

CAMERA POSE ESTIMATION FROM SURGICAL VIDEOS WITH DEEP LEARNING

PROJECT N°7

STUDENTS : Andrea Naclerio

Ali Shadman Yazdi

Siavash Taleghani

TUTOR : Alberto Rota

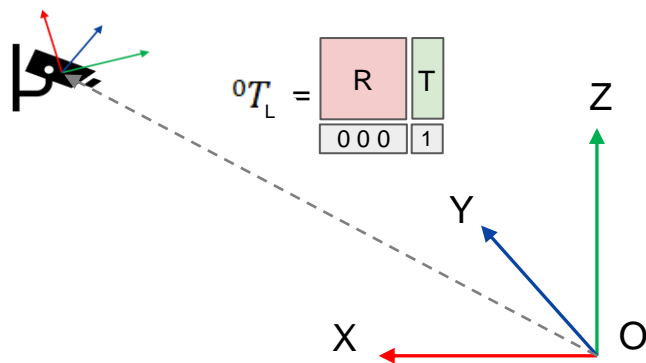
Camera Pose Estimation

What is camera pose estimation?

A homogeneous transformation matrix that describes position and orientation of a camera

Why do we need camera pose?

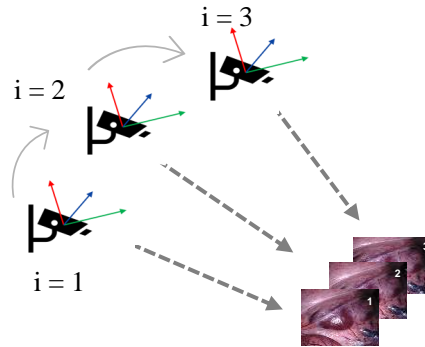
- 3D reconstruction
- Enhanced perception
- Organs and lesions localization



State of the art

Visual SLAM

- Tracking points of interest
- 3D positioning through triangulation
- Developing 3D map
- Mainly used in car localization

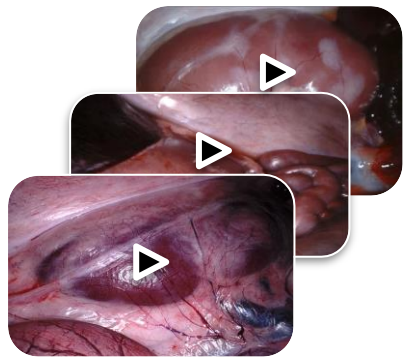


Deep Learning

- Feature extraction through DCNN
- Camera pose and image information relationship establishment
- Predicting camera pose

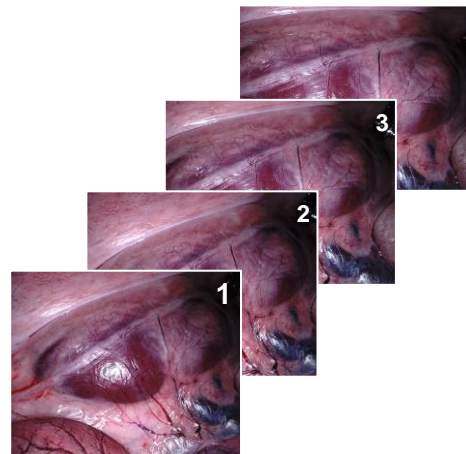
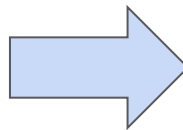
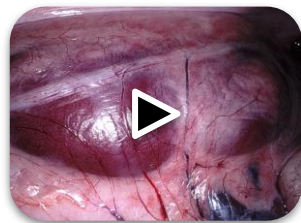
$${}^0T_i = \begin{bmatrix} R & T \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

A 3D coordinate system is shown with origin O . The Z -axis is a green arrow pointing upwards. The Y -axis is a blue arrow pointing towards the upper-left. The X -axis is a red arrow pointing towards the lower-left.

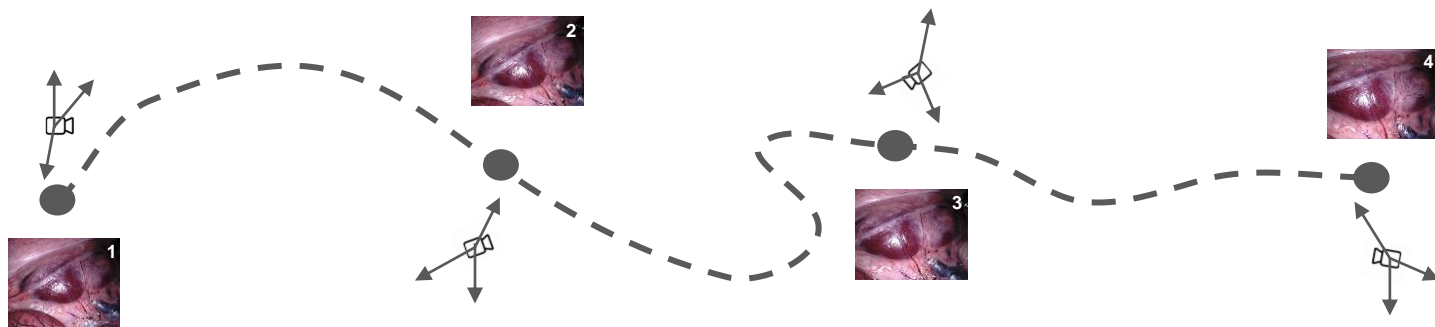


DATASET

27 monocular videos
annotated



FRAMES



TRAINING

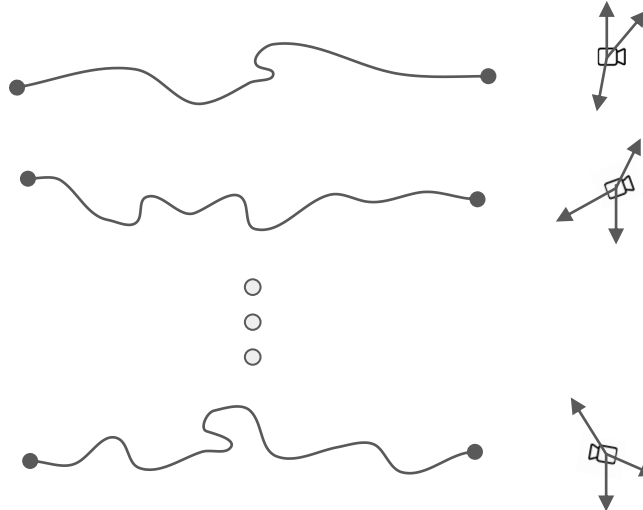
VIDEOS



FULLY SUPERVISED
MODEL



TRAJECTORY + ORIENTATION
(ground truth)



TEST

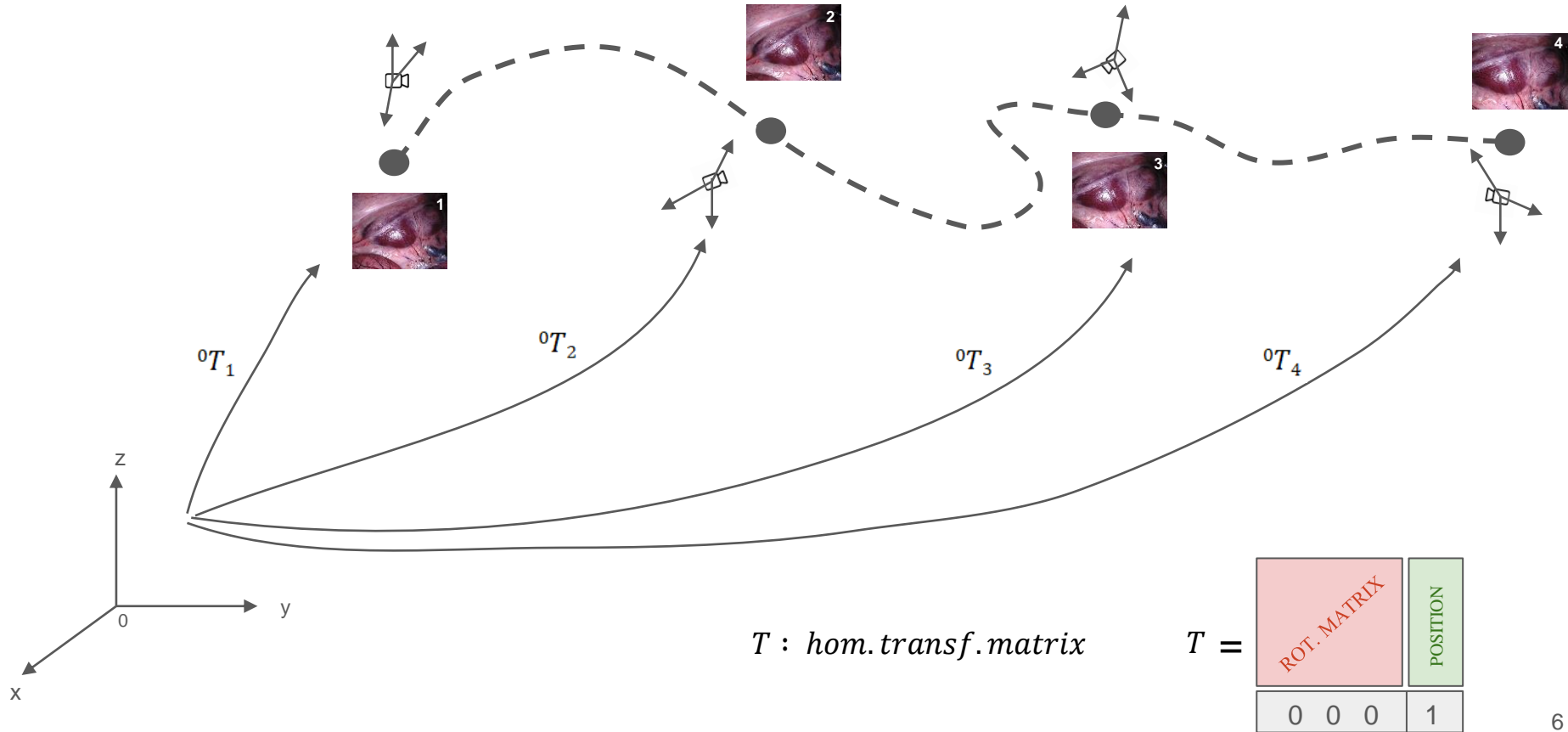


FULLY SUPERVISED
MODEL

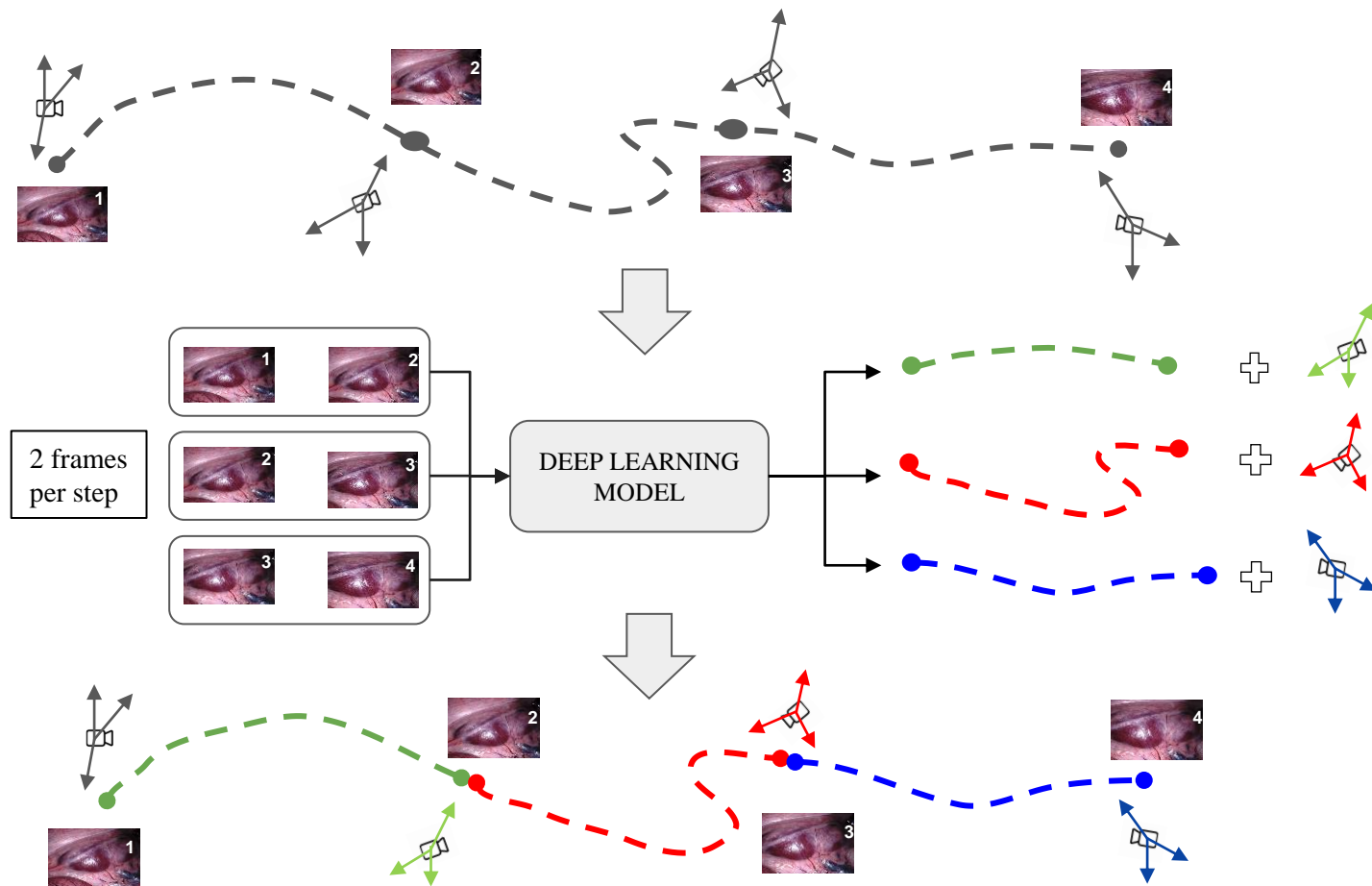


?

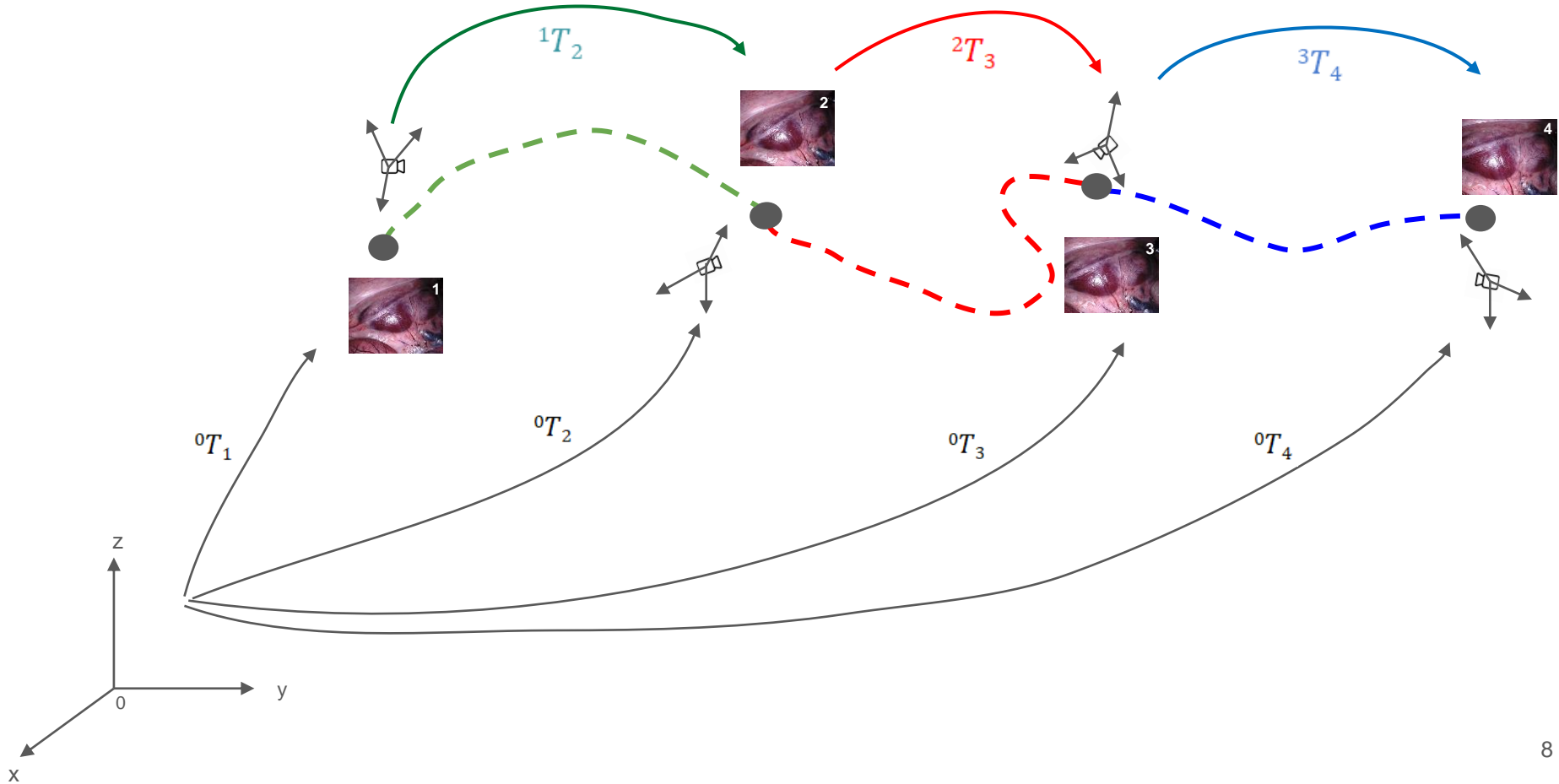
Ground-Truth data



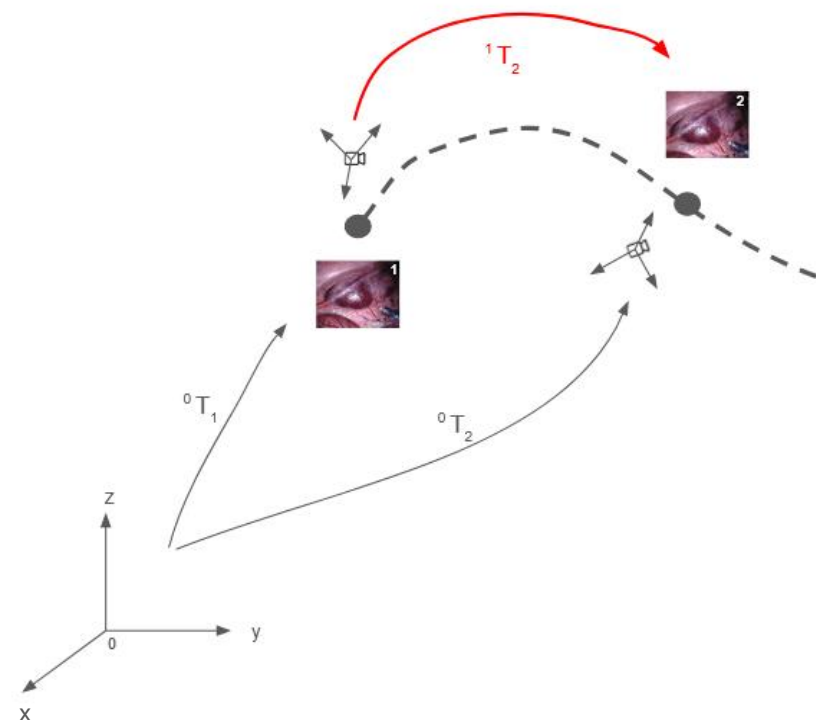
Problem decomposition



New transformation matrix



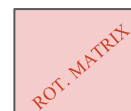
Output computation



$${}^1T_2 = {}^0T_1^{-1} * {}^0T_2$$



$${}^1T_2 = \begin{array}{|c|c|} \hline \text{ROT. MATRIX} & \text{POSITION} \\ \hline 0 & 1 \\ \hline \end{array}$$

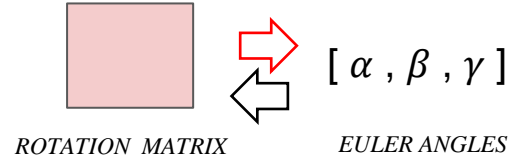


$[\alpha, \beta, \gamma]$



$[x, y, z]$

INVERSE PROBLEM



$$R_{x,y',z''}(\alpha, \beta, \gamma) = \begin{bmatrix} c\beta c\gamma & -c\beta s\gamma & s\beta \\ sa s\beta c\gamma + ca s\gamma & -sa s\beta s\gamma + ca c\gamma & -sa c\beta \\ -s\beta ca c\gamma + sa s\gamma & s\beta ca s\gamma + sa c\gamma & ca c\beta \end{bmatrix}$$

$$R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix}$$

$$\beta = \text{Atan2}(\pm \sqrt{r_{23}^2 + r_{33}^2}, r_{13})$$

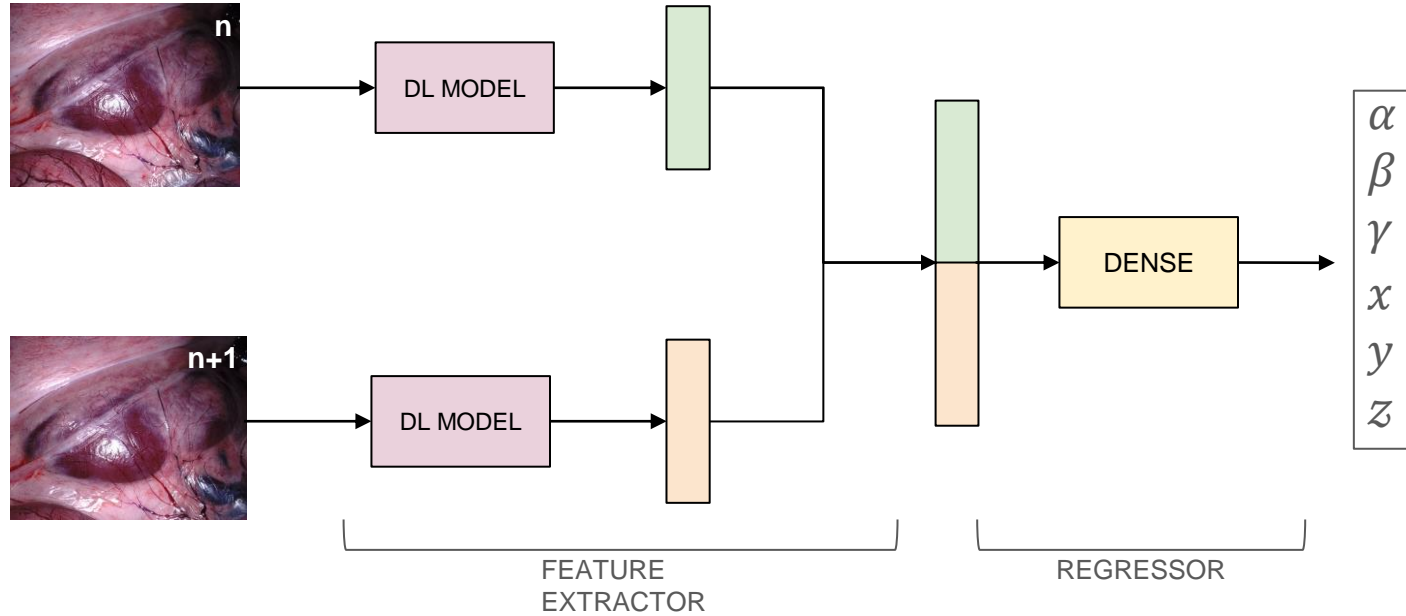
$$\alpha = \text{Atan2}\left(\frac{r_{33}}{c\beta}, \frac{-r_{32}}{c\beta}\right)$$

$$\gamma = \text{Atan2}\left(\frac{r_{11}}{c\beta}, \frac{-r_{12}}{c\beta}\right)$$

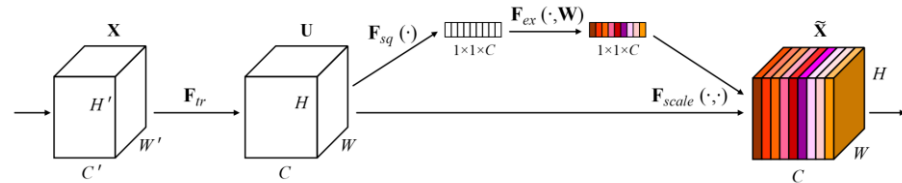
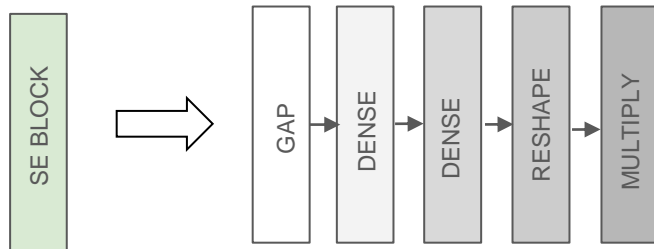
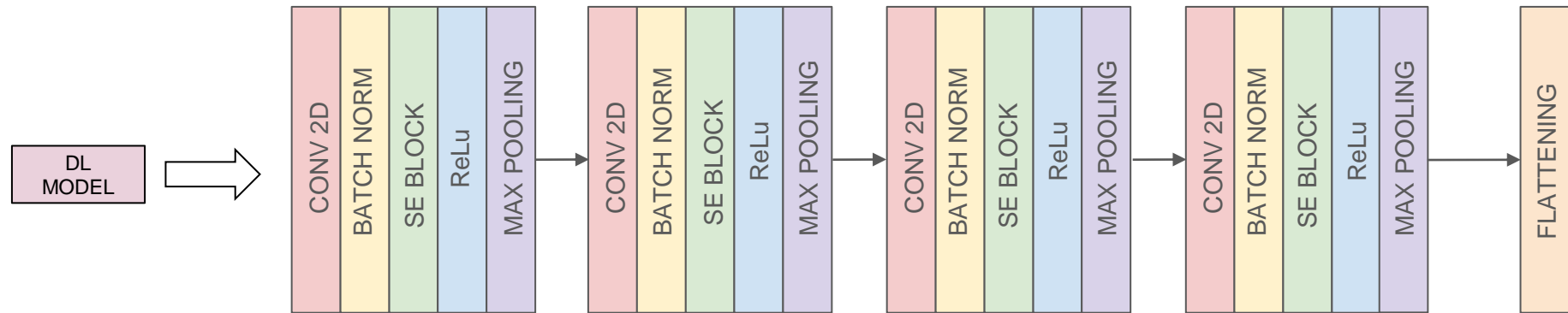


$$y = [\alpha \ \beta \ \gamma \mid x \ y \ z]$$

Two tails architecture



From scratch architecture

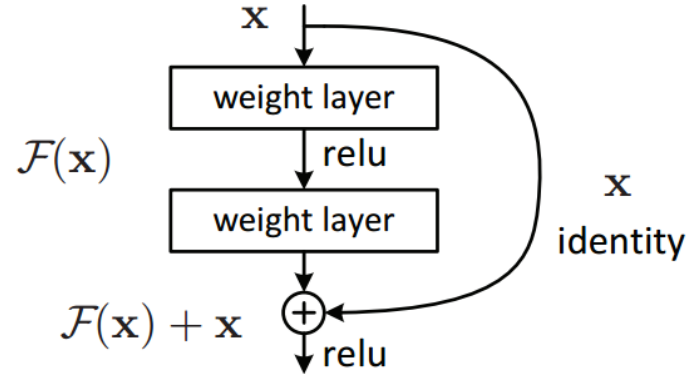


ResNet50 MODEL

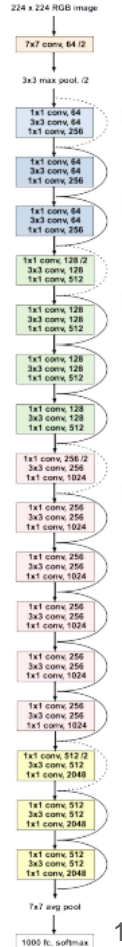
- ResNet-50 introduced the concept of residual learning to address the challenge of training very deep neural networks and mitigating the vanishing gradient problem.
- A residual block contains a shortcut or skip connection that skips one or more layers, allowing the network to learn the residual (difference) between the input and output of those skipped layers.
- Weights were initialized to pre-trained model from image-net, but adjusted to fit our dataset during training

"Instead of figuring out the entire journey, let's figure out the changes we need to make at each step, and then just add those changes to the starting point."

- chat GPT et al.(2024)

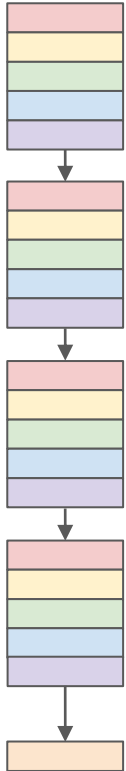


ResNet-50



Model training

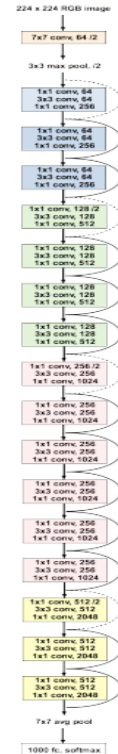
FROM SCRATCH MODEL



- TRAINING VIDEOS : 3 (v2, v3, v4)
- NUM. TRAINING IMAGES : 1784
- TEST VIDEO: 1 (v1)
- LOSS FUNCTION : **Mean Square Error (MSE)**
- BATCH : 12
- LEARNING RATE (LR) : **0.0001**
- CALLBACKS:
 - EARLY STOPPING PATIENCE : 12
 - REDUCE LR ON PLATEAU FACTOR : 6
- NUM. PARAMETERS : **40 millions**

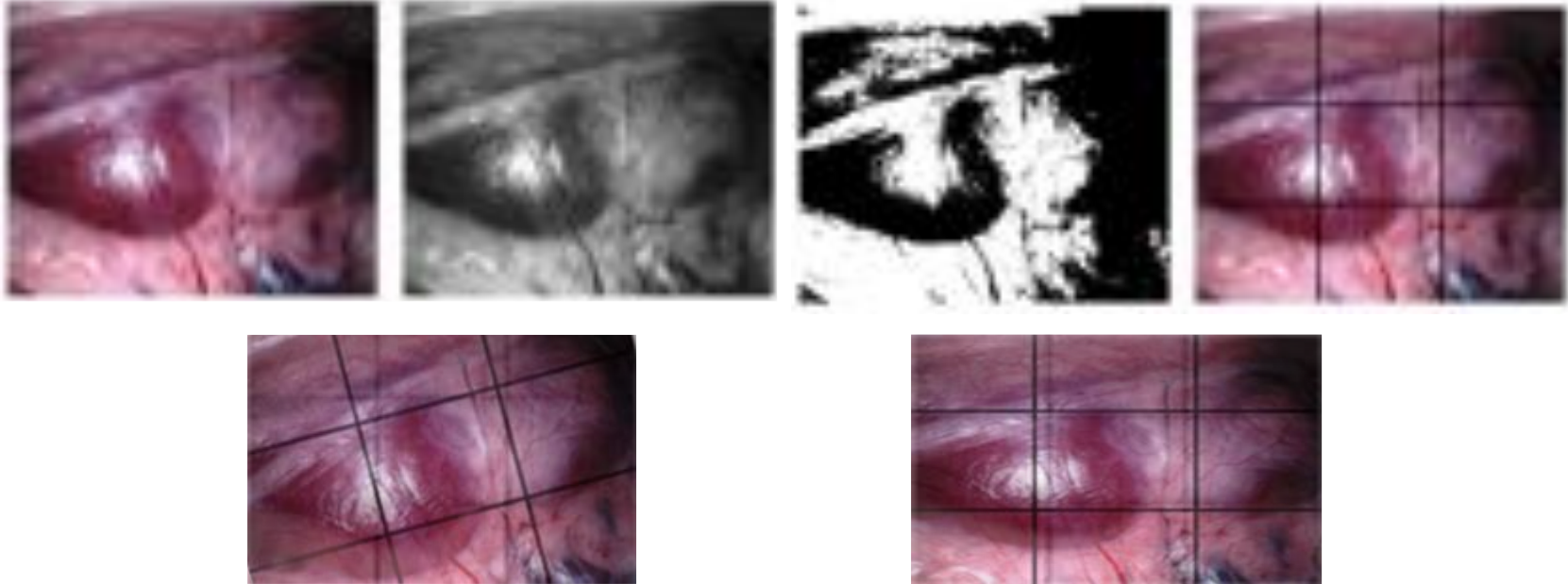
RESNET50 MODEL

ResNet-50







- TRAINING VIDEOS : 3 (v2, v3, v4)
- NUM. TRAINING IMAGES : 1784
- TEST VIDEO: 1 (v1)
- LOSS FUNCTION : **Mean Square Error (MSE)**
- BATCH : 4
- LEARNING RATE (LR) : **0.0001**
- CALLBACKS:
 - EARLY STOPPING PATIENCE : 10
 - REDUCE LR ON PLATEAU FACTOR : 3
- NUM. PARAMETERS : **340 millions**

Image preprocessing



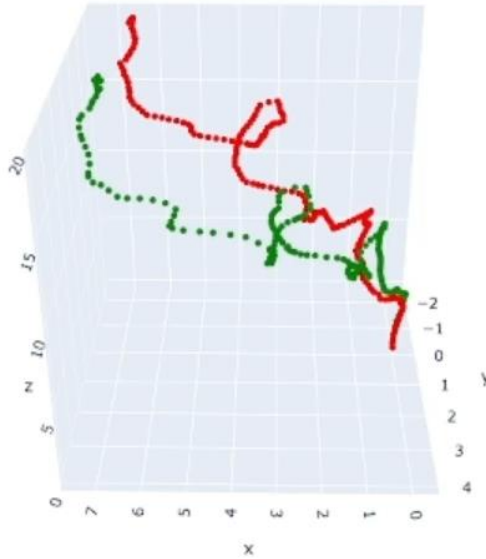
- RGB image $\rightarrow (m, 512, 640, 3) \rightarrow 983,040$ data points per image
- Binary image $\rightarrow (m, 512, 640, 1) \rightarrow 327,680$ data points per image

Effects of preprocessing

		Translation error	Rotation error
Base line		--	--
Grayscale		↑	↓
Binary		↑	↑
Grid		↑	↓

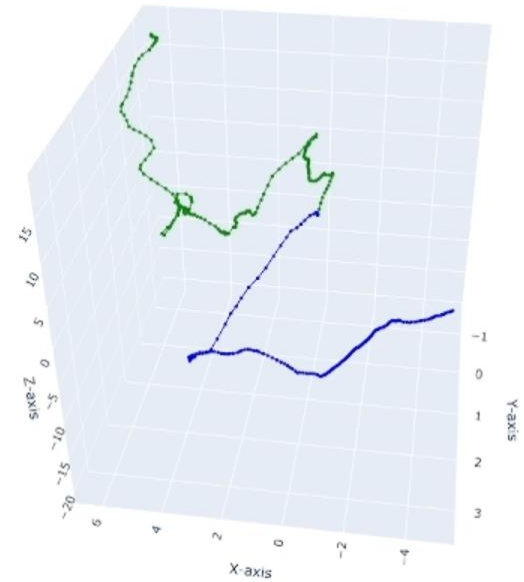
FINAL RESULTS

FROM SCRATCH MODEL



FROM SCRATCH MODEL VS TARGET TRAJECTORY

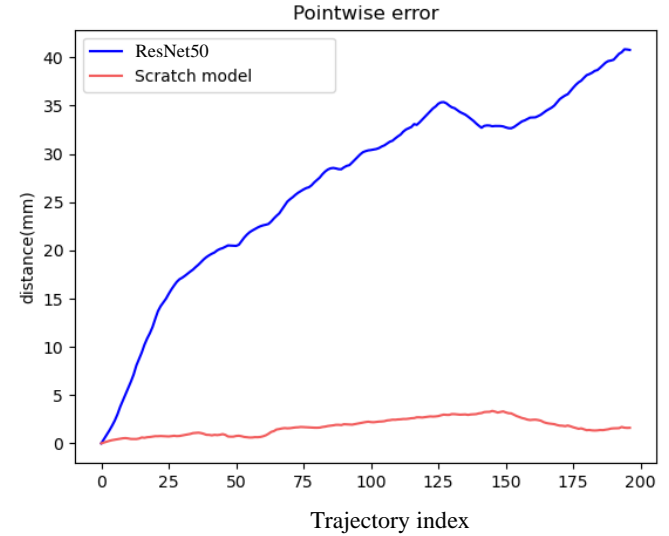
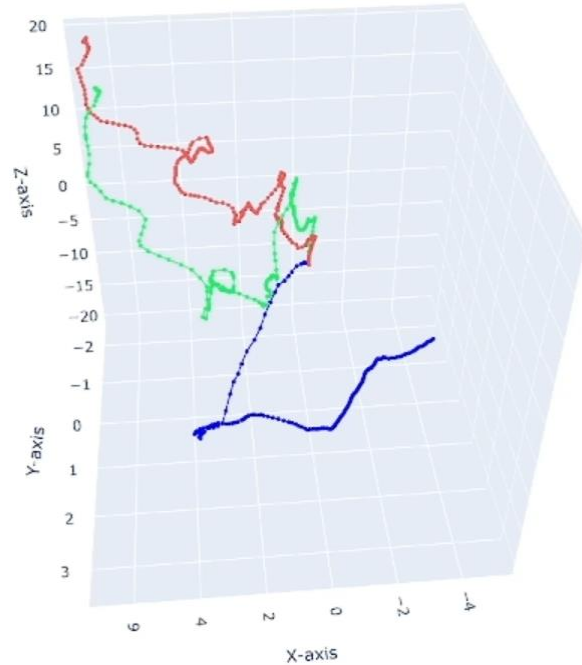
RESNET50 MODEL



RESNET 50 VS TARGET TRAJECTORY

RESULTS COMPARISON

FROM SCRATCH MODEL VS RESNET 50 VS TARGET TRAJECTORY



NUMERICAL METRICS

	Scratch model	Resnet50
Mean Angular Error (degree)	10.7	12.4
Euclidian Distance (mm)	1.7	26.9

CONCLUSION

Limitation

- Texture-less tissues and Lambertian Reflection
- GPUs limitation
- Computationally intensive
- Overfitting

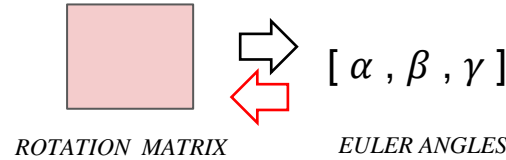
Future work

- Structural similarity
- Larger sequential
- Non sub-sequence frames
- Pretrained model on endoscopic camera dataset
- Ensembled light weight models

*THANK YOU
FOR YOUR ATTENTION*

QUESTIONS ?

DIRECT PROBLEM



R : rotation matrix

$$R_{x,y',z''}(a, \beta, \gamma) = R_x * R_{y'} * R_{z''}$$

$$R_x = \begin{bmatrix} 1 & 0 & 0 \\ 0 & ca & -sa \\ 0 & sa & ca \end{bmatrix}$$

$$R_{y'} = \begin{bmatrix} c\beta & 0 & s\beta \\ 0 & 1 & 0 \\ -s\beta & 0 & c\beta \end{bmatrix}$$

$$R_{z''} = \begin{bmatrix} c\gamma & -s\gamma & 0 \\ s\gamma & c\gamma & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R_{x,y',z''}(a, \beta, \gamma) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & ca & -sa \\ 0 & sa & ca \end{bmatrix} * \begin{bmatrix} c\beta & 0 & s\beta \\ 0 & 1 & 0 \\ -s\beta & 0 & c\beta \end{bmatrix} * \begin{bmatrix} c\gamma & -s\gamma & 0 \\ s\gamma & c\gamma & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R_{x,y',z''}(a, \beta, \gamma) = \begin{bmatrix} c\beta & 0 & s\beta \\ sa s\beta & ca & -sa c\beta \\ -s\beta ca & sa & ca c\beta \end{bmatrix} * \begin{bmatrix} c\gamma & -s\gamma & 0 \\ s\gamma & c\gamma & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R_{x,y',z''}(a, \beta, \gamma) = \begin{bmatrix} c\beta c\gamma & -c\beta s\gamma & s\beta \\ sa s\beta c\gamma + ca s\gamma & -sa s\beta s\gamma + ca c\gamma & -sa c\beta \\ -s\beta ca c\gamma + sa s\gamma & s\beta ca s\gamma + sa c\gamma & ca c\beta \end{bmatrix}$$