



## High-Order Statistics for Skeleton-Based Hand Gesture Recognition

Xuan Son Nguyen<sup>†</sup>, Luc Brun<sup>‡</sup>, Olivier Lezoray<sup>†</sup>, Sébastien Bougleux<sup>‡</sup>

<sup>†</sup> ETIS, Univ. Paris Seine, Univ. Cergy-Pontoise, ENSEA, CNRS, Cergy-Pontoise

<sup>‡</sup> Normandie Univ, ENSICAEN, CNRS, UNICAEN, GREYC, Caen  
France



## Operate on matrices

### Pros

- Richer and more aggregated data

### Cons

- More complex data structure
- Less operators are available
- Local information may be hidden.



## SPD Operators

• *GaussAvg*( $\{P_1, \dots, P_n\}$ ) : From Data to SPD.

$$GaussAvg(\{P_1, \dots, P_n\}) = \begin{bmatrix} \Sigma + \mu\mu^T & \mu^T \\ \mu^T & 1 \end{bmatrix}$$

where  $\Sigma$  is the covariance matrix and  $\mu$  the mean.

- *ReEig*( $X$ )  $\approx$  *ReLu*

$$ReEig(X) = U \max(\epsilon I, \Lambda) U^T \text{ where } X = U \Lambda U^T$$

- *LogEig*( $X$ ): Preprocessing to the new operator.

$$logEig(X) = U \log(\Lambda) U^T \text{ where } X = U \Lambda U^T$$

- *VecMat*( $X$ ): From SPD to vectors

$$VecMat(X) = [X_{1,1}, \sqrt{2}X_{1,2}, \dots, \sqrt{2}X_{1,n}, X_{2,2}, \sqrt{2}X_{2,3} \dots \sqrt{2}X_{2,n}, X_{3,3} \dots]$$

- *vl*( $X$ )

$$vl(X) = VecMat(logEig(ReEig(X)))$$



## SPD operators: a last one

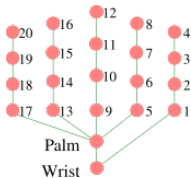
- $SPDAgg(X_1, \dots, X_n)$ : Attention network

$$SPDAgg(X_1, \dots, X_n) = \sum_{i=1}^n W_i X_i W_i^T$$

Backpropagation of  $W_i$  should insure that  $SPDAgg(x_1, \dots, X_n)$  is SPD.



## Inputs



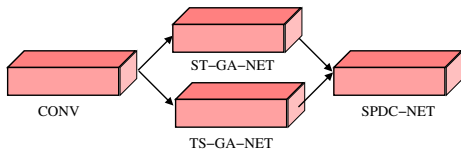
(a) Initial Graph

20	16	12	8	4
19	15	11	7	3
18	14	10	6	2
17	13	9	5	1

(b) Image encoding of joints. Each joint as a dim equal to 3.



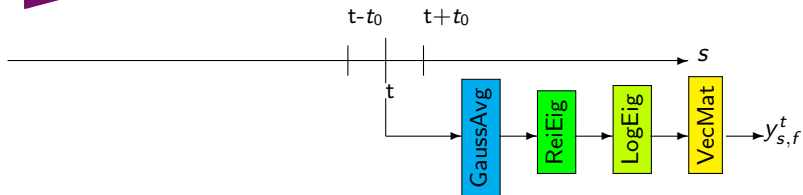
## Our convolutional network



- Conv: A classical 2D Image convolution.
- ST a spatio-temporal network
- TS a temporal-spatial network
- SPDC : a fusion step.



## ST network

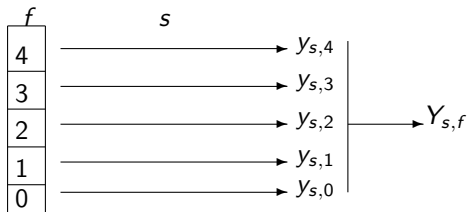


- 1 For a sequence  $s$ , each finger  $f$ , each time step  $t$  we compute

$$Y_{s,f}^t = \text{GaussAvg}(\{p_{s,j,i} \text{ with } j \in J_f, i \in [t - t_0, t + t_0]\}),$$

$J_f$  set of joints of finger  $f$ .

- 2 we compute  $y_{s,f}^t = \text{vl}(Y_{s,f}^t) = \text{VecMat}(\text{LogEig}(\text{ReEig}(Y_{s,f}^t)))$
- 3 We compute  $Y_{s,f} = \text{GaussAvg}(\{y_{s,f}^t, t \in s\})$



- 1 For each joint  $j$  let:

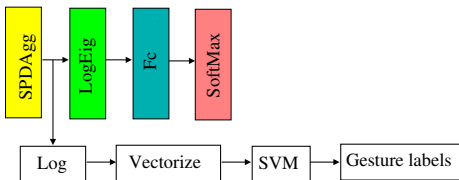
$$Y_{s,j} = \text{GaussAvg}(\{p_{s,j}^t, t \in s\})$$

- 2 Let  $y_{s,j} = vl(Y_{s,j})$
- 3 compute for each finger  $f$   $Y_{s,f} = \text{GaussAvg}(\{y_{s,j}, j \in J_f\})$





- for ST the sequence  $s$  is subdivided into 6 sub-sequences hence providing  $5 \times 6 = 30$  SPD matrices.
- For TS  $s$  is subdivided into 20 subsequences of equal length providing  $5 \times 20 = 100$  SPD matrices.
- we apply SPDAgg to merge all these matrices into a single one.



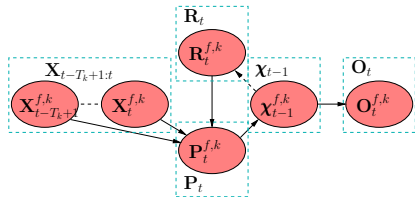


## A recurrent Network

- Main Aim: Better take into account the time dimension.



Statistical Recurrent Unit (SRU):



$$P_t^{f,k} = \text{ReEig} \left( \frac{(w_p^k)^2}{(w_p^k)^2 + (w_x^k)^2} R_t^{f,k} + \frac{(w_x^k)^2}{(w_p^k)^2 + (w_x^k)^2} h^k(X_t^f) \right) \quad (1)$$

Finger  $f$ , time  $t$ , statistics of order  $k$ .

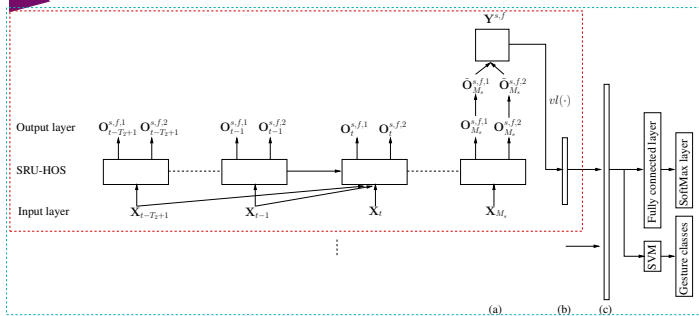
$$h^1(X_t^f) = \begin{bmatrix} \Sigma_t^{f,1} + \mu_t^{f,1}(\mu_t^{f,1})^T & \mu_t^{f,1} \\ (\mu_t^{f,1})^T & 1 \end{bmatrix} \quad \text{over the interval } [t - t_1, t].$$

$$h^2(X_t^f) = \text{GaussAvg}(vl(h^1(X_t^f))) \quad \text{over the interval } [t - t_2, t].$$

$$\forall (f, k) \in \{1, \dots, 5\} \times \{1, 2\} \quad O_t^{f,k} = \text{SRU\_HOS}_k(X_t^f)$$



# Our recurrent network



For all  $(s, f) \in \{1, \dots, 6\} \times \{1, \dots, 5\}$ :

$$Y^{s,f} = \begin{bmatrix} \tilde{O}_{M_s}^{s,f,2} + \text{vl}(\tilde{O}_{M_s}^{s,f,1})\text{vl}(\tilde{O}_{M_s}^{s,f,1})^T & \text{vl}(\tilde{O}_{M_s}^{s,f,1}) \\ \text{vl}(\tilde{O}_{M_s}^{s,f,1})^T & 1 \end{bmatrix},$$

Global representation:

$$[\text{vl}(Y^{1,1})^T, \dots, \text{vl}(Y^{6,5})^T]^T$$



## Experiments : Ablation study

- Recognition accuracy (%) of sub-networks ST-GA-NET and TS-GA-NET.

Network	FPHA	DHG (14 gestures)	DHG (28 gestures)
ST-HGR-NET	91.83	93.21	89.29
TS-HGR-NET	90.96	93.33	88.21
ST-TS-HGR-NET	<b>93.22</b>	<b>94.29</b>	<b>89.40</b>

- Relevance of  $h^1()$  and  $h^2()$  statistics

Statistics	DHG (14 gestures)	DHG (28 gestures)	FPHA
only $h^1(.)$	85.00	76.43	77.04
only $h^2(.)$	89.29	86.07	93.57
<b>Full</b>	<b>94.4</b>	<b>89.52</b>	<b>94.61</b>

- # of parameters

Model	Number of parameters
ST-TS-HGR-NET	672,243
SRU-HOS-NET	<b>18,894</b>



- Performance of our method and state-of-the-art methods on DHG dataset.

Method	Year	Color	Depth	Pose	RNN/LSTM	Accuracy (%)	
						14 gestures	28 gestures
HON4D [Oreifej and Liu, 2013]	2013	X	✓	X	X	78.53	74.03
Devanne et al. [Devanne et al., 2015]	2015	X	X	✓	X	79.61	62.00
Huang et al. [Huang and Gool, 2017]	2017	X	X	✓	X	75.24	69.64
De Smedt et al. [Smedt et al., 2016]	2016	X	X	✓	X	88.24	81.90
Devineau et al. [Devineau et al., 2018]	2018	X	X	✓	X	91.28	84.35
SRU [Oliva et al., 2017]	2018	X	X	✓	✓	82.02	76.31
SRU-SPD [Chakraborty et al., 2018]	2018	X	X	✓	✓	86.31	80.83
ST-TS-HGR-NET [Nguyen et al., 2019]	2019	X	X	✓	X	94.29	89.40
SRU-HOS-NET		X	X	✓	✓	<b>94.40</b>	<b>89.52</b>



## Experiments: Comparison with state of the art

FPHA dataset.

Method	Year	Color	Depth	Pose	RNN/LSTM	Accuracy (%)
HON4D [Oreifej and Liu, 2013]	2013	X	✓	X	X	70.61
Novel View [Rahmani and Mian, 2016]	2016	X	✓	X	X	69.21
1-layer LSTM [Zhu et al., 2016]	2016	X	X	✓	✓	78.73
2-layer LSTM [Zhu et al., 2016]	2016	X	X	✓	✓	80.14
Moving Pose [Zanfir et al., 2013]	2013	X	X	✓	X	56.34
Lie Group [Vemulapalli et al., 2014]	2014	X	X	✓	X	82.69
HBRNN [Du et al., 2015]	2015	X	X	✓	✓	77.40
Gram Matrix [Zhang et al., 2016]	2016	X	X	✓	X	85.39
TF [Garcia-Hernando and Kim, 2017]	2017	X	X	✓	X	80.69
JOULE-color [Hu et al., 2015]	2015	✓	X	X	X	66.78
JOULE-depth [Hu et al., 2015]	2015	X	✓	X	X	60.17
JOULE-pose [Hu et al., 2015]	2015	X	X	✓	X	74.60
JOULE-all [Hu et al., 2015]	2015	✓	✓	✓	X	78.78
Huang et al. [Huang and Gool, 2017]	2017	X	X	✓	X	84.35
Huang et al. [Huang et al., 2018]	2018	X	X	✓	X	77.57
SRU [Oliva et al., 2017]	2018	X	X	✓	✓	72.17
SRU-SPD [Chakraborty et al., 2018]	2018	X	X	✓	✓	78.96
ST-TS-HGR-NET [Nguyen et al., 2019]	2019	X	X	✓	X	93.22
SRU-HOS-NET		X	X	✓	✓	<b>94.61</b>



## Conclusion

- Neural Network on SPD manifold provide promising results.
- How to generalize it to arbitrary graphs ?
  - We need SPD matrices (not a big deal but to study anyway),
  - We need meaningful subgraphs to capture local information (much harder).





Questions ?





# Bibliography I



Chakraborty, R., Yang, C.-H., Zhen, X., Banerjee, M., Archer, D., Vaillancourt, D. E., Singh, V., and Vemuri, B. C. (2018).

A Statistical Recurrent Model on the Manifold of Symmetric Positive Definite Matrices. In *NeurIPS*, pages 8897–8908.



Devanne, M., Wannous, H., Berretti, S., Pala, P., Daoudi, M., and Bimbo, A. D. (2015). 3-D Human Action Recognition by Shape Analysis of Motion Trajectories on Riemannian Manifold.

*IEEE Transactions on Cybernetics*, 45(7):1340–1352.



Devineau, G., Moutarde, F., Xi, W., and Yang, J. (2018).

Deep Learning for Hand Gesture Recognition on Skeletal Data.

In *IEEE International Conference on Automatic Face Gesture Recognition*, pages 106–113.



Du, Y., Wang, W., and Wang, L. (2015).

Hierarchical Recurrent Neural Network for Skeleton Based Action Recognition.

In *CVPR*, pages 1110–1118.



Garcia-Hernando, G. and Kim, T.-K. (2017).

Transition Forests: Learning Discriminative Temporal Transitions for Action Recognition.

In *CVPR*, pages 407–415.



## Bibliography II



Hu, J., Zheng, W., Lai, J., and Zhang, J. (2015).  
Jointly Learning Heterogeneous Features for RGB-D Activity Recognition.  
In *CVPR*, pages 5344–5352.



Huang, Z. and Gool, L. V. (2017).  
A Riemannian Network for SPD Matrix Learning.  
In *AAAI*, pages 2036–2042.



Huang, Z., Wu, J., and Gool, L. V. (2018).  
Building Deep Networks on Grassmann Manifolds.  
In *AAAI*, pages 3279–3286.



Nguyen, X., Brun, L., Lézoray, O., and Bougleux, S. (2019).  
A Neural Network Based on SPD Manifold Learning for Skeleton-based Hand Gesture Recognition.  
In *CVPR*.



Oliva, J. B., Póczos, B., and Schneider, J. (2017).  
The Statistical Recurrent Unit.  
In *ICML*, pages 2671–2680.



## Bibliography III



Oreifej, O. and Liu, Z. (2013).

HON4D: Histogram of Oriented 4D Normals for Activity Recognition from Depth Sequences.

In *CVPR*, pages 716–723.



Rahmani, H. and Mian, A. (2016).

3D Action Recognition from Novel Viewpoints.

In *CVPR*, pages 1506–1515.



Smedt, Q. D., Wannous, H., and Vandeborre, J. (2016).

Skeleton-Based Dynamic Hand Gesture Recognition.

In *CVPRW*, pages 1206–1214.



Vemulapalli, R., Arrate, F., and Chellappa, R. (2014).

Human Action Recognition by Representing 3D Skeletons as Points in a Lie Group.

In *CVPR*, pages 588–595.



Zanfir, M., Leordeanu, M., and Sminchisescu, C. (2013).

The Moving Pose: An Efficient 3D Kinematics Descriptor for Low-Latency Action Recognition and Detection.

In *ICCV*, pages 2752–2759.



## Bibliography IV



Zhang, X., Wang, Y., Gou, M., Sznaiier, M., and Camps, O. (2016).

Efficient Temporal Sequence Comparison and Classification Using Gram Matrix Embeddings on a Riemannian Manifold.

In *CVPR*, pages 4498–4507.



Zhu, W., Lan, C., Xing, J., Zeng, W., Li, Y., Shen, L., and Xie, X. (2016).

Co-occurrence Feature Learning for Skeleton Based Action Recognition Using Regularized Deep LSTM Networks.

In *AAAI*, pages 3697–3703.