

PNAS



1

2 **Supporting Information for**

3 **Even Lawyers Don't Like Legalese**

4 **Eric Martínez, Francis Mollica and Edward Gibson**

5 **Eric Martínez.**

6 **E-mail: ericmartmit.edu**

7 **This PDF file includes:**

8 Supporting text

9 Figs. S1 to S5

10 Table S1

11 SI References

12 Supporting Information Text

13 1. Experiment 1

14 **A. Author Recognition Test Analyses.** As noted in the main text, for Experiment 1 we administered the Author Recognition
15 Test (a validated proxy for reading ability) to participants as part of each trial. Although this was mainly used as a filler task,
16 and we made no pre-registered predictions regarding the results, for transparency we report the results of its potential effect on
17 comprehension and recall data. The results of all of these analyses were convergent with the findings in Martinez, Mollica &
18 Gibson's sample of laypeople (1).

19 When adding ART score as a fixed-effect predictor to our comprehension model, we found a main effect of ART score on
20 comprehension ($\beta = .221, SE = .078, p = .005$), such that those who had higher ART scores had higher comprehension scores.
21 However, we did not find a significant interaction between ART and register ($p = .075$). When adding ART score as a fixed-effect
22 predictor to our recall model, we did not find a main effect of ART score on recall ($p = .787$).

23 **B. Recall annotation details.** Two trained research assistants coded whether a proposition was successfully recalled, using the
24 same method as (1). In particular, they were given a participant's retelling of a passage and then asked whether each legally
25 relevant proposition of the passage was (a) fully recalled; (b) partially recalled; or (c) not recalled. Coders were told that for a
26 response to count as "fully recalled," it did not have to be recalled verbatim (i.e. they can use their own words or syntax), so
27 long as they were confident that the meaning of what subject wrote is the same as the proposition.

28 For example, suppose the original text said "A court in Boston will resolve the dispute," and the participant wrote "something
29 will be resolved by a court." When coding responses, a coder might see three propositions that say: (i) "A court in Boston," (ii)
30 "will resolve," and (iii) "the dispute."

31 For (i), the coder would put a 0.5 for "partially recalled" (since "in Boston" was missing from "a court"); for (ii), the coder
32 would put a 1 for "fully recalled" (since "will be resolved" means basically the same thing as "will resolve"), and for (iii), the
33 coder would put a 0 for "not recalled" (since "the dispute" was not in the response).

34 To reduce potential bias, coders were unaware of whether a participant had seen or recalled the simple or legalese version of
35 a text. Moreover, the rubric that the coders used to score participants' responses was the legalese version: for each proposition,
36 they were given the language of that proposition in legalese, and were told to score a participants' response as having recalled
37 that proposition if it had language that had the same meaning as that proposition. Thus, any differences in recall favoring the
38 Plain English version would arise in spite of the coding bias, which was towards the Legalese version.

39 Of the roughly 650 retellings within the lawyer data, one coder was responsible for coding 100 percent of the retellings,
40 while the second coder was responsible for coding a random subset (20%) of the retellings, so as to assess inter-rater reliability.
41 Coder reliability was assessed with Cohen's kappa coefficient (2, 3).

42 We adjudicated ties as follows: (i) a tie between one "fully recalled" judgment and one "not recalled" judgment resulted in
43 a final "partially recalled" judgment; (ii) a tie between one "fully recalled" judgment and one "partially recalled" judgment
44 resulted in a final "fully recalled" judgment for a given proposition; and (iii) a tie between one "partially recalled" judgment
45 and one "not recalled" judgment resulted in a final "not recalled" judgment. For our regression analyses, we perform both a
46 conservative analysis (recoding "partially recalled" as "not recalled") and an anti-conservative analysis (recoding "partially
47 recalled" as "fully recalled"). Our results do not qualitatively change, so we will only report the conservative analysis in the
48 main writeup.

49 **C. Anti-conservative recall analysis.** As noted in the main writeup, for our recall analyses, we performed both a conservative
50 analysis and an anti-conservative analysis. For the conservative analysis, As our results did not qualitatively change, so we only
51 reported the conservative analysis in text.

52 For both the conservative and anti-conservative analyses, we conducted a mixed effect logistic regression with register, legal
53 training and the interaction between the two as fixed effects, and participant as random effects, with register as a random
54 slope for each. As in the conservative analysis, for the anti-conservative analysis, we found a main effect of register ($\beta = .225$,
55 $SE = .098, p = .022$) and legal training ($\beta = -.340, SE = .153, p = .026$) on recall. As in the conservative analysis, we did not
56 find a main effect of the interaction between register and condition ($p = .174$).

57 **D. Subjective rating analyses.** As noted in the main writeup, in addition to recall and comprehension analyses, we also asked
58 participants to rate how difficult a text they found the text (a) for themselves; (b) for the average layperson; and (c) for the
59 average lawyer.

60
61 Below is the wording of each of the prompts:

- 62 • (you): "How complex/difficult do you find this text to understand?"
- 63 • (lawyer): "How complex/difficult do you think the average layperson/non-lawyer would find this text to understand?"
- 64 • (layperson): "How complex/difficult do you think the average lawyer would find this text to understand?"

65 The 5 answer choices for each prompt were as follows:

- 66 • extremely simple/easy

- 67 • somewhat simple/easy
- 68 • neither complex/difficult nor simple/easy
- 69 • somewhat complex/difficult
- 70 • extremely complex/difficult

71 Results are visualized in the main text.

72 To analyze these results, we ran three different models.

73 First, we conducted a model that compared how difficult lawyers and laypeople predicted texts would be for the average
74 layperson.

75 Second, we conducted a model that compared (a) how difficult lawyers predicted texts would be for the average layperson
76 compared with (b) how difficult lay participants perceived the texts to be for themselves.

77 Third, we conducted a model that compared (a) how difficult laypeople predicted texts would be for the average lawyer
78 compared with (b) how difficult lay participants perceived the texts to be for themselves.

79 This model's predictions are less relevant for the curse of knowledge hypothesis, but for robustness purposes we report it
80 anyway.

81 The model for all three models was as follows:

- 82 • `clmm(as.factor(Response) ~ condition*training + (1 + condition | subject) + (1 + condition | item), data = .)`

83 The only difference among the three models was that the “response” variable was filtered to include a different subset of the
84 data (to include the relevant conditions)

85 For the first model, we found a main effect of condition ($\beta = -1.784$, $SE = .161$, $p < .001$), but not training ($p = .974$), nor the
86 interaction between condition and training ($p = .127$).

87 That is, contrary to the predictions of the curse of knowledge hypothesis, we did not find evidence that lawyers underestimated
88 the difficulty of legal texts for non-lawyers, nor did they particularly underestimate the difficulty of legal texts written in
89 legalese.

90 For the second model, we found a main effect of condition ($\beta = -1.784$, $SE = .161$, $p < .001$). We also found an effect of training
91 ($\beta = -.938$, $SE = .293$, $p = .001$), with laypeople's own subjective ratings being significantly easier than lawyers' predictions of the
92 average layperson's subjective ratings.

93 We also found an interaction between training and condition ($\beta = .562$, $SE = .250$, $p = .025$), such that laypeople's ratings of
94 simple texts were disproportionately high relative to lawyers' predictions of those texts relative to the groups' legalese ratings.

95 For the third model, we found a main effect of condition ($\beta = -2.235$, $SE = 8.718$, $p < .0001$). We also found an effect of
96 training ($\beta = -3.094$, $SE = .339$, $p < .0001$), with laypeople's predictions of lawyers' ratings being significantly easier than lawyers'
97 subjective ratings of the texts.

98 We also found an interaction between training and condition ($\beta = .792$, $SE = .266$, $p = .003$), such that laypeople's predictions
99 of lawyers' ratings of simple texts were disproportionately high relative to lawyers' ratings of those texts relative to the
100 legalese ratings.

101 2. Experiment 2

102 **Hypotheses and Predictions.** In Experiment 2, we aimed to test the following hypotheses and associated predictions. All of
103 these were pre-registered on OSF.

104 Hypothesis I: Lawyers simply write in a complex register out of “habit, laziness” (4) or “tradition” (5); they “copy and
105 paste” (Adams, 2022) from existing templates with old, complicated terms because that's the “quickest and cheapest way to
106 produce a contract” (6), not out of any preference.

- 107 • Prediction 1: Lawyers will rate plain English contracts as of equal quality as legalese contracts.
- 108 • Prediction 2: Lawyers will agree to sign off on a contract written in Plain English.

109 Hypothesis II: Lawyers write in legalese in order to be accepted by peers. The legalese signals in-group membership (5).

- 110 • Prediction 1: Lawyers will rate contracts written in legalese as sounding more “lawyerly” (more appropriate/suitable for
111 a lawyer) than those written in plain English.
- 112 • Prediction 2: Lawyers will rate authors of contracts written in legalese more hireable than authors of contracts written in
113 plain English.

114 Hypothesis III: Lawyers write in legalese as a way of “preserving their monopoly” (7) on legal services and “justifying fees” (5)

- 115 • Prediction: Lawyers will predict contracts written in legalese as being more likely to be signed by clients than contracts
116 written in plain English.

117 Hypothesis IV: Contractual language needs to be complex in order to convey complex legal concepts in a way that “is far more
118 precise than ordinary language” (4) and/or to be enforceable

- 119 • Prediction: Plain English contracts will be rated as unenforceable by legal experts

120 **Materials.** To evaluate our predictions, we measured two sets of outcome variables. The first set of outcome variables were
121 measured individually for each text and were as follows:

- 122 • text quality (“How would you rate the overall quality of the above contract excerpt?”)
- 123 • enforceability (“Suppose two parties signed a contract that included the above excerpt. Would the excerpt likely be
124 legally enforceable (assuming the rest of the contract was enforceable)?”)
- 125 • useability (“Suppose someone at your firm drafted a contract that included the above excerpt. Would you and your firm
126 agree to execute it as currently written (assuming the rest of the contract was okay)?”)
- 127 • hireability (“Suppose the excerpt was drafted by someone outside your firm. Would your firm be likely to hire them to
128 draft future contracts, all else equal?”)
- 129 • lawyerliness (“Does the style/tone of the excerpt sound appropriate for a lawyer?”)
- 130 • likelihood of being signed (“Suppose you drafted this excerpt for a client as part of a larger contract. Would a client be
131 likely to sign this contract (assuming the rest of the contract was written in a similar style)?”)

132 Text quality was measured on a scale of 1-5 (1 being “extremely low-quality” and 5 being “extremely high-quality”). All
133 other outcome variables in this set were measured on a yes-no scale.

134 The second set of outcome variables were measured for each contract pair as opposed to each individual contract. These
135 variables were as follows:

- 136 • more usable (“Which of the two versions would you be more likely to execute, given the choice?”)
- 137 • more likely to be signed (“Which of the two versions would a client be more likely to agree to sign?”)
- 138 • more lawyerly (“Which of the two versions sounds more appropriate for a lawyer?”)
- 139 • more hireable (“Suppose the two versions were drafted by two different authors. Which of the two would your firm be
140 more likely to hire to draft future contracts, all else equal?”)

141 All outcome variables in this set were measured on a two-point scale (version 1 or version 2).

142 **A. Analysis Plan.** To evaluate Hypothesis I, we conducted two regressions.

- 143 • An ordinal regression with the following syntax: `clmm(text quality ~ condition + (1 + condition | item) + (1`
144 `+ condition | subject), data = .)`
- 145 • A logistic regression with the following syntax: `glmer(is usable ~ condition + (1 + condition | item) + (1 +`
146 `condition | subject), data = ., family = binomial(link = "logit"))`

147 For Hypothesis II we we conducted exact binomial tests for the more lawyerly and more hireable variables.

148
149 For Hypothesis III we conducted the following logistic regression, as well as an exact binomial test:

- 150 • `glmer(client would sign ~ condition + (1 + condition | item) + (1 + condition | subject), data = ., family`
151 `= binomial(link = "logit"))`

152 For Hypothesis IV, we conducted the following regression:

- 153 • `glmer(is enforceable ~ condition + (1 + condition | item) + (1 + condition | subject), data = ., family`
154 `= binomial(link = "logit"))`

155 In our pre-registration, we stated that if we encountered issues fitting models, we would use Bayesian regression techniques
156 with similar syntax. We did not encounter issues fitting models, and therefore will report our pre-registered models.

157 **B. Supplementary Results.** Results are visualized in Table S1 and Figure S1. As noted in the main text, all of the predictions
158 of hypotheses 1-3 were disconfirmed, and all of the predictions of hypothesis 4 were confirmed.

159 With regard to hypothesis 1, contrary to the first prediction, we found that lawyers were more likely to say that they would
160 use simple contracts over legalese contracts ($\beta = 1.432, SE = .270, p < .001$), and rated simple texts as higher quality than legalese
161 texts ($\beta = 1.705, SE = .329, p < .001$).

162 With regard to hypothesis 2, contrary to both predictions, participants were more likely to rate the authors of simple texts
163 as hireable compared to the authors of legalese texts ($\beta = 1.835, SE = .318, p < .001$), and we did not find participants to be
164 more likely to rate legalese texts as sounding more lawyerly than simple texts ($p = .692$).

165 With regard to hypothesis 3, contrary to the pre-registered prediction, we found that participants were more likely to predict
166 that clients would sign a simple contract compared to a legalese contract ($\beta = 1.232, SE = .338, p < .001$).

167 With regard to hypothesis 4, in line with the predictions, we found that participants were more likely to say that they would
168 agree to use the simple contracts as written ($\beta = 1.705, SE = .329, p < .001$), and we did not find a significant difference in how
169 enforceable the different contracts were rated as ($p = .717$).

170 **C. Exploratory Demographic Analyses.** Our main analyses, predictions, and conclusions drawn on the basis of our analyses
171 were limited to those that we included in our pre-registration on OSF.

172 Although our pre-registered statistical models did not include any demographic analyses, one might wonder whether our
173 results may have been driven by the demographic composition of our lawyer sample.

174 Below is a description of each of these results as applied to our hypotheses. To help lend a visual sense of the robustness of
175 the results, results for Experiment 2 limited to “experienced” attorneys (those with 10 or more years of practice experience)
176 are visualized in Figure S4. Results limited to “fancy” attorneys (those who attended a top-25 law school or work at a top-200
177 law firm) are visualized in Figure S5.

178 To more rigorously account for the possibility of results being driven by demographic factors, we conducted additional
179 versions of our pre-registered analyses, adding each of our demographic variables as fixed-effect predictors. Doing so did not
180 alter our results.

181 Specifically, in cases where our pre-registered models found a main effect of a given predictor variable, the same was true
182 when adding all of the demographic variables as additional fixed-effect predictors.

183 Conversely, in cases where our pre-registered models did not find a main effect of a given predictor variable, the same was
184 true when adding all of the demographic variables as additional fixed-effect predictors.

185 **C.1. Curse of Knowledge.** With regard to the curse-of-knowledge hypothesis, We added the demographic factors to a modified
186 model of comprehension and recall models that (a) were limited to lawyers; and (b) did not contain fixed-effects related to
187 legal training. These models were as follows:

188 • `glmer(comprehension ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition`
189 `| item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

190 • `glmer(recall ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition`
191 `| item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

192 As with our pre-registered models, these models revealed a main effect of register, such that lawyers were significantly
193 more likely to comprehend ($\beta = .373, SE = .094, p < .0001$) and recall ($\beta = .376, SE = .132, p < .0001$) legal content written in simple
194 contracts relative to legalese contracts.

195 The comprehension model revealed a main effect of “fanciness,” such that fancy lawyers had significantly higher comprehension
196 overall than non-fancy lawyers ($\beta = .433, SE = .158, p = .006$). There were no other main effects in the two models.

197 **C.2. In-Group Signaling.** With regard to the in-group signaling hypothesis, We added the demographic factors to our model of
198 hireability, as follows:

199 • `clmm(hireability ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition`
200 `| item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

201 As with our pre-registered model, this model revealed a main effect of register, such that lawyers rated authors of simple
202 contracts as significantly more hireable than authors of legalese contracts ($\beta = 1.900, SE = .322, p < .0001$).

203 **C.3. It's Just Business.** With regard to the it's just business hypothesis, We added the demographic factors to our model of
204 willingness to sign, as follows:

205 • `clmm(client would sign ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition`
206 `| item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

207 As with our pre-registered model, this model revealed a main effect of register, such that lawyers predicted that clients
208 would be significantly more likely to sign simple contracts than legalese contracts ($\beta = 1.208, SE = .370, p = .001$).

209 **C.4. Complexity of Information.** With regard to the complexity of information hypothesis, We added the demographic factors to
210 our model of enforceability, as follows:

211 • `clmm(enforceability ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition`
212 `| item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

213 As with our pre-registered model, this model revealed no effect of register. That is, we did not find evidence that lawyers
214 rated legalese contracts as more enforceable than legalese contracts ($p = .156$).

215 **C.5. Copy and Paste.** We added the demographic factors to our models of quality and usability, as follows:

216 • `clmm(quality ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition |`
217 `item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

218 • `clmm(would use ~ condition + age + is fancy + gender + ethnicity + practice experience + (1 + condition`
219 `| item) + (1 + condition | subject), data = ., family = binomial(link = "logit"))`

220 In both cases, we still found a main effect of register on responses, such that lawyers were significantly more likely to
221 say that they would use simple contracts over legalese contracts ($\beta = 1.533, SE = .285, p < .0001$), and rated simple contracts as
222 significantly higher quality than legalese contracts ($\beta = 1.807, SE = .338, p < .0001$).

Table S1. Endorsement rates by desiderata

Desiderata	Legalese			Simple		
	endorsement %	lower CI	upper CI	endorsement %	lower CI	upper CI
hireability	31.4	28.8	33.9	59.4	56.7	62.1
likelihood of being signed	69.0	66.4	71.4	82.2	80.2	84.4
enforceability	82.3	80.2	84.4	84.3	82.2	86.2
quality	2.61	2.55	2.67	3.33	3.27	3.39

Consider the below contract excerpt, written in blue.

This agreement, by whose terms hereinafter set forth Acoustic Acapella and Elmer's Entertainment Entity, said parties being hereinafter referred to as "Artists" and "Tour," respectively, hereby agree to be bound, has been formed by both parties on this date of december 1, 2019. It is assented to by both parties that a series of twelve concerts, the percentage of revenue of which being divided evenly between the two parties and Artists' share being apportioned among members pro rata with existing shareholdings, shall be performed by Artists on the third Saturday of every month commencing January 2020 through December 2020.

How would you rate the overall quality of the above contract excerpt?

extremely high-quality

somewhat high-quality

neither high-quality nor low-quality

somewhat low-quality

extremely low-quality

Suppose two parties signed a contract that included the above excerpt. Would the excerpt likely be legally enforceable (assuming the rest of the contract was enforceable)?

Yes

No

Fig. S1. Interface of Experiment II.

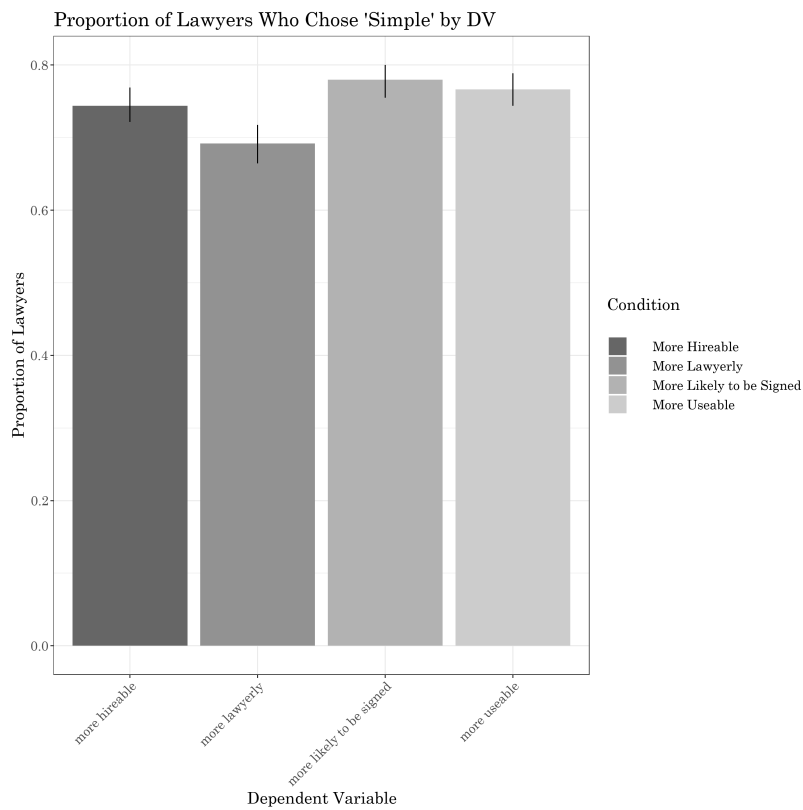


Fig. S2. Proportion of lawyers who endorsed simple version over legalese version according to different desiderata.

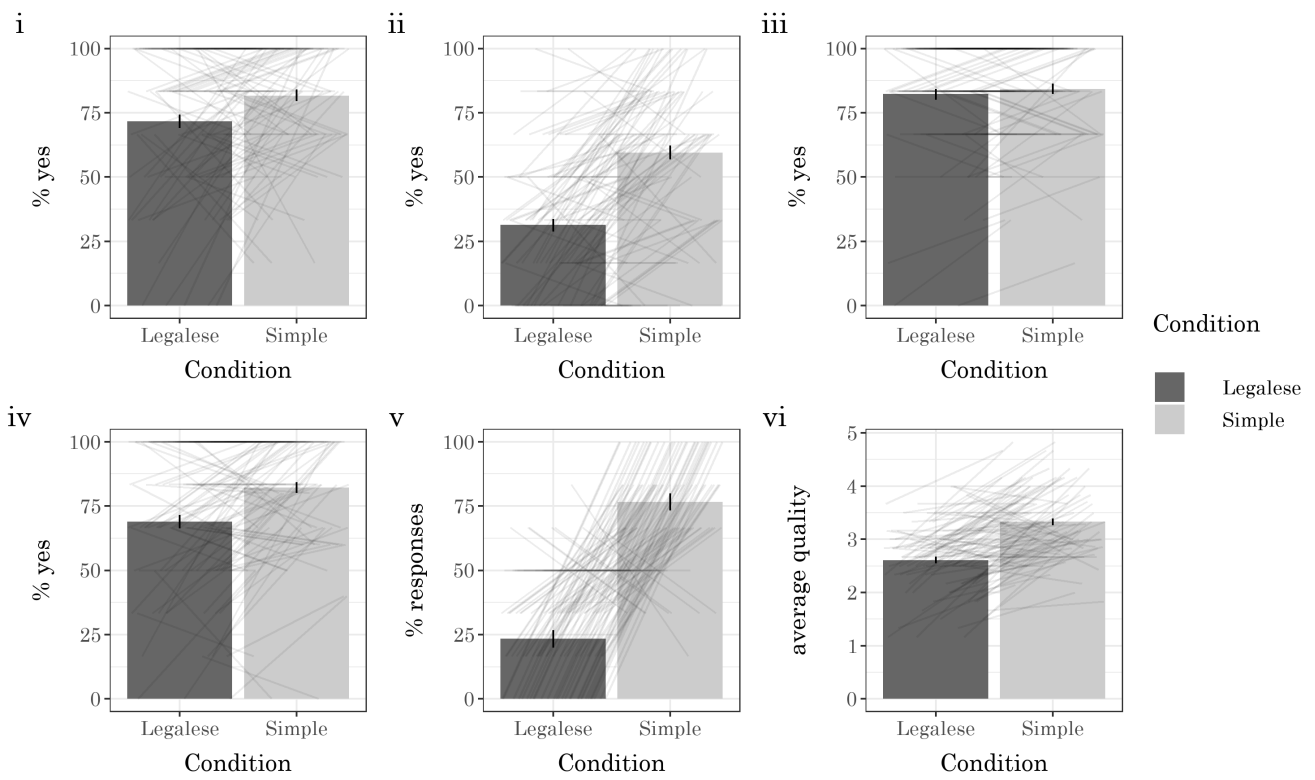


Fig. S3. Proportion of lawyers who endorsed simple and legalese contracts according to different desiderata.

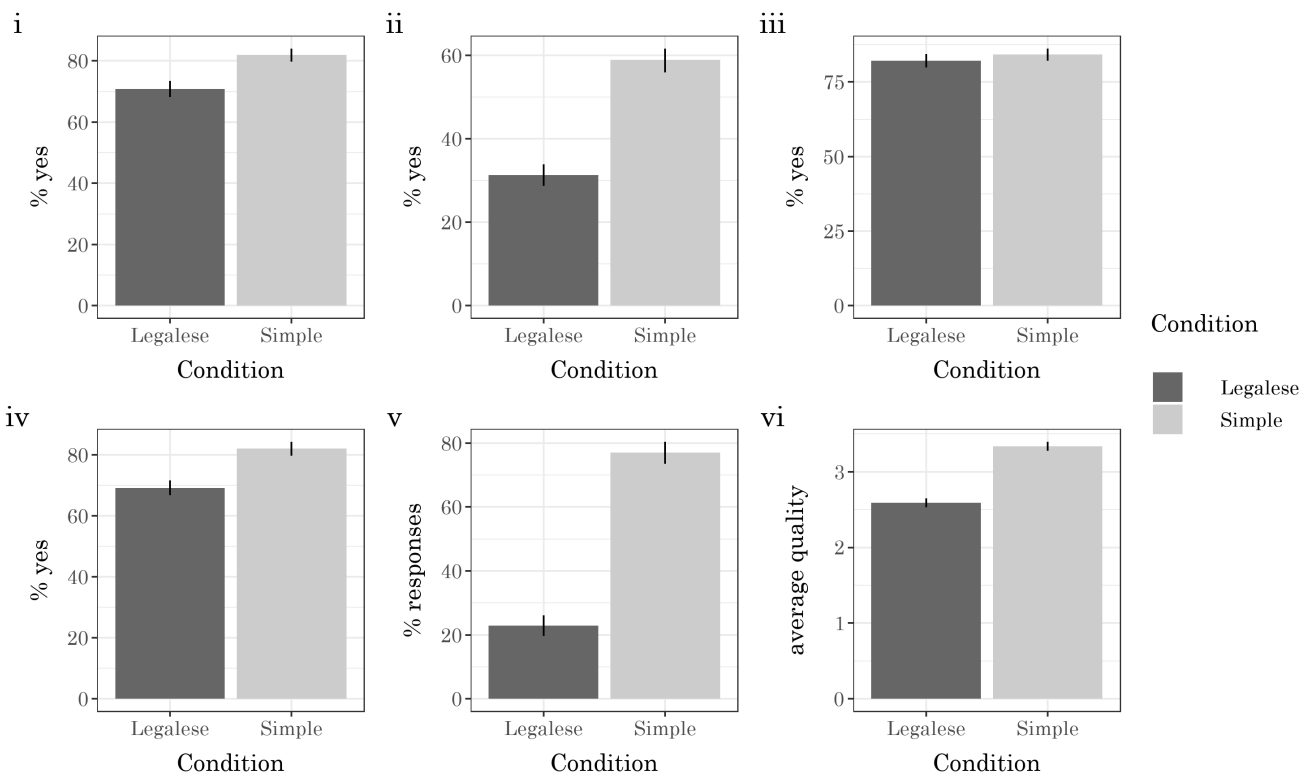


Fig. S4. Proportion of experienced lawyers who endorsed simple and legalese contracts according to different desiderata.

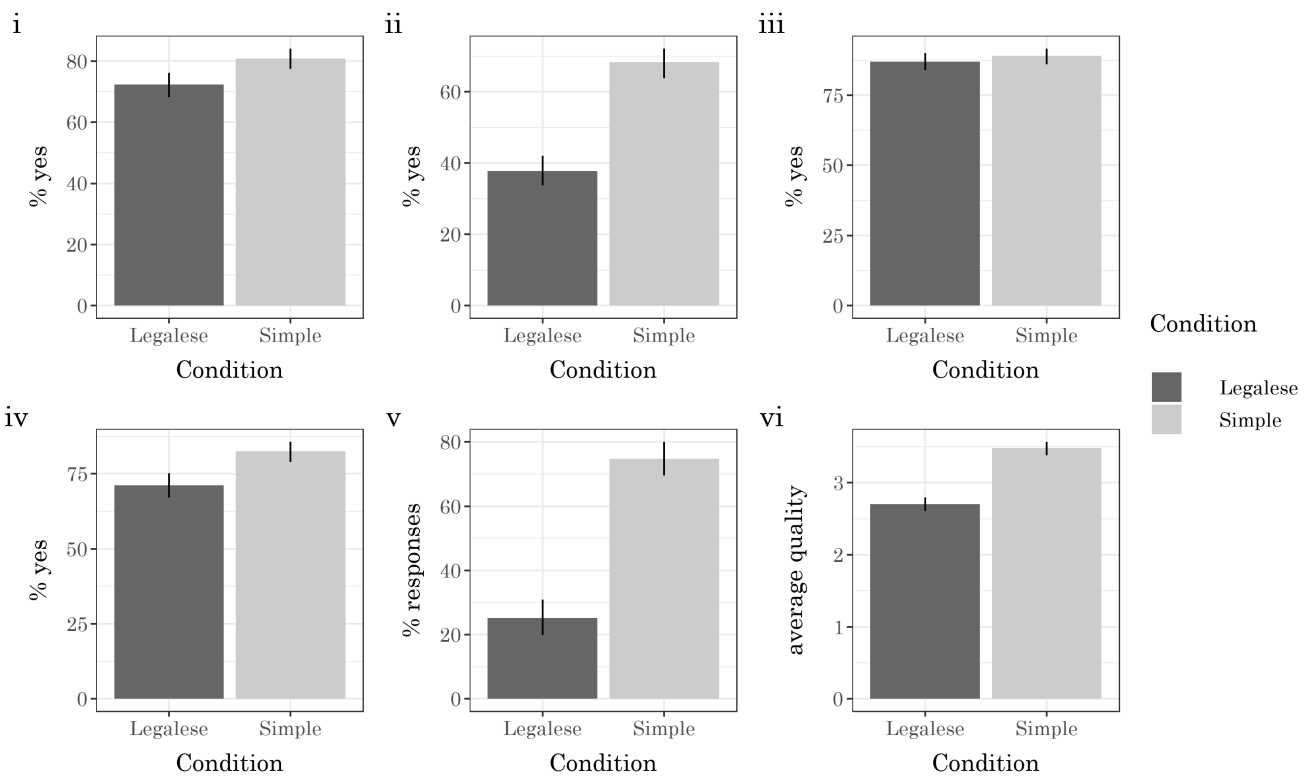


Fig. S5. Proportion of fancy lawyers who endorsed simple and legalese contracts according to different desiderata.

223 **References**

- 224 1. E Martínez, F Mollica, E Gibson, Poor writing, not specialized concepts, drives processing difficulty in legal language.
225 *Cognition* **224**, 105070 (2022).
- 226 2. J Cohen, A coefficient of agreement for nominal scales. *Educ. psychological measurement* **20**, 37–46 (1960).
- 227 3. ML McHugh, Interrater reliability: the kappa statistic. *Biochem. medica* **22**, 276–282 (2012).
- 228 4. P Tiersma, Some myths about legal language. *J. Law, Cult. Humanit. Forthcoming, Loyola-LA Leg. Stud. Pap.* (2005).
- 229 5. P Tiersma, The nature of legal language in *AILA applied linguistics series: Vol. 5. Dimensions of forensic linguistics.*
230 (John Benjamins Publishing Company), pp. 7–25 (2008).
- 231 6. CA Hill, A comment on language and norms in complex business contracting. *Chi.-Kent L. Rev.* **77**, 29 (2001).
- 232 7. D Mellinkoff, *The language of the law.* (Wipf and Stock Publishers), (2004).