



AIBOX-1688

| 16T Large Model AI Box



V1.0 2024-8-9

T-CHIP INTELLIGENCE TECHNOLOGY

Product features



32T INT4/16T INT8 Computing power

Equipped with SOPHON computing AI processor BM1688, it has an octa-core ARM Cortex-A53, with a maximum frequency of 1.6GHz, and a built-in neural network acceleration engine TPU, with 32T@INT4 peak computing power, 16T@INT8 peak computing power, 4T@FP16/BF16 computing power, and 0.5T@FP32 computing power.



Powerful multi-channel video AI performance

The AI box supports up to 32 channels of 1080P H.264/H.265 video decoding and 32 channels of 1080P HD video processing (decoding + AI analysis), making it ideal for various AI applications such as face detection and license plate recognition on video streaming.



The private deployment of large language models

Support the private deployment of ultra-large-scale parameter models under the Transformer architecture, including large language models such as LLaMa2, ChatGLM, and Qwen, as well as large vision models like ViT, Grounding DINO, and SAM.



Multiple deep learning frameworks

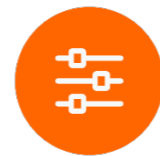
Support traditional network architectures such as CNN, RNN, and LSTM; a variety of deep learning frameworks, including TensorFlow, PyTorch, MXNet, PaddlePaddle, and ONNX, as well as custom operator development.

Product features



All-aluminum alloy enclosure for heat dissipation

The industrial-grade all-metal enclosure with aluminum alloy structure for thermal conduction. The side of the top cover features a grille design for efficient heat dissipation. Its top cover is a porous hexagonal design, combining elegance with high efficiency. The compact, exquisite device operates stably and meets the needs of various industrial-grade applications.



Abundant expansion interfaces

It has dual 1000Mbps Ethernet, 2 * USB3.0, 1 * TF Card, 1 * Type-C, 1 * HDMI2.0, 1 * Console and other expansion interfaces, which are convenient to connect various peripherals and realize multi-field product applications.



A wide range of applications

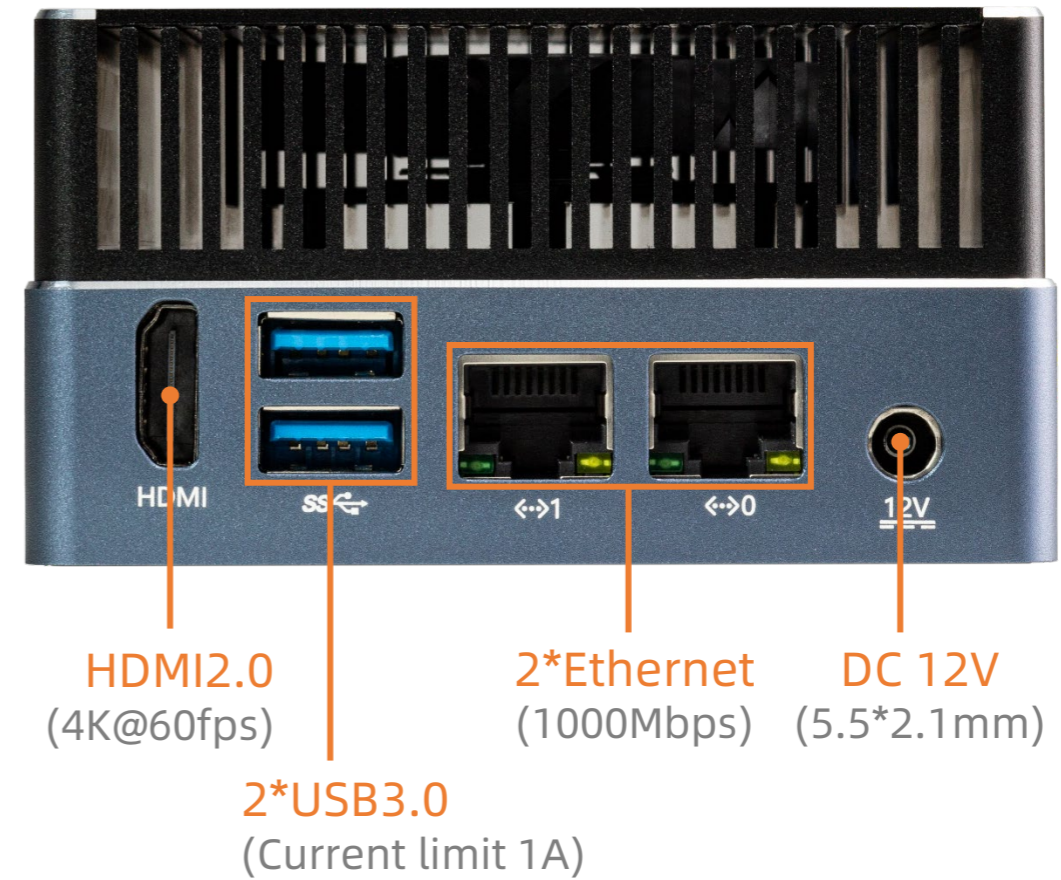
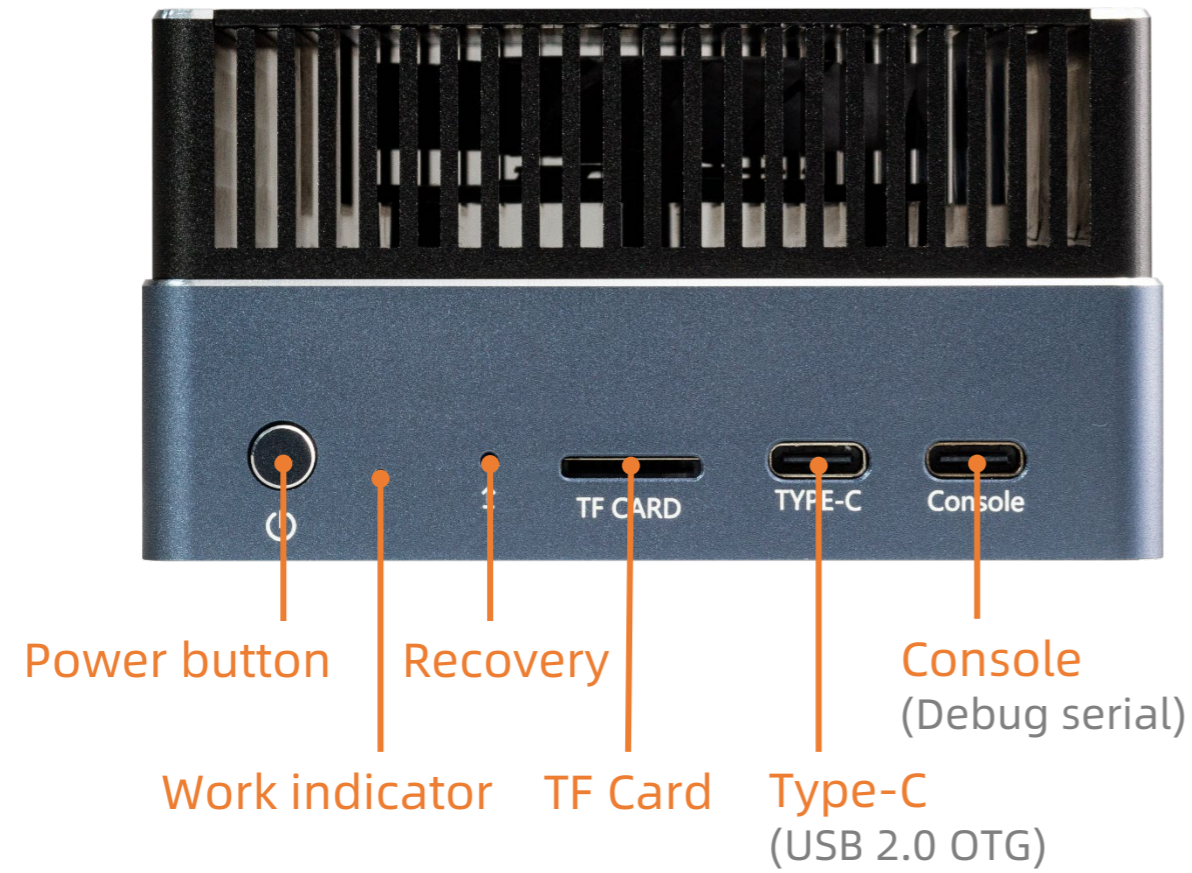
The device is widely used in intelligent surveillance, AI education, services based on computing power, edge computing, private deployment of large models, and data security and privacy protection.

Specifications

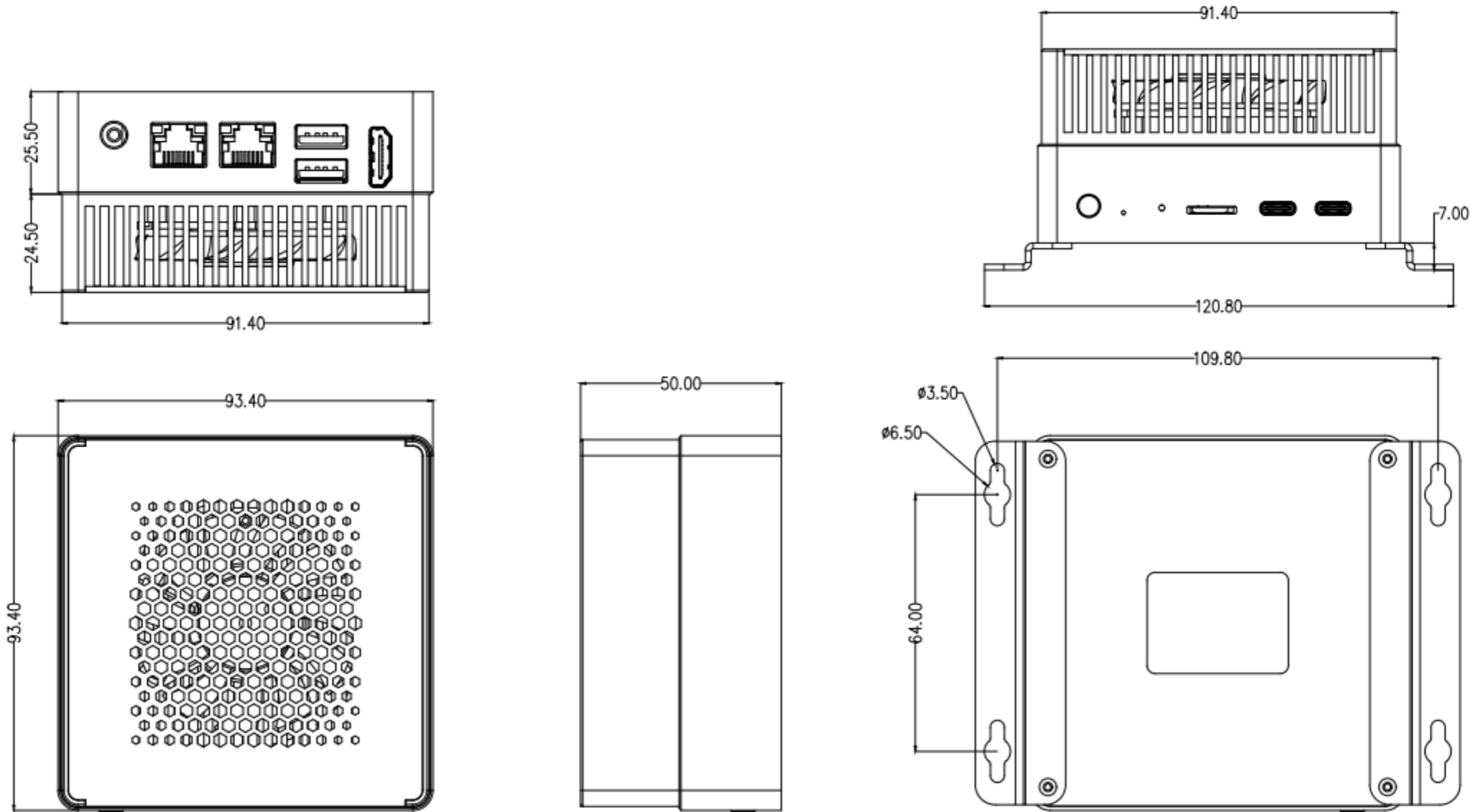


Specification		
Basic Specifications	SOC	SOPHON BM1688
	CPU	Octa-core ARM Cortex-A53 @ 1.6GHz
	TPU	Built-in SOPHGO neural network acceleration engine TPU, 32T@INT4 peak computing power, 16T@INT8 peak computing power, 4T@FP16/BF16 computing power, 0.5T@FP32 computing power
	Decoding/Encoding	Video decoding: H.264 / H.265 decoding (Max performance: 1920 * 1080@480FPS or 3840 * 2160 @120FPS) Video encoding: H.264 / H.265 encoding (Max performance: 1920 * 1080@300FPS or 3840 * 2160 @75 FPS) Image codec: JPEG/MJPEG Baseline codec (JPEG codec with a maximum resolution of 1080P@480 FPS)
	RAM	8GB LPDDR4 (4GB/8GB/16GB optional)
	Storage	32GB eMMC (32GB/64GB/128GB/256GB optional)
	Storage Expansion	1*M.2 (Expandable PCIe NVMe SSD(default support)/ SATA SSD(supported after software update), supports 2242/2260/2280) (inside the device), 1*TF Card
	Power	DC 12V/3A (DC 5.5*2.1mm)
	Power consumption	Normal: 7.2W(12V/600mA), Max: 14.4W(12V/1200mA)
	OS	Linux
	Software support	<ul style="list-style-type: none"> The private deployment of ultra-large-scale parameter models under the Transformer architecture, including large language models such as Gemma-2B, LLaMa2-7B, ChatGLM3-6B, Qwen1.5-1.8B. Traditional network architectures such as CNN, RNN, and LSTM; a variety of deep learning frameworks, including TensorFlow, PyTorch, MXNet, PaddlePaddle, and ONNX, as well as custom operator development Docker container management technology
	Size	93.4mm * 93.4mm * 50 mm
	Weight	≈ 500g
	Environment	Operating Temperature: -20°C ~ 60°C, Storage Temperature: -20°C ~ 70°C, Storage Humidity: 10% ~ 90%RH (non-condensing)
Interface Specifications	Ethernet	2*1000Mbps Ethernet
	Video output	1*HDMI2.0 (4K@60fps)
	USB	2*USB3.0 (Current Limit: 1A)
	Other interfaces	1*Type-C (USB 2.0 OTG), 1*Console (Debug serial), 1*Power button, 1*Recovery

Interface description



Dimension





T-CHIP INTELLIGENCE TECHNOLOGY



Contact Us
(+86)18688117175



E-mail
global@t-firefly.com



Website
<https://en.t-firefly.com/>



Address
Room 2101, Hongyu Building, #57 Zhongshan 4Rd, East District,
Zhongshan, Guangdong, China.