

Using a Generative Lexicon Resource to Compute Bridging Anaphora in Italian.*

Utilización de un recurso de léxico generativo para calcular Anáfora asociativas en Italiano.

Tommaso Caselli

ILC- CNR and Dip. Linguistica "T.Bolelli", Università degli Studi di Pisa
Via Moruzzi, 1 56124 Pisa, Italy
tommaso.caselli@ilc.cnr.it

Resumen: Este artículo presenta un trabajo preliminar sobre el uso de un recurso léxico basado en la teoría del léxico generativo para resolver las anáforas asociativas en italiano. Los resultados obtenidos, a pesar de no ser demasiado satisfactorios, parecen respaldar el uso de un recurso de este tipo respecto a los recursos de tipo WordNet debido al mayor número de anáforas asociativas que puede tratar.

Palabras clave: léxico generativo, resolución de anáfora, *bridging*, anáforas asociativas

Abstract: This article reports on a preliminary work on the use of a Generative Lexicon based lexical resource to resolve bridging anaphors in Italian. The results obtained, though not very satisfying, seem to support the use of such a resource with respect to WordNet-like ones due to the wider range of bridging anaphors which can be treated.

Keywords: generative lexicon, anaphora resolution, bridging

1 Introduction

Anaphora resolution is essential to capture the knowledge encoded in text. Bridging anaphora are a very challenging phenomenon because they are a "type of indirect textual reference whereby a new referent is introduced as an anaphoric not of but via the referent of an antecedent expression" (Kleiber, 1999, 339), as in the following example (bridging NPs are in bold):

- (1) Maria ha comprato *una macchina nuova*, ma **il motore** si è rotto dopo due giorni.

*Maria bought a new car, but **the engine** broke down two days later.*

Bridging anaphors are constrained to a set of semantic and pragmatic conditions. The aim of this paper is to present a preliminary study on the use of a Generative Lexicon based lexical resource (SIMPLE) as a source of these constraints to automatically resolve this kind of anaphoric definites. In order to develop the system, we have preliminary

conducted a corpus study on the identification and classification of bridging anaphors in Italian. The corpus study has been grounded on a set of theoretical statements describing the phenomenon of bridging, providing empirical evidences of their validity and also further information on their organization.

The paper is organized as follows: in section 2, we will present the semantic and pragmatic constraints underlying the phenomenon of bridging anaphora. The corpus study and its results are illustrated in section 3. We will then describe how the lexical resource is structured and what levels of semantic information encoded in it are the most relevant to accomplish the task of resolving bridging anaphors in section 4. Finally in section 5, we will describe the results obtained from the use of SIMPLE and compare its performance with that of a WordNet-based resource, namely ItalWordNet, and present our concluding remarks and observations in section 6.

2 Theoretical background

A trend in linguistic theories, which has counterparts in computational frameworks, tends to emphasize the idea that Full Definite Noun Phrases (FDNPs henceforth) are a matter of the global discourse focus, i.e.

* A preliminary version of this work has been presented at the CBA Workshop at the Universitat de Barcelona, Barcelona, 13-15 November 2008. The author wants to thank the organizers and participants for the useful comments and discussion.

they are used to retrieve a referent which is no longer accessible or to construct a conceptual representation which uniquely identifies a referent. On the contrary, empirical studies provided evidence in favor of Sidner (1979)'s hypothesis that bridging FDNPs are different from other occurrences of anaphoric FDNPs, since, in the process of identification of their antecedents, they are more sensitive to the local focus. In addition to this, bridging FDNPs trigger an inferential presupposition of the kind:

$$the[N1]R[N2] \quad (1)$$

where N1 represents the FDNP, i.e. the bridging anaphor, R is the inferential relation or bridge the interpreter has to perform in order to interpret correctly its occurrence¹, and N2 is the antecedent or anchor. Applying the formula in 1 to the example in 1 we obtain the following paraphrasis “*the [engine]_{N1} is a part_of_R [a car]_{N2}” which justifies the occurrence of the FDNP.*

Kleiber (1999) identifies some semantic restrictions on what kinds of FDNPs can enter a bridging relation. Drawing on the notion of functional nouns², he identifies two very general, language-independent factors which are at work in the mechanism of the bridging relation between the referents involved: a condition of alienation and the principle of ontological congruence. A bridging description can be conceived of as a Functional Concept of type 2 (FC2), with an implicit argument. This type of semantic definite NP introduces the referent by means of the sole sortal predicate N, without semantic subordination to another individual. In other words, the head noun looks as semantically autonomous or alienated.

Next to these semantic restrictions, a couple of pragmatic constraints can be identified. We propose to use the following pragmatic restrictions on inferencing: an Effort Condition and a Plausibility Condition as suggested by Krahmer and Piwek (2000). The two constraints can be represented by the following maxims:

- use your informational resources as little as possible (Effort Condition);
- make as few assumptions as possible (Plausibility Condition).

The Effort Condition has to do with the mental capacity the interpreter needs to resort to in order to construct a “bridge”. In particular, it states that the less time consuming inference to retrieve the right anchor should be preferred over the others. The Plausibility Condition, on the other hand, has to do with the admissibility of the constructed bridges. It is a simple consistency condition, with relevance as a side effect. The Plausibility Condition plays a major role in selecting the most plausible reading among those which passed the Effort Condition, helping us to determine the bridge and avoid ambiguity. Obviously, if the Effort Condition selects only one reading, this is considered the most plausible by definition.

The inference the hearer has to perform in order to bridge the gap from what s/he knows to the intended antecedent, bears on the possible relation(s) between the referent of the antecedent and the referent of the anaphor. The existence of such a relation is necessary for the speaker to create the bridge and for the hearer to resolve it. Most classifications of bridging anaphoras are all based on this idea (Hawkins, 1978) (Sidner, 1979). The relations that link the anaphor to the antecedent can be of various types, but they can be reduced to three pragma-cognitive dimensions: a lexical semantic dimension, a contextual, or textual, dimension and a contextual, or extralinguistic, dimension.

These elements represent the theoretical background which we have used both in the corpus-study and in the development of the automatic procedure to resolve bridging anaphors. In particular, the identification of the *R* relation between the bridging definite and its anchor has been used to identify the various classes of bridging anaphors, and the Effort and Plausibility conditions have been exploited to restrict the type and number of NPs which could be identified as anchors.

3 *Bridging Anaphora in Italian: a corpus study*

In order to verify the realizations of bridging anaphors in Italian, we have conducted a corpus study on 17 randomly chosen arti-

¹The R relation can be thought as deriving from Chierchia (1995)'s compositional semantics of FDNPs, according to which “the + N” denotes a noun N which is related in an anaphorically undetermined way B to an antecedent *u*.

²By functional nouns we intend NPs denoting a non-ambiguous interpretation, or a functional concept (FC), as proposed by Lobner (1985).

cles from the Italian financial newspaper “*il Sole-24 Ore*”, a workpackage of the SI-TAL Project, the syntactic-semantic Treebank of Italian (Montemagni et al., 2003).

The texts considered contain a total number of 1412 full definite noun phrases (FDNPs) of the form “definite article + (possessive) + N”, which represent 31.54% of all the occurrences of FDNPs in the corpus. Each newspaper article was first read entirely, and only after it was divided into segments of five sentence windows which is an arbitrary strategy to give an account of the local focus of the text i.e. the most probable place to look for anchors for bridging FDNPs.

In the classification exercise we have used an operational device such as processing requirements³ since when a FDNP is encountered in a discourse can be reduced to one of these four cases:

- it is used to pick up an entity mentioned before in the text, which, in our experiment, could be either directly or indirectly realized;
- it is not mentioned before, but its interpretation depends on , is based on, or is related in some way to an entity already present in the discourse (directly or indirectly realized);
- it is not mentioned before and is not related to any previous mentioned entity, but it refers to something which is part of the common shared knowledge of the writer and reader;
- it is self-explanatory or it is given together with its own identification.

These four types of FDNPs use reflect the classes of *Direct Anaphora*, *Bridging* and *First Mention*, respectively. The same operational device i.e. processing requirements, was used for the analysis and classification of bridging anaphors.

The classification task has led to the identification of 6 main classes of FDNPs (Table 1)⁴. One of the main interesting results deriving from the classification in 1 is represented by the class of Bridging which represents the 63.88% (299/469) of all anaphoric FDNPs,

³See also Vieira and Poesio (2000).

⁴For detailed figures and comments on the corpus study readers are referred to Caselli (2007).

FDNPs Classes	Figures
First Mention	833 (58.61%)
Possessives	36 (2.54%)
Direct Anaphora	170 (12.03%)
Bridging	299 (21.17%)
Idiom	25 (1.62%)
Doubt	49 (3.47%)
Total	1412 (100%)

Table 1: Classes of FDNPs.

thus suggesting that bridging is a more productive cohesive strategy in Italian with respect to other languages, i.e. English (Vieira and Poesio, 2000).

Five subclasses of bridging anaphors have been identified, in particular:

- Lexical: (199/299 - 39.79%) those instances of bridging descriptions whose link with the antecedent is clearly based on lexical semantics, e.g.: *la pistola - l’arma* (*the gun - the weapon*);
- Event: (18/299 - 6.02%) the antecedent is represented by a verb or a VP; it contains what Clark categorizes as indirect reference by necessary roles and optional roles, and Strand’s event-argument relations, e.g.: *fece esplodere - le macerie* (*exploded - the debris*);
- Rhetorical Relation⁵: (27/299 - 9.03%) it includes bridging anaphors whose antecedent can be identified through discourse relations, e.g.: *l’elezione - i componenti* (*the election - the members*);
- Discourse Topic: (26/299 - 8.69%) this kind of bridging is related on implicit way to the main discourse topic of a text, rather than to a specific NP or VP;
- Inferential: (109/29 - 36.45%) all cases of bridging based on complex inferential reasoning which entails use of encyclopedic, background or common shared knowledge, e.g.: *la Cina - Pechino* (*China - Beijing*).

As the classes show, different sources of information (lexical, encyclopedic and discourse structure) have important roles for

⁵It contains Clark (1997)’s relations of *reasons*, *causes* and *consequences*, part of Vieira and Poesio (2000)’s *inferential* bridging and Strand (1997)’s *argument-event*.

the resolution of these kinds of anaphoric relations. The results also suggest a preference order for the different sources of bridging anaphora: lexical semantic relations are preferred over the use of common sense inferencing and background knowledge i.e. pragmatics, which is preferred over discourse structure. Nevertheless, as it emerged from the corpus study, more than the 45% of the *R* relations needed to resolve bridging anaphors are based on commonsense knowledge (the Inferential class) and on general discourse structure (the Rhetorical Relation class).

Different strategies have been proposed to automatically resolve bridging anaphors. Most of them rely on the use of lexical resources like WordNet or WordNet-like. However, the results obtained are not very satisfactory for two main reasons: on the one hand, lexical resources have limits due to the fact that they represent closed representations of natural language and could present mistakes and missing information due to their human-based nature, and, on the other hand, the theoretical background behind their construction is unable to deal with lots of instances of *R* relations, as we have called them, which govern the ways in which bridging anaphors can be retrieved and inferred by the interpreters.

In this work we propose to use a lexical resource as well, namely PAROLE/SIMPLE/CLIPS (henceforth SIMPLE) (Ruimy et al., 2003), but the novelty of our proposal does not rely in the use of a lexical resource *per sè*, but in the use of a resource grounded on a robust lexical theory like that of Generative Lexicon (Pustejovsky, 1995). Generative Lexicon, and its developments, represents a device to model and deal both with classical lexical semantic relations, like merological relations, synonymy and others, and also with encyclopedic knowledge and even some kinds of discourse relations. The use of this lexical theory to retrieve the *R* relation responsible for the building of the bridge between the anaphoric element and its anchor will broaden the view of bridging anaphora resolution as a general problem of how much of background knowledge can be coded as part of the meaning of linguistic constituents. In the next sections, after having introduced SIMPLE, we will present the results of the performance of a semi-automatic algorithm for resolving

bridging anaphors which uses SIMPLE as its knowledge base.

4 *SIMPLE: a Generative Lexicon Resource for Italian*

The SIMPLE lexicon⁶ is a four-layered⁷ computational lexicon developed under two EU-sponsored project (PAROLE and SIMPLE) and extended under the Italian government founded project CLIPS. It represents the largest computational lexical knowledge base of Italian language, containing over 45 thousand lemmas and more that 57 thousand word senses, or semantic units.

At the semantic layer of information, lexical units are structured in terms of a semantic type system and are characterized and interconnected by means of a rich set of semantic features and relations. Combining both top-down and bottom-up approaches, the SIMPLE ontology has been elaborated in such a way as to permit an exhaustive characterization of different levels of complexity of lexical meanings.

The SIMPLE type system reflects the G.L. assumption that lexical items are multidimensional entities which present various degrees of internal complexity and thus call for a lexical semantic description able to account for different ranges of meaning components. Accordingly, a semantic type is not simply a label to be associated to a word meaning, it is rather the repository of a structured set of semantic information. Therefore, the membership of a word sense in a semantic type inherently triggers the instantiation of a rich bundle of semantic features and relations that represent the type-defining information that intrinsically characterizes the ontological type.

The core of the SIMPLE semantic relations rely on the *Qualia Structure*, which is one of the four representational level proposed by the G.L. framework. Qualia structure consists of four roles (Agentive, Telic, Formal and Constitutive) encoding the multifaceted nature of word meaning. Qualia relations enable capturing orthogonal relations existing between semantic units, regardless of their ontological classification. Querying the whole set of semantic relations in which a single keyword is involved throughout the

⁶http://www.ilc.cnr.it/clips/CLIPS_ENGLISH.htm

⁷Phonological, morphological, syntactic and semantic levels.

lexicon allows retrieving and extracting a set of semantic units belonging to different semantic types forming a semantic network. Moreover, qualia relations enable to establish a connection between a word sense and a number of events or entities strictly related to its meaning and to define the role of those events/entities in the lexical semantics of the word itself. In SIMPLE a revision of the original qualia structure was undertaken which led to the design of the *Extended Qualia Structure* whereby each of the four roles subsumes a set of semantic relations. Sixty extended qualia relations were therefore created, which allow to model the componential aspect of a word's meaning and to structure its relationships to other lexical units, on both the paradigmatic and syntagmatic axes.

However, the semantic relations are not exhausted by the (extended) qualia structure. Each semantic unit has three more relations such as synonymy, derivation, which allows a further type of connection between lexical items, and regular polisemy.

4.1 Exploiting qualia relations to resolve bridging anaphors

The core of our proposal is based on the idea that the qualia relations encoded in SIMPLE can be used to represent the *R* relations between a bridging element and its antecedent. To illustrate how to exploit qualia consider the examples from 2 to 7, all extracted from our corpus, which can only be resolved by making use of non-classical semantic relations; the anchor is in italics, the bridging element in bold and, in capital letters, the processing requirements (i.e. the *R* relations) needed to resolve the anaphoric link:

- (2) *i prezzi* – **al consumatore** [*the prices* – **the customer**]; INFERENCE
- (3) *il processo* – **gli imputati** [*the trial* – **the convicted**]; INFERENCE
- (4) *essersi sparato* – **il suicidio** [*to shoot himself* – **the suicide**]; EVENT
- (5) *fatto esplodere* – **le macerie** [*exploded* – **the debris**]; EVENT
- (6) *condannare* – **il pubblico ministero** [*to condemn* – *the attorney*]; EVENT

- (7) *il voto* – **l'elezione** [*the vote* – **the election**]; RHET. RELATION

The use of a G.L. approach allows us to claim that the *R* relations to resolve these cases of bridging are already encoded in the meanings of the lexical items themselves. Thus, for instance, in 3, the fact that a trial involves a convicted is formalized by exploiting a qualia relation between the two words, namely the constitutive “*member_of*”. In 7, the fact that if there is a vote, then there is an election (*cause/consequence*), can be formalized by exploiting the extended telic quale “*purpose*”. Moreover, bridging relations which take as anchor a verb (examples 4, 5 and 6) could as well be resolved by exploiting the extended qualia in SIMPLE. For instance, in 5, the FDNP *le macerie* can be resolved by exploiting the extended agentive quale “*result_of*”. It is quite trivial to remark that bridging relations classified as Lexical can be easily resolved as well by means of the qualia structure, including both classical lexical semantic relations and more fine-grained ones, like the one illustrated in 8, where the *R* relation can be expressed by the telic quale “*is_the_activity_of*”:

- (8) *l'attentato* – **i terroristi** [*the attack* – **the terrorists**]; LEXICAL

Before presenting the experimental data, another remark is necessary. The use of SIMPLE qualia relations has the further advantage of making explicit also what is the semantic relation which connects the bridging element to its antecedent, thus overcoming the shortcomings of machine learning approaches like Market, Nissim, and Modjeska (2003), which remain silent on this issue, i.e. do not specify *what is* the relation between the bridging anaphor and its antecedent.

5 Preliminary Experiments and Evaluation

To evaluate the reliability of the resource we have conducted an experiment on a subset⁸ of 129 bridging anaphors from our corpus. We have developed a semi-automatic procedure to query the resource. The workflow is the following: we manually provided to the system both the bridging anaphor and its an-

⁸All bridging relations which involved either as anchors or anaphoric elements named entities have been eliminated (144/299 - 48.16%), as well as those for the Discourse Topic class.

tecedent. The system, then, looks for a semantic relation between the two, either by looking for a direct connection between the two words, i.e. semantic units, or by looking for a common semantic type between the two entities. If more than a semantic relation between the two words is identified, the one with the shortest lexical distance (i.e. the one with the shortest semantic path) is selected. In case that more than a semantic relations with same lexical distance between the anaphor and the anchor is identify, both relations are considered as valid. This choice is a device to reflect the fact that even human beings when resolving bridging anaphors may agree on the anchor, but disagree on the type of relation, i.e. allow more than one relation. The maximum number of arcs allowed has been set to two. This is due to the fact that a wider range would result into inappropriate relations since the two semantic units may be linked at a very abstract level.

In order to verify our claim that a G.L. based resource should perform better in resolving bridging anaphors respect to WordNet-like ones, we have performed a comparative evaluation (by applying the same procedure) using ItalWordNet (IWN). In Table 2 we report the overall results of the two resources in terms of matching an existing semantic relation for the 129 couples of bridging anaphors and anchor, which corresponds to the number of possible bridging anaphors which could be resolved using these resources. The results are not very good, since only 22

Lexical Resource	Bridging
SIMPLE	22 (17.05%)
IWN	19 (14.72%)

Table 2: Numbers of correctly matched bridging anaphors.

couples of anchor-bridging anaphor can be resolved by using SIMPLE, a figure which is not so bigger than those which can be resolved by using IWN. The very low results are essentially due to (unexpected) missing relations and lexical entries in the SIMPLE resource. The low values for IWN are due to the absence of the necessary semantic relations, as expected and in compliance with its theoretical background. It is also interesting to notice that of the 19 correct relations which can

be retrieved by using IWN, only 11 of them cannot be identified by SIMPLE and this is due to missing information in the resource (5 over 11 couples cannot be identified because the proper semantic relations have not been introduced by the compilers of the resource) and not to theoretical shortcomings of the resource itself. Moreover, 13 of the 22 relations identified by using SIMPLE are completely out of reach for IWN, since they correspond to extended qualia.

Going into the details of the various subclasses of bridging relations the results are quite encouraging. What emerges is that the two resources can be thought as being specialized for the identification of particular subclasses of bridging anaphors. As the data in Table 3 show there is a relative high competition only for the subclass of Lexical bridging. The relative high performance of IWN in Inferential subclass is attributable to an extension of its original semantic relations as proposed by the EuroWordNet Project, of which IWN is a part. However, it is interesting to notice that all 5 Inferential bridging retrieved with IWN are identified by SIMPLE as well. The same observations hold for the class of Event as well. Finally, it is interesting to point out the fact that the subclasses of Rhetorical Relation and Inferential in SIMPLE are mainly resolved by two types of qualia (and their extensions) that is *Constitutive* and *Telic*.

Subclass	SIMPLE	IWN
Lexical	11 (50%)	12 (63.2%)
Inferential	7 (31.82%)	5 (26.31%)
Rhet. Relation	2 (9.09%)	0 (0%)
Event	2 (9.09%)	2 (10.52%)

Table 3: Subclasses of bridging matched.

6 Conclusion

The approach we have proposed is still a work-in progress and more refinements are needed. Of course a large-scale evaluation is compelling in order to provide further evidences of our proposal and a better evaluation of the SIMPLE lexicon. However, we would like to point out and emphasize some interesting aspects of this proposal:

- the use of a G.L. based resource can be seen as a way of reducing the influence

of extralinguistic knowledge;

- bridging can be used as a way of discovering semantic relations among linguistic entities and can be used to improve both the creation and maintenance of linguistic resources like SIMPLE. In particular, G.L. pattern induction from a corpus-based study can improve the resource by adding missing relations;
- the problem of bridging anaphora resolution becomes part of a more general problem of identification of semantic relations between linguistic elements;
- a resource with G.L. qualia relations encoded in it should not be compared with a world-knowledge database or similar (effort expensive and difficult) resources. G.L.-based relations are dynamic, in the sense that they allow to discover new relations between lexical items and can provide an account for the creative use of language;
- qualia relations can represent new features for machine learning approaches; considering an annotation task for anaphora resolution, it would be very useful to introduce a new attribute which expresses the qualia relation between the anchor and the anaphoric element, thus providing information to a learner to resolve also difficult (i.e. non strictly lexical) cases of bridging anaphors.

The results obtained are not very satisfying and seem to support criticisms to the use of lexical resources in tasks of anaphora resolution. We agree on some of this criticism, but we would like to point out that the resolution of bridging anaphors is not a trivial task and the use of lexical resources like SIMPLE can represent a useful strategy for the development of robust algorithms for anaphora resolution. As for SIMPLE an extended work of revision and correction of the various mistakes and missing elements is compelling in order to be used reliably. A further point which emerges from this work is represented by the observation that SIMPLE and IWN are not competitive resources, i.e. one being the extension of the other, but more complementary ones. The final proposal we suggest is a call for a new generation of lexical resources. Resources whose scope is that

of being specialized in restricted sets of lexical relations. This could result in better resources with less mistakes and missing information and easier to be integrated in NLP algorithms.

References

- Caselli, T. 2007. An annotation scheme for bridging anaphors and its evaluation. In Andrea Sansò, editor, *Language Resources and Linguistic Theory*, volume 59 of *Materiali Linguistici*. Franco Angeli, Milano, pages 149–166.
- Chierchia, G. 1995. *Dynamics of Meaning: anaphora, presuppositions and the Theory of Grammar*. University of Chicago Press, Chicago.
- Clark, H. 1997. Bridging. In P.N. Johnson-Laird and P.C. Wason, editors, *Thinking: Readings in Cognitive Science*. Cambridge University Press, Cambridge and London.
- Hawkins, J.A. 1978. *Definiteness and Indefiniteness*. Croom Helm, London.
- Kleiber, G. 1999. Associative anaphora and part-whole relationship: the condition of alienation and the principle of ontological congruence. *Journal of Pragmatics*, 31:339–362.
- Krahmer, E. and P. Piwek. 2000. Varieties of Anaphora. Course Notes, ESSLLI00, Birmingham, August 11-23.
- Lobner, S. 1985. Definites. *Journal of Semantics*, 4:297–326.
- Market, K., M. Nissim, and N. Modjeska. 2003. Using the Web for nominal anaphora resolution. In *EACL Workshop on the Computational Treatment of Anaphora*.
- Montemagni, S., F. Barsotti, M. Battista, N. Calzolari, O. Corazzari, A. Lenci, V. Pirelli, A. Zampolli, F. Fanciulli, M. Massetani, R. Raffaelli, R. Basili, M. T. Pazienza, D. Saracino, F. Zanzotto, N. Mana, F. Pianesi, and R. Delmonte. 2003. The syntactic-semantic Treebank of Italian. An Overview. *Linguistica Computazionale, Computational Linguistics in Pisa, special Issue, XVI-XVII*:461–493.
- Pustejovsky, J. 1995. *The Generative Lexicon*. MIT Press, Cambridge, MA, USA.

- Ruimy, N., M. Monachini, E. Gola, A. Spanu, N. Calzolari, M.C. Del Fiorentino, M. Olivieri, and S. Rossi. 2003. A computational semantic lexicon of Italian: SIMPLE. *Linguistica Computazionale, Computational Linguistics in Pisa, special Issue*, XVI-XVII:821–864.
- Sidner, C.L. 1979. *Towards a computational theory of definite anaphora comprehension in English discourse*. Ph.D. thesis, MIT.
- Strand, K. 1997. A taxonomy of Linking Relations. Manuscript.
- Vieira, R. and M. Poesio. 2000. An Empirically-Based System for Processing FDNPs. *Computational Linguistics*, 26(4):539–593.