

## 1. Motivation

Human mesh recovery from a single image is highly ambiguous, as an infinite set of predictions match the 2D observations equally.





 Multi-output methods propose diverse solutions given an image but are less accurate when making a single prediction.

# **2. MEGA**

Masked generative modeling enables efficient single- and multi-output predictions by leveraging a tokenized latent representation of human meshes.



MEGA relies on Mesh-VQ-VAE, an autoencoder that represents a human mesh as a sequence of 54 tokens.

# **CALLES PORT MEGA: Masked Generative Autoencoder for Human Mesh Recovery**

# Guénolé Fiche<sup>1,2</sup>

<sup>1</sup>CentraleSupélec <sup>2</sup>Naver Labs Europe <sup>3</sup>Inria Grenoble <sup>4</sup>Institut de Robòtica i Informàtica Industrial <sup>5</sup>Amazon

# **3. Training procedure**



• MEGA is first pre-trained for masked token prediction on a large-scale motion capture dataset without paired images  $\rightarrow$ random human mesh generation.

• Then, MEGA is fine-tuned for the same task but with additional conditioning from RGB images  $\rightarrow$  HMR.

• In both training stages, only the cross-entropy loss is used.

# 4. Deterministic mode

MEGA can be used as a single-output HMR model by predicting all tokens simultaneously given an image and a fully-masked sequence of human mesh tokens.

		3DPW			EMDB		
Method	Backbone	$PVE\downarrow$	$\mathbf{MPJPE}\downarrow$	PA-MPJPE $\downarrow$	$PVE\downarrow$	$\mathbf{MPJPE} \downarrow$	<b>PA-MPJPE</b> $\downarrow$
FastMETRO [13]	HRNet-w64	121.6	109.0	65.7	119.2	108.1	72.7
PARE [42]	HRNet-w32	97.9	82.0	50.9	133.2	113.9	72.2
Virtual Marker [61]	HRNet-w48	93.8	80.5	48.9	-	-	-0
CLIFF [52]	HRNet-w48	87.6	73.9	46.4	122.9	103.1	68.8
VQ-HPS [24]	HRNet-w48	84.8	71.1	45.2	112.9	99.9	65.2
MEGA (ours)	HRNet-w48	81.6	68.5	44.1	107.9	90.5	58.7
linear masking	HRNet-w48	86.5	72.6	45.9	118.7	100.1	63.3
full mask	HRNet-w48	<u>81.8</u>	68.5	44.4	<u>110.3</u>	92.7	59.2
w/o pre-training + full mask	HRNet-w48	84.1	<u>70.5</u>	46.2	113.9	95.9	62.0
HMR2.0 [27]	ViT-H	84.1	70.0	44.5	120.1	97.8	61.5
TokenHMR <sup>†</sup> [22]	ViT-H	88.1	76.2	49.3	124.4	102.4	67.5
TokenHMR <sup>∓</sup> [22]	ViT-H	84.6	71.0	44.3	109.4	91.7	55.6
MEGA (ours)	ViT-H	80.0	67.5	41.0	108.6	92.4	52.5

Table 1. Evaluation in deterministic mode. We evaluate MEGA on the 3DPW and EMDB datasets and compare it to the SOTA methods using metrics defined in Sec. 4.1 given in mm.  $\dagger$  stands for additionally using 2D training data, and  $\mp$  for additionally using 2D data and BEDLAM [7]. Methods in italic below the row "MEGA" indicate the results of the ablation study.

We compare MEGA to single-output methods on in-the-wild datasets without finetuning on the 3DPW dataset. MEGA is also SOTA on the occlusion dataset 3DPW-OCC despite not being designed specifically to handle occlusions.

## Simon Leglaive<sup>1</sup> Xavier Alameda-Pineda<sup>3</sup> Francesc Moreno-Noguer<sup>4,5</sup>



We can generate diverse human meshes iteratively. MEGA is then a multi-output HMR method.





# per-vertex uncertainty:

![](_page_0_Picture_36.jpeg)

![](_page_0_Picture_37.jpeg)

![](_page_0_Picture_38.jpeg)

![](_page_0_Picture_40.jpeg)

### **5. Stochastic mode**

The iterative process can be visualized as follows:

The diversity in prediction can be interpreted as a measure of the