

Transformation rules and Monte-Carlo sampling: a different approach for statistical paraphrase generation

Jonathan Chevelu^{1,2} Thomas Lavergne Yves Lepage¹ Thierry Moudenc²

(1) GREYC, universit  de Caen Basse-Normandie

(2) Orange Labs; 2, avenue Pierre Marzin, 22307 Lannion
{jonathan.chevelu,thierry.moudenc}@orange-ftgroup.com,
thomas.lavergne@reveurs.org, yves.lepage@info.unicaen.fr

Abstract

Paraphrase generation is often presented as a monolingual statistical machine translation problem. This approach cannot take advantage of paraphrases particularities by transforming only parts of sentences. We propose a different paradigm for statistical paraphrase generation where a paraphrase is seen as the application of a set of transformation rules on a sentence. We propose a new method, adapted to this point of view, based on *Monte-Carlo* sampling and show how this algorithm is suitable for paraphrase generation. Moreover, the basic algorithm presented here leaves a lot of opportunities for future improvement. In particular, our algorithm does not constraint the scoring function in opposite to *Viterbi* based decoders. It is now possible to use some global features in paraphrase scoring functions. This algorithm opens new outlooks for paraphrase generation and other natural language processing applications like statistical machine translation.

1 Introduction

A paraphrase generation system is a program which, given a source sentence, produces a different sentence with almost the same meaning.

Paraphrase generation is useful in applications to choose between different forms to keep the most appropriate one. For instance, automatic summary can be seen as a particular paraphrasing task (Barzilay and Lee, 2003) with the aim of selecting the shortest paraphrase. Paraphrases can also help human writers by suggesting possible alternatives and having them choose the most appropriate ones (Max and Zock, 2008).

Paraphrases can also be used to improve natural language processing (NLP) systems. Callison-

Burch et al. (2006) improved machine translations by augmenting the coverage of patterns that can be translated. Similarly, Sekine (2005) and Duclaye et al. (2003) improved information retrieval and question-answering systems based on pattern recognition by introducing paraphrase generation in such systems.

In order to produce paraphrases, a promising approach is to see the paraphrase generation problem as a translation problem, where the target language is the same as the source language (Quirk et al., 2004; Bannard and Callison-Burch, 2005; Max and Zock, 2008).

In the frame of statistical machine translation (SMT), the first problem for this approach is the need of a paraphrase table. Such a table is the monolingual version of the bilingual translation table. In SMT, aligned multilingual corpora like Europarl (Koehn, 2005) are used to produce translation tables.

The lack of equivalent monolingual corpora leads to the use of heuristics to produce them (Barzilay and Lee, 2003; Quirk et al., 2004). Bannard and Callison-Burch (2005) proposed to produce paraphrase tables using translation tables learnt from bilingual corpora. They used a standard heuristic (Lepage and Denoual, 2005) that consists in calling paraphrases a pair of phrases with the same translation.

A second problem that has drawn less attention is the generation step which corresponds to the decoding step in SMT. Most paraphrase generation tools use some standard SMT decoding algorithms (Quirk et al., 2004) or some off-the-shelf decoding tools like MOSES (Koehn et al., 2007). The goal of a decoder is to find the best path in the lattice produced from a paraphrase table. This is basically achieved by using dynamic programming and especially the *Viterbi* algorithm associated with beam searching (Koehn et al., 2007).

However decoding algorithms were designed

for translation, not for paraphrase generation. Although left-to-right decoding is justified for translation, it may not be necessary for paraphrase generation. A paraphrase generation tool usually starts with a sentence which may be very similar to some potential solution. In other words, there is no need to "translate" all of the sentences. Moreover, decoding may not be suitable for non-contiguous transformation rules.

In addition, dynamic programming imposes an incremental scoring function to evaluate the quality of each hypothesis (Langlais et al., 2008). It is easy to imagine translation or paraphrase problems where this kind of function is not appropriate. For instance, it cannot capture some scattered syntactical dependencies. Improving on this major issue is a key point to improve paraphrase generation systems.

All these problems lead us to propose a different paradigm for statistical paraphrase generation and an associated generation algorithm.

This paper first presents an alternative to decoding that is based on transformation rule application in section 2. In section 3 we propose a paraphrase generation method for this paradigm based on an algorithm used in two-player games. Section 4 briefly explain experimental context and its associated protocol for evaluation of the proposed system. We compare the proposed algorithm with a baseline system in section 5. Finally, in section 6, we point to future research tracks to improve paraphrase generation tools.

2 Statistical paraphrase generation using transformation rules

The paraphrase generation problem can be seen as an exploration problem. We seek the best paraphrase according to a scoring function in a space to search by applying successive transformations. This space is composed of states connected by actions. An action is a transformation rule with a place where it applies in the sentence. States are a sentence with a set of possible actions. Applying an action in a given state consists in transforming the sentence of the state and removing all rules that are no more applicable. In our framework, each state, except the root, can be a final state. This is modeled by adding a stop rule as a particular action. We impose the constraint that any transformed part of the source sentence cannot be transformed anymore.

This paradigm is more appropriate for paraphrase generation than the standard SMT approach in respect to several points:

- there is no need for left-to-right decoding: a transformation can be applied anywhere without order;
- there is no need to transform the whole of a sentence: each state is a final state;
- there is no need to keep the identity transformation for each phrase in the paraphrase table;
- the only domain knowledge needed is a generative model and a scoring function for final states;
- it is possible to mix different generative models: a statistical paraphrase table, an analogical solver and a paraphrase memory for instance;
- there is no constraint on the scoring function: it only scores final states.

Note that the branching factor with a paraphrase table can be around thousand actions per states which makes the generation problem a difficult computational problem. Hence we need an efficient generation algorithm.

3 Monte-Carlo based Paraphrase Generation

UCT (Kocsis and Szepesvári, 2006) (*Upper Confidence bound applied to Tree*), a *Monte-Carlo* planning algorithm, has recently become popular in two-player game problems.

UCT has some interesting properties:

- it grows the search tree non-uniformly and favours the most promising sequences, without pruning branch;
- it can deal with high branching factor;
- it is an any-time algorithm and returns best solution found so far when interrupted;
- it does not require expert domain knowledge to evaluate states.

These properties make it ideally suited for problems with high branching factor and for which there is no strong evaluation function.

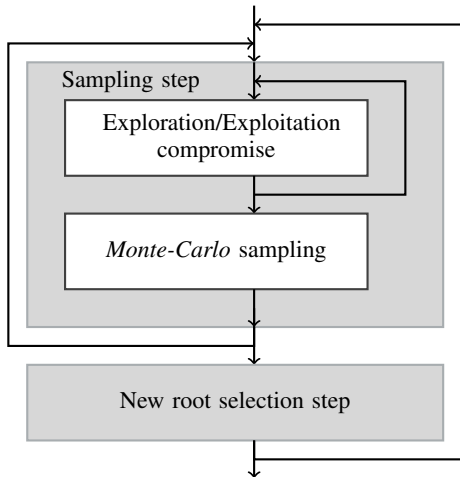


Figure 1: The MCPG algorithm.

For the same reasons, this algorithm sounds interesting for paraphrase generation. In particular, it does not put constraint on the scoring function. We propose a variation of the UCT algorithm for paraphrase generation named MCPG for *Monte-Carlo based Paraphrase Generation*. A diagram of the MCPG algorithm is presented figure 1.

The main part of the algorithm is the sampling step. An episode of this step is a sequence of states and actions, $s_1, a_1, s_2, a_2, \dots, s_T$, from the root state to a final state. During an episode construction, there are two ways to select the action a_i to perform from a state s_i .

If the current state was already explored in a previous episode, the action is selected according to a compromise between exploration and exploitation. This compromise is computed using the UCB-Tunned formula (Auer et al., 2001) associated with the RAVE heuristic (Gelly and Silver, 2007). If the current state is explored for the first time, its score is estimated using *Monte-Carlo* sampling. In other word, to complete the episode, the actions $a_i, a_{i+1}, \dots, a_{T-1}, a_T$ are selected randomly until a stop rule is drawn.

At the end of each episode, a reward is computed for the final state s_T using a scoring function and the value of each (state, action) pair of the episode is updated. Then, the algorithm computes another episode with the new values.

Periodically, the sampling step is stopped and the best action at the root state is selected. This action is then definitely applied and a sampling is restarted from the new root state. The action sequence is built incrementally and selected af-

ter being enough sampled. For our experiments, we have chosen to stop sampling regularly after a fixed amount η of episodes.

Our main adaptation of the original algorithm is in the (state, action) value updating procedure. Since the goal of the algorithm is to maximise a scoring function, we use the maximum reachable score from a state as value instead of the score expectation. This algorithm suits the paradigm proposed for paraphrase generation.

4 Experimental context

This section describes the experimental context and the methodology followed to evaluate our statistical paraphrase generation tool. Section 4.1 describes the corpus used in the experiments. Section 4.2 briefly presents the construction of the language model and the paraphrase table. Finally, Section 4.3 deals with the evaluation protocol.

4.1 Data

For the experiment reported in section 5, we use one of the largest, multi-lingual, freely available aligned corpus, Europarl (Koehn, 2005). It consists of European parliament debates. We choose French as the language for paraphrases and English as the pivot language. For this pair of languages, the corpus consists of 1,487,459 French sentences aligned with 1,461,429 English sentences. Note that the sentences in this corpus are long, with an average length of 30 words per French sentence and 27.1 for English. We randomly extracted 100 French sentences as a test corpus.

4.2 Language model and paraphrase table

Paraphrase generation tools based on SMT methods need a language model and a paraphrase table. Both are computed on a training corpus.

The language models we use are n-gram language models with back-off. We use SRILM (Stolcke, 2002) with its default parameters for this purpose. The length of the n-grams is five.

To build a paraphrase table, we use a variant of the construction method via a pivot language proposed in Bannard and Callison-Burch (2005). The first step consists in building a bilingual translation table from the aligned corpus. Given a source phrase f^i and another phrase e^i in a different language, a bilingual translation table provides the two probabilities $p(f^i|e^i)$ and $p(e^i|f^i)$. We use

GIZA++ (Och and Ney, 2003) with its default parameters to produce phrase alignments. The paraphrase table is then built from the phrase translation table. The probability for a phrase f^i to be paraphrased by a phrase f'^i in the same language is estimated by the sum of each round-trip from f^i to f'^i through any phrase e^i of a pivot language.

The construction of this table is very simple. Given a bilingual translation table sorted by pivot phrases, the algorithm retrieves all the phrases linked with the same pivot (named a *pivot cluster*). For each ordered pair of phrases, the program assigns a probability that is the product of these probabilities. This process realizes a self-join of the bilingual translation table. It produces a paraphrase table composed of tokens, instead of items. The program just needs to sum up all probabilities for all entries with identical paraphrase tokens to produce the final paraphrase table.

Three heuristics are used to prune the paraphrase table. The first heuristic prunes any entry in the paraphrase table composed of tokens with a probability lower than a threshold ϵ . The second, called *pruning pivot heuristic*, consists in deleting all pivot clusters larger than a threshold τ . The last heuristic keeps only the κ most probable paraphrases for each source phrase in the final paraphrase table. For this study, we empirically fix $\epsilon = 10^{-5}$, $\tau = 200$ and $\kappa = 10$.

Due to our proposed paradigm, we delete the rephrasing of a phrase into itself, called *identity paraphrase*, from the paraphrase table for our system. However, for each entry, we add the probability of the associated identity paraphrase. To be exact, this probability is biased to be the same as the most probable transformation of the associated source phrase. The reason is that identity paraphrase should be the best "transformation" in terms of preservation of the meaning.

4.3 Evaluation Protocol

As for assessment, automatic paraphrase evaluation is a difficult problem. Most of evaluation procedures are based on a human based judgement. Unfortunately, there is no consensus on a universally accepted human-based evaluation protocol. This state of affairs led us to develop our own protocol.

We developed a dedicated website to allow the human judges with some flexibility in workplaces and evaluation periods. We retain the principle of

the two-step evaluation, common in the machine translation domain and already used for paraphrase evaluation (Bannard and Callison-Burch, 2005). One step evaluates the syntax of paraphrases to determine if they are well formed. The other step evaluates the semantics of paraphrases. It is designed to know if the meaning of the source sentence is preserved in the paraphrase.

The question asked to the human evaluator for the syntactic task is:

Is the following sentence in good French?

The question asked to the human evaluator for the semantic task is:

Do the following two sentences express the same thing?

In our experiments, each paraphrase was evaluated by two native French evaluators.

5 Comparison with a SMT decoder

In order to validate our algorithm for paraphrase generation, we compare it with an off-the-shelf SMT decoder.

We use the MOSES decoder (Koehn et al., 2007) as a baseline. The MOSES scoring function is set by four weighting factors α_Φ , α_{LM} , α_D , α_W . Conventionally, these four weights are adjusted during a tuning step on a training corpus. The tuning step is inappropriate for paraphrase because there is no such tuning corpus available. We empirically set $\alpha_\Phi = 1$, $\alpha_{LM} = 1$, $\alpha_D = 10$ and $\alpha_W = 0$. This means that the paraphrase table and the language model are given the same weight, no reordering is allowed and no specific sentence length is favored. Hence, the scoring function (or reward function for MCPG) is equivalent to:

$$R(f'|f, I) = p(f') \times \Phi(f|f', I)$$

where f and f' are the source and target sentences, I a segmentation in phrases of f , $p(f')$ the language model score and $\Phi(f|f', I) = \prod_{i \in I} p(f^i|f'^i)$ the paraphrase table score.

The MCPG algorithm needs two parameters. One is the number of episodes η done before selecting the best action at root state. The other is k , an *equivalence parameter* which balances the exploration/exploitation compromise (Auer et al., 2001). We empirically set $\eta = 1,000,000$ and $k = 1,000$.

For our algorithm, note that identity paraphrase probabilities are biased: for each phrase it is

System	MOSES	MCPG
Well formed (Kappa)	64%(0.57)	63%(0.84)
Meaning preserved (Kappa)	58%(0.48)	55%(0.64)
Well formed and meaning preserved (Kappa)	50%(0.54)	49%(0.59)

Table 1: Results of paraphrases evaluation for 100 sentences in French using English as the pivot language. Comparison between the baseline system MOSES and our algorithm MCPG.

equal to the probability of the most probable paraphrase. Moreover, as the source sentence is the best meaning preserved "paraphrase", a sentence cannot have a better score. Hence, we use a slightly different scoring function:

$$\begin{aligned}
 R(f'|f, I) &= \min \left(\frac{p(f') \times \Phi(f|f', I)}{p(f) \times \Phi(f|f, I)}, 1 \right) \\
 &= \min \left(\frac{p(f')}{p(f)} \prod_{\substack{i \in I \\ f^i \neq f'^i}} \frac{p(f^i|f'^i)}{p(f^i|f^i)}, 1 \right)
 \end{aligned}$$

Note that for this model, there is no need to know the identity transformations probability for unchanged part of the sentence.

Results are presented in Table 1. For our algorithm, 63 sentences out of the 100 paraphrases were rated syntactically correct by both judges. Judges found 55 paraphrases which preserve the source sentence meaning. On the whole, 49% of the paraphrases are rated both syntactically and semantically correct by both judges. The Kappa statistics (Cohen, 1960) associated with the syntactic evaluation is 0.84, which is usually considered as a "perfect" agreement. For the evaluation of the meaning preservation, the Kappa statistic is 0.64, which is usually considered as a "substantial" agreement. On the whole, the Kappa statistic is 0.59, which is usually considered as a "moderate" agreement.

Results are close to evaluations from the baseline system. Because of the test corpus size, it is not statistically possible to difference the two systems. The probability for MOSES to be better than MCPG is only 0.54 with a 95% level of confidence. The main differences are from Kappa statistics which are lower for the MOSES system evaluation. Judges changed between the two experiments. We may wonder whether an evaluation with only two judges is reliable. This points to the ambiguity of any paraphrase definition.

By doing this experiment, we have shown that our algorithm with a biased paraphrase table is

state-of-the-art to generate paraphrases.

6 Conclusions and further research

In this paper, we have proposed a different paradigm and a new algorithm in NLP field adapted for statistical paraphrases generation. This method, based on large graph exploration by *Monte-Carlo* sampling, produces results comparable with state-of-the-art paraphrase generation tools based on SMT decoders.

This basic algorithm can be improved by easily adding new features hardly compatible with standard SMT decoders. The algorithm structure is flexible and generic enough to easily work with discontinuous patterns. It is also possible to mix various transformation methods to increase paraphrase variability.

Unlike SMT decoders, the output of the n-best solutions for MCPG is not more expensive than the output of the one-best solution.

The rate of ill-formed paraphrase is high at 37%. The result analysis suggests an involvement of the non-preservation of the original meaning when a paraphrase is evaluated ill-formed. Although the measure is not statistically significant because the test corpus is too small, the same trend is also observed in other experiments. Improving on the language model issue is a key point to improve paraphrase generation systems. Our algorithm can work with unconstrained scoring functions, in particular, there is no need for the scoring function to be incremental as for *Viterbi* based decoders. We are working to add, in the scoring function, a linguistic knowledge based analyzer to solve this problem.

Because MCPG is based on a different paradigm, its output scores cannot be directly compared to MOSES scores. In order to prove the optimisation qualities of MCPG versus state-of-the-art decoders, we are transforming our paraphrase generation tool into a translation tool.

This first introduction of MCPG lets some open research works. For now, the algorithm is using a

fixed number of iterations η before stopping each sampling step. A more in-depth study of the number of simulations needed, and possibly a better policy to stop sampling can improve its speed. Finally, the algorithm allows some local weighting in the exploration graph, for instance by introducing a prior on the action value function.

References

- P. Auer, N. Cesa-Bianchi, and C. Gentile. 2001. Adaptive and self-confident on-line learning algorithms. *Machine Learning*.
- Colin Bannard and Chris Callison-Burch. 2005. Paraphrasing with bilingual parallel corpora. In *ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 597–604. Association for Computational Linguistics, Morristown, NJ, USA.
- Regina Barzilay and Lillian Lee. 2003. Learning to paraphrase: An unsupervised approach using multiple-sequence alignment. In *HLT-NAACL 2003: Main Proceedings*, pages 16–23.
- Chris Callison-Burch, Philipp Koehn, and Miles Osborne. 2006. Improved statistical machine translation using paraphrases. In *Proceedings of the main conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, pages 17–24. Association for Computational Linguistics, Morristown, NJ, USA.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.
- Florence Duclaye, François Yvon, and Olivier Collin. 2003. Learning paraphrases to improve a question-answering system. In *Proceedings of the 10th Conference of EACL Workshop Natural Language Processing for Question-Answering*, page 3541.
- Sylvain Gelly and David Silver. 2007. Combining online and offline knowledge in UCT. In *24th International Conference on Machine Learning (ICML '07)*, pages 273–280.
- Levente Kocsis and Csaba Szepesvári. 2006. Bandit based monte-carlo planning. In *17th European Conference on Machine Learning, (ECML '06)*, pages 282–293.
- Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of MT Summit*.
- Philipp Koehn, Hieu Hoang, Alexandra Birch Mayne, Christopher Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Annual Meeting of the Association for Computational Linguistics (ACL), Demonstration Session*, pages 177–180.
- Philippe Langlais, Alexandre Patry, and Fabrizio Gotti. 2008. Recherche locale pour la traduction statistique par segments. In *TALN 2008*, pages 119–128. ATALA, Avignon, France.
- Yves Lepage and Etienne Denoual. 2005. Automatic generation of paraphrases to be used as translation references in objective evaluation measures of machine translation. In *IWP2005*.
- Aurélien Max and Michael Zock. 2008. Looking up phrase rephrasings via a pivot language. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 97–104. Coling 2008 Organizing Committee, Manchester, UK.
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.
- Chris Quirk, Chris Brockett, and Bill Dolan. 2004. Monolingual machine translation for paraphrase generation. In Dekang Lin and Dekai Wu, editors, *the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 142–149. Association for Computational Linguistics, Barcelona, Spain.
- Satoshi Sekine. 2005. Automatic paraphrase discovery based on context and keywords between ne pairs. In *Proceedings of International Workshop on Paraphrase (IWP2005)*.
- Andreas Stolcke. 2002. Srilmm – an extensible language modeling toolkit. In *Proceedings of International Conference on Spoken Language Processing*.