Sparse Encoding of Grammatical Gender in LSTM Language Models

Priyanka Sukumaran (psukumaran23@gmail.com)

Conor Houghton (conor.houghton@bristol.ac.uk)

School of Engineering Mathematics and Technology, University of Bristol

Nina Kazanina (nina.kazanina@unige.ch)

Department of Basic Neurosciences, University of Geneva

46

76

77

78

Abstract

1

33

Neural network language models excel at capturing the 47 2 complexities of natural language, yet their internal rep-48 3 resentations remain poorly understood. A key question 49 4 is whether such models form structured, human-like ab- 50 5 stractions that support generalization. We investigate ⁵¹ 6 how LSTM language models encode grammatical gen-52 7 der-an ideal test case, as gender is lexically fixed and 53 8 generally not inferable from semantics. We focus on 54 9 long-distance dependencies and various gender agree- 55 10 ment configurations. 56 11

We conduct single-unit ablation to identify neurons 57 12 critical for grammatical gender agreement. Across eight 58 13 LSTM models, we find between one and five units whose 59 14 removal significantly disrupts performance-by over 60 15 40% in some constructions involving gender-interfering 61 16 nouns. These units are essential for both noun-adjective 62 17 and noun-past-participle gender agreement. Neuron ac-63 18 tivity analyses reveal that these units exhibit category-64 19 specific effects, with some showing a preference for de-65 20 fault gender forms, such as masculine nouns. 66 21

Our findings show that LSTMs develop sparse and 67 22 structured representations of grammatical gender, rem-68 23 iniscent of grandmother cells in neuroscience. These re- 69 24 sults suggest that abstract grammatical categories can 70 25 emerge naturally in LSTM training. More broadly, this 71 26 work contributes to our understanding of how language 72 27 models encode linguistic structure, with implications for 73 28 model interpretability and parallels between artificial and 74 29 biological computation. 75 30

 Keywords: neural network language models; linguistic generalization; mechanistic interpretability; encoding mechanisms

Introduction

LANGUAGE, a uniquely human cognitive function, is now re-79 34 markably imitated by neural network language models. Lan- 80 35 guage models perform well across a range of language com- 81 36 prehension and generation benchmarks, sometimes reaching 82 37 human-level performance. Yet, their errors remain revealing: 83 38 where they fail, we gain insight into the distinctiveness of hu-84 39 man linguistic competence, offering valuable comparisons for 85 40 cognitive science (Chowdhury & Zamparelli, 2018; Chaves, 86 41 2020; Lan et al., 2024). Moreover, understanding how lan- 87 42 guage models perform linguistic tasks offers cognitive science 88 43 new comparative tools for studying the neural basis of lan-89 44 guage. In particular, the implicit grammatical representations 90 45

in machine models may provide testable hypotheses about how grammar is represented neurally (Houghton et al., 2023).

This study focuses on LSTMs to investigate the mechanisms of grammatical processing using the targeted syntactic evaluation approach Linzen et al. (2016). LSTMs demonstrate strong grammatical competence and greater biological plausibility than more powerful transformer models, positioning them as effective minimal models for studying neural grammar representation (Linzen et al., 2016). We investigate the internal mechanisms that enable LSTMs to encode and represent grammatical gender to maintain agreement in complex contexts. We use single-neuron ablation: systematically removing or 'ablating' neurons from a neural network to understand their function. In addition, we ask whether LSTM language models encode the grammatical gender as an *abstract category*, thus enabling linguistic generalization.

This paper confirms the emergence of specialized single neurons and sparse networks in LSTMs for encoding grammatical gender, using a test set we developed in French. We investigate two gender agreement contexts: noun-adjective and noun-past-participle. We evaluate the role of genderspecific units in these contexts. For each, we test both adjacent and long-distance agreement in the presence of attractors. Expanding on (Lakretz et al., 2019), which focused on constructions with minimal intervening words, we evaluate individual neurons' roles across varying dependency lengths. We assess whether LSTMs generalize grammatical gender across four head noun categories: singular/plural and feminine/masculine. Our findings provide evidence of sparse and selective representations in LSTMs, highlighting parallels with biological systems such as the grandmother cell phenomenon in neuroscience (Gross, 2002; Quiroga et al., 2005).

Methods

We trained a two-layer LSTM with 650 hidden units on French Wikipedia data (Mueller et al., 2020) for next-word prediction. We trained both tied and untied versions of the model, with the tied variant sharing input, output, and embedding weights (Press & Wolf, 2017) to facilitate interpretability. For robustness, we trained five untied and three tied models with different random initialization seeds.

The LSTM models were tested on grammatical agreement, which involves coordinating the noun with other elements such as the verb, determiner or adjective based on the noun's properties such as nominal number (singular/plural) and gender (masculine/feminine/neutral). Traditional psycholinguistics studies have used agreement tasks to probe hierarchicah26
syntactic knowledge that humans employ to parse language127

93 (Bock & Miller, 1991; Franck et al., 2002).

We evaluate LSTM performance before and after single-129 94 neuron ablation for two common contexts of gender agree-130 95 ment in French: noun-adjective (NA) and noun-participle (NP)131 96 agreement, see Table 2. The simplest case is (A) adjacent₁₃₂ 97 agreement, where no intervening words separate the nounisi 98 and its agreement target, making agreement straightforward,134 99 For instance, in the NA construction in Table 2, las, f robes, f estiss 100 bleues, f. To systematically test agreement, we varied the gen-136 101 der and number of the subject noun. Secondly, we also testi37 102 (B) long-distance agreement by including 1-11 words between 138 103 the noun and its agreement target using prepositional phrases 139 104 and subject-relative clauses. These constructions are labelled 140 105 NA-n and NP-n, where n is the number of intervening words as141 106 an example for NA-5: la robes, f [que j' aime beaucoup]rel esti42 107 bleue_{s f}/bleu_{m f}. Next, we test the more complex (**C**) agree-143 108 ment across an attractor by introducing another noun, with144 109 varying number and gender, using a prepositional phrase that 145 110 could potentially interfere with gender agreement. We test this146 111 using conditions NN_{num.gen}A and NN_{num.gen}P, see Table 2. 147 112

Table 1: Gender/syntax units in different models with a ran- $_{149}$ dom seed, labeled as M_{seed} . This result section focuses₁₅₀ on M_{528} . Each unit reduces gender agreement performance₁₅₁ when: SG: the head noun is singular, PG: the head noun is₁₅₂ plural, G: the head noun is either singular or plural, S: there is₁₅₃ an interfering attractor or in longer dependencies.

	Gender/syntax units			
Untied Model	SG	PG	G	S
M ₅₂₈	-	870	1098	958
M ₁₅₇₁	914	1013	-	-
M ₇₀₄	863	1012	-	-
M ₂₂₀	-	-	937	-
M ₇₃	-	-	1269	-
Tied Model				
M_tied ₅₂₈	-	-	1174	-
M_tied ₇₀₄	-	1262	930	
M_tied ₇₃	764	768	-	-

113

Results & Discussion

169

128

148

155 156

157

158

159

160

161

162

163

164

165

166

167

168

Across eight LSTM initialisations, ablating just 2-3 out of 1300₁₇₀ 114 units resulted in significant drops in grammatical agreement₁₇₁ 115 performance. We focus on model M₅₂₈, with similar patterns₁₇₂ 116 observed across other runs (Table 1). Three units consis-173 117 tently emerged as critical: Unit 870 (PG-870), which selec-174 118 tively encoded plural gender agreement; Unit 1098 (G-1098),175 119 encoding gender more broadly; and Unit 958 (S-958), which₁₇₆ 120 tracked syntactic structure. Performance on long-distance de-177 121 pendencies-especially those with intervening tokens or at-178 122 tractors-dropped to near chance after ablation. 123 179

The single-unit gate and cell state activity analyses con-180 firmed that these units encoded abstract grammatical fea-

tures. Gender units activated at the head noun and maintained gender information through to the agreement target, often via sustained cell state values. Syntax units preserved structural information across the clause, suggesting internal representations of sentence depth. t-SNE projections of unit activations revealed clear separability by gender and number, further supporting the presence of abstract, categorysensitive encodings (Figures 7 and 8).

Interestingly, the effect of ablation was asymmetric: agreement with feminine nouns was far more disrupted than masculine, aligning with the default reasoning strategy proposed by (Jumelet et al., 2019), whereby default grammatical features—such as singular number or masculine gender—are encoded implicitly in model weights, while non-default features like plural or feminine depend more on explicit, input-driven encoding. Our results support this distinction: ablations of gender-selective units impaired feminine agreement but left masculine largely intact, indicating that masculine may be redundantly or diffusely encoded.

These findings mirror a broader distinction in neural coding between localist and distributed representations. The emergence of highly selective units in our LSTM aligns with the grandmother cell hypothesis in neuroscience—the idea that individual neurons (or units) can become tuned to specific, abstract categories. Notable examples include the "Jennifer Aniston neuron" in the human hippocampus, which responds exclusively to images of that individual (Quiroga et al., 2005). Similarly, Konorski (1967) proposed gnostic units for recognizing object categories. Our gender and syntax units exhibit similar behavior in a computational setting, selectively activating for abstract grammatical roles and sustaining them across syntactic contexts.

However, localist units do not exclude distributed coding. Sparse, localist encodings (as in our identified units) can coexist with more distributed representations—particularly for default categories. This balance reflects ongoing debates in neuroscience: while localist encoding offers interpretability, distributed representations are thought to be more robust and biologically plausible (Bowers, 2017; Rolls, 2017). Importantly, distributed representations can also be sparse—involving only a few active units, each of which may not encode easily interpretable features.

Thus, LSTMs, while simplistic compared to the brain, offer a computational testbed for examining these dual encoding strategies. These findings open several avenues for future work. Do the identified units also encode related features like animacy, case, or number? Are they robust across languages with richer agreement systems (e.g., German, Bantu)? Can similar sparse mechanisms be observed in transformer models, or are they unique to sequential architectures like LSTMs? Most importantly, ablation studies like ours help demarcate the boundary between what linguistic phenomena can be explained by statistical learning in neural networks, and what might be uniquely human. As noted by Surendra et al. (2023), this boundary is often subtle.

References

234

- Bock, K., & Miller, C. A. (1991). Broken agreement. Cognitive235 182 Psychology, 23(1). doi: 10.1016/0010-0285(91)90003-7 236 183
- Bowers, J. S. (2017). Grandmother cells and localist repre-237 184

sentations: a review of current thinking (Vol. 32) (No. 3). doi.238 185 10.1080/23273798.2016.1267782 186

- Chaves, R. P. (2020). What don't RNN language models learn²⁴⁰ 187
- about filler-gap dependencies? . In Proceedings of the so-241 188 ciety for computation in linguistics. 242

189

- Chowdhury, S. A., & Zamparelli, R. (2018). RNN simula-243 190 tions of grammaticality judgments on long-distance depen-244 191
- dencies. In Coling 2018 27th international conference orf245 192 computational linguistics, proceedings. 193
- Franck, J., Vigliocco, G., & Nicol, J. (2002). Subject-verb 194 agreement errors in French and English: The role of syn-195 tactic hierarchy. Language and Cognitive Processes, 17(4). 196
- doi: 10.1080/01690960143000254 197

181

210

- Gross, C. G. (2002). Genealogy of the "grandmother cell" 198 (Vol. 8) (No. 5). doi: 10.1177/107385802237175 199
- Houghton, C., Kazanina, N., & Sukumaran, P. (2023,200 Beyond the limitations of any imaginable mecha-12). 201 nism: Large language models and psycholinguistics. Be-202 havioral and Brain Sciences, 46, e395. doi: 10.1017/ 203 S0140525X23001693 204
- Jumelet, J., Zuidema, W., & Hupkes, D. (2019). Analysing 205 neural language models: Contextual decomposition reveals 206 default reasoning in number and gender assignment. In 207 Conll 2019 - 23rd conference on computational natural 208 language learning, proceedings of the conference. doi: 209 10.18653/v1/k19-1001
- Konorski, J. (1967). Integrative activity of the brain. Chicago: 211 University of Chicago Press. 212
- Lakretz, Y., Kruszewski, G., Desbordes, T., Hupkes, D., De-213
- haene, S., & Baroni, M. (2019). The emergence of number 214 and syntax units in LSTM language models. In Naacl hlt
- 215 2019 - 2019 conference of the north american chapter of the
- 216 association for computational linguistics: Human language
- 217 technologies - proceedings of the conference (Vol. 1). 218
- Lan, N., Chemla, E., & Katzir, R. (2024). Large language 219 models and the argument from the poverty of the stimulus. 220 Linguistic Inquiry. 221
- Linzen, T., Dupoux, E., & Goldberg, Y. (2016). Assessing the 222 Ability of LSTMs to Learn Syntax-Sensitive Dependencies. 223
- Transactions of the Association for Computational Linguis-224 *tics*, 4. doi: 10.1162/tacl{_}a{_}00115 225
- Mueller, A., Nicolai, G., Petrou-Zeniou, P., Talmina, N., & 226 Linzen, T. (2020). Cross-Linguistic Syntactic Evaluation of 227
- Word Prediction Models. In Proceedings of the 58th annual 228
- meeting of the association for computational linguistics (pp. 229 5523-5539). doi: 10.18653/v1/2020.acl-main.490 230
- Press, O., & Wolf, L. (2017). Using the output embedding to 231 improve language models. In 15th conference of the euro-232
- pean chapter of the association for computational linguis-233

tics, eacl 2017 - proceedings of conference (Vol. 2). doi: 10.18653/v1/e17-2025

- Quiroga, R. Q., Reddy, L., Kreiman, G., Koch, C., & Fried, I. (2005). Invariant visual representation by single neurons in the human brain. Nature, 435(7045). doi: 10.1038/ nature03687
- Rolls, E. T. (2017). Cortical coding (Vol. 32) (No. 3). doi: 10.1080/23273798.2016.1203443
- Surendra, K., Schilling, A., Stoewer, P., Maier, A., & Krauss, P. (2023). Word class representations spontaneously emerge in a deep neural network trained on next word prediction. In 2023 international conference on machine learning and applications (icmla) (pp. 1481-1486).

Table 2: Example phrases for adjacent noun-adjective and noun-past-participle gender agreement. We tested both singular and plural noun phrases for each condition. s.f: singular feminine, s.m: singular masculine, p.f: plural feminine, p.m: plural masculine.

	Singular	Plural		
Noun-adjective				
NA	la robe est <u>bleue</u> /bleu	les robes sont <u>bleues</u> /bleus		
(No attractor)	the dress _{s.f} is <u>blue_{s.f}/blue_{s.m}</u>	the dresses_{p.f} are <u>blue_{p.f}/blue_{p.m}</u>		
NN _{gen} A	la robe avec le <i>sac</i> est <u>bleue</u> /bleu	les robes avec les <i>sacs</i> sont <u>bleues</u> /bleus		
(Gender attractor)	the $dress_{s.f}$ with the $bag_{s.m}$ is $\underline{blue}_{s.f}$ /blue $_{s.m}$	the $\textit{dresses}_{p.f}$ with the $\textit{bags}_{p.m}$ are $\underline{\textit{blue}}_{p.f}/\textit{blue}_{p.m}$		
NN _{num.gen} A	la robe avec les <i>sacs</i> est <u>bleue</u> /bleu	les robes avec le sac sont <u>bleues</u> /bleus		
(Number/gender attractor)	the $\textbf{dress}_{s.f}$ with the $\textit{bags}_{p.m}$ is $\underline{\textit{blue}_{s.f}}/\textit{blue}_{s.m}$	the $\textbf{dresses}_{p.f}$ with the $\textit{bag}_{s.m}$ are $\underline{\textit{blue}}_{p.f}/\textit{blue}_{p.m}$		
Noun-participle				
NP	la robe est <u>tombée</u> /tombé	les robes sont tombées/tombés		
(No attractor)	the dress _{s.f} <u>fell_{s.f}/fell_{s.m}</u>	the dresses _{p.f} <u>fell_{p.f}/fell_{p.m}</u>		
NN _{gen} P	la robe avec le <i>sac</i> est <u>tombée</u> /tombé	les robes avec les <i>sacs</i> sont <u>tombées</u> /tombés		
(Gender attractor)	the $dress_{s.f}$ with the $bag_{s.m}$ $\underline{fell}_{s.f}$ /fell _{s.m}	the $\textsc{dresses}_{p.f}$ with the $\textit{bags}_{p.m}$ $\underline{\text{fell}}_{p.f}/\text{fell}_{p.m}$		
NN _{num.gen} P	la robe avec les <i>sacs</i> est <u>tombée</u> /tombé	les robes avec le <i>sac</i> sont <u>tombées</u> /tombés		
(Number/gender attractor)	the $\text{dress}_{s.f}$ with the $\textit{bags}_{p.m}$ $\underline{fell}_{s.f}/fell_{s.m}$	the $dresses_{p.f}$ with the $bag_{s.m} \underline{fell}_{p.f}/fell_{p.m}$		



Figure 1: Performance on simple agreement with 0, 1, 5 and 10 intervening tokens for noun-adjective (NA) and noun-participle (NP) conditions after ablation of each LSTM unit (*x*-axis). Each panel shows mean agreement accuracy (*y*-axis) after ablating individual units (*x*-axis) for singular/plural and masculine/feminine nouns. Performance is split by head noun category: singular (top), plural (bottom), masculine (left) or feminine (right). Performance is further broken down by constructions with zero (green), one (red),five (yellow), and ten (blue) intervening tokens. Each dot represents performance after ablating a unit, with significant drops (*z*-score < -3) highlighted.



Figure 2: Input gate, forget gate, cell state and hidden state activity for gender units for simple agreement, that is, no attractor.



Figure 3: Input gate, forget gate, cell and hidden state activity for agreement across a gender attractor (NNgenA/ NNgenP)



Figure 4: Performance after ablation of each LSTM unit (*x*-axis) on noun-adjective agreement with attractor nouns of varying gender. Each panel corresponds to test sentences with different combinations of subject and attractor noun categories. For example, the SM_SF indicates Singular-Masculine subject noun and Singular-Feminine attractor.



Figure 5: Performance after ablation of each LSTM unit (*x*-axis) on noun-adjective agreement with attractor nouns of varying gender and number.



Figure 6: Input gate, forget gate, cell and hidden state activity for syntax unit S-958 across conditions



Figure 7: t-SNE projections of LSTM cell activations for selected gender units across gender agreement conditions (FS = feminine singular, FP = feminine plural, MS = masculine singular, MP = masculine plural). Activations across six syntactic configurations (simple agreement, agreement across gender and number attractors), both for noun-adjective and noun-participle across short and longer-range dependencies, were aggregated and flattened across time steps to capture the full temporal profile. Gender **Unit-1098** (left) and Gender **Unit-870** (right) show clear clustering by grammatical gender, suggesting category-specific encoding across sentence constructions. In contrast, syntax unit in Figure 8 shows less separability.



Figure 8: t-SNE projections of LSTM cell activations for the syntax-related **Unit-958** (left) and a randomly selected unit **Unit-624** (right). Compared to the gender units (Figure 7), these units show less separable patterns across gender-number categories, highlighting the contrast between specialized and non-specialized representations.