# Struck: Structured Output Tracking with Kernels

Sam Hare, Amir Saffari, And Philip H. S. Torr
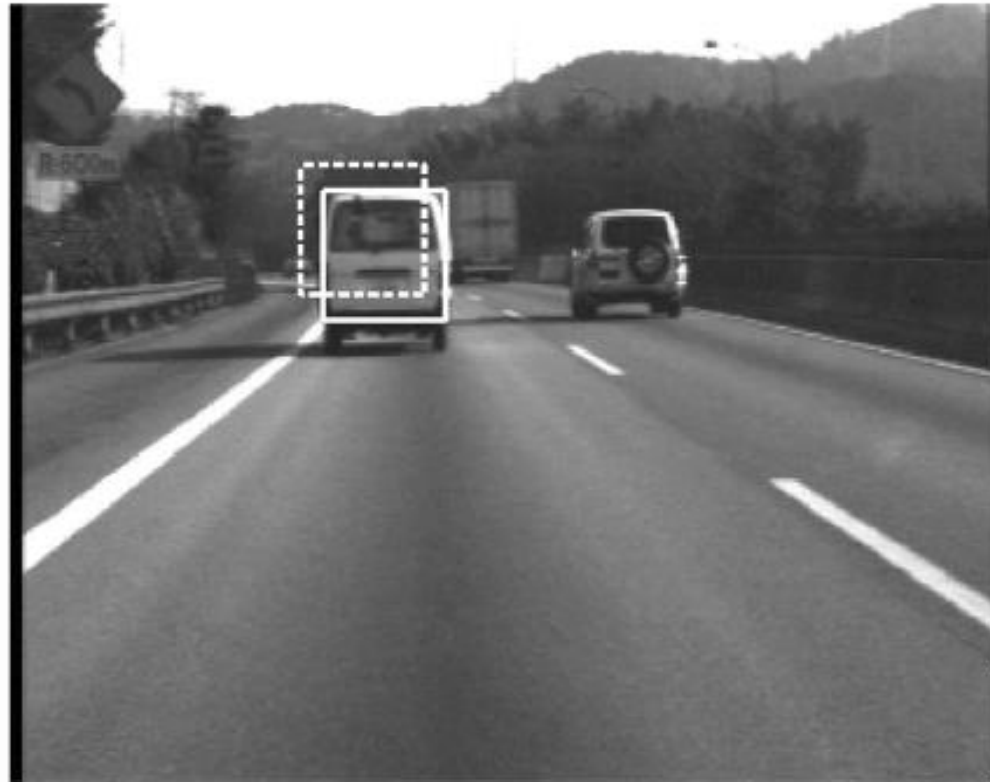
# Motivations

Problem: tracking-by-detection

Input:    target

Output: locations over times
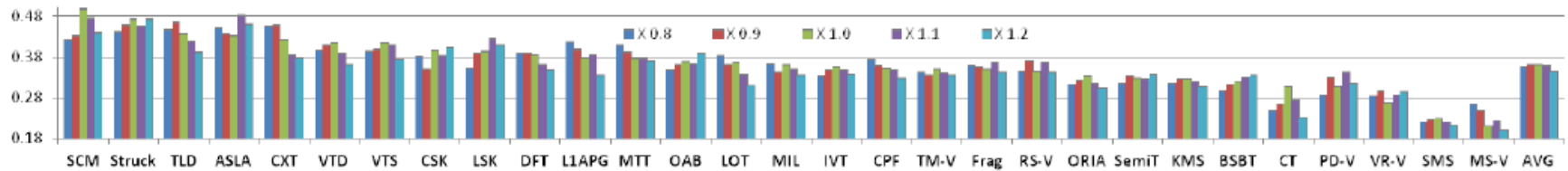
# Performance summary

Struck

TLD

MIL



Figure 6. Performance summary for the trackers initialized with different size of bounding box. AVG (the last one) illustrates the average performance over all trackers for each scale.

Y Wu, J Lim, MH Yang "Online Object Tracking: A Benchmark", Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on

# Outline

Previous works

- Tracking-by-detection
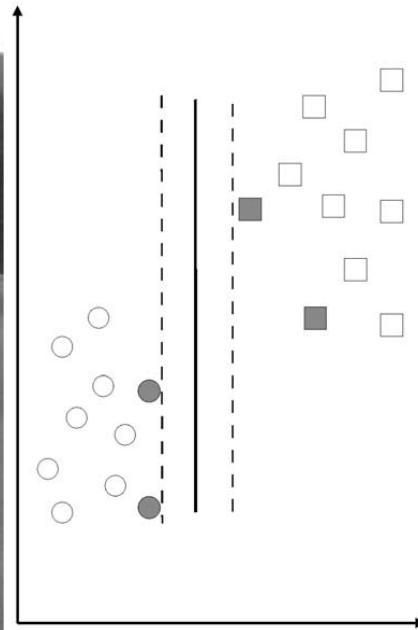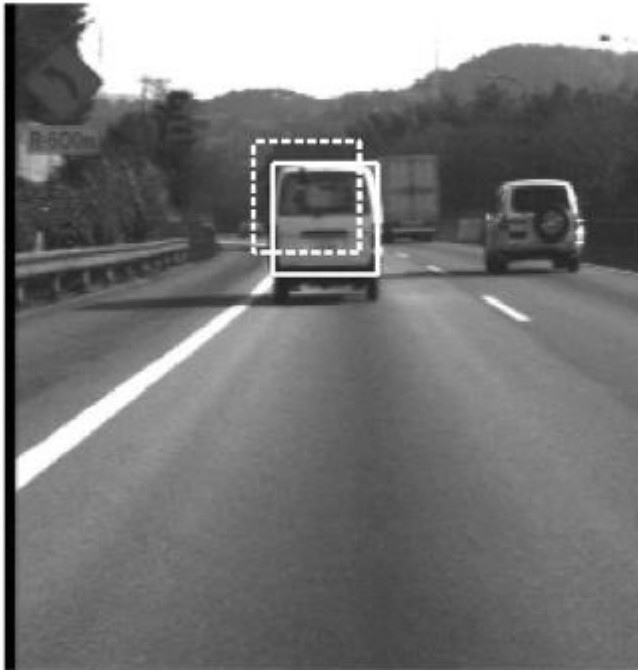- Adaptive tracking-by-detection

Methods

- Structured output tracking
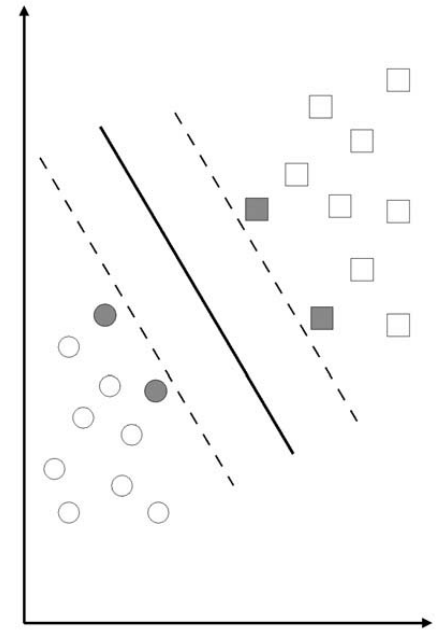- Online optimization and budget mechanism

Experiments and results

# Previous Works

Tracking problem as a detection task applied over time
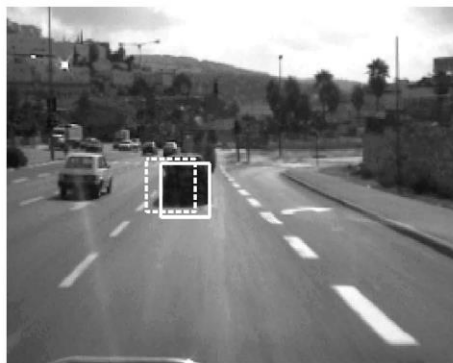


(a)                          (b)

Separating hyperplanes with different margins.

S. Avidan. Support Vector Tracking. IEEE Trans. on PAMI, 26:1064–1072, 2004.

# Previous Works

Tracking problem as a detection task applied over time



look for the image region with the highest SVM score

S. Avidan. Support Vector Tracking. IEEE Trans. on PAMI, 26:1064–1072, 2004.

# Previous Works

## Adaptive tracking-by-detection



B. Babenko, M. H. Yang, and S. Belongie. Visual Tracking with Online Multiple Instance Learning. In Proc. CVPR, 2009.

# Previous Works – Adaptive Tracking-by-detection

# Previous Works – Adaptive Tracking-by-detection

Adaptive tracking-by-detection

Tracking: A classification task

Learning: A update the object model.



B. Babenko, M. H. Yang, and S. Belongie. Visual Tracking with Online Multiple Instance Learning. In Proc. CVPR, 2009.

# Previous Works – Adaptive Tracking-by-detection



Tracking

Sampler

Learner

Problem 1

What is the best way to generate labelled samples?

+  -

Sampler

Tracking

Labeller

+ -

## Problem 2

Label prediction and position estimation are different objectives.

# Main Idea



structured output prediction

Tracking

Learner

# Main Contributions

Structured output tracking

      Avoid the intermediate classification step

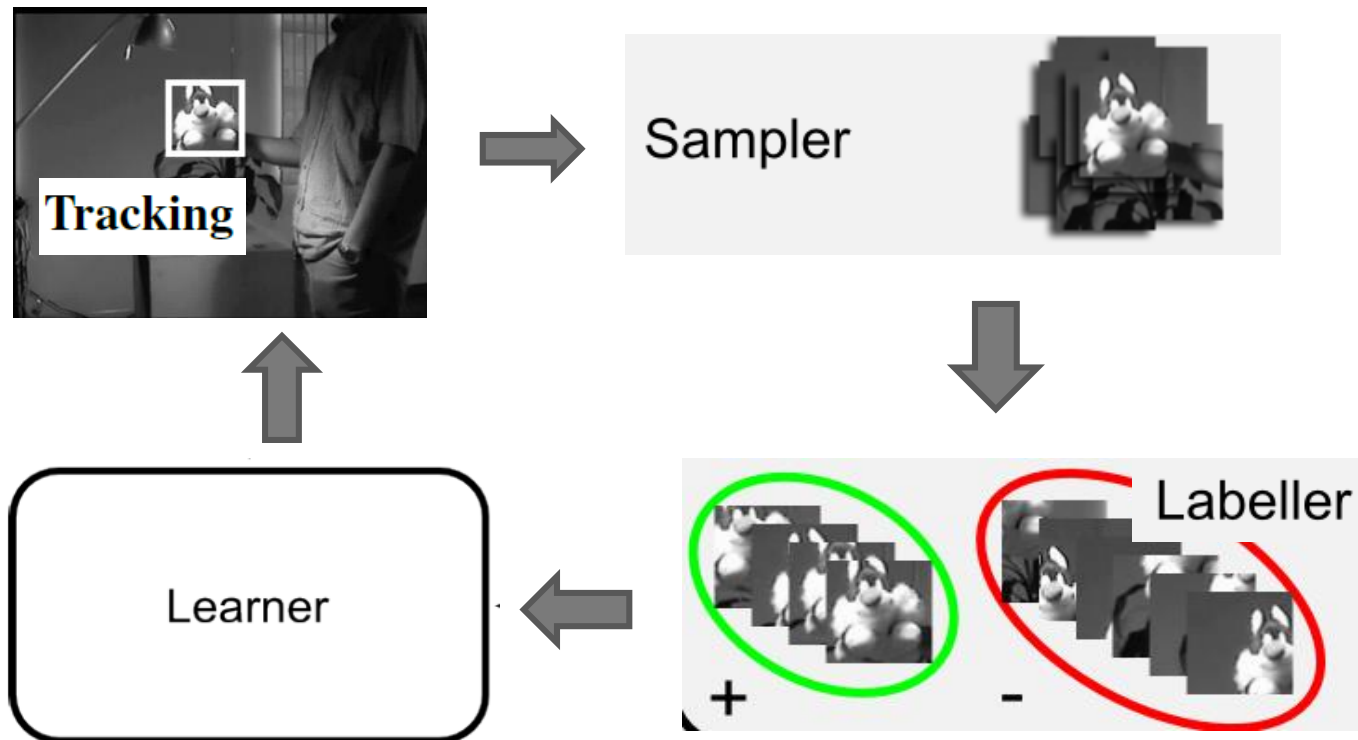Online learning and  budgeting mechanism

      Prevents too many training data

# Outline

Previous work

- Tracking-by-detection
- Adaptive tracking-by-detection

**Methods**

- **Structured output tracking**
- Online optimization and budget mechanism

Experiments and results

# Structured Output Tracking

tracker position

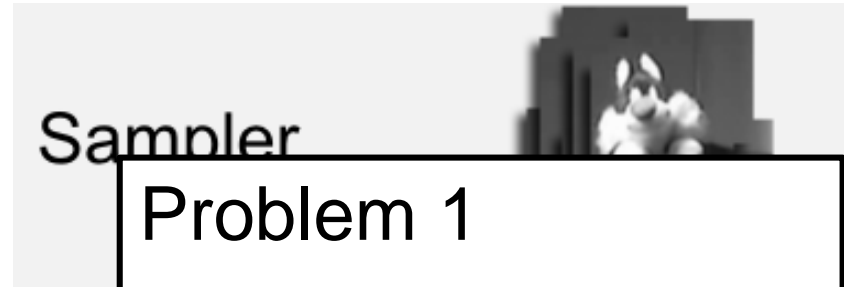$$\mathbf{y}_t = f(\mathbf{x}_t^{\mathbf{P}_{t-1}}) = \arg\max_{\mathbf{y} \in \mathcal{Y}} F(\mathbf{x}_t^{\mathbf{P}_{t-1}}, \mathbf{y})$$

Best motions

search window

image patch

M. B. Blaschko and C. H. Lampert. Learning to Localize Objects with Structured Output Regression. In Proc. ECCV, 2008.

# Structured Output Tracking

The output space is all transformations instead of the binary labels.

$$\mathbf{y}_t = f(\mathbf{x}_t^{\mathbf{P}_{t-1}}) = \underset{\mathbf{y} \in \mathcal{Y}}{\arg\max}\, F(\mathbf{x}_t^{\mathbf{P}_{t-1}}, \mathbf{y})$$

M. B. Blaschko and C. H. Lampert. Learning to Localize Objects with Structured Output Regression. In Proc. ECCV, 2008.

# Structured SVM Model



(a)

(b)

(c)

(d)

The SVM score should correlate with overlapping size with the best tracking bounding box.

S. Avidan. Support Vector Tracking. IEEE Trans. on PAMI, 26:1064–1072, 2004.

# Structured Output Tracking

**Algorithm 2** Struck: Structured Output Tracking

**Require:** $f_t, p_{t-1}, \mathcal{S}_{t-1}$

1: *Estimate change in object location*
2: $y_t = \arg\max_{y \in \mathcal{Y}} F(x_t^{P_{t-1}}, y)$
3: $p_t = p_{t-1} \circ y_t$
4: *Update discriminant function*
5: $(i, y_+, y_-) \leftarrow$ PROCESSNEW$(x_t^{P_t}, y^0)$
6: SMOSTEP$(i, y_+, y_-)$
7: BUDGETMAINTENANCE$()$
8: **for** $j = 1$ to $n_R$ **do**
9: $\quad (i, y_+, y_-) \leftarrow$ PROCESSOLD$()$
10: $\quad$ SMOSTEP$(i, y_+, y_-)$
11: $\quad$ BUDGETMAINTENANCE$()$
12: $\quad$ **for** $k = 1$ to $n_O$ **do**
13: $\quad\quad (i, y_+, y_-) \leftarrow$ OPTIMIZE$()$
14: $\quad\quad$ SMOSTEP$(i, y_+, y_-)$
15: $\quad$ **end for**
16: **end for**
17: **return** $p_t, \mathcal{S}_t$

# Structured Output Tracking



**Tracking**

---

**Algorithm 2** Struck: Structured Output Tracking

**Require:** $f_t$, $p_{t-1}$, $\mathcal{S}_{t-1}$

1: *Estimate change in object location*
2: $y_t = \underset{y \in \mathcal{Y}}{\arg\max}\, F(x_t^{P_{t-1}}, y)$
3: $p_t = p_{t-1} \circ y_t$
4: *Update discriminant function*
5: $(i, y_+, y_-) \leftarrow \text{PROCESSNEW}(x_t^{P_t}, y^0)$
6: $\text{SMOSTEP}(i, y_+, y_-)$
7: $\text{BUDGETMAINTENANCE}()$
8: **for** $j = 1$ to $n_R$ **do**
9: $\quad (i, y_+, y_-) \leftarrow \text{PROCESSOLD}()$
10: $\quad \text{SMOSTEP}(i, y_+, y_-)$
11: $\quad \text{BUDGETMAINTENANCE}()$
12: $\quad$ **for** $k = 1$ to $n_O$ **do**
13: $\quad\quad (i, y_+, y_-) \leftarrow \text{OPTIMIZE}()$
14: $\quad\quad \text{SMOSTEP}(i, y_+, y_-)$
15: $\quad$ **end for**
16: **end for**
17: **return** $p_t$, $\mathcal{S}_t$

# Structured Output Tracking
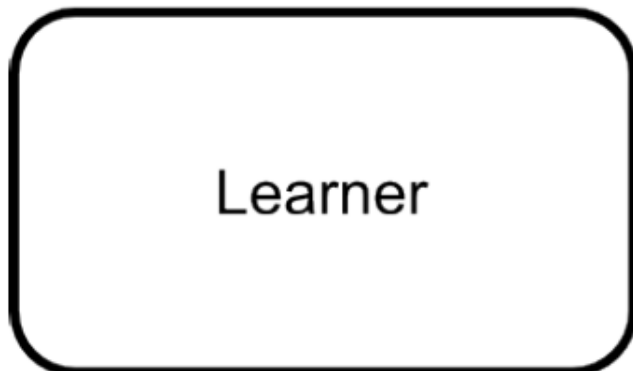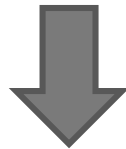


**Tracking**

Learner

**Algorithm 2** Struck: Structured Output Tracking

**Require:** $f_t, p_{t-1}, \mathcal{S}_{t-1}$

1: *Estimate change in object location*
2: $y_t = \arg\max_{y \in \mathcal{Y}} F(x_t^{P_{t-1}}, y)$
3: $p_t = p_{t-1} \circ y_t$
4: *Update discriminant function*
5: $(i, y_+, y_-) \leftarrow \text{PROCESSNEW}(x_t^{P_t}, y^0)$
6: $\text{SMOSTEP}(i, y_+, y_-)$
7: $\text{BUDGETMAINTENANCE}()$
8: **for**
9:
10:
11:
12:
13:
14: $\text{SMOSTEP}(i, y_+, y_-)$
15: **end for**
16: **end for**
17: **return** $p_t, \mathcal{S}_t$

Come back later

# Structured output SVM



$$\mathbf{y}_t = f(\mathbf{x}_t^{\mathbf{P}_{t-1}}) = \arg\max_{\mathbf{y}\in\mathcal{Y}} F(\mathbf{x}_t^{\mathbf{P}_{t-1}}, \mathbf{y})$$

Learner

$$\min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\xi_i$$

$$\text{s.t.} \quad \forall i: \xi_i \geq 0$$

$$\forall i, \forall \mathbf{y} \neq \mathbf{y}_i: \langle \mathbf{w}, \delta\mathbf{\Phi}_i(\mathbf{y})\rangle \geq \Delta(\mathbf{y}_i, \mathbf{y}) - \xi_i$$

Efficient SMO optimization (CS229, EE364)
Kernels (CS229)

A. Bordes, L. Bottou, P. Gallinari, and J. Weston. Solving multiclass support vector machines with LaRank. In Proc. ICML, 2007.

# Structured output SVM

Gaussian kernel between image feature vectors (CS229)

$$k(\mathbf{x}, \bar{\mathbf{x}}) = \exp(-\sigma\|\mathbf{x} - \bar{\mathbf{x}}\|^2),$$

Haar-like features (CS231A, CS232)

The responses of the Haar features are the input vectors of the kernel
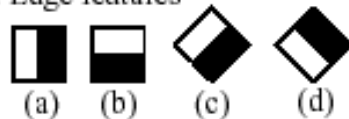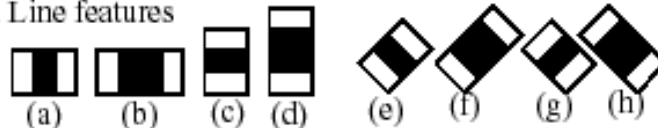
# Online optimization
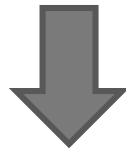


**Algorithm 2** Struck: Structured Output Tracking

**Require:** $\mathbf{f}_t, \mathbf{p}_{t-1}, \mathcal{S}_{t-1}$

1: *Estimate change in object location*
2: $\mathbf{y}_t = \arg\max_{\mathbf{y} \in \mathcal{Y}} F(\mathbf{x}_t^{\mathbf{P}_{t-1}}, \mathbf{y})$
3: $\mathbf{p}_t = \mathbf{p}_{t-1} \circ \mathbf{y}_t$
4: *Update discriminant function*
5: $(i, \mathbf{y}_+, \mathbf{y}_-) \leftarrow \text{PROCESSNEW}(\mathbf{x}_t^{\mathbf{P}_t}, \mathbf{y}^0)$
6: $\text{SMOSTEP}(i, \mathbf{y}_+, \mathbf{y}_-)$
7: $\text{BUDGETMAINTENANCE}()$
8: **for** $j = 1$ to $n_R$ **do**
9: $\quad (i, \mathbf{y}_+, \mathbf{y}_-) \leftarrow \text{PROCESSOLD}()$
10: $\quad \text{SMOSTEP}(i, \mathbf{y}_+, \mathbf{y}_-)$
11: $\quad \text{BUDGETMAINTENANCE}()$
12: $\quad$ **for** $k = 1$ to $n_O$ **do**
13: $\quad\quad (i, \mathbf{y}_+, \mathbf{y}_-) \leftarrow \text{OPTIMIZE}()$
14: $\quad\quad \text{SMOSTEP}(i, \mathbf{y}_+, \mathbf{y}_-)$
15: $\quad$ **end for**
16: **end for**
17: **return** $\mathbf{p}_t, \mathcal{S}_t$

# Outline

Previous work

- Tracking-by-detection
- Adaptive tracking-by-detection

**Methods**

- Structured output tracking
- **Online optimization and budget mechanism**

Experiments and results
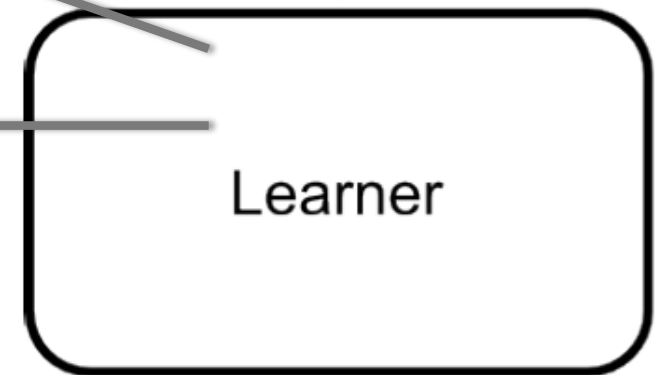
# Online optimization

PROCESSNEW():

- Processes a new example

PROCESSOLD():

- Processes an existing support pattern

OPTIMIZE():

- Processes an existing support pattern chosen at random



Learner

# Budget mechanism

The number of support vectors increase over time.

Computational and storage costs grow linearly with the number of support vectors.

# Incorporating a budget

A budget (limit) of the number of supporting vectors.

Remove the support vector which results in the smallest change to the weight vector

K. Crammer, J. Kandola, R. Holloway, and Y. Singer. Online Classification on a Budget. In NIPS, 2003.
Z. Wang, K. Crammer, and S. Vucetic. Multi-Class Pegasos on a Budget. In Proc. ICML, 2010. 2

# Outline

Previous works

- Tracking-by-detection
- Adaptive tracking-by-detection

Methods

- Structured output tracking
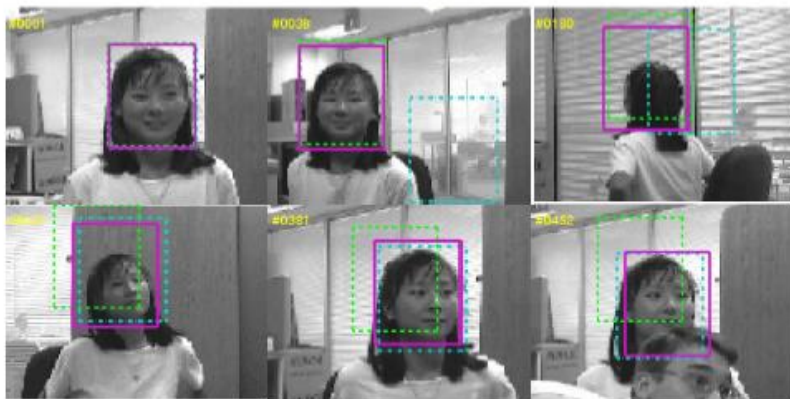- Online optimization and budget mechanism

**Experiments and results**

# Experiments

- Haar-like features

  - 6 different types arranged on a grid at 2 scales on a 4 x 4 grid, resulting in 192 features

- Search radius 60, 5 radial and 16 angular divisions.

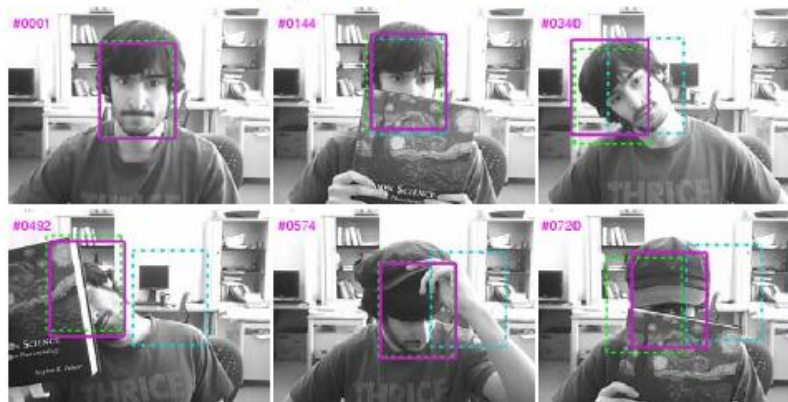- Budget size is as low as B = 20, 50, 100, inf.

# Dataset



(A) Girl (B) Tiger 2 (C) David Indoor (D) Occluded Face 2

http://vision.ucsd.edu/~bbabenko/project_miltrack.shtml;
B. Babenko, M. H. Yang, and S. Belongie. Visual Tracking with Online Multiple Instance
Learning. In Proc. CVPR, 2009.

# Overlap criterion

Jaccard similarity of bounding boxes

$$s^o_{\mathbf{p}_t}(\mathbf{y}^i_t, \mathbf{y}^j_t) = \frac{(\mathbf{p}_t \circ \mathbf{y}^i_t) \cap (\mathbf{p}_t \circ \mathbf{y}^j_t)}{(\mathbf{p}_t \circ \mathbf{y}^i_t) \cup (\mathbf{p}_t \circ \mathbf{y}^j_t)}$$

## Metric 2: Jaccard Similarity

$$100 \cdot \frac{\sum_{i=1}^{n}[\alpha_i = 1 \wedge f_i = 1]}{\sum_{i=1}^{n}[\alpha_i = 1 \vee f_i = 1]}$$

$\alpha_i$: estimated foreground/background label
$f_i$: ground truth foreground/background label

# Results



http://www.samhare.net/research/struck

# Visualization of the support vector set



(a) *girl*

(b) *david*

# Comparison



http://www.samhare.net/research/struck

# Results

Struck with the smallest budget size (B = 20) outperforms the state-of-the-art.

Average frames per second: 12 – 21.

# Extensions

- Used more objection representations

  - Haar-like features
  - Raw pixel features
  - Histogram features
- Combining multiple kernels seems to improve results, but not significantly.


- Use key points and associated descriptors for object detection.
- Consider other machine learning algorithms.

# Main Contributions

Structured output tracking

      Avoid the intermediate classification step


Online learning and  budgeting mechanism

      Prevents too many training data

# References

Sam Hare, Amir Saffari Philip H. S. Torr  Struck: Structured Output Tracking with Kernels International Conference on Computer Vision (ICCV), 2011

A. Bordes, L. Bottou, P. Gallinari, and J. Weston. Solving multiclass support vector machines with LaRank. In Proc. ICML, 2007.

I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun. Large Margin Methods for Structured and Interdependent Output Variables. JMLR, 6:1453–1484, Dec. 2005.

K. Crammer, J. Kandola, R. Holloway, and Y. Singer. Online Classification on a Budget. In NIPS, 2003.

P. Viola and M. J. Jones. Robust real-time face detection. IJCV, 57:137–154, 2004.

# Thank you!