

Human Motion Signatures: Analysis, Synthesis, Recognition

M. Alex O. Vasilescu

Department of Computer Science, University of Toronto
Toronto, ON M5S 3G4, Canada

Abstract

Human motion is the composite consequence of multiple elements, including the action performed and a motion signature that captures the distinctive pattern of movement of a particular individual. We develop a new algorithm that is capable of extracting these motion elements and recombining them in novel ways. The algorithm analyzes motion data spanning multiple subjects performing different actions. The analysis yields a generative motion model that can synthesize new motions in the distinctive styles of these individuals. Our algorithms can also recognize people and actions from new motions by comparing motion signatures and action parameters.

1. Introduction and Background

In analogy with handwritten signatures, do people have characteristic motion signatures that individualize their movements? If so, can these signatures be extracted from example motions? Can extracted signatures be used to recognize, say, a particular individual's walk subsequent to observing examples of other movements produced by this individual?

The ability to perceive motion signatures seems well-grounded from an evolutionary perspective, since survival depends on recognizing the movements of predator or prey, or of friend or foe. In the 1960s, the psychologist Gunnar Johansson performed a series of famous experiments in which he attached lights to people's limbs and recorded videos of them performing different activities, such as walking, running, and dancing [4]. Observers of these moving light displays, videos in which only the lights are visible, were asked to classify the activity performed and to note certain characteristics of the movements, such as a limp or an energetic/tired walk. Observers can usually perform this task with ease and they could sometimes determine gender and even recognize specific individuals in this way. This may corroborate the hypothesis that the motion signature is a perceptible element of human motion.

Our research [11, 12] has three goals. The first is to model human motions as the composite consequence of multiple elements—the action performed and a motion signature. The second is to determine if people have motion

signatures that are invariant of action classes. Therefore, we extract a motion signature from a subset of actions for a new individual and synthesize the remainder of the actions using the extracted motion signature. The synthetic motions are then validated by classifying against a database of all the real motions. Once our motion model has been validated our third goal is to recognize specific individuals and actions. Our algorithm exploits corpora of motion data which are now reasonably easy to acquire through a variety of modern motion capture technologies developed for use in the entertainment industry [3]. Motion synthesis through the analysis of motion capture data is currently attracting a great deal of attention within the computer graphics community as a means of animating graphical characters. Several authors have introduced generative motion models for this purpose. Recent papers report the use of hidden Markov models [1]. and neural network learning models [2].

We address the motion analysis/synthesis/recognition problem using techniques from numerical statistics. The mathematical basis of our approach is a technique known as *n-mode analysis*, which was first proposed by Tucker [10] and subsequently developed by Kapteyn *et al.* [7, 8], among others. This multilinear analysis subsumes as special cases the simple, linear (1-factor) analysis associated with conventional SVD and principal components analysis (PCA), as well as the incrementally more general bilinear (2-factor) analysis that has recently been investigated in the context of computer vision [9]. Subsuming conventional linear analysis as a special case, multilinear analysis emerges as a unifying mathematical framework suitable for addressing a variety of computer vision problems [14].

Within our framework, corpora of motion capture data spanning multiple people and actions are best organized as higher-order arrays or tensors which define multilinear operators over a *set* of vector spaces. Unlike the matrix case for which the existence and uniqueness of the singular value decomposition (SVD) is assured, the situation for higher-order tensors is not as simple. There are multiple ways to orthogonally decompose tensors [5]. However, one multilinear extension of the matrix SVD to tensors is most natural. We apply this *N-mode SVD* to extract human motion signatures among the other constitutive factors inherent to human movement.

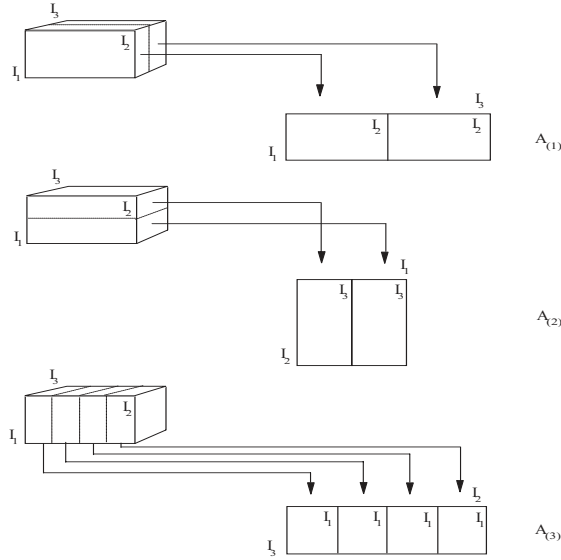


Figure 1: Flattening a (3rd-order) tensor. The tensor can be flattened in 3 ways to obtain matrices comprising its mode-1, mode-2, and mode-3 vectors.

2. Tensors and Decomposition

A *tensor* is a higher order generalization of a vector (first order tensor) and a matrix (second order tensor). Tensors are multilinear mappings over a set of vector spaces. The *order* of tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is N . An element of \mathcal{A} is denoted as $A_{i_1 \dots i_n \dots i_N}$ or $a_{i_1 \dots i_n \dots i_N}$ or where $1 \leq i_n \leq I_n$.¹ In tensor terminology, column vectors are referred to as mode-1 vectors and row vectors as mode-2 vectors. The mode- n vectors of an N^{th} order tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ are the I_n -dimensional vectors obtained from \mathcal{A} by varying index i_n while keeping the other indices fixed. The mode- n vectors are the column vectors of matrix $\mathbf{A}_{(n)} \in \mathbb{R}^{I_n \times (I_1 I_2 \dots I_{n-1} I_{n+1} \dots I_N)}$ that results from *flattening* the tensor \mathcal{A} , as shown in Fig. 1.

A generalization of the product of two matrices is the product of a tensor and a matrix. The *mode- n product* of a tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_n \times \dots \times I_N}$ by a matrix $\mathbf{M} \in \mathbb{R}^{J_n \times I_n}$, denoted by $\mathcal{A} \times_n \mathbf{M}$, is a tensor $\mathcal{B} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J_n \times I_{n+1} \times \dots \times I_N}$, whose entries are $B_{i_1 \dots i_{n-1} j_n i_{n+1} \dots i_N} = \sum_{i_n} a_{i_1 \dots i_{n-1} i_n i_{n+1} \dots i_N} m_{j_n i_n}$. The mode- n product can be expressed in terms of flattened matrices as $\mathbf{B}_{(n)} = \mathbf{M} \mathbf{A}_{(n)}$.²

¹We denote scalars by lower case letters (a, b, \dots), vectors by bold lower case letters ($\mathbf{a}, \mathbf{b}, \dots$), matrices by bold upper-case letters ($\mathbf{A}, \mathbf{B}, \dots$), and higher-order tensors by calligraphic upper-case letters ($\mathcal{A}, \mathcal{B}, \dots$).

²The mode- n product of a tensor and a matrix is a special case of the inner product in multilinear algebra and tensor analysis. Note that for tensors and matrices of the appropriate sizes, $\mathcal{A} \times_m \mathbf{U} \times_n \mathbf{V} = \mathcal{A} \times_n \mathbf{V} \times_m \mathbf{U}$

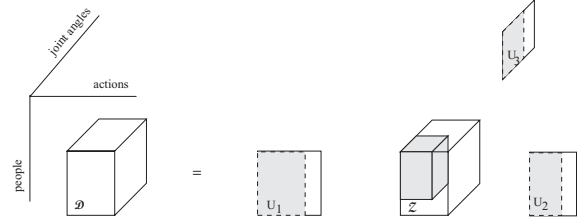


Figure 2: An N -mode SVD orthogonalizes the N vector spaces associated with an order- N tensor (the case $N = 3$ is illustrated).

A matrix $\mathbf{D} \in \mathbb{R}^{I_1 \times I_2}$ is a two-mode mathematical object that has two associated vector spaces, a row space and a column space. SVD orthogonalizes these two spaces and decomposes the matrix as $\mathbf{D} = \mathbf{U}_1 \mathbf{\Sigma} \mathbf{U}_2^T$, the product of an orthogonal column-space represented by the left matrix $\mathbf{U}_1 \in \mathbb{R}^{I_1 \times J_1}$, a diagonal singular value matrix $\mathbf{\Sigma} \in \mathbb{R}^{J_1 \times J_2}$, and an orthogonal row space represented by the right matrix $\mathbf{U}_2 \in \mathbb{R}^{I_2 \times J_2}$. In terms of the mode- n products defined above, this matrix product can be rewritten as $\mathbf{D} = \mathbf{\Sigma} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2$.

By extension, an order $N > 2$ tensor \mathcal{D} is an N -dimensional matrix comprising N spaces. “ N -mode SVD” is a “generalization” of SVD that orthogonalizes these N spaces and decomposes the tensor as the mode- n product of N -orthogonal spaces³

$$\mathcal{D} = \mathcal{Z} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_n \mathbf{U}_n \dots \times_N \mathbf{U}_N, \quad (1)$$

as illustrated in Fig. 2 for the case $N = 3$. Tensor \mathcal{Z} , known as the *core tensor*, is analogous to the diagonal singular value matrix in conventional SVD. It is important to realize, however, that the core tensor does not have a diagonal structure; rather, \mathcal{Z} is in general a full tensor [5]. The core tensor governs the interaction between the *mode matrices* \mathbf{U}_n , for $n = 1, \dots, N$. Mode matrix \mathbf{U}_n contains the orthonormal vectors spanning the column space of the matrix $\mathbf{D}_{(n)}$ that results from the mode- n flattening of \mathcal{D} , as was illustrated in Fig. 1.

Our N -mode SVD algorithm for decomposing \mathcal{D} is:

1. For $n = 1, \dots, N$, compute matrix \mathbf{U}_n in (1) by computing the SVD of the flattened matrix $\mathbf{D}_{(n)}$ and setting \mathbf{U}_n to be the left matrix of the SVD.⁴

and $(\mathcal{A} \times_n \mathbf{U}) \times_n \mathbf{V} = \mathcal{A} \times_n (\mathbf{V} \mathbf{U})$.

³A matrix representation of the N -mode SVD can be obtained by: $\mathbf{D}_{(n)} = \mathbf{U}_n \mathbf{Z}_{(n)} (\mathbf{U}_{n-1} \otimes \dots \otimes \mathbf{U}_1 \otimes \mathbf{U}_N \otimes \dots \otimes \mathbf{U}_{n+2} \otimes \mathbf{U}_{n+1})^T$, where \otimes is the matrix Kronecker product.

⁴When $\mathbf{D}_{(n)}$ is a non-square matrix, the computation of \mathbf{U}_n in the singular value decomposition $\mathbf{D}_{(n)} = \mathbf{U}_n \mathbf{\Sigma} \mathbf{V}_n^T$ can be performed efficiently, depending on which dimension of $\mathbf{D}_{(n)}$ is smaller, by decomposing either $\mathbf{D}_{(n)} \mathbf{D}_{(n)}^T = \mathbf{U}_n \mathbf{\Sigma}^2 \mathbf{U}_n^T$ and then computing $\mathbf{V}_n^T = \mathbf{\Sigma}^+ \mathbf{U}_n^T \mathbf{D}_{(n)}$ or by decomposing $\mathbf{D}_{(n)}^T \mathbf{D}_{(n)} = \mathbf{V}_n \mathbf{\Sigma}^2 \mathbf{V}_n^T$ and then computing $\mathbf{U}_n = \mathbf{D}_{(n)} \mathbf{V}_n \mathbf{\Sigma}^+$.

2. Solve for the core tensor as follows

$$\mathcal{Z} = \mathcal{D} \times_1 \mathbf{U}_1^T \times_2 \mathbf{U}_2^T \dots \times_n \mathbf{U}_n^T \dots \times_N \mathbf{U}_N^T. \quad (2)$$

3. Analysis

Given motion sequences of several people, we define a data set \mathcal{D} which takes the form of a $\mathbb{R}^{N \times M \times T}$ tensor, where N is the number of people, M is the number of action classes, and T is the number of joint angle time samples. We apply the N -mode SVD algorithm given at the end of the previous section to decompose this tensor as follows:

$$\mathcal{D} = \mathcal{Z} \times_1 \mathbf{P} \times_2 \mathbf{A} \times_3 \mathbf{J}, \quad (3)$$

into the product of a core tensor \mathcal{Z} , and three orthogonal matrices. The people matrix $\mathbf{P} = [\mathbf{p}_1 \dots \mathbf{p}_n \dots \mathbf{p}_N]^T$, whose person specific row vectors \mathbf{p}_n^T span the space of person parameters, encodes the per-person invariances across actions. Thus \mathbf{P} contains the human motion signatures. The action matrix $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_m \dots \mathbf{a}_M]^T$, whose action specific row vectors \mathbf{a}_n^T span the space of action parameters, encodes the invariances for each action across different people. The joint angle matrix \mathbf{J} whose row vectors span the space of joint angles are the *eigenmotions* that are normally computed by PCA.

The product $\mathcal{Z} \times_3 \mathbf{J}$ transforms the eigenmotions into a tensorial representation of the variation and co-variation of modes (people and action classes) and characterizes how people parameters and action parameters interact with each other. The tensor

$$\mathcal{B} = \mathcal{Z} \times_2 \mathbf{A} \times_3 \mathbf{J} \quad (4)$$

contains a set of basis matrices for all the motions associated with particular actions. The tensor

$$\mathcal{C} = \mathcal{Z} \times_1 \mathbf{P} \times_3 \mathbf{J} \quad (5)$$

contains a set of basis matrices for all the motions associated with particular people.

4. Synthesis

By performing the decomposition (3), our motion synthesis algorithm first analyzes a corpus of motion data \mathcal{D} for a group of subjects to extract \mathcal{Z} , \mathbf{A} , and \mathbf{J} . This analysis defines a *generative model* that can observe motion data $\mathcal{D}_{p,a}$ of a new subject performing one of these actions (action a) and synthesize the remaining actions, which were never before seen, for this new individual. The algorithm solves for the signature \mathbf{p} of the new individual in the equation $\mathcal{D}_{p,a} = \mathcal{B}_a \times_1 \mathbf{p}^T$, where $\mathcal{B}_a = \mathcal{Z} \times_2 \mathbf{a}_a^T \times_3 \mathbf{J}$. Note that $\mathcal{D}_{p,a}$ is a $1 \times 1 \times T$ tensor. Flattening this tensor in the people mode yields the matrix $\mathbf{D}_{p,a(\text{people})}$, actually

a row vector which we can denote as \mathbf{d}_a^T . Therefore, in terms of flattened tensors, the above equation can be written $\mathbf{d}_a^T = \mathbf{p}^T \mathbf{B}_{a(\text{people})}$ or $\mathbf{d}_a = \mathbf{B}_{a(\text{people})}^T \mathbf{p}$. A complete set of motions for the new individual is synthesized as follows:

$$\mathcal{D}_p = \mathcal{B} \times_1 \mathbf{p}^T, \quad (6)$$

where \mathcal{B} is defined in (4) and the motion signature for the new individual is given by $\mathbf{p}^T = \mathbf{d}_a^T \mathbf{B}_{a(\text{people})}^{-1}$. If several different actions \mathbf{d}_{a_k} are observed, the motion signature is computed as follows:

$$\mathbf{p}^T = [\dots \mathbf{d}_{a_k}^T \dots] \begin{bmatrix} \vdots \\ \mathbf{B}_{a_k(\text{people})}^{-1} \\ \vdots \end{bmatrix}. \quad (7)$$

Similarly, if we observe a known person (one who is already recorded in the motion database) performing a new type of action \mathbf{d}_p , we can compute the associated action parameters $\mathbf{a}^T = \mathbf{d}_p^T \mathbf{C}_{p(\text{actions})}^{-1}$ and use them to synthesize that new action for all the people in the database as follows: $\mathcal{D}_a = \mathcal{C} \times_2 \mathbf{a}^T$, where \mathcal{C} is given in (5). If several different people are observed performing the same new action \mathbf{d}_{p_k} , the action parameters are computed as follows:

$$\mathbf{a}^T = [\dots \mathbf{d}_{p_k}^T \dots] \begin{bmatrix} \vdots \\ \mathbf{C}_{p_k(\text{actions})}^{-1} \\ \vdots \end{bmatrix}. \quad (8)$$

5. Recognition

Multilinear analysis yields basis tensors that map observed motions either into the space of people parameters or the space of action parameters, thereby enabling the recognition of actions or people from motion data.

To recognize the identity of an unknown person from motion data \mathbf{d} of a known action a , we map the motion into the people signature space, by computing the projection $\mathbf{p} = \mathbf{B}_{a(\text{people})}^{-T} \mathbf{d}$. Our nearest neighbor recognition algorithm compares this signature against the person-specific signatures \mathbf{p}_n in \mathbf{P} . The best matching signature vector \mathbf{p}_p —i.e., the one that yields the smallest value of $\|\mathbf{p} - \mathbf{p}_n\|$ among all the people $n = 1, \dots, N$ —recognizes the motion \mathbf{d} as having been produced by person p .

Similarly, to recognize the action depicted in motion data \mathbf{d} generated by a known person p , we map the motion into the action parameter space, by computing the projection $\mathbf{a} = \mathbf{C}_{p(\text{actions})}^{-T} \mathbf{d}$. Our nearest neighbor recognition algorithm compares \mathbf{a} against the action parameter vectors \mathbf{a}_m in \mathbf{A} . The best matching action parameter vector \mathbf{a}_a —i.e., the one that yields the smallest value of $\|\mathbf{a} - \mathbf{a}_m\|$ among all the actions $m = 1, \dots, M$ —recognizes the motion \mathbf{d} as depicting action a .



Figure 3: The motion capture facility

Comparing our multilinear technique to conventional PCA, the latter would decompose a motion data matrix whose columns are observed motions \mathbf{d}_i into a reduced-dimensional basis matrix \mathbf{B}_{PCA} made up of the most significant eigenvectors times a matrix \mathbf{C} containing a vector of coefficients c_i for every observed motion. PCA represents each person as a set of M coefficient vectors, one for each action class. By contrast, our multilinear analysis enables us to represent each person with a single vector of coefficients relative to the bases comprising the tensor \mathcal{B} defined in (4).

6. Motion Data Acquisition

Human limb motion was recorded using a VICON system that employs four video cameras. The cameras detect infrared light reflected from 18 markers, 9 placed on each leg of a human subject. The system computes the 3D position of the markers relative to a fixed lab coordinate frame. The video cameras are positioned on one side of a 12 meter long walkway such that each marker can be observed by at least two cameras during motion. To extract the three angles spanned by a human joint, we must define a plane for each limb whose motion can be measured relative to the sagittal, frontal and transverse planes through the body.

A corpus of motion data was collected from 6 subjects. Three motions were collected for each person: walk, ascend-stairs, descend stairs. Each motion was repeated 10 times. A motion cycle was segmented from each motion sequence. To suppress noise, the collected motion data were low-pass filtered by a fourth-order Butterworth filter at a cut off frequency of 6 Hz and missing data were interpolated with a cubic spline. To compute the joint angles, we first calculate the frame coordinate transformation for each limb with respect to the lab, next we calculate the relative orientation of each limb in the kinematic chain, and finally we solve for inverse kinematic equations.

7. Results

First we model human motions as the composite consequence of the action performed and a motion signature, according to (3). Given a sufficient quantity of motion data,



Figure 5: A synthesized stair-ascending motion.

our human motion signature extraction algorithm can consistently produce walks and stair ascend/descend motions in the styles of individuals.

Next, we determine if people have motion signatures that are invariant of action classes. Therefore, we extract a motion signature from a subset of actions for a new individual (7) and synthesize the remainder of the actions using the extracted motion signature (6). The synthetic motions are then validated by classifying them against a database of all the real motions.

In a “leave-one-out” validation study, we verified that our algorithm was able to compute motion signatures sufficiently well to synthesize all three types of motions in the distinctive style of each individual compared against ground-truth motion capture data of that individual. If the motion signature \mathbf{p}_{new} captures the distinctive pattern of movement, the synthesized walk would best match the actual walk of the new person. Using a nearest neighbor classifier, the synthesized walk was indeed recognized against a complete database that includes the actual walk data for the new person.

Fig. 4(a) shows, in frontal view, the synthesis of three different styles of walking motion given only examples of descending stairs in those corresponding styles. Note that the walking styles differ subtly: The woman on the left walks in a pigeon-toed style, the clown struts, and the skeleton on the right walks with knocked knees. Fig. 4(b) shows a side view of the motions; the figures animated using synthesized motions are in the foreground. Fig. 5 shows a stair ascending motion synthesized for one of the individuals. Our algorithm extracted the motion signature from a sample walk from this individual. We then used the extracted motion signature to synthesize the stair-ascending motion for this individual. The motion signature was combined with general stair ascending parameters which were extracted a priori from our database.

In [11] we presented an animation short that was created using motion data synthesized by our algorithm. The



Figure 4: Synthesizing 3 styles of walking motions from example motions of ascending stairs in those corresponding styles. (a) Comparing synthesized walking motion data against ground truth (the synthesized data is depicted by the characters without hair), our method captures stylistic differences in motion such as pigeon-toed walking, knocked-knees or strutting. (a) The synthesized motions are depicted by the characters in the foreground and, for comparison, the captured walking motions are depicted by the characters in the background.

graphical characters shown are modeled and rendered by the *MetaCreations Poser* system.

8. Conclusion

We have introduced the notion of decomposing motion data into primitives such as action parameters, and most importantly a motion signature. To achieve such a decomposition, we have proposed an algorithm which is based on a numerical statistical analysis technique called n -mode analysis. It takes advantage of multilinear algebra in which motion data ensembles are represented as higher-dimensional tensors and an “ N -mode SVD” algorithm is applied to decompose the tensor.

Our tensor decomposition approach shows promise as a unifying mathematical framework for a variety of computer vision problems [13]. Our completely general multilinear approach accommodates any number of factors by taking advantage of the mathematical machinery of tensors. Our algorithm robustly extracts signature parameters from a corpus of motion data spanning multiple subjects performing different types of motions. We have shown that the extracted signatures are useful for the synthesis of novel motions for animating articulated characters for motion recognition.

In future work, will explore the simultaneous recognition of actions and people. We also plan to apply our approach to video input of human movement.

Acknowledgments

Professor Demetri Terzopoulos provided invaluable guidance during the writing of this paper. Motion data were collected at the Gait Laboratory of the Bloorview MacMillan Medical Centre in Toronto. The data acquisition work was done with the permission of Professor Stephen Naumann, Director of the Rehabilitation Engineering Department and with the helpful assistance of Mr. Alan Morris.

References

- [1] M. Brand and A. Hertzmann. Style machines. *Proc. ACM SIGGRAPH 2000*, New Orleans, LA, July 2000, 183–192.
- [2] R. Grzeszczuk, D. Terzopoulos, and G. Hinton. NeuroAnimator: Fast neural network emulation and control of physics-based models. *Proc. ACM SIGGRAPH 98*, July 1998, 9–20.
- [3] M. Gleicher (ed.) Making motion capture useful. *ACM SIGGRAPH 2001, Course 51*, Los Angeles, CA, August, 2001.
- [4] G. Johansson. Visual motion perception. *Scientific American*, June 1974, 76–88.
- [5] T. G. Kolda. Orthogonal tensor decompositions. *SIAM J. on Matrix Analysis and Applications*, 23(1):243–255, 2001.
- [6] L. de Lathauwer, B. de Moor, J. Vandewalle. A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.*, 21(4):1253–1278. On the best rank-1 and rank- (R_1, R_2, \dots, R_N) approximation of higher-order tensors. *SIAM J. Matrix Anal. Appl.*, 21(4):1324–1342.
- [7] A. Kapteyn, H. Neudecker, and T. Wansbeek. An approach to n -mode component analysis. *Psychometrika*, 51(2):269–275, June 1986.
- [8] J.R. Magnus and H. Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Wiley, New York, 1999.
- [9] J.B. Tenenbaum and W.T. Freeman. Separating style and content. *Advances in Neural Information Processing Systems 10*, MIT Press, 1997, 662–668.
- [10] L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31:279–311, 1966.
- [11] M.A.O. Vasilescu, Human motion signatures for character animation. *ACM SIGGRAPH 2001 Conf. Abstracts and Applications*, August, 2001, pg. 200.
- [12] M.A.O. Vasilescu. An algorithm for extracting human motion signatures. In *Proc. Tech. Sketches, IEEE Conf. Computer Vision and Pattern Recognition*, Kauai, HI, Dec., 2001.
- [13] M.A.O. Vasilescu. Multilinear image analysis for facial recognition. In *Proc. Int. Conf. on Pattern Recognition*, Quebec City, August 2002. These proceedings.
- [14] M.A.O. Vasilescu and D. Terzopoulos. Multilinear analysis of image ensembles: Tensorfaces. In *Proc. European Conf. on Computer Vision (ECCV 2002)*, Copenhagen, Denmark, May 2002. In press.