# T-HEAD

# TH1520 NPU
# User Manual

**Revision**     1.0.0

**Security**     Secret

**Date**        2023-08-26

# Revisions

| Rev | Description | Author(s) | Date |
|:---:|:---|:---:|:---:|
| V1.0.0 | Initial version | T-Head | 2023-08-26 |

# Contents

# Figures & Tables

# List of Abbreviations

| Abbreviations | Full Spelling | Chinese Explanation |
|---|---|---|
| NPU | Neural-network Processing Unit | 神经网络处理单元 |

IV

# 1 Overview

NPU is a hardware based neural network accelerator which achieves high performance with low power. NPU is a key component for SoC targeting neural networks inference acceleration with support for variable precisions for data and weight. NPU supports weight compression and flexible low precision which allow neural networks to be run quickly. Figure & Table 1-1 shows the high level block diagram of NPU.

Figure & Table 1-1 NPU block function diagram

# 2 Main Features

NPU key features are:
- Acceleration of the most common layers of neural networks, listed in Figure & Table 2-1
- Low bandwidth operation
  - Support for a large range of low precision data formats
  - Grouping layers together to reduce memory bandwidth
  - Lossless weight data compression
- DRM security
- Interoperability
  - Support for a variety of memory formats designed to share data with a CPU, GPU or other modules

Figure & Table 2-1 Supported layers

| Layer | Supported Via |
|---|---|
| **Convolution** | |
| Normal Convolution | Hardware |
| Dilated/Atrous Convolution | Hardware |
| Grouped Convolution | Hardware |
| Depthwise Convolution | Hardware |
| Convolution Transpose (Deconvolution) | Hardware |
| **Fully Connected** | |
| Fully Connected | Hardware |
| **Normalization** | |
| Batch normalization | Hardware |
| Local Response Normalization | Hardware |
| L2 Normalization | Software |
| **Activation** | |
| ReLU | Hardware |
| ReLU1 | Hardware |
| ReLU6 | Hardware |
| PReLU | Hardware |
| Clamped ReLU | Hardware |

| | |
|---|---|
| Leaky ReLU | Hardware |
| Tanh | Hardware |
| Sigmoid | Hardware |
| Logistic | Hardware |
| **Pooling** | |
| Max pooling | Hardware |
| Mean pooling | Hardware |
| Min pooling | Hardware |
| **Elementwise Operation** | |
| Negate | Hardware |
| Add/Subtract | Hardware |
| Multiply | Hardware |
| Max/min | Hardware |
| **Memory Operation** | |
| Permute | Hardware |
| Transpose | Hardware |
| Reshape | Hardware |
| Squeeze | Hardware |
| Flatten | Hardware |
| Space to Batch | Hardware |
| Batch to Space | Hardware |
| Depth to Space | Hardware |
| Space to Depth | Hardware |
| **Spatial Resize Operation** | |
| Pad | Hardware |
| Crop | Hardware |
| Bilinear resize | Hardware |
| Nearest neighbor resize | Hardware |
| **Pre-processing** | |
| Mean subtraction | Hardware |

| Post-processing | |
|---|---|
| Softmax | Software |

# 3 Function Description

## 3.1 NPU Processing Order

Figure & Table 3-1 shows the order that layers are processed by NPU. The different layers that can be combined together in single pass through the hardware is called a "layer group". If the order of processing in the target network does not match the order shown in Figure & Table 3-1, then the operations will be split into different layer groups. For example, to perform local response normalization after pooling, the first layer group would perform the pooling layer, and the second would contain the local response.
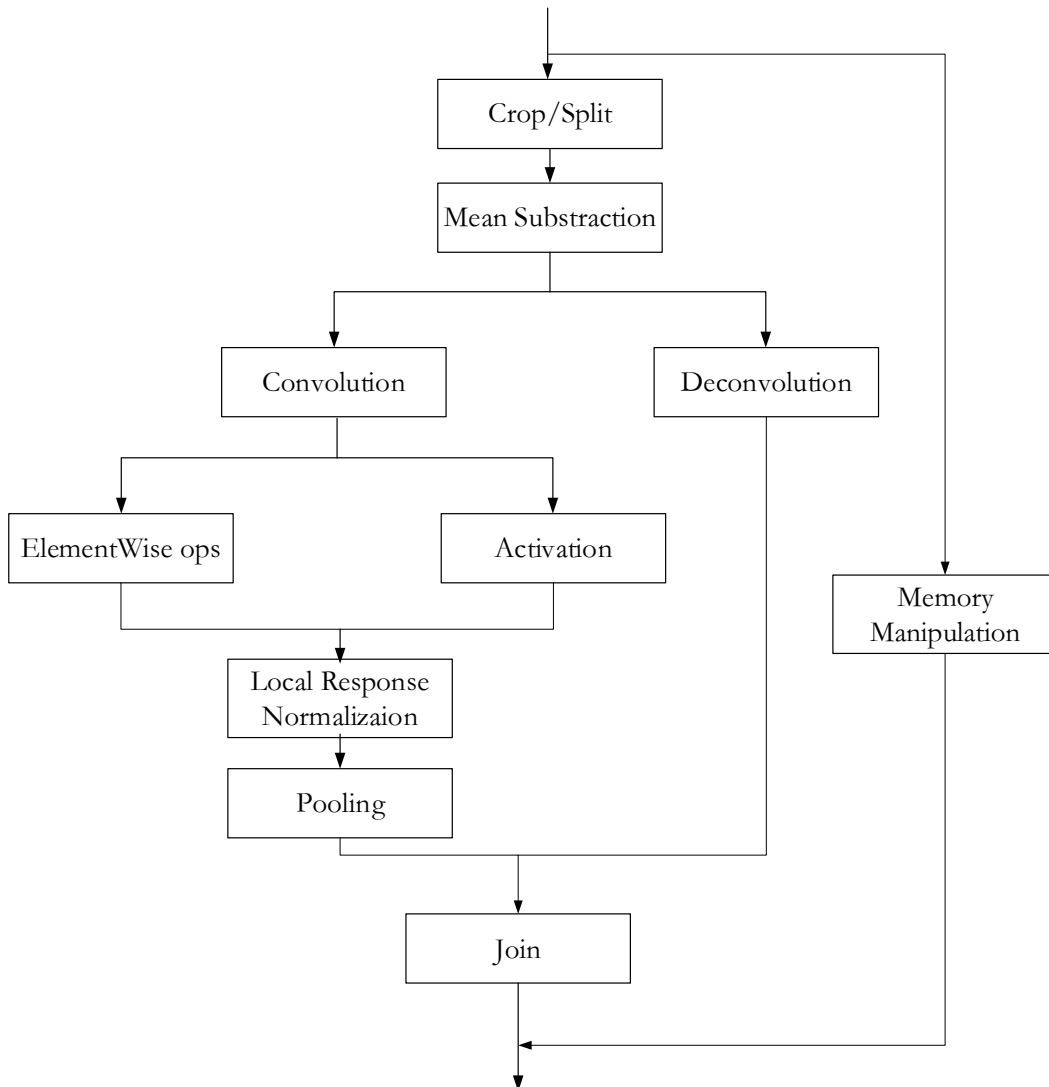
Figure & Table 3-1 NPU flow diagram

# 4 Usage

The power-up procedure of the core is as follows:

1. De-assert NPU resetn.
2. Power on NNA.
3. Enable NNA clock.
4. Wait at least 16 cycles.
5. Assert NPU resetn.
6. Wait at least 8 cycles.
7. Kick off NPU.