

Segmenting Video Into Classes of Algorithm-Suitability

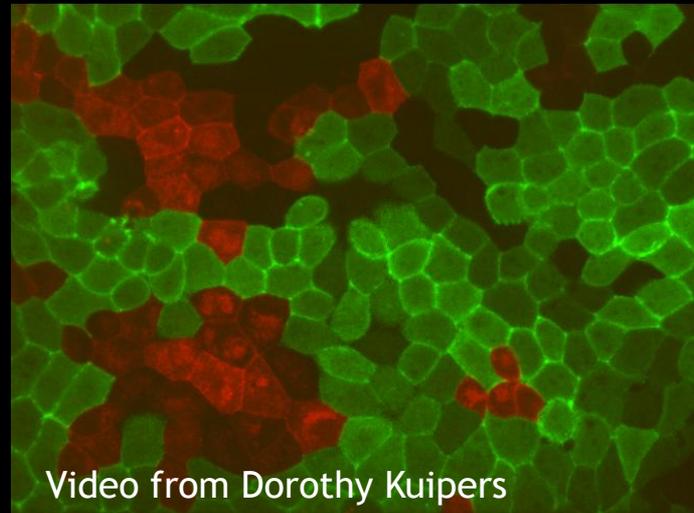
Oisin Mac Aodha (UCL)

Gabriel Brostow (UCL)

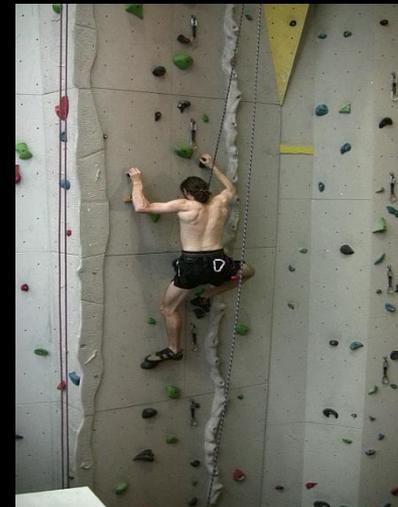
Marc Pollefeys (ETH)



Which algorithm should I (use / download / implement) to track things in this video?



Video from Dorothy Kuipers



The Optical Flow Problem

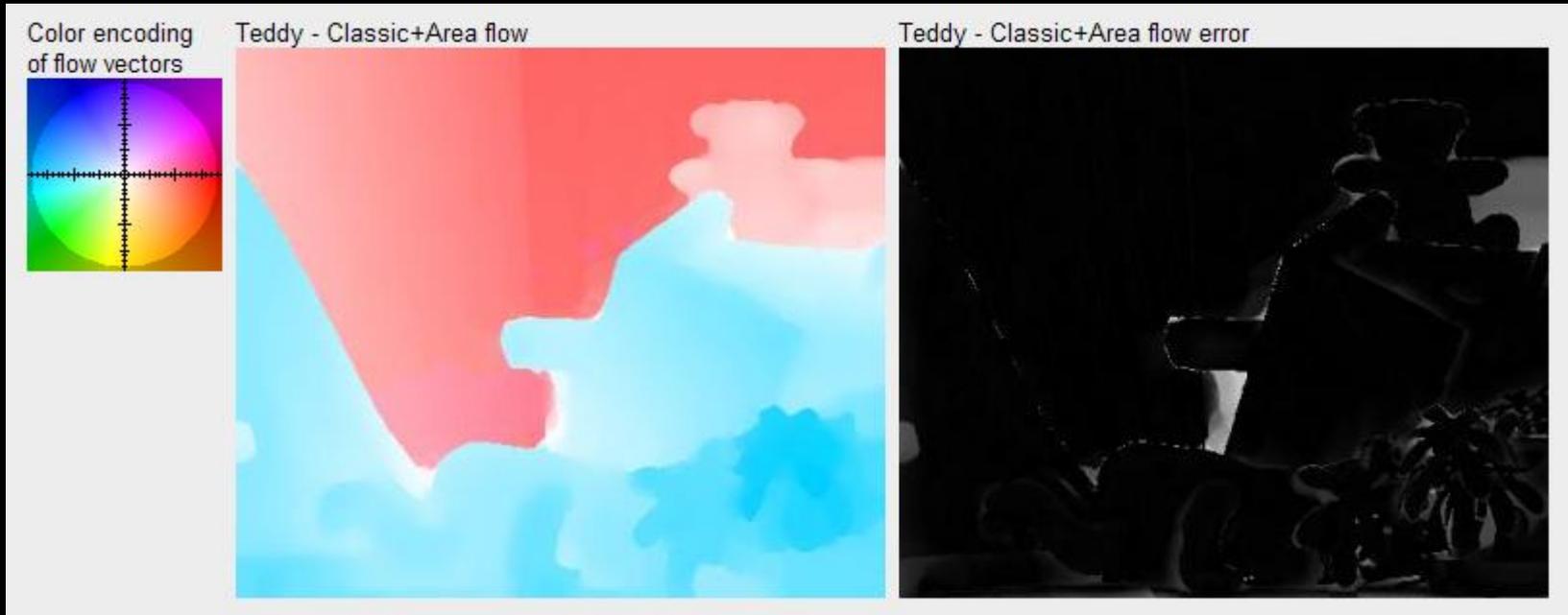
- #2 all-time Computer Vision problem (disputable)
- “Where did each pixel go?”





1st Best Algorithm (7th overall as of 17-12-2009)

DPOF [18]: C. Lei and Y.-H. Yang. [Optical flow estimation on coarse-to-fine region-trees using discrete optimization. ICCV 2009](#)



(3rd overall as of 17-12-2009) **2nd Best Algorithm**

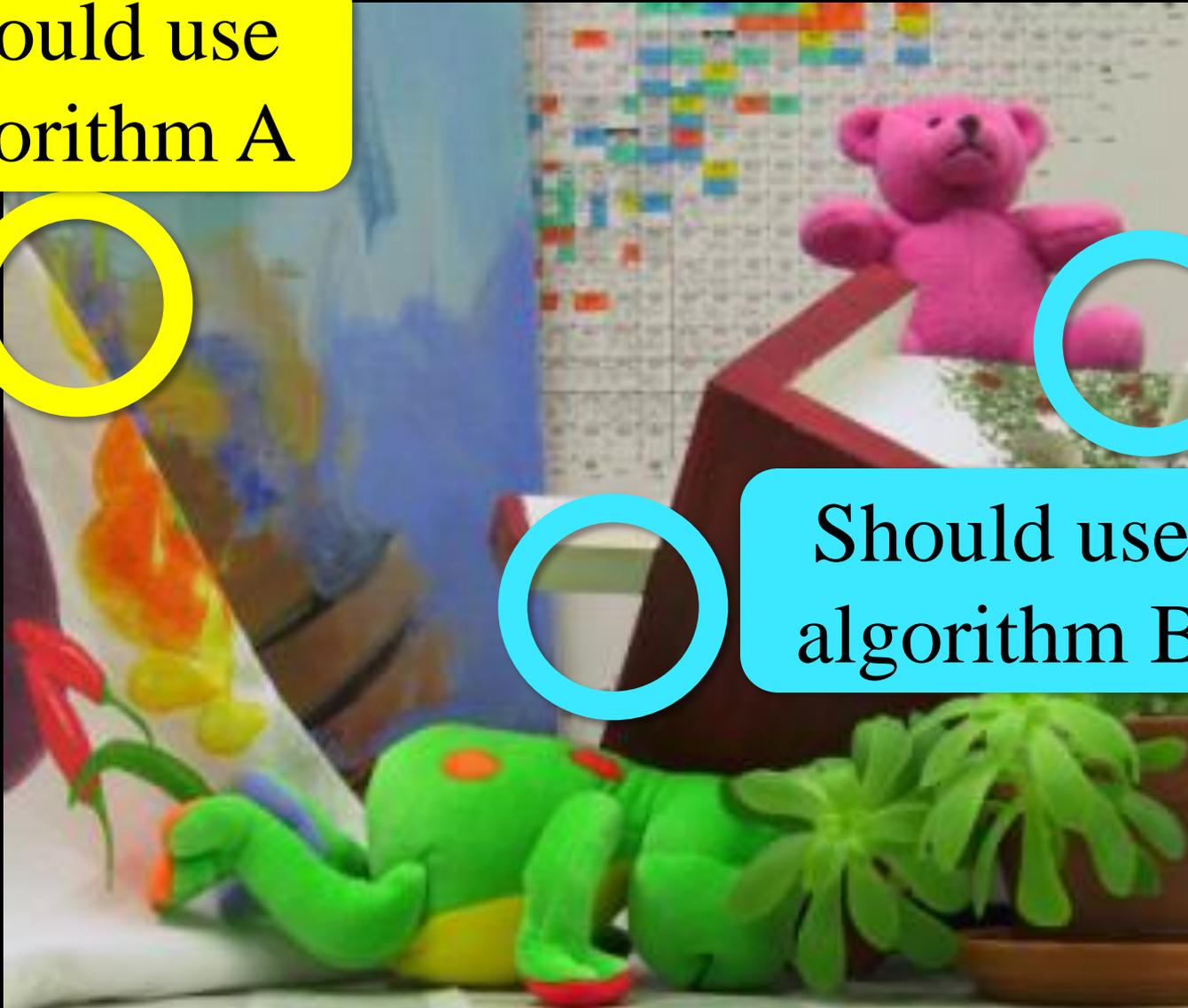
Classic+Area [31]: Anonymous. Secrets of optical flow estimation and their principles.
CVPR 2010 submission 477

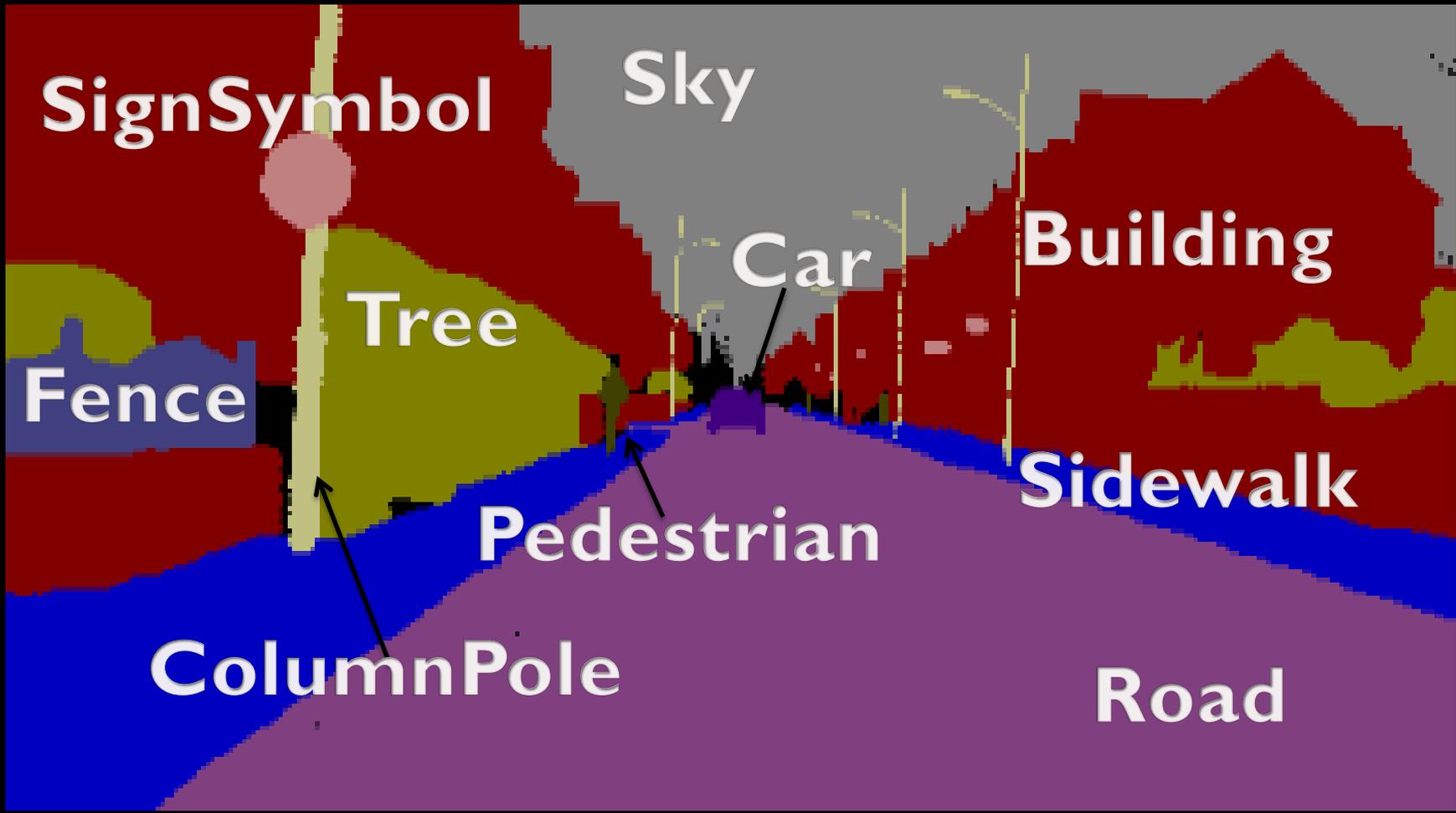


Should use
algorithm A



Should use
algorithm B





Sign

Symbol

Sky

Building

Car

Tree

Fence

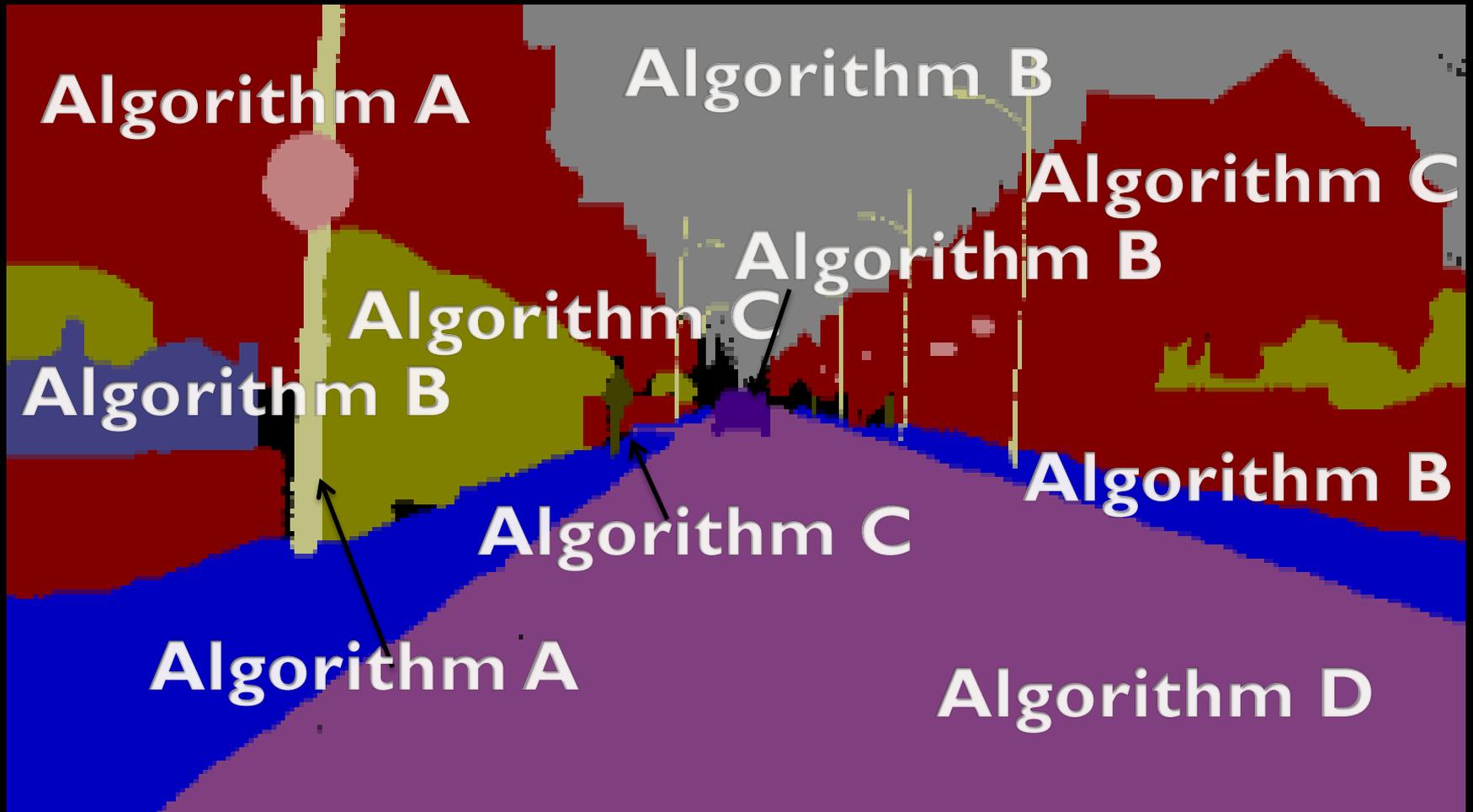
Pedestrian

Sidewalk

Column

Pole

Road



(Artistic version; object-boundaries don't interest us)

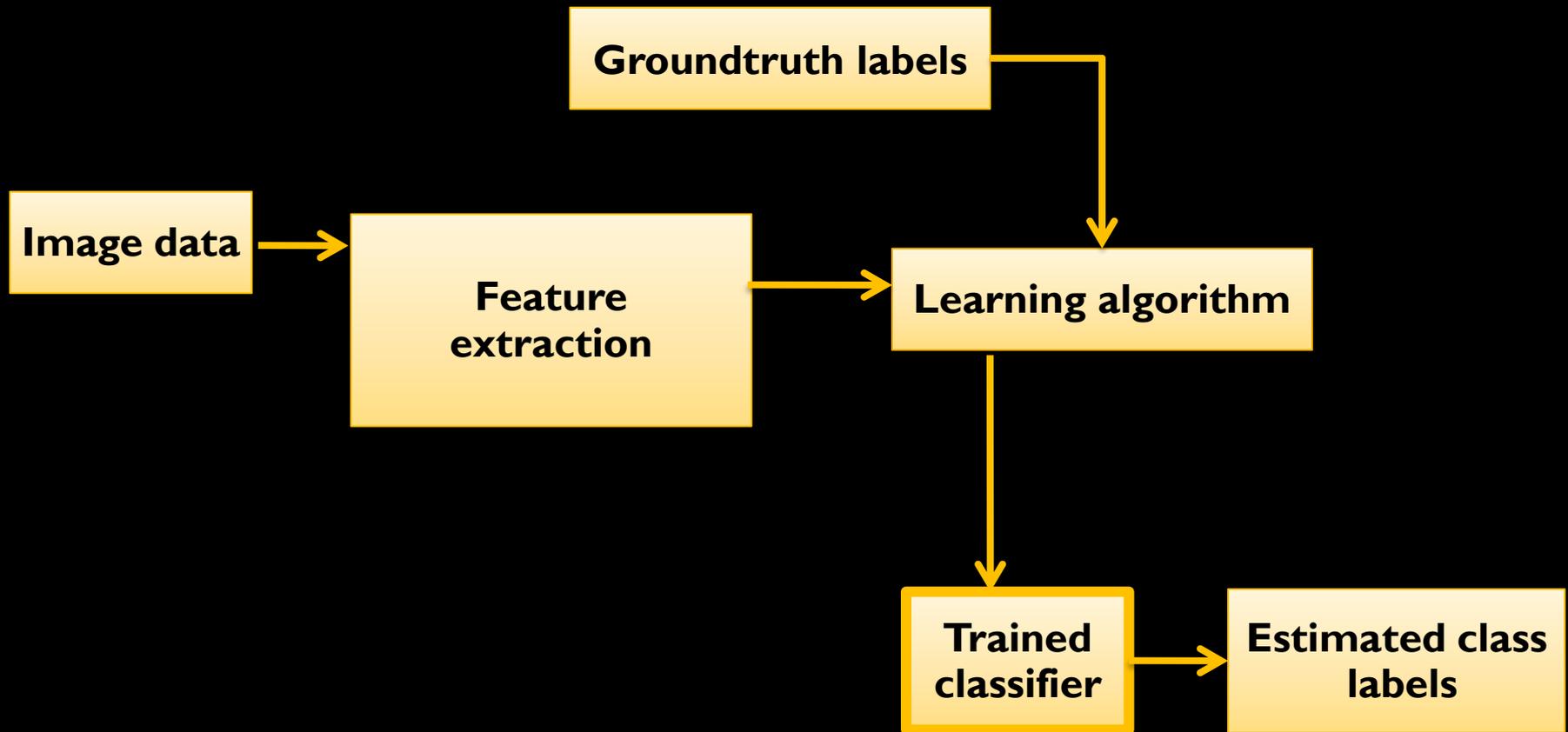
Hypothesis:

- that the most suitable algorithm can be chosen for each video automatically, through supervised training of a classifier

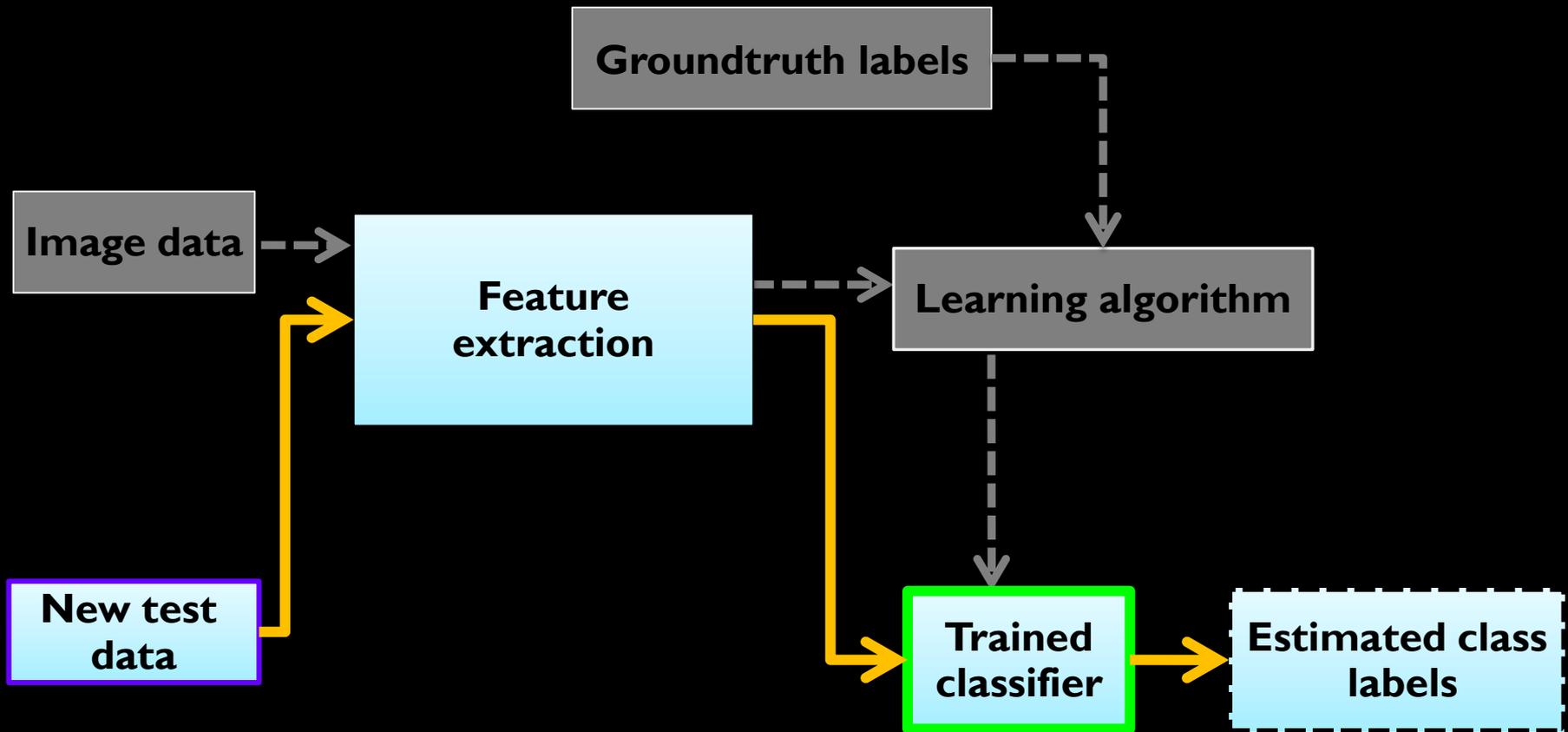
Hypothesis:

- ~~that the most suitable algorithm can be chosen for each video automatically, through supervised training of a classifier~~
- that one can predict the space-time segments of the video that are best-served by each available algorithm
 - (Can always come back to choose a per-frame or per-video algorithm)

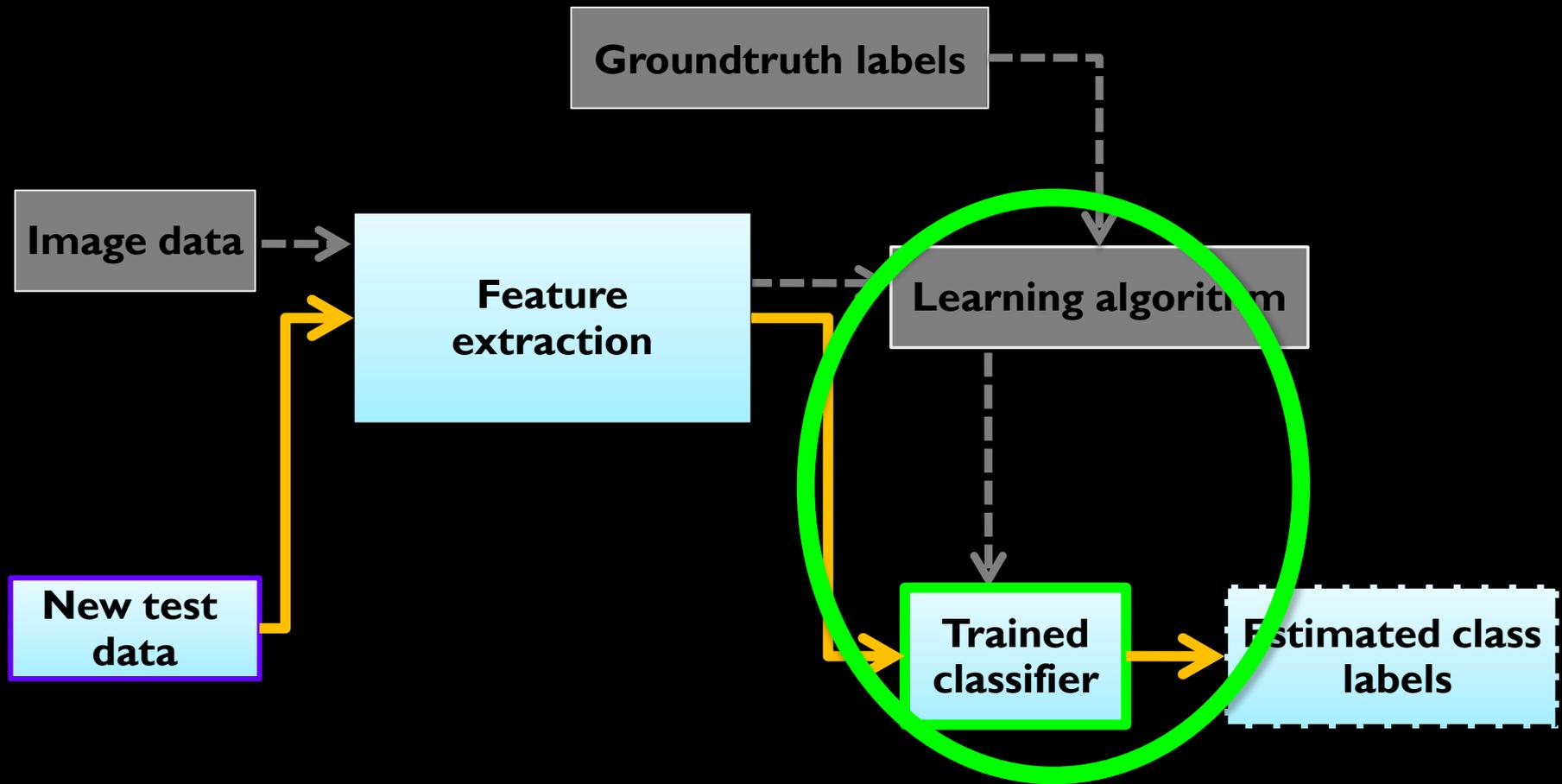
Experimental Framework



Experimental Framework



Experimental Framework

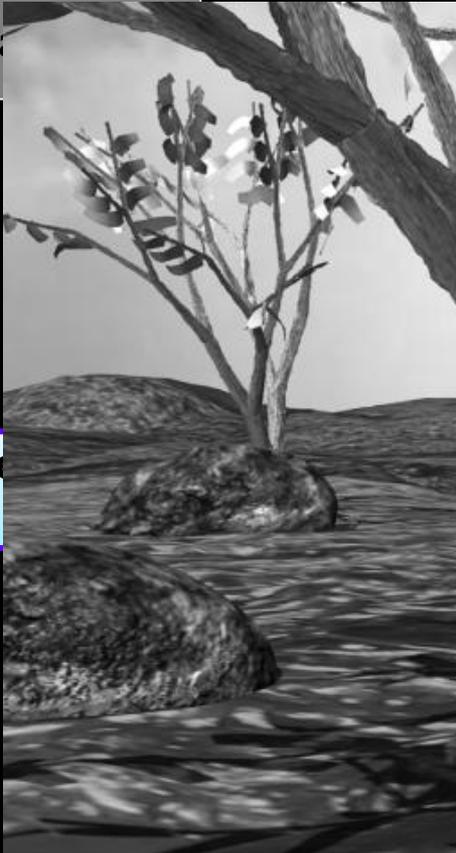


Random Forests
Breiman, 2001

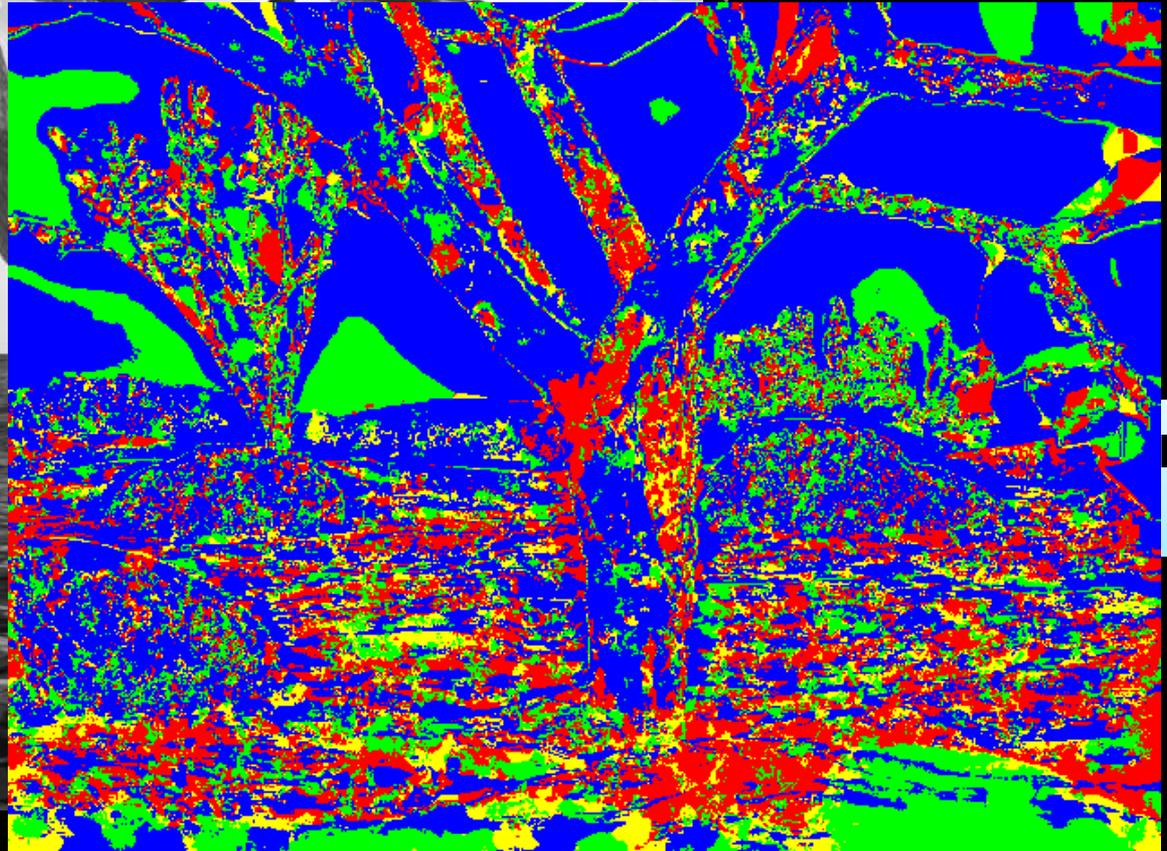
Experimental Framework

Groundtruth labels

Im

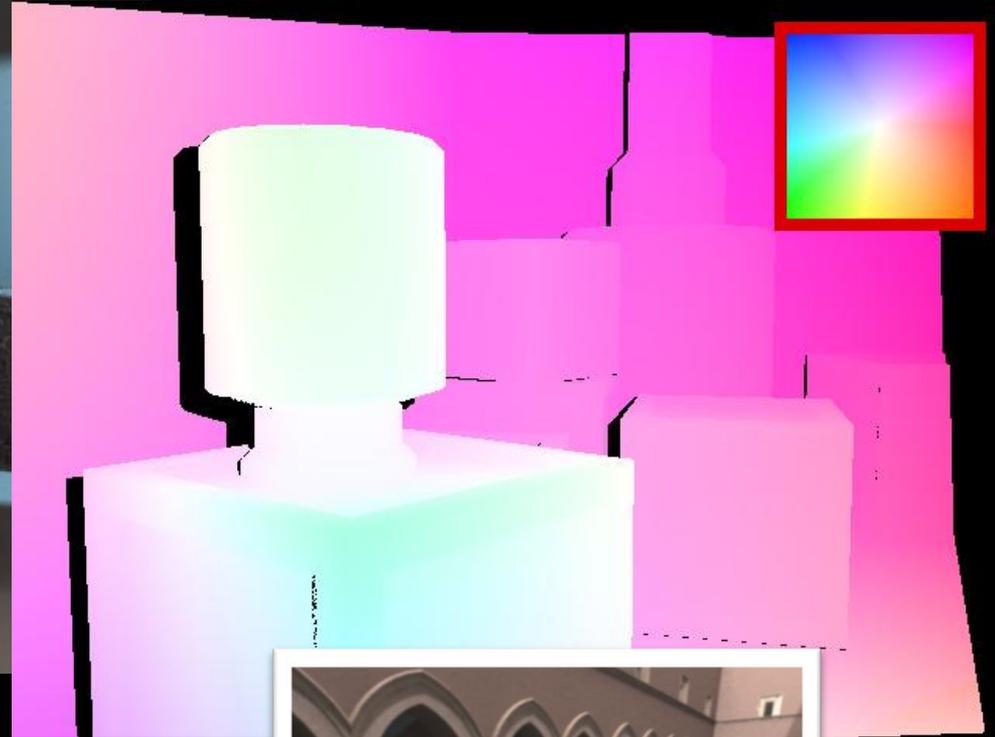
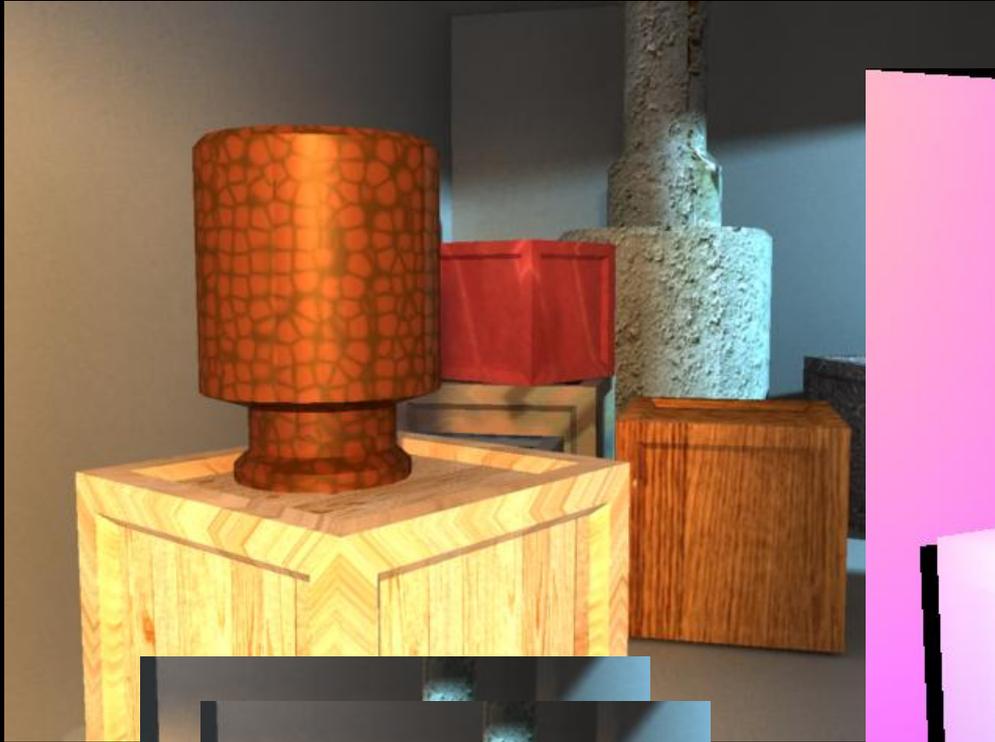


No



lass

“Making” more data



Formulation

$$\mathcal{D} = \{(\mathbf{x}_i, c_i) \mid \mathbf{x}_i \in \mathbb{R}^d, c_i \in \mathbb{Z}^k\}_{i=1}^n$$

$$\mathbf{x}_i = \{g(x, y, [1, z]), d(x, y, [1, z]), t_x(x, y, [1, z]), t_y(x, y, [1, z]), r(x, y, [1, k])\}$$

- Training data \mathcal{D} consists of feature vectors \mathbf{x} and class labels c (i.e. best-algorithm per pixel)
- Feature vector \mathbf{x} is multi-scale, and includes:
 - Spatial Gradient
 - Distance Transform
 - Temporal Gradient
 - Residual Error (after bicubic reconstruction)

Formulation

$$\mathcal{D} = \{(\mathbf{x}_i, c_i) \mid \mathbf{x}_i \in \mathbb{R}^d, c_i \in \mathbb{Z}^k\}_{i=1}^n$$

$$\mathbf{x}_i = \{g(x, y, [1, z]), d(x, y, [1, z]), t_x(x, y, [1, z]), t_y(x, y, [1, z]), r(x, y, [1, k])\}$$

- Training data \mathcal{D} consists of feature vectors \mathbf{x}_i and class labels c_i

$$g(x, y, z) = \|\nabla I_1\|$$

- Feature vector \mathbf{x} is multi-scale, and includes.
 - Spatial Gradient
 - Distance Transform
 - Temporal Gradient
 - Residual Error (after bicubic reconstruction)

Formulation

$$\mathcal{D} = \{(\mathbf{x}_i, c_i) \mid \mathbf{x}_i \in \mathbb{R}^d, c_i \in \mathbb{Z}^k\}_{i=1}^n$$

$$\mathbf{x}_i = \{g(x, y, [1, z]), d(x, y, [1, z]), t_x(x, y, [1, z]), t_y(x, y, [1, z]), r(x, y, [1, k])\}$$

$$d(x, y, z) = \text{disTrans}(\|\nabla I_1\| > \tau)$$

- Feature vector \mathbf{x} is multi-scale, and includes.
 - Spatial Gradient
 - Distance Transform
 - Temporal Gradient
 - Residual Error (after bicubic reconstruction)

Formulation Details

- Temporal Gradient

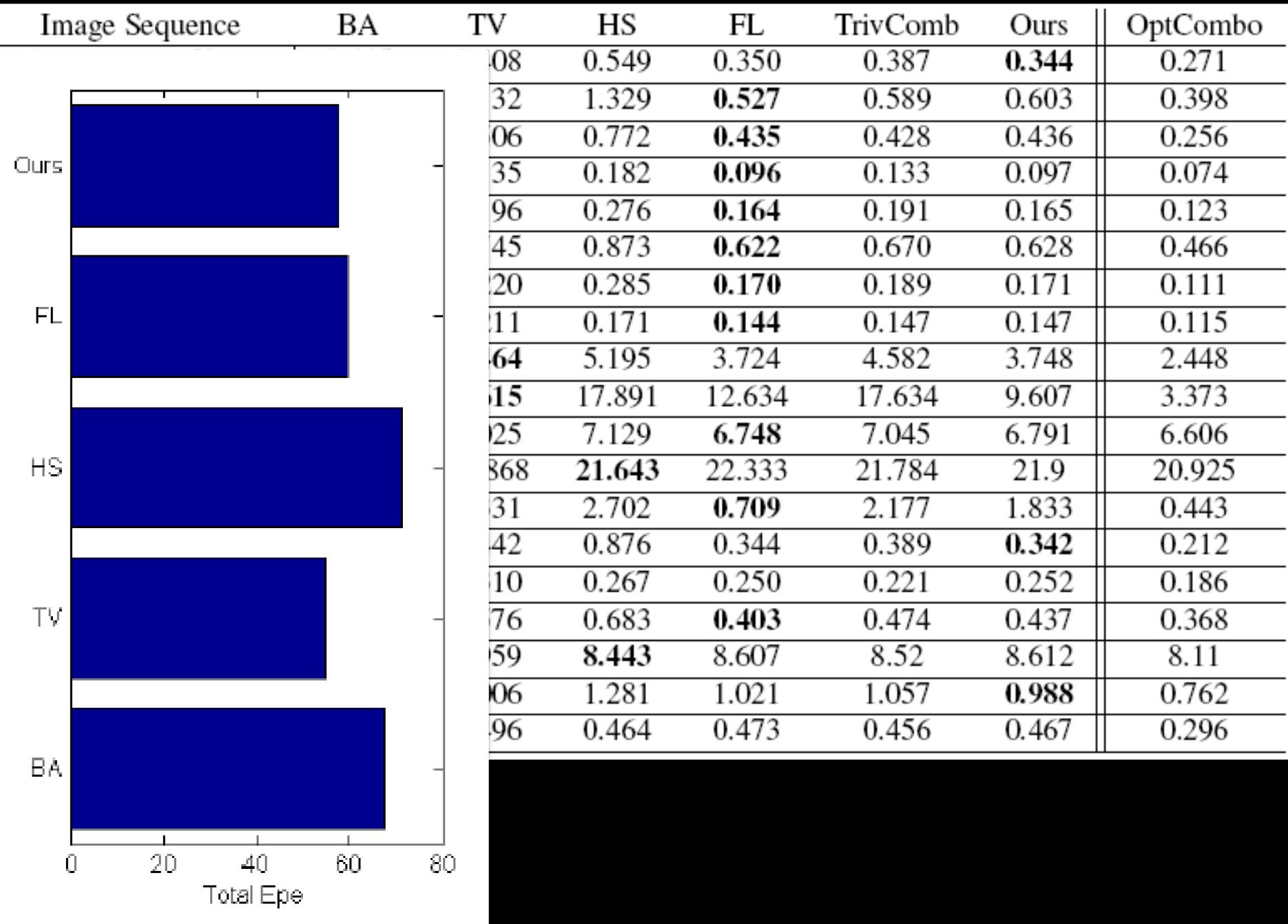
$$t_x = \|\nabla(x + \bar{u})\|$$

$$t_y = \|\nabla(y + \bar{v})\|$$

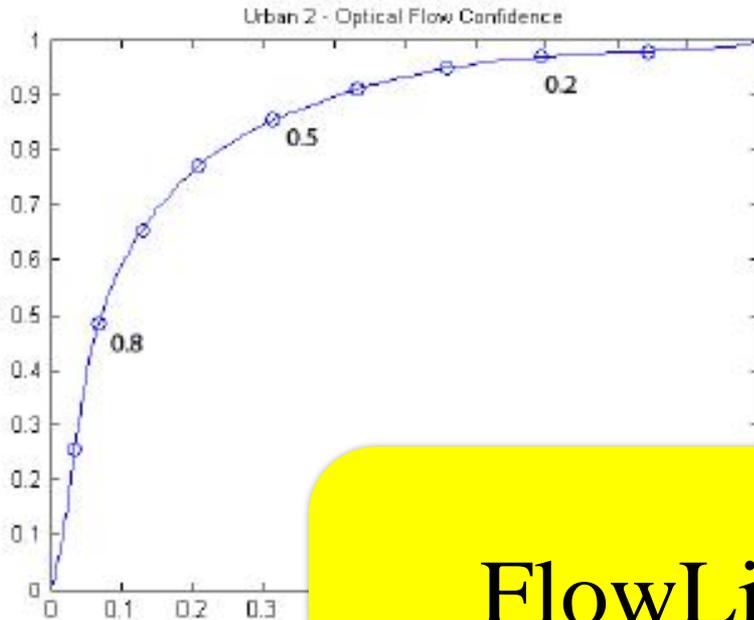
- Residual Error

$$r_i(x, y, k) = I_1(x, y) - \text{bicubic}(I_2(x + u_i(k), y + v_i(k)))$$

Application I: Optical Flow



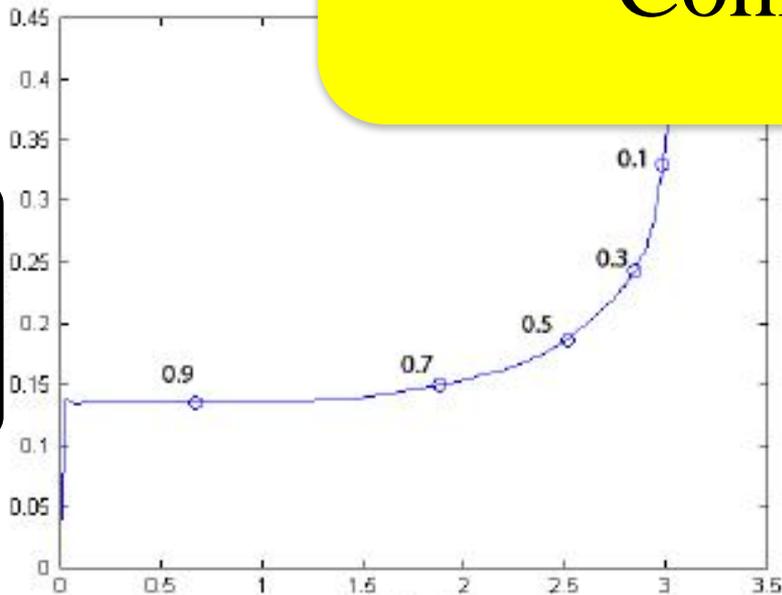
True Positive Rate



False Positive Rate

FlowLib Decision Confidence

Ave EPE



Number of pixels $\times 10^5$



Application II: Feature Matching

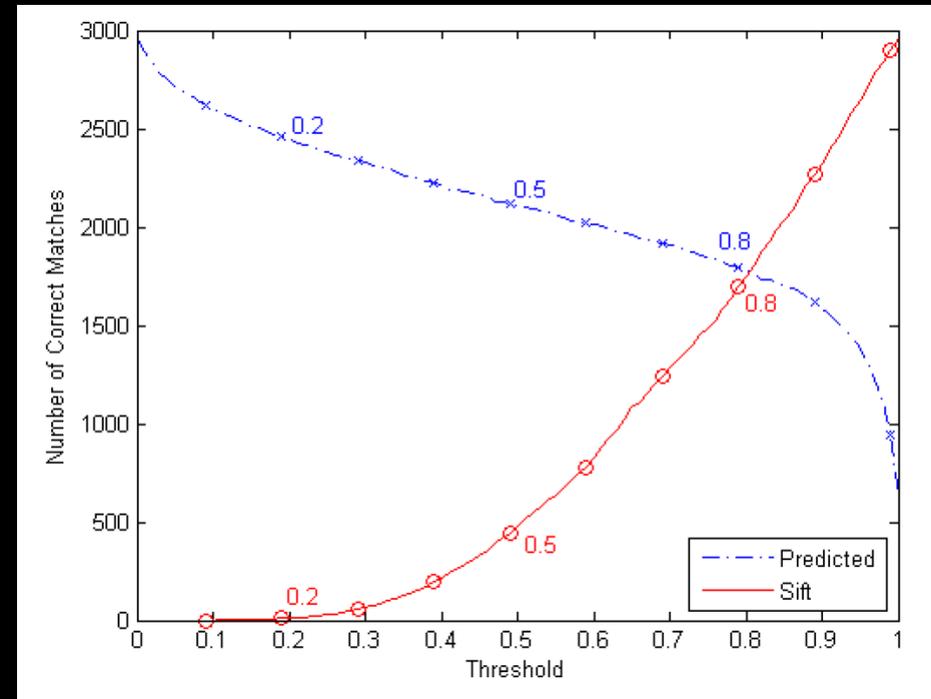
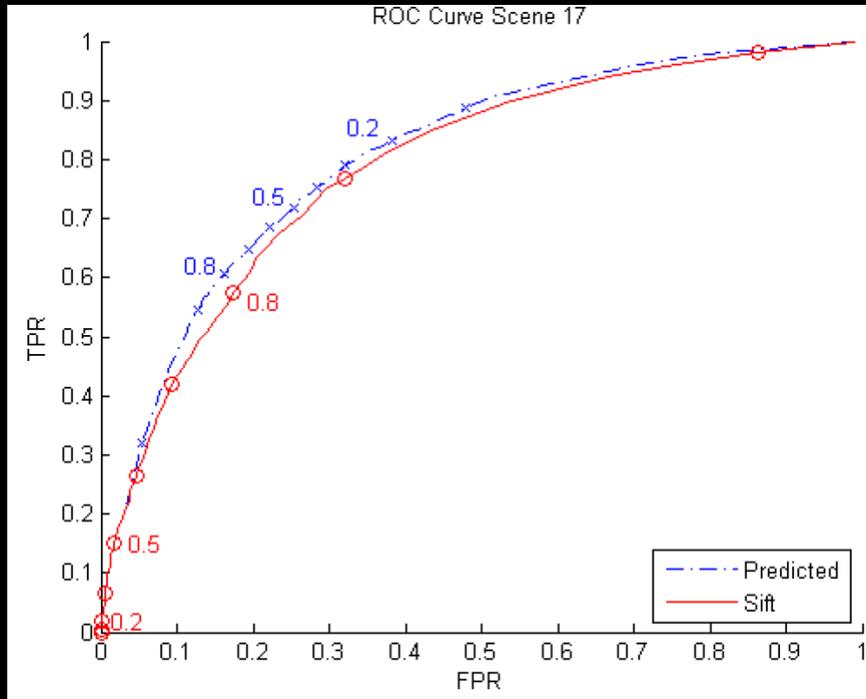
Comparing 2 Descriptions

- What is a match? Details are important...
 - Nearest neighbor (see also [FLANN](#))
 - Distance Ratio
 - PCA
- Evaluation: density, # correct matches, tolerance

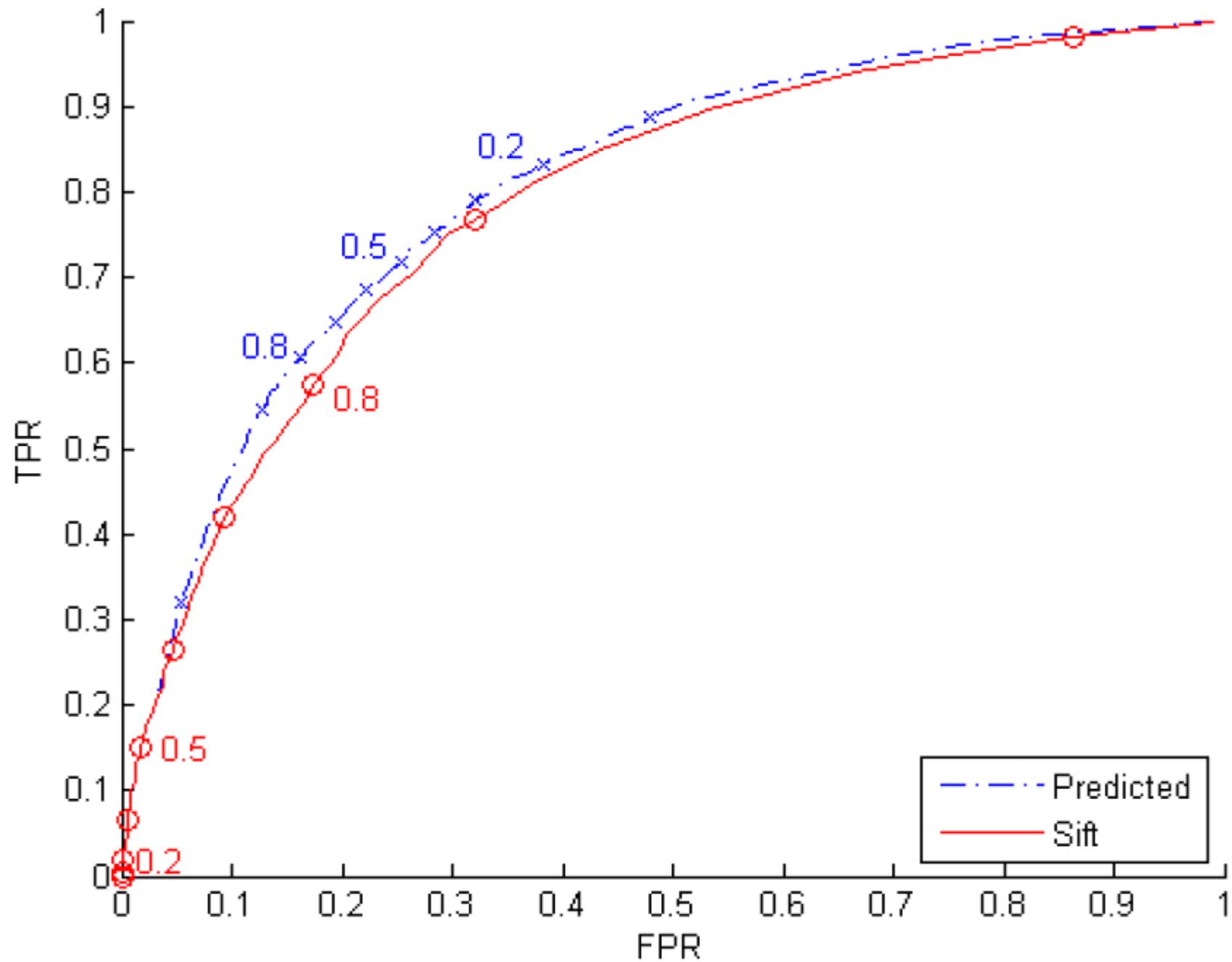


“192 correct matches (yellow) and 208 false matches (blue)”

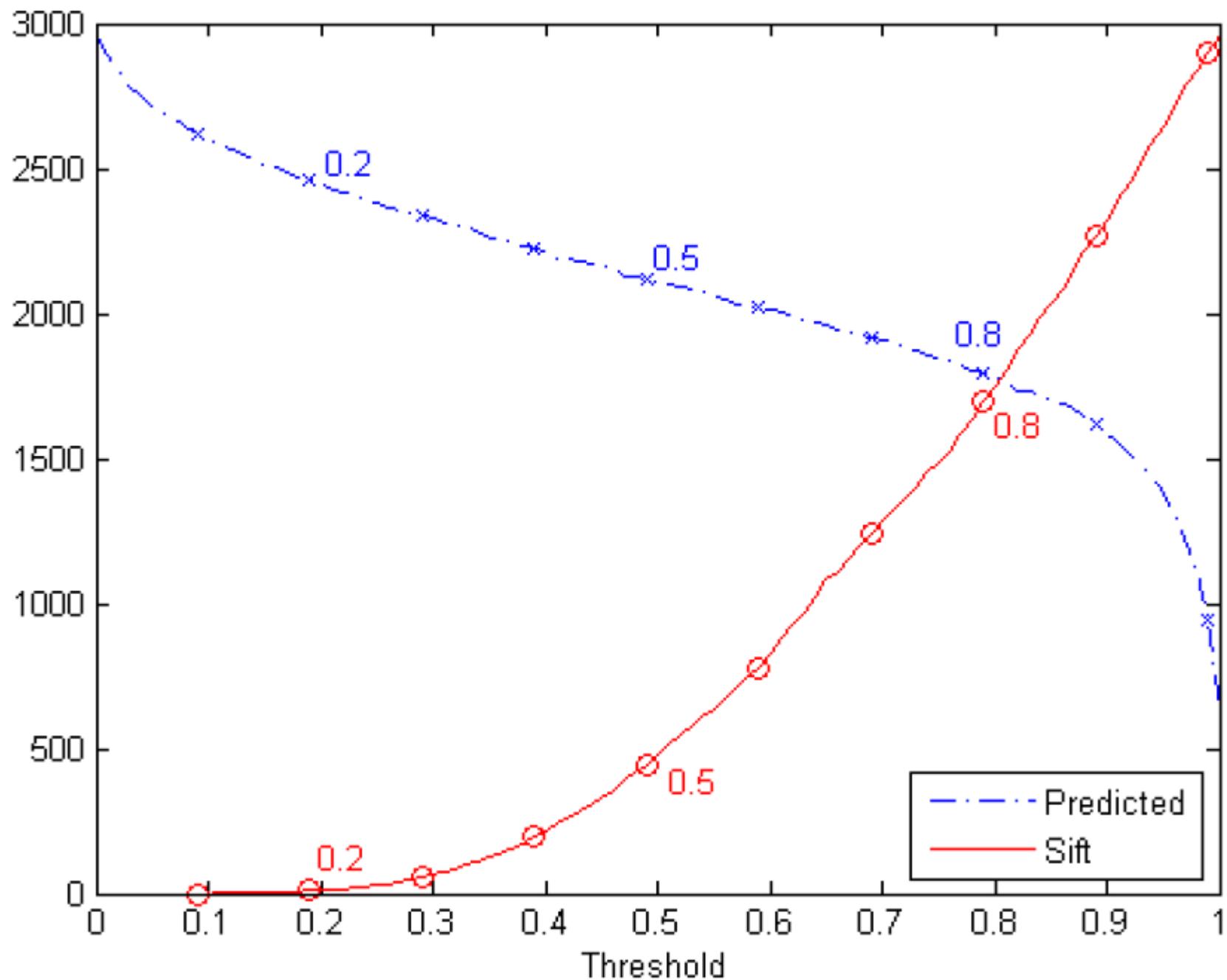
SIFT Decision Confidence



ROC Curve Scene 17



Number of Correct Matches



— · — Predicted
— Sift

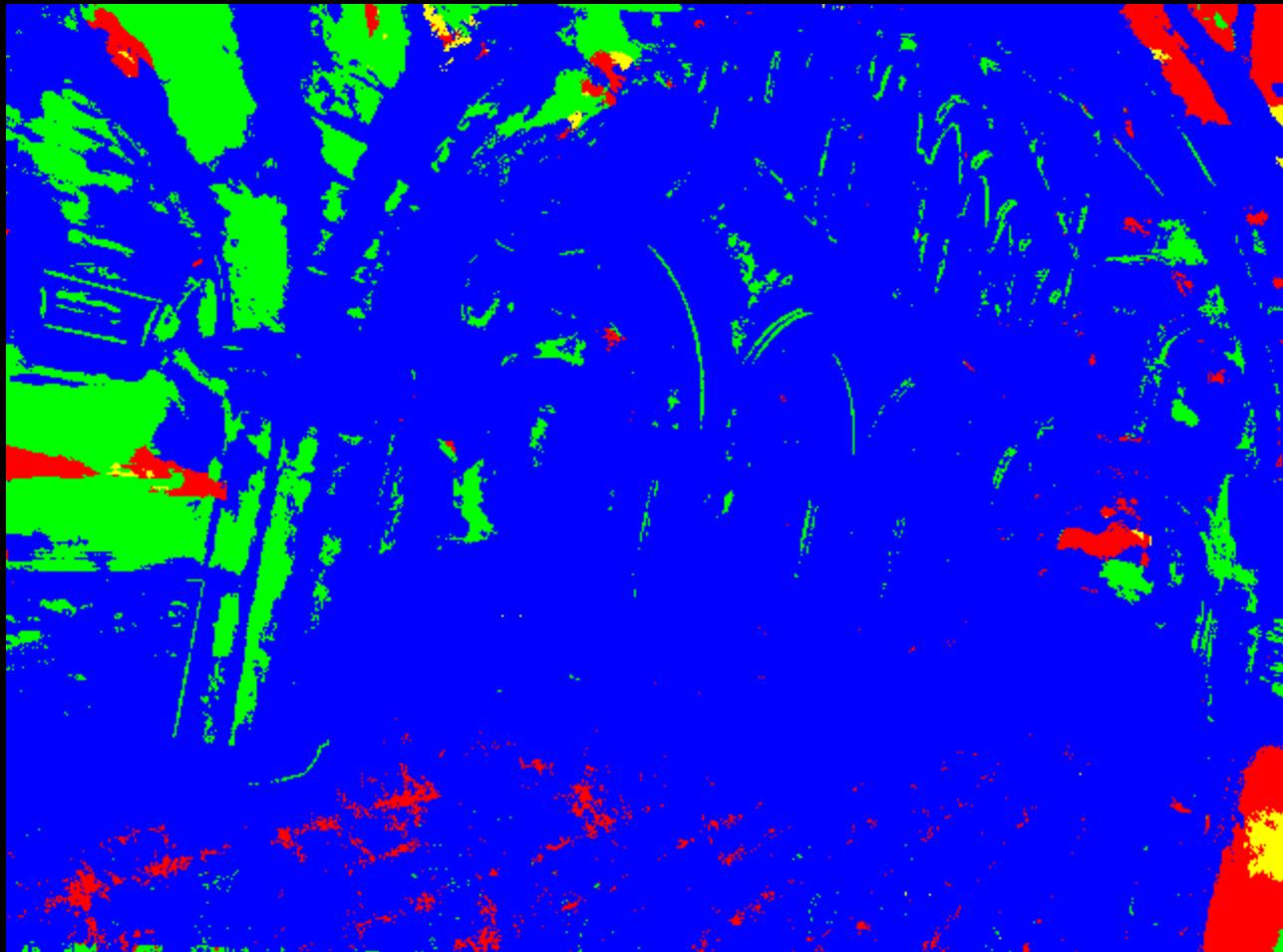
Hindsight / Future Work

- Current results don't quite live up to the theory:
 - Flaws of best-algorithm are the upper bound (ok)
 - Training data does not fit in memory (fixable)
 - “Winning” the race is more than rank (problem!)

Summary

- Overall, predictions are correlated with the best algorithm for each segment (**expressed as Pr!**)
- Training data where one class dominates is dangerous – needs improvement
- Other features could help make better predictions
 - Results don't yet do the idea justice
- One size does NOT fit all
 - At least in terms of algorithm suitability
 - Could use “bad” algorithms!





Ground Truth Best

Based on Prediction



FlowLib

Based on Prediction

White = 30 pixel end point error



FlowLib

Based on Prediction

(Contrast enhanced)