

A French Fairy Tale Corpus syntactically and semantically annotated

El Maarouf Ismail, Villaneau Jeanne

IRISA-UBS (UEB)

Rue de Saint Maudé, BP 92116 56321 Lorient France
elmaarouf.ismail@yahoo.fr, jeanne.villaneau@univ.ubs.fr

Abstract

Fairy tales, folktales and more generally children stories have lately attracted the Natural Language Processing (NLP) community. As such, very few corpora exist and linguistic resources are lacking. The work presented in this paper aims at filling this gap by presenting a syntactically and semantically annotated corpus. It focuses on the linguistic analysis of a Fairy Tales Corpus, and provides the description of the syntactic and semantic resources developed for Information Extraction. Resources include syntactic dependency relation annotation for 120 verbs; referential annotation, which is concerned with annotating each anaphoric occurrence and Proper Name with the most specific noun in the text; ontology matching for a substantial part of the nouns in the corpus; semantic role labelling for 41 verbs using the FrameNet database. The article also sums up previous analyses of this corpus and indicates possible uses of this corpus for the NLP community.

Keywords: Syntactic Annotation; Semantic annotation; Fairy Tales

1. Introduction

Fairy tales, folktales and more generally children stories have lately attracted the Natural Language Processing (NLP) community. For example, the LREC 2010 conference¹ welcomed at least three papers on such corpora; Fairy Tales markup Language have been defined to tag text sequences according to Propp's theory (see (Scheidel and Declerck, 2010) and references therein) and NLP research projects have recently been launched (e.g. FACT²).

Considered applications include Text classification, referring to work in Literature such as (Propp, 1968) or (Aarne and Thompson, 1973) and robot story tellers (see (Gelin et al., 2010; Theune et al., 2003) for instance), with a focus on expressive reading (Volkova et al., 2010). Such applications may benefit from Information Extraction (IE), a NLP task aimed at extracting entities, semantic and coreference relations from text. Children stories is a new domain of application for IE (mainly focused on newspapers and medical corpora; (Nadeau and Sekine, 2007)) and may reveal different or specific problems for NLP systems. One interesting issue concerns the specificity of their content: children stories are not always set in a "real" environment (such as one would expect in newspapers corpora), but feature magical beings as well as extraordinary events (motifs). This paper focuses on the linguistic analysis of a Fairy Tales Corpus, and provides description for the syntactic and semantic resources developed for IE in terms of classification and relation extraction.

To our knowledge, NLP ontologies such as Cyc³ or Wordnet⁴ do not provide a detailed classification (or "micro-theories") for fictional entities or events. The work pre-

sented in this paper aims at filling this gap by presenting a syntactically and semantically annotated corpus. Section 2 introduces the project and describes the corpus. Sections 3 and 4 introduce the syntactic and semantic annotations.

2. The Fairy Tales Corpus

2.1. Research project

The Fairy Tales Corpus (FTC) was originally collected in a project of the French National Research Agency, EmotiRob (Saint-Aimé et al., 2007). The goal of this project was to design an interactive companion robot for fragile children. The part of the project which is concerned here involves detecting emotion (happiness, sadness, etc.) through linguistic analysis (ASR transcripts): this is the task of EmoLogus (Le Tallec et al., 2010), a symbolic system which computes emotions on top of a semantic representation.

EmoLogus requires semantic knowledge (concepts and relations) to generate a semantic representation for each speech act. An emotion lexicon was created to cover words which could be used in children interactions, based on previous lexica and experiments in schools (see (Le Tallec et al., 2010) for further details). A corpus was then needed to extract sufficient context (linguistic) information for these words. The FTC was chosen as the best alternative since it is directed towards children and because it contained a large proportion of the targeted words (76% of nouns and 90% of verbs; (El Maarouf et al., 2009a)).

2.2. Objectives

The purpose of the annotation was to prepare and extract semantic information needed by EmoLogus. This involved the creation of a verb database where each meaning is connected to contextual patterns (see (Hanks, 2008) for a similar perspective). In order to preserve the specificity of the research context and of the corpus, the annotation was carried out in a corpus-based fashion (Sinclair, 1991): verbs were analysed one by one through their concordances and patterns were ranked according to their frequency. This

¹LREC 2010 (www.lrec-conf.org/lrec2010/); see also AMI-CUS (ilk.uvt.nl/amicus/) and the 2010 Symposium on Computational Models of Narratives (narrative.csail.mit.edu/fs10/)

²wwwhome.cs.utwente.nl/~hiemstra/2011/folktales-as-classifiable-texts.html

³www.cyc.com

⁴wordnet.princeton.edu

methodology allows to identify corpus-based semantic patterns for verbs, i.e. patterns as they appear in corpus. However, the task of syntactic relation extraction was led separately from pattern analysis in order to analyse what was involved in the creation of patterns: patterns were only merged in a second step. The most frequent full verbs (120 verbs of frequency >30) were selected for manual annotation (a sample is shown in Table 1). Only one annotator (linguist) took part in the annotation, therefore it was not possible to test the agreement rate.

Verbs	Freq.	trans.
dire	813	<i>to say</i>
trouver	299	<i>to find</i>
savoir	292	<i>to know</i>
venir	288	<i>to come</i>
prendre	275	<i>to take</i>
passer	254	<i>to pass</i>
demander	244	<i>to ask</i>
appeler	225	<i>to call</i>
partir	224	<i>to leave</i>
donner	199	<i>to give</i>
regarder	197	<i>to look</i>
entendre	181	<i>to hear</i>
sortir	169	<i>to go out</i>
répondre	165	<i>to answer</i>
manger	159	<i>to eat</i>
rester	149	<i>to stay</i>
décider	141	<i>to decide</i>
chercher	141	<i>to look for</i>
devenir	129	<i>to become</i>
penser	123	<i>to think</i>
commencer	121	<i>to begin</i>
tomber	121	<i>to fall</i>

Table 1: Annotated verbs wrt frequency

2.3. Corpus description

The FTC is a collection of tales extracted from a website and the text is copyrighted. 139 tales were manually collected, cleaned and checked. The corpus contains about 160 000 running words (the number of words per tale vary from 120 to 17000 words). As the website provided information regarding authorship (mainly age and place) for most of the tales, this information was conserved and used to classify the corpus (see Table 2).

Author	Freq.	%	Tales	%
Modern Adult	63217	39%	24	17%
Children	53109	34%	70	51%
Unknown	34314	21%	37	27%
Classic Adult	9900	6%	7	5%

Table 2: Author Categories wrt frequency and nb of Tales

The proportion of adult writing (either professional writers or not) is slightly greater than the proportion of children writing (e.g. in classroom activities) in terms of total

frequency, especially when Classic storytellers are taken into account (like the Little Red Ridding Hood). Children tend to write shorter stories, if we consider story mean length (758, as opposed to 2634 words per story for modern adults). The corpus also shows variety in terms of content: some stories involve the ‘ordinary’ (non-magical) everyday life of children while others focus on animal protagonists, fairies, witches or even aliens. In conclusion, the FTC is heterogeneous in terms of authorship and content, but its single audience (children) provides for a dimension of homogeneity.

3. Syntactic Annotation

The corpus was first automatically tokenized, lemmatized and Part-Of-Speech-tagged with the Tree-Tagger (Schmid, 1994). Tag errors were corrected as the annotation progressed.

3.1. Annotation scheme

All the words (headwords) syntactically linked to each verb were extracted and labelled using a set of categories. Each relation consists in a triple <R,V,A > where R stands for the name of the relation and V and A stand for the verb and its argument respectively. Full syntactic annotation was performed: both arguments and adjuncts were annotated and non-finite verb forms were included.

The scheme shares similarities with PropBank (Palmer et al., 2005) and the Stanford Typed Dependency ((de Marneffe and Manning, 2008)) conventions. General syntactic categories (subject, object, indirect object) are used when possible. When prepositions (simple or compounded) introduce a complement (nominal or verbal), they are used as the category label. Adjectives were generally labelled as Qualitative (QUAL).

More specific labels were added to account for subtle distinctions (especially for Adverbs and Pronouns). Instead of labelling adverbs with a broad Adjunct category, they were discriminated thanks to general semantic categories. Example (1) describes a manner adverb and example (2) illustrates a quantity adverb.

- (1) Duchesse s’approcha doucement.
trans. Duchesse approached slowly.
 (2) Biribi s’avance un peu plus.
trans. Biribi approaches a little more.

French makes great use of pronouns (Blanche-Benveniste, 1990), and the same pronoun may function differently according to context: in (3), the word “en” (which may be loosely translated as a contraction of “of it”) plays the role of a location source whereas it functions as the syntactic object in (4).

- (3) Les bisons ne peuvent presque plus en sortir.
trans. buffaloes could hardly get out of it.
 (4) On en prenait des quantités raisonnables.
trans. We took reasonable quantities of it.

Other conventions include:

- Phrasal verbs with specific meaning are given a separate index: pronominal verbs like “se pousser” (5), causatives like “faire remarquer” (6) and combinations like “se faire remarquer” (7).

(5) “Pousse-toi, triple idiot !” dit Jennyfer.
trans. “move over, you prize idiot !” said Jennyfer.

(6) Mais il m’a fait remarquer que l’on pourrait toujours le faire plus tard.
trans. But he pointed out to me that it could be done later.

(7) Il faut toujours que tu te fasses remarquer.
trans. You always behave conspicuously.

- As in PropBank, quotatives (not headwords) such as “pousse toi, triple idiot !” illustrated in (5) were annotated.
- In ambiguous contexts, the most informative head rather than the syntactic head (first governing noun) was chosen as the argument. This often occurs with collective nouns (8).

(8) Il prit un morceau de bois.
trans. He took a piece of wood.

- In special cases where verbs lack an overt subject (e.g. in imperative mood), the argument of the relation is the verb itself (9).

(9) Prends le couloir et la première porte à droite.
trans. Take the corridor and the first door to your right.

- Every headword of coordinated groups should be selected on an equal basis (duplication of syntactic relations). In (10), the Objects are both “ballon” and “peluche”.

(10) Il prit le ballon et par la même occasion la peluche.
trans. He took the ball and the toy at the same time.

- Since non-finite forms are included, the same argument may be linked to more than one verb (as in control and raising verbs), coordinated forms included. In (11), the pronoun “J” (I) is annotated as the subject of all the verbs (“aimer”, “aller”, “se coucher”, “déranger”).

(11)- J’aime autant aller me coucher que de déranger les gens.
I prefer going to sleep than disturbing people.

The relation tagset covers 219 relations when taking prepositions into account (94 hapaxs) and 32 labels excluding prepositions (5 hapaxs). A sample is provided in Table 3.

3.2. Syntactic Patterns

Syntactic patterns can be restored by adding up all syntactic relations for each verb occurrence. To show the kind of patterns obtained from the annotation, we chose the verb “répondre” (*to answer*). This verb occurs 165 times in the FTC for only one index (not a phrasal verb) and 22 syntactic patterns were collected. Table 4 shows non-hapax patterns which account for nearly 90% of the data.

Syntactic Patterns	Freq.	Prop.
SUJ-DirectSpeech	83	50%
SUJ-IOBJ-DirectSpeech	23	14%
SUJ-IOBJ	9	5%
SUJ	9	5%
SUJ-à	7	4%
SUJ-OBJ	6	4%
SUJ-IOBJ-ThatClause	6	4%
SUJ-MAN-DirectSpeech	4	2%
SUJ-en-DirectSpeech	3	2%
SUJ-avec-DirectSpeech	2	1%
SUJ-à-DirectSpeech	2	1%
...		

Table 4: Patterns of the verb “répondre”.

As can be seen, two syntactic relations play a predominant role in the FTC: Subject (SUJ) and Direct Speech (cf. Table 3). Example 12 illustrates both syntactic relations.// (12) -D’accord, merci beaucoup, répondit Jérémie.// *trans. -Ok, thanks a lot, Jérémie answered.*

Two other facts must be pointed:

- The high number of those “surface” patterns can be explained by the absence of a syntactic relation. For instance, the absence/presence of IOBJ (the addressee as pronoun) is responsible for the split of pattern 1 and 2. Since an answer is in most examples addressed to someone, the addressee is eluded because the reference can be inferred from context. Ellipsis does not necessarily influence verb meaning.
- Different syntactic relations may happen to split identical functions. For example the syntactic relation realized by the preposition “à” (*to*) often designates the addressee (IOBJ).

The conclusion which can be drawn from this verb (other experiments have confirmed this point) is that this syntactic scheme tends to create (unnecessary) scattering.

4. Semantic Annotation

The semantic annotation performed on the FTC answered two different needs:

- Limiting syntactic scattering by merging similar syntactic relations: this is the task of semantic role annotation and will be discussed in subsection 4.3.
- Restraining syntactic relations by introducing semantic features or categories (e.g. selectional restrictions), which is the task presented in the next subsection.

RELATION	FREQ.	EXAMPLES
SUJ	12110	Oui, {je} [sais]. <i>Yes, I know.</i>
OBJ	4668	{Fouad} s'est fait [arrêter]. <i>Fouad has been arrested.</i>
Direct Speech	1030(3)	{“Pousse-toi, triple idiot !”} [dit] Jennyfer. <i>“move over, you prize idiot !” said Jennyfer.</i>
I OBJ	855	je ne {vous} le [dirai] pas. <i>I won't tell it to you.</i>
INF	686	Les enfants, [courez] faire {sonner} le tocsin. <i>Children, run and ring the bell.</i>
That-clause	593	On [raconte] {que} les festivités durèrent trois mois. <i>It is told that the party lasted 3 months.</i>
...		

Table 3: Syntactic relations wrt frequency; verb between square brackets and argument between curly brackets.

4.1. Referential Annotation

A common claim in linguistics is that there exists a strong relationship between verb meaning (predicate) and the semantic categories or types of its arguments. For example, it could be proposed that only members of the Human category may be thinkers, or that only liquids can be drunk. This is only part of the story: however strong these relations are in the real world, their application to texts and words is problematic. Words regularly refer to different things and shift meaning. Pustejovsky has, among others, proposed to account for regular sense alternations, called regular polysemy or logical metonymy (Pustejovsky, 1998). The model, named the Generative Lexicon, helps to tackle cases such as container/content alternations (as in *drink tea/a cup*) and has recently been confronted to corpus data (Pustejovsky and Ježek, 2008). To do so, the possible set of categories that a verb argument belongs to (that “tea” is a liquid or that “cup” is a kind of container) needs to be known: identifying strong semantic relations is a first step, resolving metonymy is the next.

If a semantic category can be easily picked from nouns in the context of a sentence, pronouns, on the contrary, do not convey semantic information other than gender or number: they refer. In the FTC for example, pronouns account for more than a quarter of the subjects and of the indirect objects of the verb *dire* (*to say*). In a corpus-based framework, before identifying relevant semantic relations, anaphoric references can either be resolved or explicitly discarded. A similar reasoning could be applied to proper names: our analysis of the name “Christophe” in the FTC revealed that more than a third of its 18 occurrences referred to an animal. For the FTC, it was decided that each anaphoric reference should be annotated with a semantic category.

It is worth mentioning that reference resolution is generally approached as co-reference resolution (Orăsan et al., 2008): linking co-referential entities, that is by assigning the same identifier to various linguistic expressions (13).

(13)<entity id='1'>Harold</entity> came back. <entity id='1'>He</entity> had forgotten his hat.

This was not our aim, since it does not provide for a semantic characterization of the reference.

Pronouns, in their variety of forms (demonstrative, possessive, relative and personal) as well Proper Names were

annotated with a word corresponding to the most specific category expressed in the text (called a referential category). For instance, if, in the text, a referent is introduced and is afterwards referred to with a pronoun, the most specific linguistic description was used to annotate the pronoun. Example 13 would turn out as Example 14.

(14)<entity class='prince'>Harold</entity> came back. <entity class='prince'>He</entity> forgot his hat.

When a word refers to more than one entity, all the possible categories are included in the annotation (e.g. plural pronoun “les” in Example 15).

(15) <entity class='éléphant'>Tu</entity> es le seul qui sois assez costaud pour <entity class='garçon;fille'>les</entity> porter. *trans. You are the only one strong enough to carry them.*

Since referential categories are nouns, the main benefit of this scheme is to provide a common and comparable basis for pronoun verb arguments, Proper Name arguments and (regular) noun arguments at the semantic level. All in all, 24668 anaphoric occurrences have been labelled according to these conventions. The most frequent observed categories are shown in Table 5.

4.2. Ontological categories

Referential categories were then classed into ontological categories (Human, Animal, Imaginary, etc.) and could belong to only one of them. For example, the prince in (14) is human because princes only belong to this category. When animal princes were found, the Animal category was kept. These ontological categories do not cover all the words in the FTC but only those arguments syntactically linked with previously selected verbs (cf. section 3). A sample is provided in Table 6.

40 general ontological categories were selected (corresponding to a total of 17151 verb argument occurrences) from the Brandeis Linguistic Ontology (Hanks, 2008). They were tested as selectional restrictions by combining them with syntactic relations (El Maarouf, 2009). For example the subject relation of the verb “dire” (*to say* combines with Humans (60%), but also Animals (15%) and Imaginary creatures (13%). This behaviour was found to be common with other speech verbs and cognitive verbs (*to know, to think, etc.*), hence, pointing out the fact that:

Referential Category	Freq.	trans.
homme	3907	man
enfant	3451	child
fille	2161	girl
femme	1546	woman
garçon	734	boy
lutin	521	goblin
animal	448	animal
prince	429	prince
chat	421	cat
bûcheron	412	lumberjack
princesse	411	princess
cheval	354	horse
tigre	335	tiger
souris	282	mouse
roi	282	king
extra-terrestre	280	alien
fée	274	fairy
sorcière	270	witch
hérisson	238	hedgehog
...		

Table 5: Referential Category wrt frequency.

Ontological Category	Freq.
HUMAN	7723
ANIMAL	2926
IMAGINARY-CREATURE	1457
PLACE	960
OBJECT	901
VEGETAL	376
BODY-PART	353
FOOD	224
EVENT	221
INFORMATION	192
TIME-PERIOD	180
...	

Table 6: Ontological Category wrt frequency.

- Selectional restrictions were only partially useful in the FTC corpus.
- Ontological categories could be merged for the FTC corpus.
- Ontological categories could also entail scattering in the FTC corpus.

Whether these alternations (some of which could be interpreted as personifications) are cases of regular polysemy or genre-specific regular polysemy is open to discussion.

4.3. Semantic Role Annotation

In order to reduce syntactic and semantic scattering, work has also been initiated to identify semantic roles for the most frequent verbs. Semantic roles refers to functional categories regardless of their syntactic realization in the clause and each verb sense is associated with a small set of roles.

The FrameNet database⁵(Baker et al., 1998; Fillmore, 1982) was used to annotate roles (or frame elements) and predicates (or frames). For this task, syntactic patterns were mapped onto predicate-role tuples, in order to reveal the discrepancies between both levels. For example, the verb “appeler” (*to call*) in its non-pronominal form occurs in 17 patterns whereas it only concerns two frames, namely CONTACTING (establishing a communication) and BEING-NAMED (giving a name to an entity). Sometimes the mapping between semantic roles and syntactic relations is straightforward: in the patterns SUJ-OBJ (16) and SUJ-DirectSpeech (17), SUJ, OBJ and DirectSpeech always stands for the Communicator, the Addressee and the Communication, respectively.

(16) Si seulement je pouvais appeler cet oiseau.

trans. If only I could call this bird.

(17) “Ça y est Kléber! Je suis prêt”, appela son père.

trans. “All right Kléber! I am ready”, his father called.

However, the same syntactic relation may in some cases correspond to more than one semantic role. For example, prepositional phrases introduced by “de” may either be Depictive (18) or Source (19).

(18) “Tirelouï, Tirelouï”, appela-t-elle de toutes ses forces.

trans. “Tirelouï, Tirelouï”, she called with all her strength.

(19) Il entend Mélisa l’appeler de la cuisine.

trans. He heard Mélisa calling him from the kitchen.

As discussed earlier, different syntactic relations may perform the same function, and thus be labelled with the same semantic role. Currently, 63 Frames have been defined for 41 verbs, either directly linked to FrameNet data, or inspired from this resource when a given meaning pattern did not exist.

5. Conclusions and Perspectives

This paper describes a French corpus of fairy tales and the schemes used for syntactic and semantic annotation. Syntactic dependency annotation focused on 120 verbs and semantic role labelling on 41 verbs. The corpus is also referentially annotated: each anaphoric occurrence is linked to a referential category. Since only one annotator took part in the annotation process, some decisions may be subject to discussions. In order to provide a sounder ground for the annotations, another round of annotation could be forecast and the agreement rate tested.

The Fairy Tales Corpus and its detailed syntactic and semantic annotation has been used in work on the interface between discourse and semantics (El Maarouf, 2009). (El Maarouf et al., 2009b) have also compared it to a press corpus to analyse the distribution of semantic categories according to text genre. Perspectives on using this corpus by NLP applications include evaluating systems on syntactic relations (Dependency parsing), semantic classification (Information Extraction) and semantic frames (Semantic Role

⁵framenet.icsi.berkeley.edu

Labelling). With the growing interest on fairy tales, this corpus may become an asset for the NLP community and facilitate research on high level semantic analysis. The resources will be freely distributed by early 2012.

6. References

- A. Aarne and S. Thompson. 1973. *The types of the folktale: a classification and bibliography*. FF communications. Suomalainen Tiedeakatemia.
- C F Baker, C J Fillmore, and J B Lowe. 1998. The berkeley framenet project. In *Proceedings of the 17th international conference on Computational linguistics - Volume 1*, COLING '98, pages 86–90, Stroudsburg, PA, USA. Association for Computational Linguistics.
- C Blanche-Benveniste. 1990. *Le français parlé: études grammaticales*. Sciences du langage. Editions du Centre national de la recherche scientifique.
- M-C de Marneffe and C D Manning. 2008. The stanford typed dependencies representation. In *COLING Workshop on Cross-framework and Cross-domain Parser Evaluation*.
- I El Maarouf, M Le Tallec, and J Villaneau. 2009a. Ontologies naturelles et coercion : formalisation de connaissances à partir d'observations en corpus. In *Journées de Linguistique de Corpus*, Lorient.
- I El Maarouf, F Saïd, J Villaneau, and D Duhaut. 2009b. Comparing child and adult language: Exploring semantic constraints. In *The 2nd Workshop on Child, Computer and Interaction*, Cambridge, MA, USA, November. ACM.
- I El Maarouf. 2009. Natural ontologies at work: investigating fairy tales. In *Corpus Linguistics Conference*, Liverpool.
- C J Fillmore. 1982. Frame semantics. In *Linguistics in the Morning Calm*, pages 111–137, Séoul. The Linguistic Society of Korea.
- R Gelin, C d'Alessandro, A. Le Quoc, O Deroo, D Doukhan, J-C Martin, C Pelachaud, A Rilliard, and S Rosset. 2010. Towards a storytelling humanoid robot. In *Dialog with Robots – 2010 AAI Fall Symposium, November 11-13, 2010 Arlington, VA, USA*.
- P W Hanks. 2008. Lexical patterns: From hornby to hunston and beyond. In *Euralex*, pages 89–129, Barcelone.
- M Le Tallec, J Villaneau, J-Y Antoine, A Savary, and A Syssau. 2010. Emologus: a compositional model of emotion detection based on the propositional content of spoken utterances. In *13th international conference on Text, speech and dialogue, TSD'10*, page 361–368, Berlin/Heidelberg. Springer-Verlag.
- D Nadeau and S Sekine. 2007. A survey of named entity recognition and classification. *Lingvisticae Investigationes*, 30(1):3–26, January.
- C Orăsan, D Cristea, R Mitkov, and A Branco. 2008. Anaphora resolution exercise: An overview. In *Proceedings of the Sixth International Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco, May, 28 – 30.
- M Palmer, D Gildea, and P Kingsbury. 2005. The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106, March.
- V I A Propp. 1968. *Morphology of the folktale*. University of Texas Press.
- J Pustejovsky and E Ježek. 2008. Semantic coercion in language: Beyond distributional analysis. *Italian Journal of Linguistics*, 20(1):181–214.
- J Pustejovsky. 1998. *The generative lexicon*. MIT Press, Cambridge (Ma).
- S Saint-Aimé, B Le Pévédic, D Duhaut, and T Shibata. 2007. EmotiRob: companion robot project. In *IEEE International Symposium on Robots and Human Communications RO-MAN*, pages 919–924.
- A Scheidel and T Declerck. 2010. Apftml - augmented proppian fairy tale markup language. In Darányi S and Lendvai P, editors, *First International AMICUS Workshop on Automated Motif Discovery in Cultural Heritage and Scientific Communication Texts: Poster session*. Szeged University, 10.
- H Schmid. 1994. Probabilistic part-of-speech tagging using decision trees. In *Proceedings of the International Conference on New Methods in Language Processing*, Manchester, UK.
- J McH Sinclair. 1991. *Corpus, Concordance, Collocation*. Oxford University Press, Oxford.
- M Theune, S Faas, E Faas, A Nijholt, and D Heylen. 2003. The virtual storyteller: Story creation by intelligent agents. In *Proceedings of the Technologies for Interactive Digital Storytelling and Entertainment (TIDSE) Conference*, pages 204–215.
- E P Volkova, B J Mohler, D Meurers, D Gerdemann, and H H Bülthoff. 2010. Emotional perception of fairy tales: achieving agreement in emotion annotation of text. In *Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text, CAAGET '10*, pages 98–106, Stroudsburg, PA, USA. Association for Computational Linguistics.