# UNTETHER AI

## tsunAImi®
## accelerator cards
### tsn800

### Key Features

| | | | |
|---|---|---|---|
| **2 PetaOps of inference performance** | **800MB of on-chip SRAM** | **Low latency, batch = 1** | **80,000 fps ResNet-50 V1.5** |
| **12,000 qps Bert-base** | **32GB/s PCIe bandwidth** | **Multi-chip partitioning of large neural networks** | **Thermal monitoring** |

## Overview

The tsunAImi® tsn800 accelerator card is built for high performance servers and data centers, delivering an industry best compute density of over 2 PetaOps of INT8 performance. The tsn800 is powered by four runAI200™ devices, and due to their superior power efficiency, remains within a 300W Thermal Design Point (TDP). The x16 PCI-Express Gen4 interface supports up to 32 GB/s of bandwidth, enough for the most demanding AI applications.

## Applications

The on-board runAI200 devices are designed to accelerate a multiplicity of AI workloads, such as vision-based convolutional networks, transformer networks for natural language processing, and time-series analysis for financial applications.

| Markets | Application | Networks |
|---|---|---|
| Vision | Classification, object detection, semantic segmentation | ResNets, YOLO, SSD, Unets |
| Natural language processing | Text-to-speech, speech-to-text, chatbots | RNNs, Attention, BERT |
| Financial technology | X-Value adjustments, credit risk, portfolio balancing | TCNs, LSTMs |

## imAIgine Software Development Kit

The imAIgine SDK gives developers powerful automated tools and supporting software to quickly go from pilot model to production. It is organized into three parts.
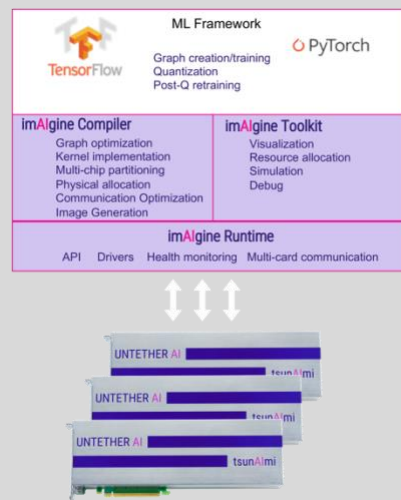
### The imAIgine Compiler
- Import TensorFlow, PyTorch, or ONNX graphs directly
- Automated quantizer and extracts performance without sacrificing accuracy
- Specify performance levels, silicon utilization, and power consumption targets

### The imAIgine Toolkit
- Evaluate functionality and performance using the extensive profiling and simulation tools

### The imAIgine Runtime
- Provides C-based API for integration into your deep learning environment
- Monitor the health and temperature of the tsunAImi® acceleration cards to ensure proper operation and prevent thermal damage
.



**ML Framework**
TensorFlow    PyTorch
Graph creation/training
Quantization
Post-Q retraining

**imAIgine Compiler**
Graph optimization
Kernel implementation
Multi-chip partitioning
Physical allocation
Communication Optimization
Image Generation

**imAIgine Toolkit**
Visualization
Resource allocation
Simulation
Debug

**imAIgine Runtime**
API    Drivers    Health monitoring    Multi-card communication

**Familiar frameworks**
Quantization and layer optimization done in familiar ML framework

**Automated graph lowering**
Optimization and allocation algorithms

**Extensive feedback**
Resource allocations, congestions, cycle-accurate simulation

**Easily integrated runtime**
Hardware abstraction, communication, and monitoring

## Product Specification

| Specification | tsun**AI**mi® tsn800 accelerator card |
|---|---|
| Form factor | Double-wide, full height, full length PCIe |
| PCIe Interface | X16 PCIe Gen4 |
| Clock Frequency | Variable, Up to 840 MHz |
| Memory | 800MB on-chip SRAM |

## Thermal Specification

| Parameter | tsun**AI**mi® tsn800 accelerator card |
|---|---|
| Total board power | TDP 300W, typical application power ~200W |
| Cooling | Passive or active heatsink options available |
| runAI200 maximum operating temperature | 85°C Junction |

## Environmental

| Parameter | tsun**AI**mi® tsn800 accelerator card |
|---|---|
| Operating temperature | 0°C to 55°C |
| Storage temperature | -40°C to 75°C |
| Operating humidity | 5% to 90% relative humidity |
| Storage humidity | 5% to 95% relative humidity |

## Power Connector
8-pin CPU power connector, capable of suppling 300W

### Trademarks
Untether AI, tsunAImi, runAI, imAIgine SDK are trademarks and/or registered trademarks of Untether AI Corporation in the U.S and other countries. Other company names may be trademarks of the respective companies.

**UNTETHER AI**