

## Observer motion problem

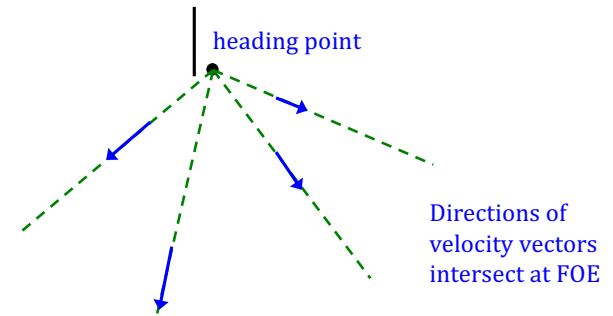


From image motion, compute:

- 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

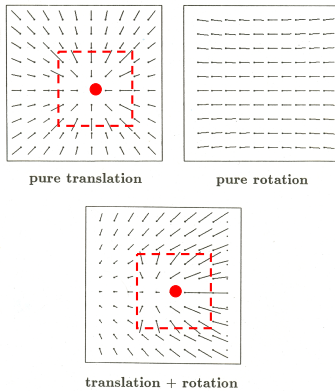
## Observer just translates toward FOE



**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

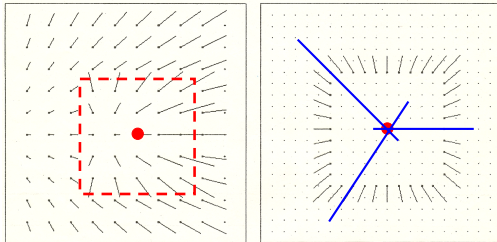
1-3

## Equations of observer motion

Translation $(T_x, T_y, T_z)$	Rotation $(R_x, R_y, R_z)$	Depth $Z(x, y)$
$V_x = (-T_x + xT_z)/Z$	$+ R_xxy - R_y(x^2+1) + R_zy$	
$V_y = (-T_y + yT_z)/Z$	$+ R_x(y^2+1) - R_yxy - R_zx$	
↓ <b>Translational Component</b>	↓ <b>Rotational Component</b>	

1-4

## Longuet-Higgins & Prazdny



- Along a depth discontinuity, *velocity differences* depend only on observer translation
- Velocity differences point to the focus of expansion

1-5

## What is a chair?



## Alignment methods

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

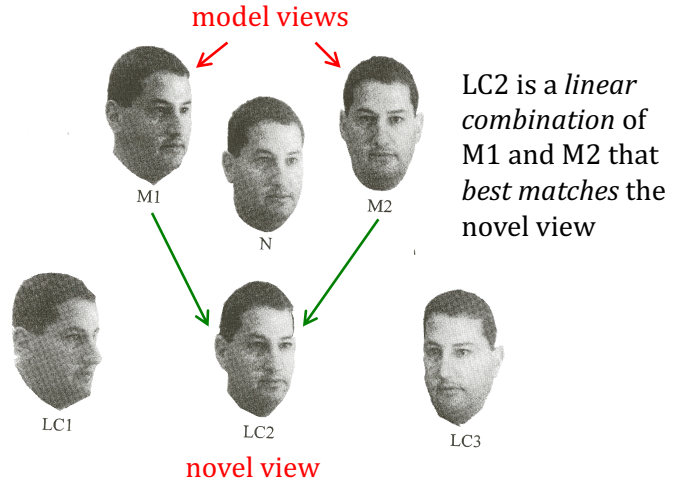
- V viewed object (image)
- $M_i$  object models
- $T_{ij}$  allowable transformations between viewed object and models
- F measure of fit between V and the expected appearance of model  $M_i$  under the transformation  $T_{ij}$

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

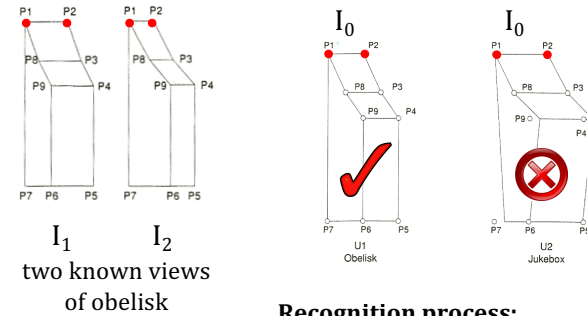
## Alignment method: recognition process

- (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

## Recognition by linear combination of views



## Predicting object appearance



$$X_{P1I0} = \alpha X_{P1I1} + \beta X_{P1I2}$$

$$X_{P2I0} = \alpha X_{P2I1} + \beta X_{P2I2}$$

- (1) compute  $\alpha, \beta$  that predict P1 & P2
- (2) use  $\alpha, \beta$  to predict other points
- (3) evaluate fit of model to image

## Why is face recognition hard?



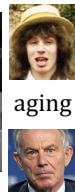
changing pose



changing illumination



changing expression



aging

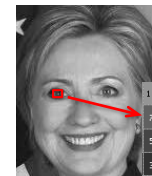


clutter  
occlusion

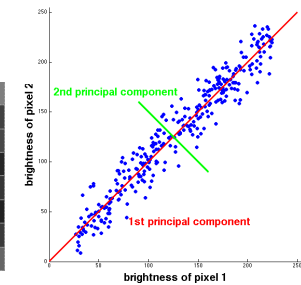
## Eigenfaces for recognition (Turk & Pentland) Principal Components Analysis (PCA)

**Goal:** reduce the dimensionality of the data while retaining as much information as possible in the original dataset

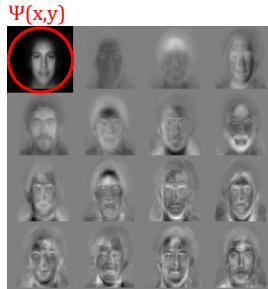
PCA allows us to compute a linear transformation that maps data from a high dimensional space to a lower dimensional subspace



131	103	79	75	75	73	77	86	78	108
77	64	52	47	44	41	44	43	48	50
53	45	41	50	62	72	80	86	93	56
38	26	32	60	76	69	62	65	58	29
41	39	71	118	121	84	66	79	66	41
44	51	89	123	118	77	75	107	165	25
91	84	102	120	109	79	73	88	94	44
62	67	83	105	116	102	73	50	73	83



## Eigenfaces for recognition (Turk & Pentland)



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  $E_i(x,y)$ :

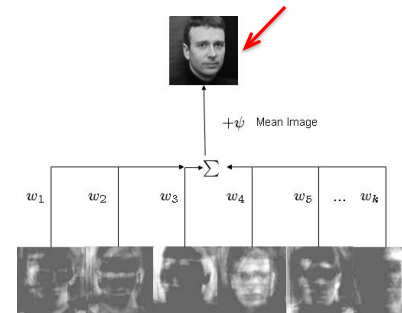
$$F(x,y) = \Psi(x,y) + \sum_i w_i * E_i(x,y)$$

1-13

## 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



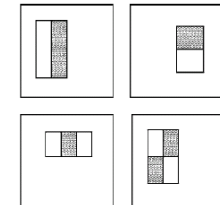
1-15

## Viola & Jones use simple features

Use simple *rectangle features*:

$\sum I(x,y)$  in gray area -  $\sum I(x,y)$  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

1-16

## Learning the best features

weak classifier using one feature:

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

$x$  = image window

$f$  = feature

$p = +1$  or  $-1$

$\theta$  = threshold



$(x_1, w_1, 1)$



$(x_n, w_n, 0)$

$n$  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

$$\epsilon_i = \min_{f,p,\theta} \sum_i w_i |h(x_i, f, p, \theta) - y_i|$$

final classifier

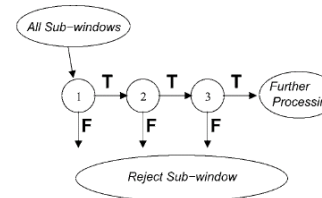
AdaBoost

use classification errors to update weights

~ 200 features yields good results for "monolithic" classifier

1-17

## "Attentional cascade" of increasingly discriminating classifiers



Early classifiers use a few highly discriminating features, low threshold

- 1<sup>st</sup> classifier uses two features, removes 50% non-face windows



- later classifiers distinguish harder examples

- Increases efficiency

- Allows use of many more features

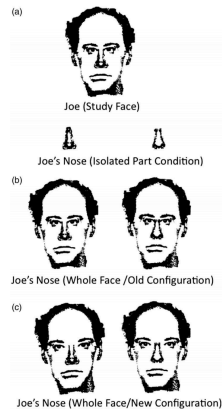
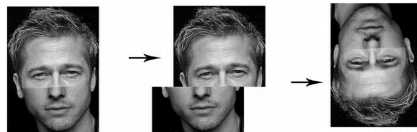
→ Cascade of 38 classifiers, using ~6000 features

1-18

## Feature based vs. holistic processing

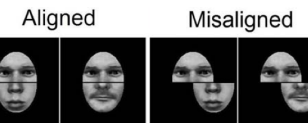
Tanaka & Simonyi (2016), Sinha et al. (2006)

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



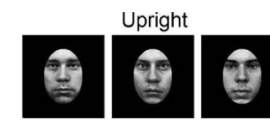
## Feature based vs. holistic processing

composite face effect

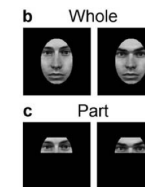


- identical top halves seen as different when aligned with different bottom halves
- when misaligned, top halves perceived as identical

face inversion effect



whole-part effect

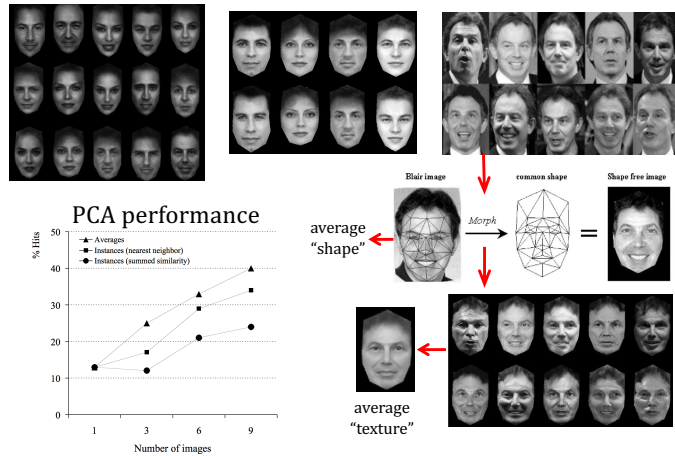


Identification of "studied" face is significantly better in whole vs. part condition

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

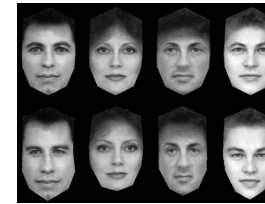


## The power of averages, Burton et al. (2005)

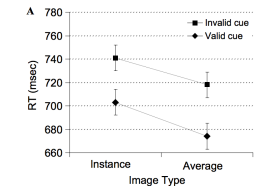


## Human recognition of average faces

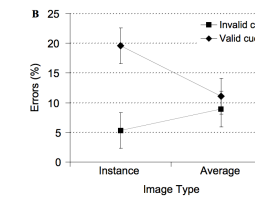
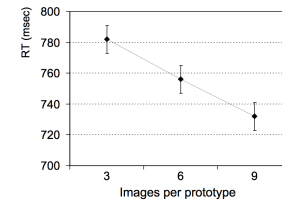
Burton et al. (2005)



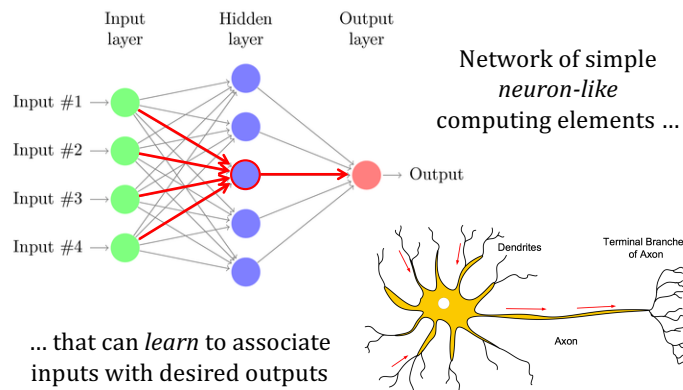
Performance: texture + shape images



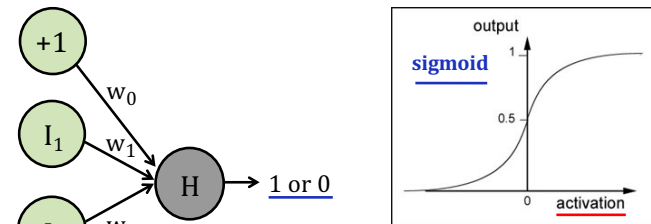
Performance: shape-free images



## What is an artificial neural network?



## Computing in a "typical" neural network

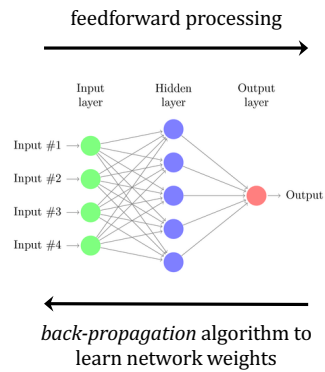


$$w_0 \cdot 1 + w_1 \cdot I_1 + w_2 \cdot I_2 + \dots + w_n \cdot I_n > 0$$

**activation**

How does each unit integrate its inputs to produce an output?  
sum of weighted inputs → sigmoid function → output between 0 and 1

## Learning to recognize input patterns



network weights can be **learned**  
from training examples  
(map inputs to correct outputs)

### back-propagation:

*iterative algorithm* progressively reduces error between computed and desired output until performance is satisfactory

*on each iteration:*

- compute output of current network and assess performance
- compute weight adjustments from hidden to output layer that reduce output errors
- compute weight adjustments from input to hidden units that improve hidden layer
- change network weights, incorporating a rate parameter