

Avat3r

Large Animatable Reconstruction Model for High-fidelity 3D Head Avatars

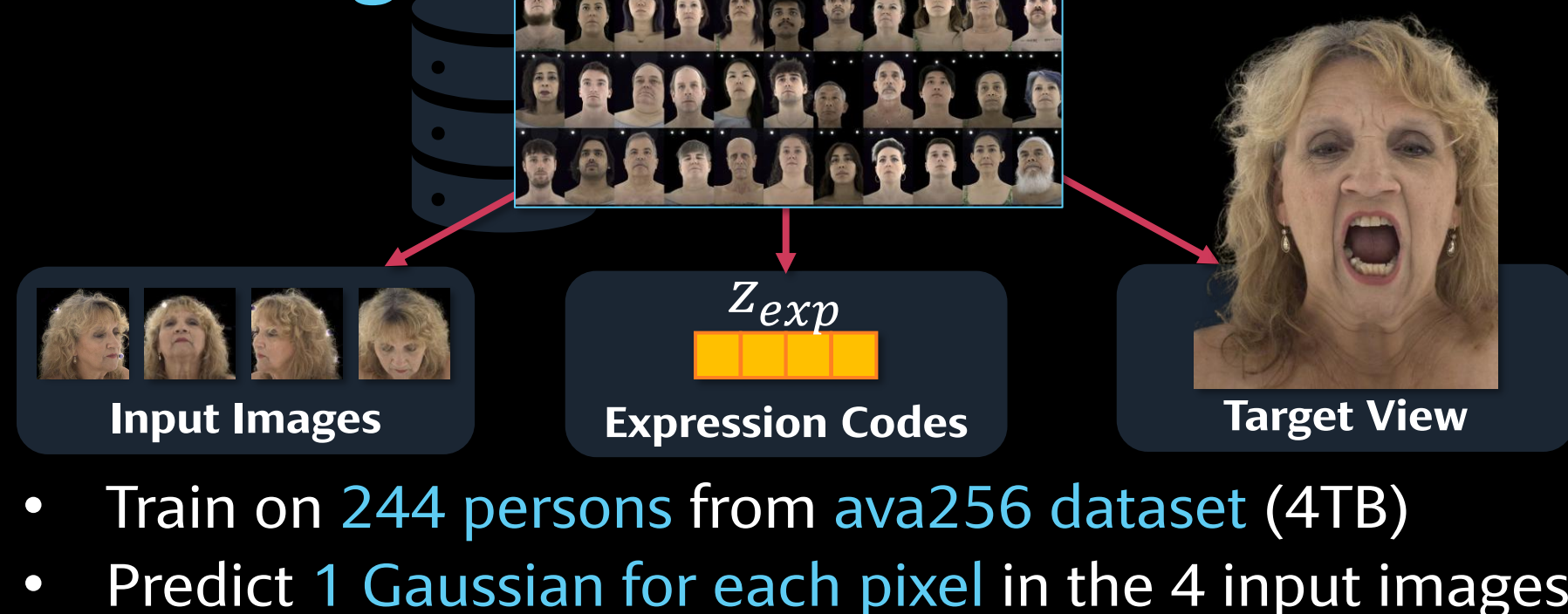
Tobias Kirschstein^{1,2} - Javier Romero² - Artem Sevastopolsky^{1,2} - Matthias Nießner¹ - Shunsuke Saito²

¹Technical University of Munich
²Meta Reality Labs

Highlights

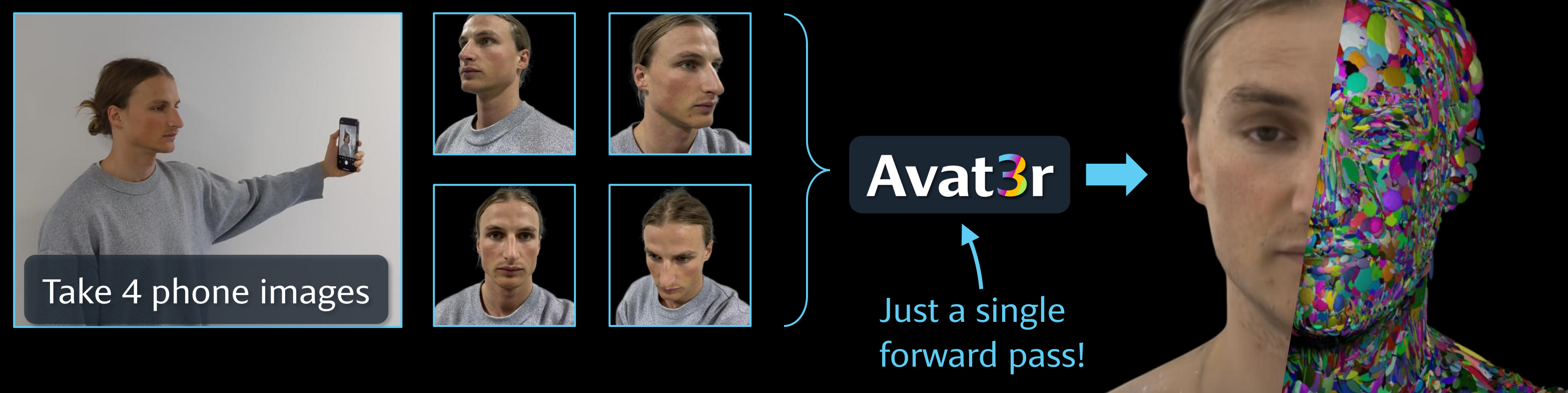
- High-quality 3D Head Avatars from just 4 images
- Make Large Reconstruction Models (LRMs) animatable via simple cross-attention to the animation signal (facial expression codes)
- No reliance on mesh-based 3D face model: Learn facial animations entirely from data
- Improve per-pixel 3D Gaussian predictions with position maps from DUST3R and feature maps from Sapiens

Training



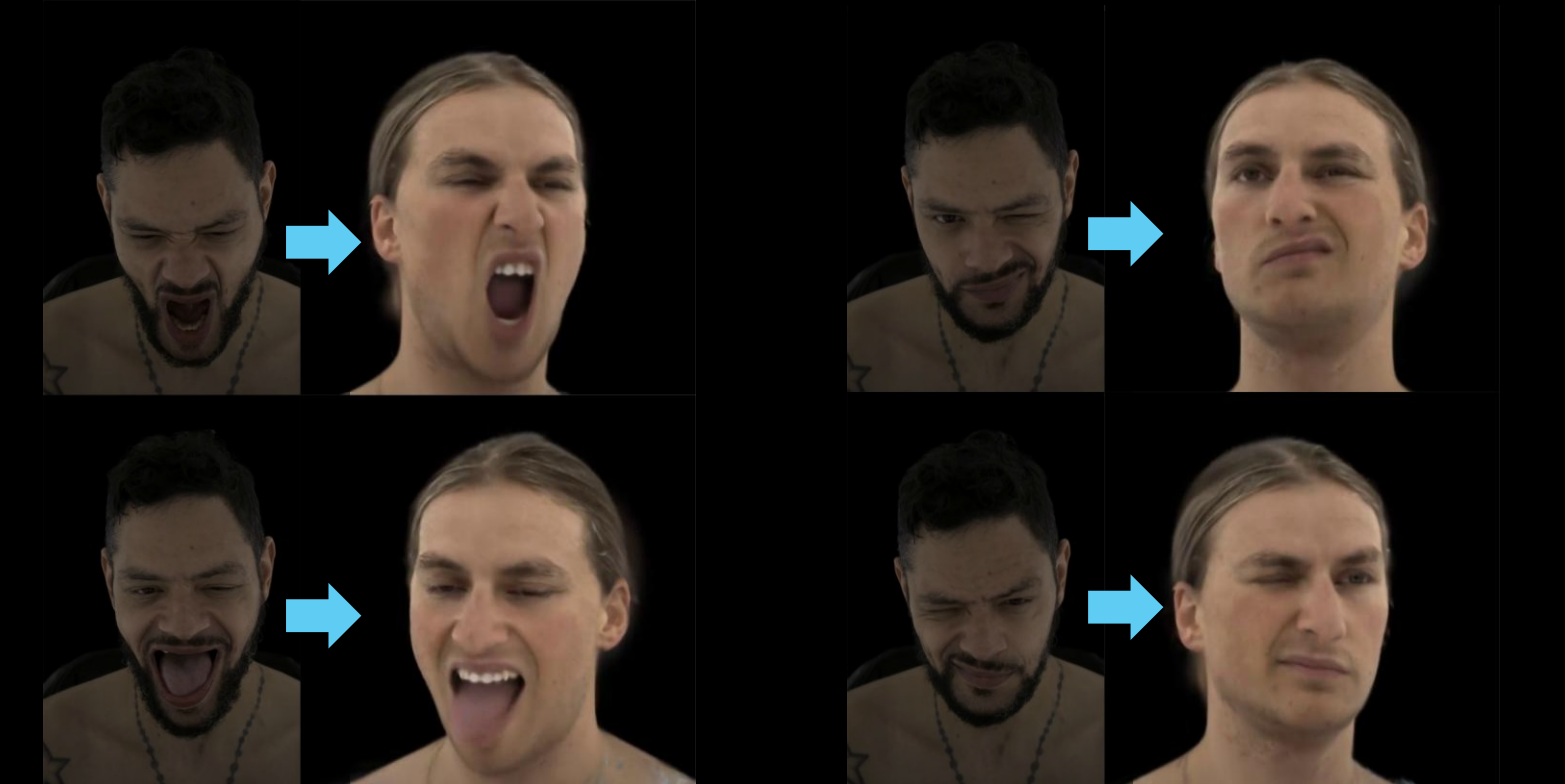
- Train on 244 persons from ava256 dataset (4TB)
- Predict 1 Gaussian for each pixel in the 4 input images

Few-shot 3D Head Avatar Creation

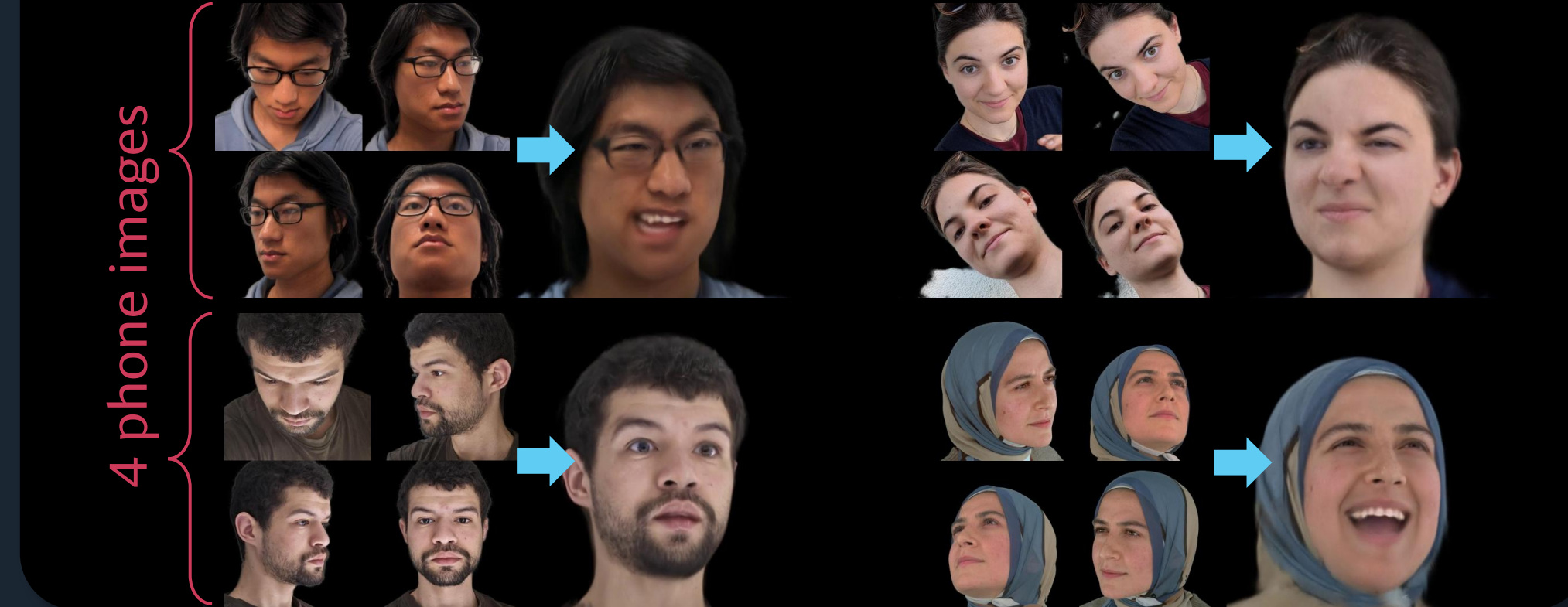


Take 4 phone images → Avat3r → Just a single forward pass!

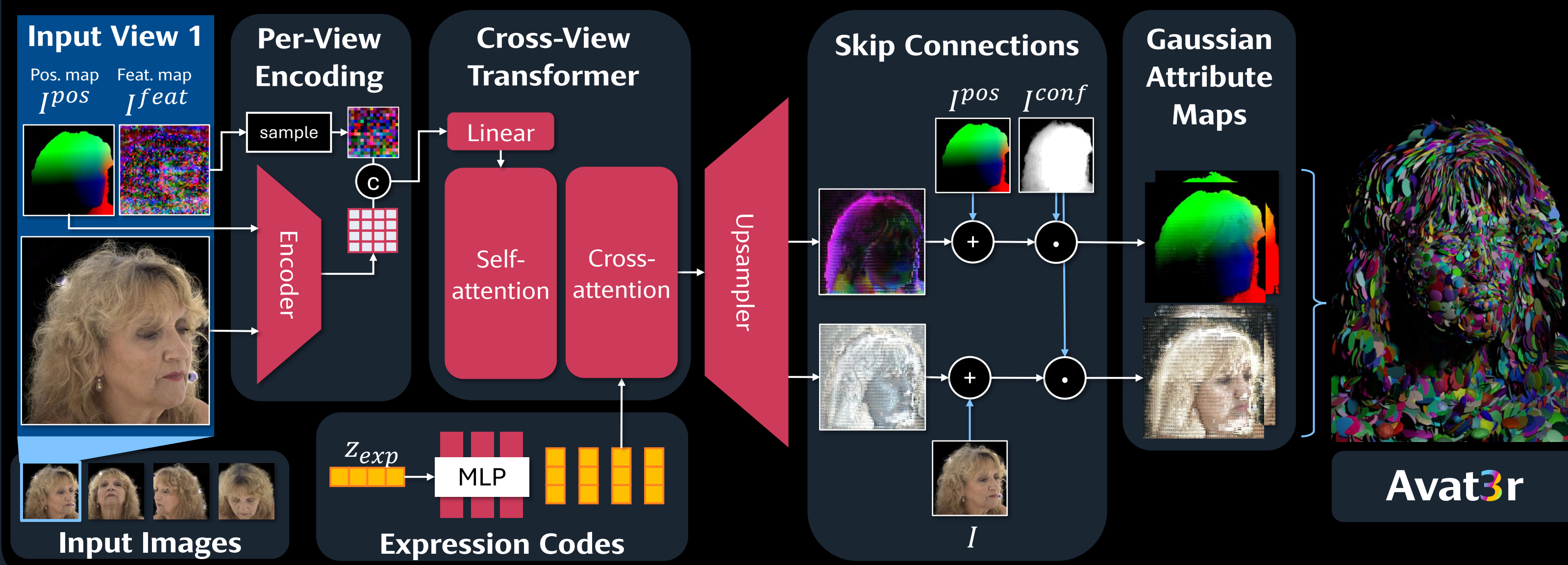
Expressive Animation



More few-shot Avat3rs



Method Overview



Single-image Pipeline



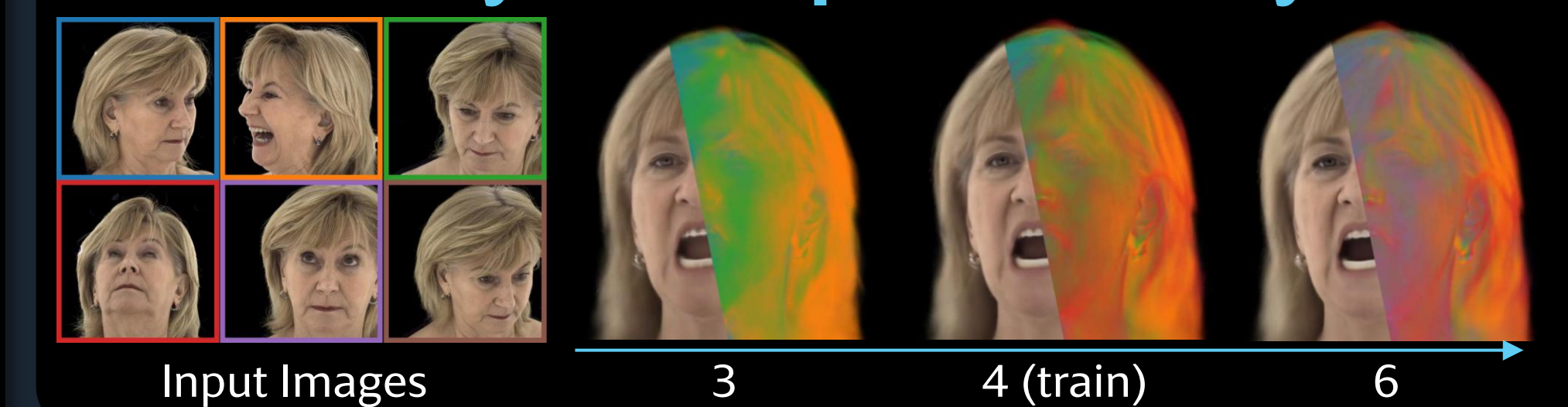
Single Input Image → 3D GAN Inversion → Avat3r → 3D head sampled from 3D GAN

Comparisons



Inputs, InvertAvatar, GPAvatar, Avat3r, GT

Analysis of Input Flexibility



Input Images, 3, 4 (train), 6

Ablations



Input Images, w/o DUST3R, w/o Sapiens, w/o rand. timesteps, Ours full

Conclusion

- Feed-forward 3D head avatar reconstruction from 4 images (no optimization needed)
- Surprisingly good generalization from just 244 persons
- Cross-attention is enough to model complex facial expressions
- No 3D face model needed → Can animate tongue
- Model is robust to inconsistent expressions in a phone scan

	\mathcal{S}	\mathcal{D}	\mathcal{T}	PSNR↑	LPIPS↓	AKD↓
Just cross-attention		☒		20.8	0.439	8.60
w/o Sapiens		☒	☒	20.9	0.434	8.59
w/o Dust3r		☒	☒	21.1	0.429	9.60
w/o rand. timesteps		☒	☒	21.3	0.409	8.86
Avat3r	☒	☒	☒	21.6	0.410	8.08

\mathcal{S} : Use Sapiens feature maps
 \mathcal{D} : Use DUST3R position maps
 \mathcal{T} : Sample inconsistent expressions during training