

MetaML: Automating Customizable Cross-Stage Design-Flow for Deep Learning Acceleration

Zhiqiang Que, Shuo Liu, Markus Rognlien, Ce Guo, Jose G. F. Coutinho, Wayne Luk
Department of Computing, Imperial College London, UK. {z.que, c.guo, jgfc, w.luk}@imperial.ac.uk

Abstract—This paper introduces a novel optimization framework for deep neural network (DNN) hardware accelerators, enabling the rapid development of customized and automated design flows. More specifically, our approach aims to automate the selection and configuration of low-level optimization techniques, encompassing DNN and FPGA low-level optimizations. We introduce novel optimization and transformation tasks for building design-flow architectures, which are highly customizable and flexible, thereby enhancing the performance and efficiency of DNN accelerators. Our results demonstrate considerable reductions of up to 92% in DSP usage and 89% in LUT usage for two networks, while maintaining accuracy and eliminating the need for human effort or domain expertise. In comparison to state-of-the-art approaches, our design achieves higher accuracy and utilizes three times fewer DSP resources, underscoring the advantages of our proposed framework.

I. INTRODUCTION

The field of deep learning has witnessed unprecedented growth in recent years, driven by the increasing demand for efficient and high-performance applications [1]. Consequently, FPGA-based deep neural network (DNN) accelerator design and optimization have gained significant attention [2, 3, 4]. The development of an efficient FPGA-based DNN design requires a diverse skill set that combines expertise in machine learning with low-level knowledge of the target hardware architecture [5]. Optimizing these DNN designs is a complex process, as it involves balancing competing objectives. On one hand, high accuracy during inference is crucial from an application perspective. On the other hand, the design must be optimized for the underlying hardware architecture, meeting power, latency and throughput requirements while fitting into the FPGA device [6, 7]. Achieving an optimal balance between these conflicting objectives requires careful consideration and effective optimization strategies [8].

In this paper, we focus on the key challenge of codifying optimization strategies that automate the selection and configuration of low-level optimization techniques covering multiple abstraction levels from DNN optimizations to FPGA optimizations. Achieving optimal results for a given problem is complex [9]. It is challenging to identify the most effective combination, order and tuning of optimization techniques to ensure optimal outcomes [4, 10]. In order to address this challenge, we propose the following contributions:

- 1) A co-optimization framework for FPGA-based DNN accelerators, which includes novel building tasks that enable

rapid development of customized design flows, automating the entire design iteration process (Section III).

- 2) A library of reusable optimization and transformation tasks designed to be customizable and flexible, and that can be easily integrated into our co-optimization framework. Some of the tasks in our library are specific to certain applications and/or target technologies, while others are agnostic, providing versatility and adaptability to the framework (Section IV).
- 3) The evaluation of the proposed framework using multiple benchmarks and different optimization strategies. This evaluation provides insights into the effectiveness of the framework and its optimization modules under different scenarios (Section V).

We present related work in Section II and conclude this paper and present future work in Section VI.

II. RELATED WORK

The field of FPGA-based DNN acceleration has rapidly expanded, leading to the development of numerous optimization techniques and tools. Several co-optimization techniques have been proposed that optimize both algorithm and hardware stages for DNNs on FPGAs, as discussed in various papers [11, 12, 13, 14, 15, 16, 17] and in hardware-aware neural architecture search studies like [18, 19]. These optimization strategies are often coupled with design space exploration, but are typically hardcoded and cannot be easily changed or customized. Other approaches offer end-to-end software frameworks, such as Xilinx’s Vitis AI [20] and Intel’s OpenVINO [21], which optimize DNNs with pre-built optimizations for deployment on specific target technologies. Moreover, frameworks, such as FINN [22], HLS4ML [23], and fpgaConvNet [24], provide optimized hardware building blocks for FPGA-based DNN accelerators. However, they do not support automated cross-stage optimization strategies.

Our framework addresses the limitations of current optimization techniques for deep neural networks (DNNs) and hardware by supporting the design of fully automated optimization flows. These flows can be described by reusable tasks that are specific to the target platform as well as ones that can be used across different platforms. Additionally, our approach seamlessly integrates both new and existing optimization techniques, catering to different levels of abstraction. For instance, it includes graph optimizations for neural networks and source-to-source optimizations for HLS C++ as described in [25].

III. OUR APPROACH

This paper introduces a novel framework that simplifies the development of customizable design flows for optimizing deep neural networks (DNNs) on FPGA platforms. A design flow consists of a series of tasks that aim to convert a high-level specification into a final hardware design. Typically, a design flow is implemented as a multi-stage pipeline, where each stage operates on a specific model abstraction. As the pipeline progresses, the model abstraction is gradually refined and optimized, taking into account the key features of the target device. For instance, to illustrate the process, let us consider a DNN model described in Tensorflow. The framework utilizes the HLS4ML tool to translate this model into a C++ HLS model in the initial stage. Subsequently, the resulting C++ HLS model can be further transformed into a Register Transfer Level (RTL) model using Vivado HLS.

Our approach enables the optimization strategies to be codified, automating the selection and configuration of DNN and hardware optimizations. It is designed to fulfill the following requirements:

Customizable: Users have the freedom to customize, extend, or modify the design-flow to meet specific needs and support experimentation. They can select a set of design-flow tasks, that implement specific optimizations or transformations, arrange them in a desired order, and fine-tune their parameters to create a specific optimization flow. Additionally, users can develop their own tasks and integrate them into the design-flow.

Cross-stage: This refers to the breadth of optimizations and the ability to target multiple levels of abstraction, typically associated with different stages of a design-flow. More specifically, optimizations performed at the neural network level operate at the graph-level, while optimizations at the HLS C++ level involve source-level transformations.

Using our framework, design flows are programmatically generated and consist of two types of components (Fig. 1): the pipe task and the meta-model. The pipe task serves as the basic unit of the design flow, executing specific optimizations or transformations. By interconnecting these tasks, we construct a complete design flow. The architecture of the design-flow is depicted as a cyclic directed graph where nodes symbolize tasks and edges signify dependencies between tasks, denoting the need to complete one task before initiating another.

The meta-model, on the other hand, serves as a shared space for storing the states of the design flow. This model consists of three sections: configuration, log, and model space. The configuration section (CFG) acts as a key-value store, holding the parameters of all pipe tasks in the design flow. The log section (LOG) records the runtime execution trace, aiding in debugging. Furthermore, the model space can store the generated models obtained during the execution of the design flow. In the example presented in Fig. 1, six models are stored, covering DNN, HLS C++, and RTL abstraction levels. Each model includes supporting files, tool reports, and computed metrics.

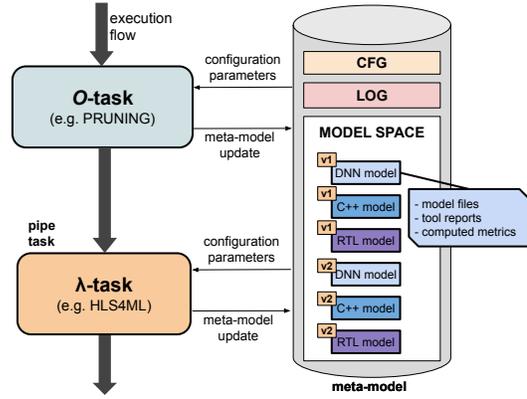


Fig. 1. A connection between an O -task and a λ -task. Each connection defines a unidirectional flow between a source and a target task.

IV. REUSABLE PIPE TASKS

Table I provides a list of pipe tasks that are currently implemented, along with their roles, multiplicity, and parameters. The multiplicity indicates the number of input and output connections that a task can handle. We consider two types of tasks:

- **O -task:** These are self-contained optimization tasks that enhance a given model based on specific objectives and constraints. Our current pipe task repository includes PRUNING and SCALING, which are implemented using the Keras API (version 2.9.0), and QUANTIZATION, which employs C++ source-to-source transformations via the Artisan framework [25];
- **λ -task:** These tasks perform functional transformations on the model space, such as compilation and synthesis. Examples include HLS4ML (version 0.6.0), which translates a DNN model into an HLS C++ model, and Vivado HLS (version 20.1), which translates an HLS C++ model into an RTL model.

Our framework is customizable. Different parameters (see Table I) can be used to customise a design flow, and new design-flows can be built from existing ones to tailor to specific needs.

TABLE I
A LIST OF IMPLEMENTED PIPE TASKS.

Type	Role	Multiplicity	Parameters
HLS4ML	λ	1-to-1	default_precision IOType FPGA_part_number clock_period test_dataset
VIVADO-HLS	λ	1-to-1	project_dir
KERAS-MODEL-GEN	λ	0-to-1	train_en train_test_dataset train_epochs
PRUNING	O	1-to-1	tolerate_acc_loss (α_p) pruning_rate_thresh (β_p) train_test_dataset train_epochs
SCALING	O	1-to-1	default_scale_factor tolerate_acc_loss (α_s) scale_auto max_trials_num train_test_dataset train_epochs
QUANTIZATION	O	1-to-1	tolerate_acc_loss (α_q) train_test_dataset

V. EVALUATION

In this section, we demonstrate how optimization strategies can be built by revising design-flow architectures, combining and reusing pipe tasks, and modifying their configuration.

A. Experimental Setup

Experiments were conducted in Python 3.9.15 with benchmark workloads from typical DNN applications, including jet identification [23, 26] (Jet-DNN), image classification using VGG7 [27] and ResNet9 [28] networks. The datasets used are: Jet-HLF [23], MNIST [29] and SVHN [30], respectively. The jet identification task targeted FPGA-based CERN Large Hadron Collider (LHC) triggers with a 40 MHz input rate and a response latency of less than 1 microsecond. Default frequencies were 100MHz for Zynq 7020 and 200MHz for Alveo U250 and VU9P, and the HLS4ML task used 18-bit fixed-point precision with 8 integer bits.

B. Optimization Strategies

In this subsection, we initially focus on three strategies, each supported by a single O -task. Then, we shift our focus towards combined strategies that utilize multiple O -tasks.

Pruning strategy. Pruning is a technique that improves the performance and efficiency of neural networks by removing insignificant weights. Our framework includes an implementation of this technique through the PRUNING O -task. A design-flow which employs the PRUNING O -task is illustrated in Fig. 2(a). This optimization task gradually zeroes out weights during training to create a more compact and efficient network while maintaining accuracy. In addition, it supports auto-pruning, which automatically determines the highest pruning rate while maintaining a given level of accuracy loss. Formally, the objective of this O -task is defined as:

$$\begin{aligned} & \text{maximum } Pruning_rate \\ & \text{subject to } Accuracy_loss(Pruning_rate) \leq \alpha_p \end{aligned} \quad (1)$$

Starting at 0% pruning rate, the auto-pruning algorithm obtains initial accuracy Acc_{p0} at step 1 (s1). It then uses a binary search approach, increasing or decreasing the pruning rate based on whether the accuracy loss is within a user-defined tolerance ($\leq \alpha_p$). The algorithm terminates when the rate difference is below a threshold (β_p). The number of steps is determined by $1 + \log_2(1/\beta_p)$. Algorithm search steps and direction are shown in Fig. 3. Both the tolerance ($\leq \alpha_p$) and threshold (β_p) values are initially set to 2% in this work.

The effectiveness of the auto-pruning algorithm is demonstrated in Fig. 4. Fig. 4(a) and (b) depict the pruning rate and accuracy for Jet-DNN and ResNet9 in each step, while Fig. 4(c) and (d) show the resources utilization. As the pruning rate increases, hardware resource requirements, particularly DSPs and LUTs, decrease, leading to improved FPGA performance. The design candidate with the highest pruning rate within the allowed tolerance is selected.

Scaling strategy. To accommodate a large DNN design on an FPGA, our framework supports the SCALING O -task which automatically reduces the layer size while tracking the

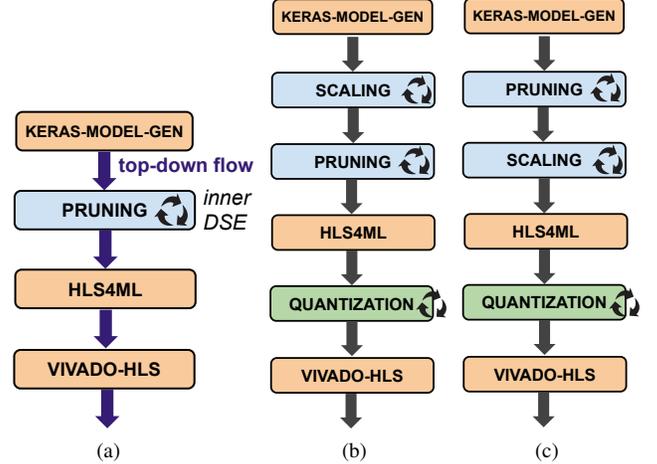


Fig. 2. (a) Pruning strategy. (b) The combined strategy of scaling, pruning and quantization. (c) The combined strategy with a different order of O -tasks.

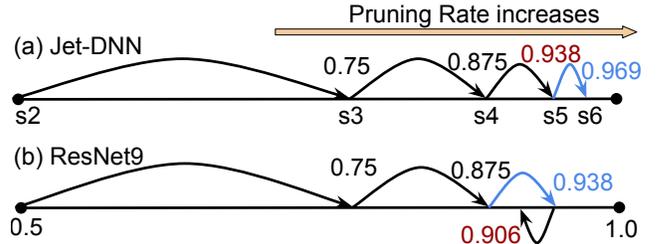


Fig. 3. The auto-pruning algorithm applied to models, (a) Jet-DNN and (b) ResNet9, with binary search direction shown. Omitting step s1 for visibility. The blue arrow indicates an accuracy loss > user threshold; red denotes the optimal pruning rate.

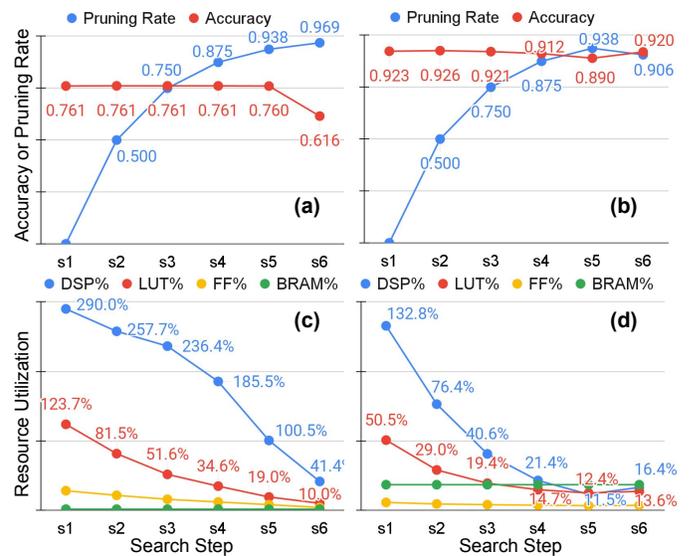


Fig. 4. (a) Pruning rates and accuracy of Jet-DNN. (b) Resource utilization of Jet-DNN design candidates with pruning on Xilinx Zynq 7020. (c) Pruning rates and accuracy of ResNet9. (d) Resource utilization of ResNet9 design candidates with pruning on Xilinx U250.

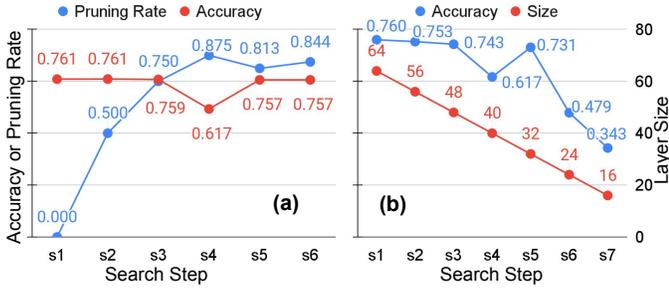


Fig. 5. (a) Jet-DNN accuracy and pruning rates with scaling then pruning. (b) Jet-DNN accuracy and layer size with pruning then scaling.

accuracy loss α_s . The search stops when the loss exceeds α_s . This parameter can be adjusted to achieve further size reduction with minimal impact on accuracy. This work sets α_s to 0.05%, which allows for model size reduction with negligible accuracy loss.

Quantization strategy. Our framework supports the QUANTIZATION O -task to automate mixed-precision quantization for networks. It operates at the HLS C++ level, providing more direct control over hardware optimizations and reducing unintended side effects when translating DNN models to HLS C++ using tools such as HLS4ML. The resulting precision configuration is directly instrumented into the C++ kernel, and a co-design simulation evaluates the accuracy of the quantized model. If the accuracy loss is within tolerance ($< \alpha_q$), this process is repeated. This work sets α_q to 1%.

Combining O -tasks. With our framework, new strategies can be derived by building and revising the design-flow architecture. For instance, by inserting the SCALING O -task before the PRUNING O -task in Fig. 2(a), a custom combined strategy is created with results shown in Fig. 5(a). The new optimal pruning rate is 84.4%, lower than the previous 93.8%, due to reduced redundancy from the preceding scaling task. By switching the order of the O -tasks, a different optimization is achieved, resulting in a 0.7% accuracy drop after one scaling step, as seen in Fig. 5(b). Moreover, the three optimization O -tasks, pruning, scaling, and quantization, can be integrated into a single automated cross-stage strategy to enhance both performance and hardware efficiency, as illustrated in Fig. 2(b) and (c).

Discussion and Comparison. Our evaluation results indicate that our combined O -task optimization strategy typically outperforms single O -task techniques. Furthermore, the order in which these optimization techniques are applied plays a crucial role, as different orders produce varying final results.

To highlight the advantages of our framework, we compared our design flow results to those of other studies targeting low-latency, low-resource, fully unfolded FPGA implementations, including original Jet-DNN [23] using HLS4ML, LogicNets [31], QKeras-based Q6 [6], and AutoQKeras-based QE and QB [6]. All designs use the same architecture, except for JSC-L, which employs a larger architecture. Our "S+P+Q" design with α_q set to 0.01 achieves an accuracy of 74.1% for JSC-S, outperforming JSC-M and JSC-L, which have

TABLE II
PERFORMANCE COMPARISON WITH THE FPGA DESIGNS OF JET-DNN NETWORK USING OTHER APPROACHES ON XILINX FPGAS

Model	α_q	FPGA	Acc. (%)	Lat. (ns)	Lat. (cycles)	DSP (%)	LUT (%)	Power (W)
HLS4ML Jet-DNN [23]	-	KU115	75	75	15	954 (17.3)	-	-
LogicNets JSC-M [31]	-	VU9P	70.6	NA	NA	0 (0)	14,428 (1.2)	-
LogicNets JSC-L [31]	-	VU9P	71.8	13 ^a	5	0 (0)	37,931 (3.2)	-
Qkeras Q6 [6]	-	VU9P	74.8	55	11	124 (1.8)	39,782 (3.4)	-
AutoQkeras QE [6]	-	VU9P	72.3	55	11	66 (1.0)	9,149 (0.8)	-
AutoQkeras QB [6]	-	VU9P	71.9	70	14	69 (1.0)	11,193 (0.9)	-
This work (same to [23])	1%	VU9P	76.1	70	14	638 (9.3)	69,751 (5.9)	2.51 ^b
This work S→P→Q	1%	VU9P	75.6	45	9	50 (0.7)	6,698 (0.6)	0.199 ^b
This work S→P→Q	4%	VU9P	72.8	40	8	23 (0.2)	7,224 (0.6)	0.166 ^b

^a A clock frequency of 384 MHz is used and the final softmax layer is removed.

^b Dynamic power reported by Vivado. The static power is about 2.5W for all these designs.

accuracies of 70.6% and 71.8%, respectively. Compared to Q6, our design has 0.8% higher accuracy, uses 2.5 times fewer DSP blocks, and 5.7 times fewer LUTs. Furthermore, our design outperforms both QE and QB, achieving over 3.3% higher accuracy and lower resource usage. Our design also boasts lower latency than all Q6, QE, and QB designs, demonstrating the benefits of our framework. Inspired by the AutoQKeras design QB, which minimizes model bit consumption, we further optimize our design by tuning the parameters, such as α_q . Increasing the quantization O -task's tolerant accuracy loss (α_q) from 1% to 4% results in a smaller model with DSP usage 3 times lower than AutoQkeras' most efficient model, QE (see Table II). Although the model accuracy decreases to 72.8%, it remains higher than the optimized AutoQkeras designs QB and QE.

VI. CONCLUSION

This paper presents a novel co-optimization framework for FPGA-based DNN accelerators, enabling efficient development of customized cross-stage design flows. Our approach reduces DSP resource usage by up to 92% and LUT usage by up to 89% while maintaining accuracy. Compared to existing research, our approach achieves higher model accuracy and uses 3 times fewer DSP resources, highlighting the benefits of our framework. Future work includes exploring more efficient search techniques with larger numbers of O -tasks, introducing control tasks to cover bottom-up and parallel flows, supporting more types of neural networks such as graph [32] and transformer [33] networks, and utilizing new FPGA resources such as AI Engines [34] and AI Tensor Blocks [35].

ACKNOWLEDGEMENT

The support of the United Kingdom EPSRC (grant numbers EP/V028251/1, EP/L016796/1, EP/N031768/1, EP/P010040/1, and EP/S030069/1), CERN, AMD and SRC is gratefully acknowledged.

REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *Proceedings of the 2015 ACM/SIGDA international symposium on field-programmable gate arrays*, 2015, pp. 161–170.
- [3] Z. Que, H. Nakahara, E. Nurvitadhi, A. Boutros, H. Fan, C. Zeng, J. Meng, K. H. Tsoi, X. Niu, and W. Luk, "Recurrent Neural Networks With Column-Wise Matrix-Vector Multiplication on FPGAs," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 30, no. 2, pp. 227–237, 2021.
- [4] X. Zhang, H. Ye, J. Wang, Y. Lin, J. Xiong, W.-m. Hwu, and D. Chen, "DNNExplorer: a framework for modeling and exploring a novel paradigm of FPGA-based DNN accelerator," in *Proceedings of the 39th International Conference on Computer-Aided Design*, 2020, pp. 1–9.
- [5] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [6] C. N. Coelho, A. Kuusela, S. Li, H. Zhuang, J. Ngadiuba, T. K. Aarrestad, V. Loncar, M. Pierini, A. A. Pol, and S. Summers, "Automatic heterogeneous quantization of deep neural networks for low-latency inference on the edge for particle detectors," *Nature Machine Intelligence*, vol. 3, no. 8, pp. 675–686, 2021.
- [7] Z. Que, E. Wang, U. Marikar, E. Moreno, J. Ngadiuba, H. Javed, B. Borzyszkowski, T. Aarrestad, V. Loncar, S. Summers *et al.*, "Accelerating recurrent neural networks for gravitational wave experiments," in *2021 IEEE 32nd International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, 2021, pp. 117–124.
- [8] E. Wang, J. J. Davis, R. Zhao, H.-C. Ng, X. Niu, W. Luk, P. Y. Cheung, and G. A. Constantinides, "Deep neural network approximation for custom hardware: Where we've been, where we're going," *ACM Computing Surveys (CSUR)*, vol. 52, no. 2, pp. 1–39, 2019.
- [9] L. Deng, G. Li, S. Han, L. Shi, and Y. Xie, "Model compression and hardware acceleration for neural networks: A comprehensive survey," *Proceedings of the IEEE*, vol. 108, no. 4, pp. 485–532, 2020.
- [10] P. Xu, X. Zhang, C. Hao, Y. Zhao, Y. Zhang, Y. Wang, C. Li, Z. Guan, D. Chen, and Y. Lin, "AutoDNNchip: An automated DNN chip predictor and builder for both FPGAs and ASICs," in *Proceedings of the 2020 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2020, pp. 40–50.
- [11] Y. Yang, Q. Huang, B. Wu, T. Zhang, L. Ma, G. Gambardella, M. Blott, L. Lavagno, K. Vissers, J. Wawrzyniek *et al.*, "Synetgy: Algorithm-hardware co-design for convnet accelerators on embedded fpgas," in *Proceedings of the 2019 ACM/SIGDA international symposium on field-programmable gate arrays*, 2019, pp. 23–32.
- [12] C. Hao, X. Zhang, Y. Li, S. Huang, J. Xiong, K. Rupnow, W.-m. Hwu, and D. Chen, "FPGA/DNN co-design: An efficient design methodology for IoT intelligence on the edge," in *2019 56th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2019, pp. 1–6.
- [13] C. Hao, Y. Chen, X. Zhang, Y. Li, J. Xiong, W.-m. Hwu, and D. Chen, "Effective algorithm-accelerator co-design for AI solutions on edge devices," in *Proceedings of the 2020 on Great Lakes Symposium on VLSI*, 2020, pp. 283–290.
- [14] W. Jiang, L. Yang, E. H.-M. Sha, Q. Zhuge, S. Gu, S. Dasgupta, Y. Shi, and J. Hu, "Hardware/software co-exploration of neural architectures," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 12, pp. 4805–4815, 2020.
- [15] Z. Dong, Y. Gao, Q. Huang, J. Wawrzyniek, H. K. So, and K. Keutzer, "Hao: Hardware-aware neural architecture optimization for efficient inference," in *2021 IEEE 29th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 2021, pp. 50–59.
- [16] H. Fan, M. Ferianc, Z. Que, H. Li, S. Liu, X. Niu, and W. Luk, "Algorithm and hardware co-design for reconfigurable CNN accelerator," in *2022 27th Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2022, pp. 250–255.
- [17] X. Zhang, Y. Li, J. Pan, and D. Chen, "Algorithm/Accelerator Co-Design and Co-Search for Edge AI," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 69, no. 7, pp. 3064–3070, 2022.
- [18] J. Ney, D. Lorocho, V. Rybalkin, N. Weber, J. Krüger, and N. Wehn, "HALF: Holistic auto machine learning for FPGAs," in *2021 31st International Conference on Field-Programmable Logic and Applications (FPL)*. IEEE, 2021, pp. 363–368.
- [19] M. S. Abdelfattah, L. Dudziak, T. Chau, R. Lee, H. Kim, and N. D. Lane, "Best of both worlds: Automl codesign of a CNN and its hardware accelerator," in *2020 57th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2020, pp. 1–6.
- [20] V. Kathail, "Xilinx Vitis unified software platform," in *Proceedings of the 2020 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2020, pp. 173–174.
- [21] "OpenVINO toolkit," <https://github.com/openvinotoolkit/openvino>, 2023.
- [22] Y. Umuroglu, N. J. Fraser, G. Gambardella, M. Blott, P. Leong, M. Jahre, and K. Vissers, "FINN: A framework for fast, scalable binarized neural network inference," in *Proceedings of the 2017 ACM/SIGDA international symposium on field-programmable gate arrays*, 2017, pp. 65–74.
- [23] J. Duarte, S. Han, P. Harris, S. Jindariani, E. Kreinar, B. Kreis, J. Ngadiuba, M. Pierini, R. Rivera, N. Tran *et al.*, "Fast inference of deep neural networks in FPGAs for particle physics," *Journal of Instrumentation*, vol. 13, no. 07, p. P07027, 2018.
- [24] S. I. Venieris and C.-S. Bouganis, "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, 2016, pp. 40–47.
- [25] J. Vandebon, J. G. F. Coutinho, W. Luk, E. Nurvitadhi, and T. Todman, "Artisan: a Meta-Programming Approach For Codifying Optimisation Strategies," in *2020 IEEE 28th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, 2020.
- [26] E. A. Moreno, O. Cerri, J. M. Duarte, H. B. Newman, T. Q. Nguyen, A. Periwal, M. Pierini, A. Serikova, M. Spiropulu, and J.-R. Vlimant, "JEDI-net: a jet identification algorithm based on interaction networks," *The European Physical Journal C*, vol. 80, no. 1, pp. 1–15, 2020.
- [27] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [29] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [30] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," 2011.
- [31] Y. Umuroglu, Y. Akhauri, N. J. Fraser, and M. Blott, "LogicNets: co-designed neural networks and circuits for extreme-throughput applications," in *2020 30th International Conference on Field-Programmable Logic and Applications (FPL)*. IEEE, 2020, pp. 291–297.
- [32] Z. Que, M. Loo, H. Fan, M. Blott, M. Pierini, A. D. Tapper, and W. Luk, "LL-GNN: Low Latency Graph Neural Networks on FPGAs for Particle Detectors," *arXiv preprint arXiv:2209.14065*, 2022.
- [33] F. Wojcicki, Z. Que, A. D. Tapper, and W. Luk, "Accelerating Transformer Neural Networks on FPGAs for High Energy Physics Experiments," in *2022 International Conference on Field-Programmable Technology (ICFPT)*. IEEE, 2022, pp. 1–8.
- [34] "Xilinx AI Engines and Their Applications," in *WP506(v1.1)*, July 10, 2020.
- [35] M. Langhammer, E. Nurvitadhi, B. Pasca, and S. Gribok, "Stratix 10 NX Architecture and Applications," in *The 2021 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2021, pp. 57–67.