

Online, semi-supervised learning for long-term interaction with object recognition systems

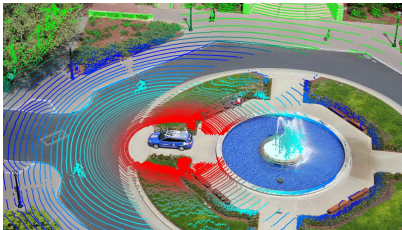
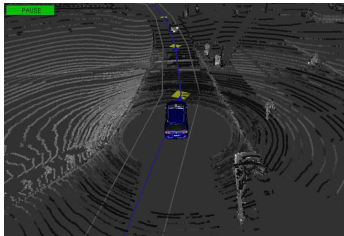
Alex Teichman and Sebastian Thrun

Department of Computer Science
Stanford University

July 12, 2012

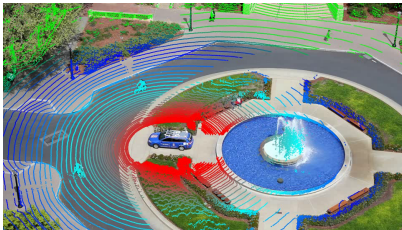
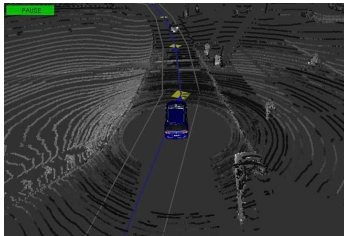


The big picture



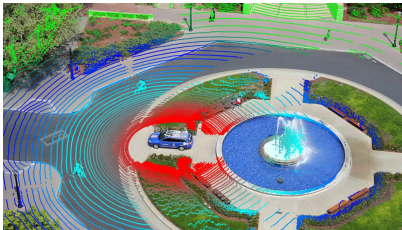
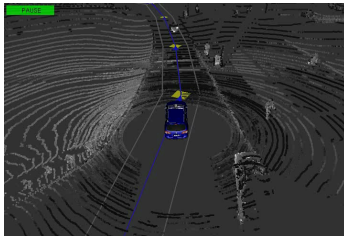
- What is the desired user interface for object recognition?
- Want autonomy with the option for user input.
- Online, active, semi-supervised learning ...

The big picture



- What is the desired user interface for object recognition?
- Want autonomy with the option for user input.
- Online, active, semi-supervised learning ...

The big picture



- What is the desired user interface for object recognition?
- Want autonomy with the option for user input.
- Online, active, semi-supervised learning ...

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Static train/test framework

Table 1. Breakdown of the dataset by class. Tracks were collected from busy streets and intersections.

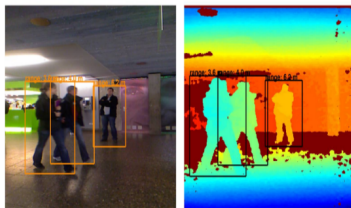
Number of tracks					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	904	205	187	6585	7881
Testing	847	112	140	4936	6035
Total	1751	317	327	11,521	13,916

Number of frames					
Set	Car	Pedestrian	Bicyclist	Background	All
Training	92,255	32,281	31,165	532,760	688,461
Testing	59,173	22,203	25,410	530,917	637,703
Total	151,428	54,484	56,575	1,063,677	1,326,164

- Rigorous evaluation and comparison
- Experimental setup
- Occasional user interaction
- Infinite unlabeled data stream

We don't want to overfit to this framework!

Object recognition approaches - sliding window & tracking-by-detection



Spinello and Arras

Spinello, Stachniss, and Burgard

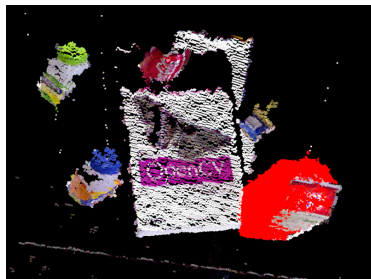
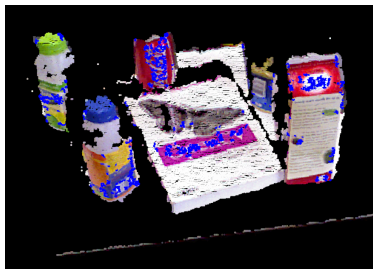
Object recognition approaches - semantic segmentation



Douillard et al.

Combining sliding windows and semantic segmentation: Lai et al.

Object recognition approaches - keypoint matching



Solutions in Perception Challenge

Collet et al.

Problem decomposition

Segmentation — Tracking — Track classification



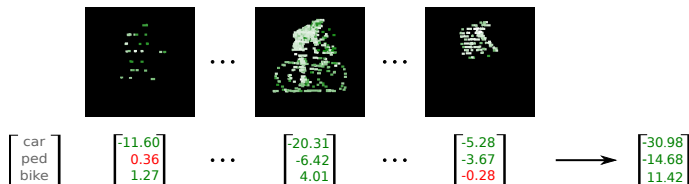
Problem decomposition

Segmentation — Tracking — Track classification



Problem decomposition

Segmentation — Tracking — Track classification



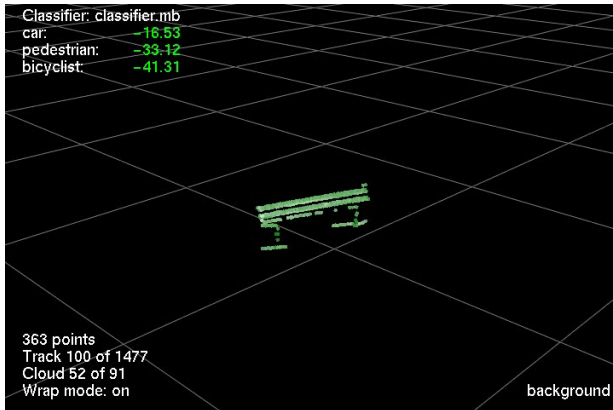
Descriptors



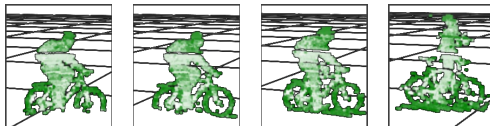
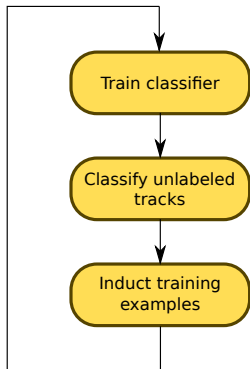
- 29 different descriptor spaces
- $x \in \mathbb{R}^{\sim 4000}$

- Oriented bounding box size
- Spin images
- HOG descriptors computed on virtual orthographic camera images

Tracks

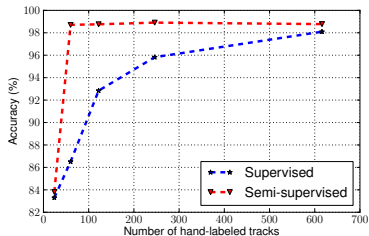


Tracking-based semi-supervised learning



- Large, automatically-labeled background dataset is provided. Often this is easy to collect.
- Only positive examples are inducted during semi-supervised learning.

Tracking-based semi-supervised learning



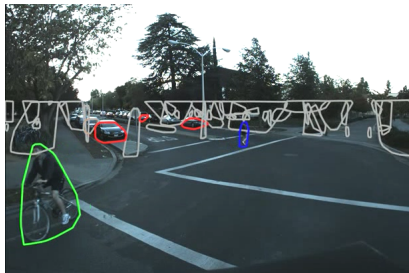
Labels

car	792	0	0	17
pedestrian	0	86	0	0
bicyclist	0	4	138	2
background	55	22	2	4917
	car	pedestrian	bicyclist	background

Predictions

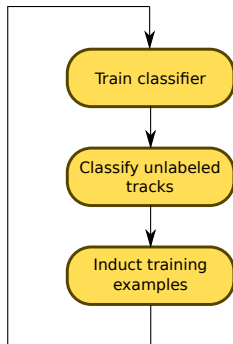
- Unsupervised method given millions of additional unlabeled examples.
- Track classification accuracy is reported. (This does not include segmentation and tracking errors.)

Tracking-based semi-supervised learning



- Three hand-labeled training examples of each class + millions of unlabeled examples used to generate these results.
- Outlines are tracked objects. Track classifications are computed offline.
- White outlines are tracked objects classified as neither person, bicyclist, or car.

Offline to online



Algorithm 1 Tracking-based semi-supervised learning

τ is a confidence threshold chosen by hand
 \mathcal{S} is a small set of seed tracks, labeled by hand
 \mathcal{U} is a large set of unlabeled tracks
 \mathcal{B} is a large set of background tracks
 \mathcal{W} is a working set, initially empty

```
 $\mathcal{W} := \mathcal{S} \cup \mathcal{B}$   
repeat  
  Train frame classifier  $\mathcal{C}$  on  $\mathcal{W}$   
   $\mathcal{W} := \mathcal{S} \cup \mathcal{B}$   
  for  $u \in \mathcal{U}$  do  
    Classify track  $u$  using  $\mathcal{C}$   
     $c := \text{confidence}(u)$   
     $l := \text{classification}(u)$   
    if  $c \geq \tau$  and  $l \neq \text{"background"}$  then  
      Add  $u$  to  $\mathcal{W}$  with label  $l$   
    end if  
  end for  
until converged
```

Modularity

Segmentation

- Connected components
- Background subtraction

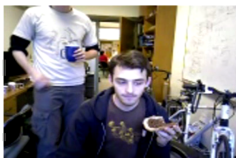
Tracking

- Kalman filters

Classification

- Boosting
- Logistic regression,
stochastic gradient descent

- Discriminative segmentation
and tracking



WAFR2012

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Logistic regression & stochastic gradient descent

- Parametric
- Fast to train and evaluate
- Easy to incrementally train

$$x \in \mathbb{R}^n, y \in \{-1, +1\}$$

$$\mathbb{P}(y|x) = \frac{1}{1 + \exp(-yw^T x)}$$

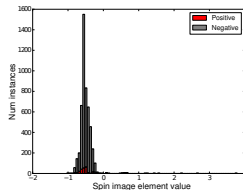
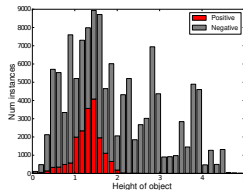
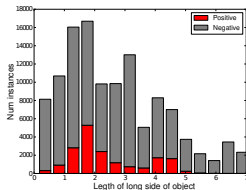
$$\underset{w}{\text{maximize}} \quad \prod_m \mathbb{P}(y^{(m)}|x^{(m)})$$

$$\underset{w}{\text{minimize}} \quad \sum_{m=1}^M \log(1 + \exp(-y^{(m)} w^T x^{(m)}))$$

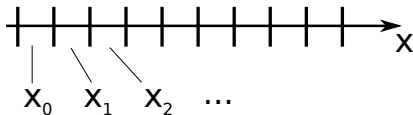
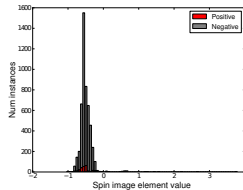
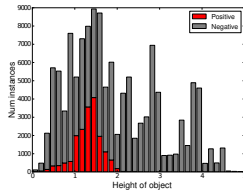
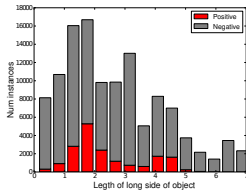
- M might be giant, or you might not have access to them all at one time.
- Stochastic gradient descent: take gradient steps using just small subsets of the data.
- ... but this fails badly if applied without thinking.

Linear models

$$\log \frac{\mathbb{P}(Y = 1|x)}{\mathbb{P}(Y = -1|x)} \approx w^T x$$

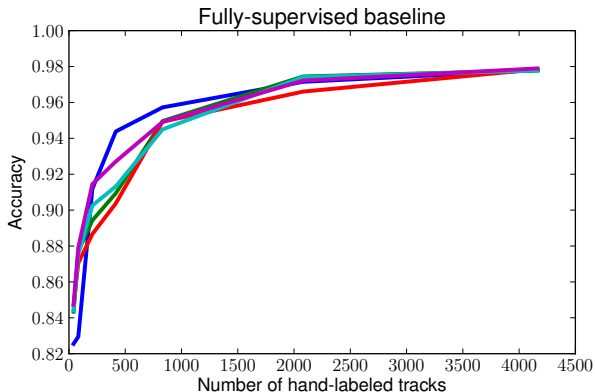


Feature transform



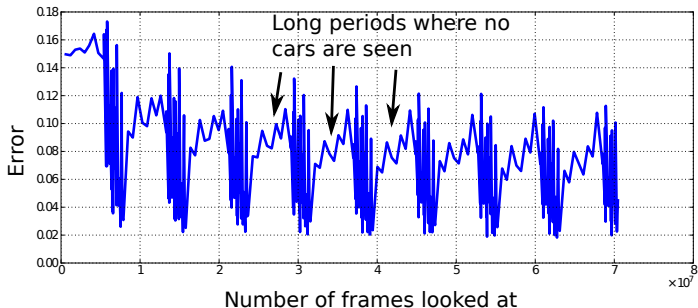
$$\mathbb{R}^{4000} \rightarrow \{0, 1\}^{400000}$$

Supervised performance



Linear model reaches a maximum of 94.0%, fully-supervised boosting 98.7%.

Prediction stability



- Fully-supervised, looping through $\sim 7\text{M}$ training examples.
- Can't do semi-supervised learning if you forget about objects after not seeing them for a while!

Training buffers



- D_S is the stream of examples seen so far.
- D_C is a new chunk of data.
- Want to maintain D_B , a fixed-size buffer of examples which is representative of D_S .
- Resample from D_B and D_C proportionally, relative to how much of the total stream they represent.

Training buffers



- D_S is the stream of examples seen so far.
- D_C is a new chunk of data.
- Want to maintain D_B , a fixed-size buffer of examples which is representative of D_S .
- Resample from D_B and D_C proportionally, relative to how much of the total stream they represent.

Training buffers



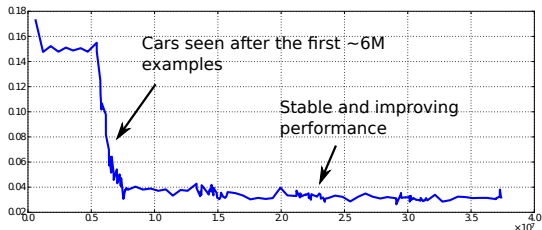
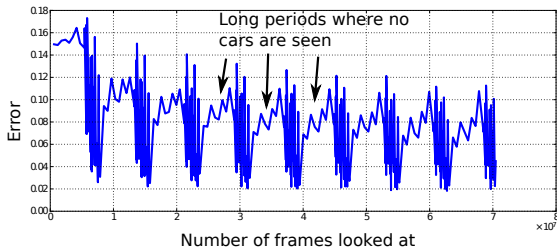
- D_S is the stream of examples seen so far.
- D_C is a new chunk of data.
- Want to maintain D_B , a fixed-size buffer of examples which is representative of D_S .
- Resample from D_B and D_C proportionally, relative to how much of the total stream they represent.

Training buffers

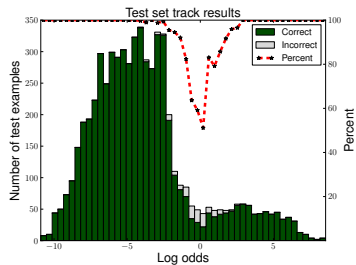


- D_S is the stream of examples seen so far.
- D_C is a new chunk of data.
- Want to maintain D_B , a fixed-size buffer of examples which is representative of D_S .
- Resample from D_B and D_C proportionally, relative to how much of the total stream they represent.

Training buffers

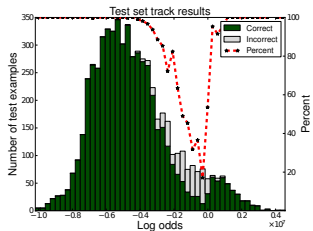
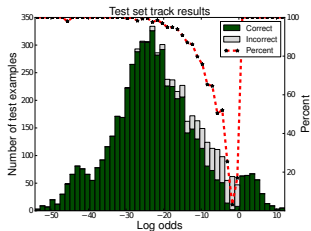


Confidence thresholds

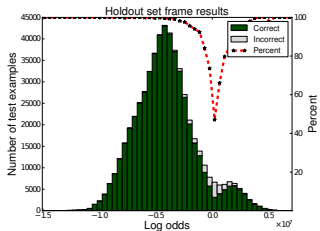
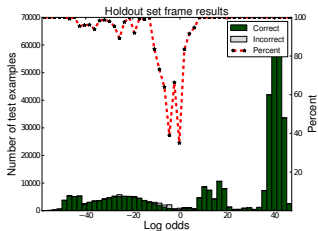
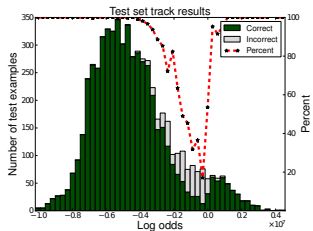
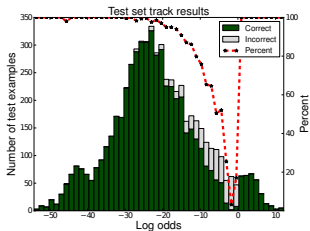


- Need to decide when to induct new tracks as positive examples of objects.

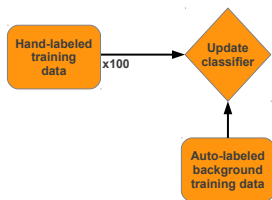
Variable confidences



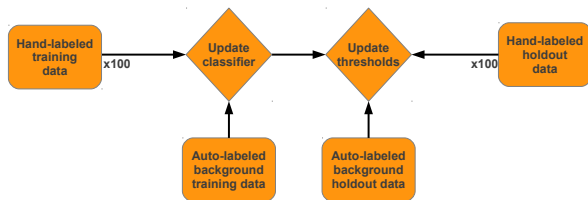
Confidence threshold learning



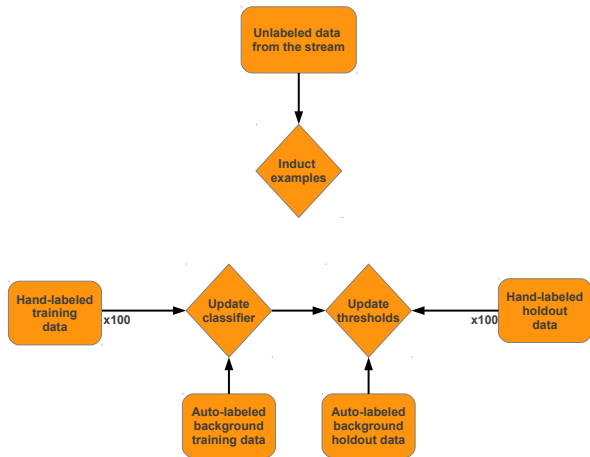
Algorithm sketch



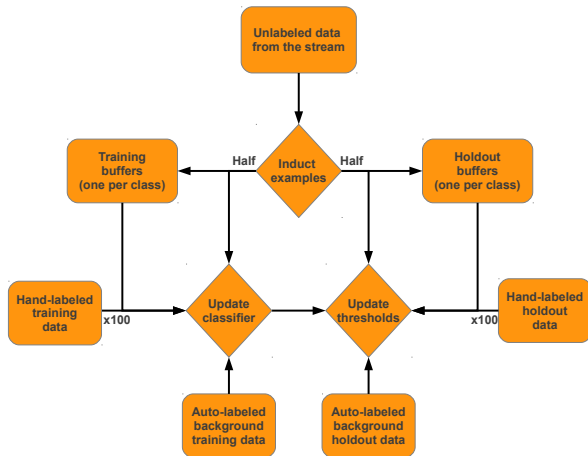
Algorithm sketch



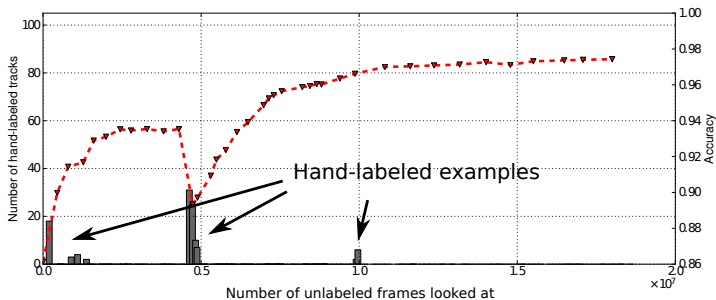
Algorithm sketch



Algorithm sketch

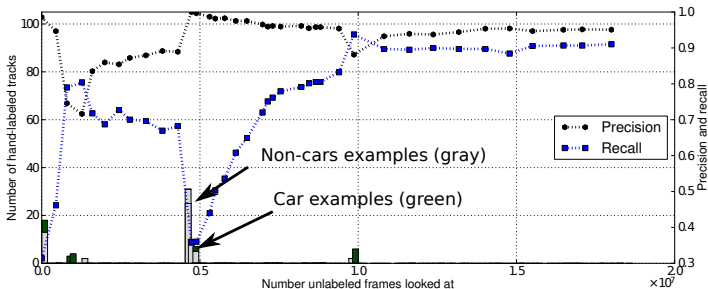


Online tracking-based semi-supervised learning



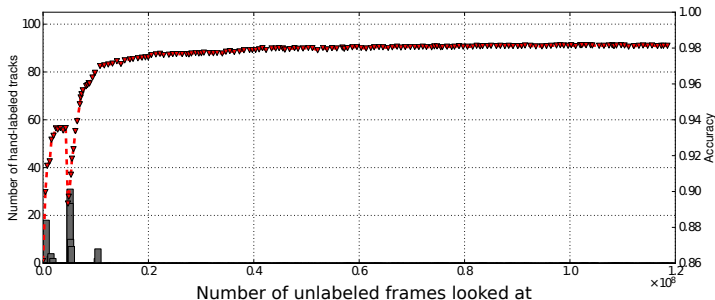
- Additional hand-labeled examples can break it out of local minima.
- $\sim 8M$ unique unlabeled examples.

Online tracking-based semi-supervised learning



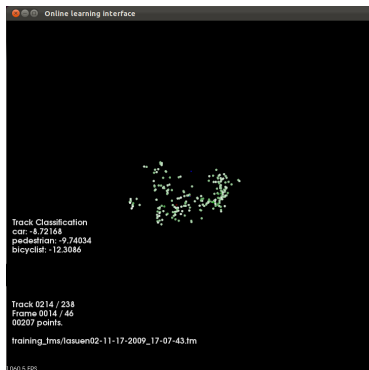
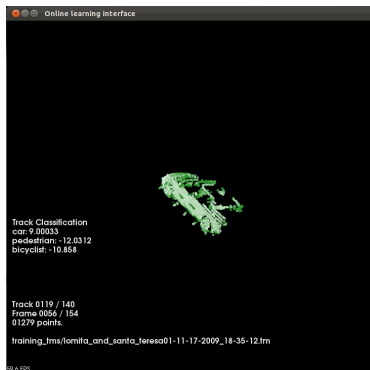
- Given lots of negative examples, recall initially drops, then recovers; overall accuracy improves.

Online tracking-based semi-supervised learning

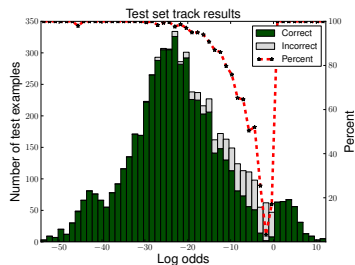
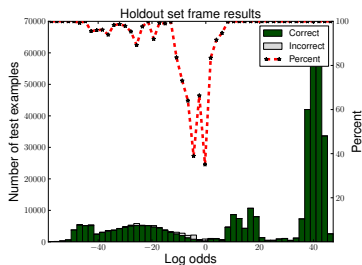


- Results after running for ~ 1 week. Total hand-labeled tracks: 108, vs ~ 4000 needed for good performance in fully-supervised case.
- Max accuracy when training on automatically-labeled background and all hand-labeled tracks: 90.1%.

Annotating

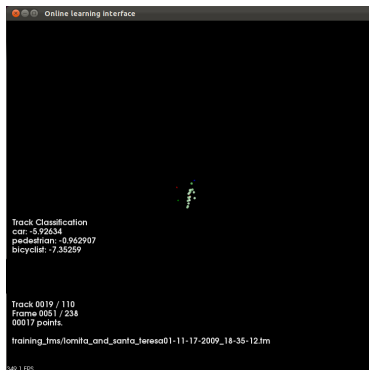
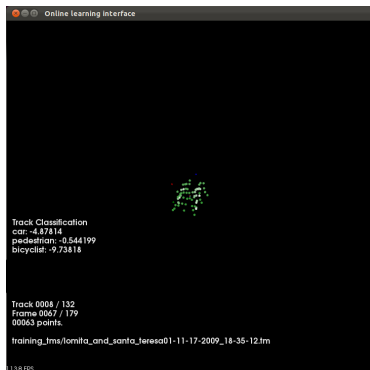


Annotating

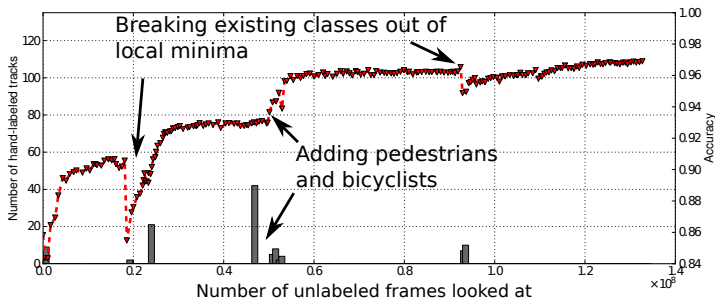


- The holdout set can tell you where to look for incorrect examples.

Annotating

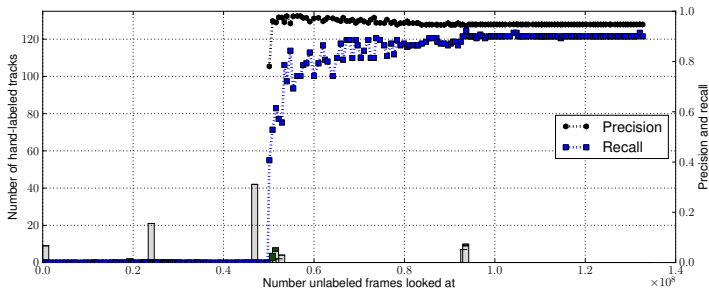


Adding classes later



- Max accuracy when training on automatically-labeled background and all hand-labeled tracks: 90.8%.

Adding classes later



Causes of failure while developing this

- Memory fragmentation
- Combined training buffer rather than one per class
- Stochastic gradient constant step size
- Not weighting the hand-labeled data

Causes of failure while developing this

- Memory fragmentation
- Combined training buffer rather than one per class
- Stochastic gradient constant step size
- Not weighting the hand-labeled data

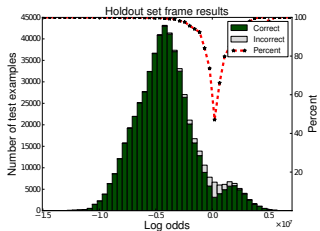
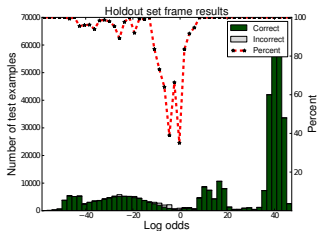
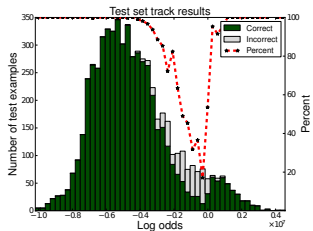
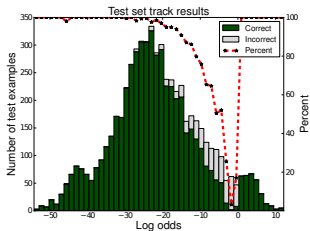
Causes of failure while developing this

- Memory fragmentation
- Combined training buffer rather than one per class
- Stochastic gradient constant step size
- Not weighting the hand-labeled data

Causes of failure while developing this

- Memory fragmentation
- Combined training buffer rather than one per class
- Stochastic gradient constant step size
- Not weighting the hand-labeled data

Future work: dual induction



Future work

Segmentation

- Connected components
- Background subtraction

Tracking

- Kalman filters

Classification

- Boosting
- Logistic regression, stochastic gradient descent

- Discriminative segmentation and tracking

