

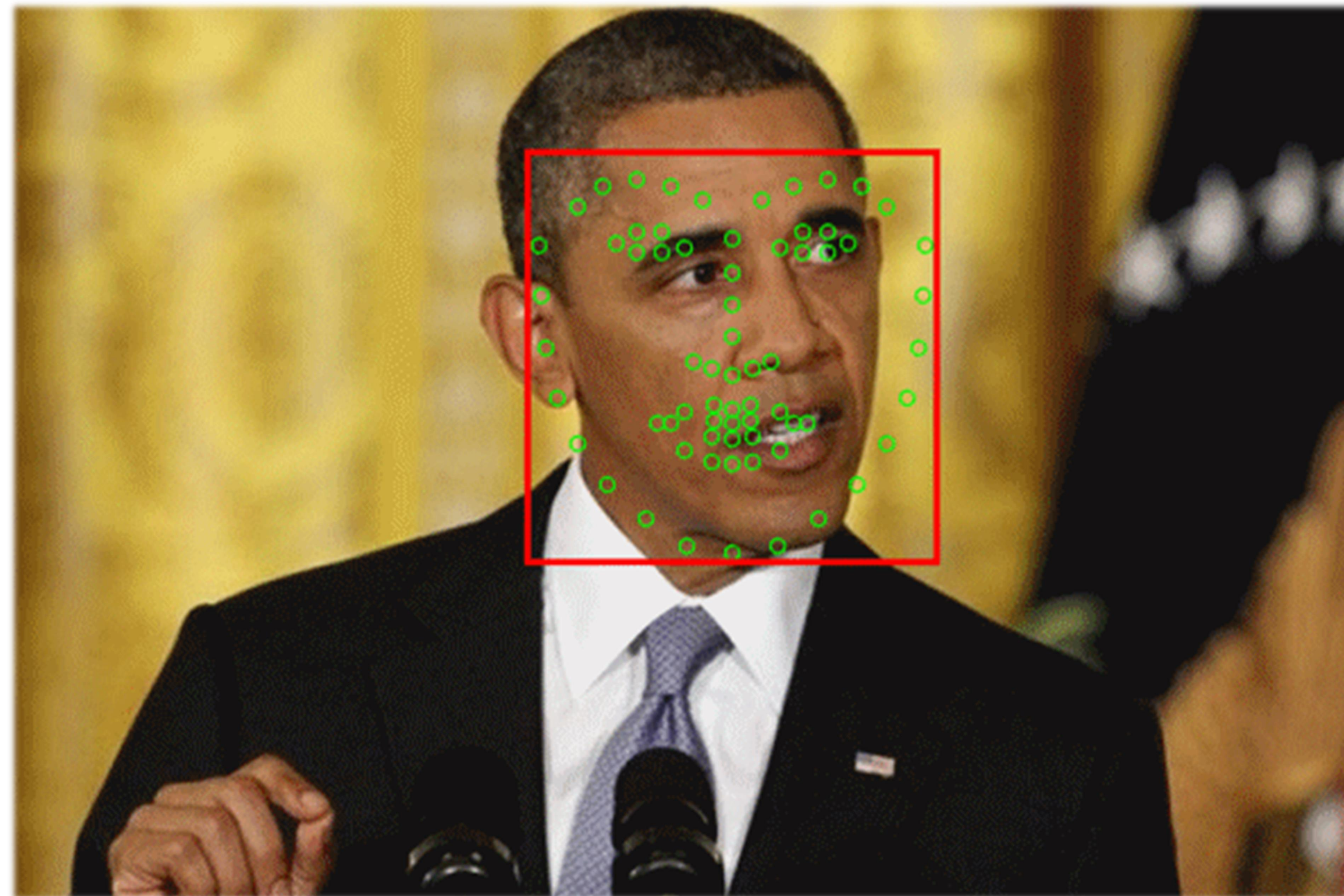
# Robust Facial Landmark Detection via Recurrent **A**ttentive-**R**efinement Networks

Shengtao XIAO, Jiashi FENG, Junliang XING, Hanjiang LAI,  
Shuicheng YAN, Ashraf KASSIM

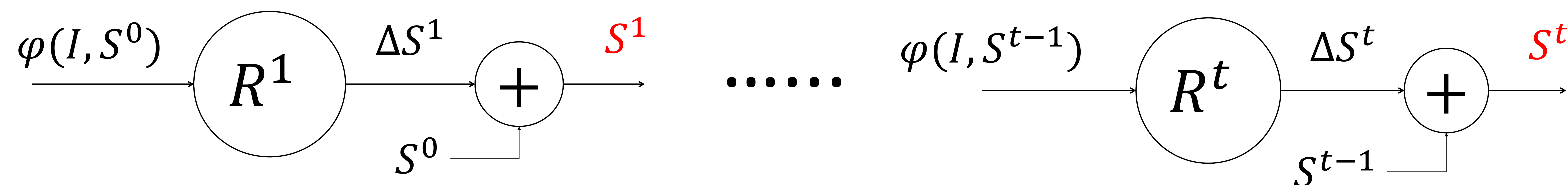


# Problem Introduction

- Obtain face shape by locating pre-defined facial landmarks.



- Challenges: face occlusions, pose variations, expressions, etc.
- Solutions: cascaded face shape regression

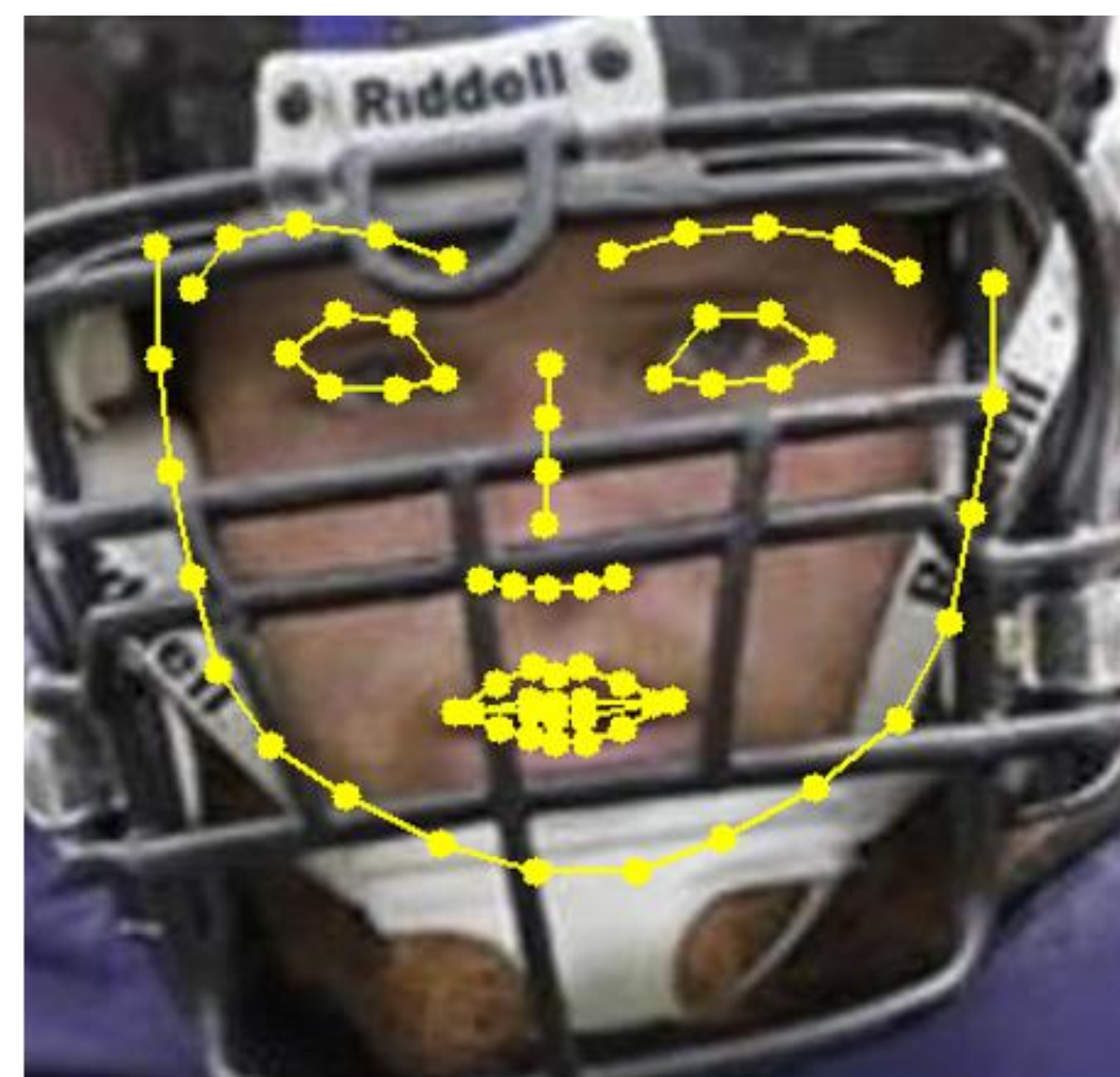


# Recurrent Attentive-Refinement (RAR) Network

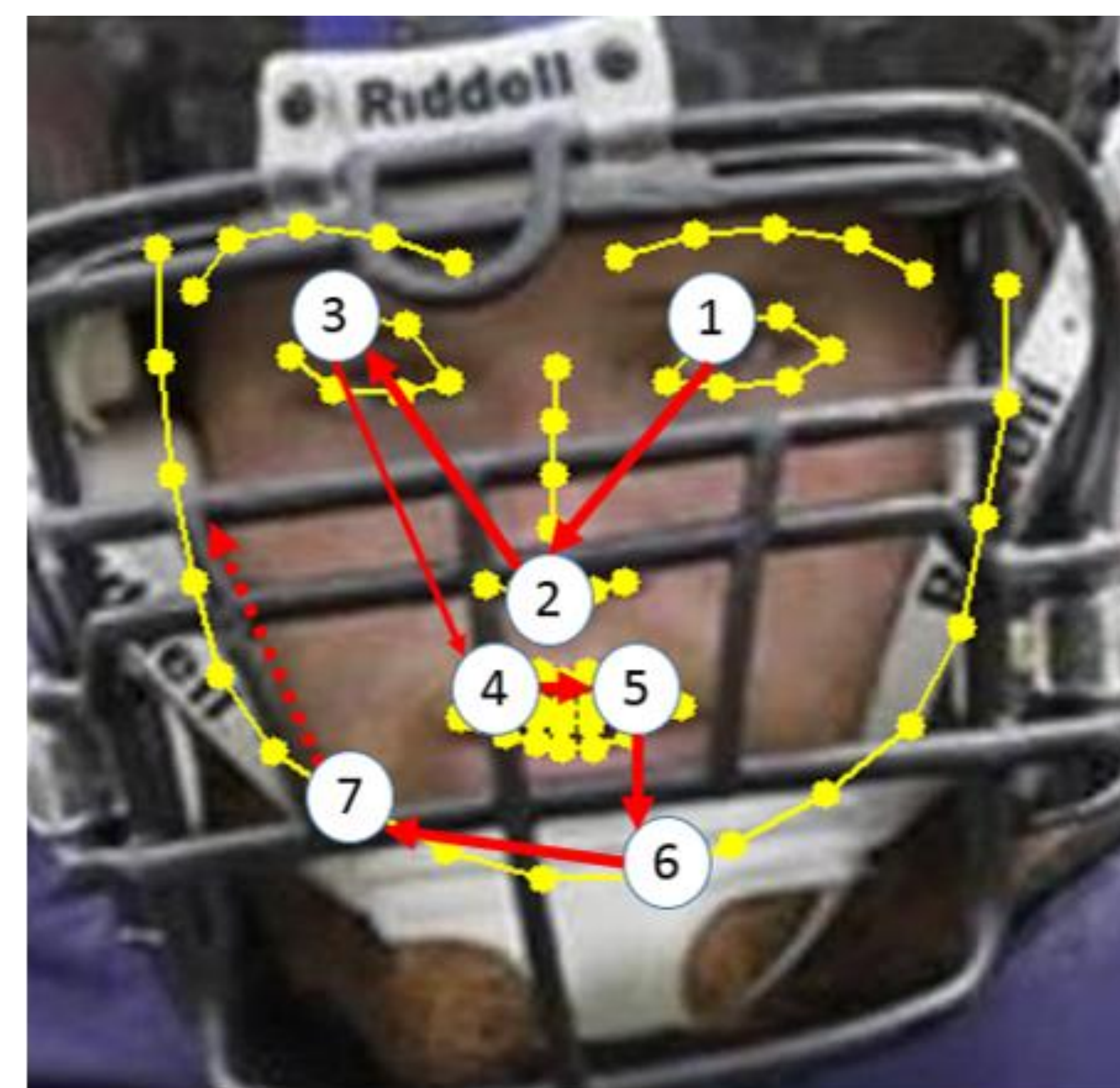
- Deep Feature Learning
- Robust Shape Initialization
- Recurrent-Attentive Refinement
  - Attention module
  - Refinement module



Input



Robust Initialization



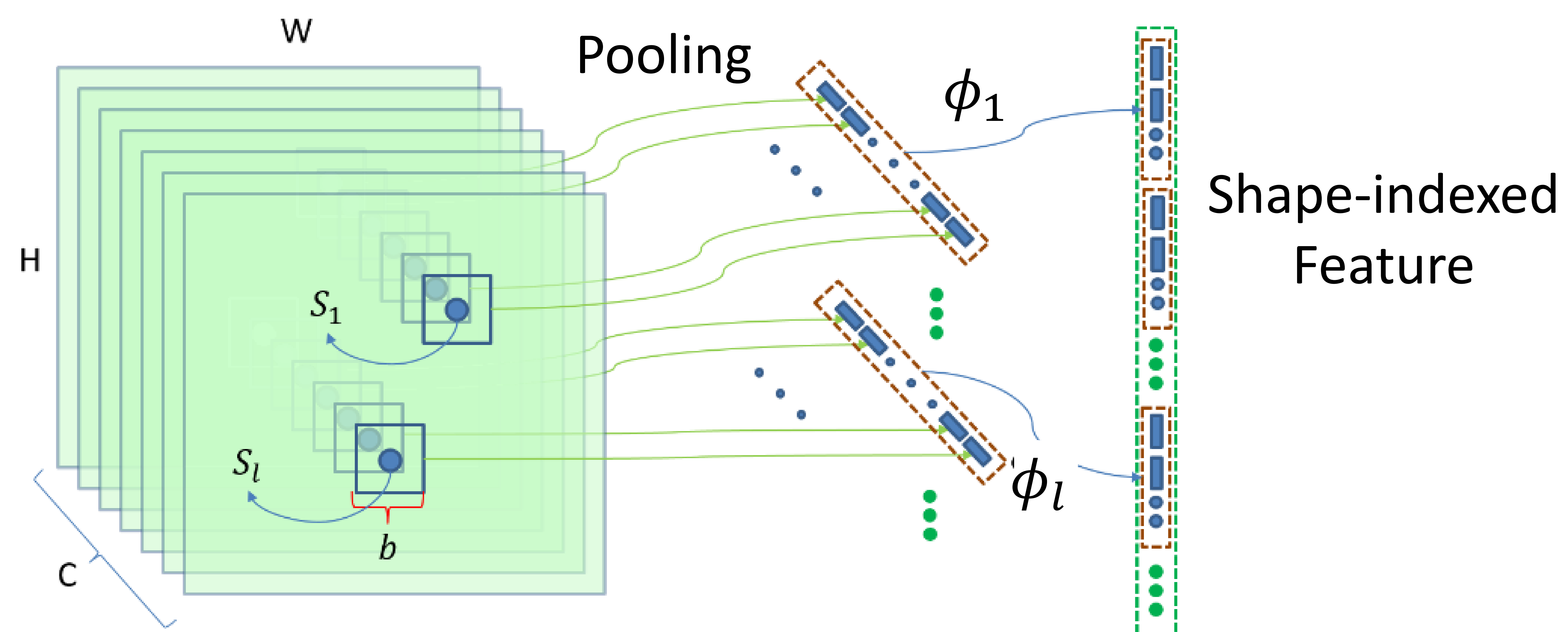
Attention-driven Refinement



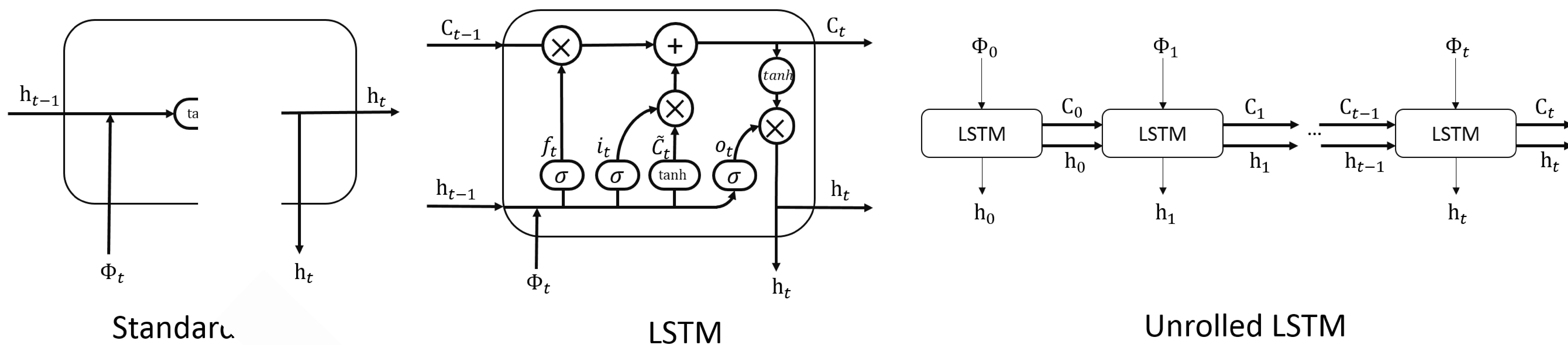
Results

# Background

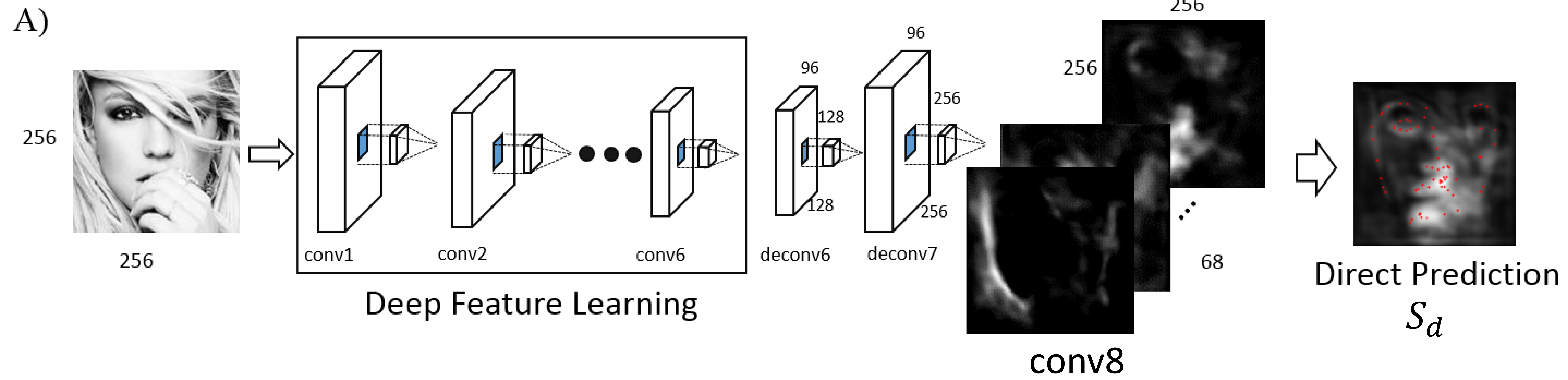
- CNN model with Shape-Indexed Pooling (SIP)



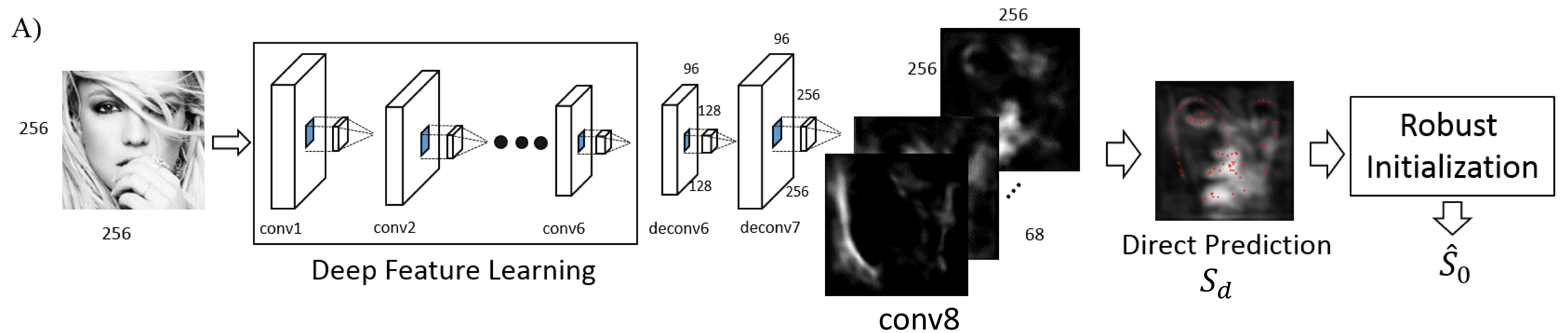
- RNN model and Long Short-Term Memory (LSTM)



# RAR Networks

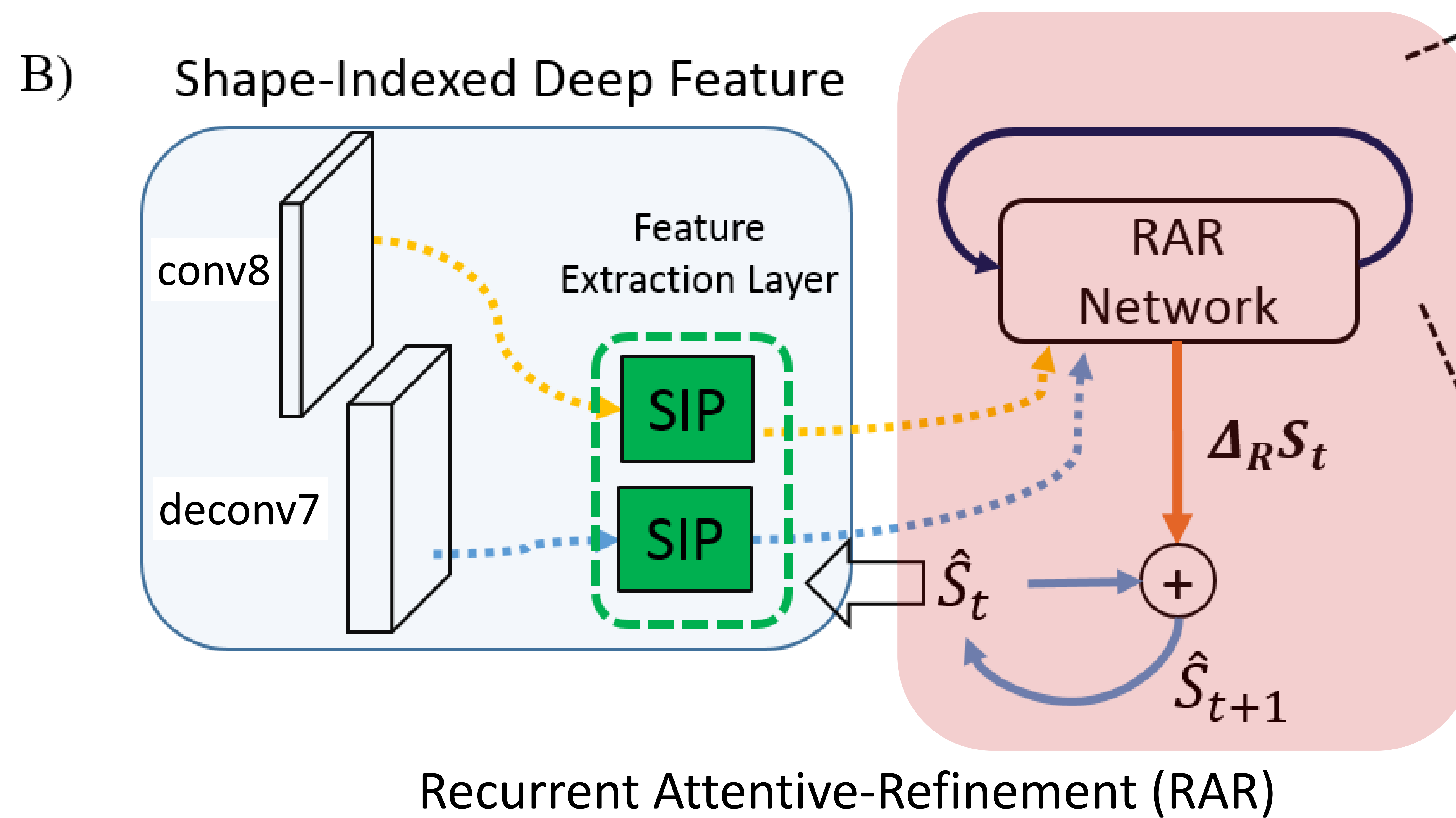
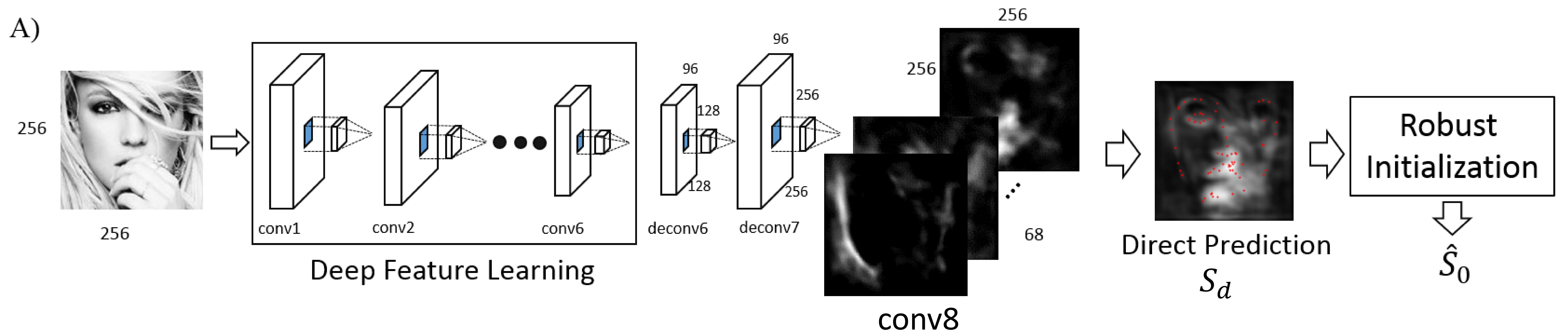


# RAR Networks



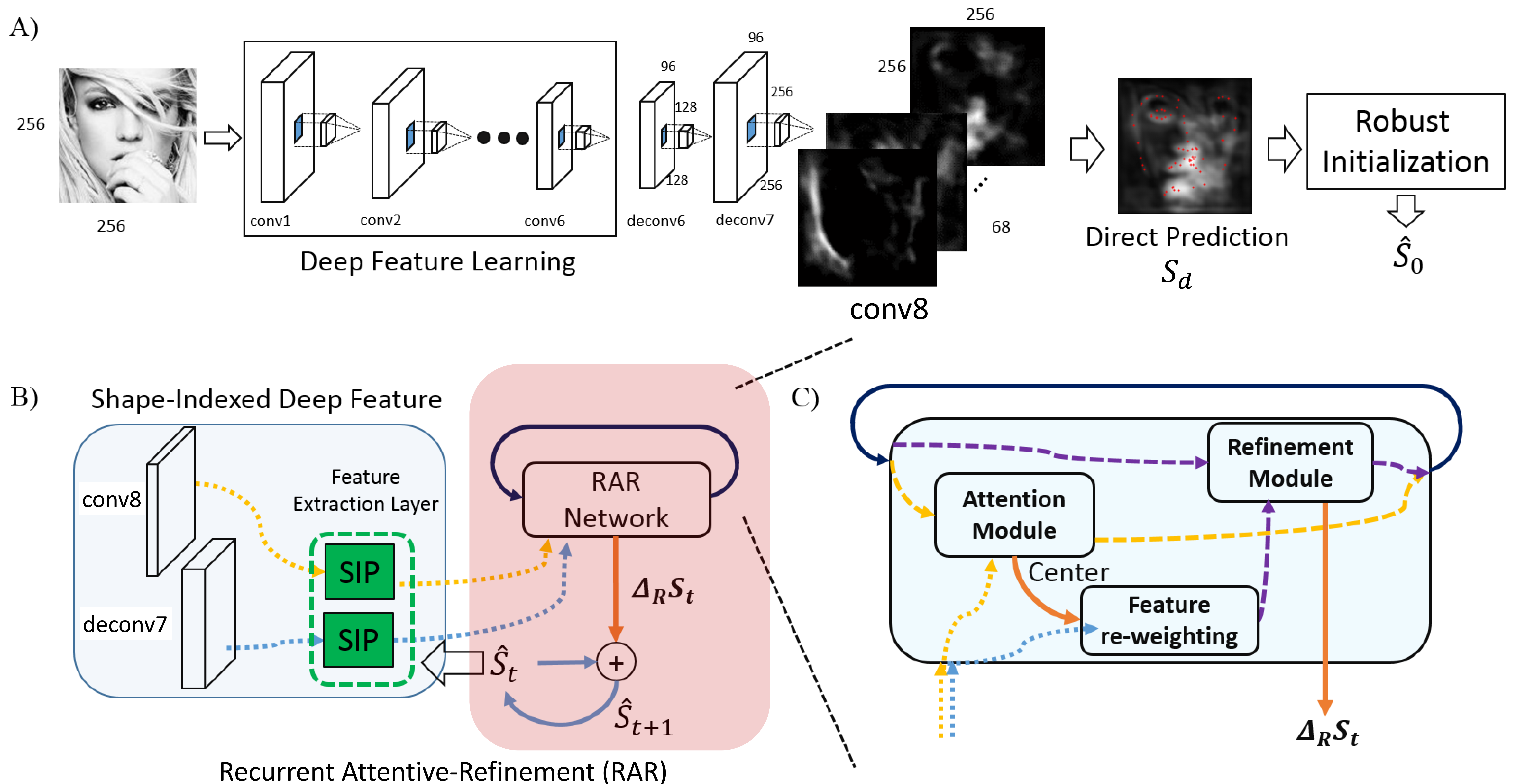
A). Deep **feature extraction**, landmark regression and robust **initialization**.

# RAR Networks



- A). Deep **feature extraction**, landmark regression and robust **initialization**.  
B). RAR **sequentially refines** the landmark estimation.

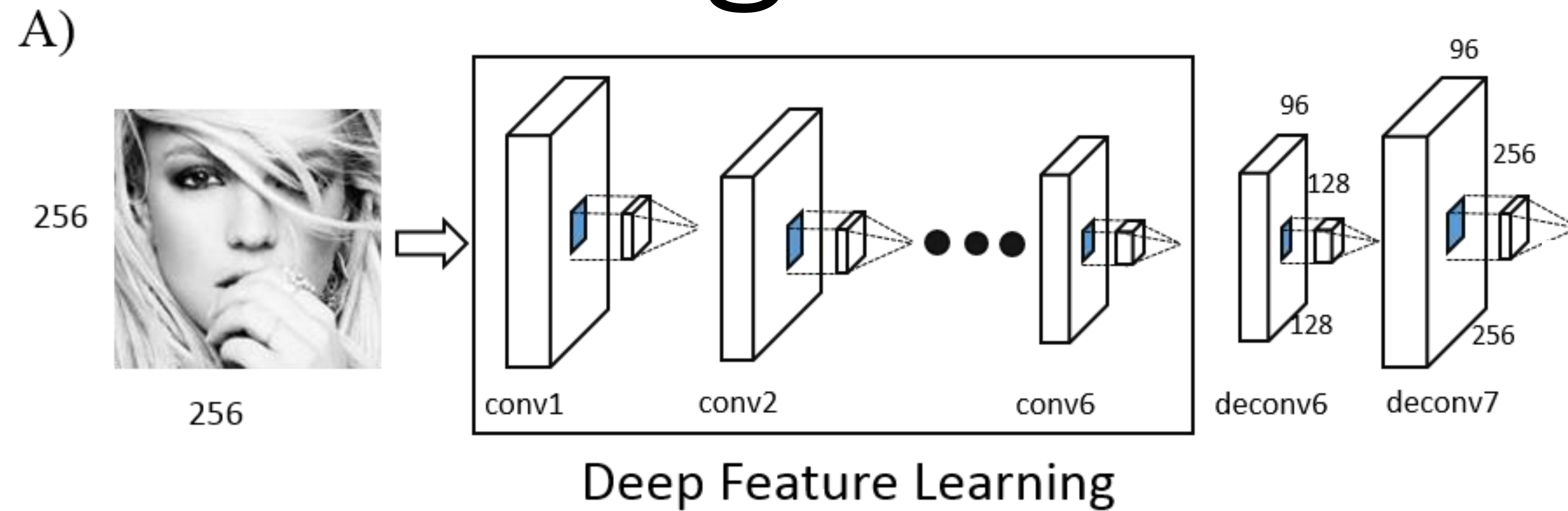
# RAR Networks



- A). Deep **feature extraction**, landmark regression and robust **initialization**.
- B). RAR **sequentially refines** the landmark estimation.
- C). An attention model in RAR for **adaptively selecting key landmark points**.

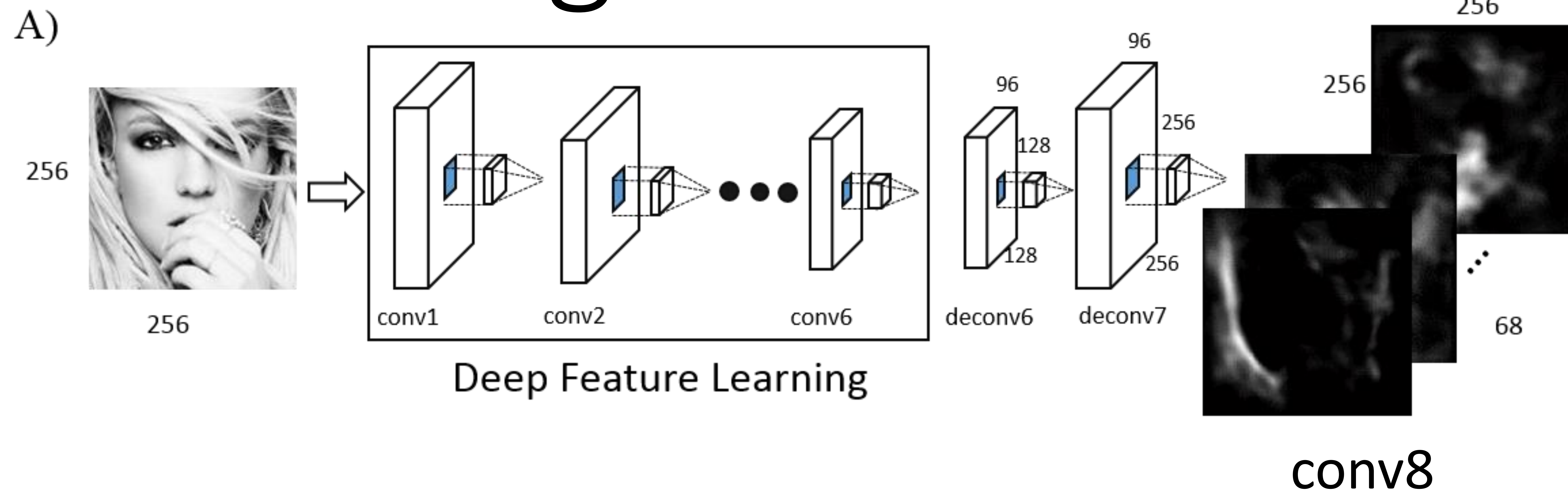


# RAR Networks: Deep Feature Learning



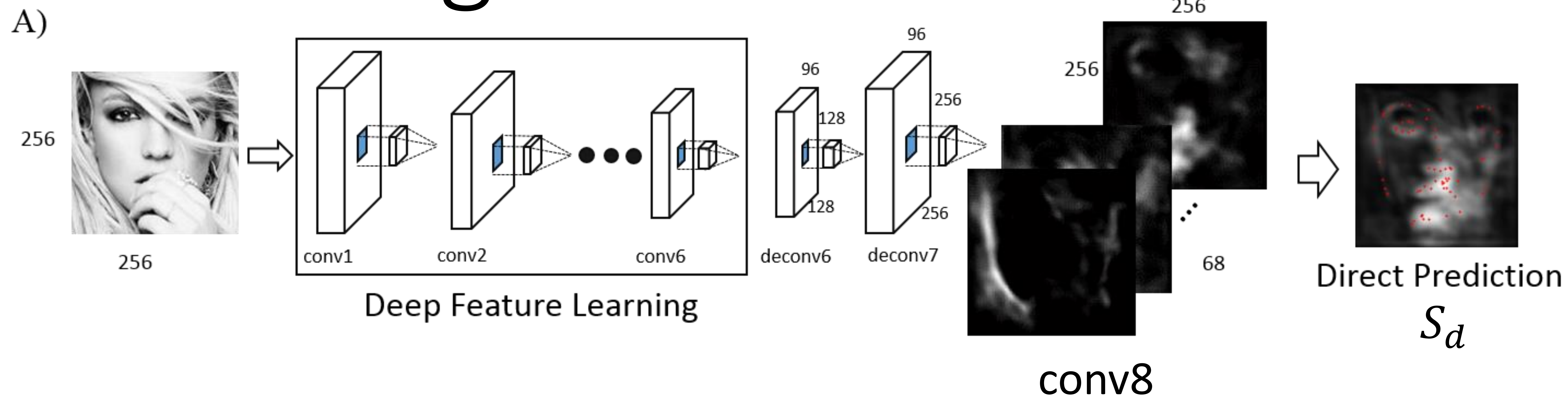
- Modified VGG19 Network + two Deconvolution layers to ensure pixel-to-pixel correspondence

# RAR Networks: Deep Feature Learning



- Modified VGG19 Network + two Deconvolution layers to ensure pixel-to-pixel correspondence
- SoftMax regression loss on **conv8**

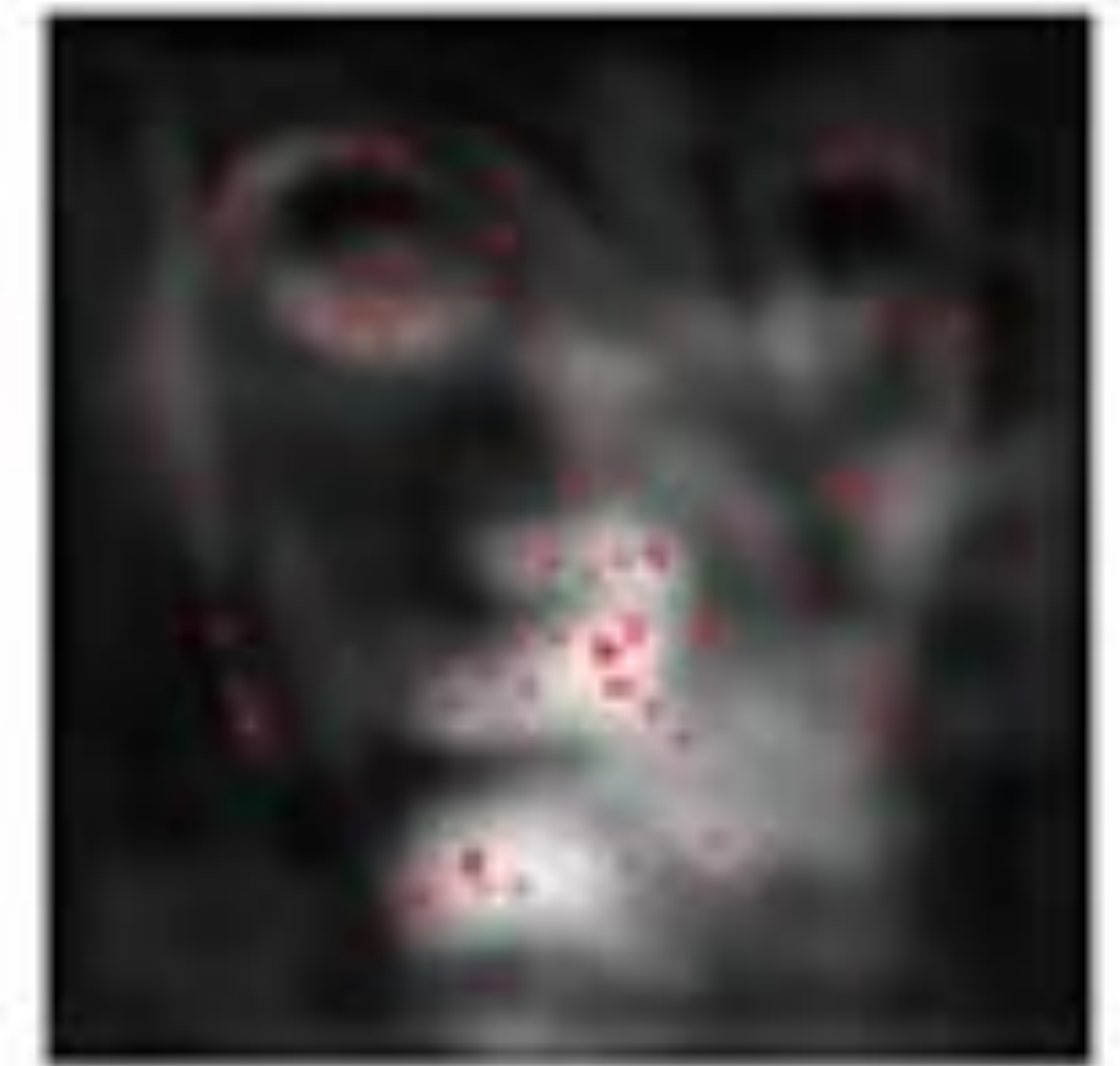
# RAR Networks: Deep Feature Learning



- Modified VGG19 Network + two Deconvolution lby selecting location of maximum response from  $v$ -th channel of **conv8**
- SoftMax regression loss on **conv8**
- Directly estimate landmark location  $S_d v d v d v d v d v$  by selecting location of maximum response from  $v$ -th channel of **conv8**

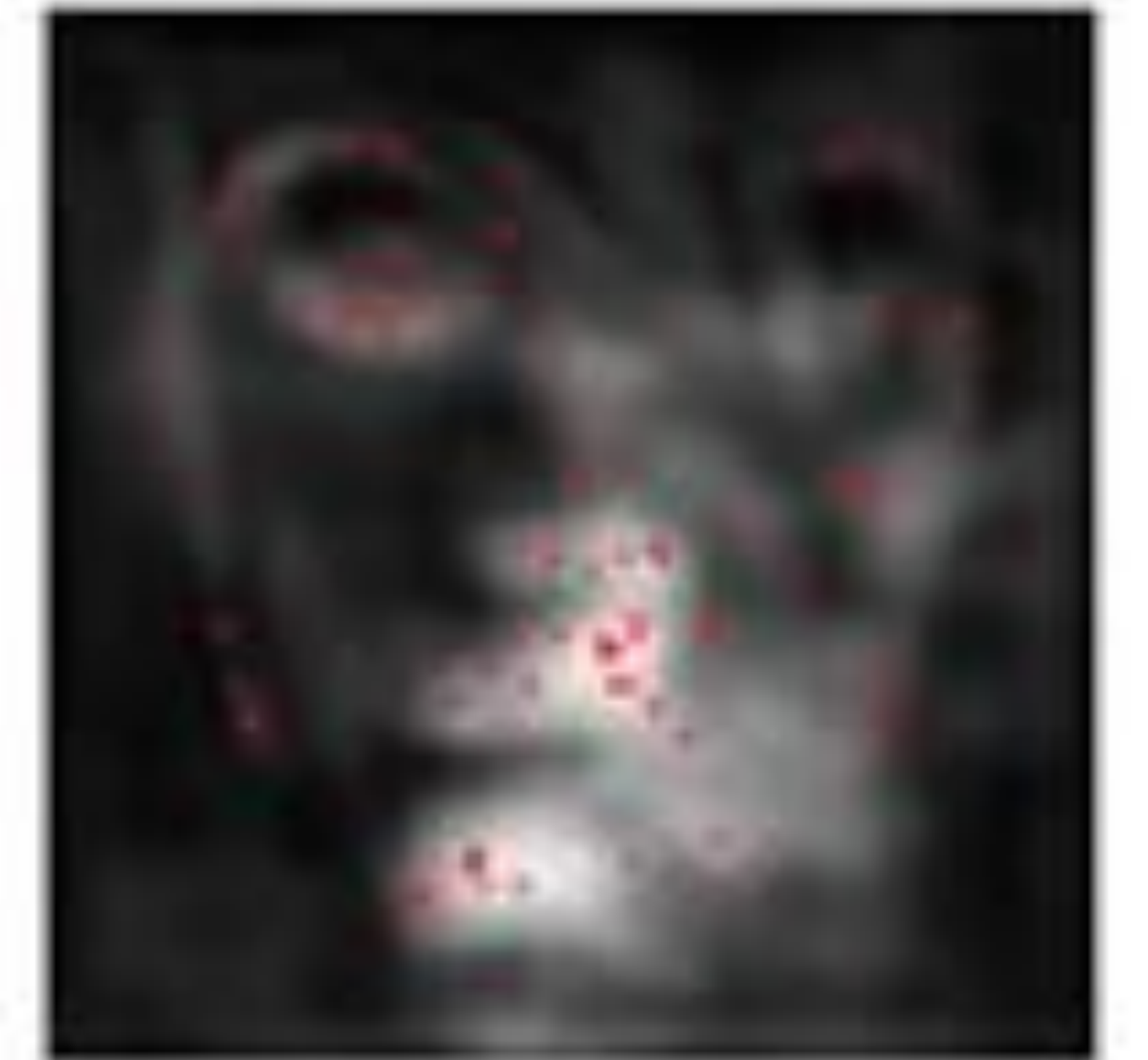
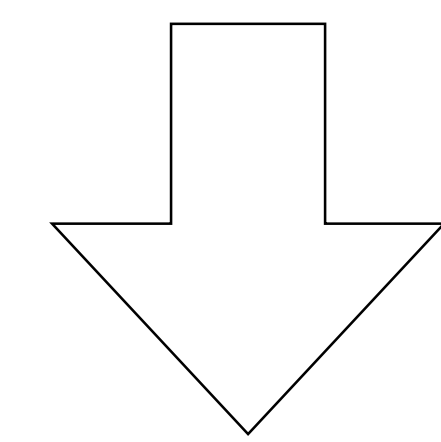
# RAR Networks: Robust Shape Initialization

However, detected shape is sensitive to occlusion



# RAR Networks: Robust Shape Initialization

However, detected shape is sensitive to occlusion



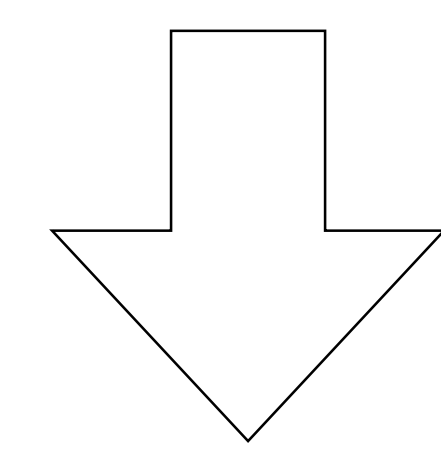
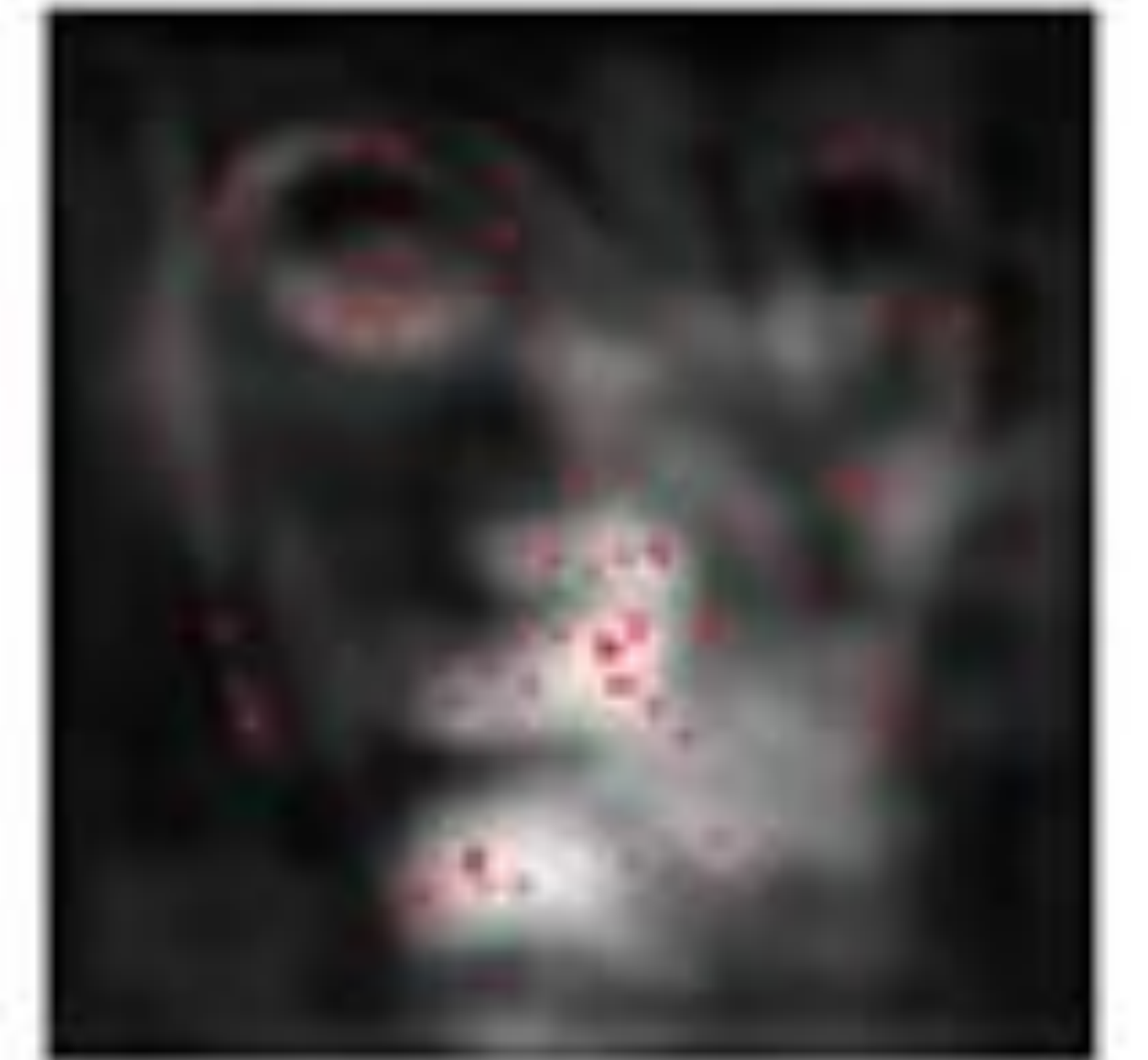
Robust Initial Shape Selection:

$$S_0 = \arg \min_S \|S - S_d\|, \text{ s. t. } S \in \mathcal{F}$$

$S_d$  → Detected Shape

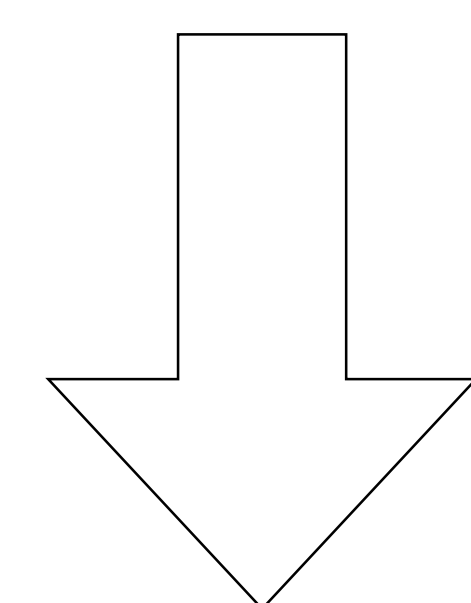
# RAR Networks: Robust Shape Initialization

However, detected shape is sensitive to occlusion



Robust Initial Shape Selection:

$$S_0 = \arg \min_S \|S - S_d\|, \text{ s. t. } S \in \mathcal{F}$$

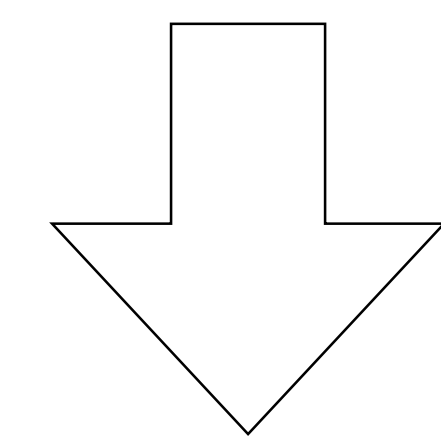
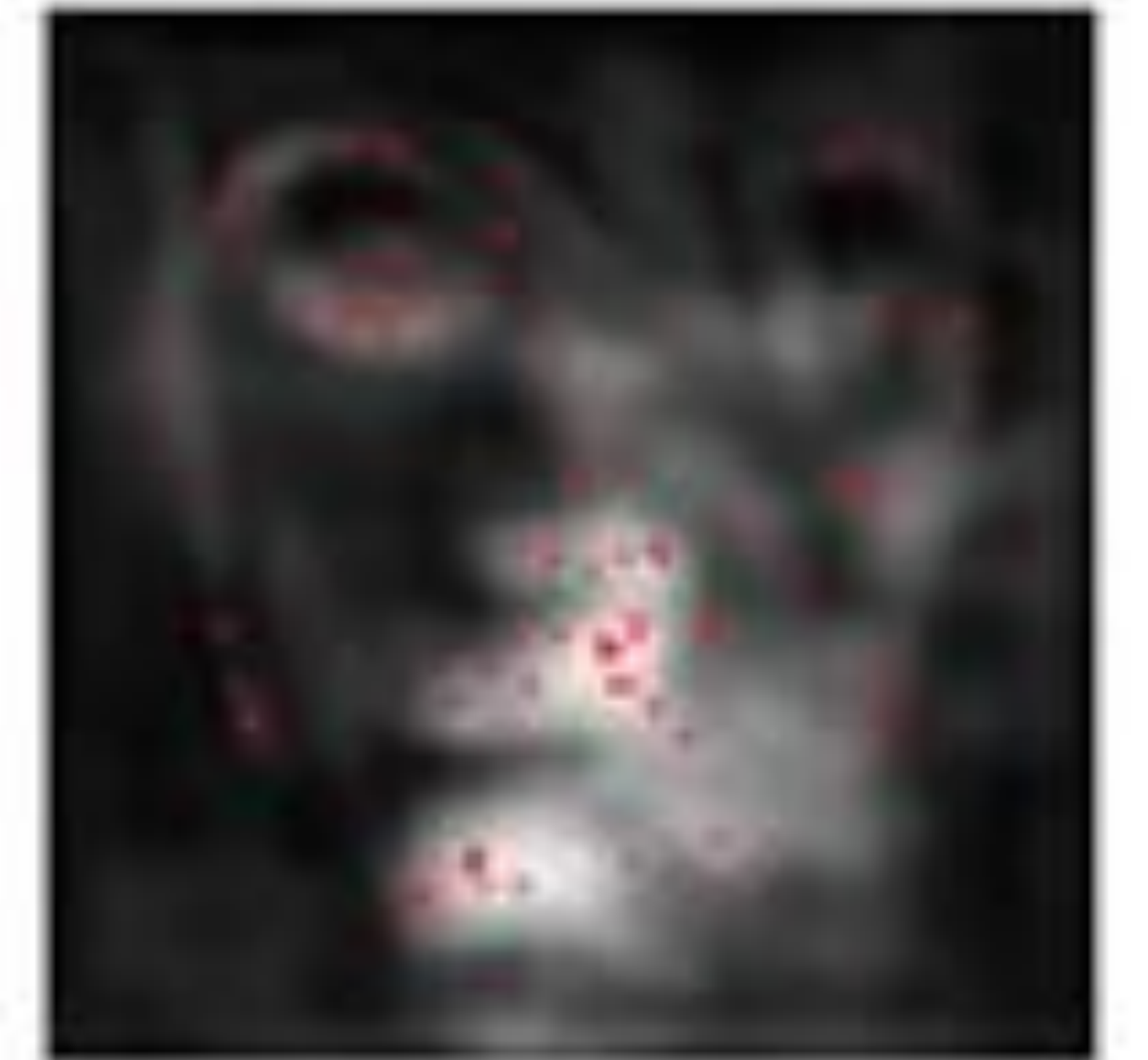


$$S_0 = \arg \min_{S, c \doteq [c_i]} \|S - S_d\|_0 + \lambda \|c\|_0, \text{ s. t. } S = \sum_i^m c_i S_i$$

GT Shapes

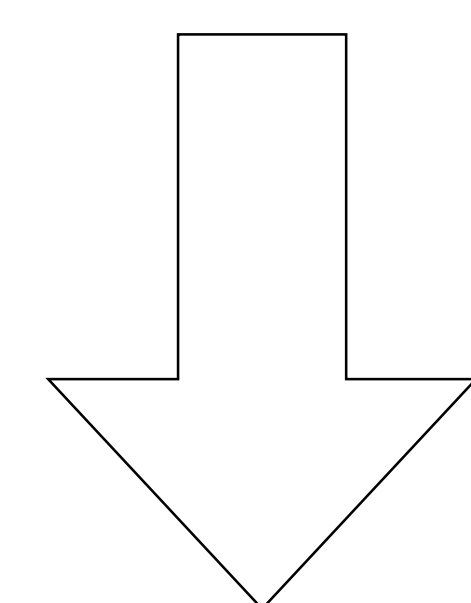
# RAR Networks: Robust Shape Initialization

However, detected shape is sensitive to occlusion



Robust Initial Shape Selection:

$$S_0 = \arg \min_S \|S - S_d\|, \text{ s. t. } S \in \mathcal{F}$$

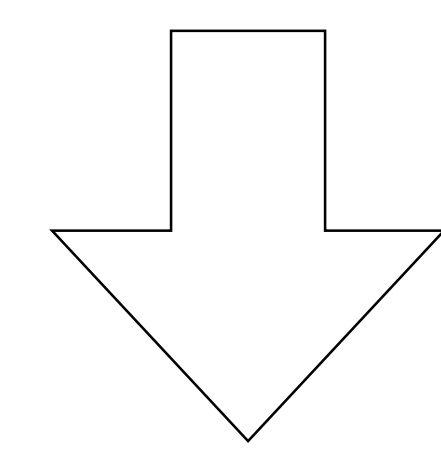
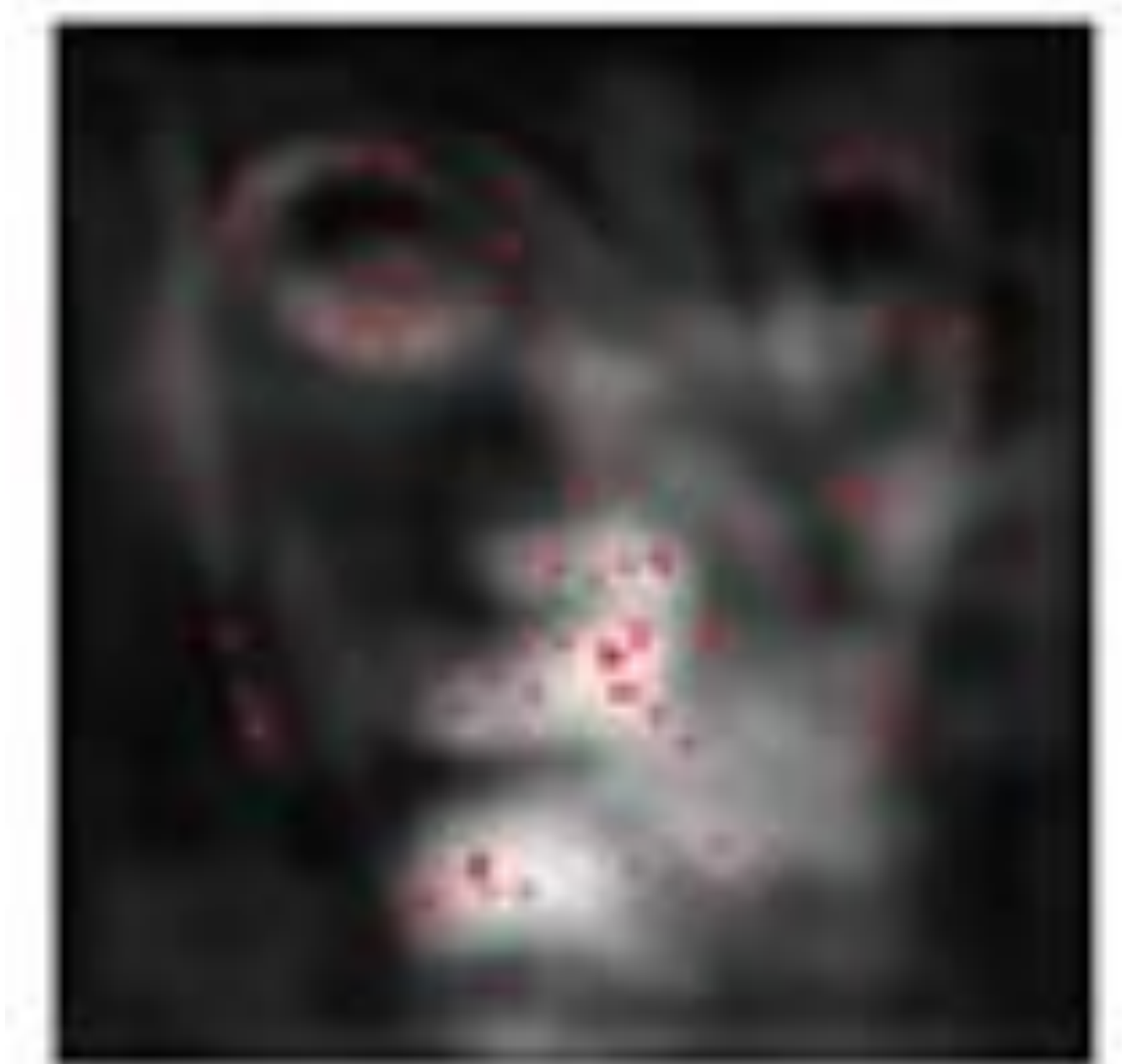


$$S_0 = \arg \min_{S, c \doteq [c_i]} \|S - S_d\|_0 + \lambda \|c\|_0, \text{ s. t. } S = \sum_i^m c_i S_i$$

Solve: get K representative shapes via K-means clustering + RANSAC method to filter out significant outliers

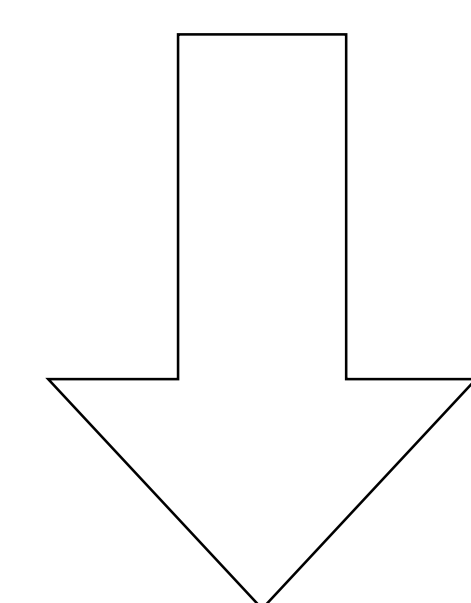
# RAR Networks: Robust Shape Initialization

However, detected shape is sensitive to occlusion

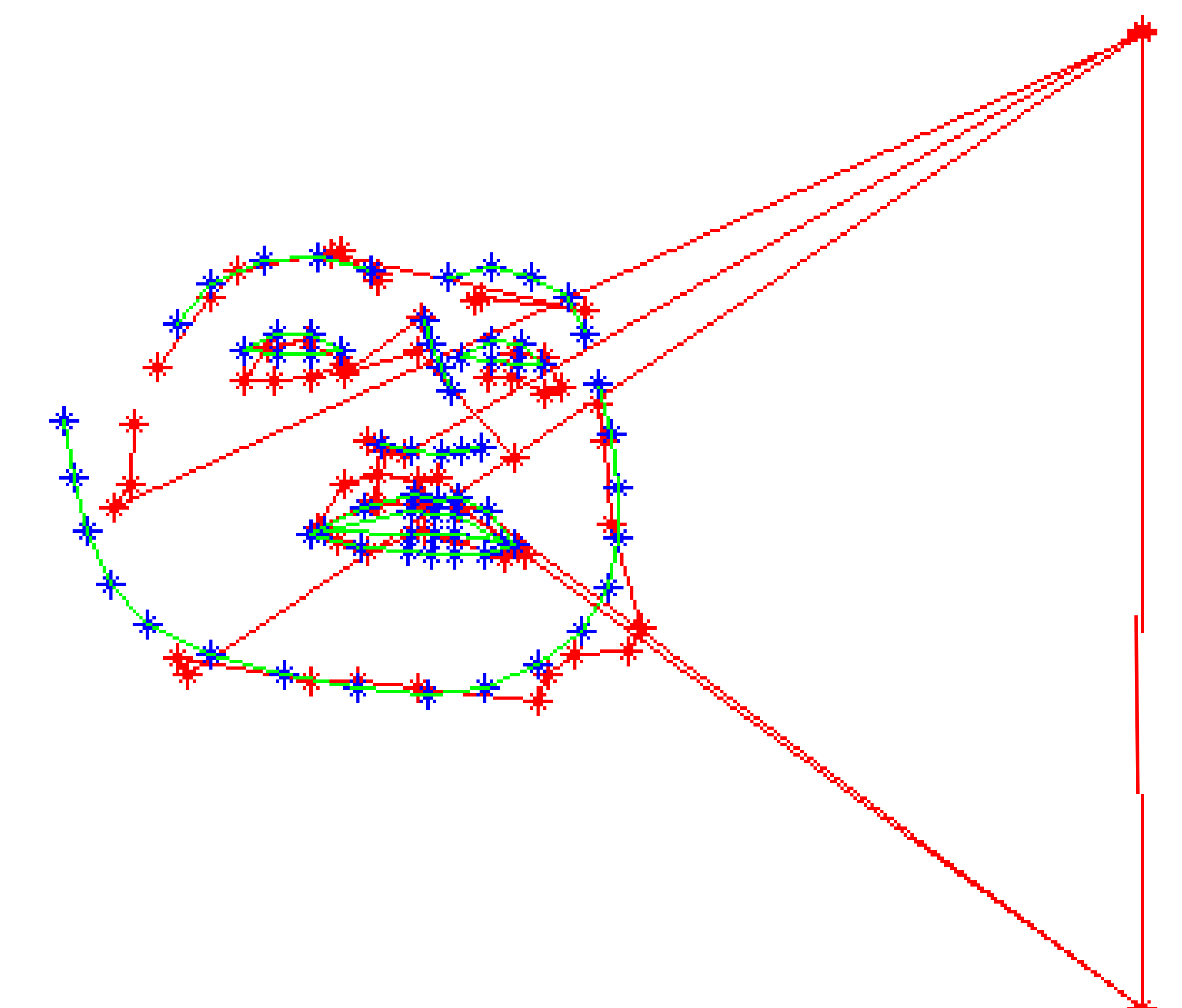
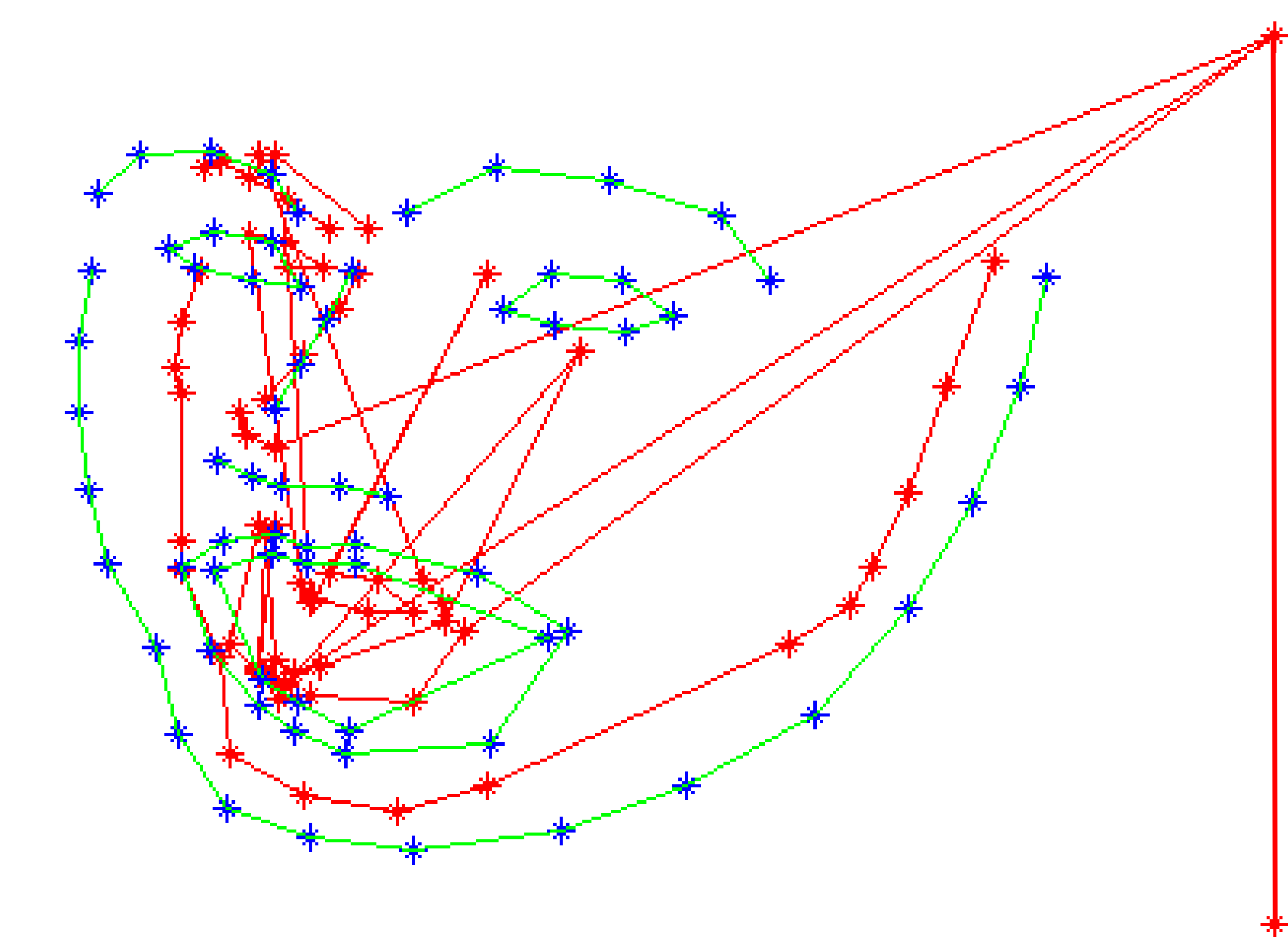


Robust Initial Shape Selection:

$$S_0 = \arg \min_S \|S - S_d\|, \text{ s. t. } S \in \mathcal{F}$$



$$S_0 = \arg \min_{S, c \doteq [c_i]} \|S - S_d\|_0 + \lambda \|c\|_0, \text{ s. t. } S = \sum_i^m c_i S_i$$



Solve: get K representative shapes via K-means clustering + RANSAC method to filter out significant outliers



# RAR Networks: Attention Module

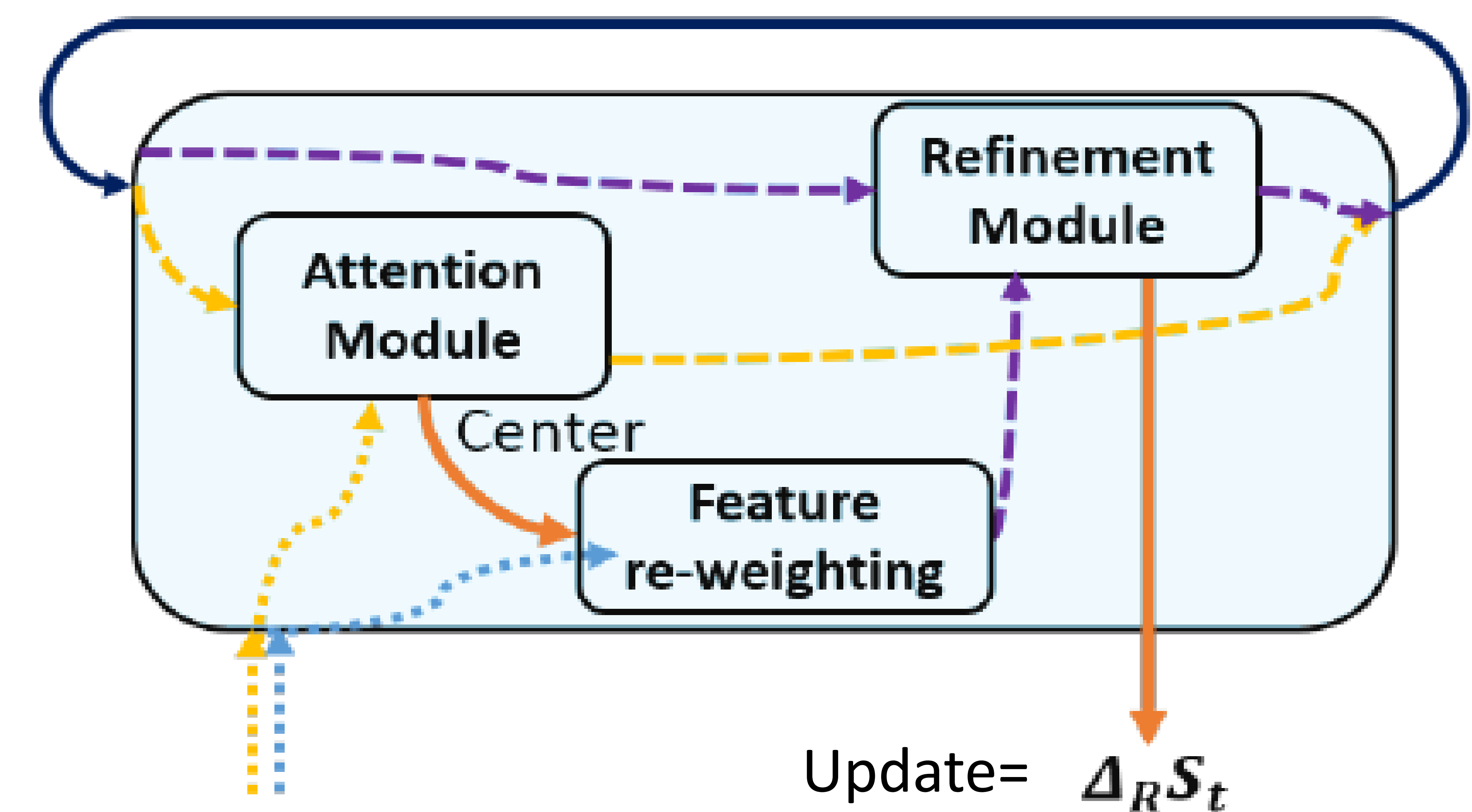
-

# RAR Networks: Attention Module

- A-LSTM (attention module) selects attention center with top confidence at each recurrent stage

$$C^* = \operatorname{argmax}_{c \in \{1, \dots, L\}} \text{A-LSTM}(\Phi_a(I_t, \hat{S}_t); W_a, c)$$

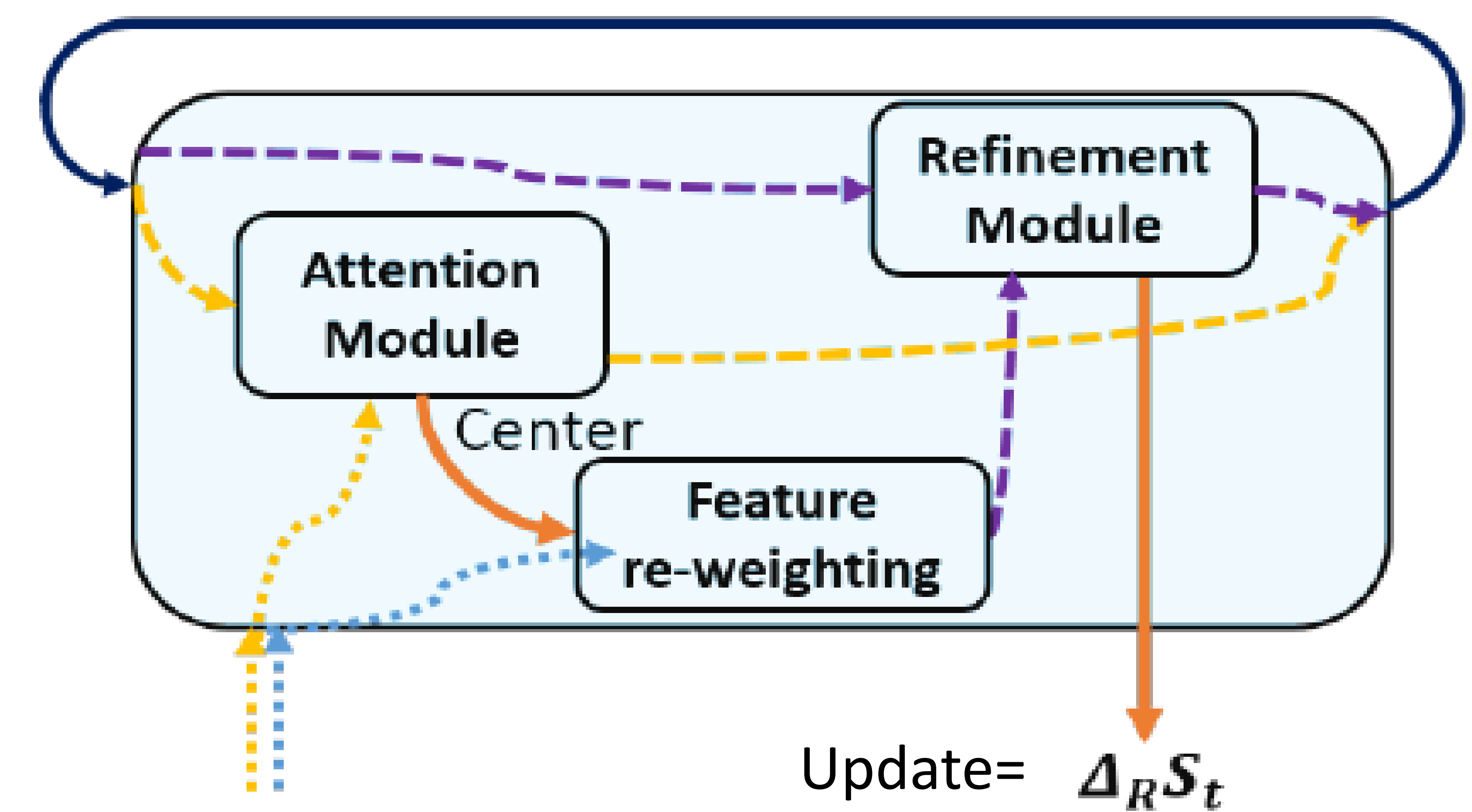
Deep SIP  
Feature



# RAR Networks: Attention Module

- $\mathcal{R} a a \mathcal{R} a$
- A-LSTM (attention module) selects attention center with top confidence at each recurrent stage  

$$c = \operatorname{argmax}_{c \in \{1, \dots, L\}} A\text{-LSTM}(\Phi_a(I_t, S_t), W_a, c)$$



- A typical attention center is selected based on maximize reward  $\mathcal{R} a$

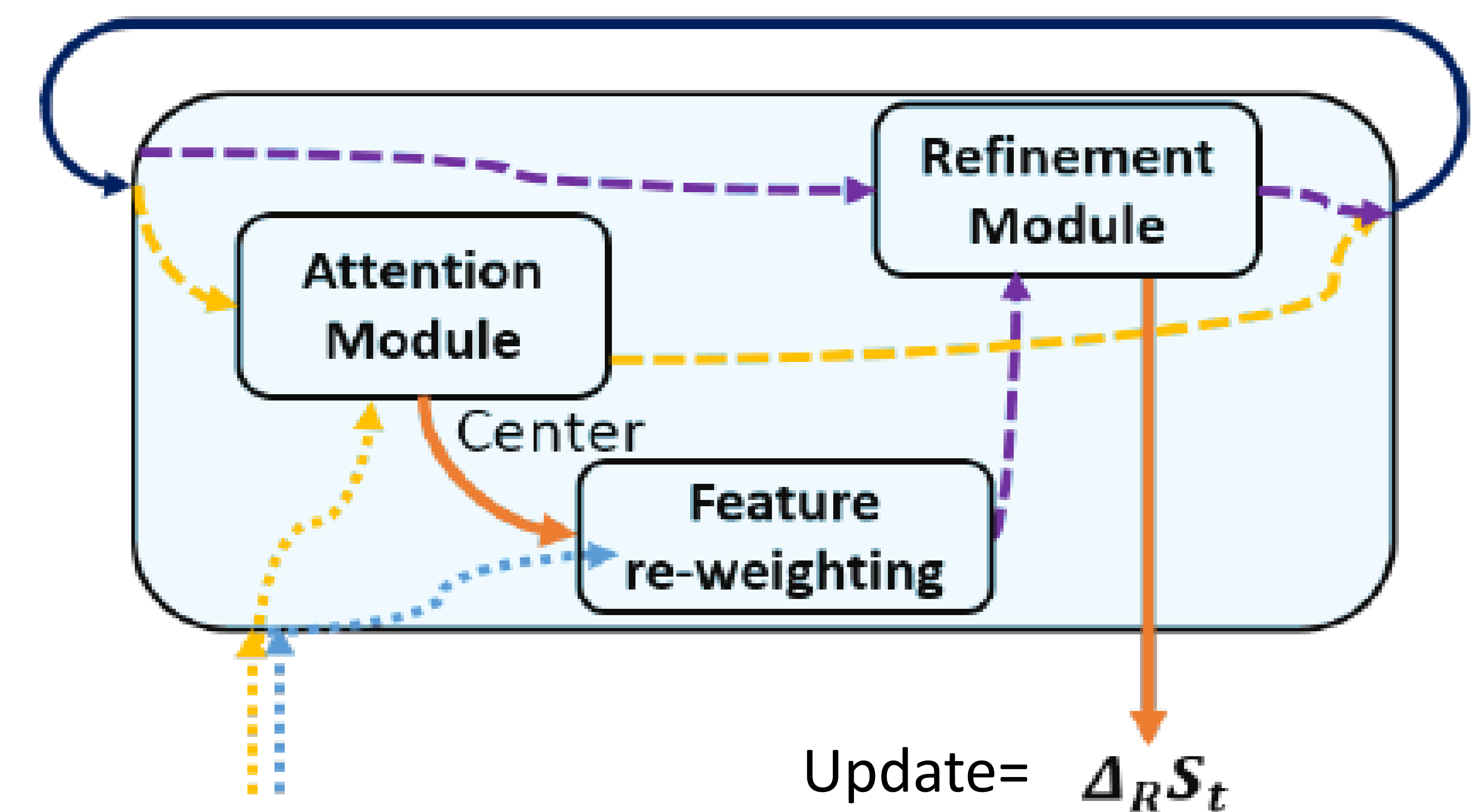
$$\mathcal{R}_a = \sum_{t=1}^{\infty} \eta^{t-1} R(\hat{S}_{t-1}, \hat{S}_t)$$

Discount Factor  $\eta$  (indicated by a red arrow pointing to  $\eta^{t-1}$ )

Intermediate Reward  $R(\hat{S}_{t-1}, \hat{S}_t)$  (indicated by a blue arrow pointing to the reward function)

# RAR Networks: Attention Module

- $\mathcal{R} a a a \mathcal{R} a$
- A-LSTM (attention module) selects attention center with top confidence at each recurrent stage
 
$$c = \operatorname{argmax}_{c \in \{1, \dots, L\}} A\text{-LSTM}(\Phi_a(I_t, S_t), W_a, c)$$



- A typical attention center is selected based on maximize reward  $\mathcal{R} a$

$$\mathcal{R}_a = \sum_{t=1}^{\infty} \eta^{t-1} R(\hat{S}_{t-1}, \hat{S}_t)$$

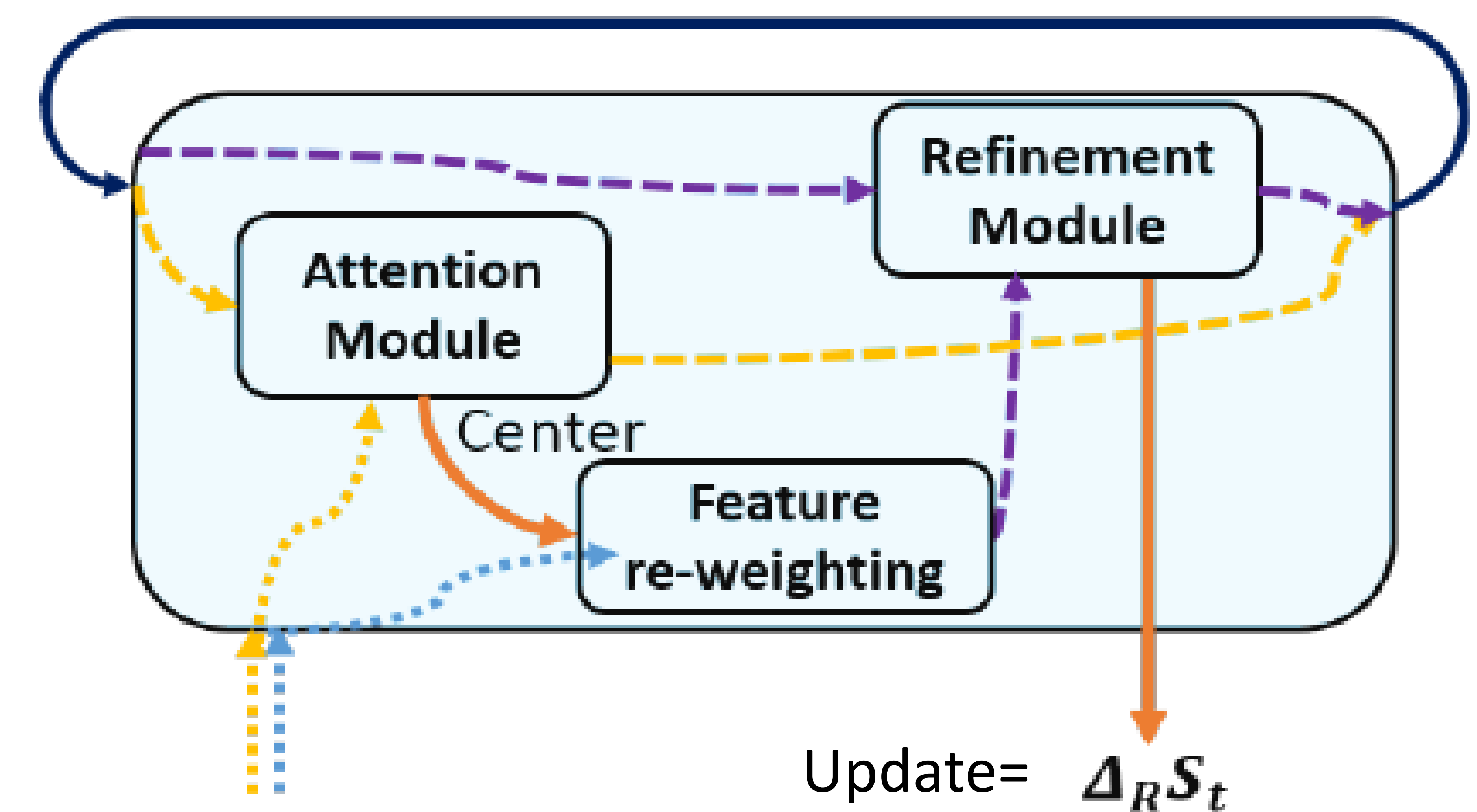
$$R(\hat{S}_{t-1}, \hat{S}_t) = \|\Gamma_t \Delta S_{t-1}\|_2^2 - \|\Gamma_t \Delta S_t\|_2^2$$

Weighting Factor (pointing to  $\Gamma_t$ )  
Ground-Truth Residue (pointing to  $\Delta S_t$ )

# RAR Networks: Attention Module

- $\mathcal{R} a a \mathcal{R} a$
- A-LSTM (attention module) selects attention center with top confidence at each recurrent stage  

$$c = \operatorname{argmax}_{c \in \{1, \dots, L\}} A\text{-LSTM}(\Phi_a(I_t, S_t), W_a, c)$$



- A typical attention center is selected based on maximize reward  $\mathcal{R} a$

$$\mathcal{R} a = \sum_{t=1}^{\infty} \eta^{t-1} R(\hat{S}_{t-1}, \hat{S}_t)$$

$$R(\hat{S}_{t-1}, \hat{S}_t) = \|\Gamma_t \Delta S_{t-1}\|_2^2 - \|\Gamma_t \Delta S_t\|_2^2$$

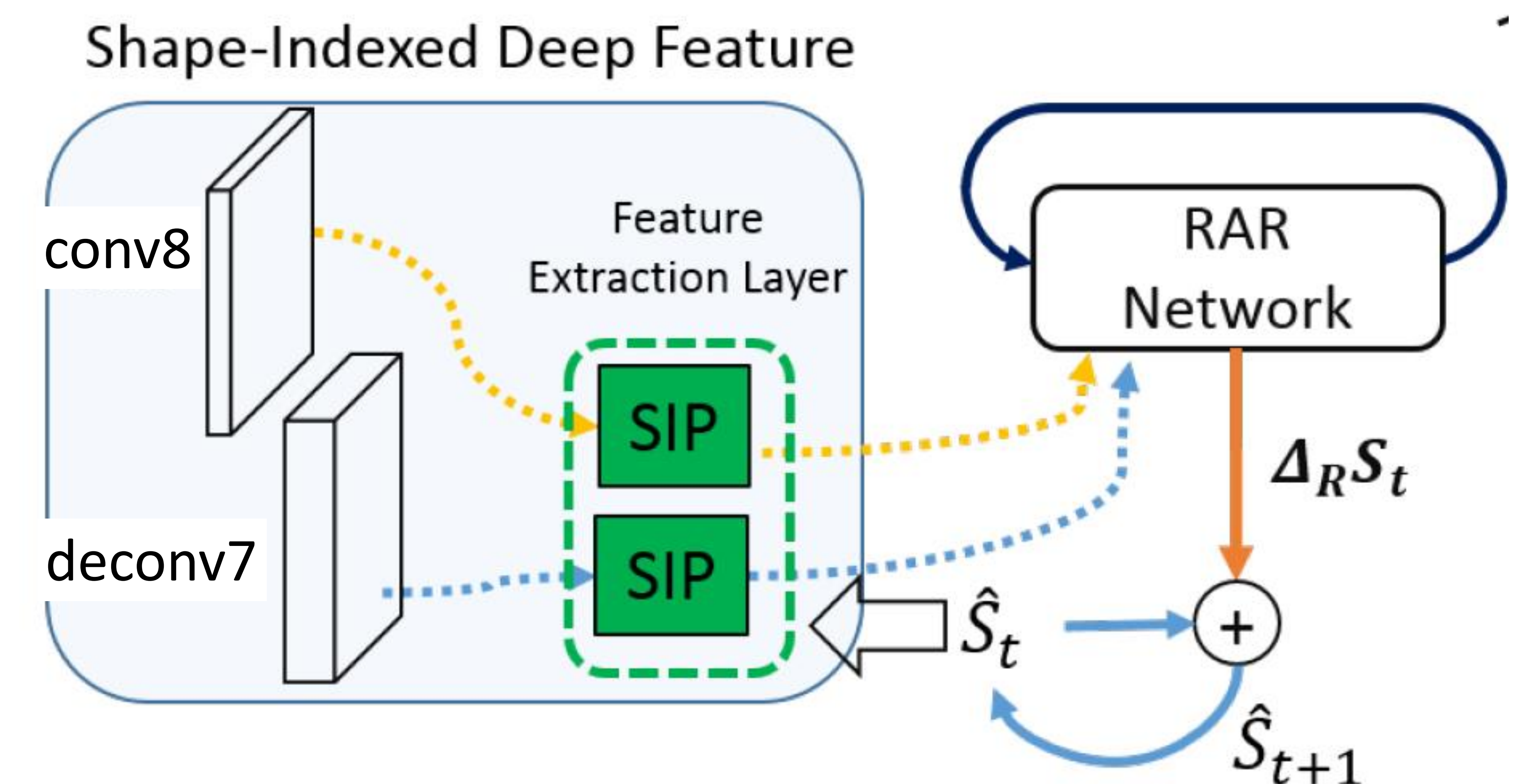
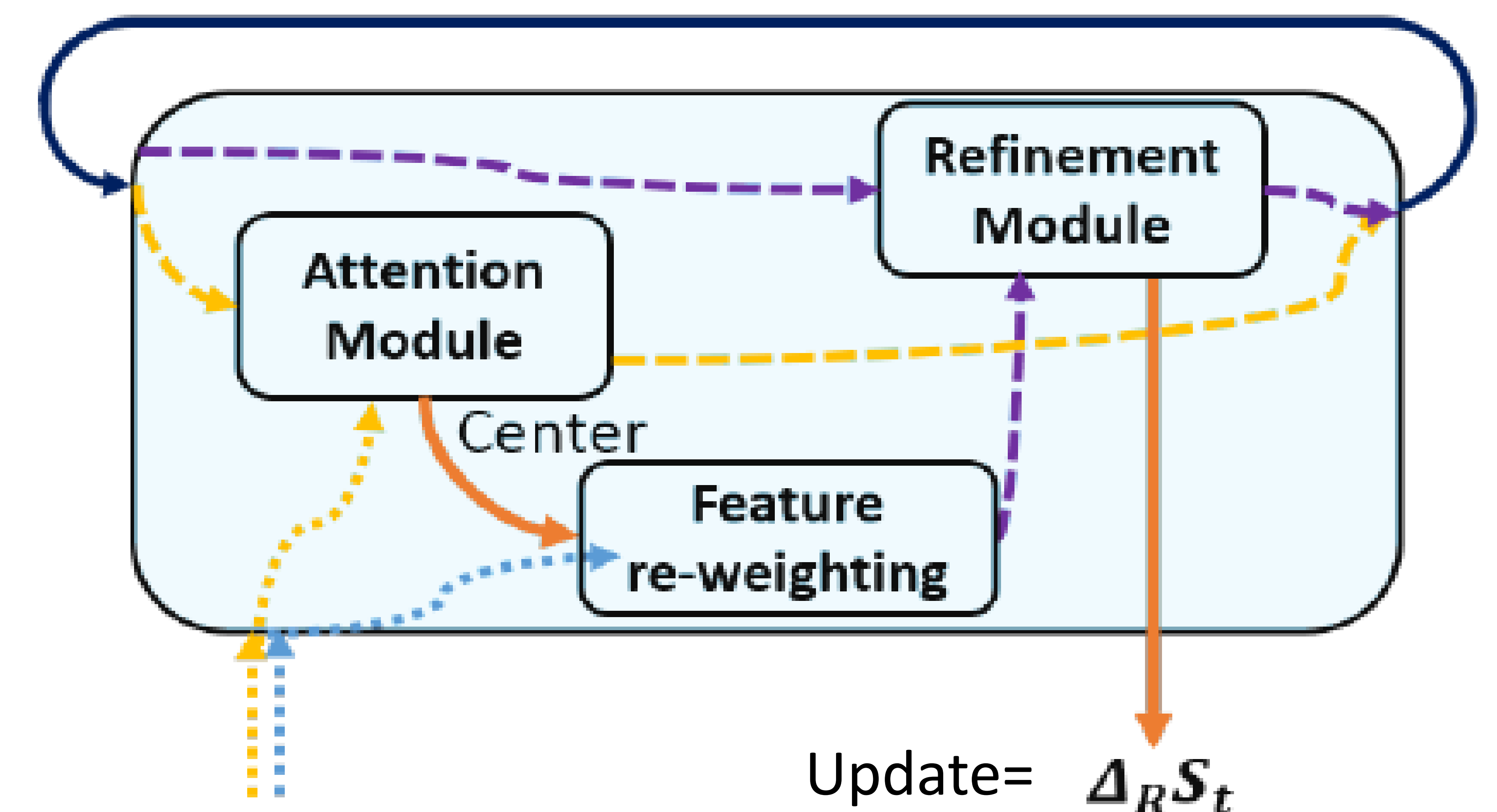
$$\Gamma_t = [\gamma_t^1, \gamma_t^2, \dots, \gamma_t^L], \text{ with } \gamma_t^l = \kappa \exp\left(\frac{-\|\hat{S}_t^l - \hat{S}_t^{c^*}\|_{l_2}^2}{4D_t^2}\right)$$

# RAR Networks: Refinement Module

- Feature re-weighting based on distance to attention center:

$$\Phi_r(I_t, \hat{S}_t) = [\gamma_t^1 \phi_t^1, \gamma_t^2 \phi_t^2, \dots, \gamma_t^L \phi_t^L]$$

Weighting Factor      Deep SIP Feature at  $S_t^2$



# RAR Networks: Refinement Module

- Feature re-weighting based on distance to attention center:

$$\Phi_r(I_t, \hat{S}_t) = [\gamma_t^1 \phi_t^1, \gamma_t^2 \phi_t^2, \dots, \gamma_t^L \phi_t^L]$$

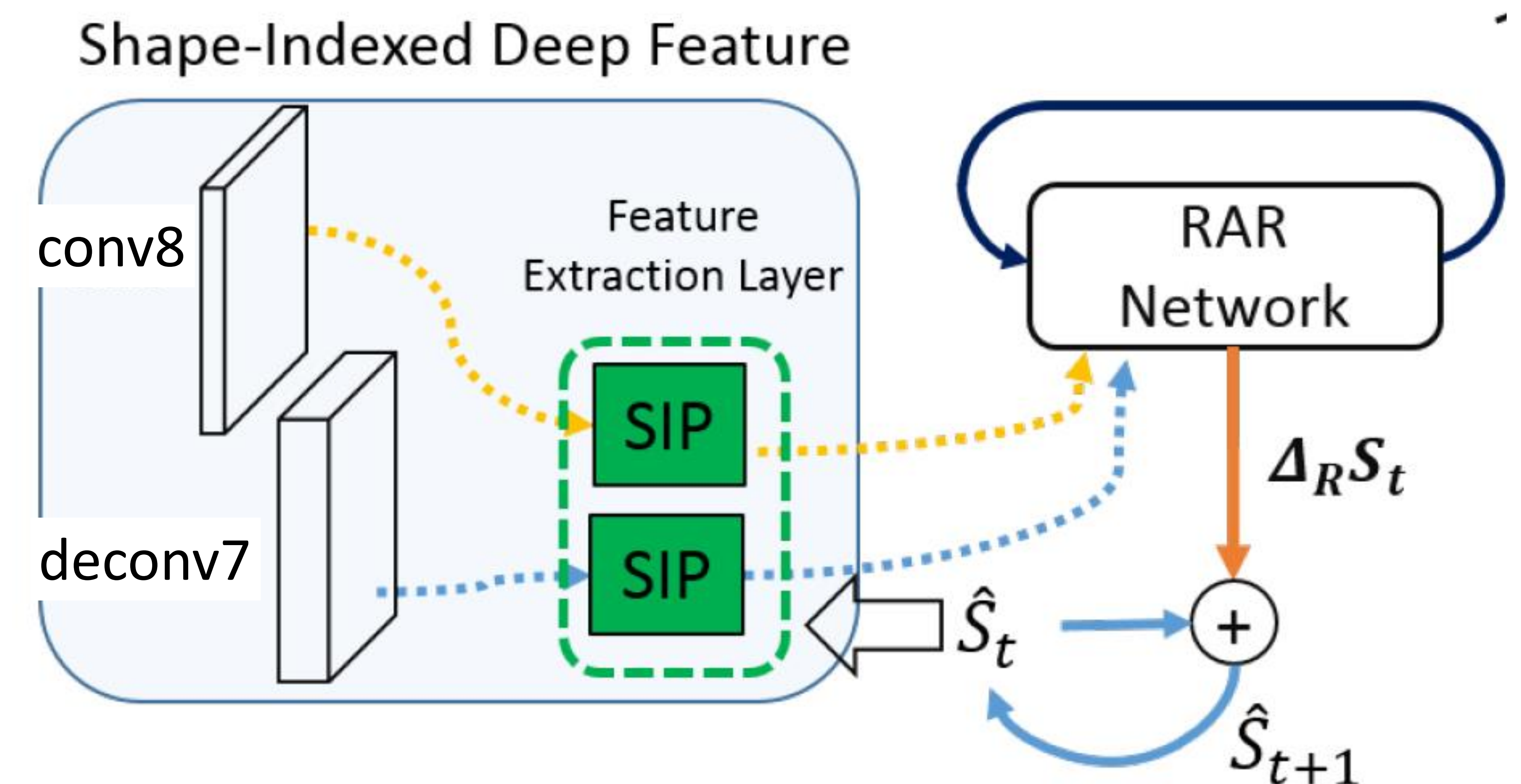
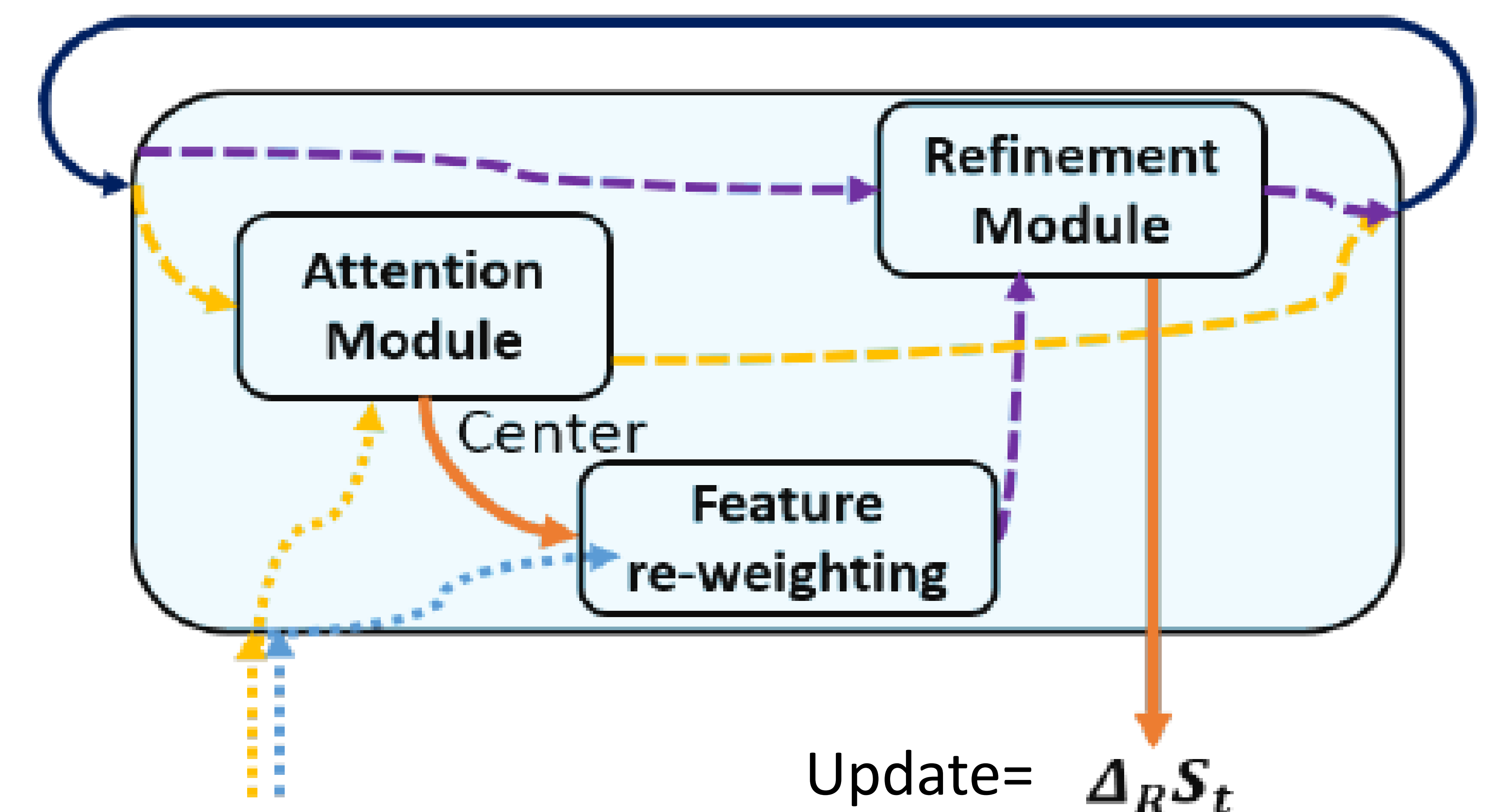
- Refinement Module to get shape update such:

$$\mathcal{L}_R^t = \|\Gamma_t(\Delta_R S_t - \Delta S_t)\|_2^2$$

Update

Ground-truth Residue

$\Delta_R S_t$  is the R-LSTM output  $\Delta_R S_t = \alpha \Gamma_t \text{R-LSTM}(\Phi_r)$



# RAR Networks: Refinement Module

- Feature re-weighting based on distance to attention center:

$$\Phi_r(I_t, \hat{S}_t) = [\gamma_t^1 \phi_t^1, \gamma_t^2 \phi_t^2, \dots, \gamma_t^L \phi_t^L]$$

- Refinement Module to get shape update such:

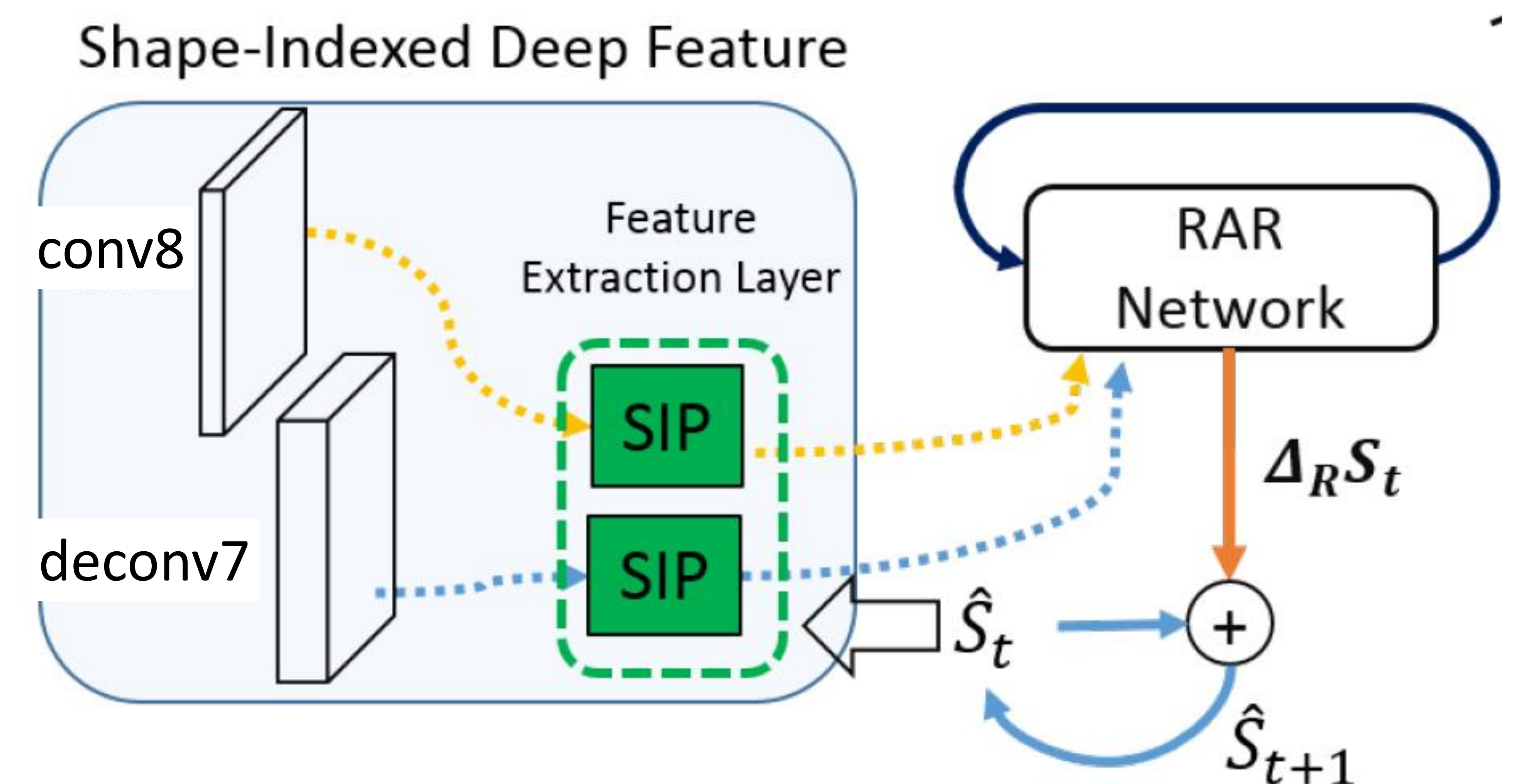
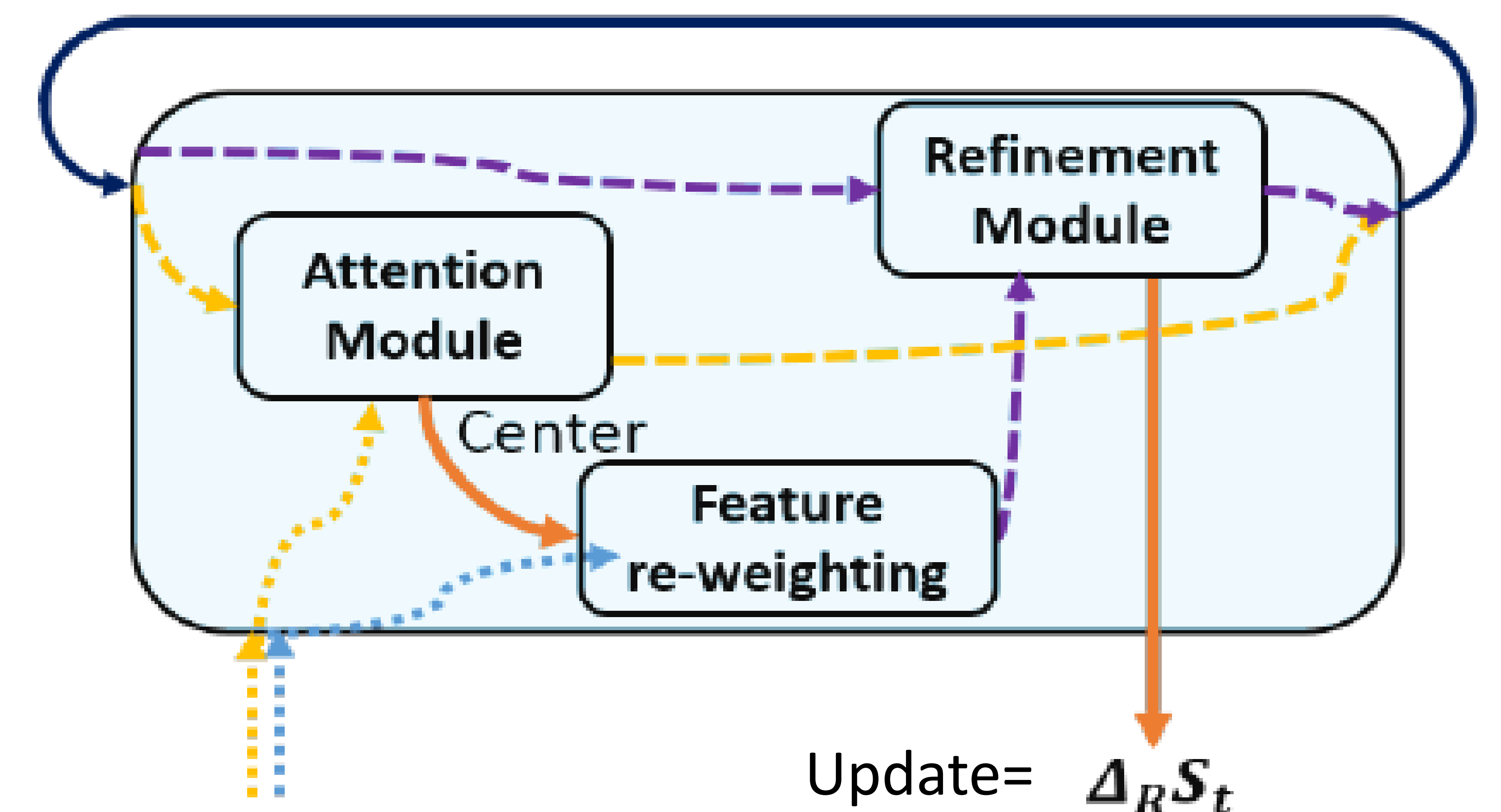
$$\mathcal{L}_R^t = || \Gamma_t(\Delta_R S_t - \Delta S_t) ||_2^2$$

$\Delta_R S_t$  is the R-LSTM output  $\Delta_R S_t = \alpha \Gamma_t \text{R-LSTM}(\Phi_r)$

Overall Training Objective of RAR:

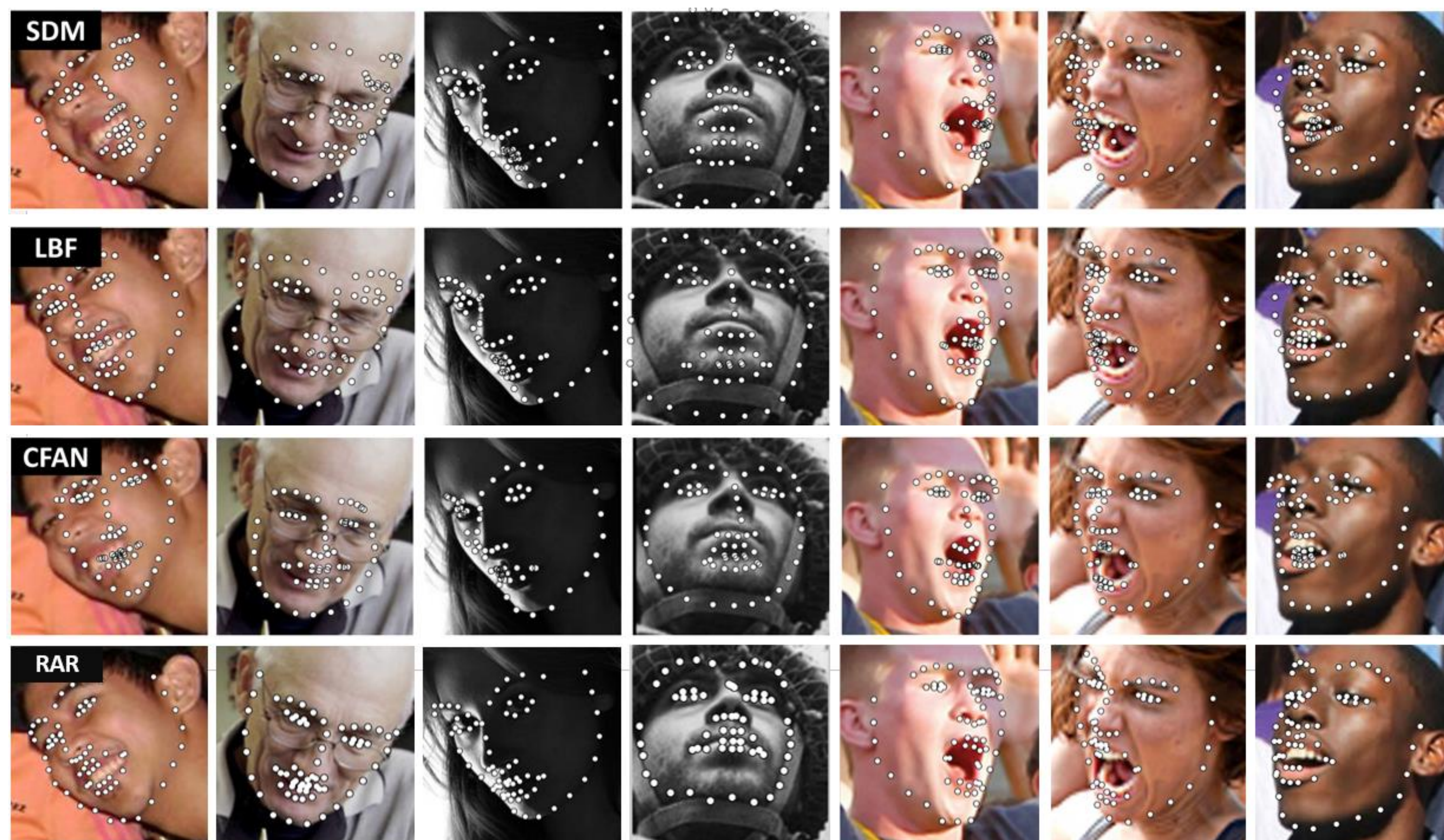
$$\sum_{t=1}^T \sum_{n=1}^N -\eta^{t-1} \mathcal{R}_a(\hat{S}_{t-1,n}, \hat{S}_{t,n}) + \mathcal{L}_{R,n}^t$$

Attention Loss
Regression Loss





# Results on 300W



Methods	300-W Dataset		
	Common	Challenging	Full
Zhu et.al [2012]	8.22	18.33	12.0
RCPR [Burgos,2013]	6.18	17.26	8.35
SDM [Xiong,2013]	5.57	15.40	7.50
LBF [Ren,2014]	4.95	11.98	6.32
LBF Fast [Ren,2014]	5.38	15.50	7.37
CFAN[Zhang, 2014]	5.50	-	-
CFSS [Zhu, 2015]	4.73	9.98	5.76
<b>Ours (RAR)</b>	<b>4.12</b>	<b>8.35</b>	<b>4.94</b>

Zhu: Face detection, pose estimation, and landmark localization in the wild. CVPR 2012

RCPR: Robust face landmark estimation under occlusion. ICCV 2013

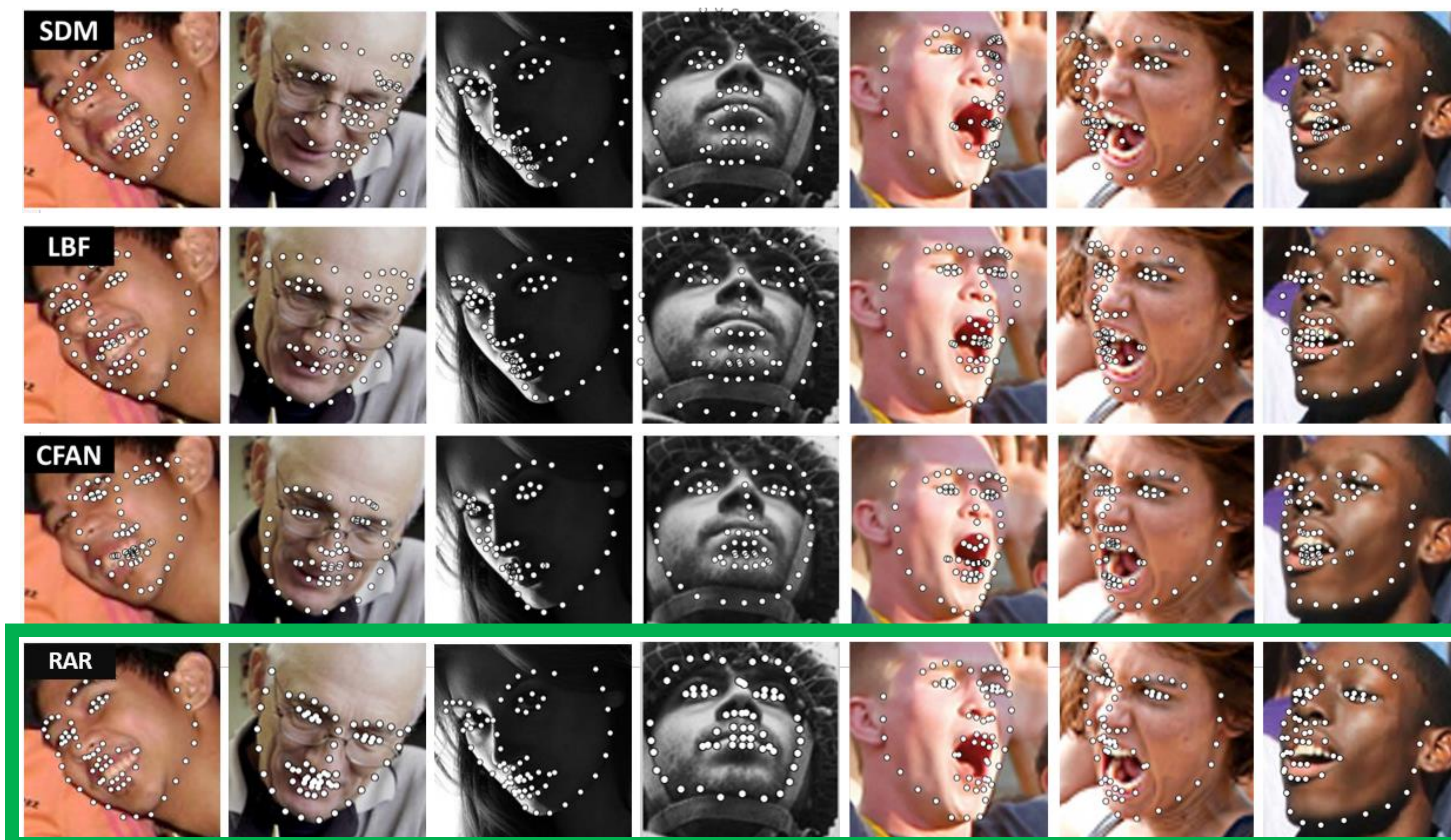
SDM: Supervised descent method and its applications to face alignment. CVPR 2013

LBF: Face alignment at 3000 fps via regressing local binary features. ECCV 2014

CFAN: Coarse-to-ne auto-encoder networks(cfan) for real-time face alignment. ECCV2014

CFSS: Face alignment by coarse-to-fine shape searching. CVPR 2015

# Results on 300W



Methods	300-W Dataset		
	Common	Challenging	Full
Zhu et.al [2012]	8.22	18.33	12.0
RCPR [Burgos,2013]	6.18	17.26	8.35
SDM [Xiong,2013]	5.57	15.40	7.50
LBF [Ren,2014]	4.95	11.98	6.32
LBF Fast [Ren,2014]	5.38	15.50	7.37
CFAN[Zhang, 2014]	5.50	-	-
CFSS [Zhu, 2015]	4.73	9.98	5.76
<b>Ours (RAR)</b>	<b>4.12</b>	<b>8.35</b>	<b>4.94</b>

Zhu: Face detection, pose estimation, and landmark localization in the wild. CVPR 2012

RCPR: Robust face landmark estimation under occlusion. ICCV 2013

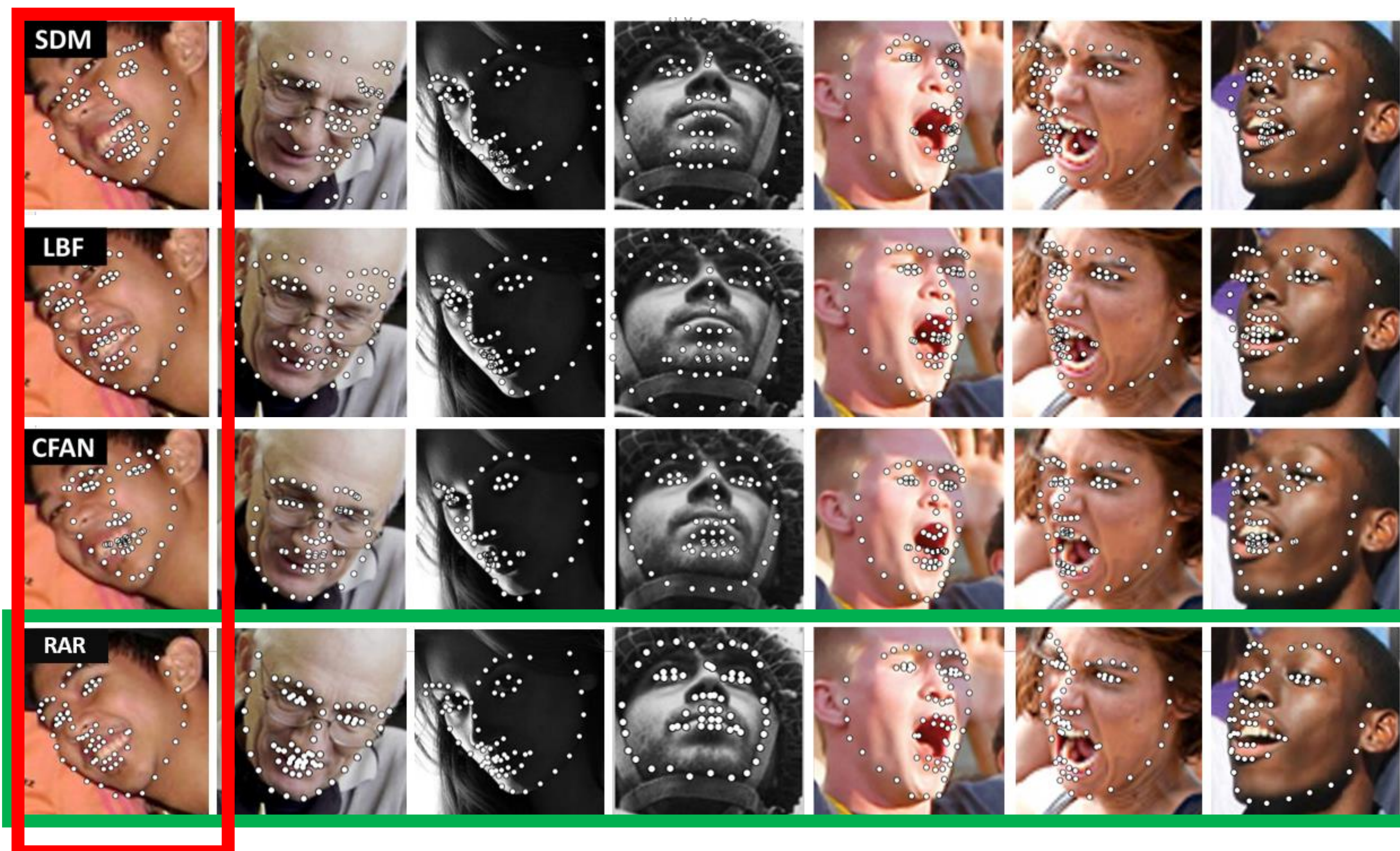
SDM: Supervised descent method and its applications to face alignment. CVPR 2013

LBF: Face alignment at 3000 fps via regressing local binary features. ECCV 2014

CFAN: Coarse-to-ne auto-encoder networks(cfan) for real-time face alignment. ECCV2014

CFSS: Face alignment by coarse-to-fine shape searching. CVPR 2015

# Results on 300W



Methods	300-W Dataset		
	Common	Challenging	Full
Zhu et.al [2012]	8.22	18.33	12.0
RCPR [Burgos,2013]	6.18	17.26	8.35
SDM [Xiong,2013]	5.57	15.40	7.50
LBF [Ren,2014]	4.95	11.98	6.32
LBF Fast [Ren,2014]	5.38	15.50	7.37
CFAN[Zhang, 2014]	5.50	-	-
CFSS [Zhu, 2015]	4.73	9.98	5.76
<b>Ours (RAR)</b>	<b>4.12</b>	<b>8.35</b>	<b>4.94</b>

Zhu: Face detection, pose estimation, and landmark localization in the wild. CVPR 2012

RCPR: Robust face landmark estimation under occlusion. ICCV 2013

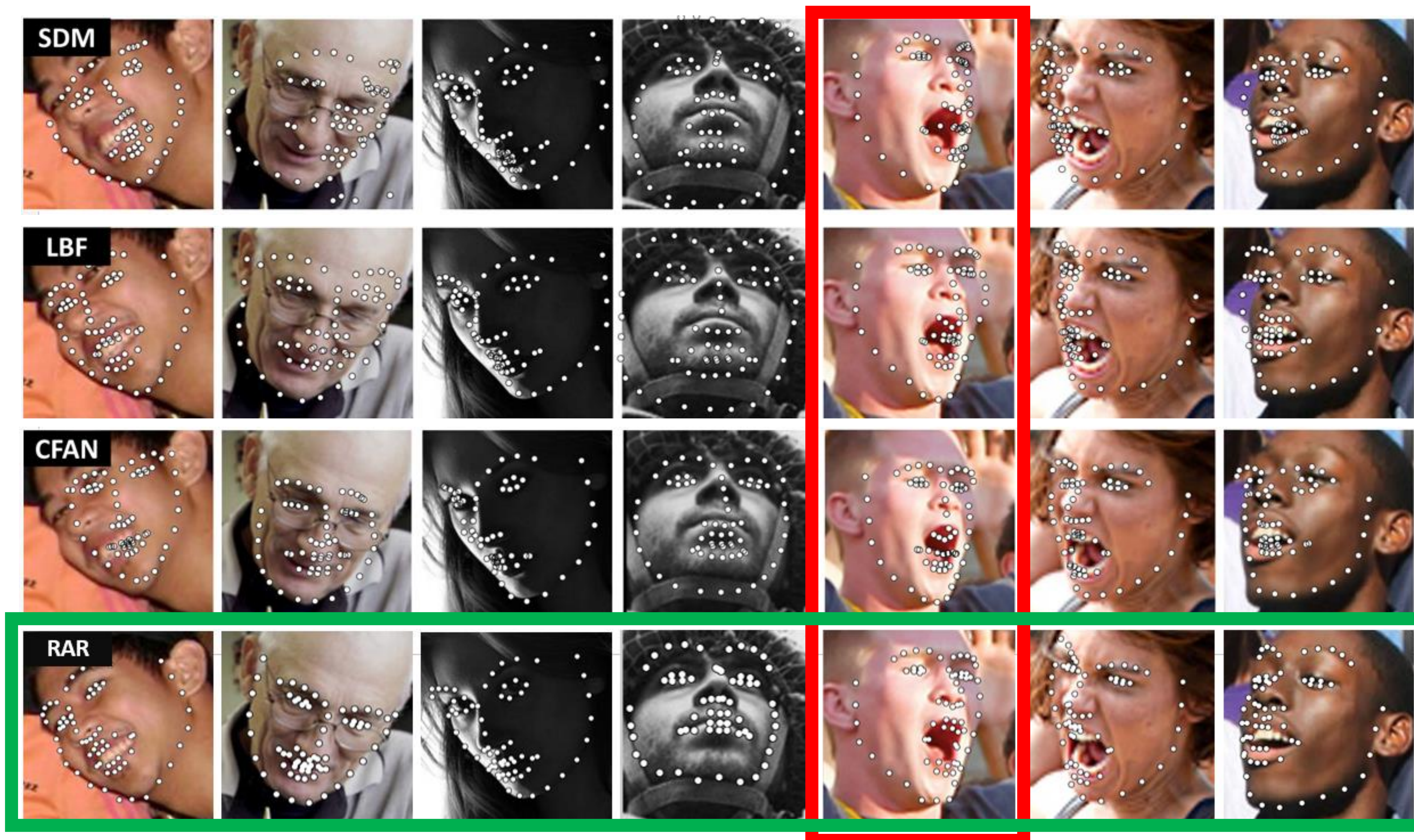
SDM: Supervised descent method and its applications to face alignment. CVPR 2013

LBF: Face alignment at 3000 fps via regressing local binary features. ECCV 2014

CFAN: Coarse-to-ne auto-encoder networks(cfan) for real-time face alignment. ECCV2014

CFSS: Face alignment by coarse-to-fine shape searching. CVPR 2015

# Results on 300W



Methods	300-W Dataset		
	Common	Challenging	Full
Zhu et.al [2012]	8.22	18.33	12.0
RCPR [Burgos,2013]	6.18	17.26	8.35
SDM [Xiong,2013]	5.57	15.40	7.50
LBF [Ren,2014]	4.95	11.98	6.32
LBF Fast [Ren,2014]	5.38	15.50	7.37
CFAN[Zhang, 2014]	5.50	-	-
CFSS [Zhu, 2015]	4.73	9.98	5.76
<b>Ours (RAR)</b>	<b>4.12</b>	<b>8.35</b>	<b>4.94</b>

Zhu: Face detection, pose estimation, and landmark localization in the wild. CVPR 2012

RCPR: Robust face landmark estimation under occlusion. ICCV 2013

SDM: Supervised descent method and its applications to face alignment. CVPR 2013

LBF: Face alignment at 3000 fps via regressing local binary features. ECCV 2014

CFAN: Coarse-to-ne auto-encoder networks(cfan) for real-time face alignment. ECCV2014

CFSS: Face alignment by coarse-to-fine shape searching. CVPR 2015

# Results on COFW and AFLW

COFW Dataset

Methods	Normalized ME	Failure Rate
RCPR	8.50	20.00%
HPM	7.46	13.24%
RPP	7.52	16.20%
TCDCN	8.05	-
<b>RAR</b>	<b>6.03</b>	<b>4.14%</b>

AFLW Dataset

Methods	Normalized ME
RCPR	11.6
SDM	8.50
CFAN	10.95
TCDCN	7.60
<b>RAR</b>	<b>7.23</b>

RCPR: Robust face landmark estimation under occlusion. ICCV 2013

HPM: Hierarchical part model. CVPR 2014.

RPP: Regional Predictive Power. TIP 2015.

TCDCN: Task constraint deep convolutional nets. PAMI 2015.

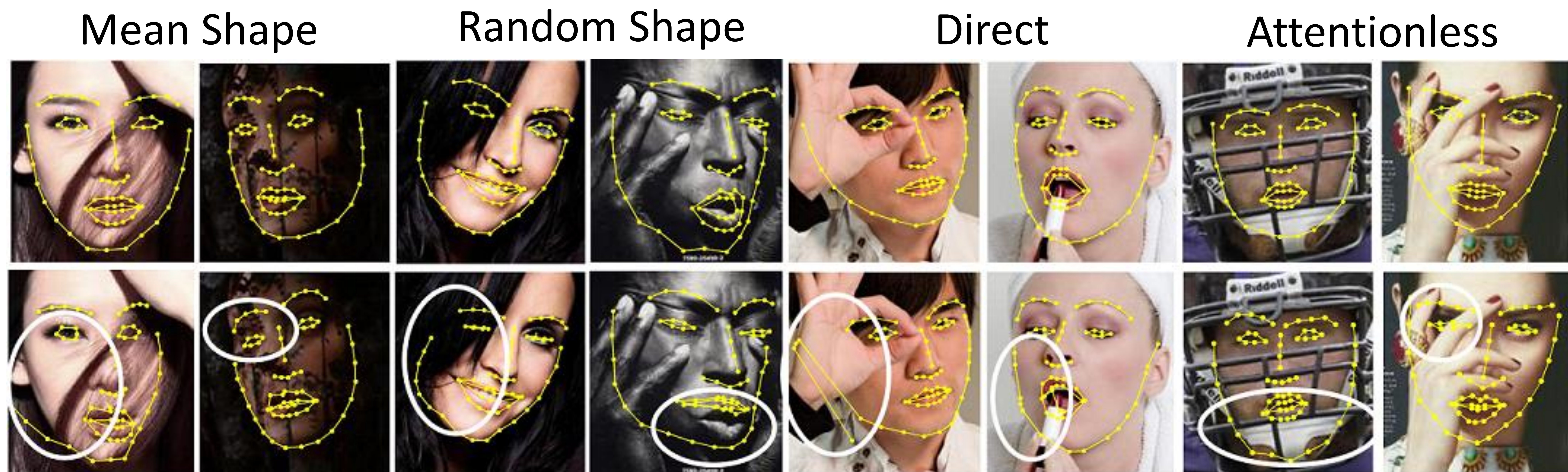
SDM: Supervised descent method and its applications to face alignment. CVPR 2013

CFAN: Coarse-to-ne auto-encoder networks(cfan) for real-time face alignment. ECCV2014

CFSS: Face alignment by coarse-to-fine shape searching. CVPR 2015

# Comparison Studies

Dataset	Conv8	Mean Shape	Random Shape	Direct	Robust Initialization
300-W	6.24	5.26	5.22	6.66	<b>4.94</b>
COFW	30.14	6.24	6.12	11.52	<b>6.03</b>
AFLW	8.14	7.36	7.42	8.15	<b>7.23</b>



Conv8:

Direct:

Mean Shape:

Random Shape:

Robust:

Prediction from conv8

RAR trained with initial shape as conv8

RAR trained with mean shape as initial shape

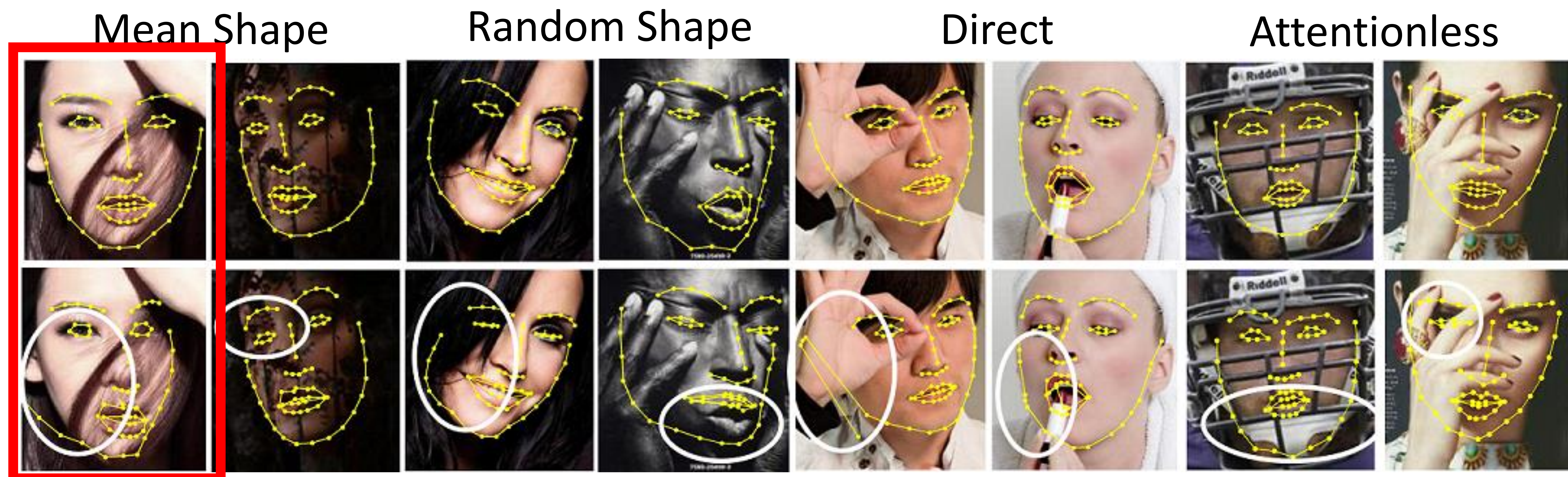
RAR trained with random shape as initial shape

RAR trained with the proposed robust initialization

RAR

# Comparison Studies

Dataset	Conv8	Mean Shape	Random Shape	Direct	Robust Initialization
300-W	6.24	5.26	5.22	6.66	<b>4.94</b>
COFW	30.14	6.24	6.12	11.52	<b>6.03</b>
AFLW	8.14	7.36	7.42	8.15	<b>7.23</b>



Conv8:

Direct:

Mean Shape:

Random Shape:

Robust:

Prediction from conv8

RAR trained with initial shape as conv8

RAR trained with mean shape as initial shape

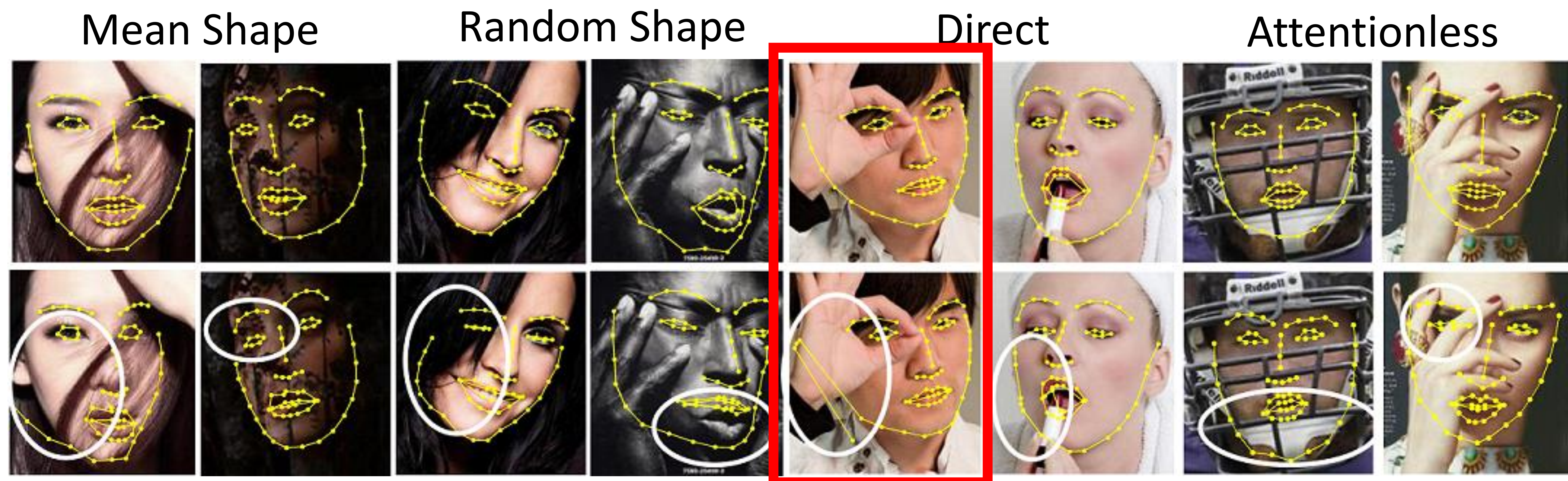
RAR trained with random shape as initial shape

RAR trained with the proposed robust initialization

RAR

# Comparison Studies

Dataset	Conv8	Mean Shape	Random Shape	Direct	Robust Initialization
300-W	6.24	5.26	5.22	6.66	<b>4.94</b>
COFW	30.14	6.24	6.12	11.52	<b>6.03</b>
AFLW	8.14	7.36	7.42	8.15	<b>7.23</b>



Conv8:

Direct:

Mean Shape:

Random Shape:

Robust:

Prediction from conv8

RAR trained with initial shape as conv8

RAR trained with mean shape as initial shape

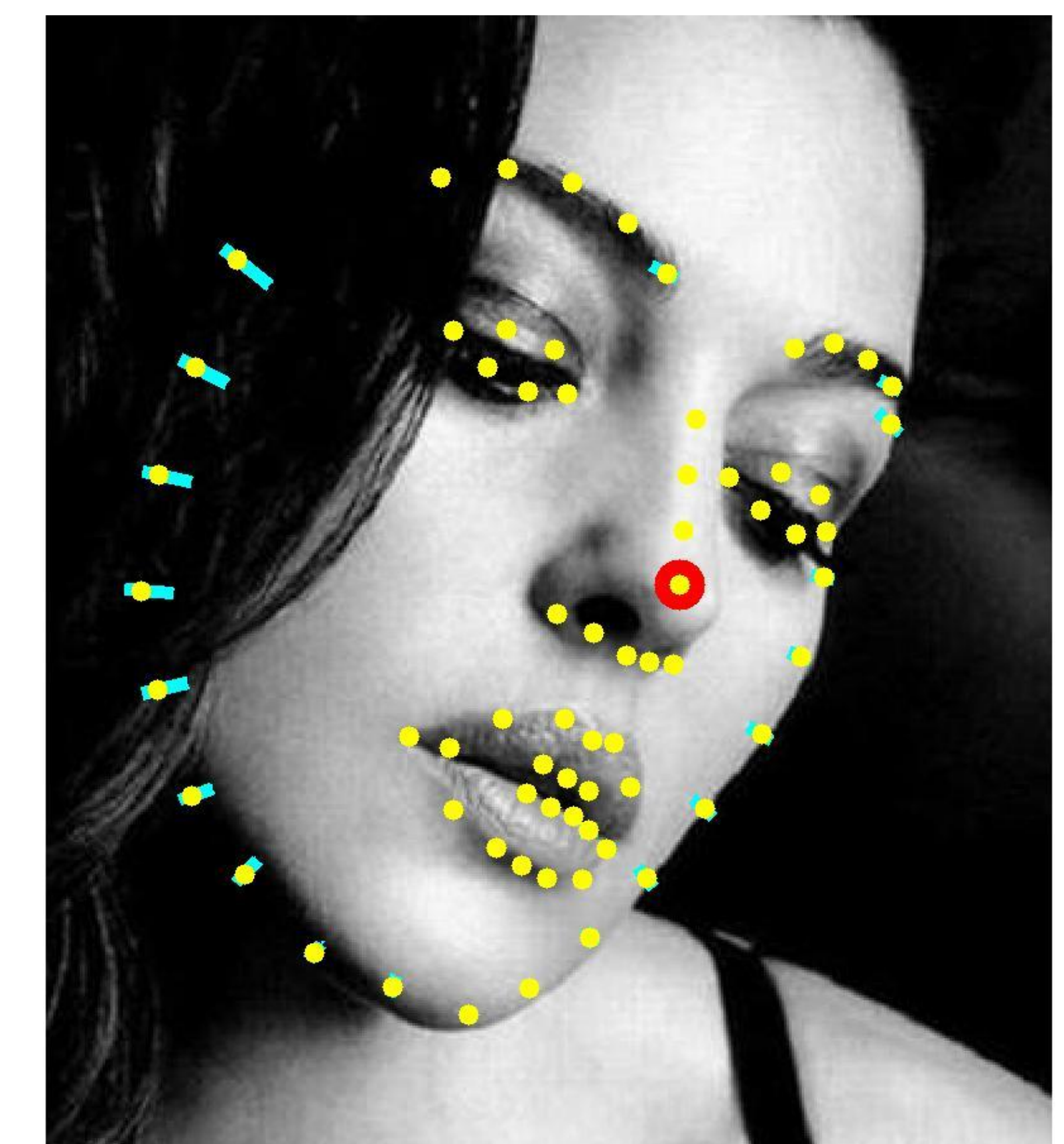
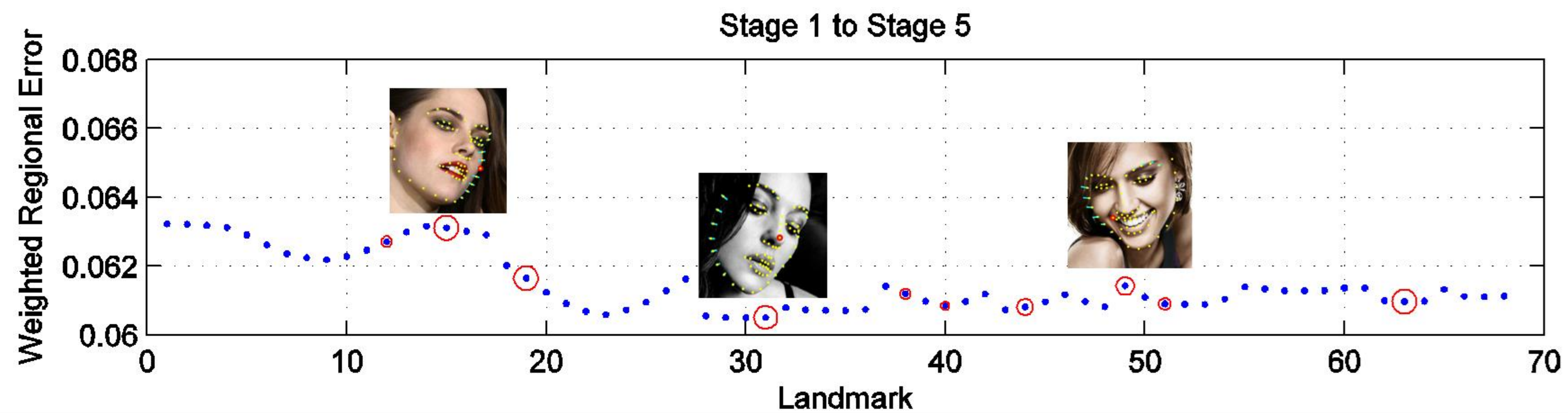
RAR trained with random shape as initial shape

RAR trained with the proposed robust initialization

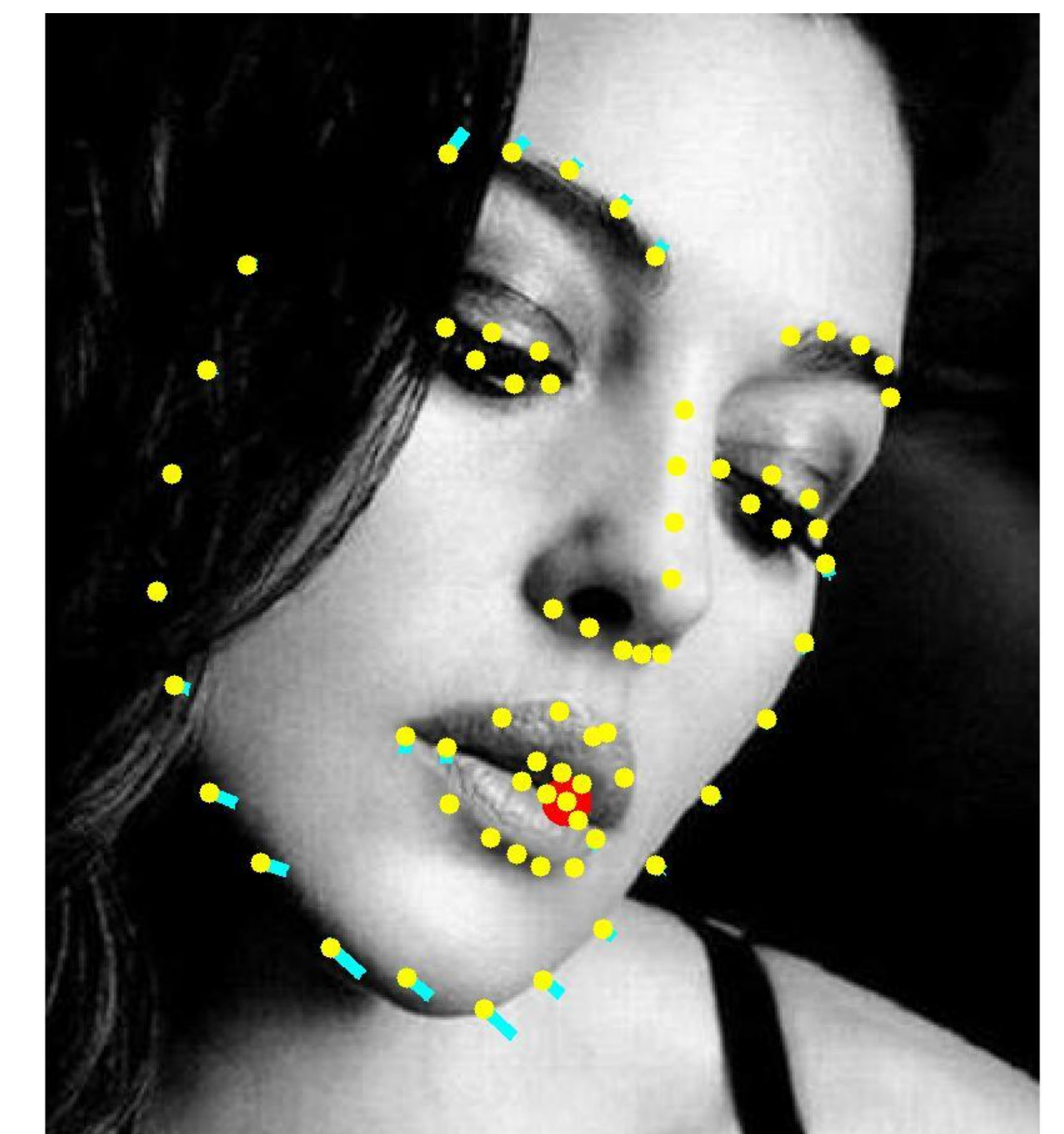
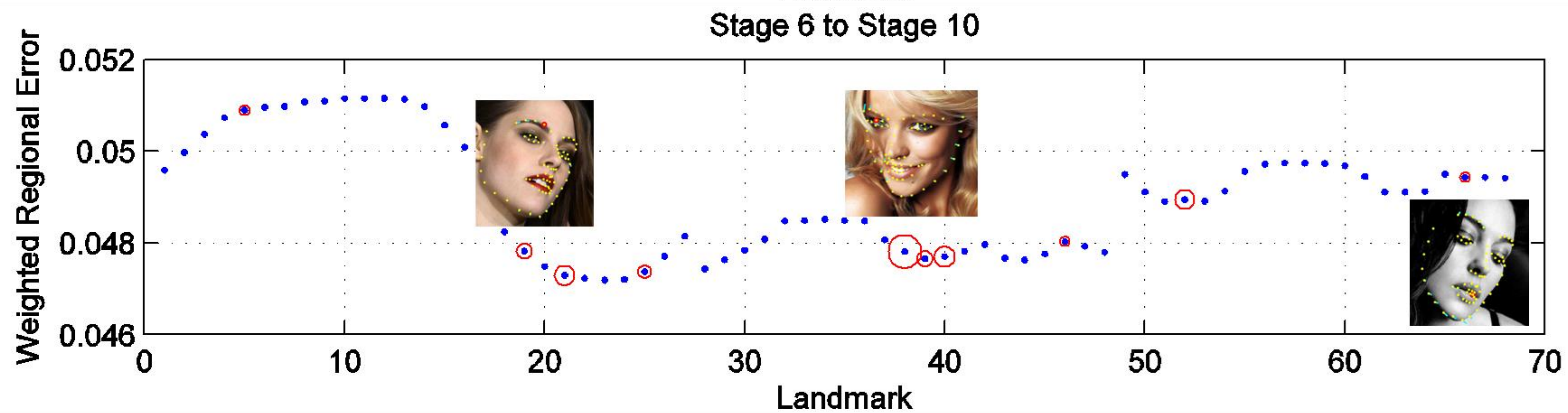
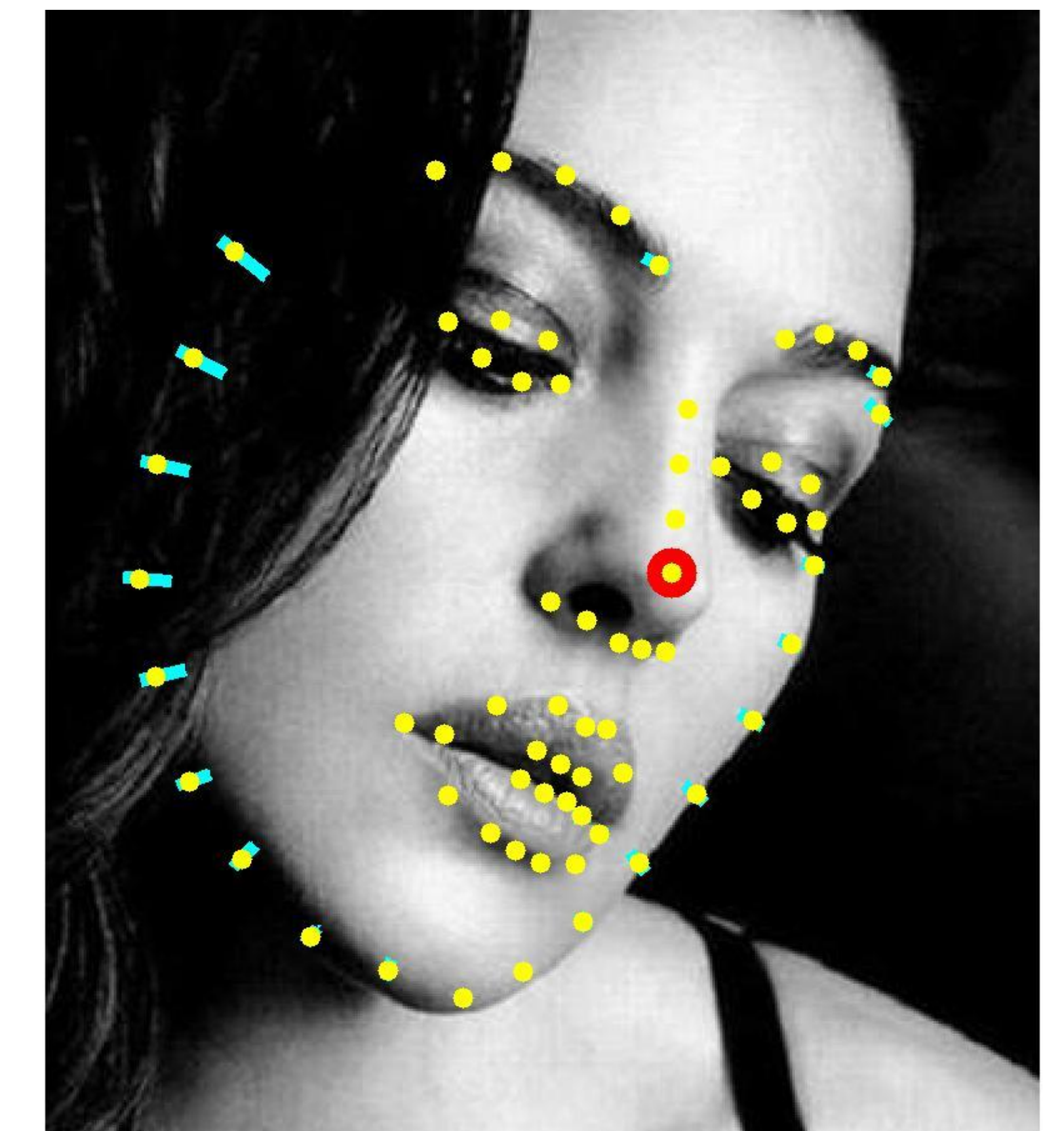
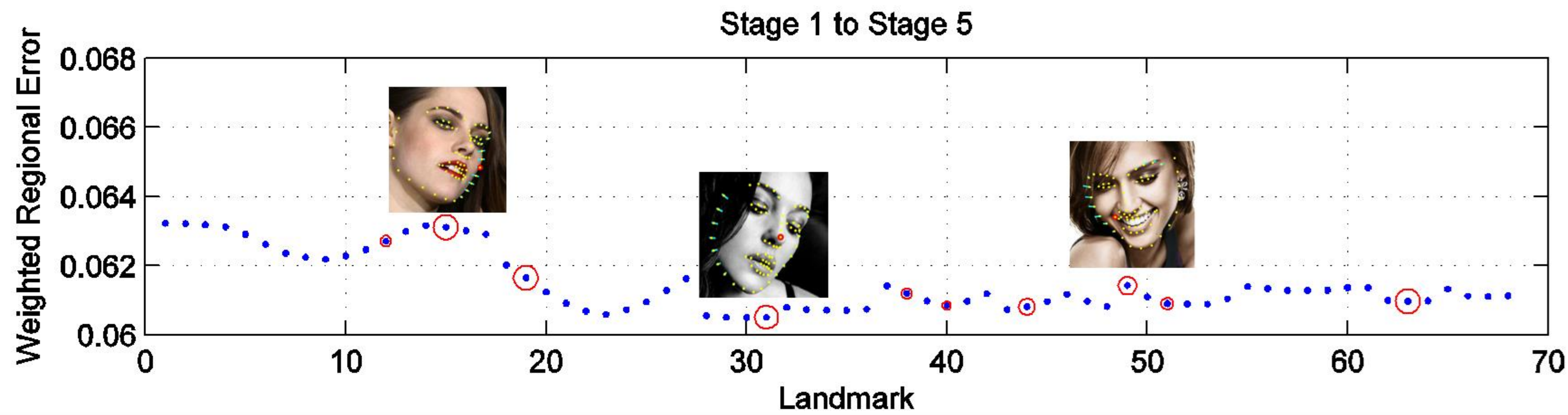
RAR



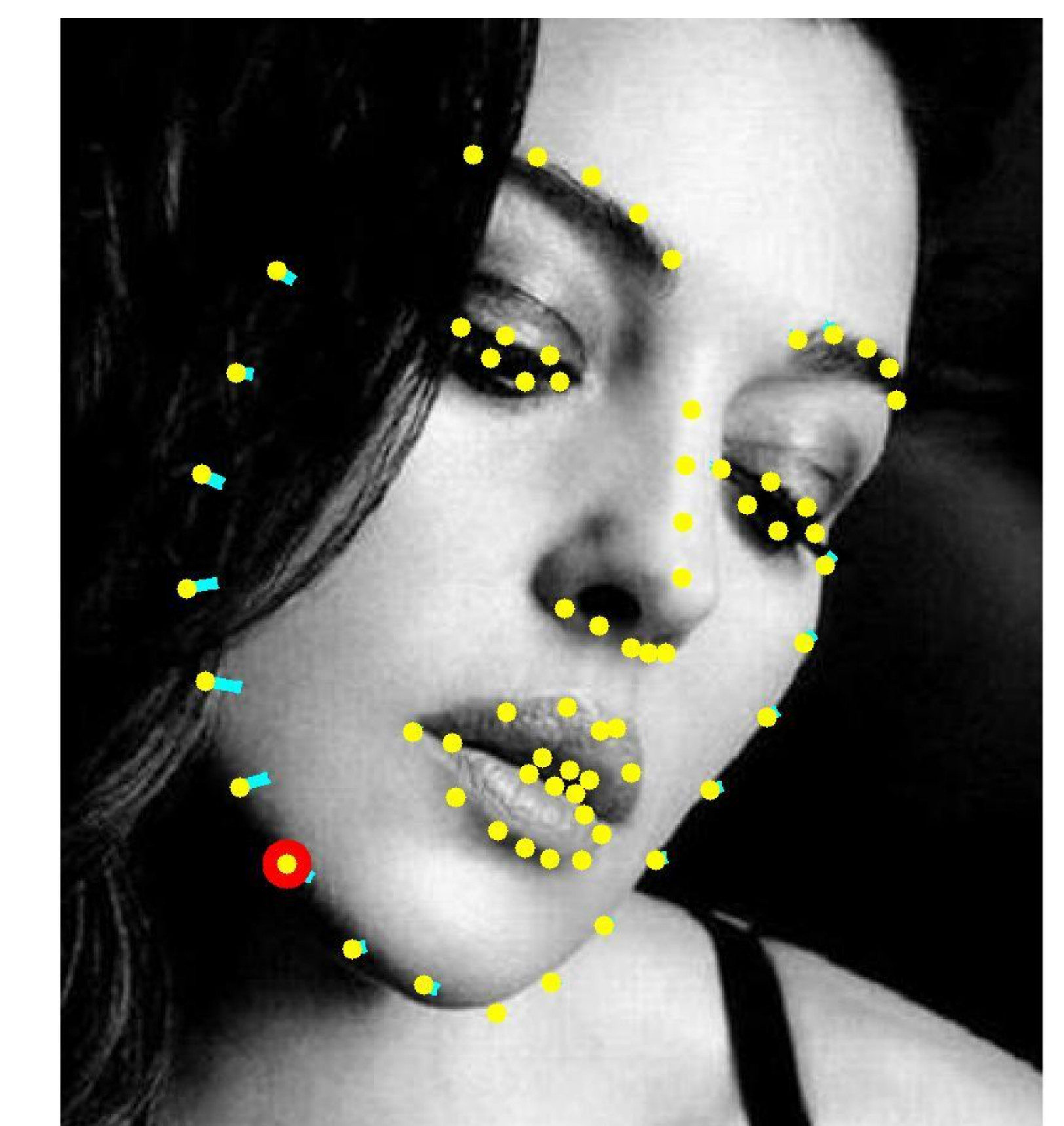
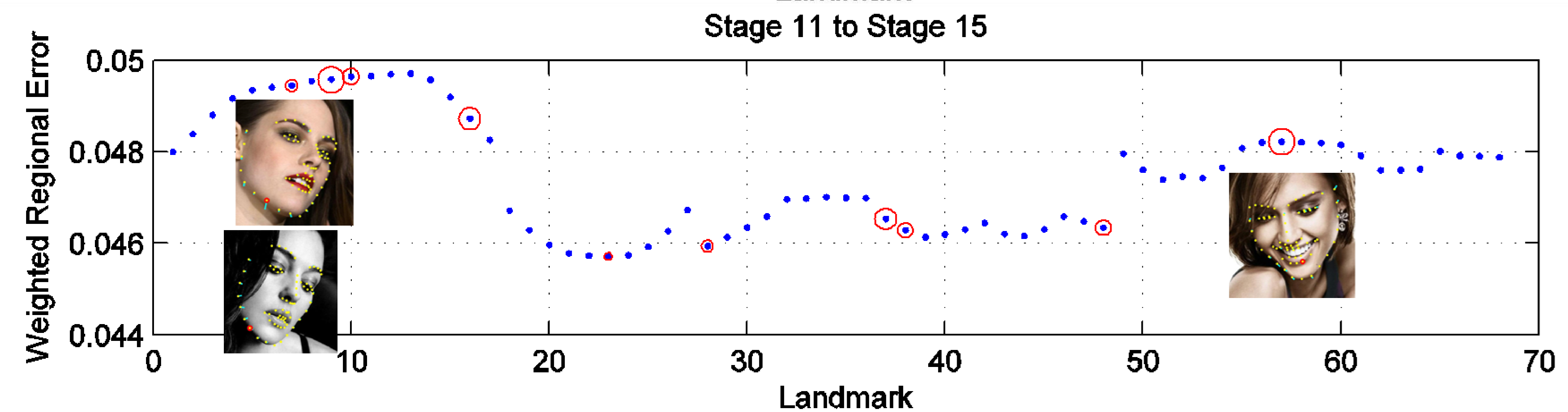
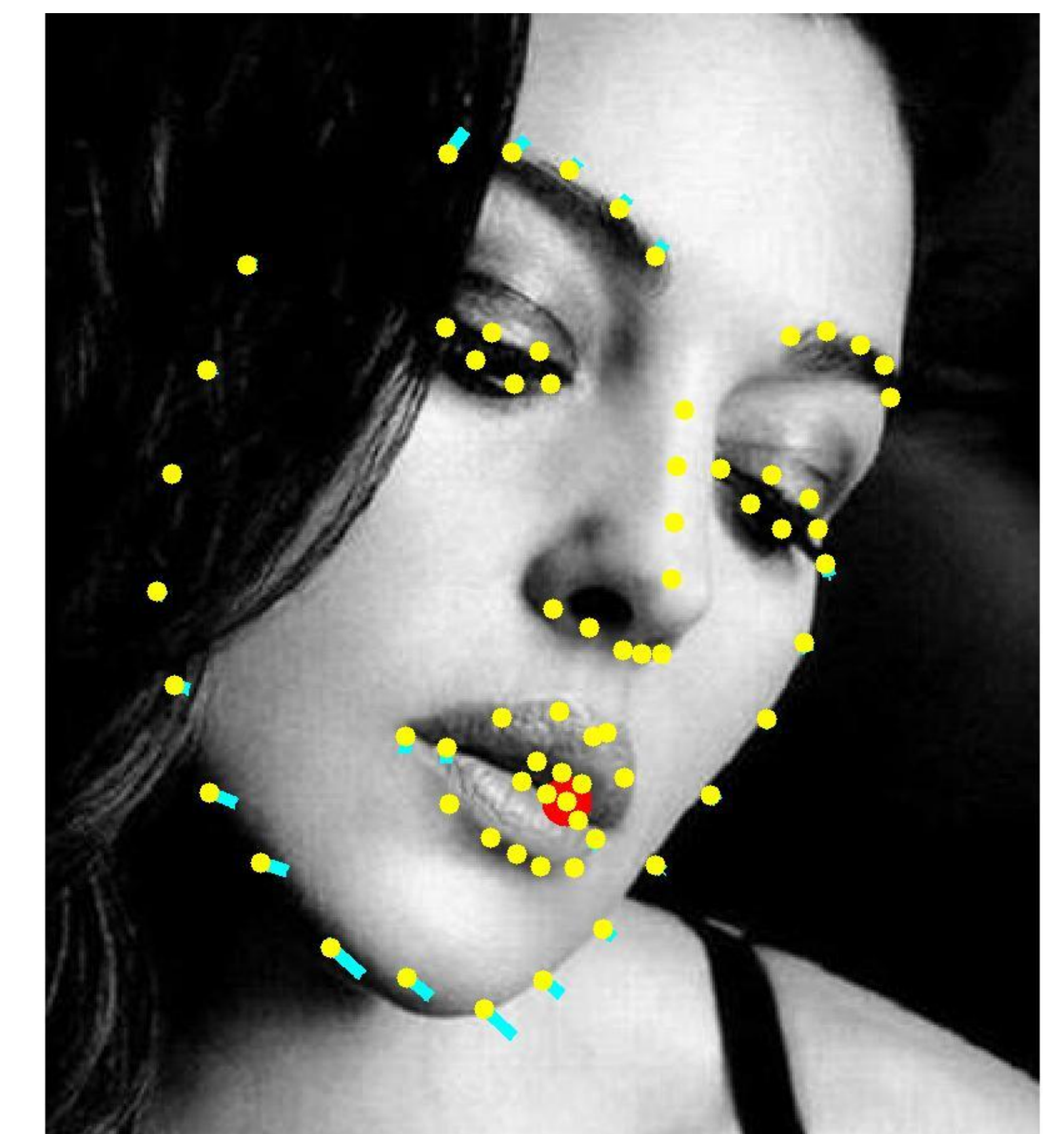
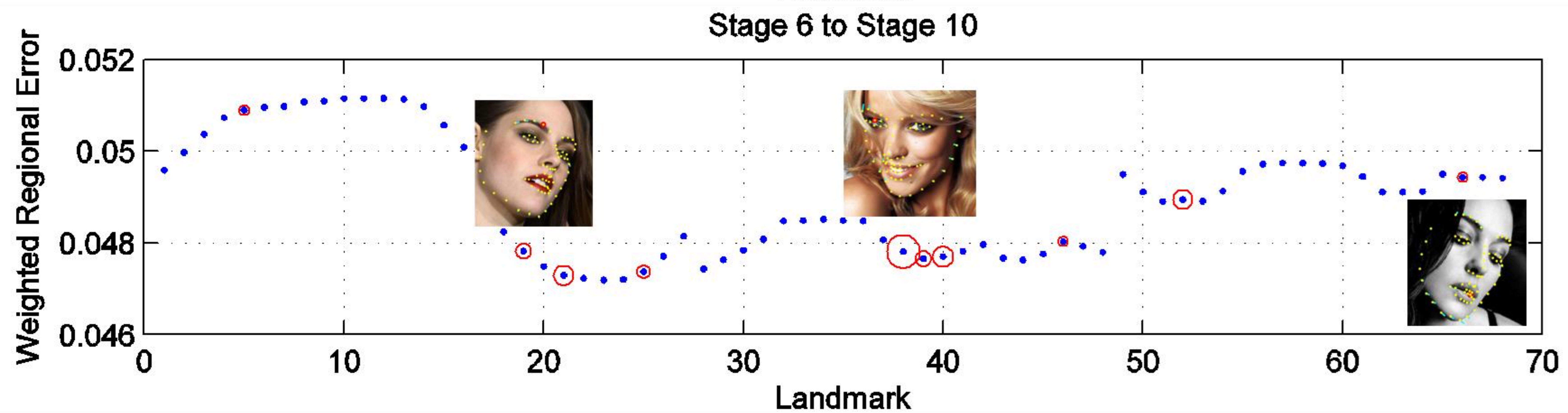
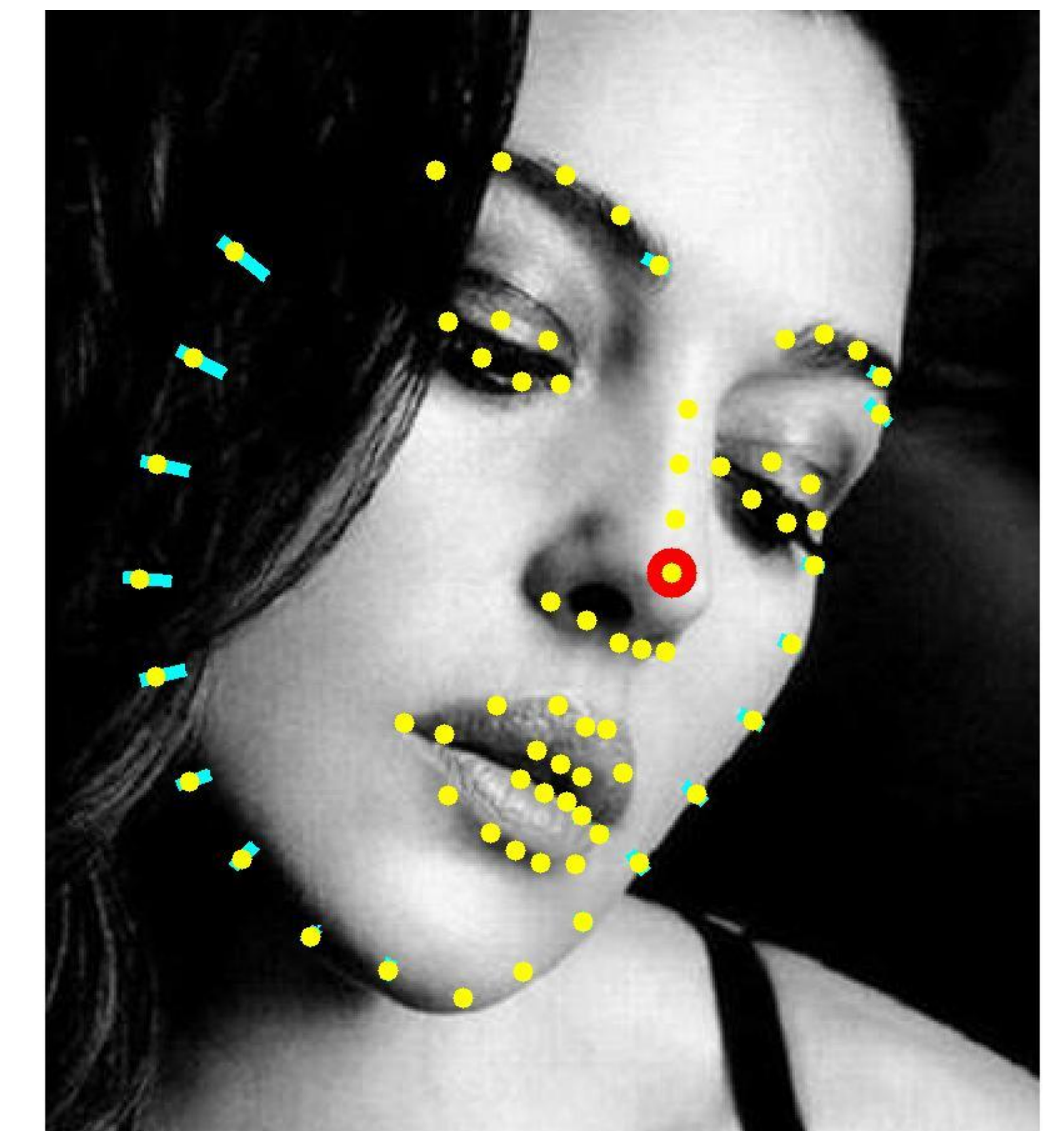
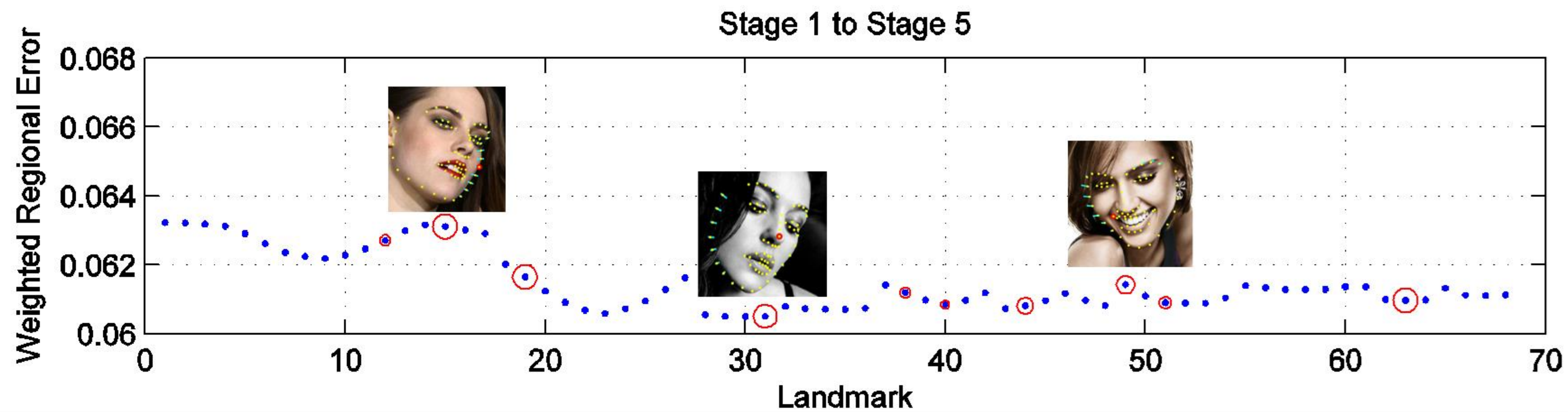
# Attention Center Selection Frequencies



# Attention Center Selection Frequencies

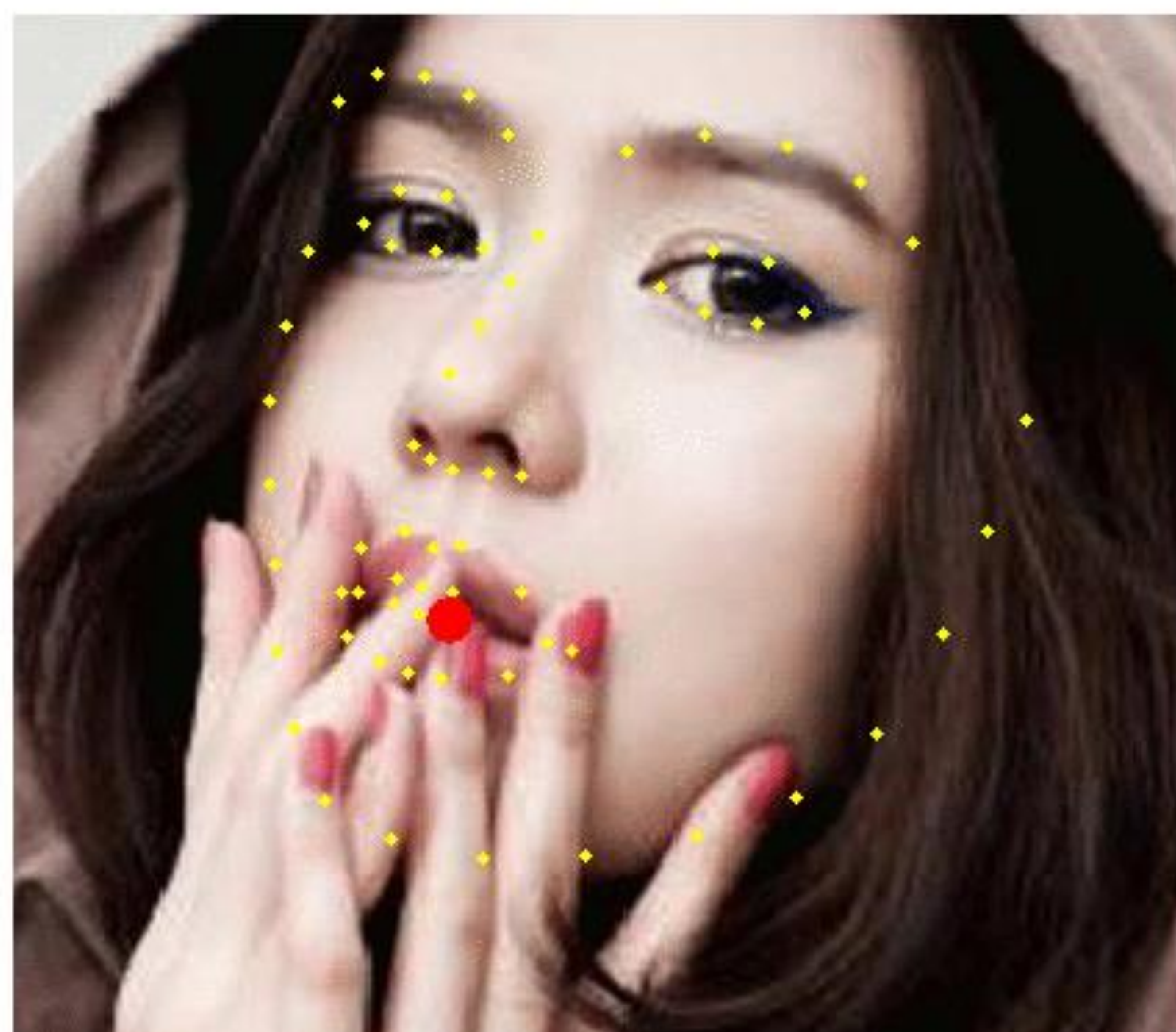


# Attention Center Selection Frequencies

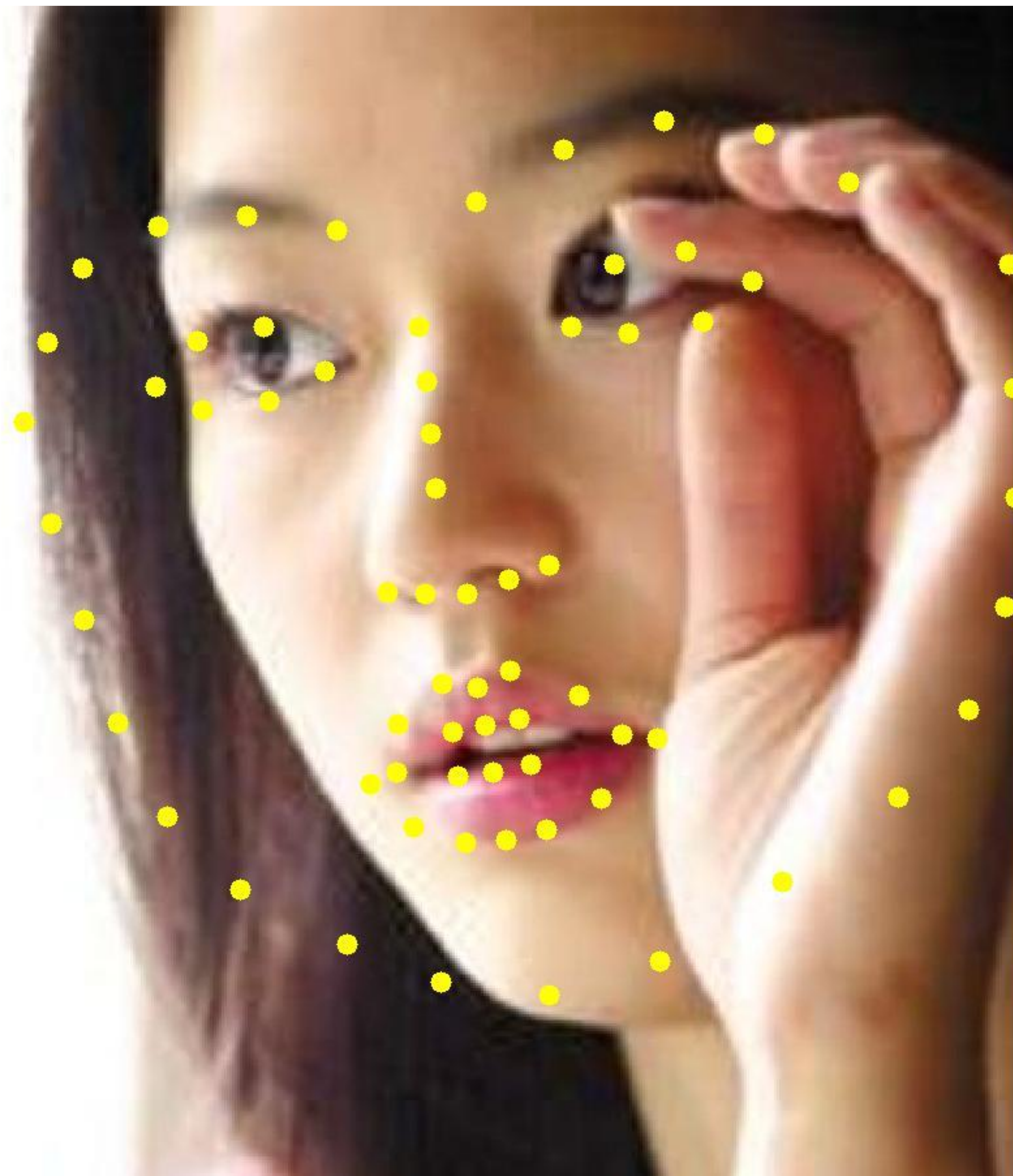


# Sample Attentive Refinement

Iter 01

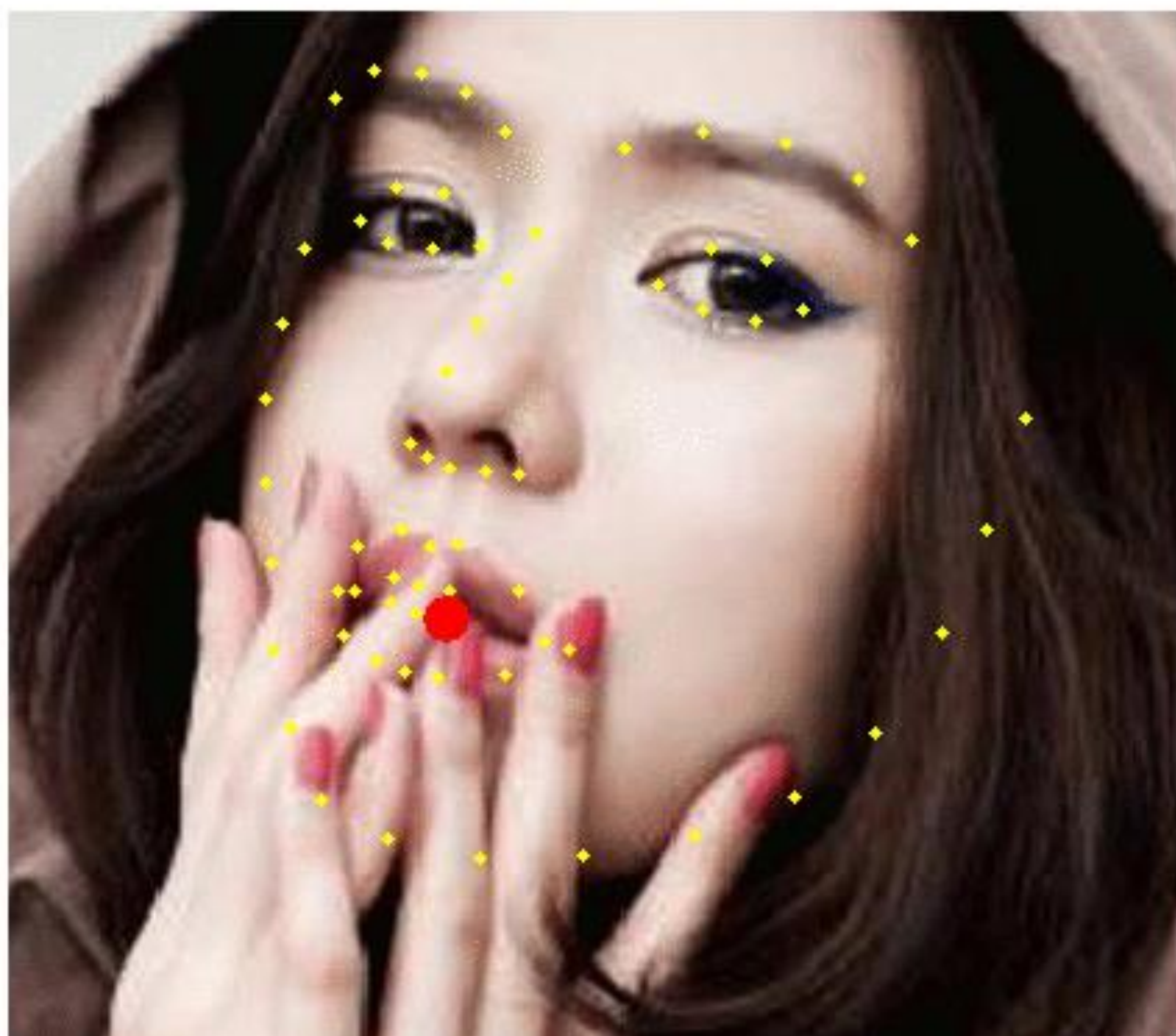


Iter 01

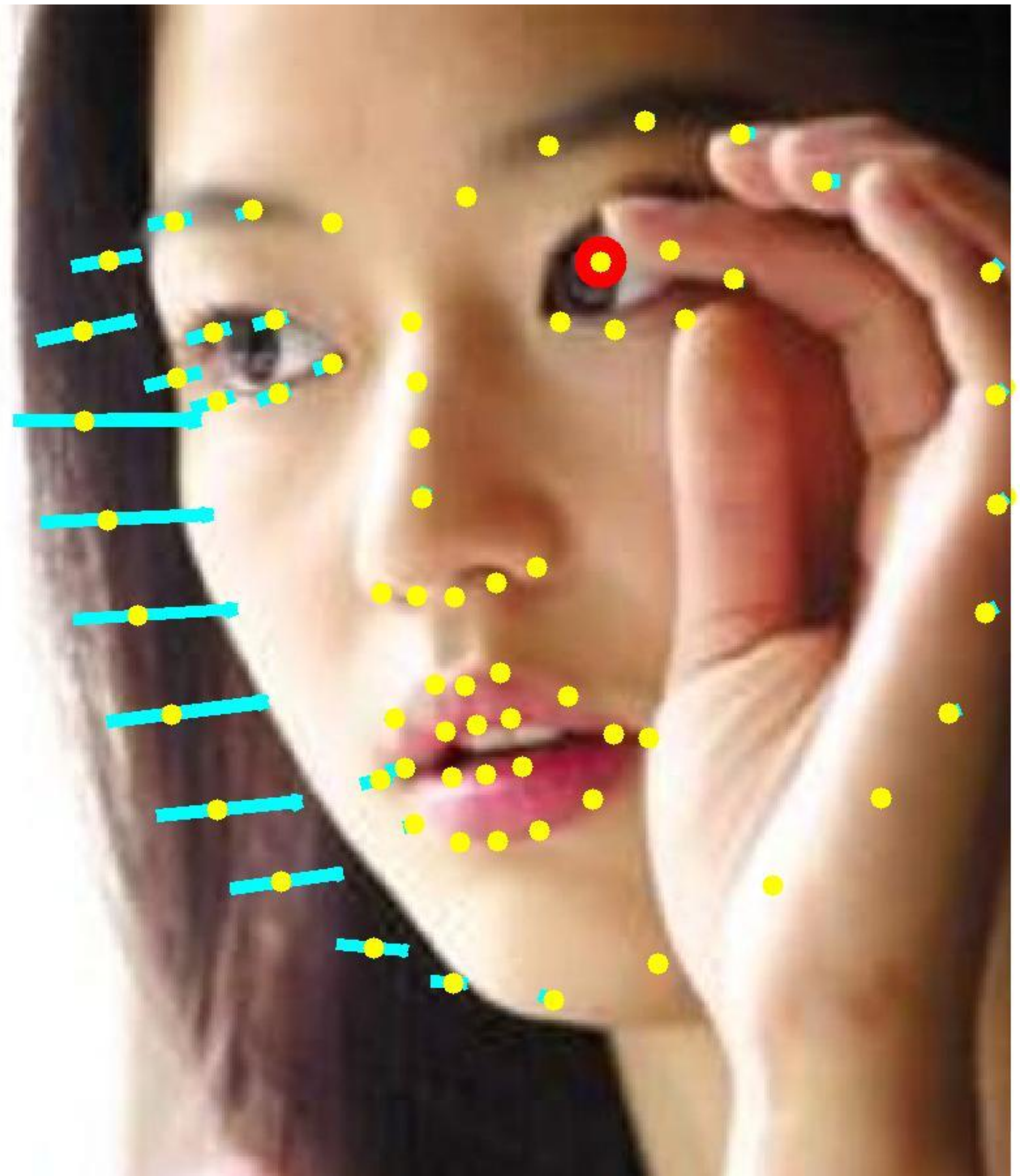
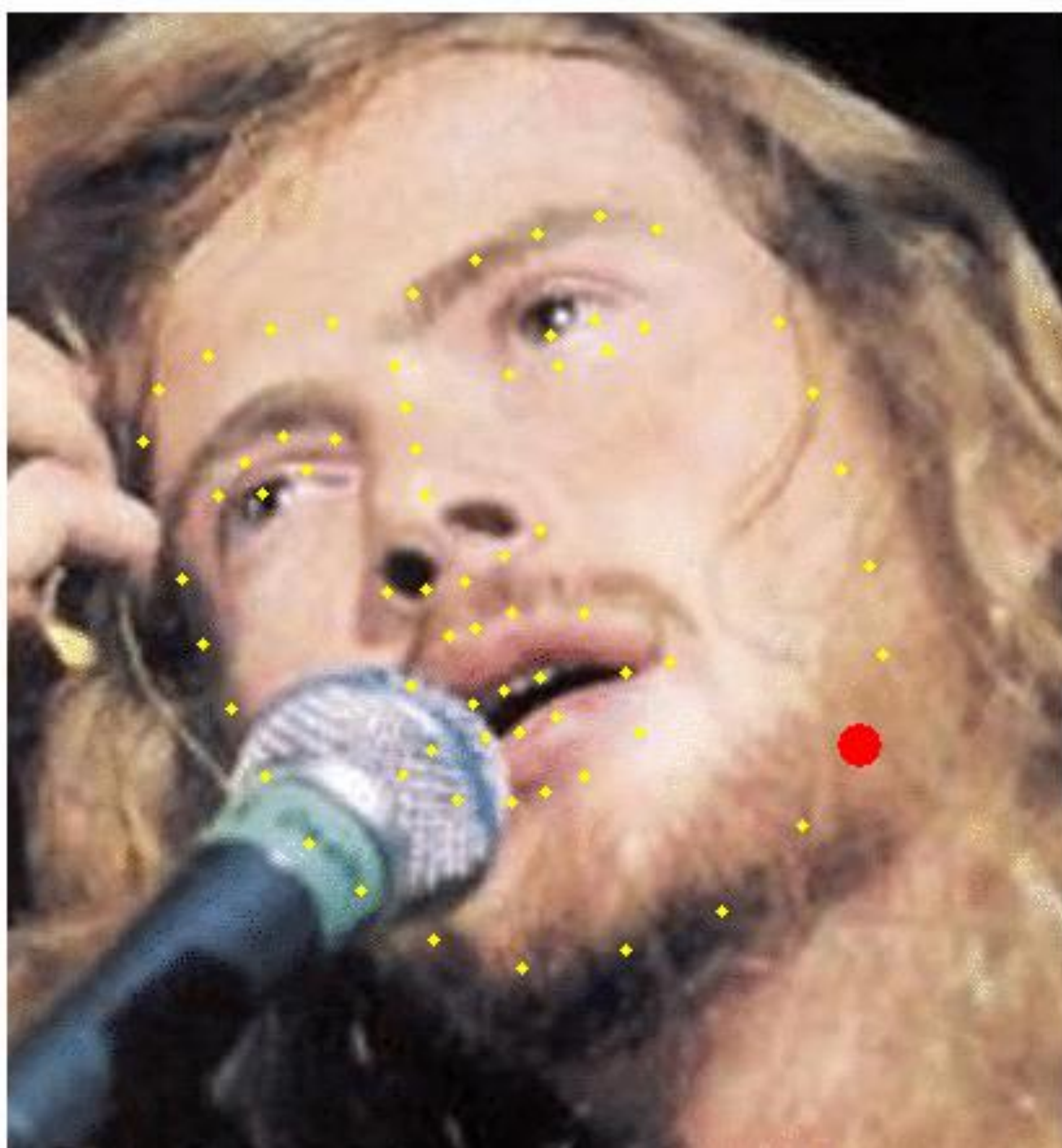


# Sample Attentive Refinement

Iter 01

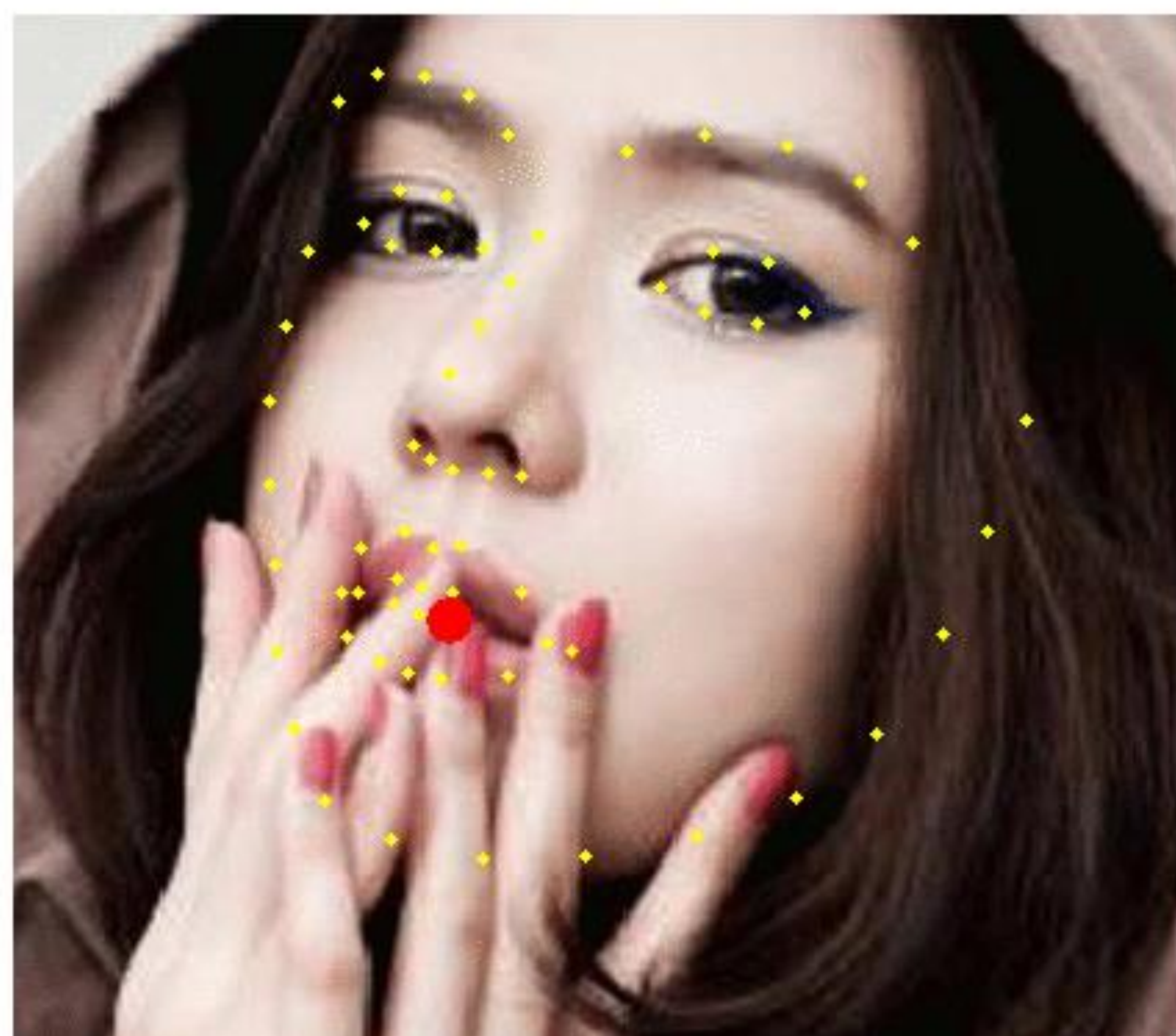


Iter 01

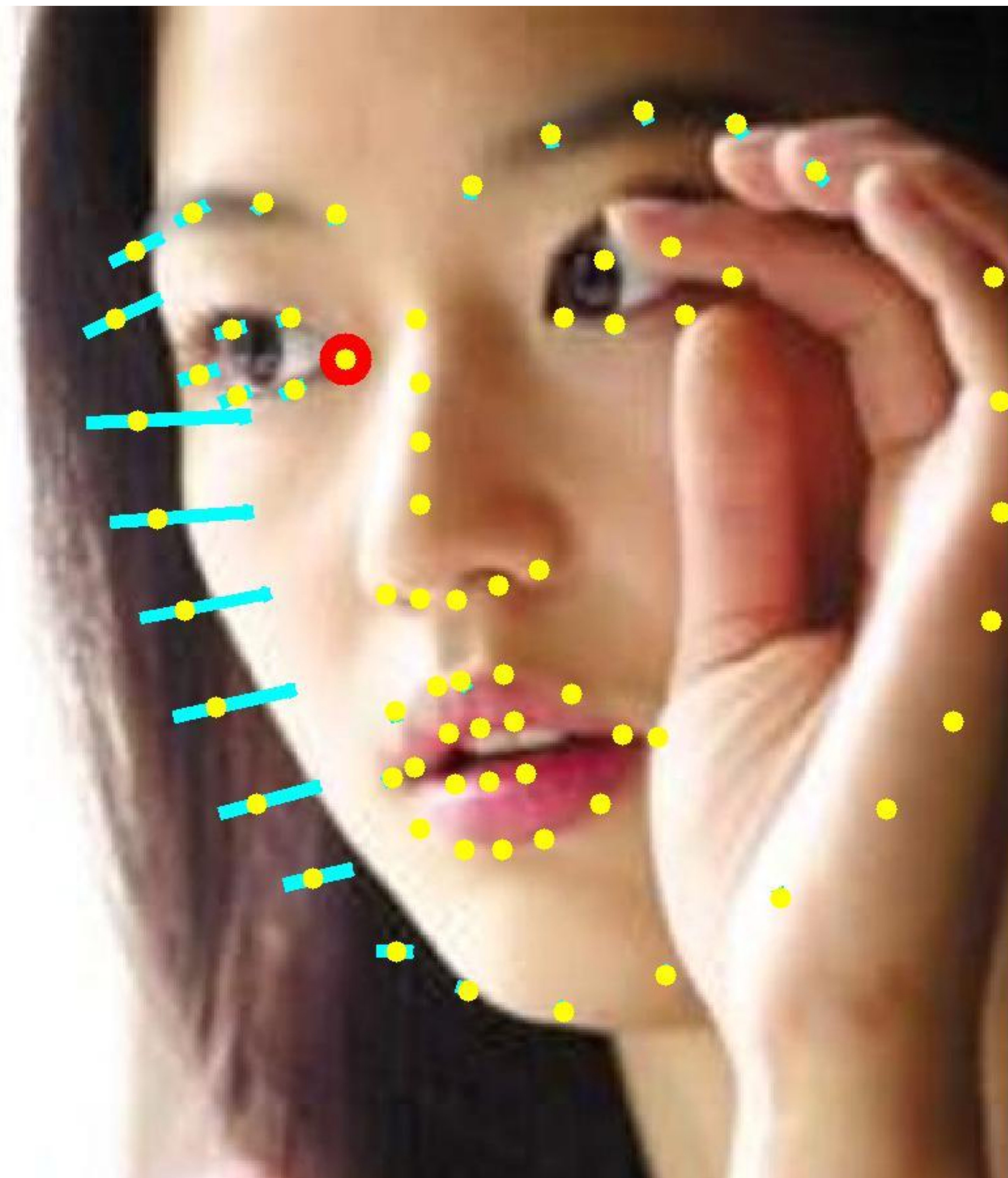
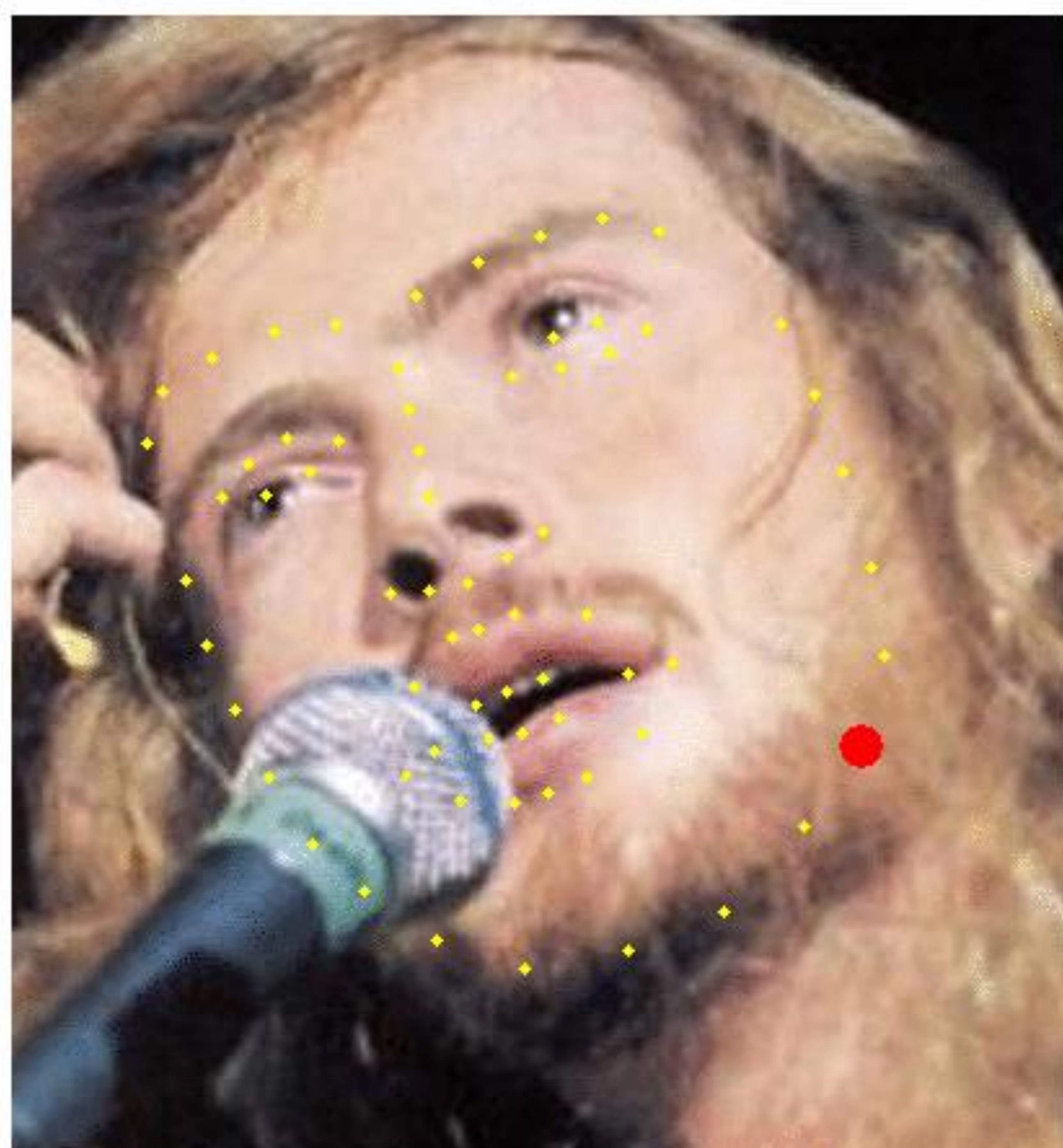


# Sample Attentive Refinement

Iter 01

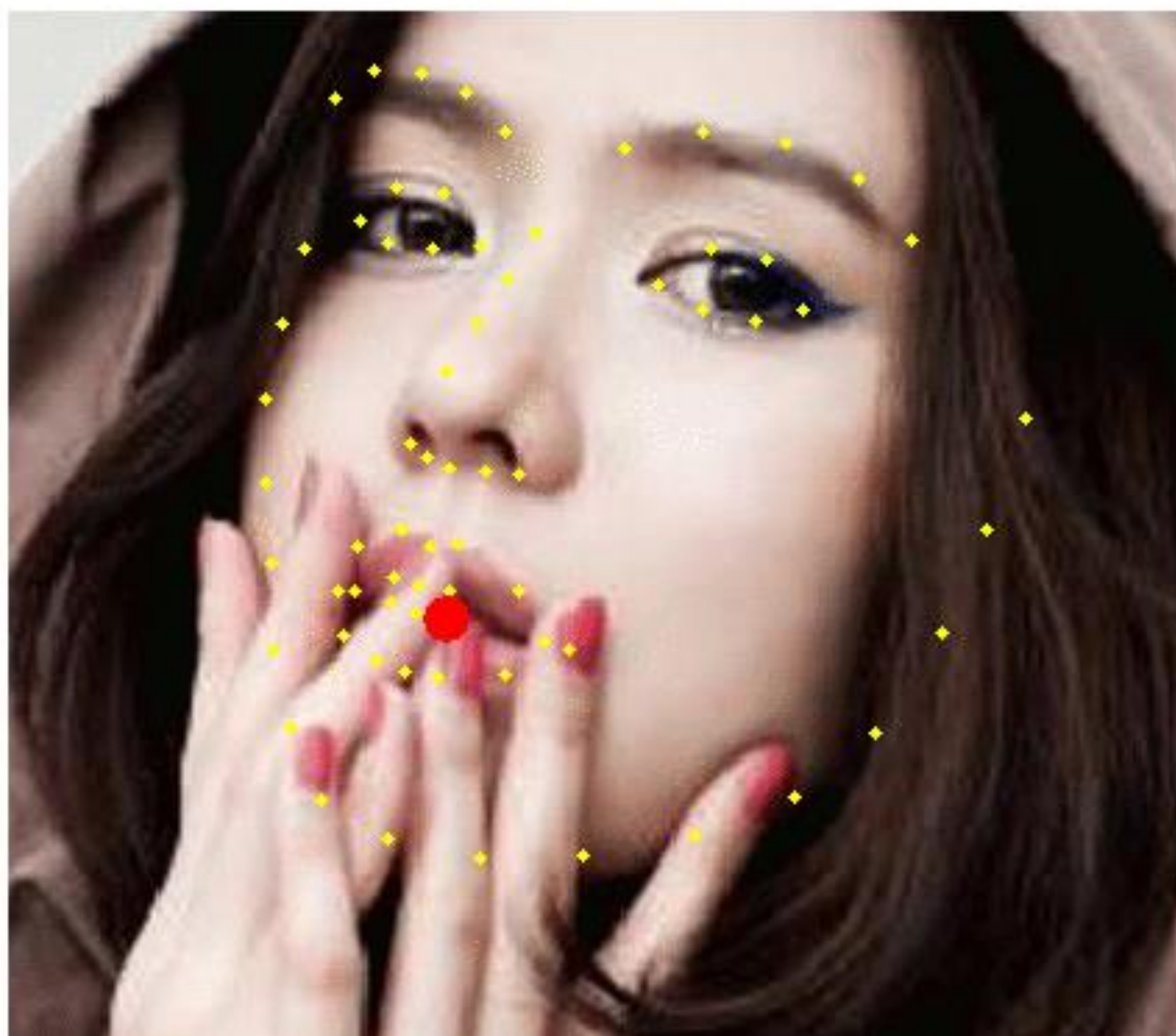


Iter 01

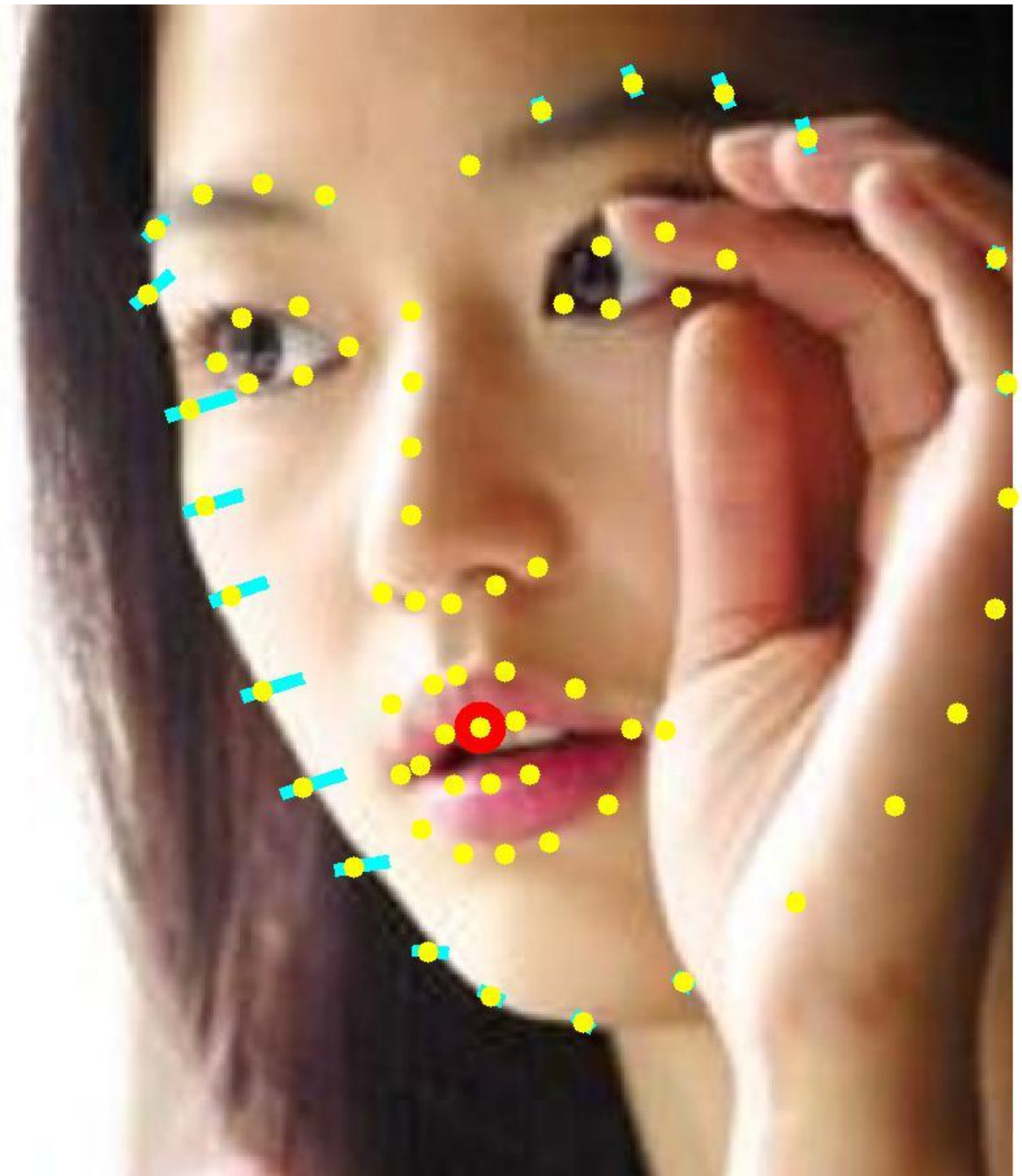
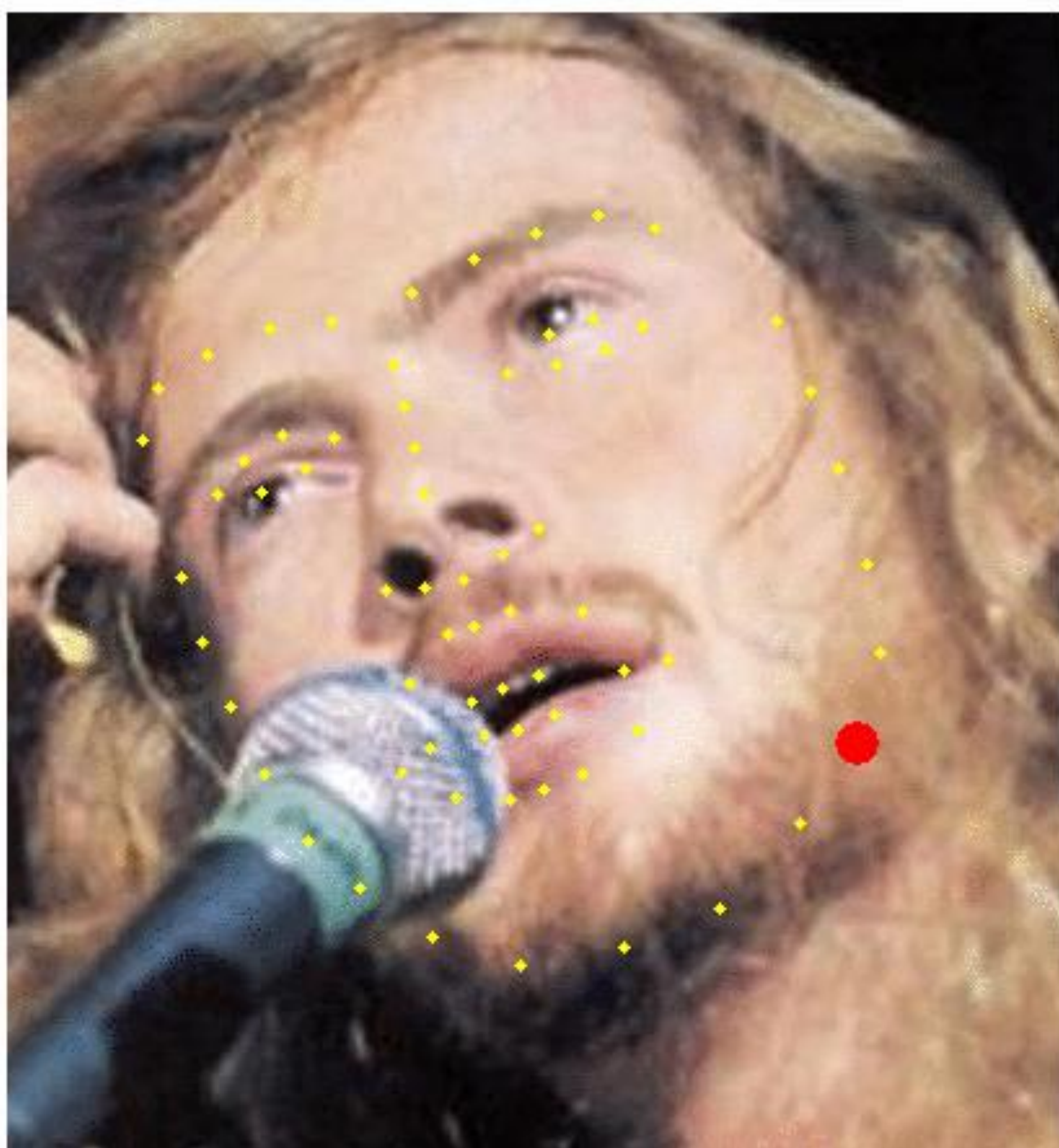


# Sample Attentive Refinement

Iter 01

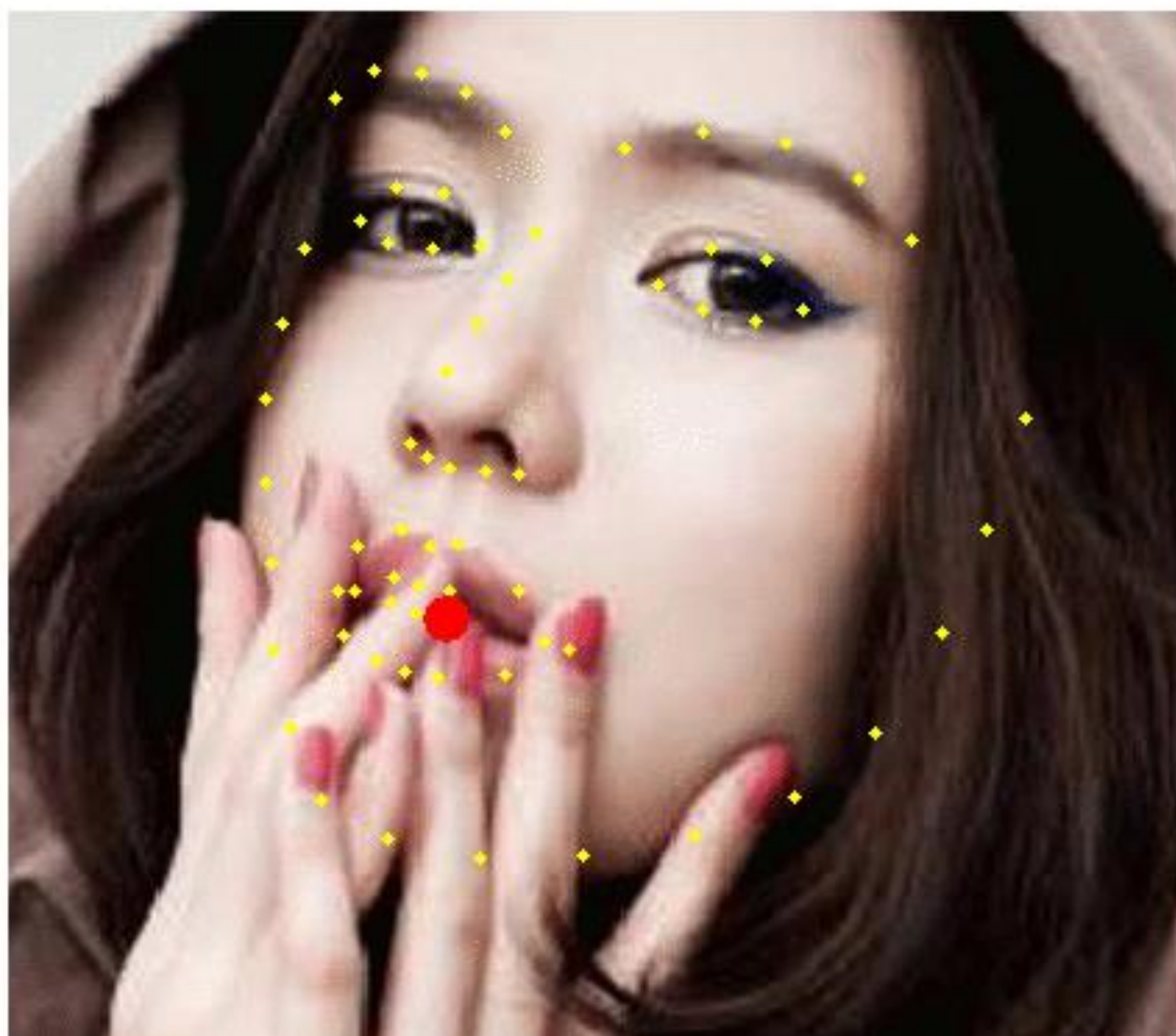


Iter 01

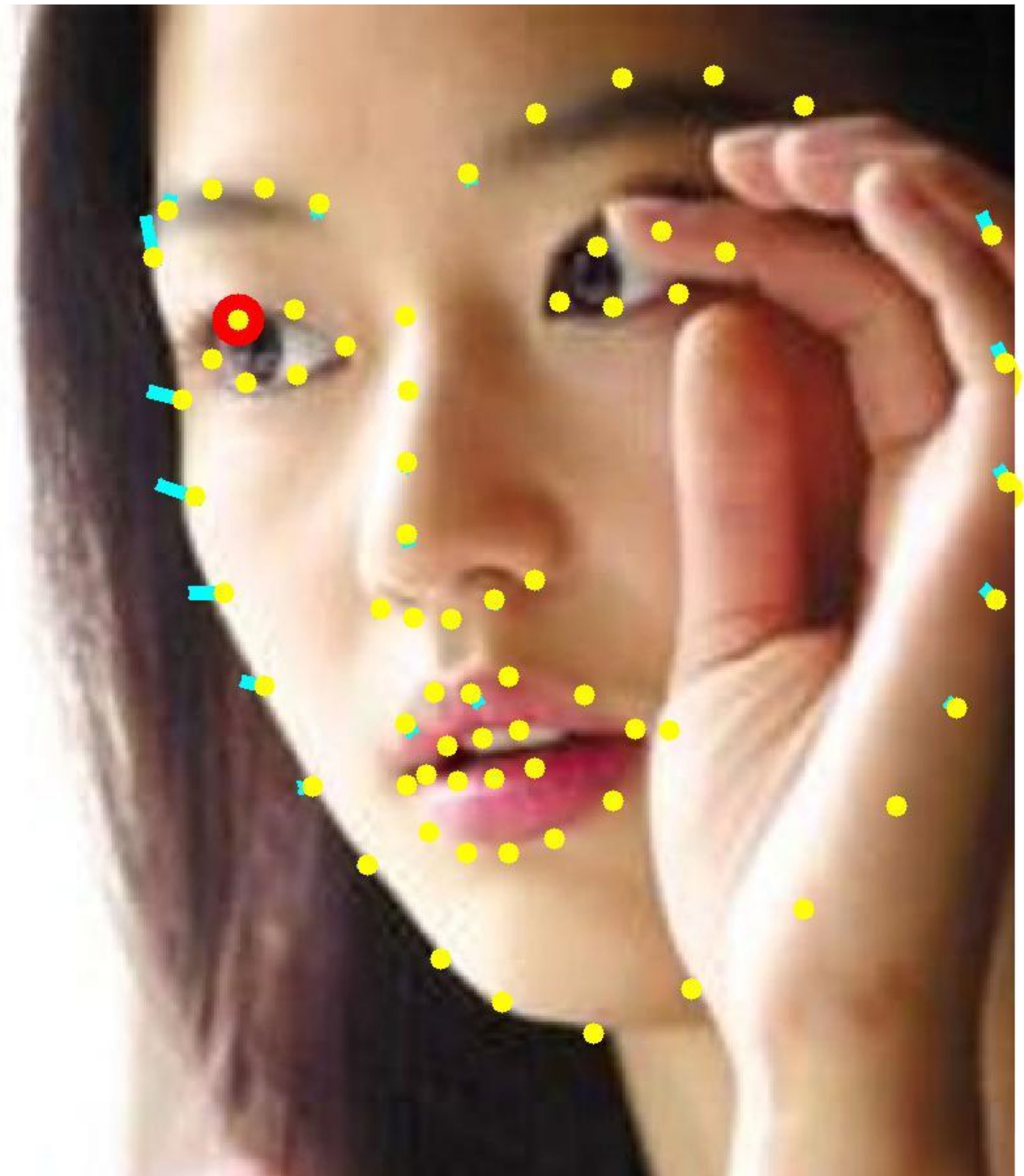


# Sample Attentive Refinement

Iter 01



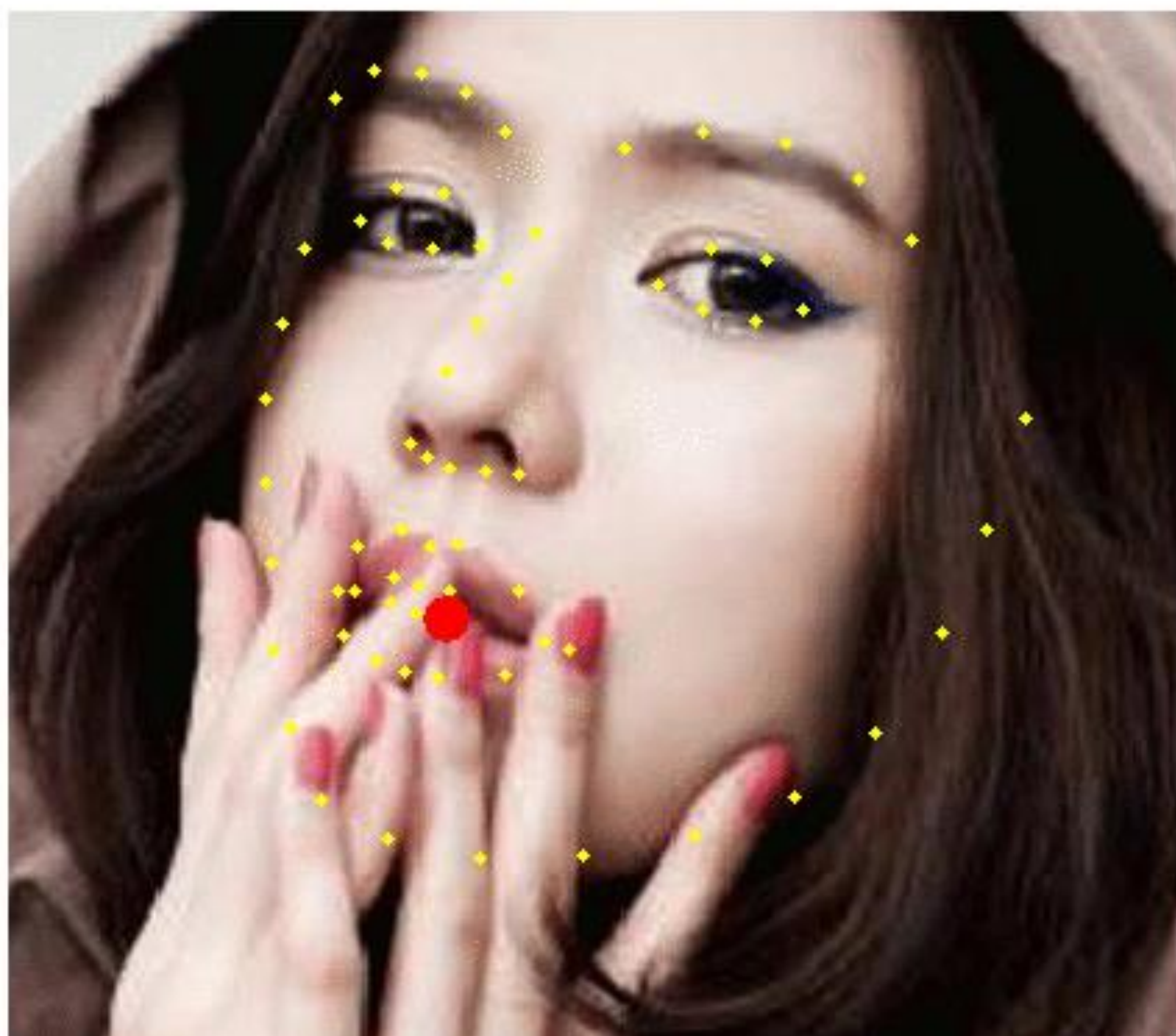
Iter 01



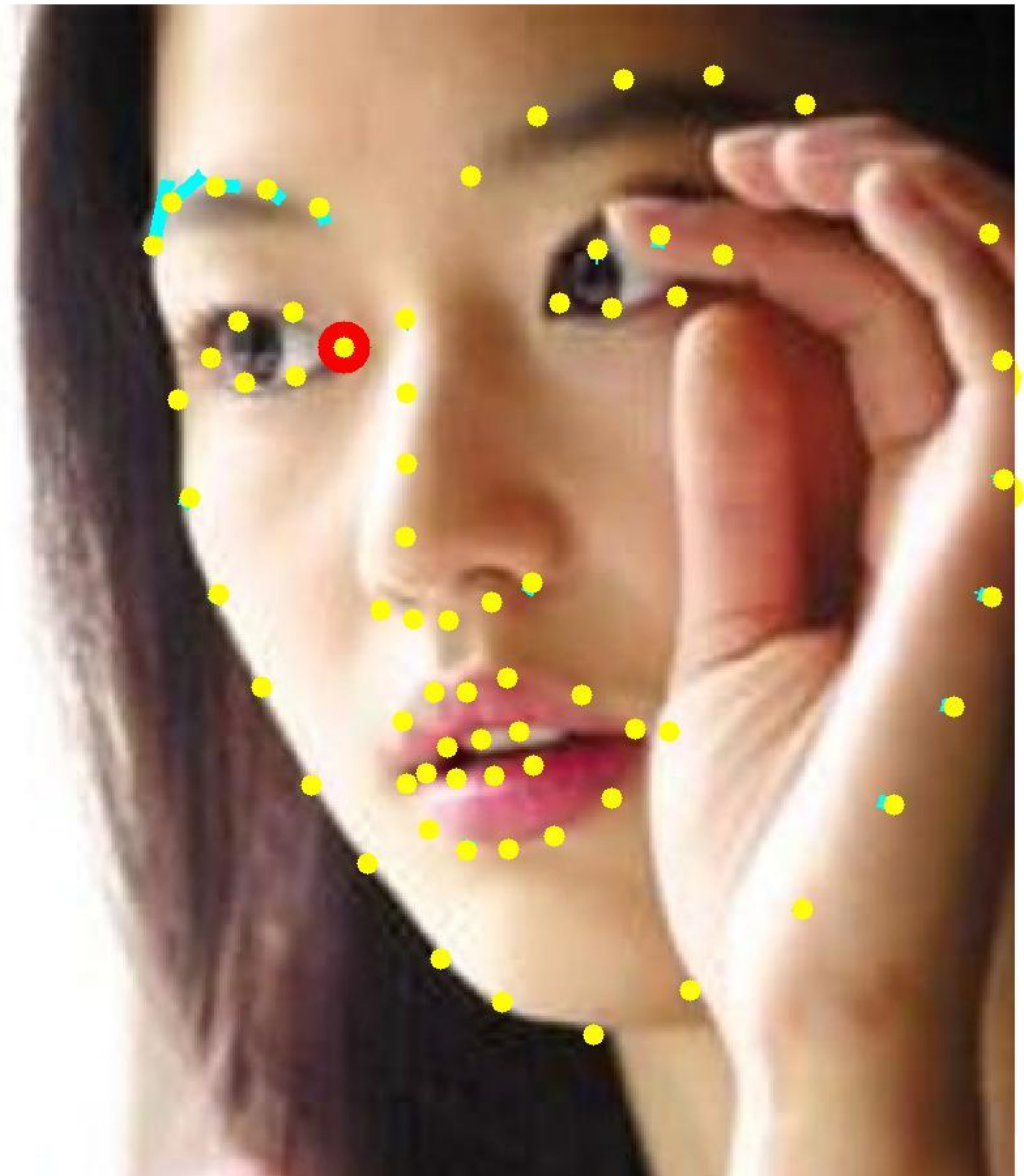


# Sample Attentive Refinement

Iter 01



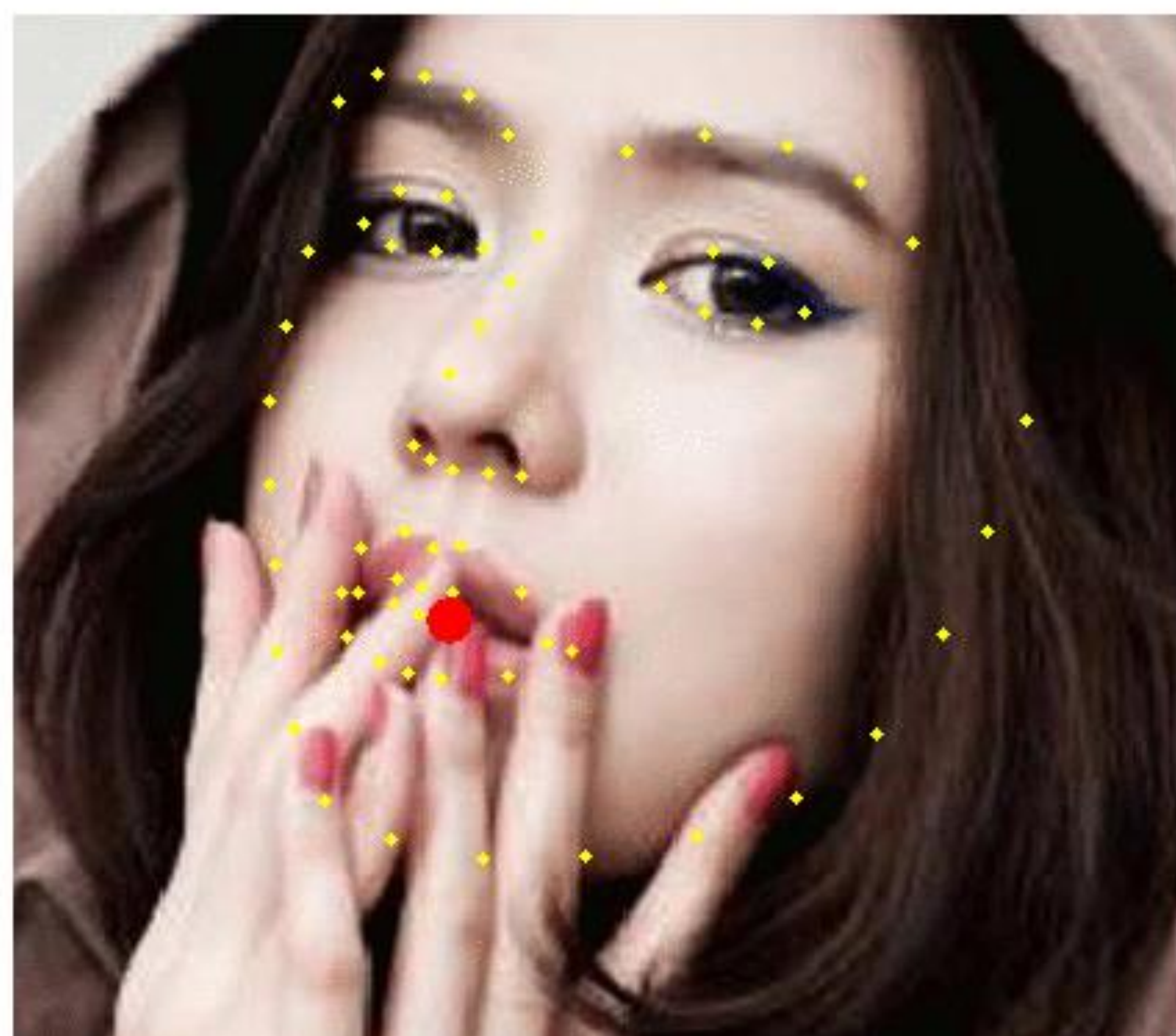
Iter 01



1 2 3 4 5

# Sample Attentive Refinement

Iter 01



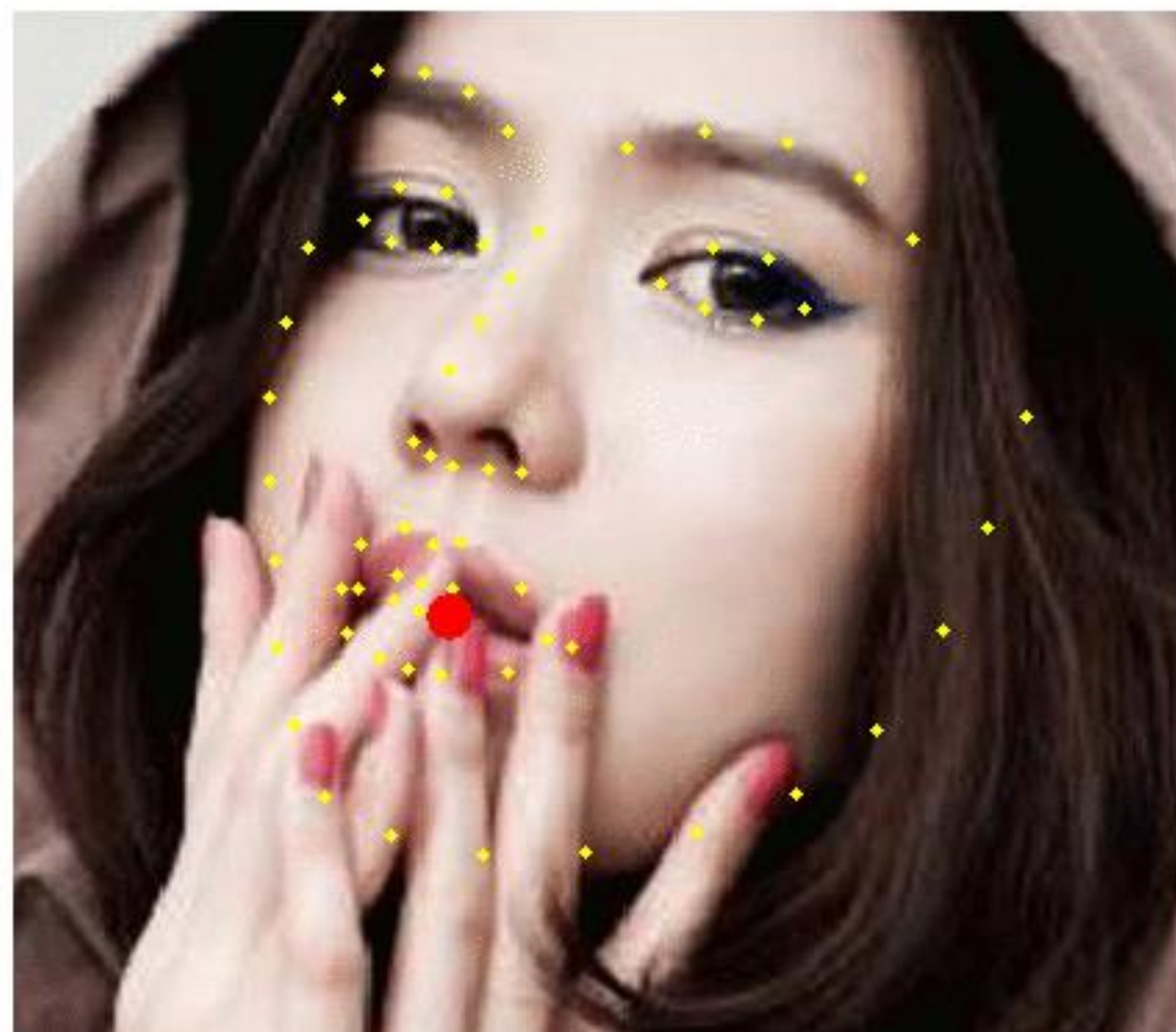
Iter 01



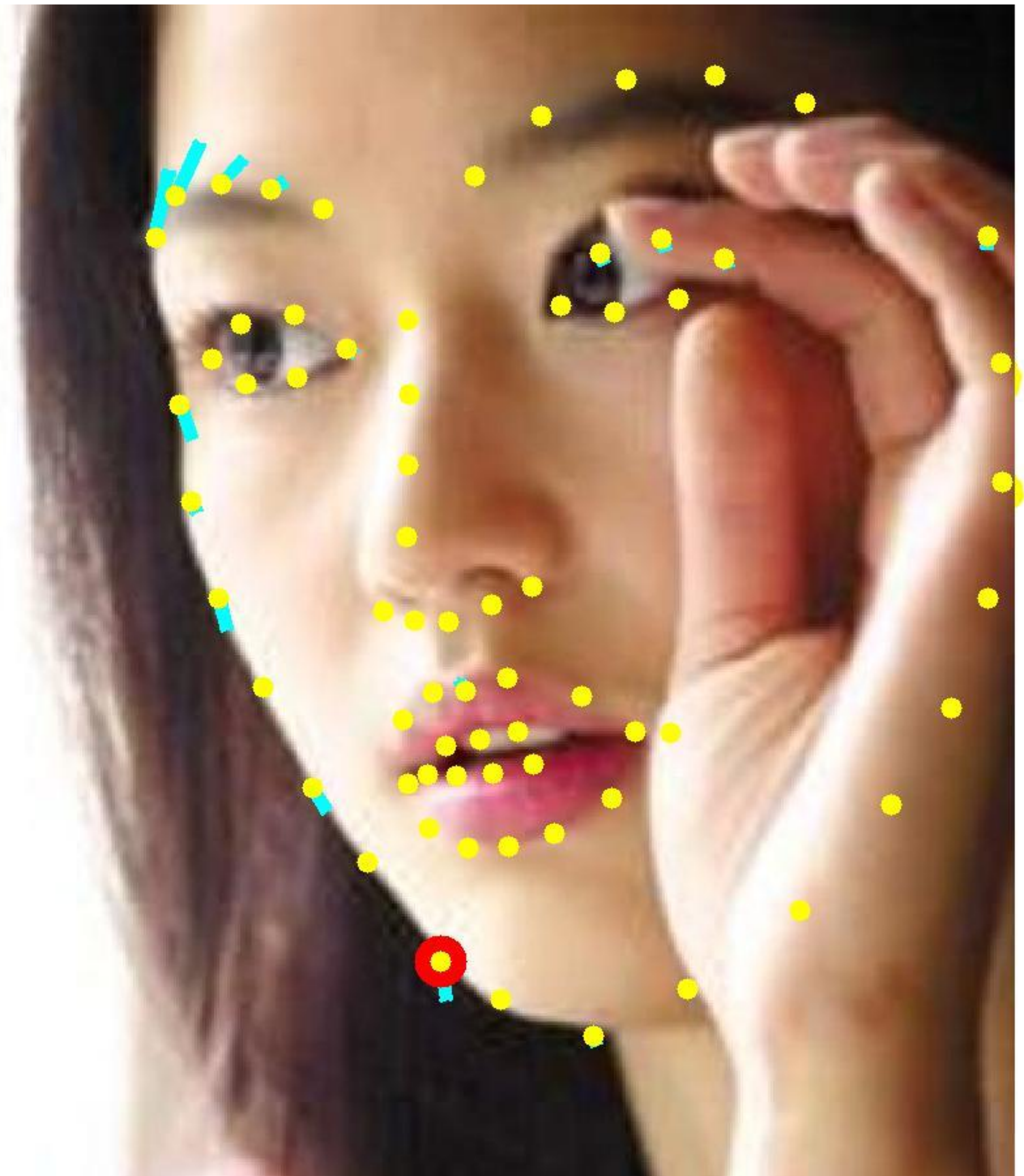
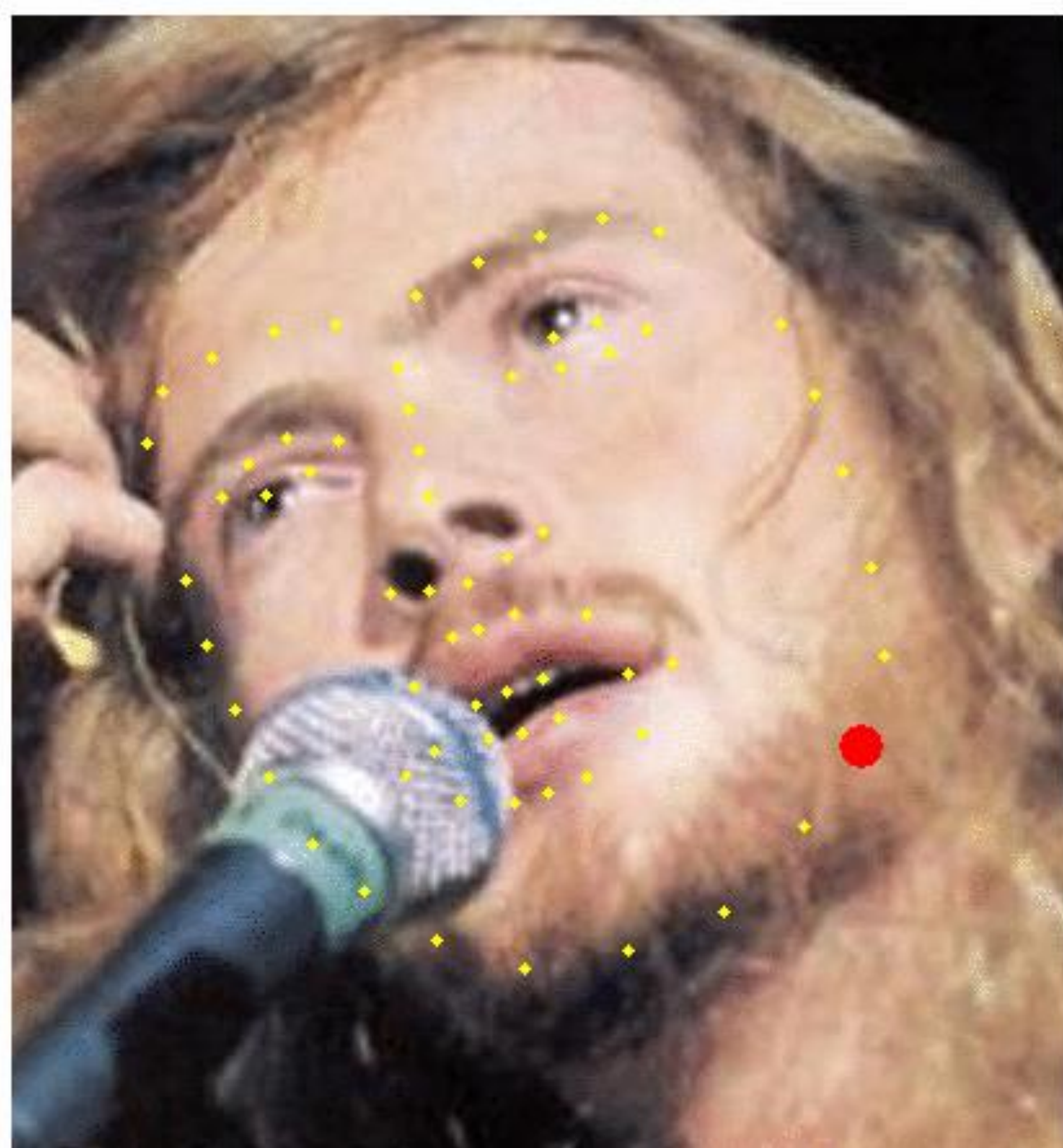
1 2 3 4 5 6

# Sample Attentive Refinement

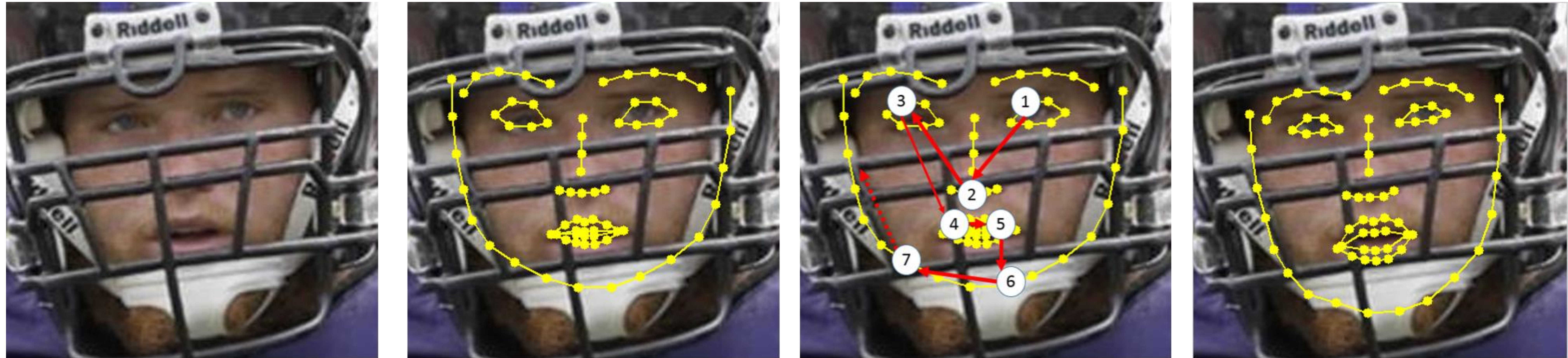
Iter 01



Iter 01



1 2 3 4 5 6 7 8 9 10 11 12 13 14 15



# Q&A

**Poster Session: O-1A-04**