

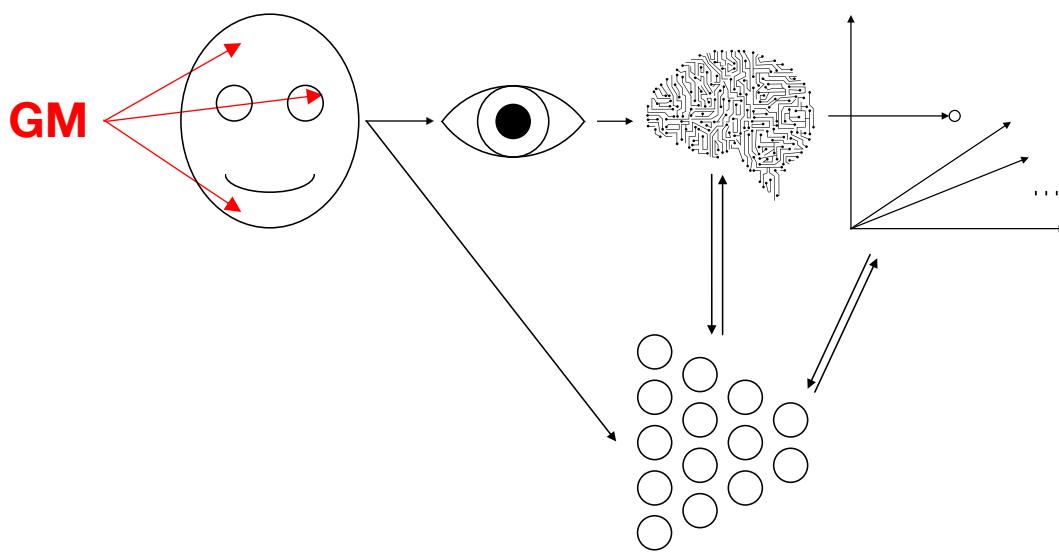
Understanding how DCNNs predict human behaviour using an interpretable generative model of 3D faces

Christoph Daube, Jiayu Zhan, Tian Xu, Andrew Webb, Robin A. A. Ince, Oliver B. Garrod, Philippe G. Schyns



@realdaubman

DCNNs are popular models of human vision – but do they really help us *understanding* it?



 **Federico claudi** @Federico_claudi · 26. März 
It's not necessarily the case that a better performing model affords a better understanding of how the brain does it. The model might be very good at the task, but using a solution radically different from the brain's.

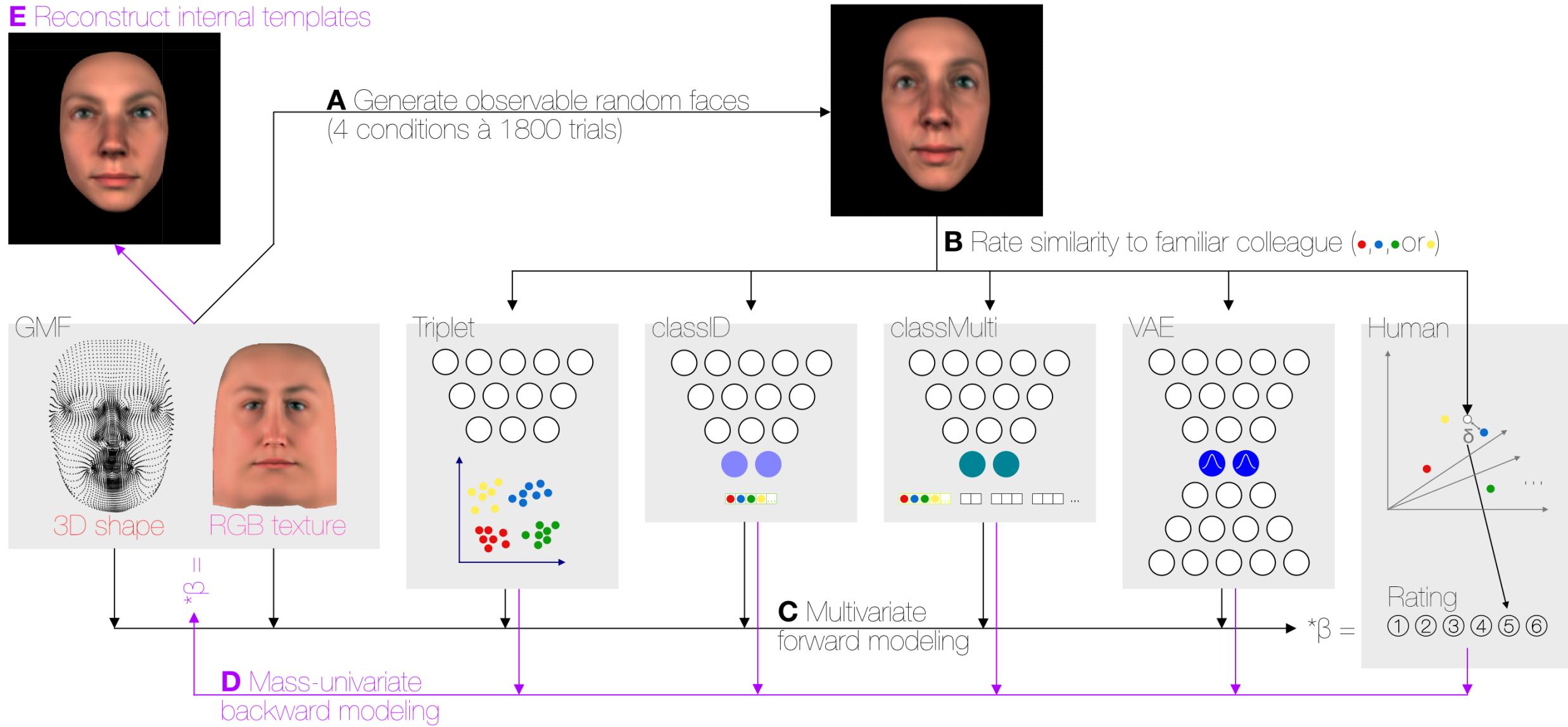
 1   4 

 **KordingLab**  @Kord... · 26. März 
sure. But, arguably, a model that well approximates reality needs to do both: (1) actually work (2) actually describe the neural data. So the idea that you should demand both of models seems pretty attractive.

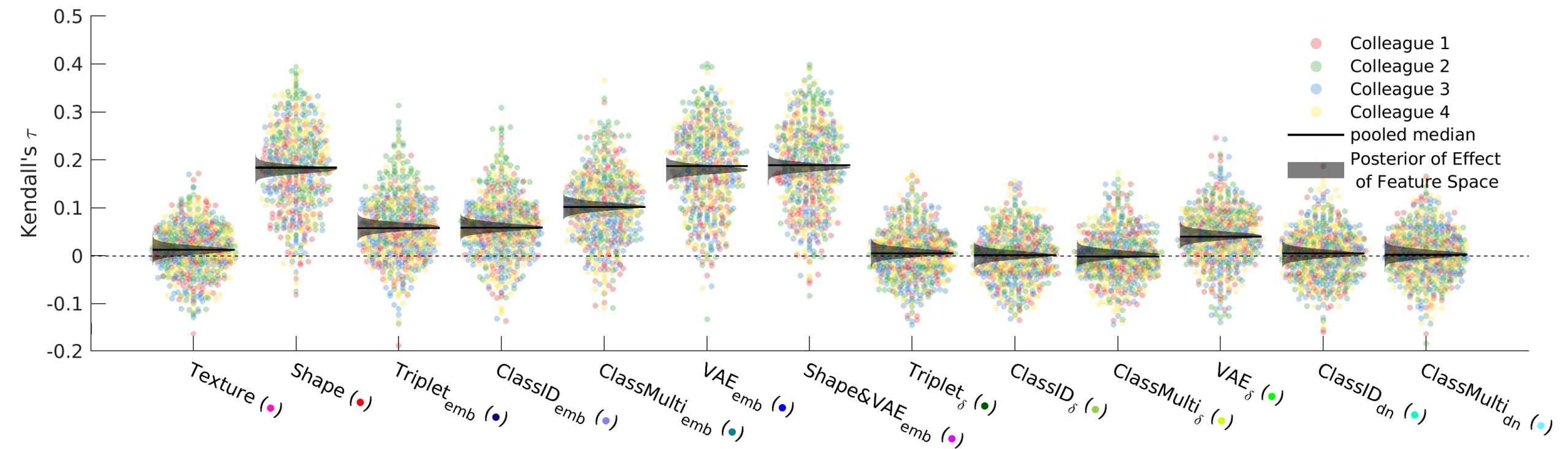
 1   1 

(3): If we want to understand what humans and models are doing and what they have in common, we have to specify what we want to understand – with a **generative model** of the stimuli that allows experimental control!

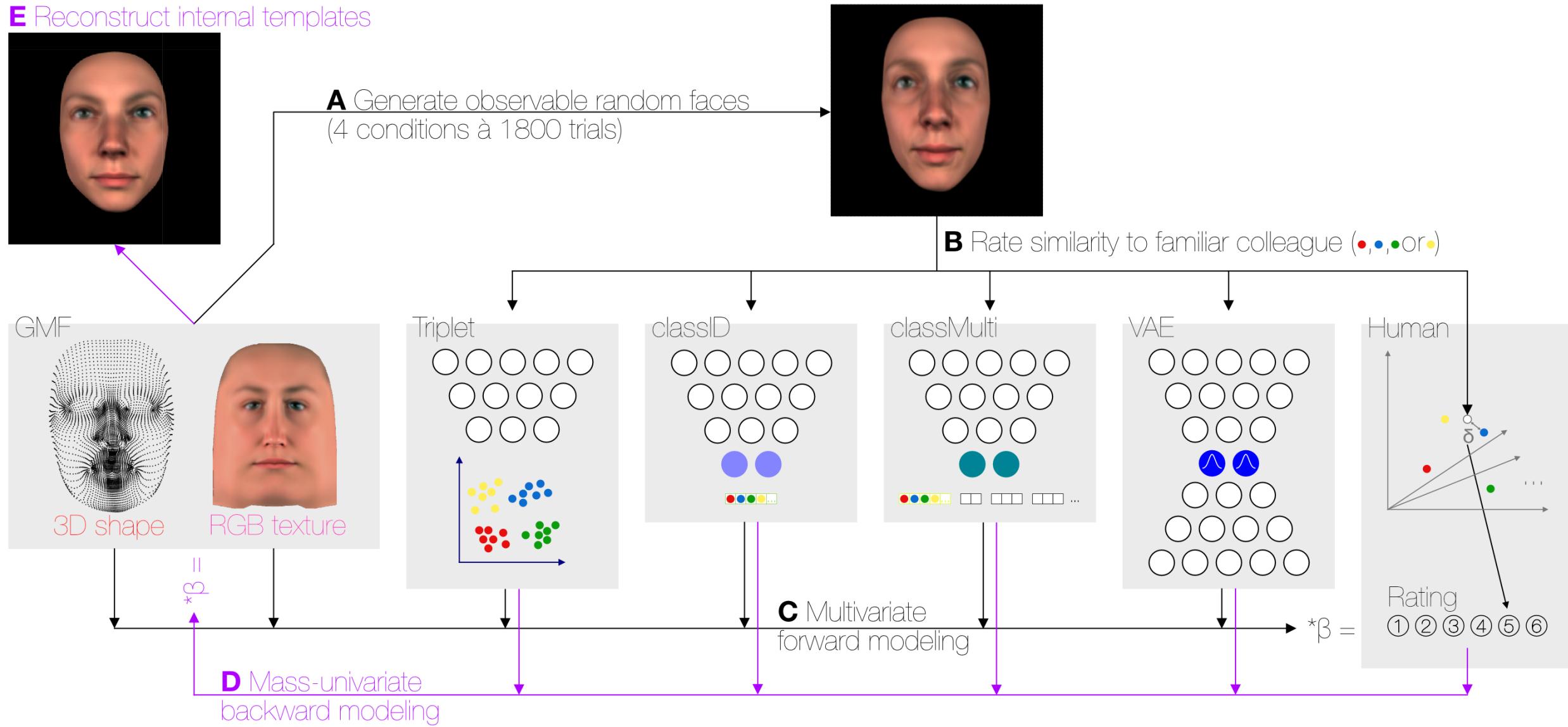
E Reconstruct internal templates



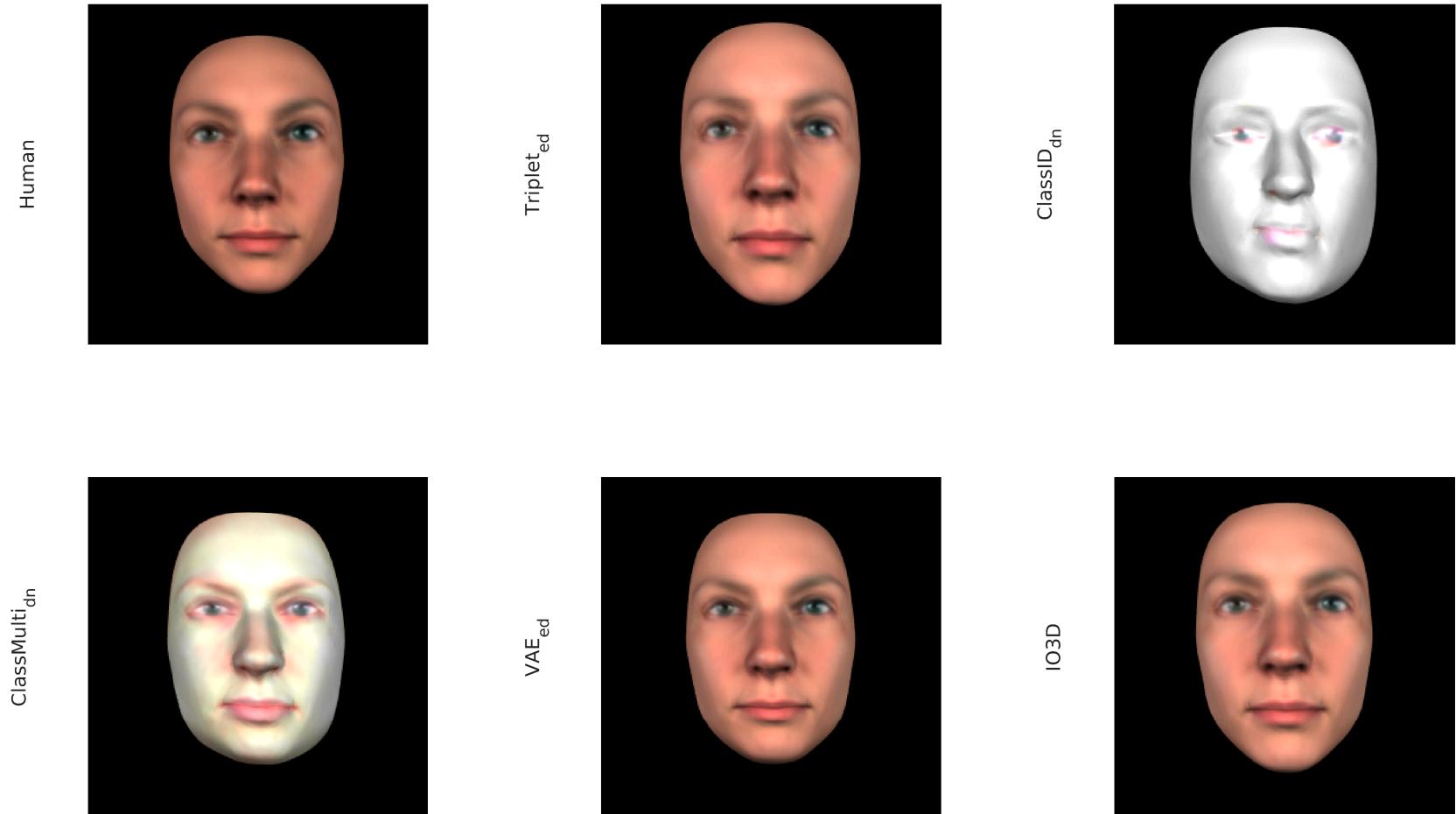
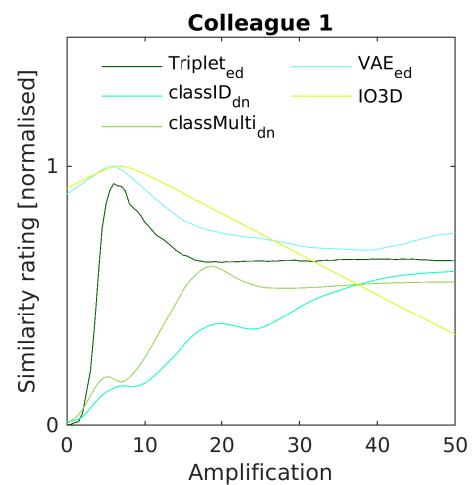
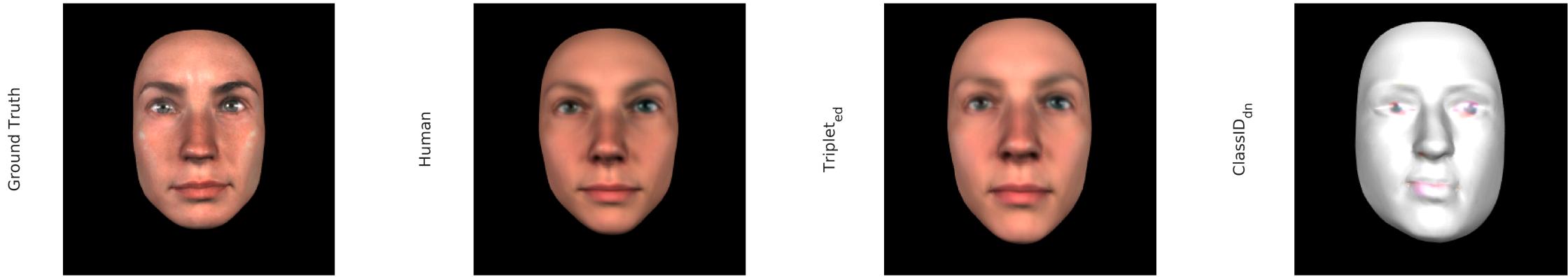
Forward modelling: Comparing predictive power of feature spaces



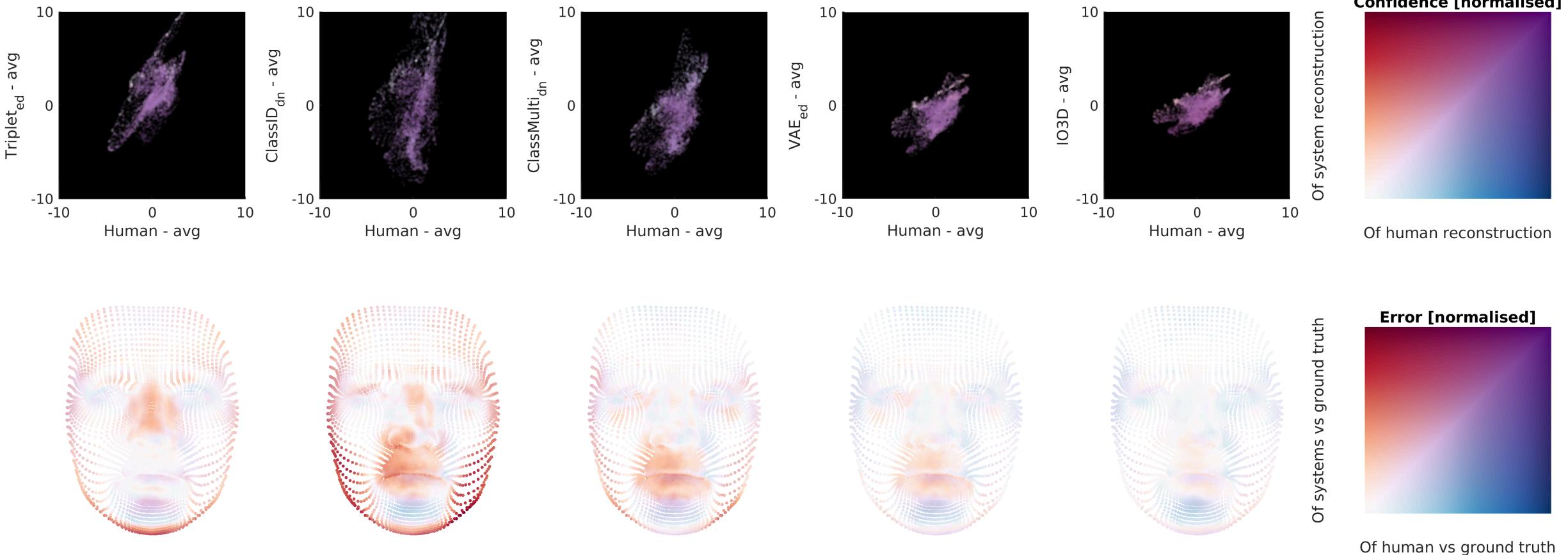
E Reconstruct internal templates



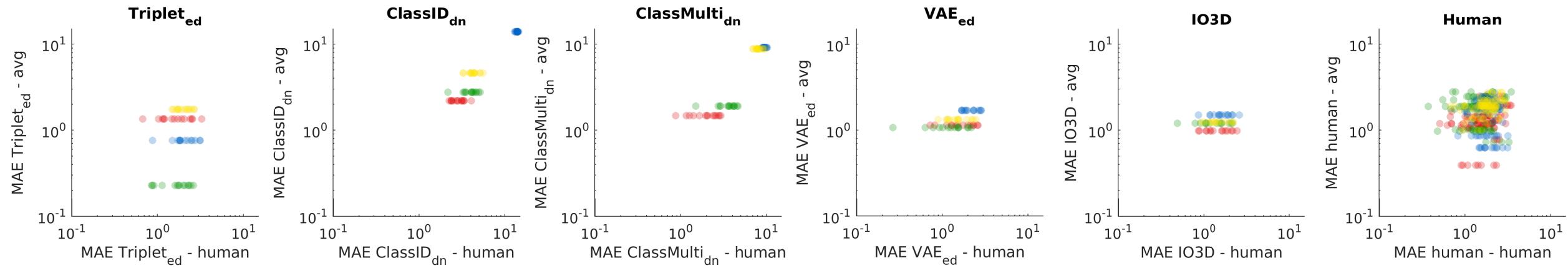
Reconstructions from human and artificial systems



Quantifying human and DNN shape perception w.r.t. generative model



Quantifying human and DNN shape perception w.r.t. generative model



Conclusions

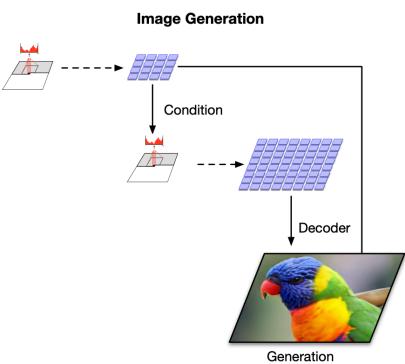
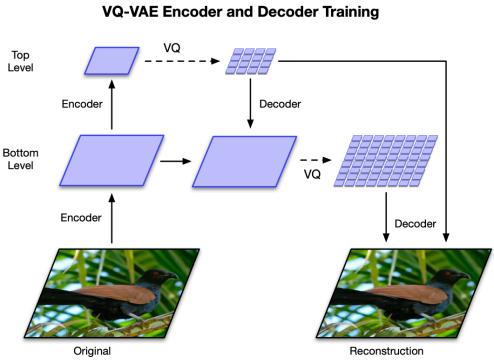
- Human rating behaviour can be better predicted from **VAE** features than from classifier DNN features
- Importantly, we also know **how** the VAE achieves this: By representing 3D shape features more like humans do

**This is a group effort. Thanks to
SchynsLab @ Uni Glasgow, especially**

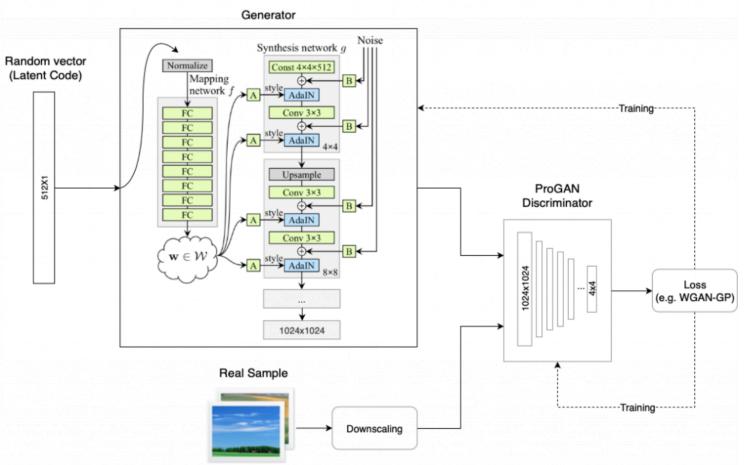
- Jiayu Zhan
- Tian Xu
- Andrew Webb
- Robin A. A. Ince
- Oliver B. Garrod
- Philippe G. Schyns

Discussion: Do this but sample from VAE vs GAN?

VQ-VAE
Razavi
et al 2019

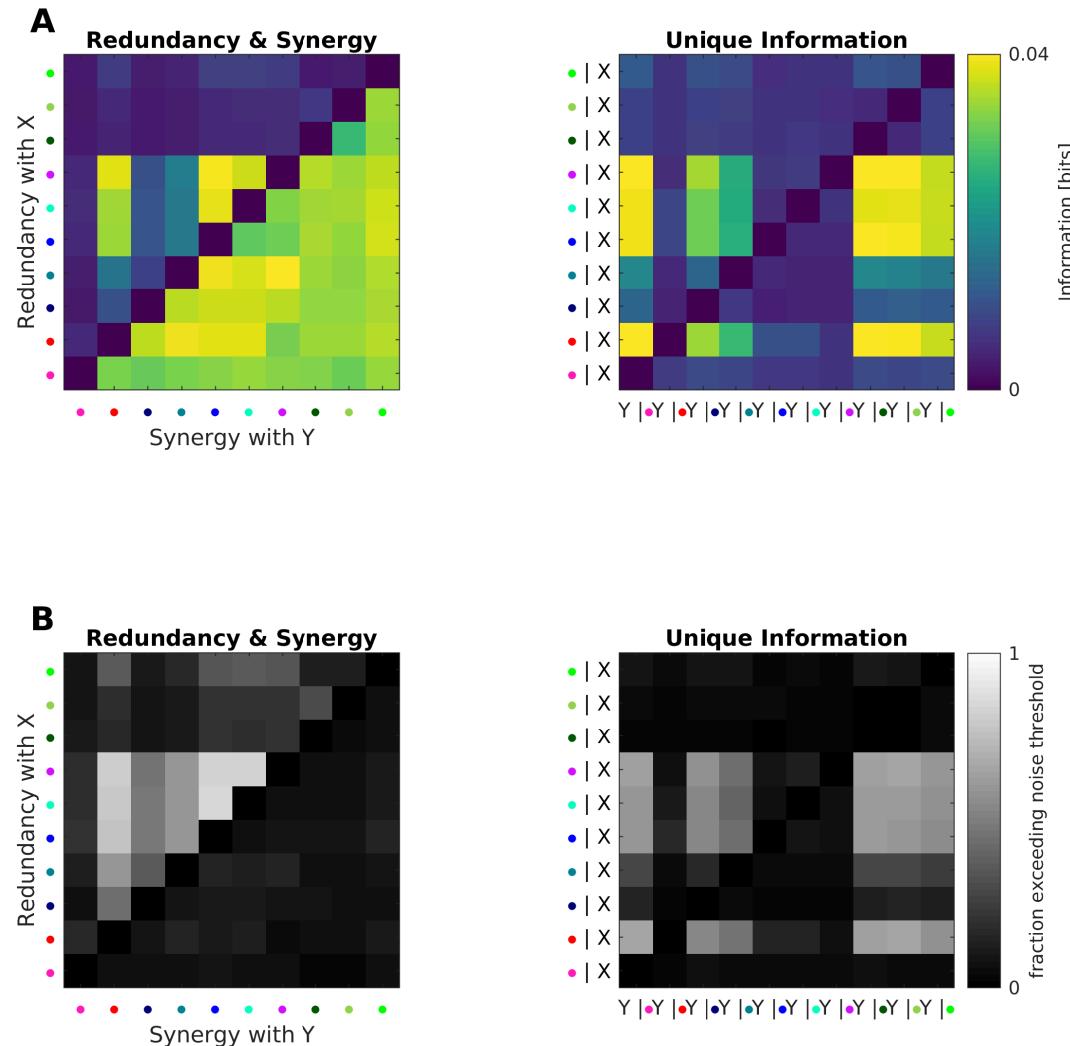
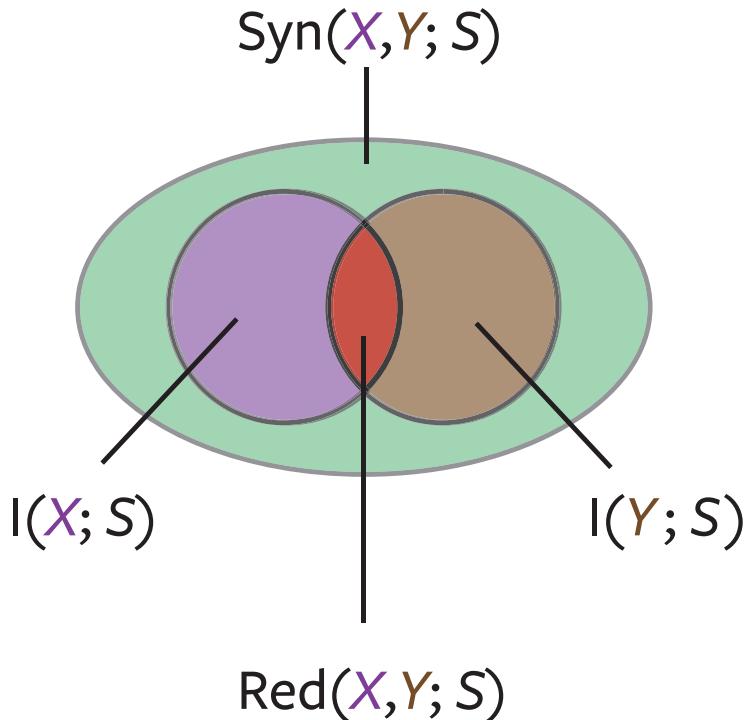


Style GAN (2)
Karras et al
2018,
(Karras et al
2020)



- would any of these models be better for reverse correlation? Why?
- with regards to the initial problem:
“What can VAEs (or GANs) do for Vision Sciences”?
What *new* insights could it provide beyond engineering?

Forward modelling: Comparing predictive power of feature spaces



Ince, 2017
Daube et al, 2019a
Daube et al. 2019b

Forward modelling: Weights are task-dependent

