



# On-board FPGA- based Deep Neural Networks processing unit

Krzysztof Czyz

[fpspace.com](http://fpspace.com)

2 On-board FPGA based Deep Neural Networks processing unit

# FP SPACE - WHO WE ARE AND WHAT WE DO...

 **Future Processing**

Earth Observation  
Application Development

 **FP Space**

 **KP LABS**

Flight Software, Big Data  
and Machine Learning

 **FPINSTRUMENTS**

Electronics Design  
and Production



## INTUITION 1

### HYPERCAM INSTRUMENT DEMONSTRATION MISSION

Hyperspectral instrument for Earth observation with increased spectral resolution enabling in orbit autonomous operation and automatic processing and classification of satellite data based on new algorithms for segmentation and classification of satellite images using deep convolutional networks.





## INTUITION 1

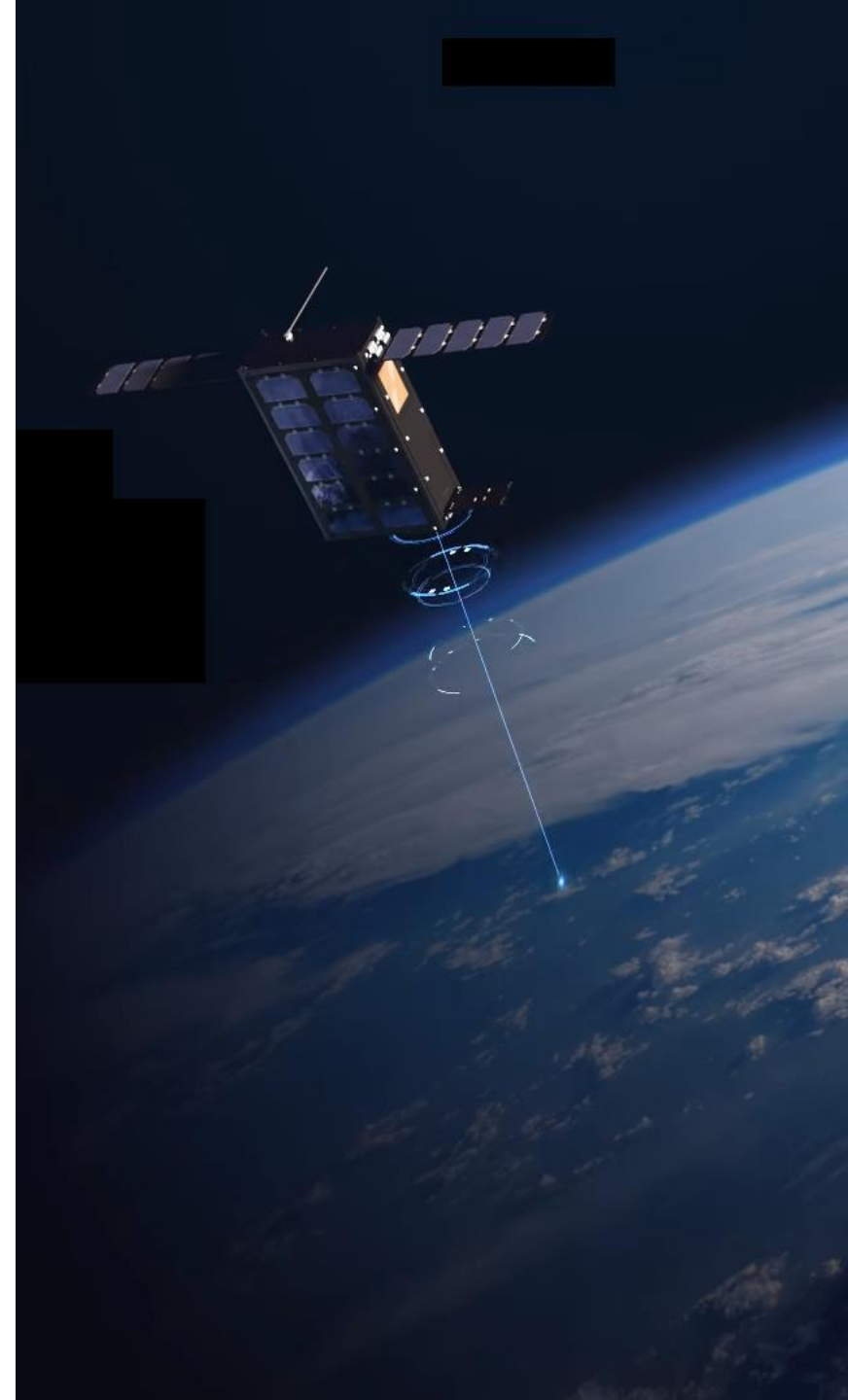
### HYPERCAM INSTRUMENT DEMONSTRATION MISSION

A high resolution hyperspectral instrument for EO operating at 600 km SSO:

- CMOS sensor with embedded Fabry-Perot spectral filters
- ground sampling distance: 25m
- spectral range: 450 nm to 900 nm
- no. of spectral bands: >75
- swath: 15 km
- RAW data per second: 0.7 GiB

A high performance processing unit for Deep Learning Acceleration:

- Scalable architecture providing overall performance >1 TOPS
- On-board data segmentation and classification
- HSI data compression
- Dynamic schedule optimization



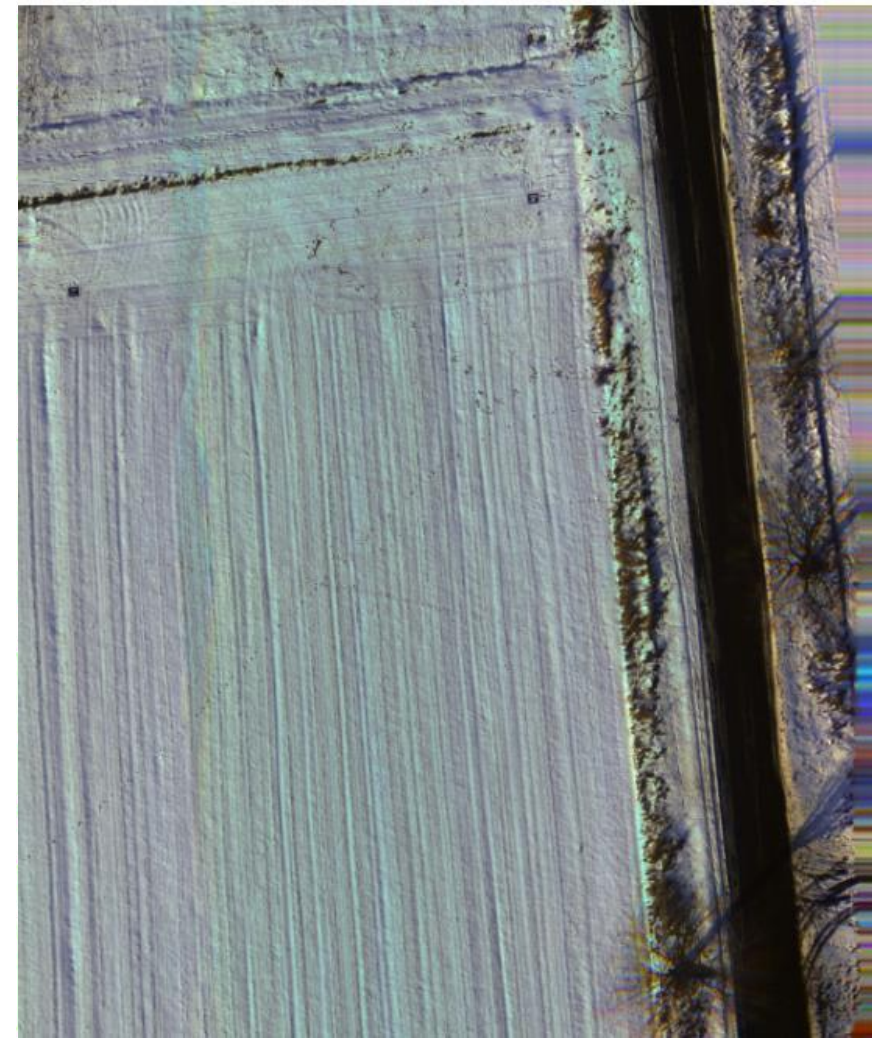
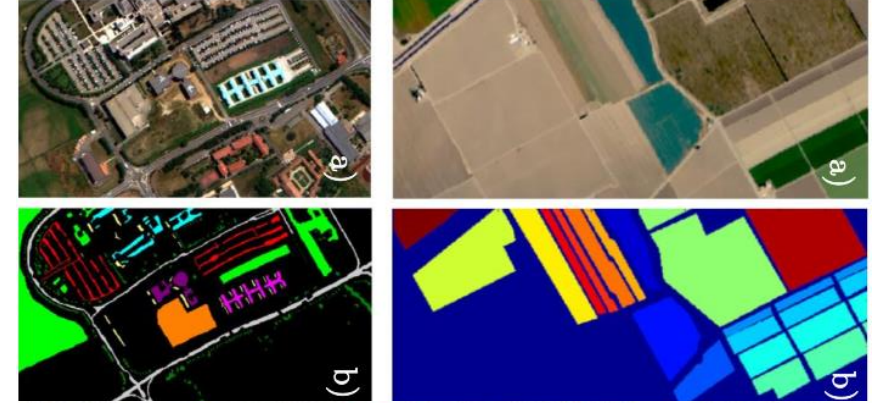
## INTUITION 1

### HYPERCAM INSTRUMENT DEMO STRATION MISSION

A algorithms for on-board HSI segmentation and classification:

- Segmentation of multiband and HSI
- Reduction of HSI by effectively select most-important bands
- Handling difficult HSI (data imbalance, representativeness)
- Automated design of deep neural nets for HSI segmentation
- HSI ground-truth data using UAV

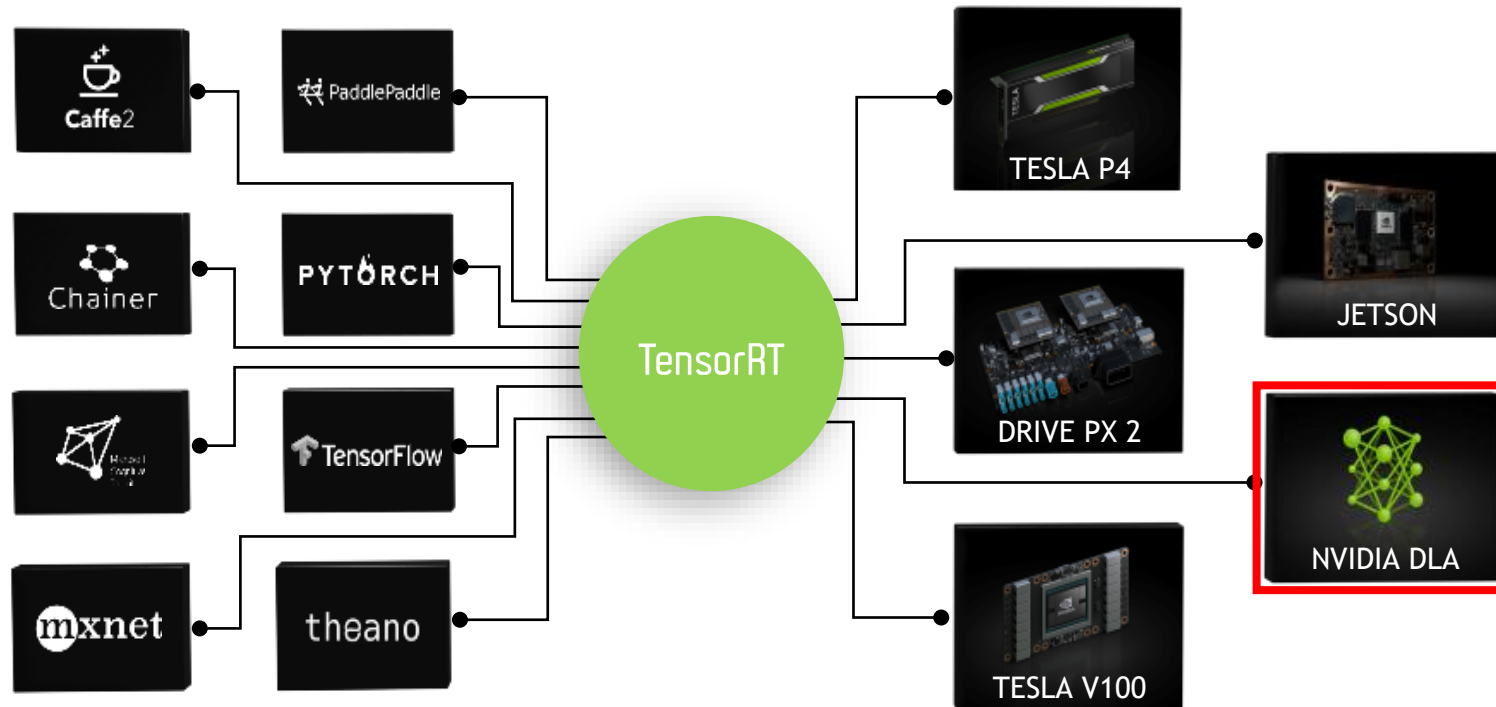
But, there is a price to be paid, these methods are computationally intensive and require supercomputing resources - that can be challenging, especially on-board of a space craft.





# FRAMEWORKS AND HARDWARE

ANYONE SEEN A SUPERCOMPUTER?

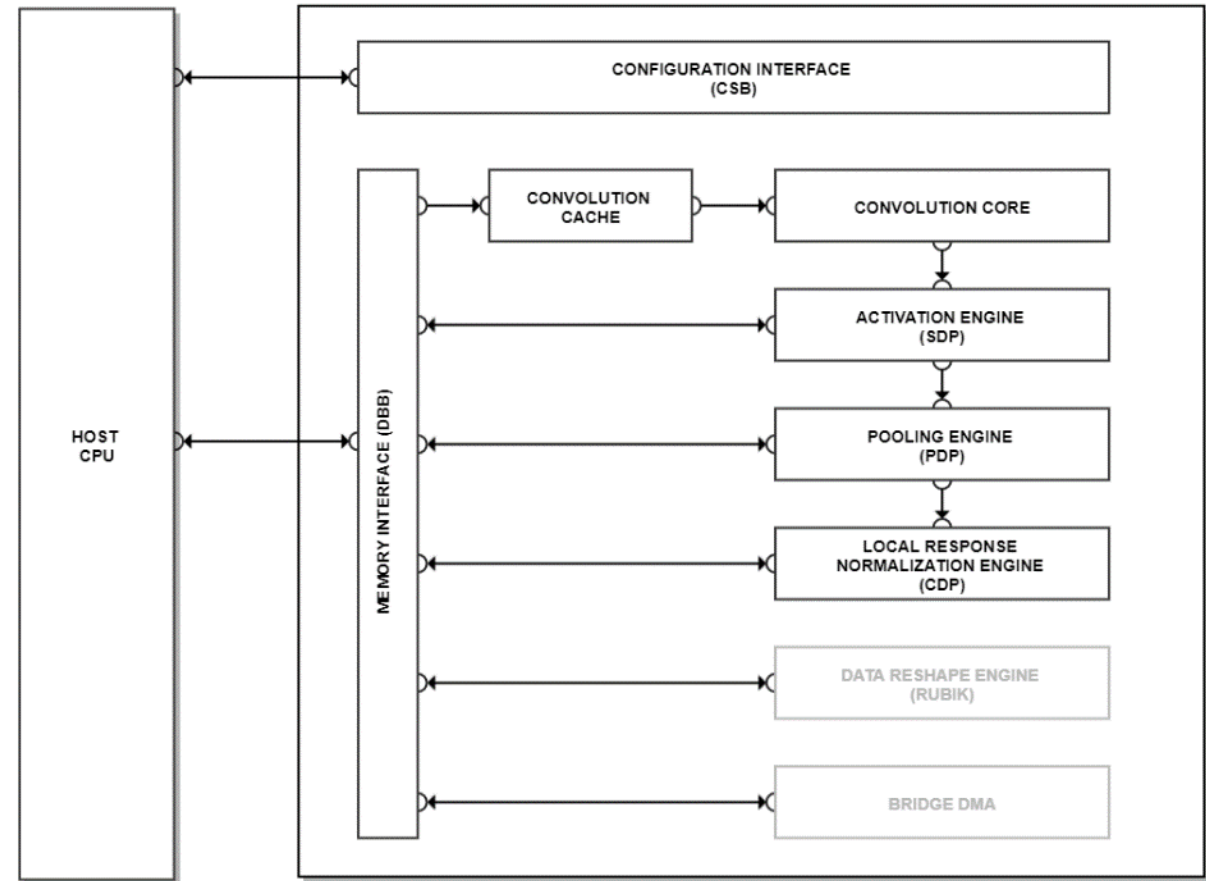


# NVIDIA DEEP LEARNING ACCELERATOR

WHAT IS IT? A HARDWARE PART

An open source neural network highly parallel deep learning accelerator created by NVIDIA providing:

- **Matrix Convolution Cores** (64-2048 cores)  
 Direct/Winograd/Multi-Batch Convolution FC layers
- **Single Data Processors** (1-16 cores)  
 Activation functions (from ReLU to non-linear)  
 Bias addition  
 Batch normalization  
 Element-wise operations
- **Planar Data Processors** (1-16 cores)  
 Pooling (min, max, average)
- **Channel Data Processors** (1-16 cores)  
 Local normalization functions

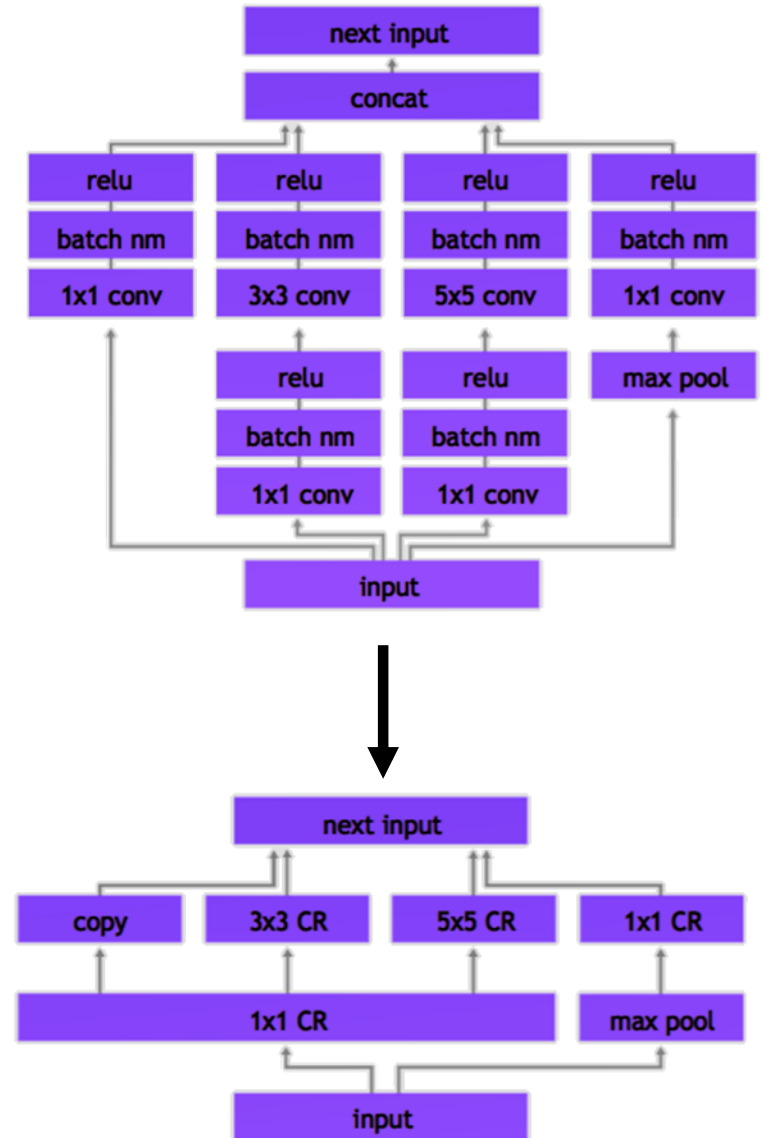


# NVIDIA DEEP LEARNING ACCELERATOR

WHAT IS IT? A HARDWARE PART

An open source neural network highly parallel deep learning accelerator created by NVIDIA providing:

- Matrix Convolution Cores** (64-2048 cores)  
 Direct/Winograd/Multi-Batch Convolution FC layers
- Single Data Processors** (1-16 cores)  
 Activation functions (from ReLU to non-linear)  
 Bias addition  
 Batch normalization  
 Element-wise operations
- Planar Data Processors** (1-16 cores)  
 Pooling (min, max, average)
- Channel Data Processors** (1-16 cores)  
 Local normalization functions



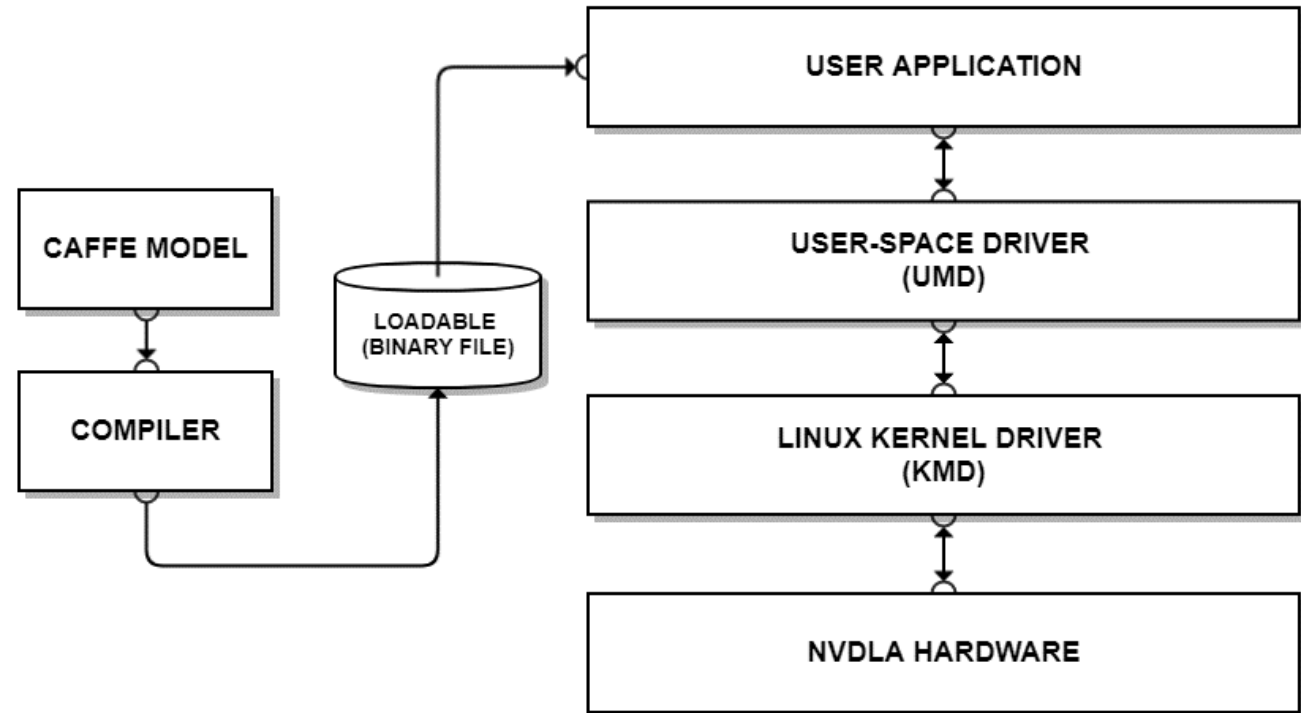


# NVIDIA DEEP LEARNING ACCELERATOR

WHAT IS IT? A SOFTWARE PART

nvDLA software ecosystem:

- **Caffe model**  
Use well established framework and deep learning data representations
- **Compiler**  
Translates DNN layers to list of low-level operations.
- **User Application**  
Run inference
- **UMD/KMD drivers**  
Hardware Abstraction Layer to the deep learning accelerator



10 On-board FPGA based Deep Neural Networks processing unit

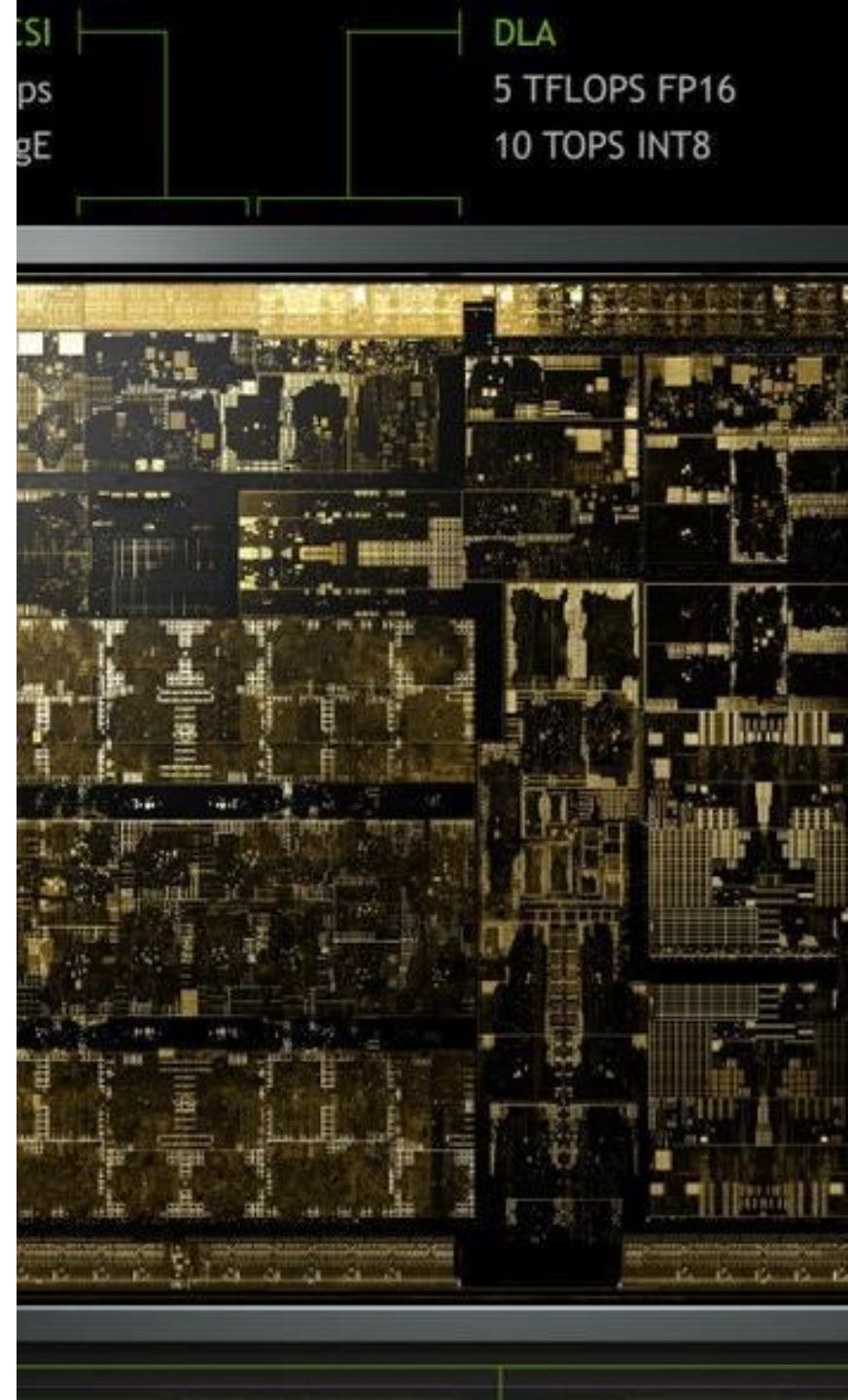
## NVIDIA DEEP LEARNING ACCELERATOR

WHAT IS IT? HARDWARE IMPLEMENTATION

nvDLA is meant to be implemented in hardware i.e. Xavier SOC.

NVIDIA has decided to go to open source with nvDLA and has released Verilog RTL:

- End of 2017  
fully featured, floating point, non-configurable accelerator (NV\_FULL)
- 2018 Q1/Q2  
small, fixed-point (INT8), non-configurable accelerator (NV\_SMALL)
- 2018 Q3  
added degrees of freedom in architecture configuration (NV\_CUSTOM)
- 2019 Q1  
release of DNN compiler for fixed point version of nvDLA





# NVIDIA DEEP LEARNING ACCELERATOR

## FPGA IMPLEMENTATION OF NVDLA

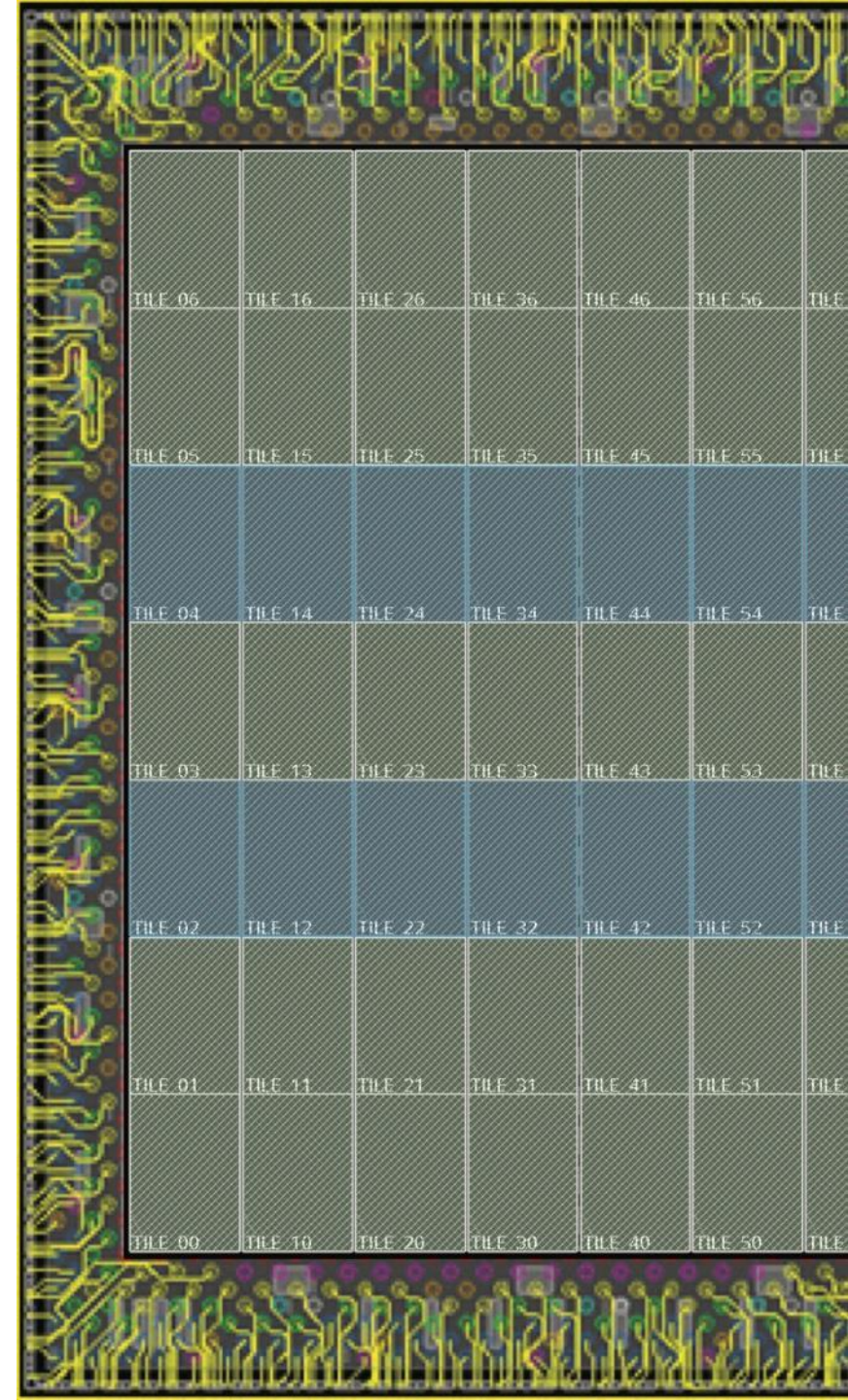
Goal: fit triple redundant nvDLA engine into FPGA

### PROS:

- very good customization
- reasonable computing power of single nvDLA instance  
AlexNet 150 FPS (~33% of TX2i), GooqLeNet 75 FPS (38% of TX2i)  
for 512 MAC cores running @ 150 MHz
- FPGA mature technology in space

### CONS:

- DNN compiler for fixed point nvDLA to be published in Q1 2019,
- the only reasonable data representation is fixed-point INT8
- higher power consumption than Jetson TX2i
- nvDLA is a complex IP, consuming a lot of FPGA area





## NVIDIA DEEP LEARNING ACCELERATOR

FPGA IMPLEMENTATION OF NVDLA

### Xilinx ZU15EG FPGA vs KU040

#### ZU15EG

- PMU with TMR, 2xR5 Cores with lock-step, 4xA53 Cores
- Good SEU performance, but exposed to SEL
- Good power efficiency

#### KU040

- No hardcoded cores (LEON seems to be a good choice)
- Higher sensitivity to SEU, but no SEL issues



# NVIDIA DEEP LEARNING ACCELERATOR

## FPGA OPTIMISATION ISSUES

### Goal:

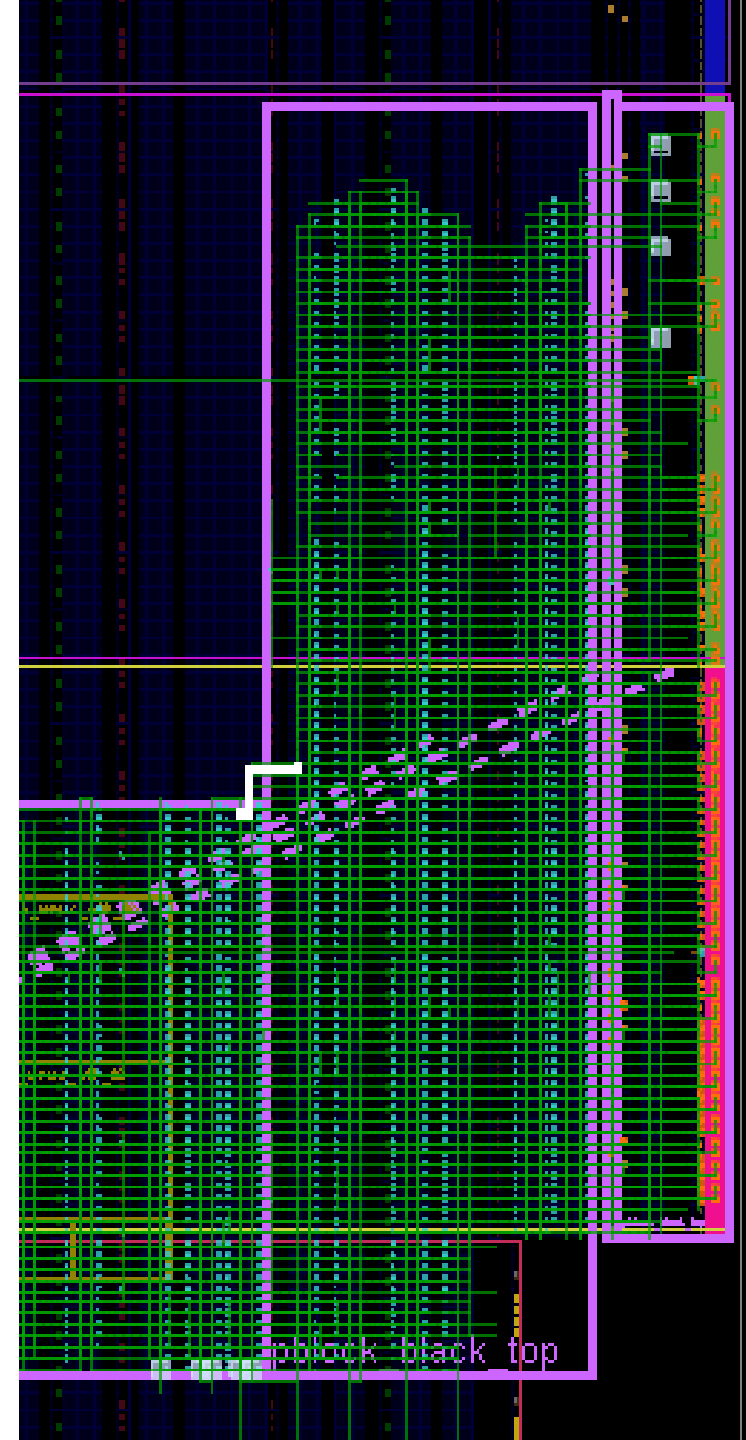
fit 3x nv\_11 in TMR mode on ZU15EG device

OR

fit 2x nv\_06 in DMR mode on ZU15EG device  
clock frequency at least 150 MHz

### Possible solution:

- very high utilization of generic LUT is observed in the design, but almost no utilization of special function LUT like Distributed RAM (few hundreds) or Shift Register (literally three).
- a lot of FIFO queues are implemented as register chain – the reimplementaion to use more efficient FPGA structures (BRAM, DRAM or SRL) can gain 3% up to 7% (depending on architecture model) extra space!



# NVIDIA DEEP LEARNING ACCELERATOR

## FPGA RESOURCE UTILIZATION

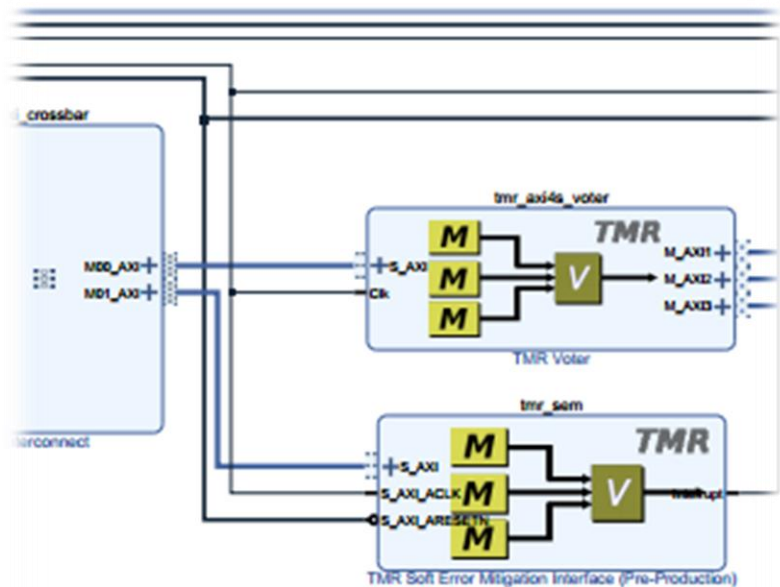
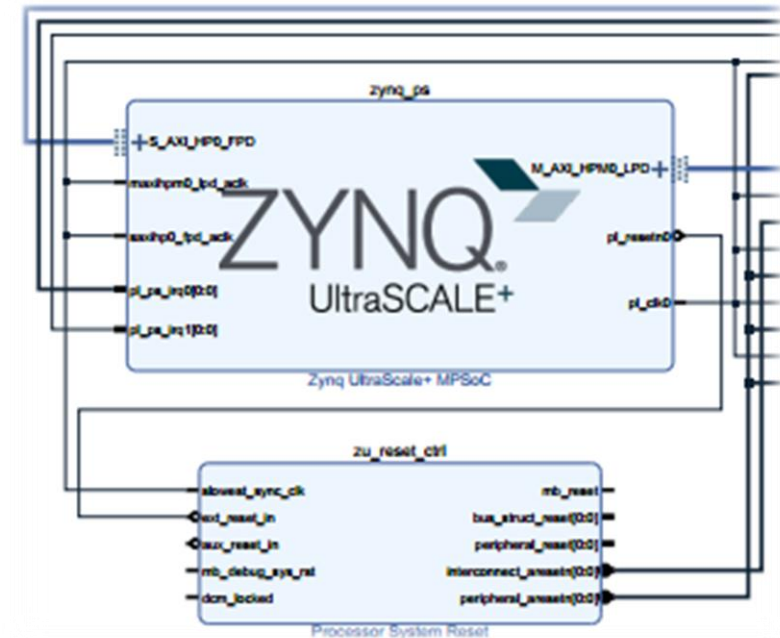
MODEL	Weight Compression	Winograd Convolution	Batch processing	SDP Non-linear	SDP Bias addition	SDP Batch norm.	SDP Element-wise	PDP Engine	CDP Engine	# of MAC	Conv. cache (kB)	SDP BS Throughput	SDP BN Throughput	SDP EW Throughput	PDP Throughput	CDP Throughput	Memory latency	Memory max. burst	Total LUTs	BRAM	DSP48	LUT usage (% ZU15EG)	FF usage (% ZU15EG)	LOCKSTEP possible	TMR possible	LUT usage (% ZU17EG)	FF usage (% ZU17EG)	LOCKSTEP possible	TMR possible
nv_01	■	■	■	■	■	■	■	■	■	512	512	16	16	4	8	8	1024	4	283495	209	239	<b>83%</b>	<b>32%</b>			<b>67%</b>	<b>26%</b>		
nv_02	■	■	■	■	■	■	■	■	■	512	512	4	4	1	2	2	1024	4	179960	197	83	<b>53%</b>	<b>25%</b>			<b>43%</b>	<b>20%</b>	✓	
nv_03			■	■	■	■	■	■	■	512	512	4	4	1	2	2	1024	4	171645	197	79	<b>50%</b>	<b>25%</b>			<b>41%</b>	<b>20%</b>	✓	
nv_04			■		■	■		■	■	512	512	4	4	-	2	2	1024	4	156351	194	65	<b>46%</b>	<b>22%</b>			<b>37%</b>	<b>18%</b>	✓	
nv_05			■		■	■		■	■	512	128	4	4	-	2	2	1024	4	155878	194	65	<b>46%</b>	<b>22%</b>			<b>37%</b>	<b>18%</b>	✓	
nv_06			■		■	■		■	■	512	128	1	1	-	1	1	1024	4	143205	194	40	<b>42%</b>	<b>21%</b>	✓		<b>34%</b>	<b>17%</b>	✓	
nv_07			■		■	■		■		512	128	1	1	-	1	1	1024	4	124613	192	35	<b>37%</b>	<b>18%</b>	✓		<b>29%</b>	<b>14%</b>	✓	✓
nv_08			■		■	■		■		512	128	1	1	-	1	1	64	4	124620	187	35	<b>37%</b>	<b>18%</b>	✓		<b>29%</b>	<b>14%</b>	✓	✓
nv_09			■		■	■		■		512	128	1	1	-	1	1	64	1	124087	187	35	<b>36%</b>	<b>18%</b>	✓		<b>29%</b>	<b>14%</b>	✓	✓
nv_10			■		■	■		■		512	32	1	1	-	1	1	64	1	118654	123	35	<b>35%</b>	<b>16%</b>	✓		<b>28%</b>	<b>13%</b>	✓	✓
nv_11			■		■	■				512	32	1	1	-	1	1	64	1	101371	119	32	<b>30%</b>	<b>12%</b>	✓	✓	<b>24%</b>	<b>10%</b>	✓	✓
nv_12			■			■				512	32	1	1	-	1	1	64	1	95772	117	30	<b>28%</b>	<b>11%</b>	✓	✓	<b>23%</b>	<b>9%</b>	✓	✓
nv_13						■				512	32	1	1	-	1	1	64	1	95541	117	30	<b>28%</b>	<b>11%</b>	✓	✓	<b>23%</b>	<b>9%</b>	✓	✓



# NVIDIA DEEP LEARNING ACCELERATOR

## TRIPLE MODULAR REDUNDANCY

- NVDLA is a complex IP, consuming a lot of physical area and processing a lot of data for extended period.
- This increases the likelihood of Single Event Upset (SEU) which can lead to corrupted output data or functional failure.



# NVIDIA DEEP LEARNING ACCELERATOR

TRIPLE MODULAR REDUNDANCY  
MULTIPLE MEMORY INTERFACES AND ADVANCED TMR VOTER

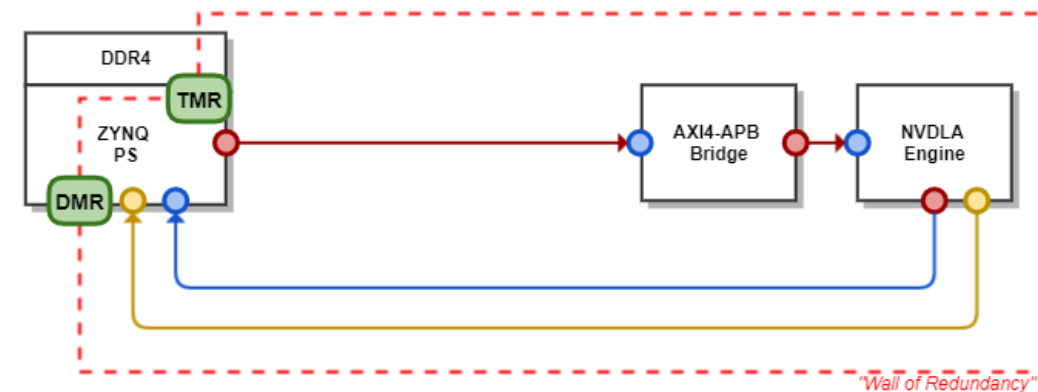
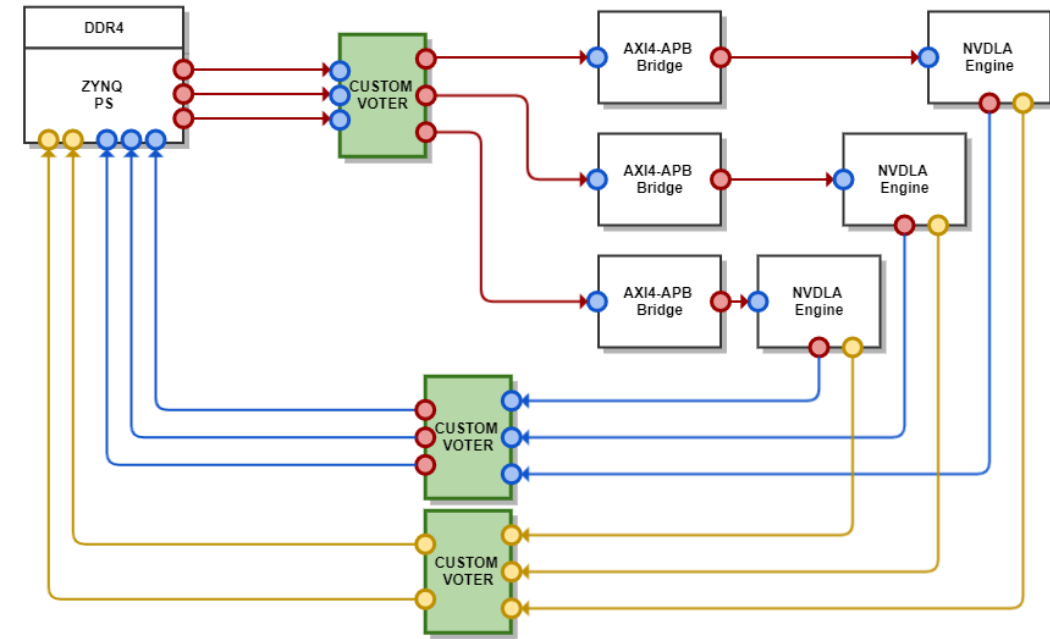
## Block TMR

### Pros:

- Error detection mechanism (well, FPGA part at least)
- Input/output data copys reside in DDR4

### Cons:

- Does not correct nvDLA operation (only partial results)
- Complicated TMR design – AXI4 may have distinct behaviour on each channel due to accessing the same DDR4 controller
- DDR4 bandwidth and latencies might be an issue (concurrent)
- DDR4 increased usage (3x)



# NVIDIA DEEP LEARNING ACCELERATOR

TRIPLE MODULAR REDUNDANCY  
MULTIPLE MEMORY INTERFACES AND ADVANCED TMR VOTER

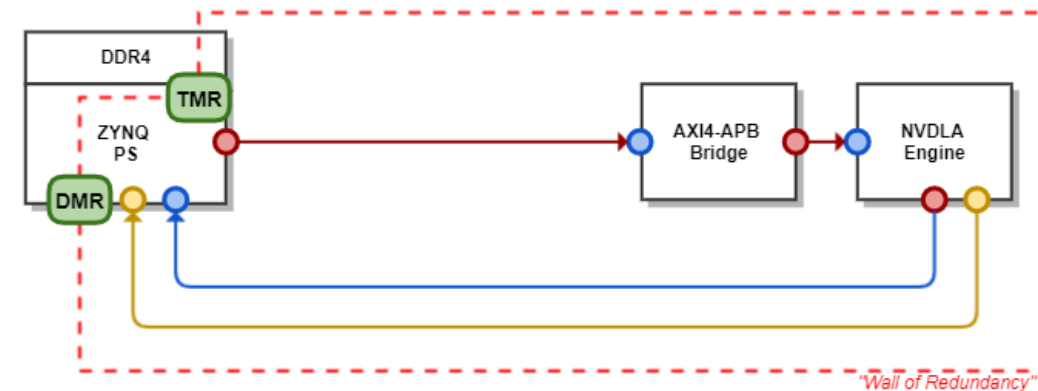
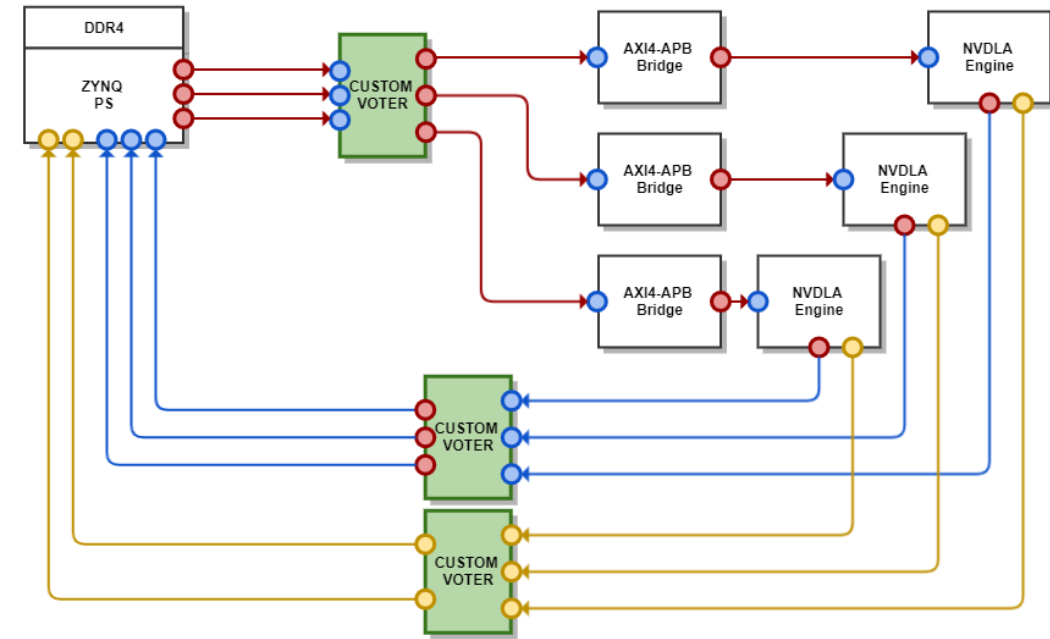
Distributed TMR (to be implemented?)

Pros:

- Implemented for individual Processors/Cores of nvDLA
- Improved resistance to SEE
- Masks configuration bit errors

Cons:

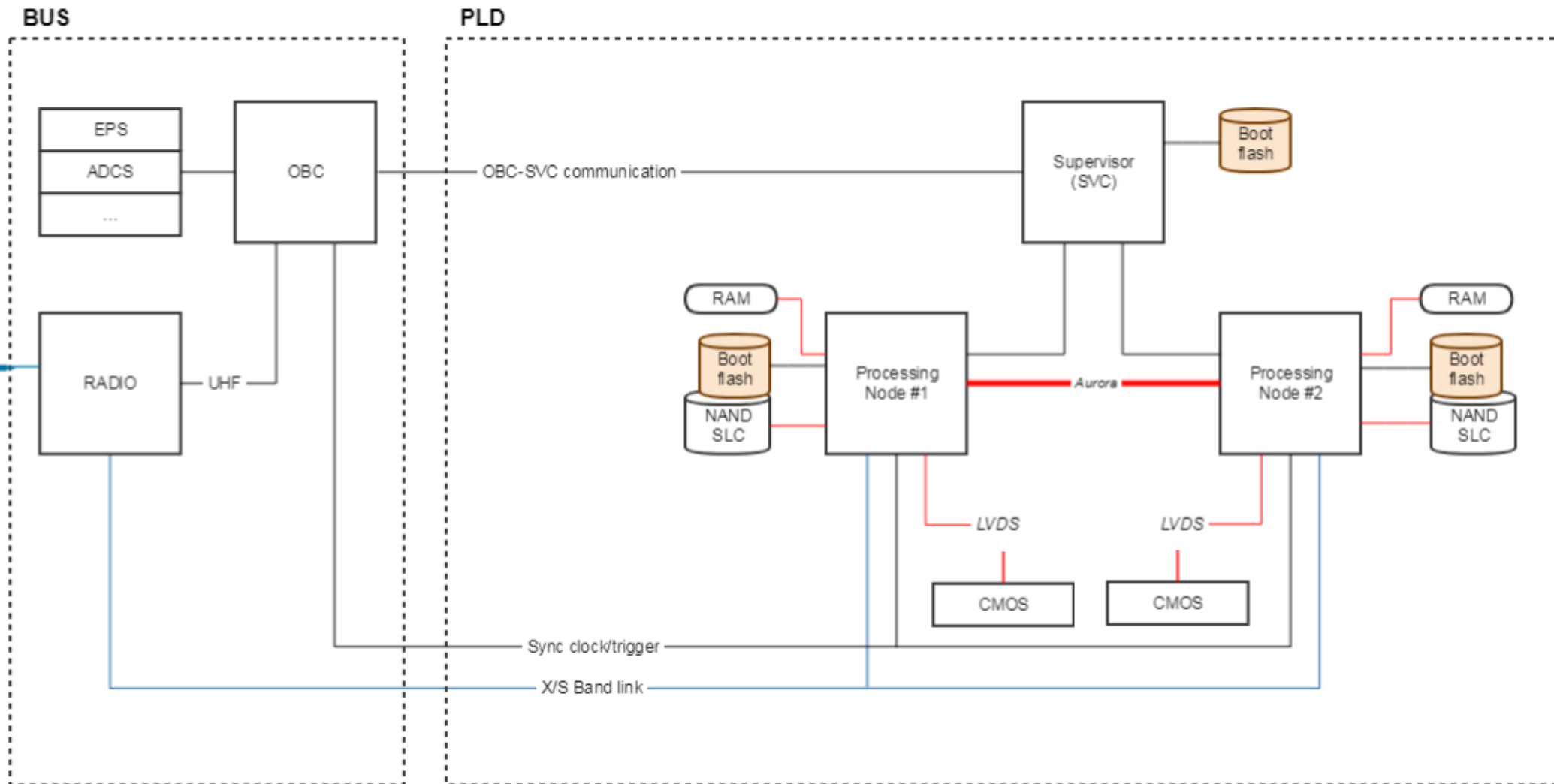
- Complicated design
- High resource utilisation





# HYPERCAM PAYLOAD

## BLOCK DIAGRAM

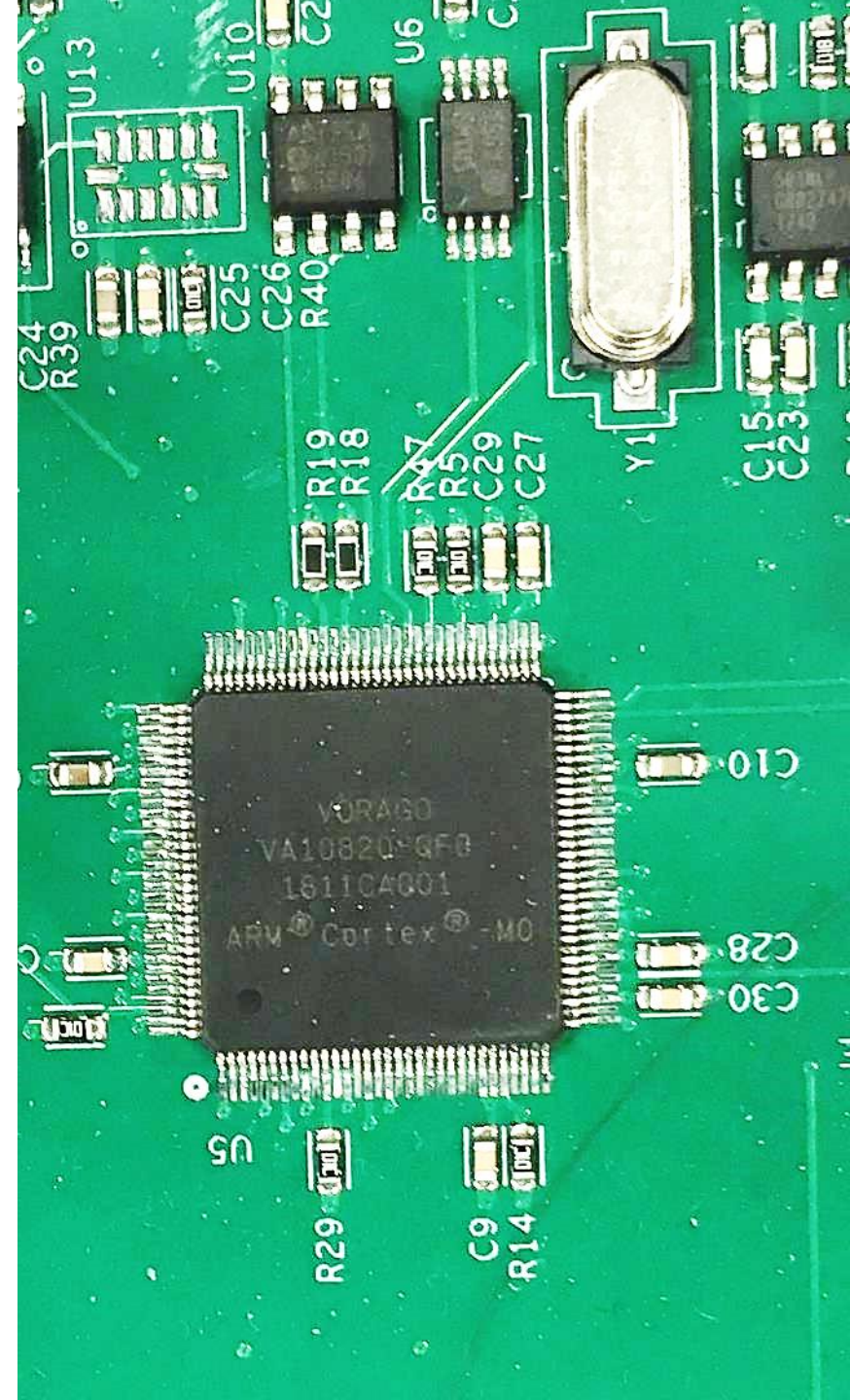


## SUPERVISORY CIRCUIT

VORAGO RADHARD MCU

Supervisory circuit / payload controller functions:

- Communication bridge between PLD and OBC
- Task scheduling for computing unit
- Power, thermal and fault management
- Error logging
- Firmware management
- Safe mode



20 On-board FPGA based Deep Neural Networks processing unit

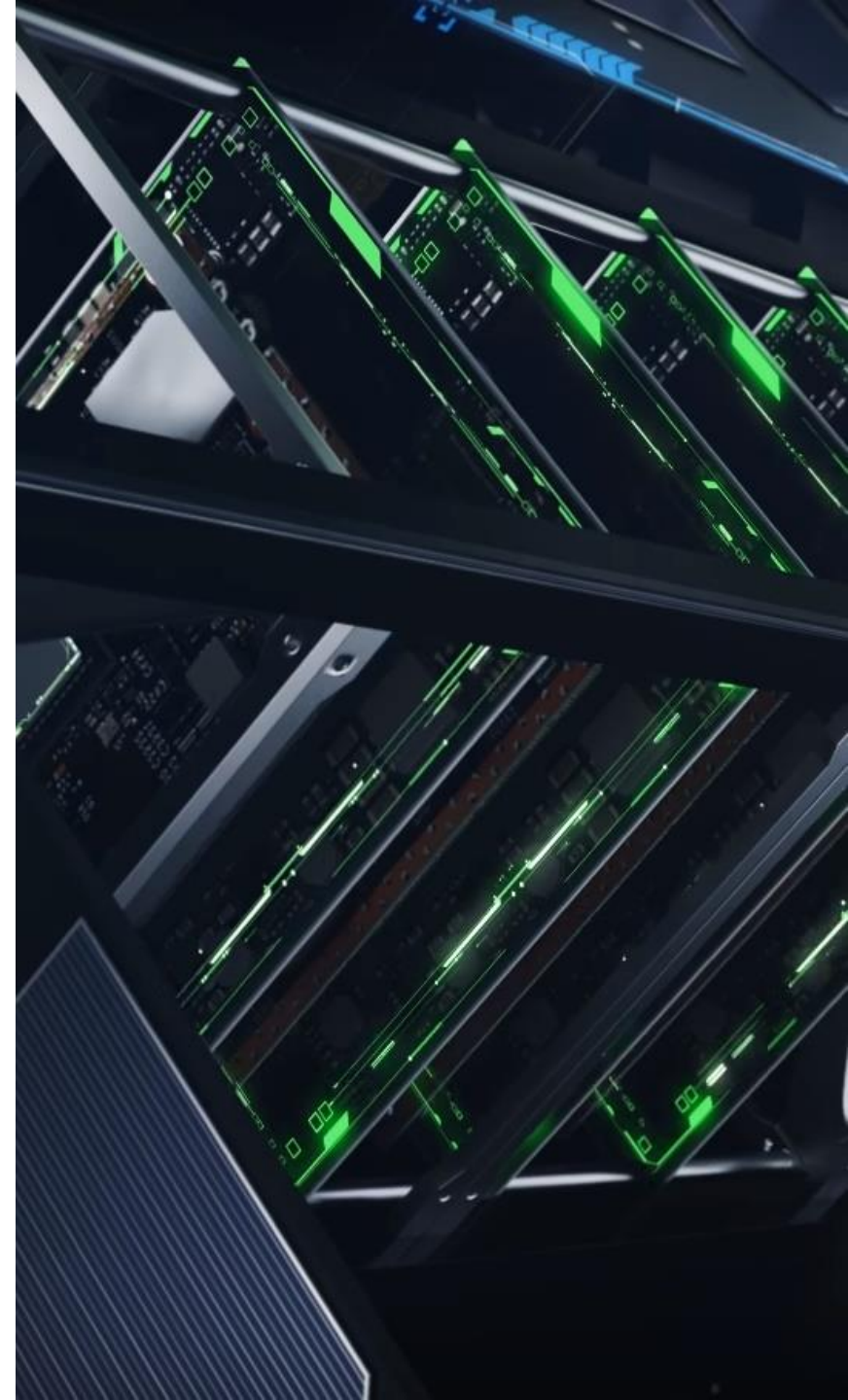
## COGNITION

ON-BOARD PROCESSING DEMONSTRATION...

A rideshare mission, initial on orbit tests of DNN processing unit.

- Architecture providing performance >200 GOPS
- 8 GiB of RAM memory
- Panchromatic camera
- Power consumption < 20 W
- Size U1

... to be launched in early 2020





## INTUITION 1

### HYPERCAM INSTRUMENT DEMONSTRATION MISSION

A high performance processing unit dedicated for segmentation of the hyperspectral images in orbit:

- Scalable architecture providing overall performance >1 TOPS
- 16 GiB of RAM memory
- Power consumption < 30 W
- Flash based data storage 0.5 TiB
- Triple Modular Redundancy

...supercomputing powers for a small sat



## INTUITION 1

### HYPERCAM INSTRUMENT DEMONSTRATION MISSION

Hyperspectral system for observing the Earth with increased spectral resolution enabling automatic processing and selection of satellite data in orbit based on new algorithms for segmentation and classification of satellite images using deep convolutional networks.

... to be launched in 2023







**THANK YOU**

---

Krzysztof Czyż,  
FP Instruments sp. z o.o.  
44-100 Gliwice,  
[kczyz@fp-instruments.com](mailto:kczyz@fp-instruments.com)

[fpspace.com](http://fpspace.com)