

Observer motion problem

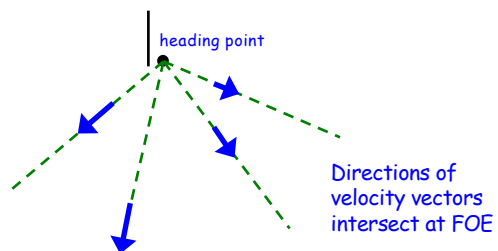


From image motion, compute:

- Observer translation
(T_x T_y T_z)
- Observer rotation
(R_x R_y R_z)
- Depth at every location
 $Z(x,y)$

1-1

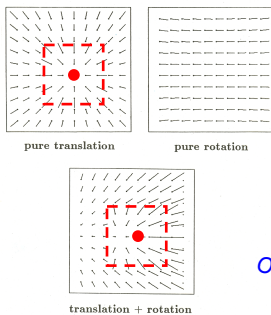
Observer translates toward heading point



But... simple strategy doesn't work if observer also rotates

1-2

Observer motion problem, revisited



From image motion, compute:

- Observer translation
(T_x T_y T_z)
- Observer rotation
(R_x R_y R_z)
- Depth at every location
 $Z(x,y)$

Observer undergoes both translation + rotation

1-3

Equations of observer motion

Translation (T_x, T_y, T_z)	Rotation (R_x, R_y, R_z)	Depth $Z(x,y)$
---	--	--------------------------

$$V_x = (-T_x + xT_z)/Z + R_xxy - R_y(x^2+1) + R_zy$$

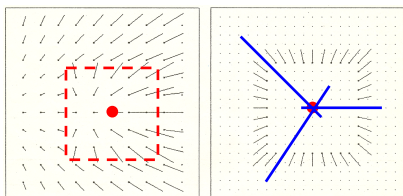
$$V_y = (-T_y + yT_z)/Z + R_x(y^2+1) - R_yxy - R_zx$$

↓
Translational Component

↓
Rotational Component

1-4

Longuet-Higgins & Prazdny



- Along a depth discontinuity, *velocity differences* depend only on observer translation
- Velocity differences point to the observer's heading point (FOE of the translational component of motion)

1-5

What is a chair?



Alignment methods

Find an object model and geometric transformation that *best match* the viewed image

- V viewed object (image)
- M_i object models
- T_{ij} allowable transformations between viewed object and models
- F measure of fit between V and the expected appearance of model M_i under the transformation T_{ij}

GOAL: Find a combination of M_i and T_{ij} that maximizes the fit F

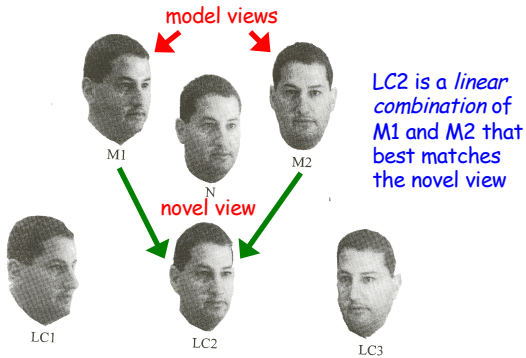
1-7

Alignment method: recognition process

- (1) Find best transformation T_{ij} for each model M_i (optimizing over possible views)
- (2) Find M_i whose best T_{ij} gives the best match to image V

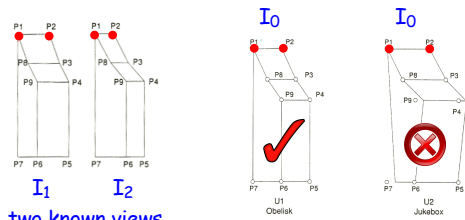
1-8

Recognition by linear combination of views



1-9

Predicting object appearance



$$X_{P_1 I_0} = \alpha X_{P_1 I_1} + \beta X_{P_1 I_2}$$

$$X_{P_2 I_0} = \alpha X_{P_2 I_1} + \beta X_{P_2 I_2}$$

Recognition process:

- (1) compute α, β that predict P1 and P2
- (2) use α, β to predict other points
- (3) evaluate fit of model to image

1-10

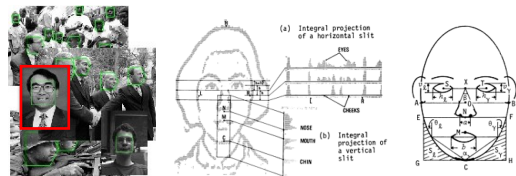
Why is face recognition hard?



It all began with Takeo Kanade (1973)...

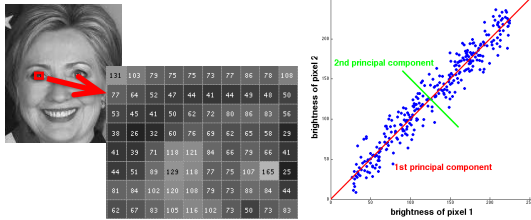
PhD thesis, *Picture Processing System by Computer Complex and Recognition of Human Faces*

- Special purpose algorithms to locate eyes, nose, mouth, boundaries of face
- ~ 40 geometric features, e.g. ratios of distances and angles between features



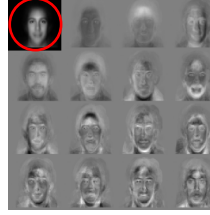
Eigenfaces for recognition (Turk & Pentland) Principal Components Analysis (PCA)

Goal: reduce the dimensionality of the data while retaining as much information as possible in the original dataset
PCA allows us to compute a linear transformation that maps data from a high dimensional space to a lower dimensional subspace



Eigenfaces for recognition (Turk & Pentland)

$\Psi(x,y)$



Perform **PCA** on a large set of training images, to create a set of *eigenfaces*, $E_i(x,y)$, that span the data set

First components capture most of the variation across the data set, later components capture subtle variations

$\Psi(x,y)$: average face (across all faces)

<http://vismod.media.mit.edu/vismod/demos/facerec/basic.html>

Each face image $F(x,y)$ can be expressed as a weighted combination of the eigenfaces $E_i(x,y)$:

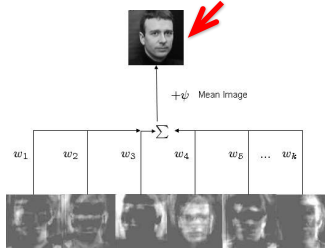
$$F(x,y) = \Psi(x,y) + \sum_i w_i * E_i(x,y)$$

1-14

Representing individual faces

Each face image $F(x,y)$ can be expressed as a weighted combination of the eigenfaces $E_i(x,y)$:

$$F(x,y) = \Psi(x,y) + \sum_i w_i * E_i(x,y)$$



Recognition process:

- (1) Compute weights w_i for novel face image
- (2) Find image m in face database with most similar weights, e.g.

$$\min \sum_{i=1}^k (w_i - w_i^m)^2$$

Face detection: Viola & Jones

Multiple view-based classifiers based on simple features that best discriminate faces vs. non-faces

Most discriminating features **learned** from thousands of samples of face and non-face image windows

Attentional mechanism: cascade of increasingly discriminating classifiers improves performance



1-16

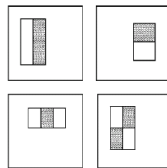
Viola & Jones use simple features

Use simple *rectangle features*:

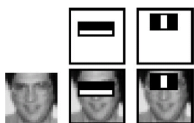
$$\sum I(x,y) \text{ in gray area} - \sum I(x,y) \text{ in white area}$$

within 24 x 24 image sub-windows

- initially consider 160,000 potential features per sub-window!
- features computed very efficiently



Which features best distinguish face vs. non-face?



Learn most discriminating features from thousands of samples of face and non-face image windows

1-17

Learning the best features

weak classifier using one feature:

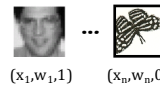
x = image window

f = feature

$p = +1$ or -1

θ = threshold

$$h(x, f, p, \theta) = \begin{cases} 1 & \text{if } pf(x) < p\theta \\ 0 & \text{otherwise} \end{cases}$$



n training samples, equal weights, known classes

$$C(x) = \begin{cases} 1 & \sum_{i=1}^T \alpha_i h_i(x) \geq \tau \\ 0 & \text{otherwise} \end{cases}$$

normalize weights → find next best weak classifier $\epsilon_t = \min_{f,p,\theta} \sum_i w_i |h(x_i, f, p, \theta) - y_i|$ → final classifier

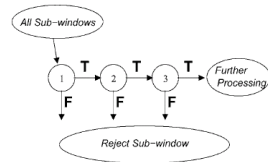
AdaBoost

use classification errors to update weights

~ 200 features yields good results for "monolithic" classifier

1-18

"Attentional cascade" of increasingly discriminating classifiers



Early classifiers use a few highly discriminating features, low threshold

• 1st classifier uses two features, removes 50% non-face windows



• later classifiers distinguish harder examples

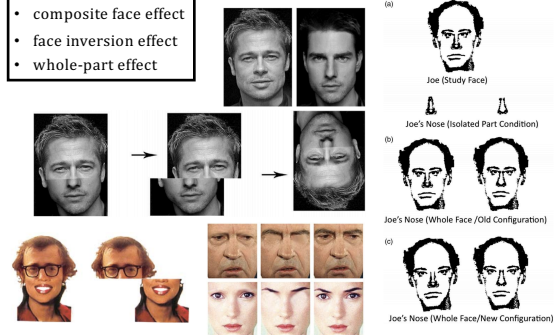
- Increases efficiency
 - Allows use of many more features
- Cascade of 38 classifiers, using ~6000 features

1-19

Feature based vs. holistic processing

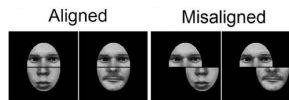
Tanaka & Simonyi (2016), Sinha et al. (2006)

- composite face effect
- face inversion effect
- whole-part effect



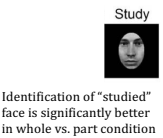
Feature based vs. holistic processing

composite face effect



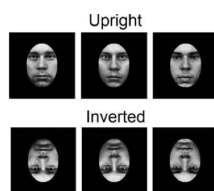
- identical top halves seen as different when aligned with different bottom halves
- when misaligned, top halves perceived as identical

whole-part effect



Identification of "studied" face is significantly better in whole vs. part condition

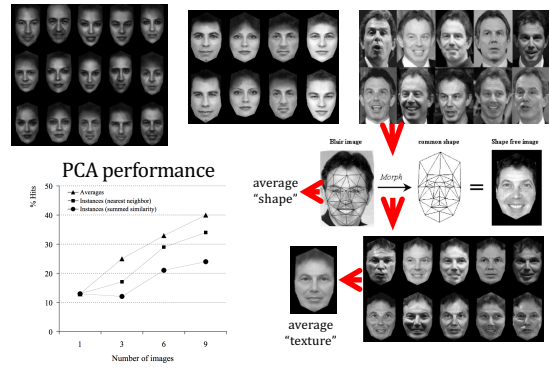
face inversion effect



- inversion disrupts recognition of faces more than other objects
- prosopagnosics do not show effect

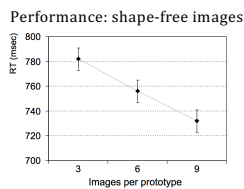
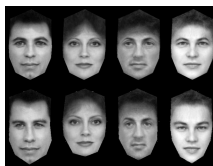


The power of averages, Burton et al. (2005)

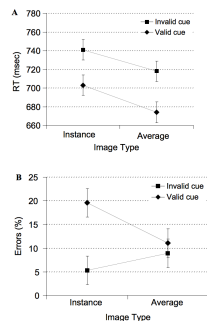


Human recognition of average faces

Burton et al. (2005)

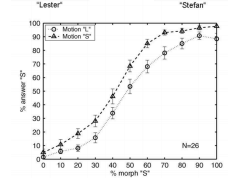
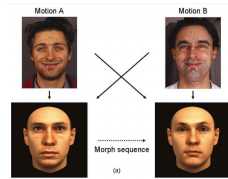
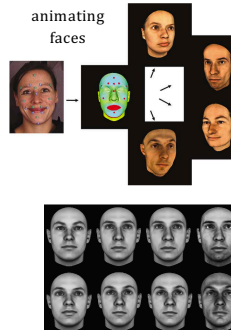


Performance: texture + shape images

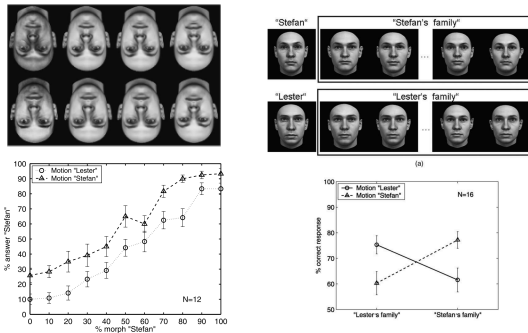


Role of facial motion, Knappmeyer et al. 2003

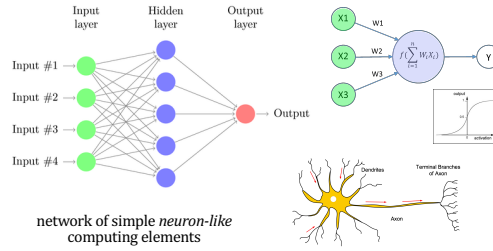
animating faces



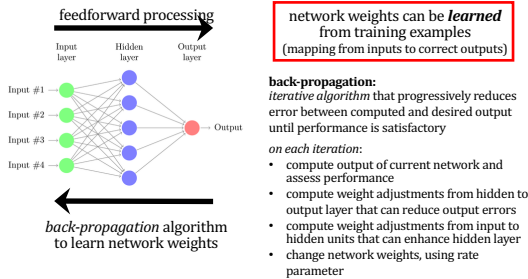
Role of facial motion, Knappmeyer et al. 2003



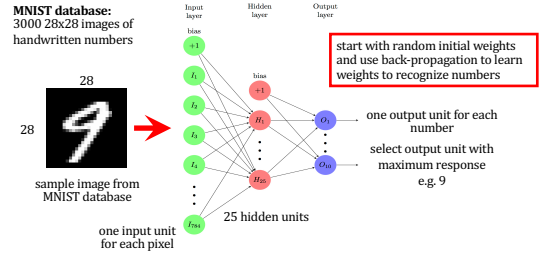
Artificial Neural Networks



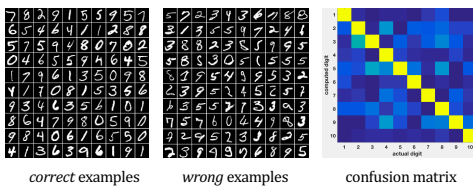
Learning to Recognize Input Patterns



Example: Learning Handwritten Digits



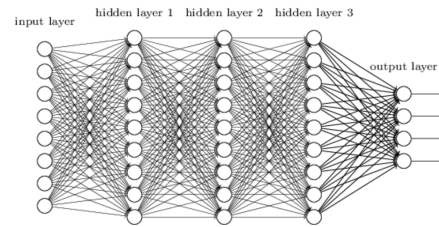
Results: Learning Handwritten Digits



overall classification accuracy: 87.1%

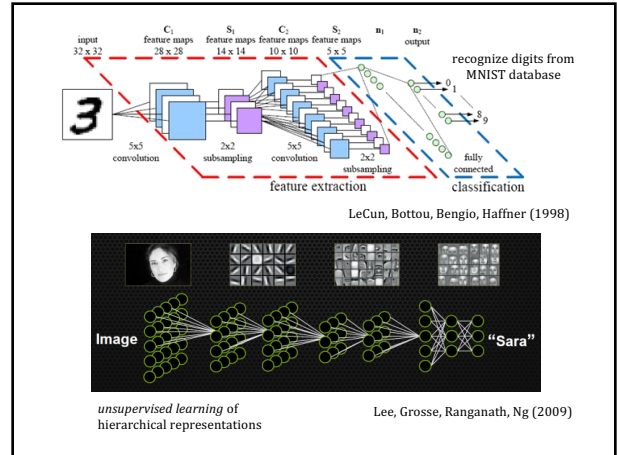
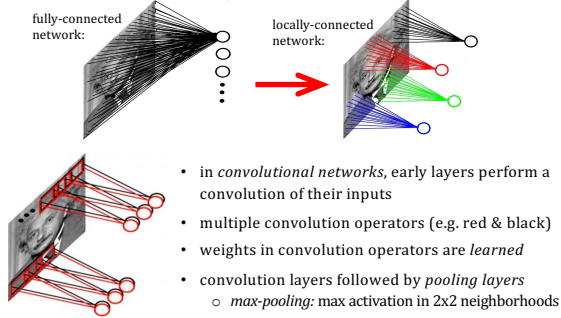
"Deep" neural networks

- early work extended simple neural networks to have multiple, highly-connected hidden layers
- **if** such networks could be trained, they would be much more powerful than "shallow" neural nets
- **but** multi-layer networks are much harder to train!!



Deep convolutional networks

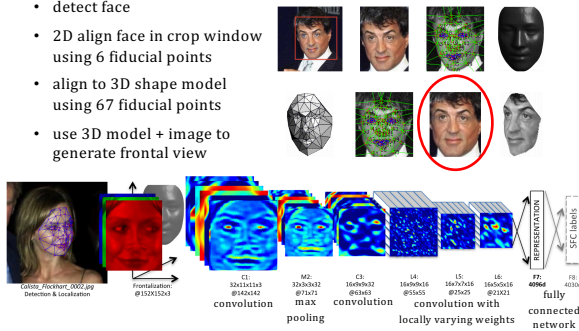
- for image processing, if network inputs are individual pixels, it does not make sense for early layers to be *fully connected*



Facebook's DeepFace system

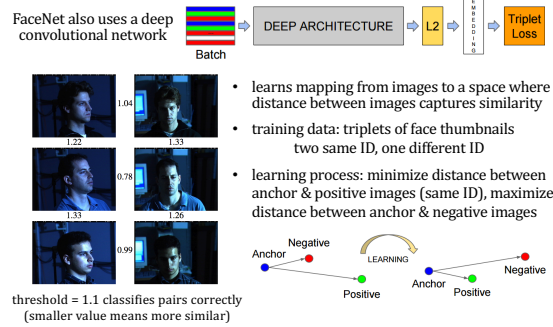
Taigman et al., 2014

- detect face
- 2D align face in crop window using 6 fiducial points
- align to 3D shape model using 67 fiducial points
- use 3D model + image to generate frontal view



Google's FaceNet system

Schroff et al., 2015



Neural Processing in the Ventral Visual Pathway

