

Generating Determiners and Quantifiers in Hebrew

Yael Dahan Netzer and Michael Elhadad
Ben Gurion University

Department of Mathematics and Computer Science, Beer Sheva, 84105, Israel
(yaeln|elhadad)@cs.bgu.ac.il

Abstract

This paper presents a specific part of HUGG, a generation grammar for Hebrew. This part deals with determiners and quantifiers. Our main goal is to determine which set of features must be present in the input to the generation grammar to control the generation of complex determiners and quantifiers.

Hebrew determiners are characterized by two properties: (1) definiteness is marked in several places in the NP and it interacts with compound constructs; (2) the order of appearance of determiners and quantifiers within a complex NP is flexible but still restricted.

HUGG is developed with the goal to design an input specification language for syntactic realization as close as possible for English and Hebrew, to allow easy development of bilingual generation applications. We show in this paper how the original set of features controlling the generation of determiners in SURGE, a large generation grammar for English, has been enriched to account for the specificities of Hebrew determiners and how, in the process, SURGE has been modified and improved.

We present a set of functional features – organized around the categories of amount, partitive and identity – and show how these features determine a variety of syntactic constructs.

1 Introduction

A text generation system generally includes two main components: content determination (that decides what to say) and surface realization (that decides how to say it). Recent research has led to the development of reusable surface realization systems which encapsulate knowledge of syntax for specific languages. Most of this research has been conducted for English, and large generation grammars for English are now available [4], [9].

Multi-lingual generation now appears as the “next frontier” in generation research: it is often easier to develop applications to generate text in different languages than it is to translate text. When developing a multi-lingual application, the engineering strategy is to develop a single content determination module and to plug it – more or less unchanged – into several surface realization modules for different languages – each encapsulating knowledge of a different language.

The main problem to address in this context is to determine which interface is appropriate between the content determination module and the linguistic modules. We are investigating this problem for the case of bilingual generation English-Hebrew. We start from an existing surface realization system for English, SURGE, and have developed a generation grammar for Hebrew called HUGG with the goal that the input specification language for both grammars be as close as possible.

This paper illustrates the general issue of designing a bilingual generation system by focusing on the specific part of HUGG that deals with determiners and quantifiers. Our main goal is to

determine which set of features must be present in the input to the generation grammar to control the generation of complex determiners and quantifiers. Because the syntactic structure of English and Hebrew determiners is very different, we use functional-semantic features in the input (as opposed to syntactic features). We present a set of such features and evaluate (1) how the features can account for specific syntactic constructs in Hebrew and (2) how the same features can be used in English.

Syntactically, Hebrew determiners are characterized by two properties: (1) definiteness is marked in several places in the NP and it interacts with compound constructs; (2) the order of appearance of determiners and quantifiers within a complex NP is flexible but still restricted. The demands that Hebrew syntax puts on the generation system are, therefore, quite different from those put by the English syntax.

In the rest of the paper, we first present previous work in the description of Hebrew determiners, and present the treatment of the determiner sequence in English for generation as implemented in SURGE. We then focus on the two issues of definiteness and order of constituents and present our treatment of the Hebrew determiner construction using a set of functional features.

We derive from this analysis an input specification for noun phrases for a generation grammar. This input specification extends the set of features used in SURGE and remains compatible with the English grammar, while providing a wide coverage of the Hebrew constructions. This case study of the determiner system illustrates how a single set of functional features can drive both the English and Hebrew generation systems despite their syntactic differences. More generally, it demonstrates the methodology followed to develop a sound bilingual generation system.

1.1 Definiteness

As a polydefinite language, Hebrew marks definiteness in several places in the NP, unlike English which only marks it in the definite/indefinite articles. Because of the complex nature of definiteness marking in Hebrew, we have found that the marking of definiteness as realized in SURGE had to be extended.

Non-Definite	Definite	Definite with smixut
(d1) <i>sefer meyuHad</i> <i>book special</i> <i>a special book</i>	(d2) <i>ha-sefer ha-meyuHad</i> <i>the-book the-special</i> <i>the special book</i>	(d3) <i>sefer ha-madaw ha-meyuHad</i> <i>book the-science the-special</i> <i>the special science book</i>

Examples d1/d2 show how the definiteness marking (ha-) appears on both the head and the modifier. Example d3 shows that when a noun-compound construction (called smixut) is used, then the definiteness marking disappears on the head. The head, however, remains semantically definite. Hebrew demonstratives can also be marked with the definiteness ha-:

Unmarked	<i>sfarym eleh</i>	<i>books those</i>
Marked	<i>ha-sfarym ha-eleh</i>	<i>the-books the-those</i>

The two forms are definite in their meaning (demonstrative implies definite), but the additional definite marking is still possible in Hebrew. The two forms differ in their pragmatic usage (the unmarked form tends to refer to previous text and the marked form to deixis proper [11] p.117)

The distinction between semantic definiteness and marked definiteness is reflected as well for the indefinite, as shown in the following examples:

Non-Definite - unmarked	Non-Definite - marked	Non-definite selective
<i>sefer</i>	<i>eyzeh sefer</i>	<i>sefer eHad</i>
<i>book</i>	<i>which book</i>	<i>book one</i>
<i>a book</i>	<i>a book (undetermined)</i>	<i>a book, some book</i>

From these examples, it appears that a single binary feature *definite*, as used in SURGE, cannot account for the syntactic complexity of the definiteness. We had, therefore, to extend the input specification for definiteness.

1.2 Order of constituents

Unlike the “determiner sequence” defined in [8] Hebrew determiners and quantifiers can be positioned before and after the head of the NP¹:

(o1) <i>col ha-yeladym</i>	(o3) <i>harbe yeladym</i>	(o5) * <i>harbe ha-yeladym</i>
<i>all the-children</i>	<i>many children</i>	* <i>many the-children</i>
<i>All the children</i>	<i>Many children</i>	

(o2) <i>ha-yeladym culam</i>	(o4) <i>yeladym rabym</i>	(o6) <i>ha-yeladym ha-rabym</i>
<i>the-children all-them</i>	<i>children numerous</i>	<i>the-children the-numerous</i>
<i>All the children</i>	<i>Many children</i>	<i>The many children</i>

Determiners and quantifiers positioned after the head behave more like adjectives: they agree in number, gender and definiteness with the head. Determiners with the same meaning can appear before or after the head, as shown in examples o1/o2 and o3/o4 above. In some cases, only one of the two options is available, as shown in examples o5/o6: because the determiner for *many* appearing before the head is lexically marked as indefinite, it cannot be used for the NP *the many children*.

In addition to the choice before/after head, when several determiners and quantifiers co-occur, they must be ordered in the right sequence. The distinction between determiners (D) and quantifiers (Q) has been proposed [3] to explain, among other facts, the order generally used:

<i>col me'ot ha-mafginym</i>	* <i>me'ot col ha-mafginym</i>
<i>all hundreds the-demonstrators</i>	* <i>hundreds all the-demonstrators</i>
<i>All of the hundreds of demonstrators</i>	* <i>Hundreds of all of the demonstrators</i>

“Determiners” (like *col / all*) generally appear before “quantifiers” (like *me'ot / hundreds*). By this definition, determiners are those words which, among other characteristics, can enter in recursive constructs. According to this model, there are no restrictions on the order of occurrence of the D determiners:

<i>col S'ar ha-mafginym</i>	<i>S'ar col ha-mafginym</i>
<i>all the-rest-of the-demonstrators</i>	<i>the-rest-of all the-demonstrators</i>
<i>All the rest of the demonstrators</i>	<i>The rest of all the demonstrators</i>

The different orders, naturally, have different meanings.

We have found, however, that there are cases where the order between D determiners can lead to syntactically unacceptable NPs:

<i>col ot-am ha-yeladym</i>	* <i>ot-am col ha-yeladym</i>
<i>all those-them the-children</i>	* <i>those-them all children</i>
<i>all those children</i>	* <i>those all children</i>

The structural distinction Determiner/Quantifier proposed in [3] does not help to explain this restriction. The problem with the last example is related to the functional incompatibility of the two determiners. We instead present a functional account of these restrictions, based on a set of features derived from those used in SURGE.

¹Determiners in English can occur after the head in the very restricted cases where a pronoun is modified by a “total” meaning – e.g., *we both, them all*.

2 Previous Work

Descriptive studies of the syntax of the Hebrew noun phrase do not generally distinguish between different functions of determiners but focus rather on the syntactic structure of the NP.

Glinert [7] adopts a functional perspective, more appropriate to the needs of a generation system, and identifies a general pattern of the NP that we use as a basis:

[partitive determiner amount head classifiers describers post-det/quant qualifiers]

In this pattern, the slots can be filled by any word that fits the function: partitive expressions identify a subset of the set referred to by the rest of the NP; amount expressions identify the quantity or number, exact or inexact, of the thing referred to; determiners identify the thing referred to more or less precisely or specifically. The general picture of the input specification is that of a reference set from which a subset is identified by different complementary means. The function of each word realizing this input specification determines where it appears in the syntactic structure.

Yizhar and Doron [3] [12] distinguish between two sets of determiners, that they call D and Q quantifiers. The distinction is based on syntactic features, such as location, ability to be modified, ability to participate in partitive structures, requirement to agree in number and gender with the head. This distinction is used to explain cooccurrence restrictions, the order of appearance of D vs Q quantifiers and the recursive structure of D determiners: D determiners can be layered on top of other D determiners. A single Q quantifier can occur in an NP and it remains attached closest to the head.

As discussed in the introduction, the distinction D/Q is not sufficient to explain some restrictions on the cooccurrence of several Ds in the same NP, and does not help to predict the restrictions on occurrence of determiners in different definiteness and partitive contexts.

In [1] we have, therefore, refined the D/Q classification using functional criteria: we map the Q quantifiers to the “amount” category defined by Glinert, and the D set is split into the partitive and determiner categories – each with a different function. Of these, only partitive are recursive.

Another controversial issue within Hebrew studies is whether determiners are types of nouns [2] [6]. This is an important issue to explain in a uniform manner why constructs like *smixut* (noun compounding) are used for certain determiners and why certain determiners determine the number and gender of the whole NP. We have implemented some of these constraints by defining lexical features for determiners which determine this behavior.

2.1 Determiner Sequence in Generation

Our goal is to define a set of functional features to include in the input to a generator. These features must provide enough information to construct a wide variety of determiner constructions both in Hebrew and in English. For the comparison with English, we based our work on the SURGE generation grammar [5].

The determiner sequence in SURGE is based on [8] and has the following pattern:

[pre-det (of) det deictic2 ordinal cardinal quantifier (describer classifier head)].

Where *pre-det* can be any one of *all*, *both*, *half*, multipliers and fractions, *det* is a deictic determiner, *deictic2* is an adjective from a restricted class that expresses the anaphoric status of the thing referred to (*e.g.*, *above*, *same*, *different*); and *quantifier* expresses the amount or quantity of the thing referred to.

In [4], we gathered a minimal set of 24 features to define the determiner sequence of the NP in English. Complex systems in SURGE implement the cooccurrence restrictions among these features, definiteness and partitive constructs.

Feature	Example	Feature	Example
Definite yes/no	<i>The/a book</i>	degree +	<i>The many cars</i>
Distance far/near/none	<i>This/That/The book</i>	degree -	<i>A little butter</i>
Selective yes/no	<i>Some/Ø children</i>	degree none	<i>Some butter</i>
Total +/-/none	<i>All/None/Ø of the children</i>	comparative yes	<i>More cars</i>
Cardinal	<i>The three children</i>	superlative yes	<i>The most cars</i>
Fraction	<i>One-third of the children</i>	evaluative yes	<i>Too many cars</i>
Multiplier	<i>Twice his weight</i>	orientation -	<i>Few cars</i>

Table 1: Features used by the determiner system of the SURGE English Generator

When studying the Hebrew determiner system, we found that this set of features had to be extended and refined. As a consequence, both the implementations of SURGE and of HUGG have been made compatible with the same set of features, which we motivate in the rest of the paper.

3 Definiteness

Definiteness in NP has two aspects: it corresponds to a complex semantic property of the thing referred to by the phrase and it is syntactically marked in different ways. We will not discuss how a noun phrase is defined as semantically definite (a development on this issue is provided in [1] Chap.5), and we assume that the decision that a referent is definite is taken by the content determination module and that the semantic definite feature is given in the input to the generation grammar. We keep the feature *definite yes* for this purpose.

In a Hebrew definite NP, the definite marker (ha-) is agglutinated to the head and to most modifiers (adjectives and quantifiers). Not all subconstituents of the NP are marked, however: in a noun-compound construction the head loses the marking and the compounded noun keeps it. When demonstratives are used, a marked and an unmarked form can be used. In addition, pronouns and proper nouns are intrinsically definite and do not receive the ha- marking.

To explain the distribution of the ha- marker, we have added the feature *mark-definite* which determines whether an article (definite or indefinite) is to be used for each subconstituent of the NP.

By default, the *mark-definite* of a subconstituent is equal to the *definite* of its mother constituent. However, certain constructions and certain lexical items block this propagation and can force a *mark-definite no* (for example, proper nouns have a lexical feature of *mark-definite no* and the noun-compounding construction enforces a *mark-definite no* on the head).

The following example illustrates how the feature *mark-definite* contrasts among the two distinct phrases:

Mark-definite no: (anaphoric)	$\left[\begin{array}{l} \text{cat} \\ \text{lex} \\ \text{definite} \\ \text{distance} \\ \text{classifier} \\ \\ \text{describer} \\ \text{mark-definite} \end{array} \left[\begin{array}{l} \text{common} \\ \text{"nawar"} \\ \text{yes} \\ \text{near} \\ \left[\begin{array}{l} \text{lex} \quad \text{"cfar"} \\ \text{lex} \quad \text{"Tov"} \end{array} \right] \\ \left[\begin{array}{l} \text{classifier} \\ \left[\begin{array}{l} \text{cat} \quad \text{common} \\ \text{lex} \quad \text{"lev"} \end{array} \right] \end{array} \right] \\ \text{yes no} \end{array} \right]$
<i>nawar cfar Tov lev ze</i>	
<i>boy village good heart this</i>	
<i>This good-hearted country boy</i>	
Mark-definite yes: (deictic)	
<i>nawar ha-cfar Tov ha-lev ha-ze</i>	
<i>boy the-village good the-heart the-this</i>	

In this example, *nawar/boy* is not marked with ha- because it enters in a compound construction with *cfar*. Similarly *Tov/good* is not marked because it is an adjective compound of *lev/heart*.

While there is agreement as to the behavior of the marking of definiteness, the marking of indefinite nouns is quite controversial. Some researchers claim that definite is not a binary but that an additional unmarked value is necessary (cf [2] for a review). We mainly adopt Ornan's view [10].

In this view, an indefinite NP is marked with a null article by default, but the indefinite can also be marked. In this case, we use the feature (mark-definite indef). This feature explains the contrast between:

<i>eyze delet</i>	<i>some door</i>	mark-definite	indef
<i>delet</i>	<i>door</i>	mark-definite	no
<i>delet aHat</i>	<i>door one</i>	selective	yes

We treat *delet aHat* as marked functionally as a selective phrase, following Halliday's feature as used in SURGE. In our analysis, *selective yes* implies *mark-definite indef*. Of the two marked indefinite articles, *eyze* is non selective while *eHad* is.

The following example, taken from an example by Winograd, illustrates the behavior of the *selective yes* feature:

Examining the cabinet, we noticed that a door was marred

The Hebrew for it will be:

? *cSe-badaknu et ha-aron, ra'ynu S-delet Svura*
while-checked-us OM the-closet, saw-us that-door broken

cSebadaknu et ha-aron, ra'ynu S-delet aHat Svura
while-checked-us OM the-closet, saw-us that-door one broken

In this context, the existence of a set of 2 or 4 doors is known pragmatically (from world knowledge about closets). The function of the indefinite article is then to select one door out of this known set. In that case, the unmarked indefinite is quite inappropriate, because it would refer to a door that cannot be related to the closet.

In summary, the complex marking of definiteness in Hebrew has led us to extend the single semantic feature *definite yes/no* used in SURGE, with the additional feature *mark-definite yes/no/indef*.

4 Determining the Order

We have mostly investigated the order of determiners and quantifiers in Hebrew that appear before the head of the NP. We have found that semantics plays a main role in determining the relative order of determiners. This follows Glinert ([7] Chap.8 and 9), who distinguishes among the different functions fulfilled by determiners, and predicts the ordering accordingly, as discussed in Sect.2.

The primary functional distinction is between *quantifiers* – that express amount or portion – and *determiners* – that express identity. Accordingly, the word *col/all* for example, has two senses: one as a quantifier (with the *all* or collective meaning) and one as a determiner (with the *every* or distributive meaning).

Quantifiers are classified along two main dimensions: exact/inexact (*e.g., three vs. many*) and portion/amount. The portion/amount system has the most effect on the structure of the NP, as it can be realized either lexically by certain quantifiers, or by using an explicit partitive syntactic structure *X m-Y/X of Y*.

We define amount quantifiers (cardinals, quantifiers like *harbeh/many, yoter/more*) as quantifiers that express quantity without explicit relation to the “total” cardinality of the reference set²

²Although amount quantifiers are always semantically related to some reference set, this relation is presented

In contrast, partitive quantifiers express a quantity that gets its meaning only in relation to the cardinality of a reference set. For example, *rov/most* refers to at least half of the reference set. Similarly, *col/all* refers to all of the reference set, and is also partitive.

The distinction among partitive quantifiers, amount quantifiers and determiners is important as it predicts the order of the words in the NP. The standard order is:
[partitive determiner amount head]

We also note that only partitives can enter into recursive structures.

Accordingly, our input specification language enforces the constraint that only a single amount and a single identification feature can be present simultaneously. The realization grammar also uses the knowledge of which word realizes which function to determine the ordering. This is illustrated in the following example, where the standard order is obtained for the combination *partitive amount head*:

<i>col waseret ha-mafginym</i>	[cat	np]	
<i>all ten the-demonstrators</i>		total	+		
<i>All the ten demonstrators</i>		ref-set	lex		"mafgin"
			definite		yes
		cardinal	[value 10]		

Note that whenever a partitive quantifier is desired, the input specification must include a ref-set construct. This enforces the constraint that partitives yield recursive constructs, as shown in the following example:

<i>wasarah me-col ha-mafginym</i>	[cat	np]	
<i>ten of-all the-demonstrators</i>		cardinal	[value 10]		
<i>Ten of all the demonstrators</i>		ref-set	total		+
			ref-set		[lex "mafgin"]

The decision to build an explicitly partitive construction *X m-Y/X of Y* is left to the realization grammar and can be quite complex. If a quantifier can be found that has the lexical property of being marked as partitive, then a non-explicitly partitive construction can be used, as in the first example above *all the ten demonstrators*. Note that in this case, in contrast to English, in Hebrew, the explicitly partitive construction is not possible * *col me-waseret ha-mafginym*. In contrast, if a feature modifying the ref-set is to be realized by a non-lexically partitive quantifier, then an explicitly partitive construction must be used.

When the "portion" part is modified with adjectives, then an explicitly partitive construction must be used:

ha-rov ha-gadol mi-beyn ha-yeladym
the-most the-big of-from the-children
The vast majority of the children

In that case, the portion functions as the logical head of the NP with the realization of the reference set as a PP – which affects agreement with the verb.

In summary, we propose a refined classification of the determiners/quantifiers in Hebrew into amount, partitive and determiner. This functional classification determines the relative order of the determiners within the NP and also determines which recursive structures can be constructed. Finally, it is important to take into account the lexical properties of specific determiners which can require or forbid the use of explicitly partitive constructions. Our input specification language encompasses all of the features found in SURGE but organizes them into a new framework which enforces these constraints (with the explicit specification of the ref-set in the input).

implicitly for amounts: *3 children* refers to 3 children out of some reference set, but the reference set is deliberately left implicit.

5 Conclusion

We have presented in this paper observations on the syntax of the Hebrew NP and deduced from its properties an input specification language for complex NPs to a generation grammar. The case study of the determiner system illustrates the general methodology we apply to develop a multi-lingual generator. The key point is to encapsulate knowledge of syntax in separate modules driven by a common set of functional features.

The input specification language extends the one used in SURGE for English to provide enough information to enforce the constraints observed in Hebrew. We have specifically extended the treatment of definiteness (and indefiniteness) with the additional feature mark-definite and the treatment of partitives with a systematic classification of quantifiers as partitives, amounts and determiners.

The resulting formalism provides a good basis for a bilingual generation grammar, and is being used to that end in practical projects.

References

- [1] Yael Dahan-Netzer. Hugg - unification-based grammar for the generation of hebrew noun phrases. Master's thesis, Ben Gurion University, Beer Sheva Israel, 1997. (In Hebrew).
- [2] G. Danon. The syntax of determiners in hebrew. Master's thesis, Tel Aviv University, May 1996.
- [3] E.Doron. The np structure. In U.Ornan E.Doron A.Ariely, editor, *Hebrew Computational Linguistics*. Ministry of Science, 1991. In Hebrew.
- [4] Michael Elhadad. *Using argumentation to control lexical choice: a unification-based implementation*. PhD thesis, Computer Science Department, Columbia University, 1992.
- [5] Michael Elhadad. Lexical choice for complex noun phrases: Structure, modifiers and determiners. *Machine Translation*, 11:159–184, 1996.
- [6] L. Glinert. Shadow noun phrases in spoken hebrew. *Hebrew Computational Linguistics*, (13), 1978. In Hebrew.
- [7] L. Glinert. *The Grammar of Modern Hebrew*. Cambridge University, 1989.
- [8] Michael A. K. Halliday. *An Introduction to Functional Grammar*. Edward Arnold, London, second edition, 1994.
- [9] William C. Mann. An overview of the Penman text generation system. In *Proceedings of the Third National Conference on Artificial Intelligence (AAAI-83)*, pages 261–265, Washington, DC, August 22-26, 1983. Also appears as USC/Information Sciences Institute Tech Report RR-83-114.
- [10] Uzi Ornan. *The Nominal Phrase in Modern Hebrew*. PhD thesis, Hebrew University, Jerusalem, 1964. (in Hebrew).
- [11] H. Rosen. *Contemporary Hebrew*. Mouton, 1977.
- [12] D. Yzhar. Computational grammar for noun phrases in hebrew. Master's thesis, Hebrew University, Jerusalem, 1993. In Hebrew.