# TensorFlow Lite to ONNX Conversion
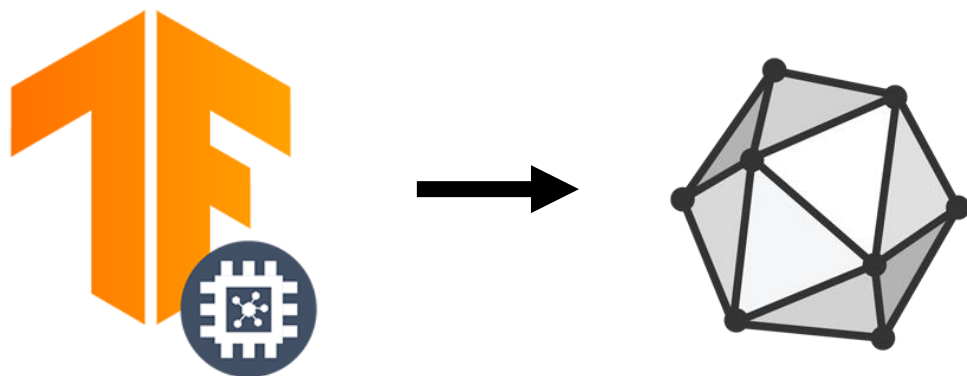
# ONNX Runtime Mobile

TOM WILDENHAIN

SCOTT MCKAY

MICROSOFT

# TensorFlow Lite to ONNX

TOM WILDENHAIN       SOFTWARE ENGINEER, ONNX CONVERTERS TEAM

MICROSOFT

# TensorFlow

- Used for inference and training
- > 1,000 ops
- Already have conversions for many ops

| | | | | |
|---|---|---|---|---|
| Abort | CountUpTo | Greater | NonSerializableDataset | ResourceApplyFtrl |
| Abs | CreateSummaryDbWriter | GreaterEqual | NotEqual | ResourceApplyFtrlV2 |
| AccumulateNV2 | CreateSummaryFileWriter | GroupByReducerDataset | NthElement | ResourceApplyGradientDescent |
| AccumulatorApplyGradient | CropAndResize | GroupByWindowDataset | OneHot | ResourceApplyKerasMomentum |
| AccumulatorNumAccumulated | CropAndResizeGradBoxes | GuaranteeConst | OneShotIterator | ResourceApplyMomentum |
| AccumulatorSetGlobalStep | CropAndResizeGradImage | HSVToRGB | OptimizeDataset | ResourceApplyPowerSign |
| AccumulatorTakeGradient | Cross | HashTable | OptimizeDatasetV2 | ResourceApplyProximalAdagrad |
| Acos | CrossReplicaSum | HashTableV2 | OptionalFromValue | ResourceApplyProximalGradientDe |
| Acosh | CudnnRNN | HistogramFixedWidth | OptionalGetValue | scent |
| Add | CudnnRNNBackprop | HistogramSummary | OptionalHasValue | ResourceApplyRMSProp |
| AddManySparseToTensorsMap | CudnnRNNBackpropV2 | IFFT | OptionalNone | ResourceConditionalAccumulator |
| AddN | CudnnRNNBackpropV3 | IFFT2D | OptionalNone | ResourceCountUpTo |
| AddSparseToTensorsMap | CudnnRNNCanonicalToParams | IFFT3D | OrderedMapClear | ResourceGather |
| AddV2 | CudnnRNNCanonicalToParamsV2 | IRFFT | OrderedMapIncompleteSize | ResourceGatherNd |
| AdjustContrast | CudnnRNNParamsSize | IRFFT2D | OrderedMapPeek | ResourceScatterAdd |
| AdjustContrastv2 | CudnnRNNParamsToCanonical | IRFFT3D | OrderedMapSize | ResourceScatterDiv |
| AdjustHue | CudnnRNNParamsToCanonicalV2 | Identity | OrderedMapStage | ResourceScatterMax |
| AdjustSaturation | CudnnRNNV2 | IdentityN | OrderedMapUnstage | ResourceScatterMin |
| All | CudnnRNNV3 | IdentityReader | OrderedMapUnstageNoKey | ResourceScatterMul |
| AllCandidateSampler | Cumprod | IdentityReaderV2 | OutfeedDequeue | ResourceScatterNdAdd |
| AllToAll | Cumsum | If | OutfeedDequeueTuple | ResourceScatterNdMax |
| Angle | CumulativeLogsumexp | Igamma | OutfeedDequeueTupleV2 | ResourceScatterNdMin |
| AnonymousIterator | DataFormatDimMap | IgammaGradA | OutfeedDequeueV2 | ResourceScatterNdSub |
| AnonymousIteratorV2 | DataFormatVecPermute | Igammac | OutfeedEnqueue | ResourceScatterNdUpdate |
| AnonymousMemoryCache | DataServiceDataset | IgnoreErrorsDataset | OutfeedEnqueueTuple | ResourceScatterSub |
| AnonymousMultiDeviceIterator | DatasetCardinality | Imag | Pack | ResourceScatterUpdate |
| AnonymousRandomSeedGenerator | DatasetFromGraph | ImageProjectiveTransformV2 | Pad | ResourceSparseApplyAdadelta |
| AnonymousSeedGenerator | DatasetToGraph | ImageProjectiveTransformV3 | PadV2 | ResourceSparseApplyAdagrad |
| Any | DatasetToGraphV2 | ImageSummary | PaddedBatchDataset | ResourceSparseApplyAdagradDA |
| ApplyAdaMax | DatasetToSingleElement | ImmutableConst | PaddedBatchDatasetV2 | ResourceSparseApplyAdagradV2 |
| ApplyAdadelta | DatasetToTFRecord | ImportEvent | PaddingFIFOQueue | ResourceSparseApplyCenteredRMS |
| ApplyAdagrad | Dawsn | InTopK | PaddingFIFOQueueV2 | Prop |
| ApplyAdagradDA | DebugGradientIdentity | InTopKV2 | ParallelConcat | ResourceSparseApplyFtrl |
| ApplyAdagradV2 | DebugGradientRefIdentity | InfeedDequeue | ParallelDynamicStitch | ResourceSparseApplyFtrlV2 |
| ApplyAdam | DebugIdentity | InfeedDequeueTuple | ParallelInterleaveDataset | ResourceSparseApplyKerasMomen |
| ApplyAddSign | DebugIdentityV2 | InfeedEnqueue | ParallelInterleaveDatasetV2 | tum |
| ApplyCenteredRMSProp | DebugNanCount | InfeedEnqueuePrelinearizedBuffer | ParallelInterleaveDatasetV3 | ResourceSparseApplyMomentum |
| ApplyFtrl | DebugNumericSummary | InfeedEnqueueTuple | ParallelInterleaveDatasetV4 | ResourceSparseApplyProximalAda |
| ApplyFtrlV2 | DebugNumericSummaryV2 | InitializeTable | ParallelMapDataset | grad |
| ApplyGradientDescent | DecodeAndCropJpeg | InitializeTableFromDataset | ParallelMapDatasetV2 | ResourceSparseApplyProximalGrad |
| ApplyMomentum | DecodeBase64 | InitializeTableFromTextFile | ParameterizedTruncatedNormal | ientDescent |
| ApplyPowerSign | DecodeBmp | InitializeTableFromTextFileV2 | ParseExample | ResourceSparseApplyRMSProp |
| ApplyProximalAdagrad | DecodeCSV | InitializeTableV2 | ParseExampleDataset | ResourceStridedSliceAssign |
| ApplyProximalGradientDescent | DecodeCompressed | InplaceAdd | ParseExampleDatasetV2 | Restore |
| ApplyRMSProp | DecodeGif | InplaceSub | ParseExampleV2 | RestoreSlice |
| ApproximateEqual | DecodeImage | InplaceUpdate | ParseSequenceExample | RestoreV2 |
| ArgMax | DecodeJSONExample | InterleaveDataset | ParseSequenceExampleV2 | RetrieveTPUEmbeddingADAMPara |
| ArgMin | DecodeJPEG | | | |

...

# TFLite

- Lightweight runtime used for inference
- ~130 ops
- Models created from TensorFlow

| | | |
|---|---|---|
| ABS | NEG | UNPACK |
| ADD_N | NON_MAX_SUPPRESSION_V4 | WHERE |
| ARG_MAX | NON_MAX_SUPPRESSION_V5 | WHILE |
| ARG_MIN | NOT_EQUAL | ZEROS_LIKE |
| AVERAGE_POOL_2D | ONE_HOT | FULLY_CONNECTED |
| BATCH_TO_SPACE_ND | PACK | ADD |
| CAST | PAD | DIV |
| CEIL | PADV2 | MUL |
| CONCATENATION | POW | SUB |
| CONV_2D | QUANTIZE | BATCH_MATMUL |
| COS | RANGE | BIDIRECTIONAL_SEQUENCE_LSTM |
| CUMSUM | RANK | BIDIRECTIONAL_SEQUENCE_RNN |
| DEPTH_TO_SPACE | REDUCE_ANY | BROADCAST_TO |
| DEPTHWISE_CONV_2D | REDUCE_MAX | CALL |
| DEQUANTIZE | REDUCE_PROD | CALL_ONCE |
| ELU | RELU | CONCAT_EMBEDDINGS |
| EQUAL | RELU6 | CUSTOM |
| EXP | RESHAPE | DELEGATE |
| EXPAND_DIMS | RESIZE_BILINEAR | DENSIFY |
| FILL | RESIZE_NEAREST_NEIGHBOR | EMBEDDING_LOOKUP |
| FLOOR | REVERSE_SEQUENCE | EMBEDDING_LOOKUP_SPARSE |
| FLOOR_DIV | REVERSE_V2 | FAKE_QUANT |
| FLOOR_MOD | ROUND | HARD_SWISH |
| GATHER | RSQRT | HASHTABLE_LOOKUP |

...

# TFLite Conversion
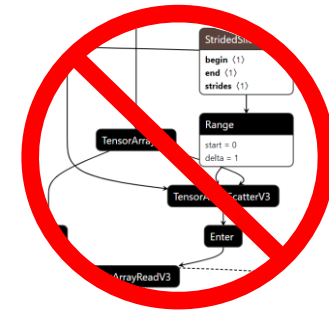
```
pip install tf2onnx
python -m tf2onnx.convert --tflite ssdmobilenet.tflite --output ssdmobilenet.onnx --opset 13
```



name: **normalized_input_image_tensor**

type: **uint8[1,300,300,3]**

quantization: **-1 ≤ 0.0078125 * (q - 128) ≤ 0.9921875**

location: **260**



Some models are only available for TFLite

Automatic quantization support!

TFLite models are often cleaner

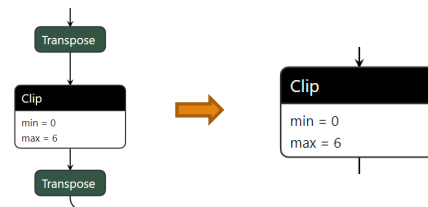# Conversion Process

1. Rewriters
   ◦ Convert op patterns

2. Handlers
   ◦ Convert individual ops

3. Optimizers
   ◦ Remove unnecessary ops

# Conversion Process

**1. Rewriters**
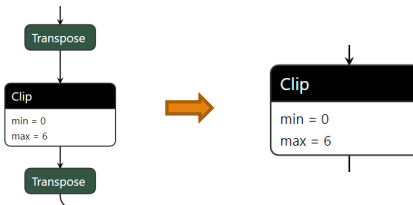- **Convert op patterns**

**2.** Handlers
- Convert individual ops
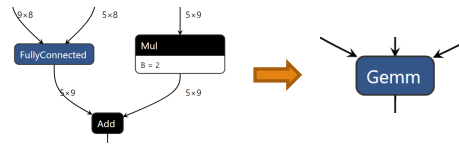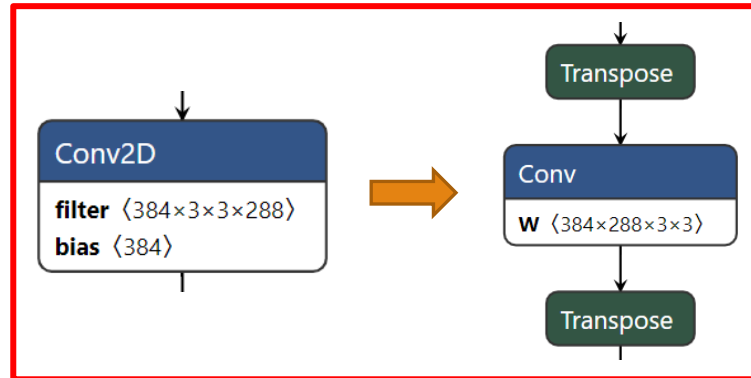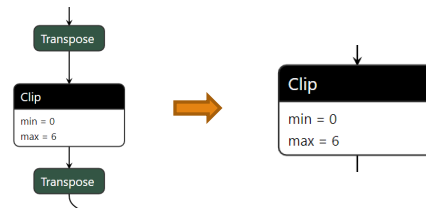
**3.** Optimizers
- Remove unnecessary ops

# Conversion Process

1. Rewriters
   - Convert op patterns

2. **Handlers**
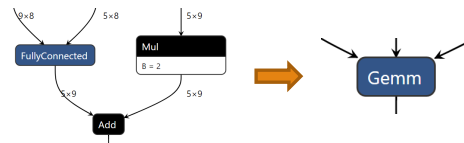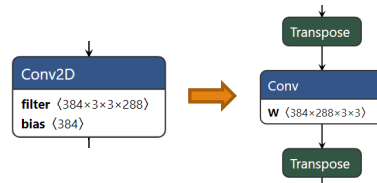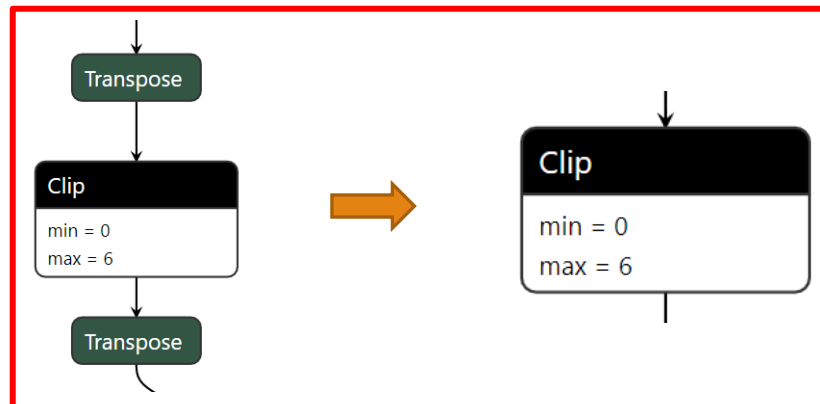   - **Convert individual ops**

3. Optimizers
   - Remove unnecessary ops

# Conversion Process

1. Rewriters
   ◦ Convert op patterns



2. Handlers
   ◦ Convert individual ops



3. **Optimizers**
   ◦ **Remove unnecessary ops**

# Quantization



tf2onnx adds q/dq

ORT recognizes op with quantized inputs/outputs

quantization: **0.033 * (q - 132)**

**TFLite graph**

**ONNX graph**

**(in ORT)**
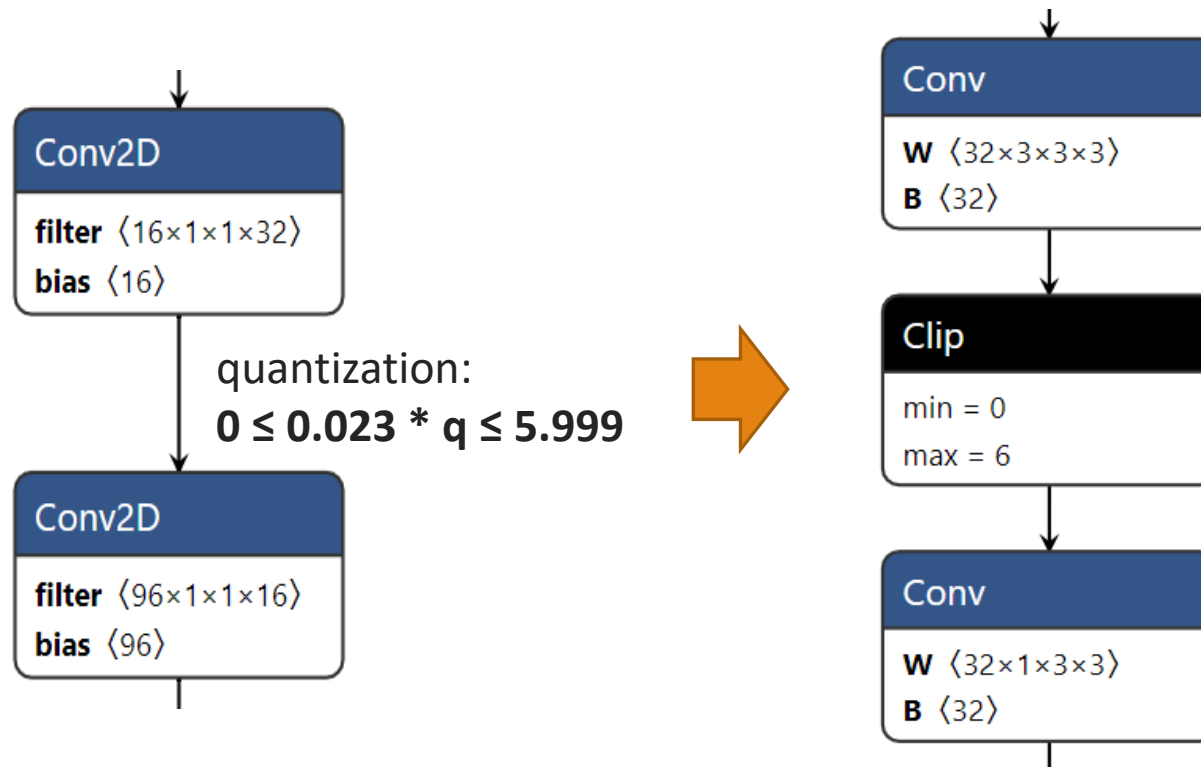
# Dequantizing Models `--dequantize`

Detect ReLU and ReLU6 ops from quantization range

# Support and Feature Requests

Please submit feature requests to GitHub

TFLite -> ONNX conversion is new, expect improvements as we support more ops

**github.com/onnx/tensorflow-onnx**

# ONNX Runtime Mobile

SCOTT MCKAY

MICROSOFT

ONNX RUNTIME MOBILE TECHNICAL LEAD

# ONNX Runtime Mobile

ONNX Runtime Mobile is a variant of ONNX Runtime that minimizes binary size for mobile and edge scenarios

- ◦ Same codebase as ONNX Runtime
- ◦ Available since ONNX Runtime v1.5, Sept 2020

Includes only required operator kernels in the build

- ◦ Can also reduce types supported by operator kernels

Custom format for the model file

# ONNX Runtime Mobile

Runtime usage of ONNX Runtime Mobile is the same as regular ONNX Runtime

◦ C, C++, Python and Java APIs are available

Supports NNAPI Execution Provider on Android

Supports CoreML Execution Provider on iOS (preview)

Documentation:

◦ ONNX_Runtime_for_Mobile_Platforms.md

# ORT format model

Created from an ONNX model
- ◦ Python script handles conversion

During conversion:
- ◦ ONNX Runtime optimizations are applied
  - ◦ e.g. constant folding
- ◦ Nodes are assigned to kernels
  - ◦ No ONNX schema dependency
    - ◦ Significant binary size and memory usage saving

Uses google::flatbuffers

# Operator Kernel selection

Configuration file specifies the kernels to include in the build

- ◦ Model conversion script will automatically generate configuration file when converting models
- ◦ Configuration file can also be manually created/edited

Example config:

- ◦ `ai.onnx;11;AveragePool,Conv,Reshape,Shape,Softmax,Squeeze,Transpose`

# Reduced Type Support

Can limit types that operator kernels support

◦ Model conversion script can automatically detect required types on a per-operator basis
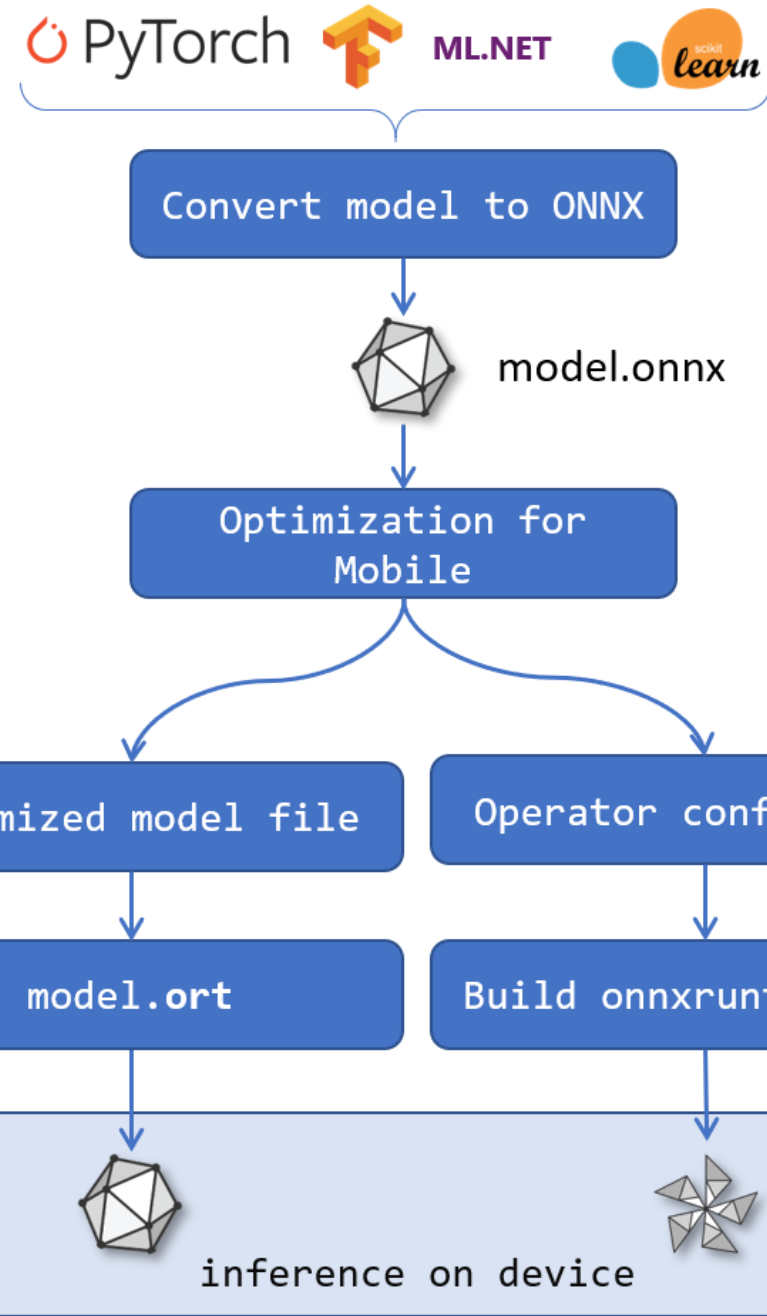◦ Alternatively, can specify a global list of types to support

Model based type reduction generally reduces kernel binary size by 25 - 33%

Available in ONNX Runtime v1.7

◦ March 2021

# ORT Mobile Usage

# Binary size

Primary choices that determine binary size:

- Operators and types to include
- Enable/disable exceptions
- Enable/disable support for traditional ML operators
- Use static or shared libc++ on Android

| Base build size for Android ARM64<br>*NDK 21.1, no operator kernels, shared libc++,*<br>*exceptions and traditional ML support disabled* | libonnxruntime.so: 755KB (280KB in AAR) |
|---|---|
| With operator kernels required by Mobilenet | libonnxruntime.so: 895KB (342KB in AAR) |
| With reduced type support enabled | libonnxruntime.so: 851KB (325KB in AAR)<br>31% reduction in size of kernels |

# NNAPI Support

Usage of NNAPI is determined at runtime
- based on whether NNAPI is available and device capabilities
  - e.g. older version of NNAPI may not support as many operators

Fallback to CPU execution if node cannot be run using NNAPI

Available in ORT v1.6
- December 2020

# Questions and Feature Requests

Please reach out to the ONNX Runtime team

- ◦ https://github.com/microsoft/onnxruntime/discussions