# ID-to-3D

## Expressive ID-guided 3D Heads via Score Distillation Sampling .

F. Babiloni, A.Lattas, J.Deng, S.Zafeiriou

Imperial College London

# ID-to-3D: Expressive ID-guided 3D Heads via Score Distillation Sampling.

**ArcFace-Conditioned 3D Head Asset Generation using SDS**

- **Text-to-2D-Normals and Albedo Models**:
  We present a novel approach to creating ID-conditioned and expression-conditioned text-to-image models capable of generating realistically plausible normals and albedo images from a small set of 3D assets.

- **ID-Conditioned Expressive Heads**:
  We propose a neural parametric representation for the expressions that creates up to 13 unique and ID-consistent expressions captured by latent codes and associated with a set of 3D assets with separate geometric, albedo, and material information.

"In-the-wild" images

ID-consistent 3D assets

Score Distillation Sampling Pipeline Optimization Objective

$$\min_{\theta} D_{KL}(q^{\theta}(\mathbf{x}_0|y_{\text{text}}) \,\|\, p(\mathbf{x}_0|y_{\text{text}})).$$

$\underbrace{\qquad\qquad}_{\text{distribution of renderings}}$  $\underbrace{\qquad\qquad}_{\text{target distribution}}$

⊗ Target distribution drift.
⊗ General guidance model does not separate texture/geometry
⊗ Using textual prompt lacks granularity.

ID to 3D Optimization Objective

$$\min_{\theta_g, \theta_a} D_{KL}\left(q^{\theta_g}(\mathbf{z}_0^n | \mathbf{c}, y_{\text{text}}, y_{\text{exp}}, y_{\text{id}}) \,\|\, p(\mathbf{z}_0^n | \mathbf{c}, y_{\text{text}}, y_{\text{exp}}, y_{\text{id}})\right)$$

*geometry generation objective*

$$+ D_{KL}\left(q^{\theta_a}(\mathbf{z}_0^a | \mathbf{c}, \mathbf{l}, y_{\text{text}}, y_{\text{exp}}, y_{\text{id}}) \,\|\, p(\mathbf{z}_0^a | \mathbf{c}, y_{\text{text}}, y_{\text{exp}}, y_{\text{id}})\right).$$

*texture generation objective*

$\theta_g$ parameterization of the 3D geometry,
$\theta_a$ parameterization of the 3D textures
$\mathbf{z}_0^a$ albedo textures
$\mathbf{z}_0^n$ normal maps
$\mathbf{l}\ \mathbf{c}$ lighting camera
$y_{\text{text}}$ textual
$y_{\text{id}}$ identity
$y_{\text{exp}}$ expression

- Geometry / Texture adapted self-attention:

$$\mathbf{Z}_{\mathbf{SA}}^n = \text{Att}(\mathbf{Q}^n, \mathbf{K}^n, \mathbf{V}^n), \quad \mathbf{Q}^n = \mathbf{X}\mathbf{W}_Q + \mathbf{X}\mathbf{W}_Q^n, \mathbf{K}^n = \mathbf{X}\mathbf{W}_K + \mathbf{X}\mathbf{W}_K^n, \mathbf{V}^n = \mathbf{X}\mathbf{W}_V + \mathbf{X}\mathbf{W}_V^n$$

- Identity and expression adapted cross-attention

$$\mathbf{Z}_{\mathbf{CA}}^n = \text{Att}(\mathbf{Q}, \mathbf{K}^{\mathbf{text}}, \mathbf{V}^{\mathbf{text}}) + \lambda_{id} \cdot \text{Att}(\mathbf{Q}, \mathbf{K}^{\mathbf{id}}, \mathbf{V}^{\mathbf{id}}) + \lambda_{exp} \cdot \text{Att}(\mathbf{Q}, \mathbf{K}^{\mathbf{exp}}, \mathbf{V}^{\mathbf{exp}})$$

- Efficient Finetuning

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{z}_0^n, \boldsymbol{\epsilon}, \mathbf{c}, t, \mathbf{y}_{\text{text}}, \mathbf{y}_{\text{id}}, \mathbf{y}_{\text{exp}}} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\phi_g}\left(\mathbf{z}_t^n, t, \mathbf{c}, \mathbf{y}_{\text{text}}, \mathbf{y}_{\text{id}}, \mathbf{y}_{\text{exp}}\right) \right\|^2$$
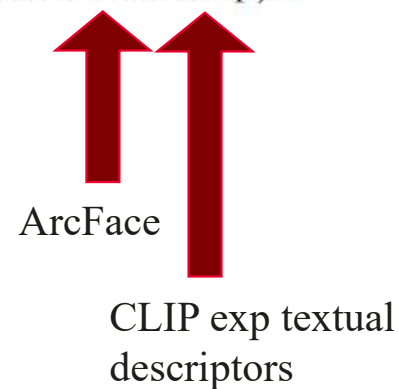
$\phi_g$   2D prior for geometry
$\phi_a$   2D prior for textures
$\mathbf{y}_{\text{id}}$   Arcface identity embeddings
$\mathbf{y}_{\text{exp}}$   CLIP embeddings textual descriptor of FaceWareHouse
$\mathbf{c}$   textual descriptor of the camera poses

ArcFace

CLIP exp textual descriptors

- Geometry Generation

$$\nabla \mathcal{L}_{SDS}(\theta_g) = \mathbb{E}_{\mathbf{c},t,\epsilon}\left[\omega(t)(\epsilon_{\phi_g}(\mathbf{z}_t^n, \mathbf{y}_{id}, \mathbf{y}_{exp}, \mathbf{y}_{text}, t) - \epsilon)\frac{\partial g(\theta_g, \mathbf{c})}{\partial \theta_g}\right]$$

$$\theta_g = [\mathbf{k}_{exp}^n, \psi_g]$$

$\theta_g$ geometry model

$\Psi_g(\Gamma, \mathbf{k}_{exp}^n)$ Transformer DMTET

$\mathbf{k}_{exp}^n \in \mathbb{R}^{d_{exp}}$ latent code

$\Gamma$ deformable tetrahedral grid

$g(\theta_g, \mathbf{c})$ rendered normals

$\phi_g$ 2D prior for geometry

- Appearance Generation

$$\nabla \mathcal{L}_{SDS}(\theta_a) = \mathbb{E}_{\mathbf{c},t,\epsilon}\left[\omega(t)(\epsilon_{\phi_a}(\mathbf{z}_t^a, \mathbf{y}_{id}, \mathbf{y}_{exp}, \mathbf{y}_{text}, t) - \epsilon)\frac{\partial g(\theta_a, \mathbf{c}, \mathbf{l})}{\partial \theta_a}\right]$$

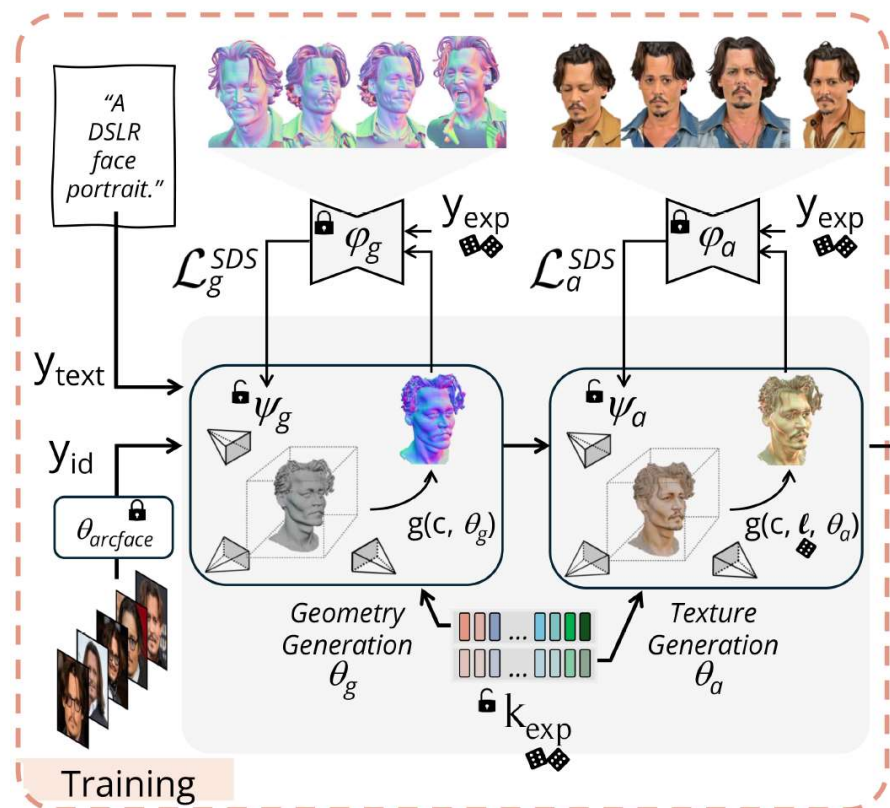$$\theta_a = [\mathbf{k}_{exp}^a, \psi_a]$$

$\theta_a$ texture appearance

$\Psi_a(\mathbf{k}_{exp}^a)$ Transformer  roughness albedo specularity

$\mathbf{k}_{exp}^a \in R^{d_{exp}}$ latent code

$g(\theta_a, \mathbf{c}, \mathbf{l})$ rendered pseudo-albedo

$\phi_a$ 2D prior for textures

# Results: Comparison with SotA SDS methods
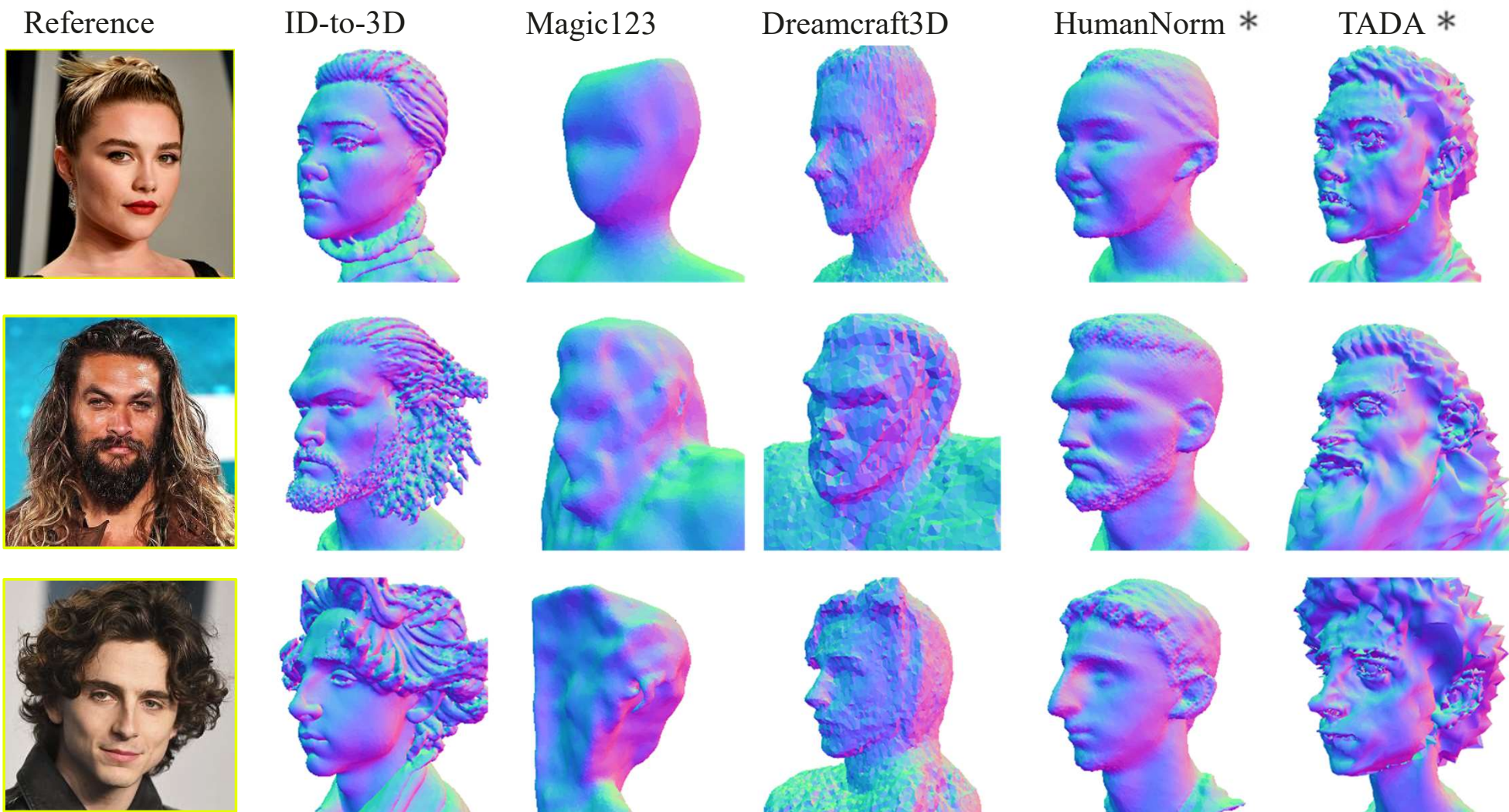


Reference     ID-to-3D     Magic123     Dreamcraft3D     HumanNorm ✳     TADA ✳

# Results: Comparison with SotA SDS methods



| Reference | ID-to-3D | Magic123 | Dreamcraft3D | HumanNorm ✳ | TADA ✳ |

Benedict
Cumberbatch

Will
Smith

Donald
Trump

Kim
Kardashian

Lucy
Liu

"...as a cute baby"

*...as a male*    *...as a male*

*...as a female*    *...as a female*

Geometry Editing

*...blonde hair*    *...brown hair*

*...green hair*    *...green hair*

Texture Editing

(a) Eyes closed    (b) Brow raised    (c) Brow lowerer    (d) Cheeks puffed

(e) Grin    (f) Dimpler    (g) Lip roll    (h) Squeeze

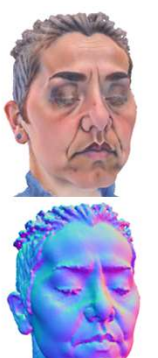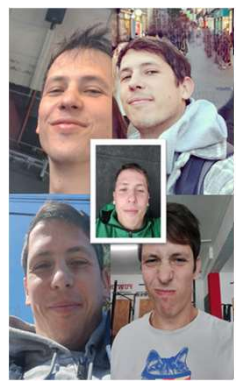(i) Eyes closed    (j) Brow raised    (k) Brow lowerer    (l) Cheeks puffed

(m) Grin    (n) Dimpler    (o) Lip roll    (p) Squeeze

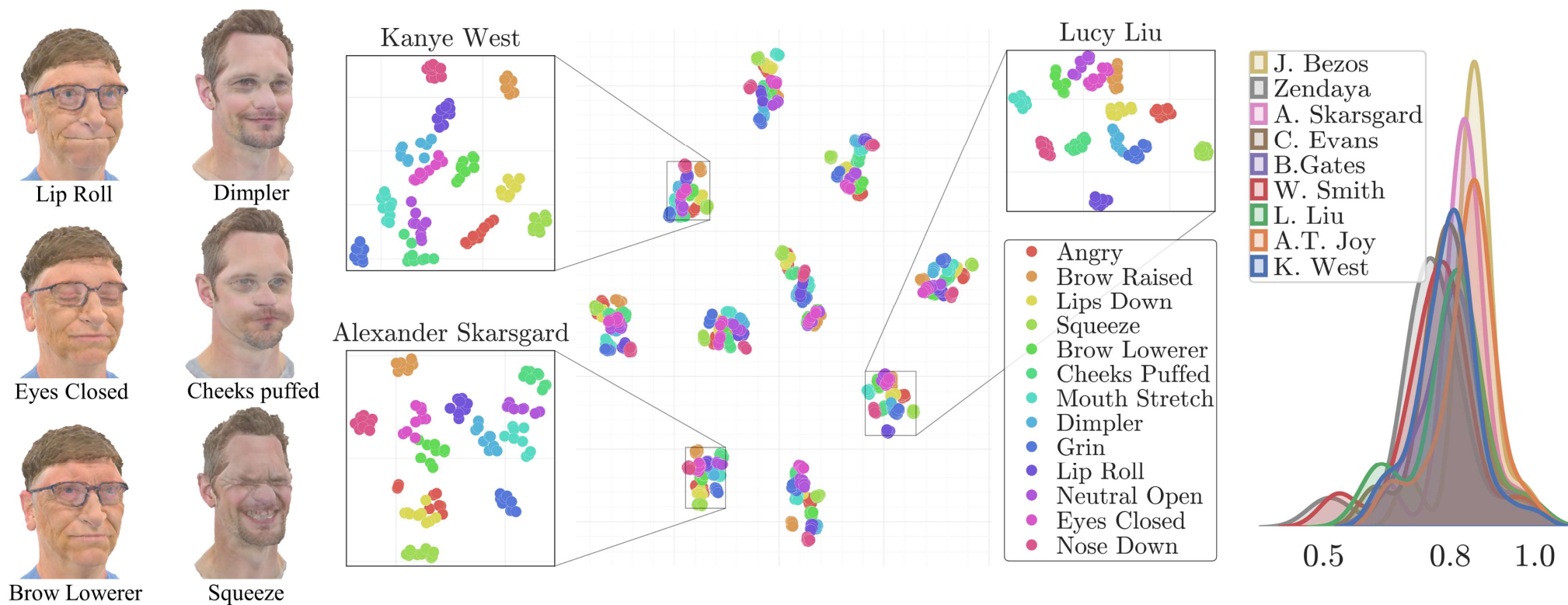"with short buzzcut hairstyle"

# Results: Expressive ID-conditioned Generation

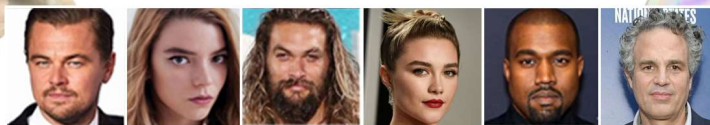# Results: Relighting



Front Light · Right Light · Back Light · Left Light

Direct · Diffused

THANKS!