FDLS: A Deep Learning Approach to Production Quality, Controllable, and Retargetable Facial Performances.

Wan-Duo Kurt Ma Weta Digital + Unity Wellington, New Zealand kma@wetafx.co.nz

Muhammad Ghifary* Weta Digital Wellington, New Zealand mghifary@gmail.com

Byungkuk Choi Weta Digital + Unity Seoul, Korea bchoi@wetafx.co.nz

Haekwang Eom Weta Digital + Unity Seoul, Korea heom@wetafx.co.nz

and difficult cases. The system has been under development for several years and has been used in major movies.

J.P. Lewis[†]

Weta Digital

Wellington, New Zealand

noisebrain@gmail.com

CCS CONCEPTS

• Computing methodologies \rightarrow Animation.

KEYWORDS

Facial animation, Deep learning, Motion capture, Optimization

ACM Reference Format:

Wan-Duo Kurt Ma, Muhammad Ghifary, J.P. Lewis, Byungkuk Choi, and Haekwang Eom. 2022. FDLS: A Deep Learning Approach to Production Quality, Controllable, and Retargetable Facial Performances.. In The Digital Production Symposium (DigiPro '22), August 7, 2022, Vancouver, BC, Canada. ACM, New York, NY, USA, 9 pages. https://doi.org/10.1145/3543664.3543672

INTRODUCTION 1



Figure 1: Facial Deep Learning Solve (FDLS) is an animator oriented tool system for solving facial animation given motion captured sparse marker set and limited training data.

Digital characters representing either humans or creatures are ubiquitous in the film and video game industries. Creating believable performances with these characters is one of the grand challenges of computer graphics. Realistic digital characters are often animated using performances captured from actors, as this allows the characters to inherit the "personality" of selected actors. On the other hand, performance capture presents its own difficulties arising from

ABSTRACT

Visual effects commonly requires both the creation of realistic synthetic humans as well as retargeting actors' performances to humanoid characters such as aliens and monsters. Achieving the expressive performances demanded in entertainment requires manipulating complex models with hundreds of parameters. Full creative control requires the freedom to make edits at any stage of the production, which prohibits the use of a fully automatic "black box" solution with uninterpretable parameters. On the other hand, producing realistic animation with these sophisticated models is difficult and laborious.

This paper describes FDLS (Facial Deep Learning Solver), which is Weta Digital's solution to these challenges. FDLS adopts a coarse-tofine and human-in-the-loop strategy, allowing a solved performance to be verified and (if needed) edited at several stages in the solving process. To train FDLS, we first transform the raw motion-captured data into robust graph features. The feature extraction algorithms were devised after carefully observing the artists' interpretation of the 3d facial landmarks. Secondly, based on the observation that the artists typically finalize the jaw pass animation before proceeding to finer detail, we solve for the jaw motion first and predict fine expressions with region-based networks conditioned on the jaw position. Finally, artists can optionally invoke a non-linear finetuning process on top of the FDLS solution to follow the motion-captured virtual markers as closely as possible. FDLS supports editing if needed to improve the results of the deep learning solution and it can handle small daily changes in the actor's face shape.

FDLS permits reliable and production-quality performance solving with minimal training and little or no manual effort in many cases, while also allowing the solve to be guided and edited in unusual

DigiPro '22, August 7, 2022, Vancouver, BC, Canada

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9418-5/22/08...\$15.00

https://doi.org/10.1145/3543664.3543672

^{*}currently with PT. Bank Rakyat Indonesia (Persero), Tbk. [†]currently with Google Research, Mountain View, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

the complexity of the models, the limited amounts of data available (to retain fidelity to a particular actor's performance of a particular character, it may not be possible to use data from other actors, nor even from the same actor performing a different character), and the need to potentially edit the resulting performances.

This paper describes FDLS, Weta Digital's solution to this challenge. FDLS has been in development since 2016 and has been used on movies such as *Gemini Man*. Our system has the following assumptions and requirements:

- **Production quality**, with high and verifiable accuracy. For example we assume that an actor's face shape may change slightly across days, and account for this. A stereo or mono head-mounted camera (HMC) with virtual (painted) markers is used in order to allow simultaneous capture of the head and body motion, since marker-free tracking methods are considered to not yet be sufficiently accurate.
- **Perceptual evaluation** is used (as is commonly the case in image quality assessment as well [Zhang et al. 2018]), with guidance provided by the VFX supervisor.
- Editable and retargetable results. Editing is needed to retain full creative control. A closed "black box" solver is not suitable because the supervisor may *require* that the performance be edited for various reasons. The actor's performance often needs to be re-targeted to other characters, such as humanoid aliens and monsters. Occasionally a particular expression may not be conveyed correctly on a character with very different proportions (e.g. a dragon). FDLS excels at post editing on rig (aka *puppet*) parameters such as FACS (Facial Action Coding System) action units [Ekman and Friesen 1978].
- We target **expert animators** rather than casual or novice users, while greatly reducing the effort needed to manipulate complex facial models,
- The ability to handle a rig with **nonlinear effects**, including intermediate and corrective blendshapes [Osipa 2010] and additional deformers e.g. for handling collisions between the lips.
- As mentioned above, we assume **limited training data** relative to the massive datasets commonly used to train deep neural networks. To address this issue we introduce strong inductive biases by way of specifically designed vertex graph features, that better capture the facial expression while remaining insensitive to irrelevant factors.
- We require performance capture that has been **stabilized** to remove rigid motion of the skull [Beeler and Bradley 2014; Lamarre et al. 2018]. We use a stereo or mono head-mounted camera, however the use of a head-mounted camera is not itself sufficient for stabilization since the camera may slip or the actor's scalp may move with respect to the skull.

This paper describes the "face solving" component of an overall character pipeline. Other components such as modeling, rigging, adaptation of the rig to specific actors [Ma et al. 2016; Seol et al. 2016], tracking, stabilization, muscle simulation, and rendering are separate topics not covered here.

The primary contributions of this paper are:

- A presentation of the motivation and design decisions for Weta Digital's Facial Deep Learning Solver (FDLS) system, a novel *multi-stage face solving* approach for production quality performance-driven facial animation using deep neural networks.
- Graph features that allow accurate animation solving with limited training data.
- A nonlinear post-finetuning method conditioned on jaw animation to allow better matching of some difficult face expressions.
- A method to accommodate the changes in the marker layout or actor's face shape without requiring re-training of the network.

2 RELATED WORK

The literature on facial performance capture dates from [Williams 1990] and is too large to fully review, for example [Zollhöfer et al. 2018] primarily targets the subset of monocular and optimizationbased approaches yet contains more than 200 references. We will focus on methods that using deep learning as well as those that give professional quality and editable results, while providing a few pointers to other literature. Broader surveys of facial performance capture and related areas include [Klehm et al. 2015; Tagliasacchi et al. 2016; Zollhöfer et al. 2018].

Computer vision research often targets fully automatic approaches that produce non-editable representations such as dense meshes or uninterpretable (e.g. PCA) parameters. [Blanz and Vetter 1999] introduced the morphable model approach in which scanned facial geometry and texture are approximated with PCA linear subspaces. This approach has been extended in a large body of research over the past two decades [Egger et al. 2020]. Performance capture of dense facial meshes usually requires complex camera setups [Bhat et al. 2013; Bradley et al. 2010; Fyffe et al. 2015] and multi-view stereo or photogrammetric reconstruction [Furukawa and Ponce 2009; Fyffe et al. 2015; Zhang et al. 2004]. Other research introduces algorithms for specific facial regions such as the lips [Garrido et al. 2016] and eyelids [Bermano et al. 2015].

Real-time tracking and re-targeting from monocular RGB video has seen widespread use in face "filters" for video conferencing systems. While these methods have impressive ease of use and robustness, they do not offer *verifiable* accuracy and may not capture fine facial detail, as well as being prone to errors from significant lighting changes and occlusions. See [Zollhöfer et al. 2018] for a recent survey of the literature in this area.

The facial animation solvers used in visual effects production and games often use physical or virtual (painted) markers as input, due to accuracy requirements and the reliability of marker-based systems [Bickel et al. 2007; Huang et al. 2011]. However, these approaches have less ability to capture fine detail in complex regions of the face, in particular around the lips. One of the accommodations in these marker-driven methods is to add additional information extracted from the captured footage. For instance, [Bhat et al. 2013] adds separately tracked contours of the lips and eyelid to the solver's input.

Keyframe animation is widely used as it allows artists to specify intuitively understandable parameters at a subset of frames, and then interpolate these parameters to the remaining frames using (for example) Catmull-Rom [Li and Deng 2008] or B-spline [Choe et al. 2001] curves. On the other hand, keyframe animation is quite laborious, and various techniques have been proposed to speed up the artist's workflow. For example [Seol et al. 2011] introduced a successive refinement scheme in which large motions (e.g. the jaw) are solved before the fine details. Temporal animation editing has been developed in several ways. [Seol et al. 2012] introduce a gradient-domain space-time editing method. [Dinev et al. 2018] use user guidance in the form of 2D sketches for generating plausible lip corrections on top of keyframes. [Berson et al. 2019, 2020] use a learning based approach to perform temporal animation editing.

Performance capture systems developed in production settings target high quality and (in some cases) editable rig representations. [Smith et al. 2017] demonstrates high-quality blendshape solving using an optimization approach. Convolutional neural networks have been employed to regress directly from video to dense meshes [Laine et al. 2017]. [Moser et al. 2018] describes a production-proven system for regressing 3D marker positions. In addition to the research reviewed above, there are also several commercial systems that address aspects of professional performance capture, including DI4D, Dynamixyz, Synthesia, Medusa, and Imagemetrics [Various 2022].

In summary, there is a large body of research on facial performance capture, including both fully automatic neural approaches, and approaches that target editable high quality performance capture using classic methods. However, there is relatively little existing research on deep learning methods that target artist-in-the-loop, editable, and high-quality performance capture.

3 HYBRID METHOD FOR FACIAL ANIMATION CAPTURE

We focus on solving performance-driven facial animation problems with blendshape models, i.e., inferring the blendshape weights $\hat{\mathbf{w}}$ such that the corresponding blendshape expression matches the captured actor's performance. Specifically, we denote a face model by $g: \mathcal{W} \to \mathcal{X}$, where $\mathcal{W} \subseteq \mathbb{R}^D$ is the blendshape weight space and $\mathcal{X} \subseteq \mathbb{R}^{3n}$ is the face expression space, each element of which is the *n* vertices with the coordinates vectorized as $[x_1, y_1, z_1, \dots, x_n, y_n, z_n]$. Correspondingly, each blendshape is given by a vector $\mathbf{b}_k \in \mathbb{R}^{3n}$. In the linear case a blendshape model in the "delta" formulation is:

$$g(\mathbf{w}) = \mathbf{b}_0 + \sum_{k=1}^{D} w_k (\mathbf{b}_k - \mathbf{b}_0) = \mathbf{b}_0 + \mathbf{B}\mathbf{w},$$
 (1)

where w_k are the blendshape weights (typically $0 \le w_k \le 1$), \mathbf{b}_0 corresponds to the neutral shape, and $\mathbf{B} \in \mathbb{R}^{3n \times D}$ contains the "delta" blendshape targets $\mathbf{b}_k - \mathbf{b}_0$.

Let us define a distance measure $\mathcal{M} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. Given a captured actor expression at a particular frame $\mathbf{x} \in \mathcal{X}$, we seek $\hat{\mathbf{w}}$ such that $\mathcal{M}(g(\hat{\mathbf{w}}), \mathbf{x})$ is small. In our case, $g(\mathbf{w})$ and \mathbf{x} represent sparse markers attached on either the face puppet or the real actor's face. In general, we can frame the problem either as a face matching problem or a regression problem.

Matching approach. A prevalent approach to solving the blendshape weights is to frame it as a *face matching* optimization problem [Choe et al. 2001]. Choosing Euclidean distance as $\mathcal{M}(\cdot, \cdot)$ and setting $g(\mathbf{w})$ as in (1), we can express the problem as the following optimization:

$$\hat{\mathbf{w}} := \arg\min_{\mathbf{w}} \|\mathbf{b}_0 + \mathbf{B}\mathbf{w} - \mathbf{x}\|_2^2 \quad \text{s.t.} \quad 0 \le \mathbf{w} \le 1,$$
(2)

which is a constrained quadratic programming (QP) problem. The solution can be obtained by applying a standard QP solver run on each frame [Boyd and Vandenberghe 2004; Fascione 2017].

Learning-based approach. Instead of solving the face matching problem by optimization, one can also apply a regression approach to inferring the blendshape weights $\hat{\mathbf{w}}$ from the captured expression \mathbf{x} through a multivariate regression function $f_{\theta}: X \to \mathcal{W}$. The regression function is often decomposed as $f = h \circ \phi$, where we call $\phi : X \to \mathcal{H}$ the *feature extractor* and $h : \mathcal{H} \to \mathcal{W}$ the *regressor.* A machine learning technique can identify the function $f_{\hat{\theta}}$ given a set of labeled examples $\mathcal{D} = \{(\mathbf{x}_v^{(i)}, \mathbf{w}_v^{(i)})\}_{i=1}^N$ [Vapnik 1998]. This trained regressor is expected to predict valid blendshape weights given the *previously unseen* captured expression \mathbf{x}_u , i.e., $\hat{\mathbf{w}}_u = f_{\hat{\theta}}(\mathbf{x}_u)$, such that the 3d model expression $g(\hat{\mathbf{w}}_u)$ is correctly generated.

Nonlinear rig. For simplicity, the preceding description presents the matching and learning-based approaches in terms of a linear blendshape model. Our blendshape puppets incorporate natural non-linear expressions, for example using in-between and corrective shapes [Osipa 2010; Seo et al. 2011]. In addition, the overall puppet consists of a large deformation chain with numerous tweaks, deformers and skin clusters applied on top of the blendshape system. A full description of the puppet is outside the scope of this paper. The details of the puppet are not needed to understand the application of the learning and matching approaches, however: the learning-based approach captures the puppets's behavior (including nonlinearities) since the puppet itself is used to produce training data (Section 4.1). For the matching approach, rather than seek a global optimum, our optimization fine-tunes a local linear subset of parameters that are selected by the artist (see Hybrid approach, next, and Section 4.5).

Hybrid approach. In FDLS, we apply the learning-based approach as the main animation solver, specifically using deep learning regression [He et al. 2016], followed by the matching approach to fine-tune some facial parameters. This hybrid approach has two main advantages. First, the main animation solving is reduced to the forward pass of the network f_{θ} , which is time-efficient. Second, the forward pass provides an initialization for the nonlinear matching optimization that is in the neighborhood of the correct local minimum and is close to the desired solution, resulting in temporally consistent and efficient fine tuning to generate high-fidelity final animations. In the actual implementation, FDLS comprises several components to enable the human-in-the-loop solving process. These are fully elaborated in the next section. DigiPro '22, August 7, 2022, Vancouver, BC, Canada



Figure 2: The Facial Deep Learning Solver (FDLS) comprises three main stages: data preparation, training and solving, where the letters denote: (A) salient sample selection (Section 4.2); (B) feature engineering (Section 4.3); (C) jaw kinematics conditional training (Section 4.4.1); (D) shape alignment (Section 4.4.2); (E) finetuning (Section 4.5); (F) anchor poses (Section 4.4.2).

4 THE COMPLETE FDLS PIPELINE

As shown in Fig. 2, the end-to-end FDLS pipeline consists of two main stages: i) the development stage, and ii) the operation stage. This partitioning forms a *human-in-the-loop* MLOps lifecycle in Weta Digital's production that ensures reproducibility, testability, and maintainability of our deep learning system. In the development stage, data preparation steps are performed, including synthetic data generation (Section 4.1), salient sample selection (Section 4.2), feature engineering (Section 4.3), and training the deep learning models (Section 4.4). The animation solving happens in the operation stage, which makes use of the trained deep learning models as the main solvers complemented by a few other techniques: shape alignment with anchor poses (Section 4.4) and fine-tuning optimization (Section 4.5).

4.1 Synthetic Data Generation and Augmentation

The training data for FDLS are in the form of 3d virtual markers created from the actor's facial landmarks of his/her neutral expression using photogrammetry – we call these "neutral markers" throughout the paper for brevity. We then attach the neutral markers to the corresponding actor's blendshape puppet and synthetically generate the training markers driven by the puppet motion $\mathbf{x}_v^{(i)} = g(\mathbf{w}_v^{(i)})$ forming a set of training tuples $\mathcal{D} = \{(\mathbf{x}_v^{(i)}, \mathbf{w}_v^{(i)})\}_{i=1}^N$. Note that the training blendshape weights $\mathbf{w}^{(i)}$ are designed and selected by animation experts, typically containing both FACS and some actor-specific expressions. Since the training markers are driven by the puppet, the training data thus implicitly includes nonlinear effects from the puppet without requiring a differentiable chain of deformations.

Additionally, we also utilize a data augmentation or oversampling technique to further increase the variation of the synthetic markers \mathbf{x}_v and therefore avoid overfitting of the trained deep learning models. This is achieved by randomly injecting Gaussian noise in the original markers $\mathbf{x}_e = \mathbf{x}_v + \mathcal{N}(0, \Sigma)$. We empirically determine the covariance Σ by analyzing the jitter of the tracked markers, for instance, a marker at the nose bridge has less variance than those on the chin region. This can produce a significant quality gain in the solved animations.

4.2 Salient Training Sample Selection

As is mentioned in Section 4.1, the training data for FDLS can be flexibly generated with various combinations of simple "one-hot" single action unit FACS expressions (producing a minimum viable solution) and more complex and realistic data from various sources such as dynamic scans of the actor, previously animated or solved performances, etc.

While we generally expect that more realistic and complex training data should result in better performance, counterintuitively we observed that this is not always the case. We found that the root cause is data imbalance – for example, some solved performances may consist mostly of a neutral expression. This issue can be avoided if the additional training examples are carefully selected by the animators, but manually selecting from among many training expressions is laborious.

To simultaneously encourage data diversity and salience, we introduce a salient sample selection technique to automatically remove redundant examples from the training set. This eliminates the need for tedious manual sample selection. We reduce the number of training examples from N to M < N by sequentially applying the following sample selection rule

$$s_i = \frac{1}{N} \left(\sum_{j=1}^N k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right) < \sigma, \tag{3}$$

where $s_i \in \{0, 1\}$ is a binary variable indicating whether the shape i is selected, σ is a tunable threshold, and $k : X \times X \rightarrow [0, 1]$ is the similarity measure. In the extreme case, this salient sample selection picks only a single datapoint from the dataset \mathcal{D} if it contains mostly similar expressions according to Eq. (3), e.g., neutral non-dialogue expressions. We choose an exponential Radial Basis Function (RBF) [Broomhead and Lowe 1988] representing $k(\cdot, \cdot)$ which provides a good set of salient shapes due to its effectiveness in capturing similarities in a non-linear feature space.

4.3 Feature Engineering

Solving animation using solely sparse marker coordinates $\mathbf{x} = [\mathbf{m}_1, \ldots, \mathbf{m}_n]$, where $\mathbf{m}_i \in \mathbb{R}^3$ to represent faces has its own limitations (e.g., [Seol et al. 2016]). Sparse markers are less informative than the full face representation (e.g., facial depth from LiDAR),

and it is not possible to fully characterize certain facial expressions using only marker coordinates. To partially address this and simultaneously help with the limited amount of training data available in this domain, we developed several *graph features* to increase the information richness of the input representations, thereby inducing robust and highly accurate deep learning solves.

In the initial enthusiasm following the success of [Krizhevsky et al. 2012], it was argued that features emerging from neural net training should outperform manually designed features aka "feature engineering". While this is arguably true especially in computer vision and natural language processing domains, engineered features continue to be used in state-of-the-art deep learning – the positional encoding in transformers and Fourier/positional encoding that distinguishes NeRF [Mildenhall et al. 2020] are prominent examples. In practice, the need for engineered features can be justified when it is impractical to search across a sufficient variety of architectures to find those that result in ideal features, or to impose an inductive bias in cases where "correct" features might not be discovered due to limited data – as is often the case in visual effects!



Figure 3: FDLS graph features. From left to right: pairwise distance, pairwise directions, delta from neural to the target expression.

FDLS benefits from carefully engineered features extracted from the marker coordinates. Fig. 3 illustrates the variants of the feature extractor $\phi : X \to \mathcal{H}$ that we design as part of the learning-based approach (recall Section 3). These are described next.

Pairwise distance. The pairwise distance features $\boldsymbol{\phi}^{\text{dist}} \in \mathbb{R}^{n^2}$ are a feature vector, whose values represent the Euclidean distances over all possible combinations of two marker coordinates in the input representation $\{\mathbf{m}_i\}_{i=1}^n$. The *k*-th element of the pairwise distance feature $\boldsymbol{\phi}_k^{\text{dist}} \in \mathbb{R}$ is given by $\boldsymbol{\phi}_k^{\text{dist}} = \|\mathbf{m}_i - \mathbf{m}_j\|_2$.

Pairwise direction. The pairwise direction features $\phi^{dir} \in \mathbb{R}^{3n^2}$ are a collection of the directional vectors from one coordinate to another. This gives complementary information to eliminate the inherent ambiguity in ϕ^{dist} . For instance, the Euclidean distance of two particular coordinates between different expressions, e.g. neutral and smiling, could be indistinguishable due to symmetry.

Analogous to the pairwise distance computation, the *k*-th coordinate $\boldsymbol{\phi}_k^{\text{dir}} \in \mathbb{R}^3$ is computed as $\boldsymbol{\phi}_k^{\text{dir}} = (\boldsymbol{\phi}_k^{\text{dist}})^{-1}(\mathbf{m}_i - \mathbf{m}_j)$. Note that the pairwise distance and direction features are both invariant to small slippage of the head camera.

Delta pose. Lastly, we define the delta pose features $\phi^{\text{delta}} \in \mathbb{R}^{3n}$, which are the coordinate differences between a particular face expression and its neutral pose. These features are invariant to the global displacement of the marker coordinates and have become a common representation used in traditional blendshape systems

in the graphics community [Lewis et al. 2014]. Given the specific expression **x** and the corresponding neutral pose \mathbf{x}_0 , the delta pose feature vector $\boldsymbol{\phi}^{\text{delta}} \in \mathbb{R}^{3n}$ is calculated as $\boldsymbol{\phi}^{\text{delta}} = \mathbf{x} - \mathbf{x}_0$.

The complete set of input features for FDLS is the concatenation of these three features, i.e., $\boldsymbol{\phi} = [\boldsymbol{\phi}^{\text{dist}}, \boldsymbol{\phi}^{\text{dir}}, \boldsymbol{\phi}^{\text{delta}}] \in \mathcal{H}$. Intuitively, these features are a graph-like representation of facial expressions that contains not only the positional information but also pairwise relationships over the marker coordinates. These features can be directly applied in our region-based training by extracting the features only in a specific face region, as will be shown later. Graph features have been used in concurrent work such as [Qi et al. 2017]. Our use is somewhat distinguished in that we intentionally do not introduce convolution. The translational equivariance of convolution is appropriate for recognition of general point clouds (for example an edge may appear at any location) however the features of a canonically positioned face are in relatively fixed locations and we seek only to capture changes due to expression.

4.4 Region-based Training and Solving

The deep learning regression f_{θ} defined in Section 3 is the core of FDLS. In its implementation, we define seven deep residual networks, each responsible for a particular face region containing a subset of the sparse markers. We denote the face regions by $R \in \{upper-face, lower-face, jaw, lips, cheek, eye-lids, eyeballs\}$. The marker subset is selected according to the semantic meaning of the face region, but a marker on the nose bridge is always included regardless of the region to insure sufficiently diverse and spatially supported features. Each face region is also linked only to a subset of blendshape weights or channels that are relevant to that region. In other words, we train $\{f_{\theta_r} : X_r \to W_r; \forall r \in R\}$.

Solving with the trained deep learning models is straightforward, i.e., we simply run the forward pass $\hat{\mathbf{w}}_r = f_{\hat{\theta}_r}(\mathbf{x}_r)$ for each face region on a per frame basis. The full solved weights are the concatenation of $\{\hat{\mathbf{w}}_r\}_{r \in \mathbb{R}}$, which we refer to as the *raw animation* – later this will be used as the basis for editing or fine-tuning. As described next, a substantial improvement in the quality of the raw animation is obtained by extending the forward pass with two additional techniques.



Figure 4: The network $f_{\hat{\theta}_r}(\mathbf{x}_r)$ for a regular face region (left) and a jaw-conditioned network $f_{\hat{\theta}_r}(\mathbf{x}_l; \hat{\mathbf{w}}_{jaw})$ (right).

4.4.1 Jaw Kinematics Conditional Training. Jaw kinematics plays an important role in affecting the shape and expressions of the face, especially the lower half [Yang et al. 2019]. For instance, one can have different smiling expressions under different jaw conditions (open, close, sideway movement, etc). In other words, the effect of face muscle activations controlling the lips and cheeks strongly depend on the jaw position. For this reason, in our pipeline the jaw motions are established *before* solving other regions. FDLS supports the artist in this process by providing a regression specifically for the jaw position. Specifically, we denote the jaw activation weights by $\mathbf{w}_{jaw} \in \mathbb{R}^z$ where in our production face models, normally z = 3, namely, jaw0pen, jawThrust, and jawSideways. The jaw motions $\hat{\mathbf{w}}_{jaw}$ are obtained from computing $f_{\hat{\theta}_{jaw}}(\mathbf{x}_{jaw})$ and are used in the following regions $l \in \{lower-face, lips, cheek\}$. We solve the animation by evaluating $f_{\hat{\theta}_i}(\mathbf{x}_l; \hat{\mathbf{w}}_{jaw})$, which is trained by the risk minimization:

$$\hat{\theta}_{l} := \arg\min_{\theta_{l}} \frac{1}{N} \sum_{i=1}^{N} \ell\left(f_{\theta_{l}}(\mathbf{x}_{l}^{(i)}; \mathbf{w}_{\text{jaw}}^{(i)}), \mathbf{w}_{l}^{(i)}\right), \qquad (4)$$

where $\ell : W x W \to \mathbb{R}$ is the standard regression loss function. Note the use of \mathbf{w}_{jaw} (the groundtruth jaw motion) in the training, versus $\hat{\mathbf{w}}_{jaw}$.

As we can see the solved jaw weight vector $\hat{\mathbf{w}}_{jaw}$ is part of the inputs to the deep learning model f_{θ_l} . However, due to its low dimensionality, we embed it as additional nodes in a hidden layer after the series of deep residual blocks rather than as part of the input layer, as shown in Fig. 4.

4.4.2 Shape Alignment with Anchor Poses. In real productions lasting several days or more, the positions of markers on a single actor's face may shift over time. There are several reasons for this. First, although the landmarks are painted using a perforated mask and their relative positions do not change, the mask may be placed in slightly different positions on different days. Second, the accuracy of the tracked and 3d-reconstructed markers can be inconsistent across different shots, depending on their complexity. Lastly, the actor may slightly gain or lose weight. These changes can introduce a small shape mismatch between the puppet-snapped markers and the performance-captured markers.



Figure 5: Shape alignment using an anchor pose. The green and red circles represent the puppet-snapped markers and the performance-captured markers, respectively.

To mitigate this problem, we introduce a simple shape alignment method that can be performed directly by an animator before (re)running the deep learning solver $f_{\hat{\theta}}(\cdot)$. The basic idea is to first calculate the offsets between the synthetic markers and the shot-captured markers producing an *anchor pose*, i.e., a hand-picked expression at a particular frame in a shot, and then propagate the offsets to displace the other face markers within the entire frame range. Specifically, suppose that we have a performance shot with *T* frames and an anchor pose at frame *a* is chosen. The artist then

sets the anchor weights w_a such that $g(w_a)$ portrays the desirable anchor pose. The shape alignment is executed as follows:

$$\mathbf{d}_{a} = \mathbf{x}_{a} - g(\mathbf{w}_{a}),$$

$$\tilde{\mathbf{x}}_{f} = \mathbf{x}_{f} - \mathbf{q}_{fa} \odot \mathbf{d}_{a}, f \in \{1, ..., T\},$$
(5)

where \mathbf{d}_a are the marker offsets and $\mathbf{q}_{fa} \propto \exp\left(-\|\mathbf{x}_a - \mathbf{x}_f\|_2^2\right)$ denotes the weighting factors of the alignment by calculating the pose similarity. After this we can rerun the deep learning solver $\tilde{w} = f_{\hat{\theta}}(\tilde{\mathbf{x}}_f), \forall f = \{1, ..., T\}$ without retraining to produce a more accurate animation. Note that this shape alignment process can be done iteratively by choosing a few other anchor poses as needed – for the first anchor pose, all elements of \mathbf{q} are set to one.

4.5 Fine-tuning Optimization

Editing can optionally be performed once the raw animation has been solved by the deep learning model $\hat{\mathbf{w}}_u = f_{\hat{\theta}}(\mathbf{x}_u)$. In cases with unusual expressions, we sometime see a small discrepancy between the solved expression of the actor puppet $g(\hat{\mathbf{w}}_u)$ and the actual captured expression \mathbf{x}_u , mainly because the shape \mathbf{x}_u has not been well represented in the training dataset for constructing $f_{\hat{\theta}}$. The artists can simply adjust the solved weights (most of the time only a few channels are needed) to get more accurate expressions.

To make the aforementioned process more scalable, we provide a *fine-tuning optimization* as a post-processing step in FDLS. It essentially further minimizes the error $||g(\hat{\mathbf{w}}_u) - \mathbf{x}_u||_2$. We run the minimization by solving the face matching problem in Eq. (2), in which the deep learning solved weights are set as the starting point for a minimization using the L-BFGS algorithm as the optimizer [Liu and Nocedal 1989]. We found this provides better results than starting the minimization from zero or a random point. Note that we also form a restricted blendshape basis (matrix **B** in Eq.(2)) that only includes artist-selected degrees of freedom to be finetuned while the remaining weights remain frozen.

With this implementation, we provide the artists with flexible options such as choosing only a specific frame range, partially selecting the face markers as objects of comparison, and selecting only a few blendshape weights to be finetuned. This makes the fine-tuning optimization an interactive tool as part of an overall editing flow to efficiently finalize high quality animations.

5 EXPERIMENTS

In this section we review the end-to-end FDLS pipeline and present the main results as well as ablation tests. While the use case shown here is for a specific actor and a single shot performance, the workflow is generic for any actor and performance capture.

Next we describe the implementation details of the FDLS training and solving stages as depicted in Fig. 2.

5.1 Training

As described in Section 4.1, we synthetically manufacture the training data in the form of a series of sparse markers by assigning animation weights to the blendshape channels. In the FDLS training, two types of training animations are utilized: "one-hot" FACS expressions and the range of motion (ROM) of an actor. The onehot FACS animation has individual FACS expressions such as jaw drop (AU26) and lip corner puller (AU12) at successive frames, thus FDLS: A Deep Learning Approach to Production Quality, Controllable, and Retargetable Facial Performances.



Figure 6: FDLS workspace and GUI. The system is implemented using the commercial Maya API [Autodesk, Inc. 2019].

representing the full extent of the blendshape basis. By using only the one-hot FACS training animation, we can already construct a *minimum viable solver* before more advanced training animation is available.

The ROM captures the motion of the specific actor. This is more realistic than the one-hot animation, as some people have difficulty activating individual muscles in isolation. One can think of the FACS-based animation as determining the range and extremes of individual muscles, while the ROM provides information about their distribution, for example, what muscle combinations actually occur and do not occur for the particular actor. In addition, the ROM allows an expert artist to specify exactly which blendshape weights contribute to an expression, thus guiding the solver to disambiguate cases where nearly (but not exactly) identical marker positions can be produced with different weight configurations.

5.2 Solving

Fig. 6 shows the FDLS workspace and GUI. Note that we do not require the artist to specify any parameters to perform solving, thus the tool is convenient to use without significant artist training. Our users have reported that FDLS system takes only about one week to train a new animator, which boosts the production speed.

After the workspace has been setup (Fig. 6), FDLS allows the artist to pick a neutral-ish frame as an anchor frame (Section 4.4.2), and correct the jaw animation given the raw jaw solve, then run the jaw-conditioned networks (Section 4.4.1). After raw solve, one or more additional anchor poses can be added to refine the projected markers. Finally, the post-finetuning optimizes the face conditioned on jaw channels to match the performance markers.

Fig. 7(a) shows a complex FDLS solve that illustrates most aspects of the system. The top row depicts the raw solved expressions induced by the region-based models $\{f_{\theta_r}(\cdot)\}_{r \in \mathbb{R}}$ after shape alignment using one neutral-ish frame. The middle row displays the solved expressions with shape alignment using the additional anchor pose highlighted in red followed by finetuning. This shows a substantial improvement in that the result is close to the final desired expressions (bottom).

Fig. 7(b) compares the performance of post-finetuning versus the raw solve, with or without an anchor pose at the marked frame. The minimization ensures the RMSE loss value decreases (e.g., curves A and B), however the error remains high across frames in this case.

By adding one anchor pose at around frame 55 (shown in (a) with the red box), the error is globally decreased with no further effort.

Fig. 7(c) shows the results of salient frame selection by varying the threshold σ in (3). Specifically, it compares the test performance with $\sigma = 0.1$, $\sigma = 0.3$, $\sigma = 0.5$, and the full dataset, corresponding to using 50%, 66%, 83% and 100% of the original dataset, respectively. One can see that the network trained with data selected using $\sigma = 0.3$ (orange) has better performance than the others, including the network trained with the full dataset (red).

Many performances include a certain amount of repeated information, for instance, in some performances a large percentage of the frames show a mostly neutral expression. The experiment in Fig. 7(c) provides evidence of an important point, which is that the solve actually *improves* when this redundant data is removed. However, a limitation of this approach is that the σ must be empirically tuned for the particular dataset. If σ is set to an extreme (e.g., $\sigma = 0.1$) or moderate value (e.g., $\sigma = 0.5$) it may give poor results.

6 CONCLUSION AND FUTURE WORK

This paper presents a human-in-the-loop deep learning based approach to facial animation solving. The FDLS system was initiated in 2016 and has seen ongoing development and improvement since then. It presents a simple and light-weight training and solving pipeline to our artists. The approach is suitable for driving sophisticated and non-linear blendshape rigs with hundreds of parameters. FDLS produces production quality results while requiring limited training data and artist effort. Both the design trade-offs and several specific features of the system are guided by consideration of artist practice and differ from other published work.

The graph features have proven successful at allowing accurate solves with limited training data. The quality of the intermediate solution can be verified and adjusted if needed following both the jaw estimation and prior to the post fine-tuning, and since FDLS produces interpretable parameters, the freedom to do further editing on the final solution is fully preserved. The anchor pose mechanism allows the network to adapt to slight changes in both the marker positions and the actor's face shape without re-training.

In our experience animators can learn to use the system in approximately a week and a large majority of shots require little or no human editing. While we could cherry pick results to provide an "objective" characterization of the production performance, in all honesty such a characterization is almost meaningless. This is due to wide variety of different performance and retargeting situations, the different ways that FDLS can be used, the quality of training data, choice of anchor poses, and other factors, and particularly because of the ultimate perceptual evaluation of the results. Instead, we highlight that FDLS succeeds as a system that can be employed in *all* these scenarios – from fully automatic solves to extremely challenging shots that require iterative solution and evolving guidance from a supervisor.

One of the disadvantages of using a sparse marker set is the ambiguity of the lip positions in certain expressions. In the future we would like to enhance the solver with additional information such the lip contours or LIDAR depth.



(a) FDLS main result

(c) Salient sample performance

Figure 7: A complex FDLS solve example. (a) raw solve with shape alignment from a single neutral-ish frame (top), finetune with the additional anchor pose outlined in red (middle), and the final manual edits from artist on top of the FDLS solve (bottom). (b) and (c) are the RMSE performance of ablation experiments at the frame marked with star with/without an anchor pose and finetuning, and on the salient data selection, respectively. Note that the RMSE generally includes an irreducible offset due to slight mismatch between the actor and puppet virtual marker positions.

ACKNOWLEDGMENTS

We thank Joe Letteri, Marco Revelant, Luca Fascione, Dejan Momcilovic, Stephen Cullingford, Stuart Adcock, Allison Orr, David Luke, Millie Maier, Andrew Moffat, Kenneth Gimpelson, Nivedita Goswami, and Zhicheng Ye for supporting this project. Moreover, we appreciate the anonymous reviewers for their suggestions.

REFERENCES

Autodesk, Inc. 2019. Maya. https://autodesk.com/maya

- Thabo Beeler and Derek Bradley. 2014. Rigid Stabilization of Facial Expressions. ACM Trans. Graph. 33, 4, Article 44 (2014), 44:1–44:9 pages.
- Amit Bermano, Thabo Beeler, Yeara Kozlov, Derek Bradley, Bernd Bickel, and Markus Gross. 2015. Detailed Spatio-Temporal Reconstruction of Eyelids. ACM Trans. Graph. 34, 4, Article 44 (jul 2015).
- Eloïse Berson, Catherine Soladié, Vincent Barrielle, and Nicolas Stoiber. 2019. A Robust Interactive Facial Animation Editing System. *Motion, Interaction and Games* (Oct 2019).
- Eloïse Berson, Catherine Soladié, and Nicolas Stoiber. 2020. Intuitive Facial Animation Editing Based On A Generative RNN Framework. *Computer Graphics Forum* 39, 8 (Nov 2020), 241–251.
- Kiran S. Bhat, Rony Goldenthal, Yuting Ye, Ronald Mallet, and Michael Koperwas. 2013. High Fidelity Facial Animation Capture and Retargeting with Contours. In Proceedings of the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation (Anaheim, California) (SCA '13). Association for Computing Machinery, New York, NY, USA, 7–14.
- Bernd Bickel, Mario Botsch, Roland Angst, Wojciech Matusik, Miguel Otaduy, Hanspeter Pfister, and Markus Gross. 2007. Multi-Scale Capture of Facial Geometry and Motion. In ACM SIGGRAPH 2007 Papers (San Diego, California) (SIGGRAPH '07). Association for Computing Machinery, 33–es.
- Volker Blanz and Thomas Vetter. 1999. A Morphable Model for the Synthesis of 3D Faces. In Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '99). ACM Press/Addison-Wesley Publishing Co., USA, 187–194.
- Stephen Boyd and Lieven Vandenberghe. 2004. *Convex Optimization*. Cambridge University Press.

- Derek Bradley, Wolfgang Heidrich, Tiberiu Popa, and Alla Sheffer. 2010. High Resolution Passive Facial Performance Capture. ACM Trans. Graph. 29, 4, Article 41 (jul 2010).
- D.S. Broomhead and D. Lowe. 1988. Multivariable Functional Interpolation and Adaptive Networks. *Complex Systems* 2 (1988), 321–355.
- Byoungwon Choe, Hanook Lee, and Hyeong seok Ko. 2001. Performance-Driven Muscle-Based Facial Animation. The Journal of Visualization and Computer Animation 12 (2001), 67–79.
- Dimitar Dinev, Thabo Beeler, Derek Bradley, Moritz Bächer, Hongyi Xu, and Ladislav Kavan. 2018. User-Guided Lip Correction for Facial Performance Capture. *Computer Graphics Forum* 37 (2018).
- Bernhard Egger, William A. P. Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhöfer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, Christian Theobalt, Volker Blanz, and Thomas Vetter. 2020. 3D Morphable Face Models - Past, Present, and Future. ACM Trans. Graph. 39, 5 (2020), 157:1–157:38.
- Paul Ekman and Wallace V Friesen. 1978. Facial action coding system. Environmental Psychology & Nonverbal Behavior (1978).
- Luca Fascione. 2017. FACETS Sci-Tech Academy Award. http://oscars.org.
- Yasutaka Furukawa and Jean Ponce. 2009. Dense 3D motion capture for human faces. IEEE Conference on Computer Vision and Pattern Recognition (Jun 2009).
- Graham Fyffe, Andrew Jones, Oleg Alexander, Ryosuke Ichikari, and Paul Debevec. 2015. Driving High-Resolution Facial Scans with Video Performance Capture. ACM Trans. Graph. 34, 1, Article 8 (2015).
- Pablo Garrido, Michael Zollhöfer, Chenglei Wu, Derek Bradley, Patrick Pérez, Thabo Beeler, and Christian Theobalt. 2016. Corrective 3D Reconstruction of Lips from Monocular Video. ACM Trans. Graph. 35, 6, Article 219 (2016).
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 770–778.
- Haoda Huang, Jinxiang Chai, Xin Tong, and Hsiang-Tao Wu. 2011. Leveraging Motion Capture and 3D Scanning for High-Fidelity Facial Performance Acquisition. In *ACM SIGGRAPH 2011 Papers* (Vancouver, British Columbia, Canada) (*SIGGRAPH* '11). Association for Computing Machinery, New York, NY, USA, Article 74.
- Oliver Klehm, Fabrice Rousselle, Marios Papas, Derek Bradley, Christophe Hery, Bernd Bickel, Wojciech Jarosz, and Thabo Beeler. 2015. Recent Advances in Facial Appearance Capture. *Comput. Graph. Forum* 34, 2 (2015), 709–733.

FDLS: A Deep Learning Approach to Production Quality, Controllable, and Retargetable Facial Performances.

- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1 (Lake Tahoe, Nevada) (NIPS'12). Curran Associates Inc., Red Hook, NY, USA, 1097–1105.
- Samuli Laine, Tero Karras, Timo Aila, Antti Herva, Shunsuke Saito, Ronald Yu, Hao Li, and Jaakko Lehtinen. 2017. Production-level facial performance capture using deep convolutional neural networks. In Proceedings of the ACM SIGGRAPH / Eurographics Symposium on Computer Animation, Los Angeles, CA, USA, July 28-30, 2017, Joseph Teran, Changxi Zheng, Stephen N. Spencer, Bernhard Thomaszewski, and KangKang Yin (Eds.). Eurographics Association / ACM, 10:1–10:10.
- Mathieu Lamarre, J.P. Lewis, and Etienne Danvoye. 2018. Face Stabilization by Mode Pursuit for Avatar Construction in the Universe. In 2018 International Conference on Image and Vision Computing, IVCNZ 2018. IEEE, 1–6.
- J. P. Lewis, Ken Anjyo, Taehyun Rhee, Mengjie Zhang, Fred Pighin, and Zhigang Deng. 2014. Practice and Theory of Blendshape Facial Models. In *Eurographics 2014 - State* of the Art Reports, S. Lefebvre and M. Spagnuolo (Eds.).
- Qing Li and Zhigang Deng. 2008. Orthogonal-Blendshape-Based Editing System for Facial Motion Capture Data. IEEE Comput. Graph. Appl. 28, 6 (nov 2008), 76–82.
- Dong C. Liu and Jorge Nocedal. 1989. On the limited memory BFGS method for large scale optimization. *Mathematical Programming* 45, 1 (01 Aug 1989), 503–528.
- Wan-Chun Ma, Mathieu Lamarre, Etienne Danvoye, Chongyang Ma, Manny Ko, Javier von der Pahlen, and Cyrus A. Wilson. 2016. Semantically-aware blendshape rigs from facial performance measurements. In SIGGRAPH ASIA 2016, Macao, December 5-8, 2016 - Technical Briefs, Johannes Kopf and Phillip Chi-Wing Fu (Eds.). ACM, 3.
- Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In ECCV.
- Lucio Moser, Mark Williams, Darren Hendler, and Doug Roble. 2018. High-Quality, Cost-Effective Facial Motion Capture Pipeline with 3D Regression. In ACM SIG-GRAPH 2018 Talks (Vancouver, British Columbia, Canada) (SIGGRAPH '18). 2 pages.
- Jason Osipa. 2010. Stop Staring: Facial Modeling and Animation Done Right, 3rd Ed. Sybex.
- Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. 2017. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA. 5099–5108.
- Jaewoo Seo, Geoffrey Irving, J.P. Lewis, and Junyong Noh. 2011. Compression and direct manipulation of complex blendshape models. ACM Trans. Graph. 30, 6,

- Article 164 (Dec. 2011), 164:1-164:10 pages.
- Yeongho Seol, J.P. Lewis, Jaewoo Seo, Byungkuk Choi, Ken Anjyo, and Junyong Noh. 2012. Spacetime expression cloning for blendshapes. ACM Trans. Graph. 31, 2, Article 14 (April 2012), 14:1–14:12 pages.
- Yeongho Seol, Wan-Chun Ma, and J. P. Lewis. 2016. Creating an Actor-Specific Facial Rig from Performance Capture. In Proceedings of the 2016 Symposium on Digital Production (Anaheim, California) (DigiPro '16). Association for Computing Machinery, New York, NY, USA, 13–17. https://doi.org/10.1145/2947688.2947693
- Yeongho Seol, Jaewoo Seo, Paul Hyunjin Kim, John P. Lewis, and Junyong Noh. 2011. Artist friendly facial animation retargeting. ACM Trans. Graph. 30, 6 (2011), 162.
- Alex Smith, Sven Pohle, Wan-Chun Ma, Chongyang Ma, Xian-Chun Wu, Yanbing Chen, Etienne Danvoye, Jorge Jimenez, Sanjit Patel, Mike Sanders, and Cyrus A. Wilson. 2017. Emotion challenge: building a new photoreal facial performance pipeline for games. In Proc. ACM SIGGRAPH Digital Production Symposium, Christopher Horvath, Cary B. Phillips, Andrew Pearce, Corban Gossett, and Stephen N. Spencer (Eds.). ACM, 8:1–8:2.
- Andrea Tagliasacchi, Sofien Bouaziz, Mark Pauly, and Hao Li. 2016. Modern techniques and applications for real-time non-rigid registration. In SIGGRAPH ASIA 2016, Macao, December 5-8, 2016 - Courses. 11:1–11:25.

Vladimir Vapnik. 1998. Statistical learning theory. Wiley.

- Various 2022. Commercial performance capture systems. https://di4d.com, https: //www.dynamixyz.com, https://www.synthesia.io, https://studios.disneyresearch. com/medusa, https://image-metrics.com.
- Lance Williams. 1990. Performance-Driven Facial Animation. SIGGRAPH Comput. Graph. 24, 4 (sep 1990), 235–242.
- Wenwu Yang, Nathan Marshak, Daniel Sýkora, Srikumar Ramalingam, and Ladislav Kavan. 2019. Building anatomically realistic jaw kinematics model from data. *The Visual Computer* 35, 6-8 (2019), 1105–1118.
- Li Zhang, Noah Snavely, Brian Curless, and Steven M. Seitz. 2004. Spacetime Faces: High Resolution Capture for Modeling and Animation. *ACM Trans. Graph.* 23, 3 (aug 2004), 548–558.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*.
- M. Zollhöfer, J. Thies, P. Garrido, D. Bradley, T. Beeler, P. Pérez, M. Stamminger, M. Nießner, and C. Theobalt. 2018. State of the Art on Monocular 3D Face Reconstruction, Tracking, and Applications. *Computer Graphics Forum* 37, 2 (2018), 523–550.