

Characterizing CNN Throughput and Energy Under Multithreaded and Multiaccelerator Execution

M A Muneeb^{1b}, *Member, IEEE*, and Rajesh Kedia^{1b}, *Senior Member, IEEE*

Abstract—Emerging applications and batch processing convolutional neural network (CNN) workloads require executing multiple CNNs concurrently. A wide variety of CNN accelerators are available today and we characterize the support for concurrency for CNNs in such accelerators. We use a commercial-off-the-shelf CNN accelerator in multithreading and multiaccelerator modes and identify that upto 3.98× improvement in throughput and 3.20× improvement in energy per inference can be obtained even with just a single accelerator. Our detailed characterization of 104 CNN models, for three different sizes of accelerator, reveals many insights that connect CNN characteristics to improvement in throughput and energy. We also present a design space and a low error throughput estimation model to explore such a design space.

Index Terms—Accelerator, convolutional neural network (CNN), embedded system, field programmable gate array (FPGA).

I. INTRODUCTION

THE SUCCESS of convolutional neural networks (CNNs) in applications like object detection, classification, and segmentation has led to a wide variety of CNN accelerators [1], [2], [3]. Most CNN accelerators for embedded applications are designed for executing one CNN at a time. However, emerging applications like autonomous driving assistant systems [4] and batch processing workloads require executing multiple CNNs concurrently. Thus, it is important to evaluate CNN accelerators for concurrent execution which has not been considered in prior CNN characterization [5], [6].

Typically, a host CPU invokes the CNN execution on such accelerators and also executes any nonstandard CNN layer not supported by the accelerator. Concurrent execution can be realized in two ways: 1) when only one accelerator is available, but multiple CNNs execute as separate threads on the host processor invoking the accelerator in-turns, or 2) when multiple accelerator instances are available and multiple CNNs executing as separate threads use them concurrently. We refer to these scenarios as multithreaded and multiaccelerator mode of executions, respectively. We use a commercial-off-the-shelf (COTS) platform for experiments (Section III-A1) which has multiple host CPU cores to realize multithreading and is configurable to implement multiple accelerator instances. It also supports different sizes of accelerator.

While our characterization is specific to CNNs and for a specific platform, many of the observations are generic in nature

and applicable to other systems as well. To support design space exploration (DSE), we also propose an analytical model to quickly estimate the throughput under concurrent execution. Specifically, we claim the following key contributions.

- 1) We perform a detailed experimentation and measure the throughput and energy consumption for 104 different CNNs, in various multithreading and multiaccelerator modes, for three different sizes of accelerator.
- 2) We develop various insights into functioning of CNN accelerators leading to further research problems.
- 3) We propose an analytical model (< 1% average error) to estimate the throughput for the multithreading scenario.

II. RELATED WORK

Many research works have proposed CNN accelerators [1], [2], [3]. However, most of their design and analysis is focused on standalone performance of the accelerator without any consideration for concurrent CNNs. Recently, a few works analysed and improved the speed of concurrent CNN accelerators [7], [8], [9]. However, their concurrency implies use of more than one accelerator and they do not consider using only one accelerator in a multithreaded context. Another work considers time multiplexing of an field programmable gate array (FPGA)-based accelerator [10] in a server-client setup but focuses on improving scheduling and utilization of FPGA accelerator across various users. PERSEUS [11] evaluates performance and cost for multitenant CNN execution for CPU-GPU servers but uses only two CNN models. Shafi et al. [5] characterized CNNs for GPU-based devices and frameworks while Hadidi et al. [6] characterized CNN execution on different edge devices. Both these works do not consider the concurrent execution scenarios. We characterize and analyze the behavior of a wide range of CNN models for concurrent execution (multithreading and multiaccelerator) scenarios.

III. CHARACTERIZING CNN EXECUTION

A. Experimental Setup

1) Hardware Platform: We use Xilinx ZCU102 board having four ARM CPU cores and a programmable FPGA logic which is used to implement a COTS CNN accelerator named deep learning processor unit (DPU). We use v4.1.0 of the DPU IP (DPUCZDX8G) and the Vitis AI 3.0 framework for compiling standard CNN models for execution on DPU. We use three different sizes of DPUs—B4096, B2304, and B512; which we refer as large, medium, and small sized accelerator, respectively. ZCU102 board can support three instances of large, four instances of medium, or eight instances of small sized accelerator. Based on the CNN characteristics, the Vitis AI compiler maps a layer on DPU if it is supported by the DPU variant or otherwise maps it to the CPU.

2) Workloads: We use 104 standard CNN models from Model Zoo [12] to form the workloads. A workload

Manuscript received 5 August 2024; accepted 9 August 2024. This work was supported by the Science and Engineering Research Board (SERB) under Project SRG/2022/001083. This manuscript was recommended for publication by A. Shrivastava. (Corresponding author: Rajesh Kedia.)

The authors are with the Department of Computer Science and Engineering, Indian Institute of Technology Hyderabad, Hyderabad 502285, India (e-mail: muneeb.bagdali@cse.iith.ac.in; rkedia@cse.iith.ac.in).

Digital Object Identifier 10.1109/LES.2024.3446896

TABLE I
SELECTED 10 CNNs DISCUSSED IN DETAIL (OUT OF 104 CNNs CHARACTERIZED)

Model Name (as in Model Zoo repository)	Wokload Name	Single thread CPU time (ms)			Single thread Acc. time (ms)		
		Large	Med.	Small	Large	Med.	Small
pt_fadnet_sceneflow_576_960_0.65_154G_3.0	fad	8.0	8.0	7.9	16.2	27.2	120.3
pt_psmnet_sceneflow_576_960_0.68_696G_3.0	psm	106.0	106.8	105.5	62.9	83.5	191.5
pt_salsanextv2_semantic-kitti_64_2048_0.75_33.27G_3.0	salsa	9.5	9.2	9.4	120.4	143.2	260.8
tf_superpoint_mixed_480_640_52.4G_3.0	super	454.3	456.4	445.4	67.4	116.8	365.4
tf_yolov4_coco_416_416_60.3G_3.0	yolo	3.5	3.7	3.0	71.9	119.7	441.4
tf_vgg16_imagenet_224_224_0.43_17.67G_3.0	vgg16	0.2	0.2	42.7	24.9	36.9	119.9
tf_vgg19_imagenet_224_224_0.24_29.79G_3.0	vgg19	0.2	116.6	123.2	51.5	52.9	206.4
tf_inceptionresnetv2_imagenet_299_299_26.35G_3.0	in2	0.4	0.4	1474.9	48.6	64.5	191.7
tf_inceptionv3_imagenet_299_299_0.2_9.1G_3.0	in3	0.3	0.4	0.4	15.7	22.6	74.1
tf_mobilenetv2_1.0_imagenet_224_224_602M_3.0	mob2	0.3	0.3	0.3	3.9	4.0	9.9

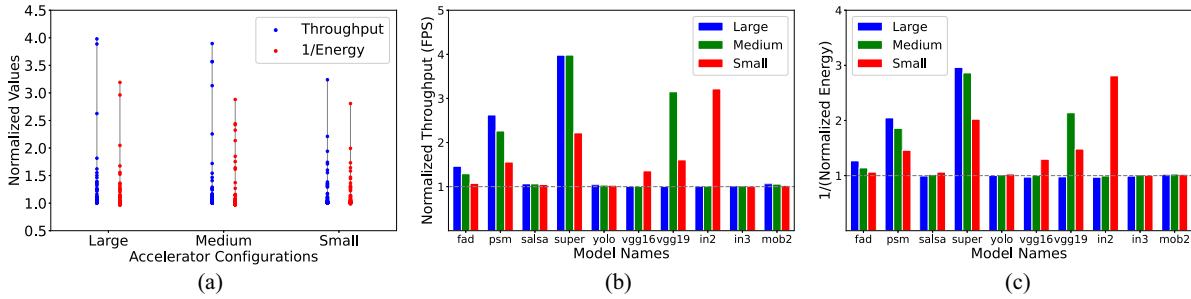


Fig. 1. Throughput and energy improvement with 4 threads, 1 accelerator (normalized to single thread). (a) Throughput and energy improvement. Each point refers to one out of 104 workloads. (b) Throughput improvement for selected CNNs. (c) Energy improvement for selected CNNs.

95 consists of one CNN being continuously executed back
 96 to back for a fixed duration. These models span various
 97 CNNs like FadNet, Inception (v1, v2, v3, v4), ResNet (v1,
 98 v2), SalsaNext, VoxelNet, YOLO (v4), RefineDet, VGG16,
 99 VGG19, EfficientNet, SqueezeNet, and MobileNet (v1, v2,
 100 v3) and include their varying resolutions and pruning ratios.
 101 While we characterize all CNNs, we choose ten of them
 102 (Table I) based on their CPU and accelerator times to illustrate
 103 important observations.

104 3) *Execution Flow and Measurement Setup*: We generate
 105 three different platform configurations, each containing one
 106 instance of the large, medium, and small sized accelerator,
 107 to be used for the multithreaded scenario. Three additional
 108 configurations containing three, four, and eight instances of the
 109 large, medium, and small sized accelerators are created for the
 110 multiaccelerator scenario. We execute each workload for 60 s
 111 which allows a large number of executions of the CNN. We
 112 specify the number of threads when invoking the workloads.
 113 Based on the platform configuration, the Vitis runtime system
 114 uses the available accelerator instances, thereby supporting
 115 multithreading and multiaccelerator scenarios.

116 We obtain the average throughput in the form of frames per
 117 second (FPS) and the average times taken by the CPU and
 118 the accelerator using the *xdputil* utility and *profiler* provided
 119 within Vitis AI framework. The power is measured using the
 120 software accessible on-board current sensors at every 1 ms and
 121 accumulated to calculate the energy consumption.

122 B. Observations and Insights

123 1) *Multithreading With Single Accelerator Instance*: Since
 124 the number of CPU cores on this platform is fixed as four, we
 125 execute various workloads with a multithreading level of four
 126 to answer the following question—How much throughput can
 127 be improved with only a single CNN accelerator instance?

128 Fig. 1(a) shows the normalized throughput for all the 104
 129 workloads for different sized accelerators. The throughput
 130 increases by upto $3.98\times$ even with only one accelerator

131 instance. CNNs having a larger fraction of the execution
 132 time on CPU core compared to the accelerator show larger
 133 improvements. This improvement is primarily because when
 134 different CNN instances are executing as different threads,
 135 one CNN might be using the DPU while some other can
 136 execute its CPU portion. Additionally, multiple CPU cores can
 137 enable upto four CNN instances to execute their CPU portions
 138 concurrently.

139 Fig. 1(b) shows the normalized performance under
 140 multithreaded execution for ten selected workloads (Table I).
 141 The workload *super* achieves a speedup close to 4 for large and
 142 medium size of the accelerator because of a significantly larger
 143 CPU time than accelerator time. For most of the workloads,
 144 the speedup decreases when moving from a bigger to a smaller
 145 accelerator. This is because the accelerator time increases
 146 from bigger to smaller accelerator causing the fraction of
 147 CPU time to decrease. However, for *vgg16*, *vgg19*, and *in2*
 148 workloads, smaller accelerators depict more speedup. As
 149 shown in Table I, the CPU time is drastically higher for these
 150 workloads for smaller sized accelerators as some layers that
 151 cannot fit on the smaller accelerator are executed on the CPU
 152 core.

153 Multithreading not only improves the throughput, but also
 154 improves (reduces) the energy consumption per inference.¹
 155 Across the workloads, Fig. 1(a) shows that multithreading
 156 reduces the energy consumption per inference by a factor of
 157 upto $3.20\times$. Such a reduction happens because the power
 158 consumption does not increase much with multithreading and
 159 therefore, more inferences can be executed for the same energy
 160 due to increased throughput. However, for a few CNNs having
 161 a negligible CPU time compared to the accelerator time, the
 162 energy per inference increases slightly (about 1%–3%). This
 163 is most likely due to switching overheads associated with
 164 multithreading. Fig. 1(c) shows the energy reduction normalized

¹We plot energy reduction as an energy improvement (>1) by taking the inverse so that it is easier to visually correlate throughput and energy.

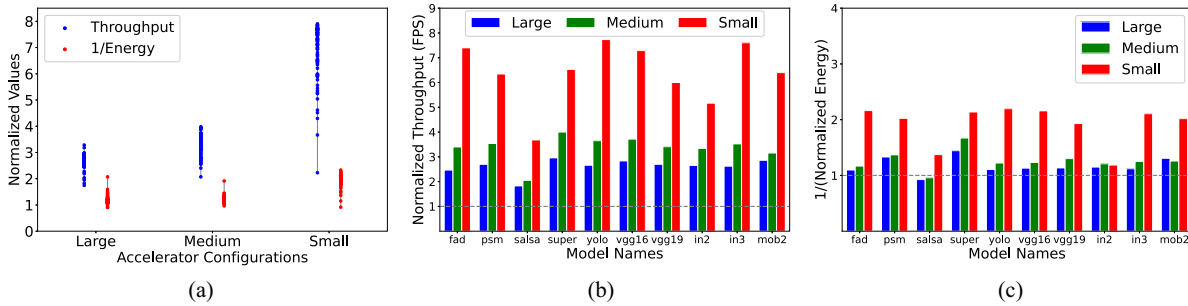


Fig. 2. Throughput and energy improvement with multi-accelerators (normalized to single thread). (a) Throughput and energy improvement. Each point refers to one out of 104 workloads. (b) Throughput improvement for selected CNNs. (c) Energy improvement for selected CNNs.

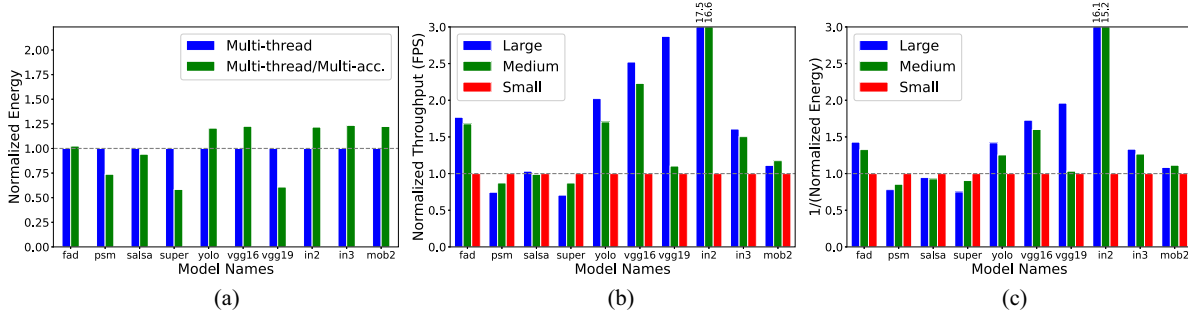


Fig. 3. Throughput and energy comparison across different implementations. (a) Energy comparison of multi-accelerator and multi-threading for Medium-sized accelerator. (b) Throughput improvement with multi-accelerators (normalized to Small-sized). (c) Energy improvement with multi-accelerators (normalized to Small-sized).

165 to a single-thread execution for the ten selected CNNs. The
166 energy reduction follow a similar trend as the throughput
167 improvement, which was discussed earlier.

168 2) *Using Multiple Instances of CNN Accelerator*: We now
169 consider platform configurations with multiple accelerator
170 instances implemented and accessed by different threads ex-
171 ecuting on the CPU cores. While the previous multithreading
172 scenario improved throughput and energy for CNNs having
173 considerable CPU time in comparison to the accelerator time,
174 the multiaccelerator mode improves throughput for almost all
175 workloads, as shown in Fig. 2(a). The improvement depends
176 upon the number of accelerator instances, which is three, four,
177 and eight for the large, medium, and small sized accelerators,
178 respectively.

179 Fig. 2(b) shows that the throughput improves by a factor
180 close to, but smaller than the number of accelerator instances.
181 This is because multiple accelerators share internal buses
182 and memory bandwidth which might slow down their execu-
183 tion. Fig. 2(c) shows the improvement in energy for the ten
184 selected models. There are two contending factors determining
185 the energy per inference: 1) the power consumption, which
186 increases due to multiple accelerator instances and 2) the time
187 per inference, which decreases due to improved throughput.
188 In comparison to the throughput, the energy per inference
189 improves (reduces) only by a small amount.

190 3) *Comparing Multithreading With Multiaccelerator*:
191 While multiaccelerator mode provides similar or better
192 throughput than the multithreaded mode, it does not apply
193 to energy consumption. Fig. 3(a) shows the relative energy
194 consumption for the medium sized accelerator. Some
195 workloads have lesser energy for multithreaded mode while
196 others have lesser energy for multiaccelerator mode depend-
197 ing upon their time on CPU core and accelerator. This opens up
198 a DSE problem to identify appropriate level of multithreading
199 and number of accelerators to use.

4) *Effect of Accelerator Size*: A larger accelerator exe- 200
cutes a workload faster, but only a fewer instances can be 201
implemented compared to a smaller accelerator which takes 202
longer to execute a workload but allows more number of 203
concurrent instances. We study the effect of accelerator size 204
with workloads executing with three, four, or eight 205
instances of large, medium, or small sized accelerator, respec- 206

207 Fig. 3(b) and (c) shows the throughput and energy per 207
inference for each workload normalized to that of the small 208
sized accelerator. For most of the workloads, throughput 209
increases and energy decreases with an increase in accelera- 210
tor size (similar to prior works [9]). This is primarily be- 211
cause the number of concurrent multiply and accumulate (MAC) 212
operations within a larger accelerator is much larger than the 213
ratio of number of instances compared to a smaller accelera- 214
tor. However, our characterization identifies that a few workloads 215
deviate from this trend which we discuss in detail. 216

217 For the *psm* and *super* workloads, the highest throughput 217
as well as lowest energy is obtained for the small sized 218
accelerator (instead of the large). This is because apart from 219
the accelerator, these models execute on the CPU core for a 220
significant amount of time (see Table I). Hence, they benefit 221
more from larger number of accelerator instances than from 222
faster accelerators. The *mob2* workload achieves the best 223
throughput and energy consumption with the medium sized 224
accelerator. This CNN has execution time of 3.9, 4.0, and 225
9.9 ms on the large, medium, and small sized accelerators, 226
respectively (Table I). Since the time for large and medium 227
accelerators is almost similar, *mob2* gets benefit from having 228
more instances of medium sized accelerator. However, its ex- 229
ecution time is drastically higher for the small sized accelera- 230
tor, making the medium sized accelerator to be best for it. While 231
prior works [9] proposed a larger accelerator to be always used 232
if available, we observe that a medium or small accelerator 233
could also be better based on the workload characteristics. 234

235 This presents another DSE problem to identify a suitable size
236 of accelerator for a given workload.

237 The *vgg19* and *in2* workloads experience a very high
238 improvement for larger accelerators as some of the layers
239 cannot execute on a smaller DPU and get mapped to CPU.
240 This reduces the throughput for smaller accelerators and hence
241 higher improvement is seen for larger accelerators. This can
242 be considered as one limitation of the end-to-end framework
243 associated with this accelerator (DPU). Ideally, the compiler
244 framework should have split the layer into smaller units to
245 execute on accelerator rather than on the CPU core. However,
246 with fast emergence of newer types of layers and CNN
247 architectures, it is also a practical approach followed by most
248 frameworks to migrate unsupported layers to the host CPU.

249 IV. SYSTEM MODELING FOR THROUGHPUT

250 Having seen that multithreading and multiaccelerator exe-
251 cution of CNNs can improve the throughput, and based on
252 the need for DSE to identify appropriate level of multi-
253 threading and count of accelerators, we propose an analytical
254 model to estimate the throughput for multithreading scenario.
255 Multiaccelerator scenario involves complex modeling to esti-
256 mate contention among accelerators [8] and is left as future
257 work.

258 For a system with n CPU cores, we want to estimate the
259 throughput for a multithreading level of n and with only 1
260 instance of the DPU. We assume that the CPU time (c) and
261 DPU time (d) for a single-thread execution are already known.
262 DPU execution can completely hide the CPU execution time
263 if $c \leq (n - 1) \times d$, as illustrated in Fig. 4 for $n = 4$ and
264 $c = 3 \times d$. Therefore, in such cases, the estimated throughput
265 is determined by only the DPU time (d).

266 However, when CPU time becomes further long, i.e., $c >$
267 $(n - 1) \times d$; then, the total time for DPU execution for other
268 threads is not sufficient for a CPU to complete its processing
269 to again invoke the DPU. Therefore, DPU's utilization reduces
270 and even a single DPU instance can be available to all threads
271 on need. In such cases, the single-thread throughput can
272 increase by a factor n . Therefore, a simple analytical model
273 (under ideal conditions) for the throughput (T) is

$$274 T = \begin{cases} \frac{1}{d}, & \text{if } c \leq (n - 1) \times d \\ \frac{1}{c + d}, & \text{if } c > (n - 1) \times d. \end{cases} \quad (1)$$

275 *Evaluating Throughput Estimation:* Fig. 5 shows the per-
276 centage error between the measured and estimated throughput.
277 The error values are positive indicating that our predictions
278 are an upper bound. Even with a simple model, the estimation
279 error is within 5% (except for 6 outlier points out of 312).
280 The mean error is $< 1\%$ (0.65%, 0.56%, and 0.67% for the
281 large, medium, and small sized accelerators, respectively).
282 The high error for some workloads (e.g., *vgg16*, *vgg19* for
283 medium sized accelerator) is because their CPU time is very
284 large and concurrent execution on four CPU cores experiences
285 slowdown due to shared resources (e.g., LLC).

286 V. CONCLUSION

287 We presented detailed characterization of 104 different
288 CNN workloads executing on an FPGA-based COTS CNN
289 accelerator named DPU. Throughput and energy improves
290 with multithreading and multiple instances of the accelerator.
291 Our study also included three different sizes of accelerator
292 and derived key insights: 1) how different CNNs are affected

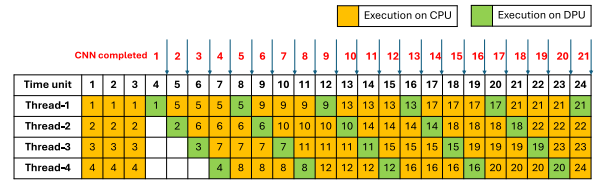


Fig. 4. CNN execution with four threads and one accelerator with $c = 3 \times d$.

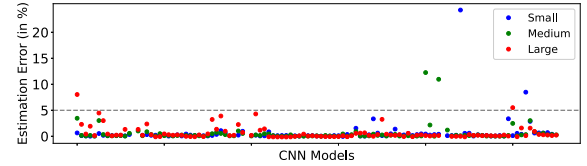


Fig. 5. Error in throughput estimation for 104 CNNs and 3 sizes.

293 with multithreading and multiaccelerators; 2) how the size of
294 accelerator affects throughput and energy; and 3) need for DSE
295 as per the workload. We also developed a model to estimate
296 the throughput improvements due to the use of multithreading,
297 which matches the measurements quite closely.

298 In the future, we would like to develop estimation models
299 for multiaccelerator mode, evaluate other platforms and archi-
300 tectures, and expand the framework to support DSE.

301 REFERENCES

- 302 [1] Y. Chen, Y. Xie, L. Song, F. Chen, and T. Tang, "A survey of accelerator
303 architectures for deep neural networks," *Engineering*, vol. 6, no. 3,
304 pp. 264–274, 2020.
- 305 [2] K. Guo et al., "Angel-eye: A complete design flow for
306 mapping CNN onto embedded FPGA," *IEEE Trans. Comput.-
307 Aided Design Integr. Circuits Syst.*, vol. 37, no. 1, pp. 35–47,
308 Jan. 2018.
- 309 [3] K. Guo, S. Zeng, J. Yu, Y. Wang, and H. Yang, "A survey of FPGA-based
310 neural network inference accelerators," *ACM Trans. Reconfig. Technol.
311 Syst.*, vol. 12, no. 1, pp. 1–26, 2019.
- 312 [4] J. Peng et al., "Multi-task ADAS system on FPGA," in *Proc. IEEE
313 Int. Conf. Artif. Intell. Circuits Syst. (AICAS)*, Hsinchu, Taiwan, 2019,
314 pp. 171–174.
- 315 [5] O. Shafi, C. Rai, R. Sen, and G. Ananthanarayanan, "Demystifying
316 TensorRT: Characterizing neural network inference engine on Nvidia
317 edge devices," in *Proc. IEEE Int. Symp. Workload Character. (IISWC)*,
318 Storrs, CT, USA, 2021, pp. 226–237.
- 319 [6] R. Hadidi, J. Cao, Y. Xie, B. Asgari, T. Krishna, and H. Kim,
320 "Characterizing the deployment of deep neural networks on commercial
321 edge devices," in *Proc. IEEE Int. Symp. Workload Character. (IISWC)*,
322 2019, pp. 35–48.
- 323 [7] Z. Du, W. Zhang, Z. Zhou, Z. Shao, and L. Ju, "Accelerating DNN infer-
324 ence with heterogeneous multi-DPU engines," in *Proc. 60th ACM/IEEE
325 Design Autom. Conf. (DAC)*, 2023, pp. 1–6.
- 326 [8] S. Goel, R. Kedia, M. Balakrishnan, and R. Sen, "INFER: INterFERENCE-
327 aware estimation of runtime for concurrent CNN execution on
328 DPUs," in *Proc. Int. Conf. Field-Program. Technol. (ICFPT)*, 2020,
329 pp. 66–71.
- 330 [9] R. Kedia, S. Goel, M. Balakrishnan, K. Paul, and R. Sen, "Design
331 space exploration of FPGA-based system with multiple DNN accel-
332 erators," *IEEE Embed. Syst. Lett.*, vol. 13, no. 3, pp. 114–117,
333 Sep. 2021.
- 334 [10] T. Nakamura, S. Saito, K. Fujimoto, M. Kaneko, and A. Shiraga,
335 "Spatial- and time- division multiplexing in CNN accelerator," *Parallel
336 Comput.*, vol. 111, Jul. 2022, Art. no. 102922.
- 337 [11] M. LeMay, S. Li, and T. Guo, "PERSEUS: Characterizing performance
338 and cost of multi-tenant serving for CNN models," in *Proc. Int. Conf.
339 Cloud Eng. (IC2E)*, 2020, pp. 66–72.
- 340 [12] (Xilinx Inc., San Jose, CA, USA). *Vitis AI Model Zoo*. Accessed: Jan. 31,
341 2024. [Online]. Available: [https://xilinx.github.io/Vitis-AI/3.0/html/docs/
342 workflow-model-zoo](https://xilinx.github.io/Vitis-AI/3.0/html/docs/workflow-model-zoo)