

A Model for Multimodal Reference Resolution

Luis Pineda*
National Autonomous University of
Mexico (UNAM)

Gabriela Garza

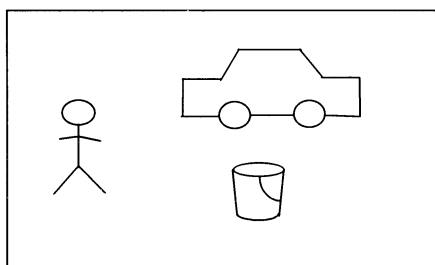
An important aspect of the interpretation of multimodal messages is the ability to identify when the same object in the world is the referent of symbols in different modalities. To understand the caption of a picture, for instance, one needs to identify the graphical symbols that are referred to by names and pronouns in the natural language text. One way to think of this problem is in terms of the notion of anaphora; however, unlike linguistic anaphoric inference, in which antecedents for pronouns are selected from a linguistic context, in the interpretation of the textual part of multimodal messages the antecedents are selected from a graphical context. Under this view, resolving multimodal references is like resolving anaphora across modalities. Another way to see the same problem is to look at pronouns in texts about drawings as deictic. In this second view, the context of interpretation of a natural language term is defined as a set of expressions of a graphical language with well-defined syntax and semantics. Natural language and graphical terms are thought of as standing in a relation of translation similar to the translation relation that holds between natural languages. In this paper a theory based on this second view is presented. In this theory, the relations between multimodal representation and spatial deixis, on the one hand, and multimodal reasoning and deictic inference, on the other, are discussed. An integrated model of anaphoric and deictic resolution in the context of the interpretation of multimodal discourse is also advanced.

1. Reference, Spatial Deixis, and Modality

In this paper a model for the resolution of multimodal references is presented. This is the problem of finding the referent of a symbol in one modality using information present either in the same or in other modalities. A model of this kind can be useful both for implementing intelligent multimodal tools (e.g., authoring tools to input natural language and graphics interactively for the automatic construction of tutorials or manuals) and from the point of view of human-computer interaction (HCI) where it can help in the design of computer interfaces in which the interpretation constraints of multimodal messages should be taken into account.

Consider Figure 1 (adapted from Rist [1996]) in which a message is expressed through two different modalities, namely text and graphics. The figure illustrates a kind of reasoning required to understand multimodal presentations: in order to make sense of the message, the interpreter must realize what individuals are referred to by the pronouns *he* and *it* in the text. For the sake of argument, it is assumed that the graphical symbols in the figure are understood directly in terms of a graphical lexicon, in the same way that the words *he*, *it*, and *washed* are understood in terms of the textual

* Department of Computer Science, Institute for Applied Mathematics and Systems (IIMAS), National Autonomous University of Mexico (UNAM), Mexico. E-mail: luis@leibniz.iimas.unam.mx.



"He washed it"

Figure 1

Instance of linguistic anaphor with pictorial antecedent.

"Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt."



Figure 2

Instance of a pictorial anaphor with linguistic antecedent.

lexicon. It can easily be seen that given the graphical context, *he* should resolve to the man, and *it* should resolve to the car. However, this inference is not valid since the information inferred is not contained in the overt graphical context and the meaning of the words involved.

One way to look at this problem is as a case of anaphoric inference. Consider that the information provided by graphical means can also be expressed through the following piece of discourse: *There is a man, a car, and a bucket. He washed it.* With Kamp's discourse representation theory (DRT) (Kamp 1981; Kamp and Reyle 1993) a discourse representation structure (DRS) in which the reference to the pronoun *he* is constrained to be the man can be built. However, the pronoun *it* has two possible antecedents, and conceptual knowledge is required to select the appropriate one. In particular, the knowledge that a man can wash objects with water, and that water is carried in buckets, must be employed. If these concepts are included in the interpretation context like DRT conditions (which should be retrieved from memory rather than from the normal flow of discourse), the anaphora can be solved. By analogy, situations like the one illustrated in Figure 1 have been considered problems of anaphors with pictorial antecedents in which the interpretation context is built not from a preceding text but from a graphical representation that is introduced with the text (André and Rist 1994).

Consider now the converse situation shown in Figure 2 (adapted from Rist [1996]), in which a drawing is interpreted as a map in the context of the preceding text. The dots and lines in the drawing, and their properties, do not have an interpretation and the picture in itself is meaningless. However, given the context introduced by the text, and also considering the common knowledge that Paris is a city in France, and Frankfurt a city in Germany, and that Germany lies to the east of France (to the right),

it is possible to infer that the denotations of the dots to the left, middle, and right in the picture are Paris, Saarbrücken, and Frankfurt, respectively, and that the dotted lines denote borders between countries, and in particular, the lower segment denotes the border between France and Germany. In this example, graphical symbols can be thought of as “variables” of the graphical representation or “graphical pronouns” that can be resolved in terms of the textual antecedent. Here again, the inference is not valid, as the graphical symbols could be given other interpretations or none at all.

The situation in Figure 2 has been characterized as an instance of a pictorial anaphor with linguistic antecedent, and further related examples can be found in André and Rist (1994). This situation, however, cannot be modeled very easily in terms of Kamp’s DRT because the “pronouns” are not linguistic objects, and lacking a proper formalization of the graphical information, there is no straightforward way to express in a discourse representation structure that a dot representing “a variable” in the graphical domain has the same denotation as a natural language name or description introduced from text in a DRS. Furthermore, the situation in Figure 1 can be thought of as anaphoric only if we ignore the modality of the graphics, as was done above; but if the notion of modality is to be considered at all in the analysis, then the situation in Figure 1 poses the same kinds of problems as the one in Figure 2. In general, graphical objects, functioning as constant terms or as variables, introduced as antecedents or as pronouns, cannot be expressed in a DRS, since the rules constructing these structures are triggered by specific syntactic configurations of the natural language in which the information is expressed. However, this limitation can be overcome if graphical information can be expressed in a language with well-defined syntax and semantics.

An alternative is to look at these kinds of problems in terms of the traditional linguistic notion of deixis (Lyons 1968). Deixis has to do with the orientational features of language, which are relative to the spatio-temporal situation of an utterance. Under this view, and in connection with the notion of graphical anaphora discussed above, it is possible to mention the deictic category of demonstrative pronouns: words like *this* and *that* permit us to make reference to extralinguistic objects. In Figure 1, for instance, the pronouns *he* and *it* can be supported by overt pointing acts at the time the expression *he washed it* is uttered. Note that the purpose of the pointing act is to provide the referents for the pronouns directly, greatly simplifying the resolution process. However, the deictic use of a pronoun does not necessarily have to be supported by a physical gesture, because deictic use is characterized, more generally, by the identification of the referent in a metalinguistic context. Ambiguity in such words is not unusual, as they can also function as anaphors if they are preceded by a linguistic context, and even as determiners with a deictic component (e.g., *this car*). Additionally, not only demonstratives and pronouns but also proper names, definite descriptions, and even indefinites can be used deictically. As a great variety of contextual factors are conceivably involved in the interpretation of a deictic expression, gestures, although prominent, should be thought of only as one particular kind of contextual factor. In summary, the denotation of a deictic term is the individual that is picked out by the human interpreter in relation to the interpretation context.¹ Consider that in the same way that an anaphoric inference is required for identifying the antecedent of an anaphoric term, an inference process is required for interpreting a term used deictically. We refer to this process as a **deictic inference**. The inference by

¹ An operator called DTHAT for mapping deictic terms into their referents in an interpretation context is introduced in Kaplan’s logic of demonstratives (Kaplan 1978).

which one determines that *he* and *it* are the man and the car is, accordingly, a deictic inference.

For our purposes, it is important to investigate the nature of the relation between the notions of deixis and modality, on the one hand, and multimodal reasoning and inference, either deictic or anaphoric, on the other. According to Kamp (1981, 283), the difference between deictic and anaphoric pronouns is that,

deictic and anaphoric pronouns select their referents from certain sets of antecedently available entities. The two pronoun's uses differ with regard to the nature of these sets. In the case of a deictic pronoun the set contains entities that belong to the real world, whereas the selection set for an anaphoric pronoun is made up of constituents of the representation that has been constructed in response to antecedent discourse.

Our concern here is how “the set of entities that belong to the real world” is accessible to the interpreter. In normal deictic spatial situations the referent of a deictic term is perceived directly through the visual modality, and as a result of such a visual interpretation process, the object is represented by the subject. The question is how the information can be expressed in this intermediate “visual” representation. A plausible answer is that there is a coding system and a medium associated with each particular modality. Our suggestion is that the notion of modality is a representational notion, and not a sensory one as normally assumed in psychological discussion. In our sense, a modality is a formal language, with a lexicon and well-defined syntactic and semantic structures, with an associated medium in which the expressions of the modality are written. Multimodal reasoning is a process involving information expressed in the languages associated with different modalities, and is achieved with the help of a translation relation similar to the relation of translation between natural languages. Performing a multimodal reasoning process is possible if the translation relation between expressions of different modalities is available. However, for particular multimodal reasoning tasks, the translation relation between individual constants of different modalities cannot be stated beforehand and has to be worked out dynamically through a deictic inferential process, as will be argued in the rest of this paper.

1.1 A Model for Multimodal Representation

This view of multimodal representation and reasoning can be formalized in terms of Montague's general semiotic program (Dowty, Wall, and Peters 1985). Each modality in the system can be captured through a particular language, and relations between expressions of different modalities can be modeled in terms of translation functions from basic and composite expressions of the source modality into expressions of the target modality. In a system of this kind, interpreting examples in Figures 1 and 2 in relation to the linguistic modality is a matter of interpreting the information expressed through natural language directly when enough information is available, and completing the interpretation process by means of translating expressions of the graphical modality into the linguistic one, and vice versa. Consider Figure 3—developing from previous work (Pineda 1989, 1998; Klein and Pineda 1990; Santana 1999)—in which a multimodal representational system for linguistic and graphical modalities is illustrated.

The circles labeled *L* and *G* in Figure 3 stand for sets of expressions of the natural language (e.g., English) and the graphical language, respectively, and the circle labeled

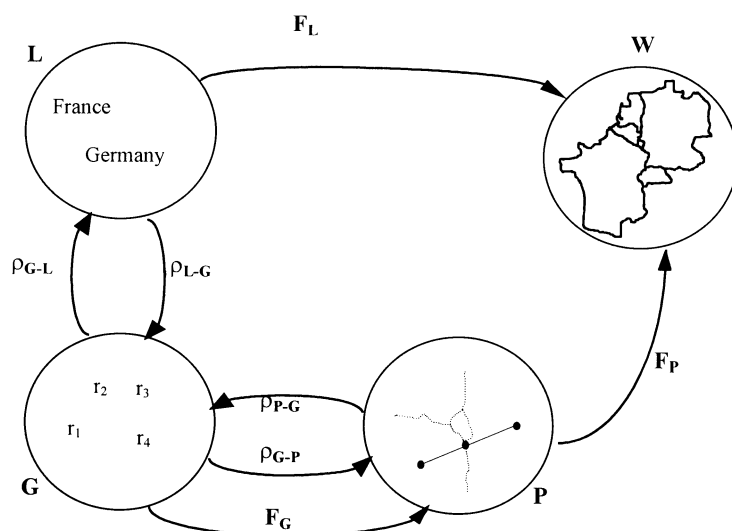


Figure 3
Multimodal representational system for linguistic and graphical modalities.

P stands for the set of graphical symbols constituting the graphical modality proper (i.e., the actual symbols on a piece of paper or on the screen). Note that two sets of expressions are considered for the graphical modality: the expressions in **G** belong to a formal language in which the geometry of pictures is represented and reasoned about, and **P** contains the overt graphical symbols that can be seen and drawn but cannot be manipulated directly. The functions ρ_{L-G} and ρ_{G-L} stand for the translation mappings between the languages **L** and **G**, and the functions ρ_{P-G} and ρ_{G-P} stand for the corresponding translations between **G** and **P**. The translation function ρ_{P-G} maps well-defined objects of the graphical modality into expressions of **G** where the interpretation process is performed. The translation ρ_{G-P} , on the other hand, maps geometrical expressions of **G** into pictures; for every well-defined term of **G** of a graphical type (e.g., *dot*, *line*, etc.) there is a graphical object or a graphical composition that can be drawn or highlighted with the application of geometrical algorithms associated to operators of **G** in a systematic fashion. The circle labeled **W** stands for the world and together with the functions F_L and F_P constitutes a multimodal system of interpretation. The ordered pair $\langle W, F_L \rangle$ defines the model M_L for the natural language, and the ordered pair $\langle W, F_P \rangle$ defines the model M_P for the interpretation of drawings. The interpretation of expressions in **G** in relation to the world is defined either by the composition $F_L \circ \rho_{G-L}$ or, alternatively, by $F_P \circ \rho_{G-P}$. The denotation of the word *France* in **L**, for instance, is the same as the denotation of the corresponding region of the map of Europe that denotes France, the country, since both refer to the same individual. The denotation of the symbol r_1 in **G** that is related to the word *France* in **L** through ρ_{G-L} , and to a particular region in **P** through ρ_{G-P} , is also France, as translation is a meaning-preserving relation between expressions. The interpretation functions F_L and F_P relate basic expressions, either graphical or linguistic, to the objects or relations of the world that these expressions happen to represent, and the definition of a semantic algebra for computing the denotation of composite graphical and linguistic expressions is required.

An important consideration for the scheme in Figure 3 is that the symbols of **P** have two roles: on the one hand, they are representational objects (e.g., a region of

the drawing represents a country), but on the other, they are also geometrical objects that can be talked about as geometrical entities. The geometrical region of the map representing France, for instance, is itself represented by the constant r_1 in \mathbf{G} . In this second view, geometrical entities are individual objects in the world of geometry, and as such they have a number of geometrical properties that are independent of whether we think of graphical symbols as objects in themselves or as symbols representing something else. The same duality can be stated from the point of view of the expressions of \mathbf{G} , since the set of individual geometrical objects (i.e., \mathbf{P}) constitutes a domain of interpretation for the language \mathbf{G} . This is to say that expressions of \mathbf{G} have two interpretations: they represent geometrical objects, properties, and relations directly, but they also represent the objects of the world (e.g., France, Germany, etc.) indirectly through the translation relation and interpretation of symbols in \mathbf{P} taken as a language (i.e., the composition $\mathbf{F}_P \circ \rho_{G-P}$). The ordered pair $\langle \mathbf{P}, \mathbf{F}_G \rangle$ defines the model \mathbf{M}_G for the geometrical interpretation of \mathbf{G} as geometrical objects; the geometrical interpretation function \mathbf{F}_G assigns a denotation for every constant of \mathbf{G} ; the denotation of individual constants of \mathbf{G} are the graphical symbols themselves, and the denotation of operators and function symbols of \mathbf{G} denoting graphical properties and relations will be given by predefined geometrical algorithms commonly used in computational geometry and computer graphics—see, for instance, Shamos (1978). The semantic interpretation of composite expressions of \mathbf{G} , on the other hand, is defined through a semantic algebra, as will be shown below in Section 2.3.2. The definition of this geometrical interpreter will allow us to perform inferences about the geometry of the drawing in a very effective fashion. Consider that to state explicitly all true and false geometrical statements about a drawing would be a very cumbersome task, as the number of statements that would have to be made even for small drawings would be very large. Note also that although a map can be an incomplete representation of the world (e.g., some cities might have been omitted), the geometrical algorithms associated with operators of \mathbf{G} will always provide complete information on the map as a geometrical object.

1.2 Multimodal interpretation

For the kind of problem exemplified in Figures 1 and 2, the objects in \mathbf{L} , \mathbf{P} , and \mathbf{G} are given, and the function \mathbf{F}_L establishes the relation between linguistic constants and the objects of the world that such constants happen to refer to. To interpret these multimodal messages, \mathbf{F}_P must be made explicit. If one asks *who is he?* looking at Figure 1, for instance, the answer is found by computing $\rho_{G-P}(\rho_{L-G}(he))$, whose value is the picture of the man on the drawing. Once this computation is performed, the picture can be highlighted or signaled by other graphical means. However, in other kinds of situations the knowledge of \mathbf{F}_P might be available and the purpose of the interpretation process could be to identify \mathbf{F}_L . If one points out the middle dot in Figure 2 at the time the question *what is this?* is asked, the answer can be found by applying the function $\rho_{G-L} \circ \rho_{P-G}$ to the dot indicated (i.e., $\rho_{G-L}(\rho_{P-G}(\bullet))$), whose value would be the word *Saarbrücken*. A similar situation arises in the interpretation of multimodal referring expressions. Consider the following example—also from André and Rist (1994)—in which a multimodal message is constituted by a picture of an espresso machine that has two switches, and by the textual expression *the temperature control*. In this scenario, the denotation of the natural language expression can be found by the human interpreter if the corresponding switch is identified in the picture through visual inspection (e.g., if the switch is highlighted). In general, multimodal coreference can be established if ρ_{L-G} and ρ_{G-L} are defined, as \mathbf{F}_P can be made explicit in terms of \mathbf{F}_L and vice versa.

In situations in which all theoretical elements illustrated in Figure 3 are given, questions about multimodal scenarios can be answered through the evaluation of expressions of a given modality in terms of the interpreters of the languages involved and the translation functions. However, when one is instructed to interpret a multimodal message, like Figures 1 and 2, not all information in the scheme of Figure 3 is available. In particular, the translation functions ρ_{L-G} and ρ_{G-L} of the graphical and linguistic individual constants mentioned in the texts and the pictures of the multimodal messages are not known, and the crucial inference of the interpretation process has as its goal to find out the definition of these functions (i.e., to establish the relations between names of **L** and **G**). It is important to emphasize that in order to find out ρ_{L-G} and ρ_{G-L} , the information overtly provided in the multimodal message is usually not enough, and in order to carry out such an interpretation process it will be necessary to consider the grammatical structure of the languages involved, the definition of translations rules between languages, and also conceptual knowledge stored in memory about the interpretation domain.

An additional consideration regarding the scheme in Figure 3 is related to the problem of ambiguity in the interpretation of multimodal messages. In the literature of intelligent multimodal systems, ambiguity is commonly seen from the perspective of human users. A multimodal referring expression constituted by the text *the temperature control* and a drawing with two switches is said to be ambiguous, for instance, if the human user is not able to tell which one is the temperature control. A well-designed presentation should avoid this kind of ambiguity by providing additional information either in a textual form (e.g., the temperature control is the switch on the left) or by a graphical focusing technique (e.g., highlighting the left switch). An important motivation in the design of intelligent presentation systems like WIP (Wahlster et al. 1993) and COMET (Feiner and McKeown 1993) is to generate graphical and linguistic explanations in which these kinds of ambiguities are avoided.² Note, however, that such situations are better characterized as problems of underspecification, rather than as problems of ambiguity, since the expression *the temperature control* has only one syntactic structure and one meaning, and the referent can be identified in a given context if enough information is available.

Ambiguity in multimodal systems has also been related to the granularity of graphical pointing acts. A map, for instance, can be represented by an expression of **G** that translates into a graphical composition in **P** denoting a single individual (e.g., Europe) or by a number of expressions of **G** that refer to the minimal graphical partitions in **P** (e.g., the countries of Europe) depending on whether the focus of the interpretation process is the whole of the drawing or its constituent parts. This problem has also been addressed in a number of intelligent multimodal systems like XTRA (Wahlster 1991) and AlFresco (Stock et al. 1993), but the lack of a formalized notion of graphical language (and also a better understanding of indexical expressions), has prevented a deeper analysis of this kind of ambiguity.

These notions of “ambiguity” in multimodal systems contrast with the traditional notion of ambiguity in natural language in which an ambiguous expression has several interpretations. The formalization of graphical representations through the definition of graphical languages with well-defined syntax and semantics allows us to face the problem of ambiguity directly in terms of the relation of translation between natural and graphical languages, and the semantics of expressions of both modal-

² It is also worth noticing that systems like WIP and COMET do not interpret multimodal messages input by human users through the interaction and, therefore, there is no ambiguity to be resolved.

ities. An interesting question is whether the graphical context offers clues that the parser can use to resolve lexical and structural ambiguity. Although we have yet to explore this issue, there are some antecedents in this regard. In Steedman's theory of incremental interpretation in dialogue, for instance, the rules of syntax, semantics, and processing are very closely linked (Steedman 1986) and local ambiguities may be resolved by taking into account their appropriateness to the context, which can be graphical. Structural ambiguity in G can be appreciated, for instance, in relation to the granularity of graphical objects, as the same drawing will have different syntactic analysis depending on whether it is interpreted as a whole or as an aggregation of parts. It is likely that the resolution of this latter kind of ambiguity is also influenced by pragmatic factors concerning the purpose of the task, the interpretation domain, and the attentional state of the interpreter, but this investigation is also pending.

We do, however, address issues of ambiguity related to the resolution of spatial indexical terms and anaphoric references in an integrated fashion. In Section 3, an incremental constraint satisfaction algorithm for resolving referential terms in relation to the graphical domain is presented. This algorithm relies on spatial constraints of drawings and general knowledge about the interpretation domain, and its computation is performed during the construction of multimodal discourse representation structures (MDRSs), which are extensions of DRSs in DRT (Kamp and Reyle 1993) as illustrated in Section 4. In the same way that DRT makes no provision for ambiguity resolution and alternative DRSs are constructed for different readings of a sentence, several MDRSs would have to be constructed in our approach for ambiguous multimodal messages.³ However, as natural language terms in L in our simplified domain refer to graphical objects, indefinites are very unlikely to have specific readings (e.g., "a city" normally refers to any city) and a simple heuristic in which indefinites are within the scope of definite descriptions and proper names can be used to obtain the preferred reading of sentences such as the one in Figures 2. Nevertheless, even if only this reading is considered, and the interpreter knows that the drawing is a map and is aware of the interpretation conventions of this kind of graphical representations (i.e., countries are represented by regions, cities by dots, etc.), drawings can still be ambiguous. In Figure 2, for instance, there are four possible interpretations for the graphical symbols that are consistent with the text if no knowledge of the geography of Europe is assumed. Our algorithm is designed to resolve reference for spatial referential and anaphoric terms in the course of the multimodal discourse interpretation, and the graphical ambiguity is resolved in the course of this process, as will be shown in detail in Sections 3 and 4.

To conclude this section, we believe the formalization of the syntax and semantics of graphical representations in a form compatible with the syntax and semantics of natural language, as in the scheme in Figure 3, may be a point of departure for investigating how the graphical or visual context helps to resolve natural language ambiguities at different levels of representation and processing.

³ A question for further research is whether our approach can be generalized to address problems of ambiguity by means of underspecified representations (e.g., van Deemter and Peters 1995). These representations result from the lexical and syntactic disambiguation process, but leave unspecified some information, like the interpretation of indexical references, the resolution of anaphoric expressions and the semantic scope of operators. A relevant antecedent related to our extension of multimodal DRSs is Poesio's extension of DRT into the so-called Conversational Representation Theory (Poesio 1994).

1.3 Multimodal Generation

An important motivation for the study of the interpretation of multimodal messages is the definition of multimodal presentation or explanation systems in which users are able to identify the referent of graphical and linguistic expressions easily. In WIP, for instance, a central concern is whether the human user is able to “activate” the relevant “representations” (presumably in his or her mind) and resolve the graphical and linguistic ambiguities and anaphors (using WIP’s terminology) present in multimodal messages. This is possible, in general, if the message conveys to the human user explicit interpretation paths from the information that is available overtly to the information that the user is expected to infer. The production of multimodal referring expressions in this kind of system depends on the use of presentation strategies defined in terms of rhetorical structures and intentional goals—e.g., along the lines of Rhetorical Structure Theory (RST) (Mann and Thompson 1988), and its computational implementation (Moore 1995). The use of a particular presentation strategy in a multimodal explanation (e.g., in WIP) depends crucially on whether the expressions generated on the basis of such a strategy satisfy the conditions defined to activate the expected representations in the user’s mind (an intentional goal). Furthermore, some rhetorical structures are designed explicitly to provide additional information to activate the expected representations if the conditions for the identification of the referent of an expression are not met. Consider again the resolution of the “ambiguity” in the interpretation of *the temperature control* example in WIP in which the presentation strategy provides the information required by the human user to identify the referent, either through the text *the temperature control is the switch on the left* or highlighting or pointing to the corresponding switch in the drawing. WIP is able to tell whether the presentation would be ambiguous for the human user if additional information were not provided because it has a representation of the actual situation and a simple model of the user’s beliefs.

Although the main representation structure of multimodal presentation and explanation systems is defined at a rhetorical level, the use of presentation strategies relies on algorithms for the generation of graphical and linguistic referring expressions. For instance, the “activate” presentation strategy of WIP (André and Rist 1994), the purpose of which is to establish a mutual belief between the human user and the system about the identity of an object, employs an algorithm for the generation of referring expressions based on an incremental interpretation algorithm proposed by Reiter and Dale (1992). It is interesting to note that presentations generated by WIP and other multimodal explanation systems like COMET (Feiner and McKeown 1993), or TEXPLAN (Maybury 1993), are limited to the production of definite descriptions only, even though the use of indefinite descriptions can be natural in multimodal communication. However, this restriction can be overcome with a more solid representational framework such as the one illustrated in Figure 3. Consider that basic or composite expressions of the languages G and L can be translated to basic or composite expressions of the other language, depending on the definition of the translation function. So, to refer linguistically to a graphical configuration, for instance, it would only be necessary to find an expression of G that succinctly expresses the relevant graphical properties of the desired object, and then translate it to its corresponding expression in L . The resulting natural language expression could be used directly or embedded in a larger natural language expression containing words that refer to abstract objects or properties. The descriptions obtained through this strategy explicitly employ the concrete and graphical properties of the representation, since expressions of G are

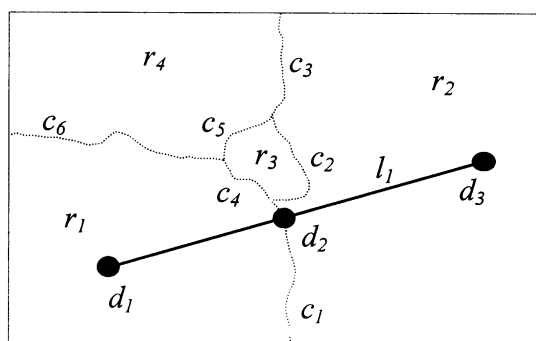


Figure 4
Labeling the graphical objects in Figure 2.

made up of constants and operators that directly describe the geometry of objects and configurations.

Consider the natural language text: *Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt*. This sentence contains the definite description *the intersection between the border between France and Germany and a line from Paris to Frankfurt*, which in turn contains *the border between France and Germany* and *a line from Paris to Frankfurt*. Finding the graphical referents of these expressions requires the identification of a dot, a curve, and a line on the map (i.e., the corresponding graphical objects). These graphical objects can be referred to directly through language; however, there are additional graphical entities on the map in Figure 2 that have an interpretation but are not mentioned explicitly in the text of the multimodal message. In Figure 4, for instance, Belgium is represented by the region r_4 , and the curve c_6 represents the border between France and Belgium. Once a picture has been interpreted, one would be entitled to ask not only for graphical objects that have been mentioned in the textual part of the message, but also for any meaningful graphical object. So, if one points to the curve c_6 in Figure 2 at the time the question *What is this?* is asked, the answer could be *the border between France and Belgium*, or alternatively, *the indefinite a border*. As some graphical objects named by constants of the graphical language do not have a proper name in natural language, the translation function ρ_{G-L} must associate a basic constant of G with a composite expression of L . The process of inducing such a translation function is closely related to the process of generating the corresponding natural language descriptions, and this relation will be explored further in Section 3.

In the rest of this paper, we discuss in more detail how the scheme for multimodal representation and interpretation in Figure 3 can be carried out. In Section 2, we present a formalization of the languages L , P , and G with their corresponding translation functions, along the lines of Montague's general semiotic program. The process of multimodal interpretation is explained, and the translation of expressions of one modality into expressions of another modality is illustrated. However, such a process can be carried out only if the translation functions are known, which is not normally the case in the interpretation of multimodal messages (as noted above). In Section 3, we offer an account of how such functions can be induced in terms of the message, constraints on the interpretation conventions of the modalities, and constraints on general knowledge of the domain. In this section we also illustrate the process of generating graphical and linguistic descriptions, which is associated with the induction of the translation functions. In Section 4, we discuss how to ex-

tend Kamp's DRS with multimodal structures. Finally, in Section 5, some concluding remarks and some directions for further work are presented.

2. A Multimodal System of Representation

In this section, we present the definition of the syntax and semantics of languages **L**, **P**, and **G**, illustrating the theory with the multimodal message of Figure 2. Language **L** is a segment of English designed to produce expressions useful for referring to objects, properties, and relations commonly found in discourse about maps. In particular, the natural language expression of Figure 2 can be constructed in a compositional fashion. The syntactic structure of **P**, on the other hand, imposes a restriction on the possible geometries of the family of drawings in the interpretation domain. Language **G** is a logical language in which interpretation and reasoning about geometrical configurations can be carried out. It is an interlingua representation for information expressed in both of the modalities.

The definitions of **L**, **P**, and **G** closely follow the general guidelines of Montague's semiotic program. As a first step in the syntactic definition of a language, the set of categories or types is stated. A number of constants—or basic expressions—for each type is defined, and the combination rules for producing composite expressions are stated. For each type of a source language, a corresponding type in the target language is assigned. Basic expressions of the source language can be mapped either to basic or to composite expressions of the corresponding type in the target language and vice versa. For each syntactic rule of a source language, a translation rule for mapping the expression formed by the rule into its translation in the target language is defined.

2.1 Definition of Language L

Language **L** contains the textual part of multimodal messages in the domain of maps. An expression of **L** is, for instance, *Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt*, which is the natural language part of Figure 2. Constants like *France* and *Germany*, and all subexpressions of the former sentence, like *the border between France and Germany* or *a line from Paris to Frankfurt* are also in **L**. In addition, **L** contains expressions like *France is a country*, *Frankfurt is a city of Germany* or *Germany is to the east of France*, which express general knowledge required in the interpretation of maps.

2.1.1 Syntactic Definition of L. The set of syntactic categories of **L** is as follows:

1. The basic syntactic categories of **L** are *t*, *IV*, *ADJ*, *CN*, and *CN'* where *t* is the category of sentences, *IV* is the category of intransitive verbs, *ADJ* is the category of adjectives, and *CN* and *CN'* are two categories of common nouns.
2. If *A* and *B* are syntactic categories then *A/B* is a category.⁴

Traditional syntactic categories of natural language like transitive verbs (*TV*), terms (*T*), prepositional phrases (*PP*), and determiners (*T/CN*) can be derived from the basic categories.

⁴ An expression of category *A/B* combines with an expression of category *B* to give an expression of category *A*.

Constant	Category name	Category definition
<i>Paris, Frankfurt, Saarbrücken, France, Germany</i>	<i>T</i>	<i>t/IV</i>
<i>city, country, border, line, intersection</i>	<i>CN</i>	<i>CN</i>
<i>east</i>	<i>CN'</i>	<i>CN'</i>
<i>big</i>	<i>ADJ</i>	<i>ADJ</i>
<i>be, lie at, be to</i>	<i>TV</i>	<i>IV/(t/IV)</i>
<i>be</i>	<i>IV/ADJ</i>	<i>IV/ADJ</i>
<i>a, the</i>	<i>T/CN</i>	<i>(t/IV)/CN</i>
	<i>PP</i>	<i>CN/CN</i>
	<i>PP'</i>	<i>CN/CN'</i>
	<i>IV</i>	<i>IV</i>

Figure 5
Constants of language *L*.

The table in Figure 5 illustrates the constants of *L* with their category names and category definitions. Common nouns are divided into *CN* and *CN'*. Expressions of category *CN* translate into graphical predicates (sets of graphical objects) while expressions of category *CN'* translate into abstract concepts. For instance, *city* translates into a set of dots representing cities, but *east* translates into a geometrical function from regions to zones (e.g., if the region representing France is the argument of this function, the zone to the right of that region is the function value). Prepositional phrases are divided into *PP* and *PP'* due to the classification of common nouns into *CN* and *CN'*. There are no basic constants of categories *PP*, *PP'*, and *IV*, as prepositional words are introduced syncategorematically and intransitive verb phrases are always composite expressions in this grammar. Transitive verbs are defined in a standard fashion, and the constant *be* of category *IV/ADJ* is used to form attributive sentences.

Next, the syntactic rules of *L* are presented. Each rule is shown in a separate item containing the purpose of the rule, the syntactic rule itself, and some examples of expressions that can be formed with the rule. Following Montague, syntactic rules and syntactic operations for combining symbols (for instance, F_{L1}) associated to each rule are separated. In the following, P_C is the set of expressions of category *C*.

SENTENCES

S1_L. If $\alpha \in P_T$ and $\beta \in P_{IV}$, then $F_{L1}(\alpha, \beta) \in P_t$, where $F_{L1}(\alpha, \beta) = \alpha\beta^*$, and β^* is the result of replacing the first *verb* in β by its third person singular present form.

Examples: -*Paris is a city of France*
 -*Germany is to the east of France*
 -*a country is big*
 -*Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt*

TRANSITIVE VERB PHRASES

S2_L. If $\alpha \in P_{TV}$ and $\beta \in P_T$, then $F_{L2}(\alpha, \beta) \in P_{IV}$, where $F_{L2}(\alpha, \beta) = \alpha\beta$.

Examples: -*be a city*
 -*be to the east of France*

ATTRIBUTIVE VERB PHRASES

S3_L. If $\alpha \in P_{IV/ADJ}$ and $\beta \in P_{ADJ}$, then $F_{L2}(\alpha, \beta) \in P_{IV}$.

Examples: *-be big*

TERMS

S4_L. If $\alpha \in P_{T/CN}$ and $\beta \in P_{CN}$ or $P_{CN'}$, then $F_{L3}(\alpha, \beta) \in P_T$, where $F_{L3}(\alpha, \beta) = \alpha^* \beta$, and α^* is α except in the case where α is *a* and the first word in β begins with a vowel; here, α^* is *an*.

Examples: *-a city*
-a city of France
-the border between France and Germany
-a line from Paris to Frankfurt
-the east of France

COMMON NOUNS

S5_L. If $\alpha \in P_{CN}$ and $\beta \in P_{PP}$, or $\alpha \in P_{CN'}$ and $\beta \in P_{PP'}$, then $F_{L2}(\alpha, \beta) \in P_{CN}$.

Examples: *-city of France*
-east of France
-border between France and Germany
-intersection between the border between France and Germany and a line from Paris to Frankfurt

of PREPOSITIONAL PHRASES⁵

S6_L. If $\alpha \in P_T$, then $F_{L4}(\alpha) \in P_{PP}$ or $P_{PP'}$, where $F_{L4}(\alpha) = \textit{of} \alpha$.

Examples: *-of France*
-of Germany
-of a country

between PREPOSITIONAL PHRASES

S7_L. If $\alpha, \beta \in P_T$, then $F_{L5}(\alpha, \beta) \in P_{PP}$, where $F_{L5}(\alpha, \beta) = \textit{between} \alpha \textit{ and } \beta$.

Examples: *-between France and Germany*
-between France and a country
-between the border between France and Germany and a line from Paris to Frankfurt

from-to PREPOSITIONAL PHRASES

S8_L. If $\alpha, \beta \in P_T$, then $F_{L6}(\alpha, \beta) \in P_{PP}$, where $F_{L6}(\alpha, \beta) = \textit{from} \alpha \textit{ to } \beta$.

Example: *-from Paris to Frankfurt*

⁵ Although *of*, *between*, and *from* have been introduced syncategorematically in **L** for simplicity, they could have been defined as constants of some category of **L**, and their translations into **G** would have been a composite expression of some graphical type.

Constant α	$F_L(\alpha)$
<i>Paris, Frankfurt, Saarbrücken, France, Germany</i>	Paris, Frankfurt, Saarbrücken, France, Germany
<i>city</i>	{Paris, Frankfurt, Saarbrücken, ...}
<i>country</i>	{France, Germany, ...}
<i>border</i>	{border between France and Germany, ...}
<i>line</i>	{line from Paris to Frankfurt, ...}
<i>intersection</i>	{intersection between the border between France and Germany and a line from Paris to Frankfurt, ...}
<i>east</i>	
<i>be, lie at, be to</i>	
<i>a, the</i>	

Figure 6
Interpretation of constants of language L .

2.1.2 Semantic Definition of L . The semantics of L is given in a model-theoretic fashion as follows: The interpretation domain is the world $\mathbf{W} = \{\text{Paris, Saarbrücken, Frankfurt, France, Germany, the border between France and Germany, ...}\}$. Let D_x be the set of possible denotations for expressions of type x , and for any types A and B , $D_{A/B} = D_A^{D_B}$ (i.e., the set of all functions from D_B to D_A). Let F_L be an interpretation function that assigns to each constant of type A a member of D_A . For the example in Figure 3, F_L is defined as shown in Figure 6.

Not every constant of L has an interpretation assigned by F_L ; in particular, words like *east*, *be*, *lie at*, and *be to* have no interpretation defined directly in L . In principle the definition of these constants could be stated as an object of the appropriate semantic type but this is not a straightforward enterprise. Consider, for instance, that the constant *east* of category CN' is a basic object (a kind of predicate), but the individual objects in its extension are not overtly defined in the interpretation domain. Furthermore, it is more natural to talk about the interpretation of composite predicates, like *east of France*, of which *east* is a part. However, even the interpretation of such composite predicates is problematic, as they have a vague spatial meaning. For these reasons, the interpretation of these constants is not defined explicitly as a part of the function F_L , but in terms of their translation into G , where a spatial meaning can be formally defined, as will be shown below. A similar strategy is used for the interpretation of spatial prepositions; although *of*, *between*, and *from-to* were introduced syncategorematically in the syntax of L , they could have been defined as objects of an appropriate category and their semantics could have been given explicitly through F_L or, alternatively, through their translation into intensional logic along the lines of PTQ. However, the semantic type of such objects is extraordinarily complex, and the actual definition of these constants is seldom seen in the literature.⁶ In our system the interpretation of spatial prepositions will also be given in terms of the translation into G and the interpretation of P . Note also that no interpretation has been defined for the determiners *a* and *the*. One strategy for assigning a denotation would be to translate these constants into intensional logic, but this would be required only for a larger fragment of English in which reference

⁶ In PTQ, prepositions—of category $(IV/IV)/T$ —are treated semantically as functions that apply to sets of properties to give functions from properties to properties, but no explicit example of the actual semantic value of any of these constants is provided. In our system it will be possible to compute the semantic value of spatial prepositional phrases in an effective manner, yet the approach is fully compatible with intensional logic.

to space was not the focus of study. In our approach the determiners will be interpreted in terms of their translations into **G** in which high-order functions can be expressed.

In summary, the semantics of some constants and all composite expressions of **L** will be given in terms of their translations into **G** and **P**. Note that according to the scheme in Figure 3, if the translations between **L** and **G**, and **G** and **P** are defined, and the semantic interpretation of **P** is overtly defined, the interpretation of the natural language expressions can be found. Although the semantics of **L** is not further discussed in this paper, we consider that the interpretation of linguistic expressions referring to spatial situations could be embedded in a larger fragment of English, and a full semantic interpretation would have to be given by translating English into intensional logic. In such a model the semantic value of spatial prepositions would be left undefined, expressions referring to spatial configurations would be translated into **G**, and the interpretation of expressions of **G** would be embedded within the interpretation of intensional logic.

2.2 Definition of Language **P**

In this section, the syntax and semantics of language **P** are formally defined. The purpose of these definitions is to characterize the family of drawings that can be interpreted as maps, and to discriminate these drawings from other kinds of graphical configurations constituted by dots, curves, and regions. This notion of a multimodal system of representation in which objects in the graphical modality are formalized through a well-defined language is similar to the notion of graphical language introduced by Mackinlay for the automatic design of graphical presentations (Mackinlay 1987), where a number of graphical languages (e.g., the languages of bar charts, area and position graphs, scatter plots, etc.) are formally specified. In Mackinlay's work, expressions of graphical languages are related to the objects of the world that they represent through an *encodes* relation with three arguments: the graphical constant or expression performing the representation, the object of the world that is represented through the graphical expression, and the graphical language to which the graphical expression belongs.⁷ The formalization of **P** permits us to define a precise statement of expressiveness of a graphical language, as follows: "a set of facts is expressible in a language (graphical) if the language contains a sentence that encodes every fact in the set and does not encode any additional facts" (Mackinlay 1987, 54). The formalization additionally allows empirical studies to determine how effectively a human user can interpret expressions of a particular graphical language in relation to another in which the same set of facts is encoded. Although all graphical languages studied by Mackinlay are conventional and have a precise geomet-

⁷ Incidentally, a similar encoding relation *encodes* is used in the WIP system to relate the representational object to the object that it represents, but the third argument of this relation in WIP is a **context space** that allows use of the same presentation in different perspectives (e.g., an espresso machine may refer to an individual machine in a context space, or alternatively it can be seen as the prototype of espresso machines in a different context space). The *encodes* relation in WIP and in Mackinlay is similar to the translation relation between objects of **P** (or **G**) and **L** in our system, and we can think of a graphical language as a language encoding the information that is intended to be communicated. However, it is interesting to note that the status of the "linguistic" argument of the *encodes* relation is different in WIP and in Mackinlay's system. In the former, it is an "internal representation"—a psychological notion—while in the latter it stands for an object or a relation in the world itself—a semantic notion. In our approach, on the other hand, there are no "internal representations" and the translation relates graphical and linguistic expressions that are both "external" and that both refer to the world through a well-defined semantics.

Constant	Type
d_1, d_2, d_3, \dots	<i>dot</i>
l_1, l_2, l_3, \dots	<i>line</i>
c_1, c_2, c_3, \dots	<i>curve</i>
r_1, r_2, r_3, \dots	<i>region</i>
z_1, z_2, z_3, \dots	<i>zone</i>
cr_1, cr_2, cr_3, \dots	<i>composite_region</i>
$\emptyset, ds_1, ds_2, \dots$	<i>dot_set</i>
$\emptyset, ls_1, ls_2, \dots$	<i>line_set</i>
m_1, m_2, m_3, \dots	<i>map</i>

Figure 7
Constants of language **P**.

rical characterization, the notions of expressiveness and effectiveness of graphical languages can be applied to more unruly graphical domains (e.g., maps are analogical representations with a diagrammatic conventional component) as long as a formalization for the family of drawings can be approximated. Here, the question of whether arbitrary families of graphical objects can be formalized through a well-defined syntax is left open, and although it is possible to think of many families of drawings with very arbitrary geometries, some important efforts have been made in the characterization of design and other kinds of objects—see, for instance, shape grammars (Stiny 1975). Another related issue that is relevant for the construction of multimodal interactive systems is whether it is possible and useful to input expressions of **P** directly, and to obtain their syntactic structure through graphical parsing techniques (Wittenburg 1998). In summary, the purpose of formalizing **P** is to be able to talk about maps as a modality, where a modality, in our sense, is a code system for the symbols expressed in a medium, and a multimodal system of representation relates information expressed through different code systems in a systematic fashion.


2.2.1 Syntactic Definition of P. The types of **P** are *dot*, *line*, *curve*, *region*, *zone*, *composite_region*, *dot_set*, *line_set*, and *map*. Let C_s be the set of constants of type s , and E_s the set of well-formed expressions of graphical type s . Although the constants of **P** are the actual graphical marks on the screen or a piece of paper, a number of labels for facilitating the presentation are illustrated in Figure 7.

For the syntactic definition of **P** we capitalize on the distinction introduced by Montague between syntactic rules and syntactic operations. This distinction is based on the observation that “syntactic rules can be thought of as comprising two parts: one which specifies under what conditions the rule is to be applied, and the other which specifies what operation to perform under those conditions” (Dowty, Wall, and Peters 1985, 254). While a syntactic rule comprises both parts and defines the syntactic structure of an expression, the syntactic operation is a rule that depends on—or at least takes into account—the shape of the symbols and the medium in which the symbols are substantially realized. For instance, the syntactic operation F_{L5} in the rule $S7_L$ (i.e., $F_{L5}(\alpha, \beta) = \textit{between } \alpha \textit{ and } \beta$) combines the symbols *between* and *and* with the arguments to form the linear string indicated by the operation. For the definition of syntactic operations of **P** we generalize the operations that manipulate strings of symbols into general geometrical operations on the shapes of the graphical symbols on the paper or the screen, and these manipulations are defined according to certain geometrical conditions.

The definition of well-formed expressions of **P** is as follows:

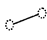
CONSTANT

S1_P. If $\alpha \in C_s$ then $\alpha \in E_s$.

Examples: •, /, 


LINE

S2_P. If $\alpha, \beta \in E_{dot}$ then $F_{P1}(\alpha, \beta) \in E_{line}$ where $F_{P1}(\alpha, \beta)$ is a line from α to β .

Example:  (the resulting graphical expression is only the line)


CURVE

S3_P. If $\alpha, \beta \in E_{region}$ such that α and β are adjacent then $F_{P2}(\alpha, \beta) \in E_{curve}$ where $F_{P2}(\alpha, \beta)$ is the curve between α and β .

Example:  (the resulting graphical expression is only the curve)

INTERSECTION

S4_P. If $\alpha \in E_{curve}$ and $\beta \in E_{line}$ then $F_{P3}(\alpha, \beta) \in E_{dot}$ where $F_{P3}(\alpha, \beta)$ is the dot in the intersection between α and β .

Example:  (the resulting graphical expression is only the dot)

RIGHT

S5_P. If $\alpha \in E_{region}$ then $F_{P4}(\alpha) \in E_{zone}$ where $F_{P4}(\alpha)$ is the zone to the right of the region α (the interpretation of “right” will be given below in the semantics of language **G**).

Example:  (the resulting graphical expression is only the gray zone)

DOT INSIDE A REGION

S6_P. If $\alpha \in E_{region}$ then $F_{P5}(\alpha) \in E_{dot}$ where $F_{P5}(\alpha)$ is the drawing of a dot inside α .

Example:  (the resulting graphical expression is only the dot)

COMPOSITE REGION (1)⁸

S7_P. If $\alpha, \beta \in C_{region}$ such that α and β are adjacent then
 $F_{P6}(\alpha, \beta) \in E_{composite_region}$ where $F_{P6}(\alpha, \beta)$ is the drawing of α and β .

COMPOSITE REGION (2)

S8_P. If $\alpha \in C_{region}$ and $\beta \in E_{composite_region}$ such that α and β are adjacent then
 $F_{P6}(\alpha, \beta) \in E_{composite_region}$.

SET OF DOTS

S9_P. If $\alpha \in E_{dot_set}$ and $\beta \in C_{dot}$ then $F_{P6}(\alpha, \beta) \in E_{dot_set}$.

SET OF LINES

S10_P. If $\alpha \in E_{line_set}$ and $\beta \in C_{line}$ then $F_{P6}(\alpha, \beta) \in E_{line_set}$.

MAP

S11_P. If $\alpha \in E_{composite_region}$, $\beta \in E_{dot_set}$ and $\delta \in E_{line_set}$ then $F_{P7}(\alpha, \beta, \delta) \in E_{map}$
 where $F_{P7}(\alpha, \beta, \delta)$ is the drawing of α , β and δ .

With the help of this grammar it is possible to draw maps like the one illustrated in Figure 2. Note that the basic object in this particular graphical construction is the region. The idea is to successfully construct a map from its constituting regions (i.e., as in a jigsaw puzzle) until the full map is produced. Once the map is constructed, other kinds of objects with conventional meanings, like dots and lines, can be drawn upon the assembly of regions. Consider Figure 8 in which the syntactic structure of the map in Figure 4 is shown. Note that the decision to use regions as basic objects in the graphical composition is not mandatory, and alternative constructions are possible; for instance, we could have designated curves as basic objects and obtained regions as compositions made out of curves. The set of graphical symbols included in a graphical syntactic tree of a map will be called the **base**. For instance, the base of the map in Figure 8 is the set $\{d_1, d_2, d_3, l_1, r_1, r_2, r_3, r_4\}$. The base is just the set of graphical objects that are taken as the **atoms** of the graphical composition in each particular interpretation task, and different graphical grammars would select different types of graphical objects for the base.

The purpose of this grammar is illustrative; we make no claims about what constitutes a map. **P** imposes very few constraints on graphical expressions, and many configurations that can be produced with these rules might not count as maps; in addition, **P** is not expressive enough to characterize a large number of objects that would be normally interpreted as maps. Another consideration is that graphical objects can be used either as basic building blocks of the construction, or as objects produced by graphical compositions (which we call **emergent objects**); for instance, in the grammar of **P**, regions are basic objects but curves are produced by graphical compositions. Additionally, in some contexts the interpretation of the graphical expression as a whole

⁸ Examples for the rules S7_P to S11_P are included in the construction of the map in Figure 8, as explained below.

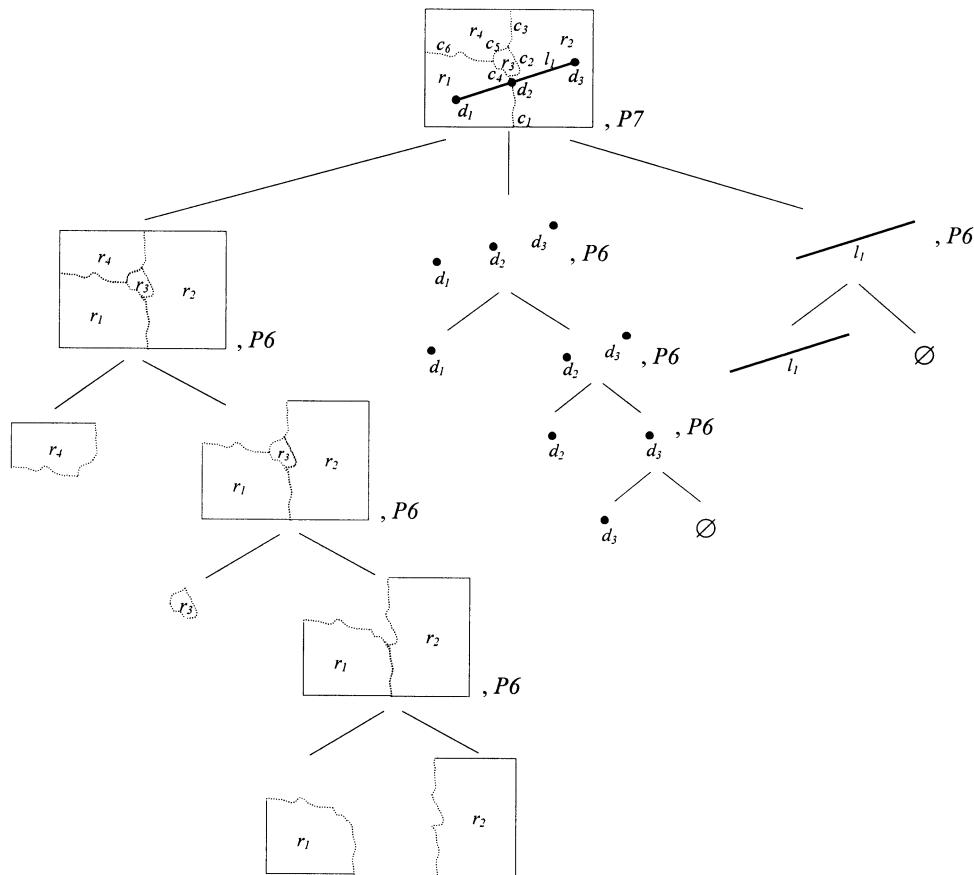


Figure 8
Construction of a map.

may be required but in others only the interpretation of some of the parts may be relevant; for instance, although curves are not a part of the syntactic tree in Figure 8 they can be generated and translated into \mathbf{G} when required through rules $S3_P$ and $T3_{P-G}$ as long as the composition is made out of regions included in the base of the map. Had the grammar allowed the generation of composite regions out of regions of the base, these emergent objects could also be used for the generation of curves. Another consideration is that expressions of type *map* are in general ambiguous as they have several syntactic analyses, but since this feature is harmless for the current discussion we do not pursue the issue further. A final remark is that alternative grammars could be defined for characterizing the same class of drawings with different consequences in the syntax and the semantics. One possibility, for instance, is to define a syntactic operation that takes two adjacent regions and produces the union of the regions as one single emerging region, instead of the set of the two regions as currently defined. Such a rule would be similar to the rule that combines two regions to produce a curve, and it would be useful in applications like XTRA (Wahlster 1991), in which the ambiguity of pointing to a part or the whole is intended to be resolved.

2.2.2 Semantic Definition of P. The semantics of \mathbf{P} is given in a model-theoretic fashion as follows: Let $\mathbf{W} = A_{\text{city}} \cup A_{\text{line}} \cup A_{\text{border}} \cup A_{\text{country}} \cup A_{\text{zone}}$ be the world. Let D_x

Constant α	$F_P(\alpha)$
d_1, d_2, d_3, \dots	Paris, Saarbrücken, Frankfurt, ...
l_1, l_2, l_3, \dots	line from Paris to Frankfurt, ...
c_1, c_2, c_3, \dots	border between France and Germany, ...
r_1, r_2, r_3, \dots	France, Germany, ...
z_1, z_2, z_3, \dots	east of France, east of Germany, ...
cr_1, cr_2, cr_3, \dots	region formed by France and Germany, ...
$\emptyset, ds_1, ds_2, \dots$	sets of cities
$\emptyset, ls_1, ls_2, \dots$	set of lines
m_1, m_2, m_3, \dots	maps

Figure 9
Semantics of constants of **P**.

be the set of possible denotations for expressions of type x , such that $D_{dot} = A_{city}$, $D_{line} = A_{line}$, $D_{curve} = A_{border}$, $D_{region} = A_{country}$, $D_{zone} = A_{zone}$, and, for any types a and b , $D_{\langle a,b \rangle} = D_b^{D_a}$ (i.e., the set of all functions from D_a to D_b). Let F_P be an interpretation function that assigns to each constant of type a a member of D_a . The interpretations of the constants are presented in Figure 9.

Following Montague, we adopt the notational convention by which the semantic value or denotation of an expression α with respect to a model M is expressed as $[[\alpha]]^M$. The semantic rules for interpreting language **L** are the following:

CONSTANT

M1_P. If $\alpha \in C_s$ then $[[\alpha]]^M = F_P(\alpha)$.

LINE

M2_P. If $\alpha, \beta \in E_{dot}$ then $[[F_{P1}(\alpha, \beta)]]^M$ is a line from $[[\alpha]]^M$ to $[[\beta]]^M$.

CURVE

M3_P. If $\alpha, \beta \in E_{region}$ such that α and β are adjacent then $[[F_{P2}(\alpha, \beta)]]^M$ is the border between $[[\alpha]]^M$ and $[[\beta]]^M$.

INTERSECTION

M4_P. If $\alpha \in E_{curve}$ and $\beta \in E_{line}$ then $[[F_{P3}(\alpha, \beta)]]^M$ is the intersection between $[[\alpha]]^M$ and $[[\beta]]^M$.

RIGHT

M5_P. If $\alpha \in E_{region}$ then $[[F_{P4}(\alpha)]]^M$ is the east of $[[\alpha]]^M$.

DOT INSIDE A REGION

M6_P. If $\alpha \in E_{region}$ then $[[F_{P5}(\alpha)]]^M$ is a city of $[[\alpha]]^M$.

COMPOSITE REGION (1)

M7p. If $\alpha, \beta \in C_{region}$ such that α and β are adjacent then $[[F_{P5}(\alpha, \beta)]]^M$ is the union of $\{[[\alpha]]^M\}$ and $\{[[\beta]]^M\}$.

COMPOSITE REGION (2)

M8p. If $\alpha \in C_{region}$ and $\beta \in E_{composite_region}$ such that α and β are adjacent then $[[F_{P5}(\alpha, \beta)]]^M$ is the union of the sets $\{[[\alpha]]^M\}$ and $[[\beta]]^M$.

SET OF DOTS

M9p. If $\alpha \in E_{dot_set}$ and $\beta \in C_{dot}$ then $[[F_{P5}(\alpha, \beta)]]^M$ is the union of the sets $[[\alpha]]^M$ and $\{[[\beta]]^M\}$.

SET OF LINES

M10p. If $\alpha \in E_{line_set}$ and $\beta \in C_{line}$ then $[[F_{P5}(\alpha, \beta)]]^M$ is the union of the sets $[[\alpha]]^M$ and $\{[[\beta]]^M\}$.

MAP

M11p. If $\alpha \in E_{composite_region}$, $\beta \in E_{dot_set}$ and $\delta \in E_{line_set}$ then $[[F_{P6}(\alpha, \beta, \delta)]]^M$ is the union of the sets $[[\alpha]]^M$, $[[\beta]]^M$ and $[[\delta]]^M$.

2.3 Definition of Language G

In this section the syntax and semantics of the graphical language **G** are formally stated. **G** is defined along the lines of intensional logic, and it is expressive enough to refer to graphical symbols and configurations, on the one hand, and to express the translation of quantified expressions of **L**, on the other.

2.3.1 Syntactic Definition of G. The types of the language **G** are as follows:⁹

1. e is a type (graphical objects).
2. t is a type (truth values).
3. If a and b are any types, then $\langle a, b \rangle$ is a type.¹⁰
4. Nothing else is a type.

Let V_s be the set of variables of type s , C_s the set of constants of type s , and E_s the set of well-formed expressions of graphical type s . The constants of **G** are presented in Figure 10. Note that constants like *right*, *curve_between*, etc. have an

⁹ A simplifying assumption rests on the consideration that the interpretations of all expressions included in these languages depend only on the current graphical state and no intensional types are included in the system. However, this analysis can be extended along the lines of intensional logic to be able to deal with a more comprehensive fragment of English.

¹⁰ An expression of type $\langle a, b \rangle$ combines with an expression of type a to give an expression of type b .

Constant	Type
$d_1, d_2, d_3, r_1, r_2, r_3, r_4, l_1$	e
<i>dot, region, curve, line, intersection</i>	$\langle e, t \rangle$
<i>right</i>	$\langle \langle \langle e, t \rangle, t \rangle, \langle e, t \rangle \rangle$
<i>lie_at, be_in_zone, inside</i>	$\langle \langle \langle e, t \rangle, t \rangle, \langle e, t \rangle \rangle$
$=$	$\langle e, \langle e, t \rangle \rangle$
$\wedge, \vee, \leftrightarrow$	$\langle t, \langle t, t \rangle \rangle$
<i>curve_between, intersection_between, line_from_to</i>	$\langle \langle \langle e, t \rangle, t \rangle, \langle \langle \langle e, t \rangle, t \rangle, \langle e, t \rangle \rangle \rangle$
<i>right*</i>	$\langle e, e \rangle$
<i>lie_at_*, be_in_zone_*, inside_*</i>	$\langle e, \langle e, t \rangle \rangle$
<i>curve_between_*, intersection_between_*, line_from_to*</i>	$\langle e, \langle e, e \rangle \rangle$

Figure 10
Constants of language **G**.

associated *right**, *curve_between**, etc. The unsubscripted version of these constants denotes a relation between sets of properties of graphical individuals and the subscripted version denotes the corresponding geometrical relation between individuals; the type-raised version is used for preserving quantification properties in the translation process from **L** into **G**, while the subscripted version is used for computing the geometry associated with the corresponding relation, as will be shown below in Section 2.3.2.

G is a formal language with constants and variables for all types, functional abstraction and application, and existential and universal quantification. The syntactic rules of **G** are as follows:

1. If $\alpha \in C_s$, then $\alpha \in E_s$.
2. If $\mu \in V_s$, then $\mu \in E_s$.
3. If $\alpha \in E_{\langle a, b \rangle}$ and $\beta \in E_a$, then $\alpha(\beta) \in E_b$.
4. If $\alpha \in E_a$ and $u \in V_b$, then $\lambda u[\alpha] \in E_{\langle b, a \rangle}$.
5. If $\mu \in V_s$ and $\beta \in E_t$ then $\exists \mu(\beta) \in E_t$.
6. If $\mu \in V_s$ and $\beta \in E_t$ then $\forall \mu(\beta) \in E_t$.

G is a very expressive language and not every well-formed expression has a translation into **L** as will be further discussed in Section 2.5. Useful translations are, for instance, names and descriptions of geometrical objects and configurations. Next, the definition of expressions of **G** that have a translation into **L** is presented. For clarity, the abbreviations in Figure 11 are used.

Two geometrical interpretations are given for the spatial prepositions *of* and *between*. Although the characterization of the meaning of these words is a very complex problem that is beyond the scope of this paper, we allow that spatial prepositions can be interpreted in more than one way, as long as each interpretation is stated in terms of a geometrical algorithm explicitly defined in **G**. For instance, the spatial meaning of *of* is different in *city of France* and *east of France*. In the former, *of* denotes a spatial inclusion relation (OF_a), but in the latter it denotes a relation of adjacency (OF_b). Similarly, the spatial meaning of *between* in *border between France and Germany* and its first occurrence in *intersection between the border between France and Germany and a line from Paris to Frankfurt* is different, as it denotes a curve in the first case ($BETWEEN_a$) and a dot in the second ($BETWEEN_b$).

Abbreviation	Formal expression	Type
A	$\lambda P \lambda Q \exists x [P(x) \wedge Q(x)]$	$\langle \langle e, t \rangle, \langle \langle e, t \rangle, t \rangle \rangle$
THE	$\lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)]$	$\langle \langle e, t \rangle, \langle \langle e, t \rangle, t \rangle \rangle$
BE _a	$\lambda P \lambda x P(\lambda y [x = y])$	$\langle \langle \langle e, t \rangle, t \rangle, \langle e, t \rangle \rangle$
BE _b	$\lambda P \lambda x P(x)$	$\langle \langle e, t \rangle, \langle e, t \rangle \rangle$
D _i	$\lambda P [P(d_i)]$	$\langle \langle e, t \rangle, t \rangle$
R _i	$\lambda P [P(r_i)]$	$\langle \langle e, t \rangle, t \rangle$
OF _a	$\lambda x_{\langle \langle e, t \rangle, t \rangle} \lambda y_{\langle e, t \rangle} \lambda z_e [y(z) \wedge \textit{inside}(x)(z)]$	$\langle \langle \langle e, t \rangle, t \rangle, \langle \langle e, t \rangle, \langle e, t \rangle \rangle \rangle$
OF _b	$\lambda x_{\langle \langle e, t \rangle, t \rangle} \lambda y_{\langle \langle \langle e, t \rangle, t \rangle, e \rangle} \lambda z_e [y(x)(z)]$	$\langle \langle \langle e, t \rangle, t \rangle, \langle \langle \langle \langle e, t \rangle, t \rangle, e \rangle, \langle e, t \rangle \rangle \rangle$
BETWEEN _a	$\lambda x_{\langle \langle e, t \rangle, t \rangle} \lambda y_{\langle \langle e, t \rangle, t \rangle} \lambda z_{\langle e, t \rangle} \lambda u_e [z(u) \wedge \textit{curve_between}(x)(y)(u)]$	$\langle \langle \langle e, t \rangle, t \rangle, \langle \langle \langle e, t \rangle, t \rangle, \langle \langle e, t \rangle, \langle e, t \rangle \rangle \rangle \rangle$
BETWEEN _b	$\lambda x_{\langle \langle e, t \rangle, t \rangle} \lambda y_{\langle \langle e, t \rangle, t \rangle} \lambda z_{\langle e, t \rangle} \lambda u_e [z(u) \wedge \textit{intersection_between}(x)(y)(u)]$	$\langle \langle \langle e, t \rangle, t \rangle, \langle \langle \langle e, t \rangle, t \rangle, \langle \langle e, t \rangle, \langle e, t \rangle \rangle \rangle \rangle$
FROM_TO	$\lambda x_{\langle \langle e, t \rangle, t \rangle} \lambda y_{\langle \langle e, t \rangle, t \rangle} \lambda z_{\langle e, t \rangle} \lambda u_e [z(u) \wedge \textit{line_from_to}(x)(y)(u)]$	$\langle \langle \langle e, t \rangle, t \rangle, \langle \langle \langle e, t \rangle, t \rangle, \langle \langle e, t \rangle, \langle e, t \rangle \rangle \rangle \rangle$

Figure 11
Shorthand definitions.

The restrictions for the expressions of **G** that can be translated into **L** are given below. In rules S6_G to S8_G, Q stands for either the quantifier **A** or **THE**.

SENTENCES

S1_G. If $\alpha \in E_{\langle \langle e, t \rangle, t \rangle}$ and $\beta \in E_{\langle e, t \rangle}$, then $F_{G1}(\alpha, \beta) \in E_t$, where $F_{G1}(\alpha, \beta) = \alpha(\beta)$.

Examples:

- D₁ (BE_a (A (OF_a(R₁) (dot))))
- R₃ (be_in_zone (THE (OF_b(R₁) (right))))
- A (region) (BE_b(big))
- D₃ (lie_at (THE (BETWEEN_b(THE (BETWEEN_a(R₁) (R₃) (curve))) (A (FROM_TO(D₁) (D₃) (line))) (intersection))))

TRANSITIVE VERB PHRASES

S2_G. If $\alpha \in E_{\langle \langle \langle e, t \rangle, t \rangle, \langle e, t \rangle \rangle}$ and $\beta \in E_{\langle \langle e, t \rangle, t \rangle}$ then $F_{G1}(\alpha, \beta) \in E_{\langle e, t \rangle}$.

Examples:

- BE_a (A (dot))
- be_in_zone (THE (OF_b(R₁)(right)))

ATTRIBUTIVE VERB PHRASES

S3_G. If $\alpha \in E_{\langle \langle e, t \rangle, \langle e, t \rangle \rangle}$ and $\beta \in E_{\langle e, t \rangle}$ then $F_{G1}(\alpha, \beta) \in E_{\langle e, t \rangle}$.

Example:

- BE_b (big)

TERMS

S4_G. If $\alpha \in E_{\langle \langle e, t \rangle, \langle \langle e, t \rangle, t \rangle \rangle}$ and $\beta \in E_{\langle e, t \rangle}$, then $F_{G1}(\alpha, \beta) \in E_{\langle \langle e, t \rangle, t \rangle}$.

Examples:

- A (dot)
- A (OF_a(R₁)(dot))
- THE (BETWEEN_a(R₁) (R₂) (curve))
- A (FROM_TO(D₁) (D₃) (line))
- THE (OF_b(R₁)(right))

COMMON NOUNS

S5_G. If $\alpha \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$ and $\beta \in E_{\langle e,t \rangle}$, or $\alpha \in E_{\langle\langle\langle e,t \rangle, t \rangle, e \rangle, \langle e,t \rangle\rangle}$ and $\beta \in E_{\langle\langle\langle e,t \rangle, t \rangle, e \rangle}$, then $F_{G1}(\alpha, \beta) \in E_{\langle e,t \rangle}$.

Examples: - $OF_a(R_1)(dot)$
 - $OF_b(R_1)(right)$
 - $BETWEEN_a(R_1)(R_2)(curve)$
 - $BETWEEN_b((THE(BETWEEN_a(R_1)(R_2)(curve))) (A(FROM_TO(D_1)(D_3)(line)))(intersection)$

of PREPOSITIONAL PHRASES

S6_G. If $\alpha \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α is either R_i or $Q(region)$, then $F_{G2}(\alpha) \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$ and $F_{G3}(\alpha) \in E_{\langle\langle\langle e,t \rangle, t \rangle, \langle e,t \rangle\rangle, \langle e,t \rangle\rangle}$, where $F_{G2}(\alpha) = OF_a(\alpha)$ and $F_{G3}(\alpha) = OF_b(\alpha)$

Examples: - $OF_a(R_1)$
 - $OF_b(R_2)$
 - $OF_a(A(region))$

between PREPOSITIONAL PHRASES

S7_G. (a) If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α, β are either R_i or $Q(region)$, then $F_{G4}(\alpha, \beta) \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$, where $F_{G4}(\alpha, \beta) = BETWEEN_a(\alpha)(\beta)$.
 (b) If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α is either C_i or $Q(curve)$ and β is either L_i or $Q(line)$, then $F_{G5}(\alpha, \beta) \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$, where $F_{G5}(\alpha, \beta) = BETWEEN_b(\alpha)(\beta)$.

Examples: - $BETWEEN_a(R_1)(R_2)$
 - $BETWEEN_a(R_1)(A(region))$
 - $BETWEEN_b((THE (BETWEEN_a(R_1)(R_2)(curve))) (A (FROM_TO(D_1)(D_3)(line)))$

from-to PREPOSITIONAL PHRASES

S8_G. If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α, β are either D_i or $Q(dot)$, then $F_{G5}(\alpha, \beta) \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$, where $F_{G5}(\alpha, \beta) = FROM_TO(\alpha)(\beta)$

Example: - $FROM_TO(D_1)(D_3)$

2.3.2 Semantic Definition of G. The interpretation of expressions of **G** is defined in relation not to the world **W** but to a domain constituted by the graphical objects in **P**. For this reason, we refer to the interpreter of **G** as a geometrical interpreter, and to the process of interpreting expressions of **G** as a geometrical interpretation process. The semantics of **G** is given in a model-theoretic fashion as follows: Let $\mathbf{P}_{base} = \{d_1, d_2, d_3, r_1, r_2, r_3, r_4, l_1\}$ be the set of basic graphical objects shown in Figure 4. Let **P** be the union of \mathbf{P}_{base} and all graphical objects that can be produced from \mathbf{P}_{base} with the help of geometrical functions: the emergent objects. Emergent objects can also be produced on the basis of other emergent objects previously generated. A particular

kind of emergent object that is interesting for the current discussion is the zone of a map that is considered to be the east of a region. For the production of emergent objects in \mathbf{P} there is a well-defined computational geometry algorithm associated with an operator symbol of \mathbf{G} , as will be seen below.

Let D_x be the set of possible denotations for expressions of type x , such that $D_e = \mathbf{P}$, $D_t = \{1, 0\}$, and, for any types a and b , $D_{\langle a,b \rangle} = D_b^{D_a}$. Let \mathbf{F}_G be an interpretation function that assigns to each constant of type a a member of D_a . For every graphical object φ in \mathbf{P}_{base} there is a constant α of type e such that $\mathbf{F}_G(\alpha) = \varphi$; for our example, \mathbf{F}_G assigns the objects $d_1, d_2, d_3, r_1, r_2, r_3, r_4$, and l_1 to the constants $d_1, d_2, d_3, r_1, r_2, r_3, r_4$, and l_1 , respectively. The interpretation (assigned by \mathbf{F}_G) of the geometrical-type predicates *dot*, *region*, *curve*, *line*, *intersection* are the sets containing the corresponding graphical objects. The constants *right**, *lie_at**, *be_in_zone**, *inside**, *curve_between**, *intersection_between**, and *line_from_to** are interpreted as geometrical functions. If the arguments of these geometrical functions are of an appropriate type, expressions containing these constants can be properly interpreted through geometrical algorithms; otherwise, these expressions have no denotation in \mathbf{G} and, as a consequence, their translations into \mathbf{L} also lack denotation. For further discussion of the interpretation of graphical expressions that have no proper graphical referent in the interpretation state, see Pineda (1992).

Following Montague, the interpretation of variables is defined in terms of an assignment function g . We adopt the notational convention by which the semantic value or denotation of an expression α with respect to a model M and a value assignment g is expressed as $[[\alpha]]^{M,g}$.

The semantic rules for interpreting expressions of \mathbf{G} are the following:

1. If $\alpha \in C_s$, then $[[\alpha]]^M = \mathbf{F}_G(\alpha)$.
2. If $\mu \in V_s$, then $[[\mu]]^{M,g} = g(\mu)$.
3. If $\alpha \in E_{\langle a,b \rangle}$, and $\beta \in E_a$, then $[[\alpha(\beta)]]^{M,g} = [[\alpha]]^{M,g}([[\beta]])^{M,g}$
4. If $\alpha \in E_a$ and $u \in V_b$, then $[[\lambda u[\alpha]]]^{M,g}$ is that function h from D_b into D_a such that for all objects k in D_b , $h(k)$ is equal to $[[\alpha]]^{M,g'}$, where g' is exactly like g except that $g'(u) = k$.
5. If $\mu \in V_s$ and $\beta \in E_t$ then $[[\exists \mu(\beta)]]^{M,g} = 1$ iff for some value assignment g' such that g' is exactly like g except possibly for the individual assigned to μ by g' , $[[\beta]]^{M,g'} = 1$.
6. If $\mu \in V_s$ and $\beta \in E_t$ then $[[\forall \mu(\beta)]]^{M,g} = 1$ iff for every value assignment g' such that g' is exactly like g except possibly for the individual assigned to μ by g' , $[[\beta]]^{M,g'} = 1$.

In order to capture the translation of expressions of \mathbf{L} into \mathbf{G} compositionally, while preserving the quantificational properties of the original source natural language expression, terms in \mathbf{G} referring to graphical objects are type-raised; consequently, graphical predicates like *be_in_zone*, *curve_between*, and *inside* have type-raised arguments. The expression *curve_between*($\lambda P[P(r_1)]$)($\lambda P[P(r_2)]$)(x), for instance, refers to the curve x between regions r_1 and r_2 ; the first two arguments refer not to the regions themselves, but to the set of properties that such regions have. Similarly, the expression *inside*($\lambda P \exists y[\text{region}(y) \wedge P(y)]$)(z) denotes that the dot z is inside a region y , but the first argument denotes the set of properties P that the region has, rather

than denoting y directly. However, whenever the full interpretation of these expressions in relation to a finite domain of graphical objects is required, they must be transformed into equivalent first-order expressions. This transformation is achieved through **meaning postulates**. The result of these transformations for the examples above are $curve_between_*(r_1, r_2) = x$ and $\exists y[region(y) \wedge inside_*(z, y)]$, where $curve_between_*$ and $inside_*$ denote geometrical functions whose arguments are graphical entities. The meaning postulates are defined as follows:

MP1. $\forall x \forall P[\delta(P)(x) \leftrightarrow P(\lambda y[\delta_*(x, y)])]$ where $\delta \in \{lie_at, be_in_zone, inside\}$

MP2. $\forall x \forall P[\delta(P)(x) \leftrightarrow P(\lambda y[\delta_*(y) = x])]$ where $\delta \in \{right\}$

MP3. $\forall x \forall P_1 \forall P_2[\delta(P_1)(P_2)(x) \leftrightarrow P_2(\lambda u[P_1(\lambda v[\delta_*(v, u) = x])])]$ where $\delta \in \{curve_between, intersection_between, line_from_to\}$

where $P, P_1,$ and P_2 are variables ranging over sets of properties (i.e., of type $\langle\langle e, t \rangle, t \rangle$), and $x, y, u,$ and v are variables ranging over individuals. Meaning postulate MP1 establishes, for instance, that a geometrical relation that holds between a set of properties of an individual a and an individual b stands in one-to-one correspondence with the relation that holds between the individuals a and b themselves, since the only property of a that is relevant for the geometrical interpretation process is the property of being in such a geometrical relation with the object b (i.e., that the object a lies at, is in a zone of, or is inside the object b). Similarly for meaning postulates MP2 and MP3.

The five examples that follow illustrate how the graphical interpreter works.

Example 1

Consider the interpretation of the expression $\Lambda(region)$ (BE (*big*)), which is the translation of *a country is big*. The expression can be reduced as follows:

1. $\lambda P \lambda Q \exists x[P(x) \wedge Q(x)]$ (*region*) ($\lambda P \lambda z P(z)$) (*big*)
2. $\lambda P \lambda Q \exists x[P(x) \wedge Q(x)]$ (*region*) ($\lambda z big(z)$)
3. $\lambda Q \exists x[region(x) \wedge Q(x)]$ ($\lambda z big(z)$)
4. $\exists x[region(x) \wedge \lambda z big(z)(x)]$
5. $\exists x[region(x) \wedge big(x)]$

Expression (5) is interpreted through the standard quantification rules of the geometrical interpreter without the help of meaning postulates. The interpretation of *big* is an algorithm that computes the average area of all regions in the map and returns the set of all regions whose area is larger than the average. This is a simple convention for illustrative purposes and alternative conventions could be chosen. Although the purpose of this paper is not to explore issues related to the interpretation of vague terms, it is interesting to note that within the present framework specific algorithms related to specific application domains that take into account the graphical context could be defined for the construction of practical applications.

Example 2

Consider the interpretation of THE (BETWEEN_a (R_1) (R_2) (*curve*))—which is the translation of *the border between France and Germany*, as will be shown in Section 2.4.1. The

expression without the abbreviations is:

1. $\lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)]$
 $(\lambda x_{\langle e,t \rangle} \lambda y_{\langle e,t \rangle} \lambda z_{\langle e,t \rangle} \lambda u_e [z(u) \wedge \text{curve_between}(x)(y)(u)]$
 $(\lambda P[P(r_1)])(\lambda P[P(r_2)])(\text{curve}))$

which can be reduced as follows:

2. $\lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)]$
 $(\lambda u [\text{curve}(u) \wedge \text{curve_between}(\lambda P[P(r_1)])(\lambda P[P(r_2)])(u)](u))$
3. $\lambda Q \exists y [\forall x [\lambda u [\text{curve}(u) \wedge \text{curve_between}(\lambda P[P(r_1)])(\lambda P[P(r_2)])(u)](x) \leftrightarrow$
 $x = y] \wedge Q(y)]$
4. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \text{curve_between}(\lambda P[P(r_1)])(\lambda P[P(r_2)])(x)) \leftrightarrow x = y] \wedge Q(y)]$
5. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \text{curve_between}(\lambda P[P(r_1)])(\lambda P[P(r_2)])(x)) \leftrightarrow x = y] \wedge Q(y)]$
6. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda P[P(r_2)])(\lambda u [\lambda P[P(r_1)])(\lambda v [\text{curve_between}_*(v, u) =$
 $x])]) \leftrightarrow x = y] \wedge Q(y)]$
7. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda P[P(r_2)])(\lambda u [\lambda v [\text{curve_between}_*(v, u) = x](r_1)]) \leftrightarrow x = y]$
 $\wedge Q(y)]$
8. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda P[P(r_2)])(\lambda u [\text{curve_between}_*(r_1, u) = x]) \leftrightarrow x = y]$
 $\wedge Q(y)]$
9. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda u [\text{curve_between}_*(r_1, u) = x](r_2)) \leftrightarrow x = y] \wedge Q(y)]$
10. $\lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \text{curve_between}_*(r_1, r_2) = x) \leftrightarrow x = y] \wedge Q(y)]$

Note that Expression (5) cannot be further reduced unless the types of the arguments of the predicate *curve_between* are lowered with the help of meaning postulate MP3. The geometrical functions in Expression (10) can be evaluated directly. Expression (10) is a denoting concept that refers to the curve between the regions r_1 and r_2 and cannot be further reduced. Consider that the expression *the border between France and Germany* is a definite description and, in order to obtain a truth value, must be combined with a predicate. The graphical object referred to by (10), on the other hand, could be identified regardless of the nature of the predicate Q , as this predicate is not used for picking out the object referred to by the definite description.¹¹ We call the object referred to by the denoting concept its **concrete extension**. The concrete extension of (10) can be identified, for instance, by interpreting the denoting concept without using the predicative abstraction Q (i.e., $\exists y [\forall x [(\text{curve}(x) \wedge \text{curve_between}_*(r_1, r_2) = x) \leftrightarrow x = y]]$) in relation to the graphical domain; if the denoting concept is indefinite, we take any object satisfying the expression as its concrete extension.

¹¹ As argued by Kaplan, contextual factors have to be considered for the identification of the referent of a definite description used referentially rather than attributively (Kaplan 1978). If the referent is identified deictically, as in the current example, the referent is found through the translation of the definite description into the graphical language, where the shape of the object is available directly. Note as well that as expressions of \mathbf{G} have an interpretation not only in relation to the graphical domain but also in relation to the world, through the translation into \mathbf{P} and the semantics of \mathbf{P} , the referent of a definite description in \mathbf{L} can be found by computing the geometrical interpretation of its translation into \mathbf{G} .

Example 3

Consider the interpretation of an expression similar to the one in Example 2, but in which an indefinite is included. The expression is $\text{THE}(\text{BETWEEN}_a(\text{R}_1)(\text{A}(\text{region}))(\text{curve}))$, which is the translation of *the border between France and a country*. The full expression is:

$$1. \quad \lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)] \\ (\lambda x_{\langle \langle e,t \rangle, t \rangle} \lambda y_{\langle \langle e,t \rangle, t \rangle} \lambda z_{\langle e,t \rangle} \lambda u_e [z(u) \wedge \text{curve_between}(x)(y)(u)] \\ (\lambda P [P(r_1)]) (\lambda P \exists z [\text{region}(z) \wedge P(z)]) \\ (\text{curve}))$$

the reduction is as follows:

$$2. \quad \lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)] \\ (\lambda u [\text{curve}(u) \wedge \text{curve_between}(\lambda P [P(r_1)]) (\lambda P \exists z [\text{region}(z) \wedge P(z)]) (u)]) \\ 3. \quad \lambda Q \exists y [\forall x [\lambda u [\text{curve}(u) \wedge \text{curve_between}(\lambda P [P(r_1)]) (\lambda P \exists z [\text{region}(z) \wedge \\ P(z)]) (u)] (x) \leftrightarrow x = y] \wedge Q(y)] \\ 4. \quad \lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \text{curve_between}(\lambda P [P(r_1)]) (\lambda P \exists z [\text{region}(z) \wedge P(z)]) (x)) \leftrightarrow \\ x = y] \wedge Q(y)] \\ 5. \quad \lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda P \exists z [\text{region}(z) \wedge \\ P(z)] (\lambda u [\lambda P [P(r_1)] (\lambda v [\text{curve_between}_*(v, u) = x])]) \leftrightarrow x = y] \wedge Q(y)] \\ 6. \quad \lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda P \exists z [\text{region}(z) \wedge P(z)] (\lambda u [\lambda v [\text{curve_between}_*(v, u) = \\ x] (r_1)]) \leftrightarrow x = y] \wedge Q(y)] \\ 7. \quad \lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \lambda P \exists z [\text{region}(z) \wedge P(z)] (\lambda u [\text{curve_between}_*(r_1, u) = \\ x]) \leftrightarrow x = y] \wedge Q(y)] \\ 8. \quad \lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \exists z [\text{region}(z) \wedge \lambda u [(\text{curve_between}_*(r_1, u) = x] (z))] \leftrightarrow x = \\ y] \wedge Q(y)] \\ 9. \quad \lambda Q \exists y [\forall x [(\text{curve}(x) \wedge \exists z [\text{region}(z) \wedge \text{curve_between}_*(r_1, z) = x] \leftrightarrow x = y] \wedge Q(y)]$$

Meaning postulate MP3 is used for reducing from (4) to (5). Expression (9) is a denoting concept similar to the final expression in Example 2, but one which has an embedded quantified expression. Meaning postulates MP1 to MP3 are defined in such a way that terms preserve quantificational properties through the reduction process.

Example 4

Consider the expression $\text{R}_2(\text{be_in_zone}(\text{THE}(\text{OF}_b(\text{R}_1)(\text{right})))$ —which is the translation of *Germany is to the east of France*. The reduced expression is the following:

$$1. \quad \text{be_in_zone}(\lambda Q \exists y [\forall x [\text{right}(\lambda P [P(r_1)]) (x) \leftrightarrow x = y] \wedge Q(y)])(r_2)$$

by meaning postulate MP2:

$$2. \quad \text{be_in_zone}(\lambda Q \exists y [\forall x [\lambda P [P(r_1)] (\lambda z [\text{right}_*(z) = x] \leftrightarrow x = y] \wedge Q(y)])(r_2) \\ 3. \quad \text{be_in_zone}(\lambda Q \exists y [\forall x [\lambda z [\text{right}_*(z) = x] (r_1) \leftrightarrow x = y] \wedge Q(y)])(r_2) \\ 4. \quad \text{be_in_zone}(\lambda Q \exists y [\forall x [\text{right}_*(r_1) = x \leftrightarrow x = y] \wedge Q(y)])(r_2)$$

by meaning postulate MP1:

5. $\lambda Q \exists y [\forall x [\text{right}_*(r_1) = x \leftrightarrow x = y] \wedge Q(y)] (\lambda z [\text{be_in_zone}_*(r_2, z)])$
6. $\exists y [\forall x [\text{right}_*(r_1) = x \leftrightarrow x = y] \wedge \lambda z [\text{be_in_zone}_*(r_2, z)](y)]$
7. $\exists y [\forall x [\text{right}_*(r_1) = x \leftrightarrow x = y] \wedge \text{be_in_zone}_*(r_2, y)]$.

Expression (7) is a first-order formula that can be directly evaluated by the interpreter of **G**. The operator right_* is interpreted as a geometrical algorithm that computes the centroid (x_c, y_c) of a region r and returns the semiplane to the right of the centroid of r (i.e., the set of all ordered pairs of reals $\langle x_i, y_i \rangle$ such that $x_i > x_c$). This convention captures objects that are to the right of a region, or those in the right part of a region.¹² The graphical predicate be_in_zone_* checks whether r_2 is within y —i.e., the zone to the right of r_1 .

Example 5

Consider the interpretation of the translation into **G** of the textual part of the multimodal message in Figure 2. The translation of *Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt* is shown in (1), its reduction in (2), and its final reduction applying the meaning postulates in (3):

1. $D_3 (\text{lie_at} (\text{THE} (\text{BETWEEN}_b (\text{THE} (\text{BETWEEN}_a (R_1) (R_3) (\text{curve}))) (\text{A} (\text{FROM_TO}(D_1) (D_3) (\text{line}))) (\text{intersection})))))$
2. $\text{lie_at} (\lambda Q \exists y [\forall x [\text{intersection}(x) \wedge \text{intersection_between} (\lambda Q \exists u [\forall v [(\text{curve}(v) \wedge \text{curve_between} (\lambda P [P(r_1)]) (\lambda P [P(r_2)]) (v)) \leftrightarrow v = u] \wedge Q(y)]) (\lambda Q \exists z [\text{line}(z) \wedge \text{line_from_to} (\lambda P [P(d_1)]) (\lambda P [P(d_3)]) (z) \wedge Q(z)]) = x \leftrightarrow x = y] \wedge Q(y)]) (D_3)$
3. $\exists y [\forall x [(\text{intersection}(x) \wedge \exists z [\text{line}(z) \wedge \text{line_from_to}_*(d_1, d_3) = z] \wedge \exists u [\forall v [(\text{curve}(v) \wedge \text{curve_between}_*(r_1, r_2) = v) \leftrightarrow v = u] \wedge \text{intersection_between}_*(u, z) = x]) \leftrightarrow x = y] \wedge \text{lie_at}_*(d_2, y)]$.

Expression (3) is true if the position of dot d_3 is the same as the position of the intersection between the curve between r_1 and r_2 and the line from d_1 to d_3 , as is the case in Figure 2.

It is worth emphasizing that as the five examples illustrate, the reason for type-raising graphical terms is to be able to translate natural language quantified expression into the graphical domain compositionally in a rather elegant way. The scheme provides a clear specification strategy; however, in a practical implementation, it would

¹² This is an arbitrary convention defined for illustrative purposes and alternative conventions could be chosen. Similar conventions could be used to interpret whether other kinds of graphical objects stand in a right-of relation. Furthermore, several conventions for the interpretation of such words can be used and a particular geometric algorithm can be defined for each interpretation. These algorithms need not be fully quantitative; more qualitative approaches can be employed as long as the computation returns a semantic value of an appropriate kind.

constant of L: α	Category name	Category definition	Translation into G: $\rho_{L-G}(\alpha)$	Corresponding type in G
<i>Paris</i>	T	t/IV	$\lambda P[P(d_1)]$	$\langle\langle e, t \rangle, t\rangle$
<i>Frankfurt</i>	T	t/IV	$\lambda P[P(d_3)]$	$\langle\langle e, t \rangle, t\rangle$
<i>Saarbrücken</i>	T	t/IV	$\lambda P[P(d_2)]$	$\langle\langle e, t \rangle, t\rangle$
<i>France</i>	T	t/IV	$\lambda P[P(r_1)]$	$\langle\langle e, t \rangle, t\rangle$
<i>Germany</i>	T	t/IV	$\lambda P[P(r_2)]$	$\langle\langle e, t \rangle, t\rangle$
<i>city</i>	CN	CN	<i>dot</i>	$\langle e, t \rangle$
<i>country</i>	CN	CN	<i>region</i>	$\langle e, t \rangle$
<i>border</i>	CN	CN	<i>curve</i>	$\langle e, t \rangle$
<i>line</i>	CN	CN	<i>line</i>	$\langle e, t \rangle$
<i>intersection</i>	CN	CN	<i>intersection</i>	$\langle e, t \rangle$
<i>east</i>	CN'	CN'	<i>right</i>	$\langle\langle\langle e, t \rangle, t \rangle, e\rangle$
<i>big</i>	ADJ	ADJ	<i>big</i>	$\langle e, t \rangle$
<i>be</i>	TV	IV/(t/IV)	$\lambda P \lambda x P(\lambda y [x = y])$	$\langle\langle\langle e, t \rangle, t \rangle, \langle e, t \rangle\rangle$
<i>be</i>	IV/ADJ	IV/ADJ	$\lambda P \lambda x P(x)$	$\langle\langle e, t \rangle, \langle e, t \rangle\rangle$
<i>lie at</i>	TV	IV/(t/IV)	<i>lie_at</i>	$\langle\langle\langle e, t \rangle, t \rangle, \langle e, t \rangle\rangle$
<i>be to</i>	TV	IV/(t/IV)	<i>be_in_zone</i>	$\langle\langle\langle e, t \rangle, t \rangle, \langle e, t \rangle\rangle$
<i>a</i>	T/CN	(t/IV)/CN	$\lambda P \lambda Q \exists x [P(x) \wedge Q(x)]$	$\langle\langle e, t \rangle, \langle\langle e, t \rangle, t \rangle\rangle$
<i>the</i>	T/CN	(t/IV)/CN	$\lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)]$	$\langle\langle\langle e, t \rangle, \langle\langle e, t \rangle, t \rangle\rangle\rangle$

Figure 12
Translation of constants of L into G.

be convenient to limit the expressive power of G and to define it as a first-order language.

2.4 Translations between L and G

In this section, the translation functions ρ_{L-G} and ρ_{G-L} are defined. As discussed in Section 1, the goal in interpreting a multimodal message like the one in Figure 2 is to find the translations of individual constants, which are not known. In this section, however, we assume that the translation is fully defined in order to illustrate all theoretical elements of the scheme in Figure 3. The induction of the translation of individual constants, on the other hand, will be shown in Section 3.

For each syntactic category of L there is a corresponding type in G. The correspondence between linguistic categories and geometrical types resembles the translation from English to intensional logic (Dowty, Wall, and Peters 1985) and is defined in terms of the function f as follows:

1. $f(t) = t$.
2. $f(CN) = f(IV) = f(ADJ) = \langle e, t \rangle$.
3. For any categories A and B , $f(A/B) = \langle f(B), f(A) \rangle$.

2.4.1 Translation from L into G. Figure 12 shows the translation of constants of L. Simple terms, such as the names of cities and countries, translate into expressions denoting characteristic functions of sets of graphical entities. This graphical type is interpreted as the set of “properties” that an individual named by the term has (for the purpose of this discussion a property is just the set of individuals, as no intensional types are considered). So, as a city is represented by a dot in the graphical domain, the translation of *Paris*, for instance, is the set of geometrical properties that the dot representing Paris has in the interpretation state. Common nouns of category CN and CN' translate into predicates and functions from sets of properties to individuals, respectively. Adjectives occurring in attributive sentences are translated as sets of individuals. Note that there are two constants *be*: one combines with a term and

the other with an adjective and both combinations produce intransitive verbs. The translations corresponding to these constants are functions from sets of properties to sets of individuals, and from sets of individuals to sets of individuals, respectively. Transitive verbs like *lie_at* and *be_to* translate into geometrical operators whose type is a function from sets of properties to sets of individuals. Determiners are translated in a standard fashion.

The translation rules for composite expressions are as follows:

SENTENCES

T1_{L-G}. If $\alpha \in P_T$ and $\beta \in P_{IV}$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then $\rho_{L-G}(F_{L1}(\alpha, \beta)) = \alpha'(\beta')$, that is to say, the function α' applied to the argument β' .

Examples: $\rho_{L-G}(\textit{Paris is a city of France}) = D_1 (\text{BE}_a (\text{A} (\text{OF}_a(\text{R}_1) (\textit{dot}))))$
 $\rho_{L-G}(\textit{Germany is to the east of France}) =$
 $R_3 (\textit{be_in_zone} (\text{THE} (\text{OF}_b(\text{R}_1) (\textit{right}))))$
 $\rho_{L-G}(\textit{a country is big}) = \text{A}(\textit{region}) (\text{BE}_b(\textit{big}))$
 $\rho_{L-G}(\textit{Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt}) =$
 $D_3 (\textit{lie_at} (\text{THE} (\text{BETWEEN}_b (\text{THE} (\text{BETWEEN}_a (\text{R}_1) (\text{R}_3)(\textit{curve})))$
 $(\text{A} (\text{FROM_TO}(\text{D}_1) (\text{D}_3) (\textit{line})))$
 $(\textit{intersection}))))$

TRANSITIVE VERB PHRASES

T2_{L-G}. If $\alpha \in P_{TV}$ and $\beta \in P_T$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then $\rho_{L-G}(F_{L2}(\alpha, \beta)) = \alpha'(\beta')$.

Examples: $\rho_{L-G}(\textit{be a city}) = \text{BE}_a (\text{A} (\textit{dot}))$
 $\rho_{L-G}(\textit{be to the east of France}) = \textit{be_in_zone} (\text{THE} (\text{OF}_b(\text{R}_1)(\textit{right})))$

ATTRIBUTIVE VERB PHRASES

T3_{L-G}. If $\alpha \in P_{IV/ADJ}$ and $\beta \in P_{ADJ}$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then $\rho_{L-G}(F_{L2}(\alpha, \beta)) = \alpha'(\beta')$.

Example: $\rho_{L-G}(\textit{be big}) = \text{BE}_b(\textit{big})$

TERMS

T4_{L-G}. If $\alpha \in P_{T/CN}$ and $\beta \in P_{CN}$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then $\rho_{L-G}(F_{L3}(\alpha, \beta)) = \alpha'(\beta')$.

Examples: $\rho_{L-G}(\textit{a city}) = \text{A} (\textit{dot})$
 $\rho_{L-G}(\textit{a city of France}) = \text{A} (\text{OF}_a(\text{R}_1)(\textit{dot}))$
 $\rho_{L-G}(\textit{the border between France and Germany}) = \text{THE} (\text{BETWEEN}_a$
 $(\text{R}_1) (\text{R}_2) (\textit{curve}))$
 $\rho_{L-G}(\textit{a line from Paris to Frankfurt}) = \text{A} (\text{FROM_TO}(\text{D}_1) (\text{D}_3) (\textit{line}))$
 $\rho_{L-G}(\textit{the east of France}) = \text{THE} (\text{OF}_b(\text{R}_1)(\textit{right}))$

Note that the term *the east* can be formed by the rule $S4_L$, but it cannot be translated into \mathbf{G} because there is a type restriction in the definition of $T4_{L-G}$ (i.e., $\beta \in P_{CN}$, but $east \in P_{CN'}$). This restriction prevents the translation of terms like *the east* as these expressions have no concrete graphical representation; however, *the east of France* can be generated, translated into \mathbf{G} and interpreted through the geometry as shown in Section 2.3.2. In general, natural language expressions denoting abstract concepts do not have a graphical representation (i.e., *the population of France*), and although in this grammar we have focused on expressions that can be translated into \mathbf{G} , the language can be extended with linguistic terms that would be interpreted only in the linguistic modality.

COMMON NOUNS

$T5_{L-G}$. If $\alpha \in P_{CN}$ and $\beta \in P_{PP}$, or $\alpha \in P_{CN'}$ and $\beta \in P_{PP'}$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then $\rho_{L-G}(F_{L2}(\alpha, \beta)) = \beta'(\alpha')$.

Examples: $\rho_{L-G}(\text{city of France}) = OF_a(R_1)(\text{dot})$
 $\rho_{L-G}(\text{east of France}) = OF_b(R_1)(\text{right})$
 $\rho_{L-G}(\text{border between France and Germany}) = BETWEEN_a(R_1)(R_2)$
 (curve)
 $\rho_{L-G}(\text{intersection between the border between France and Germany and a line from Paris to Frankfurt}) = BETWEEN_b(\text{THE}(\text{BETWEEN}_a(R_1)(R_2)(\text{curve}))) (\text{A}(\text{FROM_TO}(D_1)(D_3)(\text{line}))) (\text{intersection})$

of PREPOSITIONAL PHRASES

$T6_{L-G}$. If $\alpha \in P_T$, and $\rho_{L-G}(\alpha) = \alpha'$, then $\rho_{L-G}(F_{L4}(\alpha))$ is either $OF_a(\alpha')$ or $OF_b(\alpha')$.

Examples: $\rho_{L-G}(\text{of France}) = OF_a(R_1)$
 $\rho_{L-G}(\text{of Germany}) = OF_b(R_2)$

between PREPOSITIONAL PHRASES

$T7_{L-G}$. If $\alpha, \beta \in P_T$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then

Examples: $\rho_{L-G}(F_{L5}(\alpha, \beta))$ is either $BETWEEN_a(\alpha')(\beta')$ or $BETWEEN_b(\alpha')(\beta')$
 $\rho_{L-G}(\text{between France and Germany}) = BETWEEN_a(R_1)(R_2)$
 $\rho_{L-G}(\text{between the border between France and Germany and a line from Paris to Frankfurt}) = BETWEEN_b(\text{THE}(\text{BETWEEN}_a(R_1)(R_2)(\text{curve}))) (\text{A}(\text{FROM_TO}(D_1)(D_3)(\text{line})))$

From-to PREPOSITIONAL PHRASES

$T8_{L-G}$. If $\alpha, \beta \in P_T$, and $\rho_{L-G}(\alpha) = \alpha'$, $\rho_{L-G}(\beta) = \beta'$ then $\rho_{L-G}(F_{L6}(\alpha, \beta)) = FROM_TO(\alpha')(\beta')$.

Example: $\rho_{L-G}(\text{from Paris to Frankfurt}) = FROM_TO(D_1)(D_3)$

Constant of G:	Translation into L:
α	$\rho_{G-L}(\alpha)$
<i>dot</i>	<i>city</i>
<i>region</i>	<i>country</i>
<i>curve</i>	<i>border</i>
<i>line</i>	<i>line</i>
<i>intersection</i>	<i>intersection</i>
<i>right</i>	<i>east</i>
<i>big</i>	<i>big</i>
<i>lie_at</i>	<i>lie at</i>
<i>be_in_zone</i>	<i>be to</i>

Figure 13
Translation of constants of language G into L.

2.4.2 Translation from G into L. In this section, the translation function ρ_{G-L} is defined. The translation of expressions of G into L are shown in Figures 13 and 14. Note that constants of G in Figure 13 translate into constants of L; however, the translations shown in Figure 14 are more complex, since composite expressions of G can translate into basic or composite expressions of L.

The translation from G into L is shown below. In rules T6_{G-L} to T8_{G-L} Q stands for either the quantifier A or THE.

SENTENCES

T1_{G-L}. If $\alpha \in E_{\langle(e,t),t\rangle}$ and $\beta \in E_{\langle e,t\rangle}$, and $\rho_{G-L}(\alpha) = \alpha'$, $\rho_{G-L}(\beta) = \beta'$ then $\rho_{G-L}(F_{G1}(\alpha, \beta)) = \alpha' \beta''$ (the concatenation), where β'' is the result of replacing the first *verb* in β' with its third person singular present form.

Examples: $\rho_{G-L}(D_1(BE_a(A(OF_a(R_1)(dot)))))) = Paris\ is\ a\ city\ of\ France$
 $\rho_{G-L}(R_3(be_in_zone(THE(OF_b(R_1)(right)))))) = Germany\ is\ to\ the\ east\ of\ France$
 $\rho_{G-L}(A(region)(BE_b(big))) = a\ country\ is\ big$
 $\rho_{G-L}(D_3(lie_at(THE(BETWEEN_b(THE(BETWEEN_a(R_1)(R_3)(curve))))(A(FROM_TO(D_1)(D_3)(line))))(intersection)))) =$
Saarbrücken lies at the intersection between the border between France and Germany and a line from Paris to Frankfurt

Expression of G: δ	Translation into L: $\rho_{G-L}(\delta)$
$\lambda P[P(d_1)], \lambda P[P(d_2)], \lambda P[P(d_3)]$	<i>Paris, Frankfurt, Saarbrücken, respectively</i>
$\lambda P[P(r_1)], \lambda P[P(r_2)]$	<i>France, Germany, respectively</i>
$\lambda P[P(c_1)]$	<i>the border between France and Germany</i>
$\lambda P \lambda x P(\lambda y[x = y])$	<i>be</i>
$\lambda P \lambda x P(x)$	<i>be</i>
$\lambda P \lambda Q \exists x [P(x) \wedge Q(x)]$	<i>a</i>
$\lambda P \lambda Q \exists y [\forall x [P(x) \leftrightarrow x = y] \wedge Q(y)]$	<i>the</i>

Figure 14
Translation of some composite expressions of G into constants of L.

TRANSITIVE VERB PHRASES

T2_{G-L}. If $\alpha \in E_{\langle\langle e,t \rangle, t \rangle, \langle e,t \rangle}$ and $\beta \in E_{\langle\langle e,t \rangle, t \rangle}$, and $\rho_{G-L}(\alpha) = \alpha'$, $\rho_{G-L}(\beta) = \beta'$ then $\rho_{G-L}(F_{G1}(\alpha, \beta)) = \alpha' \beta'$.

Examples: $\rho_{G-L}(\text{BE}_a(\text{A}(\text{dot}))) = \text{be a city}$
 $\rho_{G-L}(\text{be_in_zone}(\text{THE}(\text{OF}_b(\text{R}_1)(\text{right})))) = \text{be to the east of France}$

ATTRIBUTIVE VERB PHRASES

T3_{G-L}. If $\alpha \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$ and $\beta \in E_{\langle e,t \rangle}$, and $\rho_{G-L}(\alpha) = \alpha'$, $\rho_{G-L}(\beta) = \beta'$ then $\rho_{G-L}(F_{G1}(\alpha, \beta)) = \alpha' \beta'$.

Example: $\rho_{G-L}(\text{BE}_b(\text{big})) = \text{be big}$

TERMS

T4_{G-L}. If $\alpha \in E_{\langle\langle e,t \rangle, \langle\langle e,t \rangle, t \rangle\rangle}$, $\beta \in E_{\langle e,t \rangle}$, and $\rho_{G-L}(\alpha) = \alpha'$, $\rho_{G-L}(\beta) = \beta'$ then $\rho_{G-L}(F_{G1}(\alpha, \beta)) = \alpha'' \beta'$, where α'' is α' except in the case where α' is *a* and the first word in β begins with a vowel; here, α'' is *an*.

Examples: $\rho_{G-L}(\text{A}(\text{dot})) = \text{a city}$
 $\rho_{G-L}(\text{A}(\text{OF}_a(\text{R}_1)(\text{dot}))) = \text{a city of France}$
 $\rho_{G-L}(\text{THE}(\text{BETWEEN}_a(\text{R}_1)(\text{R}_2)(\text{curve}))) = \text{the border between France and Germany}$
 $\rho_{G-L}(\text{A}(\text{FROM_TO}(\text{D}_1)(\text{D}_3)(\text{line}))) = \text{a line from Paris to Frankfurt}$
 $\rho_{G-L}(\text{THE}(\text{OF}_b(\text{R}_1)(\text{right}))) = \text{the east of France}$

COMMON NOUNS

T5_{G-L}. If $\alpha \in E_{\langle\langle e,t \rangle, \langle e,t \rangle\rangle}$ and $\beta \in E_{\langle e,t \rangle}$, or $\alpha \in E_{\langle\langle\langle\langle e,t \rangle, t \rangle, e \rangle, \langle e,t \rangle\rangle}$ and $\beta \in E_{\langle\langle\langle e,t \rangle, t \rangle, e \rangle}$, and $\rho_{G-L}(\alpha) = \alpha'$, $\rho_{G-L}(\beta) = \beta'$ then $\rho_{G-L}(F_{G1}(\alpha, \beta)) = \beta' \alpha'$.

Examples: $\rho_{G-L}(\text{OF}_a(\text{R}_1)(\text{dot})) = \text{city of France}$
 $\rho_{G-L}(\text{OF}_b(\text{R}_1)(\text{right})) = \text{east of France}$
 $\rho_{G-L}(\text{BETWEEN}_a(\text{R}_1)(\text{R}_2)(\text{curve})) = \text{border between France and Germany}$
 $\rho_{G-L}(\text{BETWEEN}_b(\text{THE}(\text{BETWEEN}_a(\text{R}_1)(\text{R}_2)(\text{curve})))$
 $(\text{A}(\text{FROM_TO}(\text{D}_1)(\text{D}_3)(\text{line}))) (\text{intersection})) =$
intersection between the border between France and Germany and a line from Paris to Frankfurt

of PREPOSITIONAL PHRASES

T6_{G-L}. If $\alpha \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α is either R_i or $Q(\text{region})$ and $\rho_{G-L}(\alpha) = \alpha'$ then $\rho_{G-L}(F_{G2}(\alpha)) = \rho_{G-L}(F_{G3}(\alpha)) = \text{of } \alpha'$

Examples: $\rho_{G-L}(\text{OF}_a(\text{R}_1)) = \text{of France}$
 $\rho_{G-L}(\text{OF}_b(\text{R}_2)) = \text{of Germany}$
 $\rho_{G-L}(\text{OF}_a(\text{A}(\text{region}))) = \text{of a country}$

between PREPOSITIONAL PHRASES

- T7_{G-L}. If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that
- (a) α, β are either R_i or $Q(\textit{region})$ or
 - (b) α is either C_i or $Q(\textit{curve})$ and β is either L_i or $Q(\textit{line})$,
- and $\rho_{G-L}(\alpha) = \alpha', \rho_{G-L}(\beta) = \beta'$ then
- $$\rho_{G-L}(F_{G4}(\alpha, \beta)) = \textit{between } \alpha' \textit{ and } \beta'.$$

Examples: $\rho_{G-L}(\text{BETWEEN}_a(R_1)(R_2)) = \textit{between France and Germany}$
 $\rho_{G-L}(\text{BETWEEN}_a(R_1)(\text{OF}_a(A(\textit{region})))) = \textit{between France and a country}$
 $\rho_{G-L}(\text{BETWEEN}_b(\text{THE}(\text{BETWEEN}_a(R_1)(R_2)(\textit{curve})))$
 $(A(\text{FROM_TO}(D_1)(D_3)(\textit{line})))) =$
between the border between France and Germany and a line from Paris to Frankfurt

from-to PREPOSITIONAL PHRASES

- T8_{G-L}. If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α, β are either D_i or $Q(\textit{dot})$ and
- $$\rho_{G-L}(\alpha) = \alpha', \rho_{G-L}(\beta) = \beta' \text{ then } \rho_{G-L}(F_{G5}(\alpha, \beta)) = \textit{from } \alpha' \textit{ to } \beta'.$$

Example: $\rho_{G-L}(\text{FROM_TO}(D_1)(D_3)) = \textit{from Paris to Frankfurt}$

As mentioned above, **G** is a very expressive language; not all expressions of **G** can be translated into expressions of **L**. Rules T1_{G-L} to T8_{G-L} define the expressions that do have a translation. Instances of expressions that cannot be translated are individual constants (e.g., d_1), equality relations between individuals (e.g., $d_1 = d_2$), and conjunctions or disjunctions (e.g., $\textit{dot}(d_1) \wedge \textit{dot}(d_2)$). Other examples are expressions of the form $\lambda P[P(e_1) \vee P(e_2) \vee \dots \vee P(e_n)]$, where e_i is an individual constant, which denote the set of properties that one or another individual has. However, this latter kind of expression could be translated if the expressiveness of **L** were augmented by allowing conjoined term phrases in the grammar.

2.5 Translations between **G** and **P**

The translation functions ρ_{G-P} and ρ_{P-G} are defined in this section, concluding the presentation of the theoretical elements of the system of multimodal representation. For each type of **P** there is a corresponding type in **G** and it is defined in terms of the function f_{P-G} as follows:

1. $f_{P-G}(\textit{dot}) = f_{P-G}(\textit{line}) = f_{P-G}(\textit{curve}) = f_{P-G}(\textit{region}) = f_{P-G}(\textit{zone}) =$
 $f_{P-G}(\textit{composite_region}) = f_{P-G}(\textit{dot_set}) = f_{P-G}(\textit{line_set}) = f_{P-G}(\textit{map}) = \langle e, t \rangle.$
2. For any types a and b , $f_{P-G}(\langle a, b \rangle) = \langle f_{P-G}(a), f_{P-G}(b) \rangle.$

2.5.1 Translation from **P into **G**.** The translations of the constants of **P** into **G** are presented in Figure 15. In the following definitions, **Q** stands for either the quantifier

Constant of P:	Translation into G:
α	$\rho_{P-G}(\alpha)$
d_1, d_2, d_3, \dots	$\lambda P[P(d_1)], \lambda P[P(d_2)], \lambda P[P(d_3)], \dots$
l_1, l_2, l_3, \dots	$\lambda P[P(l_1)], \lambda P[P(l_2)], \lambda P[P(l_3)], \dots$
c_1, c_2, c_3, \dots	$\lambda P[P(c_1)], \lambda P[P(c_2)], \lambda P[P(c_3)], \dots$
r_1, r_2, r_3, \dots	$\lambda P[P(r_1)], \lambda P[P(r_2)], \lambda P[P(r_3)], \dots$
z_1, z_2, z_3, \dots	$\lambda P[P(z_1)], \lambda P[P(z_2)], \lambda P[P(z_3)], \dots$

Figure 15
Translation of constants of language P into G.

A or THE. The translation rules are as follows:

CONSTANT¹³

T1_{P-G}. If $\alpha \in C_s$ where $s \in \{dot, line, curve, region, zone\}$ then
 (a) $\rho_{P-G}(\alpha)$ is as shown in Figure 15.
 (b) $\rho_{P-G}(\beta) = Q(s)$.

Examples: $\rho_{P-G}(\bullet) = \lambda P[P(d_1)]$
 $\rho_{P-G}(/) = \lambda P[P(l_1)]$
 $\rho_{P-G}(\text{rectangle}) = \lambda P[P(r_1)]$
 $\rho_{P-G}(\text{region}) = A(\text{region})$

LINE

T2_{P-G}. If $\alpha, \beta \in E_{dot}$, and $\rho_{P-G}(\alpha) = \alpha'$ and $\rho_{P-G}(\beta) = \beta'$ then
 $\rho_{P-G}(F_{P1}(\alpha, \beta)) = Q(\text{FROM_TO}(\alpha')(\beta')(line))$.

Example: $\rho_{P-G}(\text{arrow}) = A(\text{FROM_TO}(\lambda P[P(d_1)])(\lambda P[P(d_3)])(line))$

CURVE¹⁴

T3_{P-G}. If $\alpha, \beta \in E_{region}$ such that α and β are adjacent, and $\rho_{P-G}(\alpha) = \alpha'$ and $\rho_{P-G}(\beta) = \beta'$ then $\rho_{P-G}(F_{P2}(\alpha, \beta)) = Q(\text{BETWEEN}_a(\alpha')(\beta')(curve))$.

Examples¹⁴: $\rho_{P-G}(\text{adjacent regions}) = \text{THE}(\text{BETWEEN}_a(\lambda P[P(r_1)])(\lambda P[P(r_2)])(curve))$
 $\rho_{P-G}(\text{region with curve}) = \text{THE}(\text{BETWEEN}_a(\lambda P[P(r_1)])(A(\text{region}))(curve))$

INTERSECTION

T4_{P-G}. If $\alpha \in E_{curve}$ and $\beta \in E_{line}$, and $\rho_{P-G}(\alpha) = \alpha'$ and $\rho_{P-G}(\beta) = \beta'$ then
 $\rho_{P-G}(F_{P3}(\alpha, \beta)) = Q(\text{BETWEEN}_b(\alpha')(\beta')(intersection))$.

13 Rule (b) allows the concrete extension of a graphical object in P to be represented as its corresponding denoting concept in G.

14 These two example expressions correspond to the abbreviated expressions in Examples 2 and 3, respectively, presented in Section 2.3.2.

Example: $\rho_{P-G}(\text{img}) = \text{THE}(\text{BETWEEN}_b(\text{THE}(\text{BETWEEN}_a(\lambda P[P(r_1)])(\lambda P[P(r_2)])(\text{curve}))(A(\text{FROM_TO}(\lambda P[P(d_1)])(\lambda P[P(d_3)])(\text{line})))$
(*intersection*))

RIGHT

T5_{P-G}. If $\alpha \in E_{region}$ and $\rho_{P-G}(\alpha) = \alpha'$ then $\rho_{P-G}(F_{P4}(\alpha)) = Q(\text{OF}_b(\text{right})(\alpha'))$.

Example: $\rho_{P-G}(\text{img}) = \text{THE}(\text{OF}_b(\text{right})(\lambda P[P(r_3)]))$

DOT INSIDE A REGION

T6_{P-G}. If $\alpha \in E_{region}$ and $\rho_{P-G}(\alpha) = \alpha'$ then $\rho_{P-G}(F_{P5}(\alpha)) = Q(\text{OF}_a(\alpha')(\text{dot}))$.

Example: $\rho_{P-G}(\text{img}) = A(\text{OF}_a(\lambda P[P(r_1)])(\text{dot}))$

COMPOSITE REGION (1)¹⁵

T7_{P-G}. If $\alpha, \beta \in C_{region}$ such that α and β are adjacent, and $\rho_{P-G}(\alpha) = \lambda P[\alpha']$ and $\rho_{P-G}(\beta) = \lambda P[\beta']$ then $\rho_{P-G}(F_{P6}(\alpha, \beta)) = \lambda P[\alpha' \vee \beta']$.

COMPOSITE REGION (2)

T8_{P-G}. If $\alpha \in C_{region}$ and $\beta \in E_{composite_region}$ such that α and β are adjacent, and $\rho_{P-G}(\alpha) = \lambda P[\alpha']$ and $\rho_{P-G}(\beta) = \lambda P[\beta']$ then $\rho_{P-G}(F_{P6}(\alpha, \beta)) = \lambda P[\alpha' \vee \beta']$.

SET OF DOTS

T9_{P-G}. (a) If $\alpha \in E_{dot_set} = \emptyset$ and $\beta \in C_{dot}$, and $\rho_{P-G}(\beta) = \lambda P[\beta']$ then $\rho_{P-G}(F_{P6}(\alpha, \beta)) = \lambda P[\beta']$.

(b) If $\alpha \in E_{dot_set} \neq \emptyset$ and $\beta \in C_{dot}$, and $\rho_{P-G}(\alpha) = \lambda P[\alpha']$ and $\rho_{P-G}(\beta) = \lambda P[\beta']$ then $\rho_{P-G}(F_{P6}(\alpha, \beta)) = \lambda P[\alpha' \vee \beta']$.

SET OF LINES

T10_{P-G}. (a) If $\alpha \in E_{line_set} = \emptyset$ and $\beta \in C_{line}$, and $\rho_{P-G}(\beta) = \lambda P[\beta']$ then $\rho_{P-G}(F_{P6}(\alpha, \beta)) = \lambda P[\beta']$.

(b) If $\alpha \in E_{line_set} \neq \emptyset$ and $\beta \in C_{line}$, and $\rho_{P-G}(\alpha) = \lambda P[\alpha']$ and $\rho_{P-G}(\beta) = \lambda P[\beta']$ then $\rho_{P-G}(F_{P6}(\alpha, \beta)) = \lambda P[\alpha' \vee \beta']$.

¹⁵ Examples of the application of the rules T7_{P-G} to T11_{P-G} are included in the translation of a map shown in Figure 16, as explained below.

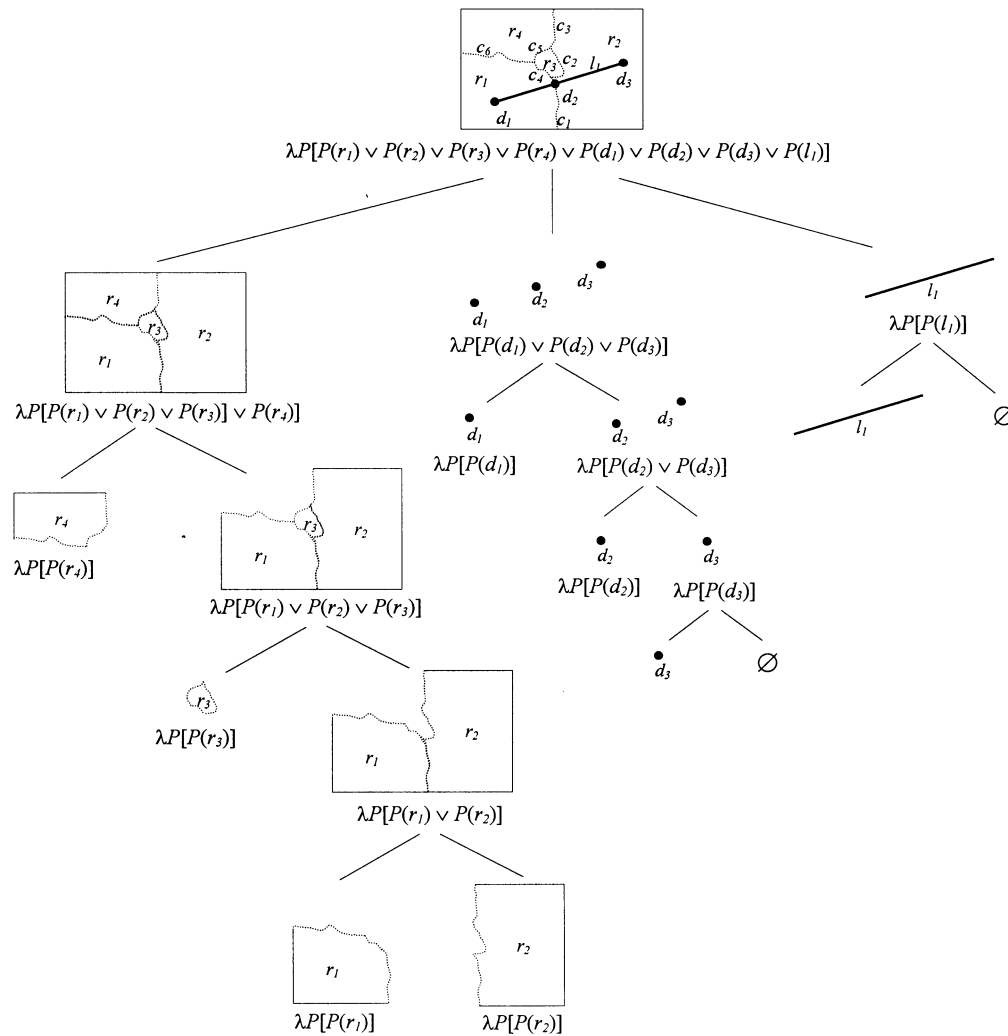


Figure 16
Translation into \mathbf{G} of a map.

MAP

$T11_{P-G}$. If $\alpha \in E_{composite_region}$, $\beta \in E_{dot_set}$ and $\delta \in E_{line_set}$, and
 $\rho_{P-G}(\alpha) = \lambda P[\alpha']$, $\rho_{P-G}(\beta) = \lambda P[\beta']$, $\rho_{P-G}(\delta) = \lambda P[\delta']$, then
 $\rho_{P-G}(F_{P7}(\alpha, \beta, \delta)) = \lambda P[\alpha' \vee \beta' \vee \delta']$.

An example of the translation of a map from \mathbf{P} into \mathbf{G} by rule $T11_{P-G}$ is shown in Figure 16. A map is interpreted in \mathbf{G} as the set of properties that one or another graphical object in the base of the map has. Computing and translating all possible syntactic structures that can be generated in \mathbf{P} on the basis of the overt graphical symbols of the drawing is not required for the interpretation of the picture in Figure 4. The translation rules permit mapping a large number of syntactic structures into \mathbf{G} , and they can be used as necessary. However, for the interpretation of a map we will only translate a designated expression ζ of type *map* that results from parsing a full drawing in terms of the graphical objects in the base. ζ

will be called the map. This criterion ensures that the drawing belongs to the map modality \mathbf{P} . In addition, the graphical terms in the disjunction of the body of expressions of type *map* are used in \mathbf{G} to define the interpretation domain \mathbf{P}_{base} . When this set is defined the semantic rules to interpret expressions of \mathbf{G} can be evaluated.

2.5.2 Translation from \mathbf{G} into \mathbf{P} . As mentioned in Section 1 in relation to the scheme in Figure 3, the purpose of this translation is to draw the graphical symbols that are referred to in \mathbf{G} . To picture the full map, the only symbols that must be drawn are the symbols of the base (\mathbf{P}_{base}), as emerging symbols do not have an independent pictorial realization. Thus, the only translations that have to be defined are the translations of the symbols contained in the expression ζ (i.e., the map). We also have to consider that graphical terms occurring in expressions of \mathbf{G} can have a graphical realization, which may be required for specific purposes. For instance, if one needs to highlight the region to the east of France the term of \mathbf{G} denoting that region should be translated and depicted in \mathbf{P} . In the definition of the rules below, \mathbf{Q} stands for either the quantifier **A** or **THE**.

CONSTANT¹⁶

- T1_{G-P}.** (a) If $\alpha = \lambda P[P(\alpha^*)]$ and $\alpha^* \in C_e$ then $\rho_{\mathbf{G-P}}(\alpha)$ is the drawing of α^* .
 (b) If $\alpha = \mathbf{Q}(s)$ where $s \in \{\text{dot}, \text{line}, \text{curve}, \text{region}, \text{zone}\}$ then $\rho_{\mathbf{G-P}}(\alpha)$ is the drawing of whatever graphical object in C_s .

Examples: $\rho_{\mathbf{G-P}}(\lambda P[P(d_1)]) = \bullet$
 $\rho_{\mathbf{G-P}}(\lambda P[P(l_1)]) = /$
 $\rho_{\mathbf{G-P}}(\lambda P[P(r_1)]) = \text{[Dashed box with a notch]}$
 $\rho_{\mathbf{G-P}}(\mathbf{A}(\text{region})) = \text{[Dashed box]}$

LINE

- T2_{G-P}.** If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α, β are either \mathbf{D}_i or $\mathbf{Q}(\text{dot})$, and $\rho_{\mathbf{G-P}}(\alpha) = \alpha'$ and $\rho_{\mathbf{G-P}}(\beta) = \beta'$ then $\rho_{\mathbf{G-P}}(\mathbf{Q}(\text{FROM_TO}(\alpha)(\beta)(\text{line}))) = F_{P1}(\alpha', \beta')$.

Example: $\rho_{\mathbf{G-P}}(\mathbf{A}(\text{FROM_TO}(\lambda P[P(d_1)])(\lambda P[P(d_3)])(\text{line}))) = \text{[Line with a hook]}$

CURVE

- T3_{G-P}.** If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α, β are either \mathbf{R}_i or $\mathbf{Q}(\text{region})$, and $\rho_{\mathbf{G-P}}(\alpha) = \alpha'$ and $\rho_{\mathbf{G-P}}(\beta) = \beta'$ then $\rho_{\mathbf{G-P}}(\mathbf{Q}(\text{BETWEEN}_a(\alpha)(\beta)(\text{curve}))) = F_{P2}(\alpha', \beta')$.


Examples:¹⁷ $\rho_{\mathbf{G-P}}(\mathbf{THE}(\text{BETWEEN}_a(\lambda P[P(r_1)])(\lambda P[P(r_2)])(\text{curve}))) = \text{[Dashed box with a curve inside]}$
 $\rho_{\mathbf{G-P}}(\mathbf{THE}(\text{BETWEEN}_a(\lambda P[P(r_1)])(\mathbf{A}(\text{region}))(\text{curve}))) = \text{[Dashed box with a curve inside]}$

¹⁶ Rule (b) allows a graphical denoting concept in \mathbf{G} to be represented in \mathbf{P} as its concrete extension.

¹⁷ These two example expressions correspond to the abbreviated expressions in Examples 2 and 3, respectively, presented in Section 2.3.2.

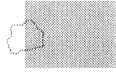
INTERSECTION

T4_{G-P}. If $\alpha, \beta \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α is either C_i or $Q(\text{curve})$ and β is either L_i or $Q(\text{line})$, and $\rho_{G-P}(\alpha) = \alpha'$ and $\rho_{G-P}(\beta) = \beta'$ then $\rho_{G-P}(Q(\text{BETWEEN}_b(\alpha)(\beta)(\text{intersection}))) = F_{P3}(\alpha', \beta')$.

Example: $\rho_{G-P}(\text{THE}(\text{BETWEEN}_b(\text{THE}(\text{BETWEEN}_a(\lambda P[P(r_1)])(\lambda P[P(r_2)])(\text{curve})))$
 $(\text{A}(\text{FROM_TO}(\lambda P[P(d_1)])(\lambda P[P(d_3)]$
 $(\text{line}))) (\text{intersection}))) =$ 


RIGHT

T5_{G-P}. If $\alpha \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α is either R_i or $Q(\text{region})$, and $\rho_{G-P}(\alpha) = \alpha'$ then $\rho_{G-P}(Q(\text{OF}_b(\text{right})(\alpha))) = F_{P4}(\alpha')$.

Example: $\rho_{G-P}(\text{THE}(\text{OF}_b(\text{right})(\lambda P[P(r_3)]))) =$ 

DOT INSIDE A REGION

T6_{G-P}. If $\alpha \in E_{\langle\langle e,t \rangle, t \rangle}$ such that α is either R_i or $Q(\text{region})$, and $\rho_{G-P}(\alpha) = \alpha'$ then $\rho_{G-P}(Q(\text{OF}_a(\alpha)(\text{dot}))) = F_{P5}(\alpha')$.

Example: $\rho_{G-P}(\text{A}(\text{OF}_a(\lambda P[P(r_1)])(\text{dot}))) =$ 

COMPOSITE REGION (1)¹⁸

T7_{G-P}. If $\alpha = \lambda P[P(\alpha^*)]$ and $\beta = \lambda P[P(\beta^*)]$ such that $\text{region}(\alpha^*)$ and $\text{region}(\beta^*)$, and $\rho_{G-P}(\alpha) = \alpha'$, $\rho_{G-P}(\beta) = \beta'$, and α' and β' are adjacent then $\rho_{G-P}(\lambda P[P(\alpha^*) \wedge P(\beta^*)]) = F_{P6}(\alpha', \beta')$.

COMPOSITE REGION (2)

T8_{G-P}. If $\alpha = \lambda P[P(\alpha^*)]$ and $\beta = \lambda P[\beta''] = \lambda P[P(\beta_1) \vee P(\beta_2) \vee \dots \vee P(\beta_n)]$ such that $\text{region}(\alpha^*)$ and $\text{region}(\beta_i)$, and $\rho_{G-P}(\alpha) = \alpha'$, $\rho_{G-P}(\beta) = \beta'$, and α' and β' are adjacent then $\rho_{G-P}(\lambda P[P(\alpha^*) \vee \beta'']) = F_{P6}(\alpha', \beta')$.

SET OF DOTS

T9_{G-P}. (a) If $\beta = \lambda P[P(\beta^*)]$ such that $\text{dot}(\beta^*)$ then $\rho_{G-P}(\beta) = F_{P6}(\emptyset, \beta)$.
 (b) If $\alpha = \lambda P[\alpha''] = \lambda P[P(\alpha_1) \vee P(\alpha_2) \vee \dots \vee P(\alpha_n)]$ and $\beta = \lambda P[P(\beta^*)]$ such that $\text{dot}(\alpha_i)$ and $\text{dot}(\beta^*)$, and $\rho_{G-P}(\alpha) = \alpha'$ and $\rho_{G-P}(\beta) = \beta'$ then $\rho_{G-P}(\lambda P[\alpha'' \vee P(\beta^*)]) = F_{P6}(\alpha', \beta')$

¹⁸ Examples for rules T6_{G-P} to T11_{G-P} are included in Figure 16 above.

SET OF LINES

- T10_{G-P}. (a) If $\beta = \lambda P[P(\beta^*)]$ such that $line(\beta^*)$ then $\rho_{G-P}(\beta) = F_{P6}(\emptyset, \beta)$.
 (b) If $\alpha = \lambda P[\alpha''] = \lambda P[P(\alpha_1) \vee P(\alpha_2) \vee \dots \vee P(\alpha_n)]$ and $\beta = \lambda P[P(\beta^*)]$ such that $line(\alpha_i)$ and $line(\beta^*)$, and $\rho_{G-P}(\alpha) = \alpha'$ and $\rho_{G-P}(\beta) = \beta'$ then $\rho_{G-P}(\lambda P[\alpha'' \vee P(\beta^*)]) = F_{P6}(\alpha', \beta')$

MAP

- T11_{G-P}. If $\alpha \in E_{\langle\langle e,t \rangle, t \rangle} = \lambda P[\alpha''] = \lambda P[P(\alpha_1) \vee P(\alpha_2) \vee \dots \vee P(\alpha_m)]$, $\beta \in E_{\langle\langle e,t \rangle, t \rangle} = \lambda P[\beta''] = \lambda P[P(\beta_1) \vee P(\beta_2) \vee \dots \vee P(\beta_n)]$ and $\delta \in E_{\langle\langle e,t \rangle, t \rangle} = \lambda P[\delta''] = \lambda P[P(\delta_1) \vee P(\delta_2) \vee \dots \vee P(\delta_r)]$ such that $region(\alpha_i)$, $dot(\beta_i)$ and $line(\delta_i)$, and $\rho_{G-P}(\alpha) = \alpha'$, $\rho_{G-P}(\beta) = \beta'$ and $\rho_{G-P}(\delta) = \delta'$ then $\rho_{G-P}(\lambda P[\alpha'' \vee \beta'' \vee \delta'']) = F_{P7}(\alpha', \beta', \delta')$.

This completes the specification of the system of multimodal representation in Figure 3. In this system, it is possible to express natural language and graphical information about maps and translate expressions between these two modalities. Natural language can be seen as stating or imposing an interpretation upon graphical representations, making the graphics meaningful. Alternatively, graphics can be seen as representing knowledge in an effective fashion. Expressions of the languages **L** and **P** can be translated through the interface language **G** in which both the semantics of **L** and the geometrical structure of **P** can be represented and reasoned about in an integrated fashion.

The system provides solid semantic ground on which to state and resolve problems of reference in multimodal scenarios. The syntactic and semantic structures of the three languages permit expression and interpretation of information in each of the modalities, and the ability to systematically find correlated expressions in different modalities with the same semantic values. As a consequence, it is possible to state formally what it means to resolve a multimodal reference: according to this theory, to resolve a multimodal reference is to find the semantic value of an expression using either the information expressed in the modality or information expressed through other modalities with the help of the translation functions. In a fully interpreted multimodal system such as the one illustrated in this section, interpreting a multimodal message is a matter of evaluating the multimodal expression. However, as argued in Section 1, the relationship between individual constants input through different modalities must be established before multimodal expressions can be evaluated. How to establish this relationship, the crucial part of the interpretation process, is illustrated in Section 3.

3. Resolution of Deictic Inference by Constraint Satisfaction

In the theory developed in Section 2, it was assumed that the translations of constants of all categories from **L** into **G** and vice versa were available, and then multimodal interpretation could be carried out; however, in the interpretation of multimodal messages, natural language and graphics are input from different sources, and working out the meaning of a multimodal message is by no means trivial. As discussed in Section 1, resolving the references and inducing the translation between graphical and linguistic terms can be thought of as the same problem. Consider, for instance, reading a book with words and pictures: when the associations between textual and graphical

symbols are realized by the reader, the message as a whole has been properly understood. However, it cannot be expected that such an association can be known in advance.

The process of inducing the translation functions for constants of **G** and **L** is similar to the computer vision problem of interpreting drawings. A related antecedent is the work on the logic of depiction (Reiter and Mackworth 1987) in which a logic for the interpretation of maps, to be applied to computer vision and intelligent graphics, is developed. It is argued that any adequate representation scheme for visual (and computer graphics) knowledge must make a distinction between knowledge of the image (the geometry) and knowledge of the scene (its linguistic interpretation), and about the relation between symbols at these two levels of representation; following Reiter and Mackworth (1987) we call this the depiction relation. In Reiter's system, two sets of first-order logic representing the scene and the image are employed. They express, respectively, the conceptual and geometrical knowledge about handdrawn sketch maps of geographical regions. In the view adopted here, the depiction relation corresponds to the translation function between constants of **L** and **G** as discussed above. An interpretation in Reiter's system is defined as a model, in the logical sense, of both sets of sentences and the depiction relation, and interpreting a drawing is a matter of finding all possible models of such sets of sentences. The domain for these models is determined by the set of individuals in the image and the scene of the picture that is being interpreted. Although computing the set of models of a set of first-order logical formulae is a very hard computational problem, the entities constituting a drawing normally form a finite set, which is often small. So, whether it is possible to compute the set of models of a given drawing is an empirical question. In particular, Reiter's system employs a constraint satisfaction algorithm to find all possible interpretations of maps, and the output of his system is a set of labels for such as "river", "road", or "shore" for curves or chains, and "land region" or "water region" for areas. As mentioned above, finding the translation functions between **G** and **L** is a similar problem, with the same level of complexity. In Section 3.1, we present a constraint satisfaction algorithm for the induction of the translation into **G** of individual constants of **L** mentioned explicitly in the text of a multimodal message. We also show how composite terms of **L** can be translated into their corresponding graphical expressions of **G** (and subsequently of **P**).

A second consideration in this section is that working out the translation between graphical and linguistic individual constants suggests a method for generating natural language expressions that refer to graphical objects and configurations. Note that inducing the linguistic translation of a graphical term that has not been mentioned overtly in the textual part of a multimodal message is the same as generating a linguistic description for the object denoted by the corresponding graphical term: once one knows the translation between individual constants of both of the modalities, the generation of multimodal descriptions can be achieved through the translation rules. For instance, in the map of Figure 4, if one points to the curve c_1 once the translation of individual constants has been found, the expression *the border between France and Germany* can be generated. This strategy for producing natural language descriptions is discussed further in Section 3.2.

3.1 Resolution of Spatial Deixis

From the point of view of our system, in interpreting multimodal messages like Figures 1 and 2, what is given are expressions of **L** and expressions of **P** and what has to be worked out is the composition $\rho_{G-P} \circ \rho_{L-G}$ and the reciprocal function $\rho_{G-L} \circ \rho_{P-G}$. However, note that the expressions of **P** are the graphical symbols on the drawings

and parsing a drawing (an expression of type *map*) produces a syntactic structure of \mathbf{P} whose translation into \mathbf{G} is the expression ζ (which we called the map). Emergent objects can also be represented in \mathbf{G} as long as they can be produced from the base through syntactic rules of \mathbf{P} and their translations into \mathbf{G} . Consequently, expressions of \mathbf{G} that refer to graphical objects stand in a one-to-one relation with the corresponding objects in \mathbf{P} . Although, theoretically, expressions in \mathbf{G} and \mathbf{P} are different representational objects, in actual interpretation processes they always come packed together. The relation between expressions of \mathbf{G} and \mathbf{L} , on the other hand, has to be worked out. For this purpose we present an algorithm for establishing a relationship between the individual constants of \mathbf{L} and the graphical constants included in the expression ζ (the map), which correspond to the interpretation domain \mathbf{P}_{base} . The algorithm for computing the translation function assigns a graphical constant to all proper names overtly mentioned in the linguistic part of a multimodal message (e.g., the graphical symbols d_1, d_2, d_3, r_1 , and r_2 to the linguistic constants *Paris, Saarbrücken, Frankfurt, France, and Germany*, respectively). The set of proper names appearing in a particular multimodal message will be referred to as *Names*. As the translations for linguistic constants of other types are given beforehand, once the translations for proper names are available, it is possible to find the graphical symbols and configurations that corefer with composite natural language descriptions through the translation rules between \mathbf{L} , \mathbf{G} , and \mathbf{P} . For instance, once the regions representing France and Germany have been identified, the term *the border between France and Germany* can be translated into an expression of \mathbf{G} , which denotes the corresponding curve, and also into the drawing of the curve in \mathbf{P} , which denotes the border between France and Germany itself. Here, it is important to highlight that the translation for individual constants cannot normally be found with the overt information expressed through the multimodal message only. For working out the interpretation of Figure 2, for instance, we need, in addition to the text and graphics, knowledge about the geography of Europe and also knowledge about the interpretation conventions of maps.

For the definition of the algorithm, a table representing the set of possible functions from linguistics predicates (e.g., *city, country*, etc.) to their corresponding graphical types (e.g., *dot, region*, etc.) is defined. This table will be referred to as a **function table**. For each particular interpretation task, a set of appropriate function tables is defined according to the following rule: For each $\delta \in C_{CN}$ of \mathbf{L} and $\delta' \in C_{(e,t)}$ of \mathbf{G} such that $\rho_{\mathbf{L}-\mathbf{G}}(\delta) = \delta'$, create a function table $(X_\delta, Y_{\delta'})$ such that:

$$\begin{aligned} X_\delta &= \{x \in C_T \mid [[x \text{ is a } \delta]]^M \text{ is true and } x \in \text{Names}\} \\ Y_{\delta'} &= \{y \in C_e \mid [[\delta'(y)]]^M \text{ is true}\} \end{aligned}$$

where X_δ and $Y_{\delta'}$ are not empty. In case either of these two sets is empty no function table for the corresponding pair is defined.

The function tables for our example are illustrated in Figure 17.

Note that if only one cell of each column of a function table is filled in, a function from proper names to graphical constants is defined. Furthermore, if the result of this process is a table in which only one cell of each row is also marked, the function is one-to-one. Accordingly, if there are n names and m graphical objects, the first column of a function table can be filled up in m different ways, the second in $m - 1$ different ways, and so on, until n graphical objects have been assigned. As a consequence, each function table with n names and m graphical objects defines $m!/(m - n)!$ possible translation functions.¹⁹ In the example, $(X_{\text{city}}, Y_{\text{dot}})$ and $(X_{\text{country}}, Y_{\text{region}})$ define 6 and 12

¹⁹ In general, if graphical objects can receive more than one name—e.g., as in the multimodal

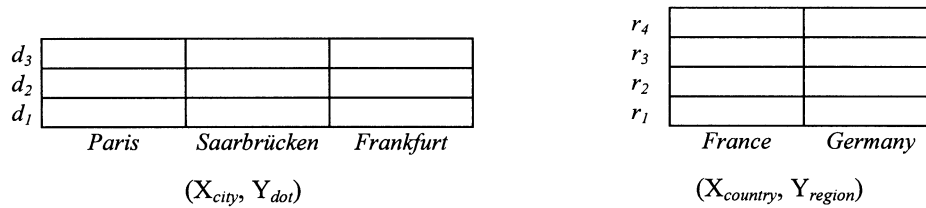


Figure 17
Function tables for the message in Figure 2.

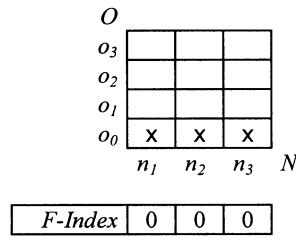


Figure 18
Set of functions associated with a function table.

possible functions, respectively. Let T_δ be the set of possible translation functions for the function table $(X_\delta, Y_{\delta'})$, and let Γ be the cross product of all T_δ in an interpretation context (i.e., the set of possible translation models). For our example, $\Gamma = T_{city} \times T_{country}$, where $|\Gamma| = 72$. This set contains 72 ordered pairs of functions, and each one represents a possible translation model for the multimodal message. Translation models can be enumerated by assigning a natural number to every cell in the array Γ . We give the following enumeration for bidimensional translation models: let $\gamma_n = (f_i, g_j)$ be the n th translation model in $\Gamma = T_x \times T_y$, where $0 \leq n < |\Gamma|$ and $f_i \in T_x, g_j \in T_y$. For every n , if $\text{mod } |T_x| \neq 0$ then $i = (n - 1 \text{ mod } |T_x|) + 1$ and $j = (n - 1 \text{ div } |T_x|) + 1$. Similar expressions can be defined for higher dimensions.

To enumerate the set of possible functions from n names to m graphical objects we use the following procedure: Let N be a list of names and O a list of graphical objects, and let $F\text{-INDEX}$ be an n -digit string containing the n digits of a numeral in base m . Every string in $F\text{-INDEX}$ codifies a total function in which the j th graphical object m_j in O (where $0 \leq j < m$) is assigned to the n_i th name in N by the rule $F\text{-INDEX}(i) = j$. The set of possible entries in $F\text{-INDEX}$ codes the m^n possible functions from n names to m graphical objects. One-to-one functions are those in which no m_j occurs more than once in a given value of $F\text{-INDEX}$. The functions sought are the $m! / (m - n)!$ one-to-one functions that result from enumerating in base m all possible values for $F\text{-INDEX}$ from 0 to $m^n - 1$, and filtering out all numbers in which the same digit occurs more than once in the enumeration order. Consider the graphical illustration of the $F\text{-INDEX}$ scheme for identifying the functions corresponding to a function table with three names and four graphical objects in Figure 18. This function table has $4^3 = 64$ possible functions out of which $4! / (4 - 3)! = 24$ are one-to-one. The graphical object m_j in O is associated to the name n_i in N by marking the corresponding cell in the table, where j is placed in the corresponding cell of $F\text{-INDEX}$. The string in $F\text{-INDEX}$ is the numeral 000 in

interpretation scenarios related to the Hyperproof system (Barwise and Etchemendy 1994)—the number of possible translation functions will be m^n , where m is the number of graphical objects and n is the number of names.

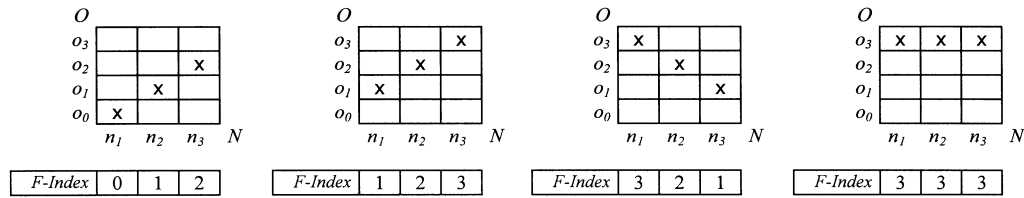


Figure 19
Examples of function index.

base 4 and represents the function in which the graphical object o_0 is assigned to all three names.

Some examples of the enumeration of functions are shown, in Figure 19. The first table shows function 12 (base 4), which is the smallest index for a one-to-one function; the second table shows function 123, which associates names $n_1, n_2,$ and n_3 to the graphical objects $o_1, o_2,$ and $o_3,$ respectively; the third table illustrates the function 321, which is the largest index for a total function in the set; and finally, the fourth table illustrates the function 333, which is a constant function assigning the object o_3 to all three names.

Armed with these concepts, we can define an algorithm for working out the interpretation of a multimodal message, as follows: Let $message_L$ be a sentence of L (the textual part of the multimodal message), θ_G an empty set of expressions of G , and Γ the set of possible translation models for $message_L$. Then, for each $\gamma_i \in \Gamma$ assume that γ_i is a translation model for $message_L$ and include its translation $message_G$ under γ_i in θ_G —i.e., $\rho_{L-G}(message_L) = message_G$. If the semantic value of all expressions θ_G in relation to the geometrical domain P_{base} is true, then γ_i is a translation model for $message_L$; otherwise, exclude γ_i from Γ . Once all translation models have been tested, check whether there is only one γ_j in Γ . If so, that γ_j is the translation function; otherwise, select a new appropriate expression of L (a general knowledge constraint) and include its translation into G in θ_G , and repeat the process until there is only one γ_j in Γ .

For our example, 4 translations out of the 72 γ 's in Γ will come out true for the first cycle of the algorithm in which the multimodal message is used as the only constraint (Example 5 in Section 2.3), as shown in Figure 20.

To continue with the algorithm, some knowledge of the geography of Europe is required. For our problem the constraints relevant to interpreting the message are illustrated in Figure 21.

The idea of the algorithm is simply to take constraints one at a time and produce the interpretation of the message incrementally. Considering constraint 1 in Figure 21, the translation functions (2) and (4) in Figure 20 can be removed; the translation function (3), in turn, can be ruled out either through constraints 2 or 3 (the interpretation of the translation of constraint 1 into G is shown in Example 4 in Section 2.3). For the example, only three cycles of the algorithm are required to rule out all but the correct translation model in Γ , which is the translation function (1) in Figure 20.

This concludes the presentation of the procedure for interpreting proper names deictically in relation to a graphical context. Although only the interpretation of this kind of constant was required for our example, the interpretation of other kinds of terms, e.g. pronouns, can be carried out in the present framework. Consider that to be able to cope with multimodal messages in which pronouns were included in the textual part, as in Figure 1, a more general definition of the language L would be required, but in such an extension both proper names and pronouns would be con-

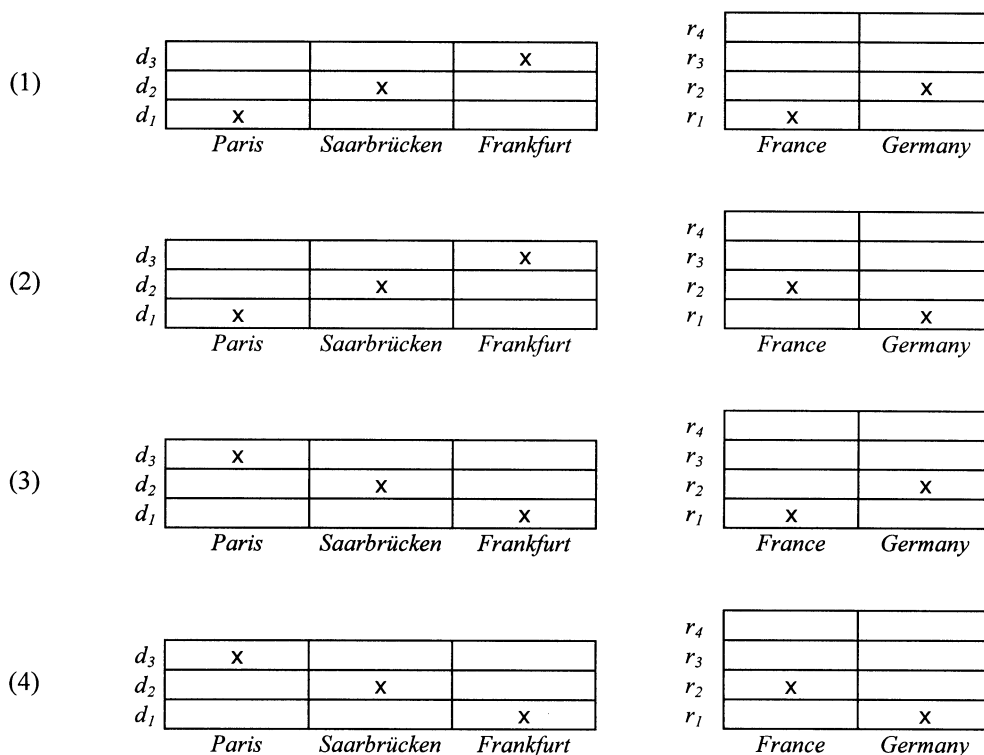


Figure 20
Possible translation models without additional constraints.

1.	<i>Germany is to the east of France.</i>
2.	<i>Paris is a city of France.</i>
3.	<i>Frankfurt is a city of Germany.</i>

Figure 21
General knowledge of geography.

stants of category T in the grammar. In the present framework, pronouns would be interpreted along the lines of proper names. For the definition of function tables, each pronoun present in a multimodal text would be included in the set $Names$, and as a first approximation, it would be a member of the domain of all function tables; different instances of the same pronoun would be considered two different objects in the interpretation process (e.g., he_0, he_1, \dots , etc.), and the interpretation would be worked out as shown above. It is also possible to think of a situation in which there are two or more graphical objects with the same name; in this context, proper names would be considered kinds of pronouns, and from the point of view of L , a different subscripted constant of category T ($name_0, name_1, \dots$) would be assigned to each such graphical object. To differentiate these objects, alternative definite descriptions could be obtained through the translation from constants of P into expressions of G , as will be argued in Section 3.2, and such descriptions could be used in the context of the particular rhetorical structures and communicative purposes of multimodal messages. A further consideration is that not only the constants of category T in a grammar can be used deictically; definite and indefinite descriptions can also be interpreted in this way. Consider that the textual part of the multimodal message in Figure 1 could have been *John washed it, the man washed it, or even a man washed it* and all three terms *John,*

the man, and *a man* would have to be interpreted deictically in relation to the graphical context. To be able to deal with this latter situation, descriptions can be interpreted deictically in our approach if terms of this kind are also included in the set *Names* for the construction of function tables. More generally, our interpretation procedure defines a function from terms into individuals of the world through the graphical context. This is because although function tables define translation models between linguistic and graphical terms, graphical objects in P denote the corresponding individuals in the world. We can think of our deictic interpretation procedure as a specific implementation for our graphical domain of Kaplan's operator DTHAT—in our simplified extensional language—which takes a term and maps it into an individual of the world in the interpretation context whenever the term is used deictically (Kaplan 1978).

The interpretation of proper names and definite descriptions has long been a source of interesting semantic problems. Consider that linguistic terms serve to identify individuals, and whenever they are used, the individual they denote should exist. However, as pointed out by Donnellan and commented on by Kaplan, “using a definite description referentially a speaker may say something true even though the description correctly applies to nothing” (Kaplan 1978). For example, suppose a bachelor enters a room accompanied by a woman who is misintroduced as his wife. Someone who notices the woman's solicitous attention to the man, says *His wife is kind to him*. The speaker uses the description *his wife* to refer to the woman, which implies that the bachelor has a wife (!), and nevertheless, what the speaker says is true. Here, one might be inclined to say that *his wife* applies to nothing, but if the woman is in the visual field of the speaker, it would be more proper to say that *his wife* applies deictically to the woman. If the expression *she is kind to him* had been used instead, or simply, the speaker had pointed to the woman at the time the expression *is kind to him* were uttered, the deictic nature of the reference would be easily revealed. According to Kaplan, whenever a description is used referentially (as opposed to attributively), describing can be taken as a form of pointing, and as he suggested, instead of taking the sense of a description as the subject of a proposition, the sense is used only to fix the denotation, which is then taken directly as the subject of the proposition. Similarly, although a proper name is usually thought of as related to an individual (the bearer) in an intimate fashion through an interpretation function in model theory, and it is often stated that proper names are related to the same individual through all world and time indices (i.e., as rigid designators [Kripke 1972]), we would argue proper names can be also interpreted deictically to fix a referent, which can then be taken directly as the subject of a proposition.

3.2 Generation of Natural Language Descriptions

The multimodal representation scheme and the resolution of deictic inferences presented above permit the generation of multimodal descriptions in a simple and systematic fashion. Once a multimodal representational system is fully defined, the generation of graphical and linguistic expressions can be achieved directly through the translation rules. As the crucial piece of knowledge required for use of the translation rules is the translation model, the deictic inference required to identify an individual and the inference required to generate a description for such an individual are but two sides of the same coin.

If a graphical object is pointed out on the screen, a number of natural language descriptions to refer to it can be produced. Several strategies for finding an appropriate description are available, depending on whether the object pointed at is in P_{base} or whether it is an emergent object. Another consideration is whether the object has

a proper name or can be referred to either by a definite or an indefinite description. Suppose that a graphical cursor for pointing to graphical objects is available. The cursor itself is modeled as a graphical object of type *dot* within the graphical domain. With this interactive device we can identify dots if the position of a dot in a drawing is the same as the position of the cursor, lines and curves if the cursor lies on the line or the curve, and regions if the cursor is inside or on the border of a region. With this device, we can select all basic or emergent graphical objects that are identified by an individual pointing act. Basic objects will be identified directly, and the emergent objects selected by a pointing act will be those that can be produced by the grammar of **P** and satisfy the geometrical conditions associated with the cursor. Objects identified by a cursor will be terms of **P** that can be translated into **G** as proper names, definite descriptions, or indefinite descriptions. To refer to a graphical object we can use the following simple strategy: if a graphical object can be translated into **L** as a proper name, use the proper name; if the graphical object can be translated as a description and it is the only one satisfying such a description, use a definite description; otherwise, use an indefinite description.

Consider a pointing action in which the dot d_1 or the region r_1 is selected. The translations from **P** into **G** of these objects, according to rule $T1_{P-G}$ (a), are $\lambda P[P(d_1)]$ and $\lambda P[P(r_1)]$, respectively. As these expressions can be translated into **L** as shown in Figure 14, these objects can be referred to in **L** as *Paris* and *France*, respectively. However, consider that as these objects are the concrete extension of a number of denoting concepts of **G**, these objects can also be translated into **G** through rule $T1_{P-G}$ (b) as $A/THE(dot)$, and $A/THE(region)$, respectively, and they can be translated into **L** as *a/the city* and *a/the country*. Another possible translation for the dot d_1 whenever it is produced by rule $S4_P$ and translated into **G** through $T4_{P-G}$ is $A/THE(OF_a(R_1)(dot))$, which in turn can be translated into **L** as *a/the city of France*.

There may be constants of **P** that cannot be translated into **L** as proper names. Consider, for instance, the line l_1 in Figure 4; the translation of this line into **G** is $\lambda P[P(l_1)]$, but as can be seen in Figure 14, there is no proper name that corresponds to this expression in **L**. However, this constant could be translated by rule $T1_{P-G}$ (b) as the denoting concepts $A/THE(line)$. As the line is also an emergent object that can be produced via the syntactic rule $S2_P$, its translations according to rule $T2_{P-G}$ are $A/THE(FROM_TO(D_1)(D_3)(line))$; subsequently, such an object can be referred to in **L** as *a/the line* and *a/the line from Paris to Frankfurt*. A similar example is the description of the curve c_1 , which is produced by rule $S3_P$ and can be translated into **L** through **G** as *a/the border between France and Germany*.

As mentioned in Section 1, the generation of descriptions is required within the context of specific rhetorical and intentional structures, such as the *activate* structure of the WIP system, which employs Reiter and Dale's algorithm for the production of definite descriptions on demand. Our system can be used to support the generation of descriptions either definite or indefinite, and even pronouns used deictically, in multimodal generation systems with a solid semantic base. These descriptions could be used according to particular rhetorical and intentional structures related to specific application domains. The advantage of such an approach is that the choice of the expressions to be used in multimodal presentations could be made not only on the basis of predefined heuristics, but also on the basis of the semantic value of these expressions in the context of use. In addition, the decision about what kind of knowledge is expressed through either modality for the production of coordinated natural language and graphical explanations could take into account not only the kind of heuristics that are currently employed in systems like WIP and COMET, but also the expressiveness and effectiveness criteria of natural and graphical languages.

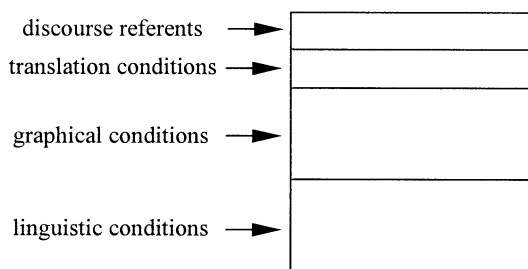


Figure 22
Components of the multimodal discourse representation structure.

4. Multimodal Discourse Representation Theory

The ability to interpret individual multimodal messages is a prerequisite for interpreting sequences of multimodal messages occurring in the normal flow of interactive conversations. In the same sense that discourse theories, like DRT, are designed to interpret sequences of sentences, it is desirable to have a theory in which sequences of multimodal messages can be interpreted. Such a theory would have to support anaphoric and deictic resolution models in an integrated fashion, and would have to be placed in a larger pragmatic setting in which intentions and presuppositions are considered, and in which mechanisms to retrieve knowledge from memory are also taken into account. To work out such a theory is quite an ambitious goal; however, in the same way that DRT focuses in internal structural processes that govern anaphoric resolution, it is plausible to consider a multimodal discourse representation theory (MDRT) to cope with the resolution of spatial deictic inferences. In the same way that DRT postulates discourse representation structures in which referents and conditions are introduced incrementally through the interpretation of the incoming natural language discourse by means of the application of construction rules, it is plausible to conceive similar multimodal discourse representation structures (MDRS) whose referents and conditions would be introduced by modality-dependent construction rules acting upon the expressions of the corresponding modality. In these structures, DRS conditions extracted from different modalities would be kept in separate partitions, but discourse referents would be abstract objects common to the whole MDRS. In particular, MDRS's could help to specify accessibility relations between anaphoric and deictic terms and their antecedents and interpretation context, imposing severe constraints on the possible interpretations, as is normal in DRT. The resolution process itself would be accomplished by incremental constraint satisfaction, as shown for deictic inferences. In the rest of this section, we present a schematic picture of how an MDRS can be developed, and illustrate using the interpretation of the multimodal message in Figure 2. Consider first the empty MDRS in Figure 22.

The MDRS is a structure with four partitions; it extends traditional DRS with one partition for graphical conditions and another to store the translation models that hold in a particular interpretation state. The partition for linguistic conditions is used as in normal DRS, and the top partition for referents includes a variable for every individual that is referred to in the multimodal message in either of the modalities. Figure 23 illustrates the initial state for the interpretation process of the multimodal message in Figure 2. Graphical expressions of G (the map) are included in the graphical section of the MDRS, and textual conditions, with the associated type information, are included in the linguistic section as in normal DRS. A refer-

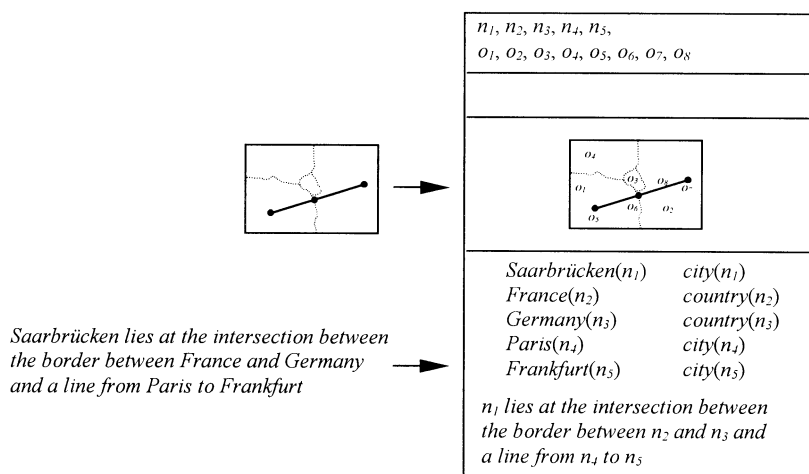


Figure 23
Initial MDRS for the interpretation of the multimodal message.

ent is included in the corresponding partition for every individual that has been introduced through either modality, as referents are considered medium-independent abstractions. In the same way that the order of processing of linguistic information is not crucial for the definition of the linguistic conditions, we abstract over scanning considerations and assume that graphical expressions are introduced as a single “sentence.”²⁰ Finally, the partition for the possible translation models is empty at this stage, as the coreference relation between text and graphics has not yet been established.

The interpretation process by constraint satisfaction is illustrated in Figure 24. Figure 24(a) illustrates the interpretation state after the first cycle of the constraint satisfaction algorithm presented in Section 3.1 has been carried out. In this state, the partition for the translation conditions contains the disjunction of the four possible translation models that are consistent with the message, taking the message itself as the only interpretation constraint. Figure 24(b) illustrates the interpretation state once the additional constraint that Germany is to the east of France has been considered.²¹ The interpretation of the corresponding expression introduces two additional discourse referents (n_6 and n_7), as the terms *Germany* and *France* in the textual part of the message should be resolved anaphorically in relation to the context previously built. However, this anaphoric resolution process is kept within the linguistic section of the MDRS and should take into account the accessibility constraints between anaphor and antecedent, as commonly done in DRT. The result of this anaphoric inference is reflected in the equality conditions $n_6 = n_3$ and $n_7 = n_2$. The inclusion of the constraint *Germany is to the east of France* permits us to rule out two possible translation models, and the result of the second cycle of the constraint satisfaction

²⁰ We leave for further research whether the analysis of scanning protocols by means of eye-tracking techniques can provide information for imposing additional constraints on accessibility relations and possible translation models for the construction of MDRS's. An interesting antecedent for the definition of such constraints can be found in Faraday and Sutcliffe (1998).

²¹ How this constraint is selected is beyond the scope of this paper and we only make the assumption that the symbols in the graphical and linguistic partitions of the MDRS form a part of the indexing scheme required to retrieve the information from memory. For a prototype implementation, this kind of constraint could be provided by the human user directly.

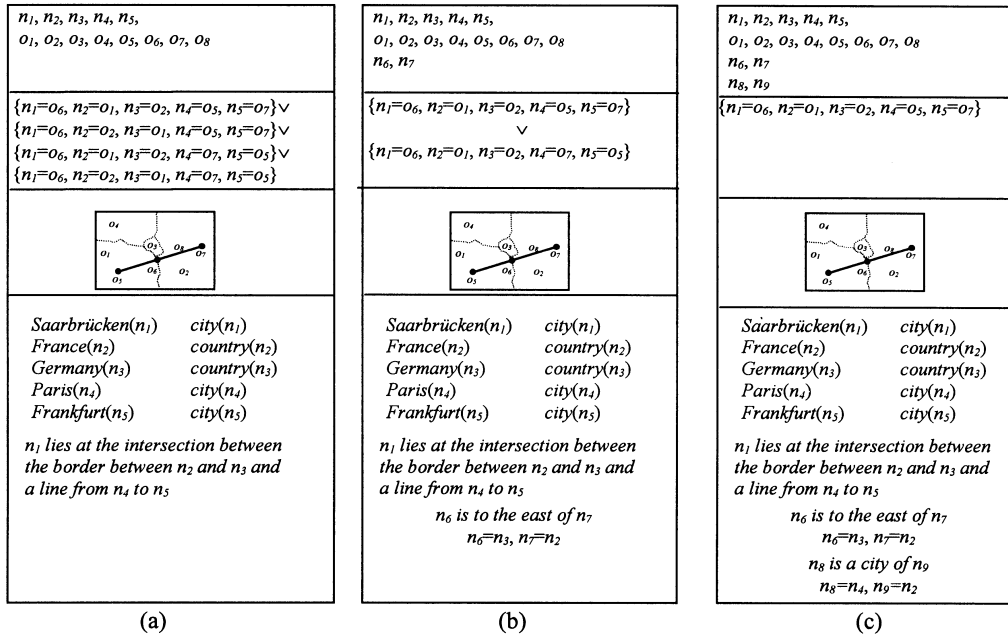


Figure 24
 Interpretation of multimodal message by constraint satisfaction.

algorithm is reflected in the new state of the partition for translation conditions of the MDRS. Figure 24(c) illustrates the final interpretation state in which the constraint that Paris is a city of France has been considered and involves anaphoric and deictic resolution inferences as in the interpretation of the previous constraint. As a result of this last constraint satisfaction cycle, only one translation model is left in the partition for translation conditions and reflects the correct interpretation of the multimodal message.

As a last example of the integrated anaphoric and deictic interpretation, consider a situation in which the natural language expression *It is big* is mentioned after the multimodal message in Figure 2 has been interpreted, as illustrated in Figure 25. In this situation, the natural language information would enter into the partition for linguistic conditions and the pronoun *it* should be interpreted anaphorically in relation to the context currently provided by the MDRS and could resolve to *Saarbrücken* (although there are several possibilities). However, if the expression is supported by an overt gesture indicating the city of Paris, for instance, it would be deictic and its interpretation would have to be worked out with the same machinery; although in this latter situation the translation relation between graphical and linguistic referents could be asserted directly in the translation model as the gesture would render unnecessary the constraint satisfaction part of the deictic inference.

With this we conclude the presentation of our model for integrated deictic and anaphoric inferences. The distinction between anaphora and deixis is clearly demarcated. The antecedent for a pronoun, a proper name, or a description used anaphorically is provided by the discourse interpretation context, while the referent for a deictic pronoun, proper name, definite or indefinite description, or a demonstrative word like *this* or *that* is taken from an intermediate representation of a nonlinguistic modality such as the graphical context, and denotes an individual of the world directly, a view that is consistent with Kamp's distinction quoted in Section 1.

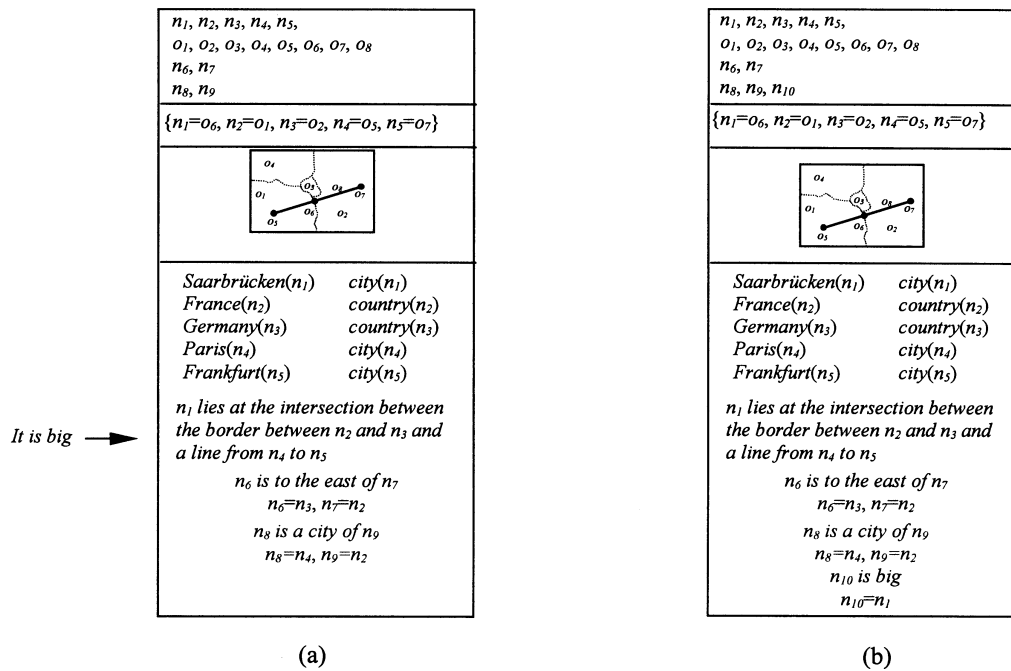


Figure 25
 Integrated interpretation of anaphora and spatial deixis.

5. Conclusions

In this paper, we have presented a theory of representation and interpretation for multimodal messages and a model for multimodal reference resolution. The model is based on the view that a modality is a code system on a medium that can be characterized by well-defined syntax and semantics. Multimodal interpretation is a matter of working out coreference relations between terms of different modalities. A central concern in articulating this theory is a clear characterization of how spatial deictic reference is resolved and of how spatial deictic reference relates to the resolution of anaphors in the normal flow of discourse. A key theoretical assumption we make is that graphics are interpreted deictically, which is in opposition to the view graphical representations are interpreted anaphorically.

The theoretical machinery for the definition of the syntax and semantics is formally developed along the lines of Montague’s semiotic program and its associated general theory of translation. We have also illustrated an algorithm for finding the translation between texts and graphics, as messages in these modalities are introduced through independent input channels, and the translation between linguistic terms and their corresponding graphical expressions must be induced dynamically. We also suggested an extension of Kamp’s DRT with multimodal discourse structures (MDRS). This model defines an integrated interpretation model for multimodal messages while maintaining a clear demarcation between indexical and anaphoric inferential processes. Natural language terms, like proper names, pronouns, and descriptions can be interpreted in relation to a model; however, these linguistic terms also admit anaphoric and deictic interpretations.

It is important to note that although we used a simplified extensional definition for the semantics of natural language and graphical expressions, the system was carefully

designed to move smoothly into the intensional domain. Consider that the extensional formulation used in the semantic definition of the graphical language can be easily extended into an intensional one by changing the types of constants, predicates, and sentences from individuals, sets of individuals, and truth values, into individual concepts, properties of individuals, and propositions, respectively. This is achieved simply by indexing the interpretation of expressions in terms of a possible world and time, and all definitions presented above could be considered relative to the current world and time. The move to the intensional domain would allow the definition of the interpretation of more comprehensive natural language segments.

Intensionality is also relevant for the interpretation of graphical languages, in general, and for the definition of graphical and linguistic interactive systems, in particular. In interactive sessions with a computer graphic interface, the interaction states can be considered as possible worlds and the interpretation of graphical constants would depend on particular graphical states. If a graphical object like a dot, for instance, is moved from one position to another in an interactive transaction, we have the intuition that the object before and after the change is the same and denotes the same object of the world and yet not even its position, which one could think of as an essential property of a dot, is the same. Accordingly, in the intensional setting the semantic value of a graphical constant is not an individual but an individual concept. Consider as well that the same graphical description can have different semantic values in different interactive states; for instance, the value of the expression *position of d_1* will be a different ordered pair before and after dot d_1 is moved. According to this, the interpretation of graphical operators at every index will be a function from sequences of graphical objects of the proper kind into graphical objects; however, unlike normal linguistic situations in which different functions at different indices are assigned to operators and predicates, the same function at every index has to be assigned to geometrical operators and predicates, as the geometry is always the same. Moving into the intensional setting is also relevant for our treatment of indexicals. In our current approach, the interpretation of a term used deictically is an individual of the world; in the intensional context, the interpretation of the same term in one particular interaction state will be the same in every state despite the fact that the description for referring to such an object in the state in which it was selected might pick up a different individual in a different state.

In the future, it would be interesting to deal with a more general fragment of natural language that includes temporal expressions. In the same way that the language **G** provided a finite and small domain for the interpretation of linguistic spatial prepositions, a similar language **T** for the interpretation of temporal prepositions could be defined. Temporal predicates and operators of this language would be interpreted in terms of arithmetic functions like those presented, for instance, in Allen's temporal logic (Allen 1983). In the same way that the constraint satisfaction algorithm for the definition of the translation between graphical and linguistic terms helped to solve deictic inferences, a constraint satisfaction algorithm for resolving temporal deictic references in relation to a finite and small domain of actions and events is conceivable. The definition for such a spatial and temporal indexical model could be quite helpful for the implementation of natural language and graphics systems in which actions and events are mentioned in the course of interactive conversations.

6. Implementation

Although a prototype system for the theory presented in this paper has not been implemented, several aspects of the theory have been implemented in relation to simpler systems. A simpler version of the strategy for multimodal interpretation of the scheme

in Figure 3 was implemented in the first version of the Graflog system (Pineda 1989). Several versions of the graphical language and its geometrical interpreter have been implemented in relation to different application domains (Morales 1994; Masse 1994; Santana 1999; Garza 1995) with BinProlog and the TCL/TK programming environment. The geometrical interpreter and the strategy of evaluating a set of geometrical constraints incrementally in relation to a graphical domain was used in a later version of Graflog to solve and generate graphical explanations of geometrical constraint satisfaction problems (Pineda 1992, 1998), and also for the definition of a model (not yet fully implemented) for the production of solids from orthogonal views of polyhedra (Garza and Pineda 1998). We also implemented the scheme for enumerating functions used in the definition of translation models for a semantic theorem-proving system written in Prolog, in order to find the possible models satisfying logical theories about graphical scenarios of the Hyperproof system (Barwise and Etchemendy, 1994).

Acknowledgments

The authors gratefully acknowledge support to Luis Pineda from the Institute for Applied Mathematics and Systems (IIMAS) at the National University of Mexico (UNAM) and Conacyt grant 400316-5-27948-A and to Luis Pineda and Gabriela Garza from the Institute for Electrical Research (IIE), Mexico. We are grateful also to the anonymous reviewers, and for helpful discussions with numerous people including James Allen, Elisabeth André, Kees Van Deemter, John Lee, Sergio Santana, Oliviero Stock, Thomas Rist, Henk Zeevat, and very specially to Ewan Klein.

References

- Allen, James. 1983. Maintaining knowledge about temporal intervals. *Communications of the ACM*. 26(11): 832–843.
- André, Elisabeth and Thomas Rist. 1994. Referring to world objects with text and pictures. In *Proceedings of COLING 94*, pages 530–534.
- Barwise, Jon and John Etchemendy. 1994. *Hyperproof*. CSLI.
- van Deemter, Kees and Stanley Peters, editors. 1995. *Semantic Ambiguity and Underspecification*. CSLI Publications 1996, Stanford, CA.
- Dowty, David R., Robert E. Wall, and Stanley Peters. 1985. *Introduction to Montague Semantics*. D. Reidel Publishing Company, Dordrecht, Holland.
- Faraday, Pete and Alistair Sutcliffe. 1998. Providing advice for multimedia designers. In L. Pineda, T. Rist, and J. Lee, editors, *Proceedings of the Workshop on Interpretation and Generation in Intelligent Multimodal Systems and Graphical Reasonings in Expert Systems*, pages 72–79, Fourth World Congress on Expert Systems. ITESM Mexico City Campus, Mexico.
- Feiner, Steven and Kathleen McKeown. 1993. Automating the generation of coordinated multimedia explanations. In M. Maybury, editor, *Intelligent Multimedia Interfaces*. The MIT Press and AAAI Press, Menlo Park, CA, pages 117–239.
- Garza, Gabriela. 1995. *Síntesis de Poliedros a partir de sus Vistas Ortogonales: Un Caso de Estudio acerca del Razonamiento Gráfico*. M. Sc. thesis, ITESM, Campus Morelos, Mexico.
- Garza, Gabriela and Luis Pineda. 1998. Synthesis of solid models of polyhedra from their orthogonal views using logical representations. In *Expert Systems with Applications*. Elsevier Science Ltd., Volume 14, pages 91–108.
- Kamp, Hans. 1981. A theory of truth and semantic representation. *Formal Methods in the Study of Language*, 136: 277–322. Mathematical Centre Tracts.
- Kamp, Hans and Uwe Reyle. 1993. *From Discourse to Logic*. Kluwer Academic Publisher, Dordrecht, Holland.
- Kaplan, David. 1978. DTHAT. *Syntax and Semantics*. Volume 9, pages 383–399.
- Klein, Ewan and Luis Pineda. 1990. Semantics and graphical information. In Diaper, Gilmore, Cockton, and Shackel, editors, *Human-Computer Interaction, Interact'90*. IFIP, North-Holland, pages 485–491.
- Kripke, Saul. 1972. *Naming and Necessity*. Basil Blackwell, Oxford.
- Lyons, John. 1968. *Introduction to Theoretical Linguistics*, Cambridge University Press, Cambridge.
- Mackinlay, Jock Douglas. 1987. *Automatic Design of Graphical Presentations*. Ph.D. thesis, Stanford University. University Microfilms International.
- Mann, William C. and Sandra A. Thompson. 1988. *Rhetorical Structure*.

- Theory: Toward a functional theory of text organization. *Text*, 8(3): 243–281.
- Masse, J. Antonio. 1994. Satisfacción de Restricciones por Referencia Simbólica en Dibujos Geométricos. B. Sc. thesis, ENEP Aragón, UNAM, Mexico.
- Maybury, Mark. 1993. Planning multimedia explanations using communicative acts. In M. Maybury, editor, *Intelligent Multimedia Interfaces*. The MIT Press and AAAI Press, Menlo Park, CA, pages 59–74.
- Moore, Johanna. 1995. *Participating in Explanatory Dialogues: Interpreting and Responding to Questions*. The MIT Press, Cambridge. A Bradford Book.
- Morales, Rafael. 1994. Pizarrones Interactivos Multimodales para la Enseñanza de Conceptos Matemáticos, M. Sc. thesis, ITESM, Campus Morelos, Mexico.
- Pineda, Luis. 1989. *Graflog: A Theory of Semantics for Graphics with Applications to Human-Computer Interaction and CAD Systems*. Ph.D. thesis, University of Edinburgh.
- Pineda, Luis. 1992. Reference, synthesis and constraint satisfaction. *Computer Graphics Forum*. 11(3): 333–344.
- Pineda, Luis. 1998. Graphical and linguistic dialogue for intelligent multimodal systems. In *Expert Systems with Applications*. Elsevier Science Ltd., Volume 14, pages 149–157.
- Poesio, Massimo. 1994. *Discourse Interpretation and the Scope of Operators*. Ph.D. thesis, University of Rochester.
- Reiter, Ehud and Robert Dale. 1992. A fast algorithm for the generation of referring expressions. In *Proceedings of the COLING'92*. Volume 1, pages 232–238, Nantes, France.
- Reiter, Raymond and Alan K. Mackworth. 1987. The logic of depiction. Technical Reports on Research in Biological and Computational Vision at the University of Toronto. RCBV-TR-87-18.
- Rist, Thomas. 1996. Current state of the reference model for intelligent multimedia presentation systems. Paper presented in the workshop “Towards a Standard Reference Model for Intelligent Presentation Systems” at the 12th European Conference on Artificial Intelligence. Budapest. August.
- Santana, Sergio. 1999. *The Generation of Coordinated Natural and Graphical Explanations in Design Environments*. Ph.D. thesis, Universidad de Salford.
- Shamos, M. I. 1978. *Computational Geometry*. Ph.D. thesis, Yale University. University Microfilms International.
- Steedman, Mark J. 1986. Incremental interpretation in dialogue. ACORD Project Deliverable T2.4. Department of Artificial Intelligence and Centre for Cognitive Science. University of Edinburgh.
- Stiny, G. 1975. *Pictorial and Formal Aspects of Shape Grammars*. Birkhauser Verlag, Basel.
- Stock, Oliviero and the AlFresco Project Team. 1993. AlFresco: Enjoying the combination of natural language processing and hypermedia for information exploration. In M. Maybury, editor, *Intelligent Multimedia Interfaces*. The MIT Press and AAAI Press, Menlo Park, CA, pages 197–224.
- Wahlster, Wolfgang. 1991. User and discourse models for multimodal communication. In J. W. Sullivan and S. W. Tyler, editors, *Intelligent User Interfaces*. ACM Press, New York, pages 45–67.
- Wahlster, Wolfgang, Elisabeth André, Wolfgang Finkler, Hans-Jürgen Profitlich, and Thomas Rist. 1993. Plan-based integration of natural language and graphics generation. *Artificial Intelligence* 63: 387–427.
- Wittenburg, Kent. 1998. Visual language parsing: If I had a hammer In Harry Bunt, Robbert-Jan Beun, and Tijn Borghuis, editors, *Multimodal Human-Computer Communication: Systems, Techniques and Experiments*. Springer-Verlag, pages 231–249.