



**MLU370-M8 Intelligent Accelerating
Card Product Manual**

Issue 0.9.1

Cambricon

2021.06.22

1. Preface.....	3
1.1. Copyright Notice.....	3
1.2. Version record.....	4
1.3. Update history.....	4
2. Overview.....	5
3. Product specification overview.....	6
3.1 Overview of the product specifications and parameters.....	6
3.2 Overview of structural specifications.....	6
3.3 Overview of power specifications.....	7
3.4 Overview of heat dissipation specifications.....	7
3.5 Interface specification overview.....	7
4. Electrical specifications.....	9
4.1 Connector pin description.....	9
4.2 Power demand.....	13
4.2.1 Enter the power supply.....	13
4.2.2 Peak power supply current.....	13
4.2.3 HSC protection circuit.....	13
4.2.4 Power timing.....	14
4.3 Signal description.....	14
4.3.1 Clock signal.....	14
4.3.2 PCIe signal.....	15
4.3.3 The MLU-Link signal.....	15

4.3.4 Other signal instructions.....	16
5. Heat dissipation specifications.....	24
5.1 Radiation description.....	24
5.2 Speed reduction protection temperature.....	24
5.3 Turn off temperature.....	24
5.4 Working environment air volume requirements.....	25
5.5 Thermal simulation model.....	26
6. The Cambricon NeuWare development environment.....	27
7. Compliance	28



1. Preface

1.1. Copyright Notice

DISCLAIMER

CAMBRICON MAKES NO REPRESENTATION, WARRANTY (EXPRESS, IMPLIED, OR STATUTORY) OR GUARANTEE REGARDING THE INFORMATION CONTAINED HEREIN, AND EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES OF MERCHANTABILITY, TITLE, NONINFRINGEMENT OF INTELLECTUAL PROPERTY OR FITNESS FOR A PARTICULAR PURPOSE, AND CAMBRICON DOES NOT ASSUME ANY LIABILITY ARISING OUT OF THE APPLICATION OR USE OF ANY PRODUCT OR SERVICES. CAMBRICON SHALL HAVE NO LIABILITY RELATED TO ANY DEFAULTS, DAMAGES, COSTS OR PROBLEMS WHICH MAY BE BASED ON OR ATTRIBUTABLE TO: (I) THE USE OF THE CAMBRICON PRODUCT IN ANY MANNER THAT IS CONTRARY TO THIS Manual, OR (II) CUSTOMER PRODUCT DESIGNS.

LIMITATION OF LIABILITY

In no event shall Cambricon be liable for any damages whatsoever (including, without limitation, damages for loss of profits, business interruption and loss of information) arising out of the use of or inability to use this Manual, even if Cambricon has been advised of the possibility of such damages. Notwithstanding any damages that customer might incur for any reason whatsoever, Cambricon's aggregate and cumulative liability towards customer for the product described in this Manual shall be limited in accordance with the Cambricon terms and conditions of sale for the product.

ACCURACY OF INFORMATION

Information provided in this document is proprietary to Cambricon, and Cambricon reserves the right to make any changes to the information in this document or to any products and services at any time without notice. The information contained in this Manual and all other information contained in Cambricon documentation referenced in this Manual is provided "AS IS." Cambricon does not warrant the accuracy or completeness of the information, text, graphics, links or other items contained within this Manual. Cambricon may make changes to this Manual, or to the products described therein, at any time without notice, but makes no commitment to update this Manual.

Performance tests and ratings set forth in this Manual are measured using specific chips or computer systems or components. The results shown in this Manual reflect approximate performance of Cambricon products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance. As set forth above, Cambricon makes no representation, warranty or guarantee that the product described in this Manual will be suitable for any specified use. Cambricon does not represent or warrant that it tests all parameters of each product. It is customer's sole responsibility to ensure that the product is suitable and fit for the application planned by the customer and to do the necessary testing for the application in order to avoid a default of the application or the product.

Weaknesses in customer's product designs may affect the quality and reliability of Cambricon product and may result in additional or different conditions and/or requirements beyond those

contained in this Manual.

IP NOTICES

Cambricon and the Cambricon logo are trademarks and/or registered trademarks of Cambricon Corporation in China and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

This Manual is copyrighted and is protected by worldwide copyright laws and treaty provisions. This Manual may not be copied, reproduced, modified, published, uploaded, posted, transmitted, or distributed in any way, without Cambricon's prior written permission. Other than the right for customer to use the information in this Manual with the product, no other right or license, either express or implied, is granted by Cambricon under this Manual. For the avoidance of doubt, Cambricon does not grant any right or license (express or implied) to customer under any patents, copyrights, trademarks, trade secret or any other intellectual property or proprietary rights of Cambricon.

COPYRIGHT NOTICE

© Cambricon Corporation. All rights reserved.

1.2. Version record

Table 1.1 Version record

Document name	MLU370-M8 Intelligent Accelerating Card Product Manual
Version number	V0.9.1
The author	Cambricon
Create date	2021.06.22

1.3. Update history

V0.9.0

Update time: 6.06.17,2021

Update content:

- Initial version

V0.9.1

Update time: 26.22.2021

Update content:

- Modified Section 4.3.4.1: Register read operation process



Figure 21 The Cambricon MLU370-M8 Intelligent Accelerating Card

Fully upgraded data center training and push integrated AI acceleration card

The MLU370-M8 intelligent accelerating card is based on the Cambricon new generation MLU370 chip. The interface for PCIe 4.0 X16, is the OPEN Accelerator Module (OAM) standard size acceleration card, which is suitable for the latest CPU platform in the industry and can quickly realize the deployment of AI computing power. The MLU370-M8 acceleration card power consumption of 300W, can provide powerful computing support for highly diverse AI applications such as computer vision, natural language processing, and voice.

The new MLU-Link™ and ROCE v2, flexibly form training clusters

Cambricon MLU-Link™ multi-core interconnection technology supports inter-chip interconnection and inter-system interconnection, which can realize computing center-scale vertical expansion and meet the needs of ultra-large AI model training. MLU 370-M 8 supports the total bandwidth of 200G B/s, training cluster without switch dependence, or ROCEv2 network with 1 * 100 Gbps bandwidth, and MLU-Link™. The hybrid network with ROCE v2 realizes the large-scale expansion of the training cluster.



3. Product specification overview

3.1 Overview of the product specifications and parameters

The MLU370-M8 intelligent accelerating card specification parameters are listed in the following table:

Table 3.1 MLU370-M8 specification parameters

Specification indicators	Description
Plate card type	MLU370-M8
Core architecture	Cambricon MLUv03
Frequency	1.3GHz
Video decoding	Support
Memory capacity	48GB
Memory bit wide	384-bit
Memory bandwidth	307.2 GB/s
System interface	PCIe 4.0 x16, supports lane reversal
The PCI identifier	PCIe Vendor ID 0xCABC PCIe Device ID 0x0370 PCIe Sub-Vendor ID 0xCABC PCIe Sub-System ID 0x0055
The MLU-L ink interface	4 Ports
Total MLU-L ink bandwidth	200GB /s BI-Direction
TDP power consumption	300W
Radiation scheme	Passive

3.2 Overview of structural specifications

The MLU370-M8 intelligent accelerating card structure specifications are listed in the following table:

Table 3.2 Structural specifications of MLU370-M8

Specification indicators	Description
Board card, appearance	The 102mm*165mm*108.24mm,OAM specification
Board card weight	1.41KG
Minimum chip pressure	30PSI
Maximum chip pressure	60PSI
The Partner Label Area	16.1mm*40.2mm

3.3 Overview of power specifications

MLU370-M8 intelligent accelerating card power specifications are in the following table:

Table 3.3 MLU370-M8 power supply specifications

Specification indicators	Description
Enter voltage	54V±5% , 5.56A
Input power peak (EDPp)	1.6X TDP ≤2ms
	1.5X TDP ≤5ms
	1.2X TDP ≤10ms
	1.1X TDP ≤20ms

3.4 Overview of heat dissipation specifications

The MLU370-M8 intelligent accelerating card cooling specifications are listed in the following table:

Table 3.4 MLU370-M8 cooling specifications

Specification indicators	Description
MLU Down Operating Temperature (Tj)	92°C
MLU off temperature (Tj)	95°C
MLU power reduction ratio	From 1 / 2 to 1 / 8

3.5 Interface specification overview

The MLU370-M8 intelligent accelerating card interface specifications are listed in the following table:

Table 3.5 MLU370-M8 interface specifications

Interface	Description
PCIe Base address	PF (1 individual, 64bit): BAR0: 256MB prefetchable BAR2: 256MB prefetchable BAR4: 256MB prefetchable VF (4s, 64bit): BAR0: 256MB prefetchable BAR2: 256MB prefetchable BAR4: 256MB prefetchable
SMBus (8bit address)	0x8E(Write) 0x8F (Read)



4. Electrical specifications

4.1 Connector pin description

The MLU370-M8 intelligent accelerating card uses two 688pin Molex Mirror Mezz buckboard connectors. The connector is designed to impedance $90 \Omega \pm 5\%$ and is compatible with both protocols of 85Ω and 100Ω protocols. The signal and power supply are connected to the motherboard through the connector, and the connector pin area is divided as follows:

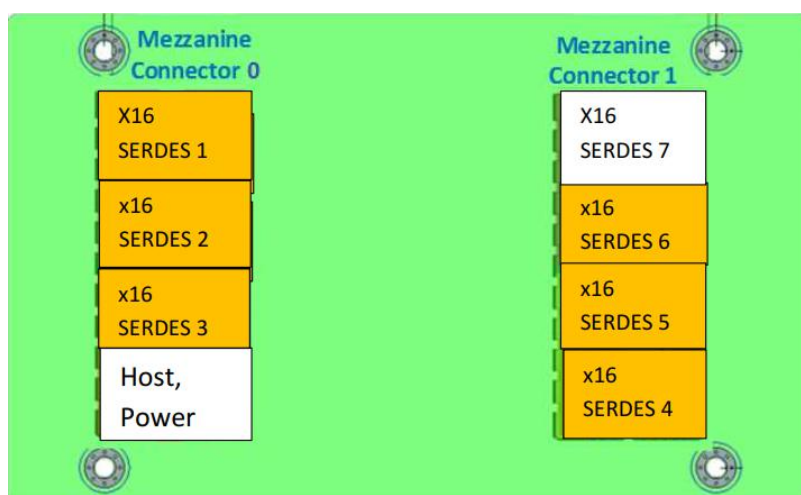


Figure 4.1 Division of the connector pin region

Connector pin arrangement is shown in the following table:

Table 4.1 MLU370-M8 Connector 0 pin arrangement

Signal	Signal direction	Signal description	Voltage
PWR_54V	I	MLU370-M8 module power input pin, supporting $54V \pm 5\%$ power supply	54V
PVREF[1:0]	O	MLU370-M8 module JTAG interface, I/O signal voltage indication signal, 1.8V	1.8V
PETp/n [15:0]	O	The PCIe sends a signal. The MLU370-M8 module is sent, and the motherboard receives it. Note the	/

		AC coupling capacitance near the buckle board connector of the motherboard, recommended 220nF	
PERp/n [15:0]	I	The PCIe receives the signal.MLU370-M8 module receives, motherboard send.Note placing the AC coupling capacitance near the motherboard sending chip with a recommended value of 220nF	/
PE_REFCLKp/n	I	The PCIe 100MHz reference clock	/
PERST#	I	MLU370-M8 module reset signal, low level effective	3.3V
HOST_PWRGD	I	Main board power supply GOOD indicates the signal	3.3V
MODULE_PWRGD	O	The MLU370-M8 module power supply GOOD indicates the signal	3.3V
PWRBRK#	I	Power brake power consumption reduces frequency to 1 / 4 of current frequency, low level effective	3.3V
MODULE_ID[4:0]	I	MLU370-M8 module slot bit ID.MLU370-M8 module internal default 10K pull-up	3.3V
I2C_SLV_D	I/O	I2C data signal, the MLU370-M8 module works in slave mode	3.3V
I2C_SLV_CLK	I	I2C clock signal, the MLU370-M8 module works in slave mode	3.3V
I2C_SLV_ALERT#	O	Reserved, NC processing	3.3V
UART_TXD	O	MLU370-M8 module MCU UART serial port output	3.3V
UART_RXD	I	MLU370-M8 module MCU UART serial port input	3.3V
JTAG0_TRST	I	The MLU370-M8 module JTAG0 TRST reset signal	1.8V
JTAG0_TMS	I	The MLU370-M8 module JTAG0 TMS mode selection signal	1.8V
JTAG0_TCK	I	The MLU370-M8 module JTAG0 TCK clock signal	1.8V
JTAG0_TDO	O	The MLU370-M8 module JTAG0 TDO data output signal	1.8V
JTAG0_TDI	I	The MLU370-M8 module JTAG0 TDI data input signal	1.8V
PRSNT0#	O	MLU370-M8 module buckle connector 0 position signal, and MLU370-M8 module default 1 pull down K.It is recommended to pull 10K	1.8V OR 3.3V
MANF_MODE#	I	No feature defined, NC processing	3.3V
FW_RECOVERY#	I	No feature defined, NC processing	3.3V
TEST_MODE#	I	Test mode.NC processing is available	1.8V OR 3.3V
RFU	/	Keep the pipe feet in advance	

Table 4.2 MLU370-M8 Connector 1 pin arrangement

Signal	Signal direction	Signal description	Voltage
MLU-Link _0Tp/n [7:0]	O	The MLU-L ink0[7:0] sends signals and defines the pin SERDES_4Tp/n [15:8] corresponding to the OAM specification	/
MLU-Link _0Rp/n [7:0]	I	The MLU-L ink0[7:0] receives the signals and defines the pin SERDES_4R p/n [15:8] corresponding to the OAM specification	/
MLU-Link _1Tp/n [7:0]	O	The MLU-L ink1[7:0] sends signals and defines the pin SERDES_4Tp/n [7:0] corresponding to the OAM specification	/
MLU-Link _1Rp/n [7:0]	I	The MLU-L ink1[7:0] receives the signals and defines the pin SERDES_4R p/n [7:0] corresponding to the OAM specification	/
MLU-Link _2Tp/n [7:0]	O	The MLU-L ink2[7:0] sends signals and defines the pin SERDES_5Tp/n [7:0] corresponding to the OAM specification	/
MLU-Link _2Rp/n [7:0]	I	The MLU-L ink2[7:0] receives the signals and defines the pin SERDES_5R p/n [7:0] corresponding to the OAM specification	/
MLU-Link _3Tp/n [7:0]	O	The MLU-L ink3[7:0] sends signals and defines the pin SERDES_7Tp/n [15:8] corresponding to the OAM specification	/
MLU-Link _3Rp/n [7:0]	I	The MLU-L ink3[7:0] receives the signals and defines the pin SERDES_7R p/n [15:8] corresponding to the OAM specification	/
MLU-Link _4Tp/n [7:0]	O	The MLU-L ink4[7:0] sends signals and defines the pin SERDES_6Tp/n [15:8] corresponding to the OAM specification	/
MLU-Link _4Rp/n [7:0]	I	The MLU-L ink4[7:0] receives the signals and defines the pin SERDES_6R p/n [15:8] corresponding to the OAM specification	/
MLU-Link _5Tp/n [7:0]	O	The MLU-L ink5[7:0] sends signals and defines the pin SERDES_6Tp/n [7:0] corresponding to the OAM specification	/
MLU-Link _5Rp/n [7:0]	I	The MLU-L ink5[7:0] receives the signals and defines the pin SERDES_6R p/n [7:0] corresponding to the OAM specification	/
AUX_156M_REFCLKp/n	I	The MLU-Link 156.25MHz reference clock	/

PWRRDT#[1:0]	I	TDP power consumption set tube feet, the motherboard needs to provide the default 3.3V pull-up 11-L0 level, Normal TDP power consumption is 300W, default value The 10-L1 level, TDP power consumption was reduced to 250W The 01-L2 level, TDP power consumption was reduced to 200W The 00-L3 level, TDP power consumption was reduced to 150W	3.3V
THERMTRIP#	O	The MLU370-M8 module overtemperature alarm will trigger the MLU370-M8 module to drop power automatically. Please check the chassis fault (such as a fan fault) and then restart the equipment recovery. The low level is valid	3.3V
LINK_CONFIG[4:0]	I	serdes link configuration topology, default 10 K pull-up inside the MLU370-M8 module	3.3V
PE_BIF[1:0]	O	PCI e interface bit width indication: the MLU370-M8 module default is 00 00-1 x16(defaults) 01-Reserve 10-Reserve 11-Reserve	1.8V OR 3.3V
PLINK_CAP	O	"PCIe" port protocol support: MLU370-M8 module default is 0 0 = supports the PCI e protocol only (default) 1 = Other protocols are supported, reserved	1.8V OR 3.3V
PRSNT1#	O	MLU370-M8 module buckle connector 1 position signal, and MLU370-M8 module default 1 pull down K. It is recommended to pull 10K	1.8V OR 3.3V
SCALE_DEBUG_EN	O	The function is not defined. NC processing is available	3.3V
DEBUG_PORT_PRSNT#	I	The commissioning signal of the main board. NC processing is available	1.8V
RFU	/	Keep the pipe feet in advance	

Connector signals Pinlist and Pinmap are shown in MLU370-M8 Connector Pinmap » .

4.2 Power demand

4.2.1 Enter the power supply

MLU370-M8 intelligent accelerating card input power requirements in the following table:

Table 4.3 Input power specification for MLU370-M8

Enter voltage	Enter the current
54V±5%	5.56A

Note:

- 1、 The voltage value shall be the test value at the connector;
- 2、 If the input voltage is low, the current value needs to be raised to meet the 300W TDP specification;

4.2.2 Peak power supply current

The MLU370-M8 Smart Accelerator enables power reduction regulation of instantaneous power changes above μ s levels, and power regulators can support power fluctuations within the ms level (e.g., 1.2x TDP).

Table 4.4 MLU370-M8 EDPp Specifications

EDP	The duration
1.6X TDP	≤ 2 m s
1.5X TDP	≤ 5 m s
1.2X TDP	≤ 10 m s
1.1X TDP	≤ 20 m s

4.2.3 HSC protection circuit

MLU370-M8 intelligent accelerating card input voltage is 54V, motherboard needs to design Hot Swap Controllers (HSC) circuit to provide slow opening, short circuit protection, overflow voltage protection for the MLU370-M8 module, each MLU370-M8 module recommended separate 1 H SC circuit, requiring its transient power consumption support of more than 2x TDP, HSC circuit box diagram is as follows:

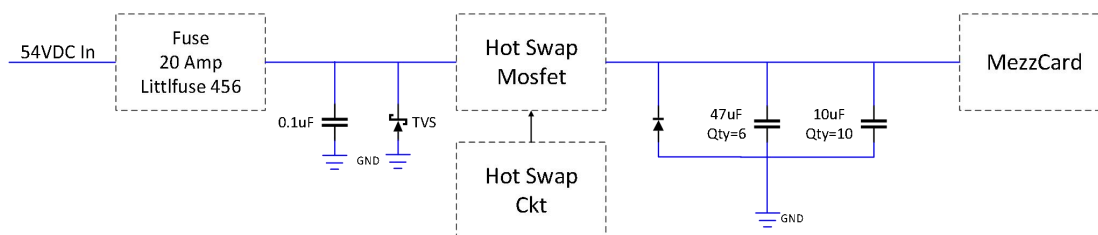


Figure 4.2 HSC circuit block diagram

The HSC Controller recommends T I's L M5069 and L M5066l. When selecting M OS type, focus on S OA curve, recommend IPB017N10N5LF of Infineon and PSMN4R8-100BSE. of Nexperia

4.2.4 Power timing

The MLU370-M8 intelligent accelerating card 54V is powered up normally, and after the 156.25MHz reference clock is stable, the H OST_PWRGD signal is sent out. The detailed timing diagram is as follows:

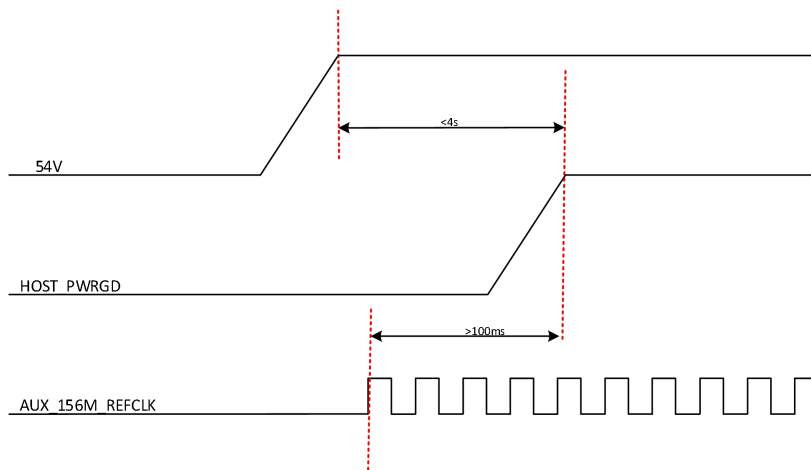


Figure 4.3 Power timing

4.3Signal description

4.3.1 Clock signal

AUX_156M_REFCLKp/n receives a common differential level signal input such as LVPECL, LVDS, CML, HCSL, since its phase noise jitter (12 KHz – 20 MHz) requires less than 270 fs. as a MLU-LINK high-speed SERDES reference clock

Table 4.5 AUX_156M_REFCLK Specifications

Parameters	Conditions	Minimum value	Typical value	Maximum value	Unit
Frequency	--		156.25		MHz
Frequency deviation	--	-100	--	100	PPM
Differential swing	--	0.6	0.8	1.0	V
Occupy empty ratio	--	45	50	55	%
Phase noise jitter	12 KHz – 20 MHz			270	fs RMS

4.3.2 PCIe signal

MLU370-M8 intelligent accelerating card supports PCIe GEN4.0 rate, default x 16 bit width, no AC coupling capacitance placed in the MLU370-M8 module, AC coupling capacitances sent and received are placed on the main board, AC coupling capacitance value range 176n F-265nF, recommended using 220n F, placement position reference figure below:

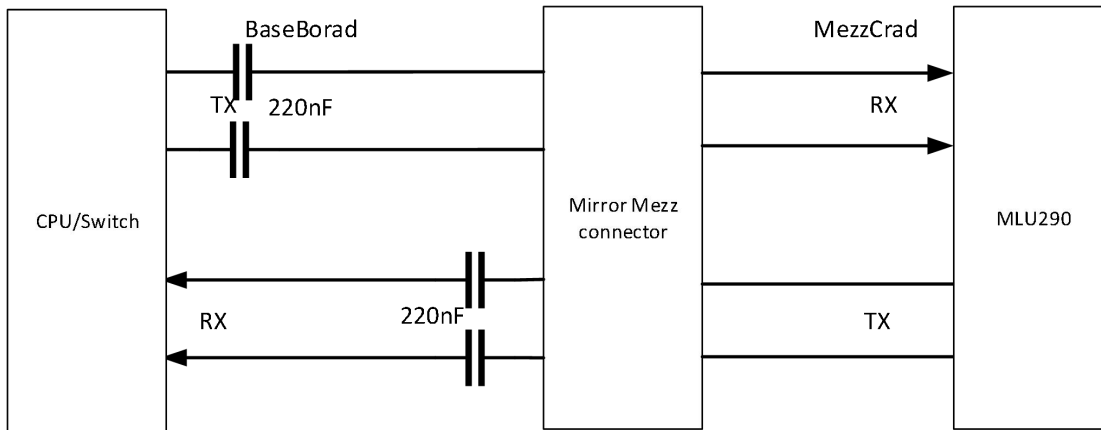


Figure 4.4 The PCIe AC coupling capacitance

The PE_REFCLK reference clock meets the requirements of the PCIe GEN4.0 specification, and the phase noise jitter (12 KHz – 20 MHz) is less than 0.5ps.

Table 4.6 Description of the PCIe signal

Signal	Signal direction	Signal description
PETp/n [15:0]	O	The PCIe sends a signal. The MLU370-M8 module is sent, and the motherboard receives it
PERp/n [15:0]	I	The PCIe receives the signal. MLU370-M8 module receives, motherboard send
PE_REFCLKp/n	I	The PCIe 100MHz reference clock

4.3.3 The MLU-Link signal

MLU-Link is the MLU370-M8 intelligent accelerating card interconnection link. A total of 4 port, per port has 4 lane, maximum rate support of 50Gbps PAM4, feet compatible with OAM specification definition, and the module board meets the 14dB requirements.

The 16 pairs defined by MLU-Link and serdes correspond to the feet defined by OAM as shown in the following table, and the other lane are not used in MLU370-M8 and suspended. See the reference topology chapter in the detailed interconnection topology.

Table 4.7 Description of the MLU-L ink signal

Signal	Signal direction	Signal description
MLU-Link_0Tp/n [3:0]	O	The MLU-L ink0[3:0] sends signals and defines the pin SERDES_1Tp/n [3:0] corresponding to the OAM specification
MLU-Link_0Rp/n [3:0]	I	The MLU-L ink0[3:0] receives the signals and defines the pin SERDES_1R p/n [3:0] corresponding to the OAM specification
MLU-Link_1Tp/n [3:0]	O	The MLU-L ink1[3:0] sends signals and defines the pin SERDES_3Tp/n [3:0] corresponding to the OAM specification
MLU-Link_1Rp/n [3:0]	I	The MLU-L ink1[3:0] receives the signals and defines the pin SERDES_3R p/n [3:0] corresponding to the OAM specification
MLU-Link_2Tp/n [3:0]	O	The MLU-L ink2[3:0] sends signals and defines the pin SERDES_4Tp/n [15:12] corresponding to the OAM specification
MLU-Link_2Rp/n [3:0]	I	The MLU-L ink2[3:0] receives the signals and defines the pin SERDES_4R p/n [15:12] corresponding to the OAM specification
MLU-Link_3Tp/n [3:0]	O	The MLU-L ink3[3:0] sends signals and defines the pin SERDES_6Tp/n [15:12] corresponding to the OAM specification
MLU-Link_3Rp/n [3:0]	I	The MLU-LI ink3[3:0] receives the signals and defines the pin SERDES_6R p/n [15:12] corresponding to the OAM specification

4.3.4 Other signal instructions

4.3.4.1 I2C signal and register description

I2C_SLV_D/CLK # meets the I2C protocol with a maximum rate of 400KHz, with an 8bit address of 0x 8E(write) / 0x 8F(read) and the MLU370-M8 module works in Slave mode.

Table 4.8 Description of the I2C signal

Signal	Signal direction	Signal description
I2C_SLV_D	I/O	The I2C data signal is SDA
I2C_SLV_CLK	I	The I2C clock signal is SCL
I2C_SLV_ALERT#	O	Reserved, NC processing

Register Read, operation process:

1	7	1	1	8	1	1	1	7	1	1	8	1	8	1	...	8	1	1
S	S	l	A	Reg	A	P	S	Slave	R	A	Data	A	Data	A	...	Data	N	P
	Address	W	A	ister				Addr			1		2			N		
	s			adres				ess										
				s														

Register writing, operating process:

1	7	1	1	8	1	8	1	8	1	...	8	1	1
S	S	l	A	Reg	A	Data1	A	Data2	A	...	DataN	A	P
	Address	W	A	ister									
				adres									

Among them, S (start), A (ACK), Sr (restart), N (NACK), P (STOP), W (write), and R (read).

The I2C register is as follows:

Table 4.9 Description of the I2C Register

Register definition	Address	Read and write	Description
Board card power consumption	0x 01	RO	[15:7] Power consumption, data type unsigned integer, in W [6:0] Reserve [31:16] Reserve
Plate card temperature	0x02	RO	[7:0] Card temperature + 100, data type unsigned integer, in °C [31:8] Reserve
Chip temperature	0x 03	RO	[7:0] Chip temperature + 100, data type unsigned integer, in °C [31:8] Reserve
Memory temperature	0x04	RO	[7:0] Memory temperature + 100, data type unsigned integer, in °C [31:8] Reserve
power brake	0x 05	WO	[7:0] The write 0x04 primary frequency drops to the current 25%; the write 0x01 returns to the pre-drop level [31:8] Reserved
Card power consumption set state	0x19	RO	[15:7] Card TDP power consumption data type unsigned integer, in W [31:23] Current power capping value is unsigned integer in W [6:0] Reserve

			[22:16] Reserve
IPU frequency status	0x1A	RO	[15:4] Current frequency capping value data type unsigned integer in MHz [31:20] The current IPU frequency data type is an unsigned integer, in MHz [3:0] Reserve [19:16] Reserve
The Frequency capping sets the range	0x1B	RO	[15:4] frequency capping can set the maximum value data type unsigned integer, in MHz [31:20] frequency capping can set the minimum data type unsigned integer, in MHz [3:0] Reserve [19:16] Reserve
Utilization Information 1	0x1E	RO	[7:0]: PCIe bandwidth utilization; LSB [15:8]: IPU utilization; LSB [23:16]: memery5 bandwidth utilization; LSB [31:24]: memery6 bandwidth utilization; LSB
Utilization Information 2	0x1F	RO	[7:0]: memery1 bandwidth utilization; LSB [15:8]: memery2 bandwidth utilization; LSB [23:16]: memery3 bandwidth utilization; LSB [31:24]: memery4 bandwidth utilization; LSB
Enabling status information	0x20	RO	Does b it0:power brake hold power position Is bit1:thermal alert, enabling position [5:2] Reserved Is bit6:power capping, enabling position Does bit 7:frequency capping hold power position [31:8] Reserved
Temperature, threshold, information	0x23	RO	[7:0] Chip super-temperature and frequency reduction temperature [15:8] Chip over-temperature and

			power-off temperature [23:16] Maximum operating temperature of the chip [31:24], HBM maximum operating temperature
Power capping control	0x29	RW	[15:0] Board card power capping power consumption control value [31:16] Reserved
Frequency capping control	0x2A	RW	[15:0] frequency capping control value [31:16] Reserved
PCI e Vendor ID and Device ID	0xA0	RO	[15:0] Vendor ID :0xCABC [31:16] Device ID :0x0370
PCI e Sub-Vendor ID and Sub-System ID	0xA1	RO	[15:0] Sub-Vendor ID :0xCABC [31:16] Sub-System ID : 0x0055
PCI e_negotiated_speed	0xA2	RO	[7:0] Shows the PCI e negotiation rate, for example, 0x 04 indicates gen4 16GT/s [31:8] Reserved
PCI e_negotiated_link_width	0xA3	RO	[7:0] Displays the PCI e negotiation width, for example, 0x16 indicates X16 [31:8] Reserved
Data-readable flags	0x B8	RO	Bit0: card inherent information readable flag (0x F0-0xFF), 1: readable, 0: unreadable Bit 1: card real-time information readable flag (0x01-0x EF), 1: readable, 0: unreadable [31:2] Reserve
Plate card type	0xF0	RO	[7:0] Display the card type, with a default value of 0x55 [31:8] Reserved
Equipment manufacturer	0xF1	RO	[4:0] Displays the device manufacturer number, with a default value of 0x 3 [31:5] Reserved
Hardware version number	0xF2	RO	[7:0] Displays the hardware version number, such as 0x10 indicates that the hardware version is V1.0 [31:8] Reserved
The firmware version number	0xF3	RO	[11:0] Displays the firmware version number, for example, the 0x100 main

			version number 0x01 subversion number 0x0 patch number is 0x0 [31:12] Reserved
Manufacturing time	0xF4	RO	[15:0] Show manufacturing years, for example 0x2006 indicates June 2020 production [31:16] Reserved
Device sequence number	0xF5	RO	[19:0] Displays the serial number of the device, for example, 0x00018 indicates that the serial number is 00018 [31:20] Reserved
SN, Low	0xF6	RO	[31:0] The SN number has low 8 bit data, for example the low 8 bit data for SN:552006300018 was saved as 0630001 8
High SN	0xF7	RO	[15:0] The SN number high 4-bit data, for example the SN:552006300018 high 4-bit data was saved as 5520 [31:16] Reserved
P art _number_1	0xF8	RO	[31:0] part_number Low 8-bit data ("ASCII code corresponding to MLU3" characters)
P art _number_2	0xF9	RO	[31:0] Eight-bit data in part_number ("ASCII code corresponding to 70-M" characters)
P art _number_3	0xF A	RO	[7:0] part_number High 8-bit data (ASCII code corresponding to the "8" characters) [31:8] Reserved

4.3.4.2 THERMTRIP# Over temperature alarm signal

The THERMTRIP# super temperature alarm is triggered at the MLU370-M8 chip 95°C temperature and is irreversible after taking effect. The system heat dissipation design recommends the MLU370-M8 chip temperature 88°C trigger fan full speed, 92°C trigger power consumption to 1 / 2, 93°C power consumption to 1 / 4, 94°C power consumption to 1 / 8. When the temperature rises to 95°C indicates that the power reduction measures fail, there may be fan failure or other catastrophic failure in the system, then the MLU370-M8 intelligent accelerating card will pull down the THERMTRIP# signal and alarm the substrate management system. After 1S, the MLU370-M8 intelligent accelerating card will automatically cut off the power supply for power loss protection.

The THERMTRIP# signal will continue to be low after taking effect. It is recommended to confirm on site and restart equipment recovery after troubleshooting.

4.3.4.3 Power consumption, configuration, signal

The MLU370-M8 Smart Accelerator provides power configuration pins on both PWRBRK# and PWRRDT#[1:0] hardware. The PWRBRK# (power brake) is based on the current power consumption, and the low levels are effectively reduced to 1 / 4 of the current frequency. PWRRDT#[1:0] is the hardware implementation of power capping, which can set four TDP power consumption such as 300W,250W,200W, 150W. It is recommended to pull up to 3.3V, on the motherboard for better performance.

At the same time, the MLU370-M8 intelligent accelerating card supports setting the TDP power consumption upper limit by issuing power capping instructions through the I2C interface, and the setting range is 150W-300W. Both power consumption configurations can take effect, with the I 2C configuration prevail when configured at the same time.

Table 4.10 Description of the power consumption configuration signal

Signal	Signal direction	Signal description
PWRBRK#	I	Power brake power consumption reduces frequency to 1 / 4 of current frequency, low level effective
PWRRDT#[1:0]	I	<p>TDP power consumption set tube foot, the motherboard needs to provide a default 3.3V high level</p> <p>11-L0 level, Normal TDP power consumption is 300W, default value</p> <p>The 10-L1 level, TDP power consumption was reduced to 250W</p> <p>The 01-L2 level, TDP power consumption was reduced to 200W</p> <p>The 00-L3 level, TDP power consumption was reduced to 150W</p>

4.3.4.4 Other configuration signal

Table 4.11 Description of the other configuration signals

Signal	Signal direction	Signal description
--------	------------------	--------------------

MODULE_ID[4:0]	I	MLU370-M8 module slot bit ID.MLU370-M8 module internal default 10K pull-up
LINK_CONFIG[4:0]	I	serdes link configuration topology, default 10 K pull-up inside the MLU370-M8 module
PE_BIF[1:0]	O	PCI e interface bit width indication: the MLU370-M8 module default is 00 00-1 x16(defaults) 01-Reserve 10-Reserve 11-Reserve
PLINK_CAP	O	"PCI e" port Protocol support: The MLU370-M8 module default value is 0 0 = supports the PCI e protocol only (default) 1 = Other protocols are supported, reserved

4.3.4.5 Reserved signals

The following signal is defined in OAM but not used in MLU370-M8 intelligent accelerating card only reserved, recommended for NC processing.

Table 4.12 Description of the reserved signal

Signal	Signal direction	Signal description
SCALE_DEBUG_EN	O	The function is not defined
DEBUG_PORT_PRSENT#	I	The commissioning signal of the main board.NC processing is available
MANF_MODE#	I	The function is not defined
FW_RECOVERY#	I	The function is not defined
TEST_MODE#	I	Test mode.NC processing is available

RFU	/	Keep the pipe feet in advance
NC	/	Hanging pipe feet



5. Heat dissipation specifications

5.1 Radiation description

The MLU370-M8 intelligent accelerating card uses passive heat dissipation in a windy environment of the system.

5.2 Speed reduction protection temperature

Reducing working temperature refers to the MLU 370 chip temperature T_j reaches this temperature point to prevent the temperature by reducing the working clock frequency to ensure the reliability of the calculation card operation, but will cause a reduction in the performance of the calculation card. MLU370 intelligent accelerating card main chip junction temperature reaches 92°C will trigger a deceleration.

Table 5.1 Description of deceleration protection

temperature	Power consumption after speed reduction
92°C	1/2 TDP
93°C	1/4 TDP
94°C	1/8 TDP

5.3 Turn off temperature

Turning off the temperature means that the MLU370 chip junction temperature T_j will prevent the main chip from continuing to heat up by cutting off the power supply to ensure that no permanent damage will be caused. The calculation card should not trigger the off temperature during normal operation, and the system may have a catastrophic error during triggering the off temperature. MLU370-M8 chip warming to 95°C will trigger a shutdown.

When the MLU370-M8 junction temperature reaches 95°C , the THERMTRIP# signal is triggered and the MLU370-M8 Intelligent Accelerating card power supply is turned off after 1S. After the power supply is turned off, it is recommended to manually eliminate the heat dissipation problem

on site before adding power.

5.4 Working environment air volume requirements

The MLU370-M8 intelligent accelerating card supports the use requirements for 0-45°C ambient temperature (inlet temperature of board card radiator) (TDP mode), and the minimum radiator air capacity requirements for each main temperature conditions are shown in the following table:

Table 5.2 VS ambient temperature gauge for radiator

Air inlet temperature is (°C)	Radiator Minimum air requirement (CFM)
25	22
30	31
35	40
40	60
45	94

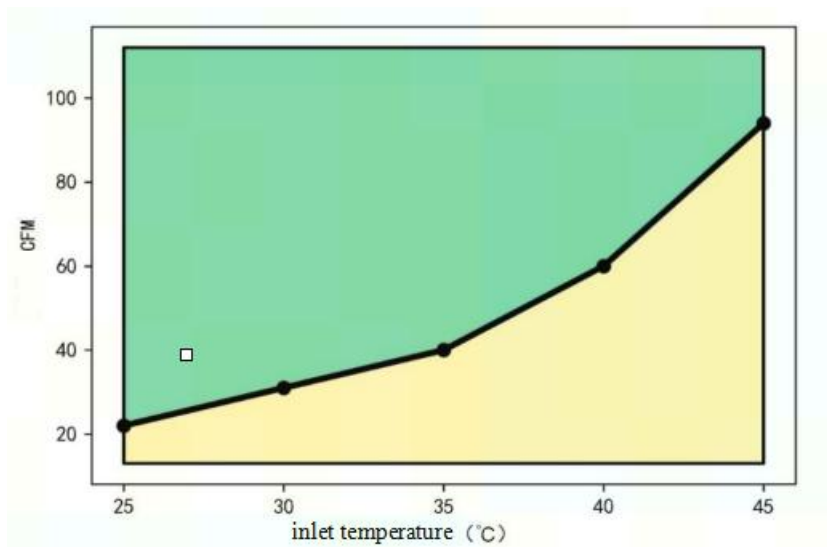


Figure 5.1 Radiator minimum air volume demand vs air inlet temperature

The MLU370-M8 intelligent accelerating card is recommended in the green area environmental conditions.

5.5 Thermal simulation model

Thermal simulation models are shown in *MLU370-M8_PDML_V1_0*.



6. The Cambricon NeuWare development environment

NeuWare fully supports all kinds of mainstream programming frameworks (such as TensorFlow, Caffe, PyTorch and MXNet, etc.).Users can easily develop and deploy deep learning applications on the Cambricon MLU370-M8.At the same time, NeuWare provides full runtime systems and drive software for rapid system integration.

NeuWare also offers a range of tools including application development, feature debugging, performance tuning, and more.Application development tools include machine learning library, runtime library, compiler, model retraining tools and specific fields (such as video analysis field) SDK, etc.; functional debugging tools can meet debugging needs at different levels such as programming framework and function library; performance tuning tools include performance analysis tools and system monitoring tools.

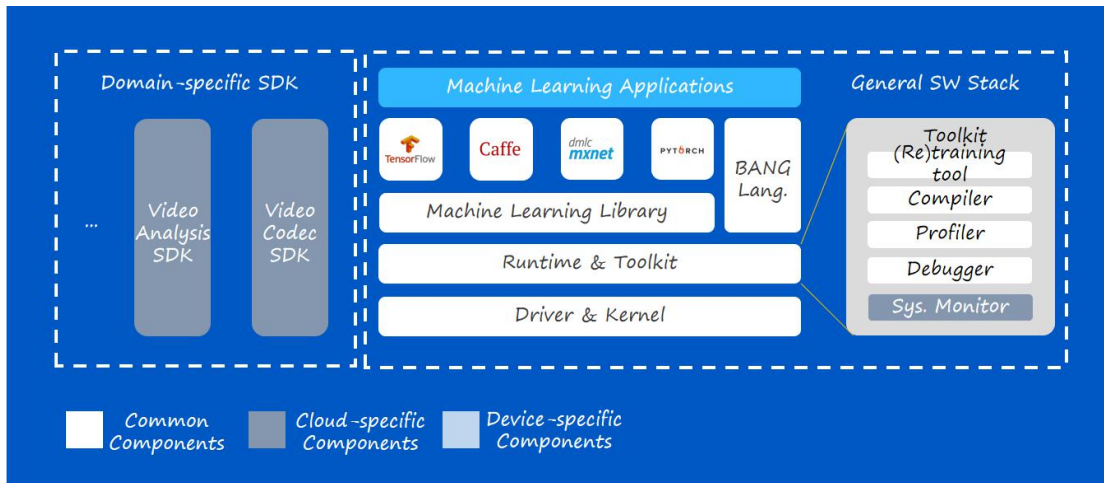


Figure 6.1 Cambricon NeuWare



7. Compliance

The MLU370-M8 is compliant with the regulations listed in this section. Compliance marks, including the FCC ID numbers, can be found on the labels of each devices.

United States

Federal Communications Commission (FCC)

This device complies with Part 15 of the FCC Rules.

Operation is subject to the following two conditions: (1) This device may not cause harmful interference, and (2) this device must accept any interference received, including interference that may cause undesired operation.

This equipment has been tested and found to comply with the limits for a Class A digital device, pursuant to part 15 of the FCC Rules. These limits are designed to provide reasonable protection against harmful interference when the equipment is operated in a commercial environment. This equipment generates, uses, and can radiate radio frequency energy and, if not installed and used in accordance with the instruction manual, may cause harmful interference to radio communications. Operation of this equipment in a residential area is likely to cause harmful interference in which case the user will be required to correct the interference at his own expense.

Caution: Any changes or modifications not expressly approved by the party responsible for compliance could void the user's authority to operate this equipment.

Underwriters Laboratories (UL)

UL Listed Product Logo for MLU370-M8 Intelligent Processing Cards , model name MLU370-M8.



The surface temperature of the product is very high during operation.

Please do not touch the surface without taking care.