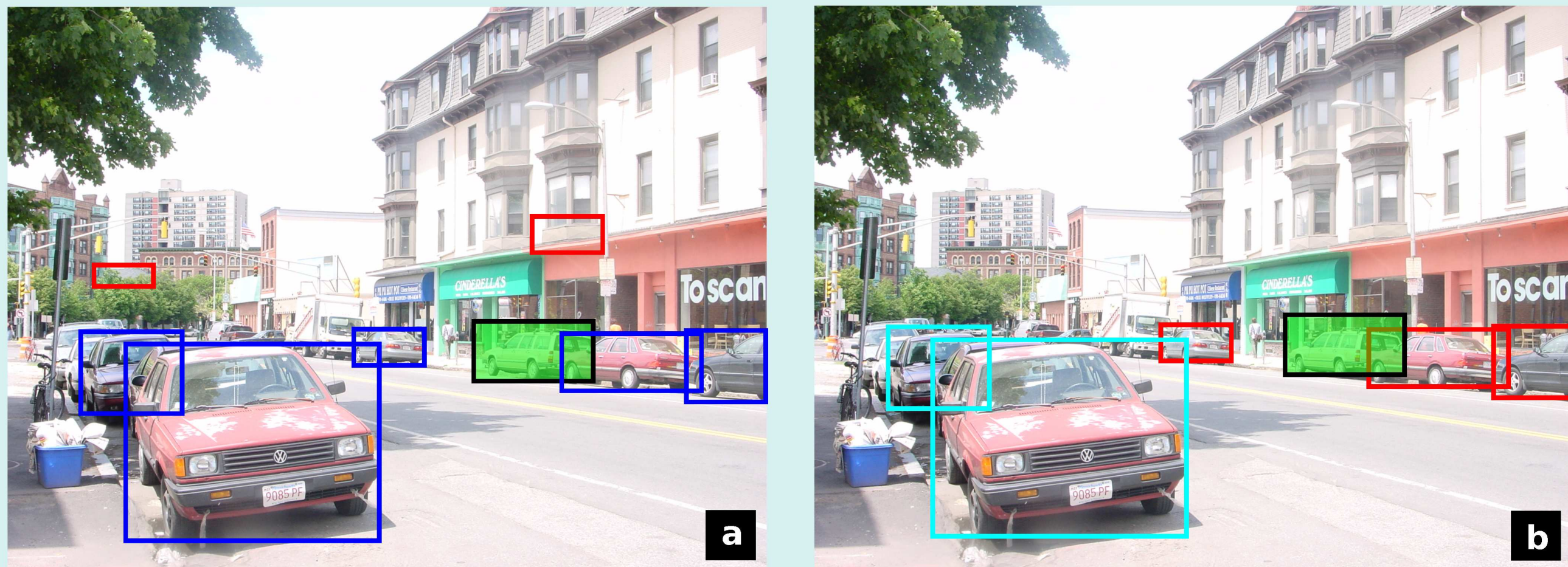


# Towards Cautious Collective Inference for Object Verification

José Oramas M., Luc De Raedt, Tinne Tuytelaars

## Abstract

- Use most certain object hypotheses to bootstrap the prediction of the other objects.
- Consider that objects associate with each other based on underlying "relationships" rather than the class they belong to.



## Motivation

- Contextual cues have shown to be of great benefit for object detection.
- Top-ranking object hypotheses from state-of-the-art detectors are usually correct --> good seed hypotheses.
- "Relationships" --> group behaviour in nature and imposed or desired rules on man-made objects.

## Challenges

- Defining relations between objects.
- How to select the contextual objects.
- How to define the influence of contextual objects.
- How to combine the local and contextual sources of information.

## Defining relations between objects

### Object:

#### Format 1

$x_i$ : X coordinate in the 2D image space.  
 $y_i$ : Y coordinate in the image space.  
 $\theta_i$ : Azimuth Pose angle  
 $s_i$ : Detection score

#### Format 2 (based on Li et al's, CVPR 2012)

$x_i$ : X coordinate in the 2D image space.  
 $y_i$ : Y coordinate in the image space.  
 $\beta_i$ : Scale of the of the object bounding box  
 $a_i$ : Aspect Ratio  
 $s_i$ : Detection score

#### Format 3 (based on Perko et al's, CVIU 2010)

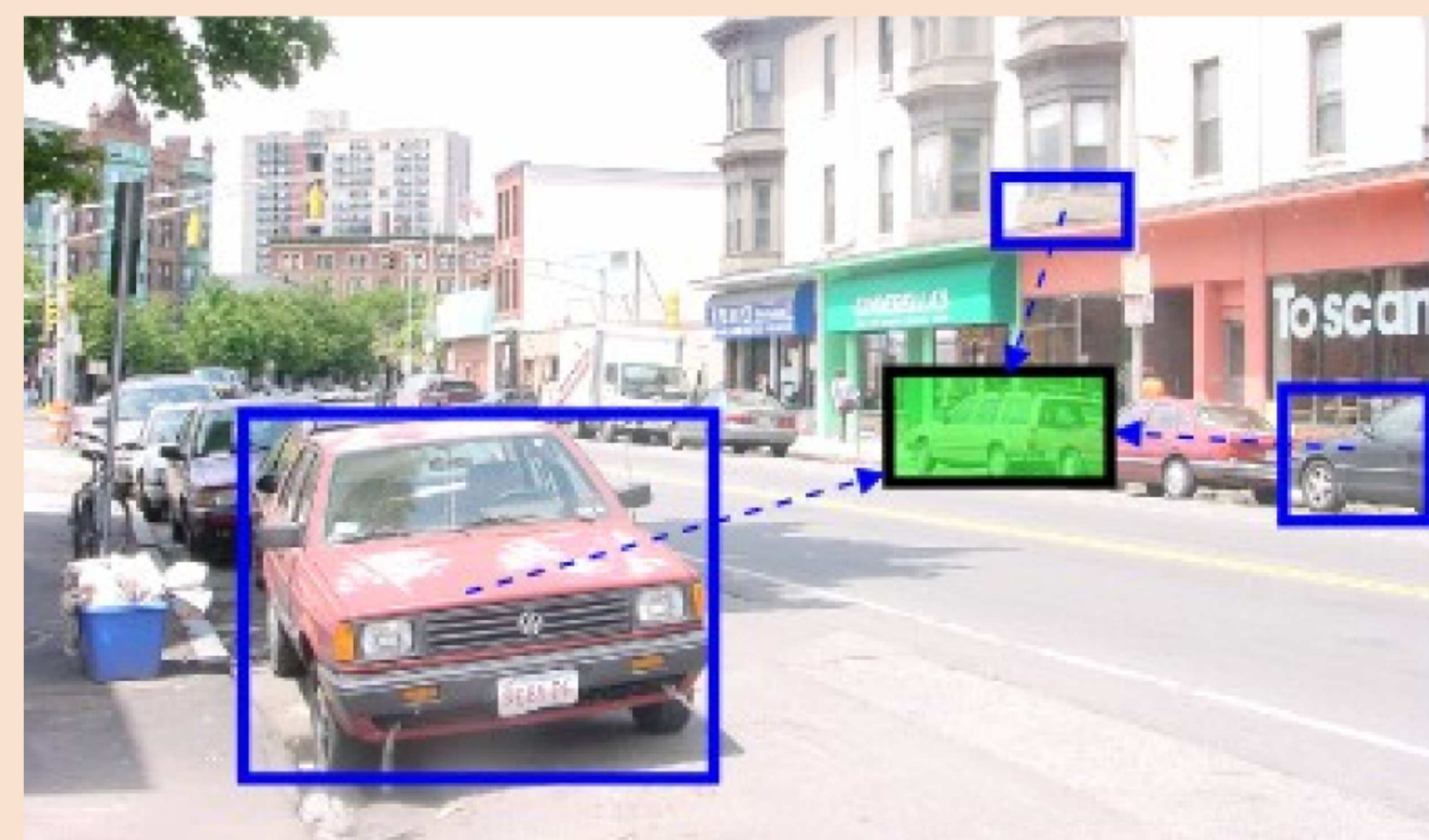
$x_i$ : X coordinate in the 2D image space.  
 $y_i$ : Y coordinate in the image space.  
 $s_i$ : Detection score

### Relations:

$$r_{ij}^{(RF1)} = (\Delta x_{ij}, \Delta y_{ij}, \Delta \theta_{ij})$$

$$r_{ij}^{(RF2)} = (rx_{ij}, ry_{ij}, r\rho_{ij}, rai_{ij})$$

$$r_{ij}^{(RF3)} = (\Delta x_{ij}, \Delta y_{ij})$$



## How to Select the Contextual Objects

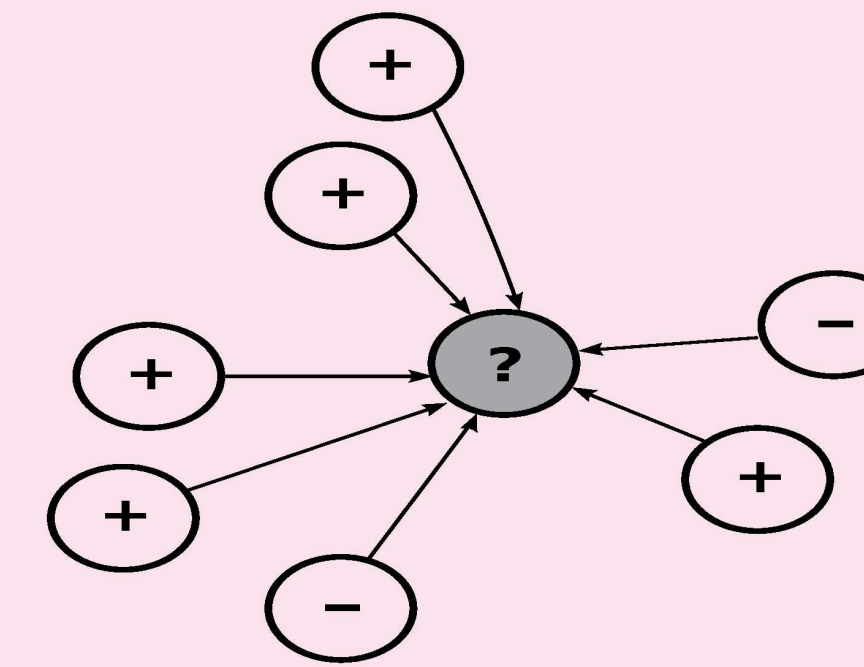
### Aggressive Inference

Weighted-Vote Relational Neighbor Classifier [1]

$$wvRN(o_i|N_i) = \frac{1}{z} \sum_{o_j \in N_i} p(o_i|r_{ij}) \cdot w_j$$

$$p(o_i|r_{ij}) = \frac{p(r_{ij}|o_i)p(o_i)}{p(r_{ij}|o_i)p(o_i) + p(r_{ij}|-o_i)p(-o_i)}$$

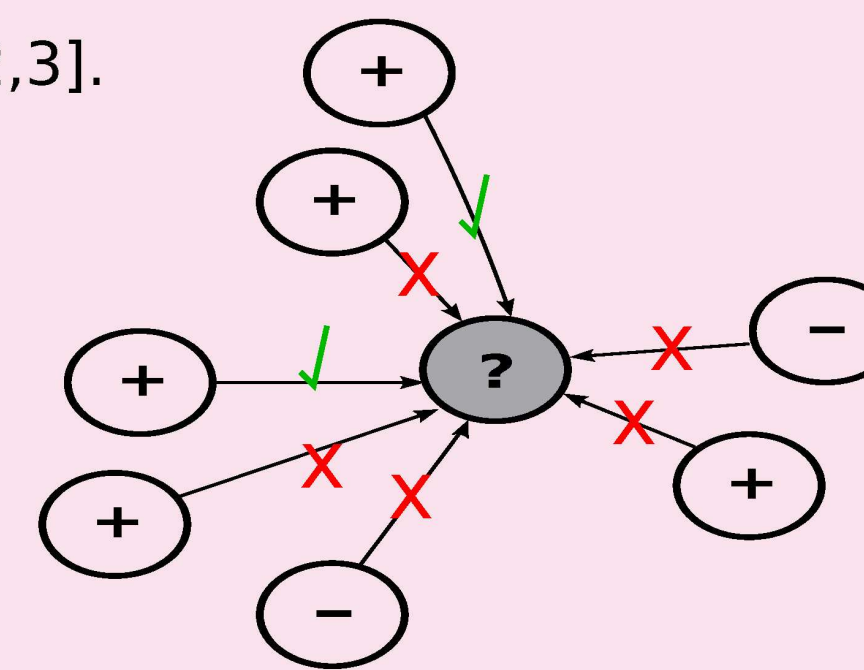
$$w_j = p(o_j|s_j) = \frac{p(s_j|o_j)p(o_j)}{p(s_j|o_j)p(o_j) + p(s_j|-o_j)p(-o_j)}$$



### Cautious Inference

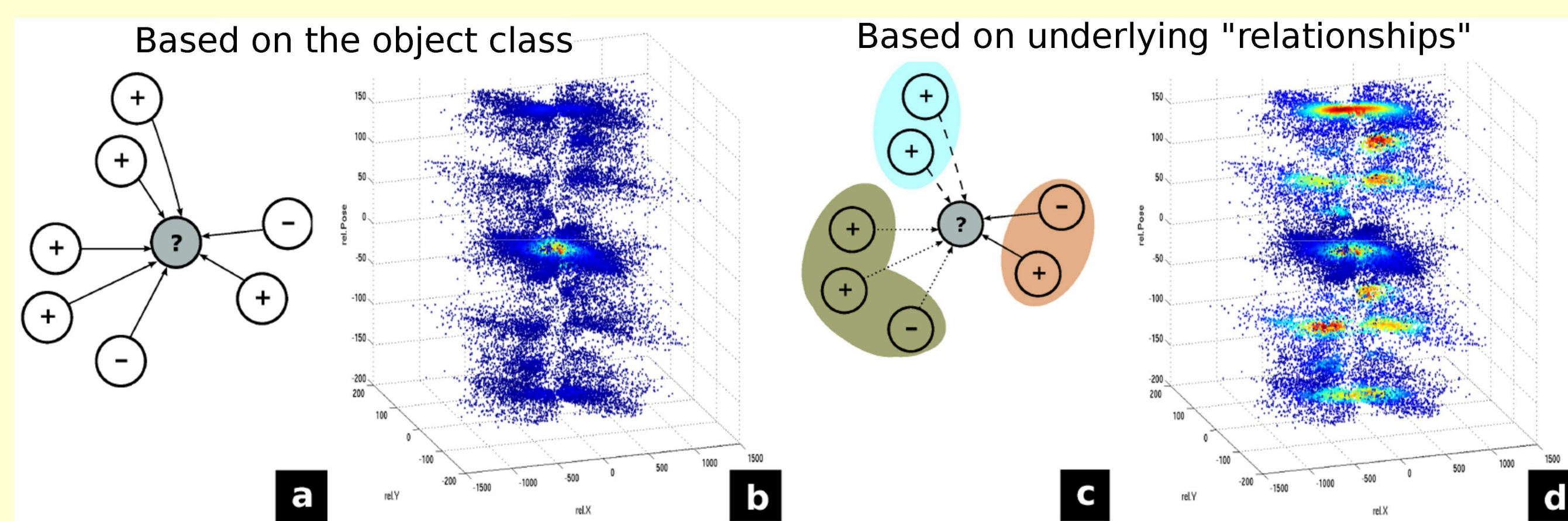
Inspired in the cautious and iterative methods from [2,3].

$$wvRN(o_i^u|N_i) = \frac{1}{z} \sum_{o_j^k \in (N_i \cap O^k)} p(o_i^u|r_{ij}) \cdot w_j$$



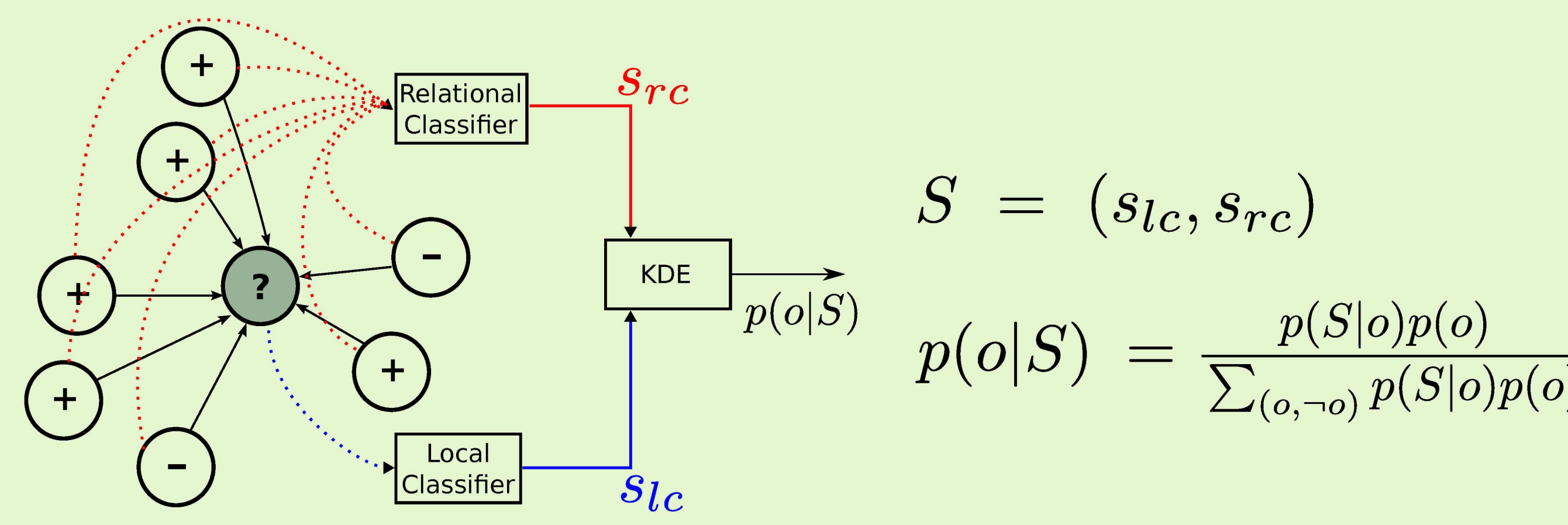
## How to define the influence of contextual objects

Relationships are extracted via XMeans clustering [4].



## How to combine local and contextual sources of information

Inspired in the score combination method from [2].



## Experiments

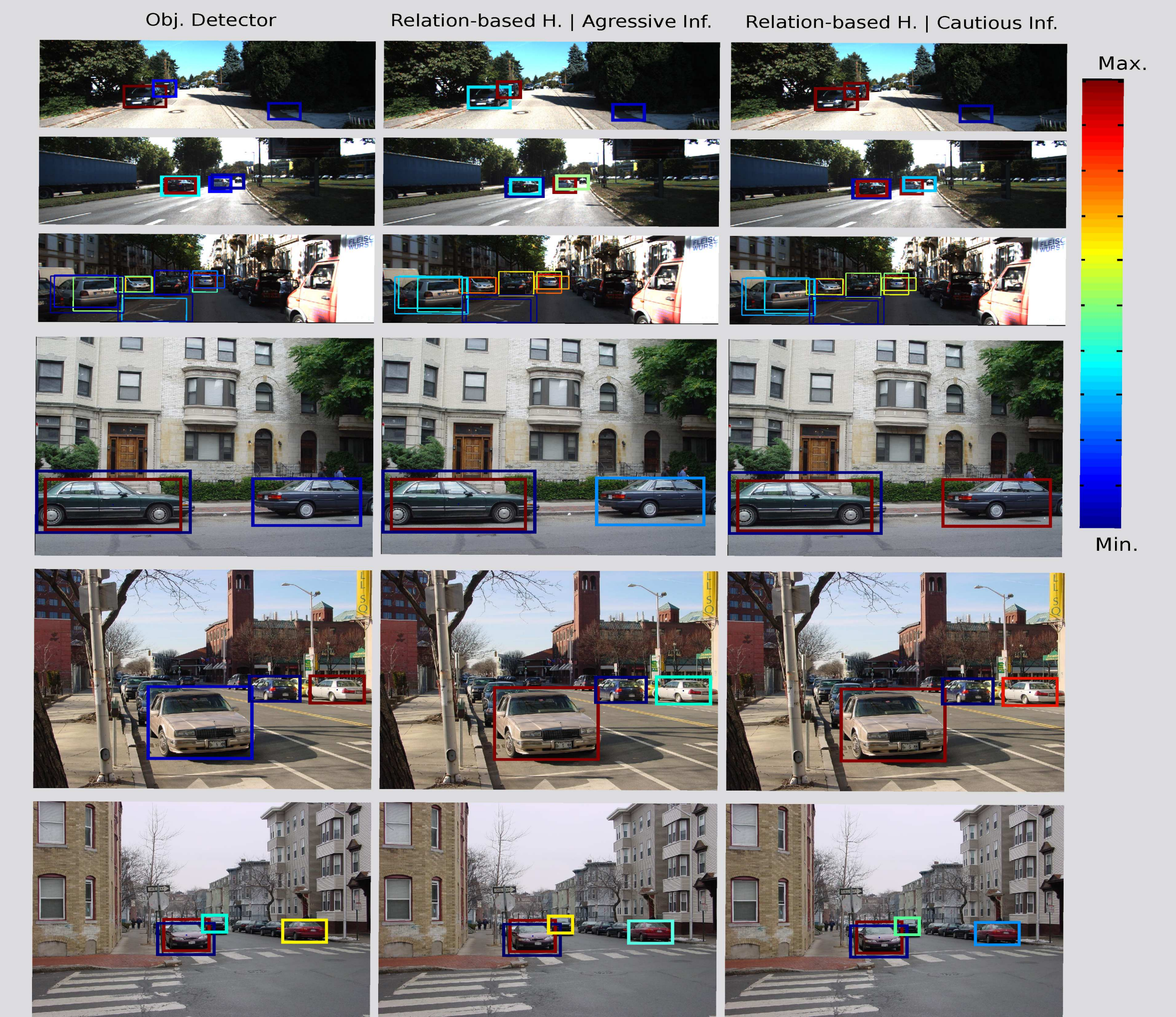
### Settings

#### Dataset :

- KITTI object detection benchmark.
- MIT StreetScenes.

#### Local Classifiers

- Viewpoint-aware DPM detector from [6] (8 views).
- DPM detector from [7] (No viewpoint information).



### Using only contextual information

Dataset	Relations Representation : RF1				Relations Representation : RF2			
	Class-based Hom.		Relation-based Hom.		Class-based Hom.		Relation-based Hom.	
	Global	Global	Global	Global	Global	Global	Global	Global
Set	aggre.	caut.	aggre.	caut.	aggre.	caut.	aggre.	caut.
KITTI benchmark	0.29	0.38	0.28	0.37	0.32	0.40	0.41	0.50
all								
Set	aggre.	caut.	aggre.	caut.	aggre.	caut.	aggre.	caut.
MIT StreetScenes	0.54	0.63	0.51	0.59	0.51	0.56	0.49	0.55
all								

### Local + contextual Information [6]

Dataset	RF1		RF2	
	Class-based Homophily	Relation-based Homophily	Class-based Homophily	Relation-based Homophily
	Global	Global	Global	Global
Set	aggre.	caut.	aggre.	caut.
KITTI benchmark	0.61±0.011	0.61±0.009	0.63±0.007	0.65±0.011
all				0.68±0.003
Set	aggre.	caut.	aggre.	caut.
MIT StreetScenes	0.77±0.001	0.80±0.028	0.73±0.011	0.76±0.014
all				

### Local + contextual Information [7]

Dataset	RF3		RF2	
	Class-based Homophily	Relation-based Homophily	Class-based Homophily	Relation-based Homophily
	Global	Global	Global	Global
Set	aggre.	caut.	aggre.	caut.
KITTI benchmark	0.68±0.003	0.71±0.007	0.72±0.009	0.75±0.003
all				
Set	aggre.	caut.	aggre.	caut.
MIT StreetScenes	0.66±0.011	0.71±0.012	0.65±0.026	0.69±0.014
all				

## Conclusions

- In a Global Neighborhood, Cautious Inference > Aggressive Inference.
- Considering underlying "relationships" is useful for cases where the local information about the unknown object is unavailable.
- Cautious Inference + "relationships" in RF2 format are better for constrained camera settings.
- In most cases, Global Neighborhood > Near Neighborhood.

[1] S. Mackassy and F. Provost, JMLR 2007.  
 [2] L. McDowell et al., JMLR 2009.  
 [3] J. Neville et al., SRL WS@AAAI 2000.  
 [4] A. Pelleg, ICML 2000.

[5] R. Perko and A. Leonardis, CVIU 2010.  
 [6] R. Lopez et al., WS@ICCV 2011.  
 [7] P. Felzenszwalb et al., TPAMI 2010.