# 3D Shape Regression for Real-time Facial Animation

Chen Cao, Yanlin Weng, Stephen Lin, Kun Zhou

# FaceWarehouse: a 3D facial Expression Database for Visual Computing

Chen Cao, Yanlin Weng, Shun Zhou, Yiying Tong, Kun Zhou

Zhejiang University     MSRA

Presented by Shu Liang

(Black-on-white slides are Shu's)

# Facial Animation

- Facial animation is widely used in films & games
- Performance-based facial animation



Avatar 2009
© 21st Century Fox

L.A. Noire 2011
© Team Bondi

# Related Work

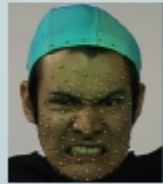- ## Performance-based Facial Animation

Quality

**Special Equipment**

Facial Markers       Camera arrays       Structured light
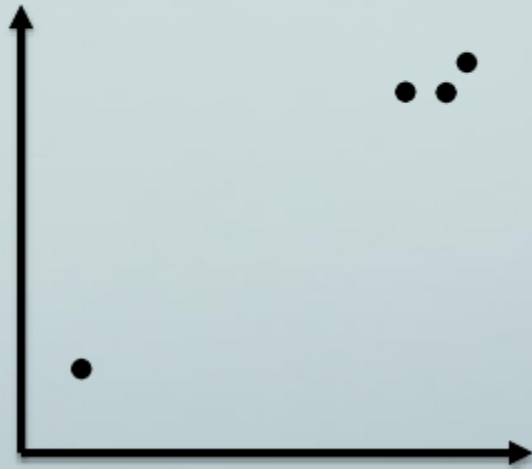


[Huang et al. 2011]       [Beeler et al. 2011]       [Weise et al. 2009]
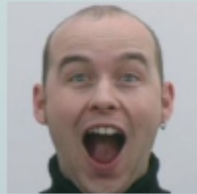
Device complexity

# Related Work

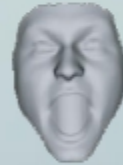- ## Performance-based Facial Animation



Quality

**Single Camera**

Optical Flow

ASM & AAM

Regression-based alignment

[Vlasic et al. 2005]

[Cootes et al. 1992-2001]

[Cao et al. 2012]

Device complexity
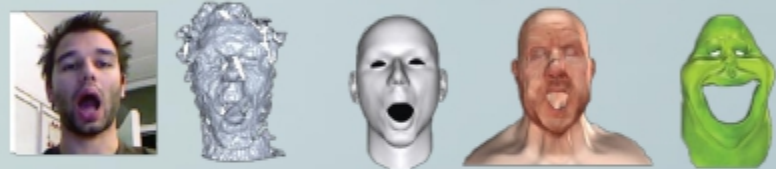
# Related Work

- ## Performance-based Facial Animation

Quality

Device complexity

**Consumer RGBD Camera**
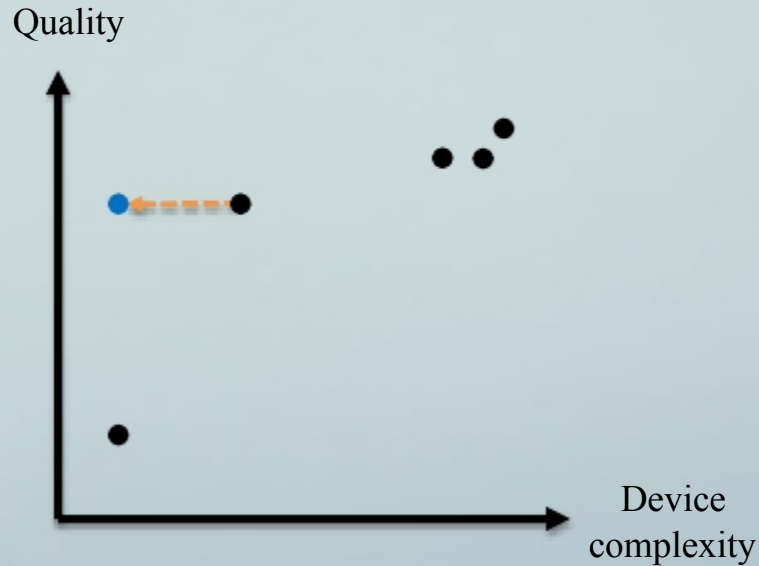


[Weise et al. 2011]



[Bouaziz et al. 2013]

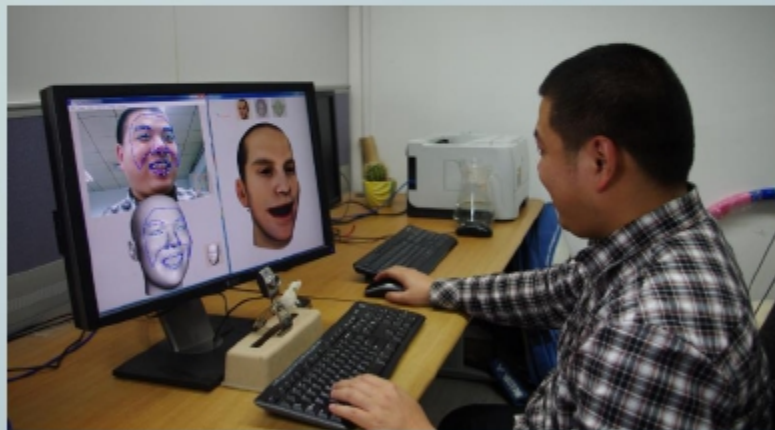

[Li et al. 2013]

# Our Goal

- Real-time facial animation for average users
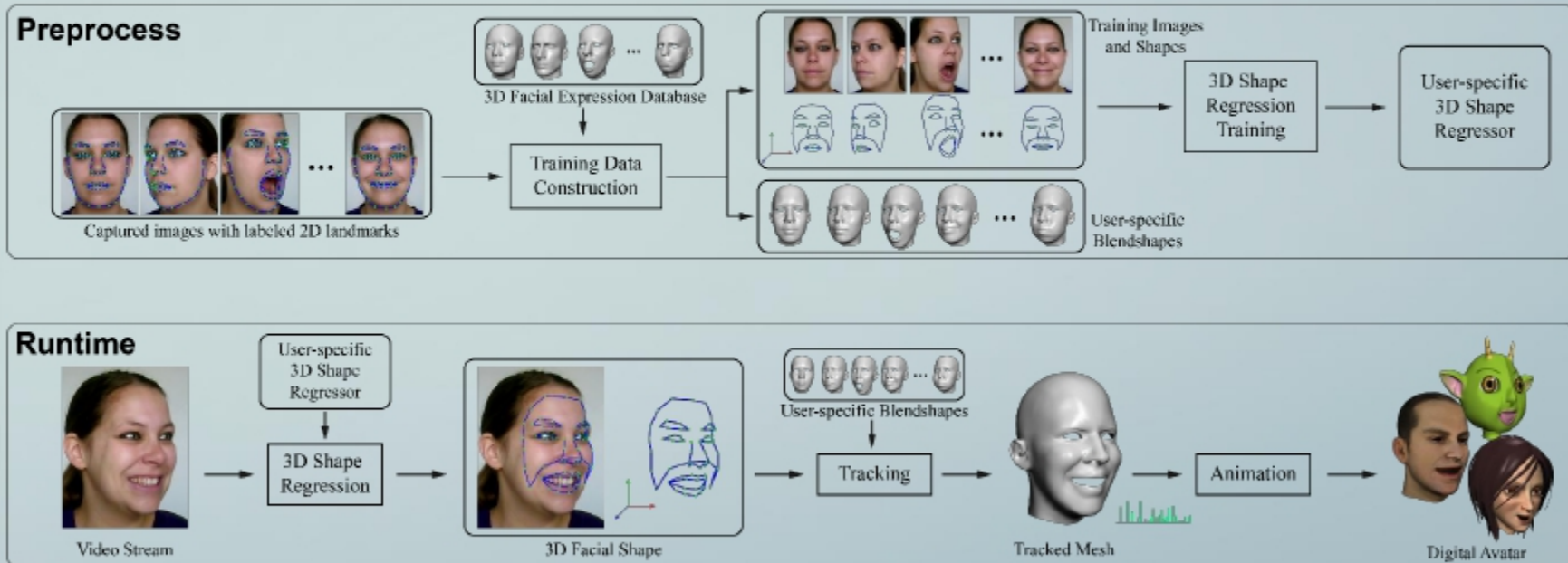
# Our Goal

- Real-time facial animation for ordinary users
  - Single web camera
  - Robust
    - Fast motions
    - Large rotations
    - Exaggerated expressions
  - General environments
    - Indoors and outdoors
  - High performance
    - Mobile devices

# Our Pipeline



**Preprocess**

Captured images with labeled 2D landmarks

3D Facial Expression Database

Training Data Construction

Training Images and Shapes

3D Shape Regression Training

User-specific 3D Shape Regressor

User-specific Blendshapes

**Runtime**

Video Stream

User-specific 3D Shape Regressor

3D Shape Regression

3D Facial Shape

User-specific Blendshapes

Tracking

Tracked Mesh

Animation

Digital Avatar

# Our Pipeline

- One-time Preprocess



3D Facial Expression Database

Training Data Construction

Captured images with labeled 2D landmarks

Training Images and Shapes

User-specific Blendshapes

3D Shape Regression Training

User-specific 3D Shape Regressor
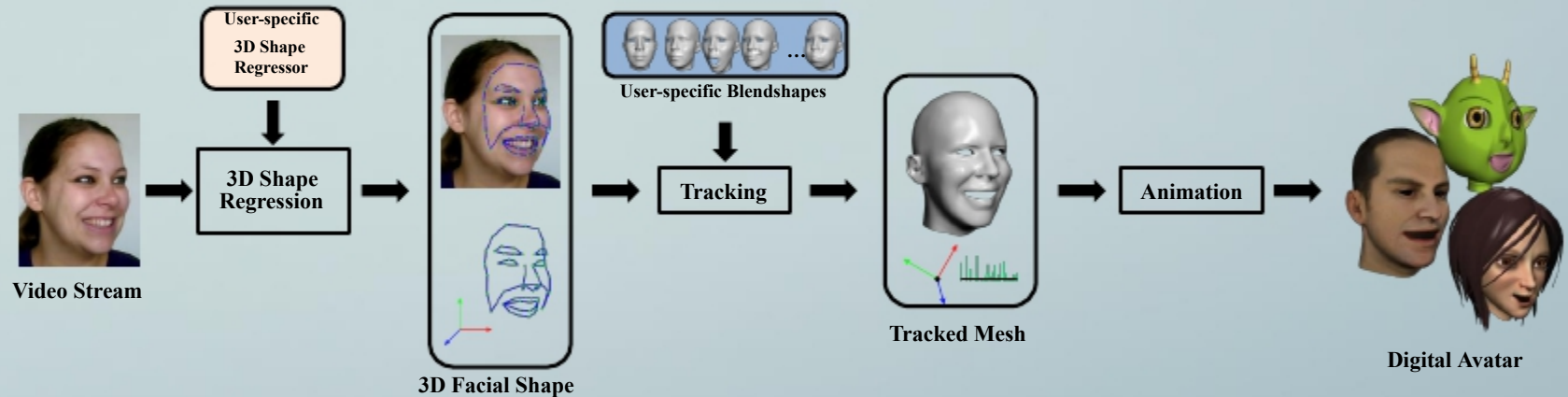
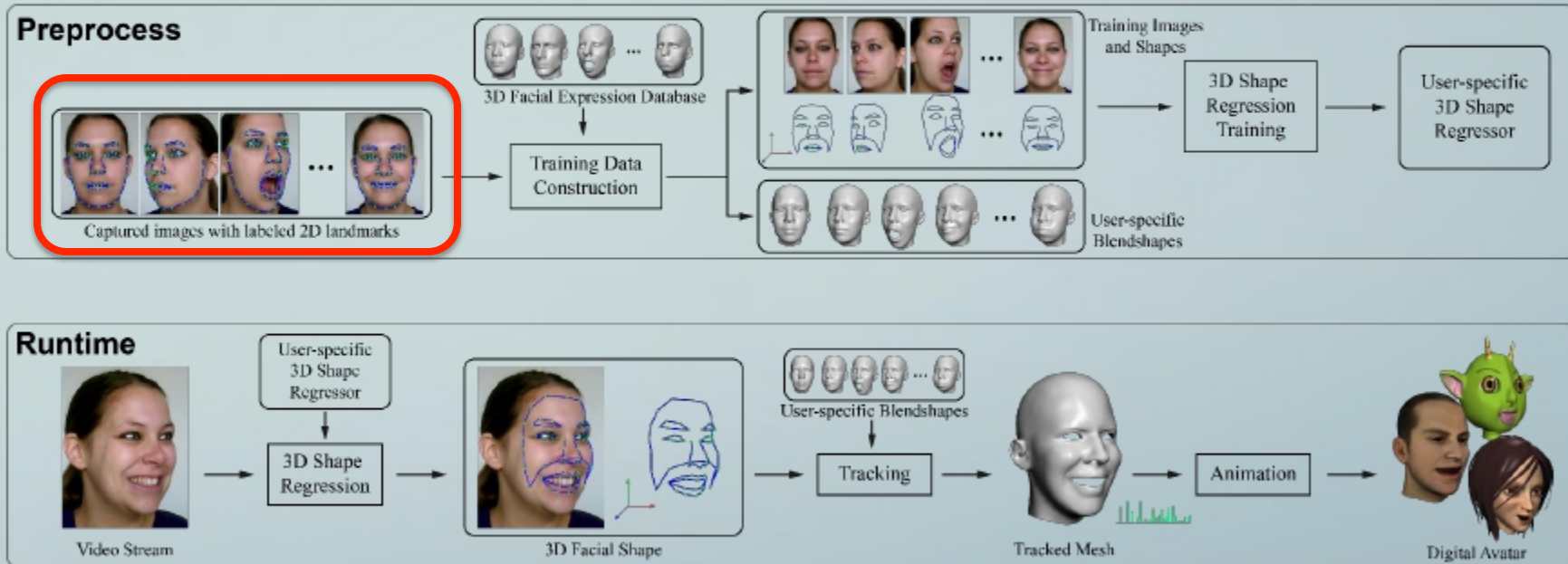# Our Pipeline

- Runtime computation

# 3D Face Shape Regression: Preprocess

- Data Collection
  - Image capturing & labeling
  - Blendshapes generation
  - Shape reconstruction
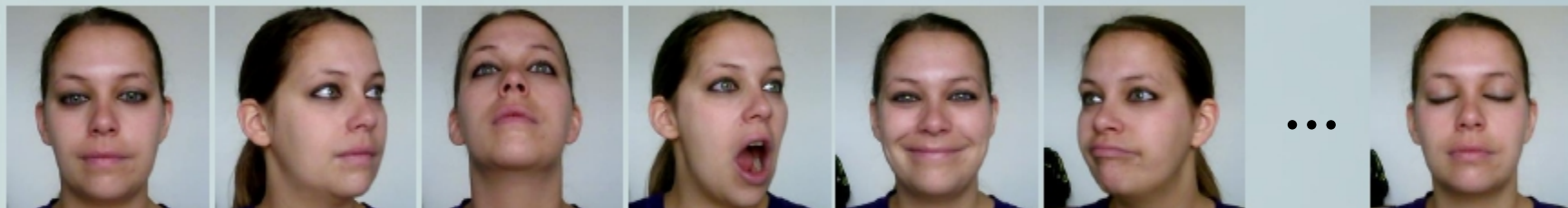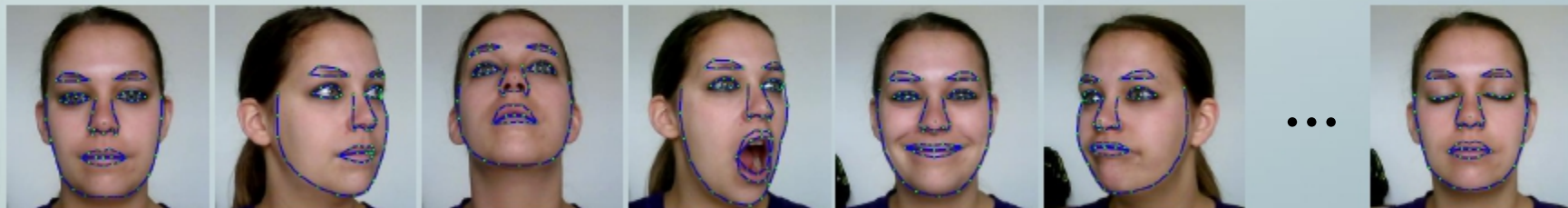  - Training data generation
- Training

# Our Pipeline

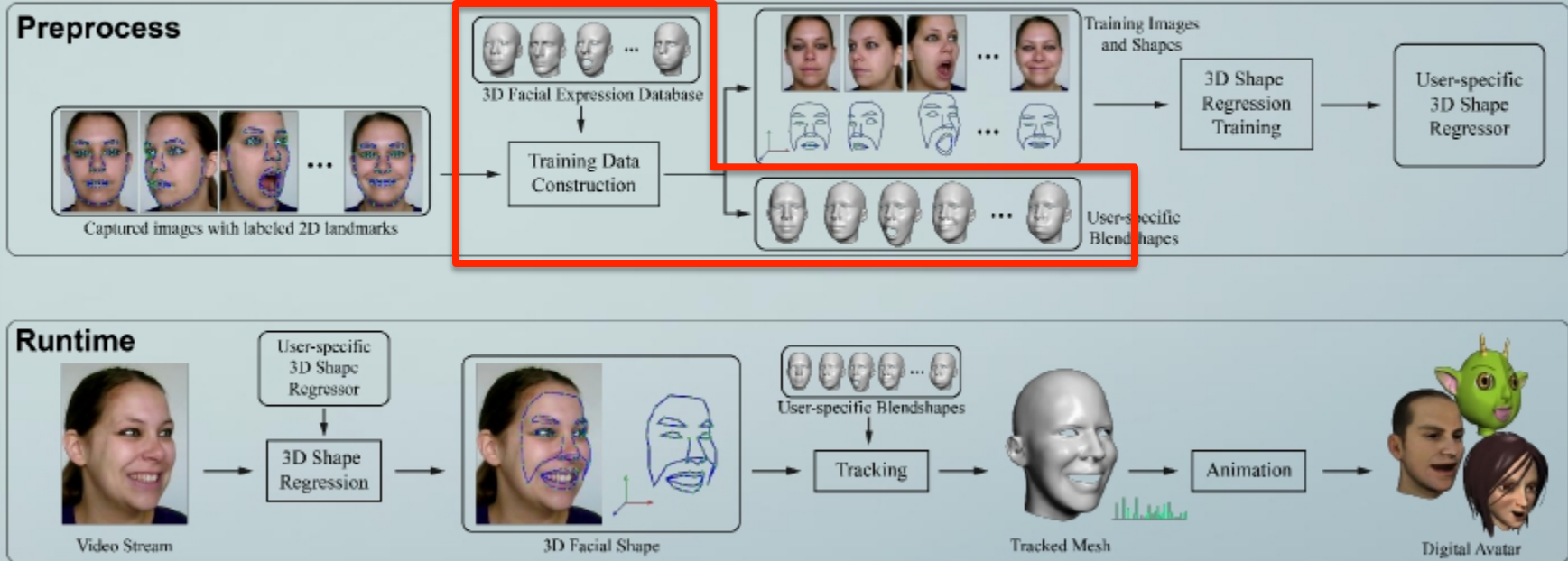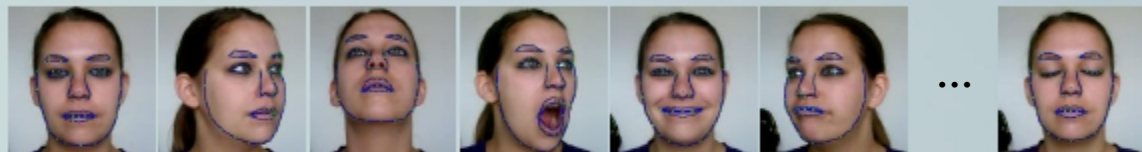# Preprocess: Image Capturing & Labeling



Captured Images

**[Cao et al. 2012] + Manual Adjustment**
Labeled 2D Feature Points

# Our Pipeline



**Preprocess**

Captured images with labeled 2D landmarks

3D Facial Expression Database

Training Data Construction

Training Images and Shapes

3D Shape Regression Training

User-specific 3D Shape Regressor

User-specific Blendshapes

**Runtime**

Video Stream

User-specific 3D Shape Regressor

3D Shape Regression

3D Facial Shape

User-specific Blendshapes

Tracking

Tracked Mesh

Animation

Digital Avatar

# Preprocess: Blendshapes Generation



Labeled 2D Feature Points

FaceWarehouse
[Cao et al. 2013]

**150** identities × **47** expressions

Fitting

User-specific Blendshapes

# Our Pipeline



FaceWarehouse

**Preprocess**

Captured images with labeled 2D landmarks

3D Facial Expression Database

Training Data Construction

Training Images and Shapes

3D Shape Regression Training

User-specific 3D Shape Regressor

User-specific Blendshapes

**Runtime**

Video Stream

User-specific 3D Shape Regressor

3D Shape Regression

3D Facial Shape

User-specific Blendshapes

Tracking

Tracked Mesh

Animation

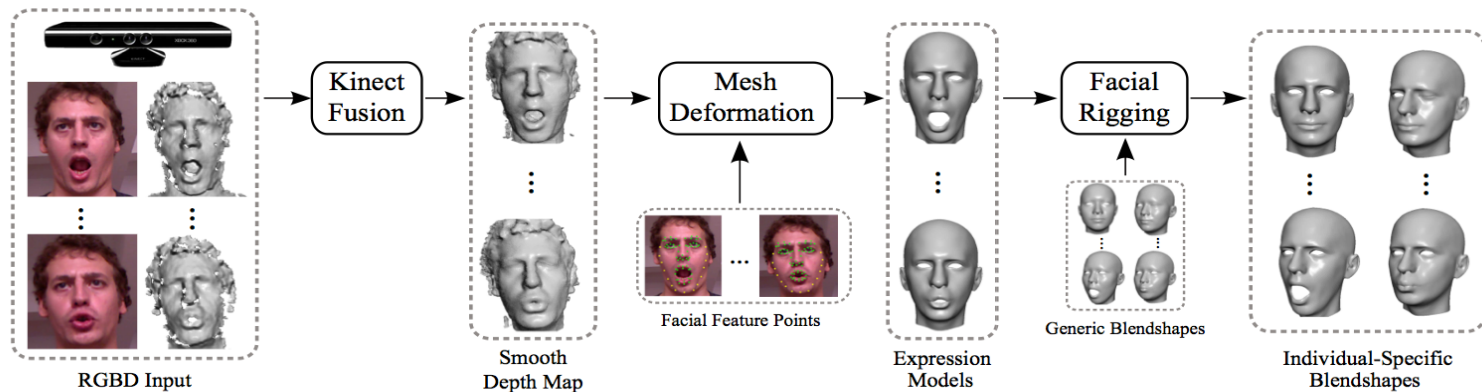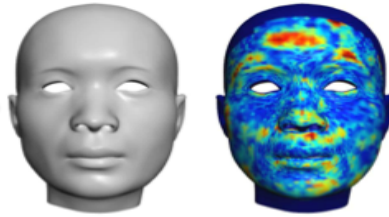Digital Avatar

# FaceWarehouse

- RGBD images of 150 individuals captured by Kinect
- Aged 7-80 from various ethnic backgrounds
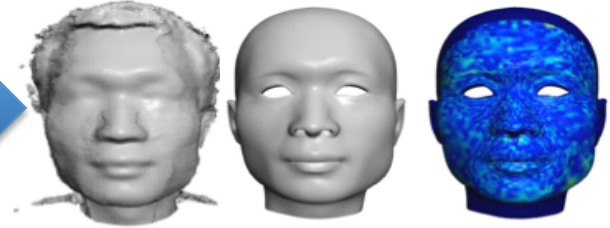- Different expressions, one neutral and 19 other expressions.



RGBD Input → Kinect Fusion → Smooth Depth Map → Mesh Deformation (Facial Feature Points) → Expression Models → Facial Rigging (Generic Blendshapes) → Individual-Specific Blendshapes

# FaceWarehouse

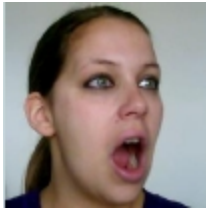- Mesh deformation

Neutral expression



Other expressions

[Blanz et. al 2004]        [Huang et. al 2006]

[Sumner et. al 2006]        [Huang et. al 2006]

$S_0, S_1, S_2 \ldots S_{19}$ for 20 expressions.

# FaceWarehouse

− Individual-specific expression blendshapes

- Example-based facial rigging algorithm:

An expression H of the person can be：

$$H = B_0 + \sum_{i=1}^{46} \alpha_i (B_i - B_0)$$

{$B_1$,$B_2$,...$B_{46}$} 46 FACS blendshapes

- Begins with a generic blendshape model $\mathbf{A} = \{A_0, A_1, ..., A_{46}\}$
- Optimized by minimizing the difference between $S_j$ and linear combination of $B_i$ with known weight for expression j, the difference between the deformation from $B_0$ to $B_i$ and that from $A_0$ to $A_i$.
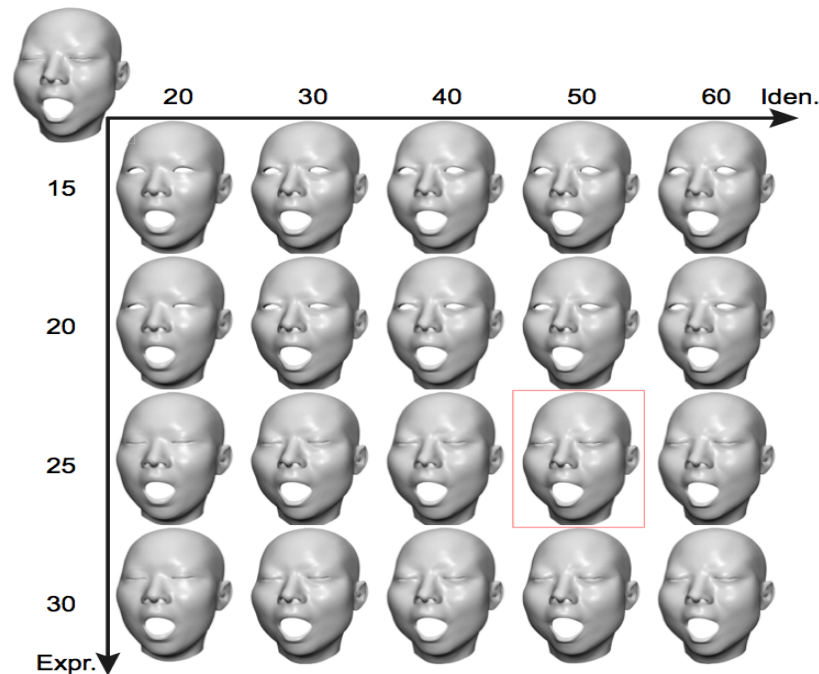
# FaceWarehouse

- ## Bilinear face model

A rank-three data tensor T.
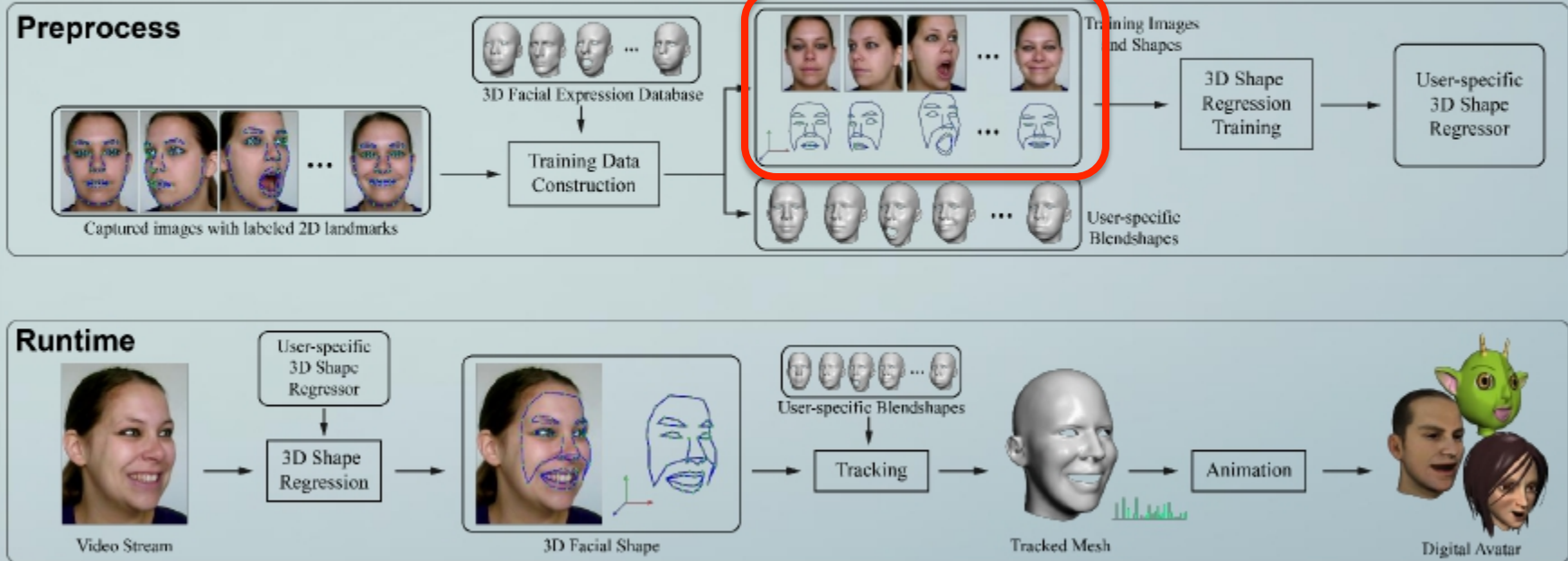(11K vertices × 150 identities ×
47 expressions）
Used N-mode SVD to
decompose the tensor.

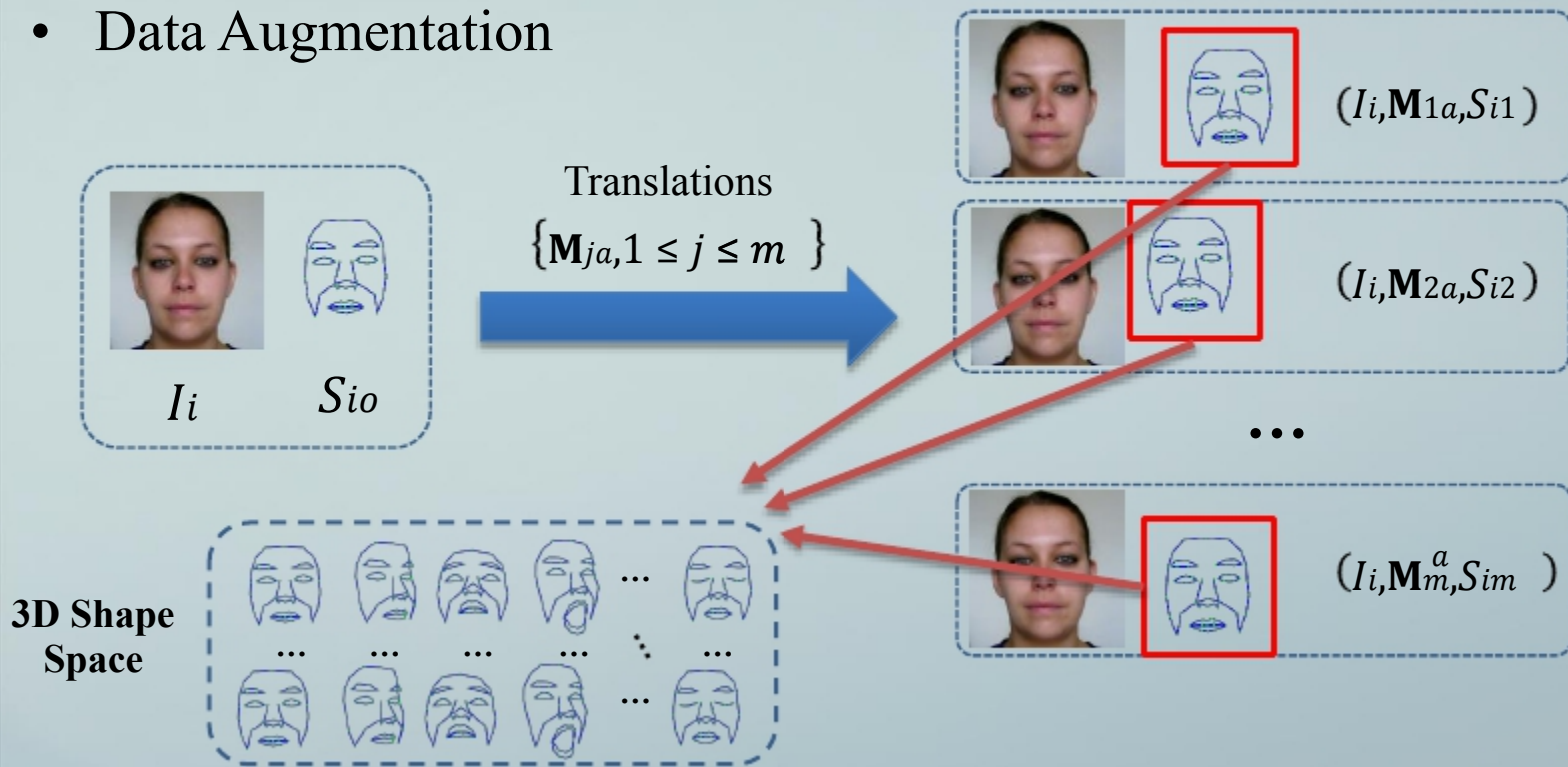$$V = C_r \times_2 \mathbf{w}_{id}^T \times_3 \mathbf{w}_{exp}^T,$$

# Our Pipeline

# Preprocess: 3D Shape Reconstruction



Image
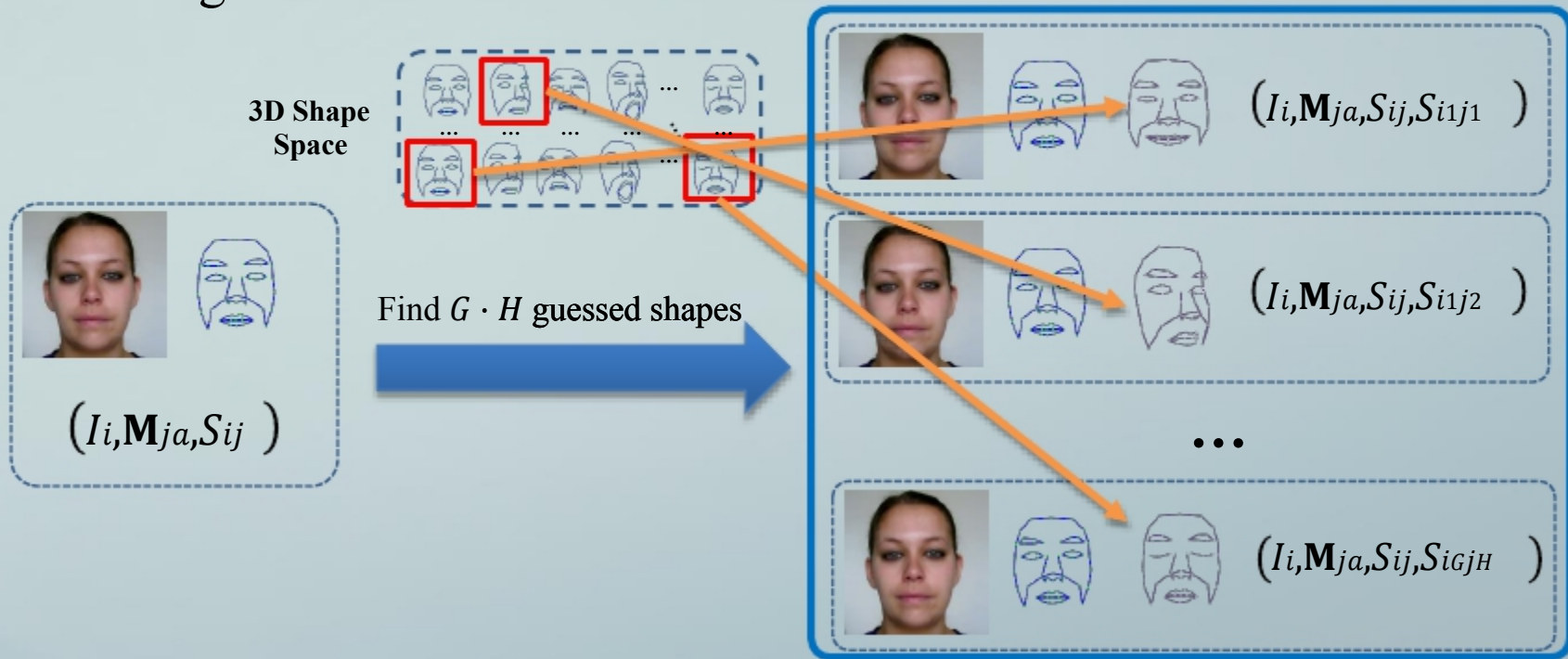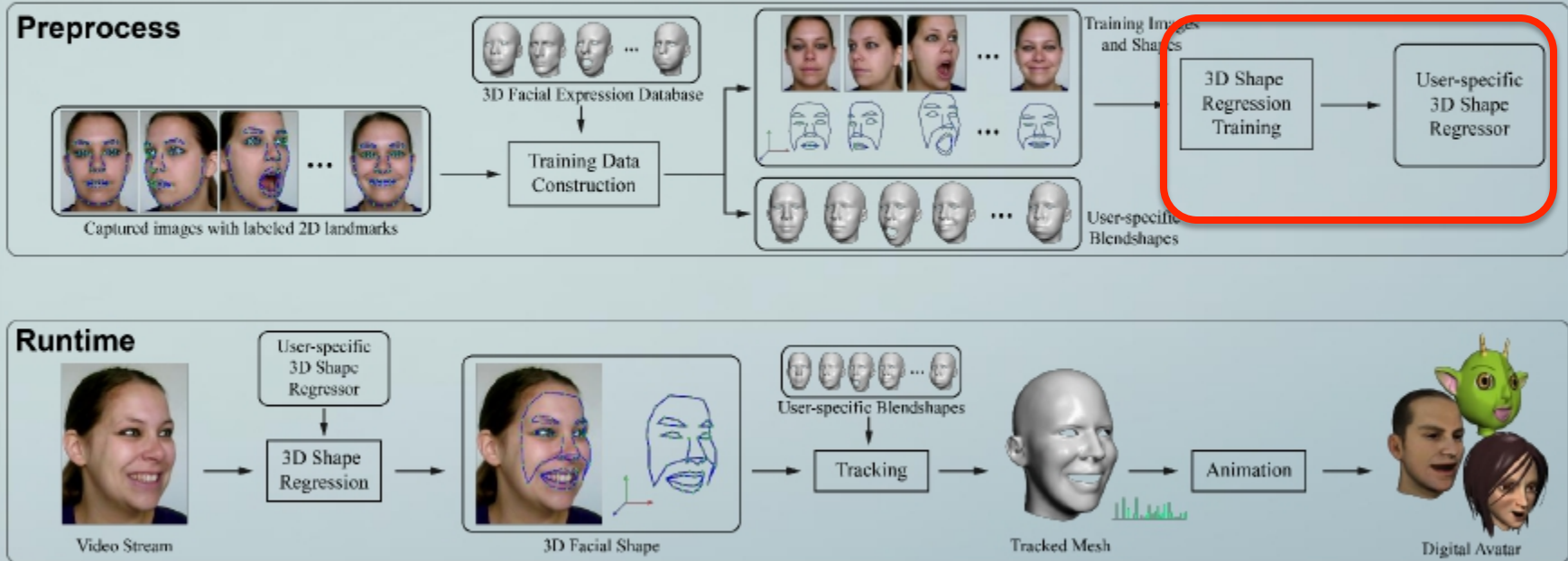
Mesh

Shape

# **Preprocess: Training Data Generation**

- Data Augmentation



$I_i$  $S_{io}$

Translations
$\{\mathbf{M}_{ja}, 1 \leq j \leq m\}$

$(I_i, \mathbf{M}_{1a}, S_{i1})$

$(I_i, \mathbf{M}_{2a}, S_{i2})$

$\cdots$

$(I_i, \mathbf{M}_m^a, S_{im})$

**3D Shape Space**

# Preprocess: Training Data Generation

- Training Set Construction

Training Data



3D Shape Space

Find $G \cdot H$ guessed shapes

$(I_i, \mathbf{M}_{ja}, S_{ij})$

$(I_i, \mathbf{M}_{ja}, S_{ij}, S_{i1j1})$

$(I_i, \mathbf{M}_{ja}, S_{ij}, S_{i1j2})$

$\cdots$
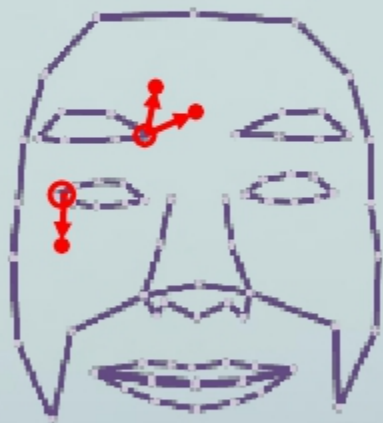
$(I_i, \mathbf{M}_{ja}, S_{ij}, S_{iGjH})$

# Our Pipeline

# Preprocess: Training

- Appearance vector
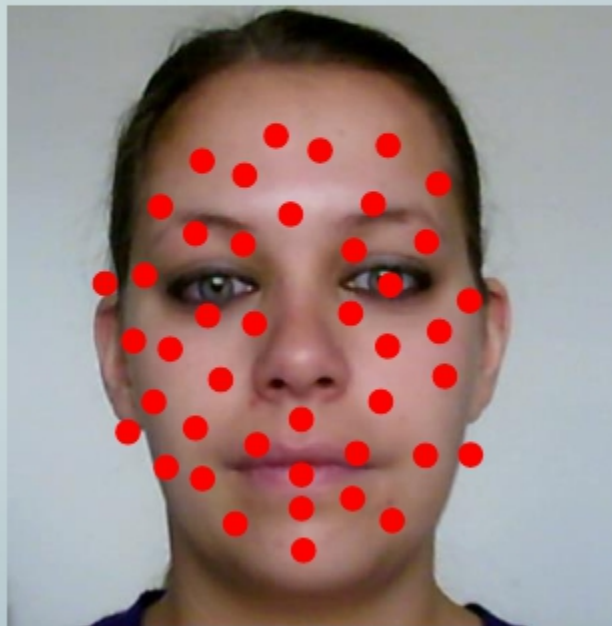- Primitive regressor: fern

- Summary: two-level boosted regressor

# Preprocess: Training

- Appearance vector



$S_{ic}$

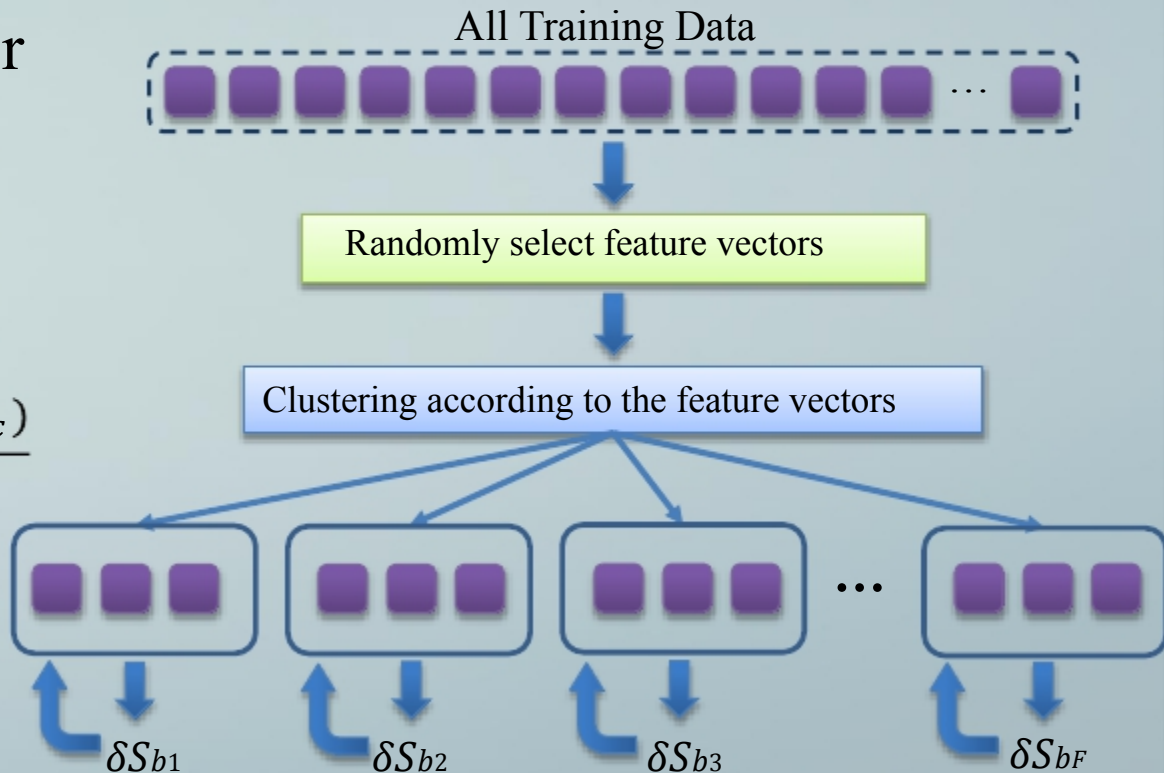$(I_i, \mathbf{M}_{ia})$

Appearance Vector

# Preprocess: Training

- Primitive regressor



$$\delta S_b = \frac{1}{1 + \beta / |\Omega_b|} \frac{\sum_{i \in \Omega_b} (S_i - S_{ic})}{|\Omega_b|}$$

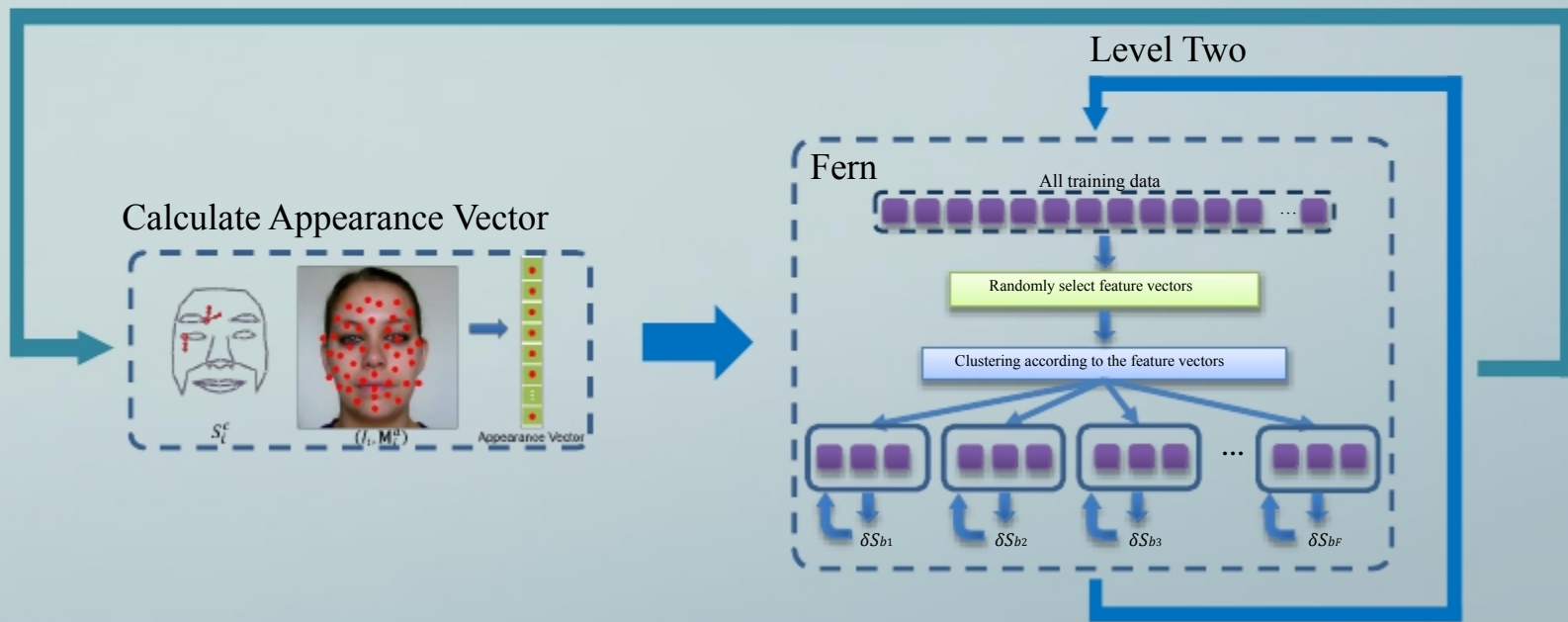$$S_{ic} = S_{ic} + \delta S_b, i \in \Omega_b$$

All Training Data

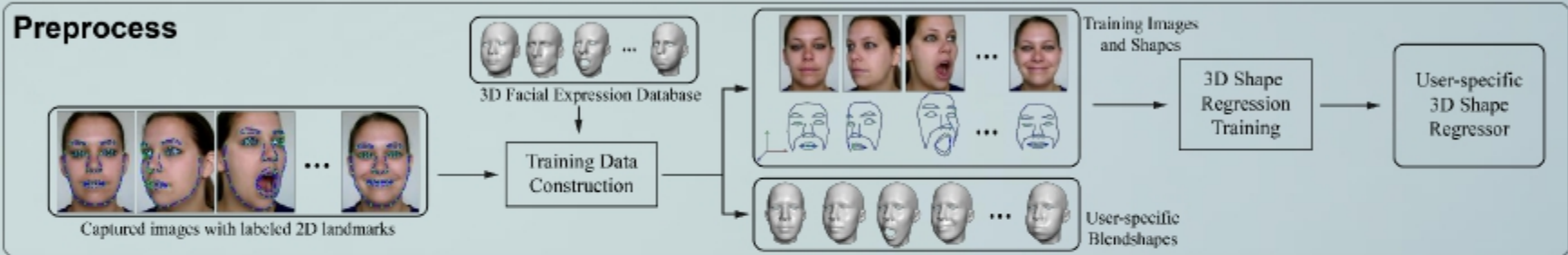Randomly select feature vectors

Clustering according to the feature vectors

$\delta S_{b1}$   $\delta S_{b2}$   $\delta S_{b3}$   $\cdots$   $\delta S_{bF}$

# Our Pipeline



**Preprocess**

Captured images with labeled 2D landmarks

3D Facial Expression Database

Training Data Construction

Training Images and Shapes

3D Shape Regression Training

User-specific 3D Shape Regressor

User-specific Blendshapes

**Runtime**

Video Stream

User-specific 3D Shape Regressor

3D Shape Regression

3D Facial Shape

User-specific Blendshapes

Tracking

Tracked Mesh

Animation

Digital Avatar

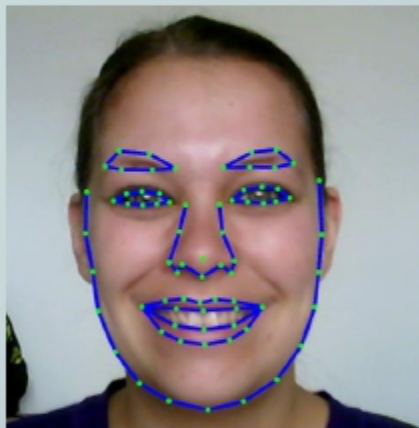# 3D Face Shape Regression: Runtime

- Initialization: first frame
- Following frames
  - Find guessed shapes
  - Two-level boosted regression
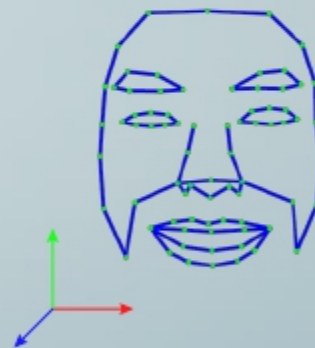
# Runtime: Initialization

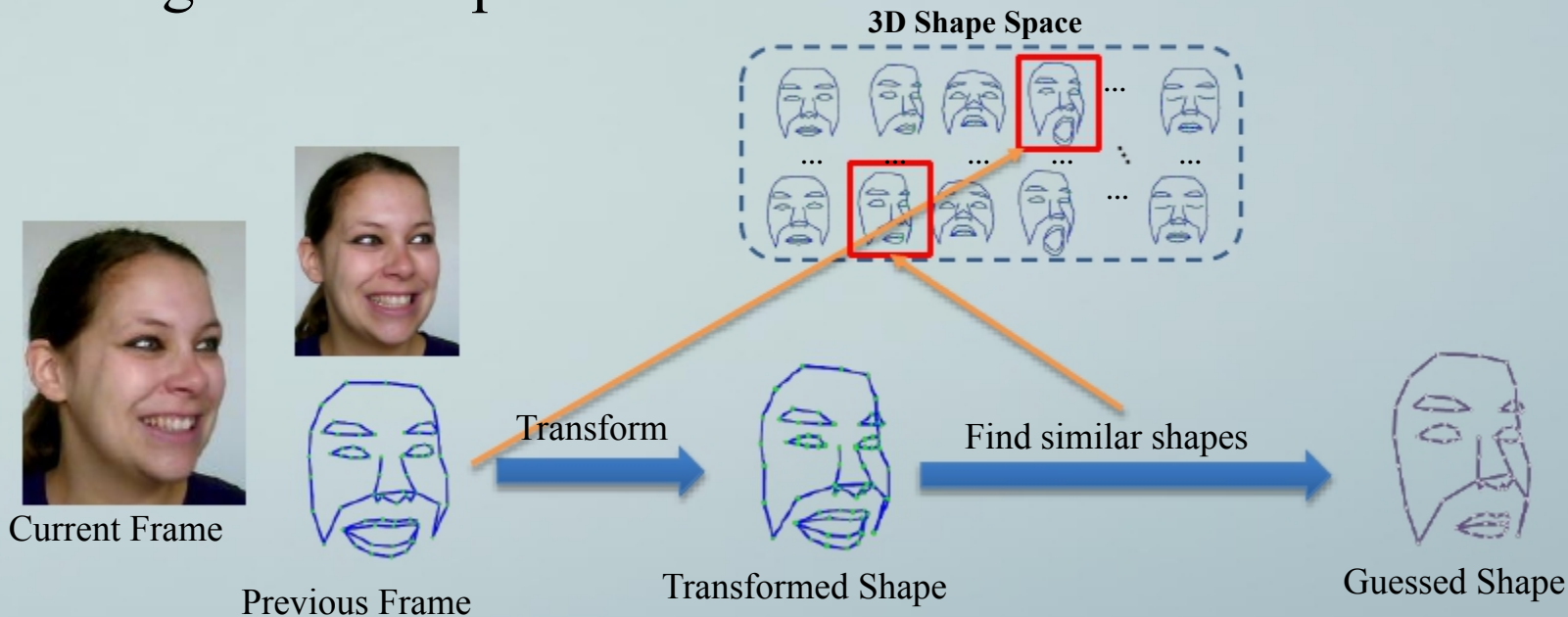- First frame



Face Detection
[Viola and Jones 2001]
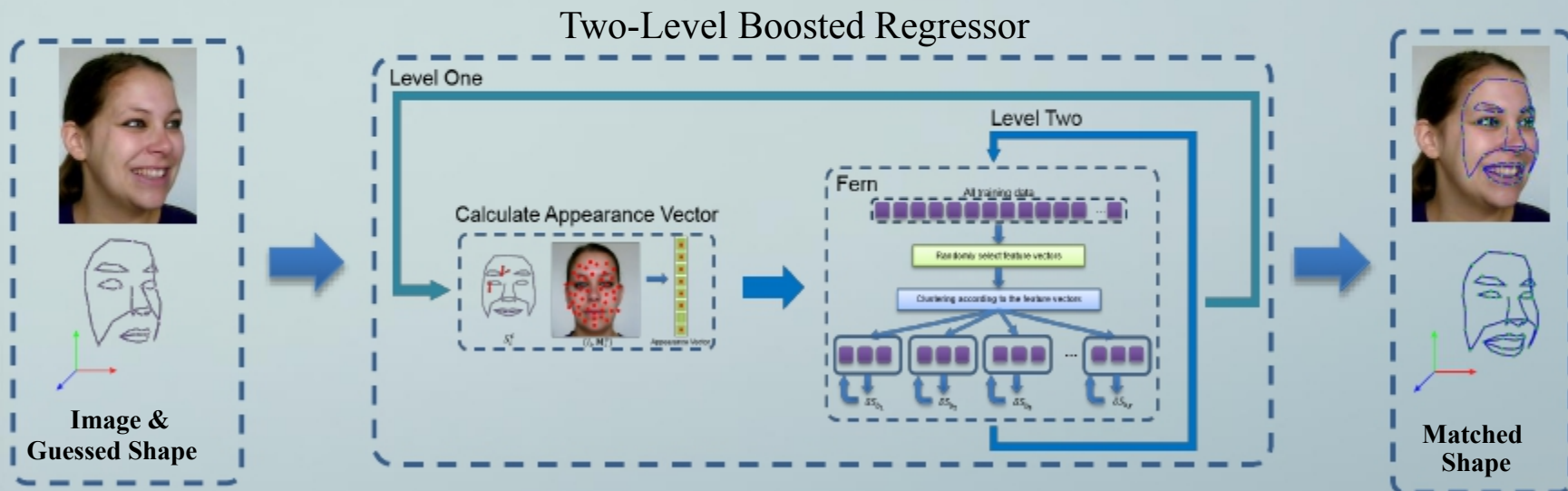
2D Feature Alignment
[Cao et al. 2012]

3D Shape Recovery
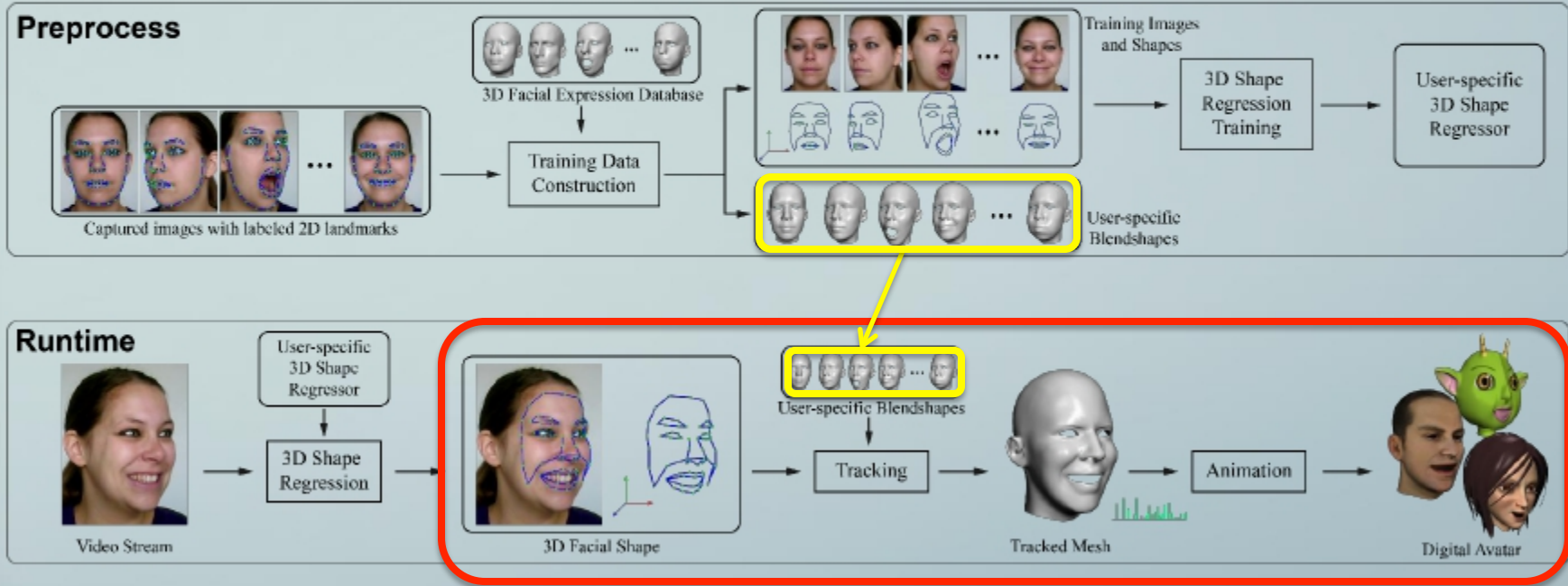
# Runtime: Following Frames

- Find guessed shape



**3D Shape Space**

Current Frame

Previous Frame

Transform

Transformed Shape

Find similar shapes

Guessed Shape

# Runtime: Following Frames

- Two-level boosted regression



Two-Level Boosted Regressor

Level One

Level Two

Calculate Appearance Vector

Fern

All training data

Randomly select feature vectors

Clustering according to the feature vectors

Image & Guessed Shape

Matched Shape

# Our Pipeline

# Tracking & Animation

Similar to [Weise et al. 2011]

User-specific Blendshapes

Matched Shape

Tracking

Rigid

Non-Rigid

Tracked Mesh

Rigid Transformation

Blendshape Weight

Digital Avatar

# Evaluation: Regressed shape vs. Kinect 3D vs. 2D vs. Optical Flow



**Figure 8:** *Comparison of depth from 3D shape regression and ground truth from Kinect.*
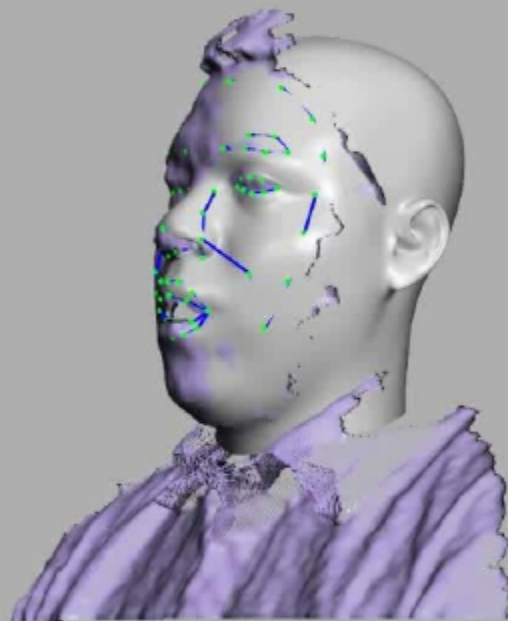
| RMSE | < 3 pixels | < 4.5 pixels | < 6 pixels |
|---|---|---|---|
| 3D Regression | 73.3% | 80.8% | 100% |
| 2D Regression | 50.8% | 64.2% | 72.5% |
| Optical Flow | 20.8% | 24.2% | 41.7% |

**Table 1:** *Percentages of frames with RMSE less than given thresholds for the tested video sequence.*
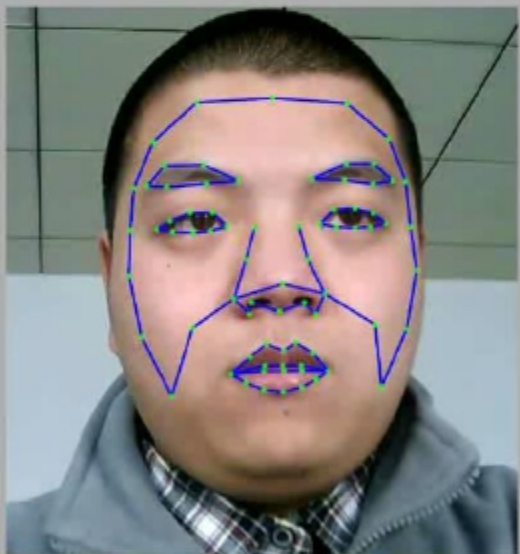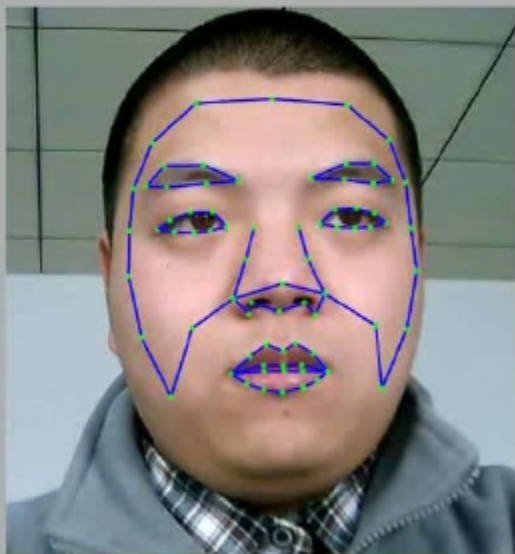
# Live Demo
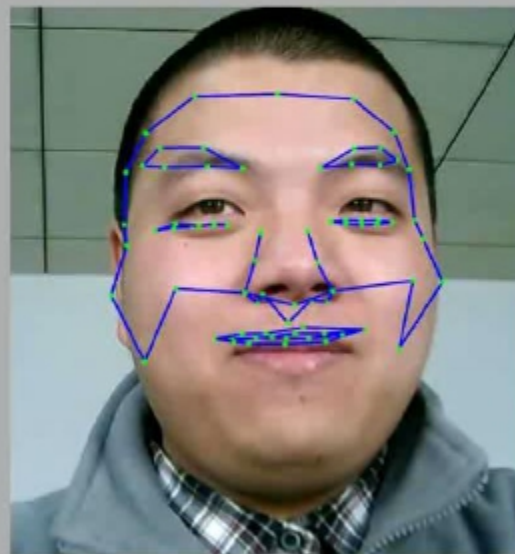
- [Demo](#)

# Evaluation: Regressed shape vs. Kinect

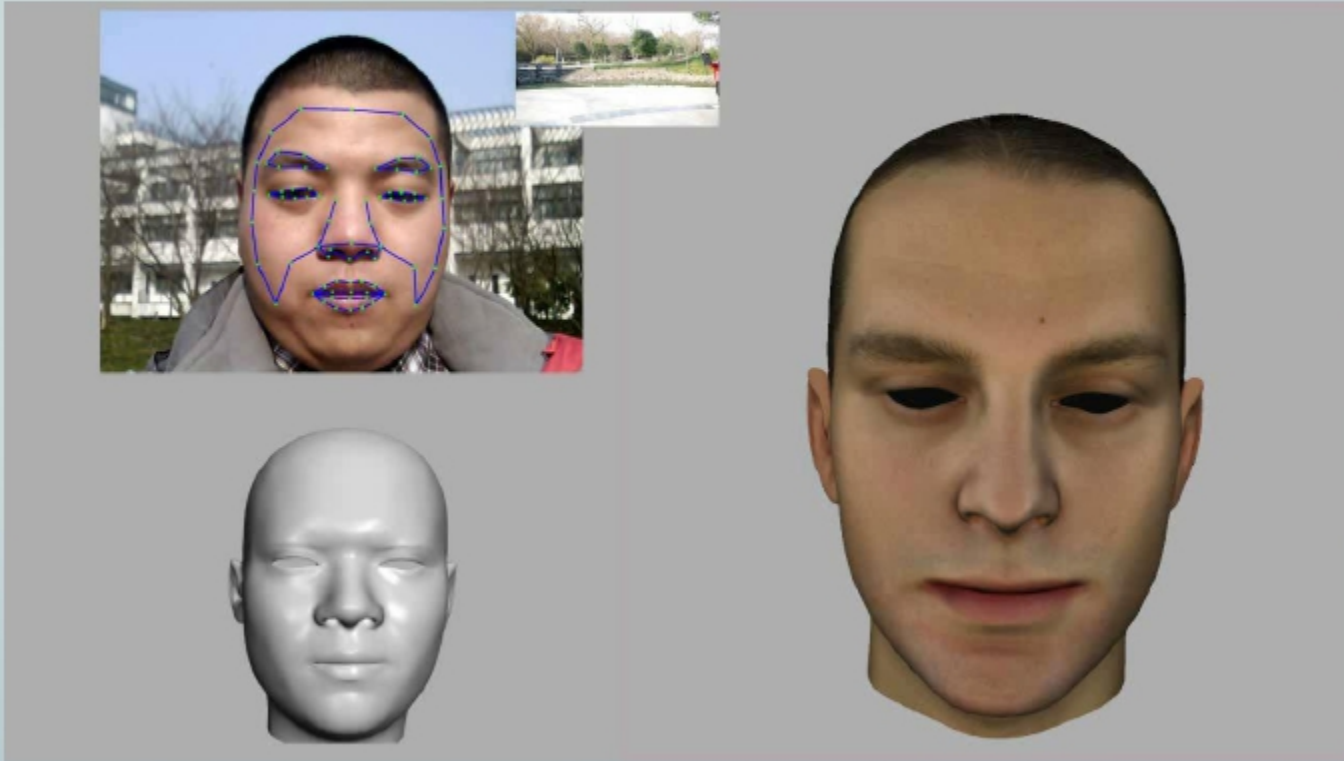# Evaluation: 3D vs. 2D vs. Optical Flow



Our 3D Regression  2D Regression  Optical Flow Based

# More Results: Outdoor
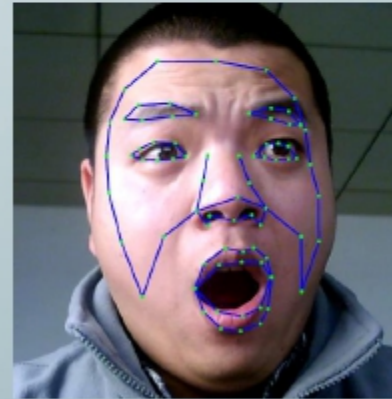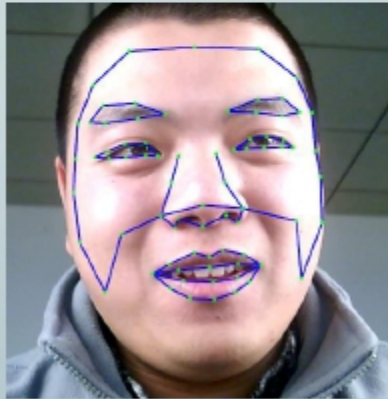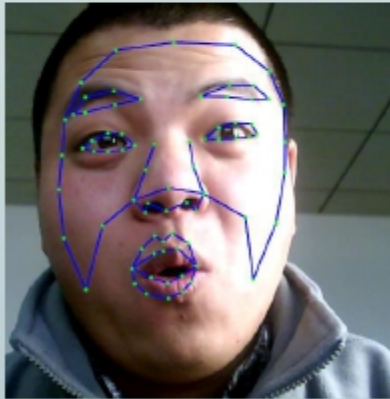
# Our System on Mobile Device

# Timings

- Preprocess: 45 mins
  - Capture: 10 mins
  - User interaction: 25 mins
  - Training: 10 mins
- Runtime: less than 15 ms
  - Regression: 5 ms
  - Tracking & Animation: 8 ms

# Limitations

- Much <span style="color:red">training data</span>
  - **60** head poses and facial expressions
- Dramatic <span style="color:red">lighting changes</span>

# Summary

- 3D facial performance capture from 2D video
  - Regression-based approach
  - Robust: fast motions, large rotations, exaggerated expressions
  - General environments: indoors and outdoors
  - High performance: real-time

- Future work
  - Handle lighting variations
  - Reduce training data

SIGGRAPH 2013

# Acknowledgement

- Face capture: Marion Blatt, Steffen Toborg
- Anonymous reviewers
- Funding
  - NSFC (No.61003145 and No.61272305)
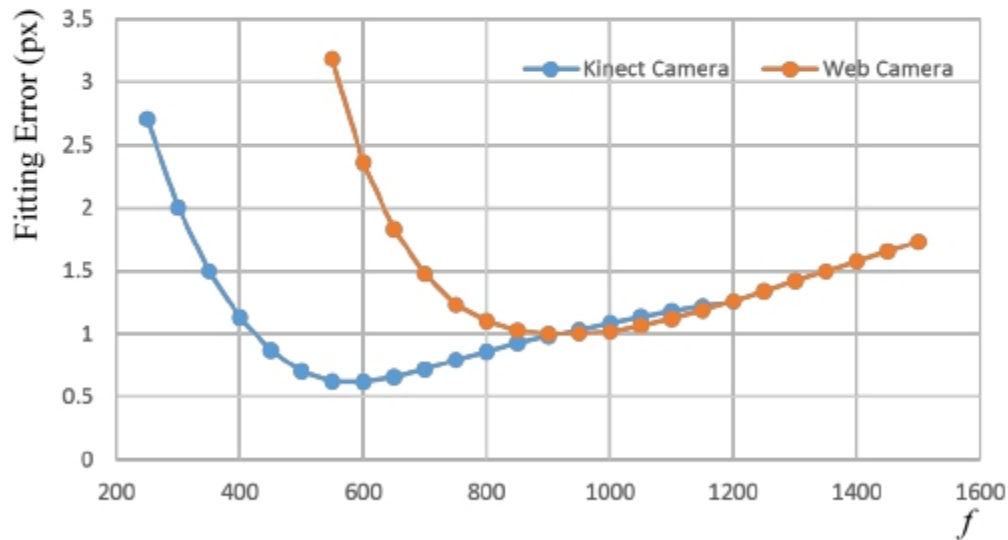  - 973 program of China (No.2009CV320801)
- FaceWarehouse Data: http://gaps-zju.org/facewarehouse/

# Thank you!

# Preprocess: Camera Calibration

Blendshape Generations:

$$E = \sum_{i=1}^{n} \sum_{k=1}^{75} \left\| \boxed{\Pi \mathbf{Q}} \left( \mathbf{M}_i \left( C_r \times_2 \mathbf{w}_{id}^T \times_3 \mathbf{w}_{exp,i}^T \right)^{(v_k)} - \mathbf{u}_i^{(k)} \right) \right\|^2$$

$$\mathbf{Q} = \begin{pmatrix} f & 0 & u_0 \\ 0 & f & v_0 \\ 0 & 0 & 1 \end{pmatrix}$$

# Why not directly use previous shape?

- Error accumulation

# Why not directly regress parameters?

- Expression coefficients in [0:1]
- Animation prior
  - Temporal coherence
- Rigid transformation & expression coefficients
  - Different spaces