



Vision-Based Gesture-Controlled Drone: A Comprehensive Python Package to Deploy Embedded Pose Recognition Systems

Summer 2021 Research Project – Final Report

August 11, 2021

Arthur FINDELAIR

Advisor: Dr. Jafar Saniie

Embedded Computing and Signal Processing Research Laboratory

Department of Electrical and Computer Engineering

Illinois Institute of Technology, Chicago, Illinois, USA

Objective of the project

Develop a **vision-based gesture control** system on an autonomous drone using industry standards to facilitate **portability and extensibility**.



Other applications of a gesture control system

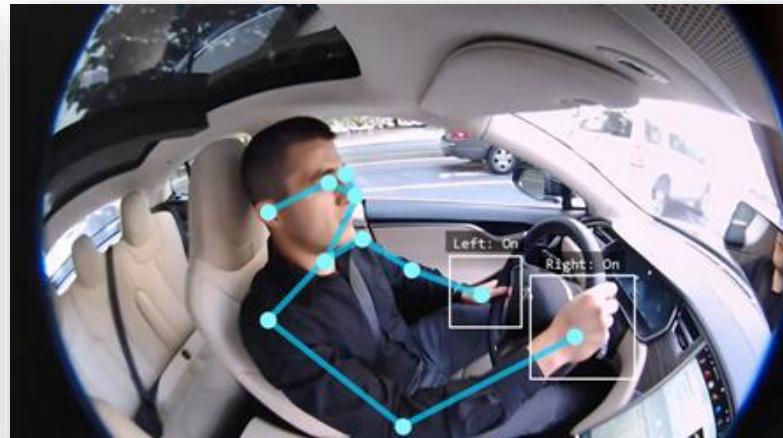
Continuous tracking

Augmented reality
VFX animations
Robotic



Gesture recognition

Seniors & patients monitoring
Body language interpretation
Basic HMI



Hybrid

Advanced HMI



Table of content

- Motivation & Overview of the project
- Hardware components
- Pose estimation
- Pose classification
- Deployment
- Demonstration
- Future works



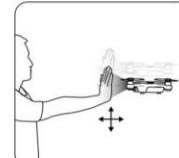
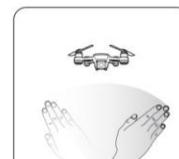
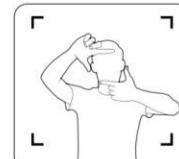
Motivation of the project – Current uses of drones

APPLIED DRONE-BASED METHODS PER INDUSTRY

Industry Vertical	Mapping & Surveying	Inspection	Photography & Filming	Localization & Detection	Other ¹	Delivery	Spraying & Seeding
Energy (Utilities)	14%	83%	0%	1%	1%	0%	0%
Construction	80%	15%	0%	0%	2%	2%	0%
Transportation and Warehousing	48%	34%	0%	6%	5%	7%	0%
Agriculture	59%	3%	0%	8%	0%	0%	30%
Mining, Quarrying, and Oil & Gas Extraction	55%	38%	0%	0%	6%	0%	2%
Public Administration	28%	5%	15%	32%	14%	0%	6%
Information	3%	6%	74%	0%	16%	0%	0%
Real Estate, Rental & Leasing, and Industrial Plants	16%	67%	17%	0%	0%	0%	0%
Arts, Entertainment, and Recreation	5%	0%	60%	0%	30%	5%	0%
Insurance	46%	10%	22%	14%	8%	0%	0%
Health Care and Social Assistance	35%	0%	2%	17%	8%	38%	0%
Professional, Scientific, and Technical Services	51%	6%	2%	15%	24%	2%	0%
Safety & Security	11%	0%	0%	69%	20%	0%	0%
Educational Services	53%	18%	0%	12%	17%	0%	0%
Waste Management & Remediation Services	34%	14%	5%	17%	30%	0%	0%

Motivation of the project

– Commercial solutions

Adjusting Position		
1. Move your palm up or down slowly to control the aircraft's altitude while maintaining a constant distance between your palm and the aircraft.	 Solid green	
2. Keep your palm at a constant distance from your body (as shown), then move your arm left or right to control the aircraft's orientation.		
3. Keep your palm at a constant distance from your body, move forward or backward to fly forward or backward.		
Follow		
1. Stand in front of the aircraft then raise one of your arms and quickly wave your hand at the camera. Position your palm about 0.7 m in front of the aircraft's nose for about two seconds.	 Blinks green twice	
2. The front LEDs will blink green twice if the gesture is recognized successfully. The aircraft will ascend and fly backward, then hover in place 3 m from where you're standing, 2.3 m above ground. Do not move your body until the aircraft hovers in place.		
3. The front LEDs will glow solid green and the aircraft will start tracking automatically.		
4. If the wave gesture fails or the aircraft loses track of you, you can activate Follow by raising both arms above your head in a Y shape, then holding for two seconds.	 Solid green	
Taking Selfies		
1. Make a frame with your hands within 23 ft (7 m) of the aircraft while facing toward the camera and it will begin taking a selfie.	 Blinks red slowly	
2. Selfie gesture has been recognized successfully if the Front LEDs blink red slowly. Wait for the selfie count down for three seconds. The Front LEDs will blink red quickly, indicating the camera is about to take a selfie.		

Spark Gesture Controls

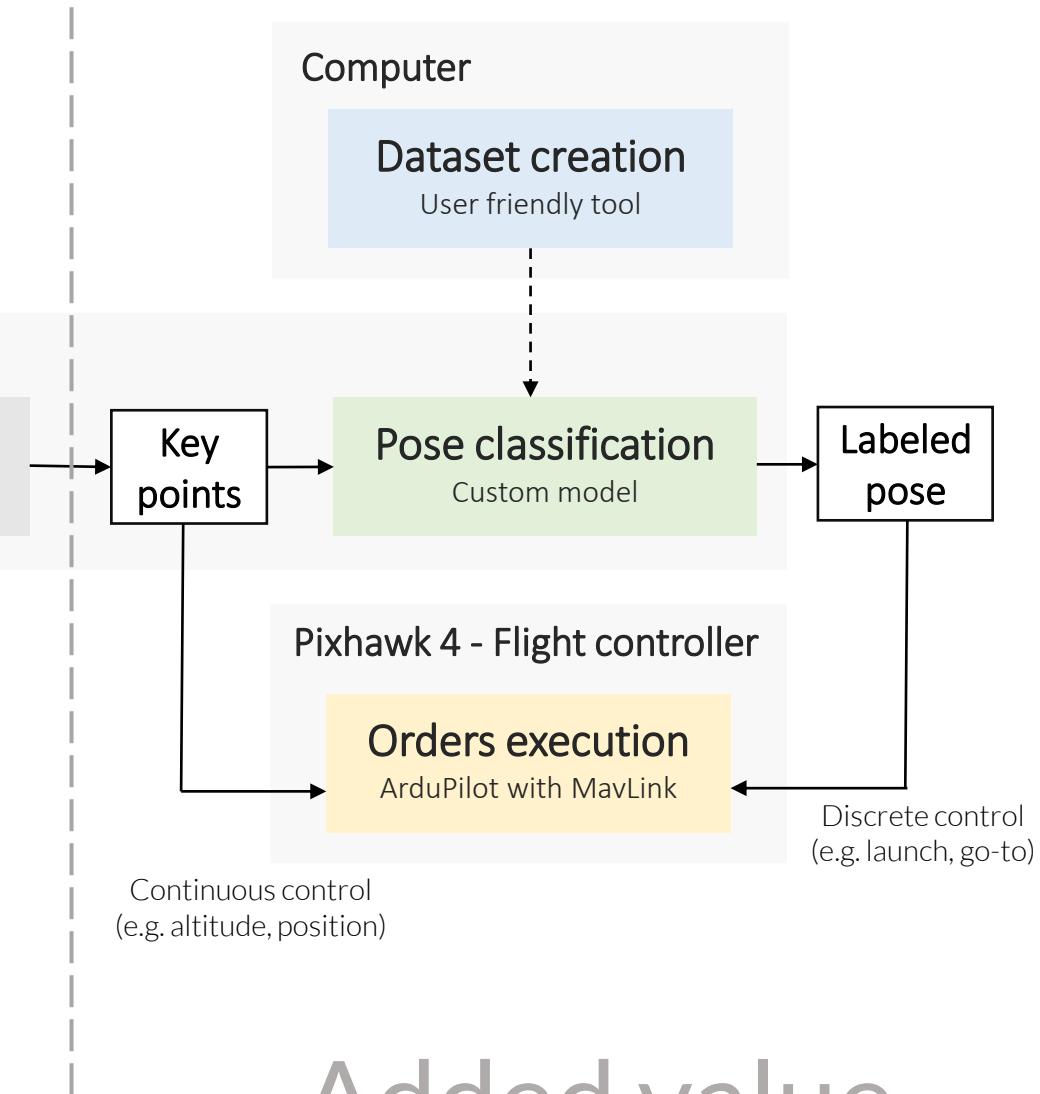
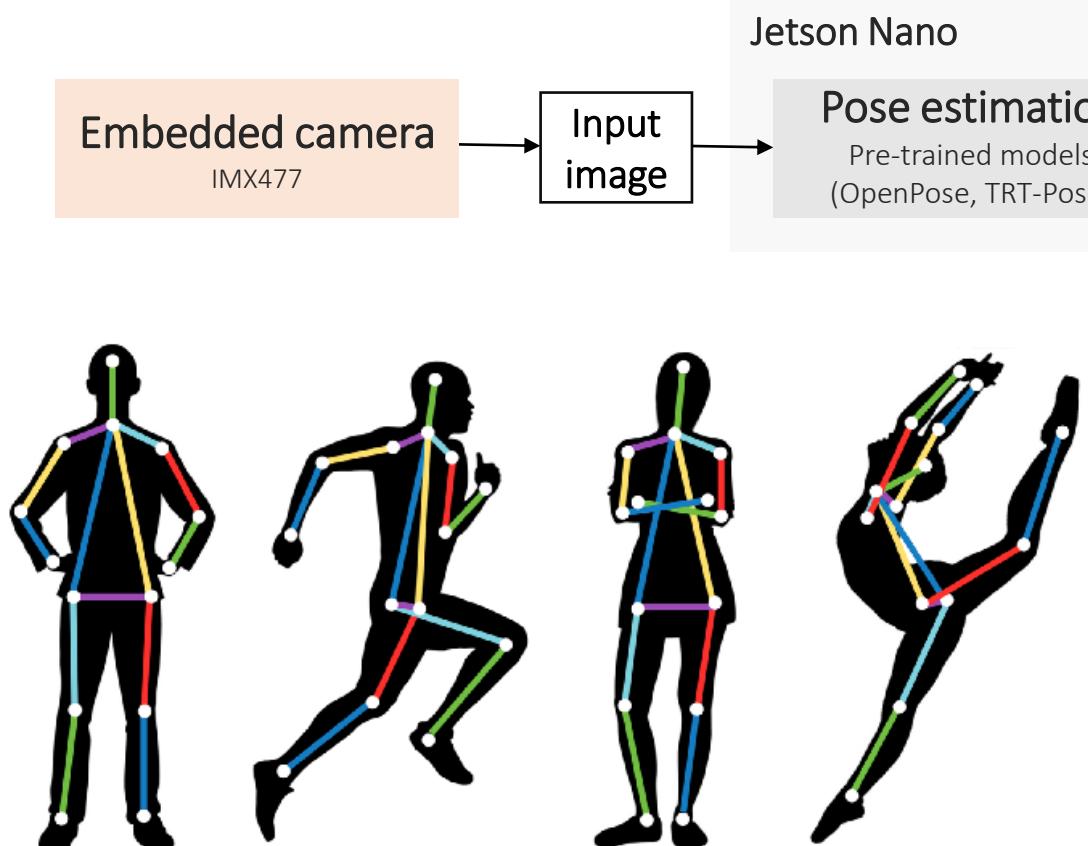


DJI Motion
Controller



Overview of the project – The pipeline

Existing tools



Added value

Overview of the project – The deployment platform



DJI Mavic Air 2



Ryze Tello



Holybro X500



Abstraction

Flexibility

Skydio X2



Parrot Anafi



Hardware – The deployment platform

Kit Contents

- Pixhawk4 Flight controller
- Power Management PM07
- Pixhawk4 GPS
- Motors – 2216 KV880
- BLHeli S ESC 20A
- Propeller 1045
- 915MHz Telemetry Radio
- Power and Radio Cables
- Battery Straps

Characteristics

- Dimensions: 410*410*300mm
- Wheelbase: 500mm
- Weight: 978g
- Max payload: 2.1kg (at 14.8V)

Additional parts

- 6000 mAh, 4 cells LiPo battery
- FrSky X8R 2.4GHz receiver
- Taranis X9D+ controller



Hardware

– The deployment platform

- Main FMU processor: STM32F765
- IO Processor: STM32F100
- Main sensors:
ICM-20689 (IMU), BMI055 (IMU), IST8310 (Mag), MS5611 (Baro)
- GNSS module sensors:
Ublox Neo-M8N (GPS/GLONASS receiver), IST8310 (Mag)

Transceiver Telemetry



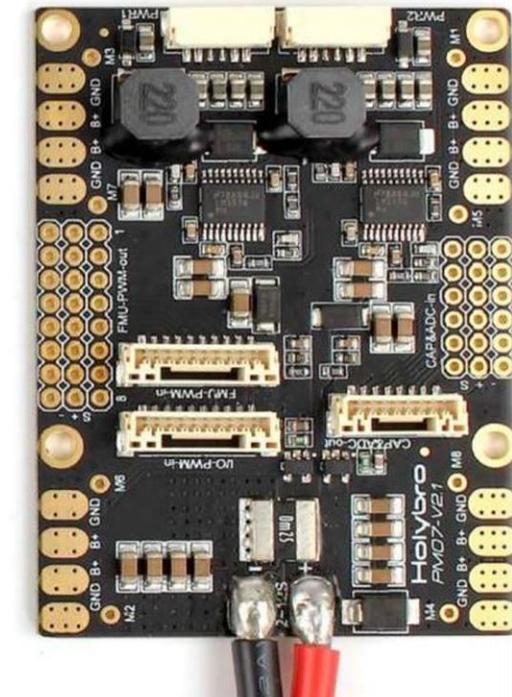
GNSS module



Pixhawk 4 (FMUv5)



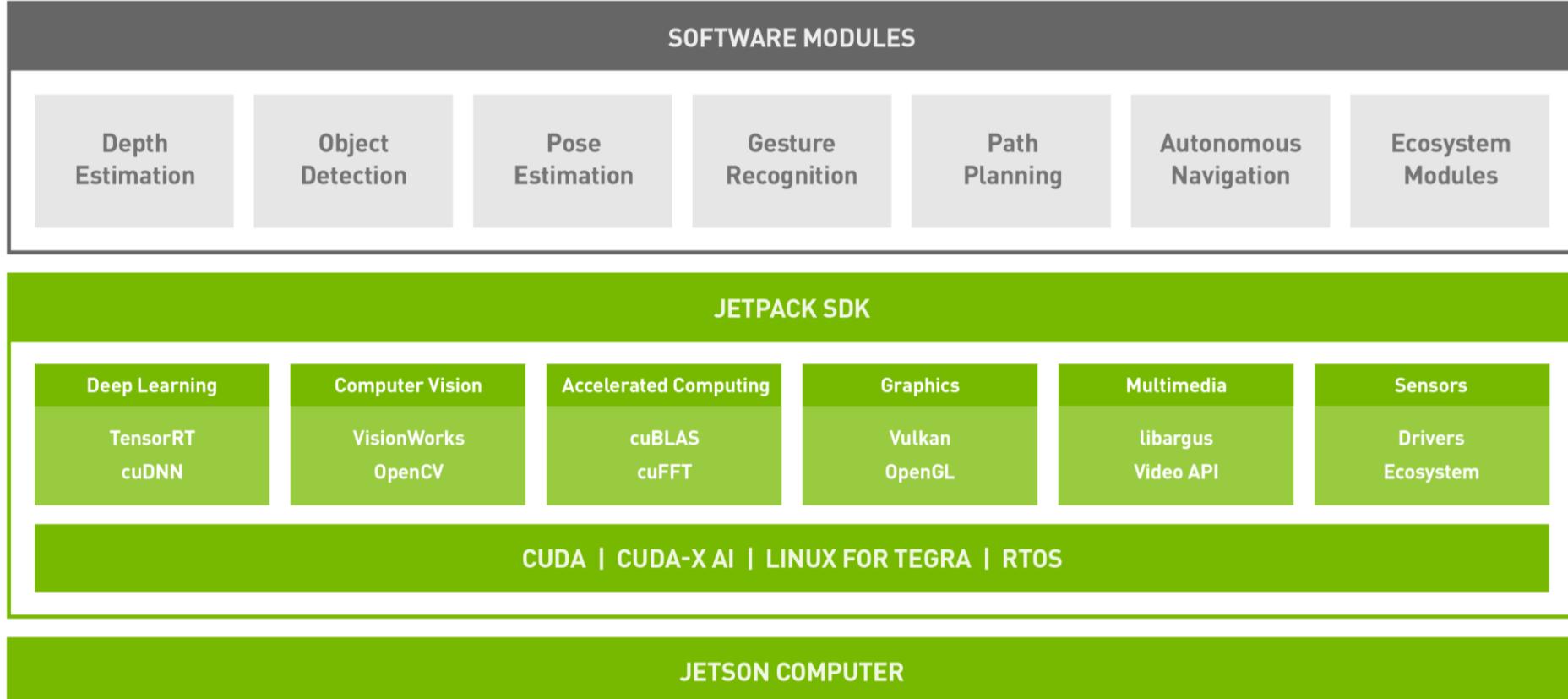
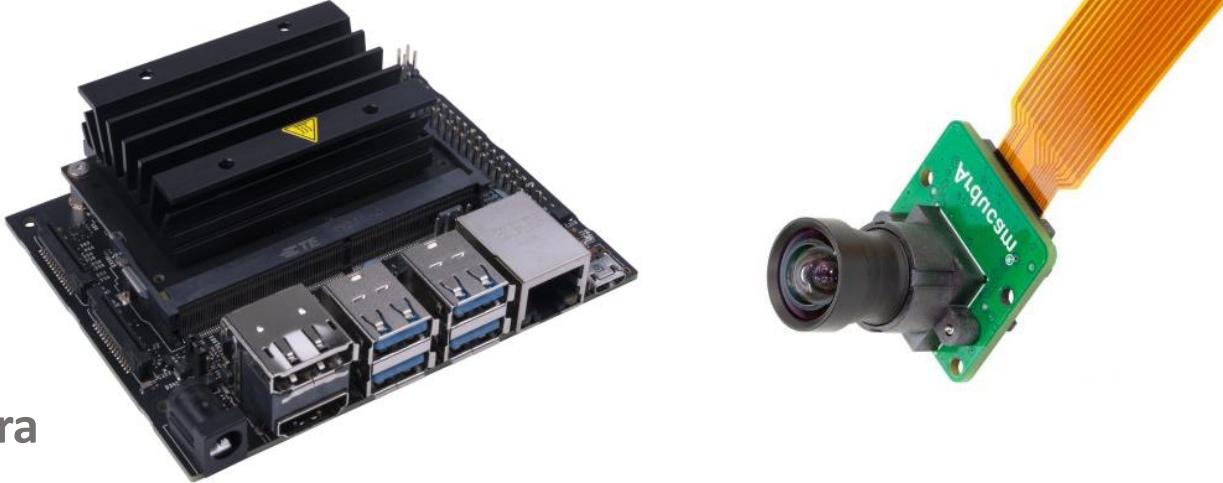
Power distribution board



Hardware

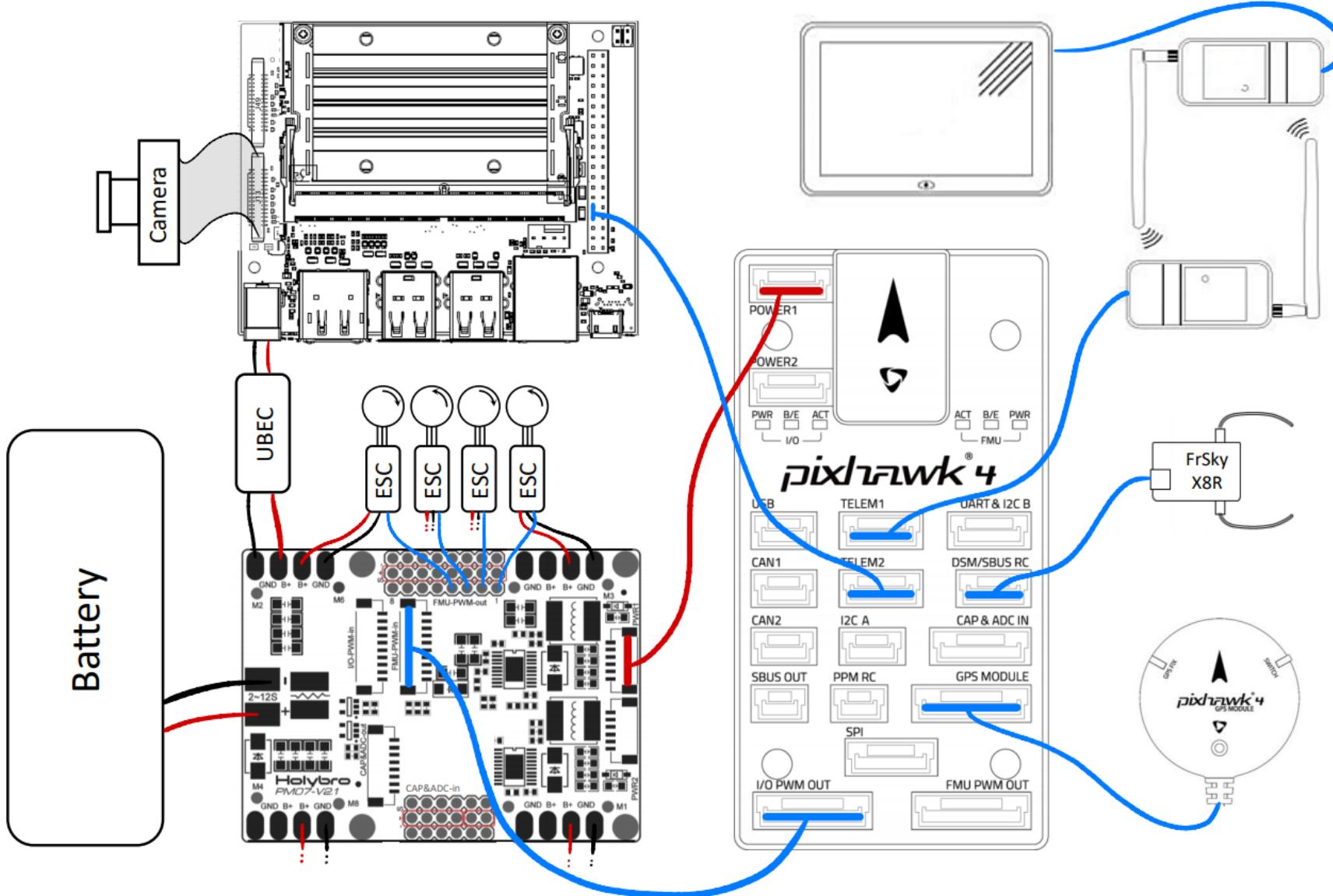
– The deployment platform

Nvidia Jetson Nano Development Kit B01
& ArduCam IMX477 Jetson compact camera



Jetson Nano	
AI Performance	472 GFLOPS
GPU	128-core NVIDIA Maxwell™ GPU
CPU	Quad-core ARM® Cortex®-A57 MPCore processor
Memory	4 GB 64-bit LPDDR4 25.6GB/s
Storage	16 GB eMMC 5.1
Power	5W 10W

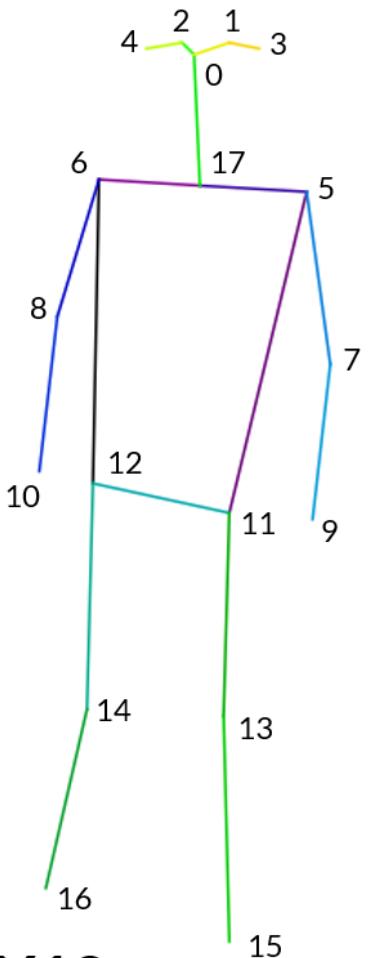
Hardware – Connections



Pose estimation – What is it?



- Torso (right)
- Torso (left)
- Shoulder (right)
- Shoulder (left)
- Arm (right)
- Forearm (right)
- Arm (left)
- Forearm (left)
- Hip
- Thigh (right)
- Leg (right)
- Thigh (left)
- Leg (left)
- Neck
- Eye (right)
- Ear (right)
- Eye (left)
- Ear (left)



Microsoft COCO 2020 Keypoint Detection Task
(250,000 person instances labeled)

BODY18

Source: Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. CoRR, abs/1405.0312, 2014.

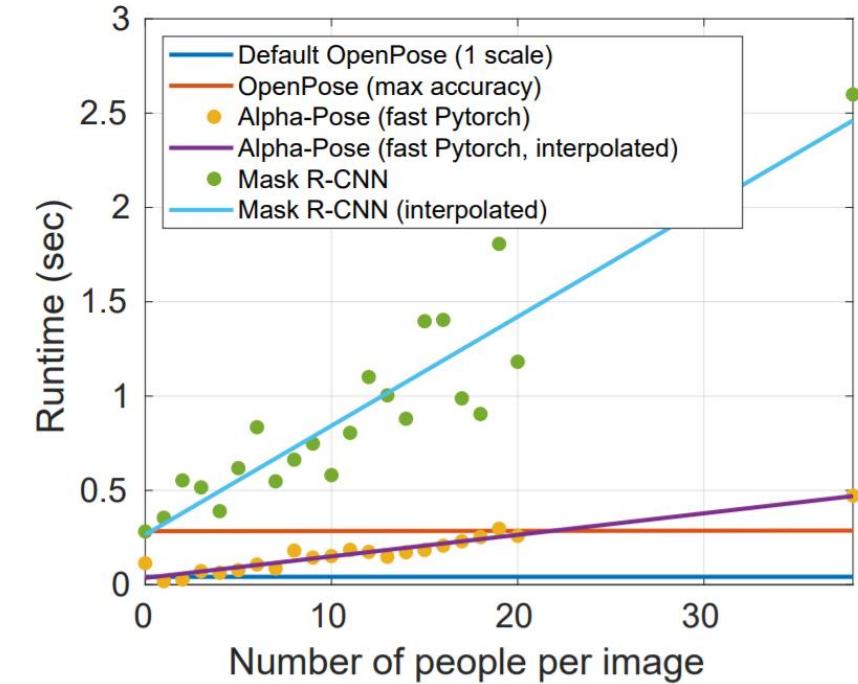
Pose estimation

– The Bottom-Up approach

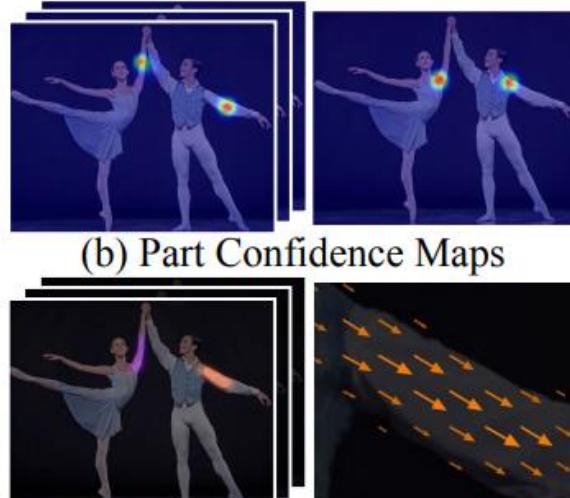
Part Confidence Maps – CMap
for body joints detection

Part Affinity Fields – PAF
for body joints association

Source: Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields, 2017. (arXiv:1611.08050)



(a) Input Image



(b) Part Confidence Maps



(c) Part Affinity Fields

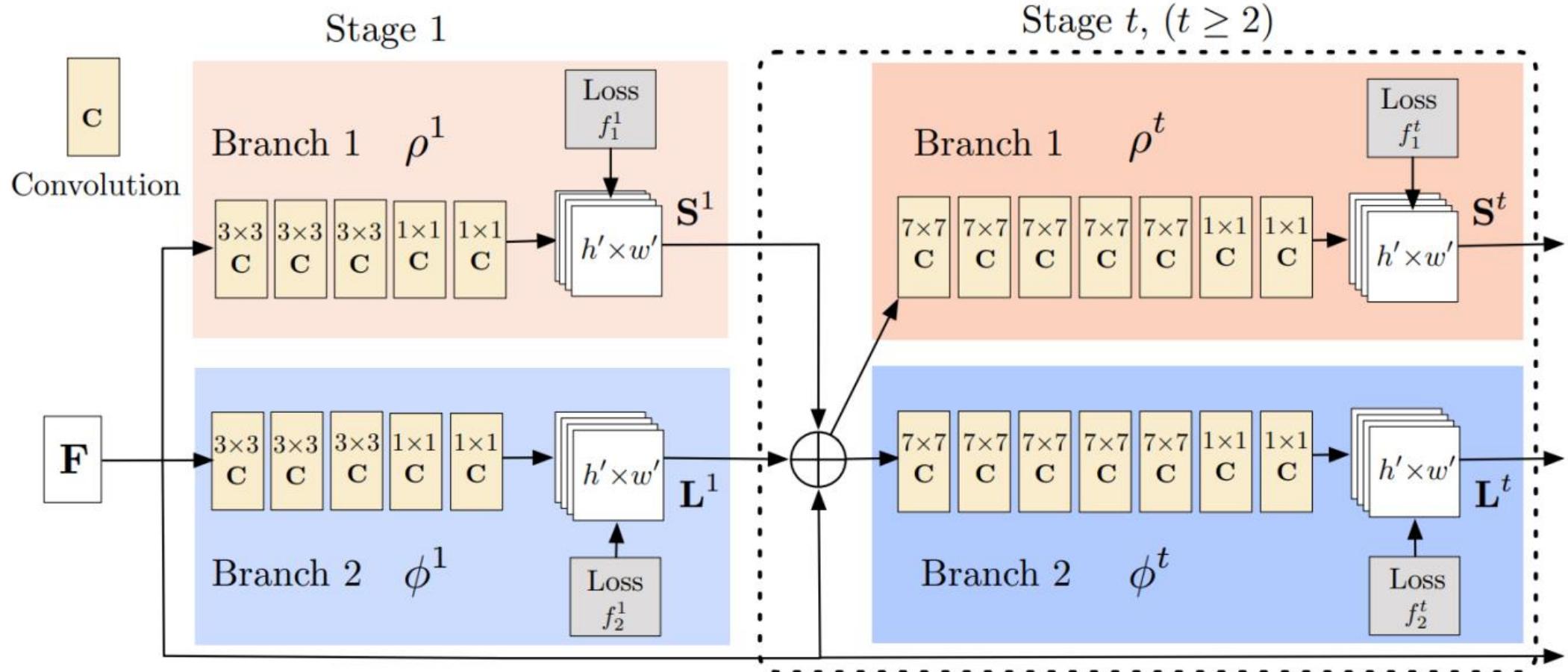


(d) Bipartite Matching

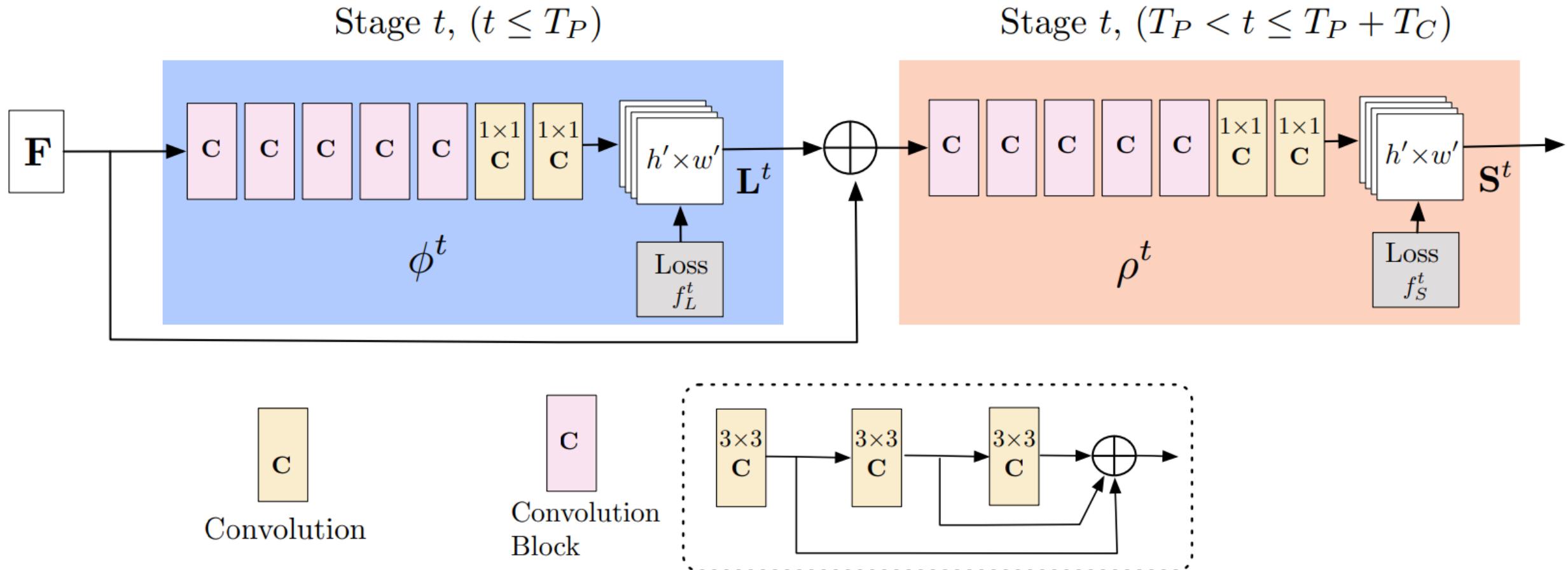


(e) Parsing Results

Pose estimation – Original architecture



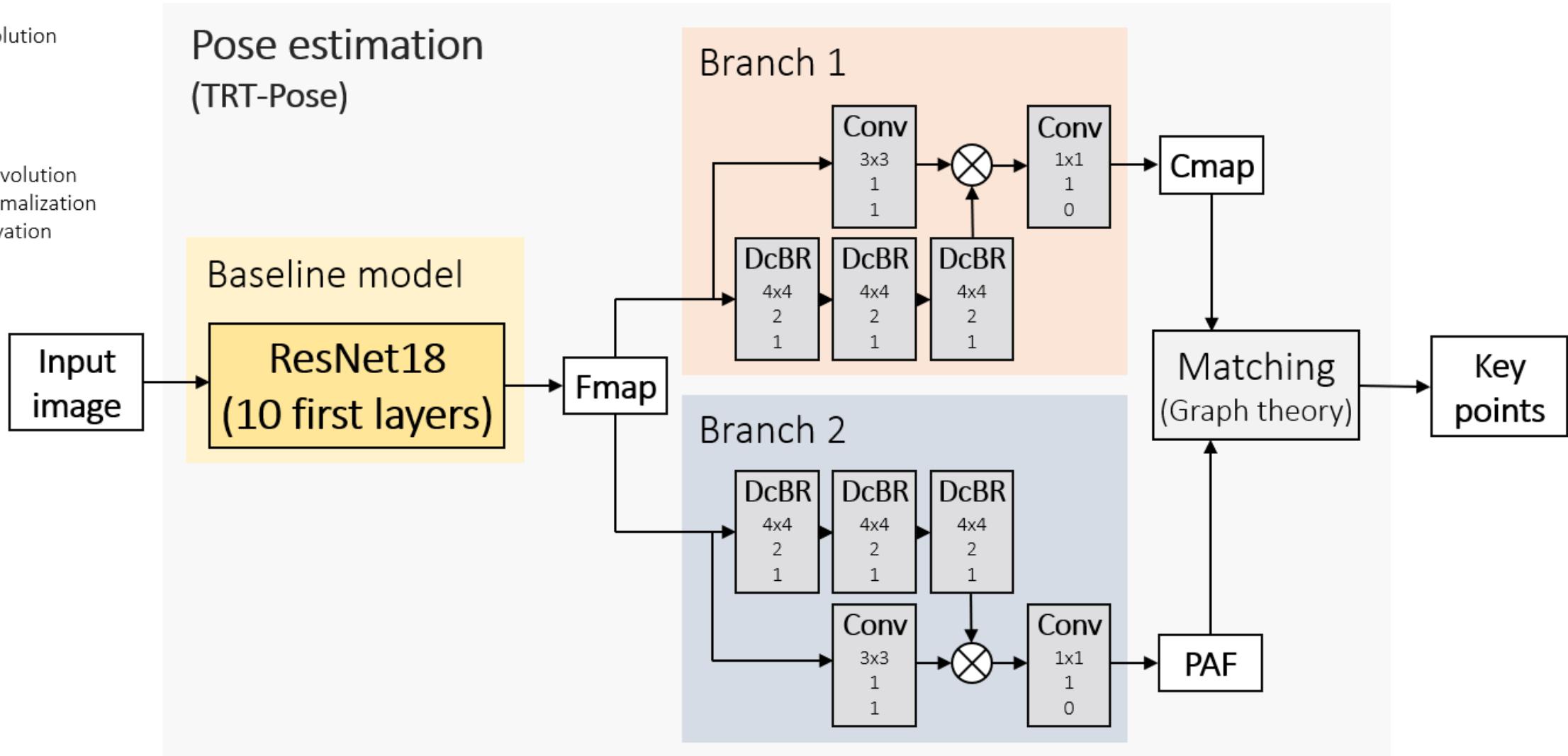
Pose estimation – OpenPose architecture



Pose estimation – TensorRT-Pose architecture

Conv
2D Convolution
Kernel
Stride
Padding

DcBR
2D Deconvolution
Batch normalization
ReLU activation
Kernel
Stride
Padding



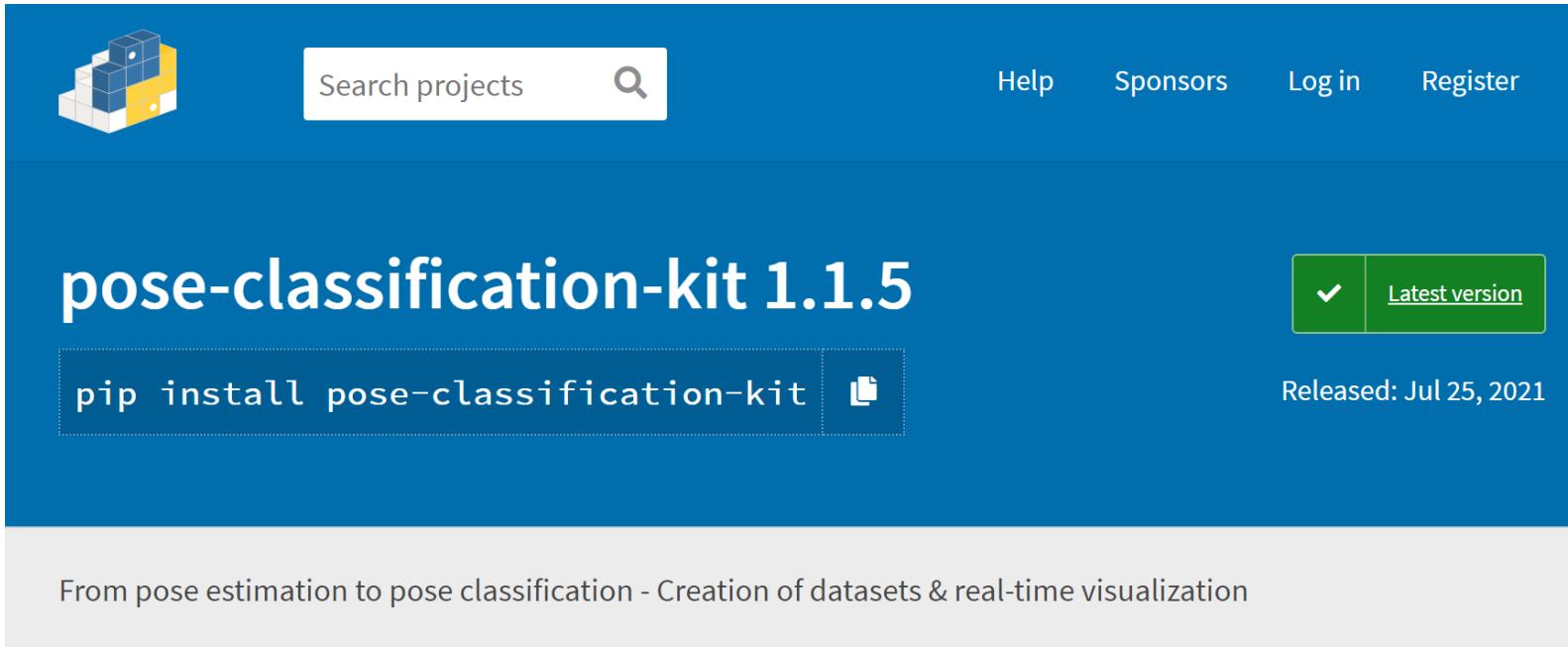
Pose estimation – TensorRT-Pose source code

```
1  class UpsampleCBR(torch.nn.Sequential):
2      def __init__(self, input_channels, output_channels):
3          layers = []
4          for i in range(3):
5              if i == 0:
6                  inch = input_channels
7              else:
8                  inch = output_channels
9          layers += [
10              torch.nn.ConvTranspose2d(inch, output_channels,
11                                     kernel_size=4, stride=2, padding=1),
12              torch.nn.BatchNorm2d(output_channels),
13              torch.nn.ReLU()
14          ]
15      super(UpsampleCBR, self).__init__(*layers)
16
17 class CmapPafHeadAttention(torch.nn.Module):
18     def __init__(self, input_channels, cmap_channels, paf_channels, upsample_channels=256):
19         super(CmapPafHeadAttention, self).__init__()
20         self.cmap_up = UpsampleCBR(input_channels, upsample_channels)
21         self.paf_up = UpsampleCBR(input_channels, upsample_channels)
22         self.cmap_att = torch.nn.Conv2d(upsample_channels, upsample_channels,
23                                         kernel_size=3, stride=1, padding=1)
24         self.paf_att = torch.nn.Conv2d(upsample_channels, upsample_channels,
25                                         kernel_size=3, stride=1, padding=1)
26
27         self.cmap_conv = torch.nn.Conv2d(upsample_channels, cmap_channels,
28                                         kernel_size=1, stride=1, padding=0)
29         self.paf_conv = torch.nn.Conv2d(upsample_channels, paf_channels,
30                                         kernel_size=1, stride=1, padding=0)
31
32     def forward(self, x):
33         xc = self.cmap_up(x)
34         ac = torch.sigmoid(self.cmap_att(xc))
35
36         xp = self.paf_up(x)
37         ap = torch.tanh(self.paf_att(xp))
38
39         return self.cmap_conv(xc * ac), self.paf_conv(xp * ap)
```

```
1  class ResNetBackbone(torch.nn.Module):
2
3      def __init__(self, resnet):
4          super(ResNetBackbone, self).__init__()
5          self.resnet = resnet
6
7      def forward(self, x):
8          x = self.resnet.conv1(x)
9          x = self.resnet.bn1(x)
10         x = self.resnet.relu(x)
11         x = self.resnet.maxpool(x)
12         x = self.resnet.layer1(x) # /4
13         x = self.resnet.layer2(x) # /8
14         x = self.resnet.layer3(x) # /16
15         x = self.resnet.layer4(x) # /32
16
17         return x
18
19     def resnet18_baseline_att(cmap_channels, paf_channels, upsample_channels=256, pretrained=True):
20         resnet = torchvision.models.resnet18(pretrained=pretrained)
21         return _resnet_pose_att(cmap_channels, paf_channels, upsample_channels, resnet, 512)
```

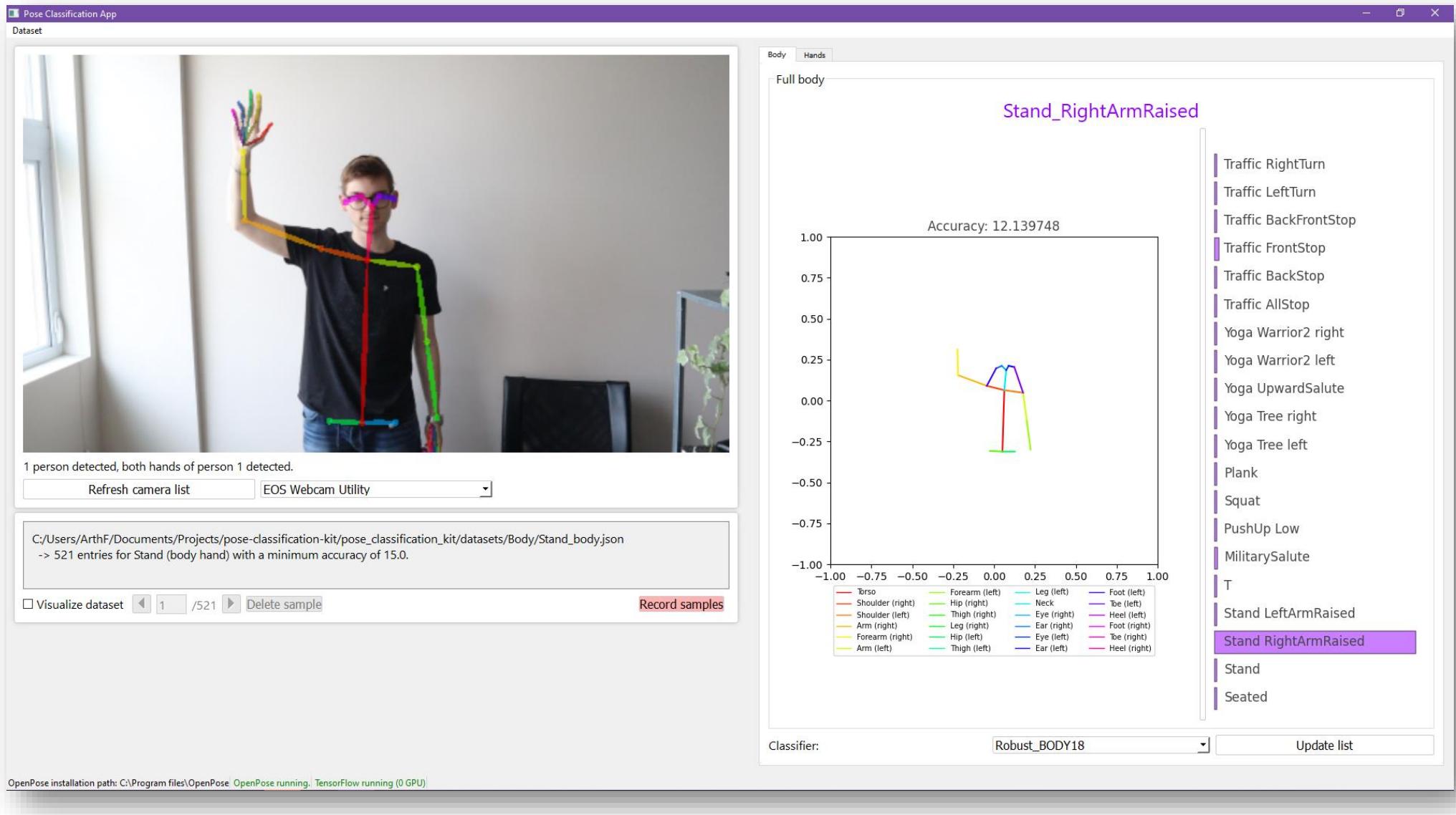
Pose classification

– Pose-Classification-Kit (PCK) package



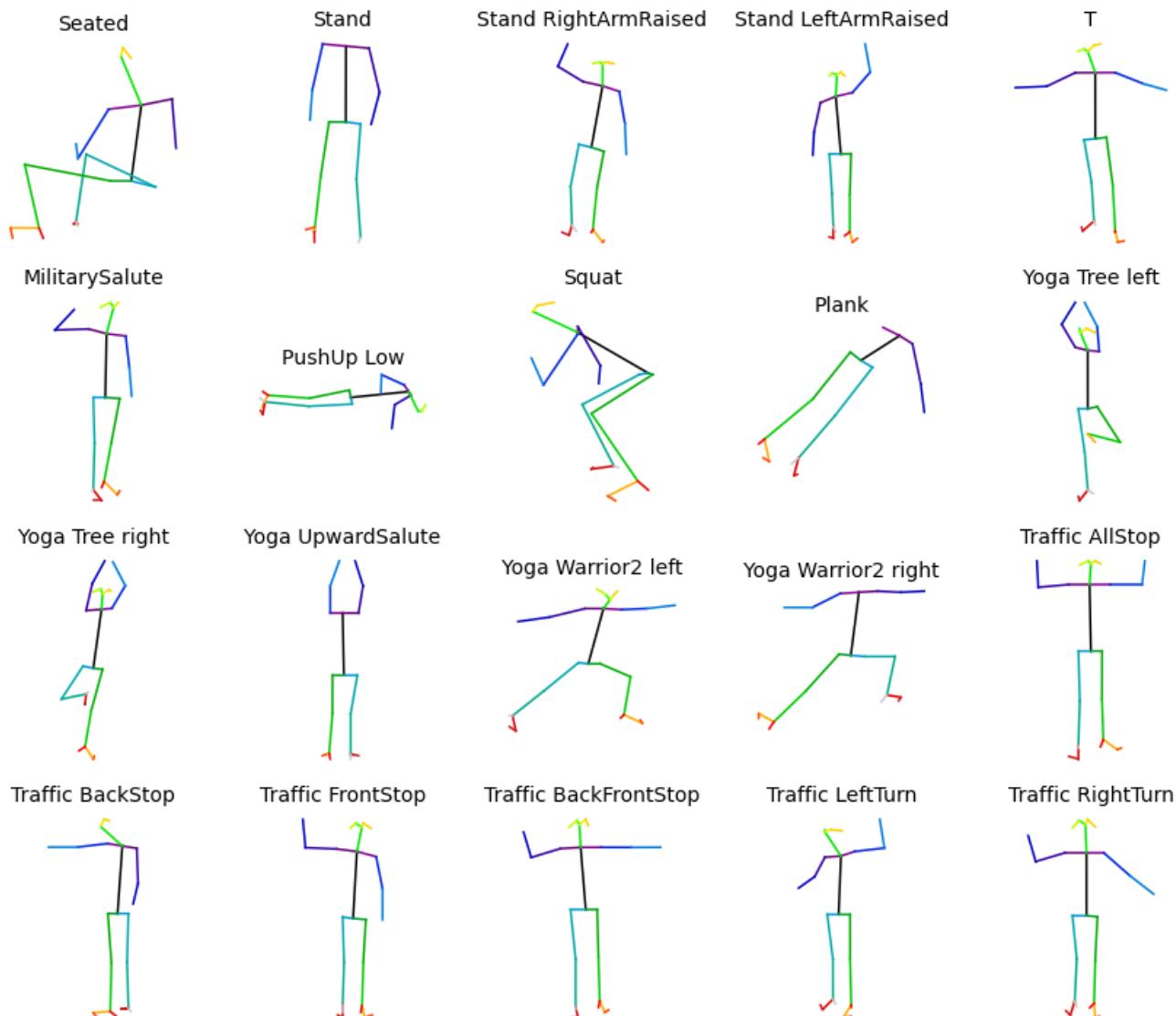
- **Dataset creation tool**
- **Two datasets** for body and hand gesture classification
- Several **pre-trained models**
- Utility functions for **data augmentation** and video analysis
- Example Jupyter Notebooks to ease **DL model creation**

Pose classification – PCK Application



Number of samples per class

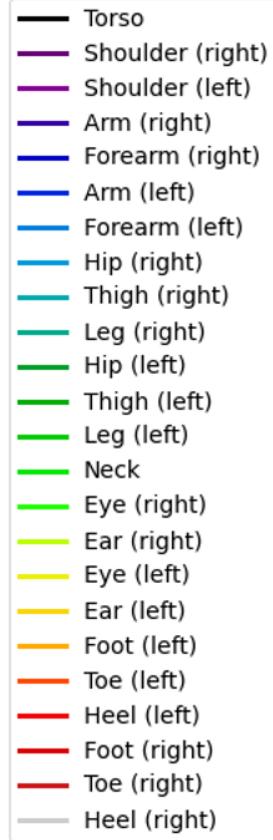
Pose classification – Body dataset



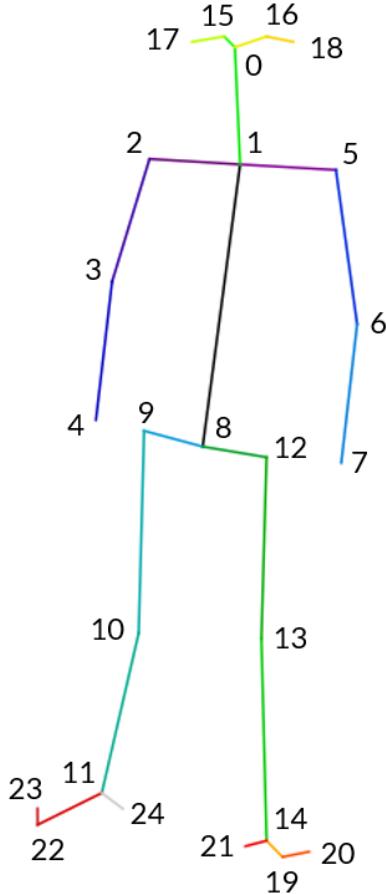
- Torso
- Shoulder (right)
- Shoulder (left)
- Arm (right)
- Forearm (right)
- Arm (left)
- Forearm (left)
- Hip (right)
- Leg (right)
- Hip (left)
- Thigh (left)
- Leg (left)
- Neck
- Eye (right)
- Ear (right)
- Eye (left)
- Ear (left)
- Foot (left)
- Toe (left)
- Heel (left)
- Foot (right)
- Toe (right)
- Heel (right)

	Number of samples per class
Yoga_Warrior2_left	521
Yoga_Warrior2_right	515
Yoga_Tree_left	580
Yoga_Tree_right	502
Traffic_AllStop	571
Traffic_BackStop	560
Traffic_FrontStop	528
Traffic_BackFrontStop	516
Traffic_LeftTurn	555
Traffic_RightTurn	521
Seated	606
Stand	521
Stand_RightArmRaised	503
Stand_LeftArmRaised	534
MilitarySalute	524
T	503
PushUp_Low	508
Squat	552
Plank	535
Yoga_UpwardSalute	526

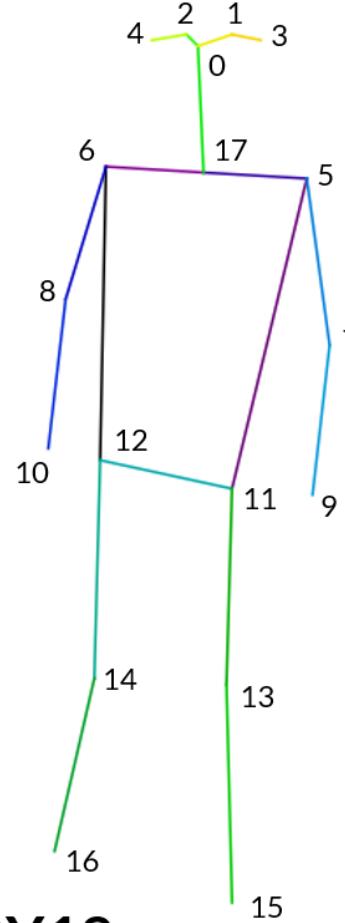
Pose classification – Body models & Normalization



BODY25



BODY18



Consider a sample x of format
BODY{I}, $\forall i \in [0, I - 1]$:

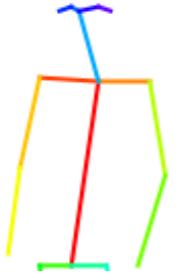
$$f(\vec{x}_i) = \frac{(\vec{x}_i - \langle \vec{x} \rangle)}{\max_{j \in \Delta} \delta_j \cdot \gamma_j}$$

Where δ_j is the length and γ_j the scaling factor of limb j

$$\Delta = \left\{ \begin{array}{l} \text{Torso, Neck, Shoulders,} \\ \text{Thigh (right), Thigh (left),} \\ \text{Leg (right), Leg (left)} \end{array} \right\}$$

Keypoints normalization

Pose classification – Dataset augmentation



Remove keypoints

Pre-defined or a random list of keypoints are removed (coordinates set to 0) from the sample.



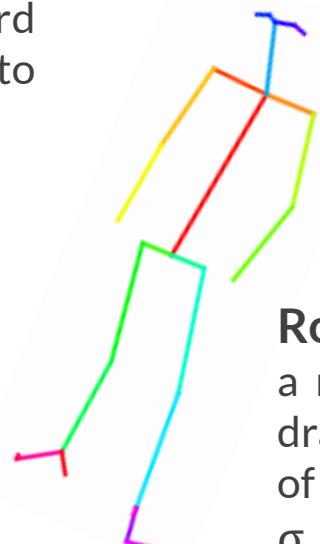
Noise

Gaussian noise of standard deviation σ_{noise} is added to coordinates of the sample.



Scaling

A random scaling factor drawn from a normal distribution of mean 0 and standard deviation σ_{scale} is applied to all sample coordinates.

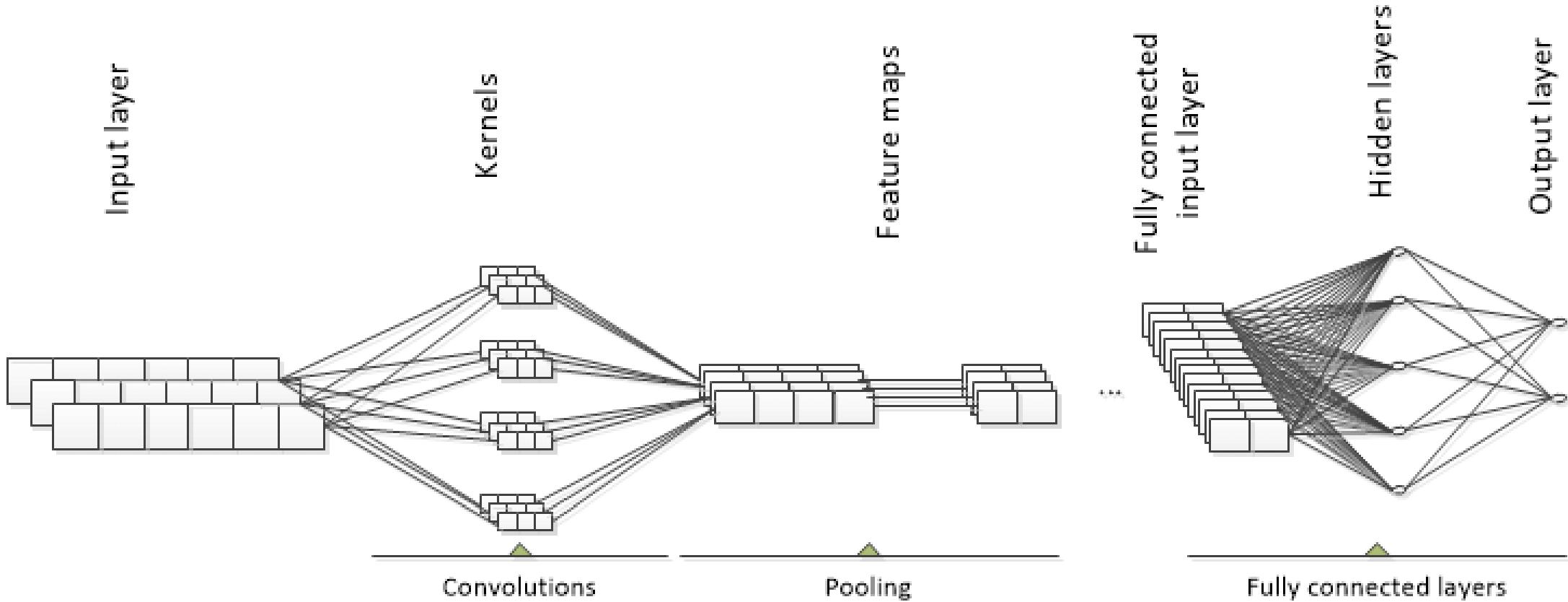


Rotation

a rotation of an angle randomly drawn from a normal distribution of mean 0 and standard deviation $\sigma_{rotation}$ is applied to the sample.

Augmentation Ratio	σ_{scale}	$\sigma_{rotation}$	σ_{noise}	Remove keypoints
10%	0.08	0.0	0.0	None
10%	0.0	10.0	0.0	None
15%	0.0	0.0	0.03	Legs
15%	0.0	0.0	0.03	Legs & Hip
20%	0.0	0.0	0.03	2 random

Pose classification – One dimensional convolution

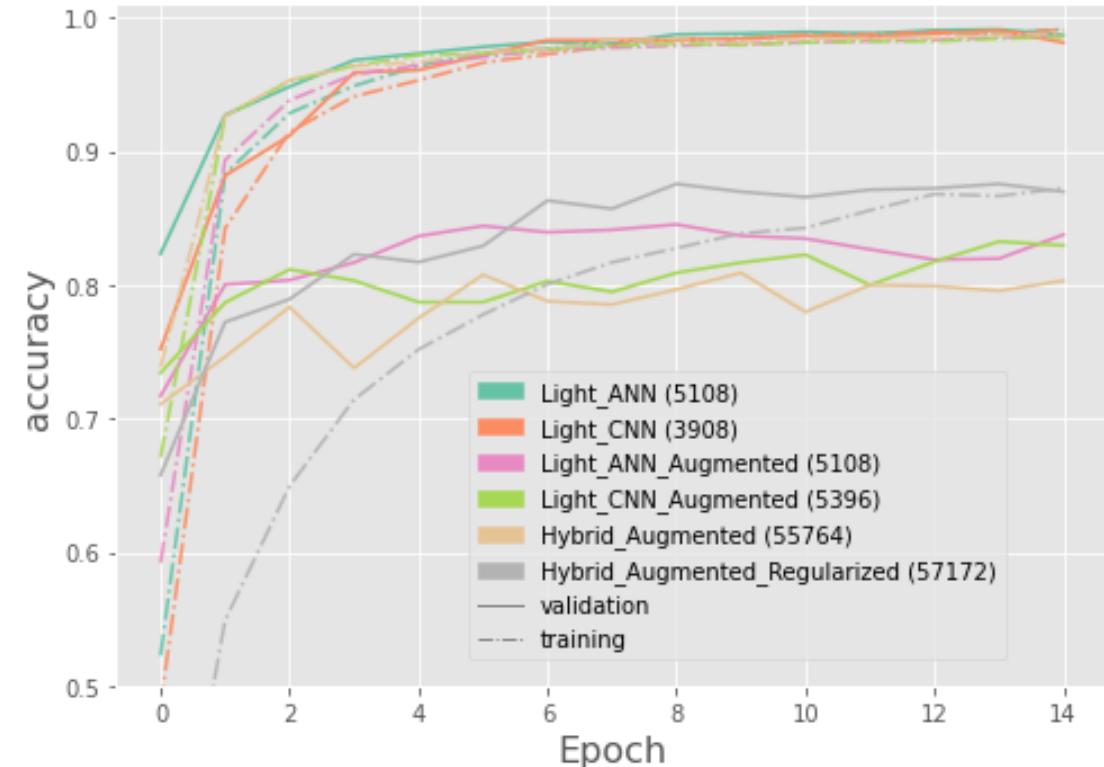
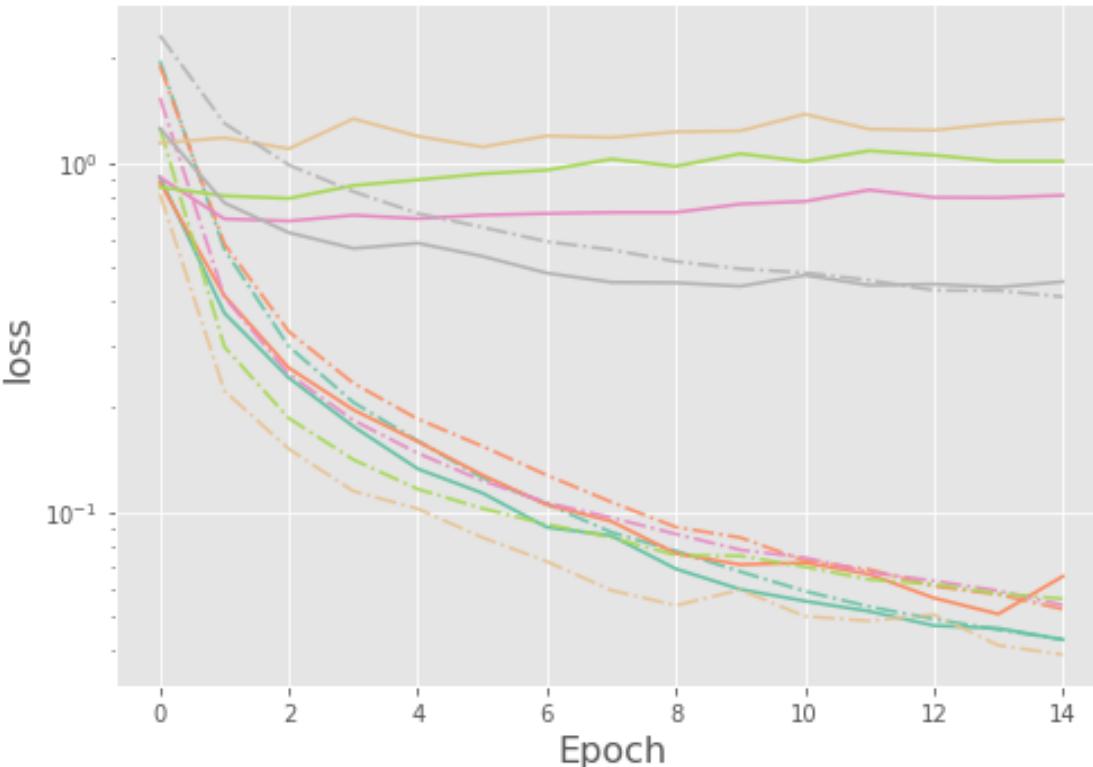
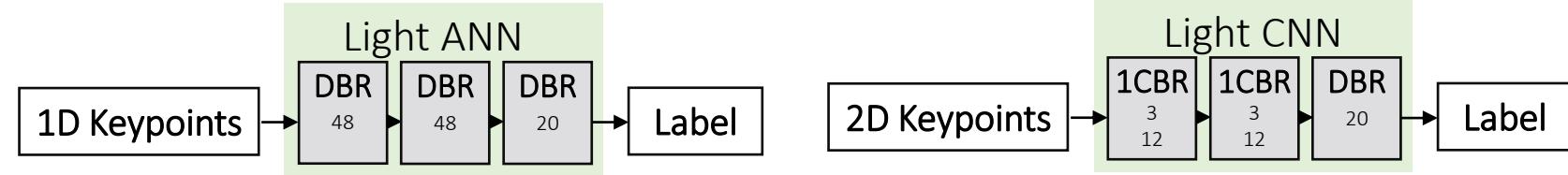


Pose classification

– Models training

1CBR
Kernel
Filters
1D Convolution
Batch normalization
Drop-out
ReLU activation

DBR
Units
Dense layer
Batch normalization
Drop-out
ReLU activation



Pose classification - Light CNN results

Original dataset

Testing accuracy: 98.25%

Augmented dataset

Testing accuracy: 50.95%

Pose classification - Hybrid model results

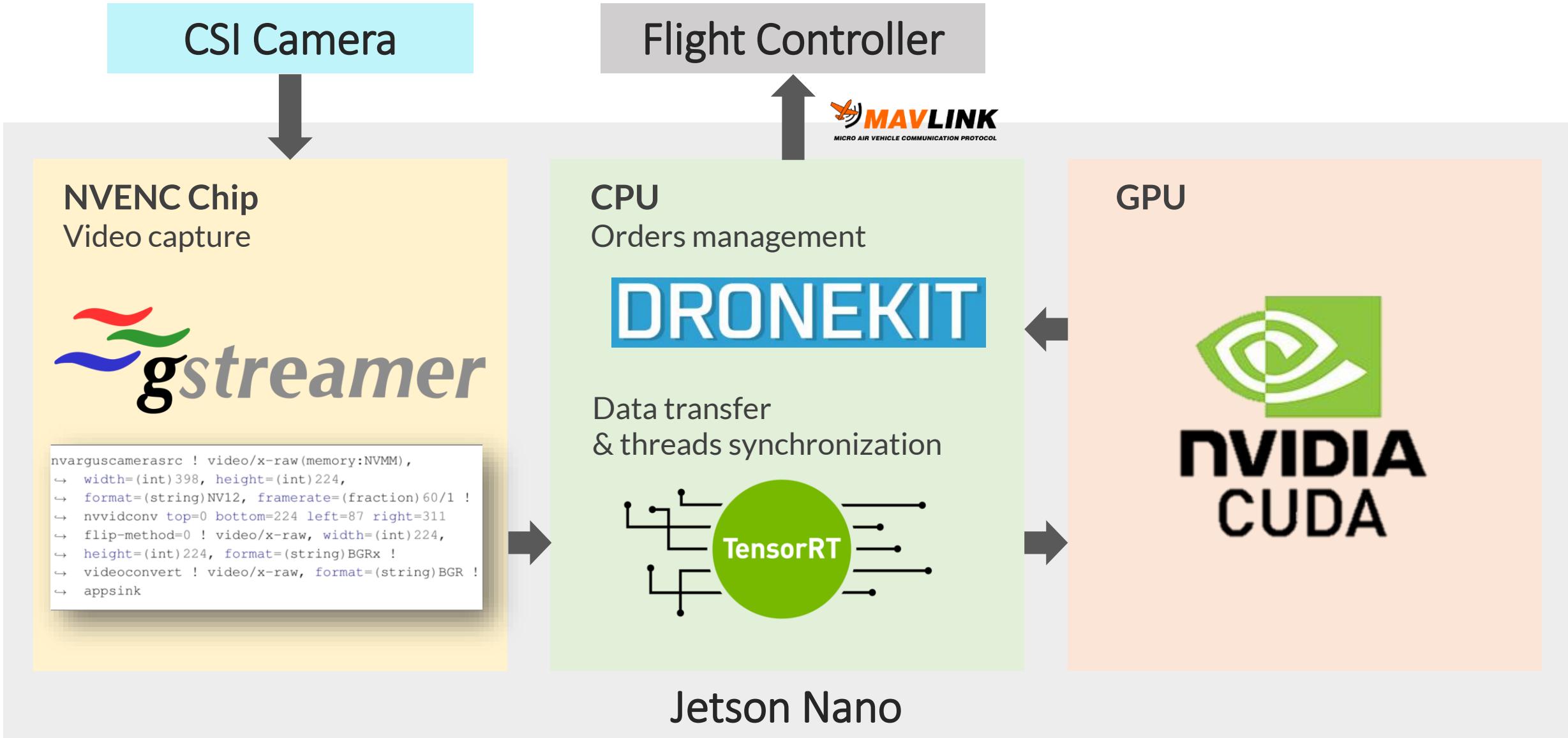
Original dataset

Testing accuracy: 98.30%

Augmented dataset

Testing accuracy: 95.05%

Deployment – Software stack



Deployment – Results

Pipeline inference time [ms] after optimization and compilation:

Original:	Data length	
	FP32	FP16
453		
Size	1GB	89 68
	33MB	86 67

Pose estimation model

Original:	Data length	
	FP32	FP16
184		
Size	1GB	6 6
	33MB	5 5

Pose classification model

Fail-safe systems:

- Output buffer to **filter-out misclassified labels**.
- Check the status of the drone to **ensure the feasibility of actions**.
- **Block the transmission of orders with a switch** on the radio controller and allow manual order override.
- **Ignore critical commands** on the serial port of the flight controller.

Orders manager: Gesture controls currently supported

Pose	Action
T	Arm the drone if it is disarmed and landed; Disarm the drone if it is armed and landed
Traffic_AllStop	Take-off at an altitude of 2.5m if the drone is armed and landed; Land if the drone is in flight
Traffic_RightTurn	Move 2m to the right if the drone is armed
Traffic_LeftTurn	Move 2m to the left if the drone is armed
Yoga_UpwardSalute	Return to Launch (RTL)

Demonstration



Future works

Technical improvements:

- **Concatenate both DL models** (estimation & classification) of the processing pipeline as a single model using a common framework.
- Deploy the pipeline on a **more efficient processing platform** such as Coral devices.
- **Increase input image resolution** to reach better accuracy and optimize multi-video feed simultaneous analysis.

Additional features:

- Augment the order manager with **new gestures** to broaden control possibilities.
- Create a new dynamic dataset to train recurrent neural network and allow **dynamic gestures detection**.
- **Add user tracking** to improve robustness and usefulness.
- **Add continuous controls** to the order manager to allow precise control.

