

Direct Object Omission as a Sign of Conceptual Defaultness

Louis Hickman, Julia Taylor, and Victor Raskin

Purdue University, West Lafayette, IN, USA
{lchickma, jtaylor1, vraskin}@purdue.edu

Abstract

This paper is another step towards the recognition of conceptual defaults based on the presence or absence of modification for the nominal direct object of semi-transitive verb events. A default is defined as any piece of information that is obvious enough to be omitted in every day communication or within a domain of interest familiar to the speakers. The insufficiency of a seemingly reliable non-semantic algorithm is exposed and corrective semantic measures discussed.

Introduction

Current military, governmental, corporate, and societal needs increasingly require applications based on computer understanding of information in natural language texts. Machine learning and other statistical methods that avoid direct representation of meaning have satisfied this need indirectly and to a limited extent. Attempts have intensified to base semantic technologies on direct representation of meaning/knowledge, as exemplified by recent work on semantic resources(see, for instance, Palmer et al. 2005; Ruppenhofer et al. 2010; Schular et al. 2009; Raskin et al. 2010a; Taylor et al. 2010a; Cimiano et al. 2014) that have made considerable strides in representing explicit meaning of each sentence. What has not yet been addressed directly, let alone comprehensively, is the information that is inferred by humans from every sentence but never mentioned explicitly in the text, that is, implicit meaning.

In this paper, we are addressing one simple element of implicit meaning, namely, something so obvious that it is commonly omitted from being mentioned in the text. In general, omitted words characterize ellipsis. Semantic ellipsis has been addressed rather rarely (see Baltes 1995—cf. the more basic and syntactic McShane 2005), even though ellipsis occurs quite commonly in writing and much more so in speech. Semantic ellipsis often indicates the failure to provide details, as in *Bob was killed a couple of*

years ago, and more detailed information on the killer may vary between people, natural disasters, wild animals, avalanche and such. On the other hand, *Bob was murdered/assassinated...* limits the perpetrator to humans, even if a venomous snake was used as an instrument.

Different from semantic ellipsis, lies default – something that is understood by people without an associated range of missing details. The defaults are omitted not because speakers fail to provide details per se, but because these details, if provided, would be unnecessary and, perhaps, awkward. For example, *John was eating* is more appropriate that *John was eating food* because humans and animals always eat food, so food does not need to be mentioned. In fact, mentioning it modifies the meaning to something atypical, for instance, to emphasize the difference between home-cooked by Grandma and fast food, in which case the modifying adjective *real* is strongly implied. Note that *John was eating some food* or *Thai food* is perfectly okay, so modification makes the verbalization perfectly normal.

The goal of this project is to enable computer understanding of natural language at the level of human comprehension, relative to the omitted but obvious and implicitly present information in independent sentences—that is, sentences without preceding text. This means that omitted information that could be recovered in previous text (as in ellipsis, coreference, etc.) is not of interest here. What is of interest is the algorithm of recovering common knowledge, a special case of it. We suspect that such algorithm can be useful to understanding what is common and acceptable for any given human. Thus, it could be useful not only in communication between robots and humans, but also in information security (for example, for insider threat or social engineering detection—cf. Raskin et al. 2010b, Taylor et al. 2010b) or in health-care communication under the purview of the Health Insurance Portability and Accountability Act (HIPAA).

Background

Mutual knowledge and shared experience are topics most frequently covered in the areas of psychology and perhaps less so in most of linguistics. Unfortunately, little work has been done on implementing access to and use of mutual knowledge in computation.

Mutual knowledge “assumes that listeners use the knowledge and beliefs they share with speakers in the process of interpreting utterances” (Gibbs, 1987). One of the explanations for the omitted information is that assumptions made between speaker and listener allow humans to skip the infinite exchange of background knowledge in order to draw conclusions more easily (Clark & Marshall, 1981). Moreover, Grice’s (1975) Cooperative Principle is a series of maxims that outline the rules by which humans seek to communicate information. In the Cooperative Principle, both speakers and listeners are said to make as much of a contribution to the conversation as is required. Both parties shall also make the contribution of information when it is necessary within the conversation. Sperber and Wilson (1994) claim that humans focus most on what is relevant in a conversation. As such, humans tend to both state and interpret things in a way that maximizes relevance.

Research on referring expressions is one area in which mutual knowledge is somewhat redefined following the work of Grice. Mainly, the implementations have focused on natural language generation of referring expressions (Dale & Reiter 1995; Appelt & Kronfield 1987; Dale 1992; Viethen & Dale 2006). Outside of ontologies, some research has been done in the area of databases to represent knowledge and use it to create richer queries. For instance, Feldman and Hirsh (1997) created a system that examines keyword labels in text documents. The system views background knowledge as constraints to a query.

A lot of research has been done on automatic acquisition of ontologies. In most cases, the researchers are interested in recovering as much information from text as possible to construct new ontological concepts or properties between existing concepts. This is loosely relevant to our purpose as the omitted information we get will be entered into the ontology, not for the purpose of the ontology building, but as data that could be retrieved later for better understanding of text. What is important here is that the omitted information is usually generic and, in many instances, is applicable to many situations. While machine learning has been used for ontology acquisition *per se*, we are not aware of any machine learning applications to this kind of problem—when the needed information is not contained in text.

The general acknowledgement of the existence of mutual knowledge exchange in combination with the varying views of how this information is obtained

demonstrates a need for research on computational implementation. According to Krauss and Fussell (1990), people have been asking how to address common background knowledge in a global workforce for years. As technology has continued to improve over the years, the number of virtual teams within a company has also increased. Thus, there is a need to collect mutual knowledge within a particular domain. Cramton (2001) notes that some of the consequences resulting from inconsistencies in background knowledge include the hesitation of individuals on a team to mention relevant and unique information assumed known by others, bringing about the rapid deterioration of working relationships. The situation can only get worse in human-robot communication.

Raskin et al (2010b) discusses the potential for verbalized information that should have been omitted within Information Security and specifically for insider threat detection. Insider threat has many definitions but essentially refers to “a breach of trust by people within an organization or system” (Bishop et al 2008). Verbalized/omitted information switching appears to be particularly useful for the detection of lies and unintentional inference which pairs well with insider threat detection. As lying is not a form of bona-fide communication, a person generally violates defaults in some way when they lie. Whether intentional or unintentional, violations can be identified if the “common” knowledge of the individual is recovered.

Non-Semantic Work on Defaults

Defaults, as described by Taylor et al. (2010b), refer to that information which is assumed to be known and is no longer salient to the speaker. Because this information is no longer salient to the speaker, it is not brought up in conversation. However, that unspoken piece of information is necessary for understanding the meaning of a statement.

We are interested in a situation in which the use of the unmodified noun—a noun phrase (NP) consisting of a single noun—is inappropriate as a direct object (considered stating the obvious) while, with a modifier for that noun, the NP and the phrase in which it occurs lose their default status. For example, the object represented by *words* is defaulted in *He wrote words* but removed from its default status by modifiers in *He wrote 5,000 words* and *He wrote polysyllabic words*.

Ringenberg (2015) conducted an experiment with the goal of extracting candidate direct and indirect object defaults (without any semantic knowledge) from unstructured text by examining the relationships between verb and the objects that relate to them. Specifically, the experiment compiled and examined events that occur as

verbs, the modifier-noun, and noun combinations that are associated with them within a verb phrase. These events, modifiers, and nouns are used to identify when information about an event is both stated and omitted.

The datasets used were Brown corpus and Wikipedia for Schools. Wikipedia for Schools is a collection of over 6,000 Wikipedia documents with more than 26 million (non-unique) words, which all pertain to the subjects taught in United Kingdom curriculum. The 200 most frequent verbs in Brown corpus were selected. From those, verbs that could not have defaults were removed. The number of instances in both corpora is shown in Table 1.

Metric	Brown Corpus	Wikipedia
Unique verbs chosen	141	141
Verb instances with no modifiers	1339	272216
Candidate instances	869	205774
Unique verbs with instances	106	141

Table 1: number of instances

A candidate default was identified as any noun that:

- Was a direct object of a verb in question, and
- Occurred when modified, and
- Never occurred as an unmodified direct object.

The results show 141 verbs had at least one candidate default. Table 2 shows verbs that had more than 5000 instances in the Wikipedia dataset.

Verb	Count of Candidate Default Instances	Count of Unique Candidates
include	21473	5740
Use	20187	4178
produce	7869	2112
Form	7763	1708
contain	7662	2034
Create	7510	2039
receive	6252	1163
Cause	5746	1566
Play	5343	843
develop	5044	1165

Table 2: Verbs with high number of default candidates

However, without semantic information, it was impossible to detect what was a default and what was a regular undefaulted direct object. Many verbs had more than one default, which correspond to different properties. Some instances were not as clear and we conjecture that they relate to the individual defaults rather than commonly shared defaults.

Working With Data

Analyzing News – the meaning of kill

To further illustrate the distinction between semantic ellipsis, on the one hand, and conceptual defaults, on the other, let us take real-life text from this day’s CNN online news:

Investigations into the series of terrorist attacks that killed more than 120 people in Paris are moving forward, with people taken into custody and two of the gun-wielding suicide bombers identified. ISIS claimed responsibility for the massacres in a statement. In response, France has carried out air strikes on targets in the militant organization’s stronghold in Raqqa, Syria.

Semantic ellipsis takes place whenever information not present in the current text, including its previous parts is evoked, and it is impossible to understand the text otherwise. In this text, the very first usage of the definite article *the* before *series* refers to a series of attacks never mentioned earlier because this sentence is the first one in a new text. True enough, *series of terrorist attacks that killed more than 120 people in Paris* is hypertext, and the link provides some of the information necessary for reconstructing the ellipsis. People taken into custody and terrorist identification must be understood as results of investigations, and this comes from the general knowledge of the world (even as gun-wielding dead people requires an amount of indulgence). An understanding of war, never mentioned in the text, as attacks and counterattacks has to be imputed also, as is, in fact, the knowledge that ISIS (ISIL, Daesh) is a terrorist organization claiming territory in Syria and Iraq.

There is only one case of default in the example. The phrasal verb *taken into custody* has *human* with quantity of one as default. However, due to the importance of informing the public when suspects are in custody, this is a situation where default is frequently violated. The defaultness is evident because one can say, *three are in custody*, and everyone understands that means three humans. Here there was, perhaps, uncertainty as to the number of people taken into custody and hence, the somewhat awkward, unspecific usage.

Analyzing Wikipedia – the meaning of play

Let us consider another example of a verb: *play*. *Play* is the only frequent verb that has more than 5000 instances of candidate defaults with less than 1000 unique candidates in Ringenber (2015).

Our goal is to define information about each meaning of the verb *play* that could be filled automatically by the system and used for the semantic candidate default

generation. In this section, we will look at 3 systems that describe different meanings of the verb *play*: FrameNet (Ruppenhofer et al. 2010), PropBank (Kingsbury and Palmer, 2002), and Ontological Semantics Technology (OST) (Raskin, Taylor, and Hempelmann, 2010b). We manually identify the meaning of the verb *play*, the meanings of its direct objects, the semantic roles between the verb meanings and the objects, as well as the meaning of the modifiers. We are thus interested in the meanings of *play* in each of the three systems.

Meaning, according to FrameNet

According to FrameNet, the verb *play* has 6 different senses, or can be used in 6 different frames, with the core frame elements described in the sub-bullets:

- Being_relevant (play into)
 - Cognizer
 - Endeavor
 - Phenomenon
- Compliance (play by the rules)
 - Act
 - Norm
 - Protagonist
 - State_of_affairs
- Performers_and_roles
 - Audience
 - Medium
 - Performance
 - Performer (multiple)
 - Role
 - Score
 - Script
 - Type
- Competition
 - Competition
 - Participant (multiple)
- Cause_to_make_noise
 - Agent
 - Cause
 - Sound-maker
- Make_noise
 - Noisy_event
 - Sound
 - Sound_source

Meaning, according to PropBank

PropBank recognizes the following meanings of *play*, with the roles of each frame listed as sub-bullets:

- Play a game
 - Player
 - Game
 - Instrument/equipment
 - Opponent, play against whom
- Play a role
 - Actor

- Role
- Play into, be a factor
 - Thing factoring in
 - Thing being factored into
- Play a trick on someone
 - Trickster
 - Mention of trick
 - Tricked, who trick was played on
- Perform music
 - Performer, player
 - Thing performed
 - Musical instrument/style

Meaning, according to OST

- Play (in sports) Anchoring concept: sports-event
 - Agent: human
 - Instrument: sports-object
 - Place: sport-facility
 - Theme: game-rule
- Play (music) Anchoring-concept: make-noise
 - Agent: human
 - Beneficiary: animal
 - Instrument: musical-instrument
 - Theme: piece-of-music
- Play (theater) Anchoring-concept: portray
 - Agent: human
 - Beneficiary: human
 - Theme: entertainment-role
 - Place: theater
- Play (initiate playback of recorded media) Anchoring-concept: media-playback
 - Agent: human
 - Theme: media-object
 - Instrument: electronics
- Play (not work) Anchoring-concept: work
 - Epistemological modality: 0
 - Agent: human
- Play: Anchoring-concept: work
 - Agent: human
 - Manner: life-role

Meaning of the verb *play* in Wikipedia

Semantic analysis of sentences from the Wikipedia dataset used by Ringenber reveal that meaning corresponding to playing music occurs with unique direct objects 172 times (474 instances). The meaning corresponding to performing/acting occurs with unique direct objects 323 times (425 instances). The meaning corresponding to participation in a game uniquely occurs with unique direct objects 209 times (997 instances). The meaning corresponding to media playback occurs with unique direct objects 21 times (35 instances). Other sentences corresponded to infrequent sentences or used idiomatic expressions.

Semantic architecture/meaning of play		to play an instrument	to perform or act	to participate in a game
FrameNet	<i>Frame name:</i>	<i>cause_to_make_noise</i>	<i>performers_and_roles</i>	<i>Competition</i>
	slot:	sound-maker	role	comp
	filler (unique/total number)	38/135	278/355	64/743
PropBank	<i>Roleset:</i>	<i>play.11</i>	<i>play.2</i>	<i>play.01</i>
	arg:	thing performed	role	game
	filler (unique/total number)	71/214	278/355	64/743
	arg:	instrument/style		
OST	<i>Event:</i>	<i>make-noise</i>	<i>pretend</i>	<i>sports-event</i>
	property:	Instrument	theme	N/A
	filler (unique/total number)	38/135	278/355	47/722
	property:	Manner		instrument
	filler (unique/total number)	71/214		8/19

Table 3: Semantic default candidates

By further analyzing the direct objects based not by their words but solely by their meaning, we enhance the analysis. Data was examined both at its current grain size of semantic detail as well as at larger grain-sizes (going ‘up’ in the ontology, as it were) to identify common ontological ancestors. Thus we can identify potential semantic defaults at the grain size in which they exist.

It is noteworthy that these meanings can be encoded by each of the systems described above and all of them have the necessary properties/slots to fill the necessary information (see Table 3). Upon further analysis, the meaning playing music has the following frequent meanings from the candidate object defaults: music-instrument, music-piece, music-event, and music-object (such as rhythm, note, scale, etc).

The meaning of performing/acting has only one frequent meaning of direct object – that of a human, albeit of a different grain size. The common performances are of fictional characters, professional roles such as artist or ballerina, and named persons such as Bach or Kennedy.

The meaning of participation in a game has 64 different sports-events mentioned with the total number of 743 occurrences. Thus, these sports-events (such as soccer, match, and Wimbledon) could play a role of default as well. When sports-events are the object of play, OST actually processes this differently than the other systems. Specifically, it replaces the anchoring concept with the finer grain size of the specific sport. OST also has different numbers than the rest of the systems for playing games because OST draws a distinction between playing sports versus types of games. Thus, the reduction in quantity for OST correlates exactly with the instances relating to card games, board games, and the like.

There exists a previously unmentioned meaning, that of contributing to an event. Approximately 1,000 instances occur with the direct object *role* and almost 200 have *part*. The only other occurrence has object *influence*, which can roughly be replaced in the sentence by *part* or *role*. This is a phrasal with multiple wordings and is not covered here because it does not allow for a variety of direct objects (i.e., all possible objects have equivalent semantics).

And finally, the meaning of initiating playback of recorded media has as candidate default concept media-object. The difference between this sense and the playing an instrument sense is whether the agent is initiating playback through some electronic medium or if the agent is performing the piece themselves, either via an instrument or at a concert. All objects that only occurred modified were some form of recorded media.

Discussion and Conclusion

The results show that the syntactic property of modification and lack thereof reliably follows from defaultness: if the concept anchoring a word is the default for a certain sense of a verb then it will be omitted, except perhaps in a number of exceptional situations when used without any modification. One example of including the default is in questions requesting a finer grain size to be specified, e.g., *what instrument do you play?* It is, apparently, not a sufficient condition to only occur modified as the verb object, and those are much harder to come by for natural language phenomena or other complex entities such as love, life, humor. Most of them are constitutive (see Searle 1969) in that they actually create a phenomenon the same way rules of chess create the game.

To access the constitutive sufficient conditions of defaultness, one needs to apply semantic rules. One of them, a pretty trivial one, follows immediately from OST and—somewhat less clearly—from other semantic resources: if a concept is the filler of an essential semantic property for an event, then any unmodified word anchored directly in this concept, is the default for the conceptual sense of the verb. This describes food for eat/ingest or text for write.

In a formal semantic system like OST, the ontology is accessed algorithmically in the course of routine semantic analysis. So, linguistic semantics can indeed be done formally, contrary to Chomsky's recently reiterated aberration (2015), and it should not be replaced by syntax in search of algorithmicity.

References

- Appelt, D. E., & Kronfeld, A. (1987). A computational model of referring. In *IJCAI 87*: 640-647)
- Baltes, P. 1995. Discourse Reduction and Ellipsis: A Semantic Theory of Interpretation and Recovery. Unpublished Ph.D. thesis, Department of English, Purdue University, West Lafayette, IN.
- Bishop, M., Engle, S., Peisert, S., Whalen, S., Gates, C., Probst, C.W. & Somayaji, A.. (2008). We have met the enemy and he is us. *New Security Paradigms Proceedings of the 2008 Workshop*, 1-12.
- Cimiano, P., Unger, C., & McCrae, J. (2014). *Ontology-Based Interpretation of Natural Language*. Morgan & Claypool Publishers.
- Chomsky, N. (2015). Semantics. BISC-Group, 10/14.
- Clark, H. H., & Marshall, C. R. (1981). Definite reference and mutual knowledge. In A. K. Koshi, B. Webber, & I. A. Sag (eds.) *Elements of discourse understanding*. Cambridge, Cambridge University Press, pp. 10-63.
- Cramton, C. D. (2001). The mutual knowledge problem and its consequences for dispersed collaboration. *Organization Science*, 12(3), 346-371.
- Dale, R. (1992). *Generating referring expressions: Building descriptions in a domain of objects and processes*. Cambridge, MA: MIT Press.
- Dale, R., & Reiter, E. (1995). Computational Interpretations of the Gricean Maxims in the Generation of Referring Expressions. *Cognitive Science* 19: 233-263.
- Feldman, R., & Hirsh, H. (1997). Exploiting background information in knowledge discovery from text. *Journal of Intelligent Information Systems*, 9(1): 83-97
- Gibbs Jr., R. W. (1987). Mutual knowledge and the psychology of conversational inference. *Journal of Pragmatics* 11(5): 561-588
- Grice, H.P. (1975). Logic and conversation. In: Cole, P., and L.J. Morgan (eds.) *Syntax and Semantics.Vol.3. Speech Acts*. New York: Academic Press, pp. 41-58.
- Kingsbury, P., & Palmer, M. (2002, May). From TreeBank to PropBank. *LREC*, pp. 1989-1993.
- Krauss, R. M., & Fussell, S. R. (1990). Mutual knowledge and communicative effectiveness. *Intellectual teamwork: Social and technological foundations of cooperative work*, 111-146.
- McShane, M. J. (2005). *A Theory of Ellipsis*, New York: Oxford University Press.
- Palmer, M., Gildea, D., & Kingsbury, P. (2005). The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics* 31(1):71-106
- Raskin, V., Hempelmann, C. F., & Taylor, J. M. (2010a). Guessing vs. knowing: The two approaches to semantics in natural language processing. *Annual International Conference on Artificial Intelligence Dialogue 2010*, Bekasovo (Moscow), Russia.
- Raskin, V., Taylor, J. M., & Hempelmann, C. F. (2010b). Ontological semantic technology for detecting insider threat and social engineering. *Proceedings of the 2010 workshop on NewSecurityParadigms*, 115-128). ACM.
- Ringenberg, T. (2015). *Creating, testing and implementing a method for retrieving conversational inference with ontological semantics and defaults*. Unpublished MS Thesis, Computer and Information Technology, Purdue University.
- Ruppenhofer, J., Ellsworth, M., Petruck, M. R. L., Johnson, C. R., & Scheffczyk, J. (2010). FrameNET II: Extended Theory and Practice. <https://framenet2.icsi.berkeley.edu/docs/r1.5/book.pdf>.
- Schular, K., Korhonen, A., & Brown, S. W. (2009). VerbNet overview, extensions, mappings and apps, Tutorial, NAACL-HLT 2009, Boulder, Colorado.
- Searle, J. R. (1969). *Speech Acts*. Cambridge, UK: Cambridge University Press.
- Sperber, D., & Wilson, D. (1994). Outline of relevance theory. *Links and Letters*, 85-106.
- Taylor, J. M., Hempelmann, C., & Raskin, V. (2010a). On an automatic acquisition toolbox for ontologies and lexicons in ontological semantics. *ICAL*, 863-869).
- Taylor, J. M., Raskin, V., Hempelmann, C., & Attardo, S. (2010b). An unintentional inference and ontological property defaults. In *Systems Man and Cybernetics (SMC), 2010 IEEE International Conference*, 3333-3339). IEEE.
- Viethen, J., & Dale, R. (2006). Algorithms for generating referring expressions: do they do what people do?. *Proceedings of the Fourth International Natural Language Generation Conference*, (3-70). Association for Computational Linguistics.