# Switch reference and its role in referential choice in Mbyá Guaraní narratives

Guillaume Thomas        Gregory Antono        Laurestine Bradford

Angelika Kiss        Darragh Winkelman

November 23, 2020

**Abstract:**    Switch reference has been analyzed as a reference tracking mechanism, whose main function is to avoid ambiguity of reference. One domain where this function has been argued to manifest itself is referential choice. Kibrik (2011) notably proposed that switch reference marking plays the role of a referential aid, which helps to prevent referential conflict, thereby enabling the production of reduced referential expressions like pronouns and zeroes. The present study probes this theory through an analysis of the role of switch reference marking in multifactorial models of referential choice in Mbyá Guaraní. We show that while switch reference increases the likelihood of mention reduction in Mbyá Guaraní, this effect is marginal relative to other predictors of referential choice. We argue that this result is compatible with the analysis of switch reference as a referential aid, but also supports analyses that emphasize the multiplicity of its functions, beyond the disambiguation of reference.

**Keywords:**    Switch reference, referential choice, Mbyá Guaraní, recursive partitioning.

# 1 Introduction

Switch reference (SR) is a family of grammatical devices that relate two clauses and indicate whether specific arguments in each of them have the same referent. Since Jacobsen's (1967) seminal study, it has been established that SR is a widespread phenomenon of central interest to linguistic typology and grammatical theory, as evidenced by recent overviews of this topic (McKenzie, 2015; van Gijn, 2016; Roberts, 2017; Baker and Camargo Souza, 2019). An outstanding question in this literature is whether SR can be said to have a core function. An influential proposal in this respect, originally formulated by Haiman and Munro (1983), is that SR is a reference tracking mechanism, whose main function is to avoid ambiguity of reference. The present study probes this theory through a quantitative investigation of the role of SR marking in referential choice in Mbyá Guaraní, a Tupí-Guaraní language spoken by approximately 30,000 speakers in Argentina, Brazil and Paraguay (Maria Inês Ladeira, 2018).

Referential choice is a process that speakers go through in the production of referential expressions. Having chosen to mention (i.e., make reference to) a certain entity in an utterance, the speaker must select the form of the expression that she will use to refer to it, such as a pronoun or a proper name. Referential choice is the process of selecting this linguistic form. The most fundamental split in referential choice is widely recognized to be the distinction between reduced referential devices, which include pronouns and zeroes (null arguments), and lexical referential devices, which include descriptions and proper names (Givón, 1983, 2017; Ariel, 1990; Arnold, 1998; Gundel et al., 1993; Kibrik, 2011). The choice between these two types of referential devices has been argued to be governed by the degree of accessibility or salience of the representation of referents in discourse (Ariel, 1990; Arnold, 1998). Referents that are more accessible tend to be realized by reduced devices, and less accessible referents tend to be realized by lexical devices.

While accessibility facilitates the use of reduced referential devices, the decision to refer to an entity with a reduced device may be hindered by the presence of other highly salient entities that the addressee may misinterpret as the intended referents of the reduced device. Kibrik (2011) dubs this phenomenon *referential conflict*. In this perspective, Kibrik analyzes SR marking as a *referential*

*aid*, which helps to preclude referential conflict by discriminating between two or more competing referents.

Against this backdrop, we ask whether there is evidence that Mbyá Guaraní speakers use SR marking as a disambiguation device in the course of referential choice. We expect this function to manifest itself as an interaction between SR marking and the number of referents that can serve as potential antecedents for a referential expression. More specifically, we expect that while the presence of multiple activated referents would decrease the likelihood of using a reduced referential device, this effect would be neutralized or lessened in the presence of SR marking. We test this hypothesis through multifactoral modelling of referential choice in a corpus of Mbyá Guaraní narratives authored by speakers from the state of Paraná, Brazil.

Our article is structured as follows. In section 2, we give an overview of Switch Reference and Referential Choice, and lay out the theoretical background for our study. In section 3, we give an overview of aspects of Mbyá Guaraní grammar that are relevant to our research question. Section 4 describes the data and variables used in the study. In section 5, referential choice in Mbyá Guaraní is analyzed using two recursive partitioning models. The results and implications of these models are discussed in section 6, which concludes the paper.

## 2   Referential choice and switch reference

### 2.1   Theories of switch reference

The term "switch reference" was coined by Jacobsen (1967) in his seminal study of the phenomenon in three Hokan-Coahuiltecan languages. The most cited definition of the phenomenon, however, is due to Haiman and Munro (1983):

> "Canonical switch reference is an inflectional category of the verb, which indicates whether or not its subject is identical with the subject of some other verb." (Haiman and Munro, 1983: ix)

Canonical SR is well illustrated by the following examples from Usan:

(1)  a.  *ye  nam  su-ab  isomei*

      I   tree  cut-SS  I.went.down

      'I cut the tree and went down.'                   (Haiman and Munro, 1983: ix)

   b.  *ye  nam  su-ine  isorei*

      I   tree  cut-DS  it.went.down

      'I cut the tree down.'                           (Haiman and Munro, 1983: ix)

In both examples, the SR construction relates two clauses. Haiman and Munro (1983) refer to the clause that contains the SR marker as the *marking clause*. SR marking indicates a dependency between an argument in the marking clause and some argument in another clause, which they call the *reference clause*. Following van Gijn (2016), we call the two arguments related by SR marking the *relative* and *controller* participants, respectively. In (1a) and (1b), the relative participant is the subject of the marking clause, namely the first person pronoun *ye*. The controller participant is the subject of the reference clause, which is first person in (1a) and third person in (1b). The relative and controller participants are also known as the *pivots* of the SR construction (Stirling, 1993).

In a canonical SR system, the relative and controller participants are the grammatical subjects of the marking and reference clauses, and SR marking indicates whether these arguments are coreferential. In this canonical form, SR has been analyzed as a reference tracking mechanism, whose function is to avoid ambiguity of reference. In support of this view, Haiman and Munro (1983) note that although there exist systems in which SR marking is restricted to third person pivots, there are no systems in which it is restricted to first and/or second person pivots. Systems in which SR marking is attested with all persons are viewed as grammaticalizations of the core reference tracking function.

While influential, Haiman and Munro's (1983) view of SR was criticized as overly restrictive. In particular Stirling (1993) objects to the characterization of SR as a reference tracking mechanism.

Building on earlier work by Roberts (1988), Stirling establishes that in some languages SR marking may signal whether the events described by the marking and reference clauses share the same location, time or purpose, regardless of the identity of their participants. Based on this observation, Stirling (1993) proposes that the general function of SR is to mark (dis)agreement between the situations described by the marking and reference clauses. The identity of their most prominent participant is only one dimension along which these situations can be compared, other dimensions including spatial and temporal locations. Uses of SR marking that do not exclusively track the pivots' identity have become known as non-canonical SR.

Given this tension in the literature, one may wonder to what extent speakers use SR marking in order to disambiguate referential expressions, and how important this function is to the organization of canonical SR systems. One way to address these questions would be to investigate linguistic processes that manipulate referential expressions, and whose outcome is expected to be affected by disambiguation mechanisms. If SR is used to disambiguate reference in such a process, we may then be able to determine the size and scope of its effect, which would in turn help us to assess the centrality of reference disambiguation among the different functions of SR marking. One domain in which one may hope to find such evidence is referential choice, to which we now turn.

## 2.2 Referential choice and referential aids

The fundamental distinction in the study of referential choice is that between expressions whose form encodes non-grammatical information about their referents, and expressions whose form only encodes grammatical information such as person, number and gender, if they encode any information at all about their referents. Following Kibrik (2011), we call these two categories lexical referential devices and reduced referential devices, respectively. Experimental and corpus based studies of referential choice have identified several linguistic variables that tend to be good predictors of the choice between reduced and referential devices across languages. In particular, the importance of the following variables has been widely recognized in the literature (cf. overviews in Almor and Nair, 2007; Arnold, 2010; Kibrik, 2011; Gatt et al., 2014):

(2)   Predictors of Referential Choice

  a.  *Anaphoric Distance* Among mentions of discourse old referents, those whose antecedent
      is more recent in discourse are more likely to be realized as reduced devices (Givón, 1983;
      Ariel, 1990; Arnold, 1998).

  b.  *Grammatical Function* Several studies have shown that referential mentions whose an-
      tecedent is a subject are more likely to be realized by reduced devices than those whose
      antecedent is an object (Brennan, 1995; Arnold, 1998).[1] In addition, referential mentions
      that have the same grammatical function as their antecedent tend to be more reduced
      (Arnold, 1998, 2003).

  c.  *Animacy* Reduced referential devices are more commonly produced for animate than for
      inanimate referents (Fukumara and van Gompel, 2011).

  d.  *Referential Conflict* Referential mentions are less likely to be realized by reduced devices
      when the discourse context includes other referents (Francik, 1985; Arnold and Griffin,
      2007; Fukumara and van Gompel, 2011).

These variables will be at the core of our model of referential choice in Mbyá Guaraní. Non-
linguistic variables have also been shown to have an effect on referential choice, in particular visual
salience (Fukumara et al., 2010) and different types of processing load (for a review of the relevant
literature in this latter domain, see Arnold, 2010: §3). Since our study is based on written narratives,
we will only consider linguistic variables.

  In the present study, our interest does not lie exclusively in referential choice, but rather in the

---

[1]An anonymous reviewer observed that Du Bois (1987) argued that argument realization is governed by ergative
alignment. More recent work by Haig and Schnell (2016) argues that Du Bois' conclusions, which were based on
a study of Sakapultek Mayan, don't generalize well crosslinguistically. Guaraní Mbyá is a split-S languages, and
our study takes into account split-S ('active-inactive') alignment as well as grammatical function (subject, object and
other), as predictors of referential choice, see section 4.3. The question whether ergative alignment has a role to play
in referential choice in Mbyá Guaraní is left to future studies.

role that SR marking may play in this process. Our starting point is a set of hypotheses on the function of SR that were articulated by Kibrik (2011). Kibrik proposes that the form of referential expressions is primarily a function of the degree of activation of the representation of their referents in the speaker's working memory. This degree of activation is itself determined by certain properties of the referent, such as animacy, and of the discourse in which its mention is embedded, such as the antecedent's grammatical function and distance to the mention. If a referent's degree of activation passes a certain threshold, the speaker will tend to produce a reduced device.

Crucially, the production of a reduced device for a referent with a high activation score may still be blocked by additional factors that Kibrik calls *filters*. The most important of these is *referential conflict*, which occurs when more than one referent is activated, and the production of a reduced device may cause the listener to associate it with a wrong referent.

Referential conflict may be precluded by linguistic resources that help to discriminate between several activated referents, and which Kibrik calls *referential aids*.[2] Canonical SR marking is argued to be such a device. More specifically, Kibrik (2011) argues that SR marking helps to identify the referent of a mention used as relative pivot by comparing it to the referent of its controller. Note that by the same logic, we expect that SR may also serve to disambiguate controllers themselves. In particular, when the relative pivot of a SS construction precedes its controller, once the relative pivot's reference is known, SS marking should help to identify the controller's referent.

In sections 4 and 5, we try to determine whether SR marking can be analyzed as a referential aid in Mbyá Guaraní, through an analysis of two models of referential choice in this language. In the next section, we give an overview of aspects of the grammar of Mbyá Guaraní that are relevant to the design and interpretation of these models. We focus on the grammar of argument realization,

---

[2]It may at first seem inconsistent to assume that referential aids should affect referential choice, since the choice of referential expression by the speaker belongs to language production, while the notion of referential aid may appear to belong to language comprehension. The inconsistency is only apparent, since speakers are known to take into account their addressee's perspective in language production, although the extent to which they do is a matter of debate. For a discussion of addressee oriented processes in referential choice, see notably Arnold (2008).

and the use of SR markers.

# 3 The Mbyá Guaraní language

## 3.1 General observations

Mbyá Guaraní is a head-marking language, core arguments being cross-referenced on the verb, as illustrated by the following example:[3,4]

(3)  *Ava  o-o     ramo  mboi   o-exa.*
     man  A3-go   DS    snake  A3-see

     'When the man went, the snake saw him.'                                    ([Dooley, 1989](): 97)

As this example also illustrates, there are no articles and no definiteness marking on nouns. Plural marking is optional and only attested with a subset of higher animate nouns. To illustrate, the noun *ava* ('man') has a plural form *ava-kue* ('men'), but the noun *mboi* ('snake(s)') is number neutral. Personal pronouns encode person, number and clusivity features, but not animacy or gender:

(4)  *Ha'e  o-u.*
     3      A3-come

     'He/she/it came.'                                                          ([Dooley, 2016](): 53)

---

[3]Glosses: A1SG: first person singular active inflection; B1SG: first person singular inactive inflection; BDY: information structure boundary marker; CONV: converbial marker; CNTX: counterexpectational; DEM: demonstrative; DS: different subject switch reference marker; HSY: hearsay evidential; NPOSSD non-possessed; OPT: optative; PL: plural; R: linking morpheme; SS: same subject switch reference marker; RESP: response particle;

[4]The glosses of Mbyá Guaraní examples from Dooley (1989) and Florentino (2011) were adjusted to fit our glossing conventions. The glosses of examples from Dooley (2015, 2016) and Veríssimo (2002a,b) are our own, and these examples have been retranslated into English, based on the original Portuguese translation.

Subjects and objects can be cross-referenced on verbs, in the form of prefixes or proclitics that encode their person, clusivity and number. This cross-referencing system is sensitive to the transitivity of the verb, as well as to the lexical class of intransitive verbs. Intransitive verbs belong to one of two classes, called active and inactive,[5] which use different paradigms of prefixes to cross-reference their subject,[6] as illustrated by the following examples:

(5) *A-nha.*

    A1SG-run

    'I ran.'

(6) *Xe-kane'õ.*

    B1SG-tired

    'I am tired.'

With transitive verbs, the active paradigm is used to cross-reference subjects, and the inactive paradigm is used to cross-reference objects. However, only one argument can be cross-referenced.[7] If both arguments are third persons, the subject is cross-referenced. Otherwise, the highest argument on the person hierarchy 1 > 2 > 3 is cross-referenced. In the following example, the verb *a-exa* cross-references its $1^{st}$ person subject:

---

[5]'Active-inactive' and 'active-stative' are common labels for this morphosyntactic alignment in the literature on Tupí-Guaraní languages, see e.g. Seki (1990); Velázquez-Castillo (2002). In the broader typological literature, this alignment is also known as 'agent-patient' or 'Split-S', among other labels, see Mithun (1991).

[6]Contrary to what is observed Paraguayan Guaraní (Velázquez-Castillo, 2002: §3.1), fluid intransitivity has not been reported in Mbyá Guaraní, see (Dooley, 2015: §10) and (Martins, 2003: §2.4).

[7]With the exception of combinations of $1^{st}$ person subject and $2^{nd}$ person object, which are cross-referenced with a portmanteau prefix *ro-*.

(7)  *Ava   a-exa.*

man  A1SG-see

'I saw the man.'

In example (8), the verb *xe-r-exa* cross-references its $1^{st}$ person object. Its implicit subject must be $2^{nd}$ or $3^{rd}$ person:

(8)  *Xe-r-exa.*

B1SG-R-see

'He/she/they/it/you saw me.'

As the previous examples illustrate, all core arguments of the verb can be implicit. In example (3), the object of *o-exa* is not realized as a free argument, and is not cross-referenced on the verb either. In example (8), the subject argument is left implicit.

It is worth emphasizing that there are two senses in which arguments can be said to be implicit. Firstly, there may be no word or phrase that realizes the argument independently of the verb, regardless of cross-reference marking. In this sense, the subjects of the verb *anha* and *xekane'õ* in examples (5) and (6) can be said to be implicit, although they are cross-referenced. Secondly, an argument may be realized neither as a morphologically independent word or phrase, nor as a cross-reference marker. Importantly, independent argument realization and cross-reference marking are different types of processes: while independent argument realization is grammatically optional, the use of cross-reference markers on verbs is subject to deterministic grammatical rules: one and only one argument of the verb must be cross-referenced, the choice of argument being governed by a competition in person and grammatical function, as described above.

## 3.2   Switch reference marking

In Mbyá Guaraní, SR is marked by the particles *vy* (Same Subjects) and *ramo* (Different Subjects) or its reduced form *rã*, both of which occur in the right periphery of the predicate of the marking

clause:[8]

(9) *Ava  o-o    vy  mboi  o-exa.*

man A3-go SS snake A3-see

'When the man went, he saw the snake.'                                    (Dooley, 1989: 97)

(10) *Ava  o-o    ramo mboi  o-exa.*

man A3-go DS    snake A3-see

'When the man went, the snake saw him.'                                   (Dooley, 1989: 97)

This subsection discusses aspects of this construction that are relevant to the study of its effects on referential choice: (i) nature of the pivots, (ii) grammatical and semantic relations between marking and reference clauses and (iii) non-canonical uses of SR markers.

*(i) Nature of the pivots.* Dooley (1989) argues that the pivots of SR in Mbyá Guaraní are subjects, which can be defined as the only cross-referenced argument of intransitive verbs, and the argument of transitive verbs that is cross-referenced with active markers. While the relevance of a grammatical opposition between subjects and objects has been questioned in the analysis of Paraguayan Guaraní (Velázquez-Castillo, 2002), Dooley (2015: §7.1.1) provides evidence that this opposition is active in the grammar of Mbyá Guaraní independently of SR marking. In partic- ular, reflexive possessives are controlled by subject arguments in both intransitive and transitive clauses, the impersonal voice targets subjects, and the subject of an intransitive converb must be coreferential with the subject of a main verb regardless of the transitivity of the latter.

*(ii) Relations between marking and reference clauses.* Dooley (2010, 2015) argues that marking clauses relate to reference clauses by peripheral subordination or ad-clausal modification, rather

---

[8]An anonymous reviewer observes that Tupí-Guaraní languages do not normally display switch reference. Indeed, Mbyá Guarani is the only Tupí-Guaraní language for which Jensen (1998) reports SR marking. Jensen proposes that the DS marker *ramo* is derived from the Proto-Tupí-Guaraní simultaneous/conditional morpheme *-(r)VmV*, while the SS marker *vy* is derived from the PTG serial verb suffix *-áβo*.

than by coordination. Dooley notably observes that (i) the order of the marking and reference clauses does not necessarily reflect the temporal order of the events they describe, (ii) tense, aspect, mood and polarity modifiers have a restricted distribution in the marking clause and (iii) SR constructions are not subject to the coordinate structure constraint on question formation. All of these properties are taken to be characteristic of subordinating relations.

The range of semantic relations that are attested between the marking and reference clauses is also indicative of adverbial subordination. In examples (9) and (10) the semantic relation between the marking clause and the reference clause is one of temporal overlap. However, SR markers are also compatible with causal and conditional interpretations, as illustrated by examples (11) and (12):

(11) *Urutau ma je nd-o-vy'a-i i-juru guaxu vaipa vy .*
     potoo BDY HSY NEG-A3-happy-NEG B3-mouth big very SS

     'The potoo was not happy because he had such a large mouth.'          (Florentino, 2011)

(12) *Ko'ẽ teĩ kuaray nd-o-jope-i rã nd-a'eve-i.*
     dawn CNTX sun NEG-A3-heat-NEG DS NEG-good-NEG

     'If the sun doesn't rise one day, it won't be good.'          (Veríssimo, 2002a)

Another important aspect of SR in Mbyá Guaraní that is illustrated in the previous examples is the fact that both anticipatory and non-anticipatory uses of SR marking are attested in the language. In anticipatory uses, such as (9), (10) and (12), the relative pivot precedes its controller. By contrast, example (11) illustrates the non-anticipatory use of SR marking, in which the relative pivot follows its controller.

*(iii) Noncanonical uses of switch reference.* The reader may wonder whether non-canonical uses of SR markers are attested in Mbyá. With respect to interclausal uses of SR, which we discussed in this section, Dooley (1989, 1992) observes that SR is overwhelmingly used canonically, stating that 98% of occurrences in his corpus indicate whether the marking and reference clause

subjects corefer. The residue of cases that deviate from this pattern are restricted to pairs of pivots whose referents are neither identical nor disjoint, namely when the subjects' referents overlap, stand in a relation of proper inclusion, or when one of the subjects is expletive (has zero reference). In other words, Dooley (1989, 1992) presents interclausal SR in Mbyá as a canonical system. However, SR is also attested in sentence initial connectives, as illustrated in example (13):

(13) *Peteĩ-gue je ava o-o o-i-ny t-ape r-upi. Ha'e vy je o-exa*
    one-time HSY man A3-go A3-be-CONV NPOSSD-path R-on 3 SS HSY A3-see

    *apere'a.*
    guinea.pig

    'Once there was a man who was going on a dirt path. Then, he saw a guinea pig.'

    (Veríssimo, 2002a)

In this example, SR marking is attested in a sentence initial connective *ha'e vy*, which Dooley (2015: §21.5) analyzes as consisting of a third person pronoun *ha'e* interpreted by propositional anaphora, and the SR marker *vy* (SS). Dooley (1992) argues that in sentence initial connectives, the primary use of SR markers is non-canonical: the Same Subject maker indicates a predictable continuation of a line of action, regardless of subject reference, while Different Subject markers indicate a continuation that violates prior expectations. It is important to emphasize that these primarily non-canonical uses of SR marking are restricted to inter-sentential discourse connectives. By contrast, our study is concerned exclusively with post-verbal SR markers in intra-sentential clause chaining constructions.

# 4   Data set and variables

## 4.1   Corpus

The corpus used in the present study consists of 81 narratives written between 1976 and 2002 by 9 Mbyá speakers from the Rio das Cobras community in Paraná, Brazil. This corpus contain 1313

sentences and 14575 tokens. The narratives in the corpus come from two different sources. The first part, which consist of 33 narratives written by two authors between 1976 and 1990, were produced during literacy workshops organized by Robert Dooley and the Summer Institute of Linguistics in Brazil. An interlinearized version of the corpus with a translation into English is available on the Archive of the Indigenous Languages of the America (Dooley, 2011). The second part consists of 48 narratives written by 7 authors and published in 2002 as Veríssimo (2002a) and Veríssimo (2002b) together with their translation into Brazilian Portuguese.

We added four layers of annotation to this corpus: (i) morpheme-by-morpheme interlinear glosses[9], (ii) syntactic dependencies, (iii) coreference relations between referential expressions and (iv) animacy tags on referential expressions. Interlinearization was produced in SIL FieldWorks (Black and Simons, 2008). Interlinearized files were then converted to the CoNLL-U format for dependency annotation in Arborator (Gerdes, 2013). Coreference and animacy anotations were then added in the Webanno software (de Castilho et al., 2016).

Syntactic annotation was done in dependency grammar, in the framework of Universal Dependencies (Nivre et al., 2019). Author (xxxx) describes the annotation principles adopted in the construction of the corpus.

A layer of coreference annotation adds tags to referential expressions, as well as anaphoric relations between referential expressions. Our coreference annotation guidelines follow Komen (2009), with minor adjustments. Overt referential expressions were tagged directly, while null arguments were tagged on their predicate.

Referential expressions were also tagged for their ontological category, which was later mapped to two animacy categories: animate and inanimate.

Our corpus imposes some limitations on the present study. Firstly, insofar as our study is based exclusively on written texts, we acknowledge that its results may not generalize to spoken language production. Note however that the study of written language production has played an important

---

[9]For narratives in Dooley's (2011) corpus, we used Dooley's interlinearization, with minor modifications. The interlinearization of narratives from Veríssimo (2002a) and Veríssimo (2002b) is our own.

role in previous studies of referential choice, notably in the work of Fox (1987), Ariel (1990), Arnold (1998) and Kibrik (1996, 2011), among others. In particular, Kibrik's (2011) model of referential choice, from which this study borrows heavily, is intended as a model of both spoken and written language production. Secondly, our study does not address variation in the grammar of switch reference in Mbyá Guaraní. Mbyá Guaraní is spoken on a large territory that includes parts of Argentina, Brazil and Paraguay, and it goes without saying that we expect dialectal and sociolinguistic variation within this territory. Our study is based exclusively on texts produced by speakers from the state of Paraná in Brazil between 1976 and 2002, and it is possible that its results will not generalize to other variants of Mbyá Guaraní. These limitations will have to be addressed in future research.

## 4.2   Data set construction

Corpus annotation files were exported from Webanno, and processed with a Python script, in order to create the data set used for our analysis. The data set was manually checked for errors and inconsistencies. Each observation in the data set corresponds to a referential mention, together with its associated values for the set of variables we describe in the next section.

Among all observations of referential expressions, we only considered $3^{rd}$ person mentions of discourse old referents for our analysis. Exclusion of $1^{st}$ and $2^{nd}$ person mentions is motivated by the fact that they are all realized by reduced referential expressions in the corpus. Their relevance to the study of referential choice is therefore limited. This is expected under the assumption that referents of $1^{st}$ and $2^{nd}$ person pronouns are always maximally accessible/activated in discourse (Ariel, 1990: §2.1; Kibrik, 2011: §2.3.2). We also excluded mentions of new referents from the analysis, as we are interested in studying the use of switch-reference as a referential aid in the face of referential conflict, which only arises with mentions of discourse old referents.

Finally, we excluded all pronominal referential expressions except personal and demonstrative pronouns. In particular, indefinite pronouns (e.g. *amongue* 'someone') and quantificational pronouns were excluded from the set of observations, since the latter are not referential and the former

introduce new discourse referents.

Note that referential mentions that were excluded from the set of observations were still taken into account as potential antecedents, and in the identification of referential conflict.

## 4.3 Variables

We define as our dependent variable the form of the observed referential mention, `MentionForm`, with two values: `Lexical`, or `Reduced`. A first set of predictors consists of activation factors that encode properties of the context in which the referent is mentioned, among which we include the grammatical function of the mention itself (`MentionFunction`), due to possible parallelism effects:

(14) Activation Factors (Discourse Context):

    a. `ClauseDistance`: distance to the referential mention's antecedent in number of clauses.

    b. `AntecedentForm`: form of the referential mention's antecedent, with two levels: `Lexical` (description or proper name) or `Reduced` (pronoun or zero).

    c. `AntecedentFunction`: grammatical function of the referential mention's antecedent, with three levels: `Subject`, `Object` or `Other` (including obliques and possessors).

    d. `MentionFunction`: grammatical function of the referential mention: with three levels: `Subject`, `Object` or `Other` (including obliques and possessors).

Our model only includes one activation factor that encodes an internal property of the referent, namely animacy:

(15) Activation Factors (Referent's internal properties):

    `MentionAnimacy`: animacy of the referential mention: `Animate` or `Inanimate`.

Referential conflict is encoded in the predictor `Competitors`. Note that when calculating the number of competitors of an observed mention, we included all referents mentioned between the

mention and its antecedent regardless of gender and number. This is motivated by the fact that Mbyá Guaraní does not mark gender on nouns and pronouns, and that nouns are number neutral in the language. We also included referents of first and second person mentions as potential competitors of observed mentions. This is justified by the fact that the texts in our corpus are narratives in which first and second person pronouns most frequently occur in direct reported speech, in which case they often corefer with third person mentions.

(16)  Referential Conflict Filter:

   `Competitors`: number of distinct referents mentioned between the observed mention and its closest antecedent.[10]

The next predictors in our model encode the presence of referential aids in the discourse context, more specifically SR marking. We distinguish two subsets of predictors in this category: those that encode properties of relative pivots, and those that encode properties of controller pivots.

At least three different features of relative pivots and their context are relevant to our study: (i) whether a Same Subject or a Different Subject marker is used, (ii) whether the relative pivot

---

[10]We decided to exclude referents mentioned before the closest antecedent of the referential mention. An alternative would have been to include all referents mentioned since the beginning of a text. We rejected this alternative since we hypothesize following Kibrik (2011) that referential choice is only sensitive to referents that have a sufficient degree of activation in the speaker's working memory. The closest antecedent provides a lower bound on the span of discourse inside which referential conflict may in principle take place: if this antecedent is close enough for its referent to have a degree of activation that warrants mention reduction, then everything else being equal, referents of intervening mentions should be sufficiently activated to generate referential conflict. By contrast, referents that were last mentioned at an arbitrary distance from a target mention may not be close enough to compete with the mention's referent in terms of activation in working memory. A second alternative would have been to select a fixed window, e.g. 10 clauses before a target mention. Since we are uncertain what window length should be chosen in that case, we decided to use the closest antecedent as a lower bound. A third alternative, adopted by Kibrik (2011), would be to calculate a degree of activation for all previously mentioned referents, and to include as competitors all referents whose degree of activation is deemed high enough. Since our model does not rely on the calculation of degrees of activation, this last alternative is unavailable to us.

follows or precedes its controller and (iii) whether the relative pivot's controller is a reduced or a lexical mention. Following Kibrik (2011), we assume that when a referential mention is used as relative pivot, its controller may help to disambiguate its reference. Both the form and position of the controller are relevant in this respect. Indeed, not only do we expect that a controller might help to disambiguate the reference of a relative pivot that follows it, we also expect that this effect might be stronger when the controller is lexical (rather than reduced) and therefore encodes more information about its referent, which may facilitate the disambiguation process.

(17) Referential Aids (Relative Pivots):

    a. `RelativePivot`: indicates whether the referential mention is marked as identical to or different from a controlling subject in a SR construction. There are three levels: `None` (mention unmarked), `DS` (mention marked as different from the controlling subject) or `SS` (mention marked as identical to the controlling subject).

    b. `PivotPosition`: indicates whether the referential mention is marked as a relative participant in a SR construction, and if so, whether it precedes or follows its controller. There are three levels: `None` (mention unmarked), `Pre` (mention is marked as a relative participant and precedes its controller) and `Post` (mention is marked as a relative participant and follows its controller).

    c. `ControllerForm`: indicates whether the referential mention is marked as a relative participant in a SR construction, and if so, what is the form of its controller. There are four levels: `None` (mention unmarked), `Expletive` (mention is marked and controller is an expletive subject), `Reduced` (mention is marked and controller is a reduced mention) and `Lexical` (mention is marked and controller is a lexical mention).

Note that the value `None` for these variables indicates that the referential mention is not a relative pivot, i.e. that it does not participate in a SR construction, or that it does so only as a controller pivot.

    SR marking may also serve to track the reference of controller pivots. When thinking about this

phenomenon, it must be kept in mind that a single subject may serve as the controller of several relative pivots, which may differ from one another with respect to their form and position relative to the controller. Consequently, we only included two predictors that encode whether observations participate in SR constructions as contoller pivots:

(18)  Referential Aids (Controller Pivots):

   a.  `ControllerSS`: indicates whether the referential mention is used as the controller of one or more relative pivots in SS constructions. The predictor is numeric, its possible values being the number of SS relative pivots controlled by the referential mention (0 or more).

   b.  `ControllerDS`: indicates whether the referential mention is used as the controller of one or more relative pivots in DS constructions. The predictor is numeric, its possible values being the number of DS relative pivots controlled by the referential mention (0 or more).

Finally, we include a predictor that encodes whether a referential mention is cross-referenced on its predicate:

(19)  `MentionCrossreference`: indicates whether the referential mention is cross-referenced on its predicate. Three levels: `None` (mention is not cross-referenced), `A` (mention is cross-referenced with an active marker) and `B` (mention is cross-referenced with an inactive marker).

Note that it would not be appropriate to take cross-reference marking into account as part of the encoding of our dependent variable. Firstly, cross-reference marking cuts across the distinction between reduced and lexical referential devices, which is the focus of our study. Secondly, since cross-reference marking is determined categorically based on lexical features of the predicate and grammatical properties of the referential mention (grammatical function and person), it cannot be said that a speaker chooses to use cross-reference marking in the same sense that she chooses to use a full or reduced expression to mention a referent.

| | levels | Lexical | Reduced | p |
|---|---|---|---|---|
| ClauseDistance | Median (IQR) | 2.0 (4.0) | 1.0 (1.0) | <0.001 |
| Competitors | Median (IQR) | 2.0 (3.0) | 1.0 (2.0) | <0.001 |
| MentionFunction | Subject | 435 (27.1) | 1171 (72.9) | <0.001 |
| | Object | 144 (40.3) | 213 (59.7) | |
| | Other | 370 (79.2) | 97 (20.8) | |
| MentionAnimacy | Animate | 668 (33.6) | 1318 (66.4) | <0.001 |
| | Inanimate | 281 (63.3) | 163 (36.7) | |
| MentionCrossreference | None | 512 (59.9) | 343 (40.1) | <0.001 |
| | A | 372 (26.5) | 1032 (73.5) | |
| | B | 65 (38.0) | 106 (62.0) | |
| AntecedentForm | Lexical | 645 (49.1) | 668 (50.9) | <0.001 |
| | Reduced | 304 (27.2) | 813 (72.8) | |
| AntecedentFunction | Subject | 471 (29.1) | 1149 (70.9) | <0.001 |
| | Object | 172 (45.7) | 204 (54.3) | |
| | Other | 306 (70.5) | 128 (29.5) | |
| PivotPosition | None | 892 (41.1) | 1279 (58.9) | <0.001 |
| | Post | 3 (6.7) | 42 (93.3) | |
| | Pre | 54 (25.2) | 160 (74.8) | |
| RelativePivot | None | 892 (41.1) | 1279 (58.9) | <0.001 |
| | DS | 28 (25.7) | 81 (74.3) | |
| | SS | 29 (19.3) | 121 (80.7) | |
| ControllerForm | None | 892 (41.1) | 1279 (58.9) | <0.001 |
| | Expletive | 1 (50.0) | 1 (50.0) | |
| | Lexical | 11 (15.1) | 62 (84.9) | |
| | Reduced | 45 (24.5) | 139 (75.5) | |
| ControllerSS | 0 | 935 (41.0) | 1348 (59.0) | <0.001 |
| | 1 | 13 (9.1) | 130 (90.9) | |
| | 2 | 1 (25.0) | 3 (75.0) | |
| ControllerDS | 0 | 906 (38.7) | 1437 (61.3) | 0.110 |
| | 1 | 42 (50.0) | 42 (50.0) | |
| | 2 | 1 (33.3) | 2 (66.7) | |
| All observations | | 949 (39.1) | 1481 (60.9) | |

Table 1: Summary of the variables in the final data set (n = 2430)

## 4.4   Data set exploration

Our final data set includes 2430 observations, and is summarized in table 1. The last column reports p-values of Wilcoxon rank sum tests for continuous predictors, and of Chi-squared tests for categorical predictors. We observe that 60.9% of all mentions are reduced. Reduced mentions tend to be closer to their antecedents than lexical ones, and they tend to have less competitors. A greater proportion of reduced devices is observed among subjects than non-subjects, and among animate referents than inanimate ones.

Regarding relative pivots of SR constructions, we observe a greater proportion of reduced devices among referential mentions that are marked as relative participants in SS or DS constructions. Among marked mentions, those that follow their controller are more frequently reduced than those that precede it.

The greatest number of relative pivots controlled by a subject in SS constructions is 2, and likewise for DS constructions. We observe a greater proportion of reduced devices among referential mentions that are controllers in SS marking constructions.

Let us now explore the relations between predictors that track referential conflict and SR marking. For this purpose, we will restrict our attention to subjects. Figure 1 displays the proportion of reduced subjects by number of competitors. The grouping variable indicates whether referential mentions are relative pivots in a SR construction, and if so whether a SS or DS marker is used. This plot only considers referential mentions with 12 competitors or less, since there are no reduced mentions among observations with a greater number of competitors. In figure 2, the grouping variable indicates whether referential mentions are relative pivots in a SR construction, and if so whether they precede or follow their controller.

Within the group of subjects that are not marked as relative pivots, the proportion of reduced mentions decreases as the numbers of competitors increases. Within the group of relative pivots on the other hand, the proportion of reduced mentions is greater with one competitor than with no competitors, and for pivots that follow their controller, it does not start to decrease until three competitors are present.
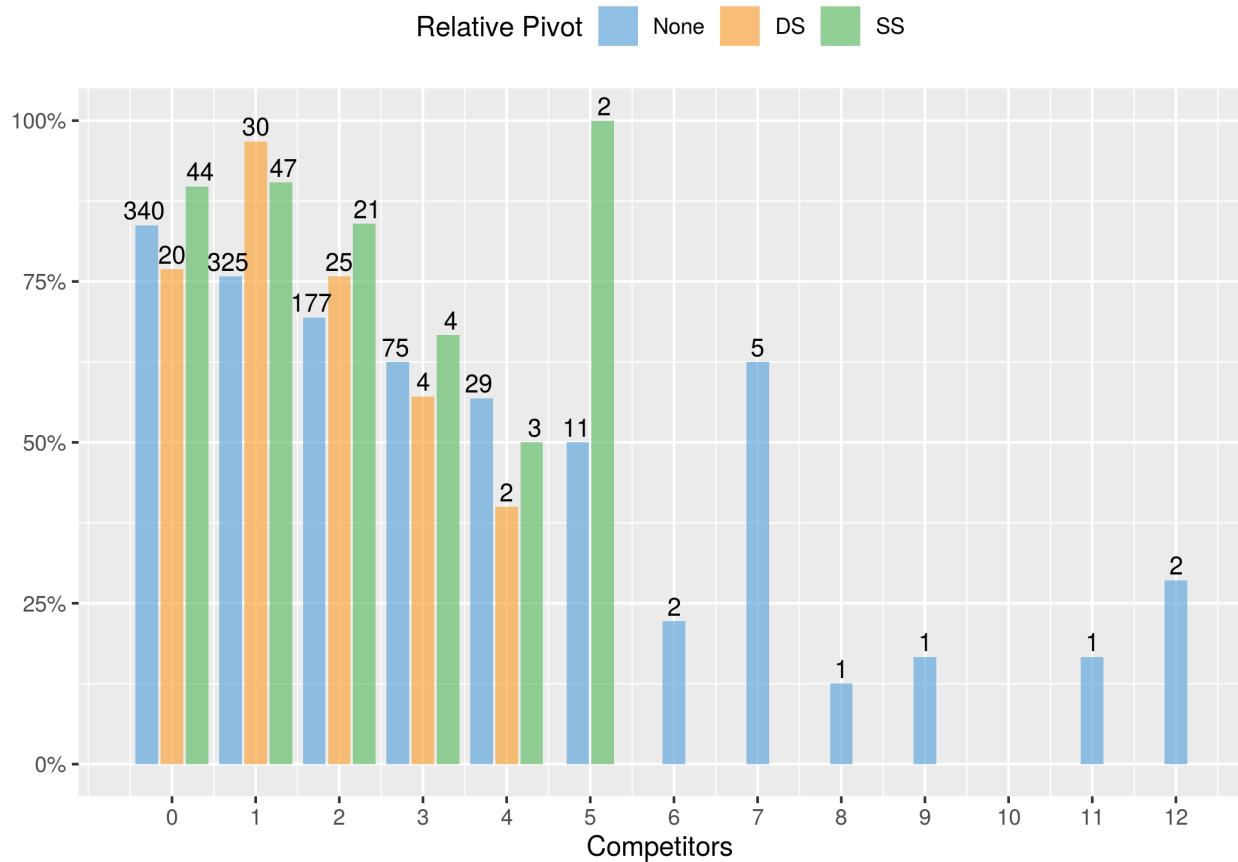
Figure 1: Proportion of reduced subjects by numbers of competitors and pivot marking

Figure 3 displays the proportion of reduced subjects by number of competitors, grouped by number of SS relative pivots controlled by the observed subject. We observe that with zero or one competitor, subjects that are controllers of SS constructions are more frequently reduced than other subjects.

These observations suggest that SR marking and referential conflict may be interacting in a way that is consistent with the analysis of SR as a referential aid: while the relative frequency of reduced mentions generally decreases as the number of competitors increases, this trend is not observed for mentions that are relative pivots or controllers of SS constructions in the presence of a small number of competitors.

In the next section, we move on to a multifactorial analysis of our data. We consider two recursive partitioning models of referential choice, and ask whether the role of predictors that encode
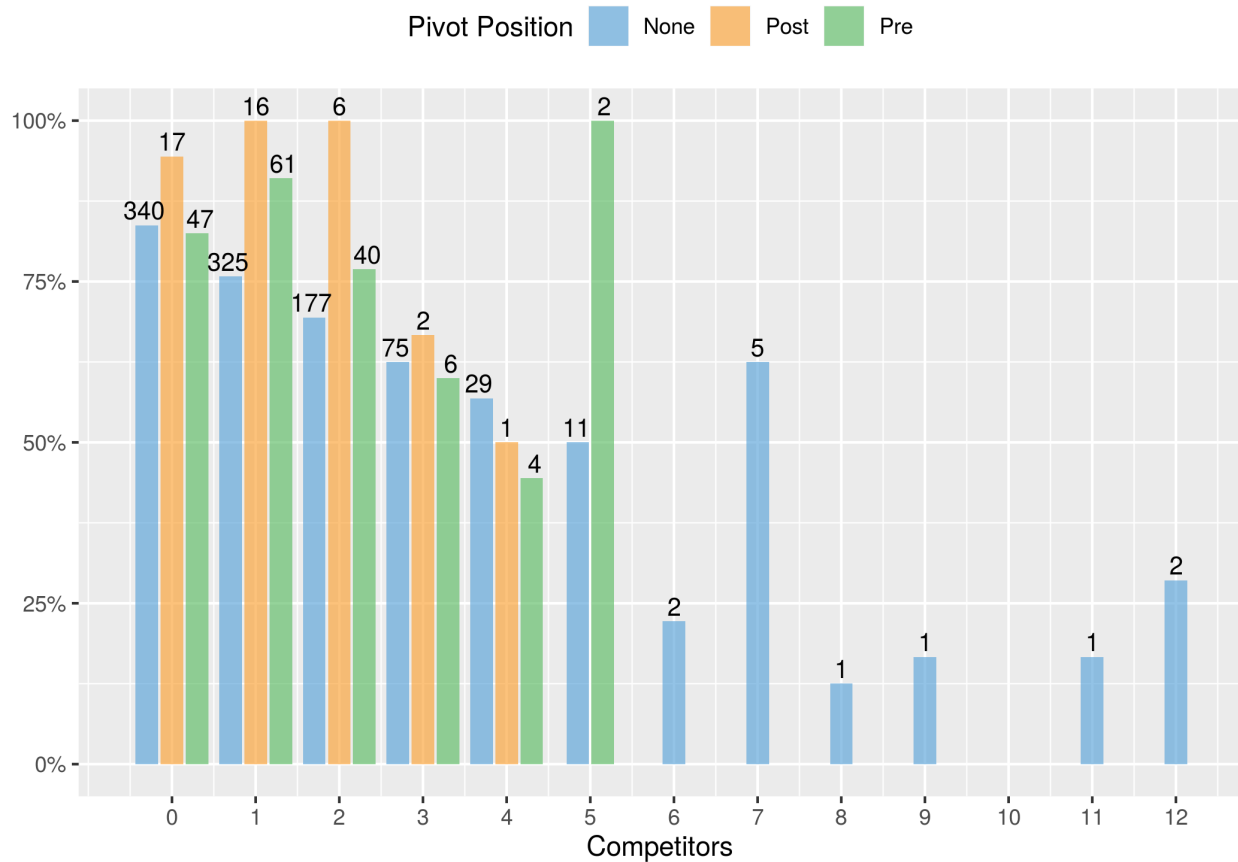
Figure 2: Proportion of reduced subjects by numbers of competitors and pivot position

SR marking in these models supports the analysis of SR as a referential aid.

# 5  Multifactorial analysis

We analyzed referential choice using two recursive partitioning models, conditional inference trees and random forests (Hothorn et al., 2006; Strobl et al., 2009). Both type of models have been argued to be tolerant to unbalanced data sets with multicollinearity (Tagliamonte and Baayen, 2012). For purposes of cross-validation, we partitioned our data into a training set and a test set with the `caret` package (Kuhn, 2008), using a 80%/20% split.

We fitted a conditional inference tree to our training set using the `ctree()` method from the `partykit` library (Hothorn and Zeileis, 2015) in R (R Core Team, 2013). In a conditional
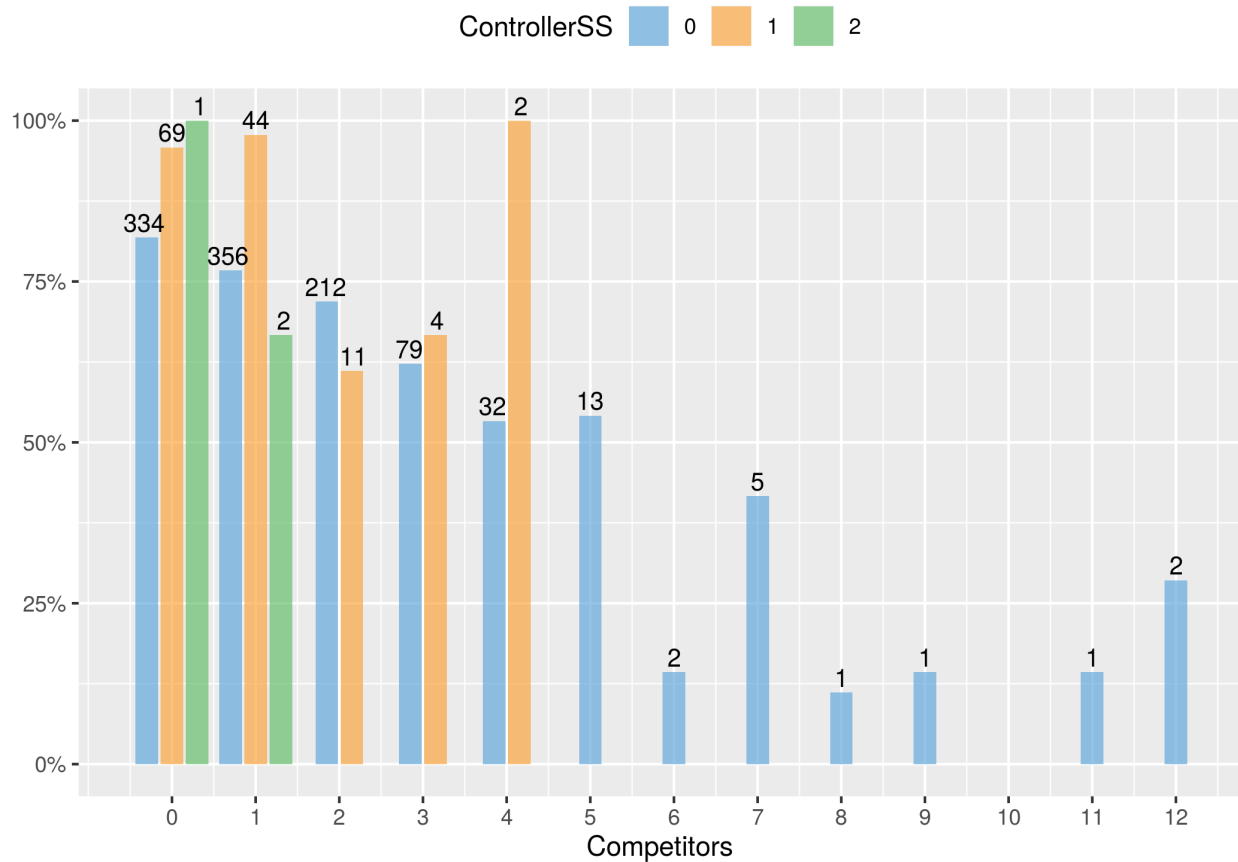
Figure 3: Proportion of reduced subjects by numbers of competitors and number of controlled pivots in SS constructions

inference tree, the data is subjected to successive binary splits. At each split, the predictor that is most significantly associated with the response variable in a series of permutation tests is selected, and the data is split in two subsets, each associated with different levels of the selected variable. The process is then applied recursively to each subset until a stopping criterion is met, which in our case was set as a significance level of 0.05.

On the test set, our conditional inference tree has a classification accuracy of 79.38% (F1 = 70.59%), which is significantly better than the 61.03% baseline (p = < 2.2e-16). This tree, which is represented in figure 4, determines whether a referential mention should be classified as reduced or lexical, depending on the values that it takes for the variables represented on the nodes of the tree. In addition, each terminal node indicates the number of observations it includes from the training set, as well as the relative frequencies of lexical and reduced mentions among these observations.

To illustrate, node 35 consists of observations that are non-core arguments (i.e. neither subjects nor objects) with a reduced antecedent. The model classifies observations of this type as lexical (rather than reduced), i.e. the model predicts that such observations will be realized by lexical referential devices. In the training set, there were 110 observations of this type, of which 60.9% were indeed lexical, and 39,1% were reduced. As Strobl et al. (2009: 328) observe, while the classification of new observations by the model is not a probabilistic process, one may estimate predicted class probabilities based on the relative frequencies in the terminal nodes. Using node 35 as an example again, one may estimate that non-core mentions with a reduced antecedent have a predicted 39,1% probability of reduction. Note that confidence intervals are not available for such estimates (*ibid.*).

Let us now look at the model in more detail. We observe that:

1. Non-core arguments (i.e. mentions that are neither subjects nor objects) are classified as lexical mentions (nodes 33, 34, 35).

2. Among core arguments (i.e. subjects and objects) with more than two competitors (node 24):

   (a) Mentions that are more than four clauses away from their antecedent are classified as lexical (node 30).

   (b) Mentions that are two to four clauses away from their antecedent are classified as reduced (node 29).

   (c) Mentions that are at most one clause away from their antecedent are classified as reduced if they are animate (node 27), and lexical if they are inanimate (node 28).

3. Among core arguments with at most two competitors and that are more than one clause away from their antecedent (node 15):

   (a) Mentions whose antecedent is not a subject are classified as reduced if they are not cross-referenced on their verb (node 22), and they are classified as lexical otherwise (node 23).

(b) Mentions whose antecedent is a subject that is more than three clauses away are classified as lexical (node 20).

(c) Mentions whose antecedent is a subject that is at most three clauses away are classified as reduced (nodes 18 and 19).

4. Core arguments with at most two competitors and that are at most one clause away from their antecedent are classified as reduced (nodes 8, 9, 10, 12, 13, 14).

Overall, the model predicts that only core arguments (subjects and objects) are reduced. Core arguments with a large number of competitors ($> 2$) and that are distant from their antecedent ($> 4$ clauses) are predicted to be realized by lexical devices (node 30). Core arguments with a small number of competitors ($\leq 2$) and that are close to their antecedent ($\leq 1$ clause) are predicted to be realized by reduced referential devices (nodes 8, 9, 10, 12, 13, 14).

The fact that cross-referenced mentions in node 23 are classified as lexical while non-cross-referenced mentions in node 22 are classified as reduced may seem surprising at first. This pattern can actually be interpreted as the manifestation of a grammatical function parallelism effect. Indeed, since we are looking only at third person mentions, mentions that are cross-referenced on their verb must be subjects. Hence, the observations in node 23 are subjects whose antecedents are objects or non-core arguments, while most of the observations in node 22 are likely to be objects, and have object or non-core antecedents. Under the assumption that a mismatch between the grammatical function of a mention and that of this antecedent decreases the likelihood of the mention's reduction, we can make linguistic sense of the model's predictions in nodes 22 and 23.

If we now estimate class probabilities based on relative frequencies in terminal nodes, we observe that among core-arguments with at most two competitors and that are at most one clause away from their antecedent (node 4), referential mentions are more likely to be reduced when they have the same grammatical function as their antecedent (cf. nodes 12 versus 13 and 8/9 versus 10). This can also be interpreted as a manifestation of the grammatical function parallelism effect.

Importantly, this model suggests that SR marking only has a very circumscribed effect on refer-

ential choice, since predictors that track SR marking are only used to split two subsets of observations that correspond to mentions with highly activated referents and a small number of competitors. The first one consists of subjects that are at most one clause away from a subject antecedent, and have at most two competitors (node 7). Among these observations, referential mentions that are relative pivots and whose controller is not expletive (node 9) are slightly more likely to be reduced (95.9%) than other mentions (node 8; 92%). The second subset of observations in which SR comes into play consists of core arguments that are two or three clauses away from a subject antecedent and have at most two competitors (node 17). Among these, referential mentions that are relative pivots in a SR construction (node 19) are more likely to be reduced (88.6%) than other mentions (node 18; 66.4%).

In order to better understand the importance of SR marking relative to other predictors of referential choice, we fitted a random forest model to our data, using the `partykit` library (Hothorn and Zeileis, 2015) in `R`. Random forests (Breiman, 2001) are an ensemble method in which multiple decision trees are fitted to randomly sampled subsets of the training data, and predictions of individual trees are combined to return the prediction of the ensemble. Random forests introduce an additional layer of randomness when fitting individual trees, by restricting the set of possible predictors to be selected at each split to a randomly sampled subset of all predictors (Strobl et al., 2009). Besides being less prone to over-fitting than individual decision trees, random forests allow us to measure the relative importance of each variable in the model, and to estimate their marginal effect using partial dependence plots.

A random forest model was fitted to our training set, growing 1000 trees and using unbiased variable selection (Hothorn et al., 2006) with four candidate variables at each split (mtry = 4). The model also has a classification accuracy of 79,38% (F1 = 70.76%), significantly higher than the 61.03% baseline ($p < 2e-16$). Note that this accuracy was calculated on our test set, rather than using the out-of-bag predictions of the model.

Figure 5 presents measures of importance for each predictor in the model. Horizontal bars represent the loss of classification accuracy that results from permuting the values of a predictor over

Figure 4: Conditional Inference Tree model of Referential Choice

all observations, so that the association between this predictor and the response variable is broken (Strobl et al., 2008: 335). The most important predictor in the model is `MentionFunction`, and the least important one is `MentionAnimacy`. Note that the absolute values of the permutation importance score should not be interpreted, since they depend on specific characteristic of the data set (*ibid.*). What figure 5 contributes to the interpretation of the model is a ranking of predictors' importance. In addition, it provides an indication of which predictors contribute to the model classification accuracy: predictors that have a null or negative importance, or whose positive importance is in the same range as the negative predictors can be discarded (Strobl et al., 2008: 336). According to this criterion, `MentionAnimacy` and `ControllerDS` can be deemed irrelevant.

We observe that three of the four most important predictors in the model (`ClauseDistance`, `AntecedentFunction` and `AntecedentForm`) are activation factors in Kibrik's (2011) model of referential choice. `MentionFunction` is known to interact with `AntecedentFunction` in referential choice across languages, in the form of parallelism effects. Among all other predictors, `ControllerSS`, which encodes whether a mention is used as a controller in a SS construction, is ranked relatively high (0.014). Two other predictors that encode different dimensions of SR marking, `ControllerForm` (0.005) and `PivotPosition` (0.005), are ranked higher than `Competitors` (0.002), which encodes referential conflict. This suggests that SR marking has a non-negligible effect on referential choice.

While figure 5 provides information on the relative importance of the predictors in the model, it does not indicate how each of them affects the model's predictions. In order to investigate these effects, we constructed partial dependence plots in figure 6, using the `pdp` library in R (Greenwell, 2017). Each plot represents the predicted probability of mention reduction for different levels of a predictor in the random forest model, when all observations in the data set are forced to assume these levels. To illustrate, the partial dependence plot for `MentionFunction` shows that while the predicted probability of reduction for non-core arguments (`Other`) is about 39%, it increases to 64% for objects and 66% for subjects. More generally, we observe that the closer a mention is to its antecedent, and the less competitors it has, the more likely it is to be reduced. In addition, being
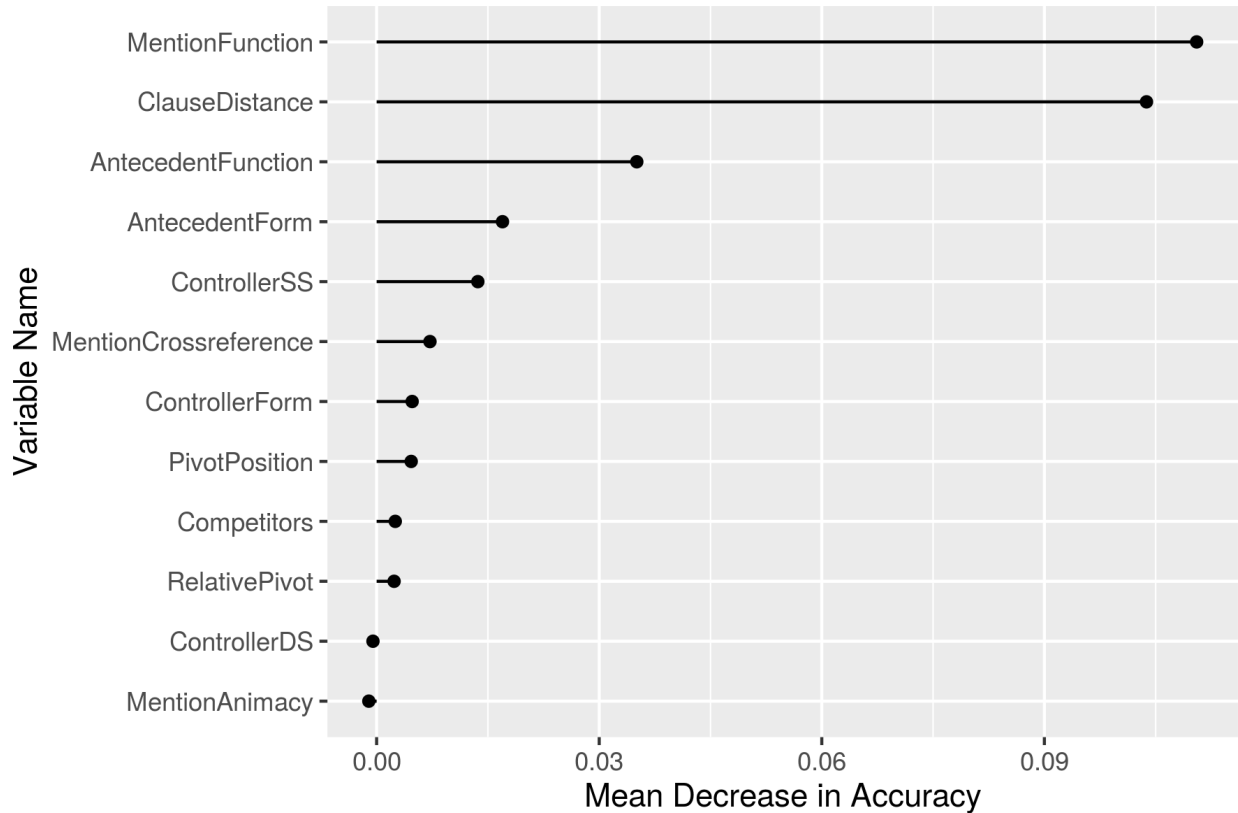
Figure 5: Variable importance in the random forest model.

animate, having a more prominent grammatical function, having an antecedent that is reduced and having a more prominent antecedent all make it more likely for a mention to be reduced.

Predicted probabilities of reduction for different levels of predictors encoding SR are presented in more detail in table 2. Mentions that are marked as relative pivots in a SR construction or a controller pivot in a SS construction are more likely to be reduced, although the observed effect is small. It is interesting to note that relative pivots are not more likely to be reduced than non-pivots when their controller is expletive (see `ControllerForm: None` versus `Expletive`). This is consistent with the view of SR marking as a referential aid: if the controller of a relative pivot is an expletive subject, SR marking (more specifically, Different Subject marking) will not help to disambiguate the relative pivot's reference, and should therefore have no effect on referential choice. Likewise, increasing the number of DS marked pivots of a controller pivot does not make it more likely that the latter will be reduced (see `ControllerDS`). This is also consistent with the

analysis of SR marking as a referential aid: since the relative pivots are marked as different from their controller, SR marking may fail to act as a referential aid in that case.

|  | levels | Reduction probability |
|---|---|---|
| ControllerForm | None | 60.58% |
|  | Expletive | 60.09% |
|  | Lexical | 63.33% |
|  | Reduced | 62.77% |
| PivotPosition | None | 60.56% |
|  | Post | 62.81% |
|  | Pre | 62.99% |
| RelativePivot | None | 60.64% |
|  | DS | 62.33% |
|  | SS | 62.87% |
| ControllerSS | 0 | 60.65% |
|  | 1 | 63.95% |
|  | 2 | 62.95% |
| ControllerDS | 0 | 60.98% |
|  | 1 | 59.36% |
|  | 2 | 59.36% |

Table 2: Predicted probabilities of mention reduction for predictors encoding dimensions of SR

# 6 Discussion and Concluding Remarks

## 6.1 Referential choice in Mbyá Guaraní

Our two models are consistent with existing theories of referential choice. The most important predictors in these models are the mention's distance to its antecedent, and the grammatical functions of the mention and its antecedent. A rich literature has argued that these predictors have an important effect on the degree of accessibility of referents in discourse, which in turn governs the form of referential expressions. Mentions that are closer to their antecedent have been argued to have more accessible referents (Givón, 1983; Ariel, 1990; Arnold, 1998). Grammatical function has also been shown to affect referential choice, mentions being more likely to be reduced when
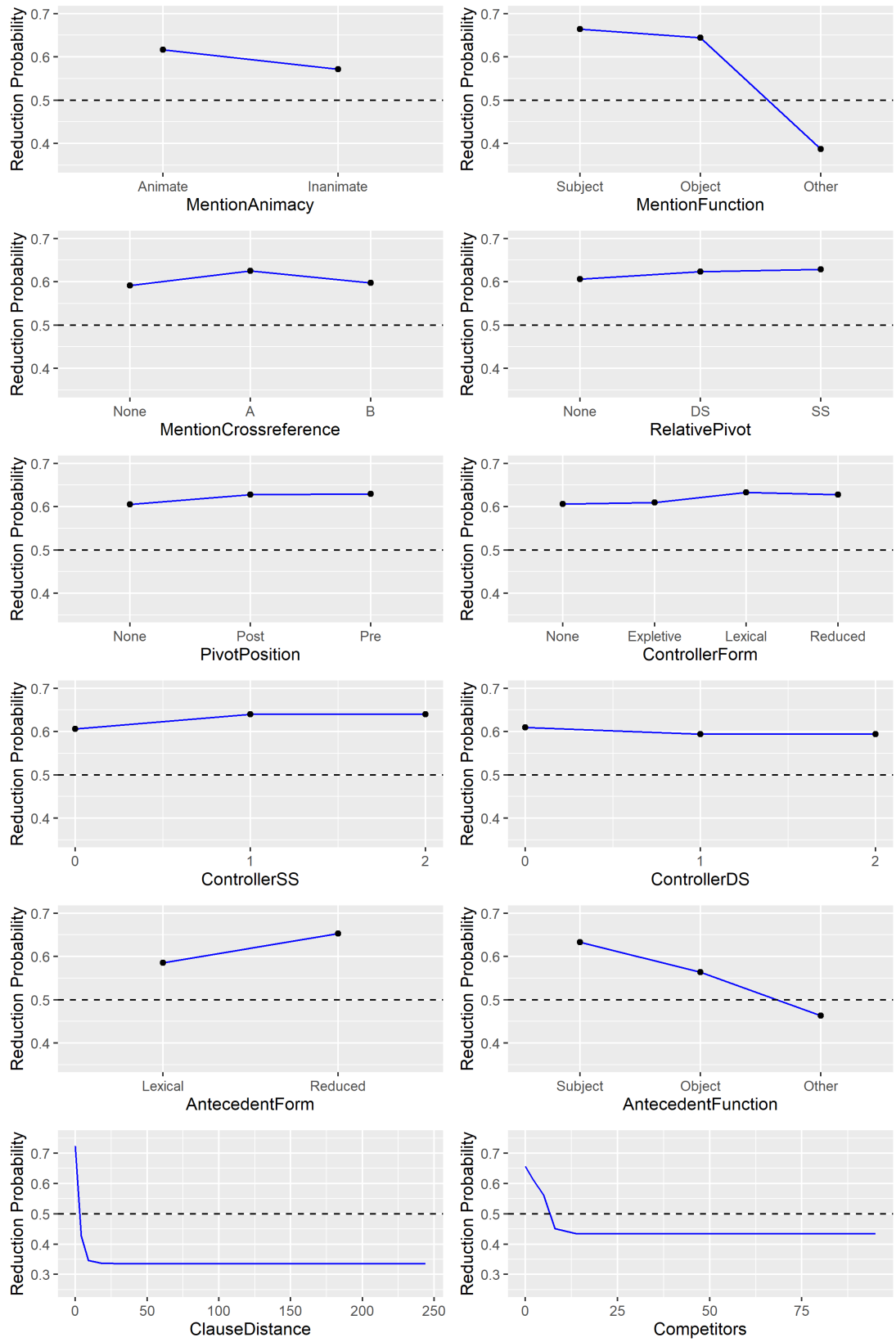
Figure 6: Predicted probabilities of mention reduction for all predictors

their antecedent is a subject (Brennan et al., 1987; Brennan, 1995; Arnold, 2001, 2010).

Our decision tree model also highlighted the effect of grammatical function parallelism on referential choice: mentions are more likely to be reduced when they have the same grammatical function as their antecedent. Such effects have been documented cross-linguistically in production (Arnold, 1998, 2003) and in interpretation (Chambers and Smyth, 1998). In our decision tree, among the set of core arguments with less than three competitors and that are at most one clause away from their antecedents, we observed that the relative frequency of reduction was higher for mentions that have the same grammatical function as their antecedent. A parallelism effect was also observed among referential mentions that are more than one clause away from their antecedent and that have less than three competitors. In this group, the relative frequency of reduction was lower for cross-referenced mentions whose antecedent is an object or a non-core argument. Since third person arguments are only cross-referenced if they are subjects, this can also be interpreted as a manifestation of the parallelism effect.

Note that the most important predictors of mention reduction in our models correspond to activation factors in Kibrik's (2011) model of referential choice, which determine whether a referent is sufficiently activated to be mentioned with a reduced device. In Kibrik's model, the grammatical function of the antecedent and its distance to the mention contribute directly to the degree of activation of the referent in working memory. While the degree of activation of a referent is not affected by the grammatical function of its prospective mention, Kibrik (2011: 55, 453) argues following Tomlin (1995) that referents that are in the focus of the speaker's attention tend to be mentioned in subject position, and that attention is more likely to be focused on highly activated referents, which explains the greater relative frequency of reduced devices in subject position (Chafe, 1994).

Finally, it is worth noting that our decision tree and random forest models are consistent with Dooley's (1976) analysis of referential choice in Mbyá Guaraní, which he calls the "naturalness problem" for participants in narratives. Although referential choice is not Dooley's (1976) focus, Dooley identifies a parallelism effect on the referential form of core-arguments. He also observes that participants who are re-introduced in the narrative after having been eclipsed by other partici-

pants must be referred to with lexical devices, which points to the effect of referential conflict and distance to antecedent on referential choice.

## 6.2  Switch reference as a referential aid

While SR marking increases the likelihood of mention reduction both in our decision tree model and in our random forest model, this effect appears to be very circumscribed and to have a small size. Firstly, figures 1, 2 and 3 show that many subjects are reduced in the presence of competitors, even in the absence of SR marking. Secondly, variable selection in our decision tree model suggests that SR marking only affects referential choice for a very circumscribed subset of subjects, which are at most three clauses away from an antecedent that is itself a subject. This suggests that the effect of SR marking on referential choice is restricted to mentions with highly activated referents. Finally, our random forest model suggests that the importance of SR marking in referential choice is marginal relative to the grammatical function of mentions, the distance to their antecedents and the grammatical function of their antecedents. This was observed in the ranking of predictors by conditional variable importance. In addition, partial dependence plots for the model suggest that the effect of SR marking on referential choice is small relative to that of the most important predictors. We believe that while these results are compatible with Kibrik's (2011) analysis of SR marking as a referential aid, they also highlight the fact that this function might be peripheral in the grammar of canonical SR in Mbyá Guaraní.

On the one hand, it is expected in Kibrik's (2011) theory that referential aids should only affect referential choice for mentions whose referent has a degree of activation that is high enough to warrant mention reduction, and even then, only when referential conflict would prevent the use of a reduced mention if left unchecked. The role played by SR marking in our decision tree model is consistent with these expectations. Predictors that encode SR marking are selected only for a subset of mentions whose referent is highly activated, and has a low number of competitors. In this case, SR marking increases the probability of mention reduction.

On the other hand, our models of referential choice suggest that if Mbyá Guaraní speakers

refrained from using SR marking altogether, there would be little impact on the relative frequency of reduced mentions in narratives. In our decision tree model, in the largest subset of observations that is split by SR marking (node 7, n = 572), the relative frequency of reduced mentions among marked subjects is only 3.9% greater than among unmarked subjects, 92% of which are reduced. In the second largest subset of observations split by SR marking (node 17, n = 296), this difference is much greater at 22.2%, but the majority of unmarked mentions (66.4%) is still reduced. In our random forest model, among predictors that track SR marking, the largest effect was observed for `ControllerSS`, and only corresponds to an increase of 3.30% in the probability of mention reduction.

## 6.3   Concluding remarks: switch reference and its functions

If the use of SR marking as a referential aid is taken as a proxy for its general functionality as a disambiguation mechanism, these results call into question analyses that take the latter to be the primary function of SR marking. Indeed, if it has such a small and restricted effect on one of the processes that it is expected to affect the most, is it reasonable to assume that the disambiguation of reference is the primary function of SR marking? One may suspect that Mbyá Guaraní speakers have other reasons to use SR markers, beyond their use as referential aids, and that these reasons may be as important as the need to disambiguate referential mentions.

This raises the question of what may be the functions of canonical SR in Mbyá Guaraní beyond reference tracking. In addition to its use in the disambiguation of reference, canonical SR in Mbyá Guaraní also serves to indicate clause linkage. SR markers belong to a paradigm of post-verbal particles that mark (co)subordination relations. While subordinating conjunctions encode specific semantic relations between events or propositions, SR markers are underspecified in this respect, and are compatible with temporal, causal and non-conterfactual conditional relations (Dooley, 1999, 2015). Dooley (2010) argues furthermore that while all (co)subordinating conjunctions can be used to form clause chains in Mbyá Guaraní, SR markers differs from other conjunctions insofar as clauses with SR markers can be either backgrounded or foregrounded, while clauses with

conjunctions that encode specific inter-clausal relations can only be backgrounded. Viewed from this perspective, the use of SR markers should be analyzed as a solution to two decision problems that speakers must face when combining clauses to build a larger stretch of discourse: the decision to combine two clauses by linking them rather than by expressing them as independent sentences, and the decision to use an underspecified SR marker rather than a more specific conjunction when linking clauses. This suggests that SR markers in Mbyá Guaraní may serve a number of functions related to discourse coherence and information packaging, even when they are interpreted canonically as markers of referential identity between pivots, a point emphasized by van Gijn (2012) for other South American SR systems. Note that this observation is not inconsistent with the analysis of SR as a referential aid. As Kibrik observes:

> "[...] referential aids are not necessarily inherent and dedicated disambiguation devices, they mostly exist in languages for other purposes, and their usefulness in the preclusion of referential conflicts is a by-product of their use with separate and specialized semantic functions." (Kibrik, 2011: p. 65)

In sum, although the present study supports the analysis of SR marking as a referential aid in Mbyá Guaraní, it also suggests that its effect on referential choice is marginal and does not exhaust its functionality. These results support studies of SR marking which emphasize that the disambiguation of reference may not be the primary function of all SR systems, without however denying the importance of reference tracking and disambiguation in the characterization of SR (Givón, 1983; Stirling, 1993; van Gijn, 2012).

# References

Almor, A. and Nair, V. A. (2007). The form of referential expressions in discourse. *Language and Linguistics Compass*, 1(1-2):84–99.

Ariel, M. (1990). *Accessing NP Antecedents*. Routledge, London.

Arnold, J. E. (1998). *Reference form and discourse patterns*. PhD thesis, Stanford University.

Arnold, J. E. (2001). The effect of thematic roles on pronoun use and frequency of reference continuation. *Discourse Processes*, 31:137–162.

Arnold, J. E. (2003). Multiple constraints on reference form: null, pronominal, and full reference in Mapudungun. In Bois, J. W. D., Kumpf, L. E., and Ashby, W. J., editors, *Preferred argument structure: grammar as architecture for function*, pages 225–245. John Benjamins, Amsterdam.

Arnold, J. E. (2008). Reference production: production-internal and addressee-oriented processes. *Language and Cognitive Processes*, 23(4):495–527.

Arnold, J. E. (2010). How speakers refer: the role of accessibility. *Language and Linguistics Compass*, 4(4):187–203.

Arnold, J. E. and Griffin, Z. M. (2007). The effect of additional characters on choice of referring expression: Everyone counts. *Journal of Memory and Language*, 56:521–536.

Author, A. (xxxx). Blinded.

Baker, M. C. and Camargo Souza, L. (2019). Switch-reference in American languages: A synthetic overview. In *The Routledge Handbook of North American Languages*, pages 210–232. Taylor and Francis.

Black, A. and Simons, G. (2008). The SIL Fieldworks Language Explorer Approach to Morphological Parsing. In *Computational Linguistics for Less Studied Languages: Texas Linguistics Society, 10*, pages 37–55. CSLI Publications.

Breiman, L. (2001). Random forests. *Machine Learning*, 45:5–32.

Brennan, S., Friedman, M. W., and Pollard, C. (1987). A centering approach to pronouns. In *Proceedings from the 25th Annual Meeting of the Association for Computational Linguistics*, pages 155–162. CSLI Publications.

Brennan, S. E. (1995). Centering attention in discourse. *Language and Cognitive Processes*, 10:137–167.

Chafe, W. (1994). *Discourse, consciousness, and time: The flow and displacement of conscious experience in speaking and writing*. University of Chicago Press, Chicago.

Chambers, C. G. and Smyth, R. (1998). Structural parallelism and discourse coherence: a test of centering theory. *Journal of Memory and Language*, 39(4):593–608.

de Castilho, R. E., Mújdricza-Maydt, E., Yimam, S. M., Hartmann, S., Gurevych, I., Frank, A., and Biemann, C. (2016). A Web-based Tool for the Integrated Annotation of Semantic and Syntactic Structures. In *Proceedings of the Workshop on Language Technology Resources and Tools for Digital Humanities LT4DH*, pages 11–17, Osaka, Japan. https://webanno.github.io/webanno/.

Dooley, R. A. (1976). Participants in Guaraní narrative. Technical Report Arquivo Linguístico no. 035, Sociedade Internacional de Linguística, Porto Velho.

Dooley, R. A. (1989). Switch reference in Mbyá Guaraní: a fair-weather phenomenon. *Work Papers of the Summer Institute of Linguistics, U. of North Dakota Session*, 33:93–119.

Dooley, R. A. (1992). When switch reference moves to discourse: developmental markers in Mbyá Guarani. In Hwang, S. J. J. and Merrifield, W. R., editors, *Language in context: essays for Robert E. Longacre.*, pages 97–108. University of Texas at Arlington, Arlington.

Dooley, R. A. (1999). A noncategorial approach to coherence relations: switch reference constructions in Mbyá Guaraní. In Loos, E., editor, *Logical relations in discourse*, pages 219–242. Summer Institute of Linguistics, Dallas.

Dooley, R. A. (2010). Foreground and background in Mbyá Guaraní clause chaining. In McElhanon, K. A. and Reesink, G., editors, *A mosaic of languages and cultures: studies celebrating the career of Karl J. Franklin*, volume 19 of *SIL e-books*, pages 90–110. SIL International, Dallas.

Dooley, R. A. (2011). Mbyá Guaraní collection of robert dooley. The Archive of the Indigenous Languages of Latin America: www.ailla.utexas.org. Media: text. Access: 100% restricted. PID ailla:119734.

Dooley, R. A. (2015). Léxico Guarani, dialeto Mbyá: introdução. Summer Institute of Linguistics.

Dooley, R. A. (2016). Léxico Guarani, dialeto Mbyá. Summer Institute of Linguistics.

Du Bois, J. (1987). The discourse basis of ergativity. *Language*, 63:805–855.

Florentino, N. (2011). Ka'aguy regua ndovy'ai okuapyague. In Dooley, R., editor, *Mbyá Guaraní Collection of Robert Dooley*. The Archive of the Indigenous Languages of Latin America: www.ailla.utexas.org. Media: text. Access: 100% restricted. PID ailla:119734.

Fox, B. (1987). Anaphora in popular written English narratives. In Tomlin, R., editor, *Coherence and Grounding in Discourse*, pages 157–174. John Benjamins, Amsterdam/Philadelphia.

Francik, E. P. (1985). *Referential choice and focus of attention in narratives (discourse anaphora, topic continuity, language production)*. PhD thesis, Stanford University.

Fukumara, K. and van Gompel, R. P. G. (2011). The effect of animacy on the choice of referring expression. *Language and Cognitive Processes*, 26(10):1472–1504.

Fukumara, K., van Gompel, R. P. G., and Pickering, M. J. (2010). The use of visual context during the production of referring expressions. *The Quarterly Journal of Experimental Psychology*, 63:1700–1715.

Gatt, A., Krahmer, E., van Deemter, K., and van Gompel, R. P. G. (2014). Models and empirical data for the production of referring expressions. *Language, Cognition and Neuroscience*, 29(8):899–911.

Gerdes, K. (2013). Collaborative dependency annotation. In *Proceedings of the Second International Conference on Dependency Linguistics (DepLing 2013)*, pages 88–97.

Givón, T. (1983). Topic continuity in discourse: an introduction. In Givón, T., editor, *Topic Continuity in Discourse: A Quantitative Cross-Language Study*, pages 3–41. John Benjamins, Amsterdam.

Givón, T. (2017). *The Story of Zero*. John Benjamins Publishing Company, Amsterdam/Philadelphia.

Greenwell, B. M. (2017). pdp: An r package for constructing partial dependence plots. *The R Journal*, 9(1):421–436.

Gundel, J. K., Hedberg, N., and Zacharski, R. (1993). Cognitive status and the form of referring expressions in discourse. *Language*, 69:274–307.

Haig, G. and Schnell, S. (2016). The discourse basis of ergativity revisited. *Language*, 92:591–618.

Haiman, J. and Munro, P. (1983). Introduction. In Haiman, J. and Munro, P., editors, *Switch-reference and universal grammar*, pages ix–xv. John Benjamins, Amsterdam.

Hothorn, T., Hornik, K., and Zeileis, A. (2006). Unbiased recursive partitioning: a conditional inference framework. *Journal of Computational and Graphical Statistics*, 15:651–674.

Hothorn, T. and Zeileis, A. (2015). partykit: A modular toolkit for recursive partytioning in R. *Journal of Machine Learning Research*, 16:3905–3909.

Jacobsen, W. H. (1967). Switch-reference in Hokan-Coahuiltecan. In Hymes, D. and Bittle, W., editors, *Studies in Southwestern ethnolinguistics*, pages 238–263. Mouton, The Hague.

Jensen, C. (1998). The use of coreferential and reflexive markers in Tupí-Guaraní languages. *Journal of Amazonian Languages*, 1(2):1–49.

Kibrik, A. (2011). *Reference in Discourse*. Oxford University Press, Oxford.

Kibrik, A. A. (1996). Anaphora in russian narrative discourse: a cognitive calculative account. In Fox, B., editor, *Studies in Anaphora*, pages 255–304. John Benjamins, Amsterdam/Philadelphia.

Komen, E. R. (2009). Coreference annotation guidelines. http://repository.ubn.ru.nl/bitstream/handle/2066/78810/78810.pdf.

Kuhn, M. (2008). Building predictive models in R using the caret package. *Journal of Statistical Software, Articles*, 28(5):1–26.

Maria Inês Ladeira (2018). Guarani Mbya. In Ricardo, F. P., editor, *Povos Indígenas no Brasil*. Instituto Socioambiental.

Martins, M. F. (2003). *Descrição e Análise de Aspectos de Gramática do Guarani Mbyá*. PhD thesis, State University of Campinas.

McKenzie, A. (2015). A survey of switch-reference in North America. *International Journal of American Linguistics*, 81(3):409–448.

Mithun, M. (1991). Active/agentive case marking and its motivation. *Language*, 67(3):510–554.

Nivre, J., Abrams, M., and Željko et al., A. (2019). Universal Dependencies 2.4. LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

R Core Team (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna.

Roberts, J. R. (1988). Amele switch-reference and the theory of grammar. *Linguistic Inquiry*, 19(1):45–63.

Roberts, J. R. (2017). A typology of switch reference. In Aikhenvald, A. Y. and Dixon, R. M. W., editors, *The Cambridge Handbook of Linguistic Typology*, Cambridge Handbooks in Language and Linguistics, pages 538–573. Cambridge University Press.

Seki, L. (1990). Kamaiurá (Tupí-Guaraní) as an activestative language. In Payne, D. L., editor, *Amazonian linguistics: Studies in lowland South American languages*, pages 367–391. University of Texas Press, Austin.

Stirling, L. (1993). *Switch-reference and Discourse Representation*. Cambridge University Press, Cambridge.

Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., and Zeileis, A. (2008). Conditional variable importance for random forests. *BMC Bioinformatics*, 307(9):1471–2105.

Strobl, C., Malley, J., and Tutz, G. (2009). An introduction to recursive partitioning: Rationale, application, and characteristics of classification and regression trees, bagging, and random forests. *Psychological Methods*, 14(4):323–348.

Tagliamonte, S. A. and Baayen, R. H. (2012). Models, forests, and trees of york english: *Was/were* variation as a case study for statistical practice. *Language Variation and Change*, 24:135–178.

Tomlin, R. (1995). Focal attention, voice and word order: An experimental cross-linguistic study. In Downing, P. and Noonan, M., editors, *Word Order in Discourse*, pages 517–554. John Benjamins, Amsterdam/Philadelphia.

van Gijn, R. (2012). Switch-attention (aka switch-reference) in south-american temporal clauses: Facilitating oral transmission. *Linguistic Discovery*, 1(10):112–127.

van Gijn, R. (2016). Switch-reference: An overview. In van Gijn, R. and Hammond, J., editors, *Switch-reference 2.0*, volume 114 of *Typological Studies in Language*, pages 1–54. John Benjamins.

Velázquez-Castillo, M. (2002). Grammatical relations in active systems. The case of Guaraní. *Functions of Language*, 9:133–167.

Veríssimo, A. T. (2002a). *Opa mba'e re nhanhembo'e aguã 1*. Nhombo'ea Guarani: Mundo Indígena, Laranjeiras do Sul.

Veríssimo, A. T. (2002b). *Opa mba'e re nhanhembo'e aguã 2*. Nhombo'ea Guarani: Mundo Indígena, Laranjeiras do Sul.