# Converter between formats of Deep Neural Network models on mobile platforms

**Martin Pavella**, xpavel39@stud.fit.vutbr.cz
supervisor: Ing. **Radek Kočí**, PhD.
consultant: Ing. **Róbert Kalmár** of NXP

## 1. The need for conversion

- High popularity and availability of Deep Neural Network (DNN) models in the **ONNX** format
- Superior HW accelerator support on mobile platforms for models in the **TFLite** format
- Existing converters produce sub-optimal models with unnecessary operators due to indirect approach

## 2. Proposed solution

A **direct** converter from ONNX to TFLite

- Represent an ONNX model using a hierarchy of objects
- Convert it to an equivalent TFLite object model (Fig.1)
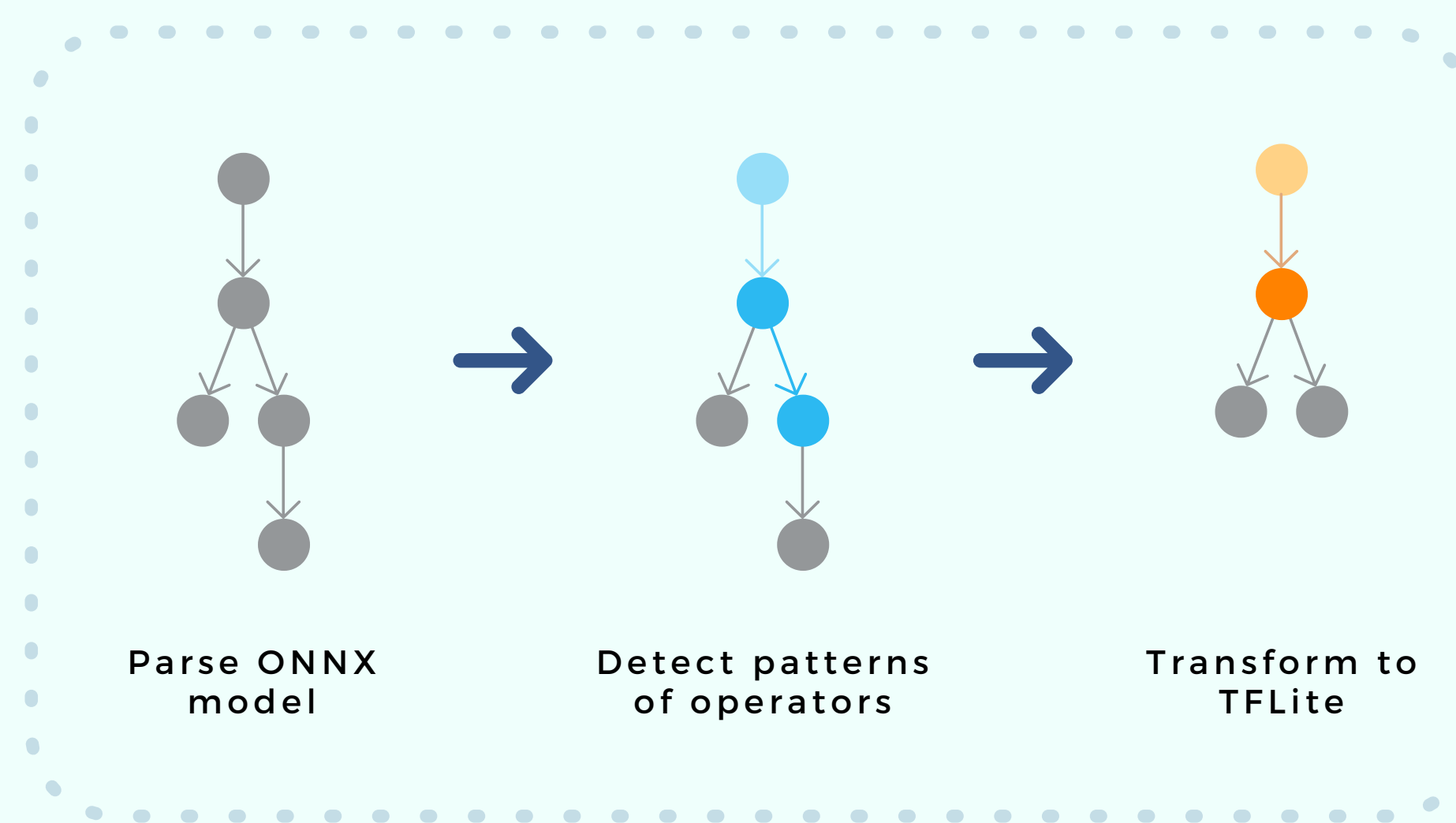- Serialize the model to the output TFLite file

Parse ONNX model → Detect patterns of operators → Transform to TFLite

**Fig. 1**
Process of model conversion

## 3. Results

- Operator conversion is a complex and evolving problem
- Conversion of all operators is not feasible -> Focus on a subset of commonly used operators
- Successful conversion of models used for classification, object detection, segmentation and analysis of acoustic data
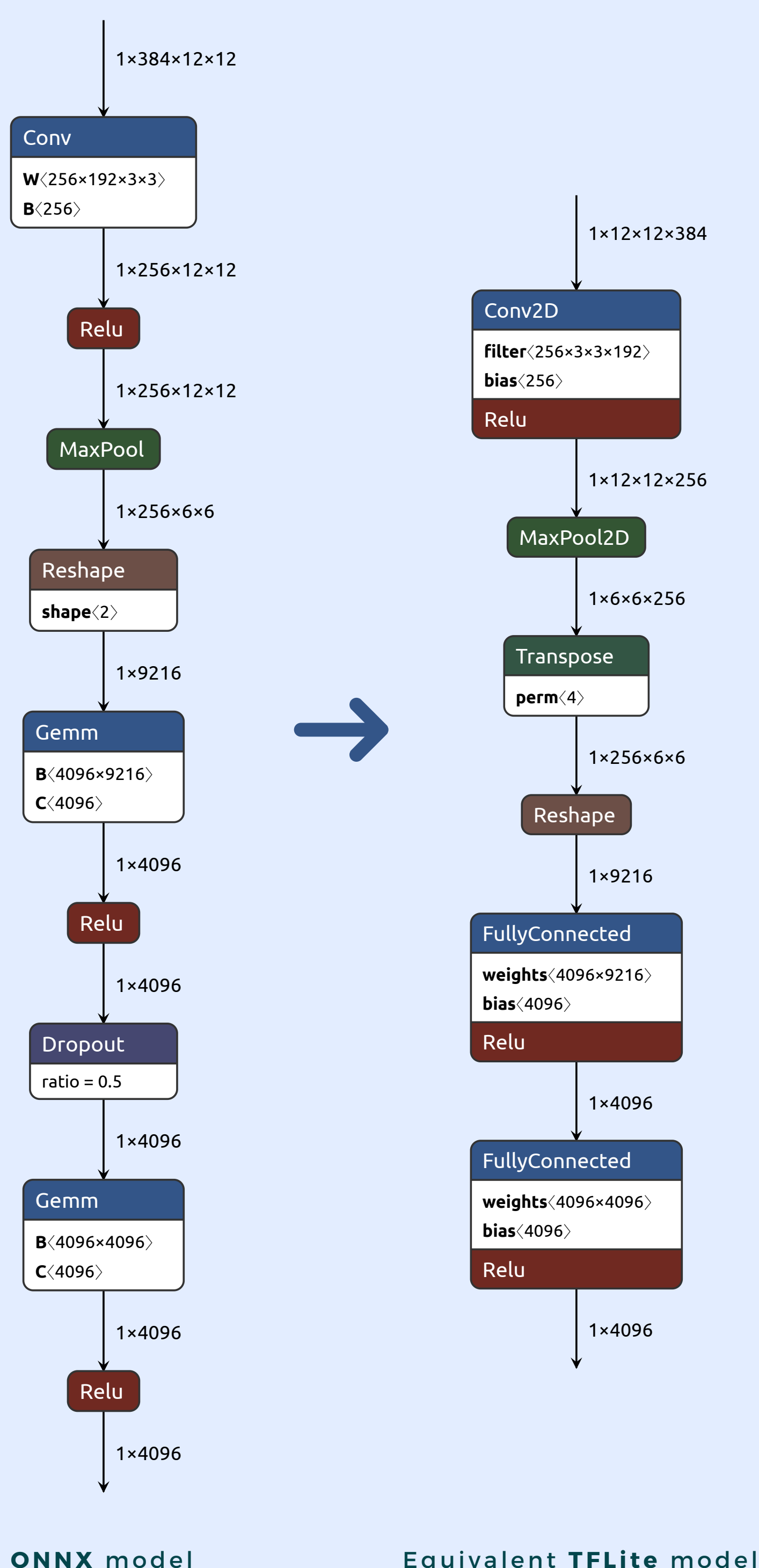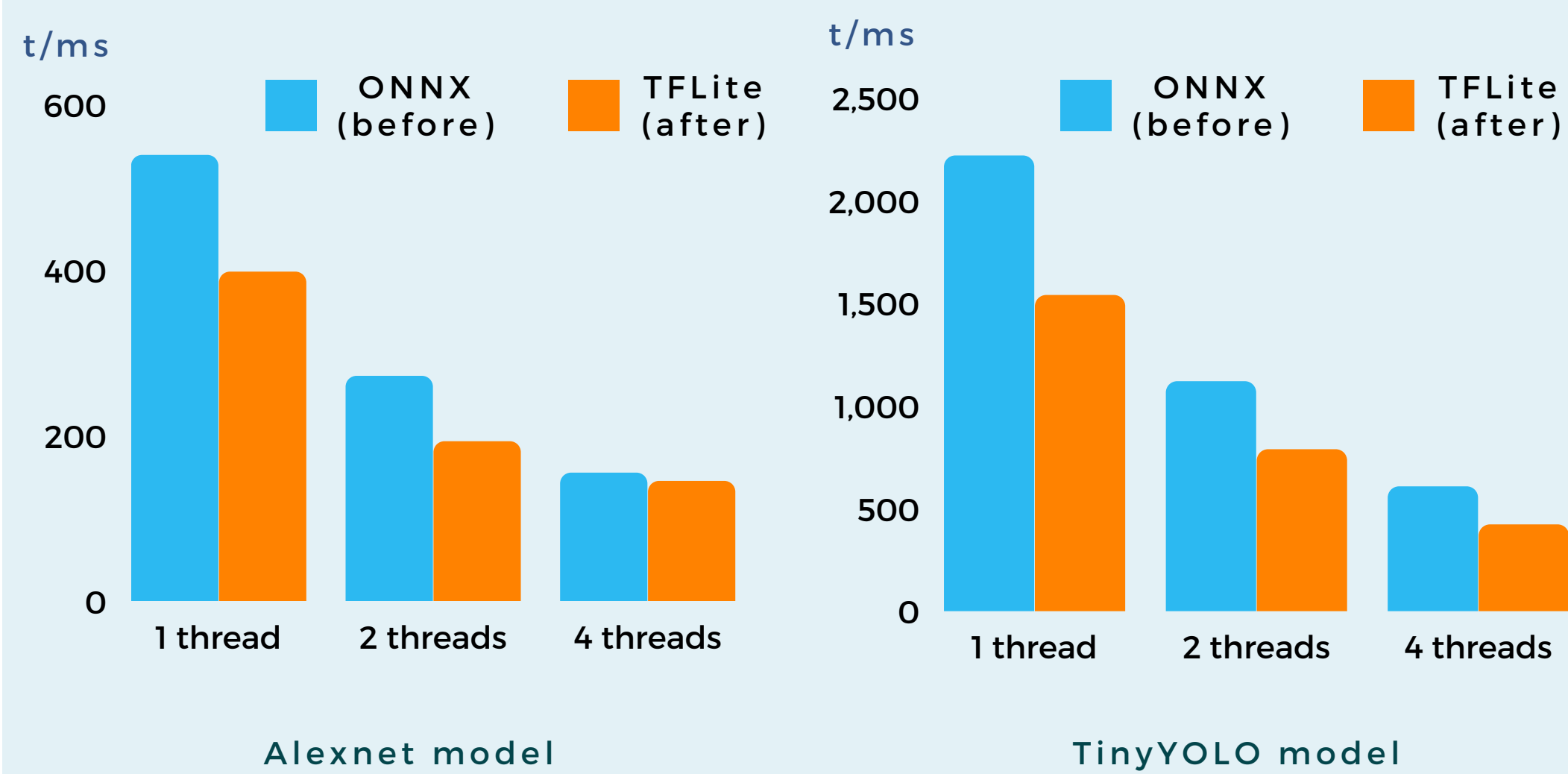- Model **size reduction** by up to 420kB

**ONNX model**

1×384×12×12
Conv
**W**⟨256×192×3×3⟩
**B**⟨256⟩
1×256×12×12
Relu
1×256×12×12
MaxPool
1×256×6×6
Reshape
**shape**⟨2⟩
1×9216
Gemm
**B**⟨4096×9216⟩
**C**⟨4096⟩
1×4096
Relu
1×4096
Dropout
ratio = 0.5
1×4096
Gemm
**B**⟨4096×4096⟩
**C**⟨4096⟩
1×4096
Relu
1×4096

**Equivalent TFLite model**

1×12×12×384
Conv2D
**filter**⟨256×3×3×192⟩
**bias**⟨256⟩
Relu
1×12×12×256
MaxPool2D
1×6×6×256
Transpose
**perm**⟨4⟩
1×256×6×6
Reshape
1×9216
FullyConnected
**weights**⟨4096×9216⟩
**bias**⟨4096⟩
Relu
1×4096
FullyConnected
**weights**⟨4096×4096⟩
**bias**⟨4096⟩
Relu
1×4096

**ONNX** model        Equivalent **TFLite** model

**Fig. 2**
Example of a section of a convolutional Alexnet model before and after conversion

## 4. Impact on inference

- Converted models produce identical outputs as the original ones
- Experiments in collaboration with the NXP company show a significant *improvement of inference speed* on target platforms

Alexnet model
t/ms — ONNX (before) / TFLite (after)
600, 400, 200, 0
1 thread, 2 threads, 4 threads

TinyYOLO model
t/ms — ONNX (before) / TFLite (after)
2,500, 2,000, 1,500, 1,000, 500, 0
1 thread, 2 threads, 4 threads

**Tab. 1**
The time duration of DNN model inference on target platforms before and after conversion

## 5. Limitations

- Limited subset of supported operators
- Conversion is not always possible
- Some models can be converted, but are not efficient