

# Beyond Covariance: Higher-order Tensor Descriptors and Applications in Computer Vision

P. Koniusz

National ICT Australia (NICTA)

In collaboration with:

A. Cherian, F. Porikli, P. H. Gosselin, F. Yan, K. Mikolajczyk.

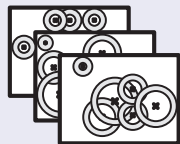
*piotr.koniusz@nicta.com.au*

July 26, 2016

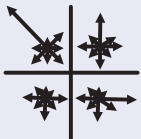
# The Content

- Recap of Bag-of-Words (universal baseline to study).
- Second- and higher-order occurrence pooling.
- Intuitive explanation: uncertainty in max-pooling.
- Evaluations and comparisons to fisher vector encoding.
- Region covariance descriptors.
- Embedding into RKHS+Third-order Super-symmetric Tensor descriptors.
- Details of linearization process.
- Sparse coding for TOSST descriptor.
- Evaluations of TOSST on texture recognition.
- Action recognition from 3D skeletons: Sequence and Dynamics Compatibility Kernels.
- Details of eigenvalue Power Normalisation.
- Results on 3D skeleton action classification.
- Natural Inner Product on Gaussians and Deep Architectures.

## Bag-of-Words (First-order Approaches)



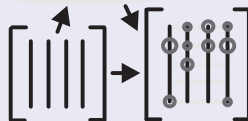
keypoints



descriptors



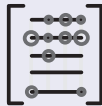
dictionary learning



descriptor coding



pyramid  
matching



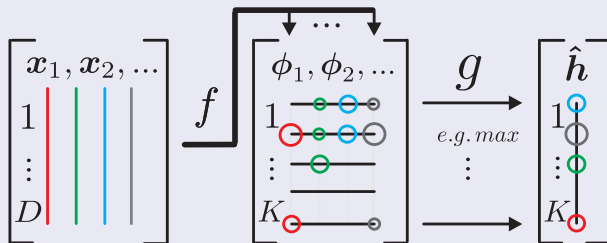
feature pooling



the image  
signature

# First-order Occurrence Pooling

- The local descriptors  $\mathbf{x}$  are extracted from an image and coded by  $f$  that operates on columns.
- Pooling  $g$  aggregates visual words from the mid-level features  $\phi$  along rows:



- Let me remind the following three steps (without Pyramid Matching):

$$\phi_n = f(\mathbf{x}_n, \mathbf{D}), \quad \forall n \in \mathcal{N} \quad (\text{encode}) \quad (1)$$

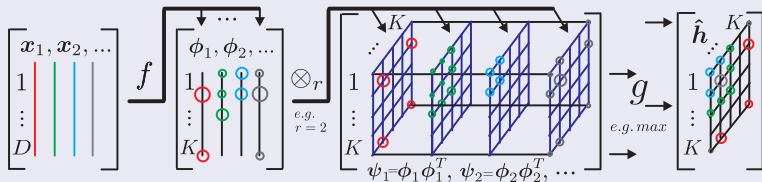
$$\hat{\mathbf{h}}_k = g(\{\phi_{kn}\}_{n \in \mathcal{N}}) \quad (\text{pool}) \quad (2)$$

$$\mathbf{h} = \hat{\mathbf{h}} / \|\hat{\mathbf{h}}\|_2 \quad (\text{normalise}) \quad (3)$$

# Higher-order Occurrence Pooling

- Note that Fisher Vector Encoding and Vector of Locally Aggregated Tensors use the second-order statistics and Power Normalisation.
- BoW can employ the second-order statistics with  $\uparrow \otimes_r$ , e.g.

$$\uparrow \otimes_2 \phi = \phi \phi^T:$$



- Formally, this can be expressed in four steps:

$$\phi_n = f(\mathbf{x}_n, \mathbf{D}), \quad \forall n \in \mathcal{N} \quad (\text{encode}) \quad (4)$$

$$\psi_n = \otimes_r \phi_n \quad (\text{co-occurrences}) \quad (5)$$

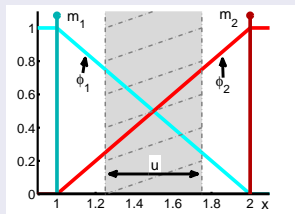
$$\psi_n := u: (\psi_n) \quad (\text{vectorise}) \quad (6)$$

$$\hat{h}_k = g(\{\psi_{kn}\}_{n \in \mathcal{N}}) \quad (\text{pool}) \quad (7)$$

$$\mathbf{h} = \hat{\mathbf{h}} / \|\hat{\mathbf{h}}\|_2 \quad (\text{normalise}) \quad (8)$$

# Uncertainty in max-pooling

- Two linear slopes (LLC) - coding values  $\phi_1$  and  $\phi_2$  for any  $1 \leq x \leq 2$ .
- Draw randomly descriptors from this interval and apply max-pooling.
- If we were to draw several times



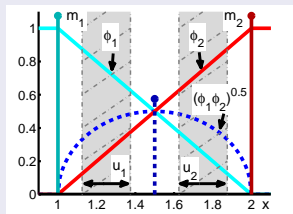
(a)  $x=1.5$ , we would obtain: For (b)  $x_1=1, x_2=2, x_3=1.5$ , we get:

	$\phi_1$	$\phi_2$
0.5	0.5	
...		
0.5	0.5	
max	0.5	0.5

	$\phi_1$	$\phi_2$
0	1	
1	0	
0.5	0.5	
max	1	1

- Position of the descriptor  $x=1.5$  in (a) can be uniquely retrieved from  $\phi$  because  $f^{-1}([0.5, 0.5]^T) = 1.5$ .
- Position of  $x_3=1.5$  in (b) is lost as  $f^{-1}([1, 1]^T) \in [1; 2]$ .

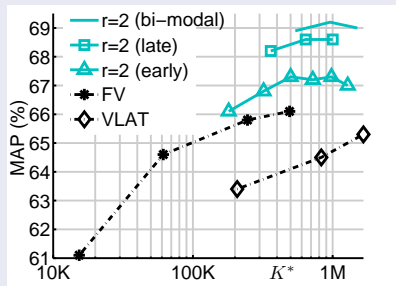
# Uncertainty in max-pooling



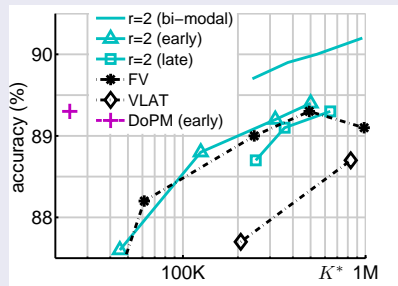
- Two linear slopes (LLC) display coding values  $\phi_1$  and  $\phi_2$  for any  $1 \leq x \leq 2$ .
- Co-occurrences improve on the making effect, e.g. take co-occurrence  $\phi_1\phi_2$ .
- It results in a new maximum for  $x=1.5$  - masking region  $u$  is now split in two smaller regions  $u_1$  and  $u_2$ .

# Second-order Occurrence Pooling: results

## PascalVOC07, Flower102, Opponent SIFT, Sparse Coding



(a) VOC07



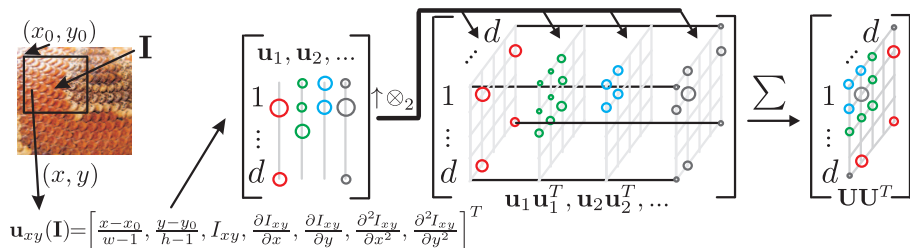
(b) Flower102

- $K^*$  - length of the image signature.
- $r=2$  - Second-order Occurrence Pooling
- FV/VLAT - Fisher Kernels/Vector of Locally Aggregated Tensors.
- Fusions: early - descriptor level, late - kernel level, bi-modal - tensor.
- DoPM - a first-order method.



# TOSST Texture Descriptors

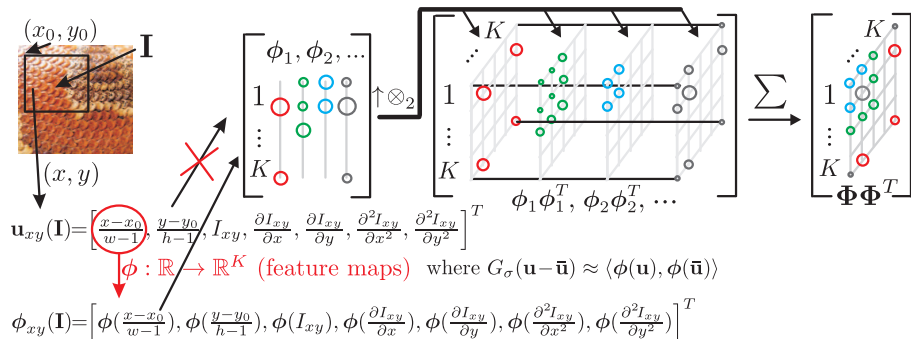
- Region covariance descriptors (co-occurrences).
- They typically use outer-product of low-level feature vectors  
 $\uparrow \otimes_2 \mathbf{u} = \mathbf{u}\mathbf{u}^T$ .



- Low-level features  $\mathbf{u}$  (linear). Non-linear features are better.  
We embed  $\mathbf{u}$  into RKHS/linearise the RBF kernel by feature maps.

# TOSST Texture Descriptors

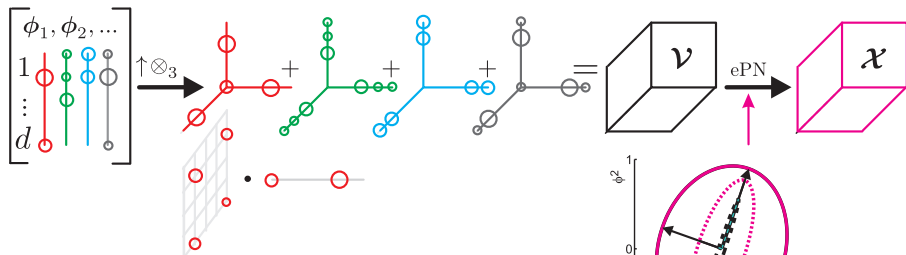
- Non-linear co-occurrence descriptors.



- Can we form more informative co-occurrences?  
 Yes, we extend  $\uparrow \otimes_2$  to third-order outer product  $\uparrow \otimes_3$ .

# TOSST Texture Descriptors

- Non-linear third-order descriptors  
+ eigenvalue Power Normalization (ePN).



- Eigenvalue Power Normalisation prevents correlated signal bursts. Imagine the largest eigenvalue represents the count of pattern of brick and the 2<sup>nd</sup> represents the tree bark. Surely, the amount of brick and bark patterns should not affect the prediction. But it does!
- Higher-order models can be derived analytically.

# Higher-order Occurrence Pooling: derivation

- Assume a kernel, e.g. RBF and its linearisation given by:  
 $ker(\mathbf{u}, \bar{\mathbf{u}}) \approx \langle \phi, \bar{\phi} \rangle$ .
- Assume the dot product  $\langle \phi, \bar{\phi} \rangle$  on a pair of features and polynomial kernel:  $\langle \phi, \bar{\phi} \rangle^r, r \geq 2$ .
- Define a sum kernel between two sets of features  $\mathbf{U} = \{\mathbf{u}_n\}_{n \in \mathcal{N}}$  and  $\bar{\mathbf{U}} = \{\bar{\mathbf{u}}_{\bar{n}}\}_{\bar{n} \in \bar{\mathcal{N}}}$  for two images/regions/sequences (anything you like):

$$\begin{aligned} Ker(\mathbf{U}, \bar{\mathbf{U}}) &= \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \frac{1}{|\bar{\mathcal{N}}|} \sum_{\bar{n} \in \bar{\mathcal{N}}} ker(\mathbf{u}_n, \bar{\mathbf{u}}_{\bar{n}})^r \\ &\approx \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \frac{1}{|\bar{\mathcal{N}}|} \sum_{\bar{n} \in \bar{\mathcal{N}}} \langle \phi_n, \bar{\phi}_{\bar{n}} \rangle^r \\ &= \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \frac{1}{|\bar{\mathcal{N}}|} \sum_{\bar{n} \in \bar{\mathcal{N}}} \left( \sum_{k=1}^K \phi_{kn} \bar{\phi}_{k\bar{n}} \right)^r \end{aligned} \quad (9)$$

# Higher-order Occurrence Pooling: derivation

- The rightmost summation can be re-expressed as a dot product of two outer-products of order  $r$  on  $\phi$ :

$$\left( \sum_{k=1}^K \phi_{kn} \bar{\phi}_{k\bar{n}} \right)^r = \sum_{k^{(1)}=1}^K \dots \sum_{k^{(r)}=1}^K \phi_{k^{(1)}} \bar{\phi}_{k^{(1)}} \cdot \dots \cdot \phi_{k^{(r)}} \bar{\phi}_{k^{(r)}} = \langle \otimes_r \phi_n, \otimes_r \bar{\phi}_{\bar{n}} \rangle_F \quad (10)$$

- Now, the problem is further simplified:

$$\begin{aligned} \text{Ker}(\mathbf{U}, \bar{\mathbf{U}}) &\approx \text{Ker}'(\Phi, \bar{\Phi}) = \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \frac{1}{|\bar{\mathcal{N}}|} \sum_{\bar{n} \in \bar{\mathcal{N}}} \langle \otimes_r \phi_n, \otimes_r \bar{\phi}_{\bar{n}} \rangle_F \\ &= \left\langle \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \otimes_r \phi_n, \frac{1}{|\bar{\mathcal{N}}|} \sum_{\bar{n} \in \bar{\mathcal{N}}} \otimes_r \bar{\phi}_{\bar{n}} \right\rangle_F = \left\langle \text{Avg}_{n \in \mathcal{N}}(\otimes_r \phi_n), \text{Avg}_{\bar{n} \in \bar{\mathcal{N}}}(\otimes_r \bar{\phi}_{\bar{n}}) \right\rangle_F \end{aligned}$$

- We introduce operator  $\mathcal{G}$  (similarity for matrices/tensors, e.g. ePN):

$$\text{Ker}^*(\Phi, \bar{\Phi}) = \left\langle \mathcal{G}\left(\frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \otimes_r \phi_n\right), \mathcal{G}\left(\frac{1}{|\bar{\mathcal{N}}|} \sum_{\bar{n} \in \bar{\mathcal{N}}} \otimes_r \bar{\phi}_{\bar{n}}\right) \right\rangle_F \quad (11)$$

# Sparse Coding for Third-order Tensor Descriptors (TSC)

- However,  $\mathcal{X}$  is cubic w.r.t. size of features.  
We propose Sparse Coding for Third-order Tensor Descriptors.
- We can learn a dictionary to encode TOSST:

$$\arg \min_{\substack{\mathbf{B}_1, \dots, \mathbf{B}_K \\ \alpha^1, \dots, \alpha^N}} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K \mathbf{B}_k \alpha_k^n \right\|_F^2 + \lambda \|\alpha^n\|_1. \quad (12)$$

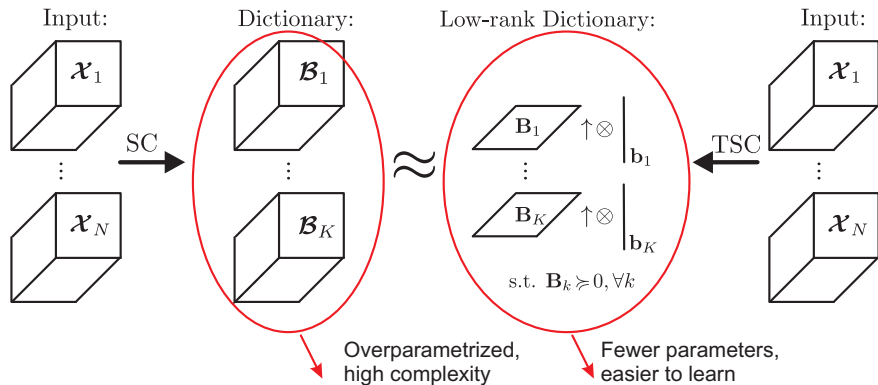
- However,  $\mathbf{B}$  have three modes (overparametrised model), so we learn instead low-rank dictionary:

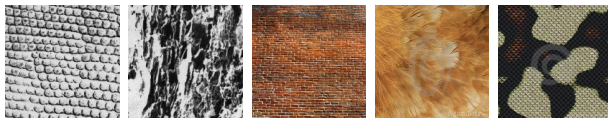
$$\arg \min_{\substack{\mathbf{B}_1, \dots, \mathbf{B}_K \\ \mathbf{b}_1, \dots, \mathbf{b}_K \\ \alpha^1, \dots, \alpha^N}} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K (\mathbf{B}_k \uparrow \otimes \mathbf{b}_k) \alpha_k^n \right\|_F^2 + \lambda \|\alpha^n\|_1. \quad (13)$$

- Resulting sparse codes  $\alpha$  are pooled and used for SVM training.

# Sparse Coding for Third-order Tensor Descriptors (TSC)

- We use training set of TOSST descriptors  $\mathcal{X}_1, \dots, \mathcal{X}_N$ .
- We learn low-rank dictionary atoms  $\mathbf{B}_1 \uparrow \otimes \mathbf{b}_1, \dots, \mathbf{B}_K \uparrow \otimes \mathbf{b}_K$  (outer product of matrices with vectors).
- They approximate full-rank tensor atoms  $\mathcal{B}_1, \dots, \mathcal{B}_K$ .



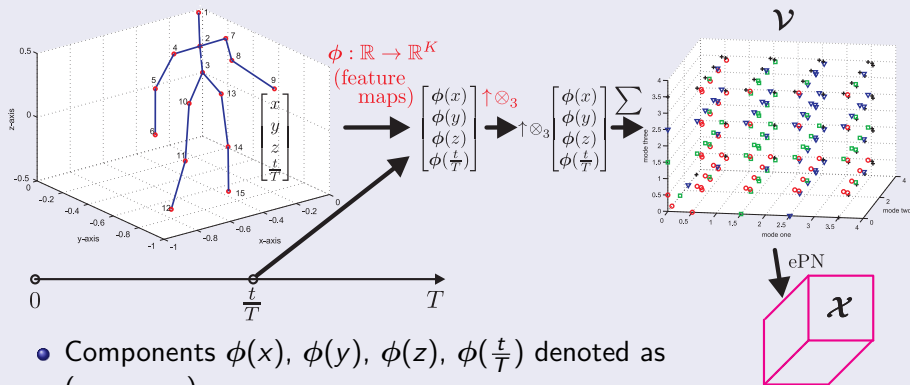


- Brodatz textures; 99.9% accuracy (the state of the art); others score  $\sim 98.72\%$ .
- UIUC materials recognition; 58.0% accuracy.
- PASCAL VOC07 descriptor compression: 61.2% mAP (25K signature) vs. 61.3% mAP (176K signature).



# Action Recognition from 3D Skeletons

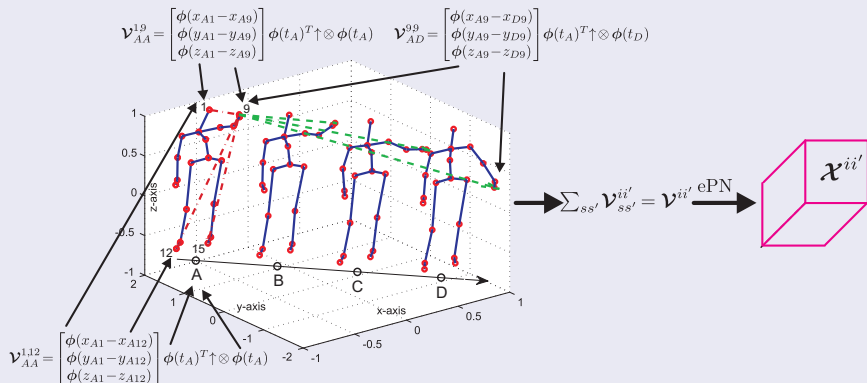
## Sequence Compatibility Kernel



- Components  $\phi(x)$ ,  $\phi(y)$ ,  $\phi(z)$ ,  $\phi(\frac{t}{T})$  denoted as  $(\circ, \square, \nabla, +)$ .
- $\mathcal{V}$  captures all triplets:  $(\circ\square\nabla)$ ,  $(\circ\square+)$ ,  $(\circ\nabla+)$ ,  $(\square\nabla+)$ .
- ePN evens out counts of these co-occurrences.
- Tensors  $\mathcal{X}$  are the samples for training SVM.

# Action Recognition from 3D Skeletons

## Dynamics Compatibility Kernel



- Enumerate all unique joint displacement vectors  $\mathbf{x}_{it} - \mathbf{x}_{jt'}, i \leq j, t \leq t'$ .
- Embed displacements into RKHS and linearise to obtain  $\phi(\mathbf{x}_{it} - \mathbf{x}_{jt'})$ .
- Embed start-/end-times into RKHS, linearise to obtain  $\phi(\frac{t}{T}), \phi(\frac{t'}{T})$ .
- Take outer products  $\phi(\mathbf{x}_{it} - \mathbf{x}_{jt'}) \phi(\frac{t}{T}) \uparrow \otimes \phi(\frac{t'}{T})$ , aggregate+ePN.

- Four simple steps in MATLAB:

$$(\mathcal{E}; \mathbf{A}_1, \dots, \mathbf{A}_r) = \text{HOSVD}(\mathcal{V}) \quad (14)$$

$$\hat{\mathcal{E}} = \text{Sgn}(\mathcal{E}) |\mathcal{E}|^\gamma \quad (15)$$

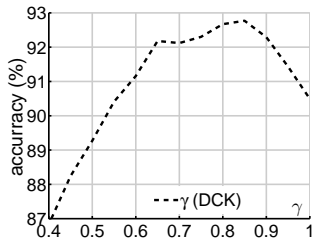
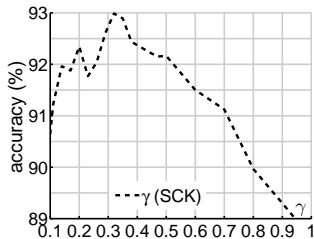
$$\hat{\mathcal{V}} = ((\hat{\mathcal{E}} \otimes_1 \mathbf{A}_1) \dots) \otimes_r \mathbf{A}_r \quad (16)$$

$$\mathcal{X} = \mathcal{G}(\mathcal{V}) = \text{Sgn}(\hat{\mathcal{V}}) |\hat{\mathcal{V}}|^{\gamma^*} \quad (17)$$

- Perform Higher Order SVD (equivalent of SVD for more than 2 modes).
- Obtain the core tensor  $\mathcal{E}$  (equivalent of singular values).
- Power-normalise this spectrum (values  $\mathcal{E}$  can be negative).
- Assemble back tensor, if needed, perform additionally standard PN.

# Action Recognition from 3D Skeletons: results

- Eigenvalue Power Normalisation w.r.t.  $\gamma$ :



# Action Recognition from 3D Skeletons: results

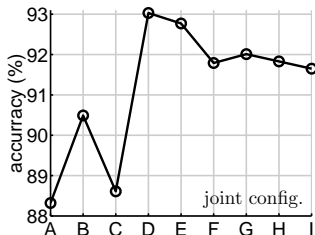
- Florence3D-Action (state-of-the-art):

	SCK	DCK		SCK+DCK
accuracy	92.98%	93.03%	92.77%	<b>95.47%</b>
size	26565	9450	16920	43485

Bag-of-Poses	82.00%	SE(3)	90.88%
--------------	--------	-------	--------

- Using fewer less noisy key-joints may be better (SCK):



# Action Recognition from 3D Skeletons: results

- UTKinect-Action (state-of-the-art):

	SCK	DCK	SCK+DCK
accuracy	96.08%	97.69%	<b>98.39%</b>
size	40480	16920	57400
<hr/>			
3D joints hist.	90.92%	$SE(3)$	97.08%

- MSR-Action3D:

	SCK+DCK	$SE(3)$
accuracy, standard protocol	<b>92.7%</b>	89.48%
accuracy, specific classes/subjects	<b>96%</b>	92.46%
size	57400	-

# Natural Inner Product on Gaussians

- Gaussian kernel between  $\mathbf{u} \in \mathbb{R}^{d'}$  and  $\bar{\mathbf{u}} \in \mathbb{R}^{d'}$  can be simply rewritten as:

$$G_\sigma(\mathbf{u} - \bar{\mathbf{u}}) = e^{-\|\mathbf{u} - \bar{\mathbf{u}}\|_2^2 / 2\sigma^2} = \left( \frac{2}{\pi\sigma^2} \right)^{\frac{d'}{2}} \int_{\zeta \in \mathbb{R}^{d'}} G_{\sigma/\sqrt{2}}(\mathbf{u} - \zeta) G_{\sigma/\sqrt{2}}(\bar{\mathbf{u}} - \zeta) d\zeta. \quad (18)$$

- Finite approximation by  $\zeta_1, \dots, \zeta_Z$  pivots is given by:

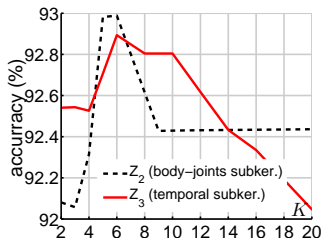
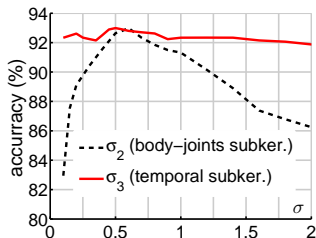
$$\phi(\mathbf{u}) = \left[ G_{\sigma/\sqrt{2}}(\mathbf{u} - \zeta_1), \dots, G_{\sigma/\sqrt{2}}(\mathbf{u} - \zeta_Z) \right]^T, \quad (19)$$

$$\text{and } G_\sigma(\mathbf{u} - \bar{\mathbf{u}}) \approx \left\langle \frac{\phi(\mathbf{u})}{\|\phi(\mathbf{u})\|_2}, \frac{\phi(\bar{\mathbf{u}})}{\|\phi(\bar{\mathbf{u}})\|_2} \right\rangle. \quad (20)$$

- As few as 6 pivots yield  $\leq 0.8\%$  approximation error.

# Action Recognition from 3D Skeletons: results

- Florence3D-Action w.r.t. kernel radii and pivot numbers:

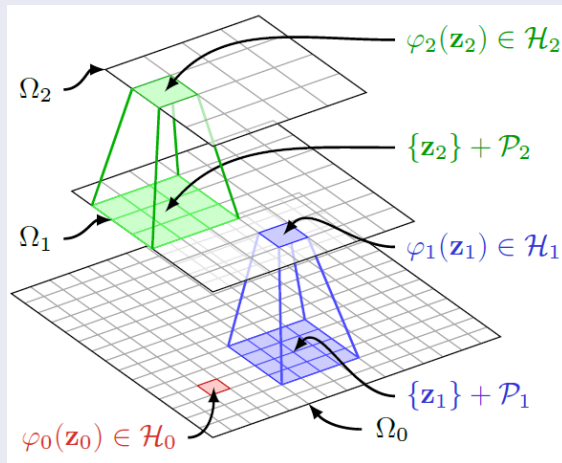


- Not too sensitive to parameter variations.
- More pivots provide better RBF approximation but SVM has to learn more parameters (overfitting).



# Extensions to Deep Architectures

- Convolutional Kernel Networks: future talk.



- The simplest extension - instead of Fisher Vectors apply HOSVD+ePN to aggregate over multiple CNN-based sub-patches.

- It may look difficult, but it requires few easy steps only in practice:
  - A. Extract your favourite features.
  - B. Embed them into RKHS/linearise.
  - C. Form outer products of desired order.
  - D. Aggregate and apply ePN.
  - E. Train SVM (or any favourite classifier).
- Interested in any related ideas? Talk to me to see if we can collaborate :-)
- Ideas take know-how, time and people to develop them.
- References:



P. Koniusz, F. Yan, P. H. Gosselin, K. Mikolajczyk.

Higher-order Occurrence Pooling for Bags-of-Words: Visual Concept Detection, TPAMI, 2016.



P. Koniusz, A. Cherian.

Sparse Coding for Third-order Super-symmetric Tensor Descriptors with Application to Texture Recognition, CVPR, 2016.



P. Koniusz, A. Cherian, F. Porikli.

Tensor Representations via Kernel Linearization for Action Recognition from 3D Skeletons, ECCV, 2016.



J. Mairal, P. Koniusz, Z. Harchaoui, C. Schmid.

Convolutional Kernel Networks, NIPS, 2014.

# Thank You

