# REVISITING THE COCO PANOPTIC METRIC
## TO ENABLE VISUAL AND QUALITATIVE ANALYSIS OF HISTORICAL MAP INSTANCE SEGMENTATION

### E. CARLINET AND J. CHAZALON
EPITA Research & Development Laboratory (LRDE)

## MOTIVATION

- Evaluate historical map (dense instance) segmentation (ICDAR21 MapSeg competition)
- Identify uses for major document segmentation metrics, and compare them.
- Propose a unifying framework for these metrics.
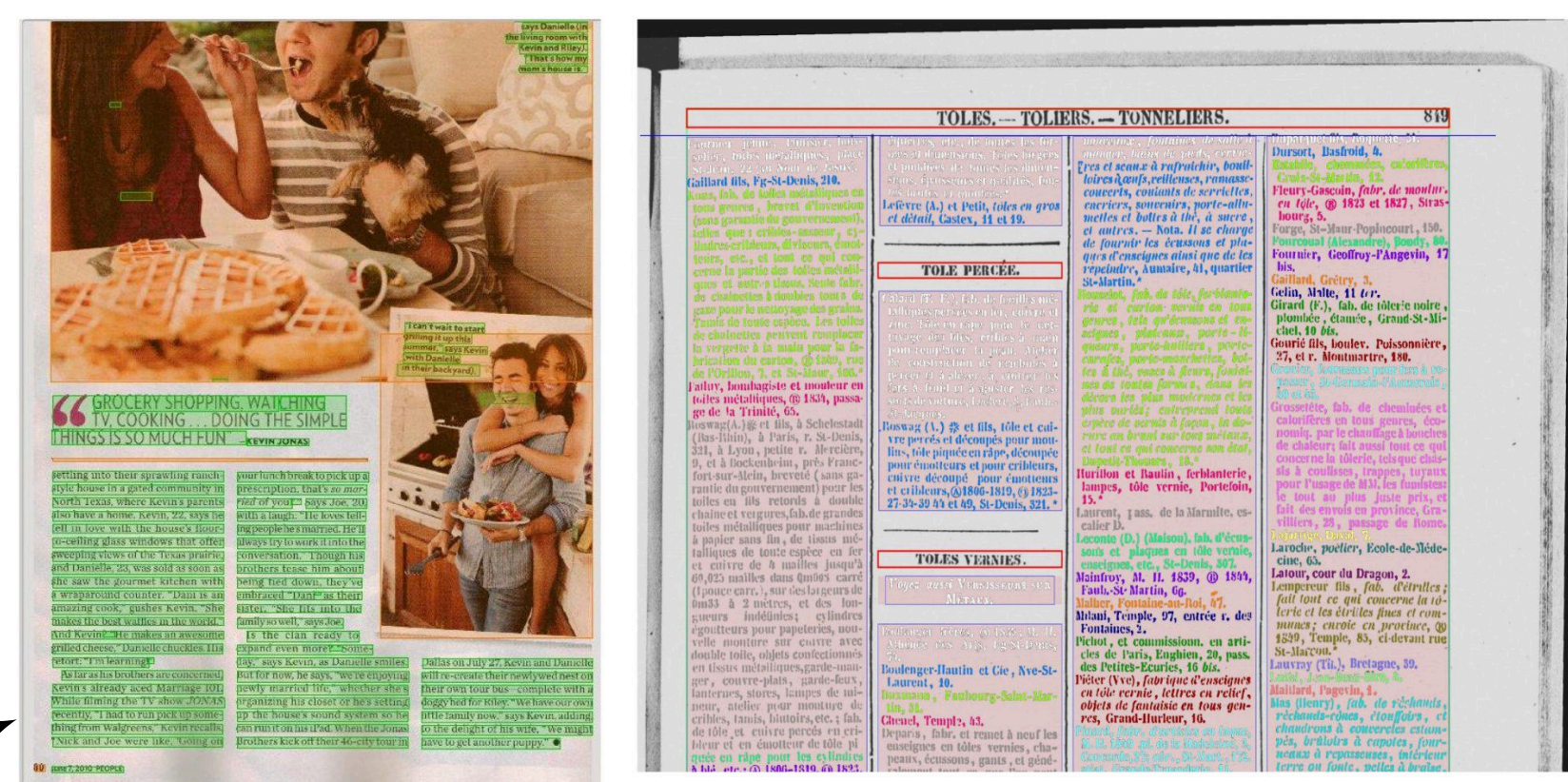- Provide new qualitative and quantitative tools to leverage these metrics.

https://github.com/icdar21-mapseg/coco-pano-ext-demo

pip install coco-pano-ext-demo

## SEGMENTATION METRICS

What makes a metric bad ?
1. Not normalized/with no bounds (*"Congrats you've got 641! All contenders are between 111 and 4747"*)
2. Does not fit the human perception of "good" (*"The method B totally messed up and got the score 0.951"*)
3. Many thresholds/parameters (*"for our experiments, we use $\alpha_1 = 0.51$ and $\zeta_{26} = 3.29$"*)
4. Hard to explain (*"because a log makes it look smarter"*)
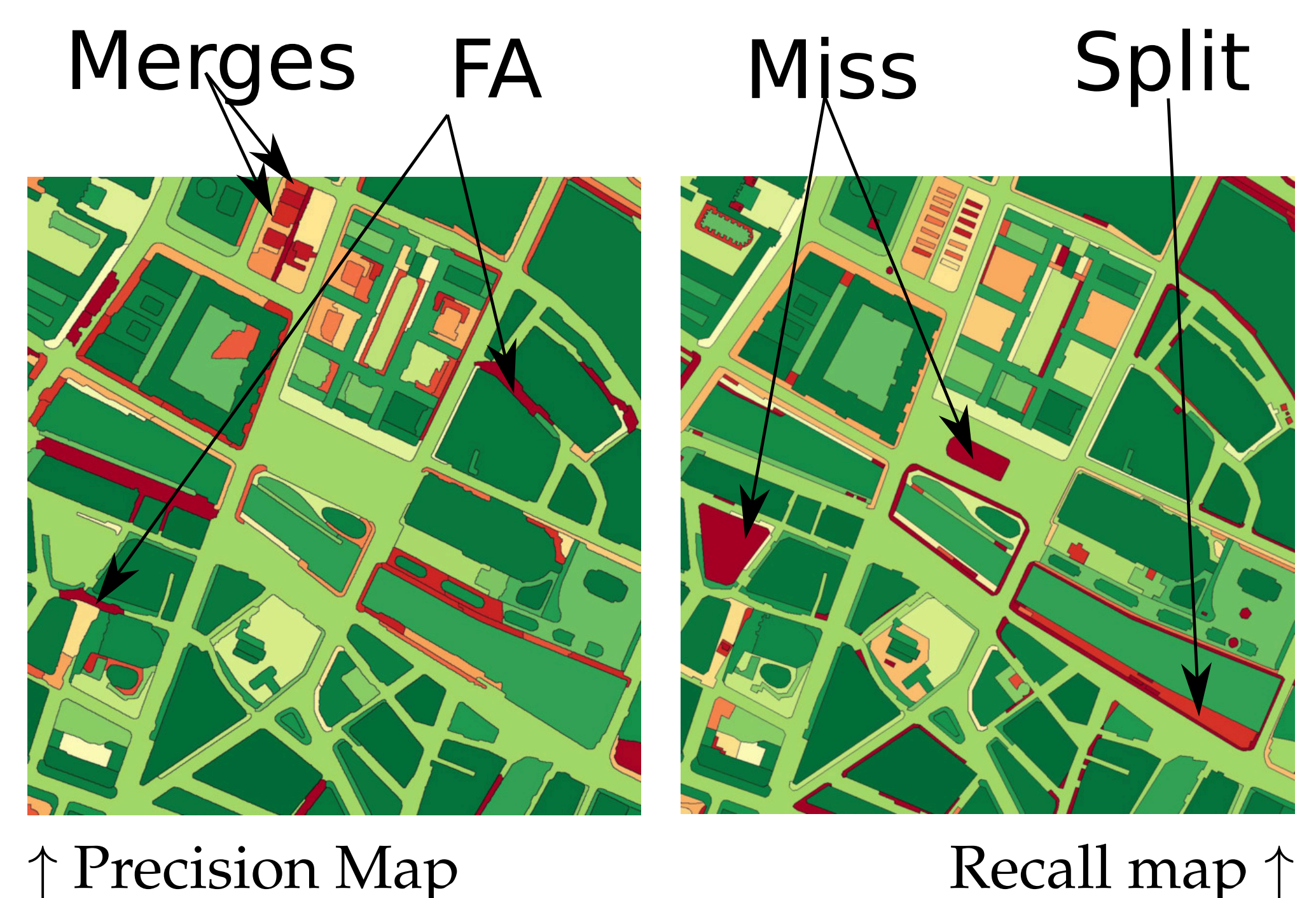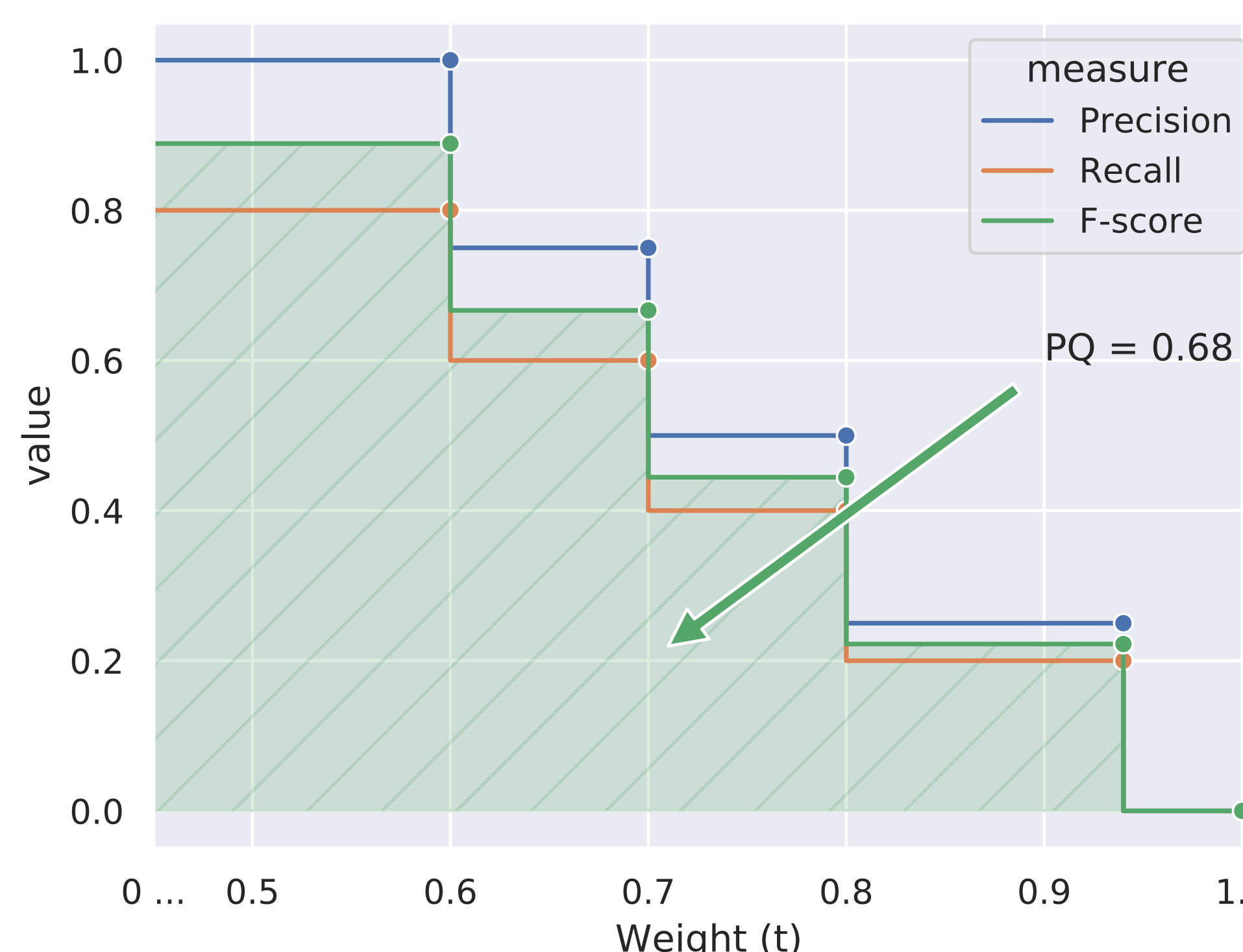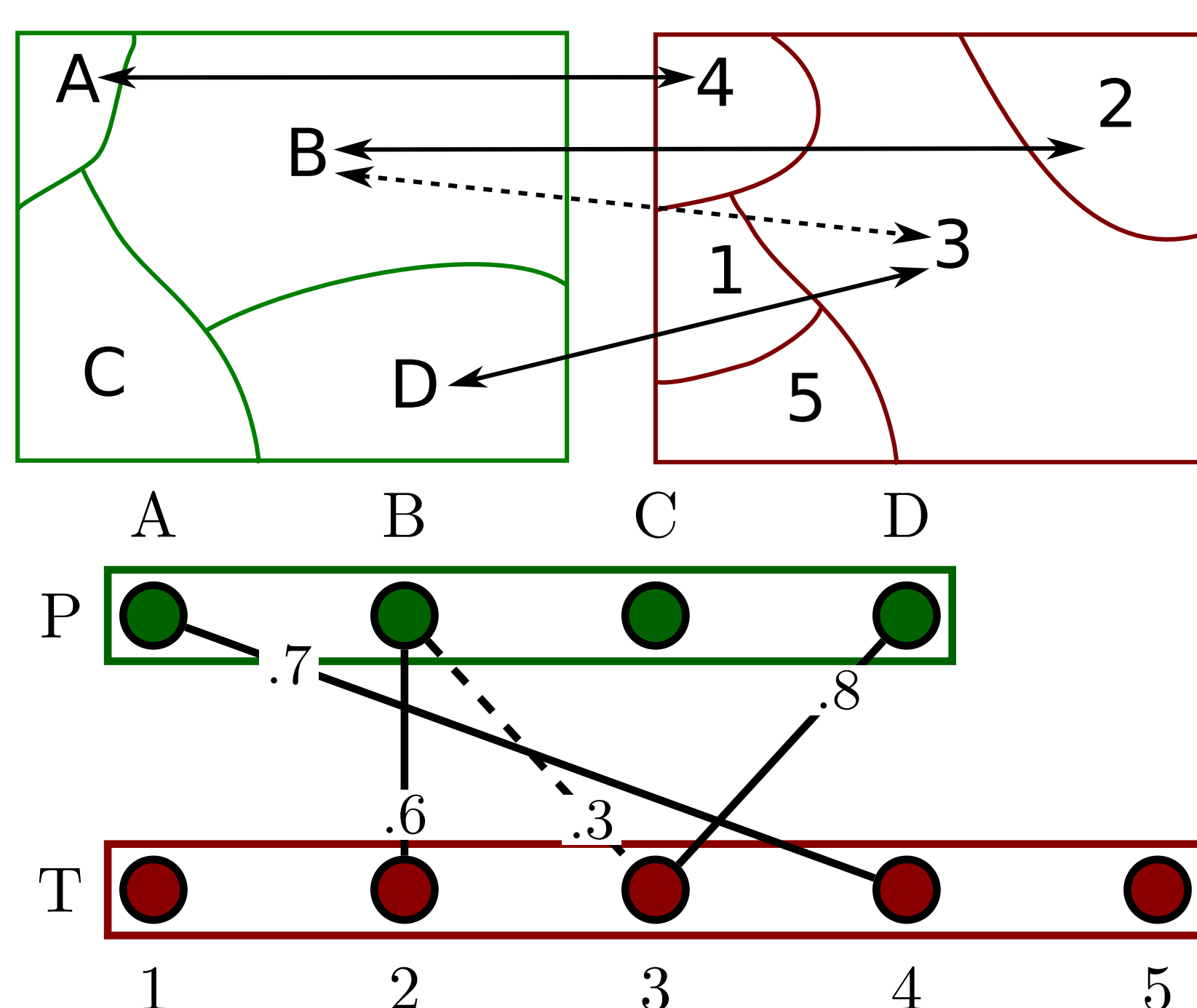
### Document-oriented: DetEval, ZoneMap

- ✓ Consider fragmentation (1-to-many relations)
- ✗ Under/over-segmentation is somehow over-scoring
- ✗ ZoneMap: Non-normalized scores
- ✗ DetEval: calibration of several thresholds

### Computer Vision-related: COCO Panoptic

- ✓ Simple: a combination of "segmentation accuracy" and "detection performance"
- ✓ A single parameter, easy to tune and sensible (the IoU min)
- ✗ Single global score (ok for ranking but bad for a system performance analysis)

## A FRAMEWORK FOR WEIGHTED PAIRING EVALUATION



PQ = 0.68

Merges  FA  Miss  Split

↑ Precision Map          Recall map ↑

### 1-1 Pairing

- Edges weighted with IoU (or other)
- A single and meaningful parameter $\alpha$
  - With IoU, $\alpha = 0.5 \rightarrow$ pairing
  - **Tolerance Threshold**: the level over which no human correction is required
- Compute the number of remaining pairing $T_c$ as the IoU increases.

### Toward COCO PQ - Quantitative Analysis

$$prec. = \frac{T_c}{|P|} \quad recall = \frac{T_c}{|T|} \quad \text{F-score} = \frac{2.T_c}{|P| + |T|}$$

- Precision/Recall/F-Score Curves to analyze systems performances  *New!*
- The COCO PQ is the AUC of the F-Score curve ! *Bang!*

### Prec./Recall Maps - Qualitative Analysis

*New!* A region is valued by the best match score in the pairing.

- Precision Map. Assess the quality of detections (objects from the system)
- Recall Map. Detect the missed components (objects from the GT)