

# Machine Learning for Computer Vision

[MVA 2015/2016]

## Final assignment

Maha ELBAYAD

### 1 Detector training

We trained different classifiers: Linear regression, L2-logistic regression (c.f. assignment 1), linear SVM and RBF-SVM (c.f. assignment 2) to detect particular points of a facial image namely the left eye, right eye, left mouth, right mouth and nose. Each of these classifier is either trained on SIFT features or CNN features extracted with an Alexnet pre-trained network<sup>1</sup> shown in figure (1) where the features are either the output of the relu2 layer denoted CNN2 or that of the relu5 layer denoted CNN5.

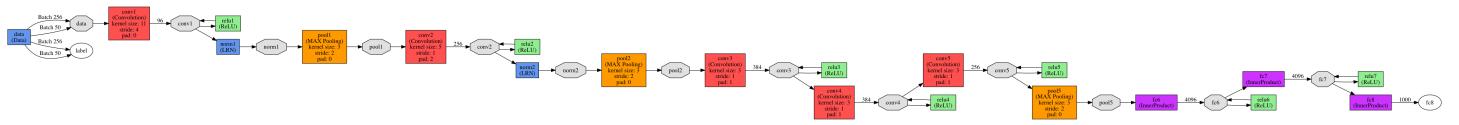


Figure 1: Alexnet architecture

If needed we choose the classifier's optimal parameters with cross-validation (5-folds).

#### 1.1 Precision recall curves:

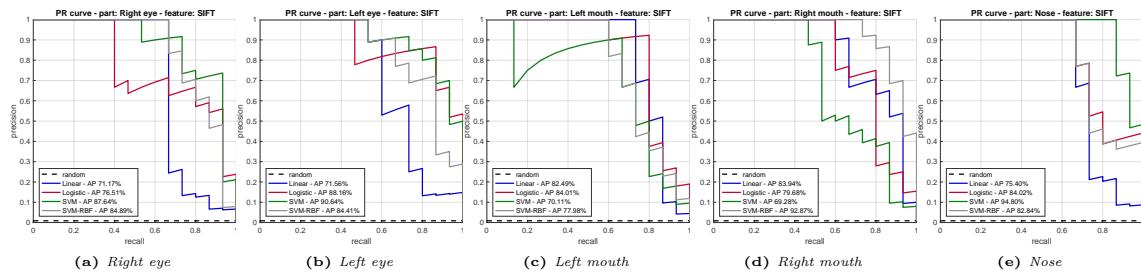


Figure 2: SIFT features

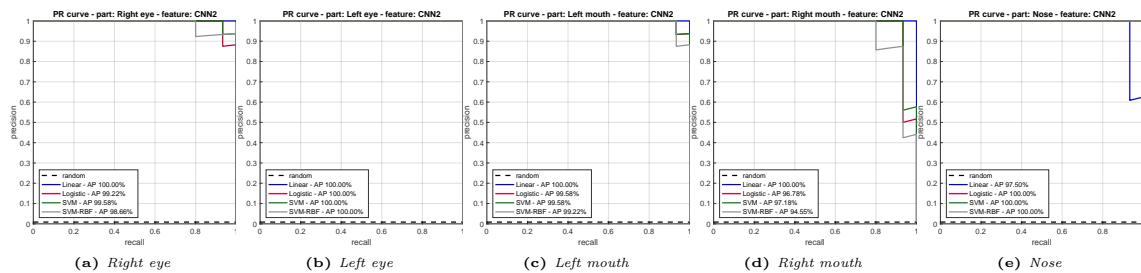
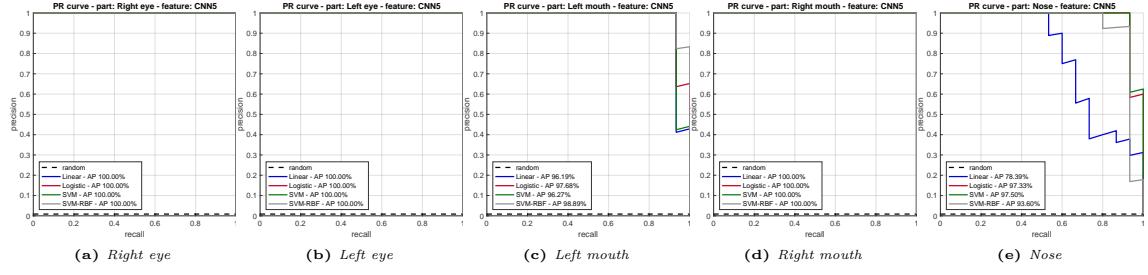


Figure 3: CNN2 features

<sup>1</sup><http://www.vlfeat.org/matconvnet/models/beta16/imagenet-caffe-alex.mat>



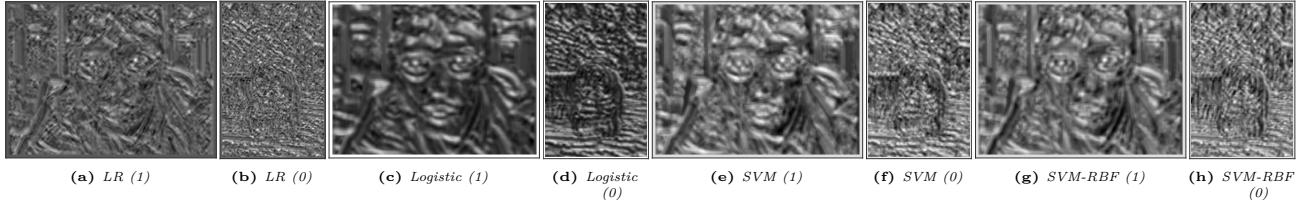
**Figure 4:** CNN5 features

## 1.2 SIFT features - Dense scores

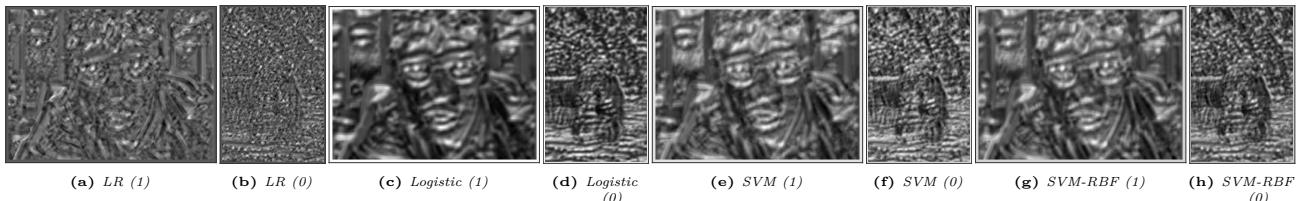
For two samples from different categories we show the dense scores of each classifier.



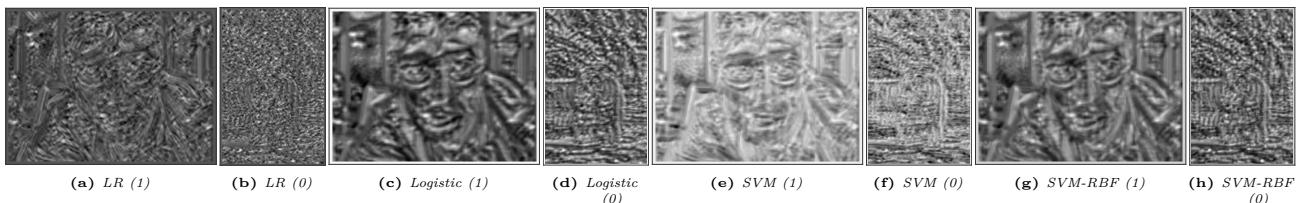
**Figure 5:** Input images



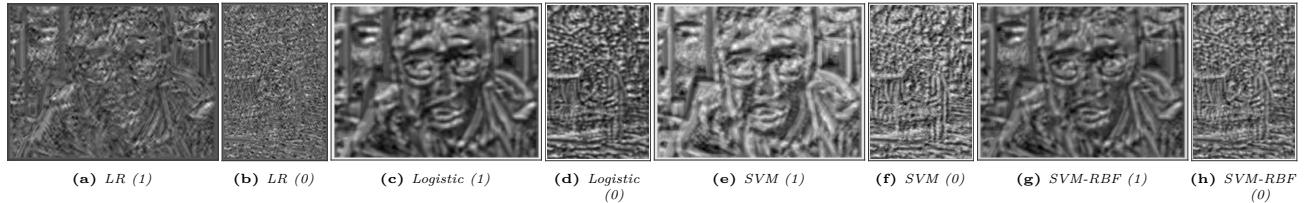
**Figure 6:** Right eye



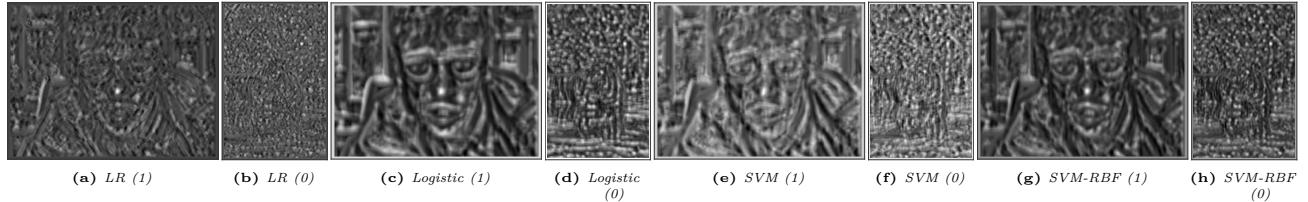
**Figure 7:** Left eye



**Figure 8:** Left mouth

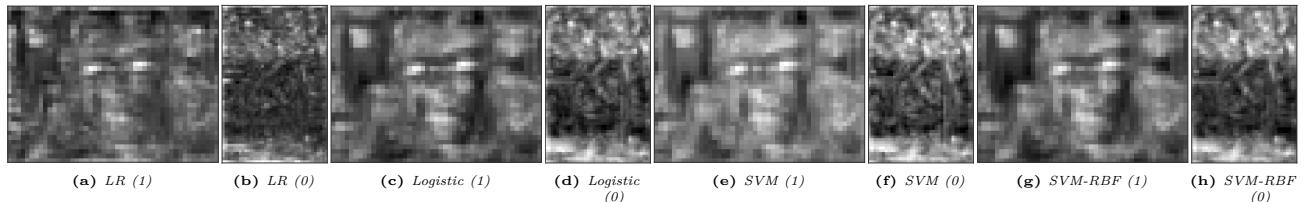


**Figure 9:** Right mouth

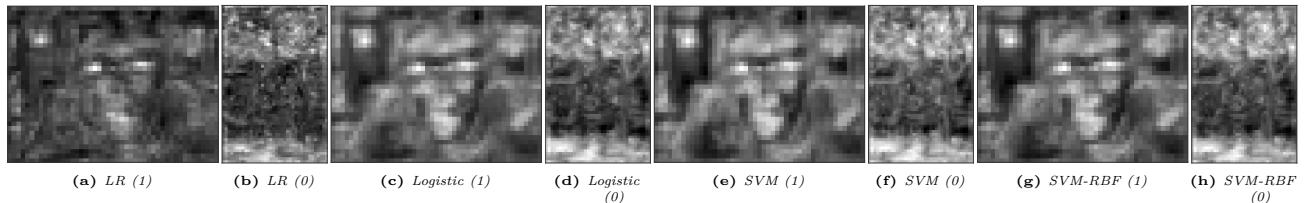


**Figure 10:** Nose

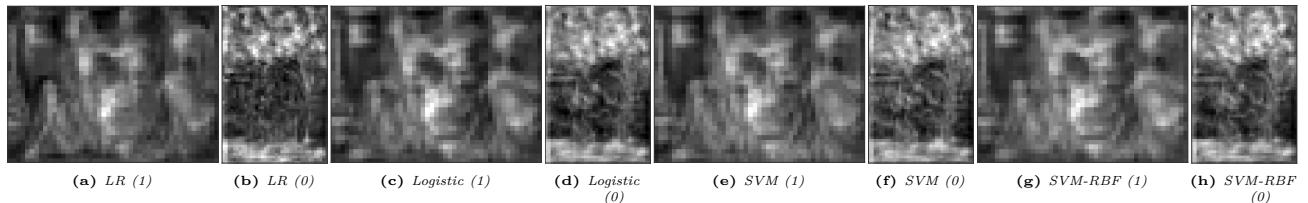
### 1.3 CNN2 features - Dense scores



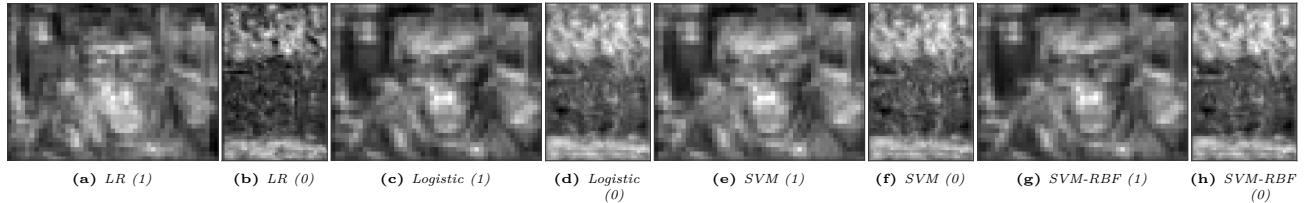
**Figure 11:** Right eye



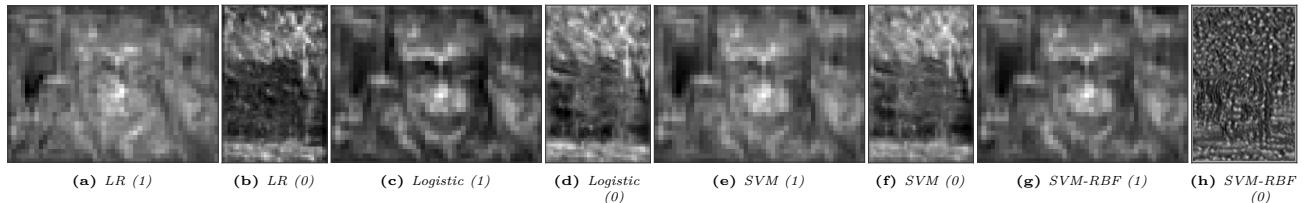
**Figure 12:** Left eye



**Figure 13:** Left mouth

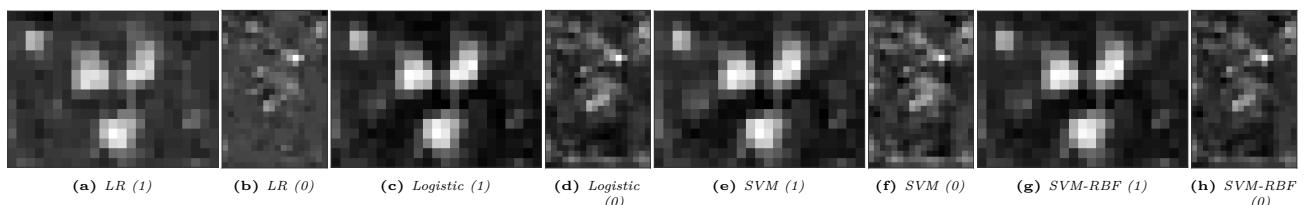


**Figure 14:** Right mouth

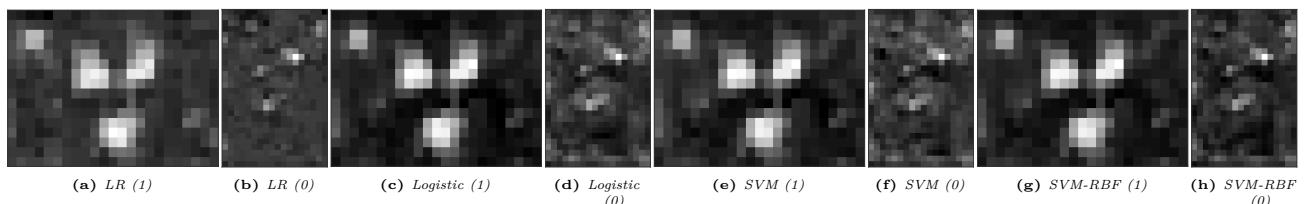


**Figure 15:** Nose

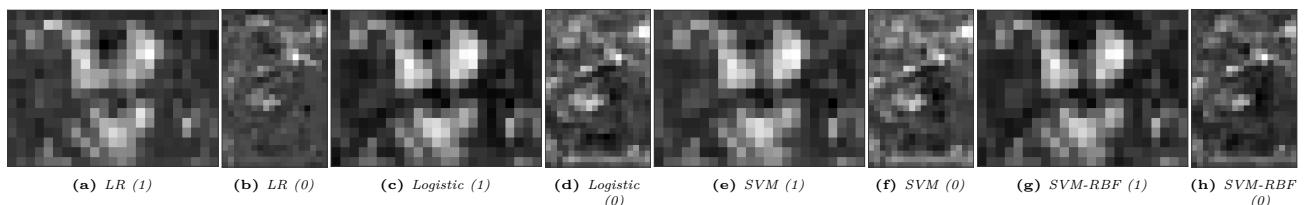
#### 1.4 CNN5 features - Dense scores



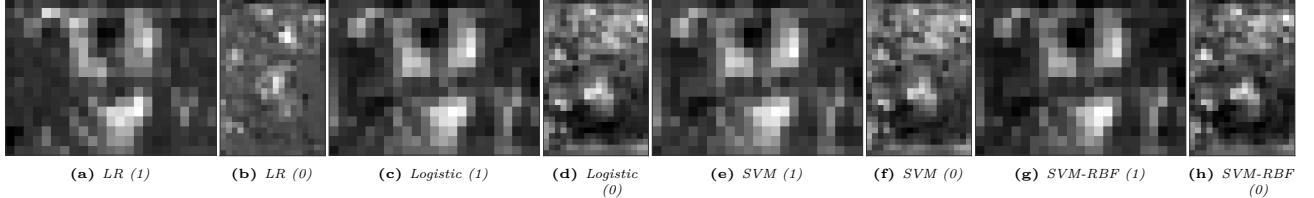
**Figure 16:** Right eye



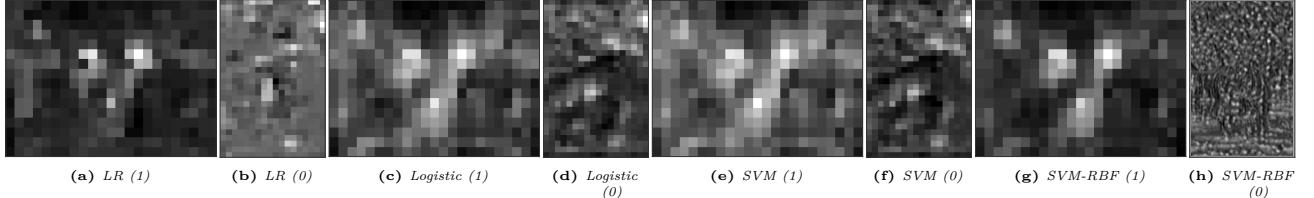
**Figure 17:** Left eye



**Figure 18:** Left mouth

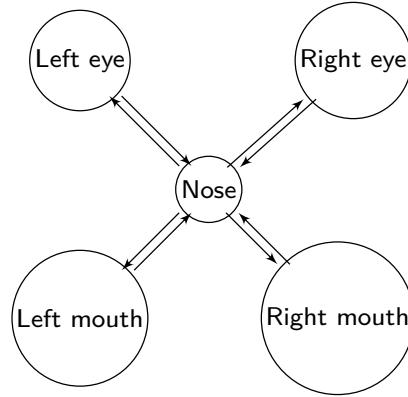


**Figure 19:** Right mouth



**Figure 20:** Nose

## 2 Combination: Max-product for DPM detection



We consider the probabilistic model above to which we apply the max-product algorithm. Each leaf node {Left eye, right eye, left mouth, right mouth} (index  $p$ ) passes a message to the root (index  $r$ ) of the form:

$$\begin{aligned}
 m_{p \rightarrow r}(X) &= \log \left( \max_{X_p} \Phi_p(X_p) \Psi_{p,r}(X_p, X_r) \right) \\
 &= \max_{X_p} \phi_p(X_p) + a(X_p(1) - X_r(1))^2 + b(X_p(1) - X_r(1)) + c(X_p(2) - X_r(2))^2 + d(X_p(2) - X_r(2)) + e
 \end{aligned}$$

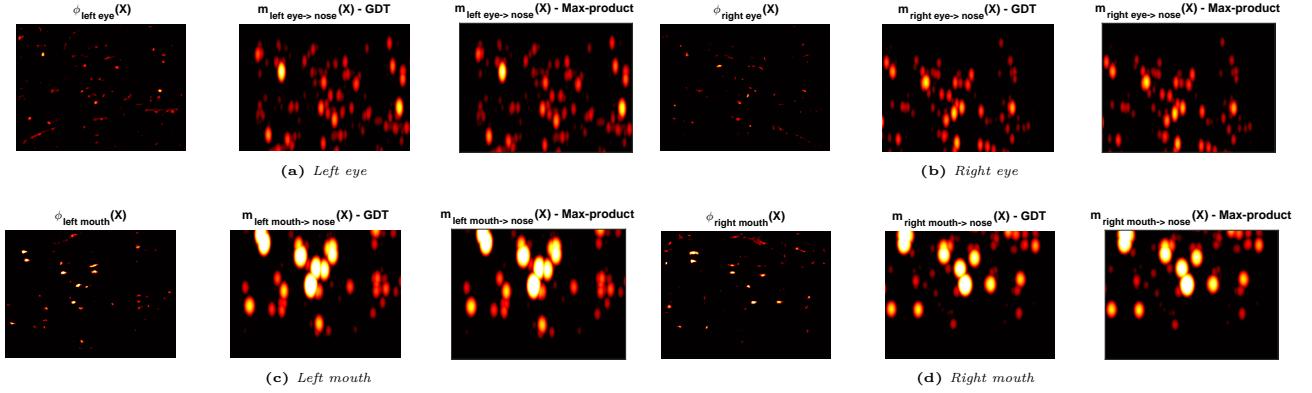
where  $\phi = \log \Phi, \psi = \log \Psi$  and  $a, b, c, d$  and  $e$  are parameters characterizing the distribution of the displacement vector  $X_P - X_r$  assumed Gaussian.

We have:

$$\begin{cases} a = -\frac{1}{2\sigma_p(1)^2} \\ b = \frac{\mu_p(1)}{\sigma_p(1)^2} \\ c = -\frac{1}{2\sigma_p(2)^2} \\ d = \frac{\mu_p(2)}{\sigma_p(2)^2} \\ e = -\frac{\mu_p(1)^2}{2\sigma_p(1)^2} - \frac{\mu_p(2)^2}{2\sigma_p(2)^2} - \log(2\pi\sigma_1\sigma_2) \end{cases}$$

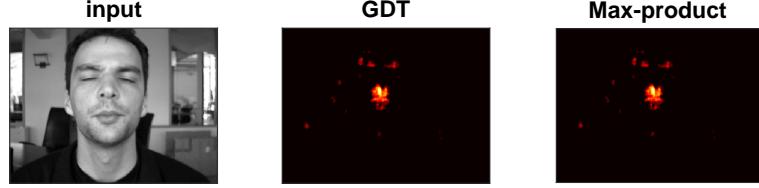
The unary potentials are computed first with the provided SVM weights using SIFT features and then with an SVM learnt using the CNN2 features with the Alexnet architecture.

## 2.1 Provided SVM model with SIFT



**Figure 21:** Message passing : GDT vs Max-product (image 3)

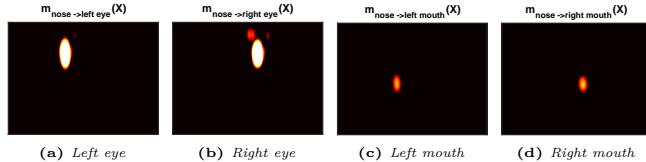
At the root we collect the messages to compute the belief function  $b(X) = \phi_r(X) + \sum_{i \neq r} m_{i \rightarrow r}(X)$



**Figure 22:** Belief function at the root ( $\equiv$  nose) - (image 3)

In a similar fashion, we send messages back to the parts:

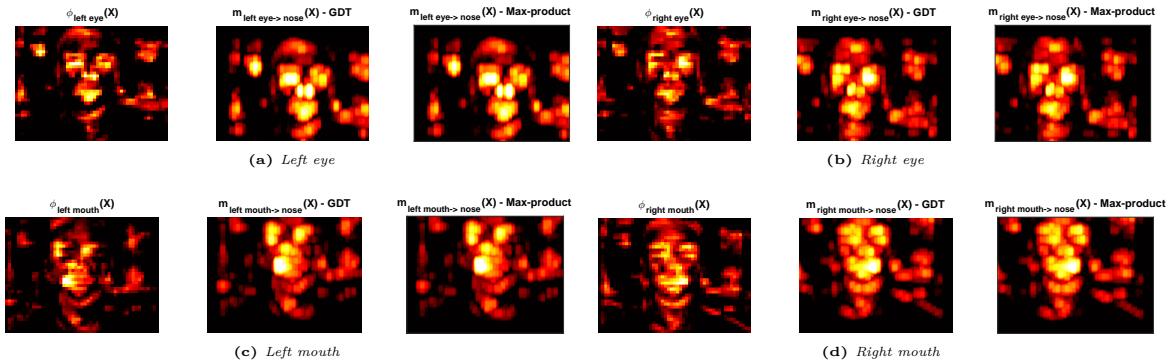
$$\begin{aligned} m_{r \rightarrow p}(X) &= \log \left( \max_{X_r} \left[ \Phi_r(X_r) \Psi_{r,p}(X_r, X_p) \prod_{k \in \mathcal{N}(r) \setminus p} M_{k \rightarrow r}(X_r) \right] \right) \\ &= \max_{X_r} \left[ \phi_r(X_r) + \psi_{r,p}(X_r, X_p) + \sum_{k \in \mathcal{N}(r) \setminus p} m_{k \rightarrow r}(X_r) \right] \end{aligned}$$



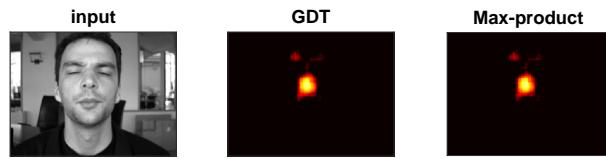
**Figure 23:** Messages to the parts (image 3)

## 2.2 Part 1' SVM model with CNN2

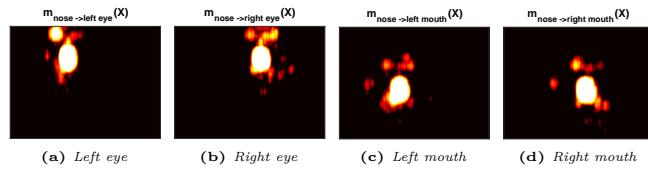
Image 3:



**Figure 24:** Message passing : GDT vs Max-product (image 3)

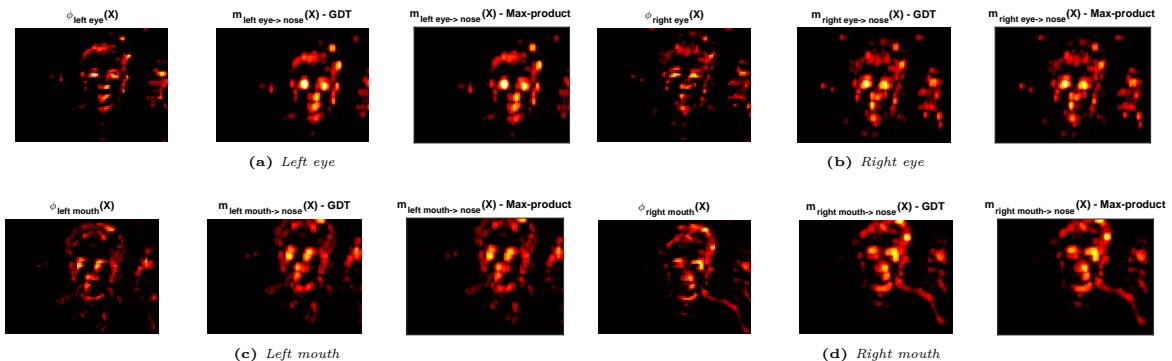


**Figure 25:** Belief function at the root ( $\equiv$  nose) - (image 3)

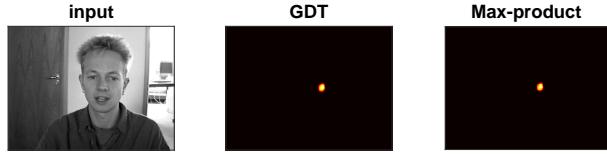


**Figure 26:** Messages to the parts (image 3)

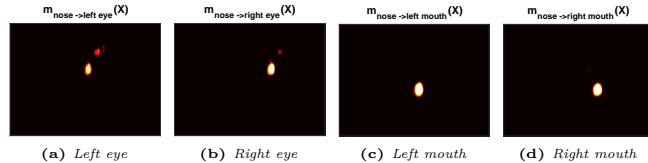
Image 231:



**Figure 27:** Message passing : GDT vs Max-product (image 231)

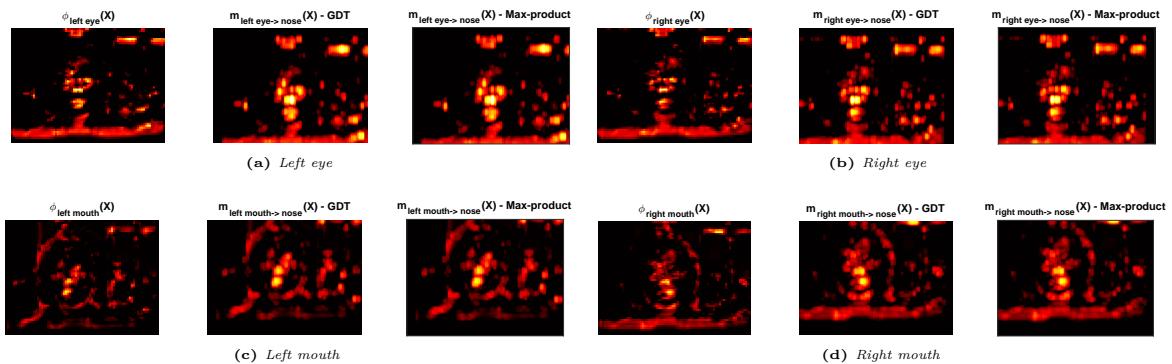


**Figure 28:** Belief function at the root ( $\equiv$  nose) - (image 231)

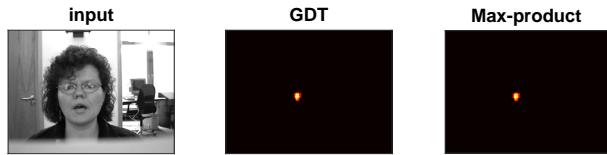


**Figure 29:** Messages to the parts (image 231)

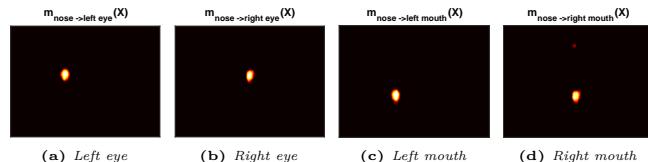
**Image 507:**



**Figure 30:** Message passing : GDT vs Max-product (image 507)



**Figure 31:** Belief function at the root ( $\equiv$  nose) - (image 507)

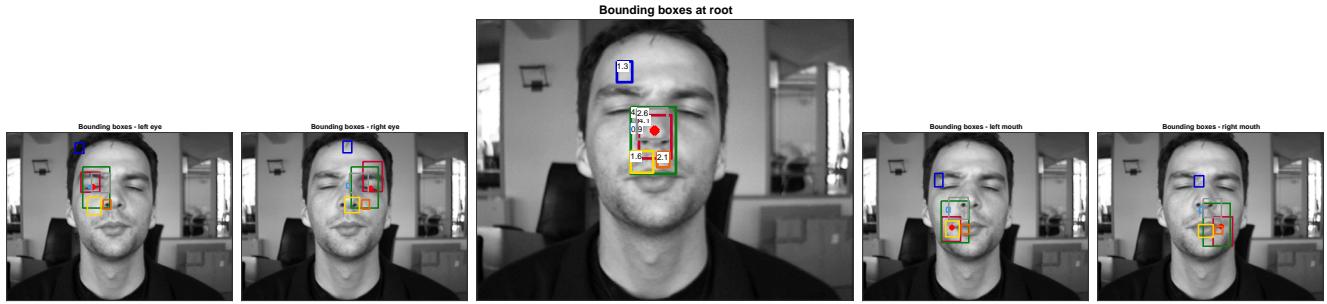


**Figure 32:** Messages to the parts (image 507)

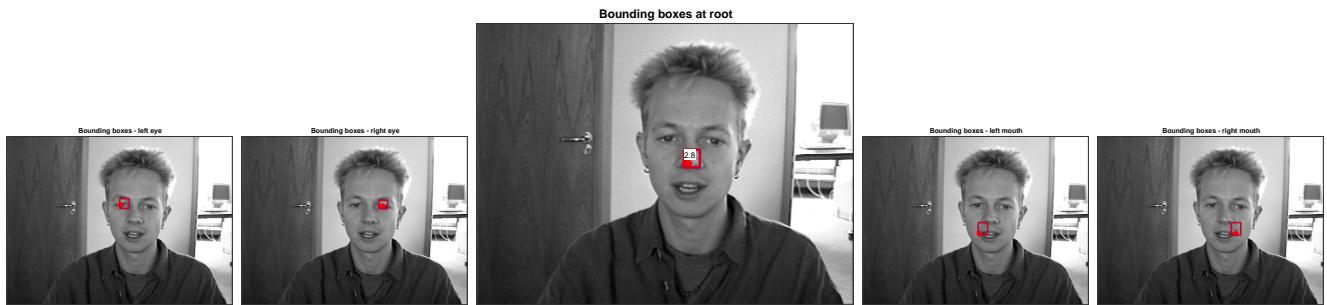
### 3 Multi-scale detection , non-maximum suppression, benchmarking

We perform feature extraction with the previously trained SVM on the CNN2 features then compute the unary scores and run GDT to compute the belief at the root at different scales  $[2^{-1}, 2^{-0.5}, 1, 2^{0.5}, 2]$ . The input image remains unnormalized i.e. we

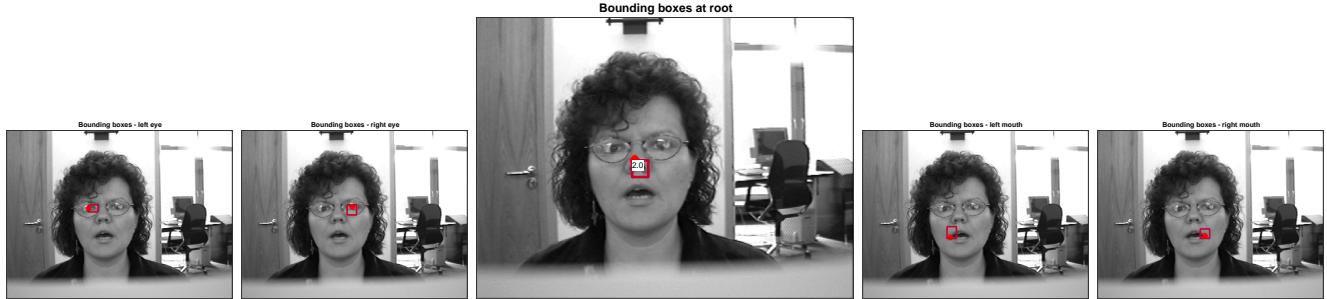
do not correct the distance between the eyes. The candidate roots are selected when the belief is above a selection threshold and the remaining parts are located using the *ix* and *iy* maps outputted by GDT given the root location.



**Figure 33:** Threshold = 0.5 - image 3

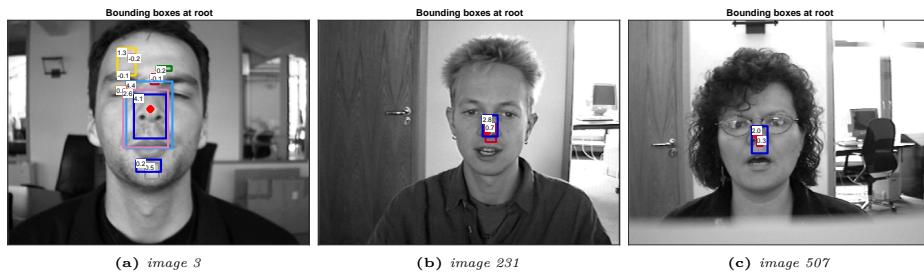


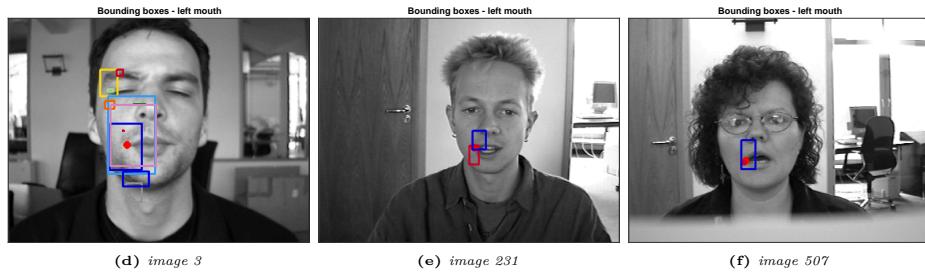
**Figure 34:** Threshold = 0.5 - image 231



**Figure 35:** Threshold = 0.5 - image 507

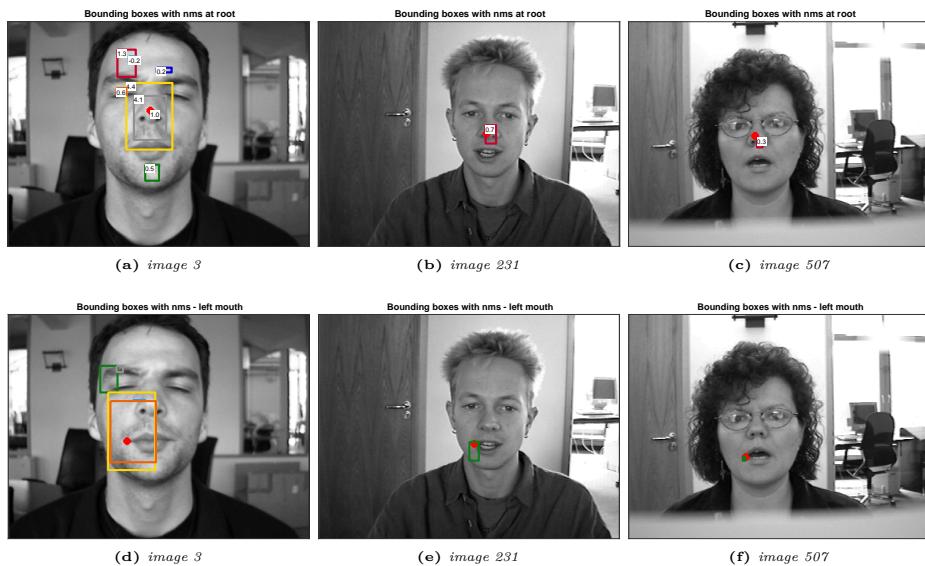
With a lower threshold we can suggest more locations:





**Figure 36:**  $\text{Threshold} = -0.5$

To filter the suggested boxes we apply non-maximum suppression with an overlap threshold of 0.7



**Figure 37:** NMS - box selection threshold -0.5 - box overlap threshold 0.7