

# Comparison of Spatio-Temporal Hand Pose Denoising Models

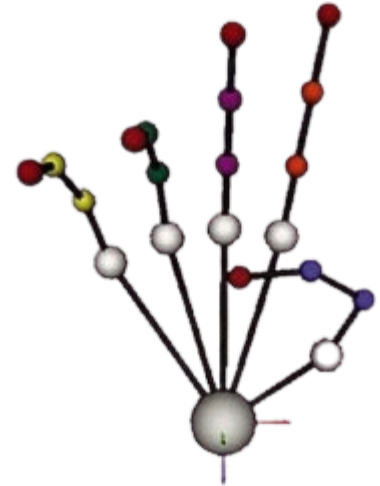
Johnny Núñez Cano

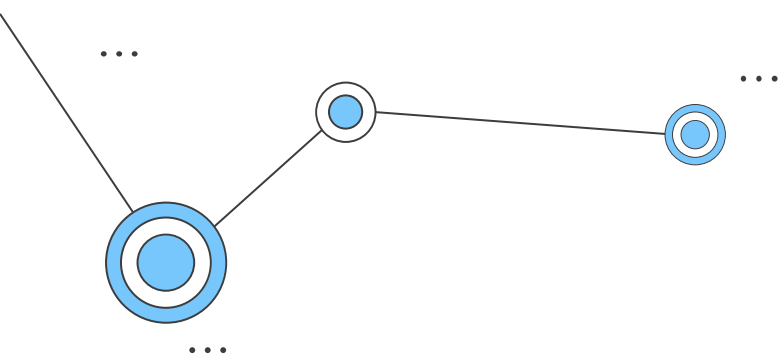
Directors:

Dr. Sergio Escalera

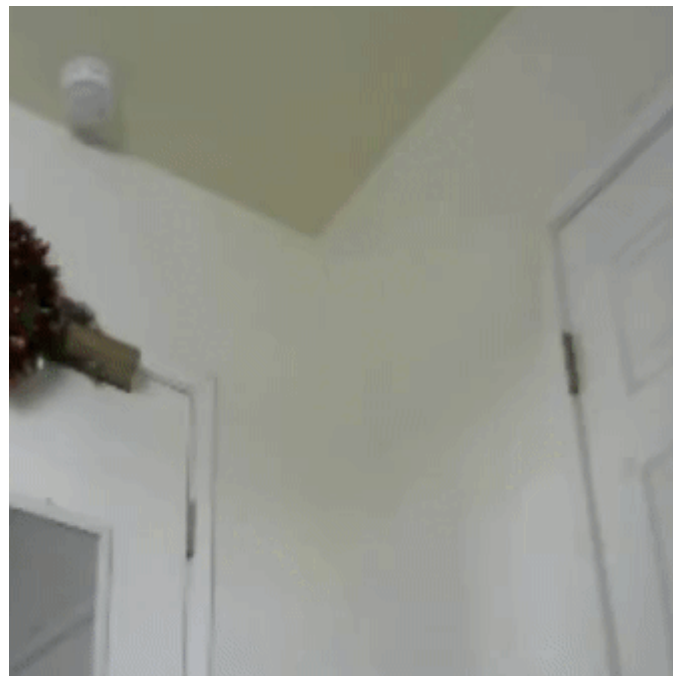
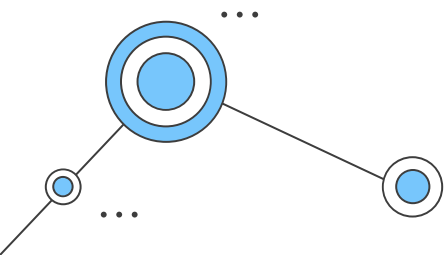
MSc. German Barquero

MSc. Cristina Palmero





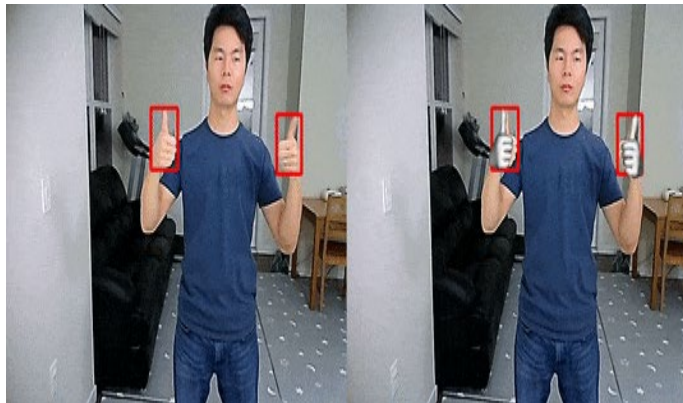
# Hand Pose Estimation



[https://github.com/NVIDIA-AI-IOT/trt\\_pose\\_hand](https://github.com/NVIDIA-AI-IOT/trt_pose_hand)

# Estimation Algorithms

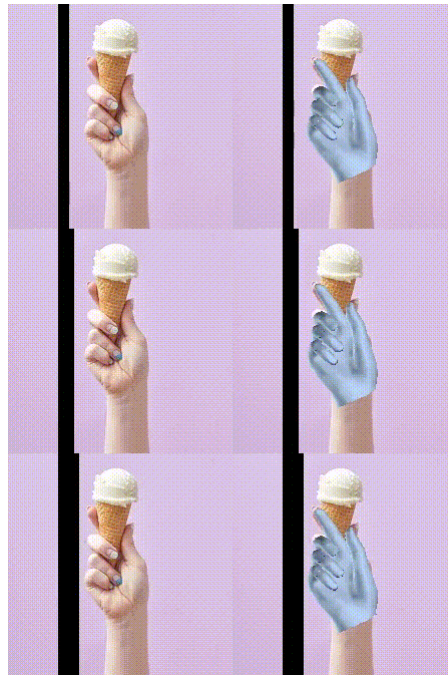
- FrankMocap



Rong, Yu, et al.. Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. In IEEE ICCV Workshops, 2021.

# Estimation Algorithms

- FrankMocap
- Mesh Graphormer



Rong, Yu, et al.. Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. In IEEE ICCV Workshops, 2021.  
Lin, Kevin and Wang, Lijuan and Liu, Zicheng.. Mesh Graphormer

# Typical Errors

01  
...

## Jitter

Small displacement.

02  
...

## Inversion

Error on the same instance.

03  
...

## Swap

Error on different instances.

04  
...

## Miss

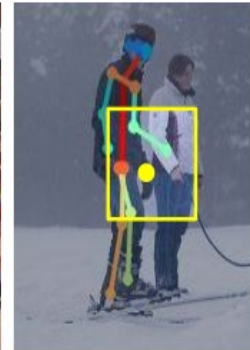
Not exists keypoint.



(a) Jitter



(b) Inversion

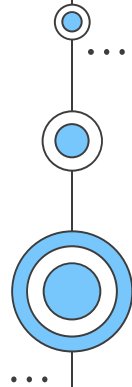
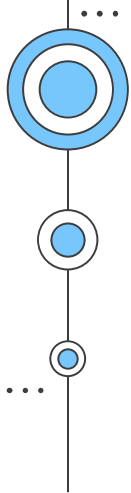
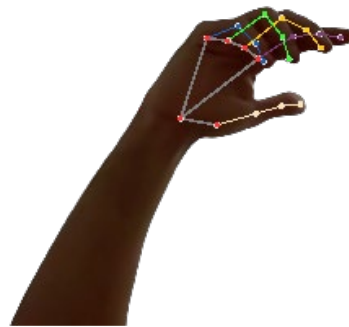
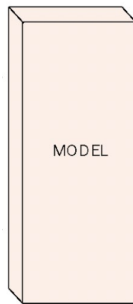
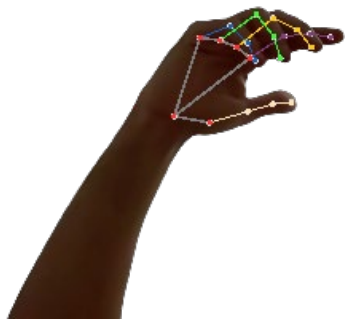


(c) Swap

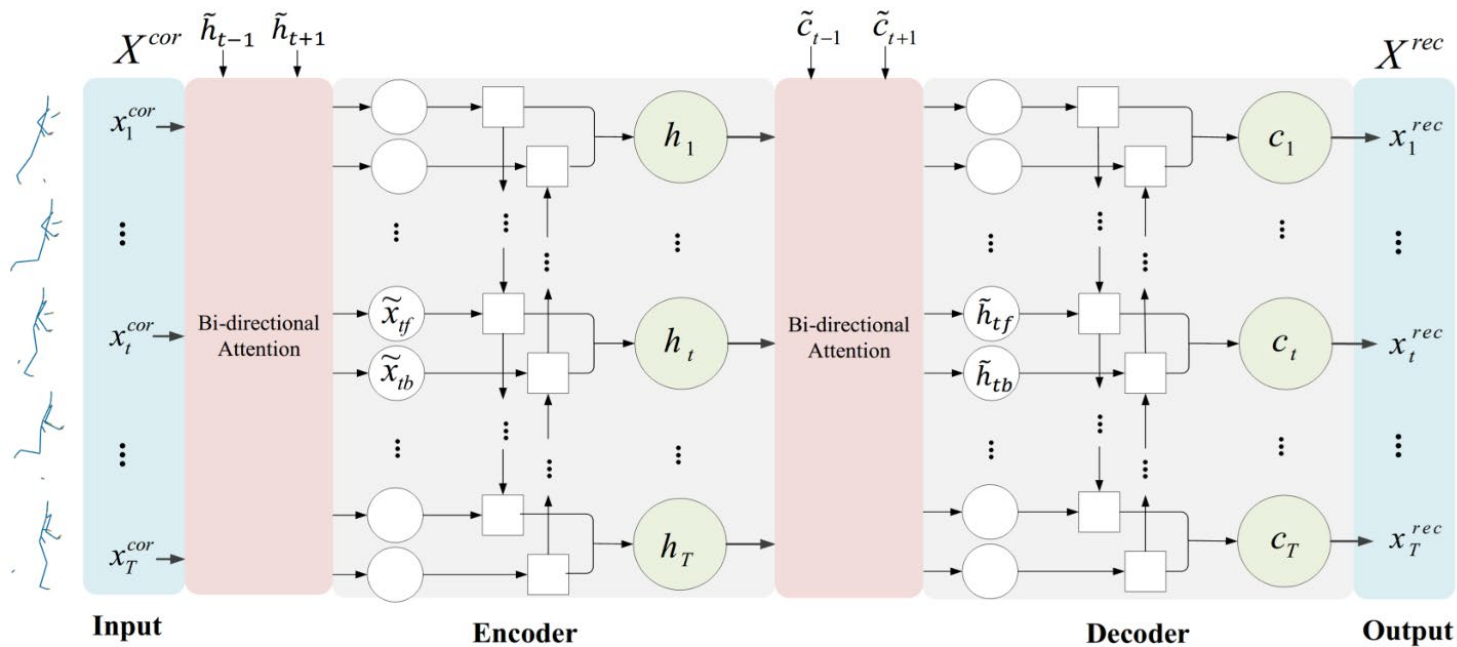


(d) Miss

# Our Goals

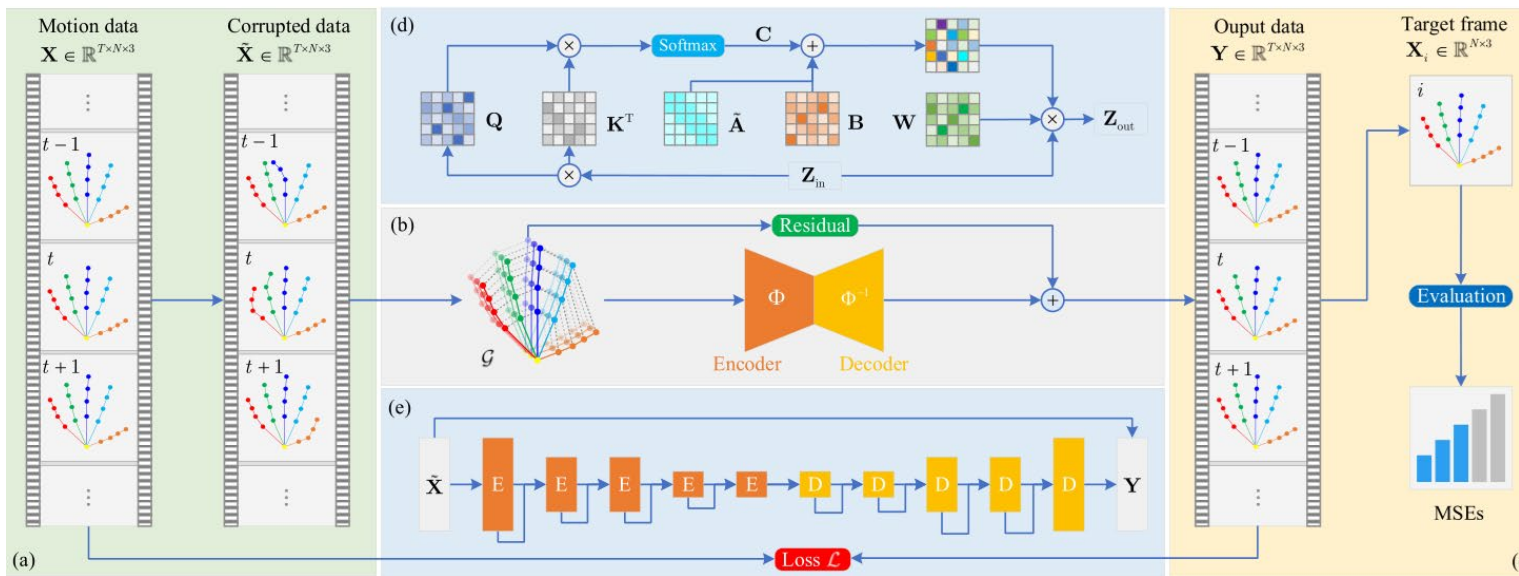


# Related Work



Qiongjie Cui, Huaijiang Sun, Yupeng Li, Yue Kong.. A Deep Bi-directional Attention Network for Human Motion Recovery. IJCAI, 2019

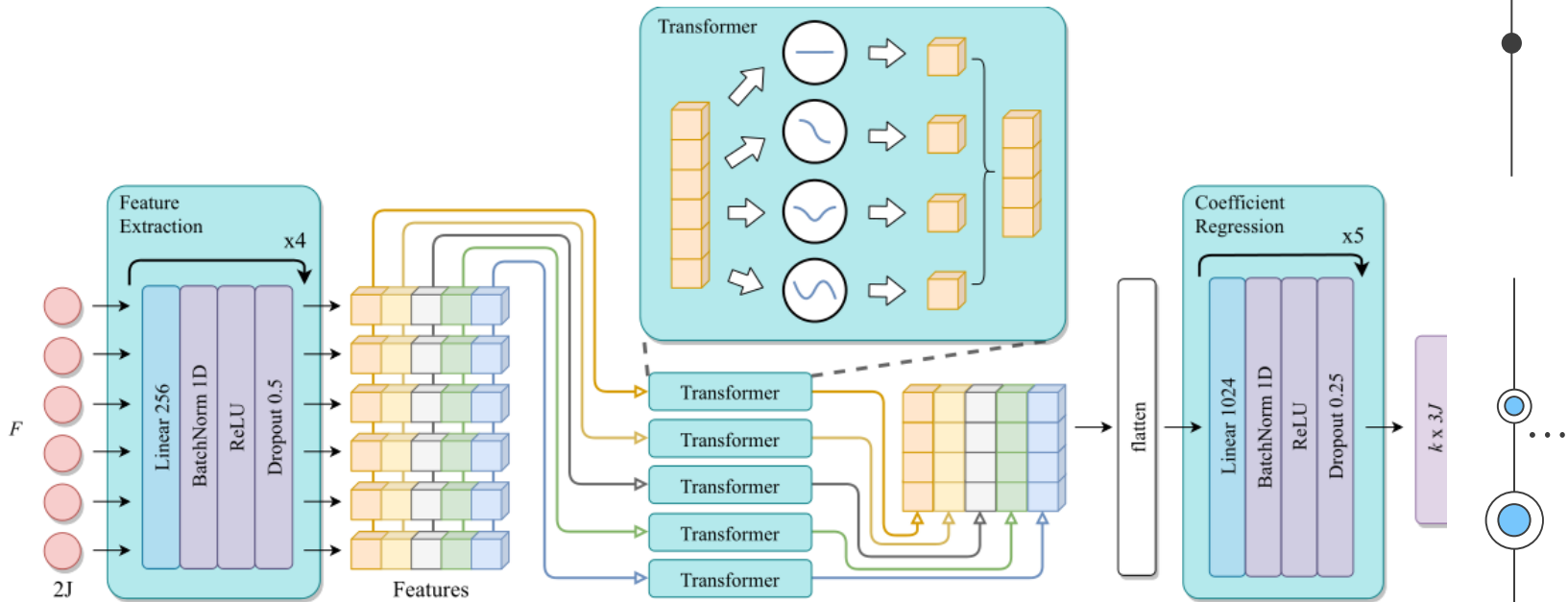
# Related Work

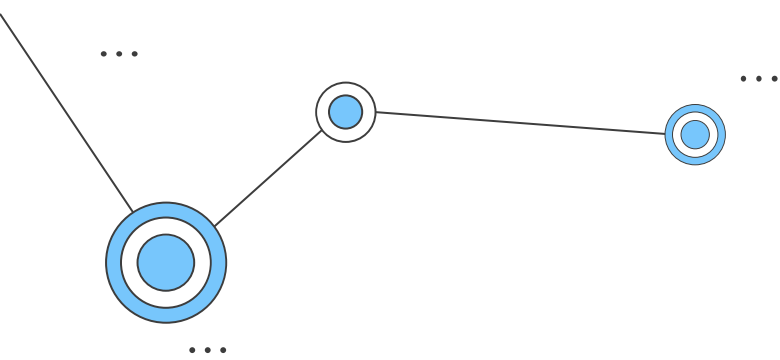


Durham Research, Kanglei Zhou, Zhiyuan Cheng, Hubert P H Shum, Frederick W B Li, Liang, Xiaohui Liang. STGAE: Spatial-Temporal Graph Auto-Encoder for Hand Motion Denoising. ISMAR, 2021.



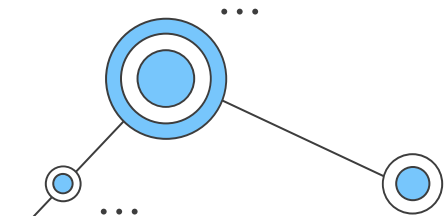
# Related Work





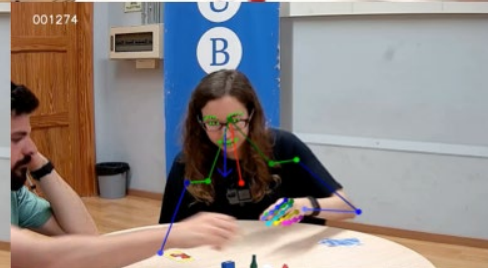
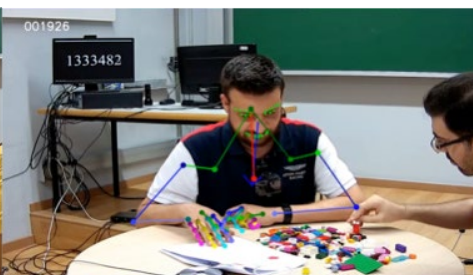
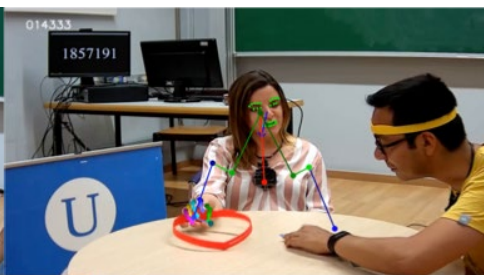
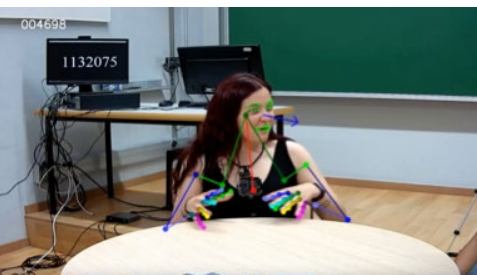
UDIVA

**U**nderstanding  
**D**yadic  
Interactions from  
**V**ideo and  
**A**udio signals



# UDIVA v0.5

- 4 tasks from 145 sessions: 116 for training (raw annotations), 18 for validation and 11 for testing (cleaned annotations).



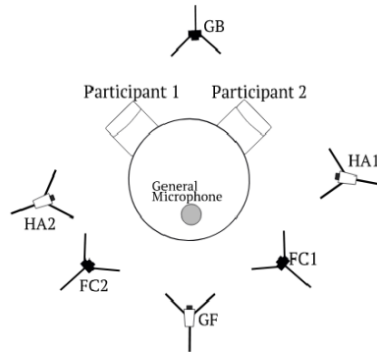
*Talk*

*Animals*

*Lego*

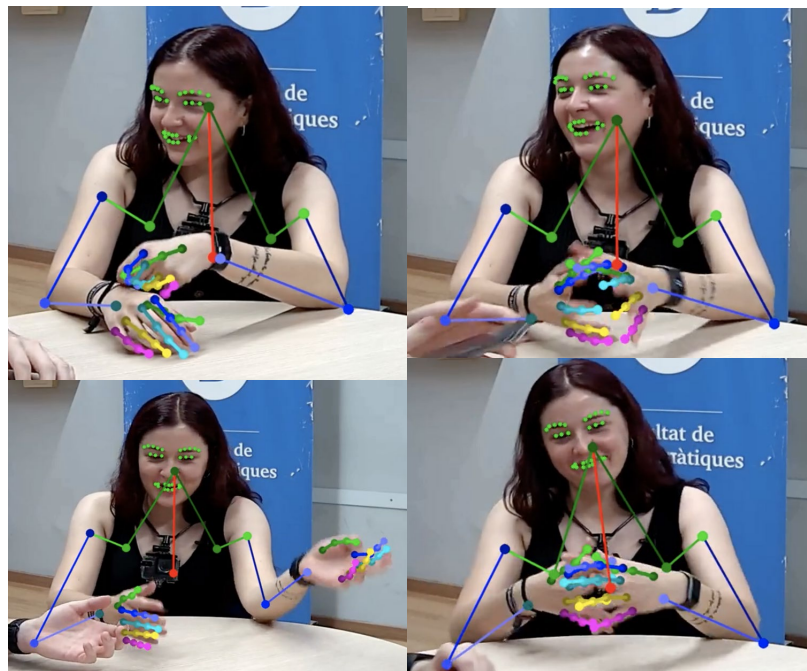
*Ghost*

# UDIVA v0.5



...

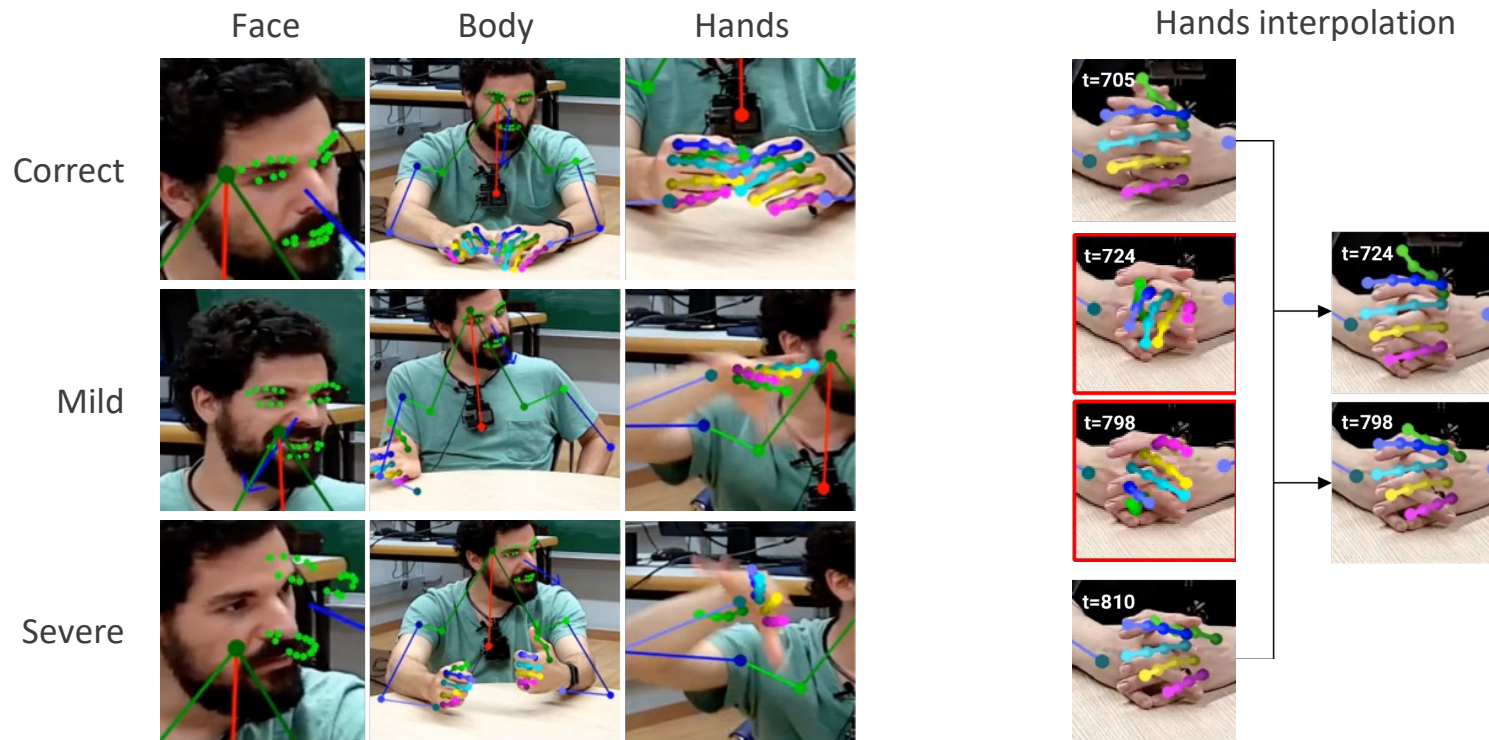
# Landmarks extraction



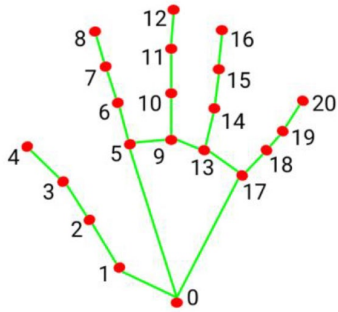
Hand Pose estimation Algorithm:  
**FrankMocap**

# Visual inspection

- Validation (18 sessions) and test (11 sessions) sets underwent a visual inspection process.

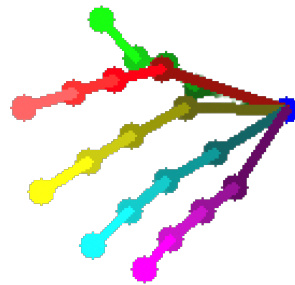


# Hand Representation

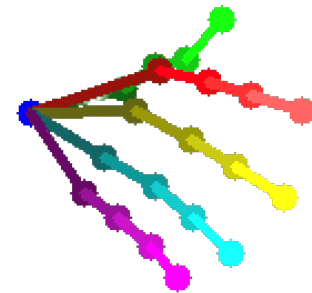


- 1x20 right hand landmarks
- 1x20 left hand landmarks (flip horizontally)
- Absolute coordinates were transformed to root-relative coordinates.

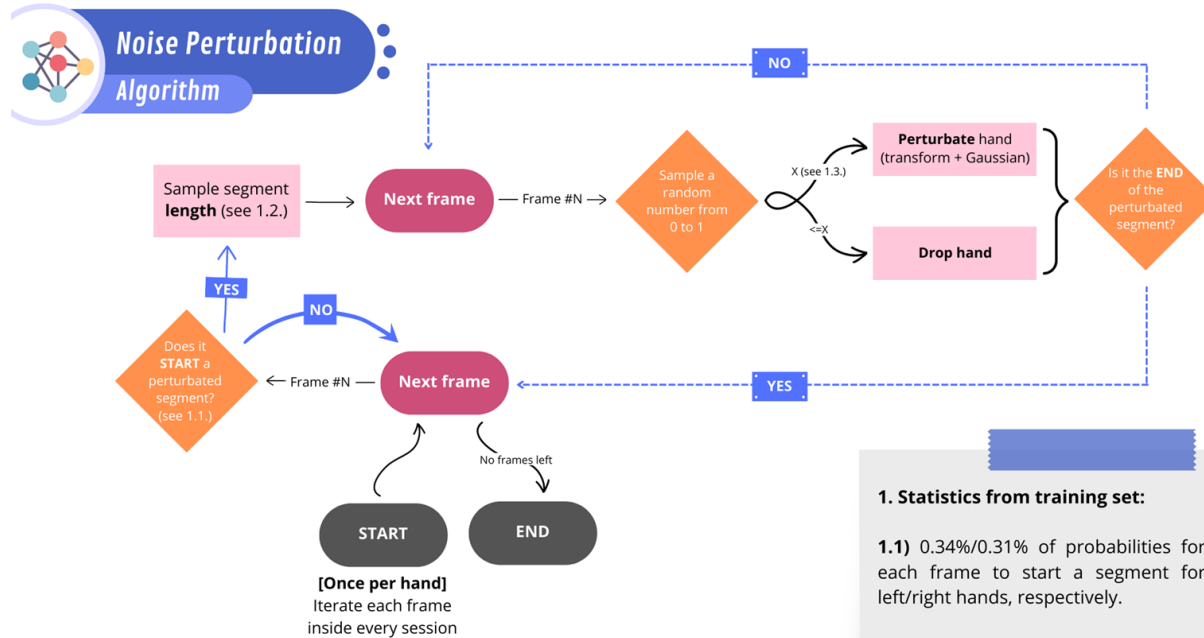
Left Hand



Left Hand Transformed



# Perturbation Algorithm



## 1. Statistics from training set:

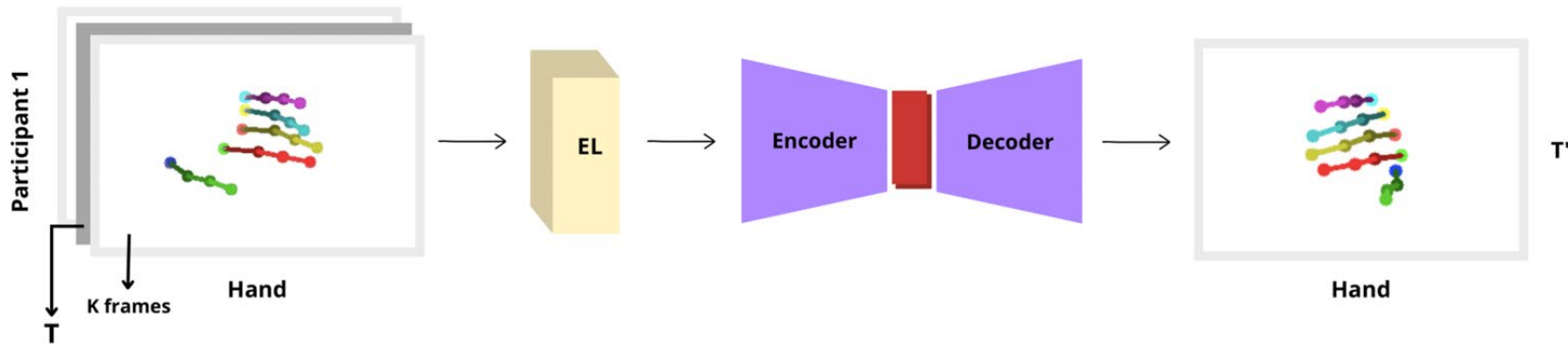
**1.1)** 0,34%/0,31% of probabilities for each frame to start a segment for left/right hands, respectively.

**1.2)** The length of the segments follow a 'lognorm' distribution.

**1.3)** There are 14,4%/8,0% of left/right hands missing inside the segments (joints equal to 0).



# Methodologies



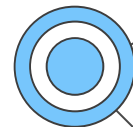
- **Architectures:**

- MLP  $\longleftrightarrow$  MLP Frequential
- Seq2Seq  $\longleftrightarrow$  Seq2Seq-Bidirectional
- Spatio-Temporal Transformer

- **Data:**

- Participant 1
- Two hands consecutively
- Ground Truth: Cleaned Frames

# Visual Results



- MLP

Sample 21 - [H] Frame 0

Context

Raw

Groundtruth

Difference

algorithm\_0



MSE\_AVG: 622.09  
MSE\_ind\_AVG: 1786.13  
MSE\_root\_AVG: 335.56  
Dist\_Euclidean\_AVG: 0.15

MSE: 329.17  
MSE\_ind: 742.25  
MSE\_root: 282.50  
Dist\_Euclidean: 22.27

MSE: 552.05  
MSE\_ind: 1622.81  
MSE\_root: 556.12  
Dist\_Euclidean: 32.00

- MLP Frequential

Sample 31 - [H] Frame 0

Context

Raw

Groundtruth

Difference

algorithm\_0



MSE\_AVG: 677.77  
MSE\_ind\_AVG: 1968.40  
MSE\_root\_AVG: 277.70  
Dist\_Euclidean\_AVG: 0.46

MSE: 319.17  
MSE\_ind: 864.95  
MSE\_root: 102.50  
Dist\_Euclidean: 22.82

MSE: 185.96  
MSE\_ind: 464.31  
MSE\_root: 103.56  
Dist\_Euclidean: 16.89

# Visual Results

- RNN

Sample 2 - [H] Frame 0

Context

Raw

Groundtruth

Difference

algorithm\_0



MSE\_AVG: 356.57  
MSE\_ind\_AVG: 671.89  
MSE\_root\_AVG: 90.99  
Dist\_Euclidean\_AVG: 0.16

MSE: 189.70  
MSE\_ind: 195.30  
MSE\_root: 32.00  
Dist\_Euclidean: 15.62

MSE: 421.59  
MSE\_ind: 1395.90  
MSE\_root: 795.49  
Dist\_Euclidean: 27.90

- Transformer

Sample 36 - [H] Frame 0

Context

Raw

Groundtruth

Difference

algorithm\_0

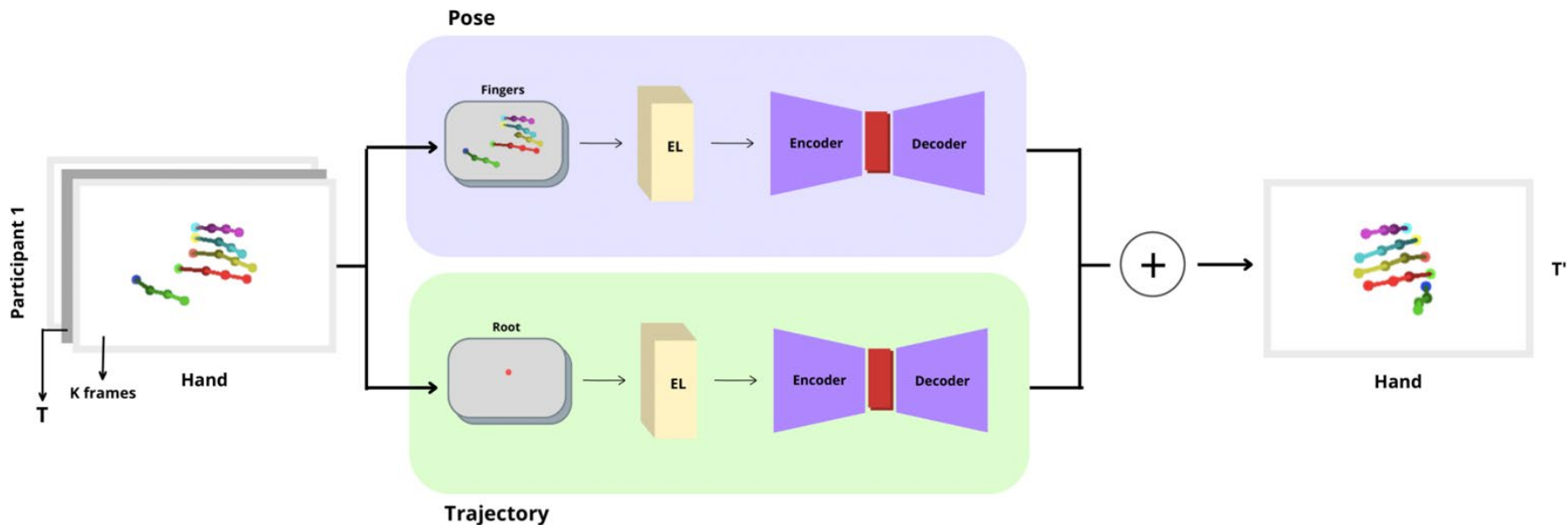


MSE\_AVG: 1127.01  
MSE\_ind\_AVG: 728.35  
MSE\_root\_AVG: 79.86  
Dist\_Euclidean\_AVG: 0.40

MSE: 1116.70  
MSE\_ind: 791.17  
MSE\_root: 74.50  
Dist\_Euclidean: 41.17

MSE: 393.35  
MSE\_ind: 424.97  
MSE\_root: 218.67  
Dist\_Euclidean: 27.41

# Divide and Conquer



- **Pose Architectures:**

- MLP
- MLP Frequential
- Transformer

- **Trajectory Architectures:**

- Seq2Seq
  - LSTM

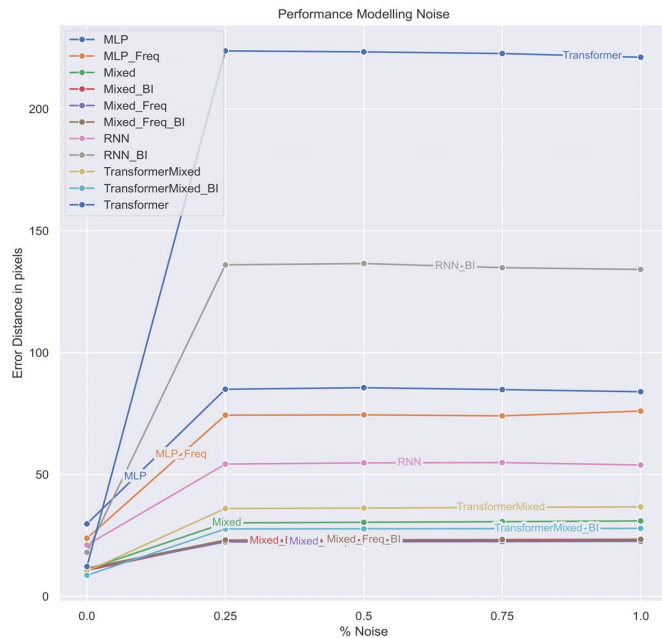
# Ranking

- Baseline: 9.138526

Model	Prediction
TransformerMixed_BI_100	<b>8.649507</b>
TransformerMixed_100	10.361105
Mixed_BI_100	10.661939
Mixed_100	11.080834
Mixed_Freq_100	11.377095
Mixed_Freq_BI_100	11.522652
Transformer_100	12.230602
RNN_BI_100	18.057422
RNN_100	21.049572
MLP_Freq_100	23.856506
MLP_100	29.756484

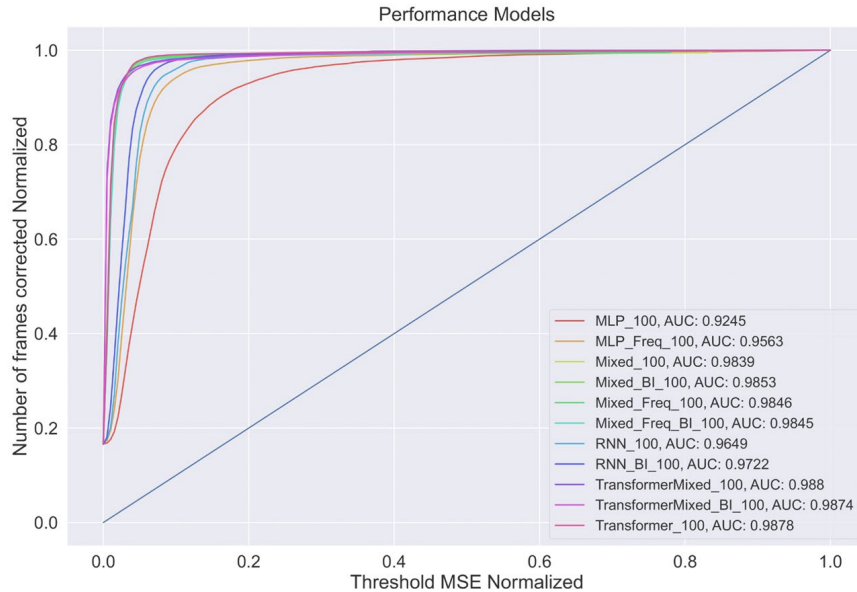
Mean Error of Euclidean Distance: Two hands

# More Evaluation



Performance of the 100 observation length models respect to noise application.

# More Evaluation

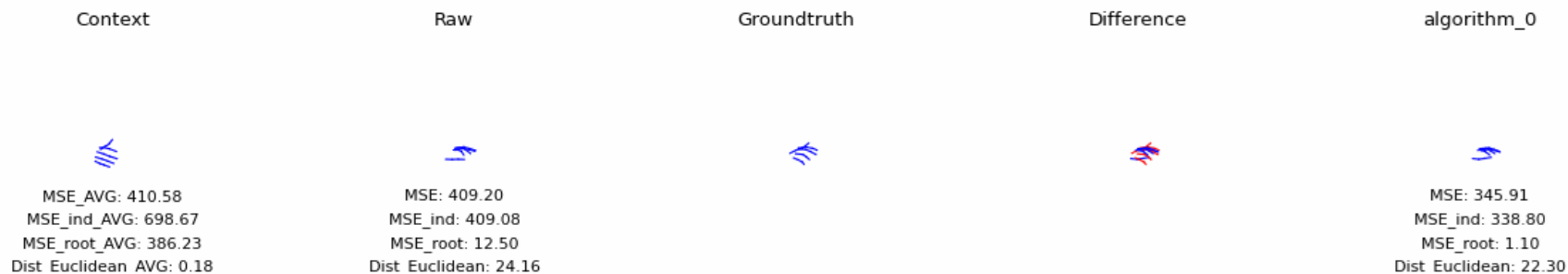


Counting frames with MSE less than the threshold. Normalized values.

# Best Algorithm Visual Result

- Mixed Transformed Bidirectional

Sample 66 - [H] Frame 0




Sample 3 - [H] Frame 0





## First Case

Raw



MSE: 0.00  
MSE\_ind: 0.00  
MSE\_root: 0.00  
Dist\_Euclidean: 0.00


Groundtruth



Difference



algorithm\_0



MSE: 15.61  
MSE\_ind: 6.29  
MSE\_root: 3.40  
Dist\_Euclidean: 5.04

## Second Case

Raw



MSE: 503450.00  
MSE\_ind: 463797.50  
MSE\_root: 495378.00  
Dist\_Euclidean: 1003.06

Groundtruth



Difference



algorithm\_0



MSE: 24067.67  
MSE\_ind: 7004.61  
MSE\_root: 10887.46  
Dist\_Euclidean: 217.08

## Third Case

Raw



MSE: 167.95  
MSE\_ind: 363.33  
MSE\_root: 32.50  
Dist\_Euclidean: 16.35

Groundtruth



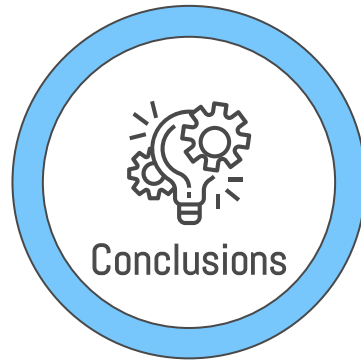
Difference



algorithm\_0



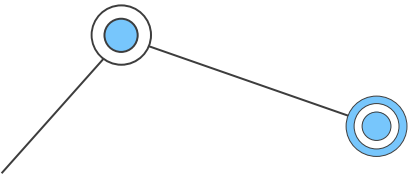
MSE: 36.89  
MSE\_ind: 64.71  
MSE\_root: 5.68  
Dist\_Euclidean: 7.78



- We carefully reviewed all related state-of-the-art literature on pose denoising.
- **The innovative Data Perturbation** algorithm allows to worsen the hand poses.
- The frequency space **improves** results in simple models over the original space, but not is representative.
- Mixed model **are more robust** applying more noise.
- **Pose and Trajectory disentanglement** is a good solution approach to find a robust denoising architecture algorithms.

# Future Work

- How the observation window affects model performance.
- Not add noisy into cleaned frames.
- New recent Deep Learning Models.
- Use frequency domain only in trajectory.
- Instead of denoising a unique pose, denoising a sequence of poses.



# Thanks!

Do you have any questions?

