

The OPU ISA Specification

Version 0.1-draft, July 2020

Edited by Louis Delhez, Lei He, Yunxuan Yu, and Tiandong Zhao.
opu-admin@googlegroups.com

Contents

1	Introduction	2
1.1	Data types	2
1.2	Instruction Fundamentals	2
1.3	Memory	3
2	OPU ISA	4
2.1	Registers	4
2.2	Buffers	5
2.3	Instructions	5
2.3.1	Configuration Instructions	5
2.3.2	Imperative Instructions	9
3	Arithmetic rules	14
3.1	Representable set	14
3.2	Effective width	14
3.3	Conversion from \mathbb{R}	14
A	Instruction Listing	15

Chapter 1

Introduction

Academic research of CPU has been facilitated greatly by open source toolchain and resources such as instruction set (ISA), compiler, and instruction/microarchitecture level simulation. Examples include SimpleScalar, GEM5, and more recently, RISC-V toolset. Yet, such open-source and complete eco-system does not exist for general machine learning algorithms. OPU, short for *open processing unit*, is meant to be an open-source and complete eco-system for machine learning hardware research, including: ISA with executable specifications, compiler with formal verification, instruction level (functional) and microarchitecture level (cycle-accurate) simulation, parametrized modules in RTL and Chisel, and FPGA emulation and development boards.

This documents describes the OPU Instruction Set Architecture.

1.1 Data types

The OPU ISA is designed to be polymorphic: different implementations operating on different data types can execute the same binary. Implementations are characterized by the four data types that they use for different kinds of values:

- ITYPE: type of feature values in memory
- BTYPE: type of bias values in memory
- KTYPE: type of kernel weights in memory
- OTYPE: type of partial sums in the processor

Chapter 3 defines how the various combinations of datatypes should be handled by the implementations.

1.2 Instruction Fundamentals

An OPU *program* consists of a sequence of *instructions* stored contiguously in memory. Instructions are always encoded on 32 bits, following the little-endian convention in memory. Instructions are executed one after the other in program order. The execution stops when an end instruction is executed. The address of the first instruction of a program must be 64-byte aligned. There is no such restriction regarding the last instruction of the program.

1.3 Memory

OPU has an address space of 2^{32} bytes. The little-endian convention is always assumed by OPU when reading and writing multi-bytes values in memory. Instructions and data share a same unique address space.

Memory accesses are always performed in program order. A load operation always sees the effect of all previous store operations.

Chapter 2

OPU ISA

2.1 Registers

Contrary to conventional ISAs that define registers for storing temporary values on which operations can be performed, the registers of the OPU ISA are exclusively utilized for configuration of subsequent operations.

The registers of the OPU ISA are listed in Table 2.1.

Registers are set by *configuration instructions*, whose names start with the symbol @.

<i>Name</i>	<i>Description</i>	<i>Possible values</i>
ifm_h	first dimension of the ifm buffer	$\{1, \dots, 2^7 - 1\}$
ifm_w	second dimension of the ifm buffer	$\{1, \dots, 2^7 - 1\}$
ifm_c	third dimension of the ifm buffer	$2^i, i \in \{4, \dots, 6\}$
ofm_h	first dimension of the ofm buffer	$\{1, \dots, 2^7 - 1\}$
ofm_w	second dimension of the ofm buffer	$\{1, \dots, 2^7 - 1\}$
ofm_c	third dimension of the ofm buffer	$2^i, i \in \{1, \dots, 6\}$
ker_n	first dimension of the ker buffer	$\{1, \dots, 36\}$
stride_h	stride along first dimension for conv instructions	$\{1, \dots, 7\}$
stride_w	stride along second dimension for conv instructions	$\{1, \dots, 7\}$
ifm_addr	address offset for ld.ifm instructions	$a \times 2^{28}, a \in \{0, \dots, 2^4 - 1\}$
ker_addr	address offset for ld.ker instructions	$a \times 2^{28}, a \in \{0, \dots, 2^4 - 1\}$
bias_addr	address offset for ld.bias instructions	$a \times 2^{28}, a \in \{0, \dots, 2^4 - 1\}$
ifm_mem_w	second dimension of array for ld.ifm instructions	$\{1, \dots, 2^{10} - 1\}$
ofm_mem_w	second dimension of array for ld.ofm instructions	$\{1, \dots, 2^{10} - 1\}$
ofm_mem_h	first dimension of array for pad instructions	$\{1, \dots, 2^{10} - 1\}$
ifm_shift	scale factor of the ifm buffer	$\{-2^7, \dots, 2^7 - 1\}$
bias_shift	scale factor of the bias buffer	$\{-2^7, \dots, 2^7 - 1\}$
post_order	order of operations for store instructions	$\{0, 1, 2\}$
add_ifm	switch for ifm addition for store instructions	$\{0, 1\}$
act	activation function to use for store instructions	$\{0, 1, 2\}$
pool_h	first dimension of the pooling filter	$\{1, \dots, 2^4 - 1\}$
pool_w	second dimension of the pooling filter	$\{1, \dots, 2^4 - 1\}$
pool_h_stride	stride along first dimension for the pooling filter	$\{1, \dots, 2^3 - 1\}$
pool_w_stride	stride along first dimension for the pooling filter	$\{1, \dots, 2^3 - 1\}$

Table 2.1: Registers of the OPU ISA.

2.2 Buffers

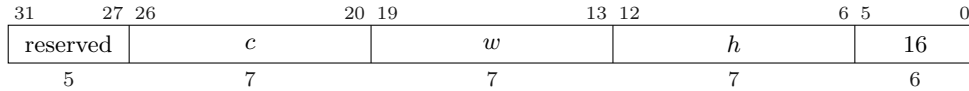
In OPU, the role of traditional processor registers is fulfilled by large flexible multi-dimensional memories called buffers. Four different buffers are defined: `ifm`, `ker`, `bias`, and `ofm`. Table 2.2 details their characteristics. While the dimensions of these buffers can be modified at runtime by setting the appropriate registers, the types of the elements of the buffers are constant for a given implementation.

Buffer	Dimensions	Data type
<code>ifm</code>	$\text{ifm_h} \times \text{ifm_w} \times \text{ifm_c}$	ITYPE
<code>ofm</code>	$\text{ofm_h} \times \text{ofm_w} \times \text{ofm_c}$	OTYPE
<code>ker</code>	$\text{ker_n} \times \text{ofm_c} \times \text{ifm_c}$	KTYPE
<code>bias</code>	ofm_c	BTYPE

2.3 Instructions

2.3.1 Configuration Instructions

@shape.ifm



The `@shape.ifm` instruction sets the size of the `ifm` buffer as

$$\begin{aligned} \text{ifm_h} &\leftarrow h \\ \text{ifm_w} &\leftarrow w \quad \text{where} \quad \gamma = 2^c. \\ \text{ifm_c} &\leftarrow \gamma \end{aligned}$$

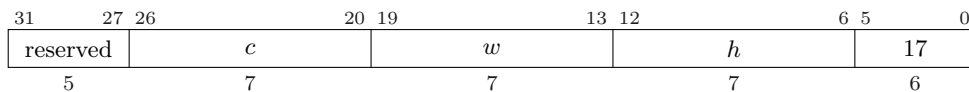
When this instruction is executed, the content of the `ifm` and `ker` buffers is invalidated. The instruction is only valid if the following condition is satisfied:

$$1 \leq h \times w \leq 2048 \wedge 16 \leq \gamma \leq 64.$$

In assembly, this instruction is represented by

`@shape.ifm [h,w, γ]`

@shape.ofm



The `@shape.ofm` instruction sets the size of the `ofm` buffer as

$$\begin{aligned} \text{ofm_h} &\leftarrow h \\ \text{ofm_w} &\leftarrow w \quad \text{where} \quad \gamma = 2^c. \\ \text{ofm_c} &\leftarrow \gamma \end{aligned}$$

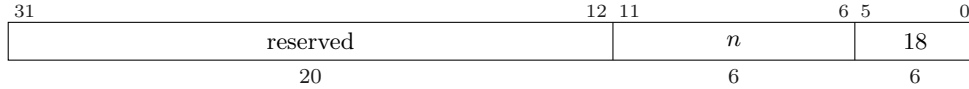
When this instruction is executed, the content of the `ofm` and `ker` buffers is invalidated. The instruction is only valid if the following condition is satisfied:

$$1 \leq h \times w \leq 2048 \wedge 2 \leq \gamma \leq 64.$$

In assembly, this instruction is represented by

`@shape.ofm [h,w, γ]`

@shape.ker



The `@shape.ker` instruction sets the first dimension of the `ker` buffer as

$$\text{ker}_n \leftarrow n.$$

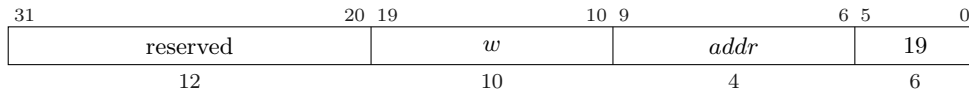
When this instruction is executed, the content of the `ker` buffer is invalidated. The instruction is only valid if the following condition is satisfied:

$$1 \leq n \leq 36.$$

In assembly, this instruction is represented by

`@shape.ker n`

@mem.ifm



The `@mem.ifm` instruction specifies the location of the data to be loaded by the subsequent `ld.ifm` instructions as

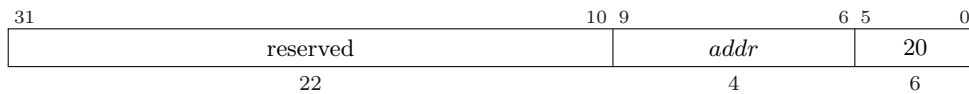
$$\text{ifm_addr} \leftarrow \text{addr} \ll 28$$

$$\text{ifm_mem}_w \leftarrow w.$$

In assembly, this instruction is represented by

`@mem.ifm addr, w`

@mem.ker



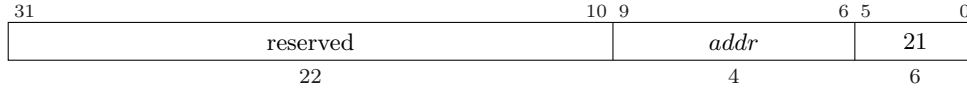
The `@mem.ker` instruction specifies the location of the data to be loaded by the subsequent `ld.ker` instructions as

$$\text{ker_addr} \leftarrow \text{addr} \ll 28.$$

In assembly, this instruction is represented by

`@mem.ker addr`

@mem.bias



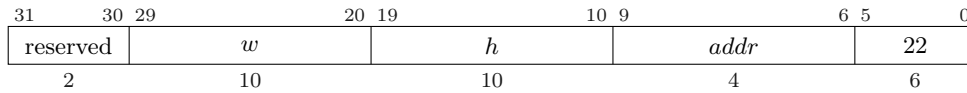
The @mem.bias instruction specifies the location of the data to be loaded by the subsequent ld.bias instructions as

$$\text{bias_addr} \leftarrow \text{addr} \ll 28.$$

In assembly, this instruction is represented by

@mem.bias *addr*

@mem.ofm



The @mem.ofm instruction specifies the location where data will be stored by the subsequent store instructions as

$$\text{ofm_addr} \leftarrow \text{addr} \ll 28$$

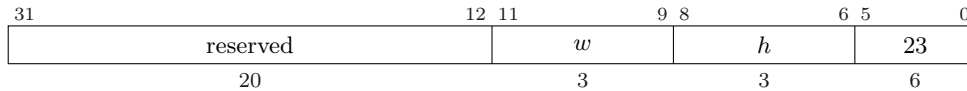
$$\text{ofm_mem_h} \leftarrow h$$

$$\text{ofm_mem_w} \leftarrow w.$$

In assembly, this instruction is represented by

@mem.ofm *addr*, [*h*, *w*]

@stride



The @stride instruction sets the strides for the subsequent conv instructions with

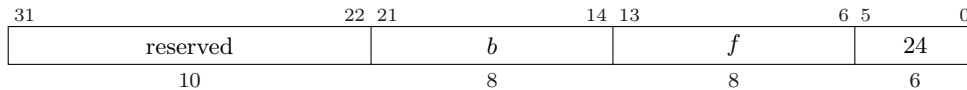
$$\text{stride_h} \leftarrow h$$

$$\text{stride_w} \leftarrow w.$$

In assembly, this instruction is represented by

@stride [*h*, *w*]

@shift



The @shift sets the scale factors for the values of the ifm and bias buffers:

$$\text{ifm_shift} \leftarrow f$$

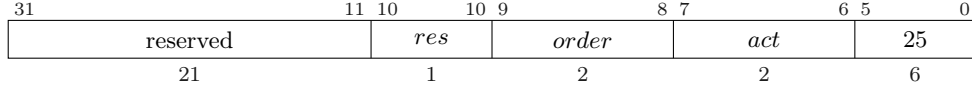
$$\text{bias_shift} \leftarrow b$$

The fields *f* and *b* are signed integers encoded using two's complement representation.

In assembly, this instruction is represented by

@shift *f*, *b*

@post



The `post` instruction specifies what operation must be performed during subsequent store as

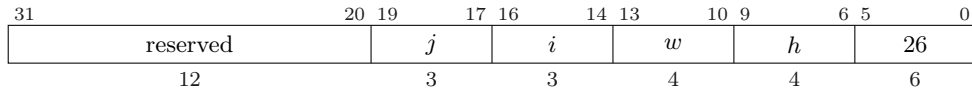
$$\begin{aligned} \text{act} &\leftarrow \text{act} \\ \text{post_order} &\leftarrow \text{order} \\ \text{add_ifm} &\leftarrow \text{res} \end{aligned}$$

This instruction is only valid for the values listed in Table 2.2. As shown in this table, the representation of this instruction in assembly depends on the values of *res*, *order*, and *act*.

<i>order</i>	<i>act</i>	<i>res</i>	Assembly
0	0	0	@post pool
0	0	1	@post res, pool
0	1	0	@post act.relu, pool
0	1	1	@post act.relu, res, pool
0	2	0	@post act.leaky, pool
0	2	1	@post act.leaky, res, pool
1	1	1	@post res, act.relu, pool
1	2	1	@post res, act.leaky, pool
2	0	1	@post pool, res
2	1	1	@post act.relu, pool, res
2	2	1	@post act.leaky, pool, res

Table 2.2: Assembly representation of @post.

@pool



The `pool` instruction characterizes the pooling filter applied in store instructions:

$$\begin{aligned} \text{pool_h} &\leftarrow h \\ \text{post_w} &\leftarrow w \\ \text{pool_h_stride} &\leftarrow i \\ \text{pool_w_stride} &\leftarrow j \end{aligned}$$

The instruction is only valid if the following condition is satisfied:

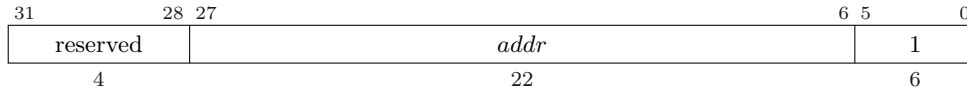
$$(1 \leq h) \wedge (1 \leq w) \wedge (1 \leq i) \wedge (1 \leq j).$$

In assembly, this instruction is represented by

$$\text{@pool } [h, w], [i, j]$$

2.3.2 Imperative Instructions

ld.ifm



The `ld.ifm` instruction loads a 3-dimensional array of type `ITYPE` from memory into the `ifm` buffer.

The shape of the loaded array is $\text{ifm_h} \times \text{ifm_w} \times \text{ifm_c}$ (i.e. the `ifm` buffer is filled entirely). The first element of this array is assumed to be stored at address

$$\text{ifm_addr} + \text{addr} \times 64.$$

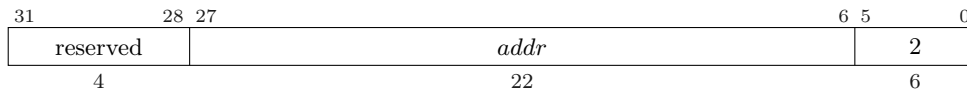
The access stride along the three dimensions of the array are respectively

$$\text{ifm_mem_w} \times 64, 64, \text{ and } 1.$$

In assembly, this instruction is represented by

`ld.ifm addr`

ld.ker



The `ld.ker` instruction loads a 3-dimensional array of type `KTYPE` from memory into the `ker` buffer.

The shape of the loaded array is $\text{ker_n} \times \text{ofm_c} \times \text{ifm_c}$ (i.e. the `ker` buffer is filled entirely). The first element of this array is assumed to be stored at address

$$\text{ker_addr} + \text{addr} \times 64.$$

The array is assumed to be stored contiguously in memory in row-major order, i.e. the access strides along the three dimensions of the array are respectively

$$\text{ofm_c} \times \text{ifm_c}, \text{ofm_c}, \text{ and } 1.$$

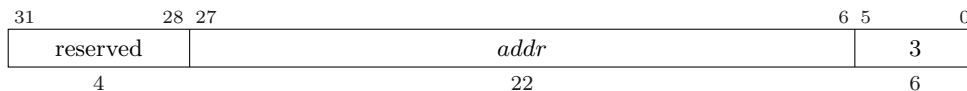
The `ld.ker` instruction is only valid if

$$n \times \max \left\{ \frac{\text{ifm_c} \times \text{ofm_c}}{1024}, 1 \right\} \leq 36.$$

In assembly, this instruction is represented by

`ld.ker addr`

ld.bias



The `ld.bias` instruction loads a 1-dimensional array of type `BTYPE` from memory into the `bias` buffer.

The length of the loaded array is equal to `ofm_c` (i.e. the `bias` buffer is filled entirely). The first element of this array is assumed to be stored at address

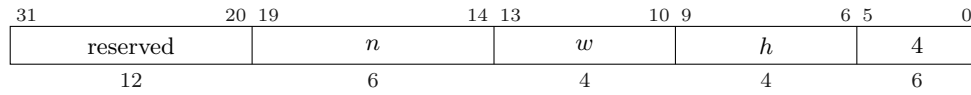
$$\text{bias_addr} + \text{addr} \times 64.$$

The array is assumed to be stored contiguously in memory.

In assembly, this instruction is represented by

`ld.bias addr`

conv



The `conv` instruction performs a pointwise convolution between the `ifm` buffer and one slice of the `ker` buffer. The result of this operation is written to the `ofm` buffer:

$$\forall (i, j, k) \in \{0, \dots, \text{ofm_h} - 1\} \times \{0, \dots, \text{ofm_w} - 1\} \times \{0, \dots, \text{ofm_c} - 1\} :$$

$$\text{ofm}[i, j, k] := \kappa_{\text{OTYPE}} \left(2^{\text{ifm_shift}} \times \sum_{\ell=0}^{\text{ifm_c}-1} A[i, j, k, \ell] \right)$$

where

$$A[i, j, k, \ell] = \rho_{\text{KTYPE}}(\text{ker}[\text{ker_n}, k, \ell]) \times \rho_{\text{ITYPE}}(\text{ifm}[h + \text{stride_h} \times i, w + \text{stride_w} \times j, \ell]).$$

The type casting functions ρ and κ are defined in Chapter 3.

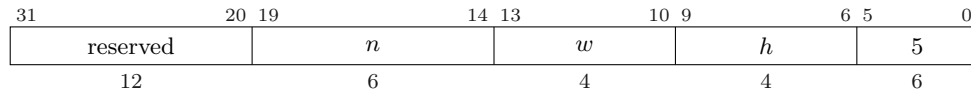
The instruction is only valid if the following condition is satisfied:

$$h + \text{ofm_h} \times (\text{stride_h} - 1) < \text{ifm_h} \wedge w + \text{ofm_w} \times (\text{stride_w} - 1) < \text{ifm_w}.$$

In assembly, this instruction is represented by

`conv ifm:[h,w], ker:n`

conv.bias



Similarly to `conv`, the `conv.bias` instruction performs a pointwise convolution between the `ifm` buffer and one slice of the `ker` buffer. In addition, the content of the `bias` buffer is added to the result. The result of this operation is written to the `ofm` buffer:

$$\forall (i, j, k) \in \{0, \dots, \text{ofm_h} - 1\} \times \{0, \dots, \text{ofm_w} - 1\} \times \{0, \dots, \text{ofm_c} - 1\} :$$

$$\text{ofm}[i, j, k] := \kappa_{\text{OTYPE}} \left(2^{\text{bias_shift}} \times \rho_{\text{BTYPE}}(\text{bias}[k]) + 2^{\text{ifm_shift}} \times \sum_{\ell=0}^{\text{ifm_c}-1} A[i, j, k, \ell] \right)$$

where

$$A[i, j, k, \ell] = \rho_{\text{KTYPE}}(\text{ker}[\text{ker_n}, k, \ell]) \times \rho_{\text{ITYPE}}(\text{ifm}[h + \text{stride_h} \times i, w + \text{stride_w} \times j, \ell]).$$

The type casting functions ρ and κ are defined in Chapter 3.

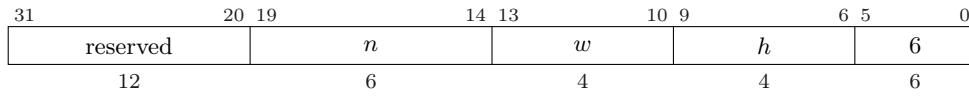
The instruction is only valid if the following condition is satisfied:

$$h + \text{ofm_h} \times (\text{stride_h} - 1) < \text{ifm_h} \wedge h + \text{ofm_w} \times (\text{stride_w} - 1) < \text{ifm_w}.$$

In assembly, this instruction is represented by

`conv.bias ifm:[h,w], ker:n`

conv.acc



Similarly to `conv`, the `conv.acc` instruction performs a pointwise convolution between the `ifm` buffer and one slice of the `ker` buffer. In addition, the content of the `ofm` buffer is added to the result:

$$\forall (i, j, k) \in \{0, \dots, \text{ofm_h} - 1\} \times \{0, \dots, \text{ofm_w} - 1\} \times \{0, \dots, \text{ofm_c} - 1\} :$$

$$\text{ofm}[i, j, k] := \kappa_{\text{OTYPE}} \left(\rho_{\text{OTYPE}}(\text{ofm}[i, j, k]) + 2^{\text{ifm_shift}} \times \sum_{\ell=0}^{\text{ifm_c}-1} A[i, j, k, \ell] \right)$$

where

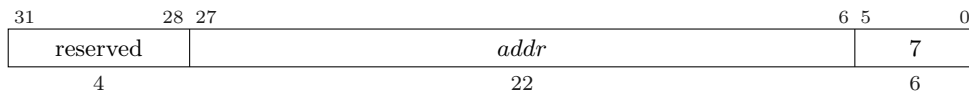
$$A[i, j, k, \ell] = \rho_{\text{KTYPE}}(\text{ker}[\text{ker_n}, k, \ell]) \times \rho_{\text{ITYPE}}(\text{ifm}[h + \text{stride_h} \times i, w + \text{stride_w} \times j, \ell]).$$

The type casting functions ρ and κ are defined in Chapter 3.

In assembly, this instruction is represented by

`conv.acc ifm:[h,w], ker:n`

store



Let A be the 3-dimensional array of shape $(\text{ofm_h}, \text{ofm_w}, \text{ofm_c})$ obtained by casting the elements of the `ofm` buffer to the type `ITYPE` as follows:

$$A[i, j, k] := \kappa_{\text{ITYPE}} \left(2^{\delta_{\text{ITYPE}} - \delta_{\text{OTYPE}}} \times \rho_{\text{OTYPE}}(\text{ofm}[i, j, k]) \right)$$

where δ , κ , and ρ are defined in Chapter 3.

Let *activation* be a function that accepts as argument a 3-dimensional array X of shape $(\text{ofm_h}, \text{ofm_w}, \text{ofm_c})$ and type `ITYPE` and that returns an array of same shape and type according to

$$\text{activation}(X)[i, j, k] := \begin{cases} X[i, k, k] & \text{if act} = 0 \\ \kappa_{\text{ITYPE}}(\max\{\rho_{\text{ITYPE}}(X[i, j, k]), 0\}) & \text{if act} = 1 \\ \kappa_{\text{ITYPE}}(\max\{\rho_{\text{ITYPE}}(X[i, j, k]), \rho_{\text{ITYPE}}(X[i, j, k])/8\}) & \text{if act} = 2 \end{cases}$$

Let *residual* be a function that accepts as argument a 3-dimensional array X of type ITYPE and returns an array of same shape and type according to

$$\text{residual}(X)[i, j, k] := \kappa_{\text{ITYPE}} \left(\rho_{\text{ITYPE}}(X[i, j, k]) + \rho_{\text{ITYPE}}(\text{ifm}[i, j, k]) \right).$$

Let *pooling* be a function that accepts as argument a 3-dimensional array X of shape $(\text{ofm_h}, \text{ofm_w}, \text{ofm_c})$ and type ITYPE and returns an array of same type and shape

$$\left(\left\lfloor \frac{\text{ofm_h} - \text{pool_h}}{\text{pool_h_stride}} \right\rfloor + 1, \left\lfloor \frac{\text{ofm_w} - \text{pool_w}}{\text{pool_w_stride}} \right\rfloor + 1, \text{ofm_c} \right)$$

according to

$$\text{pooling}(X)[i, j, k] := \kappa_{\text{ITYPE}} \left(\max_{(a,b) \in \mathcal{A} \times \mathcal{B}} \rho_{\text{ITYPE}}(X[i \times \text{pool_h_stride} + a, j \times \text{pool_w_stride} + b, k]) \right)$$

where where $\mathcal{A} = \{0, 1, \dots, \text{pool_h} - 1\}$ and $\mathcal{B} = \{0, 1, \dots, \text{pool_w} - 1\}$.

The `store` instruction produces a 3-dimensional array B of type ITYPE as given by

$$B := \begin{cases} \text{pooling}(\text{residual}(\text{activation}(A))) & \text{if } \text{pos_order} = 0 \\ \text{pooling}(\text{activation}(\text{residual}(A))) & \text{if } \text{pos_order} = 1 \\ \text{residual}(\text{pooling}(\text{activation}(A))) & \text{if } \text{pos_order} = 2 \end{cases}$$

The array B is then stored to memory starting at address

$$\text{ofm_addr} + \text{addr} \times 64.$$

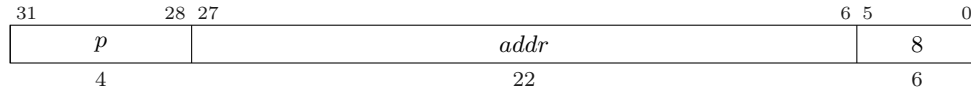
The strides along the three dimensions are respectively

$$\text{ofm_mem_w} \times 64, 64, \text{ and } 1.$$

In assembly, this instruction is represented by

`store addr`

pad



The `pad` instruction pads a 3-dimensional array in memory along its first two dimension. More specifically, assuming that *rows* and *cols* refer respectively to the first and second dimensions of the array, this instruction fills with zeros (i.e. $\kappa_{\text{ITYPE}}(0)$) the first p rows, first p cols, last p rows, and last p columns.

The shape of the considered array is

`[ofm_mem_h, ofm_mem_w, 64]`.

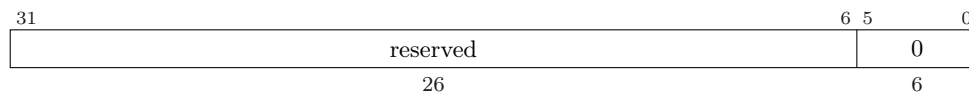
This array is assumed to be stored contiguously in memory in row-major order, i.e. its strides are respectively

`ofm_mem_w, 64, and 1.`

In assembly, this instruction is represented by

`pad addr, p`

end



The end instruction terminates the execution of the current program. All programs must end with an end instruction.

In assembly, this instruction is simply represented by

`end`

Chapter 3

Arithmetic rules

In OPU, a data type T is defined by

- a set \mathcal{S}_T of all valid representations
- a function $\rho_T : \mathcal{S}_T \rightarrow \mathbb{R}$ that maps a represented number to its actual value

3.1 Representable set

The set $\mathcal{R}_T \subset \mathbb{R}$ of all representable values using a data type T is defined as

$$\mathcal{R}_T \equiv \{x \in \mathbb{R} \mid (\exists y \in \mathcal{S}_T)[\rho_T(y) = x]\}.$$

3.2 Effective width

The effective width δ_T of a data type T is defined as

$$\delta_T \equiv \left\lceil \log_2 \left(\max_{x \in \mathcal{R}_T} |x| \right) \right\rceil.$$

3.3 Conversion from \mathbb{R}

The function $\kappa_T : \mathbb{R} \rightarrow \mathcal{S}_T$ maps any real number to its nearest representation. For all $x \in \mathbb{R}$, this function must satisfy

$$\rho_T(\kappa_T(x)) = \max_{a \in \mathcal{A}} (\rho_T(a))$$

where

$$\mathcal{A} = \arg \min_{y \in \mathcal{S}_T} |\rho_T(y) - x|.$$

Appendix A

Instruction Listing

end	31	-														6	5	0	0												
ld.ifm	31	-	28	27	<i>addr</i>														6	5	1	0									
ld.ker	31	-	28	27	<i>addr</i>														6	5	2	0									
ld.bias	31	-	28	27	<i>addr</i>														6	5	3	0									
conv	31	-														20	19	<i>n</i>	14	13	<i>w</i>	10	9	<i>h</i>	6	5	4	0			
conv.bias	31	-														20	19	<i>n</i>	14	13	<i>w</i>	10	9	<i>h</i>	6	5	5	0			
conv.acc	31	-														20	19	<i>n</i>	14	13	<i>w</i>	10	9	<i>h</i>	6	5	6	0			
store	31	-	28	27	<i>addr</i>														6	5	7	0									
pad	31	<i>p</i>	28	27	<i>addr</i>														6	5	8	0									
@shape.ifm	31	-	27	26	<i>c</i>	20	19	<i>w</i>	13	12	<i>h</i>	6	5	16	0																
@shape.ofm	31	-	27	26	<i>c</i>	20	19	<i>w</i>	13	12	<i>h</i>	6	5	17	0																
@shape.ker	31	-											12	11	<i>n</i>	6	5	18	0												
@mem.ifm	31	-														20	19	<i>w</i>	10	9	<i>addr</i>	6	5	19	0						
@mem.ker	31	-														10	9	<i>addr</i>	6	5	20	0									
@mem.bias	31	-														10	9	<i>addr</i>	6	5	21	0									
@mem.ofm	31	-	30	29	<i>w</i>	20	19	<i>h</i>	10	9	<i>addr</i>	6	5	22	0																
@stride	31	-											12	11	<i>w</i>	9	8	<i>h</i>	6	5	23	0									
@shift	31	-														22	21	<i>b</i>	14	13	<i>f</i>	6	5	24	0						
@post	31	-											11	10	9	<i>r</i>	8	7	<i>ord</i>	6	5	25	0								
@pool	31	-														20	19	<i>j</i>	17	16	<i>i</i>	14	13	<i>w</i>	10	9	<i>h</i>	6	5	26	0