

UNIVERSITY OF OTAGO

SCHOOL OF COMPUTING

COSC385 PROJECT REPORT

Talking in French Like Academia

Machine Learning Powered Verlan Identification

Author:
Yitian LI (4556502)

Supervisor(s):
Dr. Lech SZYMANSKI
Dr. Veronica
LIESAPUTRA

October 11, 2025



Abstract

something.

1 Introduction

1.1 Context and Motivation

Since the early 19th century, the French people have started to talk using verlan. Just like Pig Latin¹ exists in English culture, verlan is an unusual and creative form of *argot* (slang) that is formed by flipping the syllables around in a word.²[1, 2] Time flies, verlan has become more and more popular, and it is now widely used amongst teens and young people in francophone societies³[3]. Examples of verlan can be as follows:

- bite = bi + te → te + bi → tebie (penis)
- shit = shi + t → t + shi → teuchi[3]
- bonjour = bon + jour → jour + bon → jourbon (greetings)

In real-life conversations, such can be used as in the example sentences below:

- *Le graff géant représente une tebie pixel art.*
(The giant graffiti depicts a pixel art penis.)
- *Il a du bon teuchi du bled.*
(He's got some good shit from the countryside.)
- *Un p'tit⁴ jourbon et tout le monde sourit.*
(A quick hello and everyone smiles.)

Indeed, verlan can be formed with different original languages, not only French, but also English and other languages. However, it always follows the same rule of flipping syllables, although, for better pronunciation reasons, certain minor amendments such as dropping unnecessary letters and applying accents (e.g., é, è) can be used from time to time[1]. Besides, due to the

¹en.wikipedia.org/wiki/Pig_Latin

²In fact, the word *verlan* is a verlan from the word *l'inver* (the inversion).

³Such as France, Belgium, Switzerland, Luxembourg, and Canada.

⁴Standard spelling: petit.

universal trait of slang being used more often phonetically instead of written, verlan users tend to spell them differently when writing them down. As technology develops, this has been occurring more frequently than ever in daily texting[4].

Thinking internationally, when people are communicating with translators, it is possible that slang in their mother language can be brought to the conversation, which could be tricky for translators to translate[5]. Using translators such as DeepL⁵ and Google Translate⁶ to translate sentences that contain verlan from French to English can be a specific example to prove this. Furthermore, although both of the translators above are using Machine Learning (ML) for translation, their results of translating verlan are not ideal[6, 7]. For example, when attempting to translate the sentence above, *Le graff géant représente une tebie pixel art.*, both Google Translate¹ and DeepL² cannot translate the word *tebie* correctly. Specifically, for DeepL, there is no desired translation as *penis* in its alternative word list for *tebie*³.

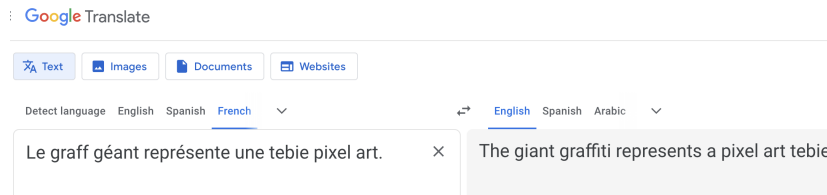


Figure 1: Google Translate cannot translate the verlan *tebie* correctly.

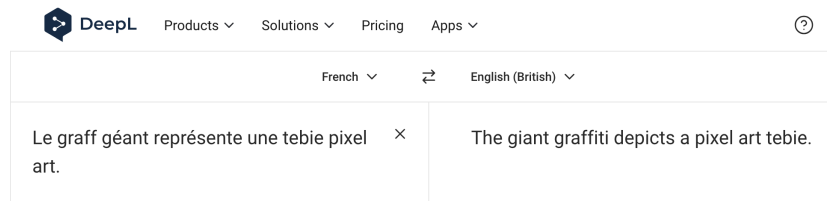


Figure 2: DeepL cannot translate the verlan *tebie* correctly.

⁵www.deepl.com

⁶translate.google.com

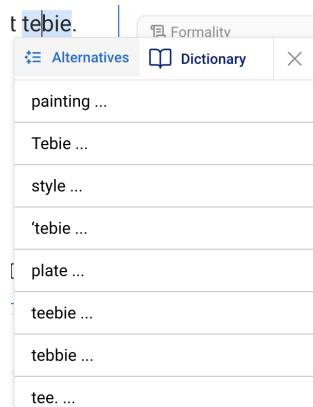


Figure 3: No desired translation for verlan *tebie* in DeepL’s alternative word list.

Thus, a question shall naturally arise: Can we improve translators’ performance in translating slang by improving the ML model? The answer is undoubtedly ‘yes’ in an era where artificial intelligence research is expanding rapidly. Researchers have been making progress in identifying slang using ML[13] and, moreover, in translating noisy text, of which slang is a part[8].

But what about verlan? There is no known ongoing or completed research on identifying *such* slang or their translations⁷, nor does a proper dataset exist. The only work similar to this is an assignment published at the University of Toronto⁸, asking students to train a Neural Machine Translation (NMT) model to transform standard English into Pig Latin. It is not only the other way around; instead of identifying Pig Latin and transforming it back to standard English, it is also more of an example for students to practice using NMT than a discussion on its identification and translation. Shouldn’t we do something?

This report aims to change that.

1.2 Objective

The purpose of the project is to create two verlan datasets: one functioning as a dictionary, containing the verlan words and their normalised standard French equivalents; the other a dataset of sentences that contain verlan,

⁷Until September 2025.

⁸<https://uoft-csc413.github.io/2022/assets/assignments/PA03.pdf>

paired with the same sentences containing normalised words, with labels indicating whether a sentence contains verlan. After that, the project embeds and classifies verlan using Large Language Models (LLMs) and analyses the results.

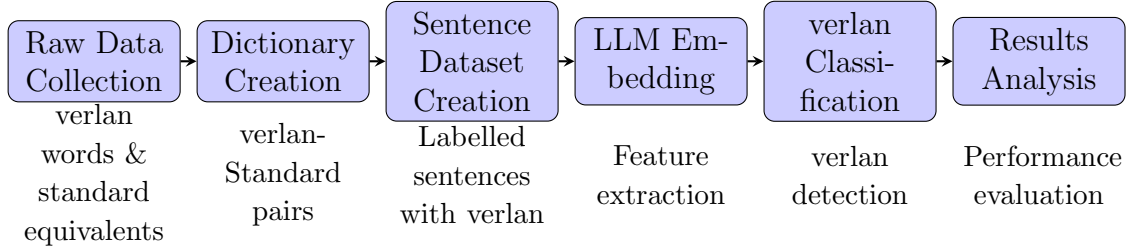


Figure 4: A visulisation of the objectives.

With the purpose above, the report contributes to the linguistics and the AI researchers two verlan datasets, for dictionary making or LLMs training. The report also evaluates how good we can achieve the identification of verlan with ML, to benefit machine translation in the future.

The code and the unannotated, un peer-reviewed dataset developed as part of the project are released under openlicences and aligns with open science best practices, with the usage of a version controlled software development platform (GitHub)⁹. The annotated, peer-reviewed dataset will be published shortly after this report, aiming by the end of 2025.

2 Background

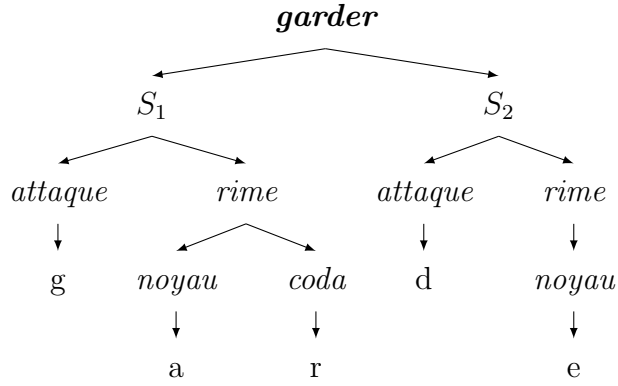
2.1 A Living Verlan

Vivienne Véla, a former scholar from Université Paris 8, poetically captured one of Verlan’s most important traits: it pursues confusion instead of clarity[21]. One reason is that it is widely used among lower-class people, drug users, gangs, or those in jail. Thus, making the context unidentifiable is important — certain phenomena such as reverlanisation (flipping the Verlan again if it becomes too popular) and truncation are therefore applied.

⁹github.com/greateden/verlan-Identification-Normalisation

However, although Verlan is used for concealing meaning, it still follows certain rules. The most general rule is syllabic reversal, as mentioned in the introduction chapter of this report.

Specifically, to delve into the linguistic rules, V  la pointed out that the analytic model proposed by Kaye and Lowenstamm provides the best description[22]. The syllable can be disassembled into *attaque* (onset), *rime* (rhyme), *noyau* (nucleus), and *coda*. For example, here is a representation of the word *garder*, IPA¹⁰ [garde].

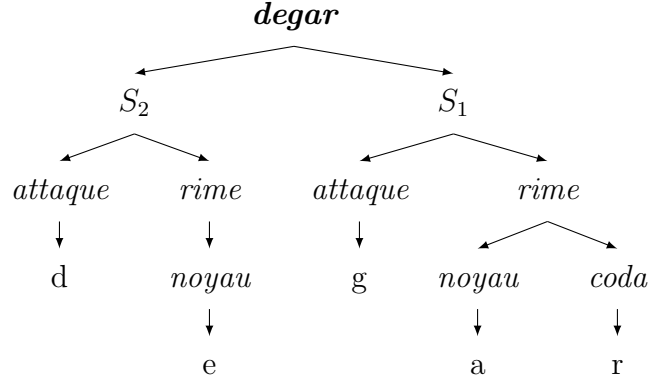


It has two syllables, S_1 and S_2 . To create the Verlan form, we follow the permutation equation below:

$$(S_1S_2) \rightarrow (S_2S_1) \quad (1)$$

After the permutation, we obtain the Verlan form of *garder* as *degar*, represented below.

¹⁰International Phonetic Alphabet, https://en.wikipedia.org/wiki/International_Phonetic_Alphabet



Notably, the permutation occurs only at the syllable level (i.e., between S_1 and S_2); it does not affect the internal structure of each syllable tree, although in some cases, certain letters (such as *e*) might be dropped after permutation. That said, the example above is not an exhaustive explanation of forming a Verlan. To avoid confusing the readers, this report suggests that this example perfectly illustrates its regular rule. For further details, readers are advised to consult V  la’s paper.

With such a sub-word permutation, researchers can not only discuss it within the linguistic realm, but it is also intriguing for computer scientists to explore how machines, such as LLMs, perceive this kind of difference. Just as V  la describes Verlan — ambiguous, sometimes violent, sometimes amazing, and always vivid.

2.2 Detecting Slang

To the best of our knowledge, there is no existing computational research¹¹ on the *detection* of Verlan — this particular form of French slang. However, there are a few scholars who have included Verlan in their research[9, 10, 11, 12]. Yet, these studies commonly included Verlan as a type of slang in their datasets or corpora. Moreover, they did not specifically focus on how to detect this particular type of slang, but rather approached it in a broader sense — they created slang datasets that contain Verlan, and some of them employed computational approaches to detect such slang.

Fortunately, there are several papers related to computational slang detection, and their approaches could contribute to Verlan detection to a large

¹¹As of September 2025.

extent[13, 14, 15, 18]. These studies are not limited to French but also cover other Indo-European languages¹².

Therefore, regarding the history of Verlan detection, this report first generalises the task as slang detection, and then discusses possible methods that could be implemented for Verlan identification, in order to provide readers with a general and useful background.

2.2.1 1910s-2016: A Super-Condensed History of Slang Detection

The background of traditional slang detection often leverages fuzzy-matching methods. Two main methods were introduced and widely cited: Soundex, a phonetic indexing system for names introduced by Russell in 1918, and the edit-distance-based spelling-correction method introduced by Levenshtein in 1966[23, 24]. Afterwards, scholars introduced more algorithms, such as Philips’s Metaphone and Double Metaphone, which improved on Russell’s Soundex; Kukich’s methods for detecting and correcting spelling; Sproat’s normalisation of Non-Standard Words (NSW); and Aw et al.’s phrase-based Machine Translation (MT) approach for standardising SMS messages[25, 26, 27, 28, 29]. While these are not directly slang-detection research, over time their methodology became increasingly related to slang — some slang can be treated as misspelling or NSW, and people frequently use slang in text messages.

2.2.2 2016-2019: Dictionary Search

The easiest way we can think of dealing with slang is to use a dictionary — just like how we look up a word that we do not know. The pros and cons are highly similar to consulting a dictionary. It is fast (if using a digital one) and accurate. On the other hand, because it is purely fixed data, it only works with existing words and thus cannot identify newly invented ones.

Examples of existing slang dictionaries include SlangNet, SlangSD, and SLANGZY[17, 18, 19]. As for French slang dictionaries, we have, for example, *Dictionnaire du chilleur*[20]. Specifically for Verlan, the report identifies several online dictionaries, including *Dictionnaire Interactif du Verlan*¹³,

¹²For example, English, German, and Russian. For more information, please refer to: https://en.wikipedia.org/wiki/Indo-European_languages.

¹³<https://ecoleng.com/verlan-comprendre-argot-francais-parler/dictionnaire-interactif-du-verlan>

Wiktionary¹⁴, and *Dictionnaire Verlan*¹⁵.

With these existing dictionaries, implementing a tool to identify Verlan should be straightforward. However, two major issues limit the possibility of directly using these dictionaries for Verlan identification: they lack comprehensive coverage, and some are fan-made, which neither captures the full extent of this slang nor guarantees accuracy. Licensing for certain dictionaries could also be a concern.

Although dictionaries have the drawbacks mentioned above, they remain essential resources for implementing LLM-based approaches, as discussed later. Consequently, new dictionaries continue to be produced.

2.2.3 Meanwhile, for Fuzzy Search

The 2010s belonged to social media and research on user-generated text. Representative work includes Beaufort et al.’s hybrid finite-state framework for SMS normalisation, Han and Baldwin’s lexical normalisation for Twitter, and the W-NUT shared tasks on Twitter message normalisation[30, 31, 32].

While these works are not directly about slang recognition, they provided immensely useful background for the research specifically on slang that followed.

2.2.4 2020-2025: Fuzzy Search + Slang Corpus = BOOM

In the 2020s, everyone tended to check what could be done with Machine Learning (ML) for this task, using Natural Language Processing (NLP). Researchers started to apply NLP to slang detection. Wilson’s paper used two million entries from *Urban Dictionary*, with terms, definitions, examples, and tags[33]. They pre-processed the dataset with techniques like lowercasing and removal of punctuation, followed by training a fastText¹⁶ skip-gram for 10 epochs with a 300-dimensional vector space. Using a fastText classifier, they analysed properties such as sentiment and sarcasm. For evaluation, they used accuracy, precision, recall, and F1 score.

Notably, the report has found two theses highly related to this project, *Slang or not?* and *Toward Informal Language Processing*[14, 15]. Both cre-

¹⁴<https://en.wiktionary.org/wiki/Category%3AVerlan>

¹⁵https://zlang.fandom.com/fr/wiki/Dictionnaire_Verlan

¹⁶A library for learning of word embeddings and text classification created by Facebook’s AI Research (FAIR) lab.

ated their own datasets that were manually annotated and validated. The former compared the performance of traditional ML (SVM¹⁷-linear with TF-IDF¹⁸ + n-grams), Convolutional Neural Network (CNN)¹⁹ / Bidirectional Long-Short Term Memory (BiLSTM)²⁰ with Bidirectional Encoder Representations from Transformers (BERT)²¹ embeddings, Transformer models (e.g., BERT-large-uncased), and Large Language Models (LLMs) (GPT-4o and GPT-4o-mini), finding that a fine-tuned Transformer performed best. The latter compared traditional baselines, Language Models (LMs), and LLMs.

2.2.5 Detecting Verlan?

The results from the last section provide this report with a clear guideline regarding Verlan identification. They have absorbed and adapted the historical development of slang detection into a modern, up-to-date framework. Building upon these insights, this report argues that BERT and contemporary LLMs represent the most effective tools for the Verlan detection task.

3 Dataset

¹⁷Support Vector Machine, https://en.wikipedia.org/wiki/Support_vector_machine

¹⁸<https://en.wikipedia.org/wiki/Tf%E2%80%93idf>

¹⁹https://en.wikipedia.org/wiki/Convolutional_neural_network

²⁰https://en.wikipedia.org/wiki/Long_short-term_memory

²¹[https://en.wikipedia.org/wiki/BERT_\(language_model\)](https://en.wikipedia.org/wiki/BERT_(language_model))

4 Conclusion

References

- [1] Radjabov, Ruslan Rajabmurodovich. *Understanding "verlan" in the French Language*. Web of Scientist: International Scientific Research Journal, vol. 6, no. 3, 2025, pp. 368-372. Available at: <https://webofjournals.com/index.php/3/article/view/3264>.
- [2] Bach, Xavier. *Tracing the origins of verlan in an early nineteenth century text*. Journal of French Language Studies, vol. 28, no. 1, 2018, pp. 1-18. Cambridge University Press. doi:10.1017/S0959269516000221.
- [3] Olivier Sécardin. *Évolution du verlan, marqueur social et identitaire, comme reflet de la langue et de la société françaises*. Synergies Europe, no. 3, 2008, pp. 223-232. Available at: <https://journal.lib.uoguelph.ca/index.php/synergies/article/download/1037/1859?inline=1>.
- [4] Rúa, Paula López. "Shortening Devices in Text Messaging." *Journal of Computer-Mediated Communication*, vol. 10, no. 4, July 2005. Wiley. doi:10.1111/j.1083-6101.2005.tb00268.x.
- [5] Hajiyevea, Bulbul. "Translating Idioms and Slang: Problems, Strategies, and Cultural Implications." *Acta Globalis Humanitatis et Linguarum*, vol. 2, no. 2, 2025, pp. 284-293. doi:10.69760/aghel.025002123.
- [6] DeepL. "DeepL Translator translates texts using artificial neural networks. These networks are trained on many millions of translated texts." *DeepL Blog*, 2020. Available at: <https://www.deepl.com/en/blog/how-does-deepl-work>.
- [7] Wu, Yonghui, et al. "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation." arXiv preprint arXiv:1609.08144, 2016. Available at: <https://arxiv.org/abs/1609.08144>.
- [8] Michel, Paul, and Graham Neubig. "MTNT: A Testbed for Machine Translation of Noisy Text." *Proceedings of EMNLP*, 2018. Available at: <https://aclanthology.org/D18-1050/>.

- [9] Zurbuchen, Lucas, and Rob Voigt. *A Computational Analysis and Exploration of Linguistic Borrowings in French Rap Lyrics*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics — Student Research Workshop (ACL SRW 2024)*, 2024, pp. 200-208. DOI: 10.18653/v1/2024.acl-srw.27.
- [10] Podhorná-Polická, Alena. *RapCor, Francophone Rap Songs Text Corpus*. In *Proceedings of the Fourteenth Workshop on Recent Advances in Slavonic Natural Language Processing (RASLAN 2020)*, 2020, pp. 95-102. Available at: <https://nlp.fi.muni.cz/raslan/raslan20.pdf#page=95>.
- [11] Mekki, Jade; Lecorvé, Gwénolé; Battistelli, Delphine; Béchet, Nicolas. *TREMoLo-Tweets: A Multi-Label Corpus of French Tweets for Language Register Characterization*. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, Held Online, INCOMA Ltd., Sep 1-3, 2021, pp. 950-958. DOI: 10.26615/978-954-452-072-4_108.
- [12] Panckhurst, Rachel; Lopez, Cédric; Roche, Mathieu. *A French text-message corpus: 88milSMS. Synthesis and usage*. Corpus [En ligne], 20 — 2020 (mis en ligne le 28 janvier 2020). DOI: 10.4000/corpus.4852.
- [13] Pei, Zhengqi, Zhewei Sun, and Yang Xu. *Slang Detection and Identification*. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL 2019)*, Hong Kong, China, 2019, pp. 881-889. Available at: <https://aclanthology.org/K19-1082/>.
- [14] Sun, Zhewei, Qian Hu, et al. *Toward Informal Language Processing: Knowledge of Slang in Large Language Models*. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2024)*, 2024. DOI: 10.18653/v1/2024.naacl-long.94.
- [15] Anonymous. *Slang or Not? Exploring NLP Techniques for Slang Detection Using the SlangTrack Dataset*. ACL ARR (OpenReview) submission, December 2024 (ACL ARR 2024 December). Available at: <https://openreview.net/forum?id=bIS03DD8sU>.

- [16] Wu, Tianyang; Morstatter, Fred; Liu, Huan; et al. *SlangSD: Building, Expanding, and Using a Sentiment Dictionary of Slang Words for Short-Text Sentiment Classification*. Language Resources and Evaluation (2018). DOI: 10.1007/s10579-018-9416-0.
- [17] Dhuliawala, Shehzaad; Kanojia, Diptesh; Bhattacharyya, Pushpak. *SlangNet: A WordNet like Resource for Slang Words*. In: Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016). Portorož, Slovenia (2016). Available at: <https://www.cse.iitb.ac.in/~pb/papers/lrec16-slangnet.pdf>.
- [18] Wu, Tianyang; Morstatter, Fred; Liu, Huan; et al. *SlangSD: Building, Expanding, and Using a Sentiment Dictionary of Slang Words for Short-Text Sentiment Classification*. Language Resources and Evaluation (2018). DOI: 10.1007/s10579-018-9416-0.
- [19] Gupta, Vishal; Rani, Rekha; et al. *SLANGZY: A Slang Word Recognition System for Hindi-English Code-Mixed Social Media Text*. In: Proceedings of the 6th Workshop on South and Southeast Asian Natural Language Processing (WSSANLP 2019). Kolkata, India (2019). Available at: <https://aclanthology.org/K19-1082.pdf>.
- [20] Parent, Philippe; Parent, André. *Dictionnaire du chilleur*. Éditions Somme toute (2024). ISBN: 9782925124351.
- [21] Méla, Vivienne. *Le verlan ou le langage du miroir*. Langages, No. 101, Les javanais (Mars 1991), pp. 73–94. Published by Armand Colin. Available at: <https://www.jstor.org/stable/23906698>.
- [22] Kaye, Jonathan D.; Lowenstamm, Jean. *De la syllabité*. In: Dell, François; Hirst, Daniel; Vergnaud, Jean-Roger (eds.), *Forme sonore du langage*. Hermann, Paris (1984), pp. 123–159. Available at: <https://archive.org/details/formesonoredulangage>.
- [23] Russell, Robert C.; Odell, Margaret K. *Soundex system of indexing names*. U.S. Patent 1,261,167, filed June 3, 1918, and issued April 2, 1918. Available at: <https://patents.google.com/patent/US1261167A/en>.
- [24] Levenshtein, Vladimir I. *Binary codes capable of correcting deletions, insertions, and reversals*. Soviet Physics Doklady, vol. 10, no. 8, 1966,

- pp. 707-710. Available at: <https://nymity.ch/sybilhunting/pdf/Levenshtein1966a.pdf>.
- [25] Philips, Lawrence. *Hanging on the Metaphone*. Computer Language (1990). Available at: <https://aspell.net/metaphone/>.
 - [26] Philips, Lawrence. *The Double Metaphone Search Algorithm*. C/C++ Users Journal (June 2000). Available at: <https://xlinux.nist.gov/dads/HTML/doubleMetaphone.html>.
 - [27] Kukich, Karen. *Techniques for automatically correcting words in text*. ACM Computing Surveys (1992). DOI: 10.1145/146370.146380.
 - [28] Sproat, Richard; Black, Alan W.; Chen, Stanley; Kumar, Shankar; Ostendorf, Mari; Richards, Christopher. *Normalization of non-standard words*. Computer Speech & Language, 15(3):287–333 (2001). DOI: 10.1006/csla.2001.0169.
 - [29] Aw, AiTi; Zhang, Min; Xiao, Juan; Su, Jian. *A Phrase-Based Statistical Model for SMS Text Normalization*. Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions (2006), pages 33–40. Available at: <https://aclanthology.org/P06-2005/>.
 - [30] Beaufort, Richard; Roekhaut, Sophie; Cougnon, Louise-Amélie; Fa-iron, Cédric. *A Hybrid Rule/Model-Based Finite-State Framework for Normalizing SMS Messages*. Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics (ACL 2010). Available at: <https://aclanthology.org/P10-1079.pdf>.
 - [31] Han, Bo; Baldwin, Timothy. *Lexical Normalisation of Short Text Mes-sages: Makn Sens a #twitter*. Proceedings of the 49th Annual Meet-ing of the Association for Computational Linguistics: Human Language Technologies (ACL-HLT 2011), pages 368–378. Available at: <https://aclanthology.org/P11-1038/>
 - [32] Baldwin, Timothy; de Marneffe, Marie Catherine; Han, Bo; Kim, Young-Bum; Ritter, Alan; Xu, Wei. *Shared Tasks of the 2015 Work-shop on Noisy User-generated Text: Twitter Lexical Normalization and Named Entity Recognition*. Proceedings of the Workshop on Noisy User-generated Text (W-NUT 2015). DOI: 10.18653/v1/W15-4319.

- [33] Urban Dictionary Embeddings. *Urban Dictionary Embeddings for Slang NLP Applications*. LREC / ACL Anthology (2020). Available at: <https://aclanthology.org/2020.lrec-1.586/>.
- [34] Sun, X.; et al. *Knowledge of Slang in Large Language Models*. Proceedings of the 2024 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-Long 2024). Available at: <https://aclanthology.org/2024.naacl-long.94/>.

Appendix A Some extra things

If you have anything more to add such as:

- not essential details - things that might be too much for first time reading, or could be distracting from the main points...but are still important for reproducibility or deeper understanding
- work that was done in the project but doesn't go with the main work, or detracts/is not essential for the main narrative.

Appendix B Aims and Objectives

Interim report only! – you do not need to include this appendix in the final report. However, in your interim the last appendix should include your original Aims and Objectives, and, if the things have changed, the revised Aims and Objectives. If you used the L^AT_EX template provided for your Aims and objectives document, just copy the `\paragraph{Aims}` and `\paragraph{Objectives}` sections and paste them here.

Original

Aims Here you are describing the term goal of the project. What do you want to achieve by the end? What is the ultimate goal of this work? For example, the primary aim of this document is to have students produce suitable aims and objectives for their COSC480/490 project. While the aims and objectives document is not an assessed deliverable, a clear definition of what is to be done, and a bit of planning of how it is to be accomplished is

paramount to the project's success. It is important to establish the scope of the project.

Objectives Objectives list the milestones that you need to achieve in order to achieve the project's aim(s). It's a rough plan for what needs to happen in what order. It's best to list the objectives in bullet point form. For many projects the structure to these objectives might follow the following pattern (objective names are just examples – you can have different objective names):

- background reading; going through the literature; learning about the research field;
- setting up of some kind of system for the project; getting the environment for experiments working;
- conducting preliminary experiments; implementation of a basic/simple approach; producing base case results;
- trying method 1; recording the results;
- trying method 2; recording the results.

Revised

Aims Here you are describing the term goal of the project. What do you want to achieve by the end? What is the ultimate goal of this work? For example, the primary aim of this document is to have students produce suitable aims and objectives for their COSC480/490 project. While the aims and objectives document is not an assessed deliverable, a clear definition of what is to be done, and a bit of planning of how it is to be accomplished is paramount to the project's success. It is important to establish the scope of the project.

Objectives Objectives list the milestones that you need to achieve in order to achieve the project's aim(s). It's a rough plan for what needs to happen in what order. It's best to list the objectives in bullet point form. For many projects the structure to these objectives might follow the following pattern (objective names are just examples – you can have different objective names):

- background reading; going through the literature; learning about the research field;
- setting up of some kind of system for the project; getting the environment for experiments working;

- conducting preliminary experiments; implementation of a basic/simple approach; producing base case results;
- trying method 1; recording the results;
- trying method 2; recording the results.