

# tinyML<sup>®</sup> Foundation

*Enabling Ultra-low Power Machine Learning at the Edge*

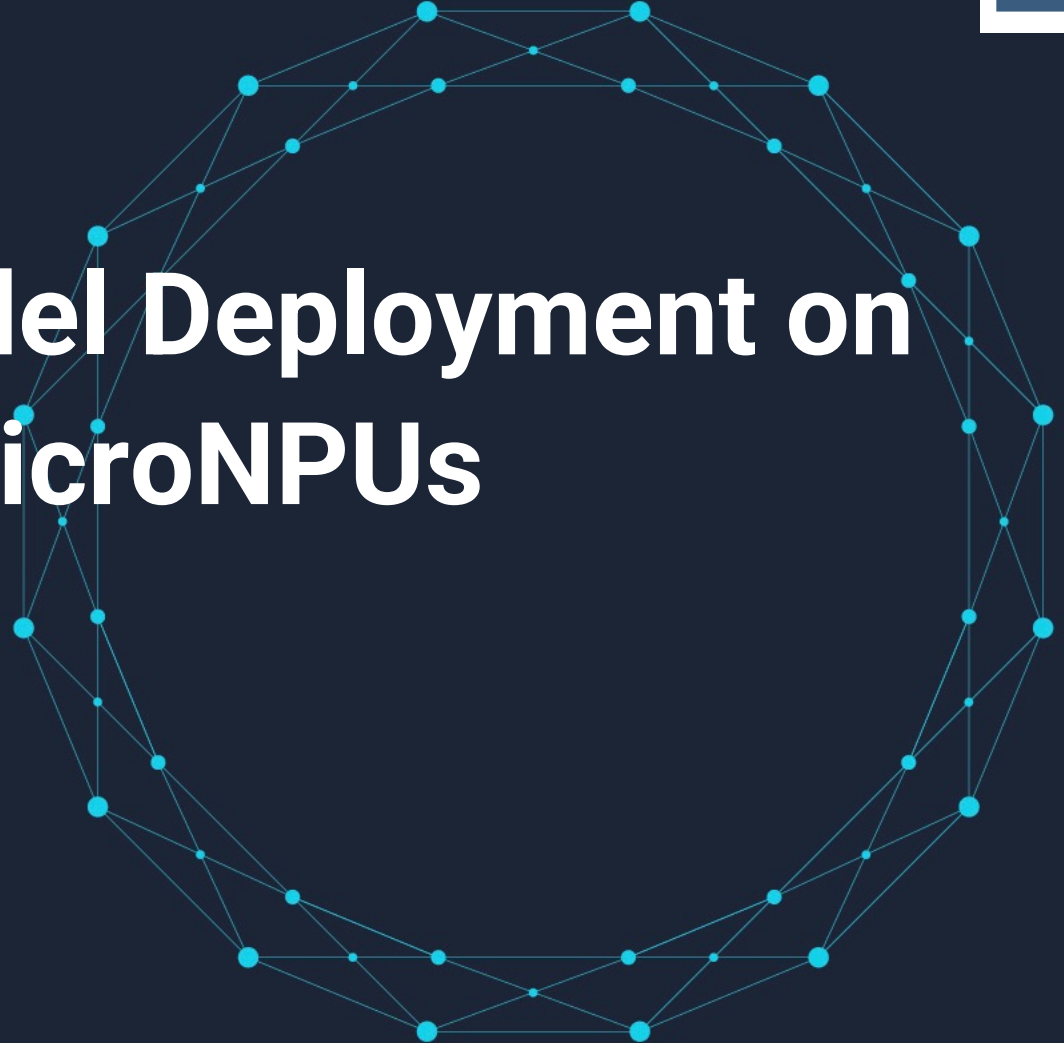
**tinyML Summit April 22 - 24, 2024**



[www.tinyML.org](http://www.tinyML.org)



# Transformer-Based Model Deployment on Edge Devices through MicroNPUs Operator Converter





**Nota AI<sup>®</sup>**

# Nota AI®

Corporate Journey: Pioneering the AI Industry

- Nota AI®'s journey showcases growth, partnerships, events, and awards since inception.



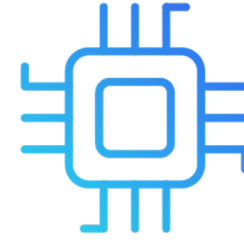
Investment Funding

**US\$  
22.7M**



Global Customers

**101+**



Supported Devices

**30+**



IP Rights

**69+**

# Nota AI<sup>®</sup>

Global Partnerships: Collaborating with Tech Leaders

- Strategic alliances with industry leaders drive Nota AI<sup>®</sup>'s global expansion.

## Partners



arm



RENESAS



seed studio

aetina



## Clients



SAMSUNG SDS



NAVER



LINE



emart

Telechips

coway



# Nota AI®

Platform & Edge Solutions: Elevating Excellence

- Continuously advancing technologically through the utilization of NetsPresso®, Nota ITS, and Nota DMS.



# Table of Contents

## Problem Statement

- Compiler
- Transformer

## Operator Converter

- Operator Converter
- Model Editor

## Results

- Ethos-U NPU
- Conformer, ViT

The background is a solid blue color. Overlaid on this background are several thin, light blue lines that form a complex, abstract pattern. These lines connect small, light blue circular dots, creating a series of interconnected paths that resemble a stylized map or a network diagram. The lines and dots are scattered across the page, with a higher concentration around the central text.

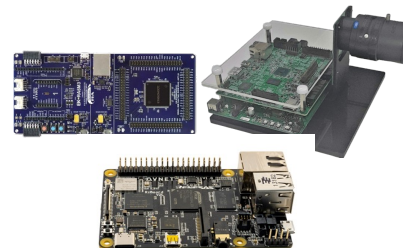
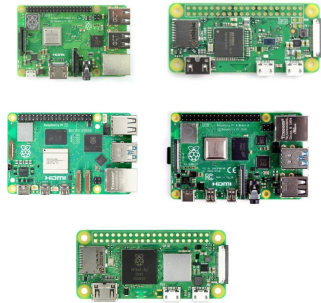
# Problem Statement



# Problem statement

Various device support

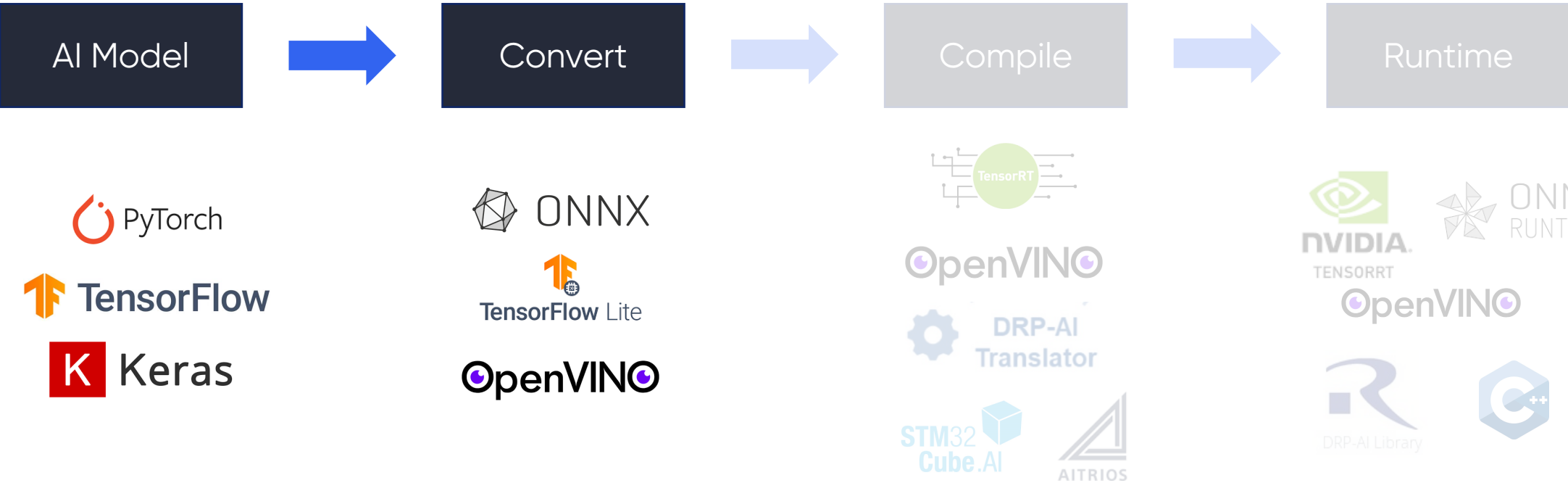
- Nota AI® supports a diverse range of devices such as CPUs, GPUs, NPUs, MCUs, etc.



# Problem statement

## Runtime-specific IR Converter

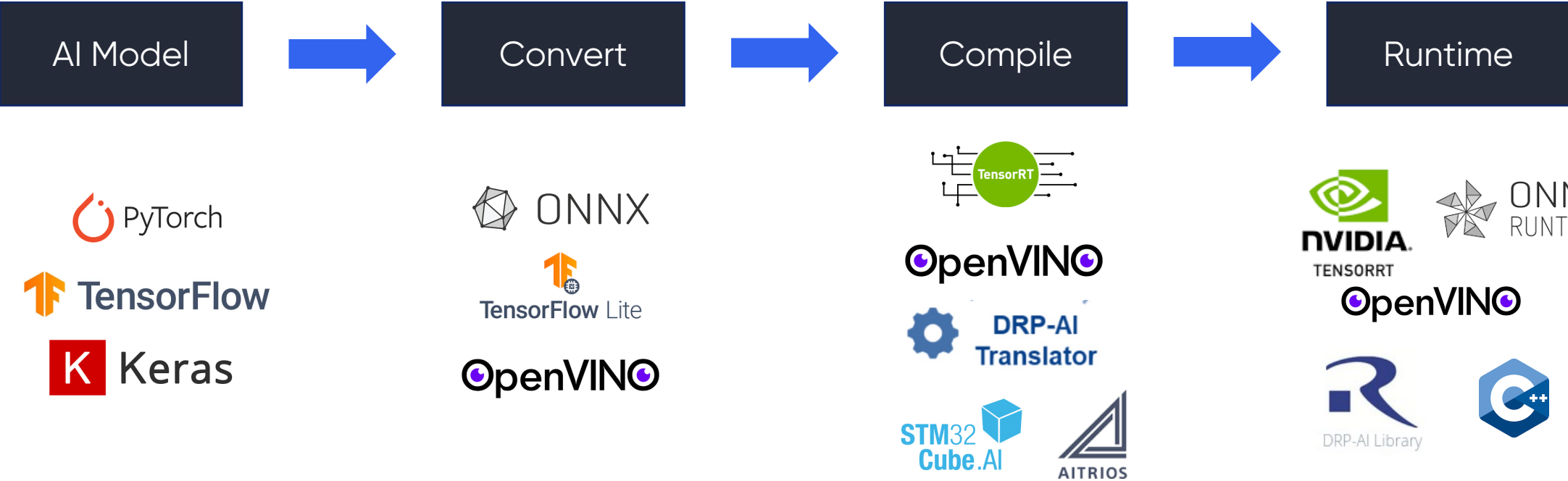
- To deploy AI models onto various devices, converting and compiling the AI model from deep learning frameworks into a runtime-specific Intermediate Representation(IR) is necessary.



# Problem statement

Limitation of compiler

- Each compiler can only support the operators that it has implemented.
- If an AI model contains unsupported operators, compilation may fail, or the operations may be executed on the CPU or MCU, resulting in reduced performance.

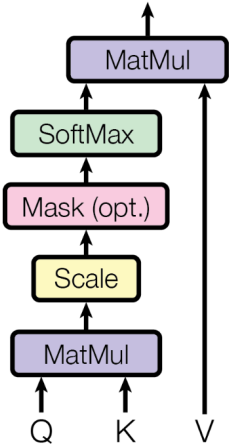


# Problem statement

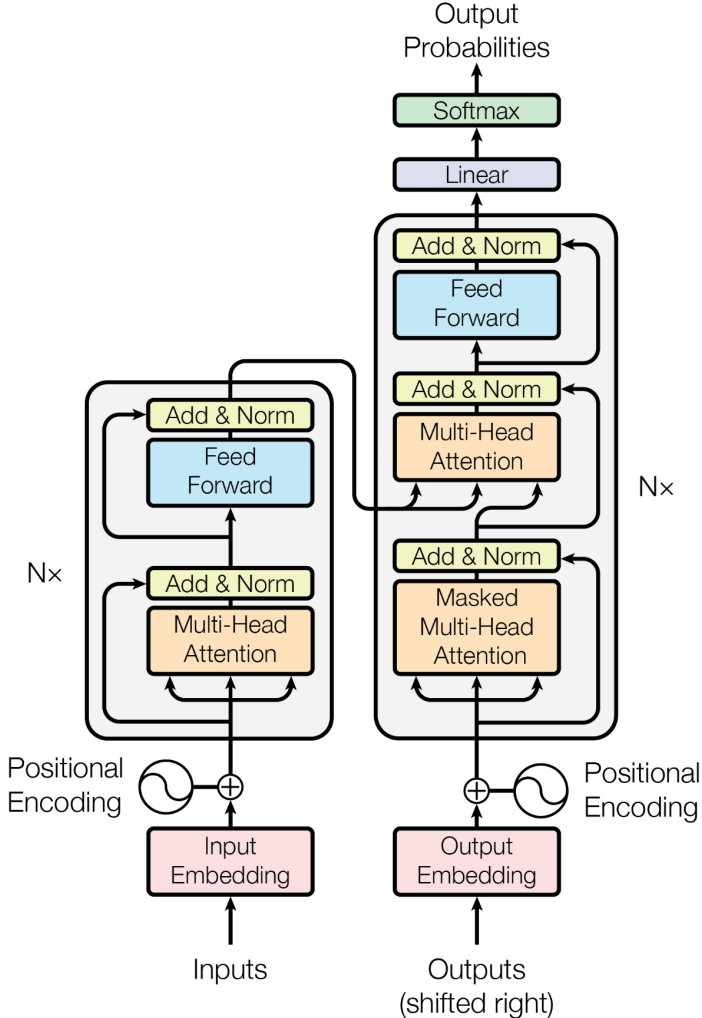
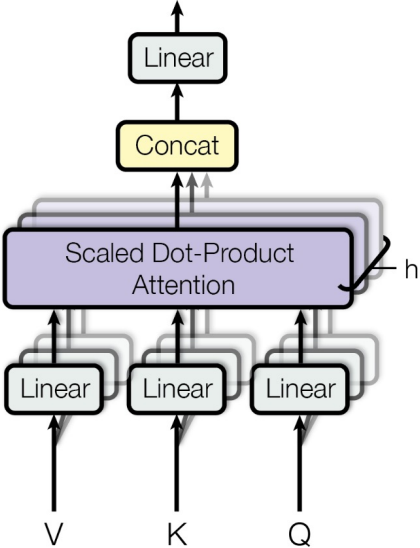
## Transformer

- Transformers have revolutionized deep learning, dominating natural language processing (NLP).
- Recent advancements have shown their remarkable performance in diverse domains, marking a paradigm shift in AI.

Scaled Dot-Product Attention



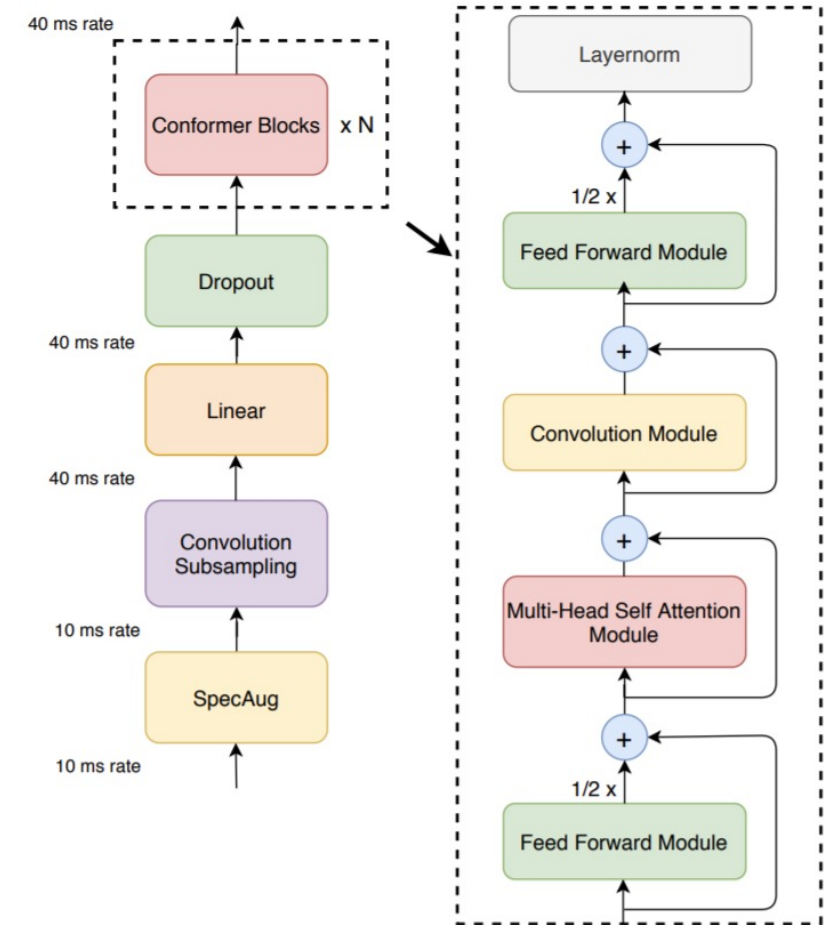
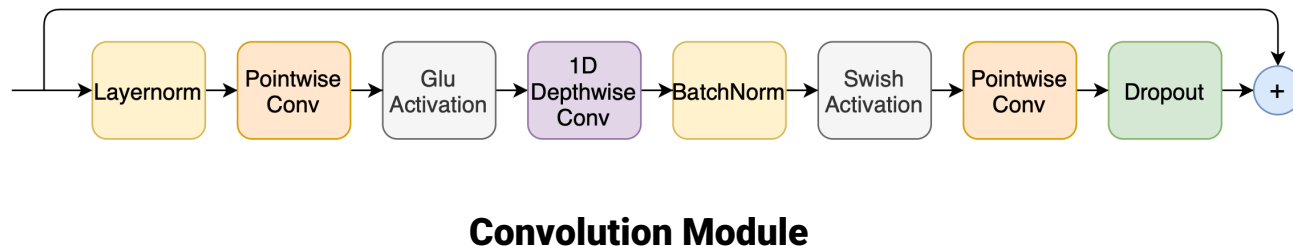
Multi-Head Attention



# Problem statement

## Conformer

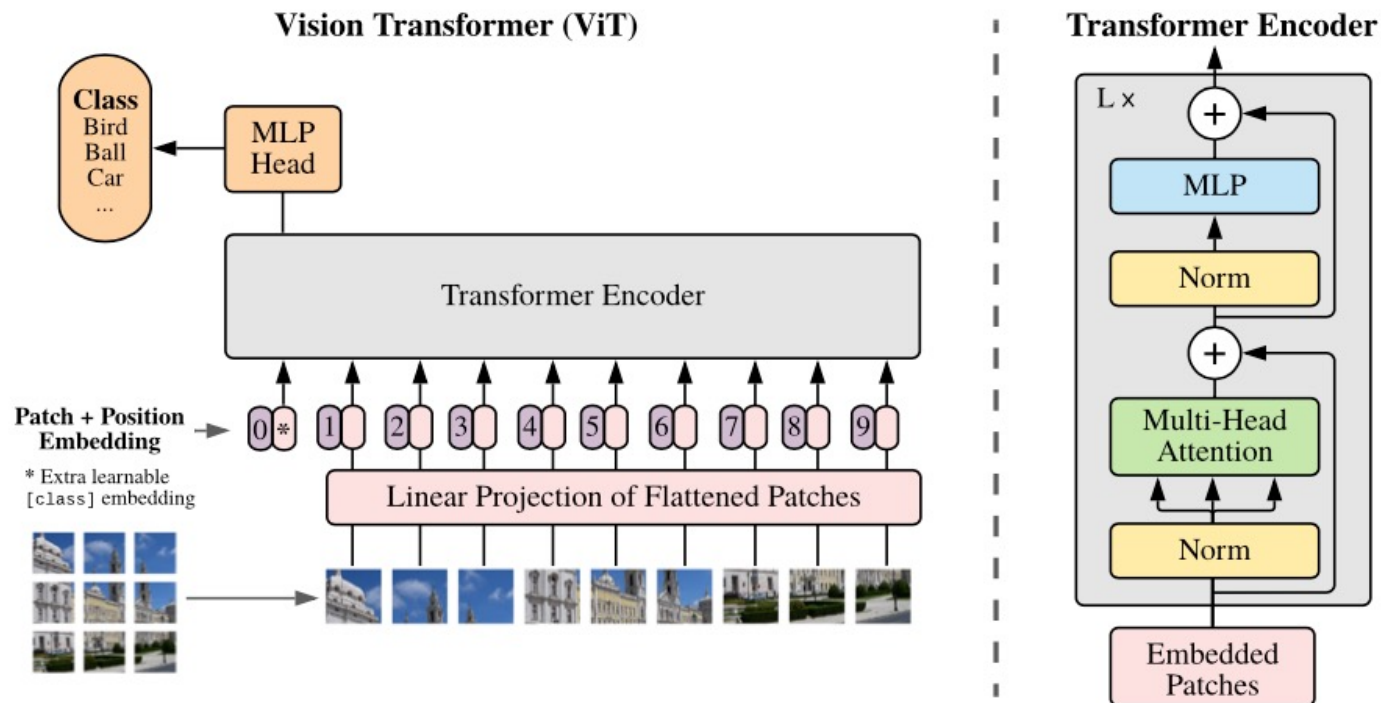
- Conformer enhances Transformers for speech and audio tasks, providing effective solutions in voice recognition and audio processing.
- Conformer combines CNNs with self-attention to improve speech-to-text accuracy, harnessing local and global data patterns.



# Problem statement

## Vision Transformer (ViT)

- Vision Transformer (ViT) expands Transformer's capabilities to computer vision tasks, achieving remarkable results in image classification and object detection.
- ViT processes images by dividing them into patches.



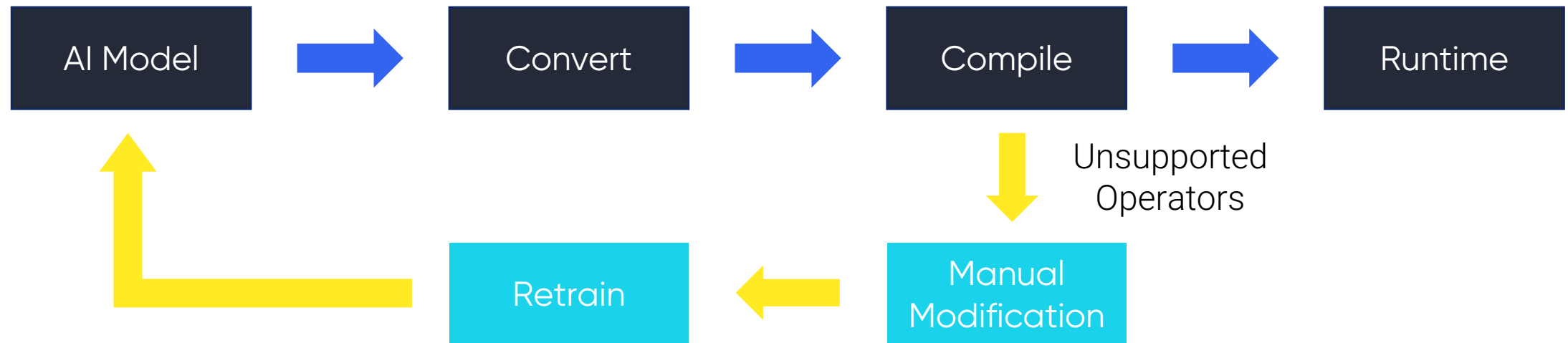
The background features a complex, abstract pattern of light blue lines and dots. The lines form a series of interconnected, irregular shapes that resemble a stylized map or a network diagram. The dots are placed at various points along these lines, creating a sense of movement and structure. The overall aesthetic is clean and modern, with a strong emphasis on geometric forms and a limited color palette.

# Operator Converter

# Operator Converter

## Challenges of Manual Model Modification

- AI models with unsupported operators require replacement or removal for deployment.
- Manual model modification in deep learning frameworks is cumbersome.

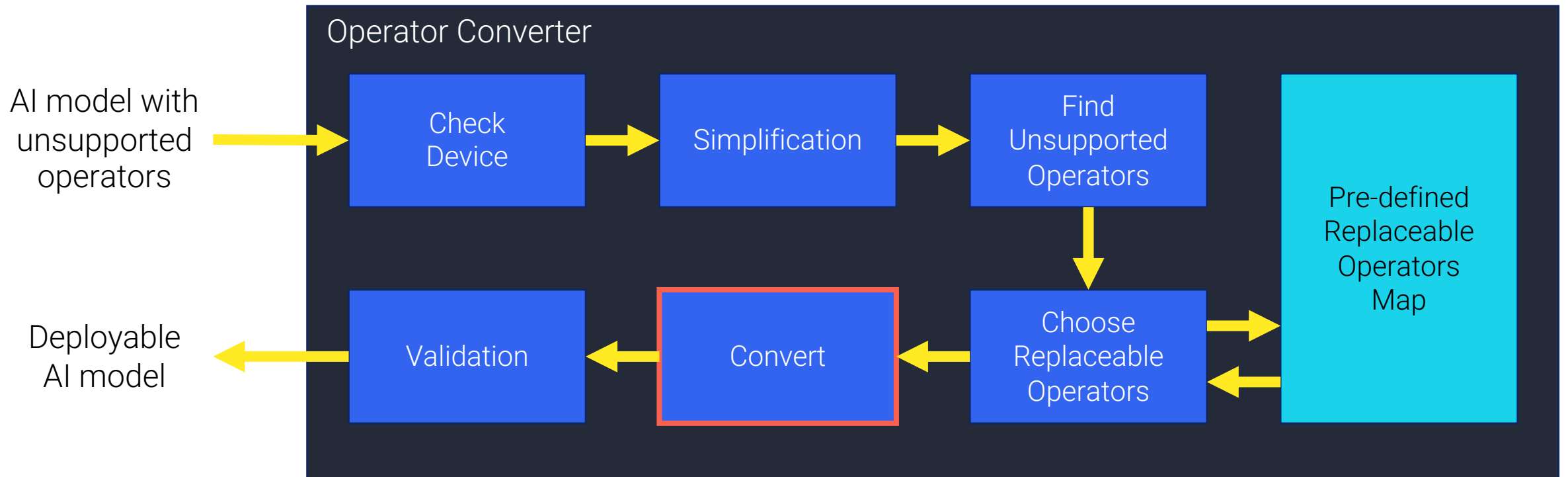




# Operator Converter

## Operator Converter

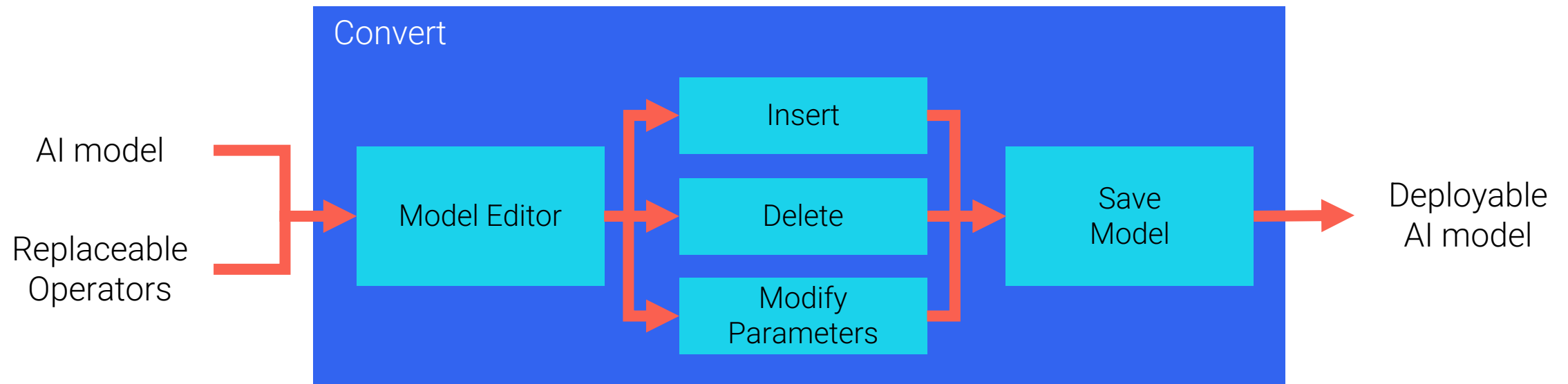
- The Operator Converter automates unsupported operation replacement, easing deployment.



# Operator Converter

Model Editor

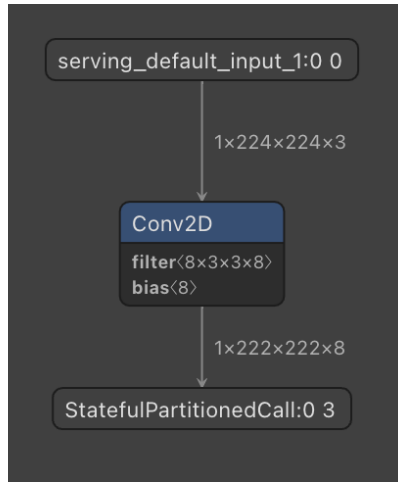
- The convert functionality within the operator converter utilizes a model parser equipped with insert, delete, and modify functions to alter the graph.



# Operator Converter

Model Editor

- Example



Subtract layer



Original

The screenshot shows the Model Editor interface. On the left, a diagram of the model structure is visible, with the second Conv2D layer highlighted. On the right, the 'NODE PROPERTIES' panel is open for this layer. The 'type' is 'Conv2D' and the 'location' is '1'. The 'ATTRIBUTES' section shows 'dilation\_h\_factor: 1', 'dilation\_w\_factor: 1', 'padding: VALID', 'stride\_h: 1', and 'stride\_w: 1'. The 'INPUTS' section shows the 'input' and 'filter' tensors. The 'filter' tensor is highlighted with a red box, showing its quantization parameters: 'quantization: linear' and a list of 8 values: 0: 0.001570379827171564 \* q, 1: 0.001570379827171564 \* q, 2: 0.001570379827171564 \* q, 3: 0.001570379827171564 \* q, 4: 0.001570379827171564 \* q, 5: 0.001570379827171564 \* q, 6: 0.00153042480815202 \* q, 7: 0.0016055121086537838 \* q.

Modify quantization params

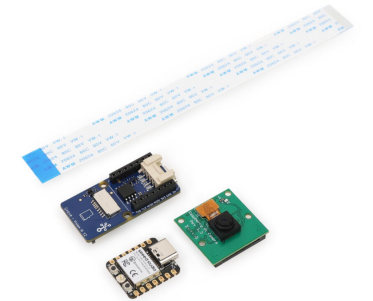
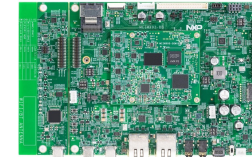
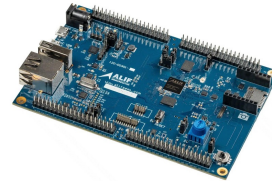


# Results

# Results

Target Device

- Alif Ensemble E7 DevKit Gen2
  - Arm Cortex-M55 + Arm Ethos-U55 NPU
- NXP i.MX 93 Evaluation Kit
  - Arm Cortex-M33 + Arm Ethos-U65 NPU
- Seeed Studio Grove Vision AI V2
  - Arm Cortex-M55 + Arm Ethos-U55 NPU



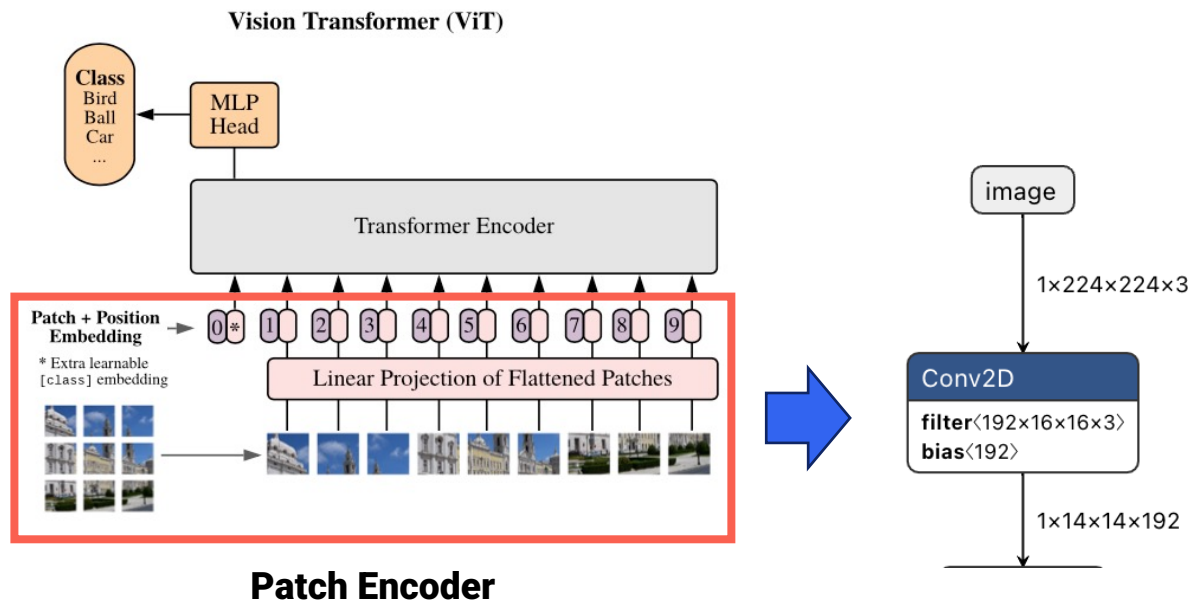
- Convert the TensorFlow Lite model with full INT8 quantization using the Vela compiler



# Results

## Constraints with Vela Compiler

- The patch encoder in ViT uses a convolution with kernel size and stride matching the patch size.
- This particular convolution configuration is not currently supported by the Vela compiler.



### TFLite CONV\_2D Constraints

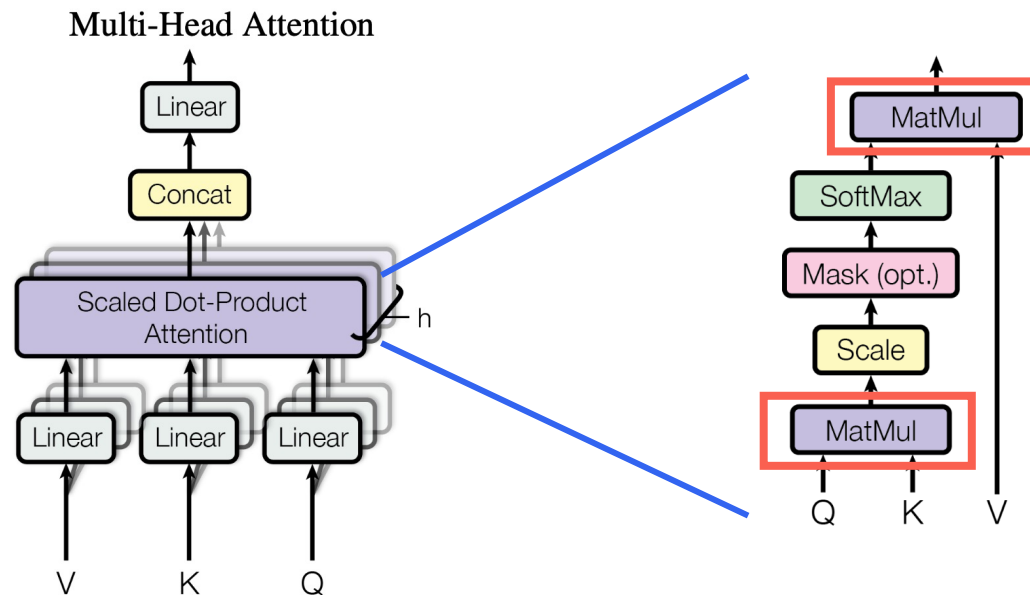
This is a list of constraints that the CONV\_2D operator must satisfy in order to be scheduled on the NPU.

- Stride values for both width and height must be integer types
- IFM depth must be a whole multiple of the filter kernel depth
- Number of filter kernels must be equally divisible by the number of convolution groups
- Dilation factor values for both width and height must be integer types
- Strides must fulfil the following criteria:
  - Stride h must be between 1 and 3 when ofm height is greater than 1
  - Stride w must be between 1 and 3 when ofm height is greater than 1 or stride w must be divisible by 2 or 3 and ifm width must be divisible by stride\_w/2 or stride\_w/3
- Dilated kernel height must be in the range [1, 64]
- Product of dilated kernel width and height must be in the range [1, 4096]
- Weight tensor must be 8-bit
- Weight tensor must be constant
- The sum of the weights cannot exceed 8323072
- Optional Bias tensor must be of shape: 1D
- Optional Bias tensor must be of type: int32, int64
- Optional Bias tensor values must fit within 40-bits

# Results

## Constraints with Vela Compiler

- MatMul in Multi-Head Attention is not supported by the Vela compiler.
- Addressing NPU compatibility issues with Transformer blocks by representing unsupported operators with supported operator combinations.

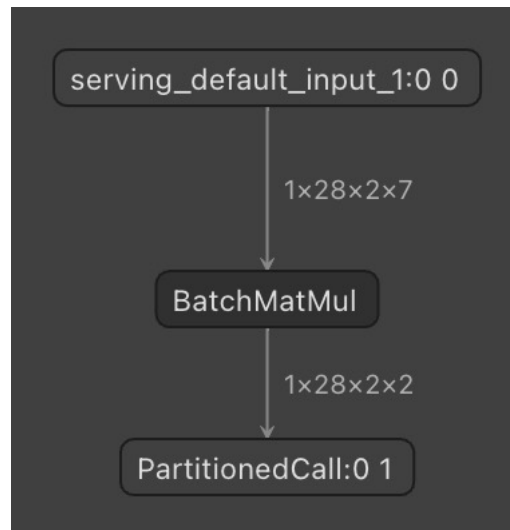


Operator	TFLite Constraints
ABS	Generic, Specific
ADD	Generic, Specific
ARG_MAX	Generic, Specific
AVERAGE_POOL_2D	Generic, Specific
CONCATENATION	Generic, Specific
CONV_2D	Generic, Specific
DEPTHWISE_CONV_2D	Generic, Specific
EXP	Generic, Specific
EXPAND_DIMS	Generic, Specific
FULLY_CONNECTED	Generic, Specific
HARD_SWISH	Generic, Specific
LEAKY_RELU	Generic, Specific
LOGISTIC	Generic
MAXIMUM	Generic, Specific
MAX_POOL_2D	Generic, Specific
MEAN	Generic, Specific
MINIMUM	Generic, Specific
MIRROR_PAD	Generic, Specific
MUL	Generic, Specific

# Results

## Constraints with Vela Compiler

### BatchMatmul



### Unsupported case in Ethos-U NPU

```
Warning: Using internal-default values for system configuration
Warning: Using internal-default values for memory mode
Warning: Unsupported TensorFlow Lite semantics for BATCH_MATMUL 'PartitionedCall:0'. Placing on CPU instead
- Input(s), Output and Weight tensors must have quantization parameters
  Op has tensors with missing quantization parameters: PartitionedCall:0

Network summary for only_matmul
Accelerator configuration      Ethos_U55_256
System configuration          internal-default
Memory mode                   internal-default
Accelerator clock              500 MHz

CPU operators = 1 (100.0%)
NPU operators = 0 (0.0%)

Neural network macs          0 MACs/batch
Network Tops/s               nan Tops/s

NPU cycles                   0 cycles/batch
SRAM Access cycles           0 cycles/batch
DRAM Access cycles           0 cycles/batch
On-chip Flash Access cycles  0 cycles/batch
Off-chip Flash Access cycles  0 cycles/batch
Total cycles                  0 cycles/batch

Batch Inference time         0.00 ms,    nan inferences/s (batch size 1)
```

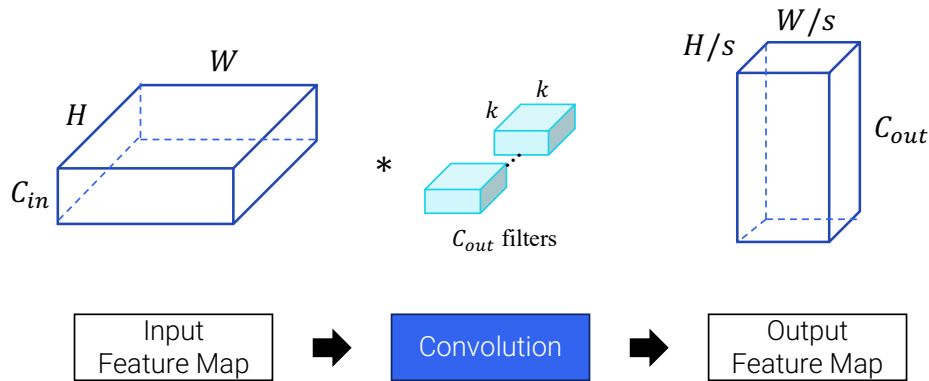


# Results

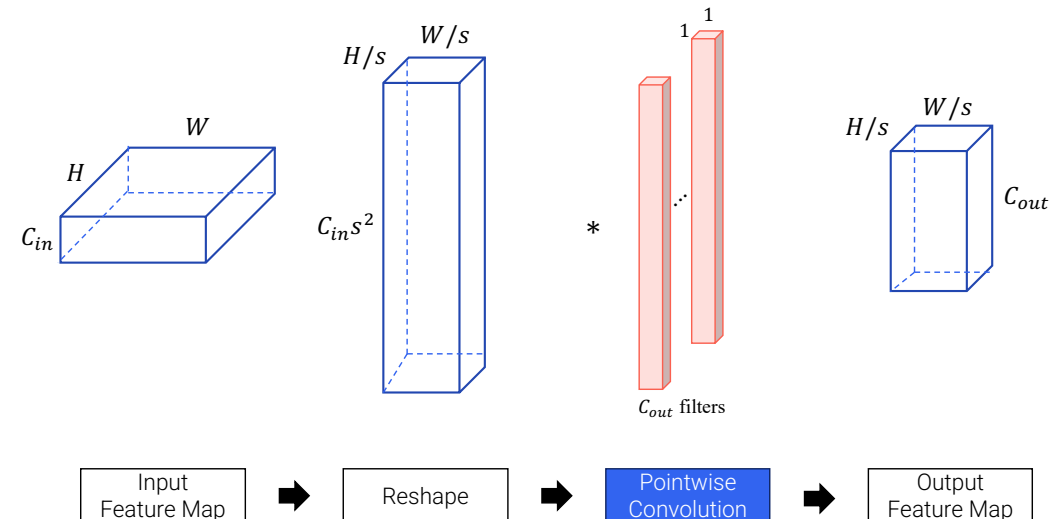
## Operator Converter for Patch Encoder

- Patch Encoder can be expressed using combinations of supported operators such as Slice, Reshape, and Pointwise Convolution.

### Without Operator Converter



### With Operator Converter

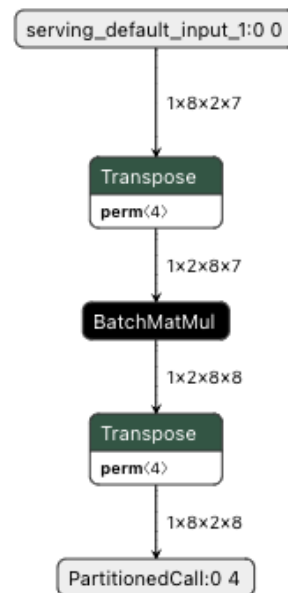


# Results

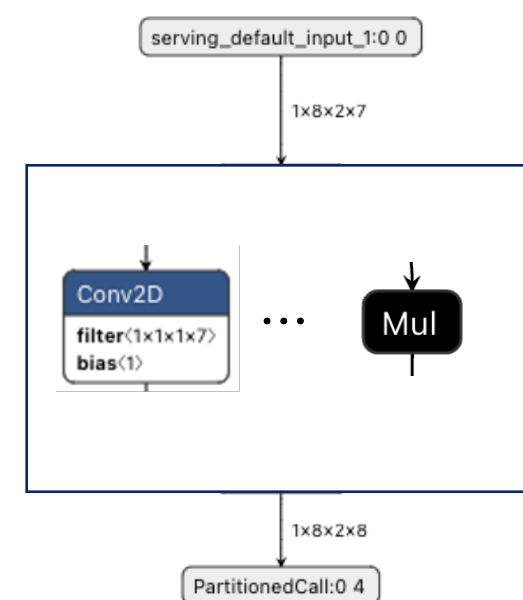
## Operator Converter for BatchMatMul

- Batchmatmul can be expressed using combinations of supported operators such as Convolution, Multiplication, and so on.
- Quantization values have also been appropriately filled in.

### Without Operator Converter



### With Operator Converter

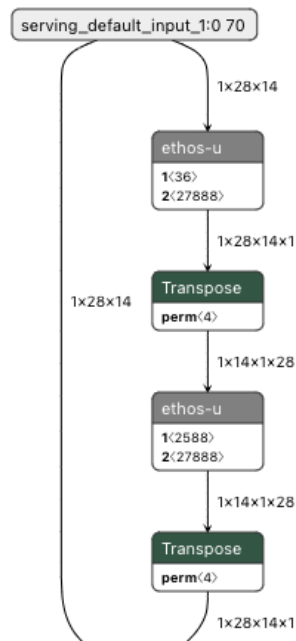


# Results

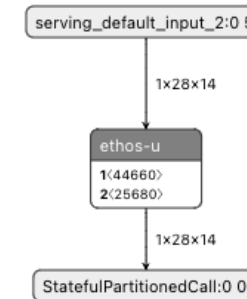
## Optimizing NPU Utilization

- After the operator converter, the model consists solely of supported operators, condensing into a single layer through the Vela compiler.
- This setup facilitates full NPU acceleration.

### Without Operator Converter



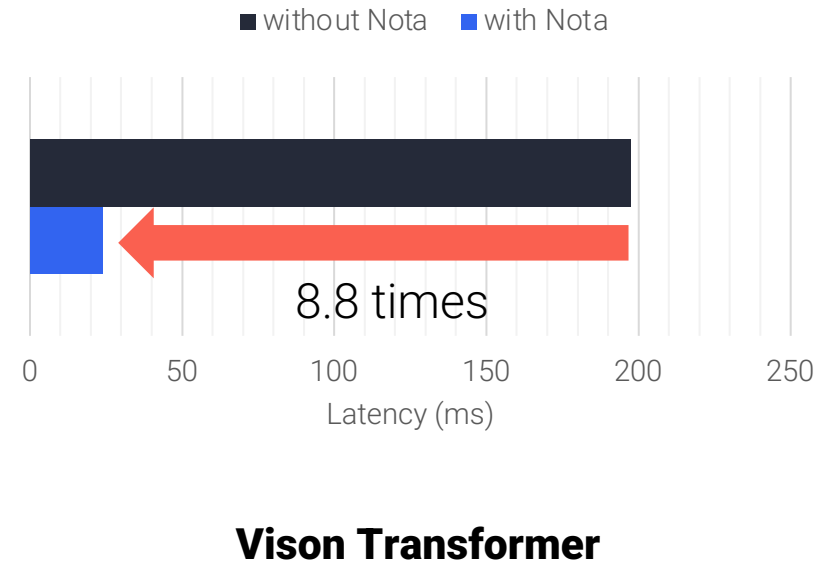
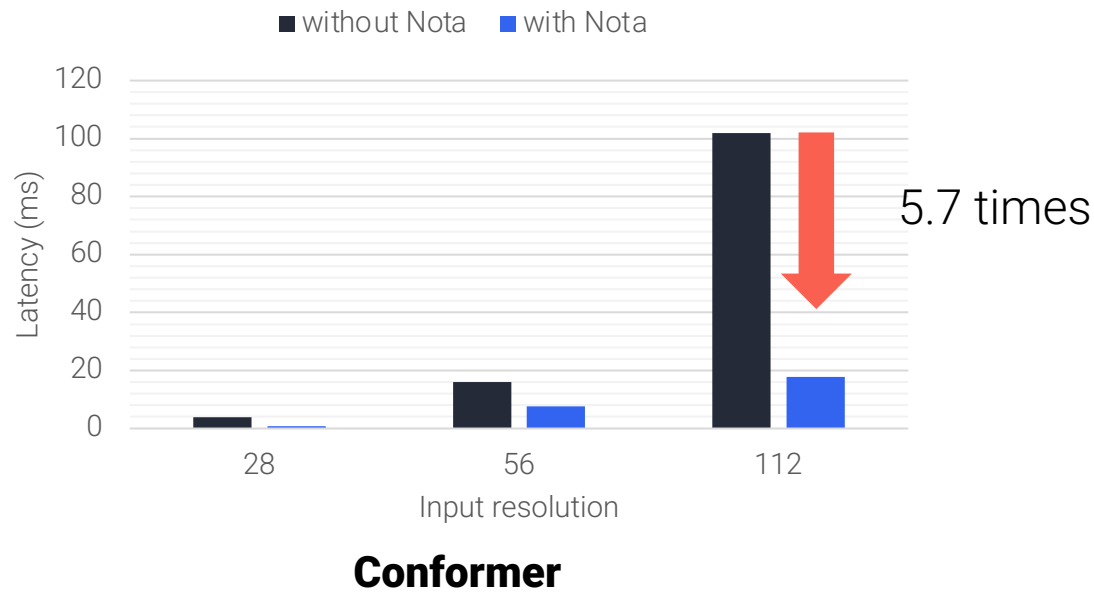
### With Operator Converter



# Results

## Benchmark

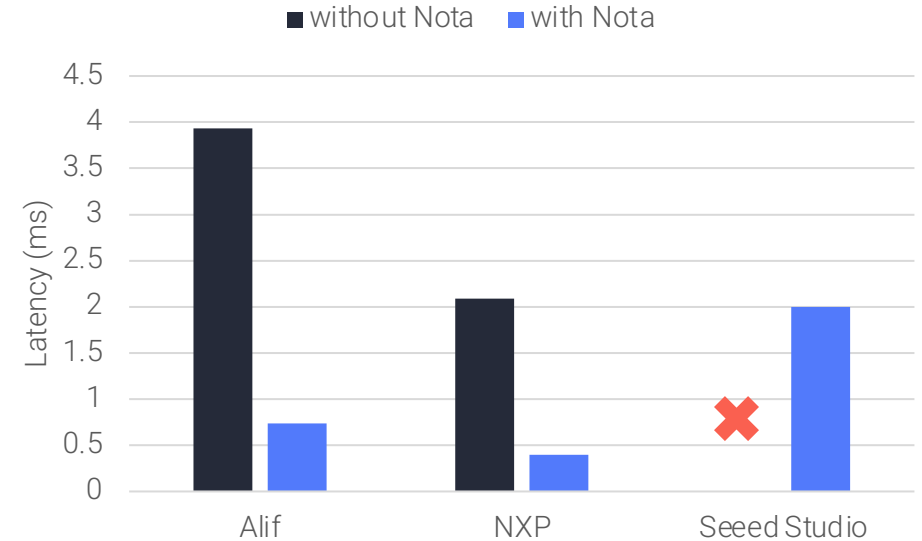
- Ethos-U NPUs don't currently support Conformer & Vision Transformer blocks.
- Ran AI model with Nota AI<sup>®</sup> using only Ethos-U NPUs
- On Alif Ensemble E7 DevKit Gen2, Conformer is 5.7 times faster, and ViT is 8.8 times faster



# Results

## Benchmark

- All devices with Arm Ethos NPUs (Alif, NXP, and Seeed Studio) demonstrated accelerated model inference speeds with Nota AI®.
- Using only NPUs, the Conformer model achieves notable performance improvements on each device.



Block	Type	Latency (ms)		
		Alif Ensemble E7 DevKit Gen2	NXP i.MX 93 Evaluation Kit	Seeed Studio Grove Vision AI V2
Conformer	Original	3.93	2.09	X
	with Nota	0.73	0.40	2.00



# Copyright Notice

This presentation in this publication was presented at the tinyML<sup>®</sup> Summit 2024. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

**[www.tinyml.org](http://www.tinyml.org)**