

Repetition Facilitates Processing: The Processing Advantage of Construction Repetition in Dialogue

Anonymous ACL submission

Abstract

Repetitions occur frequently in dialogue. This study focuses on the repetition of lexicalised constructions—i.e., recurring multi-word units—in English open domain spoken dialogues. We hypothesise that construction repetition is an efficient communication strategy that reduces processing effort, and we make three predictions based on this hypothesis. We conduct a quantitative analysis, measuring reduction in processing effort via two surprisal-based measures and estimating surprisal with an adaptive neural language model. Our three predictions are confirmed: (i) repetitions facilitate the processing of constructions and of their linguistic context; (ii) facilitating effects are higher when repetitions accumulate, (iii) and lower when repetitions are less locally distributed. Our findings suggest that human-like patterns of repetitions can be learned implicitly by utterance generation models equipped with psycholinguistically motivated learning objectives and adaptation mechanisms.

1 Introduction

In language production, speakers select—among a set of possible realisations—the lexical, syntactic, and semantic alternatives they deem most appropriate to verbalise their communicative intents. For instance, speakers can choose to precede reported speech with ‘*I said*’ or ‘*I was like*’: ‘*I was like where is this going?*’, ‘*I said you don’t have to love each other*’. Given such sets of alternatives, speakers’ choices are influenced, among other things, by their recent linguistic experience. In a dialogue, a speaker may be more prone to choose ‘*I was like*’ if they or their conversational partner have already used it. This is an example of priming: under the influence of previous mentions, ‘*I was like*’ is repeated more often than expected by chance.

Most studies on priming have targeted the repetition of syntactic structures (Levelt and Kelter, 1982; Bock, 1986; Branigan et al., 2000; Reit-

ter et al., 2006b, 2011), often explaining them within the framework of the interactive alignment model (Pickering and Garrod, 2004). Lexical repetitions have also been investigated (e.g., Brennan, 1996) and they have been typically explained as the result of collaborative mechanisms (Brennan and Clark, 1996) or social pressures (Danescu-Niculescu-Mizil et al., 2012; Noble and Fernández, 2015; Doyle and Frank, 2016). Less is known about the mechanisms underlying speakers’ repetition of particular configurations of structures and lexemes, *constructions*, a pervasive phenomenon in conversational language use (Tomasello, 2003; Goldberg, 2006; Sinclair and Fernández, 2021). The reuse of constructions has been analysed by Fusaroli et al. (2014) as part of a process of ‘interpersonal synergy’ between conversational partners. In this study, we investigate whether speakers repeat lexicalised constructions (such as ‘*I was like*’) throughout a dialogue as a result of two information processing mechanisms traditionally argued to affect priming: 1) *residual activations* due to exposure to local context (Pickering and Branigan, 1998; Cleland and Pickering, 2003) and 2) *implicit learning* of the global statistics of expressions and structures (Bock and Griffin, 2000; Fine and Florian Jaeger, 2013). We use a computational model to approximate these mechanisms, hypothesising that, if they are in place, construction repetition becomes a rational strategy of information transmission (Gibson, 1998; Levy, 2008): processing effort is reduced when speakers follow this strategy.

We use *surprisal* to operationalise the processing advantage of construction repetition, estimated with a neural language model. Surprisal measures the unpredictability of a linguistic signal, which can be taken as an estimate of the amount of effort required to process the signal (e.g., Jelinek et al., 1975; Keller, 2004; Levy, 2008). We predict (i) that construction repetition has a facilitating effect on processing, observable in the form of a surprisal

reduction both for the construction itself and for its linguistic context. To further understand the nature of the processing advantage, we study how it varies across different types of repetition. We predict (ii) that the processing advantage of construction repetition increases with the total number of repetitions made in a dialogue, and (iii) that it decreases with the distance between repetitions. Our experiments confirm these three predictions, providing new empirical evidence that dialogue partners use repetitions as a communication strategy due to it leading to higher information processing efficiency.

Our findings inform the development of better dialogue models. They indicate that avoiding repetitions in utterance generation (Li et al., 2016; Welleck et al., 2019) may not be the most appropriate strategy. Instead, models should be encouraged to follow human-like patterns of repetitions to be successfully deployed in conversational settings.

2 Background

2.1 Constructions

This work focuses on *constructions*, seen as particular configurations of structures and lexemes in usage-based accounts of natural language (Tomasello, 2003; Bybee, 2006, 2010; Goldberg, 2006). According to these accounts, models of language processing must consider not only individual lexical elements according to their syntactic roles, but also more complex form-function units, which can break regular phrasal structures—e.g., ‘*I know I*’, ‘*something out of*’. We further focus on fully lexicalised constructions (sometimes called *formulaic expressions*, or *multi-word expressions*). Commonly studied types of constructions are idioms (‘*break the ice*’), collocations (‘*pay attention to*’), phrasal verbs (‘*make up*’), and lexical bundles (‘*a lot of the*’). In Section 5, we explain how the notion of lexicalised construction is operationalised in the current study; Table 1 shows some examples.

A common property of constructions is their frequent occurrence in natural language. As such, they possess what in usage-based accounts is sometimes referred to as ‘processing advantage’ (Conklin and Schmitt, 2012; Carrol and Conklin, 2020). Evidence for the processing advantage of construction *usage* has been found in reading (Arnon and Snider, 2010; Tremblay et al., 2011), naming latency (Bannard and Matthews, 2008; Janssen and Barber, 2012), eye-tracking (Underwood, 2004; Siyanova-Chanturia et al., 2011), and electrophys-

SYXU	S7ZG	SVPK
<i>had a few</i>	<i>if you look at</i>	<i>I think it was just</i>
<i>it I was</i>	<i>yes of course</i>	<i>like this is</i>
<i>I’d be like</i>	<i>look at what</i>	<i>like you’re not</i>
<i>were like oh</i>	<i>if you give</i>	<i>so I didn’t</i>
<i>do you get</i>	<i>and all of that</i>	<i>that I know</i>
<i>and I went</i>	<i>it doesn’t have to</i>	<i>it’s not even</i>
<i>I don’t like</i>	<i>right okay so</i>	<i>and I was kind of</i>
<i>a bit more</i>	<i>something out of</i>	<i>and it was like oh</i>
<i>I know I</i>	<i>that in itself</i>	<i>think of it like</i>
<i>I was like</i>	<i>yeah that’s fine</i>	<i>kind of thing where</i>

Table 1: Top 10 constructions from three dialogues of the Spoken BNC (Love et al., 2017). Constructions are sorted according to the PMI between a construction and its dialogue (see Section 5 for extraction procedure). Headers correspond to the dialogues’ IDs in the corpus.

iology (Tremblay and Baayen, 2010; Siyanova-Chanturia et al., 2017). In this paper, we study the processing advantage of the *repetition* of lexicalised constructions.

2.2 Surprisal and Processing Effort

Estimates of surprisal have been shown to be good predictors of processing effort in perception (Jelinek et al., 1975; Clayards et al., 2008), reading (Keller, 2004; Demberg and Keller, 2008; Levy et al., 2009), and sentence interpretation (Levy, 2008; Gibson et al., 2013). Because speakers take into consideration their addressee’s processing effort (Clark and Wilkes-Gibbs, 1986; Clark and Schaefer, 1989), their linguistic choices can often be explained as an optimal strategy to manage the fluctuations of surprisal levels over time. Surprisal-based accounts have indeed been successful at explaining various aspects of language production: speakers tend to reduce the duration of less surprising sounds (Aylett and Turk, 2004, 2006; Bell et al., 2003; Demberg et al., 2012); they are more likely to drop sentential material within less surprising scenarios (Jaeger and Levy, 2007; Frank and Jaeger, 2008; Jaeger, 2010); they tend to overlap at low-surprisal dialogue turn transitions (Dethlefs et al., 2016); and they produce sentences at a uniform surprisal rate in texts (Genzel and Charniak, 2002, 2003; Qian and Jaeger, 2011).

To estimate surprisal, we use GPT-2 (Radford et al., 2019), a neural language model. Using language models to approximate surprisal is an established approach (e.g., Genzel and Charniak, 2002; Keller, 2004; Xu and Reitter, 2018) and *neural models’* surprisal estimates in particular have been shown to be good predictors of processing effort,

168 measured as reading time, gaze duration, and N400
169 response (van Schijndel and Linzen, 2018; Merx
170 and Frank, 2021).

171 2.3 Priming Mechanisms

172 Priming has been widely studied through the anal-
173 ysis of structural repetitions, whether densely clus-
174 tered (e.g., Branigan et al., 1999; Wheeldon and
175 Smith, 2003), or occurring across multiple utter-
176 ances and interactions (e.g., Branigan et al., 2000;
177 Kaschak et al., 2014). These two types of prim-
178 ing (often called *short-term priming* and *long-term*
179 *priming*, respectively) are thought to be the result of
180 different underlying mechanisms (for a review see,
181 e.g., Hartsuiker et al., 2008). Quickly decaying,
182 short-term priming effects rely on an activation-
183 based mechanism dependent on residual traces
184 left by lexical material (Pickering and Branigan,
185 1998; Cleland and Pickering, 2003). Slowly decay-
186 ing, long-term priming effects are independent of
187 lexical material and rely on an implicit learning
188 mechanism (Bock and Griffin, 2000; Fine and Flo-
189 rian Jaeger, 2013). In the current study, we model
190 both mechanisms so that we do not limit a priori
191 the space of possible processes underlying priming.

192 3 Hypotheses

193 Does construction repetition come with a process-
194 ing advantage? Is this advantage due to the mecha-
195 nisms underlying priming? To answer these ques-
196 tions, we formulate the following three hypotheses.

197 **H1** *Repetition facilitates processing.* We predict
198 1) a construction has lower surprisal when
199 repeated than when first produced, and 2) rep-
200 etitions of a construction (i.e., the occurrences
201 that follow its first mention) have a stronger
202 reduction effect on the surprisal of the dia-
203 logue turn (i.e., a stronger *facilitating effect*)
204 than first mentions.

205 **H2** *The processing advantage of repetition is cu-*
206 *mulative.* We predict multiple repetitions of a
207 construction contribute 1) to a stronger reduc-
208 tion in the surprisal of the construction itself,
209 and 2) to a stronger facilitating effect.

210 **H3** *The processing advantage of repetition decays*
211 *as a function of the distance between repeti-*
212 *tions.* We predict that a larger distance be-
213 tween a construction repetition and its previ-
214 ous mention results 1) in a weaker reduction
215 in the surprisal of the construction, and 2) in
216 a weaker facilitating effect.

217 **H1** tests whether repeating a construction re-
218 duces processing effort. Comprehenders are known
219 to process written and spoken words more rapidly
220 when they are repeated (for a review, see Bigand
221 et al., 2005), suggesting increased expectation for
222 these words. An increase in expectation (hence
223 reduction in surprisal) due to repetition is compat-
224 ible with the implicit learning account of priming
225 (Kaschak et al., 2006; Reitter et al., 2011; Fine
226 et al., 2013). However, if repetitions are closely
227 clustered, any surprisal reduction could also be the
228 result of residual activations from previous men-
229 tions, in line with the activation-based account.

230 Because **H1** does not distinguish between differ-
231 ent repetitions of a construction and their distribu-
232 tion across time, **H2** tests how surprisal reduction
233 effects vary along chains of repetitions in terms
234 of cumulation (Table 4 shows an example chain).
235 Changes in the magnitude of the processing advan-
236 tage of construction repetition may interact with
237 the number of times the construction has already
238 been repeated (Jaeger and Snider, 2008; Fine and
239 Jaeger, 2016). Cumulative effects propagating over
240 distant repetitions would be evidence in favour of
241 the implicit learning account, whereas cumulative
242 effects taking place locally are compatible with the
243 activation-based account.

244 The processing advantage of construction rep-
245 etition may also be determined by the distance
246 between mentions. Inspired by earlier analyses
247 conducted for lexical and syntactic priming with
248 varying results (Reitter et al., 2011; Howes et al.,
249 2010; Healey et al., 2014), **H3** investigates the in-
250 fluence of recency of previous mention on a rep-
251 etition’s processing advantage. Fast decay effects
252 could be taken in support of the activation-based
253 account, whereas slow decay effects would suggest
254 reduction in surprisal is due to sensitivity to the
255 global statistics of expressions and structures in a
256 dialogue, in line with the implicit learning account.

257 4 Data

258 We test our hypotheses on the Spoken British Na-
259 tional Corpus¹ (Love et al., 2017), a dataset of tran-
260 scribed spoken open domain dialogues containing
261 1,251 contemporary British English conversations,
262 collected in a range of real-life contexts. We focus
263 on the 622 dialogues that feature only two speakers,
264 and randomly split them into a 70% finetuning set
265 (to be used as described in Section 6) and a 30%

¹<http://www.natcorp.ox.ac.uk>.

analysis set. Table 2 shows basic statistics for the dialogues used in this study.

	Mean \pm Std	Median	Min	Max
Dialogue length (# turns)	736 \pm 599	541.5	67	4859
Dialogue length (# words)	7753 \pm 5596	6102	819	39575
Turn length (# words)	11 \pm 15	6	1	982

Table 2: Two-speaker dialogue statistics, Spoken BNC.

5 Extracting Repeated Constructions

We define constructions as multi-word sequences that are repeated within a dialogue. We analyse constructions produced by only one of the dialogue participants as well as those produced by both speakers. To extract a set of constructions from each dialogue, we use the sequential pattern mining method proposed by Duplessis et al. (2017a,b, 2021), which treats the extraction task as an instance of the longest common subsequence problem (Hirschberg, 1977; Bergroth et al., 2000).² We modify it to not discard multiple repetitions of a construction that occur in the same dialogue turn. We focus on constructions of at least three tokens, uttered at least three times in a dialogue. Repeated sequences that mostly appear as a sub-part of a larger repeated construction are discarded.³

We apply the following further constraints. First, we exclude topic-determined constructions and referential expressions in order to disentangle priming effects from topic coherence effects. To this end, we filter out constructions that include nouns, unless the nouns are highly generic.⁴ For example, we discard sequences such as ‘playing table tennis’ and ‘a woolly jumper’ and retain constructions such as ‘a lot of’ and ‘the thing is’. Second, we filter out repetitions that are simply due to a high base frequency rate and not to the speakers’ self and mutual priming effects. We measure the association strength between a construction c and a dialogue d as the pointwise mutual information (PMI) between the two:

$$PMI(c, d) = \log_2 \frac{P(c|d)}{P(c)} \quad [1]$$

²Their code is freely available at <https://github.com/GuillaumeDD/dialign>.

³We discard constructions that appear less than twice outside of a larger repeated construction in a given dialogue (e.g., ‘think of it’ vs. ‘think of it like’).

⁴We define a limited specific vocabulary of generic nouns (e.g., ‘thing’, ‘fact’, ‘time’); full vocabulary in Appendix B.

which measures how unusually frequent a construction is in a given dialogue, compared to the rest of the corpus. We discard all constructions that have a PMI score lower than 1 in their respective dialogue. The probabilities in Eq. 1 are obtained using maximum likelihood estimation over the analysis split of the Spoken BNC. Finally, we exclude sequences containing punctuation marks or which consist of more than 50% filled pauses (e.g., ‘mm’, ‘erm’).⁵

Applying the described extraction procedure to the 187 dialogues in the analysis split of the Spoken BNC, we obtain a total of 3,676 unique constructions and 33,103 occurrences. Further statistics on the extracted constructions are presented in Table 3. Table 1 shows examples of the top 10 constructions extracted from three dialogues, ranked according to their PMI score.

	Mean \pm Std	Median	Min	Max
Construction length	3.23 \pm 0.52	3	3	7
Construction frequency	3.87 \pm 1.93	3	3	58
Constructions per dialogue	206 \pm 307	100	3	2023
Words per dialogue turn	31 \pm 37	21	3	959

Table 3: Construction statistics for the analysis split of the Spoken BNC. *Construction frequency* is the number of occurrences of a given construction in a dialogue, *Constructions per dialogue* is the number of occurrences of all constructions in a dialogue, *Words per dialogue turn* is computed on turns containing a construction.

6 Experimental Setup

In this section, we present two surprisal-based measures of processing advantage, the adaptive language model that produces surprisal estimates, and statistical tests used to confirm our hypotheses.⁶

6.1 Measures of Processing Advantage

The *surprisal* of a word choice w_i is the negative logarithm of the corresponding word probability, conditioned on the dialogue turn context t (i.e., the words that precede w_i in the dialogue turn) and on the local dialogue context l :

$$H(w_i|t, l) = -\log_2 P(w_i|t, l) \quad [2]$$

We define the local dialogue context l as the 50 tokens that precede the first word in the dialogue turn.⁷ We use tokens as a unit of context size, rather

⁵The full list of filled pauses can be found in Appendix B.

⁶Data and code will be made public upon acceptance.

⁷Building on prior work (Reitter et al., 2006a) that uses a window of 15 seconds of spoken dialogue as the locus of

Speaker	RI	RI Turn	Dist	Turn	<i>S</i>	<i>FE</i>
A	0	0	-	Drink? that was what he did yeah just just to just to know that I he might not be a complete twat but just a fyi	4.73	0.40
B	1	0	1586	Especially for my birthday mind you I might not be here for mine and I went what do you mean you might not be here?	4.01	0.53
	2	1	14		2.70	0.90

Table 4: Repetition chain for the construction ‘*might not be*’ in dialogue SXWH, Spoken BNC, annotated with repetition index (RI), RI within dialogue turn (RI Turn), and distance from previous mention (Dist; in tokens).

than dialogue turns, since they more closely correspond to the temporal units used in previous work (e.g., Reitter et al., 2006a), and since the length of dialogue turns can vary significantly (see Table 2). To measure the surprisal of a construction c , we average over word-level surprisal values:

$$S(c; t, l) = \frac{1}{|c|} \sum_{w_i \in c} H(w_i | t, l) \quad [3]$$

Surprisal estimates provide a computational approximation of the effort required to process a construction in context. We also measure the surprisal change (increase or reduction in processing effort) contributed by a construction c to its dialogue turn context, which we call the *facilitating effect* of a construction. The facilitating effect is positive when the construction has lower surprisal than its context, and negative when it has higher surprisal:

$$FE(c; t, l) = \log_2 \frac{\frac{1}{|s|-|c|} \sum_{w_j \in s, w_j \notin c} H(w_j | t, l)}{\frac{1}{|c|} \sum_{w_i \in c} H(w_i | t, l)} \quad [4]$$

Due to human memory constraints, the facilitating effect of constructions is more likely to affect the processing of words that are produced immediately before and after the construction itself. We define the locus of the facilitating effect (s in Eq. 4) as the 10 tokens preceding and the 10 tokens following the construction.⁸ The tokens exceeding the limits of the current dialogue turn are discarded. When the locus s corresponds to the construction itself, the facilitating effect equals 0.

6.2 Estimates of Surprisal

To produce surprisal estimates, we use a computational model of next word prediction which implements approximations of both the activation-based

local priming effects, we compute the average speech rate in the Spoken BNC (3.16 tokens/second) and multiply it by 15; we then round up the result (47.4) to 50 tokens.

⁸This is motivated by the fact that the average length of turns containing a construction is 31 tokens (median length is 21), with constructions being 3 to 7 tokens long—see Table 3.

and the implicit learning mechanism: it is conditioned on local contextual cues while it learns from exposure to the global dialogue context. We use GPT-2 (Radford et al., 2019), a pre-trained autoregressive Transformer language model. We take GPT-2’s attention mechanism (Vaswani et al., 2017) over the preceding context of a word as a proxy for the local activation-based mechanism: words in the more proximate dialogue context shape the model’s expectations for next words, and thus their contextualised surprisal. As an implicit learning mechanism, we use the Transformer’s standard learning rule, back-propagation on the cross-entropy next word prediction error, which has been successful at modelling a wide range of linguistic phenomena (Rumelhart and McClelland, 1986; Elman, 1991; Cleeremans and Elman, 1993; van Schijndel and Linzen, 2018). We rely on HuggingFace’s implementation of GPT-2 with default tokenizers and parameters (Wolf et al., 2020), and finetune the pre-trained model on a 70% training split of the Spoken BNC in order to adapt it to the idiosyncrasies of spoken dialogic data.⁹ We refer to this finetuned version as the *frozen* model. We use an attention window of length 50, i.e., the size of the local dialogue context, which may span over multiple dialogue turns (see Section 6.1).

Adaptive language model When estimating surprisal for a dialogue, we begin by processing the first turn using the frozen language model and then gradually update the model parameters after each turn, using back-propagation with cross-entropy loss. The magnitude of the learning rate is important for these updates to have the desired effect. The learning rate should be sufficiently high for the language model to adapt during a single dialogue, yet an excessively high learning rate can cause the language model to lose its ability to generalise across dialogues. To find the appropriate learning rate, we randomly select 18 dialogues from

⁹More details on finetuning can be found in Appendix C.1.

the analysis split of the Spoken BNC¹⁰ and run an 18-fold cross-validation for a set of six candidate learning rates: $1e - 5$, $1e - 4$, \dots , 1 . We fine-tune the model on each dialogue using one of these learning rates, and compute perplexity reduction 1) on the dialogue itself (*adaptation*) as well as 2) on the remaining 17 dialogues (*generalisation*). We select the learning rate yielding the best adaptation over cross-validation folds ($1e - 3$), while still improving the model’s generalisation ability. See Appendix C.2 for further details.

6.3 Statistical Modelling

To test **H1**, we split all occurrences of constructions by whether they are the first mention in a dialogue or a repetition. Our dataset consists of 8,562 first mentions and 24,541 repetitions. Using a Two Sample Bayesian t-test,¹¹ we compare the S distribution of first mentions to that of repetitions. We perform the same analysis for FE values.

H2 and **H3** focus on analysing repetitions only. We label each occurrence with a *repetition index* (the first repetition of a construction has an index of 1, the second, 2, etc.), and with the *distance from the previous mention* in a dialogue, measured as the number of words between the first word of the current occurrence and the first word of the previous occurrence (see Table 4). We fit two linear mixed effect models using S and FE as response variables, and include multi-level random effects grouped by dialogue and individual speaker ID. To select the models’ fixed effects, we start with a collection of motivated features—including repetition index and distance from previous mention—and perform an ablation selection procedure, iteratively removing features with the lowest significance, keeping only those that yield a p -value lower than 0.05.¹²

7 Results

We now present the results of our experiments, testing three hypotheses on the processing advantage (surprisal reduction and facilitating effect) of construction repetition. The final linear mixed effect models for both construction surprisal S and facilitating effect FE include repetition index and dis-

tance from the previous mention, which are directly related to our hypotheses, as well as construction length and repetition index within the current turn. The full specification of the best models, with fixed and random effect coefficients, is in Appendix D.

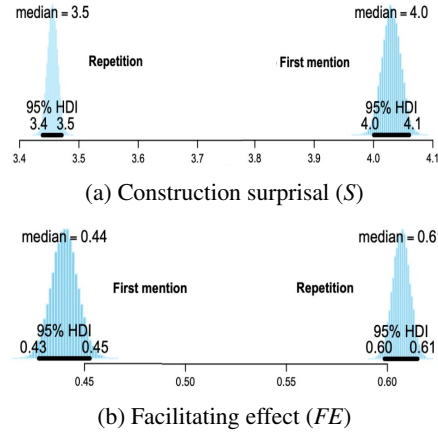


Figure 1: Posterior predictive distributions for the mean S and FE according to the Bayesian t-test between first mentions and repetitions.

Repetition facilitates processing (H1) Figures 1a and 1b show that the posterior distributions of the mean S and FE do not overlap between groups. For both metrics, highest density intervals of difference between means do not include 0. In sum, we find surprisal of construction repetitions is lower than that of first mentions, and repetitions have a stronger facilitating effect than first mentions. Our first two predictions are thus confirmed.

The processing advantage of repetition is cumulative (H2) The effect of repetition index is negative on S ($-24.85e - 2$, $p < 2e - 16$) and positive on FE ($7.57e - 2$, $p < 2e - 16$). Figures 2a and 2b show the opposite trajectories of the measures, with a stronger effect of repetition index on construction surprisal. In sum, we find that the surprisal of construction decreases, and their facilitating effect increases, as previous mentions accumulate. This confirms our second pair of predictions.

The processing advantage of repetition decays (H3) The distance of a construction from its previous mention has a positive effect on S ($9.66e - 2$, $p < 2e - 16$) and a negative effect on FE ($-4.29e - 2$, $p < 2e - 16$), also shown in Figures 2c and 2d. Surprisal increases, and facilitating effect decreases, as the current usage of a construction gets further away from its previous mention. Our third pair of predictions is thus confirmed.

¹⁰This amounts to ca. 10% of the analysis split. We use the analysis split because there is no risk of “overfitting” with respect to our main analyses.

¹¹We use the t-test implemented in the ‘Bayesian First Aid’ R-JAGS package (https://github.com/rasmusab/bayesian_first_aid) with the default uninformative priors and a credible interval of 95%.

¹²The full list of features can be found in Appendix D.

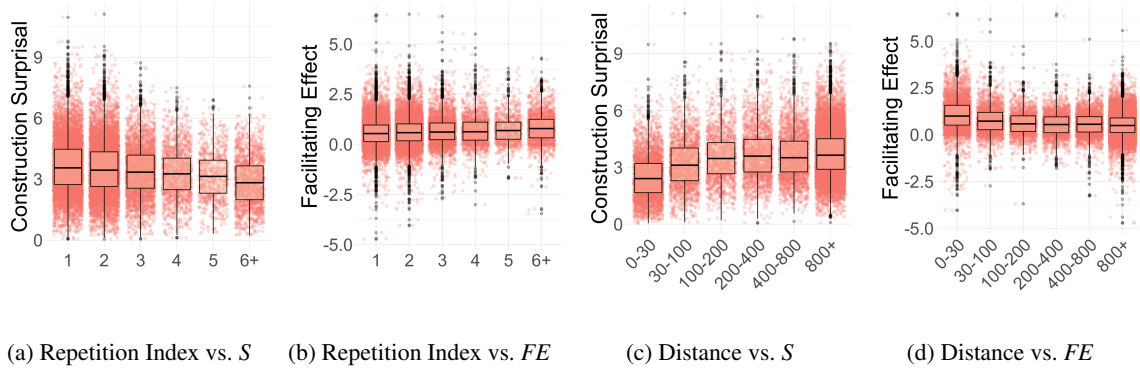


Figure 2: Construction surprisal (S , bits) and facilitating effect (FE) vs. repetition index and distance from previous mention (number of words). The first distance bin is the mean length of a turn containing a construction (Table 3).

8 Analysis

Having confirmed our hypotheses, we now further analyse the distribution of FE and S estimates, their relationship, and how their values across repetitions are influenced by additional factors.

8.1 Measures of Processing Advantage

Our first observation is that not only construction *repetition* but also construction *usage* comes with a processing advantage, as measured with both S and FE —a finding in line with prior work (e.g., Arnon and Snider, 2010; Bannard and Matthews, 2008; Tremblay et al., 2011; Janssen and Barber, 2012). On the one hand, as shown in Figure 1b, the posterior distribution of the mean FE spans over positive values for both first mentions and repetitions. The estimated mean FE of constructions is higher than the mean (0.07 ± 0.82) and median (0.01) FE of non-construction sequences in the Spoken BNC dialogues.¹³ On the other hand, the posterior predictive mean value of S for constructions (Figure 1a) does not include the mean (5.59 ± 2.36) nor the median (5.36) S of non-construction sequences.

Our second observation is that the two metrics show similar but opposite patterns in our results. Based on the definition of the two metrics (Section 6.1)—these trends can be predicted a priori: it is more likely for a construction to have a facilitating effect if its surprisal is low; if construction surprisal is high, the context of the construction must be even more surprising for facilitating effect

¹³We calculate S and FE of all 3- to 7-grams in our analysis split of the Spoken BNC, excluding all n -grams that are equal to extracted constructions. We then sample, for each length n from 3 to 7, s_n non-construction sequence occurrences—where s_n is the number of occurrences of n -tokens-long extracted constructions. The length distributions should match because length has an effect on S and FE (see Section 8.2).

to occur. Empirically, we find that the Kendall’s rank-correlation between facilitating effect and surprisal is -0.569 ($p < 2e - 16$): although this is a rather strong negative correlation, the fact that the score is not closer to -1 indicates that there are cases where the two values do not follow the predicted pattern. Some constructions have high surprisal and high facilitating effect:

A: So what have you got? what have you got going on with enrichments?
 B: **I have to do** drama enrichment ($S = 5.46$ $FE = 1.32$)

While there are cases where construction surprisal is low and facilitating effect is low or negative:¹⁴

A: But like I always really love strawberries but hate strawberry-flavoured things so I don’t
 B: I don’t like strawberries **but I like** strawberry-flavoured things ($S = 2.24$ $FE = -0.70$)

These examples show that our measures capture different types of context-dependent processing advantage.¹⁵

8.2 Other Predictors of Processing Advantage

Other factors that influence facilitating effect and surprisal beyond those directly related to our hypotheses are construction length and repetition index within a dialogue turn. Construction length has the strongest effect on both metrics (S : $-110.90e - 2$, $p < 2e - 16$; FE : $30.16e - 2$, $p < 2e - 16$): the longer the construction the lower its surprisal and the stronger its facilitating effect. Table 4 shows a full repetition chain for a construction of length 3; Table 5 (Appendix B) shows a chain for one of length 6. Because constructions, per se, have a processing advantage, and their repetition facilitates processing (see Section 7), construction repetition is more advantageous when constructions occupy

¹⁴A negative facilitating effect indicates that the surprisal of the construction is higher than the surprisal of its context.

¹⁵The examples have been selected among occurrences with S and FE higher or lower than the mean $S / FE \pm$ std.

540 a larger portion of processing time (which is pro- 587
541 portional to the number of words). 588

542 The repetition index of a construction mention 589
543 *within a dialogue turn* also has an effect on both 590
544 metrics of processing advantage ($S: -29.48e - 2, p < 0.05; FE: 14.38e - 2, p < 2e - 16$); we 591
545 find strong cumulativity effects for self-repetitions 592
546 within the current dialogue turn.¹⁶ Only 6.46% of 593
547 the total construction occurrences have at least one 594
548 previous mention in the same turn; yet when this 595
549 is the case, the magnitude of S and FE increases 596
550 with the number of previous local mentions. This 597
551 interaction between cumulativity and recency (me- 598
552 dian distance between repetitions in the same turn 599
553 is 7 words; across turns is 1208 words) indicates 600
554 that processing advantage accumulates faster when 601
555 repetitions are densely clustered. See Appendix E. 602
556

557 9 Conclusion 604

558 We have hypothesised that speakers repeat lexi- 605
559 calised constructions in dialogues because repeti- 606
560 tion eases information processing, and have formu- 607
561 lated concrete predictions that follow from this 608
562 hypothesis. To quantify the processing advantage 609
563 of constructions we have proposed two surprisal- 610
564 based measures, facilitating effect and construc- 611
565 tion surprisal, and have analysed how the values of 612
566 these measures—estimated with a neural language 613
567 model—vary as constructions are repeated. Al- 614
568 though our experiments do not rely on direct mea- 615
569 surements of the processing effort of human sub- 616
570 jects, there is evidence that neural language models 617
571 produce reliable estimates (Goodkind and Bicknell, 618
572 2018; Linzen, 2019; Schrimpf et al., 2021). 619

573 Our experiments on English spoken open do- 620
574 main dialogues confirmed our three predictions: 621
575 (i) construction repetition reduces processing ef- 622
576 fort; (ii) the effort reduction increases with the 623
577 frequency of repetitions and (iii) decreases with 624
578 the distance between repetitions. These empiri- 625
579 cal results provide new evidence that construction 626
580 repetition in dialogue is an efficient communica- 627
581 tion strategy. They thus complement prior work 628
582 on the processing advantage of construction usage 629
583 (cf. Section 2.1) and contribute to an understudied 630
584 type of priming, with priming research tradition- 631
585 ally focusing on repetitions of syntactic structures 632
586 and lexical elements (cf. Section 1). Our findings reveal 633

¹⁶The identity of the speaker producing previous mentions does not influence FE or S . All fixed effects related to speaker identity are discarded during the ablation procedure; see Section 6.3 and Appendix D. 634

587 that the information processing efficiency of con- 588
589 struction repetition results from a combination of 590
591 the activation-based and implicit learning priming 592
593 mechanisms. In line with activation-based accounts 594
595 of priming, we find that the processing advantage 596
597 of repetitions accumulates faster when repetitions 598
599 are densely clustered, and it decays faster within 600
601 more local distances. However, implicit learning is 602
603 necessary to explain the fact that both cumulativity 604
605 and decay effects are still present across distant 606
607 repetitions. The discovered decreasing patterns of 608
609 surprisal may seem to contradict the entropy rate 609
610 constancy principle (Genzel and Charniak, 2002) 610
611 and the principle of uniform information density 611
612 (Jaeger and Levy, 2007), according to which sur- 612
613 prisal remains stable over consecutive utterances. 613
614 Yet we believe our findings can help explain these 614
615 principles by providing insights into the informa- 615
616 tion structure of individual utterances: the process- 616
617 ing advantage of repeated constructions, which are 617
618 not topic related, allows for progressively more 618
619 information-dense topical and referential expres- 619
620 sions. We conjecture that it is as a result of this bal- 620
621 ance that surprisal remains stable over utterances. 621
622

623 Besides contributing new empirical evidence 624
625 on construction usage and repetition in dialogue, 625
626 this study highlights the importance of a few key 626
627 desiderata for the design of human-compatible 627
628 computational dialogue models. First, models 628
629 should both attend to the local dialogue context 629
630 and use the global statistics collected throughout 630
631 a dialogue for on-the-fly adaptation. This would 631
632 have the natural effect of models being more likely 632
633 to repeat constructions established as part of the 633
634 dialogue lexicon. Second, although excessive and 634
635 unnatural repetitions should be avoided in machine- 635
636 generated utterances (Li et al., 2016; Holtzman 636
637 et al., 2019), a certain degree of repetition makes 637
638 a dialogue sound more natural. Human-like repeti- 638
639 tion patterns can be explicitly learned by auxiliary 639
640 modules (Holtzman et al., 2018) or, as our study 640
641 suggests, they may be implicitly acquired if next- 641
642 word surprisal training and decoding objectives are 642
643 complemented with context-dependent surprisal- 643
644 based objectives. Simple techniques such as those 644
645 proposed by Wei et al. (2021) and Meister et al. 645
646 (2020) could be used to operationalise facilitating 646
647 effect as a psycholinguistically motivated inductive 647
648 bias to be used in training, and as a word choice 648
649 criterion in decoding. 649
650

References

Inbal Arnon and Neal Snider. 2010. More than words: Frequency effects for multi-word phrases. *Journal of memory and language*, 62(1):67–82.

Matthew Aylett and Alice Turk. 2004. The smooth signal redundancy hypothesis: A functional explanation for relationships between redundancy, prosodic prominence, and duration in spontaneous speech. *Language and Speech*, 47(1):31–56.

Matthew Aylett and Alice Turk. 2006. Language redundancy predicts syllabic duration and the spectral characteristics of vocalic syllable nuclei. *The Journal of the Acoustical Society of America*, 119(5):3048–3058.

Colin Bannard and Danielle Matthews. 2008. Stored word sequences in language learning: The effect of familiarity on children’s repetition of four-word combinations. *Psychological science*, 19(3):241–248.

Alan Bell, Daniel Jurafsky, Eric Fosler-Lussier, Cynthia Girand, Michelle Gregory, and Daniel Gildea. 2003. Effects of disfluencies, predictability, and utterance position on word form variation in English conversation. *The Journal of the Acoustical Society of America*, 113(2):1001–1024.

Lasse Bergroth, Harri Hakonen, and Timo Raita. 2000. A survey of longest common subsequence algorithms. In *Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000*, pages 39–48. IEEE.

Douglas Biber and Federica Barbieri. 2007. Lexical bundles in university spoken and written registers. *English for specific purposes*, 26(3):263–286.

Douglas Biber, Susan Conrad, and Viviana Cortes. 2004. If you look at...: Lexical bundles in university teaching and textbooks. *Applied linguistics*, 25(3):371–405.

Emmanuel Bigand, Barbara Tillmann, Bénédicte Poulin-Charronnat, and D Manderlier. 2005. Repetition priming: Is music special? *The Quarterly Journal of Experimental Psychology Section A*, 58(8):1347–1375.

J Kathryn Bock. 1986. Syntactic persistence in language production. *Cognitive Psychology*, 18(3):355–387.

Kathryn Bock and Zenzi M Griffin. 2000. The persistence of structural priming: Transient activation or implicit learning? *Journal of Experimental Psychology: General*, 129(2):177.

Holly P Branigan, Martin J Pickering, and Alexandra A Cleland. 1999. Syntactic priming in written production: Evidence for rapid decay. *Psychonomic Bulletin & Review*, 6(4):635–640.

Holly P Branigan, Martin J Pickering, Andrew J Stewart, and Janet F McLean. 2000. Syntactic priming in spoken production: Linguistic and temporal interference. *Memory & Cognition*, 28(8):1297–1302. 689
690
691
692

Susan E Brennan. 1996. Lexical entrainment in spontaneous dialog. *Proceedings of ISSD*, 96:41–44. 693
694

Susan E. Brennan and Herbert H. Clark. 1996. Conceptual pacts and lexical choice in conversation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 22:1482–1493. 695
696
697
698

Joan Bybee. 2006. From usage to grammar: The mind’s response to repetition. *Language*, pages 711–733. 699
700

Joan Bybee. 2010. *Language, usage and cognition*. Cambridge University Press. 701
702

Joan Bybee and Joanne Scheibman. 1999. The effect of usage on degrees of constituency: The reduction of don’t in English. *Linguistics*, 37(4):575–596. 703
704
705

Gareth Carrol and Kathy Conklin. 2020. Is all formulaic language created equal? Unpacking the processing advantage for different types of formulaic sequences. *Language and Speech*, 63(1):95–122. 706
707
708
709

Herbert H. Clark and Edward F. Schaefer. 1989. Contributing to discourse. *Cognitive Science*, 13(2):259–294. 710
711
712

Herbert H. Clark and Deanna Wilkes-Gibbs. 1986. Referring as a collaborative process. *Cognition*, 22(1):1–39. 713
714
715

Meghan Clayards, Michael K. Tanenhaus, Richard N. Aslin, and Robert A. Jacobs. 2008. Perception of speech reflects optimal use of probabilistic speech cues. *Cognition*, 108(3):804–809. 716
717
718
719

Axel Cleeremans and Jeffrey Elman. 1993. *Mechanisms of implicit learning: Connectionist models of sequence processing*. MIT press. 720
721
722

Alexandra A Cleland and Martin J Pickering. 2003. The use of lexical and syntactic information in language production: Evidence from the priming of noun-phrase structure. *Journal of Memory and Language*, 49(2):214–230. 723
724
725
726
727

Georgie Columbus. 2013. In support of multiword unit classifications: Corpus and human rating data validate phraseological classifications of three different multiword unit types. *Yearbook of Phraseology*, 4(1):23–44. 728
729
730
731
732

Kathy Conklin and Norbert Schmitt. 2012. The processing of formulaic language. *Annual Review of Applied Linguistics*, 32:45–61. 733
734
735

Cristian Danescu-Niculescu-Mizil, Lillian Lee, Bo Pang, and Jon Kleinberg. 2012. Echoes of power: Language effects and power differences in social interaction. In *Proceedings of the 21st international conference on World Wide Web*, pages 699–708. 736
737
738
739
740

741	Vera Demberg and Frank Keller. 2008. Data from eye-tracking corpora as evidence for theories of syntactic processing complexity. <i>Cognition</i> , 109(2):193–210.	Riccardo Fusaroli, Joanna Rączaszek-Leonardi, and Kristian Tylén. 2014. Dialog as interpersonal synergy. <i>New Ideas in Psychology</i> , 32:147–157.	795 796 797
744	Vera Demberg, Asad Sayeed, Philip Gorinski, and Nikolaos Engonopoulos. 2012. Syntactic surprisal affects spoken word duration in conversational contexts. In <i>Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning</i> , pages 356–367.	Dmitriy Genzel and Eugene Charniak. 2002. Entropy rate constancy in text. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 199–206.	798 799 800 801
751	Nina Dethlefs, Helen Hastie, Heriberto Cuayáhuítl, Yanchao Yu, Verena Rieser, and Oliver Lemon. 2016. Information density and overlap in spoken dialogue. <i>Computer speech & language</i> , 37:82–97.	Dmitriy Genzel and Eugene Charniak. 2003. Variation of entropy and parse trees of sentences as a function of the sentence number. In <i>Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 65–72.	802 803 804 805 806
755	Gabriel Doyle and Michael C Frank. 2016. Investigating the sources of linguistic alignment in conversation. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 526–536.	Edward Gibson. 1998. Linguistic complexity: Locality of syntactic dependencies. <i>Cognition</i> , 68(1):1–76.	807 808
760	Guillaume Dubuisson Duplessis, Franck Charras, Vincent Letard, Anne-Laure Ligozat, and Sophie Rosset. 2017a. Utterance retrieval based on recurrent surface text patterns. In <i>European Conference on Information Retrieval</i> , pages 199–211. Springer.	Edward Gibson, Leon Bergen, and Steven T. Piantadosi. 2013. Rational integration of noisy evidence and prior semantic expectations in sentence interpretation. <i>Proceedings of the National Academy of Sciences</i> , 110(20):8051–8056.	809 810 811 812 813
765	Guillaume Dubuisson Duplessis, Chloé Clavel, and Frédéric Landragin. 2017b. Automatic measures to characterise verbal alignment in human-agent interaction. In <i>18th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)</i> , pages 71–81.	Adele E Goldberg. 2006. <i>Constructions at work: The nature of generalization in language</i> . Oxford University Press on Demand.	814 815 816
771	Guillaume Dubuisson Duplessis, Caroline Langlet, Chloé Clavel, and Frédéric Landragin. 2021. Towards alignment strategies in human-agent interactions based on measures of lexical repetitions. <i>Language Resources and Evaluation</i> , 55(2):353–388.	Adam Goodkind and Klinton Bicknell. 2018. Predictive power of word surprisal for reading times is a linear function of language model quality . In <i>Proceedings of the 8th Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2018)</i> , pages 10–18, Salt Lake City, Utah. Association for Computational Linguistics.	817 818 819 820 821 822 823
776	Jeffrey L Elman. 1991. Distributed representations, simple recurrent networks, and grammatical structure. <i>Machine learning</i> , 7(2):195–225.	Robert J Hartsuiker, Sarah Bernolet, Sofie Schoonbaert, Sara Speybroeck, and Dieter Vanderelst. 2008. Syntactic priming persists while the lexical boost decays: Evidence from written and spoken dialogue. <i>Journal of Memory and Language</i> , 58(2):214–238.	824 825 826 827 828
779	Alex B Fine and T Florian Jaeger. 2013. Evidence for implicit learning in syntactic comprehension. <i>Cognitive Science</i> , 37(3):578–591.	Patrick GT Healey, Matthew Purver, and Christine Howes. 2014. Divergence in dialogue. <i>PloS one</i> , 9(6):e98598.	829 830 831
782	Alex B Fine and T Florian Jaeger. 2016. The role of verb repetition in cumulative structural priming in comprehension. <i>Journal of Experimental Psychology: Learning, Memory, and Cognition</i> , 42(9):1362.	Daniel S Hirschberg. 1977. Algorithms for the longest common subsequence problem. <i>Journal of the ACM (JACM)</i> , 24(4):664–675.	832 833 834
786	Alex B Fine, T Florian Jaeger, Thomas A Farmer, and Ting Qian. 2013. Rapid expectation adaptation during syntactic comprehension. <i>PloS one</i> , 8(10):e77661.	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In <i>International Conference on Learning Representations</i> .	835 836 837 838
790	Austin F. Frank and T. Florian Jaeger. 2008. Speaking rationally: Uniform information density as an optimal strategy for language production. In <i>Proceedings of the Annual Meeting of the Cognitive Science Society</i> .	Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. 2018. Learning to write with cooperative discriminators. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1638–1649.	839 840 841 842 843 844
794		Christine Howes, Patrick GT Healey, and Matthew Purver. 2010. Tracking lexical and syntactic alignment in conversation. In <i>Proceedings of the Annual Meeting of the Cognitive Science Society</i> .	845 846 847 848

849	T. Florian Jaeger. 2010. Redundancy and reduction: Speakers manage syntactic information density. <i>Cognitive Psychology</i> , 61(1):23–62.	905
850		906
851		907
852	T. Florian Jaeger and Roger P. Levy. 2007. Speakers optimize information density through syntactic reduction. In <i>Advances in neural information processing systems</i> , pages 849–856.	908
853		909
854		910
855		911
856	T Florian Jaeger and Neal Snider. 2008. Implicit learning and syntactic persistence: Surprisal and cumulatvity. In <i>Proceedings of the 30th Annual Conference of the Cognitive Science Society</i> , volume 827812. Cognitive Science Society Austin, TX.	912
857		913
858		914
859		915
860		916
861	Niels Janssen and Horacio A Barber. 2012. Phrase frequency effects in language production. <i>PLoS one</i> , 7(3):e33202.	917
862		918
863		919
864	Frederick Jelinek, Lalit Bahl, and Robert Mercer. 1975. Design of a linguistic statistical decoder for the recognition of continuous speech. <i>IEEE Transactions on Information Theory</i> , 21(3):250–256.	920
865		
866		
867		
868	Hajnal Jolsvai, Stewart M McCauley, and Morten H Christiansen. 2013. Meaning overrides frequency in idiomatic and compositional multiword chunks. In <i>Proceedings of the Annual Meeting of the Cognitive Science Society</i> .	921
869		922
870		923
871		924
872		925
873	Michael P Kaschak, Timothy J Kutta, and Jacqueline M Coyle. 2014. Long and short term cumulative structural priming effects. <i>Language, cognition and neuroscience</i> , 29(6):728–743.	926
874		927
875		928
876		
877	Michael P Kaschak, Renrick A Loney, and Kristin L Borreggine. 2006. Recent experience affects the strength of structural priming. <i>Cognition</i> , 99(3):B73–B82.	929
878		930
879		931
880		932
881	Frank Keller. 2004. The entropy rate principle as a predictor of processing effort: An evaluation against eye-tracking data. In <i>Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 317–324.	933
882		934
883		935
884		
885		
886	Willem JM Levelt and Stephanie Kelter. 1982. Surface form and memory in question answering. <i>Cognitive Psychology</i> , 14(1):78–106.	936
887		937
888		938
889	Roger Levy. 2008. A noisy-channel model of human sentence comprehension under uncertain input. In <i>Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 234–243.	939
890		940
891		941
892		942
893		943
894	Roger Levy, Klinton Bicknell, Tim Slattery, and Keith Rayner. 2009. Eye movement evidence that readers maintain and act on uncertainty about past linguistic input. <i>Proceedings of the National Academy of Sciences</i> , 106(50):21086–21090.	944
895		945
896		946
897		947
898		
899	Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. 2016. Deep reinforcement learning for dialogue generation. In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1192–1202.	948
900		949
901		950
902		951
903		952
904		
	Tal Linzen. 2019. What can linguistics and deep learning contribute to each other? response to pater. <i>Language</i> , 95(1):e99–e108.	953
		954
		955
		956
		957
	Robbie Love, Claire Dembry, Andrew Hardie, Vaclav Brezina, and Tony McEnery. 2017. The spoken BNC2014. <i>International Journal of Corpus Linguistics</i> , 22(3):319–344.	
	Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2173–2185.	
	Danny Merckx and Stefan L Frank. 2021. Human sentence processing: Recurrence or attention? In <i>Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics</i> , pages 12–22.	
	Bill Noble and Raquel Fernández. 2015. Centre stage: How social network position shapes linguistic coordination. In <i>Proceedings of the 6th workshop on cognitive modeling and computational linguistics</i> , pages 29–38.	
	M. J. Pickering and S. Garrod. 2004. Toward a mechanistic psychology of dialogue. <i>Behavioral and Brain Sciences</i> , 27(02):169–190.	
	Martin J Pickering and Holly P Branigan. 1998. The representation of verbs: Evidence from syntactic priming in language production. <i>Journal of Memory and language</i> , 39(4):633–651.	
	Ting Qian and T. Florian Jaeger. 2011. Topic shift in efficient discourse production. In <i>Proceedings of the Annual Meeting of the Cognitive Science Society</i> .	
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. Technical report, OpenAI.	
	David Reitter, Frank Keller, and Johanna D Moore. 2006a. Computational modelling of structural priming in dialogue. In <i>Proceedings of the Human Language Technology Conference of the NAACL, companion volume: Short papers</i> , pages 121–124.	
	David Reitter, Frank Keller, and Johanna D Moore. 2011. A computational cognitive model of syntactic priming. <i>Cognitive science</i> , 35(4):587–637.	
	David Reitter, Johanna D Moore, and Frank Keller. 2006b. Priming of syntactic rules in task-oriented dialogue and spontaneous conversation. <i>Proceedings of the 28th Annual Conference of the Cognitive Science Society</i> .	
	DE Rumelhart and JL McClelland. 1986. On learning the past tenses of English verbs. In <i>Parallel distributed processing: explorations in the microstructure, vol. 2: psychological and biological models</i> , pages 216–271. MIT press Cambridge, MA.	

958	Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2021. The neural architecture of language: Integrative modeling converges on predictive processing. <i>Proceedings of the National Academy of Sciences</i> , 118(45).	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. <i>Advances in Neural Information Processing Systems</i> , 30:5998–6008.	1011
959			1012
960			1013
961			1014
962			1015
963			
964	Arabella Sinclair and Raquel Fernández. 2021. Construction coordination in first and second language acquisition . In <i>Proceedings of the 25th Workshop on the Semantics and Pragmatics of Dialogue - Full Papers</i> , Potsdam, Germany. SEMDIAL.	Jason Wei, Clara Meister, and Ryan Cotterell. 2021. A cognitive regularizer for language modeling. <i>arXiv preprint arXiv:2105.07144</i> .	1016
965			1017
966			1018
967		Sean Welleck, Iliia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2019. Neural text generation with unlikelihood training. In <i>International Conference on Learning Representations</i> .	1019
968			1020
969	Anna Siyanova-Chanturia, Kathy Conklin, Sendy Cafarra, Edith Kaan, and Walter JB van Heuven. 2017. Representation and processing of multi-word expressions in the brain. <i>Brain and language</i> , 175:111–122.		1021
970			1022
971		Linda Wheeldon and Mark Smith. 2003. Phrase structure priming: A short-lived effect. <i>Language and Cognitive Processes</i> , 18(4):431–442.	1023
972			1024
973	Anna Siyanova-Chanturia, Kathy Conklin, and Walter JB Van Heuven. 2011. Seeing a phrase “time and again” matters: The role of phrasal frequency in the processing of multiword sequences. <i>Journal of Experimental Psychology: Learning, Memory, and Cognition</i> , 37(3):776.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	1025
974			1026
975			1027
976			1028
977			1029
978			1030
979	Patrizia Tabossi, Rachele Fanari, and Kinou Wolf. 2009. Why are idioms recognized fast? <i>Memory & Cognition</i> , 37(4):529–540.		1031
980			1032
981			1033
982	Debra Titone and Maya Libben. 2014. Time-dependent effects of decomposability, familiarity and literal plausibility on idiom priming: A cross-modal priming investigation. <i>The Mental Lexicon</i> , 9(3):473–496.		1034
983			1035
984			1036
985		Alison Wray. 2002. <i>Formulaic Language and the Lexicon</i> . Cambridge, UK: Cambridge University Press.	1037
986	Debra A Titone and Cynthia M Connine. 1994. Descriptive norms for 171 idiomatic expressions: Familiarity, compositionality, predictability, and literality. <i>Metaphor and Symbol</i> , 9(4):247–270.		1038
987			1039
988		Yang Xu and David Reitter. 2018. Information density converges in dialogue: Towards an information-theoretic model. <i>Cognition</i> , 170:147–163.	1040
989			1041
990	Michael Tomasello. 2003. <i>Constructing a language: A usage-based theory of language acquisition</i> . Harvard University Press.		1042
991			1043
992		Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. Dialogpt: Large-scale generative pre-training for conversational response generation. In <i>ACL, system demonstration</i> .	1044
993	Antoine Tremblay and R Harald Baayen. 2010. Holistic processing of regular four-word sequences: A behavioral and ERP study of the effects of structure, frequency, and probability on immediate free recall. <i>Perspectives on formulaic language: Acquisition and communication</i> , pages 151–173.		1045
994			1046
995			1047
996			1048
997			
998			
999	Antoine Tremblay, Bruce Derwing, Gary Libben, and Chris Westbury. 2011. Processing advantages of lexical bundles: Evidence from self-paced reading and sentence recall tasks. <i>Language learning</i> , 61(2):569–613.		1049
1000			1050
1001			1051
1002			1052
1003			1053
1004	G Underwood. 2004. The eyes have it. An eye-movement study into the processing of formulaic sequences.		1054
1005			1055
1006			1056
1007	Marten van Schijndel and Tal Linzen. 2018. A neural model of adaptation in reading. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4704–4710.		1057
1008			1058
1009			1059
1010			1060
			1061
			1062

1063	phrases (Tabossi et al., 2009; Jolsvai et al.,	communicative functions. Biber et al. (2004),	1111
1064	2013; Carrol and Conklin, 2020). Therefore	e.g., distinguish between stance expressions	1112
1065	we ignore this criterion in the current study.	(attitude, certainty with respect to a proposi-	1113
1066	• Literal plausibility This criterion is typically	tion), discourse organisers (connecting prior	1114
1067	used to discriminate among different types	and forthcoming discourse), and referential	1115
1068	of idioms (Titone and Connine, 1994; Titone	expressions; and for each of these three pri-	1116
1069	and Libben, 2014)—as compositional phrases	mary discourse functions, more specific sub-	1117
1070	are literally plausible by definition. Because	categories are defined. This type of classi-	1118
1071	we ignore distinctions made on the basis of	fication is typically done a posteriori—i.e.,	1119
1072	compositionality, we do not use this criterion.	after a manual analysis of the expressions re-	1120
1073	• Meaningfulness Meaningful expressions are	trieved from a corpus according to other cri-	1121
1074	idioms and compositional phrases (e.g. ‘on	teria (Biber and Barbieri, 2007). In the BNC,	1122
1075	my mind’, ‘had a dream’) whereas sentence	for example, we find epistemic lexical bun-	1123
1076	fragments that break constituency boundaries	dles (‘I don’t know’, ‘I don’t think’), desire	1124
1077	(e.g., ‘of a heavy’, ‘by the postal’) are consid-	bundles (‘do you want to’, ‘I don’t want to’),	1125
1078	ered less meaningful (as measured in norming	obligation/directive bundles (‘you don’t have	1126
1079	studies, e.g., by Jolsvai et al., 2013). There	to’), and intention/prediction bundles (‘I’m	1127
1080	is some evidence that the meaningfulness of	going to’, ‘it’s gonna be’). We do not use this	1128
1081	multi-word expressions correlates with their	criterion to avoid an a priori selection of the	1129
1082	processing advantage even more than their	constructions.	1130
1083	frequency (Jolsvai et al., 2013); yet expres-	B Extraction of Repeated Constructions	1131
1084	sions are particularly frequent, they present	We define a limited specific vocabulary of generic	1132
1085	processing advantages even if they break reg-	nouns to filter out topical and referential construc-	1133
1086	ular phrasal structures (Bybee and Scheibman,	tion. The vocabulary includes: <i>bit, bunch, day,</i>	1134
1087	1999; Tremblay et al., 2011). Moreover, ut-	<i>days, fact, god, idea, ideas, kind, kinds, loads, lot,</i>	1135
1088	terances that break regular constituency rules	<i>lots, middle, ones, part, problem, problems, reason,</i>	1136
1089	are particularly frequent in spoken dialogue	<i>reasons, rest, side, sort, sorts, stuff, thanks, thing,</i>	1137
1090	data (e.g., ‘if you could search for job and	<i>things, time, times, way, ways, week, weeks, year,</i>	1138
1091	that’s not’, ‘you don’t wanna damage your	<i>years.</i>	1139
1092	relationship with’). For these reasons, we do	We also find all the filled pauses and exclude	1140
1093	not exclude constructions that span multiple	word sequences that consist for more than 50% of	1141
1094	constituents from our analysis.	filled pauses. Filled pauses in the Spoken BNC are	1142
1095	• Schematicity This criterion distinguishes ex-	transcribed as: <i>huh, uh, erm, hm, mm, er.</i>	1143
1096	pressions where all the lexical elements are	Table 5 shows a whole construction chain (from	1144
1097	fixed from expressions “with slots” that can be	the first mention to the last repetition) for a con-	1145
1098	filled by varying lexical elements. In this study,	struction of length 6.	1146
1099	we focus on fully lexicalised constructions.	C Language Model	1147
1100	• Familiarity This is a subjective criterion that	C.1 Finetuning	1148
1101	strongly correlates with objective frequency	We finetune the ‘small’ variant of GPT-2 (Radford	1149
1102	(Carrol and Conklin, 2020). Human experi-	et al., 2019) and DialoGPT (Zhang et al., 2020)	1150
1103	ments would be required to obtain familiarity	on our finetuning split of the Spoken BNC (see	1151
1104	norms for our target data, and the resulting	Section 4) using HuggingFace’s implementation of	1152
1105	norms would only be an approximation of the	the models with default tokenizers and parameters	1153
1106	familiarity judgements of the true speakers we	(Wolf et al., 2020). Dialogue turns are simply con-	1154
1107	analyse the language of. Therefore, we ignore	catenated; we have experimented with labelling the	1155
1108	this criterion in the current study.	dialogue turns (i.e., A: utterance 1, B: utterance 2	1156
1109	• Communicative function Formulaic expres-	and found that this leads to higher perplexity. The	1157
1110	sions can fulfil a variety of discourse and	finetuning results for both models are presented in	1158
		Table 6. We finetune the models and measure their	1159

Speaker	RI	RI Turn	Dist	Turn	<i>S</i>	<i>FE</i>
A	0	0	-	[...] I think that everyone should have the same opportunities and I don't think you should be proud or ashamed of what your you know what your situation is whether you what your what your race is whether you're a woman or a man whether you live from this pl whether you're in this place [...]	1.90	1.21
A	1	0	80	I well I th I don't think it should I don't think you should be	1.73	1.40
A	2	0	19	Well yes perhaps but I don't think you should be like um embarrassed about it or I think I think you should just sort of	1.06	2.48

Table 5: A chain of repetitions of the construction ‘*I don't think you should be*’ in dialogue S2AX of the Spoken BNC, annotated with repetition index (RI), repetition index within dialogue turn (RI Turn), and distance from previous mention (Dist; in tokens).

perplexity using Huggingface’s finetuning script. We use early stopping over 5 epochs.¹⁷ Sequence length and batch size vary together because they together determine the amount of memory required; more expensive combinations (e.g., 256 tokens with batch size 16) require an exceedingly high amount of GPU memory. Reducing the maximum sequence length has limited impact: 99.90% of dialogue turns have at most 128 words.

DialoGPT starts from extremely high perplexity values but catches up quickly with finetuning. GPT-2 starts from much lower perplexity values and reaches virtually the same perplexity as DialoGPT after finetuning. For the pre-trained DialoGPT perplexity is extremely high, and the perplexity trend against maximum sequence length is surprisingly upward. These two behaviours indicate that the pre-trained DialoGPT is less accustomed than GPT-2 to the characteristics of our dialogue data. DialoGPT is trained on written online group conversations, while we use a corpus of transcribed spoken conversations between two speakers. In contrast, GPT-2 has been exposed to the genre of fiction, which contains scripted dialogues, and thus to a sufficiently similar language use. We select GPT-2 finetuned with a maximum sequence length of 128 and 512 as our best two models; these two models (which we now refer to as *frozen*) are used for the adaptive learning rate selection (Section C.2).

¹⁷The number of epochs (5) has been selected in preliminary experiments together with the learning rate ($1e-4$). In these experiments—which we ran for 40 epochs—we noticed that the $1e-4$ learning rate offers the best tradeoff of training time and perplexity out of four possible values: $1e-2$, $1e-3$, $1e-4$, $1e-5$. We obtained insignificantly lower perplexity values with a learning rate of $1e-5$, with significantly longer training time: 20 epochs for GPT-2 and 28 epochs for DialoGPT.

C.2 Learning rate selection

To find the appropriate learning rate for on-the-fly adaptation (see Section 6.2), we randomly select 18 dialogues D from the analysis split of the Spoken BNC and run an 18-fold cross-validation for a set of six candidate learning rates: $1e-5$, $1e-4$, ..., 1. We finetune the model on each dialogue using one of these learning rate values, and compute perplexity change 1) on the dialogue itself (to measure *adaptation*) as well as 2) on the remaining 17 dialogues (to measure *generalisation*). We set the Transformer’s context window to 50 to reproduce the experimental conditions presented in Section 6.1.

More precisely, for each dialogue $d \in D$, we calculate the perplexity of our two frozen models (Section C.1) on d and D ($ppl_{before}(d)$ and $ppl_{before}(D)$, respectively). Then, we finetune the models on d using the six candidate learning rates, and measure again the perplexity over d and D ($ppl_{after}(d)$ and $ppl_{after}(D)$). The change in performance is evaluated according to two metrics: $\frac{ppl_{after}(d) - ppl_{before}(d)}{ppl_{before}(d)}$ measures the degree to which the model has successfully adapted to the target dialogue; $\frac{ppl_{after}(D) - ppl_{before}(D)}{ppl_{before}(D)}$ measures whether finetuning on the target dialogue has caused any loss of generalisation.

The learning rate selection results are presented in Figure 3. We select $1e-3$ as the best learning rate and pick the model finetuned with a maximum sequence length of 512 as our best model. The difference in perplexity reduction (both adaptation and generalisation) is minimal with respect to the model finetuned with a maximum sequence length of 128, but since the analysis split of the Spoken

Model	Learning rate	Max sequence length	Batch size	Best epoch	Perplexity finetuned	Perplexity pretrained
DialoGPT	0.0001	128	16	3	23.21	7091.38
DialoGPT	0.0001	256	8	4	22.26	12886.92
DialoGPT	0.0001	512	4	4	21.73	21408.32
GPT-2	0.0001	128	16	4	23.32	173.76
GPT-2	0.0001	256	8	3	22.21	159.23
GPT-2	0.0001	512	4	3	21.55	149.82

Table 6: Finetuning results for GPT-2 and DialoGPT on our finetuning split of the Spoken BNC.

BNC contains turns longer than 128 tokens, we select the 512 version. Similarly to van Schijndel and Linzen (2018), we find that finetuning on a dialogue does not cause a loss in generalisation but instead helps the model generalise to other dialogues. Unlike (2018), who used LSTM language models, we find that learning rates larger than $1e-1$ cause backpropagation to overshoot, even within a single dialogue. In Figure 3, the bars for $1e-1$ and 1 are not plotted because the corresponding data contains infinite perplexity values (due to numerical overflow). The selected learning rate, $1e-3$, is a relatively low learning rate for on-the-fly adaptation but it is still higher than the best learning rate for the entire dataset by a factor of 10.

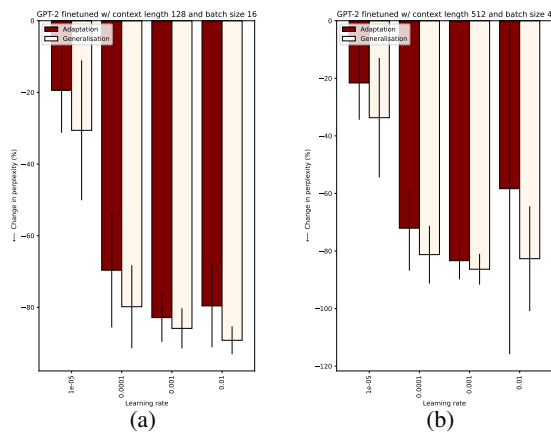


Figure 3: The adaptation and generalisation performance (defined in Section C.2) with varying learning rate.

D Linear Mixed Effect Models

As explained in Section 6.3 of the main paper, we fit linear mixed effect models using facilitating effect and construction surprisal as response variables and including multilevel random effects grouped by dialogues and individual speakers.¹⁸ To select

¹⁸We also try grouping observations only by dialogue and only by individual speakers. The amount of variance explained

the fixed effects of the models, we start with a collection of motivated features and perform an ablation selection procedure, iteratively removing features with the lowest significance, and keeping only those that yield a p -value lower than 0.05. We start with the following features: the logarithm of the repetition index, the logarithm of the repetition index *within the current turn*, the logarithm of the distance¹⁹ from the previous mention (computed in three ways: with respect to the previous mention of any speaker, of the current speaker, and of the other speaker), the logarithm of construction length (measures as the number of tokens in a construction), the logarithm of the number of tokens between the current occurrence and the first mention of a construction, and binary features indicating whether the previous mention is by the current speaker, whether it is produced by the initiator of the construction, whether the construction has been already uttered by both speakers, and whether the previous mention is in the current dialogue turn.

The ablation selection procedure yields two models with the following fixed effects: log repetition index, log repetition index within the current dialogue turn, log distance from the previous mention (of any speaker), and log construction length. The best model for facilitating effect is summarised in Listing 1 and the best model for construction surprisal in Listing 2.

E Local Effects of Processing Advantage

Table 7 shows the distribution of repetition indices within the dialogue turn. An index of n indicates that n previous mentions of the construction take place in the current dialogue turn. Figures 4a

(but unaccounted for by the fixed effects) decreases, so we keep the two-level random effects.

¹⁹Distance is measured as the number of words between the first word of the current occurrence and the first word of the previous occurrence. We choose this strategy as there exist overlapping constructions and the distance values would be negative if we used the last word of the previous occurrence as a starting point to compute the distance.

Listing 1: Best linear mixed effect model for Facilitating Effect

```

Linear mixed model fit by REML. t-tests use Satterthwaite's method [
lmerModLmerTest]
Formula:
logFE10 ~ 1 + logLength + logRepIndexInTurn + logRepetitionIndex +
  logDistance + (1 | `Dialogue ID`/Speaker)
Data: data

REML criterion at convergence: 51869.1

Scaled residuals:
  Min       1Q   Median       3Q      Max
-7.3884 -0.6125 -0.0438  0.5574  8.4443

Random effects:
 Groups                Name                Variance Std.Dev.
 Speaker:`Dialogue ID` (Intercept) 0.006503 0.08064
 Dialogue ID           (Intercept) 0.006100 0.07810
 Residual              0.478766 0.69193
Number of obs: 24540, groups:
Speaker:`Dialogue ID`, 364; Dialogue ID, 185

Fixed effects:
              Estimate Std. Error      df t value Pr(>|t|)
(Intercept)  4.056e-01  5.335e-02 2.036e+04  7.603 3.02e-14
logLength    3.016e-01  2.901e-02 2.452e+04 10.394 < 2e-16
logRepIndexInTurn 1.438e-01  1.709e-02 2.451e+04  8.416 < 2e-16
logRepetitionIndex 7.569e-02  6.902e-03 2.360e+04 10.965 < 2e-16
logDistance  -4.290e-02  1.741e-03 2.309e+04 -24.638 < 2e-16

(Intercept)    ***
logLength      ***
logRepIndexInTurn ***
logRepetitionIndex ***
logDistance    ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Correlation of Fixed Effects:
      (Intr) lgLngt lgRIIT lgRptI
logLength -0.909
lgRpIndxInT -0.177 -0.008
lgRpttnIndx -0.291  0.067 -0.031
logDistance -0.342  0.030  0.563  0.095

```

Listing 2: Best linear mixed effect model for Construction Surprisal

Linear mixed model fit by REML. t-tests use Satterthwaite's method [lmerModLmerTest]
 Formula: S ~ 1 + logLength + logRepIndexInTurn + logRepetitionIndex + logDistance + (1 | 'Dialogue ID'/Speaker)
 Data: data

REML criterion at convergence: 78900.3

Scaled residuals:

Min	1Q	Median	3Q	Max
-3.0885	-0.6807	-0.0779	0.6062	6.5359

Random effects:

Groups	Name	Variance	Std.Dev.
Speaker: 'Dialogue ID'	(Intercept)	0.01282	0.1132
Dialogue ID	(Intercept)	0.04292	0.2072
Residual		1.43852	1.1994

Number of obs: 24540, groups:

Speaker: 'Dialogue ID', 364; Dialogue ID, 185

Fixed effects:

	Estimate	Std. Error	df	t value	Pr(> t)
(Intercept)	4.866e+00	9.319e-02	1.810e+04	52.215	<2e-16
logLength	-1.109e+00	5.033e-02	2.451e+04	-22.042	<2e-16
logRepIndexInTurn	-2.948e-01	2.964e-02	2.452e+04	-9.943	<2e-16
logRepetitionIndex	-2.485e-01	1.197e-02	2.346e+04	-20.761	<2e-16
logDistance	9.657e-02	3.028e-03	2.408e+04	31.889	<2e-16

(Intercept) ***
 logLength ***
 logRepIndexInTurn ***
 logRepetitionIndex ***
 logDistance ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Correlation of Fixed Effects:

	(Intr)	lgLngt	lgRIIT	lgRptI
logLength	-0.903			
lgRpIndxInT	-0.176	-0.007		
lgRpttnIndx	-0.289	0.068	-0.030	
logDistance	-0.339	0.031	0.563	0.096

Previous mentions in the current dialogue turn									
Tot	0	1	2	3	4	5	6	7	8
33103	30965	1872	188	46	16	11	3	1	1

Table 7: The distribution of repetition indices *within the dialogue turn*.

1279 and 4b show how facilitating effect and construc-
 1280 tion surprisal vary locally, for repetitions occurring
 within the same dialogue turn.

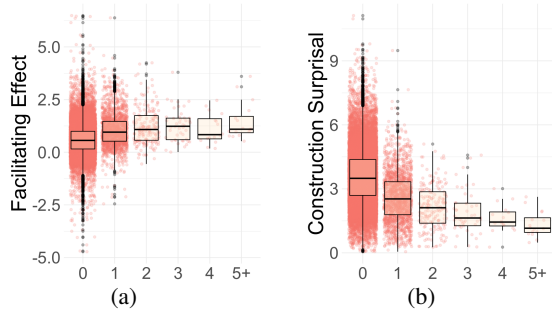


Figure 4: Facilitating effect and construction surprisal (bits) against repetition index *within the current dialogue turn*.

1281 **F Computing Infrastructure and Budget**

1282 Our experiments were carried out using a single
 1283 GPU on a computer cluster with Debian Linux OS.
 1284 The GPU nodes on the cluster are GPU GeForce
 1285 1001 1080Ti, 11GB GDDR5X, with NVIDIA
 1286 driver version 418.56 and CUDA version 10.1. The
 1287 total computational budget required to finetune the
 1288 language model amounts to 45 minutes; obtaining
 1289 surprisal estimates requires 4 hours, and selecting
 1290 the adaptation learning rate requires 9 hours.
 1291