

## KunTai A722推理服务器：多种推理卡灵活配置满足多样需求



Atlas 300I Pro  
Atlas 300V Pro  
Atlas 300V

支持 1 ~ 7 张卡

插卡式



Atlas 300I DUO

支持 1 ~ 4 张卡



KunTai A722

\*推理卡不能混插使用

**应用场景** 部署于数据中心机房中，进行AI推理或小规模训练

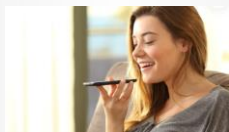
搜索推荐



金融大脑



语音识别



内容审核



### 关键特性

### 规格描述

型号	KunTai A722
形态	2U机架服务器 86.1mm × 447mm × 790mm
CPU	2 x 鲲鹏920
内存	最多32个DDR4内存插槽
AI加速卡	最大支持 7 个Atlas 300I/V Pro/ Atlas 300V 或 4 个Atlas 300I DUO
散热	风冷散热

## Atlas 300I Duo 推理卡：双芯配置，性能翻倍



### 应用场景

- 集成于服务器，进行AI推理
- 主要用于互联网场景
- 受友商引导时也可用于其他推理场景

#### 互联网(主要场景)

推荐、检索聚类



#### 智慧城市

城市治理



#### 智慧金融

智慧营业厅、推荐



#### 智慧园区

目标识别、结构化



### 关键特性

### 规格描述

形态	单槽位全高全长PCIe卡
尺寸	266.7mm (长) × 18.46mm (宽) × 111.15mm (高)
处理器核	内置16个Arm Core (最大主频1.9GHz)
内存	容量：96GB，总带宽：408 GB/s；支持ECC
AI算力	整数精度 (INT8)：280 TOPS 半精度 (FP16)：140 TFLOPS
编解码能力	内置DVPP (数字视觉预处理) 单元 视频 256路 1080P 30FPS (硬件解码能力)
功耗	150 W

### 优势

- 单卡对业界厂商时，Atlas有效算力更高
- 视频解码能力为业界厂商的2倍，适合视频内容审核场景

1

### 模型性能好

ResNet-50 images/s

Atlas 300I Duo vs 业界厂商 1.3X ↑

2

### 视频解析强

1080P 30FPS

Atlas 300I Duo vs 业界厂商 2X ↑



# Atlas 300I Pro 推理卡：扩大算力优势，提升典型模型性能



**应用场景** 集成于服务器中，进行AI推理

关键特性	规格描述
形态	单槽位半高半长PCIe卡
尺寸	169.5mm（长） x 18.45mm（宽） x 68.9mm（高）
处理器核	内置8个Arm Core（最大主频1.9GHz）
内存	容量：24GB； 支持ECC
AI算力	整数精度（INT8）： <b>140</b> TOPS 半精度（FP16）： <b>70</b> TFLOPS
功耗	72 W

## 1 稳定算力优势显著

Atlas 300I Pro **VS** 业界厂商

**2.1X** ↑  
TOPS (INT8)

## 2 性能功耗比优

Atlas 300I Pro **VS** 业界厂商

**2.1X** ↑  
TOPS/W (INT8)

## 3 典型模型性能好

Atlas 300I Pro **VS** 业界厂商

**1.7X** ↑  
BERT large  
Sentences/Second

搜索推荐



金融大脑



语音识别



内容审核



## Atlas 300V Pro 视频解析卡：业界极致性能视频分析卡



Atlas 300V Pro 视频解析卡

插卡式

7卡



KunTai A722 推理服务器

**应用场景** 集成于服务器中，提供超强AI推理、视频图片编解码等功能

关键特性	规格描述
形态	单槽位半高半长PCIe卡
尺寸	169.5mm (长) x 18.45mm (宽) x 68.9mm (高)
处理器核	内置8个Arm Core (最大主频1.9GHz)
内存	容量：48G；支持ECC
AI算力	整数精度 (INT8) : <b>140</b> TOPS 半精度 (FP16) : <b>70</b> TFLOPS
编解码能力	内置DVPP (数字视觉预处理) 单元 视频 <b>128路</b> 1080P 30FPS (硬件解码能力)
功耗	72 W

### 智慧城市



### 智慧交通

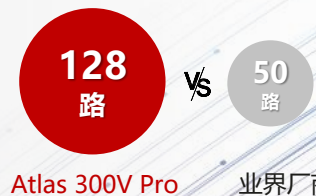


### 智慧园区

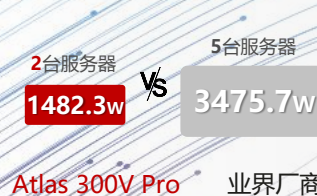


### 支持更多高清视频解析路数

单卡支持视频解析路数对比



支持1000路数视频解析 功耗对比





## Atlas 300V 视频解析卡：极高性价比视频分析卡



Atlas 300V 视频解析卡

插卡式

7卡



KunTai A722 推理服务器

### 应用场景

- 集成于服务器中，提供超强AI推理
- 以CV（Computer Vision-计算机视觉）为主
- 高性价比

#### 智慧城市



#### 智慧交通



#### 智慧园区



关键特性	规格描述
形态	单槽位半高半长PCIe卡
尺寸	169.5mm（长） x 18.45mm（宽） x 68.9mm（高）
处理器核	内置8个Arm Core（最大主频1.9GHz）
内存	容量： <b>24G</b> ；支持ECC
AI算力	整数精度（INT8）： <b>100</b> TOPS 半精度（FP16）： <b>50</b> TFLOPS
编解码能力	内置DVPP预处理单元 视频 <b>100路</b> 1080P 25FPS
功耗	72 W

### 支持更多高清视频解析路数

单卡支持视频解析路数对比

能效（路/W）对比

