Nvidia Jetson TX2 and its Software Toolset



~

João Fernandes 2017/2018



In this presentation

- Nvidia Jetson TX2: Hardware
- Nvidia Jetson TX2: Software
- Machine Learning: Neural Networks
 - Convolutional Neural Networks
- Tensorflow
- A case study





NVIDIA Jetson TX2



NVIDIA Jetson TX2 Developer Kit



- "Jetson is the world's leading low-power embedded platform, enabling server-class AI compute performance for edge devices everywhere"
- "Edge computing is an emerging paradigm which uses local computing to enable analytics at the source of the data"
- "Jetson TX2 accelerates cutting-edge deep neural network (DNN) architectures using the NVIDIA cuDNN and TensorRT libraries"



Image Recognition Classification



Object Detection Localization

NVIDIA Jetson TX2



Segmentation Free Space



CPU	
GPU	
Memory	8G
Storage	
I/O	UART, SI

Hardware Specifications



NVIDIA Jetson TX2

ARM Cortex-A57 (quad-core) @ 2GHz +

NVIDIA Denver2 (dual-core) @ 2GHz

256-core Pascal @ 1300MHz (2 SM's)

B 128-bit LPDDR4 @ 1866Mhz | 59.7 GB/s

32GB eMMC 5.1

PI, I2C, I2S, GPIOs, HDMI 2.0, USB, Ethernet, CAN

Heterogeneous Multi-Processing: Cortex-A57 vs Denver2

	4xCortex - A57	2xDenver2				
ARM ISA	ARMv8 (32/64 bits)	ARMv8 (32/64 bits)				
Frequency	Up to 2 GHz	Up to 2 GHz				
L1 Cache	48 KB (Inst) + 32 KB (Data)	128KB (Inst) + 64 KB (Data)				
L2 Cache	2MB (shared between the 2 Chips)	2MB (shared between the 2 Chips)				
GFLOPS(FP32)	750	437				
Full coherency quaranteed across the 2 CPU Clusters						





Performance Modes

Mode	Mode Name	Denver2	ARM A57	GPU Freq.	
0	Max-N	2 @ 2.0 GHz	4 @ 2.0 GHz	1.30 GHz	
1	Max-Q	0	4 @ 1.2 GHz	0.85 GHz	
2	Max-P Core All	2 @ 1.4 GHz	4 @ 1.4 GHz	1.12 GHz	
3	Max-P ARM	0	4 @ 2.0 GHz	1.12 GHz	
4	Max-P Denver	2 @ 2.0 GHz	0	1.12 GHz	

Can be changed at run time using:

sudo nvpmodel -m <desired mode>

		N
GoogLeNet	Perf	
	Power Consumption	
	Efficiency (FPS/W)	
AlexNet	Perf	
	Power Consumption	
	Efficiency	

Max-Q vs Max-P

Max-N	Max-P	Max-Q			
290 FPS	253 FPS	196 FPS			
12.8 W	8.9 W	5.9 W			
22.7	28.5	33.2			
692 FPS	601 FPS	463 FPS			
12.4 W	8.6 W	5.6 W			
55.8	69.9	82.7			

The Software: NVIDIA JetPack

- A set of software tools to ease the development for the Jetson platform:
 - Flash Jetson Developer Kit with the latest OS image
 - Install developer tools for both host PC and Developer Kit
 - Install libraries and APIs
 - Samples and documentation
- Key Software:
 - OS: L4T (Based on Ubuntu)
 - CUDA 8
 - TensorRT 2.1
 - cuDNN 6.0

Requires Ubuntu 14 on the Host!

The Software

CUDA

and optimizing the performance of your applications

cuDNN

includes support for convolutions, activation functions and tensor transformations.

TensorRT 2.1

memory footprint for convolutional and deconv neural networks.

Multimedia API

Ο video decode, encode, format conversion and scaling functionality

• CUDA Toolkit provides a comprehensive development environment for C and C++ developers building GPU-accelerated applications. It includes a compiler for NVIDIA GPUs, math libraries, and tools for debugging

• CUDA Deep Neural Network library provides high-performance primitives for all deep learning frameworks. It

• TensorRT is a high performance deep learning inference runtime for image classification, segmentation, and object detection neural networks. It speeds up deep learning inference as well as reducing the runtime

The Jetson Multimedia API package provides low level APIs for image acquisition (via camers). It enables

Real World Utilization

Jetson Al Pipeline

Relative Performance (GTX 1070) Matrix-Mul FP32

Machine Learning

- training phase).
 - Decision tree learning
 - Clustering
 - Neural Networks

• Machine learning is the science of getting the computer to perform a certain task without explicitly programming it. Normally, this requires a sample dataset from which the algorithm can infer the model parameters (commonly known as the

input layer

hidden layer 1

Neural Networks

hidden layer 2

Neural Networks: The neuron

Convolutional Neural Networks

Inp	Input Volume (+pad 1) (7x7x3)) () () () () () () () () () (Filter W			
x[]	1,2,	.01	()	i Terre				w0	1.	1,0
0	0	0	0	0	0	8		L	1	-1
0	2	2	0	1	1	0		1	1	0
0.	0	1	1	t	0	0		1	1	1
0	2	1	1	1	1	0		w0	÷.,	,1
0	2	2	г	0	2	0	/	1	1	T
0	0	0	1	1	a	0		-0	0	4
0	0	0	15	0	0	-0	- 1999 - 1999	0	-1	0
x1.	1	.11	1	/	-	1.00	/	w0	20	72
0	0	0	0	0	Q.	-0	1	10	-1	V
0	ī	1	12	2	1	0	1	0	X	1
0	+	2	0	1	0	x		10	1	-1
0	1	0	0	0	h	0		Bia	e har	60
0	1	1	2/	6	0	0/	/	boy	.,	.,0
0	2	2/	1	0	4	0		1	1	
0	9	0	0	8/	6	0	/			
×C.	1, 1,	21	1	200		1				
0	0	2	16	0	0/	0				
0	2/	0	0	2	1	0				
9	1	2	0	2	0	0				
0	1	2	1	1	2	0				
0	0	0	2	2_{i}	1	0				
0	1	0	2.	1	0	0				
0	0	0	0	0	0	0				

Filter W1 (3x3x3) Output Volume (3x3x2) x3x3) w1[:,:,0] 0[1,1,0] 1 1 0 3 4 5 -1 -1 1 10 12 10 0 1.25 W1[:,:,1] 0[:,:,1] 6 0 0 -1 3 3 w1[1,1,2] -1 0 1 0 1 x1) 1 Bias bI (1x1x1) b1[:,:,0] 0 toggle movement

Convolutional Neural Networks

State of the art: Object Classification (VGG16)

Others: Inception, ResNet

Tensorflow

TensorFlow is an open source software library for numerical computation using data flow graphs. Nodes in the graph represent mathematical operations, while the graph edges represent the multidimensional data arrays (tensors) communicated between them.

Computational Graph

with tf.Session() as session:

dot operation = tf.matmul(r2, r1)

start time = time.time()

result = session.run(dot_operation)

r1 = tf.random uniform(shape=shape, minval=0, maxval=1, dtype=data_type) r2 = tf.random_uniform(shape=shape, minval=0, maxval=1, dtype=data_type)

The Case Study

Naive Implementation

while(True):

image = CaptureFrame()

(boxes, scores, classes, num) = sess.run(

[detection_boxes, detection_scores, detection_classes, num_detections],

feed_dict={image_tensor: image_np_expanded})

DrawBBoxes(image, boxes, scores, ...)
ShowFrame(image)

Naive Implementation: Computational Graph

Implementation: Computational Graph

Conclusion

- same:
 - Measure/Profile
 - Identify and Understand the Bottleneck
 - Start the iterative tuning process
 - Validade
 - Be scientific

• Performance Engineering can look like it's changing... But the basics are always the

