

Capturing and Synthesis of 3D Digital Humans

Wojciech Zielonka

Committee:

Prof. Justus Thies

Prof. Matthias Niessner

Prof. Jan Gugenheimer

Prof. Simone Schaub-Meyer



MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



**TECHNISCHE
UNIVERSITÄT
DARMSTADT**

Meta Codec Avatars



CODEC AVATAR



CODEC AVATAR





Potential Applications of Digital Avatars

🧠 Research & Simulation [1]



🛍️ Retail & Fashion [2]



🤖 Human-Computer Interaction [3]

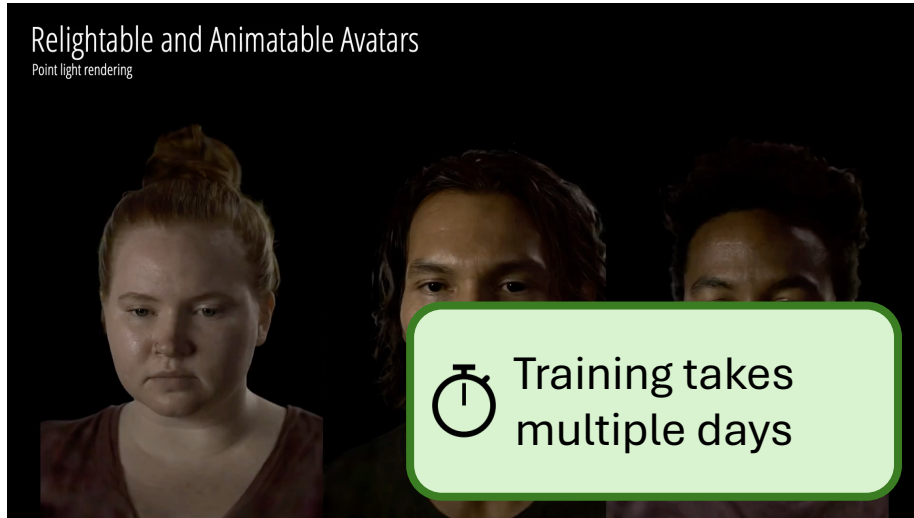


🎮 Entertainment & Media [4]

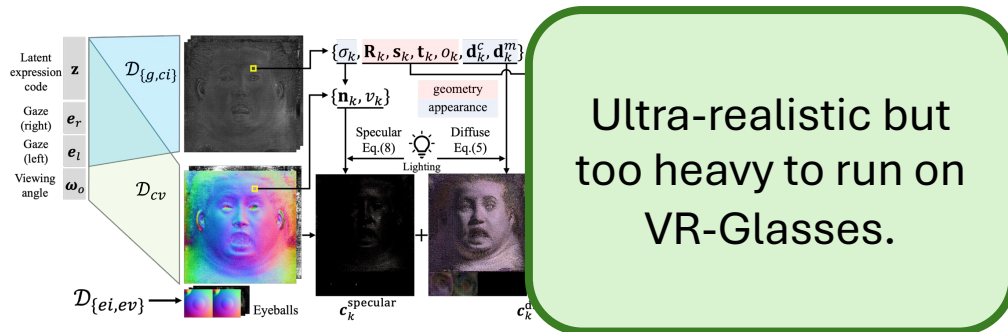


- 1) [Human behavior modeling](#)
- 2) [Google Virtual Try-on](#)
- 3) [Apple Persona](#)
- 4) [Runways ACT 2](#)

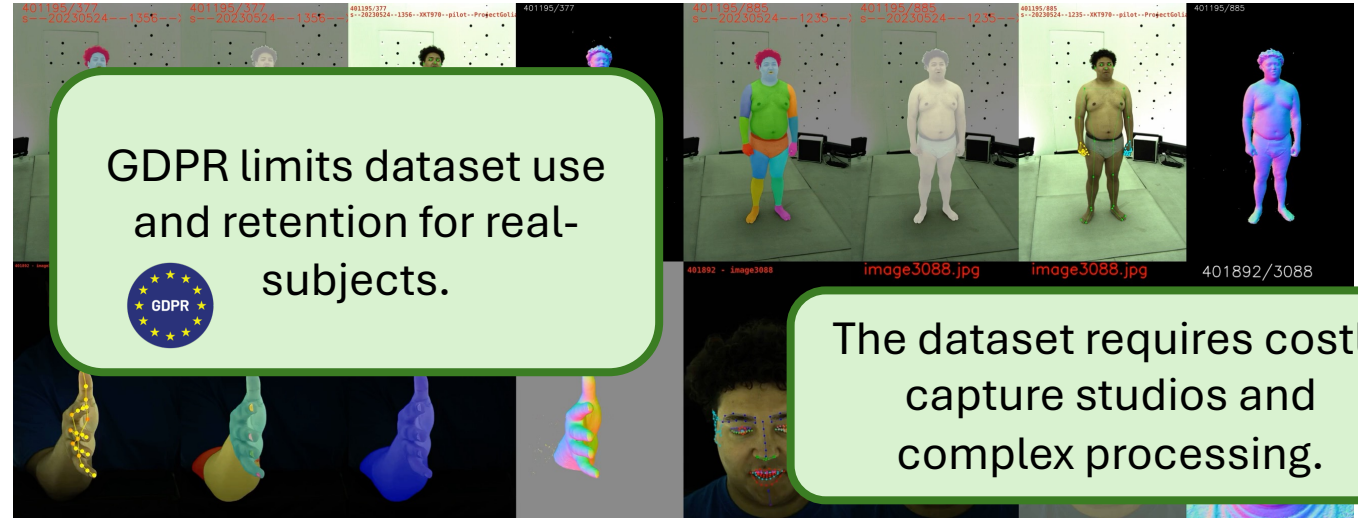
Why is it a Hard Problem?



Photorealism comes at the cost of sophisticated and powerful networks. [2]



Ultra-realistic but too heavy to run on VR-Glasses.



Creating photorealistic 3D digital avatars requires assets that take up hundreds of gigabytes and take weeks to process. [3]



The avatar must be photorealistic.

Any lack of photorealism immediately breaks immersion and makes the avatars look artificial. [1]

- 1) A 3D Face Model for Pose and Illumination Invariant Face Recognition, Paysan *et al.*
- 2) Relightable Gaussian Codec Avatars, Saito *et al.*
- 3) Codec Avatar Studio, Martinez *et al.*

Challenges Addressed in this PhD Work



☑ Speed of the avatar's creation.



☑ Few-shot inversion by utilizing synthetic prior.



☑ Distillation to lightweight representation.

Challenges Addressed in this PhD Work



☑ Speed of the avatar's creation.

- Many methods need days (even up to 5 days) to train a single avatar.
- Appearance changes were not address due to compute heavy training.
- They often lack real-time capabilities.



☑ Few-shot inversion by utilizing synthetic prior.



☑ Distillation to lightweight representation.

Challenges Addressed in this PhD Work



☑ Speed of the avatar's creation.

- Many methods need days (even up to 5 days) to train a single avatar.
- Appearance changes were not address due to compute heavy training.
- They often lack real-time capabilities.



☑ Few-shot inversion by utilizing synthetic prior.

- Monocular avatars don't generalize well to new views or expressions.
- 3D models need diverse paired data, which is scarce.
- GDPR limits dataset use and retention.



☑ Distillation to lightweight representation.

Challenges Addressed in this PhD Work



☑ Speed of the avatar's creation.

- Many methods need days (even up to 5 days) to train a single avatar.
- Appearance changes were not address due to compute heavy training.
- They often lack real-time capabilities.



☑ Few-shot inversion by utilizing synthetic prior.

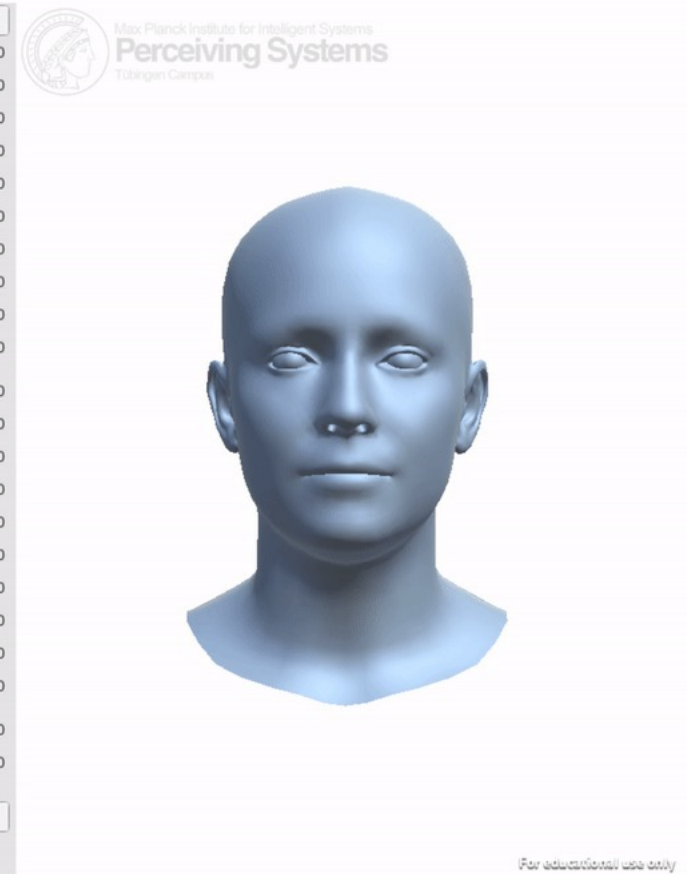
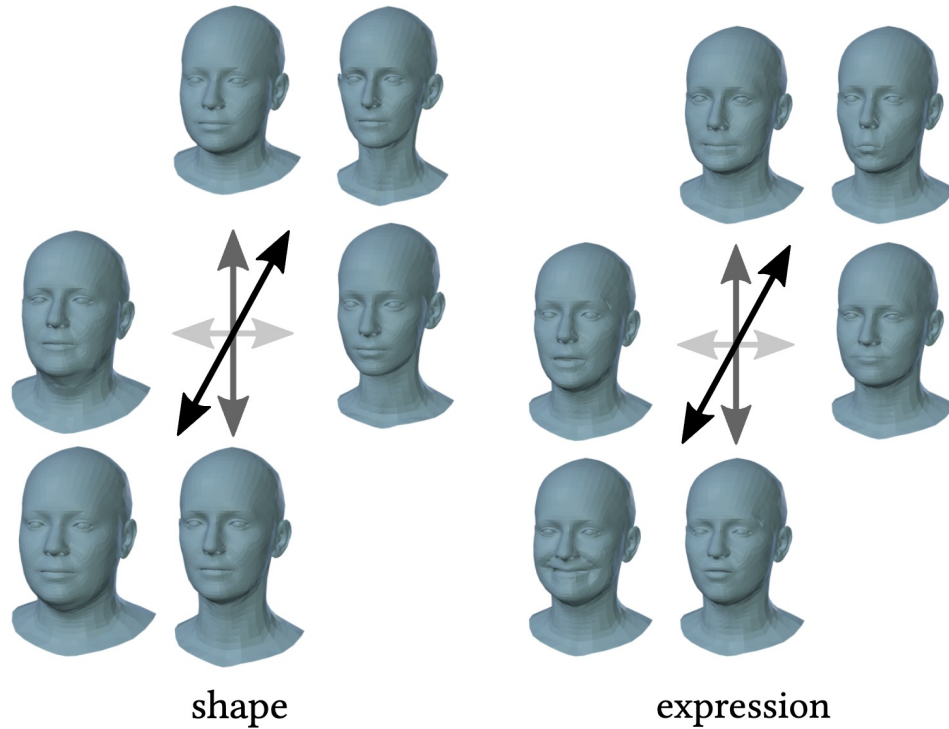
- Monocular avatars don't generalize well to new views or expressions.
- 3D models need diverse paired data, which is scarce.
- GDPR limits dataset use and retention.



☑ Distillation to lightweight representation.

- Lightweight models look unrealistic.
- High-quality avatars are too resource-heavy for common devices.

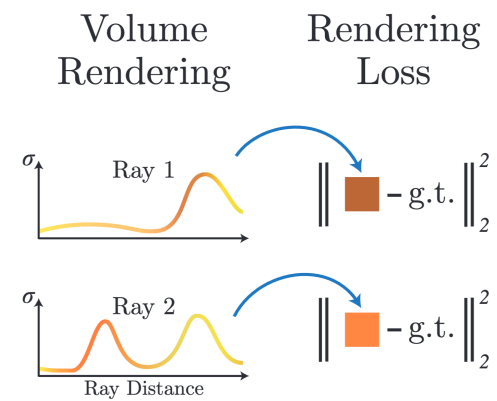
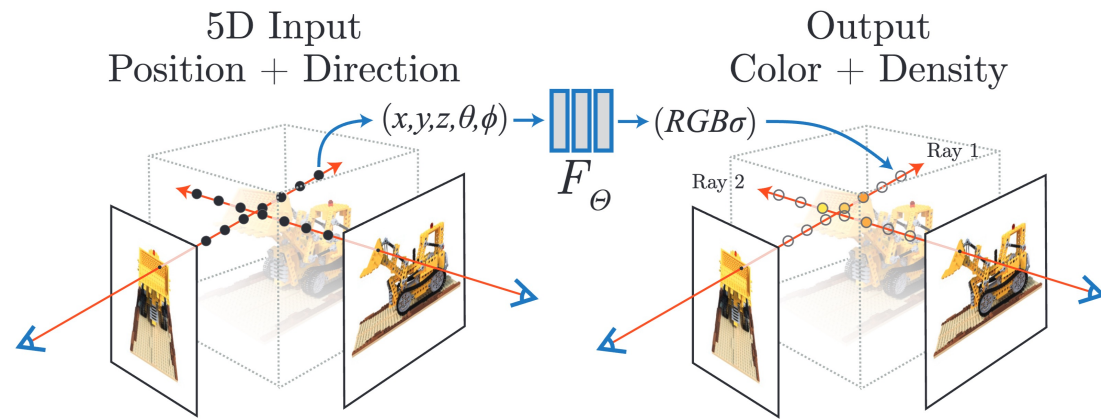
FLAME*



FLAME is a common 3DMM model used as a prior in most of 3D face avatar methods

$$\mathbf{S} = \bar{\mathbf{S}} + \delta \mathbf{B}_{\text{id}} + \gamma \mathbf{B}_{\text{expr}}$$
$$\mathbf{C} = \bar{\mathbf{C}} + \sigma \mathbf{D}_{\text{id}}$$

NeRF*



$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i$$



3DGS*



Primitives as ellipsis



Splatted Gaussians

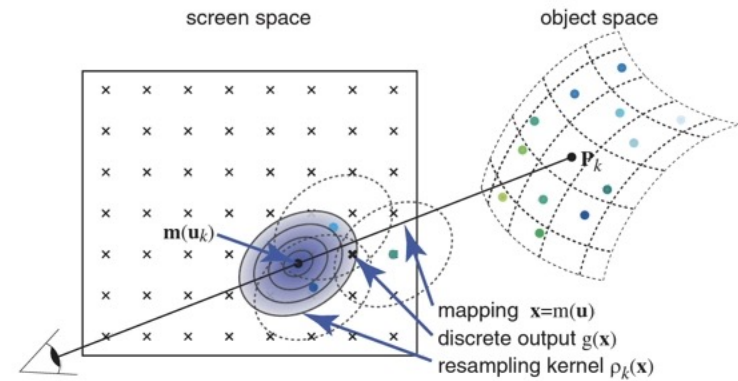
Anisotropic Volumetric 3D Gaussians



Final Rendering

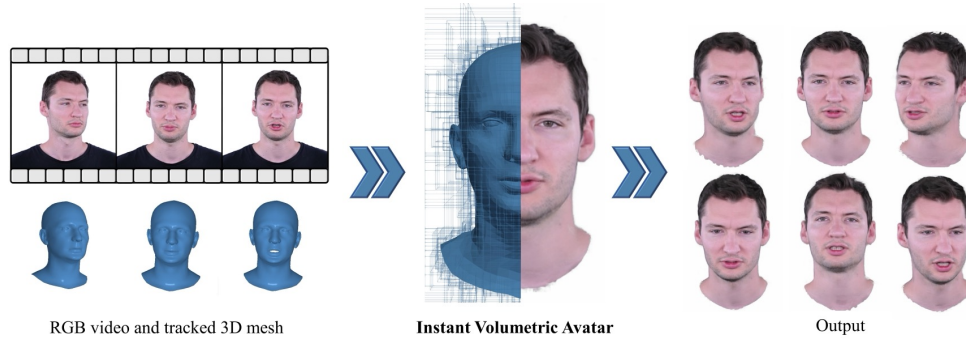
3D Gaussian Visualization

3D Gaussian Splatting by Kerbl et al.

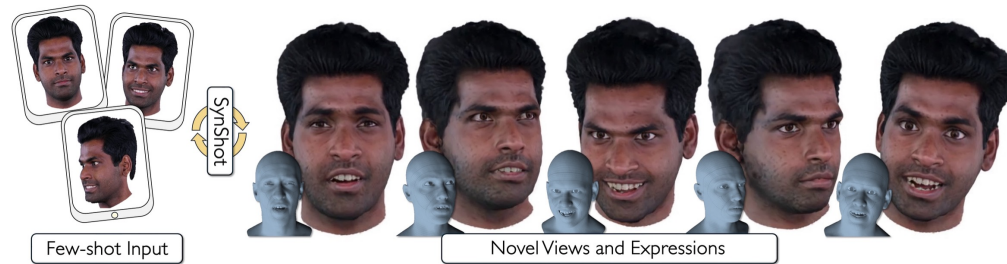


Surface Splatting by Zwicker et al.

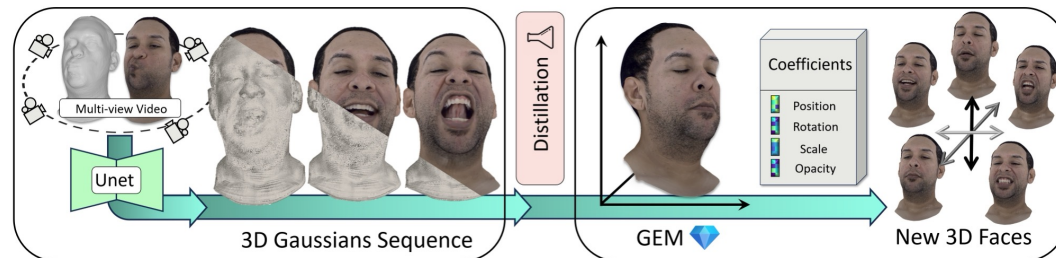
Presentation Outline



INSTA - Instant Volumetric Head Avatars
[Zielonka, Bolkart, Thies]
CVPR'23

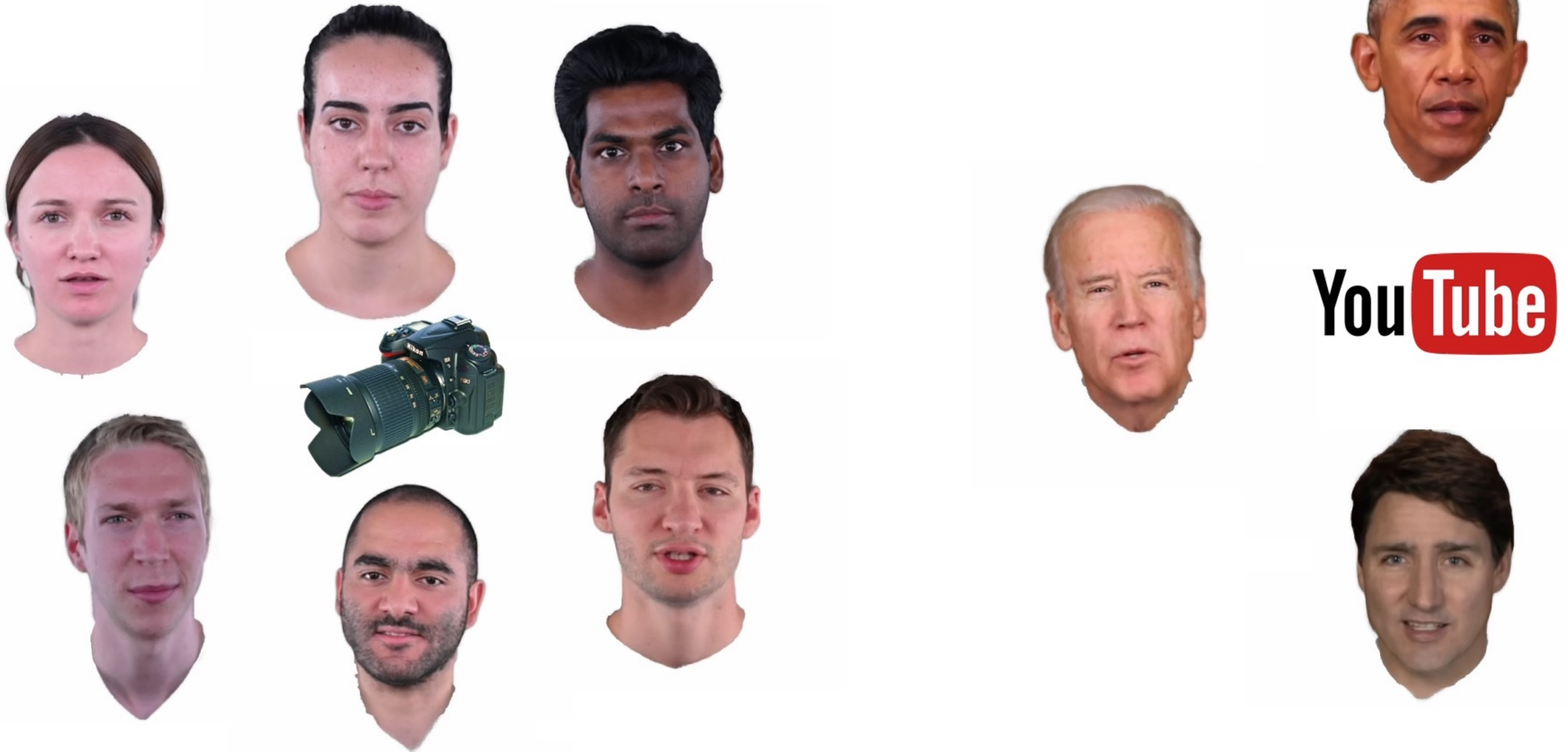


SynShot - Synthetic Prior for Few-Shot Drivable Avatar Inversion
[Zielonka, Garbin, Lattas, Kopanas, Gotardo, Beeler, Thies, Bolkart]
CVPR'25



GEM - Gaussian Eigen Models for Human Heads
[Zielonka, Bolkart, Beeler, Thies]
CVPR'25

Motivation: Avatar from a monocular video



Motivation: Optimization Time



~ 2 days

1 GPUs

NHA [Grassal *et al.*]



~ 3-4 days

1 GPU

NeRF [Gafni *et al.*]



~ 4-5 days

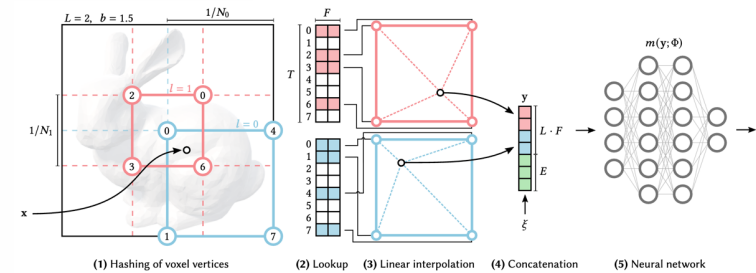
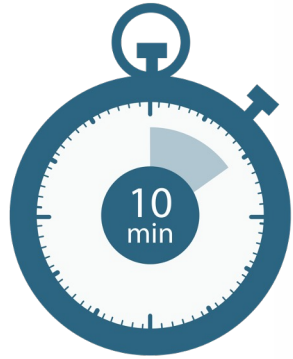
1 GPU

IMAvatar [Zheng *et al.*]



Ground truth

Solution: NeRF embedded on the surface of a mesh



Instant Neural Graphics Primitives with a Multiresolution Hash Encoding by Muller *et al.*

Solution: Order of magnitude time improvement



~ 2 days
1 GPU

NHA [Grassal *et al.*]



~ 3-4 days
1 GPU

NeRFace [Gafni *et al.*]



~ 4-5 days
1 GPU

IMAvatar [Zheng *et al.*]



~ 10 minutes
1 GPU
Ours



Ground truth

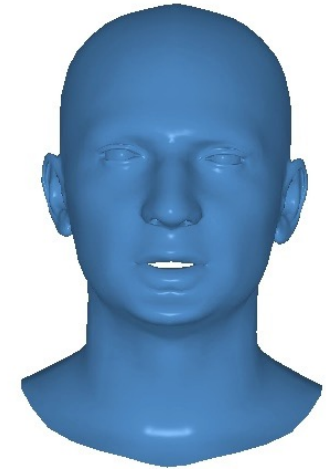
Method: Input (monocular video + tracked mesh)



RGB video

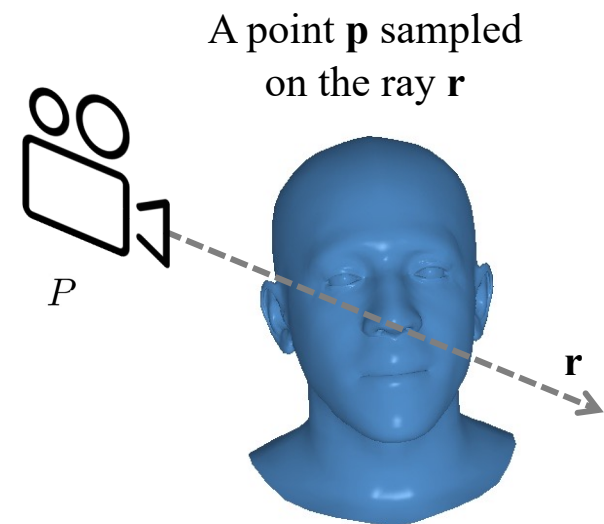
Landmarks

Statistical
texture



Geometry

Method: Pipeline



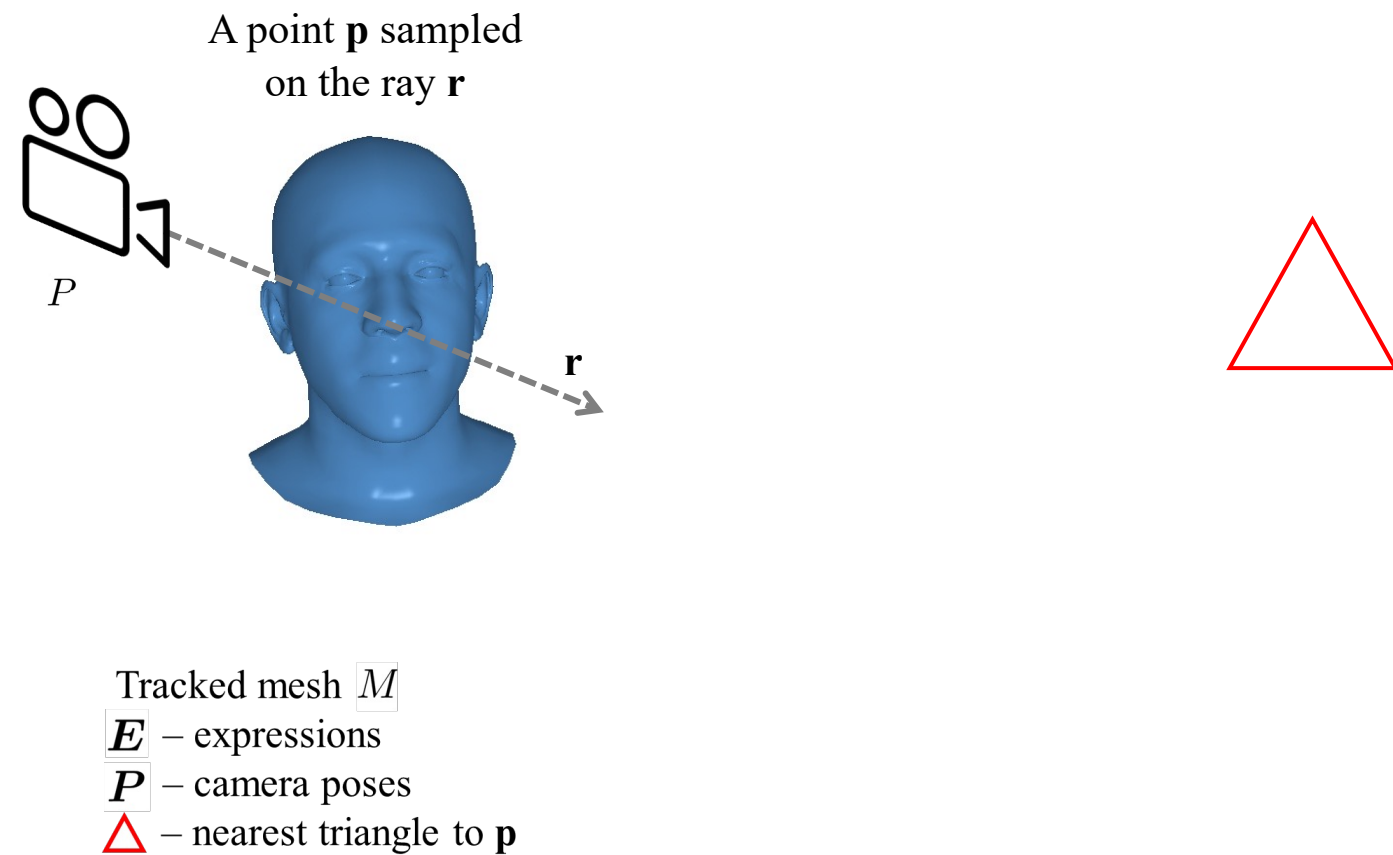
Tracked mesh M

E – expressions

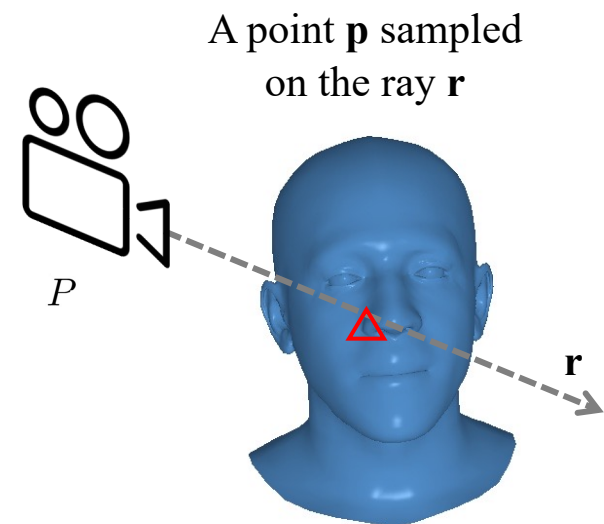
P – camera poses

\triangle – nearest triangle to \mathbf{p}

Method: Pipeline



Method: Pipeline



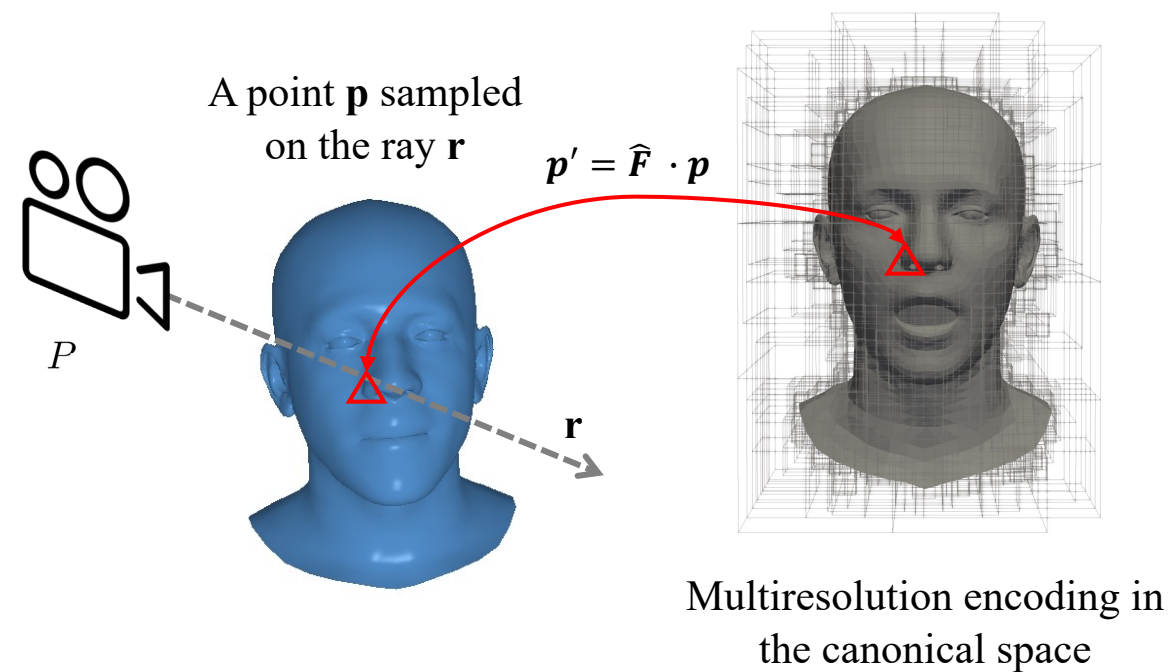
Tracked mesh \bar{M}

\bar{E} – expressions

\bar{P} – camera poses

\triangle – nearest triangle to \mathbf{p}

Method: Pipeline

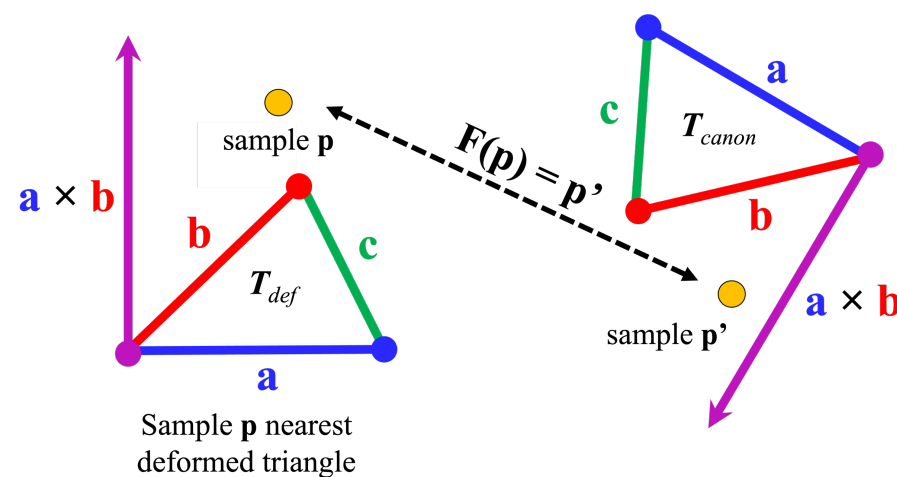


Tracked mesh \mathbf{M}

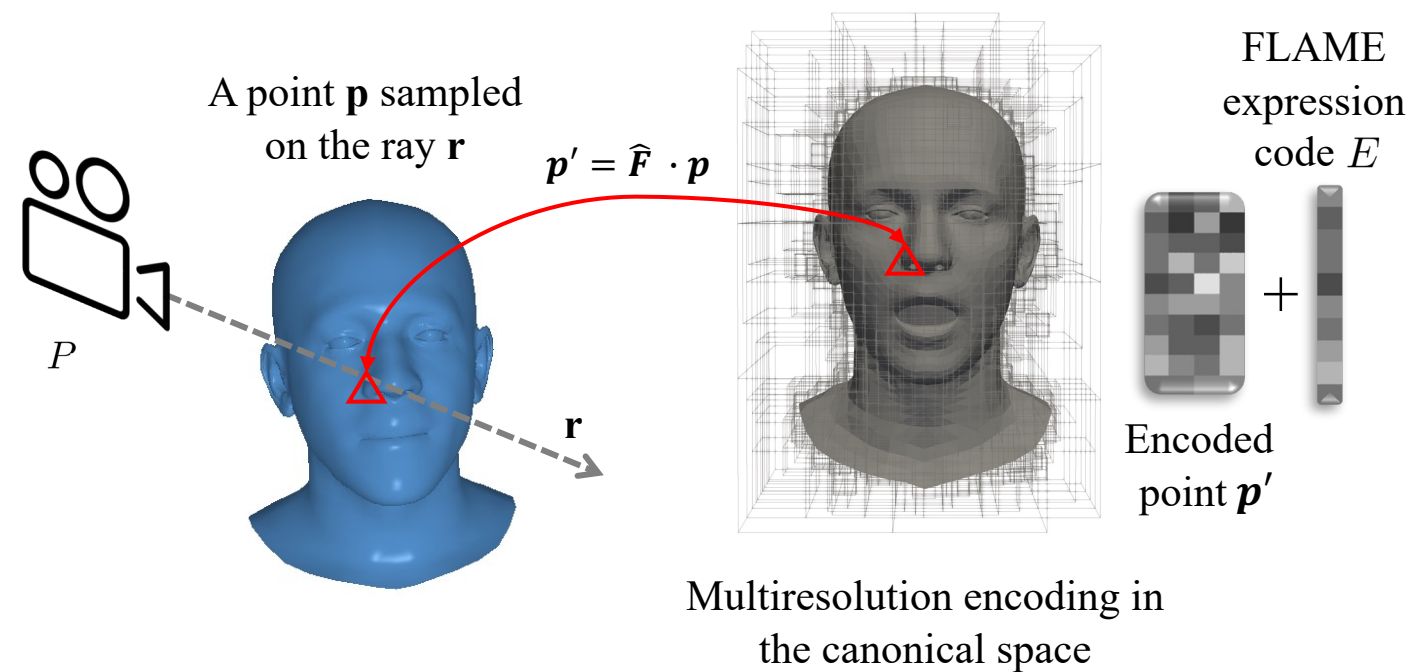
\mathbf{E} – expressions

\mathbf{P} – camera poses

\triangle – nearest triangle to \mathbf{p}



Method: Pipeline



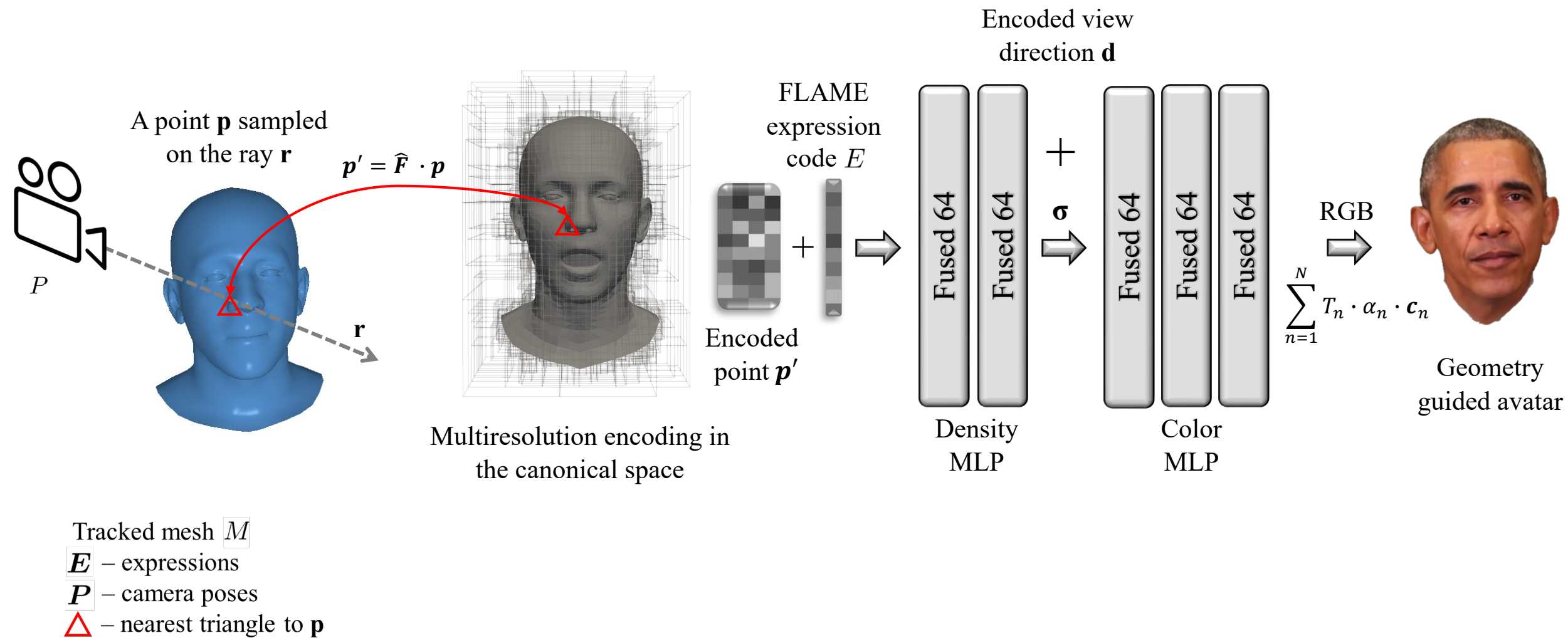
Tracked mesh M

E – expressions

P – camera poses

\triangle – nearest triangle to \mathbf{p}

Method: Pipeline



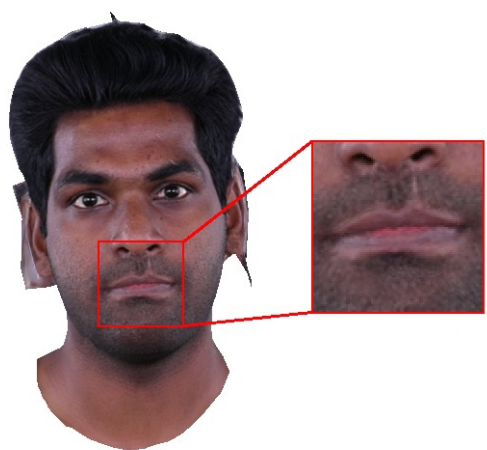
Short Real-time Demo



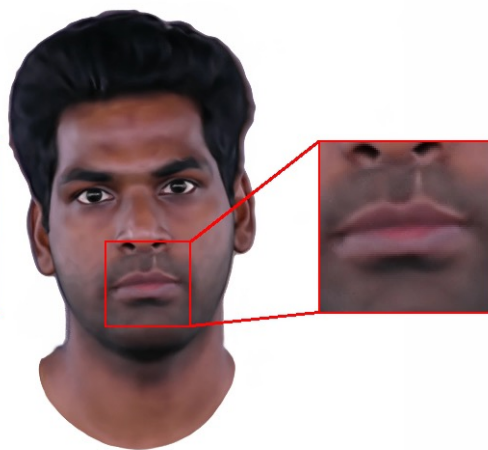
▶ instant-ngp v1.0dev

▶ Camera path

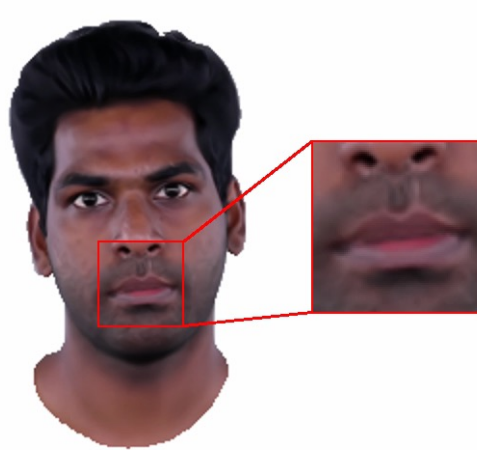
Results: Qualitative Comparison



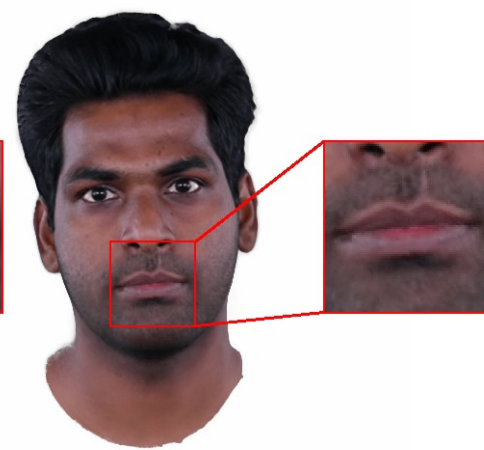
NHA
[Grassal et al.]



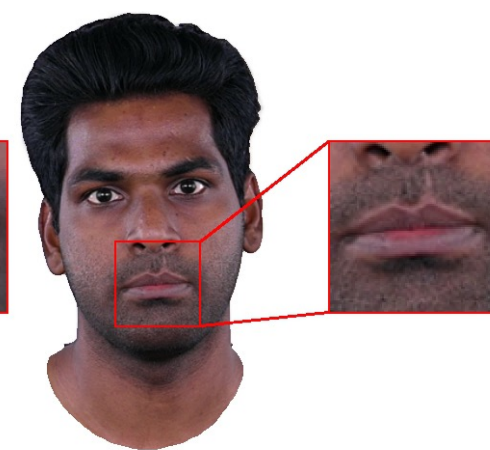
NeRF face
[Gafni et al.]



IMAvatar
[Zheng et al.]



Ours



Ground truth

Results: Quantitative Comparison

Method	L2 ↓	PSNR ↑	SSIM ↑	LPIPS ↓
NHA [Grassal et al.]	0.0018	28.65	0.96	0.03
IMAvatar [Zheng et al.]	0.0014	29.10	0.95	0.06
NeRFace [Gafni et al.]	0.0010	30.87	0.96	0.05
Ours	0.0010	30.51	0.96	0.03

Results: Retargeting

INSTA works well for mostly small deformations but lacks a prior to properly generalize to other identities.

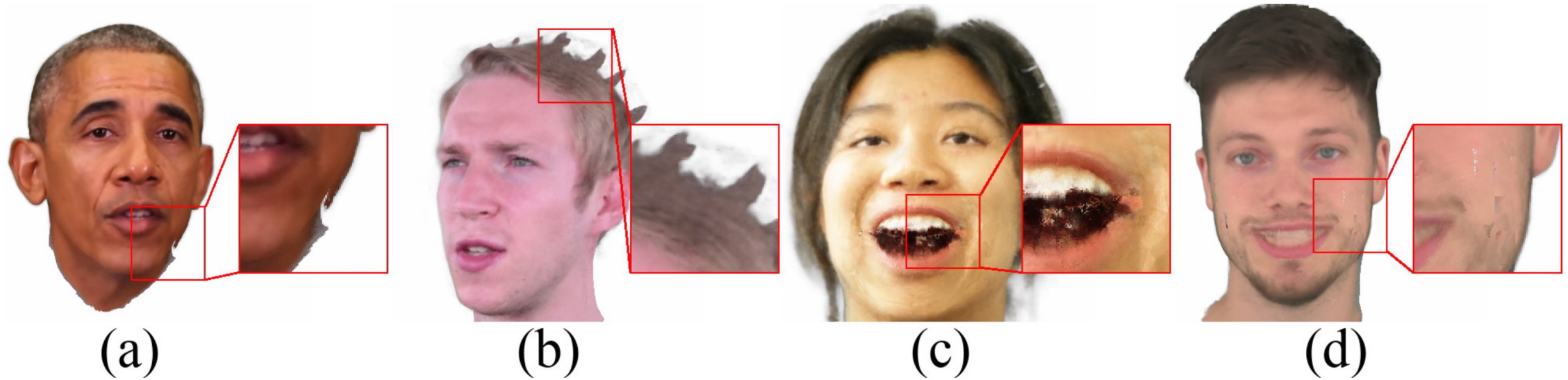


Source



Targets

Limitations



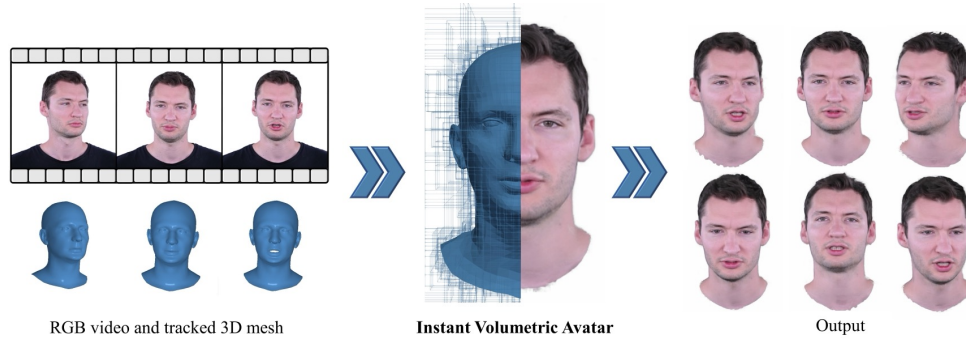
Failure cases: **(a)** and **(b)** exhibits outline artifacts at the chin and hair which stem from geometry misalignment of the tracker, **(c)** extreme expressions can cause artifacts in the mouth region, and **(d)** extrapolation of expressions can lead to artifacts.

Take-home Messages of INSTA

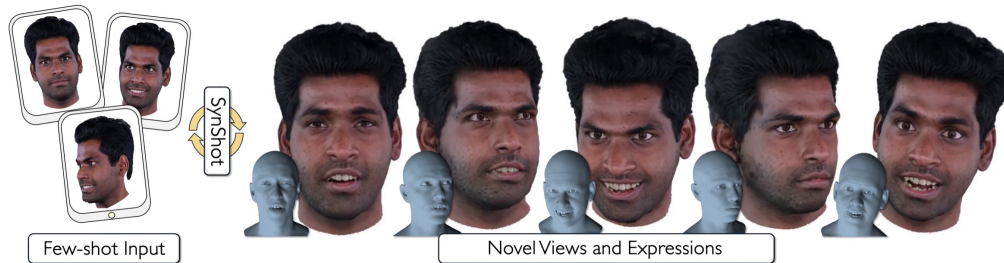
1. Given a monocular RGB video as input, we optimize a controllable avatar in less than 10 minutes.
2. In this way, we can create a new avatar almost on the fly that reflects the current appearance instead of a prerecorded one.



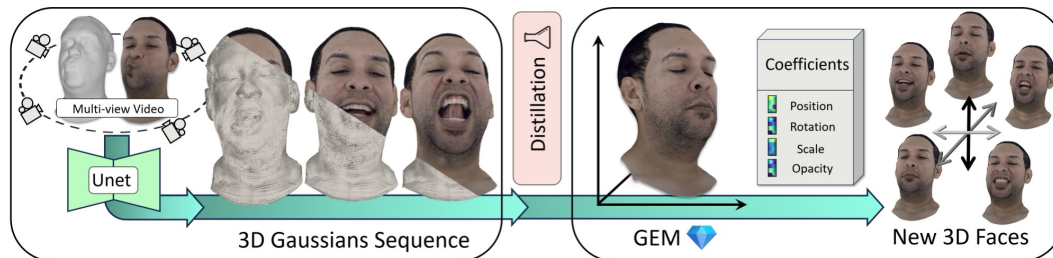
Presentation Outline



INSTA - Instant Volumetric Head Avatars
[Zielonka, Bolkart, Thies]
CVPR'23



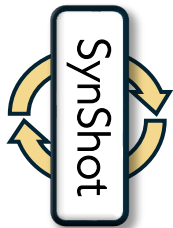
SynShot - Synthetic Prior for Few-Shot Drivable Avatar Inversion
[Zielonka, Garbin, Lattas, Kopanas, Gotardo, Beeler, Thies, Bolkart]
CVPR'25



GEM - Gaussian Eigen Models for Human Heads
[Zielonka, Bolkart, Beeler, Thies]
CVPR'25

Motivation: Few-shot avatar inversion

Few-shot Input



Novel Views and Expressions

Motivation: Monocular Methods



Source



INSTA¹



Splatting Avatar²



Flash Avatar³

1) Zielonka *et al.* Instant Volumetric Head Avatars

2) Xiang *et al.* FlashAvatar: High-fidelity Head Avatar with Efficient Gaussian Embedding

3) Zhijing *et al.* SplattingAvatar: Realistic Real-Time Human Avatars with Mesh-Embedded Gaussian Splatting

Motivation: GAN-based Methods



Source



InvertAvatar¹



Portrait4D²



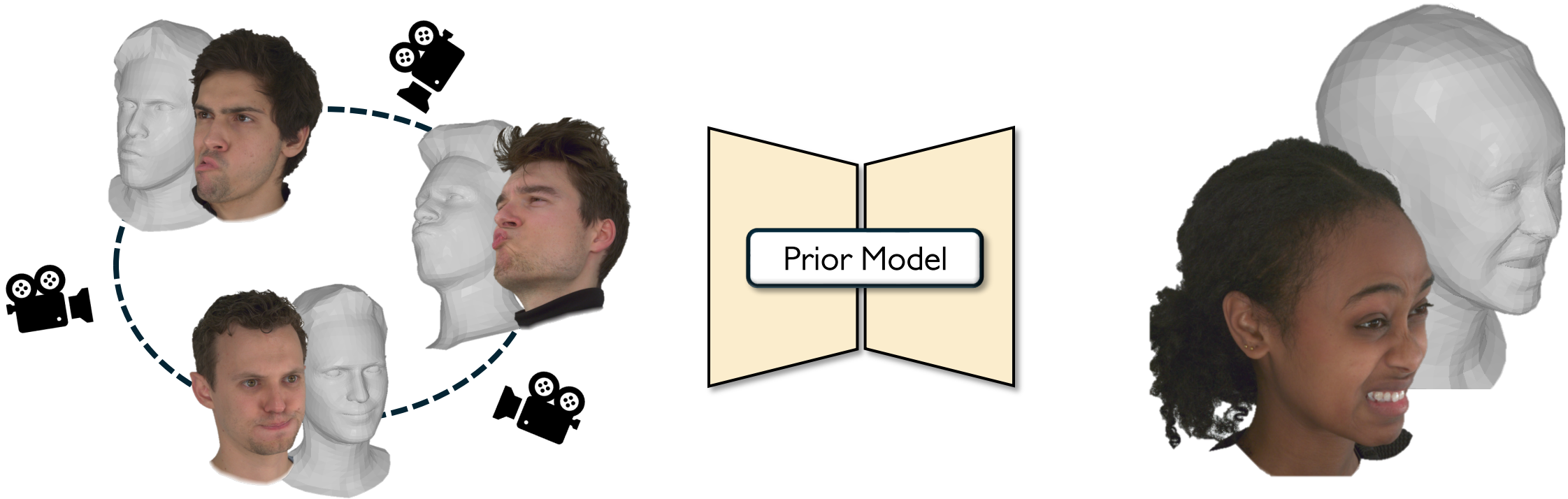
Next3D³

1) Zhao *et al.* InvertAvatar: Incremental GAN Inversion for Generalized Head Avatars

2) Deng *et al.* Portrait4D: Learning One-Shot 4D Head Avatar Synthesis using Synthetic Data

3) Sun *et al.* Next3D: Generative Neural Texture Rasterization for 3D-Aware Head Avatars

Solution: Train Prior Model



Use multi-view dataset with tracked meshes to build a prior model used for inversion and driving the avatars.

Solution: Real Datasets?



Nersemble¹



Multiface²



FaceScape³

1) Kirschstein *et al.* NeRSemble: Multi-View Radiance Field Reconstruction of Human Heads

2) Wu *et al.* Multiface: A Dataset for Neural Face Rendering

3) Zhu *et al.* FaceScape: 3D Facial Dataset and Benchmark for Single-View 3D Face Reconstruction

Solution: Real Datasets?



The **General Data Protection Regulation (GDPR)** is an EU law that protects individuals' personal data and privacy, enforced since May 25, 2018.

What does it mean for **Digital Humans research**:

1. Dataset derivatives must be frequently deleted e.g., each 30 days.
2. Trained models the same, periodically removed.



Nersemble¹

Multiface²

FaceScape³

1) Kirschstein *et al.* NeRSemble: Multi-View Radiance Field Reconstruction of Human Heads

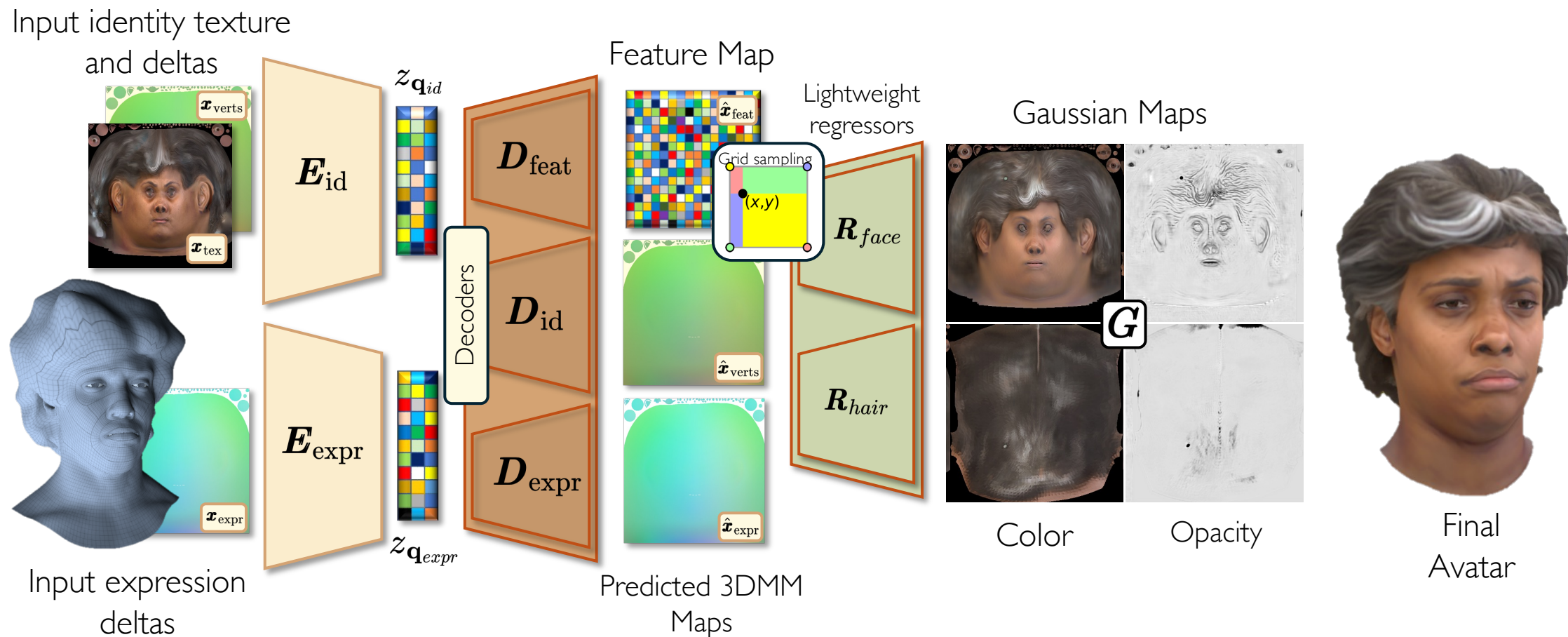
2) Wu *et al.* Multiface: A Dataset for Neural Face Rendering

3) Zhu *et al.* FaceScape: 3D Facial Dataset and Benchmark for Single-View 3D Face Reconstruction

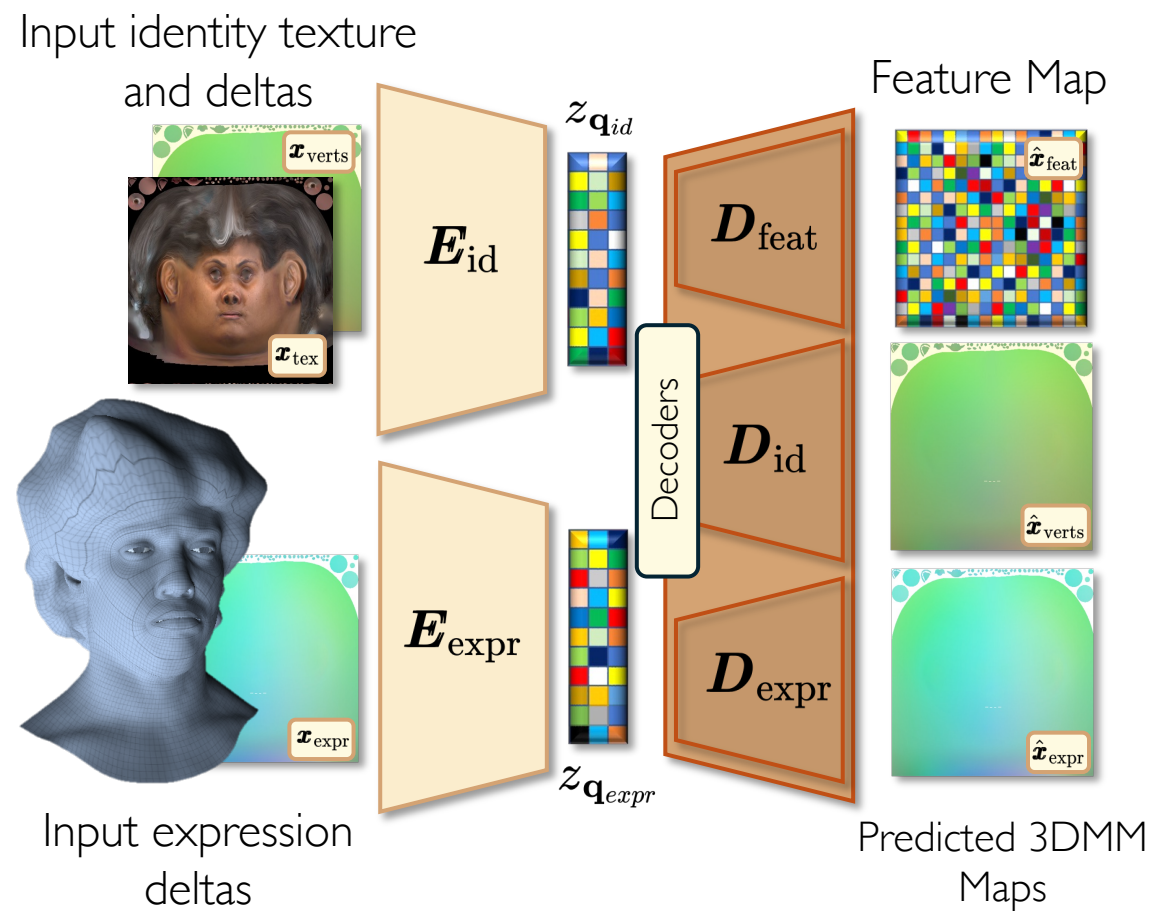
Solution: Synthetic Datasets ✓



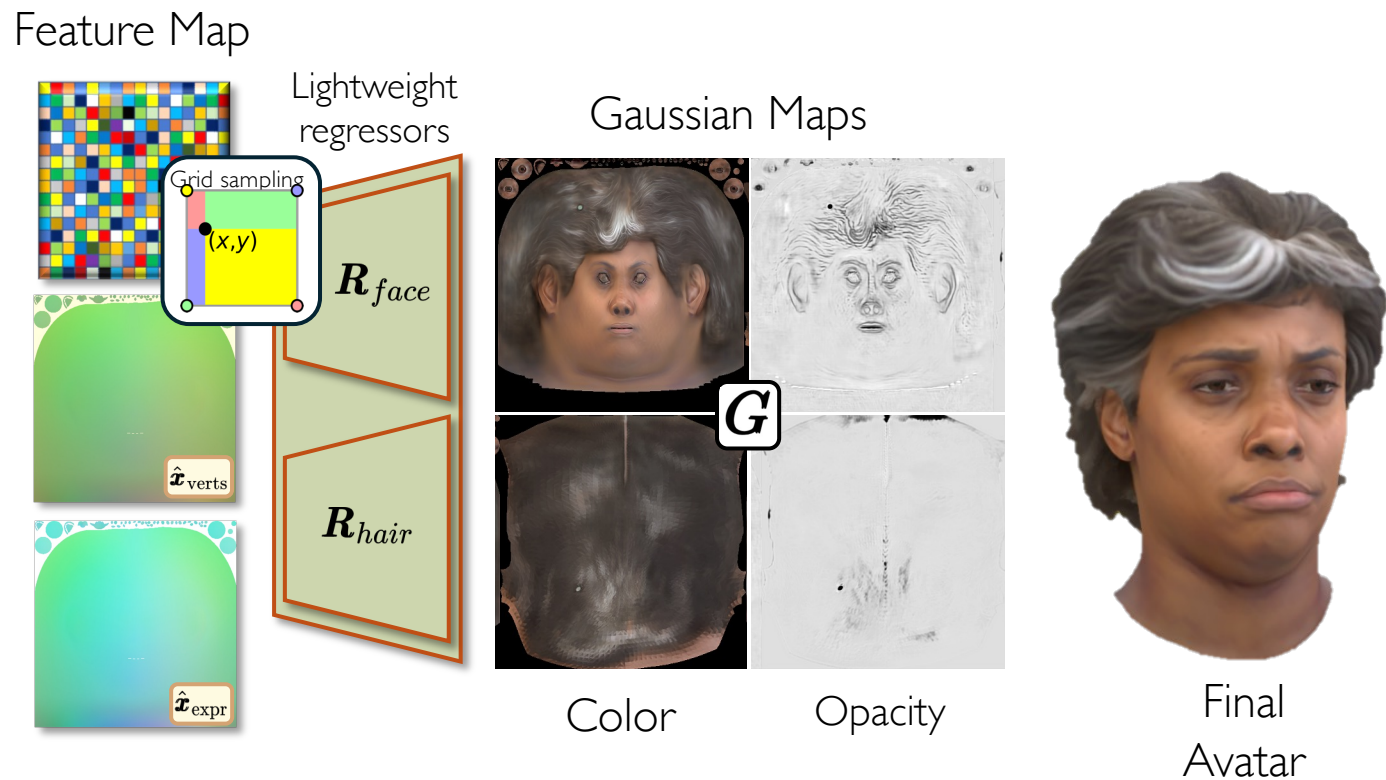
Method: Pipeline



Method: Encoder



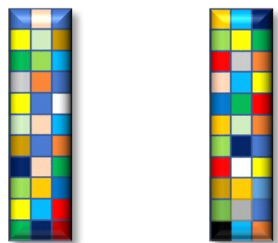
Method: Decoder



Method:

Latent Space

$z_{q_{id}}$ $z_{q_{expr}}$



Interpolation



Method: Avatar Inversion



Few-shot Input

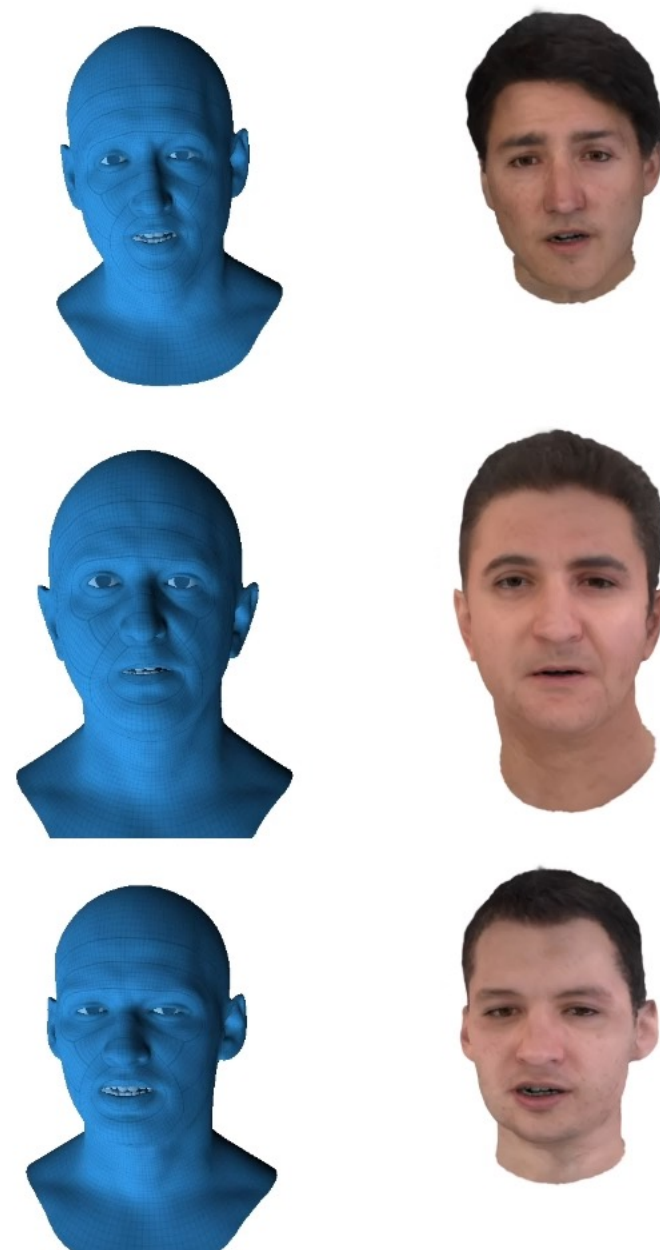
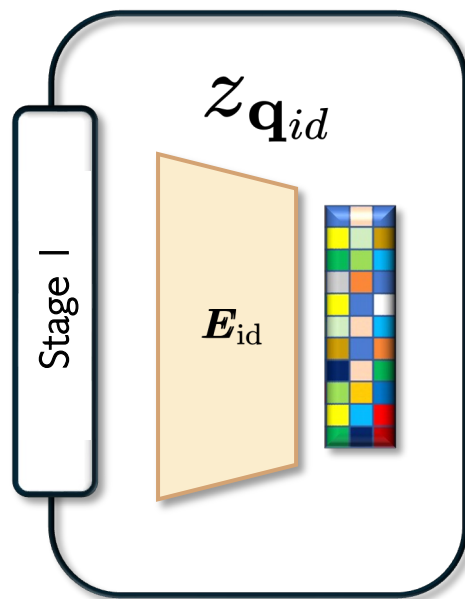


Inverted Avatar

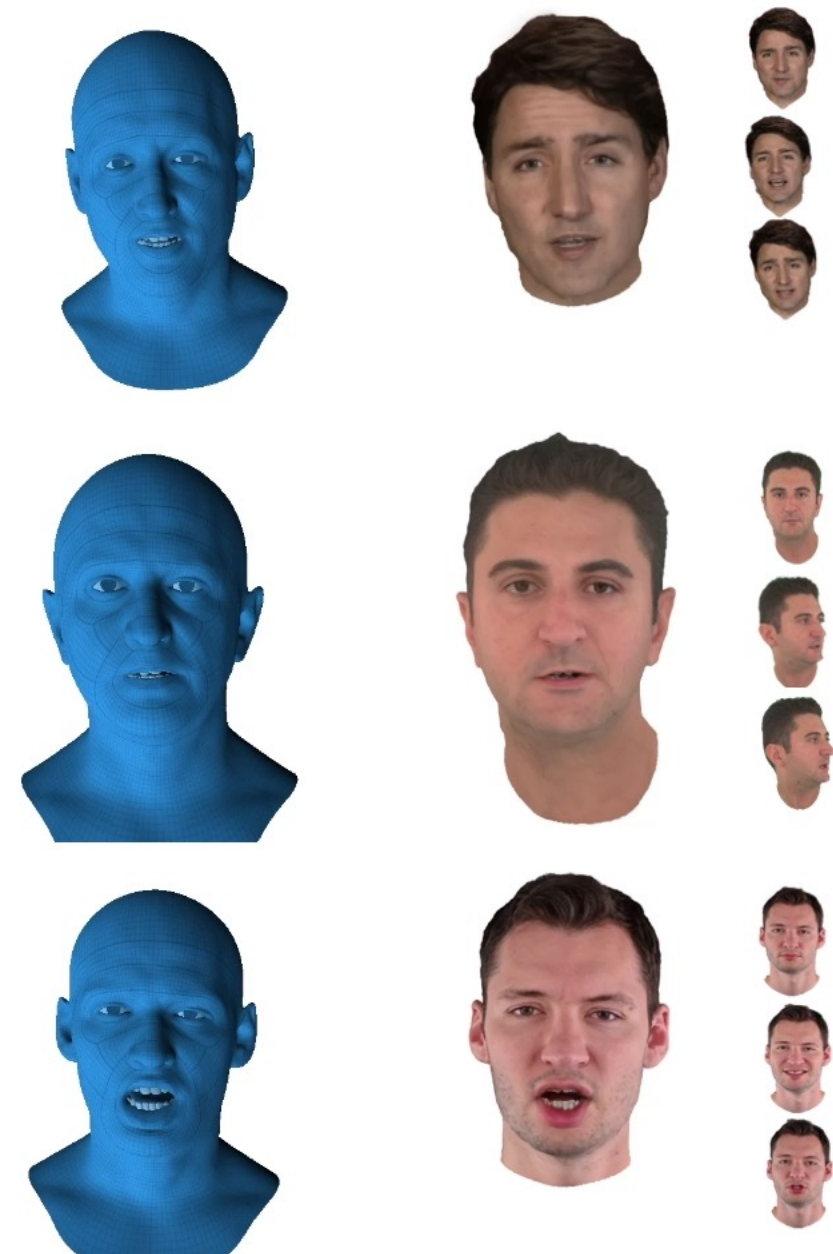
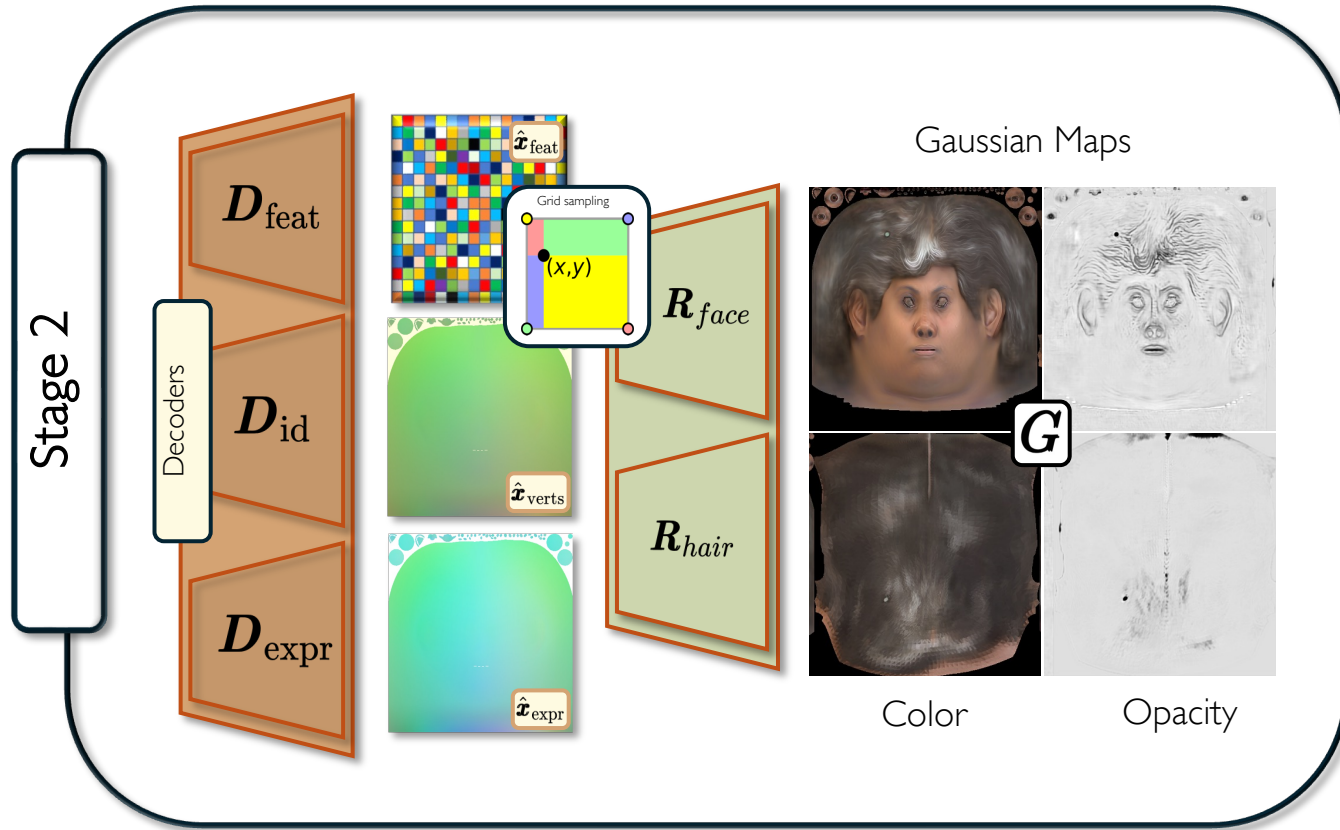
Input



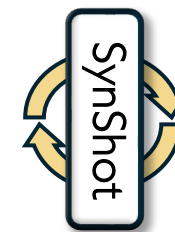
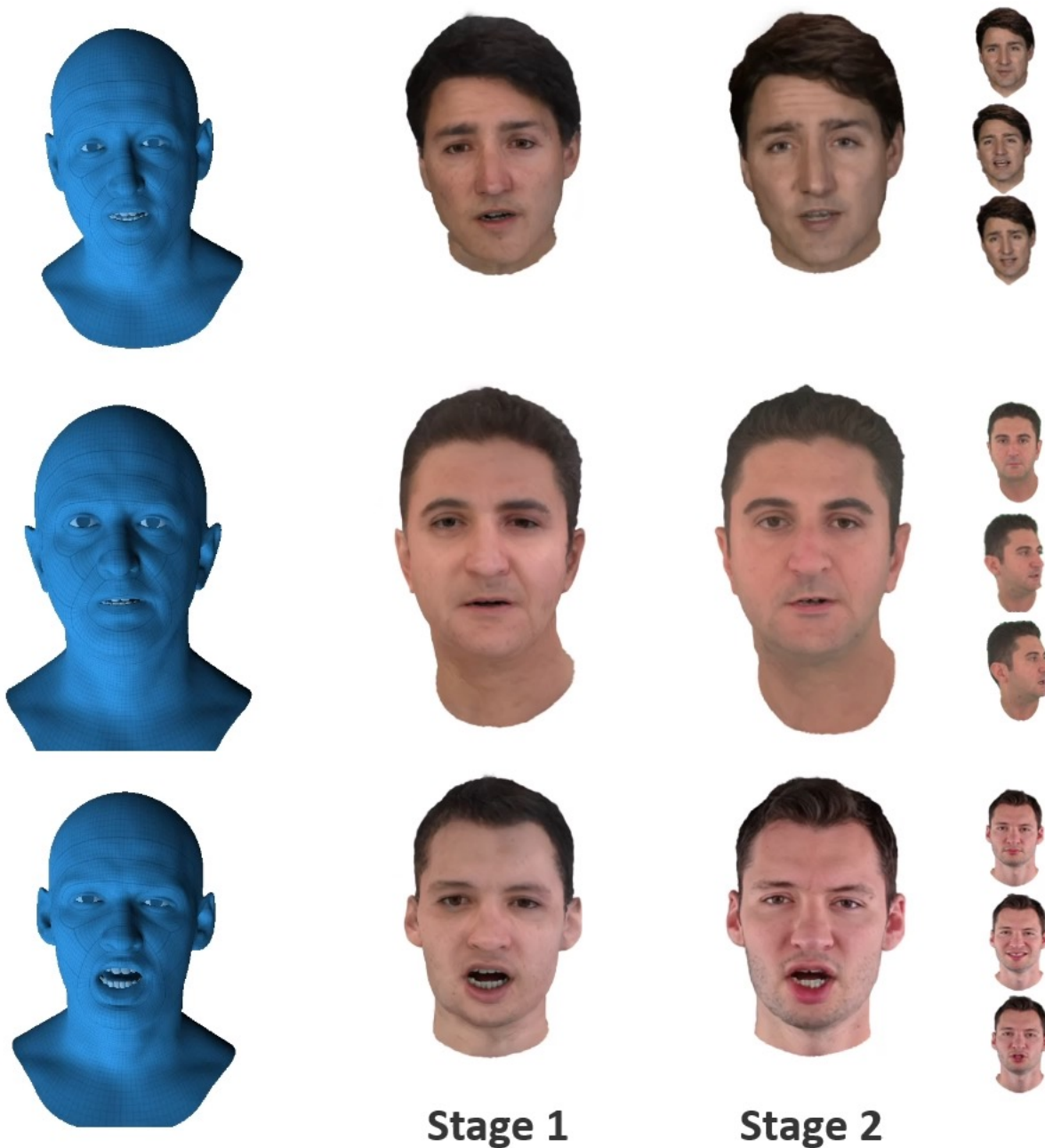
Inversion – Stage 1



Inversion – Stage 2



Inversion – Stage 1-2



Results: Personalized Baselines



Source



Ours



INSTA¹



Splatting Avatar²



Flash Avatar²

1) Zielonka *et al.* Instant Volumetric Head Avatars

2) Xiang *et al.* FlashAvatar: High-fidelity Head Avatar with Efficient Gaussian Embedding

3) Zhijing *et al.* SplattingAvatar: Realistic Real-Time Human Avatars with Mesh-Embedded Gaussian Splatting

Results: GAN-based Baselines



Source



Ours



InvertAvatar¹



Portrait4D²



Next3D³

1) Zhao *et al.* InvertAvatar: Incremental GAN Inversion for Generalized Head Avatars

2) Deng *et al.* Portrait4D: Learning One-Shot 4D Head Avatar Synthesis using Synthetic Data

3) Sun *et al.* Next3D: Generative Neural Texture Rasterization for 3D-Aware Head Avatars

Results: Three-shot Inversions



Results: Three-shot Inversions



Limitations



Missing hair



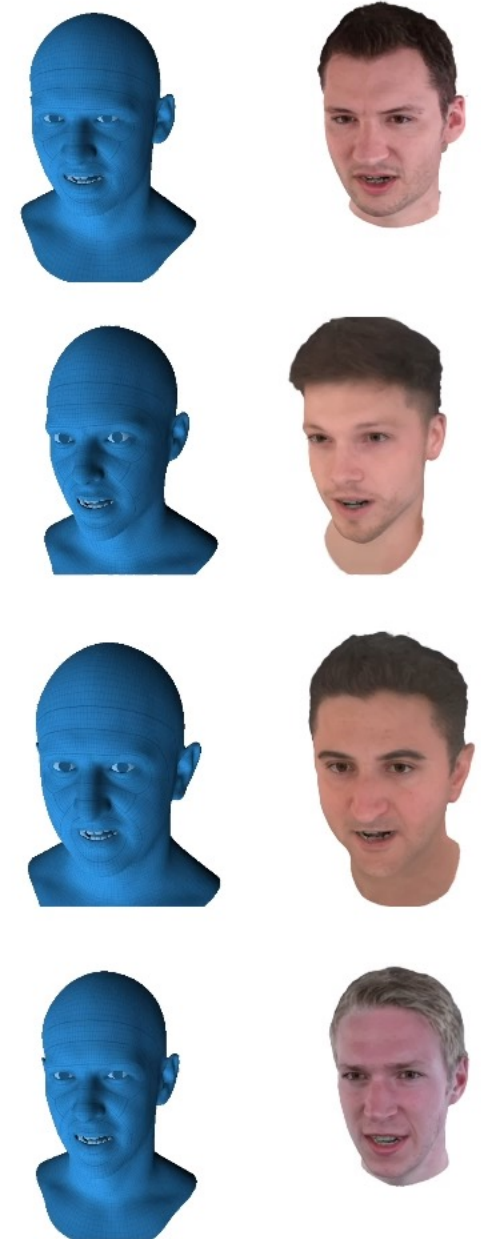
No glasses in the dataset



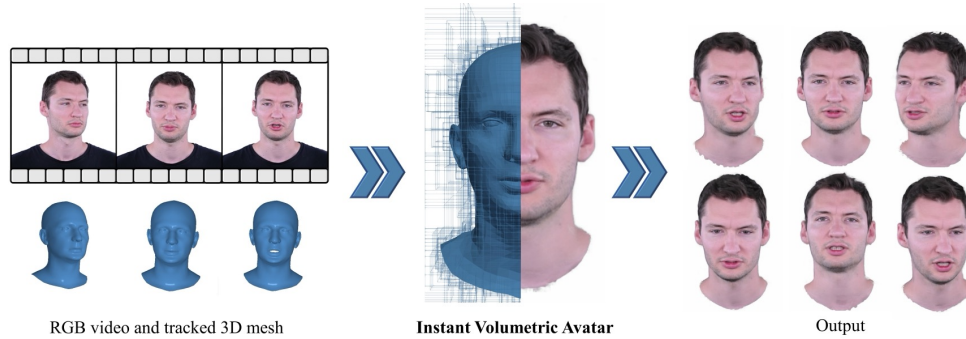
Wrong illumination

Take-home Messages of SynShot

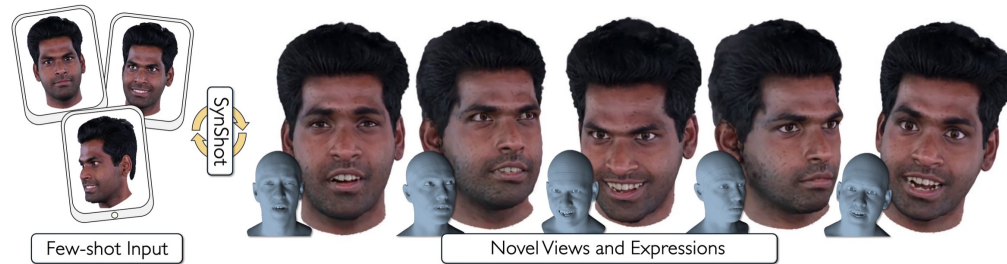
1. **SynShot** builds a multi-view prior using only *synthetic* data.
2. Enables cross-reenactment and outperforms monocular personalized methods like **INSTA**.
3. Pivotal fine-tuning bridges the real-synthetic gap, enabling drivable 4D avatars from synthetic-only priors.



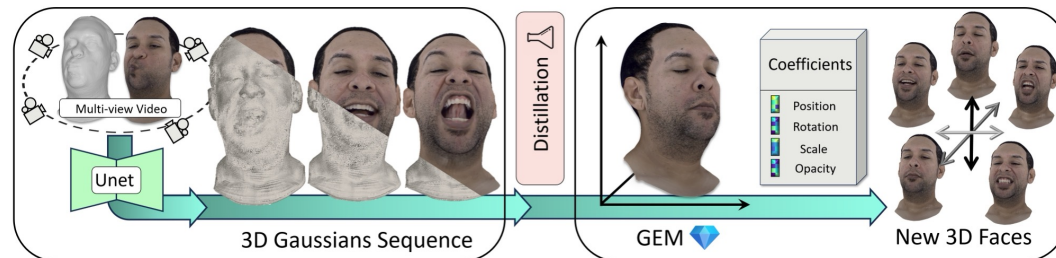
Presentation Outline



INSTA - Instant Volumetric Head Avatars
[Zielonka, Bolkart, Thies]
CVPR'23



SynShot - Synthetic Prior for Few-Shot Drivable Avatar Inversion
[Zielonka, Garbin, Lattas, Kopanas, Gotardo, Beeler, Thies, Bolkart]
CVPR'25



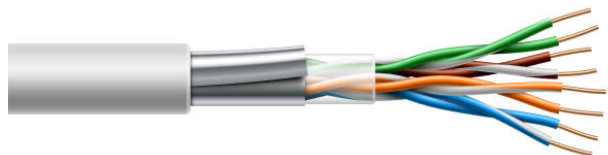
GEM - Gaussian Eigen Models for Human Heads
[Zielonka, Bolkart, Beeler, Thies]
CVPR'25

Motivation: High-Quality Avatars on VR-Glasses

Place A



GEM  Avatar

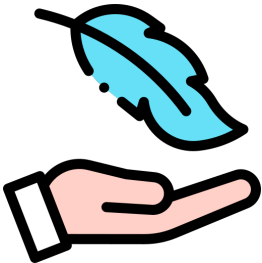
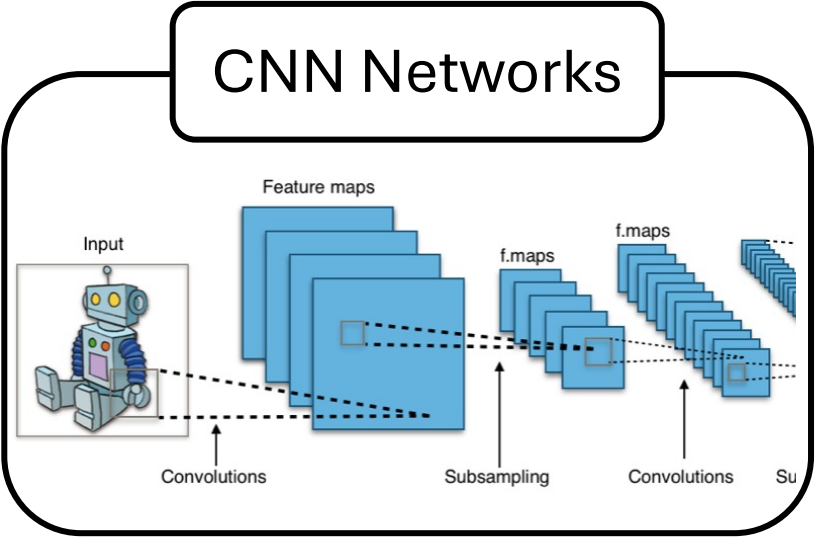


Place B



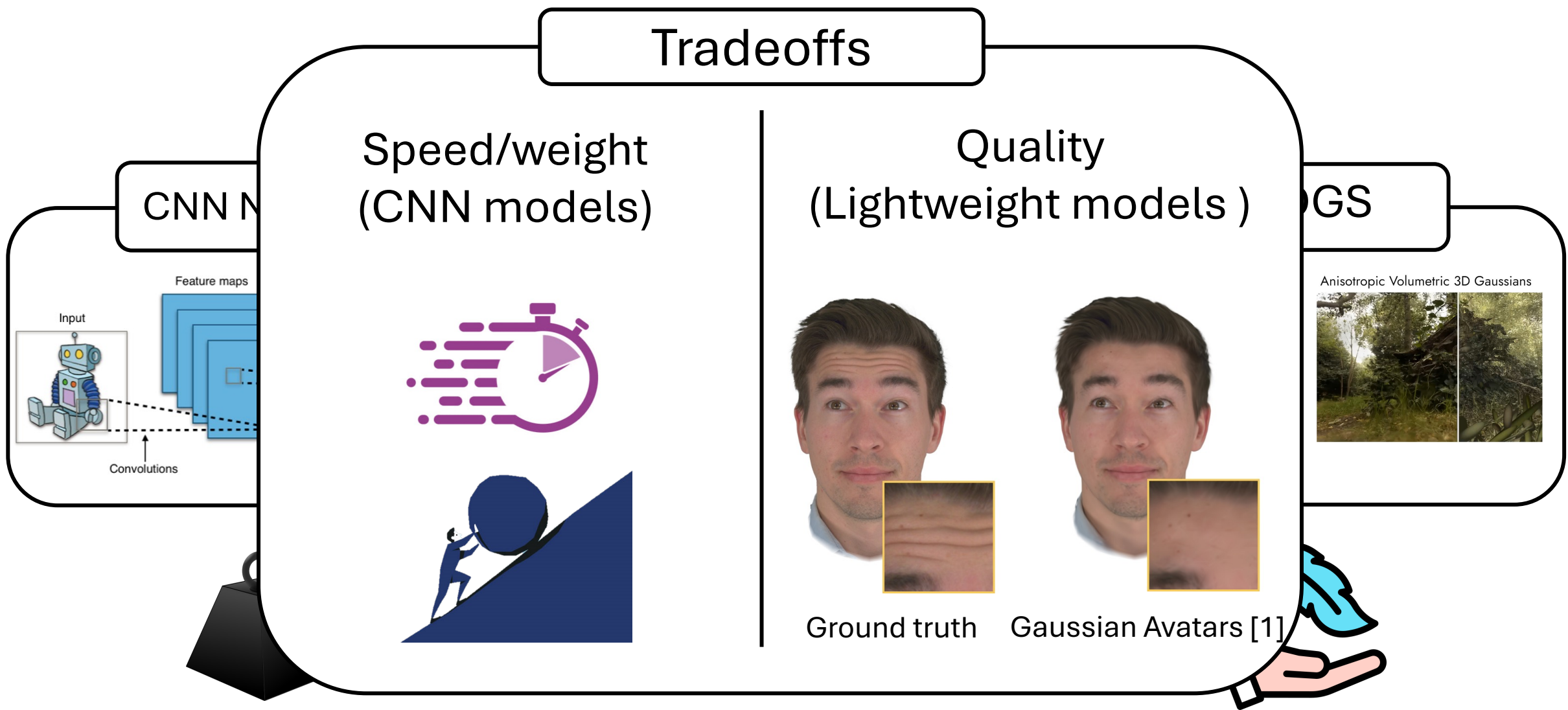
GEM  Avatar

Motivation: There is no Free Lunch



1) Qian et al., *Photorealistic Head Avatars with Rigged 3D Gaussians*, CVPR 2024

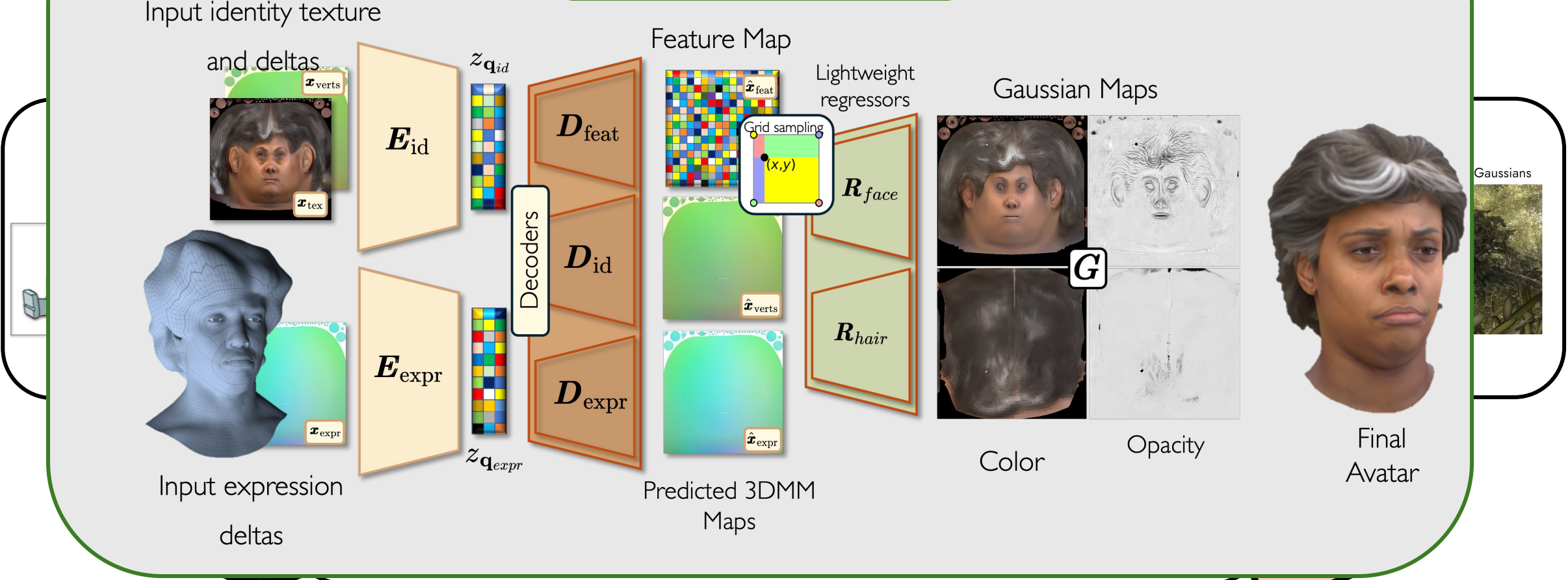
Motivation: There is no Free Lunch



1) Qian et al., *Photorealistic Head Avatars with Rigged 3D Gaussians*, CVPR 2024

Motivation: There is no Free Lunch

SynShot

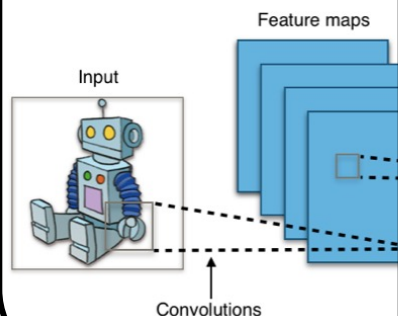


1) Qian et al., *Photorealistic Head Avatars with Rigged 3D Gaussians*, CVPR 2024





Motivation

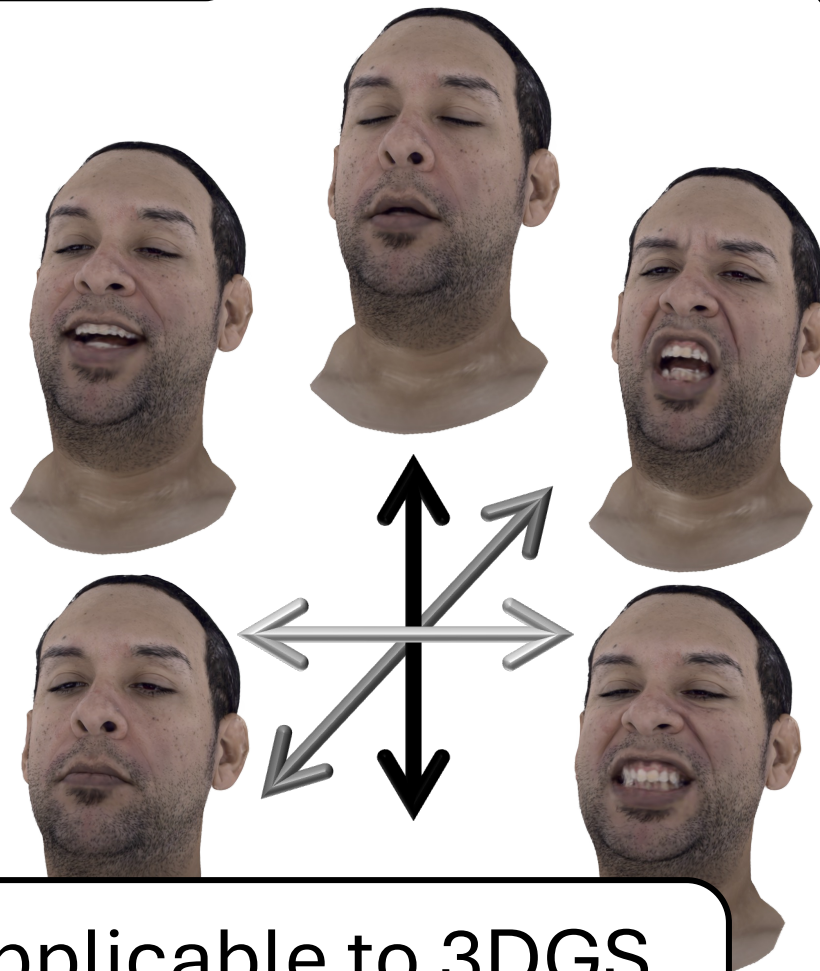
GEM 

CNN N



Coefficients

-  Position
-  Rotation
-  Scale
-  Opacity



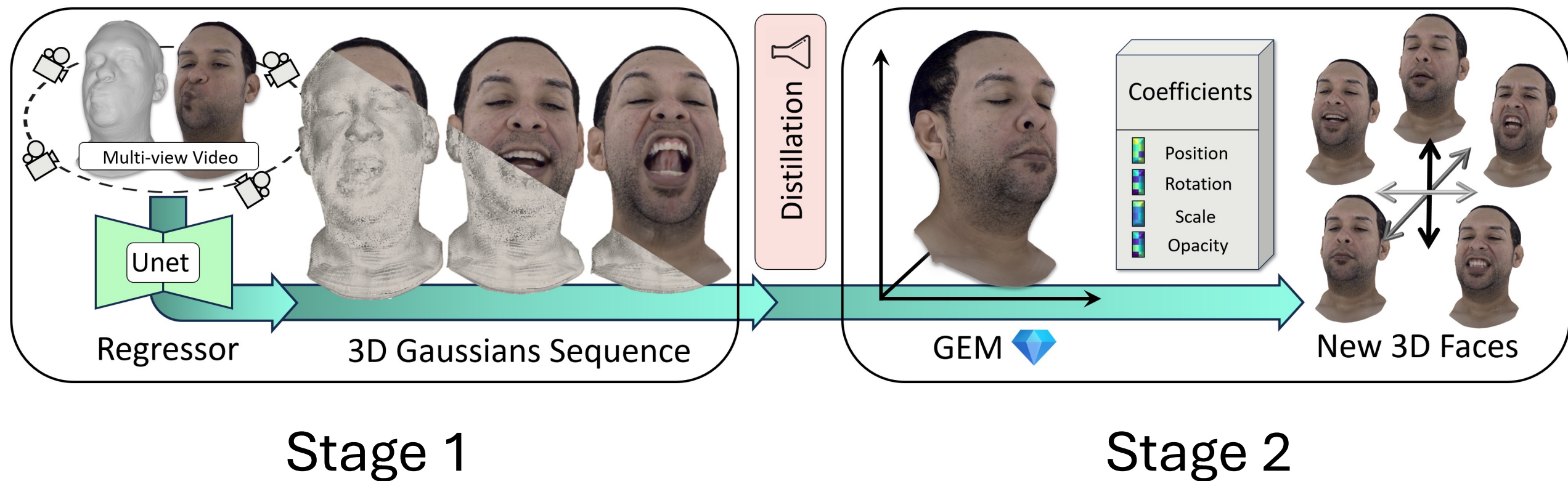
DGS



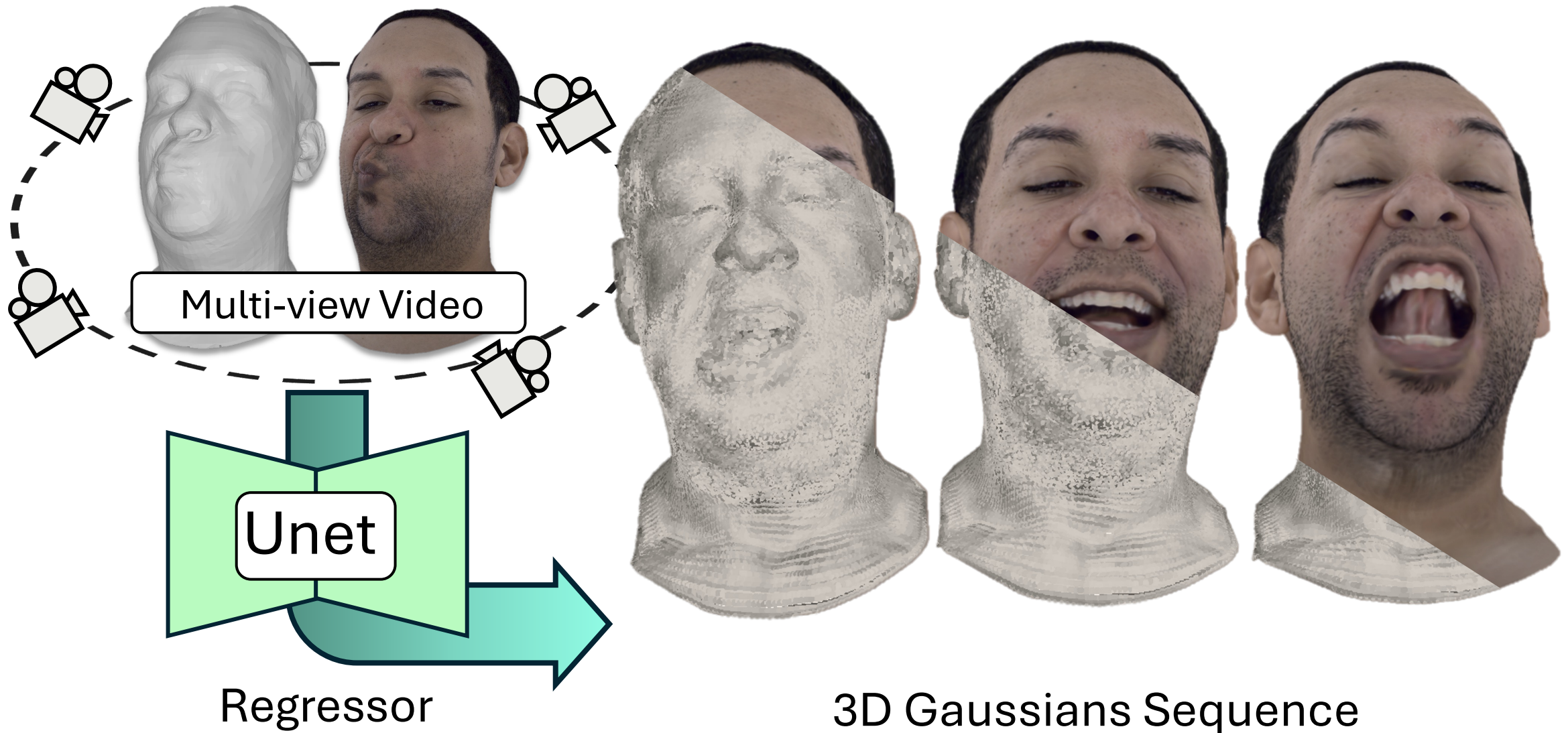
Distillation applicable to 3DGS
head avatar methods

1) Qian et al., *Photorealistic Head Avatars*

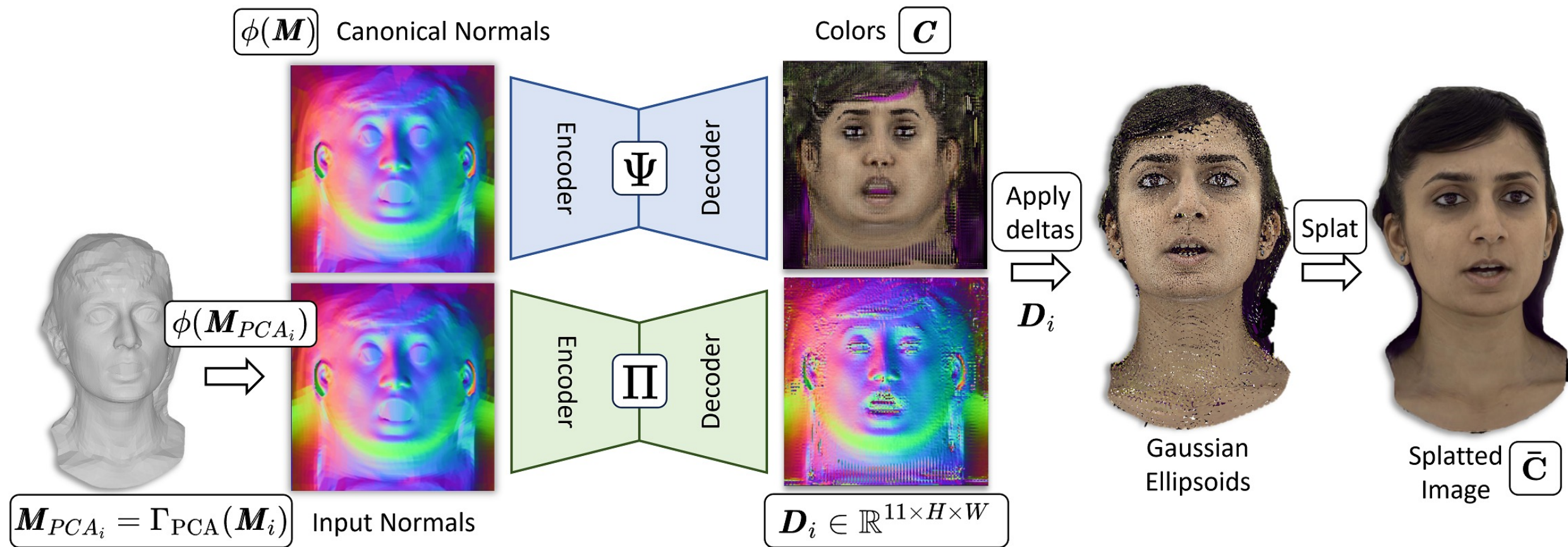
Method: Two Stage Approach



Method: Stage One (Generator)



Method: Stage One (Sequence of Gaussians)



Animatable Gaussians [Li et al.]

Method: Stage One (Sequence of Gaussians)

SynShot as Generator

Input identity texture

and deltas



E_{id}

$z_{q_{id}}$



D_{feat}

D_{id}

D_{expr}

Decoders

E_{expr}

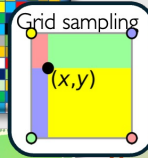


$z_{q_{expr}}$

Feature Map

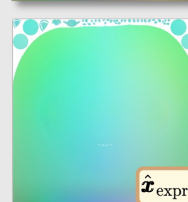
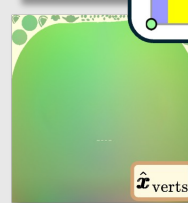


Lightweight regressors



R_{face}

R_{hair}



Predicted 3DMM Maps

Gaussian Maps



G



Color

Opacity



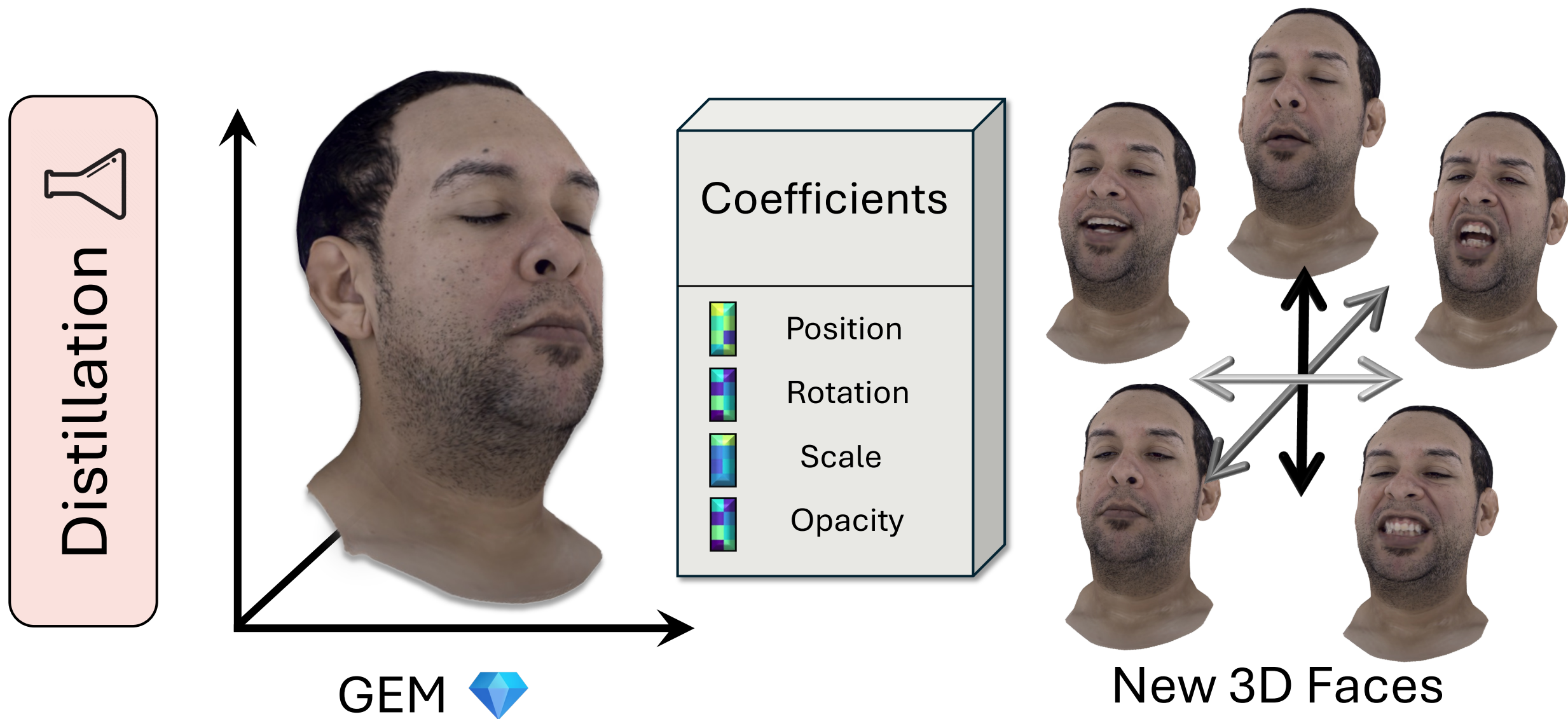
Final Avatar

\bar{C}

M_{PCA}

Input expression
deltas

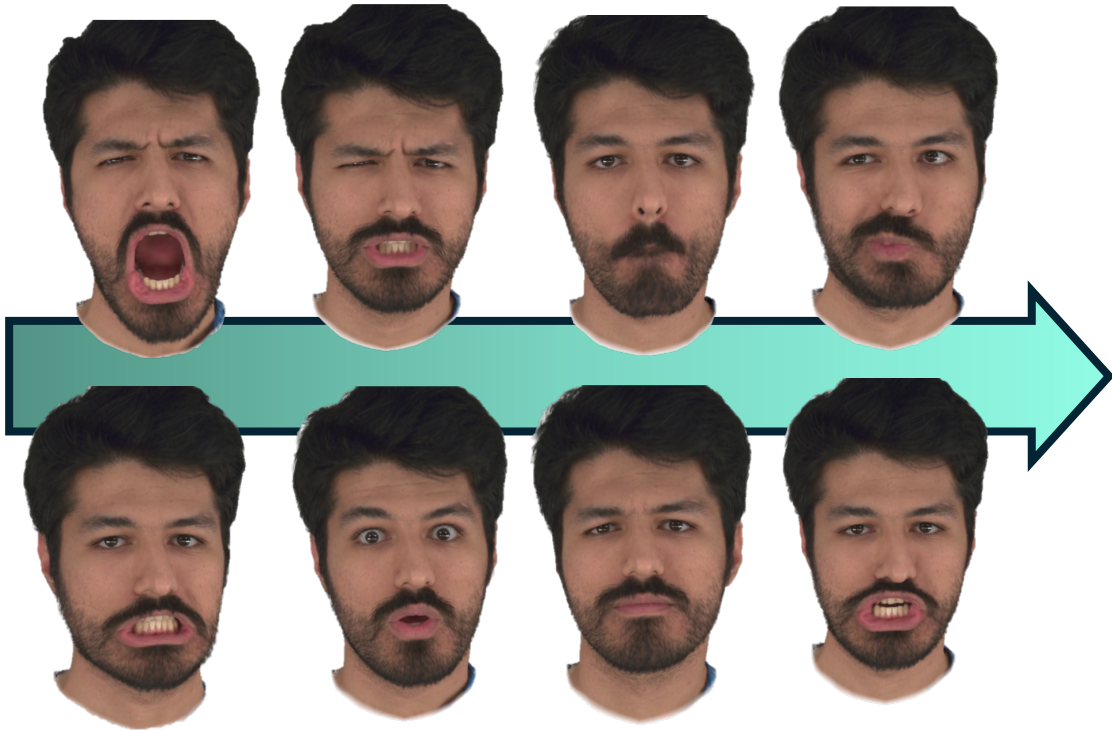
Method: Stage Two (Distillation)



Method: Stage Two

Distilling

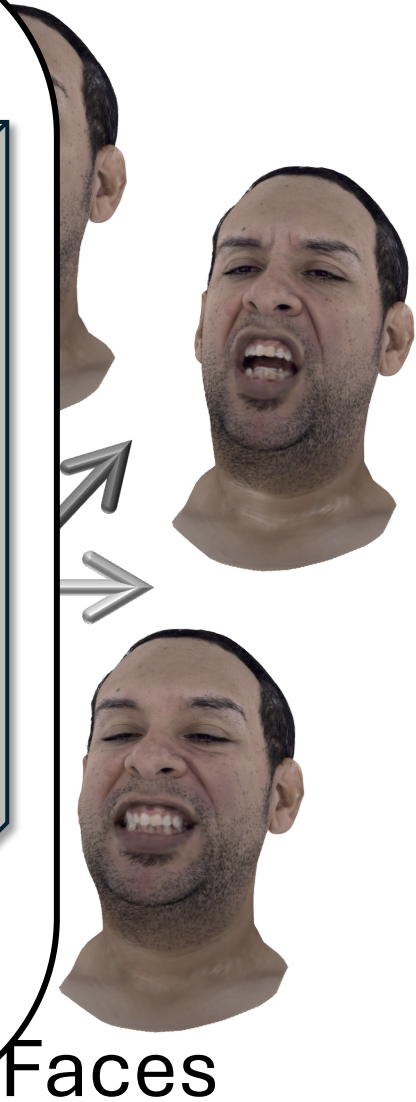
Distillation



Dataset of regressed
Gaussians primitives

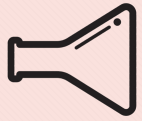
Coefficients

	Position
	Rotation
	Scale
	Opacity



Method: Stage Two

Ensamble of eigen basis



Distillation



GEOMETRY COMP 00 = -2.6



Position



OPACITY COMP 00 = -2.6



Opacity



ROTATION COMP 00 = -2.6



Rotation



SCALES COMP 01 = -2.6



Scale

Space traversal $[-3\sigma, 3\sigma]$



Faces

Results: Novel Expressions

AG – Animatable Gaussians [Li et al.]

GA – Gaussian Avatars [Qian et al.]

INSTA – [Zielonka et al.]

🔥 Needs a 3DMM like FLAME



Ground Truth



Ours  GEM



Ours Net 🔥



AG 🔥

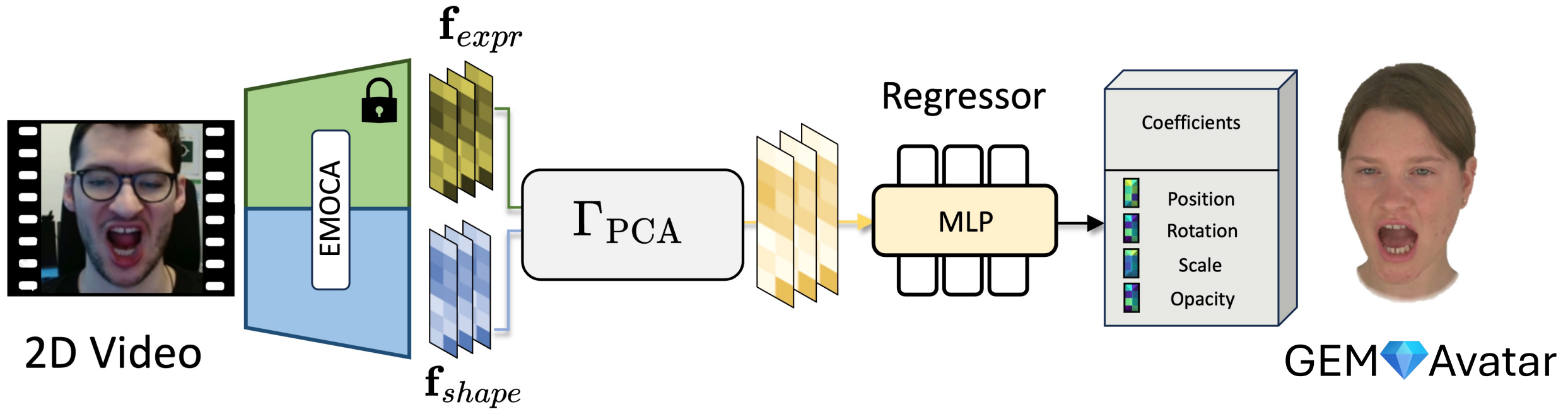


GA 🔥

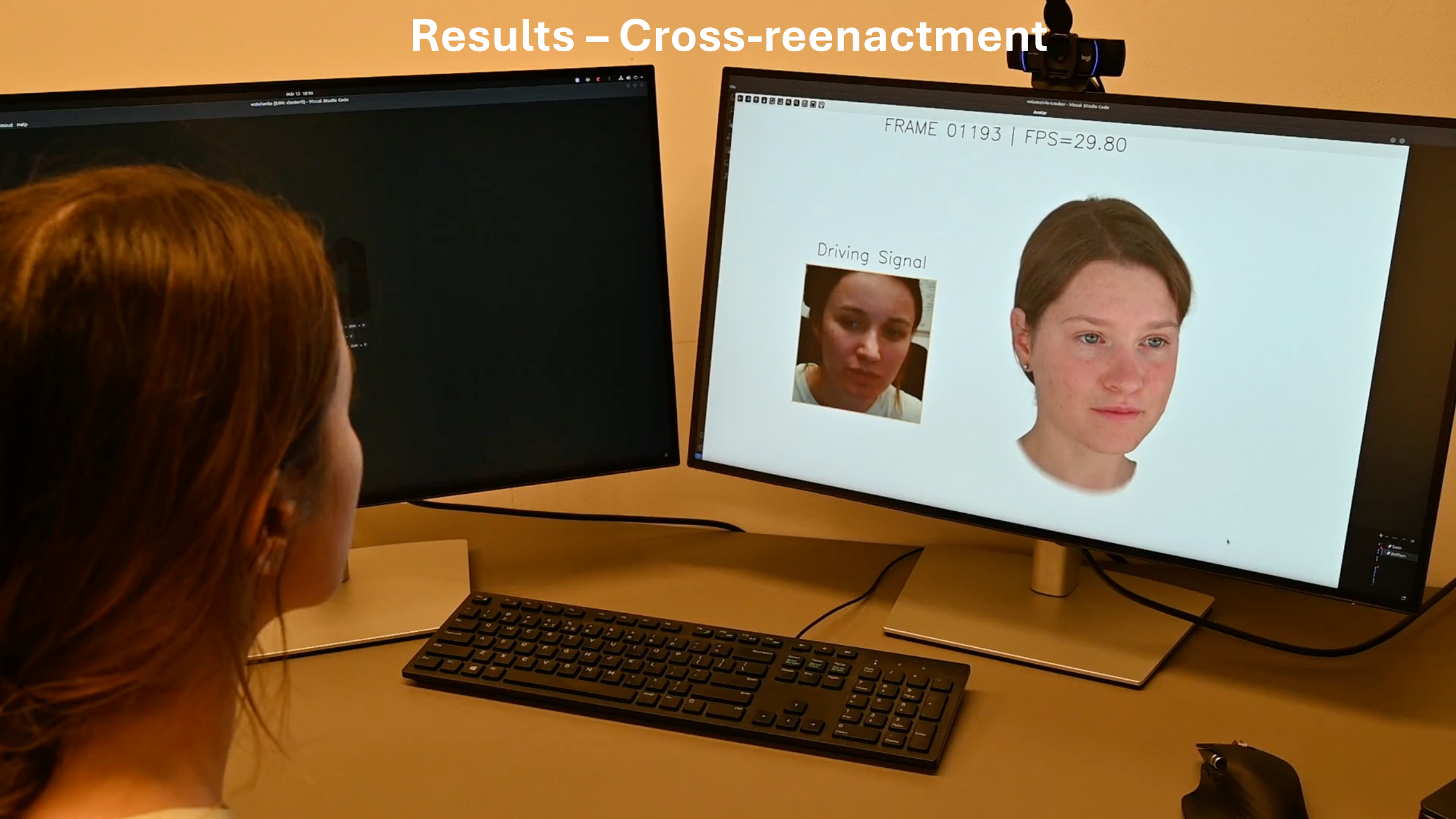


INSTA 🔥

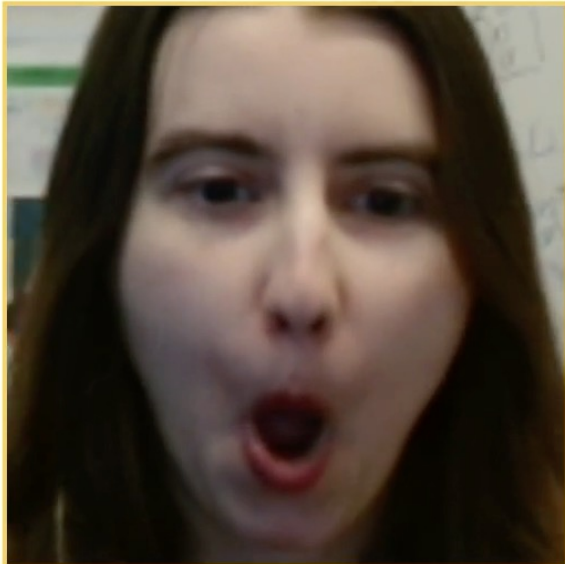
Results: Image-base Cross-reenactment



Results – Cross-reenactment



Driving Signal



Driving Signal



Limitations: PCA's Global Extend



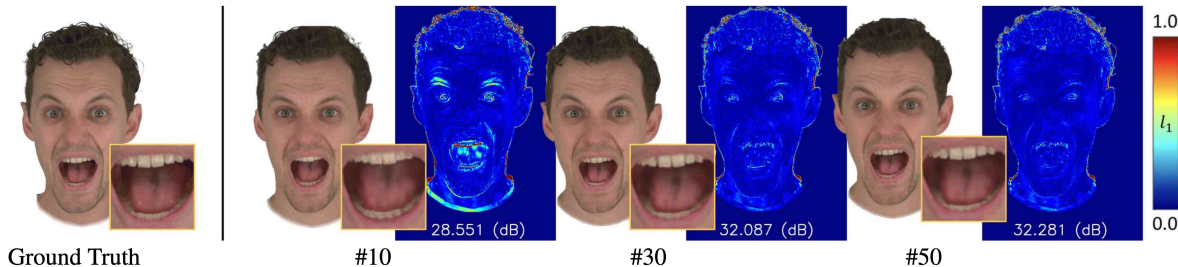
Source



GEM

Take-home Messages of GEM

1. GEM distills heavy CNN-based avatars into efficient linear Eigenbases without needing 3DMM.
2. Quality and memory are adjustable via the number of bases.
3. Enables real-time image-based cross-reenactment with a pretrained ResNet.



All Projects During my PhD



SynShot - Synthetic Prior for Few-Shot Drivable Head Avatar Inversion [CVPR25]



RGB video and



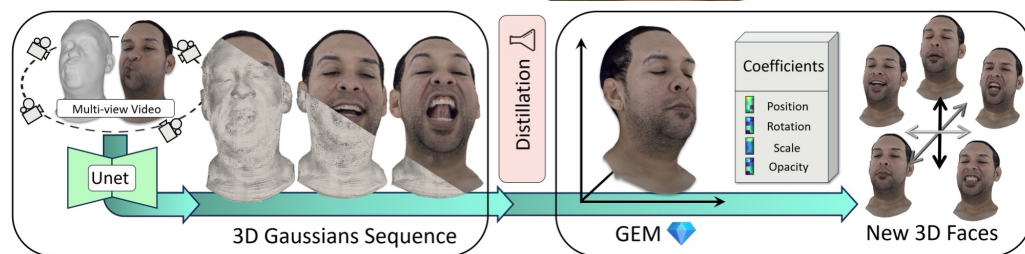
INSTA - Instant Volumetric Head Avatars [CVPR23]



MICA - Towards Metrical Reconstruction of Human Faces [ECCV22]

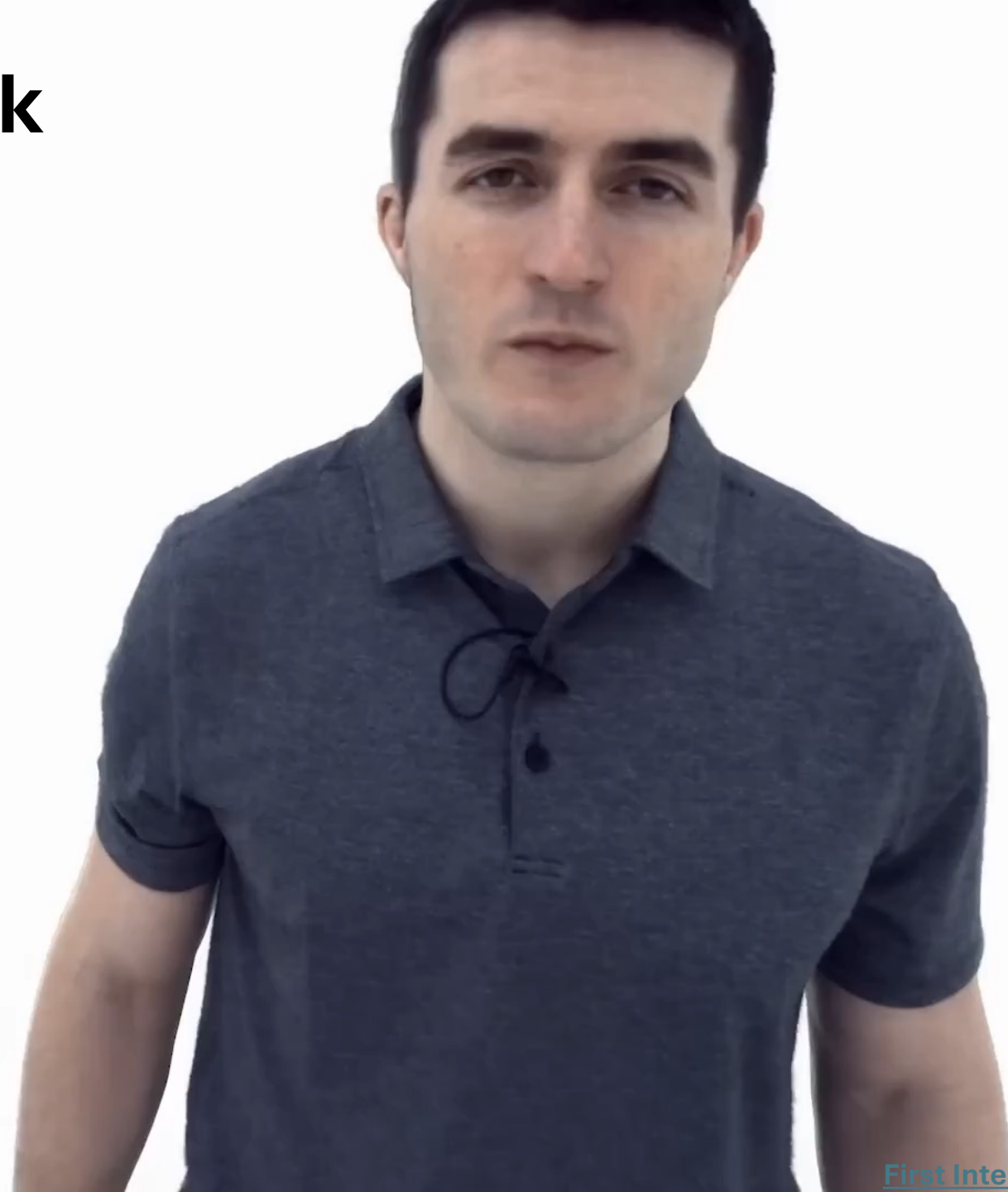


D3GA - Drivable 3D Gaussian Avatars [3DV25]



GEM - Gaussian Eigen Models for Human Heads [CVPR25]

Future Work







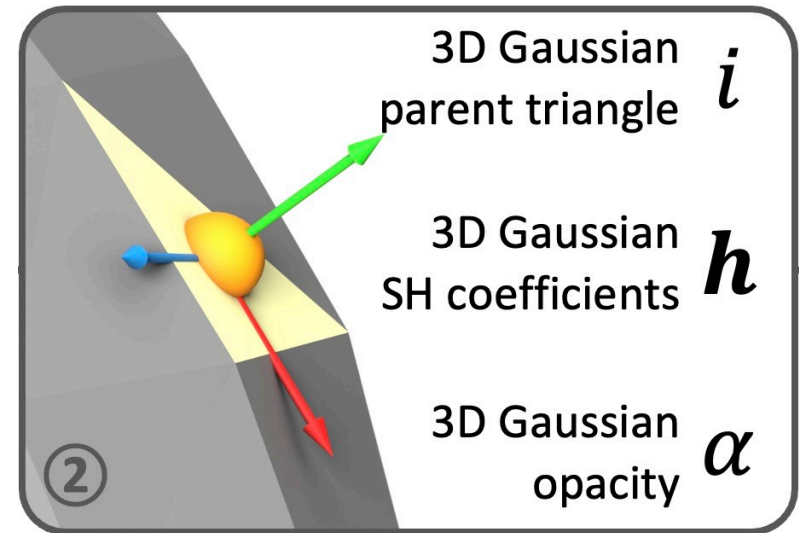
Thank you!





Additional Slides

(3DGS + FLAME)*



assign a 3D Gaussian at the center of each triangle

Results: Qualitative Comparison (Novel View)



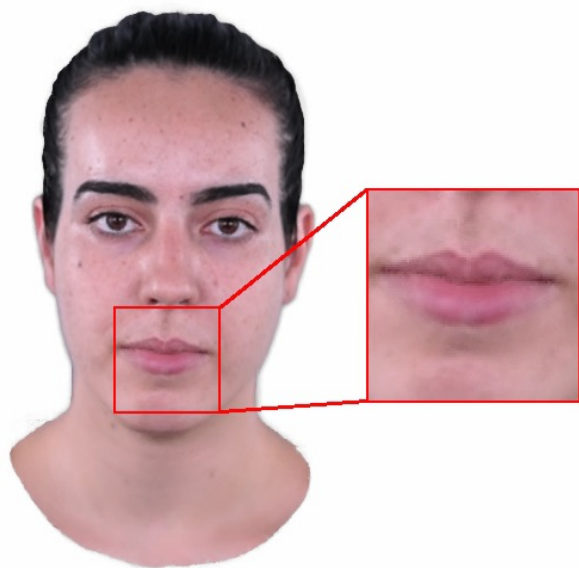
NHA [Grassal et al.]

NeRFace [Gafni et al.]

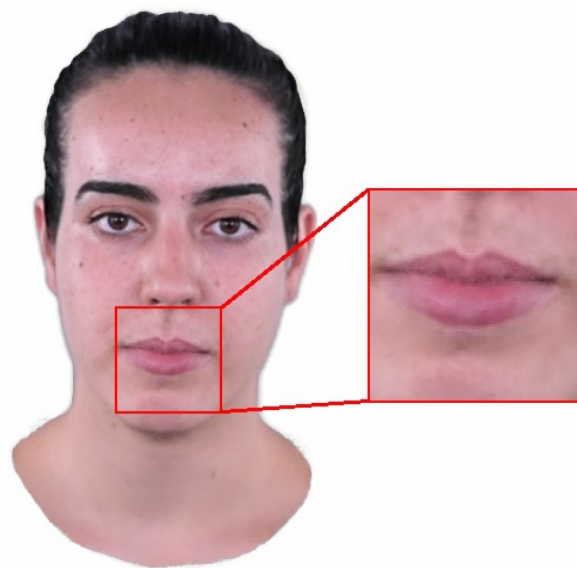
IMAvatar [Zheng et al.]

Ours

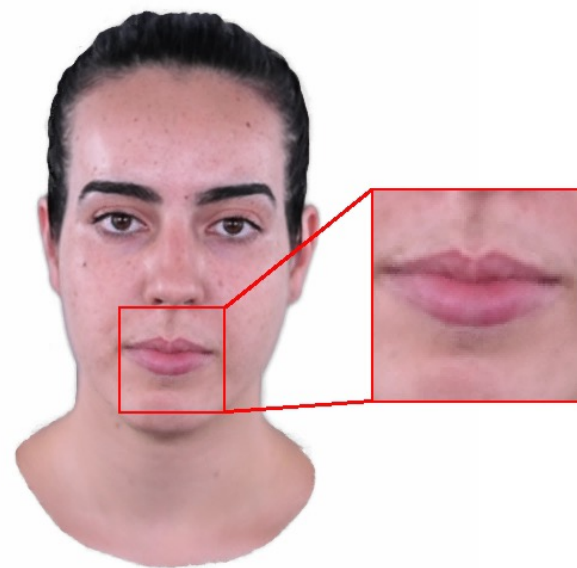
Results: Ablation Studies



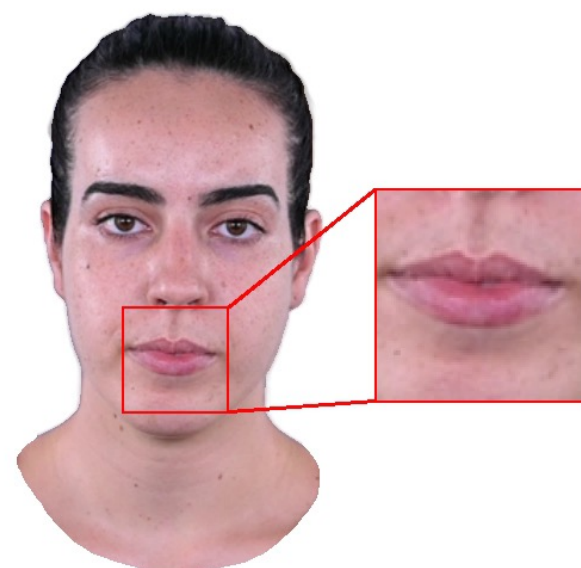
w/o mouth
weighting



w/o conditioning

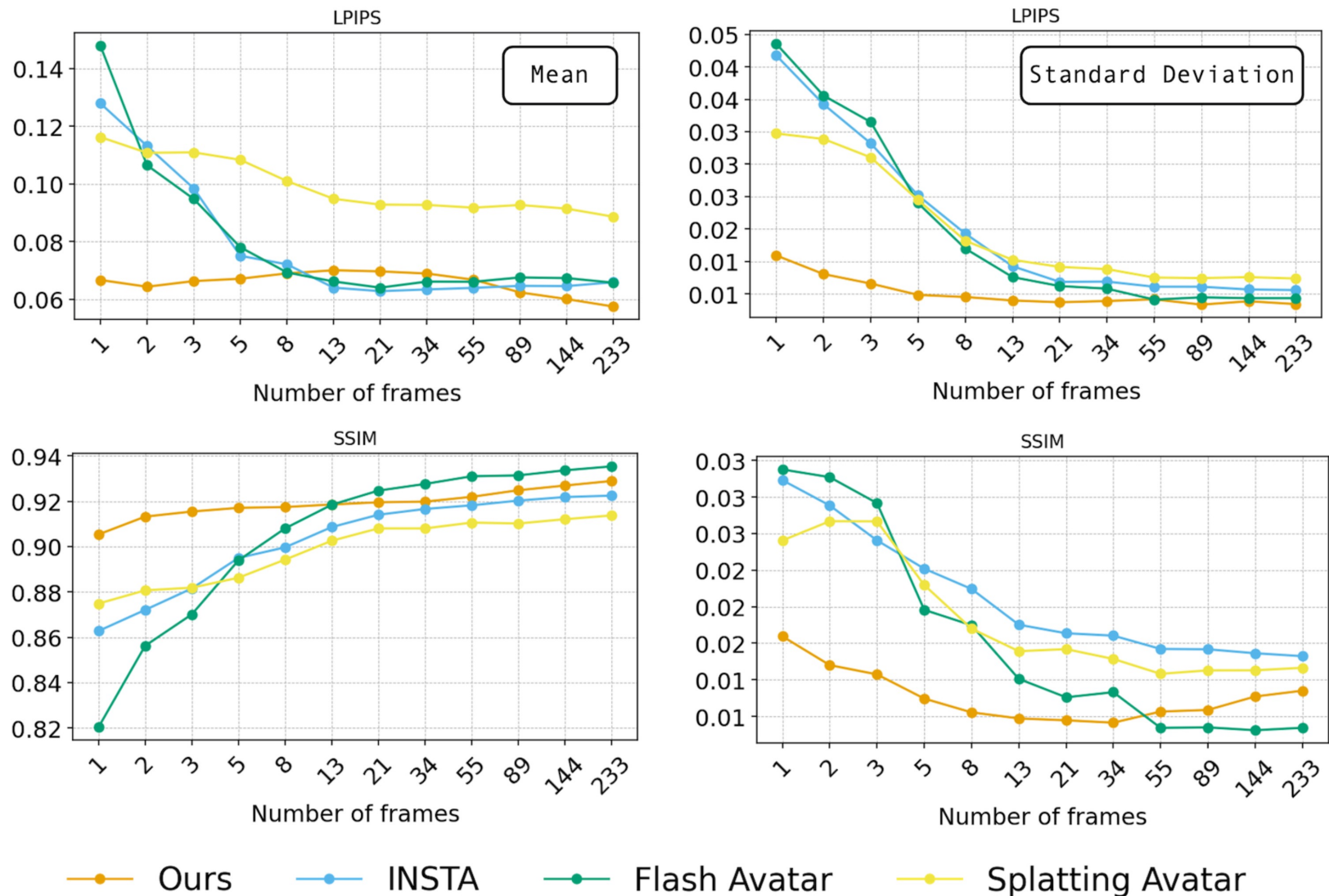


Final



Ground Truth

Results: Quantitative Comparison



Results: Ablation Study (Architecture)

Architecture	L1 ↓	LPIPS ↓	SSIM ↑	PSNR ↑
$F = 128$	0.0356	0.2686	0.8189	20.1536
Tex. up-sampling	0.0352	0.2695	0.8196	20.1909
Single Layer	0.0369	0.2702	0.8177	19.8871
$F = 32$	0.0375	0.2732	0.8146	19.7002
w/o VQ	0.0396	0.2747	0.8122	19.2861
$F = 64$	0.0400	0.2765	0.8104	19.2731
No Sampling	0.0403	0.2853	0.8158	19.9787
256×256	0.0365	0.2865	0.8194	20.4010

Results: Novel Expressions

Method	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	L1 \downarrow
AG	29.0114	0.0812	0.9429	0.0099
GA	28.3137	0.0815	0.9433	0.0102
INSTA	27.9181	0.1153	0.9340	0.0128
Ours Net	29.2454	0.0777	0.9448	0.0096
Ours GEM	32.6781	0.0675	0.9633	0.0069

AG – Animatable Gaussians [Li et al.]
GA – Gaussian Avatars [Qian et al.]
INSTA – [Zielonka et al.]

Results: Novel Viewpoint

Method	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	L1 \downarrow
AG	32.4166	0.0712	0.9614	0.0066
GA	31.3197	0.0786	0.9567	0.0075
INSTA	27.7786	0.1232	0.9294	0.0163
Ours Net	32.4622	0.0713	0.9617	0.0067
Ours GEM	33.5528	0.0678	0.9662	0.0061

AG – Animatable Gaussians [Li et al.]
GA – Gaussian Avatars [Qian et al.]
INSTA – [Zielonka et al.]

GEM: Localized PCA



GEOMETRY/eyeballs COMP 00 = -3.0

GEM: Localized PCA



Source



GEM



Localized GEM