# Face Landmarks

**and stuff you can do with them**

CS448V — Computational Video Manipulation

May 2019

Constrained Local Models

Improving Portraits

Editing Video

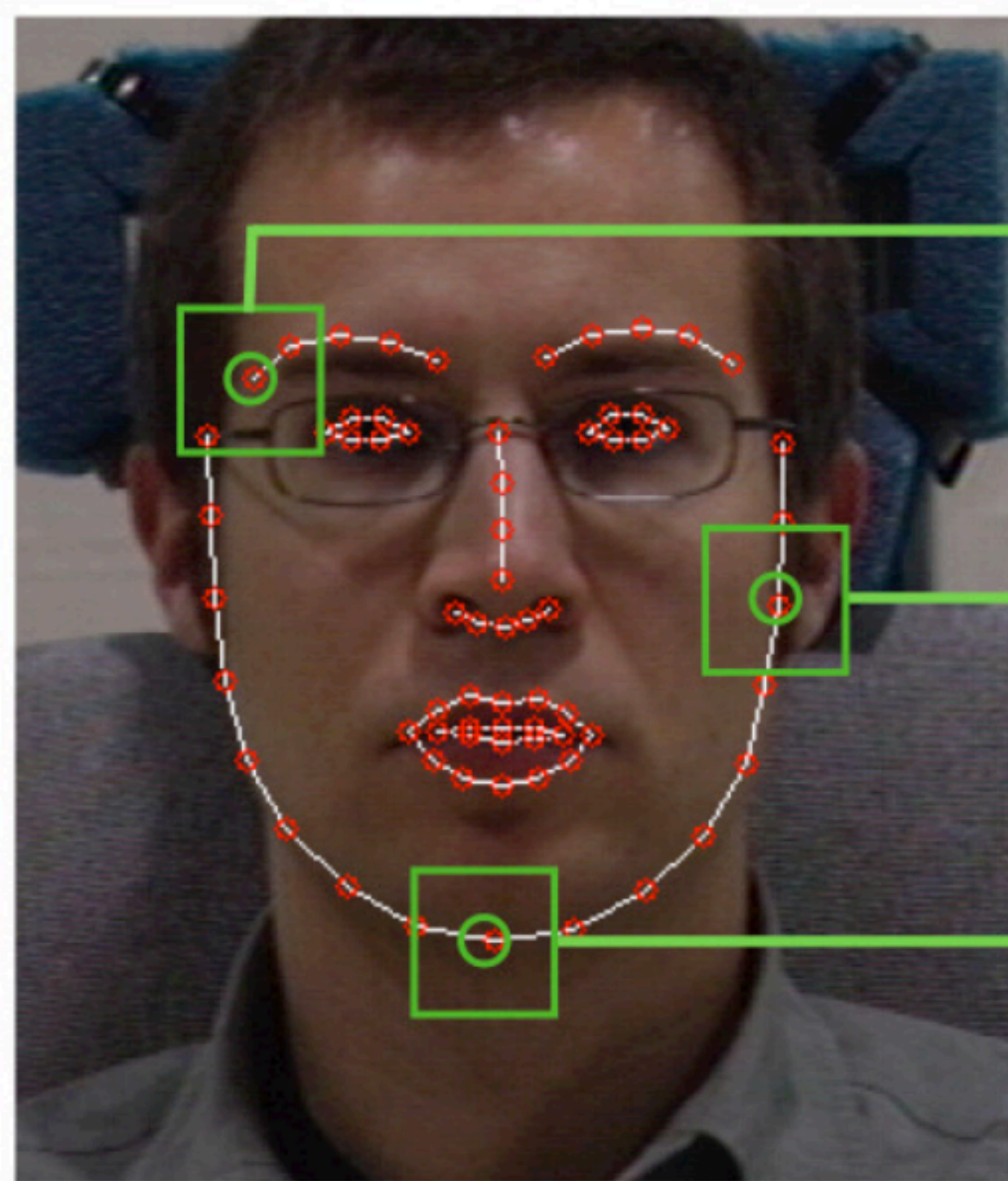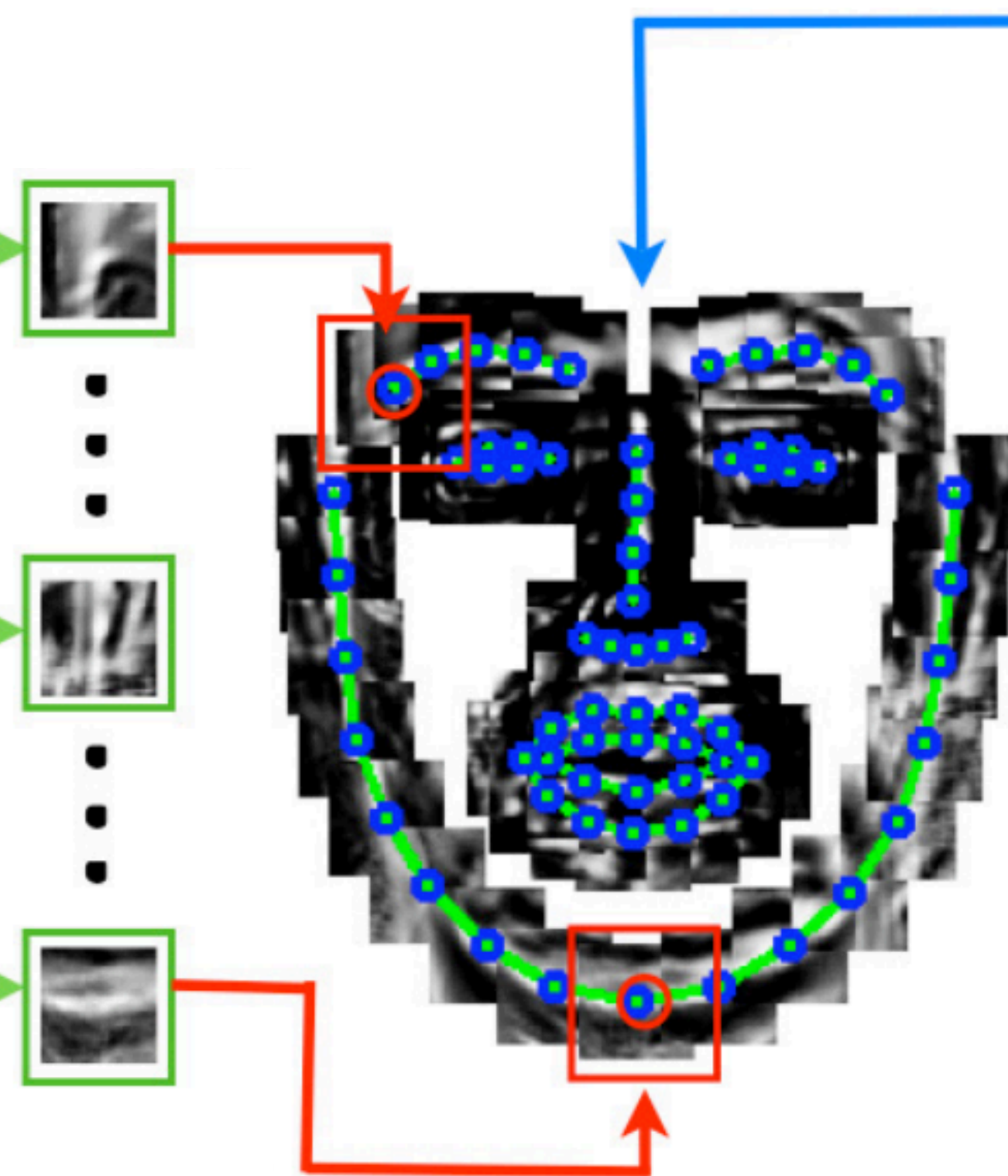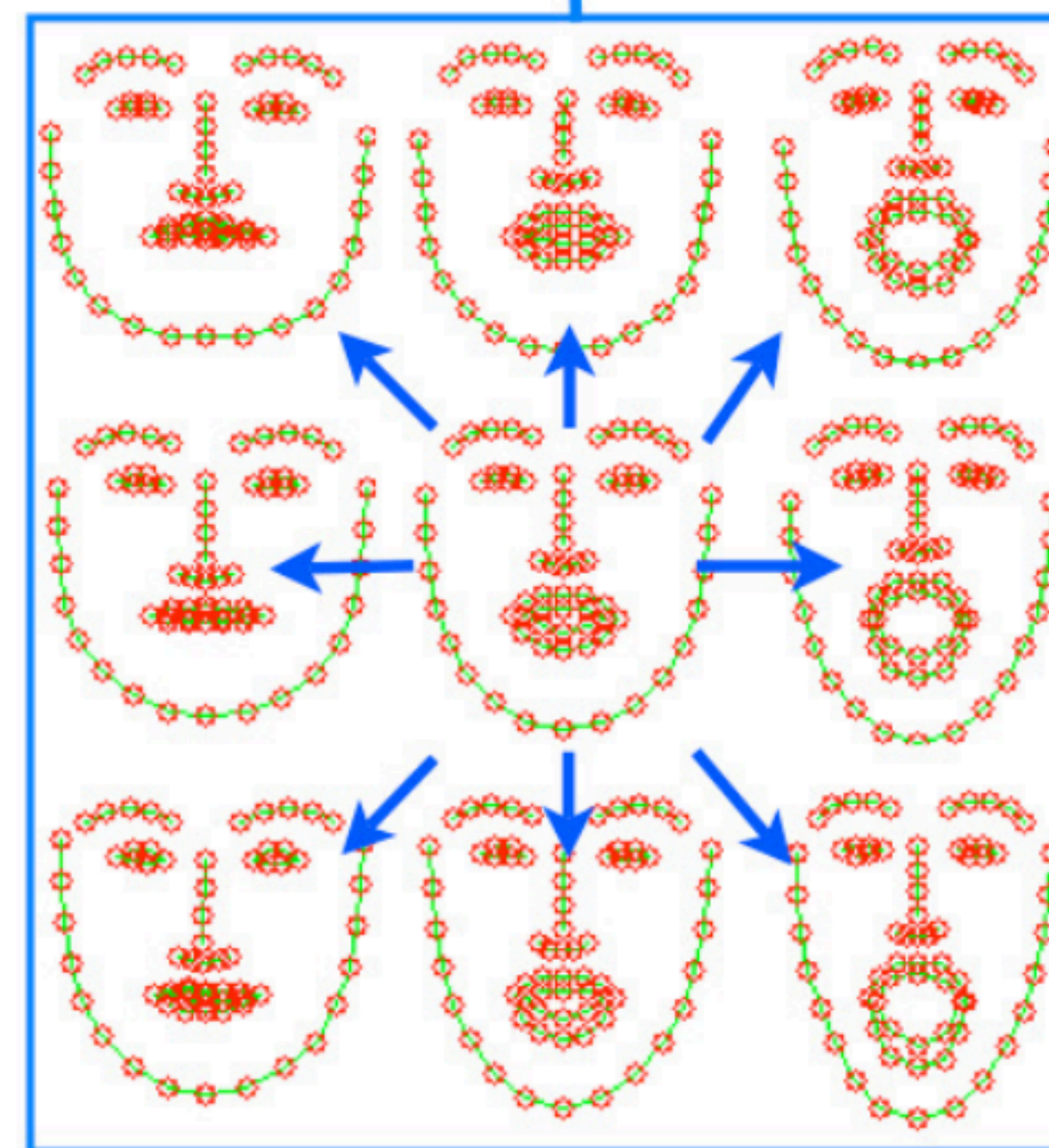# Constrained Local Models

# Constrained Local Models

Image and Search Windows      Optimization      Point Distribution Model

Probability that a landmark's true location is x

Image and Search Windows

Optimization

Point Distribution Model

Probability that a landmark's true location is x

Change current (independent) estimation for landmark location
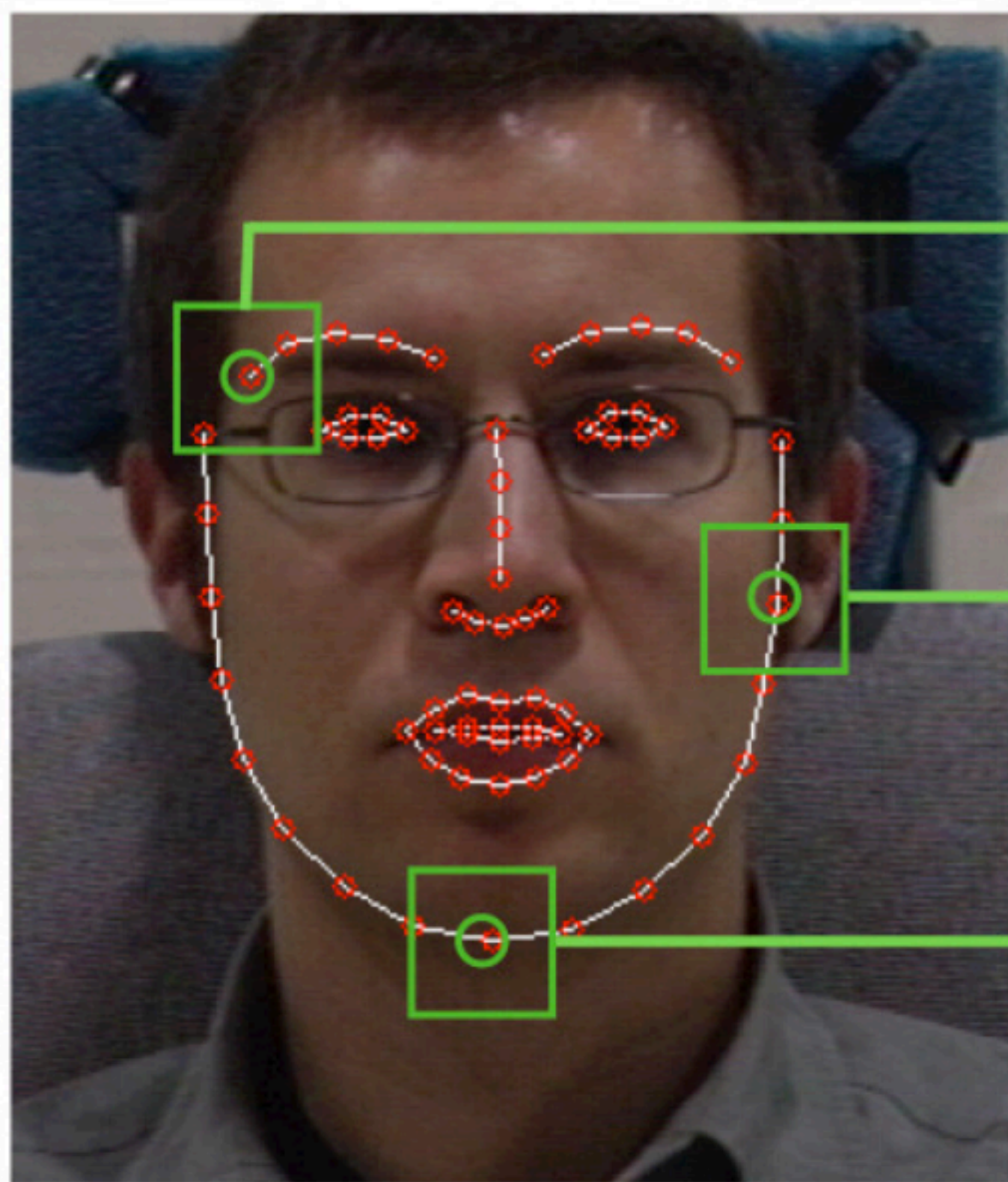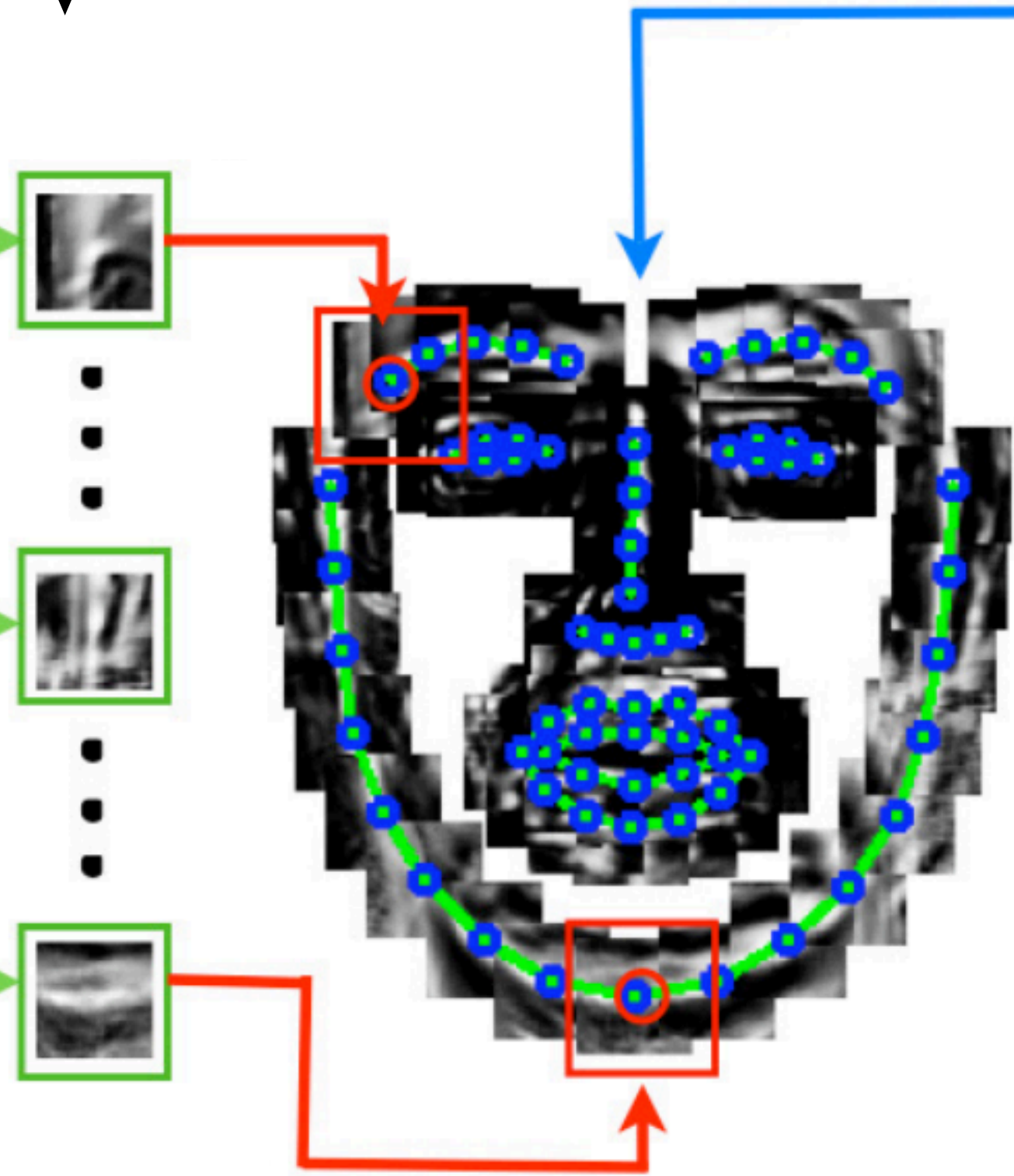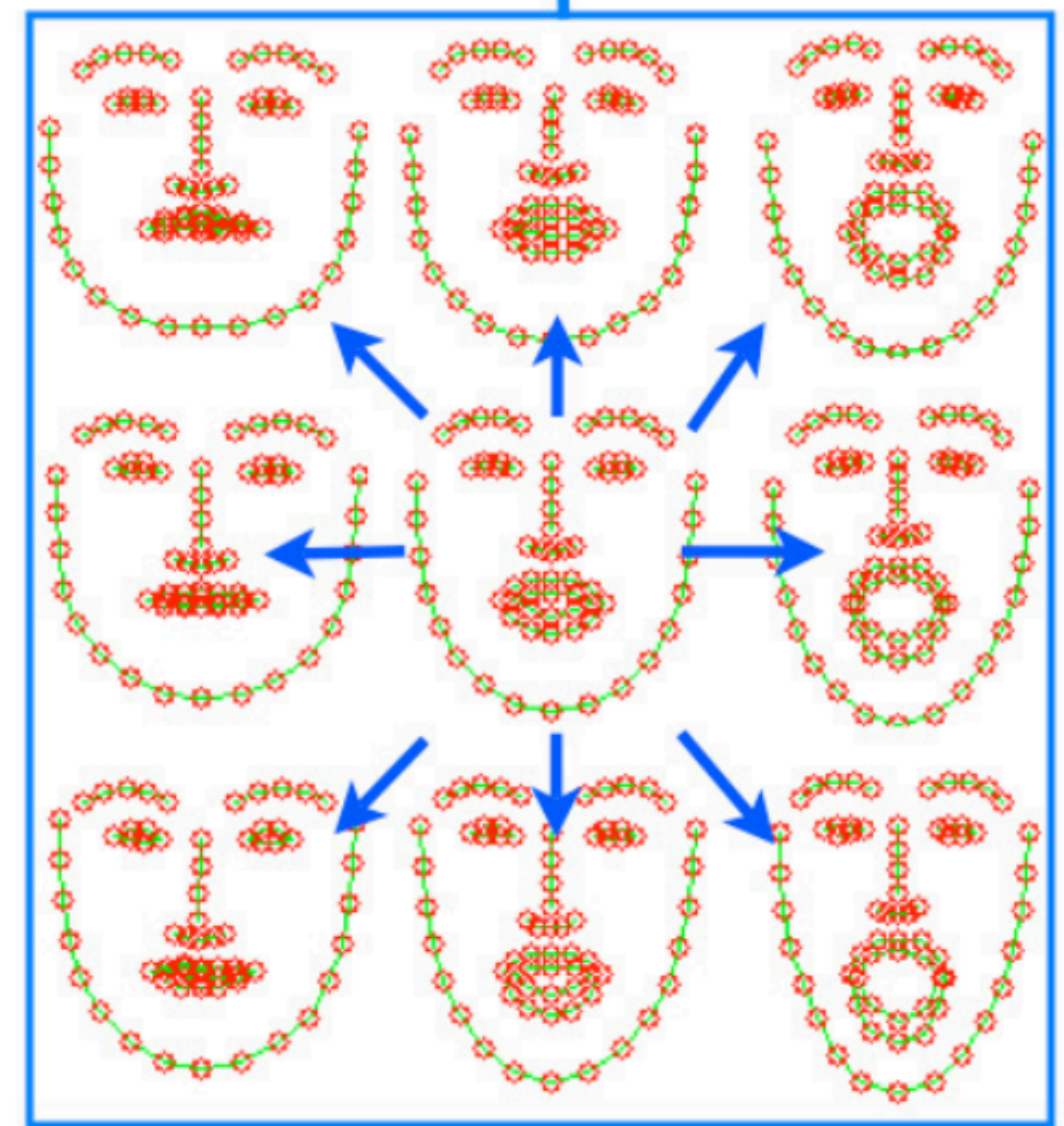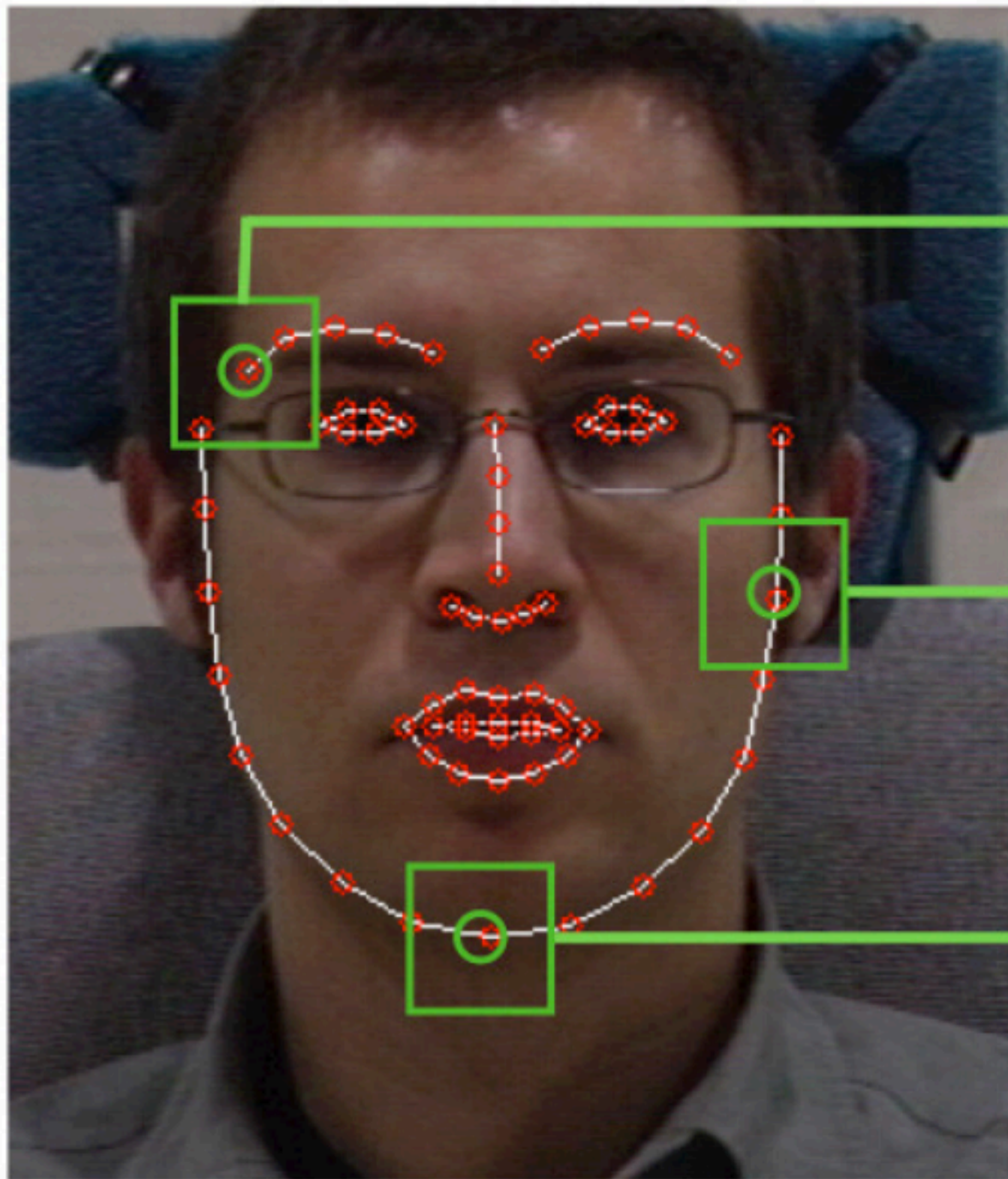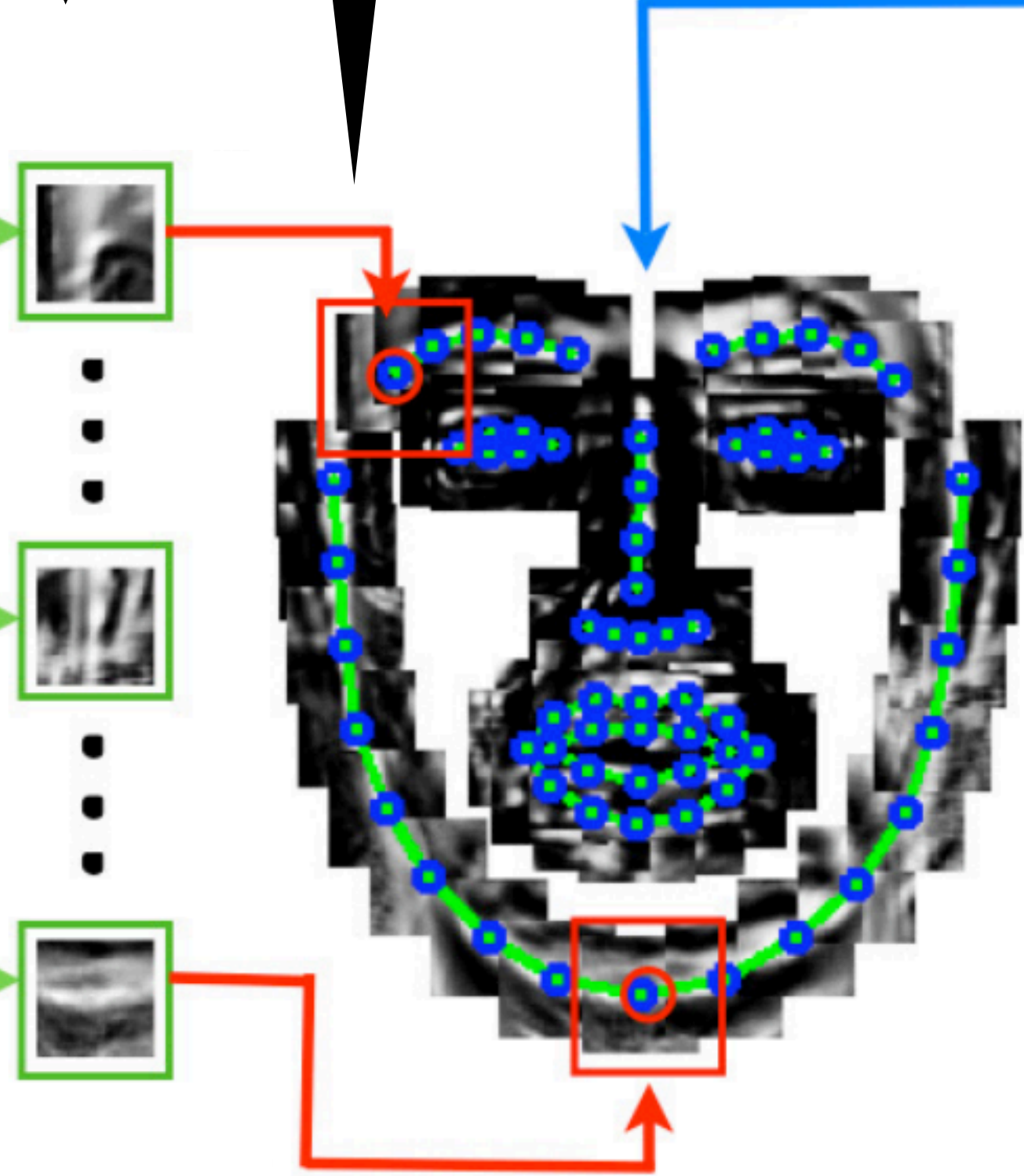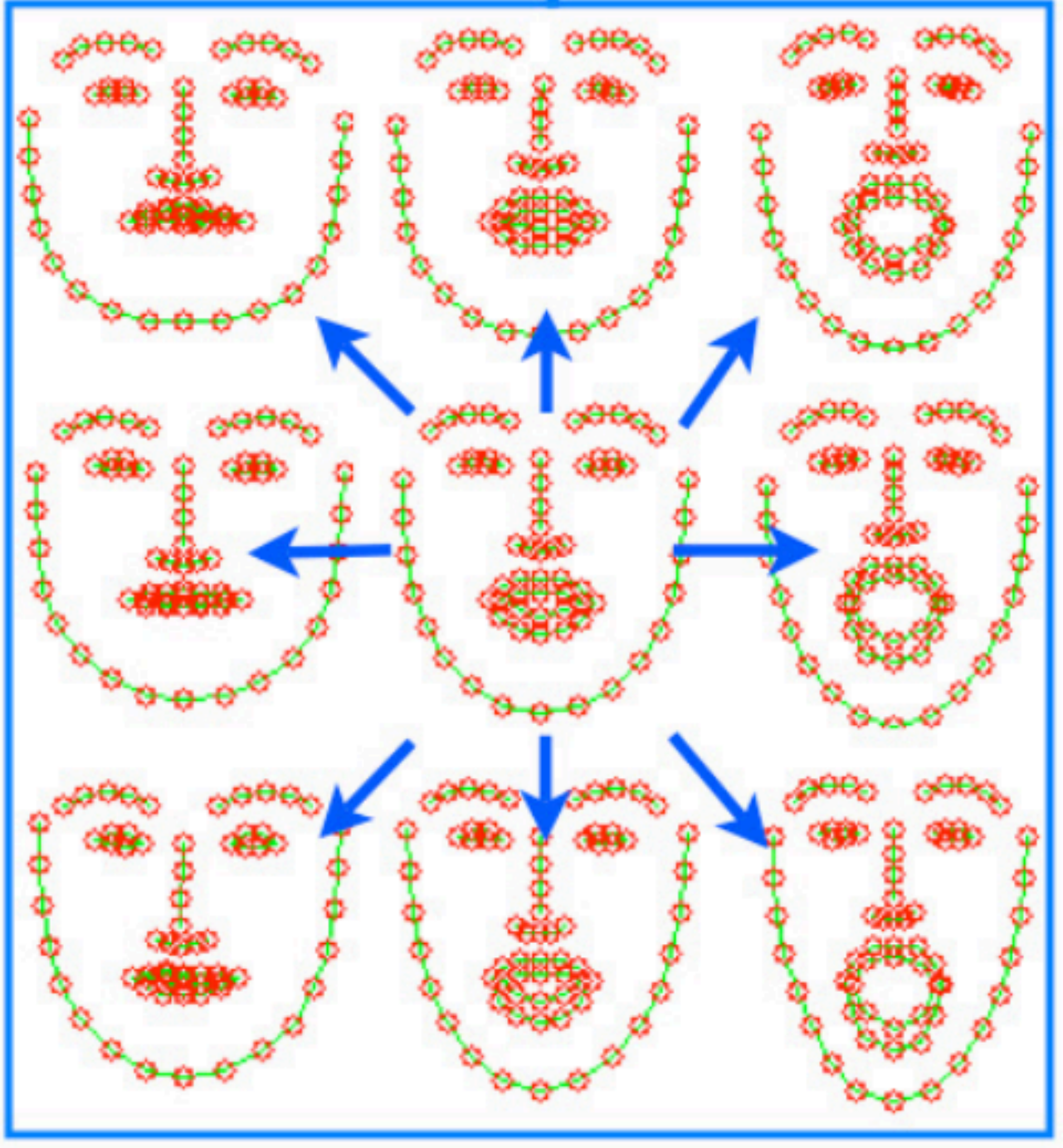
Image and Search Windows

Optimization

Point Distribution Model

Probability that a landmark's true location is x

Change current (independent) estimation for landmark location

Change current (global) estimation for parameterized landmark model

Image and Search Windows

Optimization

Point Distribution Model

# Landmark probabilities



$$p(l_i|\mathbf{x})$$

# Landmark probabilities



$p(l_i|\mathbf{x})$  ASM  CQF  GMM$_5$  KDE$_{20}$  KDE$_5$  KDE$_1$

# Landmark probabilities



$p(l_i|\mathbf{x})$    ASM    CQF    $\text{GMM}_5$    $\text{KDE}_{20}$    $\text{KDE}_5$    $\text{KDE}_1$

# Mean shift

# Mean shift

# Mean shift

Probability that a landmark's true location is x

Change current (independent) estimation for landmark location
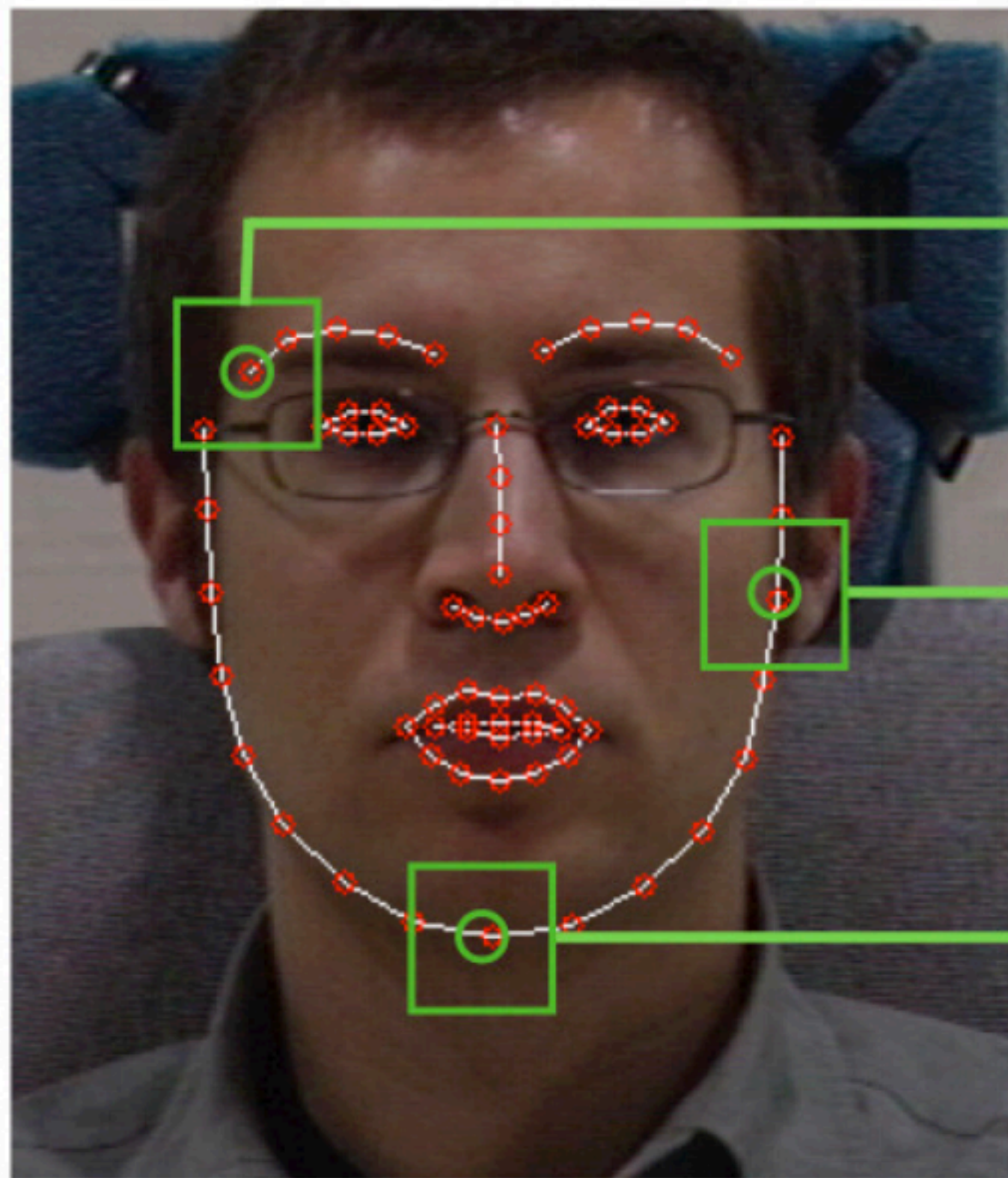
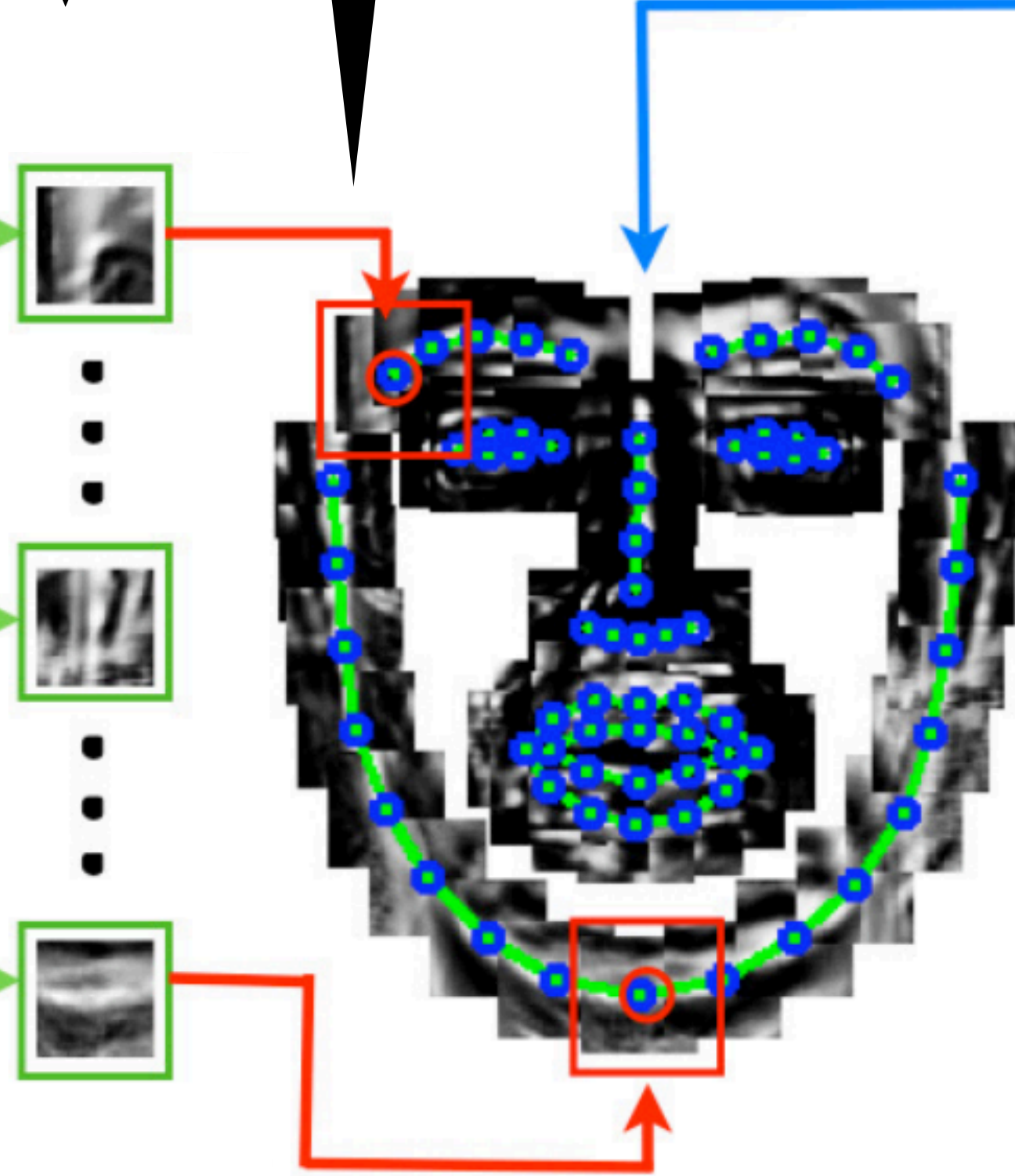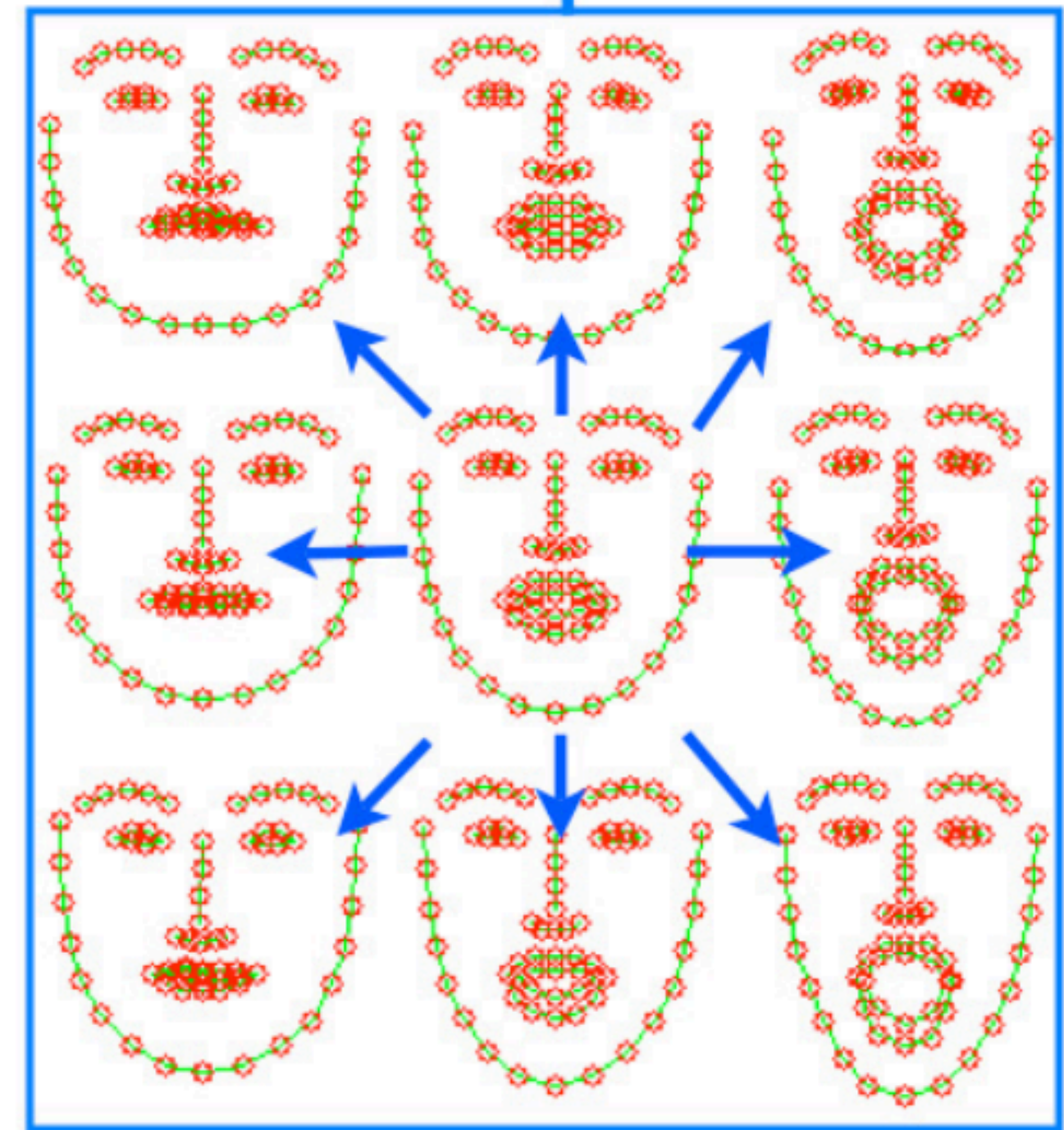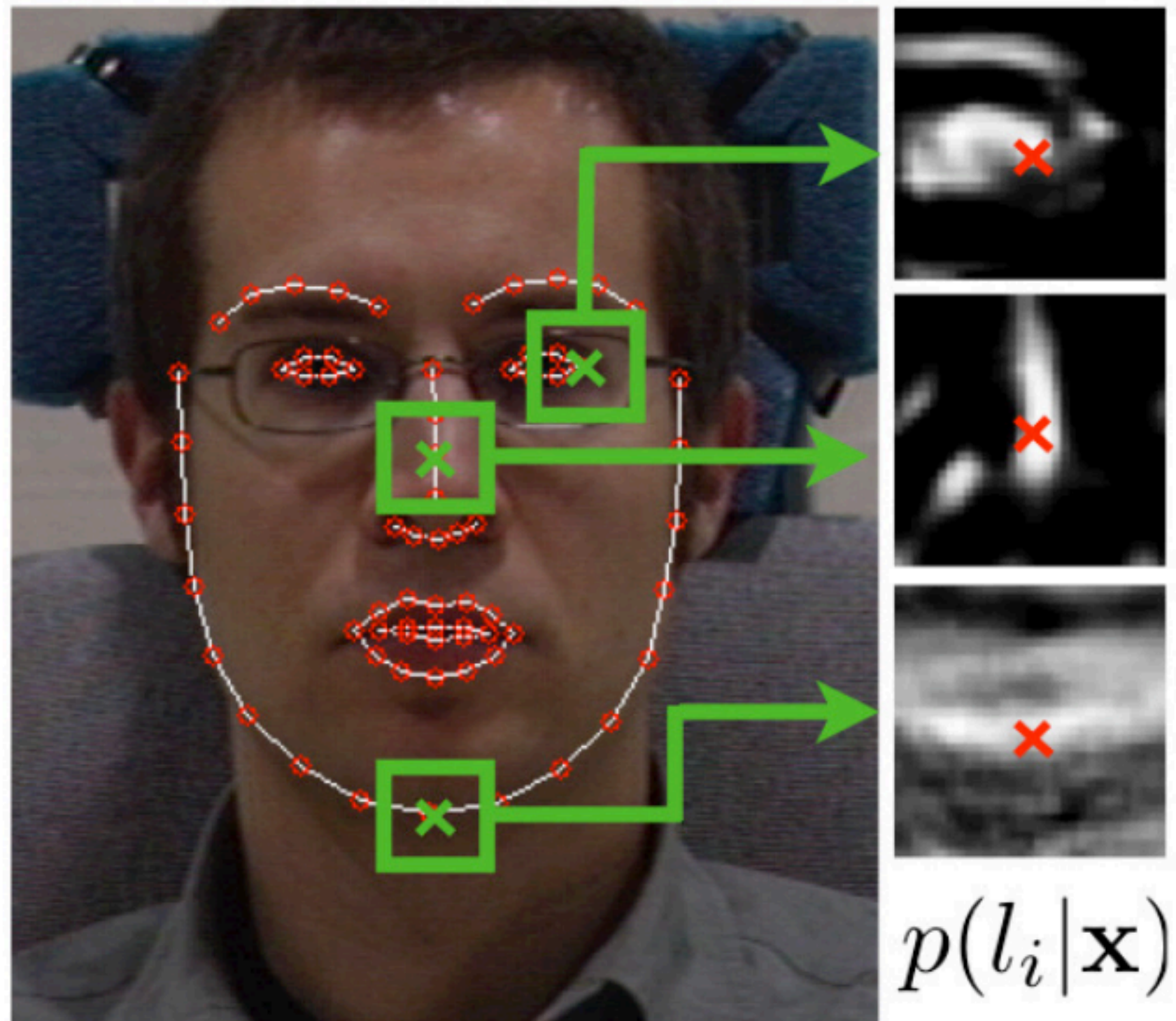Change current (global) estimation for parameterized landmark model
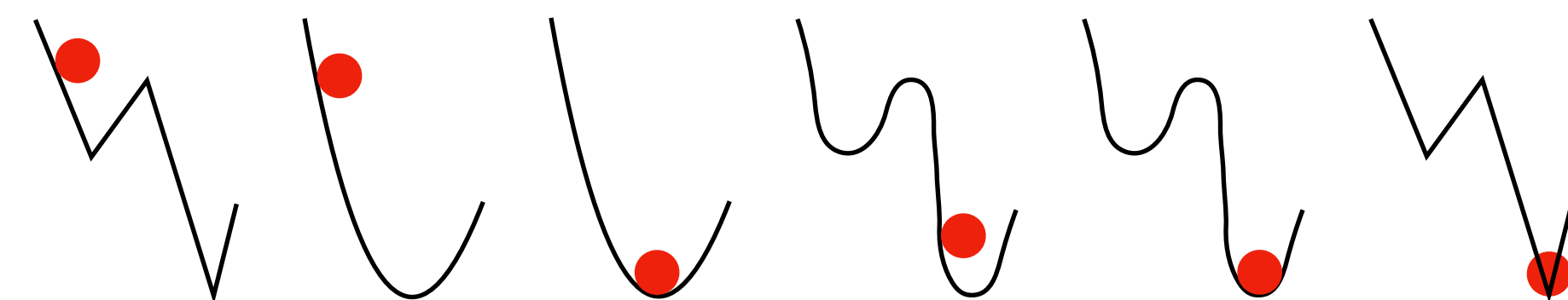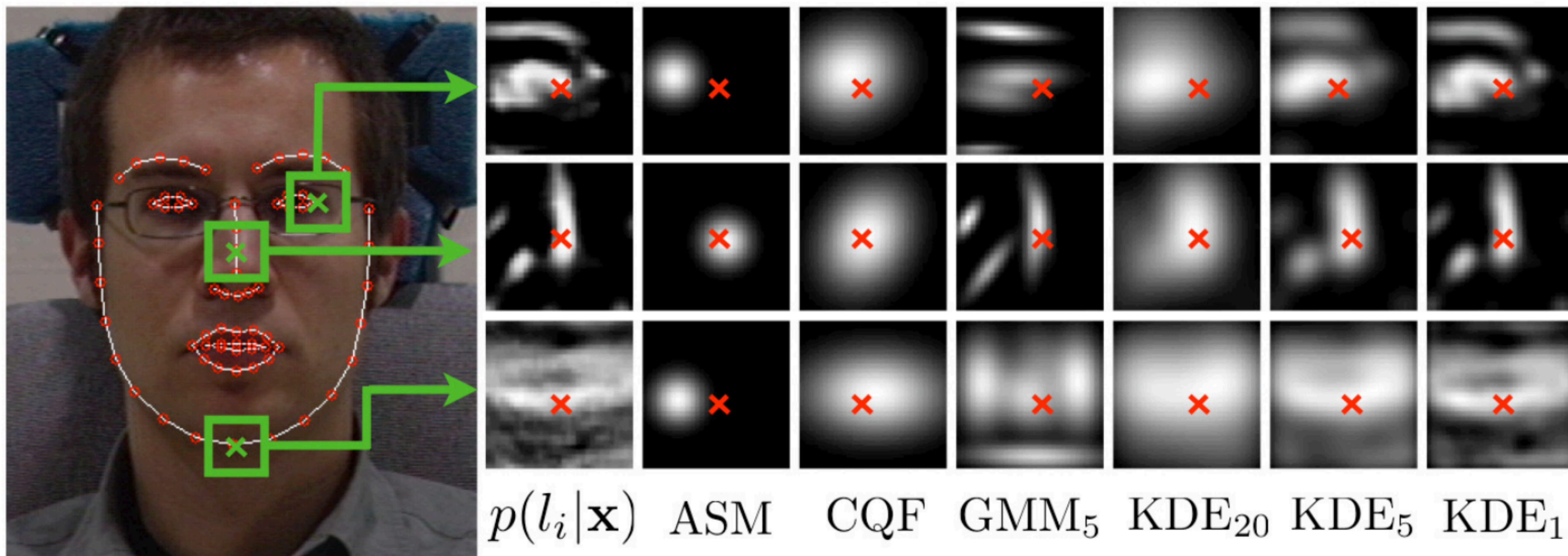
Image and Search Windows

Optimization

Point Distribution Model

# What can you do with sparse face landmarks?

# What can you do with sparse face landmarks?

# What can you do with sparse face landmarks?

## Bootstrap a **better** model

# What can you do with sparse face landmarks?

Bootstrap a **better** model

**dense**

# What can you do with sparse face landmarks?

Bootstrap a **better** model

**dense**

**3D**

# What can you do with sparse face landmarks?

Bootstrap a **better** model

**dense**

**3D**

**richer** (albedo, lighting, …)

# What can you do with sparse face landmarks?

Bootstrap a **better** model

**dense**

**3D**

**richer** **(albedo, lighting, …)**

Sounds familiar?

→ **Structure from motion** → **Multi-view stereo** →

$\longrightarrow$ **Structure from motion** $\longrightarrow$ **Multi-view stereo** $\longrightarrow$

Sparse points that we
"trust"

Dense reconstruction
using information from
previous step

You can also use deep learning…

Constrained Local Models

Improving Portraits

Editing Video

# Improving Portraits

# Improving Portraits

Perspective-aware Manipulation of Portrait Photos

Ohad Fried · Eli Shechtman · Dan B Goldman · Adam Finkelstein

# Distance Matters

# Distance Matters



Pro Photographer
Uses telephoto lens

Selfie
Limited by arm length

# Distance Matters



Pro Photographer — →

Uses telephoto lens

← Selfie →

Limited by arm length

Close ➞ "approachable"

Far ➞ "impressive"

[Perona 2007]

# Camera placement is hard

# Camera placement is hard



- Lack of equipment

# Camera placement is hard

- Lack of equipment

- Lack of expertise

# Camera placement is hard

- Lack of equipment

- Lack of expertise

- Ephemeral moment

# Camera placement is hard

- Lack of equipment

- Lack of expertise

- Ephemeral moment

- Physical constraints

# It is hard to capture the perfect photo...

- Lack of equipment

- Momentary event

- Passersby

- Lack of knowledge

We would like to move the camera in
post processing

# Demo

# Goals

Virtually move perspective camera

- From single image

- Results similar to ground truth

- No background artifacts

- No seams

# Goals

Virtually move perspective camera

- From single image

- Results similar to ground truth

- No background artifacts

- No seams

60cm

480cm

Input 60cm

Simulated 480cm

60cm

480cm

Ground Truth 480cm

Input 60cm        Simulated 480cm        Ground Truth 480cm

# Key Contributions

# Key Contributions

- New application

  - Timely: "year of the selfie"

# Key Contributions

- New application

  - Timely: "year of the selfie"

- New optimization framework

  - Full perspective model

  - Robust for single input photo

# Key Contributions

- New application

  - Timely: "year of the selfie"

- New optimization framework

  - Full perspective model

  - Robust for single input photo

- Fit in 3D, warp in 2D (following, e.g. [Yang '11])

  - Warp field works for perspective model

# Part I: Fitting



Input

Single Image

Head Model

Fiducial
Detection

Parameter Optimization

# Part II: Warping



Change
Model
Parameters

Before     After

Generate
Warp Field

Warp

# Head Model

We want to support **diverse input photos**

We need an **expressive model**

# Head Model

We want to support **diverse input photos**

We need an **expressive model**

What should it include?

# Head Model

# Head Model



- Head shape ("identity")

# Head Model

- Head shape ("identity")

- Bone/muscle layout ("expression")

# Head Model

- Head shape ("identity")

- Bone/muscle layout ("expression")

- Location and pose, relative to camera

# Head Model

- Head shape ("identity")

- Bone/muscle layout ("expression")

- Location and pose, relative to camera

- Internal camera parameters

FaceWarehouse

[Cao et al. 2014]

# Aligned models

# Aligned models



[Cao et al. 2014]

# Aligned models



[Cao et al. 2014]

# Aligned models

identities

expressions

[Cao et al. 2014]

11K vertices
x
150 identities
x
47 expressions

11K vertices
x
150 identities
x
47 expressions

We can combine these
to create new heads!

11K vertices
x
150 identities
x
47 expressions

We can combine these
to create new heads!

HOSVD ⬠

# Short detour — SVD

$$M_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}$$

$$M_{3 \times 4} = \begin{bmatrix} | & | & | \\ u_1 & u_2 & u_3 \\ | & | & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \end{bmatrix} \begin{bmatrix} - & v_1 & - \\ - & v_2 & - \\ - & v_3 & - \\ - & v_4 & - \end{bmatrix}$$

# Short detour — SVD

$$M_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}$$

$$M_{3 \times 4} = \begin{bmatrix} | & | & | \\ u_1 & u_2 & u_3 \\ | & | & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \end{bmatrix} \begin{bmatrix} - & v_1 & - \\ - & v_2 & - \\ - & v_3 & - \\ - & v_4 & - \end{bmatrix}$$

# Short detour — SVD

$$M_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}$$

$$M_{3 \times 4} = \begin{bmatrix} | & | & | \\ u_1 & u_2 & u_3 \\ | & | & | \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \end{bmatrix} \begin{bmatrix} - & v_1 & - \\ - & v_2 & - \\ - & v_3 & - \\ - & v_4 & - \end{bmatrix}$$

11K vertices
150 identities
47 expressions

34530x150x47

HOSVD

Core Tensor $\otimes_1$ $A_{34530 \times 40}$ $\otimes_2$ $B_{150 \times 50}$ $\otimes_3$ $\Gamma_{47 \times 25}$

40 x 50 x 25

11K vertices
150 identities
47 expressions

34530x150x47

HOSVD

Low rank approximation:
✓ less space
✓ less noise

Core Tensor

$\otimes_1$   $A_{34530x40}$   $\otimes_2$   $B_{150x50}$   $\otimes_3$   $\Gamma_{47x25}$

40 x 50 x 25

11K vertices
150 identities
47 expressions

34530x150x47

HOSVD

Low rank approximation:
✓ less space
✓ less noise

Core Tensor $\otimes_1$ $A_{34530 \times 40}$ $\otimes_2$ $B_{150 \times 50}$ $\otimes_3$ $\Gamma_{47 \times 25}$

40 x 50 x 25

And for a single head:

Core Tensor $\otimes_1$ $A_{34530 \times 40}$ $\otimes_2$ $\beta_{1 \times 50}$ $\otimes_3$ $\gamma_{1 \times 25}$

40 x 50 x 25

$$\begin{bmatrix} f_x & s & x_0 \\ 0 & f_y & y_0 \\ 0 & 0 & 1 \end{bmatrix}$$

Camera Calibration Matrix $*$ Rotation $*$ Translation $*$

HOSVD

Core Tensor $\otimes_1 A \otimes_2 \beta \otimes_3 \gamma$

Specific Head

$$\begin{bmatrix} f_x & s & x_0 \\ 0 & f_y & y_0 \\ 0 & 0 & 1 \end{bmatrix}$$

Camera Calibration Matrix * Rotation * Translation * HOSVD Specific Head

Core Tensor $\otimes_1 A \otimes_2 \beta \otimes_3 \gamma$

(3 of 5)    (3)    (3)    (50)  (25)

3 + 3 + 3 + 50 + 25 = 84 parameters

# Fitting

Input

Fiducial Detection

Single Image

Head Model

Parameter Optimization

# Fitting

Input



Single Image

Head Model

Fiducial Detection

Parameter Optimization

# Fiducial points



66 points

[Saragih et al. 2009]

# Fiducial points



66 points

[Saragih et al. 2009]

3 points

(manually marked)

# Fitting

Input

Fiducial Detection

Single Image

Head Model

Parameter Optimization

# Fitting



good     bad

Plug in parameters

# Optimization



Solve rotation & translation → Solve identity → Solve expression → Solve camera → Update 3D landmarks

Coordinate descent

# Fitting

Input



Single Image

Head Model

Fiducial Detection

Parameter Optimization

# Part II: Warping



Change Model Parameters

Before    After

Generate Warp Field

Warp

# Warping



Change
Model
Parameters

Before     After

Generate
Warp Field

Warp

# Warp field



Triangulation based interpolation

# Warp field



Discrete Laplace operator

Averaging filter

# Warping



Change Model Parameters

Before    After

Generate Warp Field

Warp

# Recap

# Recap

# Recap

# Recap

# Recap

# Recap

# Evaluation

Radial | Face only | Mean head | Ours

Radial Face only Mean head Ours

Radial

Face only

vs.

Mean head

Ours

Radial

Face only

Mean head

Ours

vs.

[Vlasic et al. 2005]

[Yang et al. 2011]

...

Radial

Face only

VS.

[Vlasic et al. 2005]

[Yang et al. 2011]

...

Mean head

no expression
or identity

Ours

# Radial

# Face only

vs.

[Vlasic et al. 2005]

[Yang et al. 2011]

...

# Mean head

no expression
or identity

[Hassner et al. 2015]

# Ours

Radial

Face only

Mean head

Ours

vs.

[Vlasic et al. 2005]

[Yang et al. 2011]

...

no expression
or identity

[Hassner et al. 2015]

Full model

| Input | Ground truth | Radial | Face only | Mean head | Ours |

no error  high error

Input                    Mean head                    Ours

# Rotation



Input

# Rotation



Left 4°

# Rotation



Left 7°

# Rotation



Input

# Rotation



Right 4°

# Rotation



Right 7°

# Rotation

Useful for
Face
Identification!

(e.g. DeepFace '14)

Right 7°

# Rotation



Not Perfect

Useful for Face Identification!

(e.g. DeepFace '14)

Right 7°

# Limitations & Future Work

# Limitations & Future Work

- Fitting may fail

# Limitations & Future Work

- Fitting may fail

  - Use more data — RGB values

# Limitations & Future Work

- Fitting may fail

  - Use more data — RGB values

- Does not support extreme rotations

# Limitations & Future Work

- Fitting may fail

  - Use more data — RGB values

- Does not support extreme rotations

  - Hallucinate missing features

More results and demo:

# http://faces.cs.princeton.edu

Constrained Local Models

Improving Portraits

Editing Video

# Editing
# Video

# Editing Video

**Text-based Editing of Talking-head Video**

Ohad Fried · Ayush Tewari · Michael Zollhöfer · Adam Finkelstein · Eli Shechtman ·
Dan Goldman · Kyle Genova · Zeyu Jin · Christian Theobalt · Maneesh Agrawala

Video Blogs                Interviews                Online Courses




Speeches                   Commercials                          ...

In these videos we care mostly about the spoken transcript

In these videos we care mostly about the spoken transcript

But we edit them just like any other video…

Bin: CNMT_Mom_Never_Wrong_31312 × | Project: CMNT test | Source: Clip #14 ▾ × | Effect Controls | Audio Mixer: test sequence | Metadata | Program: test sequence ▾ × 

CMNT test.prproj\CNMT_Mom_Never_Wrong_31312    61 Items

In: All ▾

Clip #14    19;17
Clip #15    4;27
Clip #19    2;15
Clip #3    37;08
Clip #4    2;27
Clip #7    6;01

SODA SPECIAL $3⁰⁰
CHIPS 2 LESS 1.00
FRESCAS
Los BAJOS Precios
VERDURAS    FRUTAS
BOYLE HEIGHTS

Mott St

00;11;09;24    Fit ▾    1/4 ▾    00;00;19;17

00;01;30;16    Fit ▾    1/4 ▾    00;00;15;24

Media Browser | Info | Effects × | Markers | History

Presets
Audio Effects
Audio Transitions
Video Effects
Video Transitions

test sequence × | rough cut | Sequence 01

00;01;30;16

01;04;02  00;01;08;02  00;01;12;02  00;01;16;02  00;01;20;02  00;01;24;02  00;01;28;02  00;01;32;02  00;01;36;02  00;01;40;02  00;01;44;02  00;01;48;02  00;01;52;02  00;01;5

Video 4
Video 3
Video 2    Clip #14 .city:Opacity ▾
Video 1    Clip #6 Opacity:Opacity ▾  Clip #6 [V] Opacity:Opacity ▾  Clip #  Clip #10 [' Clip #9 [V]  Clip #15  Clip #15 [V] y ▾ Clip #15 [V]  Clip #  Clip Clip #5 [264%] :ity ▾ Clip #8 [V] Opacity:Opacity ▾

A1    Audio 1    Clip #6 Volume:Level ▾  Clip #6 [A] Volume:Level ▾  Clip #  Clip #10 [  Clip #9 [A]  Clip #15  Clip #15 [A] ' ▾ Clip #15 [A]  Clip #  C    Clip #8 [A] Volume:Level ▾

A2    Audio 2    Clip #6 Volume:Level ▾    10-09 Grocery Store Ambience Near Cashiers; Medium Elect

Audio 3    Airport Lounge.aif Volume:Level ▾    10-09 Grocery Store Ambience Near Cashiers; Medium Elect

Audio 4
Audio 5

| Task | Current | Proposed |
| --- | --- | --- |
|  |  |  |
|  |  |  |
|  |  |  |

| Task | Current | Proposed |
| --- | --- | --- |
| Compose a sentence from multiple takes | | |

| Task | Current | Proposed |
| --- | --- | --- |
| Compose a sentence from multiple takes | Jump cuts | |

| Task | Current | Proposed |
| --- | --- | --- |
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |

| Task | Current | Proposed |
|---|---|---|
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |
| Delete words | | |
| | | |

| Task | Current | Proposed |
|------|---------|----------|
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |
| Delete words | Jump cuts | |

| Task | Current | Proposed |
|---|---|---|
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |
| Delete words | Jump cuts | Seamless transitions |

| Task | Current | Proposed |
|------|---------|----------|
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |
| Delete words | Jump cuts | Seamless transitions |
| Change or add words | | |

| Task | Current | Proposed |
| --- | --- | --- |
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |
| Delete words | Jump cuts | Seamless transitions |
| Change or add words | New recording session | |

| Task | Current | Proposed |
|------|---------|----------|
| Compose a sentence from multiple takes | Jump cuts | Seamless transitions |
| Delete words | Jump cuts | Seamless transitions |
| Change or add words | New recording session | Synthesize new video |

"The market closed today with apple stock price at one hundred and ninety one point four five dollars per share"

"The market closed today with apple stock price at
one hundred and ~~ninety one point four~~ five dollars per share"
eighty two point two

# Idea: use existing snippets to construct new words

# Reframed as a selective interpolation problem

Only use expressions

# Reframed as a selective interpolation problem

# Method overview

*The quick brown ...*

Input video

Preprocessing

Input video

*The quick brown ...*

DH IY0 K W IH1 K sp B R AW1 N

Phoneme Alignment

**Input video**

*The quick brown ...*

Preprocessing

DH IY0 K W IH1 K sp B R AW1 N

**Phoneme Alignment**

→

Head
Parameters

Model

**Tracking & Reconstruction**

Input video

*The quick brown ...*

Preprocessing

Phoneme Alignment

DH IY0 K W IH1 K sp B R AW1 N

Tracking & Reconstruction

Head Parameters

Model

*The quick brown spider jumped fox*

Edit Operation

**Preprocessing**

**Input video**

*The quick brown ...*

**Phoneme Alignment**

DH IY0 K W IH1 K sp B R AW1 N

**Tracking & Reconstruction**

Head Parameters

Model

**Edit Operation**

*The quick brown spider jumped* ~~*fox*~~

**Viseme Search**

*The quick brown spider ...*
*... viper ... ox ...*
f ox

Preprocessing

Input video

The quick brown ...

DH IY0 K W IH1 K sp B R AW1 N

Phoneme Alignment

Head Parameters

Model

Tracking & Reconstruction

Edit Operation

The quick brown spider jumped fox

The quick brown spider ...
... viper ... ox ...

f ox

Viseme Search

brown    v    ox    jumped

Parameter Blending

Preprocessing

Background Retiming

**Input video**

*The quick brown ...*

**Phoneme Alignment**

DH IY0 K W IH1 K sp B R AW1 N

**Tracking & Reconstruction**

Head Parameters

Model

**Edit Operation**

*The quick brown spider jumped fox*

**Viseme Search**

*The quick brown spider ...*
*... viper ... ox ...*

f ox

**Parameter Blending**

*brown*   *v*   *ox*   *jumped*

**Rendering**

Preprocessing

Background Retiming

**Input video**

*The quick brown ...*

**Edit Operation**

*The quick brown spider jumped fox*

**Phoneme Alignment**

DH IY0 K W IH1 K sp B R AW1 N

**Viseme Search**

*The quick brown spider ...*
*... viper ... ox ...*
f ox

**Tracking & Reconstruction**

Head Parameters

Model

**Parameter Blending**

brown          v          ox          jumped

**Rendering**

Preprocessing

Background Retiming

The quick brown ...

Input video

DH IY0 K W IH1 K sp B R AW1 N

Phoneme Alignment

Head Parameters

Model

Tracking & Reconstruction

The quick brown
spider jumped
fox

Edit Operation

The quick brown spider ...
... viper ... ox ...

f ox

Viseme Search

brown        v        ox        jumped

Parameter Blending

Rendering

New stuff

# Viseme search

# Viseme search

- Identical phonemes are likely to be visually similar

# Viseme search

- Identical phonemes are likely to be visually similar

- Same for visemes (but less so)

# Viseme search

- Identical phonemes are likely to be visually similar

- Same for visemes (but less so)

- **Cannot** expect to find a good coherent viseme sequence for long edits

# Viseme search

- Identical phonemes are likely to be visually similar

- Same for visemes (but less so)

- **Cannot** expect to find a good coherent viseme sequence for long edits

  - Instead, find several matching subsequences and combine

# Viseme search



Query

Split

Match

# Viseme search

# Viseme search

# Viseme search

# Viseme search

**Modified Levenshtein edit distance**

# Viseme search

**Modified Levenshtein edit distance**

$$C_{insert} = 1$$

$$C_{delete} = 1$$

# Viseme search

**Modified Levenshtein edit distance**

$$C_{insert} = 1$$

$$C_{delete} = 1$$

$$C_{swap} =$$

# Viseme search

## Modified Levenshtein edit distance

$$C_{insert} = 1$$

$$C_{delete} = 1$$

$$C_{swap} = C_{vis}(\blacksquare, \blacksquare)$$

# Viseme search

**Modified Levenshtein edit distance**

$$C_{insert} = 1$$

$$C_{delete} = 1$$

$$C_{swap} = C_{vis}(\blacksquare, \blacksquare)(|\blacksquare| + |\blacksquare|)$$

# Viseme search

**Modified Levenshtein edit distance**

$$C_{insert} = 1$$

$$C_{delete} = 1$$

$$C_{swap} = C_{vis}(\blacksquare, \blacksquare\blacksquare)(|\blacksquare| + |\blacksquare\blacksquare|) + \chi \left| |\blacksquare| - |\blacksquare\blacksquare| \right|$$

# Viseme search

# Parameter blending

# Parameter blending



Head Parameters    Model

# Parameter blending

- Geometry


Head Parameters   Model

# Parameter blending



Head Parameters — Model

- Geometry

- Albedo

# Parameter blending


Head Parameters → Model

- Geometry

- Albedo

- Illumination

# Parameter blending



Head Parameters → Model

- Geometry

- Albedo

- Illumination

- Pose

# Parameter blending



Head Parameters → Model

- Geometry

- Albedo

- Illumination

- Pose

- Expression

# Parameter blending



Head Parameters → Model

- Geometry

- Albedo

**Constant**

- Illumination

- Pose

- Expression

# Parameter blending



Head Parameters    Model

- Geometry

- Albedo

**Constant**

- Illumination

**Linear interpolation in new region**

- Pose

- Expression

# Parameter blending



Head Parameters    Model

- Geometry

- Albedo

**Constant**

- Illumination

**Linear interpolation in new region**

- Pose

**Later…**

- Expression

# Parameter blending



Head Parameters → Model

- Geometry

- Albedo

**Constant**

- Illumination

**Linear interpolation in new region**

- Pose

**Later…**

- Expression

**Linear interpolation between snippets**

# Background retiming

# Background retiming

**… The quick brown spider jumped …**

**… The quick brown fox jumped …**

# Background retiming

**… The quick brown spider jumped …**

Longer

Shorter

**… The quick brown fox jumped …**

# Background retiming

… **The** **quick brown spider** **jumped** …

Longer

Shorter

… **The** **quick brown fox** **jumped** …

We want localized edits. Everything else should stay the same

# Background retiming

# Background retiming

- Use longer sequence (even if edit is short)

# Background retiming

- Use longer sequence (even if edit is short)

- Calculate number of frames to add / remove

# Background retiming

- Use longer sequence (even if edit is short)

- Calculate number of frames to add / remove

  - Spread equally

# Background retiming

- Use longer sequence (even if edit is short)

- Calculate number of frames to add / remove

  - Spread equally

- Long enough —> no retiming artifacts

# Background retiming

- Use longer sequence (even if edit is short)

- Calculate number of frames to add / remove

  - Spread equally

- Long enough —> no retiming artifacts

- Pose parameters taken from retimed background

# Neural Face Rendering

# Neural Face Rendering

# Neural Face Rendering

# Neural Face Rendering

**Loss** = Photometric + Spatial adversarial + Temporal adversarial

Input video

*The quick brown ...*

Preprocessing

Phoneme Alignment

DH IY0 K W IH1 K sp B R AW1 N

Tracking & Reconstruction

Head Parameters

Model

Background Retiming

Edit Operation

*The quick brown spider jumped fox*

Viseme Search

*The quick brown spider ...*
*... viper ... ox ...*

f ox

Parameter Blending

*brown*    *v*    *ox*    *jumped*

Rendering

# Results

Input

Input

Visual voice assistant

Visual voice assistant

**Translation**

Blend with input (temporal)

Real

Blend with input
(spatial)

Voice by VoCo
[Jin et al. '17]

**Voice by VoCo**
**[Jin et al. '17]**

# Evaluation & comparisons

# Parameter blending



**Without blending**

**With blending**

# Parameter blending



**Without blending**                                    **With blending**

# Parameter blending



Without blending                    With blending

# Parameter blending



Without blending

With blending

# Parameter blending



**Without blending**                    **With blending**

# Parameter blending



**Without blending**          **With blending**

# Dataset size



5% data                10% data                50% data                100% data

# Dataset size



5% data           10% data           50% data           100% data

# Dataset size



5% data       10% data       50% data       100% data

# Dataset size



5% data      10% data      50% data      100% data

# Dataset size



5% data          10% data          50% data          100% data

**Ground truth**

**Deep Video Portraits
[Kim et al. '18]**

**Ours**

**Ground truth**

**Deep Video Portraits**
**[Kim et al. '18]**

**Ours**

Ground truth

Deep Video Portraits
[Kim et al. '18]

Ours

**Ground truth**

**Deep Video Portraits
[Kim et al. '18]**

**Ours**

**Ground truth**

**Deep Video Portraits**
**[Kim et al. '18]**

**Ours**

**Ground truth**

**Deep Video Portraits**
**[Kim et al. '18]**

**Ours**

**Ground truth**

**Deep Video Portraits**
**[Kim et al. '18]**

**Ours**

Ground truth

Deep Video Portraits
[Kim et al. '18]

Ours

**Face2Face**
**[Thies et al. 18]**

**Ours**

Face2Face
[Thies et al. 18]

Ours

**Morph Cut**
**[Berthouzoz et al. 12]**

**Ours**

"It's one of the approaches ~~to~~ to machine learning"

**Morph Cut**
**[Berthouzoz et al. 12]**

**Ours**

"It's one of the approaches ~~to~~ to machine learning"

**Morph Cut**
**[Berthouzoz et al. 12]**

**Ours**

"It's one of the approaches ~~to~~ to machine learning"

**Morph Cut**
**[Berthouzoz et al. 12]**

**Ours**

"… learning from examples ~~and~~ and scientists …"

**Morph Cut**
**[Berthouzoz et al. 12]**

**Ours**

"… learning from examples ~~and~~ and scientists …"

**Morph Cut**
**[Berthouzoz et al. 12]**

**Ours**

"… learning from examples ~~and~~ and scientists …"

# User study

"This video clip looks real to me"

% of scores

**Legend:** GT Base Videos · GT Target Videos · Our Modified Videos

X-axis: 5 (strongly agree), 4, 3, 2, 1 (strongly disagree)

% of scores

"This video clip looks real to me"

Our results rated 'real' in 59.6% of cases

50

37.5

25

12.5

0

5 (strongly agree)   4   3   2   1 (strongly disagree)

■ GT Base Videos   ■ GT Target Videos   ■ Our Modified Videos

# Limitations & future work

# Limitations & future work

- Moods and facial expressions

  - We might blend incompatible sequences

  - Can we control these when synthesizing?

# Limitations & future work

- Moods and facial expressions

  - We might blend incompatible sequences

  - Can we control these when synthesizing?

- Viseme search is slow

  - Can speed up with some relaxations

# Limitations & future work

- Moods and facial expressions

  - We might blend incompatible sequences

  - Can we control these when synthesizing?

- Viseme search is slow

  - Can speed up with some relaxations

- Interactivity

  - Algorithm speedups

  - Editing UI

Constrained Local Models

Improving Portraits

Editing Video

Image and Search Windows    Optimization    Point Distribution Model

Constrained Local Models

Improving Portraits

Editing Video



Scene Description

Illumination $\gamma$

Pose $(R, t)$

Fine Layer

Detail p

Coarse Layer

Shape $\alpha$     Albedo $\beta$     Expression $\delta$

Medium Layer

Correctives $\tau$



Image and Search Windows          Optimization          Point Distribution Model

Constrained Local Models

Improving Portraits

Editing Video

Image and Search Windows    Optimization    Point Distribution Model

Scene Description

Illumination $\gamma$    Pose $(R, t)$    Fine Layer

Coarse Layer    Detail p

Shape $\alpha$    Albedo $\beta$    Expression $\delta$    Medium Layer

Correctives $\tau$

**Scene Description**

Illumination $\gamma$ — Pose $(R, t)$

**Fine Layer** — Detail p

**Coarse Layer**

Shape $\alpha$ — Albedo $\beta$ — Expression $\delta$

**Medium Layer** — Correctives $\tau$

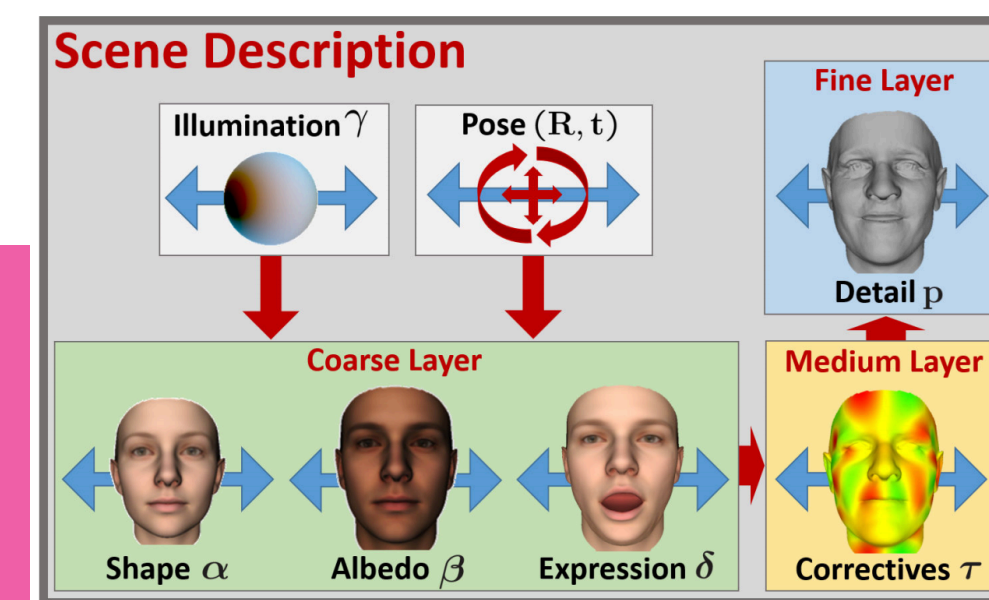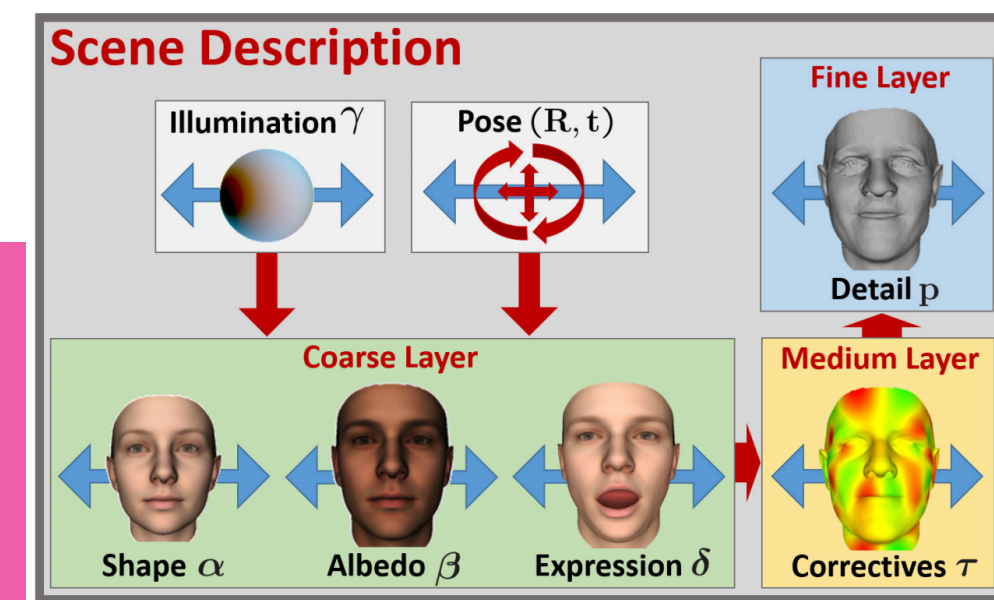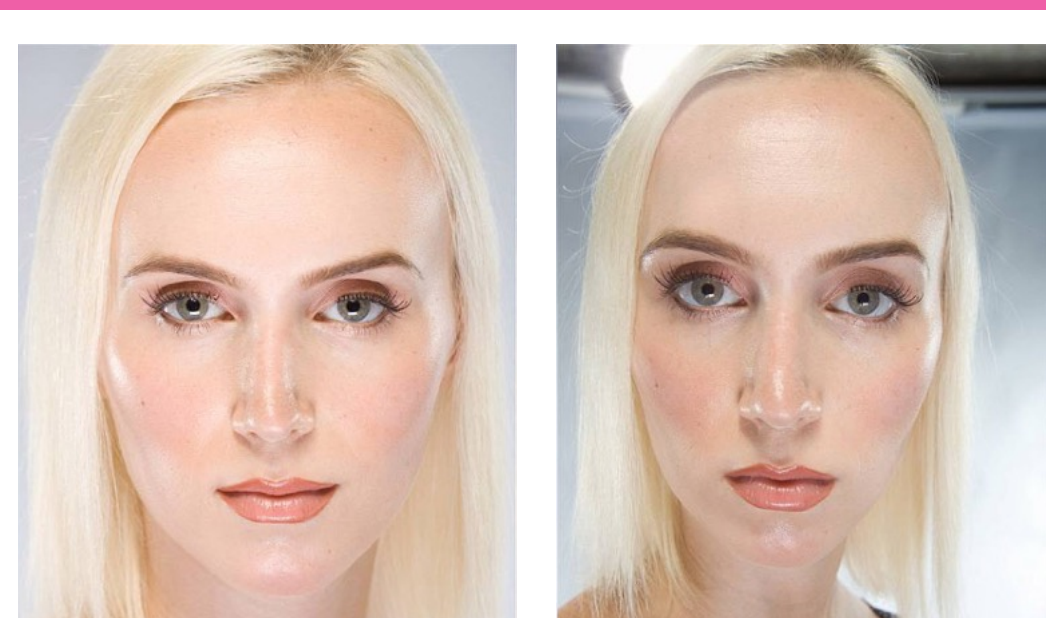# Constrained Local Models

# Improving Portraits

# Editing Video

Image and Search Windows — Optimization — Point Distribution Model

The quick brown ... — Input video

DH IY0 K W IH1 K sp B R AW1 N — Phoneme Alignment

Head Parameters — Model — Tracking & Reconstruction

Background Retiming — Rendering

The quick brown spider jumped fox — Edit Operation

The quick brown spider ... ... viper ... ox ... f ox — Viseme Search

brown — v — ox — jumped — Parameter Blending