# TransCasm: A Bilingual Corpus of Sarcastic Tweets

**Desline Simon[2], Sheila Castilho[2], Pintu Lohar[2] and Haithem Afli[1]**
ADAPT Centre
[1]Department of Computer Sciences, Munster Technological University, Cork, Ireland
[2]School of Computing, Dublin City University, Dublin, Ireland
{firstname.lastname}@adaptcentre.ie

## Abstract

Sarcasm is extensively used in User Generated Content (UGC) in order to express one's discontent, especially through blogs, forums, or social media such as Twitter. Several works have attempted to detect and analyse sarcasm in UGC. However, the lack of freely available corpora in this field makes the task even more difficult. In this work, we present "TransCasm" corpus, a parallel corpus of sarcastic tweets translated from English into French along with their non-sarcastic representations. To build the bilingual corpus of sarcasm, we select the "SIGN" corpus, a monolingual data set of sarcastic tweets and their non-sarcastic interpretations, created by (Peled and Reichart, 2017). We propose to define linguistic guidelines for developing "TransCasm" which is the first ever bilingual corpus of sarcastic tweets. In addition, we utilise "TransCasm" for building a binary *sarcasm* classifier in order to identify whether a tweet is sarcastic or not. Our experiment reveals that the *sarcasm* classifier achieves 61% accuracy on detecting *sarcasm* in tweets. "TransCasm" is now freely available online and is ready to be explored for further research.

**Keywords:** Sarcasm, parallel corpus, translation

## 1. Introduction

Sarcasm is an extremely ambiguous form of wit, extremely hard to analyse for a machine but as well as for the average human reader. Sarcasm is even harder to be analysed in written form, as the tone of voice (pitch, heavy stress), gestures or facial clues (hand movements, rolling the eyes, etc.) which are important to detect sarcasm in spoken form, cannot be considered in textual communication. Sarcasm is an indirect way to communicate negation using a contradiction between a sentiment and a situation. It is usually used to express the opposite of one truly wants to say. In fact, its structure very often consists of a contrast between positive or intensified sentiment words to convey a negative feeling - a very strong one most of the time - such as insult, irritation, hostility, disagreement, mockery, etc. Sarcasm is a particular form of irony, but with the intention and consciousness not to be explicit to express a negative feeling or aggressive attitude (acrimony, bitterness). In the example "I love being ignored", a very positive word *love* is used in a negative context *being ignored*. The particular structure of a positive word followed by a negative situation is often a strong indicator of sarcasm.

Currently, sarcasm is extensively used in UGC in order to express one's discontent, especially through blogs, forums, or social media such as Twitter. The importance and need to detect sarcasm is real. It is attracting interest in many application domains. In fact, BBC reported[1] on January $5^{th}$, 2014 that the US Secret Service is seeking a Twitter sarcasm detector. The motivation of our present work comes from the challenges in processing sarcastic tweets. We develop a first ever corpus of bilingual sarcastic tweets for the English–French pair. On top of it, we also build a binary sarcasm classifier using the corpus in order to mitigate the difficulties in detecting sarcastic tweets.

## 2. Related work

Twitter sentiment analysis has attracted many researchers during the last few years (Agarwal et al., 2011; Zimbra et al., 2018; Branz and Brockmann, 2018). Parallel twitter corpus has also been developed for Twitter sentiment translation (Afli et al., 2017; Lohar et al., 2017). The sarcasm phenomena has been well-studied in linguistics, psychology and cognitive science (González-Ibáñez et al., 2011; Gibbs, 1986). But in the natural language processing (NLP) literature, automatic detection of sarcasm is considered a difficult problem (Pang and Lee, 2008) and has been addressed in several works. Bouazizi and Ohtsuki (2015); Reyes and Rosso (2012) use sarcasm detection to enhance the efficiency of after-sales services and consumer assistance through understanding the intentions and real opinions of customers when browsing their feedback or complaints.

Inspired by the advances on sentiment analysis (SA) (O'Connor et al., 2010) and figurative language processing research (Reyes et al., 2012), the field of sarcasm detection has started to benefit from using SA (Ghosh et al., 2015). To the best of our knowledge, no parallel corpus of sarcastic tweets is available till date. In this work, we propose to define linguistic guidelines for building a first ever bilingual corpus of sarcastic tweets, based on the extension of the unique corpora created by Peled and Reichart (2017). In addition, we train a sarcasm classifier using this corpus in order to

---

[1]https://www.bbc.com/news/technology-27711109

detect whether a tweet is sarcastic or not. We conduct initial experiment on sarcasm detection and measure its performance.

## 3. Method

### 3.1. SIGN Corpus

To build the bilingual corpus of sarcasm, we have selected the SIGN corpus, a monolingual data set of sarcastic tweets and their non-sarcastic interpretations. The SIGN corpus was created by Peled and Reichart (2017) and consists of $3,000$ sarcastic tweets containing text only, where the average sarcastic tweet length is 13.87 utterances and the vocabulary size is $8,788$ unique words. The authors used Twitter API[2], and collected tweets with the hashtag *#sarcasm* from January 2016 to June 2016. Following the collection, the authors asked five human judges to write honest, non-sarcastic interpretations of those tweets to capture the original meaning behind them. Therefore, a monolingual parallel corpus was created in which sarcastic English tweets were translated into non-sarcastic English. For example, a sentence "How I love Mondays. #sarcasm" would generate interpretations such as "how I hate Mondays' or "I really hate Mondays". Table 1 shows examples of sarcastic tweets and their interpretations.

Quite often, the human judges would interpret the sarcastic tweets in the same manner, getting the same meaning behind the sarcasm, but at times, the interpretation would greatly differ from each other, which is a strong indicator that sarcasm is extremely complex to analyse and interpret, even for humans. As the non-sarcastic interpretations composed for each tweet might help the task of sentiment analysis - since those honest interpretations capture the meaning behind the original sarcastic utterances - we also decided to translate a few of those interpretations. However, due to time constraints, no more than three interpretation per tweets were translated. The decision on what interpretation to translate depended on the vocabulary encountered, and the overall relevance of those interpretations.

### 3.2. Creating the TransCasm Parallel Corpus

The corpus is being translated from English into French with the help of MateCat[3] , a web-based open source tool. Figure 1 shows the translation set-up. In the corpus, the same tweet is presented each time their interpretations are presented. Both the original and the interpretations are separated by a comma. As the SIGN corpus is normalised, that is, the hashtag (#sarcasm), capital letters and all punctuations of any kind are cleaned, the interpretation of the tweets was challenging, and in consequence, their translation. The topics in the SIGN corpus are greatly eclectic considering that the method to extract them was only based on the

query #sarcasm. The topics we were able to identify while translating include:

- weather

  e.g.1 ' really looking forward to more rain today', " i'm hoping that there won t be more rain today"

  e.g.2 ' 96 degrees one of the great benefits of living in california', ' 96 degrees california is too hot for me'

- traffic and transport

  e.g.1 ' delays on the piccadilly line yayyyyyyy', ' oh no there are delays on the piccadilly line'

  e.g.2 ' ah yes i love replacing a tire at 9 in the morning', ' i hate replacing tires especially early in the day'

- work

  e.g.1 ' almost lunchtime i get a half hour away from this paradise', ' almost lunchtime i get a half hour close far this terrible workplace'

  e.g.2 " don't you just love when people takes the credit for something you did ? yeah ? me too", ' i hate people taking credit for something i did'

- school

  e.g.1 ' finals are so nice like omg best thing ever to happen to my life like love studying', ' i hate studying for the final exams'

  e.g.2 " making flash cards for my exams tomorrow i'm having fun", ' i am not enjoying making these flash cards for my exams tomorrow'

  e.g.3 " it's teacher appreciation week and the lovefest is in full swing at my school", " it's teacher appreciation week but they are not getting any love at my school"

- health

  e.g.1 ' dentists make money off of people with bad teeth so should we really trust the toothpaste they recommend ?', " we shouldn't trust dentist"

  e.g.2 ' a nice long wait in the doctors office should calm my nerves', ' waiting at the doctor office will stress me out'

- sports (especially American football, baseball, ice hockey and football (soccer))

  e.g.1 ' did bartolo colon hit a hr tonight ? i didnt see it mentioned anywhere on twitter ?', ' did bartolo colon hit a hr tonight ? it was mentioned everywhere on twitter'

  e.g.2 ' terry got sent off ?  maybe hiddink needs to have more control over his players', ' hiddink already has a lot of control over his players'

  e.g.3 " cavs aren't getting any calls this is new", " cavs aren't getting any calls as usual"

| Sarcastic Tweets | Honest Interpretations |
|---|---|
| What a great way to end my night #sarcasm | 1. Such a bad ending to my night<br>2. Oh what a great way to ruin my night<br>3. What a horrible way to end a night<br>4. Not a good way to end a night<br>5. Well that wasn't the night I was hoping for |
| Staying up till 2 : 30 was a brilliant idea, very productive #sarcasm | 1. Bad idea staying up late, not very productive<br>2. It was not smart or productive for me to stay up so late<br>3. Staying up till 2 : 30 was not a brilliant idea<br>4. I need to go to bed on time<br>5. Staying up till 2 : 30 was completely useless |

Table 1: Example of two sarcastic tweets augmented by five non-sarcastic interpretations



Figure 1: Translation process using the Matecat tool

- politics (especially the American Republican Party and the then candidate Donald Trump)

  e.g.1 ' now trump is bringing up bill clinton 90 s affairs because we all know there are no current day pressing matters to focus on', ' there are important matters to focus on instead of bringing up affairs of bill clinton from the 1990 s'

  e.g.2 " so trump won the republican nominee so that's that good job america", " so trump won the republican nominee so that's that shameful job america"

  e.g.3 ' obviously a well prepared speech by trump', ' obviously a poorly prepared speech by trump'

- social relationships (love, friends, family and co-workers)

  e.g.1 ' being left out is such an amazing feeling', ' being left out is such a painful feeling'

  e.g.2 " really don't know what i'd do without you", ' i am doing great without you'

  e.g.3 ' some people are just really brilliant', ' some people are just really dumb'

  e.g.4 " many people don't know this but you can actually read a book or go to the gym without announcing it on facebook", ' people need to understand that nobody cares of what you do stop trying hard on facebook'

As sarcasm is very often, if not all the time, used to express frustration about being in unpleasant situations, undesirable states, and unenjoyable activities, it is not surprising that the corpus contained the topics listed above.

### 3.2.1. Complexity of the Translation Task

Each social media platform has its own idiosyncrasies, and occasionally even new dialects. Twitter is no exception to that rule, as writing a tweet is limited to 280 characters. The length constraint has created a new way of communicating, especially considering the massive use of acronyms, abbreviations, slang, phonetisation, onomatopoeia and interjections, words contractions, etc. Moreover, the informal nature of tweets which of-

100

ten lack grammatical structure (spelling errors, syntax problems) frequently leads to misunderstanding issues, which in turn, leads to translation issues.

The lack of context, and the lack of knowledge about the real intentions when writing the tweet, is another issue faced. The collected tweets are not linked to each other and the order they are presented are random. This issue was also shared with Peled and Reichart (2017), as at times, the human judges could not understand the meaning of sarcasm and, therefore, unable to write an interpretation. In this case, they were told to skip the tweet. These issues described above had an impact on the translation as each time an unknown acronym, or the use of technical vocabulary used by people to comment were found, added to the lack of context, it meant that the translator needed to research those specific terminology, which in turn, affected the translator's efficiency.

### 3.2.2. Guidelines

Considering the issues we found during translation, a set of rules was agreed upon in order to keep the translation to the same standard:

1. Length: even if the translation exceed the 280 characters, find the best equivalent in French.

2. Structure: keep the original structure of the original corpus. The original tweet is displayed first between quotations followed by its interpretation, separated by a comma [' ', ' '].

3. Capitalization: keep the original capitalization found int he corpus.

4. Accents: write the translations using the French diacritics/accents, even if the tweets appear unstructured or ungrammatical, implying that the language register used by the author is casual or informal.

5. Slangs or Onomatopoeia: find the best French equivalent.

   'ass hat' = 'tête de nœud'
   'that sucks' = 'ça craint' or 'c'est nul'
   'wow' = 'ouah'
   'yuck' = 'beurk' or 'pouah'
   'boo' = 'bouh'

6. Acronyms and abbreviations: find the best French equivalent, but if not possible, leave the original acronym untranslated.

   'bc' for 'because' = 'pcq' for 'parce que'
   'omg' for 'oh my god' = 'omd' for 'oh mon dieu'
   'lol' for 'lot of laugh' = 'mrd' for 'mort de rire'
   'thk' for 'thanks' = 'mrc' for 'merci'
   'bter' or 'btr' for 'better' = 'mx' for 'mieux'

7. Metaphors and idioms: translate metaphors and idioms focusing on the adequacy/equivalence of the translation, by finding the best French equivalent.

   doublespeak' = 'langue de bois'
   'wipe the floor' = 'mordre la poussière'
   'like taking candy from a baby' = 'c'est enfantin' or 'c'est un jeu d'enfant'
   'bet the farm' = 'tout miser'

8. graphemic stretching: keep the same number of repeated letters in the best French equivalent translation.

   'greeeat' = 'géniallll'
   'realllly' = 'vraimentttt'

These rules were applied for the set of the corpus that has been translated so far. We believe that further rules will be added as the translation advances.

## 4.   Initial experiments with TransCasm

In addition to developing the "TransCasm" corpus, we also utilise it in sarcasm detection. As mentioned earlier, each sarcastic tweet in English is transformed into at most 5 non-sarcastic interpretations. Some of the tweets contain less than 5 interpretations. Sometimes it is very difficult to transform a sarcastic tweet into many different interpretations, a single non-sarcastic representation is enough in these cases. TransCasm consists of $1,831$ different non-sarcastic interpretations of $860$ unique sarcastic tweets in English. These unique $860$ sarcastic tweets and their $1,831$ non-sarcastic interpretations are translated into French. The corpus statistics is shown in Table 2.

| Language | #Sarcastic Tweets | #non-sarcastic interpretations |
|---|---|---|
| English | 860 | $1,831$ |
| French | 860 | $1,831$ |

Table 2: corpus statistics

In our experiment of binary sarcasm classification, we use $610$ out of $860$ unique sarcastic tweets as the training data set and held out the remaining $250$ as the test data set. These $610$ sarcastic tweets are then mixed with randomly selected $610$ non sarcastic representations so that our training data becomes an equal distribution of sarcastic and non-sarcastic tweets. Similarly, the test data set is created by mixing the $250$ non-sarcastic tweets with $250$ sarcastic tweets. We use multinomial naive bayes (MNB) classification approach to train our sarcasm detection model. The model is automatically tuned by randomly selecting the $20\%$ of the training data itself during the training process. The data distribution of the whole experiment is shown in Table 3.

## 5.   Results

Although the main objective of this work is to develop the first ever bilingual corpus of English–French sar-

| Distribution | #sarcastic | #non-sarcastic |
|---|---|---|
| train | 610 | 610 |
| test | 250 | 250 |

Table 3: Data distribution

castic tweets, we conduct an initial experiment on sarcasm detection using this data set and obtain interesting results. We apply our binary sarcasm detection model to the test data set of 500 mixed tweets (both sarcastic and non-sarcastic) to see how the system performs in identifying sarcastic tweets. Our system achieves an accuracy of 61% in identifying whether a tweet is sarcastic or not. Table 4 shows the detailed results with accuracy for each category. We can observe that the

| Tweet type | #of tweets | #Correctly classified | Accuracy |
|---|---|---|---|
| sarcastic | 250 | 169 | 67.6% |
| non-sarcastic | 250 | 136 | 54.4% |
| all | 500 | 305 | 61% |

Table 4: Experimental results

sarcasm classifier obtains an accuracy of 67.6% for sarcastic and 54.4% for non-sarcastic tweets, respectively. However, we note that this is a preliminary phase of our experiments and further research directions with this data set is planned.

## 6. Conclusion and Future Work

The main contribution of this work is to present an on-going translation project that aims at building the first ever parallel sarcasm corpus - "TransCasm" corpus - by fully translating the SIGN corpus, a freely available corpus of sarcastic tweets and their non-sarcastic interpretations, from English into French.

We translated 860 unique sarcastic tweets in English into French. Each of the tweets, having at least 1 and at most 5 interpretations was translated, which amounted to 1, 831 translations of the interpretations. Due to the informal nature of the tweets and the issues found during translation, we developed guidelines in order to aid the translation process. In addition to developing TransCasm, we utilised this corpus for sarcasm detection in Tweets.

We built a binary classifier using the MNB classification algorithm. In the initial experiments, our classifier achieved an accuracy of 61% for the test data of 500 tweets, which can be considered a good mark given that this is the beginning of our work on "TransCasm". There is no such corpus available according to the best of our knowledge, therefore, "TransCasm" may become very useful resource in the NLP community, especially those working with the twitter data.

For future work, we intend to extend the corpus, as well as the translation of the interpretations. Upon extension we will also build the first sarcasm translation system (TransCasm MT). In addition, we will aim to build

more robust sarcasm classifier with the extended data set and compare with the state-of-the-art sarcasm detectors. We will also apply deep learning techniques to further refine our sarcasm detection model. Moreover, we believe that annotating the text features in which sarcasm can be identified, would be beneficial and make the corpus useful for multilingual sarcasm detection as well sentiment translation research. We have released the "TransCasm" corpus online[4] to facilitate further research on this data set.

## 8. Bibliographical References

Afli, H., McGuire, S., and Way, A. (2017). Sentiment translation for low resourced languages: Experiments on irish general election tweets. In *Proceedings of the 18th International Conference on Computational Linguistics and Intelligent Text Processing*, Budapest, Hungary.

Agarwal, A., Xie, B., Vovsha, I., Rambow, O., and Passonneau, R. (2011). Sentiment analysis of twitter data. In *Proceedings of the Workshop on Languages in Social Media*, LSM '11, pages 30–38, Portland, Oregon, USA.

Bouazizi, M. and Ohtsuki, T. (2015). Sarcasm detection in twitter: "all your products are incredibly amazing!!!" - are they really? In *GLOBECOM*, pages 1–6. IEEE.

Branz, L. and Brockmann, P. (2018). Sentiment analysis of twitter data: Towards filtering, analyzing and interpreting social network data. In *Proceedings of the 12th ACM International Conference on Distributed and Event-based Systems*, DEBS '18, pages 238–241, Hamilton, New Zealand.

Ghosh, A., Li, G., Veale, T., Rosso, P., Shutova, E., Barnden, J. A., and Reyes, A. (2015). Semeval-2015 task 11: Sentiment analysis of figurative language in twitter. In *SemEval@NAACL-HLT*, pages 470–478, Denver, Colorado, USA.

Gibbs, R. W. (1986). On the psycholinguistics of sarcasm. *Journal of Experimental Psychology: General*, 115(1):3 – 15.

---

[4]https://github.com/HAfli/TransCasm_Corpus

González-Ibáñez, R. I., Muresan, S., and Wacholder, N. (2011). Identifying sarcasm in twitter: A closer look. In *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA - Short Papers*, pages 581–586, Portland, Oregon, USA.

Lohar, P., Afli, H., and Way, A. (2017). Maintaining sentiment polarity in translation of user-generated content. *The Prague Bulletin of Mathematical Linguistics*, 108(1).

O'Connor, B., Balasubramanyan, R., Routledge, B. R., and Smith, N. A. (2010). From tweets to polls: Linking text sentiment to public opinion time series. In *Proceedings of the Fourth International AAAI Conference on Weblogs and Social Media*, pages 122–129, Washington D.C., USA.

Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis. *Found. Trends Inf. Retr.*, 2(1-2):1–135, January.

Peled, L. and Reichart, R. (2017). Sarcasm sign: Interpreting sarcasm with sentiment based monolingual machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1690–1700, Vancouver, Canada.

Reyes, A. and Rosso, P. (2012). Making objective decisions from subjective data: Detecting irony in customer reviews. *Decision Support Systems*, 53(4):754 – 760.

Reyes, A., Rosso, P., and Buscaldi, D. (2012). From humor recognition to irony detection: The figurative language of social media. *Data Knowl. Eng.*, 74:1–12, April.

Zimbra, D., Abbasi, A., Zeng, D., and Chen, H. (2018). The state-of-the-art in twitter sentiment analysis: A review and benchmark evaluation. *ACM Trans. Manage. Inf. Syst.*, 9(2):5:1–5:29, August.