

Automatic Classification of Communicative Functions of Definiteness

Archna Bhatia* Chu-Cheng Lin* Nathan Schneider* Yulia Tsvetkov*
Fatima Talib Al-Raisi* Laleh Roostapour* Jordan Bender† Abhimanu Kumar*
Lori Levin* Mandy Simons* Chris Dyer*

*Carnegie Mellon University
Pittsburgh, PA 15213

Abstract

Whether an NP is realized grammatically as definite or not depends on a variety of semantic, pragmatic, and discourse criteria, or *communicative functions*, the interaction of which varies from language to language. We present supervised classifiers for English that use lexical, morphological, and syntactic features to predict communicative functions of definiteness. The benefits of this work are twofold: linguistically, the classifiers' features and parameters model the *grammaticalization* of definiteness in English, not all of which are obvious. Computationally, it presents a framework to predict *semantic and pragmatic* communicative functions of definiteness which, unlike lexical and morphosyntactic features, are preserved in translation. The classifiers may therefore be useful for coreference resolution, MT, and other mono- and cross-lingual NLP tasks.

1 Introduction

Definiteness is a morphosyntactic category associated with some aspects of meaning (Lyons, 1999). While some morphosyntactic forms of definiteness are employed by all the languages, namely, demonstratives, personal pronouns and possessives, languages display a vast range of variation with respect to the form and meaning of definiteness. For example, while languages like English make use of definite and indefinite articles to distinguish between the discourse status of various entities (*the car* vs. *a car* vs. *cars*), many other languages—including Czech, Indonesian, and Russian—do not have articles (although they do have demonstrative determiners). Sometimes definiteness is marked with affixes or clitics, as in Arabic. Sometimes it is expressed with other constructions, as in Chinese (a language without articles), where the existential construction is used to express indefinite subjects and the *ba*- construction to express definite direct objects (Chen, 2004).

Aside from this variation in the form of (in)definite noun phrases (NPs) within and across languages, there is also variability in the semantic, pragmatic, and discourse-related functions expressed by (in)definites. We will refer to these as *communicative functions* of (in)definiteness. Croft (2003, pp. 6–7), for instance, shows that the conditions under which an NP is marked as definite or indefinite (or not at all) are language-specific by contrasting English and French translations such as:

- (1) He showed **extreme care**. (unmarked)
Il montra **un soin extrême**. (indef.)
- (2) I love **artichokes** and asparagus. (unmarked)
J'aime **les artichauts** et les asperges. (def.)
- (3) His brother became **a soldier**. (indef.)
Son frère est devenu **soldat**. (unmarked)

A cross-linguistic classification of communicative functions should be able to characterize the aspects of meaning that account for the different patterns of definiteness marking exhibited in (1–3): e.g., that (2) concerns a generic class of entities while (3) concerns a role filled by an individual. For more on communicative functions, see §2.

• NONANAPHORA $[-A, -B]$	999	• ANAPHORA $[+A]$	1574
- UNIQUE $[+U]$	287	- BASIC_ANAPHORA $[-B, +F]$	795
* UNIQUE_HEARER_OLD $[+F, -G, +S]$	251	* SAME_HEAD	556
• UNIQUE_PHYSICAL_COPRESENCE $[+R]$	13	* DIFFERENT_HEAD	329
• UNIQUE_LARGER_SITUATION $[+R]$	237	- EXTENDED_ANAPHORA $[+B]$	779
• UNIQUE_PREDICATIVE_IDENTITY $[+P]$	1	* BRIDGING_NOMINAL $[-G, +R, +S]$	43
* UNIQUE_HEARER_NEW $[-F]$	36	* BRIDGING_EVENT $[+R, +S]$	10
- NONUNIQUE $[-U]$	581	* BRIDGING_RESTRICTIVE_MODIFIER $[-G, +S]$	614
* NONUNIQUE_HEARER_OLD $[+F]$	169	614[^A _B this is going in the next row, need to make everything small so it fits in the row.]	
• NONUNIQUE_PHYSICAL_COPRESENCE		* BRIDGING_SUBTYPE_INSTANCE $[-G]$	0
$[-G, +R, +S]$	39	* BRIDGING_OTHER_CONTEXT $[+F]$	112
• NONUNIQUE_LARGER_SITUATION $[-G, +R, +S]$	117	• MISCELLANEOUS $[-R]$	732
• NONUNIQUE_PREDICATIVE_IDENTITY $[+P]$	13	- PLEONASTIC $[-B, -P]$	53
* NONUNIQUE_HEARER_NEW_SPEC $[-F, -G, +R, +S]$	231	- QUANTIFIED	248
* NONUNIQUE_NONSPEC $[-G, -S]$	181	- PREDICATIVE_EQUATIVE_ROLE $[-B, +P]$	58
- GENERIC $[+G, -R]$	131	- PART_OF_NONCOMPOSITIONAL_MWE	100
* GENERIC_KIND_LEVEL	0	- MEASURE_NONREFERENTIAL	125
* GENERIC_INDIVIDUAL_LEVEL	131	- OTHER_NONREFERENTIAL	148

	+	-	0		+	-	0		+	-	0		+	-	0
<u>Anaphoric</u>	1574	999	732	<u>Generic</u>	131	1476	1698	<u>Predicative</u>	72	53	3180	<u>Specific</u>	1305	181	1819
<u>Bridging</u>	779	1905	621	<u>Familiar</u>	1327	267	1711	<u>Referential</u>	690	863	1752	<u>Unique</u>	287	581	2437

Figure 1: CFD (Communicative Functions of Definiteness) annotation scheme, with frequencies in the corpus. Internal (non-leaf) labels are in bold; these are not annotated or predicted. +/− values are shown for ternary attributes Anaphoric, Bridging, Familiar, Generic, Predicative, Referential, Specific, and Unique; these are inherited from supercategories, but otherwise default to 0. Thus, for example, the full attribute specification for **UNIQUE_PHYSICAL_COPRESENCE** is $[-A, -B, +F, -G, 0P, +R, +S, +U]$. Counts for these attributes are shown in the table at bottom.

This paper describes classifiers that predict communicative function labels for English NPs using lexical, morphological, and syntactic features. The contribution of our work is in both the output of the classifiers and the models themselves (features and weights). Each classifier predicts communicative function labels that capture aspects of discourse-newness, uniqueness, specificity, and so forth. Such functions are usually preserved in translation, even when the grammatical mechanisms for expressing them are different. Indeed, previous work has noted that machine translation systems face problems while translating from one language to another when the languages use different grammatical strategies (see §7). The communicative function labels also represent the discourse status of entities, making them relevant for entity tracking, knowledge base construction, and information extraction.

Our log-linear model is a form-meaning mapping, consisting of the syntactic, lexical, and morphological features and weights that are predictive of communicative functions. We are using this model to generate plausible hypotheses regarding the form-meaning relationship which can then be tested rigorously through controlled experiments. This hypothesis generation is linguistically significant as it indicates new grammatical mechanisms beyond the articles *the* and *a* that are used for expressing definiteness in English.

To build our models, we leverage a cross-lingual definiteness annotation scheme (§2) and annotated English corpus (§3) from prior work (Bhatia et al., 2014). Our classifiers, §4, are supervised models with features that combine lexical and morphosyntactic information with prespecified attributes or groupings of the communicative function labels (such as Anaphoric, Bridging, Specific in fig. 1) to predict leaf node labels; the evaluation measures (§5) include one that exploits these label groupings to award partial credit according to relatedness. §6 presents experiments comparing several models and discussing their strengths and weaknesses; computational work and applications related to definiteness are addressed in §7.

2 Annotation scheme

The literature on definiteness describes functions such as uniqueness, familiarity, identifiability, anaphoricity, specificity, and referentiality (Birner and Ward, 1994; Condoravdi, 1992; Evans, 1977, 1980; Gundel et al., 1988, 1993; Heim, 1990; Kadmon, 1987, 1990; Lyons, 1999; Prince, 1992; Roberts,

2003; Russell, 1905, *inter alia*) as being related to definiteness. The reductionist approaches to definiteness try to define it in terms of one or two of above mentioned communicative functions. For example, Roberts (2003) proposes that the combination of uniqueness and a presupposition of familiarity underlie all definite descriptions. However, possessive definite descriptions (*John's daughter*) and the weak definites (*the son of Queen Juliana of the Netherlands*) are neither unique nor necessarily familiar to the listener before they are spoken. In contrast to the reductionist approaches are approaches to grammaticalization (Hopper and Traugott, 2003) in which grammar develops over time in such a way that each grammatical construction has some prototypical communicative functions, but may also have many non-prototypical communicative functions. The scheme we are adopting for this work, the annotation scheme for Communicative Functions of Definiteness (CFD), as described in Bhatia et al. (2014), assumes that there may be multiple functions to definiteness. It is based on a combination of these functions and is summarized in fig. 1. It was developed by annotating texts in two languages (English and Hindi) for four different genres, namely TED talks, the presidential inaugural speech, news articles and fictional narratives, keeping in mind the communicative functions that have been associated with definiteness in the linguistic literature.

The CFD is hierarchically organized. This hierarchical organization serves to reduce the number of decisions that an annotator needs to make for speed and consistency. At the highest level, the distinction is made between **Anaphora**, **Nonanaphora** and **Miscellaneous** functions of an NP (the annotatable unit). While the communicative functions **Anaphora** and **Nonanaphora** determine whether an entity is old or new in discourse respectively, the **Miscellaneous** function is mainly assigned to various kinds of non referential NPs. Within the **Anaphora** category, further distinctions are made based on communicative functions, such as **Basic_Anaphora** vs. **Extended_Anaphora**. **Basic_Anaphora** refers to entities that have been mentioned before. **Extended_Anaphora** refers to entities that have not been mentioned themselves but are evoked by a previously mentioned entity. For example, after mentioning a wedding, *the bride*, *the groom*, and *the cake* are considered to be **Extended_Anaphora**. Within the **Nonanaphora** category, the distinction is made between a **Unique** vs. **Nonunique** vs. **Generic** communicative function. The function **Unique** refers to entities that become unique in a context due to various reasons. For example, *Obama* can safely be considered unique in a political discourse in US contemporarily. The function **Nonunique** refers to entities that can possible have multiple referents which may or may not become identifiable in a speech situation. For example, *a little riding hood of red velvet* in fig. 2 could be annotated as a **Nonunique** entity. The function **Generic** refers to classes or types of entities rather than specific entities. For example, *Dinosaurs* in *Dinosaurs are extinct.* is a **Generic**. Another important distinction CFD makes is between the communicative functions **Hearer_Old** and **Hearer_New** which refer to entities which are familiar to the hearer (e.g. if they are physically present in the speech situation) or non familiar respectively. However this distinction cuts across the two subparts of the hierarchy, **Anaphora** and **Nonanaphora**. While we only express **Hearer_Old** through the attribute value in the **Anaphora** category, in the **Nonanaphora**, there are specific labels used corresponding to this distinction, e.g. **Unique_Hearer_Old** and **Unique_Hearer_New**.

There are multiple ways to organize the hierarchy based on which of these distinctions are made first and which ones are made subsequently. For example, Komen (2013) has also proposed a hierarchy with similar leaf nodes, but different internal structure. Since it is possible that some natural groupings of labels are not reflected in the hierarchy used, we have also associated several attributes with each label type, viz. Anaphoric, Bridging, Familiar, Generic, Predicative, Referential, Specific, and Unique. These attributes can have values of +, -, or 0, see fig. 1. For instance, with the Anaphoric attribute, a value of + applies to labels that can never mark NPs new to the discourse, - applies to labels that can *only* apply if the NP is new in the discourse, and 0 applies to labels such as **Pleonastic** (where anaphoricity is not applicable because there is no discourse referent). For further details on the scheme, see Bhatia et al. (2014).

Once upon a time there was a dear little girl who was loved by everyone who looked at her, but most of all by her grandmother, and there was nothing that she would not have given to the child.

Once she gave her a little riding hood of red velvet, which suited her so well that
SAME_HEAD DIFFERENT_HEAD OTHER_NONREFERENTIAL SAME_HEAD
NONUNIQ_HEARER_NEW_SPEC
she would never wear anything else; so she was always called 'Little Red Riding Hood.'
SAME_HEAD QUANTIFIED SAME_HEAD UNIQ_HEARER_NEW

Figure 2: An annotated sentence from “Little Red Riding Hood.” The previous sentence is shown for context.

3 Data

We use the English definiteness corpus of Bhatia et al. (2014), which consists of texts from multiple genres annotated with the scheme described in §2.¹ The 17 documents consist of prepared speeches (TED talks and a presidential address), published news articles, and fictional narratives. The TED data predominates (75% of the corpus)²; the presidential speech represents about 16%, fictional narratives 5%, and news articles 4%. All told, the corpus contains 13,860 words (868 sentences), with 3,422 NPs (the annotatable units). Bhatia et al. (2014) report high inter-annotator agreement, estimating Cohen’s $\kappa = 0.89$ within the TED genre as well as for all genres.

Figure 2 is an excerpt from the “Little Red Riding Hood” annotated with the CFD scheme.

4 Classification framework

To model the relationship between the grammar of definiteness and its communicative functions in a data-driven fashion, we work within the supervised framework of feature-rich discriminative classification, treating the functional categories from §2 as output labels y and various lexical, morphological, and syntactic characteristics of the language as features of the input x . Specifically, we learn two kinds of probabilistic models. The first is a log-linear model similar to multiclass logistic regression, but deviating in that logistic regression treats each output label (response) as atomic, whereas we decompose each into *attributes* based on their linguistic definitions, enabling commonalities between related labels to be recognized. Each weight in the model corresponds to a feature that mediates between *percepts* (characteristics of the input NP) and attributes (characteristics of the label). This is aimed at attaining better predictive accuracy as well as feature weights that better describe the form–function interactions we are interested in recovering. We also train a random forest model, which sacrifices interpretability of the learned parameters for predictive accuracy.

Our setup is formalized below, where we discuss the mathematical models and linguistically motivated features.

4.1 Models

We experiment with two classification methods: a log-linear model and a nonlinear tree-based ensemble model. Due to their consistency and interpretability, linear models are a valuable tool for quantifying and analyzing the effects of individual features. Non-linear models, while less interpretable, often outperform logistic regression (Perlich et al., 2003), and thus could be desirable when the predictions are needed for a downstream task.

4.1.1 Log-linear model

At test time, we model the probability of communicative function label y conditional on an NP x as follows:

$$p_{\theta}(y|x) = \log \frac{\exp \theta^{\top} \mathbf{f}(x, y)}{\sum_{y' \in \mathcal{Y}} \exp \theta^{\top} \mathbf{f}(x, y')} \quad (1)$$

¹The data can be obtained from http://www.cs.cmu.edu/~ytsvetko/definiteness_corpus.

²The TED talks are from a large parallel corpus obtained from <http://www.ted.com/talks/>.

where $\boldsymbol{\theta} \in \mathbb{R}^d$ is a vector of parameters (feature weights), and $\mathbf{f}: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ is the feature function over input–label pairs. The feature function is defined as follows:

$$\mathbf{f}(x, y) = \boldsymbol{\phi}(x) \times \tilde{\boldsymbol{\omega}}(y) \quad (2)$$

where the percept function $\boldsymbol{\phi}: \mathcal{X} \rightarrow \mathbb{R}^c$ produces a vector of real-valued characteristics of the input, and the attribute function $\tilde{\boldsymbol{\omega}}: \mathcal{Y} \rightarrow \{0, 1\}^a$ encodes characteristics of each label. There is a feature for every percept–attribute pairing: so $d = c \cdot a$ and $f_{(i-1)a+j}(x, y) = \phi_i(x) \tilde{\omega}_j(y)$, $1 \leq i \leq c$, $1 \leq j \leq a$. The contents of the percept and attribute functions are detailed in §4.2 and §4.3 respectively.

For prediction, having learned weights $\hat{\boldsymbol{\theta}}$ we use the Bayes-optimal decision rule for minimizing misclassification error, selecting the y that maximizes this probability:

$$\hat{y} \leftarrow \arg \max_{y \in \mathcal{Y}} p_{\hat{\boldsymbol{\theta}}}(y|x) \quad (3)$$

Training optimizes $\hat{\boldsymbol{\theta}}$ so as to maximize a convex L_2 -regularized³ learning objective over the training data \mathcal{D} :

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} -\lambda \|\boldsymbol{\theta}\|_2^2 + \sum_{\langle x, y \rangle \in \mathcal{D}} \log \frac{\exp \boldsymbol{\theta}^\top \mathbf{f}(x, y)}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta}^\top \mathbf{f}(x, y'))} \quad (4)$$

With $\tilde{\boldsymbol{\omega}}(y) = \text{the identity of the label}$, this reduces to standard logistic regression.

4.1.2 Non-linear model

We employ a random forest classifier (Breiman, 2001), an ensemble of decision tree classifiers learned from many independent subsamples of the training data. Given an input, each tree classifier assigns a probability to each label; those probabilities are averaged to compute the probability distribution across the ensemble.

An important property of the random forests, in addition to being an effective tool in prediction, is their immunity to overfitting: as the number of trees increases, they produce a limiting value of the generalization error.⁴ Thus, no hyperparameter tuning is required. Random forests are known to be robust to sparse data and to label imbalance (?), both of which are challenges with the definiteness dataset.

4.2 Percepts

The characteristics of the input that are incorporated in the model, which we call *percepts* to distinguish them from model features linking inputs to outputs,⁵ are intended to capture the aspects of English morphosyntax that may be relevant to the communicative functions of definiteness. [B^A should we explain the difference between percepts and features - may be when the 1st time we mention them in section 4 intro above?- Nathan??? there we say features are lexical, morphological, syntactic characteristics of the input x, here we say these are percepts, and features link input with output.]

After preprocessing the text with a dependency parser and coreference resolver, which is described in §6.1, we extract several kinds of percepts for each NP.

4.2.1 Basic

Words of interest. These are the *head* within the NP, all of its *dependents*, and its *governor* (external to the NP). We are also interested in the *attached verb*, which is the first verb one encounters when traversing the dependency path upward from the head. For each of these words, we have separate percepts capturing: the token, the part-of-speech (POS) tag, the lemma, the dependency relation, and (for the head only) a binary indicator of plurality (determined from the POS tag). As there may be multiple dependents, we have additional features specific to the first and the last one. Moreover, to better capture tense, aspect and modality, we collect the attached verb’s *auxiliaries*. We also make note of *neg* if it is attached to the verb.

³ As is standard practice with these models, bias parameters (which capture the overall frequency of percepts/attributes) are excluded from regularization.

⁴ See Theorem 1.2 in Breiman (2001) for details.

⁵ See above.

Structural. The structural percepts are: the *path length* from the head up to the root, and to the attached verb. We also have percepts for the number of dependents, and the number of dependency relations that link non-neighbors. Integer values were binarized with thresholding.

Positional. These percepts are the *token length* of the NP, the NP’s *location* in the sentence (first or second half), and the *attached verb’s position* relative to the head (left or right). 12 additional percept templates record the POS and lemma of the left and right neighbors of the head, governor, and attached verb.

4.2.2 Contextual NPs

When extracting features for a given NP (call it the “target”), we also consider NPs in the following relationship with the target NP: its *immediate parent*, which is the smallest NP whose span fully subsumes that of the target; the *immediate child*, which is the largest NP subsumed within the target; the *immediate precedent* and *immediate successor* within the sentence; and the *nearest preceding coreferent mention*.

For each of these related NPs, we include all of their basic percepts conjoined with the nature of the relation to the target.

4.3 Attributes

As noted above, though labels are organized into a tree hierarchy, there are actually several dimensions of commonality that suggest different groupings. These attributes are encoded as ternary characteristics; for each label (including internal labels), every one of the 8 attributes is assigned a value of +, −, or 0 (refer to fig. 1). In light of sparse data, we design features to exploit these similarities via the attribute vector function

$$\omega(y) = [y, A(y), B(y), G(y), O(y), P(y), R(y), S(y), U(y)]^T \quad (5)$$

where $A: \mathcal{Y} \rightarrow \{+, -, 0\}$ returns the value for Anaphoric, $B(y)$ for Bridging, etc. The identity of the label is also included in the vector so that different labels are always recognized as different by the attribute function. The categorical components of this vector are then binarized to form $\tilde{\omega}(y)$; however, instead of a binary component that fires for the 0 value of each ternary attribute, there is a component that fires for *any* value of the attribute—a sort of bias term. The weights assigned to features incorporating + or − attribute values, then, are easily interpreted as deviations relative to the bias.

5 Evaluation

The following measures will be used to evaluate our predictor against the gold standard for the held-out evaluation (dev or test) set \mathcal{E} :

- **Exact match:** This accuracy measure gives credit only where the predicted and gold labels are identical.
- **By leaf label:** We also compute precision and recall of each leaf label to determine which categories are reliably predicted.
- **Soft match:** This accuracy measure gives partial credit where the predicted and gold labels are related. It is computed as the proportion of attributes whose (categorical) values match: $|\omega(y) \cap \omega(y')|/9$.

[^A_B TODO: We need to enter these numbers. Chu-Cheng?]

6 Experiments

6.1 Experimental Setup

The annotated corpus of Bhatia et al. (2014) (§3) contains 17 documents in 3 genres: 13 prepared speeches (mostly TED talks), 2 newspaper articles, and 2 fictional narratives. We arbitrarily choose some documents to hold out from each genre; the resulting test set consists of 2 TED talks (“Alisa_News”, “RobertHammond_park”), 1 newspaper article (“crime1_iPad_E”), and 1 narrative (“Little Red Riding Hood”). The test set then contains 19,28 tokens (111 sentences), in which there are 511 annotated NPs; while the training set contains 2,911 NPs among 11,932 tokens (757 sentences). Gold NP boundaries are assumed throughout our experiments.

Condition	$ \theta $	λ	Exact Match Acc.	Soft Match Acc.
Majority baseline	—	—	12.1	25.5
Log-linear classifier, attributes only	473,064	100	38.7	52.0
Log-linear classifier, labels only	413,931	100	40.8	52.1
Full log-linear classifier (labels + attributes)	926,417	100	43.7	55.6
Random forest classifier	20,363	—	49.7	60.1

Table 1: Classifiers and baseline, as measured on the test set. The first two columns give the number of parameters and the tuned regularization hyperparameter, respectively; the third and fourth columns give accuracies as percentages; the best in each column is bolded. [A Chu-Cheng will add precision and recall for the exact match. We can’t do that for soft-match. Also Chu-Cheng will add another table where we have accuracies corresponding to each leaf label.]

The log-linear model variants are trained with an in-house implementation of supervised learning with L_2 -regularized AdaGrad (Duchi et al., 2011). Hyperparameters are tuned on a development set formed by holding out every tenth instance from the training set (test set experiments use the full training set): the power of 10 giving the highest soft match accuracy was chosen for λ .⁶ The Python `scikit-learn` toolkit (Pedregosa et al., 2011) was used for the random forest classifier.⁷ Automatic dependency parses and coreference information were obtained with the parser and coreference resolution system in Stanford CoreNLP v. 3.3.0 (Socher et al., 2013; Recasens et al., 2013) for use in features (§4.2). From the parser output, the basic dependencies were used to extract the percepts used in this work. On two of the reviewers’ suggestions, to evaluate the performance of Stanford system on our data, we manually inspected the dependencies and the coreference for a subset of sentences from our corpus (using texts from TED talks and fictional narratives genres) processed using Stanford CoreNLP v. 3.3.0 and recorded the errors. We found that the Stanford parser had an accuracy of about 70% on this subset of the corpus and coreference system had an accuracy of 75% [A TODO: enter correct numbers after we are completely done with our Gold Standard- Archana,Chu-Cheng,Jordan].

6.2 Results

Measurements of overall classification performance appear in table 1. While far from perfect, our classifiers achieve promising accuracy levels given the small size of the training data and the number of labels in the annotation scheme. The random forest classifier is the most accurate, likely due to the robustness of that technique under conditions where the data are small and the frequencies of individual labels are imbalanced.

Among the log-linear models, the most successful is the richest model, which combines the fine-grained communicative function labels with higher-level attributes of those labels. This is encouraging because it suggests that the model has correctly exploited known linguistic generalizations to account for the grammaticalization of definiteness in English.

An advantage of log-linear models is that feature weights offer insights into the model’s behavior. Figure 3 lists the 10 features that received the highest positive weights in the full model for the + and – values of the Specific attribute. Reassuringly, the definite article, possessives (PRP\$), proper nouns (NNP), and the second person pronoun are associated with specific NPs, while the indefinite article is associated with nonspecific NPs. The model also seems to have picked up on the less obvious but well-attested tendency of objects to be nonspecific (Aissen, 2003).⁸

In addition to confirming known grammaticalization patterns of definiteness, we can mine the highly-weighted features for new hypotheses: here the model thinks that objects of “from” are especially likely to be specific, and that NPs with comparative adjectives (JJR) are especially likely to be nonspecific. Whether these are general trends, or just an artifact of the sentences that happened to be in the training data, will require further investigation, ideally with additional datasets and more rigorous hypothesis testing.

Finally, we can remove features to test their impact on predictive performance. Notably, in experiments

⁶Preliminary experiments with cross-validation on the training data showed that the value of λ was stable across folds.

⁷Because it is a randomized algorithm, the results may vary slightly between runs; however, a cross-validation experiment on the training data found very little variance in accuracy.

⁸The percept VB_pos_outerhead_left fires when the NP is governed by a verb to its left.

+Specific	-Specific
PRP\$_first_pos_dependents	a_first_lemma_dependents
PRP\$_pos_head_left	a_last_lemma_dependents
you_last_lemma_dependents	a_lemma_dependents_1
you_lemma_dependents_1	JJR_pos_head_left
PRP\$_pos_dependents_1	JJR_last_pos_dependents
PRP\$_pos_outerhead_right	a_lemma_dependents_2
NNP_last_pos_dependents	new_first_lemma_dependents
PRP\$_last_pos_dependents	new_last_lemma_dependents
the_first_lemma_dependents	JJR_pos_dependents_2
from_lemma_outerhead	VB_pos_outerhead_left

Figure 3: Percepts receiving highest positive weights in association with attribute values. [^A maybe we should add pointers to the not duplicated features, although am not sure if that is necessary.]

ablating features indicating articles—the most obvious exponents of definiteness in English—we see a decrease in performance, but not a drastic one. This suggests that the expression of communicative functions of definiteness is in fact much richer than morphological definiteness.

Errors. Several labels are unattested or virtually unattested in the training data, so the models unsurprisingly fail to predict them correctly at test time. SAME_HEAD and DIFFERENT_HEAD, though both common, are confused quite frequently. Whether the previous coreferent mention has the same or different head is a simple distinction for humans; low model accuracy is likely due to errors propagated from coreference resolution. This problem is so frequent that merging these two categories and retraining the random forest model improves exact match accuracy by 8% absolute and soft match accuracy by 5% absolute. Another common confusion is between the highly frequent category UNIQUE_LARGER_SITUATION and the rarer category UNIQUE_HEARER_NEW; the latter is supposed to occur only for the first occurrence of a proper name referring to an entity that is not already part of the knowledge of the larger community. In other words, this distinction requires world knowledge about well-known entities, which could perhaps be mined from the Web or other sources.

7 Related Work

Because semantic/pragmatic analysis of referring expressions is important for many NLP tasks, a computational model of the communicative functions of definiteness has the potential to leverage diverse lexical and grammatical cues to facilitate deeper inferences about the meaning of linguistic input. We have used a coreference resolution system to extract features for modeling definiteness, but an alternative would be to predict definiteness functions as input to (or jointly with) the coreference task. Applications such as information extraction and dialogue processing could be expected to benefit not only from coreference information, but also from some of the semantic distinctions made in our framework, including specificity and genericity.

Better computational processing of definiteness in different languages stands to help machine translation systems. It has been noted that machine translation systems face problems when the source and the target language use different grammatical strategies to express the same information (Stymne, 2009; Tsvetkov et al., 2013). Previous work on machine translation has attempted to deal with this in terms of either (a) preprocessing the source language to make it look more like the target language (Collins et al., 2005; Habash, 2007; Nießen and Ney, 2000; Stymne, 2009, *inter alia*); or (b) post-processing the machine translation output to match the target language, (e.g., Popović et al., 2006). Attempts have also been made to use syntax on the source and/or the target sides to capture the syntactic differences between languages (Liu et al., 2006; Yamada and Knight, 2002; Zhang et al., 2007). Automated prediction of (in)definite articles has been found beneficial in a variety of applications, including postediting of MT output (Knight and Chander, 1994), text generation (Elhadad, 1993; Minnen et al., 2000), and identification and correction of ESL errors (Han et al., 2006; De Felice and Pulman, 2008; Gamon et al., 2008; Rozovskaya and Roth, 2010). More recently, Tsvetkov et al. (2013) trained a classifier to predict where English articles might plausibly be added or removed in a phrase, and used this classifier to improve the

quality of statistical machine translation.

While definiteness morpheme prediction has been thoroughly studied in computational linguistics, studies on additional, more complex aspects of definiteness are limited. Reiter and Frank (2010) exploit linguistically-motivated features in a supervised approach to distinguish between generic and specific NPs. To the best of our knowledge, no studies have been conducted on automatic prediction of semantic and pragmatic communicative functions of definiteness more broadly.

Our work is related to research in linguistics on the modeling of syntactic constructions such as dative shift and the expression of possession with “of” or “’s”. Bresnan and Ford (2010) used logistic regression with semantic features to predict syntactic constructions. Although we are doing the opposite (using syntactic features to predict semantic categories), we share the assumption that reductionist approaches (as mentioned earlier) are not able to capture all the nuances of a linguistic phenomenon. Following Hopper and Traugott (2003) we observe that grammaticalization is accompanied by *function drift*, resulting in multiple communicative functions for each grammatical construction.

8 Conclusion

We have presented a data-driven approach to modeling the relationship between communicative functions associated with (in)definiteness and their lexical/grammatical realization in a particular language. Our feature-rich classifiers can give insight into this relationship as well as predict communicative functions for the benefit of NLP systems. This work has focused on English, but in future work we will build similar models for other languages—including languages without articles, under the hypothesis that such languages will rely on other, subtler devices to encode many of the functions of definiteness.

References

- Judith Aissen. 2003. Differential object marking: iconicity vs. economy. *Natural Language & Linguistic Theory*, 21(3):435–483.
- Archana Bhatia, Mandy Simons, Lori Levin, Yulia Tsvetkov, Chris Dyer, and Jordan Bender. 2014. A unified annotation scheme for the semantic/pragmatic components of definiteness. In *Proc. of LREC*. Reykjavík, Iceland.
- Betty Birner and Gregory Ward. 1994. Uniqueness, familiarity and the definite article in English. In *Proc. of the Twentieth Annual Meeting of the Berkeley Linguistics Society*, pages 93–102.
- Leo Breiman. 2001. Random forests. *Machine Learning*, 45(1):5–32.
- Joan Bresnan and Marilyn Ford. 2010. Predicting syntax: Processing dative constructions in American and Australian varieties of English. *Language*, 86(1):168–213.
- Ping Chen. 2004. Identifiability and definiteness in Chinese. *Linguistics*, 42:1129–1184.
- Herbert H. Clark. 1977. Bridging. In P.N. Johnson-Laird and P.C. Wason, editors, *Thinking: Readings in Cognitive Science*. Cambridge University Press.
- Michael Collins, Philipp Koehn, and Iлона Kucerova. 2005. Clause restructuring for statistical machine translation. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 531–540. Association for Computational Linguistics, Ann Arbor, Michigan.
- Cleo Condoravdi. 1992. Strong and weak novelty and familiarity. In *Proc. of SALT II*, pages 17–37.
- William Croft. 2003. *Typology and Universals*. Cambridge University Press.
- Rachele De Felice and Stephen G Pulman. 2008. A classifier-based approach to preposition and determiner error correction in L2 English. In *Proc. of the 22nd International Conference on Computational Linguistics*, pages 169–176. Association for Computational Linguistics.
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159.

- Michael Elhadad. 1993. Generating argumentative judgment determiners. In *Proc. of AAAI*, pages 344–349.
- Gareth Evans. 1977. Pronouns, quantifiers and relative clauses. *Canadian Journal of Philosophy*, 7(3):46.
- Gareth Evans. 1980. Pronouns. *Linguistic Inquiry*, 11.
- Michael Gamon, Jianfeng Gao, Chris Brockett, Alexandre Klementiev, William B Dolan, Dmitriy Belenko, and Lucy Vanderwende. 2008. Using contextual speller techniques and language modeling for ESL error correction. In *Proc. of IJCNLP*, volume 8, pages 449–456.
- Jeanette K. Gundel, Nancy Hedberg, and Ron Zacharski. 1988. The generation and interpretation of demonstrative expressions. In *Proc. of XIIth International Conference on Computational Linguistics*, pages 216–221.
- Jeanette K. Gundel, Nancy Hedberg, and Ron Zacharski. 1993. Cognitive status and the form of referring expressions in discourse. *Language*, 69:274–307.
- Nizar Habash. 2007. Syntactic preprocessing for statistical machine translation. In *MT Summit XI*, pages 215–222. Copenhagen.
- Na-Rae Han, Martin Chodorow, and Claudia Leacock. 2006. Detecting errors in English article usage by non-native speakers.
- Irene Heim. 1990. E-type pronouns and donkey anaphora. *Linguistics and Philosophy*, 13:137–177.
- Paul J. Hopper and Elizabeth Closs Traugott. 2003. *Grammaticalization*. Cambridge University Press.
- Nirit Kadmon. 1987. *On unique and non-unique reference and asymmetric quantification*. Ph.D. thesis, University of Massachusetts.
- Nirit Kadmon. 1990. Uniqueness. *Linguistics and Philosophy*, 13:273–324.
- Kevin Knight and Ishwar Chander. 1994. Automated postediting of documents. In *Proc. of the National Conference on Artificial Intelligence*, pages 779–779. Seattle, WA.
- Erwin Ronald Komen. 2013. *Finding focus: a study of the historical development of focus in English*. LOT, Utrecht.
- Yang Liu, Qun Liu, and Shouxun Lin. 2006. Tree-to-string alignment template for statistical machine translation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 609–616. Association for Computational Linguistics, Sydney, Australia.
- Christopher Lyons. 1999. *Definiteness*. Cambridge University Press.
- Guido Minnen, Francis Bond, and Ann Copestake. 2000. Memory-based learning for article generation. In *Proc. of the 2nd Workshop on Learning Language in Logic and the 4th Conference on Computational Natural Language Learning*, pages 43–48. Association for Computational Linguistics.
- Sonja Nießen and Hermann Ney. 2000. Improving SMT quality with morpho-syntactic analysis. In *Proceedings of the 18th Conference on Computational Linguistics*, pages 1081–1085. Association for Computational Linguistics.
- Fabian Pedregosa, Gael Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, M. Perrot, and Edouard Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Claudia Perlich, Foster Provost, and Jeffrey S. Simonoff. 2003. Tree induction vs. logistic regression: a learning-curve analysis. *Journal of Machine Learning Research*, 4:211–255.
- Massimo Poesio and Renata Vieira. 1998. A corpus based investigation of definite description use. *Computational*

Linguistics, 24:183–216.

- Maja Popović, Daniel Stein, and Hermann Ney. 2006. Statistical machine translation of German compound words. In *Advances in Natural Language Processing*, pages 616–624. Springer.
- Ellen F. Prince. 1992. The ZPG letter: Subjects, definiteness and information status. In S. Thompson and W. Mann, editors, *Discourse description: diverse analyses of a fund raising text*, pages 295–325. John Benjamins.
- Marta Recasens, Marie-Catherine de Marneffe, and Christopher Potts. 2013. The life and death of discourse entities: identifying singleton mentions. In *Proc. of NAACL-HLT*, pages 627–633. Atlanta, Georgia, USA.
- Nils Reiter and Anette Frank. 2010. Identifying generic noun phrases. In *Proc. of ACL*, pages 40–49. Uppsala, Sweden.
- Craig Roberts. 2003. Uniqueness in definite noun phrases. *Linguistics and Philosophy*, 26:287–350.
- Alla Rozovskaya and Dan Roth. 2010. Training paradigms for correcting errors in grammar and usage. In *Proc. of NAACL-HLT*, pages 154–162. Association for Computational Linguistics.
- Bertrand Russell. 1905. On denoting. *Mind, New Series*, 14:479–493.
- Richard Socher, John Bauer, Christopher D. Manning, and Andrew Y. Ng. 2013. Parsing with compositional vector grammars. In *Proc. of ACL*, pages 455–465. Sofia, Bulgaria.
- Sara Stymne. 2009. Definite noun phrases in statistical machine translation into Danish. In *Proc. of Workshop on Extracting and Using Constructions in NLP*, pages 4–9.
- Yulia Tsvetkov, Chris Dyer, Lori Levi, and Archana Bhatia. 2013. Generating English determiners in phrase-based translation with synthetic translation options. In *Proc. of WMT*. Association for Computational Linguistics.
- Kenji Yamada and Kevin Knight. 2002. A decoder for syntax-based statistical MT. In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 303–310. Association for Computational Linguistics, Philadelphia, Pennsylvania, USA.
- Yuqi Zhang, Richard Zens, and Hermann Ney. 2007. Improved chunk-level reordering for statistical machine translation. In *IWSLT 2007: International Workshop on Spoken Language Translation*, pages 21–28.