

Using vision to understand the world by the combination of learning and innate mechanisms.



Understanding the world through vision



'Digital Baby'

Input: image sequences

Innate capacities



197	178	172	127	165	173
218	204	202	196	193	194
215	198	186	180	187	184
218	199	203	195	191	156
167	170	134	193	106	110
95	157	160	204	168	151
192	197	203	197	187	175



219	188	218	204	202	196
190	235	215	198	186	180
163	223	218	199	203	195
210	224	167	170	134	193
226	179	95	157	160	204
216	193	192	197	203	197
218	221	204	203	186	218



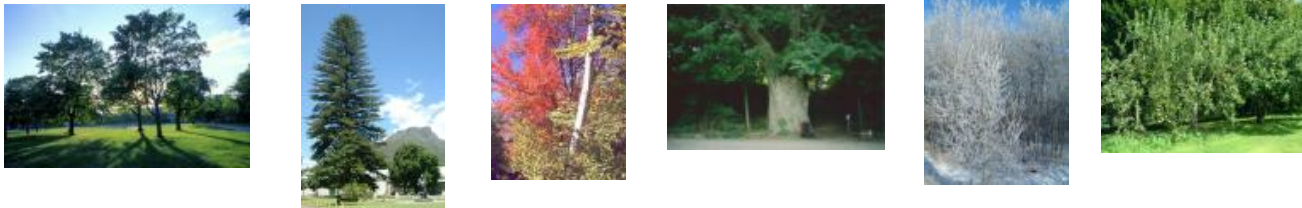
184	113	118	105	117	82
151	95	122	131	87	100
160	156	159	197	178	172
136	219	188	218	204	202
184	190	235	215	198	186
175	163	223	218	199	203
221	210	224	167	170	134



Objects ←
Actions
Agents
Goals
Tools
Interactions

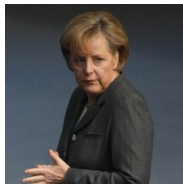
Learn visual interpretation

Object Categories

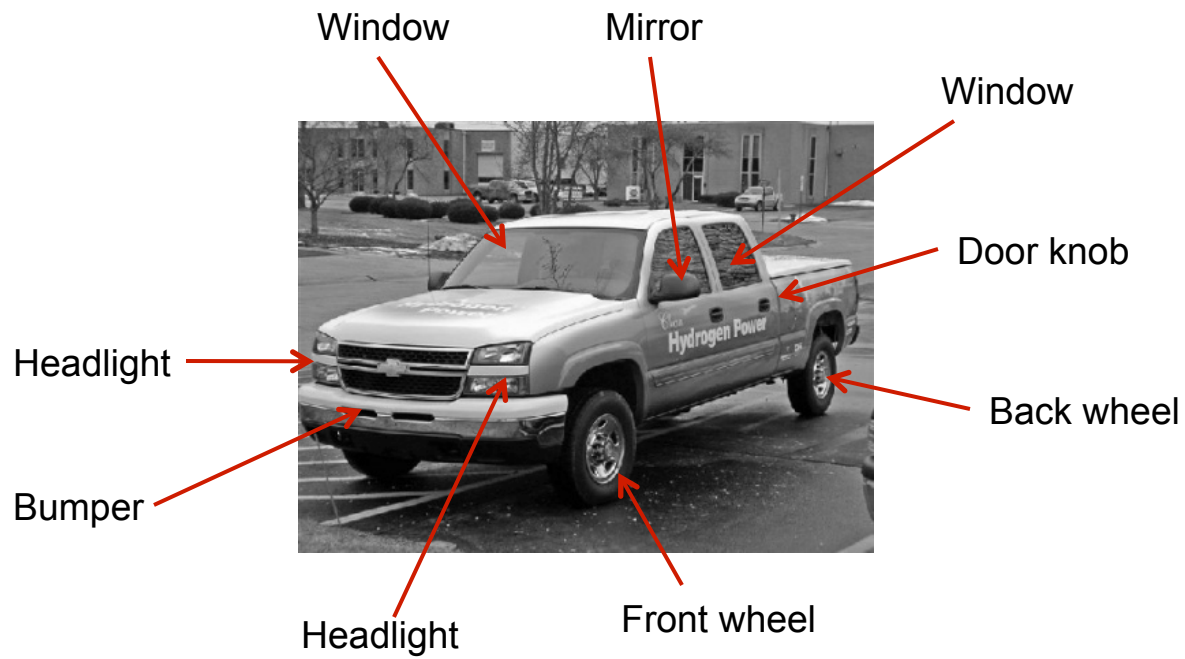


- We perceive the world in term of objects and classes
- Large variability within a each class

Individual Recognition



Object parts

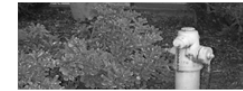


I. Categorization: dealing with class variability

Class



Non-class



We learn categories from examples

Class

81	82	85	70	195	247	91
133	82	85	88	76	89	88
89	110	92	85	85	85	89
160	69	129	87	85	85	88
92	82	125	79	89	89	89
128	129	106	99	93	93	90
185	127	122	135	162	93	90

197	178	172	127	165	173
218	204	202	196	193	194
215	198	186	180	187	184
218	199	203	195	191	156
167	170	134	193	106	110
95	157	160	204	168	151
192	197	203	197	187	175

219	188	218	204	202	196
190	235	215	198	186	180
163	223	218	199	203	195
210	224	167	170	134	193
226	179	95	157	160	204
216	193	192	197	203	197
218	221	204	203	186	218

Non-class

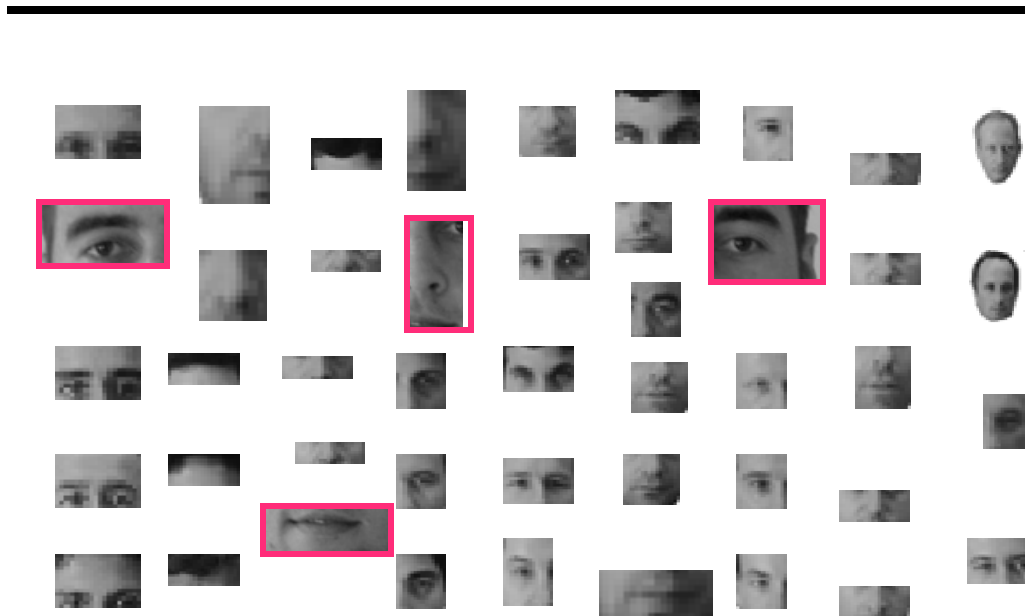
164	171	161	215	88	96	89
174	178	165	255	123	91	89
189	179	189	255	113	88	89
206	151	170	0	86	96	91
81	162	251	252	14	38	97
154	236	252	214	8	33	20
249	244	249	224	16	24	5

184	113	118	105	117	82	:
151	95	122	131	87	100	:
160	156	159	197	178	172	:
136	219	188	218	204	202	:
184	190	235	215	198	186	:
175	163	223	218	199	203	:
221	210	224	167	170	134	:

202	196	193	194	185	127
186	180	187	184	190	199
203	195	191	156	176	196
134	193	106	110	76	174
160	204	168	151	176	195
203	197	187	175	169	185
186	218	183	156	179	201

Natural for the brain, difficult computationally

Visual Class: Similar Configurations of Shared Components

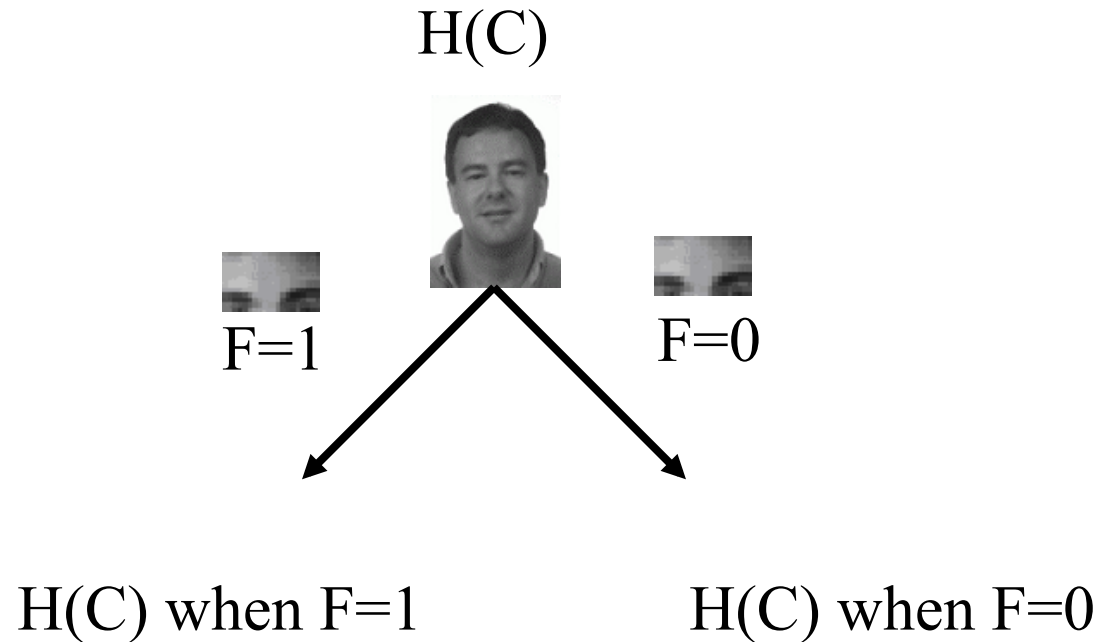




What will be optimal building-blocks for the class?

Mutual information


$$H(c) = -\sum P(c) \text{Log}(P(c))$$



$$I(C;F) = H(C) - H(C/F)$$

Mutual Information $I(C,F)$



Class:	1	1	0	1	0	1	0	0
Feature 	1	0	0	1	1	1	0	0

$$I(F,C) = H(C) - H(C|F)$$

$$\sum_{C,F} p(C,F) \log \frac{p(C,F)}{p(C)p(F)}$$

Optimal classification features

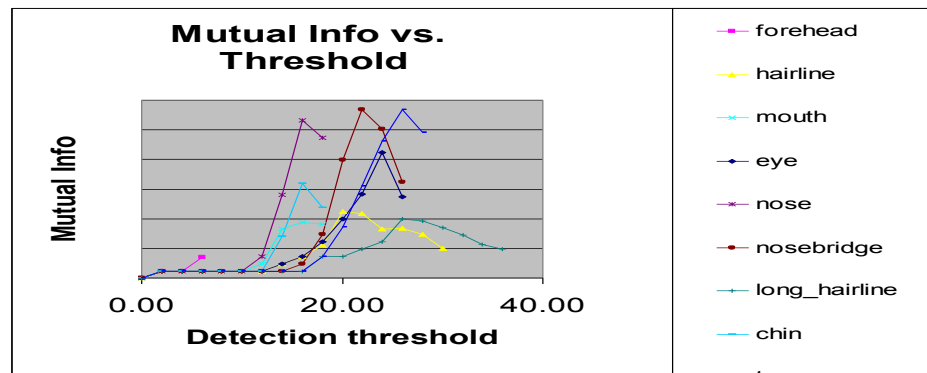
- Theoretically: maximizing delivered information minimizes classification error

$$\text{Error} = H - I(C;F)$$

- In practice: informative object components can be identified in training images



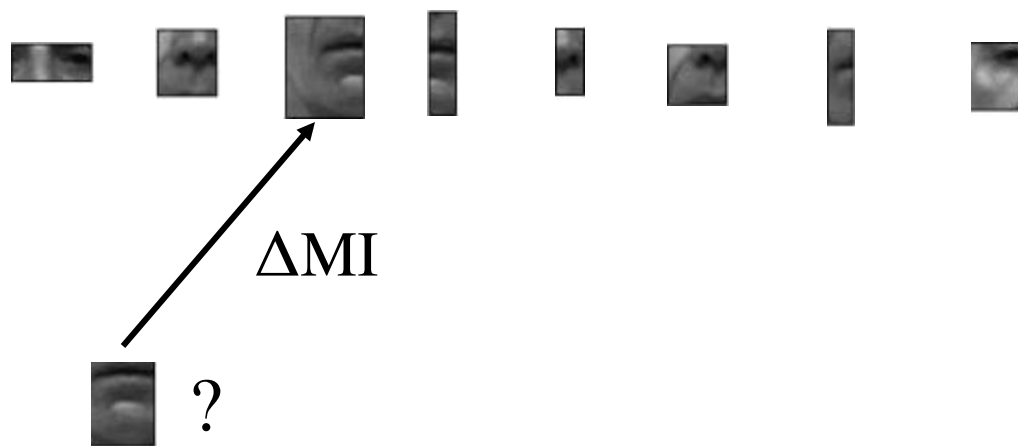
Selecting Fragments



'Imprinting' many receptive fields and selecting a subset

Adding a New Fragment

(Avoiding redundancy by max-min selection)



Compare new fragments F_i to all the previous ones.

Select F which maximizes the additional information

$$\text{Max}_i \text{Min}_k \Delta MI (F_i, F_k)$$

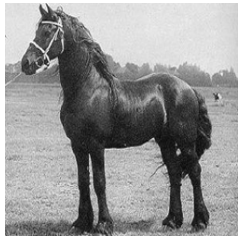
Competition between units with similar responses

Highly Informative Face Fragments



Optimal receptive fields for Faces

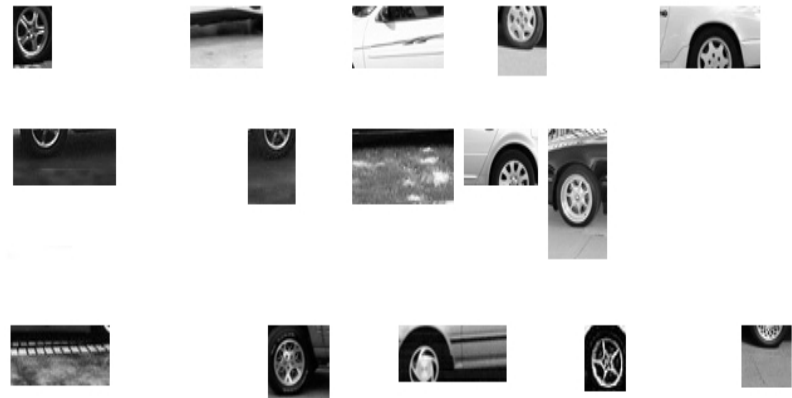
Informative class features



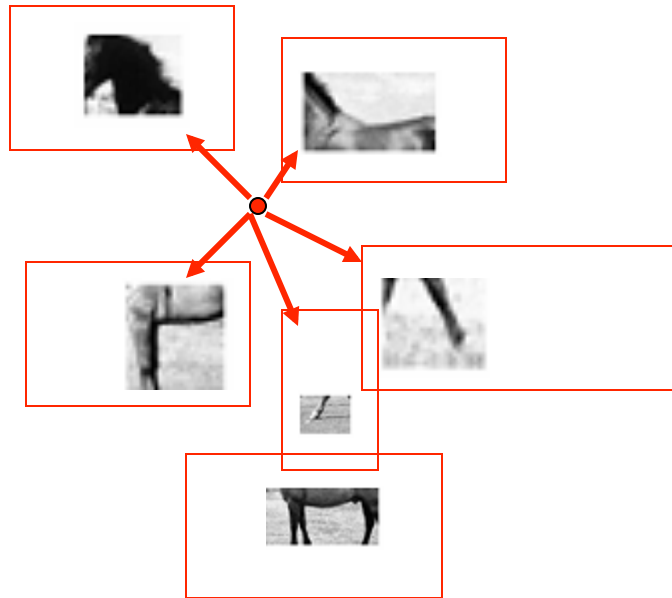
Horse-class features



Car-class features



Fragments with positions



$$\sum w_k F_k > \theta$$

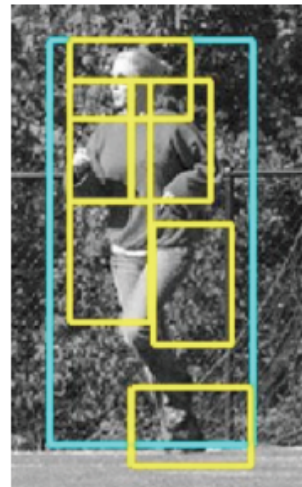
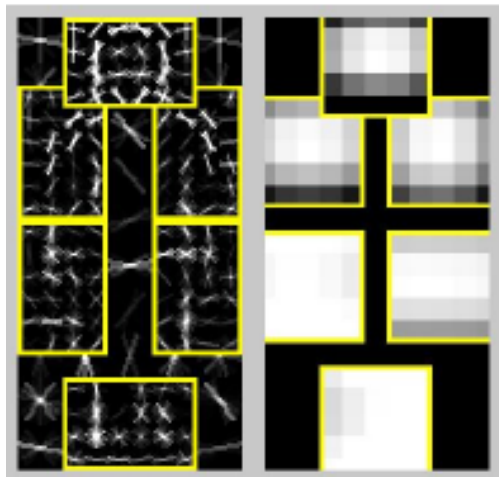
On all detected fragments within their regions

Variability of Horses Detected

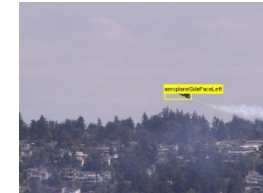
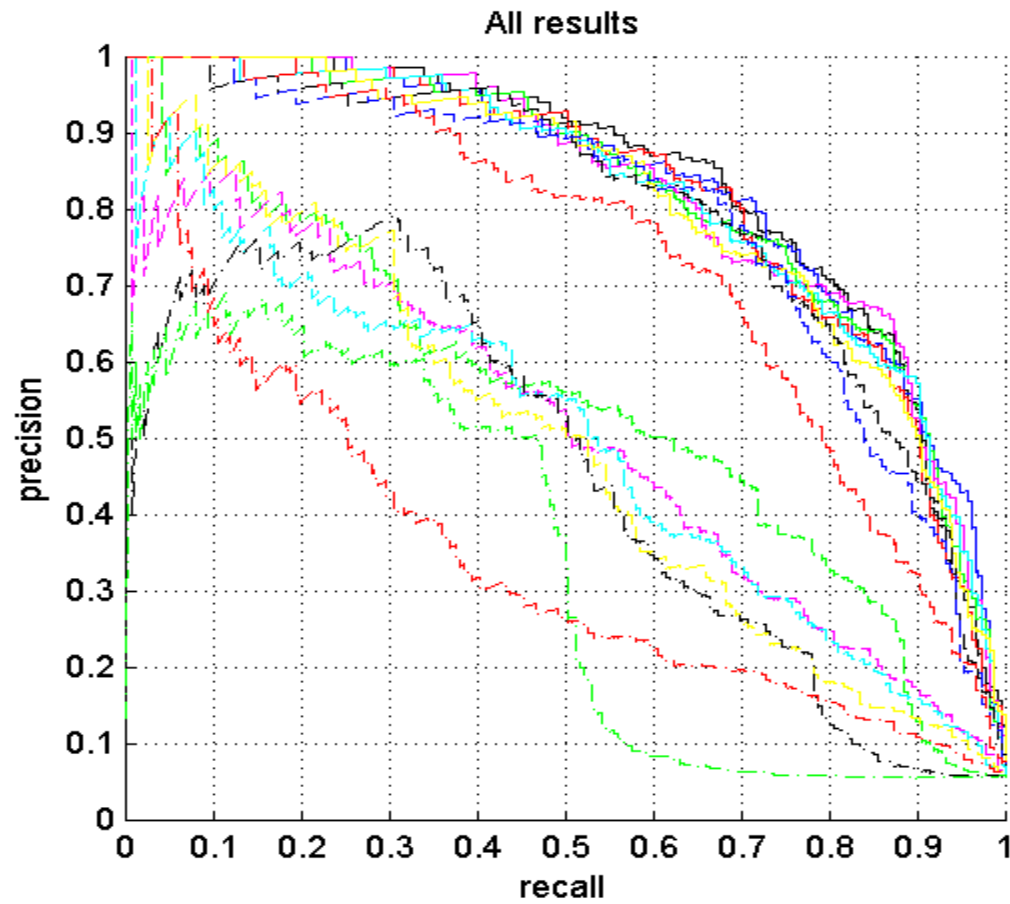


Using informative class-fragments to deal with variability

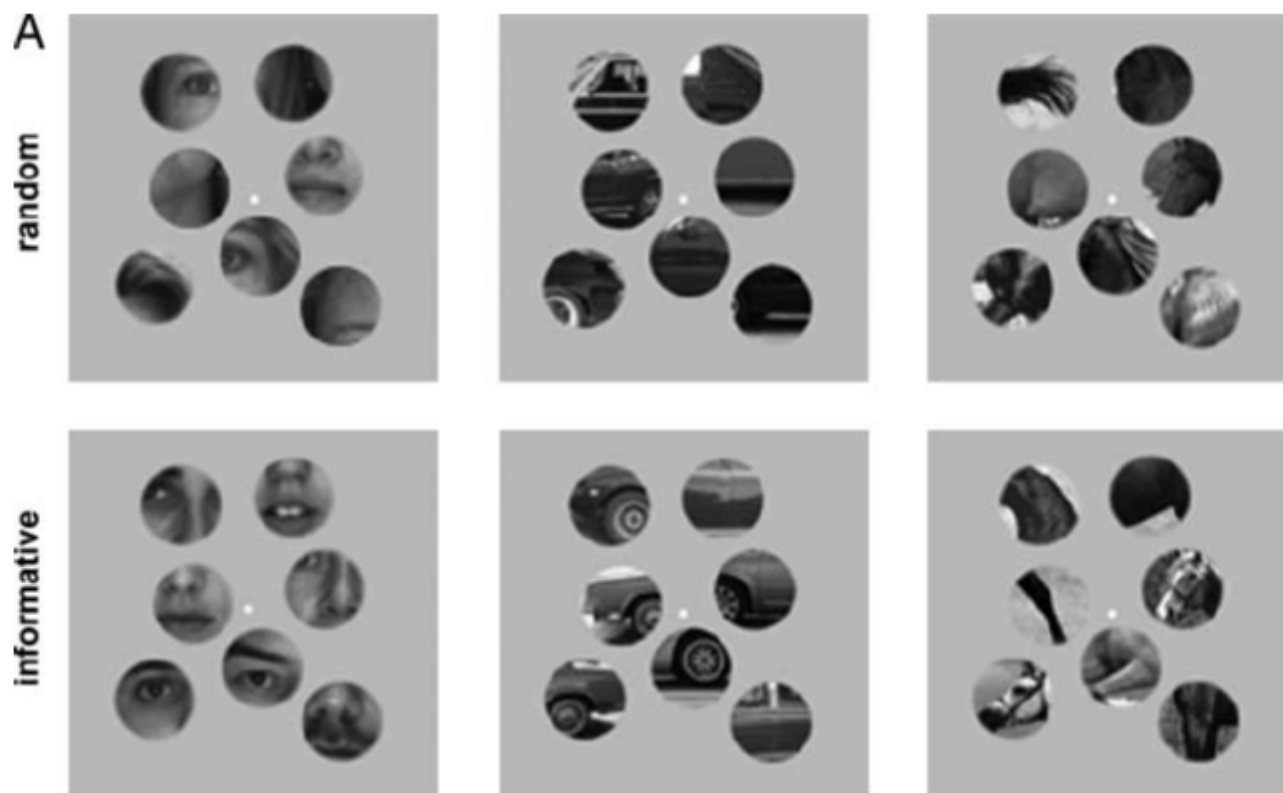
Felzenszwalb et al (Pascal 2007)



'Pascal Challenge' Airplanes

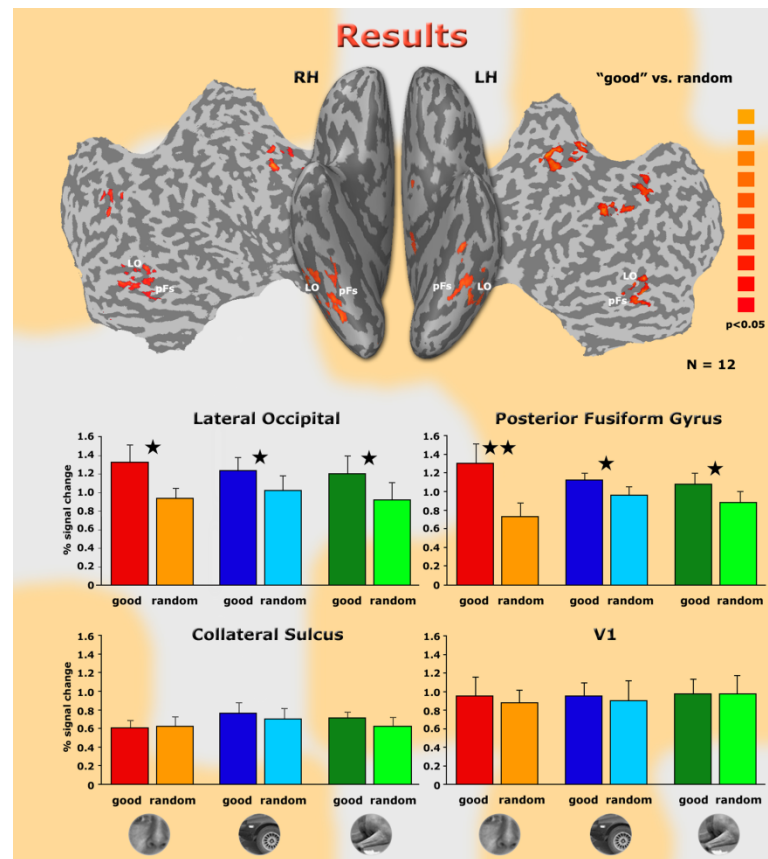


Looking for Class Features in the Brain: fMRI



Lerner, Epshtein Ullman Malach JCON 2008

Class-fragments and Activation



Malach et al 2008

ERP



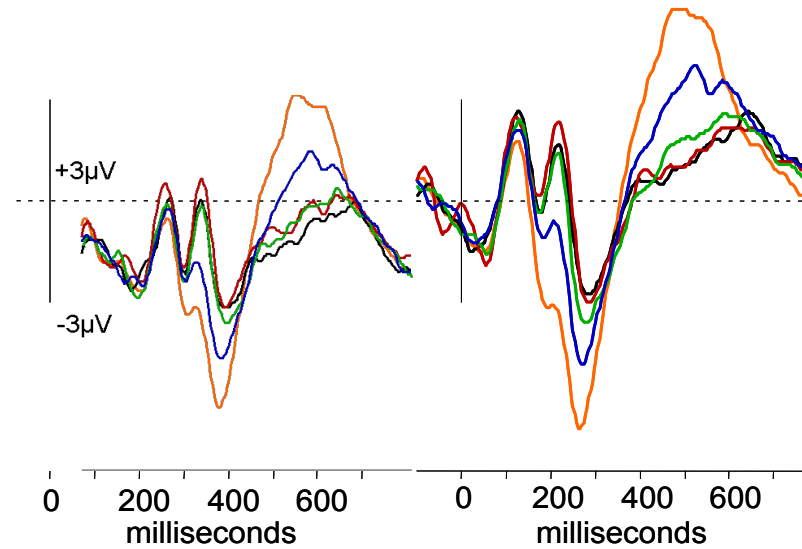
FACE FEATURES

Posterior-Temporal sites

Left Hemisphere

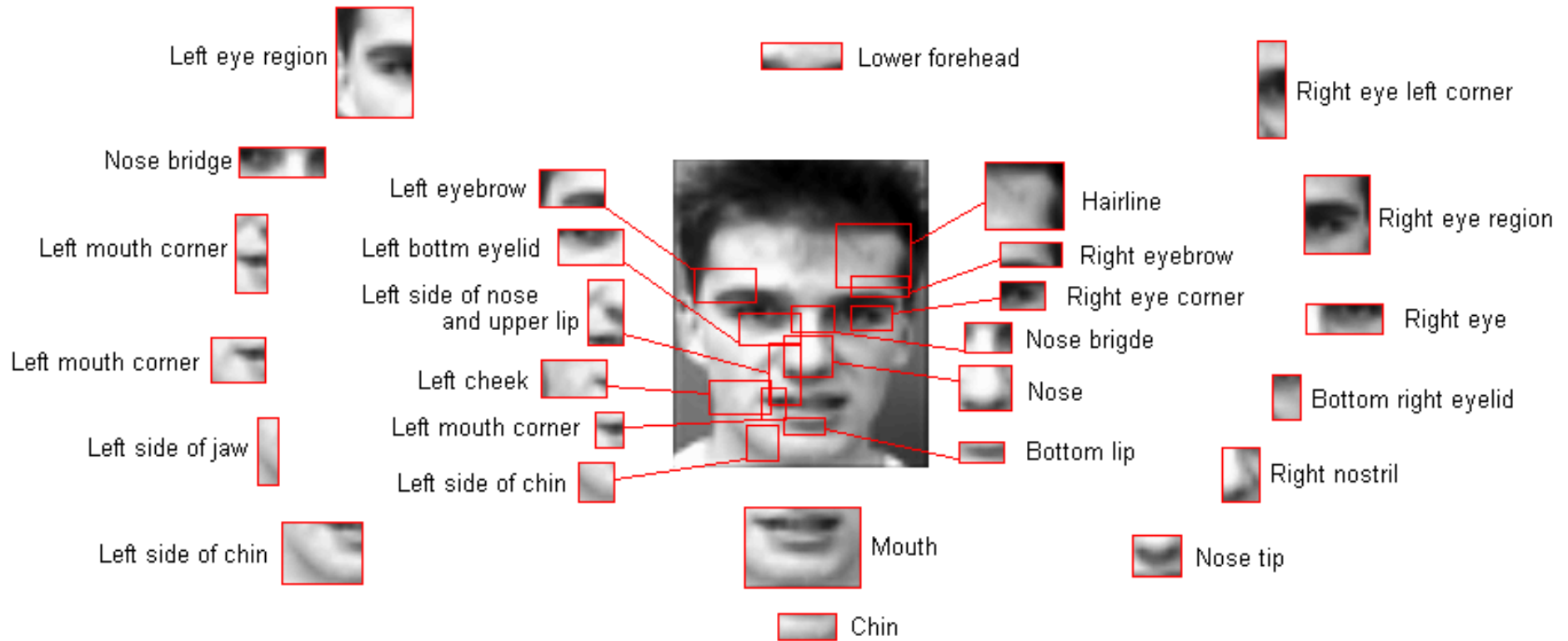
Right Hemisphere

MI 1 —
MI 2 —
MI 3 —
MI 4 —
MI 5 —



‘Object interpretation’:

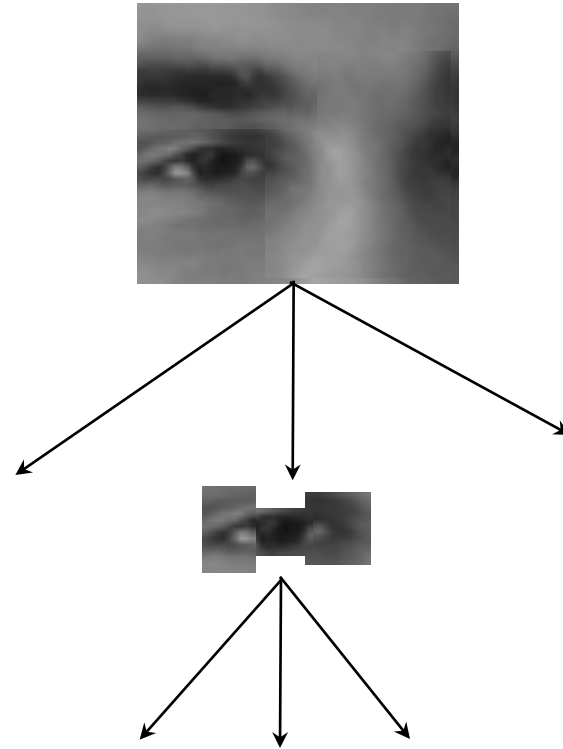
Full recognition of parts and sub-parts at many levels



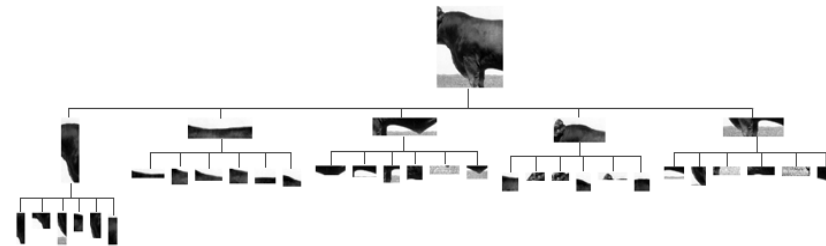
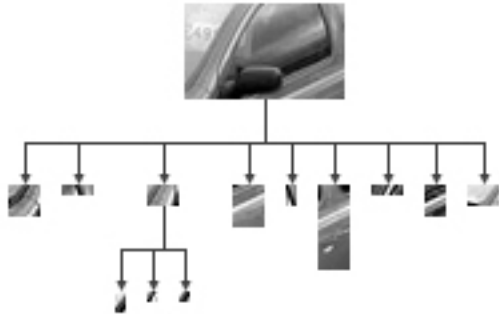
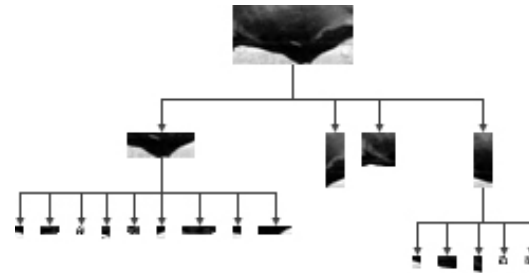
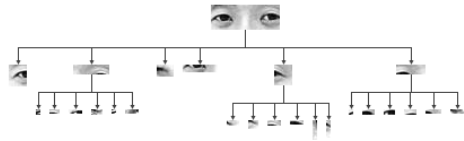
Hierarchies of sub-fragments

Detect the part by simpler sub-parts

Repeat at multiple levels, to obtain a hierarchy of parts and sub-parts

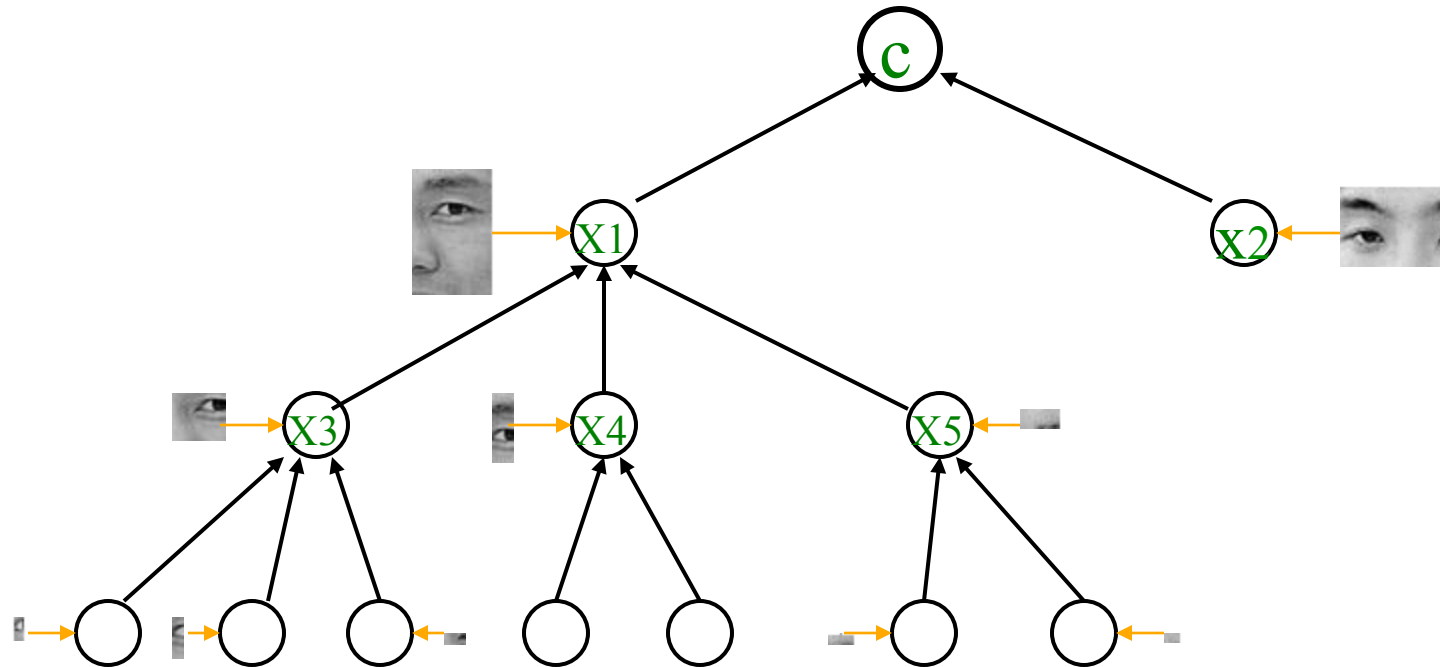


Example Hierarchies



Classification and detecting all the parts

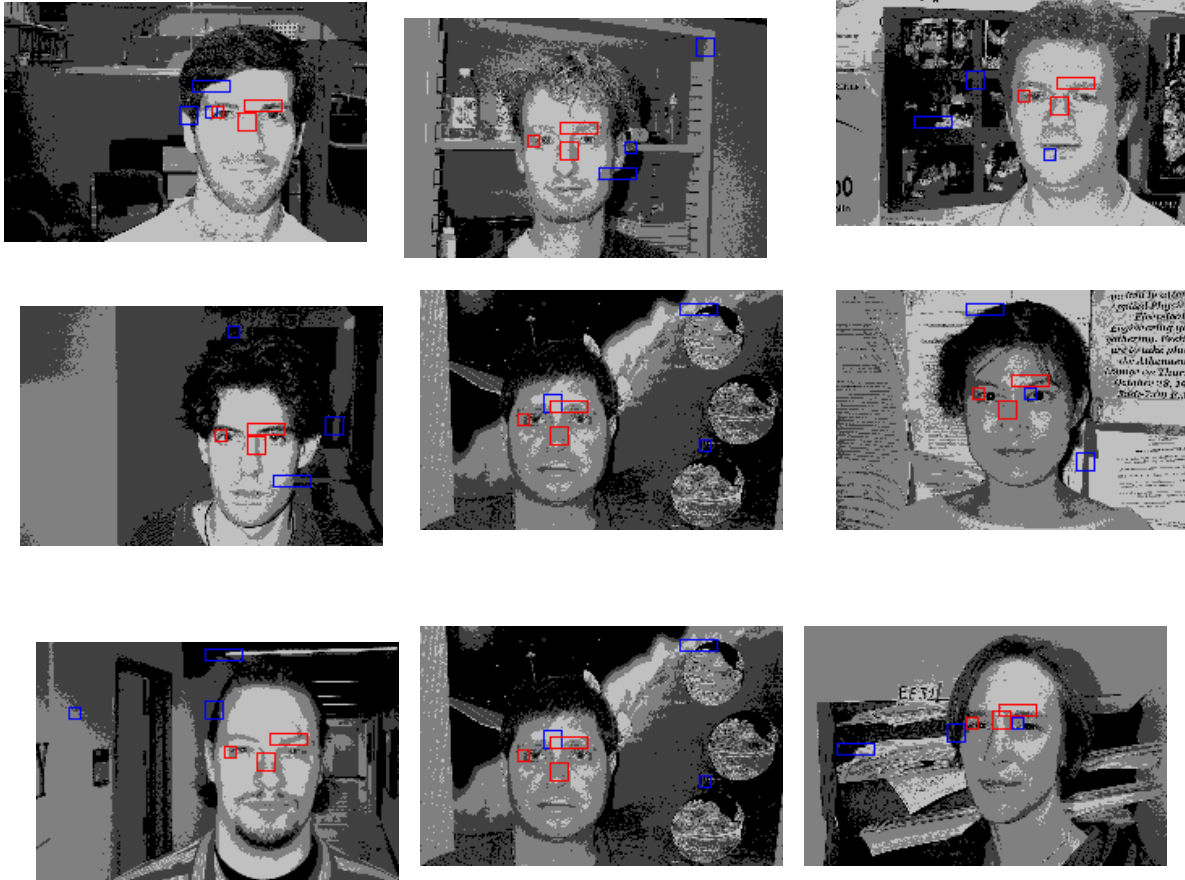
Two pass computation: BU then TD



- The full hierarchy is used for recognition
- A bottom-up then top-down pass

$$p(c, \underline{x}, \underline{F}) = p(c) \prod p(x_i | x_{i^-}) p(F_i | x_i)$$

Difficult parts



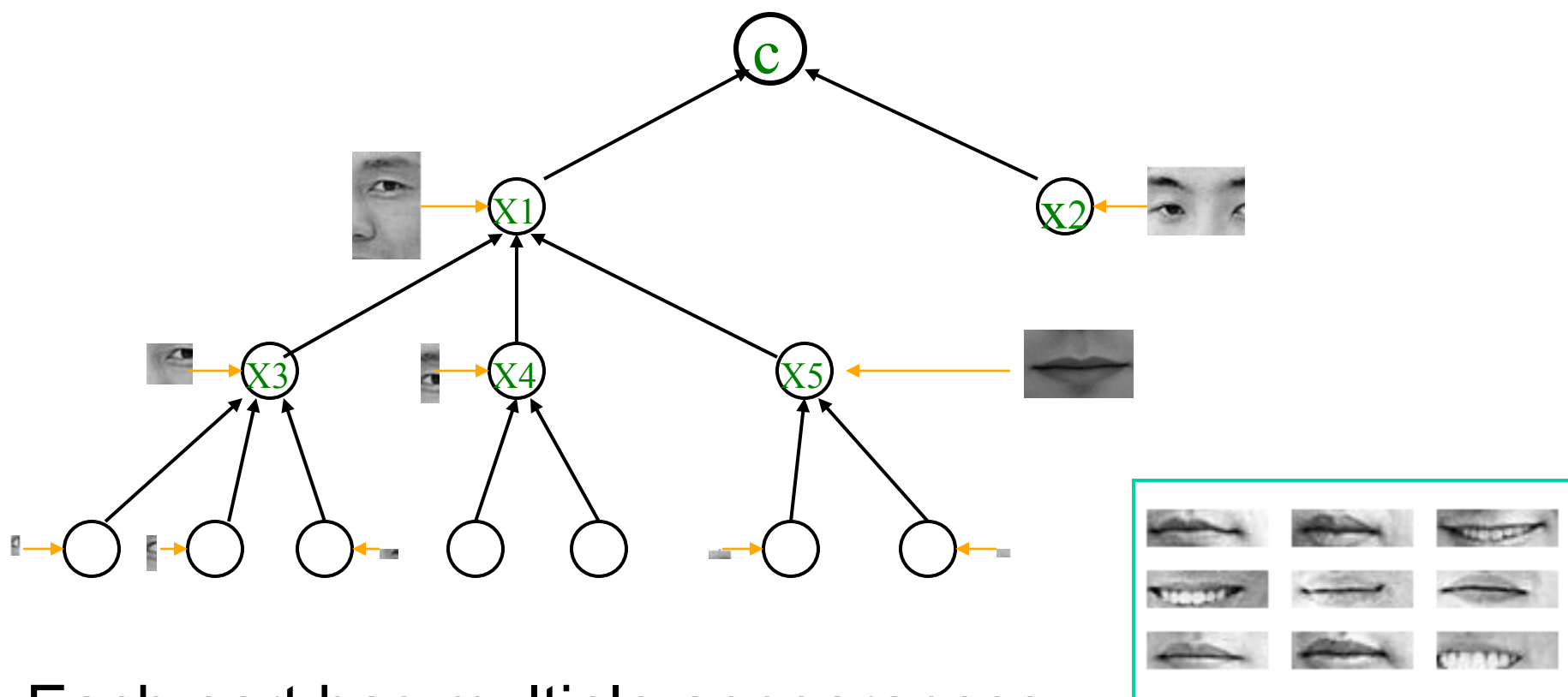
Eyebrow, nose-bridge, eye-corner
All missed by the feed-forward pass

III. Individual Recognition

Same Object Under Different Views



From Appearance to Semantic Features



Each part has multiple appearances

Semantic Features

- Group together alternative appearances:

- *Dynamic*: by motion

- *Static*: by common context

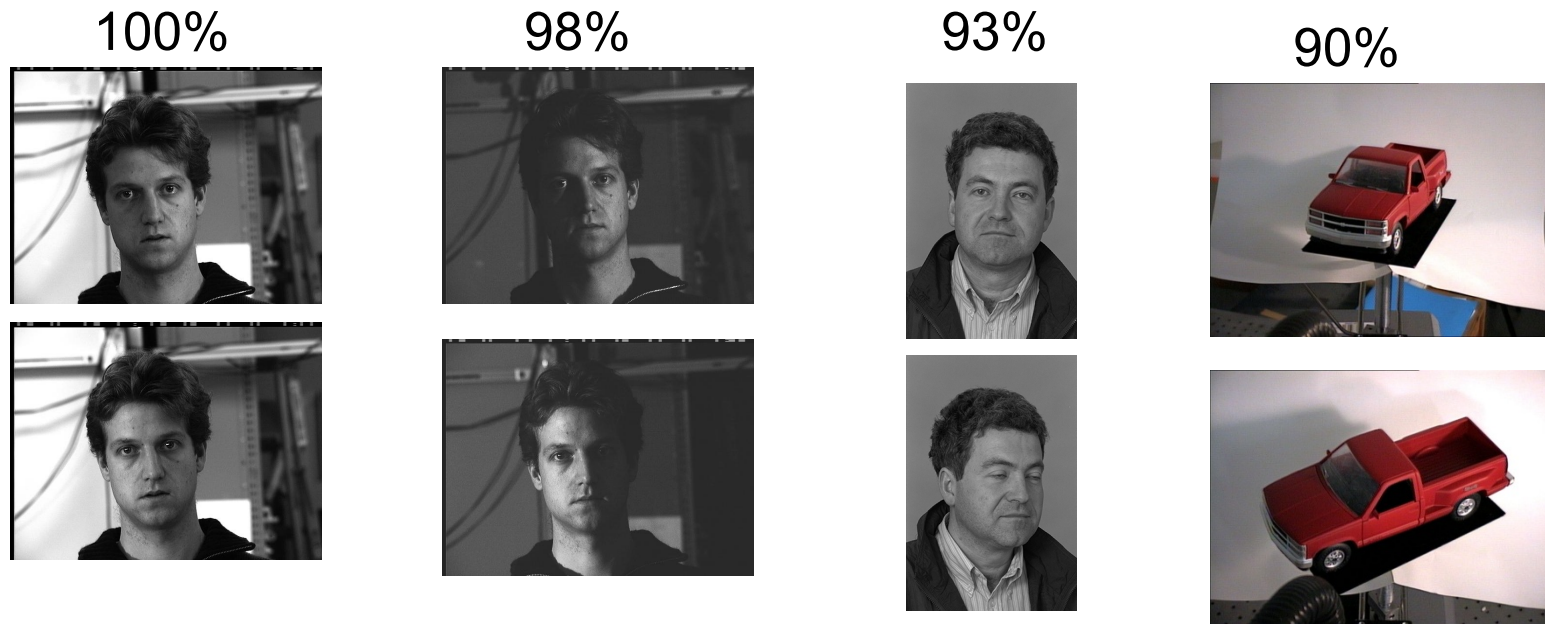


Semantic Features: Examples

Emergence of abstract representation

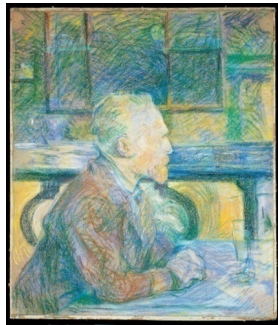


Individual Recognition by Semantic Fragments

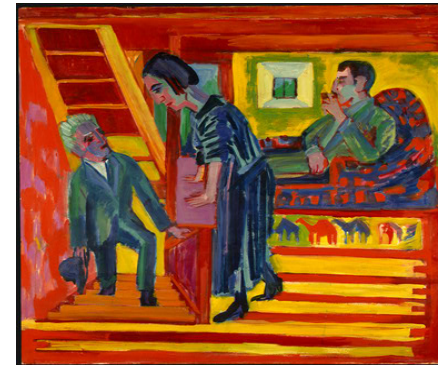


Invariance to individual objects is inherited from the learned invariance of semantic fragments

Full Interpretation, Difficult Parts: Finding Hand by Context



Van Gogh

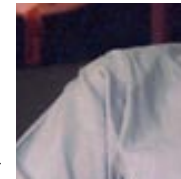


Kirchner

Hand by Context



Face



Face → Shoulder → Upper-arm → Lower-arm → Hand

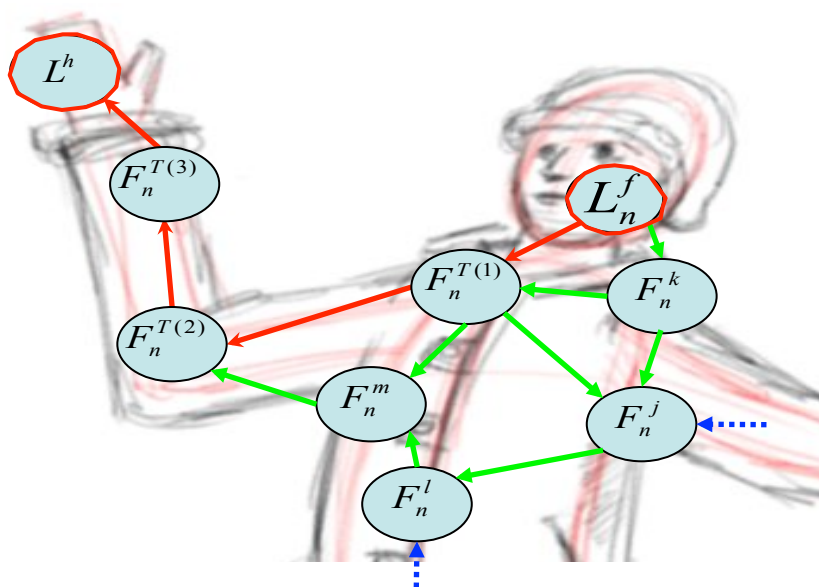
Use appearance to learn the context

The 'Chains Model' for learning useful chains from sources to targets

Chains model

$$\underline{E} \rightarrow \underline{X}$$

$$p(\underline{E}, \underline{X})$$



$$p_{\tau}(\text{Hand, Face, Features}) =$$

$$p(\text{Face}) p(F_1|\text{Face})p(F_2|F_1)\dots p(\text{Hand}|F_k) \quad \prod p(F_j)$$

(On-chain)

(Off-chain)

$$p(\text{Face}) p(F_1|\text{Face})p(F_2|F_1)\dots p(\text{Hand}|F_k) * \prod p(F_j)$$

After dividing by the product $\prod p(F_j)$:

$p(\text{entire configuration}) :$

$$p_T(\text{Hand, Face, Features}) \sim \frac{P_c(F_i, F_{i+1})}{\prod P(F_i) \cdot P(F_{i+1})} \quad \text{Links along the chain}$$

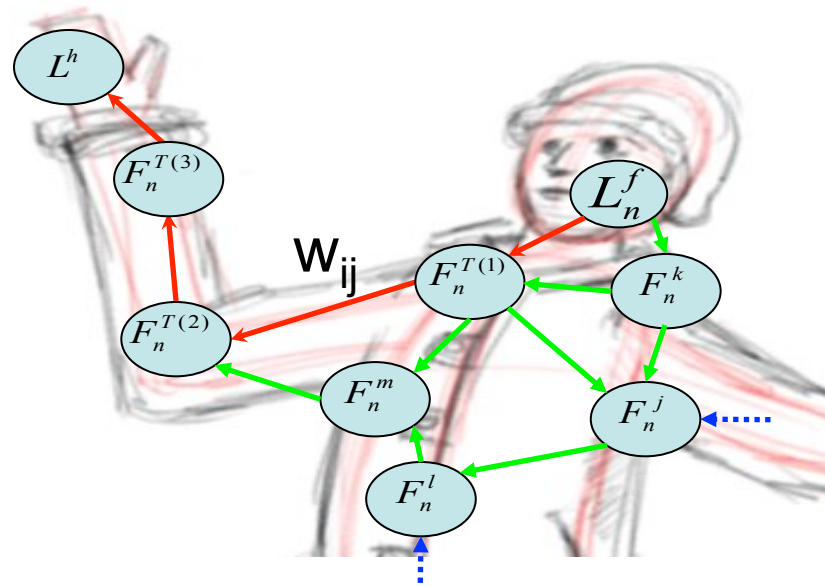
Or: \prod 'Suspicious coincidence' (F_i, F_j)

Finally, finding the hand -- combine all chains:

$$p(\text{Hand} | \text{Face}, \underline{F}_k) \sim \sum_T p_T(\text{hand, Face}, \underline{F}_k) = \sum_T \prod_T w_{ij}$$

The final computation:

Chains
model



$$\frac{P_c(F_i, F_{i+1})}{P(F_i) \cdot P(F_{i+1})}$$

- Features are connected by a lateral weight W_{ij}
- Features that occur frequently together will have high W_{ij}
- The probability of a chain is $\prod w_{ij}$ along the chain.
- The computation finds high likelihood chains between the source and the target (face and hand).
- Propagate from known parts to new parts

The chains model inference

$$P_{CH}(L_H | L_F, \{F_i\}) \propto \sum_T P_{CH}(T, L_H, L_F, \{F_i\}) = D \cdot C^* \cdot S \approx D \cdot C \cdot S$$

$$D = [P_{obj}(L_H, F_1), \dots, P_{obj}(L_H, F_K)], C = \left(\sum_{m=0}^{M-1} A^m \right), S = \left[\frac{P_{obj}(L_F, F_1)}{P_{obj+bg}(F_1)}, \dots, \frac{P_{obj}(L_F, F_K)}{P_{obj+bg}(F_K)} \right]^T$$

Hand

Internal

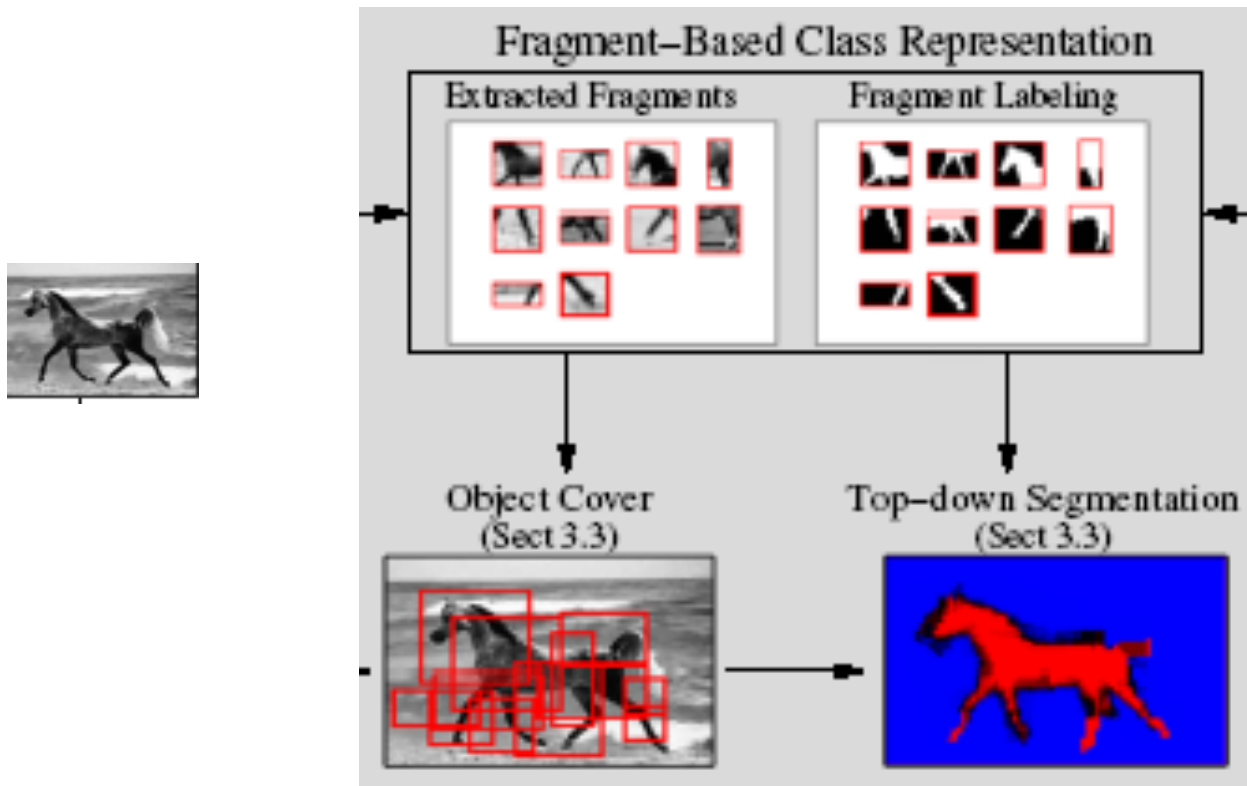
Fac

e

- The inference boils down to matrix multiplication approximating marginalization on chains of length $\leq m$.
- A is the weighted adjacency matrix of the feature graph, w_{ij} .
- C^* is the marginalization over only the simple chains (NP-hard), approximated by C which is the marginalization over all chains.



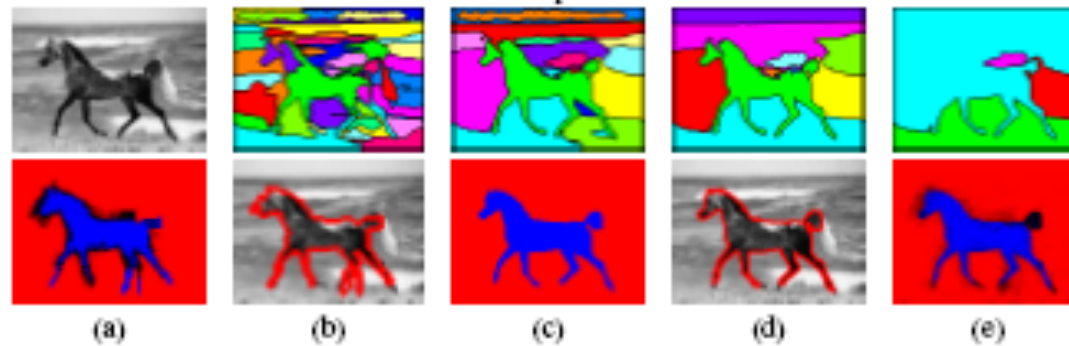
Class-based Segmentation



Top-down Segmentation



Combining TD with BU



Segmentation: combining TD with BU

Summary

- Informative object fragments
- Class specific
- Continuously extracted from examples
- Hierarchy of informative parts
- Grouped into semantic features (motion, context)
- Categorization, identification, segmentation
- Bi-directional interpretation process BU-TD