

# Generating animation from natural language texts and semantic analysis for motion search and scheduling

Masaki Oshita

Published online: 16 February 2010  
© Springer-Verlag 2010

**Abstract** This paper presents an animation system that generates an animation from natural language texts such as movie scripts or stories. It also proposes a framework for a motion database that stores numerous motion clips for various characters. We have developed semantic analysis methods to extract information for motion search and scheduling from script-like input texts. Given an input text, the system searches for an appropriate motion clip in the database for each verb in the input text. Temporal constraints between verbs are also extracted from the input text and are used to schedule the motion clips found. In addition, when necessary, certain automatic motions such as locomotion, taking an instrument, changing posture, and cooperative motions are searched for in the database. An animation is then generated using an external motion synthesis system. With our system, users can make use of existing motion clips. Moreover, because it takes natural language text as input, even novice users can use our system.

**Keywords** Computer animation · Motion database · Natural language processing

## 1 Introduction

Recently, computer animation has been widely used in movies, video games, TV programs, web graphics, etc. Because computer animation is a very powerful tool to present a story, drama, or instruction, there are demands from non-professional people to create computer animation. However,

it is a difficult task because of two main issues. The first issue is the difficulty of making and reusing motion data. Currently, motion data are mainly created using motion capture or keyframe techniques. Either way, they are very time consuming and require professional skills. Although there is demand for reusing existing motion data, this is difficult because of the lack of a system for storing and searching large amounts of motion data. Because there can be various motions of various characters, it is difficult to manage them in a standard file system or database. Currently, most motion data are created from scratch for individual scenes and are thrown away without reuse. The second issue is the limitation of current animation systems. A computer animation can be created by combining a number of existing motion clips using animation software such as MotionBuilder, Maya, 3ds Max, etc. However, it is difficult for novice users to utilize such software, because handling motion data is tricky and these systems require training.

To address these issues, we developed an animation system that generates an animation from natural language texts such as movie scripts or stories (Fig. 1). We also developed a motion database that stores many motion clips for different characters. When an input text is given, the system searches for an appropriate motion clip from the database for each verb. Temporal constraints between verbs are also extracted from the input text. The searched motion clips are scheduled based on the temporal constraints. In addition, when necessary, some automatic motions such as locomotion, taking an instrument, changing posture, and cooperative motions are searched from the database. The system outputs a motion timetable which consists of motion clips and their execution timings. An animation is then generated using an external motion synthesis system. Using our system, even novice users can create animation by making use of existing motion clips.

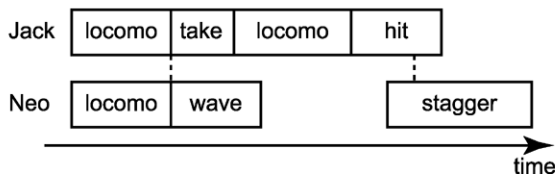
---

M. Oshita (✉)  
Kyushu Institute of Technology, 680-4 Kawazu, Iizuka,  
Fukuoka 820-8502, Japan  
e-mail: [oshita@ces.kyutech.ac.jp](mailto:oshita@ces.kyutech.ac.jp)

## (a) Input text

*Neo waves to Jack. At the same time, Jack takes the red bottle. Jack hits Neo with it.*

## (b) Output motion timetable



## (c) Generated animation



**Fig. 1** Example of our system. **a** Input text. **b** Searched motion clips and their execution timings. **c** Generated animation

There are many possible applications of our system. Recently, in movie production, simple animations are created before production to check camerawork, screenplay, necessary visual effects, etc. These animations are called “previsualization” or “animatics”. They are also often created for the scenes in which no computer graphics are involved. Using our system, even directors or writers who are not professional animators can create an animation very quickly. Moreover, our system can be used by non-professional people who want to make an animation but do not have professional skills. It can also be used for children to visualize a story to make it interesting and easy to understand. Our system can be used for movie production. Even though animators want to add more details to the output of our system, our method is much easier than making animations from scratch.

In this paper, we propose a motion frame that contains meta-information about a motion clip, an object-oriented database framework for storing a number of motions of a number of characters in a hierarchical structure, natural language analysis methods that are specialized for extracting motion related descriptions from an input text, and scheduling of multiple motions based on the temporal constraints in an input text. In addition, we have done preliminary experiments which showed that our system generates expected results from various input texts.

This paper is an extended version of our previous work [1]. As explained in Sect. 5, we have mainly extended our natural language analysis methods to enable our system to handle various expressions in input texts. Based on the experiments presented in Sect. 8, 87% of the verbs in a sample movie script can be dealt with using our methods and represented as motions, although 78% of these were handled by the system before the extension [1].

The rest of this paper is organized as follows. Section 2 reviews related work in the literature. Section 3 gives an overview of our system. Sections 4, 5, 6, and 7 describe our methods used in the framework of the motion database, natural language analysis, motion search and motion scheduling, respectively. In Sect. 8, some experimental results are

presented together with a discussion thereof. Finally, Sect. 9 concludes the paper.

## 2 Related work

Generating animation from natural language texts has been a challenge. Many research groups have tackled this problem. The SHRDLU system, which was developed by Winograd [2], is known as the pioneer. Using SHRDLU, a user can give commands to a robot using English in an interactive manner, and make it arrange objects in a scene. However, the types of commands were very limited.

Badler et al. [3, 4] developed virtual agents that follow natural language interactions. They proposed Parameterized Action Representation (PAR), which has a similar purpose to the motion frame in our research. The PAR has more complex information such as pre-condition and achievement. The motion generator of each PAR is programmed using a state machine. It can use motion data or any motion generation methods. However, specifying detailed information and constructing motion generators are very time consuming.

Tokunaga et al. [5] developed the K2 system, which has similar goals to Badler et al. In their system, agents are controlled via spoken language. Their research is rather focused on solving the vagueness of natural language instructions. They use case frames [6] to search for motions. Unlike our work, they use all cases that are used in linguistic analysis. The interpretation of each case is left to the user who adds the case frame handler. The motion generator for each case frame must be manually programmed by the user.

These previous works aim at developing intelligent agents that understand natural language instructions and make plans to execute them. However, the systems are very complex, and many rules are required. On the other hand, our system aims to reuse existing motion data easily and efficiently. The motion frame in our work contains just enough information to search for appropriate motions that match natural language texts and it is easy to describe. We believe that our system is more practical.

Lu and Zhan [7] developed an animation production system that includes story understanding, plot planning, act planning, camera planning, etc. Although their system takes simple Chinese as input, it requires a great deal of additional knowledge, including not only case frames but also many dictionaries, templates and rules.

Sumi et al. [8] developed a system for visualizing short stories for children. The system extracts keywords from an input text, and chooses an appropriate scene, characters, and motions from a database. It simply plays a motion that matches the keywords. Although a user can add motion data to the system, the system cannot select motions appropriate for the objects or characters and cannot generate interactions between characters and the scene.

There is very little research that deals with motion scheduling from natural language texts. The above systems simply execute motions as instructions are given or events happen, and no scheduling is considered. However, in order to execute multiple motions of multiple characters as instructed by an input text, the execution timing of the motions must be coordinated. Baba et al. [9] developed a system for generating an animation that satisfies temporal and spatial constraints given by natural language texts. The system determines appropriate initial positions of the agents and objects that are specified in the input text. However, the motions of the agents and motion scheduling were not considered.

Coyne and Sproat [10] developed WordsEye, which converts natural language texts to a scene. Because their purpose is to generate a still image, when a character motion is indicated in a given text, the system simply chooses a pose for the action from the database.

There have been various studies on generating a character's gestures for a monologue or conversation [11]. These methods generate motions by composing short fragments of motions based on signal processing of the input speech rather than by interpreting the meaning of the speech.

There are also animation engines that support some script language such as Improv [12] and Alice [13]. However, it is still difficult to program the agents and to make use of a large amount of existing motion data. In addition, markup language formats for describing animation including scenes, characters and actions have been proposed [14, 15]. However, they are difficult to describe by hand. The animation files should be created by using specific authoring software. Moreover, it is difficult to add and reuse motion data using such file formats and authoring software.

There are many motion synthesis methods which generate new motions from a small number of motions [17, 18]. However, they require a manual setup for each motion module. It is difficult for end users to add new motion modules. Although currently our system selects one motion from the

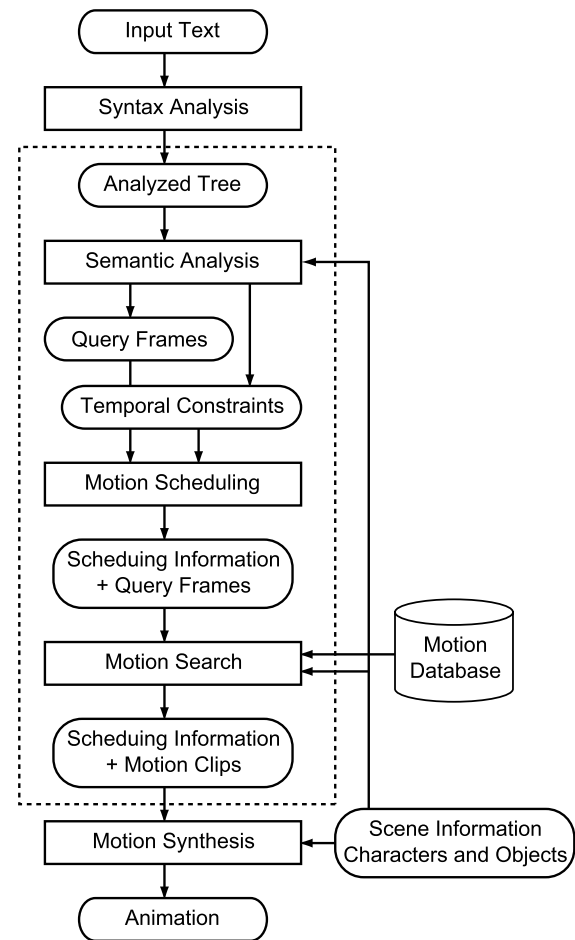


Fig. 2 System overview

database, it is possible to extend our system to blend a number of selected motions based on quantitative motion query parameters such as contact position.

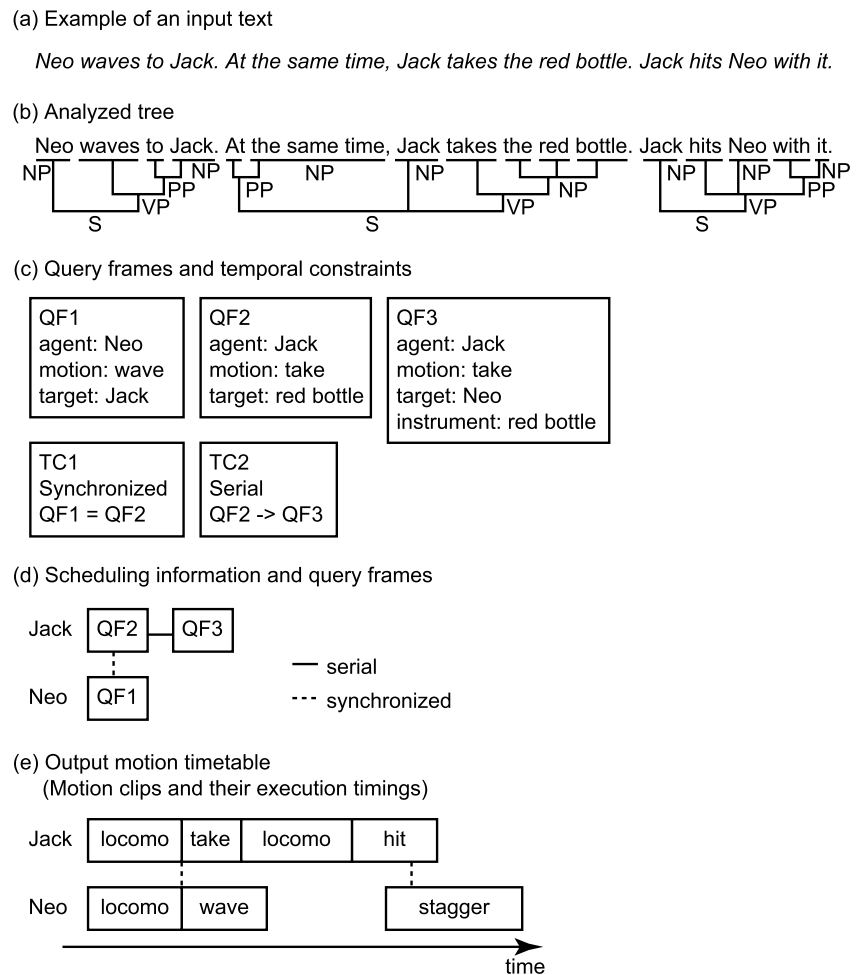
### 3 System overview

In this section, we explain the overview of our system (Fig. 2) and data representation (Fig. 3).

When an input text is given to the system, natural language processes (syntax analysis and semantic analysis) are applied first. The syntax analysis is the process of converting a plain text to a tree structure with phrase tags and dependencies. Figure 3(b) is an example of the analyzed tree which is computed from an input text (Fig. 3(a)). The type of each phrase and the dependency between phrases are determined. For example, S, NP, VP and PR in Fig. 3(b) represent sentence, noun phrase, verb phrase and preposition, respectively.

The semantic analysis extracts information about motions described in the input text from the tree structure. A query frame contains information for the motion search.

**Fig. 3** Example of data representation



One is generated for each verb in the text. The temporal constraints contain information about execution timing between verbs. For example, QF1~QF3 and TC1~TC2 in Fig. 3(c) represent query frames and temporal constraints, respectively.

Based on the temporal constraints, motion scheduling determines the execution order of each motion clip, which corresponds to each query frame as shown in Fig. 3(d). Note that exact execution times are not decided at this point, because the duration of each motion is not known until motion clips are searched from the database and automatic motions are added later.

The motion search is applied for each query frame. In addition, when it is necessary, automatic motions are inserted before the motion. Finally, motion clips and their execution timings are passed to the motion synthesis module as a motion timetable, as shown in Fig. 3(e).

The motion synthesis generates an animation by smoothly connecting given motion clips. The interactions between characters and between a character and objects are handled by this module based on the information that the motion clips have.

The scene information contains characters and objects and their initial states, including postures, positions, and orientation. Each object has certain object information including names, a default contact point, alternative contact points and their names. For example, a desk object has “desk”, “table”, etc. as its names. An object also has sets of pairs consisting of a part name and its position (e.g., “above”, “under”, “side”, etc.). This information is used to search for appropriate motions and determine appropriate contact positions according to an adjective that is used in the input text. In addition, an object has a default position which is used when no adjective is specified. This kind of object information is commonly used in similar approaches [5, 10]. In addition, a scene also has default entering and leaving points as the default goal locations for locomotive motions (see Sect. 5.5). Currently, our system assumes that the scene information is provided in advance by the user.

The scope of this paper is the components in the dotted box in Fig. 2. There are many tools for syntax analysis that can be used with our system. The Stanford parser [19] is used for our implementation. For motion synthesis, our system uses an external animation system [20]. The sys-

tem generates continuous motions from given motion clips and their execution timings. The system determines an appropriate synthesis method for each transition based on the constraints between the foot and the ground during motions. Alternatively, another commercial animation system such as MotionBuilder, Maya, 3ds Max, etc. can be used.

#### 4 Motion database

In this section, we describe the representation of motion data. We first explain the case frame that is used in natural language processing. Then, we explain our motion frame, which is inspired by the case frame. We also describe our database of characters and motions.

##### 4.1 Case frame

The idea of a case frame was proposed by Fillmore [6]. A case frame represents the role of a verb. Each case of a case frame is a phrase that represents an aspect of the verb. Typically a case frame has the following cases:

- Agent: the person who performs the motion.
- Experiencer: the person who experiences something.
- Object: the object that an effect is caused to during the motion.
- Instrument: the object that causes an effect during the motion.
- Source: the source or origin of the motion.
- Goal: the goal or target of the motion.
- Time: the time when the motion is performed.
- Location: the location where the motion is performed.

Each case needs to be a specific type of entity. Some cases are mandatory for some verbs. A verb that has different roles depending on context has multiple case frames.

In general natural language processing systems, a procedure to select a case frame for an input text is as follows. First, based on the types and dependency of phrases in the analyzed tree, candidate cases of each phrase are determined. By searching for case frames that match the candidate cases, the most appropriate case frame and all its cases are determined.

The case frame is a good way to extract and represent the meanings of texts. The case frame is widely used in many research papers such as [5, 10]. However, the case frame is not suitable for representation of motion data for animation. From the view point of motion representation, each case has different roles depending on case frames. For example, the “object” case of a case frame could be an object that the character uses or another character that the character’s motion causes an effect on. Moreover, the case frame does not contain information about postures and contact positions, which are important for selecting motions.

Item	Value
Agent	human
Names of Motion	take, pick up, get
Instrument	NULL
Target	appropriate size and weight ranges
Contact Position	hand position of contact
Target Direction	NULL
Initial Posture	standing
Adverbs	slowly

Fig. 4 Example motion frame of “taking-an-object”

##### 4.2 Motion frame

We propose a motion frame which contains the information about a motion clip. The motion frame is inspired by the case frame. However, we define the items of the motion frame based on importance when we search for a motion according to input texts.

There are many kinds of verbs in general English. However, our system handles only action verbs that involve a physical motion, in other words, verbs that can be visualized as an animation. Other kinds of verbs such as non-action verbs (e.g., “think”, “believe”) or state verbs (e.g., “know”, “exist”) are ignored in our system because they are difficult to represent by a motion clip. Action verbs are categorized into intransitive, transitive, and ditransitive verbs. Intransitive verbs involve no other object (e.g., “he runs”). Transitive verbs include one target object/character/position (e.g., “he opens the door”, “he hits her”, “he walks to the door”). Ditransitive verbs include two target objects (e.g., “he gives her the book”, “he cuts the bread with a knife”). For distractive verbs, one of the two target objects should be the object that the character possesses. We call such objects “instruments”. Therefore, action verbs have at most one “target” object /character/position and at most one “instrument” object. We use them as items of a motion frame instead of cases in a case frame. In addition, contact position is used to select a motion that fits the environment and previous motions.

The items of the motion frame are as follows. An example of a motion frame is shown in Fig. 4. Note that some items may not have any value depending on the motion.

- Agent  $M_{agent\_ref}$ : The reference to the character in the database who performs the motion.
- Names of motion  $M_{motion\_strings}$ : The set of verbs that represent the motion. When a verb in the input text matches one of the motion names, the motion frame will be a candidate for the verb. To handle ambiguity, a motion frame may have multiple names. For example, a “taking-an-object” motion may have “take” and “pick up” as its names.
- Instrument  $M_{instrument\_ref}$ ,  $M_{instrument\_params}$ : The object that the character uses in the motion. This is either a reference to an object in the database  $M_{instrument\_ref}$  or the



- size and weight ranges of an object  $M_{instrument\_params}$ . If the motion requires a specific object such as “cutting with a knife”, the object should be specified as a reference to the instrument. Otherwise abstract conditions of an object are specified. For example, if the motion is “poking something with a long object”, then appropriate size and weight ranges of the object are specified.
- Target: The reference to an object  $M_{target\_ref}$  or the size and weight ranges  $M_{target\_params}$  are specified in the same way as the instrument. If the target is a character, the reference to the character is specified in  $M_{target\_ref}$ .
  - Contact position  $M_{contact\_vertical}$ ,  $M_{contact\_horizontal}$ : the position of the end-effector when it makes contact with the target. A contact position is specified when the motion involves contact with a target character or object. Vertical and horizontal positions are handled differently. Because the horizontal position can be adjusted by lateral movement (see Sect. 7.2), the vertical position is more important for motion selection. For example, if multiple “taking an object” motions are in the database and an input text “he takes the bottle on the ground” is given, then based on the position of the bottle, the appropriate taking motion (e.g., “taking an object with squatting”) will be selected. The contact position is automatically computed from the contact information (see Sect. 4.3) of the motion data. The contact position is expressed in the local coordinates of the motion data.
  - Target direction  $M_{target\_direction}$ : The direction of the target. For some motion, even though the motion does not involve contact with the target, the target direction is important. For example, when “waving to a person” or “shooting a target” motion is executed, the character should face the right direction. For some motion, both contact position and target direction are specified. For example, “sitting down on a sofa” motion should make contact with a sofa from the front of the sofa.
  - Initial posture  $M_{initial\_posture\_flag}$ : the character’s posture when the motion begins. Currently, it is represented as one of three states: standing, sitting, or lying down. The initial posture is used to select a motion that matches the terminal posture of the previous motion. In cases where no such motion is in the database, an automatic changing posture motion will be added (see Sect. 7.2).
  - Adverbs  $M_{adverb\_strings}$ : The set of adverbs represent the style of the motion such as “slowly” or “happily”.

Each item of motion frames must be specified by a user. However, this is not such a difficult task for users. For each motion frame (each motion clip), the user is asked to specify the agent, verbs, target, and instrument. The agent is selected from the character database. For the target and instrument, it is either an appropriate object or agent that is selected from the database or the size and weight range of an object. When the motion involves a specific object (e.g., “cutting with a

sword”), the object should be selected. Otherwise, object conditions are specified (e.g., “lifting up a light object using one hand”). The contact position is automatically computed from the motion and its contact information (see Sect. 4.3). The initial posture is also automatically computed from the motion clip. As a result, specifying the items of a motion frame is very easy.

#### 4.3 Motion data

Our system supposes that each motion is short and simple. A complex motion is difficult to represent by a motion frame. If a user wants to add a long motion to the database, the motion should be divided into pieces.

Some motions involve an interaction with an object or a character. This information is very important for generating animation and for selecting motions. Therefore, it is specified on the motion frame. The contact information consists of the contact type (hold, release, or hit), contact time (local time in the motion clip) and the end-effector (e.g., right hand). This information is also necessary for generating animation in the motion synthesis module (see Sect. 6.2).

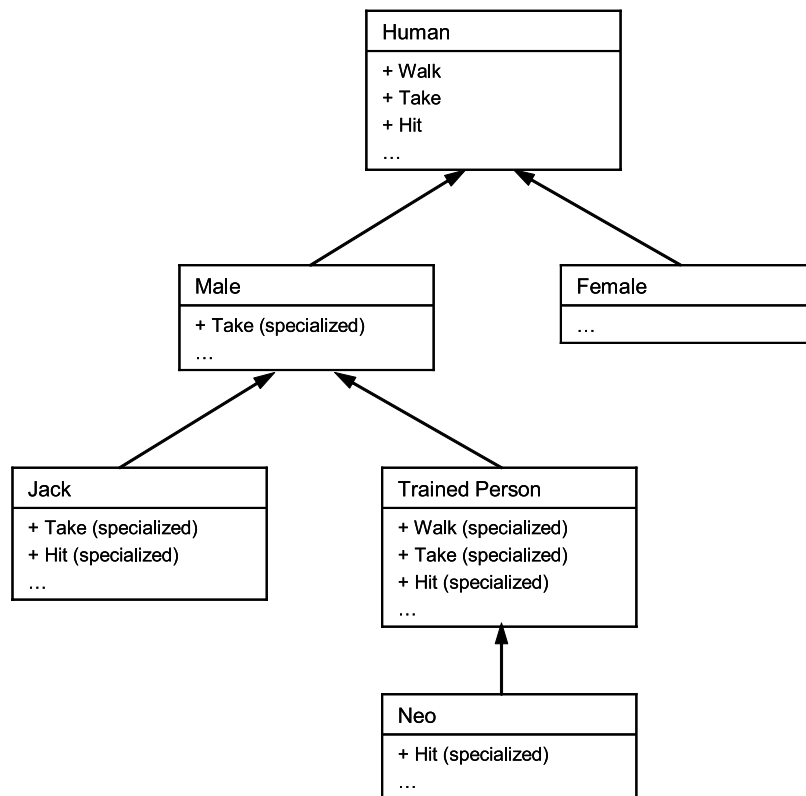
Some motions that interact with another character cause the reaction of the other character (e.g., “Jack hits Neo. Neo falls”). Usually such cooperative motions are captured or created at the same time but are stored as separate motion clips. In our system, such cooperative motions are specified on the motion frame. If a motion has cooperative motions and no cooperative motion is indicated in the input text, the system automatically executes a cooperative motion (see Sect. 7.2). In addition, when two cooperative motions include physical contact, the timings and the initial positions of these motions are coordinated (see Sect. 7.1).

#### 4.4 Character and motion database

We use an object-oriented framework for the character and motion database. As shown in Fig. 5, each character is considered to be an object that has various motions as its methods. A character inherits from a base character. A motion of the base character can be overridden by another motion. The motions that are not overridden are used as the motions for the derived character. In this way, the hierarchy of characters and their motions are efficiently managed. A character can inherit from multiple base characters. All motions that the base characters have are used for the derived character. Since the motion that most closely matches an input sentence is selected from the available motions, even if there are multiple motions with the same name, there is no problem with conflicts caused by the multiple inheritance.

If a user wishes to create a new character, she/he simply adds the new character that inherits from a base character or multiple base characters to the database and adds character-specific motions to that character. Even if there are not many

**Fig. 5** Example of a hierarchical database of characters



new motions for the new character, the motions of the base characters are used. In this way, users can add new characters very easily.

The database can be implemented in various ways. If the characters and motions are implemented using an object-oriented programming language (e.g., C++ or Java), we would represent motions as objects rather than methods and implement a mechanism of motion inheritance on the character class, because it is practically difficult to handle motions as methods using such programming languages.

### 5 Natural language analysis

Although natural language processing techniques have advanced in recent years, it is still a challenge to understand general texts because it requires not only language processing but also a large knowledge of the world. However, our system is supposed to take script-like text and only motion-related descriptions in the text matter. This makes the natural language analysis much easier than general natural language processing systems such as machine translation or summarization systems. Moreover, because scene information, such as characters and objects, is given in advance, we do not need the same large dictionary required by general natural language processing systems.

As explained in Sect. 3, the semantic analysis takes an analyzed tree and generates query frames and temporal constraints. A query frame contains information of a verb for the motion search. The temporal constraints contain information about the execution timing between verbs. In the followings of this subsections, we explain how the semantic analysis works.

#### 5.1 Query frame

To select a motion that matches an input text, we use a query frame, which has the same items as the motion frame, and whose items are determined by analyzing the syntax tree of the input text (see Fig. 3(b)). Scene information is also used to determine some items.

As explained in Sect. 4.2, unlike generic semantic analysis, motion searches only need a target and an instrument for each verb. Therefore, we determine these by applying the following rules to each verb in the input text.

- A verb is used as the name of motion of the query frame  $Q_{motion\_strings}$ . If the verb is followed by a preposition or noun, then all sets of the verb and the following word are also set to  $Q_{motion\_strings}$  because this could represent an idiom. Therefore,  $Q_{motion\_strings}$  can contain multiple phases. For example, in “Jack falls back”, both the phrases “fall” and “fall back” are set to  $Q_{motion\_strings}$ .

- If a noun represents a character in the scene and the verb is dependent on the noun, the character is considered as the agent (subject) of the query frame  $Q_{agent\_ref}$ .
- If two nouns are dependent on the subject that the verb is related to, they are considered as the target  $Q_{target\_ref}$  and the instrument  $Q_{instrument\_ref}$ . (E.g., in “Jack gives Neo the book”, “Neo” is the target and “the book” is the instrument.)
- If only one noun is dependent on the subject, it is considered as the target  $Q_{target\_ref}$ .
- If a preposition phrase (e.g., “to Neo”) is dependent on the subject, it is considered as the target  $Q_{target\_ref}$  or the instrument  $Q_{instrument\_ref}$  depending on the preposition. If the preposition is “with” and the noun in the phrase represents an object, the object is used as the instrument. Otherwise, the noun is used as the target.
- If the character is holding an object, the object is also used as  $Q_{instrument\_ref}$ , even if it is not specified in the input text.
- If a phrase considered to be an adverb in the syntax analysis is dependent on the subject, the phrase is used as one of the adverbs  $Q_{adverb\_strings}$  that can contain multiple phrases.

After the names of the target and instrument are determined, we obtain the reference or value of each item from the scene information. We suppose that the characters or objects in input texts always exist in the scene. Therefore, unlike general semantic analysis, by looking up the scene information all nouns in input texts are determined.

The target character or object that is indicated in the input text is searched from the scene information and the reference and position are set to the query frame. When the target is a character and a body part is indicated in the text such as “She hit him in the head”, the reference and position of the body part is set. When the target is an object in the scene, the target size and weight are set in the query frame. The contact position  $Q_{contact\_vertical}$ ,  $Q_{contact\_horizontal}$  and the target direction  $Q_{target\_direction}$  are set based on the position and direction, respectively, of the target character or object. If an adjective is used in the input text (e.g., “top of the table”, “under the table”) and the object has the corresponding part (see Sect. 3), the position of the corresponding part assigned to the object is used as the contact position. If there is no adjective, the default position specified for the object is used. The instrument object that is indicated in the input text is also set to the query frame.

## 5.2 Temporal constraints

Temporal constraints are extracted from input texts. The types of temporal constraint are serial execution or synchronized execution between two verbs. A serial execution

constraint has the execution order of two motions. A synchronized execution constraint has relative execution timing. Temporal constraints are generated from a syntax tree as follows:

1. For all pairs of sequential verbs in the input text, serial execution constraints are assigned. For example, when the input text “Jack walks in the room. Neo stands up.” is given to the system, a serial execution section constraint (Jack, walk) to (Neo, stands up) is generated.
2. When a word that indicates a reverse order exists in the input text (e.g., “after”), the order of the serial execution constraint is reversed. If a serial execution constraint is already created, the old constraint is overridden. For example, when the input text “Jack walks in the room after Neo stands up.” is given to the system, a serial execution constraint (Neo, stands up) to (Jack, walk) is generated.
3. When a word that indicates synchronization exists in the input text (e.g., “at the same time” or “while”), a synchronized execution constraint is added. If there is a conflicting constraint, it is overridden. For example, when the input text “Jack walks in the room. At the same time, Neo stands up.” is given to the system, a synchronized execution constraint (Neo, stands up) and (Jack, walk) is generated. The relative timings between two motions are set to zero so that they start at the same time.
4. When the motions of two characters are cooperative motions and they include contact with each other, a synchronized execution constraint is added and the relative execution timings of the two motions are determined based on their contact information (Sect. 4.3). For example, when the input text “Jack hits Neo. Neo falls” is given to the system, a synchronized execution constraint (Jack, hit) and (Neo, fall) is generated. At this point, the relative timings are not set. They will be set based on the contact times in the searched motion data, when the motions are searched later.

## 5.3 Adjective and pronouns

Sometimes a character or an object is referred to by a combination of adjectives and a noun instead of its name. In this case, the system has to determine to which character or object the expression refers. Moreover, when a pronoun (e.g., “he”, “she”, “it”, “they”) is used in the input text, the system has to determine to which character or object the pronoun refers. This process can be difficult especially when ambiguous or euphemistic expressions are used. However, since our system is meant to take simple script-like texts, we handle these problems using the following method.

Each character and object in our database has a list of adjectives, nouns, and pronouns by which the character can be referred to. For example, a male soldier character would have “he”, “they”, “man”, “guy”, “soldier”, etc. as its list of



adjectives and pronouns. The list is inherited from the base character. If the adjective, noun, or pronoun appearing in the input text represents a single character or object (e.g., “he”, “the soldier”, “it”), the system searches for a character or an object that matches the words.

However, if the words are ambiguous, meaning that there are multiple characters or objects matching the given adjective, noun, or pronoun in the scene, the system has to choose an appropriate character or object. Basically, if the noun or pronoun represents a character or an object, this should be mentioned in the previous sentences. A character can be the agent  $Q_{agent\_ref}$  or the target  $Q_{target\_ref}$  of the query frame, while an object (e.g., “it”) can be the target  $Q_{target\_ref}$  or the instrument  $Q_{instrument\_ref}$  thereof. By using this constraint, the noun or pronoun is determined as follows.

1. If one of two items is clearly mentioned in the input text, the other should be a different character or object. For example, if the input text is “Neo comes to Jack. Jack gives him a book”, the word “him” cannot represent Jack since Jack is already the agent  $Q_{agent\_ref}$  of the query frame of the second sentence. Therefore, in this case, the system searches for characters in the previous sentence and uses the other character Neo as the target  $Q_{target\_ref}$ .
2. If both of the items for characters are pronouns, the agent of the previous sentence is also used as the agent  $Q_{agent\_ref}$ . For example, if the input text is “Neo comes to Jack. He gives him a book”, the system decides that the first “he” in the second sentence represents Neo. In the same way, if both the items for objects are pronouns, although this is not common, the target of the previous sentence is used as the current target and the instrument of the previous sentence is used as the current instrument.

If the noun or pronoun represents multiple characters (e.g., “they”, “soldiers”), the sentence should be represented by multiple motions. Therefore, in such a case, multiple query frames are generated. For example, if the input text is “Neo hits the soldiers” and there are two soldiers A and B in the scene, two query frames are generated with all items except the target character the same. The same rule is applied when an item of the query frame represents multiple characters (e.g., “Neo and Jack walk.”). However, if a motion frame matching the motion name has a target direction, but not a contact position, the center position of all characters referred to is used as the target direction, instead of generating multiple query frames. For example, with the input “Neo shoots the soldiers”, the “shooting” motion frame has only a target direction, and a query frame whose target position is the center of the soldiers is generated.

#### 5.4 Infinitives and gerunds

Infinitives and gerunds are often used with a verb. In this case, the system generates appropriate query frames and temporal constraints depending on the verb.

- If the verb is “do”, “perform”, etc. (e.g., “Neo performs dancing”), the infinitive or the gerund is represented as the motion and is used to generate a corresponding query frame. In this case, the verb is not involved in the query frame.
- If the verb is “start”, “try”, etc., the infinitive or the gerund is used to generate a corresponding query frame in the same way. However, in this case, the next event is considered to happen before the motion finishes. Therefore, a temporal constraint is generated to execute the next motion just after this motion starts.
- If the verb is “repeat”, “keep”, etc., the infinitive or gerund is used to generate a corresponding query frame. The motion is repeated before the next motion starts, and is therefore, specified in the query frame. This information is used for motion scheduling to duplicate a motion when it can be executed more than once.
- If the verb does not fall into any of the above cases, both the verb and the infinitive or gerund are represented as motions (e.g., “Neo walks waving to Jack”). In this case, multiple query frames are generated. In addition, temporal constraints are generated to execute all motions at the same time.

The system uses a dictionary of pairs consisting of the verb and the corresponding method, to determine which method should be applied.

In addition, if a gerund exists on its own in the input text (“Neo walks to the door, waving to Jack”), a query frame is generated for the gerund and a temporal constraint is generated to execute the gerund and the verb in the sentence at the same time. In this case, the agent of the verb becomes the agent of the gerund as well.

#### 5.5 Locomotive motions

Locomotive motions such as walking and running require special care because the target position and path vary depending on the situation and appropriate motions should be generated instead of simply executing a motion in the database. How to generate locomotive motions is explained in Sect. 6.3. In this section, we explain how to handle locomotive motions in natural language analysis.

As discussed in Sect. 8, natural language is not suited to specifying the locomotion path. Therefore, our system currently does not handle it and only determines the target position of locomotive motions. We categorize locomotive motions into the following types depending on how the target position is handled.

- Moving to a target position. If the verb is a locomotive motion (e.g., “walk”, “run”, “go”, etc.) and a target position is explicitly specified in the input text (e.g., “Neo walks to the door.”), the query frame includes the target

position  $Q_{contact\_vertical}$ ,  $Q_{contact\_horizontal}$  and a flag indicating that this is a locomotive motion.

- Entering and leaving. Sometimes the target position is not specified in the input text. In this case, if the verb is a specific verb, the locomotive motion is handled as an entering or leaving motion. For example, if the verb is “leave”, “walk away”, “disappear”, etc., the verb is handled as a leaving motion and the leaving position specified in the scene information is used as the target position. The query frame also includes a flag indicating locomotive motion.
- Simple walking. If a target position is not specified and the verb represents a locomotive motion excluding an entering or leaving motion (e.g., “Neo walks.”), a walking motion is simply executed. In this case, the query frame does not include a flag indicating locomotive motion. This query frame is handled in the same way as the other query frames. As a result, a motion of walking forward from the character’s current position is selected and executed.

## 5.6 Adverbs

Adverbs are handled in different ways depending on the word. As explained in Sect. 5.2, if the adverb represents temporal information, an appropriate temporal constraint is generated. If the adverb represents the frequency or timing of executing a verb, the adverb is handled in the same way as infinitives and gerunds in Sect. 5.4. For example, if an adverb such as “repeatedly”, “twice”, etc. is specified, the third option in Sect. 5.4 is applied. The system has a dictionary of adverbs for these cases. If the adverb is not found in the dictionary, it is assigned to a query frame to search for an appropriate motion as explained in Sect. 5.1.

## 6 Motion search

In this section, we explain how to search for an appropriate motion for each verb in the input text. Handling multiple verbs and motions is dealt with in the next section. A query frame is generated for each verb as explained in the previous section. Based on the query frame, a motion is selected from the database.

### 6.1 Evaluation of motion frame

A motion frame that best matches the query frame is searched for in the database. This search is performed in three steps.

In the first step, all candidate motion frames in which the motion name and agent match the query frame are selected from the database. All motion frames with the agent character or its base characters are potential candidates.

In the second step, the motion frames whose items do not match the query frame are excluded as candidates. If the

query frame has a target  $Q_{target\_ref}$ , or  $Q_{target\_params}$  and/or an instrument  $Q_{instrument\_ref}$ , or  $Q_{instrument\_params}$  but the motion frame does not, then it is excluded. Moreover, if a motion frame has target parameters, instrument parameters, or the vertical contact position, and the values of the query frame exceed the specified ranges, then that motion frame is also excluded.

In the third step, all candidate motion frames are evaluated based on the similarity between the motion frame and the query frame items using the following equation:

$$\begin{aligned}
 E = & w_0 R(M_{target\_params}, Q_{target\_params}) \\
 & + w_1 R(M_{instrument\_params}, Q_{instrument\_params}) \\
 & + w_2 D(M_{contact\_vertical}, Q_{contact\_vertical}) \\
 & + w_3 D(M_{contact\_horizontal}, Q_{contact\_horizontal}) \\
 & + w_4 D(M_{target\_direction}, Q_{target\_direction}) \\
 & + w_5 F(M_{initial\_posture\_flag}, Q_{initial\_posture\_flag}) \\
 & + w_6 A(M_{adverb\_strings}, Q_{adverb\_strings}) \\
 & + w_7 H(M_{agent\_ref}, Q_{agent\_ref}), \quad (1)
 \end{aligned}$$

where  $R(M, Q)$ ,  $D(M, Q)$ ,  $F(M, Q)$ ,  $A(M, Q)$ ,  $H(M, Q)$  are the functions that compute normalized distance (0.0 ~ 1.0) between size and weight parameters, contact positions, posture flags, adverbs, and hierarchical positions, respectively. The distances between the size and weight range of the motion frame and the object size and weight of the query frame are computed so that the distance becomes zero when the values are at the center of the range and the distance becomes one when the values are at the edge of the range. The distance between posture flags is computed in such a way that the distance is zero when they match and otherwise the distance is one. The distance between adverbs is computed so that the distance is zero when there is at least one pair of matching adverb between the motion frame and query frame and otherwise the distance is one. The distance between hierarchical positions of the characters is computed from the number of inheritances between them (see Fig. 5). The candidate motion frame whose evaluation is the smallest will be selected and used for animation.  $w_0 \sim w_7$  are weight parameters. They can be set for each motion frame in the case that some items are important for the motion. In our current experiments, we used 1.0 for all weights on all motions.

### 6.2 Motion modification

The motion clip of the selected motion frame is used for animation. However, even if the closest motion frame is selected, the contact position may not exactly match the query frame. In that case, the motion clip is modified using inverse

kinematics. The posture of the character during the motion is modified so that the contact position of the end-effector (e.g., hand) matches the target position in the query frame.

When the character is far from the target, changing the end-effector position is not enough. In addition, when the character executes the selected motion it may need to first take an instrumental object or change its posture (e.g., standing up). These cases are handled by adding automatic motions before the selected motion instead of modifying the selected motion. Automatic motions are explained in Sect. 7.2.

### 6.3 Locomotive motion

When a query frame indicates a locomotive motion (see Sect. 5.5), appropriate motion enabling the character to move to the target position must be generated. Several methods, such as [16], have been developed to generate walking motions. In our system, the character should not only walk, but also turn and step, in order to move to the appropriate position and direction. Therefore, we generate locomotive motions based on the target position and/or target direction according to the steps below using the set of motion data that the character has.

1. If the target direction is specified and the target position is not, an appropriate turning motion is generated. An appropriate motion based on the target direction is selected from the motions with ‘turn’ as their motion name in the database. If only the target direction is specified, the process stops here.
2. If the target position is specified and it is not in front of the character, a turning motion is added in the same way as in the first step so that the agent faces the target position.
3. If the target position is within one step, a stepping motion is added in the same way as the turning motion. The motion is selected from all ‘step’ motions.
4. If the target position is more than one step in the distance, a walking motion is added. The walking motion is repeated until the agent reaches the target position. The step length in each walking cycle is adjusted so that the walking cycle ends at the target position. The motion is selected from ‘walk’ motions. Currently, our system generates a straight path to the target position even if there are obstacles.
5. If the target direction is specified and it does not match the character’s direction at the end point of the walking motion, a turning motion is once again added.

As explained above, the system uses the “turn”, “step” and “walk” motions that the character has. Currently, the system selects an appropriate motion and modifies it if necessary. Alternatively, motion blending can be used to generate more appropriate motions by using multiple motions [16, 18].

## 7 Motion scheduling

In this section, we explain how our system handles multiple motions from an input text. Basically, the system searches for a motion for each verb in the input text. However, in order to make an animation, the execution timing of each motion must also be determined. Moreover, the continuity of motions should be considered. For example, when a character makes contact with an object in the scene, the character must first move close to the object. Our system takes care of this kind of continuity of motions.

When multiple characters perform multiple motions the motions should be scheduled. However, an exact execution time for each motion is not usually specified in the input text. In order to determine the motion schedule, we need information about the motions such as duration and contact information.

Our motion schedule works as follows. First, temporal constraints are extracted from input texts in addition to query frames (Sect. 5). Second, query frames are roughly scheduled based on the temporal constraints (Sect. 7.1). Note that at this point, only process orders of query frames are determined. Finally, by searching for a motion frame that matches each query frame in order of process, the execution timing of each motion is determined. When automatic motions are required to be executed before a motion, they are added incrementally (Sect. 7.2). By repeating this process for all query frames, the motion clips and their execution timings are determined.

### 7.1 Scheduling query frames

Based on temporal constraints, the query frames are scheduled roughly at first. After that, the process order of all query frames (verbs) is determined. For motions that have a synchronized execution constraint, their process orders are temporarily set as one of them being processed first. The exact timings of all query frames are determined in the process order.

For each query frame, a motion clip is searched from the database as explained in Sect. 6.1. Before searching each motion, the scene condition is set to the time when the motion is executed because the selected motion may change depending on the position of the character or object that the motion involves. The execution timing of the motion is determined based on the duration of the selected motion. The next motion is started just after the previous motion is finished if they have a serial execution constraint. If they have a synchronized executing constraint, their execution timings are determined based on the contact timings of the selected motions.

This process is repeated from the first motion to the last. When multiple query frames are synchronized based on the

temporal constraints, the motions for all query frames are searched and their execution timings are delayed until all constraints are satisfied.

## 7.2 Automatic motions

During the motion scheduling and motion search, a searched motion can sometimes not be executed. In that case, automatic motions are generated and added before the searched motion. As explained earlier, the purpose of our system is to reuse motion data without complex motion planning which may require additional programming for each motion. Therefore, our system deals with minimum automatic motions. The additional motions are also selected from the database. Therefore, each character is easily customized by adding specific kinds of motion to the database without adding any rules or modules.

### 7.2.1 Locomotive motion

If a motion includes interaction with another character or an object in the scene (i.e., a query frame contains a target object or character), the character has to be in the right place to make contact with the object or character. If not, the system automatically adds locomotive motions for the character to move to the right place and to face the right direction.

If the motion frame has a contact position and target direction (e.g., “sitting on a chair” motion should be executed in the right position and direction to the chair), an appropriate locomotive motion is generated so that the character approaches the right point and turns in the right direction. The method for generating locomotive motions explained in Sect. 6.3 is used. If the motion frame has a contact position, an appropriate locomotive motion is generated so that the character approaches the right point. If the motion frame only has a direction (e.g., “shooting toward the target” motion), the character merely turns without walking. As explained in Sect. 6.3, our current system has no path planning; the character merely moves in a straight line to the target position.

### 7.2.2 Taking an instrument

When a character uses an instrument in a motion (i.e., a query frame contains an instrument and the character does not hold it), the character must pick up the instrument object before they use it. When a motion to take the instrument is not explicit in the input text, a ‘take’ motion is selected from the database. When the character is away from the instrument, locomotive motions are also added before the taking motion.

### 7.2.3 Changing posture

For motion searches, if there is no candidate motion whose initial posture matches the terminal posture of the previous motion (i.e., the initial posture of a query frame does not match any of the candidate motion frames), a changing posture motion such as standing up is added. In this case, all motions that include a state change will be candidate motions.

### 7.2.4 Cooperative motion

As explained in Sect. 4.3, when a motion involves interaction with another character, a cooperative motion of the other character follows. When a selected motion frame has cooperative motions and any of them are not indicated in the input text, the default cooperative motion and a temporal constraint of the motion frame are automatically added.

## 8 Experiment and discussion

We have implemented our method and motion database. Currently, the system has six characters as shown in Fig. 5 and about 50 motions that are collected from a commercially available motion capture library. We have tested our system with some short sentences and found that an appropriate motion was selected from each sentence even though the same verb is used in different sentences. An example of the generated animation is available from the author’s web site (<http://www.cg.ces.kyutech.ac.jp/research/modb/index.html>).

To evaluate our framework, we tested it with a published movie script (The Matrix, 1999). Because our motion database does not yet have enough data, we checked whether our methods could handle the descriptions in the movie script and output appropriate query frames. There were about 830 actions (verbs) in the script. We found that about 87% of these were processed by our system without any problems. However, 4% were complex expressions that are difficult to handle using simple rules, such as a sentence with the subject being a character’s body (e.g., “His elbow hits the enemy.”, “His body jumps.”), vague representation (e.g., “he stares into the darkness”), indirect expression (e.g., “He has no answer.”), and ambiguous nouns or pronouns. 4% were verbs that cannot be represented by a motion including non-action or state verbs explained in Sect. 4.2, such as “He feels that ~”, a verb representing a result of a motion such as “miss” in “he shoots her and misses”. 5% were verbs representing initial states in the scene but not actions (e.g., “they are dead”, “he stands in the room”). As discussed later, a non-text-based interface is suitable for specifying initial states or positions of locomotion. According to the above results, 9% of the verbs in the sample script were actually



verbs that cannot be represented as motions. This shows that 95% of verbs that can be represented as motion were handled by our methods. Although it is possible to extend our semantic analysis to support more complex expressions, this will require a great deal of knowledge and rules, which is contrary to the aims of this research. If a complex expression cannot be handled by the system, the user should rephrase it as a plain expression rather than adding more knowledge and rules to the system.

However, even if an animation is generated from a given text, since there is limited information in an input text, a user may not be satisfied with the motions that are found in the database. Moreover, since we use a simple method for generating locomotive motion, motion modification, and motion synthesis, a user may not be satisfied with the synthesized animation. To evaluate the effectiveness of our system, we intend conducting a practical user study in a future work. Improvement of motion generation and the external motion synthesis system [20] is also a future work.

Our current system cannot handle object motions. However, as they are also important for animation, it is easily possible to extend our system to handle them, as they tend to be simpler than human motions.

The fundamental principle of our framework is to make use of motion data without requiring any additional motion specific rules. Currently, our system does not support high-level motion planning such as automatically dividing complex motion into small motions or path planning with object avoidance. Because we use simple rules for automatic locomotion, the resulting animations are not so natural. This can be solved by adding more motion data and some sophisticated modules that generate new motion from a number of motion data sources such as [17, 18].

Our system supposes that scene information, such as the positions of objects and characters, is provided by the user. The existing text-to-scene system [10] can be integrated with our system. However, specifying the positions using natural language can be harder than using a conventional mouse-based interface. So can specifying locomotion path. From a practical viewpoint, a hybrid of a text-based interface and a conventional interface might be more useful.

With our current system, if the user is not satisfied with or wants to change an output motion, they must change the input text and they cannot change the output motions directly. To address this, we are going to develop a natural language-based motion editing interface with which a user can change generated motions interactively by giving instructions to agents, as real directors do with actors.

## 9 Conclusion

We have proposed an animation system that generates animation from natural language text such as movie scripts or

stories. Our future work includes the expansion of both the system and the motion database. Currently, creating animations is very difficult, especially for nonprofessionals. We believe that our system will alleviate this and provide many creators with a means of expressing their stories as animation.

## References

- Oshita, M.: Generating animation from natural language texts and framework of motion database. In: Proc. of International Conference on Cyberworlds 2009, pp. 146–153, Bradford, UK (2009)
- Winograd, T.: Understanding Natural Language. Academic Press, San Diego (1972)
- Badler, N., Bindiganavale, R., Allbeck, J., Schuler, W., Zhao, L., Palmer, M.: Parameterized action representation for virtual human agents. In: Embodied Conversational Agents, pp. 256–284 (2000)
- Bindiganavale, R., Schuler, W., Allbeck, J., Badler, N., Joshi, A., Palmer, M.: Dynamically altering agent behaviors using natural language instructions. In: Proc. of Autonomous Agents 2000, pp. 293–300 (2000)
- Tokunaga, T., Funakoshi, K., Tanaka, H.: K2: animated agents that understand speech commands and perform actions. In: Proc. of 8th Pacific Rim International Conference on Artificial Intelligence 2004, pp. 635–643 (2004)
- Fillmore, C.J.: The case for case. In: Universals in Linguistic Theory, pp. 1–88 (1968)
- Lu, R., Zhan, S.: Automatic Generation of Computer Animation: Using AI for Movie Animation. Springer, Berlin (2002)
- Sumi, K., Nagata, M.: Animated storytelling system via text. In: Proc. of International Conference on Advances in Computer Entertainment Technology (2006)
- Baba, H., Noma, T., Okada, N.: Visualization of temporal and spatial information in natural language descriptions. Trans. Inf. Syst. **E79-D**(5), 591–599 (1996)
- Coyne, B., Sproat, R.: Wordseye: an automatic text-to-scene conversion system. In: Proc. of SIGGRAPH 2001, pp. 487–496 (2000)
- Levine, S., Theobalt, C., Koltun, V.: Real-time prosody-driven synthesis of body language. ACM Trans. Graph. **28**(5), 172 (2009) (In: Proc. of ACM SIGGRAPH Asia 2009)
- Perlin, K., Goldberg, A.: Improv: a system for scripting interactive actors in virtual worlds. In: Proc. of SIGGRAPH '96 Proceedings, pp. 205–216 (1996)
- Conway, M.J.: Alice: easy-to-learn 3D scripting for novices. PhD Dissertation, University of Virginia (1997)
- Hayashi, M., Ueda, H., Kurihara, T., Yasumura, M.: TVML (TV program Making Language)—automatic TV program generation from text-based script. In: Proc. of Imagina '99, pp. 84–89 (1999)
- Shim, H., Kang, B.G.: CAMEO—camera, audio and motion with emotion orchestration for immersive cinematography. In: Proc. of International Conference on Advances in Computer Entertainment Technology (ACE) 2008, pp. 115–118 (2008)
- Park, S.I., Shin, H.J., Shin, S.Y.: On-line locomotion generation based on motion blending. In: Proc. of ACM SIGGRAPH Symposium on Computer Animation 2002, pp. 105–111 (2002)
- Rose, C., Cohen, M.F., Bodenheimer, B.: Verbs and adverbs: Multidimensional motion interpolation. IEEE Comput. Graph. Appl. **18**(5), 32–40 (1998)
- Kovar, L., Gleicher, M.: Automated extraction and parameterization of motions in large data sets. ACM Trans. Graph. **23**(3), 559–568 (2004)



19. Klein, D., Manning, C.D.: Fast exact inference with a factored model for natural language parsing. In: *Advances in Neural Information Processing Systems 15 (NIPS 2002)*, pp. 3–10 (2003)
20. Oshita, M.: Smart motion synthesis. *Comput. Graph. Forum* **27**(7), 1909–1918 (2008)



**Masaki Oshita** received his BS, MS, and PhD degrees from Kyushu University in 1998, 2000, and 2003, respectively. Since 2003, he has been an associate professor in the Department of Systems Innovation and Informatics at Kyushu Institute of Technology. His research interest is on-line computer animation which includes interactive motion control of human figures, motion control interface, motion synthesis, physics-based simulation, human–computer interaction, computer vision, etc.