
Information-theoretic analysis of disfluencies in speech

Shiva Upadhye and Richard Futrell
{upadhye, rfutrell}@uci.edu
Department of Language Science
University of California, Irvine

Abstract

1 This study proposes and examines an information-theoretic measure of planning in
2 incremental speech production, and investigates the effects of planning, predictabil-
3 ity, and interference-based measures on lexical substitution errors. We then present
4 a rate-distortion theoretic model of speech production that explicates how these
5 factors affect the production of lexical substitution errors.

6 1 Introduction

7 Spontaneous speech is punctuated with disfluencies that reflect delays or interruptions in language
8 production mechanisms [1, 2, 3, 4]. One common kind of disfluency is a lexical substitution, where a
9 speaker says one word (which we term a **distractor**) and then corrects it to another word (which we
10 term a **target**) [5]. For example, see the following naturally-occurring utterance from the Switchboard
11 corpus [6]:

12 (1) and I believe that mental acuity is easy to *sustain* **maintain** if you just simply continue to
13 exercise your mind

14 Here the distractor is *sustain*, which the speaker corrects to *maintain*. More generally, several
15 prominent models of speech production have attributed hesitations and disfluencies of various kinds
16 to difficulties in selecting what word to say next [7, 8, 9].

17 Here we develop an information-theoretic analysis of lexical substitution errors in a dataset of
18 naturally-occurring disfluencies in the Switchboard corpus, applying a new measure of planning. We
19 leverage large language models to estimate the relevant information quantities, and perform a targeted
20 analysis that predicts *which word* is likely to appear as a distractor in each context. Finally, we
21 explicate our results within the framework of a rate–distortion theoretic model of speech production.

22 1.1 Background and related work

23 Computational characterizations of lexical selection have emphasized the effects of word frequency
24 and predictability on the ease of retrieving the appropriate word from memory. Low frequency words
25 were not only associated with delayed production in naming experiments [10, 11, 12], but were also
26 more likely to be preceded by disfluencies in naturalistic speech [13, 14]. Similarly, words that
27 were less predictable conditional on the past (low *forward* predictability) or within the surrounding
28 context (low *contextual* predictability) were also more likely to be preceded by hesitations, repetitions,
29 and corrections [1, 15, 13, 16, 17]. Intriguingly, disfluencies and slowdowns in speech were also
30 associated with *backward* transitional probability: the probability of a word given its *following*
31 context [2, 18]; this is usually considered to be related to a speaker’s plans for future production. In
32 addition to being a strong predictor of fillers and corrections [19, 20], backward predictability has

33 also been associated with delays caused by the reactivation of past material to aid retrieval of a target
34 compatible with the planned future [21].

35 A complementary line of experimental research in word production has also demonstrated that
36 delays or errors in selecting a target can arise due to the activation accrued by its semantic and/or
37 phonological neighbors [22, 23, 24]. The semantic interference effect, in particular, has been a
38 robust finding in the Picture Word Interference paradigm, where the degree of semantic relatedness
39 between a pictorial target and an orthographic distractor predicts a delay in retrieving the target [25].
40 Similarly, the presentation of a phonologically-related distractor has been associated with an increase
41 in activation of the segments that the distractor shares in common with the target [26]. Hence, higher
42 proximity between the target and distractor along semantic and/or phonological dimensions correlates
43 to a greater degree of co-activation during processing.

44 1.2 Planning measure

45 Integrating perspectives from both lines of inquiry, our analyses of lexical substitution errors in-
46 corporate the effects of frequency, incremental production, and planning along with distance-based
47 metrics such as phonological and semantic distance. We propose an information-theoretic measure of
48 planning based on Pointwise Mutual Information (PMI) [27]. In particular, we use the PMI of a word
49 x with the future context c_f given the past context c_p : $\text{pmiFP} = \ln \frac{p(x|c_f, c_p)}{p(x|c_p)}$, where $p(x | c_f, c_p)$
50 is the probability of x given both the preceding and following contexts and $p(x|c_p)$ is its forward
51 predictability. Unlike backward predictability, pmiFP does not assume that the future is planned
52 independent of the past. Rather it serves as a measure of planning that (i) does not commit to a
53 particular planning order and (ii) uses the information provided by the past context to estimate the
54 association between the word and the future context. Therefore, pmiFP can accommodate both linear
55 [28, 29] and hierarchical planning [30, 31]. Previous works investigating backward predictability
56 have not considered the integration of past and future context represented in pmiFP . In Section 4, we
57 give a more detailed theoretical justification for the use of this measure.

58 2 Methods

59 We develop a regression model that predicts whether a given word x , as opposed to any other word
60 in the vocabulary, is the distractor that the speaker selects and eventually overrides after production.
61 Specifically, we examine how proximity to the target, frequency, incremental or forward predictability,
62 and planned production guide the selection of a particular distractor given an utterance context.

63 *Measures*

64 **Frequency:** We use the frequencies from the SUBTLEXus corpus [32] to estimate unigram probabil-
65 ity $p(x)$ of a given word x

66 **Predictability-based measures:** We estimate the forward, backward, and contextual probabilities
67 of targets and distractors using XLNet [33], a large transformer-based model trained on both causal
68 and masked modeling objectives. We select utterances where the length of both the preceding and
69 following contexts exceed one word. To estimate $p(x | c_p)$ and $p(x | c_f)$ for a continuation x , we
70 provide as input the past context c_p and the future context c_f (in reverse) respectively. To estimate
71 $p(x | c_p, c_f)$, we provide the entire bidirectional utterance context (c_p, c_f) to the model with a *mask*
72 applied to x in order to indicate the word to be predicted (see Appendix A for examples).

73 **Distance measures:** We estimate phonological distance between the target and distractor based on
74 the Soundex algorithm [34] and semantic distance using cosine distance between pretrained GloVe
75 embeddings [35].

76 *Materials*

77 For our analyses, we use Switchboard Annotations [36], a human-parsed subset of the Switchboard
78 corpus of conversational speech [37] with annotations marking reparanda and repairs. In order to
79 extract utterances with lexical substitutions, we identify two distinct signatures characteristic of these
80 errors, namely (i) where the distractor or *reparandum* is immediately followed by the target or *repair*
81 (1) and (ii) where a single word is substituted within a repeated phrase (Appendix A). We restrict
82 our analyses to utterances with an equal number of reparanda and repairs in order to exclude cases
83 of additions, deletions, or structural revisions. Utterances with multiple substitution errors were
84 preprocessed into context frames with single substitution errors (see Appendix A for examples).

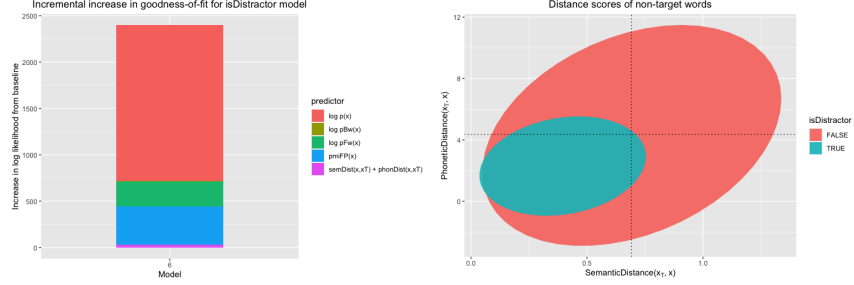


Figure 1: **Left:** Increase in log likelihood of isDistractor model with incrementally added predictors, in the order determined by goodness of fit in single-predictor models: $\log p$, distances, $\log pFw$, pmiFP, and $\log pBw$. The contribution of log backward probability is not visible. **Right:** Semantic and phonological distance to target for actual distractor words (blue) vs. all words in the vocabulary (red). Actual distractors are both semantically and phonologically close to the targets.

85 Model

86 We use a generalized linear mixed model (GLMM) [38] with a binary dependent variable *isDistractor*
 87 and the target identity as a random effect:

$$\text{isDistractor}(x) \sim \ln p(x) + \ln p(x | c_p) + \text{pmiFP}(x) + \text{phonDist}(x, x_T) + \text{semDist}(x, x_T) + (1 | x_T)$$

88 For each instance of a lexical substitution error (N = 1368), the model above is fit to predict which
 89 word out of a subset of nontarget words (N = 150) selected from the vocabulary with varying degrees
 90 of semantic and phonetic relatedness to the target is the chosen distractor.

91 We hypothesize that a nontarget x that is more frequent, and has higher forward predictability and
 92 *pmiFP* will have a higher likelihood of being the distractor that receives the activation required for
 93 selection. A word likely to be the distractor will also be co-activated with the target x_T by virtue
 94 of their common semantic and/or phonological features. Hence, we predict that the most probable
 95 candidate for the distractor will also have lower semantic and/or phonological distance from x_T .

96 3 Results

97 Both predictability and distance-based metrics emerged as significant predictors in the isDistractor
 98 model. As hypothesized, we observe a strong positive effect of frequency or unigram probability
 99 ($\beta = 0.240, p < 0.001$), forward predictability ($\beta = 0.269, p < 0.001$), and pmiFP ($\beta = 0.223, p <$
 100 0.001). In contrast, we find a strong negative effect of semantic ($\beta = -0.695, p < 0.001$) and
 101 phonological distance ($\beta = -0.103, p < 0.001$).

102 We analyze the effects of individual predictors by examining the goodness-of-fit of the isDistractor
 103 model with incrementally added features. We observe that a model with pmiFP as the planning
 104 measure provides a better fit than one with backward predictability ($\chi^2 = 800.5, p < 0.001$), the
 105 measure used in previous work. We also find that when both these predictors are included as planning
 106 measures, backward predictability has a negligible contribution to the increase in log likelihood
 107 (Figure 1).

108 4 Model sketch

109 To explicate our results, we present a sketch of a model of speech production within a rate–distortion
 110 framework, deriving an optimal probabilistic policy $p_g(x | c_p)$ for the selection of a word x given
 111 a communicative goal g and current state c_p consisting of a sequence of previous words, subject to
 112 a constraint on the policy’s usage of information about the goal g [39, 40, 41, 42, 43, 44, 45, 46].
 113 Following Todorov [39]’s KL-control framework, a policy is selected to maximize an average
 114 future-discounted **value-to-go**

$$v_g(x | c_p) = \ell_g(x | c_p) + \alpha \langle v_g(x | x, c_p) \rangle_{p_g^*(x|c_p)}, \quad (1)$$

115 where $\alpha \in [0, 1]$ is a future-discount parameter, $\langle \cdot \rangle_p$ indicates an average over distribution $p, p_g^*(x | c_p)$
 116 is an optimal policy, and $\ell(x | c_p)$ is the **local value** of an action x given goal g and state c_p . The

117 local value is given by the communicative value of a word x minus a **control cost** term which reflects
 118 the KL divergence between the policy p_g and an **automatic policy** p_0 which is not conditional on the
 119 goal g :

$$\ell_g(x | c_p) = \underbrace{R_g(x | c_p)}_{\text{Communicative value}} - \underbrace{\ln \frac{p_g(x | c_p)}{p_0(x | c_p)}}_{\text{Control cost}}. \quad (2)$$

120 Communicative value R_g is meant to signify how well the word x conveys the speaker’s intended
 121 message to a listener. Viewed within a rate-distortion theoretic framework, maximizing the com-
 122 municative value corresponds to minimizing the distortion subject to the rate or control cost, which
 123 quantifies the amount of information about the goal used to determine the optimal action. Given this
 124 setting, the policy that maximizes average value-to-go as derived by [39] is

$$p_g^*(x | c_p) \propto \exp\{\ln p_0(x | c_p) + R_g(x | c_p) + \alpha \langle v(x | c_p, x) \rangle\}. \quad (3)$$

125 We propose that the selection of the next word in speech is determined by Eq. 3. The policy predicts
 126 that what matters for the selection of a word x is (1) predictability given past context c_p , (2) the
 127 communicative value R_g of the word with respect to the current goal and state, and (3) the expected
 128 value of words following x , a kind of planning effect.

129 4.1 Application to lexical substitution errors

130 We consider lexical substitution errors to reflect cases where there are two words (the target and the
 131 distractor) that both receive high probability under Eq. 3. In that case, a word is likely to appear as a
 132 distractor whenever any of the three terms inside Eq. 3 are high. We will see that these correspond to
 133 forward predictability, semantic and phonetic distance, and pmiFP respectively.

134 To understand the effect of the communicative value of word x , we consider the difference in
 135 communicative value between the distractor x and target x_T , $\Delta R_g = R_g(x | c_p) - R_g(x_T | c_p)$.
 136 This value differential ΔR_g corresponds to a negative communicative cost for saying x instead of
 137 x_T . This cost should reflect both semantic distance, because semantically similar words will share
 138 many of their features relevant for communicative goals, and phonetic distance, because phonetically
 139 similar words may be indistinguishable to a listener.

140 In order to draw out predictions from Eq. 3, we make three simplifying assumptions. First, we
 141 assume production consists of a word x followed by a second word representing the entire future of
 142 the utterance, c_f . Second, we assume that the production of the future c_f is deterministic given the
 143 communicative goal, and third, that the communicative value of c_f is independent of the choice of x .
 144 Under these assumptions, the policy in Eq. 3 simplifies to

$$p_g^*(x | c_p) \propto \exp\{\ln p_0(x | c_p) + \Delta R_g + \alpha \ln p_0(c_f | c_p, x)\}. \quad (4)$$

145 Rewriting the probability of the future c_f using Bayes’ Rule in terms of the probability of the current
 146 word x given the future c_f , $p_0(x | c_p, c_f)$, we get

$$p_g^*(x | c_p) \propto \exp \left\{ \underbrace{\ln p_0(x | c_p)}_{\text{Forward predictability}} + \underbrace{\Delta R_g}_{\text{Distance}} + \alpha \underbrace{\ln \frac{p_0(x | c_p, c_f)}{p_0(x | c_p)}}_{\text{pmiFP}} \right\}, \quad (5)$$

147 where the three terms correspond to factors that were shown to predict which words appear as
 148 distractors in our corpus studies, including pmiFP. See Appendix B for full derivations. In addition,
 149 the frequency effects that we observe may be accommodated by adding an additional control cost for
 150 use of information about the state c_p , but we leave this modeling question to future work.

151 5 Conclusion

152 We have presented an analysis of lexical substitution errors in speech that predicts which words
 153 surface as distractors, and shown how the results can be accommodated in a rate–distortion control
 154 framework. The work opens the way for information-theoretic models of speech production that are
 155 tightly linked with rate–distortion models in other fields such as neuroscience and psychophysics.

References

- 156 [1] F. Goldman-Eisler, Speech production and language statistics, *Nature* 180 (1957) 1497–1497.
- 157 [2] F. Goldman-Eisler, *Psycholinguistics: Experiments in spontaneous speech* (1968).
- 158 [3] J. E. Fox Tree, The effects of false starts and repetitions on the processing of subsequent words
159 in spontaneous speech, *Journal of Memory and Language* 34 (1995) 709–738.
- 160 [4] E. Shriberg, Disfluencies in Switchboard, in: *Proceedings of International Conference on*
161 *Spoken Language Processing*, volume 96, Citeseer, 1996, pp. 11–14.
- 162 [5] W. J. Levelt, Monitoring and self-repair in speech, *Cognition* 14 (1983) 41–104.
- 163 [6] J. J. Godfrey, E. C. Holliman, J. McDaniel, SWITCHBOARD: Telephone speech corpus for
164 research and development, in: *IEEE International Conference on Acoustics, Speech, and Signal*
165 *Processing (ICASSP-92)*, volume 1, IEEE, 1992, pp. 517–520.
- 166 [7] V. A. Fromkin, Slips of the tongue, *Scientific American* 229 (1973) 110–117.
- 167 [8] G. Kempen, E. Hoenkamp, An incremental procedural grammar for sentence formulation,
168 *Cognitive Science* 11 (1987) 201–258. URL: [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S036402138780006X)
169 [article/pii/S036402138780006X](https://www.sciencedirect.com/science/article/pii/S036402138780006X). doi:[https://doi.org/10.1016/S0364-0213\(87\)](https://doi.org/10.1016/S0364-0213(87)80006-X)
170 [80006-X](https://doi.org/10.1016/S0364-0213(87)80006-X).
- 171 [9] K. Bock, W. J. M. Levelt, *Language production : Grammatical encoding*, 1994.
- 172 [10] R. C. Oldfield, A. Wingfield, Response latencies in naming objects, *Quarterly Journal of*
173 *Experimental Psychology* 17 (1965) 273–281.
- 174 [11] Z. M. Griffin, K. Bock, Constraint, word frequency, and the relationship between lex-
175 ical processing levels in spoken word production, *Journal of Memory and Language*
176 38 (1998) 313–338. URL: [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0749596X9792547X)
177 [S0749596X9792547X](https://www.sciencedirect.com/science/article/pii/S0749596X9792547X). doi:<https://doi.org/10.1006/jmla.1997.2547>.
- 178 [12] J. Almeida, M. Knobel, M. Finkbeiner, A. Caramazza, The locus of the frequency effect in
179 picture naming: When recognizing is not enough, *Psychonomic bulletin & review* 14 (2007)
180 1177–1182.
- 181 [13] G. Beattie, B. Butterworth, Contextual probability and word frequency as determinants of
182 pauses and errors in spontaneous speech, *Language and Speech* 22 (1979) 201 – 211.
- 183 [14] V. Kapatsinski, Frequency of use leads to automaticity of production: Evidence from repair in
184 conversation, *Language and Speech* 53 (2010) 105 – 71.
- 185 [15] P. H. Tannenbaum, F. A. Williams, C. S. Hillier, Word predictability in the environments of
186 hesitations, *Journal of Verbal Learning and Verbal Behavior* 4 (1965) 134–140.
- 187 [16] E. Shriberg, A. Stolcke, Word predictability after hesitations: a corpus-based study, *Proceeding*
188 *of Fourth International Conference on Spoken Language Processing. ICSLP '96* 3 (1996)
189 1868–1871 vol.3.
- 190 [17] Z. Harmon, V. Kapatsinski, Studying the dynamics of lexical access using disfluencies, *Papers*
191 *presented at DiSS* (2015) 41–44.
- 192 [18] A. Bell, J. M. Brenier, M. L. Gregory, C. Girand, D. Jurafsky, Predictability effects on durations
193 of content and function words in conversational english, *Journal of Memory and Language* 60
194 (2009) 92–111.
- 195 [19] S. Dammalapati, R. Rajkumar, S. Agarwal, Expectation and locality effects in the prediction of
196 disfluent fillers and repairs in English speech, in: *Proceedings of the 2019 Conference of the*
197 *North American Chapter of the Association for Computational Linguistics: Student Research*
198 *Workshop, Association for Computational Linguistics, Minneapolis, Minnesota, 2019*, pp.
199 103–109. URL: <https://aclanthology.org/N19-3015>. doi:10.18653/v1/N19-3015.
- 200

- 201 [20] S. Dammalapati, R. Rajkumar, S. Agarwal, Effects of duration, locality, and surprisal in
 202 speech disfluency prediction in English spontaneous speech, in: Proceedings of the Society for
 203 Computation in Linguistics 2021, Association for Computational Linguistics, Online, 2021, pp.
 204 91–101. URL: <https://aclanthology.org/2021.scil-1.9>.
- 205 [21] Z. Harmon, V. Kapatsinski, A theory of repetition and retrieval in language production.,
 206 Psychological review (2021).
- 207 [22] A. Roelofs, A spreading-activation theory of lemma retrieval in speaking, *Cognition* 42 (1992)
 208 107–142.
- 209 [23] G. M. Oppenheim, G. S. Dell, M. F. Schwartz, The dark side of incremental learning: A model
 210 of cumulative semantic interference during lexical access in speech production, *Cognition* 114
 211 (2010) 227–252.
- 212 [24] M. A. Goldrick, Limited interaction in speech production: Chronometric, speech error, and
 213 neuropsychological evidence, *Language and Cognitive Processes* 21 (2006) 817 – 855.
- 214 [25] R. R. Rosinski, Picture-word interference is semantically based., *Child Development* 48 (1977)
 215 643–647.
- 216 [26] A. S. Meyer, H. Schriefers, Phonological facilitation in picture-word interference experiments:
 217 Effects of stimulus onset asynchrony and types of interfering stimuli., *Journal of Experimental*
 218 *Psychology: Learning, Memory and Cognition* 17 (1991) 1146–1160.
- 219 [27] K. W. Church, P. Hanks, Word association norms, mutual information, and lexicography,
 220 in: 27th Annual Meeting of the Association for Computational Linguistics, Association for
 221 Computational Linguistics, Vancouver, British Columbia, Canada, 1989, pp. 76–83. URL:
 222 <https://aclanthology.org/P89-1010>. doi:10.3115/981623.981633.
- 223 [28] A. S. Meyer, Lexical access in phrase and sentence production: Results from picture–word
 224 interference experiments, *Journal of Memory and Language* 35 (1996) 477–496.
- 225 [29] S. Brown-Schmidt, M. K. Tanenhaus, Watching the eyes when talking about size: An investi-
 226 gation of message formulation and utterance planning, *Journal of Memory and Language* 54
 227 (2006) 592–609.
- 228 [30] K. Christianson, F. Ferreira, Conceptual accessibility and sentence production in a free word
 229 order language (odawa), *Cognition* 98 (2005) 105–135.
- 230 [31] S. Momma, V. S. Ferreira, Beyond linear order: The role of argument structure in speaking,
 231 *Cognitive Psychology* 114 (2019).
- 232 [32] M. Brysbaert, B. New, Moving beyond kučera and francis: A critical evaluation of current
 233 word frequency norms and the introduction of a new and improved word frequency measure for
 234 american english, *Behavior research methods* 41 (2009) 977–990.
- 235 [33] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, Q. V. Le, Xlnet: Generalized au-
 236 toregressive pretraining for language understanding, *Advances in neural information processing*
 237 systems 32 (2019).
- 238 [34] R. Russell, M. Odell, The soundex indexing system, National Archives and Records Adminis-
 239 tration (1918).
- 240 [35] J. Pennington, R. Socher, C. D. Manning, Glove: Global vectors for word representation,
 241 in: Proceedings of the 2014 conference on empirical methods in natural language processing
 242 (EMNLP), 2014, pp. 1532–1543.
- 243 [36] S. Calhoun, J. Carletta, J. M. Brenier, N. Mayo, D. Jurafsky, M. Steedman, D. I. Beaver,
 244 The nxt-format switchboard corpus: a rich resource for investigating the syntax, semantics,
 245 pragmatics and prosody of dialogue, *Language Resources and Evaluation* 44 (2010) 387–419.
- 246 [37] J. J. Godfrey, E. Holliman, J. McDaniel, Switchboard: telephone speech corpus for research and
 247 development, [Proceedings] ICASSP-92: 1992 IEEE International Conference on Acoustics,
 248 Speech, and Signal Processing 1 (1992) 517–520 vol.1.

- 249 [38] D. G. Clayton, Generalized linear mixed models, Markov chain Monte Carlo in practice 1
250 (1996) 275–302.
- 251 [39] E. Todorov, Efficient computation of optimal actions, Proceedings of the National Academy of
252 Sciences 106 (2009) 11478–11483.
- 253 [40] N. Tishby, D. Polani, Information theory of decisions and actions, in: Perception-action cycle,
254 Springer, 2011, pp. 601–636.
- 255 [41] S. G. Van Dijk, D. Polani, C. L. Nehaniv, Hierarchical behaviours: getting the most bang for
256 your bit, in: European Conference on Artificial Life, Springer, 2009, pp. 342–349.
- 257 [42] P. A. Ortega, D. A. Braun, Thermodynamics as a theory of decision-making with information-
258 processing costs, Proceedings of the Royal Society A: Mathematical, Physical, and Engineering
259 Sciences 469 (2013) 20120683.
- 260 [43] T. Genewein, F. Leibfried, J. Grau-Moya, D. A. Braun, Bounded rationality, abstraction,
261 and hierarchical decision-making: An information-theoretic optimality principle, Frontiers in
262 Robotics and AI 2 (2015) 27.
- 263 [44] S. J. Gershman, Origin of perseveration in the trade-off between reward and complexity, bioRxiv
264 (2020).
- 265 [45] N. Zaslavsky, J. Hu, R. P. Levy, A Rate–Distortion view of human pragmatic reasoning, arXiv
266 preprint arXiv:2005.06641 (2020).
- 267 [46] R. Futrell, An information-theoretic account of semantic interference in word production,
268 Frontiers in Psychology 12 (2021) 672408.

269 A Appendix: Lexical substitution examples and data processing

270 Example of a lexical substitution within repeated material:

271 1 so until i see the entire quote old guard of the soviet *military* of the soviet **government**
272 completely roll over and disappear preferably buried i still consider them a threat

273 Example of utterance with multiple substitution errors preprocessed into contexts with single substi-
274 tution errors:

275 A it depends on whether you whether we figure that we have a defense oriented military or an
276 *aggressive aggression* oriented military

277 a it depends on whether [you/we] figure that we have a defense oriented military or an
278 aggression oriented military

279 b it depends on whether we figure that we have a defense oriented military or an [*aggres-*
280 *sion/aggressive*] oriented military

281 Examples of preprocessed XLnet inputs for estimating forward, backward, and masked probabilities:

282 A.3 it depends on whether <mask>

283 A.4 military oriented aggression an or military oriented defense a have we that figure <mask>

284 A.5 it depends on whether <mask> figure that we have a defense oriented military or an aggres-
285 sion oriented military

286 Code for preprocessing, calculating metrics, and analysis can be found at: InfoTheoreticDisfModel

287 **B Appendix: Derivation of pmiFP from the speech production model**

288 Starting with policy of Eq. 3, we can get to Eq. 5 by assuming (1) production consists of a word
 289 x followed by a second word representing the entire future of the utterance, c_f , (2) production of
 290 the future c_f is deterministic given the communicative goal, and (3) that the communicative value
 291 of c_f is independent of the choice of x . Under these assumptions, the policy in Eq. 3. The first
 292 assumption means that we only need to consider a finite time horizon with one future action. The
 293 second assumption means that the policy $p_g^*(c | c_p, x) = \delta_{cc_f}$. The third assumption means that
 294 $R_g(c_f | x, c_p)$ is the same for all x . These assumptions represent a scenario where the speaker has
 295 high certainty about what they will say next, and where the next part of an utterance is relatively
 296 independent of the current part, for example at the end of a phrase or clause.

297 We start with the policy probability (setting $\alpha = 1$ to save writing):

$$p_g^*(x | c_p) \propto \exp\{\ln p_0(x | c_p) + R_g(x | c_p) + \langle v(x' | c_p, x) \rangle\}. \quad (6)$$

298 First, we rewrite $R_g(x | c_p) = R_g(x_T | c_p) + \Delta R_g$. Because $R_g(x_T | c_p)$ is not a function of the
 299 action under consideration x , it can be absorbed into the normalizing constant of Eq. 3, giving

$$p_g^*(x | c_p) \propto \exp\{\ln p_0(x | c_p) + \Delta R_g + \langle v(x' | c_p, x) \rangle\}. \quad (7)$$

300 Now using assumption (1), we can rewrite the policy in Eq. 3 as:

$$p_g^*(x | c_p) \propto \exp\left\{\ln p_0(x | c_p) + \Delta R_g - \left\langle R_g(c | c_p, x) + \ln \frac{p_g^*(c | c_p, x)}{p_0(c | c_p, x)} \right\rangle_{p_g^*(c | c_p, x)}\right\}. \quad (8)$$

301 Using $p_g^*(c | c_p, x) = \delta_{cc_f}$ for the expectation over future actions (assumption 2), we get

$$p_g^*(x | c_p) \propto \exp\left\{\ln p_0(x | c_p) + \Delta R_g + R_g(c_f | c_p, x) - \ln \frac{1}{p_0(c_f | c_p, x)}\right\}. \quad (9)$$

302 Because $R_g(c_f | c_p, x)$ is invariant to x (assumption 3), it can be absorbed into the implicit normaliz-
 303 ing constant of Eq. 9, giving

$$p_g^*(x | c_p) \propto \exp\{\ln p_0(x | c_p) + \Delta R_g + \ln p_0(c_f | c_p, x)\}. \quad (10)$$

304 Now applying Bayes' rule to the term $\ln p_0(c_f | c_p, x)$, and then applying logarithm rules, we get

$$p_g^*(x | c_p) \propto \exp\left\{\ln p_0(x | c_p) + \Delta R_g + \ln \frac{p_0(x | c_p, c_f)p_0(c_f | c_p)}{p_0(x | c_p)}\right\} \quad (11)$$

$$= \exp\left\{\ln p_0(x | c_p) + \Delta R_g + \ln \frac{p_0(x | c_p, c_f)}{p_0(x | c_p)} + \ln p_0(c_f | c_p)\right\}. \quad (12)$$

305 Again, the last term is invariant to x , so it can be absorbed into the normalizing constant, leaving us
 306 with the policy considered in the main text

$$p_g^*(x | c_p) \propto \exp\left\{\ln p_0(x | c_p) + \Delta R_g + \ln \frac{p_0(x | c_p, c_f)}{p_0(x | c_p)}\right\}. \quad (13)$$