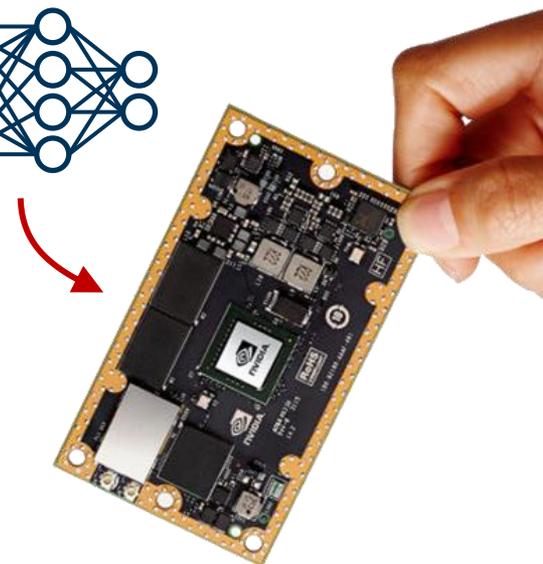
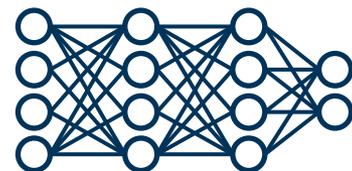


# Speeding up Deep Neural Networks on the Jetson TX1

Markus Eisenbach

Ilmenau University of Technology  
(Germany)

Neuroinformatics and Cognitive Robotics Lab  
[markus.eisenbach@tu-ilmenau.de](mailto:markus.eisenbach@tu-ilmenau.de)  
[www.tu-ilmenau.de/neurob](http://www.tu-ilmenau.de/neurob)

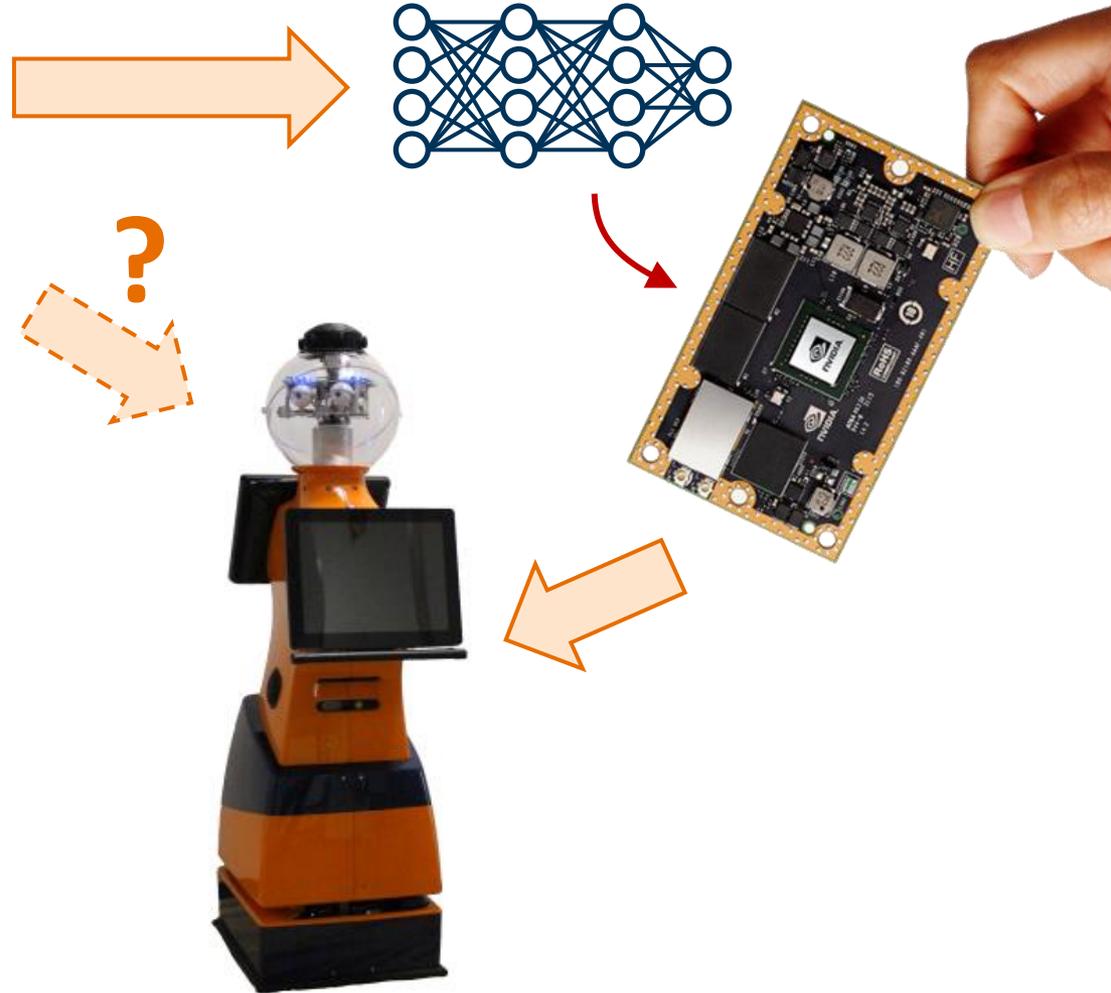


R. Stricker, D. Seichter, A. Vorndran, T. Wengefeld, H.-M. Gross  
Ilmenau University of Technology  
Germany

# Motivation

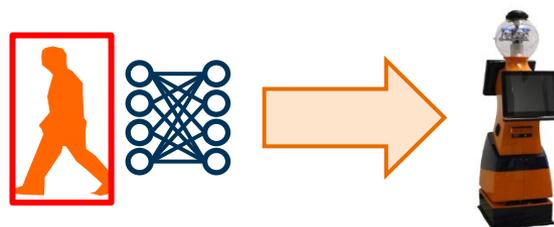


[IJCNN 2016]



# Motivation

Why do we need a DL-based person detector on a robot?



# Why a DL-based person detector?

– Motivation from a robotic viewpoint



Do state of the art person detectors perform well in real-world applications?

## Example: Navigation

Many **false detections near doors** (vertical gradients)



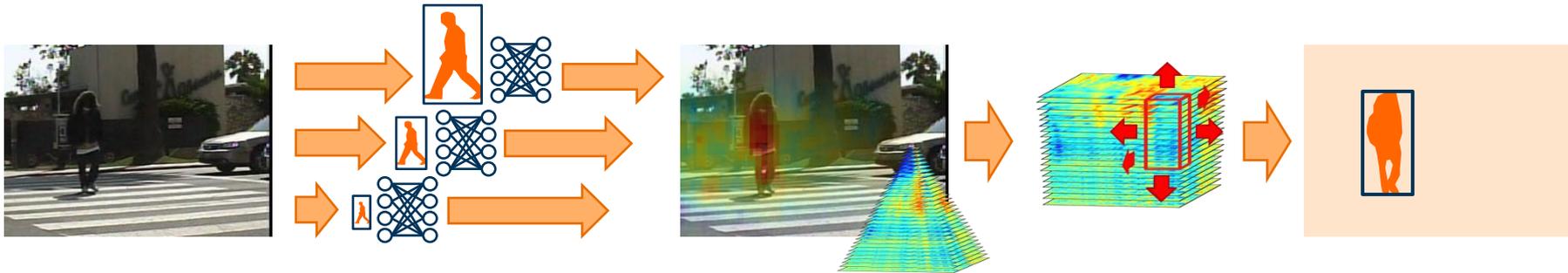
→ **Robot does not pass** the door due to possible personal space violation

# Outline

What can you expect?

# Outline

## Person detector in a nutshell

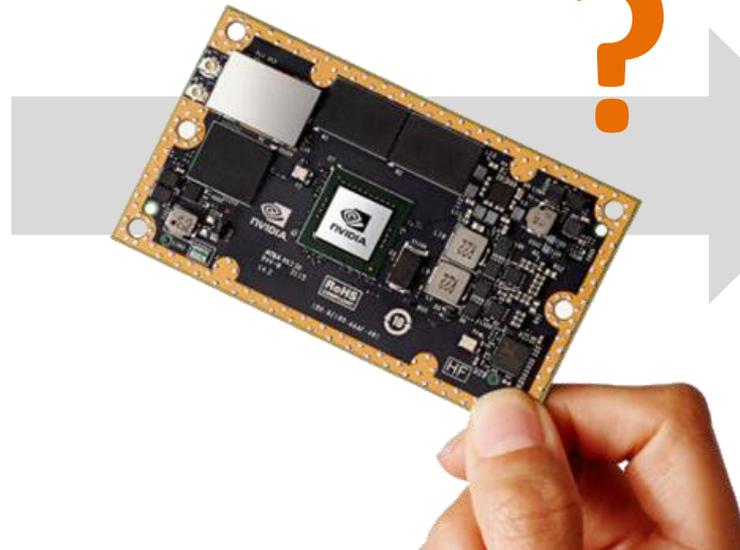


## Is the Jetson TX1 really necessary?



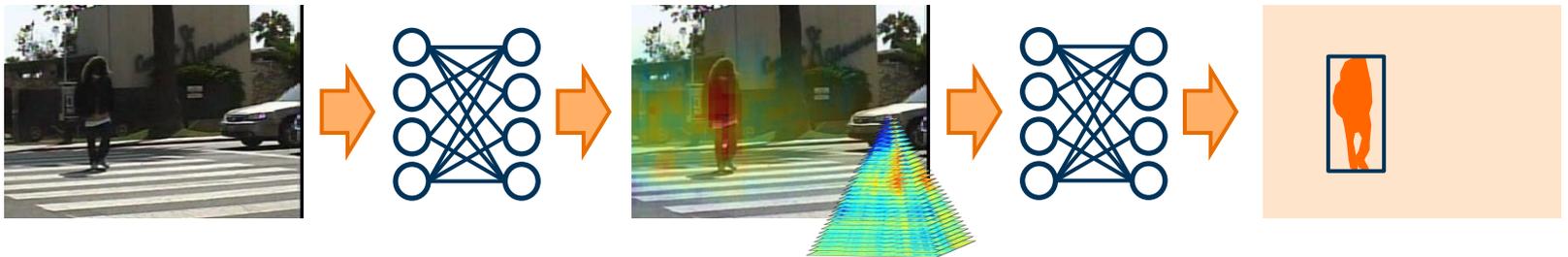
# Outline

How can the run time be improved on this device without a loss in accuracy?

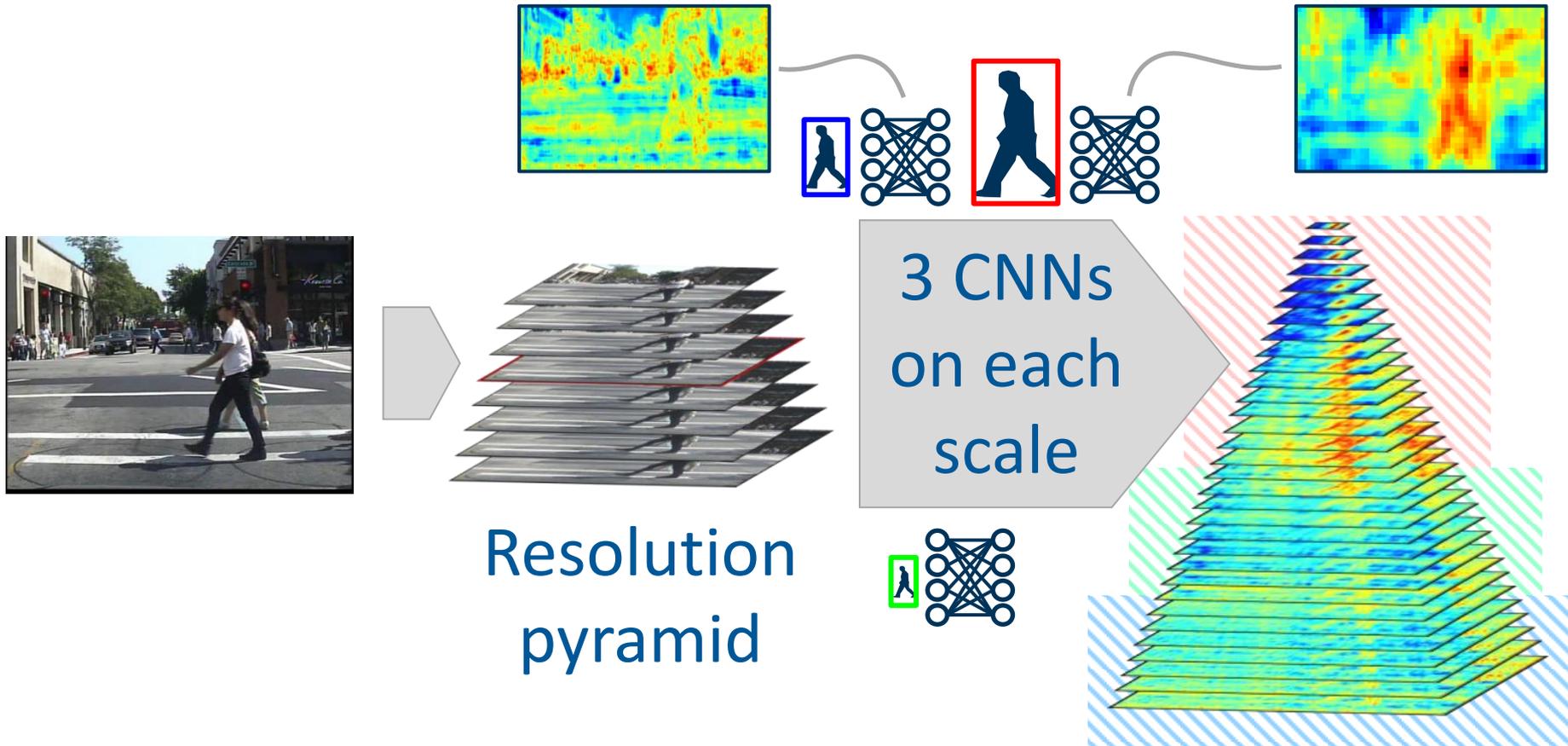


# Our person detector in a nutshell

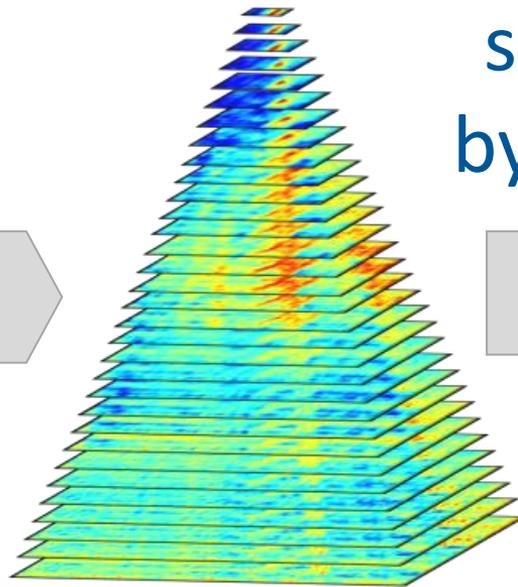
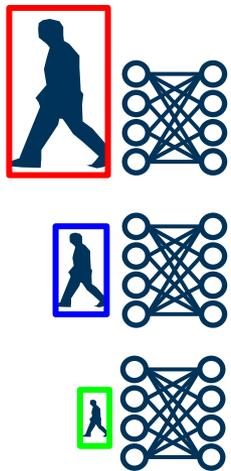
What is the basic idea?



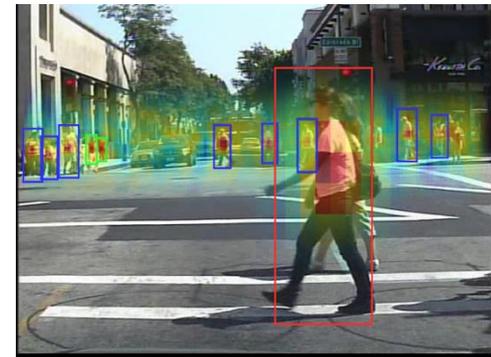
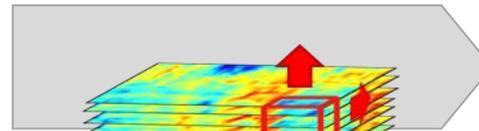
# Person detector in a nutshell



# Person detector in a nutshell



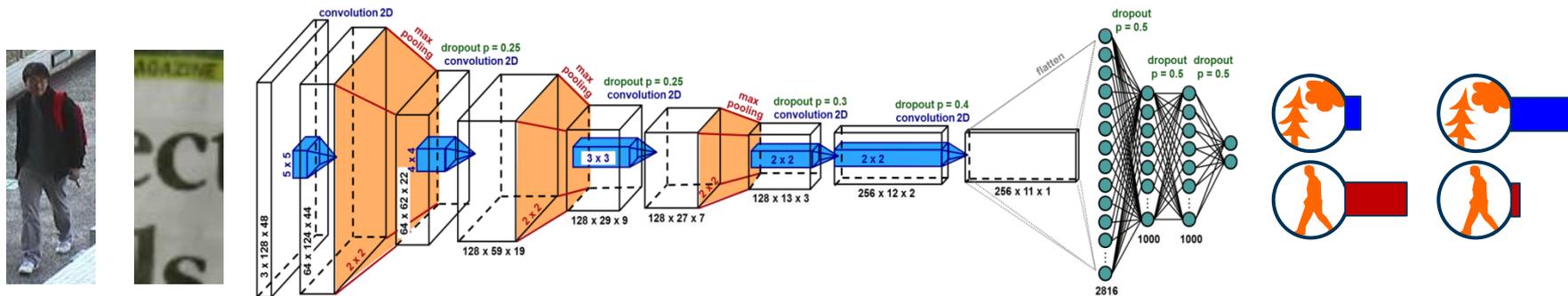
Non-maximum  
suppression  
by 3D pooling



# Person detector in a nutshell

Input 128 x 48 x 3

Output

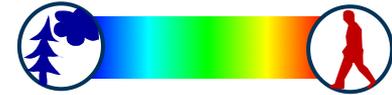


5 Conv layers: sizes 2x2 – 5x5

(1 – 3) 2x2 Max pooling layers with stride 2

2 FC layers + Softmax output layer → Conv layers

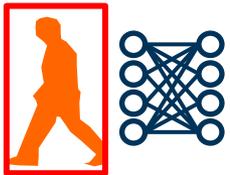
# Detector – Visual Examples



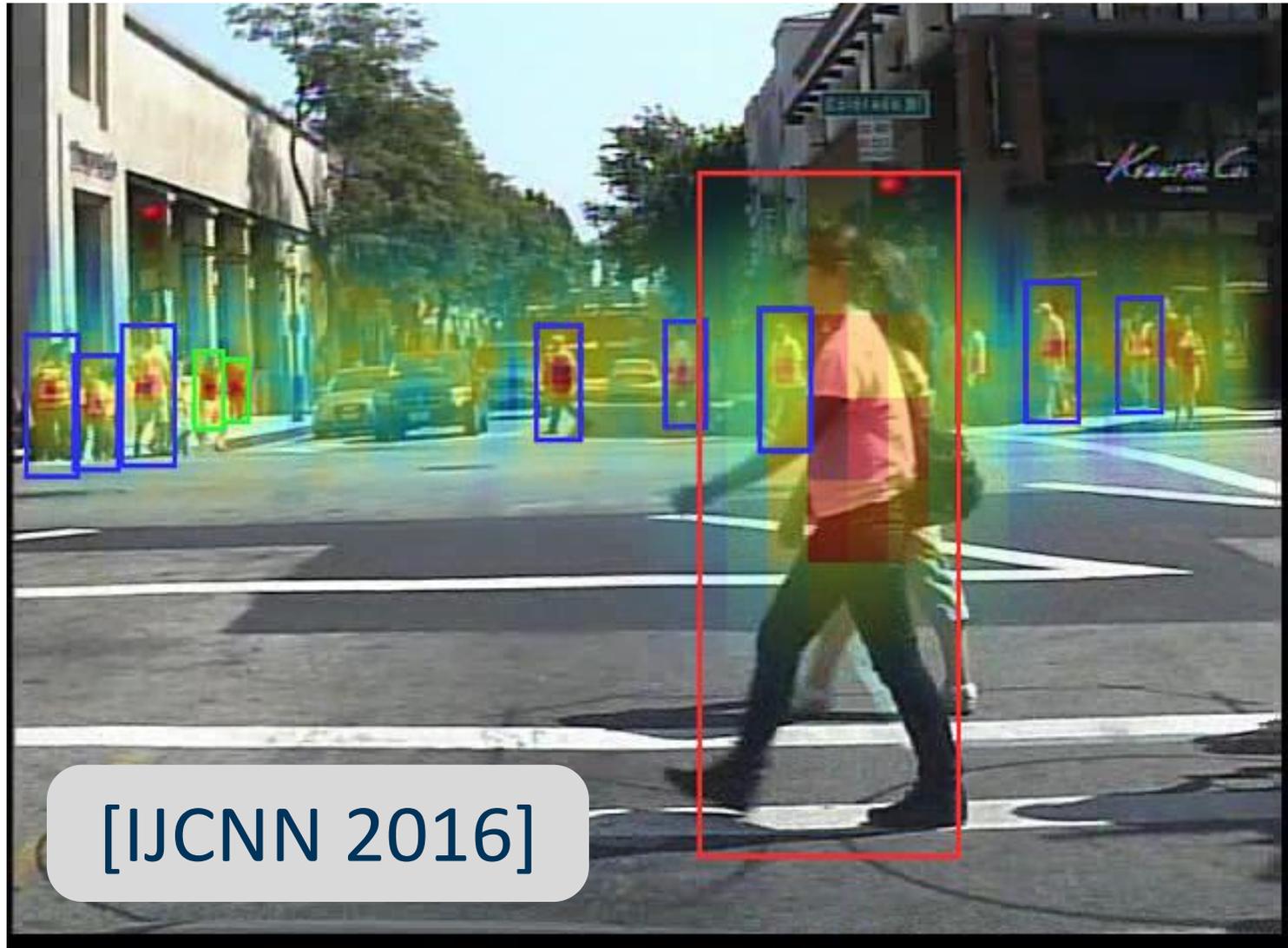
Far



Medium

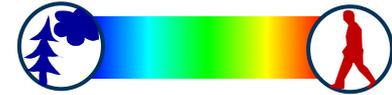


Near



[IJCNN 2016]

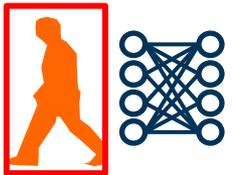
# Detector – Visual Examples



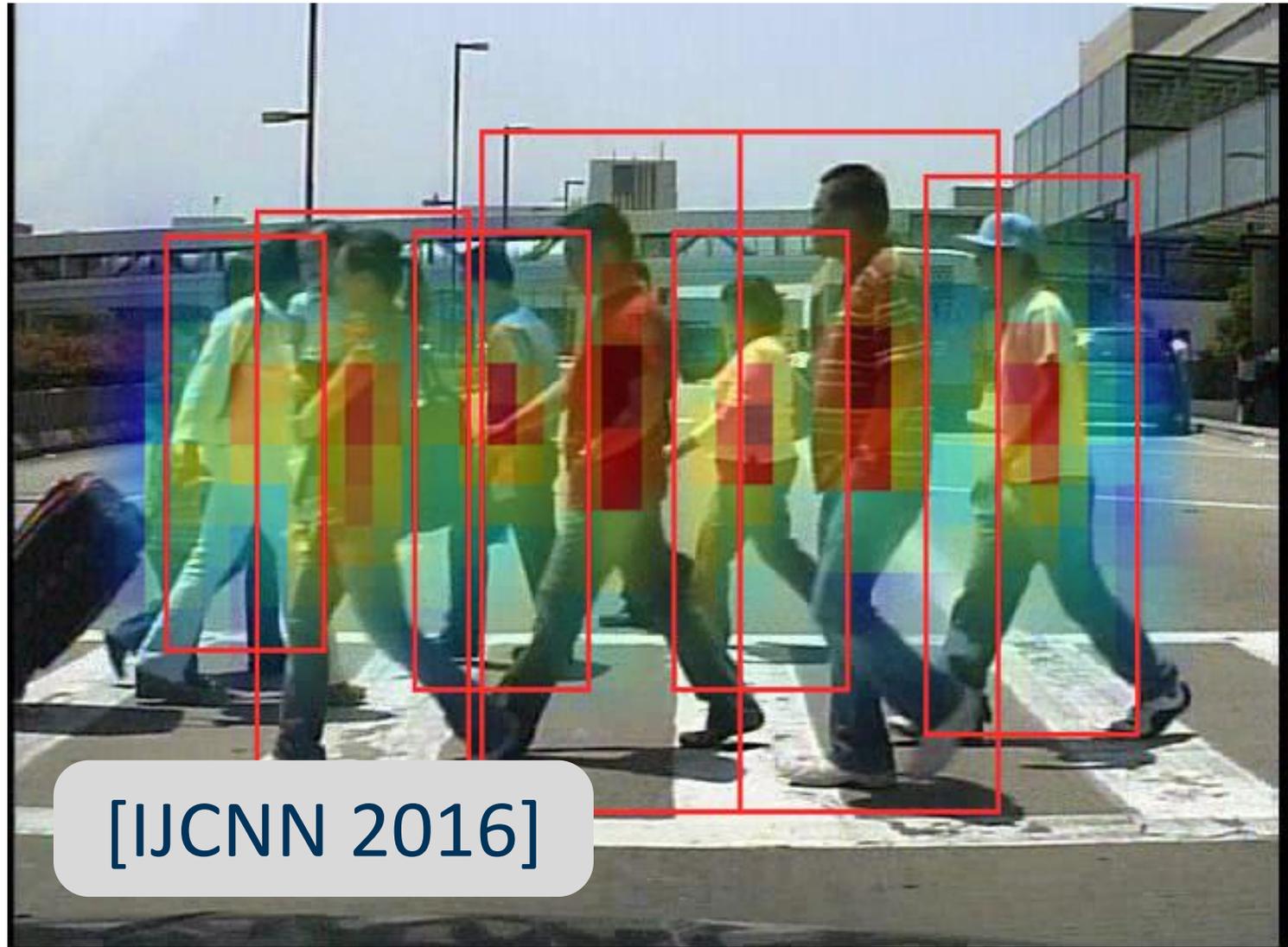
Far



Medium



Near



[IJCNN 2016]

# Why do we need the Jetson TX1?



# Why do we need the Jetson TX1?



# Why do we need the Jetson TX1?



# Why do we need the Jetson TX1?

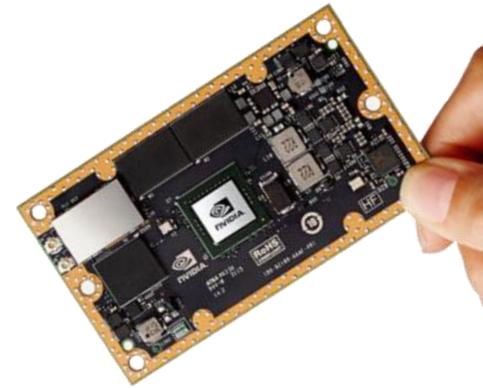


# ... and run time?



8.391 s

Jetson  
TX1



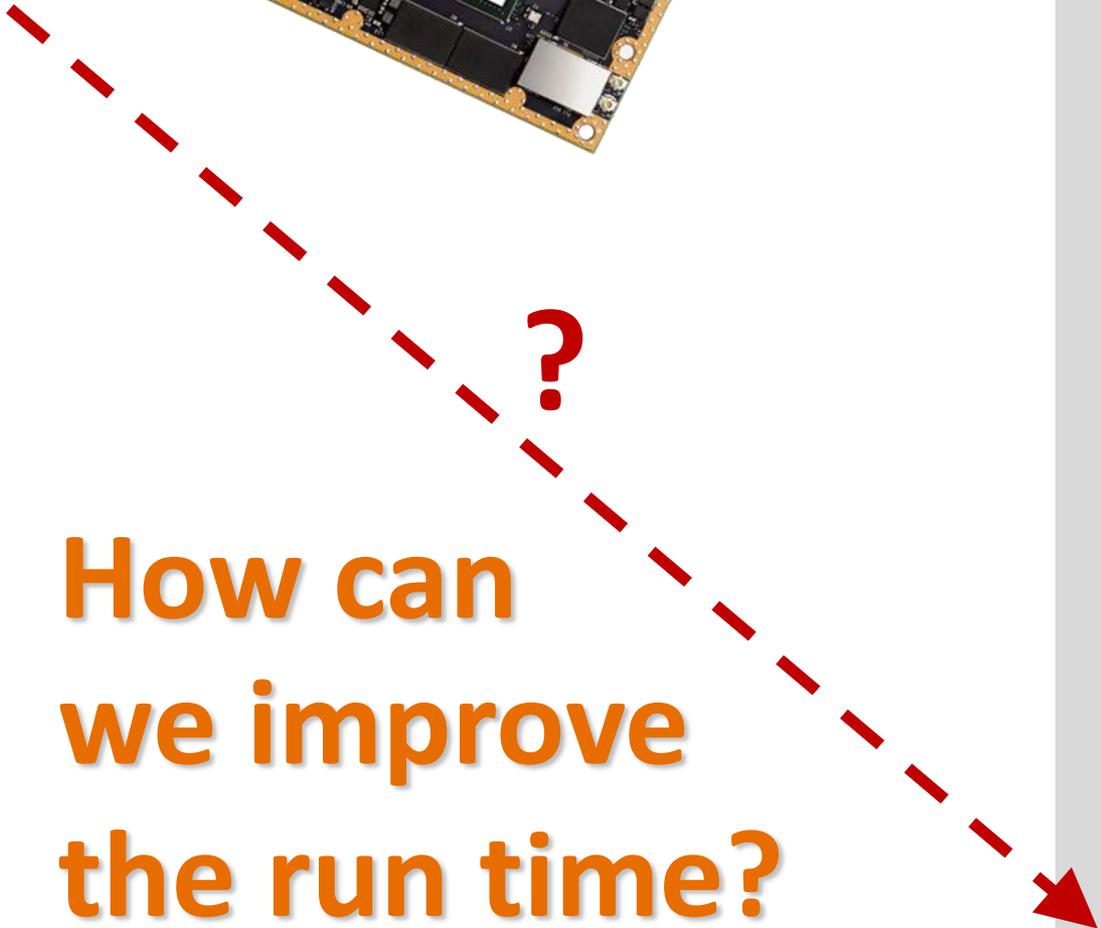
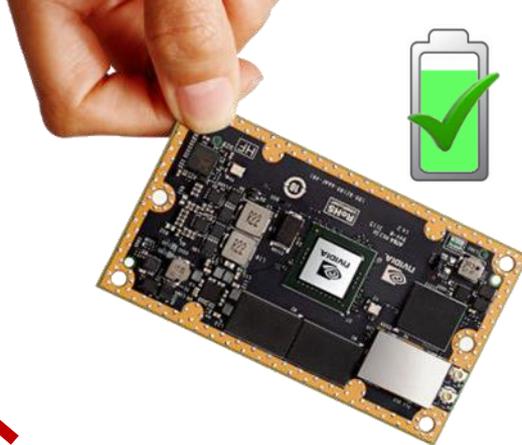
0.462 s

Titan X



8.391 s

Jetson  
TX1



How can  
we improve  
the run time?



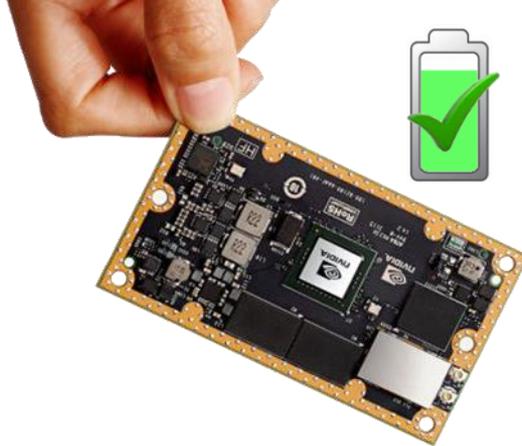
0.462 s

Titan X



8.391 s

Jetson  
TX1



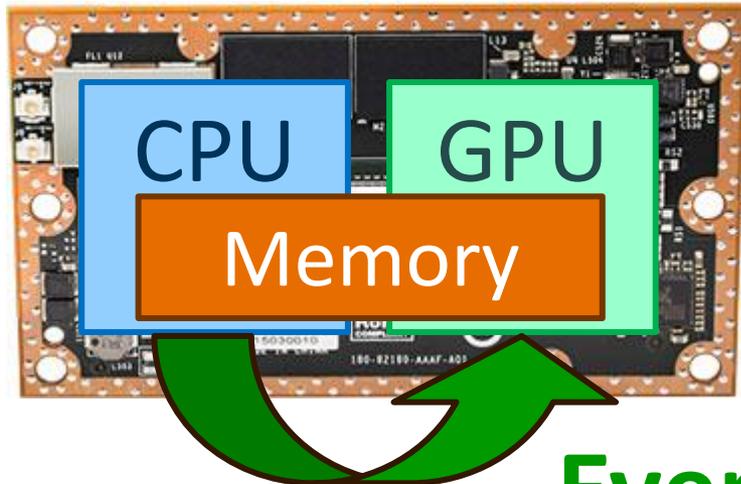
# Computation graph optimization and overhead removal



0.462 s

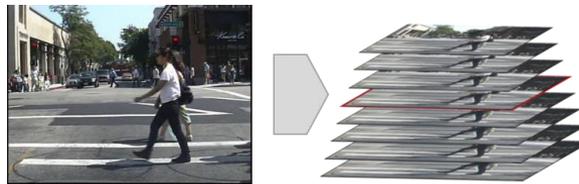
Titan X

# Computation graph optimization and overhead removal

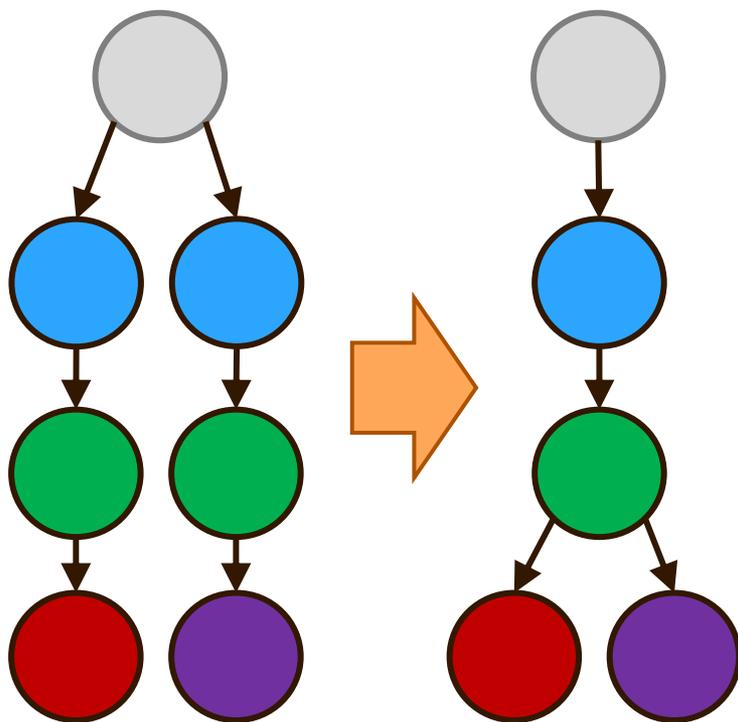


Copy operations between CPU and GPU (Frameworks)

**Every operation on GPU (incl. non-DL operations)**



# Computation graph optimization and overhead removal

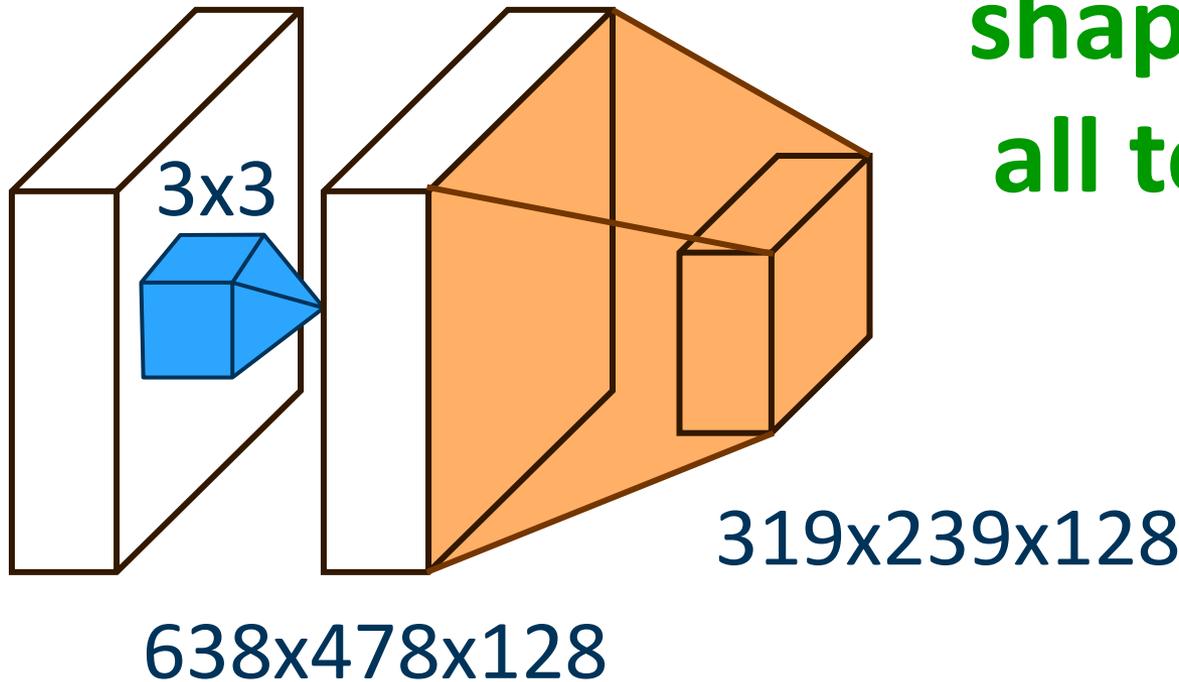


**Remove  
redundancies from  
the computation  
graph**

# Computation graph optimization and overhead removal



640x480x3

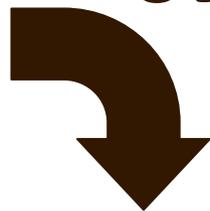


Specify exact shapes to make all tensor sizes static



8.391 s

Jetson  
TX1



67%



2.8 s

Remove  
Overhead



Reduce  
floating point  
precision



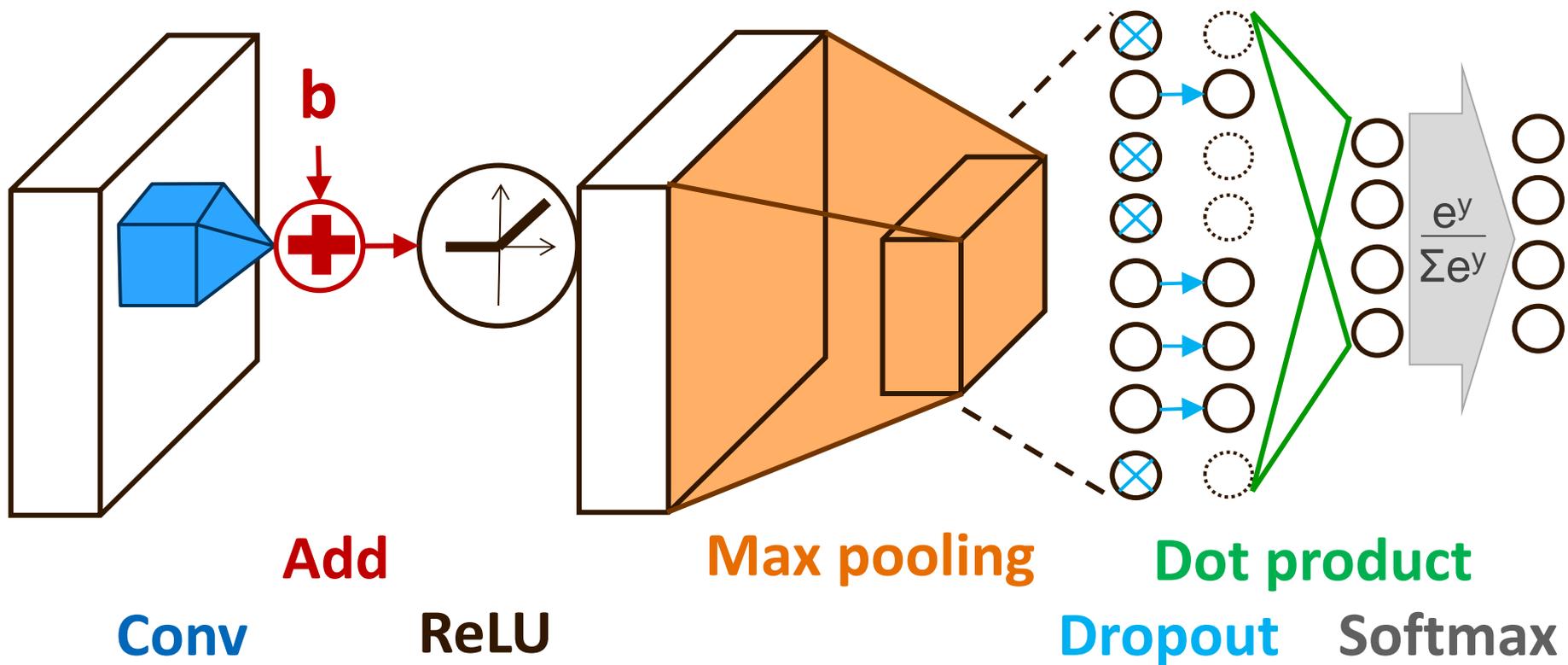
0.462 s

Titan X

# Reducing floating point precision

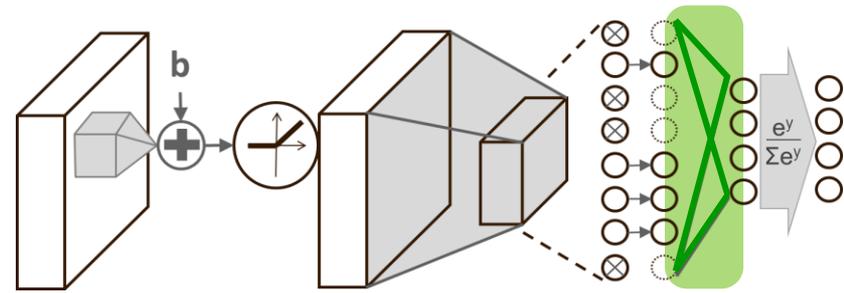
float 32 → float 16

## Typical Deep Learning operations:

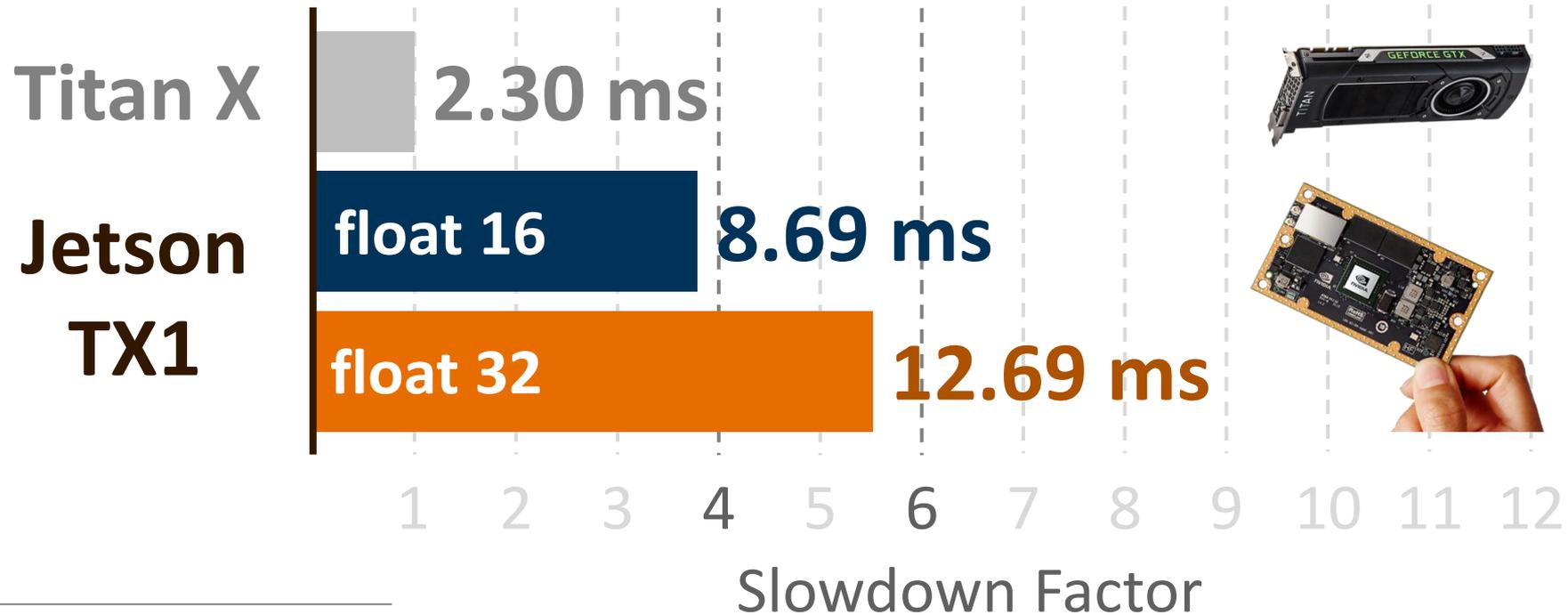


# Reducing precision

float 32 → float 16



## Dot product (fully connected layers)



(Matrix: 1000x1000)

# Reducing floating point precision

## Summary

Slowdown factor



float 32



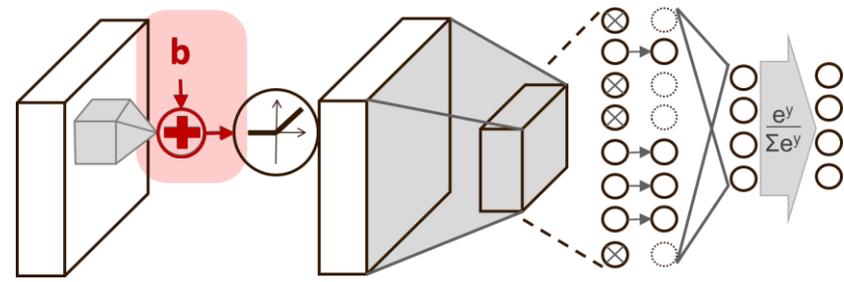
float 16

Matrix multiplication

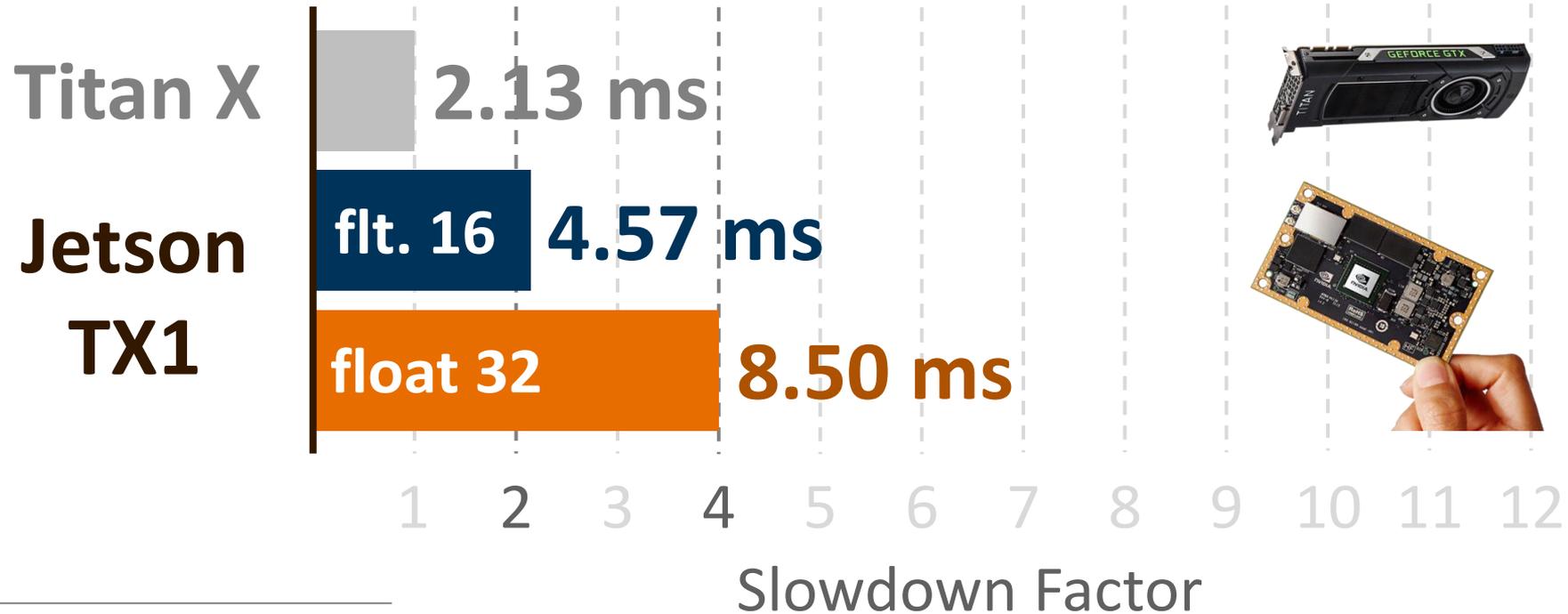
6  $\xrightarrow{-33\%}$  4

# Reducing precision

float 32 → float 16



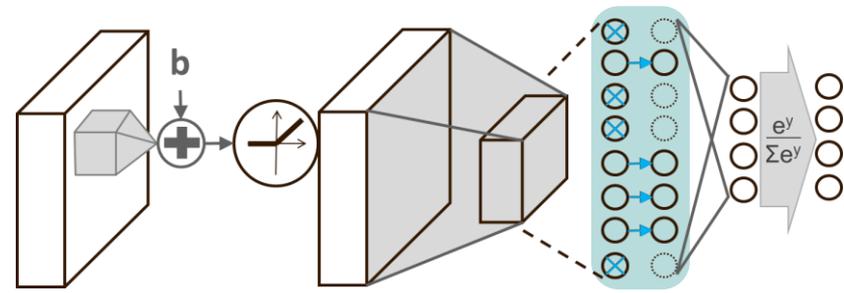
## Element-wise add (bias)



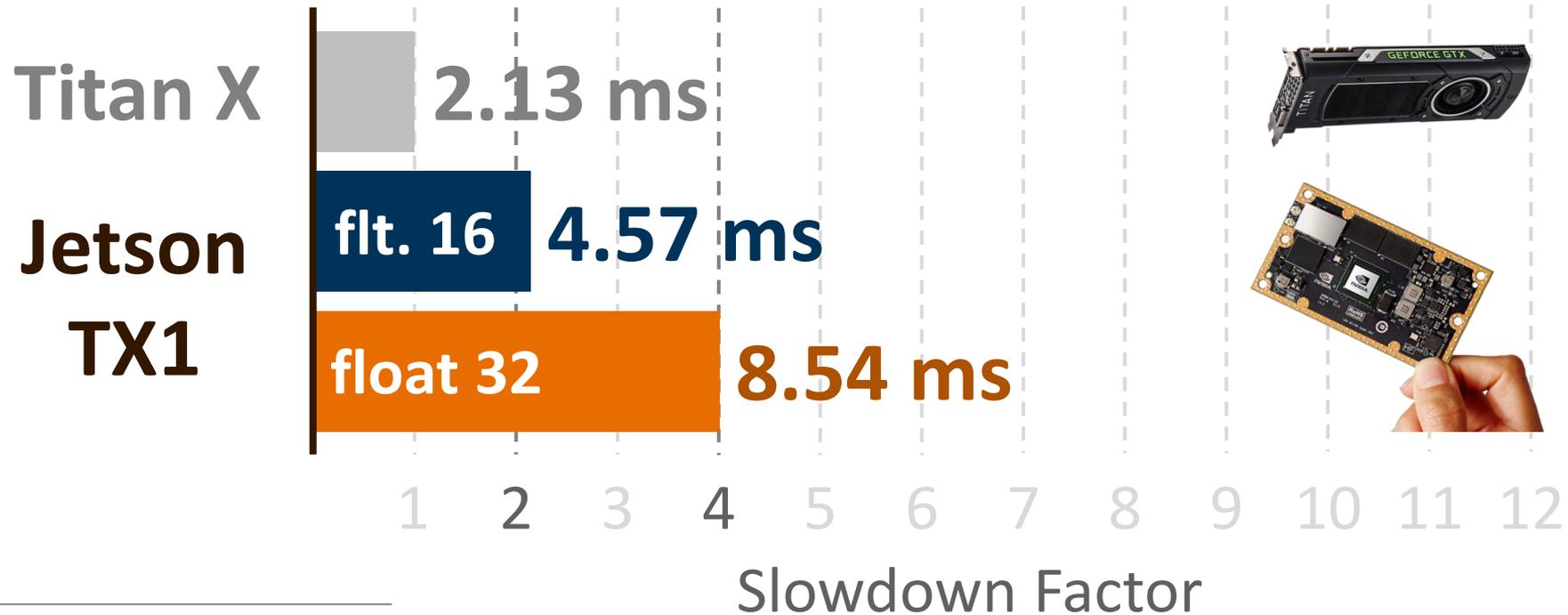
(Matrix: 1000x1000)

# Reducing precision

float 32 → float 16



## Element-wise multiplication (dropout)

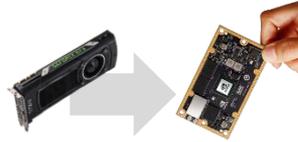


(Matrix: 1000x1000)

# Reducing floating point precision

## Summary

Slowdown factor



**float 32**



**float 16**

Matrix multiplication

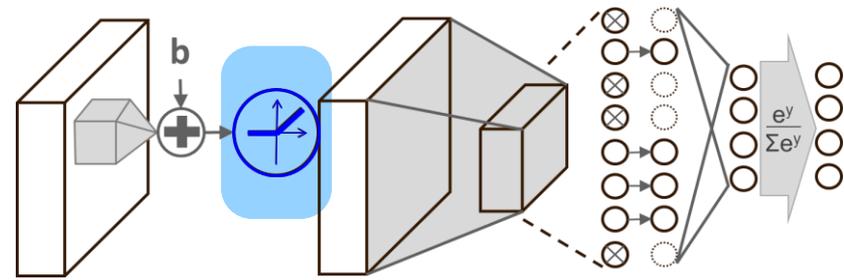
6  $\xrightarrow{-33\%}$  4

Element-wise operations

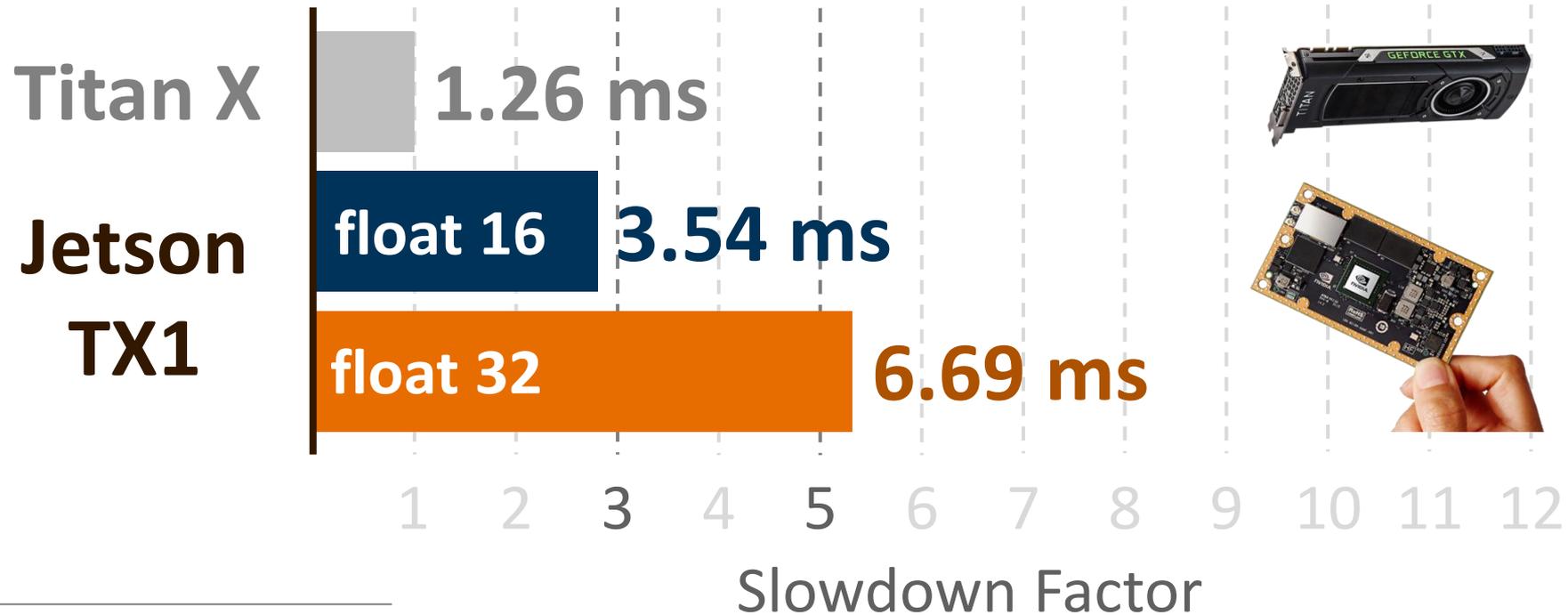
4  $\xrightarrow{-50\%}$  2

# Reducing precision

float 32 → float 16



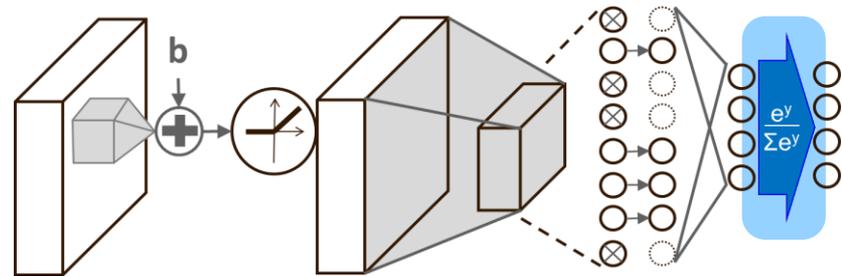
## Element-wise non-linearity (ReLU activation)



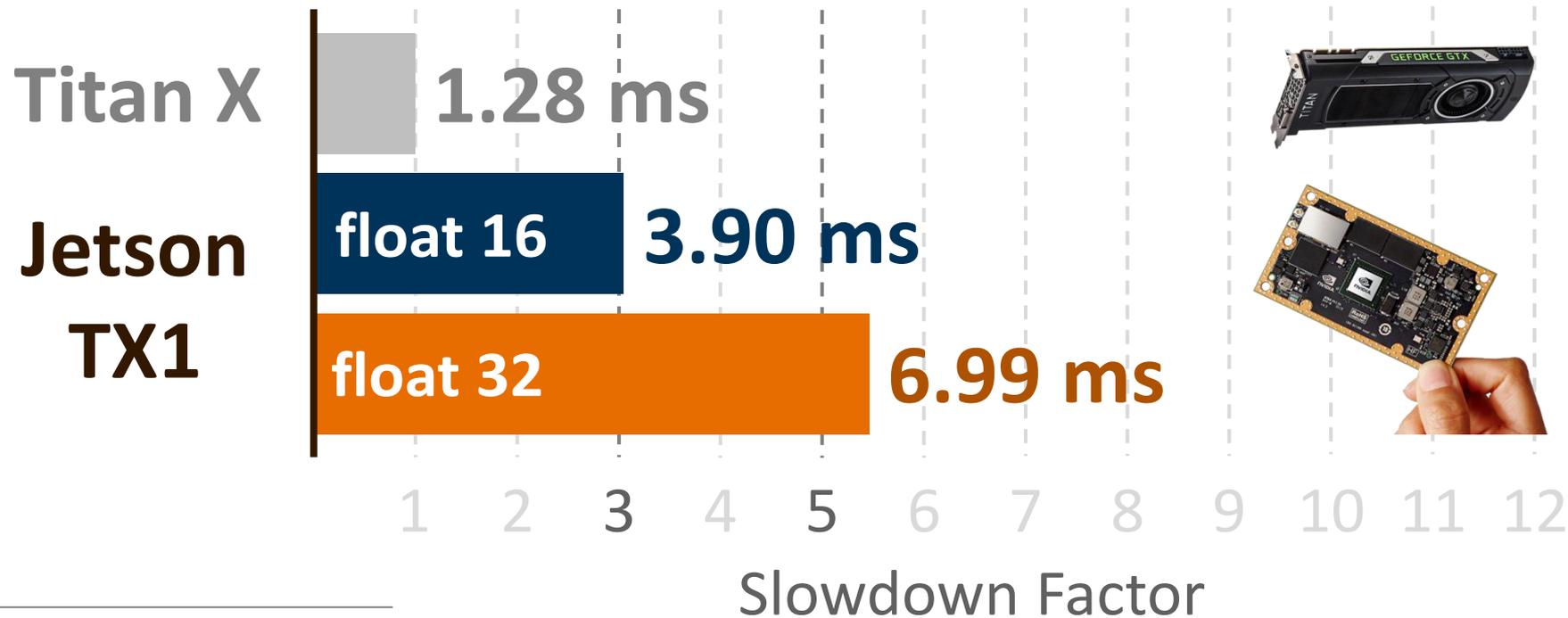
(Matrix: 1000x1000)

# Reducing precision

float 32 → float 16



## Element-wise non-linearity (Softmax output)



(Matrix: 1000x1000)

Slowdown Factor

# Reducing floating point precision

## Summary

Slowdown factor



**float 32**



**float 16**

Matrix multiplication

6  $\xrightarrow{-33\%}$  4

Element-wise operations

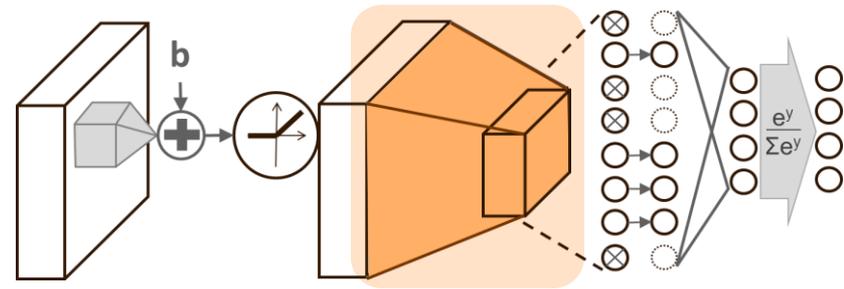
4  $\xrightarrow{-50\%}$  2

Non-linear activations

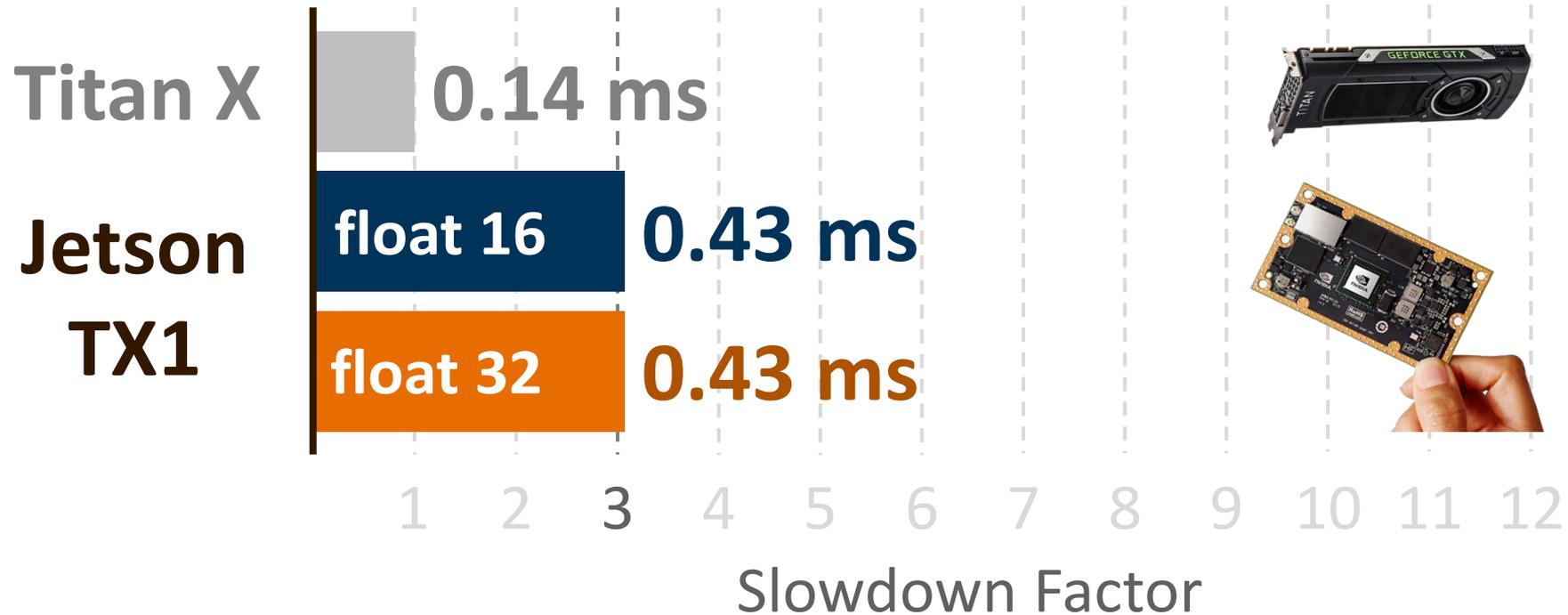
5  $\xrightarrow{-40\%}$  3

# Reducing precision

float 32 → float 16



## Max pooling



Input Volume: 3 x 320 x 240, 200 filters, pooling/stride 2x2

# Reducing floating point precision

## Summary

Slowdown factor



**float 32**



**float 16**

Matrix multiplication

6  $\xrightarrow{-33\%}$  4

Element-wise operations

4  $\xrightarrow{-50\%}$  2

Non-linear activations

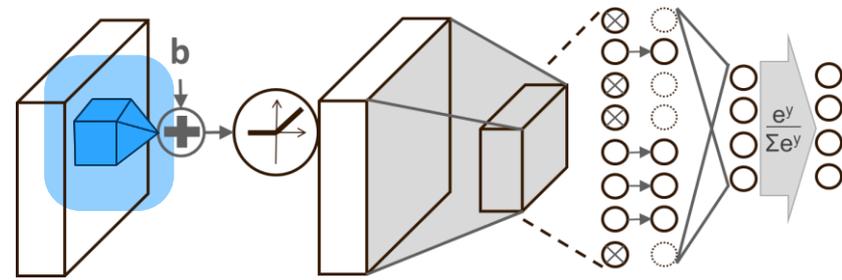
5  $\xrightarrow{-40\%}$  3

**Pooling**

**3**  $\xrightarrow{=}$  **3**

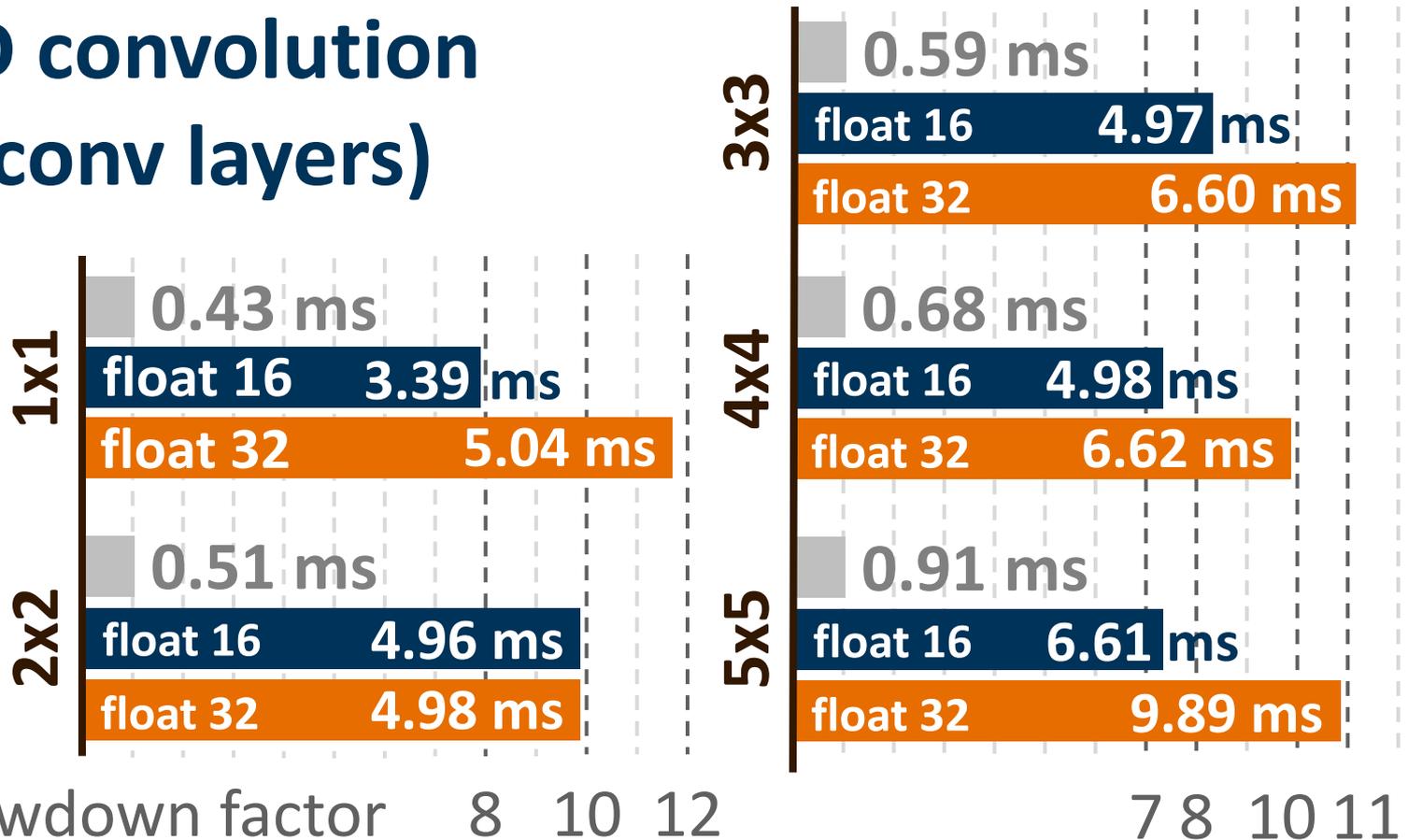
# Reducing precision

float 32 → float 16



## 2D convolution (conv layers)

Titan X  
Jetson  
TX1



Input Volume: 3 x 320 x 240, no padding, 200 filters, batch size 1

# Reducing floating point precision

## Summary

Slowdown factor



**float 32**



**float 16**

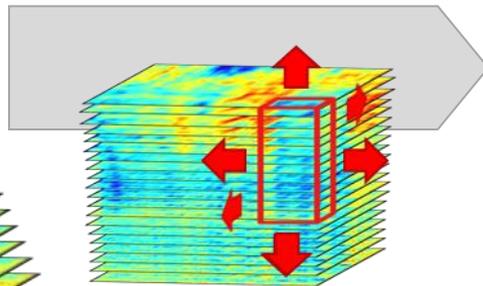
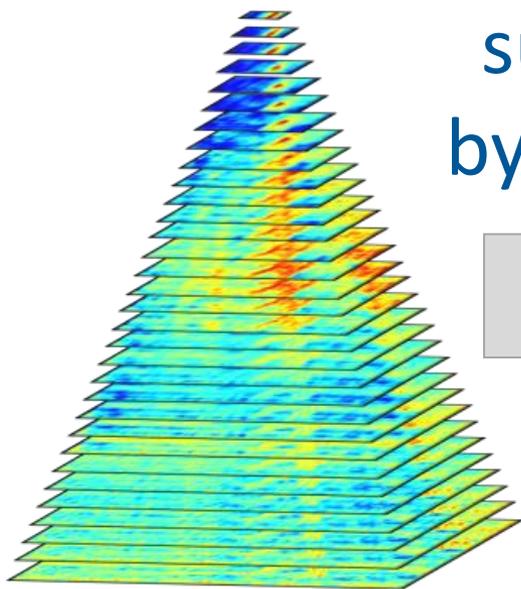
Matrix multiplication	6	$\xrightarrow{-33\%}$	4
Element-wise operations	4	$\xrightarrow{-50\%}$	2
Non-linear activations	5	$\xrightarrow{-40\%}$	3
Pooling	3	$\xrightarrow{=}$	3
<b>Convolutions</b>	<b>10-11</b>	$\xrightarrow{-25-30\%}$	<b>7-8</b>

# Reducing floating point precision

float 32 → float 16

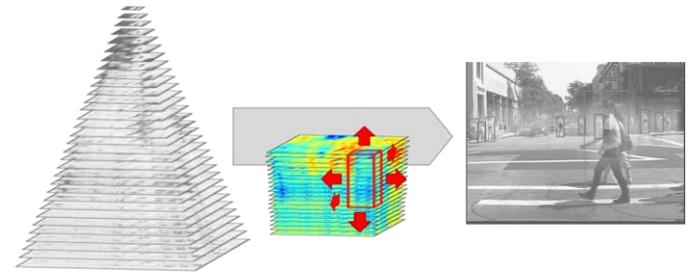
## Rarely used Deep Learning operations

Non-maximum  
suppression  
by **3D pooling**

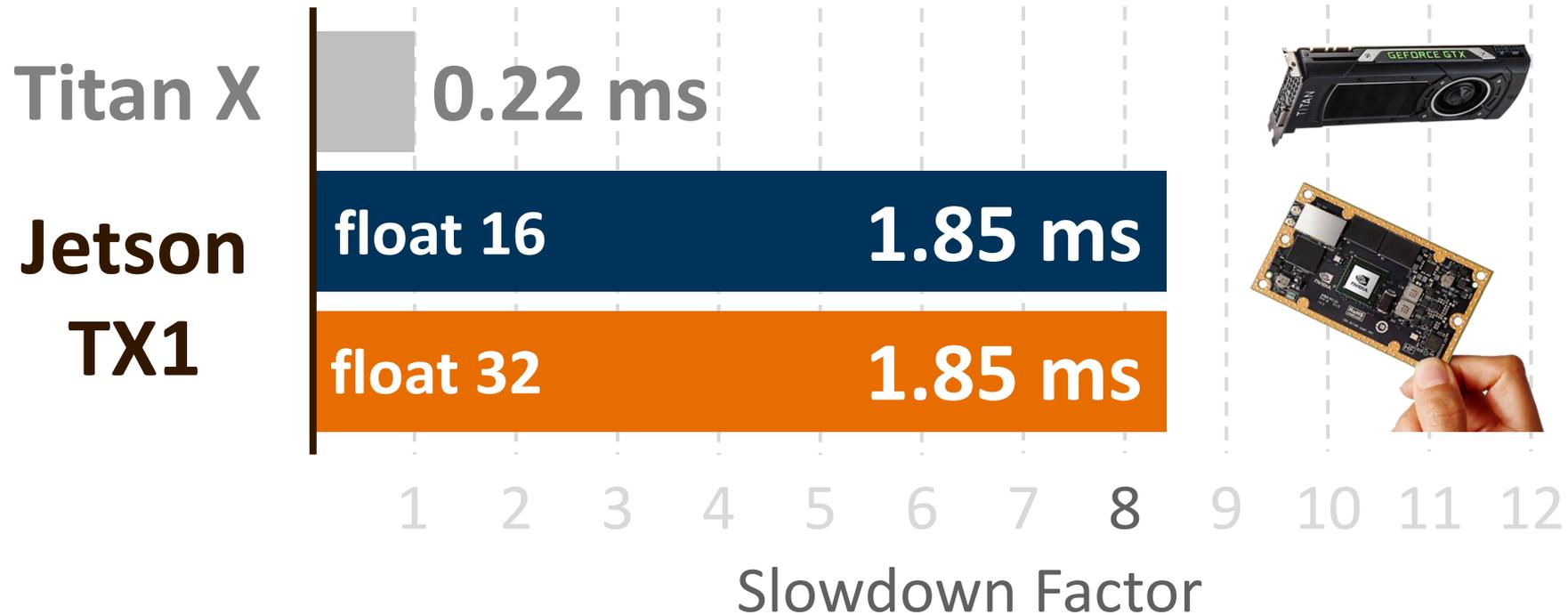


# Reducing precision

float 32 → float 16



## 3D pooling (Non-maximum suppression)



Input Volume: 3 x 320 x 240, pooling 3x7x3, stride 1x1x1

# Reducing floating point precision

## Summary

Slowdown factor



**float 32**



**float 16**

Matrix multiplication

6  $\xrightarrow{-33\%}$  4

Element-wise operations

4  $\xrightarrow{-50\%}$  2

Non-linear activations

5  $\xrightarrow{-40\%}$  3

Pooling

3  $\xrightarrow{=}$  3

Convolutions

10-11  $\xrightarrow{-25-30\%}$  7-8

**3D Pooling**

**8  $\xrightarrow{=}$  8**

# Reducing floating point precision

## Summary

Slowdown factor



float 32

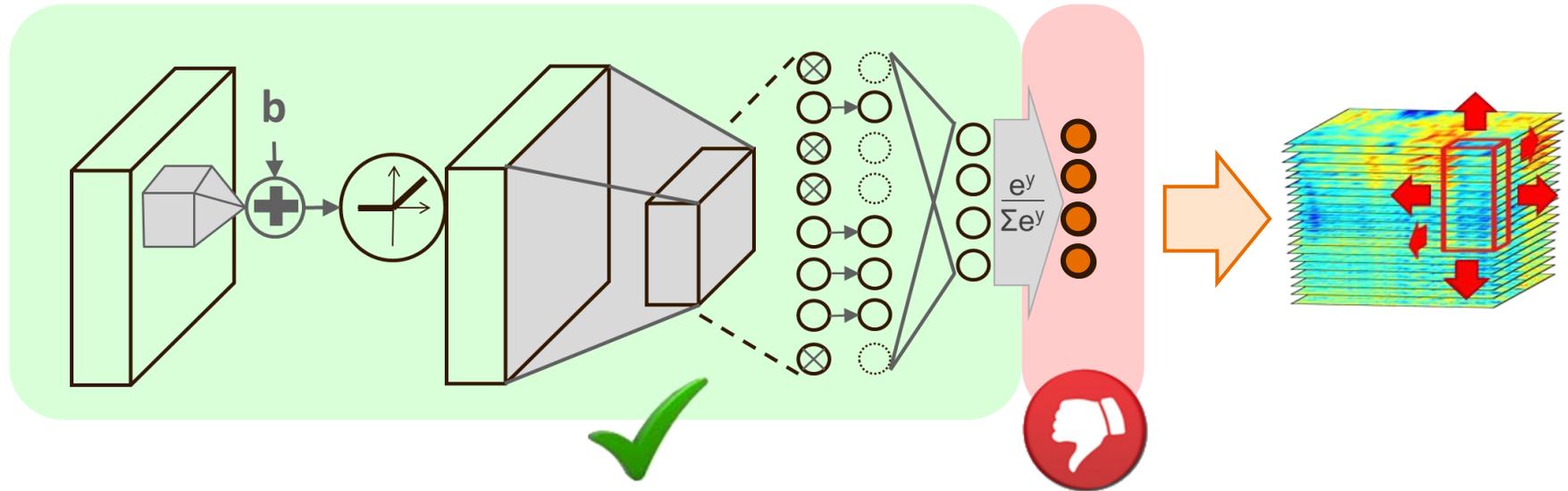


float 16

Matrix multiplication	6	- 33%	4
Element-wise operations	4	- 50%	2
Non-linear activations	5	- 40%	3
Pooling	3	=	3
Convolutions	10-11	- 25-30%	7-8
3D Pooling	8	=	8
Overall	10	- 25%	7.5

>90%

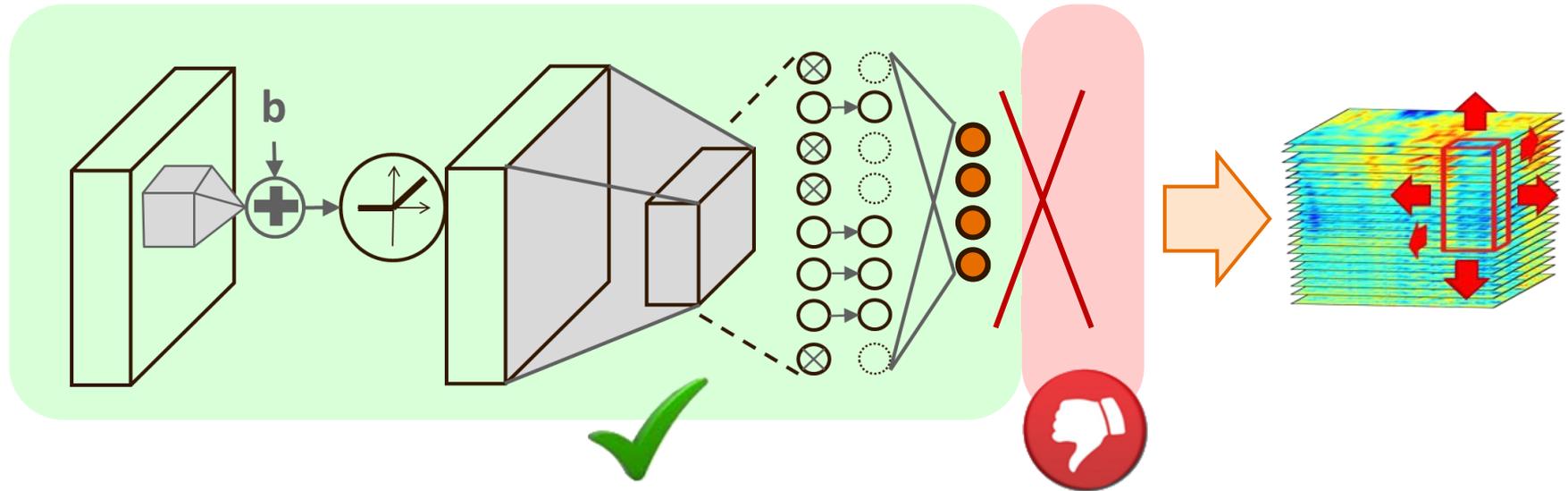
# Loss of precision / accuracy ?



Loss of precision after exponential function

Problematic for NMS

# Loss of precision / accuracy ?



Loss of precision after exponential function

Problematic for NMS

Solution: Use linear output for NMS

# Loss of precision / accuracy ?

## Which precision do you need?

→ It depends on dataset

[Rzayev et al., IJCNN 2017]

On ImageNET:



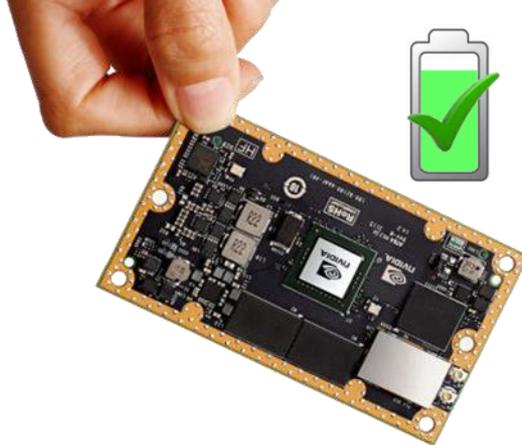
Model	32-float	12-bit	10-bit	8-bit	6-bit
VGG19	<b>71.36%</b>	71.35%	<b>71.34%</b>	70.88%	56.00%
ResNet152	<b>77.55%</b>	77.51%	<b>77.40%</b>	74.95%	9.29%
SqueezeNetV1	<b>56.52%</b>	56.52%	<b>56.24%</b>	54.56%	17.10%
InceptionV3	<b>76.41%</b>	76.43%	<b>76.44%</b>	73.67%	1.50%

<https://github.com/aaron-xichen/pytorch-playground/blob/master/README.md>



8.391 s

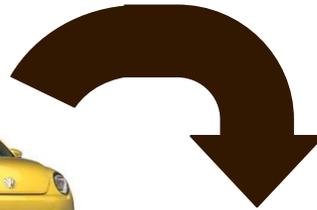
Jetson  
TX1



2.8 s

Remove  
Overhead

25%



2.103 s

Float 16

Task specific:  
ground plane  
assumption

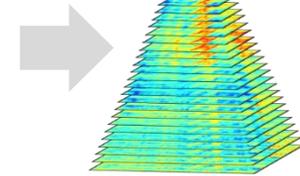
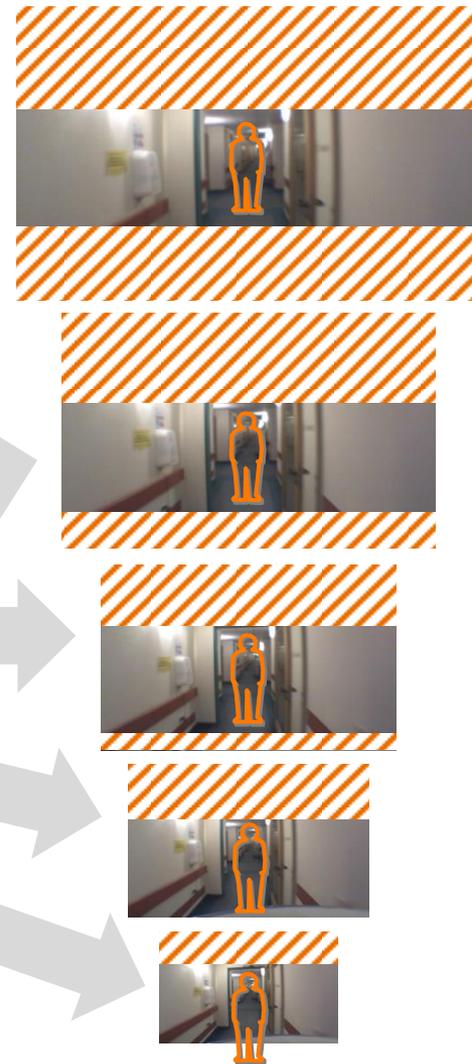
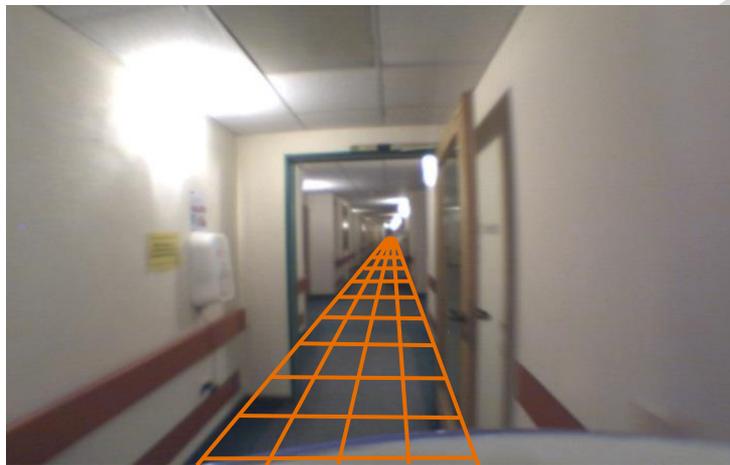


0.462 s

Titan X

# Ground plane assumption

**Assumption:**  
Every person stands  
on the ground



→ **Less computations**



8.391 s

Jetson  
TX1



2.8 s

Remove  
Overhead



2.103 s

Float 16



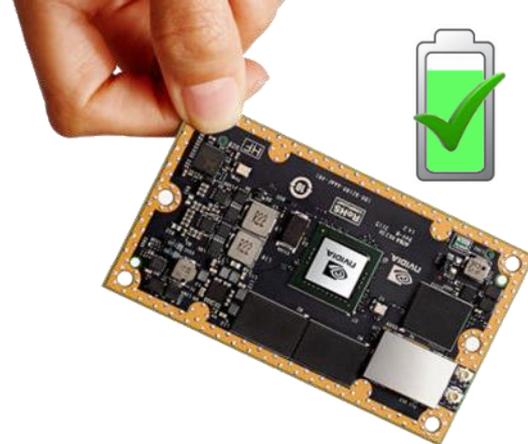
0.588 s

Ground  
Plane



0.462 s

Titan X





8.391 s

Jetson  
TX1



2.8 s

Remove  
Overhead



2.103 s

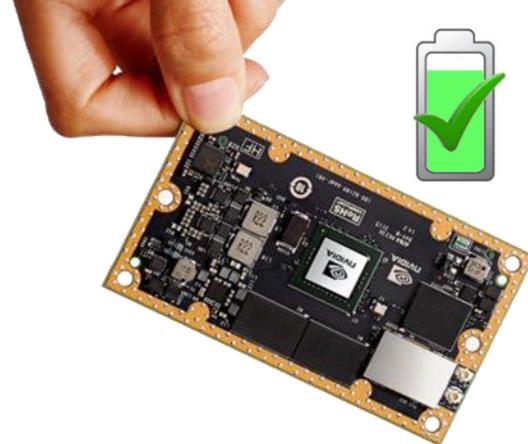
Float 16



0.588 s  
**Ground  
Plane**

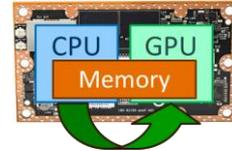


0.462 s  
**Titan X**

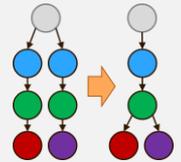


# Summarized: What can you do to improve the run time?

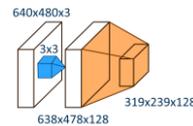
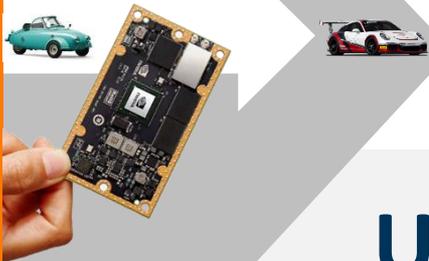
More operations on GPU



Remove redundancies



Specify exact tensor shapes



Use float 16 precision

Jetson TX1	float 16	3.90 ms
	float 32	6.99 ms

Avoid needless computations



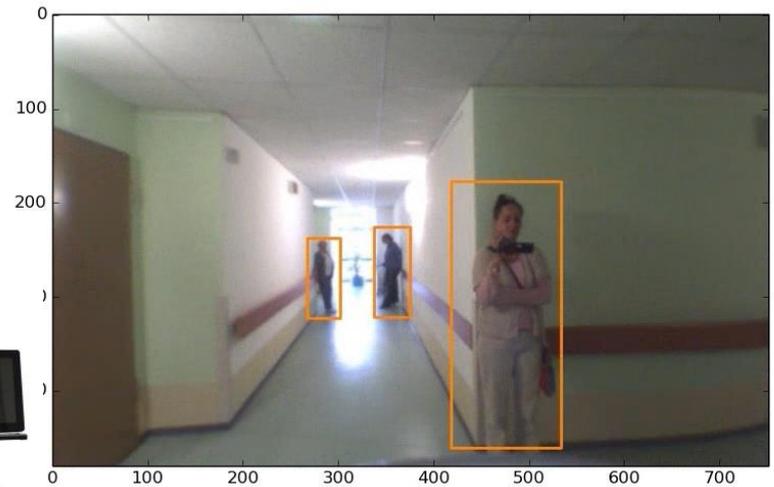
# Performance

How does the detector  
perform on a robot?

# Performance: video example



Front camera



Back camera

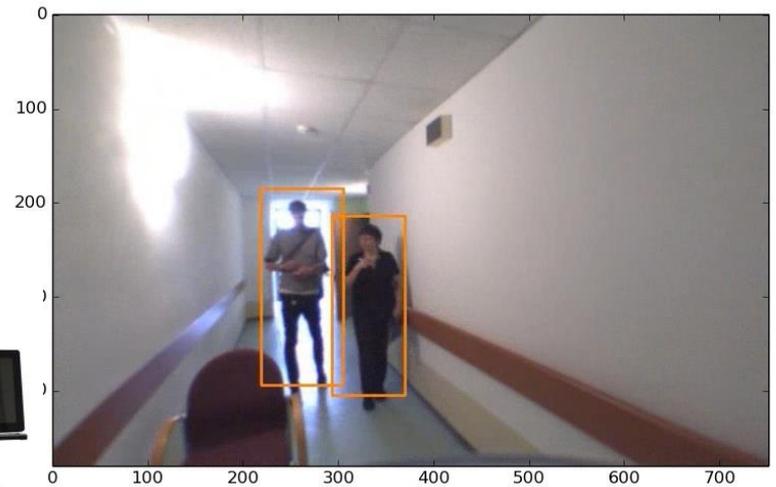
Detections



# Performance: video example



Front camera



Back camera

Detections



# Performance: video example



Front camera

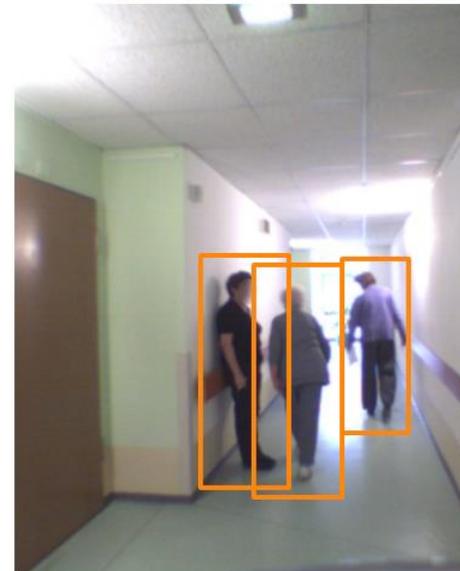
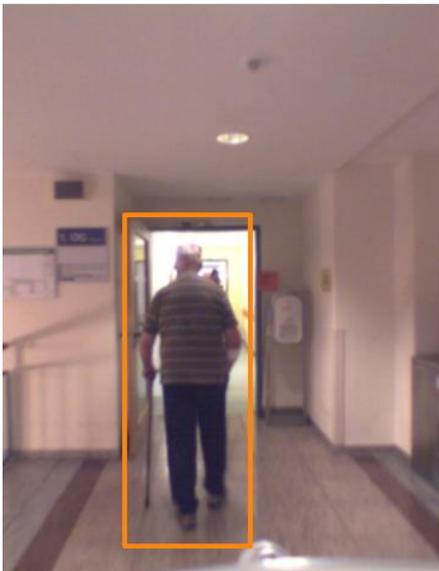
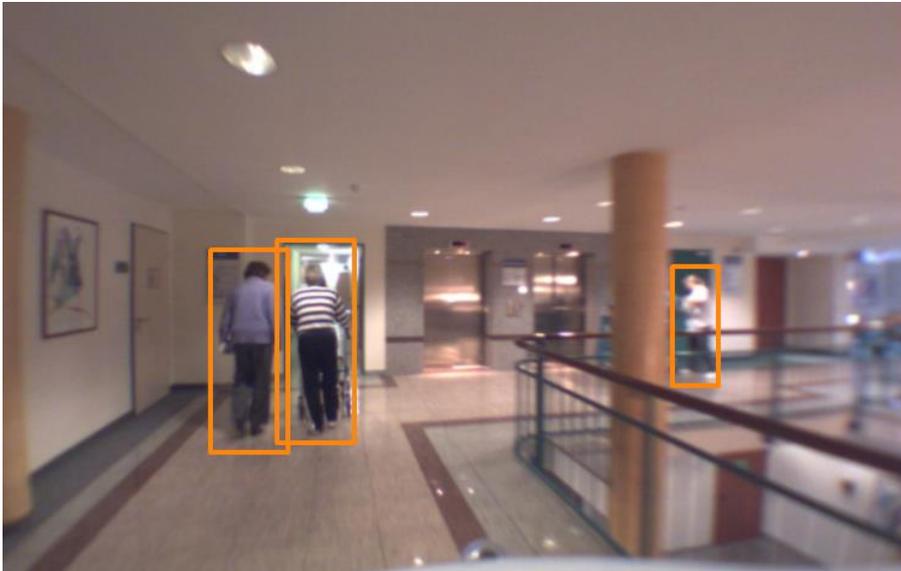


Back camera

Detections



# Performance: visual examples



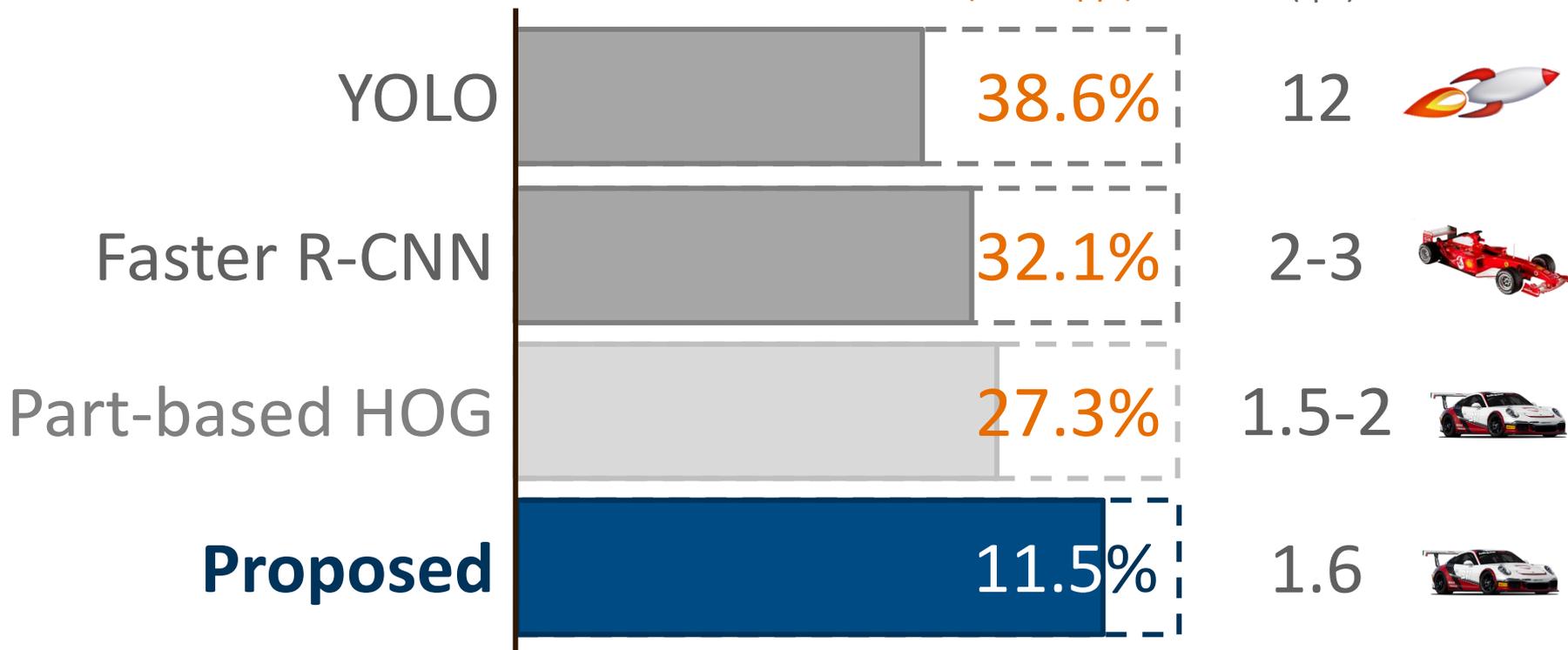
# Performance in numbers

(Robotic dataset)



Miss rate  
(@0.1 fppi)

Speed  
(fps)



# Conclusion

# Conclusion Speeding up Deep Neural Networks on the Jetson TX1

67%

Process all on GPU, optimize graph

25%

Use float16 instead of float32

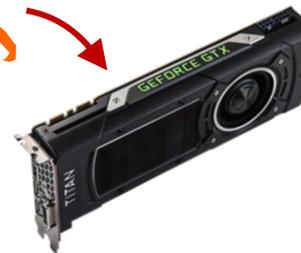
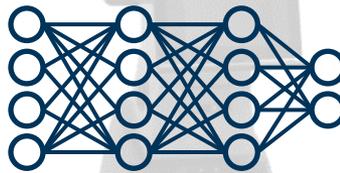
75%

Use assumptions for less computations



0.588 s

Jetson TX1



0.462 s

Titan X