

Arm China Zhouyi Compass Driver and Runtime

Version: 1.0

User Guide

Confidential

arm CHINA

Arm China Zhouyi Compass Driver and Runtime

User Guide

Copyright © 2021–2023 Arm Technology (China) Co., Ltd. All rights reserved.

Release Information

Document History

Issue	Date	Confidentiality	Change
0001-00	30 September 2023	Confidential	First release for 1.0

Confidential Proprietary Notice

This document is **CONFIDENTIAL** and any use by you is subject to the terms of the agreement between you and Arm Technology (China) Co., Ltd ("Arm China") or the terms of the agreement between you and the party authorized by Arm China to disclose this document to you.

This document is protected by copyright and other related rights and the practice or implementation of the information contained in this document may be protected by one or more patents or pending patent applications. No part of this document may be reproduced in any form by any means without the express prior written permission of Arm China. **No license, express or implied, by estoppel or otherwise to any intellectual property rights is granted by this document unless specifically stated.**

Your access to the information in this document is conditional upon your acceptance that you will not use or permit others to use the information: **(i)** for the purposes of determining whether implementations infringe any third party patents; **(ii)** for developing technology or products which avoid any of Arm China's intellectual property; or **(iii)** as a reference for modifying existing patents or patent applications or creating any continuation, continuation in part, or extension of existing patents or patent applications; or **(iv)** for generating data for publication or disclosure to third parties, which compares the performance or functionality of the Arm China technology described in this document with any other products created by you or a third party, without obtaining Arm China's prior written consent.

THIS DOCUMENT IS PROVIDED "AS IS". ARM CHINA PROVIDES NO REPRESENTATIONS AND NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, INCLUDING, WITHOUT LIMITATION, THE IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, NON-INFRINGEMENT OR FITNESS FOR A PARTICULAR PURPOSE WITH RESPECT TO THE DOCUMENT. For the avoidance of doubt, Arm China makes no representation with respect to, and has undertaken no analysis to identify or understand the scope and content of, third party patents, copyrights, trade secrets, or other rights.

This document may include technical inaccuracies or typographical errors.

TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL ARM CHINA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF ARM CHINA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

This document consists solely of commercial items. You shall be responsible for ensuring that any use, duplication or disclosure of this document complies fully with any relevant export laws and regulations to assure that this document or any portion thereof is not exported, directly or indirectly, in violation of such export laws. Use of the word "partner" in reference to Arm China's customers is not intended to create or refer to any partnership relationship with any other company. Arm China may make changes to this document at any time and without notice.

If any of the provisions contained in these terms conflict with any of the provisions of any click through or signed written agreement covering this document with Arm China, then the click through or signed written agreement prevails over and supersedes the conflicting provisions of these terms. This document may be translated into other languages for convenience, and you agree that if there is any conflict between the English version of this document and any translation, the terms of the English version of the Agreement shall prevail.

Arm China is a trading name of Arm Technology (China) Co., Ltd. The words marked with ® or ™ are registered trademarks in the People's Republic of China and/or elsewhere. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners.

Copyright © 2021–2023 Arm China (or its affiliates). All rights reserved.

Confidentiality Status

This document is Confidential. This document may only be used and distributed in accordance with the terms of the agreement entered into by Arm China and the party that Arm China delivered this document to.

Product Status

The information in this document is Final, that is a for a developed product.

Web Address

<https://www.armchina.com>

Contents

Arm China Zhouyi Compass Driver and Runtime User Guide

	Preface	5
Chapter 1	NPU driver	1-9
1.1	About the NPU driver	1-10
1.2	Linux-based driver.....	1-11
1.3	QNX-based driver (for customized solutions only).....	1-50
1.4	Bare-metal-based driver.....	1-68
1.5	RTOS-based driver (for customized solutions only).....	1-78
Chapter 2	NPU runtime	2-86
2.1	About the NPU runtime	2-87
2.2	Arm NN runtime	2-88
2.3	Android neural networks runtime	2-90
2.4	TensorFlow Lite delegate runtime	2-94

Preface

This preface introduces the *Arm China Zhouyi Compass Driver and Runtime User Guide*.

It contains the following sections:

- [About this book on page 6.](#)
- [Feedback on page 8.](#)

About this book

This book describes how the Zhouyi Compass Diver and Runtime work and how to use them.

Using this book

This book is organized into the following chapters:

Chapter 1 NPU driver

This chapter describes the driver of Zhouyi NPUs.

Chapter 2 NPU runtime

This chapter describes the runtime of Zhouyi NPUs.

Glossary

The Arm® Glossary is a list of terms used in Arm China documentation, together with definitions for those terms. The Arm Glossary does not contain terms that are industry standard unless the Arm China meaning differs from the generally accepted meaning.

See the [Arm® Glossary](#) for more information.

Typographic conventions

italic

Introduces special terminology, denotes cross-references, and citations.

bold

Highlights interface elements, such as menu names. Denotes signal names. Also used for terms in descriptive lists, where appropriate.

`monospace`

Denotes text that you can enter at the keyboard, such as commands, file and program names, and source code.

`monospace`

Denotes a permitted abbreviation for a command or option. You can enter the underlined text instead of the full command or option name.

`monospace italic`

Denotes arguments to monospace text where the argument is to be replaced by a specific value.

`monospace bold`

Denotes language keywords when used outside example code.

<and>

Encloses replaceable terms for assembler syntax where they appear in code or code fragments. For example:

```
MRC p15, 0, <Rd>, <CRn>, <CRm>, <Opcode_2>
```

SMALL CAPITALS

Used in body text for a few terms that have specific technical meanings, that are defined in the *Arm® Glossary*. For example, IMPLEMENTATION DEFINED, IMPLEMENTATION SPECIFIC, UNKNOWN, and UNPREDICTABLE.

Additional reading

This book contains information that is specific to this product. See the following documents for other relevant information.

Arm China publications

The following confidential documents are only available to licensees:

- *Arm China Zhouyi Compass Getting Started Guide.*
- *Arm China Zhouyi Compass Software Technical Overview.*
- *Arm China Zhouyi Compass NN Compiler User Guide.*

Feedback

Feedback on this product

If you have any comments or suggestions about this product, contact your supplier and give:

- The product name.
- The product revision or version.
- An explanation with as much information as you can provide. Include symptoms and diagnostic procedures if appropriate.

Feedback on content

If you have comments on content, send an e-mail to errata@armchina.com. Give:

- The title *Arm China Zhouyi Compass Driver and Runtime User Guide*.
- The number 61010023_0001_00_en.
- If applicable, the page number(s) to which your comments refer.
- A concise explanation of your comments.

Arm China also welcomes general suggestions for additions and improvements.

_____ **Note** _____

Arm China tests the PDF only in Adobe Acrobat and Acrobat Reader, and cannot guarantee the quality of the represented document when used with any other PDF reader.

Chapter 1

NPU driver

This chapter describes the driver of the Zhouyi NPU.

It contains the following sections:

- *1.1 About the NPU driver on page 1-10.*
- *1.2 Linux-based driver on page 1-11.*
- *1.3 QNX-based driver on page 1-50.*
- *1.4 Bare-metal-based driver on page 1-68.*
- *1.5 RTOS-based driver (for customized solutions only) on page 1-78.*

1.1 About the NPU driver

The Zhouyi NPU driver is responsible for parsing middleware generated by an offline NN compiler and scheduling corresponding AI acceleration tasks to the NPU accelerator. Currently, the driver can support NPU controlling on Arm Linux, bare-metal, QNX and RTOS host platforms.

The following figure shows the components of a typical NPU driver.

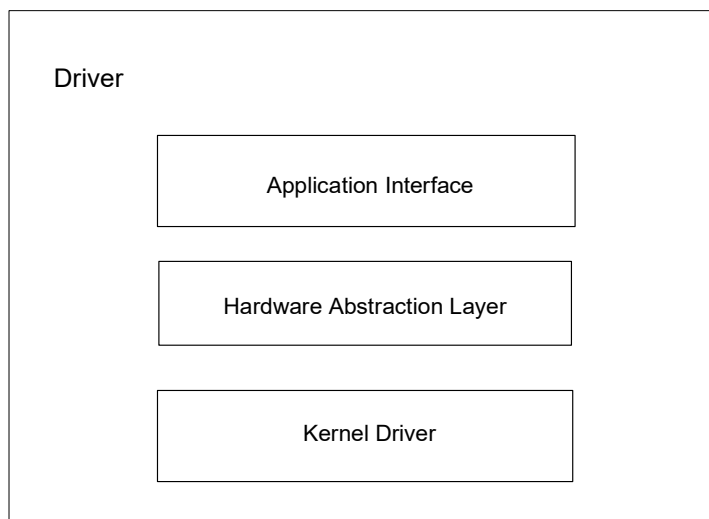


Figure 1-1 Components of a typical NPU driver

1.2 Linux-based driver

The Linux-based driver consists of two key parts—*User Mode Driver* (UMD) and *Kernel Mode Driver* (KMD). The UMD parses NPU job descriptors passed by applications, allocates resources and schedules jobs by calling KMD interfaces.

1.2.1 UMD

The UMD is compiled as a dynamic library and provides user applications a series of interfaces. Standard UMD APIs are designed for applications to schedule standard AI executable benchmarks onto the NPU.

User applications should call the APIs in a relatively fixed order and check the return values step by step to ensure correct operations.

The typical usage of the standard UMD interfaces in an application is as follows:

1. The application calls the context initialization API first to initialize UMD runtime context.
2. After initialization, the application should call the graph loading API to load an offline built binary. If successful, a graph ID is returned. Multiple graph binaries can be loaded in order.
3. After graph loading, the application can create and schedule an inference job which is bound to a previously loaded graph. All input, output, and intermediate buffers are allocated in this step.
4. The application should load the input frame data into the corresponding input buffers.
5. If the blocking scheduling API is used and returns without an error code, the application can get the output tensor directly from the tensor output buffer. If the non-blocking scheduling API is used, after querying the status of execution and no exception is found, the application can fetch the results in the same way.
6. The application should clean a finished job after it terminates, or reuse the job with new input data. The bind buffer can be reused in loading the inputs of the next frame and rescheduling the job again as described in steps 4–5.
7. Before exiting, the application should unload all graphs and destroy the context.

For more information about the corresponding data structures and interfaces, see *5.2.3 Linux driver user interface - Standard API*. Additionally some samples are supplied for reference in *driver/samples*.

1.2.2 KMD

The KMD works closely with the Linux kernel to respond to contiguous physical memory allocation requests from the UMD, controls job scheduling on the NPU, and handles NPU interrupts. The NPU KMD supports the management of both multiple NPU core instances with heterogeneous configurations (for AIPUV2/V3) and partition-cluster architecture (for AIPU V3). System software developers porting the NPU driver should implement the corresponding Linux device tree to tell the KMD the low-level hardware information. For more information about the device tree bindings and other usage/configuration options of the kernel driver, see the KMD readme document in the release package of driver source code.

To support UMD requests, like other char drivers, several standard Linux device file operation interfaces are provided to the UMD—open, release, mmap, ioctl and poll.

The ioctl options supported are as follows:

- `AIPU_IOCTL_QUERY_CAP`: Queries the common capability of NPUs.
- `AIPU_IOCTL_QUERY_PARTITION_CAP`: Queries the capability of an NPU core or partition.
- `AIPU_IOCTL_REQ_BUF`: Requests to allocate a coherent buffer.

- AIPU_IOCTL_FREE_BUF: Requests to free a coherent buffer allocated by AIPU_IOCTL_REQBUF.
- AIPU_IOCTL_DISABLE_SRAM: Disables the management of SoC SRAM in the kernel driver.
- AIPU_IOCTL_ENABLE_SRAM: Enables the management of SoC SRAM in the kernel driver disabled by AIPU_IOCTL_DISABLE_SRAM.
- AIPU_IOCTL_SCHEDULE_JOB: Schedules a user job to the kernel mode driver for execution.
- AIPU_IOCTL_QUERY_STATUS: Queries the execution status of one or multiple scheduled jobs.
- AIPU_IOCTL_KILL_TIMEOUT_JOB: Kills a timeout job and cleans it from the kernel mode driver.
- AIPU_IOCTL_REQ_IO: Reads/writes an external register of an NPU core.
- AIPU_IOCTL_GET_HW_STATUS: Gets the working status of NPU cores (for example, busy or idle).
- AIPU_IOCTL_ABORT_CMD_POOL: Aborts the NPU command pools (for AIPUv3 only).
- AIPU_IOCTL_DISABLE_TICK_COUNTER: Disables the tick counter (for AIPUv3 only).
- AIPU_IOCTL_ENABLE_TICK_COUNTER: Enables the tick counter (for AIPUv3 only).
- AIPU_IOCTL_CONFIG_CLUSTERS: Configures the NPU clusters, for example, disabling cores (for AIPUv3 only).
- AIPU_IOCTL_ALLOC_DMA_BUF: Requests to allocate a buffer using the dma-buf framework.
- AIPU_IOCTL_FREE_DMA_BUF: Frees a buffer allocated with the dma-buf allocation ioctl.
- AIPU_IOCTL_GET_DMA_BUF_INFO: Gets description of a buffer allocated with the dma-buf allocation ioctl.
- AIPU_IOCTL_GET_DRIVER_VERSION: Gets a string that contains the driver version number.

To port the KMD on a device, refer to the guide [AI610-SDK-1012/Linux-driver/driver/kmd/porting_guide.txt](#).

1.2.3 Linux driver user interface - Standard API

Data structure

Data structure	Type	Members	Comment
device_status_t	enum	DEV_IDLE (0x0)	-
		DEV_BUSY	-
		DEV_EXCEPTION	-

Data structure	Type	Values	Comment
aipu_data_type_t	enum	TENSOR_DATA_TYPE_NONE (0x0)	No data type.
		AIPU_DATA_TYPE_BOOL(0x1)	Bool type.
		TENSOR_DATA_TYPE_U8 (0x2)	Unsigned int8.
		TENSOR_DATA_TYPE_S8 (0x3)	Signed int8.
		TENSOR_DATA_TYPE_U16 (0x4)	Unsigned int16.
		TENSOR_DATA_TYPE_S16 (0x5)	Signed int16.
		AIPU_DATA_TYPE_U32 (0x6)	Unsigned int32.
		AIPU_DATA_TYPE_S32 (0x7)	Signed int32.
		AIPU_DATA_TYPE_U64 (0x8)	Unsigned int64.
		AIPU_DATA_TYPE_S64 (0x9)	Signed int64.
		AIPU_DATA_TYPE_f16 (0xA)	Float16.

		AIPU_DATA_TYPE_f32 (0xB)	Float32.
		AIPU_DATA_TYPE_f64 (0xC)	Float64.

Data structure	Type	Members	Comment
aipu_tensor_type_t	enum	AIPU_TENSOR_TYPE_INPUT(0)	-
		AIPU_TENSOR_TYPE_OUTPUT	-
		AIPU_TENSOR_TYPE_INTER_DUMP	-
		AIPU_TENSOR_TYPE_PRINTF	-
		AIPU_TENSOR_TYPE_PROFILER	-
		AIPU_TENSOR_TYPE_LAYER_COUNTER	-
		AIPU_TENSOR_TYPE_ERROR_CODE	-

Data structure	Type	Members	Member type	Comment
aipu_tensor_desc_t	struct	id	uint32_t	Tensor ID.
		size	Uint32_t	Tensor size.
		scale	float	Tensor scale parameter.
		zero_point	float	Tensor zero-point parameter.
		data_type	aipu_data_type_t	Tensor format information.

Data structure	Type	Members	Member type	Comment
aipu_global_config_simulation_t	struct	simulator	const char*	Simulator executable full name.
		x2_arch_desc	const char*	Simulated aipu_v3 ARCH.
		log_file_path	const char*	The file path to store log files.
		log_level	uint32_t	Log level.
		gm_size	uint32_t	GM size for the aipu_v3 simulator
		verbose	bool	Verbose switch.
		enable_avx	bool	Simulator ENABLE_AVX switch.
		enable_calloc	bool	Simulator ENABLE_CALLOC switch.
		en_eval	bool	Simulator ENABLE_EVAL switch.

Data structure	Type	Members	Member type	Comment
aipu_job_config_simulation_t	struct	data_dir	const char*	The data path to store simulation dumped data files.

Data structure	Type	Members	Member type	Comment
aipu_job_config_dump_t	struct	dump_dir	const char*	The data path to store dumped data files.
		prefix	const char *	The file name prefix for dumped files.
		output_prefix	const char*	The file name prefix for dumped output files (can be ignored).

		misc_prefix	const char*	The file name prefix for dumped profile files.
--	--	-------------	-------------	--

Data structure	Type	Members	Member type	Comment
aipu_create_job_cfg_t	struct	partition_id	uint8_t	Indicates which cluster partition the created job is committed to. Only for aipu_v3.
		dbg_dispatch	uint8_t	Debug dispatch flag. When set to 1, it indicates that the job is assigned to the debug core to run.
		dbg_core_id	uint8_t	Specifies the debug core ID [0, max_core_id in the cluster].
		qos_level	uint8_t	Indicates which QoS level the created job is assigned with (low or high). Only for aipu_v3.
		fm_mem_region	uint8_t	Controls which region the feature map buffer memory is allocated from. By default, it is allocated from DDR. If configured as AIPU_MEM_REGION_SRAM, it is allocated from SRAM (if available and configured on DTS). If configured as AIPU_MEM_REGION_DTCM, it is allocated from DTCM (if available and configured on DTS and the CPU can access).
		wt_mem_region	uint8_t	Controls which region the weight buffer memory is allocated from. By default, it is allocated from DDR. If configured as AIPU_MEM_REGION_SRAM, it is allocated from SRAM (if available and configured on DTS). If configured as AIPU_MEM_REGION_DTCM, it is allocated from DTCM (if available and configured on DTS and the CPU can access).
		fm_idxes	int32_t *	Specifies feature maps allocated from fm_mem_region.
		fm_idxes_cnt	int32_t	The element number in fm_idxes.
		wt_idxes	int32_t *	Specifies weights allocated from wt_mem_region.
		wt_idxes_cnt	int32_t	The element number in wt_idxes.

Data structure	Type	Values	Comment
aipu_job_status_t	enum	AIPU_JOB_STATUS_NO_STATUS (0x0)	-
		AIPU_JOB_STATUS_DONE (0x1)	-
		AIPU_JOB_STATUS_EXCEPTION (0x2)	-

Data structure	Type	Members	Member type	Comment
aipu_debugger_job_info_t	struct	instr_base	uint64_t	Instruction section base address (physical).
		simulation_aipu	void *	NPU simulator object.
		simulation_mem_engine	void *	Simulation mem_engine object.

Data structure	Type	Members	Member type	Comment
aipu_shared_tensor_info_t	struct	type	aipu_tensor_type_t	Type of the shared tensor: input or output.
		tensor_idx	uint32_t	The index number of the shared tensor.
		id	uint64_t	Job ID: on marking one IO buffer as shared.

				Graph ID: on specifying the shared buffer to a new graph. (Ignored for dma_buf)
		pa	uint64_t	<ul style="list-style-type: none"> On marking the shared buffer: return the buffer address of the marked shared IO buffer. On specifying the shared buffer: transfer the shared buffer address to a new graph. (Ignored for dma_buf)
		dmabuf_fd	int	The file ID of a certain dma_buf.
		offset_in_dmabuf	uint32_t	The valid offset in dma_buf referenced by dmabuf_fd.

Data structure	Type	Members	Member type	Comment
aipu_dmabuf_op	struct	dmabuf_fd	int	The file fd of dma_buf.
		offset_in_dmabuf	uint32_t	The dma_buf referenced by dmabuf_fd.
		size	uint32_t	Filled data size or fetched data size.
		data	char *	The buffer of data that is written to dma_buf or that is read back from dma_buf.

Data structure	Type	Members	Member type	Comment
aipu_core_info_t	struct	umd_version[16]	char	The buffer to store the UMD version number.

Data structure	Type	Members	Member type	Comment
aipu_core_info_t	struct	reg_base	uint64_t	Base address of the NPU core external register.

Data structure	Type	Values	Comment
aipu_ioctl_cmd_t	enum	AIPU_IOCTL_MARK_SHARED_TENSOR	Mark one IO buffer of the job as a shared one.
		AIPU_IOCTL_SET_SHARED_TENSOR	Specify the previously marked shared buffer to replace the IO buffer of the new graph.
		AIPU_IOCTL_SET_PROFILE	Dynamically turn on or off the profiling feature in simulation.
		AIPU_IOCTL_ALLOC_DMABUF	Request one dma_buf from KMD.
		AIPU_IOCTL_FREE_DMABUF	Release one dma_buf to KMD.
		AIPU_IOCTL_WRITE_DMABUF	Write data to one dma_buf.
		AIPU_IOCTL_READ_DMABUF	Read data from one dma_buf.
		AIPU_IOCTL_GET_VERSION	Get the version number of UMD and KMD.

Data structure	Type	Values	Comment
aipu_config_type_t	enum	AIPU_JOB_CONFIG_TYPE_DUMP_TEXT	Dump text section(s).

		AIPU_JOB_CONFIG_TYPE_DUMP_WEIGHT	Dump weight section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_RODATA	Dump rodata section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_DESCRIPTOR	Dump descriptor section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_INPUT	Dump input data section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_OUTPUT	Dump output data tensor(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_TCB_CHAIN	Dump TCB chain tensor(s). Not used for Zhouyi aipu_v1 and aipu_v2.
		AIPU_JOB_CONFIG_TYPE_DUMP_EMULATION	Dump emulation files. Internally used.
		AIPU_CONFIG_TYPE_SIMULATION	Simulation configuration switch.
		AIPU_CONFIG_TYPE_HW	The switch of configuration on hardware.
		AIPU_GLOBAL_CONFIG_TYPE_DISABLE_VER_CHECK	Disable graph bin version check.
		AIPU_GLOBAL_CONFIG_TYPE_ENABLE_VER_CHECK	Enable graph bin version check.

Data structure	Type	Values	Comment
aipu_status_t	enum	AIPU_STATUS_SUCCESS (0x0)	The execution of this API returns successfully without any error.
		AIPU_STATUS_ERROR_NULL_PTR	The arguments passed by user applications to UMD API are NULL pointers which are invalid.
		AIPU_STATUS_ERROR_INVALID_CTX	The context pointer passed by user applications to UMD API is invalid.
		AIPU_STATUS_ERROR_OPEN_DEV_FAIL	UMD fails in opening the device file /dev/aipu which may result from the failure in probing the KMD module.
		AIPU_STATUS_ERROR_DEV_ABNORMAL	The NPU device is in abnormal state.
		AIPU_STATUS_ERROR_DEINIT_FAIL	UMD API fails in de-initiating a context.
		AIPU_STATUS_ERROR_UNKNOWN_BIN	The type of the binary loaded is unknown. It cannot be loaded or executed.
		AIPU_STATUS_ERROR_INVALID_CONFIG	The configuration passed by a user application to UMD API is invalid.
		AIPU_STATUS_ERROR_GVERSION_UNSUPPORTED	The version of an executable graph binary passed by a user application is not supported on the current UMD.
		AIPU_STATUS_ERROR_TARGET_NOT_FOUND	UMD fails in finding a matching NPU hardware target for the executable graph passed by a user application to execute.
		AIPU_STATUS_ERROR_INVALID_GBIN	The executable graph binary file passed by a user application to UMD API contains invalid items which cannot be parsed or executed.
		AIPU_STATUS_ERROR_INVALID_GRAPH_ID	Graph ID provided is an invalid one which has been unloaded or never existed.
		AIPU_STATUS_ERROR_OPEN_FILE_FAIL	UMD fails in opening a file passed by a user application.

	AIPU_STATUS_ERROR_MAP_FILE_FAIL	UMD fails in mapping a file passed by a user application.
	AIPU_STATUS_ERROR_READ_FILE_FAIL	UMD fails in reading a file passed by a user application.
	AIPU_STATUS_ERROR_WRITE_FILE_FAIL	UMD fails in writing a file passed by a user application.
	AIPU_STATUS_ERROR_INVALID_JOB_ID	Job ID provided is an invalid one which has been cleaned or never existed.
	AIPU_STATUS_ERROR_JOB_EXCEPTION	The execution of a job ID passed by a user application to UMD API ends with an exception.
	AIPU_STATUS_ERROR_JOB_TIMEOUT	The execution of a job ID passed by a user application to UMD API timed out.
	AIPU_STATUS_ERROR_OP_NOT_SUPPORTED	The operation is not supported in the current UMD version or system environment.
	AIPU_STATUS_ERROR_INVALID_OP	The operation is invalid and not allowed in the current system environment.
	AIPU_STATUS_ERROR_INVALID_SIZE	The size argument is invalid.
	AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	The buffer allocation requests fail because of system memory limitations.
	AIPU_STATUS_ERROR_BUF_FREE_FAIL	The buffer free requests fail because the buffers are in busy state.
	AIPU_STATUS_ERROR_INVALID_CORE_ID	The AIPU core ID that the application provides is invalid and cannot be found in the system.
	AIPU_STATUS_ERROR_RESERVE_SRAM_FAIL	UMD fails in reserving SRAM as the executable binary requested because there is no SoC SRAM or SRAM is busy.
	AIPU_STATUS_ERROR_INVALID_TENSOR_ID	The tensor ID that the application provides is invalid.
	AIPU_STATUS_ERROR_INVALID_CLUSTER_ID	The AIPU cluster ID that the application provides is invalid and cannot be found in the system.
	AIPU_STATUS_ERROR_INVALID_PARTITION_ID	The AIPU partition ID that the application provides is invalid.
	AIPU_STATUS_ERROR_PRINTF_FAIL	UMD fails in parsing the printf buffer and printing corresponding logs.
	AIPU_STATUS_ERROR_INVALID_TENSOR_TYPE	The specified tensor type is invalid for this operation.
	AIPU_STATUS_ERROR_INVALID_GM	The GM is invalid.
	AIPU_STATUS_ERROR_INVALID_SEGMMU	The segmented MMU is invalid.
	AIPU_STATUS_ERROR_INVALID_QOS	The specified QoS level is invalid.
	AIPU_STATUS_ERROR_INVALID_TENSOR_CNT	The specified tensor count is invalid.
	AIPU_STATUS_ERROR_TIMEOUT	Timeout on the polling job's status.
	AIPU_STATUS_ERROR_NO_BATCH_QUEUE	There is no specific batch queue.

		AIPU_STATUS_ERROR_MARK_SHARED_TENSOR	Mark shared tensor: no corresponding tensor.
		AIPU_STATUS_ERROR_SET_SHARED_TENSOR	Set shared tensor: no corresponding tensor.
		AIPU_STATUS_MAX	Maximum error code value.

Data structure	Type	Members	Member type	Comment
aipu_ctx_handle_t	struct	opaque	-	The pointer to this opaque context is returned by the init API.

User application interfaces

Function declaration	aipu_status_t aipu_init_context(aipu_ctx_handle_t** ctx);	
Parameter	ctx	Pointer to a memory location allocated by the application where the UMD stores the opaque context handle struct.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_OPEN_DEV_FAIL AIPU_STATUS_ERROR_DEV_ABNORMAL	
Comments	This API is used to initialize the NPU UMD context.	

Function declaration	aipu_status_t aipu_get_error_message (const aipu_ctx_handle_t* ctx, aipu_status_t status, const char** msg);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	status	Status returned by UMD standard API.
	msg	Pointer to a memory location allocated by the application where the UMD stores the message string pointer.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX	
Comments	This API is used to query additional information about a status returned by UMD API.	

Function declaration	aipu_status_t aipu_deinit_context(const aipu_ctx_handle_t* ctx);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_DEINIT_FAIL	
Comments	This API is used to destroy the NPU UMD context.	

Function declaration	aipu_status_t aipu_config_global(const aipu_ctx_handle_t* ctx, uint64_t types, void* config)		
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.	
	types	Configuration type(s)	
	config	Pointer to a memory location allocated by the application where stores the configurations.	

Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_CONFIG
Comments	This API is used to configure a specified global option.

Function declaration	aipu_status_t aipu_load_graph(const aipu_ctx_handle_t* ctx, const char* graph, uint64_t* id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph	Pointer to a memory location allocated by the application where stores the graph.
	id	Graph ID returned to the caller.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_OPEN_FILE_FAIL AIPU_STATUS_ERROR_MAP_FILE_FAIL AIPU_STATUS_ERROR_UNKNOWN_BIN AIPU_STATUS_ERROR_GVERSION_UNSUPPORTE AIPU_STATUS_ERROR_TARGET_NOT_FOUND AIPU_STATUS_ERROR_INVALID_GBIN AIPU_STATUS_ERROR_BUF_ALLOC_FAIL AIPU_STATUS_ERROR_RESERVE_SRAM_FAIL AIPU_STATUS_ERROR_INVALID_GM	
Comments	This API loads an offline built NPU executable graph binary from file system.	

Function declaration	aipu_status_t aipu_load_graph_helper(const aipu_ctx_handle_t* ctx, const char* graph_buf, uint32_t graph_size, uint64_t* id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph_buf	Loadable graph binary file data.
	graph_size	The byte size of loadable graph binary data.
	id	Graph ID returned to the caller.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_OPEN_GBIN_FAIL AIPU_STATUS_ERROR_MAP_GBIN_FAIL Other values returned by AIPU_load_graph	
Comments	This API loads a graph from the corresponding data buffer.	

Function declaration	aipu_status_t aipu_unload_graph(const aipu_ctx_handle_t* ctx, uint64_t id);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	id	Graph ID
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX	

	AIPU_STATUS_ERROR_INVALID_GRAPH_ID
Comments	This API is used to unload a loaded graph.

Function declaration	aipu_status_t aipu_create_job(const aipu_ctx_handle_t* ctx, uint64_t graph, uint64_t* job, aipu_create_job_cfg_t *config = nullptr)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph	Graph ID returned by aipu_load_graph.
	job	Pointer to a memory location allocated by the application where the UMD stores the created job ID.
	config	<ol style="list-style-type: none"> 1. Specify the partition ID and QoS level of the job, only for Zhouyi X2. 2. Specify the memory region for the feature map and weight buffer. 3. Bind the job to one core in the cluster. [aipu_v3 only] 4. Specify which weight buffer or feature map buffer to be allocated from the specified weight region or feature map region.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	
Comments	This API is used to create a job.	

Function declaration	aipu_status_t aipu_finish_job(const aipu_ctx_handle_t* ctx, uint64_t job_id, int32_t time_out)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	job_id	Job ID returned by aipu_create_job.
	time_out	The timeout parameter for the poll system call.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_OP AIPU_STATUS_ERROR_JOB_EXCEPTION AIPU_STATUS_ERROR_JOB_TIMEOUT	
Comments	This API is used to flush a new computation job onto the NPU. This API is a blocking one which will schedule a job and block in waiting for it to be done. The application can safely fetch the inference result data if this API returns successfully.	

Function declaration	aipu_status_t aipu_flush_job(const aipu_ctx_handle_t* ctx, uint64_t job_id, aipu_job_callback_func_t cb_func = nullptr)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	job_id	Job ID returned by aipu_create_job.
	cb_func	Callback function to handle this job. Prototype: int (*aipu_job_callback_func_t)(uint64_t job_id, aipu_job_status_t job_state)

Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_OP
Comments	This API is used to flush a new computation job onto the NPU. Unlike <code>aipu_finish_job</code> , this API returns as soon as a job is scheduled. By using this API with multiple tensor buffers' ping-pong operation, the pipeline feature will be enabled. The application should query the job done status before fetching the inference result data. It is optional to assign a callback function to handle the done job timely.

Function declaration	<code>aiipu_status_t aipu_get_job_status(const aipu_ctx_handle_t* ctx, uint64_t job_id, aipu_job_status_t* status, int32_t timeout = 0)</code>	
Parameter	<code>ctx</code>	Pointer to a context handle struct returned by <code>aipu_init_context</code> .
	<code>job_id</code>	Job ID returned by <code>aipu_create_job</code> .
	<code>status</code>	Pointer to a memory location allocated by the application where the UMD stores the job status.
	<code>timeout</code>	timeout value (ms) to poll the job's status.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID	
Comments	<p>This API is used to get the execution status of a job scheduled via <code>aipu_flush_job</code>. This API is a blocking one which will block in waiting for the job to be done if it is still running.</p> <ul style="list-style-type: none"> <code>timeout > 0</code>: The max polling time window is 'timeout'. <code>timeout = 0</code>: Non-blocking and return the job's status immediately. <code>timeout = -1</code>: Blocking until the job is really done or an exception occurs. 	

Function declaration	<code>aiipu_status_t aipu_clean_job(const aipu_ctx_handle_t* ctx, uint64_t id)</code>	
Parameter	<code>ctx</code>	Pointer to a context handle struct returned by <code>aipu_init_context</code> .
	<code>id</code>	Job ID returned by <code>aipu_create_job</code> .
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID	
Comments	This API is used to clean a finished job object in UMD. After this API successfully returns, the job ID immediately becomes invalid and cannot be used again. After this API successfully returns, the buffer handle used in creating this job will immediately be free to be used to create another job with the same graph ID.	

Function declaration	<code>aiipu_status_t aipu_get_tensor_count(const aipu_ctx_handle_t* ctx, uint64_t id, aipu_tensor_type_t type, uint32_t* cnt)</code>	
Parameter	<code>ctx</code>	Pointer to a context handle struct returned by <code>aipu_init_context</code> .
	<code>id</code>	Job ID returned by <code>aipu_create_job</code> .

	type	Tensor type.
	cnt	Pointer to a memory location allocated by the application where the UMD stores the count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to get the tensor count of specified type.	

Function declaration	aipu_status_t aipu_get_tensor_descriptor(const aipu_ctx_handle_t* ctx, uint64_t id, aipu_tensor_type_t type, uint32_t tensor, aipu_tensor_desc_t* desc)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	id	Job ID returned by aipu_create_job.
	type	Tensor type.
	tensor	Tensor ID.
	desc	Pointer to a memory location allocated by the application where the UMD stores the tensor descriptor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID	
Comments	This API is used to get the tensor descriptor of specified type.	

Function declaration	aipu_status_t aipu_load_tensor(const aipu_ctx_handle_t* ctx, uint64_t id, uint32_t tensor, const void* data)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	id	Job ID returned by aipu_create_job.
	tensor	Tensor ID.
	data	Data of the input tensor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to load the input tensor data.	

Function declaration	aipu_status_t aipu_get_tensor(const aipu_ctx_handle_t* ctx, uint64_t job, aipu_tensor_type_t type, uint32_t tensor, void* buf);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	id	Job ID returned by aipu_create_job.
	type	Tensor type.

	tensor	Tensor ID.
	buf	Pointer to a memory location allocated by the application where the UMD stores the tensor data.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get the tensor data of specified type.	

Function declaration	aipu_status_t aipu_config_job(const aipu_ctx_handle_t* ctx, uint64_t id, uint64_t types, void* config)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	id	Job ID returned by aipu_create_job.
	type	Configuration type(s)
	config	Pointer to a memory location allocated by the application where the application stores the configuration data struct.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_CONFIG	
Comments	This API is used to configure a specified option of a job.	

Function declaration	aipu_status_t aipu_specify_iobuf(const aipu_ctx_handle_t* ctx, uint64_t id, aipu_shared_tensor_info_t *shared_buf)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	id	Job ID returned by aipu_create_job.
	shared_buf	Specify one dma_buf or non-dma_buf as the IO buffer of the job.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID AIPU_STATUS_ERROR_INVALID_SIZE AIPU_STATUS_ERROR_DMABUF_SHARED_IO	
Comments	This API is used to specify a shared buffer as the input or output buffer of the job. The shared buffer has to be one dma_buf currently.	

Function declaration	aipu_status_t aipu_get_partition_count(const aipu_ctx_handle_t* ctx, uint32_t* cnt)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.

	cnt	Pointer to a memory location allocated by the application where the UMD stores the partition count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get the NPU partition count.	

Function declaration	aipu_status_t aipu_get_cluster_count(const aipu_ctx_handle_t* ctx, uint32_t partition_id, uint32_t* cnt)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	partition_id	Indicates which partition it gets the cluster count from. Only one partition currently.
	cnt	Pointer to a memory location allocated by the application where the UMD stores the cluster count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get the NPU cluster count.	

Function declaration	aipu_status_t aipu_get_core_count(const aipu_ctx_handle_t* ctx, uint32_t partition_id, uint32_t cluster, uint32_t* cnt)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	partition_id	Partition ID.
	cluster	Cluster ID.
	cnt	Pointer to a memory location allocated by the application where the UMD stores the core count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_OP AIPU_STATUS_ERROR_INVALID_CLUSTER_ID	
Comments	This API is used to get the NPU core count in a specific cluster.	

Function declaration	aipu_status_t aipu_debugger_get_core_info(const aipu_ctx_handle_t* ctx, uint32_t core_id, aipu_core_info_t* info)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	core_id	ID of an NPU core.
	info	Pointer to a memory location allocated by the application where the UMD stores the core information.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_CORE_ID	

	AIPU_STATUS_ERROR_DEV_ABNORMAL
Comments	This API is used to get information about an NPU core for the debugger to use.

Function declaration	aipu_status_t aipu_debugger_get_job_info(const aipu_ctx_handle_t* ctx, uint64_t job, aipu_debugger_job_info_t* info)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	job_id	Job ID.
	info	Pointer to a memory location allocated by the application where the UMD stores a pointer to the job information.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID	
Comments	This API is used to get the number of NPU cores in a system for the debugger to use. The NPU core IDs are numbered as [0, cnt - 1), respectively. This API should be called by the debugger before calling any other APIs for the debugger to use.	

Function declaration	aipu_status_t aipu_debugger_bind_job(const aipu_ctx_handle_t* ctx, uint32_t core_id, uint64_t job_id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	core_id	ID of an NPU core.
	job_id	Job ID returned by aipu_create_job.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_CORE_ID AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API binds a created job to an idle NPU core for execution later. External registers of the specified NPU core is written after this API returns, but the start PC register is not triggered to run. For the same core or job, it should only be bound once, unless the job is done. The core to be bound should be in idle state, otherwise UMD returns error code	

Function declaration	aipu_status_t aipu_debugger_run_job(const aipu_ctx_handle_t* ctx, uint32_t job_id);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	job_id	Job ID returned by aipu_create_job.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_JOB_NOT_EXIST AIPU_STATUS_ERROR_JOB_EXCEPTION AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API triggers a previously bound job to run on a target NPU core. A debugger job must be bound before running.	

	<p>After this API successfully returns, the start PC register of the NPU core previously bound with the job specified by <code>job_id</code> is triggered.</p> <p>This API is a blocking API which returns after the job execution ends on hardware.</p> <p>A debugger job should only be scheduled by <code>aipu_debugger_run_job()</code>. You cannot bind a job by <code>aipu_debugger_bind_job()</code> but schedule it by calling <code>aipu_flush_job()</code> or <code>aipu_finish_job()</code>.</p>
--	---

Function declaration	<code>aipu_status_t aipu_debugger_malloc(const aipu_ctx_handle_t* ctx, uint32_t size, void** va)</code>	
Parameter	<code>ctx</code>	Pointer to a context handle struct returned by <code>aipu_init_context</code> .
	<code>size</code>	Size of the requested buffer in bytes.
	<code>va</code>	Pointer to a virtual address where stores the base address of the buffer.
Return value	<p>AIPU_STATUS_SUCCESS</p> <p>AIPU_STATUS_ERROR_NULL_PTR</p> <p>AIPU_STATUS_ERROR_INVALID_CTX</p> <p>AIPU_STATUS_ERROR_INVALID_SIZE</p> <p>AIPU_STATUS_ERROR_BUF_ALLOC_FAIL</p>	
Comments	<p>This API allocates a buffer for the NPU debugger to use.</p> <p>This API shall be used after <code>aipu_load_graph</code> and before calling <code>aipu_debugger_bind_job</code>.</p>	

Function declaration	<code>aipu_status_t aipu_debugger_free(const aipu_ctx_handle_t* ctx, void* va)</code>	
Parameter	<code>ctx</code>	Pointer to a context handle struct returned by <code>aipu_init_context</code> .
	<code>va</code>	Virtual address to be freed.
Return value	<p>AIPU_STATUS_SUCCESS</p> <p>AIPU_STATUS_ERROR_NULL_PTR</p> <p>AIPU_STATUS_ERROR_INVALID_CTX</p> <p>AIPU_STATUS_ERROR_BUF_FREE_FAIL</p>	
Comments	This API frees a buffer allocated by <code>aipu_debugger_malloc</code> .	

Function declaration	<code>aipu_status_t aipu_printf(char* printf_base, char* redirect_file)</code>	
Parameter	<code>printf_base</code>	Pointer to a tensor buffer where stores the printf log data.
	<code>redirect_file</code>	Printf output redirect file path.
Return value	<p>AIPU_STATUS_SUCCESS</p> <p>AIPU_STATUS_ERROR_NULL_PTR</p> <p>AIPU_STATUS_ERROR_PRINTF_FAIL</p>	
Comments	<p>This API prints NPU execution log information after the corresponding job ends. The application can choose to redirect the printf log to the terminal or a file. It passes NULL to the <code>redirect_file</code> results in printing logs to the terminal directly. Nothing will be printed if the printf function is not used in the NPU executable binary loaded by <code>aipu_load_graph</code>.</p>	

Function declaration	<code>aipu_status_t aipu_create_batch_queue(const aipu_ctx_handle_t* ctx, uint64_t graph_id, uint32_t *queue_id)</code>	
Parameter	<code>ctx</code>	Pointer to a context handle struct returned by <code>aipu_init_context</code> .
	<code>graph_id</code>	Graph ID returned by <code>aipu_load_graph</code> .

	queue_id	Pointer to store the batch queue ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to create a batch queue.	

Function declaration	aipu_status_t aipu_clean_batch_queue(const aipu_ctx_handle_t *ctx, uint64_t graph_id, uint32_t queue_id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph_id	Graph ID returned by aipu_load_graph.
	queue_id	Queue ID returned by aipu_create_batch_queue.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to clean a specific batch queue.	

Function declaration	aipu_status_t aipu_config_batch_dump(const aipu_ctx_handle_t *ctx, uint64_t graph_id, uint32_t queue_id, uint64_t types, aipu_job_config_dump_t *dump_cfg)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph_id	Graph ID returned by aipu_load_graph.
	queue_id	Queue ID returned by aipu_create_batch_queue.
	types	Dump options for each batch job.
	dump_cfg	The root path to store dump files of each batch.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_NO_BATCH_QUEUE	
Comments	This API is used to perform basic configuration for batch inference on the simulator and hardware.	

Function declaration	aipu_status_t aipu_add_batch(const aipu_ctx_handle_t *ctx, uint64_t graph_id, uint32_t queue_id, char *inputs[], char *outputs[])	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph_id	Graph ID returned by aipu_load_graph.
	queue_id	Queue ID returned by aipu_create_batch_queue.
	inputs	Buffer pointers for input tensors.
	outputs	Buffer pointers for output tensors.
Return value	AIPU_STATUS_SUCCESS	

	AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_NO_BATCH_QUEUE
Comments	This API is used to add a batch buffer for one frame inference.

Function declaration	aipu_status_t aipu_finish_batch(const aipu_ctx_handle_t *ctx, uint64_t graph_id, uint32_t queue_id, aipu_create_job_cfg_t *create_cfg);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	graph_id	Graph ID returned by aipu_load_graph.
	queue_id	Queue ID returned by aipu_create_batch_queue.
	create_cfg	Config for all batches in one queue.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_NO_BATCH_QUEUE	
Comments	This API is used to run multiple batch inferences.	

Function declaration	aipu_status_t aipu_ioctl(aipu_ctx_handle_t *ctx, uint32_t cmd, void *arg = nullptr);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_context.
	cmd	<p>AIPU_IOCTL_MARK_SHARED_TENSOR: Mark a shared buffer from a belonged job, arg {aipu_shared_tensor_info_t}. The marked buffer is not freed on destroying a job and directly used by a new job.</p> <p>AIPU_IOCTL_SET_SHARED_TENSOR: Specify a marked shared buffer to a new job, arg {aipu_shared_tensor_info_t}. Marking a shared buffer and setting a shared buffer to a new job must be performed in the same one process context. See the sample in <code>samples/src/sharebuffer_test</code>.</p> <p>AIPU_IOCTL_ENABLE_TICK_COUNTER: Enable the performance counter, no arg.</p> <p>AIPU_IOCTL_DISABLE_TICK_COUNTER: Disable the performance counter, no arg.</p> <p>AIPU_IOCTL_CONFIG_CLUSTERS: Configure the number of enabled cores in a cluster, arg {struct aipu_config_clusters}. It is used to disable the clock of some cores to reduce power consumption.</p> <p>AIPU_IOCTL_SET_PROFILE: Dynamically enable or disable the profiling feature of AIPUv3 simulation, arg {1/0}. 1: Enable profiling. 0: Disable profiling.</p> <p>AIPU_IOCTL_ALLOC_DMABUF: Request dma_buf from KMD, arg {struct aipu_dma_buf_request}. aipu_dma_buf_request->bytes: Request size (filled by UMD). aipu_dma_buf_request->fd: fd corresponding to dma_buf (filled by KMD).</p>

		AIPU_IOCTL_FREE_DMABUF: Free a dma_buf with its fd, arg { fd }. Refer to the samples in samples/src/{dmabuf_mmap_test/dmabuf_vmap_test /dmabuf_dma_test}.
	arg	Input or output argument.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_DEV_ABNORMAL	
Comments	This API is used to send a specific command to the driver. [aipu_v3 only]	

1.2.4 Linux driver user interface - Python API

Public APIs

Function	NPU()
Return value	Returns an NPU class object.
Comments	This API should be called first of all.

Class NPU methods

Function	aipu_init_context()	
Parameter	None.	
Return value	An integer.	
Comments	This method is used to open a device and initialize the UMD context.	

Function	aipu_deinit_context()	
Parameter	None.	
Return value	An integer.	
Comments	This method is used to destroy the UMD context.	

Function	aipu_get_error_message(status)	
Parameter	status	Status returned by other NPU methods.
Return value	A string.	
Comments	This method is used to query additional information about a status returned by NPU methods.	

Function	aipu_config_global(types, global_config_simulation)	
Parameter	types	Configuration type.
	global_config_simulation	A dictionary contains the configurations.
Return value	An integer.	
Comments	This method is used to configure a specified global option.	

Function	aipu_load_graph(bin_file)	
Parameter	bin_file	The path of the AIPU bin file.
Return value	A dictionary contains the returned status and graph ID.	
Comments	This method loads an offline built NPU executable graph binary from the file system.	

Function	aipu_load_graph(graph_data, length)	
Parameter	graph_data	Returned value by file.read().
	length	The length of the graph_data.
Return value	A dictionary contains the returned status and graph ID.	
Comments	This method loads an offline built NPU executable graph binary from the buffer.	

Function	aipu_unload_graph(graph_id)	
Parameter	graph_id	Returned value by method aipu_load_graph.
Return value	An integer.	
Comments	This method is used to unload a loaded graph.	

Function	aipu_create_job(graph_id, job_config, fm_idxe, wt_idxe)	
Parameter	graph_id	Returned value by method aipu_load_graph.
	job_config	1. Specify the partition ID and QoS level of the job (only for aipu v3).
	fm_idxe	2. Specify the memory region for the feature map and weight buffer.
	wt_idxe	3. Bind the job to one core in the cluster. [aipu_v3 only] 4. Specify which weight buffer or feature map buffer to be allocated from the specified weight region or feature map region.
Return value	A dictionary contains the returned status and job ID.	
Comments	This method is used to create a job.	

Function	aipu_config_job(job_id, type, job_config_dump)	
Parameter	job_id	Returned value by method aipu_create_job.
	type	Configuration type.
	job_config_dump	A dictionary contains the configuration information.
Return value	A dictionary contains the returned status and job ID.	
Comments	This method is used to create a job.	

Function	aipu_finish_job(job_id, timeout)	
Parameter	job_id	Returned value by method aipu_create_job.
	timeout	The number of milliseconds that the method should block in waiting for the job to complete.
Return value	An integer.	

Comments	This method is used to flush a new computation job onto the NPU. This method is a blocking one which will schedule a job and block in waiting for it to be done. The application can safely fetch the inference result data if this method returns successfully.	
-----------------	--	--

Function	aipu_flush_job(job_id, py_cb)	
Parameter	job_id	Returned value by method aipu_create_job.
	py_cb	Callback function to handle this job.
Return value	An integer.	
Comments	This method is used to flush a new computation job onto the NPU. Unlike aipu_finish_job, it returns as soon as a job is scheduled. By using this method with multiple tensor buffers' ping-pong operation, the pipeline feature will be enabled. The application should query the job done status before fetching the inference result data. It is optional to assign a callback function to handle the done job timely.	

Function	aipu_get_job_status(job_id, timeout)	
Parameter	job_id	Returned value by method aipu_create_job.
	timeout	The number of milliseconds that the method should block in waiting for the job status.
Return value	A dictionary contains the return status and job status.	
Comments	This method is used to get the execution status of a job scheduled by aipu_flush_job. This method is a blocking one which will block in waiting for the job to be done if it is still running.	

Function	aipu_clean_job(job_id)	
Parameter	job_id	Returned value by method aipu_create_job.
Return value	An integer.	
Comments	This method is used to clean a finished job object in UMD. After this API successfully returns, the job ID immediately becomes invalid and cannot be used again. After this API successfully returns, the buffer handle used in creating this job will immediately be free to be used to create another job with the same graph ID.	

Function	aipu_get_tensor_count(graph_id, type)	
Parameter	graph_id	Returned value by method aipu_load_graph.
	type	Tensor type.
Return value	A dictionary contains the return status and tensor count.	
Comments	This method is used to get the tensor count of the specified type	

Function	aipu_get_tensor_descriptor(graph_id, type, tensor_id)	
Parameter	graph_id	Returned value by method aipu_load_graph.
	type	Tensor type.
	tensor_id	The sequence number of the specified tensor type.
Return value	A data structure contains the tensor information.	
Comments	This method is used to get the specified tensor descriptor.	

Function	aipu_load_tensor_from_file(job_id, tensor_id, file_path)	
Parameter	job_id	Returned value by method aipu_create_job.
	tensor_id	The sequence number of the input tensor.
	file_path	Input tensor file path.
Return value	An integer.	
Comments	This method is used to load the input tensor data.	

Function	aipu_load_tensor(job_id, tensor_id, data)	
Parameter	job_id	Returned value by method aipu_create_job.
	tensor_id	The sequence number of the input tensor.
	data	Returned value of file.read() or numpy array. Only supports int8 numpy array.
Return value	An integer.	
Comments	This method is used to load the input tensor data.	

Function	aipu_get_tensor(job_id, type, tensor_id)	
Parameter	job_id	Returned value by method aipu_create_job.
	type	Tensor type.
	tensor_id	The sequence number of the specified tensor.
Return value	A dictionary contains the return status and tensor data.	
Comments	This method is used to get specified tensor data. Only supports int8 tensor type.	

Function	aipu_ioctl(cmd, py_arg)	
Parameter	cmd	<p>AIPU_IOCTL_MARK_SHARED_TENSOR: Mark a shared buffer from a belonged job. The marked buffer is not freed on destroying a job and directly used by a new job.</p> <p>AIPU_IOCTL_SET_SHARED_TENSOR: Specify a marked shared buffer to a new job. Marking a shared buffer and setting a shared buffer to a new job must be performed in the same one process context.</p> <p>AIPU_IOCTL_EN_TICK_COUNTER: Enable the performance counter, no arg.</p> <p>AIPU_IOCTL_DI_COUNTER: Disable the performance counter, no arg.</p> <p>AIPU_IOCTL_SET_PROFILE: Dynamically enable or disable the profiling feature of AIPUv3 simulation.</p> <p>AIPU_IOCTL_ALLOC_DMABUF: Request dma_buf from KMD.</p> <p>AIPU_IOCTL_FREE_DMABUF: Free a dma_buf with its fd.</p>
	py_arg	A dictionary contains the args.
Return value	A dictionary contains the return status and cmd required data.	

Comments	This method is used to send a specific command to the driver. [aipu_v3 only]	
Function	aipu_specify_jobbuf(job_id, py_arg)	
Parameter	job_id	Returned value by method aipu_create_job.
	py_arg	Specify one dma_buf or non-dma_buf as the job's IO buffer.
Return value	An integer.	
Comments	This method is used to specify a shared buffer as the job's IO buffer. The shared buffer can be one dma_buf or non-dma_buf.	

1.2.5 Driver features

Adapt the SMMU

Except for the device specific reserved memory, the NPU driver also allows to use the physically discontinuous memory regions with the mapping of Arm SMMU v3. To enable an SMMU for the NPU, the corresponding configurations of SMMU v3, the stream ID for the NPU, and DMA ranges should be configured in the device tree. For detailed configurations, see Section 3.1.4.2 in `driver/kmd/porting_guide.txt`. If correctly configured, the KMD will allocate the virtually contiguous addresses for user threads, and both the UMD and the NPU can use these addresses in the same way as they handle the reserved physically contiguous ones.

Specify memory regions

If you configure the DDR, SRAM and DTCM (Zhouyi X1 only) memory regions for the NPU in the device tree file (dts), you need to specify memory regions. It is important to note that the address range of SRAM should be located at ASID0. Refer to the KMD porting guide for detailed configuration.

When creating a job, you can choose which region the weight buffer or feature map buffer is allocated from. The three macros `AIPU_MEM_REGION_DEFAULT`, `AIPU_MEM_REGION_SRAM`, and `AIPU_MEM_REGION_DTCM` represent the DDR, SRAM and DTCM regions, respectively. You can set them to field `fm_mem_region` or `wt_mem_region` in struct `aipu_create_job_cfg`. If there is not enough free space in the specified region, the allocation logic will fall back in the order of DTCM->SRAM->DDR. Actually, you only need to set the most expected region macro other than multiple macros. Refer to the `benchmark_test` sample.

Share IO buffer via DMA-BUF

The driver already supports sharing input/output buffer among modules through the Linux DMA-BUF framework. To apply this scheme, it needs to confirm that the input/output buffer is non-shared with other feature map buffers. In other words, when generating the model binary with the `aipugb` offline tool, it is required to specify such parameters `--disable_input_buffer_reuse` or `--disable_output_buffer_reuse` to keep the input or output as an independent one. When it allocates one `dma_buf` successfully, this `dma_buf` can be used as the input or output buffer of a certain job by configuring struct `aipu_shared_tensor_info` {type, tensor_idx, dmabuf_fd, offset_in_dmabuf}. Refer to samples `dmabuf_mmap_test`, `dmabuf_vmap_test`, and `dmabuf_dma_test`.

Share IO buffer via common buffer

This scheme is used to share IO buffer among different jobs which are created by the same graph or different graphs. It can only be applied in one process context. For example, the output buffer of the first job is shared with the input buffer of another job, or several jobs share one common input buffer.

Specific utilization steps:

1. Be sure to mark a certain IO buffer as a shared one via macro (`AIPU_IOCTL_MARK_SHARED_TENSOR`).

2. Replace the IO buffer of another job with the previously marked buffer via macro (AIPU_IOCTL_SET_SHARED_TENSOR). Refer to the sharebuffer_test sample.

Pipeline mode

This mode allows the application to flush many jobs to the KMD in advance. This is implemented by calling `aipu_flush_job` other than `aipu_finish_job`, where the latter API is for a blocking mode. The committed jobs are queued in the KMD. After the hardware resource is released, the queued job can be dispatched immediately to the NPU. It can save some schedule cost between jobs. Refer to the `flush_job_test` sample.

Job reclaiming thread

The UMD supports two ways to reclaim the job status. The default is the reclaiming thread that is identical to the committing thread of the job. In other words, the job is polled right after the job is committed via `aipu_flush_job` in the same thread context. To adapt to more flexible application scenarios, actually it can create a specific thread to reclaim the job. This is implemented by specifying `poll_in_commit_thread` as false in struct `aipu_global_config_hw_t` and creating the reclaiming thread to call `aipu_get_job_status` and execute other related logic.

Job callback function

If you want to apply different handling logic for separate jobs, you can specify a custom callback function for each job. The callback function is assigned through the third parameter of `aipu_flush_job`. It will be invoked after the job status is polled successfully by the internal `aipu_get_job_status`.

About readme

To fully understand the whole NPU driver implementation and utilization, it is strongly recommended that you first read the readme files in the driver package. These files provide detailed information about compiling the KMD and UMD, the device tree configuration and usages of samples.

1.2.6 Linux driver programming model – Standard API

The following code samples are just for quick reference. The detailed samples are located in the sample folder of the UMD driver source code package. Using the pipeline, multi_batch, multi thread, multi process, and benchmark can run on multiple cores (if any).

Code 1 shows the init & de-init API calls that Linux applications should follow before scheduling specific inference tasks and after exiting inference on the NPU.

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int main(int argc, char* argv[])
6.  {
7.      /* ... */
8.      aipu_status_t status = AIPU_STATUS_SUCCESS;
9.      aipu_ctx_handle_t *ctx = nullptr;
10.     const char* status_msg = nullptr;
11.
12.     status = aipu_init_context(&ctx);
13.     if (AIPU_STATUS_SUCCESS != status)
14.     {
15.         aipu_get_error_message(ctx, status, &status_msg);
16.         fprintf(stderr, "Error: %s\n", status_msg);
17.         /* error handling or quitting */
18.     }
19.
20.     /* load graphs and schedule jobs */

```

```

21.
22.     status = aipu_deinit_context(ctx);
23.     if (AIPU_STATUS_SUCCESS != status)
24.     {
25.         aipu_get_error_message(ctx, status, &status_msg);
26.         fprintf(stderr, "Error: %s\n", status_msg);
27.         /* error handling or quitting */
28.     }
29.
30.     /* ... */
31. }

```

Code 1. Initialization and De-Initialization

Code 2 shows the config API call that Linux applications should follow before scheduling specific inference tasks on the NPU simulator.

```

1. /**
2.  * @file main thread
3.  */
4.
5. int main(int argc, char* argv[])
6. {
7.     /* ... */
8.     aipu_global_config_simulation_t sim_glb_config;
9.
10.    /* for aipu_v1/v2 */
11.    /* z2/z3/x1_aipu_simulator optionally */
12.    sim_glb_config.simulator = "./x1_aipu_simulator";
13.
14.    sim_glb_config.log_level = 0;
15.
16.    /* for aipu_v3 */
17.    sim_glb_config.x2_arch_desc = "X2_1204";
18.
19.    /* initialize context */
20.
21.    /* config simulation global options */
22.    ret = aipu_config_global(ctx, AIPU_CONFIG_TYPE_SIMULATION,
23.        &sim_glb_config);
24.    if (AIPU_STATUS_SUCCESS != status)
25.    {
26.        aipu_get_error_message(ctx, status, &status_msg);
27.        fprintf(stderr, "Error: %s\n", status_msg);
28.        /* error handling or quitting */
29.    }
30.
31.    /* load graphs and schedule jobs */
32.
33.    /* deinit */
34.    /* ... */
35. }

```

Code 2. Simulation Configuration

Code 3 shows an example of getting the partition number, cluster number, and core number of the NPU before specifying the job to a specific partition.

```

1. /**
2.  * @file main thread
3.  */

```

```

4.
5. int main(int argc, char* argv[])
6. {
7.     uint32_t part_idx = 0, part_cnt = 0, cluster_cnt = 0, core_cnt = 0;
8.
9.     /* ... */
10.
11.    /* initialize context */
12.
13.    /* config simulation global options */
14.
15.    /* get partition count */
16.    ret = aipu_get_partition_count(ctx, &part_cnt);
17.    if (ret != AIPU_STATUS_SUCCESS)
18.    {
19.        aipu_get_error_message(ctx, ret, &msg);
20.        fprintf(stderr, "aipu_get_partition_count: %s\n", msg);
21.        return ret;
22.    }
23.
24.    for (uint32_t i = 0; i < part_cnt; i++)
25.    {
26.        /* get cluster count of a partition */
27.        ret = aipu_get_cluster_count(ctx, i, &cluster_cnt);
28.        if (ret != AIPU_STATUS_SUCCESS)
29.        {
30.            aipu_get_error_message(ctx, ret, &msg);
31.            fprintf(stderr, "aipu_get_cluster_count: %s\n", msg);
32.            return ret;
33.        }
34.
35.        for (uint32_t j = 0; j < cluster_cnt; j++)
36.        {
37.            /* get core count of a cluster within a partition */
38.            ret = aipu_get_core_count(ctx, i, j, &core_cnt);
39.            if (ret != AIPU_STATUS_SUCCESS)
40.            {
41.                aipu_get_error_message(ctx, ret, &msg);
42.                fprintf(stderr, "aipu_get_core_count: %s\n", msg);
43.                return ret;
44.            }
45.            fprintf(stdout, "<part_idx, cluster_idx, core_cnt> =
46.            <%u, %u, %u>\n", i, j, core_cnt);
47.        }
48.
49.        /* load graphs and schedule jobs */
50.
51.        /* deinit */
52.        /* ... */
53.    }

```

Code 3. Getting partition, cluster, and core information

Code 4 shows the graph loading and unloading API calls that Linux applications should follow before scheduling specific inference tasks.

```

1. /**
2.  * @file main thread
3.  */
4.
5. int main(int argc, char* argv[])

```

```

6. {
7.     /* ... */
8.     uint64 graph_id;
9.     Const char *graph = "./aipu.bin";
10.
11.    /* initialize context */
12.
13.    /* load one or multiple graphs from buffer */
14.    status = aipu_load_graph(ctx, graph, &graph_id);
15.    if (AIPU_STATUS_SUCCESS != status)
16.    {
17.        aipu_get_error_message(ctx, status, &status_msg);
18.        fprintf(stderr, "Error: %s\n", status_msg);
19.        /* error handling or quitting */
20.    }
21.
22.    /* alloc tensor buffer, create and schedule jobs in this single thread */
23.
24.    /* unload one or multiple graphs */
25.    status = aipu_unload_graph(ctx, graph_id);
26.    if (AIPU_STATUS_SUCCESS != status)
27.    {
28.        aipu_get_error_message(ctx, status, &status_msg);
29.        fprintf(stderr, "Error: %s\n", status_msg);
30.        /* error handling or quitting */
31.    }
32.
33.
34.    /* deinit */
35.    /* ... */
36. }

```

Code 4. Graph(s) Loading

Code 5, code 6, and code 7 show the inference job creating, scheduling and execution status checking model that single thread Linux applications should follow.

```

1. /**
2.  * @file main thread
3.  */
4.
5. int non_pipeline_scheduling()
6. {
7.     uint64_t job_id = 0;
8.     aipu_create_job_cfg_t config = {0};
9.
10.    /* load graph and allocate buffers */
11.    /* load input tensors data */
12.
13.    create_job_cfg.fm_mem_region = 0;
14.    create_job_cfg.wt_mem_region = 0;
15.    /* the below 2 lines just for X2, if run for Z2/Z3/X1, it doesn't need to
    set this parameter */
16.    config.partition_id = 0;
17.    config.qos_level = 0;
18.
19.    status = aipu_create_job(ctx, graph_id, &job_id, &config);
20.    if (AIPU_STATUS_SUCCESS != status)
21.    {
22.        aipu_get_status_msg(ctx, status, &status_msg);
23.        fprintf(stderr, "Error: %s\n", status_msg);
24.        /* error handling or quitting */

```

```

25.     }
26.     status = aipu_finish_job(ctx, job_id, time_out);
27.     if (AIPU_STATUS_SUCCESS != status)
28.     {
29.         aipu_get_error_message(ctx, status, &status_msg);
30.         fprintf(stderr, "Error: %s\n", status_msg);
31.         /* error handling or quitting */
32.     }
33.     else
34.     {
35.         /* job results fetching */
36.     }
37.
38.     /* ... */
39.
40.     status = aipu_clean_job(ctx, job_id);
41.     if (AIPU_STATUS_SUCCESS != status)
42.     {
43.         aipu_get_error_message(ctx, status, &status_msg);
44.         fprintf(stderr, "Error: %s\n", status_msg);
45.         /* error handling or quitting */
46.     }
47. }
48.
49. int main(int argc, char* argv[])
50. {
51.     /* ... */
52.     /* initialize */
53.
54.     ret = non_pipeline_scheduling();
55.     /* ... */
56.
57.     /* deinit */
58.     /* ... */
59. }

```

Code 5. Non-pipeline Scheduling in Single Thread Application

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int pipeline_scheduling()
6.  {
7.      aipu_create_job_cfg_t config = {0};
8.
9.      /* ... */
10.     int pipe_cnt = 3;
11.     Uint64_t job_id[pipe_cnt];
12.     int time_out = 1000;
13.     aipu_job_status_t job_status = AIPU_JOB_STATUS_NO_STATUS;
14.
15.     /* load graph and allocate buffers */
16.     /* load input tensors data */
17.
18.     /* create multiple jobs and then schedule one by one */
19.     for (uint32_t i = 0; i < pipe_cnt; i++)
20.     {
21.         /* load inputs */
22.
23.         create_job_cfg.fm_mem_region = 0;
24.         create_job_cfg.wt_mem_region = 0;

```

```

25.      /* the below 2 lines just for X2 running, if run for Z2/Z3/X1, it
doesn't need to set this parameter */
26.      config.partition_id = 0;
27.      config.qos_level = 0;
28.
29.      /* create job with corresponding buffer handle */
30.      ret = aipu_create_job(ctx, graph_id, &job_id[i], &config);
31.      if (ret != AIPU_STATUS_SUCCESS)
32.      {
33.          aipu_get_error_message(ctx, ret, &status_msg);
34.          fprintf(stderr, "Error: %s\n", status_msg);
35.          /* error handle or quit */
36.      }
37.
38.      ret = aipu_flush_job(ctx, job_id[i]);
39.      if (ret != AIPU_STATUS_SUCCESS)
40.      {
41.          aipu_get_error_message(ctx, ret, &status_msg);
42.          fprintf(stderr, "Error: %s\n", status_msg);
43.          /* error handle or quit */
44.      }
45.  }
46.  /* application can do some other tasks if any before getting job results */
47.  /* get job results */
48.  for (uint32_t i = 0; i < pipe_cnt; i++)
49.  {
50.      ret = aipu_get_job_status(ctx, job_id[i], &job_status);
51.      if ((ret != AIPU_STATUS_SUCCESS) || (job_status != AIPU_JOB_STATUS_DONE
))
52.      {
53.          aipu_get_error_message(ctx, ret, &status_msg);
54.          fprintf(stderr, "Error: %s\n", status_msg);
55.          /* error handle or quit */
56.      }
57.
58.      /* fetch job #i result */
59.
60.      ret = aipu_clean_job(ctx, job_id[i]);
61.      if (ret != AIPU_STATUS_SUCCESS)
62.      {
63.          aipu_get_error_message(ctx, ret, &status_msg);
64.          fprintf(stderr, "Error: %s\n", status_msg);
65.          /* error handle or quit */
66.      }
67.  }
68. }
69.
70. int main(int argc, char* argv[])
71. {
72.     /* ... */
73.     /* initialize */
74.
75.     ret = pipeline_scheduling();
76.     /* ... */
77.
78.     /* free buffers, unload and deinit */
79.     /* ... */
80. }

```

Code 6. Pipeline Scheduling in Single Thread Application

```

1.  /**

```

```

2.  * @file main thread
3.  */
4.  #define MAX_BATCH 4
5.
6.  int multi_batch_scheduling()
7.  {
8.      /* ... */
9.      uint32_t input_cnt, output_cnt;
10.     char **input_buf = nullptr, **output_buf[MAX_BATCH];
11.     uint32_t queue_id = 0;
12.     uint32_t batch_loop_cnt = 2;
13.     uint64_t graph_id;
14.     const char* msg = nullptr;
15.     vector<aipu_tensor_desc_t> input_desc;
16.     vector<aipu_tensor_desc_t> output_desc;
17.     aipu_status_t ret = AIPU_STATUS_SUCCESS;
18.
19.     /* load graph and allocate buffers */
20.     /* get input and output count */
21.
22.     /* found input_buf and output_buf */
23.     input_buf = new char* [input_cnt];
24.     for (uint32_t i = 0; i < MAX_BATCH; i++)
25.         output_buf[i] = new char* [output_cnt];
26.
27.     for (uint32_t i = 0; i < input_cnt; i++)
28.         input_buf[i] = opt.inputs[i];
29.
30.     for (uint32_t k = 0; k < MAX_BATCH; k++)
31.     {
32.         for (uint32_t i = 0; i < output_cnt; i++)
33.         {
34.             char* output = new char[output_desc[i].size];
35.             output_data[k].push_back(output);
36.             output_buf[k][i] = output;
37.         }
38.     }
39.
40.     /* create multi batch queue */
41.     ret = aipu_create_batch_queue(ctx, graph_id, &queue_id);
42.     if (ret != AIPU_STATUS_SUCCESS)
43.     {
44.         aipu_get_error_message(ctx, ret, &msg);
45.         fprintf(stderr, "aipu_create_batch_queue: %s\n", msg);
46.         /* error handle or quit */
47.     }
48.
49.     ret = aipu_config_batch_dump(ctx, graph_id, queue_id,
50.         AIPU_CONFIG_TYPE_SIMULATION, &mem_dump_config);
51.     if (ret != AIPU_STATUS_SUCCESS)
52.     {
53.         aipu_get_error_message(ctx, ret, &msg);
54.         fprintf(stderr, "aipu_config_batch_dump: %s\n", msg);
55.         /* error handle or quit */
56.     }
57.
58.     for (uint32_t round = 0; round < batch_loop_cnt; round++)
59.     {
60.         AIPU_INFO() << "Batch round #" << round;
61.         for (uint32_t i = 0; i < MAX_BATCH; i++)
62.         {
63.             /* add batch */

```



```

64.         ret = aipu_add_batch(ctx, graph_id, queue_id, input_buf,
65.                             output_buf[i]);
66.         if (ret != AIPU_STATUS_SUCCESS)
67.         {
68.             aipu_get_error_message(ctx, ret, &msg);
69.             fprintf(stderr, "aipu_add_batch: %s\n", msg);
70.             /* error handle or quit */
71.         }
72.     }
73.
74.     ret = aipu_finish_batch(ctx, graph_id, queue_id, &create_job_cfg);
75.     if (ret != AIPU_STATUS_SUCCESS)
76.     {
77.         aipu_get_error_message(ctx, ret, &msg);
78.         fprintf(stderr, "aipu_finish_batch: %s\n", msg);
79.         /* error handle or quit */
80.     }
81.     /* fetch batch #i result */
82. }
83.
84. ret = aipu_clean_batch_queue(ctx, graph_id, queue_id);
85. if (ret != AIPU_STATUS_SUCCESS)
86. {
87.     aipu_get_error_message(ctx, ret, &msg);
88.     fprintf(stderr, "aipu_clean_batch_queue: %s\n", msg);
89.     /* error handle or quit */
90. }
91. }
92.
93. int main(int argc, char* argv[])
94. {
95.     /* ... */
96.     /* initialize */
97.
98.     ret = multi_batch_scheduling ();
99.     /* free buffers, unload and deinit */
100. }

```

Code 7. Multi-batch Scheduling in Single Thread Application

Code 8 shows the pipeline scheduling in multithread application.

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int main(int argc, char* argv[])
6.  {
7.
8.      /* initialize context */
9.      void *data1 = nullptr, *data2 = nullptr;
10.
11.     /* initialize data1 & data2 */
12.
13.     /* create more scheduling threads to load graphs independently and run jobs
14.      */
15.     std::thread t1(job_scheduling_thread, data1);
16.     std::thread t2(job_scheduling_thread, data2);
17.
18.     t1.join();
19.     t2.join();

```

```

20.  /* deinit */
21.  /* ... */
22. }
23.
24. /**
25.  * @file scheduling thread
26.  */
27.
28. void* job_scheduling_thread(void* data)
29. {
30.  /* ... */
31.  ret = pipeline_scheduling(); /* the same as the scheduling func in Code 6
32.  */
33.  /* ... */
33. }

```

Code 8. Pipeline Scheduling in Multithread Application

Code 9 shows the usage of the NPU printf and profiling dump features in a single thread non-pipeline application. In real applications, you can decide whether to use these features. Note that the profiling data generated by the NPU profiler will be ready in the profiling buffer automatically after the profiling job ends. The printf data is also stored to the corresponding buffer if the printf API in the NPU binary is called. Note that currently the printf feature is not supported on aipu_v3.

```

1.  /**
2.  * @file main thread
3.  */
4.
5.  int dump_file_helper(const char* fname, const void* src, unsigned int size)
6.  {
7.  /* write <size> bytes from pointer <src> into a file named <fname> */
8.  }
9.
10. int main(int argc, char* argv[])
11. {
12.  /* ... */
13.  char fname[4096];
14.  aipu_status_t status = AIPU_STATUS_SUCCESS;
15.  const char* status_msg = nullptr;
16.  aipu_tensor_desc_t desc;
17.  char *buff = nullptr;
18.
19.  /* init context and load graphs */
20.  /* ... */
21.
22.  /* buffer allocations */
23.  status = aipu_get_tensor_descriptor(ctx, graph_id, AIPU_TENSOR_TYPE_PRINTF,
24.  0, &desc);
25.  if (AIPU_STATUS_SUCCESS != status)
26.  {
27.      aipu_get_error_message(ctx, status, &status_msg);
28.      fprintf(stderr, "Alloc Error: %s\n", status_msg);
29.      /* error handling or quitting */
30.  }
31.  buff = new char[desc.size];
32.
33.  /* create, schedule a job and wait for it ends */
34.  /* ... */
35.
36.  /* Job Done Now */
37.  /* Case #9.1: print AIPU log information to terminal */
38.  status = aipu_get_tensor(m_ctx, job_id, AIPU_TENSOR_TYPE_PRINTF, 0, buff);

```

```

38.     status = aipu_printf(buff, "print.log");
39.     if (AIPU_STATUS_SUCCESS != status)
40.     {
41.         aipu_get_error_message(ctx, status, &status_msg);
42.         fprintf(stderr, "Printf Error: %s\n", status_msg);
43.         /* error handling or quitting */
44.     }
45.
46.     /* Case #9.2: get data in profiling-dump buffer(s) */
47.
48.     /* buffer allocations */
49.     status = sts = aipu_get_tensor_descriptor(ctx, graph_id,
AIPU_TENSOR_TYPE_PROFILER, 0, &desc);
50.     if (AIPU_STATUS_SUCCESS != status)
51.     {
52.         aipu_get_error_message(ctx, status, &status_msg);
53.         fprintf(stderr, "Alloc Error: %s\n", status_msg);
54.         /* error handling or quitting */
55.     }
56.     buff = new char[desc.size];
57.
58.     status = aipu_get_tensor(m_ctx, job_id, AIPU_TENSOR_TYPE_PROFILER, 0,
buff);
59.     if (AIPU_STATUS_SUCCESS != status)
60.     {
61.         aipu_get_error_message(ctx, status, &status_msg);
62.         fprintf(stderr, "Printf Error: %s\n", status_msg);
63.         /* error handling or quitting */
64.     }
65.
66.     /* dump the profiling data generated by profiler as a file */
67.     /* or do any other operations application would like to */
68.     snprintf(fname, sizeof(fname), "AIPU_Profiler_Data.bin");
69.     dump_file_helper(fname, buff, desc.size);
70.
71.
72.     /* clean jobs and graphs, and deinit */
73.     /* ... */
74. }

```

Code 9. Usage of NPU Printf and Profiling Dump Features

Code 10 shows the usage of NPU debugger APIs.

```

1.  /**
2.   * @file main thread
3.   */
4.
5.
6.  int main(int argc, char* argv[])
7.  {
8.      /* ... */
9.      aipu_status_t status = AIPU_STATUS_SUCCESS;
10.     const char* status_msg = nullptr;
11.     uint64_t job_id = 0;
12.     uint32_t core_cnt = 0;
13.     aipu_core_info_t* core_info = nullptr;
14.     aipu_debugger_job_info_t dbg_job_info;
15.
16.     /* init context, load graphs, allocate buffers */
17.     /* ... */
18.

```

```

19.  /* create an AIPU debugger job */
20.  status = aipu_create_job(ctx, graph_id, &job_id);
21.  if (AIPU_STATUS_SUCCESS != status)
22.  {
23.      aipu_get_error_message(ctx, status, &status_msg);
24.      fprintf(stderr, "Create Job Error: %s\n", status_msg);
25.      /* error handling or quitting */
26.  }
27.  /* The following 5 APIs should be called after aipu_create_job and
    before AIPU_clean_job */
28.
29.  status = aipu_get_core_cnt(ctx, 0, 0, &core_cnt);
30.  if (AIPU_STATUS_SUCCESS != status)
31.  {
32.      aipu_get_error_message(ctx, status, &status_msg);
33.      fprintf(stderr, "Error: %s\n", status_msg);
34.      /* error handling or quitting */
35.  }
36.  core_info = new aipu_core_info_t[core_cnt];
37.  for (uint32_t i = 0; i < core_cnt; i++)
38.  {
39.      status = aipu_debugger_get_core_info(ctx, i, &core_info[i]);
40.      if (ret != AIPU_STATUS_SUCCESS)
41.      {
42.          aipu_get_error_message(ctx, status, &status_msg);
43.          fprintf(stderr, "Error: %s\n", status_msg);
44.          /* error handling or quitting */
45.      }
46.  }
47.
48.  status = aipu_debugger_get_job_info(ctx, job_id, &dbg_job_info);
49.  if (AIPU_STATUS_SUCCESS != status)
50.  {
51.      aipu_get_error_message(ctx, status, &status_msg);
52.      fprintf(stderr, "Error: %s\n", status_msg);
53.      /* error handling or quitting */
54.  }
55.
56.  /* bind this job to AIPU core 0 (or another one if applicable) */
57.  status = aipu_debugger_bind_job(ctx, core_selected, job_id);
58.  if (AIPU_STATUS_SUCCESS != status)
59.  {
60.      aipu_get_error_message(ctx, status, &status_msg);
61.      fprintf(stderr, "Error: %s\n", status_msg);
62.      /* error handling or quitting */
63.  }
64.
65.  /* trigger this job to run on the bind core */
66.  status = aipu_debugger_run_job(ctx, job_id);
67.  if (AIPU_STATUS_SUCCESS != status)
68.  {
69.      aipu_get_error_message(ctx, status, &status_msg);
70.      fprintf(stderr, "Error: %s\n", status_msg);
71.      /* error handling or quitting */
72.  }
73.
74.  /* get results and clean all */
75. }

```

Code 10. Usage of NPU Debugger Features

Code 11 shows the usage of accessing dma_buf via mmap in user mode.

Note that the `dma_buf` mechanism is mainly for sharing IO buffers across multiple modules or processes. Currently it is only implemented for `aipuv1/v2/v3`. Refer to the detailed implementation related to the `mmap` way in the sample `samples/src/dmabuf_mmap_test`. In addition, if you want to use `dma_buf` in kernel mode via `vmap` or DMA APIs, refer to the samples in `samples/src/{dmabuf_vmap_test, dmabuf_dma_test}`.

```
1. #define DEV_EXPORTER "/dev/aipu"
2. #define DMABUF_SZ 0x100000
3.
4. int dmabuf_malloc(uint64_t size)
5. {
6.     int ret = 0;
7.     int fd = 0, dmabuf_fd = 0;;
8.     struct aipu_dma_buf_request dma_buf_req = {0};
9.
10.    fd = open(DEV_EXPORTER, O_RDWR);
11.    if (fd < 0)
12.    {
13.        ret = -1;
14.        fprintf(stderr, " open %s [fail]\n", DEV_EXPORTER);
15.        goto out;
16.    }
17.
18.    dma_buf_req.bytes = size;
19.    ret = ioctl(fd, AIPU_IOCTL_ALLOC_DMA_BUF, &dma_buf_req);
20.    if (ret < 0)
21.    {
22.        fprintf(stderr, "ioctl %s [fail]\n", DEV_EXPORTER);
23.        goto out;
24.    }
25.
26.    dmabuf_fd = dma_buf_req.fd;
27.
28.out:
29.    close(fd);
30.    return dmabuf_fd;
31. }
32.
33. int dmabuf_free(int _fd)
34. {
35.     int ret = 0;
36.     int fd = 0;
37.
38.     fd = open(DEV_EXPORTER, O_RDWR);
39.     if (fd < 0)
40.     {
41.         ret = -1;
42.         fprintf(stderr, "open %s [fail]\n", DEV_EXPORTER);
43.         goto out;
44.     }
45.
46.     ret = ioctl(fd, AIPU_IOCTL_FREE_DMA_BUF, &_fd);
47.     if (ret < 0)
48.     {
49.         fprintf(stderr, "ioctl %s [fail]\n", DEV_EXPORTER);
50.         goto out;
51.     }
52.
53.out:
54.    close(fd);
55.    return ret;
56. }
```

```

57.
58. int dmabuf_fill(int fd, char *data, uint32_t size)
59. {
60.     int ret = 0;
61.     char *va = nullptr;
62.
63.     va = (char *)mmap(NULL, DMABUF_SZ, PROT_READ|PROT_WRITE, MAP_SHARED, fd,
64. 0);
65.     if (va == MAP_FAILED)
66.     {
67.         ret = -1;
68.         fprintf(stderr, "mmap dmabuf [fail]\n");
69.         goto out;
70.     }
71.     memcpy(va, data, size);
72.     munmap(va, DMABUF_SZ);
73.
74. out:
75.     return ret;
76. }
77.
78. int main(int argc, char* argv[])
79. {
80.     /* ... */
81.     aipu_status_t status = AIPU_STATUS_SUCCESS;
82.     const char* status_msg = nullptr;
83.     uint64_t graph_id, job_id;
84.     aipu_shared_tensor_info_t share_tensor;
85.     dmabuf_fd = 0;
86.
87.     /* init context, load graphs, allocate buffers */
88.     /* ... */
89.
90.     /* request dma_buf */
91.     dmabuf_fd = dmabuf_malloc(DMABUF_SZ);
92.
93.     /* fill dma_buf with input data */
94.     dmabuf_fill(dmabuf_fd, data, size);
95.
96.     /*
97.      * construct share buffer arguments
98.      * @dmabuf_fd: the fd of dma_buf
99.      * @offset_in_dmabuf: the start offset of valid data in dma_buf
100.      * @tensor_idx: the tensor index which the dma_buf will replace
101.      * @type: the replaced tensor type(input or output)
102.      */
103.     share_tensor.dmabuf_fd = dmabuf_fd;
104.     share_tensor.offset_in_dmabuf = 0;
105.     share_tensor.tensor_idx = 0;
106.     share_tensor.type = AIPU_TENSOR_TYPE_INPUT;
107.
108.     status = aipu_create_job(ctx, graph_id, &job_id);
109.     if (status != AIPU_STATUS_SUCCESS)
110.     {
111.         aipu_get_error_message(ctx, status, &status_msg);
112.         fprintf(stderr, "Error: %s\n", status_msg);
113.     }
114.
115.     /* specify dma_buf as input tensor buffer-0 */
116.     status = aipu_specify_iobuf(ctx, job_id, &share_tensor);
117.     if (status != AIPU_STATUS_SUCCESS)

```

```

118.     {
119.         aipu_get_error_message(ctx, status, &status_msg);
120.         fprintf(stderr, "Error: %s\n", status_msg);
121.     }
122.
123.     /* aipu_finish_job ... */
124.
125.     /* aipu_clean_job, aipu_deinit_ctx ... */
126. }

```

Code 11. Usage of dma_buf via the mmap method in user mode

Code 12 shows the usage of aipu_ioctl to handle some miscellaneous requests with aipu_v3.

The aipu_ioctl API supports many miscellaneous controls for aipu_v3, for example, request or free dma_buf, enable or disable some idle cores in a cluster, and enable or disable the performance counter. In addition, it also can be used to share some IO buffers in the same process context.

```

1. int main(int argc, char* argv[])
2. {
3.     /* ... */
4.     aipu_status_t status = AIPU_STATUS_SUCCESS;
5.     const char* status_msg = nullptr;
6.     aipu_ctx_handle_t ctx;
7.
8.     status = aipu_init_context(&ctx);
9.
10.    /* load graphs, allocate buffers */
11.    /* ... */
12.
13.    /* demo1: enable performance tick counter */
14.    status = aipu_ioctl(ctx, AIPU_IOCTL_ENABLE_TICK_COUNTER, nullptr);
15.
16.    /* demo2: enable performance tick counter
17.     *
18.     * Just enable one core in cluster 0
19.     */
20.    struct aipu_config_clusters config_clusters;
21.    memset(&config_clusters, 0, sizeof(config_clusters));
22.    config_clusters.clusters[0].en_core_cnt = 1;
23.    status = aipu_ioctl(ctx, AIPU_IOCTL_CONFIG_CLUSTERS, &config_clusters);
24.
25.    /* demo3: alloc/free dma_buf */
26.    struct aipu_dma_buf_request dma_buf_req = {0};
27.    dma_buf_req.bytes = 4096; /* request 4KB */
28.    status = aipu_ioctl(ctx, AIPU_IOCTL_ALLOC_DMABUF, &dma_buf_req);
29.
30.    /* free dma_buf through its fd */
31.    status = aipu_ioctl(ctx, AIPU_IOCTL_FREE_DMABUF, dma_buf_req.fd);
32.
33.    if (status != AIPU_STATUS_SUCCESS)
34.    {
35.        aipu_get_error_message(ctx, status, &status_msg);
36.        fprintf(stderr, "Error: %s\n", status_msg);
37.    }
38.
39.    /* aipu_finish_job ... */
40.    /* get result ... */
41.
42.    /* aipu_clean_job, ... */
43.    status = aipu_init_context(ctx);
44.    if (status != AIPU_STATUS_SUCCESS)

```

```

45.     {
46.         aipu_get_error_message(ctx, status, &status_msg);
47.         fprintf(stderr, "Error: %s\n", status_msg);
48.     }
49. }

```

Code 12. Usage of dma_buf via the mmap method in user mode

1.2.7 Linux driver programming model – Python APIs

Code 13 shows the usage of the NPU Linux driver Python APIs.

```

1.
2. from libaipudrv import *
3.
4. # create a npu object and init UMD context.
5. npu = NPU()
6. ret = npu.aipu_init_context()
7. if ret != 0:
8.     print("failed")
9.     exit(-1)
10.
11. # load an AIPU executable binary file
12. retmap = npu.aipu_load_graph("path/aipu.bin")
13. if retmap["ret"] == 0:
14.     graph_id = retmap["data"]
15. else:
16.     exit(-1)
17.
18. # get input tensor count
19. retmap = npu.aipu_get_tensor_count(graph_id, AIPU_TENSOR_TYPE_INPUT)
20. if retmap["ret"] == 0:
21.     input_cnt = retmap["data"]
22. else:
23.     exit(-1)
24.
25. # get output tensor count
26. retmap = npu.aipu_get_tensor_count(graph_id, AIPU_TENSOR_TYPE_OUTPUT)
27. if retmap["ret"] == 0:
28.     output_cnt = retmap["data"]
29. else:
30.     exit(-1)
31.
32. # create an AIPU job
33. fm_idxes = []
34. wt_idxes = []
35. job_cfg = {"partition_id":0, "dbg_dispatch":0,
36.            "dbg_core_id":0, "qos_level":0}
37. retmap = npu.aipu_create_job(graph_id, job_cfg, fm_idxes, wt_idxes)
38. if retmap["ret"] == 0:
39.     job_id = retmap["data"]
40. else:
41.     exit(-1)
42.
43. # load input tensor
44. npu.aipu_load_tensor_from_file(job_id, 0, "path/input0.bin")
45.
46. # flush job to AIPU HW and waiting, -1 indicate wait forever
47. ret = npu.aipu_finish_job(job_id, -1)
48. if ret != 0:
49.     exit(-1)
50.

```



```
51. # get output tensor data
52. retmap = npu.aipu_get_tensor(job_id, AIPU_TENSOR_TYPE_OUTPUT, 0)
53. if retmap["ret"][0] == 0:
54.     output_data = retmap["data"]
55. else:
56.     exit(-1)
57.
58. # clean job and unload graph
59. ret = npu.aipu_clean_job(job_id)
60. ret = npu.aipu_unload_graph(graph_id)
61.
62. # destroy UMD context
63. ret = npu.aipu_deinit_context()
64.
```

Code 13. Usage of NPU Driver Python APIs

1.3 QNX-based driver (for customized solutions only)

The QNX-based driver consists of two key parts—*User Mode Driver* (UMD) and *Resource Manager Driver* (RMD). The UMD parses NPU job descriptors passed by applications, allocates resources and schedules jobs by calling RMD interfaces.

1.3.1 UMD

The UMD is compiled as a dynamic library or static library and provides user applications a series of interfaces. Standard UMD APIs are designed for applications to schedule standard AI executable benchmarks onto the NPU.

User applications should call the APIs in a relatively fixed order and check the return values step by step to ensure correct operations.

The typical usage of the standard UMD interfaces in an application is as follows:

1. The application calls the context initialization API first to initialize UMD runtime context.
2. After initialization, the application should call the graph loading API to load an offline built binary. If successful, a graph ID is returned. Multiple graph binaries can be loaded in order.
3. After graph loading, the application can create and schedule an inference job which is bound to a previously loaded graph. All input, output, and intermediate buffers are allocated in this step.
4. The application should load the input frame data into the corresponding input buffers.
5. If the blocking scheduling API is used and returns without an error code, the application can get the output tensor directly from the tensor output buffer. If the non-blocking scheduling API is used, after querying the status of execution and no exception is found, the application can fetch the results in the same way.
6. The application should clean a finished job after it terminates and then the bind buffer can be reused in loading the next frame's inputs and rescheduling the job again as described in steps 4–5.
7. Before exiting, the application should unload all graphs and destroy the context.

For more information about the corresponding data structures and interfaces, see 5.3.3 *QNX driver user interface - Standard API*.

1.3.2 Resource Manager Driver

The Resource Manager Driver is a server program that receives messages from applications. It is running in the user space and communicates with the NPU hardware. When the driver is running, it will generate a device file for each device such as the Linux char devices.

The Resource Manager Driver works as a user space application in QNX to respond to contiguous physical memory allocation requests from the UMD, controls job scheduling on the NPU, and handles NPU interrupts.

To support UMD requests, like other Linux char drivers, several standard file operation interfaces are provided to the UMD—open, close, mmap_device_memory, devctl and poll.

The devctl options supported are as follows:

- AIPU_IOCTL_QUERY_CAP: Queries the common capability of NPUs.
- AIPU_IOCTL_QUERY_PARTITION_CAP: Queries the capability of an NPU core.

- AIPU_IOCTL_REQ_BUF: Requests to allocate a coherent buffer.
- AIPU_IOCTL_FREE_BUF: Requests to free a coherent buffer allocated by AIPU_IOCTL_REQBUF.
- AIPU_IOCTL_SCHEDULE_JOB: Schedules a user job to the kernel mode driver for execution.
- AIPU_IOCTL_QUERY_STATUS: Queries the execution status of one or multiple scheduled jobs.
- AIPU_IOCTL_REQ_IO: Reads/Writes an external register of an NPU core.

1.3.3 QNX driver user interface - Standard API

Data structure

Data structure	Type	Values	Comment
aipu_data_type_t	enum	TENSOR_DATA_TYPE_NONE (0x0)	No data type.
		AIPU_DATA_TYPE_BOOL(0x1)	-
		TENSOR_DATA_TYPE_U8 (0x2)	Unsigned int8.
		TENSOR_DATA_TYPE_S8 (0x3)	Signed int8.
		TENSOR_DATA_TYPE_U16 (0x4)	Unsigned int16.
		TENSOR_DATA_TYPE_S16 (0x5)	Signed int16.
		TENSOR_DATA_TYPE_U32 (0x6)	Unsigned int32.
		TENSOR_DATA_TYPE_S32 (0x7)	Signed int32.
		TENSOR_DATA_TYPE_U64 (0x8)	Unsigned int64.
		TENSOR_DATA_TYPE_S64 (0x9)	Signed int64.
		TENSOR_DATA_TYPE_f16 (0xa)	Float 16.
		TENSOR_DATA_TYPE_f32 (0xb)	Float 32.
		TENSOR_DATA_TYPE_f64 (0xc)	Float 64.

Data structure	Type	Members	Comment
aipu_tensor_type_t	enum	AIPU_TENSOR_TYPE_INPUT	-
		AIPU_TENSOR_TYPE_OUTPUT	-
		AIPU_TENSOR_TYPE_PRINTF	-
		AIPU_TENSOR_TYPE_PROFILER	-
		AIPU_TENSOR_TYPE_INTER_DUMP	-
		AIPU_TENSOR_TYPE_LAYER_COUNTER	-
		AIPU_TENSOR_TYPE_ERROR_CODE	-

Data structure	Type	Members	Member type	Comment
aipu_tensor_desc_t	struct	id	uint32_t	Tensor ID.
		size	Uint32_t	Tensor size.
		scale	float	Tensor scale parameter.
		zero_point	float	Tensor zero-point parameter.

		data_type	aipu_data_type_t	Tensor format information.
--	--	-----------	------------------	----------------------------

Data structure	Type	Members	Member type	Comment
aipu_job_config_dump_t	struct	dump_dir	const char*	The data path to store dumped data files.
		prefix	const char*	The file name prefix for dumped files.
		output_prefix	const char*	The name prefix of output dump files.

Data structure	Type	Members	Member type	Comment
aipu_create_job_cfg_t	struct	misc	unsigned int	-
		partition_id	unsigned char	The default value is 0, that is, in partition-0.
		qos_level	unsigned char	The default value is 0, that is, low priority.

Data structure	Type	Values	Comment
aipu_job_status_t	enum	AIPU_JOB_STATUS_NO_STATUS (0x0)	-
		AIPU_JOB_STATUS_DONE (0x1)	-
		AIPU_JOB_STATUS_EXCEPTION (0x2)	-

Data structure	Type	Members	Member type	Comment
aipu_core_info_t	struct	reg_base	uint64_t	Base address of the NPU core external register.

Data structure	Type	Values	Comment
aipu_config_type_t	enum	AIPU_JOB_CONFIG_TYPE_DUMP_TEXT	Dump text section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_WEIGHT	Dump weight section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_RODATA	Dump rodata section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_DESCRIPTOR	Dump descriptor section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_INPUT	Dump input data section(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_OUTPUT	Dump output data tensor(s).
		AIPU_JOB_CONFIG_TYPE_DUMP_TCB_CHAIN	Dump TCB chain tensor(s). Not used for aipu v1/v2.
		AIPU_JOB_CONFIG_TYPE_DUMP_EMULATION	Dump emulation files. Internally used.
		AIPU_CONFIG_TYPE_SIMULATION	QNX does not support simulation.
		AIPU_GLOBAL_CONFIG_TYPE_DISABLE_VER_CHECK	Disable graph bin version check.
		AIPU_GLOBAL_CONFIG_TYPE_ENABLE_VER_CHECK	Enable graph bin version check.

Data structure	Type	Values	Comment
aipu_status_t	enum	AIPU_STATUS_SUCCESS (0x0)	The execution of this API returns successfully without any error.

		AIPU_STATUS_ERROR_NULL_PTR	The arguments passed by user applications to UMD API are NULL pointers which are invalid.
		AIPU_STATUS_ERROR_INVALID_CTX	The context pointer passed by user applications to UMD API is invalid.
		AIPU_STATUS_ERROR_OPEN_DEV_FAIL	UMD fails in opening the device file <code>/dev/aipu</code> which may result from the failure in probing the KMD module.
		AIPU_STATUS_ERROR_DEV_ABNORMAL	The NPU device is in abnormal state.
		AIPU_STATUS_ERROR_DEINIT_FAIL	UMD API fails in de-initiating a context.
		AIPU_STATUS_ERROR_INVALID_CONFIG	The configuration passed by a user application to UMD API is invalid.
		AIPU_STATUS_ERROR_GVERSION_UNSUPPORTED	The version of an executable graph binary passed by a user application is not supported on the current UMD.
		AIPU_STATUS_ERROR_TARGET_NOT_FOUND	UMD fails in finding a matching NPU hardware target for the executable graph passed by a user application to execute.
		AIPU_STATUS_ERROR_INVALID_GBIN	The executable graph binary file passed by a user application to UMD API contains invalid items which cannot be parsed or executed.
		AIPU_STATUS_ERROR_GRAPH_NOT_EXIST	The graph description pointer passed by a user application to UMD API cannot be found in the current context.
		AIPU_STATUS_ERROR_OPEN_FILE_FAIL	UMD fails in opening a file passed by a user application.
		AIPU_STATUS_ERROR_MAP_FILE_FAIL	UMD fails in mapping a file passed by a user application.
		AIPU_STATUS_ERROR_READ_FILE_FAIL	UMD fails in reading a file passed by a user application.
		AIPU_STATUS_ERROR_WRITE_FILE_FAIL	UMD fails in writing a file passed by a user application.
		AIPU_STATUS_ERROR_JOB_NOT_EXIST	The job ID passed by a user application to UMD API cannot be found in the current context.
		AIPU_STATUS_ERROR_JOB_NOT_SCHEDULED	The operation regarding to the job ID passed by a user application to UMD API cannot be executed because it should be done after the job is scheduled.
		AIPU_STATUS_ERROR_JOB_SCHEDULED	The operation regarding to the job ID passed by a user application to UMD API cannot be executed because it should be done before the job is scheduled.
		AIPU_STATUS_ERROR_JOB_NOT_END	The operation regarding to the job ID passed by a user application to UMD API cannot be executed because it should be done after the job execution terminates.
		AIPU_STATUS_ERROR_JOB_EXCEPTION	The execution of a job ID passed by a user application to UMD API ends with an exception.

	AIPU_STATUS_ERROR_JOB_TIMEOUT	The execution of a job ID passed by a user application to UMD API timed out.
	AIPU_STATUS_ERROR_INVALID_OPTIONS	The options passed by a user application to UMD are invalid and cannot be reflected.
	AIPU_STATUS_ERROR_INVALID_PATH	The path passed by a user application to UMD is invalid and cannot be operated.
	AIPU_STATUS_ERROR_OP_NOT_SUPPORTED	The operation is not supported in the current UMD version or system environment.
	AIPU_STATUS_ERROR_INVALID_OP	The operation is invalid and not allowed in the current system environment.
	AIPU_STATUS_ERROR_INVALID_SIZE	The size argument is invalid.
	AIPU_STATUS_ERROR_INVALID_HANDLE	The buffer handle argument is invalid which cannot be found in the current context.
	AIPU_STATUS_ERROR_BUSY_HANDLE	The buffer handle argument cannot be used in this API because it is in busy state and is used by another job.
	AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	The buffer allocation requests fail because of system memory limitations.
	AIPU_STATUS_ERROR_BUF_FREE_FAIL	The buffer free requests fail because the buffers are in busy state.
	AIPU_STATUS_ERROR_RESERVE_SRAM_FAIL	UMD fails in reserving SRAM as the executable binary requested because there is no SoC SRAM or SRAM is busy.
	AIPU_STATUS_ERROR_NO_BATCH_QUEUE	Transfer an incorrect batch queue ID.
	AIPU_STATUS_MAX	Maximum error code value.

Data structure	Type	Members	Member type	Comment
aipu_ctx_handle_t	struct	opaque	-	The pointer to this opaque context is returned by the init API.

User application interfaces

Function declaration	aipu_status_t aipu_init_ctx(aipu_ctx_handle_t** ctx);	
Parameter	ctx	Pointer to a memory location allocated by the application where the UMD stores the opaque context handle struct.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_OPEN_DEV_FAIL AIPU_STATUS_ERROR_DEV_ABNORMAL	
Comments	This API is used to initialize the NPU UMD context.	

Function declaration	aipu_status_t aipu_get_error_message(const aipu_ctx_handle_t* ctx, aipu_status_t status, const char** msg);		
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.	
	status	Status returned by UMD standard API.	

	msg	Pointer to a memory location allocated by the application where the UMD stores the message string pointer.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX	
Comments	This API is used to query additional information about a status returned by UMD API.	

Function declaration	aipu_status_t aipu_deinit_ctx(const aipu_ctx_handle_t* ctx);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_DEINIT_FAIL	
Comments	This API is used to destroy the NPU UMD context.	

Function declaration	aipu_status_t aipu_config_global(const aipu_ctx_handle_t* ctx, uint64_t types, void* config)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	types	Configuration type(s)
	config	Pointer to a memory location allocated by the application where stores the configurations.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX	
Comments	This API is used to configure a specified global option.	

Function declaration	aipu_status_t aipu_load_graph(const aipu_ctx_handle_t* ctx, const char* graph, uint64_t* id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	graph	Pointer to a memory location allocated by the application where stores the graph.
	id	Graph ID returned to the caller.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_OPEN_FILE_FAIL AIPU_STATUS_ERROR_MAP_FILE_FAIL AIPU_STATUS_ERROR_UNKNOWN_BIN AIPU_STATUS_ERROR_GVERSION_UNSUPPORTED AIPU_STATUS_ERROR_TARGET_NOT_FOUND AIPU_STATUS_ERROR_INVALID_GBIN AIPU_STATUS_ERROR_BUF_ALLOC_FAIL AIPU_STATUS_ERROR_RESERVE_SRAM_FAIL	
Comments	This API loads an offline built NPU executable graph binary from the file system.	

Function declaration	aipu_status_t aipu_load_graph_helper(const aipu_ctx_handle_t* ctx, const char* graph_buf, uint32_t graph_size, uint64_t* id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	graph_buf	Loadable graph binary file data.
	graph_size	The byte size of loadable graph binary data.
	id	Graph ID returned to the caller.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_OPEN_GBIN_FAIL AIPU_STATUS_ERROR_MAP_GBIN_FAIL Other values returned by AIPU_load_graph	
Comments	This API loads a graph from the corresponding data buffer.	

Function declaration	aipu_status_t aipu_unload_graph(const aipu_ctx_handle_t* ctx, uint64_t id);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Graph ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to unload a loaded graph.	

Function declaration	aipu_status_t aipu_create_job(const aipu_ctx_handle_t* ctx, uint64_t id, uint64_t* job, aipu_create_job_cfg_t* config)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Pointer to a graph descriptor returned by aipu_load_graph.
	job	Pointer to a memory location allocated by the application where the UMD stores the created job ID.
	config	Specify the partition ID and QoS level of the job, only for AIPUv3.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	
Comments	This API is used to create a job.	

Function declaration	aipu_status_t aipu_finish_job(const aipu_ctx_handle_t* ctx, uint64_t job_id, int32_t time_out)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	job_id	Job ID returned by aipu_create_job.
	time_out	The timeout parameter for the poll system call.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR	

	AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_OP AIPU_STATUS_ERROR_JOB_EXCEPTION AIPU_STATUS_ERROR_JOB_TIMEOUT
Comments	This API is used to flush a new computation job onto the NPU. This API is a blocking one which will schedule a job and block in waiting for it to be done. The application can safely fetch the inference result data if this API returns successfully.

Function declaration	aipu_status_t aipu_flush_job(const aipu_ctx_handle_t* ctx, uint64_t job_id, void* priv)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	job_id	Job ID returned by aipu_create_job.
	priv	Pointer to the private data structure of the UMD application that is used as an argument of the callback when it is called.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to flush a new computation job onto the NPU. Unlike aipu_finish_job, this API returns as soon as a job is scheduled. By using this API with the ping-pong operation of multiple tensor buffers, the pipeline feature will be enabled. The application should query the job done status before fetching the inference result data.	

Function declaration	aipu_status_t aipu_get_job_status(const aipu_ctx_handle_t* ctx, uint64_t job_id, aipu_job_status_t* status)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	job_id	Job ID returned by aipu_create_job.
	status	Pointer to a memory location allocated by the application where the UMD stores the job status.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID	
Comments	This API is used to get the execution status of a job scheduled via aipu_flush_job. This API is a blocking one which will block in waiting for the job to be done if it is still running.	

Function declaration	aipu_status_t aipu_clean_job(const aipu_ctx_handle_t* ctx, uint64_t id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Job ID returned by aipu_create_job.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID	

Comments	This API is used to clean a finished job object in UMD. After this API successfully returns, the job ID immediately becomes invalid and cannot be used again. After this API successfully returns, the buffer handle used in creating this job will immediately be free to be used to create another job with the same graph ID.
-----------------	--

Function declaration	aipu_status_t aipu_get_tensor_count(const aipu_ctx_handle_t* ctx, uint64_t id, aipu_tensor_type_t type, uint32_t* cnt)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Job ID returned by aipu_create_job.
	type	Tensor type.
	cnt	Pointer to a memory location allocated by the application where the UMD stores the count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to get the tensor count of specified type.	

Function declaration	aipu_status_t aipu_get_tensor_descriptor(const aipu_ctx_handle_t* ctx, uint64_t id, aipu_tensor_type_t type, uint32_t tensor, aipu_tensor_desc_t* desc)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Job ID returned by aipu_create_job.
	type	Tensor type.
	tensor	Tensor ID.
	desc	Pointer to a memory location allocated by the application where the UMD stores the tensor descriptor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID	
Comments	This API is used to get the tensor descriptor of specified type.	

Function declaration	aipu_status_t aipu_load_tensor(const aipu_ctx_handle_t* ctx, uint64_t id, uint32_t tensor, const void* data)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Job ID returned by aipu_create_job.
	tensor	Tensor ID.
	ddata	Data of the input tensor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID	

	AIPU_STATUS_ERROR_INVALID_OP
Comments	This API is used to load the input tensor data.

Function declaration	aipu_status_t aipu_get_tensor(const aipu_ctx_handle_t* ctx, uint64_t job, aipu_tensor_type_t type, uint32_t tensor, void* buf);	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Job ID returned by aipu_create_job.
	type	Tensor type.
	tensor	Tensor ID.
	buf	Pointer to a memory location allocated by the application where the UMD stores the tensor data.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_TENSOR_ID AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get the tensor data of specified type.	

Function declaration	aipu_status_t aipu_config_job(const aipu_ctx_handle_t* ctx, uint64_t id, uint64_t types, void* config)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	id	Job ID returned by aipu_create_job.
	type	Configuration type(s).
	config	Pointer to a memory location allocated by the application where the application stores the configuration data struct.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_JOB_ID AIPU_STATUS_ERROR_INVALID_CONFIG	
Comments	This API is used to configure a specified option of a job.	

Function declaration	aipu_status_t aipu_get_partition_count(const aipu_ctx_handle_t* ctx, uint32_t* cnt)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	cnt	Pointer to a memory location where stores the partition count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get AIPU partition count, only for AIPUv3.	

Function declaration	aipu_status_t aipu_get_cluster_count(const aipu_ctx_handle_t* ctx, uint32_t partition_id, uint32_t* cnt)	
-----------------------------	--	--

Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	partition_id	Indicates which partition it gets the cluster count from.
	cnt	Pointer to a memory location where stores the cluster count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get AIPU cluster count, only for AIPUv3.	

Function declaration	aipu_status_t aipu_get_core_count(const aipu_ctx_handle_t* ctx, uint32_t partition_id, uint32_t cluster, uint32_t* cnt)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	partition_id	Partition ID.
	cluster	Cluster ID.
	cnt	Pointer to a memory location where stores the core count.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_OP	
Comments	This API is used to get AIPU core count, only for AIPUv3.	

Function declaration	aipu_status_t aipu_get_target(const aipu_ctx_handle_t *ctx, char *target)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	target	Pointer to a memory where stores the hardware ARCH information.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX	
Comments	This API is used to get NPU ARCH information.	

Function declaration	aipu_status_t aipu_get_device_status(const aipu_ctx_handle_t* ctx, uint32_t *status)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	status	Pointer to a memory where stores the hardware status.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX	
Comments	This API is used to get NPU status.	

Function declaration	aipu_status_t aipu_create_batch_queue(const aipu_ctx_handle_t* ctx, uint64_t graph_id, uint32_t *queue_id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.

	graph_id	Graph ID.
	queue_id	The batch queue ID will be saved in pointer queue_id.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to create a batch queue.	

Function declaration	aipu_status_t aipu_clean_batch_queue(const aipu_ctx_handle_t* ctx, uint64_t graph_id, uint32_t queue_id)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	graph_id	Graph ID.
	queue_id	The batch queue ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID	
Comments	This API is used to clean a specific batch queue.	

Function declaration	aipu_status_t aipu_add_batch(const aipu_ctx_handle_t* ctx, uint64_t graph_id, uint32_t queue_id, char *inputs[], char *outputs[])	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	graph_id	Graph ID.
	queue_id	The batch queue ID.
	inputs	Buffer pointers for input tensors.
	outputs	Buffer pointers for output tensors.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_NO_BATCH_QUEUE	
Comments	This API is used to add group buffer of one benchmark to the batch queue.	

Function declaration	aipu_status_t aipu_finish_batch(const aipu_ctx_handle_t* ctx, uint64_t graph_id, uint32_t queue_id, aipu_create_job_cfg_t *create_cfg)	
Parameter	ctx	Pointer to a context handle struct returned by aipu_init_ctx.
	graph_id	Graph ID.
	queue_id	The batch queue ID.
	create_cfg	Configuration for all batches in one queue.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR	

	AIPU_STATUS_ERROR_INVALID_CTX AIPU_STATUS_ERROR_INVALID_GRAPH_ID AIPU_STATUS_ERROR_NO_BATCH_QUEUE
Comments	This API is used to run multiple batches.

Function declaration	aipu_status_t aipu_printf(char* printf_base, char* redirect_file)	
Parameter	printf_base	Pointer to a tensor buffer where stores the printf log data.
	redirect_file	Printf output redirect file path.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_NULL_PTR AIPU_STATUS_ERROR_PRINTF_FAIL	
Comments	This API prints NPU execution log information after the corresponding job ends. The application can choose to redirect the printf log to the terminal or a file. It passes NULL to the redirect_file results in printing logs to the terminal directly. Nothing will be printed if the printf function is not used in the NPU executable binary loaded by aipu_load_graph.	

1.3.4 QNX driver programming model – Standard API

Code 14 shows the init & de-init API calls that QNX applications should follow before scheduling specific inference tasks and after exiting inference on the NPU.

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int main(int argc, char* argv[])
6.  {
7.      /* ... */
8.      aipu_status_t status = AIPU_STATUS_SUCCESS;
9.      aipu_ctx_handle_t *ctx = nullptr;
10.     const char* status_msg = nullptr;
11.
12.     status = aipu_init_ctx(&ctx);
13.     if (AIPU_STATUS_SUCCESS != status)
14.     {
15.         aipu_get_status_msg(ctx, status, &status_msg);
16.         fprintf(stderr, "Error: %s\n", status_msg);
17.         /* error handling or quitting */
18.     }
19.
20.     /* load graphs and schedule jobs */
21.
22.     status = aipu_deinit_ctx(ctx);
23.     if (AIPU_STATUS_SUCCESS != status)
24.     {
25.         aipu_get_status_msg(ctx, status, &status_msg);
26.         fprintf(stderr, "Error: %s\n", status_msg);
27.         /* error handling or quitting */
28.     }
29.
30.     /* ... */
31. }
```

Code 14. Initialization and De-Initialization

Code 15 shows the graph loading and unloading API calls that QNX applications should follow before scheduling specific inference tasks.

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int main(int argc, char* argv[])
6.  {
7.      /* ... */
8.      aipu_global_config_simulation_t sim_glb_config;
9.
10.
11.     /* initialize context */
12.
13.     /* config simulation global options */
14.     ret = aipu_config_global(ctx, AIPU_CONFIG_TYPE_SIMULATION,
15.         &sim_glb_config);
16.     if (AIPU_STATUS_SUCCESS != status)
17.     {
18.         aipu_get_status_msg(ctx, status, &status_msg);
19.         fprintf(stderr, "Error: %s\n", status_msg);
20.         /* error handling or quitting */
21.     }
22.     /* load graphs and schedule jobs */
23.
24.     /* deinit */
25.     /* ... */
26. }
```

Code 15. Graph(s) Loading

Code 16 and 17 show the inference job creating, scheduling and execution status checking model that single thread QNX applications should follow.

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int non_pipeline_scheduling()
6.  {
7.      uint64_t job_id = 0;
8.      /* load graph and allocate buffers */
9.      /* load input tensors data */
10.
11.     status = aipu_create_job(ctx, graph_id, &job_id);
12.     if (AIPU_STATUS_SUCCESS != status)
13.     {
14.         aipu_get_status_msg(ctx, status, &status_msg);
15.         fprintf(stderr, "Error: %s\n", status_msg);
16.         /* error handling or quitting */
17.     }
18.     status = aipu_finish_job(ctx, job_id, time_out);
19.     if (AIPU_STATUS_SUCCESS != status)
20.     {
21.         aipu_get_status_msg(ctx, status, &status_msg);
22.         fprintf(stderr, "Error: %s\n", status_msg);
23.         /* error handling or quitting */
24.     }
25.     else
```

```

26.  {
27.      /* job results fetching */
28.  }
29.
30.  /* ... */
31.
32.  status = aipu_clean_job(ctx, job_id);
33.  if (AIPU_STATUS_SUCCESS != status)
34.  {
35.      aipu_get_status_msg(ctx, status, &status_msg);
36.      fprintf(stderr, "Error: %s\n", status_msg);
37.      /* error handling or quitting */
38.  }
39. }
40.
41. int main(int argc, char* argv[])
42. {
43.     /* ... */
44.     /* initialize */
45.
46.     ret = non_pipeline_scheduling();
47.     /* ... */
48.
49.     /* deinit */
50.     /* ... */
51. }

```

Code 16. Non-pipeline Scheduling in Single Thread Application

```

1.  /**
2.   * @file main thread
3.   */
4.
5.  int pipeline_scheduling()
6.  {
7.      /* ... */
8.      int pipe_cnt = 3;
9.      Uint64_t job_id[pipe_cnt];
10.     int time_out = 1000;
11.     aipu_job_status_t job_status = AIPU_JOB_STATUS_NO_STATUS;
12.
13.     /* load graph and allocate buffers */
14.     /* load input tensors data */
15.
16.     /* create multiple jobs and then schedule one by one */
17.     for (uint32_t i = 0; i < pipe_cnt; i++)
18.     {
19.         /* load inputs */
20.
21.         /* create job with corresponding buffer handle */
22.         ret = aipu_create_job(ctx, graph_id, &job_id[i]);
23.         if (ret != AIPU_STATUS_SUCCESS)
24.         {
25.             aipu_get_status_msg(ctx, ret, &status_msg);
26.             fprintf(stderr, "Error: %s\n", status_msg);
27.             /* error handle or quit */
28.         }
29.
30.         ret = aipu_flush_job(ctx, job_id[i]);
31.         if (ret != AIPU_STATUS_SUCCESS)
32.         {
33.             aipu_get_status_msg(ctx, ret, &status_msg);

```



```

34.         fprintf(stderr, "Error: %s\n", status_msg);
35.         /* error handle or quit */
36.     }
37. }
38. /* application can do some other tasks if any before getting job results */
39. /* get job results */
40. for (uint32_t i = 0; i < pipe_cnt; i++)
41. {
42.     ret = aipu_get_job_status(ctx, job_id[i], &job_status);
43.     if ((ret != AIPU_STATUS_SUCCESS) || (job_status != AIPU_JOB_STATUS_DONE
44. ))
45.     {
46.         aipu_get_status_msg(ctx, ret, &status_msg);
47.         fprintf(stderr, "Error: %s\n", status_msg);
48.         /* error handle or quit */
49.     }
50.     /* fetch job #i result */
51.
52.     ret = aipu_clean_job(ctx, job_id[i]);
53.     if (ret != AIPU_STATUS_SUCCESS)
54.     {
55.         aipu_get_status_msg(ctx, ret, &status_msg);
56.         fprintf(stderr, "Error: %s\n", status_msg);
57.         /* error handle or quit */
58.     }
59. }
60. }
61.
62. int main(int argc, char* argv[])
63. {
64.     /* ... */
65.     /* initialize */
66.
67.     ret = pipeline_scheduling();
68.     /* ... */
69.
70.     /* free buffers, unload and deinit */
71.     /* ... */
72. }

```

Code 17. Pipeline Scheduling in Single Thread Application

Code 18 shows the pipeline scheduling in multithread application.

```

1. /**
2.  * @file main thread
3.  */
4.
5. int main(int argc, char* argv[])
6. {
7.
8.     /* initialize context */
9.     void *data1 = nullptr, *data2 = nullptr;
10.
11.     /* initialize data1 & data2 */
12.
13.     /* create more scheduling threads to load graphs independently and run jobs
14.      */
15.     std::thread t1(job_scheduling_thread, data1);
16.     std::thread t2(job_scheduling_thread, data2);

```

```

16.
17.     t1.join();
18.     t2.join();
19.
20.     /* deinit */
21.     /* ... */
22. }
23.
24. /**
25.  * @file scheduling_thread
26.  */
27.
28. void* job_scheduling_thread(void* data)
29. {
30.     /* ... */
31.     ret = pipeline_scheduling(); /* the same as the scheduling func in Code 5 */
32.     /* ... */
33. }

```

Code 18. Pipeline Scheduling in Multithread Application

Code 19 shows the usage of the NPU printf and profiling dump features in a single thread non-pipeline application. In real applications, you can decide whether to use these features. Note that the profiling data generated by the NPU profiler will be ready in the profiling buffer automatically after the profiling job ends. The printf data is also stored to the corresponding buffer if the printf API in the NPU binary is called.

```

1. /**
2.  * @file main_thread
3.  */
4.
5. int dump_file_helper(const char* fname, const void* src, unsigned int size)
6. {
7.     /* write <size> bytes from pointer <src> into a file named <fname> */
8. }
9.
10. int main(int argc, char* argv[])
11. {
12.     /* ... */
13.     char fname[4096];
14.     aipu_status_t status = AIPU_STATUS_SUCCESS;
15.     const char* status_msg = nullptr;
16.     aipu_tensor_desc_t desc;
17.     char *buff = nullptr;
18.
19.     /* init context and load graphs */
20.     /* ... */
21.
22.     /* buffer allocations */
23.     status = sts = aipu_get_tensor_descriptor(ctx, graph_id,
24. AIPU_TENSOR_TYPE_PRINTF, 0, &desc);
25.     if (AIPU_STATUS_SUCCESS != status)
26.     {
27.         aipu_get_status_msg(ctx, status, &status_msg);
28.         fprintf(stderr, "Alloc Error: %s\n", status_msg);
29.         /* error handling or quitting */
30.     }
31.     buff = new char[desc.size];
32.
33.     /* create, schedule a job and wait for it ends */
34.     /* ... */

```

```

34.
35.     /* Job Done Now */
36.     /* Case #9.1: print AIPU log information to terminal */
37.     status = aipu_get_tensor(m_ctx, job_id, AIPU_TENSOR_TYPE_PRINTF, 0, buff);
38.     status = aipu_printf(buff, "print.log");
39.     if (AIPU_STATUS_SUCCESS != status)
40.     {
41.         aipu_get_status_msg(ctx, status, &status_msg);
42.         fprintf(stderr, "Printf Error: %s\n", status_msg);
43.         /* error handling or quitting */
44.     }
45.
46.
47.     /* Case #9.2: get data in profiling-dump buffer(s) */
48.
49.     /* buffer allocations */
50.     status = sts = aipu_get_tensor_descriptor(ctx, graph_id,
AIPU_TENSOR_TYPE_PROFILER, 0, &desc);
51.     if (AIPU_STATUS_SUCCESS != status)
52.     {
53.         aipu_get_status_msg(ctx, status, &status_msg);
54.         fprintf(stderr, "Alloc Error: %s\n", status_msg);
55.         /* error handling or quitting */
56.     }
57.     buff = new char[desc.size];
58.
59.     status = aipu_get_tensor(m_ctx, job_id, AIPU_TENSOR_TYPE_PROFILER, 0,
buff);
60.     if (AIPU_STATUS_SUCCESS != status)
61.     {
62.         aipu_get_status_msg(ctx, status, &status_msg);
63.         fprintf(stderr, "Printf Error: %s\n", status_msg);
64.         /* error handling or quitting */
65.     }
66.
67.     /* dump the profiling data generated by profiler as a file */
68.     /* or do any other operations application would like to */
69.     snprintf(fname, sizeof(fname), "AIPU_Profiler_Data.bin");
70.     dump_file_helper(fname, buff, desc.size);
71.
72.
73.     /* clean jobs and graphs, and deinit */
74.     /* ... */
75. }

```

Code 19. Usage of NPU Printf and Profiling Dump Features

1.4 Bare-metal-based driver

The bare-metal driver is compiled as a static library, which uses similar user interfaces and workflow as its Linux counterparts, except that no standard operating system APIs are used. The interfaces are listed in *5.4.2 Bare-metal platform programming model*.

User applications should call the APIs in a relatively fixed order and check the return values step by step to ensure correct operations.

The typical usage of bare-metal driver interfaces in an application is as follows:

1. The application calls the address config API first to initialize the NPU External Control Register group base address on the SoC. The memory offset between the SoC and the NPU will also be updated in this API.
2. After the address is set, the application should call the context initialization API to initialize bare-metal runtime context.
3. After context initialization, the application should call the graph loading API to load an offline built binary. If successful, a graph descriptor will be returned. Multiple graph binaries can load in order before next steps.
4. After graph loading, the application should allocate at least one copy of tensor buffers for a graph to be executed. For every graph, the application can allocate multiple copies of tensor buffers to realize executions in the pipeline.
5. According to the graph descriptor information in step 3, the application should allocate the corresponding space for IO tensor descriptors.
6. After IO tensor descriptors are allocated, the application should call the tensor descriptor get API to pass the descriptors to the driver to get specified IO buffer address allocated by the bare-metal runtime driver.
7. The application should load the input frame data into the corresponding buffers obtained in step 6.
8. After the preceding steps finish, the application can call the start API and specify the pipeline index to trigger NPU running.
9. The application should call the NPU status get API to check NPU running status in polling mode. If there are more than one copy of tensor buffers, the application can load the next input into another buffer while the NPU is running.
10. If the application wants to run multiple input frames, it should follow steps 7-9 repeatedly.
11. Before exiting, the application should call the tensor buffer free API, graph unload API and context de-init API in order to free all buffers.

1.4.1 Bare-metal user interface

Data structure

Data structure	Type	Values	Comment
aipu_tensor_layout_t	enum	TENSOR_LAYOUT_NONE (0x0)	-
		TENSOR_LAYOUT_NHWC (0x1)	-
		TENSOR_LAYOUT_NCHW (0x2)	-
		TENSOR_LAYOUT_NWH (0x3)	-
		TENSOR_LAYOUT_NC (0x4)	-

Data structure	Type	Members	Member type	Comment
aipu_tensor_shape_t	struct	N	unsigned int	-
		H		-
		W		-
		C		-

Data structure	Type	Values	Comment
aipu_data_type_t	enum	TENSOR_DATA_TYPE_NONE (0x0)	Unsigned char.
		TENSOR_DATA_TYPE_U8 (0x2)	Unsigned char.
		TENSOR_DATA_TYPE_S8 (0x3)	Signed char.
		TENSOR_DATA_TYPE_U16 (0x4)	Unsigned short.
		TENSOR_DATA_TYPE_S16 (0x5)	Signed short.

Data structure	Type	Members	Member type	Comment
aipu_tensor_fmt_t	struct	layout	aipu_tensor_layout_t	-
		shape	aipu_tensor_shape_t	-
		data_type	aipu_data_type_t	-

Data structure	Type	Values	Comment
aipu_tensor_io_type_t	enum	TENSOR_IO_TYPE_IN(0x0)	-
		TENSOR_IO_TYPE_OUT(0x1)	-
		TENSOR_IO_TYPE_DUMP(0x2)	-

Data structure	Type	Members	Member type	Comment
aipu_tensor_desc_t	struct	fmt	aipu_tensor_fmt_t	-
		io_type	aipu_tensor_io_type_t	-
		buffer_desc	aipu_buffer_desc_t*	-

Data structure	Type	Members	Member type	Comment
aipu_buffer_desc_t	struct	index	unsigned int	Tensor ID.
		size	unsigned int	Tensor size.
		addr	void*	Tensor buffer address.

Data structure	Type	Members	Member type	Comment
aipu_graph_desc_t	struct	graph_id	int	Graph ID.
		input_num	unsigned int	Input number of the graph ID.
		output_num	unsigned int	Output number of the graph ID.
		dump_num	unsigned int	Dump number of the graph ID.

Data structure	Type	Values	Comment
aipu_task_status_t	enum	AIPU_TASK_STATUS_NO_STATUS(0x0)	-
		AIPU_TASK_STATUS_RUNNING(0x1)	-
		AIPU_TASK_STATUS_DONE(0x2)	-
		AIPU_TASK_STATUS_EXCEPTION(0x3)	-

Data structure	Type	Values	Comment
aipu_status_t	enum	AIPU_STATUS_SUCCESS	-
		AIPU_STATUS_ERROR	-
		AIPU_STATUS_ERROR_INVALID_ARGS	-
		AIPU_STATUS_ERROR_GRAPH_NOT_EXIST	-
		AIPU_STATUS_ERROR_GRAPH_INIT_FAIL	-
		AIPU_STATUS_ERROR_INVALID_BIN	-
		AIPU_STATUS_ERROR_INVALID_BIN_HEAD_SIZE	-
		AIPU_STATUS_ERROR_INVALID_BIN_MAGIC	-
		AIPU_STATUS_ERROR_PARSE_BSS_STATIC_FAIL	-
		AIPU_STATUS_ERROR_PARSE_BSS_REUSE_FAIL	-
		AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	-
		AIPU_STATUS_ERROR_INVALID_BUFFER	-
		AIPU_STATUS_ERROR_INVALID_GRAPH_VERSION	-
		AIPU_STATUS_ERROR_INIT_CONTEXT_FAIL	-
		AIPU_STATUS_ERROR_DEINIT_CONTEXT_FAIL	-
		AIPU_STATUS_ERROR_CONTEXT_NOT_INIT	-
		AIPU_STATUS_ERROR_TENSOR_NOT_ALLOC	-
		AIPU_STATUS_ERROR_TENSOR_NOT_EXIST	-
		AIPU_STATUS_ERROR_RUN_BUSY	-
		AIPU_STATUS_ERROR_RUN_EXCEPTION	-

		AIPU_STATUS_ERROR_RUN_SUCCESS_BUT_WARNING	-
		AIPU_STATUS_ERROR_GRAPH_IS_RUNNING	-
		AIPU_STATUS_ERROR_GRAPH_ID_NOT_RUN	-
		AIPU_STATUS_ERROR_GRAPH_NOT_MATCH_HARDWARE	-
		AIPU_STATUS_ERROR_PARSE_BSS_LINKED_FAIL	-
		AIPU_STATUS_ERROR_INVALID_IO_HANDLE	-
		AIPU_STATUS_MAX	-

User application interfaces

Function declaration	void aipu_config_address(unsigned long ctrl_reg_base_addr, unsigned long memory_addr_offset);	
Parameter	ctrl_reg_base_addr	External control register base address.
	memory_addr_offset	Memory offset between the CPU and NPU.
Return value	void	
Comments	This API is used to set address information between the CPU and NPU.	

Function declaration	aipu_status_t aipu_init_ctx(void);	
Parameter	void	
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_INIT_CONTEXT_FAIL	
Comments	This API is used to initialize the NPU context.	

Function declaration	aipu_status_t aipu_deinit_ctx(void);	
Parameter	void	
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_DEINIT_CONTEXT_FAIL	
Comments	This API is used to destroy the NPU context.	

Function declaration	aipu_status_t aipu_load_graph(void *graph, aipu_graph_desc_t *graph_desc);	
Parameter	graph	Graph binary buffer pointer.
	graph_desc	Graph descriptor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_BUF_ALLOC_FAIL AIPU_STATUS_ERROR_GRAPH_INIT_FAIL AIPU_STATUS_ERROR_INVALID_ARGS AIPU_STATUS_ERROR_INVALID_BIN_MAGIC AIPU_STATUS_ERROR_INVALID_BIN_HEAD_SIZE AIPU_STATUS_ERROR_INVALID_GRAPH_VERSION AIPU_STATUS_ERROR_GRAPH_NOT_MATCH_HARDWARE AIPU_STATUS_ERROR_PARSE_BSS_STATIC_FAIL AIPU_STATUS_ERROR_PARSE_BSS_REUSE_FAIL	

	AIPU_STATUS_ERROR_PARSE_BSS_LINKED_FAIL
Comments	This API is used to load a graph binary for the driver to parse and allocate static buffers.

Function declaration	aipu_status_t aipu_unload_graph(int graph_id);	
Parameter	graph_id	Graph ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_INVALID_ARGS AIPU_STATUS_ERROR_GRAPH_IS_RUNNING	
Comments	This API is used to unload a loaded graph.	

Function declaration	aipu_status_t aipu_alloc_tensor_buffer(int graph_id, int io_buff_num);	
Parameter	graph_id	Graph ID.
	io_buff_num	Copies of tensor buffers.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_INVALID_ARGS AIPU_STATUS_ERROR_INVALID_BUFFER AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	
Comments	This API is used to allocate buffers for all input and output tensors of a graph. The io_buff_num must be >=1, which indicates that how many copies of tensor buffers are allocated by runtime. After the API returns successfully, the valid tensor buffer handle is >=0 and < io_buff_num. If the io_buff_num is set to 1, the pipeline feature will be disabled.	

Function declaration	aipu_status_t aipu_get_tensor_desc(int graph_id, int io_buff_handle, aipu_tensor_io_type_t io_type, aipu_tensor_desc_t* tdesc);	
Parameter	graph_id	Graph ID.
	io_buff_handle	The valid tensor buffer handle which is >=0 and < io_buff_num.
	io_type	Tensor IO type.
	tdesc	Tensor descriptor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_TENSOR_NOT_ALLOC AIPU_STATUS_ERROR_INVALID_IO_HANDLE AIPU_STATUS_ERROR_TENSOR_INVALID_ARGS	
Comments	This API is used to get input/output/dump tensor descriptor of a graph. Applications load inputs and get outputs through these descriptors.	

Function declaration	aipu_status_t aipu_start(int graph_id, int io_buff_handle);	
-----------------------------	---	--

Parameter	graph_id	Graph ID.
	io_buff_handle	Valid buffer handle which has loaded inputs for aipu_v1v2. It is useless for aipu_v3, and recommended to use value 0.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_TENSOR_NOT_ALLOC AIPU_STATUS_ERROR_INVALID_IO_HANDLE AIPU_STATUS_ERROR_RUN_BUSY AIPU_STATUS_ERROR_RUN_EXCEPTION AIPU_STATUS_ERROR_RUN_SUCCESS_BUT_WARNING	
Comments	<ul style="list-style-type: none"> Before calling this API, the aipu_alloc_tensor_buffer() and tensor data should be put into tdesc. aiPU_start() and aipu_get_status() should be one-to-one, otherwise, the status will be covered by the next running graph. If this API returns the exception state, the application can call the aipu_get_status() API to clear the exception status. 	

Function declaration	aiPU_status_t aipu_get_status(int graph_id, aipu_task_status_t *status);	
Parameter	graph_id	Graph ID.
	status	NPU task status.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_GRAPH_ID_NOT_RUN AIPU_STATUS_ERROR_RUN_EXCEPTION	
Comments	<ul style="list-style-type: none"> This API is non-blocking, so the application needs to poll this API to get status. When this API finds that the NPU runs into exception state, the API will clear the exception state, and then return the exception state value. 	

Function declaration	aiPU_status_t aipu_free_tensor_buffer(int graph_id);	
Parameter	graph_id	Graph ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST	
Comments	This API is used to free the tensor buffer.	

1.4.2 Bare-metal platform programming model

Code 20 shows the NPU address configuration API usage. The API should be called before any operation. At the same time, the default values of AIPU_CTRL_REG_BASE_ADDR and MEMORY_ADDR_OFFSET are set as (0x64000000u, 0x80000000u) in the bare-metal runtime driver.

```

1. #define AIPU_CTRL_REG_BASE_ADDR (0x64000000u)
2. #define MEMORY_ADDR_OFFSET (0x80000000u)
3.
4. aipu_config_address(AIPU_CTRL_REG_BASE_ADDR, MEMORY_ADDR_OFFSET);

```

```
5.
6. /* ... */
```

Code 20. NPU Driver Library Use Case – for Bare-metal (Part 1)

Code 21 shows the context init & de-init API call. The init API should be called after `aipu_config_address()` to create an environment for next APIs. The de-init API should be called when the application tries exit on the NPU.

```
1. aipu_status_t ret = AIPU_STATUS_SUCCESS;
2.
3. ret = aipu_init_ctx();
4. if(AIPU_STATUS_SUCCESS != ret){
5.     printf("aipu context init fail\n");
6.     /* error handling or quitting */
7. }
8.
9. /* ... */
10.
11. ret = aipu_deinit_ctx();
12. if (AIPU_STATUS_SUCCESS != ret)
13. {
14.     printf("aipu context deinit fail\n");
15.     /* error handling or quitting */
16. }
17.
18. /* ... */
```

Code 21. NPU Driver Library Use Case – for Bare-metal (Part 2)

Code 22 shows the graph loading and unloading. After the loading graph is called, the application can get the graph ID and tensor number information. If the application will not use the graph, it can call the unloading API. Before the unloading graph API is called, the application should call the tensor free API to free tensor buffer allocated by the driver at first.

```
1. aipu_graph_desc_t graph_desc = {0};
2.
3. /* config address of cpu and aipu */
4. /* initialize context */
5.
6. /* graph_addr is pointer of graph binary */
7. ret = aipu_load_graph(graph_addr, &graph_desc);
8. if(AIPU_STATUS_SUCCESS != ret){
9.     printf("aipu load graph fail\n");
10.    /* error handling or quitting */
11. }
12. /* continue loading if multigraphs used */
13. /* ... */
14.
15. /* unload graph */
16. ret = aipu_unload_graph(graph_desc.graph_id);
17. if (AIPU_STATUS_SUCCESS != ret)
18. {
19.     printf("aipu unload graph fail\n");
20.     /* error handling or quitting */
21. }
22.
23. /* ... */
```

Code 22. NPU Driver Library Use Case – for Bare-metal (Part 3)

Code 23 shows the memory allocation and free API. After loading graph, the application can call the tensor allocate buffer API which supports one or more copies of tensor buffers allocated. These buffers will provide the space for the application to load inputs and get outputs.

```

1. /* config address of cpu and aipu */
2. /* initialize context */
3. /* load graph */
4.
5. /* alloc tensor buffer, only one copy of buffer, no use of pipeline */
6. ret = aipu_alloc_tensor_buffer(graph_desc.graph_id, 1);
7. if(AIPU_STATUS_SUCCESS != ret){
8.     printf("aipu alloc tensor buffer fail\n");
9.     /* error handling or quitting */
10. }
11.
12. /* alloc tensor descriptor buffer */
13. /* get tensor descriptor entity from driver */
14. /* load inputs */
15. /* start aipu and get status*/
16. /* ... */
17.
18. /* free all tensor buffer */
19. ret = aipu_free_tensor_buffer(graph_desc.graph_id);
20. if (AIPU_STATUS_SUCCESS != ret)
21. {
22.     printf("aipu free tensor buffer fail\n");
23.     /* error handling or quitting */
24. }
25.
26. /* ... */

```

Code 23. NPU Driver Library Use Case – for Bare-metal (Part 4)

Code 24 shows the entity getting API of descriptors. Before calling this API, the application should create input and output tensor descriptor entities, and pass them to the driver level through the following API.

```

1. /* config address of cpu and aipu */
2. /* initialize context */
3. /* load graph */
4. /* alloc tensor buffer */
5. /* alloc tensor descriptor buffer */
6.
7. /* get input tensor descriptor entity through the io buffer handle */
8. ret = aipu_get_tensor_desc(graph_desc.graph_id,
9.     0, TENSOR_IO_TYPE_IN, &input_tensor_desc);
10. if(AIPU_STATUS_SUCCESS != ret){
11.     printf("aipu get tensor desc fail\n");
12.     /* error handling or quitting */
13. }
14. /* get output tensor descriptor entity through the io buffer handle */
15. ret = aipu_get_tensor_desc(graph_desc.graph_id, 0, TENSOR_IO_TYPE_OUT, &output
16.     t_tensor_desc);
17. if(AIPU_STATUS_SUCCESS != ret){
18.     printf("aipu get tensor desc fail\n");
19.     /* error handling or quitting */
20. }
21. /* load inputs */
22. /* start aipu and get status */
23. /* ... */

```

Code 24. NPU Driver Library Use Case – for Bare-metal (Part 5)

Code 25 shows the action of loading inputs. `load_inputs` should be realized by the application itself. It is a process of copying input frames from the application layer to the driver layer.

```

1.  /* config address of cpu and aipu */
2.  /* initialize context */
3.  /* load graph */
4.  /* alloc tensor buffer */
5.  /* alloc tensor descriptor buffer */
6.  /** <load_inputs> is aim to copy input_src buffer to input tensor  buffer al
    located by runtime driver
7.    <input_src> is pointer of source input frame buffer
8.  */
9.  if(1 == graph_desc.input_num){
10.     ret = load_inputs(&input_tensor_desc.buffer_desc[0], input0_src);
11.     if(AIPU_STATUS_SUCCESS != ret){
12.         printf("load input fail\n");
13.         /* error handling or quitting */
14.     }
15. }
16. else if(2 == graph_desc.input_num){
17.     ret = load_inputs(&input_tensor_desc.buffer_desc[0], input0_src);
18.     if(AIPU_STATUS_SUCCESS != ret){
19.         printf("load input fail\n");
20.         /* error handling or quitting */
21.     }
22.     ret = load_inputs(&input_tensor_desc.buffer_desc[1], input1_src);
23.     if(AIPU_STATUS_SUCCESS != ret){
24.         printf("load input fail\n");
25.         /* error handling or quitting */
26.     }
27. }
28. /* start aipu and get status */
29. /* ... */

```

Code 25. NPU Driver Library Use Case – for Bare-metal (Part 6)

Code 26 shows the NPU start & status get API. The application can get status in polling mode.

```

1.  /* config address of cpu and aipu */
2.  /* initialize context */
3.  /* load graph */
4.  /* alloc tensor buffer */
5.  /* alloc tensor descriptor */
6.  /* load inputs */
7.  /* start aipu and specify which tensor buffer aipu uses */
8.  ret = aipu_start(graph_desc.graph_id, 0);
9.  if(!((AIPU_STATUS_SUCCESS == ret) || (AIPU_STATUS_ERROR_RUN_SUCCESS_BUT_WAR
NING == ret))){
10.     printf("aipu start fail\n");
11.     /* error handling or quitting */
12. }
13. /* check aipu running status, while(1) just for an example */
14. while(1){
15.     aipu_get_status(graph_desc.graph_id, &aipu_status);
16.     if(AIPU_TASK_STATUS_RUNNING != aipu_status){
17.         break;
18.     }
19. }
20. /* ... */

```

Code 26. NPU Driver Library Use Case – for Bare-metal (Part 7)

Code 27 shows how to use the pipeline feature. The application will apply two copies of tensor buffers, and complete the ping-pong operation. The `get_next_ready_pipeline` and `get_next_idle_pipeline` should be realized by the application itself. They are in charge of managing the pipeline and recording which pipeline is idle or ready.

```

1.  /* config address of cpu and aipu */
2.  /* initialize context */
3.  /* load graph */
4.  /* alloc two copies of tensor buffer */
5.  aipu_alloc_tensor_buffer(graph_desc.graph_id, 0);
6.
7.  /* alloc tensor descriptor buffer */
8.  /* ... */
9.  /* get input tensor descriptor entity through the io buffer handle one by
   one */
10. ret = aipu_get_tensor_desc(graph_desc.graph_id,0, TENSOR_IO_TYPE_IN, &input_
    tensor_desc0);
11.
12. ret = aipu_get_tensor_desc(graph_desc.graph_id,1, TENSOR_IO_TYPE_IN, &input_
    tensor_desc1);
13.
14. while(1){
15.     /* get the ready pipeline, if none, load input */
16.     pipeline = get_next_ready_pipeline();
17.     /* no ready pipeline, load input into it */
18.     /* load input */
19.     /* ... */
20.
21.     /* start aipu and specify which tensor buffer aipu uses */
22.     aipu_start(graph_desc.graph_id, pipeline);
23.
24.     while(1){
25.         aipu_get_status(graph_desc.graph_id, &aipu_status);
26.         /* when aipu is ongoing, load input to an idle pipeline */
27.         if(AIPU_TASK_STATUS_RUNNING == aipu_status){
28.             index = get_next_idle_pipeline();
29.             /* load input */
30.             /* ... */
31.         }
32.     }
33. }
34. /* free tensor buffer */
35. /* unload graph */
36. /* context deinit */
37.

```

Code 27. NPU Driver Library Use Case – for Bare-metal (Part 8)

1.5 RTOS-based driver (for customized solutions only)

The RTOS driver is also compiled as a static library, which keeps the same user interfaces and workflow as the bare-metal driver. Therefore, all programming models in bare metal can be used in RTOS. The RTOS driver provided by Arm China is implemented based on FreeRTOS v10.1.1, but it is easy to be ported to other platforms, like RTX, QNX and RT-Thread.

User applications should call the APIs in a relatively fixed order and check the return values step by step to ensure correct operations.

The typical usage of RTOS driver interfaces in an application is as follows:

1. The application calls the address config API first to initialize the NPU External Control Register group base address on the SoC. The memory offset between the SoC and the NPU will also be updated in this API.
2. After the address is set, the application should call the context initialization API to initialize RTOS runtime context.
3. After context initialization, the application should call the graph loading API to load an offline built binary. If successful, a graph descriptor will be returned. Multiple graph binaries can load in order before next steps.
4. After graph loading, the application should allocate at least one copy of tensor buffers for a graph to be executed. For every graph, the application can allocate multiple copies of tensor buffers to realize executions in the pipeline.
5. According to the graph descriptor information in step 3, the application should allocate the corresponding space for IO tensor descriptors.
6. After IO tensor descriptors are allocated, the application should call the tensor descriptor get API to pass the descriptors to the driver to get specified IO buffer address allocated by the RTOS runtime driver.
7. The application should load the input frame data into the corresponding buffers obtained in step 6.
8. After the preceding steps finish, the application can call the start API and specify the pipeline index to trigger NPU running.
9. The application should call the NPU status get API to check NPU running status in polling mode. If there are more than one copy of tensor buffers, the application can load the next input into another buffer while the NPU is running.
10. If the application wants to run multiple input frames, it should follow steps 7-9 repeatedly.
11. Before exiting, the application should call the tensor buffer free API, graph unload API and context de-init API in order to free all buffers.

1.5.1 FreeRTOS demo

Arm China provides a RTOS demo based on FreeRTOS v10.1.1 in AI610-SDK-1015-xxxx-xxx/FreeRTOS-driver/demo/. It is compiled into an executable OS ELF file, which consists of the FreeRTOS lib, the AIPU bare-metal driver and the NPU application.

The typical usage of the demo is as follows:

1. Build the FreeRTOS kernel lib. The lib must include basic APIs of FreeRTOS, such as the memory management component and task scheduler component.
2. According to different environments, choose the file system component, the Ethernet component and other components. They are determined by the users and development boards. It is not a mandatory option.
3. Put the compiled lib into the bsp/lib/ folder and put the FreeRTOS OS related head files into the bsp/include/ folder.
4. Create the main() function, which will be called by the FreeRTOS OS to set up an NPU application (for example, run a neural network model). In this main() function, create a task (using FreeRTOS API xTaskCreate()) to call the bare-metal driver APIs to set up an application. For more information about the FreeRTOS APIs, see the FreeRTOS documentation or source code.

Currently the driver supports a single task only.

1.5.2 RTOS user interface

Data structure

Data structure	Type	Values	Comment
aipu_tensor_layout_t	enum	TENSOR_LAYOUT_NONE (0x0)	-
		TENSOR_LAYOUT_NHWC (0x1)	-
		TENSOR_LAYOUT_NCHW (0x2)	-
		TENSOR_LAYOUT_NWH (0x3)	-
		TENSOR_LAYOUT_NC (0x4)	-

Data structure	Type	Members	Member type	Comment
aipu_tensor_shape_t	struct	N	unsigned int	-
		H		-
		W		-
		C		-

Data structure	Type	Values	Comment
aipu_data_type_t	enum	TENSOR_DATA_TYPE_NONE (0x0)	Unsigned char.
		TENSOR_DATA_TYPE_U8 (0x2)	Unsigned char.
		TENSOR_DATA_TYPE_S8 (0x3)	Signed char.
		TENSOR_DATA_TYPE_U16 (0x4)	Unsigned short.
		TENSOR_DATA_TYPE_S16 (0x5)	Signed short.

Data structure	Type	Members	Member type	Comment
aipu_tensor_fmt_t	struct	layout	aipu_tensor_layout_t	-
		shape	aipu_tensor_shape_t	-
		data_type	aipu_data_type_t	-

Data structure	Type	Values	Comment
aipu_tensor_io_type_t	enum	TENSOR_IO_TYPE_IN(0x0)	-
		TENSOR_IO_TYPE_OUT(0x1)	-
		TENSOR_IO_TYPE_DUMP(0x2)	-

Data structure	Type	Members	Member type	Comment
aipu_tensor_desc_t	struct	fmt	aipu_tensor_fmt_t	-
		io_type	aipu_tensor_io_type_t	-
		buffer_desc	aipu_buffer_desc_t*	-

Data structure	Type	Members	Member type	Comment
aipu_buffer_desc_t	struct	index	unsigned int	Tensor ID.
		size	unsigned int	Tensor size.
		addr	void*	Tensor buffer address.

Data structure	Type	Members	Member type	Comment
aipu_graph_desc_t	struct	graph_id	int	Graph ID.
		input_num	unsigned int	Input number of the graph ID.
		output_num	unsigned int	Output number of the graph ID.
		dump_num	unsigned int	Dump number of the graph ID.

Data structure	Type	Values	Comment
aipu_task_status_t	enum	AIPU_TASK_STATUS_NO_STATUS(0x0)	-
		AIPU_TASK_STATUS_RUNNING(0x1)	-
		AIPU_TASK_STATUS_DONE(0x2)	-
		AIPU_TASK_STATUS_EXCEPTION(0x3)	-

Data structure	Type	Values	Comment
aipu_status_t	enum	AIPU_STATUS_SUCCESS	-
		AIPU_STATUS_ERROR	-
		AIPU_STATUS_ERROR_INVALID_ARGS	-
		AIPU_STATUS_ERROR_GRAPH_NOT_EXIST	-
		AIPU_STATUS_ERROR_GRAPH_INIT_FAIL	-

	AIPU_STATUS_ERROR_INVALID_BIN	-
	AIPU_STATUS_ERROR_INVALID_BIN_HEAD_SIZE	-
	AIPU_STATUS_ERROR_INVALID_BIN_MAGIC	-
	AIPU_STATUS_ERROR_PARSE_BSS_STATIC_FAIL	-
	AIPU_STATUS_ERROR_PARSE_BSS_REUSE_FAIL	-
	AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	-
	AIPU_STATUS_ERROR_INVALID_BUFFER	-
	AIPU_STATUS_ERROR_INVALID_GRAPH_VERSION	-
	AIPU_STATUS_ERROR_INIT_CONTEXT_FAIL	-
	AIPU_STATUS_ERROR_DEINIT_CONTEXT_FAIL	-
	AIPU_STATUS_ERROR_CONTEXT_NOT_INIT	-
	AIPU_STATUS_ERROR_TENSOR_NOT_ALLOC	-
	AIPU_STATUS_ERROR_TENSOR_NOT_EXIST	-
	AIPU_STATUS_ERROR_RUN_BUSY	-
	AIPU_STATUS_ERROR_RUN_EXCEPTION	-
	AIPU_STATUS_ERROR_RUN_SUCCESS_BUT_WARNING	-
	AIPU_STATUS_ERROR_GRAPH_IS_RUNNING	-
	AIPU_STATUS_ERROR_GRAPH_ID_NOT_RUN	-
	AIPU_STATUS_ERROR_GRAPH_NOT_MATCH_HARDWARE	-
	AIPU_STATUS_ERROR_PARSE_BSS_LINKED_FAIL	-
	AIPU_STATUS_ERROR_INVALID_IO_HANDLE	-
	AIPU_STATUS_MAX	-

User application interfaces

Function declaration	void aipu_config_address(unsigned long ctrl_reg_base_addr, unsigned long memory_addr_offset);	
Parameter	ctrl_reg_base_addr	External control register base address.
	memory_addr_offset	Memory offset between the CPU and NPU.
Return value	void	
Comments	This API is used to set address information between the CPU and NPU. It is not thread safe.	

Function declaration	aipu_status_t aipu_init_ctx(void);	
Parameter	void	
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_INIT_CONTEXT_FAIL	
Comments	This API is used to initialize the NPU context. It is not thread safe.	

Function declaration	aipu_status_t aipu_deinit_ctx(void);	
Parameter	void	

Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_DEINIT_CONTEXT_FAIL
Comments	This API is used to destroy the NPU context. It is not thread safe.

Function declaration	aipu_status_t aipu_load_graph(void *graph, aipu_graph_desc_t *graph_desc);	
Parameter	graph	Graph binary buffer pointer.
	graph_desc	Graph descriptor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_BUF_ALLOC_FAIL AIPU_STATUS_ERROR_GRAPH_INIT_FAIL AIPU_STATUS_ERROR_INVALID_ARGS AIPU_STATUS_ERROR_INVALID_BIN_MAGIC AIPU_STATUS_ERROR_INVALID_BIN_HEAD_SIZE AIPU_STATUS_ERROR_INVALID_GRAPH_VERSION AIPU_STATUS_ERROR_GRAPH_NOT_MATCH_HARDWARE AIPU_STATUS_ERROR_PARSE_BSS_STATIC_FAIL AIPU_STATUS_ERROR_PARSE_BSS_REUSE_FAIL AIPU_STATUS_ERROR_PARSE_BSS_LINKED_FAIL	
Comments	This API is used to load a graph binary for the driver to parse and allocate static buffers. It is thread safe.	

Function declaration	aipu_status_t aipu_unload_graph(int graph_id);	
Parameter	graph_id	Graph ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_INVALID_ARGS AIPU_STATUS_ERROR_GRAPH_IS_RUNNING	
Comments	This API is used to unload a loaded graph. It is thread safe.	

Function declaration	aipu_status_t aipu_alloc_tensor_buffer(int graph_id, int io_buff_num);	
Parameter	graph_id	Graph ID.
	io_buff_num	Copies of tensor buffers.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_INVALID_ARGS AIPU_STATUS_ERROR_INVALID_BUFFER AIPU_STATUS_ERROR_BUF_ALLOC_FAIL	
Comments	This API is used to allocate buffers for all input and output tensors of a graph. The io_buff_num must be >=1, which indicates that how many copies of tensor buffers are allocated by runtime. After the API returns successfully, the valid tensor buffer handle is >=0 and < io_buff_num. If the io_buff_num is set to 1, the pipeline feature will be disabled. It is thread safe.	

Function declaration	aipu_status_t aipu_get_tensor_desc(int graph_id, int io_buff_handle, aipu_tensor_io_type_t io_type, aipu_tensor_desc_t* tdesc);	
Parameter	graph_id	Graph ID.
	io_buff_handle	The valid tensor buffer handle which is ≥ 0 and $< io_buff_num$.
	io_type	Tensor IO type.
	tdesc	Tensor descriptor.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_TENSOR_NOT_ALLOC AIPU_STATUS_ERROR_INVALID_IO_HANDLE AIPU_STATUS_ERROR_TENSOR_INVALID_ARGS	
Comments	This API is used to get input/output/dump tensor descriptor of a graph. Applications load inputs and get outputs through these descriptors. It is thread safe.	

Function declaration	aipu_status_t aipu_start(int graph_id, int io_buff_handle);	
Parameter	graph_id	Graph ID.
	io_buff_handle	Valid buffer handle which has loaded inputs for aipu_v1/v2. It is useless for aipu_v3, and recommended to use value 0.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_TENSOR_NOT_ALLOC AIPU_STATUS_ERROR_INVALID_IO_HANDLE AIPU_STATUS_ERROR_RUN_BUSY AIPU_STATUS_ERROR_RUN_EXCEPTION AIPU_STATUS_ERROR_RUN_SUCCESS_BUT_WARNING	
Comments	<ul style="list-style-type: none"> Before calling this API, putting aipu_alloc_tensor_buffer() and tensor data into tdes should be done. aiipu_start() and aipu_get_status() should be one-to-one, otherwise, the status will be covered by the next running graph. If this API returns the exception state, the application can call the aipu_get_status() API to clear the exception status. It is thread safe. 	

Function declaration	aipu_status_t aipu_get_status(int graph_id, aipu_task_status_t *status);	
Parameter	graph_id	Graph ID.
	status	NPU task status.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST AIPU_STATUS_ERROR_GRAPH_ID_NOT_RUN AIPU_STATUS_ERROR_RUN_EXCEPTION	
Comments	<ul style="list-style-type: none"> This API is non-blocking, so the application needs to poll this API to get status. 	

	<ul style="list-style-type: none"> When this API finds that the NPU runs into exception state, the API will clear the exception state, and then return the exception state value. It is thread safe.
--	--

Function declaration	aipu_status_t aipu_free_tensor_buffer(int graph_id);	
Parameter	graph_id	Graph ID.
Return value	AIPU_STATUS_SUCCESS AIPU_STATUS_ERROR_CONTEXT_NOT_INIT AIPU_STATUS_ERROR_GRAPH_NOT_EXIST	
Comments	This API is used to free the tensor buffer. It is thread safe.	

1.5.3 RTOS platform programming model

A programming model in bare metal is also suitable for RTOS. For the usage of the interfaces, see 5.4.2 *Bare-metal platform programming model*. In addition, RTOS is an operating system which supports the multi-task feature that is not supported in bare metal. Therefore, a multi-task programming model is only applicable to RTOS.

Code 28 shows a multi-task programming model. The main task is responsible for the NPU context initialization and de-initialization, and each subtask loads the graph and runs the inference job. Before the subtask loads the graph, the main task must complete the context initialization first.

```

1. #define AIPU_CTRL_REG_BASE_ADDR (0x64000000u)
2. #define MEMORY_ADDR_OFFSET (0x80000000u)
3. /* main task */
4. void main()
5. {
6.     /* config address of cpu and aipu */
7.     aipu_config_address(AIPU_CTRL_REG_BASE_ADDR, MEMORY_ADDR_OFFSET);
8.     aipu_status_t ret = AIPU_STATUS_SUCCESS;
9.     /* initialize context */
10.    ret = aipu_init_ctx();
11.
12.    /* waiting all subtasks complete and exit */
13.    /* ... */
14.
15.    /* de-init context */
16.    ret = aipu_deinit_ctx();
17. }
18.
19. /* task 1 */
20. void task1()
21. {
22.     /* load graph */
23.     /* alloc tensor buffer */
24.     /* get tensor descriptor entity from driver */
25.     /* load inputs */
26.     /* start aipu and get status*/
27.     /* ... */
28.     /* free all tensor buffer */
29.     /* unload graph */
30. }
31.
32. /* task 2 */
33. void task2()
34. {
35.     /* load graph */

```

```
36.    /* alloc tensor buffer */
37.    /* get tensor descriptor entity from driver */
38.    /* load inputs */
39.    /* start aipu and get status*/
40.    /* ... */
41.    /* free all tensor buffer */
42.    /* unload graph */
43. }
44.
45. /* task n */
46. /* ... */
```

Code 28. NPU Driver Library Use Case – for RTOS (Part 1)

Chapter 2

NPU runtime

This chapter describes the Zhouyi NPU runtime.

It contains the following sections:

- [2.1 About the NPU runtime on page 2-87.](#)
- [2.2 Arm NN runtime on page 2-88.](#)
- [2.3 Android neural networks runtime on page 2-90.](#)
- [2.4 TensorFlow Lite delegate runtime on page 2-94.](#)

2.1 About the NPU runtime

The Zhouyi NPU runtime is a program that provides a runtime system for the *Neural Network* (NN). It sets up the running environment, with the handle of the NPU and maybe some other hardware, along with their drivers. In the NPU runtime, it has a hardware abstraction for the NPU and other hardware, to treat different hardware with the same interfaces as the NPU.

The NPU runtime also provides a set of application interfaces, with which engineers can build their neural network applications.

The Zhouyi Compass NPU runtime is developed based on the Arm NN inference engine and Android NNAPI.

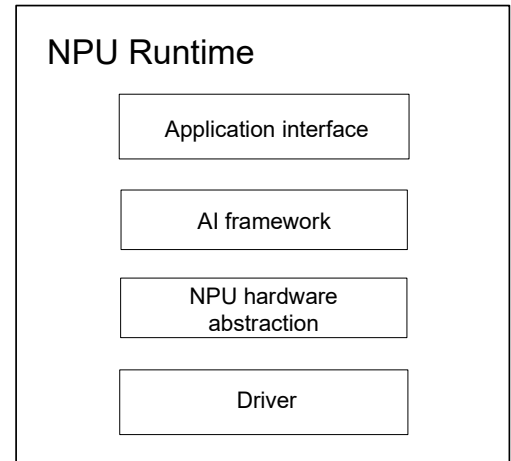


Figure 2-1 NPU runtime architecture

2.2 Arm NN runtime

The Arm NN runtime for Zhouyi NPUs contains the original Arm NN framework and the Zhouyi backend.

Arm NN is an inference engine for CPUs, GPUs and NPUs. It bridges the gap between existing NN frameworks and the underlying IPs. It enables efficient translation of existing NN frameworks such as TFLite and ONNX, allowing them to run efficiently, without modification, across Arm Cortex-A CPUs, Arm Mali GPUs, and third-party NPUs. Based on the Arm NN framework, the Zhouyi backend is added to support the inferences on Zhouyi NPUs. The Zhouyi backend translates and schedules Arm NN inference tasks by calling the APIs provided by the NPU build tool and driver. Currently, the Arm NN Zhouyi backend supports AIPU V2 and V3 NPUs.

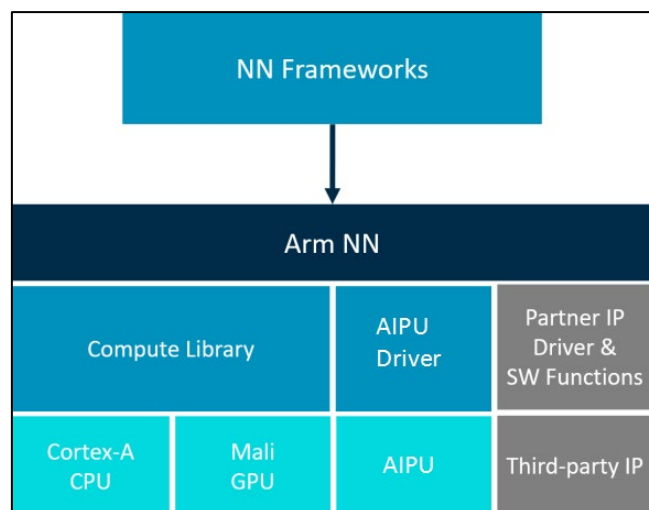


Figure 2-2 Arm NN architecture

2.2.1 Building Arm NN with the Zhouyi backend

The Arm NN version used is 23.08 with no modification. You can download the Arm NN source (and dependencies of it) from GitHub and get the Zhouyi NPU backend from the release package.

The Zhouyi NPU backend supports both x86-linux and arm64-linux environments. The x86-linux Arm NN is only suggested to be used in your software bring-up stage and it depends on the Zhouyi NPU simulators.

Generally, Arm NN with the Zhouyi backend depends on the binary modules as shown in the following table.

Table 2-1 Dependent software modules

Module	Description	Source
libaipubuildtool.so	Zhouyi build tool lib	Arm China Compass Software release package
libaipuopplugin.so	Zhouyi operator plugin lib	
libaipu_xlib.so	Zhouyi NPU operator layer lib	
aiffcllib.a & tpccllib.a	Zhouyi AIPU V3 NPU operator layer lib	
aipuoas, aipuold, libAIPU0GB.so, libaiputoolchain.so	Zhouyi compiler toolchains (for AIPUv3 only)	
Zhouyi NPU simulator binaries and libraries (for x86-linux simulation only)	-	

Module	Description	Source
aipu.ko	Zhouyi NPU KMD driver (for arm64-linux only)	Built from the released Compass driver source code
libaipudrv.so	Zhouyi NPU UMD driver	
Dependencies of Arm NN (such as TensorFlow, and FlatBuffer)	-	Downloaded as described in the Arm NN documentation

To run an application with Arm NN, you also need to build an Arm NN TFLite parser and the test applications of it. You can download them from the Arm NN GitHub repository. Note that currently only quantized TFLite models are accepted.

For detailed instructions on how to integrate the Zhouyi NPU backend into Arm NN and test the applications, see the README.md file in the Zhouyi NPU backend source directory. To simplify the integration flow, Arm China also provides an example build.sh script in the release package. You can find how to configure and use the build script in the readme file and the script itself.

2.2.2 Running the sample application

You can find three mobilenet-v2 sample applications (one of them is a multi-batch case) in AI610-SDK-1014/armnn-runtime/sample for x86-linux environments.

Follow the steps in AI610-SDK-1014/armnn-runtime/README to run the sample applications on different target platforms.

If you want to use the auto-tiling feature provided by the Zhouyi build tool, enable it by configuring the environment variable TILING_KEY to be one of the following tiling methods before executing your Arm NN application:

```
$ export TILING_KEY=fps
```

2.3 Android neural networks runtime

The *Neural Network* (NN) HAL defines an abstraction of various devices, such as *Graphics Processing Units* (GPUs) and *Digital Signal Processors* (DSPs), that are in a product (for example, a phone or a tablet). The drivers for these devices must conform to the NN HAL. As a vendor device, the Zhouyi NPU complies with the API of Android NN HAL from HIDL 1.0~1.3, and AIDL from Android 12.

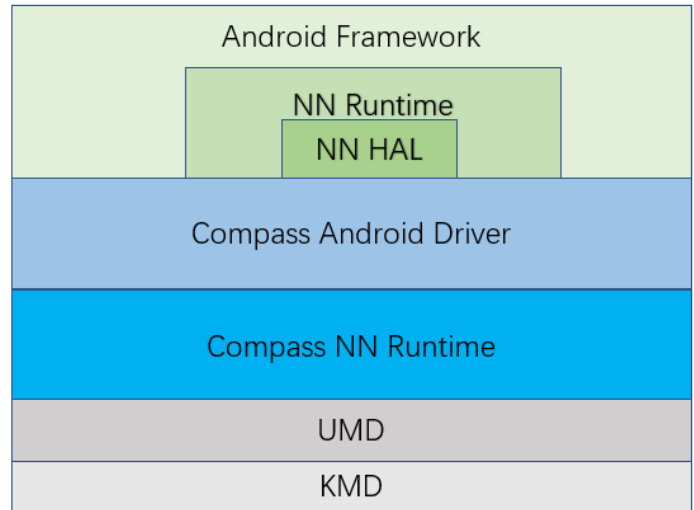


Figure 2-3 Android NN Runtime architecture

2.3.1 Using the Android Neural Networks driver

This section describes how to use the NPU driver for the Android Neural Networks API to implement the `android.hardware.neuralnetworks` HALs.

Before you begin

Check the AOSP version in P/Q/R/S/T.

Files provided

You can find the following file structure in the software package:

```

Compass_Android_Driver/nnapi/
|-- hidl
|-- aidl (from Android12)
|-- LICENSE.TXT
|-- NnapiSupport.txt (support operator list)
|-- README.md
|-- sepolicy_sample
`-- test (test apk)

Compass_NN_Runtime/
| |-- aipulib
| |-- Android.mk
| |-- build_libs_ndk.sh
| |-- CMakeLists.txt

```

```
| |-- config.mk
| |-- include
| |-- LICENSE.TXT
| |-- OperatorsSupport.txt
| |-- README.md
| |-- samples
| `-- src
Dependencies/
Simulator/ (executable files only for V2)

Prebuilt aipulib/
| |-- libaipu_buildingtool.so
| |-- v2
| | |-- libaipu_driver.so
| | `-- libaipu_layerlib.so
| `-- v3
|     |-- libaipu_driver.so
|     |-- libaipu_layerlib.so
|     |-- libaiputoolchain_core.so
|     `-- libaiputoolchain.so
```

Procedure

1. Prepare AOSP source tree <ANDROID_ROOT>.
2. Create the driver directory <AOSP_ROOT>/hardware/armchina/neuralnetworks/, and then copy <RELEASE_PACKAGE>/Compass_Android_Driver into it. If you do not attempt to build AIDL, you can delete Compass_Android_Driver/nnapi/aidl.
3. Copy <RELEASE_PACKAGE>/Compass_NN_Runtime into <AOSP_ROOT>/hardware/armchina/neuralnetworks/.
4. Update the Android build environment to add the Zhouyi driver.

Add PRODUCT_PACKAGES to vendor.mk in

<ANDROID_ROOT>/device/<manufacturer>/<product>, for example, if you try to run on the x86_64 emulator, it is <ANDROID_ROOT>/device/generic/goldfish/vendor.mk.

For HAL HIDL 1.3:

```
PRODUCT_PACKAGES += android.hardware.neuralnetworks@1.3-service-zhouyi \
    libcompass_nn_runtime \
    <prebuilt shared libraries>
```

For AIDL:

```
PRODUCT_PACKAGES += android.hardware.neuralnetworks-service-zhouyi \
    libcompass_nn_runtime \
    <prebuilt shared libraries>
```

5. Set the configuration file under <RELEASE_PACKAGE>/Compass_NN_Runtime/config.mk.

- **AIPU_VERSION:** If you set it to v2, the default target will be X1_1204.
If you set it to v3, the default target is X2_1204. If you want to use X2_1204 in a 3-core configuration, the target name is X2_1204MP3.
When you try to specify other targets such as Z3_1204, you can use the following:

```
$ adb -s <emulator_name> shell setprop vendor.nn.zhouyi.simulator_target Z3_1204
```


For hardware users, AIPU_VERSION must be consistent with the actual hardware, and cannot be modified by setprop.
 - **USE_AIPULIB_LOCAL:** Set to 1.
 - **BUILD_STANDALONE_LIB:** If you try to build libcompass_nn_runtime.so separately such as for delegate, set it to 1, otherwise set it to 0.
6. Build AOSP. For details, visit <https://source.android.com/setup/build/building>.
 7. To confirm that the Zhouyi driver has been built, check the driver service at `<ANDROID_ROOT>/out/target/product/<product>/vendor/bin/hw/android.hardware.neuralnetworks@1.3-service-zhouyi` or `android.hardware.neuralnetworks-service-zhouyi` for AIDL, and check the shared library in `<ANDROID_ROOT>/out/target/product/<product>/vendor/lib64/libcompass_nn_runtime.so` and shared libraries in the prebuilt list.

Note

For Android emulator + AIPU simulator users, the AIPUv2 simulator is an executable file, while the AIPUv3 simulator is a shared library.

When the lunched TARGET_ARCH is x86_64, the built files will run on the AIPU simulator by default. If it is arm64, they will run on hardware.

For HAL users, it is best to restart (stop|start) the service if you want to enable the new property:

- hidl: neuralnetworks_hal_service_zhouyi
- aidl: neuralnetworks_hal_service_aidl_zhouyi

2.3.2 Running the driver on the emulator

1. Launch the emulator with writable, then root, and remount with adb.
2. Push relative files into the Android emulator.
 - For the AIPUv2 architecture:


```
$ adb -s <emulator_name> shell mkdir -p /vendor/etc/npu/
$ adb -s <emulator_name> push <RELEASE_PACKAGE>/simulator /vendor/etc/npu/
$ adb -s <emulator_name> shell chmod 777 -R /vendor/etc/npu/simulator
```
 - For the AIPUv3 architecture:


```
$ adb -s <emulator_name> shell mkdir -p /vendor/etc/npu/operators/
$ adb -s <emulator_name> push <RELEASE_PACKAGE>/dependencies/*.a
/vendor/etc/npu/operators/
```

Where, <RELEASE_PACKAGE> is the package directory. <emulator_name> is the current emulator.
3. The service will check the above files at startup. If this is the first time you prepare them, you can reboot the emulator or pack these files into the Android image directly.

2.3.3 Using the profiler feature on Android

This section describes how to use the profiler feature on Android. Before you proceed, ensure that you have completed the instructions in *6.3.1 Using the Android Neural Networks driver*.

Note

You cannot enable/disable the profiler when you are running it, so you need to enable/disable the profiler before you start up the driver service. The profiler is disabled by default.

Perform the following steps to use the profiler feature:

1. Set properties before running the model.

```
$ adb -s <emulator_name> setprop vendor.nn.zhouyi.profile <on/off>
```
2. Run the model, such as `pineapple_test.apk` as described in *2.3.5 Testing*.
3. Fetch profiler files on the simulator (the simulator shares one `aiff_dumps` directory).

```
$ adb -s <emulator_name> pull /data/user/nnapi/aiff_dumps
```

```
$ adb -s <emulator_name> pull /data/user/nnapi/graph_1/graph.json
```
4. Fetch profiler files on the hardware.

```
$ adb -s <device_name> pull /data/user/nnapi/graph_<n>/graph.json
```

```
$ adb -s <device_name> pull /data/user/nnapi/graph_<n>/temp.profile
```

For information on how to analyze the profile data, see *6.2 Using the aipu_profiler* in the *Arm China Zhouyi Compass NN Compiler User Guide*.

2.3.4 Using the tiling feature on Android

1. Run the following command to enable/disable the tiling feature.

```
$ adb -s <emulator_name> shell setprop vendor.nn.zhouyi.tiling <fps/disable>
```

2.3.5 Testing

1. Install the application.

```
$adb -s <emulator_name> install -t <RELEASE_PACKAGE>/test/pineapple_test.apk
```

```
$adb -s <emulator_name> shell am start
```

```
com.example.tflite_imageclassifier/.MainActivity
```
2. Execute.
3. Check whether Android NN uses Zhouyi as the backend.

```
$ adb -s <emulator_name> shell setprop vendor.nn.zhouyi.vlog on
```
4. Click **Detect** to begin inferring.
A moment later, you can see the result on the application.
5. After inference, run `logcat` as follows:

```
$ adb -s <emulator_name> logcat
```

You should find ‘ZhouyiDriver’ and ‘EXECUTION WAS SUCCESSFUL’ in the log.

2.4 TensorFlow Lite delegate runtime

TensorFlow Lite delegates enable hardware acceleration of TensorFlow Lite models by leveraging on-device accelerators such as the NPU, GPU and *Digital Signal Processor* (DSP).

By default, TensorFlow Lite utilizes CPU kernels that are optimized for the Arm Neon instruction set. However, the CPU is a multi-purpose processor that is not necessarily optimized for the heavy arithmetic typically found in Machine Learning models (for example, the matrix math involved in convolution and dense layers). As most modern mobile phones contain chips that are better at handling these heavy operations, utilizing them for neural network operations provides huge benefits in terms of latency and power efficiency. Each of these accelerators have associated APIs that enable custom computations. TensorFlow Lite's Delegate API acts as a bridge between the TFLite runtime and these lower-level APIs.

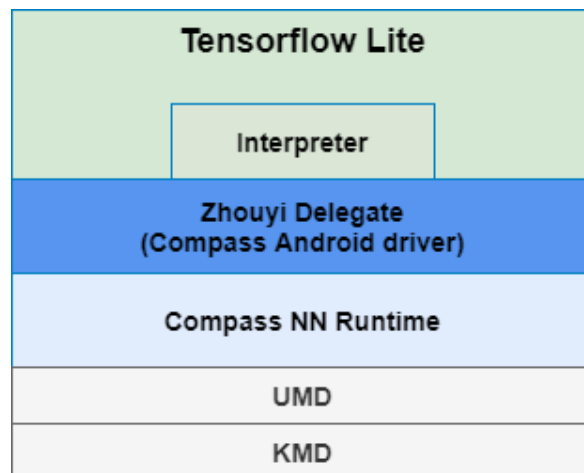


Figure 2-4 TensorFlow Lite Runtime architecture

This section describes how to use the TensorFlow Lite Zhouyi Delegate in your application using Java and/or C APIs.

The delegate leverages the Zhouyi library to execute quantized kernels on the NPU.

Currently the following devices are tested:

- Android x86_64 emulator: Runs on the NPU simulator.
- Juno: Runs on the NPU hardware.

Currently the following Zhouyi architecture is tested:

- v3

The Zhouyi delegate supports all models that conform to Arm China quantization specification, including symmetric quantization models and asymmetric quantization models.

2.4.1 How to build

- Set up the environment.
 - a. Get the TensorFlow code from GitHub and check out the branch to r2.12.
 - b. Ensure that tensorflow-lite.aar can be compiled.
 - c. Get the Zhouyi delegate code from the release package. You can find the following file structure in the software package:

```

Compass_Android_Driver/zhouyi
|-- BUILD
|-- builders
|-- compass_nn_runtime

```

```
|-- docs
|-- java
|-- utils.cc
|-- utils.h
|-- version_script.lds
|-- zhouyi_delegate.cc
|-- zhouyi_delegate.h
|-- zhouyi_delegate_kernel.cc
|-- zhouyi_delegate_kernel.h
|-- zhouyi_implementation.cc
|-- zhouyi_implementation.h
```

- d. Put the zhouyi folder in the following path:

tensorflow/lite/delegates/.

- Build the AAR package.

After Bazel is properly configured, you can build the Zhouyi delegate AAR from the root checkout directory as follows:

```
1. bazel build -c opt --config=<arch>
   tensorflow/lite/delegates/zhouyi/java:tensorflow-lite-zhouyi
```

This will generate an AAR file in bazel-bin/tensorflow/lite/delegates/zhouyi/java/.

- Build the delegate.

```
1. bazel build -c opt --config=<arch>
   tensorflow/lite/delegates/zhouyi:zhouyi_delegate
```

This will generate a libzhouyi_delegate.so file in bazel-bin/tensorflow/lite/delegates/zhouyi.

- Build tests.

```
1. bazel --bazelrc=/dev/null build -c opt --config=<arch>
   tensorflow/lite/delegates/zhouyi/builders/tests:all
```

This will generate the *zhouyi_xxx* test file in bazel-bin/tensorflow/lite/delegates/zhouyi/builders/tests.

- Implement the architecture support:

- android_x86_64
- android_arm64

2.4.2 Zhouyi delegate Java API

```
2. public class ZhouyiDelegate implements Delegate, Closeable {
3.     /**
4.      * Zhouyi Delegate DMA buffer
5.      * Use this object to access zhouyi's dma buffer.
6.      */
7.     public static final class DMABuffer implements Closeable {
8.         /**
9.          * init DMABuffer
10.        */
11.        public DMABuffer(long size)
12.
13.        /**
```

```

14.     * get buffByteBuffer
15.     */
16.     public ByteBuffer bytebuffer()
17.
18.     /**
19.     * get buffer size
20.     */
21.     public long size
22.
23.     /**
24.     * Bind buffer handle to tensor
25.     */
26.     public long bind(long delegateHandle, int tensor_index)
27.
28.     /**
29.     * Explicit call to release resources.
30.     */
31.     @Override
32.     public void close()
33. }
34.
35. /** Delegate options. */
36. public static final class Options {
37.     /**
38.     * Enable or disable the Zhouyi runtime dump file.
39.     * default settings: false.
40.     */
41.     public Options setEnableDump(boolean enable)
42.
43.     /**
44.     * Enable or disable the Zhouyi runtime profile.
45.     * It generates performance files
46.     * default settings: false.
47.     */
48.     public Options setEnableProfile(boolean enable)
49.
50.     /**
51.     * Set the Zhouyi runtime tiling method, it support "fps" "footprint".
52.     * Titing is a model segmentation function.
53.     * default settings: none.
54.     */
55.     public Options setTiling(String tilingMethod)
56.
57.     /**
58.     * Configure the location to be used to store model compilation cache
59.     entries.
60.     * If not set, the default is "/data/tata/{app package name}"
61.     */
62.     public Options setCacheDir(String cacheDir)
63.
64.     /**
65.     * This corresponds to the log level in the Zhouyi Runtime.
66.     * 0(defalut): LOG_ERROR
67.     * 1: LOG_WARN
68.     * 2: LOG_INFO
69.     * 3: LOG_DEBUG
70.     * 4: LOG_VERBOSE
71.     */
72.     public Options setLogLevel(int level)
73.
74.     /**
75.     * Returns whether the Zhouyi runtime "dump file" function is enabled.

```



```

75.     */
76.     public Boolean getEnableDump()
77. }
78.
79. /*
80.  * Creates a new ZhouyiDelegate object without param.
81.  * Throws UnsupportedOperationException if Zhouyi AIPU delegation is not
    available on this device.
82.  */
83.     public ZhouyiDelegate() throws UnsupportedOperationException
84.
85.     /*
86.     * Creates a new ZhouyiDelegate object given the current 'options'.
87.     * Throws UnsupportedOperationException if Zhouyi AIPU delegation is not
    available
88.     * on this device.
89.     */
90.     public ZhouyiDelegate(Options setOptions) throws
    UnsupportedOperationException
91.
92.     /**
93.     * Frees TFLite resources in C runtime.
94.     * User is expected to call this method explicitly.
95.     */
96.     @Override
97.     public void close();
98. }

```

Example usage

1. Add the AAR file to the application.
2. Add tensorflow-lite-zhouyi-xxx.aar in the release package to the app/libs/ directory of the application.
3. Edit app/build.gradle to use the Zhouyi delegate AAR.

```

1. dependencies {
2.     ...
3.     implementation 'org.tensorflow:tensorflow-lite:xx.xx.xx'
4.     implementation files('libs/tensorflow-lite-zhouyi-xxx.aar')
5. }

```

Note

The AAR file can also be added using the Android Studio tool.

4. Add Zhouyi libraries to your Android app. See *2.4.5 Adding shared libraries to an app*.

Note

The Zhouyi shared library must be used with the AAR.

5. Create a delegate and initialize a TensorFlow Lite Interpreter.

```

1. import org.tensorflow.lite.ZhouyiDelegate;
2. private ZhouyiDelegate zhouyidelegate = null;
3. private Interpreter tfLiteInterpreter;
4. private Interpreter.Options tfLiteOptions = new Interpreter.Options();

```

```

5. // Create the Delegate instance.
6. try {
7.     ZhouyiDelegate.Options zhouyiDelegateOptions = new
        ZhouyiDelegate.Options();
8.     zhouyiDelegateOptions.setLogLevel(4);
9.     zhouyiDelegateOptions.setEnableProfile(false);
10.    zhouyiDelegateOptions.setEnableDump(false);
11.    zhouyiDelegateOptions.setTiling("fps");
12.    if(zhouyidelegate == null) {
13.        zhouyidelegate = new ZhouyiDelegate(zhouyiDelegateOptions);
14.    }
15.    tfliteOptions.addDelegate(zhouyidelegate);
16. } catch (UnsupportedOperationException e) {
17.     // Zhouyi delegate is not supported on this device.
18. }
19.
20. tfliteInterpreter = new Interpreter(tfliteModel, tfliteOptions);
21.
22. //...
23. //work
24. //...
25.
26.
27. // Dispose after finished with inference.
28. tfliteInterpreter.close();
29. if (zhouyiDelegate != null) {
30.     zhouyiDelegate.close();
31. }

```

2.4.3 Zhouyi delegate C API

```

1. // Use TfliteZhouyiDelegateOptionsDefault() for Default options.
2. struct TFL_CAPI_EXPORT TfliteZhouyiDelegateOptions {
3.     // This determines whether Zhouyi runtime dumps intermediate files.
4.     // false(defalut)
5.     bool enableDump;
6.
7.     // This corresponds to the opening of the Zhouyi runtime profile function.
8.     // false(defalut)
9.     bool enableProfile;
10.
11.    // This corresponds to the log level in the Zhouyi Runtime.
12.    // 0(defalut): LOG_ERROR
13.    // 1: LOG_WARN
14.    // 2: LOG_INFO
15.    // 3: LOG_DEBUG
16.    // 4: LOG_VERBOSE
17.    int logLevel;
18.
19.    // This corresponds to the opening of the Zhouyi runtime tiling function.
20.    // Only support "footprint", "fps".
21.    char tilingMethod[512];
22.
23.    // Current,the compilation process needs to generate intermediate files,
24.    // it is necessary to specify a directory to store intermediate files.
25.    // This directory usually is "/data/data/{package name}".
26.    char cacheDir[512];
27.
28.    // This sets the maximum number of Zhouyi graphs created with
29.    // zhouyi_nn_init. Each graph corresponds to one delegated node subset in
    the

```

```

30. // TFLite model.
31. int max_delegated_partitions;
32. // This sets the minimum number of nodes per graph created with
33. // zhouyi_nn_init. Defaults to 2.
34. int min_nodes_per_partition;
35. };
36.
37. // Returns TfLiteZhouyiDelegateOptions populated with default values.
38. TFL_CAPI_EXPORT TfLiteZhouyiDelegateOptions
39. TfLiteZhouyiDelegateOptionsDefault();
40.
41. // Same as above method but doesn't accept the path params.
42. // Return a delegate that uses Zhouyi SDK for ops execution.
43. // Must outlive the interpreter.
44. TfLiteDelegate* TFL_CAPI_EXPORT
45. TfLiteZhouyiDelegateCreate(const TfLiteZhouyiDelegateOptions* options);
46.
47. // Do any needed cleanup and delete 'delegate'.
48. void TFL_CAPI_EXPORT TfLiteZhouyiDelegateDelete(TfLiteDelegate* delegate);
49.
50. // Assumes the environment setup is already done. Only initialize Zhouyi.
51. void TFL_CAPI_EXPORT TfLiteZhouyiInit();
52.
53. // Clean up and switch off the AIPU.
54. // This should be called after all processing is done and delegate is
   deleted.
55. void TFL_CAPI_EXPORT TfLiteZhouyiTearDown();
56.
57. //Allocate zhouyi dma buffer.
58. int TFL_CAPI_EXPORT
59. TfLiteZhouyiAllocatedDMABuffer(int size);
60.
61. //Get buffer addr.
62. void* TFL_CAPI_EXPORT
63. TfLiteZhouyiGetAddr(int buffer_handle);
64.
65. //Free zhouyi dma buffer by handle.
66. void TFL_CAPI_EXPORT
67. TfLiteZhouyiFreeDMABuffer(int buffer_handle);
68.
69. //Free zhouyi dma buffer by pointer.
70. void TFL_CAPI_EXPORT
71. TfLiteZhouyiFreeDMABufferByAddr(void* buffer);
72.
73. // Bind DMA buffer to I/O tensor.
74. // Note that this function must be called after
   Interpreter::ModifyGraphWithDelegate().
75. TfLiteStatus TFL_CAPI_EXPORT
76. TfLiteZhouyiBindBufferToTensor(
77.     TfLiteDelegate* delegate, int buffer_handle, int tensor_index);
78.
79. // Bind DMA buffer to I/O tensor by tensor pointer.
80. // Note that this function must be called after
   Interpreter::ModifyGraphWithDelegate().
81. TfLiteStatus TFL_CAPI_EXPORT
82. TfLiteZhouyiBindBufferToTensorP(
83.     TfLiteDelegate* delegate, int buffer_handle, TfLiteTensor* tensor);
84.

```

Example usage

1. Add the AAR file to the application.

2. Add tensorflow-lite-zhouyi-xxx.aar in the release package to the app/libs/ directory of the application.
3. Edit app/build.gradle to use the Zhouyi delegate AAR.

```
6. dependencies {
7.     ...
8.     implementation 'org.tensorflow:tensorflow-lite:xx.xx.xx'
9.     implementation files('libs/tensorflow-lite-zhouyi-xxx.aar')
10. }
```

Note

The AAR file can also be added using the Android Studio tool.

4. Add Zhouyi libraries to your Android app. See *2.4.5 Adding shared libraries to an app*.

Note

The Zhouyi shared library must be used with the AAR.

5. Include the C header.

The header file zhouyi_delegate.h can be extracted from the Zhouyi delegate AAR.

6. Create a delegate and initialize a TensorFlow Lite Interpreter.

In your code, ensure that the native Zhouyi library is loaded. This can be done by calling in your Activity or Java entry-point.

```
1. System.loadLibrary("tensorflowlite_zhouyi_jni");
```

Example of creating a delegate:

```
1. #include "tensorflow/lite/delegates/zhouyi/zhouyi_delegate.h"
2.
3. // If files are packaged with native lib in android App then it
4. // will typically be equivalent to the path provided by
5. // "getContext().getApplicationInfo().nativeLibraryDir"
6. const char[] library_directory_path = "/data/data/{package name}";
7. TfLiteZhouyiInit();
8.
9. ::tflite::TfLiteZhouyiDelegateOptions options = {
10.     .cacheDir = library_directory_path,
11. };
12. // 'delegate_ptr' Need to outlive the interpreter. For example,
13. // If use case will need to resize input or anything that can trigger
14. // re-applying delegates then 'delegate_ptr' need to outlive the interpreter.
15. auto* delegate_ptr = ::tflite::TfLiteZhouyiDelegateCreate(&options);
16. Interpreter::TfLiteDelegatePtr delegate(delegate_ptr,
17.     [](TfLiteDelegate* delegate) {
18.         ::tflite::TfLiteZhouyiDelegateDelete(delegate);
19.     });
20. interpreter->ModifyGraphWithDelegate(delegate.get());
21. // After usage of delegate.
22. TfLiteZhouyiTearDown(); // Needed once at end of app/Zhouyi usage.
```

2.4.4 Supported ops

Refer to the Supported_Ops.md file in the code package.

2.4.5 Adding shared libraries to an app

Generally, the Zhouyi delegate depends on the binary modules as shown in the following table.

Table 2-2 Delegate dependencies

Module	Description	Source
libaipu_buildingtool.so	Zhouyi build tool lib	Arm China Compass Software release package
libaiputoolchain.so	Zhouyi operator plugin lib	
libaiputoolchain_core.so	Zhouyi NPU operator layer lib	
libaipugenerate_afile.so	Zhouyi AIPU V3 NPU operator layer lib	
libaipu_layerlib.so		
libaipu_simulator_x2.so	For Android x86_64 emulator only	
libcompass_nn_runtime.so	Compass NN runtime lib	Arm China Compass Software release package or built from the released Compass NN runtime source code
libaipu_driver.so	Zhouyi NPU UMD driver	Arm China Compass Software release package or built from the released Compass driver source code
aipu.ko	Zhouyi NPU KMD driver (for arm64 only)	Built from the released Compass driver source code

The release package provides libraries for the following two target architectures:

- x86_64
- arm64-v8a

To add shared libraries to an application

1. Find the file of the corresponding architecture, and add it to the app/libs/{arch} path.
2. Modify build.gradle of the application to add the following code:

```

1. sourceSets {
2.     main {
3.         jniLibs.srcDirs = ['libs']
4.     }
5. }
```

2.4.6 Input/Output buffers

The Zhouyi delegate provides a set of C++/Java APIs to directly read from and write to the NPU hardware buffer, bypassing avoidable memory copies.

C++ API usage

The following techniques are only available when you are using Bazel or building TensorFlow Lite yourself.

- The function must be called after `Interpreter::ModifyGraphWithDelegate()`. In addition, the inference output is copied from NPU memory to CPU memory by default. You can turn this behavior off by calling `Interpreter::SetAllowBufferHandleOutput(true)` during initialization.

- If the input of the network is already loaded in the NPU memory, the input buffer handle can be directly bound to the input tensor.
- The output data of the network can also be used directly as the input of the next network.
- The buffer data is released when the subgraph is destroyed, and it also can be released by calling interface `TfLiteZhouyiFreeDMABuffer`.

C++ API example usage

```

1. #include "tensorflow/lite/delegates/zhouyi/zhouyi_delegate.h"
2. // ...
3. // Prepare NPU delegate.
4.
5. // If files are packaged with native lib in android App then it
6. // will typically be equivalent to the path provided by
7. // "getContext().getApplicationInfo().nativeLibraryDir"
8. const char[] library_directory_path = "/data/local/tmp";
9. TfLiteZhouyiInit();
10.
11. ::tflite::TfLiteZhouyiDelegateOptions options = {
12.     .cacheDir = library_directory_path,
13. };
14. auto* delegate_ptr = ::tflite::TfLiteZhouyiDelegateCreate(&options);
15. Interpreter::TfLiteDelegatePtr delegate(delegate_ptr,
16.     [](TfLiteDelegate* delegate) {
17.         ::tflite::TfLiteZhouyiDelegateDelete(delegate);
18.     });
19. interpreter->ModifyGraphWithDelegate(delegate.get());
20.
21. interpreter->SetAllowBufferHandleOutput(true); // disable default npu->cpu
    copy
22. int user_provided_input_buffer = TfLiteZhouyiAllocatedDMABuffer(input_size);
23. int user_provided_output_buffer =
    TfLiteZhouyiAllocatedDMABuffer(output_size);
24. if (!TfLiteZhouyiBindBufferToTensor(
25.     delegate, interpreter->inputs()[0], user_provided_input_buffer)) {
26.     return false;
27. }
28. if (!TfLiteZhouyiBindBufferToTensor(
29.     delegate, interpreter->outputs()[0], user_provided_output_buffer)) {
30.     return false;
31. }
32.
33. // Run inference.
34. if (interpreter->Invoke() != kTfLiteOk) return false;

```

Java API usage

By default, the inference output is copied from NPU memory to CPU memory. You can turn this behavior off by calling `Interpreter.Options.setAllowBufferHandleOutput(true)` during initialization.

JAVA API example usage

```

35. // interpreter prepare
36. // delegate prepare
37.
38. // ...
39.
40. // input
41. ZhouyiDelegate.DMABuffer dmaBufferIn = new
    ZhouyiDelegate.DMABuffer(input_size);

```

```

42. ByteBuffer bytebuffer = dmaBufferIn.bytebuffer();
43. dmaBufferIn.bind(delegate.getNativeHandle(),
    interpreter.getInputTensor(0).index());
44. // copy input data to bytebuffer
45.
46.
47. // output
48. ZhouyiDelegate.DMABuffer dmaBufferOut = new
    ZhouyiDelegate.DMABuffer(output_size);
49. dmaBufferOut.bind(delegate.getNativeHandle(),
    interpreter.getOutputTensor(0).index());
50. // ...
51.
52. // the output data:
53. ByteBuffer output_bytebuffer = dmaBufferOut.bytebuffer();
54.
55. // after invoke
56. dmaBufferIn.close();
57. dmaBufferOut.close();

```

2.4.7 Testing

- OP test
 1. Push shared libraries to device path /vendor/lib64/.
 2. Push *zhouyi_xxx_test* to device path /data/local/tmp.
 3. Run *zhouyi_xxx_test*.
For the compilation instructions of *zhouyi_xxx_test*, see 2.4.2 *Zhouyi delegate Java API*.
- Sample app test
 1. Get the `delegateJavaSample` code from the release package.
 2. Build the app and install it onto the device.
 3. Run the app. If you use the Juno device, install the driver first, that is, `aipu.ko`.

2.4.8 FAQs

The following information provides answers to general frequently asked questions.

How can I tell that the model is using Zhouyi when I enable the delegate?

Two log messages will be displayed when you enable the delegate:

- One is to indicate whether the delegate was created.
- The other is to indicate how many nodes are running using the delegate.

```

1. Created TensorFlow Lite delegate for Zhouyi.
2. Zhouyi delegate: X nodes delegated out of Y nodes.

```

Do all Ops in the model need to be supported to run the delegate?

No. The model will be partitioned into subgraphs based on the supported ops. Any unsupported ops will run on the CPU.

System permission: I have rooted my phone but still cannot run the delegate, what should I do?

Be sure to disable SELinux enforcing by running `adb shell setenforce 0`.