

25 15. The application of experimental methods in semantics

26

27 1. Introduction

28 2. The stumbling blocks

29 3. Off-line evidence for scope interpretation

30 4. Underspecification vs. full interpretation

31 5. On-line evidence for representation of scope

32 6. Conclusions

33 7. References

34

35 Abstract

36 *The purpose of this paper is twofold. On the methodological side, we shall attempt to*
37 *show that even relatively simple and accessible experimental methods can yield*
38 *significant insights into semantic issues. At the same time, we argue that experimental*
39 *evidence, both the type collected in simple questionnaires and measures of on-line*
40 *processing, can inform semantic theories. The specific case that we address here*
41 *concerns the investigation of quantifier scope. In this area, where judgments are often*
42 *subtle and controversial, the gradient data that psycholinguistic experiments provide*
43 *can be a useful tool to distinguish between competing approaches, as we demonstrate*
44 *with a case study. Furthermore, we describe how a modification of existing*
45 *experimental methods can be used to test predictions of underspecification theories. The*
46 *programme of research we outline here is not intended to be a prescriptive set of*
47 *instructions for researchers, telling them what they should do; rather it is intended to*
48 *illustrate some problems an experimental semanticist may encounter but also the profit*

49 *of this enterprise.*

50

51 1. Introduction

52 A wide range of data types and sources are used in the field of semantics, as is
53 demonstrated by the related article 12 (Krifka) *Varieties of semantic evidence* in this
54 volume. The aim of this article is to show with an example research study series what
55 sort of questions can be addressed with experimental tools and suggest that these
56 methods can deliver valuable data which is relevant to basic assumptions in semantics.
57 This text also attempts to address the constraints on and limits to such an approach.
58 These are both methodological and theoretical: it has long been recognized that links
59 between empirical measures and theoretical constructs require careful argumentation to
60 establish.

61 The authors therefore have two aims: one related to experimental methodologies and the
62 other to do with the value of processing data. They first seek to show that even
63 relatively simple and accessible experimental methods can yield significant insights into
64 semantic issues. They second wish to illustrate that experimental evidence such as that
65 gathered in their eye-tracking study has the potential to inform semantic theory.

66 Semanticists have of course always sought confirmatory evidence to support their
67 analyses. There is, on the one hand, fairly extensive use of computational techniques
68 and corpus data in the field, and a growing body of experimental work on semantic
69 processing, language acquisition, and pragmatics, but in the area of theoretical and
70 formal semantics the experimental methods are less frequently employed.

71 Now there are good reasons for this. There are inherent factors related to the
72 accessibility of the relevant measures why controlled data gathering techniques are still

73 somewhat less frequent in this field than in some others. We shall discuss what these
74 reasons are and demonstrate with a case study what constraints they place on empirical
75 studies, particularly experimental studies. The example research program that we shall
76 report is thus not simply a recipe for others for what should be done, rather it is an
77 illustration of the difficulties involved, which aims to explore some of the boundaries of
78 what is accessible to experimental studies.

79 The specific case that we address here concerns the investigation of quantifier scope, a
80 perennial issue in semantics. Previous attempts to account for the complex data patterns
81 to be found in natural languages have met with the difficulty that the causal factors and
82 preferences need first to be identified before a realistic model can be developed. This
83 requires as an initial step the capture and measurement of the relevant effects and their
84 interactions, which is no trivial task.

85 The next section lays out a range of reasons why semanticists do not routinely seek to
86 test the empirical bases of their theories with simple experiments. Section 3 reports the
87 series of empirical investigations on quantifier scope carried out by Bott and Radó in
88 on-going research. Section 4 lays out some of the theoretical background and
89 importance of these studies for current theory (the underspecification debate). The final
90 section takes as a starting point Bott and Radó (2009) to suggest how some of the
91 problems noted in section 3 may be overcome with a more sophisticated experimental
92 procedure.

93

94 2. The stumbling blocks

95 As Manfred Krifka notes in his neighbouring article 12 (Krifka) *Varieties of semantic*
96 *evidence*, a major problem with investigating meaning is that we cannot yet fully define

97 what it is. This is indeed a root cause of difficulty, but here we shall attempt to illustrate
98 in more practical detail what effects this has on attempts to conduct experiments in this
99 field.

100

101 2.1. Specifying meaning without using language

102 The essential feature distinguishing experiment procedure is control. In language
103 experiments we may distinguish three (sets of) variables: linguistic form, context, and
104 meaning. In the typical experiment we will keep two of them constant and
105 systematically vary the other. Much semantic research concerns the systematic
106 interdependence of form, context, and meaning. These issues can be investigated for
107 example by:

- 108 a) keeping form and context constant, manipulating meaning systematically, and
109 measuring the *felicity* of the outcome (in judgements, or reaction times, or processing
110 effort), or
- 111 b) manipulating (at least one of) form and context, and measuring perceived meaning.

112 The first requires the experimenter to *manipulate* meaning as a variable, which entails
113 expressing meaning in a form other than language, (pictures, situation descriptions, etc);
114 the second requires the experimenter to *measure* perceived meaning, which again
115 normally demands reference to meanings captured in non-linguistic form. But precisely
116 this expression of tightly constrained meaning in non-linguistic form is very difficult.

117 To show how this factor affects studies in semantics disproportionately, it is worth
118 noting how this makes controlled studies in semantics more challenging than in syntax.

119 Work in experimental syntax is often interested in addressing precisely those effects of
120 form change which are *independent* of meaning. The variable meaning can thus be held

121 constant, but this does not require it to be exactly specified. It often does not much
122 matter exactly what interpretation subjects assign to the example structures as long as it
123 is the same for all of them. Thus only the syntactic analysis need be controlled, not the
124 meaning that this analysis gives rise to. This makes empirical studies in syntax much
125 less difficult than those in semantics.

126

127 2.2. The boundaries of form, context, and meaning

128 A further problem of exact studies concerning meaning is that the three variables are not
129 always clearly distinguished, in part because they systematically covary, but also in part
130 because linguists do not always agree about the boundaries. This is particularly visible
131 when we seek to identify where an anomaly lies. Views have changed over time in
132 linguistics about the nature and location of ill-formedness (e.g. the discussion of the
133 status of *I am lurking in a culvert* in Ross 1970) but the fundamental ambiguity is still
134 with us. For example, Weskott & Fanselow (2009) give the following examples and
135 judgements of syntactic and semantic well-formedness: (1a) is syntactically ill-formed
136 (*), (1b) is semantically ill-formed (#), and (1c) is ill-formed on both accounts (*#).

137 (1) a. *Die Suppe wurde gegen versalzen.

138 the soup was against oversalted

139 b. #Der Zug wurde gekaut.

140 the train was chewed

141 c. *#Das Eis wurde seit entzündet.

142 the ice was since inflamed

143 Our own judgements suggest that the structures in (1-a) and (1-c) have no acceptable
144 syntactic analysis, and therefore no semantic analysis can be constructed -- they are thus

145 both syntactically and semantically ill-formed. Crucially, the semantic anomaly is
146 dependent upon the syntactic problem; the lack of a recognizable compositional
147 interpretation is a result of the lack of a possible structural analysis. We would
148 therefore regard these examples as primarily syntactically unacceptable. This contrasts
149 with (1-b), which we regard as well-formed on both parameters, being merely
150 implausible, except in a small child's playroom, where a train being chewed is an
151 entirely normal situation (cf. Hahne & Friederici 2002).

152

153 2.3. Plausibility

154 Such examples highlight another problem in manipulating meaning as an experimental
155 variable: the human demand to make sense of linguistic forms. We associate possible
156 meanings with things that we can accept as being true or plausible. So 'the third-floor
157 apartment reappeared today', which is both syntactically and semantically flawless,
158 will cause irrelevant experimental effects since subjects will find it difficult to fit the
159 meaning into their mental model of the world. Zhou & Gao (2009) for example argue
160 that participants interpret *Every robber robbed a bank* in the surface scope reading
161 because it is more *plausible* that each robber robbed a different bank.

162 This links in to a wider discussion of the role of plausibility as a factor in semantic
163 processing and as a filter on possible readings. Zhou & Gao (2009) claim that such
164 doubly quantified sentences are ambiguous in Mandarin, since their experimental
165 evidence suggests that both interpretations are built up in parallel, but one reading is
166 subsequently filtered out by plausibility, which accounts for the contrary judgements in
167 work on semantic theory (e.g. Huang 1982, Aoun & Li 1989).

168

169 2.4. Meaning as a complex measure

170 The meaning of a structure is not fixed or unique, even when linguistic, social, and
171 discourse context are fixed. First, a single expression may have multiple readings,
172 which compete for dominance. Often a specific relevant reading of a structure needs to
173 be forced in an experiment. Some readings of theoretical interest may be quite
174 inaccessible, though nevertheless real. This raises the issue of expert knowledge, which
175 again contrasts with the situation in syntax. Syntactic well-formedness judgements are
176 generally available and accessible to any native speaker and require no expertise. On
177 the other hand, it can require specialist knowledge to ‘get’ some readings since the
178 access to variant readings is usually via different analyses. This is a crucial point in
179 semantics, since it reduces the likelihood that the intuitions of the naïve native speaker
180 can be the final arbiter in this field, as they can reasonably be argued to be in syntax
181 (Chomsky 1965). A fine example of this is from Hobbs & Schieber (1987):

182 (2) Two representatives of three companies saw most samples.

183 They claim that this sentence is five-ways ambiguous. Park (1995) however denies the
184 existence of one of these readings (*three > most > two*). It is doubtful whether this
185 question is solvable by asking naïve informants.

186 Even within a given analysis of a construction, the meaning may not be fully
187 determined. Aspects of meaning are left unspecified, which means that two different
188 perceivers can interpret a single structure in different ways. This too requires great care
189 and attention to detail when designing experiments which aim to be exact.

190

191 2.5. The observer’s paradox

192 A frequent aim in semantic experiments is to discover how subjects interpret linguistic

193 input under normal conditions. A constant problem is how experimenters can access
194 this information, because whatever additional task we instruct the subjects to carry out
195 renders the conditions abnormal. For example, if we ask them to choose which one of a
196 pair of pictures illustrates the interpretation that they have gathered, or even if we just
197 observe their eye movements, the very presence of two pictures is likely to make them
198 more aware that more than one interpretation is possible, thus biasing the results. Even a
199 single picture can alter or trigger the accessibility of a reading.

200

201 2.6. Inherent meaning and inferred meaning

202 One last linguistic distinction which we should note here is that between the inherent
203 meaning of an expression ("what is said") and the inferred meaning of a given utterance
204 of an expression. This distinction is fundamental in the division of research into
205 meaning into separate fields, but it is in practice very difficult to apply in experimental
206 work, since naïve informants do not naturally differentiate the two. The recent 'literal
207 Lucy' approach of Larson et al. (2010) is a promising solution to this problem; in this
208 paradigm participants must report how 'literal Lucy', who only ever perceives the
209 narrowly inherent meaning of utterances and makes no inferences, would understand
210 example sentences. This distinction is particularly important when an experimental
211 design requires a disambiguation, and extreme care must be taken that its content is not
212 only inferred. For example, in (3), it is implicated that every rugby player broke one of
213 their own fingers, but this is not necessarily the case. This example cannot thus offer
214 watertight disambiguation.

215 (3) Every rugby player broke a finger.

216 *Implication:* Every rugby player broke one of their own fingers.

217

218 2.7. Experimental measures and the object of theory

219 As a rule, semantic theory makes no predictions about semantic processing. Instead it
220 concerns itself with the final stable interpretation which is achieved after a whole
221 linguistic expression, usually at the sentence level, has been processed and all
222 reanalyses, for example as a result of garden paths, have been resolved. It
223 fundamentally concerns the stative, holistic result of the processing of an expression,
224 indeed many theoretical approaches regard meaning as only coming about in a full
225 sentence (cf. article 8 (Meier-Oeser) *Emergence of linguistic semantics*). But the
226 processing of a sentence is made up of many steps which are incremental and which
227 interact strongly with each other, partly predicting, partly parsing input as it arrives,
228 partly confirming or revising previous analyses. Much of the experimental evidence
229 available to us provides direct evidence only of these processing steps.

230 It thus follows that for many semantics practioners much of the empirical evidence
231 which we can gather concerns at best our *predictions* about what the sentence is going
232 to mean, not really aspects of its actual meaning. The time course of our arriving at a
233 particular reading, whether it be remote or readily accessible, has no direct implications
234 for the theory, since this makes no predictions about processing speed (cf. Phillips &
235 Wagers 2007). One aim of this article is to show that experimental techniques can
236 deliver data which can contribute to theory building.

237

238 2.8. Categorical predictions and gradient data

239 Predictions of semantic theories typically concern the *availability* of particular
240 interpretations. Experiments deliver more fine-grained data that reflect the relative

241 preferences among the interpretations. Mapping these gradient data onto the categorical
242 predictions, that is, drawing the line between still available and impossible readings is a
243 non-trivial task. At the same time, the ability to distinguish preferences among the
244 “intermediate” interpretations may be highly relevant for testing predictions concerning
245 readings that fall between the clearly available and the clearly impossible.

246

247 2.9. Outlook

248 In the remainder of this paper we will discuss two ways in which systematically
249 collected experimental data can contribute to semantic theorizing. We will use
250 quantifier scope as an example of a phenomenon where results of psycholinguistic
251 experiments can make significant contributions to the theoretical discussions. We will
252 not attempt to review here the considerable psycholinguistic literature on the processing
253 of quantifiers (for a comprehensive survey cf. article 103 (Frazier) *Meaning in*
254 *psycholinguistics*). Instead we will concentrate on a small set of studies that show the
255 usefulness of end-of-sentence judgements in establishing the available interpretations of
256 quantified sentences. Then we will sketch an experiment to address aspects of the
257 unfolding interpretation of quantifier scope which are of interest to theoretical
258 semanticists as well.

259

260 3. Off-line evidence for scope interpretation

261 Semantic theories are typically based on introspective judgements of a handful of
262 theoreticians. The judgements concern available readings of a sentence, possibly
263 ranked as to how easily available these readings are. Not surprisingly, judgements of
264 this sort are subtle and often controversial. For instance, the sentence *Everyone loves*

265 *someone* has been alternately considered to only allow the wide-scope universal reading
266 (e.g. Hornstein 1995; Beghelli & Stowell 1997) or to be fully ambiguous (May 1977,
267 1985; Hornstein 1984; Higginbotham 1985). Example (2) above illustrates the same
268 point. Park (1995) and Hobbs & Shieber (1987) disagree about the number of available
269 readings.

270 The data problem has been known for a long time. Studies as early as Ioup (1975) and
271 VanLehn (1978) tried to consider the intuitions of naïve speakers in developing an
272 empirically motivated theory. However, it has been clear from the beginning that
273 "obvious" tasks such as paraphrasing a presumably ambiguous doubly-quantified
274 sentence or asking informants to choose a (preferred) paraphrase is rather complex and
275 that linguistically untrained participants may not be able to carry them out reliably.

276 Another purely linguistic task has been problematic for a different reason. Researchers
277 have tried to combine the quantified sentence with a disambiguating continuation, as in
278 (4).

279 (4) Every kid climbed a tree.

280 (a) The tree was full of apples.

281 (b) The trees were full of apples.

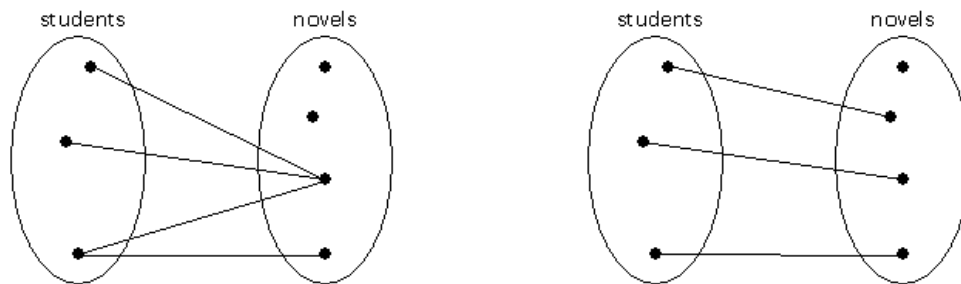
282 Disambiguation of this type was used by Gillen (1991), Kurtzman & MacDonald
283 (1993), Tunstall (1998) and Filik, Paterson & Liversedge (2004), for instance. Here the
284 plural continuation is only acceptable if multiple trees are instantiated, that is, the wide-
285 scope universal interpretation is chosen, whereas the singular continuation is intended to
286 only fit the wide-scope existential interpretation. Unfortunately the singular
287 continuation fails to disambiguate the sentence, as Tunstall (1998) points out: *the tree*
288 (4b) can easily be taken to mean *the tree the kid climbed*, thus making it compatible

289 with the wide-scope universal interpretation as well (see also Bott & Radó 2007 and
290 article 103 (Frazier) *Meaning in psycholinguistics*).

291 Problems of these kinds have prompted researchers to look for non-linguistic methods
292 of disambiguation. Gillen (1991) used, among other methods, simple pictures
293 resembling set diagrams. In her experiments subjects either drew diagrams to represent
294 the meaning of quantified sentences, chose the diagram that corresponded to the
295 (preferred) reading or judged how well the situation depicted in the diagram fitted the
296 sentence. Bott & Radó (2007) tested a somewhat modified form of the last of these
297 methods using diagrams like those in Figure 15.1. to see whether they constitute a
298 reliable mode of disambiguation that naïve informants can use easily. They found that
299 participants consistently delivered the expected judgements both for scopally
300 unambiguous quantified sentences (i.e. sentences where one scope reading was
301 excluded due to an intervening clause boundary) and for ambiguous quantified
302 sentences where expected preferences could be determined based on theoretical
303 considerations and corpus studies. These results show that there is no a priori reason to
304 exclude the judgements of non-linguist informants from consideration.

A) exactly one > each

B) each > exactly one



305 Figure 15.1: *DISAMBIGUATING DIAGRAMS FOR THE SENTENCE Exactly one novel was read*
306 *by each student.*

307

308 For informative experiments, however, we need to be able to derive testable hypotheses
309 based on existing semantic proposals. Although semantic theories are not formulated to
310 make predictions about processing, it is still possible to identify areas where different
311 approaches lead to different predictions concerning the judgement of particular
312 constructions. The interpretation of quantifiers provides an example here as well.

313 One possible way of classifying theories of quantifier scope has to do with the way
314 different factors are supposed to affect the scope properties of quantifiers. In
315 configurational models such as Reinhart (1976, 1978, 1983, 1995) and Beghelli &
316 Stowell (1997), quantifiers move to/are interpreted in different structural positions. A
317 quantifier higher in the (syntactic) tree will always outscope lower ones. The absolute
318 position in the tree is irrelevant; what matters is the position relative to the other
319 quantifier(s). While earlier proposals only considered syntactic properties of quantifiers,
320 Beghelli and Stowell also include semantic factors in the hierarchy of quantifier
321 positions. Taking *distributivity* as an example, assuming that a +dist quantifier is
322 interpreted in Spec,QP which is the highest position available for quantifiers, Q1 will

323 outscope Q2 if only Q1 is +dist, regardless of what other properties Q1 or Q2 may have.
324 An effect of other factors will only become apparent if neither of the quantifiers is +dist.
325 By contrast, the basic assumption in multi-factor theories of quantifier scope is that each
326 factor has a certain amount of influence on quantifier scope regardless of the presence
327 or absence of other factors (cf. Ioup 1975; Kurtzman & MacDonald 1993; Kuno 1991
328 and Pafel 2005). The effects of different factors can be combined, resulting in greater
329 or lesser preference for a particular interpretation. Theories differ in whether one of the
330 readings disappears when it is below some threshold, or whether sentences with
331 multiple quantifiers are always necessarily ambiguous.

332 Let us assume that the two scope-relevant factors we are interested in are distributivity
333 and discourse-binding, the latter indicated by the partitive NP *one of these N*, see (6).
334 Crossing these factors yields four possible combinations: +dist/+d-bound, +dist/-d-
335 bound, -dist/+d-bound, and -dist/-d-bound. In a configurational theory presumably there
336 will be a structural position reserved for discourse-bound phrases. Let us consider the
337 case where this position is lower than that for +dist, but higher than the lowest scope
338 position available for quantifiers. Thus Q1 should outscope Q2 in the first two
339 configurations, Q2 should outscope Q1 in the third, and the last one may in fact be fully
340 scope ambiguous unless some additional factors are at play as well. Moreover, as
341 configurational theories of scope have no mechanism to predict relative strength of
342 scope preference, the first two configurations should show the same size preference for
343 a wide-scope interpretation of Q1. In statistical terms, we expect an interaction: d-
344 binding should have an effect when Q1 is -dist, but not when it is +dist.

345 In multi-factor theories, on the other hand, the prediction would usually be that the
346 effects of the different factors should add up. That is, the difference in scope bias

347 between a d-bound and a non-d-bound +dist quantifier should be the same as between a
348 d-bound and a non-d-bound -dist quantifier. A given factor should be able to exert its
349 influence regardless of the other factors present.

350 Bott and Radó have been testing these predictions in on-going work. In two
351 questionnaire studies subjects read doubly-quantified German sentences and used
352 magnitude estimation to indicate how well disambiguating set diagrams fitted the
353 interpretation of the sentence. Experiment 1 manipulated distributivity and linear order
354 and used materials like (5). Experiment 2 tested the factors distributivity and d-binding
355 using sentences like (6).

356 (5) a. Genau einen dieser Professoren haben alle Studentinnen verehrt.
357 Exactly one these professors_{acc} have all female students adored.
358 *All female students adored exactly one of these professors.*

359 b. Genau einen dieser Professoren hat jede Studentin verehrt.
360 Exactly one these professors_{acc} has each female students adored.
361 *Each female student adored exactly one of these professors.*

362 c. Alle Studentinnen haben genau einen dieser Professoren verehrt.
363 All female students have exactly one these professors_{acc} adored.
364 *All female students adored exactly one of these professors.*

365 d. Jede Studentin hat genau einen dieser Professoren verehrt.
366 Each female student has exactly one these professors_{acc} adored.
367 *Each female student adored exactly one of these professors.*

368 (6) a. Genau einen Professor haben alle diese Studentinnen verehrt.
369 Exactly one professor_{acc} have all these female students adored.
370 *All of these female students adored exactly one professor.*

371 b. Genau einen dieser Professoren haben alle Studentinnen verehrt.

372 Exactly one these professors_{acc} have all female students adored.

373 *All female students adored exactly one of these professors.*

374 c. Genau einen Professor hat jede dieser Studentinnen verehrt.

375 Exactly one professor_{acc} has each these female students adored.

376 *Each of these female students adored exactly one professor.*

377 d. Genau einen dieser Professoren hat jede Studentin verehrt.

378 Exactly one these professors_{acc} has each female student adored.

379 *Each female student adored exactly one of these professors.*

380 Bott and Radó found clear evidence for the influence of all three factors. The
381 distributive quantifier *jeder* took scope more easily than *alle*, d-binding of a quantifier
382 and linear precedence both resulted in a greater tendency to take wide scope. Crucially,
383 the effects were additive, which is compatible with the predictions of multi-factor
384 theories but unexpected under configurational approaches.

385 These results show that even simple questionnaire studies can deliver theoretically
386 highly relevant data. This is particularly important in an area like quantifier scope,
387 where the judgements are typically subtle and not always accessible to introspection.

388 Of course the study reported here cannot address all possible questions concerning the
389 interpretation of quantified sentences like those in (5)-(6). It cannot for example clarify
390 whether the processor initially constructs a fully specified representation of quantifier
391 scope or whether it first builds only a underspecified structure which is compatible with
392 both possible readings, an outstanding question of much current interest in semantics.

393 The data that we have presented so far is off-line, in that it measures preferences only at
394 the end of the sentence, when its content has been disambiguated. In section 5 we

395 present an experimental design which will allow investigating the on-going (on-line)
396 processing of scope ambiguities. In the next section we relate the semantic issue of
397 underspecification to experimental data and predictions for on-line processing.

398

399 4. Underspecification vs. full interpretation

400 It is generally agreed that syntactic processing is *incremental* in nature (e.g. van Gompel
401 & Pickering 2007) i.e. a full-fledged syntactic representation is assigned to every
402 incoming word. Whether semantic processing is incremental in the strict sense, is far
403 from beyond dispute and still an empirical question. To formulate hypotheses about the
404 time-course of semantic processing, we will now look at the on-going debate in
405 semantic theory on underspecification in semantic representations.

406 Underspecified semantic representations are a tool intended to handle the problem of
407 ambiguity. The omission of parts of the semantic information allows one single
408 representation to be compatible with a whole set of different meanings (for an overview
409 of underspecification approaches, see e.g. Pinkal, 1999; articles 24 (Egg) *Semantic*
410 *underspecification* and 110 (Pinkal & Koller) *Semantics in computational linguistics*). It
411 is thus an economic method of dealing with ambiguity in that it avoids costly reanalysis,
412 used above all in computational applications.

413 Taking the psycholinguistic perspective, one would predict that constructing
414 underspecified representations in semantically ambiguous regions of a sentence avoids
415 processing difficulties in ambiguous regions and at the point of disambiguation.

416 Underspecification can be contrasted with an approach that assumes strict
417 incrementality and thus immediate full interpretation even in ambiguous regions. This
418 would predict processing difficulties in cases of disambiguations to non-preferred

419 readings. A candidate for a semantic processing principle guiding the choice of one
420 specified semantic representation would be a complexity-sensitive one (for example:
421 "Avoid quantifier raising" captured in Tunstall's *Principle of Scope Interpretation* 1998
422 and Anderson's 2004 *Processing Scope Economy*).

423 In the psycholinguistic investigation of coercion phenomena, the experimental evidence
424 is interpreted along these lines. Processing difficulties at the point of disambiguation are
425 taken as evidence for full semantic interpretation (see e.g. Piñango, Zurif & Jackendoff
426 1999; Todorova, Straub, Badecker & Frank 2000) whereas the lack of measurable
427 effects is seen as support for an underspecified semantic representation (see e.g.
428 Pyllkkänen & McElree 2006; Pickering, McElree, Frisson, Chen & Traxler 2006).

429 Analogously, in the processing of quantifier scope ambiguities, experimental evidence
430 for processing difficulties at the point of disambiguation will be interpreted as support
431 for full interpretation. However, this need not be taken as final. If we look at
432 underspecification approaches in semantics, non-semantic factors are mentioned which
433 might explain (and predict) difficulties in processing local scope ambiguities (see article
434 24 (Egg) *Semantic underspecification*, section 6.4.1.). And these are exactly the factors
435 which are assumed by multi-factor theories to have an impact on quantifier scope:
436 syntactic structure and function, context, and type of quantifier. The relative weighting
437 and interaction of these factors are not made fully explicit, however.

438 For the full picture, it would be necessary to examine not only the point of
439 disambiguation but also the ambiguous part of the input, for it is there that the effects of
440 these factors might be identified. Underspecification is normally only temporary,
441 however, and a full interpretation will presumably be constructed at some stage. This
442 might be recognizable for example in behavioural measures, but the precise predictions

443 of underspecification theory are not always clear. For example, it might be assumed
444 that even representations which are never fully specified by the input signal (or context)
445 do receive more specific interpretations at some later stage. This of course raises the
446 question what domains of interpretation are relevant here (sentence boundary, utterance,
447 ...). In the next section we present experimental work which may offer a starting point
448 for the empirical investigation of such issues.

449

450 5. On-line evidence for representation of scope

451 Given the underspecification view, relative scope should remain underspecified as long
452 as neither interpretation is forced. Indeed there should not even be any preference for
453 one reading. The results of the questionnaire studies reported in Section 3 already
454 indicate that this view cannot be right: A particular combination of factors was found to
455 systematically support a certain reading. Furthermore it is unlikely that the task itself
456 introduced a preference towards one interpretation -- although the diagram representing
457 the wide-scope existential reading was somewhat more complex, this did not seem to
458 interfere with participants' performance. The observed preferences must thus be due to
459 the experimental manipulation. That is, even if all possible interpretations are available
460 up to the point where disambiguating information arrives, there must be some inherent
461 ranking of the various scope-determining factors that results in certain interpretations
462 being more activated than others.

463 Off-line results such as those discussed above are thus equally compatible with two
464 different explanations; one where quantifier scope is fully determined (at least) by the
465 end of the sentence, and another one where several (presumably all combinatorially
466 possible) interpretations are available but weighted differently. A different

467 methodology is needed to find out whether there is any psycholinguistic support for an
468 underspecified view of quantifier scope.

469 As it turns out, the currently existing results of on-line studies cannot distinguish the
470 two alternatives, either. In on-line experiments a scope-ambiguous initial clause is
471 followed by a second one that is only compatible with one scope reading. An indication
472 of difficulty during the processing of the second sentence is typically taken as evidence
473 that the disambiguation is incompatible with the (sole) interpretation that had been
474 entertained up to that point. However, there is another way to look at such effects.
475 When the disambiguation is encountered, the underspecified representation needs to be
476 enriched to allow only one reading and exclude all others. It is conceivable that
477 updating the representation may require more or less effort depending on the ultimate
478 interpretation that is required.

479 This situation poses a dilemma for researchers investigating the interpretation of
480 quantifier scope. If explicit disambiguation is provided we can only test how easily the
481 required reading is available -- the results don't tell us what other reading(s) may have
482 been constructed. Without explicit disambiguation, however, reading time (or other)
483 data cannot be interpreted, since we do not know what reading(s) the participants had in
484 mind.

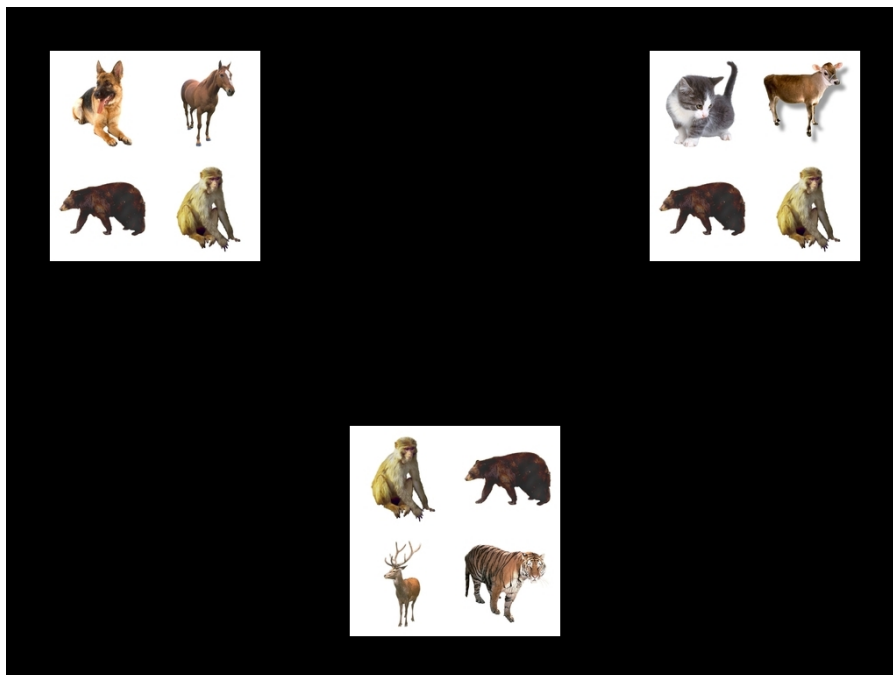
485 Bott & Radó (2009) approached this problem using eye-tracking while participants read
486 ambiguous sentences and then asking them to report the interpretation they computed.
487 Although the results they got are only partly relevant for the underspecification debate,
488 we will describe the experiment in some detail, since it provides a good starting point
489 for a more conclusive investigation. We will then sketch a modification of the method
490 that makes it possible to avoid some problems with the original study.

491 The scope-ambiguous sentences in Bott and Radó's study were instructions like those in
492 (7):

493 (7) a. Genau ein Tier auf jedem Bild sollst du nennen!
494 Exactly one animal on each picture should you name!
495 *Name exactly one animal from each picture!*

496 b. Genau ein Tier auf allen Bildern sollst du nennen!
497 Exactly one animal on all pictures should you name!
498 *Name exactly one animal from all pictures!*

499



500

501 Figure 15.2: *Display following inverse linking constructions.*

502

503 The first quantifier was always the indefinite *genau ein* “exactly one”. Q2 was either
504 distributive (*jeder*) or not (*alle*). In one set of control conditions Q1 was replaced by a
505 definite NP (*das Tier* “the animal”). In another set of control conditions the two
506 possible interpretations of (7) (one animal that is present in all fields vs. a possibly

507 different animal from each field on a display) were expressed by scope-unambiguous
508 quantified sentences, as in (8).

509 (8) a. Name exactly one animal that is found on all pictures.

510 b. From each picture name exactly one animal.

511 In each experimental trial participants first read one of these instruction sentences and
512 their eye-movements were monitored. Then the instruction sentence disappeared and a
513 picture display as in Figure 15.2. replaced it. Participants inspected this and had to
514 provide an answer within four seconds. Displays were constructed to be compatible
515 with both possible readings: a wide-scope universal one where different animals can be
516 selected from each field, as well as a wide-scope existential one where a particular
517 animal appeared in all fields (e.g. the monkey in Figure 15.2.). To make the quantifier
518 *exactly one* felicitous, the critical displays always allowed two potential answers for the
519 wide-scope existential interpretation.

520 The scope-ambiguous instructions were so-called inverse linking constructions, in
521 which the two quantifiers are contained within one NP. It has been assumed (e.g. May
522 & Bale 2006) that in inverse linking constructions the linearly second quantifier
523 preferentially takes scope over the first. The purpose of the study was to test this
524 prediction and to investigate to what extent the distributivity manipulation is able to
525 modulate it. Based on earlier results (Bott & Radó 2007) it was assumed that *jeder*
526 would prefer wide scope, which should further enhance the preference for the inverse
527 reading. When *alle* occurred as Q2, there should be a conflict between the preferences
528 inherent to the construction and those arising from the particular quantifiers.

529 The experimental setup made it possible to look at both the process of computing the
530 relative scope of the quantifiers (eye-movement behavior while reading the instructions)

531 and at the final interpretation (the answer participants gave) without providing any
532 disambiguation. Thus the answers could be taken to reflect the scope preferences at the
533 end of the sentence, whereas processing difficulty during reading would serve as an
534 indication that scope preferences are computed at a point where no decision is yet
535 required.

536 The off-line answers showed the expected effects. There was an overall preference for
537 the inverse scope reading, which was significantly stronger with *jeder* than with *alle*.
538 Crucially, the reading time data showed clear evidence of a conflict between the scope
539 factors: there was a significant slow-down at the second quantifier in (7b). The effect
540 was present already in first-pass reading times suggesting that scope preferences were
541 computed immediately. Bott and Radó interpret these results as strong indication that
542 readers regularly disambiguate sentences during normal reading.

543 However, this conclusion may be too strong. In Bott and Radó's experiment
544 participants had to choose a particular interpretation in order to carry out the
545 instructions (i.e. *name an animal*). Although they did not have to settle on that
546 interpretation while they were reading the instruction, they had to make a decision as to
547 the preferred reading immediately after the end of the sentence. This may have caused
548 them to disambiguate constructions that are typically left ambiguous during normal
549 interpretation.

550 Moreover, the instructions used in the experiment were highly predictable in structure:
551 they always contained a complex NP with two quantifiers (experimental items), a
552 definite NP1 followed by a quantified NP2 (fillers A), or else an unambiguous sentence
553 with two quantifiers. Although the content of NP1 (animal, vehicle, flag) and
554 distributivity of Q2 was varied, the rest of the instruction was the same: *sollst du nennen*

555 "you should name". This pattern was easy to recognize and may have resulted in a
556 strategy of starting to compute the scope preferences as soon as the second NP had been
557 received. To rule out this explanation Bott and Radó compared responses provided in
558 the first and the last third of each experimental session and failed to find any indication
559 of strategic behavior. Still the possibility remains that consistent early disambiguation
560 in the experiment resulted from the task of having to choose a reading quickly in order
561 to provide an answer. The ultimate test of underspecification would have to avoid such
562 pressure to disambiguate fast.

563 We propose a modification of Bott and Radó's experiment that may not only avoid this
564 pressure but actually encourage participants to delay disambiguation. In the proposed
565 experiment participants will have to judge the accuracy of sentences like those in (9):

566 (9) a. Genau eine geometrische Form auf allen Bildern ist rechteckig.

567 Exactly one geometrical shape on all pictures is rectangular.

568 *Exactly one geometrical shape on all pictures is rectangular.*

569 b. Genau eine geometrische Form auf jedem Bild ist rechteckig.

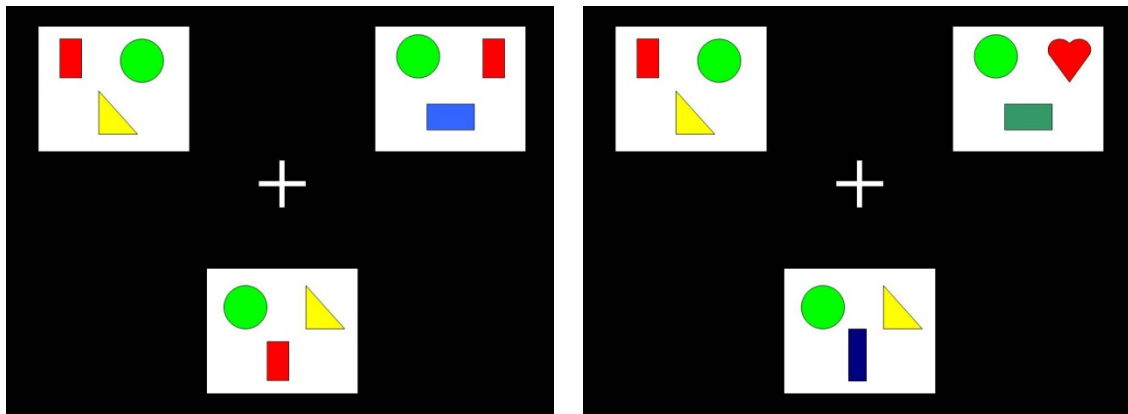
570 Exactly one geometrical shape on each picture is rectangular.

571 *Exactly one geometrical shape on each picture is rectangular.*

572

A) wide scope existential disambiguation

B) wide scope universal disambiguation



573 Figure 15.3: *Disambiguating displays in the proposed experiment*

574

575 The experiment procedure is as before. The sentences will be paired with unambiguous
 576 displays supporting either the wide-scope universal or the wide-scope existential
 577 reading (Figure 15.3.). In (9) full processing of the semantic content is not possible
 578 until the critical information (*rechteckig*) has been received. Since the display following
 579 the sentence is only compatible with one reading which the participant cannot
 580 anticipate, they are better off waiting to see which interpretation will be required for the
 581 answer. If underspecification is indeed the preferred strategy, there should be no
 582 difference in reading times across the different conditions, nor should there be any
 583 difficulty in judging any kind of sentence-display pair. Assuming immediate full
 584 specification of scope, however, we would expect the same pattern of results as in Bott
 585 and Radó's study: slower reading times in (9a) than in (9b) at the second quantifier, as
 586 well as slower responses to displays requiring the wide-scope existential interpretation,
 587 the latter presumably modulated by *distributivity* of Q2.

588 The experiment sketched above would be able to distinguish intermediate positions
 589 between the two extremes of complete underspecification and immediate full

590 interpretation. It is conceivable, for instance, that scope interpretation is only initiated
591 when the perceiver can be reasonably sure that they have received all (or at least
592 sufficient) information. This would correspond to the same reading time effects (and
593 same answering behavior) as predicted under immediate full interpretation, but the
594 effects would be somewhat delayed. Another possibility is an initial underspecification
595 of scope, but the construction of a fully specified interpretation at the boundary of some
596 interpretation domain such as the clause boundary. That would predict a complete lack
597 of reading time effects but answer times showing the same incompatibility effects as
598 under versions of the full interpretation approach.

599 It is worth emphasizing how this design differs from existing studies. First, it looks at
600 the ambiguous region and not just the disambiguation point. Second, it differs from
601 Filik, Paterson & Liversedge (2004), who also measured reading times in the
602 ambiguous region, but who used the kind of disambiguation that we criticized in section
603 3.

604

605 6. Conclusions

606 In this article we have attempted to show that experimentally obtained data can, in spite
607 of certain complicating and confounding factors, be of relevance to semantic theory and
608 provide both support for and in some cases falsification of its assumptions and
609 constructs. In section 2 we noted that the field of theoretical semantics has made less
610 use of experimental verification of its analyses and assumptions. We have seen that
611 there are some quite good reasons for this and laid out what some of the problematic
612 factors are. While some of these are shared to a greater or lesser degree with other

613 branches of linguistics, some of them are peculiar to semantics or are especially severe
614 in this case.

615 The main part of our paper reports a research programme addressing the issue of
616 relative scope in doubly quantified sentences. We present this work as an example of
617 the ways in which experimental approaches can contribute to the development of
618 theory. They also illustrate some of the practical constraints upon such studies. For
619 example, we have seen that clear disambiguation is not always easy to achieve, in
620 particular, it is difficult to achieve without biasing the interpretational choices of the
621 experiment participant. The use of eye-tracking and fully ambiguous picture displays is
622 a real advance on previous practice (Bott & Radó 2009).

623 Section 3 shows how experimental procedures which are simple enough for non-
624 specialist experimenters can nevertheless yield evidence of value for the development of
625 semantic theories: a carefully constructed and counter-balanced design can produce data
626 of sufficient quality to answer outstanding questions with some degree of finality. In
627 this particular case the configurational account of scope can be seen as failing to
628 account for data that the multi-factor account succeeds in capturing. The unsupported
629 account is demonstrated to need adaptation or development. Experimentation can make
630 the field of theory more dynamic and adaptive; an account which repeatedly fails to
631 capture evidence gathered in controlled studies and which cannot economically be
632 extended to do so will eventually need to be reconsidered.

633 In section 5 we lay out some experimental designs to provide evidence which
634 distinguishes between two accounts (section 4) of the way that perceivers deal with
635 ambiguity in the input signal: Underspecification vs. Full Interpretation. This is an
636 example of how processing data can under certain circumstances provide decisive

637 evidence which distinguishes between theoretical accounts. While it is often the case
638 that theory does not make any direct predictions about psycholinguistically testable
639 measures of processing, this is not always the case, and it may require the collaboration
640 of psycholinguists and semanticists to make these apparent.

641 We therefore argue for experimental linguists and semanticists to cooperate more and
642 take more notice of each other's work for their mutual benefit. Semanticists will gain
643 additional ways to falsify theoretical analyses or aspects of them, which can deliver a
644 boost to theory development. This will be possible, because experimenters can tailor
645 experimental methods, tasks, and designs to their specific requirements.

646 Experimenters for their part will benefit by having the questioning eye of the
647 semanticist look over their experimental materials, which will surely avoid many
648 experiments being carried out whose materials fail to uniquely fulfill the requirements
649 of the design. An example of this is the mode of disambiguation which we discussed in
650 section 3. Further to this, experimenters will doubtless be able to derive more testable
651 predictions from semantic theories, if they discuss the finer workings of these with
652 specialist semanticists. We might mention here the example of semantic
653 underspecification: can we find evidence for its psychological reality? Further
654 questions might be: if some feature of an expression remains underdetermined by the
655 input, how long can the representation remain underspecified? Is it possible for a final
656 representation of a discourse to have unspecified features and nevertheless be fully
657 meaningful?

658 We conclude, therefore, that controlled experimentation can provide a further source of
659 evidence for semantics. This data can under certain circumstances give a more detailed
660 picture of the states of affairs which theories aim to account for. This additional

661 evidence could be the catalyst for some advances in semantic theory and explanation, in
662 the same way that it has in syntactic theory.

7. References

Anderson, Catherine 2004. *The Structure and Real-time Comprehension of Quantifier Scope Ambiguity*. Ph.D. dissertation. Northwestern University.

Aoun, Joseph & Yen-hui Audrey Li 1989. Scope and constituency. *Linguistic Inquiry* 16, 623–637.

Beghelli, Filippo & Tim Stowell 1997. Distributivity and negation: The syntax of *each* and *every*. In: A. Szabolcsi (ed.). *Ways of Scope Taking*. Dordrecht: Kluwer, 71–107.

Bott, Oliver & Janina Radó 2007. Quantifying quantifier scope. In: S. Featherston & W. Sternefeld (eds.). *Roots. Linguistics in Search of its Evidential Base*. Berlin/New York: Walter de Gruyter, 53–74.

Bott, Oliver & Janina Radó 2009. How to provide exactly one interpretation for every sentence, or what eye movements reveal about quantifier scope. In: S. Winkler & S. Featherston (eds.). *The Fruits of Empirical Linguistics, Volume 1: Process*. Berlin/New York: Walter de Gruyter, 25–46.

Chomsky, Noam 1965. *Aspects of the Theory of Syntax*. Cambridge, Mass.: The MIT Press.

Filik, Ruth, Kevin B. Paterson & Simon P. Liversedge 2004. Processing doubly

quantified sentences: Evidence from eye movements. *Psychonomic Bulletin & Review* 11(5), 953–959.

Gillen, Kathryn 1991. *The Comprehension of Doubly Quantified Sentences*. Ph.D. Dissertation. University of Durham.

van Gompel, Roger P.G., & Martin J. Pickering 2007. Syntactic parsing. In: G. Gaskell (ed.). *The Oxford Handbook of Psycholinguistics*. Oxford: Oxford University Press. 455–504.

Hahne, Anja & Angela D. Friederici 2002. Differential task effects on semantic and syntactic processes as revealed by ERPs. *Cognitive Brain Research* 13, 339–356.

Higginbotham, James 1985. On semantics. *Linguistic Inquiry* 16, 547–594.

Hobbs, Jerry & Stuart M. Shieber 1987. An algorithm for generating quantifier scopings. *Computational Linguistics* 13, 47–63.

Huang, Cheng-Teh James 1982. *Logical Relations in Chinese and the Theory of Grammar*. PhD dissertation. MIT.

Hornstein, Norbert 1984. *Logic as Grammar*. Cambridge, MA: The MIT Press.

Hornstein, Norbert 1995. *Logical Form: From GB to Minimalism*. Oxford: Blackwell.

Ioup, Georgette 1975. *The Treatment of Quantifier Scope in Transformational Grammar*. Ph.D. dissertation. University of New York.

Kuno, Susumu 1991. Remarks on quantifier scope. In: H. Nakajima (ed.). *Current English Linguistics in Japan*. Berlin: Mouton de Gruyter, 261–287.

Kurtzman, Howard S. & Maryellen C. MacDonald 1993. Resolution of quantifier scope ambiguities. *Cognition* 48, 243–279.

Larson, Meredith, Ryan Doran, Yaron McNabb, Rachel Baker, Matthew Berends, Alex Djalali & Gregory Ward 2010. Distinguishing the said from the implicated using a novel experimental paradigm. In: U. Sauerland & K. Yatsushiro (eds.). *Semantics and Pragmatics. From Experiment to Theory*. Houndmills: Palgrave Macmillan.

May, Robert 1977. *The Grammar of Quantification*. Ph.D. dissertation. MIT. Reproduced in 1982. Indiana University Linguistics Club.

May, Robert 1985. *Logical Form: Its Structure and Derivation*. Cambridge, MA: The MIT Press.

May, Robert & Alan Bale 2006. Inverse linking. In: M. Everaert & H. van Riemsdijk (eds.). *Blackwell Companion to Syntax*. Oxford: Blackwell, chap. 36, 639–667.

Pafel, Jürgen 2005. *Quantifier scope in German*. Vol. 84 of *Linguistics Today*. Amsterdam: John Benjamins.

Park, Jong C. 1995. Quantifier scope and constituency. In: H. Uszkoreit (ed.). *Proceedings of the 33rd Annual Meeting of the Association of Computational Linguistics*. Boston, Palo Alto, CA: Morgan Kaufmann, 205–212.

Phillips, Colin & Matthew Wagers 2007. Relating structure and time in linguistics and psycholinguistics. In: M.G. Gaskell (ed.). *The Oxford Handbook of Psycholinguistics*. Oxford: Oxford University Press, 739–756.

Pickering, Martin J., Brian McElree, Steven Frisson, Lilian Chen & Matthew J. Traxler 2006. Underspecification and aspectual coercion. *Discourse Processes* 42(2), 131–155.

Piñango, Maria, Edgar Zurif & Ray Jackendoff 1999. Real-time processing implications of enriched composition at the syntax-semantics interface. *Journal of Psycholinguistic Research* 28, 395–414.

Pinkal, Manfred 1999. On semantic underspecification. In: H. Bunt & R. Muskens (eds.). *Computing Meaning*, Dordrecht: Kluwer, 33–55.

Pylkkänen, Liina & Brian McElree 2006. The syntax-semantics interface: On-line composition of meaning. In: M. A. Gernsbacher & M. Traxler (eds.). *Handbook of Psycholinguistics*, 2. Edition. New York: Elsevier, 537–577.

Reinhart, Tanya 1976. *The Syntactic Domain of Anaphora*. Ph.D. dissertation. MIT.

Reinhart, Tanya 1978. Syntactic domains for semantic rules. In: F. Guenther & S.J. Schmidt (eds.). *Formal Semantics and Pragmatics for Natural Language*. Dordrecht: Reidel, 107–130.

Reinhart, Tanya 1983. *Anaphora and Semantic Interpretation*. London/Sydney: Croom Helm.

Reinhart, Tanya 1995. *Interface Strategies*. OTS Working Papers in Linguistics.

Todorova, Marina, Kathy Straub, William Badecker & Robert Frank 2000. Aspectual coercion and the online computation of sentential aspect. In: L. R. Gleitman & A. K. Joshi (eds.). *Proceedings of the 22nd Annual Conference of the Cognitive Science Society*. Mahwah, N.J.: Lawrence Erlbaum Associates, 3–8.

Ross, John R. 1970. On declarative sentences. In: R. Jacobs & P. Rosenbaum (eds). *Readings in English Transformational Grammar*. Waltham, Massachusetts: Ginn, 222–272.

Tunstall, Susanne L. 1998. *The Interpretation of Quantifiers: Semantics and Processing*. Ph.D. Dissertation. University of Massachusetts Amherst.

Weskott, Thomas & Gisbert Fanselow 2009. Scaling issues in the measurement of linguistic acceptability. In: S. Winkler & S. Featherston (eds.). *The Fruits of Empirical Linguistics, Volume 1: Process*. Berlin/New York: Walther de Gruyter, 229–246.

VanLehn, Kurt A. 1978. *Determining the Scope of English Quantifiers*. Technical Report (AI-TR 483). Artificial Intelligence Laboratory, MIT.

Zhou, Peng & Liqun Gao 2009. Scope processing in Chinese. *Journal of Psycholinguistic Research* 38, 11–24.

Oliver Bott, Tübingen (Germany)

Sam Featherston, Tübingen (Germany)

Janina Radó, Tübingen (Germany)

Britta Stolterfoht, Tübingen (Germany)