

# Scalable Kernel Correlation Filter with Sparse Feature Integration

Andrés Solís Montero, Jochen Lang and Robert Laganière.



uOttawa

University of Ottawa

December 12, 2016

# Outline

## 1 Motivation, Problems and Objectives

- Motivation
- Problem
- Objectives
- Contributions
- Related Work

## 2 Algorithm Overview

- Estimation Position
- Adjustable Windows
- Estimate Scale
- Improving Performance

## 3 Evaluation Methodology

- Relevant Datasets
- Performance Measures

## 4 Results

- Speed
- Visual Tracker Benchmark
- VOT Challenges

## 5 Conclusions

- Conclusions
- Future Work

# Motivation

- Fast object tracking with live learning
- Object representation, independent of the type of object
- Live estimation of location and scale changes
- General solution for tracking objects??

## Problem

- Tracking object with a moving camera
- No information of the object except an initial selection
- Challenging scenarios and object representations, i.e., partial occlusions, noise, and small and low textured objects
- Estimating location and change of scale
- Speed performance and scalability

## Objectives

- Develop a fast and accurate tracking framework
- Estimate changes in location and scale
- Uses a general object representation
- Benchmark the solution: Visual Benchmark and VOT Challenges
  - Precision, Success, Accuracy, and Robustness

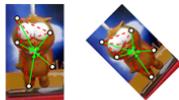
## Contributions

- Extended the KCF framework to add on-line scale estimation
- Improved object/background separation.
- Combines sparse and dense object representations to estimate location and scale on-line
- Improved real-time frame rates and low latency using fHOG (SSE2) and Intel's CCS format for Fourier spectrums
- Improved precision, success, accuracy, and robustness
- Possibility of processing high dimensional data with different feature/scale/correlation estimation methods



Tracking Learning and Detection (TLD)  
Z.Kalal, K. Mikolajczyk and J. Matas 2012.

Consensus-based Match. Track. of Keypoints  
(CMT) G. Nebehay & Roman Pflugfelder, 2014.



Structured Output Tracking with Kernels  
(Struck) S. Hare, A. Saffari and P. H. Torr, 2011

Object Tracking by Oversampling Local Features  
(Alien) F. Pernici and A. del Bimbo, 2014.



# Object Representations

## Dense Representations



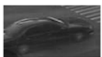
Color : RGB, HSV, HLS ....



Histogram



Frequency Domain



Single channel: gray , tir, ...



Gradient, HOG, ...

## Sparse Representations



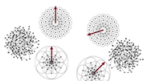
Keypoints



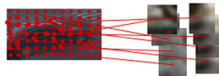
Grid points



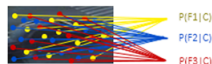
Random points



Descriptors



Patches



Classifiers



## Related Work

- Visual Tracker Benchmark: 29 Trackers.
- VOT Challenges: 27 Trackers (2013), 38 Trackers (2014) ...
- Among most relevant work:
  - TLD, SCM, Struck, CMT, Alien, KCF, CSK, SAMF, etc

## Selected Work

- Henriques, J. F. et al., High-Speed Tracking with Kernelized Correlation Filters, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2015.

# Algorithm Overview - KCF- Estimating Location

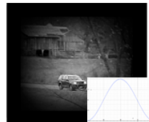
Image



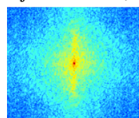
$I$



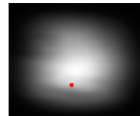
$$I_w = I * W$$



$$I_f = fft(I_w)$$



$$iffit(I_f .* conj(T_w))$$



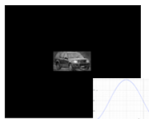
Template



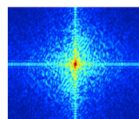
$T$



$$T_w = T * W$$



$$T_f = fft(T_w)$$

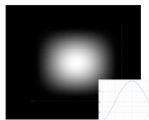


Cosine Window

$$hann(N) = \frac{(1 - \cos(\frac{2\pi x}{N-1}))}{2}$$

$$W = hann(I_h) * hann(I_w)'$$

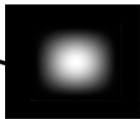
$W$



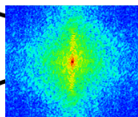
$$O(N \log N)$$

# Algorithm Overview - KCF - Estimating Location

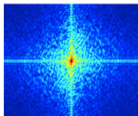
Weighted vertical and horizontal shifts



Object model with encoded shifts



Autocorrelation



Gaussian correlation

$$k^{xx'} = e\left(-\frac{1}{\sigma^2}(\|x\|^2 + \|x'\|^2 - 2F^{-1}(\sum_c x_f^c \odot (x_f^c)^*))\right)$$

$$\alpha_f = \frac{y_f}{k_f^{xx'} + \lambda}$$

Learning Formula

# Algorithm Overview - KCF

## Algorithm 1 : KCF.

Variables with subscript  $f$  are in the frequency domain. Circled operators represent element-wise operations (i.e.,  $\odot$  and  $\oslash$ ).

- $w\text{-sz}$ : size of the tracked region, ( $W \times H$ ).
- $pos$ : center location of the tracker in spatial domain.
- $patch$ : region of  $img$  centered at  $pos$  with size  $w\text{-sz}$ , ( $W \times H \times C$ ).
- $features(x)$ : extracted features (e.g., HoG), ( $m \times n \times c$ ).
- $cos\text{-window}$ : cosine window weights each feature channel, ( $m \times n \times 1$ ).

- 1: **for each**  $img$  **in sequence**:
- 2:   **if not** first image:
  - 3:      $patch \leftarrow \text{region}(img, pos, w\text{-sz})$
  - 4:      $z_f \leftarrow F(\text{features}(patch) \odot cos\text{-window})$
  - 5:      $k_f^{z\tilde{x}} \leftarrow F(\text{correlation}(z_f, \tilde{x}_f)) \triangleright \text{Eq. (2)}$
  - 6:      $pos \leftarrow pos + \underset{loc}{\text{argmax}}(r(k_f^{z\tilde{x}})) \triangleright \text{Eq. (3)}$
- 7:    $patch \leftarrow \text{region}(img, pos, w\text{-sz})$
- 8:    $x_f \leftarrow F(\text{features}(patch) \odot cos\text{-window})$
- 9:    $k_f^{z\tilde{x}} \leftarrow F(\text{correlation}(x_f, x_f)) \triangleright \text{Eq. (2)}$
- 10:    $\alpha_f \leftarrow y_f \oslash (k_f^{z\tilde{x}} + \lambda) \triangleright \text{Eq. (1)}$
- 11:   **if** first image:  $f \leftarrow 1$  **else**  $f \leftarrow \text{factor}$
- 12:    $\tilde{\alpha}_f \leftarrow f \times \alpha_f + (1 - f) \times \tilde{\alpha}_f$
- 13:    $\tilde{x}_f \leftarrow f \times x_f + (1 - f) \times \tilde{x}_f$

## Learning Formula: Eq. 1

$$\alpha_f = \frac{y_f}{k_f^{x x'} + \lambda}$$

## Gaussian Correlation: Eq. 2

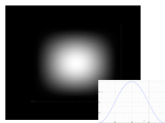
$$k^{x x'} = e\left(-\frac{1}{\sigma^2} (\|x\|^2 + \|x'\|^2 - 2F^{-1}(\sum_c x_f^c \odot (x_f^c)^*))\right)$$

## Response: Eq. 3

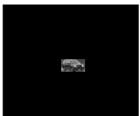
$$r(k_f^{z\tilde{x}}) = F^{-1}(k_f^{z\tilde{x}} \odot \tilde{\alpha}_f)$$

# Algorithm Overview - Adjustable Windows

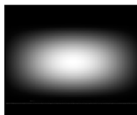
*Cosine Window*



*Scale Changes*



*Adjustable Windows*

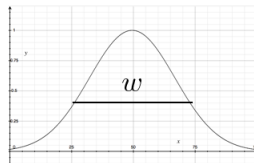


*Gaussian window*

$$g(N) = e^{-\frac{1}{2} \left( \frac{i}{\sigma(N-1)} \right)^2},$$

$$0 \leq i \leq N.$$

$$\sigma = \frac{w}{N}$$

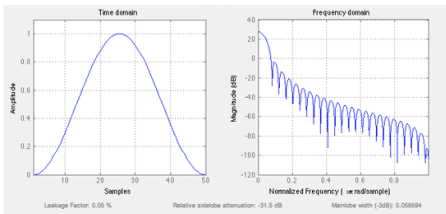


## Algorithm Overview - Adjustable Windows [examples]

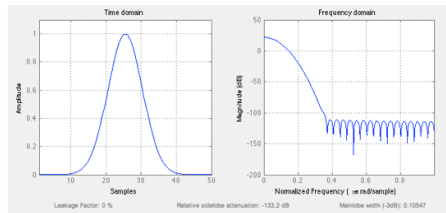
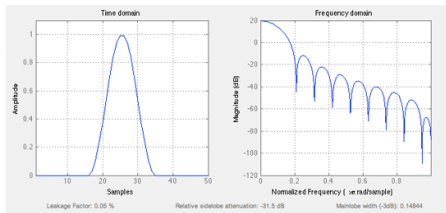
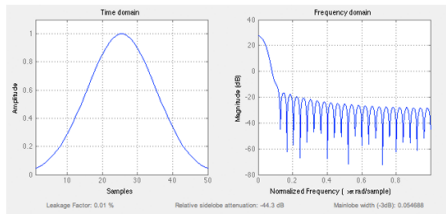


# Cosine vs Gaussian Window

Cosine window



Gaussian window



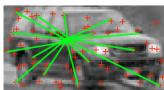
# Algorithm Overview - Estimating Scale

$$\square * scale = \square$$



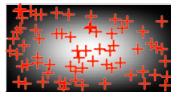
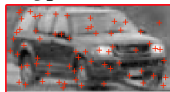
$$\|k_1 - k_j\|_2$$

$$\|m_1 - m_j\|_2$$

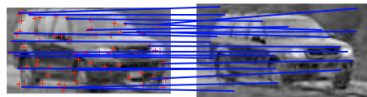


keypoints  $k_i$

weights  $w_i$



optical flow



Weighted arithmetic mean

$$scale = \frac{\sum_i \sum_j w_i w_j * \frac{\|m_i - m_j\|_2}{\|k_i - k_j\|_2}}{\sum_i \sum_j w_i w_j}$$

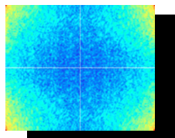
$O(KN)$



## Improving Performance

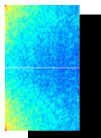
- fast HOG descriptors (SSE instructions)  
Felzenszwalb et al. Object detection with discriminatively trained part, TPAMI 2010.
- Intel's CCS packed format
- Optimal search area  $N = 2^p \times 3^q \times 5^r$   
(e.g.,  $300 \times 300 = 5^2 \times 3 \times 2^2$ , closer power of two is  $512 \times 512$ ).

Full Spectrum  
(real + imag)



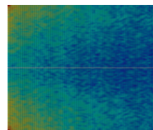
2 Channels

Half



2 Channels

CCS



1 Channel



Object Tracking

# Algorithm sKCF

## Algorithm 2 : sKCF.

Changes to the KCF pipeline are showed in different color.

- $w\_sz$ : size of the tracked region, ( $W \times H$ ).
- $t\_sz$ : size of the target, ( $w \times h$ ).
- $features(x)$ : extracted features (e.g., HoG), ( $m \times n \times c$ ).

```

1: for each img in sequence:
2:   if not first image:
3:     p2 ← region(img, pos, w_sz)
4:     zf ← F(features(p2) ⊙ cos_window)
5:     kfz $\tilde{x}$  ← F(correlation(zf,  $\tilde{x}_f$ )) ▷ Eq. (2)
6:     pos ← pos + argmaxloc(r(kfz $\tilde{x}$ ))) ▷ Eq. (3)
7:     t_sz ← t_sz * scale( $K^{p1}$ ,  $K^{p2}$ ) ▷ Eq. (7)
8:   p1 ← region(img, pos, w_sz)
9:   gauss ← G(m, n, t_sz, w_sz) ▷ Eq. (4, 5)
10:  xf ← F(features(p1) ⊙ gauss)
11:  kfx $x'$  ← F(correlation(xf, xf)) ▷ Eq. (2)
12:   $\alpha_f$  ← yf ⊙ (kfx $x'$  +  $\lambda$ ) ▷ Eq. (1)
13:  if first image: f ← 1 else f ← factor
14:   $\tilde{\alpha}_f$  ← f ×  $\alpha_f$  + (1 - f) ×  $\tilde{\alpha}_f$ 
15:   $\tilde{x}_f$  ← f × xf + (1 - f) ×  $\tilde{x}_f$ 

```

Learning Formula: Eq. 1

$$\alpha_f = \frac{y_f}{k_f^{x x'} + \lambda}$$

Gaussian Correlation: Eq. 2

$$k_f^{x x'} = e\left(-\frac{1}{\sigma^2}(\|x\|^2 + \|x'\|^2 - 2F^{-1}(\sum_c x_f^c \odot (x_f^c)^*))\right)$$

Response: Eq. 3

$$r(k_f^{z \tilde{x}}) = F^{-1}(k_f^{z \tilde{x}} \odot \tilde{\alpha}_f)$$

Gaussian Window : Eq. 4

$$g(N, \sigma) = \exp\left(-\frac{1}{2} \left(\frac{i}{\sigma(N-1)}\right)^2\right), \quad 0 \leq i \leq N.$$

Scale Estimation: Eq. 7

$$scale(K^{p1}, K^{p2}) = \frac{\sum_i^T \sum_j^T w_i w_j * \frac{\|K_i^{p2} - K_j^{p2}\|^2}{\|K_i^{p1} - K_j^{p1}\|^2}}{\sum_i^T \sum_j^T w_i w_j}.$$

# Datasets

## Tracker Benchmark v1.0 [Yi Wu et al. 2013]

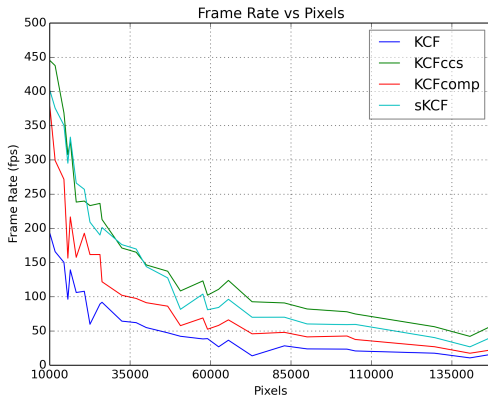
- 50 sequences with 29 trackers
- Measures: precision and success

## VOT Challenge [Kristan et al.]

- VOT2013: 16 sequences with 27 trackers
- VOT2014: 25 sequences with 37 trackers
- VOT2015: 60 sequences
- VOTTIR2015: 20 sequences
- Measures: accuracy and robustness/reliability

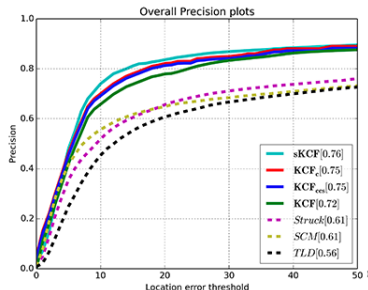
# Speed

Frame rate expressed in frames per second (y-axis of the plot) measured by the number of pixels processed (x-axis of the plot).

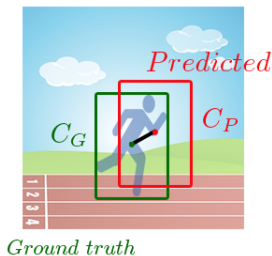


## Precision [Yi Wu et al.]

Precision plot shows the ratio of successful frames whose tracker output is within the given threshold (x-axis of the plot, in pixels) from the ground-truth, measured by the center distance between bounding boxes.

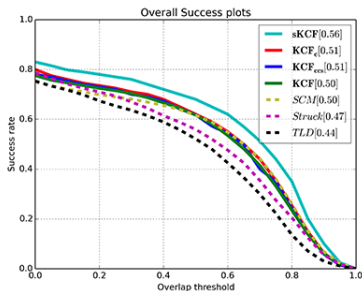


$$\|C_p - C_c\| < T_p$$

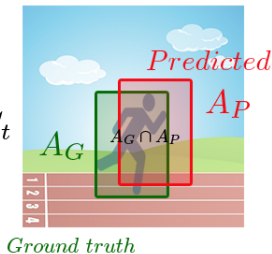


# Success [Yi Wu et al.]

For an overlap threshold (x-axis of the plot), the success ratio is the ratio of the frames whose tracked box has more overlap with the ground-truth box than the threshold.

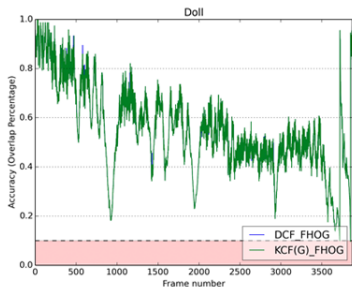


$$\frac{A_G \cap A_P}{A_G \cup A_P} \geq S_t$$

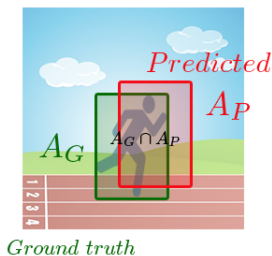


# Accuracy [Kristan et al.]

Overlap between the ground-truth  $A_G$  and the area predicted by a tracker, i.e.,  $A_P$ . The overall accuracy of a sequence is the average accuracy of all the frames in the sequence.

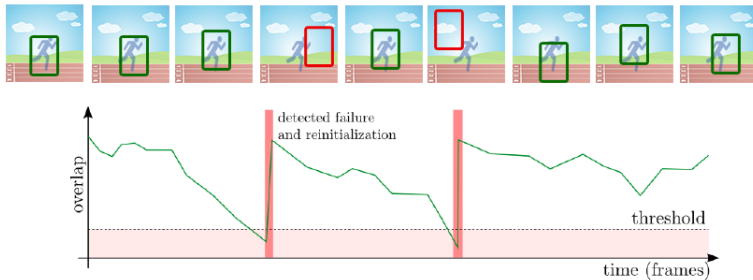


$$\frac{A_G \cap A_P}{A_G \cup A_P}$$



## Robustness/Reliability

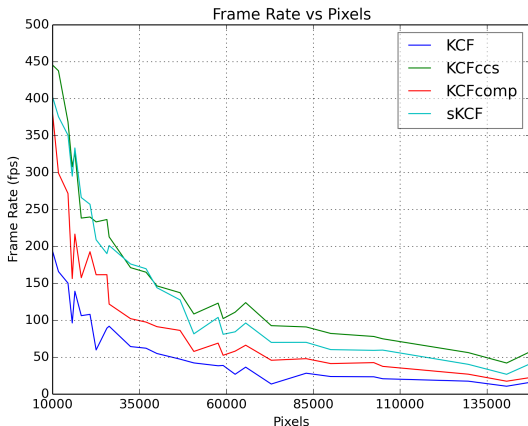
Counts the number of times the tracker failed and had to be reinitialized. Failure occurs when the overlap drops below a threshold.





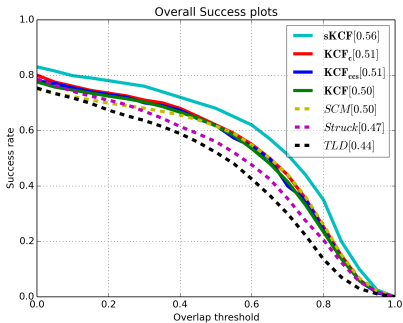
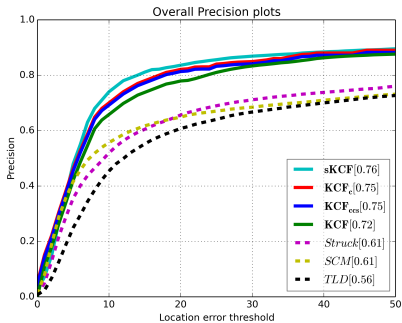
## Speed Benchmark

Comparison between KCF implementation [Henriques et al. 2015] and our solution



# Precision and Success

Dataset: Visual Tracker Benchmark [Yi Wu et al. 2013]



# VOT 2014

Table : VOT 2014 Results

	Overall		Rank			fps
	Acc.	Fail.	Acc.	Rob.	Overall	
DSST	<b>0.65</b>	<b>16.90</b>	5.44	<b>12.17</b>	<b>8.81</b>	5.8
SAMF	<b>0.65</b>	19.23	<b>5.23</b>	12.94	9.09	1.6
sKCF	0.61	18.44	7.68	13.14	10.41	<b>65.4</b>
KCF	0.56	27.14	13.14	18.02	15.58	20.3

# VOT 2015

Table : VOT and VOT TIR 2015 Results

## VOT 2015

	Overall		Rank			
	Acc.	Fail.	Acc.	Rob.	Overall	fps
sKCF	<b>0.50</b>	<b>2.49</b>	<b>2.22</b>	<b>2.60</b>	<b>2.41</b>	<b>64.5</b>
KCF	0.47	2.61	3.29	2.68	2.99	24.4

## VOT TIR 2015

sKCF	0.58	<b>5.28</b>	2.92	<b>2.50</b>	<b>2.71</b>	<b>215.0</b>
KCF	0.56	5.66	3.40	2.65	3.02	94.8

## Conclusions

- Scalable KCF solution that reacts better to object transformations and changes of scale
- Gaussian Window filtering for better object/background separation.
- Combines sparse and dense object representations to estimate location and scale on-line
- Improved real-time frame rates and low latency using fHOG (SSE2) and CCS format for Fourier spectrums
- Improved precision, success, accuracy, and robustness
- Possibility of processing high dimensional data with different scale estimation methods

## Future Work

- Including rotation
- Improve learning methodology, tracker should drop information while occluded
- Improve speed performance
- Compare adjustable filtering functions