

# A Framework for Multimodal Urban Scene Understanding

Philippe XU

Jean-Baptiste BORDES

Franck DAVOINE

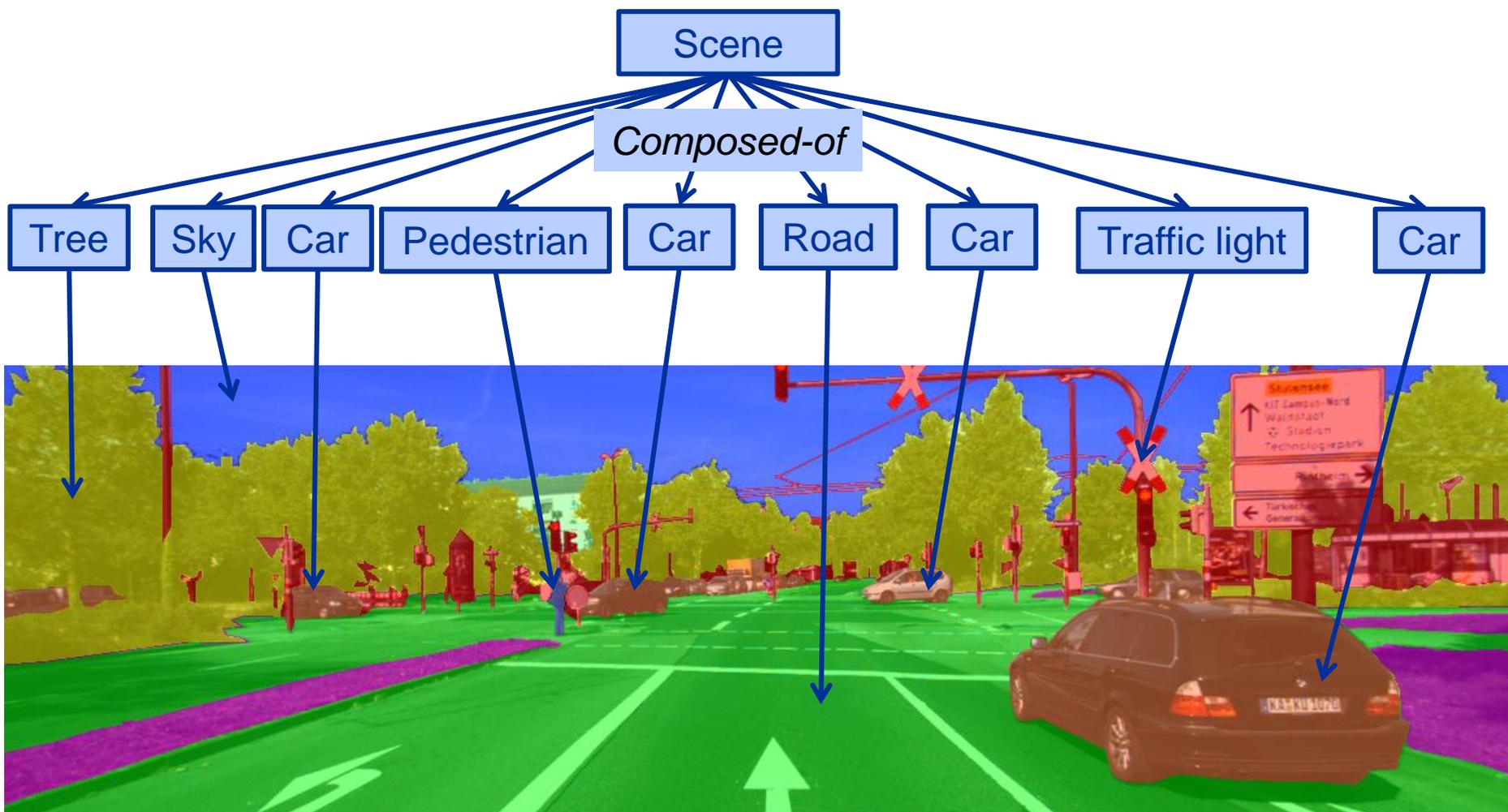
Thierry DENOEUUX

Huijing ZHAO

# Multi-sensors context



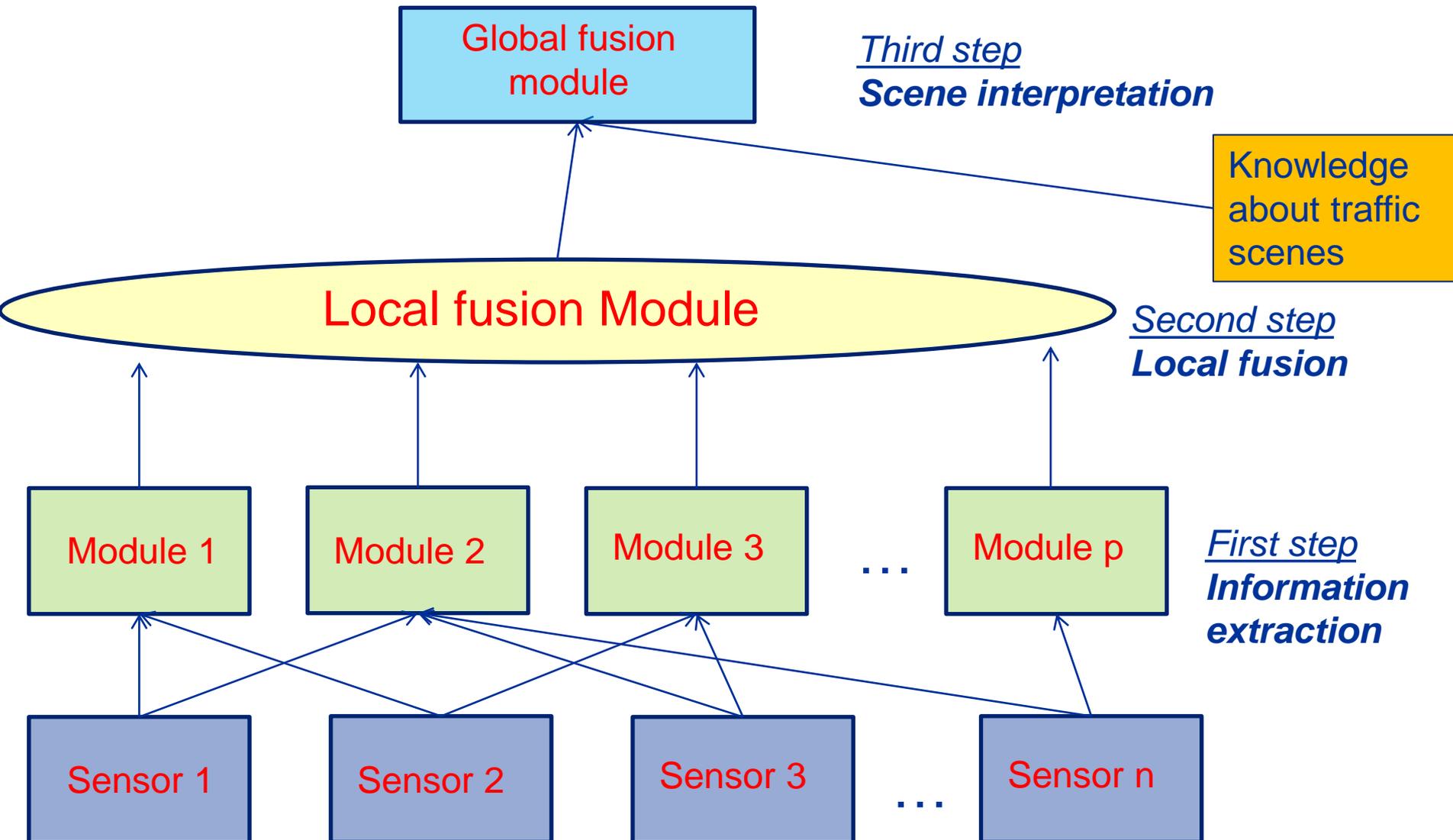
# Expected Result



# Requirements of the targeted framework

- **Flexibility:**
  - Cope with sensor failure
  - Add new classes easily
- **Multimodal:**
  - Fuse the output incoming from various kinds of sources of information
  - Take into account easily new sources of information

# Proposed global Framework



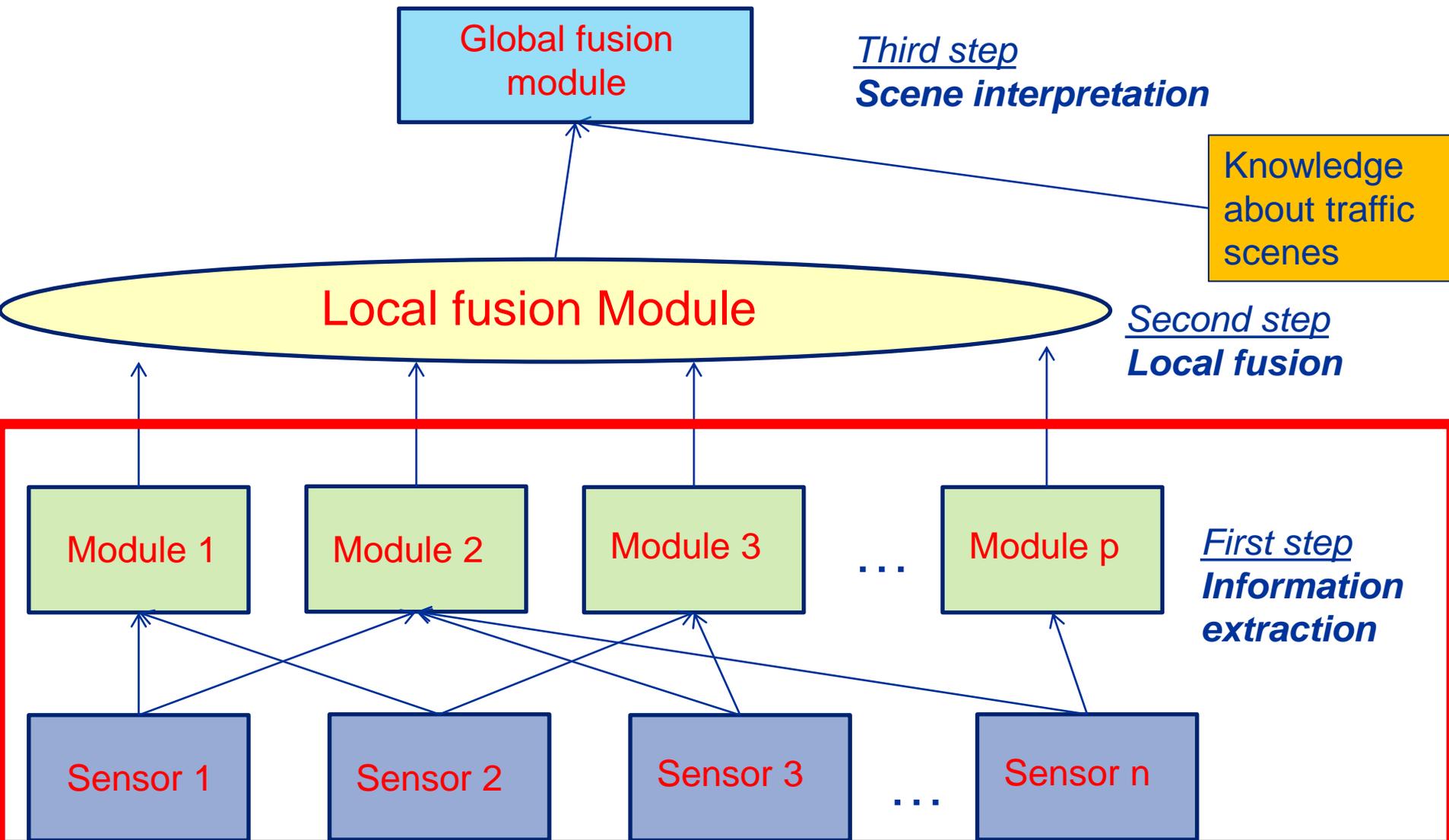
# Outlines

**1. Information extraction**

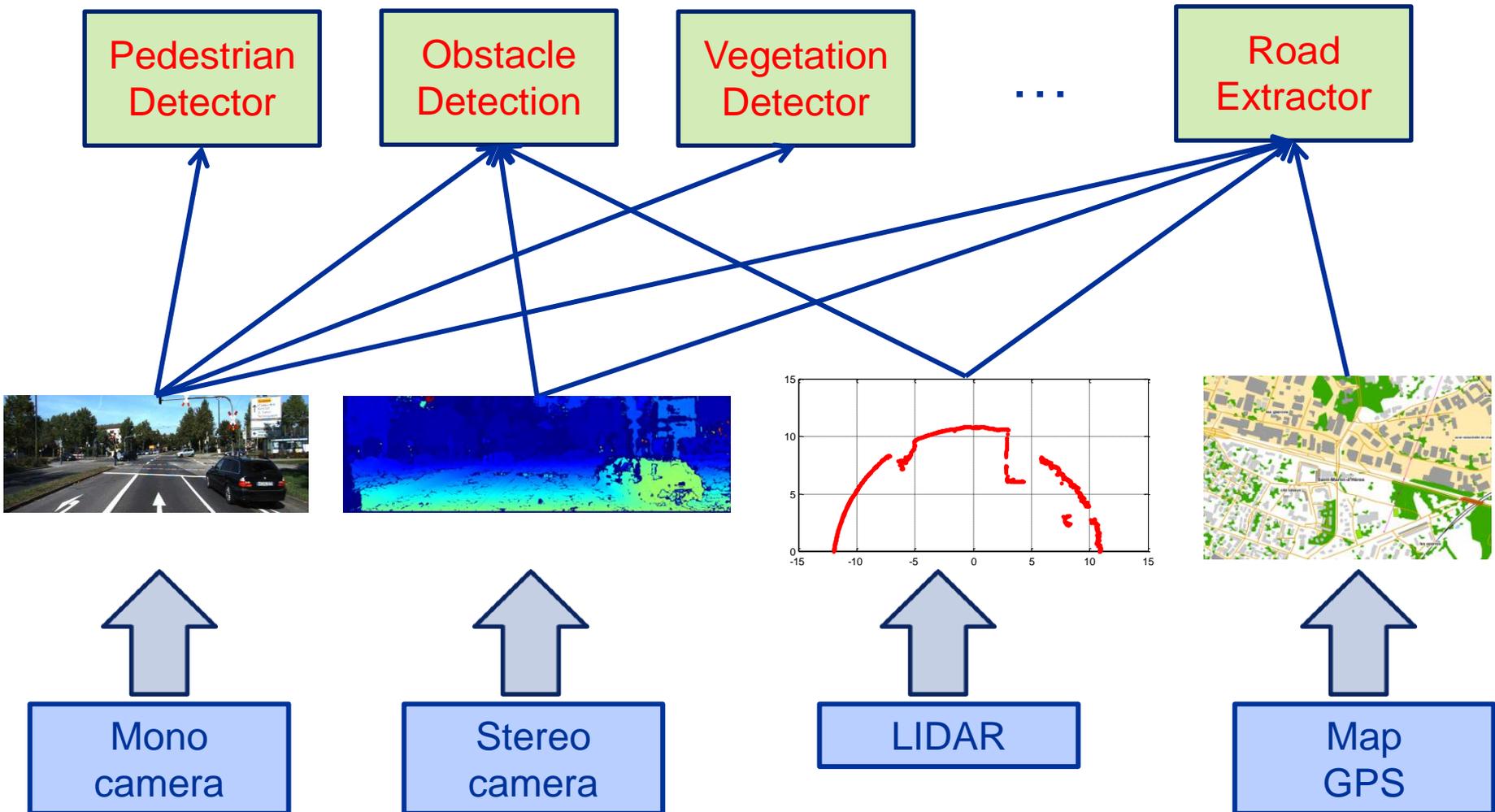
**2. Local Fusion using Dempster-Shafer theory**

**3. Global fusion using Evidential Grammar**

# Proposed global Framework

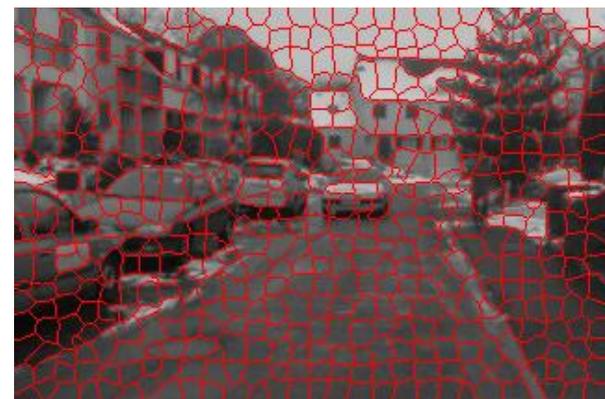


# Global Framework



# Decision space

- Reasoning at the image level:
  - It is what the driver sees
  - Adapted for driver assistance systems (may not be the case for autonomous driving)
  
- Classification over an over-segmented image
  - Intermediate level between pixel level/object level



# Outlines

**1. Information extraction**

**2. Local Fusion using Dempster-Shafer theory**

**3. Global fusion using Evidential Grammar**

# Local fusion step



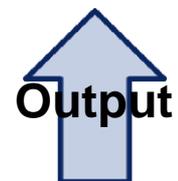
Local fusion Module



Output



Output



Output



Output

Pedestrian  
Detector

Obstacle  
Detection

Vegetation  
Detector

...

Road  
Extractor

# Class space

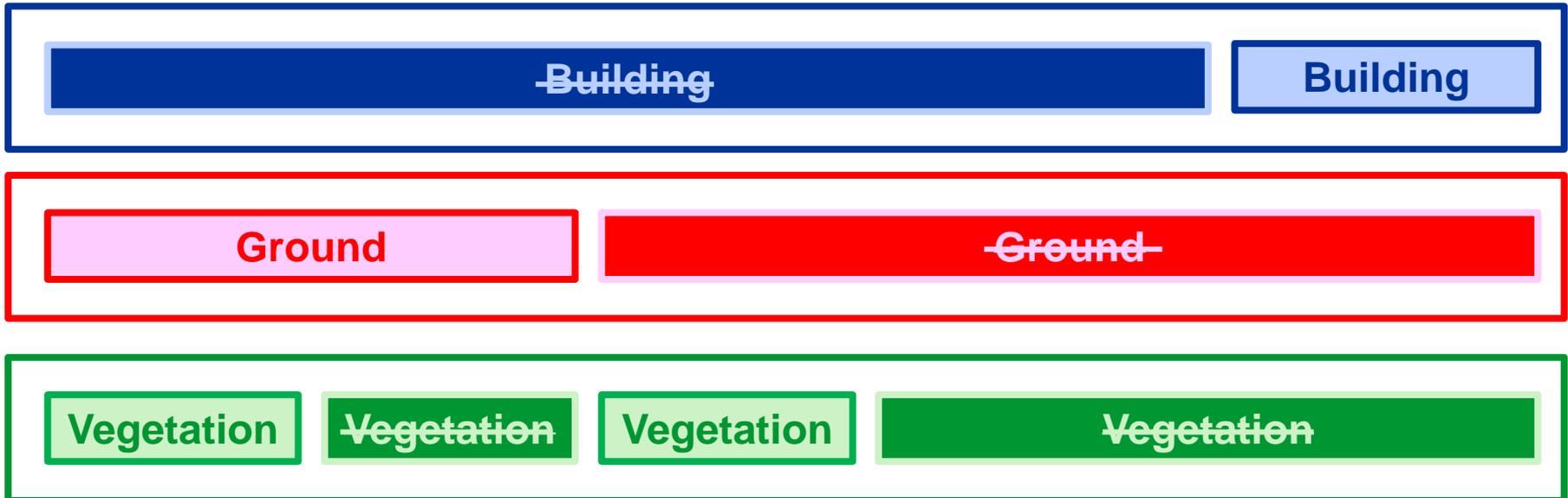
## Several classes could be considered:

- Pedestrians, cyclists
- Cars, motorbikes, trucks
- Roads, buildings, trees
- Traffic signs
- ...

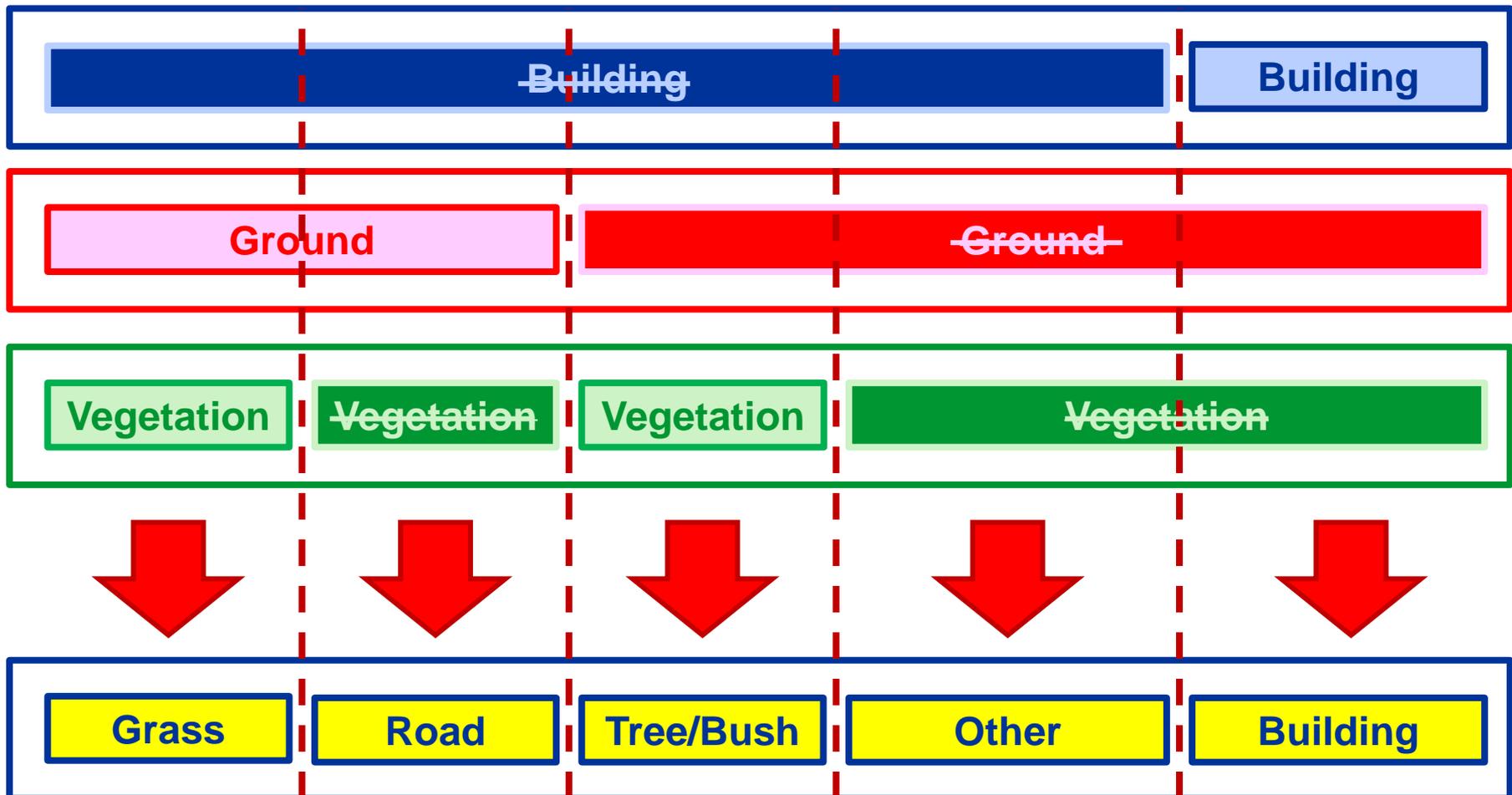
## Main issues:

- **Which** classes to choose?
- How can we make it **flexible** to add new classes?
- How to make **common decision space** to all classifiers?

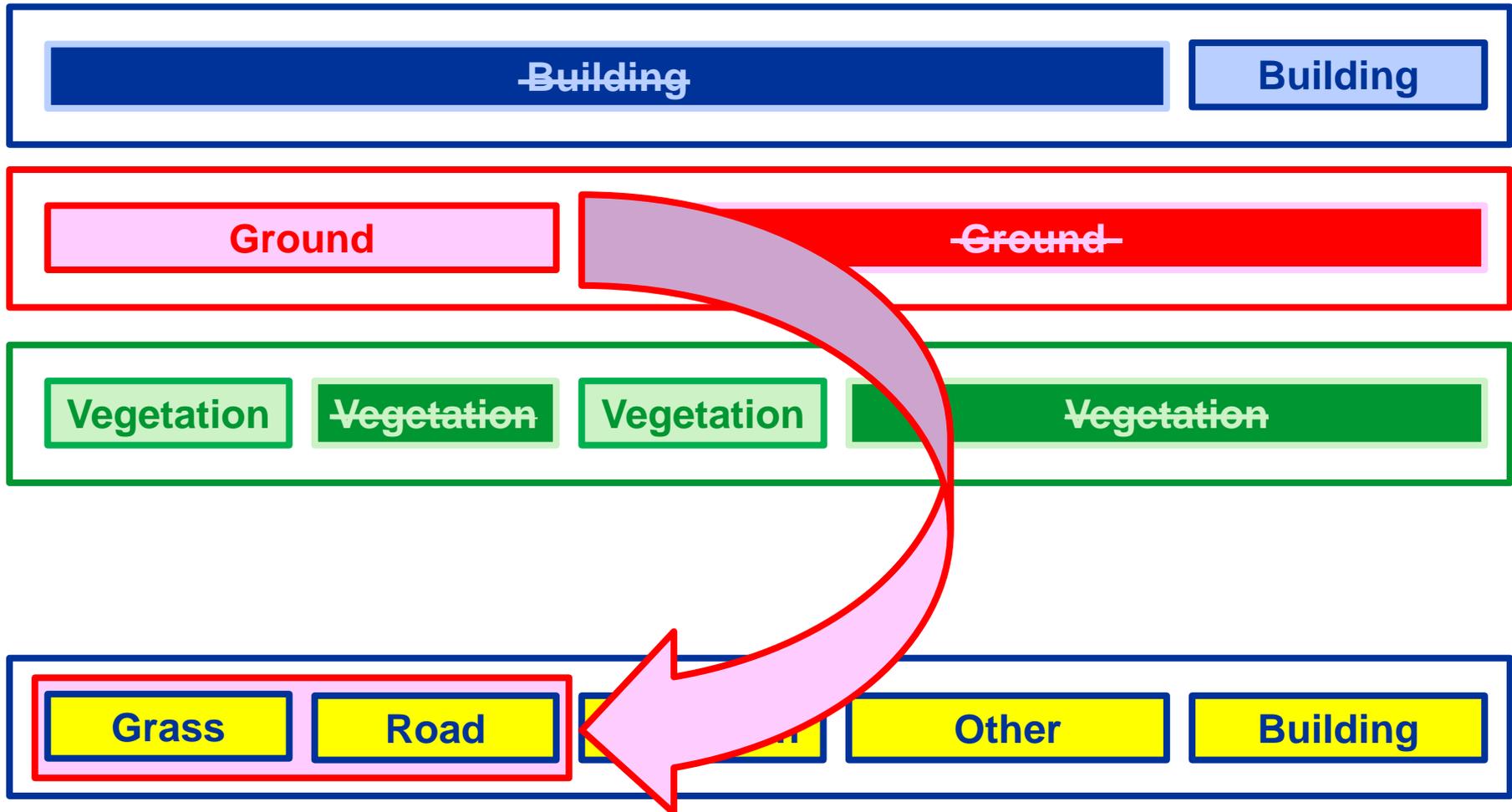
# Belief functions work on sets



# Belief functions work on sets



# Belief functions work on sets



# Ground detection

|   |  |
|---|--|
| <b>Disparity from stereo camera</b>         | <b>Stereo based ground detector</b>            |
| <b>Laser points (Velodyne)</b>              | <b>Laser based free space detector</b>         |
| <b>Optical flow from consecutive images</b> | <b>Optical flow based temporal propagation</b> |
| <b>Oversegmentation using Turbopixels</b>   | <b>Demspter-Shafer fusion</b>                  |

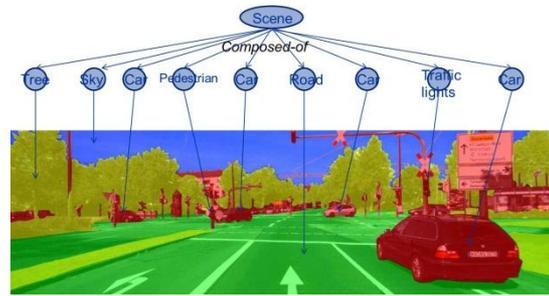
# Outlines

**1. Information extraction**

**2. Local Fusion using Dempster-Shafer theory**

**3. Global fusion using Evidential Grammar**

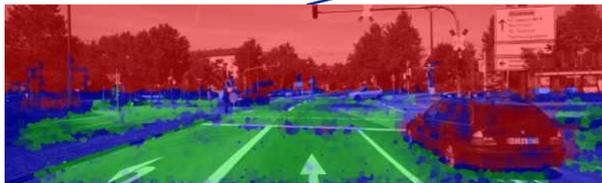
# Global fusion



Output

Global fusion module

*Third step*  
*Scene interpretation*



Output

Knowledge about traffic scenes

Local fusion Module

# Evidential grammars

- **Input : Classification at the segment level**
  - Objects or parts of objects
  - Contain uncertainty
- **Output: Image understanding**
  - Identification of the objects
  - Relationships between the objects
- **Visual Grammars**
  - Model of the decomposition of a scene into objects, parts of objects and primitives

# Stochastic grammars

- **A stochastic grammar is a 5 tuple  $(S, V_N, V_T, \Gamma, P)$  where:**
  - **$S$  is a starting symbol**
  - **$V_N$  is a set of non-terminal nodes**
  - **$V_T$  is a set of terminal nodes**
  - **$\Gamma$  is a set of production rules augmented with a set of probabilities  $P$ :**
    - **$A \rightarrow A_1$  with probability  $p_1$**
    - **$A \rightarrow A_2$  with probability  $p_2$**
    - **...**
    - **$A \rightarrow A_n$  with probability  $p_n$**

# Visual Grammars

- **Extension of the notion of stochastic grammars for the image**

- The natural left-to-right ordering of words is replaced by spatial relationships:

- Ex: Pedestrian -> Head “over” Body

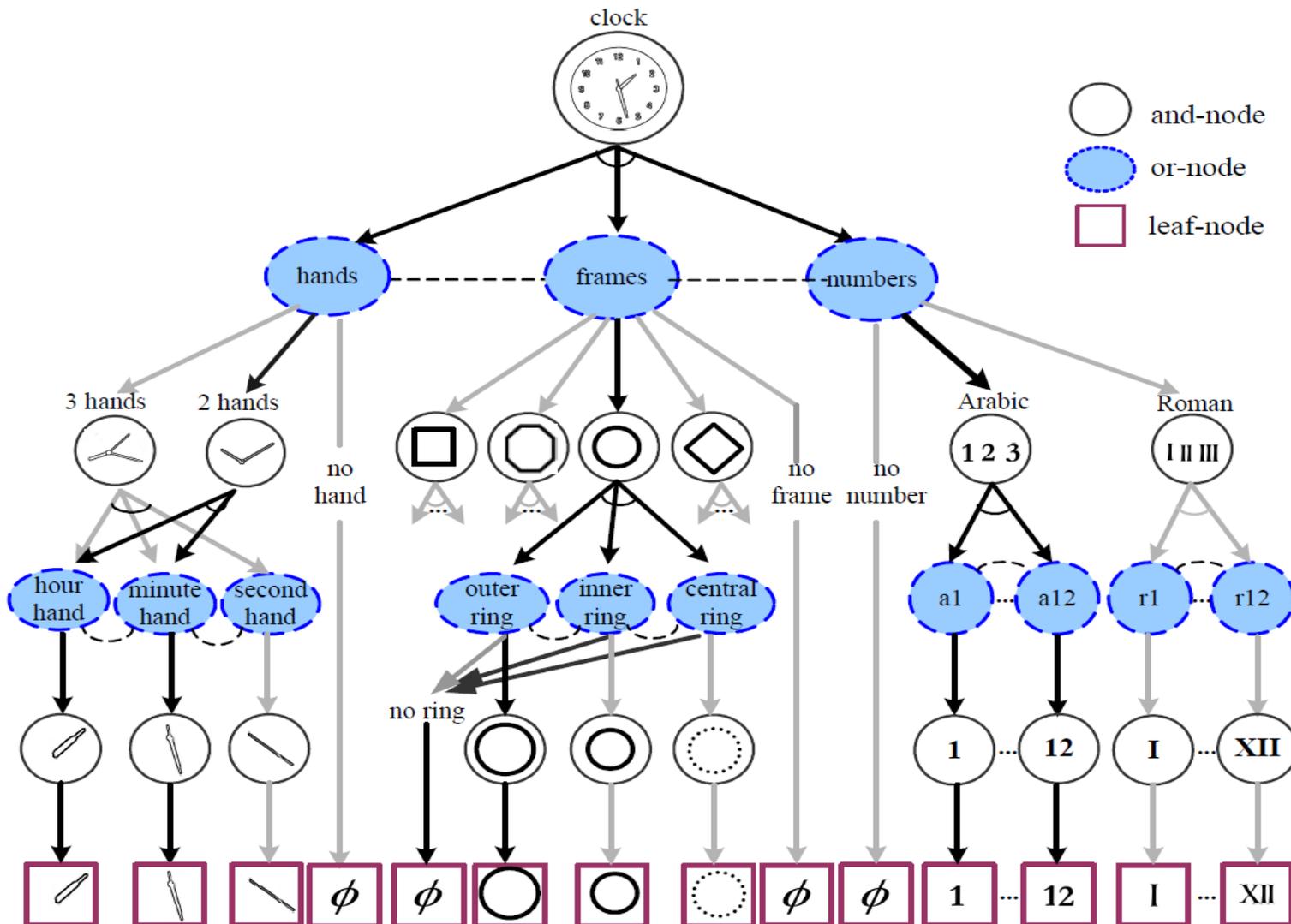
- The terminal nodes are called “visual primitives”



- The relationships may depend on the semantic level of the vocabulary

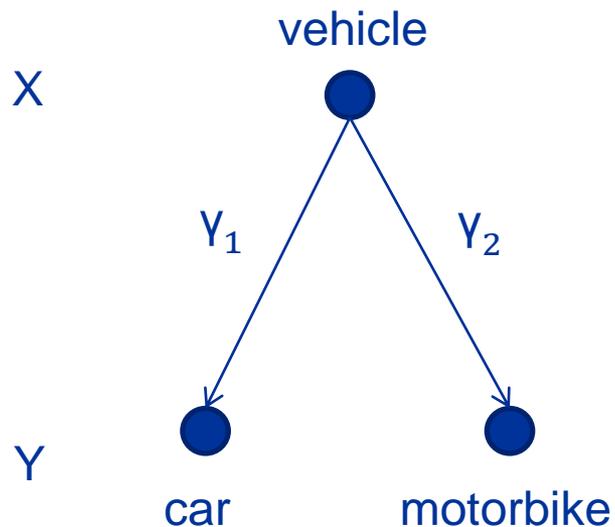
- Low-level: “adjacent”, “disjoint”
- Middle-level: “radial”, “bordering”, “hinge”
- High-level: “support”, “occlude”

# Visual grammars



# Evidential grammars

- An evidential grammar is a 5 tuple  $(S, V_N, V_T, \Gamma, \mathbf{M})$
- $\mathbf{M}$  contains a set of conditional mass functions defining the grammar rules



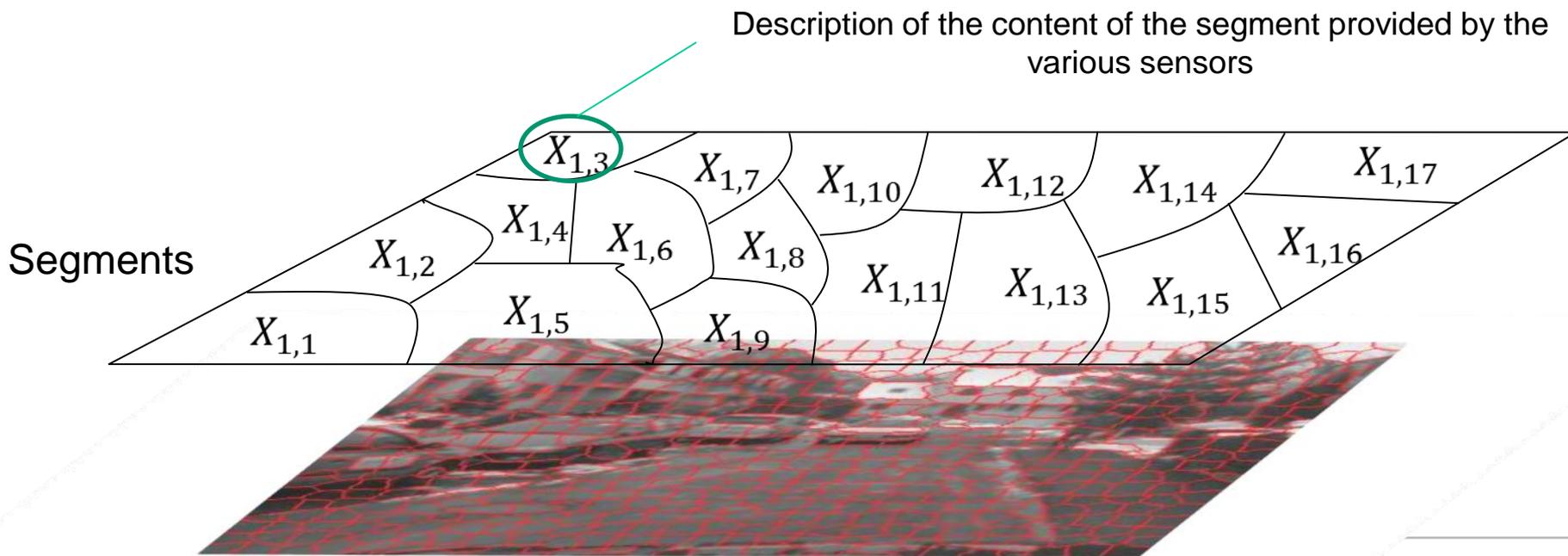
Assumption of complete ignorance:  
 $m(Y \in \{\text{car}, \text{motorbike}\} | X = \text{vehicle}) = 1$

# Why using evidential grammars?

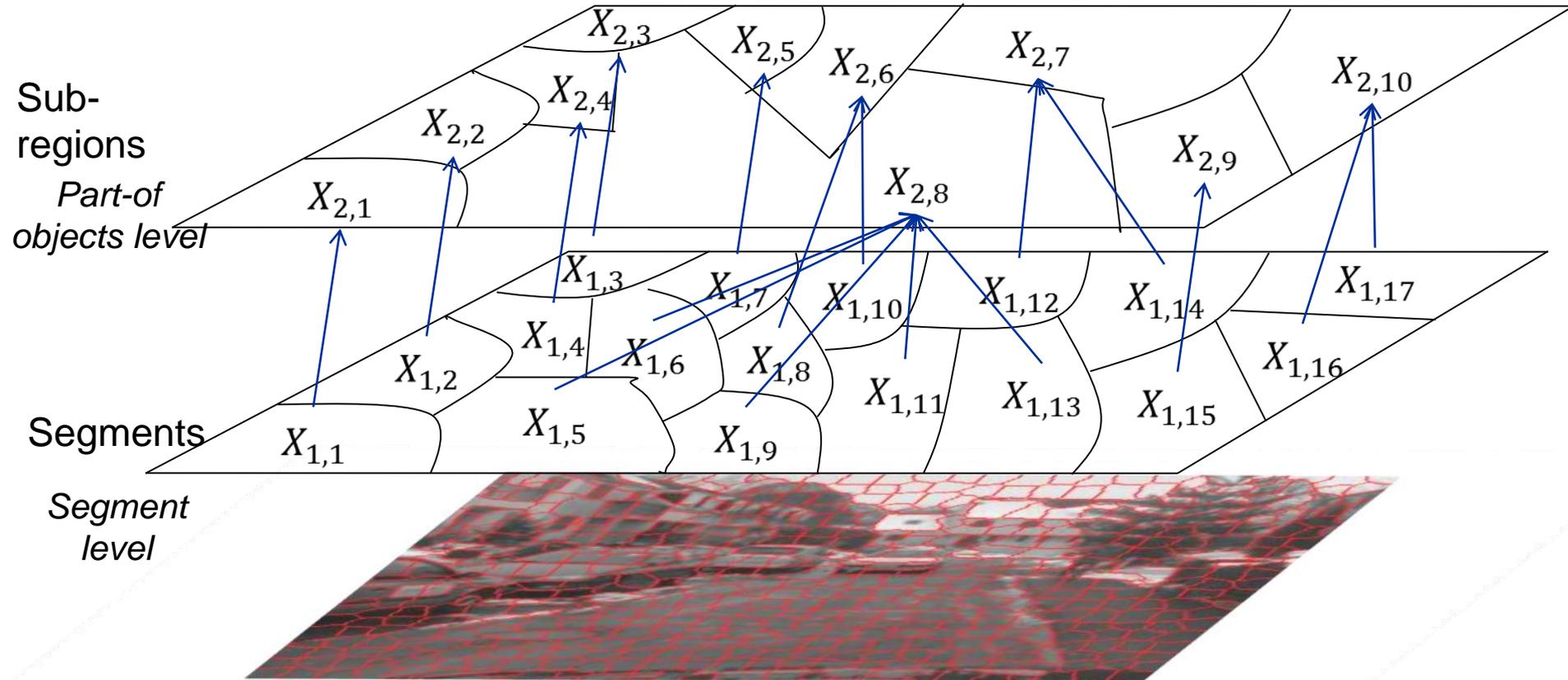
- **The knowledge about the situations which can occur in traffic scenes has to be used**
    - The lack of training data can be balanced by supplying knowledge to the system
  - **Learn part-of objects can be used to detect several objects**
    - A wheel can belong to a truck, a car, a bike, a moto ...
  - **The informations provided by the various sensors will provide belief on the different spaces**
    - It is of highest importance for us to handle the uncertainty in the interpretation process
- => A precise framework to handle visual grammars and belief functions has to be defined

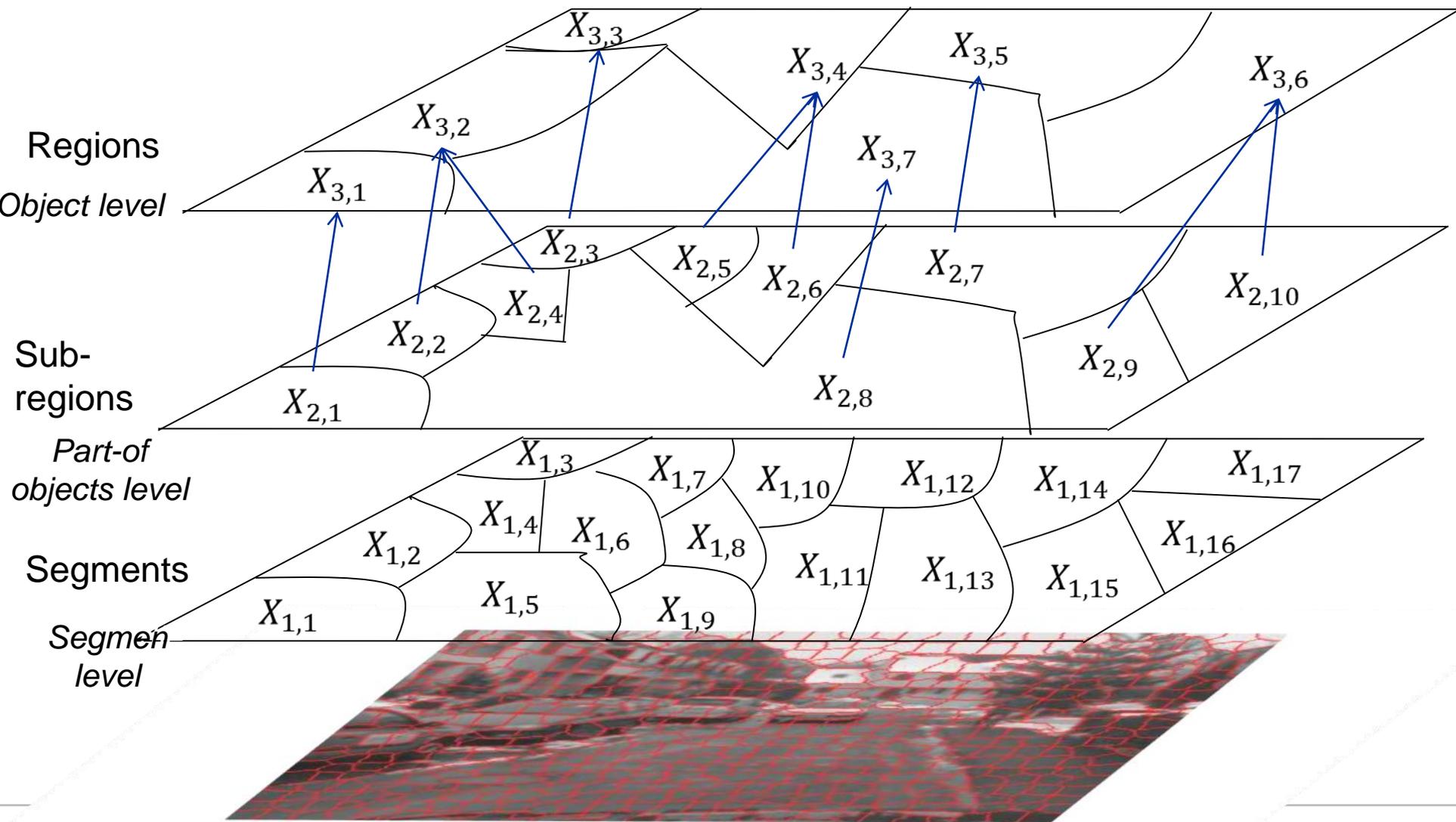
# Interpretation Tree

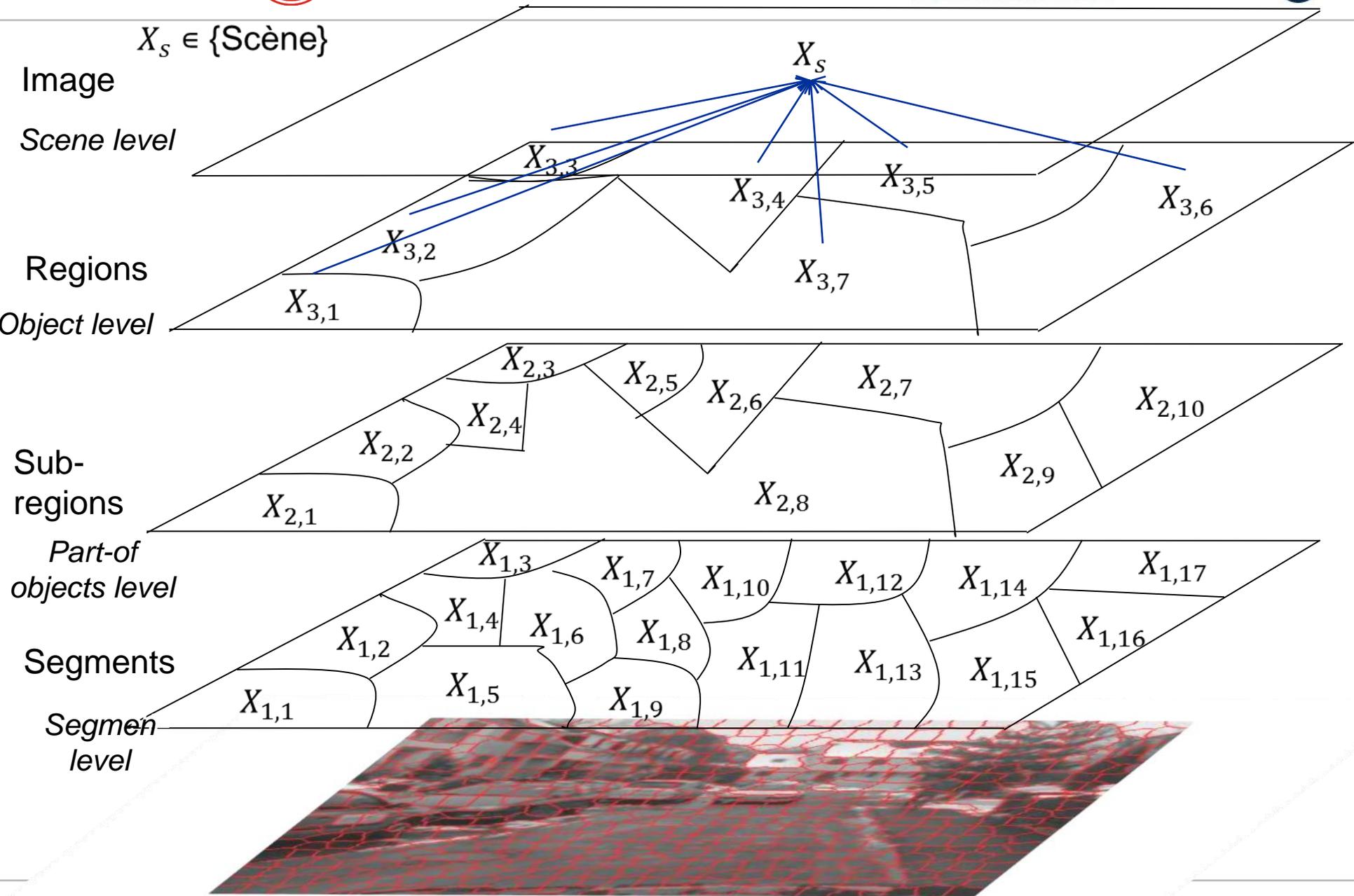
- The initial image is oversegmented and the information about the class contained in each one of these segments is described by a belief function and modelled by a random variable



# Interpretation Tree

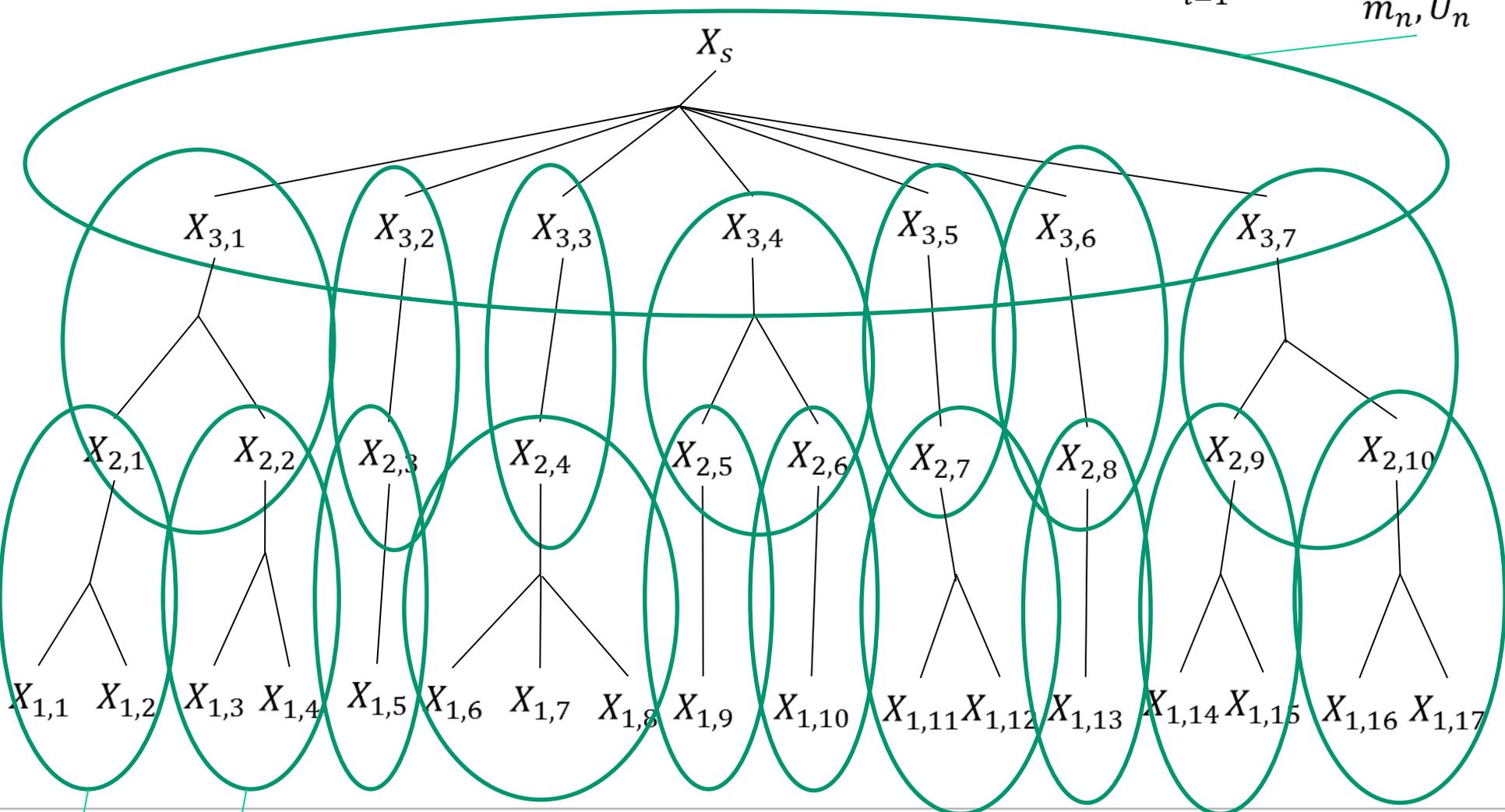






# Evidential Network

$$m_U = \bigcap_{i=1}^n m_{U_i \uparrow U} \quad m_n, U_n$$

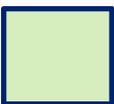


$m_1, U_1$      $m_2, U_2$

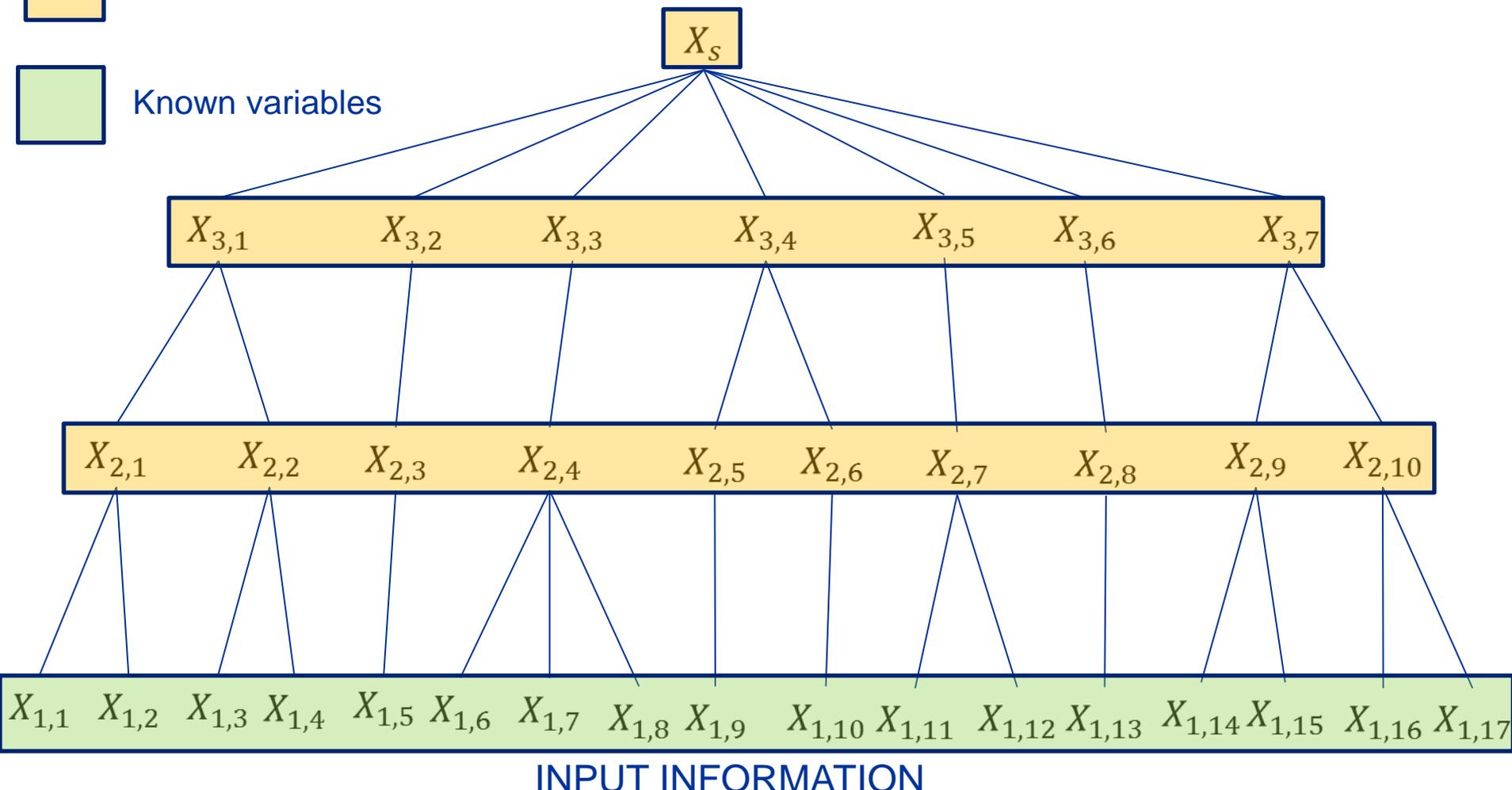
# Bottom-up Propagation of the belief



Unknown variables



Known variables

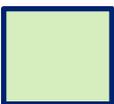


INPUT INFORMATION

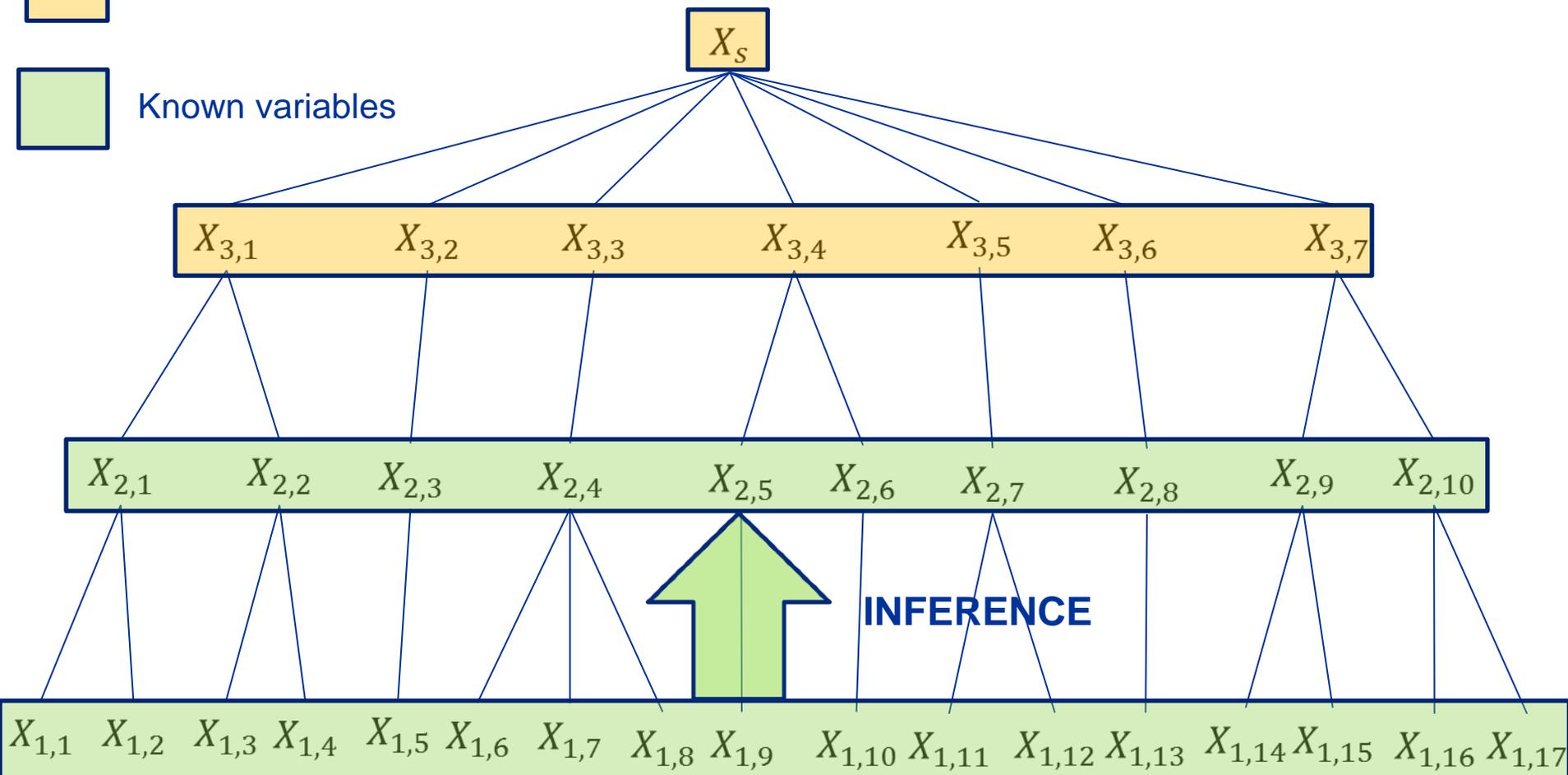
# Bottom-up Propagation of the belief



Unknown variables



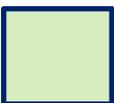
Known variables



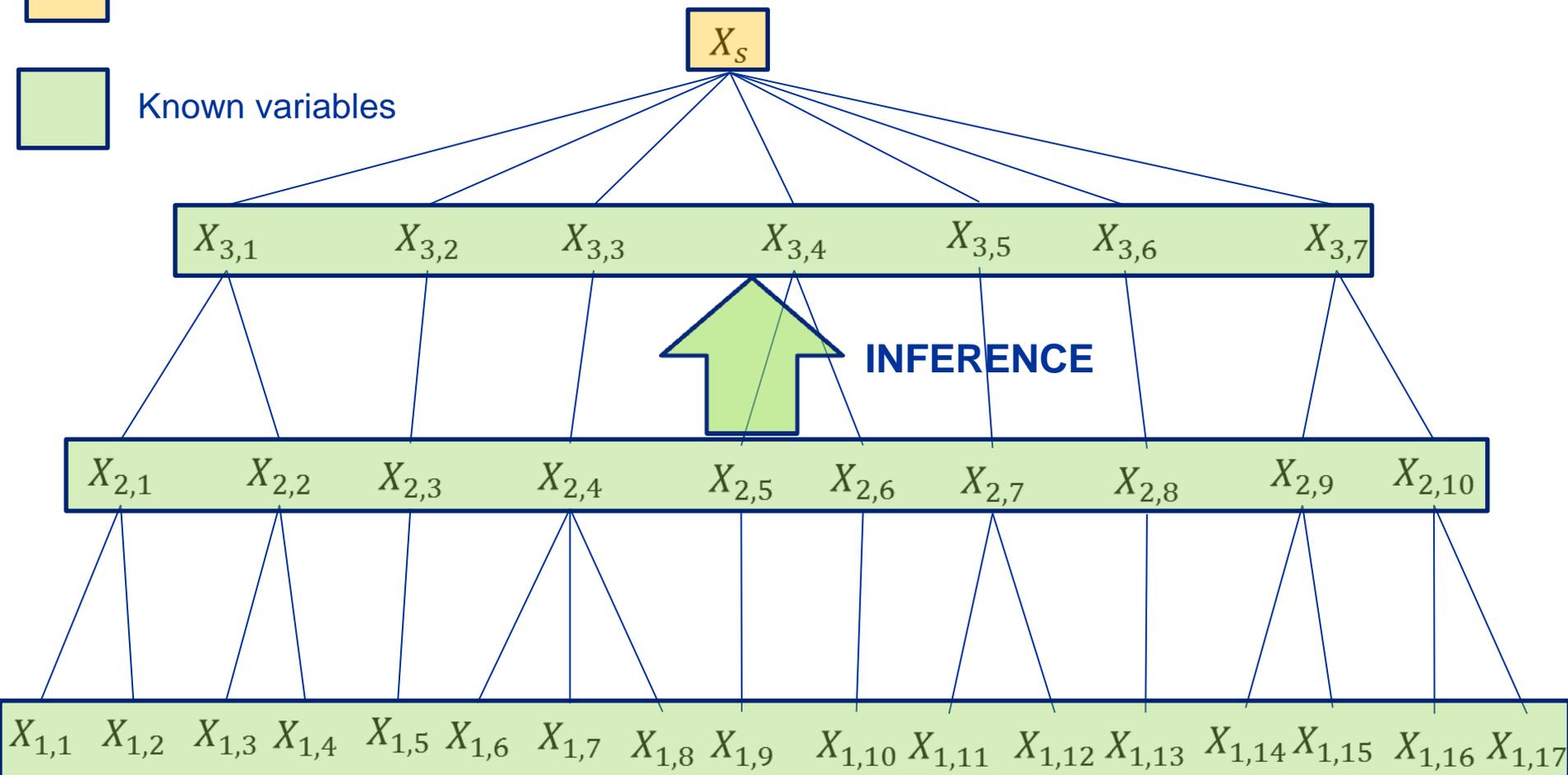
# Bottom-up Propagation of the belief



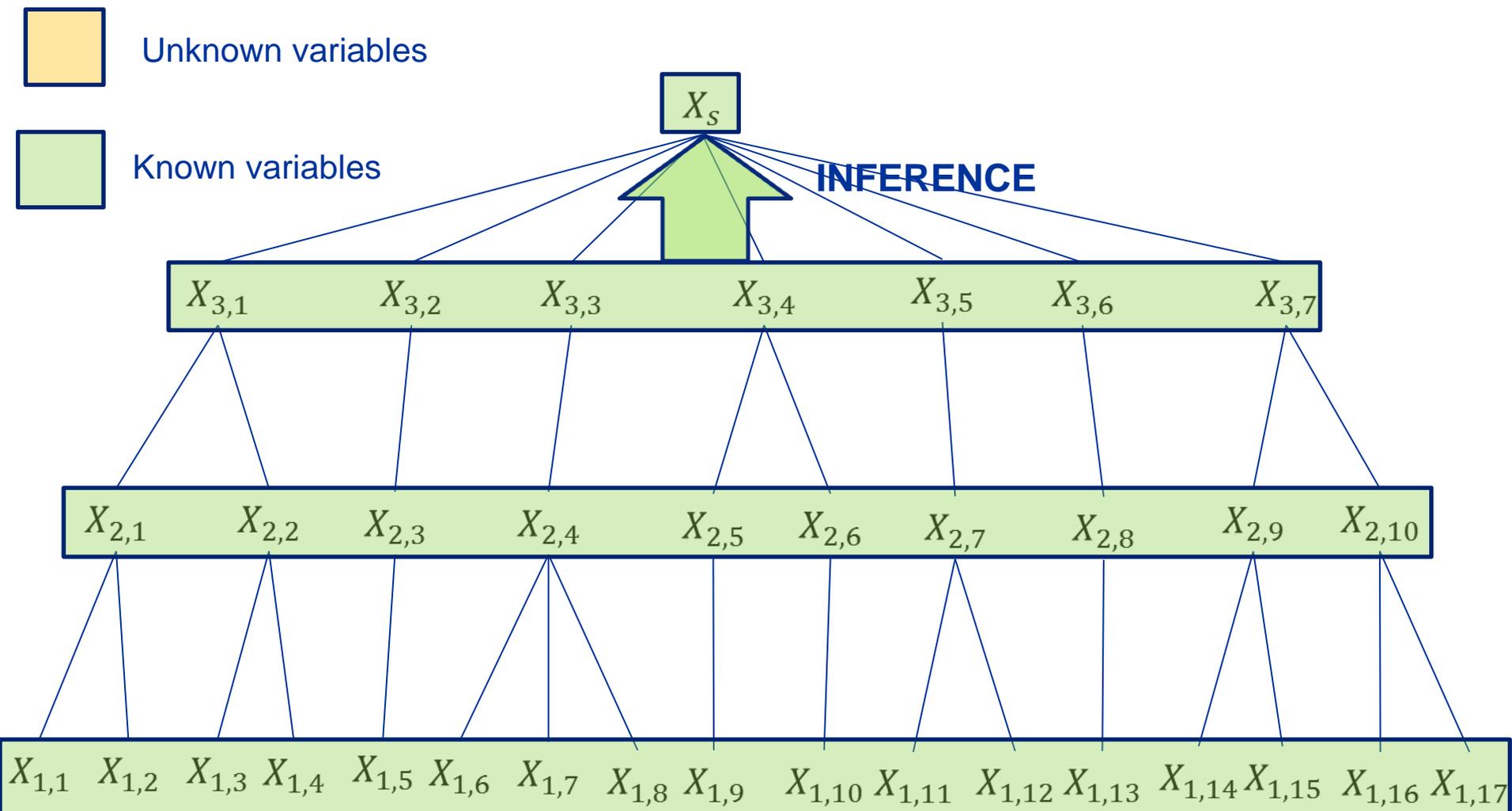
Unknown variables



Known variables



# Bottom-up Propagation of the belief



# Research of the Parse tree

- $X_S$  has only one possible value: S (Scene)
  - $m_{X_S}$  is defined on  $\{S, \emptyset\}$
  - $m_{X_S}(\emptyset)$  measures the consistency of the image interpretation
- Optimal interpretations of an image :
  - Parse tree minimizing  $m_{X_S}(\emptyset)$

# Experiments

- **Evidential grammars is a framework:**
  - The 5-tuple  $(S, V_N, V_T, \Gamma, \mathbf{M})$  has to be instantiated
  - For traffic scenes, we have to define:
    - The objects and parts of objects
    - Spatial relationships
    - The production rules
- **First step of experiments:**
  - input directly at the objects level
- **Second step of experiments:**
  - input at the part-of-objects level

# Ongoing student works

- **Semi-automatic data annotation**
  - 5 master students (Beihang University, Ecole Centrale)
- **Local feature (texture) analysis**
  - 1 master student (Peking University)
- **Part-based object detection**
  - Master internship (Spring 2013)
- **One class learning**
  - Master internship (Spring 2013)

# Ongoing and future publications

- **IAPR MVA 2013:**
  - Information Fusion on Oversegmented Images: An Application for Urban Scene Understanding
- **ORASIS 2013:**
  - Fusion d'informations sur des images sursegmentées : Une application à la compréhension de scènes routières
- **IUKM 2013:**
  - Evidential Grammars Framework for Image Interpretation. Application to Multimodal Traffic Scene Understanding