

# A Three-Layered Approach to Facade Parsing

Anđelo Martinović<sup>1</sup>   Markus Mathias<sup>1</sup>  
Julien Weissenberg<sup>2</sup>   Luc Van Gool<sup>1,2</sup>

<sup>1</sup>ESAT-PSI/VISICS, KU Leuven

<sup>2</sup>Computer Vision Laboratory, ETH Zurich

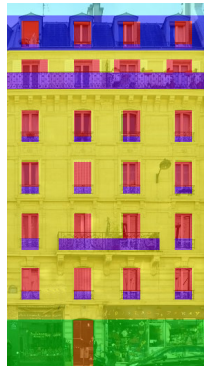


# We aim to improve the state of the art in facade parsing

From an image ...



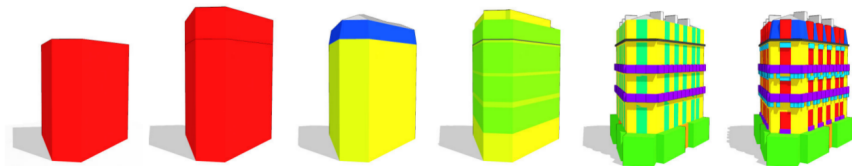
... to its labeling



- window
- wall
- balcony
- door
- roof
- sky
- shop

# We do not use shape grammars!

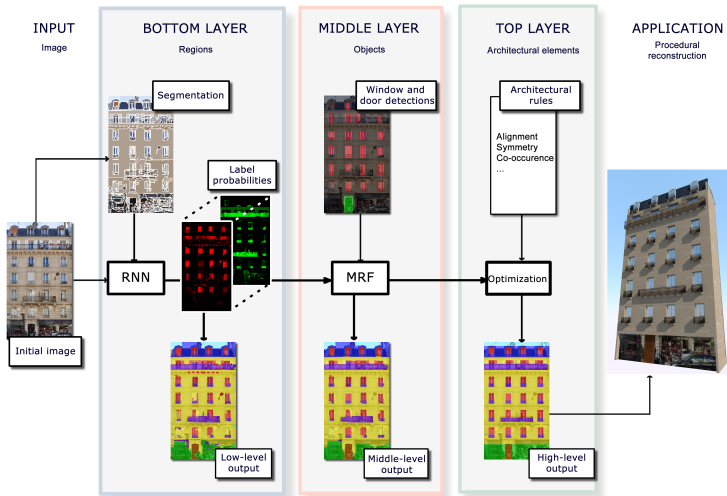
- State-of-the-art methods in facade parsing assume that an appropriate shape grammar is available [1].



- We do not use shape grammars as priors, and still achieve superior performance.

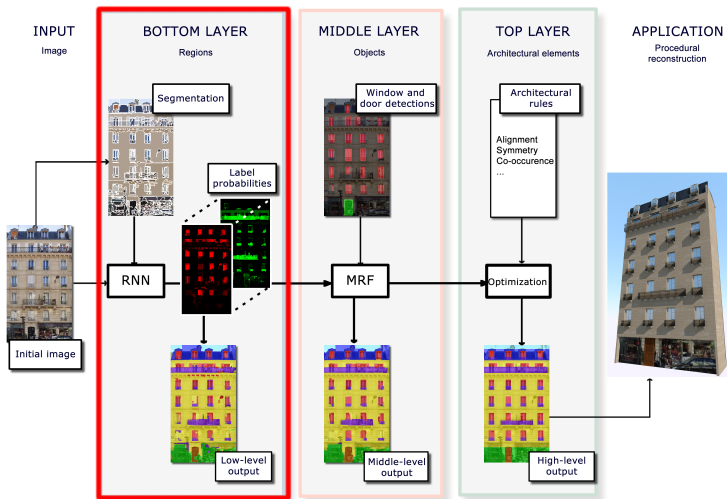
[1] Teboul, Kokkinos, Simon, Koutsourakis, Paragios: "Shape grammar parsing via Reinforcement Learning", CVPR, (2011).

# A Three-Layered Approach



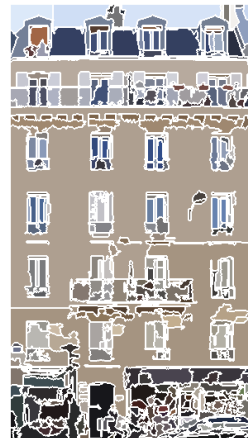


# Bottom layer - segments

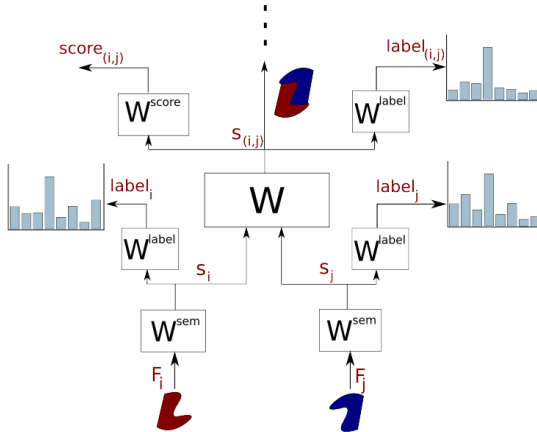


# Image preparation

- We segment the image using mean-shift.
- The appearance (color and texture), geometry, and location features are extracted for each region.
  - STAIR Vision Library
- This results in 225-dimensional feature vectors.



# Recursive Neural Network

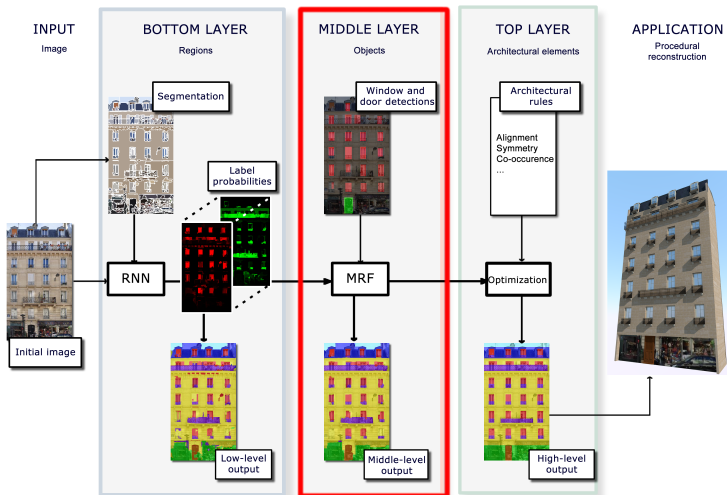


# Bottom Layer Output

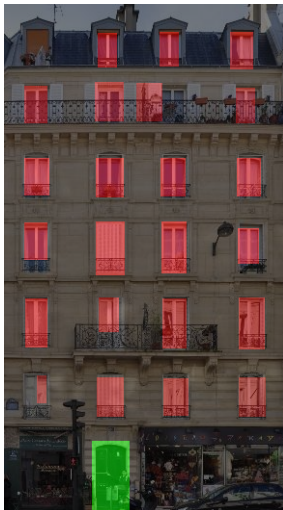


-  window
-  wall
-  balcony
-  door
-  roof
-  sky
-  shop

# Middle layer - objects



# Window and Door Detection



# Incorporating Detector Knowledge With MRFs

## Energy minimization with graph cuts

- Potts model

$$E(I) = \sum_{x_i} \phi_s(l_i | x_i) + \lambda \sum_{x_i} \sum_{x_j \sim x_i} \phi_p(l_i, l_j | x_i, x_j) \quad (1)$$

- Pairwise potentials

$$\phi_p(l_i, l_j | x_i, x_j) = \begin{cases} 0, & \text{if } l_i = l_j \\ 1, & \text{otherwise} \end{cases} \quad (2)$$

- Unary potentials

$$\phi_s(l_i | x_i) = -\log p(l_i | RNN(x_i)) - \sum_k \alpha_k \log p(l_i | D_k(x_i)) \quad (3)$$

# Incorporating Detector Knowledge With MRFs

## Energy minimization with graph cuts

- Potts model

$$E(I) = \sum_{x_i} \phi_s(l_i | x_i) + \lambda \sum_{x_i} \sum_{x_j \sim x_i} \phi_p(l_i, l_j | x_i, x_j) \quad (1)$$

- Pairwise potentials

$$\phi_p(l_i, l_j | x_i, x_j) = \begin{cases} 0, & \text{if } l_i = l_j \\ 1, & \text{otherwise} \end{cases} \quad (2)$$

- Unary potentials

$$\phi_s(l_i | x_i) = -\log p(l_i | RNN(x_i)) - \sum_k \alpha_k \log p(l_i | D_k(x_i)) \quad (3)$$



# Incorporating Detector Knowledge With MRFs

## Energy minimization with graph cuts

- Potts model

$$E(l) = \sum_{x_i} \phi_s(l_i | x_i) + \lambda \sum_{x_i} \sum_{x_j \sim x_i} \phi_p(l_i, l_j | x_i, x_j) \quad (1)$$

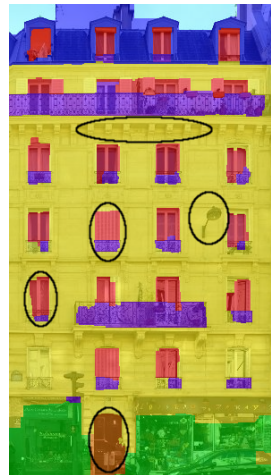
- Pairwise potentials

$$\phi_p(l_i, l_j | x_i, x_j) = \begin{cases} 0, & \text{if } l_i = l_j \\ 1, & \text{otherwise} \end{cases} \quad (2)$$

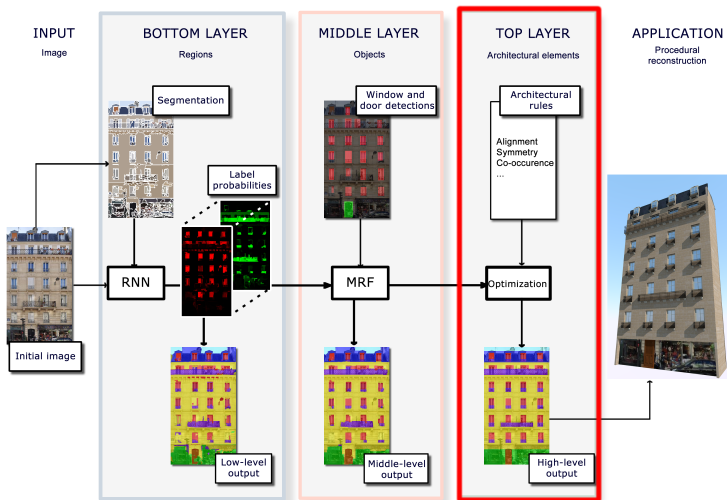
- Unary potentials

$$\phi_s(l_i | x_i) = -\log p(l_i | RNN(x_i)) - \sum_k \alpha_k \log p(l_i | D_k(x_i)) \quad (3)$$

# From Bottom To Middle Layer Output



# Top layer - architectural elements

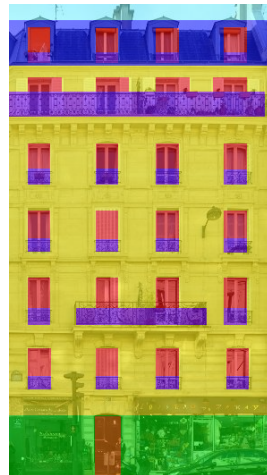
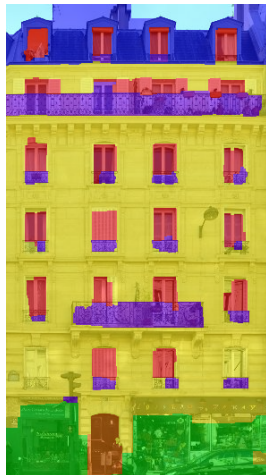


# Weak Architectural Principles

- Soft constraints instead of fixed grammar structure
- Only enforced if there is enough image support

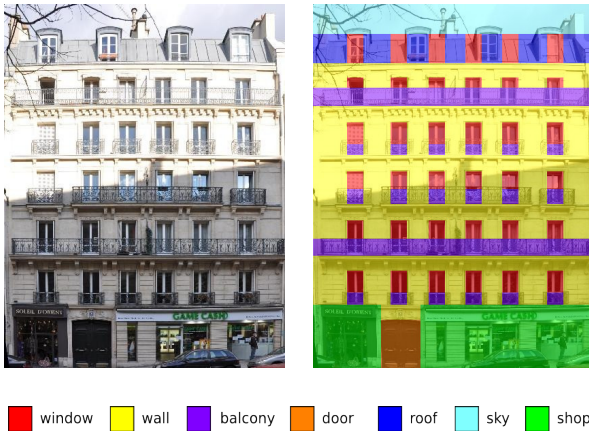
Principle	Alter	Add	Remove
Vertical and horizontal (non)alignment	✓	-	-
Window similarity	-	✓	-
Facade symmetry	-	✓	✓
Element co-occurence	-	✓	✓
Equal width/height in a row or column	✓	-	-
Door hypothesis	✓	✓	✓
Vertical region order	✓	-	-

# From Middle To Top Layer Output



# Ecole Centrale Paris Facades Database [2]

- Contains 104 rectified and cropped Haussmannian facades.



[2] Teboul, O. , "Ecole Centrale Paris Facades Database" (2010).

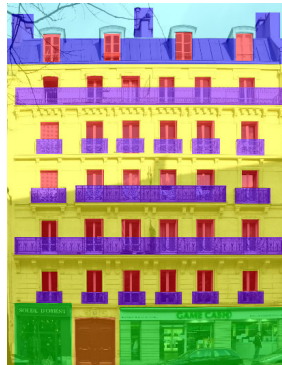
# Ecole Centrale Paris Facades Database

- Original labeling is plausible, but imprecise.
- We provide more precise annotations (available online).

Old annotation



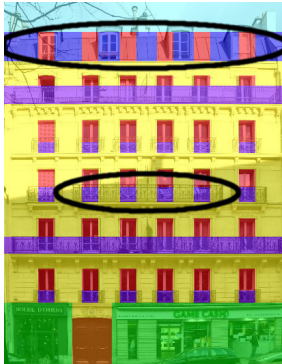
New annotation



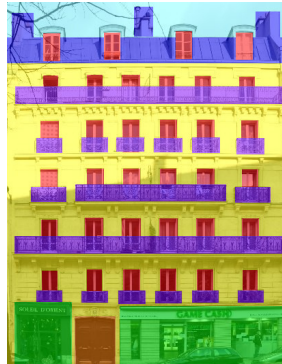
# Ecole Centrale Paris Facades Database

- Original labeling is plausible, but imprecise.
- We provide more precise annotations (available online).

Old annotation



New annotation





# Results - ECP Dataset

Class	Baseline[4]	Layer 1	Layer 2	Layer 3
<i>window</i>	62	62	69	<b>75</b>
<i>wall</i>	82	91	<b>93</b>	88
<i>balcony</i>	58	<b>74</b>	71	70
<i>door</i>	47	43	60	<b>67</b>
<i>roof</i>	66	70	73	<b>74</b>
<i>sky</i>	95	91	91	<b>97</b>
<i>shop</i>	88	79	86	<b>93</b>
Pixel acc.	74.71	82.63	<b>85.06</b>	84.17

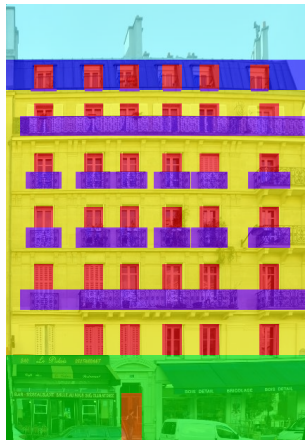
[4] Teboul, O., "Shape Grammar Parsing: Application to Image-based Modeling" (2011).

# Pixel Accuracy vs Visual Effect

Pixel accuracy: 89.48%



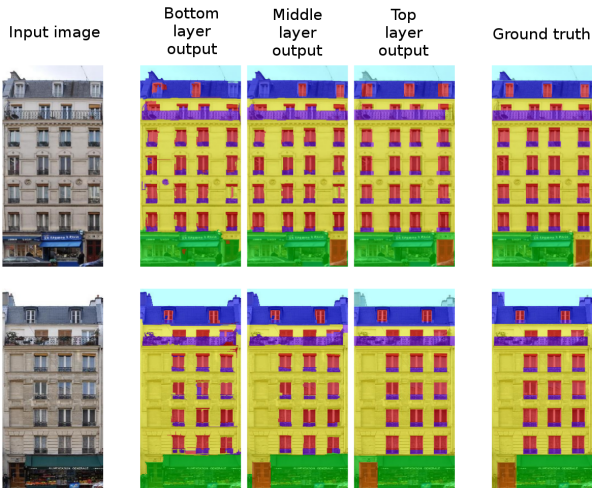
Pixel accuracy: 87.82%



# Results - ECP Dataset

Class	Baseline[4]	Layer 1	Layer 2	Layer 3
<i>window</i>	62	62	69	<b>75</b>
<i>wall</i>	82	91	<b>93</b>	88
<i>balcony</i>	58	<b>74</b>	71	70
<i>door</i>	47	43	60	<b>67</b>
<i>roof</i>	66	70	73	<b>74</b>
<i>sky</i>	95	91	91	<b>97</b>
<i>shop</i>	88	79	86	<b>93</b>
Pixel acc.	74.71	82.63	<b>85.06</b>	84.17
Class acc.	71.14	72.86	77.46	<b>80.71</b>

# Example Outputs - ECP Dataset



## eTRIMS Database [3]

- Contains 60 images of various building styles.
- We perform automatic rectification.



■ building ■ car ■ door ■ pavement ■ road ■ sky ■ vegetation ■ window

[3] Korč, F. and Förstner, W., "eTRIMS Image Database for Interpreting Images of Man-Made Scenes" (2009).

# Example Outputs - eTRIMS Dataset

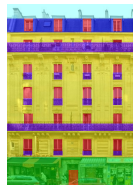
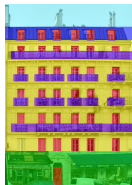
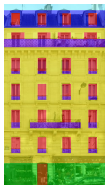
Input image

Bottom  
layer  
outputMiddle  
layer  
outputTop  
layer  
output

Ground truth



# Example Outputs - Procedural Models



# Summary

- We developed a **novel three-layer approach** for facade parsing.
- We **significantly outperform** the state-of-the-art on two facade parsing datasets.
- We utilize the concept of **weak architectural knowledge**.
- Outlook
  - So far, the inferred procedural models are instance-specific.
  - We want to generalize between buildings of the same style.
  - As we no longer depend on grammars as priors, can we instead induce them from the data?



# Questions?



Andelo Martinović

<http://homes.esat.kuleuven.be/~amartino/>

Available online: updated ECP annotations, paper manuscript, supplementary material, spotlight video








**ETH**

Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich



## References

-  [1] Teboul, O. and Kokkinos, I. and Simon, L. and Koutsourakis, P. and Paragios, N. , "Shape grammar parsing via Reinforcement Learning" (2011).
-  [2] Teboul, O. , "Ecole Centrale Paris Facades Database" (2010).
-  [3] Korč, F. and Förstner, W., "eTRIMS Image Database for Interpreting Images of Man-Made Scenes" (2009).
-  [4] Teboul, O., "Shape Grammar Parsing: Application to Image-based Modeling" (2011).
-  [5] Yang, M.Y. and Förstner, W. , "Regionwise Classification of Building Facade Images", Springer (2011).
-  [6] Socher et al. , "Parsing Natural Scenes and Natural Language with Recursive Neural Networks", ICML (2011).

## Results - eTRIMS Dataset

The results for eTrims were obtained by automatically rectifying both the input images and the ground truth labelings. Our results were computed in the rectified space. As previous work did not perform any rectification, we repeated the evaluation by “unrectifying” our output labeling and comparing to the original ground truth. The results obtained in this way are actually better by ~1% than reported in the paper.

Class	Baseline[5]	Layer 1	Layer 2	Layer 3
<i>building</i>	71	88	<b>91</b>	87
<i>car</i>	35	<b>69</b>	<b>69</b>	<b>69</b>
<i>door</i>	16	<b>25</b>	18	19
<i>pavement</i>	22	<b>34</b>	33	<b>34</b>
<i>road</i>	35	<b>56</b>	55	<b>56</b>
<i>sky</i>	78	<b>94</b>	93	<b>94</b>
<i>vegetation</i>	66	<b>89</b>	<b>89</b>	88
<i>window</i>	75	71	74	<b>79</b>
Pixel acc.	65.8	81.87	<b>83.16</b>	81.63
Class acc.	49.75	<b>65.85</b>	65.4	65.6