

# Natural Logic and Natural Language Inference

Bill MacCartney and Christopher D. Manning

**Abstract** We propose a model of natural language inference which identifies valid inferences by their lexical and syntactic features, without full semantic interpretation. We extend past work in *natural logic*, which has focused on semantic containment and monotonicity, by incorporating both semantic exclusion and implicativity. Our model decomposes an inference problem into a sequence of atomic edits linking premise to hypothesis; predicts a lexical entailment relation for each edit; propagates these relations upward through a semantic composition tree according to properties of intermediate nodes; and joins the resulting entailment relations across the edit sequence. A computational implementation of the model achieves 70 % accuracy and 89 % precision on the FraCaS test suite. Moreover, including this model as a component in an existing system yields significant performance gains on the Recognizing Textual Entailment challenge.

## 1 Introduction

Natural language inference (NLI) is the problem of determining whether a natural language hypothesis  $h$  can reasonably be inferred from a given premise  $p$ . For example:

- (1)  $p$ : *Every firm polled saw costs grow more than expected, even after adjusting for inflation.*  
 $h$ : *Every big company in the poll reported cost increases.*

A capacity for open-domain NLI is clearly necessary for full natural language understanding, and NLI can also enable more immediate applications, such as semantic search and question answering. Consequently, NLI has been the focus of intense research effort in recent years, centered around the annual Recognizing Textual Entailment (RTE) competition (Dagan et al. 2006).

---

B. MacCartney (✉) · C.D. Manning  
Stanford University, Stanford, CA, USA  
e-mail: [wcmac@cs.stanford.edu](mailto:wcmac@cs.stanford.edu)

C.D. Manning  
e-mail: [manning@cs.stanford.edu](mailto:manning@cs.stanford.edu)

For a semanticist, the most obvious approach to NLI relies on full semantic interpretation: first, translate  $p$  and  $h$  into some formal meaning representation, such as first-order logic (FOL), and then apply automated reasoning tools to determine inferential validity. While the formal approach can succeed in restricted domains, it struggles with open-domain NLI tasks such as RTE. For example, the FOL-based system of Bos and Markert (2005) was able to find a proof for less than 4 % of the problems in the RTE1 test set. The difficulty is plain: truly *natural* language is fiendishly complex. The formal approach faces countless thorny problems: idioms, ellipsis, paraphrase, ambiguity, vagueness, lexical semantics, the impact of pragmatics, and so on. Consider for a moment the difficulty of fully and accurately translating example (1) to a formal meaning representation.

Yet example (1) also demonstrates that full semantic interpretation is often not necessary to determining inferential validity. To date, the most successful NLI systems have relied on surface representations and approximate measures of lexical and syntactic similarity to ascertain whether  $p$  subsumes  $h$  (Glickman et al. 2005; MacCartney et al. 2006; Hickl et al. 2006). However, these approaches face a different problem: they lack the precision needed to properly handle such commonplace phenomena as negation, antonymy, downward-monotone quantifiers, non-factive contexts, and the like. For example, if *every* were replaced by *some* or *most* throughout (1), the lexical and syntactic similarity of  $h$  to  $p$  would be unaffected, yet the inference would be rendered invalid.

In this paper, we explore a middle way, by developing a model of what Lakoff (1970) called *natural logic*, which characterizes valid patterns of inference in terms of syntactic forms which are as close as possible to surface forms. For example, the natural logic approach might sanction (1) by observing that: in ordinary *upward monotone* contexts, deleting modifiers preserves truth; in *downward monotone* contexts, inserting modifiers preserves truth; and *every* is downward monotone in its restrictor NP. Natural logic thus achieves the semantic precision needed to handle inferences like (1), while sidestepping the difficulties of full semantic interpretation.

The natural logic approach has a very long history,<sup>1</sup> originating in the syllogisms of Aristotle (which can be seen as patterns for natural language inference) and continuing through the medieval scholastics and the work of Leibniz. It was revived in recent times by van Benthem (1988, 1991) and Sánchez Valencia (1991), whose *monotonicity calculus* explains inferences involving semantic containment and inversions of monotonicity, even when nested, as in *Nobody can enter without a valid passport*  $\models$  *Nobody can enter without a passport*. However, because the monotonicity calculus lacks any representation of semantic exclusion, it fails to license many simple inferences, such as *Stimpy is a cat*  $\models$  *Stimpy is not a poodle*.

Another model which arguably belongs to the natural logic tradition (though not presented as such) was developed by Nairn et al. (2006) to explain inferences involving implicatives and factives, even when negated or nested, as in *Ed did not forget to force Dave to leave*  $\models$  *Dave left*. While the model bears some resemblance

---

<sup>1</sup>For a useful overview of the history of natural logic, see van Benthem (2008). For recent work on theoretical aspects of natural logic, see (Fyodorov et al. 2000; Sukkarieh 2001; van Eijck 2005).

to the monotonicity calculus, it does not incorporate semantic containment or explain interactions between implicatives and monotonicity, and thus fails to license inferences such as *John refused to dance*  $\models$  *John didn't tango*.

We propose a new model of natural logic which extends the monotonicity calculus to incorporate semantic exclusion, and partly unifies it with Nairn et al.'s account of implicatives. We first define an inventory of *basic entailment relations* which includes representations of both containment and exclusion (Sect. 2). We then describe a general method for establishing the entailment relation between a premise  $p$  and a hypothesis  $h$ . Given a sequence of *atomic edits* which transforms  $p$  into  $h$ , we determine the *lexical entailment relation* generated by each edit (Sect. 4); project each lexical entailment relation into an *atomic entailment relation*, according to properties of the context in which the edit occurs (Sect. 5); and join atomic entailment relations across the edit sequence (Sect. 3). We have previously presented an implemented system based on this model (MacCartney and Manning 2008); here we offer a detailed account of its theoretical foundations.

## 2 An Inventory of Entailment Relations

The simplest formulation of the NLI task is as a binary decision problem: the relation between  $p$  and  $h$  is to be classified as either *entailment* ( $p \models h$ ) or *non-entailment* ( $p \not\models h$ ). The *three-way* formulation refines this by dividing non-entailment into *contradiction* ( $p \models \neg h$ ) and *compatibility* ( $p \not\models h \wedge p \not\models \neg h$ ).<sup>2</sup> The monotonicity calculus carves things up differently: it interprets entailment as a *semantic containment* relation  $\sqsubseteq$  analogous to the set containment relation  $\subseteq$ , and thus permits us to distinguish forward entailment ( $p \sqsubseteq h$ ) from reverse entailment ( $p \supseteq h$ ). Moreover, it defines  $\sqsubseteq$  for expressions of every semantic type, including not only complete sentences but also individual words and phrases. Unlike the three-way formulation, however, it lacks any way to represent contradiction (semantic exclusion). For our model, we want the best of both worlds: a comprehensive inventory of entailment relations that includes representations of both semantic containment and semantic exclusion.

Following Sánchez Valencia, we proceed by analogy with set relations. In a universe  $U$ , the set of ordered pairs  $\langle x, y \rangle$  of subsets of  $U$  can be partitioned into 16 equivalence classes, according to whether each of the four sets  $x \cap y$ ,  $x \cap \bar{y}$ ,  $\bar{x} \cap y$ , and  $\bar{x} \cap \bar{y}$  is empty or non-empty.<sup>3</sup> Of these 16 classes, nine represent degenerate cases in which either  $x$  or  $y$  is either empty or universal. Since expressions having empty denotations (e.g., *round square cupola*) or universal denotations (e.g., *exists*) fail to divide the world into meaningful categories, they can be regarded as semantically vacuous. Contradictions and tautologies may be common in

<sup>2</sup>The first three RTE competitions used the binary formulation, while the three-way formulation was adopted for RTE4. The three-way formulation was also employed in the FraCaS test suite (Cooper et al. 1996) and has been investigated in depth by Condoravdi et al. (2003).

<sup>3</sup>We use  $\bar{x}$  to denote the complement of set  $x$  in universe  $U$ ; thus  $x \cap \bar{x} = \emptyset$  and  $x \cup \bar{x} = U$ .

**Table 1** The set  $\mathfrak{B}$  of seven basic entailment relations

Symbol <sup>a</sup>	Name	Example	Set theoretic definition <sup>b</sup>
$x \equiv y$	equivalence	<i>couch</i> $\equiv$ <i>sofa</i>	$x = y$
$x \sqsubset y$	forward entailment	<i>crow</i> $\sqsubset$ <i>bird</i>	$x \subset y$
$x \sqsupset y$	reverse entailment	<i>European</i> $\sqsupset$ <i>French</i>	$x \supset y$
$x \wedge y$	negation	<i>human</i> $\wedge$ <i>nonhuman</i>	$x \cap y = \emptyset \wedge x \cup y = U$
$x \mid y$	alternation	<i>cat</i> $\mid$ <i>dog</i>	$x \cap y = \emptyset \wedge x \cup y \neq U$
$x \smile y$	cover	<i>animal</i> $\smile$ <i>nonhuman</i>	$x \cap y \neq \emptyset \wedge x \cup y = U$
$x \# y$	independence	<i>hungry</i> $\#$ <i>hippo</i>	(all other cases)

<sup>a</sup>Selecting an appropriate symbol to represent each relation is a vexed problem. We sought symbols which (a) are easily approximated by a single ASCII character, (b) are graphically symmetric iff the relations they represent are symmetric, and (c) do not excessively abuse accepted conventions. The  $\wedge$  symbol was chosen to evoke the logically similar bitwise XOR operator of the C programming language family; regrettably, it may also evoke the Boolean AND function. The  $\mid$  symbol was chosen to evoke the Sheffer stroke commonly used to represent the logically similar Boolean NAND function; regrettably, it may also evoke the Boolean OR function. The  $\sqsubset$  and  $\sqsupset$  symbols were obviously chosen to resemble their set-theoretic analogs, but a potential confusion arises because some logicians use the horseshoe  $\supset$  (with the *opposite* orientation) to represent material implication

<sup>b</sup>Each relation in  $\mathfrak{B}$  obeys the additional constraints that  $\emptyset \subset x \subset U$  and  $\emptyset \subset y \subset U$  (i.e.,  $x$  and  $y$  are non-vacuous)

logic textbooks, but they are rare in everyday speech. Thus, in a practical model of informal natural language inference, we will rarely go wrong by assuming the *non-vacuity* of the expressions we encounter.<sup>4</sup> We therefore focus on the remaining seven classes, which we designate as the set  $\mathfrak{B}$  of *basic entailment relations*, shown in Table 1.

First, the semantic containment relations ( $\sqsubset$  and  $\sqsupset$ ) of the monotonicity calculus are preserved, but are factored into three mutually exclusive relations: equivalence ( $\equiv$ ), (strict) forward entailment ( $\sqsubset$ ), and (strict) reverse entailment ( $\sqsupset$ ). Next, we have two relations expressing semantic exclusion: negation ( $\wedge$ ), or exhaustive exclusion, which is analogous to set complement; and alternation ( $\mid$ ), or non-exhaustive exclusion. The next relation is cover ( $\smile$ ), or non-exclusive exhaustion. Though its utility is not immediately obvious, it is the *dual under negation* of the alternation relation.<sup>5</sup> Finally, the independence relation ( $\#$ ) covers all other cases: it expresses non-equivalence, non-containment, non-exclusion, and non-exhaustion. Note that  $\#$

<sup>4</sup>Our model can easily be revised to accommodate vacuous expressions and relations between them, but then becomes somewhat unwieldy. The assumption of non-vacuity is closely related to the assumption of *existential import* in traditional logic. For a defense of existential import in natural language semantics, see (Böttner 1988).

<sup>5</sup>We describe relations  $R$  and  $S$  as *duals under negation* iff  $\forall x, y : \langle x, y \rangle \in R \Leftrightarrow \langle \bar{x}, \bar{y} \rangle \in S$ . Thus  $\sqsubset$  and  $\sqsupset$  are dual;  $\mid$  and  $\smile$  are dual; and  $\equiv$ ,  $\wedge$ , and  $\#$  are self-dual. The significance of this duality will become apparent in Sect. 5.

is the least informative relation, in that it places the fewest constraints on its arguments.<sup>6</sup>

Following Sánchez Valencia, we define the relations in  $\mathfrak{B}$  for all semantic types. For semantic types which can be interpreted as characteristic functions of sets,<sup>7</sup> the set-theoretic definitions can be applied directly. The definitions can then be extended to other types by interpreting each type as if it were a type of set. For example, propositions can be understood (per Montague) as denoting sets of possible worlds. Thus two propositions stand in the  $|$  relation iff there is no world where both hold (but there is some world where neither holds). Likewise, names can be interpreted as denoting singleton sets, with the result that two names stand in the  $\equiv$  relation iff they refer to the same entity, or the  $|$  relation otherwise.

By design, the relations in  $\mathfrak{B}$  are mutually exclusive, so that we can define a function  $\beta(x, y)$  which maps every ordered pair of expressions<sup>8</sup> to the unique relation in  $\mathfrak{B}$  to which it belongs.

### 3 Joining Entailment Relations

If we know that entailment relation  $R$  holds between  $x$  and  $y$ , and that entailment relation  $S$  holds between  $y$  and  $z$ , then what is the entailment relation between  $x$  and  $z$ ? The *join* of entailment relations  $R$  and  $S$ , which we denote  $R \bowtie S$ ,<sup>9</sup> is defined by:

$$R \bowtie S \stackrel{\text{def}}{=} \{ \langle x, z \rangle : \exists y (\langle x, y \rangle \in R \wedge \langle y, z \rangle \in S) \}$$

Some joins are quite intuitive. For example, it is immediately clear that  $\sqsubset \bowtie \sqsubset = \sqsubset$ ,  $\sqsubset \bowtie \sqsupset = \sqsupset$ ,  $\wedge \bowtie \wedge = \equiv$ , and for any  $R$ ,  $(R \bowtie \equiv) = (\equiv \bowtie R) = R$ . Other joins are less obvious, but still accessible to intuition. For example,  $| \bowtie \wedge = \sqsubset$ . This can be seen with the aid of Venn diagrams, or by considering simple examples: *fish*  $|$  *human* and *human*  $\wedge$  *nonhuman*, thus *fish*  $\sqsubset$  *nonhuman*.

But we soon stumble upon an inconvenient truth: not every join yields a relation in  $\mathfrak{B}$ . For example, if  $x | y$  and  $y | z$ , the relation between  $x$  and  $z$  is not determined. They could be equivalent, or one might contain the other. They might be independent

---

<sup>6</sup>Two sets selected uniformly at random from  $2^U$  are overwhelmingly likely to belong to  $\#$  (for large  $|U|$ ).

<sup>7</sup>That is, all functional types whose final output is a truth value. If we assume a type system whose basic types are  $e$  (entities) and  $t$  (truth values), then this includes most of the functional types encountered in semantic analysis:  $e \rightarrow t$  (common nouns, adjectives, and intransitive verbs),  $e \rightarrow e \rightarrow t$  (transitive verbs),  $(e \rightarrow t) \rightarrow (e \rightarrow t)$  (adverbs),  $(e \rightarrow t) \rightarrow (e \rightarrow t) \rightarrow t$  (binary generalized quantifiers), and so on.

<sup>8</sup>Assuming the expressions are non-vacuous, and belong to the same semantic type.

<sup>9</sup>In Tarskian relation algebra, this operation is known as *relation composition*, and is often represented by a semi-colon:  $R ; S$ . To avoid confusion with semantic composition (Sect. 5), we prefer to use the term *join* for this operation, by analogy to the database JOIN operation (also commonly represented by  $\bowtie$ ).

**Table 2** The join table for the basic entailment relations

$\bowtie$	$\equiv$	$\sqsubset$	$\sqsupset$	$\wedge$	$\mid$	$\smile$	$\#$
$\equiv$	$\equiv$	$\sqsubset$	$\sqsupset$	$\wedge$	$\mid$	$\smile$	$\#$
$\sqsubset$	$\sqsubset$	$\sqsubset$	$\equiv\sqsubset\mid\#$	$\mid$	$\mid$	$\sqsubset\wedge\smile\#$	$\sqsubset\mid\#$
$\sqsupset$	$\sqsupset$	$\equiv\sqsubset\sqsupset\smile\#$	$\sqsupset$	$\smile$	$\sqsupset\wedge\smile\#$	$\smile$	$\sqsupset\smile\#$
$\wedge$	$\wedge$	$\smile$	$\mid$	$\equiv$	$\sqsubset$	$\sqsubset$	$\#$
$\mid$	$\mid$	$\sqsubset\wedge\smile\#$	$\mid$	$\sqsubset$	$\equiv\sqsubset\sqsupset\mid\#$	$\sqsubset$	$\sqsubset\mid\#$
$\smile$	$\smile$	$\smile$	$\sqsupset\wedge\smile\#$	$\sqsupset$	$\sqsupset$	$\equiv\sqsubset\sqsupset\smile\#$	$\sqsupset\smile\#$
$\#$	$\#$	$\sqsubset\smile\#$	$\sqsupset\mid\#$	$\#$	$\sqsupset\mid\#$	$\sqsubset\smile\#$	$\bullet$

or alternative. All we can say for sure is that they are not exhaustive (since both are disjoint from  $y$ ). Thus, the result of joining  $\mid$  and  $\mid$  is not a relation in  $\mathfrak{B}$ , but a *union* of such relations, specifically  $\bigcup\{\equiv, \sqsubset, \sqsupset, \mid, \#\}$ .<sup>10</sup>

We will refer to (non-trivial) unions of relations in  $\mathfrak{B}$  as *union relations*.<sup>11</sup> Of the 49 possible joins of relations in  $\mathfrak{B}$ , 32 yield a relation in  $\mathfrak{B}$ , while 17 yield a union relation, with larger unions conveying less information. Union relations can be further joined, and we can establish that the smallest set of relations which contains  $\mathfrak{B}$  and is closed under joining contains just 16 relations.<sup>12</sup> One of these is the total relation, which contains all pairs of (non-vacuous) expressions. This relation, which we denote  $\bullet$ , is the black hole of entailment relations, in the sense that (a) it conveys zero information about pairs of expressions which belong to it, and (b) joining a chain of entailment relations will, if it contains any noise and is of sufficient length, lead inescapably to  $\bullet$ .<sup>13</sup> This tendency of joining to devolve toward less-informative entailment relations places an important limitation on the power of the inference method described in Sect. 7.

A complete join table for relations in  $\mathfrak{B}$  is shown in Table 2.<sup>14</sup>

In an implemented model, the complexity introduced by union relations is easily tamed. Every union relation which results from joining relations in  $\mathfrak{B}$  contains  $\#$ , and thus can safely be approximated by  $\#$ . After all,  $\#$  is already the least informative relation in  $\mathfrak{B}$ —loosely speaking, it indicates ignorance of the relationship between two expressions—and further joining will never serve to strengthen it. Our implemented model therefore has no need to represent union relations.

<sup>10</sup>We use this notation as shorthand for the union  $\equiv \cup \sqsubset \cup \sqsupset \cup \mid \cup \#$ . To be precise, the result of this join is not identical with this union, but is a subset of it, since the union contains some pairs of sets (e.g.  $\langle U \setminus a, U \setminus a \rangle$ , for any  $|a| = 1$ ) which cannot participate in the  $\mid$  relation. However, the approximation makes little practical difference.

<sup>11</sup>Some union relations hold intrinsic interest. For example, in the three-way formulation of the NLI task described in Sect. 2, the three classes can be identified as  $\bigcup\{\equiv, \sqsubset\}$ ,  $\bigcup\{\wedge, \mid\}$ , and  $\bigcup\{\sqsupset, \smile, \#\}$ .

<sup>12</sup>That is, the relations in  $\mathfrak{B}$  plus 9 union relations. Note that this closure fails to include most of the 120 possible union relations. Perhaps surprisingly, the unions  $\bigcup\{\equiv, \sqsubset\}$  and  $\bigcup\{\wedge, \mid\}$  mentioned in footnote 11 do not appear.

<sup>13</sup>In fact, computer experiments show that if relations are selected uniformly at random from  $\mathfrak{B}$ , it requires on average just five joins to reach  $\bullet$ .

<sup>14</sup>For compactness, we omit the union notation here; thus  $\sqsubset\mid\#$  stands for  $\bigcup\{\sqsubset, \mid, \#\}$ .

## 4 Lexical Entailment Relations

Suppose  $x$  is a compound linguistic expression, and let  $e(x)$  be the result of applying an *atomic edit*  $e$  (the deletion, insertion, or substitution of a subexpression) to  $x$ . The entailment relation which holds between  $x$  and  $e(x)$ , which we denote  $\beta(x, e(x))$ , will depend on (1) the *lexical entailment relation* generated by  $e$ , which we label  $\beta(e)$ , and (2) other properties of the context  $x$  in which  $e$  is applied (to be discussed in Sect. 5). For example, suppose  $x$  is *red car*. If  $e$  is  $\text{SUB}(\text{car}, \text{convertible})$ , then  $\beta(e)$  is  $\sqsubset$  (because *convertible* is a hyponym of *car*). On the other hand, if  $e$  is  $\text{DEL}(\text{red})$ , then  $\beta(e)$  is  $\sqsubset$  (because *red* is an intersective modifier). Crucially,  $\beta(e)$  depends solely on the lexical items involved in  $e$ , independent of context.

How are lexical entailment relations determined? Ultimately, this is the province of lexical semantics, which lies outside the scope of this work. However, the answers are fairly intuitive in most cases, and we can make a number of useful observations.

**Substitutions** The entailment relation generated by a substitution edit is simply the relation between the substituted terms:  $\beta(\text{SUB}(x, y)) = \beta(x, y)$ . For open-class terms such as nouns, adjectives, and verbs, we can often determine the appropriate relation by consulting a lexical resource such as WordNet. Synonyms belong to the  $\equiv$  relation (*sofa*  $\equiv$  *couch*, *forbid*  $\equiv$  *prohibit*); hyponym-hypernym pairs belong to the  $\sqsubset$  relation (*crow*  $\sqsubset$  *bird*, *frigid*  $\sqsubset$  *cold*, *soar*  $\sqsubset$  *rise*); and antonyms and coordinate terms generally belong to the  $|$  relation (*hot*  $|$  *cold*, *cat*  $|$  *dog*).<sup>15</sup> Proper nouns, which denote individual entities or events, will stand in the  $\equiv$  relation if they denote the same entity (*USA*  $\equiv$  *United States*), or the  $|$  relation otherwise (*JFK*  $|$  *FDR*). Pairs which cannot reliably be assigned to another entailment relation will be assigned to the  $\#$  relation (*hungry*  $\#$  *hippo*). Of course, there are many difficult cases, where the most appropriate relation will depend on subjective judgments about word sense, topical context, and so on—consider, for example, the pair *system* and *approach*. And some judgments may depend on world knowledge not readily available to an automatic system. For example, plausibly *skiing*  $|$  *sleeping*, but *skiing*  $\#$  *talking*.

Closed-class terms may require special handling. Substitutions involving generalized quantifiers generate a rich variety of entailment relations: *all*  $\equiv$  *every*, *every*  $\sqsubset$  *some*, *some*  $\wedge$  *no*, *no*  $|$  *every*, *at least four*  $\sim$  *at most six*, and *most*  $\#$  *ten or more*.<sup>16</sup> Two pronouns, or a pronoun and a noun, should ideally be assigned to the  $\equiv$  relation if it can be determined from context that they refer to the same entity, though this may be difficult for an automatic system to establish reliably. Prepositions are somewhat problematic. Some pairs of prepositions can be interpreted as antonyms, and thus

<sup>15</sup>Note that most antonym pairs do *not* belong to the  $\wedge$  relation, since they typically do not exclude the middle.

<sup>16</sup>Some of these assertions assume the non-vacuity (Sect. 2) of the predicates to which the quantifiers are applied.

assigned to the  $|$  relation (*above*  $|$  *below*), but many prepositions are used so flexibly in natural language that they are best assigned to the  $\equiv$  relation (*on [a plane]*  $\equiv$  *in [a plane]*  $\equiv$  *by [plane]*).

**Generic Deletions and Insertions** For deletion edits, the default behavior is to generate the  $\sqsubset$  relation (thus *red car*  $\sqsubset$  *car*). Insertion edits are symmetric: by default, they generate the  $\sqsupset$  relation (*sing*  $\sqsupset$  *sing off-key*). This heuristic can safely be applied whenever the affected phrase is an intersective modifier, and can usefully be applied to phrases much longer than a single word (*car which has been parked outside since last week*  $\sqsubset$  *car*). Indeed, this principle underlies most current approaches the RTE task, in which the premise  $p$  often contains much extraneous content not found in the hypothesis  $h$ . Most RTE systems try to determine whether  $p$  subsumes  $h$ : they penalize new content inserted into  $h$ , but do not penalize content deleted from  $p$ .

**Special Deletions and Insertions** However, some lexical items exhibit special behavior upon deletion or insertion. The most obvious example is negation, which generates the  $\wedge$  relation (*didn't sleep*  $\wedge$  *did sleep*). Implicatives and factives (such as *refuse to* and *admit that*) constitute another important class of exceptions, but we postpone discussion of them to Sect. 6. Then there are non-intersective adjectives such as *former* and *alleged*. These have various behavior: deleting *former* seems to generate the  $|$  relation (*former student*  $|$  *student*), while deleting *alleged* seems to generate the  $\#$  relation (*alleged spy*  $\#$  *spy*). We lack a complete typology of such cases, but consider this an interesting problem for lexical semantics. Finally, for pragmatic reasons, we typically assume that auxiliary verbs and punctuation marks are semantically vacuous, and thus generate the  $\equiv$  relation upon deletion or insertion. When combined with the assumption that morphology matters little in inference,<sup>17</sup> this allows us to establish, e.g., that *is sleeping*  $\equiv$  *sleeps* and *did sleep*  $\equiv$  *slept*.

## 5 Entailment Relations and Semantic Composition

How are entailment relations affected by semantic composition? In other words, how do the entailment relations between compound expressions depend on the entailment relations between their parts? Say we have established the value of  $\beta(x, y)$ , and let  $f$  be an expression which can take  $x$  or  $y$  as an argument. What is the value of  $\beta(f(x), f(y))$ , and how does it depend on the properties of  $f$ ?

The monotonicity calculus of Sánchez Valencia provides a partial answer. It explains the impact of semantic composition on entailment relations  $\equiv$ ,  $\sqsubset$ ,  $\sqsupset$ , and  $\#$  by assigning semantic functions to one of three monotonicity classes: UP, DOWN, and NON. If  $f$  has monotonicity UP (the default), then the entailment relation between  $x$

<sup>17</sup>Indeed, the official definition of the RTE task explicitly specifies that tense be ignored.

and  $y$  is projected through  $f$  without change:  $\beta(f(x), f(y)) = \beta(x, y)$ . Thus *some parrots talk*  $\sqsubset$  *some birds talk*. If  $f$  has monotonicity DOWN, then  $\sqsubset$  and  $\sqsupset$  are swapped. Thus *no carp talk*  $\sqsupset$  *no fish talk*. Finally, if  $f$  has monotonicity NON, then  $\sqsubset$  and  $\sqsupset$  are projected as  $\#$ . Thus *most humans talk*  $\#$  *most animals talk*.

The monotonicity calculus also provides an algorithm for computing the effect on entailment relations of multiple levels of semantic composition. Although Sánchez Valencia’s presentation of this algorithm uses a complex scheme for annotating nodes in a categorial grammar parse, the central idea can be recast in simple terms: propagate a lexical entailment relation upward through a semantic composition tree, from leaf to root, while respecting the monotonicity properties of each node along the path. Consider the sentence *Nobody can enter without pants*. A plausible semantic composition tree for this sentence could be rendered as (*nobody* (*can* ((*without pants*) *enter*))). Now consider replacing *pants* with *clothes*. We begin with the lexical entailment relation: *pants*  $\sqsubset$  *clothes*. The semantic function *without* has monotonicity DOWN, so *without pants*  $\sqsupset$  *without clothes*. Continuing up the semantic composition tree, *can* has monotonicity UP, but *nobody* has monotonicity DOWN, so we get another reversal, and find that *nobody can enter without pants*  $\sqsubset$  *nobody can enter without clothes*.

While the monotonicity calculus elegantly explains the impact of semantic composition on the containment relations (chiefly,  $\sqsubset$  and  $\sqsupset$ ), it lacks any account of the exclusion relations ( $\wedge$  and  $\mid$ , and, indirectly,  $\smile$ ). To remedy this lack, we propose to generalize the concept of monotonicity to a concept of *projectivity*. We categorize semantic functions into a number of *projectivity signatures*, which can be seen as generalizations of both the three monotonicity classes of Sánchez Valencia and the nine implication signatures of Nairn et al. (see Sect. 6). Each projectivity signature is defined by a map  $\mathfrak{B} \mapsto \mathfrak{B}$  which specifies how each entailment relation is projected by the function. (Binary functions can have different signatures for each argument.) In principle, there are up to  $7^7$  possible signatures; in practice, probably no more than a handful are realized by natural language expressions. Though we lack a complete inventory of projectivity signatures, we can describe a few important cases.

**Negation** We begin with simple negation (*not*). Like most functions, it projects  $\equiv$  and  $\#$  without change (*not happy*  $\equiv$  *not glad* and *isn’t swimming*  $\#$  *isn’t hungry*). As a downward monotone function, it swaps  $\sqsubset$  and  $\sqsupset$  (*didn’t kiss*  $\sqsupset$  *didn’t touch*). But we can also establish that it projects  $\wedge$  without change (*not human*  $\wedge$  *not nonhuman*) and swaps  $\mid$  and  $\smile$  (*not French*  $\smile$  *not German* and *not more than 4*  $\mid$  *not less than 6*). Its projectivity signature is therefore  $\{\equiv:\equiv, \sqsubset:\sqsupset, \sqsupset:\sqsubset, \wedge:\wedge, \mid:\smile, \smile:\mid, \#:\#\}$ .

**Intersective Modification** Intersective modification has monotonicity UP, but projects both  $\wedge$  and  $\mid$  as  $\mid$  (*living human*  $\mid$  *living nonhuman* and *French wine*  $\mid$  *Spanish wine*), and projects  $\smile$  as  $\#$  (*metallic pipe*  $\#$  *nonferrous pipe*). It therefore has signature  $\{\equiv:\equiv, \sqsubset:\sqsubset, \sqsupset:\sqsupset, \wedge:\mid, \mid:\mid, \smile:\#, \#:\#\}$ .<sup>18</sup>

<sup>18</sup>At least for practical purposes. The projection of  $\wedge$  and  $\mid$  as  $\mid$  depends on the assumption of non-vacuity, and  $\smile$  is actually projected as  $\bigcup\{\equiv, \sqsubset, \sqsupset, \mid, \#\}$ , which we approximate by  $\#$ , as described in Sect. 3.

**Table 3** Projectivity signatures for various quantifiers

Quantifier	Projectivity for 1st argument						Projectivity for 2nd argument							
	≡	⊂	⊃	∧		∪	#	≡	⊂	⊃	∧		∪	#
<i>some</i>	≡	⊂	⊃	∪ <sup>†</sup>	#	∪ <sup>†</sup>	#	≡	⊂	⊃	∪ <sup>†</sup>	#	∪ <sup>†</sup>	#
<i>no</i>	≡	⊃	⊂	<sup>†</sup>	#	<sup>†</sup>	#	≡	⊃	⊂	<sup>†</sup>	#	<sup>†</sup>	#
<i>every</i>	≡	⊃	⊂	<sup>‡</sup>	#	<sup>‡</sup>	#	≡	⊂	⊃	<sup>†</sup>	<sup>†</sup>	#	#
<i>not every</i>	≡	⊂	⊃	∪ <sup>‡</sup>	#	∪ <sup>‡</sup>	#	≡	⊃	⊂	∪ <sup>†</sup>	∪ <sup>†</sup>	#	#

**Quantifiers** While semanticists are well acquainted with the monotonicity properties of common quantifiers, how they project the exclusion relations may be less familiar. Table 3 summarizes the projectivity signatures of the most common binary generalized quantifiers for each argument position.

A few observations:

- All quantifiers (like most other semantic functions) project  $\equiv$  and  $\#$  without change.
- The table confirms well-known monotonicity properties: *no* is downward-monotone in both arguments, *every* in its first argument, and *not every* in its second argument.
- Relation  $|$  is frequently “blocked” by quantifiers (i.e., projected as  $\#$ ). Thus *no fish talk # no birds talk* and *someone was early # someone was late*. A notable exception is *every* in its second argument, where  $|$  is preserved: *everyone was early | everyone was late*. (Note the similarity to intersective modification.)
- Because *no* is the negation of *some*, its projectivity signature can be found by projecting the signature of *some* through the signature of *not*. Likewise for *not every* and *every*.
- Some results depend on assuming the non-vacuity of the other argument to the quantifier: those marked with <sup>†</sup> assume it to be non-empty, while those marked with <sup>‡</sup> assume it to be non-universal. Without these assumptions,  $\#$  is projected.

**Verbs** Verbs (and verb-like constructions) exhibit diverse behavior. Most verbs are upward-monotone (though not all—see Sect. 6), and many verbs project  $\wedge$ ,  $|$ , and  $\cup$  as  $\#$  (*eats humans # eats nonhumans*, *eats cats # eats dogs*, and *eats mammals # eats nonhumans*). However, verbs which encode functional relations seem to exhibit the same projectivity as intersective modifiers, projecting  $\wedge$  and  $|$  as  $|$ , and  $\cup$  as  $\#$ .<sup>19</sup> Categorizing verbs according to projectivity is an interesting problem for lexical semantics, which may involve codifying some amount of world knowledge.

<sup>19</sup>Consider the verbal construct *is married to*: *is married to a German | is married to a non-German, is married to a German | is married to an Italian, is married to a European # is married to a non-German*. The AUCONTRAIRE system (Ritter et al. 2008) includes an intriguing approach to identifying such *functional phrases* automatically.

**Table 4** Implicatives and factives

	Signature	$\beta(\text{DEL}(\cdot))$	$\beta(\text{INS}(\cdot))$	Example
implicatives (UP)	+/-	$\equiv$	$\equiv$	<i>he managed to escape <math>\equiv</math> he escaped</i>
	+/o	$\sqsubset$	$\sqsupset$	<i>he was forced to sell <math>\sqsubset</math> he sold</i>
	o/-	$\sqsupset$	$\sqsubset$	<i>he was permitted to live <math>\sqsupset</math>he lived</i>
implicatives (DOWN)	-/+	$\wedge$	$\wedge$	<i>he forgot to pay <math>\wedge</math> he paid</i>
	-/o			<i>he refused to fight   he fought</i>
	o/+	$\smile$	$\smile$	<i>he hesitated to ask <math>\smile</math> he asked</i>
factives (NON)	+/+	$\sqsubset$	$\sqsupset$	<i>he admitted that he knew <math>\sqsubset</math> he knew</i>
	-/-			<i>he pretended he was sick   he was sick</i>
	o/o	#	#	<i>he wanted to fly # he flew</i>

## 6 Implicatives and Factives

In (Nairn et al. 2006), Nairn et al. offer an elegant account of inferences involving implicatives and factives<sup>20</sup> such as *manage to*, *refuse to*, and *admit that*. Their model classifies such operators into nine *implication signatures*, according to their implications—positive (+), negative (-), or null (o)—in both positive and negative contexts. Thus *refuse to* has implication signature  $-/o$ , because it carries a negative implication in a positive context (*refused to dance* implies *didn't dance*), and no implication in a negative context (*didn't refuse to dance* implies neither *danced* nor *didn't dance*).

Most of the phenomena observed by Nairn et al. can be explained within our framework by specifying, for each implication signature, the relation generated when an operator of that signature is deleted from (or inserted into) a compound expression, as shown in Table 4.

This table invites several observations. First, as the examples make clear, there is room for variation regarding the appearance of infinitive arguments, complementizers, passivization, and morphology. An implemented model must tolerate such diversity.

Second, some of the examples may seem more intuitive when one considers their negations. For example, deleting signature  $o/-$  generates  $\sqsupset$ ; under negation, this is projected as  $\sqsubset$  (*he wasn't permitted to live  $\sqsubset$  he didn't live*). Likewise, deleting signature  $o/+$  generates  $\smile$ ; under negation, this is projected as | (*he didn't hesitate to ask | he didn't ask*).

Third, a fully satisfactory treatment of the factives (signatures  $+/+$ ,  $-/-$ , and  $o/o$ ) would require an extension to our present theory. For example, deleting signature  $+/+$  generates  $\sqsubset$ ; yet under negation, this is projected not as  $\sqsupset$ , but as | (*he*

<sup>20</sup>We use “factives” as an umbrella term embracing counterfactuals and nonfactuals along with factives proper.

*didn't admit that he knew | he didn't know*). The problem arises because the implication carried by a factive is not an entailment, but a presupposition.<sup>21</sup> As is well known, the projection behavior of presuppositions differs from that of entailments (van der Sandt 1992). It seems likely that our model could be elaborated to account for projection of presuppositions as well as entailments, but we leave this for future work.

We can further cement implicatives and factives within our model by specifying the monotonicity class for each implication signature: signatures  $+/-$ ,  $+/\circ$ , and  $\circ/-$  have monotonicity UP (*force to tango*  $\square$  *force to dance*); signatures  $-/+$ ,  $-/\circ$ , and  $\circ/+$  have monotonicity DOWN (*refuse to tango*  $\square$  *refuse to dance*); and signatures  $+/+$ ,  $-/-$ , and  $\circ/\circ$  (the propositional attitudes) have monotonicity NON (*think tangoing is fun*  $\#$  *think dancing is fun*). We are not yet able to specify the complete projectivity signature corresponding to each implication signature, but we can describe a few specific cases. For example, implication signature  $-/\circ$  seems to project  $\wedge$  as  $|$  (*refuse to stay | refuse to go*) and both  $|$  and  $\sim$  as  $\#$  (*refuse to tango # refuse to waltz*).

## 7 Putting It All Together

We now have the building blocks of a general method to establish the entailment relation between a premise  $p$  and a hypothesis  $h$ . The steps are as follows:

1. Find a sequence of atomic edits  $\langle e_1, \dots, e_n \rangle$  which transforms  $p$  into  $h$ : thus  $h = (e_n \circ \dots \circ e_1)(p)$ . For convenience, let us define  $x_0 = p$ ,  $x_n = h$ , and  $x_i = e_i(x_{i-1})$  for  $i \in [1, n]$ .
2. For each atomic edit  $e_i$ :
  - a. Determine the lexical entailment relation  $\beta(e_i)$ , as in Sect. 4.
  - b. Project  $\beta(e_i)$  upward through the semantic composition tree of expression  $x_{i-1}$  to find an *atomic entailment relation*  $\beta(x_{i-1}, x_i)$ , as in Sect. 5.
3. Join atomic entailment relations across the sequence of edits, as in Sect. 3:

$$\beta(p, h) = \beta(x_0, x_n) = \beta(x_0, e_1) \bowtie \dots \bowtie \beta(x_{i-1}, e_i) \bowtie \dots \bowtie \beta(x_{n-1}, e_n)$$

However, this inference method has several important limitations, including the need to find an appropriate edit sequence connecting  $p$  and  $h$ ;<sup>22</sup> the tendency of

<sup>21</sup>Of course, the implicatives may carry presuppositions as well (*he managed to escape*  $\rightarrow$  *it was hard to escape*), but these implications are not activated by a simple deletion, as with the factives.

<sup>22</sup>The order of edits can be significant, if one edit affects the projectivity properties of the context for another edit. In practice, we typically find that different edit orders lead to the same final result (albeit via different intermediate steps), or at worst to a result which is compatible with, though less informative than, the desired result. But in principle, edit sequences involving lexical items with unusual properties—not exhibited, so far as we are aware, by any natural language expressions—could lead to incompatible results. Thus we lack any formal guarantee of soundness.

**Table 5** An example inference involving semantic exclusion

$i$	$e_i$	$x_i = e_i(x_{i-1})$	$\beta(e_i)$	$\beta(x_{i-1}, x_i)$	$\beta(x_0, x_i)$
		<i>Stimpy is a cat</i>			
1	SUB( <i>cat</i> , <i>dog</i> )				
		<i>Stimpy is a dog</i>			
2	INS( <i>not</i> )		^	^	□
		<i>Stimpy is not a dog</i>			
3	SUB( <i>dog</i> , <i>poodle</i> )		□	□	□
		<i>Stimpy is not a poodle</i>			

the join operation toward less informative entailment relations, as described in Sect. 3; and the lack of a general mechanism for combining information from multiple premises.<sup>23</sup> Consequently, the method has less deductive power than first-order logic, and fails to sanction some fairly simple inferences, including de Morgan’s laws for quantifiers. But the method neatly explains many inferences not handled by the monotonicity calculus.

For example, while the monotonicity calculus notably fails to explain even the simplest inferences involving semantic exclusion, such examples are easily accommodated in our framework. We encountered an example of such an inference in Sect. 1: *Stimpy is a cat*  $\models$  *Stimpy is not a poodle*. Clearly, this is a valid natural language inference. To establish this using our inference method, we must begin by selecting a sequence of atomic edits which transforms the premise  $p$  into the hypothesis  $h$ . While there are several possibilities, one obvious choice is first to replace *cat* with *dog*, then to insert *not*, and finally to replace *dog* with *poodle*. An analysis of this edit sequence is shown in Table 5. In this representation (of which we will see several more examples in the following pages), we show three entailment relations associated with each edit  $e_i$ , namely:

- $\beta(e_i)$ , the lexical entailment relation generated by  $e_i$ ,
  - $\beta(x_{i-1}, x_i)$ , the atomic entailment relation which holds across  $e_i$ , and
  - $\beta(x_0, x_i)$ , the cumulative join of all atomic entailment relations up through  $e_i$ .
- This can be calculated in the table as  $\beta(x_0, x_{i-1}) \bowtie \beta(x_{i-1}, x_i)$ .

In Table 5,  $x_0$  is transformed into  $x_3$  by a sequence of three edits. First, replacing *cat* with its coordinate term *dog* generates the lexical entailment relation |. Next, inserting *not* generates ^, and | joined with ^ yields □. Finally, replacing *dog* with its hyponym *poodle* generates □. Because of the downward-monotone context created by *not*, this is projected as □, and □ joined with □ yields □. Therefore, premise  $x_0$  entails hypothesis  $x_3$ .

<sup>23</sup>However, some inferences can be enabled by auxiliary premises encoded as lexical entailment relations. For example, *men* □ *mortal* can enable the classic syllogism *Socrates is a man* □ *Socrates is mortal*.

**Table 6** An example inference involving an implicative

$i$	$e_i$	$x_i = e_i(x_{i-1})$	$\beta(e_i)$	$\beta(x_{i-1}, x_i)$	$\beta(x_0, x_i)$
<i>We were not permitted to smoke</i>					
1	DEL( <i>permitted to</i> )		$\sqsupset$	$\sqsupset$	$\sqsupset$
<i>We did not smoke</i>					
2	DEL( <i>not</i> )		$\wedge$	$\wedge$	$\sqsupset\sqsupset$
<i>We smoked</i>					
3	INS( <i>Cuban cigars</i> )		$\sqsupset$	$\sqsupset$	$\sqsupset$
<i>We smoked Cuban cigars</i>					

For an example involving an implicative, consider the inference in Table 6. Again,  $x_0$  is transformed into  $x_3$  by a sequence of three edits.<sup>24</sup> First, deleting *permitted to* generates  $\sqsupset$ , according to its implication signature; but because *not* is downward-monotone, this is projected as  $\sqsupset$ . Next, deleting *not* generates  $\wedge$ , and  $\sqsupset$  joined with  $\wedge$  yields  $\sqsupset\sqsupset$ . Finally, inserting *Cuban cigars* restricts the meaning of *smoked*, generating  $\sqsupset$ , and  $\sqsupset\sqsupset$  joined with  $\sqsupset$  yields  $\sqsupset$ . So  $x_3$  contradicts  $x_0$ .

Let’s now look at a more complex example (first presented in (MacCartney and Manning 2008)) that demonstrates the interaction of a number of aspects of the model we’ve presented. The inference is:

*p*: Jimmy Dean refused to move without blue jeans.

*h*: James Dean didn’t dance without pants.

Of course, the example is quite contrived, but it has the advantage that it compactly exhibits several phenomena of interest: semantic containment (between *move* and *dance*, and between *pants* and *jeans*); semantic exclusion (in the form of negation); an implicative (namely, *refuse to*); and nested inversions of monotonicity (created by *refuse to* and *without*). In this example, the premise *p* can be transformed into the hypothesis *h* by a sequence of seven edits, as shown in Table 7. This time we include even “light” edits yielding  $\equiv$  for the sake of completeness.

We analyze these edits as follows. The first edit simply substitutes one variant of a name for another; since both substituends denote the same entity, the edit generates the  $\equiv$  relation. The second edit deletes an implicative (*refuse to*) with implication signature  $-/o$ . As described in Sect. 6, deletions of this signature generate the  $\sqsupset$  relation, and  $\equiv$  joined with  $\sqsupset$  yields  $\sqsupset$ . The third edit inserts an auxiliary verb (*did*); since auxiliaries are more or less semantically vacuous, this generates the  $\equiv$  relation, and  $\sqsupset$  joined with  $\equiv$  yields  $\sqsupset$  again. The fourth edit inserts a negation, generating the  $\wedge$  relation. Here we encounter the first interesting join: as explained in Sect. 3,  $\sqsupset$  joined with  $\wedge$  yields  $\sqsupset$ . The fifth edit substitutes *move* with its hyponym *dance*, generating the  $\sqsupset$  relation. However, because the edit occurs within the scope of the newly-introduced negation,  $\sqsupset$  is projected as  $\sqsupset$ , and  $\sqsupset$  joined with  $\sqsupset$  yields  $\sqsupset$ . The sixth edit deletes a generic modifier (*blue*), which generates the  $\sqsupset$  relation by default. This

<sup>24</sup>We neglect edits involving auxiliaries and morphology, which simply yield the  $\equiv$  relation.

**Table 7** Analysis of a more complex inference

$i$	$e_i$	$x_i = e_i(x_{i-1})$	$\beta(e_i)$	$\beta(x_{i-1}, x_i)$	$\beta(x_0, x_i)$
		<i>Jimmy Dean refused to move without blue jeans</i>			
1	SUB( <i>Jimmy Dean, James Dean</i> )		$\equiv$	$\equiv$	$\equiv$
		<i>James Dean refused to move without blue jeans</i>			
2	DEL( <i>refused to</i> )				
		<i>James Dean moved without blue jeans</i>			
3	INS( <i>did</i> )		$\equiv$	$\equiv$	
		<i>James Dean did move without blue jeans</i>			
4	INS( <i>n't</i> )		$\wedge$	$\wedge$	$\sqsubset$
		<i>James Dean didn't move without blue jeans</i>			
5	SUB( <i>move, dance</i> )		$\sqsubset$	$\sqsubset$	$\sqsubset$
		<i>James Dean didn't dance without blue jeans</i>			
6	DEL( <i>blue</i> )		$\sqsubset$	$\sqsubset$	$\sqsubset$
		<i>James Dean didn't dance without jeans</i>			
7	SUB( <i>jeans, pants</i> )		$\sqsubset$	$\sqsubset$	$\sqsubset$
		<i>James Dean didn't dance without pants</i>			

time the edit occurs within the scope of *two* downward-monotone operators (*without* and negation), so we have two inversions of monotonicity, and  $\sqsubset$  is projected as  $\sqsubset$ . Again,  $\sqsubset$  joined with  $\sqsubset$  yields  $\sqsubset$ . Finally, the seventh edit substitutes *jeans* with its hypernym *pants*, generating the  $\sqsubset$  relation. Again, the edit occurs within the scope of two downward-monotone operators, so  $\sqsubset$  is projected as  $\sqsubset$ , and  $\sqsubset$  joined with  $\sqsubset$  yields  $\sqsubset$ . Thus  $p$  entails  $h$ .

Of course, the edit sequence shown in Table 7 is not the only sequence which can transform  $p$  into  $h$ . A different edit sequence might yield a different sequence of intermediate steps, but the same final result. Consider, for example, the edit sequence shown in Table 8. Note that the lexical entailment relation  $\beta(e_i)$  generated by each edit is the same as before. But because the edits involving downward-monotone operators (namely, INS(*n't*) and DEL(*refused to*)) now occur at different points in the edit sequence, many of the atomic entailment relations  $\beta(x_{i-1}, x_i)$  have changed, and thus the sequence of joins has changed as well. In particular, edits 3 and 4 occur within the scope of *three* downward-monotone operators (negation, *refuse*, and *without*), with the consequence that the  $\sqsubset$  relation generated by each of these lexical edits is projected as  $\sqsubset$ . Likewise, edit 5 occurs within the scope of two downward-monotone operators (negation and *refuse*), and edit 6 occurs within the scope of one downward-monotone operator (negation), so that | is projected as  $\smile$ . Nevertheless, the ultimate result is still  $\sqsubset$ .

However, it turns out not to be the case that every edit sequence which transforms  $p$  into  $h$  will yield equally satisfactory results. Consider the sequence shown in Table 9. The crucial difference in this edit sequence is that the insertion of *not*, which generates lexical entailment relation  $\wedge$ , occurs within the scope of *refuse*,

**Table 8** An alternative analysis of the inference from Table 7

$i$	$e_i$	$x_i = e_i(x_{i-1})$	$\beta(e_i)$	$\beta(x_{i-1}, x_i)$	$\beta(x_0, x_i)$
		<i>Jimmy Dean refused to move without blue jeans</i>			
1	INS( <i>did</i> )		≡	≡	≡
		<i>Jimmy Dean did refuse to move without blue jeans</i>			
2	INS( <i>n't</i> )		^	^	^
		<i>Jimmy Dean didn't refuse to move without blue jeans</i>			
3	DEL( <i>blue</i> )		⊃	⊃	
		<i>Jimmy Dean didn't refuse to move without jeans</i>			
4	SUB( <i>jeans, pants</i> )		⊃	⊃	
		<i>Jimmy Dean didn't refuse to move without pants</i>			
5	SUB( <i>move, dance</i> )		⊃	⊃	
		<i>Jimmy Dean didn't refuse to dance without pants</i>			
6	DEL( <i>refuse to</i> )			⋈	⊃
		<i>Jimmy Dean didn't dance without pants</i>			
7	SUB( <i>Jimmy, James</i> )		≡	≡	⊃
		<i>James Dean didn't dance without pants</i>			

**Table 9** A third analysis of the inference from Table 7

$i$	$e_i$	$x_i = e_i(x_{i-1})$	$\beta(e_i)$	$\beta(x_{i-1}, x_i)$	$\beta(x_0, x_i)$
		<i>Jimmy Dean refused to move without blue jeans</i>			
1	INS( <i>did</i> )		≡	≡	≡
		<i>Jimmy Dean did refuse to move without blue jeans</i>			
2	INS( <i>not</i> )		^		
		<i>Jimmy Dean did refuse not to move without blue jeans</i>			
3	DEL( <i>refuse to</i> )				≡⊃⊃ #
		<i>Jimmy Dean didn't move without blue jeans</i>			
4	DEL( <i>blue</i> )		⊃	⊃	•
		<i>Jimmy Dean didn't move without jeans</i>			
5	SUB( <i>jeans, pants</i> )		⊃	⊃	•
		<i>Jimmy Dean didn't move without pants</i>			
6	SUB( <i>move, dance</i> )		⊃	⊃	•
		<i>Jimmy Dean didn't dance without pants</i>			
7	SUB( <i>Jimmy Dean, James Dean</i> )		≡	≡	•
		<i>James Dean didn't dance without pants</i>			

so that ^ is projected as atomic entailment relation | (see Sect. 5). But the deletion of *refuse to* also produces atomic entailment relation | (see Sect. 6), and | joined with | yields a relatively uninformative union relation, namely  $\bigcup\{\equiv, \sqsupset, \sqsupset, |, \#\}$  (which

could also be described as the NON-EXHAUSTION relation). The damage has been done: further joining leads directly to the “black hole” relation  $\bullet$ , from which there is no escape. Note, however, that even for this infelicitous edit sequence, our inference method has not produced an *incorrect* answer (because the  $\bullet$  relation includes the  $\sqsubset$  relation), only an *uninformative* answer (because it includes all other relations in  $\mathfrak{B}$  as well).

Additional examples are presented in (MacCartney 2009).

## 8 Implementation and Evaluation

The model of natural logic described here has been implemented in software as the NatLog system. In previous work (MacCartney and Manning 2008), we have presented a description and evaluation of NatLog; this section summarizes the main results. NatLog faces three primary challenges:

1. *Finding an appropriate sequence of atomic edits connecting premise and hypothesis.* NatLog does not address this problem directly, but relies instead on edit sequences from other sources. We have investigated this problem separately in (MacCartney et al. 2008).
2. *Determining the lexical entailment relation for each edit.* NatLog learns to predict lexical entailment relations by using machine learning techniques and exploiting a variety of manually and automatically constructed sources of information on lexical relations.
3. *Computing the projection of each lexical entailment relation.* NatLog identifies expressions with non-default projectivity and computes the likely extent of their arguments in a syntactic parse using hand-crafted tree patterns.

We have evaluated NatLog on two different test suites. The first is the FraCaS test suite (Cooper et al. 1996), which contains 346 NLI problems, divided into nine sections, each focused on a specific category of semantic phenomena. The goal is three-way entailment classification, as described in Sect. 2. On this task, NatLog achieves an average accuracy of 70%.<sup>25</sup> In the section concerning quantifiers, which is both the largest and the most amenable to natural logic, the system answers all problems but one correctly. Unsurprisingly, performance is mediocre in four sections concerning semantic phenomena (e.g., ellipsis) not relevant to natural logic and not modeled by the system. But in the other five sections (representing about 60% of the problems), NatLog achieves accuracy of 87%. What’s more, precision is uniformly high, averaging 89% over all sections. Thus, even outside its areas of expertise, the system rarely predicts entailment when none exists.

The RTE3 test suite (Giampiccolo et al. 2007) differs from FraCaS in several important ways: the goal is binary entailment classification; the problems have much longer premises and are more “natural”; and the problems employ a diver-

---

<sup>25</sup>Our evaluation excluded multi-premise problems, which constitute about 44% of the test suite.

sity of types of inference—including paraphrase, temporal reasoning, and relation extraction—which NatLog is not designed to address. Consequently, the NatLog system by itself achieves mediocre accuracy (59 %) on RTE3 problems. However, its precision is comparatively high, which suggests a strategy of hybridizing with a broad-coverage RTE system. We were able to show that adding NatLog as a component in the Stanford RTE system (Chambers et al. 2007) led to accuracy gains of 4 %.

## 9 Conclusion

The model of natural logic presented here is by no means a universal solution to the problem of natural language inference. Many NLI problems hinge on types of inference not addressed by natural logic, and the inference method we describe faces a number of limitations on its deductive power (discussed in Sect. 7). Moreover, there is further work to be done in fleshing out our account of projectivity, particularly in establishing the proper projectivity signatures for a broader range of quantifiers, verbal constructs, implicatives and factives, logical connectives, and other semantic functions.

Nevertheless, we believe our model of natural logic fills an important niche. While approximate methods based on lexical and syntactic similarity can handle many NLI problems, they are easily confounded by inferences involving negation, antonymy, quantifiers, implicatives, and many other phenomena. Our model achieves the logical precision needed to handle such inferences without resorting to full semantic interpretation, which is in any case rarely possible. The practical value of the model is demonstrated by its success in evaluations on the FraCaS and RTE3 test suites.

## References

- Bos, J., & Markert, K. (2005). Recognising textual entailment with logical inference. In *Proceedings of the conference on human language technology and empirical methods in natural language processing* (pp. 628–635). Vancouver: Association for Computational Linguistics.
- Böttner, M. (1988). A note on existential import. *Studia Logica*, 47(1), 35–40.
- Chambers, N., Cer, D., Grenager, T., Hall, D., Kiddon, C., MacCartney, B., de Marneffe, M.-C., Ramage, D., Yeh, E., & Manning, C. D. (2007). Learning alignments and leveraging natural logic. In *Proceedings of the ACL-07 workshop on textual entailment and paraphrasing* (pp. 165–170). Prague: Association for Computational Linguistics.
- Condoravdi, C., Crouch, D., de Paiva, V., Stolle, R., & Bobrow, D. (2003). Entailment, intensionality and text understanding. In *Proceedings of the HLT-NAACL 2003 workshop on text meaning* (pp. 38–45). Morristown: Association for Computational Linguistics.
- Cooper, R., et al. (1996). *Using the framework* (Technical Report LRE 62-051 D-16). The FraCaS Consortium.
- Dagan, I., Glickman, O., & Magnini, B. (2006). The PASCAL recognising textual entailment challenge. In J. Quiñero-Candela, C. E. Rasmussen, F. Sinz, O. Bousquet, & B. Schölkopf (Eds.), *Machine learning challenges. Evaluating predictive uncertainty, visual object classification, and recognising textual entailment* (Vol. 3944, pp. 177–190). Berlin: Springer.

- Fyodorov, Y., Winter, Y., & Francez, N. (2000). A natural logic inference system. In *Proceedings of the 2nd international workshop on inference in computational semantics (ICoS-2)*, Germany: Dagstuhl.
- Giampiccolo, D., Magnini, B., Dagan, I., & Dolan, B. (2007). The third PASCAL recognizing textual entailment challenge. In *Proceedings of the ACL-07 workshop on textual entailment and paraphrasing* (pp. 1–9). Prague: Association for Computational Linguistics
- Glickman, O., Dagan, I., & Koppel, M. (2005). Web based probabilistic textual entailment. In *Proceedings of the PASCAL challenges workshop on recognizing textual entailment*. [http://u.cs.biu.ac.il/~nlp/RTE1/Proceedings/glickman\\_et\\_al.pdf](http://u.cs.biu.ac.il/~nlp/RTE1/Proceedings/glickman_et_al.pdf).
- Hickl, A., Williams, J., Bensley, J., Roberts, K., Rink, B., & Shi, Y. (2006). Recognizing textual entailment with LCC's GROUNDHOG system. In *Proceedings of the second PASCAL challenges workshop on recognizing textual entailment*, Venice, Italy, PASCAL (pp. 137–142).
- Lakoff, G. (1970). Linguistics and natural logic. *Synthese*, 22, 151–271.
- MacCartney, B. (2009). *Natural language inference*. Ph.D. thesis, Stanford University.
- MacCartney, B., & Manning, C. D. (2008). Modeling semantic containment and exclusion in natural language inference. In *Proceedings of the 22nd international conference on computational linguistics (COLING-08)* (pp. 521–528). Manchester: Association for Computational Linguistics.
- MacCartney, B., Grenager, T., de Marneffe, M.-C., Cer, D., & Manning, C. D. (2006). Learning to recognize features of valid textual entailments. In *Proceedings of the human language technology conference of the North American chapter of the Association of Computational Linguistics* (pp. 41–48). New York: Association for Computational Linguistics.
- Nairn, R., Condoravdi, C., & Karttunen, L. (2006). Computing relative polarity for textual inference. In *Proceedings of the fifth international workshop on inference in computational semantics (ICoS-5)* (pp. 67–76).
- MacCartney, B., Galley, M., & Manning, C. D. (2008). A phrase-based alignment model for natural language inference. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 802–811). Honolulu: Association for Computational Linguistics.
- Ritter, A., Downey, D., Soderland, S., & Etzioni, O. (2008). It's a contradiction—no, it's not: A case study using functional relations. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 11–20). Honolulu: Association for Computational Linguistics.
- Sánchez Valencia, V. (1991). *Studies on natural logic and categorial grammar*. Ph.D. thesis, Univ. Amsterdam.
- Sukkarieh, J. (2001). Quasi-NL knowledge representation for structurally-based inferences. In *Proceedings of the 3rd international workshop on inference in computational semantics (ICoS-3)*, Siena, Italy.
- van Benthem, J. (1988). The semantics of variety in categorial grammars. In W. Buszkowski, W. Marciszewski, & J. van Benthem (Eds.), *Categorial grammar* (pp. 33–55). Amsterdam: Benjamins.
- van Benthem, J. (1991). *Studies in logic: Vol. 130. Language in action: Categories, lambdas and dynamic logic*. Amsterdam: North-Holland.
- van Benthem, J. (2008). *A brief history of natural logic* (Technical Report PP-2008-05). Institute for Logic, Language & Computation. <http://www.illc.uva.nl/Publications/ResearchReports/PP-2008-05.text.pdf>.
- van der Sandt, R. A. (1992). Presupposition projection as anaphora resolution. *Journal of Semantics*, 9(4), 333–377.
- van Eijck, J. (2005). Natural logic for natural language. <http://homepages.cwi.nl/~jve/papers/05/nlnl/NLNL.pdf>.