**Slide 1**

# Aerial Image Labeling

Quantitative Big Imaging Course 2016

Javier Montoya, PhD
montoya@ethz.ch
PostDoc @ ScopeM

---

**Slide 2**

## What do I mean by **Aerial Image Labeling**?

- **Task:** divide a given input image into a set of semantic coherent regions: road networks and buildings.

Buildings
Road Networks

Input Image → Semantic Segmented Image

---

**Slide 3**

## Why is it **Important?**

Earth Observation & Environmental Modeling

Urban Planning

Virtual Representations / 3D city models

Location-Aware Applications / Navigation Maps

---

**Slide 4**

## Which are the **existing techniques?**

**Deterministic Representations:** *detected pieces of objects are stitched together using low-level image processing, e.g. Miao et.al. (GRSL'13), Poullis et.al. (JPRS'10).*

- ✓ Successful when objects appear more clearly (no occlusions, etc).
- ✗ Many parameters must be tuned empirically.
- ✗ Errors from each step are propagated.

**Local Statistical Approaches:** *set of features are locally extracted to train class-specific models, e.g. Dollar et.al. (CVPR'06), Mnih et.al. (ECCV'10).*
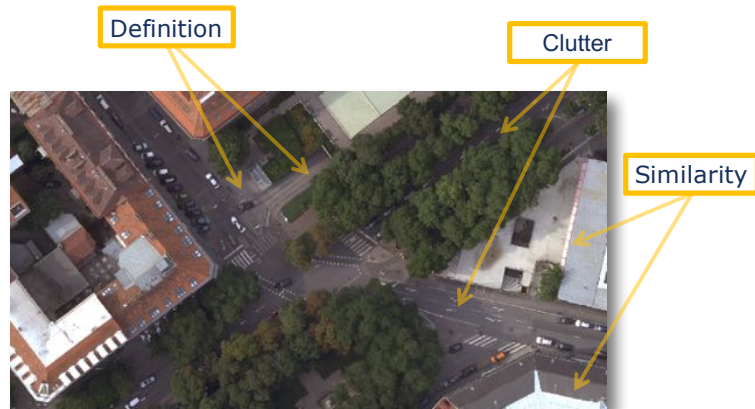
- ✓ Robust models can be obtained.
- ✗ Availability of training datasets, hand-labeling process.
- ✗ Predictions are focused on local information only.

**Probabilistic Representations of Image Context:** *high-level semantic knowledge is incorporated, e.g. Lacoste (PAMI'05), Türetken et.al. (CVPR'13).*

- ✓ Encode rich semantic level information.
- ✓ Smooth and precise segmentations.
- ✗ Inference can be time consuming, e.g. Markov Point Processes.
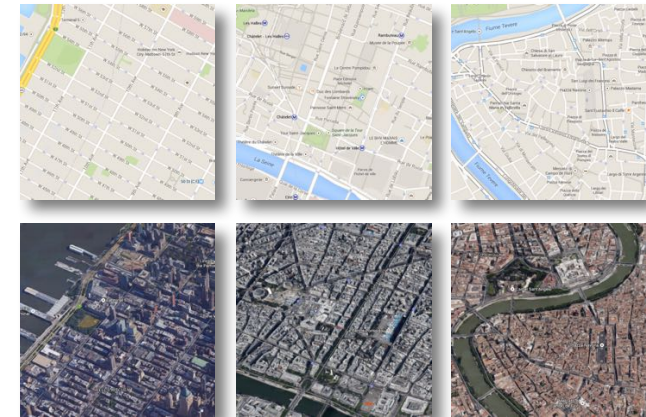
**Slide 5 — What are the Challenges?**

Definition — Clutter — Similarity

… unsolved problem since almost 40 years! *Bajcsy et al. (TSMC'76)*

---

**Slide 6 — What are the Challenges?**

New York — Paris — Rome

+ Complex Structural Prior

---

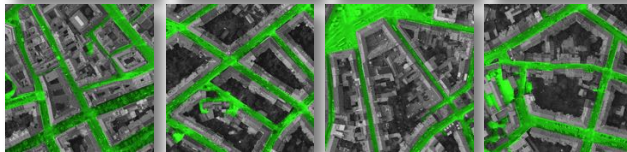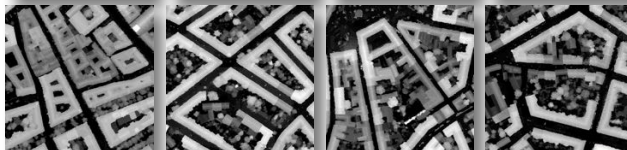**Slide 7 — Graz Road Dataset**

Image Tiles
GTs
nDSMs

(i) Urban region, two classes (road vs. background), 67 tiles of 1000x1000 pixels, manual annotations.
(ii) Road networks: major avenues + secondary streets, change slow in width/curvature.
(iii) Presence of occlusions, e.g. trees and cars.
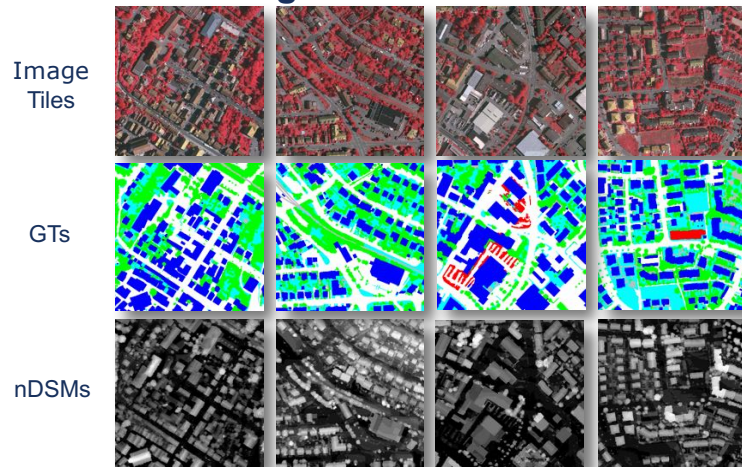
---

**Slide 8 — Vaihingen Road Dataset**

Image Tiles
GTs
nDSMs

(i) Countryside region, two classes (road vs. background), 16 tiles of 1000x1000 pixels, manual annotations.
(ii) Road networks: very irregular, mainly narrow, partially occluded trees, shadows.

## Vaihingen Multi-class Dataset



Image Tiles

GTs

nDSMs

(i) Six classes: natural ground, background, roads, trees, grass, and buildings.
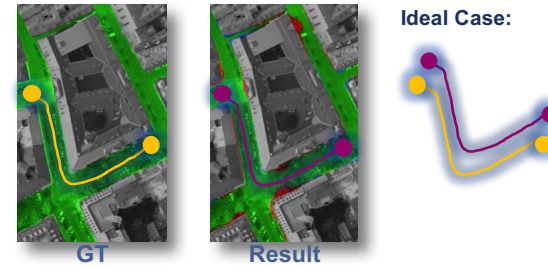(ii) Buildings: vary strongly in shape and are often densely clustered.

---

## How are the results **measured?**

**Pixel-wise classification accuracy:**
- F1-score, Precision, Recall.

**Road Network Topology:**
- What fraction of connecting paths between road seeds have the correct length within 5% tolerance, are respectively:

**Ideal Case:**



GT          Result

---

## How are the results **measured?**

**Pixel-wise classification accuracy:**
- F1-score, Precision, Recall.

**Road Network Topology:**
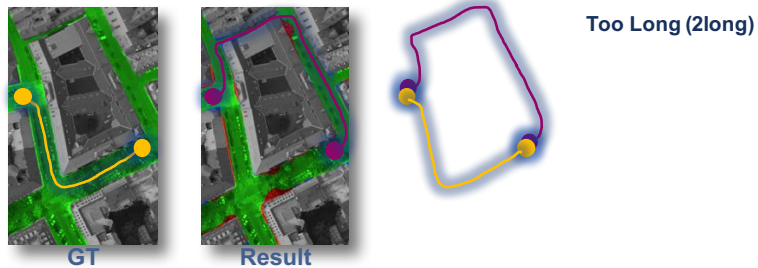- What fraction of connecting paths between road seeds have the correct length within 5% tolerance, are respectively:

**Too Long (2long)**



GT          Result

---

## How are the results **measured?**

**Pixel-wise classification accuracy:**
- F1-score, Precision, Recall.

**Road Network Topology:**
- What fraction of connecting paths between road seeds have the correct length within 5% tolerance, are respectively:
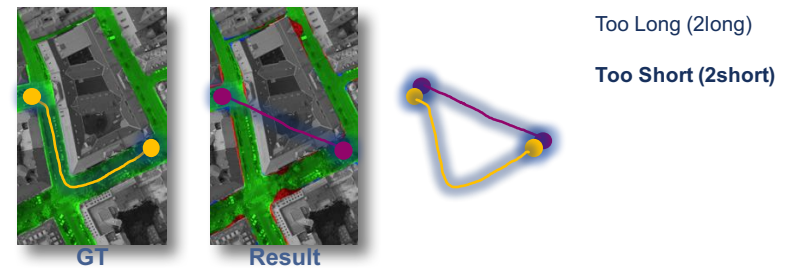
Too Long (2long)

**Too Short (2short)**



GT          Result

## How are the results **measured?**

**Pixel-wise classification accuracy:**
- F1-score, Precision, Recall.

**Road Network Topology:**
- What fraction of connecting paths between road seeds have the correct length within 5% tolerance, are respectively:
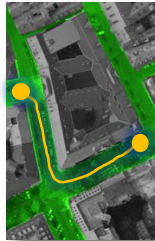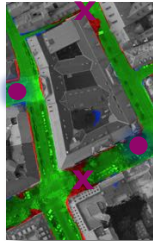


GT    Result

Too Long (2long)

Too Short (2short)

**No Connectivity (NoConn)**

---

## How are the results **measured?**
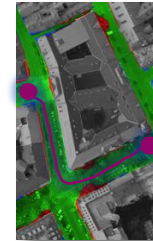
**Pixel-wise classification accuracy:**
- F1-score, Precision, Recall.

**Road Network Topology:**
- What fraction of connecting paths between road seeds have the correct length within 5% tolerance, are respectively:
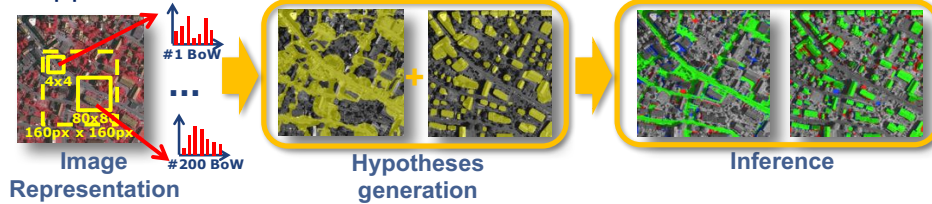


GT    Result

$$\text{TopoCorrectness} = 100\% - 2\text{Short} - 2\text{Long} - \text{noConn}$$

---

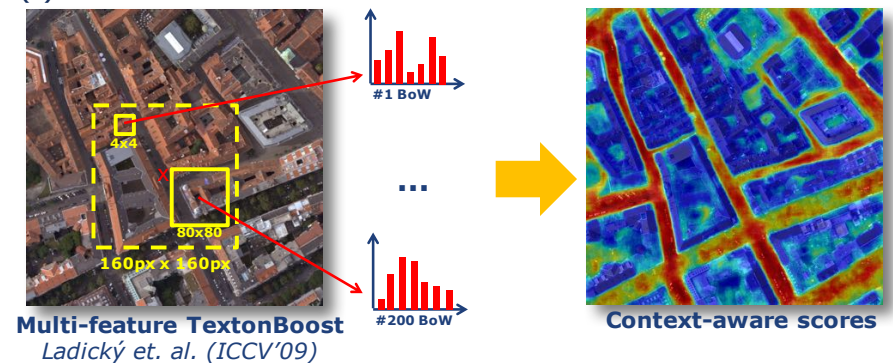## Class-Specific higher-order cliques

Approach Overview:



#1 BoW

#200 BoW

**Image Representation**

**Hypotheses generation**

**Inference**

**Model**:
(i) Multilabel pixelwise classification using powerful neighbor features.
(ii) Overcomplete representation of *building* and *road* candidates.
(iii) Candidates are prunned to optimal subset through CRF.

---

## Class-Specific higher-order cliques

**(I) Context-aware road scores**



#1 BoW

#200 BoW
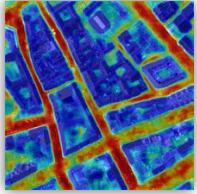
**Multi-feature TextonBoost**
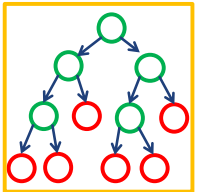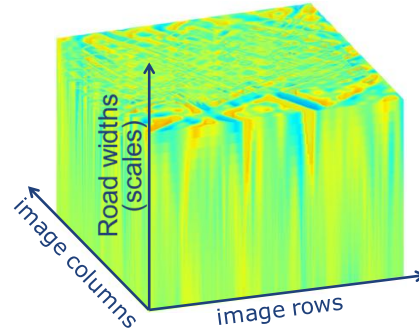*Ladický et. al. (ICCV'09)*

**Context-aware scores**

- Multi-label pixelwise classification.
- Self context/Local layout, appearance information is encoded over large spatial neighborhoods.

## Class-Specific higher-order cliques
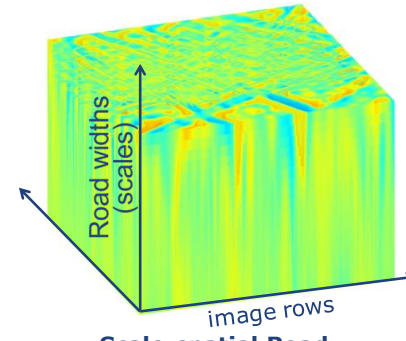
(II) Hyphoteses Generation: Roads



**Context-aware road scores**
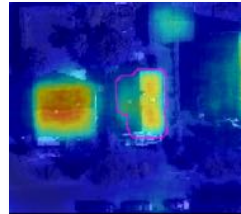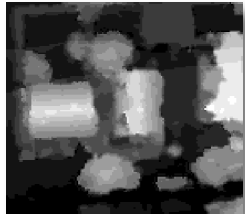
**Random Forest estimates road widths**

Road widths (scales)

image columns

image rows

**Likelihoods per road width**

---

## Class-Specific higher-order cliques

(II) Hyphoteses Generation: Roads



Road widths (scales)

image rows

**Scale-spatial Road likelihoods**

**Sampled Road candidates**

---

## Class-Specific higher-order cliques

(II) Hyphoteses Generation: Buildings



- Building candidates based on classifier scores.
- Connected components with high building likelihood.
- Building segments are approximated through alpha-shapes.

---

## Class-Specific higher-order cliques

(III) Hypheses Selection: Inference



**Sampled candidates** (high recall)    **Candidate selection** (high precision)

$$E = \sum_{pixels} E_u(x_{pix}) + \sum_{pix} \sum_{neigh} E_p(x_{pix}, x_{neigh}) + \sum_{roads} E_R(Q_m) + \sum_{build.} E_B(Q_n)$$

class evidence        soft-clique membership

## Slide 21

### CRF-Model for road superpixel segmentation

**Baselines:**

$$E = \sum_{pixels} E_u(x_{pix}) + \sum_{pix}\sum_{neigh} E_p(x_{pix}, x_{neigh}) + \sum_{roads} E_R(Q_m)$$



**Winn**

Road/Background Pixel-likelihoods
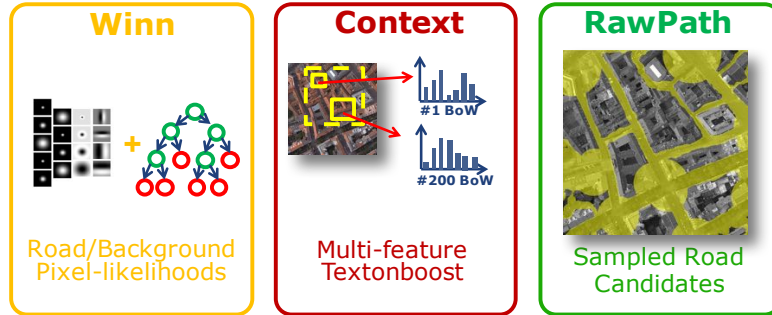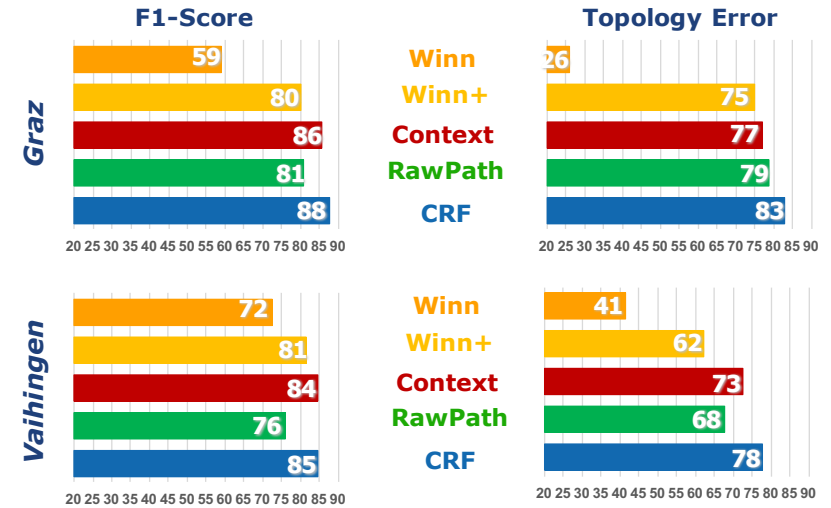
**Context**

#1 BoW
#200 BoW

Multi-feature Textonboost

**RawPath**

Sampled Road Candidates

## Slide 22

### Class-Specific higher-order cliques

**F1-Score**

**Topology Error**

*Graz*

| | Winn | Winn+ | Context | RawPath | CRF |
|---|---|---|---|---|---|
| F1-Score | 59 | 80 | 86 | 81 | 88 |
| Topology Error | 26 | 75 | 77 | 79 | 83 |

20 25 30 35 40 45 50 55 60 65 70 75 80 85 90

*Vaihingen*

| | Winn | Winn+ | Context | RawPath | CRF |
|---|---|---|---|---|---|
| F1-Score | 72 | 81 | 84 | 76 | 85 |
| Topology Error | 41 | 62 | 73 | 68 | 78 |

20 25 30 35 40 45 50 55 60 65 70 75 80 85 90

## Slide 23

### Class-Specific higher-order cliques

**Visual Results on Road Network Extraction:**

True Positives
False Negatives
False Positives



**Winn**

**Context**

**RawPath**

**Ours**

Road network on GRAZ, img6

## Slide 24

### Class-Specific higher-order cliques

**Experimental Results on Joint Road Networks + Buildings:**

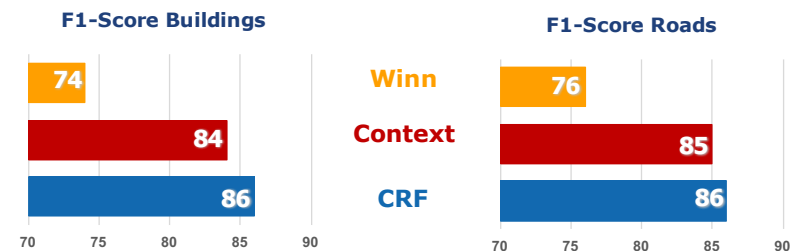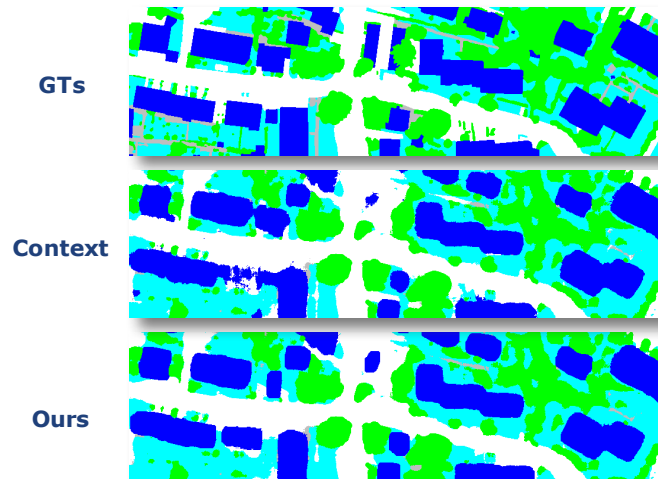$$E = \sum_{pixels} E_u(x_{pix}) + \sum_{pix}\sum_{neigh} E_p(x_{pix}, x_{neigh}) + \sum_{roads} E_R(Q_m) + \sum_{build.} E_B(Q_n)$$

**F1-Score Buildings**

**F1-Score Roads**

| | Winn | Context | CRF |
|---|---|---|---|
| F1-Score Buildings | 74 | 84 | 86 |
| F1-Score Roads | 76 | 85 | 86 |

70   75   80   85   90

Labeling Performance on the **Vaihingen Multi-class Dataset** (all numbers percentages).

**Slide 25:**

## Class-Specific higher-order cliques

Visual Results on Joint Road Networks + Buildings:



GTs

Context

Ours

---

**Slide 26:**

## Class-Specific higher-order cliques

Visual Results on Joint Roads + Buildings:

True Positives
False Negatives
False Positives



Buildings

Context

Ours

Road Networks

Context

Ours

---

**Slide 27:**

## Conclusions

- **Class-specific Priors:**
  - Higher-level representations for buildings and roads are useful multi-class segmentation.
  - Buildings are represented as a set of compact-like polygons.
  - Roads are modeled as a collection of long, narrow segments.

---

**Slide 28:**

## Outlook & Future Work

- **Generation of object candidates:**
  - Probabilistic Sampling Scheme?
  - Hypotheses parameters as regression task?

- **Training Dataset:**
  - Manual labeling of roads is time-consuming and costly
    => use publicly available data as ground truth, e.g. Open Street Map.
    => deep learning.

- **Applicability in other domains:**
  - Generic model potentially applicable to other networks such as neurons and vessel in medical imaging.