

Annotation Guidelines for Citations in Scientific Literature

November 23, 2011

1 Introduction

These guidelines outline an annotation scheme for citations in scientific literature. Citations are used in the text of scientific literature to refer to other sources (most often they refer to other published literature). For example, in the sentence below, **(Oviatt 1996)** is a citation that points to another paper.

- (1) Multimodal systems provide a natural and effective way for users to interact with computers through multiple modalities such as speech, gesture, and gaze **(Oviatt 1996)**.

For this annotation scheme we would like to consider two aspects of the citation: (1) what is the author saying about the quality of the cited work, and (2) what is the relationship of the citing work to the cited work, i.e., how is the author using the cited work.

With this in mind, we would like to annotate each citation along four dimensions or facets (taken from the classification scheme by Moravcsik and Murgesan (1975)):

- conceptual vs operational
 - *Is this an idea or a tool?*
- evolutionary vs juxtapositional
 - *Is the author building on the cited work or working in contrast to it?*
- organic vs perfunctory
 - *Is this **particular** citation necessary for understanding the paper or can the paper still be understood without it?*
- confirmative vs negational
 - *Is the cited work correct or are there some limitations to it?*

The first two dimensions correspond to the *utility* of the cited work and the last two dimensions relate to the *quality* of the cited work.

For our annotation, all citations should be completely defined, i.e., no facets left undefined. Finally, a note on terminology, **author** will be used in the guidelines to describe either the citing paper or the authors of the citing paper, i.e. that paper that makes reference to another paper. We will then use **cited work** for either the cited paper or authors of the cited papers.

2 Conceptual vs Operational

Generally, if the citation refers to the use of some tool or resource it should be labeled **operational**, otherwise if it is an idea or algorithm it should be labeled **conceptual**. Some examples of tools and resources in NLP might include: taggers, parsers, stemmers, classifiers, or corpora.

Note that there are often cases when a paper has both a conceptual and operational component. Be careful to annotate this accurately for each citation in the case where the work is cited more than once.

2.1 Operational

Label the citation **operational** if:

- the citation refers to the use of a tool (e.g. tagger, parser, stemmer, etc.), a corpus, etc.
- (2) However, most of the existing models have been developed for English and trained on the Penn Treebank (Marcus et al., 1993)

2.2 Conceptual

Otherwise, for example if the citation refers to an idea or algorithm, label **conceptual**. Some specific examples might be citations to theories, algorithms, or any abstract concept found in the cited work.

- (3) Context is typically treated as a set of unordered words, although in some cases syntactic information is taken into account (Lin, 1998; Grefenstette, 1994; Lee, 1999).

Also in the case that the author refers to implementing the cited work, use the label **conceptual**.

2.3 Possibly tricky examples?

- (4) More specifically, we combine a probabilistic topological field parser for German (Becker and Frank, 2002) with the HPSG parser of (Callmeier, 2000). <OPER >

- (5) Various parsing techniques have been developed for lexicalized grammars such as Lexicalized Tree Adjoining Grammar (LTAG) (Schabes et al., 1988), and Head-Driven Phrase Structure Grammar (HPSG) (Pollard and Sag, 1994). <CONCEPT >
- (6) In the following decade, great success in terms of parse disambiguation and even language modeling was achieved by various lexicalized PCFG models (Magerman, 1995; Charniak, 1997; Collins, 1999; Charniak, 2000; Charniak, 2001). <CONCEPT >
- (7) Table 2 compares the results of our algorithm with the results in (Och and Ney, 2000), where an HMM model is used to bootstrap IBM Model 4.<OPER >

In (7), note that the citation to “results” refers to results from a tool and are therefore labeled *operational*.

3 Evolutionary vs Juxtapositional

We define *evolutionary* to be any citation that is compatible with what is being claimed by the author, and *juxtapositional* is any citation that contradicts or contrasts the claims of the author.

Again, a cited work may be cited in one context as evolutionary and in another context as juxtapositional. For example, in a discussion of using machine learning for predicting pitch accent, one citation context may describe the common problem of predicting pitch accent, which is labeled *evolutionary*, and a second citation may describe the different machine learning approach used in the cited work, which is then labeled *juxtapositional*.

3.1 Juxtapositional

- If the author proposes an alternative to the cited work, label *juxtapositional*.
 - (8) Our approach differs from Lin (1998) in three important ways: (a) by introducing dependency paths...
- If there is any contrastive or juxtapositional element in the citation then label it *juxtapositional*.
 - (9) **Alshawi et al. (2000)** also presented a two-level arranged word ordering and chunk ordering by a hierarchically organized collection of finite state transducers. The main difference from our work is that their approach is basically deterministic, while the chunk-based translation model is non-deterministic. The former method, of course, performs more efficient decoding but requires stronger heuristics to generate a set of transducers. Although the latter approach demands a large amount of decoding time and hypothesis

space, it can operate on a very broad-coverage corpus with appropriate translation modeling.

3.2 Evolutionary

If the cited work is the basis of the author’s work, is used in the author’s work, or even if the cited work is compatible with what is being claimed by the author, we define the citation as **evolutionary**. Below are listed some typical instances where the citation should be labeled **evolutionary**. This is not, however, an exhaustive list, i.e. **evolutionary** instances are not limited to the conditions listed below.

- If the citation is used or even compatible with citing work, mark as **evolutionary**.

(10) we follow Ennis and Bi (1998) and use the identities

- If the citation refers to an agreed upon definition, term, or metric, label as **evolutionary**. For example, in example (11), although the BLEU score is not being extended, just by using it we assume it is an endorsement of the metric.

(11) We utilize BLEU (Papineni et al., 2002) for the automatic evaluation of MT quality in this paper.

- If the citation discusses a shared problem, label as **evolutionary**. This context should be labeled **evolutionary**, even if later in the paper there is a separate citation that discusses how the author distinguishes itself from the cited work.

(12) Information Extraction (IE) is the process of identifying events or actions of interest and their participating entities from a text. As the field of IE has developed, the focus of study has moved towards automatic knowledge acquisition for information extraction, including domain-specific lexicons (Riloff, 1993; Riloff and Jones, 1999) and extraction patterns (Riloff, 1996; Yangarber et al., 2000; Sudo et al., 2001).

- If it is a tool (i.e., labeled **operational**), then label **evolutionary** if the tool is simply being used (i.e., when the author uses a particular tagger) or if a series of third-party tools are being compared (i.e., a review of several different taggers). However, label as **juxtapositional** if the cited tool is being directly compared to the author’s tool being presented.

4 Organic vs Perfunctory

Generally, *organic* citations will be those that are very important for understanding the author’s work. These can be citations that form the basis of the

author's work, or any citations without which the paper would not make sense, or citations to otherwise unique work that cannot be referred to with any other citation. *Perfunctory* citations on the other hand are citations used to point to related literature, work, or authors that are not necessarily essential to understanding the author's paper.

4.1 Perfunctory

Label *perfunctory* if:

- the citation could easily be replaced by another citation (or removed altogether) and the general point could still be understood and make sense.
 - (13) Vector spaces enjoy widespread use in information retrieval (Salton and McGill, 1983; Baeza-Yates and Ribiero-Neto, 1999)...
- the citation is in a list of citations (i.e., the citation could be replaced or omitted). This is the case for explicit lists like example (14) or implicit lists like example (15). One exception to this rule may be if all of the cited work in the list has at least one common author, so that it could be considered a unique work that spans several papers. In this case the citations may be labeled *organic* (see example 18).
 - (14) Corpus-based methods and machine learning techniques have been applied to anaphora resolution in written text with considerable success (Soon et al., 2001; Ng & Cardie, 2002, among others).
 - (15) The utilization of language technology for the creation of hyperlinks has a long history (e.g., Allen et al., 1993).
- the citation comes at the end of a sentence without being specifically referred to in the text and with no further explanation. The assumption being that that citation justifies the statement.
 - (16) Even narrow-coverage context-free natural language grammars produce explosive ambiguity (**Church and Patil, 1982**).
- the citation refers to details in another paper (usually by the author)
 - (17) See Baldwin and Bond (2003) for further details.
- the citation is to a tool (i.e. a citation to a tagger or parser that could be replaced by using another tagger or parser).

4.2 Organic

If the citation refers to important work or work that is uniquely necessary for making sense of the the author's work then it should be labeled *organic*. Examples of this may be when the author's work is based on or inspired the cited

work, when the cited work is fundamental in realization of the author’s work, or when the author cites something very specific that can only be found in that particular cited work.

Note it should also be labeled **organic** if:

- there is a list of citations with the same author (or one common author). (This is an exception to one of the *perfunctory* rules above.)

(18) STILL NEED TO FIND EXAMPLE OF THIS

Typical **organic** example:

- (19) In order to exploit syntactic dependencies in a larger context, we propose a new model of supertagging based on Sparse Network of Winnow (SNoW) (Roth, 1998).

5 Confirmative vs Negational

This facet is similar to the NLP task of *sentiment analysis*, which is basically determining if the author is describing something as positive or negative. However, sentiment is manifested differently in published scientific literature with respect to product reviews for example. We should reconsider what constitutes positive and negative language in scientific literature and keep in mind that negative citations have been shown to be rare in published articles (Moravcsik and Murugesan, 1975).

Following the labels used by Moravcsik and Murugesan, we use **negational** to refer to negative citations and **confirmative** to refer to positive citations.

We will consider **negational** citations to be those where the author is critical of the cited work, highlights shortcomings or limitations of the cited work (and likely proposes solutions to it), or generally disagrees with the assertions in the cited work.

We will consider **confirmative** citations to be those where the author supports the cited work, highlights particular positive aspects of the cited approach, or generally agrees with or at least accepts the assertions in the cited work. We will also consider citations that do not seem to be positive or negative to be **confirmative**. This is because simply by citing the work we assume that the author agrees with it or thinks positively of it.

Note that different citations to the same paper can be assigned different labels. For example, a citation might be introduced and praised for its initial contribution and later criticised for its shortcomings. If this is nicely separated by having two citations, label each of them accordingly. If for one citation there is mixed positive and negative feedback, the annotator should label the citation as **negational**.

5.1 Negational

Label **negational** if:

- the citation is explicitly negative: illustrating major faults in the cited work’s methodology, results, or conclusions, etc.
- the author points out limitations in the cited work (and proposes alternative solutions). If these critical comments follow statements of praise, then a decision must be made by the annotator on the positivity/negativity of the citation. However, the annotator should lean towards **negational**.

(20) Various supervised learning methods for Named Entity (NE) tasks were successfully applied and have shown reasonably satisfiable performance.((Zhou and Su, 2002)(Borthwick et al., 1998)(Sassano and Utsuro, 2000)) However, most of these systems heavily rely on a tagged corpus for training. ...

- If the citation is marked **juxtapositional**, take care in labeling **confirmative/negational**. Mark **negational** if the the citing work fills a void or corrects something in the cited work. If the cited work is different and distinct enough (still **juxtapositional**) the citation might not necessitate a negative value and therefore be marked **confirmative**, i.e. if we are “comparing apples to oranges.”

For example, in (21), the author’s work and the cited work differ, but the objectives of each are distinct. The author doesn’t necessarily have any negative comments.

(21) This differs from the BioCreAtIvE competition tasks that aimed at classifying entities (gene products) into classes based on Gene Ontology (Ashburner et al., 2000).

Typical **negational** example:

- (22) Unlike well-known bootstrapping approaches (Yarowsky, 1995), EM and CE have the possible advantage of maintaining posteriors over hidden labels (or structure) throughout learning;

5.2 Confirmative

Remaining citations may be labeled **confirmative**.

Some more specific cases:

- if the citation refers to the use of a tool, label **confirmative**, even if there is no explicit value judgement (it is assumed that any use of the tool at all is positive).
- if the author uses the cited algorithm, technique, etc., without alteration.

Typical **confirmative** example:

- (23) A later study (Pang and Lee, 2004) found that performance increased to 87.2% when considering only those portions of the text deemed to be subjective.

References

Moravcsik, M. J. and Murugesan, P. (1975). Some results on the function and quality of citations. *Social Studies of Science*, 5:86–92.