



AN02008: Reducing Memory Pressure in Neural Networks using Operator Splitting

Publication Date: 2026/4/8

Document Number: XM-015073-AN v1.0.0

IN THIS DOCUMENT

1	Prerequisites	1
2	Introduction: the Size of a Neural Network	1
3	Analyzing the Memory Requirements of MobileNet-V2	2
4	Making the Model Use Less Memory	4
5	Running the Model on Hardware	7
6	Summary	7

This app-note discusses running neural networks on XCORE devices that require more space for tensors than is available on the device. The `xcore-ai-tools` provide a method for dealing with this, called *Operator Splitting*, which reduces the size of intermediate data stored without affecting the outcome of the computation.

1 Prerequisites

This document assumes that you have installed the `xmos-ai-tools` package and are familiar with its basic usage.

To install the required Python dependencies, run the following command in your project directory:

```
pip install -r requirements.txt
```

2 Introduction: the Size of a Neural Network

The size of a neural network typically encompasses two variables:

- ▶ The amount of memory required to store learned parameters
- ▶ The amount of memory required to store intermediate results

The former is a number that is often mentioned in papers on neural networks. For example, `resnet50` has 25,583,592 learned parameters. Note that they are also known as *weights* or *trainable parameters*. Each learned parameter typically takes up a byte of space (assuming 8-bit quantisation), therefore `ResNet50` requires approximately 25 MByte of space for learned parameters.

The second number is the amount of memory that is used to store data whilst it is processing the network. The network is typically organised as a sequence of operators, where each operator consumes and produces a tensor (a multi-dimensional vector). The input tensor may, for example, be 2D RGB image or a 1D sound sample. Intermediate tensor represents features that the Neural Network algorithm has inferred, and the final output is the desired classification or embedding.

The intermediate values are known as the activations. Unlike learned parameters, the activations are read and written and therefore need to be stored in a RAM memory. The maximum number of activations that need to be stored at any one time has to be less



than the amount of RAM memory available. Unlike weights, the number is dependent on the implementation of the neural network engine, so this is not a number that is typically published. As a first proxy one can use the largest tensor that is produced in the computation, but any forward connections in the network need to be added to that.

Given the an implementation of the neural network, one can calculate the memory required, and indeed the XMOS AI Tools can output the size required and a diagram depicting which part of the network requires this amount of memory.

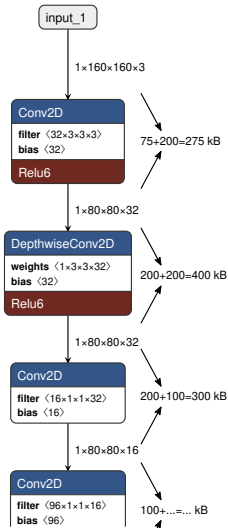


Fig. 1: First few layers of mobilenet V2

As an example we use a well known neural network that requires a large amount of memory: MobileNet-V2. It may be old and out of date, but it is well-known. We can look at a representation of the first few layers of the model as shown in [First few layers of mobilenet V2](#). The graph (drawn by Netron) is annotated with the size of each tensor, and the total size of tensors going in and coming out of an operator. Assuming that the two cannot be overlapped, we can add these together to get a first idea how much memory will be required at least.

Looking a bit further down in the network we highlight the layers that include the first bypass in [First bypass in mobilenet V2](#). Like before we have annotated the graph with the size of each tensor, and the total size of tensors being alive during a straightforward top-down execution of the graph (assuming zero overlap). The colours link the size calculation to each of the tensors.

3 Analyzing the Memory Requirements of MobileNet-V2

The particular flavour of MobileNet that we use has an input image of 160x160 and 1280 features in the last layer (an alpha of 1.0). We've chosen this network as it is too large to fit into memory. First we run the xcore-optimization tool on the file to establish what size memory is required:

```
% cd model
% xcore-opt -xp model --xcore-conv-err-threshold=2.1 mobilenetv2.tflite -o model.tflite -f params.flash
<unknown>:0: remark: Tensor arena size : 798120
%
% cd model
% xcore-opt -xp model --xcore-conv-err-threshold=2.1 mobilenetv2.tflite -o model.tflite -f params.flash
<unknown>:0: remark: Tensor arena size : 798120
%
```



and amount of memory are minimal for the dual-split, and execution time is optimal for the triple-split.

The value **10** that the xformer suggested is optimal for the single split, but when applying multiple splits it is worthwhile to cut them into fewer chunks. That reduces the overhead of the splitting in both memory and time. This is shown in the final two columns where we apply fewer splits. In this case, the triple split is the most efficient in terms of both memory and time.

Note that the execution time in this example includes the time taken to load coefficients from flash which accounts for 72 ms. This is due to the large number of parameters. As many of the parameters are used only once during the final fully connected layer, there is little scope to hide these times amongst compute. In typical networks, the final fully connected layer will be smaller reducing the parameters count and total execution time.

An excerpt of the graph produced in the right-hand column (excluding loading from flash and only the first few operators) is shown in *The split graph; the first part is split eight times, the second part is split four times, the third part if split twice, below that there are no splits.*

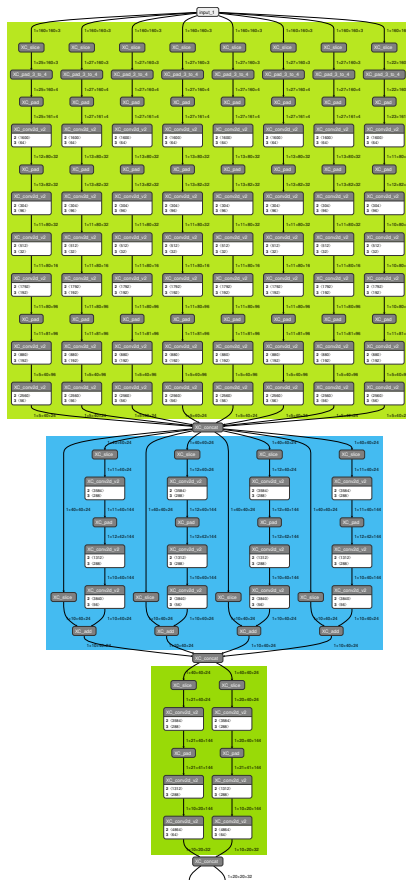


Fig. 3: The split graph; the first part is split eight times, the second part is split four times, the third part if split twice, below that there are no splits.



5 Running the Model on Hardware

This Application Note provides a source code to run the model on hardware. The required HW is the following:

- ▶ 1x *xcore.ai Explorer Development Kit* (XK-EVK-XU316).
- ▶ 2x Micro USB cable (for Power supply and xTag).

This application includes a CMake configuration that automates the build process. The CMake file sets up a Python virtual environment, installs all necessary Python packages and C libraries, and converts the TFLite model into a C++ header file for integration. Additional files required for linking and compiling the application with the XMOS AI tools are also generated automatically.

To compile and run the application on the XK-EVK-XU316, follow these steps:

```
# compile
cmake -G "Unix Makefiles" -B build
xmake -C build
# flash
xflash --target XK-EVK-XU316 --data app_an02008/src/weights.bin
# run
xrun --xscope app_an02008/bin/app_an02008.xe
```

After running the above commands, the model will run on the XK-EVK-XU316 board and a profile of each layer will be printed to the console. The output will include the time taken for each layer, the number of operations, and the memory used by each layer. Here is an example of the output:

```
Cumulative times for invoke()...
20 OP_XC_slice           18471      0.18ms
10 OP_XC_pad_3_to_4    155816     1.56ms
63 OP_XC_pad           281739     2.82ms
86 OP_XC_ld_weights    7120681    71.21ms
182 OP_XC_conv2d_v2    9134347    91.34ms
19 OP_XC_add           53282      0.53ms
5 OP_XC_concat         6913       0.07ms
1 OP_XC_mean          203804     2.04ms
1 OP_XC_no_op         23         0.00ms
1 OP_XC_softmax       11031      0.11ms

Total time for invoke() - 16986107 169.86ms
```

6 Summary

This app-note gives a brief example on how to use “operator splitting” to make larger networks fit on an xcore.



Copyright © 2026, All Rights Reserved.

XMOS Ltd. is the owner or licensee of this design, code, or Information (collectively, the “Information”) and is providing it to you “AS IS” with no warranty of any kind, express or implied and shall have no liability in relation to its use. XMOS Ltd makes no representation that the Information, or any particular implementation thereof, is or will be free from any claims of infringement and again, shall have no liability in relation to any such claims.

XMOS, XCORE, VocalFusion and the XMOS logo are registered trademarks of XMOS Ltd. in the United Kingdom and other countries and may not be used without written permission. Company and product names mentioned in this document are the trademarks or registered trademarks of their respective owners.

