

# A Common Gesture and Speech Production Framework for Virtual and Physical Agents

Quoc Anh Le  
Telecom ParisTech  
37 rue Dareau  
75014, Paris  
quoc@enst.fr

Jing Huang  
Telecom ParisTech  
37 rue Dareau  
75014, Paris  
jing.huang@enst.fr

Catherine Pelachaud  
CNRS, LTCI  
37 rue Dareau  
75014, Paris  
catherine.pelachaud@enst.fr

## ABSTRACT

We introduce a modular system to generate communicative expressive gestures accompanying speech for an agent. This system is designed as a common model for different embodiments so that its processes are independent from a specific agent. There are two main features of this system. Firstly gesture expressivity is taken into account when gesture animation are computed on the fly from abstract gesture templates. Secondly gestures are scheduled to ensure that their execution are tightly tied to speech. In this paper, we present the first implementation of this system being used to control co-verbal gestures of the Greta virtual agent and of the Nao physical robot.

## Categories and Subject Descriptors

H.5.2 [Information Interfaces and Presentation]: Miscellaneous

## General Terms

Algorithms, Design, Language

## Keywords

Gesture, Speech, Synchronization, Expressivity, HRI, HMI, BML, FML, SAIBA, GRETA, NAO

## 1. INTRODUCTION

For many years, we have developed a virtual intelligent agent (IVA) system, namely GRETA [25] that enables to produce and to respond appropriately verbal and non verbal behaviors like gaze, facial expressions, head movements and gestures to human users. The modular architecture of this system follows SAIBA (Situation, Agent, Intention, Behavior, Animation), an international standard multimodal behavior generation framework for embodied agents [29].

Recently, the advance of robotics technology bring us humanoid robots with certain behavior capacities as much as

the virtual agents have [15]. For instance the expressive anthropomorphic robot Kismet at MIT can communicate rich information through its facial expressions [2]. The ASIMO robot produces gestures accompanying speech in human communication [27]. The Nao humanoid robot can convey several emotions such as anger, happiness, sadness through its dynamic body movements [9, 20]. The approach of two domains, virtual embodied agents (e.g., embodied conversational agents) and physical embodied agents (e.g. robots) allows us to think about a common framework to control their behaviors in a same way. For this reason we aim at extending and developing our existing system to be able to handle both virtual and physical agents. The common gesture generation model for the virtual agent Greta [25] and the robot Nao [8] is our first attempt to reach this goal. In this model we focus on three main aspects of human gestures. They are the form of gestures, the expressivity of gestures and the synchronization of gestures with speech.

Since the virtual and physical agents have different motion capacities (e.g., the robot has less degrees of freedom and has some limits in its movement speed), our methodology is to control the agents' behaviors at a symbolic level through representation languages such as FML [12] and BML [29]. This solution enables using the same processes for selecting and planning gestures, and different algorithms for creating animation only.

Regarding the form of gestures, the robot and the virtual agent may not be able to display the same gestures but their selected gestures have to convey the same meaning (or at least similar meanings). For this reason, we create two repertoires of gesture templates, one for the virtual agent and another one for the robot. These two repertoires have entries for the same list of communicative intentions. Given an intent, the system selects appropriate gestures from either repertoires. For instance to point at an object, Greta can select an index gesture with one finger. Nao has only two hand configurations, open and closed. It cannot extend one finger as the virtual agent does, but it can full stretch its arm to point the object. As a result, for the same intent of object pointing, while the Nao repertoire contains a gesture of whole stretched arm, the Greta repertoire contains an index gesture with one finger.

Concerning gesture expressivity, we have designed a set of quality dimensions such as: 1) Spatial extent (SPC) determines the amplitude of movements (e.g., contracting vs. expanding); 2) Fluidity (FLD) refers to smoothness and continuity of movements (e.g., smooth vs. jerky); 3) Power (PWR) defines acceleration and dynamic properties of move-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICMI 2012 Workshop on Speech and Gesture Production in Virtually and Physically Embodied Conversational Agents, October 26, 2012, Santa Monica, CA, USA.

Copyright 2012 ACM 978-1-4503-1514-2/12/10...\$15.00..

ments (e.g., weak vs. strong); 4) Temporal extent (TMP) refers to the global duration of movements (e.g., quick vs. sustained actions); 5) Repetition (REP) defines tendency to rhythmic repeats of specific movements; 6) Tension (TEN) refers to hand-arm muscle states (e.g., relax vs. tense); 7) Openness (OPE) determines spatial relation of hand-arm positions to the body (e.g., away from body in an open gesture). These parameters have been implemented for the virtual agent Greta [11]. We want to realize such a set of expressivity parameters for the Nao robot’s gestures. From a same gesture template, an agent can animate the gesture in different ways depending on current emotion state or personality of the agent. For instance a sad agent may realize gestures slowly and weakly vs. an angry agent can gesture quickly and strongly.

In this framework, the synchronization of gestures with speech is ensured by adapting gesture movements to the speech timing. According to Kendon and McNeill [16, 21], the most meaningful part of a gesture (i.e., the stroke phase) mainly happens at the same time or lightly before the stressed syllables of speech. While a robot may potentially need longer time for execution of hand movements than a virtual agent, our synchronization engine has to be able to predict gesture duration for each agent’s embodiment type so that their gestures are scheduled correctly. In our case, the duration of gesture movements between any two positions in gesture space of the Nao robot is pre-calculated because we cannot have it on the fly.

The paper is structured as follows. The next section presents some recent initiatives in generating gestures for virtual agents and for humanoid robots and how our approach differs from these existing works. Then, Section 3 gives an overview of our system and explains how our system is designed to be common for both virtual and physical agents. Section 4 presents gesture lexicons which are elaborated to be adapted to agents’ embodiment. In Section 5 and 6, we describe the mechanism to select and plan gestures from gesture lexicons to synchronize with speech and to be rendered expressive. Section 7 shows how gestures with expressivity are produced and realized for Greta and Nao. Section 8 concludes the paper and proposes some future works.

## 2. STATE OF THE ART

This section presents some recent initiatives to generate co-verbal gestures for virtual agents and physical robots. The differences and similarities between these approaches and our system are analyzed in detail.

### *Co-verbal Gesture Production for Virtual Agents*

The first system that generates gestures for a virtual agent was proposed by Cassell et al [3]. In their system, gestures are selected and computed from gesture templates. These gesture templates are predefined and stored in a gesture repertoire called lexicon. A similar method is still used in our system. However our model takes into account a set of expressivity parameters while creating gesture animations. So that we can produce variants of a gesture from a same abstract gesture template.

Stone et al. [28] proposed a data-driven method for synchronizing small units of pre-recorded gesture animation and speech. Their approach generates gestures synchronized with each phrase of speech automatically. Different combina-

tion schemes simulate agent’s communicative style. Another data-driven method was proposed by Neff et al. [22]. In this method their model creates gesture animation based on gesturing styles extracted from gesture annotations of real human subjects. In general, both these two systems and our model create gestures from predefined gestural prototypes. In our system, gestural prototypes are abstract gesture templates that have no reference to a specific animation parameters of agents (e.g., wrist joint).

The model of Bergmann et al. [1] combines data-driven machine learning techniques and rule-based decision methods. It also introduces several contextual factors. The whole architecture is used for a computational Human-Computer Interaction simulation, focusing on the production of the speech-accompanying iconic gestures. This model allows the generation of gestures on the fly. It is one of the few models to have such a capacity. However this is a domain dependent gesture generation model. While our model can handle all types of gestures regardless specific domains, the model of Bergmann is limited to iconic gestures and it have to be re-trained with a new data corpus to be able to produce appropriate gestures for a new domain.

Concerning the expressivity of nonverbal behaviors (e.g., gesture expressivity), it exists several expressivity models either act as filter over an animation or modulate the gesture specification ahead of time. EMOTE implements the effort and shape components of the Laban Movement Analysis [4]. These parameters affect the wrist location of the humanoid. They act as a filter on the overall animation of the virtual humanoid. On the other hand, a model of nonverbal behavior expressivity has been defined that acts on the synthesis computation of a behavior [10]. It is based on perceptual studies conducted by Wallbott [30]. Among a large set of variables that are considered in the perceptual studies, six parameters [11] are retained and implemented in the Greta ECA system.

### *Speech Gesture Production for Humanoid Robots*

The most similar approach to our model is the work of Salem et al. [27]. We share the same idea of using an existing virtual agent system to control a physical humanoid robot. Both of us have to face difficulties of physical constraints while creating robot gestures (e.g., limit of space and speed robot movements). However, we have certain differences in resolving these problems. While Salem et al. fully use the MAX system to produce gesture parameters (i.e., joint angles or effector targets) which are still designed for the virtual agent, our existing GRETA system is extended and developed so that its extern parameters can be customized to produce gesture parameters for a specific agent embodiment (e.g., a virtual agent or a physical robot). For instance, the MAX system produces an iconic gesture of complicated hand shapes that is feasible for the MAX agent but have to be mapped to one of three basic hand shapes of ASIMO. In our system, we deal with this problem ahead of time when elaborating lexicon for each agent type. This allows us to ensure that both agents convey the same information. In addition, the quality of our robot’s gestures is increased with a set of expressivity parameters that is taken into account while the system generate gesture animations. This gesture expressivity has not yet been studied in Salem’s robot system although it was mentioned in development of the Max agent [1].

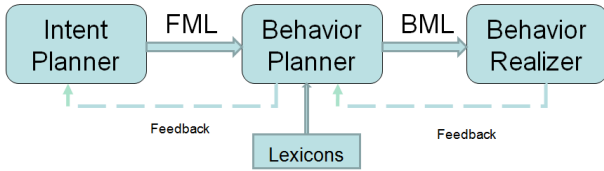


Figure 1: SAIBA framework.

An implementation and evaluation of gesture expressivity was done in the robot gesture generation system of Ng-Thow-Hing [23]. This system selects gesture types corresponding to input text through a parts-of-speech analysis. Then it schedules the gestures to be synchronized with speech using temporal information returned from a text-to-speech engine. The system calculates gesture trajectories on the fly from gesture templates while taking into account its style parameters. Differently from our model, his system was not designed as a common framework for both virtual and physical agents.

There are also other initiatives that generate gestures for a humanoid robot such as [24, 14] but they are limited in simple gestures or gestures for certain functions only. For instance pointing gestures in presentation [24].

All of the above systems have a mechanism to synchronize gestures with speech. Gesture movements are adapted to speech’s timing in [27, 23, 24]. This solution is also used in our system. Some systems have a feedback mechanism to receive and process feedback information from the robot in real-time, which is then used to improve the smoothness of gesture movements [27], or to improve the synchronization of gestures with speech [14]. They have also a common characteristic that robot gestures are driven by a script language such as MURML [27], BML [14] and MPML-HR [24].

### 3. SYSTEM OVERVIEW

Our system follows the architecture of the SAIBA framework [29] (cf. Figure 1). This architecture consists of three separated modules: (i) the first module, Intent Planner, defines the communicative intents that the agent aims to communicate to the users such as emotional states, beliefs or goals; (ii) the second, Behavior Planner, selects and plans the corresponding multi-modal behavior to be realized; (iii) and the third module, Behavior Realizer, synchronizes and realizes the planned behaviors. The results of the first module is the input of the second module through an interface described with the Function Markup Language (FML) [13]. The output of the second module is encoded the Behavior Markup Language (BML) [29], and then sent to the third module. Both languages FML and BML are XML-based and do not refer to specific animation parameters of agents (e.g. wrist joint). That means the Intent Planner and Behavior Planner modules in this platform are independent of the agent’s embodiment and the animation player technology.

The Behavior Realizer receives the BML message and instantiates the BML tags from either gesture repertoires (i.e. one repertoire for the virtual agent and another one for the physical robot) in order to schedule gesture phases and generate a set of gesture keyframes. This module is common to both agents. The next module, Animation Realizer, is responsible in generating the animation from the keyframes. Only, this module is specific to each agent. Figure 2 illus-

trates the data flow of our model. A message service system (i.e. in our case ActiveMQ) is used to exchange data in real-time between modules. The ActiveMQ facilitates us to integrate a new module into the system to send as well as receive messages from other modules.

The following subsections present in detail each process in the system.

## 4. GESTURE TEMPLATES

In our system, gestures are generated on the fly from abstract gesture templates in a gestuary that was introduced firstly by De Ruiter [5]. Each entry in a gestuary is a pair of two informations: the name of communicative intention and the description of gesture that conveys the given communicative intention. Gesture templates are described symbolically with a representation language as an extension of BML [29]. Their descriptions have no reference to specific animation parameters of agents (e.g. wrist joint).

Gesture is specified symbolically in the agent and robot lexicons. We rely on the theory of gestures of McNeill [21], the gestural hierarchy of Kendon [16] to specify a symbolic gesture. As a result, a gestural action may be divided into several phases of wrist movement, in which the obligatory phase is call stroke transmitting the meaning of the gesture. The stroke phase may be preceded by a preparatory phase which serves to take the articulatory joints (e.g. hand and wrist) to a position where the stroke occurs. After that it may be followed by a retraction phase that returns the articulatory joints to relax position or a position initialized for the next gesture. In our lexicons, only the description of the stroke phase is specified for each gesture. Other phases will be generated automatically by the system. A stroke phase is represented through a sequence of key poses, each of which is described with the information of hand shape, wrist position, palm orientation, etc. A trajectory type is declared as linear, curve, etc to indicate how to move from one key pose to another one.

## 5. FML-APML TO BML

The FML language has not yet been standardized so that we use our FML-APML language [19]. The FML-APML is based on the Affective Presentation Markup Language (APML) [6] and has similar syntax with FML [12].

A FML message includes two description parts: one for speech and another one for communicative intents. The description of speech is borrowed from the BML syntax. It indicates the text to be uttered by the agent as well as time markers for synchronization purposes. The second part is based on the work of Poggi [26]; it defines information on the world and on the speaker’s mind. In this part, each tag corresponds to one of the communicative intentions. Each intention has tag attributes to indicate its importance degree (probability to happen), timing (absolute or relative to the speech’s time markers), etc. The Behavior Planner selects from the agent’s lexicon the behaviors that convey specific communicative acts. It also calculates absolute start and end time for them, as well as values of expressivity parameters. A speech synthesizer (e.g. Acapela or OpenMary) is called in this module to create audio data and to instantiate time markers. The selected gestures and speech’s information are outputted within a BML message and sent to the Behavior Realizer module.

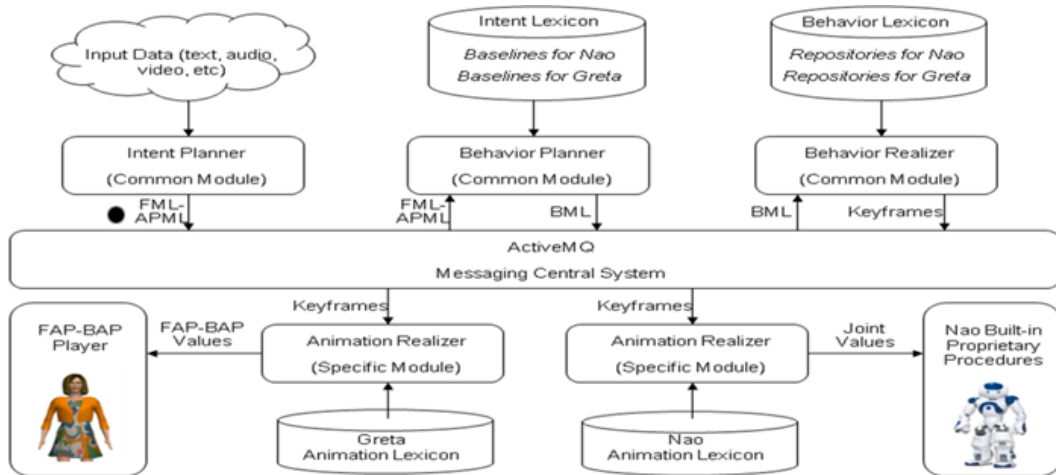


Figure 2: A Common Gesture Generation Framework for Virtual and Physical Agents.

## 6. BML TO KEYFRAMES

This process has two main tasks: scheduling gesture phases to synchronize with speech while taking into account the expressivity parameters and loading gestures from either gestural lexicons to create corresponding keyframes. Each keyframe contains the symbolic description and timing of each gesture phase. The symbolic representation of keyframes allow us to use the same algorithm for the synchronization of gestures with speech independently of the agent embodiment or animation parameters. Speech signal is also described within a keyframe. This keyframe indicates the audio source provided by the speech synthesizer as well as the start time to play this audio.

## SYNCHRONIZATION

In our system, the synchronization between gesture signal and speech is realized by adapting the gesture timing to speech. It means the temporal information of gestures within *bml* tag (i.e. for gesture phases) are relative to the speech. They are specified through time markers encoded by seven synchronization points: *start*, *ready*, *stroke-start*, *stroke*, *stroke-end*, *relax* and *end* [29]. The most meaningful part occurs between the *stroke-start* and the *stroke-end* (i.e. the stroke phase). The preparation phase goes from *start* to *ready*. In our system, the synchronization between gesture and speech is ensured by forcing the end time of the stroke phase (i.e. *stroke-end* sync point) to coincide with the stressed syllables. The duration of the preparation and stroke phase are hence pre-estimated so that the system can calculate exactly the time to start the gesture. This ensures that the stroke happens on the stressed syllables. This pre-estimation is done by calculating the distance between the current hand-arm position and the next desired position and by computing how long it takes to perform the trajectory. In case that the allocated time is not enough to do the preparation phase, the whole gesture has to be canceled, leaving free time to prepare for the next gesture. In other cases, if the allocated duration totally for a gesture is too long, a hold phase is added to keep this gesture movement more natural. The retraction phase is optional. It depends on its available time and also on the start time for the next gesture. This phase will be canceled if it has not enough time to move hands to

a defined relax position.

We apply the Fitts' Law (ie. simulating human movement law) [7] to have the natural movement speed. The parameters of Fitts' Law function is customized to adapt to each agent.

## GESTURE EXPRESSIVITY

The set of expressivity parameters is divided into two subsets. The first subset including spatial extent (SPC), temporal extent (TMP), stroke repetition (REP) is taken into account whilst the timing of gesture phases is calculated. The second subset including other parameters of the set (i.e fluidity, power, openness, tension of gesture movement) is applied when creating gesture animation. The reason is that the expressivity parameters in the second subset is dependent on the agents' embodiment. For instance the Nao robot does not support the acceleration modulation of gesture movements in real-time. In the first subset of expressivity parameters, the temporal extent(TMP) modifies the duration of a gesture. If the TMP value increases, the gesture lasts less. It means the speed of the movement is faster. However, in order to keep the synchronization with speech the time of *stroke-end* sync point can not be changed. Consequently the time of *stroke-start* and *start* sync points is later. On the contrary, their time is earlier if the TMP value decreases. Concerning spatial extent (SPC), it modulates the amplitude of gesture movements along the vertical, horizontal and depth dimensions. When a gesture is elaborated, certain dimensions are fixed to keep a gesture meaning. So that only re-sizable dimensions are affected by the SPC parameter. They are increased if the SPC value increases and vice versa. The REP parameter defines the number of repeating stroke phase in a gesture action. The duration of the complete gesture increases linearly with the REP value.

## 7. KEYFRAMES TO ANIMATION

The process to compute the animation from a given set of keyframes is specific to each embodiment. While all previous computations use the common agent framework, this stage is embodiment dependent. The following subsections present in detail how to calculate the values of the animation parameters for the Greta virtual agent and the Nao robot.

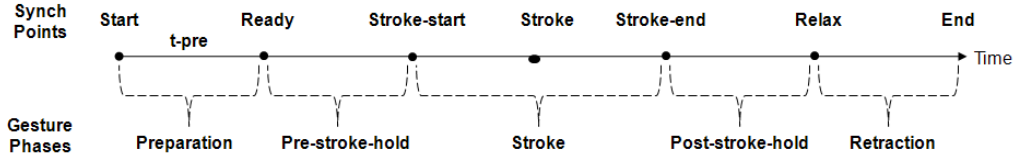


Figure 3: Standard BML synchronization points.

## 7.1 Generating Greta gesture animation

In this section, we present the implementation of our animation pipeline. It starts by receiving BML-like symbolic key frames time stamped in the motion planner. All key frames are received by streaming, and hence our animation computations need to be achieved on the fly. Each key frame includes gesture phases, expressivity parameters, gesture trajectory and the description of shape and motion for hand, torso, head, etc. We group keyframes per modalities, ie torso movements, head movements, arm gesture movements (two groups: left and right sides) in order to create a full body information. A key frame is defined by two computational attribute types: movement descriptions and targets to be reached through forward and inverse kinematics techniques. Direct movement descriptions are used to define forward kinematics (FK); the data can be abstracted from either motion capture or edited motion of different body parts. The targets will describe the gesture trajectory: we can perform a targeting process to reorganize the gesture trajectory that can take the form of line, curve, circle, and spiral. After this path targeting process, we obtain animation sequences for each body part (head, torso, gestures, etc). The next step is to gather these animation sequences into a single time stamps sequence covering the whole body. With this gathering process, we can create full body animation dependency, such as arm gestures influencing torso movements. This influence mechanism is part of the reaching model. We use forward kinematics to define the initial states for our agent skeleton system. Our IK method is applied to complete the key frames specification for the body. When the full body posture is computed, we apply re-targeting when processing the second subset of expressivity parameters (FLD, PWR, OPE, TEN) (see section Gesture Expressivity). We defined several different expressivity parameters. Using various easing functions to modulate speed and acceleration interpolation curves allows the simulation of PWR and TEN. The last process of our pipeline is to generate animation frames from key frames and finally to convert these animation frames into BAP (MPEG-4 body animation parameter) to animate our conversational virtual agent. This process is only performed in 3D rotation space. All the BAP frames are sent to the rendering and animation player.

## 7.2 Generating Nao gesture animation

Similarly to the Greta gesture animation module, this process receives and processes keyframes on the fly (through ActiveMQ). Then it translates keyframes into joint values of the robot. The second subset of expressivity parameters is applied in this stage.

To avoid singular positions in the gesture movement space of the robot, we predefine a set of wrist positions the robot can reach. In our case this set has 105 positions correspond-

ing to the key positions in McNeill’s gesture space [17]. The symbolic position of a gesture keyframe is instantiated with corresponding wrist position. From the actual position of the wrist, the palm orientation and hand shape are computed in real-time. The robot has only two hand shape configurations (i.e. open and close). The TMP value modifies the complete duration of a gesture, the PWR value modulates the acceleration of the movement of this gesture. For the Nao robot, while the movement acceleration cannot be modified, the system adjusts the duration of each phase of the gesture to simulate a change of movement speed. A hold time is also added after stroke phase when the PWR value increases to simulate a powerful movement. The Fluidity (FLD) parameter modifies the smoothness of single gesture and the continuity between consecutive gestures. It modifies the motion curve. However, the modification of the acceleration and trajectory curve is not available for the Nao robot so that we can not apply these changes. So far, the FLD value modulates the way that the robot link consecutive gestures. For instance when the FLD value increases, the movement between two consecutive gestures is smoother, the robot does a movement liaison from the first gesture without retraction phase to the second gesture.

Lastly all joint values with timing information are sent to the robot (as an animation layer). The animation is obtained by interpolating between joint values with the robot built-in proprietary procedures [8].

## Experimental results

The Nao’s gestures generation system was evaluated through perceptive tests. We wanted to evaluate how robot’s gestures were perceived by human users at the level of the expressivity, the naturalness of gestures and the synchronization of gestures with speech while the robot was telling a French tale [18]. 63 French speakers participated in our experiment. The results showed that the co-verbal expressive gestures generated by our model and displayed by the Nao robot were acceptable. 48 participants (76%) agreed that gestures were synchronized with speech and 44 participants (70%) approved that gestures were expressive. However, the naturalness of gestures were not appropriate and need to be improved in future work.

## 8. CONCLUSIONS

We have designed and implemented a framework to animate virtual and physical agents. This framework is as much as possible independent of the embodiment of the agents. Only the last step, consisting in interpolating keyframes into animation frames, is agent dependent. In our system a gesture lexicon is elaborated for each agent. It allows us to encompass variations and limitations of agent embodiments. Elements of the lexicon are stored using the same symbolic language. An extended set of expressivity parameters have

been implemented. The parameters act on the volume and dynamism of gestures production. Our gesture engine ensures also that the timing of gesture phases is synchronized with speech.

## 9. ACKNOWLEDGMENTS

The authors would like to thank André-Marie Pez for his help in implementing the system. This work has been partially supported by the French national projects ANR CECIL, GVLEX and IMMOMO.

## 10. REFERENCES

- [1] K. Bergmann and S. Kopp. Modeling the production of coverbal iconic gestures by learning bayesian decision networks. *Appl. Artif. Intell.*, 24(6):530–551, 2010.
- [2] C. Breazeal. Emotion and sociable humanoid robots. *Int. J. Hum.-Comput. Stud.*, 59(1-2):119–155, 2003.
- [3] J. Cassell, T. Bickmore, M. Billinghurst, L. Campbell, K. Chang, H. Vilhjálmsón, and H. Yan. Embodiment in conversational interfaces: Rea. In *Proceedings of the SIGCHI conference on Human factors in computing systems: the CHI is the limit*, pages 520–527. ACM, 1999.
- [4] D. Chi, M. Costa, L. Zhao, and N. Badler. The emote model for effort and shape. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 173–182. ACM Press/Addison-Wesley Publishing Co., 2000.
- [5] J. P. De Ruiter. *Gesture and Speech Production*. Doctoral dissertation at Catholic University of Nijmegen, Netherlands, 1998.
- [6] B. DeCarolis, C. Pelachaud, I. Poggi, and M. Steedman. Apml, a mark-up language for believable behavior generation. *Life-like Characters. Tools, Affective Functions and Applications*.
- [7] P. Fitts. The information capacity of the human motor system in controlling the amplitude of movement. *Journal of experimental psychology*, 47(6):381, 1954.
- [8] D. Gouaillier, V. Hugel, P. Blazevic, C. Kilner, J. Monceaux, P. Lafourcade, B. Marnier, J. Serre, and B. Maisonnier. Mechatronic design of nao humanoid. *The Int. Conf. on Robotics and Automation, 2009.*, pages 769–774, 2009.
- [9] M. Haring, N. Bee, and E. Andre. Creation and evaluation of emotion expression with body movement, sound and eye color for humanoid robots. In *RO-MAN, 2011 IEEE*, pages 204–209, 2011.
- [10] B. Hartmann, M. Mancini, and C. Pelachaud. Towards affective agent action: Modelling expressive eca gestures. In *International conference on Intelligent User Interfaces-Workshop on Affective Interaction, San Diego, CA, 2005*.
- [11] B. Hartmann, M. Mancini, and C. Pelachaud. Implementing expressive gesture synthesis for embodied conversational agents. *LNCS: Gesture in human-Computer Interaction and Simulation*, pages 188–199, 2006.
- [12] D. Heylen, S. Kopp, S. Marsella, C. Pelachaud, and H. Vilhjálmsón. The next step towards a function markup language. *Intelligent Virtual Agents*, pages 270–280, 2008.
- [13] D. Heylen, S. Kopp, S. Marsella, C. Pelachaud, and H. Vilhjálmsón. The next step towards a function markup language. pages 270–280, 2008.
- [14] A. Holroyd and C. Rich. Using the behavior markup language for human-robot interaction. In *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, pages 147–148. ACM, 2012.
- [15] T. Holz, M. Dragone, and G. O’Hare. Where robots and virtual agents meet. *International Journal of Social Robotics*, 1(1):83–93, 2009.
- [16] A. Kendon. *Gesture: Visible action as utterance*. Cambridge University Press, 2004.
- [17] Q. Le, S. Hanoune, and C. Pelachaud. Design and implementation of an expressive gesture model for a humanoid robot. *11th IEEE-RAS Humanoid Robots*, pages 134–140, 2011.
- [18] Q. A. Le and C. Pelachaud. Evaluating an expressive gesture model for a humanoid robot: Experimental results. *Submitted to 8th ACM/IEEE International Conference on Human-Robot Interaction, 2012*.
- [19] C. P. M. Mancini. The fml - apml language. *The First FML workshop*, 2008.
- [20] V. Manohar, S. al Marzooqi, and J. W. Crandall. Expressing emotions through robots: a case study using off-the-shelf programming interfaces. In *The 6th Int. Conf. on HRI*, pages 199–200. ACM, 2011.
- [21] D. McNeill. Hand and mind: What gestures reveal about thought. 1996.
- [22] M. Neff, M. Kipp, I. Albrecht, and H. Seidel. Gesture modeling and animation based on a probabilistic re-creation of speaker style. *ACM Transactions on Graphics (TOG)*, 27(1):5, 2008.
- [23] V. Ng-Thow-Hing, P. Luo, and S. Okita. Synchronized gesture and speech production for humanoid robots. *The Int. Conf. on Intelligent Robots and Systems (IROS’10). IEEE/RSJ*, 2010.
- [24] Y. Nozawa, H. Dohi, H. Iba, and M. Ishizuka. Humanoid robot presentation controlled by multimodal presentation markup language mpml. *Computer animation and virtual worlds*, pages 153–158, 2004.
- [25] C. Pelachaud. Multimodal expressive embodied conversational agents. pages 683–689, 2005.
- [26] I. Poggi, C. Pelachaud, and E. Caldognetto. Gestural mind markers in ecas. *Gesture-Based Communication in Human-Computer Interaction*, pages 481–482, 2004.
- [27] M. Salem, S. Kopp, I. Wachsmuth, K. Rohlfing, and F. Joubin. Generation and evaluation of communicative robot gesture. *International Journal of Social Robotics*, pages 1–17, 2012.
- [28] M. Stone, D. DeCarlo, I. Oh, C. Rodriguez, A. Stere, A. Lees, and C. Bregler. Speaking with hands: Creating animated conversational characters from recordings of human performance. *ACM Transactions on Graphics (TOG)*, 23(3):506–513, 2004.
- [29] H. Vilhjálmsón et al. The behavior markup language: Recent developments and challenges. *Intelligent Virtual Agents*, pages 99–111, 2007.
- [30] H. Wallbott. Bodily expression of emotion. *European journal of social psychology*, 28(6):879–896, 1998.