

# The lexical and grammatical sources of neg-raising inferences

Hannah Youngeun An

Department of Computer Science  
University of Rochester

Aaron Steven White

Department of Linguistics  
University of Rochester

## Abstract

We investigate *neg(ation)-raising* inferences, wherein negation on a predicate can be interpreted as though in that predicate’s subordinate clause. To do this, we collect a large-scale dataset of neg-raising judgments for effectively all English clause-embedding verbs and develop a model to jointly induce the semantic types of verbs and their subordinate clauses and the relationship of these types to neg-raising inferences. We find that some neg-raising inferences are attributable to properties of particular predicates, while others are attributable to subordinate clause structure.

## 1 Introduction

Inferences that are triggered (at least in part) by particular lexical items provide a rich test bed for distinguishing the relative semantic contribution of lexical items and functional structure. One class of such inferences that has garnered extended attention is *neg(ation)-raising*, wherein negation on a predicate can be interpreted as though in that predicate’s subordinate clause (Fillmore, 1963; Bartsch, 1973; Horn, 1978; Gajewski, 2007). For example, a neg-raising inference is triggered by (1) while one is not triggered by (2).

- (1) Jo doesn’t think that Bo left.  
     $\rightsquigarrow$  Jo thinks that Bo didn’t leave.
- (2) Jo doesn’t know that Bo left.  
     $\nrightarrow$  Jo knows that Bo didn’t leave.

Though accounts vary with respect to whether neg-raising inferences are explained as a syntactic or a pragmatic phenomenon, all associate these inferences with particular predicates in some way or other—e.g. *think*, *believe*, *suppose*, *imagine*, *want*, and *expect* are often taken to be associated with neg-raising inferences as a matter of knowledge one has about those predicates, while *say*, *claim*, *regret*, and *realize* are not (Horn, 1971, 1978).

One challenge for such approaches is that whether a neg-raising inference is triggered varies with aspects of the context, such as the predicate’s subject—e.g. (3a) triggers the inference that the speaker thinks Jo didn’t leave—and tense—e.g. (3b) does not trigger the same inference as (3a).

- (3) a. I don’t know that Jo left.  
    b. I didn’t know that Jo left.

While some kinds of variability can be captured by standing accounts, other kinds have yet to be discussed at all. For example, beyond a predicate’s subject and tense, the syntactic structure of its clausal complement also appears to matter: (4a) and (5a) can both trigger neg-raising interpretations, while (4b) and (5b) cannot.

- (4) a. Jo wasn’t thought to be very intelligent.  
    b. Jo didn’t think to get groceries.
- (5) a. Jo wasn’t known to be very intelligent.  
    b. Jo didn’t know to get groceries.

Should these facts be chalked up to properties of the predicates in question? Or are they general to how these predicates compose with their complements? These questions are currently difficult to answer for two reasons: (i) there are no existing, lexicon-scale datasets that measure neg-raising across a variety of contexts—e.g. manipulating subject, tense and complement type; and (ii) even if there were, no models currently exist for answering these questions given such a dataset.

We fill this lacuna by (i) collecting a large-scale dataset of neg-raising judgments for effectively all English clause-embedding verbs with a variety of both finite and non-finite complement types; and (ii) extending White and Rawlins’ (2016) model of s(ematic)-selection, which induces semantic type signatures from syntactic distribution, with a module that associates semantic types with the inferences they trigger. We use this model to jointly

induce semantic types and their relationship to neg-raising inferences, showing that the best fitting model attributes some neg-raising inferences to properties of particular predicates and others to general properties of syntactic structures.<sup>1</sup>

We begin with background on theoretical approaches to neg-raising, contrasting the two main types of accounts: syntactic and pragmatic (§2). We then present our methodology for measuring neg-raising across a variety of predicates and syntactic contexts (§3) as well as our extension of White and Rawlins’ s-selection model (§4). Finally, we discuss the results of fitting (§5) our model to our neg-raising dataset (§6).

## 2 Background

Two main types of approaches have been proposed to account for neg-raising interpretations: syntactic and pragmatic (see Zeijlstra 2018; Crowley 2019 for reviews). We do not attempt to adjudicate between the two here—rather aiming to establish the explanatory devices available to each for later interpretation relative to our modeling results.

**Syntactic Approach** In syntactic approaches, neg-raising interpretations arise from some syntactic relation between a matrix negation and an unpronounced embedded negation that is licensed by the neg-raising predicate. This is classically explained via a syntactic rule that “raises” the negation from the subordinate clause to the main clause, as in (6), though accounts using alternative syntactic relations exist (Fillmore 1963; Kiparsky 1970; Jackendoff 1971; Pollack 1976; Collins and Postal 2014, 2017, 2018; cf. Klima, 1964; Zeijlstra, 2018; see also Lasnik, 1972).

- (6) Jo does not believe Bo did            leave.

Evidence for syntactic accounts comes from the distribution of negative polarity items, Horn-clauses, and island phenomena (Horn, 1971; Collins and Postal, 2014, 2017, 2018; cf. Zwarts, 1998; Gajewski, 2011; Chierchia, 2013; Horn, 2014; Romoli and Mandelkern, 2019).

Purely syntactic approaches to neg-raising have effectively one method for explaining variability in neg-raising inferences relative to subject, tense, and subordinate clause structure (as discussed in §1): if a certain lexical item—e.g. *know*—occurs in some sentence that licenses a neg-raising

inference—e.g. (5a)—and another that doesn’t—e.g. (5b)—one must say that the structure in the first differs from the second in such a way that the first allows the relevant syntactic relation while the second does not. This implies that, even in cases like (3a) v. (3b), where there is no apparent structural difference (beyond the subject), the structures differ on some neg-raising-relevant property. This can be implemented by saying that, e.g. the same verb can select for two different structural properties—one that licenses neg-raising and one that does not—or that the verb is somehow ambiguous and its variants differ with respect to some neg-raising-relevant, syntactic property.

**Semantic/Pragmatic Approach** In semantic/pragmatic approaches, neg-raising interpretations are derived from an *excluded middle* (EM or *opinionatedness*) inference (Bartsch, 1973; Horn, 1978; Horn and Bayer, 1984; Toven, 2001; Gajewski, 2007; Romoli, 2013; Xiang, 2013; Homer, 2015). This approach posits that, anytime a neg-raising predicate *v* is used to relate entity *x* with proposition *p*, the hearer assumes that either  $x \vee p$  or  $x \vee \neg p$ . For example, in the case of *believe*, as in (7), the hearer would assume that Jo either believes that Bo left or that Bo didn’t leave.

- (7) Jo believes that Bo left.

- a. *truth conditions*:  $x \text{ BELIEVE } p$   
b. *inference*:  $x \text{ BELIEVE } p \vee x \text{ BELIEVE } \neg p$

The EM inference is impotent in the positive cases but drives further inferences in the negative, where the first EM disjunct is cancelled by the truth conditions: if Jo doesn’t believe that Bo left and Jo believes that Bo left or that Bo didn’t leave, then Jo must believe that Bo didn’t leave.

- (8) Jo doesn’t believe that Bo left.

- a. *truth conditions*:  $x \neg \text{BELIEVE } p$   
b. *inference*:  $x \neg \text{BELIEVE } p \vee x \text{ BELIEVE } \neg p$

To capture non-neg-raising predicates, one must then say that some predicates trigger the EM inference, while others don’t (Horn, 1989). However, such lexical restrictions alone cannot exhaustively explain the variability in whether verbs trigger presuppositions with certain subjects, as noted for (2) and (3a). To explain this, Gajewski (2007) posits that neg-raising predicates are soft presupposition triggers. Effectively, the EM inferences are defeasible, and when they are cancelled, the neg-raising inference does not go through (Abusch, 2002). This is supported by cases of explicit cancella-

<sup>1</sup>Data are available at [megaattitude.io](http://megaattitude.io).

tion of the EM inference—e.g. the neg-raising inference (9c) that would otherwise be triggered by (9b) does not go through in the context of (9a).

- (9) a. Bill doesn't know who killed Caesar. He isn't even sure whether or not Brutus and Caesar lived at the same time. So...  
 b. Bill doesn't believe Brutus killed Caesar.  
 c.  $\nrightarrow$  Bill believes Brutus didn't kill Caesar.

This sort of explanation relies heavily on semantic properties of particular verbs and naturally covers variability that correlates with subject and tense differences—e.g. (3a) v. (3b)—since facts about how one discusses their own belief or desire states, in contrast to others belief states, at different times plausibly matter to whether a hearer would make the EM inference. The explanation for variation relative to subordinate clause structure is less clear but roughly two routes are possible: (i) some property of the subordinate clause licenses (or blocks) EM inferences; and/or (ii) predicate ambiguity correlates with which subordinate clause structure (or property thereof) a predicate selects.

**Abstracting the Approaches** Across both approaches, there are roughly three kinds of explanations for neg-raising inferences that can be mixed-and-matched: (i) lexical properties might directly or indirectly (e.g. via an EM inference) license a neg-raising inference; (ii) properties of a subordinate clause structure might directly or indirectly license a neg-raising inference; and/or (iii) lexical and structural properties might interact—e.g. via selection—to directly or indirectly license a neg-raising inference. We incorporate these three kinds of explanation into our models (§4), which we fit to the data described in the next section.

### 3 Data

We develop a method for measuring neg-raising analogous to White and Rawlins-White et al.'s (2018) method for measuring veridicality inferences. With the aim of capturing the range of variability in neg-raising inferences across the lexicon, we deploy this method to test effectively all English clause-embedding verbs in a variety of subordinate clause types—finite and nonfinite—as well as matrix tenses—*past* and *present*—and matrix subjects—*first* and *third person*.

**Method** Participants are asked to answer questions like (10) using a 0-1 slider, wherein the

first italicized sentence has negation in the matrix clause and the second italicized sentence has negation in the subordinate.<sup>2</sup>

- (10) If I were to say *I don't think that a particular thing happened*, how likely is it that I mean *I think that that thing didn't happen*?

Because some sentences, such the italicized in (11), sound odd with negation in the matrix clause, participants are asked to answer how easy it is to imagine someone actually saying the sentence—again, on a 0-1 slider. The idea here is that the harder it is for participants to imagine hearing a sentence, the less certain they probably are about the judgment to questions like (10).

- (11) How easy is it for you to imagine someone saying *I don't announce that a particular thing happened*?

Acknowledging the abuse of terminology, we refer to responses to (11) as *acceptability responses*. We incorporate these responses into our model (§4) as weights determining how much to pay attention to the corresponding *neg-raising response*.

**Materials** We use the MegaAcceptability dataset of White and Rawlins (2016) as a basis on which to construct acceptable items for our experiment. MegaAcceptability contains ordinal acceptability judgments for 50,000 sentences, including 1,000 clause-embedding English verbs in 50 different syntactic frames. To avoid typicality effects, these frames are constructed to contain as little lexical content as possible besides the verb at hand—a method we follow here. This is done by ensuring that all NP arguments are indefinite pronouns *someone* or *something* and all verbs besides the one being tested are *do*, *have* or *happen*. We focus on the six frames in (12)–(17).

- (12) [NP \_ that S]  
 Someone knew that something happened.  
 (13) [NP \_ to VP[+EV]]  
 Someone liked to do something.  
 (14) [NP \_ to VP[-EV]]  
 Someone wanted to have something.  
 (15) [NP be \_ that S]  
 Someone was told that something happened.  
 (16) [NP be \_ to VP[+EV]]  
 Someone was ordered to do something.  
 (17) [NP be \_ to VP[-EV]]  
 Someone was believed to have something.

<sup>2</sup>The full task instructions are given in Appendix A.

These frames were chosen so as to manipulate (i) the presence and absence of tense in the subordinate clause; (ii) the presence or absence of a direct object; and (iii) the lexical aspect of the complement. The frames with direct objects were presented in passivized form so that they were acceptable with both communicative predicates—e.g. *tell*—and emotive predicates—e.g. *sadden*—the latter of which tend to occur with expletive subjects. Lexical aspect was manipulated because some verbs—e.g. *believe*—are more acceptable with nonfinite subordinate clauses headed by a stative than ones headed by an eventive, while others—e.g. *order*—show the opposite pattern.

In light of the variability in neg-raising inferences across the same verb in different tenses—compare again (3a) and (3b)—we aim to manipulate the matrix tense of each clause-taking verb in our experiment. This is problematic, because the MegaAcceptability dataset only contains items in the past tense. We could simply manipulate the tense for any acceptable sentences based on such past tense items, but some verbs do not sound natural in the present tense with some subordinate clauses—compare the sentences in (18).

- (18) a. Jo wasn’t told that Mary left.  
b. Jo isn’t told that Mary left.

To remedy this, we extend MegaAcceptability with tense/aspect information by collecting acceptability judgments for modified versions of each sentence in MegaAcceptability, where the target verb is placed in either present or past progressive.<sup>3</sup> Combined with MegaAcceptability, our extended dataset results in a total of 75,000 verb-tense-frame pairs: 50,000 from the MegaAcceptability dataset and 25,000 from our dataset. From this combined dataset, we take past and present tense items rated on average 4 out of 7 or better (after rating normalization), for our experiment. This yields 3,968 verb-tense-frame pairs and 925 unique verbs. With our subject manipulation (first v. third person), the number of items doubles, producing 7,936 items. Table 1 summarizes the distribution of verbs in each frame and tense.

To construct items, we follow the method of White et al. (2018) of “bleaching” all lexical category words in our sentences (besides the subordinate clause-taking verb) by realizing NPs as *a particular person* or *a particular thing*. Verbs are

Matrix tense	Frame	# verbs
<i>past</i>	NP _ that S	556
	NP _ to VP[+EV]	400
	NP _ to VP[-EV]	359
	NP be _ that S	255
	NP be _ to VP[+EV]	461
	NP be _ to VP[-EV]	460
<i>present</i>	NP _ that S	413
	NP _ to VP[+EV]	219
	NP _ to VP[-EV]	155
	NP be _ that S	176
	NP be _ to VP[+EV]	268
	NP be _ to VP[-EV]	246

Table 1: # of verbs acceptable in each tense-frame pair based on our extension of MegaAcceptability.

replaced with *do*, *have*, or *happen*. This method aims to avoid unwanted typicality effects that might be introduced by interactions between our predicates of interest and more contentful items elsewhere in the sentence.<sup>4</sup>

We partition items into 248 lists of 32 items. Each list is constrained such that (i) 16 items had a first person subject, and 16 items had a third person subject; (ii) 16 items contain a *low frequency* verb and 16 items contain a *high frequency* verb, based on a median split of the frequencies in the SUBTLEX\_US word frequency database (Brysbaert and New, 2009); (iii) 16 items are *low acceptability* and 16 items are *high acceptability*, based on a median split of the normalized acceptabilities for items selected from our extension of the MegaAcceptability dataset; (iv) no verb occurred more than once in the same list; (v) items containing a particular combination of matrix tense and syntactic frame occur in rough proportion to the number of verbs that are acceptable with that tense-frame combination based on our extension of the MegaAcceptability dataset (Table 1).

**Participants** 1,108 participants were recruited through Amazon Mechanical Turk to give 10 ratings per sentence in the 248 lists of 32—i.e. the end result contains 79,360 ratings for each of neg-raising and acceptability judgments. Participants were not allowed to respond to the same list more than once, though they were allowed to respond to as many lists as they liked. Each participants re-

<sup>3</sup>See Appendix B for details.

<sup>4</sup>Because this method has not been previously validated for measuring neg-raising, we report two validation experiments in Appendix C, which demonstrate that the measure accords with judgments from prior work.



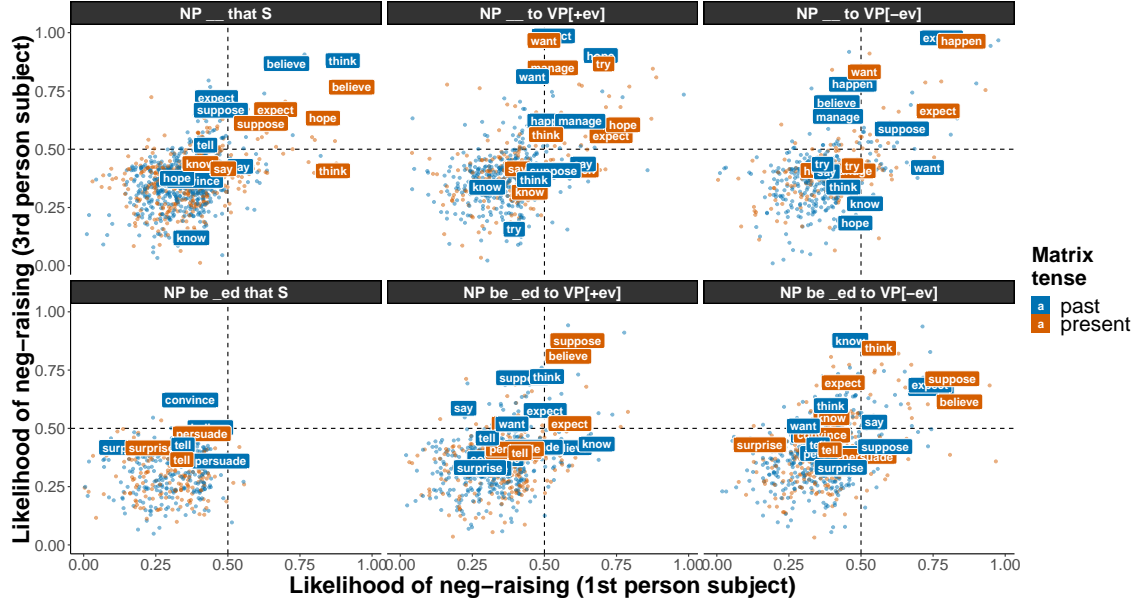


Figure 1: Normalized neg-raising scores for different subject, tense, and frame pairs.

sponded to 2.3 lists on average (min: 1, max: 16, median: 1). Of the 1,108 participants, 10 reported not speaking American English as their native language. Responses from these participants were filtered from the dataset prior to analysis. From this, responses for 27 lists were lost ( $\sim 1\%$  of the responses). This filtering removed at most two judgments for any particular item.

**Results** Figure 1 plots the normalized neg-raising scores for verbs in different subject (axes)-tense (color)-frame (block) contexts.<sup>5</sup> A verb (in some tense) being toward the top-right corner means that it shows strong neg-raising inferences with both first person and third person subjects, while a verb being towards the bottom-right corner means that it shows neg-raising behavior with first person subjects but not with third person subjects. The converse holds for the top-left corner: neg-raising behavior is seen with third person subjects but not first. We see that our method correctly captures canonical neg-raising predicates—e.g. *think* and *believe* with finite complements and *want* and *expect* with infinitival complements—as well as canonical non-neg-raising predicates—e.g. *know* and *say* with finite complements and *try* and *manage* with infinitivals.

## 4 Model

We aim to use our neg-raising dataset to assess which aspects of neg-raising inferences are due to properties of lexical items and which aspects are

due to properties of the structures they compose with. To do this, we extend White and Rawlins’ (2016) model of s(emantic)-selection, which induces semantic type signatures from syntactic distribution, with a module that associates semantic types with the *inference patterns* they trigger.

Our model has two hyperparameters that correspond to the theoretical constructs of interest: (i) the number of lexical properties relevant to neg-raising; and (ii) the number of structural properties relevant to neg-raising. In §5, we report on experiments aimed at finding the optimal setting of these two hyperparameters, and we analyze the parameters of the model fit corresponding to these hyperparameters in §6.

**S-selection Model** White and Rawlins’ (2016) model of s-selection aims to induce verbs’ semantic type signatures—e.g. that *love* can denote a relation between two entities and *think* can denote a relation between an entity and a proposition—from their syntactic distribution—e.g. that *love* is acceptable in NP \_\_ NP frames and that *think* is acceptable in NP \_\_ S frames. They formalize this task as a boolean matrix factorization (BMF) problem: given a boolean matrix  $\mathbf{D} \in \mathbb{B}^{|\mathcal{V}| \times |\mathcal{F}|} = \{0, 1\}^{|\mathcal{V}| \times |\mathcal{F}|}$ , wherein  $d_{vf} = 1$  iff verb  $v \in \mathcal{V}$  is acceptable in syntactic frame  $f \in \mathcal{F}$ , one must induce boolean matrices  $\mathbf{\Lambda} \in \mathbb{B}^{|\mathcal{V}| \times |\mathcal{T}|}$  and  $\mathbf{\Pi} \in \mathbb{B}^{|\mathcal{T}| \times |\mathcal{F}|}$ , wherein  $\lambda_{vt} = 1$  iff verb  $v$  can have semantic type signature  $t \in \mathcal{T}$  and  $\pi_{tf} = 1$  iff  $t$  can be mapped onto syntactic frame  $f$ , such that (19): verb  $v$  is acceptable in frame  $f$  iff  $v$  has some type  $t$  that can be mapped (or *projected*) onto  $f$ .

<sup>5</sup>See Appendix D for details on normalization.

$$(19) d_{vf} \approx \bigvee_t \lambda_{vt} \wedge \pi_{tf}$$

As is standard in matrix factorization, the equivalence is approximate and is only guaranteed when there are as many semantic type signatures  $\mathcal{T}$  as there are frames  $\mathcal{F}$ , in which case, the best solution is the one with  $\mathbf{\Lambda} = \mathbf{D}$  and  $\mathbf{\Pi}$  as the identity matrix of dimension  $|\mathcal{T}| = |\mathcal{F}|$ . Because this solution is trivial,  $|\mathcal{T}|$  is generally much smaller than  $|\mathcal{F}|$  and determined by fit to the data—in BMF, the count of how often  $d_{vf} \neq \bigvee_t \lambda_{vt} \wedge \pi_{tf}$ .

As an estimate of  $\mathbf{D}$ , [White and Rawlins](#) use the MegaAcceptability dataset, which we use in constructing our neg-raising dataset (§3). Instead of directly estimating the boolean matrices  $\mathbf{\Lambda}$  and  $\mathbf{\Pi}$ , they estimate a probability distribution over the two under the strong independence assumption that all values  $\lambda_{vt}$  and  $\pi_{tf}$  are pairwise independent of all other values. This implies (20).<sup>6</sup>

$$(20) \mathbb{P}(d_{vf}) = 1 - \prod_t 1 - \mathbb{P}(\lambda_{vt})\mathbb{P}(\pi_{tf})$$

[White and Rawlins](#) treat  $\mathbb{P}(d_{vf})$  as a fixed effect in an ordinal mixed effects model, which provides the loss function against which  $\mathbb{P}(\lambda_{vt})$  and  $\mathbb{P}(\pi_{tf})$  are optimized. They select the number of semantic type signatures to analyze by setting  $|\mathcal{T}|$  such that an information criterion is optimized.

**Neg-Raising Model** We retain the main components of [White and Rawlins](#)’ model but add a notion of *inference patterns* associated both with properties of verbs, on the one hand, and with semantic type signatures, on the other. In effect, this addition models inferences, such as neg-raising, as arising via a confluence of three factors: (i) properties of the relation a lexical item denotes—e.g. in a semantic/pragmatic approach, whatever property of a predicate triggers EM inferences; (ii) properties of the kinds of things that a predicate (or its denotation) relates—e.g. in a syntactic approach, whatever licenses “raising” of the negation; and (iii) whether a particular verb has a particular type signature. With respect to (ii) and (iii), it is important to note at the outset that, because we do not attempt to model acceptability, semantic type signatures play a somewhat different role in our model than in [White and Rawlins](#)’: instead of determining which structures a verb is compatible with—i.e. (non)finite subordinate clauses, presence of a direct object, etc.—our model’s type signatures control the inferences a particular verb can trigger when taking a particular structure. As such,

our model’s semantic type signatures might be more easily construed as properties of a structure that may or may not license neg-raising.<sup>7</sup> We thus refer to them as *structural properties*—in contrast to predicates’ *lexical properties*.

Our extension requires the addition of three formal components to [White and Rawlins](#)’ model: (i) a boolean matrix  $\mathbf{\Psi} \in \mathbb{B}^{|\mathcal{V}| \times |\mathcal{I}|}$ , wherein  $\psi_{vi} = 1$  iff verb  $v \in \mathcal{V}$  has property  $i \in \mathcal{I}$ ; (ii) a boolean tensor  $\mathbf{\Phi} \in \mathbb{B}^{|\mathcal{I}| \times |\mathcal{J}| \times |\mathcal{K}|}$ , wherein  $\phi_{ijk} = 1$  iff property  $i$  licenses a neg-raising inference with subject  $j \in \mathcal{J}$  and tense  $k \in \mathcal{K}$ ; and (iii) a boolean tensor  $\mathbf{\Omega} \in \mathbb{B}^{|\mathcal{T}| \times |\mathcal{J}| \times |\mathcal{K}|}$ , wherein  $\omega_{tjk} = 1$  iff semantic type signature  $t \in \mathcal{T}$  licenses a neg-raising inference with subject  $j$  and tense  $k$ .

As it stands, this formulation presupposes that there are both lexical and structural properties relevant to neg-raising. To capture the possibility that there may be only one or the other relevant to neg-raising, we additionally consider two families of *boundary models*. In the boundary models that posit no lexical properties—which (abusing notation) we refer to as  $|\mathcal{I}| = 0$ —we fix  $\mathbf{\Psi} = \mathbf{1}_{|\mathcal{V}|}$  and  $\mathbf{\Phi} = \mathbf{1}_{|\mathcal{I}|} \otimes \mathbf{1}_{|\mathcal{J}|} \otimes \mathbf{1}_{|\mathcal{K}|}$ . In the boundary models that posit no structural properties ( $|\mathcal{T}| = 0$ ) we fix  $\mathbf{\Pi} = \mathbf{1}_{|\mathcal{F}|}$ ,  $\mathbf{\Lambda} = \mathbf{1}_{|\mathcal{V}|}$ , and  $\mathbf{\Omega} = \mathbf{1}_{|\mathcal{T}|} \otimes \mathbf{1}_{|\mathcal{J}|} \otimes \mathbf{1}_{|\mathcal{K}|}$ .

Analogous to [White and Rawlins](#), we treat our task as a problem of finding  $\mathbf{\Lambda}, \mathbf{\Pi}, \mathbf{\Psi}, \mathbf{\Phi}, \mathbf{\Omega}$  that best approximate the tensor  $\mathbf{N}$ , wherein  $n_{vfjk} = 1$  iff verb  $v$  licenses neg-raising inferences in frame  $f$  with subject  $j$  and tense  $k$ . This is formalized in (21), which implies that  $n_{vfjk} = 1$  iff there is some pairing of semantic type signature  $t$  and inference pattern  $i$  such that (i) verb  $v$  has semantic type signature  $t$ ; (ii) verb  $v$  licenses inference pattern  $i$ ; (iii) semantic type signature  $t$  can map onto frame  $f$ ; and (iv) both  $t$  and  $i$  license a neg-raising inference with subject  $j$  and tense  $k$ .

$$(21) n_{vfjk} \approx \bigvee_{t,i} \lambda_{vt} \wedge \psi_{vi} \wedge \phi_{ijk} \wedge \pi_{tf} \wedge \omega_{tjk}$$

Also analogous to [White and Rawlins](#), we aim to estimate  $\mathbb{P}(n_{vfjk})$  (rather than  $n_{vfjk}$  directly) under similarly strong independence assumptions:  $\mathbb{P}(\lambda_{vt}, \psi_{vi}, \phi_{ijk}, \pi_{tf}, \omega_{tjk}) = \mathbb{P}(\lambda_{vt})\mathbb{P}(\psi_{vi})\mathbb{P}(\phi_{ijk})\mathbb{P}(\pi_{tf})\mathbb{P}(\omega_{tjk}) = \zeta_{vtifjk}$ , implying (22).

$$(22) \mathbb{P}(n_{vfjk}) = 1 - \prod_{t,i} 1 - \zeta_{vtifjk}$$

We design the loss function against which  $\mathbb{P}(\lambda_{vt})$ ,

<sup>6</sup>See Appendix E for the derivation of (20).

<sup>7</sup>Alternatively, they might be construed as (potentially cross-cutting) classes of syntactic structures and/or semantic type signatures that could be further refined by jointly modeling acceptability (e.g. as measured by MegaAcceptability) alongside our measure of neg-raising inferences.

$\mathbb{P}(\psi_{vi})$ ,  $\mathbb{P}(\phi_{ijk})$ ,  $\mathbb{P}(\pi_{tf})$ , and  $\mathbb{P}(\omega_{tjk})$  are optimized such that (i)  $\mathbb{P}(n_{vfjk})$  is monotonically related to the neg-raising response  $r_{vfjkl}$  given by participant  $l$  for an item containing verb  $v$  in frame  $f$  with subject  $j$  and tense  $k$  (if one exists); but (ii) participants may have different ways of using the response scale. For example, some participants may prefer to use only values close to 0 or 1, while others may prefer values near 0.5; or some participants may prefer lower likelihood values while others may prefer higher values. To implement this, we incorporate (i) a fixed scaling term  $\sigma_0$ ; (ii) a fixed shifting term  $\beta_0$ ; (iii) a random scaling term  $\sigma_l$  for each participant  $l$ ; and (iv) a random shifting term  $\beta_l$  for each participant  $l$ . We define the expectation for a response  $r_{vfjkl}$  as in (23).

$$(23) \quad \hat{r}_{vfjkl} = \text{logit}^{-1}(m_l \nu_{vfjk} + \beta_0 + \beta_l) \\ \text{where } \nu_{vfjk} = \text{logit}(\mathbb{P}(n_{vfjk})) \\ m_l = \exp(\sigma_0 + \sigma_l)$$

We optimize  $\mathbb{P}(\lambda_{vt})$ ,  $\mathbb{P}(\psi_{vi})$ ,  $\mathbb{P}(\phi_{ijk})$ ,  $\mathbb{P}(\pi_{tf})$ , and  $\mathbb{P}(\omega_{tjk})$  against a KL divergence loss, wherein  $r_{vfjkl}$  is taken to parameterize the true distribution and  $\hat{r}_{vfjkl}$  the approximating distribution.

$$(24) \quad D(r \parallel \hat{r}) = r \log \frac{\hat{r}}{r} + (1 - r) \log \frac{1 - \hat{r}}{1 - r}$$

To take into account that it is harder to judge the neg-raising inferences for items that one cannot imagine hearing used, we additionally weight the above-mentioned KL loss by a normalization of the acceptability responses for an item containing verb  $v$  in frame  $f$  with subject  $j$  and tense  $k$ . We infer this value from the acceptability responses for an item containing verb  $v$  in frame  $f$  with subject  $j$  and tense  $k$  given by participant  $l$ , assuming a form for the expected value of  $a_{vfjkl}$  as in (25)—analogous to (23). (Unlike  $\nu_{vfjk}$  in (23),  $\alpha_{vfjk}$  in (25) is directly optimized.)

$$(25) \quad \hat{a}_{vfjkl} = \text{logit}^{-1}(m'_l \alpha_{vfjk} + \beta'_0 + \beta'_l) \\ \text{where } m'_l = \exp(\sigma'_0 + \sigma'_l)$$

The final loss against which  $\mathbb{P}(\lambda_{vt})$ ,  $\mathbb{P}(\psi_{vi})$ ,  $\mathbb{P}(\phi_{ijk})$ ,  $\mathbb{P}(\pi_{tf})$ ,  $\mathbb{P}(\omega_{tjk})$  are optimized is (26).<sup>8</sup>

$$(26) \quad \mathcal{L} = - \sum \alpha'_{vfjk} D(r_{vfjkl} \parallel \hat{r}_{vfjkl}) \\ \text{where } \alpha'_{vfjk} = \text{logit}^{-1}(\alpha_{vfjk}).$$

## 5 Experiment

We aim to find the optimal settings, relative to our neg-raising data, for (i) the number  $|\mathcal{I}|$  of lexical

<sup>8</sup>An additional term (not shown) is added to encode the standard assumption that the random effects terms are normally distributed with mean 0 and unknown variance.

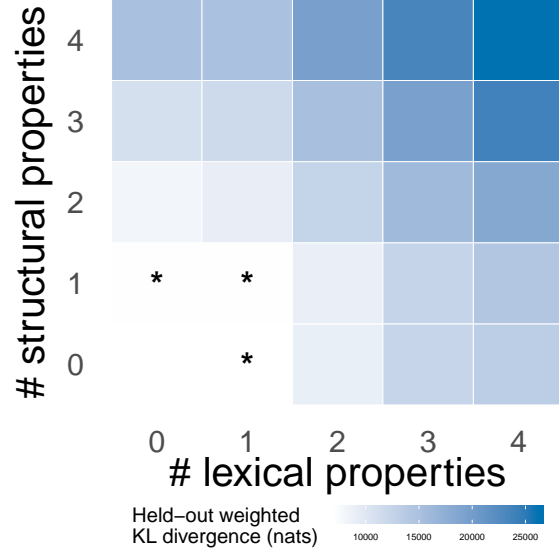


Figure 2: Sum of the weighted KL divergence loss across all five folds of the cross-validation for each setting of  $|\mathcal{I}|$  (# of lexical properties) and  $|\mathcal{T}|$  (# of structural properties).  $|\mathcal{I}| = |\mathcal{T}| = 0$  was not run.

properties relevant to neg-raising that it assumes; and (ii) the number  $|\mathcal{T}|$  of structural properties relevant to neg-raising that it assumes. As with other models based on matrix factorization, higher values for  $|\mathcal{I}|$  (with a fixed  $|\mathcal{T}|$ ) or  $|\mathcal{T}|$  (with a fixed  $|\mathcal{I}|$ ) will necessarily fit the data as well or better than lower values, since a model with larger  $|\mathcal{I}|$  or  $|\mathcal{T}|$  can embed the model with a smaller value. However, this better fit comes at the cost of increased risk of overfitting due to the inclusion of superfluous dimensions. To mitigate the effects of overfitting, we conduct a five-fold cross-validation and select the model(s) with the best performance (in terms of our weighted loss) on held-out data.

**Method** In this cross-validation, we pseudorandomly partition sentences from the neg-raising experiments into five sets (folds), fit the model with some setting of  $|\mathcal{I}|$ ,  $|\mathcal{T}|$  to the neg-raising responses for sentences in four of these sets (80% of the data), then compute the loss on the held-out set—repeating with each partition acting as the held-out set once. The assignment of items to folds is pseudorandom in that each fold is constrained to contain at least one instance of a particular verb with a particular complement type in some tense with some subject. If such a constraint were not enforced, on some folds, the model would have no data upon which to predict that verb with that complement. We consider each possible pairing of  $|\mathcal{I}|, |\mathcal{T}| \in \{0, 1, 2, 3, 4\}$ , except  $|\mathcal{I}| = |\mathcal{T}| = 0$ . The same partitioning is used for

every setting of  $|\mathcal{I}|$  and  $|\mathcal{T}|$ , enabling paired comparison by sentence.

**Implementation** We implement our model in `tensorflow 1.14.0` (Abadi et al., 2016). We use the Adam optimizer (Kingma and Ba, 2015) with a learning rate of 0.01 and default hyperparameters otherwise.

**Results** Figure 2 plots the sum of the weighted KL divergence loss across all five folds of the cross-validation for each setting of  $|\mathcal{I}|$  (number of lexical properties) and  $|\mathcal{T}|$  (number of structural properties). The best-performing models in terms of held-out loss (starred in Figure 2) are (in order): (i) one that posits one lexical property and no structural properties; (ii) one that posits no lexical properties and one structural property; and (iii) one that posits one lexical property and one structural property. None of these models’ performance is reliably different from the others—as determined by a nonparametric bootstrap computing the 95% confidence interval for the pairwise difference in held-out loss between each pairing among the three—but all three perform reliably better than all other models tested.

Among these three, the model with the best fit to the dataset has  $|\mathcal{I}| = 1$  and  $|\mathcal{T}| = 1$ . This result suggests that neg-raising is not purely a product of lexical knowledge: properties of the subordinate clause that a predicate combines with also influence whether neg-raising inferences are triggered. This is a surprising finding from the perspective of prior work, since (to our knowledge) no existing proposals posit that syntactic properties like the ones we manipulated to build our dataset—i.e. the presence or absence of tense, the presence or absence of an overt subject of the subordinate clause, and eventivity/stativity of a predicate in the subordinate clause—can influence whether neg-raising inferences are triggered. We next turn to analysis of this model fit to understand how our model captures patterns in the data.

## 6 Analysis

Table 2 gives the  $|\mathcal{I}| = |\mathcal{T}| = 1$  model’s estimate of the relationship between neg-raising inferences and lexical  $\mathbb{P}(\phi_{ijk})$  (top) and structural properties  $\mathbb{P}(\omega_{tjk})$  (bottom) with different subjects and tenses. The fact that all of the values in Table 2 are near 1 suggests that predicates having the lexical property or structures having the structural

Property	Person	Tense	
		<i>past</i>	<i>present</i>
<i>lexical</i>	<i>first</i>	0.93	0.98
	<i>third</i>	0.95	0.98
<i>structural</i>	<i>first</i>	0.93	0.98
	<i>third</i>	0.95	0.98

Table 2: Relationship between neg-raising inferences and lexical property  $\mathbb{P}(\phi_{ijk})$  (top) and structural property  $\mathbb{P}(\omega_{tjk})$  (bottom) with different subjects and tenses in  $|\mathcal{I}| = |\mathcal{T}| = 1$  model.

property will give rise to neg-raising inferences regardless of the subject and tense.<sup>9</sup>

This pattern is interesting because it suggests that the model does not capture the variability across different subjects and tenses observed in Figure 1 as a matter of either lexical or structural properties. That is, the model treats any variability in neg-raising inferences across different subjects and/or tenses as an idiosyncratic fact about the lexical item and the structure it occurs with—i.e. noise. This result makes intuitive sense insofar as such variability arises due to pragmatic reasoning that is specific to particular predicates, as opposed to some general semantic property.

But while the model does not distinguish among neg-raising inference with various subject and tense combinations, it does capture the coarser neg-raising v. non-neg-raising distinction among predicates—namely, by varying the probability that different lexical items have the lexical property  $\mathbb{P}(\psi_{vi})$  and the probability that they select the structural property  $\mathbb{P}(\lambda_{vt})$ . Figure 3 plots the distribution of  $\mathbb{P}(\psi_{vi}) \times \mathbb{P}(\lambda_{vt})$  across predicates.<sup>10</sup> We see that predicates standardly described as neg-raising (*think*, *believe*, *want*, *seem*, *feel*, etc.) fall to the right, while those standardly

<sup>9</sup>These tables appear to be copies of each other, but they are not. What is happening here is that the model is learning to associate  $\mathbb{P}(\phi_{ijk})$  and  $\mathbb{P}(\omega_{tjk})$  with (roughly) the square root of the largest expected value across all predicates for the neg-raising response to sentences with subject  $j$  and tense  $k$ . (It sets these values to the square root of the largest expected value because they will be multiplied together.) This strategy allows the model to simply vary  $\mathbb{P}(\lambda_{vt})$ ,  $\mathbb{P}(\psi_{vi})$ , and  $\mathbb{P}(\pi_{tf})$  to capture the likelihood a particular predicate or structure gives rise to neg-raising inferences, as described below.

<sup>10</sup>We plot the distribution of  $\mathbb{P}(\psi_{vi}) \times \mathbb{P}(\lambda_{vt})$ , instead of showing a scatter plot, because these probabilities show extremely high positive rank correlation—approximately 1. This happens because, when there is only one lexical property and one structural property, the lexical property and selection probabilities are effectively a single parameter  $p$ , with  $\mathbb{P}(\psi_{vi})$  and  $\mathbb{P}(\lambda_{vt})$  themselves being set to  $\sqrt{p}$  (see also Footnote 9).



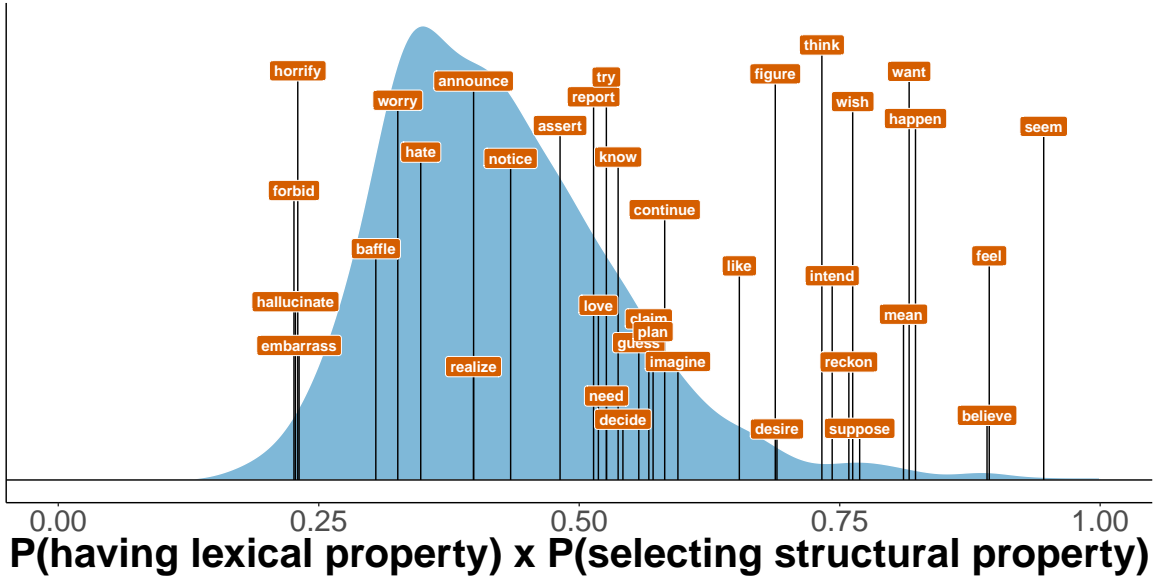


Figure 3: Distribution of  $\mathbb{P}(\psi_{vi}) \times \mathbb{P}(\lambda_{vt})$  across predicates, along with selected neg-raising (toward right) and non-neg-raising (toward left) predicates in  $|\mathcal{I}| = |\mathcal{T}| = 1$  model. (Label height is jittered to avoid overplotting.)

described as non-neg-raising (*know*, *notice*, *realize*, *love*, etc.) fall to left. Thus, in some sense, a predicate’s probability of having the model’s single lexical property (plus its probability of selecting the single structural property) appears to capture something like the probability of neg-raising.

Structure	Probability
NP __ that S	0.91
NP be _ed that S	0.84
NP __ to VP[+ev]	0.98
NP be _ed to VP[+ev]	0.93
NP __ to VP[-ev]	0.94
NP be _ed to VP[-ev]	0.98

Table 3: Relationship between structural property and structures  $\mathbb{P}(\pi_{tf})$  in  $|\mathcal{I}| = |\mathcal{T}| = 1$  model.

The model captures variability with respect to different syntactic structures by modulating  $\mathbb{P}(\pi_{tf})$ , shown in Table 3. Looking back to Figure 1, these values roughly correlate with the largest neg-raising response (across subjects and tenses) seen in that frame, with NP be \_ed that S showing the lowest such value. The value of  $\mathbb{P}(\pi_{tf})$  is not the *same* as the largest neg-raising value in Figure 1, likely due to the fact that many of the predicates that occur in that frame also have small values for  $\mathbb{P}(\psi_{vi}) \times \mathbb{P}(\lambda_{vt})$ , and thus, when  $\mathbb{P}(\pi_{tf})$  is multiplied by that values, it is small.

## 7 Conclusion

We presented a probabilistic model to induce the mappings from lexical sources and their gram-

matical sources to neg-raising inferences. We trained this model on a large-scale dataset of neg-raising judgments that we collected for 925 English clause-embedding verbs in six distinct syntactic frames as well as various matrix tenses and subjects. Our model fit the best when positing one lexical property and one structural property. This is a surprising finding from the perspective of prior work, since (to our knowledge) no existing proposals posit that syntactic properties like the ones we manipulated to build our dataset—i.e. the presence or absence of tense, the presence or absence of an overt subject of the subordinate clause, and eventivity/stativity of a predicate in the subordinate clause—can influence whether neg-raising inferences are triggered. Our findings suggest new directions for theoretical research attempting to explain the interaction between lexical and structural factors in neg-raising. Future work in this vein might extend the model proposed here to investigate the relationship between neg-raising and acceptability as well as other related phenomena with associated large-scale datasets, such as lexically triggered veridicality inferences (White and Rawlins, 2018; White et al., 2018; White, 2019).

## Acknowledgments

We would like to thank the FACTS.lab at UR as well as three anonymous reviewers for useful comments. This work was supported by an NSF grant (BCS-1748969/BCS-1749025) *The MegaAttitude Project: Investigating selection and polysemy at the scale of the lexicon*.

## References

- Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. [Tensorflow: A system for large-scale machine learning](#). In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 265–283.
- Dorit Abusch. 2002. Lexical alternatives as a source of pragmatic presuppositions. *Semantics and Linguistic Theory*, 12:1–19.
- Renate Bartsch. 1973. “Negative transportation” gibt es nicht. *Linguistische Berichte*, 27(7).
- Marc Brysbaert and Boris New. 2009. Moving beyond kucera and francis: A critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for american english. *Behavior Research Methods*, 41:977–990.
- Gennaro Chierchia. 2013. *Logic in grammar: polarity, free choice, and intervention*, first edition. Oxford University Press, Oxford.
- Chris Collins and Paul Martin Postal. 2014. *Classical NEG Raising: An Essay on the Syntax of Negation*. MIT Press.
- Chris Collins and Paul Martin Postal. 2017. Interclausal neg raising and the scope of negation. *Glossa: A Journal of General Linguistics*, 2:1–29.
- Chris Collins and Paul Martin Postal. 2018. Disentangling two distinct notions of neg raising. *Semantics and Pragmatics*, 11(5).
- Paul Crowley. 2019. Neg-raising and neg movement. *Natural Language Semantics*, 27(1):1–17.
- Mark Davies. 2017. [Corpus of Contemporary American English \(COCA\)](#).
- Charles J. Fillmore. 1963. The position of embedding transformations in a grammar. *WORD*, 19(2):208–231.
- Jon R. Gajewski. 2007. Neg-raising and polarity. *Linguistics and Philosophy*, 30(3):289–328.
- Jon R. Gajewski. 2011. Licensing strong npis. *Natural Language Semantics*, 19(2):109–148.
- Vincent Homer. 2015. Neg-raising and positive polarity: The view from modals. *Semantics and Pragmatics*, 8(4):1–88.
- Laurence Robert Horn. 1971. Negative transportation: unsafe at any speed? *Papers from the seventh regional meeting, Chicago Linguistic Society*, pages 120–133.
- Laurence Robert Horn. 1978. Remarks on neg-raising. *Syntax and Semantics*, 9:129–220.
- Laurence Robert Horn. 1989. *A Natural History of Negation*. University of Chicago Press.
- Laurence Robert Horn. 2014. The cloud of unknowing. In Jack Hoeksema and Dicky Gilbers, editors, *Black Book: A Festschrift for Frans Zwarts*, pages 178–196. University of Groningen, Groningen, the Netherlands.
- Laurence Robert Horn and Samuel Bayer. 1984. Short-circuited implicature: A negative contribution. *Linguistics and Philosophy*, 7(4):397–414.
- Ray S. Jackendoff. 1971. On some questionable arguments about quantifiers and negation. *Language*, 47(2):282–297.
- Adam Kilgariff, Vt Baisa, Jan Buta, Milo Jakubek, Vojtech Kov, Jan Michelfeit, Pavel Rychl, and Vt Suchomel. 2014. [The sketch engine: ten years on](#). *Lexicography*, pages 7–36.
- Diederik P. Kingma and Jimmy Lei Ba. 2015. Adam: A method for stochastic optimization. *Proceedings of the 3rd International Conference on Learning Representations*.
- Paul Kiparsky. 1970. Semantic rules in grammar. In Hreinn Benediktsson, editor, *The Nordic Languages and Modern Linguistics*, pages 262–285. Visindafelag Islendinga, Reykjavik.
- Edward S. Klima. 1964. *Negation in English*. Englewood Cliffs, NJ: Prentice-Hall.
- Howard Lasnik. 1972. *Analyses of Negation in English*. Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA.
- Fabian Pedregosa, Gal Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay. 2011. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12:2825–2830.
- Jay M. Pollack. 1976. A re-analysis of neg-raising in english. *Working Papers in Linguistics*, 21:189–239.
- Jacopo Romoli. 2013. A scalar implicature-based approach to neg-raising. *Linguistics and Philosophy*, 36(4):291–353.
- Jacopo Romoli and Matthew Mandelkern. 2019. Whats not to like. *Linguistic Inquiry*, 0(ja):1–21.
- Nadine Theiler, Floris Roelofsen, and Maria Aloni. 2017. Whats wrong with believing whether? *Semantics and Linguistic Theory*, 27:248–265.

- Lucia M. Tovená. 2001. Neg-raising: negation as failure. In Jack Hoeksema, Hotze Rullmann, Victor Sánchez-Valencia, and Ton van der Wouden, editors, *Perspectives on Negation and Polarity Items*, pages 331–356. John Benjamins, Amsterdam.
- Aaron Steven White. 2019. Lexically triggered veridicality inferences. To appear in *Handbook of Pragmatics*.
- Aaron Steven White and Kyle Rawlins. 2016. A computational model of s-selection. *Semantics and Linguistic Theory*, 26:641–663.
- Aaron Steven White and Kyle Rawlins. 2018. The role of veridicality and factivity in clause selection. *Proceedings of the 48th Meeting of the North East Linguistic Society*.
- Aaron Steven White and Kyle Rawlins. 2019. Frequency, acceptability, and selection: A case study of clause-embedding. Under revision for *Glossa*.
- Aaron Steven White, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2018. Lexicosyntactic inference in neural models. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4717–4724.
- Yimei Xiang. 2013. Neg-raising: Focus and implications. *Proceedings of Sinn und Bedeutung*, 18:487–503.
- Hedde Zeijlstra. 2018. Does neg-raising involve neg-raising? *Topoi*, 37(3):417–433.
- Frans Zwarts. 1998. Three types of polarity. In Fritz Hamm and Erhard Hinrichs, editors, *Plurality and Quantification*, pages 177–238. Springer Netherlands, Dordrecht.

## A Instructions

In this experiment, you will be asked to answer questions about what a person is likely to mean if they say a particular sentence.

Your task will be to respond about the likelihood on the slider that will appear under each question, where the left side corresponds to *extremely unlikely* and the right side corresponds to *extremely likely*.

For instance, you might get the question *If I were to say John has three kids, how likely is it that I mean John has exactly three kids?* with a slider. In this case you would move the slider handle fairly far to the right (toward *extremely likely*), since if someone says “John has three kids”, it’s pretty likely that they mean that John has exactly three children.

If the question were *If I were to say some of the boys left, how likely is it that I mean all of the boys*

*left?*, then you might move the slider pretty far to the left (toward *extremely unlikely*), since it would be odd if someone says “Some of the boys left and by that, I mean all of the boys left”.

And if the question were *If I were to say Ann didn’t greet everyone politely, how likely is it that I mean Ann was unwelcoming to every single person?*, you might leave the slider in the middle (which corresponds to *maybe or maybe not*), since quite often such sentence can be used to mean Ann greeted some people politely but not all, or to mean Ann was not polite to every single person.

Try to answer the questions as quickly and accurately as possible. Many of the sentences may not be sentences that you can imagine someone ever saying. Try your best to interpret what a speaker would mean in using them. After each question, you will be given a chance to tell us whether the sentence you just responded to isn’t something you can imagine a native English speaker ever saying.

Not all questions have correct answers, but a subset in each HIT do. Prior to approval, we check the answers given for this subset. We will reject work containing a substantial number of answers that do not agree with the correct answer.

When the experiment is over, a screen will appear telling you that you are done, and a submission button will be revealed.

## B Data

We extend White and Rawlins’ (2016) MegaAcceptability v1.0 dataset by collecting acceptability judgments for sentences in present and past progressive tenses—resulting in MegaAcceptability v2.0, which subsumes MegaAcceptability v1.0. To enable comparison of the judgments given in MegaAcceptability v1.0 and those we collect, we run an additional *linking experiment* with half items from MegaAcceptability v1.0 and our extension. We then normalize all three datasets separately using the procedure described in White and Rawlins 2019 and then combine them by using the linking experiment data to train a model to map them into a comparable normalized rating space. Both the extended MegaAcceptability and linking datasets are available at [megaattitude.io](https://megaattitude.io).

**Extended MegaAcceptability** Our test items are selected and modified from the top 25% most acceptable verb-frame pairs from the MegaAcceptability dataset of White and Rawlins (2016),

determined by a modified version of the normalization procedure used in [White and Rawlins 2019](#). This item set thus contains 12,500 verb-frame pairs, with 1000 unique verbs and the same 50 subcategorization frames (35 in active voice and 15 in passive voice) that are used in MegaAcceptability.

Given the 12,500 verb-frame pairs, we construct new sentences in both present and past progressive tense/aspect, resulting in a total of 25,000 items. Examples of two sentences from MegaAcceptability v1.0 are given in (27) and the corresponding present and past progressive versions are given in (28) and (29), respectively.

- (27) a. Someone knew which thing to do.  
b. Someone talked about something.
- (28) a. Someone knows which thing to do.  
b. Someone talks about something.
- (29) a. Someone is knowing which thing to do.  
b. Someone was talking about something.

All methods follow [White and Rawlins 2016](#). Sentences are partitioned into 500 lists of 50, with each list constructed such that (i) each frame shows up once in a list, making each list contain 50 unique frames, if possible; (ii) otherwise, the distribution of frames are kept as similar as possible across lists; and (iii) no verbs appear more than once in a list. We gather 5 acceptability judgments per sentence, yielding a total of 125,000 judgments for 25,000 items.

Judgments for each sentence in a list are collected on a 1-to-7 scale. To avoid typicality effects, we construct the frames to contain as little lexical content as possible besides the verb at hand. For this, we instantiate all NP arguments as indefinite pronouns *someone* or *something* and all verbs besides the one being tested as *do* or *happen*. 565 participants were recruited from Amazon Mechanical Turk, where 562 speak American English as their native language.

**Linking experiment** Because our extension of MegaAcceptability was built in such a way that it likely contains higher acceptability items, the ratings in MegaAcceptability v1.0 and the ratings in our extension are likely not comparable—i.e. a rating in MegaAcceptability v1.0 is, in some sense, a worse rating than in our extension, since our sentences are, by construction, better overall. To put the existing MegaAcceptability dataset and our extended dataset on a comparable scale,

we run another experiment to assist in mapping the two datasets to such a comparable scale. We choose 50 items, each with a unique verb, by selecting 26 items from our dataset (14 in present tense and 12 in past progressive tense) and 24 items from MegaAcceptability (all past tense).

This item selection was constrained such that half of the items chosen were below the median acceptability score and half were above, evenly split across items from our experiment and items from MegaAcceptability v1.0. The items with the lowest acceptability scores consist of 8 in the present, 6 in the past progressive, and 12 in the past tense and so do the items with the highest acceptability scores. Example items with the low acceptability scores (under this criterion) are shown in (30), and example items with high acceptability scores are shown in (31).

- (30) a. Someone demands about whether something happened.  
b. Someone was judging to someone that something happened.  
c. Someone invited which thing to do.
- (31) a. Someone is distracted.  
b. Someone was teaching.  
c. Someone dared to do something.

The linking experiment is built in a very similar manner to our extension of MegaAcceptability, described above. Ordinal acceptability judgments are collected on a 1-to-7 scale. 50 participants were recruited to rate all 50 items in the experiment. All of the 50 participants report speaking American English as their native language.

After running the linking experiment, we normalize the ratings in all three datasets separately using a modified version of the procedure described in [White and Rawlins 2019](#). Then, we construct one mapping from the normalized ratings in our extension of MegaAcceptability to the normalized ratings for the linking dataset and another mapping from the normalized ratings in the linking dataset to the normalized ratings in MegaAcceptability v1.0 with two linear regressions—implemented in `scikit-learn` ([Pedregosa et al., 2011](#)). We then compose these two regressions to map the normalized ratings in our extended MegaAcceptability dataset to those in MegaAcceptability v1.0. This gives us a combined dataset of acceptability judgments for sentences in three different tense/aspect combinations



Subordinate clause	Neg-raising	Non-neg-raising
<i>Finite</i>	think, believe, feel, reckon, figure, guess, suppose, imagine	announce, claim, assert, report, know, realize, notice, find out
<i>Infinitival</i>	want, wish, happen, seem, plan, intend, mean, turn out	love, hate, need, continue, try, like, desire, decide

Table 4: Verbs used in validation experiments

(*past*, *present*, and *past progressive*) and 50 different syntactic frames, which we use to construct our neg-raising experiment.

## C Validation Experiments

We conduct experiments aimed at validating our method for measuring neg-raising. In both experiments, we test the same set of 32 clause-embedding verbs, half of which we expect to show neg-raising behavior and the other half we do not (based on the literature discussed in §2). For neg-raising verbs, we refer to the neg-raising predicates listed in Gajewski 2007 and Collins and Postal 2018; and for non-neg-raising verbs, we choose factive verbs and those that Theiler et al. (2017) claim are not neg-raising. The experiments differ with respect to whether we employ “bleached” items (as in the data collection described in the main body of the paper) or “contentful” items, which are constructed based on sentences drawn from English corpora.

**Materials** We select neg-raising and non-neg-raising verbs such that half of each type takes infinitival subordinate clauses and half takes finite subordinate clauses. Table 4 shows the 32 verbs we choose for the pilot. Some verbs listed as taking one kind of subordinate clause can also take the other. In these cases, we only test that verb in the subordinate clause listed in Table 4.

The matrix subject (first v. third person) and matrix tense (present v. past) are manipulated for each predicate: (32) schematizes four items from our bleached experiment and (33) schematizes four items from our contentful experiment.

(32) {I, A particular person} {don’t/doesn’t, didn’t} want to do a particular thing.

(33) {I, Stephen} {don’t/doesn’t, didn’t} want to introduce new rules.

Items for the bleached experiment are constructed

automatically using the templates, which select *to have a particular thing* for *turn out* and *seem* as their subordinate clause, *to do a particular thing* for other verbs taking infinitival subordinate clauses, and *that something happened* for the verbs taking finite subordinate clauses. Items for the contentful experiment are constructed by replacing all bleached words (*a particular person*, *a particular thing*, *do*, *have*, and *happen*) from the bleached experiment items by contentful lexical words.

The high content sentences are constructed based on sentences sampled from the Corpus of Contemporary American English (Davies, 2017) and the Oxford English Corpus (Kilgarriff et al., 2014). The contentful items are modified so that third person subject is a proper name and sentences do not include any pauses or conjunctions. To allow possible item variability, we create five contentful items per each bleached item.

For the bleached experiment, four lists of 32 items each are constructed by partitioning the resulting 128 items under the constraints that (i) every list contains every verb with exactly one subject (*first*, *third*) and tense (*past*, *present*) and (ii) every subject-tense pair is seen an equal number of times across verbs. We ensure that the same level of a particular factor is never assigned to the same verb more than once in any list and that the items in a list are randomly shuffled. To construct items, we manipulate neg-raising, embedded complement, matrix subject, matrix tense. Neg-raising and embedded complements are pre-determined for each verb, while matrix subject and matrix tense are randomly selected for a verb in each task. The same constraints apply for the contentful experiment except that the test items were partitioned into 20 lists of 32 instead of four lists because the total number of sentences for the contentful experiment is five times bigger than the bleached experiment.

**Participants** For the bleached experiment, 100 participants were recruited such that each of the four lists was rated by 25 unique participants. For the contentful experiment, 100 participants were recruited as well, to have each of the 20 lists of 32 rated by five unique participants. No participant was allowed to rate more than one list. In each experiment, one participant out of 100 reported not speaking American English natively and this participant’s responses were filtered prior to analysis.

**Analysis** We test whether our task correctly captures canonical (non-)neg-raising verbs using linear mixed effects models. For both validation experiments, we start with a model containing fixed effects for NEGRAISING (*true, false*; as in Table 4), random intercepts for PARTICIPANT, VERB, and (in the contentful validation) ITEM. Nested under both verb and participant, we also included random intercepts for MATRIX SUBJECT (*1st, 3rd*) and MATRIX TENSE (*past, present*) and their interaction. We compare this against a model with the same random effects structure but no effect of NEGRAISING. We find a reliably positive effect of NEGRAISING for both the bleached experiment ( $\chi^2(1) = 34.5, p < 10^{-3}$ ) and the contentful experiment ( $\chi^2(1) = 19.8, p < 10^{-3}$ ). This suggests that participants’ responses are consistent with neg-raising inferences being more likely with verbs that have previously been claimed to give rise to such inferences.

## D Normalization

For the purposes of visualization in §3, we present normalized neg-raising scores. These scores are derived using a mixed effects robust regression with loss the same loss (26) as for the model described in Section 4, except that, unlike for the model, where  $\nu_{vfjk}$  is defined in terms of the model, for the purposes of normalization, both  $\nu_{vfjk}$  in (23) and  $\alpha_{vfjk}$  in (25) are directly optimized. Figure 1 plots  $\text{logit}^{-1}(\exp(\sigma_0)\nu_{vfjk}) + \beta_0$ .

## E Model Derivation

$$\begin{aligned}
\mathbb{P}(d_{vf}) &= \mathbb{P}\left(\bigvee_t \lambda_{vt} \wedge \pi_{tf}\right) \\
&= \mathbb{P}\left(\neg \neg \bigvee_t \lambda_{vt} \wedge \pi_{tf}\right) \\
&= \mathbb{P}\left(\neg \bigwedge_t \neg(\lambda_{vt} \wedge \pi_{tf})\right) \\
&= \mathbb{P}\left(\neg \bigwedge_t \neg(\lambda_{vt} \wedge \pi_{tf})\right) \\
&= 1 - \mathbb{P}\left(\bigwedge_t \neg(\lambda_{vt} \wedge \pi_{tf})\right) \\
&= 1 - \prod_t \mathbb{P}(\neg(\lambda_{vt} \wedge \pi_{tf})) \\
&= 1 - \prod_t 1 - \mathbb{P}(\lambda_{vt} \wedge \pi_{tf}) \\
&= 1 - \prod_t 1 - \mathbb{P}(\lambda_{vt})\mathbb{P}(\pi_{tf})
\end{aligned}$$