Therefore, it could strive to reduce latency further by increasing off-chip bandwidth or on-chip buffer size to exploit available data reuse.

The goal of this paper is to develop a framework for explainable DSE of HW/SW codesigns of DNN accelerators.

Our approach is to use bottleneck analysis to enable explainability in DSE of DNN accelerator/dataflow codesigns. Enabling explainability in DSE with bottleneck analysis requires *bottleneck models*. Conventional DSE approaches evaluate only *cost models* in DSE that provide just a single value like latency. In contrast, the domain-specific bottleneck models can provide richer information about how design parameters contribute to various execution factors like the time for computation, memory accesses, and communication via NoCs, which in turn, leads to total cost such as latency (Fig. 2). Bottleneck models also provide mitigation strategies, i.e., when one of these factors say on-chip communication gets identified as a bottleneck, how to tune different design parameters based on key execution-related characteristics of the workloads (e.g., increase bit-widths of NoCs by certain amount or increase physical links or time-shared unicast support). These bottleneck models can be developed based on domain-specific information, which is often embedded in experts-defined, domain-specific cost models (like [14, 36]) but *implicitly*. Having the *explicit* bottleneck models and their driving the DSE can help DSE explain inefficiencies of acquired designs (referred to as *bottleneck analysis*) and to make mitigating acquisition decisions.

For enabling DSE of DNN accelerators using bottleneck analysis, our approach overcomes the following shortcomings.

1) We develop a bottleneck model for deep learning accelerator design domain. Taking latency minimization as an example, we describe what execution characteristics of DNN accelerators need to be leveraged, how to construct a corresponding bottleneck model, how its bottleneck graph provides insights in execution inefficiencies of design and how to pinpoint bottlenecks with it, and what are mitigation strategies once a bottleneck is identified. By

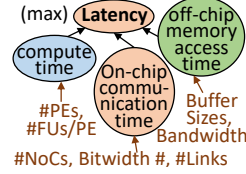


Figure 2: Example bottleneck graph of DNN accelerator latency.

applying bottleneck analysis on software-optimized executions of each hardware design, our DSE co-explores both hardware-software configurations of DNN accelerators in adaptive and tightly coupled manner.

2) We propose an API for interfacing DSE with domain-specific bottleneck models. Through proposed API, bottleneck model of a system can be described as a tree corresponding to the target cost. Navigating such tree enables DSE to analyze the bottlenecks, relate the bottlenecks with the design parameters, and reason about the desired scaling for mitigations. For instance, by parsing a latency tree (Fig. 2), DSE could reason that latency is a maximum value of the time taken for computations, on-chip communications, and memory accesses; if computational time exceed other factors by 3×, then related parameters are number of functional units in PEs and number of PEs, which may need to be scaled next.

The API can allow expert designers to systematically express their domain-specific bottleneck models and integrate them in DSE while leveraging constrained exploration framework. This helps overcome a limitation of previous DSEs using bottleneck analysis in other domains like multimedia or FPGA-HLS [20, 23, 58, 72] which lack such interface; as search mechanisms were defined in domain-specific ways for their bottleneck models, they cannot be decoupled or reused for other domains.

3) We propose a generic framework for constrained DSE using bottleneck models, with acquisitions accounting for multiple bottlenecks in multi-workload executions. Prior frameworks for DSE using bottleneck analysis in other domains optimize only a single task at a time, i.e., consider a single cost value of executing a loop-kernel or a whole task and iteratively mitigate its bottleneck. However, when workloads involve different functions of diverse execution characteristics, e.g., a DNN with multiple layers or multiple DNNs, changing a design parameter impacts their contribution to overall cost in distinct ways; considering just a total cost could not be useful. Also, mitigation strategies to address layer-wise bottlenecks can lead to range of different values for diverse parameters. So, our framework systematically aggregates parameters predicted for mitigating bottlenecks in executions of multiple functions in one or more workloads, for making next acquisitions.

Results: We demonstrate our explainable and agile DSE framework by exploring high-performance edge inference accelerators for recent computer vision and language processing models. By iteratively mitigating bottlenecks, Explainable-DSE reduces latency under constraints in mostly every attempt (1.3× on average). Thus, it explores effectual candidates and achieves efficient codesigns in *minutes*, while non-explainable optimizations may fail to obtain even a feasible solution over days. Explainable-DSE obtains codesigns of 6× lower latency in (36× less search time on average and up to 1675×) 47× fewer iterations vs. previous DSE approaches for DNN accelerators. By achieving highly efficient solutions in only 54 iterations, Explainable-DSE enables opportunities for cost-effective and dynamic explorations in vast space.

2 LIMITATIONS OF PRIOR DNN ACCELERATOR DSES

Non-feedback DSE approaches search either exhaustively over statically reduced space (e.g., grid search in [26, 41]) or randomly (e.g., in

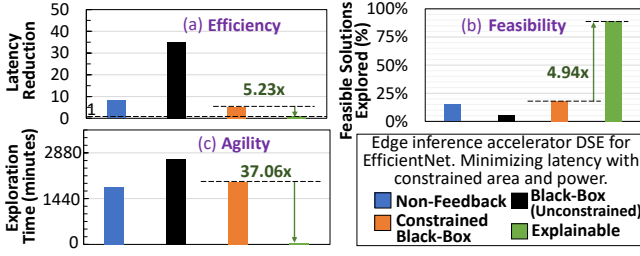


Figure 3: Effectiveness of non-explainable and explainable DSE frameworks for exploring efficient and feasible solutions in the vast space: (a) Efficiency (latency minimization); (b) Feasibility (in % of the total solutions evaluated); (c) Agility (exploration time in minutes). Analysis is shown here for exploring edge accelerator design for EfficientNetB0 model.

[45]). So, they do not consider any outputs like objective or utilized constraints and terminate after using a large exploration budget. This is illustrated in Fig. 1(a). On the other hand, black-box optimizations such as Bayesian Optimization (e.g., in [43, 46, 50, 69, 74]) consider values of objective for previously tried solutions. This is illustrated in Fig. 1(b)–(c). Considering the objective helps them predict the likelihood of where the minima may lie; they acquire a candidate for the next trial accordingly. The process repeats until convergence or the number of trials exceeding a threshold. While black-box DSE could be more efficient than non-feedback DSE, they all face the following limitations:

DSE techniques lack reasoning about bottlenecks incurring high costs: An efficient DSE mechanism should determine challenges hindering the reduction of objectives or utilized constraints. It should also determine which of the many parameters can help mitigate those inefficiencies and with what values. However, with the objective as only input, these black-box or system-oblivious DSEs can figure out only which prior trials reduced the objective. But, they are non-explainable as in they cannot reason about what costs a solution could lead to and why – a crucial aspect in exploring enormous design space. This is exacerbated by the fact that execution characteristics of different functions in workloads are diverse (e.g., memory- vs. compute-bounded DNN operators; energy consumption characteristics). By considering just total cost value, black-box DSEs cannot consider diverse bottlenecks in multi-functional or multi-workload executions that need to be systematically addressed.

Implications: A major implication of excessive sampling caused by lack of explainability is **inefficiency of obtained solutions**. Fig. 4(a) illustrates this through a toy scenario, i.e., exploring the number of PEs and global buffer size for a single ResNet layer. It shows exploration from early iterations to later iterations with HyperMapper 2.0 [43] – an efficient, Bayesian-based optimizer. The figure shows that even for a tiny space, acquired solutions are mostly inefficient (high latency), as there is no reasoning about underlying bottlenecks and their mitigation. So, even though DSE has already acquired some better solutions before, later acquisitions correspond to inefficient solutions. We find that as the space becomes vast, the non-explainable DSE techniques can require too many trials (at

least, in thousands [30, 43, 74]), and they may still not find the most efficient solutions. As an example, Fig. 3(a) shows that the latency of the solutions obtained by non-explainable DSEs can be up to 35× higher, even for 2500 trials (two days of search time). This is because practical exploration budget is typically fractional (thousands) compared to the vast design space (quadrillions). By generating trials without understanding executional bottlenecks and their mitigation, most of the search budget gets used for excessive and mostly ineffectual trials.

Lacking reasoning about design’s inefficiencies can deprive the DSE of *tightly coupled hardware/software codesign*. For instance, DSEs in [10, 30, 34, 50, 55, 70, 77] mainly explore architectural parameters with black-box DSEs and use a fixed dataflow for executions.³ Fixing the execution methods limit the effectual utilization of architectural resources when subjected to various tensor shapes and functionalities [9, 12]. Consequently, DSEs may achieve architecture designs that are either incompatible with the dataflow (**infeasible solutions**) or **inefficient**. Likewise, separate optimizations of architectural design and dataflow that are *oblivious* of each other can lead to excessive trials and inefficient solutions. Further, for these constrained optimizations, excessive trials are also caused by the fact that DSEs cannot determine which constraints are violated and how configuring different accelerator design parameters could affect that. Fig. 3(b) illustrates this for an edge accelerator DSE that is subjected to power and area constraints. Out of 2500 solutions evaluated, for constrained optimizations like HyperMapper2.0 [43], only 18% of the all solutions evaluated were feasible and up to only 52% for constrained reinforcement learning [30].

Another implication of excessive trials is **inapplicability to dynamic DSE scenarios**. Excessive trials lead to low agility, as illustrated in Fig. 3(c). Non-explainable DSEs consume very high exploration time, even *weeks*, while obtaining solutions of lower efficiency. This makes existing DSE approaches unsuitable for dynamic explorations (e.g., DSE convergence within a few tens to 100 iterations). For instance, unlike one-time ASIC designs, deploying accelerator overlays over FPGAs (edge/cloud; dedicated/multi-tenant) can benefit from dynamic DSEs, where constraints for DSE and resource budget may also become available just before deployment.

3 DSE OF DEEP LEARNING ACCELERATORS USING BOTTLENECK ANALYSIS: MOTIVATION AND CHALLENGES

3.1 Making DSE Explainable Through Bottleneck Analysis

In Fig. 4(b), we illustrate the same search problem of designing a DNN accelerator as in Fig. 4(a), but by using bottleneck analysis in the DSE. Before acquiring a candidate, the DSE analyzes the current design through the bottleneck model and pinpoints the bottleneck in achieved latency. Then, it uses mitigations suggested by the bottleneck model to make the next acquisitions. The bottleneck, in terms of latency optimization for a deep learning accelerator, can be attributed to execution factors such as time consumed by computations, communication via NoCs, and off-chip memory accesses with direct memory access (DMA) controller. For instance,

³§A.2 and §G provide background on HW/SW codesign DSE process.

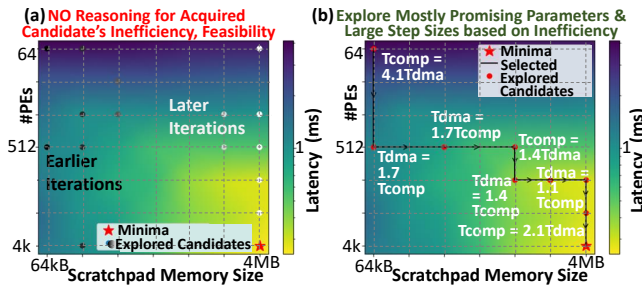


Figure 4: Example DSE of #PEs and shared memory sizes for ResNet CONV5_2b [24], with: (a) Prior techniques (HyperMapper2.0 [43]), and (b) Explainable-DSE, which reasons about inefficiencies in achieved executions, limiting the search to crucial parameters and tuning them accordingly.

after evaluating the initial point (*number of PEs, shared memory size*) = (64, 64kB), the DSE can reason that the computation time of the design is $4.14\times$ higher than the time taken by off/on-chip data communication. From the mitigation strategy, DSE concludes and communicates to the designers that it would scale the total number of PEs next by at least $4.14\times$.⁴ Since this is the only mitigation suggested, the newly acquired and optimized design becomes (512, 64kB). By repeating this process, the DSE informs that the previous bottleneck got mitigated and DMA-transfers is the new bottleneck. Using the bottleneck model, the DSE considers execution characteristics (like data accessed from off-chip memory and unexploited data reuse) and mitigation for the current design point, adjusting the size of shared on-chip memory or off-chip bandwidth. This iterative process continues. It enables the DSE to not just characterize, explain, and optimize DSE decisions and acquired designs, but also optimize objectives at almost every acquisition attempt and converge to efficient solutions quickly.

3.2 Challenges in Enabling DSE of DNN Accelerators Using Bottleneck Analysis

Need bottleneck models for DNN accelerator domain. DSE using bottleneck analysis requires bottleneck models. Unlike cost models used in black-box DSEs that provide a single value, bottleneck models can provide richer information about 1) how design parameters contribute to different factors that finally lead to the overall cost, and 2) mitigation strategies when any of those factors gets identified as a bottleneck. Such bottleneck/root-cause analysis have been developed/applied for characterizing fixed designs and finding mitigation strategies, e.g., for industry pipelines and production systems, hardware or software for specific applications [13, 61, 72], FPGA-based HLS [23, 58], overlapping microarchitectural events [19], and power outage [22]. Likewise, optimizing DNN accelerator designs with bottleneck analysis also require developing bottleneck models.

Need an interface to decouple domain-specific bottleneck models from a domain-independent exploration mechanism

⁴Just to note the power of explicit bottleneck mitigation strategies, if area constraint was unmet, DSE could intelligently let communication time increase but meet constraints first through reduced buffer/NoC sizes.

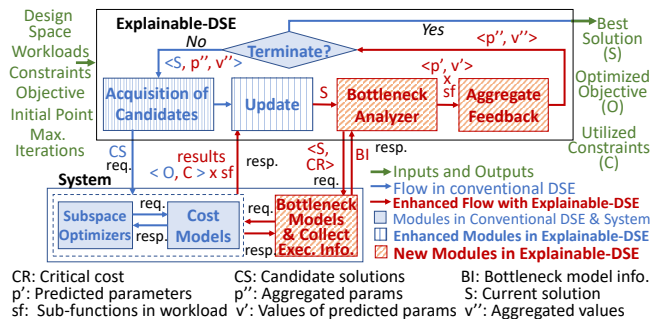


Figure 5: Explainable-DSE: A framework for exploring design space using domain-specific bottleneck models.

and express them to DSE. Once bottleneck models are developed, there needs to be a DSE framework that can integrate such a domain-specific bottleneck model to drive the iterative search. However, since bottleneck models are usually domain-specific, search mechanisms provided by prior DSE techniques using bottleneck analysis [20, 58, 72] are implemented too specifically for their domain. There needs to be an interface to decouple the domain-independent search mechanism from domain-specific bottleneck models so that designers can reuse and apply the same search mechanism for exploring designs in new domains like DNN acceleration.

Need acquisitions accounting for mitigations of multiple bottlenecks in workload executions. Prior DSE techniques using bottleneck analysis (in other domains) [20, 23, 58, 72] optimize only a single task at a time, i.e., consider a single cost value of executing a loop-kernel or whole task and iteratively mitigate arising bottleneck. However, when workloads involve different functions of diverse execution characteristics, e.g., a DNN with multiple layers or multiple DNNs, changing a design parameter impacts their contribution to the overall cost in distinct ways; considering just a total cost may not be useful. Mitigation strategies to address these layer-wise bottlenecks can lead to changing diverse parameters and a range of values possible for the same parameter. Therefore, when DSE framework makes its next acquisitions, it needs to ensure that multiple bottlenecks arising from executing different functions of target workloads are mitigated systematically and effectively.

4 EXPLAINABLE-DSE: A CONSTRAINED DSE FRAMEWORK USING BOTTLENECK ANALYSIS

This section presents Explainable-DSE – a framework for an agile and explainable DSE using bottleneck analysis for optimizing deep learning accelerator designs. First, we discuss our framework’s overall workflow and illustrate it with a walk-through example. Then, we describe how its bottleneck analyzer processes bottleneck models, i.e., determines factors incurring a high cost, parameters relevant to the bottleneck factors, and new values of parameters that can reduce the cost. We also introduce an API through which architects can specify domain-specific bottleneck models, e.g., for accelerator execution costs. For bottleneck analysis involving the execution of multiple functions within workloads, we discuss how Explainable-DSE aggregates the obtained parameters and their new

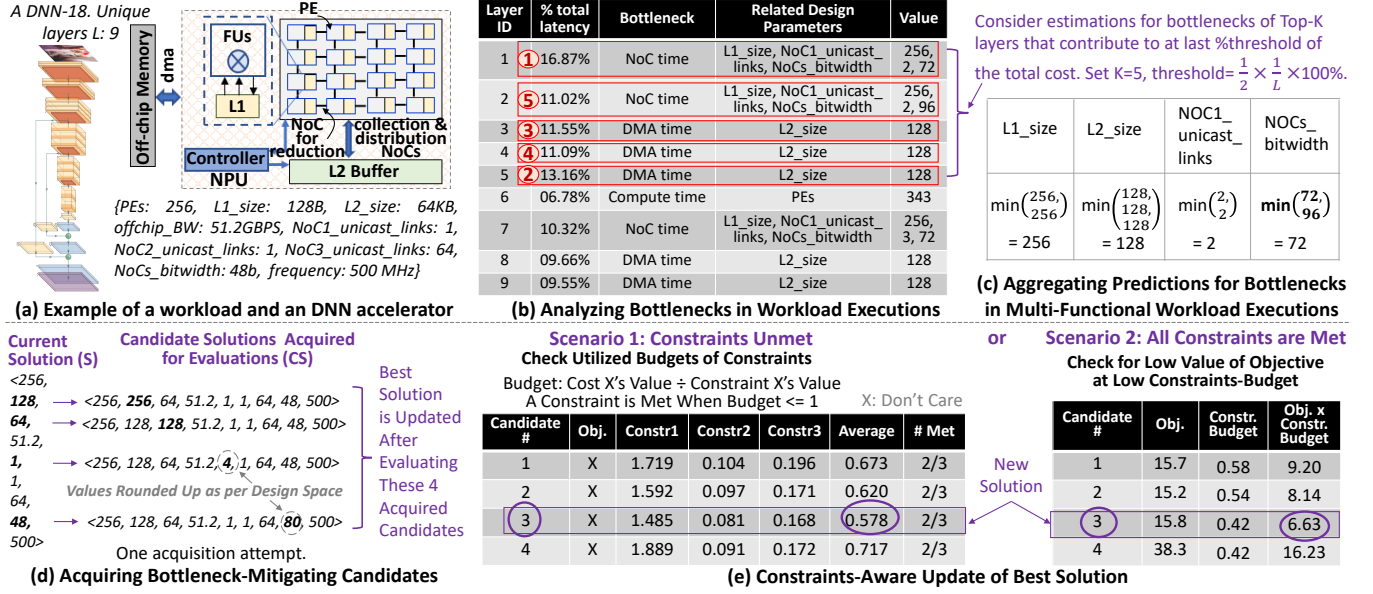


Figure 6: Example Walkthrough. a) A DNN and accelerator architecture parameters. b) Analyzing bottlenecks in multi-functional workload executions (§4.3). c) Aggregations of predictions for multiple bottlenecks (§4.4). d) Acquiring new candidates that mitigate bottlenecks (§4.5). e) Constraints-accommodating update of new solution (§4.6).

values, including considering bottlenecks of only execution-critical functions. We then describe how the proposed framework considers inequality constraints when updating the obtained solutions, prioritizing the exploration of feasible regions. We provide an in-depth bottleneck model to analyze and mitigate bottlenecks in exploring low-latency designs of DNN accelerators using Explainable-DSE. Lastly, we discuss how our approach can enable tightly coupled accelerator/mappings co-explorations.

4.1 Framework Workflow

Fig. 5 illustrates the workflow of Explainable-DSE. The DSE uses bottleneck analysis to explore solutions that reduce a critical cost, denoted as CR . Critical cost is usually an objective O that needs to be minimized and optionally an unmet inequality constraint value C . To reduce a critical cost, the bottleneck analyzer considers the current best solution (S) and analyzes cost-related bottleneck information (I). The analyzer identifies the bottleneck factors incurring higher cost value and finds the *scaling* “ s ” by which the objective/constraint value needs to be reduced (“ s ” is internal to the analyzer, so not shown in Fig. 5). Then, the analyzer determines design parameters (p') crucial for mitigating the bottleneck and their values (v'). Workloads usually involve multiple functions or sub-functions (sf), e.g., different DNNs or layers in a DNN. So, the DSE applies bottleneck analysis to the costs of each function individually and aggregates the corresponding feedback obtained. This aggregation leads to a set of predicted design parameters (p'') and their respective values (v''). It corresponds to a new set of candidate solutions (CS) for the subsequent acquisition. The process iterates, as depicted in Fig. 5. We refer to acquiring and evaluating candidates in a CS as one “*acquisition attempt*”. It is analogous to z sequential DSE iterations if there are z candidates in a CS . The best solution,

S , is updated once (from z candidates) at every acquisition attempt. When some inequality constraint is not met, the framework considers the utilized budgets of constraints for acquired candidates in updating the best solution. This approach enables the DSE to prioritize reaching feasible subspaces. In Fig. 5, the introduction of new modules for the proposed approach and corresponding information flow is illustrated through a diagonal stride pattern and a different shade (red). The workings of these modules are described next, accompanied by a walk-through example (illustrated in Fig. 6). Additional information regarding the capabilities of the framework, current limitations, and future works are discussed in §B and §C, respectively.

4.2 Framework Inputs and Outputs

Inputs: Information of design space, constraints, objective, workloads, initial point, and total iterations. **Outputs** upon convergence or termination: Optimized solution and its costs.

Design Space: It defines parameters of type integer, real, or categorical. Their possible values can be expressed as either a list or a mathematical expression.

Constraints and Objective: Users can define inequality constraints on multiple costs. Our current implementation optimizes a single objective. It can be extended for multiple objectives through existing acquisition techniques.

Target System and Cost Models: System can incorporate arbitrary cost models and subspace optimizations for populating costs. It can also provide costs at sub-functions granularity, e.g., the latency of individual DNN layers. The proposed API (§4.3) enables the seamless integration of bottleneck models.

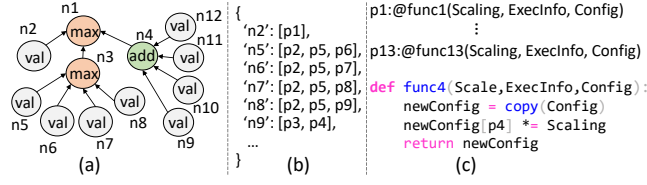


Figure 7: Proposed API through which architects can specify a bottleneck model of a system. The information can contain: (a) Bottleneck graph containing factors contributing to a cost; (b) Different parameters impacting the factors; (c) Handles to subroutines that calculate new values of the parameters.

To demonstrate DNN accelerator design explorations, we leverage existing cost models and use them to evaluate all techniques. We use Accelergy [68] to obtain the total area, energy per data access (for 45nm technology node), and maximum power. The maximum power is obtained from the maximum energy consumed by all design components in a single cycle. Accelergy provides technology-specific estimations via plugins for Aladdin [57] and CACTI [42]. We use our dMazeRunner infrastructure [14] to obtain the latency and energy consumed by mappings of DNN layers and for quick mapping optimizations for each architecture design.

4.3 Bottleneck Analyzer

Before each acquisition attempt, Explainable-DSE conducts bottleneck analysis on the previously obtained best solution. It uses the bottleneck model, which helps pinpoint the execution bottlenecks and suggests solutions to mitigate them (as detailed in §4.7), ultimately reducing costs. For instance, Fig. 6 demonstrates this exploration process for an 18-layer DNN, where nine layers have unique tensor shapes for execution-critical operators (CONV and GEMM). Fig. 6(a) shows the architectural template and parameter values of the current best solution during the DSE. Fig. 6(b) displays the bottleneck analyzer's ability to identify bottlenecks for each DNN layer and estimate which parameters should be updated with what specific values. This section further explains how the analyzer works and presents an API through which designers can specify their domain-specific bottleneck models for the DSE.

By evaluating the bottleneck model, the bottleneck analyzer determines (a) bottleneck factors, (b) parameters that are most critical for reducing the costs of these bottleneck factors, and (c) values of these critical parameters. Designers can provide the information for bottleneck models through an API that comprises three data structures, as illustrated in Fig. 7. The first is a bottleneck graph, which outlines the underlying factors contributing to the total cost. The second includes a list of related parameters for each factor. The third contains handles to subroutines that predict the next values of parameters. When some information is unavailable, such as how to predict the value of a parameter, Explainable-DSE resorts to its black-box counterpart (e.g., sampling neighboring values).

(a) Determining bottleneck factors from a bottleneck graph containing execution factors: A bottleneck graph in the bottleneck model outlines how various factors contribute to a workload's execution cost, as depicted in Fig. 7(a). It is represented as a tree

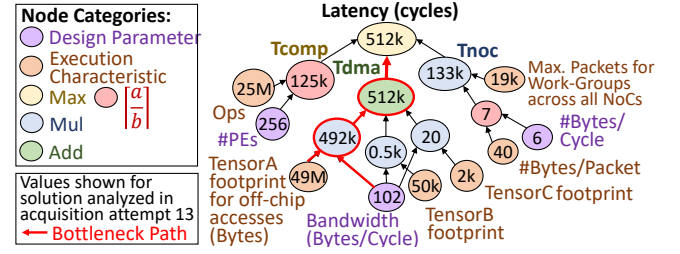


Figure 8: A simplified bottleneck graph for analyzing the latency of a DNN layer execution on a DNN accelerator.

whose nodes are mathematical functions like addition, multiplication, division, and maximum. Each node typically represents a cost factor, which is calculated from its children by applying the corresponding mathematical function. For example, Fig. 8 shows a simplified bottleneck graph for a DNN layer execution, where the root corresponds to the overall cost (e.g., latency). The total cost depends on child nodes representing underlying cost factors, such as computational time or data communication time. For example, the total latency is determined as the maximum value among the computational time, the total on-chip communication time, and the total DMA time for off-chip memory accesses. The total DMA time, in turn, is additive and depends on the off-chip footprint of different tensors and the bandwidth. Similarly, the time for communicating data from on-chip buffers to PEs via NoCs is approximated with the total data packets communicated to different workgroups and NoC bus widths. Thus, leaf nodes typically represent values of primary components, such as *design parameters* and the accelerator's *execution characteristics* for a given workload/application. The execution characteristics include data allocated to buffers, on/off-chip communication of the data, unexploited reuse, etc. (§4.7).

During each acquisition attempt, the analyzer considers the obtained solution and populates the graph with the corresponding actual values. For each cost factor, which is an intermediate node, the analyzer calculates its contribution to the total cost as the ratio of its value to the total cost. The analyzer traverses the graph and computes the contribution of each factor based on the associated mathematical operation. For instance, at a max node, it traces back to the maximum value; at an add node, it counts contributions proportionally. It identifies the factor with the highest contribution as the primary bottleneck. The analyzer then calculates the **scaling "s"**, which is the ratio by which the cost of the bottleneck factor should be reduced to alleviate the bottleneck. In Fig. 8, DMA time dominates the total latency, whereas the computational and on-chip communication time contributes to only 24.4% and 25.9% of the total latency, respectively. The analyzer finds that the later latency factors can be balanced by scaling the DMA time down, e.g., by a factor of $100\% \div 25.9\%$ or $3.85\times$. Through traversal, the analyzer identifies the memory footprint of tensor A as the primary bottleneck operand. The analyzer may also determine multiple bottlenecks (based on decreasing order of their contributions) so that the acquisition function can generate an adequate number of candidates.

(b) Selecting Parameters Associated with the Bottleneck: To determine which parameters impact specific bottleneck factors,

the analyzer can traverse the bottleneck graph, or designers can provide this information through a dictionary that maps the node names/numbers to relevant parameters (Fig. 7b). In the example bottleneck graph of Fig. 7(a), nodes 'n4' and 'n9' correspond to DMA time and the off-chip footprint of Tensor A, respectively. They are associated with parameters 'p3' and 'p4' (e.g., 'L2_size' and 'offchip_BW' in Fig. 6a). Once the bottleneck factor and mitigating parameters are identified, DSE obtains new values from supporting subroutines.

(c) Obtaining Values of Critical Parameters with Mitigation Strategies: Designers can provide handles to domain-specific subroutines that contain mitigation strategies for different design parameters, as shown in Fig. 7(c). Each subroutine calculates the new value of a parameter based on the current parameter value, the scaling s required for reducing the bottleneck factor, and the execution characteristics of the current design configuration (§4.7). For example, the function 'func4' can scale the off-chip bandwidth to reduce DMA time, and functions 'func5' to 'func8' can scale the bus width or links of NoCs to lower the on-chip communication time. The DSE leverages these subroutines to predict new values of critical parameters and evaluates the corresponding design points to identify the best configuration.

4.4 Addressing Bottlenecks in Multi-Functional Execution

As Fig. 6(b) illustrates, the analyzer performs bottleneck analysis on each sub-function of workloads (DNN layer) one by one. Due to the diverse execution characteristics of these functionalities, the predictions obtained for each sub-function can be distinct, depending on the factors like available reuse and parallelism. Additionally, predictions for mitigating multiple bottleneck factors of various DNN layers may involve the same parameter. Hence, an aggregation is required to determine the next set of parameters and their values (Fig. 6c). The DSE employs two methods for aggregation:

(i) *Aggregating different values of the same parameter:* After analyzing solution S for bottlenecks of multiple sub-functions, there can be different predicted values of the same parameter. So, the final prediction can be obtained by either iterating over some of these values or applying a function (maximum, minimum, average) on the predicted values. Choosing the maximum value can lead to faster convergence, but it can favor a single sub-function and be overly aggressive for others. For instance, selecting a new value as $16\times$ (from options like $4\times$, $8\times$, $16\times$ the current number of PEs) can significantly reduce latency of a non-performance-critical DNN layer but not of other layers, while consuming higher area and power. Thus, exploration can quickly exhaust the budget for constraints without getting a chance to explore a considerable range of intermediate candidates that could minimize the overall cost. Instead, we opt for the logical minimum of the estimated values as the final prediction (shown in Fig. 6c).

(ii) *Aggregating parameters from only bottleneck sub-functions:* Not all sub-functions or cost factors require improvement. Explainable-DSE allows focusing on only the bottleneck ones, i.e., contributing the most to the total cost. This capability is achieved through two user-tunable parameters: K and $threshold$. The DSE considers predictions from up to top- K sub-functions whose fractional

contributions to the total cost exceed a certain *threshold*. In target DNNs, the number of layers with unique tensor shapes (l) can range from a few to several tens. So, we set K to five and the *threshold* to $0.5 \cdot (1/l) \cdot 100\%$, considering predictions from layers that consume higher portions of the cost. As Fig. 6 shows, the analyzer considers mitigating bottlenecks from the top-5 layers that contribute at least 5.5% to the total latency.

4.5 Bottlenecks-Aware Acquisitions of Candidates

After aggregating parameter values for mitigating bottlenecks, the DSE populates candidate solutions CS to be acquired next. For simplicity, our acquisition function samples a candidate for each new value of a parameter. As Fig. 6(d) shows, all but one parameter of the candidate has the same value as the current solution. This mechanism naturally facilitates an iterative search that adaptively tunes among bottleneck parameters. Our acquisition avoids falling into a greedy local search [48] by the following means. i) It limits exploration parameters to only a few (critical for addressing the bottleneck); ii) It can predict values of larger step-size (non-neighbors) based on bottleneck mitigation analysis (whereas local search explores all p immediate neighboring values of all p parameters in the selected solution). Acquisitions by addressing multiple, dynamic bottlenecks (different parameters to be optimized at each DSE iteration) and exploring larger step sizes usually help avoid over-optimization of a design within the local neighborhood (converging to local optimal). §C further discusses workarounds for greediness in the search. Due to the modular design of the framework, users also may specify other acquisition/update functions that act upon bottlenecks-mitigating parameters. When acquiring a candidate, if a predicted value is not present in the defined design search space (e.g., non-power-of-2), the DSE rounds it up to the closest value.

4.6 Constraints-Budget-Aware Update of Solution

When exploring a vast space under tight constraints, initially acquired solutions usually fail to meet all constraints (e.g., low-area, high-latency region, or vice versa). To effectively explore the space, our DSE accounts for the *constraints budget* when selecting the best solution, which, in turn, impacts the acquisitions of new candidates. In determining a new solution from explored candidates, our DSE first checks whether the solutions meet all constraints and by what margin. If any candidate does not meet all constraints, it selects a candidate as the best solution that uses the least *constraints budget*. The constraints budget is calculated as the average of the utilized constraint values that are normalized to constraint thresholds. Such accounting is illustrated in Fig. 6(e) - scenario 1. Further, for monomodal cost models, when a candidate (corresponding to the new value of some parameter) violates more constraints than the obtained solution, the DSE can disable further exploration of the parameter's range. Thus, by prioritizing the feasibility of solutions, the DSE limits acquiring solutions that optimize the objective at the expense of violating constraints. When multiple candidates satisfy all constraints (scenario 2), the DSE selects the one (as the new solution) that achieves the lowest objective value with a lower

constraints budget, i.e., the smallest value for *objective* × *constraints budget*. Such a strategy can help avoid greedy optimization that chases marginal objective reduction and rather seeks more promising solutions.

4.7 Bottleneck Mitigations for DNN Accelerators

We use the latency of executing a DNN as an example cost for a bottleneck model of DNN accelerator/mapping code designs. We describe what information about latency can be analyze and how to predict parameters that mitigate various bottlenecks.

Information embedded in bottleneck model: The bottleneck model incorporates execution characteristics of an optimized mapping of a DNN layer onto an architecture design. They include:

- *T_{comp}* *T_{comm}*, *T_{dma}*: Total time consumed by computations on PEs, communicating data via NoCs, and accessing data from off-chip memory via DMA, respectively.
- *Accel_freq*: Frequency of the accelerator (MHz)
- *data_offchip*: Data (bytes) accessed from off-chip, per operand
- *data_noc*: Data (bytes) communicated via NoC, per operand
- *NoC_groups_needed*: Maximum number of concurrent links that can be provided for communicating unique data to different PE-groups; one variable per operand.
- *NoC_bytes_per_group*: Size of the data that can be broadcast to PEs within every PE-group; one variable per operand.

Using above information, a bottleneck graph can be created as illustrated in Fig. 8. Typically, this information is available from experts-defined cost models such as [14, 36, 68]. If not, it may be obtained through similar analysis, hardware counters, or ML models.

Dictionary of Affected Parameters: It contains different factors contributing to the latency as keys and a list of relevant parameters as values. For example, the computation time is affected by the number of PEs and functional units in PEs. The time consumed by NoC communication is affected by the concurrent unicast links in NoCs, bit-widths of NoCs, and size of the local buffer or RF. The buffer size impacts the exploited reuse and the size of the data to be communicated. DMA time is affected by the bandwidth for off-chip memory accesses and the size of the shared memory.

Determining Values of Accelerator Design Parameters: Analyzing the bottleneck graph of a cost provides *s*, which is the scaling to be achieved by reducing a bottleneck factor's cost. *X_{current}* and *X_{new}* indicates the current and predicted value of a parameter *X*, respectively. *X* is a parameter impacting the bottleneck factor (obtained from dictionary). We next describe the calculation for various design parameters.

- **PEs:** The number of PEs required can be calculated directly from the needed speedup. $PEs_{new} = s * PEs_{current}$.
- **Off-chip BW:** Bandwidth (BW) for off-chip and on-chip communication is obtained from the number of data elements communicated per operand and targeted speedup. E.g.,
 $scaled_T_{dma} = T_{dma} \div s$;
 $footprint = sum(data_offchip)$;
 $bytes_per_cycle = footprint \div scaled_T_{dma}$
 $offchip_BW_{new} = bytes_per_cycle * Accelerator_freq$
- **NoC Links and Bit-width:** For DNN accelerators, separate NoCs communicate different operands, each with multiple concurrent

links for various PE groups. For every NoC, the maximum number of PE-groups with simultaneous access and the total bytes broadcast to each group are obtained from the cost model [14]. If communication time is a bottleneck, the operand causing it ('*op*') is available from the bottleneck analysis of the graph. Then, for the corresponding NoC, its width (bits) is scaled to make the broadcast faster based on the needed speedup. The new value is clamped to avoid exceeding the maximum width feasible for a one-shot broadcast.

$$max_width_feasible = exec_info[noc_bytes_per_group][op] * 8$$

$$width_scaled = noc_width_current * s$$

$$noc_width_new = min(width_scaled, max_width_feasible)$$

Similarly, total unicast links needed by the NoC for *op* are calculated from required concurrent accesses by PE groups.

$$max_links_feasible = exec_info[noc_groups_needed][op]$$

$$lnk_scaled = noc_unicast_links_current[op] * s$$

$$unicast_links_new[op] = min(lnk_scaled, max_links_feasible)$$

Whenever the number of PE-groups requiring different data elements exceeds available unicast links (by *V*×), data is unicast with time-sharing (*V* times) over configurable NoC (as in Eyeriss [8]) to facilitate mapping. Parameter *virtual_unicast_links* indicates time-sharing over a unicast link, which can be set as number of time-sharing instances (*V*).

- **Sizing RFs and Memory:** The total NoC communication time can be reduced by increasing the bottleneck operand (*op*)'s reuse in the RF (local buffer) of PEs. Increasing the reuse by *R* requires (*R*×) larger chunks of non-bottleneck operands, which need to be stored in RF and communicated via other NoCs. Using the information about non-exploited (available) reuse of the bottleneck operand and the required speedup, the new RF size can be calculated as:

$$target_scaling = min(max_reuse_available_RF[op], S)$$

$$RF_size_new = \sum_{op_i} [exec_info[data_RF][op_i] * target_scaling \div reuse_available_RF[op_i]]$$

The calculation is similar for global scratchpad memory, except for targeted scaling. In off-chip data communication, multiple operands are communicated one by one via DMA (unlike simultaneously by NoCs per operand). So, the targeted speedup depends on the bottleneck operand's (with remaining reuse) contribution (*f*) to the total off-chip footprint. The speedup achievable through reuse (*A*) can be approximated with Amdahl's law as:

$$A = (s * f) \div (1 - s + (s * f))$$

$$target_scaling = min(max_reuse_available_SPM[op], A)$$

$$SPM_size_new = \sum_{op_i} [exec_info[data_SPM][op_i] * target_scaling \div reuse_available_SPM[op_i]]$$

We implemented Explainable-DSE workflow and bottleneck analysis for DNN accelerators in python. It allows easy interfacing with DNN accelerator cost models. Since the implementation of the bottleneck analysis module and multi-bottleneck DSE is external to the cost model, they could be extended to interface with other cost models like MAESTRO that make execution characteristics available (e.g., bandwidth, Ops, data packets to be communicated). §C and §D discuss such specification efforts. We plan to open-source framework.

4.8 Tightly Coupled Hardware/Software Co-Explorations

Efficient codesign requires optimizing both hardware configurations and mappings in a coordinated manner. However, when using

back-box DSEs, these configurations are typically explored in a loosely coupled manner, as in the acquired values usually do not address inefficiencies in the achieved execution with their counterpart. For example, the acquired values of off-chip/NoC bandwidth may be inefficient for the selected loop tile configuration in the same/previous trials, resulting in significantly higher communication time and total latency.

To address these inefficiencies, our DSE integrates mapping space optimizations for DNN executions, and it explores HW/SW codesign in a tightly coupled manner through our bottleneck-based exploration. It considers software optimization as a subspace, which allows tailoring hardware configurations for obtained software configurations and optimizing software configurations to utilize hardware resources effectively. For a hardware configuration, when the DSE optimizes mappings through explorations or even a fixed schema, it mostly leads to efficient executions that can adapt to the tensor shapes and workload characteristics (reuse, batching, parallelism, etc.). Then, the DSE finds bottlenecks in the optimized executions obtained. In the next acquisition attempt, the DSE acquires new hardware candidates such that they address bottlenecks in the executions optimized previously by software configurations. Once a new hardware design is updated as the solution, software configurations are optimized again in tandem. Consequently, this approach leads to an efficient codesign for diverse tensor shapes and workload characteristics.

To enable efficient exploration of hardware/mapping codesign within practical budgets, DSE needs to explore quality mappings quickly. Our approach builds on previous research on mappers for DNN accelerators that eliminate infeasible and ineffective mappings by pruning loop tilings and orderings (detailed in §7, §F). For fast mapping optimizations, we have integrated and extended dMazeRunner [14], which can find near-optimal solutions within seconds. Mappers like dMazeRunner [14], Interstellar [71], or ZigZag [41] consider comprehensive space, optimally prune loop orderings, and prune tilings based on utilization of architectural resources (PEs, buffers, non-contiguous memory accesses). However, one challenge with their fixed utilization thresholds for pruning is that they may lead to a search space that either contains too few mappings (e.g., tens) for some DNN layers or too many (many thousands) for others. To address this challenge, we automatically adjust these search hyperparameters of dMazeRunner to formulate the mapping search space that contains up to the top- N mappings based on utilization thresholds. N is the size of pruned mapping space formulated by iteratively adjusted thresholds, which must be within a user-specified range, such as [10, 10000]. These mapping trials are then evaluated linearly, as in dMazeRunner [14] or Timeloop [45]. This approach helps achieve quality mappings by pruning ineffectual methods like in dMazeRunner/Interstellar, while also ensuring reasonably large space of high-quality mappings as per user-specified exploration budget.

5 EXPERIMENTAL METHODOLOGY

• **Benchmarks:** We evaluate eleven DNNs for Computer Vision (CV) and Natural Language Processing (NLP) tasks [51]. CV models include ResNet18, MobileNetV2, and EfficientNetB0 [60] (light)

Table 1: Design Space for Edge DNN Accelerators.

Data: int16; Freq. 500 MHz; Constraints: Throughput $\geq 40/10$ FPS (vision light/heavy), 120/530/176k samples/second (NLP: Transformer/BERT/wav2vec2); Area $< 75 \text{ mm}^2$; Max. power $< 4\text{W}$.

Objective: Minimize latency.

Parameter	Values	Options
PEs	64, 128, ..., 4096	7
L1 buffer (B)	8, 16, ..., 1024	8
L2 buffer (kB)	64, 128, ..., 4096	7
Offchip bandwidth (MBPS)	1024, 2048, 4096, 6400, 8192, 12800, 19200, 25600, 38400, 51200	10
NOC datawidth	16^*i ; i : [1, 16]	16
Physical unicast ($\times 4$)	$PEs^*i / 64$; i : [1, 64]	64^4
Virtual unicast ($\times 4$)	2^{3i} ; i : [0, 3]	4^4

and VGG16, ResNet50, and Vision Transformer (heavy) for classifying ImageNet images. The light and heavy labels differentiate models based on inference latency and total computations. For object detection, we evaluated recent models FasterRCNN-MobileNetV3 [27] and YOLOv5 [4] (heavy). NLP models include Transformer for English-German sentence translation and BERT-base-uncased [18] for Q&A on SQuAD dataset. We also evaluated Facebook wav2vec 2.0. [5] for ASR. Their DNN layers are 18, 53, 82, 16, 54, 86, 79, 60, 163, 85, and 109 respectively. We obtained models from PyTorch and HuggingFace frameworks.

• **Design Space:** Table 1 lists the design space of a DNN accelerator for inference at the edge. As in existing accelerators, we considered four dedicated NoCs for a total of four read/write operands [12]. The number of links for concurrent or time-shared unicasting is per NoC. To minimize the space for related techniques, we considered expressing the number of unicast links as a fraction of total PEs. We selected execution constraints based on the requirements for ML benchmarks [51] and designs of industrial edge accelerators for ML inference e.g., [1, 2]. We set the objective as minimizing the latency of the single-stream execution [51].

• **DSE Techniques:** We evaluated Explainable-DSE against previous accelerator DSE frameworks using constrained optimizations - Hypermapper 2.0 [43] and Confucius [30] - reinforcement learning (RL). Confucius limits the total parameters to two, works with a single constraint, and requires the same number of values for all parameters. So, we generalized its implementation for evaluations. We also evaluated our approach against non-feedback or black-box approaches like Grid search, Random search, Simulated annealing (Scipy [64]), Genetic algorithm (Scikit-Opt [3]), and Bayesian optimization [44]. We evaluated all techniques on a Dell precision 5820 tower workstation. Like previous DNN accelerator DSEs, we used validated cost models [14, 68]. The system for evaluating the candidates with cost models was same for all techniques.

• **Mapping Optimizations and Codesign Explorations:** Prior works mostly used a fixed dataflow, such that exploration time is primarily spent on optimizing hardware configurations, while getting efficient mappings with fixed strategy. So, we first fixed the mapping technique as an optimized output stationary dataflow (SOC-MOP) [7] for all approaches. Then, we demonstrate the codesign with Explainable-DSE by a tightly coupled optimization of both the hardware and mapping configurations. We also compare

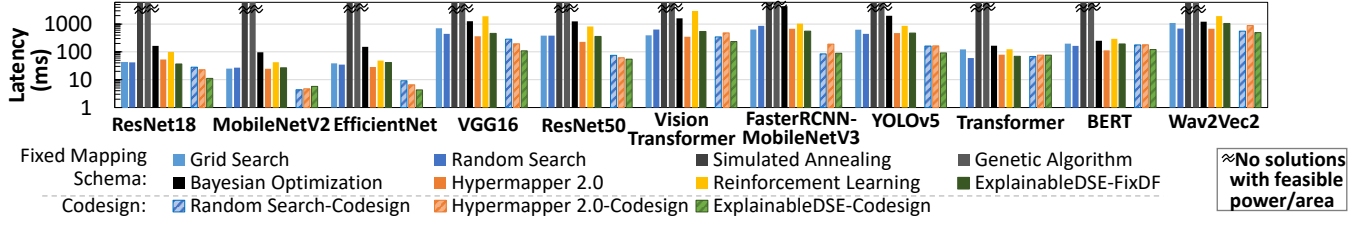


Figure 9: Explainable-DSE obtained codesigns of 6× lower latency.

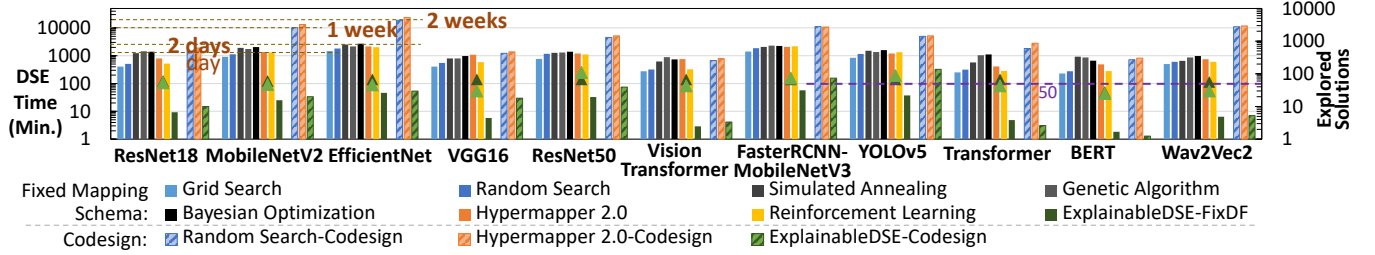


Figure 10: Explainable-DSE with fixed dataflow and codesigns reduce search time by 53× and 103× (minutes vs. days-weeks).

obtained codesigns with those obtained by black-box approaches. Black-box codesign DSE explores hardware configurations with two techniques that were found effective: random search and HyperMapper 2.0. §F details setup for effective black-box exploration of mappings in a comprehensive yet highly pruned space of feasible/effectual mappings. For mapping each DNN layer on every hardware configuration, black-box DSE uses Timeloop-like random search for 10,000 mapping trials, as it was found effective in quickly obtaining high-quality mappings (§F).

- **Exploration Budget:** We consider 2500 iterations for statically finding the best solutions. We also analyze dynamic DSE capabilities by explorations in 100 iterations.

6 RESULTS AND ANALYSIS

6.1 Explainable-DSE Obtained Codesigns of 6× Low Latency in 47× Less Iterations and 36× Less Time

Fig. 9 illustrates the latency obtained by different techniques for static exploration. By exploring among quality solutions, Explainable-DSE obtained 6× more efficient solutions, on average, as compared to previous approaches, and up to 9.6× over random search and 49.3× over Bayesian optimization. Even when dataflow (schema for optimized mappings) was fixed for all techniques, it obtained 1.77× lower latency on average and up to 7.89×. By applying bottleneck analysis on workload executions at every acquisition attempt, Explainable-DSE could determine parameters critical for improving efficiency. Thus, it can effectively navigate high-reward subspaces among the vast space. Fig. 11 illustrates this with latency reduction obtained over iterations by taking examples of two models, EfficientNet for CV and Transformer for NLP. With objective reduction at almost every attempt, the Explainable-DSE converges to quality solutions early on (some tens of iterations), and usually of better efficiency. E.g., obtained solutions have 6.6×–35.1× lower latency for EfficientNet, as compared to DSEs with fixed dataflow and 2.1×–9.7× as compared to black-box co-optimization. Overall,

at every attempt, it reduced the values of objective for feasible acquisitions by geomean 1.30× and 1.32× for fixed and co-explored mappings (as shown in Table 3). Acquisitions of non-explainable techniques, being bottlenecks-unaware, do not increasingly focus on high-reward sub-spaces. In fact, in some of our evaluations for Bayesian, random, constrained RL, overall growth was negative. For instance, many feasible candidates were acquired by these techniques without understanding bottlenecks, which corresponded to lower efficiencies than previously obtained best solutions.

Fig. 10 shows the total time (bars) taken by DSE techniques. Through constraints accommodation and systematically mitigating bottlenecks in multi-functional workload executions, explorations quickly converged or terminated while achieving even more efficient solutions. For example, Explainable-DSE with fixed and optimized mappings explored about only 59 and 54 designs, respectively (shown by triangles; ~2500 for other techniques). This translated in search time reduction of 53× and 103× on average over black-box explorations, when using fixed dataflow for all techniques and hardware/mapping co-optimization, respectively. Maximum reduction in the search time was up to 501× and 1675×, respectively. Using modest information on mitigating bottlenecks, our DSEs consumed only 21 and 64 minutes, on average. In fact, they achieved the most efficient solutions for BERT under just two minutes!

6.2 Including Software Design Space in the Exploration Enables 4.24× Better Solutions

With the availability of exploration budget (by a drastic reduction in the search time), hardware/software codesigns can truly be enabled by optimizing both of them in a tightly coupled manner. codesigns obtained with Explainable-DSE reduced objective by 4.24× on average, as compared to using a single optimized mapping per DNN operator. The higher efficiency emanates from achieving better mappings tailored for processing various DNN layers (different functionality and tensor shapes of DNN operators) on the selected hardware configuration. They leverage higher spatial

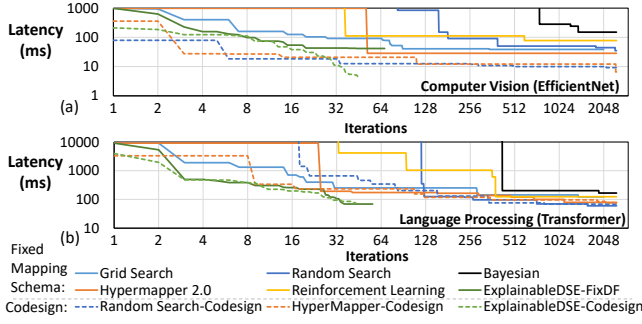


Figure 11: Latency reduced over iterations for (a) EfficientNet; (b) Transformer.

parallelism and more effectively hide data communication latency behind computations, as compared to a pre-set dataflow. Further, mapping optimizations reduce the objective considerably, without necessarily increasing hardware resources. Thus, by having a more constraints-budget on hand, DSE reduced the objective further (also evident in Fig. 11a).

For exploring comprehensively defined vast space of architectural configurations with non-explainable DSEs, presetting dataflow can lead to many infeasible solutions (§6.3). Note that infeasible solutions are not just hardware configurations with exceeding constraints like area or power. The designs can also be infeasible when generated hardware configuration is incompatible with the used software, i.e., dataflow for mapping. For instance, in configurations generated by non-explainable DSEs, the total number of links for time-shared unicast was often lower than that needed by spatial parallelism in the dataflow used for mapping. That’s exactly why a codesign or joint exploration with the software is important.

Black-box co-optimizations incorporated mapping explorations and reduced latency of obtained solutions further by 2.33× for HyperMapper 2.0 and 2.63× for random search, as compared to their DSEs using a fixed schema for optimized mappings. This is primarily because of the availability of more constraints-budget at hand, as discussed before. The co-optimizations also alleviated aforementioned challenge of mapping-hardware incompatibility. As Fig. 12 shows, the black-box co-optimizations find more feasible designs, when hardware design configurations are explored for the same number of trials. However, even after 2500 trials for exploring hardware configurations and 10,000 trials for exploring mappings of each DNN layer on every hardware configuration, the latency of codesigns obtained by black-box approaches are still 1.6× higher than the codesigns obtained by Explainable-DSE, while consuming 103× more search time (Taking 7-16 days for four workloads). Key reasons for such effective explorations by Explainable-DSE include generating fewer yet objective-reducing trials and tightly coupled codesigns. As Explainable-DSE leverages the domain knowledge, its generated designs target addressing execution inefficiencies, converging in 47× less iterations. In black-box co-optimizations, the DSE is loosely coupled, as in generated hardware configuration is not necessarily tailored to work best with the optimized mappings (from the previous/same trial for the hardware DSE). In contrast, tightly coupled codesign in Explainable-DSE explores hardware

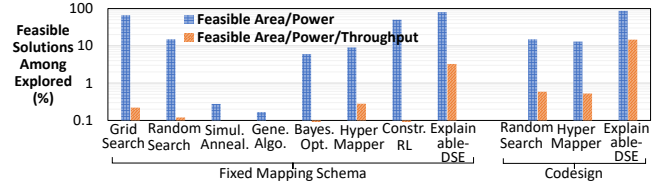


Figure 12: Most acquisitions by Explainable-DSE met area and power constraints, as compared to non-explainable techniques. Solutions obtained by all but Explainable-DSE mostly did not meet strict throughput requirements.

configurations that alleviate inefficiencies in the workload executions optimized previously by mappings; once new hardware configuration is generated, mapping exploration strives to utilize hardware resources effectively and lowers costs further. And, this repeats. Thus, optimizations for both hardware and software configurations strive to reduce inefficiencies in the execution optimized by their counterpart.

Although optimizing the mappings for every hardware design requires additional search time, the overall increase for exploring codesigns with Explainable-DSE was only 3× on average (from 21 minutes to 64). In fact, for all but heavy object detection models, the DSE time increased from 16 minutes to only 26 minutes. One reason is that the mappings can be quickly evaluated with analytical performance models (e.g., a minute each for several hundred to a few thousand mappings) and concurrent execution with multiple threads [14] (subjected to execution on 4 cores at maximum in our evaluations). Moreover, applying bottleneck analysis on efficient mappings helped obtain efficient designs faster (1.1× lower iterations for hardware designs on average, and up to 1.9×). Whenever the DSE for codesign evaluated a similar number of architecture designs as our DSE with fixed dataflow, it went on to explore even more efficient solutions (e.g., 2.33× lower latency for Vision Transformer).

6.3 By Considering the Utilization of Constraints, DSE Mostly Focuses on Feasible Solutions

Non-explainable black-box optimization approaches, e.g., with Genetic Algorithm or Bayesian Optimization, did not know which configurations could likely lead to feasible subspaces. Therefore, even after exploration over days, they almost did not obtain a single feasible solution. When considering only area and power constraints, feasibility of explored solutions was higher for mostly all techniques (Fig. 12), e.g., 15% for the random search and 50% for constraints-aware reinforcement learning. However, when considering throughput requirement for inference, the feasibility of the explored solutions was barely ~0.1%–0.3%. By exploring the mappings, the black-box codesign optimizations addressed the challenge of mappings being incompatible for the obtained hardware configurations. Thus, they improved feasibility by 2×–5×, but the overall feasibility was still ~0.6%. Such low feasibility for DSE in humongous space is presumably caused due to DSEs not accommodating constraints during exploration and bottlenecks-unaware acquisition trials. Contrarily, Explainable-DSE prioritized to meet

the constraints for its acquisitions and update of the best solutions, which helped avoid infeasible subspaces. Plus, addressing bottlenecks in executions helped acquiring high-performance solutions. Hence, 87% and 15% of solutions explored by Explainable-DSE code-signs were feasible when considering area and power constraints and all the three constraints, respectively. For DNNs like BERT and MobileNetV2, 89%–98% of the explored solutions met area and power constraints. Once Explainable-DSE achieved a solution that met all constraints, it always ensured to optimize further with a feasible solution.

6.4 Enabling Efficient Dynamic Exploration in Vast Space

Table 2 shows latency of solutions achieved in 100 iterations by different techniques. Under a short exploration budget, non-explainable techniques did not find a feasible solution (shaded values). Even after ignoring throughput requirements, most techniques could not find feasible solutions. Contrarily, by exploring spaces where candidates utilize low budget of constraints, Explainable-DSE quickly landed feasible solutions. Black-box approaches explored feasible codesigns, but they did not meet throughput requirements. On the other hand, by addressing the bottlenecks in multi-functional executions, Explainable-DSE achieved solutions of one to two orders of magnitude lower latency over other techniques.

7 ADDITIONAL RELATED WORKS

In this section, we discuss additional related work beyond the background on DNN accelerator DSE techniques described in §A.2, their limitations in §2, and previous DSEs using bottleneck analysis for different domains in §3.

• **Execution Cost Models of DNN Accelerators:** The cost models of SECDA [23] and TVM/VTA [6] support end-to-end simulation and synthesis, while faster analytical models are more commonly used to optimize mappings and accelerator design configurations. Their examples include MAESTRO [36], Accelergy [68], SCALE-Sim [53], and those of Timeloop [45], dMazeRunner [14], and Interstellar [71] infrastructures. Most of these models estimate both latency/throughput and energy. In addition to computational cycles, MAESTRO, dMazeRunner, and Timeloop account for on-chip and off-chip communication latency. Table 5 compares their execution modeling features. For the DSE, we used the cost model of our dMazeRunner infrastructure, which also considers the performance overheads of non-contiguous memory accesses and allows explicit specification of NoC bandwidths and flexibly specifying mappings through loop nest configurations.

• **Mappers for DNN Accelerators:** Mappers typically target the space of all valid loop tilings and orderings. For tensor shapes of a layer, there can be many factors of loop iteration counts, and just populating the space of valid mappings could be time-consuming (*microseconds*–*several seconds*) [52]. Table 6 compares different mappers. Timeloop [45], a commonly used mapper, explores mappings through random sampling, while GAMMA [31] uses a genetic algorithm. However, GAMMA limits the number of loops that can be executed spatially and does not prune invalid tilings before exploration, requiring several-fold more trials for convergence [32]. Without eliminating ineffectual loop tilings and orderings beforehand, black-box explorations typically require thousands of trials,

generating many invalid mappings, and take hours to map a single DNN layer once [29]. Mind Mappings [25] reduces the search time by training a surrogate model that estimates costs faster than analytical models. CoSA [29] uses a prime factorization-based approach to construct the tiling space for a mixed-integer programming solver. But, many tilings corresponding to combinations of prime factors remain unexplored, potentially resulting in sub-optimal solutions. Additionally, most mappers do not support depthwise-convolutions, invoking convolutions channel-by-channel. So, they miss opportunities for exploiting parallelism across multiple channels and reducing miss penalties for accessing contiguous data of consecutive channels from the off-chip memory.

Interstellar [71] prunes ineffectual tilings by constraining the search to pre-set resource utilization thresholds. dMazeRunner [14] goes further and prunes loop orderings for unique/maximum reuse of operands and proposes heuristics that reduce the space to highly efficient mappings, which can be explored in second(s). Hence, we utilize our dMazeRunner infrastructure in our codesign and extend it to construct the space of up to top- N mappings, where N is the maximum mapping trials allowed. ZigZag [41] and follow-up mappers build upon such pruning strategies. ZigZag allows uneven blockings of loops for processing different tensors, which may partially improve efficiency. However, ZigZag’s search time for a DNN layer is nearly hours [41]. While works such as [37, 41, 62, 75, 76] optimize DNN mappings on one or more hardware accelerators, they require exploring hardware parameters exhaustively or with black-box optimizations.

• **Hardware/Software Codesign Explorations of DNN Accelerators:** Previous DNN-accelerator DSEs, such as [38, 50, 65, 69, 74], used black-box optimizations. They incur excessive trials and ineffectual solutions, as they lack reasoning about the higher costs of obtained candidates and the potential efficiency of candidates to be acquired next (§2). Further, DSEs of [10, 30, 34, 50, 55, 70, 77] used a fixed dataflow in explorations. It obviates increasing search time further but may not lead to the most efficient solutions compared to codesigns.

Recent approaches HASCO [69] and DiGamma [33] optimize both hardware and mapping configurations in a black-box manner, encountering the same challenges of ineffectual and excessive trials due to non-explainability (§2). Second, with a loosely coupled codesign exploration (§4.8), they acquire HW/SW configurations that may not be effective or suitable for the counterpart. Furthermore, they target a limited hardware design space comprising only buffers and PEs. Finally, they typically do not explore a single accelerator design that addresses inefficiencies in executing DNNs with many layers.

• **DSE Using Bottleneck Analysis:** DSEs of [20, 72] use bottleneck analysis, but they are constraints-unaware and optimize only a single loop-kernel. Plus, they explored only neighboring values of parameters (instead of scaling them to mitigate bottleneck in one shot). It leads to search time comparable to black-box DSEs [72]. AutoDSE [58] and SECDA [23] proposed bottleneck models specific to FPGA-based HLS and their search optimizes a single loop-kernel/task of a single workload at a time. We propose bottleneck models for DNN accelerator domain; our DSE framework generalizes prior bottleneck-based DSEs to the case of multiple

Table 2: Latency minimized by DSE Techniques in 100 iterations.

Explainable-DSE evaluated ~54 solutions. Designs obtained by Non-Explainable DSEs were *low-throughput* (shaded values) and *incompatible* with dataflow (dashes). More importantly, * denotes none of the obtained candidates met even area/power constraints.

DSE Technique	ResNet18	MobileNetv2	EfficientNet	VGG16	ResNet50	Vision Transformer	FasterRCNN-MobileNetv3	YOLOv5	Transformer	BERT	Wav2Vec2
Grid Search-FixDF	278	73.4	92.0	3650	747	1973	1625	1477	251	780	1933
Random Search-FixDF	-*	197	694	41912	626	1376	3152	7754	157	1044	2357
Simulated Annealing-FixDF	-*	-*	-*	-*	-*	-*	-*	-*	-*	-*	-*
Genetic Algorithm-FixDF	-*	-*	-*	-	-*	-	-	-*	-*	-*	-*
Bayesian Optimization-FixDF	-	-	-	-	-	-	-	-	-	-	-
HyperMapper 2.0-FixDF	53.3	46.5	135	1339	493	1308	13582	1142	171	663	912
Reinforcement Learning-FixDF	-	-	360	-	-	-	21150	18082	143	1428	1428
Random Search-Codesign	69.6	12.7	9.5	870	209	857	224	218	244	240	1427
HyperMapper 2.0-Codesign	63.1	5.1	10.3	1233	87.3	1084	830	348	133	637	1945
ExplainableDSE-Codesign	11.2	5.7	4.3	109	54.9	233	89.2	92.1	76.2	121	494

loop-nests and multiple workloads through aggregation of bottleneck mitigations. Further, via proposed API and data structures, our framework decouples bottleneck models from search algorithms, allowing designers to express their bottleneck models.

8 CONCLUSIONS

Agile and efficient exploration in the vast design space, e.g., for hardware/software codesigns of DNN accelerators, require techniques that not just should consider objectives and constraints but are also explainable. They need to reason about obtained costs for acquired solutions and how to improve underlying execution inefficiencies. Non-explainable DSE with black-box optimizations (evolutionary, ML-based) lack such capability; obtaining efficient solutions even after thousands of trials or days can be challenging. To overcome such challenges, we proposed Explainable-DSE, which analyzes execution through bottleneck models and determines the bottleneck factors behind obtained costs and acquire solutions based on relevant mitigation strategies. Our demonstration of optimizing codesigns of DNN accelerators showed how Explainable-DSE could effectively explore feasible and efficient candidates (6× low-latency solutions). By obtaining most efficient solutions in short exploration budgets (47× fewer iterations or minutes/hours vs. days/weeks), it opens up cost-effective and dynamic exploration opportunities.

ACKNOWLEDGEMENTS

We thank anonymous reviewers for their valuable feedback and suggestions. This work was partially supported by National Science Foundation (NSF) grant CPS 1645578 and Graduate College Fellowship at Arizona State University. This work is done in part for Artificial Intelligence Hardware (AIHW) program of Semiconductor Research Corporation (SRC).

REFERENCES

- [1] 2018. ARM Machine Learning Processor. https://en.wikichip.org/wiki/arm_holdings/microarchitectures/mlp.
- [2] 2019. Intel Nervana NNP-I 100. https://en.wikichip.org/wiki/nervana/nnp/nnp-i_1100.
- [3] 2019. scikit-opt. github.com/guofei9987/scikit-opt/.
- [4] 2019. YOLOv5 Classification. https://pytorch.org/hub/ultralytics_yolov5.
- [5] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems* 33 (2020), 12449–12460.
- [6] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Q. Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018. {TVM}: An Automated {End-to-End} Optimizing Compiler for Deep Learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*. 578–594.
- [7] Yu-Hsin Chen, Joel Emer, and Vivienne Sze. 2016. Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*. 367–379. <https://doi.org/10.1109/ISCA.2016.40>
- [8] Yu-Hsin Chen, Tushar Krishna, Joel S. Emer, and Vivienne Sze. 2017. Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks. *IEEE Journal of Solid-State Circuits* 52, 1 (2017), 127–138. <https://doi.org/10.1109/JSSC.2016.2616357>
- [9] Yu-Hsin Chen, Tien-Ju Yang, Joel Emer, and Vivienne Sze. 2019. Eyeriss v2: A Flexible Accelerator for Emerging Deep Neural Networks on Mobile Devices. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* 9, 2 (2019), 292–308. <https://doi.org/10.1109/JETCAS.2019.2910232>
- [10] Kanghyun Choi, Deokki Hong, Hojae Yoon, Joonsang Yu, Youngsok Kim, and Jinho Lee. 2021. Dance: Differentiable accelerator/network co-exploration. In *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 337–342.
- [11] Coral. [n. d.]. Edge TPU Performance Benchmarks. <https://coral.ai/docs/edgetpu/benchmarks/>.
- [12] Shail Dave, Riyadh Baghdadi, Tony Nowatzki, Sasikanth Avancha, Aviral Shrivastava, and Baoxin Li. 2021. Hardware Acceleration of Sparse and Irregular Tensor Computations of ML Models: A Survey and Insights. *Proc. IEEE* 109, 10 (2021), 1706–1752. <https://doi.org/10.1109/JPROC.2021.3098483>
- [13] Shail Dave, Mahesh Balasubramanian, and Aviral Shrivastava. 2018. RAMP: Resource-Aware Mapping for CGRAs. In *2018 55th ACM/ESDA/IEEE Design Automation Conference (DAC)*. 1–6. <https://doi.org/10.1109/DAC.2018.8465892>
- [14] Shail Dave, Youngbin Kim, Sasikanth Avancha, Kyoungwoo Lee, and Aviral Shrivastava. 2019. DMazerunner: Executing perfectly nested loops on dataflow accelerators. *ACM Transactions on Embedded Computing Systems (TECS)* 18, 5s (2019), 1–27.
- [15] Shail Dave, Alberto Marchisio, Muhammad Abdullah Hanif, Amira Guesmi, Aviral Shrivastava, Ihcen Alouani, and Muhammad Shafique. 2022. Special Session: Towards an Agile Design Methodology for Efficient, Reliable, and Secure ML Systems. In *2022 IEEE 40th VLSI Test Symposium (VTS)*. IEEE, 1–14.
- [16] Shail Dave and Aviral Shrivastava. 2022. Design Space Description Language for Automated and Comprehensive Exploration of Next-Gen Hardware Accelerators. in *Workshop on Languages, Tools, and Techniques for Accelerator Design (LATTE'22)* (2022). co-located with the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2022).
- [17] Shail Dave, Aviral Shrivastava, Youngbin Kim, Sasikanth Avancha, and Kyoungwoo Lee. 2020. dMazeRunner: Optimizing Convolutions on Dataflow Accelerators. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1544–1548.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. [n. d.]. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.
- [19] Brian A Fields, Rastislav Bodik, Mark D Hill, and Chris J Newburn. 2003. Using interaction costs for microarchitectural bottleneck analysis. In *Proceedings. 36th Annual IEEE/ACM International Symposium on Microarchitecture, 2003. MICRO-36*. IEEE, 228–239.
- [20] Gennette Gill and Montek Singh. 2009. Bottleneck analysis and alleviation in pipelined systems: A fast hierarchical approach. In *2009 15th IEEE Symposium on Asynchronous Circuits and Systems*. IEEE, 195–205.

- [21] Soonhoi Ha, Jürgen Teich, Christian Haubelt, Michael Glaß, Tulika Mitra, Rainer Dömer, Petru Eles, Aviral Shrivastava, Andreas Gerstlauer, and Shuvra S Bhat-tacharyya. 2017. Introduction to hardware/software codesign. *Handbook of Hardware/Software Codesign* (2017), 3–26.
- [22] Di Han, Wei Chen, Bo Bai, and Yuguang Fang. 2019. Offloading optimization and bottleneck analysis for mobile cloud computing. *IEEE Transactions on Communications* 67, 9 (2019), 6153–6167.
- [23] Jude Haris, Perry Gibson, José Cano, Nicolas Bohm Agostini, and David Kaeli. 2021. SECD: Efficient Hardware/Software Co-Design of FPGA-based DNN Accelerators for Edge Inference. In *2021 IEEE 33rd International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*. IEEE, 33–43.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [25] Kartik Hegde, Po-An Tsai, Sitao Huang, Vikas Chandra, Angshuman Parashar, and Christopher W Fletcher. 2021. Mind mappings: enabling efficient algorithm-accelerator mapping space search. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*.
- [26] Christian Heidorn, Frank Hannig, and Jürgen Teich. 2020. Design space exploration for layer-parallel execution of convolutional neural networks on CGRAs. In *Proceedings of the 23th International Workshop on Software and Compilers for Embedded Systems*. 26–31.
- [27] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. [n. d.]. Searching for MobileNetV3. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, 1314–1324.
- [28] Qijing Huang, Charles Hong, John Wawrzyniec, Mahesh Subedar, and Yakun Sophia Shao. 2022. Learning A Continuous and Reconstructible Latent Space for Hardware Accelerator Design. In *2022 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 277–287.
- [29] Qijing Huang, Aravind Kalaiah, Minwoo Kang, James Demmel, Grace Dinh, John Wawrzyniec, Thomas Norell, and Yakun Sophia Shao. 2021. Cosa: Scheduling by constrained optimization for spatial accelerators. In *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 554–566.
- [30] Sheng-Chun Kao, Geonhwa Jeong, and Tushar Krishna. 2020. Confucius: Autonomous hardware resource assignment for dnn accelerators using reinforcement learning. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 622–636.
- [31] Sheng-Chun Kao and Tushar Krishna. 2020. Gamma: Automating the hw mapping of dnn models on accelerators via genetic algorithm. In *Proceedings of the 39th International Conference on Computer-Aided Design*. 1–9.
- [32] Sheng-Chun Kao, Angshuman Parashar, Po-An Tsai, and Tushar Krishna. 2022. Demystifying Map Space Exploration for NPUs. *arXiv preprint arXiv:2210.03731* (2022).
- [33] Sheng-Chun Kao, Michael Pellauer, Angshuman Parashar, and Tushar Krishna. 2022. Digamma: Domain-aware genetic algorithm for hw-mapping co-optimization for dnn accelerators. In *2022 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 232–237.
- [34] Liu Ke, Xin He, and Xuan Zhang. 2018. Nnest: Early-stage design space exploration tool for neural network inference accelerators. In *Proceedings of the International Symposium on Low Power Electronics and Design*. 1–6.
- [35] David Koeplinger, Christina Delimitrou, Raghu Prabhakar, Christos Kozyrakis, Yaqi Zhang, and Kunle Olukotun. 2016. Automatic Generation of Efficient Accelerators for Reconfigurable Hardware. In *Proceedings of the 43rd International Symposium on Computer Architecture*. IEEE Press, 115–127. <https://doi.org/10.1109/ISCA.2016.20>
- [36] Hyoukjun Kwon, Prasanth Chatarasi, Michael Pellauer, Angshuman Parashar, Vivek Sarkar, and Tushar Krishna. 2019. 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*. 754–768.
- [37] Hyoukjun Kwon, Liangzhen Lai, Michael Pellauer, Tushar Krishna, Yu-Hsin Chen, and Vikas Chandra. 2021. Heterogeneous dataflow accelerators for multi-DNN workloads. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 71–83.
- [38] Yujun Lin, Mengtian Yang, and Song Han. 2021. NAAS: Neural accelerator architecture search. In *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 1051–1056.
- [39] Google LLC. [n. d.]. Coral Edge TPU Accelerator. <https://coral.ai/products/accelerator-module>.
- [40] Atefeh Mehrabi, Aninda Manocha, Benjamin C Lee, and Daniel J Sorin. 2020. Prospector: Synthesizing efficient accelerators via statistical learning. In *2020 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 151–156.
- [41] Linyan Mei, Pouya Houshmand, Vikram Jain, Sebastian Giraldo, and Marian Verhelst. 2021. ZigZag: Enlarging joint architecture-mapping design space exploration for DNN accelerators. *IEEE Trans. Comput.* 70, 8 (2021), 1160–1174.
- [42] Naveen Muralimanohar, Rajeev Balasubramanian, and Norman P Jouppi. 2009. CACTI 6.0: A tool to model large caches. *HP laboratories* 27 (2009), 28.
- [43] Luigi Nardi, David Koeplinger, and Kunle Olukotun. 2019. Practical design space exploration. In *2019 IEEE 27th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*. IEEE, 347–358.
- [44] Fernando Nogueira. 2014–. Bayesian Optimization: Open source constrained global optimization tool for Python. <https://github.com/fmfn/BayesianOptimization>
- [45] Angshuman Parashar, Priyanka Raina, Yakun Sophia Shao, Yu-Hsin Chen, Victor A Ying, Anurag Mukkara, Rangharajan Venkatesan, Bruce Khailany, Stephen W Keckler, and Joel Emer. 2019. Timeloop: A systematic approach to dnn accelerator evaluation. In *2019 IEEE international symposium on performance analysis of systems and software (ISPASS)*. IEEE, 304–315.
- [46] Maryam Parsa, Aayush Ankit, Amirkoushyar Ziabari, and Kaushik Roy. 2019. Pabo: Pseudo agent-based multi-objective bayesian hyperparameter optimization for efficient neural accelerator design. In *2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 1–8.
- [47] Andy D. Pimentel. 2017. Exploring Exploration: A Tutorial Introduction to Embedded Systems Design Space Exploration. *IEEE Design & Test* 34, 1 (2017), 77–90. <https://doi.org/10.1109/MDAT.2016.2626445>
- [48] David L Poole and Alan K Mackworth. 2010. *Artificial Intelligence: foundations of computational agents*. Cambridge University Press.
- [49] Nirmal Prajapati, Sanjay Rajopadhye, Hristo Djidjev, Nandakishore Santhi, Tobias Grosser, and Rumen Andonov. 2019. Optimization Approach to Accelerator Codesign. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 39, 6 (2019), 1300–1313.
- [50] Brandon Reagen, José Miguel Hernández-Lobato, Robert Adolf, Michael Gelbart, Paul Whatmough, Gu-Yeon Wei, and David Brooks. 2017. A case for efficient accelerator design space exploration via bayesian optimization. In *2017 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*. IEEE, 1–6.
- [51] Vijay Janapa Reddi, Christine Cheng, David Kanter, Peter Mattson, Guenther Schmuelling, Carole-Jean Wu, Brian Anderson, Maximilien Breughe, Mark Charlebois, William Chou, Ramesh Chukka, Cody Coleman, Sam Davis, Pan Deng, Greg Damos, Jared Duke, Dave Fick, J. Scott Gardner, Itay Hubara, Sachin Idgunji, Thomas B. Jablin, Jeff Jiao, Tom St. John, Pankaj Kanwar, David Lee, Jeffery Liao, Anton Lokhmotov, Francisco Massa, Peng Meng, Paulius Mickevicius, Colin Osborne, Gennady Pekhimenko, Arun Tejusve Raghunath Rajan, Dilip Sequeira, Ashish Sirasao, Fei Sun, Hanlin Tang, Michael Thomson, Frank Wei, Ephrem Wu, Lingjie Xu, Koichi Yamada, Bing Yu, George Yuan, Aaron Zhong, Peizhao Zhang, and Yuchen Zhou. 2020. Mlperf inference benchmark. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. 446–459.
- [52] Enrico Russo, Maurizio Palesi, Davide Patti, Salvatore Monteleone, Giuseppe Ascia, and Vincenzo Catania. 2022. Multi-Objective End-to-End Design Space Exploration of Parameterized DNN Accelerators. *IEEE Internet of Things Journal* (2022).
- [53] Ananda Samajdar, Yuhao Zhu, Paul Whatmough, Matthew Mattina, and Tushar Krishna. 2018. Scale-sim: Systolic cnn accelerator simulator. *arXiv preprint arXiv:1811.02883* (2018).
- [54] Roberto Santana. 2017. Gray-box optimization and factorized distribution algorithms: where two worlds collide. *arXiv preprint arXiv:1707.03093* (2017).
- [55] Giulia Santoro, Mario R Casu, Valentino Peluso, Andrea Calimera, and Massimo Alioto. 2018. Energy-performance design exploration of a low-power microprogrammed deep-learning accelerator. In *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 1151–1154.
- [56] Kiran Seshadri, Berkin Akin, James Laudon, Ravi Narayanaswami, and Amir Yazdanbakhsh. 2022. An Evaluation of Edge TPU Accelerators for Convolutional Neural Networks. In *2022 IEEE International Symposium on Workload Characterization (IISWC)*. IEEE, 79–91.
- [57] Yakun Sophia Shao, Brandon Reagen, Gu-Yeon Wei, and David Brooks. 2014. Aladdin: A pre-rtl, power-performance accelerator simulator enabling large design space exploration of customized architectures. In *2014 ACM/IEEE 41st International Symposium on Computer Architecture (ISCA)*. IEEE, 97–108.
- [58] Atefeh Sohrabizadeh, Cody Hao Yu, Min Gao, and Jason Cong. 2022. AutoDSE: Enabling Software Programmers to Design Efficient FPGA Accelerators. *ACM Transactions on Design Automation of Electronic Systems (TODAES)* 27, 4 (2022), 1–27.
- [59] Naveen Suda, Vikas Chandra, Ganesh Dasika, Abinash Mohanty, Yufei Ma, Sarma Vrudhula, Jae-sun Seo, and Yu Cao. 2016. Throughput-Optimized OpenCL-Based FPGA Accelerator for Large-Scale Convolutional Neural Networks. In *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. Association for Computing Machinery, 16–25. <https://doi.org/10.1145/2847263.2847276>
- [60] Mingxing Tan and Quoc Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*. PMLR, 6105–6114.

- [61] Catia Trubiani, Antinisca Di Marco, Vittorio Cortellessa, Nariman Mani, and Dorina Petriu. 2014. Exploring synergies between bottleneck analysis and performance antipatterns. In *Proceedings of the 5th ACM/SPEC International Conference on Performance engineering*. 75–86.
- [62] Miheer Vaidya, Aravind Sukumaran-Rajam, Atanas Rountev, and P Sadayappan. [n. d.]. Comprehensive accelerator-dataflow co-design optimization for convolutional neural networks. In *2022 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 325–335.
- [63] Stylianos I Venieris and Christos-Savvas Bouganis. 2016. fpgaConvNet: A framework for mapping convolutional neural networks on FPGAs. In *2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 40–47.
- [64] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, Ilhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy. 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature methods* 17, 3 (2020), 261–272.
- [65] Jie Wang and Jason Cong. 2021. Search for Optimal Systolic Arrays: A Comprehensive Automated Exploration Framework and Lessons Learned. *arXiv preprint arXiv:2111.14252* (2021).
- [66] Jian Weng, Sihao Liu, Vidushi Dadu, Zhengrong Wang, Preyas Shah, and Tony Nowatzki. 2020. Dsagen: Synthesizing programmable spatial accelerators. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 268–281.
- [67] Nan Wu, Yuan Xie, and Cong Hao. 2021. Ironman: Gnn-assisted design space exploration in high-level synthesis via reinforcement learning. In *Proceedings of the 2021 on Great Lakes Symposium on VLSI*. 39–44.
- [68] Yannan Nellie Wu, Joel S. Emer, and Vivienne Sze. 2019. Accelergy: An Architecture-Level Energy Estimation Methodology for Accelerator Designs. In *2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. 1–8. <https://doi.org/10.1109/ICCAD45719.2019.8942149>
- [69] Qingcheng Xiao, Size Zheng, Bingzhe Wu, Pengcheng Xu, Xuehai Qian, and Yun Liang. 2021. Hasco: Towards agile hardware and software co-design for tensor computation. In *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 1055–1068.
- [70] Lei Yang, Zheyu Yan, Meng Li, Hyoukjun Kwon, Liangzhen Lai, Tushar Krishna, Vikas Chandra, Weiwen Jiang, and Yiyu Shi. 2020. Co-Exploration of Neural Architectures and Heterogeneous ASIC Accelerator Designs Targeting Multiple Tasks. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. 1–6. <https://doi.org/10.1109/DAC18072.2020.9218676>
- [71] Xuan Yang, Mingyu Gao, Qiaoyi Liu, Jeff Setter, Jing Pu, Ankita Nayak, Steven Bell, Kaidi Cao, Heonjae Ha, Priyanka Raina, Christos Kozyrakis, and Mark Horowitz. 2020. 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*. 369–383.
- [72] Yang Yang, Marc Geilen, Twan Basten, Sander Stuijk, and Henk Corporaal. [n. d.]. Automated bottleneck-driven design-space exploration of media processing systems. In *2010 Design, Automation & Test in Europe Conference & Exhibition (DATE 2010)*. 1041–1046.
- [73] Ye Yu, Yingmin Li, Shuai Che, Niraj K Jha, and Weifeng Zhang. 2020. Software-defined design space exploration for an efficient dnn accelerator architecture. *IEEE Trans. Comput.* 70, 1 (2020), 45–56.
- [74] Dan Zhang, Safeen Huda, Ebrahim Songhori, Kartik Prabhu, Quoc Le, Anna Goldie, and Azalia Mirhoseini. 2022. A full-stack search technique for domain optimized deep learning accelerators. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*. 27–42.
- [75] Xiaofan Zhang, Yuan Ma, Jinjun Xiong, Wen-Mei W. Hwu, Volodymyr Kindratenko, and Deming Chen. 2022. Exploring HW/SW Co-Design for Video Analysis on CPU-FPGA Heterogeneous Systems. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 41, 6 (2022), 1606–1619.
- [76] Shixuan Zheng, Xianjue Zhang, Leibo Liu, Shaojun Wei, and Shouyi Yin. [n. d.]. Atomic Dataflow based Graph-Level Workload Orchestration for Scalable DNN Accelerators. In *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. 475–489.
- [77] Yanqi Zhou, Xuanyi Dong, Tianjian Meng, Mingxing Tan, Berkin Akin, Daiyi Peng, Amir Yazdanbakhsh, Da Huang, Ravi Narayanaswami, and James Laudon. 2022. Towards the Co-design of Neural Networks and Accelerators. *Proceedings of Machine Learning and Systems* 4 (2022), 141–152.

A ACCELERATOR HARDWARE/SOFTWARE DSE

A.1 DSE Problem Formulation and Terminology

Exploration of accelerator designs is a *constrained minimization* problem, where the most efficient *solution*⁵ corresponds to minimized *objective* (e.g., latency), subjected to *inequality constraints* on some *costs* (e.g., area, power) and *parameters* p of accelerator design [47]. Fig. 13 illustrates the DSE problem formulation. During the optimization, every solution gets evaluated by *cost models* for objectives and inequality constraints. The DSE technique needs to consider only *feasible* solutions and determine the most efficient solution by processing several *iterations*. It is a discrete optimization since the *search space* is usually confined to presumably effective solutions, e.g., power-of-two or categorical values of parameters. It is also *derivative-free optimization* [43].

$$\begin{aligned} \min \quad & \text{obj}(p), \quad p = (p_1, p_2, \dots, p_n) \in \mathbb{R}^n \\ \text{subject to} \quad & \text{cost}_i(p) \leq \text{constraint}_i; \quad \text{for } i = 1, 2, \dots, m \end{aligned}$$

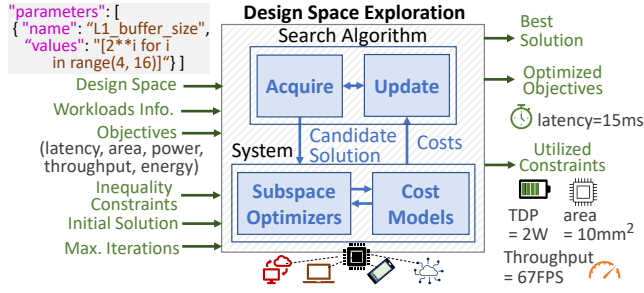


Figure 13: Accelerator DSE As Constrained Minimization.

A.2 DNN Accelerator Hardware/Software Codesigns DSE

Hardware/software codesigns can be explored by partitioning the search space and optimizing *software* space as a *subspace* in a loop. So, the DSE technique needs to find the best mapping of a task onto architecture and repeat the search with different architectural configurations [21]. Partitioning enables exploration in reduced space compared to exploring parameters from multiple spaces altogether. DSE techniques for DNN accelerators explore hardware designs through non-feedback or black-box optimizations like evolutionary or ML-based [10, 28, 30, 46, 50, 59, 65, 69, 70, 73, 74, 77]. For mapping DNNs on a design (subspace optimization), they typically fix the way of execution or dataflow, e.g., in [10, 30, 34, 50, 55, 70, 72]. Hence, for processing each functionality (nested loop such as a DNN layer), these techniques usually have just one mapping. Thus, they primarily optimize designs of accelerator architecture, i.e., parameters for buffers, processing elements (PEs), and NoCs.

B CAPABILITIES AND DISTINGUISHED FEATURES

In this section, we highlight the capabilities of the proposed DSE approach for agile and explainable explorations.

⁵In the accelerator design space exploration context, we use terms "solutions", "designs", and "configurations" interchangeably.

- **Efficient designs.** Explainable-DSE finds better solutions since it investigates costs and bottlenecks that incur higher costs; by exploring candidates that can mitigate inefficiencies in obtained designs, DSE provides efficient designs.
- **Quick DSE.** The DSE can reduce objective values at almost every acquisition attempt; it searches mostly in feasible/effectual solution spaces. Thus, DSE achieves efficient solutions quickly, which is beneficial for the early design phase and for dynamic DSEs, e.g., deployments of accelerator overlays at run time. Additionally, it can help when acquisition budgets are limited, e.g., due to evaluation of a solution consuming minutes to hours [74]. Further, when designers optimize designs offline with hybrid optimization methodologies [30] comprising multiple optimizations, quickly found efficient solutions can serve as high-quality initial points.
- **Explainability in the DSE and design process.** This work shows the need for explainability in the design process, e.g., in exploring the vast design space of deep learning accelerators, and how DSE driven by bottleneck models can achieve explainability. Exploration based on bottleneck analysis can help explain why designs perform well/poorly and which regions are well-explored/unexplored in vast space and why.
- **Generalized bottleneck-driven DSE for multiple microbenchmarks and workloads.** In acquiring new candidates, the DSE accounts for various bottlenecks in executing multiple loop nests (e.g., DNN layers) of diverse execution characteristics. Thus, the DSE can provide a single solution that is most effective overall, in contrast to previous DSEs that provide loop-kernel-specific solutions.
- **Specification for expressing domain-specific bottleneck models to the DSE.** This work proposes an API for expressing domain-specific bottleneck models so that the designers can integrate them to bottleneck-driven DSE frameworks and reuse the DSE.
- **Comprehensive design space specification.** In the DSE, appropriate values of a parameter is selected through bottleneck models. Thus, the DSE can alleviate the need for fine-tuning the design space; users can comprehensively define/explore vast space, e.g., more parameters and large ranges of values (arbitrary instead of power-of-two).
- **Bottleneck analysis for hardware/software codesign of deep learning accelerators.** By taking the latency of accelerators as an example, this work shows how to construct bottleneck models for designing deep learning accelerators and bottleneck analysis for improving the accelerator designs based on their execution characteristics.

C CURRENT LIMITATIONS AND FUTURE WORK

- **Improving efficiency further through better acquisitions:** By using bottleneck models and considering budgeted constraints in the explorations, the DSE can find outperforming solutions. However, it may converge to a suboptimal solution (e.g., for VGG-16) due to the greed for resolving the bottlenecks in the optimized designs. This can be mitigated

Table 3: At every acquisition attempt, Explainable-DSE reduces objective by 30% vs. ~1.4% by non-explainable techniques. N/A is indicated when a technique could not find a single feasible hardware solution.

DSE Technique	ResNet18	MobileNetv2	EfficientNet	VGG16	ResNet50	Vision Transformer	FasterRCNN-MobileNetv3	YOLOv5	Transformer	BERT	Wav2Vec2	Average
Grid Search-FixDF	1.71%	1.03%	1.07%	1.21%	1.25%	1.41%	0.71%	0.55%	0.98%	1.04%	1.07%	1.09%
Random Search-FixDF	0.52%	-0.87%	7.34%	-2.26%	4.69%	-4.29%	-1.41%	-0.90%	0.01%	0.97%	-1.45%	0.21%
Simulated Annealing-FixDF	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Genetic Algorithm-FixDF	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Bayesian Optimization-FixDF	11.26%	26.57%	19.57%	19.22%	-1.09%	-4.89%	-10.01%	-12.28%	10.15%	-0.27%	11.33%	6.32%
HyperMapper 2.0	5.32%	1.21%	0.44%	2.67%	4.94%	4.86%	-0.20%	0.87%	1.40%	3.35%	1.18%	2.37%
Reinforcement Learning-FixDF	-0.75%	-4.13%	5.18%	-2.51%	-2.97%	-10.47%	0.67%	1.66%	4.62%	0.50%	-0.50%	-0.79%
Random Search-Codesign	-0.07%	-0.29%	0.14%	0.44%	0.33%	0.57%	0.23%	0.91%	0.48%	-0.24%	0.02%	0.23%
HyperMapper 2.0-Codesign	0.56%	0.46%	0.59%	0.64%	0.68%	0.68%	0.72%	0.76%	0.43%	0.52%	0.73%	0.62%
ExplainableDSE-FixDF	53.54%	21.92%	20.48%	52.42%	15.32%	31.74%	23.73%	21.54%	30.96%	40.66%	21.44%	30.34%
ExplainableDSE-Codesign	30.50%	23.45%	32.10%	32.03%	18.77%	46.29%	27.03%	18.78%	26.19%	47.30%	46.70%	31.74%

by making the acquisitions more "exploratory", e.g., exploring distant high-reward sub-spaces. This can be achieved by exploring multiple spaces side-by-side by targeting a pool of various initial points [43]. Alternatively, the acquisition function can incorporate both self-supervised learning and bottleneck/constraint considerations.

- **Implications of specifying ineffectual bottleneck analysis for an arbitrary domain-specific system:** In Explainable-DSE, the bottleneck model helps explain design cost and guide the DSE. It considers which factors constitute overall cost and how design parameters and their scaling could impact each factor. For different domain-specific systems, designers usually characterize target systems and develop domain-specific bottleneck models/analyses [19, 20, 22, 58, 61, 72]. Even for designing deep learning accelerators, designers already develop cost models or closely work with them [14, 36, 45, 46, 50, 59, 63, 65, 69, 73, 74]. It enables the designers to determine bottleneck-affecting parameters and mitigation strategies. For instance, designers can obtain a bottleneck model by simplifying their full model for an analytical cost calculation, which they can provide along with the cost model. Alternately, designers could estimate bottleneck mitigation strategies through characterization or sensitivity analysis of design parameters. Designers can also opt for automation techniques, as discussed later below. A domain-specific bottleneck model constructed by domain experts or possibly learned from domain-specific data can lead to an effective DSE. However, in the absence of an effectual analysis for an arbitrary system, the resulting exploration may be slower or suboptimal. For instance, the DSE may require more acquisitions and be slower if the user scales the parameter values conservatively or associates irrelevant parameters with bottleneck factors (e.g., the total number of PEs to the DMA time). Conversely, the DSE may converge to suboptimal solutions if the user either skips associating a critical parameter with a bottleneck factor or scales values aggressively. The examples are a user not suggesting explorations of NOC links or bit-widths when on-chip communication time is the bottleneck and the DSE scaling the number of PEs by 2× or more, even if the required scaling was only 1.2×. Either can cause the DSE to miss a range of efficient and constraints-satisfying solutions.

- **Generalizing bottleneck mitigation strategies for arbitrary domain-specific systems:** For expert-defined cost models, such as those of DNN accelerators, manual bottleneck analysis is possible. In the case of accelerators described as flow graphs, the analysis of costs/bottlenecks may be automated by parsing execution information for flow graph components [15, 16, 66, 72]. However, for other systems, mitigation can be estimated using gray-box optimization functions that approximate the relevance and contributions of each parameter to the total cost [54] or surrogates [25].

Our dynamic DSE evaluation demonstrates the potential of incorporating explainability into the DSE. Efficiency and agility can be improved even further by improving the provided feedback (mitigation strategies), improving the decision mechanism that uses the feedback, and enabling parallel evaluation/exploration of candidates.

D SPECIFICATION NEEDED FOR THE DSES

Black-box optimization techniques such as random search, simulated annealing, or Bayesian optimization are easy to deploy and require minimal tuning and specification for meaningfully constraining the design space, which may take some hours. However, due to their non-explainable nature, they may mostly explore inefficient or infeasible solutions, while consuming significant search time for excessive sampling. Applying these black-box exploration approaches to a domain-specific design optimization problem may require some additional specification efforts [30, 31], depending on their implementation in prior works. For instance, Confucius [30] is a reinforcement learning-based DSE framework that uses an LSTM/MLP-based policy network. For the targeted evaluations, we extended it to allow an arbitrary number of parameters, a different environment (target setup), different numbers of possible options for different parameters (different list sizes), and an arbitrary number of constraints. We also improved the reward calculation to consider an arbitrary number of constraints. This extension required adding/modifying a few tens of lines of code (LoC) and a few days of work.

For bottleneck analysis-driven DSE, we developed a bottleneck model for the target domain of DNN accelerator's hardware/mapping codesign. This is similar to efforts made previously in other domains [19, 20, 22, 58, 72]. The bottleneck model was expressed to the DSE

framework via proposed API. It led to about tens of LoC that specified how different factors contribute to overall cost and a few to several lines that specified each parameter’s scaling. It is worth noting that this is significantly less code compared to domain-specific analytical cost models [14, 36], which typically require thousands of lines of code. The bottleneck model used in the DSE is generally simpler than the full analytical model, as it only considers major execution-related factors and incorporate simple estimates about bottleneck mitigation. Domain experts can derive the bottleneck model from either the analytical model they develop/use or the sensitivity of a cost to the design parameters, which may take several days. §C discusses how bottleneck mitigation may be generalized for arbitrary systems.

E CASE STUDY: EFFICIENCY OF THE DESIGNS ACHIEVED BY DSE

Methodology: In this study, we compare the efficiency of designs obtained by our DSE (§6; Fig. 9) to those of edge AI accelerators. The DSE can only explore designs for the spatial architecture templates used by the cost models (e.g., in [14, 36, 45, 68]). Therefore, to make a fair comparison, we only considered edge accelerators with similar architectural features. Specifically, we compare with two edge accelerators: 1) Coral Edge TPU [39], a state-of-the-art industrial edge accelerator platform developed by Google, and 2) Eyeriss [7], an efficient edge accelerator that incorporates several special optimizations for high energy efficiency and low latency.

While novel edge accelerators have been developed recently, they typically contain several (micro)architectural features for specialization (e.g., mixed-precision computations or sparsity exploiting features [12]) that are not accurately modeled by the commonly used/available latency models for edge AI accelerators. The considered architectures are suitable for the comparison, as they are commonly used as a template in the system-level tools for design and execution modeling [14, 25, 30, 36, 45, 71]. Table 4 compares the architectural features and execution optimizations for the Edge TPU, Eyeriss, and the template architecture used by the DSE cost models. We obtained the results for Edge TPU from the performance benchmarking of TPU-optimized models [11], and the results for Eyeriss from its evaluations [7].

Fig. 14 shows the performance achieved by all three designs and the resultant energy efficiency and area efficiency. The results demonstrate that, on average, the codesigns obtained by the DSE attain 3.7× higher throughput than the Edge TPU. This is due to the codesign DSE’s ability to analyze execution bottlenecks and optimize both mappings and hardware configurations such as NoC configurations/bandwidth and buffer sizes. The DSE-achieved designs also required more than an order of magnitude less on-chip area, presumably due to allocating significantly smaller buffers and fewer MAC units (e.g., considering Edge TPU configurations studied in [56]). The overall area efficiency is higher by about 14× for MobileNetV2 and 49× on average due to a high-throughput execution for VGG-16. Although the DSE was focused on minimizing latency, energy efficiency of its designs matches that of the Edge TPU-aware EfficientNet, and is even higher by 2.9× for VGG-16.

Table 4: Comparison of Architectural Features and Execution Optimization for Edge AI Acceleration.

Feature	Edge TPU	Eyeriss	DSE
Data Precision	8b*	16b	16b
Technology	-	65nm	45nm
PEs	Scalar MAC	Scalar MAC	Scalar MAC
Temporal Reuse	Yes	Yes	Yes
Spatial Reuse	Data Distribution and Reduction	Data Distribution and Reduction	Data Distribution
Mapping Optimizations	Automatic (TFLite)	Fixed-Style	Automatic
Hardware-Mapping Co-optimization	No	Yes	Yes
Accelerator-DNN Codesign	Yes†	No	No
Frequency (MHz)	125-500	100-250	500

* Edge TPU results are scaled to match 16b precision used in comparison.

† For Edge TPU, the EfficientNet-EdgeTPU model is codesigned.

‡ In lack of information about the actual power consumed by edge TPU for different models, 1.4W power is considered (as reported for MobileNetV2 in the edge TPU datasheet).

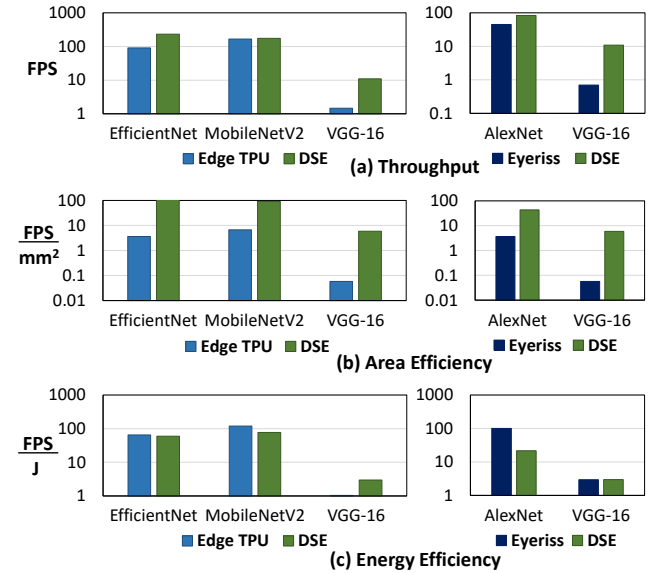


Figure 14: Comparison of efficiency of the designs obtained by the DSE with that of the edge accelerators Google Coral Edge TPU [39] and Eyeriss [8]: (a) Throughput (frames per second i.e., FPS) [3.7×, 8.7×], (b) Area efficiency (FPS/mm²) [49×, 57×], and (c) Energy efficiency (FPS/Joule) [1.5×, 0.6×].

When compared to Eyeriss, the DSE achieves designs of lower latency (with similar area/power budgets), and improves area efficiency by up to 11.84× and energy efficiency of VGG-16 by up to 1.33×—while not incorporating additional efficiency-gaining features as in Eyeriss. For instance, Eyeriss leverages frequency scaling, power gating, zero skipping, and compression (for up to 86% sparsity in AlexNet layers). Thus, Eyeriss achieves about 2×–3.75× higher energy efficiency when compared to the designs obtained by DSE that minimized latency. In most cases, the codesigns obtained

Table 5: Execution Cost Models for Deep Learning Accelerators.

Features	Timeloop [45]	dMazeRunner [14]	MAESTRO [36]	Interstellar [71]	SCALE-Sim [53]	Accelergy [68]
Performance Estimate	Yes	Yes	Yes	No	Yes	No
Energy Estimate	Yes	Yes	Yes	Yes	No	Yes
Integrated Support for ML Libraries	No	Yes	Yes	No	No	N/A
Layers Supported	GEMM, CONV	GEMM, CONV, DWCONV	GEMM, CONV, DWCONV	GEMM, CONV	GEMM, CONV	N/A
Data Precision	Variable	Variable	Variable	Variable	Fixed	Variable
Mapping Specification	Loop Nest Configuration	Loop Nest Configuration	Directives	Loop Nest Configuration	N/A	N/A
Dataflow	All	All	All	All	WS, OS, IS	N/A
Spatial and Temporal Data Reuse	Yes	Yes	Yes	Yes	Yes	N/A
Data Reuse with Striding Convolution	Yes	No	Yes	N/A	N/A	N/A
Considers On-Chip Communication Latency	Yes	Yes	Yes	No	N/A	N/A
Memory Hierarchy	Arbitrary	3-level	3-level	3/4-level	Fixed	Arbitrary
Models overhead of non-contiguous accesses	No	Yes	No	No	No	No

Table 6: Mappers for Deep Learning Accelerators.

Features	Timeloop [45]	Gamma [31]	Mind Mappings [25]	CoSA [29]	Interstellar [71]	ZigZag [41]	dMazeRunner [14]
Comprehensive Mapping Space	Yes	No	Yes	No	Yes	Yes	Yes
Discard Invalid Mappings for Mapping-Space Construction	Yes	No	No	Yes	Yes	Yes	Yes
Prune Inefficient Tilings	No	No	No	No	No	Yes	Yes
Prune Inefficient Orderings	No	No	No	No	Fixed	Yes	Yes
Layers Supported	CONV, GEMM	CONV, GEMM, DWCONV	CONV, MTTKRP	CONV, GEMM	CONV, GEMM	CONV, GEMM	CONV, GEMM, DWCONV
Exploration Approach	Random	Genetic Algorithm	Gradient Descent	Mixed Integer Programming	Heuristic	Heuristic	Heuristic
Uneven Tiling for Tensors	No	No	No	No	No	Yes	No
Search Time	Minutes	Minutes	Minutes	Seconds	Minutes	Hours	Seconds
Off-line Training	No	No	Yes	No	No	No	No

by DSE in a tightly coupled manner outperform the Eyeriss-like designs. These comparisons demonstrate the potential of the proposed capabilities for the accelerator system design. With the continued development of automated execution modeling and inclusion of more template architectures and specialization components, the methodology can be expected to enhance the efficiency further.

F CONSTRUCTING AN EFFECTUAL MAPPING SPACE AND BLACK-BOX MAPPING OPTIMIZATIONS

Background: The process of optimizing mappings of a DNN layer on an accelerator design involves exploring an enormous search space, as demonstrated by previous works [17, 31] and later summarized in Table 7. This search space is primarily composed of configurations for loop tile sizings and loop orderings. Loop tilings determine the parallel processing on PEs and the data bursts that need to be accessed via memory hierarchy and communicated, while loop orderings determine the data reuse and impact the total

memory accesses. For mapping DNN operators such as convolutions and multi-layer perceptrons on an accelerator with a 3-level buffer/memory hierarchy and 1-level spatial parallelism [8], the mapping space corresponds to 28-deep and 12-deep nested loops, respectively [14]. Tile sizings are configurations that represent the values for loop iterations at each temporal/spatial level. For a selected tiling configuration, processing of a 7-deep nested loop at a buffer level corresponds to 7! different loop orderings. Table 7 lists the DNNs and their layers with the largest search space, which can contain up to $O(10^{36})$ configurations.

Pruning Invalid Tilings: To optimize the tile sizes of a DNN layer using a black-box optimizer, designers often set the index variables of tiled loops as design parameters and specify lower/upper bounds (loop iterations) to define the range of values that the optimizer can explore [25, 31]. However, the search space for tile sizes can be enormous, containing as many as $O(10^{28})$ configurations for certain DNN layers, as shown in Table 7. A black-box optimizer is unlikely to find more than a few valid mappings under practical exploration budgets when working with such a large search space

Table 7: Analyzing Size of the Mapping Space for Deep Learning Accelerators.

Model	Layer	Tile Sizings	Tile Sizings with Valid Factors	Valid Tilings w.r.t. Hardware	Orderings at a Memory Level	Orderings with Unique/Max Data Reuse	Full Map. Space	Factorization-Constrained Mapping Space	Factorization-Constrained Reuse-Aware Map. Space
		A	B	C	D	E	F: A^2D^2	G: B^2D^2	H: B^2E^2
ResNet18	CONV_2_1a	$O(10^{25})$	$O(10^{13})$	$O(10^7)$	$O(10^4)$	15/3	$O(10^{32})$	$O(10^{20})$	$O(10^{14})$
MobileNetV2	features.2.conv.0	$O(10^{22})$	$O(10^{12})$	$O(10^6)$	$O(10^4)$	15/3	$O(10^{30})$	$O(10^{19})$	$O(10^{13})$
EfficientNetB0	blks.2.expand	$O(10^{22})$	$O(10^{12})$	$O(10^6)$	$O(10^4)$	15/3	$O(10^{29})$	$O(10^{20})$	$O(10^{13})$
VGG-16	CONV_1_2	$O(10^{28})$	$O(10^{14})$	$O(10^7)$	$O(10^4)$	15/3	$O(10^{36})$	$O(10^{21})$	$O(10^{15})$
ResNet50	CONV_2_1b	$O(10^{25})$	$O(10^{13})$	$O(10^7)$	$O(10^4)$	15/3	$O(10^{32})$	$O(10^{20})$	$O(10^{14})$
Vision Transformer	patchembeddings.CONV2D	$O(10^{25})$	$O(10^{13})$	$O(10^6)$	$O(10^4)$	15/3	$O(10^{32})$	$O(10^{20})$	$O(10^{14})$
FasterRCNN-MobileNetV3	features.12.conv2.excite	$O(10^{26})$	$O(10^{13})$	$O(10^6)$	$O(10^4)$	15/3	$O(10^{33})$	$O(10^{20})$	$O(10^{14})$
YOLOv5	features.1.conv	$O(10^{27})$	$O(10^{14})$	$O(10^7)$	$O(10^4)$	15/3	$O(10^{34})$	$O(10^{21})$	$O(10^{15})$
Transformer	decoder.output_projection	$O(10^{27})$	$O(10^9)$	$O(10^4)$	$O(10^1)$	3/3	$O(10^{28})$	$O(10^{10})$	$O(10^{10})$
BERT	encoder.layer.0.output.dense	$O(10^{26})$	$O(10^9)$	$O(10^5)$	$O(10^1)$	3/3	$O(10^{27})$	$O(10^{11})$	$O(10^{10})$
Wav2Vec2	encoder.layers.0.intermediate.dense	$O(10^{28})$	$O(10^{12})$	$O(10^6)$	$O(10^1)$	3/3	$O(10^{29})$	$O(10^{13})$	$O(10^{12})$

[29]. This is because, for each loop in a DNN operator, only a small subset of factors of loop iterations lead to valid tile sizes. For example, a set (8, 4, 2, 16) makes a valid 4-level tiling configuration for 1024 loop iterations over output activations/filters, while many more do not. Therefore, designers formulate the space of valid tiling sizes based on factorization of loop iterations [29, 45]. As shown in Table 7, this prunes the search space of tiling configurations by a square root or even a cube root, e.g., from $O(10^{22})$ – $O(10^{28})$ configurations to a much smaller range of $O(10^9)$ – $O(10^{14})$ configurations. Invalid tiling configurations, which require more architectural resources than are available in the target hardware configuration, are usually discarded by black-box optimizers during the exploration (as indicated in column C of Table 7).

Pruning Ineffectual Orderings: We further reduced the space by discarding ineffective loop orderings. Previous black-box mappers explored the search space of all orderings, resulting in a large number of configurations ($7!$ or $O(10^4)$) for processing a convolution at a memory level. Instead, our approach build upon previous research that has shown only a handful of loop-orderings having unique data reuse of tensors (15 for convolutions) and a few with maximum reuse of various tensors [14, 41].

Overall Mapping Space: The GAMMA-like mapper considers full (non-factorized) tiling space and all loop-orderings [31] (column F), while Timeloop [45] and CoSA [29] consider factorized tile sizes but all loop orderings. For black-box mapping optimizations in our evaluations, we used factorized tile sizes (all valid factorizations) and considered only loop orderings with unique data reuse. In practice, there are only a few unique data reuse scenarios (vs. 15/3) for each tiling configuration, so all of them can be explored linearly [14]. Therefore, we invoked black-box optimizations with 10,000 trials for mapping each layer on a hardware configuration. During each trial, the optimization acquired a tiling configuration and evaluated all effectual loop-orderings through the cost model, selecting the one that minimized the objective. Thus, the mapping space

formulation discarded a large number of invalid and ineffective configurations (column H) without compromising the optimality.

Selection of the Mapping Optimization Technique: The results of our experiments comparing black-box DSEs of hardware configurations with fixed mapping schema showed that random search and Bayesian optimization-based HyperMapper 2.0 [43] were the most effective, as depicted in Figure 9. As a result, we selected these two techniques for optimizing both the hardware and mapping configurations. However, when we applied them to optimize mappings from the pruned space (column H in Table 7), we found that random search obtained efficient mappings within several seconds. By contrast, HyperMapper 2.0 generated efficient mappings, but its search overhead was prohibitively high, requiring a few hours to evaluate just a single DNN layer. Consequently, we used random search for optimizing mappings during the codesign DSE.

We also compared the quality of mappings obtained by random search with those obtained by simulated annealing, genetic algorithm, and Bayesian optimization for ResNet18 layers, as shown in Fig. 15.⁶ The random search successfully achieved low-latency mappings for all layers, whereas simulated annealing [64] failed to map a few layers (in 10,000 trials), and genetic algorithm [3] led to a higher overall latency than random search and took almost four hours to optimize mappings for the nine layers, making it impractical for codesign. Therefore, we opted to use a Timeloop-like random search to quickly and efficiently explore the highly pruned mapping space.

G CODESIGN OPTIMIZATION

The optimization of hardware and software codesigns can be done either by exploring partitioned sub-spaces in a sequential manner or simultaneously. In a partitioned or a two-stage optimization, an outer loop iterates over different hardware configurations, and an

⁶The mappings were evaluated for the initial hardware configuration, corresponding to the lowest values of design parameters in Table 1.

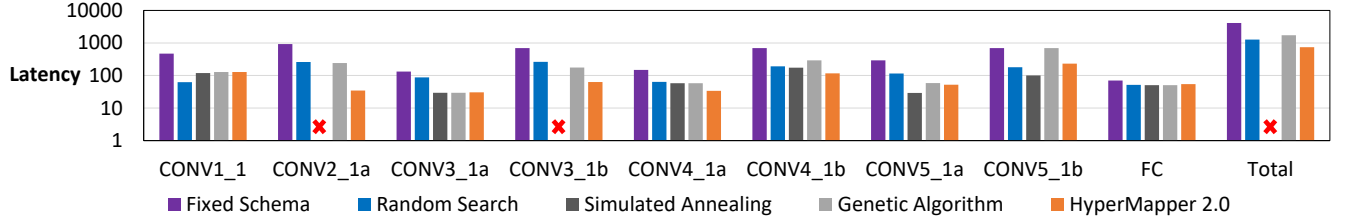


Figure 15: Efficiency of mappings obtained by different black-box optimizations.

inner loop optimizes the software space for each hardware configuration selected. On the other hand, the joint or simultaneous exploration involves finding a new configuration for both the hardware and software parameters at the same time in a trial. Although approaches using simultaneous search have been proposed, they are often infeasible to apply to a multi-workload exploration, target system with diverse and time-consuming cost evaluations, and huge collective search space. Therefore, partitioned sub-space exploration is commonly used for optimizing codesigns (§A.2). For demonstration of Explainable-DSE, all our DSE evaluations also follow two-stage optimization.

Firstly, approaches using simultaneous search [33, 43] typically optimize configurations for individual loop kernels such as a single DNN layer, as they optimize both the hardware and software parameters at every search attempt. This does not necessarily lead to a single accelerator design that is most efficient for the entire workload or a set of workloads, as layer-specific designs may not be optimal overall for the entire DNN or multiple DNNs.

Furthermore, optimizing both hardware and software parameters simultaneously can be very time-consuming. A target system often involves different cost functions or modules for different metrics that could consume different evaluation times. For example, evaluating area and power of each hardware configuration via Accelerger [68] could take a few seconds, whereas the cost models of dMazeRunner [14] or Timeloop [45] could estimate latency/energy for hundreds–thousands of mappings in a second. For exploring codesigns for a DNN with $L=50$ unique layers, consider a black-box DSE that is budgeted $H=2,500$ trials for hardware configurations and $M = 10,000$ trials for mapping each DNN layer on each hardware configuration. Simultaneous exploration of hardware and software configurations in $H \times M$ trials for each of the L layers requires the system to evaluate power/area costs for $H \times M \times L$ times, which would take more than 0.7 million hours or 79 years! In contrast, a two-stage partitioned exploration evaluates power/area costs only for H trials, and if the DSE samples infeasible mappings for a hardware configuration, they can be discarded promptly without further detailed evaluation. Our experiments show that the black-box DSEs obtained codesigns in a few days to a few weeks with the partitioned exploration approach.

Finally, in addition to the design parameters such as the total PEs or buffer sizes, hardware configurations can have various parameters, such as bandwidth, reconfiguration of NoCs (time-multiplexed communication of data packets, bus widths), and those for architectural specialization/heterogeneity, which further increase the

search space for both the hardware and software/mapping configurations. With the vast space for both the hardware and software/mapping configurations, the collective search space becomes huge, compounding the already challenging exploration of feasible and effective solutions for either of the hardware and software parameters. Additionally, in the DSE trials, simultaneously acquired hardware and software configurations may not be compatible with each other or may not mitigate execution inefficiencies corresponding to their counterpart.

Received 20 October 2022; revised 2 March 2023; accepted 27 April 2023