# Lightweight Markerless Monocular Face Capture with 3D Spatial Priors

Shridhar Ravikumar

**Figure 1:** *Results from our solver re-targeted onto multiple face models. As seen in the side views, our method is able to handle ambiguities in depth, inherent in monocular input, and achieve results that are physically plausible even in areas with occlusions like around the lips.*

### Abstract

*We present a simple lightweight markerless facial performance capture framework using just a monocular video input that combines Active Appearance Models for feature tracking and prior constraints on 3D shapes into an integrated objective function. 2D monocular inputs inherently lack information along the depth axis and can lead to physically implausible solutions. In order to address this loss of information, we enforce a constraint on our objective function within a probabilistic framework that uses preexisting animations obtained from accurate 3D tracking systems, thus achieving more plausible results. Our system fits a Blendshape model to tracked 2D features while also handling noise in estimation of features and camera parameters. We learn separate constraints for the upper and lower regions of the face thus maintaining flexibility. We show that using this approach, we can obtain significant improvement in tracking especially along the depth dimension. Our method uses easily obtainable prior animation data. We show that our method can generate convincing animations using only a monocular video input. We quantitatively evaluate our results comparing it with an approach using a monocular input without our spatial constraints and show that our results are closer to the ground-truth geometry. Finally, we also evaluate the effect that the choice of the Blendshape set has on the results of the solver by solving for a different set of Blendshapes and quantitatively comparing it with our previous results and to the ground truth. We show that while the choice of Blendshapes does make a difference, the use of our spatial constraints generates results that are closer to the ground truth.*

## 1. Introduction

Over the last few years many methods for facial performance capture have been developed. These methods range in complex-

ity from active marker based approaches with multiple cameras, head-mounted devices and controlled environments [Wil90, Vic, FHW*11, FP09, FJA*14, BHPS10, BHB*11, TlTM11, Rai04, SLS*12, GGW*98, BGY*13, VWB*12, MJC*08] to passive markerless monocular approaches [SKSS14, CET01, SLC11, CWLZ13, CHZ14, CBZB15, GVWT13]. Then there are methods that use depth sensing devices like the Microsoft Kinect and others [ZSCS08, WLGP09, WBLP11, LYYB13, BWP13]. Although many of these methods provide good quality animations, the time taken to setup the environment including the time taken to apply markers on the actor's face can be a hindrance to the quick application of these methods, similarly the capture devices used may be custom built and not easily available or cheap enough for general consumer use. That being said, there is ample data that exists from these legacy systems, freely available online, captured from multiple people performing a range of facial movements and speech. One of the problems with 2D monocular inputs is that information along the depth axis is inherently missing and this proves to be a challenge. Even though the solver may minimize the error in the objective function, the true 3D shape may be ambiguous. In this work, we present a lightweight approach that achieves good quality solves using just a monocular input while exploiting this existing data from legacy systems and using it to learn a prior in order to make sure physically implausible results aren't generated. Our contributions are as below:

- We present a lightweight monocular markerless capture method that achieves good quality animation parameters and does not require special equipment or complex training phases.
- We exploit easily available prior animation data obtained from 3D tracking systems and use this in a density estimation framework to regularize our objective function and generate more plausible results. We enable further flexibility by learning separate prior constraints for the upper and lower face regions.
- We combine initial estimates of 2D landmark points on the face based on an ensemble of regression trees, with an Active Appearance Model for improved accuracy.
- We handle noise in input 2D features and in the estimation of camera extrinsic parameters thus eliminating jitter in the resulting animation.
- Finally we evaluate the effect that the choice of Blendshapes has on the results of the solver in general and show that while it affects the results, incorporating the spatial priors generates results that are closer to the ground truth.

## 2. Related Work

Many different approaches for facial performance capture have been demonstrated in the last decade or so and can broadly be categorized into different types mainly based on the representation that is used for the underlying animation model and also based on the input devices used for the capture process.

Mesh propagation approaches deform a single high quality mesh of the actor's face through the sequence in order to generate the animation based on motion capture data. Bradley et al. [BHPS10] propagate a mesh by applying optical flow over multiple cameras and combine both 2D video and 3D point clouds in order to obtain very high resolution captures. Fyffe et al. [FHW*11, FJA*14]

use five high speed cameras combined with gradient illumination patterns in order to get extremely high resolution captures. They also make use of multiple static scans in order to account for high resolution detail. Beeler et al. [BHB*11] use anchor frames in order to reduce drift in optical flow. [SKSS14] automatically reconstruct the 3D face shape of the subject from an unconstrained data-set of images obtained over the internet and track the face through a monocular video sequence by finding the optimal scene flow parameters. They then add in fine scale details using a shape from shading approach. Other mesh propagation approaches include [GGW*98, BPL*05, VWB*12].

On the other hand, parametric model based approaches optimize for the values of the parameters through the sequence in order to best capture the movement of the face. Statistical models, such as those based on Principle Component Analysis [CET01, BV99, TlTM11] provide an orthogonal set of basis which can then be combined linearly in order to generate the animation. Blendshape models consist of a set of facial expressions as a basis and these expressions are combined in a linear fashion in order to generate different expressions. This is arguably the most popular approach for representing animation owing to its intuitiveness and relative ease of use. [LAR*14] gives an excellent in-depth analysis of Blendshapes. [WLGP09, BWP13, GVWT13, LWP10, LYYB13, CWLZ13, CHZ14, DCFN06] all make use of Blendshape models combined with either multi-view, monocular or depth sensing input devices. [TZN*15] presented a method for real-time facial reenactment using an RGB-D device. They use a statistical model for the identity, expression and albedo of the face and optimize for the best model parameters in an analysis by synthesis approach. They model the scene lighting using spherical harmonic basis by assuming that the light source is distant and the scene is predominantly lambertian. They optimize for the parameters in real-time using a data parallel GPU solver.

In this work, our objective is to capture the movement of the actor's face using a lightweight method that doesn't require expensive equipment or lighting apparatus and using only a single monocular off-the-shelf camera while simultaneously addressing the issue of missing depth information. In light of this, we will mainly discuss previous work that is similar to our approach in these regards.

Our approach is most similar in spirit to Garrido et al. [GVWT13]. They present a lightweight approach that makes use of a single monocular video input and are able to generate a high quality output animation with fine scale details using a Blendshape model. They track a few sparse landmark feature points on the face reliably using forward and backward optical flow combined with automatic key-frame selection based on local binary patterns for robustness. The pose and facial expressions are estimated from these sparsely tracked points on the face in an iterative fashion. Temporally coherent dense motion fields tracked from video combined with a smoothness constraint is then used in order to refine the pose and facial expression. Fine scale details are then added on top using a shape from shading approach. Although their method produces good results, as it is a monocular system, it doesn't account for inherent loss of information along the depth dimension. Our method uses a prior constraint in order to regularize the data which alleviates the error due to the lack of depth information. Also our method

doesn't depend on optical flow across the sequence, but only on the current and previous frames and thus can be implemented in an online fashion.

Cao et al. [CWLZ13] presented a method for obtaining real-time performance capture from monocular input. Similar to our approach, they track 2D points but instead of fitting directly to the 2D features, they train a user specific two-level boosted regressor trained on labelled 2D points and corresponding 3D shapes, in order to map from 2D to 3D features at run-time. They then fit a Blendshape model to the obtained 3D features by iteratively solving for transformation parameters and expression weights. Their method requires a training phase where images of the user in different poses and expressions are captured. [CHZ14] extend the previous work to be independent of a user and instead learn a regressor from public image databases. They infer both the 2D facial landmarks and the 3D shape of the face simultaneously. Their algorithm adapts to the user's face at run-time by solving for the user-specific Blendshapes and the expression co-efficients in an iterative manner. [CBZB15] enhance a low-resolution tracked mesh with medium-scale wrinkle details which are generated by local regressors trained on high-resolution scan data. They require a one-time training phase in order to learn the mapping from UV space to vertex offsets and can be applied to an unseen actor at runtime. While their method obtains very impressive and detailed results, the training process is quite involved and requires generation of multiple guess-truth pairs for the parameters for each training image in order to train the DDE regressor as explained in [CHZ14].

## 3. System Overview

Our pipeline fits a Blendshape model to automatically tracked 2D feature points in the video. In order to generate our Blendshapes (Sec.5), we first deform the neutral expression mesh from an existing template Blendshape model to a 3D scanned face of the actor to obtain the mesh of the actor in a desired topology. We then automatically generate the user-specific Blendshape expressions for this actor. Our 2D landmark features are initially obtained using the method of [KS14] giving us 68 distinct landmarks (Sec.4). These landmarks are detected per frame and although they provide a good starting point, the detections are not accurate enough for our purposes and give unacceptable results especially for extreme facial expressions. In order to address this, we use the Fast Simultaneous Inverse Compositional algorithm of [TP13], trained on a few images of the user which is applied on top, using the results from [KS14] as an initialization. This gives us more robust landmark detections especially in extreme facial expressions. We solve for the camera extrinsic parameters using the Perspective-n-Point approach with the Levenberg-Marquardt algorithm for non-linear optimization (Sec.6). In order to address the inherent noise from estimating landmarks and camera parameters per frame, we smooth the noise in the 2D features and the camera parameters using a Kalman filter (Sec.7). In order to regularize our solve results, we make use of a prior energy constraint (Sec.9 and Sec.10) that estimates the probability of the solution and factors it into the objective function (Sec.11) resulting in more plausible shapes. Our optimization function then solves for the optimal coefficients of the Blendshapes that minimize the re-projection error of the landmark

**Figure 2:** *Landmark points initialized using an ensemble of regression trees and updated by the AAM*

vertices whilst taking into consideration its likelihood. In Section 12.1 we quantitatively compare our results to the ground-truth and show that our method generates more accurate results. Finally in Section 13, we evaluate a completely different set of Blendshapes in order to analyze the effect that the choice of Blendshapes has on the solver result, and also quantitatively show that our method still generates more accurate results.

## 4. 2D Facial Feature Detection

To obtain our initial 2D features, we first use the algorithm of [KS14] implemented within the Dlib library [Kin09], to detect 68 landmarks on the face on a frame by frame basis. The algorithm uses an ensemble of regression trees to accurately predict the optimal displacement of each landmark at each level, based on differences in pixel intensity around that landmark. We trained the ensemble of regression trees using the images from the HE-LEN, LFPW, AFW, IBUG and 300 faces In-the-wild databases [SAT*16, STZP13b, STZP13a], to give a total of 4213 training images with a wide variety of pose, lighting and shape variations. We used a cascade depth of 10, tree depth of 10, 500 trees per cascade, a feature pool of 400 and 50 test splits (see [KS14] for details). This gives us reasonable initial detections of landmarks on the user's face but isn't accurate or robust enough for the application of performance capture as shown in the accompanying video. The landmark detections are incorrect during extreme expressions and this leads to unacceptable solve results.

In order to improve upon this, we further train an Active Appearance Model [CET01] using the Fast Simultaneous Inverse Compositional approach of [TP13], on a few select images of the user performing a few facial expressions. We used 15 facial expressions that included a few that elicit the extreme range of the user's facial movements (jaw-open, lip-swing) and also a few challenging expressions that generate lip and eye occlusions (pucker, squinch). We then use the result of the previous step as a starting point and then use the AAM to improve landmark tracking through the sequence. In our experiments, this gives us much more robust detections and also covers the full range of the user's expressions as shown in the accompanying video. The results of this landmark detection are further fed into the Kalman filter to account for discrepancies between frames and to remove noise in an online fashion.

## 5. Automatic Blendshape Generation

In order to obtain the neutral expression mesh of the actor, we scan the actor in a one-time pre-processing step. Alternatively, this can also be obtained automatically using the approach of [BV99]. The resulting scanned mesh has a random noisy topology which cannot be used directly. In order to rectify this, we then deform the neutral expression mesh from an existing template Blendshape model to this scanned mesh using the algorithm of [ARV07] to obtain a new mesh with the desired topology. We then automatically generate the expression Blendshapes for the actor's face by applying the Deformation Transfer approach of [SP04] to give us 140 unique Blendshapes. This gives us a linear Blendshape model with $N$ Blendshapes for use in our solver as below:

$$B = B_0 + \sum_{i=1}^{N} \alpha_i(B_i - B_0) \qquad (1)$$

where, $B_0$ is the neutral expression Blendshape, $\alpha_i$ is the weight associated with Blendshape $i$ and $B_i$ corresponds to the i-th Blendshape.

The vertices on the 3D mesh that correspond to the automatically detected landmarks from Sec.4, are chosen by the user in a one time manual step.

## 6. Camera Calibration for Facial Projection

We calibrate our camera using a standard checkerboard pattern in order to obtain our camera intrinsic matrix $K$ in a one time step. We then calculate the extrinsic parameters $[R|t]$ i.e. the rotation and translation that take the 3D face mesh from model coordinates to camera coordinates. This is essentially a Perspective-n-Point problem between the vertex coordinates on the mesh and the corresponding landmarks in 2D. Given these corresponding landmarks, we then solve for the camera extrinsic parameters every frame using the Levenberg-Marquadt non-linear optimization framework combined with RANSAC for robustness. This gives us the $[R|t]$ values every frame that we combine with the camera intrinsic matrix $K$ to give us our projection matrix $-K[R|t]$. We update the values of the
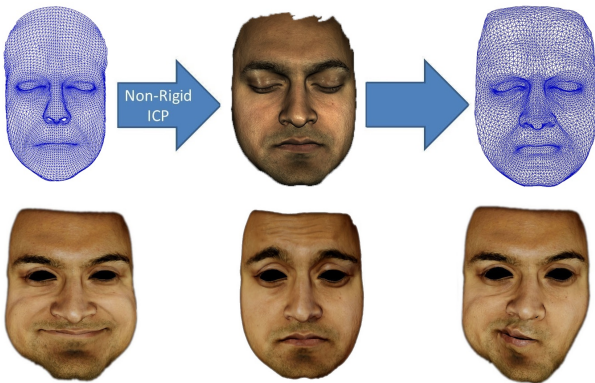


**Figure 3:** *Non-rigid ICP from template mesh to scan resulting in user specific Blendshapes*

3D landmark points on the mesh every frame by using the previous frame's shape. This ensures better projection every frame and improves the accuracy of the solver. In order to account for differences between adjacent frames, we further process these results and smooth them across frames by feeding these into the Kalman filter.

## 7. Online Facial Feature Smoothing

The 2D features and the camera extrinsic parameters so far, are obtained on a frame by frame basis. Both the 2D landmark detections and the camera extrinsic estimates are prone to noise and this independence between the frames inevitably leads to jitter in the final animation as shown in the accompanying video. This necessitates a smoothing operation in order to ensure consistency and avoid sudden changes. We use a Kalman filter in order to smoothly transition these values between frames. The Kalman filter predicts the value of the landmarks and extrinsic parameters every frame and then uses the observations to update its belief about what the parameters should actually be, based on the value of the Kalman-gain-factor, which it calculates based on the noise in observation and the noise in the process. This ensures that any updates made to the landmarks and the extrinsic parameters are updated smoothly. We use one Kalman filter for updating the changes in all the 2D landmarks and one for updating the changes in rotation and translation parameters for the camera extrinsics.

The location of the landmark point detections are predicted every frame by taking into account the velocity and acceleration along each dimension of the tracked points based on the previous frames and then combined with the observations at that frame. A similar process is applied for the translation values in the camera extrinsics. As for the rotations, the orientations in 3D have 3 degrees of freedom ( yaw, pitch and roll ) but our rotation matrix has 9 parameters. Smoothing the rotation matrix directly in this space will not constrain the values properly and will not ensure valid rotation values. Hence we perform the smoothing in quaternion space. This removes ambiguity and also constrains the rotations better thus ensuring that our rotation values are valid throughout the sequence.

This process of online smoothing using Kalman filters greatly helps with tackling noise inherent in the 2D feature detection and camera extrinsics and ensures that our objective function isn't affected by noise. The results of our updated landmarks and camera projection matrix can be seen in Figure 4 and in the accompanying video.

## 8. 2D Objective Function

Given the corrected 2D landmarks and the projection parameters from the previous step, we then solve for the facial expressions by calculating the optimal coefficient weights for the Blendshape model using the following objective function, similar to [CWLZ13].

$$E_{2D(\alpha)} = \sum_{l=1}^{n} \|\Pi_Q(M(B_0 + \sum_{i=1}^{N} \alpha_i B_i)^{(v_l)}) - q^{(l)}\|^2 \qquad (2)$$

where:

- $n$ is the number of landmarks
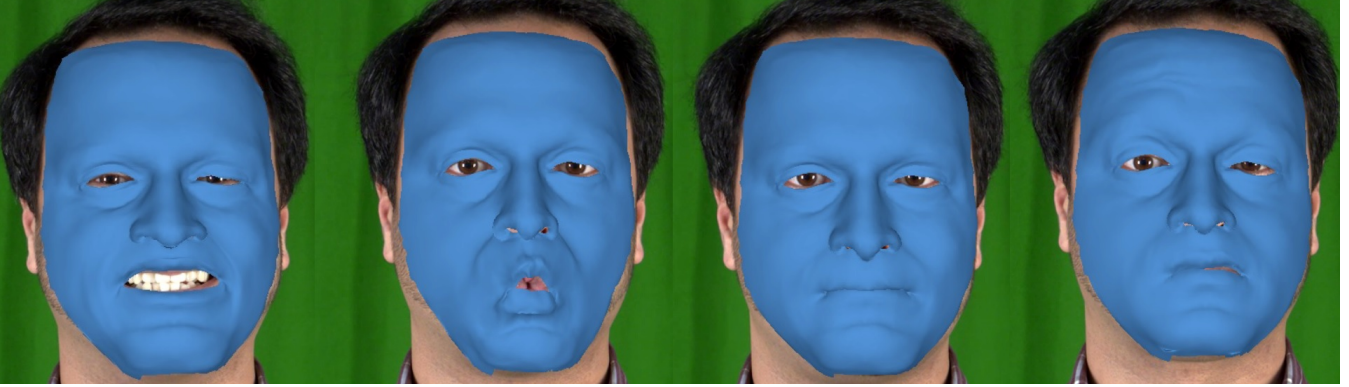- $\Pi_Q$ is the camera projection matrix

**Figure 4:** *Results from the projection of our 3D face mesh onto the video*

- *M* is the rigid transform from object space to camera coordinates
- *N* is the number of Blendshapes
- $B_0$ is the neutral expression Blendshape
- $\alpha_i$ is the weight associated with Blendshape *i* with the constraint $0 \leq \alpha_i \leq 1$
- $B_i$ corresponds to the i-th Blendshape
- $q^{(l)}$ represents the l-th 2D landmark point in the image
- $v_l$ represents the vertex corresponding to landmark *l*

This objective function by itself only considers the 2D landmarks to reduce the error and is inherently under constrained as the information along the depth axis is missing in the video. This leads to insufficient constraining and to errors that show up especially along the depth axis. In Section 9 and 10 we tackle this problem by improving on our objective function and making use of prior constraints.

## 9. 3D Spatial Constraints

One of the problems inherent in monocular capture approaches is that the information along the depth dimension is lost. This leads to situations where the optimization function described above is able to obtain coefficients that minimize the squared distance between the projected vertices and the corresponding 2D landmarks but it does so with solutions that no longer adhere to physically plausible shapes of the face. This effect is especially visible around the mouth region as there is a lot of variation in the depth dimension during speech. This can also be seen when opening the jaw as there is movement along the depth axis. In order to ensure that our objective function provides physically plausible results, we need to ensure that our results are regularized to stay within such a solution space. In general 3D face capture systems can be less flexible compared to monocular markerless systems and placing physical markers on the actor can be very time-consuming, but these systems provide accurate tracking of points. There is ample data available from legacy 3D face capture systems from multiple people performing different facial expressions and speech sequences. It makes sense to use this data in order to regularize our results. We can use the data from these accurate 3D systems and use it to constrain our results while still retaining a monocular markerless approach.

In a one time training step, we use prior data in the form of 3D marker locations over capture sequences from previous captures using the Vicon Cara 3D head-mounted device [Vic]. Our data spanned multiple sequences of speech and facial expressions from 5 different people performing diverse facial movements. In order to use this data for our purposes, we first need to solve for the Blendshape coefficient weights specific to our Blendshape model. We do this using the standard objective function [LWP10] for 3D solves, as shown below.

$$E_{3D(\alpha)} = \arg\min_{\alpha} \|B_0 + B\alpha - T\|_2^2 + \beta\|\alpha\|_1 + \alpha^T \Gamma \alpha \quad (3)$$

where:

- $B_0$ is a $3n \times 1$ vector representing the neutral face, where *n* is the number of markers.
- *B* is a matrix of size $3n \times N$, that contains the deltas for each of the Blendshapes $B_{i...N}$, where *N* is the number of Blendshapes.
- $\alpha$ is a $N \times 1$ vector of weights with the constraint $0 \leq \alpha_i \leq 1$.
- *T* is a $3n \times 1$ vector representing the target markers.
- The term $\|\alpha\|_1$ is an L1-norm on $\alpha$ that penalizes the sum of weights.
- $\beta$ is a weighting factor on the L1 regularizer.
- The $\Gamma$ term is a Tikhonov regularizer that ensures that the function is convex and has a unique global solution. $\Gamma = \varepsilon I$ where $\varepsilon$ is a very small constant and *I* is the identity matrix of size $N \times N$.

Solving this for our multiple training sequences gives us valid coefficient weights over multiple people. We use this data $D_{prior}$, in order to estimate the posterior of the solution in a probabilistic framework. Essentially given a solution vector of coefficient weights i.e. our likelihood, we need to estimate the posterior $P(\alpha|D_{prior})$, given the prior $P(\alpha)$. We can learn this prior $P(\alpha)$ using a Kernel Density Estimation method - the Parzen window with an RBF kernel as explained below:

$$P(\alpha) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(\alpha_i - \alpha)^2}{2\sigma^2}\right) \quad (4)$$

where *n* is the number of prior data points, $\alpha$ is the estimated coefficient vector and $\alpha_i$ are the prior points. The bandwidth or standard deviation $\sigma$ is obtained as shown in the next section.

## 10. Data Binning

One of the issues with using a density estimation technique directly is that it will give higher probability to points that occur more frequently in our prior data $D_{prior}$. This is not desirable as we assume that all of our prior data is accurate as we obtained it from a high accuracy 3D capture system and we want to weight them equally. So in order to overcome this, we essentially need to perform a data-binning operation where we replace multiple data points by a single vector corresponding to the mean of the cluster to which they belong. This makes sense as it allows us to give equal weight to different valid facial configurations without letting its frequency affect the probability. We use the Mean Shift algorithm [FH75] to cluster our prior data.

The mean shift algorithm is a non-parametric clustering algorithm which does not require the knowledge about the number of clusters. It works by updating candidates for centroids to be the mean of the points within a given region. Given a candidate centroid $x_i$ for iteration $t$, the candidate is updated according to the following equation:

$$x_i^{t+1} = x_i^t + m(x_i^t) \tag{5}$$

where $m$ is the mean shift vector that is computed for each centroid that points towards a region of the maximum increase in the density of points

We make use of an automatic method of bandwidth selection [CD12, HKV13, CM13] for use with the mean shift algorithm. This gives us our bandwidth, $\sigma$ in equation 4. We then use the means of the clusters thus obtained ($\alpha_i$) within the Parzen window density estimation framework (Eqn. 4) in order to obtain the prior probability of a solution.

## 11. Objective Function with 3D Spatial Prior

Finally, we factor this probability into our objective function in order to obtain optimal coefficients $\alpha$ that adhere to physically plausible face configurations. Since the probability of the upper face coefficients should be independent of the lower face coefficients — e.g. to avoid the probability of the eyebrows being raised being affected by the jaw being open — we perform the density estimation of these coefficients separately. This is essential since our prior data is not guaranteed to cover all possible combinations of the upper-face and lower-face shapes in tandem. Our updated objective function is as shown below:

$$E_{Final(\alpha)} = E_{2D} + \frac{\lambda}{P(\alpha|D_{prior})} \tag{6}$$

where $\lambda$ is the weight given to the probability term and $P(\alpha|D_{prior})$ is the posterior distribution on the coefficients. This ensures that the solution provided by our objective function lies within the valid space of facial expressions which is controlled by the prior data. As seen in the results in Figures 5,6 and 7, this yields significant improvement over the solve using just the 2D landmarks. This is especially visible in regions around the mouth as this is where majority of the variation in depth occurs.

## 12. Results and Discussion

We render our 3D model within an OpenGL framework combined with OpenCV for calculation of the projection matrix and solving the perspective-n-point problem. In order to perform our optimization, we make use of the Ceres solver, an open source C++ library for modeling and solving constrained non-linear optimization problems. Our method is online as it doesn't require any operations across frames, although it isn't real-time owing to speed bottlenecks. Our system consists of a 4th generation Intel-I7 2.80GHz quadcore processor with 16GB RAM and an NVIDIA GTX765M graphics card.

While the 2D feature detector used for initializing the landmarks per frame is independent of the user (as it's trained on different faces), in our experiments it wasn't accurate enough, especially for extreme facial expressions including ones with occlusions. The Active Appearance Model is trained on 15 images of the user performing few facial expressions that elicit the range of his expressions both on the upper and lower face regions. These training images are marked semi-automatically – initialized using the previous 2D feature tracker and then corrected where needed. Although this makes our system dependent on the specific actor, it is a small price to pay for the improvement in tracking especially considering only 15 images were sufficient to track accurately through a video sequence of over a thousand frames.

The accompanying video shows the improvements obtained owing to our Kalman filtering operation performed both on the 2D features and on the camera projection parameters.

As seen in the results in Figures 5, 6, 7 and 8, the addition of the prior makes our results more reliable. The improvements can be seen in general but specifically around the mouth region where the changes in depth are most prominent. Areas where occlusions are common, around the eyes and mouth, stand to benefit the most as these regions are prone to loss of information in the 2D input. As our density estimator involves the Radial Basis Function as a kernel in the Parzen Window scheme, our objective function becomes non-linear, but because we initialize the parameters with the solutions from the previous frame, we do not run into popping in the animations owing to local minima. Nevertheless, the Parzen window density estimator is the prime bottleneck in taking this approach from merely online to real-time. In future work we will consider more efficient mechanisms for calculating prior probabilities such as Probabilistic PCA [TB99] or Gaussian Process Latent Variable models [Law04]. Our prior data consisted of animation sequences over 5 people and multiple facial expressions and speech movements to give us a total of 45817 frames of animation. Our automatic bandwidth estimation gives us a bandwidth of 0.3066 and 81 clusters for the lower face and a bandwidth of 0.1202 and 82 clusters for the upper face.

### 12.1. Quantitative Evaluation - Comparison with Ground Truth

In order to evaluate the results of our algorithm quantitatively, we needed to collect ground truth data for comparisons. For this purpose, we recorded both of our subjects using 5 high-definition cameras with baseline offsets that capture the actor's face from 5 dif-
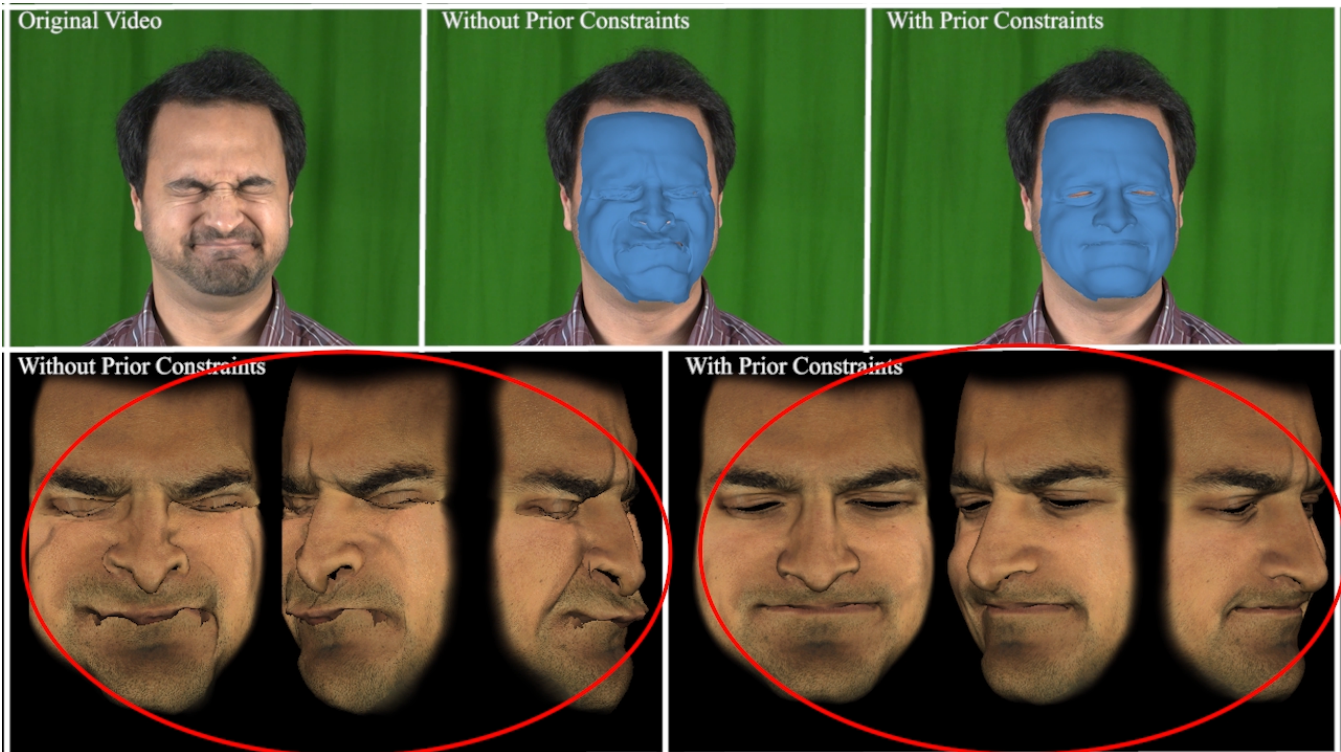
**Figure 5:** *Result showing a comparison between the solve using only 2D points vs the improvements obtained using our Prior constraint.*

ferent angles. These cameras were synced with each other using our gen-lock mechanism allowing us to obtain video frames from the performance sequence from different views that are temporally in correspondence with each other. We use the data from these cameras for a few expressions and reconstruct the 3D geometry of the actor's face for those frames using the Agisoft Photoscan software [ASP14] and obtain high detail textured meshes as shown in Figure 9. The obtained meshes have a random topology and hence we use the method of Amberg et al. [ARV07] to deform the neutral face of the actor to the reconstructed 3D meshes. This gives us ground truth geometry with the same topology as our Blendshapes which we use for quantitative evaluation. We then align the resulting shapes and the ground-truth geometry in 3D space by doing a rigid transformation after manually selecting corresponding points on the 3D meshes. We generate a heat map based on the per-vertex error of our results compared to the ground-truth. The results of this evaluation can be seen in Figure 10 and 11. The heat map was generated by using the error value per vertex (after scaling between 0 and 1) as a UV coordinate on a color gradient image. The gradient image goes from green at 0 to red at 1. Regions that have higher error show up in red and regions with lower error show up in green. The resulting quantitative errors are as explained in the image captions in Figures 10 and 11.

## 13. Comparison with reduced Blendshape Set

In the previous sections, our Blendshape model was comprised of 140 Blendshapes. We note that one of the factors affecting the

results of the solver is the fact that the Blendshapes themselves are not all compatible with each other i.e certain combinations of Blendshapes when triggered together result in a shape that is no longer within the feasible spectrum of shapes that a human face can assume. This is expected as the Blendshapes are not orthogonal to each other unlike in methods using PCA bases [TlTM11]. Our claim is that Blendshapes that are inherently more compatible with each other — i.e. Blendshape sets containing bases such that there are fewer combinations of shapes that result in implausible expressions when triggered together — will result in better solutions. Naturally, this implies that a Blendshape set having shapes that go relatively well with each other, will perform better. In order to test this hypothesis, we make use of the 'Emily' Blendshape set [ARL*10], which is comprised of 68 shapes, and solve for the Blendshape weights using our solver without the prior constraints and then quantitatively evaluate it with the results obtained using our original 140 Blendshapes, with and without the prior constraints.

Note that because the two Blendshape sets are inherently different from each other, we restrict the error calculation to the relevant facial regions comprising of the frontal face, excluding the outer vertices around the edges of the mesh. This is to ignore the error resulting from the edge vertices which occurs due to the different sizes of and the number of faces in the two Blendshape sets. Once again, we ensure that the topology of both the resulting Emily shapes and those from our original Blendshape set are identical by using the method of [ARV07] and follow it by a per-vertex error
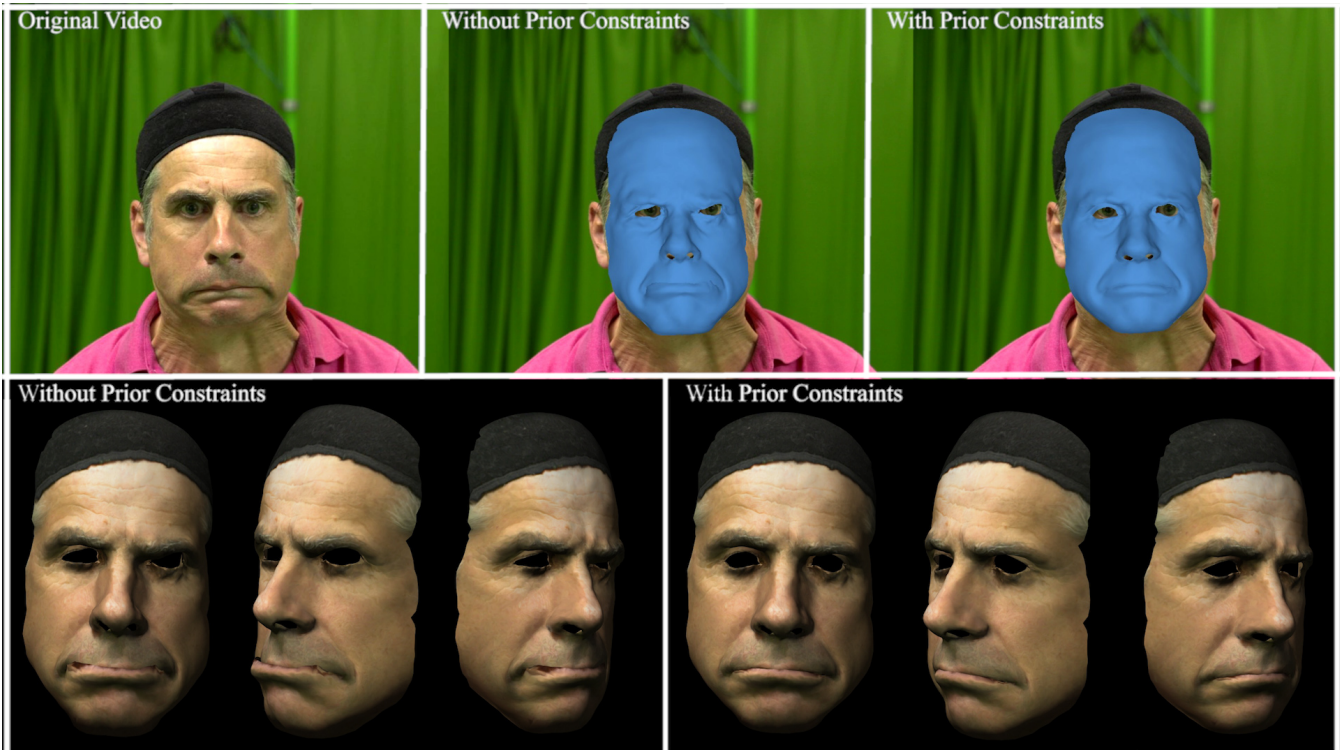
**Figure 6:** *Comparison of the frown expression using only 2D points vs the improvements obtained using our Prior constraint.*

calculation. The results of these comparisons can be seen in Figure 12. As can be seen, the error using the Emily set of Blendshapes is lower than those using our original 140 Blendshapes when we solve for both without using the prior constraints, but the results when we solve with the prior constraints are significantly lower in error when compared to the previous results. This shows that while the choice of Blendshape sets makes a difference, incorporating our prior in the solver provides significantly better solutions.

## 14. Conclusion

We have presented a lightweight markerless approach for capturing the facial movements of a user using only a 2D monocular input resulting in high quality animation parameters. We make use of prior constraints on the solve results, obtained from pre-existing accurate 3D captures and use this to improve the physical plausibility of our solve results which helps tackle the loss of information in depth that is inherent in 2D monocular inputs. We combine 2D landmark detections based on an ensemble of regression trees with an Active Appearance Model for improved accuracy. We make use of Kalman filters to handle noise across frames in an online fashion, both in 2D feature detection and in estimating the camera parameters giving us improved tracking and solves. We quantitatively compare the results of our solver which uses the spatial constraints with one that doesn't and show that our method generates better results when compared to the ground truth data. We also evaluate the effect that the choice of Blendshapes has on the results and quantitatively show that while the choice of shapes has an impact, overall

our solution using the spatial prior still generates results closer to the ground truth data.
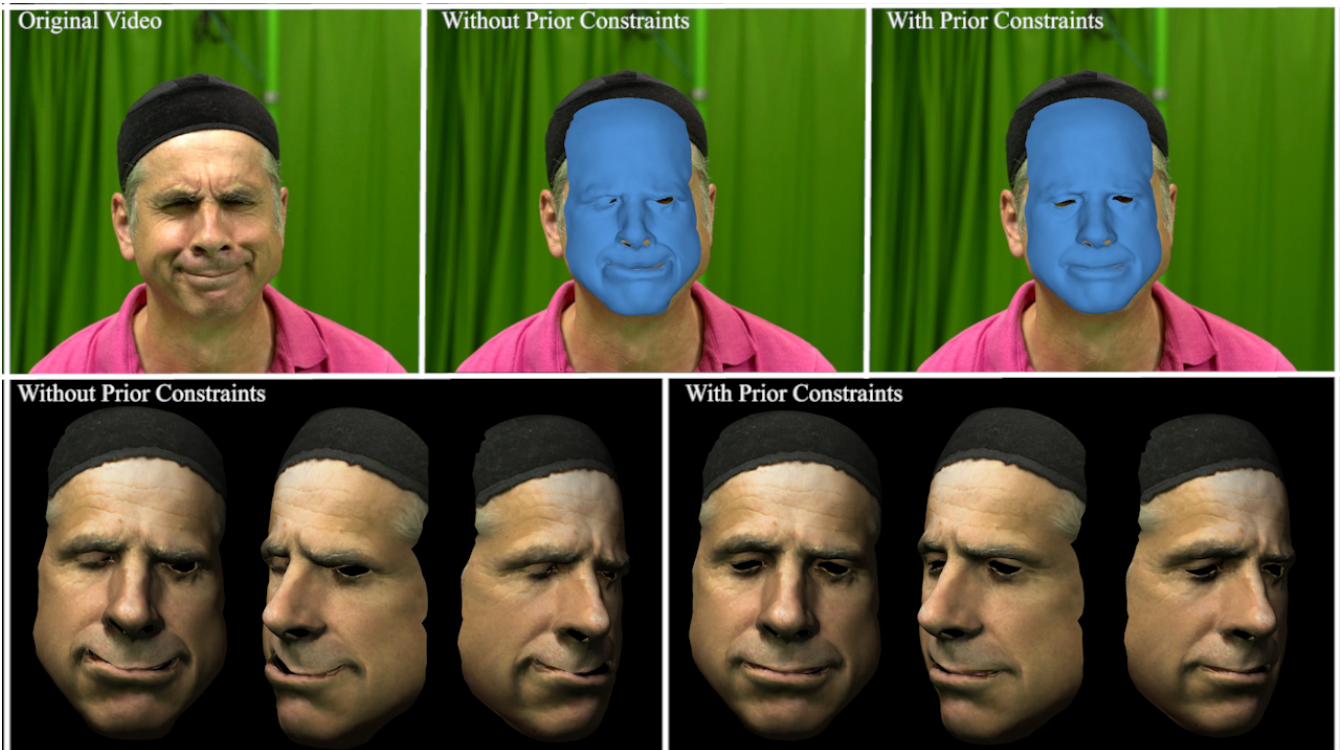
**Figure 7:** *Comparison of the angry expression using only 2D points vs the improvements obtained using our Prior constraint.*
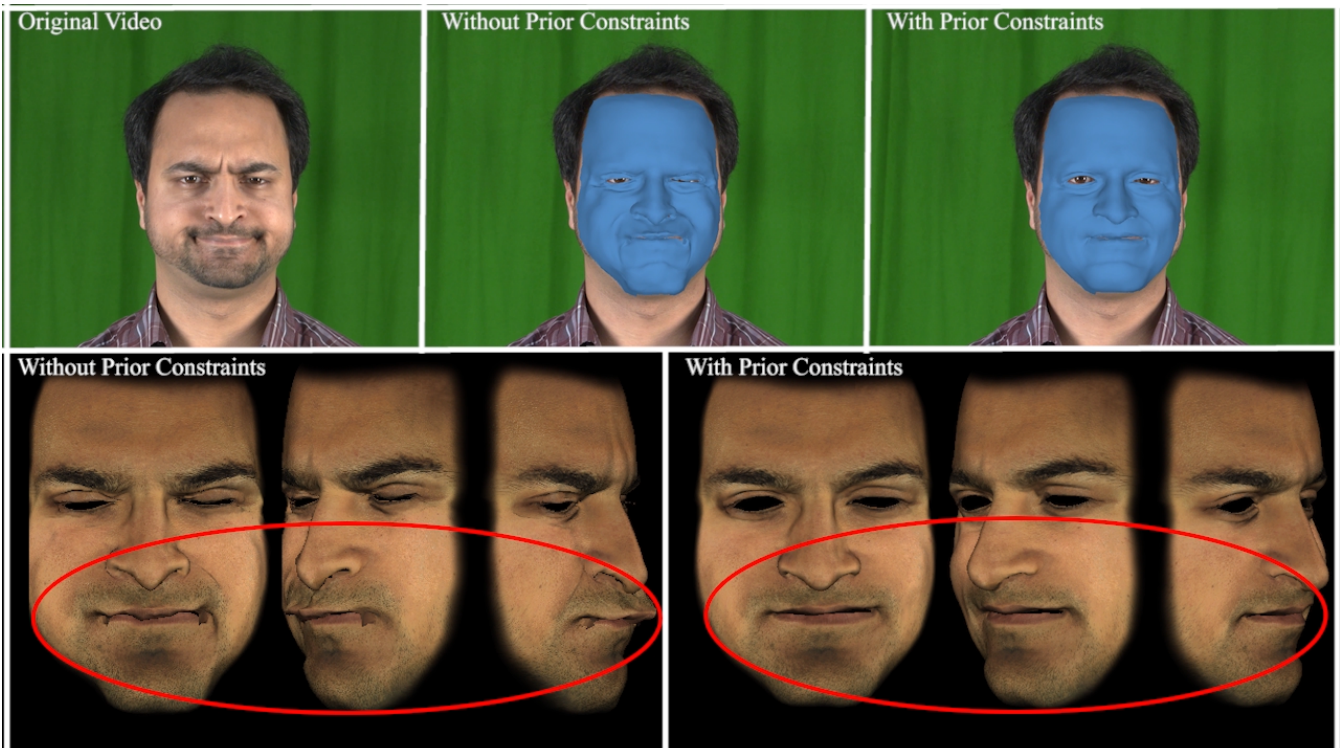


**Figure 8:** *Result showing a comparison between the solve using only 2D points vs the improvements obtained using our Prior constraint.*

**Figure 9:** *Ground-truth geometry obtained by reconstructing synchronized frames from multiple views of the actor's face.*



**Figure 10:** *Ground-truth evaluation for the lip-swing expression (top) and the frown expression (bottom). The total error from ground truth geometry for the lip-swing expression was: 4.4908e+05 cm$^2$ (without prior) and 3.3752e+05 (with prior). The total error from ground truth geometry for the frown expression was: 5.5391e+05 cm$^2$ (without prior) and 4.3288e+05 (with prior). The ground truth geometry was reconstructed from 5 different views of the actor's face.*
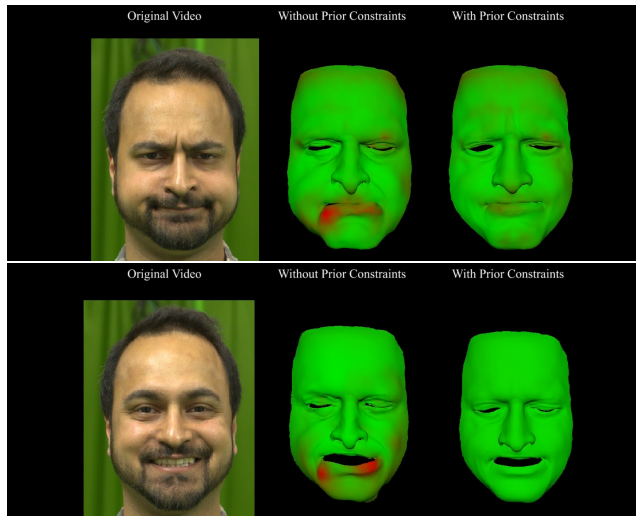


**Figure 11:** *Ground-truth evaluation for the angry expression (top) and the smile expression (bottom). The total error from ground truth geometry for the angry expression was: 2.06719e+05 cm$^2$ (without prior) and 1.78974e+05 (with prior). The total error from ground truth geometry for the smile expression was: 8.18939e+05 cm$^2$ (without prior) and 3.47592e+05 (with prior). The ground truth geometry was reconstructed from 5 different views of the actor's face.*



**Figure 12:** *Ground-truth evaluation for the OO expression (top) and the smile expression (bottom). The total error from ground truth geometry for the OO expression was: 1.53622e+05 cm$^2$ (140 Blendshapes, without prior), 4.6528+04 cm$^2$ (140 Blendshapes, with prior constraints) and 1.37559e+05 cm$^2$ (Emily Blendshapes).The total error from ground truth geometry for the smile expression was: 1.36572e+05 cm$^2$ (140 Blendshapes, without prior), 6.6960e+04 cm$^2$ (140 Blendshapes with, prior constraints) and 1.04310e+05 (Emily Blendshapes).*

# References

[ARL*10] ALEXANDER O., ROGERS M., LAMBETH W., CHIANG J.-Y., MA W.-C., WANG C.-C., DEBEVEC P.: The digital emily project: Achieving a photorealistic digital actor. *IEEE Computer Graphics and Applications 30*, 4 (2010), 20–31. 7

[ARV07] AMBERG B., ROMDHANI S., VETTER T.: Optimal step non-rigid icp algorithms for surface registration. *In Proc. of IEEE Computer Vision and Pattern Recognition (CVPR)* (2007). 4, 7

[ASP14] AGISOFT L., ST PETERSBURG R.: Agisoft photoscan. *Professional Edition* (2014). 7

[BGY*13] BHAT K. S., GOLDENTHAL R., YE Y., MALLET R., KOPERWAS M.: High fidelity facial animation capture and retargeting with contours. *Proc. 12th ACM SIGGRAPH/Eurographics Symp. Comput. Animat. (SCA)* (2013). 2

[BHB*11] BEELER T., HAHN F., BRADLEY D., BICKEL B., BEARDSLEY P., GOTSMAN C., SUMNER R. W., GROSS M.: High-quality passive facial performance capture using anchor frames. *ACM Trans. Graph 30*, 4 (2011). 2

[BHPS10] BRADLEY D., HEIDRICH W., POPA T., SHEFFER A.: High resolution passive facial performance capture. *ACM Trans. Graph. 29*, 4 (2010). 2

[BPL*05] BORSHUKOV G., PIPONI D., LARSEN O., LEWIS J. P., TEMPELAAR-LIETZ C.: Universal capture - image-based facial animation for the matrix reloaded. *ACM SIGGRAPH Courses* (2005). 2

[BV99] BLANZ V., VETTER T.: A morphable model for the synthesis of 3d faces. *ACM SIGGRAPH* (1999), 187âĂŞ194. 2, 4

[BWP13] BOUAZIZ S., WANG Y., PAULY M.: Online modeling for realtime facial animation. *ACM Trans. Graph 32*, 4 (2013). 2

[CBZB15] CAO C., BRADLEY D., ZHOU K., BEELER T.: Real-time high-fidelity facial performance capture. *ACM Transactions on Graphics (TOG) 34*, 4 (2015), 46. 2, 3

[CD12] CHACÓN J. E., DUONG T.: *Bandwidth selection for multivariate density derivative estimation, with applications to clustering and bump hunting*. Tech. rep., 2012. 6

[CET01] COOTES T. F., EDWARDS G. J., , TAYLOR C. J.: Active appearance models. *IEEE Trans. Pattern Anal. Mach. Intell. 23*, 6 (2001), 681âĂŞ685. 2, 3

[CHZ14] CAO C., HOU Q., ZHOU K.: Displaced dynamic expression regression for real-time facial tracking and animation. *ACM Transactions on Graphics (TOG) 33*, 4 (2014), 43. 2, 3

[CM13] CHACÓN J. E., MONFORT P.: A comparison of bandwidth selectors for mean shift clustering. *arXiv preprint arXiv:1310.7855* (2013). 6

[CWLZ13] CAO C., WENG Y., LIN S., ZHOU K.: 3d shape regression for real-time facial animation. *ACM Tran. Grap. 32*, 4 (2013). 2, 3, 4

[DCFN06] DENG Z., CHIANG P.-Y., FOX P., NEUMANN U.: Animating blendshape faces by cross-mapping motion capture data. *Proc. of symposium on Interactive 3D graphics and games (SI3D)* (2006). 2

[FH75] FUKUNAGA K., HOSTETLER L.: The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE Transactions on information theory 21*, 1 (1975), 32–40. 6

[FHW*11] FYFFE G., HAWKINS T., WATTS C., MA W.-C., DEBEVEC P.: Comprehensive facial performance capture. *Computer Graphics Forum 30*, 2 (2011), 425âĂŞ434. 2

[FJA*14] FYFFE G., JONES A., ALEXANDER O., ICHIKARI R., DEBEVEC P.: Driving high-resolution facial scans with video performance capture. *ACM Transactions on Graphics (TOG) 34*, 1 (2014), 8. 2

[FP09] FURUKAWA Y., PONCE J.: Dense 3d motion capture for human faces. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on* (2009), IEEE, pp. 1674–1681. 2

[GGW*98] GUENTER B., GRIMM C., WOOD D., MALVAR H., PIGHIN F.: Making faces. *Proc. 25th Annu. Conf. Comput. Graph. Interact. Tech* (1998). 2

[GVWT13] GARRIDO P., VALGAERTS L., WU C., THEOBALT C.: Reconstructing detailed dynamic face geometry from monocular video. *ACM Trans. Graph 32*, 6 (2013). 2

[HKV13] HOROVÁ I., KOLÁČEK J., VOPATOVÁ K.: Full bandwidth matrix selectors for gradient kernel density estimate. *Computational Statistics & Data Analysis 57*, 1 (2013), 364–376. 6

[Kin09] KING D. E.: Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research 10* (2009), 1755–1758. 3

[KS14] KAZEMI V., SULLIVAN J.: One millisecond face alignment with an ensemble of regression trees. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2014), pp. 1867–1874. 3

[LAR*14] LEWIS J. P., ANJYO K., RHEE T., ZHANG M., PIGHIN F., DENG Z.: Practice and Theory of Blendshape Facial Models. In *Eurographics 2014 - State of the Art Reports* (2014), Lefebvre S., Spagnuolo M., (Eds.), The Eurographics Association. doi:10.2312/egst.20141042. 2

[Law04] LAWRENCE N. D.: Gaussian process latent variable models for visualisation of high dimensional data. *Advances in neural information processing systems 16*, 3 (2004), 329–336. 6

[LWP10] LI H., WEISE T., PAULY M.: Example-based facial rigging. *ACM Trans. Graph 29*, 4 (2010). 2, 5

[LYYB13] LI H., YU J., YE Y., BREGLER C.: Realtime facial animation with on-the-fly correctives. *ACM Trans. Graph. 32*, 4 (2013). 2

[MJC*08] MA W.-C., JONES A., CHIANG J.-Y., HAWKINS T., FREDERIKSEN S., PEERS P., VUKOVIC M., OUHYOUNG M., DEBEVEC P.: Facial performance synthesis using deformation-driven polynomial displacement maps. In *ACM Transactions on Graphics (TOG)* (2008), vol. 27, ACM, p. 121. 2

[Rai04] RAITT: Presentation at u. southern california institute for creative technologiesâĂŹs frontiers of facial animation workshop, 2004. 2

[SAT*16] SAGONAS C., ANTONAKOS E., TZIMIROPOULOS G., ZAFEIRIOU S., PANTIC M.: 300 faces in-the-wild challenge: Database and results. *Image and Vision Computing 47* (2016), 3–18. 3

[SKSS14] SUWAJANAKORN S., KEMELMACHER-SHLIZERMAN I., SEITZ S.: Total moving face reconstruction. *ECCV 8692* (2014), 796–812. 2

[SLC11] SARAGIH J. M., LUCEY S., COHN J. F.: Real-time avatar animation from a single image. In *Automatic Face & Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on* (2011), IEEE, pp. 117–124. 2

[SLS*12] SEOL Y., LEWIS J., SEO J., CHOI B., ANJYO K., NOH J.: Spacetime expression cloning for blendshapes. *ACM Transactions on Graphics (TOG) 31*, 2 (2012), 14. 2

[SP04] SUMNER R. W., POPOVIC J.: Deformation transfer for triangle meshes. *ACM Trans. Graph 23*, 3 (2004). 4

[STZP13a] SAGONAS C., TZIMIROPOULOS G., ZAFEIRIOU S., PANTIC M.: 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *Proceedings of the IEEE International Conference on Computer Vision Workshops* (2013), pp. 397–403. 3

[STZP13b] SAGONAS C., TZIMIROPOULOS G., ZAFEIRIOU S., PANTIC M.: A semi-automatic methodology for facial landmark annotation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops* (2013), pp. 896–903. 3

[TB99] TIPPING M. E., BISHOP C. M.: Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology) 61*, 3 (1999), 611–622. 6

[TlTM11] TENA J. R., LA TORRE F. D., MATTHEWS I.: Interactive region-based linear 3d face models. *ACM Trans. Graph. 30*, 4 (2011). 2, 7

/

[TP13] TZIMIROPOULOS G., PANTIC M.: Optimization problems for fast aam fitting in-the-wild. In *Proceedings of the IEEE international conference on computer vision* (2013), pp. 593–600. 3

[TZN*15] THIES J., ZOLLHÖFER M., NIESSNER M., VALGAERTS L., STAMMINGER M., THEOBALT C.: Real-time expression transfer for facial reenactment. *ACM Transactions on Graphics (TOG) 34*, 6 (2015), 183. 2

[Vic] VICON: Vicon. *http://www.vicon.com*. 2, 5

[VWB*12] VALGAERTS L., WU C., BRUHN A., SEIDEL H.-P., THEOBALT C.: Lightweight binocular facial performance capture under uncontrolled lighting. *ACM Trans. Graph. 31*, 6 (2012), 187. 2

[WBLP11] WEISE T., BOUAZIZ S., LI H., PAULY M.: Realtime performance-based facial animation. In *ACM Transactions on Graphics (TOG)* (2011), vol. 30, ACM, p. 77. 2

[Wil90] WILLIAMS L.: Performance-driven facial animation. *ACM SIGGRAPH 24*, 4 (1990), 235–242. 2

[WLGP09] WEISE T., LI H., GOOL L. V., PAULY M.: Face-off: Live facial puppetry. *Proc. of ACM SIGGRAPH/Eurographics Symposium on Computer animation* (2009). 2

[ZSCS08] ZHANG L., SNAVELY N., CURLESS B., SEITZ S. M.: Space-time faces: High-resolution capture for˜ modeling and animation. In *Data-Driven 3D Facial Animation*. Springer, 2008, pp. 248–276. 2