

# Beyond Correlation Filters: Learning Continuous Convolution Operators for Visual Tracking

Martin Danelljan, Andreas Robinson,  
Fahad Shahbaz Khan, Michael Felsberg

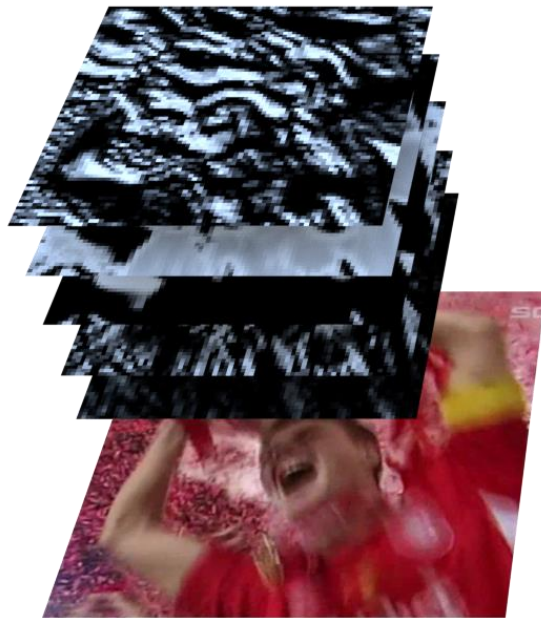
# Discriminative Correlation Filters (DCF)

## Applications

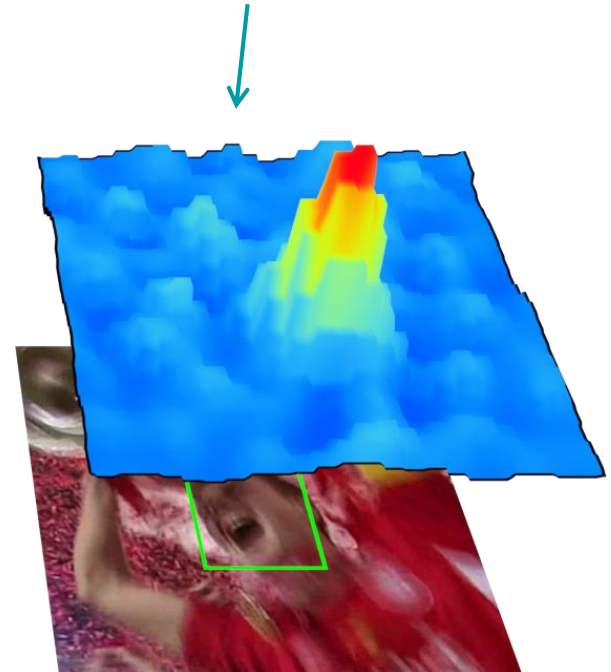
- Object recognition
- Object detection
- Object tracking
  - Among state-of-the-art since 2014
  - KCF, DSST, HCF, SRDCF, Staple ...

# Discriminative Correlation Filters (DCF)

Limitations: Single-resolution feature map



Coarse output scores



## DCF Limitations:

### 1. Single-resolution feature map

- Why a problem?
  - Combine convolutional layers of a CNN
    - Shallow layers: low invariance – high resolution
    - Deep layers: high invariance – low resolution
- How to solve?
  - Explicit resampling?
    - Artefacts, information loss, redundant data
  - Independent DCFs with late fusion?
    - Sub-optimal, correlations between layers

## DCF Limitations:

### 2. Coarse output scores

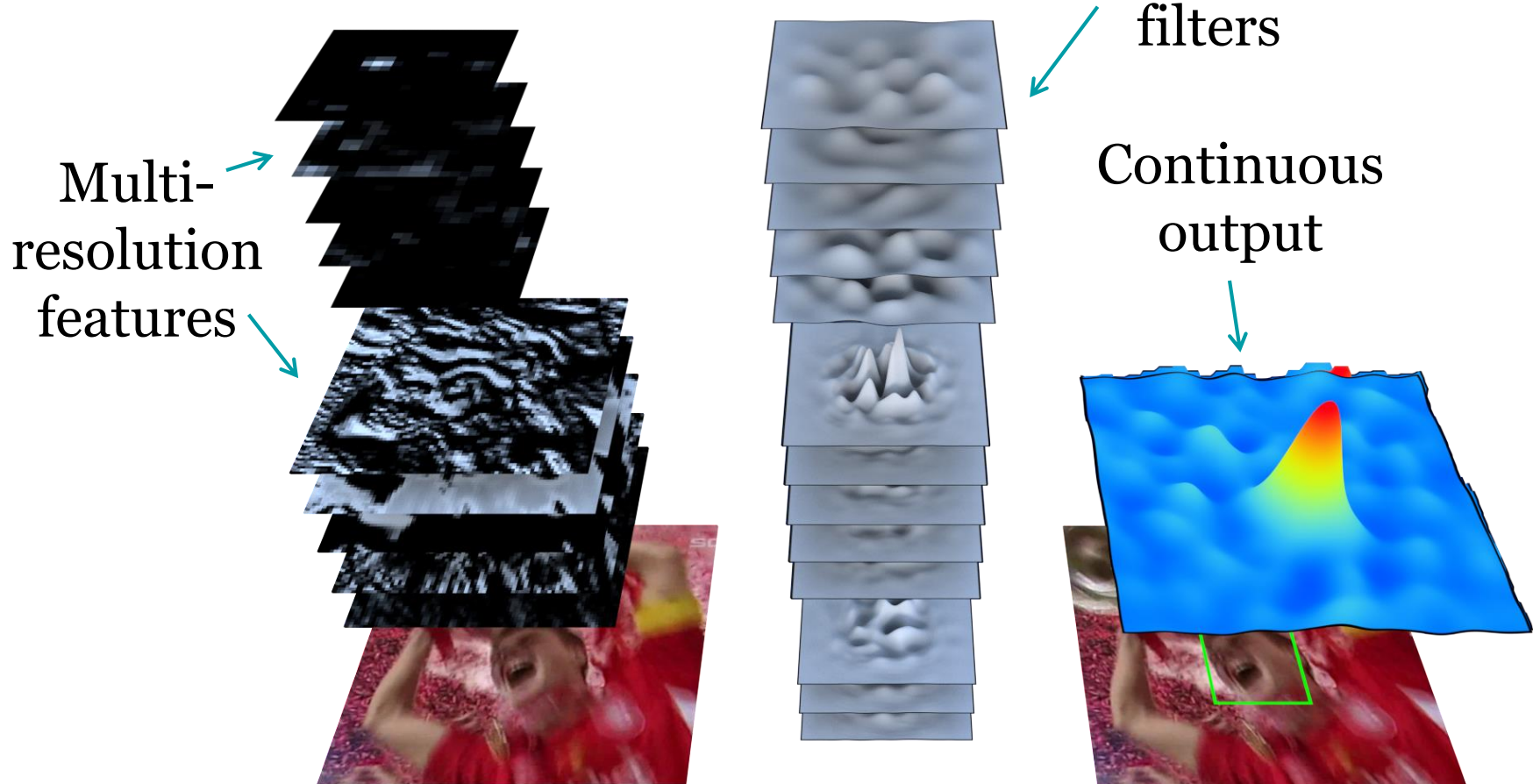
- Why a problem?
  - Accurate localization
    - Sub-grid (e.g. HOG grid) or sub-pixel accuracy
    - More accurate annotations=> less drift
- How to solve?
  - Interpolation?
    - Which interpolation strategy?
  - Interweaving?
    - Costly

## DCF Limitations:

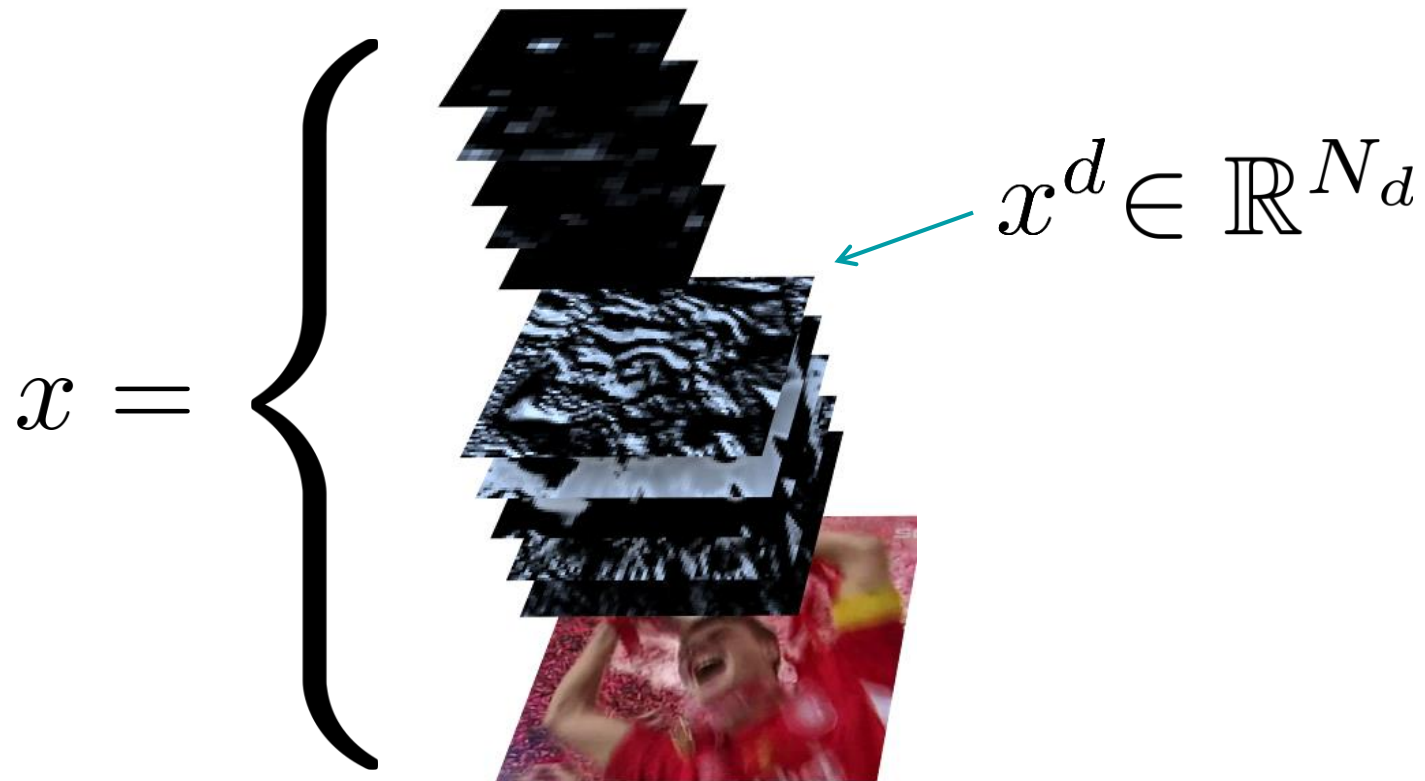
### 3. Coarse labels

- Why a problem?
  - Accurate learning
    - Sub-grid or sub-pixel supervision
- How to solve?
  - Interweaving?
    - Costly
  - Explicit interpolation of features?
    - Artefacts

# Our Approach: Overview



# Multiresolution Features



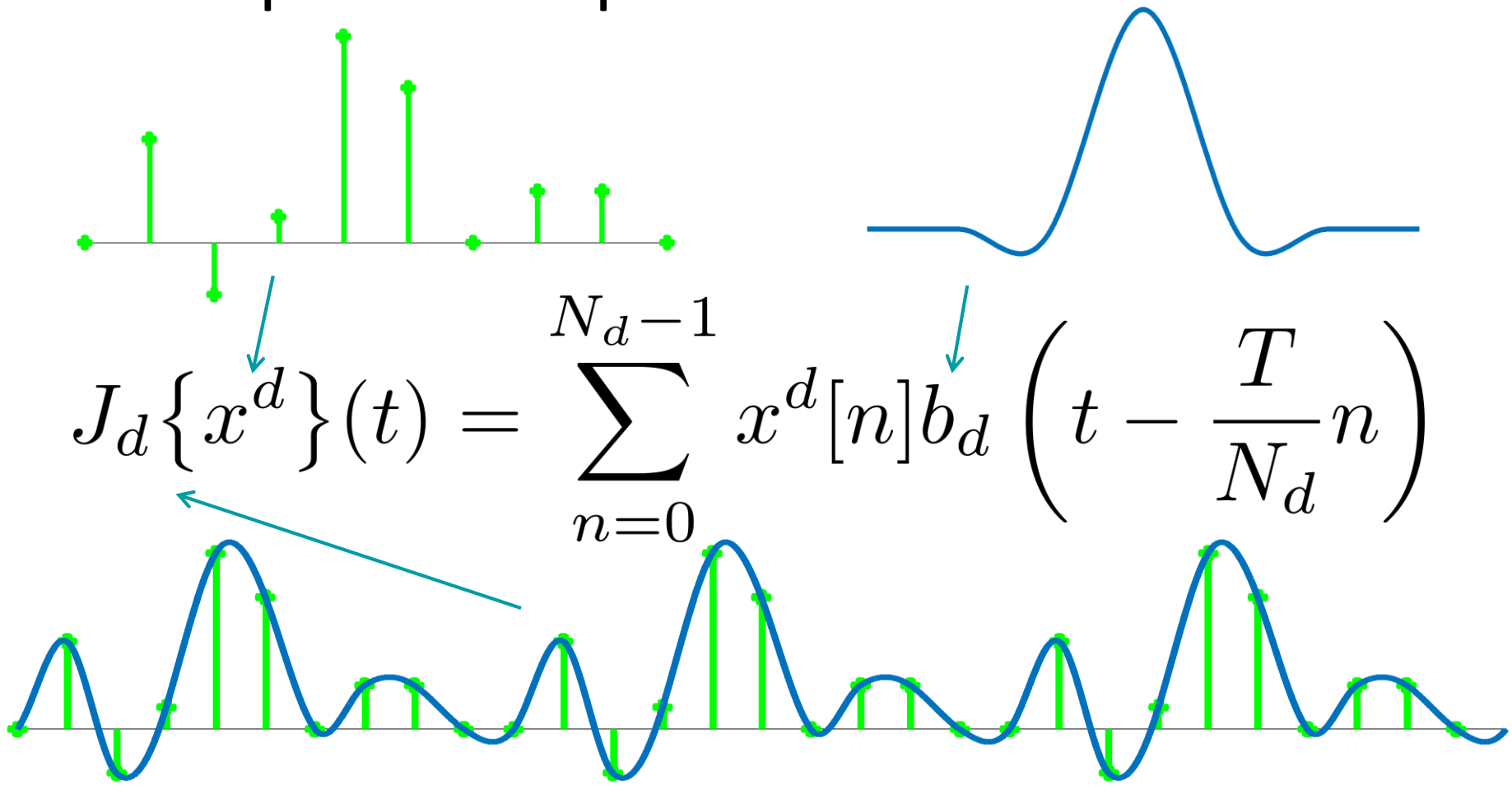


# Interpolation Operator

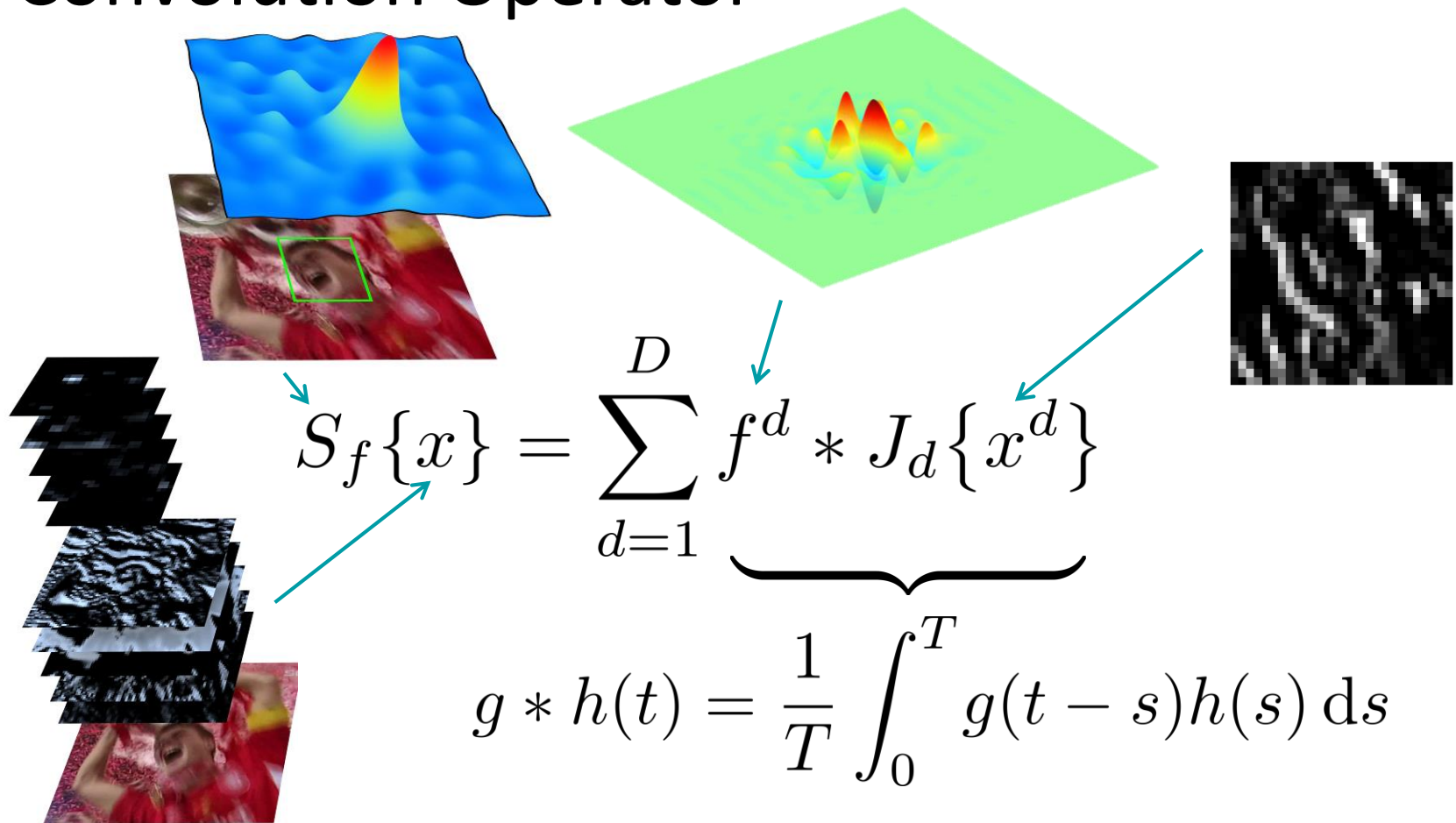
$$J_d : \mathbb{R}^{N_d} \rightarrow L^2(T)$$

$$J_d \left\{ \text{[Low-resolution image]} \right\} = \text{[High-resolution image]}$$

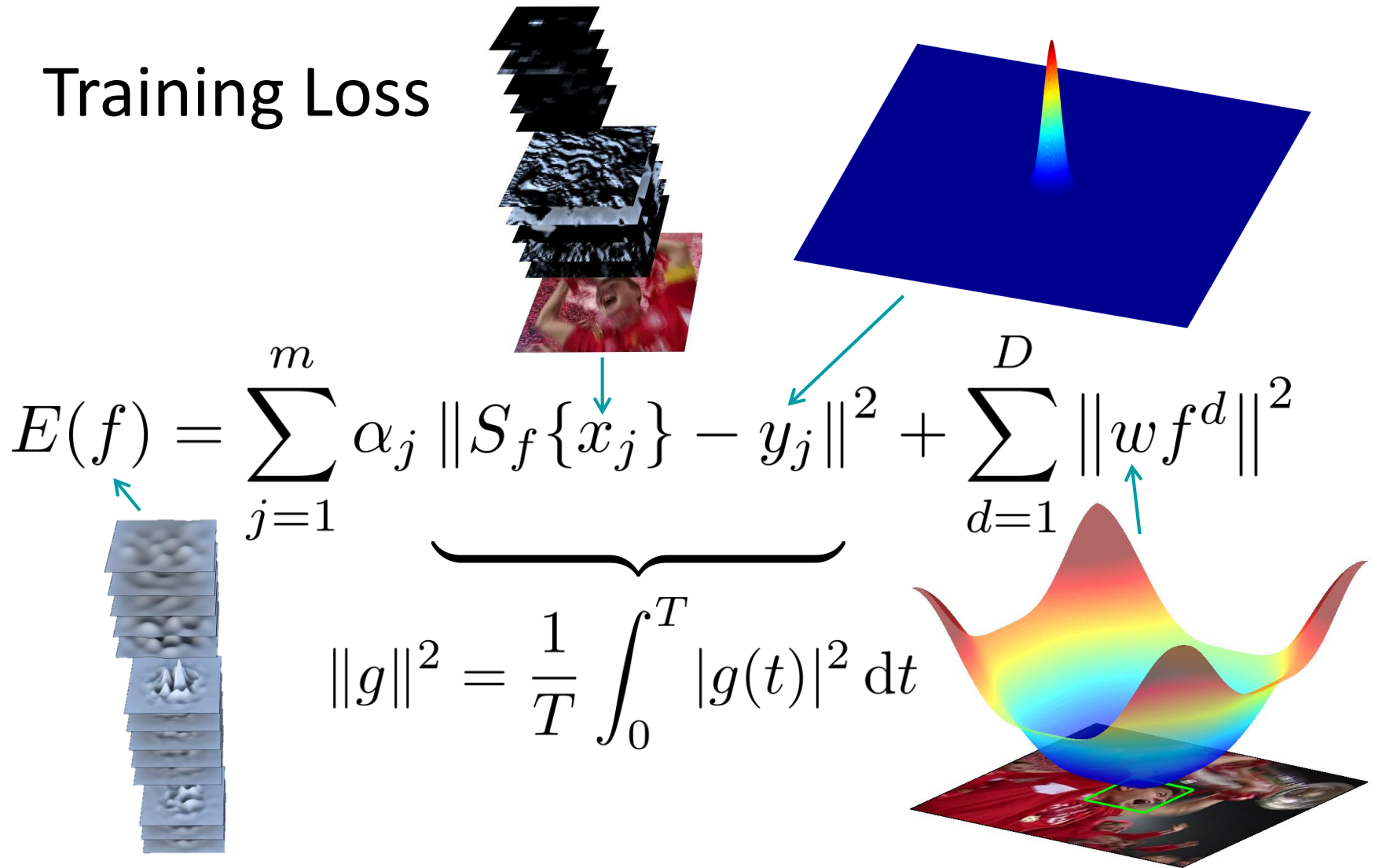
# Interpolation Operator



# Convolution Operator



# Training Loss



The diagram illustrates the training loss function  $E(f)$  with various visual components and mathematical terms:

- Top Left:** A stack of grayscale images representing input frames  $x_j$ .
- Top Right:** A 3D surface plot showing a sharp peak, representing the correlation filter  $f$ .
- Bottom Left:** A stack of 3D surface plots representing the learned filters  $f^d$ .
- Bottom Right:** A 3D surface plot showing a smooth, multi-peaked surface, representing the continuous convolution operator  $S_f$ .
- Center:** The main equation for the training loss:
 
$$E(f) = \sum_{j=1}^m \alpha_j \underbrace{\|S_f\{x_j\} - y_j\|^2}_{\|g\|^2} + \sum_{d=1}^D \|w f^d\|^2$$
- Below the Equation:** The definition of the norm  $\|g\|^2$ :
 
$$\|g\|^2 = \frac{1}{T} \int_0^T |g(t)|^2 dt$$

Arrows indicate the mapping from visual elements to terms in the equation: from the input stack to  $x_j$ , from the correlation filter to  $f$ , from the learned filters to  $f^d$ , and from the continuous convolution operator to  $S_f$ .

# Training Loss – Fourier Domain

$$E(f) = \sum_{j=1}^m \alpha_j \underbrace{\left\| \sum_{d=1}^D \hat{f}^d X_j^d \hat{b}_d - \hat{y}_j \right\|_{\ell^2}^2}_{\infty} + \sum_{d=1}^D \left\| \hat{w} * \hat{f}^d \right\|_{\ell^2}^2$$

$$\|\hat{g}\|_{\ell^2}^2 = \sum_{-\infty}^{\infty} |\hat{g}[k]|^2$$

$$\hat{g}[k] = \langle g, e_k \rangle = \frac{1}{T} \int_0^T g(t) e^{-i \frac{2\pi}{T} kt} dt$$

$$X^d[k] = \sum_{n=0}^{N_d-1} x^d[n] e^{-i \frac{2\pi}{N_d} nk}$$

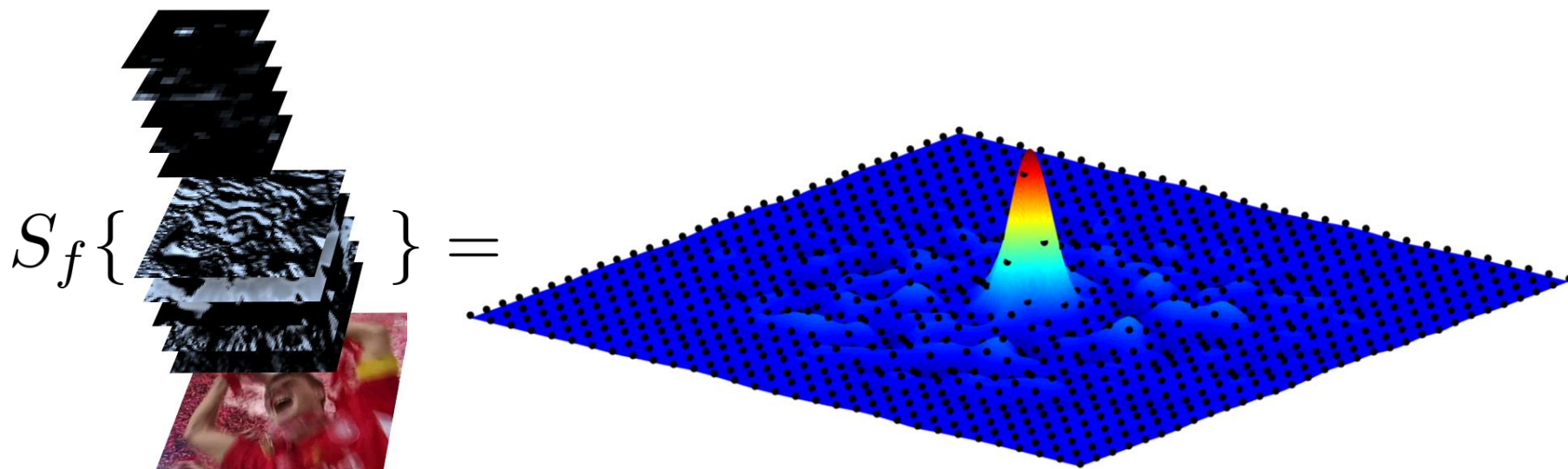
# Training Loss – Fourier Domain

$$E(f) = \sum_{j=1}^m \alpha_j \left\| \sum_{d=1}^D \hat{f}^d X_j^d \hat{b}_d - \hat{y}_j \right\|_{\ell^2}^2 + \sum_{d=1}^D \left\| \hat{w} * \hat{f}^d \right\|_{\ell^2}^2$$

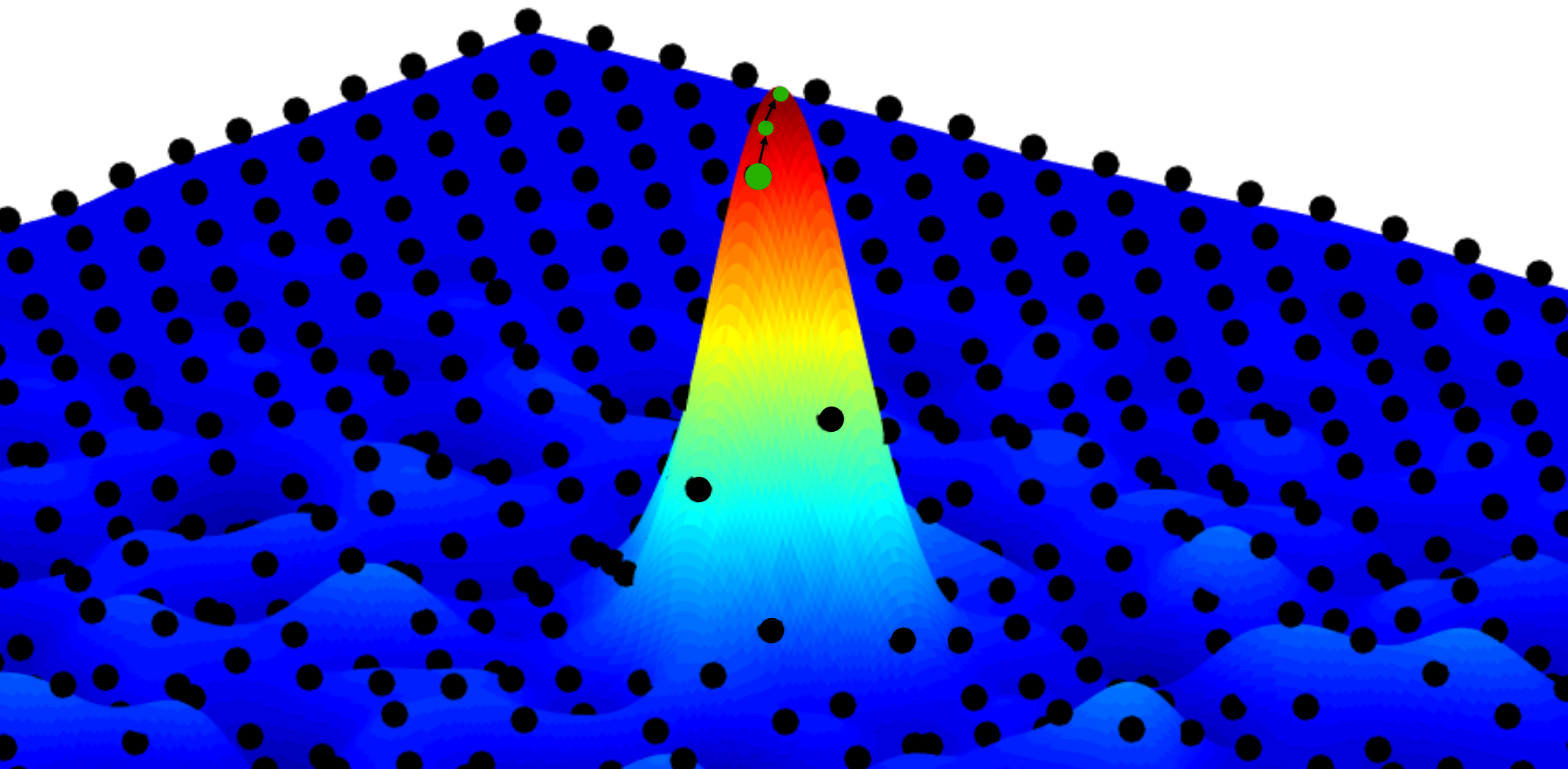


$$(A^H \Gamma A + W^H W) \hat{\mathbf{f}} = A^H \Gamma \hat{\mathbf{y}}$$

# Localization



# Localization





## How to set $y_j$ and $b_d$ ?

- Use periodic summation of functions  $g : \mathbb{R} \rightarrow \mathbb{R}$  :

$$g_T(t) = \sum_{n=-\infty}^{\infty} g(t - nT)$$

- Gaussian function for  $y_j$
- Cubic spline kernel for  $b_d$
- Fourier coefficients  $\hat{y}_j, \hat{b}_d$  with Poisson's summation formula:

$$\hat{g}_T[k] = \frac{1}{T} \hat{g}\left(\frac{k}{T}\right)$$

# Object Tracking Framework: Features

- VGG network
  - Pre-trained on ImageNet
  - No fine-tuning on application specific data

# Object Tracking Framework: Optimization

$$\text{Solving } (A^H \Gamma A + W^H W) \hat{\mathbf{f}} = A^H \Gamma \hat{\mathbf{y}}$$

## SRDCF: Gauss-Seidel

- ☹ Explicit computation of  $(A^H \Gamma A + W^H W)$
- ☹ Sparse matrix handling
- ☹  $\mathcal{O}(D^2)$
- ☹ “Infinite” memory
- ☺ Warm starting: trivial

## C-COT: Conjugate Gradient

- ☺ Only need to evaluate  $(A^H \Gamma A + W^H W) \hat{\mathbf{f}}$
- ☺ **No** sparse matrix handling
- ☺  $\mathcal{O}(D)$
- ☹ Finite memory
- ☹ Warm starting: **non**-trivial
- ☹ Tuning of pre-conditioners

# Object Tracking Framework: Pipeline

- Simple:  
... – Track – Train – Track – Train – ...
- No thresholds
- No hidden “tricks”

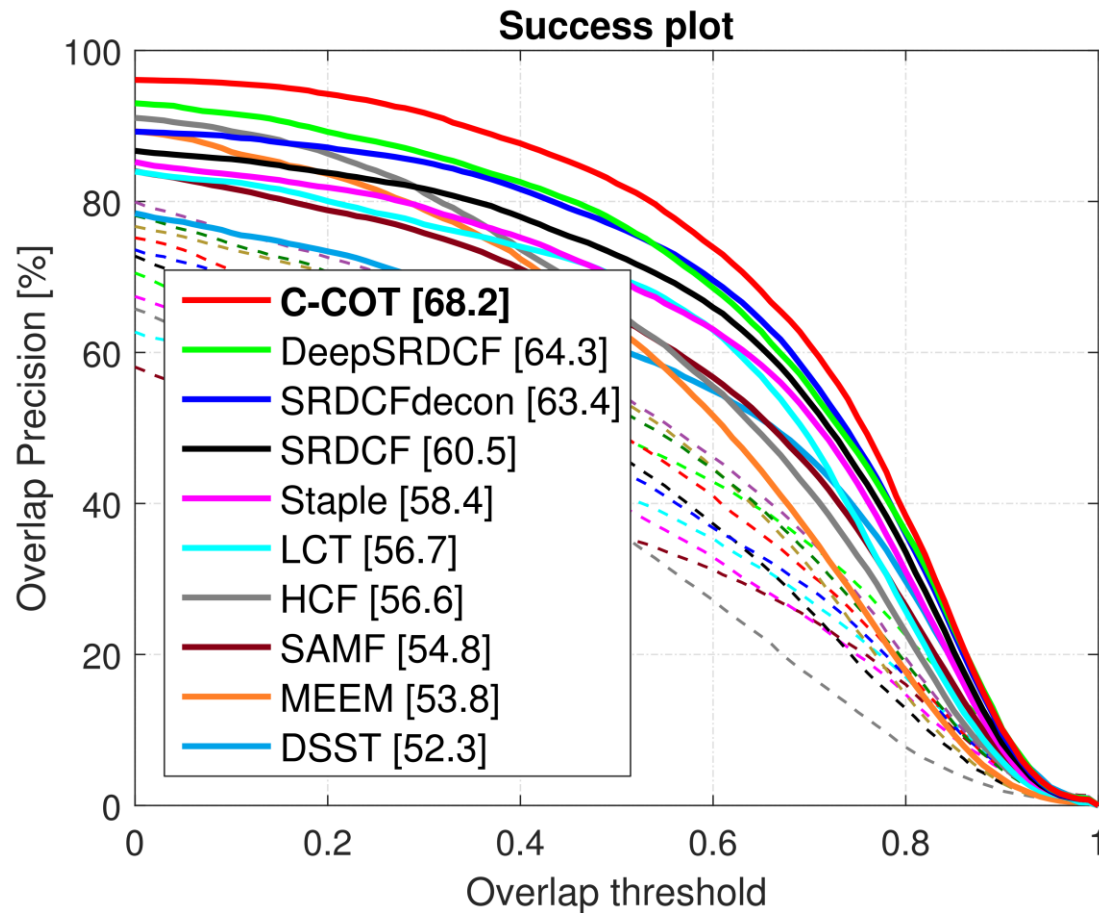
# Experiments: Object Tracking

- 3 datasets: OTB-100, TempleColor, VOT2015
- Layer fusion on OTB:

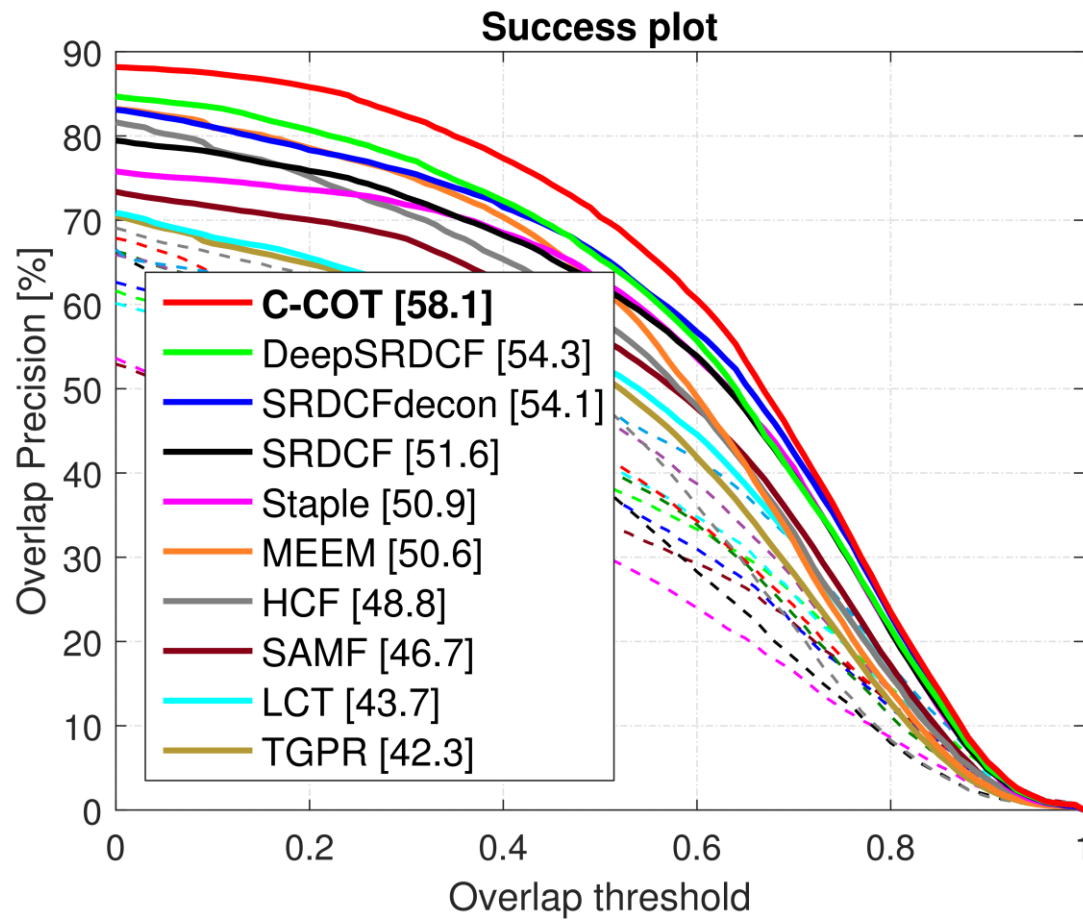
	Layer 0	Layer 1	Layer 5	Layers 0, 1	Layers 0, 5	Layers 1, 5	Layers 0, 1, 5
Mean OP	58.8	78.0	60.0	77.8	70.7	81.8	82.4
AUC	49.9	65.8	51.1	65.7	59.0	67.8	68.2

- Compared to explicit resampling in DCF
  - Performance gain: +7.4% AUC
  - Efficiency gain: –80% data size






# Experiments: OTB (100 videos)



# Experiments: Temple-Color (128 videos)



# Experiments: VOT2016

	Tracker	EAO	A	R	$A_{rank}$	$R_{rank}$	AO	EFO	Impl.
1.	 C-COT	<b>0.331</b>	0.539	<i>0.238</i>	12.000	<b>1.000</b>	<i>0.469</i>	0.507	D M
2.	 TCNN	<i>0.325</i>	0.554	0.268	4.000	<i>2.000</i>	<i>0.485</i>	1.049	S M
3.	 SSAT	<i>0.321</i>	<b>0.577</b>	0.291	<b>1.000</b>	<i>3.000</i>	<b>0.515</b>	0.475	S M
4.	 MLDF	0.311	0.490	<b>0.233</b>	36.000	<b>1.000</b>	0.428	1.483	D M
5.	 Staple	0.295	0.544	0.378	5.000	10.000	0.388	11.144	D C

[Matej et al., ECCV VOT workshop 2016]



# Object Tracking: Speed

- With CNN features: slow  $\sim 1$  FPS (no GPU)
- With HOG features:  $\sim$  real time at SRDCF performance

# Feature Point Tracking Framework

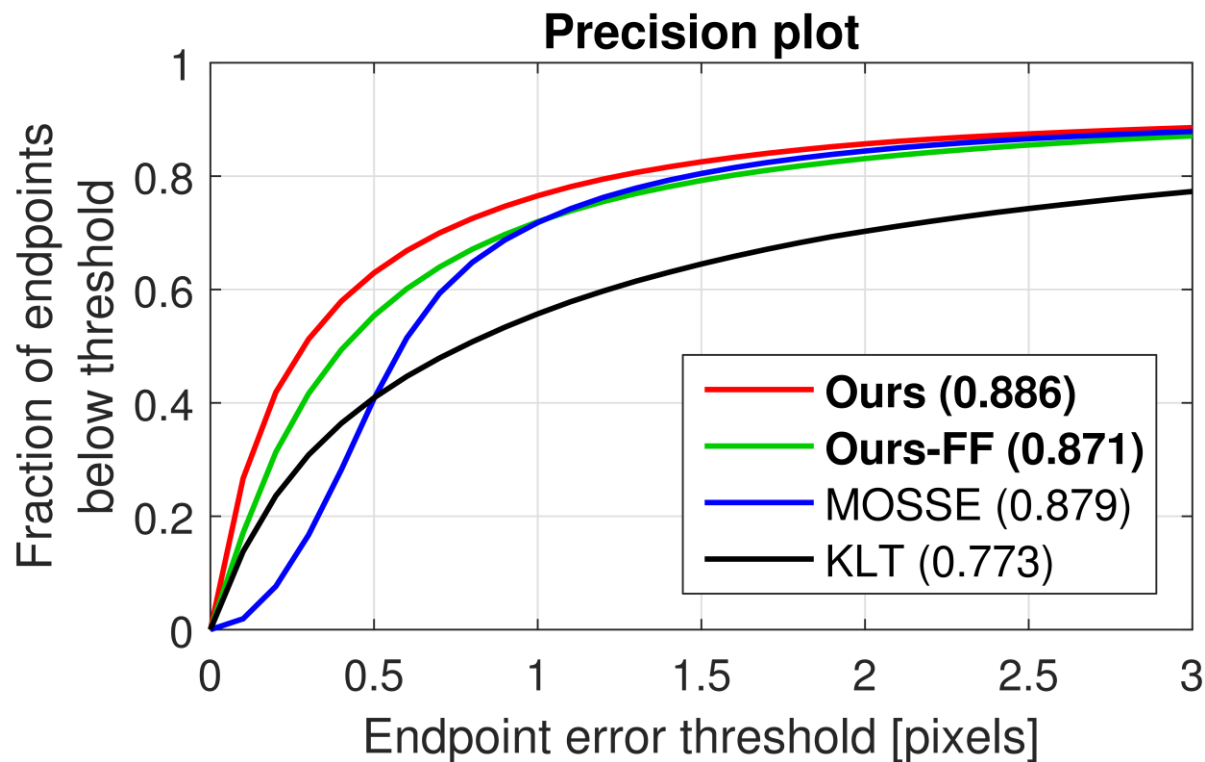
- Grayscale pixel features,  $D = 1$
- Uniform regularization,  $w(t) = \beta$



$$\hat{f}[k] = \frac{\sum_{j=1}^m \alpha_j \overline{X_j[k] \hat{b}[k]} \hat{y}_j[k]}{\sum_{j=1}^m \alpha_j |X_j[k] \hat{b}[k]|^2 + \beta^2}$$

# Experiments: Feature Point Tracking

- The Sintel dataset



# Feature Point Tracking: KITTI

**C-COT**



**KLT**



# Future Work

- Features
  - Fine tuning
  - Unsupervised learning
- Optimization
  - Warm start in CG (theory and heuristics)
  - Preconditioners
  - Implementation aspects
  - Alternative strategies or update rules
- Further explore of the continuous formulation

# Conclusions

- **Continuous domain** learning formulation
  - **Multi-resolution** deep feature maps
  - **Sub-pixel** accurate localization
  - **Sub-pixel** supervision
- Superior results for two applications
  - Object tracking
  - Feature point tracking

# Oral and poster: O-4B-03

Friday afternoon  
(last session)



Martin Danelljan

[www.liu.se](http://www.liu.se)