

Emotional Variation in Speech-Based Natural Language Generation

Michael Fleischman and Eduard Hovy

USC Information Science Institute
4676 Admiralty Way
Marina del Rey, CA 90292-6695
U.S.A.{fleisch, hovy} @ISI.edu

02/14/02

Session: Main

Topic Area: Speech-Based Natural Language Generation

Key Words: Emotion, Natural Language Generation, Speech-Based Generation, Multi-modal Virtual Environments, Empathy, Personality, Embodied Agents, MRE

Word Count: 3,196

Under consideration for other conferences (specify)? No

Abstract

We present a framework for handling emotional variations in a speech-based natural language system for use in the MRE virtual training environment. The system is a first step toward addressing issues in emotion-based modeling of verbal communicative behavior. We cast the problem of emotional generation as a distance minimization task, in which the system chooses between multiple valid realizations for a given input based on the emotional distance of each realization from the speaker's attitude toward that input.

Emotional Variation in Speech-Based Natural Language Generation

Michael Fleischman and Eduard Hovy

USC Information Science Institute
4676 Admiralty Way
Marina del Rey, CA 90292-6695
U.S.A.
{fleisch, hovy} @ISI.edu

Abstract

We present a framework for handling emotional variations in a speech-based natural language system for use in the MRE virtual training environment. The system is a first step toward addressing issues in emotion-based modeling of verbal communicative behavior. We cast the problem of emotional generation as a distance minimization task, in which the system chooses between multiple valid realizations for a given input based on the emotional distance of each realization from the speaker's attitude toward that input. We discuss evaluations of the system and future work that includes modeling personality and empathy within the same framework.

1. Introduction

Emotion is an ever-present characteristic of human experience and behavior. As fundamental to the human condition as cognition, emotion has begun to pique the interest of those researchers in the Artificial Intelligence community concerned with simulating human behavior in embodied agents. Nowhere is this interest more prominent than in the domain of multi-modal, virtual training environments. In such environments, realistic modeling of emotion

enhances the user's ability to suspend disbelief (Marsella & Gratch, 2001), and can be used as an additional parameter in creating more variable training scenarios.

While much attention has been paid to the effect of emotion on planning and non-verbal behavior (Marsella et al., 2001), little work has been done on the effects of emotion on the verbal behavior of embodied agents. Our research represents a first step toward developing an integrated framework for modeling emotion in the speech-based natural language generation of embodied agents.

2. Mission Rehearsal Exercise

The emotional NLG system that we present is designed within the Mission Rehearsal Exercise (MRE) virtual training environment (Swartout et al. 2001). The MRE is a large-scale collaborative research effort to develop a fully interactive training simulation modeled after the holodeck in Star Trek. The project brings together researchers working on graphics, 3-D audio, artificial intelligence, and Hollywood screenwriters to create a realistic virtual world in which human subjects can interact naturally with simulated agents. The agents are modeled using the Steve system of Rickel and Johnson (1999). They communicate through voice and gesture, reason about tasks and actions, and incorporate a complex model of their own emotions, as well as the emotional states of

the other agents in their environment (Gratch & Marsella, 2001). Users can query and interact with one (and eventually many) agent in real-time as they proceed through a scenario developed for the particular training mission at hand.

The scenario presently implemented is designed to train army lieutenants for eastern European peace keeping missions. The scenario centers around the trainee, a human lieutenant, who is attempting to move his platoon to a support position, when one of his drivers unexpectedly collides with a civilian car. A civilian passenger, a young boy, is critically injured and the boy's mother, as well as a crowd of local onlookers, is becoming increasingly agitated. The trainee must interact with his or her virtual platoon sergeant in order to stabilize the situation.

The MRE represents the integration of many fields in NLP. As the trainee interacts with the virtual agents in the environment, automatic speech recognition translates the user's speech into a text string that is passed to the natural language understanding module. This module uses a finite state machine to convert the string into a case frame structure that is passed to a dialogue manager. At this point, the dialogue manager interacts with the task planner, the action selector, and the emotion model to initiate a particular response. The content of this response is then passed as an impoverished case frame to the NLG system. Generation converts the input into a tree structure that contains both syntactic and semantic information. The tree is then passed to a gesture module and is tagged with non-verbal information to control gaze and body movements. Finally, the tree is flattened, the gestures and visemes are synched using the BEAT system (Cassell, 2001), and the speech is synthesized.

3. NLG in MRE

Generation in the MRE is a hybrid process. The generator can take as input both highly elaborated case frames, for scenario specific utterances, and more impoverished

frames, for use in interactive conversation. We discuss only the conversational aspect of the system.

The generator is, at this point, highly domain dependent, but has sufficient coverage to generate utterances for every task in the agents' task models. The generator is implemented in the SOAR programming language (Newell, 1990) and takes place in three stages: sentence planning, realization and ranking.

3.1 Sentence Planning

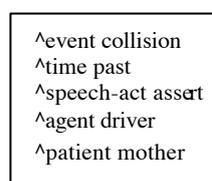


Figure 1a. Input from dialogue manager: input to sentence planning phase of generation

As seen in Figure 1(a), the inputs to this stage are received from the dialogue manager. These inputs contain minimal information about the state or event to be described, along with references to the actors and objects involved. A set of SOAR production rules converts this information into an enriched case frame structure, seen in Figure 1(b), which contains more detailed information about the events and objects in the input. The conversion process, which involves choosing the appropriate object case frames, relies heavily on the emotional decision engine.

3.2 Realization

Realization is a highly lexicalized procedure, and tree construction begins with the selection of main verbs (more on this below). Each verb in the lexicon carries with it slots for its constituents (e.g., agent, patient), which form branches in the tree. Once the verb is chosen, production rules recursively expand the nodes in the tree until no more nodes can be expanded. As each production rule fires, the relevant portion of the semantic frame is propagated down into

the expanded nodes. Thus, every node in the tree contains a pointer to the specific aspect of the semantic frame from which it was created. For example, in Figure 1(c), the NP node of “the mother” contains in it a pointer to the frame <patient> from Figure 1(b). By keeping semantic content localized in the tree, we allow the gesture and speech synthesis modules convenient access to needed semantic information. This strategy is particularly convenient in a setting such as the MRE, where modules require increasing amounts of information as research continues.

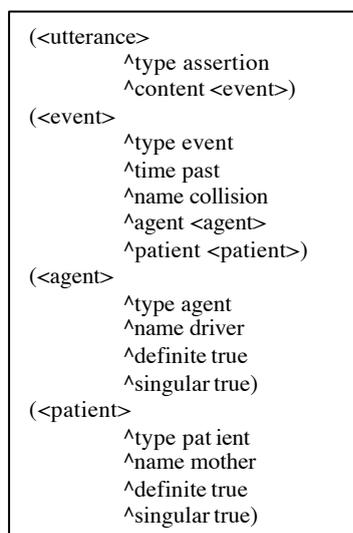


Figure 1b. Expansion of input from dialogue manager; output of sentence planning

For any given state and event, there are a number of theoretically valid realizations available in the lexicon. Instead of attempting to decide which is most appropriate at any stage, we adopt a strategy similar to that introduced by (Knight & Hatzivassiloglou, 1995), which puts off the decision until realization is complete. We realize all possible valid trees that correspond to a given semantic input, and store the fully constructed trees in a forest structure. After all such trees are constructed we move on to the final stage.

3.3 Ranking

In this stage we examine all the trees in the forest structure and decide which tree will be propagated further down the NLP pipeline. Each tree is given a rank score based upon the tree’s information content and emotional quality. The score of each tree is calculated by recursively summing the scores of the nodes along the frontiers of the tree, and then percolating that sum up to the next layer. Summing and percolating proceeds until the root node is given a score that is equivalent to the sum of the scores for the individual nodes of that tree. The tree with the highest root node score is selected.

4. Emotional Variations

We cast the problem of emotional language generation as an optimization problem in which multiple acceptable realizations of a given semantic frame are produced. Given a set of valid realizations for a given frame, we output the sentence that most closely fits the emotional state of the speaker.

4.1 Speaker’s Emotions

In the current (somewhat simplified scheme) we represent the emotional state of the speaker as a set of integer values (ranging from -5 to 5). Each value corresponds to the speaker’s emotional attitude toward a specific element of the input. For example, Figure 2(a) depicts an input describing an event (collision) with an agent (driver) and a patient (mother). Each element is further described by an emotional attitude representing how the speaker feels about each of the concepts (agent: +4, patient: +1, event: -1). These values are calculated by the emotion model and passed as input to the generator, along with the semantic input, by the dialogue manager.

4.2 Emotional Distance

We calculate the fit of a sentence to the emotional state of the speaker as the distance between the speaker’s emotional

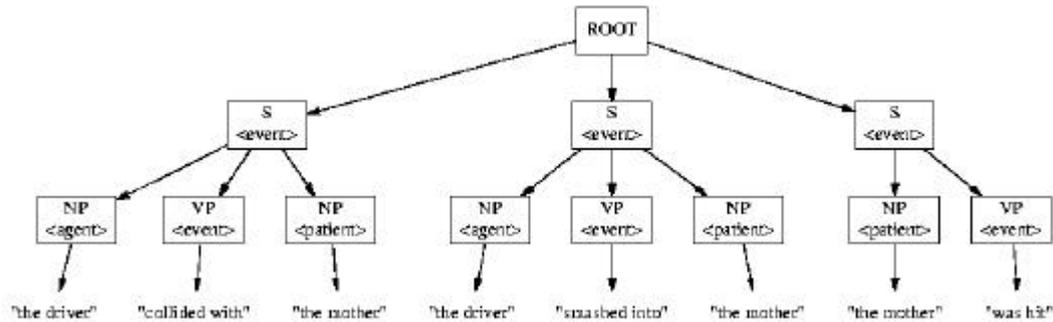


Figure 1c. A subset of the forest output of realization.

attitude toward an object and the default emotional shade of the lexical item or expression used to express that object. While the emotional attitudes of the speaker are given by the emotion module, the emotional shades for the lexical items are stored in the lexicon, as real valued scores.

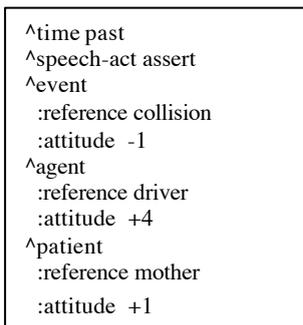


Figure 2a. Input from dialogue manager showing speaker's attitudes toward objects and events

Deciding what default value shade each lexical item is given is, at this point, a matter of linguistic intuition. However, empirical alternatives are discussed in later sections.

In order to avoid the memory explosion that comes with calculating distances for every possible valid sentence that represents a frame, we divide the task between two stages of generation: planning and ranking.

During planning, the impoverished input given by the dialogue manager is expanded into a semantic frame ready for

realization. The task of expansion involves deciding which frame is to be chosen to represent each object in the input. For example, Figure 2(b), shows a number of

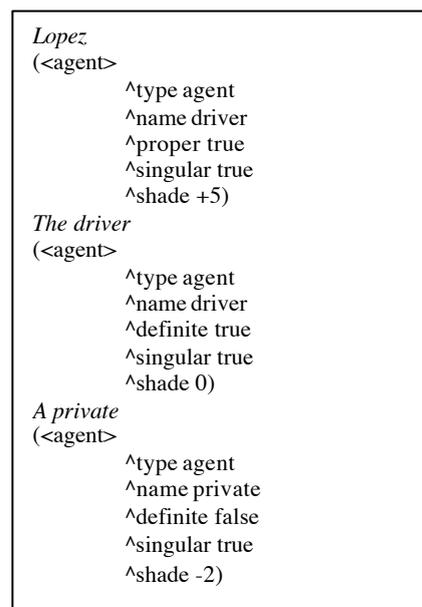


Figure 2b. Subset of possible case frame expansions for object "driver"; ^shade value represents emotional shade of using that frame to refer to "driver".

possible frames that could be used to represent the agent "driver." The decision is based on the emotional shade of each semantic option. A distance is calculated between the shade of each semantic frame that could represent the "driver" and the emotional attitude of the speaker toward the "driver."

The frame with the minimum distance is chosen for expansion. This is done for each of the objects associated with the event or state. Once all objects have been assigned a frame, planning is complete, and realization begins.

During realization, semantic frames are expanded as described in section 3.2. In this phase, all verbs in the lexicon that are valid representations of the input frame are used to create distinct trees. Each verb carries with it its emotional shade. This shade is expressed in two ways: by the overall emotional connotation of the verb itself, and by the emotional connotations that the verb imparts on its constituents. A sample of the lexicon for verbs that describe the event “collision” is shown in Figure 3.

Gloss	Agent Shade	Event Shade	Patient Shade
<i>A bumped into P</i>	0	-1	0
<i>A collided with P</i>	-1	-1	0
<i>A ran into P</i>	-2	-2	0
<i>A hit P</i>	-3	-2	+1
<i>P was hit by A</i>	-2	-2	+1
<i>P was hit</i>	Na	-2	+1
<i>A crashed into P</i>	-3	-3	+1
<i>A smashed into P</i>	-5	-4	+2
<i>There was an accident</i>	Na	-2	Na

Figure 3. Valid lexical representations for event “collision” including the shades that the verbs apply to the objects related to the event.

As seen in the entry, the verb “hit” casts a more negative emotional shade on the agent and event than the verb “smash.” However, “smash” casts a more positive shade on the patient of the event than “hit.” This effect is seen in the realizations: “The driver hit the mother’s car” and “The driver smashed into the mother’s car.” While both verbs betray negativity toward the “driver,” the latter is far more severe than the former. Further, because of the intensity of the verb (and the negativity of the event) the patient is cast as more sympathetic in the latter sentence. These sorts of differences are what the lexical entries capture through the use of different valued shades.

When the ranking phase begins, each tree formed of these verbs is ranked and

compared, as described in section 3.3. The tree finally selected is that in which the *total* emotional distance from the speaker’s attitude is minimized across the event itself, as well as across all the constituents of that event. Thus, even if the speaker feels very negatively toward the event described in Figure 2(a), because the distances for each tree are summed across all of its constituents, the generator may still opt not to use the strong lexical item “smash” if the speaker has intensely positive feelings toward the agent.

Table 1 shows example calculations for three variations of the input given in Figure 2(a). The emotion scores for each variation are computed using the formula below, where $attitude(x)$ is the speaker’s attitude toward x and $shade(x)$ is the shade of the lexical item used to represent x .

$$EmotScore(x) = Dist(verb) + \sum_i Dist(constituent_i)$$

$$Dist(x) = |attitude(x) - shade(x)|$$

Verb	D(agent)	D(verb)	D(patient)	Score
was hit	Na	 -2--1 =1	 1-1 =0	1
collided	1-4 =5	1-1 =0	0-1 =1	6
smashed	1-5-4 =9	1-4-1 =3	2-1 =1	13

Table 1. Emotional scores for input shown in Figure 2(a). “was hit” obtains minimal distance score, and is selected.

4.3 Discussion

This method of calculating emotional effect provides a great deal of variation with very little overhead. Once the lexicon is updated with items that carry emotional shadings, it is simply a matter of assigning the speaker attitudes, and applying a simple distance metric. The system will then automatically decide between possible realizations based not only on lexical choice, but also on sentence structure.

As seen in Figure 3, the passive construction of the verb “hit” shades the elements of the event differently than the active construction of the same verb. Because the agent is not realized at all, the passive will be preferred when the attitude of the speaker is very positive toward the agent. This is because the event itself is such that it always

Event	Agent	Patient	Output
-2	-3	5	<i>A private crashed into the mother</i>
-2	-3	4	<i>A private hit the mother</i>
-2	-3	3	<i>A private hit the mother</i>
-2	-3	2	<i>The mother was hit by a private</i>
-2	-3	1	<i>The mother was hit by a private</i>
-2	-3	0	<i>A private ran into a woman</i>
-2	-3	-1	<i>A private ran into a civilian</i>
-2	-3	-2	<i>A private ran into a civilian</i>
-2	-3	-3	<i>A private ran into a civilian</i>
-2	-3	-4	<i>A private ran into a civilian</i>
-2	-3	-5	<i>A private ran into one of our "responsibilities"</i>

Figure 4a. Effect of varying the speaker's attitude toward the patient of an event; attitude toward the agent and the event itself are held constant.

shades the agent negatively. Thus, by not mentioning the agent at all, the speaker avoids having to say something negative about an object it regards positively. In extreme cases, the speaker's attitude may even lead him or her to elide most of the sentence or to not speak at all.

In practice, however, the need to convey information must be taken into account as well. We therefore compute a total rank score as a linear combination of the emotional distance and information content expressed by the tree:

$$\text{Total Score}(x) = a \text{Info}(x) - (1-a)\text{EmotScore}(x)$$

Here, the $\text{Info}(x)$ is the number of slots from the input frame that are realized by x , and $\text{EmotScore}(x)$ is as above. By changing the coefficient a , different weight will be given to the information content of the utterance versus its emotional shade. One can view an aspect of the personality of the speaker as a tendency toward a certain value for a : An agent who is more interested in the facts will always use a high a , while one who is more concerned with expressing emotion will use a low value.

5. Evaluation

In evaluating this system, we were particularly concerned with two points. First, how sensitive is the system to different inputs,

and second, how much do the outputs actually mimic the emotional behavior of humans.

To determine the sensitivity of the system to different inputs we cycled through the parameters of the input space and observed the frequency of change in the output sentences. Because of the large number of possible inputs even for a simple frame such as in Figure 2(a) (i.e., the number of possible values raised to the power of the number of objects), we present results only for a subset of examples. Figure 4(a) shows the outputs of the generator when the attitude toward the patient is changed, holding all else constant; and Figure 4(b) shows the output when the attitude toward the agent is changed, holding all else constant. (Notice that the realization of the object being held constant does not change. This is because the frames that dictate the realizations are chosen at the sentence planning stage.)

It is interesting to notice the difference in sensitivity between the two cases; changing the attitude toward the agent has more effect than changing the attitude toward the patient. This is because of the nature of the event "collision." As can be seen in Figure 3, the different realizations of the event vary mostly in their effect on the agent of the sentence. Thus, changing the attitude toward the patient has an effect on the sentence only at the extremes of the range of attitudes. We conclude that using a distance measure as the basis for the emotion calculus is adequately sensitive.

Event	Agent	Patient	Output
0	5	3	<i>A woman was hit</i>
0	4	3	<i>A woman was hit</i>
0	3	3	<i>Our driver bumped into a woman</i>
0	2	3	<i>Our driver bumped into a woman</i>
0	1	3	<i>One of our drivers bumped into a woman</i>
0	0	3	<i>The driver bumped into a woman</i>
0	-1	3	<i>The driver collided with a woman</i>
0	-2	3	<i>A driver ran into a woman</i>
0	-3	3	<i>One of our privates collided with a woman</i>
0	-4	3	<i>One of our privates ran into a woman</i>
0	-5	3	<i>A damn private collided with a woman</i>

Figure 4b. Effect of varying the speaker's attitude toward the agent of an event; attitude toward the patient and the event itself are held constant.

Evaluating how a generated output correlates with human intuitions regarding the speaker's attitude is not an easy task. Judging the emotional state of someone based solely on their utterances is near impossible and presents many methodological challenges.

One way of determining such correlation is by having humans guess the attitudes of the system and comparing this to the system's emotional input. We ask subjects to rate the objects in the sentence on scales from 5 to -5 (where 5 means the speaker thinks most favorably about the object and -5 is most unfavorably). The correlation between what the subjects believed to be the attitudes of the speaker and the actual attitudes used for generation was significant even with a low power ($r=0.659$, $n=10$) indicating that the expressiveness of the system is reliable.

An interesting prospect for future work is incorporating the procedure for evaluating the system into the system's actual design. We plan to examine the feasibility of using averages of human judgments as the shades for verbs in the lexicon. This is essentially the method that is employed now (using only the authors' intuitions), but by increasing the size of the sample, we suspect even more reliable outputs can be found.

6. Related Work

We are unaware of much work on emotional variation in language production that focuses on utterances (as opposed to

intonation and non-verbal communication). The most relevant work is over 10 years old. In his thesis, Hovy (1988) implemented a 3-valued (positive, negative, neutral) system of emotional shades with a simple sign multiplication calculus to control affect laden text generation. Our distance calculus adds flexibility and also allows us to extend the emotional input from simple like/dislike to more complicated constructs.

Work by Bateman and Paris (1989), and Paris (1988) focused on variations of expert system output based on the reader's knowledge. Also here, the rules for combining ratings of sentence constituents was fairly simple and not easily extendable.

7. Discussion/Future Work

This work is a first step toward incorporating emotionality into natural language generation. The system we present, while not complete, shows much promise for future work. The notion of a speaker's attitude toward an object or event, for example, while very simple in this implementation, can easily be expanded to fit the needs of the system. Because the decision method is a simple Euclidean distance metric, the single valued attitudes that we describe can easily be converted to more complex vectors of emotions. Thus, in future implementations, it may not be satisfactory to feel -3 about the "driver," but rather, the speaker may feel toward the driver: (1:respect); (-2:blame); (-4:envy); etc. Once

the lexicon is updated for the richer format, the distance metric need only be changed to operate on vectors.

Of further interest is the possibility of incorporating empathy into the generation process. In the current system, generation is based only on the emotions of the speaker. In the future, with more information from the emotion model, we will be able to generate sentences also taking into account the emotional attitudes of the hearers.

For example, in the MRE, when the agent is asked about the status of the injured boy, it knows that the boy is critically injured and that this information will upset the boy's mother. Empathic generation takes this into account when generating a response, saying, instead of: "the boy is dying," the more appropriate: "the boy needs a doctor."

Such empathy is easily implemented in our framework by replacing the vector of the speaker's attitudes with a linear combination of the speaker's attitudes and the hearers' attitudes. Under this formulation, the personality of the speaker can be partially described in terms of the weights with which one performs this combination (much like the a used in ranking, see section 4.3). For example, a speaker who is *sensitive* is just someone who tends to give higher weight to the attitudes of the hearers than to their own. On the other hand, an *indifferent* speaker would be one who ignores the attitudes of the hearers when generating an utterance.

We believe that the framework that we have set up is a simple and convenient method for treating emotion in language. As virtual environments become more common, and the population of virtual characters in those environments explodes, the need for such emotional generation becomes more apparent. While our system is not complete, it is a simple and intuitive method for dealing with a necessary and ignored area of natural language generation.

8. References

1. Bateman, J. and Paris, C.L. (1989). Phrasing a text in terms the user can understand. In *Proceedings of the Eleventh International Joint*

Conference on Artificial Intelligence, Detroit, Michigan.

2. Cassell, J., Vilhjlmsson, H., & Bickmore, T. (2001). *BEAT: The Behavior Expression Animation Toolkit*, Proc. of SIGGRAPH, ACM Press.

3. Gratch, J. and Marsella, S. (2001). Tears and Fears: Modeling emotions and emotional behaviors in synthetic agents, in *Proc. of the 5th International Conference on Autonomous Agents*, Montreal, Canada.

4. Hovy, E. H. (1988). *Generating Natural Language under Pragmatic Constraints*. Lawrence Erlbaum, Hillsdale, New Jersey.

5. Knight, K. & Hatzivassiloglou, V. (1995). *Two-level, many-paths generation*. In Proceedings of the 33rd Annual Meeting of the Association for Computational Linguistics, Cambridge, Mass, June 1995, pp. 252-260.

6. Marsella, S. and Gratch, J. (2001) Modeling the Interplay of Emotions and Plans in Multi-Agent Simulations, in *Proc. of the 23rd Annual Conference of the Cognitive Science Society*, Edinburgh, Scotland.

7. Marsella, S., Gratch, J., and Rickel, J. (2001). The Effect of Affect: Modeling the Impact of Emotional State on the Behavior of Interactive Virtual Humans, in *Proc. of the Agents2001 Workshop on Representing, Annotating, and Evaluating Non-Verbal and Verbal Communicative Acts to Achieve Contextual Embodied Agents*, Montreal, Canada.

8. Newell, A. (1990). *Unified Theories of Cognition*. Harvard University Press.

9. Paris, C.L. (1988). Tailoring Object Descriptions to a User's Level of Expertise. *Computational Linguistics* 14(3):64-78.

10. Rickel, J., & Johnson, W. (1999). Virtual Humans for Team Training in Virtual Reality. In *Proceedings of the Ninth International Conference on AI in Education*, pp. 578-585. IOS Press.

11. Swartout, W., Hill, R., Gratch, J., Johnson, W.L., Kyriakakis, C., Labore, K., Lindheim, R., Marsella, D., Moore, B., Morie, J., Rickel, J., Thiebaut, M., Tuch, L., Whitney, R. (2001). Towards the Holodeck: Integrating Graphics, Sound, Character and Story. In *Proc. of the Fifth International Conference on Autonomous Agents*, Montreal, Canada.

