

ESP Accelerator Specifications

Copyright © 2011-2020 Columbia University
System Level Design Group

January 6, 2025



Contents

- 1 Introduction** **2**
- 1.1 Conventions 3

- 2 ESP Accelerator Specifications** **4**
- 2.1 Accelerator Model 4
- 2.1.1 Accelerator Configuration 5
- 2.1.2 Private Local Memory 5
- 2.2 DMA Transactions 6
- 2.2.1 DMA User Field 7

1. Introduction

This document describes the signal-level protocol specification of an ESP accelerator. The guide is intended for RTL designers who wish to implement a native ESP accelerator using a hardware-description language, such as SystemVerilog, or VHDL. Any accelerator that complies with the protocol specification described in this guide can be integrated in ESP and leverage all platform services through the *ESP accelerator socket*.

This document does not describe the *ESP third-party accelerator flow*. The latter enables the seamless integration of an existing accelerator IP leveraging an ARM AMBA open standard interface.



1.1 Conventions

- **bitwidth:** number of bits. This is typically associated to a signal, or to a unit of data.
- **token:** the unit of input or output data transferred between the accelerator and the ESP socket. The bitwidth of a token depends on the particular accelerator and may vary across different transactions over a bus or data channel.
- **beat:** the unit of data transferred on a bus, or a data channel. The bitwidth of one beat depends on the particular implementation of the accelerator (e.g. *dma32* or *dma64*) and not on the data type of the input or output token in a transaction. Therefore, for any given implementation of an ESP accelerator, the bitwidth of a beat is constant.
- **flit:** the unit of data transferred over a network-on-chip (NoC). For ESP accelerators, the bitwidth of a flit is equal to the bitwidth of a beat plus two bits. These additional bits indicate if the flit is the head, part of the body, or the tail of a packet.
- **packet:** a set of flits transferred in an ordered sequence across the NoC. Packets must have one header flit, one tail flit and as many body flits as necessary. Single-flit packets have just one flit with both *head* and *tail* bits set. A packet that is granted a link of the NoC will traverse such link from head to tail not interleaved with another packet.
- **initiator** or **master:** a component that can initiate a transaction over a bus, or a NoC.
- **target** or **slave:** a component that servers a transaction initiated by a master.
- **latency-insensitive channel (LIC):** a bundle of data wires and two control wires named *ready* and *valid*. During read transactions, the master drives the *ready* control signal, while the slave drives the data and the paired *valid* control signal. Roles are inverted for write transactions. A beat is transferred over a LIC when both *ready* and *valid* are set. Both master and slave have the ability to delay the transfer of a beat for as many cycles as necessary. For more on latency-insensitive channels please refer to [Carloni, 2015].
- **CSR:** configuration and/or status register.
- **DMA:** the acronym for *direct-memory access*. When referring to an ESP accelerator, the term DMA refers to the mechanism used by the accelerator to access data in the system memory hierarchy. A DMA transaction initiated by an accelerator in ESP may be accessing external memory *directly* or by mediation of the ESP cache hierarchy. The selection is managed by software at run time and is transparent to the accelerator.
- **PLM:** the accelerator's private local memory, composed of a set of SRAM bank groups customized for the accelerator's datapath.

2. ESP Accelerator Specifications

2.1 Accelerator Model

The block diagram of Figure 2.1 illustrates the ESP accelerator socket and shows the three main set of signals at the interface of an ESP accelerator: read and write port for data transfers through DMA requests, configuration port and interrupt line.

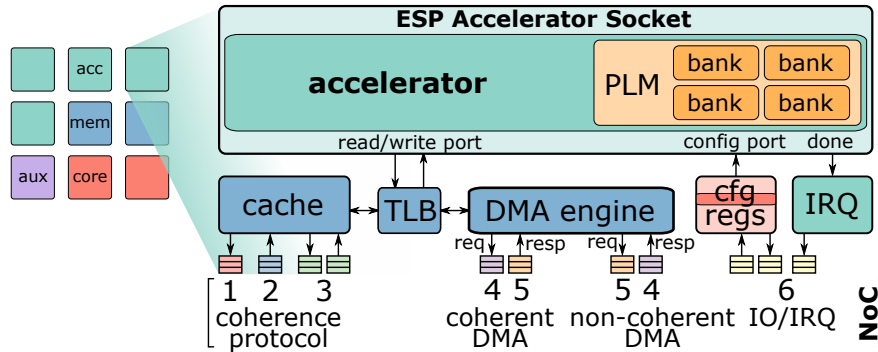


Figure 2.1: Block diagram of ESP accelerator tiles

A typical ESP accelerator is composed of three control blocks (*configuration*, *load*, *store*), one or more computation blocks and a customized private local memory (PLM).

Once configuration registers are valid, the configuration block activates the other components. The *load* module initiates the first DMA transaction to fetch the input data, or a portion of it, from the system memory hierarchy into the PLM. Next, the computation blocks process the available input and produce the corresponding output. Finally, the *store* block writes back the output to the system memory hierarchy with a DMA request. A single accelerator invocation from software typically results into multiple iterations of *load*, *compute* and *store* phases, therefore we recommend implementing a portion of the PLM as a set of ping-pong buffers to enable pipelining. Depending on the particular task and accelerator implementation, this strategy may significantly improve the overall accelerator throughput by masking most of the time for data transfers with the overlapping computation steps.

The above accelerator model corresponds to what ESP automation provides through any of the available high-level synthesis flows. RTL designers are not required to follow these directions, as long as they comply with the signal-level protocol at the interface with the ESP socket.

2.1.1 Accelerator Configuration

The configuration block regulates the accelerator execution and implements the interface with software by sampling the value of common and user-defined configuration registers located in the ESP socket.

Table 2.1: Description of the ESP accelerator configuration port.

Signal	Driver	Description
clk	socket	accelerator clock.
rst	socket	accelerator reset active low . The socket activates this reset signal when software clears the interrupt request to ensure that the accelerator is ready for a new invocation and internal state is clean. If the accelerator is expected to retain its state across different invocations, a user-defined configuration register can be used to implement a software-controlled reset signal.
conf_done	socket	Configuration registers are valid and computation can start. This signal is active high and asserted for one clk cycle to trigger the accelerator execution.
conf_info_<register_name>	socket	User-defined configuration input. The corresponding memory-mapped configuration registers are automatically generated in the ESP socket when creating the SoC instance. There can be up to 14 user-defined registers that must be listed in the accelerator definition XML file. For each register the accelerator must expose one conf_info_ input. Bitwidth must be between 1 and 32 bits. These inputs should be considered valid when the conf_done input is active high.
acc_done	accelerator	Single-cycle pulse active high. This flag indicates that the accelerator has completed its task. The pulse should occur only after the last DMA write transaction has completed and all output data have been transferred from the PLM to the memory hierarchy. Asserting acc_done will trigger an interrupt request to the interrupt controller located in the ESP auxiliary tile. The software interrupt handler is responsible for clearing the interrupt, thus resetting the state of the socket and activating the rst input of the accelerator.
debug	accelerator	32-bit debug output. The accelerator designer can use this output to encode error codes. The state of this output can be accessed through the common memory-mapped registers present in the socket.

2.1.2 Private Local Memory

The PLM can be generated with the ESP Memgen utility, which combines SRAM primitives available as part of the target technology libraries. Alternatively, the accelerator designer can manually implement the PLM in RTL.

The PLM is not memory mapped, hence it is not exposed to software. Furthermore, the PLM is not part of the SoC cache hierarchy as it is solely intended as a customized working buffer for the accelerator data path. As a result, the PLM has no external interface exposed to the ESP accelerator socket and RTL designers are not required to comply with any hierarchy convention or signal-level protocol to implement the PLM, unless they use the ESP Memgen utility (refer to the ESP Memory Generator documentation for further information). Accelerators that operate on small batches of data and don't have particular buffering requirements can also be implemented without a PLM.



2.2 DMA Transactions

The master of a direct-memory access (DMA) transaction is always an accelerator. The accelerator initiates a DMA read transfer through the `rd_ctrl` channel and a DMA write transfer through the `wr_ctrl` channel. Tables 2.2 and 2.3 describe the fields and the encoding of the two control channels.

DMA control channels are LIC that follow a simple protocol [Carloni, 2015]: when both *valid* and *ready* control signals are set, the value of the data bus is sampled by the slave. From the accelerator view point *valid* and *ready* are independent and there should be no combinational path between the two signals.

An ESP accelerator does not issue requests using physical addresses. The field *index* of the control channels indicates an offset with respect to a virtual memory region reserved for the accelerator. The ESP device driver allocates this region in virtual pages and generates a corresponding page table. The ESP accelerator socket handles address translation, therefore the accelerator can operate as if the reserved area was contiguous.

Table 2.2: Encoding of DMA size

Encoding	Name	Bitwidth
000	BYTE	8
001	HWWORD	16
010	WORD	32
011	DWORD	64

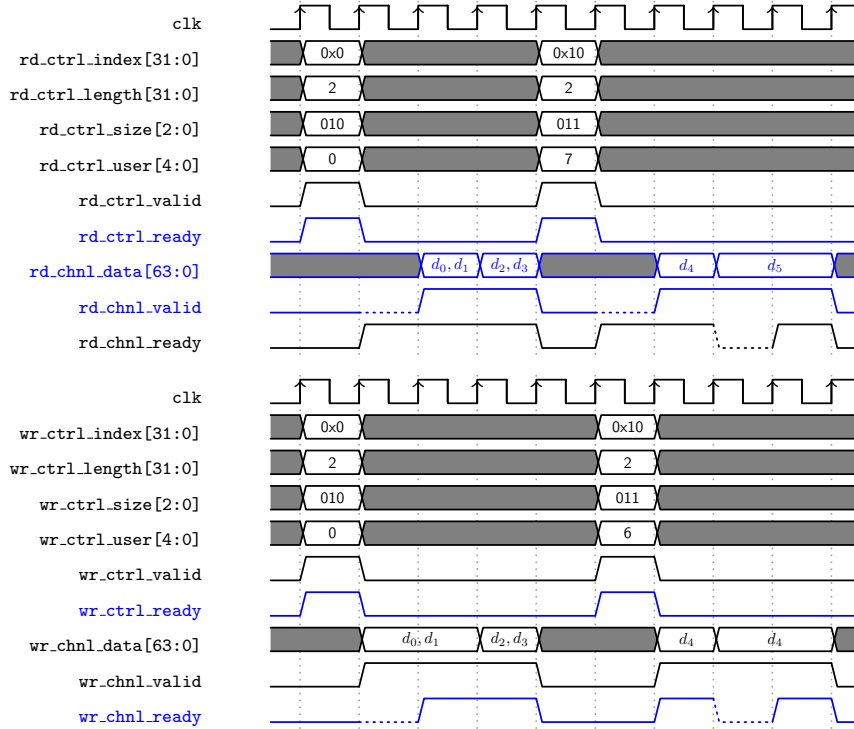
Table 2.3: Description of the DMA control channels `rd_ctrl` and `wr_ctrl`.

Signal	Driver	Description
<code>dma_[read write]_ctrl_data_index</code>	accelerator	Offset of a DMA read or write transaction expressed as number of beats. This offset is used to compute the starting address of the transaction.
<code>dma_[read write]_ctrl_data_length</code>	accelerator	Length of a DMA read or write transaction expressed as number of beats.
<code>dma_[read write]_ctrl_data_size</code>	accelerator	Bitwidth of the data token for the DMA transaction. This signal is used to correct the NoC flits when the processor architecture follows the <i>big endian</i> convention to store data in memory. This signal follows the encoding in Table 2.2.
<code>dma_[read write]_ctrl_data_user</code>	accelerator	This field gives the accelerator control over the data access mode used to load/store data from/to the SoC. On the read channel, this field specifies the source of the data. On the write channel, it specifies the number of consumers of the data. More details are presented in Section 2.2.1.
<code>dma_[read write]_ctrl_valid</code>	accelerator	Flag indicating a new DMA transaction request. When set, all data fields must be valid. This flag must not depend combinatorially on the corresponding <i>ready</i> signal.
<code>dma_[read write]_ctrl_ready</code>	socket	Flag indicating that the ESP socket is ready to accept a new DMA request. This flag must not depend combinatorially on the corresponding <i>valid</i> signal.

For an accelerator with N-bits DMA interface (e.g. 64 bits), the physical address in bytes of a DMA transaction is computed by the ESP socket as follows:

$$addr = walk_accelerator_ptable(index * N/8) \tag{2.1}$$

The user-level driver is responsible to prepare data in memory using the same offsets used for DMA transfers by the accelerator. Offset calculation can be defined at design time by hard-coding the logic to compute offsets in the accelerator. Alternatively, offsets can be computed in software and configured at run time through user-defined control and status registers (CSRs).



Symmetrically, when a DMA read transfer is configured, the accelerator must transfer the exact number of beats set with the *length* field. Data beats are transferred through the DMA write channel by setting the *valid* flag high when the corresponding data signal is valid. A beat is transferred when both *valid* and *ready* are set during the same cycle. No restriction is imposed on the throughput of the transfer. The accelerator must hold valid data on the DMA write channel when the socket is not ready to sample it. This condition may occur in case of contention for NoC links, or external memory channels.

Figure 2.2 and 2.2 show two examples of DMA read and DMA write transactions. Signals in blue are driven by the socket, while signals in black are driven by the accelerator. Dotted lines indicate back-pressure, which can be applied by either the accelerator or the socket.

2.2.1 DMA User Field

ESP supports 3 basic types of accelerator communication: 1) direct memory access (DMA), 2) point-to-point (P2P) communication [Giri et al., 2020] from one accelerator to another, and 3) multicast communication from one accelerator to several others. Traditionally, the accelerator’s communication mode is set through configuration registers and is fixed for the entire duration of an accelerator invocation. Beginning with ESP release 2025.1.0, the communication



mode instead can be selected by the accelerator itself through the *user* fields on the DMA read and write control channels [Zuckerman et al., 2024]. This feature is optional, however, and the user field can be tied to zero on both control channels to default to the settings specified by the configuration registers.

On the read control channel, the user field specifies the *source* of the data. A value of zero defaults to the configuration register settings (memory or P2P). Any value greater indexes into a small lookup table that is implemented in configuration registers in the accelerator socket. Each entry in the lookup table is an X-Y pair that specifies the tile coordinates of a source in the SoC. Entry 0 of this table is hard-coded to always hold the local tile's coordinates; other entries must be programmed by SW prior to the accelerator invocation. The coordinates obtained from using the user field to index into this table are then used to formulate a P2P data request to the selected source.

On the write control channel, the user field specifies the *number of consumers* of the data. A value of zero defaults to the configuration register settings (memory or P2P). A value of 1 indicates a standard P2P request, while any value greater specifies a multicast transfer, in which all consumers will receive the same data using a single message on the NoC. When a non-zero value is specified, the DMA controller (socket) will wait for the specified number of P2P requests to be received before sending the write data on the NoC.

For further information on this feature, please see the ESP release for code examples, the hands-on tutorials on the ESP website (**coming soon**), and the two papers cited above.

Bibliography

- [Carloni, 2015] Carloni, L. P. (2015). From latency-insensitive design to communication-based system-level design. *Proceedings of the IEEE*, 103(11):2133–2151.
- [Giri et al., 2020] Giri, D., Chiu, K.-L., Guglielmo, G. D., Mantovani, P., and Carloni, L. P. (2020). ESP4ML: Platform-based design of systems-on-chip for embedded machine learning.
- [Zuckerman et al., 2024] Zuckerman, J., Wellman, J.-D., Vanamali, A., Shankar, M., Tombesi, G., Swaminathan, K., Lee, K., Kapur, M., Philhower, R., Bose, P., and Carloni, L. P. (2024). Towards generalized on-chip communication for programmable accelerators in heterogeneous architectures.