

# **Let the Shape Speak - Discriminative Face Alignment using Conjugate Priors**

---

Pedro Martins, Rui Caseiro, João F. Henriques, Jorge Batista

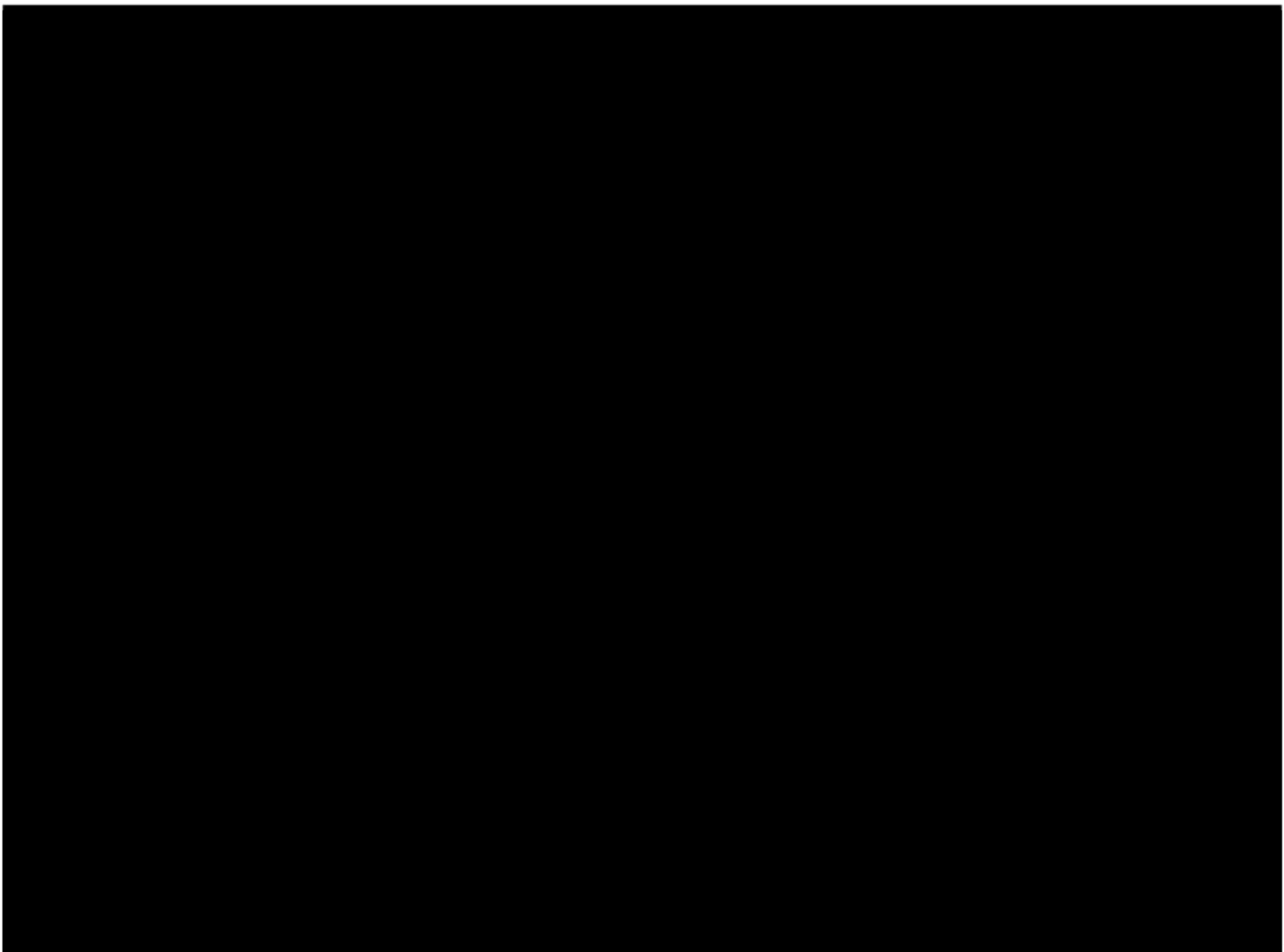
<http://www.isr.uc.pt/~pedromartins>



Institute of Systems and Robotics  
University of Coimbra  
Portugal

British Machine Vision Conference 2012

# Demo Video

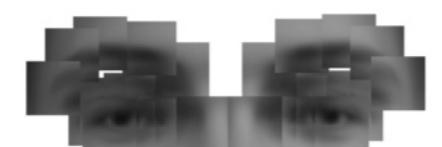
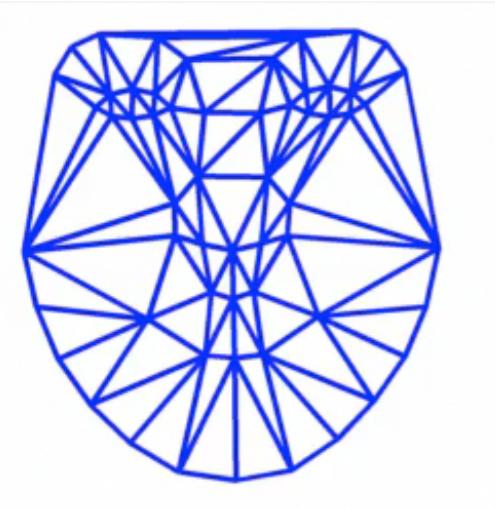


# Introduction

---

- **Goal:** Face alignment in unseen images.
- Closely related to Constrained Local Models (CLM) and Active Shape Models (ASM), where a set of local detectors is constrained to lie in the subspace spanned by a Point Distribution Model (PDM).
  - Two step fitting approach:
    - (1) Local search using the local detectors (response maps for each landmark)
    - (2) Global optimization strategy that finds the PDM parameters that jointly maximize all the detections.
- **Proposed Work:** New Bayesian global optimization strategy where the prior distribution encodes the transition of the PDM parameters.
- The prior distribution is modeled using recursive Bayesian estimation. The mean and covariance are assumed to be unknown and treated as random variables.

# Related Work - Parametric Image Alignment



**Point Distribution Model**

$$s = \mathcal{S}(s_0 + \Phi b; q)$$

↑  
Shape Parameters

Pose Parameters

- **Generative / Holistic methods**

- Active Appearance Models (AAM)  
T.F.Cootes, G.J.Edwards, C.J.Taylor - ECCV 98
- 3D Morphable Models (3DMM)  
V.Blanz, T.Vetter - SIGGRAPH 99
- Real Time Combined 2D+3D Active Appearance Models  
J.Xing, S.Baker, I.Matthews, T.Kanade - CVPR 2004

- **Discriminative / Patch-Based**

- Active Shape Models (ASM)  
T.F.Cootes, G.J.Edwards, C.J.Taylor - CVIU 95
- Constrained Local Model (CLM)  
D.Cristinacce, T.F.Cootes - BMVC 2006
- Convex Quadratic Fitting (CQF)  
Y.Wang, S.Lucey, J.Cohn - CVPR 2008
- Bayesian Constrained Local Model (BCLM)  
U.Paquet - CVPR 2009
- Subspace Constrained Mean-Shifts (SCMS)  
J.Saragih, S.Lucey, J.Cohn - ICCV 2009

# The Alignment Goal

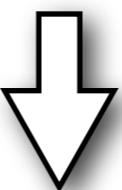
---

- Given a shape observation vector ( $\mathbf{y}$ ), find the optimal set of shape (and pose) parameters ( $\mathbf{b}$ ) that maximize the posterior probability

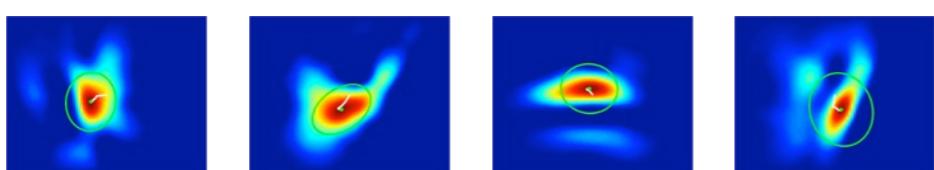
$$\mathbf{b}^* = \arg \max_{\mathbf{b}} p(\mathbf{b}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{b})p(\mathbf{b})$$

- Assuming:
  - Conditional independence between landmarks
  - Close to a solution

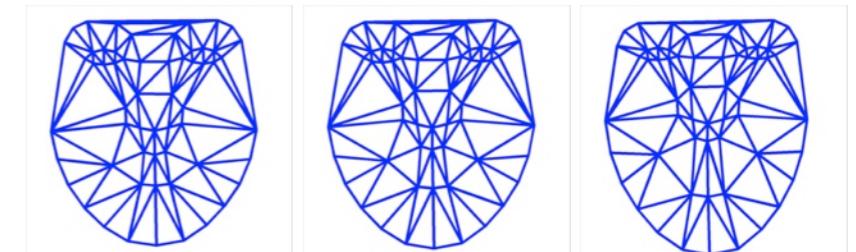
$$p(\mathbf{b}|\mathbf{y}) \propto \left( \prod_{i=1}^v p(\mathbf{y}_i|\mathbf{b}) \right) p(\mathbf{b}|\mathbf{b}_{k-1}^*)$$



**Likelihood** from the local detectors



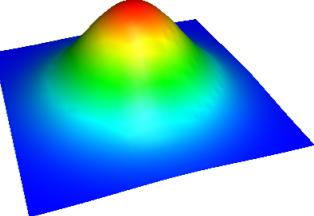
**Prior** on how parameters change

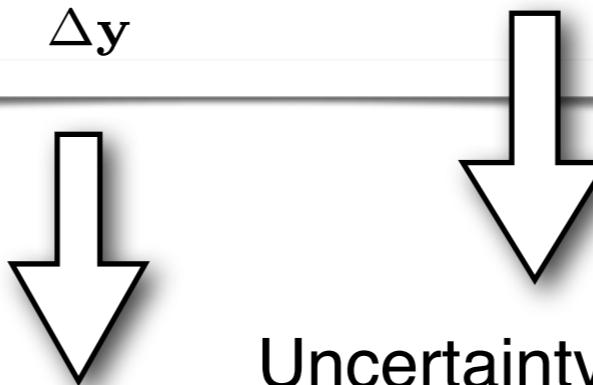


# The Likelihood Term

- Convex energy function:

Observed shape ( $\mathbf{y}$ )

$$p(\mathbf{y}|\mathbf{b}) \propto \exp \left( -\frac{1}{2} \underbrace{(\mathbf{y} - (\mathbf{s}_0 + \Phi \mathbf{b}))^T}_{\Delta \mathbf{y}} \Sigma_{\mathbf{y}}^{-1} (\mathbf{y} - (\mathbf{s}_0 + \Phi \mathbf{b})) \right)$$




Uncertainty covariance

Difference between the  
**observed** and the **mean shape**

$$\Sigma_{\mathbf{y}}$$

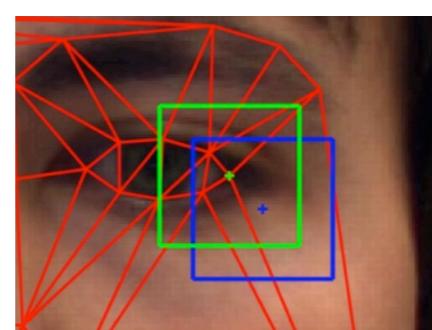
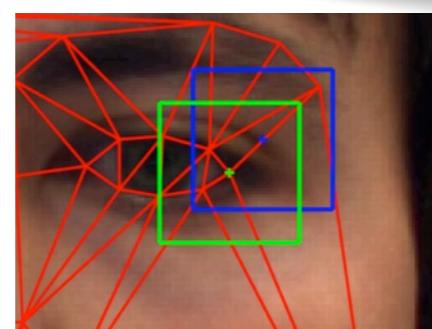
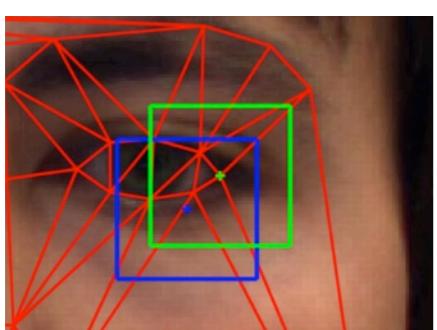
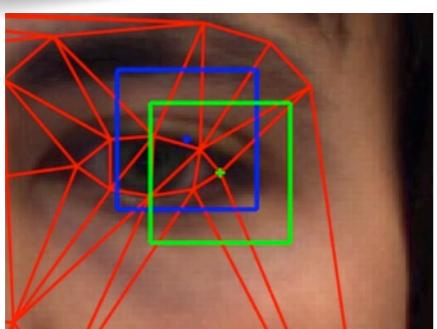
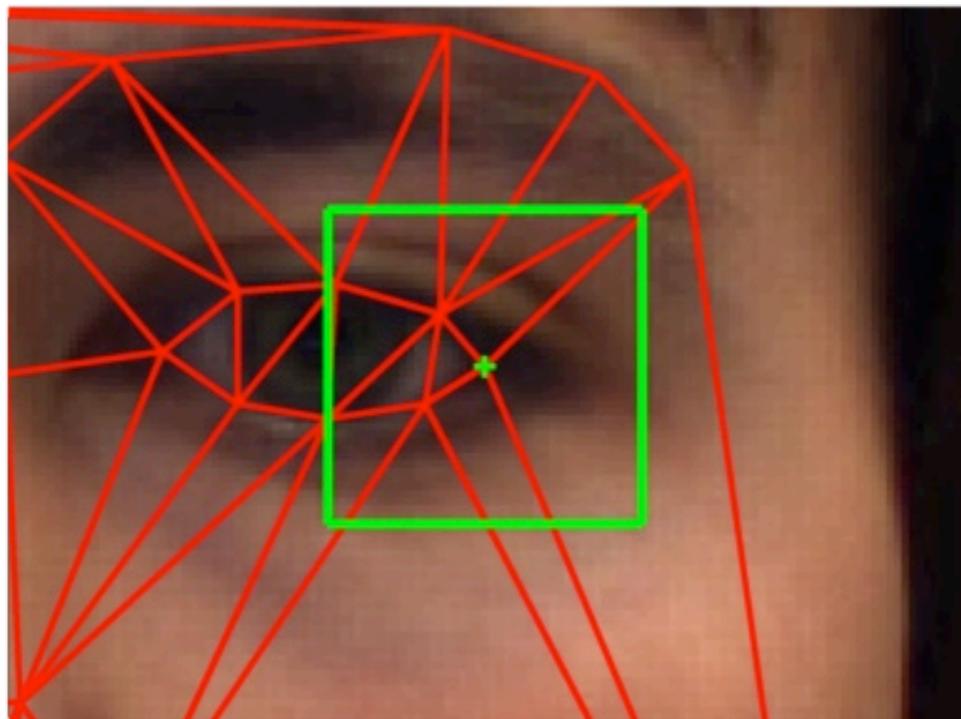
				0
				...
0				

The likelihood follow a Gaussian distribution

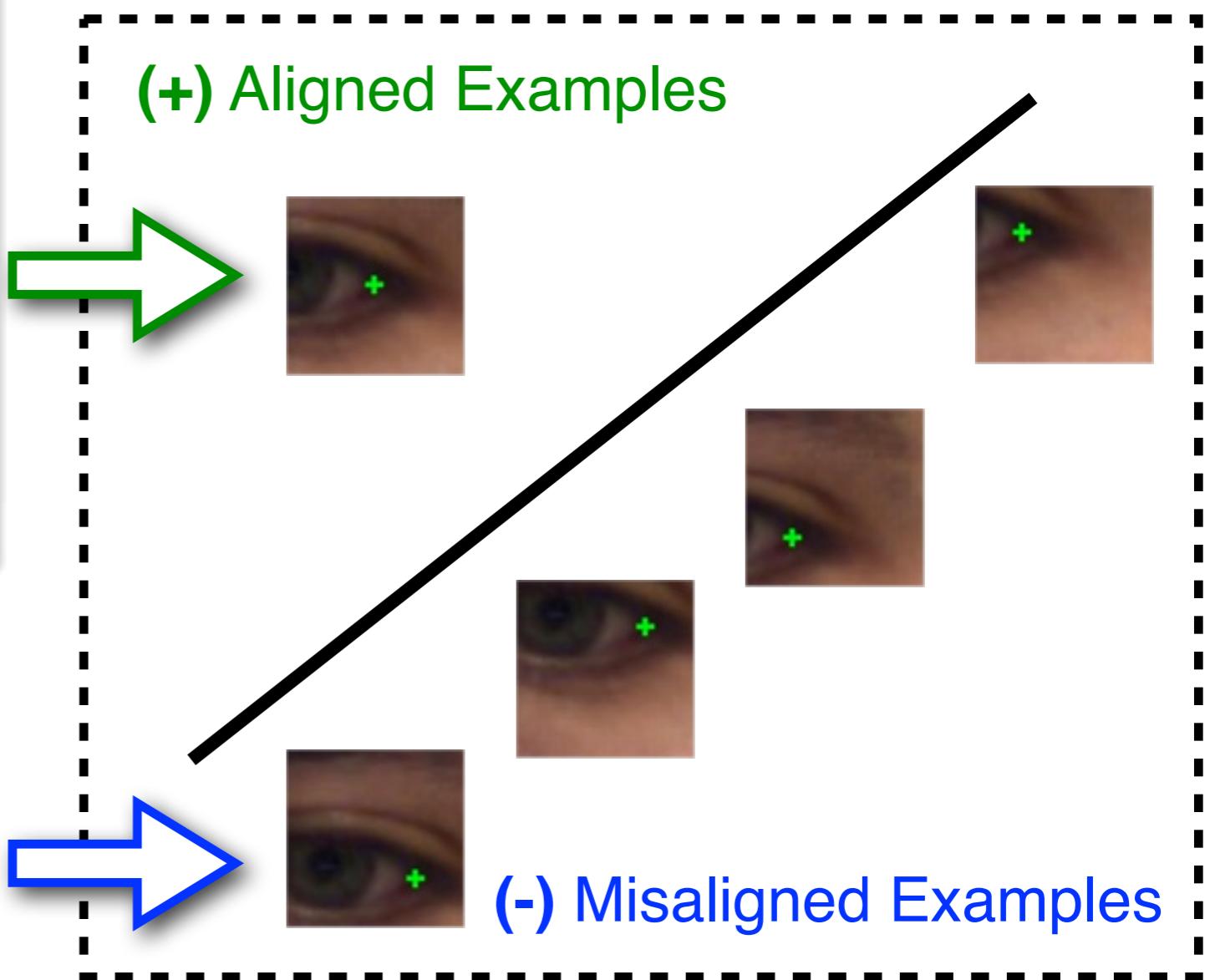
$$p(\mathbf{y}|\mathbf{b}) \propto \mathcal{N}(\Delta \mathbf{y} | \Phi \mathbf{b}, \Sigma_{\mathbf{y}})$$

$2v \times 2v$  Block diagonal

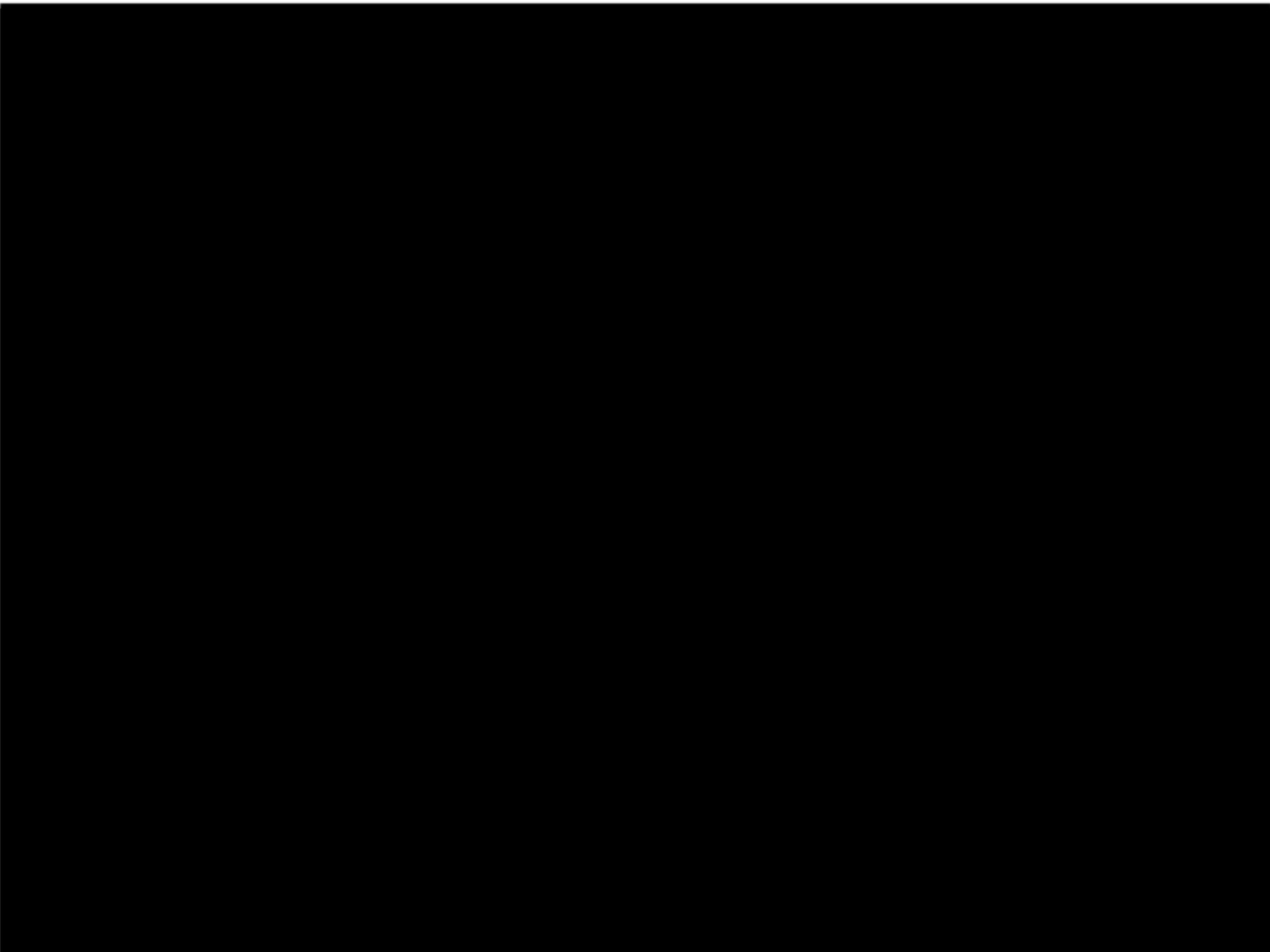
# Local Landmark Detectors



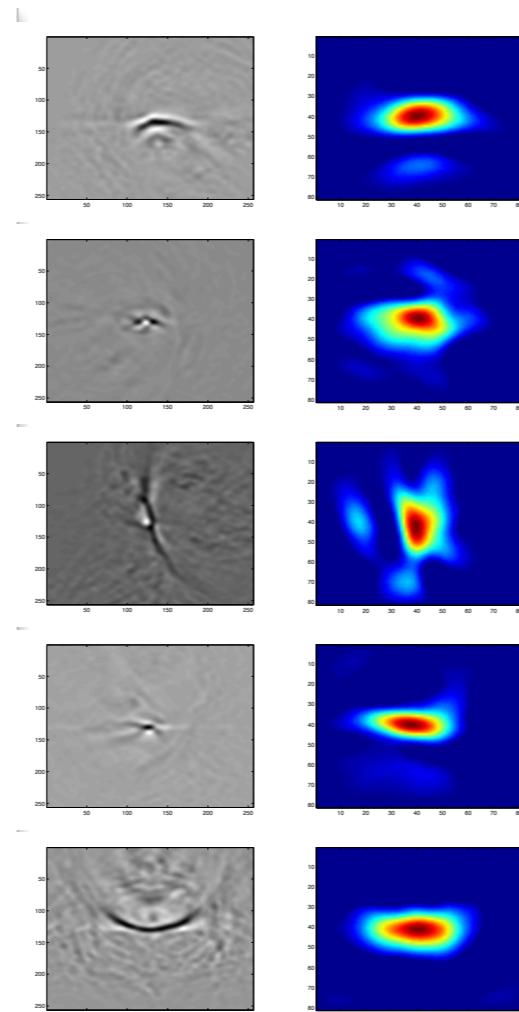
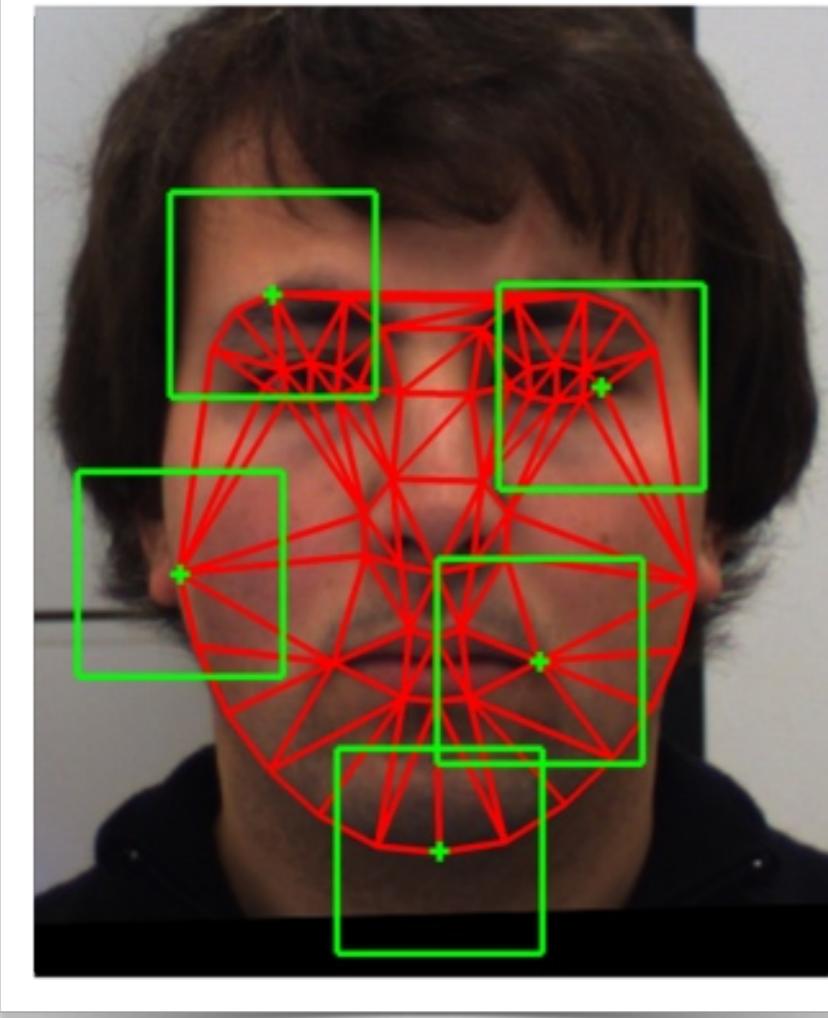
## Linear SVM



## Local Detectors



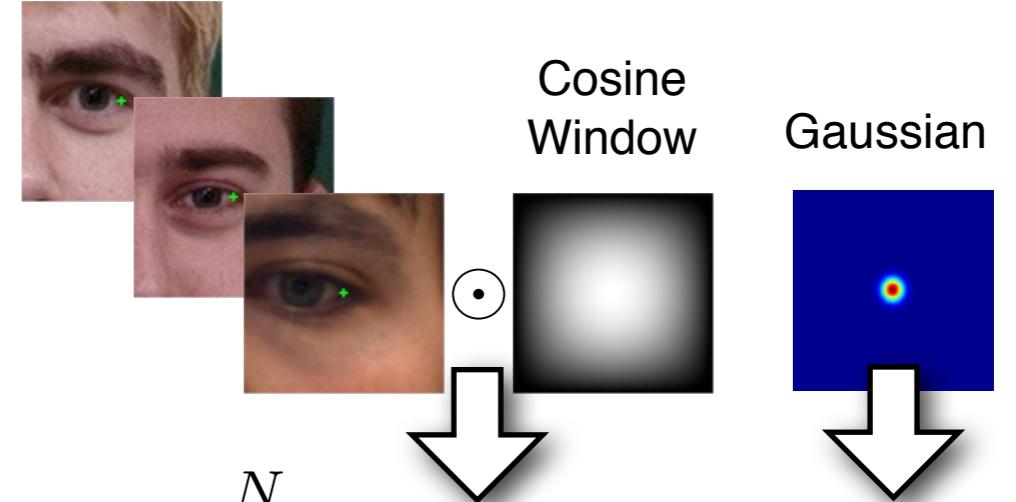
# Local Landmark Detectors - MOSSE Filters



$$\mathcal{F}^{-1}\{\mathbf{H}_i^*\} \quad \mathcal{D}_i^{\text{MOSSE}}(\mathbf{I}(\mathbf{y}_i))$$

- Correlation in Fourier Domain

$$\mathbf{G} = \mathcal{F}\{\mathbf{I}\} \odot \mathbf{H}^*$$



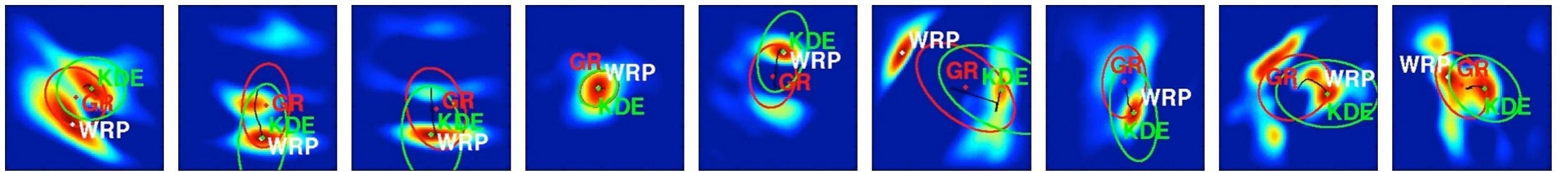
$$\min_{\mathbf{H}^*} \sum_{j=1}^N (\mathcal{F}\{\mathbf{I}_j\} \odot \mathbf{H}^* - \mathbf{G}_j)^2$$

## MOSSE Filter

$$\mathbf{H}^* = \frac{\sum_{j=1}^N \mathbf{G}_j \odot \mathcal{F}\{\mathbf{I}_j\}^*}{\sum_{j=1}^N \mathcal{F}\{\mathbf{I}_j\} \odot \mathcal{F}\{\mathbf{I}_j\}^*}$$

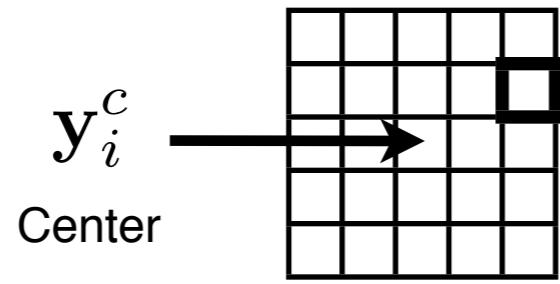
- Visual object tracking using adaptive correlation filters
- D.Bolme, J.Beveridge, B.Draper, Y.Lui, CVPR 2010

# Local Optimization Strategies



$p_i(\mathbf{z}_i)$

Prob.  $\mathbf{z}_i$  is aligned



$\mathbf{z}_i = (x_i, y_i)$

Pixel candidate to  $i^{\text{th}}$  landmark location

Patches under occlusion

## Weighted Peak Response (WPR)

$$\mathbf{y}_i^{\text{WPR}} = \max_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} (p_i(\mathbf{z}_i))$$

$$\Sigma_{\mathbf{y}_i}^{\text{WPR}} = \text{diag}(p_i(\mathbf{y}_i^{\text{WPR}})^{-1})$$

## Gaussian Response (GR)

$$\mathbf{y}_i^{\text{GR}} = \frac{1}{d} \sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} p_i(\mathbf{z}_i) \mathbf{z}_i \quad d = \sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} p_i(\mathbf{z}_i)$$

$$\Sigma_{\mathbf{y}_i}^{\text{GR}} = \frac{1}{d-1} \sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} p_i(\mathbf{z}_i) (\mathbf{z}_i - \mathbf{y}_i^{\text{GR}})(\mathbf{z}_i - \mathbf{y}_i^{\text{GR}})^T$$

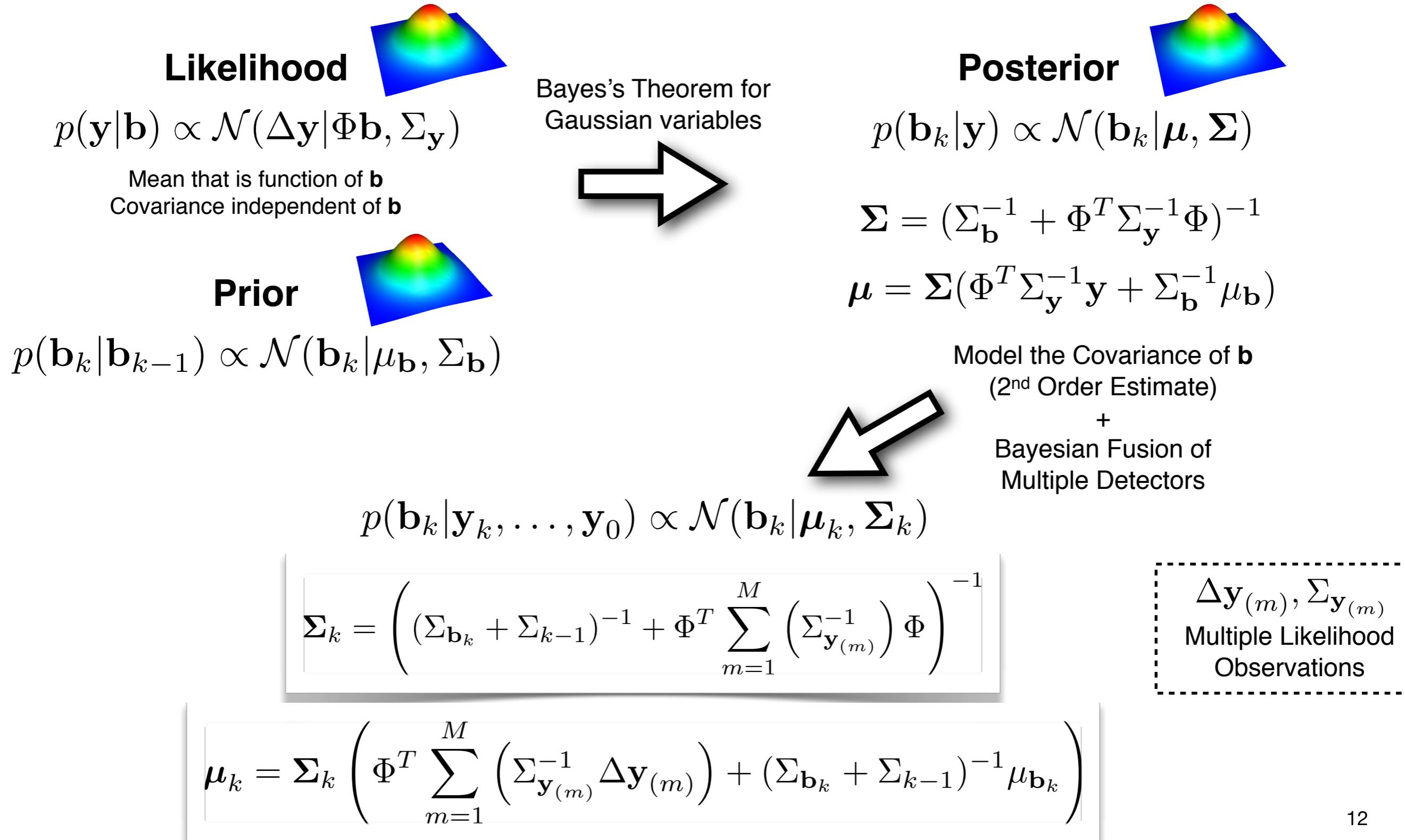
## Kernel Density Estimator (KDE)

$$\mathbf{y}_i^{\text{KDE}(\tau+1)} \leftarrow \frac{\sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} \mathbf{z}_i p_i(\mathbf{z}_i) \mathcal{N}(\mathbf{y}_i^{\text{KDE}(\tau)} | \mathbf{z}_i, \sigma_{h_j}^2 \mathbf{I}_2)}{\sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} p_i(\mathbf{z}_i) \mathcal{N}(\mathbf{y}_i^{\text{KDE}(\tau)} | \mathbf{z}_i, \sigma_{h_j}^2 \mathbf{I}_2)}$$

$$\Sigma_{\mathbf{y}_i}^{\text{KDE}} = \frac{1}{d-1} \sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} p_i(\mathbf{z}_i) (\mathbf{z}_i - \mathbf{y}_i^{\text{KDE}})(\mathbf{z}_i - \mathbf{y}_i^{\text{KDE}})^T$$

# KDE Demo Video

# MAP Global Alignment



# The Prior Term

---

- Mean and Covariance ( $\mu_b, \Sigma_b$ ) are assumed to be **unknown** and modeled as **random variables**.

$$p(\mathbf{b}_k | \mathbf{b}_{k-1}) \propto \mathcal{N}(\mathbf{b}_k | \boxed{\mu_b, \Sigma_b})$$

Unknown

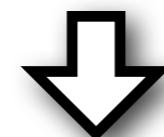
Observable vector  $\mathbf{b}$

**Bayes Theorem:**  $p(\mu_b, \Sigma_b | \mathbf{b}) \propto p(\mathbf{b} | \mu_b, \Sigma_b) p(\mu_b, \Sigma_b)$

Joint Posterior  
Normal Inverse-Wishart

Joint Prior  
Normal Inverse-Wishart

Parameters



Degrees of freedom  $v_k = v_{k-1} + m$

Number of measurements  $\kappa_k = \kappa_{k-1} + m$

Mean

$$\theta_k = \frac{\kappa_{k-1}}{\kappa_{k-1} + m} \theta_{k-1} + \frac{m}{\kappa_{k-1} + m} \bar{\mathbf{b}}$$

m - Number of samples to update the model

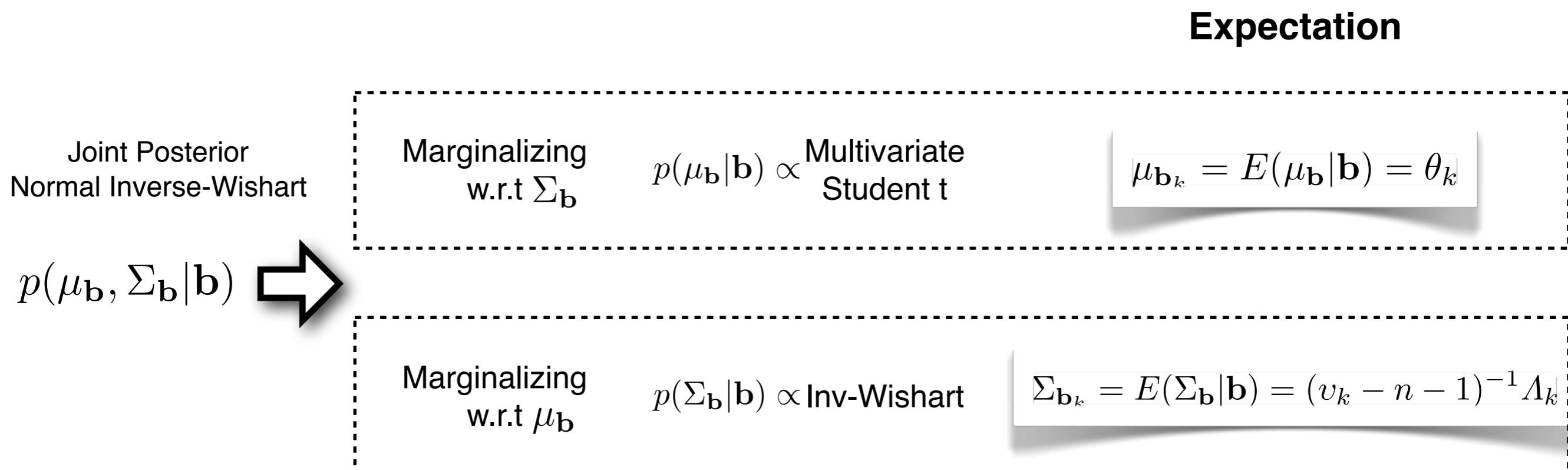
$\bar{\mathbf{b}}$  - Mean of all samples

Scale matrix

$$\Lambda_k = \Lambda_{k-1} + \frac{\kappa_{k-1}m}{\kappa_{k-1} + m} (\bar{\mathbf{b}} - \theta_{k-1})(\bar{\mathbf{b}} - \theta_{k-1})^T$$

# The Prior Term

- Using the expectation of marginal posterior distributions as the model parameters update.



The Prior distribution is **continuously kept up to date**

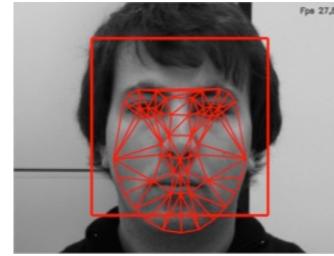
# The Algorithm

**Precompute:**

PDM:  $\mathbf{s}_0, \Phi, \Psi, \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

Initial estimate

$(\mathbf{b}_0, \Sigma_0), (\mathbf{q}_0, \Sigma_0^q)$

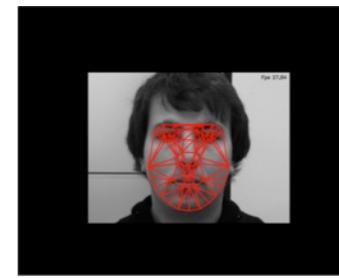
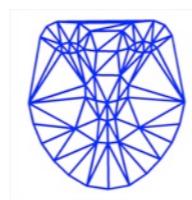


MOSSE Filters:  $\mathbf{H}_i^*$   $i=1, \dots, v$



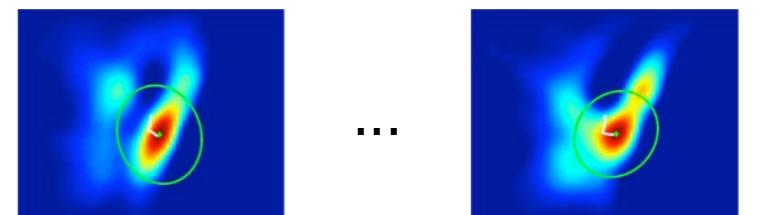
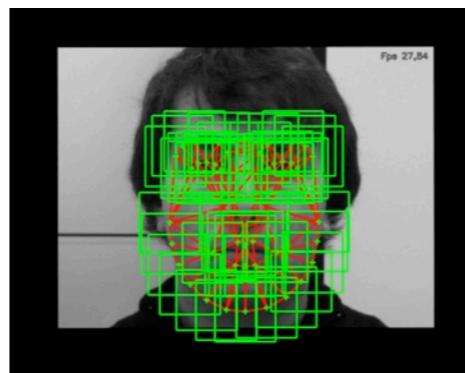
**for**  $k=1:1:\text{MaxIterations}$

Warp Image to the base mesh, using the current pose parameters



Generate current shape  $\mathbf{s} = \mathcal{S}(\mathbf{s}_0 + \Phi \mathbf{b}_k; \mathbf{q}_k)$

**WPR, GR or KDE Local Strategy**



**for**  $i=1:1:\text{Landmarks}$

Evaluate detectors response

Find the likelihood parameters  $\mathbf{y}_i, \Sigma_{\mathbf{y}_i}$

**end**

Estimate the shape/pose parameters:

Update the parameters of Normal Inv-Wishart distribution  $\rightarrow v_k, \kappa_k, \theta_k, \Lambda_k$

Expectation of the prior shape parameters  $\longrightarrow \mu_{\mathbf{b}_k} = \theta_k, \Sigma_{\mathbf{b}_k} = (v_k - n - 1)^{-1} \Lambda_k$

Evaluate the **global** shape parameters and the covariance  $\rightarrow \mu_k, \Sigma_k$

**end**

# Hierarchical Search (BASM-KDE-H)

- When response maps are approximated by KDE representations.

Mean-Shift Landmark Update

$$\mathbf{y}_i^{\text{KDE}(\tau+1)} \leftarrow \frac{\sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} \mathbf{z}_i p_i(\mathbf{z}_i) \mathcal{N}(\mathbf{y}_i^{\text{KDE}(\tau)} | \mathbf{z}_i, \sigma_{h_j}^2 \mathbf{I}_2)}{\sum_{\mathbf{z}_i \in \Omega_{\mathbf{y}_i^c}} p_i(\mathbf{z}_i) \mathcal{N}(\mathbf{y}_i^{\text{KDE}(\tau)} | \mathbf{z}_i, \sigma_{h_j}^2 \mathbf{I}_2)}$$

Bandwidth schedule

$$\sigma_h^2 = (15, 10, 5, 2)$$

## Standard Search

```
for k=1:1:MaxIterations
```

```
    for i=1:1:v (LandMarks)
```

Evaluate de Detectors Response

Mean-Shift Landmark Update  $\sigma_h^2 = (15, 10, 5, 2)$

```
end
```

```
end
```

Global Optimization

## Hierarchical Search

```
for k=1:1:MaxIterations
```

$\sigma_h^2 = (15, 10, 5, 2)$

```
    for i=1:1:v (LandMarks)
```

Evaluate de Detectors Response

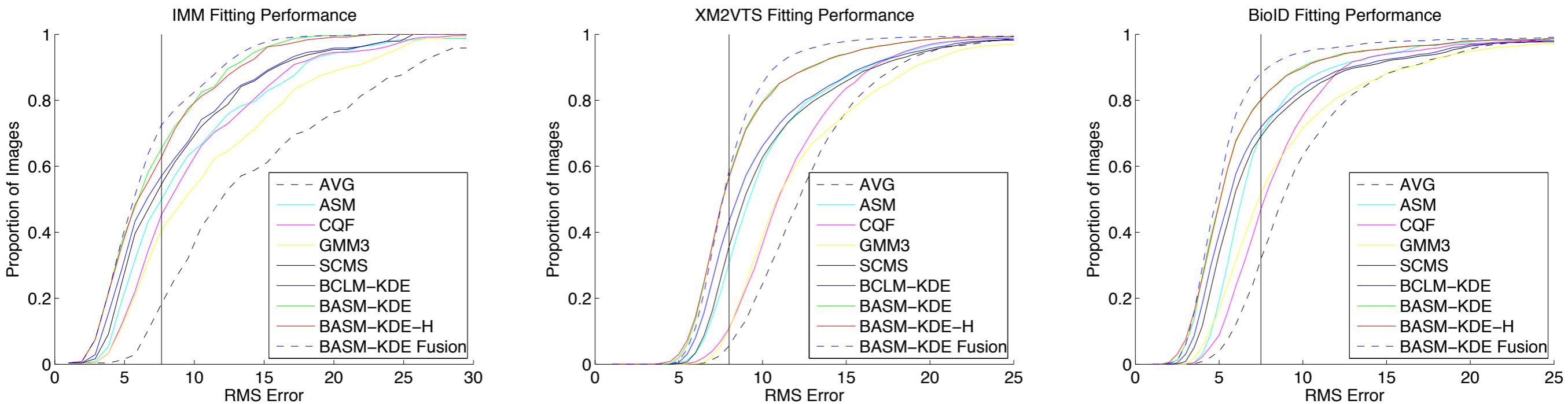
Mean-Shift Landmark Update  $\sigma_{h_j}^2$

```
end
```

```
end
```

Global Optimization

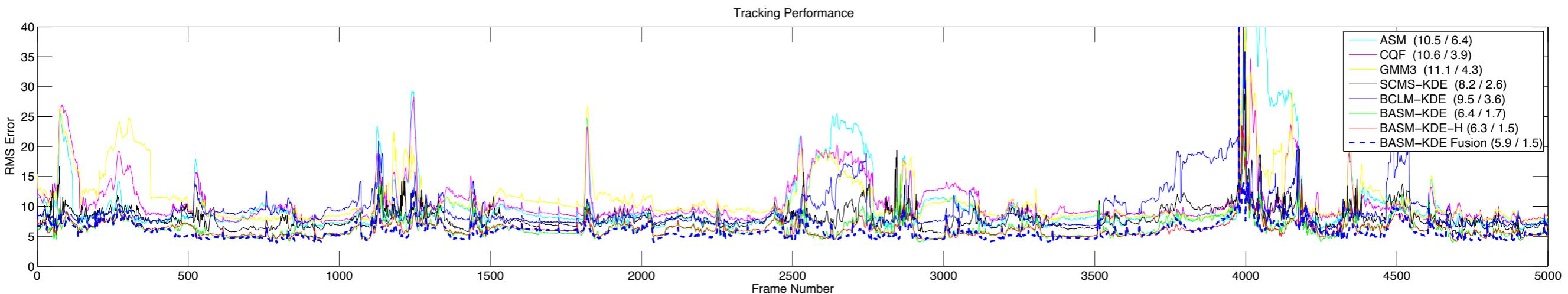
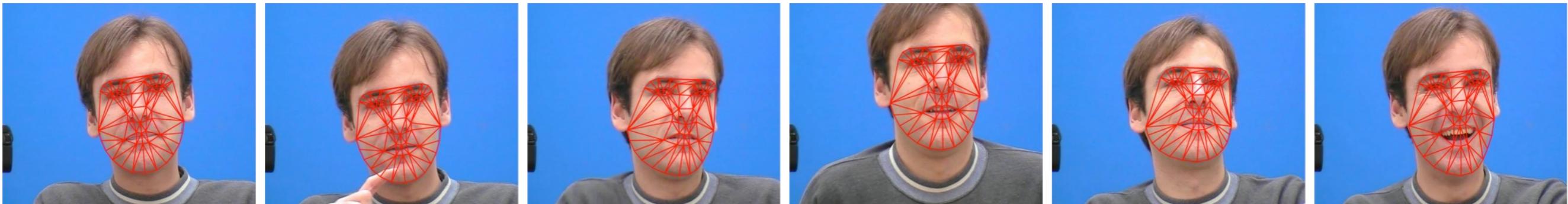
# Evaluation Results



Reference 7.5 RMS	IMM (240 images)	XM2VTS (2360 images)	BioID (1521 images)
ASM	50.0	30.7	70.0
BASM-WPR (our method)	<b>58.4</b> (+8.4)	<b>47.4</b> (+16.7)	<b>77.1</b> (+7.1)
CQF	45.4	10.9	47.0
GMM3	40.8 (-4.6)	10.4 (-0.5)	51.7 (+4.7)
BCLM-GR	48.3 (+2.9)	15.9 (+5.0)	54.2 (+7.2)
BASM-GR (our method)	<b>51.8</b> (+6.4)	<b>19.7</b> (+8.8)	<b>63.5</b> (+16.5)
SCMS-KDE	54.6	35.7	69.0
BCLM-KDE	57.1 (+2.5)	43.4 (+7.7)	71.9 (+2.9)
BASM-KDE (our method)	<b>65.4</b> (+10.8)	<b>57.0</b> (+21.3)	<b>80.3</b> (+11.3)
BASM-KDE-H (our method)	64.0 (+9.4)	56.6 (+20.9)	79.9 (+10.9)
BASM-KDE Fusion of 2 Detectors	<b>72.5</b> (+17.9)	<b>58.7</b> (+23.0)	<b>88.2</b> (+19.2)

# Tracking Performance

- FGNET Talking Face Video Sequence

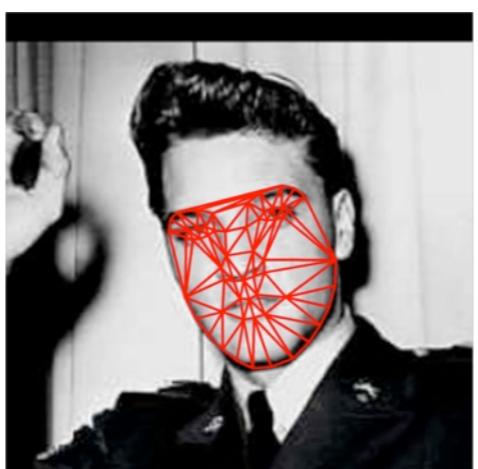
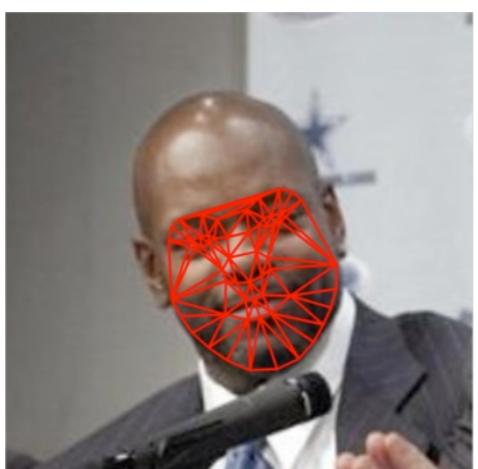
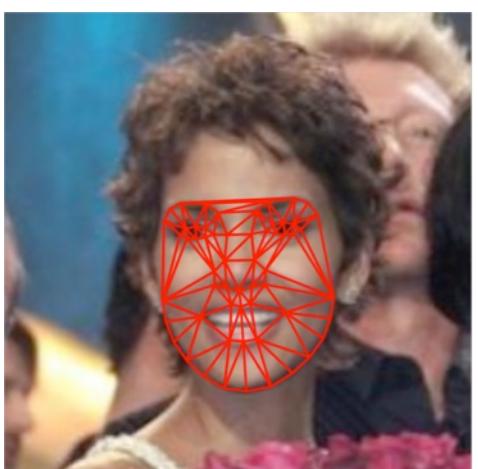
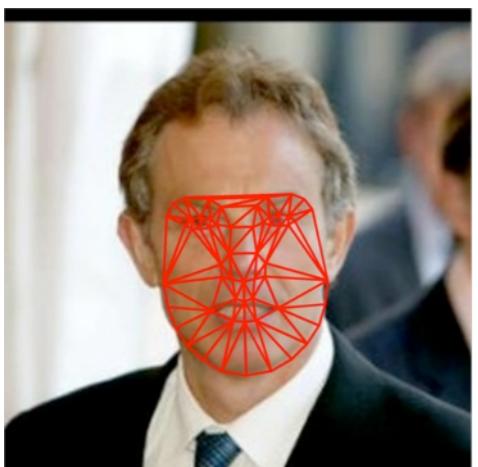
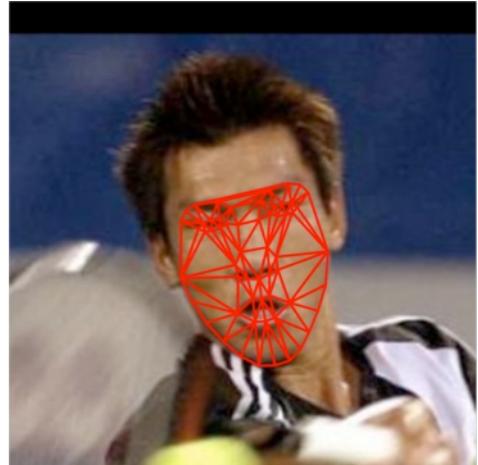
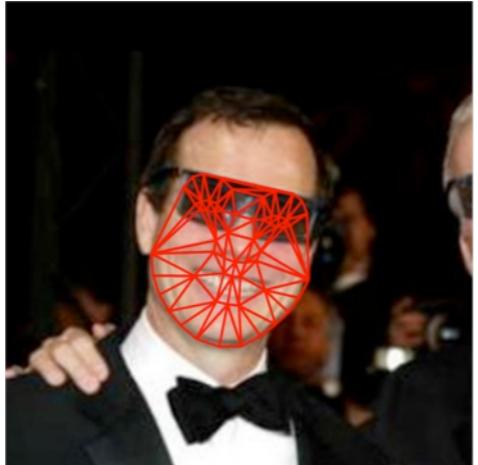
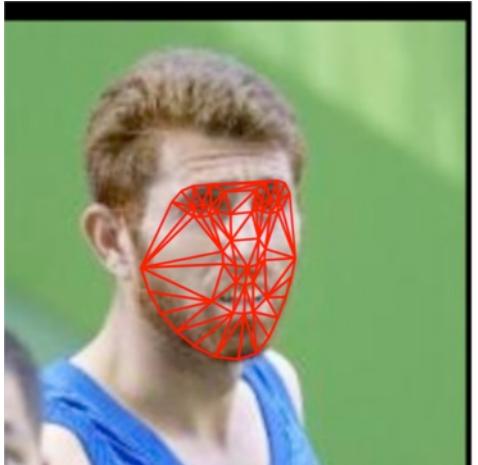


RMS Error	ASM	CQF	GMM3	SCMS-KDE	BCLM-KDE	BASM-KDE	BASM-KDE-H	BASM-KDE Fusion
Mean	10.5	10.6	11.1	8.2	9.5	6.4	6.3	5.9
Standard Deviation	6.4	3.9	4.3	2.6	3.6	1.7	1.5	1.5

# Tracking Performance

# Labeled Faces in the Wild (LFW)

---



# Conclusions

---

- Bayesian formulation for aligning faces in unseen images.
- New global optimization strategy infers both shape and pose parameters, in MAP sense, by explicitly modeling the prior distribution.
- The prior distribution is continuously kept up to date.
- Recursive Bayesian estimation is used to treat the mean and covariance as random variables.
- Extensive evaluations were performed on several standard datasets (IMM, XM2VTS, BiID, LFW and FGNET Talking Face) against state-of-the-art methods while using the same local detectors.

## Acknowledgements

- Work supported by the Portuguese Science Foundation (*Fundação para a Ciência e Tecnologia - FCT*) under the project grant PTDC/EIA-CCO/108791/2008.
- Pedro Martins, Rui Caseiro and João F. Henriques also acknowledge the FCT through the grants SFRH/BD/45178/2008, SFRH/BD74152/2010 and SFRH/BD/75459/2010, respectively.

# Qualitative Evaluation - Labeled Faces in the Wild

