# **Small Language Model Learning with Inconsistent Input**

**Anonymous ACL Submission** 

## Abstract

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

Modern language models, such as GPT-3, BERT, and LLaMA, are notoriously datahungry, requiring millions to over a trillion tokens of training data. Yet, transformer-based learning models have demonstrated a remarkable ability to learn natural languages. After sufficient training, they can consistently distinguish grammatical from ungrammatical sentences. Children as young as 14 months already have the capacity to learn abstract grammar rules from very few examples, even in the presence of non-rule-following exceptions. Yang's (2016) Tolerance Principle specifies an exact threshold for how many exceptions are permissible in a given dataset for a rule to be learnable by humans. We explore the minimal amount and quality of training data necessary for rules to be generalized by language models to test for evidence of Tolerance-Principle-like effects. We implement BabyBERTa (Huebner et al. 2021), a transformer-based language model optimized for training on smaller corpora than most LLMs. We train it on very small artificial grammar sets. For the simplest kind of rule, BabyBERTa can learn from datasets of under 1,000 tokens. The effect of type and token frequency of exemplars vs. exceptions on learning follows a gradient. We see no effect that can be related to the Tolerance Principle.

#### 1 Introduction

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

## 1.1 Tolerance Principle

Learning a rule from a set of examples, in an unsupervised (no feedback) setting, requires the ability to generalize the rule to novel instances unseen in the training set (Huebner et al. 2021). We say that a rule is *productive* if it is learnable from a set of examples. Given a sufficient set of examples, being able to determine whether a rule should or should not be productive—i.e., having a theory that can explain how it is that abstract rules are generalized—is a challenge relevant to linguists and cognitive scientists because, among other things, a good explanation of the ability to generalize would shed a great deal of light on the process of early language acquisition in humans.

One theory of rule generalization is the 47 Tolerance Principle (TP), originally proposed via 48 mathematical derivation in The Price of Linguistic 49 Productivity (Yang 2016). The Tolerance 50 Principle was derived as a necessary consequence 51 of a rule-ordering algorithm known as the 52 Elsewhere Condition (Anderson 1969, Kiparsky 53 1973), which Yang proposes as a cognitive model 54 of processing rules and exceptions. According to 55 the Elsewhere Condition, as applied to the human 56 brain, learning operates in an "exceptions-first, 57 rule-later" fashion. When encountering a new 58 exemplar and needing to decide whether to apply 59 a rule, the brain must first consider every known 60

1

exception to the rule (in order to see if this 61 exemplar is one of these exceptions) before the 62 general rule can be applied to it. When there are 63 very many exceptions and very few rule-64 following examples, it is more time-efficient to 65 just memorize each exemplar on a case-by-case 66 basis and not try to learn a rule at all. The 67 Tolerance Principle explores, mathematically, the 68 relationship between the number of exceptions 69 and the number of rule-following examples that 70 allows the brain to "optimize/minimize the time 71 complexity of language use," (Yang, 2016, p. 60). 72

The Tolerance Principle is well described in "A 73 User's Guide to the Tolerance Principle," (Yang 74 2018), but it is made most explicit in "A User's 75 Defense of the Tolerance Principle" with the 76 following excerpt: "The TP is first and foremost a 77 theory of learning. It specifies a precise threshold, 78 as a proportion of items in the learner's 79 experience, that a generalization can tolerate as 80 exceptions:  $\theta_N = N/lnN$ , where N is the 81 cardinality of the item set," (Yang, 2023, p. 2). 82 Yang makes the claim that the Tolerance Principle 83 is applicable to many kinds of learning where a 84 rule must be generalized despite the possibility of 85 exceptions, and is not explicitly limited to natural 86 language rules. 87

One key feature of the TP is that rule learning, according to the theory, does not occur gradually; it is instead quantal, meaning a rule is either productive or unproductive on a given set. In other words, given a sufficient number of examples, a learner should either be able to generalize a rule, or they will be entirely unable to do so, in which case they can only memorize the behavior of the examples they were given on a case-by-case basis, with no capacity to generalize beyond them. This applies to the learning of dominant rules over an entire set, and it also applies to the learning of subrules for subsets of a set. 100

88

89

90

91

92

93

94

95

96

97

98

99

Also crucial for a full understanding of the TP 101 is the fact that the set size N, as well as the number 102 of permissible exceptions  $e \leq \theta_N$ , both refer to 103 numbers of unique types of items, with no regard 104 to the number of repetitions of items of the same 105 type. In a linguistic context, this means the only 106 consideration that goes into productivity is the 107 type frequency (number of unique items), with no 108

109 regard to the token frequency (total number of 110 items, including repeated items of the same type). 111 So long as a learner is exposed to enough different 112 examples for rule learning to occur at all, the 113 number of repetitions of individual elements of 114 the example set will not affect the productivity of 115 the rule.

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

The Tolerance Principle's application to nonhuman learners has not been explored.

## 1.2 Rule Generalization in Human Infants

A long line of research in the laboratory of Rushen Shi has investigated the abilities of human infants to generalize abstract grammar rules. Koulaguina & Shi (2013) showed that infants as young as 14 months can generalize abstract grammar rules to novel instances from relatively little training (as few as 8 exemplar sentences, repeated four times). Koulaguina & Shi (2019) showed with 14-montholds that a training set that consisted of 50% rulefollowing and 50% non-rule- following sentences was insufficient for the word-order rule to be generalized, while a training set consisting of 80% rule-following and 20% non-rule-following was sufficient. They also found that it was the type frequency of the example set and not the token frequency that determined whether a word-order shift rule was productive.

Shi & Emond (2023) continued the above paradigm with more rigor, attempting to find a threshold of permissible exceptions beyond which generalizability would be impossible. They also investigated the gradual/quantal question. Their findings lent significant support to the Tolerance Principle.

## 1.3 Motivation

It is difficult to explain how or why 14-month-147 olds are so remarkably capable of generalizing 148 abstract rules to novel instances; however, 149 computational models are less of a black box than 150 a human brain. When a model uses unsupervised 151 learning to learn a rule from noisy or exception-152 filled data, is its learning governed by the 153 Tolerance Principle, or something like it? This 154 was the question that motivated our work. If it is 155 possible to show that models can do the same 156

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

157 thing human infants can do, examining how they158 do it might help explain how human infants do it.

159 The problem of explaining the capacity to learn 160 is also at the forefront of language model research 161 (see Contreras et al. 2023, Jawahar et al. 2019). 162 Whereas there is plenty of research on how LLMs 163 learn when they are provided with superhuman 164 amounts of data and training, there is very limited 165 research on their capacity to learn with small 166 amounts of training data. 167

## 1.4 Related Studies

168

185

169 A few efforts have been made to optimize LLMs 170 with the objective of achieving substantial 171 learning from developmentally feasible quantities 172 of training data. In the BabyLM challenge 173 (Warstadt et al. 2023), language models were 174 optimized to maximize learning with a training 175 data size of 10M words or less. Huebner et al. 176 (2021) developed the BabyBERTa transformer-177 based language model as a variation of 178 RoBERTa-base (Liu 2019) and pre-trained it on 179 as few as 5M words, simulating the input 180 available to children aged one to six years old. 181 Some of the best performers on the BabyLM 182 challenge used BabyBERTa. 183

184 **2** Implementation

## 2.1 Task

186 Our objective is to attempt to address the 187 following questions: (1) What is the minimal 188 amount of training data that our language model 189 needs in order to learn a rule? (2) How noisy can 190 this training data be? In other words, what 191 proportion of training data in the training set can 192 be non-rule-following for the rule to still be 193 learnable? What is the relation between this 194 proportion and the size of the dataset? (3) Is 195 productivity quantal or gradient? That is, if a 196 language model can generalize a rule to novel 197 instances, is there a gradient or quantal effect as 198 we move from unproductive regions of the 199 parameter space to productive regions? 200

## 2.2 Model Selection

## 2.2.1 Architecture

We implement BabyBERTa (Huebner et al. 2021), whose code is available on GitHub. BabyBERTa uses the Transformers architecture (Vaswani 2017) and is the result of a fine-tuning of the hyper-parameters of RoBERTa (Liu 2019).

BabyBERTa, in line with RoBERTa and differing from BERT (Devlin 2018), does not do next-sentence prediction. It is instead trained only on the masked language model (MLM) pretraining objective used by BERT. A new random subsample of tokens is selected for masking every epoch.

Unlike RoBERTa-base, BabyBERTa is trained exclusively on single sentences. This means that the prediction of masked tokens takes into account only the rest of the tokens in the same sentence as the masked token. The MLM procedure is a form of self-supervised learning.

## 2.2.2 Hyper-Parameters

Like the original BabyBERTa implementation, our model uses 8 layers, 8 attention heads, 256 hidden units, and an intermediate size of 1024. We use Adam optimizer (Kingma 2014) with a learning rate of 1e - 4. Batch size is set to 16. In creating a random subsample of tokens for masking, tokens are selected with a probability of 0.15.

## 2.2.3 Training Procedure

We train our model on a text (.txt) file. The primary reason we use transformers rather than another neural network architecture is to be able to train our model on sequential text data. The simplest kind of rule, with as few features as possible, is a binary rule. We trained the model on binary strings of 0's and 1's of length 16. Our rule was: the first digit of each vector should be '1'.

In all our trials, we separated sentences in the training sets by a newline character (one vector is considered a sentence), like the original BabyBERTa's training data.

We trained the model many times from scratch,245varying (a) the proportion of exceptions in the246

247 dataset, (b) the number of unique vectors 248 (sentences) in the dataset, and, (c) the number of 249 epochs of training.

#### 250 2.2.4 Evaluation Procedure 251

257

258

259

260

261

262

263

264

265

After a full training sequence was complete, we 252 tested the trained models on novel test sets, whose 253 format was inspired by the grammar test suites 254 used to evaluate BabyBERTa (Huebner et al. 255 2021). 256

Test vectors were generated in pairs. Each pair of vectors was identical, except that the first digit of the 16 digits of one of the vectors was '1'and in the other it was '0'. Each vector has its "surprisal" calculated. Surprisal is equivalent to the sum of the cross-entropy errors of each token in a given sequence. Since our sequences were only one token each, surprisal was just the cross-entropy error of that token.

266 If the model has learned a rule, then it should assign a lower surprisal score to a vector that 267 follows the rule than to a nearly identical vector 268 that breaks the rule. The model did a better job 269 predicting one sentence in each pair over the 270 other-the one with a lower surprisal score. We 271 say that the model *prefers* sentences with lower 272 surprisal scores. 273

The model's accuracy on each test set is 274 equivalent to how often, as a percentage, the 275 model prefers vectors that follow the rule, which 276 we compute by dividing the number of vector 277 pairs for which the model prefers the rule-278 following vector by the overall number of vector 279 280 pairs.

For each from-scratch model, we generated a new unique test set of 1,000 vector pairs.

**3** Trials



Figure 1: Model accuracies (represented by color of a point) for different training set sizes (x-axis), proportions of exceptions per training set (y-axis), and number of epochs (5, 10, & 20).



283

281

282



Figure 2: Effects of varying the proportion of exceptions (x-axis) on model accuracy (y-axis) for all combinations of e (# of epochs) and n (size of training dataset). Each graph contains 2 linear regressions: one on the left side of the TP threshold ( $\theta_n = n/ln(n)$ ) and one to the right.

Figures 1 and 2 show the results of training and
testing our models. If our models' learning were
governed by the Tolerance Principle, we would
expect:

288

289

290

291

292

293

294

(1) Learning should occur quantally. There should be a statistically significant jump, for each graph in Figure 2, from the regression on the left of the TP threshold to the regression on right of the TP threshold. The slope of each regression should be close to 0.

295(2) Varying the number of epochs should296have no significant effects on learning,297since token frequency is not significant in298determining whether the language model299can learn.

We observe some clear trends in the above300figures. For one, noticeable learning of a rule301appears to be possible from training sets of just a302few hundred vectors. The number of epochs in303training has a major effect on learning: increasing304the number of epochs leads to higher overall305accuracies.306

In general, we see no Tolerance-Principle-like 307 quantal effect. In Figure 2, the jump from one 308 regression to the other (at the TP threshold) was 309 only statistically significant in 2 instances out of 310 30: (e=5, n=50) and (e=20, n=450), no more than 311 we would expect by chance. In combinations of e 312 and *n* where learning occurred, we tend to see a 313 gradient decrease in model accuracy as the 314 315 proportion of exceptions in the training set316 increases.

## 4 Conclusion

317

318

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

For this machine learning architecture, datasets of 319 a few hundred examples are large enough for rule 320 learning to occur. Learning appears to follow a 321 gradient-as the proportion of exceptions is 322 increased, there is a linear, not quantal, decrease 323 in accuracy. Token frequency is significant in 324 determining whether this language model can 325 learn; training for more epochs over the same data 326 increases accuracy. The threshold predicted by 327 the Tolerance Principle seems to have no 328 significant bearing on the language model's 329 learning. 330

331 5 Limitations

This work does not reflect a broad study on many language models. It is limited in scope to the study of one model with fixed hyper-parameters.

## References

- Anderson, S. R. (1969). West Scandinavian vowel systems and the ordering of phonological rules. PhD thesis, MIT.
- Contreras Kallens, P., Kristensen-McLachlan, R. D., & Christiansen, M. H. (2023). Large language models demonstrate the potential of statistical learning in language. *Cognitive Science*, 47(3), e13256.
  - Devlin, J. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Huebner, P. A., Sulem, E., Cynthia, F., & Roth, D. (2021, November). BabyBERTa: Learning more grammar with small-scale child-directed language. In *Proceedings of the 25th conference on computational natural language learning* (pp. 624-646).

- Jawahar, G., Sagot, B., & Seddah, D. (2019). What does BERT learn about the structure of language?.
  In ACL 2019-57th Annual Meeting of the Association for Computational Linguistics.
- Kingma, D. P. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- 357 Kiparsky, P. (1973). Elsewhere in phonology. In358 Anderson, S. R. and Kiparsky, P., editors, A

festschrift for Morris Halle, pages 93-106. Holt,	359
Rinehart and Winston, New York.	360
Koulaguina, E., & Shi, R. (2013). Abstract rule	361
learning in 11-and 14-month-old infants. Journal of	362
psycholinguistic research, 42, 71-80.	363
Koulaguina, E., & Shi, R. (2019). Rule generalization	364
from inconsistent input in early infancy. Language	365
Acquisition, 26(4), 416-435.	366
Liu, Y. (2019). Roberta: A robustly optimized bert	367
pretraining approach. arXiv preprint	368
arXiv:190/.11092, 304.	260
Macwhinney, B. (2000). The CHILDES project: The	309
Shi P & Emond E (2022) The threshold of rule	370
productivity in infants. Exontiars in Psychology 14	371
1251124	372
Vaswani A (2017) Attention is all you need	373
Advances in Neural Information Processing	374
Systems002E	375
Warstadt, A., Mueller, A., Choshen, L., Wilcox, E.,	376
Zhuang, C., Ciro, J., & Cotterell, R. (2023).	377
Findings of the BabyLM Challenge: Sample-	378
efficient pretraining on developmentally plausible	379
corpora. In Proceedings of the BabyLM Challenge	380
at the 27th Conference on Computational Natural	381
Language Learning.	382
Yang, C. (2016). The Price of Linguistic Productivity.	383
Cambridge, MA: The MIT Press.	38/
Yang, C. (2018). A user's guide to the tolerance	385
principle. Unpublished work.	200
Yang, C. (2023). A User's defense of the tolerance	300
principle. University of Pennsylvania.	387
	388