

Attention Mesh: High-fidelity Face Mesh Prediction in Real-time

Ivan Grishchenko Artsiom Ablavatski Yury Kartynnik Karthik Raveendran Matthias Grundmann
Google Research
1600 Amphitheatre Pkwy, Mountain View, CA 94043, USA
{igrishchenko, artsiom, kartynnik, krav, grundman}@google.com

Abstract

We present *Attention Mesh*, a lightweight architecture for 3D face mesh prediction that uses attention to semantically meaningful regions. Our neural network is designed for real-time on-device inference and runs at over 50 FPS on a Pixel 2 phone. Our solution enables applications like AR makeup, eye tracking and AR puppeteering that rely on highly accurate landmarks for eye and lips regions. Our main contribution is a unified network architecture that achieves the same accuracy on facial landmarks as a multi-stage cascaded approach, while being 30 percent faster.

1. Introduction

In this work, we address the problem of registering a detailed 3D mesh template to a human face on an image. This registered mesh can be used for the virtual try-on of lipstick or puppeteering of virtual avatars where the accuracy of lip and eye contours is critical to realism.

In contrast to methods that use a parametric model of the human face [1], we directly predict the positions of face mesh vertices in 3D. We base our architecture on earlier efforts in this field [5] that use a two stage architecture involving a face detector followed by a landmark regression network. However, using a single regression network for the entire face leads to degraded quality in regions that are perceptually more significant (e.g. lips, eyes).

One possible way to alleviate this issue is a cascaded approach: use the initial mesh prediction to produce tight crops around these regions and pass them to specialized networks to produce higher quality landmarks. While this directly addresses the problem of accuracy, it introduces performance issues, e.g. relatively large separate models that use the original image as input, and additional synchronization steps between the GPU and CPU that are very costly on mobile phones. In this paper, we show that it is possible for a single model to achieve the same quality as the cascaded approach by employing region-specific heads that transform the feature maps with spatial transformers [4], while being

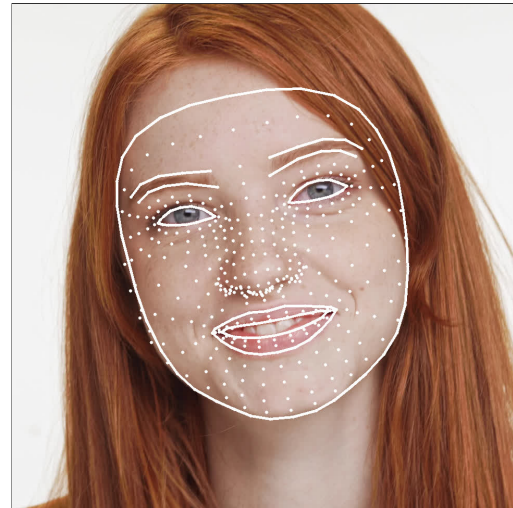


Figure 1. Salient contours predicted by Attention Mesh submodels

up to 30 percent faster during inference. We term this architecture as *attention mesh*. An added benefit is that it is easier to train and distribute since it is internally consistent compared to multiple disparate networks that are chained together.

We use an architecture similar to one described in [7], where the authors build a network that is robust to the initialization provided by different face detectors. Despite the differing goals of the two papers, it is interesting to note that both suggest that a combination of using spatial transformers with heads corresponding to salient face regions produces marked improvements over a single large network. We provide the details of our implementation for producing landmarks corresponding to eyes, irises, and lips, as well as quality and inference performance benchmarks.

2. Attention mesh

Model architecture The model accepts a 256×256 image as input. This image is provided by either the face detector or via tracking from a previous frame. After extract-

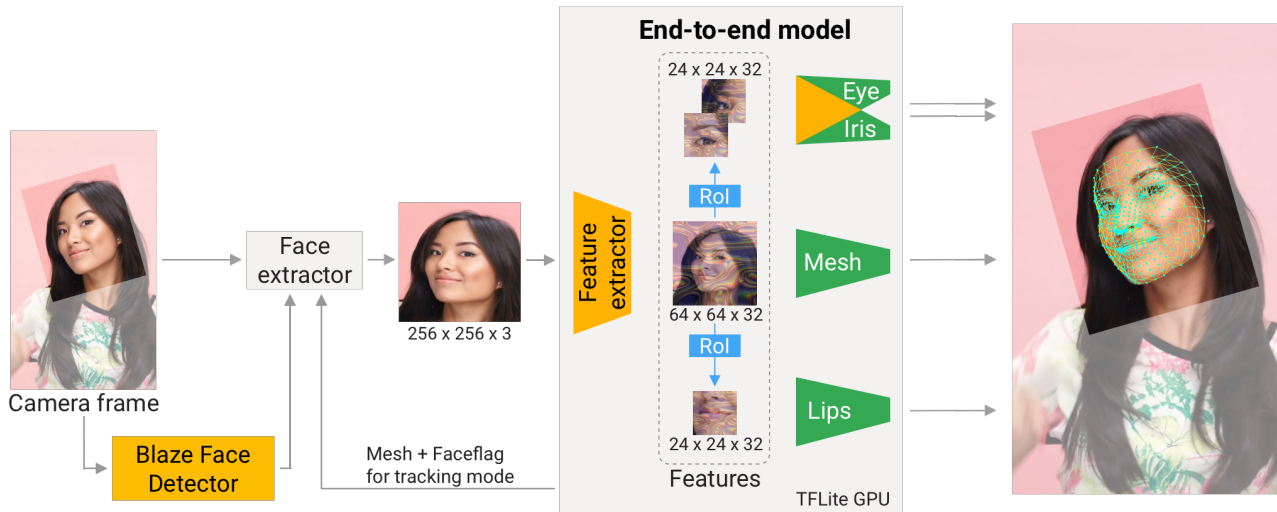


Figure 2. The inference pipeline and the model architecture overview

ing a 64×64 feature map, the model splits into several sub-models (Figure 2). One submodel predicts all 478 face mesh landmarks in 3D and defines crop bounds for each region of interest. The remaining submodels predict region landmarks from the corresponding 24×24 feature maps that are obtained via the attention mechanism.

We concentrate on three facial regions with key contours: the lips and two eyes (Figure 1). Each eye submodel predicts the iris as a separate output after reaching the spatial resolution of 6×6 . This allows the reuse of eye features while keeping dynamic iris independent from the more static eye landmarks.

Individual submodels allow us to control the network capacity dedicated to each region and boost quality where necessary. To further improve the accuracy of the predictions, we apply a set of normalizations to ensure that the eyes and lips are aligned with the horizontal and are of uniform size.

We train the attention mesh network in two phases. First, we employ ideal crops from the ground truth with slight augmentations and train all submodels independently. Then, we obtain crop locations from the model itself and train again to adapt the region submodels to them.

Attention mechanism Several attention mechanisms (soft and hard) have been developed for visual feature extraction [2, 4]. These attention mechanisms sample a grid of 2D points in feature space and extract the features under the sampled points in a differentiable manner (*e.g.* using 2D Gaussian kernels or affine transformations and differentiable interpolations). This allows to train architectures end-to-end and enrich the features that are used by the attention mechanism. Specifically, we use a spatial transformer mod-

ule [4] to extract 24×24 region features from the 64×64 feature map. The spatial transformer is controlled by an affine transformation matrix θ (Equation 1) and allows us to zoom, rotate, translate, and skew the sampled grid of points.

$$\theta = \begin{bmatrix} x_x & sh_x & t_x \\ sh_y & s_y & t_y \end{bmatrix} \quad (1)$$

This affine transformation can be constructed either via supervised prediction of matrix parameters, or by computing them from the output of the face mesh submodel.

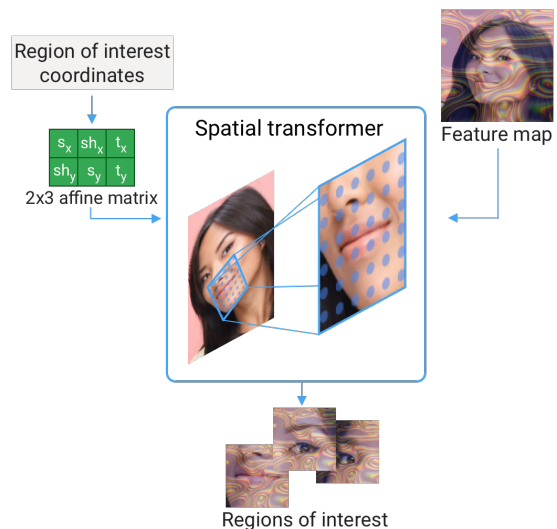


Figure 3. Spatial transformer as the attention mechanism

Dataset Our dataset contains 30K in-the-wild mobile camera photos taken with numerous camera sensors and in varied conditions. We used manual annotation with special emphasis on consistency for salient contours to obtain the ground truth mesh vertex coordinates in 2D. The Z coordinate was approximated using a synthetic model.

3. Results

To evaluate our unified approach, we compare it against the cascaded model which consists of independently trained region-specific models for the base mesh, eyes and lips that are run in succession.

Performance Table 1 demonstrates that the attention mesh runs 25%+ faster than the cascade of separate face and region models on a typical modern mobile device. The performance has been measured using the TFLite GPU inference engine [6]. An additional 5% speed-up is achieved due to the reduction of costly CPU-GPU synchronizations, since the whole attention mesh inference is performed in one pass on the GPU.

Model	Inference Time (ms)
Mesh	8.82
Lips	4.18
Eye & iris	4.70
Cascade (sum of above)	22.4
Attention Mesh	16.6

Table 1. Performance on Pixel 2XL (ms)

Mesh quality A quantitative comparison of both models is presented in Table 2. As the representative metric, we employ the mean distance between the predicted and ground truth locations of a specific subset of the points, normalized by 3D interocular distance (or the distance between the corners in the case of lips and eyes) for scale invariance. The attention mesh model outperforms the cascade of models on the eye regions and demonstrates comparable performance on the lips region.

Model	All	Lips	Eyes
Mesh	2.99	3.28	6.6
Cascade	2.99	2.70	6.28
Attention mesh	3.11	2.89	6.04

Table 2. Mean normalized error in 2D.

4. Applications

The performance of our model enables several real-time AR applications like virtual try-on of makeup and puppeteering.

AR Makeup Accurate registration of the face mesh is critical to applications like AR makeup where even small errors in alignment can drive the rendered effect into the “uncanny valley” [8]. We built a lipstick rendering solution (Figure 4) on top of our attention mesh model by using the contours provided by the lip submodel. A/B testing on 10 images and 80 people showed that 46% of AR samples were actually classified as real and 38% of real samples — as AR.



Figure 4. Virtual makeup comparison: base mesh without refinements (left) vs. attention mesh with submodels (right)

Puppeteering Our model can also be used for virtual puppeteering and facial triggers. We built a small fully connected model that predicts 10 blend shape coefficients for the mouth and 8 blend shape coefficients for each eye. We feed the output of the attention mesh submodels to this blend shape network. In order to handle differences between various human faces, we apply Laplacian mesh editing to morph a canonical mesh into the predicted mesh [3]. This lets us use the blend shape coefficients for different human faces without additional fine-tuning. We demonstrate some results in Figure 5.



Figure 5. Puppeteering

5. Conclusion

We present a unified model that enables accurate face mesh prediction in real-time. By using a differentiable attention mechanism, we are able to devote computational resources to salient face regions without incurring the performance penalty of running independent region-specific models. Our model and demos will soon be avail-

able in MediaPipe (<https://github.com/google/mediapipe>).

References

- [1] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *Proceedings of 36th International Conference and Exhibition on Computer Graphics and Interactive Techniques*, pages 187–194, 1999. 1
- [2] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. *arXiv preprint arXiv:1502.04623*, 2015. 2
- [3] Jianwei Hu, Ligang Liu, and Guozhao Wang. Dual laplacian morphing for triangular meshes. *Computer Animation and Virtual Worlds*, 18(45):271–277, 2007. 3
- [4] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. 1, 2
- [5] Yury Kartynnik, Artsiom Ablavatski, Ivan Grishchenko, and Matthias Grundmann. Real-time Facial Surface Geometry from Monocular Video on Mobile GPUs. *arXiv preprint arXiv:1502.04623*, July 2019. 1
- [6] Juhyun Lee, Nikolay Chirkov, Ekaterina Ignasheva, Yury Piskarchyk, Mogan Shieh, Fabio Riccardi, Raman Sarokin, Andrei Kulik, and Matthias Grundmann. On-device neural net inference with mobile gpus. *arXiv preprint arXiv:1907.01989*, 2019. 3
- [7] J. Lv, X. Shao, J. Xing, C. Cheng, and X. Zhou. A deep regression architecture with two-stage re-initialization for high performance facial landmark detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3691–3700, 2017. 1
- [8] Junichiro Seyama and Ruth S. Nagayama. The uncanny valley: Effect of realism on the impression of artificial human faces. *Presence: Teleoper. Virtual Environ.*, 16(4):337351, Aug. 2007. 3