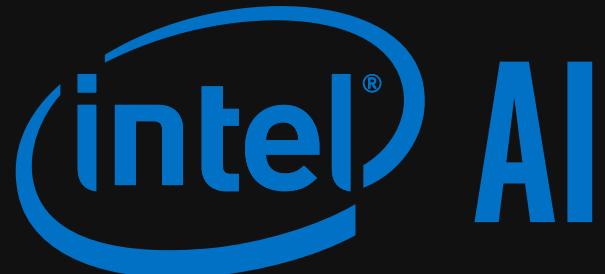


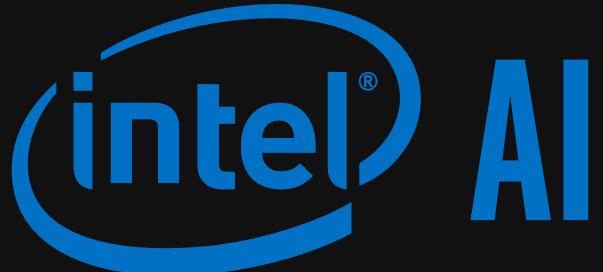


# DEEP LEARNING

---

on the inside





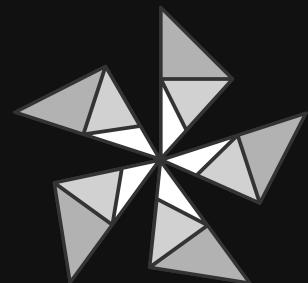
nGraph



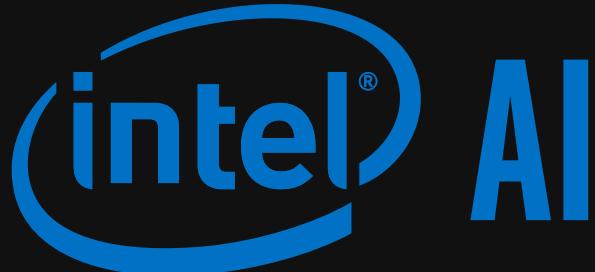
ONNX



OpenVINO



ONNXRuntime



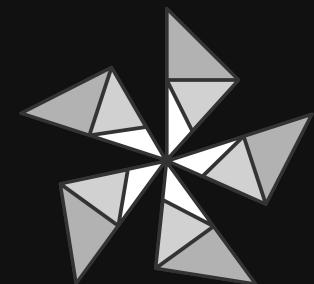
nGraph



ONNX



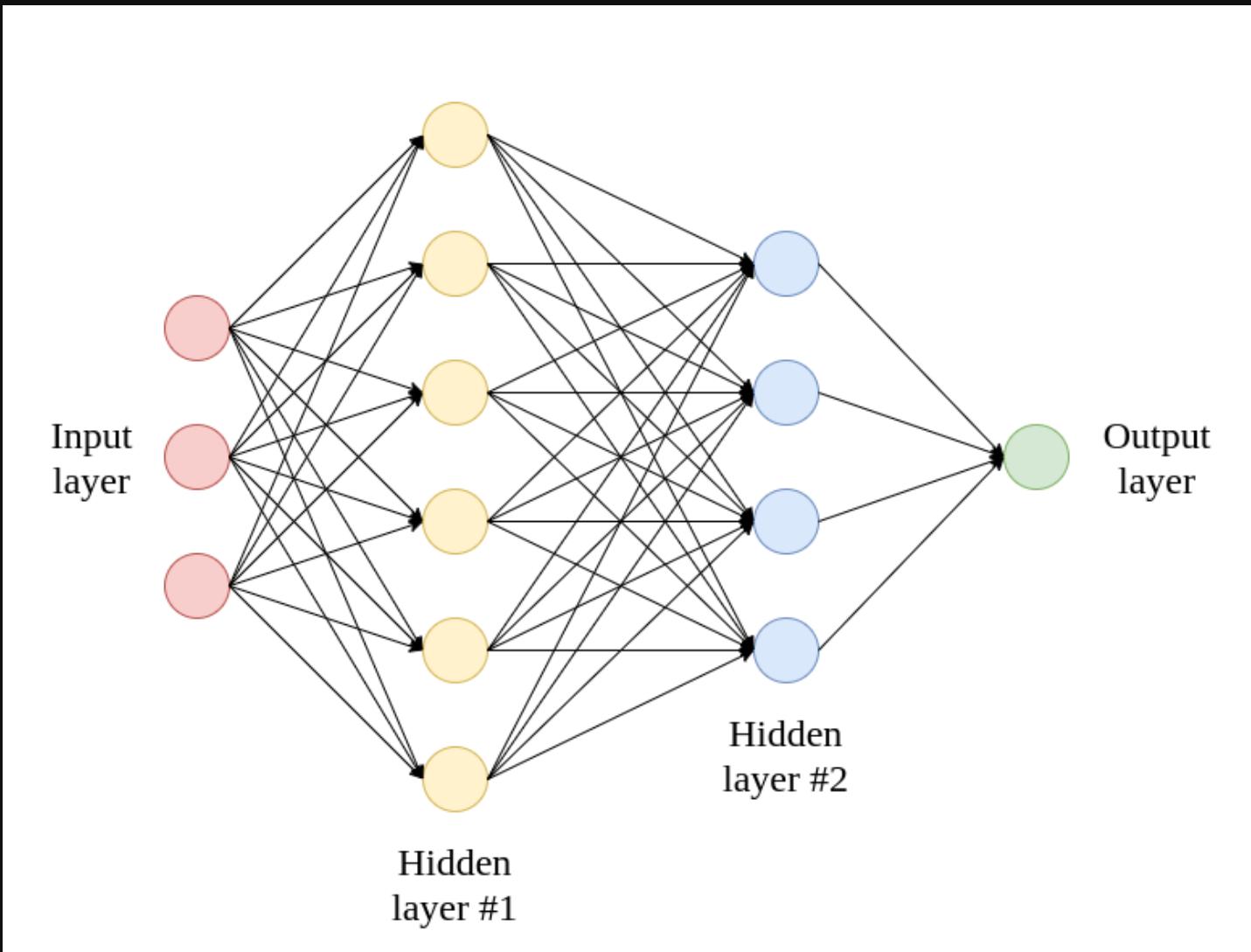
OpenVINO



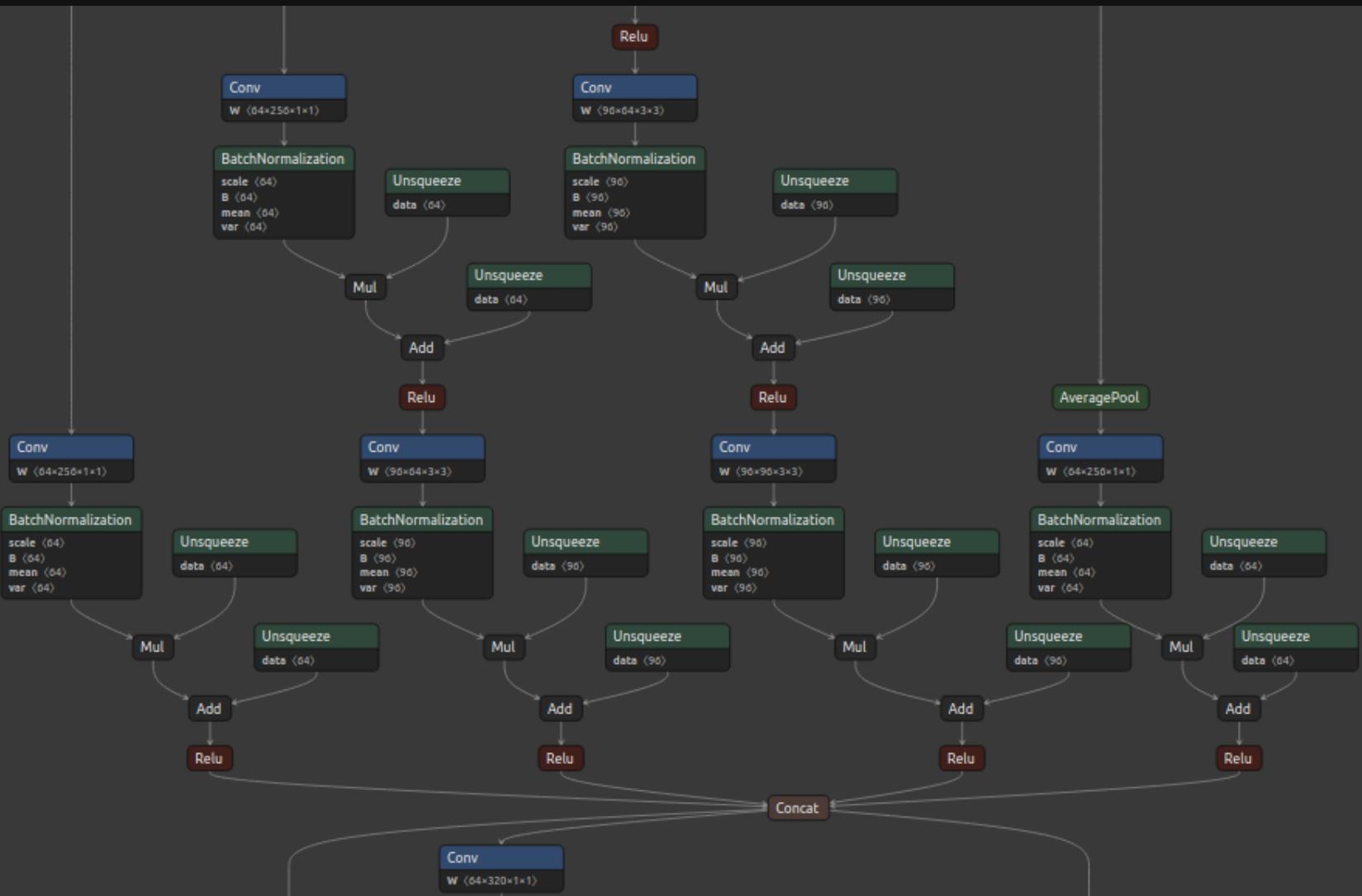
ONNXRuntime



# Artificial neural networks

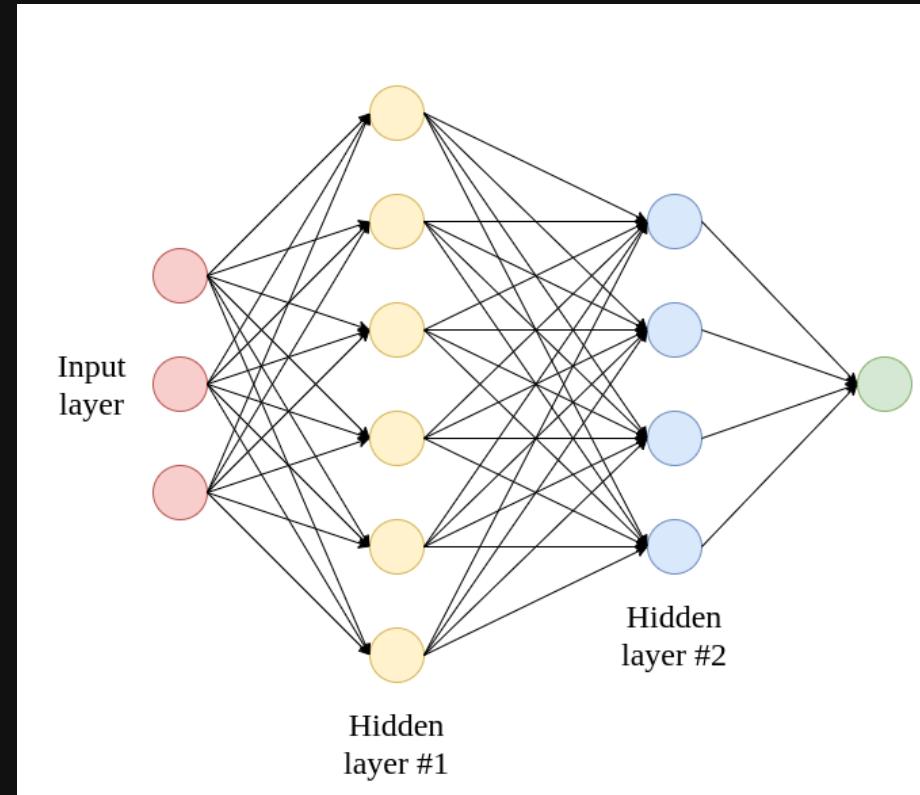
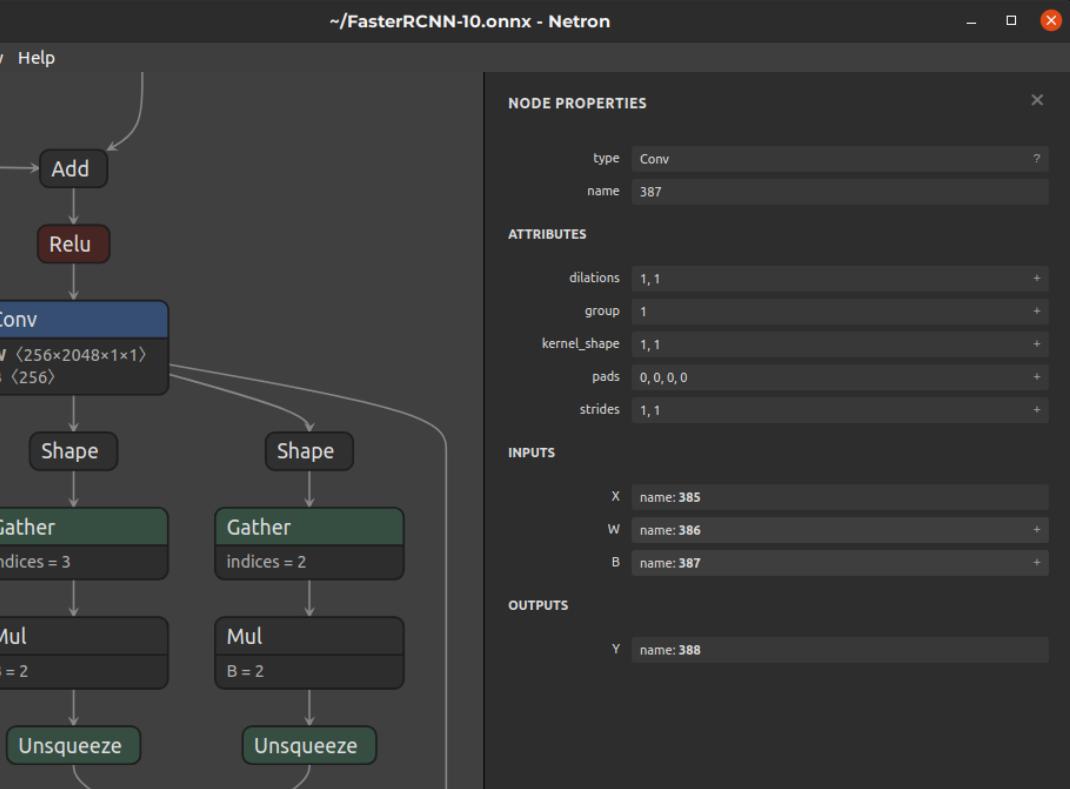


# Real-life models

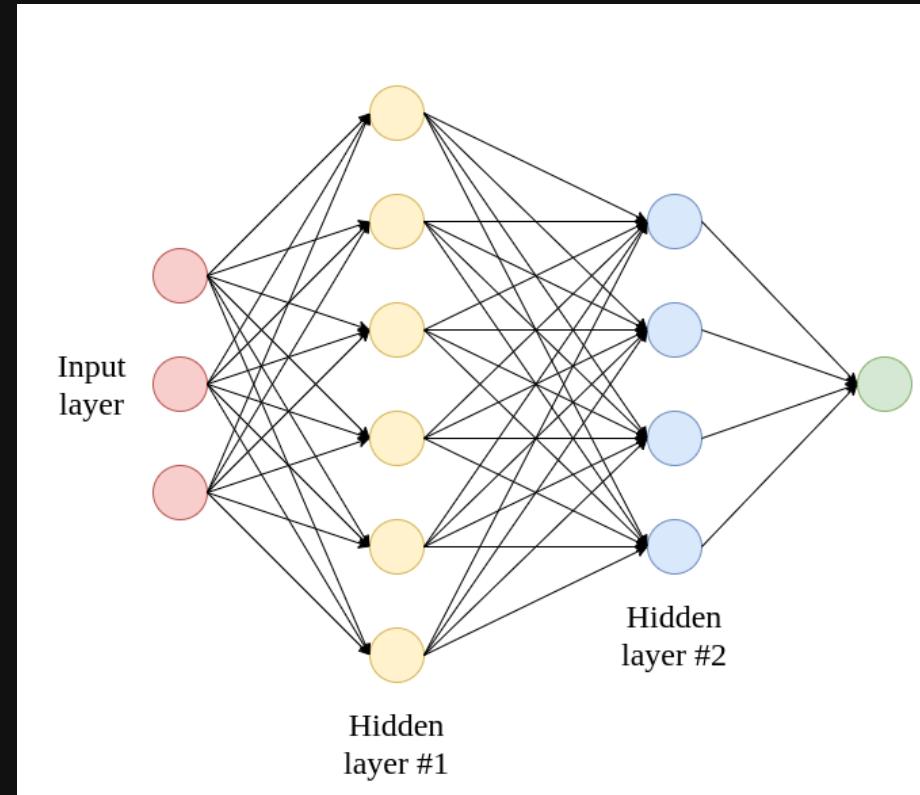
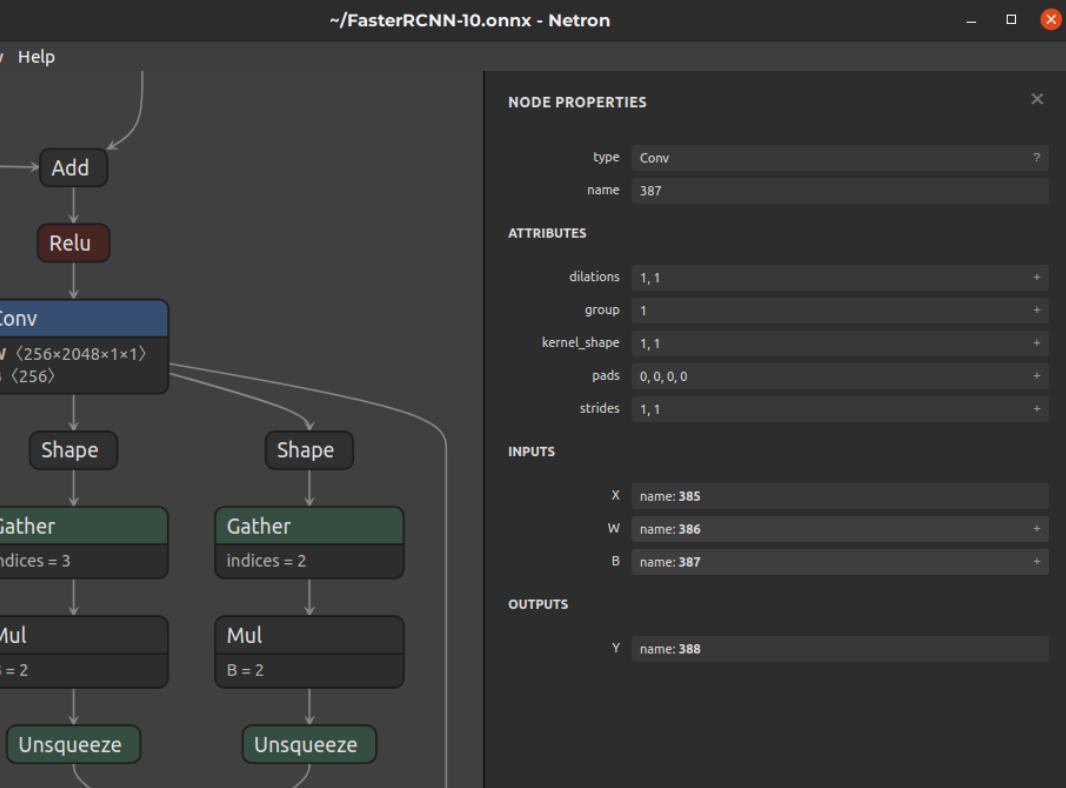


(fragment of the Faster R-CNN model)

# Graph representations

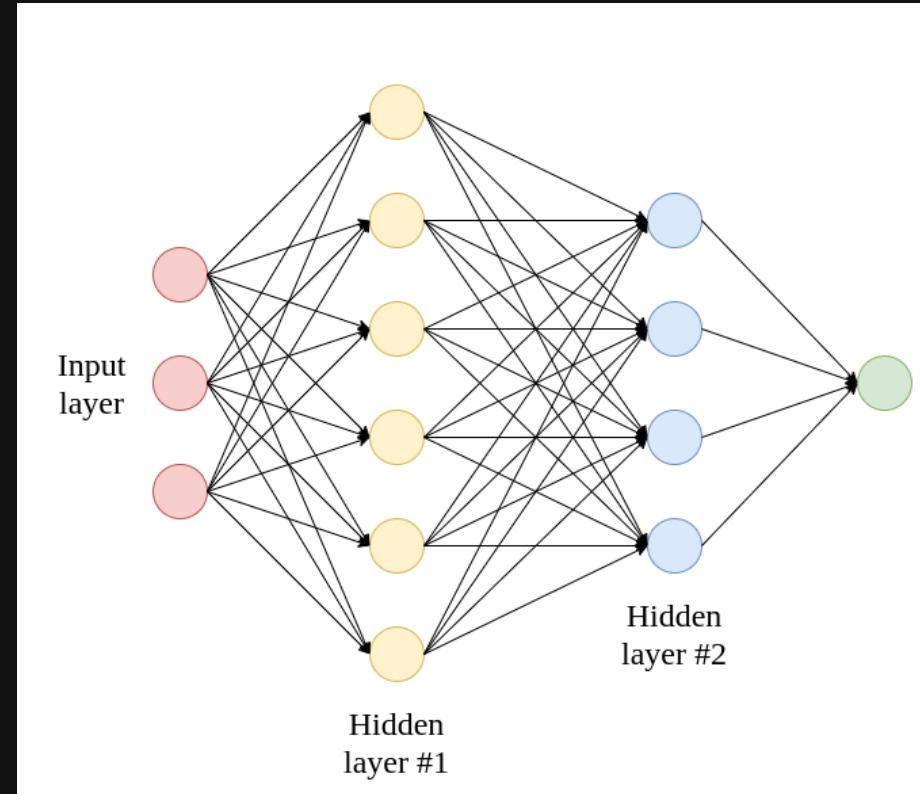
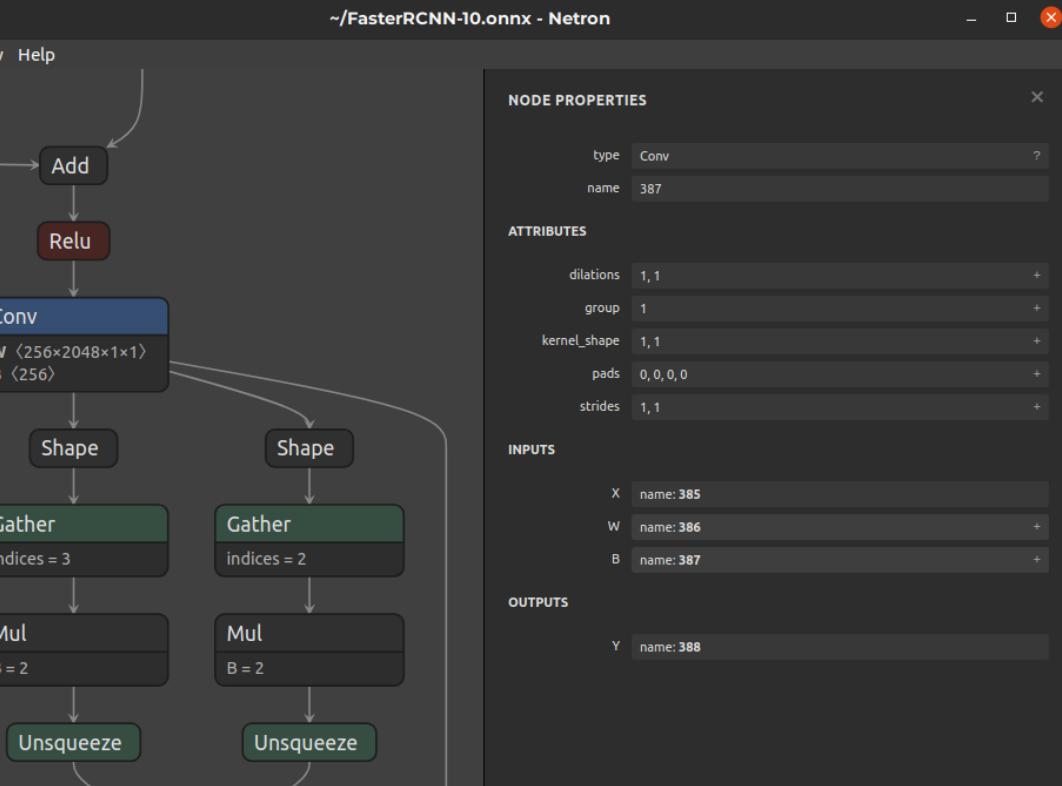


# Graph representations



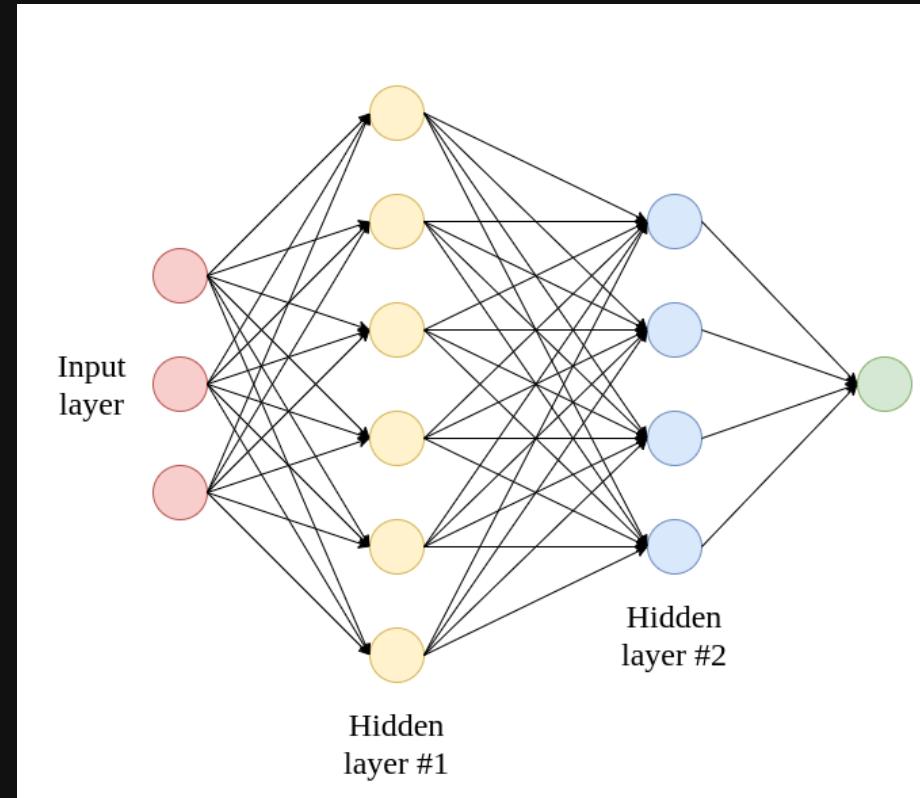
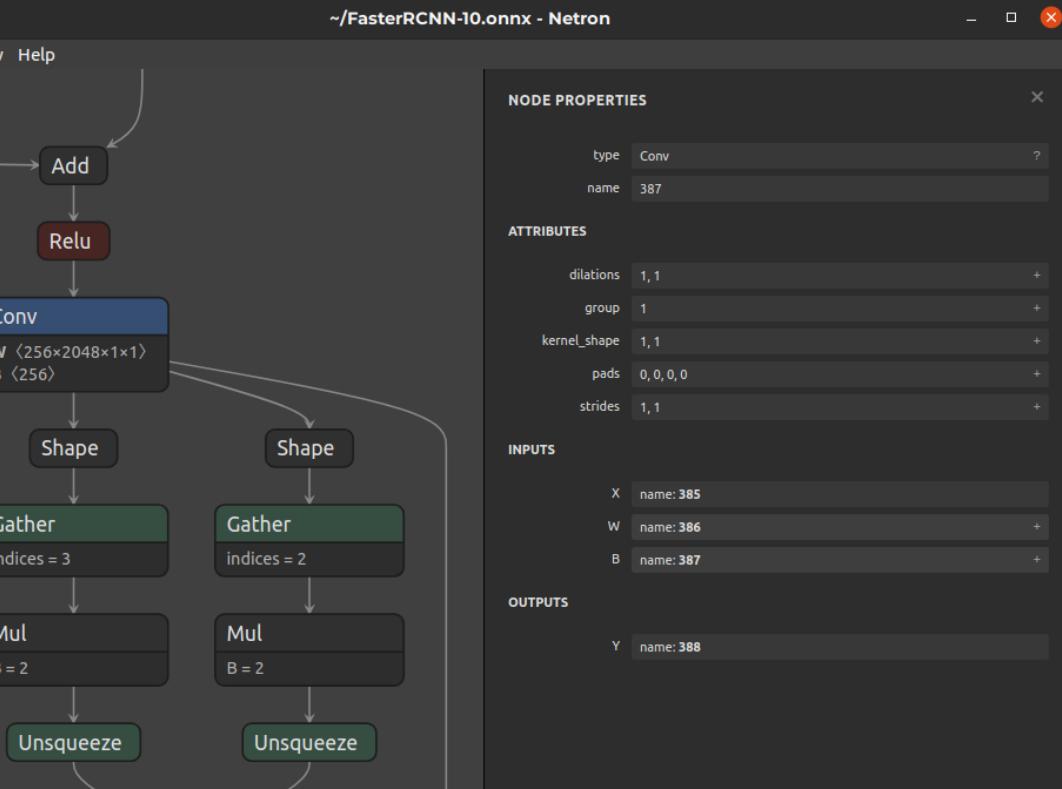
- nodes: operations *vs* individual neurons

# Graph representations



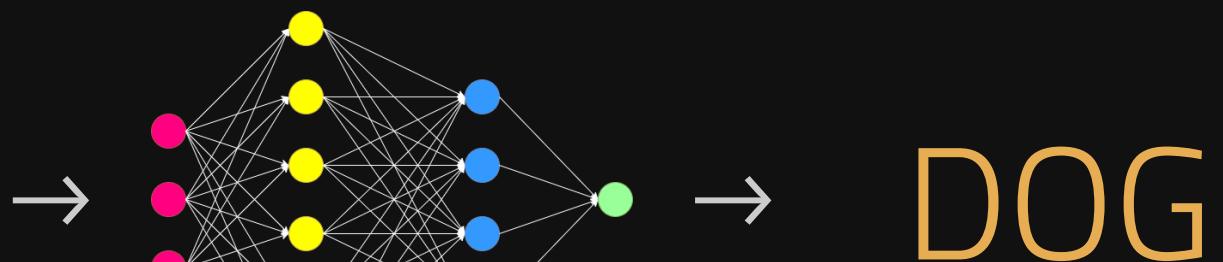
- nodes: operations *vs* individual neurons
- edges: tensors *vs* individual weights

# Graph representations



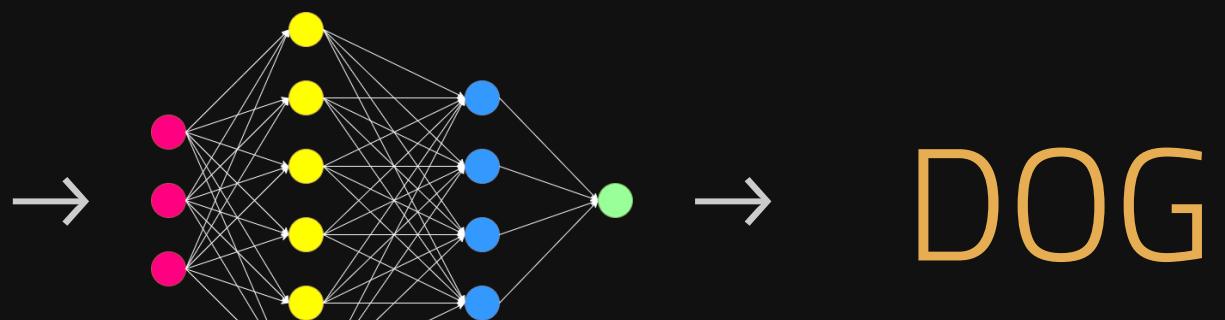
- nodes: operations *vs* individual neurons
- edges: tensors *vs* individual weights
- graph: computational graph *vs* linear graph with layers

# Model as a function





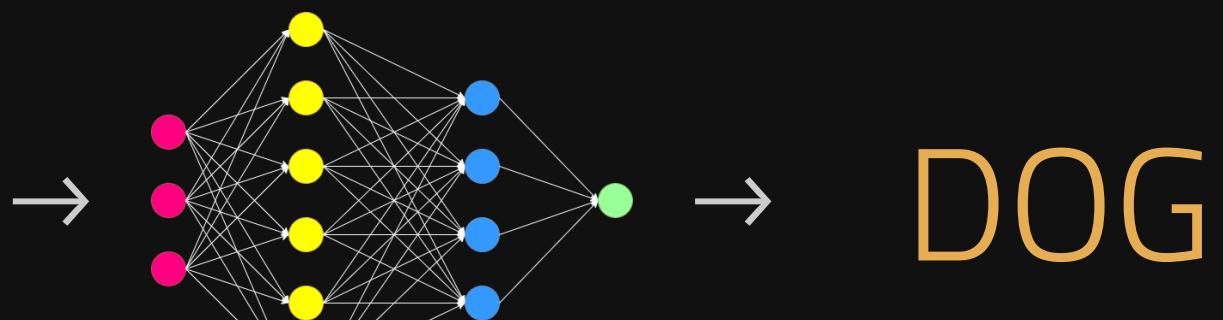
# Model as a function



$$y = f(x)$$



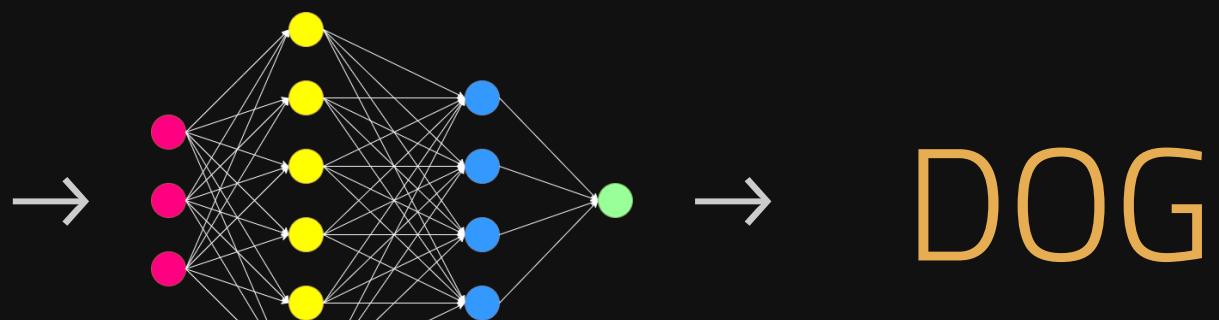
# Model as a function



$$y = f(x) \quad 1 \Rightarrow \text{DOG}$$



# Model as a function

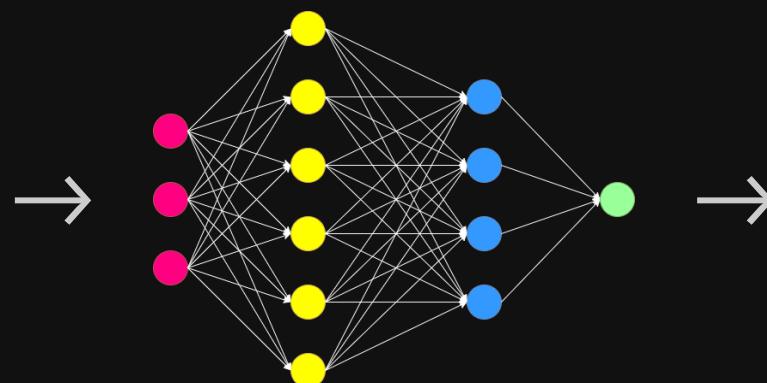


$$y = f(x)$$

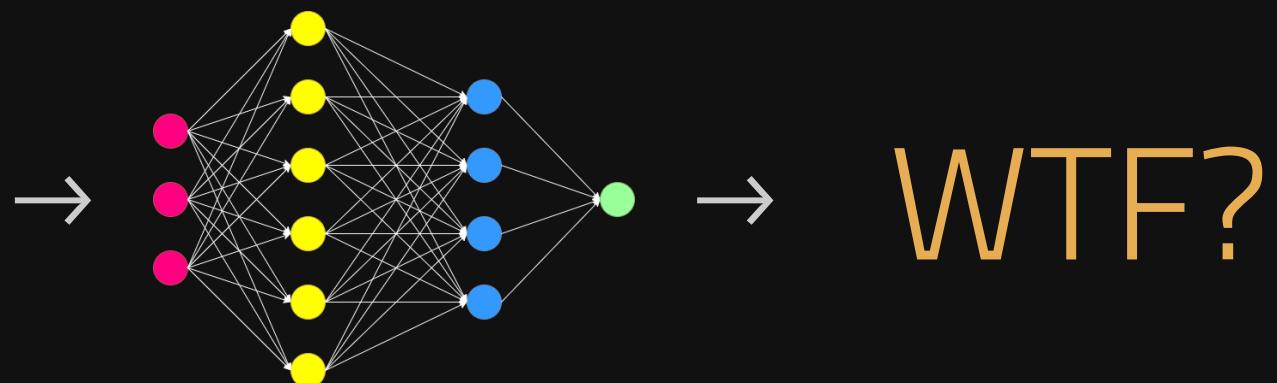
1 ⇒ DOG  
0 ⇒ NOT A DOG



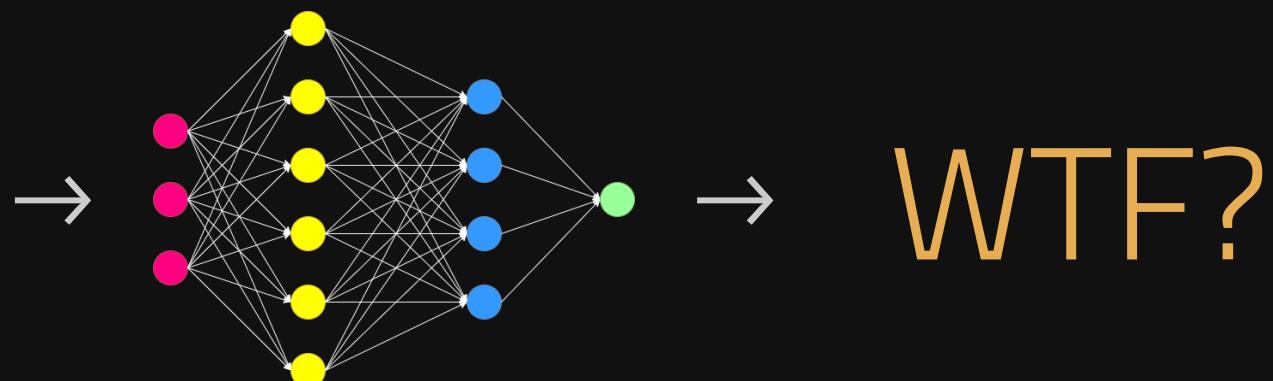
# Model as a function



# Model as a function



# Model as a function



model = approximation of a function  
with a certain accuracy

# Universal approximation theorem

*Universal approximation theorems imply that neural networks can represent a wide variety of interesting functions when given appropriate weights.*

[en.wikipedia.org/wiki/Universal\\_approximation\\_theorem](https://en.wikipedia.org/wiki/Universal_approximation_theorem)

# Deep Learning model

from a C++ developer's perspective

# Deep Learning model

from a C++ developer's perspective

```
// generates a class of functions
template <float A, float B>
float linear_function(const float x)
{
    return A * x + B;
}

float linear_function(const float x)  {
    return 3.14 * x + 1.618;
}
```

# Deep Learning model

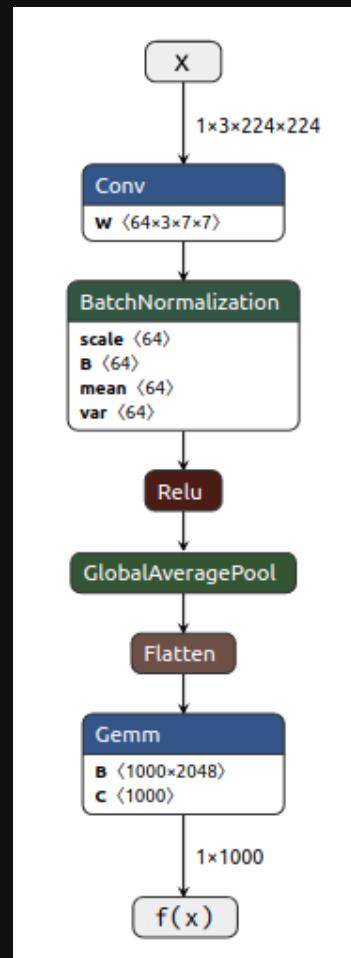
from a *generic* developer's perspective



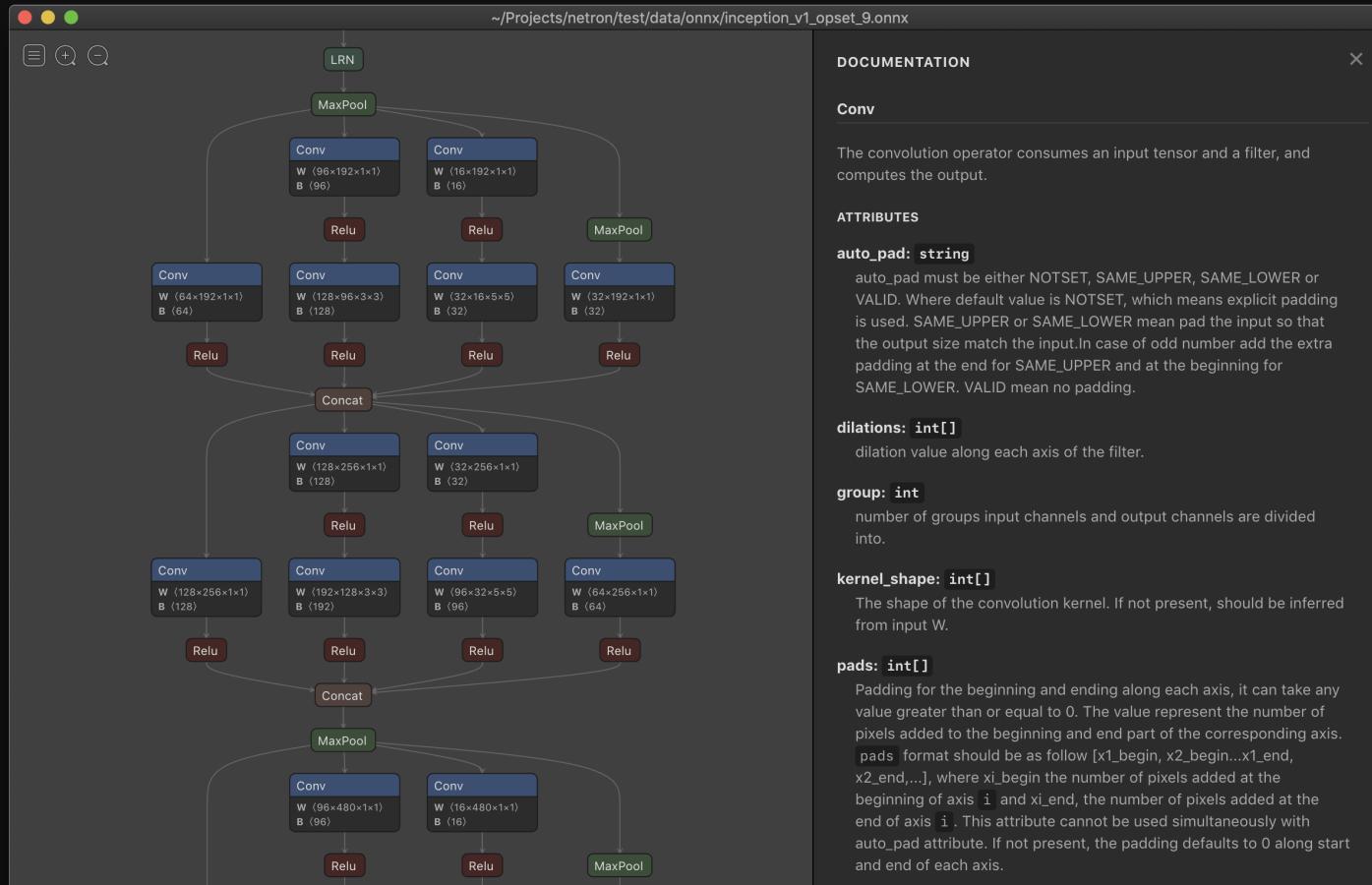
# Deep Learning model

from a *generic* developer's perspective

$$f(x) = \text{Gemm}(\text{Flatten}(\text{GlobalAveragePool}(\text{Relu}(\text{BatchNormalization}(\text{Conv}(x)$$



# Graphs in Deep Learning



lutzroeder/netron

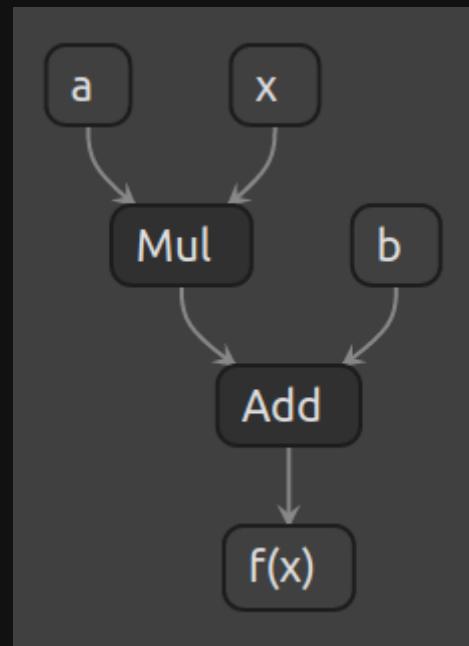
<https://netron.app>

model → function → graph

$$f(x) = ax + b$$

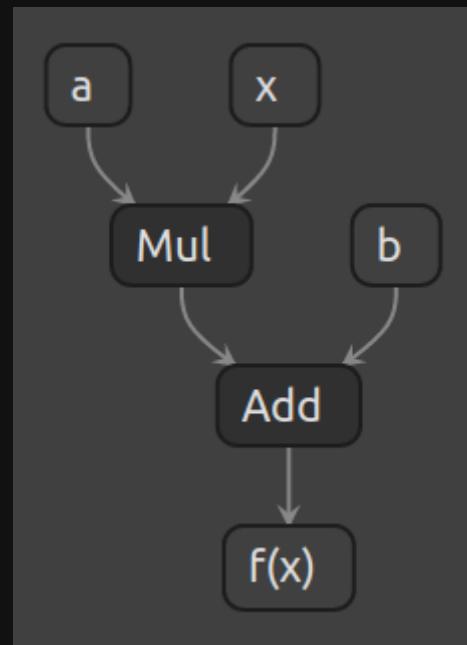
model → function → graph

$$f(x) = ax + b$$



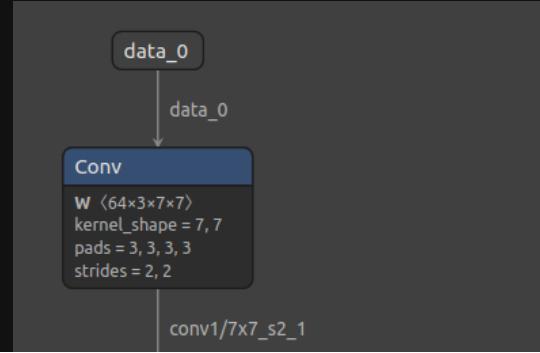
model → function → graph

$$f(x) = ax + b$$

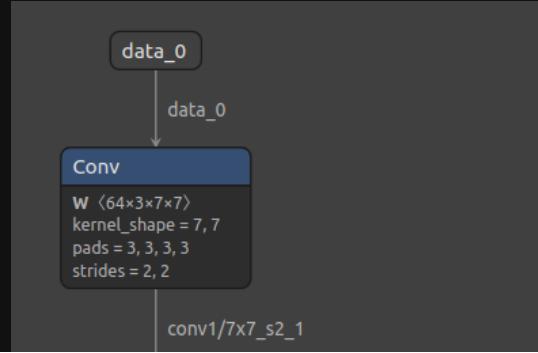


computational graphs

# Operators



# Operators



# Tensors

# Tensors

- N-dimentional data arrays (float, int, int8, int64, ...)

# Tensors

- N-dimensional data arrays (float, int, int8, int64, ...)
- Flat, contiguous blocks of allocated memory

# Tensors

- N-dimentional data arrays (float, int, int8, int64, ...)
- Flat, contiguous blocks of allocated memory
- Shape - {1, 3, 224, 224} or {3, 3}

# Tensors

- N-dimentional data arrays (float, int, int8, int64, ...)
- Flat, contiguous blocks of allocated memory
- Shape - {1, 3, 224, 224} or {3, 3}
- Dimensions (static or dynamic)

# Tensors

- N-dimentional data arrays (float, int, int8, int64, ...)
- Flat, contiguous blocks of allocated memory
- Shape - {1, 3, 224, 224} or {3, 3}
- Dimensions (static or dynamic)
- Rank

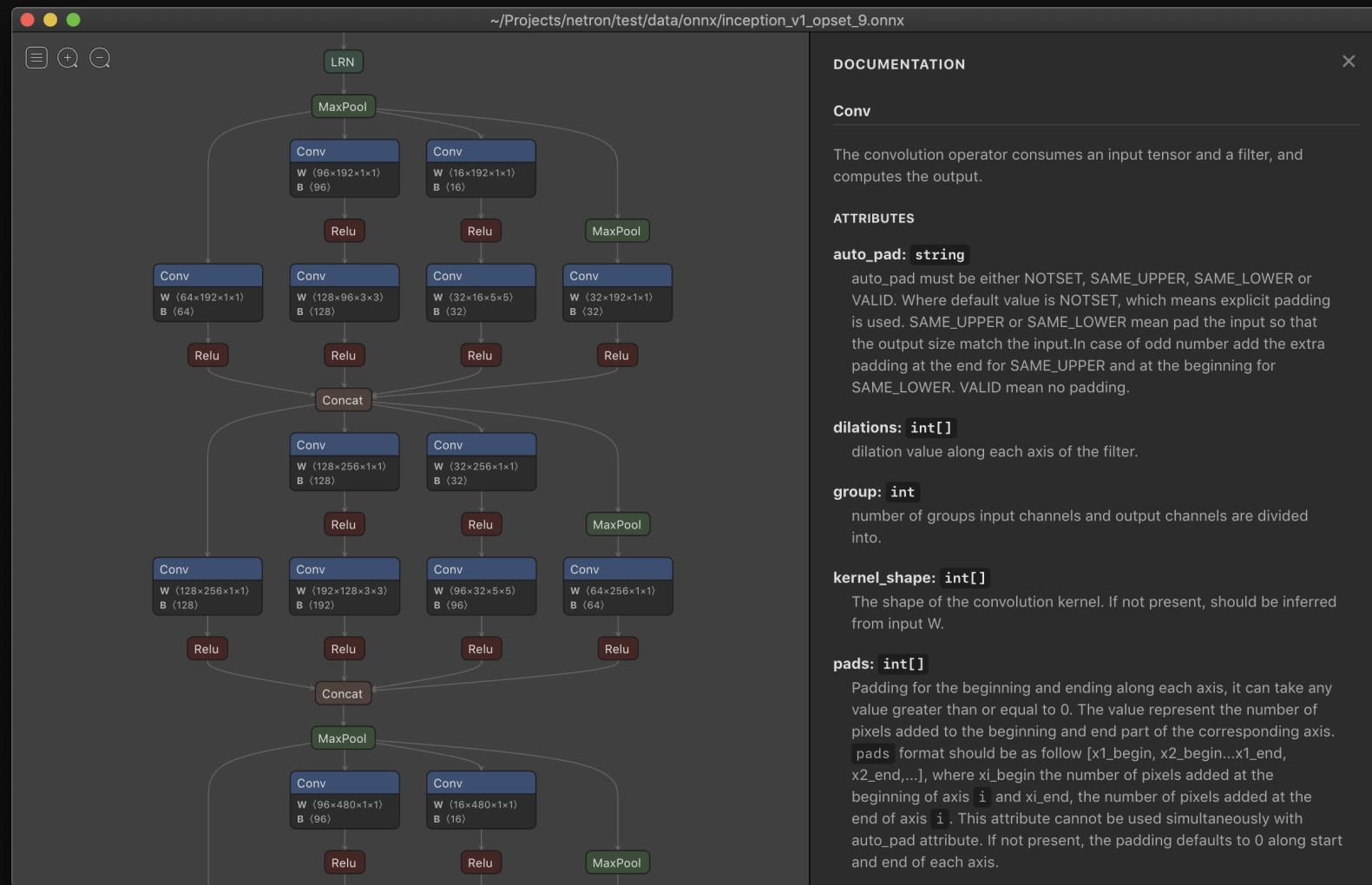
# ReduceSum of a 2D tensor

type	ReduceSum
name	reduce_sum
<strong>ATTRIBUTES</strong>	
axes	1
<strong>INPUTS</strong>	
data	name: data
<strong>OUTPUTS</strong>	
reduced	name: reduced

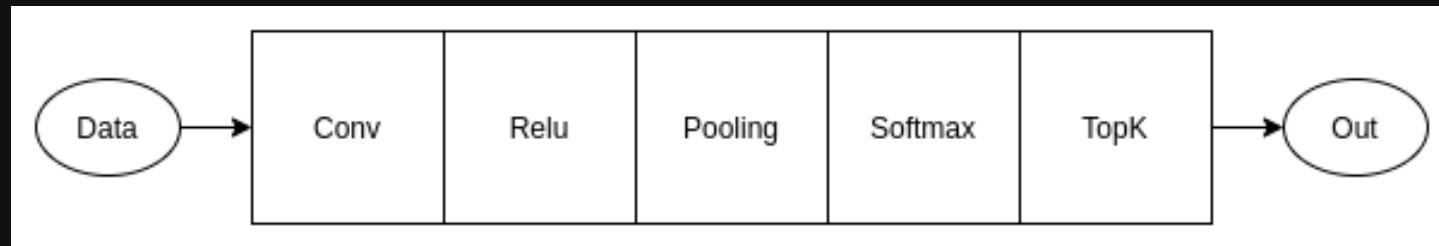


```
const auto data = Tensor({1, 2, 3, 4, 5, 6}, {  
    // [[1 2 3],  
    //  [4 5 6]]  
const auto over_axis_1 = ReduceSum(data, 1);  
// [6, 15]  
const auto over_axis_0 = ReduceSum(data, 0);  
// [5, 7, 9]
```

# Inference

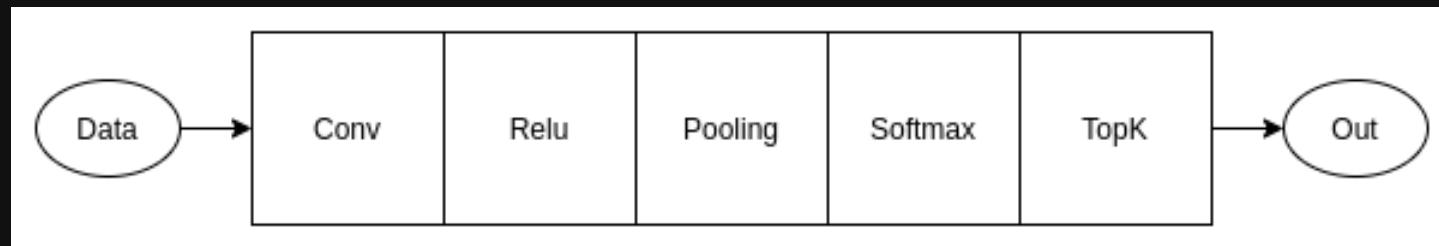


# Naive implementation



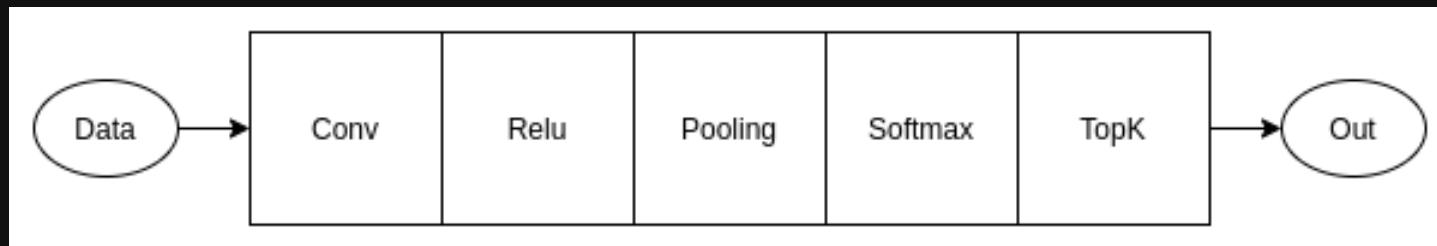
- Topological sort of the graph

# Naive implementation



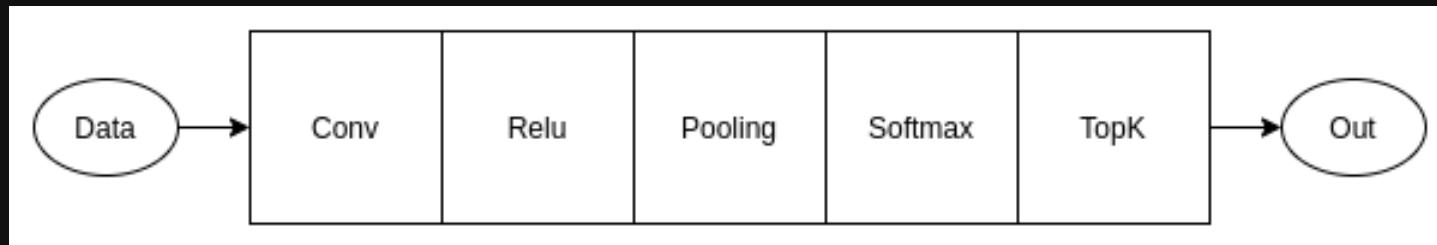
- Topological sort of the graph
- Implementation of the operators

# Naive implementation



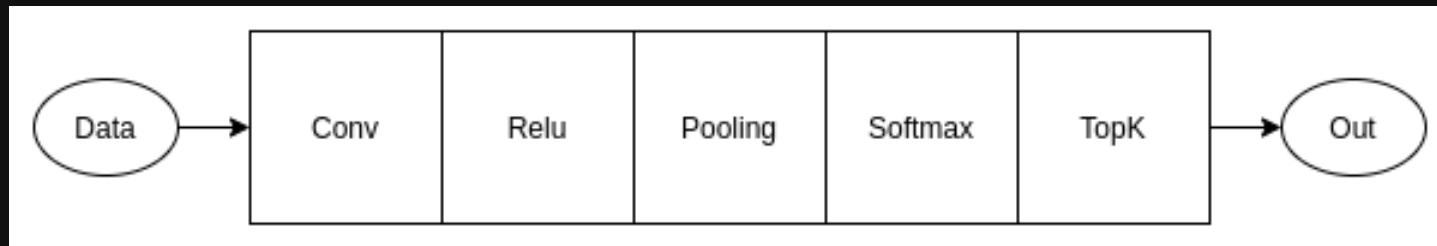
- Topological sort of the graph
- Implementation of the operators
- Sequential execution of sorted nodes

# Naive implementation



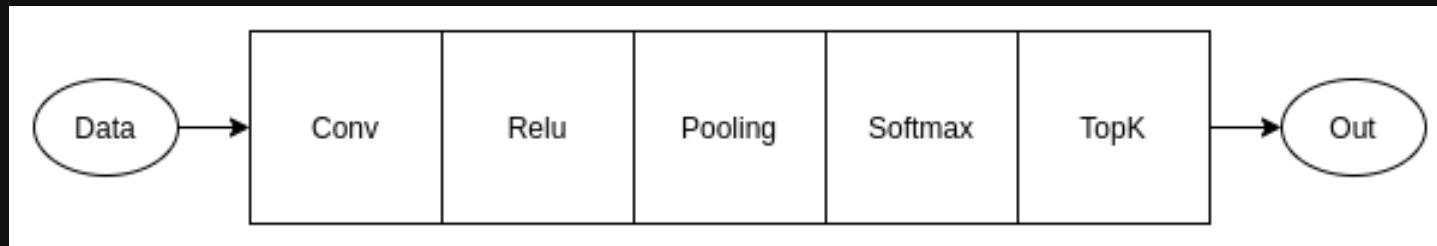
- Topological sort of the graph
- Implementation of the operators
- Sequential execution of sorted nodes
- Probably wasting a lot of HW potential

# Naive implementation



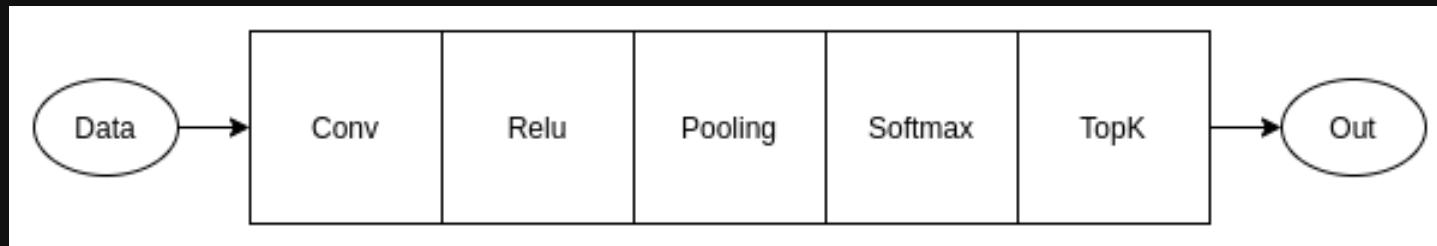
- Topological sort of the graph
- Implementation of the operators
- Sequential execution of sorted nodes
- Probably wasting a lot of HW potential
- Cores not utilized

# Naive implementation



- Topological sort of the graph
- Implementation of the operators
- Sequential execution of sorted nodes
- Probably wasting a lot of HW potential
- Cores not utilized
- No vectorization

# Naive implementation



- Topological sort of the graph
- Implementation of the operators
- Sequential execution of sorted nodes
- Probably wasting a lot of HW potential
- Cores not utilized
- No vectorization
- Not good enough for big, real-life models

# Kernels

# Kernels

- Highly specialized functions/primitives

# Kernels

- Highly specialized functions/primitives
- Hardware specific

# Kernels

- Highly specialized functions/primitives
- Hardware specific
- Architecture & available instructions dependent

# Kernels

- Highly specialized functions/primitives
- Hardware specific
- Architecture & available instructions dependent
- Data type, data shape (rank)

# Kernels

- Highly specialized functions/primitives
- Hardware specific
- Architecture & available instructions dependent
- Data type, data shape (rank)
- oneDNN - Intel CPU & GPU

# Kernels

- Highly specialized functions/primitives
- Hardware specific
- Architecture & available instructions dependent
- Data type, data shape (rank)
- oneDNN - Intel CPU & GPU



oneapi-src/oneDNN

Example: Binary Convolution

# Direct Optimization

# Direct Optimization

- DL framework + kernels = performance

# Direct Optimization

- DL framework + kernels = performance
- $N$  frameworks \*  $M$  hardware platforms = a lot of work

# Direct Optimization

- DL framework + kernels = performance
- $N$  frameworks \*  $M$  hardware platforms = a lot of work
- $K$  operations

# Direct Optimization

- DL framework + kernels = performance
- $N$  frameworks \*  $M$  hardware platforms = a lot of work
- $K$  operations in each framework

# Direct Optimization

- DL framework + kernels = performance
- $N$  frameworks \*  $M$  hardware platforms = a lot of work
- $K$  operations in each framework = even more work

# Direct Optimization

- DL framework + kernels = performance
- $N$  frameworks \*  $M$  hardware platforms = a lot of work
- $K$  operations in each framework = even more work
- HW independent optimizations?



 Intel® nGraph™

 NervanaSystems/ngraph

 openvinotoolkit/openvino



from the technical standpoint



from the technical standpoint

- a library



from the technical standpoint

- a library
-  inside



from the technical standpoint

- a library
-  inside
-  python API



from the functional standpoint



from the functional standpoint

- HW independent IR of DL models



from the functional standpoint

- HW independent IR of DL models
- graph compiler



from the functional standpoint

- HW independent IR of DL models
- graph compiler
- an optimizing compiler

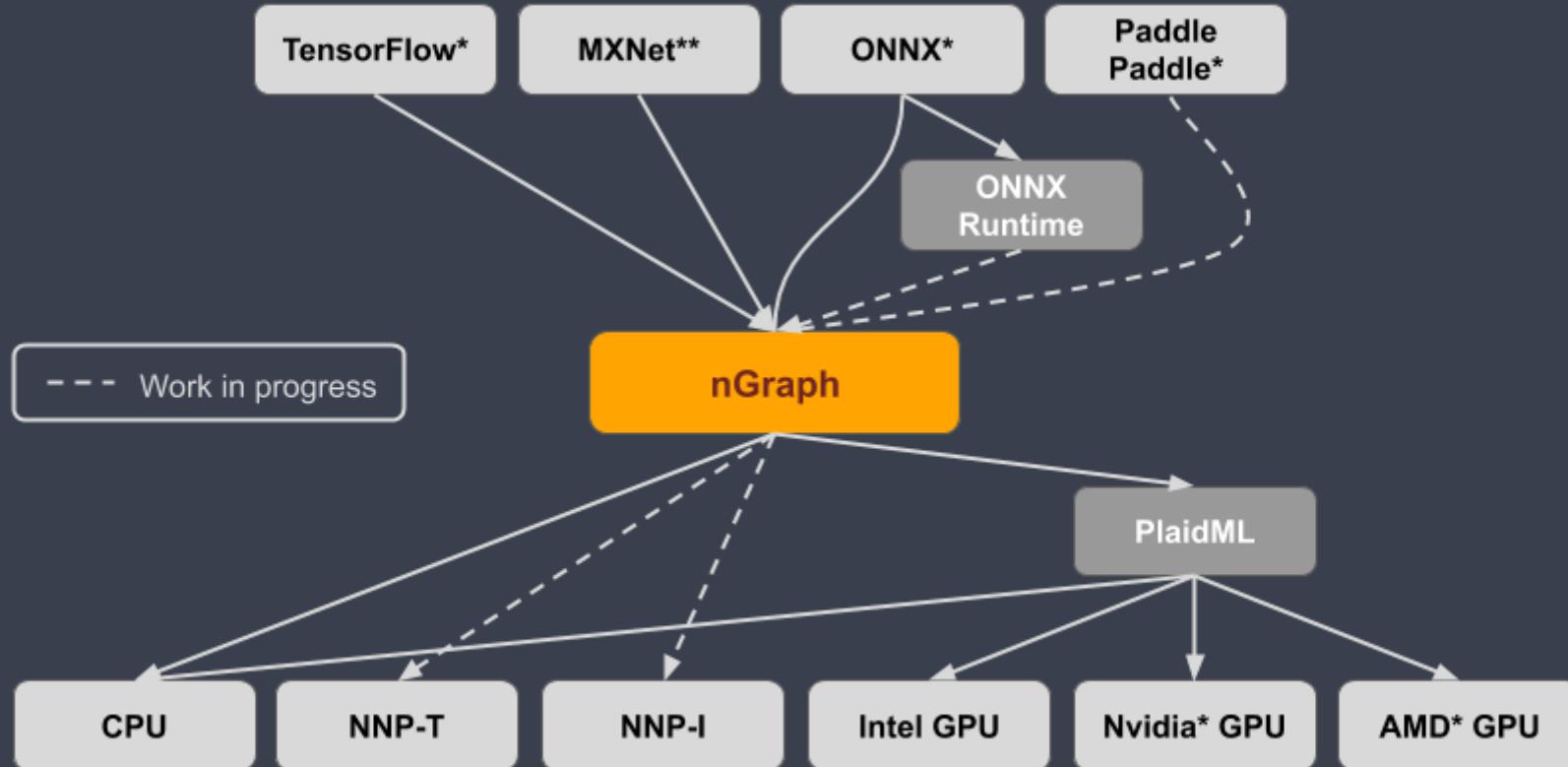


from the functional standpoint

- HW independent IR of DL models
- graph compiler
- an optimizing compiler
- `ngraph::Function`



# nGraph





# Graph optimizations



# Graph optimizations

- HW independent



# Graph optimizations

- HW independent
- Operate on an `ngraph :: Function`



# Graph optimizations

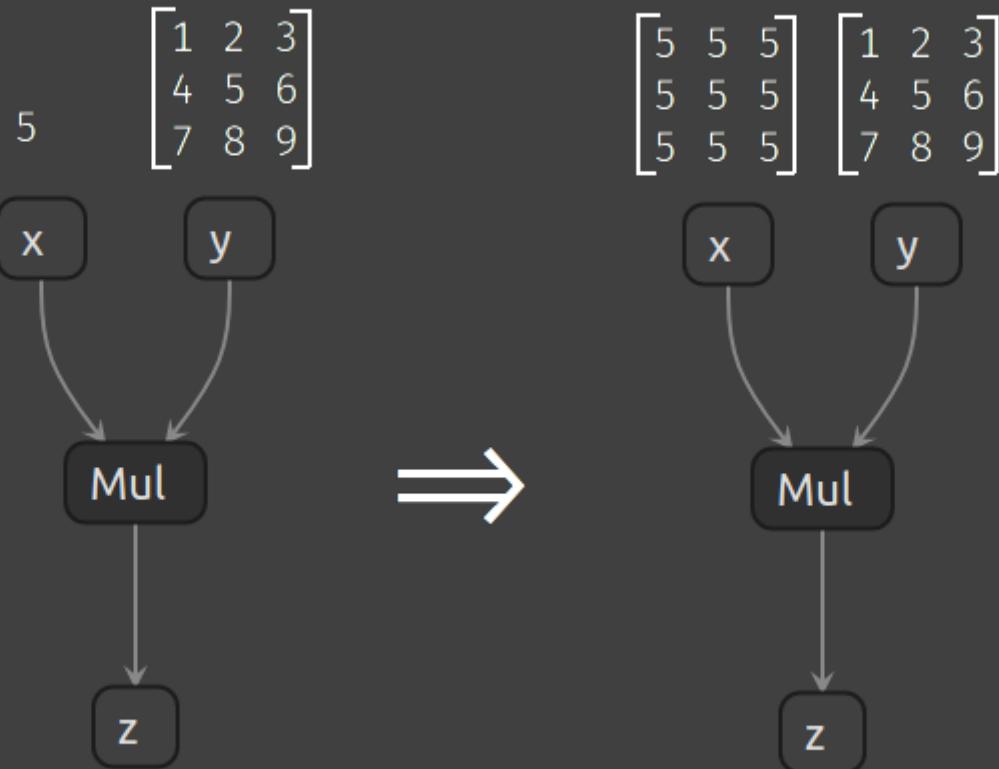
- HW independent
- Operate on an `ngraph :: Function`
- Graph pattern matching & replacement



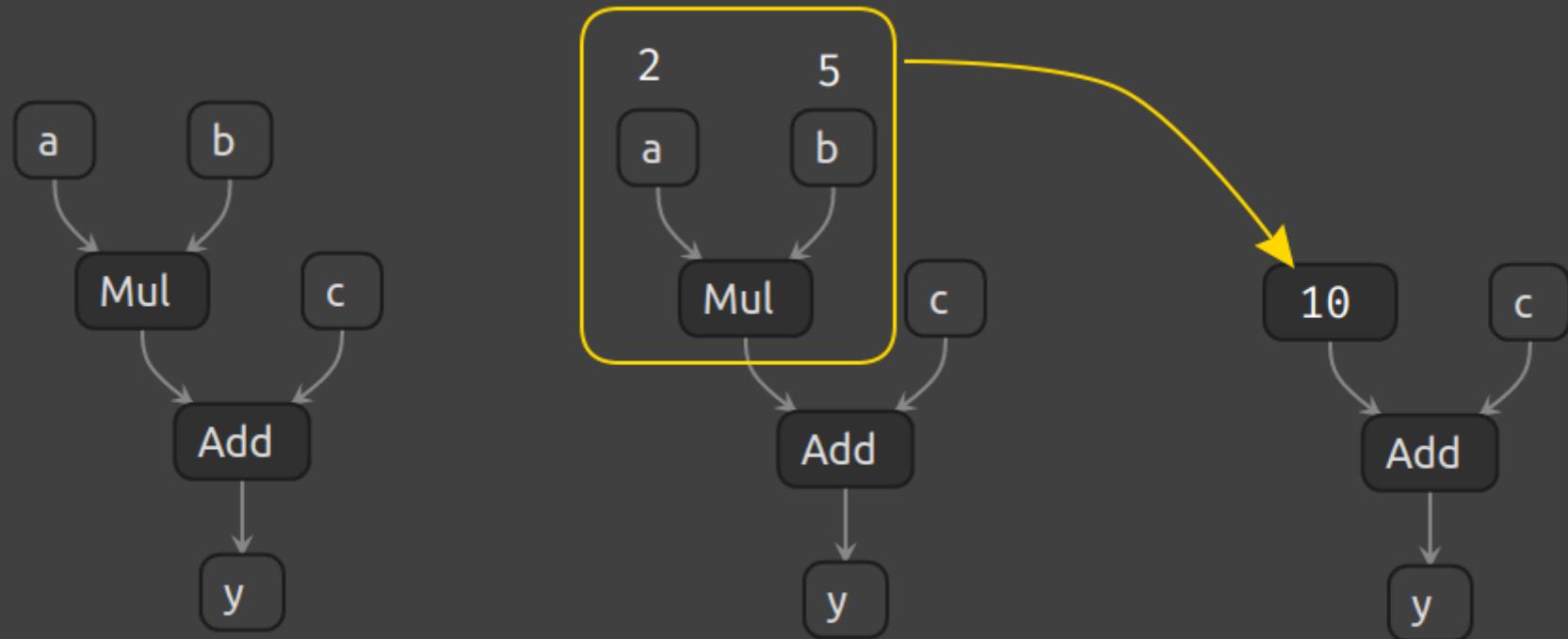
# Graph optimizations

- HW independent
- Operate on an `ngraph :: Function`
- Graph pattern matching & replacement
- "Function passes"

# Implicit broadcast



# Constant folding



# NOP elimination

- Elimination of "useless" nodes

# NOP elimination

- Elimination of "useless" nodes
- Convert  $\text{int} \rightarrow \text{int}$

# NOP elimination

- Elimination of "useless" nodes
- Convert  $int \rightarrow int$
- Broadcast or Reshape when  
 $target\_shape = input\_shape$

# NOP elimination

- Elimination of "useless" nodes
- Convert  $int \rightarrow int$
- Broadcast or Reshape when  
 $target\_shape = input\_shape$
- Reduction without a given axis

# NOP elimination

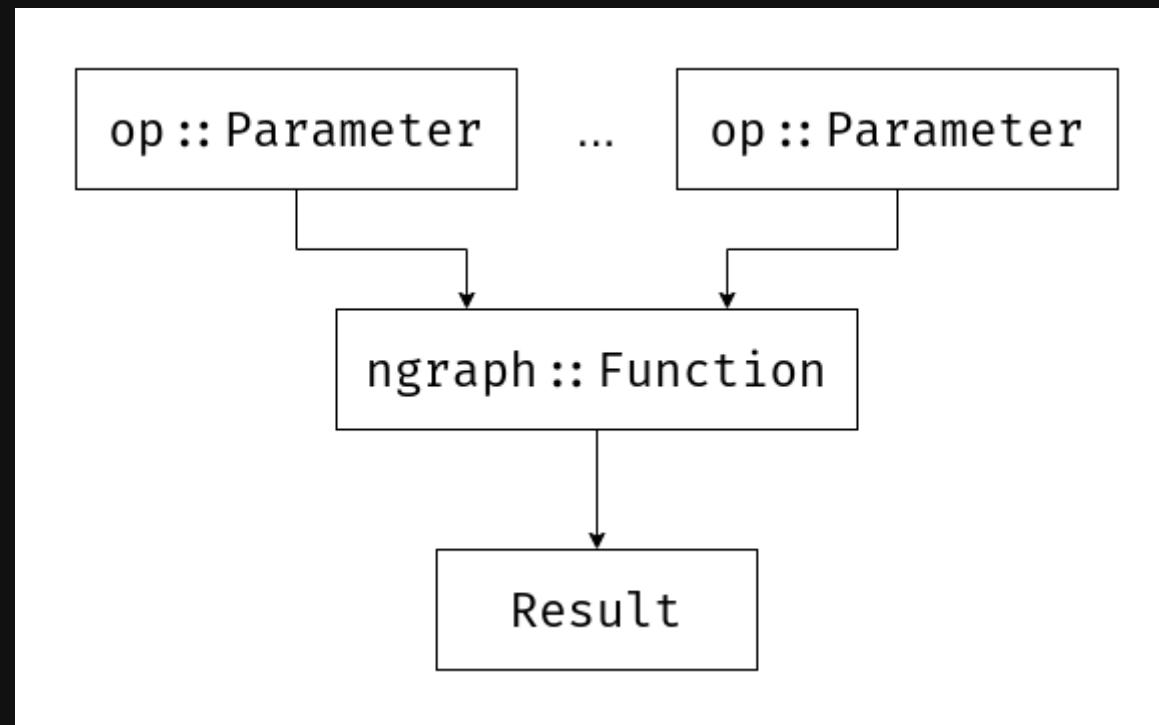
- Elimination of "useless" nodes
- Convert  $int \rightarrow int$
- Broadcast or Reshape when
$$target\_shape = input\_shape$$
- Reduction without a given axis
- Addition of zero or multiplication by one

# NOP elimination

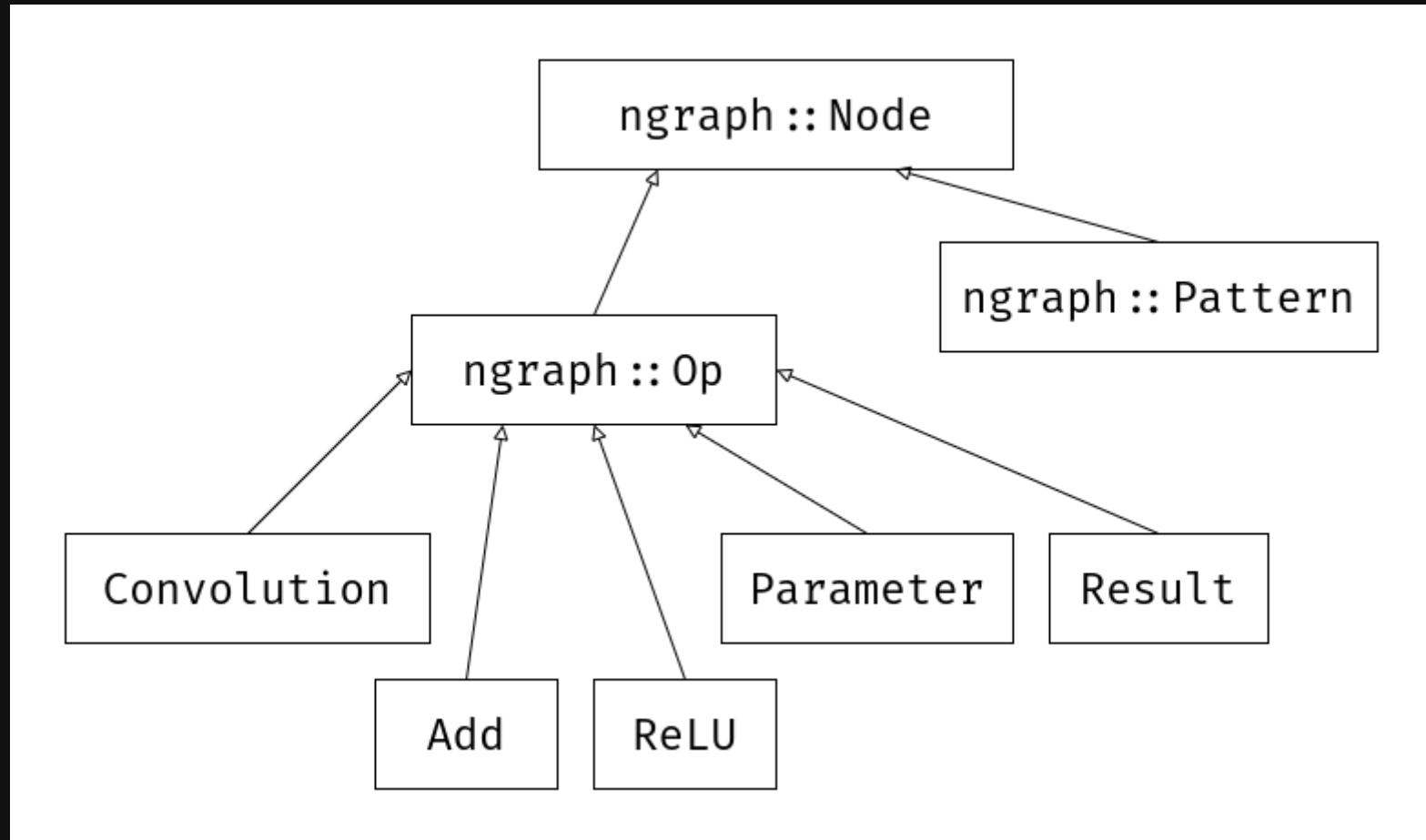
- Elimination of "useless" nodes
- Convert  $int \rightarrow int$
- Broadcast or Reshape when  
 $target\_shape = input\_shape$
- Reduction without a given axis
- Addition of zero or multiplication by one

common\_optimizations.cpp

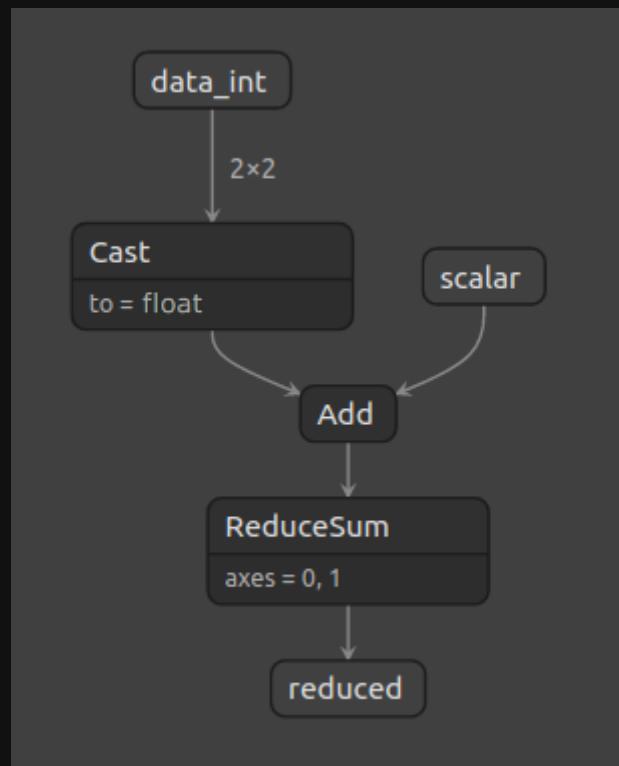
# ngraph::Function



# ngraph::Node



# Creating a simple Function



# Creating a simple Function



```
const auto data_int = make_shared<op::Parameter>(element::i32, Shape{2, 2});
const auto scalar = make_shared<op::Parameter>(element::f32, Shape{});

const auto cast = make_shared<op::Convert>(data_int, element::f32);
const auto add = make_shared<op::Add>(cast, scalar);

const auto reduction_axes = op::Constant::create(element::i32, Shape{2}, {0, 1});
const auto reduce_sum = make_shared<op::ReduceSum>(add, reduction_axes);

const auto function = make_shared<Function>(reduce_sum, ParameterVector{data_int, scalar});
```

# Creating a simple Function



```
const auto data_int = make_shared<op::Parameter>(element::i32, Shape{2, 2});
const auto scalar = make_shared<op::Parameter>(element::f32, Shape{});

const auto cast = make_shared<op::Convert>(data_int, element::f32);
const auto add = make_shared<op::Add>(cast, scalar);

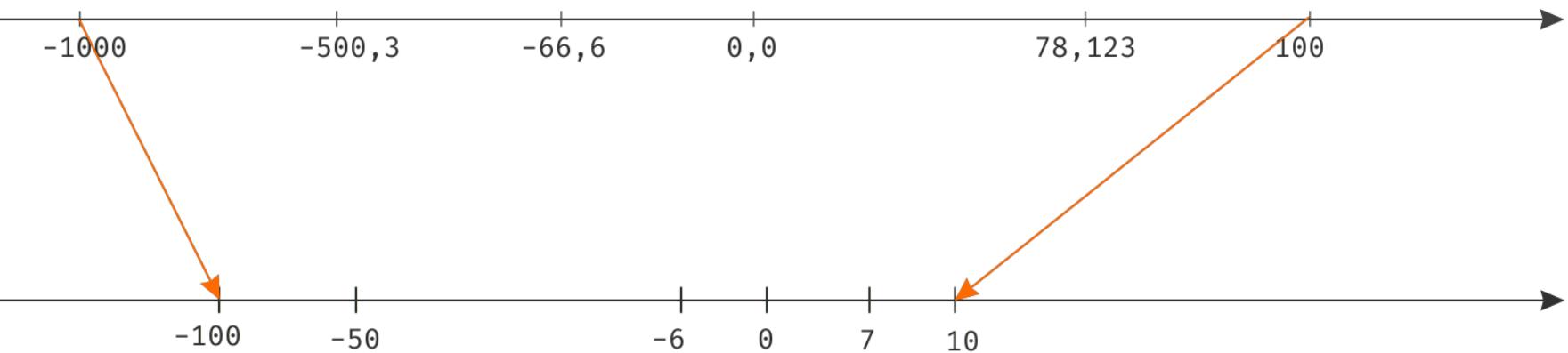
const auto reduction_axes = op::Constant::create(element::i32, Shape{2}, {0, 1});
const auto reduce_sum = make_shared<op::ReduceSum>(add, reduction_axes);

const auto function = make_shared<Function>(reduce_sum, ParameterVector{data_int, scalar});
```

[ideone.com/TJQ1IL](https://ideone.com/TJQ1IL)

# Quantization

$$q(r) = \lfloor r/s \rfloor + z$$



# #include <limits>

- `std::numeric_limits<float>`

$$1 \times 10^{-38} \longleftrightarrow 3 \times 10^{38}$$



# #include <limits>

- `std::numeric_limits<float>`

$$1 \times 10^{-38} \longleftrightarrow 3 \times 10^{38}$$

- `std::numeric_limits<int8_t>`

$$-128 \longleftrightarrow 127$$



# #include <limits>

- `std::numeric_limits<float>`

$$1 \times 10^{-38} \longleftrightarrow 3 \times 10^{38}$$

- `std::numeric_limits<int8_t>`

$$-128 \longleftrightarrow 127$$

- `std::numeric_limits<uint8_t>`

$$0 \longleftrightarrow 255$$



# Quantization

$$q(r) = \lfloor r/s \rfloor + z$$

REAL VALUE	SCALE	SCALE	SCALE	SCALE	ZERO POINT
	2	4	8	50	0
-1000,000	-500	-250	-125	-20	
-97,555	-49	-24	-12	-2	
-66,600	-33	-17	-8	-1	
0,000	0	0	0	0	
1,618	1	0	0	0	
3,142	2	1	0	0	
20,000	10	5	3	0	
1000,000	500	250	125	20	

# Joules and $\mu\text{m}^2$ per operation

Operation	Energy (pJ)	Area ( $\mu\text{m}^2$ )
int8 addition	0.03	36
int16 addition	0.05	67
int32 addition	0.1	137
float16 addition	0.4	1,360
float32 addition	0.9	4,184
int8 multiplication	0.2	282
int32 multiplication	3.1	3,495
float16 multiplication	1.1	1,640
float32 multiplication	3.7	7,700

# Quantization benefits

- Less power (edge devices)

# Quantization benefits

- Less power (edge devices)
- Smaller chip

# Quantization benefits

- Less power (edge devices)
- Smaller chip
- Better cache utilization

# Quantization benefits

- Less power (edge devices)
- Smaller chip
- Better cache utilization
- Smaller weights == smaller model

# Quantization benefits

- Less power (edge devices)
- Smaller chip
- Better cache utilization
- Smaller weights == smaller model
- Performance

# HW support - VNNI

[intel.ai/intel-deep-learning-boost](https://intel.ai/intel-deep-learning-boost)

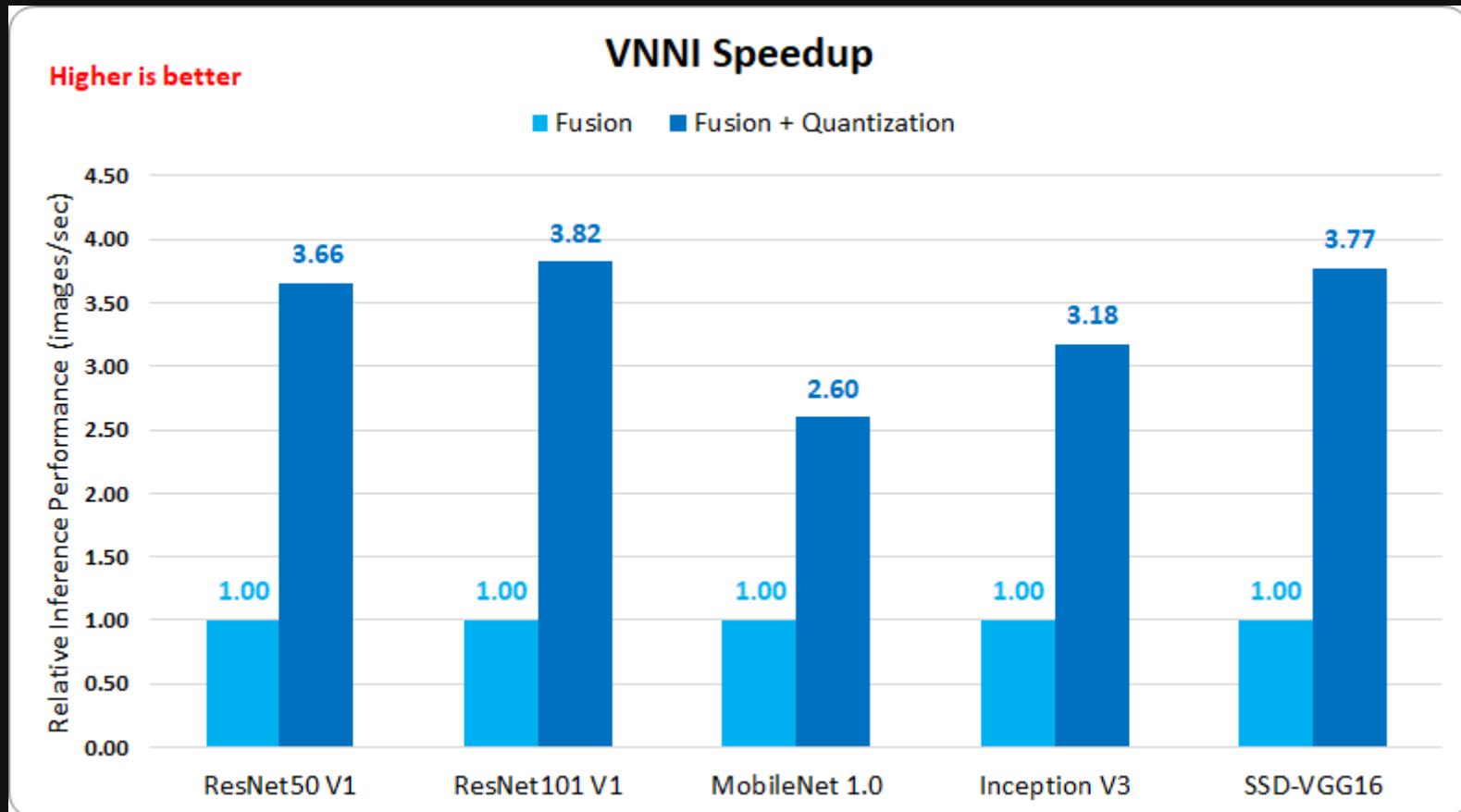
[avx-512-animation.html](https://avx-512-animation.html)

[lowering-numerical-precision-increase-deep-learning-performance.html](https://lowering-numerical-precision-increase-deep-learning-performance.html)

AVX512\_VNNI

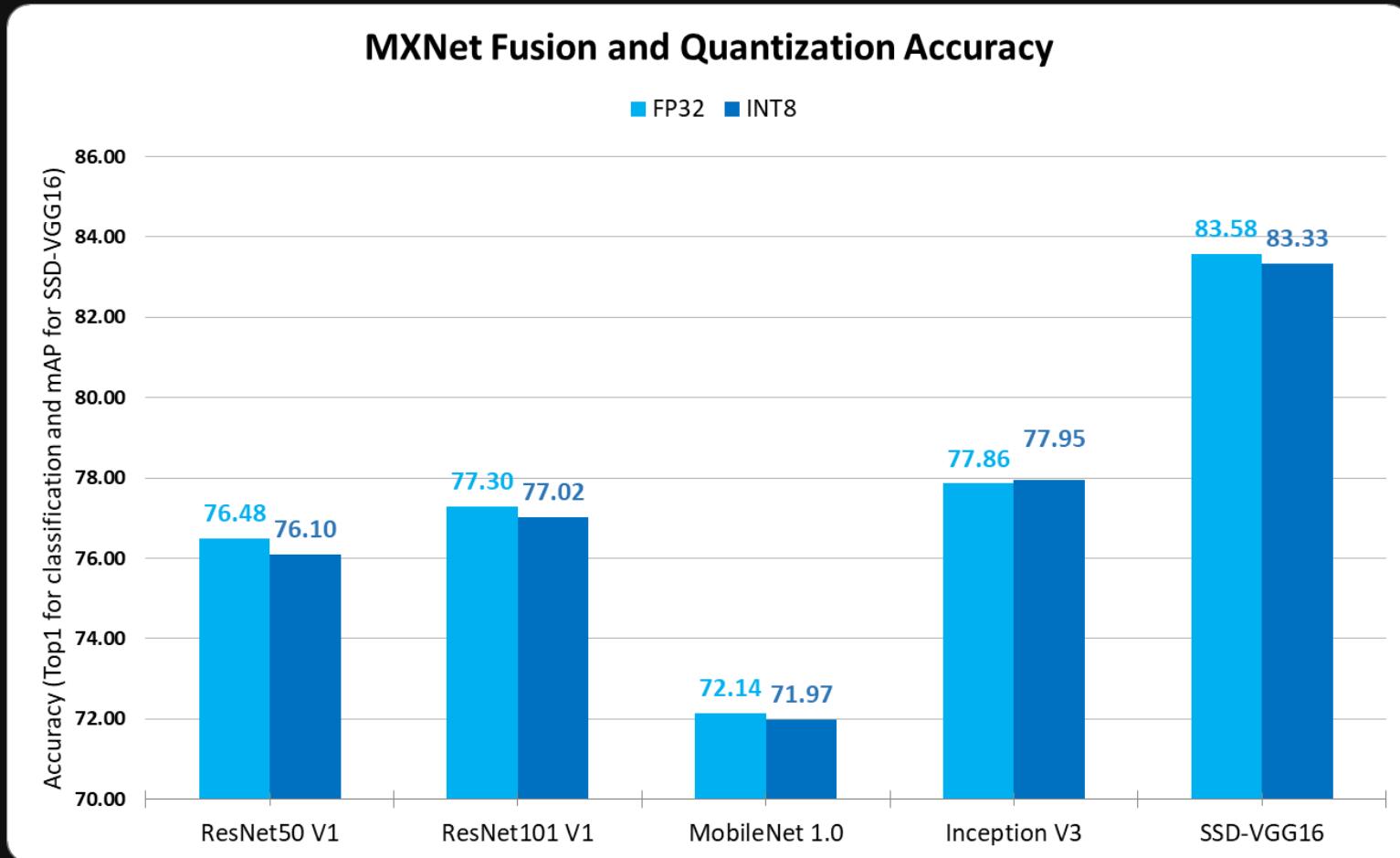


# fp32 vs int8



[medium.com/apache-mxnet/model-quantization-for-production-level-neural-network-inference-f54462ebba05](https://medium.com/apache-mxnet/model-quantization-for-production-level-neural-network-inference-f54462ebba05)

# What about accuracy?



[medium.com/apache-mxnet/model-quantization-for-production-level-neural-network-inference-f54462ebba05](https://medium.com/apache-mxnet/model-quantization-for-production-level-neural-network-inference-f54462ebba05)



# DEEP LEARNING

---

on the inside

[tomdol.github.io/dl-on-the-inside](http://tomdol.github.io/dl-on-the-inside)



[www.youtube.com/watch?v=C9S0nmNS8bQ](https://www.youtube.com/watch?v=C9S0nmNS8bQ)

# Materials

- [docs.openvinotoolkit.org/latest/index.html](https://docs.openvinotoolkit.org/latest/index.html)
- [docs.openvinotoolkit.org/latest/omz\\_models\\_public\\_index.html](https://docs.openvinotoolkit.org/latest/omz_models_public_index.html)
- [docs.openvinotoolkit.org/latest/omz\\_demos.html](https://docs.openvinotoolkit.org/latest/omz_demos.html)
- [github.com/openvinotoolkit/openvino](https://github.com/openvinotoolkit/openvino)
- [github.com/openvinotoolkit/model\\_server](https://github.com/openvinotoolkit/model_server)
- [github.com/openvinotoolkit/openvino\\_notebooks](https://github.com/openvinotoolkit/openvino_notebooks)
- [onnx.ai](https://onnx.ai)
- [github.com/onnx/models](https://github.com/onnx/models)
- [3Blue1Brown - neural nets](#)
- [coursera.org/specializations/deep-learning](https://coursera.org/specializations/deep-learning)
- [infoshare.pl](https://infoshare.pl)