2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

# Research Internship:
# 3D Pose Estimation for Driver Monitoring

## Victoria BRAMI

Master Mathématiques Vision Apprentissage (MVA)

*victoria.brami@eleves.enpc.fr*

Supervised by Patrick Pérez, advised by Souhaiel Khalfaoui and Renaud Marlet

Thursday September 29$^{th}$ 2022

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Context of work

- Distraction accounts for **20%** of car accidents in 2020.[1]
- **Driver Monitoring System (DMS)**: Set of equipment tools developed around the driver to ease his way of driving.
- EU Comission: new regulations on DMS to be introduced by 2024.

$\rightarrow$ Necessity to Improve existing systems.

École des Ponts
ParisTech

valeo.**ai**

[1]Report made by the European Comission in 2020.

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Context of work

**Motivations:**

> Get knowledge of in-car occupation to understand the occupants' behaviour while driving.

> Supply the best IMS possible (security, confort, etc.)

**Our Goal:**

> Propose **a real-time 3D Pose Estimation of the driver** to be capable to analyse his activities in a second phase.

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Outline

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Interior Monitoring Datasets Constraints

| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
|---|---|---|---|---|---|---|---|
| Scene Type | Real | Sitting, | In-Cabin | In-Cabin | Real | Real | Real |
| | Condition | Driving-like | Driving | Driving | Condition | Condition | Condition |
| Occupants | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only |
| Views | 6 | 1 | 2 | 1 | 3 | 2 | >1 |
| Nb. frames | >9.6M | 250k | 1.1M / 315k(view 1) | 119.7k / 3.3k | 4.4M | 2.1M | 10k |
| Nb. videos | 29 | 110 | 21 | | | 386 | - |
| RGB/Gray | ✓ | ✓ | ✓ | ✓ | ✓ | - | ✓ |
| IR | ✓ | - | ✓ | ✓ (6.7k) | ✓ | ✓ | - |
| Depth | ✓ | ✓ | ✓ | ✓ (6.7k) | ✓ | ✓ | - |
| Subjects[a] | 15 (4/11) | 22 (10/12) | 21 (10/11) | 13 (N/A) | 37 (10/27) | 31 (N/A) | 19 (7/12) |
| **Annotations Contents** | | | | | | | |
| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
| Activity | ✓ | ✓ | - | ✓ | ✓ | - | - |
| Nb. Activ. | 83 | 20 | - | 20 | 13 | - | - |
| 2D joints | ✓ | ✓ | - | ✓ | - | - | ✓ |
| 3D joints | ✓ | ✓ | ✓ | - | N/A | - | - |
| Format | COCO 17 | 17 Upper | Head center | COCO 17 | - | - | COCO17 |

Table: Main large-scale Driver Monitoring datasets

[a](F/M) for female / male

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Interior Monitoring Datasets Constraints

| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
|---|---|---|---|---|---|---|---|
| Scene Type | Real Condition | Sitting, Driving-like | In-Cabin Driving | In-Cabin Driving | Real Condition | Real Condition | Real Condition |
| Occupants | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only |
| Views | 6 | 1 | 2 | 1 | 3 | 2 | >1 |
| Nb. frames | >9.6M | 250k | 1.1M / 315k(view 1) | 119.7k / 3.3k | 4.4M | 2.1M | 10k |
| Nb. videos | 29 | 110 | 21 | | | 386 | - |
| RGB/Gray | ✓ | ✓ | ✓ | ✓ | ✓ | - | ✓ |
| IR | ✓ | - | ✓ | ✓ (6.7k) | ✓ | ✓ | - |
| Depth | ✓ | ✓ | ✓ | ✓ (6.7k) | ✓ | ✓ | - |
| Subjects[a] | 15 (4/11) | 22 (10/12) | 21 (10/11) | 13 (N/A) | 37 (10/27) | 31 (N/A) | 19 (7/12) |
| **Annotations Contents** | | | | | | | |
| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
| Activity | ✓ | ✓ | - | ✓ | ✓ | - | - |
| Nb. Activ. | 83 | 20 | - | 20 | 13 | - | - |
| 2D joints | ✓ | ✓ | - | ✓ | - | - | ✓ |
| 3D joints | ✓ | ✓ | ✓ | - | N/A | - | - |
| Format | COCO 17 | 17 Upper | Head center | COCO 17 | - | - | COCO17 |

Table: Main large-scale Driver Monitoring datasets

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

# Interior Monitoring Datasets Constraints

| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
|---|---|---|---|---|---|---|---|
| Scene Type | Real Condition | Sitting, Driving-like | In-Cabin Driving | In-Cabin Driving | Real Condition | Real Condition | Real Condition |
| Occupants | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only |
| Views | 6 | 1 | 2 | 1 | 3 | 2 | >1 |
| Nb. frames | >9.6M | 250k | 1.1M / 315k(view 1) | 119.7k / 3.3k | 4.4M | 2.1M | 10k |
| Nb. videos | 29 | 110 | 21 | | | 386 | - |
| RGB/Gray | ✓ | ✓ | ✓ | ✓ | ✓ | - | ✓ |
| IR | ✓ | - | ✓ | ✓(6.7k) | ✓ | ✓ | - |
| Depth | ✓ | ✓ | ✓ | ✓(6.7k) | ✓ | ✓ | - |
| Subjects[a] | 15 (4/11) | 22 (10/12) | 21 (10/11) | 13 (N/A) | 37 (10/27) | 31 (N/A) | 19 (7/12) |
| **Annotations Contents** | | | | | | | |
| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
| Activity | ✓ | ✓ | - | ✓ | ✓ | - | - |
| Nb. Activ. | 83 | 20 | - | 20 | 13 | - | - |
| 2D joints | ✓ | ✓ | - | ✓ | - | - | ✓ |
| 3D joints | ✓ | ✓ | ✓ | - | N/A | - | - |
| Format | COCO 17 | 17 Upper | Head center | COCO 17 | - | - | COCO17 |

Table: Main large-scale Driver Monitoring datasets

[a](F/M) for female / male

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

# Interior Monitoring Datasets Constraints

| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
|---------|---------|------------|-------------|----------|--------|--------|----------|
| Scene Type | Real | Sitting, | In-Cabin | In-Cabin | Real | Real | Real |
| | Condition | Driving-like | Driving | Driving | Condition | Condition | Condition |
| Occupants | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only |
| Views | 6 | 1 | 2 | 1 | 3 | 2 | >1 |
| Nb. frames | >9.6M | 250k | 1.1M / 315k(view 1) | 119.7k / 3.3k | 4.4M | 2.1M | 10k |
| Nb. videos | 29 | 110 | 21 | | | 386 | - |
| RGB/Gray | ✓ | ✓ | ✓ | ✓ | ✓ | - | ✓ |
| IR | ✓ | - | ✓ | ✓(6.7k) | ✓ | ✓ | - |
| Depth | ✓ | ✓ | ✓ | ✓(6.7k) | ✓ | ✓ | - |
| Subjects[a] | 15 (4/11) | 22 (10/12) | 21 (10/11) | 13 (N/A) | 37 (10/27) | 31 (N/A) | 19 (7/12) |
| **Annotations Contents** | | | | | | | |
| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
| Activity | ✓ | ✓ | - | ✓ | ✓ | - | - |
| Nb. Activ. | 83 | 20 | - | 20 | 13 | - | - |
| 2D joints | ✓ | ✓ | - | ✓ | - | - | ✓ |
| 3D joints | ✓ | ✓ | ✓ | - | N/A | - | - |
| Format | COCO 17 | 17 Upper | Head center | COCO 17 | - | - | COCO17 |

Table: Main large-scale Driver Monitoring datasets

[a](F/M) for female / male

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Interior Monitoring Datasets Constraints

| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
|---|---|---|---|---|---|---|---|
| Scene Type | Real Condition | Sitting, Driving-like | In-Cabin Driving | In-Cabin Driving | Real Condition | Real Condition | Real Condition |
| Occupants | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only | Driver Only |
| Views | 6 | 1 | 2 | 1 | 3 | 2 | >1 |
| Nb. frames | >9.6M | 250k | 1.1M / 315k(view 1) | 119.7k / 3.3k | 4.4M | 2.1M | 10k |
| Nb. videos | 29 | 110 | 21 | | | 386 | - |
| RGB/Gray | ✓ | ✓ | ✓ | ✓ | ✓ | - | ✓ |
| IR | ✓ | - | ✓ | ✓ (6.7k) | ✓ | ✓ | - |
| Depth | ✓ | ✓ | ✓ | ✓ (6.7k) | ✓ | ✓ | - |
| Subjects[a] | 15 (4/11) | 22 (10/12) | 21 (10/11) | 13 (N/A) | 37 (10/27) | 31 (N/A) | 19 (7/12) |
| **Annotations Contents** | | | | | | | |
| Dataset | D&A [5] | Pandora[1] | AutoPOSE[7] | TICaM[3] | DMD[6] | DAD[4] | DriPE[2] |
| Activity | ✓ | ✓ | - | ✓ | ✓ | - | - |
| Nb. Activ. | 83 | 20 | - | 20 | 13 | - | - |
| 2D joints | ✓ | ✓ | - | ✓ | - | - | ✓ |
| 3D joints | ✓ | ✓ | ✓ | - | N/A | - | - |
| Format | COCO 17 | 17 Upper | Head center | COCO 17 | - | - | COCO17 |

Table: Main large-scale Driver Monitoring datasets

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

## Drive And Act Dataset Format

- 6 views.
- 15 drivers filmed 20-30 min each (10 / 2 / 3).
- **9.6 Million** frames.
- Annotations triangulated from **OpenPose**[2]

(a) COCO17 annotation format[4]

(b) Sample from Drive&Act[5]

---

[2]Cao & al., OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, in *TPAMI*, 2019.

[4]Lin & al., Microsoft COCO: Common objects in context, in *ECCV*, 2014.

[5]Martin & al., Drive&Act: A Multi-modal Dataset for Fine-grained Driver Behavior Recognition in Autonomous Vehicles, in *ICCV*, 2019.

École des Ponts
ParisTech

valeo.**ai**

**2D Pose Estimation**
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
Experiments

## Outline

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
Experiments

# 2D Pose Models Fields

## In HR Net (Top Down)



Figure: Heatmap

## In OpenPifPaf (Bottom-Up)

(a) CIF



(b) CAF



valeo.ai

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
Experiments

## Top Down Model: HR-Net (2019)



Figure: HR-Net Model Architecture[6]

École des Ponts
ParisTech
[6]Sun & al., Deep High-Resolution Representation Learning for Human Pose
Estimation, in *CVPR*, 2019.

valeo.ai

9 / 36

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
Experiments

# Bottom-Up Model: OpenPifPaf (2019-2021)



Figure: OpenPifPaf Model Architecture[7]

[7]Kreiss & al., OpenPifPaf: Composite fields for semantic keypoint detection and spatio-temporal association, in *TITS*, 2021.

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

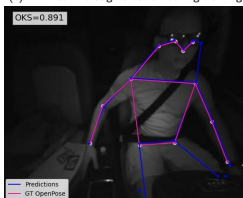Studied models
**Experiments**

## 2D Pose Models Finetuning

**Our Framework**:

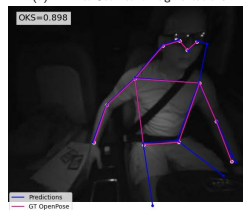- Dataset: Drive & Act.
- Metrics: AP (↑) and AR (↑).
- Finetune on **30 epochs**.
- Augmentations: scale, noise, blur.
- Specificity: Apply a **binary mask** on the joints loss **discard Feet Pose predictions**.

École des Ponts
ParisTech

valeo.**ai**

**2D Pose Estimation**
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
**Experiments**

# 2D Pose: Visual Results



Table: Visualization of the retrained models on Drive & Act test set.

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
Experiments

## 2D Pose: Quantitative Results

| HR Net | Input | AP | AP50 | AP75 | AR | AR50 | AR75 |
|---|---|---|---|---|---|---|---|
| No Finetuning | 256 x 192 | 85.0 | 96.5 | 90.2 | 90.9 | 98.7 | 93.7 |
| Finetuned (no aug.) | 256 x 192 | 87.0 | 98.1 | 90.8 | 90.3 | 98.7 | 93.9 |
| Finetuned (with geom. aug.) | 256 x 192 | 90.1 | **99.0** | **94.2** | **93.7** | 99.4 | **96.0** |
| Finetuned (with geom. aug. + noise + blur) | 256 x 192 | **90.4** | 98.6 | 92.2 | 91.2 | **99.5** | 94.2 |
| **OpenPifPaf** | **Input** | **AP** | **AP50** | **AP75** | **AR** | **AR50** | **AR75** |
| Finetuned (with geom. aug.) | 256 x 192 | 84.0 | 93.6 | 87.0 | 88.1 | 93.8 | 90.7 |

Table: AP and AR on Drive & Act test set

→ HR-Net finetuned with more augmentations **outperforms** OpenPifPaf.

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Studied models
**Experiments**

## 2D Pose Results Analysis

### Pros HR Net

1. **Better scores** obtained on Drive And Act as the model's size is $2.5\times$ bigger.

2. **More keypoints** estimated.

### Pros OpenPifPaf

1. **No need** of a prior detection step.

2. Inference time is **much lower** (almost a requirement for embedded systems).

3. **More stability and consistency** across consecutive frames.

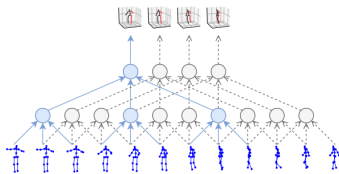- **Conclusion:** Keep working with OpenPifPaf.

valeo.**ai**

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
3D Pose Lifting Experiments

## Outline

valeo.**ai**

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
3D Pose Lifting Experiments

# 3D Pose Lifting with a CNN Model

**Idea:** From a sequence of 2D consecutive skeleton, predicts the 3D pose of the middle frame.

On Human3.6M[8]: **Mean Error is 37.2mm.**



(a) VideoPose3D[9] Model

(b) Causal Form of VideoPose3D

---

[8]Ionescu & al., Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments, *TPAMI*, 2014.

[9]Pavllo & al., 3D human pose estimation in video with temporal convolutions and semi-supervised training, in *CVPR*, 2019.

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

# 3D Pose Lifting with a CNN Model



Figure: Adaptation of VideoPose3D with the addition of joints' confidence scores $\hat{c}_i$ in input.

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

## 3D Pose Lifting with a CNN Model

| Blocks | kernel Size length | Input frames | MPJPE($\downarrow$) (mm) | P-MPJPE($\downarrow$) (mm) | N-MPJPE($\downarrow$) (mm) | MPJVE ($\downarrow$) (mm.s$^{-1}$) |
|--------|--------------------|--------------|--------------------------|-----------------------------|------------------------------|-----------------------------------|
| $B = 1$ | $K = (3, 3)$ | 9 | $34.9_{\pm 0.3}$ | $23.2_{\pm 0.1}$ | $\underline{27.5_{\pm 0.3}}$ | $6.7_{\pm 0.01}$ |
| $B = 2$ | $K = (3, 3, 3)$ | 27 | $34.6_{\pm 0.5}$ | $22.8_{\pm 0.1}$ | $28.0_{\pm 0.3}$ | $6.63_{\pm 0.03}$ |
| $B = 3$ | $K = (3, 3, 3, 3)$ | 81 | $33.5_{\pm 0.4}$ | $22.8_{\pm 0.2}$ | $27.9_{\pm 0.4}$ | $6.59_{\pm 0.02}$ |
| $B = 4$ | $K = (3, 3, 3, 3, 3)$ | 243 | $\underline{33.3_{\pm 0.3}}$ | $\underline{22.6_{\pm 0.1}}$ | $27.6_{\pm 0.3}$ | $\underline{6.55_{\pm 0.01}}$ |

Table: VideoPose3D predictions on *Drive&Act* test set. with different architectures.

$\rightarrow$ No major difference with bigger architecture.

valeo.**ai**

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

# 3D Pose Lifting with a CNN Model



(a) Angles Constraint    (b) Symmetry Constraint

Figure: Kinematics Constraints added

$$\mathcal{L}_{\text{sym}}(\hat{p}) = \sum_{((i,j),(k,l)) \in M} (\|\hat{p}_i - \hat{p}_j\|_2 - \|\hat{p}_k - \hat{p}_l\|_2)^2 \qquad (1)$$

$$\mathcal{L}_{\text{illegal}}(\hat{p}) = \exp\left(-\min(\overrightarrow{n_s^r} \cdot \overrightarrow{v_{we}}, 0)\right) \qquad (2)$$

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

# 3D Pose Lifting with a CNN Model

| Model | MPJPE ($\downarrow$) | P-MPJPE ($\downarrow$) |
|---|---|---|
| $\lambda_{sym} = 0.$ | $34.6_{\pm 0.5}$ | $22.8_{\pm 0.1}$ |
| $\lambda_{sym} = 1.10^{-4}$ | $33.9_{\pm 0.4}$ | $\underline{22.6_{\pm 0.2}}$ |
| $\lambda_{sym} = 1.10^{-3}$ | $34.5_{\pm 0.3}$ | $23.0_{\pm 0.1}$ |
| $\lambda_{sym} = 1.10^{-2}$ | $34.9_{\pm 0.4}$ | $22.9_{\pm 0.1}$ |
| $\lambda_{sym} = 1.10^{-1}$ | $\underline{33.5_{\pm 0.4}}$ | $23.8_{\pm 0.1}$ |
| $\lambda_{sym} = 1.10^{0}$ | $50.0_{\pm 0.6}$ | $43.7_{\pm 0.3}$ |
| $\lambda_{sym} = 1.10^{1}$ | $111.0_{\pm 1.8}$ | $93.1_{\pm 2.1}$ |
| $\lambda_{sym} = 1.10^{2}$ | $194.4_{\pm 19.9}$ | $156.9_{\pm 20.4}$ |

Table: Results when training with various weighted symmetry loss.

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

## 3D Pose Lifting with a CNN Model

| Model | MPJPE ($\downarrow$) | P-MPJPE ($\downarrow$) |
|---|---|---|
| $\lambda_a = 0.$ | $34.6_{\pm 0.5}$ | $\underline{22.8_{\pm 0.1}}$ |
| $\lambda_a = 1.10^{-3}$ | $34.5_{\pm 0.4}$ | $22.7_{\pm 0.1}$ |
| $\lambda_a = 1.10^{-2}$ | $34.3_{\pm 0.7}$ | $22.8_{\pm 0.3}$ |
| $\lambda_a = 1.10^{-1}$ | $34.7_{\pm 0.4}$ | $22.9_{\pm 0.0}$ |
| $\lambda_a = 1.10^{0}$ | $34.3_{\pm 0.3}$ | $23.0_{\pm 0.1}$ |
| $\lambda_a = 1.10^{1}$ | $\underline{34.3_{\pm 0.4}}$ | $23.6_{\pm 0.3}$ |
| $\lambda_a = 1.10^{2}$ | $35.3_{\pm 0.8}$ | $25.0_{\pm 1.0}$ |
| $\lambda_a = 1.10^{3}$ | $44.2_{\pm 1.9}$ | $32.4_{\pm 1.4}$ |

Table: Results when training with various weighted angle loss.

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

# 3D Pose Lifting Qualitative results



Figure: 3D Pose Prediction on Drive & Act test set.

2D Pose Estimation
**2D to 3D Pose Lifting**
Extension of the pipeline to Face and body Pose
Conclusions and future work

3D Pose Lifting Model
**3D Pose Lifting Experiments**

## 3D Pose Lifting with a CNN Model

**Conclusions:**

- CNN-based VideoPose3D lifter works well with a **Mean Error around 34.0mm.**
- Study self-supervised approaches.
- Look for lighter models using transformers like **P-STMO**[10].

[10]Shan & al., P-STMO: Pre-trained spatial temporal many-to-one model for 3d human pose estimation. in *ECCV*, 2022.

2D Pose Estimation
2D to 3D Pose Lifting
**Extension of the pipeline to Face and body Pose**
Conclusions and future work

Principles
First Experiments
Occlusion Experiments

## Outline

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
**Extension of the pipeline to Face and body Pose**
Conclusions and future work

**Principles**
First Experiments
Occlusion Experiments

## Dataset Pseudo Annotation

Motivation: Face Landmarks pose give better interpretability of the driver's state.

Goal: Incorporate the 3D face landmarks estimation.

Means: Use a pretrained network to estimate the 3D facial landmarks: **3DDFA v2 model**[11].



Refined Body Pose representation with $17 + 68$ joints.[12]

[11]Guo & al., Towards fast, accurate and stable 3D dense face alignment, in *ECCV*, 2020.

[12]Jin & al., Whole-body human pose estimation in the wild, in *ECCV*, 2020.

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Principles
First Experiments
Occlusion Experiments

# Dataset Pseudo-Labelling



68 landmarks from COCO

3DDFA V2

Predictions of the 3D face pose

Rescale Distance

3D Alignment

2D Re-projection

Figure: Face Alignment Protocol.

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Principles
First Experiments
Occlusion Experiments
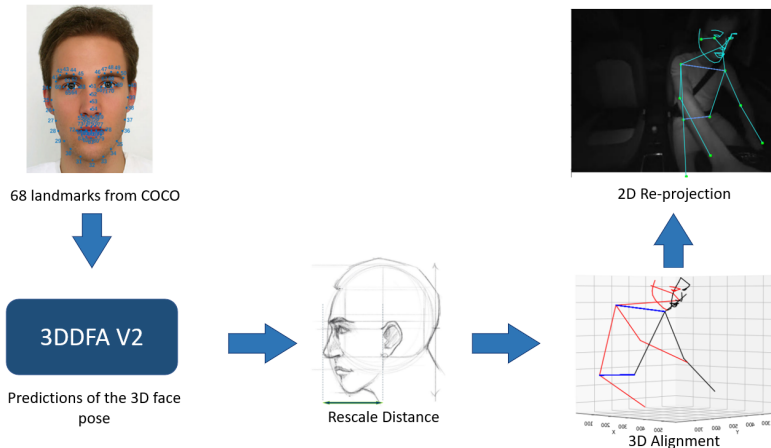
## Protocol Applied on Wholebody

**Training Framework:**

- **Input sequence:** 27 frames of $17 + 68$-joints skeletons.
- **Architecture:** 2 Blocks of Causal Convolutions with 3 dilations.
- Train on **100** epochs.
- Loss and Metric: **Mean Per Joint Error Loss**.
- Learning Rate and Batch size: $1.10^{-3}$ and 1024.
- Add some Dropout : **0.25**.

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
**Extension of the pipeline to Face and body Pose**
Conclusions and future work

Principles
**First Experiments**
Occlusion Experiments

# 3D Pose Lifting Results

Click for video

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Principles
First Experiments
Occlusion Experiments

## Occlusions Experiments



(a) VideoPose3D initial Input      (b) Added occlusions in VideoPose3D Input

Figure: Experiments on VideoPose3D's robustness, adding occlusions in the training to facilitate domain adaptation.

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
**Extension of the pipeline to Face and body Pose**
Conclusions and future work

Principles
First Experiments
**Occlusion Experiments**

# 3D Pose Lifting Comparison Results

Click for video

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
Conclusions and future work

Principles
First Experiments
Occlusion Experiments

## Occlusions: 3D Pose Lifting Results

| Model | Input frames | Occlusions ratio (%) | MPJPE($\downarrow$) (mm) | P-MPJPE($\downarrow$) (mm) | N-MPJPE($\downarrow$) (mm) | MPJVE ($\downarrow$) (mm.s$^{-1}$) |
|---|---|---|---|---|---|---|
| VideoPose3D | 27 | 0 % | $\underline{39.4}_{\pm 0.8}$ | $\underline{13.4}_{\pm 0.2}$ | $23.8_{\pm 0.2}$ | $7.28_{\pm 0.02}$ |
| | | 5 % | $39.9_{\pm 0.7}$ | $14.3_{\pm 0.5}$ | $\underline{23.6}_{\pm 0.6}$ | $7.48_{\pm 0.05}$ |
| | | 10 % | $40.7_{\pm 1.6}$ | $14.9_{\pm 0.2}$ | $24.0_{\pm 1.0}$ | $7.63_{\pm 0.02}$ |
| | | 20 % | $41.2_{\pm 0.6}$ | $15.9_{\pm 0.1}$ | $24.9_{\pm 0.5}$ | $7.88_{\pm 0.08}$ |
| | | 30 % | $42.5_{\pm 0.7}$ | $16.3_{\pm 0.1}$ | $26.8_{\pm 0.2}$ | $8.07_{\pm 0.05}$ |
| | | 40 % | $41.8_{\pm 0.3}$ | $16.7_{\pm 0.2}$ | $27.0_{\pm 0.9}$ | $8.19_{\pm 0.10}$ |
| VideoPose3D | 243 | 0 % | $37.4_{\pm 0.6}$ | $12.6_{\pm 0.4}$ | $18.2_{\pm 0.7}$ | $5.88_{\pm 0.09}$ |
| | | 5 % | $43.1_{\pm 0.3}$ | $14.2_{\pm 0.3}$ | $25.0_{\pm 0.8}$ | $\underline{6.86}_{\pm 0.08}$ |
| | | 10 % | $44.8_{\pm 0.5}$ | $16.4_{\pm 0.1}$ | $26.0_{\pm 0.7}$ | $7.34_{\pm 0.06}$ |
| | | 20 % | $43.4_{\pm 0.3}$ | $16.6_{\pm 0.2}$ | $26.6_{\pm 0.3}$ | $7.64_{\pm 0.03}$ |
| | | 30 % | $46.5_{\pm 0.5}$ | $17.5_{\pm 0.1}$ | $28.0_{\pm 0.3}$ | $7.81 \pm 0.06$ |
| | | 80 % | $75.5_{\pm 4.4}$ | $34.0_{\pm 1.2}$ | $50.9_{\pm 2.7}$ | $8.44_{\pm 0.01}$ |

Table: Errors obtained when incorporating occlusions

École des Ponts
ParisTech

valeo.ai

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
**Conclusions and future work**

Discussions
Summary
Perspectives

## Outline

1. 2D Pose Estimation

2. 2D to 3D Pose Lifting

3. Extension of the pipeline to Face and body Pose

4. Conclusions and future work
   - Discussions
   - Summary
   - Perspectives

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
**Conclusions and future work**

**Discussions**
Summary
Perspectives

## Limits of the method

- Work restricted on a single dataset.
- No real Ground Truth: Data is **pseudo-labelled** by OpenPose[13].
- Intented to minimize the errors by running each evaluation **5 times**.

[13]Cao & al., OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, in *TPAMI*, 2019.

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
**Conclusions and future work**

Discussions
**Summary**
Perspectives

## Conclusion

### Our Contributions:

- **Survey and exhaustive comparison** of Interior Monitoring datasets.
- Pseudo annotation and **3D face alignment** over Drive And Act dataset.
- **End-to-end framework** for Driver's 3D body and face landmarks pose estimation.
- Average error in 3D Pose Estimation at **34mm on average**.

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
**Conclusions and future work**

Discussions
Summary
**Perspectives**

## Perspectives

### Short term:

1. Extend the Pipeline with the **addition of 3D Hands Pseudo annotations**.
2. Smooth the pose estimation over consecutive frames.
3. Study **self-supervised methods deeper** to discard the lack of data issue.
4. Evaluate our framework on other datasets: Valeo collecting the data.

### Long term:

1. **Lighten the model** to make it embeddable.
2. Activity Recognition based on sequence of 3D Pose Estimations.

École des Ponts
ParisTech

valeo.**ai**

2D Pose Estimation
2D to 3D Pose Lifting
Extension of the pipeline to Face and body Pose
**Conclusions and future work**

Discussions
Summary
**Perspectives**

# Thank you for your Attention !

# References I

📄 Guido Borghi, Marco Venturelli, Roberto Vezzani, and Rita Cucchiara.
Poseidon: Face-from-depth for driver pose estimation.
In *CVPR*, 2017.

📄 Romain Guesdon, Carlos Crispim-Junior, and Laure Tougne.
Dripe: A dataset for human pose estimation in real-world driving settings.
In *ICCV*, 2021.

📄 Jigyasa Singh Katrolia, Bruno Mirbach, Ahmed El-Sherif, Hartmut Feld, Jason Rambach, and Didier Stricker.
TICaM: A time-of-flight in-car cabin monitoring dataset.
*arXiv preprint arXiv:2103.11719*, 2021.

valeo.**ai**

Okan Kopuklu, Jiapeng Zheng, Hang Xu, and Gerhard Rigoll.
Driver anomaly detection: A dataset and contrastive learning approach.
In *WACV*, 2021.

Manuel Martin, Alina Roitberg, Monica Haurilet, Matthias Horne, Simon Reiß, Michael Voit, and Rainer Stiefelhagen.
DriveAct: A multi-modal dataset for fine-grained driver behavior recognition in autonomous vehicles.
In *ICCV*, 2019.

École des Ponts
ParisTech

valeo.**ai**

📄 Juan Diego Ortega, Neslihan Kose, Paola Cañas, Min-An Chao, Alexander Unnervik, Marcos Nieto, Oihana Otaegui, and Luis Salgado.
Dmd: A large-scale multi-modal driver monitoring dataset for attention and alertness analysis.
In *ECCV*, 2020.

📄 Mohamed Selim, Ahmet Firintepe, Alain Pagani, and Didier Stricker.
AutoPOSE: Large-scale automotive driver head pose and gaze dataset with deep head orientation baseline.
In *VISIGRAPP (4: VISAPP)*, 2020.

École des Ponts
ParisTech

valeo.**ai**

# 3D Pose Lifting with CNN model: Semi-Supervised Approach



Figure: VideoPose3D Semi-supervised approach[14]

---

[14] Pavllo & al., 3D human pose estimation in video with temporal convolutions and semi-supervised training, in *CVPR*, 2019.

Figure: P-STMO model[15]

[15]Shan & al., P-STMO: Pre-trained spatial temporal many-to-one model for 3d human pose estimation. in *ECCV*, 2022.