

# Improved Statistical Evaluation for Grammatical Error Correction

Anonymous EMNLP submission

## Abstract

xxx

### 1 Main ideas(abstract base)

We propose two notions of conservatism: semantic conservatism and formal conservatism and state the former one is mostly the one we strive for when correcting.

We state semantic annotation is valuable for error correction and error assessment.

We present a new method to compare UCCA annotation

we show UCCA can be used for annotating ungrammatical texts

we show UCCA is stable over grammar correction

we show state of the art error correction methods are highly formally conservative (why formal and not structural or another name?)

We show current correction methods are too formally conservative, they don't change enough sentence boundaries, enough words and enough characters.

We suggest one of the factors contributing to over formal conservatism is having only a small number of references covering low percentage of the possible correction. It affects both assessment and development. In the development process we would expect good learning to understand not to try to correct complicated sentences as those will be very likely to be judged a mistake. Even when they are not.

Current correctors undercorrect, maybe due to lack of gold references.

Current assessment is an under assessment giving a much lower score than ought to be.

Significance of current scores, and what can we say about it...

### 2 Needed background

learner language is hypothesized to be a consistent language allowing us to think of error correction as a translation for closely related languages.

ungrammatical texts are being a main subject to research where learner language can be seen as one case.

### 3 Opening

In the history of error correction conservatism was considered an important trait of an error correcting system(?). This was also the reason why  $F_{0.5}$  became since conll2014(?) the measure of choice for error correction evaluation, emphasizing precision over recall. This emphasis can be understood as encouraging avoidance of wrong corrections at the cost of correcting less errors overall. The thought that stands behind such emphasis is that a user would be understanding towards errors he did, of which he is probably not even aware, not being corrected, but would not be so understanding when he sees a correction changes what he knows to be correct to a sentence saying something he did not mean it to. We want to refine this idea and suggest that there are subtleties we better address in this intuition.

We wish to say there are two different conservatism types, semantic conservatism and formal one, of which the semantic is the one we strongly need to adopt and the formal is merely a technicality for the user. in part 1\*\*make sure it is right where all is written, cross references\*\*we address the two conservatism types in part 2 we show semantics can be a consistent and measurable allowing use of it for ensuring semantic conservatism in part 3 we show current systems tend to be over too conservative when compared with human corrections in part 4 we suggest this to be an outcome of the evaluation measure used, being formally conservative and lacking.

## 4 Conservatism - not a single concept

### 4.1 what we really wish to be conservative about

In the task of grammatical error correction it is important to be conservative, not to over-correct. The user expects the minimum corrections necessary and wants no intervention in what he wishes to say. More specifically, he expects that what he has said would not be changed into something he did not. In this we may find two notions of conservatism, formal conservatism and semantic conservatism. Where any change in the original string would not be considered formal conservatism, only changes in meaning are accounted for semantic conservatism.

In many of the tasks, such as error correction for learner language, the user does understand his grammar is not perfect and would accept a change in grammar when needed. Because of this approval we also hypothesis, and it may call for a user study to prove or disprove this hypothesis, that users might accept a correct text unit of theirs to be corrected to another correct text unit with the same meaning. Maybe even more importantly, we aim to have as many correct sentences as possible, but as the grammar isn't fully correct in the first place, nor is the user's understanding of it, failing to correct grammar is acceptable. Changing meaning will be totally unacceptable, and also surely detectable by the user. In other words, the users do expect the corrector to be active and not too formally conservative,

but only as long as it is semantically conservative.

Moreover, as corrections are based on statistics, they might even just correct to a more common way of saying the same thing. Such unnecessary correction is not formally conservative, and at grammatical error correction maybe be unwanted, but not strictly unwanted as overall it is semantically conservative still. Additionally, some may even say this correction is a a needed one having a better grammar considering Fuzzy Grammar(?, ?) or a more fluent way to say the exact same thing. The latter was suggested as a necessary shift in the goals of error correction(?) Considering all this, we propose that next generation grammatical error correctors and evaluation will be focused on semantic conservatism when possible rather than on formal conservatism.

## 5 Semantics in learner language

### 5.1 Uses of semantic annotation

As semantic annotation was not used before to aid grammatical error correction, it is worth noting the a-priori reasons for developing it. The first and perhaps the most obvious use of semantic annotation would be to use it as a feature for correctors. The annotation may capture the gist that is supposed to stay the same when correcting, allowing the corrector to filter results or re-rank them based on the annotation or just to put it inside the mesh of features and learn automatically what to do with it, just as done with grammatical annotation. Later in this section we will not only discuss how to compare those features, specifically UCCA, but also why it is suspected to be a valuable feature.

Another approach of using semantic annotation would be for assessment. Reliable assessment by a gold standard might be hard to obtain (see 7), and human annotation for each output is great(?) but costly, especially considering development process. In these conditions, given a reliable semantic annotation we can enhance the reliability of our assessment. One way to do that might be to decouple the meaning from the structure. We propose a broad idea for a reduction from grammatical error detection and a comparable semantics annotation to grammatical error cor-

rection assessment. Lets assume we have both a reliable error detection tool and a good way to measure semantic changes. Then, we can transform assessment to a 3 steps assessment. First, detect errors in the original text. Assess the percentage of needed corrections that were actually corrected. Second, assess how much of was the semantics changed. Give a negative score for changing semantics. Third, use the error detection again to assess how many errors exist in the correction output, whether uncorrected by the corrector or new errors presented by the correction process itself.

This assessment was partially inspired by the WAS evaluation scheme(?), in short it states we should account in the assessment for 5 types, not only the True\False Positive\Negative but also for the case where the annotation calls for a correction, and the system did a correction, but one that is unlike the annotation's one. With the proposed assessment we can measure how many of the corrections were corrected correctly (First + Second), and how many errors do we have eventually (Third) and combine them to get something similar to the Precision Recall that is widely used. We can also account for the places where the error was detected and check if it was corrected in a way that makes it grammatical and did not change semantics, the fifth type. We do that without getting a human to confirm this is indeed a correction.

This system would be even more informative. Allowing assessment of what exactly is the part in which a corrector failed. Answering questions like: was it over formally conservative and did not make enough corrections? Was it making changes in the right places but not correcting grammar successfully? Was the system correcting grammar but changing things it was not supposed to? etc.

## 5.2 grammar can be annotated but is ill defined

Syntactic representation is very popular and useful in many NLP tasks\*\*cites\*\*. Thus, one thought that rises to mind is to use grammar annotation to evaluate corrections. While not useless, this approach is not well defined, and unclear bot practically and theoretically. One would say that the grammar would be the one induced by the actual words that appears

in the sentence, this would lead to annotation that calls for applying the syntax of Proper English to the different learner languages that just don't correspond to it. Thus, the structures may differ between different learners and they will tell us little about how to understand the sentence. This approach was being pursued in (?).

Others may suggest, and indeed they have(?), an opposite approach, saying the grammar meant by the learner is the one we should tag, but that requires having a corrected form of the sentences. Later<sup>7</sup> we show that for many sentences different corrections are possible. And where as sentences differ so does their grammar. later in this section, we propose semantic annotation as a well defined structure.

## 5.3 Learner language can be annotated by UCCA

At least theoretically, semantics are well defined even on ungrammatical text. With the right tools we might capture at least some of the semantics of sentences and use them for whatever we wanted grammar for and for other tasks. In this work we will use Universal Conceptual Cognitive Annotation (UCCA)(?), we will show that practically there are semantic annotation schemes that can be used for the purposes discussed.

But as in the syntactic representation, before we can claim anything about semantics using UCCA it is needed to show that UCCA is even consistent when applied to ungrammatical language such as learner language. To do that we used NUS(?) a parallel corpus of learner language and corrected versions which is the de facto standard since CoNLL 2013\14 shared tasks (?; ?). The NUS corpus consists of paragraphs of about 400 words each about various topics. We employed two cognition graduate students, both with background of working for a couple of years as translator. Each one had received the guidelines to read and annotated a couple of proper English paragraphs and then learner language paragraphs as an exercise. These paragraphs were compared between the annotators and each disagreement discussed in the hope of finding common annotation mistakes and choosing a methodological approach to borderline cases. After that each annotator has annotated

2 learner language paragraphs consisting of almost 800 UCCA nodes each. Over the uncoordinated paragraphs we computed the strict inter annotator agreement mentioned by (?) considering each Node in the directed acyclic graph (DAG) of UCCA annotation as agrees if and only if its label and the labels of all its span leaves were considered to have the same labels respectively, from that we derive an F1 score.

We got an F1 score for the inter annotator agreement of 0.845 with Precision 0.834 and Recall 0.857 we see that as enough to be a proof that UCCA can be applied to, especially considering those numbers are a bit higher than the inter annotator agreement reported in the reported originally for formal English(?). We explain the rise in agreement by the fact that the guidelines and procedures were refined since UCCA was first introduced and not to superiority of UCCA for annotating learner language. A similar F1 score for inter annotator agreement (0.849) over 2 corrected paragraphs suggests the same.

#### 5.4 Semantics are preserved when correcting grammar

As a next step each annotator annotated corrected paragraphs corresponding to ones he already annotated, 7 different paragraphs were annotated in this way. To avoid misleading high score due to the fact that each annotator annotated both the learner language paragraph and the corrected paragraph 3 different tuples of paragraphs were annotated by both annotators allowing a cross comparison, meaning that for each paragraph we compared the annotation of the learner language done by one annotator with the annotation of the other annotator done for the corrected paragraph and vice versa.

As a next step a comparison between the annotations was needed, but there exists no measure for how similar two different UCCA annotations of different texts are. We considered using suggested semantic measures such as SMATCH(?) but it can not work for UCCA or DAG similarity measures such as graph kernels (e.g.(?)), but those tend to work on bigger graphs and would be the wrong tool for the small UCCA DAGs. Thus, a new measure is called upon.

#### 5.5 Similarity measures

We propose several new methods to compare UCCA annotation of a learner language with UCCA annotation of corrected texts, giving a more accurate measure than the upper bound suggested by (?) for comparing two parallel texts in different languages, while keeping the essence of comparing how many of the aligned nodes conserve meaning and tag. For that we may think for a moment on error correction as translation from learner language to Proper English, and a good translation would be a translation which keeps the meaning but has the syntax of English. Considering that, just like in translation we can align words from the learner language to the corresponding words in English and keep record of how many of those nodes kept their labels.

As comparing labels is trivial between a pair, \*\*should mention somewhere weak labeling? we should focus on how we propose to align nodes. We should note first that alignment should not be at the token level, as we want to allow tokens to be replaced or removed as long as the higher structures convey the same meaning. We thus prune the labels above the leaves, the tokens of the sentence. To define an alignment of the nodes, we suggest some possible ways, all based on first aligning the words in order to give order to the DAG and then comparing the structure in one way or another.

In order to align tokens we use the fact that, unlike in translation, aligning words is a simple task as most of the words are kept unchanged, deleted fully, added, or only changed slightly. This allows us to align words well using edit distance measure, knowing that words that exists in both sentences will have low edit distance. We consider aligning to sets of words a bipartite graph matching problem, with weights according to the edit distance. For tie breaking, we add a penalty. The penalty is always smaller than 1, the minimum cost of one action, favoring a sentence order when a word occurs twice.

As to aligning nodes, we can use word spans of each node, based on the token alignment and the DAG structure, to choose how best to align. A first and most straightforward approach would be to compare all pairs of nodes in parallel paragraphs and to each node from

one paragraph assign the one most similar node, span wise from the other. That approach is quite similar to the inter annotator agreement aligning, but it has three drawbacks; it is assymetric; it may be over optimistic aligning nodes without considering the DAG structure; and second it might be slow for many nodes. Being assymetric is not much of a problem as we can compute the measure twice and use the mean of the results, that would also be the case for other assymetric methods we suggest. In order to address the other drawbacks we propose different aligning methods.

A second method driven by the assumption that nodes higher in the hierarchy are more important to the semantic representation is measuring the largest cut in which nodes are aligned (top down) to each other and have the same labels. This is expected gives a harsher lower similarity score but one of which might be more representative of the semantics that are kept and hopefully more informative for tasks that will use it.

A third type of methods were token similarity methods, these methods use one kind of aligning (top down, bottom up or all to all) and only compare the meaningful nodes. This was called upon in the (?) paper. This approach makes sense due to the fact that some labels are well defined and thought upon while others are still vague and call for future work on refining or adjusting them, moreover, some labels are more semantic while other labels are currently just a place holder as each node must get a level, and the semantic role is not always clearly defined (e.g. the word “is” in “he is walking” seem to be more syntactically related than semantically). The unused labels are center, elaborator, function, relation, linker, ground and connector.

A bit different way than all the others is to compute the labeled tree edit distance(?), for that we first needed the trees to be ordered, we did that in a top down fashion. An interesting future work would be to use unordered tree edit distance methods(?).

All of the code to implement UCCA structures, align them and evaluate them is also given as a free contribution.\*\*\*link\*\*

## 5.6 Results

We present in table \*\*\*\* the scores of the different presented methods. For each method we present the average results of 9 tuples of paragraphs annotated by the same annotator and 6 tuples where each paragraph was annotated by a different annotator.

Finally, we present as a control measure and a bound on the best score we can expect to get in such comparison the scores of 7 paragraphs in which we compare two annotations for the same paragraph using all the similarity measures discussed, it can be thought of as a different way to defining inter annotator agreement. Note that a similarity of 1 and distance of 0 is indeed reached when comparing an annotation with itself.

In table \*\*\* we present the results of the token analysis, the upper bound suggested by (?), showing similar results for learner language - corrected tuples as those seen in English - French comparison.

## 5.7 Discussion

From the result we learn a number of things, we show that the upper bounds in table\*\* suggest high stability of UCCA over grammatical error correction, and the results are similar to those shown over translation. This upper bound seem not to be very strict if the other measures are to be considered true values, we do note that because of the aligning errors those measures are actually more of a close lower bound than an exact value.

We see that measurements for symmetry that are similar to the inter annotator agreement measure also suggest high stability, achieving scores not much lower than the one different annotators get for the same paragraph. This result is quite strong as an inter annotator agreement is the upper bound being the score of comparing a paragraph to itself. Most importantly we learn from it all that even when correcting grammatical errors the semantic structure (as represented by UCCA at least) is hardly changed and thus can be used as a tool to avoid introducing semantic changes when trying to only change grammar. The symmetry measures we introduce can be used to enforce semantic conservatism. This would be a good place to remind that a direct way to measure

semantic conservatism as we have got here will allow us to be less formally conservative while focusing on the conservatism users and hence we are more interested in.

## 6 Over conservatism in current error correction attempts

In recent years a lot of research was done trying to create automatic error correction(?, ?; ?; ?) and given our research on semantic conservatism, it is reasonable to wonder whether these measures can help improving the existing correctors. To answer that we need to analyze how conservative these correctors are, something that we see as insightful and important by itself in order to improve the correctors.

The first step would be formal conservatism. If corrections are very formally conservative they are likely to be semantically so too. In addition, this analysis as will be discussed soon will show that this analysis is the one really needed at this step of development.

### 6.1 Assessing formal conservatism

Our goal was to analyze the output of all of the participants in Conll 2014 shared task(?) and of the current state of the art (?). We started at manually analyzing, our impression was that there is a real lack of corrections. Albeit important, manual analysis is not enough and we aimed for some quantitative measures. For that we first aligned each learner language text unit to a corresponding corrected text unit. We used an exact match for last words in a sentence as a boundary symbol, thus allowing a text unit to be more than one sentence. This alignment is needed because we only know the final corrections, a main obstacle that was considered in the assessment methods as well(?). Our first result to present will be how many sentences are concatenated and how many split using the different methods. Moreover, we present the same measures for the corrections done in the NUS(?). To have better evaluation of the real goal of corrections we also compute all of the measures on the TreeBank of Learner Language (?) based on the Cambridge First Certificate in English (FCE) (?), a new large parallel corpus containing language of learners native of different languages.

Next we were interested in how many words are being changed and how much word order was disrupted. We used the alignment of sentences and for each sentence we aligned words by edit distance in the same manner explained in 5.5. We calculated the number of words changed per sentence to assess how many words were edited or removed. To measure how much the word order was preserved we used spearman's rho for the indexes of the aligned words in each sentence.

### 6.2 discussion

Lets discuss the results in \*\*\*\*\* starting at what we see isn't change in the gold standard. We can call this what calls for formal conservatism. We see that it is most common for a sentence to have no change, not to concatenate two consecutive sentences and not to split a sentence into two. We also observe high correlation coefficient for most of the sentences. Summing it all together we indeed have more unchanged than changed in every measure we have. But, with a closer look we should also notice that it mostly tells us about the dataset. Specifically the level of English found in the dataset.\*\*rephrasing\*\* These measures should be seen as a gold standard of the amount of corrections to be done, and as we might wish to be a bit conservative and not exceed it, this is still where we should aim.

When we broaden our view and consider the results of the different correctors, the picture is clear. All correctors are over conservative. It is not only that correctors don't tend to over-correct, they all, to the last one, by all the measures we checked, undercorrect a lot and are over conservative. \*\*say a word on how do we see that in each graph and conclude\*\*

## 7 May lack of corrections lead to over formal conservatism?

### 7.1 The idea

Some corrections might really be too hard for current correctors to correct, leading to cautious corrections. Another cause for over conservatism might hide in the assessment methods.before we show it, lets assume that each ungrammatical sentence has some possible corrections. From the corrections only

Figure 1:  
text

2 are in the gold standard. That would lead to problems in both assessment and development process. In the assessment, results will not be reliable, the assessment will only assess whether we predict the gold standard annotation and not how many of the sentences where corrected. Perhaps less obvious will be how it affects the development process. Sentences, even if corrected well, which have more possible corrections will get lower scores on precision for correcting, while not getting scores for recall anyway. This will lead either through machine learning or algorithm development cycles to learn not to correct those sentences at all. In the rest of this section we will show that the assumptions we made are more than just assumptions.

## 7.2 Formalism

Lets denote  $X = x_1 \dots x_n \sim \text{World Sentences}$  as the considered sentences for evaluation. For each  $x_i$  there exists  $\{y_1^i, \dots, y_M^i\} \sim x_i \text{ corrections}$  the set of gold standard annotations and we consider the output of the corrector to be a function  $f(x_i)$ . A certain assessment statistic is a function  $\hat{S} = \text{Eval}(f(x_1), \{y_1^1, \dots, y_M^1\}, \dots, f(x_N), \{y_1^N, \dots, y_M^N\})$ .

## 7.3 our data

\*\*\*put it somewhere in context\*\*\* This way of measuring assumes the sentences do not need context, and while surely untrue we do assume the context will account to having a bit less possible corrections but the bigger picture will stay more or less the same. We also did not use sentences longer than 15 words, assuming these will be harder to annotate and are more likely to have independent corrections\*\*maybe explain that before?\*\*. This choice might give us a bit lower results in the number of corrections, negating the effect of the context assumption and only exclaiming the claim that there are \*\*much more corrections than we account for currently\*\*

If the NLP community has agreed one correction is not enough\*\*cite\*\*, we can now say 2 is no magic number either. We can also see what we lose and gain from different amounts of references, and may also suggest

reduction + even with indexes the problem is hard as indexes vary a lot too.

## 7.4 Lack of corrections also leads to under estimation of the statistic and hence over conservatism

## 8 On significance and variation

While  $E(\hat{S})$  vary only as we change  $M$  the number of annotations, but not  $N$  the number of corrections,  $\text{Var}(\hat{S})$  depends on both. We try to assess and give an upper bound on how much it varies for different  $M$  and  $N$ , allowing for both a smart allocation of resources when building a corpus and for assessing on given corpora whether two systems are actually different.

## 9 Other things(conclusions? discussion further work?)