# FOS: A Modular FPGA Operating System for Dynamic Workloads

**Anuj Vaishnav,**
Department of Computer Science,
The University of Manchester,
anuj.vaishnav@manchester.ac.uk

**Khoa Dang Pham,**
Department of Computer Science,
The University of Manchester,
khoa.pham@manchester.ac.uk

**Joseph Powell,**
Department of Computer Science,
The University of Manchester,
joseph.powell@manchester.ac.uk

**Dirk Koch,**
Department of Computer Science,
The University of Manchester,
dirk.koch@manchester.ac.uk

January 29, 2020

## ABSTRACT

With FPGAs now being deployed in the cloud and at the edge, there is a need for scalable design methods which can incorporate the heterogeneity present in the hardware and software components of FPGA systems. Moreover, these FPGA systems need to be maintainable and adaptable to changing workloads while improving accessibility for the application developers. However, current FPGA systems fail to achieve modularity and support for multi-tenancy due to dependencies between system components and lack of standardised abstraction layers. To solve this, we introduce a modular FPGA operating system – FOS, which adopts a modular FPGA development flow to allow each system component to be changed and be agnostic to the heterogeneity of EDA tool versions, hardware and software layers. Further, to dynamically maximise the utilisation transparently from the users, FOS employs resource-elastic scheduling to arbitrate the FPGA resources in both time and spatial domain for any type of accelerators. Our evaluation on different FPGA boards shows that FOS can provide performance improvements in both single-tenant and multi-tenant environments while substantially reducing the development time and, at the same time, improving flexibility.

***Keywords*** FPGA, Operating System, Resource-elasticity, Modular Development, Dynamic Workloads

## 1 Introduction

FPGA accelerators can provide high-performance computing at very low energy cost for applications ranging from neural-networks to network processing. This has brought FPGAs in cloud datacenters as well as in embedded systems at the edge. However, to sustain this large scope of requirements present in terms of heterogeneity of devices, environments, EDA tools, users and developer needs, we require a standardised way to manage and integrate FPGA systems components, in other words: we need an FPGA operating system.

The idea of an FPGA operating system is old and has taken many forms over time. It has evolved from providing a) OS-level APIs such as hardware threads [1, 2] and UNIX interfaces [3] for allowing hardware accelerators to use existing software OS services, to b) light-weight FPGA shells with library APIs for hardware acceleration [4–8].

We believe this evolution continues for two primary reasons: *overhead and lack of portability*. The introduction of intermediate layers and communication cost must be low in terms of latency as well as in the amount of resources required. However, to achieve this is difficult because hardware accelerators and boards require unique optimisations and implementations, leading to heterogeneity in FPGA shells, accelerators, EDA tools, and the low-level software
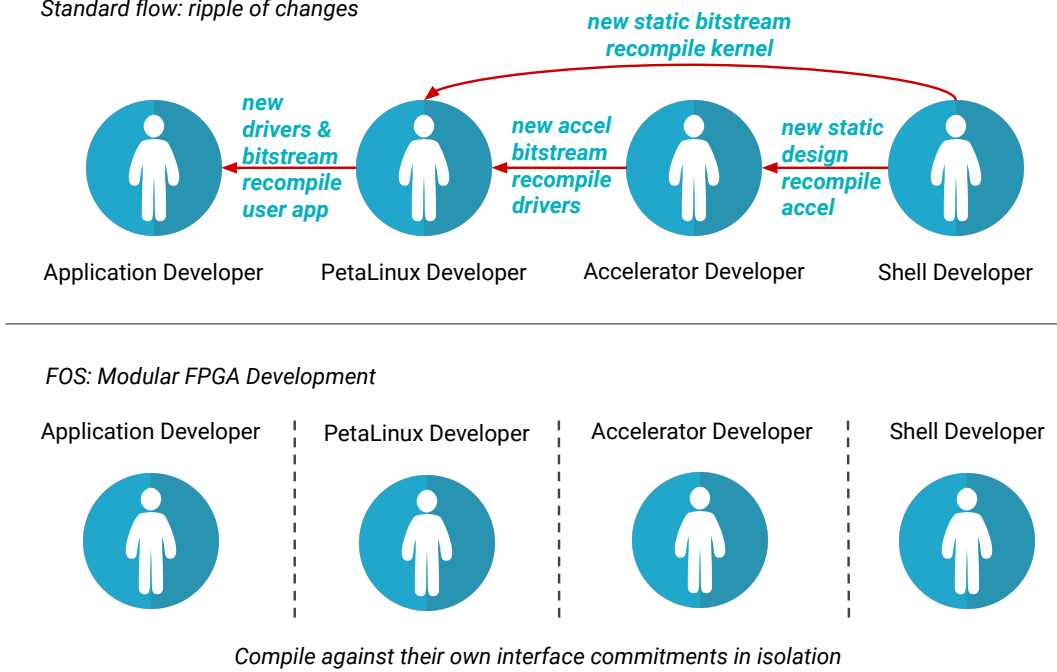
*Standard flow: ripple of changes*

*new static bitstream recompile kernel*

*new drivers & bitstream recompile user app*

*new accel bitstream recompile drivers*

*new static design recompile accel*

Application Developer          PetaLinux Developer          Accelerator Developer          Shell Developer

*FOS: Modular FPGA Development*

Application Developer          PetaLinux Developer          Accelerator Developer          Shell Developer

*Compile against their own interface commitments in isolation*

Figure 1: If we update the shell IP or EDA tool version, the rest of the system components must be recompiled because of the dependencies in standard tool flow (see Section 2.2). Ideally, we want each component to be compiled in isolation with given interfaces, as we do for FOS.

required to interface with it. In particular, this reduces not only the portability of a system but also its *maintainability*, which is a key aspect for any operating system (OS). When using, for example, the current tool flows of the major FPGA vendors, a simple change in tool version or addition of system IP can lead to re-compilation of the whole FPGA stack [9, 10], as shown in Figure 1. Essentially, implying that an operating system update means recompilation of all applications which run on top of it.

Moreover, with the need for multi-tenancy in dynamic environments like the cloud, we expect multiple types of workloads and accelerators to execute concurrently [11] (e.g., for enabling resource pooling on FPGA accelerator cards). However, most present FPGA systems operate and support either static accelerators or single accelerator execution [12–15]. Systems which support hosting multiple accelerators, do not allow resources to be re-allocated during execution or employ only time-domain multiplexing [16–18]. This commonly leads to under-utilisation and a lack of *adaptability* in the system.

At the same time, there is a need for better APIs to make FPGAs more *accessible* to software programmers or application domain experts that are commonly not FPGA experts. This should be achieved, while ensuring that hardware developers can integrate their accelerators into software stacks with high productivity, i.e. without changing their hardware designs or writing complex software wrappers.

To solve these challenges of maintainability, adaptability, and accessibility, we are introducing a modular FPGA Operating System (FOS). FOS adopts a modular FPGA development flow to allow each type of OS component (hardware or software) to be replaced, reused or ported to different FPGA systems without extensive effort or compilation time. Moreover, the runtime system provides full support for multi-tenancy with the ability to concurrently execute accelerators written in various languages (C, C++, OpenCL, or RTL) and dynamically replicate or switch to a different version of an accelerator to improve utilisation and system performance. To interact with these accelerators, application developers can use high-level APIs (available in multiple languages) to access the FPGA in three modes: 1) static acceleration for a single user, 2) dynamic acceleration for a single user and 3) dynamic acceleration in multi-tenant environments, as shown in Figure 2. At the same time, hardware developers can write light-weight interface descriptions to integrate their hardware accelerators into the FOS platform and benefit from its high-level software APIs.
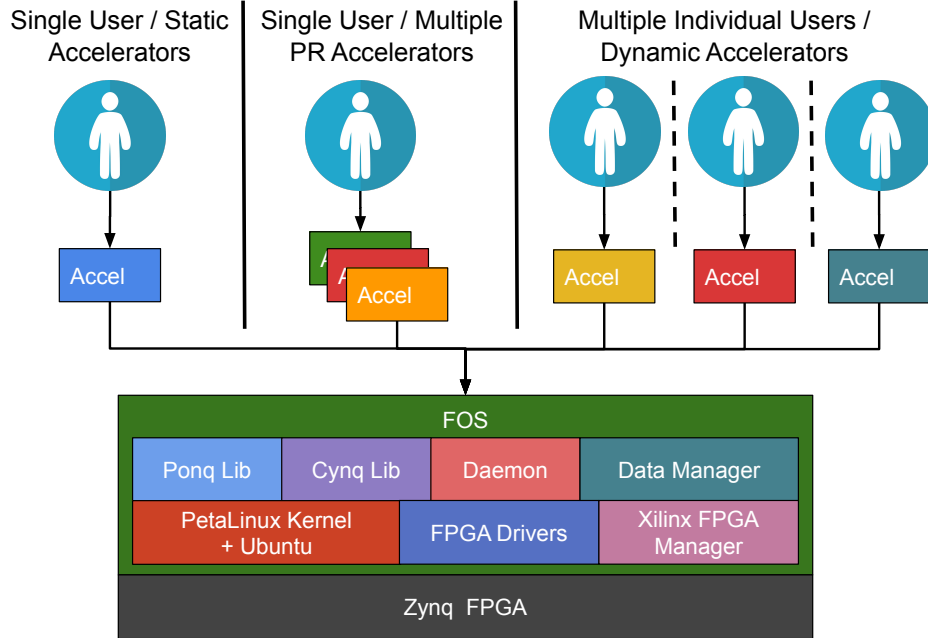
Figure 2: System components and usage mode of the FPGA Operating System – FOS.

All the research artefacts proposed in this paper were presented at the FPL 2019 Demo Night and are available online under an open-source license[1]. The key contributions of this paper are:

- A *modular FPGA development flow* to encapsulate the heterogeneity at each level of the development process (Section 2.2).

- An *open-source software platform* to improve programmability for static acceleration systems, dynamic systems for single-user as well as multi-user dynamic acceleration systems (Section 3 and 4).

- A *decoupled compilation flow for shell and modules* with support for module relocation and flexibility to combine PR regions to host accelerators of varying sizes (Section 4.1).

- A *resource-elastic runtime system for dynamic workloads* supporting virtually any type of accelerators (written in RTL, C, C++, and OpenCL) with a unified user interface and programming model (Section 4.4).

- A thorough evaluation of the FPGA Operating System (FOS) on resource overhead, compilation latency, memory throughput performance, application performance, as well as flexibility and maintainability of the system (Section 5).

## 2 Concepts

### 2.1 FPGA Operating System

Unlike standard CPU operating systems which consist only of a software infrastructure, an FPGA operating system requires an infrastructure at both the hardware and the software level (as shown in Figure 3). The following subsections describe the role of each level in detail.

#### 2.1.1 Hardware Infrastructure

The hardware infrastructure is implemented directly on top of the physical FPGA resources and is commonly built using a partial reconfiguration (PR) flow (e.g., Amazon F1 FPGA instances [12]). A PR flow generates a system in two partitions: a static system and one or more reconfigurable regions. An FPGA shell is essentially the static system of the PR flow and can be considered equivalent to the OS kernel to the hardware infrastructure. It provides common functionalities required by the accelerators such as an interconnect, network and memory access as well as

---
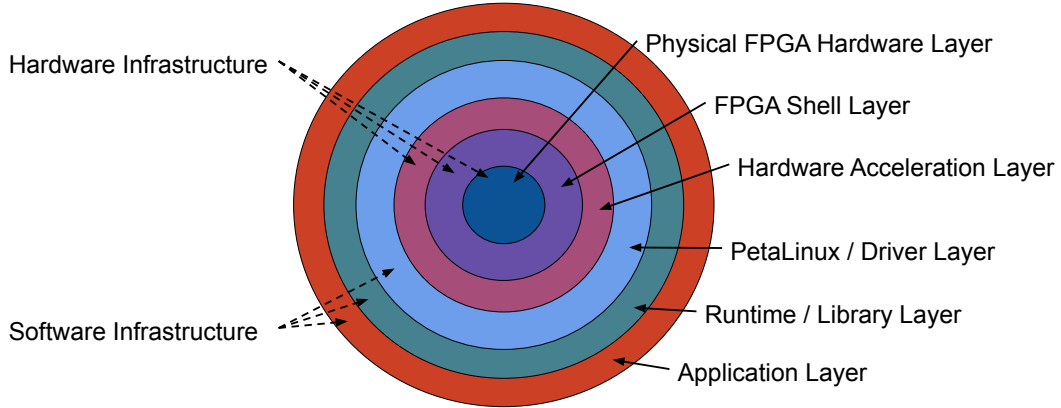
[1]https://github.com/khoapham/fos

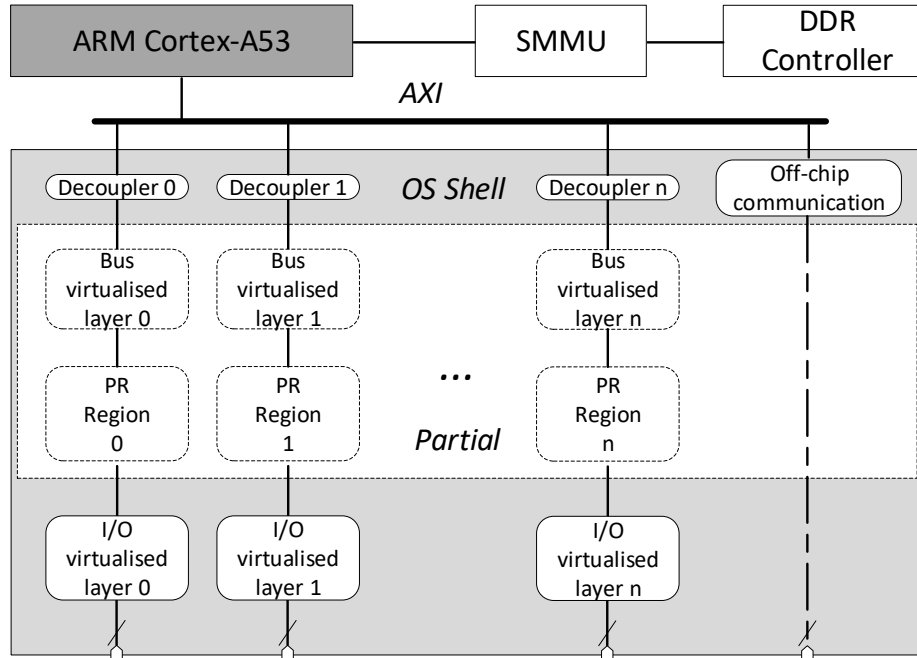Figure 3: The layered architecture of an FPGA operating system.



Figure 4: The overall organisation of an FPGA shell.

(optional) other I/O and management IPs necessary for the target system. The counterpart of the shell is the partially reconfigurable region (also called slot) which can host different hardware accelerators at runtime. A shell which supports multiple partial regions can support multi-tenancy in the spatial domain by allowing multiple accelerators to execute in parallel. Figure 4 shows an example of such a shell. However, the number of partial regions, their location and sizes depend on the system requirements and changes from system to system. Hence, an integral part of the shell is the compilation flow required to map the hardware modules onto the resources of a PR region with the required physical interface for the target system infrastructure.

Moreover, a shell often includes a CPU in the form of a soft-core (e.g., on the datacenter FPGAs [19]) or a hardened-core on MPSoC platforms used for embedded systems [19]. This CPU is used to host the software infrastructure necessary for the management of the FPGA resources (see Section 2.1.2 for details).

With this, a shell provides the basic OS functionality at the hardware level to host hardware applications (accelerators) on an FPGA.

### 2.1.2   Software Infrastructure

The software infrastructure of an FPGA operating system serves as a middle layer between the hardware accelerators and the user applications (software) while executing on a host CPU (soft or hardened). However, unlike software operating systems, the software infrastructure for an FPGA OS should be resilient to dynamic changes in the underlying hardware and be portable to different FPGA boards while hiding the heterogeneity from the software programmer. Consequently, an FPGA OS is responsible for managing the heterogeneity of the hardware resources available on the FPGA as well as to provide the high-level APIs and software integration for ease of development. This requires the software infrastructure to provide five important functionalities:

1. A boot-loader and a kernel to bring up the FPGA system in an operational state.

2. A set of drivers and hardware abstraction layer (HAL) to communicate with the hardware accelerators.

3. A scheduler to dynamically allocate FPGA resources (i.e. partial regions) and hardware accelerators to software applications.

4. High-level standard libraries to make hardware acceleration easily accessible.

5. Integration with existing CPU software stacks to benefit from legacy code.

## 2.2   Modular Development Flow

There are five primary stages of development required for a multi-tenant/cloud FPGA system, as shown in Figure 5. The bottom two stages of this stack are hardware development stages which drive the application performance and the support for multi-tenancy (via PR). In particular, 1) shell development requires designing and implementing the common system functionalities required by the user accelerators. An important step at this stage is floorplanning, which includes a decision on how many resources we allocate to the accelerators and a definition of the interface exposed to the accelerators. This stage also requires identifying the address mappings at which the driver layer will access the accelerators. While, its counterpart, 2) accelerator development, involves designing RTL or HLS accelerators with application-specific optimisations (commonly for the highest performance for the allocated resource budget). To compile the accelerators for the shell, the compiler must know the exact resources available in the partial regions and the physical locations of the interface pins. Without this information, we cannot implement partial reconfigurable accelerators. However, with the current development flow, accelerator compilation directly depends on the shell and requires knowing the implementation (internal resource allocation and routing) of a shell [20]. This is because accelerator modules are implemented as an increment to a specific shell. Hence, any update made to the shell requires recompiling of the accelerators. Note that this flow is used to build the most state-of-the-art shells [7, 15, 21–23].

The remaining three stages of the stack are software development stages which aim at programming the accelerators, runtime optimizations, and improving the developer productivity. 3) PetaLinux/Driver development involves building the embedded software required to boot up the board with necessary I/O and high-level functionalities as well as means to communicate with the accelerators to send and receive data. The current Xilinx tool flow expects the user to write a new driver for each accelerator or generate one automatically when using Vivado HLS [24]. This new driver has to be either built with the embedded Linux kernel [25] directly (for foreknown hardware) or with a device-tree overlay (for hardware known only at runtime). EDA tool vendors provide a hardware abstraction level (HAL) to make the development of the drivers easier via APIs to read and write to accelerators, to perform reconfiguration calls as well as basic C functionalities to enable debugging and fast development [26]. However, an accelerator developer or embedded Linux developer must take the responsibility to write and integrate the driver correctly in the rest of the Linux environment.

The layer above this fundamental driver and kernel layer forms 4) a runtime system or libraries for hardware acceleration. This provides a high-level API such as OpenCL to the user for serving reconfiguration and acceleration requests to the lower hardware layers [9, 10, 26]. However, currently, no FPGA vendor tools provide support for multi-tenancy.

Finally, 5) in the application development stage a developer (commonly a domain expert) performs the required task in software while using hardware acceleration where appropriate. Note that a developer at this stage does not need to possess the skills required to implement the underlying FPGA development stack.

Overall, each step requires sophisticated knowledge which makes designing systems challenging for individual or small design teams. Given the complexity and the effort required for the task involved, often the final system is application-specific and cannot be reused or ported for different needs. To avoid this, we need a standard set of APIs such that a component can be swapped at each stage without recompilation or redevelopment of the components above or below it as long as the APIs are maintained. Thus, allowing each stage to be a modular artefact in itself.
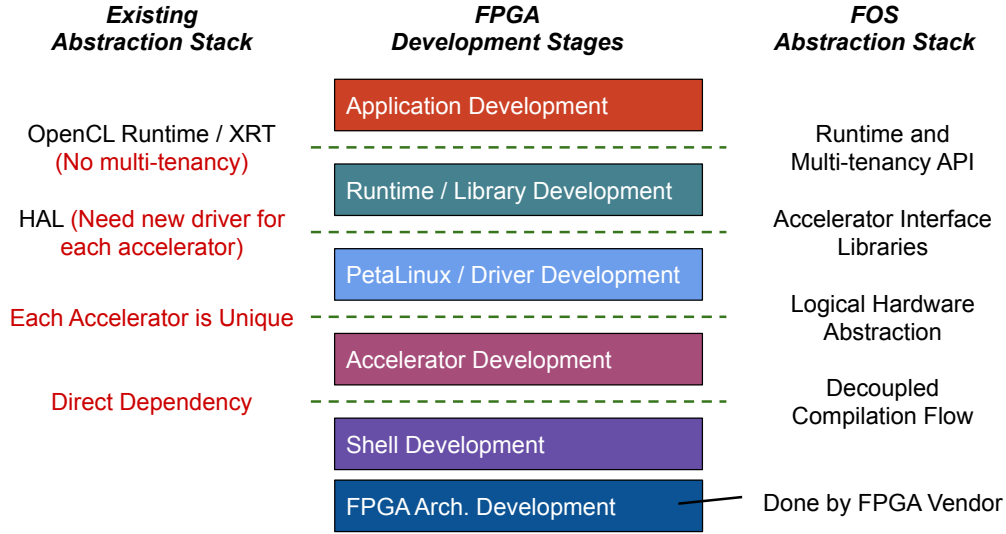
Figure 5: The abstraction layers required between FPGA development stages.

Such an abstraction stack would also allow 1) the software stack to be reusable across all types of FPGAs and FPGA boards, 2) not require the accelerator developers to write drivers, and 3) update system components in the shell with no need to propagate the changes to other stages in the stack.

We can achieve this by implementing 4 key layers of abstractions:

- **Decoupled Compilation Flow**: Ability to compile accelerators and shells in isolation from each other. Whereby the accelerator compiles against a fixed physical interface and a bounded resource region. This stops any changes in IPs or the shell from propagating to accelerators (given that we do not change the PR region and its interface).

- **Logical Hardware Abstraction**: Exposing the accelerator and the shell in a high-level format with only a minimum set of parameters required to build the drivers and perform the resource allocation. For the shell, this would include information such as how many regions it supports and what addresses they can be accessed at. Whereas, for the accelerators, it would include the i) internal hardware address register mapping (ADR map) for programming the accelerator and ii) meta-data associated with the accelerators for scheduling and management purposes (e.g., the size or maximum execution latency). Note that HLS compilation (through Vivado HLS) provides this information completely without the need of any manual step.

- **Accelerator Interface Libraries**: Using a standardised register mapping format to provide generic driver support for accelerators with streaming or master-slave interfaces (which are the most common interfaces provided by shells). Thus, it relieves the accelerator developer from the responsibility of writing drivers. In the case that the adoption of the standardised format is not feasible, a developer can use HAL APIs to build drivers.

- **Runtime and Multi-tenancy API**: Execution API which can abstract the interface at the logical level and which can span across multiple languages to provide high level integration to existing software stacks, while supporting the necessary primitives or programming models for dynamic resource allocation.

## 3   FOS Overview

The concepts for building an FPGA operating system (as introduced in Section 2.1) have been developed into FOS. FOS is a modular and lightweight FPGA Operating System which provides three primary modes of operation (as shown in Figure 2):

1. Execution on static accelerators in a single-tenant mode.
2. Using multiple partially reconfigurable accelerators in a single-tenant mode.
3. Dynamically offloading acceleration request in a multi-tenant mode.

To provide the first two-modes of operation with ease of access, FOS provides two libraries: Cynq and Ponq. They provide high-level APIs to interact with hardware from C++ and Python applications, respectively. These APIs are platform-independent and can support the traditional Xilinx PR flow as well as the decoupled compilation flow introduced with FOS.

The multi-tenant mode is provided through a daemon, which orchestrates the hardware acceleration requests from different users and schedules them in time *and* in the spatial domain. Moreover, the daemon can execute hardware acceleration requests on heterogeneous accelerators (i.e. accelerators written in different languages) concurrently while providing the same user interface.

To support the underlying hardware interaction in each mode, the libraries include generic drivers to program accelerators and the Xilinx FPGA manager [25] for partial reconfiguration. Moreover, FOS uses the PetaLinux Kernel [25] and Ubuntu rootfs, to adopt the standard functionalities provided by the Ubuntu Linux distribution such as access to file systems, network access (via host-CPU), standard libraries, debugging tools, and other forms of I/O.

The hardware infrastructure used in this paper is based on the open-source ZUCL 2.0 shell [27]. It currently supports three boards of varying FPGA capacity: ZCU102 (Xilinx MPSoC development kit), UltraZed, and Ultra-96 (suitable for IoT and edge deployment). The shell provides the ability to reallocate hardware accelerators across different partial regions as well as to combine multiple adjacent partial regions to host bigger accelerators.

With these additions to the software and hardware ecosystem for FPGAs, FOS achieves a similar level of support for rapid development, deployment, and ease of use, as known from standard operating systems for CPUs.

## 4   Implementation

### 4.1   Decoupled Compilation Flow for Shell and Modules

The primary requirement for abstracting accelerators from the shell is to decouple their compilation. However, this alone does not enable the flexible resource allocation required for the maximum utilisation of the resources. This is because of two main reasons, i) the standard vendor EDA tool flow cannot easily combine partial regions to host bigger accelerators and ii) it requires the accelerator to be compiled for each PR region it will be hosted in, hence resulting in multiple bitstreams for the same accelerator [20]. To solve this, we need a shell that can support module relocation and re-adjustable PR regions (i.e. the ability to combine PR regions). There are strict requirements to achieve development isolation with relocation and PR region flexibility:

1. PR regions should be homogeneous in terms of their resource foot-print (i.e. the relative layout of FPGA primitives) in order to allow for accelerator relocation, and regions should be located adjacent to each other to allow hosting of bigger accelerators without interfering with other system components.

2. Communication interfaces between modules and the static system must be identical in terms of both logical protocol and their physical implementation such that relative positions of connection wires are the same in all PR regions. This ensures that modules can receive operation commands from the host CPU and it provides interfaces to transfer data back and forth to the main memory, irrespective of a module placement position.

3. We must distribute clock signals in the same regular pattern across every PR region. These constrained clock routing paths will be used to provide clock signals to the resources used by the modules.

4. We must prohibit routing from the static part to pass through the reconfiguration part and vice versa (except the interface) to ensure that module relocation does not interfere with other parts of the system.[2]

The proposed design methodology depends on the standard FPGA development flow (Vivado toolchain for Xilinx FPGAs) to implement 1) the basic infrastructure, acting as the OS shell, and 2) the hardware applications (i.e. accelerator modules). Our contributions, however, include all the additional design steps and their corresponding tools that customise and adapt the default flow to automatically build the final system. This integrates the entire shell and module development process into a unified design framework, as illustrated in Figure 6.

The main steps of the compilation process are as follows:

1. *Planning*: In this step, a shell developer needs to make a series of system-level design decisions.

---

[2]This is not a strict requirement and the tool GoAhead [28] used to implement the shell can be used in a mode that allows static routing to cross reconfigurable regions.
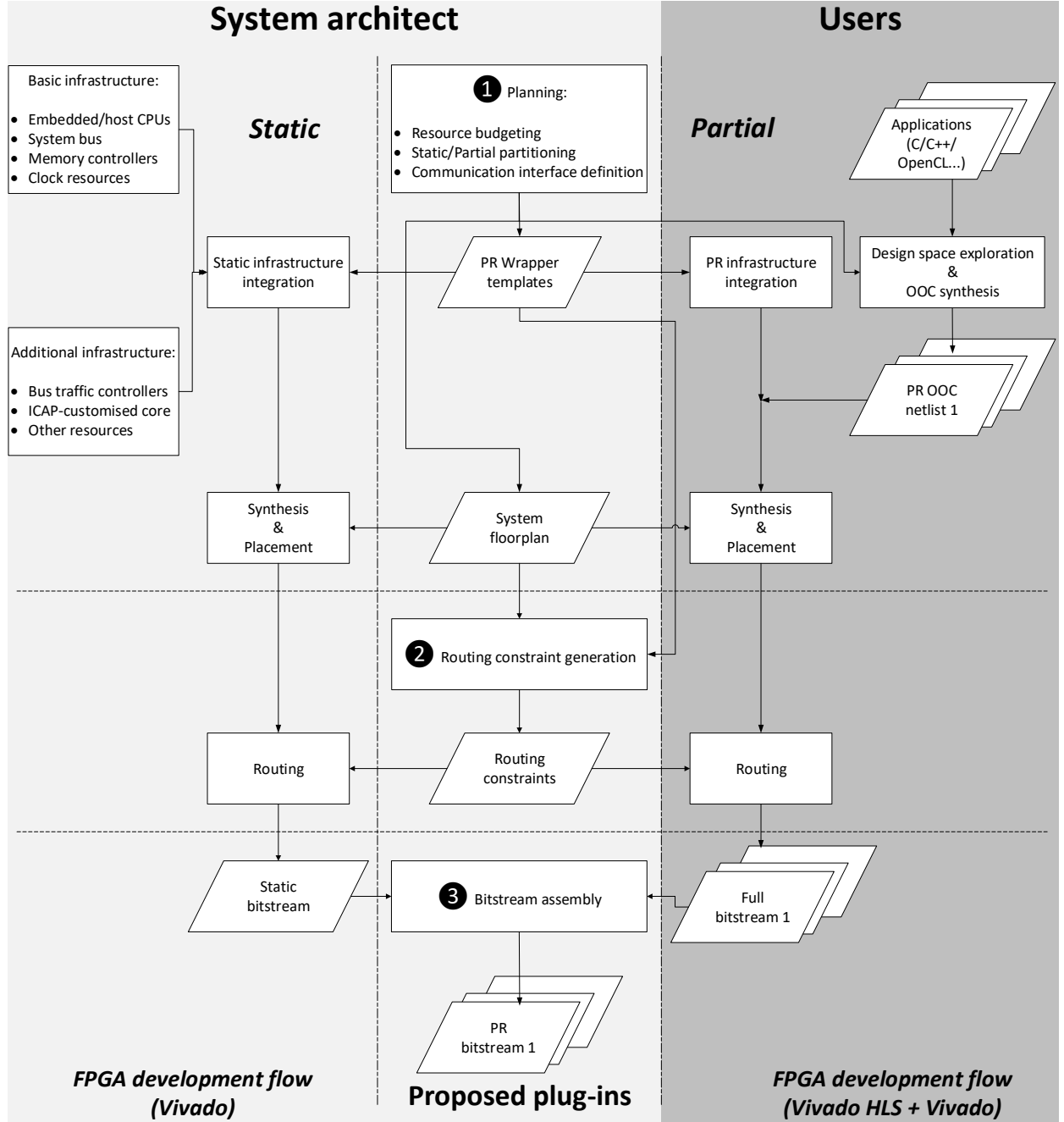
Figure 6: The proposed design methodology. The left part is performed once for a specific version of a shell by the system architect while the right part (with the darker background) is performed once for each module by FOS users.
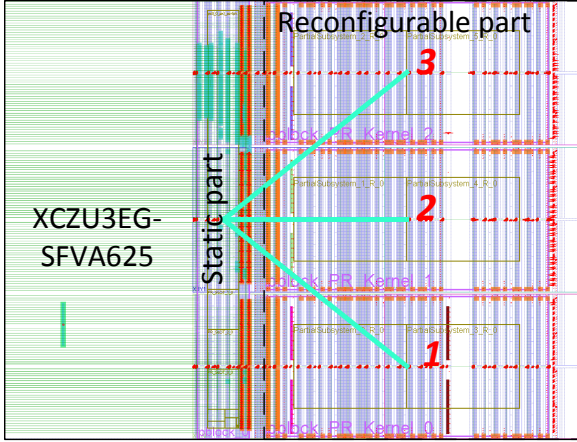
Figure 7: The physical shell implementation on the UltraZed and Ultra96 boards. This version has three PR regions which can host up to three FPGA applications simultaneously.
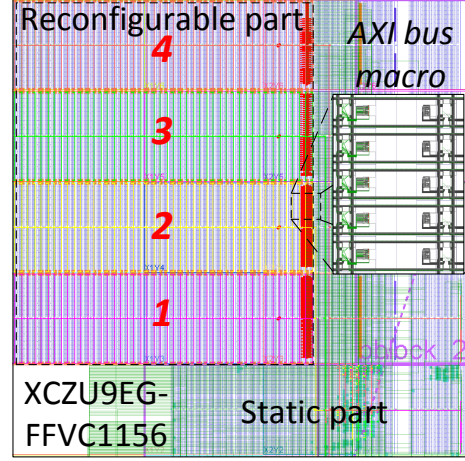


Figure 8: The physical shell implementation on the ZCU102 board. This version has four PR slots which can host up to four FPGA applications simultaneously.

(a) **Resource partitioning**: Based on the FPGA fabric layout and the number of resources available, a shell developer needs to split the FPGA fabric into two parts: 1) the static region for FPGA shell and 2) one or more reconfigurable regions for hosting hardware accelerators. This also defines the size of each PR region on which the High-Level Synthesis (HLS) tools [24] can perform the Design Space Exploration (DSE) for throughput optimisation [29]. The trade-off at this stage requires the developer to allocate the maximum amount of resources to the reconfigurable part while leaving sufficient resources to the static part for shell functionalities and future upgrades. This trade-off is case-dependent and requires a thorough understanding of the target FPGA device and system requirements. However, this step is only performed once for a particular system. The result of this process is 1) the static system floorplan and 2) bounding boxes of the reconfigurable regions.

(b) **Communication interface definition**: Select the protocol and data-width for communication between the static system and the partial regions based on the system requirements.

2. *Routing Constraint Generation*: The system floorplan and the PR interface template from the previous step are used to generate routing constraints. We can describe the implementation rules with the help academic PR frameworks (GoAhead [28] and a TCL library [30]) as TCL files for routing constraints automatically. These TCL constraints will then guide the routing stage of Vivado.

3. *Configuration Bitstream Generation*: The proposed flow results in *static bitstreams* for both shell and module designs. To compose *partial bitstreams* for accelerator modules from these bitstreams, we use the bitstream manipulation tool BitMan [31], see Section 4.1.3 for further details.

### 4.1.1 Shell Development

The static design starts with integrating basic infrastructure and additional infrastructure to the unified top-level design. In our designs the basic infrastructure includes 64-bit ARM Cortex-A53 CPU cores, AXI4 interconnects, memory controllers, Xilinx PR Decouplers for disabling/enabling static and module communication, clock management tiles for tuning module frequencies, and PR Module Interfaces. We can add other resources such as memory management or network communication IPs if required.

The PR Module Interface provides an AXI4-Lite Slave for control register access via the CPU and an AXI4 Master for memory access. This fixed interface between the hardware module and the static system is implemented using the available routing resources in the Zynq UltraScale+ FPGA devices. In particular, we keep this interface identical for all PR regions to serve relocatable hardware modules by pre-placing and pre-routing these communication signals in a constrained, predefined manner. This reassembles the bus macro approach which was very popular for Xilinx Virtex-II devices [32]. However, our flow can implement this without logic cost (in terms of LUTs used for the bus macro).

To keep the clocking resources of PR regions identical for relocatable hardware modules, we block all except for a defined subset of the `BUFCE_LEAF` primitives inside PR regions. This forces the router to use only a defined subset of these clock driver primitives that each drive a specific vertical clock spline which is ultimately connected to the flops, BRAMs, and DSPs in the region. However, we only route the clock for the PR regions this way. This means
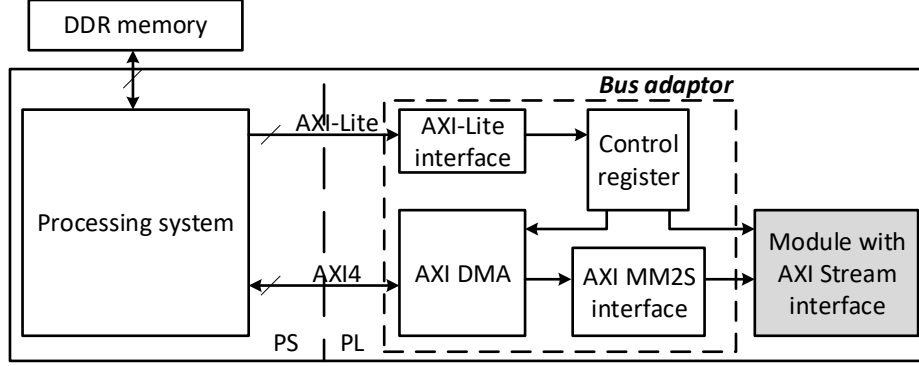
Figure 9: An example for bus virtualisation: the module has a 32-bit AXI-Lite interface and a 32-bit AXI Stream interface without DMA engine. In this case, the *bus adaptor* with AXI DMA and AXI MM2S IPs are chosen to carry out the communication with the rest of system.

when routing the static system, 1) we route the PR module clocks with `prohibit` constraints on the `BUFCE_LEAF` primitives, then 2) we remove these constraints and 3) incrementally route the rest of the system. This allows us to route additional clock nets as needed by the static system (e.g., for providing clocks to memory controllers or gigabit transceivers).

Finally, to prevent any static signal from violating PR regions, we insert a blocker macro. This blocker is non-functional, but uses all local wire resources inside the reconfigurable regions before routing the static system. We generate this blocker macro using GoAhead [28] or the TedTCL library [30] according to the system floorplanning.

The above steps result in the final static design shown in Figure 7 for the UltraZed/Ultra96 board and Figure 8 for the ZCU102 board. We can see that the PR region interfaces have the same relative physical positions and the even distribution of clock splines across all PR regions. Both systems support combining multiple adjacent regions for hosting larger monolithic modules. In this case, only one PR module interface will be used.

### 4.1.2 Bus Virtualisation

Operating a hardware accelerator needs communication with the host CPU to issue commands and to provide access to memory for data processing. A module kernel can use a wide range of bus widths, such as 32/64/128-bit width, and various bus protocols, such as AXI4 Master/Stream in the proposed approach. In particular HLS modules often, by default, include DMA engines for fetching data from memory and writing back results. However, this is not always the case with hand-crafted RTL or customised netlist accelerators.

We tackle this issue by providing another level of abstraction for bus interfaces between the FPGA applications and shells. For this, we selected the interface to provide the 32-bit AXI-Lite protocol and the 128-bit AXI4 protocol. Depending on the exact physical interface required by a module, we instantiate a module wrapper with a set of *bus adaptors* such that a module can communicate with the rest of the system as required by the individual FPGA modules. Figure 9 shows a *bus adaptor* being used to translate between different AXI bus standards. We can perform this wrapping at design time (where the wrapper is transparently instantiated when designing the module) or at runtime with partial reconfiguration (where the bus adaptor is a partially reconfigurable module located between the static system infrastructure and the accelerator module). The bus adaptor uses mostly IP components provided by FPGA vendor Xilinx [24]. These components are then automatically parameterised and integrated according to the specific interface requirements of the accelerator module. With this, shells can remain light-weight, operational, and unchanged while supporting a wide range of AXI interfaces.

The present FOS shells support up to 128-bit wide datapath for memory accesses because this is the native width to the ARM SoC. It is important to understand that a static acceleration system would also use such bus adaptors provided by the vendor, hence, our bus adaptors do not necessarily cause an additional overhead. The advantage of the here proposed *bus adaptor* concept is that an adaptor is only integrated into a module if needed and not speculatively provided by the shell.

We also use a bus adaptor to translate between AXI Master and AXI Stream protocols. We provide different versions of AXI Stream adaptors to be used depending on the AXI Stream channel width. A user can either re-compile their modules with a logical wrapper of the bus adaptor at design phase or stitch their modules with a pre-built binary of that bus adaptor at run-time. Figure 10 shows the implementation of the bus adaptor with its logic for interfacing an
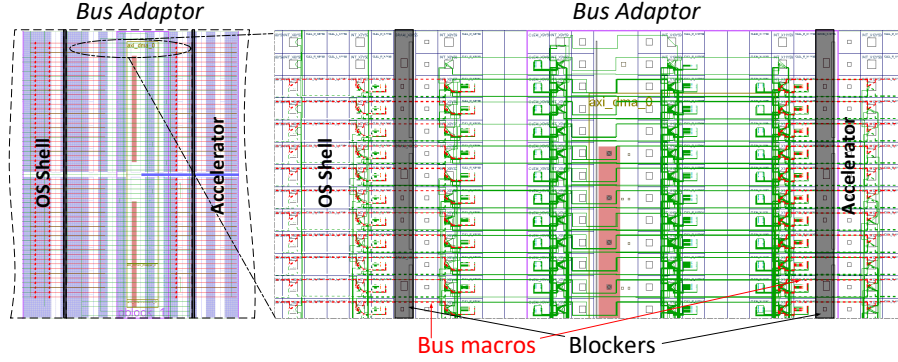
Figure 10: Implementation of a bus abstraction layer on UltraZed/Ultra96 platforms. The bus adaptor is provided as an implemented module bitstream and stitched to into the shell at run-time by using partial reconfiguration. The adaptor is a partial module that, in turn, interfaces to other partial modules. This technique avoids re-compiling bus adaptors but comprises an area overhead for a partial region to host a bus adaptor.
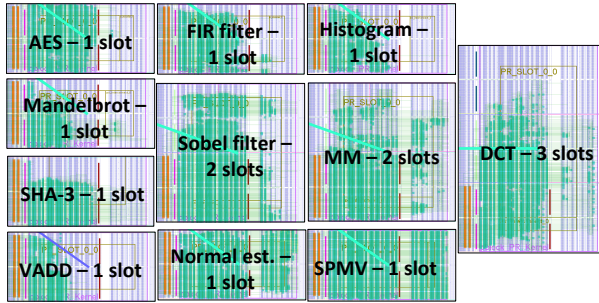


Figure 11: Compiled modules for Spector benchmark suite [33] and our in-house accelerators for Ultra-96 board.
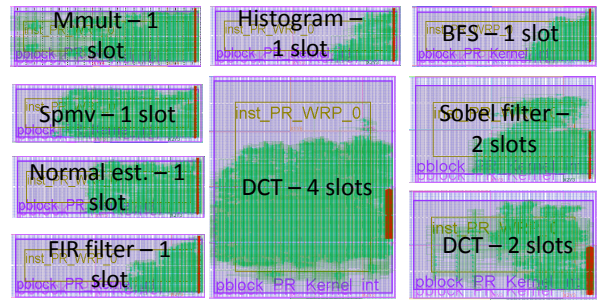


Figure 12: Compiled modules for Spector benchmark suite [33] accelerators for ZCU102 board.

accelerator with the AXI protocol, AXI MM2S, and AXI DMA services for a module which has a 32-bit AXI-Lite and 32-bit AXI Stream interface.

### 4.1.3   Module Compilation

The module design begins from either high-level language (HLL) source code (C, C++, OpenCL etc.) or a hardware description (RTL/netlist). In the case of HLL source code, we go through a High-Level Synthesis (HLS) step [24] to generate the RTL source codes. We then synthesise the resulting RTL source code in out-of-context (OOC) mode.

The PR Wrapper templates are used to create a minimal top-level placeholder for the module implementation. This temporary placeholder acts as sink/source connection points and substitutes the surrounding static system. We then integrate the module OOC netlist into this placeholder for the synthesis and placement stages.

We generate blockers (as TCL routing constraints) to enforce that all partial modules primitives and routing resources are following the strict implementation rules mentioned in Section 4.1 (by using GoAhead [28] or TedTCL library [30]). We place these blockers around the selected area to act as a fence for implementing hard module bounding box constraints. The blockers include routing tunnels for the communication to and from the temporary placeholder. The position of these tunnels matches exactly the tunnels used in the static design to implement the communication between static and partial areas.

As we implement a module in separation from the static system, the result generated by Vivado is a full configuration bitstream. We pass this full bitstream to BitMan [31] to extract the configuration data that corresponds to the module only as a partial bitstream. At run-time, BitMan manipulates those partial bitstreams to relocate modules to the desired partial region of the static system. Figure 11 and 12 shows the resulting modules for the Spector benchmark suite [33] and our in-house accelerators for Ultra-96 and ZCU102 boards, respectively.

## 4.2 Logical Hardware Abstraction

This is the primary layer between the hardware and software infrastructure. It is designed to hide the differences or changes in the hardware from the software, in order to detach the software infrastructure from the underlying hardware layer as much as possible. To achieve this conveniently, we propose describing the shell in terms of logical functionalities using a JSON file description with the following information (see listing 1 for an example):

1. Name of the shell
2. Bitstream name of the shell
3. Partial region:
   (a) Name of the partial region
   (b) Blanking bitstream for the partial region
   (c) AXI bridge decoupler address
   (d) Base address of an accelerator placed in the region

Listing 1: JSON description example of a shell.

```
1  {
2    "name": "Ultra96_100MHz_2",
3    "bitfile": "Ultra96_100MHz_2.bin",
4    "regions": [
5      {"name": "pr0", "blank": "Blanking_slot_0.bin", "bridge": "0xa0010000", "addr": "0xa0000000"},
6      {"name": "pr1", "blank": "Blanking_slot_1.bin", "bridge": "0xa0020000", "addr": "0xa0001000"},
7      {"name": "pr2", "blank": "Blanking_slot_2.bin", "bridge": "0xa0030000", "addr": "0xa0002000"}
8    ]
9  }
```

Similarly for the accelerator, we defined a JSON file description (see Listing 2) with the following details for accelerator programming and management:

1. Name of the accelerator
2. Bitstreams:
   (a) Name of the bitstream
   (b) Name of the shell it is compiled for
   (c) Type of AXI interface used (if not AXI4 master and slave)
   (d) Name and number of the PR regions it is compiled for
3. Register mappings:
   (a) Name of the HW register
   (b) Address offset at which the HW register can be access

Given that the control register map follows the standard Vivado HLS [24] interface (see Listing 3), an accelerator can have an arbitrary number of 32-bit registers. This allows us to build generic drivers for accelerators to relieve hardware developers from the responsibility of writing and integrating drivers. Hence, with this logical hardware abstraction, we can arbitrarily update the shell or accelerators with no need to recompile the Linux kernel or drivers. Section 4.3 will describe the driver and acceleration library interface.

The name of the PR region for each bitstream allows backward compatibility to the Xilinx PR flow, where we must compile an accelerator for each region (no relocation support). Moreover, the name of the accelerator is unique, but it may contain bitstreams of varying sizes corresponding to different acceleration implementations with the same functionality. The scheduler can later use these implementation alternatives to perform resource-elastic scheduling, i.e. dynamically change the resource allocation used by the accelerator based on the workload and available resources.

We then register these JSON descriptions for shell and accelerators into a JSON based registry to enable a centralised view of the available hardware to the upper software layers. This allows the application developers or the runtime system to request hardware based on just the name (a logical accelerator functionality) and corresponding input data, without needing any further information about the underlying hardware layer and accelerator implementation.

Note that we can automatically generate the JSON descriptor for the accelerator from the files generated by the Vivado HLS compilation flow [24]. Whereas, the shell developer must write the JSON description to allow the runtime system to manage the resource allocation and use generic drivers. However, this process is commonly required only once per system and can be omitted entirely when using our pre-built FOS shells.

Listing 2: JSON description of a vector add accelerator.

```
1  {
2    "name": "vadd",
3    "bitfiles": [
4      {"name": "vadd.bin", "shell": "Ultra96",
5       "region": ["pr0", "pr1"]},
6    ],
7    "registers": [
8      {"name": "control", "offset": "0"},
9      {"name": "a_op", "offset": "0x10"},
10     {"name": "b_op", "offset": "0x18"},
11     {"name": "c_out", "offset": "0x20"},
12   ]
13 }
```

Listing 3: Control bits for the accelerator.

```
1  // 0x00 : Control signals
2  //        bit 0 - ap_start (Read/Write/COH)
3  //        bit 1 - ap_done (Read/COR)
4  //        bit 2 - ap_idle (Read)
5  //        bit 3 - ap_ready (Read)
6  //        bit 7 - auto_restart (Read/Write)
7  //        others - reserved
```

### 4.3 Acceleration Interface Libraries

In general, an OS has to provide high-level APIs that can be used from existing software stacks to improve the accessibility of FPGAs to software developers (who are often non-FPGA experts). One such effort is the PYNQ [26] framework from Xilinx, which allows accessing FPGA accelerators through high-level Python APIs. However, the current implementation of PYNQ relies on a traditional development flow and contains many direct dependencies on artefacts produced by the Vivado compilation flow which restrict modular development [24, 26]. While it is possible to engineer around these shortcomings, the resulting system would suffer from code inflation and would be non-scalable due to its legacy support. Hence, we built a new light-weight acceleration interface libraries called Ponq and Cynq for Python and C++ languages, respectively. These libraries are built based on a modular development flow and provide access to static and dynamic FPGA accelerators via its generic drivers for programming accelerators, the Xilinx FPGA manager for partial reconfiguration [25], memory mapped I/O (MMIO) modules for direct access to accelerators and a data manager for contiguous physical memory allocation. Figure 13 shows the integration of Ponq and Cynq libraries into the FOS modular development flow and existing high-level software libraries.

Moreover, the libraries are backwards compatible with the standard Xilinx development flow, i.e. both the PR flow and the static acceleration environment, as we build them on top of the logical hardware abstraction layer. The libraries provide the following basic HAL functionality and generic drivers for hardware acceleration:

- Load an FPGA shell
- Load a partially reconfigurable accelerator
- Load a static accelerator
- Load and program an accelerator based on a logical function name
- Program an accelerator for execution via generic drivers
- Read and write calls to HW registers of the accelerators
- Contiguous physical memory allocation

### 4.4 Runtime and Multi-tenancy API

To support multi-tenancy, a runtime system is necessary to arbitrate access to reconfigurable resources between multiple users *transparently*. The common approach for allocation is to employ time-domain multiplexing with a run-to-completion model, but there also exists a few spatial domain scheduling mechanisms for FPGAs [34,35]. These spatial domain schedulers, however, only support a specific type of accelerator, such as OpenCL or DSL based accelerators. An operating system, in contrast, must support all types of accelerators and allow them to execute concurrently while
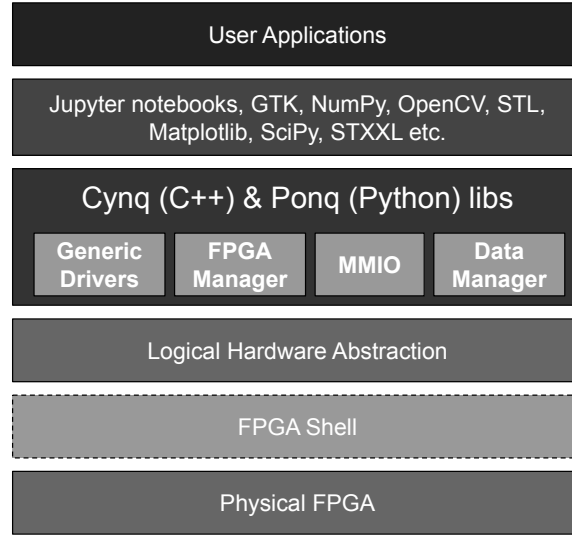
Figure 13: Cynq and Ponq libraries as an acceleration interface layer for static and dynamic acceleration on FPGAs.

using space-time domain scheduling. To make this possible, we need three main components: an API to a daemon, a programming model and a scheduler, as discussed in the following paragraphs.

### 4.4.1  Daemon API

To truly support multi-tenancy with portability, we need to design an API which can span across multiple languages and be portable to different OS kernels or a base processor system with ease. The two efficient ways to perform this Inter-process communication (IPC) are 1) message passing and 2) shared memory. In our platform we adopt the gRPC framework [36] which is a standard RPC framework with support for multi-languages. We use gRPC to send the acceleration requests from the client process to the daemon process, while the data is passed via shared memory to avoid additional latency of copying data (i.e. zero copy operation). The adoption of gRPC allows us to extend the runtime to accept acceleration requests from remote nodes in the future. The final interface exposed to the application developer in C++ and Python is shown in Listing 4 and 5. Note that each user can offload multiple data-parallel acceleration requests in a single RPC call to the daemon.

Listing 4: C++ daemon execution call example.

```
1  // Create a job
2  Job &job = jobs.emplace_back();
3  job.accname = "Partial_accel_vadd";
4
5  // Set accelerator parameters
6  job.params["a_op"] = a_op_phy_addr;
7  job.params["b_op"] = b_op_phy_addr;
8  job.params["c_out"] = c_op_phy_addr;
9
10 // Launch jobs
11 fpgaRpc.Run(job);
```

Listing 5: Python daemon execution call example.

```
1  # Create a job and set accelerator parameters
2  jobs = [{
3      "name": "Partial_accel_vadd",
4      "params": {
5          "a_op": a_op_phy_addr,
6          "b_op": b_op_phy_addr,
7          "c_out": c_op_phy_addr,
8      }]
9
10 # set accel parameters and run hardware unit
11 fpga_rpc.Run(jobs)
```

### 4.4.2  Programming Model

We achieve the dynamic resource allocation in the spatial-domain using resource-elasticity [35] in two forms: module replication and module replacement. However, to make this possible for all types of accelerators we allow applications to expose data-parallelism to the scheduler, i.e. an application developer can choose to express an acceleration job into a varying degree of parallelism appropriate for the application. A common example of this is chopping the image into
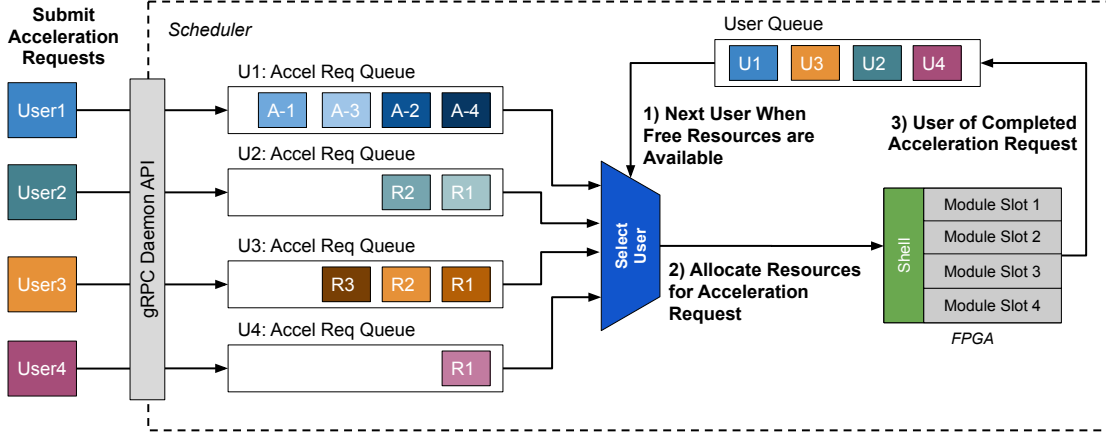
Figure 14: Scheduler organisation for resource allocation between different users and their data-parallel acceleration requests.
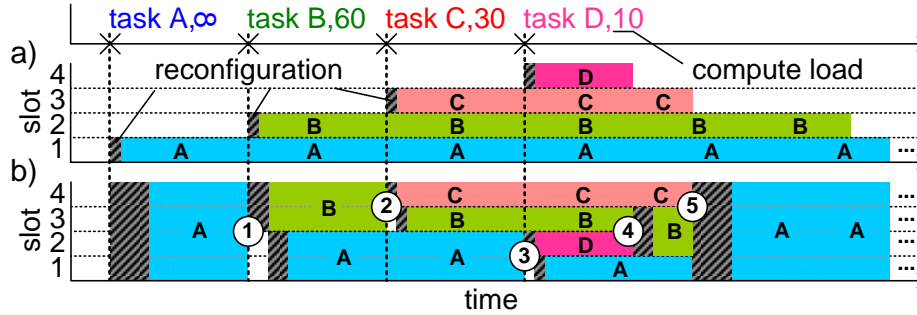


Figure 15: Resource allocation for kernels (tasks) A, B, C and D in time when using a) Standard fixed module scheduling and b) Resource-elastic scheduling on a 4 PR region FOS FPGA shell. The circled events highlight cases where resources are needed to accommodate new arriving tasks (①, ②, ③) or cases where tasks complete (④, ⑤).

multiple parts for image-processing accelerators. Note that this is analogous to a software developer deciding for the number of threads for efficient execution on a multi-core system. The runtime is then responsible for executing those parts i) in parallel, ii) use a better implementation, or iii) in case of a exceeding FPGA resource capacity perform time-domain multiplexing for the resources.

### 4.4.3  Scheduling

The scheduler maintains a queue of the users and performs round-robin scheduling between users at a coarse granularity of data-parallel acceleration requests, as shown in Figure 14. Each user has an individual queue of acceleration requests. Each request in this queue is independent of other requests in the queue and can execute in parallel and any order. At the end of each acceleration request, the scheduler relinquishes the accelerator and selects an acceleration request from the next user in the queue. This scheme implements a cooperative scheduling policy (by breaking down a job into fine-grain run-to-completion acceleration requests) where each request includes fetching operands and writing back results to main (DDR) memory. This corresponds directly to the OpenCL programming model where work-groups can be executed in any order.

In the case there is no other user, the scheduler executes requests from the same user in parallel and attempts to use the biggest module (assumed to be the fastest, i.e. Pareto-optimal) to maximise the utilisation and performance. Moreover, the scheduler avoids partial reconfiguration and reuses an accelerator if it is already available on-chip. This allows multiple different applications that require acceleration of the same functionality to share an accelerator in time without paying an additional penalty or user effort.

Consequently, a single task execution call may execute on multiple accelerators in parallel or use an implementation alternative or share an accelerator with other tasks in time, during its execution lifetime. All these modes can be arbitrarily used without an application being aware of other tasks and types of accelerators it executes on. Hence, the runtime system is responsible for dynamically arranging the loading and unloading of these heterogeneous accelerators

Table 1: Available resources for acceleration on the ZCU102 platform and the UltraZed & Ultra96 platforms. The version on ZCU102 has 4 PR regions in total, while the other platforms provide 3 PR regions in total.

| Resources on ZCU102 | Number of resources for 1 PR region | Chip utilisation per PR region (%) | Total chip utilisation for accelerators (%) |
| --- | --- | --- | --- |
| CLB LUTs | 32640 | 11.70 | 46.80 |
| CLB Regs. | 65280 | 11.90 | 47.60 |
| BRAMs | 108 | 12.10 | 48.40 |
| DSPs | 336 | 13.30 | 53.20 |
| Resources on Ultra96 & UltraZed | Number of resources for 1 PR region | Chip utilisation per PR region (%) | Total chip utilisation for accelerators (%) |
| CLB LUTs | 17760 | 25.17 | 75.51 |
| CLB Regs. | 35520 | 25.17 | 75.51 |
| BRAMs | 60 | 27.78 | 83.33 |
| DSPs | 96 | 26.67 | 80 |

(written in C, C++, OpenCL or RTL) transparently from the user. Figure 15 shows an example of how such resource allocation can allow maximising the FPGA utilisation and performance compared to standard fixed-module scheduling policies.

## 5   Evaluation

There are five primary dimensions onto which we can evaluate an FPGA operating system: 1) FPGA resource overhead, 2) software stack overhead, 3) available memory performance, 4) level of modularity and 5) application performance. We detail and discuss the performance of FOS on each dimension for Ultra-96 and ZCU102 boards.

### 5.1   FPGA Resource Overhead

#### 5.1.1   FPGA shell

The resources used by an FPGA shell have a direct impact on the FPGA resources available for user hardware accelerators. Hence, it is important to minimise the overhead as much as possible. Table 1 shows the resources available for hardware acceleration when using FOS on ZCU102, Ultra-96, and UltraZed boards. For the ZCU102 board (an MPSoC development kit from Xilinx) around 50% of the resources are available for user acceleration whereas on a small IoT category Ultra-96 board it is about 80%. This is because the layout of the ZUC102 chip is irregular, limiting the available resources when supporting relocatable modules (see Figure 8). With a regular resource layout like in Ultra-96 (see Figure 7), we can maximise the resource allocation for relocatable modules and considerably reduce the resource overhead of a FOS shell. However, it is important to note that we do not entirely waste unused resources; they are available for future extensions as well as to implement other host system components, for example static accelerators or I/O functionalities.
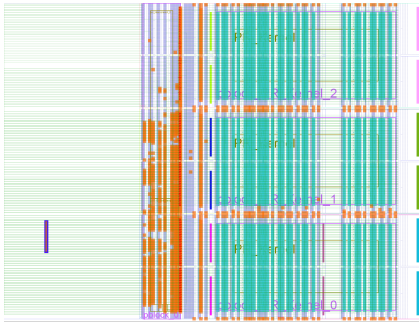
#### 5.1.2   Bus Virtualisation

In order to achieve modularity at the hardware interface layer, bus virtualisation is vital. However, to be able to dynamically load the interconnect wrappers implies pre-allocation of a partial region. Table 2 shows the overhead of this pre-allocation of resources. We can identify that the unused resources are only about 448LUTs (18% of pre-allocation) when we change the interconnection interface and protocol considerable, such as from AXI Stream to AXI master and slave. In particular, for large FPGAs this overhead is negligible. However, when using small FPGAs such as Ultra-96 or performing minor changes to the interface (e.g., changing bus width), the overhead of dynamic bus virtualisation is considerable. Hence, in such scenarios, we recommend using compile-time bus wrappers.
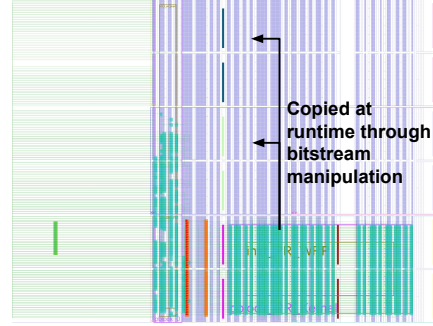
Table 2: Resource overheads for bus virtualisation at the logical and physical levels.

| Module Interface | Shell Interface | Bus adaptor's services | Primitives | Resource overhead Logical Level | Physical Level |
|---|---|---|---|---|---|
| 32-bit AXI-Lite & 32-bit AXI4 Master | 32-bit AXI-Lite & 128-bit AXI4 Master | AXI Interconnect | LUTs | 153 | 2400 |
| | | | FFs | 284 | 4800 |
| | | | BRAMs | 0 | 12 |
| 32-bit AXI -Lite & 32-bit AXI Stream | 32-bit AXI-Lite & 128-bit AXI4 Master | Control reg., AXI MM2S, & AXI DMA | LUTs | 1952 | 2400 |
| | | | FFs | 2694 | 4800 |
| | | | BRAMs | 2.5 | 12 |



(a) Xilinx PR compilation flow result.



(b) FOS compilation flow result.

Figure 16: Place and route result of Black Scholes accelerator [37] for Xilinx PR flow and FOS. Note, FOS generates relocatable bitstream with the help of BitMan [31] for other PR regions at runtime.

## 5.2 Software Stack Overhead

The software stack of an FPGA OS incurs two types of latencies: compile-time and runtime. Compile-time latency relates to the time taken to compile hardware accelerators with all the additional constraints for partial reconfiguration in the implementation phase, while the other relates to the overhead caused by intermediate layers in the software stack during accelerator execution.

### 5.2.1 Compilation Latency

The standard Xilinx partial reconfiguration (PR) flow requires compiling both static and accelerator designs together and, more importantly, it needs to perform place and route (P&R) and generate a new bitstream for each partial region. This leads to additional latency compared to our decoupled compilation flow, where we first generate a *full-static* bitstream with Vivado and then a *relocatable* partial bitstream with BitMan [31] for all regions. To measure this, we used accelerators of three different types of size: sparse (AES), medium (Normal est. [33]) and dense (Black Scholes [37]). The utilisation for each is 33%, 63% and 81%, respectively. Figure 11 shows the AES and Normal est. modules while the Black Scholes module is shown in Figure 16 with a comparison to design generated by Xilinx PR flow.

Table 3 shows the latency breakdown for place and route and bitstream generation for both Xilinx PR flow and FOS compilation flow on Ultra-96. The results show that per region P&R latency is higher for FOS as it adds additional constraints for relocatability, however, when compiling for multiple regions (i.e. 3 for Ultra-96 platform) it outperforms traditional compilation flow by up to 2.34× if a single slot accelerator can be duplicated to scale up the system. Overall, when increasing the number of partial regions, the compilation flow latency of FOS stays constant, whereas the latency of the Xilinx PR flow increases linearly. This relationship turns into an exponential increase in compilation time when compiling several applications with standard Xilinx PR flow, as it needs to compile *each module for each partial region*. Hence, it is important to adopt a scalable and modular PR flow when targeting multiple applications and shell versions which are commonly used in cloud environments.

Table 3: Total place and route (P&R), and bitstream generation latency for AES, Normal Est. [33], and Black Scholes accelerator [37] when compiling for all three partial regions on Ultra-96 shell. Evaluation is conducted using Vivado 2018.2.1 on Intel core i7-4930K CPU running at 3.4 GHz with 64 GB of RAM.

| Applications | Region Util. | Xilinx PR | | | FOS | | | Speed Up |
|---|---|---|---|---|---|---|---|---|
| | | P&R (s) | Bitgen (s) | Total (s) | P&R (s) | Bitgen (s) | Total (s) | |
| AES | 33% | 429.40 | 176.19 | 605.59 | 284.18 | 64.06 | 348.24 | 1.74 × |
| Normal Est. | 63% | 747.75 | 201.21 | 948.96 | 387.41 | 70.09 | 457.50 | 2.07 × |
| Black Scholes | 81% | 1296.26 | 231.27 | 1527.53 | 574.56 | 77.11 | 651.67 | 2.34 × |

Table 4: Execution overhead caused by various software layers.

| Software Layer | Latency (ms) |
|---|---|
| Initialize gRPC (once) | 12.20 |
| JSON parsing (once) | 2.27 |
| gRPC Call to Daemon | 0.71 |
| Scheduler | 0.02 |

### 5.2.2 Runtime Execution Overhead

The runtime overhead occurs because of the four main steps performed when using the FOS software stack: i) initialisation of gRPC server, ii) parsing of JSON files for accelerators and shells, iii) gRPC call to the daemon and iv) scheduling latency. Table 4 details the latency of each step. The first two steps (i and ii) have heavy dependencies on I/O as they use the network and file system, respectively. This leads to latencies in the range of milliseconds. However, this overhead is amortised over time as we perform steps i) and ii) only once at system start. The gRPC call to the daemon goes through many levels of the Linux stack (processes) before reaching the FOS multi-tenancy daemon, leading to a latency of about a millisecond. We can speed up the gRPC call by reducing the Linux timer interrupt period to achieve a better response time from the Linux kernel for quick turn around between processes if required. The scheduling latency is in the range of micro-seconds and comparable to standard CPU schedulers. Moreover, the scheduler is event driven and executed only when an accelerator finishes or a new acceleration request arrives (due to its cooperative nature) rather than at every timer interrupt like preemptive scheduling.

### 5.3 Memory Performance

One important characteristic of an FPGA operating system on the hardware acceleration is the memory bandwidth it can provide given the shell implementation (interconnect to the memory). Hence, we evaluated the available throughput on different AXI ports made available to partial regions as well as their combined throughput on Ultra-96 and ZCU102 boards with varying burst sizes from the PL using the memory evaluation kit [38].

Figure 17 shows the breakdown of the read and write throughput of each AXI port and the total accessible bandwidth for the Ultra-96 board, respectively. On average, there is an even split of read and write bandwidth with a throughput of 530 MB/s for each read and write operation. The aggregated read-write throughput of the individual AXIs is about 1060 MB/s and, when activated at the same time, all available AXI ports collectively achieve up to 3187 MB/s. This translates to about 25% and 74% of the theoretical DDR peak throughput when using AXI ports individually and concurrently, respectively.

Figure 18 shows the breakdown of read and write throughput on the ZCU102 board. Each AXI port achieves a throughput distribution between read and write transactions 1600 MB/s each. The total throughput is 3200 MB/s for individual AXI ports and 8804 MB/s when using all AXIs together, as shown in Figure 18. We expect that sub-linear improvement in the total throughput when using all AXI ports (HP0-3) concurrently is caused by the row pollution and AXI interconnect multiplexing in the memory controller.
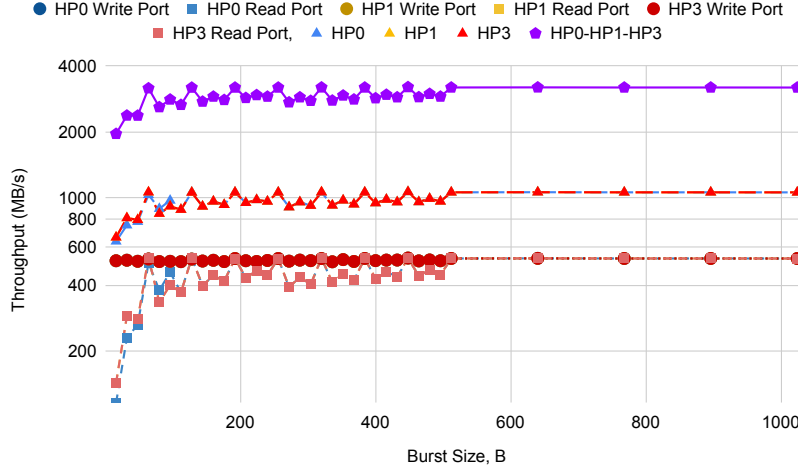
Figure 17: Memory throughput for varying burst sizes on the duplex AXI ports (HP0, HP1, and HP3) available to the hardware accelerators in PR regions 0 to 2 of Ultra-96 FOS platform at 100 MHz.
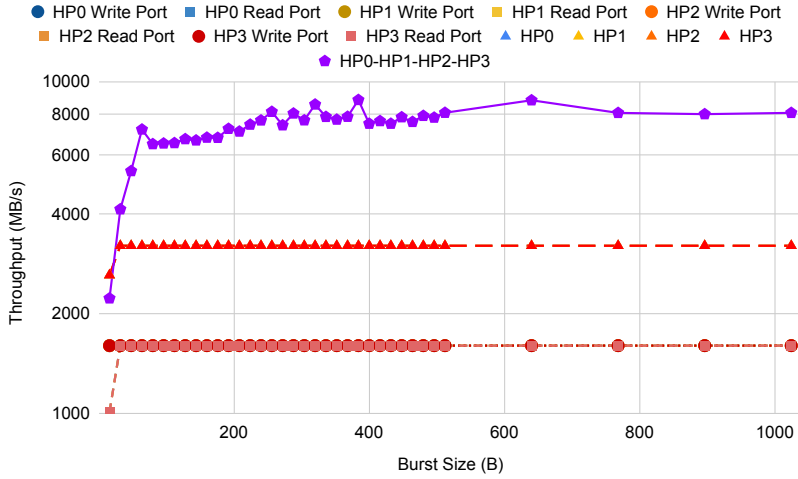


Figure 18: Memory throughput for varying burst sizes on the duplex AXI ports (HP0 to HP3) available to the hardware accelerators in PR regions 0 to 3 of ZCU102 FOS platform at 100 MHz.

## 5.4   Quantifying Modularity

Compared to the standard FPGA development flow, where a change in the shell can mean recompilation of all system components (both hardware and software), the modular FOS FPGA development flow provides the freedom to *update individual components* without recompilation of other components which may take hours [7, 24, 25]. This means that FOS only has to pay compilation and re-initialisation latency (as shown in Table 5) for the component which is being changed given that it does not change the interfaces defined between them.

From Table 5, we can identify that this avoidance of recompilation allows changing the shell at runtime with additional functionalities or bug fixes costs less than 21 ms on Ultra-96 (IoT device) and 99 ms on ZCU102 (Xilinx MPSoC development kit)[3]. Similarly, swapping an accelerator implementation is easy and includes only the partial reconfiguration latency, because FOS provides generic drivers. In contrast to FOS, the standard flow would require generating/writing new drivers and re-installing them separately. Changing the kernel involves the biggest re-initialisation latency, as this requires a system reboot which takes 66 seconds in total, including I/O setup (for keyboard, mouse, Wi-Fi, monitor and webcam) on Ultra-96. However, this still avoids the need to update user software binaries for re-integration as required in the standard PetaLinux flow [25].

---

[3]This assumes that the shell bitstream is pre-compiled as the recompilation of shell itself is unavoidable for all systems.

Table 5: Re-initialisation latencies for component change on FOS platforms.

| Component Updated | U-96 Latency (ms) | ZCU102 Latency (ms) |
|---|---|---|
| Accelerator | 3.81 | 6.77 |
| Shell | 20.74 | 98.4 |
| Runtime | 15.2 | 15.2 |
| Kernel | 66000 | 15760 |

Overall, modularity of FOS removes the recompilation latencies and re-development steps which cost in the range of hours [7] and bring down the component change latency by two-orders of magnitude compared to the standard development flow while supporting all the features required from an FPGA OS.

## 5.5  Application Case Study

We evaluated a case-study in two different environments: 1) single-tenant but multiple partial regions and 2) multi-tenant with dynamic offloading. All the accelerators used in this case study operate at 100 MHz.

### 5.5.1  Single-tenant with Multiple Partial Regions

To evaluate the benefits of multiple partial regions and the ability to replicate accelerators dynamically, we selected OpenCL accelerators from the Spector benchmark suite [33] for two reasons: 1) it comprises a wide range of application behaviour from multiple domains, and 2) it is annotated with HLS pragmas to generate implementation alternatives without modifying the source code. Figure 19 shows the results of the execution latencies with a varying amount of resources available for acceleration on the ZCU102 platform. Most of the benchmark applications show an almost *linear* performance improvement when replicated across multiple partial regions. In particular, DCT benefits from the ability to switch to bigger module implementation (by using more data buffers and a larger unrolling factor) and achieves a *super-linear* performance improvement of $3.55\times$ for $2\times$ resources.

To understand the effects of exposing more parallelism than available on the FPGA platform, we conducted the experiments on Ultra-96 using compute-bound (Mandelbrot and Black-Scholas [37]) and memory-bound (Sobel [39]) applications which are more sensitive to reconfiguration overhead due to their smaller execution latencies than Spector benchmarks [33]. Figure 20 shows the results in which the number of requests dictates the amount of parallelism exposed to the runtime. The performance improves almost linearly until it hits the number of available partial regions (3 in the case of Ultra-96 platform), after which the performance tends to stagnate (see Figure 21). This is because the scheduler uses time multiplexing to provide the illusion of an unlimited number of regions, leading to behaviour similar to multi-threading on CPUs. In particular, cases, where the number of requests is a multiple of the number of physical accelerators, perform better than others as it avoids bottlenecks caused by the (leftover) pending requests at the end of execution.

Overall, this highlights that FOS can help to improve performance even when using a single application along with its other benefits of modularity and developer productivity.

### 5.5.2  Multi-tenant Dynamic Offload

To understand the performance changes when using multiple applications at the same time, we execute the Mandelbrot and Sobel [39] applications concurrently on the Ultra-96 platform. Note that individual applications are not aware of other applications executing on the platform. In particular, the accelerators are written in C (mandelbrot) and OpenCL (sobel), demonstrate support for multiple different languages used concurrently.

Figure 22 shows the execution latencies relative to 1-Mandel$\times$1-Sobel scenario with a varying amount of acceleration requests. As we can see, the execution latencies tend to decrease with an increase in the number of requests as in the standalone case. However, here the optimal performance is achieved at 3-Mandel$\times$1-Sobel rather than 3-Mandel$\times$3-Sobel. This is because of two reasons: 1) adding more Sobel units reduces memory performance due to row-bank pollution and 2) multiple request from different application induces rapid reconfiguration latencies. Regardless of this behaviour, it is important to note that if each application takes a greedy decision to request the highest amount of parallelism suited for the application based on the standalone results, the system can still achieve a near-optimal performance resulting in 46% improvement over 1-Mandel$\times$1-Sobel by dynamically reallocating resources using the same accelerators.
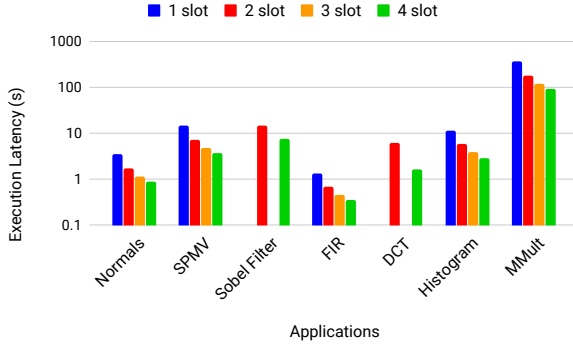
Figure 19: Execution latencies of accelerators from the Spector benchmark suite on ZCU102 platform.
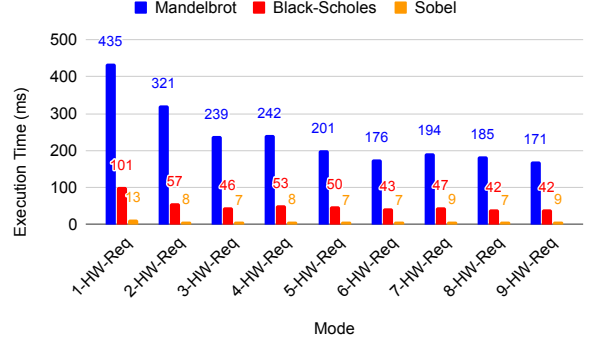


Figure 20: Execution latencies of Mandelbrot, Black Scholes (European option) [37], and Sobel [39] when executing concurrently with varying amount of hardware requests on Ultra-96 platform.
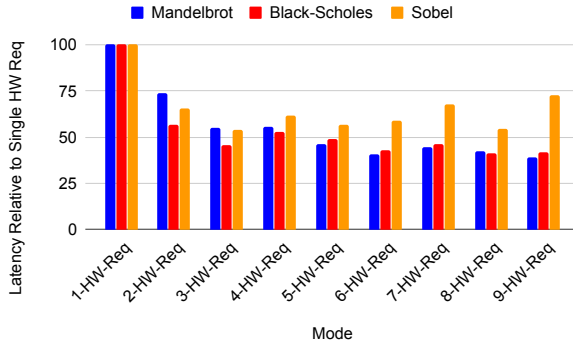


Figure 21: Relative execution latencies of Mandelbrot, Black Scholes (European option) [37], and Sobel [39] applications when exposing varying amount of parallelism to process a frame on Ultra-96 platform.
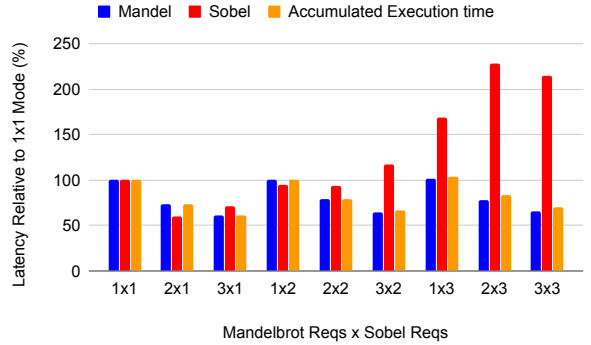


Figure 22: Relative execution latencies of Mandelbrot and Sobel [39] when executing concurrently with varying amount of HW requests on Ultra-96 platform

## 6    Conclusion

In this paper, we described the underlying concepts and abstraction layers involved in building a modular FPGA operating system that can adapt to dynamic workloads. The resulting FPGA operating system – FOS, provides support for traditional as well as multi-tenancy environments with low overhead and easy to use software interfaces. The dynamic resource allocation capabilities of FOS, allows FOS to share multiple FPGA accelerators transparently between multiple users in the time and the spatial domain as well as the ability to switch between accelerator implementations on the fly. Our evaluation shows that FOS can speed up the compilation time by up to 2.34x and avoid standard recompilation requirements to reduce the update latencies by over 100x with its modular FPGA development flow. The overheads caused by the hardware and software layers are minor and recovered by the ability to schedule resources dynamically in both single and multi-tenant environments. Overall, FOS directly caters the needs of upcoming FPGA systems which are being deployed at scale (cloud and edge) and allows systems to be more maintainable, adaptable and accessible, and benefiting both FPGA and application domain experts. As a distinct feature, FOS provides an application-centric view to the developers, by hiding most of the complexity encountered when using a heterogeneous CPU-FPGA acceleration system with Linux back-end. With this, FOS is empowering a larger FPGA user community to implement complex and scalable heterogeneous systems.

# References

[1] E. Lubbers and M. Platzner. ReconOS: An RTOS Supporting Hard-and Software Threads. In *FPL*, 2007.

[2] W. Peck et al. Hthreads: A Computational Model for Reconfigurable Devices. In *FPL*, 2006.

[3] H. So and R. Brodersen. *BORPH: An Operating System for FPGA-Based Reconfigurable Computers*. PhD thesis, EECS Department, University of California, Berkeley, Jul 2007.

[4] W. Wang et al. pvFPGA: Accessing an FPGA-based Hardware Accelerator in a Paravirtualized Environment. In *CODES+ ISSS*, 2013.

[5] J. Weerasinghe et al. Network-Attached FPGAs for Data Center Applications. In *FPT*, 2016.

[6] S. Byma et al. FPGAs in the Cloud: Booting Virtualized Hardware Accelerators with OpenStack. In *FCCM*, 2014.

[7] A. Khawaja et al. Sharing, protection, and compatibility for reconfigurable fabric with amorphos. In *OSDI*, 2018.

[8] F. Chen et. al. Enabling FPGAs in the Cloud. In *CF*, 2014.

[9] L. Wirbel. Xilinx SDAccel: A Unified Development Environment for Tomorrows Data Center. *The Linley Group Inc*, 2014.

[10] V. Kathail et al. SDSoC: A Higher-level Programming Environment for Zynq SoC and Ultrascale+ MPSoC. In *FPGA*, 2016.

[11] A. Vaishnav et al. A Survey on FPGA Virtualization. In *FPL*, 2018.

[12] Amazon Web Services. AWS EC2 FPGA Hardware and Software Development Kit, 2009. Accessed: 2017-12-04.

[13] A. Putnam et al. A reconfigurable fabric for accelerating large-scale datacenter services. In *ISCA*, 2014.

[14] Xilinx. Reconfigurable Acceleration in the Cloud, 2017.

[15] S. A. Fahmy, K. Vipin, and S. Shreejith. Virtualized FPGA Accelerators for Efficient Cloud Computing. In *CloudCom*, 2015.

[16] M. Vesper, D. Koch, and K. D. Pham. PCIeHLS: an OpenCL HLS Framework. In *FSP*, 2017.

[17] T. Xia et al. Hypervisor Mechanisms to Manage FPGA Reconfigurable Accelerators. In *FPT*, 2016.

[18] M. Happe et al. Preemptive Hardware Multitasking in ReconOS. In *ARC*, 2015.

[19] Xilinx. Platform overview — xilinx runtime 2019.1 documentation, 2019.

[20] Xilinx. UG909 - Vivado Design Suite User Guide: Partial Reconfiguration, June 2019.

[21] J. Weerasinghe et al. Enabling FPGAs in Hyperscale Data Centers. In *UIC-ATC-ScalCom*, 2015.

[22] O. Knodel et al. Virtualizing Reconfigurable Hardware to Provide Scalability in Cloud Architectures. *RECATA*, 2, 2017.

[23] Q. Zhao et al. Enabling FPGA-as-a-Service in the Cloud with hCODE Platform. *IEICE Transactions on Information and Systems*, 2018.

[24] T. Feist. Vivado design suite. *White Paper*, 5:30, 2012.

[25] Xilinx. UG1144 - PetaLinux Tools Documentation Reference Guide, 2019.

[26] Xilinx. PYNQ, 2019.

[27] K. D. Pham et al. Zucl 2.0: Virtualised memory and communication for zynq ultrascale+ fpgas. In *FSP*, 2019.

[28] Christian Beckhoff, Dirk Koch, and Jim Torresen. GoAhead: A Partial Reconfiguration Framework. In *FCCM*, 2012.

[29] F. B. Muslim, L. Ma, M. Roozmeh, and L. Lavagno. Efficient FPGA Implementation of OpenCL High-Performance Computing Applications via High-Level Synthesis. *IEEE Access*, 5:2747–2762, 2017.

[30] M. Vesper. *Dynamic Stream Processing Pipelines on FPGAs Examplified on the PostGreSQL DBMS*. PhD thesis, The University of Manchester, 2018.

[31] K. D. Pham, E. Horta, and D. Koch. BITMAN: A Tool and API for FPGA Bitstream Manipulations. In *DATE*, 2017.

[32] P. Lysaght et al. Invited Paper: Enhanced Architectures, Design Methodologies and CAD Tools for Dynamic Reconfiguration of Xilinx FPGAs. In *FPL*, 2006.

[33] Q. Gautier et al. Spector: An OpenCL FPGA Benchmark Suite. In *FPT*, 2016.

[34] M. Asiatici et al. Virtualized Execution Runtime for FPGA Accelerators in the Cloud. *IEEE Access*, 2017.

[35] A. Vaishnav et al. Resource Elastic Virtualization for FPGAs using OpenCL. In *FPL*, 2018.

[36] Google. grpc, 2019. Accessed: 2019-12-19.

[37] L. Ma, F. B. Muslim, and L. Lavagno. High performance and low power monte carlo methods to option pricing models via high level design and synthesis. In *EMS*, pages 157–162, 2016.

[38] K. Manev. Zynq ultrascale+ memory speed test, 2019.

[39] Xilinx. Xilinx sdaccel examples, 2019.