# Fisher Vector Faces (FVF) in the Wild

**Karén Simonyan**, Omkar Parkhi,
Andrea Vedaldi, Andrew Zisserman

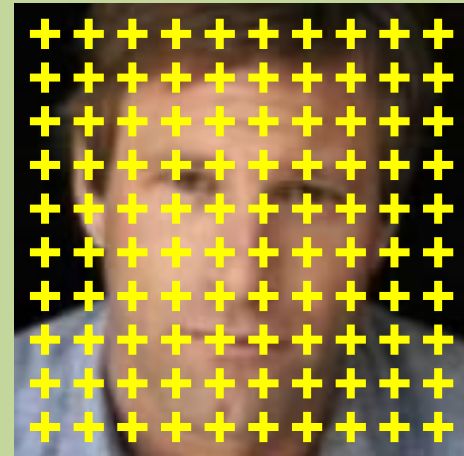Visual Geometry Group, University of Oxford

# Objective

**Face descriptor** for recognition:

- dense sampling

- relevant face parts learnt automatically

- compact and discriminative



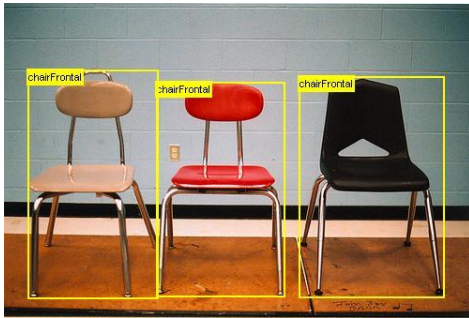**Conventional approach (describe landmarks)**



**Our approach (describe everything)**

# Motivation

- State-of-the-art image recognition pipeline:
  - **dense SIFT → Fisher vector encoding → linear SVM**
  - very competitive on (generic) image recognition tasks: Caltech 101/256, PASCAL VOC, ImageNet ILSVRC

- Can it be applied to faces? Yes!

# Application – Face Verification

## «Is it the same person in both images?»



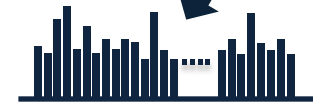SAME

DIFFERENT
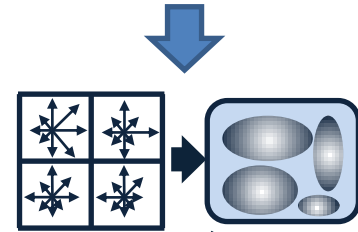
Labelled Faces in the Wild (LFW) dataset

- large-scale: 13K images, 5.7K people
- collected using Viola-Jones face detector
- high variability in appearance
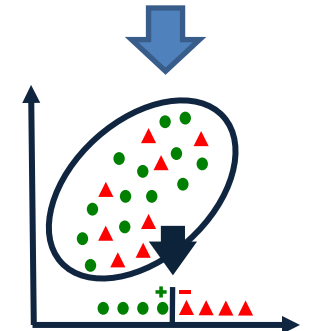- several evaluation settings (restricted, unrestricted)

# Pipeline Overview

**face image**

- Input: face image, e.g.
  - LFW + face alignment[1]
  - pre-aligned: LFW-funneled, LFW-a
  - no alignment: just Viola-Jones detection!

- Output: Fisher Vector Face descriptor (FVF)
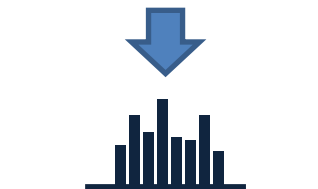  - discriminative
  - compact

**face FV extraction**

**discriminative projection**

**compact descriptor**

[1] "Taking the bite out of automatic naming of characters in TV video",
M. Everingham, J. Sivic, and A. Zisserman. IVC 2009.

# Dense Features

**face image → set of local features**



face image



face FV extraction

## Dense SIFT

- dense scale-space grid: 1 pix step, 5 scales

- 24x24 patch size

- rootSIFT[1] – explicit Hellinger kernel map

- 64-D PCA-rootSIFT

- augmented with (x,y): 66-D

discriminative projection

compact descriptor

[1] "Three things everyone should know to improve object retrieval", R. Arandjelovic and A. Zisserman. CVPR, 2012.

# Face Fisher Vector

face image

**set of local features → high-dim Fisher vector**

## Fisher Vector (FV) encoding[1]

- describes a set of local features in a single vector

- diagonal-covariance GMM as a codebook

  - appearance: SIFT

  - location: (x,y)

- GMM can be seen as a face model

**face FV extraction**

discriminative projection

ellipses – means & variances of GMM's (x,y) components

compact descriptor

[1] "Improving the Fisher kernel for large-scale image classification", Perronnin  et al., ECCV 2010

# Face Fisher Vector

face image

**set of local features → high-dim Fisher vector**

- Image FV – normalised sum of feature FVs
- Feature FV – feature space location statistics:

**1st order stats (k-th Gaussian):** $\Phi_k^{(1)} \sim \alpha_k \left( \dfrac{\mathbf{x} - \mu_k}{\sigma_k} \right)$

**2nd order stats (k-th Gaussian):** $\Phi_k^{(2)} \sim \alpha_k \left( \dfrac{(\mathbf{x} - \mu_k)^2}{\sigma_k^2} - 1 \right)$

**face FV extraction**

**soft-assignment to GMM**

$(\mu_k, \sigma_k)$

$\alpha_k$

$\mathbf{X}$

discriminative projection

compact descriptor

# Face Fisher Vector

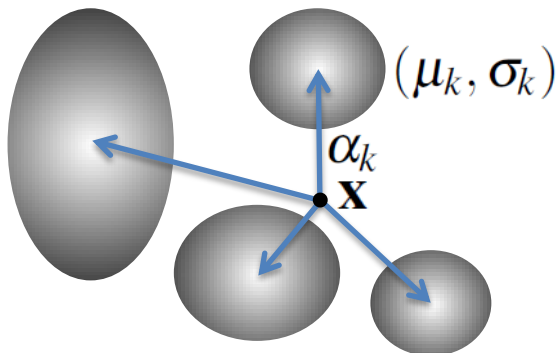**set of local features → high-dim Fisher vector**

- Image FV – normalised sum of feature FVs
- Feature FV – feature space location statistics:

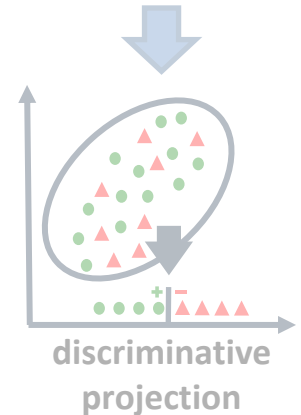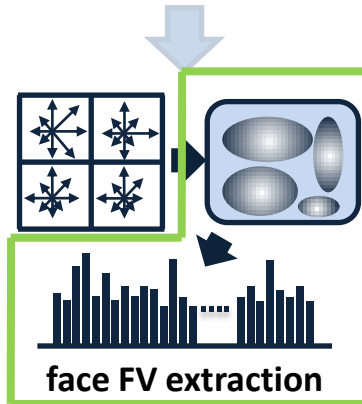**1st order stats (k-th Gaussian):** $\quad \Phi_k^{(1)} \sim \alpha_k \left( \dfrac{\mathbf{x} - \mu_k}{\sigma_k} \right)$

**2nd order stats (k-th Gaussian):** $\quad \Phi_k^{(2)} \sim \alpha_k \left( \dfrac{(\mathbf{x} - \mu_k)^2}{\sigma_k^2} - 1 \right)$

**face FV extraction**

**soft-assignment to GMM**

$(\mu_k, \sigma_k)$

$\alpha_k$

$\mathbf{X}$

**stacking**

$$\phi(\mathbf{x}) = \left[ \Phi_1^{(1)}, \Phi_1^{(2)}, \ldots, \Phi_K^{(1)}, \Phi_K^{(2)} \right]$$

**66-D**  **66-D**   **66-D**

FV dimensionality: 66×2×512=67,584
(for a mixture of 512 Gaussians)

**discriminative projection**

**compact descriptor**

# Distance Learning


face image

**high-dim FV → low-dim face descriptor**

- Large-margin distance constraints:

$$y_{ij}\left(b - d^2(\phi_i, \phi_j)\right) > 1$$

$y_{ij} = 1$ iff (i,j) is the same person, $\phi_i, \phi_j$ – FV



**same**    **different**

$0$    $b-1$    $b+1$    FV distance $d^2(\phi_i, \phi_j)$


face FV extraction


**discriminative projection**

- Distance models:
  - low-rank Mahalonobis       $W = (\equiv)$
  - joint distance-similarity    $W = (\equiv)$  $V = (\equiv)$
  - weighted Euclidean          $U = (\longrightarrow)$


compact descriptor

# Projection Learning

- Low-rank Mahalanobis distance (projection W):

$$d_W^2(\phi_i, \phi_j) = \|W\phi_i - W\phi_j\|_2^2 = (\phi_i - \phi_j)^T W^T W (\phi_i - \phi_j)$$

$$W = \begin{pmatrix} \rule{2cm}{0.4pt} \\ \rule{2cm}{0.4pt} \\ \rule{2cm}{0.4pt} \end{pmatrix}$$

- Large-margin objective: $\arg\min\limits_{W} \sum\limits_{i,j} \max\left[1 - y_{ij}\left(b - d_W^2(\phi_i, \phi_j)\right), 0\right]$

  - regularisation by $W \in \mathbb{R}^{p \times d}$, $p \ll d$
  - stochastic sub-gradient solver
  - initialised by PCA-whitening

  **Fisher vectors**

- $+$
  - Models dependencies between FV elements
  - Explicit dimensionality reduction

- $-$
  - Non-convex

# Joint Distance-Similarity Learning

- Difference of low-rank distance and inner product[1] :

$$d^2_{W,V}(\phi_i, \phi_j) = \|W\phi_i - W\phi_j\|^2_2 - \langle V\phi_i, V\phi_j \rangle =$$
$$(\phi_i - \phi_j)^T W^T W (\phi_i - \phi_j) - \phi_i V^T V \phi_j$$

$$W = \left(\begin{array}{c}\rule{3em}{0.4pt}\\\rule{3em}{0.4pt}\\\rule{3em}{0.4pt}\end{array}\right)$$

$$V = \left(\begin{array}{c}\rule{3em}{0.4pt}\\\rule{3em}{0.4pt}\\\rule{3em}{0.4pt}\end{array}\right)$$

- Large-margin objective: $\arg\min\limits_{W,V} \sum\limits_{i,j} \max\left[1 - y_{ij}\left(b - d^2_{W,V}(\phi_i, \phi_j)\right), 0\right]$

  - stochastic sub-gradient solver (as before)

**Fisher vectors**

**+**
- Models dependencies between FV elements
- More complex decision (distance) function

**-**
- Two low-dim representations (W & V projections)
- Non-convex

[1] "Blessing of dimensionality: high dimensional feature and its efficient compression for face verification", D. Chen, X. Cao, F. Wen, and J. Sun. CVPR, 2013.

# Distance Learning

- Weighted Euclidean distance (diagonal Mahalanobis)

$$d_u^2(\phi_i, \phi_j) = \sum_k u_k \left( \phi_i^{(k)} - \phi_j^{(k)} \right)^2, \quad u_k \geq 0 \,\forall k$$

$$U = (\text{———})$$

- Large-margin (SVM-like) objective:

$$\arg \min_{u_k \geq 0} \sum_{i,j} \max \left[ 1 - y_{ij} \left( b - d_u^2(\phi_i, \phi_j) \right), 0 \right]$$

**Fisher vectors**

**+**
- Convex, fast to train
- Less parameters → less training data needed

**-**
- Doesn't model dependencies between FV elements
- No dimensionality reduction

# Effect of Parameters

| SIFT density | GMM Size | Spatial Aug. | Desc. Dim. | Distance Function | Hor. Flip. | ROC-EER,% |
|---|---|---|---|---|---|---|
| 2 pix | 256 | | 32768 | diag. metric | | 89.0 |
| 2 pix | 256 | ✓ | 33792 | diag. metric | | 89.8 |
| 2 pix | 512 | ✓ | 67584 | diag. metric | | 90.6 |
| 1 pix | 512 | ✓ | 67584 | diag. metric | | 90.9 |
| 1 pix | 512 | ✓ | 128 | low-rank PCA-whitening | | 78.6 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | | 91.4 |
| 1 pix | 512 | ✓ | 256 | low-rank Mah. metric | | 91.0 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | ✓ | 92.0 |
| 1 pix | 512 | ✓ | $2 \times 128$ | low-rank joint metric-sim. | | 92.2 |
| 1 pix | 512 | ✓ | $2 \times 128$ | low-rank joint metric-sim. | ✓ | 93.1 |

**Effect of FV parameters on accuracy @ ROC-EER[1] (LFW-unrestricted)**

[1] "Is that you? Metric learning approaches for face identification", Guillaumin et al., ICCV 2009.

# Effect of Parameters

| SIFT density | GMM Size | Spatial Aug. | Desc. Dim. | Distance Function | Hor. Flip. | ROC-EER,% |
|---|---|---|---|---|---|---|
| 2 pix | 256 | | 32768 | diag. metric | | 89.0 |
| 2 pix | 256 | ✓ | 33792 | diag. metric | | 89.8 |
| 2 pix | 512 | ✓ | 67584 | diag. metric | | 90.6 |
| 1 pix | 512 | ✓ | 67584 | diag. metric | | 90.9 |
| 1 pix | 512 | ✓ | 128 | low-rank PCA-whitening | | 78.6 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | | 91.4 |
| 1 pix | 512 | ✓ | 256 | low-rank Mah. metric | | 91.0 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | ✓ | 92.0 |
| 1 pix | 512 | ✓ | $2 \times 128$ | low-rank joint metric-sim. | | 92.2 |
| 1 pix | 512 | ✓ | $2 \times 128$ | low-rank joint metric-sim. | ✓ | 93.1 |

**Effect of FV parameters on accuracy @ ROC-EER[1] (LFW-unrestricted)**

Performance increases with:

- spatial augmentation, more Gaussians, higher density

# Effect of Parameters

| SIFT density | GMM Size | Spatial Aug. | Desc. Dim. | Distance Function | Hor. Flip. | ROC-EER,% |
|---|---|---|---|---|---|---|
| 2 pix | 256 | | 32768 | diag. metric | | 89.0 |
| 2 pix | 256 | ✓ | 33792 | diag. metric | | 89.8 |
| 2 pix | 512 | ✓ | 67584 | diag. metric | | 90.6 |
| 1 pix | 512 | ✓ | 67584 | diag. metric | | 90.9 |
| 1 pix | 512 | ✓ | 128 | low-rank PCA-whitening | | 78.6 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | | 91.4 |
| 1 pix | 512 | ✓ | 256 | low-rank Mah. metric | | 91.0 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | ✓ | 92.0 |
| 1 pix | 512 | ✓ | $2\times128$ | low-rank joint metric-sim. | | 92.2 |
| 1 pix | 512 | ✓ | $2\times128$ | low-rank joint metric-sim. | ✓ | 93.1 |

**Effect of FV parameters on accuracy @ ROC-EER[1] (LFW-unrestricted)**

Performance increases with:
- spatial augmentation, more Gaussians, higher density
- discriminative projection (also **500-fold** dimensionality reduction)

# Effect of Parameters

| SIFT density | GMM Size | Spatial Aug. | Desc. Dim. | Distance Function | Hor. Flip. | ROC-EER,% |
|---|---|---|---|---|---|---|
| 2 pix | 256 | | 32768 | diag. metric | | 89.0 |
| 2 pix | 256 | ✓ | 33792 | diag. metric | | 89.8 |
| 2 pix | 512 | ✓ | 67584 | diag. metric | | 90.6 |
| 1 pix | 512 | ✓ | 67584 | diag. metric | | 90.9 |
| 1 pix | 512 | ✓ | 128 | low-rank PCA-whitening | | 78.6 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | | 91.4 |
| 1 pix | 512 | ✓ | 256 | low-rank Mah. metric | | 91.0 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | ✓ | 92.0 |
| 1 pix | 512 | ✓ | $2 \times 128$ | low-rank joint metric-sim. | | 92.2 |
| 1 pix | 512 | ✓ | $2 \times 128$ | low-rank joint metric-sim. | ✓ | 93.1 |

**Effect of FV parameters on accuracy @ ROC-EER[1] (LFW-unrestricted)**

Performance increases with:

- spatial augmentation, more Gaussians, higher density
- discriminative projection (also **500-fold** dimensionality reduction)
- averaging across 4 combinations of horizontally flipped faces

# Effect of Parameters

| SIFT density | GMM Size | Spatial Aug. | Desc. Dim. | Distance Function | Hor. Flip. | ROC-EER,% |
|---|---|---|---|---|---|---|
| 2 pix | 256 | | 32768 | diag. metric | | 89.0 |
| 2 pix | 256 | ✓ | 33792 | diag. metric | | 89.8 |
| 2 pix | 512 | ✓ | 67584 | diag. metric | | 90.6 |
| 1 pix | 512 | ✓ | 67584 | diag. metric | | 90.9 |
| 1 pix | 512 | ✓ | 128 | low-rank PCA-whitening | | 78.6 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | | 91.4 |
| 1 pix | 512 | ✓ | 256 | low-rank Mah. metric | | 91.0 |
| 1 pix | 512 | ✓ | 128 | low-rank Mah. metric | ✓ | 92.0 |
| 1 pix | 512 | ✓ | $2\times128$ | low-rank joint metric-sim. | | 92.2 |
| 1 pix | 512 | ✓ | $2\times128$ | low-rank joint metric-sim. | ✓ | 93.1 |

**Effect of FV parameters on accuracy @ ROC-EER[1] (LFW-unrestricted)**

Performance increases with:
- spatial augmentation, more Gaussians, higher density
- discriminative projection (also **500-fold** dimensionality reduction)
- averaging across 4 combinations of horizontally flipped faces
- combined distance-similarity score function

# Effect of Face Alignment

- Robust w.r.t. alignment and crop:
    - LFW → align & crop[1]:                                      92.0%
    - LFW-deep-funneled[2] → 150×150 crop:            92.0%
    - LFW-funneled[3] → 150×150 crop:                    91.7%
    - LFW → Viola-Jones crop (**no alignment**):        90.9%

- Good results without alignment
    - just run Viola-Jones and compute FVF!
    - might not hold for other datasets

- Setting: LFW-unrestricted, projection learning, horiz. flipping

[1] "Taking the bite out of automatic naming of characters in TV video",  Everingham et al., IVC 2009.
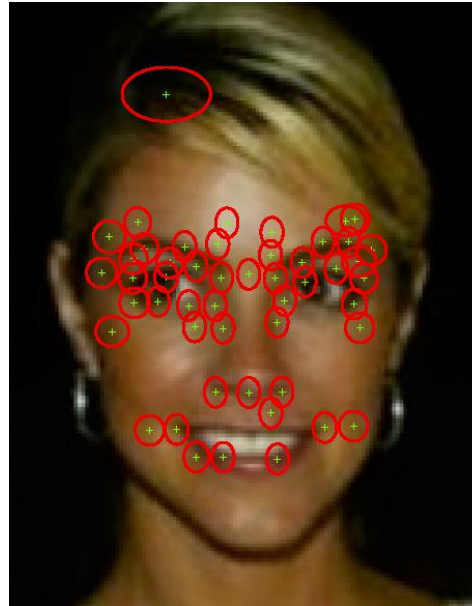[2] "Learning to align from scratch", Huang et al., NIPS 2012
[3] "Unsupervised joint alignment of complex images", Huang et al., ICCV 2007
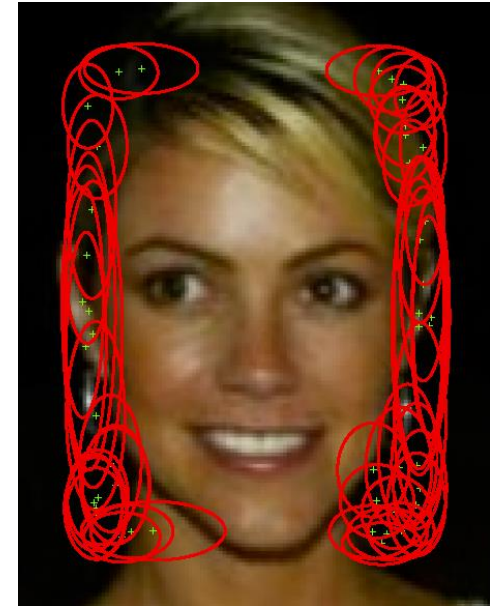
# Learnt Model Visualisation



**all Gaussians**

**important
(top-50 Gaussians)**

**irrelevant
(bottom-50 Gaussians)**

Gaussian ranking (for visualisation):

GMM component → FV sub-vector → W sub-matrix → its energy

**dimensionality
reduction projection**

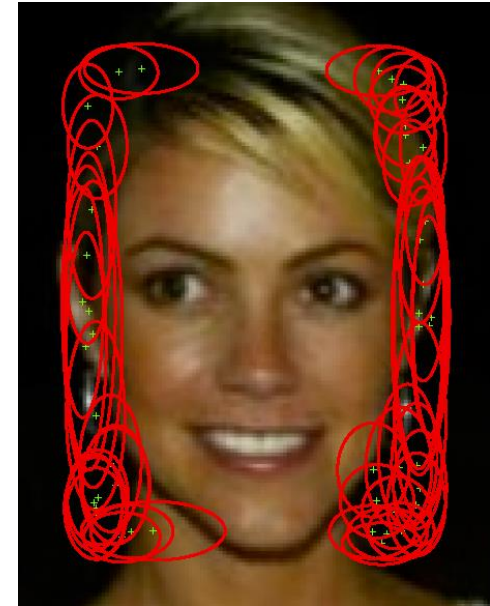$$W = \begin{array}{|c|c|c|c|} \hline \textbf{1}^{\textbf{st}}\textbf{ Gaussian} & \textbf{2}^{\textbf{nd}}\textbf{ Gaussian} & & \textbf{512}^{\textbf{th}}\textbf{ Gaussian} \\ \hline \end{array}$$

# Learnt Model Visualisation



**all Gaussians**

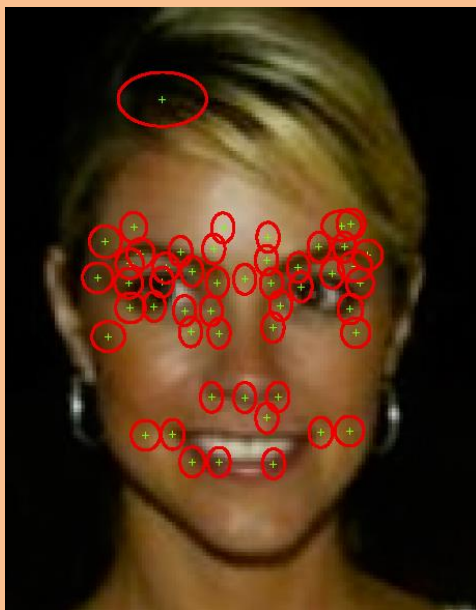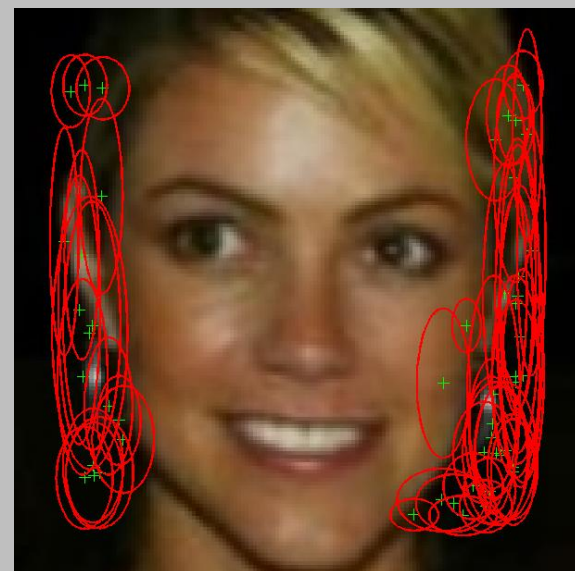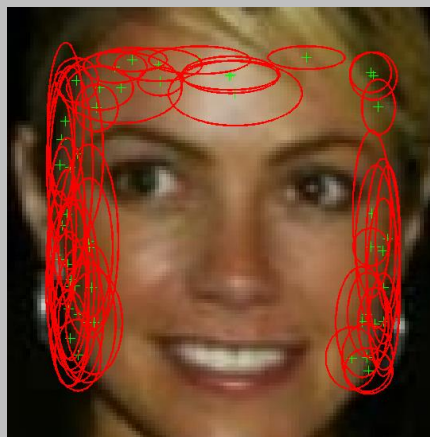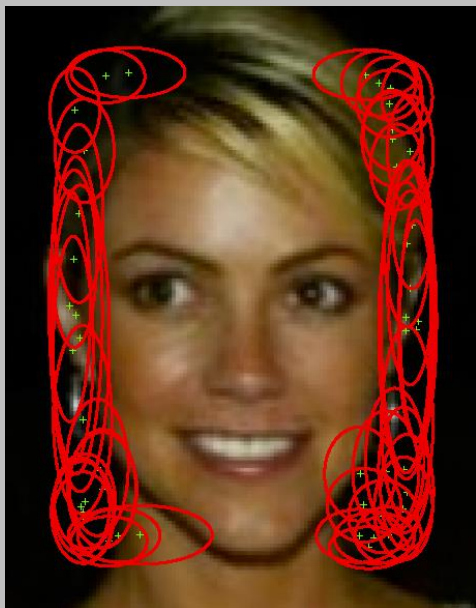**important**
**(top-50 Gaussians)**

**irrelevant**
**(bottom-50 Gaussians)**

- High-ranked Gaussians (centre)
  - **match facial features** (weren't explicitly trained to do so)
  - fine localisation (low spatial variance)
- Low-ranked Gaussians (right)
  - cover background areas
  - loose localisation (high spatial variance)

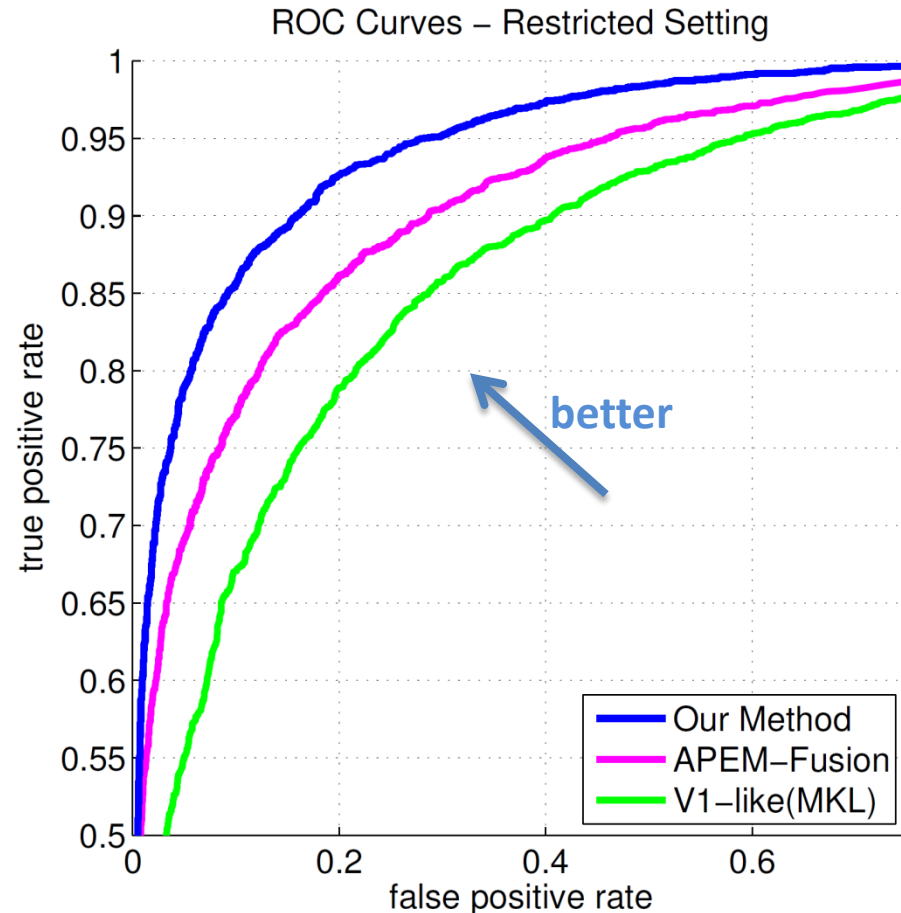LFW → alignment | LFW, no alignment (Viola-Jones box) | LFW-funneled

important (top-50 Gaussians)

irrelevant (bottom-50 Gaussians)

# Results: LFW-restricted

| Method | Mean Acc. |
|---|---|
| V1-like/MKL  [26] | $0.7935 \pm 0.0055$ |
| PEM SIFT [19] | $0.8138 \pm 0.0098$ |
| APEM Fusion [19] | $0.8408 \pm 0.0120$ |
| **Our Method** | $\mathbf{0.8747 \pm 0.0149}$ |

**verification accuracy**

- no outside training data

- LFW-funneled images
  - 150×150 centre crop

- limited training data
  - just 5400 fixed image pairs
  - used diagonal metric (SVM)
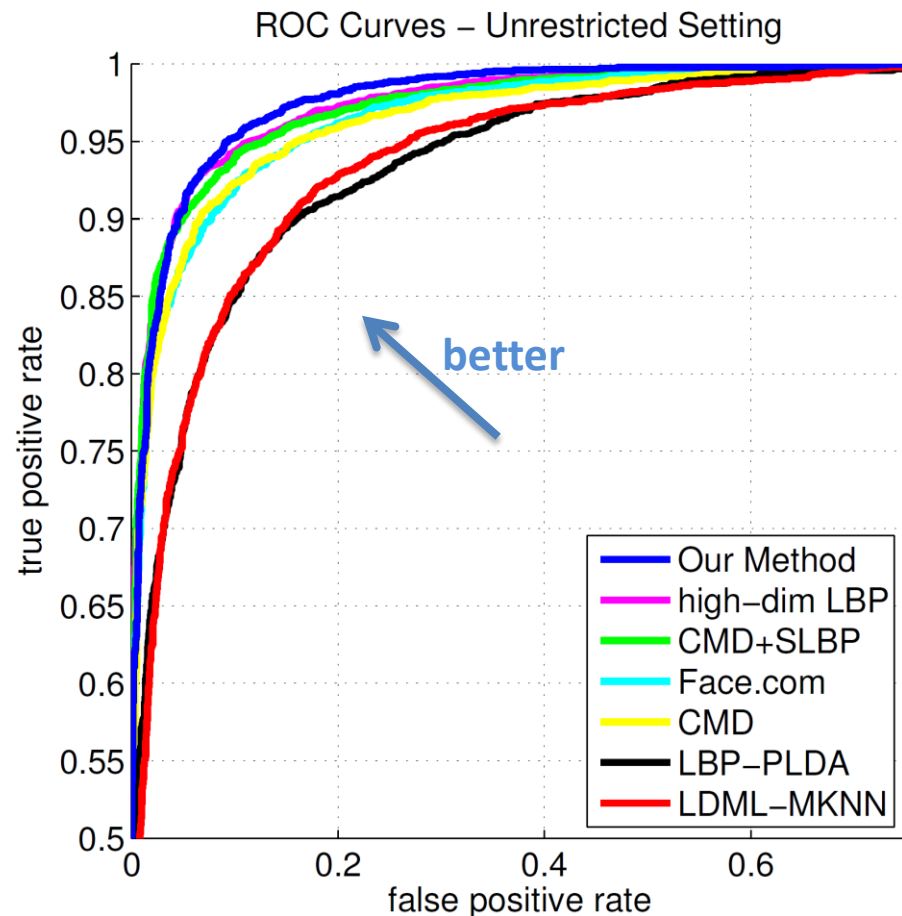
- **state-of-the-art** accuracy: 87.47% vs 84.08%[1]



ROC Curves – Restricted Setting

[1] "Probabilistic elastic matching for pose variant face verification", H. Li, G. Hua, J. Brandt, and J. Yang. CVPR 2013.

# Results: LFW-unrestricted

| Method | Mean Acc. |
|---|---|
| LDML-MkNN [10] | $0.8750 \pm 0.0040$ |
| Combined multishot [32] | $0.8950 \pm 0.0051$ |
| Combined PLDA [20] | $0.9007 \pm 0.0051$ |
| face.com [31] | $0.9130 \pm 0.0030$ |
| CMD + SLBP [12] | $0.9258 \pm 0.0136$ |
| LBP multishot [32] | $0.8517 \pm 0.0061$ |
| LBP PLDA [20] | $0.8733 \pm 0.0055$ |
| SLBP [12] | $0.9000 \pm 0.0133$ |
| CMD [12] | $0.9170 \pm 0.0110$ |
| High-dim SIFT [6] | $0.9177 \pm$ N/A |
| High-dim LBP [6] | $0.9318 \pm 0.0107$ |
| **Our Method** | $\mathbf{0.9303 \pm 0.0105}$ |

**verification accuracy**

ROC Curves – Unrestricted Setting

better

Our Method
high–dim LBP
CMD+SLBP
Face.com
CMD
LBP–PLDA
LDML–MKNN

true positive rate

false positive rate

- outside training data only for alignment [Everingham '09]
- any number of training image pairs
- matches **state-of-the-art** accuracy: 93.03% vs 93.18%[1]

[1] "Blessing of dimensionality: high dimensional feature and its efficient compression for face verification", D. Chen, X. Cao, F. Wen, and J. Sun. CVPR, 2013.

# Summary

- **Fisher Vector Face (FVF)** representation
  - achieves state-of-the-art on LFW (restricted & unrestricted)
  - performs very well on top of different alignment schemes

- FVF is based on off-the-shelf techniques
  - dense SIFT (no need for sophisticated landmark detectors)
  - Fisher vector
  - discriminative dimensionality reduction