

GMS: A BASIC FORMAT TO ENCODE LOW-LEVEL GESTURE AND MOTION SIGNALS

Annie Luciani⁺

Matthieu Evrard⁺

Nicolas Castagné^{*}

Damien Couroussé^{*}

^{*} ACROE

⁺ ICA laboratory

INPG,

46 av. Félix Viallet, 38 000 Grenoble, FRANCE

ABSTRACT

Gesture encoding may become central in the coming years, given the growing importance of gesture analysis, storage and data exchanges between heterogeneous systems and applications. This importance is particularly clear in Computer Music where gesture control and interaction have a more and more crucial place in the context of new musical performance interfaces. For the past 20 years, various gesture-related formats have been proposed. Though, these formats, at different degrees, were designed in the context of specific applications (especially character animation and motion capture), or focus on the encoding of high-level gesture feature (segmentation, etc). Hence, the question of a generic format dedicated to the encoding of the gesture signal itself (i.e. gesture at a low level) remains to be solved. The article introduces the ‘Gesture and Motion Signal’ (GMS) format. GMS has been designed with the aim of low-level and generic, by analysing the minimal features a format carrying movement/gesture data needs. Hence, GMS features the inner properties of such signals: versatile dimensionality, versatile structural features, specific types of variables, and spatial and temporal properties.

1. INTRODUCTION

Gesture, motricity, and haptic perception are deciding factors in the way we act on and we perceive our environment. During the last few years, many research centers have focused on gesture control and movement analysis or synthesis, in domains as various as Computer Music, Computer Graphics, Virtual Reality, Surgery, Aeronautics, and generally speaking in every interactive systems. More and more, applications and systems are developed in opened and versatile contexts, especially in Computer Music [1], in a such a way that gestures captured with various means (haptic systems [2], motion capture systems [3], new musical interfaces, etc.) have to be correlated to other processes, especially sound synthesis. As a consequence, the question of formats dedicated to the encoding of gesture-related data becomes crucial – a Panel Discussion on the Topic has been submitted to ICMC 2007, for example [4].

This article questions the need of a *minimal* and *low level* description of gesture and motion data, able to code and structure all their features independently of the

context in which they have been produced and will be used, and introduce the GMS format.

In Computer Music, MIDI has played an important role as for real-time performance systems. However, MIDI is dedicated to the encoding of event-based control data. It cannot be considered as a mean to encode the performed gesture itself, which is today’s needed for many new musical interfaces.

Several formats exist for encoding motion control/gesture data, especially in motion capture. However, these formats were introduced in the context of quite specific application and uses, were not designed, at first, to be used in versatile and heterogeneous contexts, or focus on the encoding of high-level gesture features (segmentation, annotation, etc.). The data they encode are thus context dependent, and hardly correspond with generic low-level gesture data in its more general meaning.

More generally, when dealing with gesture data, most researchers and designers of new interactive systems use raw, header free formats, unpublished or device dependant proprietary formats.

Hence, the situation as for gesture and motion encoding is very different from the situation for visual and acoustical and signals or event-type data. While there exists standards for sound and image, such as AVI, AIFF, WAVE, etc., while MIDI standardized event data particularly in Computer Music, there isn’t any standard that deals with low-level gesture and motion. As well, while we share a common definition of what a sound signal is (e.g. a sampled signal at 8-96KHz, made of ”tracks”, often stereo, 4D, 4+1D, etc.), we still don’t know what raw “gesture signals” or ”gesture streams” are, and how they should be structured. MIDI, AIFF, MPEG norms for instance have played an important role in the development of computer music. They have allowed researchers, community of users, etc, to exchange data, systems that support music to easy communicate with one another. While gesture is more and more an object of study, while gesture storing, exchanges and uses develop rapidly, one can assume that the definition of a minimal, generic, low-level (close to the signal, not at a symbolic level) gesture format will be of crucial importance in the coming year.

The two next sections situate the various notions at hand and discuss the versatile properties of gesture data. The section 4 reviews the various existing file formats, and discuss their advantages and lacks. The two last sections present the requirements, specification and implantation of the proposed GMS format for a low-level, minimal and generic encoding of gesture signals.

2. CLARIFICATION OF TERMS

Given the ambiguity of various terms the article focuses on, a few definitions have to be discussed.

2.1. Action / Gesture / Movement / Signal

Following [5], action is the result of the task achieved or to be achieved - for example: “to drink a glass of water”.

“Gesture” is sometimes understood as “action”. In this article, though, we consider “gesture” as the “performed gesture” – the gestures actually performed.

Movement (or motion) is the moving in space of a part or of the totality of a system. A gesture may often be seen as being homogeneous to a “movement”, but, as we will discuss, there are cases in which gesture cannot be encoded with spatial variables.

Anyhow, each “action” can be achieved by means of many different gestures or movements. And a gesture or movement can be characterized by properties, such as soft, vivid, etc.

Action can be described at a high symbolic level, for example by language, or by means of event-based parameterized representations. Conversely, movements and performed gestures are explicit temporal phenomena. They need to be represented as temporal signals.

2.2. Motion / Gesture

An evidence concerning motions is that they refer to the evolution produced by a physical system, whatever it is: human body, real mechanical objects equipped with sensors, virtual objects, etc. One can speak as well of the motion of a human body, of a leaf, of a musical instrument, a sounding source, etc. Motion is therefore considered as the result of the performance, i.e. as an output of an evolving system.

Conversely, when focusing the use of such movements as “a command” or an input to control the behavior of another system (for example: a sound synthesis system), the only usable terms are “action” or “gesture”. Unfortunately, as said in the previous paragraph, “action” refers more to the high level of the task than to way it was achieved. It is very general, and can be used for any type of systems. Conversely, “gesture” is usually reserved to humans, and some times understood as “action”. Facing this dilemma, we decide to emphasize the differentiation between input and output functionalities, more than the nature of the evolving system that produces the signals. Thus, we call “gestures” all the physical data that that can be applied as an input signal, i.e. as a “cause of a performance”, whatever the producing system is (human beings or other objects).

Anyhow, Motion/output and Gesture/input are two representatives of similar temporal signals, corresponding with an evolving physical system – see Figure 1. Thus, to simplify, the term ‘gesture’ and ‘motion data’ could be used in the following to refer to this sort of ‘gesture like’ data, no matter it is considered

as an input or an output of an evolving systems, possibly a human body.

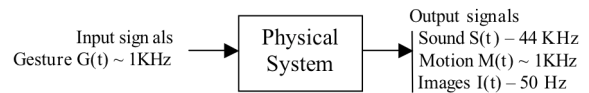


Figure 1: Gestures and Motion as input and outputs of physical systems. Gesture (input) and Motion (output) are comparable types of data.

Rather than discussing the differences between Gesture and Motion, one can emphasize their similarities. This leads to identify them as a new and unique type of data, beside visual and acoustic data, and to introduce a dedicated encoding format.

2.3. Gesture and Control

Gestures can be considered as a common feature of numerous situations, applications, systems etc., and they may be a privileged mean for allowing those systems to cooperate: sound synthesis systems, musical control and performance, animation control, multisensory control and performance such as in VR or in multimedia arts.

The figure 2 illustrates some basic cases of gesture signals and gesture control. On the figure, grey circles are placed on the streams one can consider as gesture streams (or as being equivalent to a gesture signal). Conversely, the grey squares are placed on the streams that are no more gesture streams.

In the case of gesture control of real or digital musical instruments, the gesture signals can be either the signals sensed directly on the human body (fingers, hands, arms, whole body) (figure 2, grey circles 1 and 2), or the motions produced by an object manipulated by the player (hammer, key, bow, stick, etc...) (figure 2, grey circle 3).

In the case of gesture control of 2D or 3D visual motions, the gesture signals can be either the signal sensed directly on the human body (figure 2, grey circles 4 and 5), or the motions produced by an object manipulated by the performer (sticks, force feedback devices, etc.) (figure 2, grey circles 6 and 7). Conversely to what happens in musical instruments, in which the output “sound” is usually not considered as a “gesture”, the 3D motions provided by a 3D virtual objects (Figure 2, grey circle 7) are of the same nature as the 3D motions produced by a human body in body motion capture. They may be used as an input for controlling another system (3D object, musical instrument, etc.). Consequently, a 3D virtual object produces signals that can be considered as being homogeneous with “gestures” (Figure 2, grey circles 4, 5, 6).

In the case of the gestures that produce formal and symbolic data, as when acting on a computer keyboard, the gesture signals (Figure 2, grey circles 8 and 9) are transformed in outputs (words, sentences) that can usually no more be considered as gesture signals.

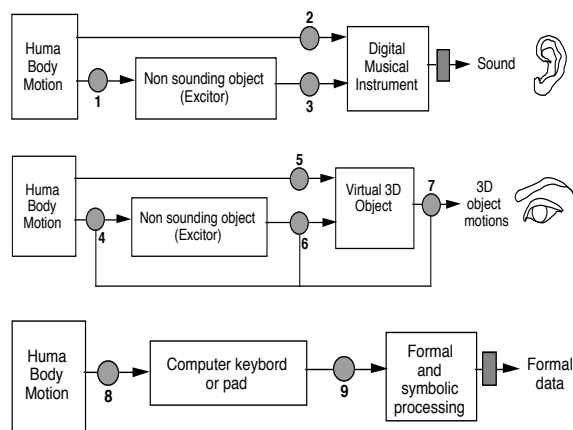


Figure 2: Various cases of gesture control. Up: gesture control of digital musical instrument. Middle: gesture control of a 3D virtual object. Down: gesture control of formal and symbolic data.

3. CHARACTERIZATION OF GESTURE SIGNALS

Gesture signals, whatever the way they are produced (objects or human motion, virtual object...), and whatever the way they are considered (as outputs or inputs of evolving systems), do present specific properties that allow distinguishing them among other temporal signals (especially aero-acoustical signals or visual signals). This section reviews these properties that will be used as a basis for defining a basic generic format.

3.1. Morphological versatility

A first evidence is the morphological versatility of gestures. While sounds and images can be displayed in predefined environments (displays or 3D Caves, with fixed sizes and resolution, stereo or quadrasonic rendering for sounds, etc.), the structure and the morphology of gesture signals are more versatile, depending on the tasks or the manipulated tools. To take into account this inner versatility, we propose to structure gesture signals along two complementary features: geometrical and structural dimensionalities.

3.1.1. Geometrical dimensionality

Geometrical dimensionality refers to the dimensionality of the space in which the gesture is evolving. It can vary from a pure scalar or a set of scalars to geometrical 1D, 2D, 3D.

For example, piano or clarinet keys are pushed or closed according to a 1D finger motion (figure 1a). More generally, controlling sounding objects or tuning parameters (for example the value of an elasticity or the amplitude of a deformation), can be made through devices that evolve in a 1D non oriented space, like set of sliders or set of knobs (figures 1c, 1d), and that can be described by a scalar or a set of scalars.

Conversely, in cartoon animation or in scrap-paper animation or animated painting, the space is reduced to a plane. Gestures and motions evolve in a 2D space (Figure 3, g), and can be described on two orthogonal oriented axes.

When we manipulate an object (real or virtual), the dimensionality of the space is obviously 3D. Thus, describing motion or gesture requires using three orthogonal oriented axes (figure 3, e, f, h).

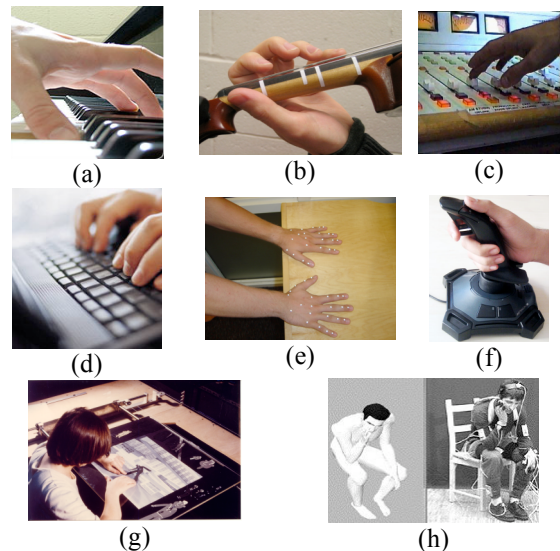


Figure 3: Versatility of the gesture morphology

3.1.2. Structural dimensionality

For a given geometrical dimensionality, gestures have another type of variability, related to their functional organization. We call this second variability structural dimensionality. In the case of articulated 3D solids, it usually corresponds to the number of degrees of freedom. A couple of examples allow illustrating the variability of structural dimensionality of gesture.

When we are acting on a keyboard of n keys (a piano keyboard, a computer keyboard, a set of buttons...), the performed gesture, and similarly the signals produced by the n keys, the geometrical dimensionality is 1, but the system can be considered in two ways either as n independent systems (for example n systems of one key), or as a single system of n degrees of freedom.

When dealing with human body motion (which is more and more the case in Computer Music) the geometrical dimensionality is 3 but the structural dimensionality (ie: the number of considered degrees of freedom) is variable: more than 200 in the real body, only N when using a motion capture systems with N sensors ($N=1$ to 16, or more).

In physical models of a bowed string, the two dimensions of the deformations (pressing the string, bowing the string) are usually decoupled, and the system can be considered as two superposed 1D gestures, thus as one system of two 1D DoF.

In particle physical model of a flow, the number N of 3D particles, i.e. the number N of 3D DoF, could be

more than some hundreds. Considering that this flow can be used to control large spatial scene (for example wind playing a musical instrument, flow controlling atmospheric or oceanic deformations), the motions can be stored as gestures of N 3D DoF.

3.1.3. To conclude

This analysis and these examples show that the geometrical dimensionality of gestures should not be fixed, since it can vary from scalar to 3D (we don't examine here space with larger geometrical dimensionality, although the case can occur). They show also that the structural dimensionality can vary widely (from 1Dof to several hundred). In addition, the structural properties of the considered gesture signal cannot be totally independent of how the system is analyzed by the designer, that must be able to declare them freely.

3.2. Quantitative Ranges

In addition to the geometrical and structural dimensionalities, gesture signals present specific quantitative spatial and temporal features.

3.2.1. Temporal ranges

As shown on Figure 4, the frequency bandwidth range is one of the quantitative features that allow distinguishing gesture signals among others (especially aero-mechanical signals, visual motions).

Although the three zones of the figure 4 are overlapped, they enable a useful categorization: visualizing motions requires a sampling rate up to 100 Hz; manipulating an object with force feedback requires a sampling rate from some Hz to some kHz; recording sounds requires a sampling rate from around 10 kHz to some 40 kHz. Indeed, gesture signals are at the middle range.

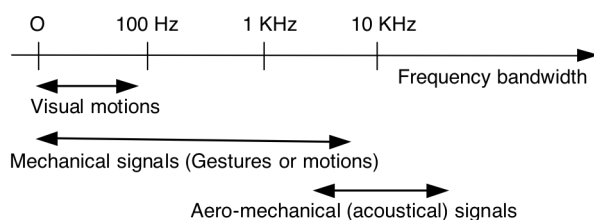


Figure 4: Temporal range of sensory signals

3.2.2. Spatial ranges

Spatial positions and ranges is another feature of gesture signals. Whereas audio signals are zero-centered deformations of about some millimeters, gestures and motions correspond with deformations and displacements that are non-necessarily centered on 0, and that are at a larger spatial scale, from centimeters to meters.

3.2.3. To conclude

One can note that the quantitative properties of gesture signals position them at middle range: (1) it is similar to visual motion as for spatiality, but it needs higher frequency rate; (2) it needs a lower frequency rate than the sound, but it runs at higher non-centered spatial range.

Indeed, one can consider the gesture signal as a pivot for multisensory platforms, including interactive sound systems. Visual motions can be produced by simply under-sampling them. A non-linear control-simulation process can then produce sounds, transforming low frequency gesture command signals in high frequency sound signals.

3.3. Type of variables

Since motions and gestures are produced by physical systems (either real or virtual), and used to control physical systems, the data can be of two different types: (1) extensive variables, such as variables that derives from spatial information (positions, velocities, angles...); and (2) intensive variables as forces, torques, etc. For instance, there are human gestures (bowing a string, acting on a pressure device, etc.) that cannot be correctly analyzed nor encoded by considering (only) motions and displacements, and that require intensive variables. Noticeably, one can note that visual and acoustical data are only extensive variables (positions and/or displacements).

In natural situations, when gestures are used for manipulating object, physical energy is exchanged between the two interacting bodies (for example object and human). To represent such interactive dynamic systems, the correlation between extensive and intensive variables must be considered, either explicitly as in Newtonian formalism, or implicitly as in energy formalisms. Both extensive and intensive variable are needed for a complete encoding of the interaction.

4. A REVIEW OF GESTURE FORMATS

Several file formats dedicated to gesture have been developed so far [6, 7] in various communities. Here, we position some of the most known (see Table 1) along the properties of gesture signals elicited in the previous paragraphs: versatility of the geometrical and structural dimensionality, spatial and temporal requirements, and multiplicity of the types of data. Noticeably, we do not discuss MIDI in this study, since MIDI is an event-based norm that does not apply to the encoding of "gesture signals" as defined in section 2.

BVA, BVH ASK, SDL	[8] Developed by Biovision, a firm providing optical motion capture systems.
HTR, TRC	[8] Developed by the Motion Analysis.
AOA	[8] Designed by Adaptive Optics, a firm that develops motion capture hardware.
ASF	[9] Developed by Acclaim, a video game

AMC	firm. Now used by Oxford Metrics (Vicon motion capture Systems). Public domain.
BRD	[10] Introduced by Lambsoft, for the "Flock of Birds" motion capture system.
CSM	[10] 3D Studio Max and Character Studio format. Recoding of characters' movements.
C3D	[11] National Institute of Health. Initially encoding movements of athletes. Largely used today (ex: as a file format for encoding basic marker data, for the Vicon Motion Systems).
GDIF	[12, 13] Currently developed within the ConGAS European project. Aims at handling different levels of movement data: from raw data to higher level descriptors.

Table 1: some of the format that relates to gesture encoding

4.1. Origin of the existing formats

The considered formats have been historically introduced in the context of specific application and uses – namely in the context of Motion Capture, Biomechanics and Computer Music. Though some have eventually been adapted to other applications and usages than their preliminary target, they had neither been designed at first to be used in other applications, nor as a data exchange mean. Consequently, the data they encode, and the manner these are encoded, depend on the context of the usages preliminary foreseen, and hardly correspond with gesture / motion data in its more general meaning.

4.2. Object dependence

Except C3D and GDIF, all the formats are based on a hierarchy depicting a hierarchical dependency over the various recorded signals. The hierarchy is inherited from the hierarchical structure of the object that produces the signals. For example, describing the movement of a humanoid is made according to a pre-defined tree-like structure of a cinematic chain of the body [14]. In that case, most often, the value of a position/rotation is given in the coordinate system of the upper hierarchical level it belongs to.

As for it, the CSM file format developed along with the Character Studio software is not explicitly hierarchy-based. It encodes the motion of a set of 3D points in absolute coordinate. However, using CSM files within Character Studio requires the names of all the points to match Character Studio setup. A skeleton is implied, and the number of points and their 'meaning' are fixed.

Such a hierarchical structure has proved to be operational in domains such as motion planning or character animation. It is most often adapted to the case of articulated objects such as human body, rigid-articulated robot, etc. However, the use of a hierarchical structure does not apply to motion/ gesture encoding in its more general meaning – especially to the case of deformable non tree-articulated objects and mechanisms.

4.3. Context Dependence

Some of the discussed formats contain raw data directly outputted from gesture devices (that are device-dependant and cannot be considered yet as "gesture data"), or additional data dealing with the context of data acquisition. The context data are typically data concerning the parameters of the capture system (number of sensors, captors used to determine the movement at a certain moment, sensitivity of sensors, etc.).

The C3D file format is the most representative example of such additional data recording. In a C3D file, each motion stream is associated with the description of the measurement process that was used. Indeed, C3D corresponds firstly with measurement or experimental setup recordings.

4.4. Predefined dimensionality

Most of the formats predefine a specific dimension for the recorded data: 3D points in AOA or CSM, "7D points" in BRD. In the BVH, ASF/AMC and HTR file formats, the dimension of each joint is declared in accordance with the constraints underlain by the hierarchical structure. Omitting a coordinate corresponds with constraints of the mechanical system. Once again, one can observe that these formats suppose that the recorded motion data are intimately associated with a tree-like structure of the object supporting the sensors. As for them, C3D files store conventional 3D positional information. 1D or 2D motion data may be recorded, by fixing one or more coordinate to zero all along the stream. Though, scalar data can also be recorded within the raw data (see §4.5).

Anyhow, one can consider that pre-defining the dimensionality of the data to be stored is a real problem since the versatility of dimensionality is a core property of gestures, and since one may often have to consider coexistence of different dimensionalities in the same gesture stream.

4.5. Type of variable

C3D files can encode two categories of data: 'analog data', i.e. scalar data acquired from any kind of sensor, and '3D data' i.e. motion-like data displayed by a motion capture system. The format implies no inner limit in the type of the 'analog data', so that any type of variable may be recorded, including hearth rate, electromyogram, etc. Additionally, a non-optimized trick allows recording 'kinematical data' (angles, moment, acceleration...) processed from the motion data. Thus, besides the fact that C3D is widely used for motion recording, it does overpass the strict scope of motion and gesture. It is a very flexible and adaptable format – but it is too general for being a Basic Low-Level format adapted to the category of data we call "gestures".

As for the type of the 'gesture/motion data' themselves, the various file formats are restricted to position-like data (C3D allow recording other types, but only within the very open 'analog data' sections). They do not natively include the possibility to encode the type

of the variables. They are especially not adapted natively to the recording of intensive variable.

The currently developed GDIF is quite flexible as for the type of the variable of the raw data. Conversely, however, this flexibility implies that the encoded data could be of any type – not only gesture. Indeed, GDIF does not aim at focusing on the question of a generic encoding of *low-level* gesture data. It rather aims at encoding in a multi-layered organization various categories of data that relates to the performance, including high-level descriptors and annotation such as Laban qualitative analysis cues, or data directly acquired on sensors (ie: not transformed into signals representing the gesture, but sensor-dependent). Conversely, we consider, a need today is to focus on the encoding of the low-level gesture data themselves – somehow as well as, for example, MPEG7 was introduced a long time after the acceptance of basic sound file formats.

4.6. Data coding

Except for the C3D format, ASCII encoding is used in all the discussed formats. This corresponds to the fact that the frequency bandwidth of motion capture systems and gesture controllers used to be on the VLF range (from 10 to 300 Hz) and that recorded gestures used to have a small number of DoF and tracks (less than one hundred). Though, today, one usually considers that the frequency bandwidth of gesture signals may be of some (up to ten, or perhaps more) kHz, especially when dealing with force feedback gesture interaction. In addition, the diversification of the objects (real or virtual) producing motions, and the decrease of the costs of sensors, imply that the number of tracks to be stored is significantly increasing. Consequently, gesture files may correspond with amount of data far more important than it had been considered so far, and ASCII encoding is no more usable. Indeed, the use of a binary format (like for sound, images or movies) is a necessity for a generic gesture/motion format.

5. THE GMS FORMAT

On the basis of the previous analysis of gesture (or more generally of motions usable as a command of an evolving system), this section introduces the proposed Low Level Gesture File Format. This format is called GMS, which stands for “Gesture and Motion Signals” - file extension is .gms.

5.1. Requirements for a gesture file

5.1.1. Sampling Rate

First of all, a gesture signal is a signal sampled at the necessary Shannon frequency. In the current 0.1 version, a unique sampling rate is used for all the data in a GMS file. This choice does not correspond to an optimal encoding, and will be reconsidered in future versions.

5.1.2. Geometry and Structure Encoding

We propose to organize the gesture data in four levels at maximum: Tracks – Channels – Units – Scenes:

- **Gesture Track.** The track i contains a sampled monodimensional scalar $a_i(t)$ corresponding to a single A-D track $a(t)$, sampled at its Shannon rate.

- **Gesture Channels.** Channels support the geometrical dimensionality of gesture signal. A Channel is composed of one or several tracks. A Channel can be 1D0 (a pure scalar), 1Dx (a vector on x axis), 1Dy, 1Dz, 2Dxy, 2Dyz, 2Dzy, or 3Dxyz. A Channel can encode either “extensive variable (EV)” i.e. homogeneous to spatial variables (positions P, Angles A, velocities V, Accelerations G) or “intensives variables (IV)”, i.e. homogeneous to forces F.

- **Gesture Units.** Units are provided for supporting the structural dimensionality of gestures. A unit is then composed of several channels in which the signals are not dynamically independent, whatever the way in which they are correlated is. This differs from the structural properties of motion capture formats that focus on tree-like structures. The user can declare units freely. Along with supporting a low-level structuring of the signal, Units allow preventing non-relevant signal processing in future uses of the data, for example breaking of the correlation.

- **Gesture Scenes.** Finally, a scene is composed of several units that are not – or can be considered as not - dynamically linked.

5.2. An example

In table 2, we illustrate how a user could structure a complex gesture scene like the one described in the ‘DESCRIPTION’ column.

UNIT	DESCRIPTION	CHANNELS
Pianist	One musician plays with a small keyboard composed of 8 piano-like keys.	8 Mono dimensional Channels “PianoKeys” PK1 [EV(P), 1Dz], ... , PK 8 [EV(P), 1Dz].
Stick Source	Another one is controlling the spatial position of a sound source via a 3D stick.	1 3Dxyz Channe 1 “SoundSource”: SS [EV(P), 3Dxyz]
Light	Another one is controlling the orientation of a light source via a 2D force sensor (a force pad).	1 2Dxy channel “LightSource” : LS [IV(F), 2Dxy]
Dancer	A dancer is equipped with a motion capture system of 16 3D points.	16 3Dxyz channels “DancerPoints”: DP1[EV(P), 3Dxyz], ..., DP16[EV(P), 3Dxyz]
Juggler	A juggler is manipulating a 6D force feedback ball to play with a 6D virtual racket.	1 3Dxyz channel “Ball1”: BL1 [EV(P), 3Dxyz] 1 3Drqf channel “Ball2”: BL2 [EV(A), 3Drqf]

Fluid	A virtual fluid of N 1D masses acts on a virtual string.	n 1D “fluidMass”: FL [EV(V), 1D0]
-------	--	--------------------------------------

Table 2: a possible structure for a gesture scene.

6. TMGMS Format Implementation

This paragraph describes the current implementation of the concepts developed previously. The GMS format, version 0.1, have been implemented in a free C/C++ read/write library, which is currently submitted to file formats authorities, and as a MIME type.

The format is built above the IFF (Interchange File Format) coding standard for binary files [15]. IFF, and its equivalent RIFF, are largely used by several formats as AIFF for sound files, AVI for movie files, ILBM or GIFF for picture files.

As required by the IFF coding standard, data in the proposed gesture file are encapsulated in structures called chunks.

Table 3 presents a description of the data contained in each of these chunks. The binary types of those data are CHAR (a byte), USHORT (a 2 bytes unsigned integer), LONG or ULONG (a 4 bytes signed or unsigned integer), FLOAT32 or FLOAT64 (a 4 or 8 bytes floating point number - IEEE 754 standard). All the data covering more than one byte in the file are encoded in BIG ENDIAN.

Header	CHAR[4] chunkId = ‘FORM’		ULONG fileSize;	
	CHAR[4] fileType = ‘GSM ’			
Version	CHAR[4] chunkId = ‘VERS’		ULONG chunkSize	
	USHORT versNum		USHORT subVersNum	
Scene	CHAR[4] chunkId = ‘SCEN’		ULONG chunkSize	
	USHORT sceneNameLength		CHAR*sceneName	
	ULONG nbFrame	FLOAT64 freq	USHORT dataType	
	FLOAT64 scale	ULONG blockSize		
Unit	CHAR[4] chunkId = ‘UNIT’		ULONG chunkSize	
	USHORT unitNameLength		CHAR* unitName	
Channel	CHAR[4] chunkId = ‘CHAN’		ULONG chunkSize	
	USHORT chanNameLength		CHAR* chanName	
	USHORT dimension		USHORT type	
Frame	CHAR[4] chunkId = ‘FRAM’		ULONG chunkSize	
	TypeOfData[nbTrack][nbFrame] frames			

Table 3: description of the data contained in each chunk

The version chunk corresponds to the chunk structure to be found in the file. Future evolutions of the gesture file are foreseen.

The scene chunk contains information about the gesture scene encoded in the file. dataType gives the type of data of the signal. The current version supports three type of data: FLOAT32, FLOAT64, and LONG. scale is a scale factor to apply on the signal. blockSize allows setting up an alignment of the frames (ie: the gesture data) on multiple of blockSize. This may be useful on some architecture to process the file in real time. If blockSize is 0 (default), the size of a block equals to the size of a frame.

The unit and channel chunks encode information on a gesture/motion unit and channel. Channel and unit chunks are interleaved: the current channel belongs to the last declared unit. The dimension integer encodes the dimension of the channel, and accordingly its number of tracks. Dimensions supported in version 0.1 are: 0D (pure scalar); a vector on 1Dx (resp. 1Dy or 1Dz), on 2Dxy (resp. 2Dyz or 2Dzx), on 3Dxyz. type encodes the type of the data. Current supported types are position, and force.

The frame chunk contains the signal itself, given frame by frame.

7. CONCLUSION

While gesture devices, especially haptic devices, develop, while applications communicate more and more though gesture-like data, the definition of a low level gesture format allowing to structure and encode performed gesture data can be seen as a major need for the near future. This is especially true in Computer Music, in which the study of instrumental gesture, and the design of new musical interfaces and instruments are important topics.

We extracted the specific properties of gesture/motion signals over other temporal signals (sounds and images): morphological versatility (decomposed in geometrical and structural properties), spatial and temporal ranges, and variety of type of data. We showed that these properties are able to explain why gestures can be considered as at “a hub place” in multisensory situations and how they can be shared by various applications. This led to introduce a basic file format, able to encode all the minimal and necessary features.

Various evolutions are foreseen, including:

- extending the format to multi-frequency gesture stream; eventually, defining a set of standard frequencies.
- extending the possible dimensions of channels, especially to 6D; although the necessity of such ND channels is not true evidence, and is still under study.
- extending the supported types of data, including for example displacement, velocities, angles, torques, etc.

These extensions, though, require further analysis for being exhaustive without damaging the ‘minimality’

and ‘simplicity’ of the proposed format. Suggestions and comments will be welcome.

Since its very first release in 2006, the proposed format has been used in various new interactive sound systems, especially including synchronous haptic interaction with physically based multisensory simulators [2]. It is also available in the GENESIS environment for musical creation by means of physical modeling [16], as a gesture input format, and in the MIMESIS environment for movement synthesis and animation with physical modeling [17], as an input and output gesture/movement format. Experiences for connecting various systems (computer animation, computer sound, multimodal systems) through the proposed gesture format are in progress in the laboratory.

The specifications of the GMS format, and a multiplatform, portable, lightweight library for handling GMS files and streams, are available freely on the web pages <http://acroe.imag.fr/gms> since January 2007.

8. ACKNOWLEDGMENTS

This work was supported by the FP6 Network of Excellence IST-2002-002114 - Enactive Interfaces, and by the French Ministry of Culture.

9. REFERENCES

Web addresses are provided for the file formats that have not been precisely described in scientific publications.

- [1] Cadoz C., Wanderley M., 2000. « *Gesture – Music* », pages 71-94, Ircam – Centre Pompidou.
- [2] Luciani, A., Cadoz, C., Florens, J.-L., 1998. “Towards a complete representation by means of computer – the instrumental communication interface concept”. In *Cyberworld*, Springer-Verlag, Tokyo.
- [3] Menache, A., 1999. *Understanding Motion Capture for computer animation and video games*. Morgan Kaufmann Ed.
- [4] Alexander Refsum Jensenius, Nicolas Castagn, Esteban Maestre, Joe Malloch, Dougie McGilvray, Diemo Schwarz, InfoMUS DIST : “The need of formats for streaming and storing music-related movement and gesture data”. Panel discussion and paper, Submitted to *ICMC07*, 2007, to appear.
- [5] Smyth, M., Wing, A. 1984. *The psychology of Human movement*. Academic Press.
- [6] Meredith, Maddock, 2000. *Motion capture file formats explained*. Department of Computer Science Technical Report CS-01-11, University of Sheffield.
- [7] Van, Van Baerle, S., 2000. “Motion Editing : Principles and Practice”. In 2000 *Game Development Conference Proceedings*. March 20-24. San Jose, CA.
- [8] AOA, BVH, HTR: <http://www.cs.wisc.edu/graphics/Courses/cs-838-1999/Jeff/{AOA,BVH,HTR}.html>
- [9] ASF: <http://www.darwin3d.com/gamedev/acclaim.zip>
- [10] BRD, CSM: <http://www.dcs.shef.ac.uk/~mikem/fileformats/{brd,csm}.html>
- [11] C3D, <http://www.c3d.org>
- [12] A. R. Jensenius, T. Kvifte, and R. I. Godøy. “Towards a gesture description interchange format”. In *Proceedings of New Interfaces for Musical Expression*, NIME 06, IRCAM - Centre Pompidou, Paris, France, June 4-8, pages 176–179, 2006.
- [13] E. Maestre, J. Janer, A. R. Jensenius, and J. Malloch. “Extending gdif for instrumental gestures: the case of violin performance”. Submitted to *ICMC07*, 2007, to appear.
- [14] HERD Herda, L., Fua, P., Plänkner, R., Boulic, R., Thalmann, D., 2000. Skeleton-Based Motion Capture for Robust Reconstruction of Human Motion. In *Computer Animation 2000 (CA00)*, p.77.
- [15] Morrison, 1985. EA IFF 85 - *Standard for interchange format files* - Electronic Arts.
- [16] N. Castagne and C. Cadoz. GENESIS: “A Friendly Musician-Oriented Environment for Mass-Interaction Physical Modeling”. In *Proceedings of the 2002 International Computer Music Conference*, Göteborg, Sweden. San Francisco, CA, 2002. International Computer Music Association.
- [17] Matthieu Evrard, Annie Luciani, Nicolas Castagné, “MIMESIS: Interactive Interface for Mass-Interaction Modeling” - In *Proceedings of CASA 2006*, Geneva, July 2006, Nadia Magnenat-Thalmann& al. editors. pp177-186.