

# High-order Tensor Pooling with Attention for Action Recognition

Lei Wang<sup>1,2</sup> Ke Sun<sup>2,1</sup> Piotr Koniusz<sup>2,1</sup>

<sup>1</sup>Australian National University

<sup>2</sup>Data61/CSIRO

April 9, 2024



Australian  
National  
University



# Motivation



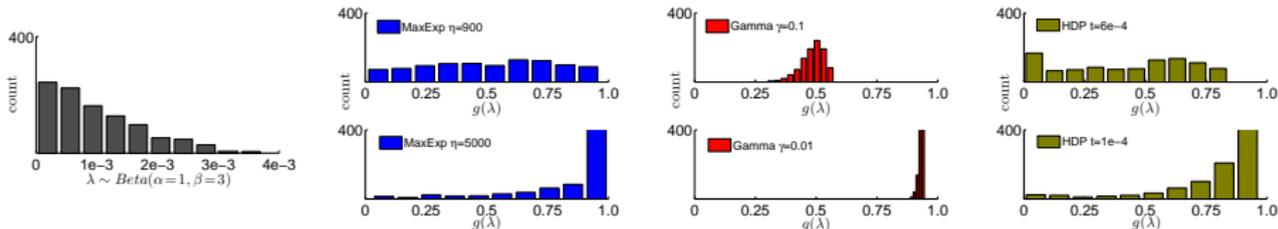
Figure 1: CNN filters respond differently to tree leaf stimuli across spatial regions. **Detecting a leaf reliably predicts a tree's presence.**



Figure 2: Differing **feature counts** challenge classifier generalization. Training with few leaves may lead to misclassification of images with thousands, as boundaries are **sensitive** to observed features.

## Motivation (cont.)

Higher-order representations undergo a non-linearity such as **Power Normalization (PN)**: reduce/boost contributions from frequent/infrequent visual stimuli in an image, respectively.



(a) Initial spectral dist.

(b) MaxExp

(c) Gamma

(d) HDP

**Figure 3:** The intuitive principle of the **Eigenvalue Power Normalization (EPN)**.

- Given a discrete eigenspectrum following a Beta distribution, the pushforward distribution of MaxExp and HDP are very similar.
- For small  $\gamma$ , Gamma is also similar to MaxExp and HDP.
- Note that all three EPN functions **whiten the spectrum** (map the majority of values to be  $\sim 1$ ) thus removing burstiness (acting as a **spectral detector**).
- As EPN prevents burstiness, it replaces counting correlated features with detecting them, thus being invariant to their spatial/temporal extent.

# HoT with EPN

EPN performs a spectrum transformation on  $\mathcal{X} \in \mathbb{R}^{d_1 \times d_2 \dots \times d_r}$ :

$$(\boldsymbol{\lambda}; \mathbf{U}_1, \dots, \mathbf{U}_r) = \text{HOSVD}(\mathcal{X}), \quad (1)$$

$$\hat{\boldsymbol{\lambda}} = g(\boldsymbol{\lambda}), \quad (2)$$

$$\mathcal{G}(\mathcal{X}) = ((\hat{\boldsymbol{\lambda}} \times_1 \mathbf{U}_1) \dots) \times_r \mathbf{U}_r, \quad (3)$$

- Let  $\Phi \equiv \{\phi_1, \dots, \phi_N \in \mathbb{R}^d\}$  be feature vectors extracted from an instance to classify, e.g., video sequences, images, text documents, etc.
- EPN retrieves factors which quantify whether there is **at least one** datapoint  $\phi_n$ ,  $n \in \mathcal{I}_N$ , projected into each subspace spanned by  $r$ -tuples of eigenvector from matrices  $\mathbf{U}_1 = \mathbf{U}_2 = \dots = \mathbf{U}_r$ .
- For brevity, assume order  $r=3$ , a super-symmetric tensor, and **any 3-tuple of eigenvectors**  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  from  $\mathbf{U}$ .
- Note that  $\mathbf{u} \perp \mathbf{v}$ ,  $\mathbf{v} \perp \mathbf{w}$  and  $\mathbf{u} \perp \mathbf{w}$  due to orthogonality of eigenvectors for super-symmetric tensors, e.g.,  $\mathbf{U}\boldsymbol{\lambda}^\ddagger\mathbf{V} = [\boldsymbol{\mathcal{X}}_{::,1}, \dots, \boldsymbol{\mathcal{X}}_{::,d}] \in \mathbb{R}^{d \times d^2}$  where  $\boldsymbol{\lambda}^\ddagger$  are eigenvalues of the **unfolded tensor**  $\mathcal{X}$ .
- If we have  $d$  unique eigenvectors, we can enumerate  $\binom{d}{r}$   $r$ -tuples and thus  $\binom{d}{r}$  subspaces  $\mathbb{R}^{d \times r} \subset \mathbb{R}^{d \times d}$ .

## HoT with EPN (cont.)

Our **super-symmetric tensor descriptor** is:

$$\mathcal{X} = \frac{1}{N} \sum_{n \in \mathcal{I}_N} \uparrow \otimes_r \phi_n, \quad (4)$$

The 'diagonalization' of  $\mathcal{X}$  by eigenvectors  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  produces core tensor:

$$\lambda_{\mathbf{u}, \mathbf{v}, \mathbf{w}} = \mathcal{X} \times_1 \mathbf{u} \times_2 \mathbf{v} \times_3 \mathbf{w}, \quad (5)$$

$\lambda_{\mathbf{u}, \mathbf{v}, \mathbf{w}}$  is a **coefficient** from the core tensor  $\lambda$ . Combining Eq. (4) & (5) yields:

$$\begin{aligned} \lambda_{\mathbf{u}, \mathbf{v}, \mathbf{w}} &= \frac{1}{N} \sum_{n \in \mathcal{I}_N} \uparrow \otimes_3 \phi_n \times_1 \mathbf{u} \times_2 \mathbf{v} \times_3 \mathbf{w} \\ &= \frac{1}{N} \sum_{n \in \mathcal{I}_N} \langle \phi_n, \mathbf{u} \rangle \langle \phi_n, \mathbf{v} \rangle \langle \phi_n, \mathbf{w} \rangle. \end{aligned} \quad (6)$$

- Let  $\phi_n$  be 'optimally' projected into subspace spanned by  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\mathbf{w}$  when  $\psi'_n = \langle \phi_n, \mathbf{u} \rangle \langle \phi_n, \mathbf{v} \rangle \langle \phi_n, \mathbf{w} \rangle$  is **maximized**.
- As our  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  are orthogonal w.r.t. each other and  $\|\phi_n\|_2 = 1$ , simple Lagrange equation  $\mathcal{L} = \sum_{i=1}^r e_i^T \phi_n + \lambda (\|\phi_n\|_2^2 - 1)$  yields **maximum of**  $\kappa = (1/\sqrt{r})^r$  **at**  $\phi_n = [(1/\sqrt{r}), \dots, (1/\sqrt{r})]^T$ .
- For **each**  $n \in \mathcal{I}_N$ , we store  $\psi_n = \psi'_n / \kappa$  in a so-called event vector  $\psi$ .

## HoT with EPN (cont.)

Assume  $\psi \in \{0, 1\}^N$  stores  $N$  outcomes of drawing from Bernoulli distribution under the i.i.d. assumption: the probability  $p$  of an event ( $\psi_n = 1$ ) &  $1-p$  for ( $\psi_n = 0$ ) are estimated by an expected value,  $p = \text{avg}_n \psi_n = \lambda_{\mathbf{u}, \mathbf{v}, \mathbf{w}} / \kappa$  ( $0 \leq \psi \leq 1$ ). The probability of at least one positive event ( $\psi_n = 1$ ) projecting into the subspace spanned by  $r$ -tuples in  $N$  trials is:

$$\hat{\lambda}_{\mathbf{u}, \mathbf{v}, \mathbf{w}} = 1 - (1-p)^N = 1 - \left(1 - \frac{\lambda_{\mathbf{u}, \mathbf{v}, \mathbf{w}}}{\kappa}\right)^N. \quad (7)$$

Each of  $\binom{d}{r}$  subspaces spanned by  $r$ -tuples acts as a detector of projections into this subspace. Eq. (7) is the spectral MaxExp pooling with  $\kappa$  normalization. Considering the dot-product between EPN-norm. tensors  $\mathcal{G}(\mathcal{X})$  and  $\mathcal{G}(\mathcal{Y})$ :

$$\begin{aligned} & \langle \mathcal{G}(\mathcal{X}), \mathcal{G}(\mathcal{Y}) \rangle \\ &= \sum_{\substack{\mathbf{u} \in \mathbf{U}(\mathcal{X}) \\ \mathbf{v} \in \mathbf{V}(\mathcal{X}) \\ \mathbf{w} \in \mathbf{W}(\mathcal{X})}} \sum_{\substack{\mathbf{u}' \in \mathbf{U}(\mathcal{Y}) \\ \mathbf{v}' \in \mathbf{V}(\mathcal{Y}) \\ \mathbf{w}' \in \mathbf{W}(\mathcal{Y})}} \hat{\lambda}_{\mathbf{u}, \mathbf{v}, \mathbf{w}} \hat{\lambda}'_{\mathbf{u}', \mathbf{v}', \mathbf{w}'} \langle \mathbf{u}, \mathbf{u}' \rangle \langle \mathbf{v}, \mathbf{v}' \rangle \langle \mathbf{w}, \mathbf{w}' \rangle. \end{aligned} \quad (8)$$

Hence, all subspaces of  $\mathcal{X}$  and  $\mathcal{Y}$  spanned by  $r$ -tuples (e.g.,  $r = 3$  as above) are compared against each other for alignment by the cosine distance.

## Backpropagating through HOSVD and/or SVD

Let  $\mathbf{M}^\# = \mathbf{M}\mathbf{M}^T = \mathbf{U}\boldsymbol{\lambda}\mathbf{U}^T$  be an SPD matrix with simple eigenvalues, *i.e.*,  $\lambda_{ii} \neq \lambda_{jj}, \forall i \neq j$ . Then  $\mathbf{U}$  coincides also with the eigenvector matrix of tensor  $\boldsymbol{\mathcal{X}}$  for the given unfolding. To compute the derivative of  $\mathbf{U}$  (we drop the index) w.r.t.  $\mathbf{M}$  (and thus  $\boldsymbol{\mathcal{X}}$ ), one has to follow the chain rule:

$$\frac{\partial \mathbf{U}}{\partial M_{kl}} = \sum_{i,j} \frac{\partial \mathbf{U}}{\partial (\mathbf{M}\mathbf{M}^T)_{ij}} \cdot \frac{\partial (\mathbf{M}\mathbf{M}^T)_{ij}}{\partial M_{kl}},$$

where  $\frac{\partial u_{ij}}{\partial \mathbf{M}^\#} = u_{ij}(\lambda_{jj}\mathbb{I} - \mathbf{M}^\#)^\dagger.$  (9)

For SVD, we simply have to backpropagate through the chain rule:

$$\frac{\partial \mathbf{U}\boldsymbol{\lambda}\mathbf{U}^T}{\partial X_{m'n'}} = 2 \text{Sym} \left( \frac{\partial \mathbf{U}}{\partial X_{m'n'}} \boldsymbol{\lambda}\mathbf{U}^T \right) + \mathbf{U} \frac{\partial \boldsymbol{\lambda}}{\partial X_{m'n'}} \mathbf{U}^T,$$

where  $\text{Sym}(\mathbf{X}) = \frac{1}{2}(\mathbf{X} + \mathbf{X}^T).$  (10)

Let  $\mathbf{X} = \mathbf{U}\boldsymbol{\lambda}\mathbf{U}^T$  be an SPD matrix with simple eigenvalues, *i.e.*,  $\lambda_{ii} \neq \lambda_{jj}, \forall i \neq j$ , and  $\mathbf{U}$  contain eigenvectors of matrix  $\mathbf{X}$ , then one can apply  $\frac{\partial \lambda_{ii}}{\partial X} = \mathbf{u}_i \mathbf{u}_i^T$  and  $\frac{\partial u_{ij}}{\partial X} = u_{ij}(\lambda_{jj}\mathbb{I} - \mathbf{X})^\dagger.$

# Application to Action Recognition

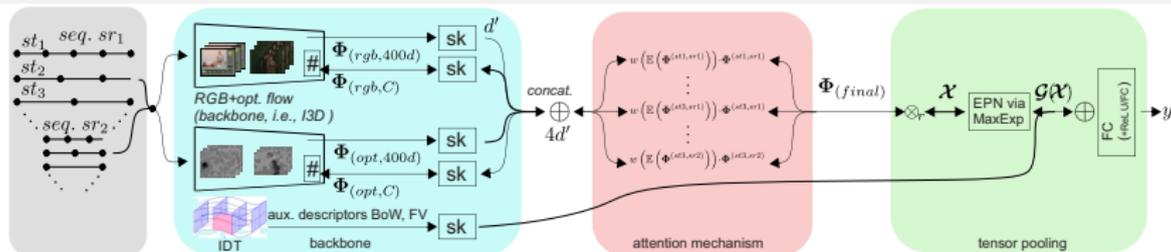


Figure 4: Our action recognition pipeline with the attention mechanism.

## Our pipeline:

- extract subsequences (invariance to action localization)
- apply various sampling rates (invariance to action speed)
- extract 400D features (I3D pretrained on Kinetics-400)
- obtain intermediate matrices with feature vectors
- use count sketching ( $sk$ ) to reduce dimensionality & concatenate features

## Attention mechanism:

- The attention network  $w : \mathbb{R}^{d'} \rightarrow \mathbb{R}$  outputs an attention score
- $\Phi_w^{(i,j)} = w(\mathbb{E}(\Phi^{(i,j)})) \cdot \Phi^{(i,j)}$ ,  $i \in \{st_1, st_2, \dots\}$  &  $j \in \{sr_1, sr_2, \dots\}$
- form final feature matrix  $\Phi_{(final)} \in \mathbb{R}^{d \times N}$ ,  $d = 4d'$ , then passed via Eq. (4).
- pass  $\mathcal{X}$  via EPN to obtain  $\mathcal{G}(\mathcal{X}) \in \mathbb{R}^{d \times d \times d}$ , one per instance to classify

# Results & Discussions

| SO+          | <i>sp1</i> | <i>sp2</i>       | <i>sp3</i> | mean               | TO+        | <i>sp1</i>         | <i>sp2</i> | <i>sp3</i> | mean        |
|--------------|------------|------------------|------------|--------------------|------------|--------------------|------------|------------|-------------|
| (no EPN)     | 76.2       | 75.3             | 76.7       | 76.1               | (no EPN)   | 75.4               | 74.0       | 75.0       | 74.8        |
| HDP          | 81.4       | 78.8             | 80.1       | 80.1               | HDP        | 81.8               | 79.6       | 81.3       | 80.9        |
| MaxExp       | 81.7       | 79.1             | 80.1       | 80.3               | MaxExp     | 82.3               | 79.9       | 81.2       | 81.1        |
| MaxExp+IDT   | 86.1       | 85.2             | 85.8       | <b>85.7</b>        | MaxExp+IDT | 87.4               | 86.7       | 87.5       | <b>87.2</b> |
| ADL+I3D 81.5 |            | Full-FT I3D 81.3 |            | SCK(SO+) +IDT 85.1 |            | SCK(TO+) +IDT 86.1 |            |            |             |

**Table 1:** (*top*) Our model vs. (*bottom*) SOTA on HMDB-51.

|               | <i>static</i> | <i>dynamic</i> | <i>mixed</i> | mean<br>stat/dyn | mean<br>all |
|---------------|---------------|----------------|--------------|------------------|-------------|
| SO+MaxExp     | 92.52         | 82.03          | 89.44        | 87.3             | 88.0        |
| SO+MaxExp+IDT | 94.92         | 86.63          | 96.02        | 90.8             | 92.5        |
| TO+MaxExp+IDT | <b>95.36</b>  | 86.90          | <b>97.04</b> | 91.1             | <b>93.1</b> |
| T-ResNet      | 92.41         | 81.50          | 89.00        | 87.0             | 87.6        |
| ADL I3D       | 95.10         | <b>88.30</b>   | -            | 91.7             | -           |

**Table 2:** (*top*) Our pipeline vs. (*bottom*) SOTA on YUP++.

|               | <i>sp1</i>      | <i>sp2</i> | <i>sp3</i> | <i>sp4</i> | <i>sp5</i>   | <i>sp6</i> | <i>sp7</i> | mAP         |
|---------------|-----------------|------------|------------|------------|--------------|------------|------------|-------------|
| SO+MaxExp+IDT | 75.7            | 82.5       | 79.4       | 75.1       | 75.7         | 76.8       | 75.9       | 77.3        |
| TO+MaxExp+IDT | 78.6            | 83.4       | 81.5       | 78.8       | 81.7         | 79.2       | 79.6       | <b>80.4</b> |
| KRP-FS 70.0   | KRP-FS+IDT 76.1 |            | GRP 68.4   |            | GRP+IDT 75.5 |            |            |             |

**Table 3:** (*top*) Our pipeline vs. (*bottom*) SOTA on MPII.

# Thank you!