# COMPUTER VISION

**AI - complete problem**
↳ signal to symbol

→ requires: (1) Find robust representations of world
(2) Maintaining + updating them (with Machine Learning
(3) Interfacing with attention, goals and plans.

Challenging as:

↳① Inverse optics ∴ 2D → 3D world. Correspondence problem.

↳② Inverse graphics ⇒ need to deal with 3D world but only 2D image with surfaces occluded, shading, etc.

↳③ Cognitive Penetrance ∨ Paradox : challenging to solve problems that are simple for humans. Since cannot reverse engineer - process of seeing things includes a highly complex model that is tough to introspect ⇒ eg Facial Recognition

↳④ Few tasks can be done bottom-up (data driven). Need top-down (prior knowledge) + model-driven.

↳⑤ Signal data is often terrible

↳⑥ Goals means of problem is often not **well posed**

↳ Solution (1) exists, (2) unique (3) depends continuously on the data.

↳⑦ Pose-invariance is often a large problem

↳⑧ we have to be able to deal with objects that haven't been seen before. (wide variety)

spatial resolution determined by CCD density and lens properties

→ pixel size limited by photon flux into small areas → per pixel

## Pixel Arrays, CCD/CMOS sensors, image coding

filter layer
sensor array

↳ CMOS
CCD cameras contain independent sensors, converting incident photons (focused by lenses) into charge proportional to light energy.
↳ charge is coupled to allow voltage to be read out easier

↳ Luminance Resolution is the number of distinguishable gray levels = number of bits per pixel
↳ colour arises from three CCD arrays (three types of sensor)

↳ Video = Luma and Chroma channels. ↳ Composite video uses high-frequency chrominance burst

↳ Framegrabber (strobed sampling block) contains high-speed ADC to discrete video into frames.
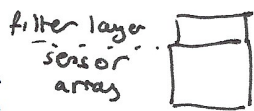↳ video frames stored in 3 byte arrays (each different colour planes)
↳ generally 8 bits / col / pixel = 24 ⊕ 16 / pixel

**IMAGE FORMATS** see slide 25:
↳ need to revert format to use image

NYQUIST SAMPLING THEOREM : highest spatial frequency component of information in an image = ½ sampling density of pixel array

② Pixel array with 640,000 columns can represent spatial frequency components of image structure no higher than 320 cycles/image. If image frames sampled at 30 fps, max temporal frequency component of information within moving sequence is 15 Hertz.
↳ can use NIR to do iris mapping ⇒ with pixel variance and mean, imaging the ratio.

# BIOLOGICAL VISUAL MECHANISMS

Neurones: richly interconnected cells (analogue + digital) with non-linear, adaptive features. Consist of enclosing membrane ⇒ voltage difference between inside and outside.
↳ lipid bilayer (100 µF) - pores that are differentially selective to ions. Cross the neural membrane through protein pores

(3)

Photo-chemical isomerisation
11-cis-retinal +
hv →all-
trans-
retinal
↓
Carbon double
bond flips
from cis to
trans +
causes pore
to close to
Na⁺ ions.

As Na⁺ ions actively
pumped (dark current),
increased resistance,
causes increased
trans-membrane
voltage.

P(photon capture)
$\lambda = c/v$

As more positive ions flow into the neurone, voltage becomes positive on the inside reducing membrane's resistance allowing more to enter. This breakdown in resistance constitutes the nerve impulse. (after refractory period (2ms) it is active again). Impulses propagate down axons at $100 ms^{-1}$. Summations of current flows into neurone from other neurones at synapses ⇒ triggering impulse.

Neural activity is asynchronous ~300 Hz. 2/3 rds of brain receives visual input.
↳ 30 different visual areas (specifically Primary Visual Cortex + Occipital Lobe)

axons of — image processing at + temporal processing
Ganglion cells → first synapse at second synapse
↳ specialised red, blue, green

Retina: 1mm thick, 120 mn light-sensitive photoreceptors ⇒ 6mn (cones) + rest are rods.
↳ part of the brain!    400 - 700 nm    good for night
120 mn photoreceptors                    vision
1mn output channels    mostly near
↳ can reliably see individual photons    fovea - central 20°
(1 to 10")

Blind spot - optic nerve

↳ digital neurones are analogue devices. Photoreceptors respond to absorption by hyperpolarisation

Rods and cones distributed in hexagonal lattices with varying relative densities
↳ imperfect, → incoherent, not crystalline

Retina network: multi-layered - 3 nuclear layer + 2 plexiform layers ⌐ synaptic interconnections
↳ photoreceptors at rear    2 directions of signal flow ⌐ bipolar cells
↳ ① longitudinal (photoreception)    ↓
    gang lion cells
↳ ② horizontal + amacrine cells, outer/inner plexiform

↳ Therefore, both convergent and divergent signals.

↳ Centre-surround comparisons implemented by bipolar cells

↳ Temporal differentiation by amacrine cells, for motion detection

↳ Separate channels for sustained w transient image information by different classes of ganglion cells

↳ Right and left visual fields project to different brain hemispheres

  ↳ at optic chiasm ⇒ crosses over to project only to contralateral brain hemisphere

  ↳ Projects only to the same brain hemisphere
    ↳ share information with Corpus Callosum

  ↳ Projections then go to Lateral Geniculate Nucleus (LGN) - in thalamus
    ↳ model builds here ⇒ neurones receive input primarily from one eye with left and right eye alternating.
      ↳ Ocular Dominance Columns have cycle of 1mm.

Orientation Selectivity : new tuning variable

  ↳ neurones in orientation columns respond to image structures in a preferred range of orientations ⇒ arises from alignment of isotropic subunits in LGN.

  ↳ Constructed into hypercolumns

## Spatial Image Encoding

  ↳ 5 main DOF in spatial structure: ⇒ position $(x, y)$, orientation, receptive field size, phase.   ↳ inferred from boundaries between excitatory and inhibitory regions — bipartite and tripartite

  ↳ Receptive field profiles well described by 2D Gabor Wavelets.


## MATHEMATICAL IMAGE OPERATIONS

Image processing is built with 2D convolutions of an image with small kernel arrays.
  ↳ eg. Edge Detection, Filtering, Feature Extraction.
  CONVOLUTION ⟺ FILTERING ⟺ FOURIER OPERATION
    ↳ convolutions in the Fourier domain are much faster
      ↳ multiplication (given FFT)

Image is superposition of many 2D Fourier components : $f(x,y) = \exp(i\pi(\mu x + \nu y))$
  ↳ 2D spatial frequency $\sqrt{\mu^2 + \nu^2}$, orientation $= \arctan(\nu/\mu)$
Adding conjugate pair is real valued wave

## Convolution Theorem

$f(x,y)$ has FFT $F(\mu, \nu)$, $g(x,y)$ has FFT $G(\mu, \nu)$

$f(x,y) * g(x,y) = h(x,y) = \int_\alpha \int_\beta f(\alpha, \beta) g(x - \alpha, y - \beta) \, d\beta \, d\alpha$

$H(\mu, \nu) = F(\mu, \nu) \, G(\mu, \nu)$

④

$h(x,y)$ is normally subjected to non-linear operations of various kinds of analysis, segmentation, pattern recognition and object classification.

↳ For explicit convolution algorithm, see (60)

Actually, $O(2(\log_2^{*}(n) + 1))$

**Differentiation Theorem:** Computing derivatives of $f(x,y)$, ~~$F(\mu,v)$~~ is equivalent to multiplying its 2DFT, $F(\mu,v)$ by the corresponding spatial frequency coordinate ($\times i$) raised to the power equal to order of differentiation

$$\left(\frac{\partial}{\partial x}\right)^m \left(\frac{\partial}{\partial y}\right)^n f(x,y) \xrightarrow{\text{2DFT}} (i\mu)^m (iv)^n F(\mu,v)$$

Notably, for Laplacian:

$$\nabla^2 f(x,y) \equiv \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}\right) f(x,y) \xrightarrow{\text{2DFT}} -(\mu^2 + v^2) F(\mu,v)$$

## EDGE DETECTION

Why? - demarkate boundaries + occlusions. Helps solve stereo correspondence problem
↳ DISCONTINUITIES = INFORMATION

$$\vec{\nabla} f(x,y) = \left(\frac{\partial f(x,y)}{\partial x}, \frac{\partial f(x,y)}{\partial y}\right)$$

— can be discretized by finite differences
↳ convolution with FINITE DIFFERENCE KERNEL [-1, 1]
↳ can concatenate to get higher derivatives + 2D kernel
↳ Effectively high pass) filters

Grad direction $\theta = \arctan\left(\frac{\partial f}{\partial y} \Big/ \frac{\partial f}{\partial x}\right)$

$$\|\vec{\nabla} f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$

Isotropic Operator - Laplacian: has no preferred orientation → approximation is:
$\begin{bmatrix} -1 & -2 & -1 \\ -2 & 12 & -2 \\ -1 & -2 & -1 \end{bmatrix}$

apply threshold on this

$-\begin{bmatrix} -1 & 2 & -1 \\ -1 & 2 & -1 \\ -1 & 2 & -1 \end{bmatrix}$

But this only works in a specific orientation Integrating vertically, second derivative horizontally

Laplacian of Gaussian works well:

$$\nabla^2 G_\sigma(x,y) = \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}\right) G_\sigma(x,y) = \frac{x^2 + y^2 - 2\sigma^2}{2\pi\sigma^6} \exp\left(-(x^2 + y^2)/2\sigma^2\right)$$

↳ smoothing parameter

↳ to extract good information, we define a scale analysis:
⌐→ multi-scale family of filters
 ↳ Laplacian pyramid extracts image structure in octave bands of spatial frequency

In 2D fourier domain, $\nabla^2$ multiplies by paraboloid $(\mu^2 + v^2)$
↳ Blurring Laplacian by a Gaussian limits the high-frequency component: ⎯ ⎯

Scale parameter $\sigma$ determines where the high-frequency cut-off occurs. Zero-crossings correspond to edge locations. Bandwidth of $\nabla^2 G_\sigma(x,y)$ filter is 1.3 octaves.

$\hookrightarrow$ Logan's Theorem shows this doesn't satisfy one-octave constraint.

$\hookrightarrow$ Doesn't matter what order Laplacian and Gaussian are applied

CANNY EDGE OPERATOR Removes spurious edges that are detected.

① Smooth image with Gaussian filter to reduce noise

② Compute $\vec{\nabla} I(x,y)$ over image

③ Non-max suppression to remove spurious edges

④ Double threshold to local gradient magnitude : strong, weak, suppressed

⑤ Impose connectivity constraint : edges that are weak + not connected to strong are eliminated

## WAVELETS & ACTIVE CONTOURS

Gabor Wavelets :- proposed as model for receptive field profiles of neurones in visual cortex. Wavelets optimal for extracting orientation, position and modulation of image structure. Achieves theoretical lower bound over variables.

position in image

Mother wavelet / family codebook :
$$f(x,y) = \exp\left(-\left[(x-x_0)^2/\alpha^2 + (y-y_0)^2/\beta^2\right]\right)\exp\left(-\left[v_0(x-x_0)\right.\right.$$
$$\left.\left. + v_0(y-y_0)\right]\right)$$

effective width and length

modulation of spatial frequency
$$w_0 = \sqrt{(u_0^2 + v_0^2)}$$
$$\theta_0 = \arctan(v_0/u_0)$$

Given a generic $\psi(x,y)$ 2D Gabor wavelet, can generate daughter wavelets :
$$F(u,v) = \exp\left(-\left[(u-u_0)^2\alpha^2 + (v-v_0)^2\beta^2\right]\right)$$
$$\times \exp\left(i\left[x_0(u-u_0) + y_0(v-v_0)\right]\right)$$

$$\Psi_{mpq\theta}(x,y) = 2^{-2m}\psi(x',y')$$

$\longrightarrow$ to incorporate dilations in size by $2^{-2m}$, translations $(p,q)$ and rotation $\theta$

$$x' = 2^{-m}\left[x\cos\theta + y\sin\theta\right] - p$$
$$y' = 2^{-m}\left[-x\sin\theta + y\cos\theta\right] - q$$

## Quadrature Wavelets : used for automatic localisation of facial features

$\hookrightarrow$ most features can easily be captured using a handful of wavelets.
$\hookrightarrow$ taking the modulus of facial image after convolving with complex 2D wavelet, find features easily.

It is possible to find only circular and parabolic boundary shapes by computing derivatives of contour integrals

$$g(x,y) = \int_\alpha \int_\beta e^{-((x-\alpha)^2 + (y-\beta)^2)/\sigma^2} \cos(w(x-\alpha))I(\alpha,\beta)\,d\beta\,d\alpha$$

$$h(x,y) = \int_\alpha \int_\beta e^{-((x-\alpha)^2 + (y-\beta)^2)/\sigma^2} \sin(w(x-\alpha))I(\alpha,\beta)\,d\beta\,d\alpha$$

$$A^2(x,y) = g^2(x,y) + h^2(x,y)$$

**Hough Transform**: find curves whose parameters wrt to increasing radious of contour integrals

  ↳ Hough Transform is a voting scheme to find instances of shapes within certain class of objects.

    ↳ accumulator space groups edge evidence — parameters of curve.     ↳ gradient magnitudes.

                 ↳ output of Canny operator.

    ↳ For each edge pixel, increment all the compatible accumulator cells. Accumulator cell for which greatest edge evidence found.

## Active Contours: deformable shape models (snakes) — by energy minimisation
(spline)

                     ↳ pull it towards object contour

Can also split or merge contours as well.

  ↳ Changes shape under competing forces : (1) Image Forces

                 (2) Internal Forces — resist excessive deformations

  ↳ External Energy — reflects how poorly snake is fitting a contour

  ↳ Internal Energy — reflects how much snake is bent or stretched

  ↳ Sum of energies minimised by: (1) Gradient descent, (2) Simulated annealing (3) PDEs

    ↳ BUT : numerical instability + stuck in local minima

## Scale - Invariant Feature Transform (SIFT)

               — does not deal with non rigid deformations

(i) Object recognition with geometric invariance

       — try to estimate a homography by identifying keypoints that correspond in different images & find transformation

    ↳ photometric invariance

(2) Matching corresponding parts of different images or objects

              ↳ Feature detectors with orientation index and scale index

(3) 3D Scene Understanding + Action Recognition

            find orientation by edge detectors ↳Eg. extrema, Gaussian image pyramid and resampling

    bins of orientation histogram normalized relative to dominant grad direction.

              ↳ for each region, orientation histogram constructed from gradient directions

↳ MATCHING PROCESS ← matches sought across wide range of scales + positions, 30° orientation bin sizes.

    ↳ compare relative configuration of groups of minutiae    ↳Best bin first priority queue

    ↳ Best candidate match determined as nearest neighbours in extracted keypoints

    ↳ use Hough Transform voting.

## PARALLEL FUNCTIONAL STREAMS

Multiple parallel functional streams in brain for specific visual subdomain : (1) form, (2) colour, etc.

dorsal stream → spatial information of visual information  dorsal & ventral hierarchies

Primary Visual input    Also, conscious and unconscious vision

      But lots of reciprocal pairwise connections between separate areas.

ventral stream → Higher level processing of object form.

# Structure From Texture

supports ← Texture is helpful to identify shapes, both 2D and 3D
figure/ground
segmentation
by dipole statistics.

↳ defined by statistical correlations across the image

Variations in the texture reveal 3D shape, slant, distance etc.

↳ Quasi-periodicity detected best by Fourier-related methods - can estimate
especially using Gabor wavelets. <u>Energy within the periodicities</u>
with modulus of Gabor wavelets coefficient
reveal texture energy variation

⌐ Phase
   Analysis
for person identification
is particularly powerful

NB. Resolving textural spectra
with location information limited
by Heisenberg's Uncertainty Principle +
optimized by Gabor Wavelets.

## Colour Information

$$R(\lambda) = I(\lambda) \odot O(\lambda)$$

Wavelength ⟷ wavelength composition
mixture received    of the illuminant
by camera at
corresponding point

Colour assignments are a matter of
calibration

↳ spectral reflectance
of the object.

## RETINEX

$\boxed{M}$

① Find max $(r_{max}, g_{max}, b_{max})$
across all pixels
② Assume scene contains
objects that reflect all red, blue, green etc.
③ $M = I(\lambda)$
④ Hence $(r, g, b) \Rightarrow (r/r_{max}, g/g_{max}, b/b_{max})$
↳ discounted the illuminant.
Can also be done in local areas rather than
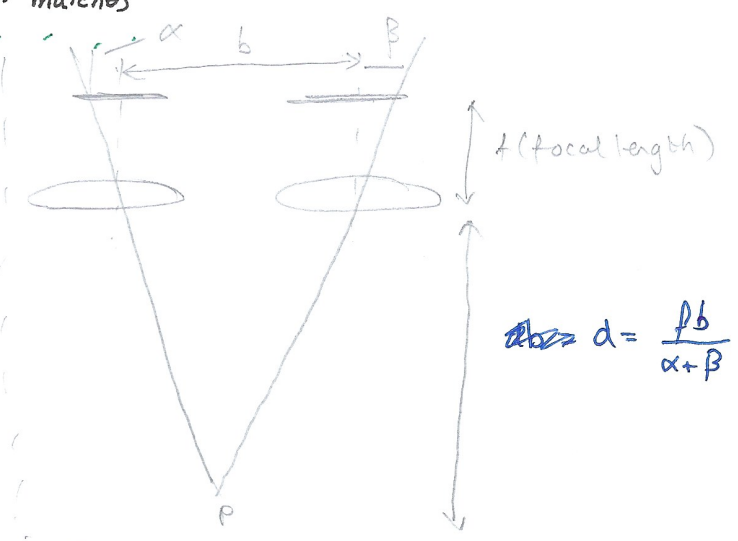just global

## Stereo Vision

2 eyes with base of separation having stereoscopic disparity, depending
on 3D geometry and camera properties. But requires solving correspondence problem.

Parallax: if objects project onto
different parts of the images.
↳ disparity ∝ distance of object
in front or behind point of fixation

Base of triangulation: ~~increased~~
distance between the two cameras

permutation-matching
space is greatly attenuated
terminating with single-
pixel precision matches

⌐ multi-scale image pyramid ⟷
steer search by course-to-fine
strategy to maximise efficiency

## Optical Flow

Apparent motion in a scene due
to relative motion between observer and the
scene (and camera)
↳ ego-motion
Motion estimation requires the solving of the
correspondence problem.

Create velocity vector field for image
↳ may be necessary to assign more than one velocity vector to any local image region
Need to detect a coherent overall motion pattern across options    ↳ motion transparency
↳ need to disambiguate object motion from contour motion

$d = \dfrac{fb}{\alpha + \beta}$

## Intensity Gradient Models

$$-\frac{dI(x,y,t)}{dt} = \vec{v} \cdot \vec{\nabla} I(x,y,t)$$

## Fourier Methods

generally looking for correlated signals across time.

$$F(\omega_x, \omega_y, \omega_t) = \int_x \int_y \int_t I(x,y,t)$$

$$\exp(-i(\omega_x x + \omega_y y + \omega_t t))$$

$$dt\, dy\, dx$$

optical flow also used for localisation through SLAM and LIDAR

## Dynamic Zero-Crossing Models: measure image velocity finding edges and contours.

$$-\frac{\partial}{\partial t}\left[\nabla^2 G_\sigma(x,y) * I(x,y,t)\right]$$

Time-derivative of Laplacian of Gaussian-convolved image

In vicinity of Laplacian zero-crossing. Amplitude is estimate of speed, sign of quantity determines direction of motion relative to normal to contour

Local spatio-temporal spectrum collapses onto 2D inclined plane.

Find motion by applying filters to image sequence observing centre frequencies are co-planar, in this 3-space. Azimuth and elevation correspond to direction and speed of motion.

① Have $I(x,y,t)$ and $F(\omega_x, \omega_y, \omega_t)$. Detecting $\vec{v} = (v_x, v_y)$

② $I(x,y,t) = I(x - v_x t_0, y - v_y t_0, t - t_0)$

③ $F(\omega_x, \omega_y, \omega_t) = \exp(-i(\omega_x v_x t_0 + \omega_y v_y t_0 + \omega_t t_0)) F(\omega_x, \omega_y, \omega_t)$

④ ③ only true if $F() = 0$ where exp factor $\neq 1$     $speed = \sqrt{(v_x^2 + v_y^2)}$

⑤ ∴ $F(\cdots) \neq 0$ only on 3D plane $\omega_x v_x + \omega_y v_y + \omega_t = 0$

$azimuth = direction = \arctan\left(\frac{v_y}{v_x}\right)$

---

## SURFACE AND REFLECTANCE MAPS

Albedo: fraction of illuminant re-emitted from a surface in all directions
↳ Light reflectance is dependent on albedo and geometric factor based on angle

① LAMBERTIAN SURFACES: $\phi(i,e,g) = \cos(i)$ Amount of light (diffuse/matte) reflected dependent on angle of incidence (Lambert's Law) not on angle of emission

② SPECULAR SURFACES: $\phi(i,e,g) = 1$ when $i = e$ ($g =$, $+e$) else $= 0$ Snell's law, perfect reflections

Most surfaces on continuum between Lambertian and specular.

↳ reflection depending on ratio of cosines of angle of incidence and angle of emission ⇒ $\phi(i,e,g) = \frac{\cos(i)}{\cos(e)}$

③ LUNAR SURFACES
↳ Why looks spherical.

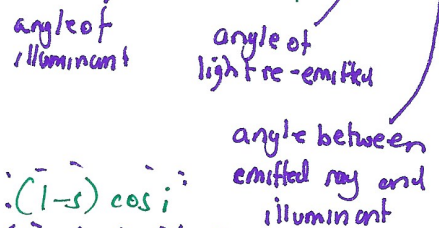Faces: $\phi(i,e,g) = \frac{1}{2}(s(nH)(2\cos i \cos e - \cos g))^n + (1-s)\cos i$

specular: $n$ reflects sharpness       Lambertian

Reflectance Maps: $\phi(i,e,g)$ fraction of incident light reflected per unit solid angle in direction of camera. Otherwise flux/steradian

angle of illuminant    angle of light re-emitted

angle between emitted ray and illuminant

## Shape-from-shading: Requires disambiguation of

↳ (1) Illumination geometry
↳ (2) Reflectance properties of surface (and variations)
↳ (3) Geometry of surface
↳ (4) Rotations of surface
↳ (5) Variations in surface albedo

] all not well known so that the problem is well-posed.

# SHAPE REPRESENTATION & CODON SHAPE GRAMMARS

Curvature map: $\Theta(s) = \lim\limits_{\Delta s \to 0} \dfrac{1}{r(s)}$ where local radius of curvature defined as 'limiting' radius of circle that best fits contour at position $s$.

↳ Curvature sign depends on if circle is inside or outside the figure

↗ can result from active contour

    ↳ Concavities linked with minima
    ↳ convexities linked with maxima

Properties of curvature-map descriptions:
    ↳ (1) Position-independent
    ↳ (2) Orientation-independent
    ↳ (3) Perimeter traversed in opposite direction by changing sign of $s$.
    ↳ (4) Scaling property: $\Theta(s) \to K\Theta(s)$ to scale an object.

## Codon Grammar:

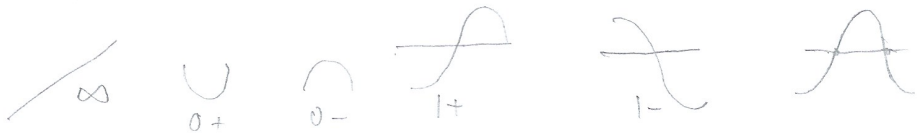used to create a taxonomy of closed shapes. Number of zero crossings as well as before or after maximum

Therefore, object recognition and classification as follows:

↳ Active contours to fit deformable

↳ Extract codon string from $\Theta(s)$ by traversing outline


0+    0-    1+    1-   

Therefore, can generate 3 codon pairs, 5 codon triples, 9 codon quads
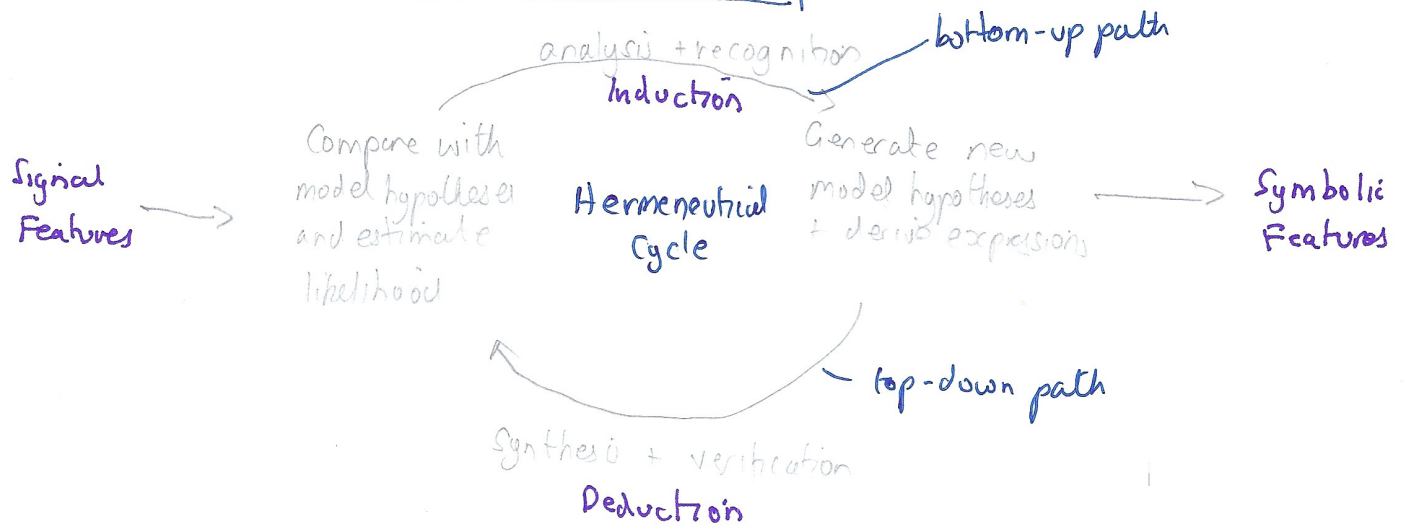
### Description of 3D shape

↳ Use codon string as index to lexicon

Superquadratics: represent objects as union a/o intersections of generalised superquadratic closed surfaces, loci of points in $(x, y, z)$ space:

↳ Object then classified by shape with lots of invariance

$\boxed{\text{Spheres: } A > B = C}$   $Ax^{\alpha} + By^{\beta} + Cz^{\gamma} = R$

Rotations produce cross terms in $(xy, xz, yz)$. Parameters define object dimensions.

# VISION AS MODEL BUILDING

analysis + recognition → bottom-up path
Induction

Signal Features →

Compare with model hypotheses and estimate likelihood

Hermeneutical Cycle

Generate new model hypotheses + derive expressions → Symbolic Features

~ top-down path

Synthesis + verification
Deduction

↳ human vision not veridical — illusions expected
↳ can learn from neurological traumas (aphasias and agnosias)

# BAYESIAN INFERENCE

Impossible to perform computer vision tasks in a bottom up fashion.
Can use bayesian method to use priors

 ↳① Some events more probable than others
 ↳② Matter doesn't disappear
 ↳③ Objects rarely change surface colour
 ↳④ Uniform texturing much more likely
 ↳⑤ Rigid rotation more likely then boundary deformation

 ↳ can apply the rule recursively using latest posterior as the new prior.

$$p(H|D) = \frac{p(D|H)p(H)}{p(D)}$$

<u>Statistical Decision Theory</u> : Pattern classification on basis of vector of acquired features.
 ↳ decide whether feature vector is consistent with a particular class.
 ↳ in 2-state decision problem, feature vectors arise from overlapping probability distributions.

For OCR, slide 161 onwards

$$\text{detectability} = d' = \frac{|\mu_2 - \mu_1|}{\sqrt{\frac{1}{2}(\sigma_2^2 + \sigma_1^2)}} \quad (\geqslant 3 \text{ is normally considered good})$$
(discriminability)

For each class separately, measure how likely any sample value is: $P(x|C_k)$

$$P(x) = \sum_k P(x|C_k)P(C_k)$$

Posterior $P(C_k|x) = \frac{1}{P(x)} \underbrace{P(x|C_k)}_{\substack{\text{class conditional} \\ \text{likelihoods}}} \underbrace{P(C_k)}_{\text{Priors}}$

Minimise total probability if assign each observation to class with highest posterior
 ↳ can rewrite minimum by assuming denominator in Bayes' probability rule is independent of $C_k$

$$P(x|C_k)P(C_k) > P(x|C_j)P(C_j) \quad \forall j \neq k$$

<u>Discriminant Functions</u> : construct [set of functions] $y_k(x)$ of data $x$, one function for each class $C_k$, st classification decisions made by assigning $x$ to
$C_k$ if : $y_k(x) > y_j(x) \quad \forall j \neq k$
 ↳ discriminant functions ⇒ normally posterior prob. functions : $P(C_k|x)$
    or : $P(x|C_k)P(C_k)$

<u>Discriminative Methods</u>: Learn function $y_k(x) = P(C_k|x)$ that maps features $x$ to class labels $C_k$.
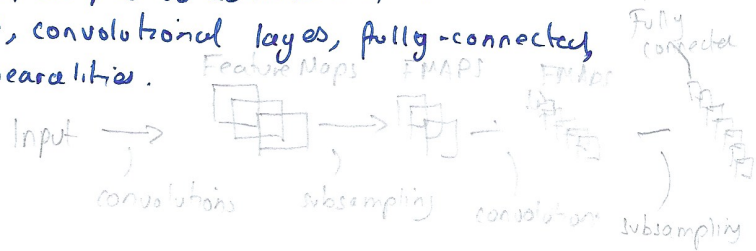
<u>Generative Methods</u>: learn likelihood model expressing prob data features $x$ would be observed in instances $C_k$, which can be used for classification using Bayes' rule. Have predictive power as allow samples from joint distribution $P(x, C_k)$

## Convolutional Neural Networks : Feed forward artificial neural networks.
- ↳ Multiple layers of small collections of neurons
- ↳ Tiling and overlapping of outputs to achieve shift invarious
  - ↳ pooling layers, convolutional layers, fully-connected, point non-linearalities.
- ↳ Little pre-processing

## OCR
CNN designed by Yann LeCun (slide 170)
- ↳ input is 32×32.
- ↳ Trained with 100,000+ examples, using supervised back-propogation. Target output +1, rest to -1. Error back propogate to produce feature maps. Neurons have 5×5 kernels, convolved with input
- ↳ Trained to extract visual feature. Subsequent feature maps achieve size, slant and style invariances. Neurons in final layer identify input as a target.

Output of each neuron at $(i,j)$ applies non-linear activation function $f_{act}$ to sum of ib input pixels × weights $w_{mn}$ and bias term:

$$ O_{ij} = f_{act} \left( w_o + \sum_m \sum_n w_{mn} \, I_{(i-m),(j-n)} \right) $$

## FACE DETECTION, RECOGNITION, IDENTIFICATION

Facial detection very challenging : is within-class variation > between-class variation
- ↳ pose, illumination, family, time
- ↳ Treat as 3D problem or 2D problem.

### Viola-Jones Face Detection:  , 30± layers
Use cascade of weak classifiers to build a strong detector.
- ↳ feature detector with 2D Haar Wavelets. - multiplication not required, therefore quicker

$$ h_j(x) = \begin{cases} -p_j & \text{if } f_j < \theta_i \\ p_j & \text{else} \end{cases} \quad , \quad h(x) = \text{sign} \left( \sum_j \alpha_j h_j \right) $$

At intermediate point, face provisionally accepted if $h(x) > 0$. Only those accepts passed onto next layer

AdaBoost: supervised, adapt weights such that early layer have high accept rates and later one more discriminatory
- ↳ cascade evaluated using sliding window approach

### Gabor Wavelets : act as effective compact code
- ↳ features represented with handful of wavelets

## Eigenfaces

↳ database of pre-normalised for size, position, foontal pose, decomposed into Principal Components as sequence of orthonormal eigenfunctions with descending eigenvalues

↳ Extract 20 most eigenfaces and for presentry photos, project onto eigenfaces and store coefficients

↳ **Accurate (90 - 95%) + fast to use**

↳ But pose and illumination
↳ deal with this by brute force, having lots of cameras.

## 3D Approaches

Need shape model     and     texture model
↳ laser, LIDAR, stereo cameras,     Project texture onto shape model
multiple images

↳ Can then be used by extracting correct pose to do 2D comparisons

2017 Face
Recognition Comp
tested with non-ideal
images, etc

## FaceNet : CNN with 22 layers and 140 mn parameters  ↗ using back-propagation
↳ trained on 200 mn face images ⇒ 8mn identities ∴ 2,000 hours training
↳ use Euclidean distance as metric.
↳ Use triplets of images – one pair from same person, minimise loss function $L$

$$L = \sum_i \left[ \| f(x_i^a) - f(x_i^p) \|^2 - \| f(x_i^a) - f(x_i^n) \|^2 \right]$$

↳ Embedding create compact code for each face
↳ Euclidean distance gives decision of same vs different

## Affective Computing : faces used for emotions – lots of brain to interpret other's faces
↳ use MRIs to show brain areas as interpretting different facial expressions

## Facial Action Coding System : taxonomy of facial expressions
↳ 32 Action Units by 7 muscles

↳ 14 Action Descriptors – use message judgement to use AUs and ADs to get meaning
↳ ① Pre-processing – face detection + normalisation
↳ ② Feature Extraction (appearance based or using (2) spatio temporal ideas
↳ ③ AU temporal segmentation, classification, intensity estimation.

generative models ⤴     ⤴ deformable discriminate methods     , Issues
infer state from muscular     fit a deformable model.
actions     ① Small dataset, often not reliable as well
② Manual scoring required.