

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047

# The emergence of number and syntax units in LSTM language models

Anonymous NAACL submission

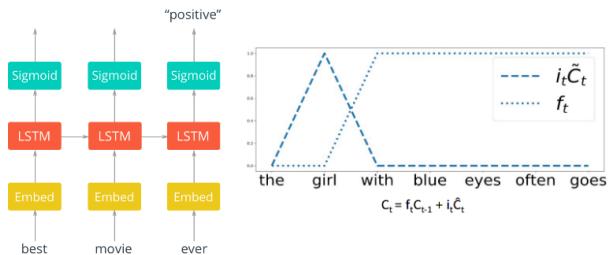
## Abstract

Recent work has shown that LSTMs trained on a generic language modeling objective capture syntax-sensitive generalizations such as long-distance number agreement. We have however no mechanistic understanding of how they accomplish this remarkable feat. Some have conjectured it depends on heuristics that do not truly take hierarchical structure into account. We present here a detailed study of the inner mechanics of number tracking in LSTMs at the single neuron level. We discover that number information is managed by very few “grandmother cells” in a localist fashion. Importantly, the behaviour of the number cells is partially controlled by other units that are independently shown to track the syntactic structure of sentences. We conclude that LSTMs are, to some extent, implementing genuinely syntactic processing mechanisms, paving the way to a more general understanding of grammatical encoding in LSTMs.

## 1 Introduction

[1-6]

Studies showing that LSTMs trained on language modeling do well on the agreement task: (?; ?), to a lesser extent: (?; ?). Studies conjecturing this is just heuristics: (?; ?). Grandma cells: (?).



**Figure 1:** Caption.

## 2 Related literature

### 2.1 Interpreting LSTM networks

Short survey of works that: - Interpret neural networks in general (e.g., CNNs in vision)  
- Interpret LSTMs (e.g. Karpathy’s) - Interpret LSTM language models

Besides ?), remember to discuss ?), as they find a “sentiment” grandma cell. ?) found a character-level RNN to track morpheme boundaries in a localist fashion, through a single cell.

Radford:etal:2017,

### 2.2 Subject-verb agreement in English

- Subject-verb agreement from psycholinguistics (e.g., Miller and Bock, Franck and Rizzi) - SV agreemnt in LSTMs (?; ?; ?; ?; ?) - Relate to Nelson et. al 2017 PNAS, an intracranial study that identifies electrodes whose high-gamma activity correlates with syntactic tree-depth (numnber of open nodes)

048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095

096	Simple	the boy greets the guy
097	Adv	the boy probably greets the guy
098	2Adv	the boy most probably greets the guy
099	CoAdv	the boy openly and deliberately greets the guy
100	NamePP	the boy near Pat greets the guy
101	NounPP	the boy near the car greets the guy
102	SubjRel	the boy that avoids the girl greets the guy
103	ObjRel	the boy that the girl avoids greets the guy
104	ObjRel0	the boy the girl avoids greets the guy

Table 1: NA tasks illustrated by representative singular sentences.

### 3 Data

#### 3.1 Number-agreement tasks

We generated number-agreement tasks (NA-tasks) with fixed syntactic structures and varied lexical material that probe subject-verb number agreement in increasingly challenging set-ups. The different structures are illustrated in Table 1 by examples where all forms are in the singular. Distinct sentences were randomly generated by selecting (different) words from pools of 20 subject/object nouns, 15 verbs, 10 adverbs, 5 prepositions, 10 proper nouns and 10 location nouns. The items were selected so that their combination would not lead to semantic anomalies. For each NA-task, we generated singular and plural versions of each sentence. We refer to each such version as a *condition*. For NA-tasks that have other nouns occurring between subject and main verb, we also systematically varied their number, resulting in four conditions. For example, the NounPP sentence in the table illustrates the SS (singular-singular) condition. The corresponding sentences in the other conditions are: “the boy near the cars greets the guy” (SP), “the boys near the car greet the guy” (PS), and “the boys near the cars greet the guy” (PP).

Possibly move this to Models or Results section. For each NA-task, model performance was evaluated for each of its conditions separately, by computing its accuracy in predicting the correct form of the main verb. Specifically, for each sentence from a condition, the model was presented with all words up to the one immediately preceding the main verb. If the model assigned higher likelihood to the correct form of the verb than to its counterpart with the wrong number,

then the trial was marked as a hit.

Yair: consider move these explanations about contrasts to the results section: Note that 2Adv features the same subject-verb distance as NamePP, but without gender-carrying words acting as possible distractors in the middle. Similarly, CoAdv can serve as a distractor-free control for NounPP.

Task data-set sizes across conditions ranged from 600 (Simple) to 18k sentences (relatives), based on the possibilities for variation allowed by the combinatorics given each structure.<sup>1</sup>

Finally, we also used the naturalistic, corpus-derived agreement test set of (?), in the version made available by (?).

#### 3.2 Number of open nodes

### 4 The models

The recurrent models that are the focus of our studies are LSTM models (?). To inspect the internal dynamics of these models, we use additional linear models trained to predict information from the hidden state activations of the language models. (DH: Or something like this. I'd put (?) and (?), but @yair perhaps we have also some neuroscience ref? Would you mind calling it diagnostic classifiers?) In this section, we describe how we train the language models we study, as well as our diagnostic classifier protocol.

#### 4.1 LSTM language model

Following (?), we study LSTM language models with two layers containing 650 nodes and an embedding size of 650. We consider the pretrained model that was made available by (?), as well as 3 (?) additional models that were trained following the same protocol, but with a different seed and/or different dropout rate.

This leaves us with 4 (?) models that ... (DH: One sentence describing the number of modes that we used and how they differed. Here also name the models. I believe this can be very short, for details we can refer to the code/appendix, we shouldn't waste our space

<sup>1</sup>We uploaded the data-set generation script as supplementary material.

192 here with too much information unimportant  
 193 to the story. E.g., LSTM equations, the ex-  
 194 act training data...these details can be found  
 195 in earlier papers.)

196 The models we consider were comparable in  
 197 terms of perplexity on a heldout dataset (with  
 198 perplexities of W, X, Y and Z, respectively) but  
 199 differ in their accuracy on the *agreement task*  
 200 ... (DH: I have to find these numbers, then I'll  
 201 finish this subsection.)

## 202 4.2 Regression model

203 - Explain that the model was used to predict the  
 204 depth of the syntactic tree from network activity

### 205 4.2.1 Model description

- 206 - Describe the Features: network activi-  
 207 ty (hidden/cell activity)
- Describe the la-  
 208 bel: Tree depth (refer to the section describ-  
 209 ing the synthetic data)
- Describe the model: A  
 210 Ridge/LASSO model.

### 211 4.2.2 Model training and evalution

- Model training: nested 5-fold CV procedure.  
 Train/val/test. Optimal regularization size was  
 estimated from the validation set (report optimal  
 lamda - figures in FAIRNS.pdf on slack)
- model evaluation: R-squared on test set. Report  
 resulting values (text+figures in FAIRNS.pdf).

## 212 5 Results

213 To successfully perform the NA-task, it seems  
 214 that the LSTM-LM should encode and store at  
 215 least two types of informations: (1) the gram-  
 216 matical number of the subject; and (2) the syn-  
 217 tactic structure of the sentence. The latter in-  
 218 formation is important for knowing at which time  
 219 point the network should output and update the  
 220 stored number information. This section de-  
 221 scribes the detailed mechanism by which units  
 222 in the network encode and store these types of  
 223 information, and the interaction between these  
 224 units. Section 5.1 identifies and describes 'num-  
 225 ber units', units in the network that specializes  
 226 in encoding and storing the grammatical num-  
 227 ber of the subject for short- and long-range de-  
 228 pendencies; Section 5.2 describes 'syntax units',  
 229 units whose activity correlates with the depth of

240 the syntactic tree; Then, section 5.3 explores the  
 241 interactions among syntax and number units, re-  
 242 quired for accomplishing the NA-task. Finally,  
 243 section 5.4 looks into the dynamics of these  
 244 units during the processing of complex syntactic  
 245 structures with single or double relative clauses.

### 246 5.1 Long-range number-units

247 To accomplish the NA-task, the LSTM network  
 248 should encode and store the grammatical num-  
 249 ber of the subject up to one step before the  
 250 verb, when prediction of the verb form (sin-  
 251 gular or plural) occurs. In some cases, this  
 252 may be quite challenging, in particular in the  
 253 case of a long-range dependency between sub-  
 254 ject and verb, and in particular, when another  
 255 noun with an opposite number appears before  
 256 the verb (cite). This section explores the un-  
 257 derlying mechanism that enables the network to  
 258 encode and store number information in various  
 259 syntactic constructions, including ones with an  
 260 interfering noun. Subsection 5.1.1 describes an  
 261 ablation study, which reveals *long-range num-  
 262 ber units (LR-number units)* that can store and  
 263 carry number information from subject to verb,  
 264 also across interfering nouns. Subsection 5.1.2  
 265 describes the intricate gate and state dynam-  
 266 ics of LR-number units during the processing  
 267 of sentences with long-range dependencies. Sec-  
 268 tion 5.1.3 describes a different type of number  
 269 units that encode grammatical number for only  
 270 short-range dependencies, such as in the case  
 271 of a verb immediately following a subject. Fi-  
 272 nally, subsection 5.1.4 looks into the structure of  
 273 the efferent weights of number units, by which  
 274 number information is propagated to the output  
 275 layer.

#### 276 5.1.1 Local vs. distributed code - an 277 ablation study

278 Generally, number information may be stored  
 279 in the network in either a local, sparse, or a dis-  
 280 tributed way, depending on the fraction of ac-  
 281 tive units that carry number information. We  
 282 hypothesized that if the network uses a local  
 283 or sparse coding, meaning that there's a small  
 284 set of units that encode number information,  
 285 then ablating these units would lead to a drastic

decrease in performance on the NA-task, compared to when ablating other units. To test this, we conducted ablation experiments in which each time a single unit of the network is ablated and the resulting model is then evaluated on several NA-tasks. Each NA-task contained sentences with a fixed syntactic structure, such as "Det Noun Adv Verb" or "Det Noun P Det Noun Verb", and each task was composed of several conditions depending on the possible assignments of grammatical number to the noun(s) in the sentence (see section 3.1 for details about all NA-tasks). In addition, we also evaluated each ablated model on the Linzen task (cite). Tables 1 summarizes the results from all ablation experiments, showing units whose ablation resulted in a performance decrease of more than 10% (TODO: choose a non-arbitrary threshold by looking at the distribution). For each NA-task, the performance of the full, non-ablated, model is also reported.

The discussion of Table 2 should lead more clearly to the following discussion in which we focus on 776 and 988. We should say at a certain point that the 770 and 1283 might be complementary cells, but clearly the main job of tracking LR number is left to 776 and 988.

It's still unclear to me that the relative clause results really fit here.

We first highlight several aspects of the behavioral results of the full network (table 1 - right column) before describing in more details the ablation results. First, some NA-tasks and conditions are clearly more difficult for the network than others. For example, performance on the simple NA-task is better than that on the nounPP NA-task, which in turn is better than that of the objrel task. This matches previously reported results in humans and LSTM-LMs (cite). Second, having an interfering noun before the verb, with an opposite number than that of the subject, is clearly a more challenging task for the network - we find that for the nounPP, subjrel and objrel tasks:  $ACC_{SP} > ACC_{SS}$  and  $ACC_{PS} > ACC_{PP}$ . We return

NA task	C	770	776	988	1283	Full	
Simple	S	-	-	-	-	100	337
Adv	S	-	-	-	-	100	338
2Adv	S	-	-	-	-	99.8	339
Co-Adv	S	-	-	84.0	84.0	98.8	340
namePP	SS	-	-	-	-	98.9	341
nounPP	SS	-	-	-	-	97.5	342
nounPP	SP	-	-	58.8	-	88.5	343
subjrel	SS	-	-	88.0	-	97.0	344
subjrel	SP	-	-	-	-	58.8	345
objrel	SS	-	-	-	-	64.7	346
objrel	SP	-	-	-	-	45.7	347
Simple	P	-	-	-	-	100	348
Adv	P	-	-	-	-	99.6	349
2Adv	P	-	-	-	-	99	350
Co-Adv	P	-	78.9	-	-	99.7	351
namePP	PS	-	57.6	-	-	66.8	352
nounPP	PS	85.2	49.7	-	-	93.2	353
nounPP	PP	-	81.7	-	-	98.3	354
subjrel	PS	85.8	58.6	-	-	87.8	355
subjrel	PP	-	88.1	-	-	99.3	356
objrel	PS	-	-	-	-	69.0	357
objrel	PP	-	-	-	-	81.0	358
Linzen	-	?	?	?	?	?	359

Table 2: Ablation experiments results: Percentage of correct subject-verb agreements in all NA-tasks (section 3.1). Full - non-ablated model, C - condition, S - singular, P - plural. For task with two nouns, SS - singular-singular, SP - singular-plural, PS - plural-singular, PP - plural-plural. Red: singular number units, Blue: Plural number units.

to this point in section 5.4. Finally, for long-range dependencies, reliably encoding that the subject is singular is in most cases more difficult than plural. For example, in all the above tasks:  $ACC_{SS} < ACC_{PP}$  and  $ACC_{SP} < ACC_{PS}$ . Interestingly, this singular-plural asymmetry has been reported also in humans (cite). We elaborate on this point in the discussion section.

We next describe several important aspects of the ablation-experiment results. First, in all NA-tasks, only four units from the entire network (1300 LSTM units in total) had a significant effect on task performance. This result suggests a local coding scheme for long-range grammatical-number information

(TODO: quantify a 'significant' reduction, and perhaps link to the emergence of local coding in neural-network simulations (Bowers) and to findings about grandmother neurons in humans (e.g., Fried)). Second, we note that all number units emerged at the second layer of the network. This seems appropriate if number information needs to be directly projected to the output layer for correct verb-form prediction. In section 5.1.4 we further explore the projection weights from number units. Third, for simple, 1Adv and 2Adv NA-tasks, none of the units had a significant effect on task performance. This suggests that for short-range dependencies number information may be also encoded elsewhere in the network, perhaps via a more distributed code. We therefore make a distinction between long-range (LR) and short-range (SR) number units in what follows. We return to this point in section 5.1.3 (TODO: complete the identification of short-range number units from the resulting weights of the classifier in the generalization-across-time experiment). Fourth, LR-number units can be further divided into two types, depending on the grammatical number of the subject. Units 770 and 776 had a significant effect only when the first noun was plural, but not singular, and vice versa for units 988 and 1283 (blue and red in table 1, respectively). We therefore refer to the former as *plural units* and to the latter as *singular units*. Finally, we note that two of the number units (776 & 988) had an exceptional effect on network performance in both nounPP-SP&PS conditions. These two conditions are in particular revealing since they involve both a long-range dependency (over a prepositional phrase) and an interfering noun before the verb, while performance of the non-ablated network is still relatively high (88.5%&93.2%, respectively) in contrast to these conditions in subjrel and objrel. Ablating one of these two units brought the network from high performance on the NA-task to around chance-level performance (58.8%&49.7%, respectively). In the next section, we therefore focus on these two units when exemplifying gate and state dynamics of number units.

### 5.1.2 Visualizing gate and cell-state dynamics

Results from the ablation study suggest that there's a small set of units that encodes number information for long-range dependencies, in particular, we find that in some conditions two units can bring the network from relatively high performance to around chance-level performance on the NA-task (section 5.1.1). However, it remains unclear what is the exact mechanism underlying successful trials in the NA-task, and what goes wrong in unsuccessful ones. To better understand this, we now look into gate and state dynamics of these units during the processing of sentences from the nounPP NA-task. **Why do we focus on nounPP: because it is the minimal case with the four conditions, but without the complications of tracking another agreement relation as in the relatives.**

To anticipate the results and facilitate their interpretations, we begin by discussing what could be a solution to the NA-task implemented by number units in their gate and state activity. We recall that the update rule of the LSTM cell has two terms (equation 1.x). In the first term  $f_t * C_{t-1}$ , the forget gate controls whether to keep the previous content  $C_{t-1}$  stored in the cell ( $f_t = 1$  - perfect remembering), or forget it ( $f_t = 0$  - complete forgetting). In the second term  $i_t * \tilde{C}_t$ , the input gate controls whether the information currently presented to the network could be updated onto the cell state:  $i_t = 1$  - full access,  $i_t = 0$  - no access. Therefore, to produce correct number agreement, it seems that number units should at least have the following three properties: (1) The grammatical number of the subject should first be encoded by  $\tilde{C}_{t_{subject}}$ , encoding singular and plural with different values. (2) To grant the encoded grammatical number  $\tilde{C}_{t_{subject}}$  access to the cell, the input gate should be open at the time when the subject is presented:  $i_{t_{subject}} > 0$ , and ideally  $i_{t_{subject}} = 1$ . In addition, to protect the stored grammatical number from interfering information updating onto the cell, such as in the case of an interfering noun, the input gate should be

closed during all successive time steps until the verb:  $i_t = 0, t < t_{verb}$ ; (3) Finally, to successfully store number information in the cell for a long-range dependency, the forget gate should be in a remembering state, starting one time step after the subject:  $f_t = 1, t > t_{subject}$ . In addition, to clean up the cell from previously stored information, the forget gate should reset when the subject is presented:  $f_{t_{subject}} = 0$ . Figure 1B summarizes these three presumably desired properties.

Figure 2 presents the actual gate and state dynamics of units 776 and 988 during the processing of sentences from the nounPP NA-task. For each unit, we draw the dynamics of the suggestion  $\tilde{C}_t$  (panels A-B), input-gate (panels C-D), forget-gate (panels E-F) and the cell variable (G-H). For each of these cases, the four condition (SS, SP, PS and PP) are described in separate curves. Error-bars represent standard deviation across 1000 sentences in each condition.

We describe the results along the order of the properties discussed above. First, the values of the cell suggestions  $\tilde{C}_t$  of both units seem to obey the first property. For the singular unit 988, we find that singular nouns are encoded with negative values  $\tilde{C}_{t_{subject}} = -1$ , and plurals with positive  $\tilde{C}_{t_{subject}} = 1$  (panel A), and similarly for unit 776 (panel B). This shows that singular and plural nouns are indeed encoded differently by these units, in accordance with the results of the ablation study that suggested the labeling of units 988 and 776 as singular and plural units, respectively.

Second, input-gate dynamics of both number units seem to correspond to the second property described above. Input-gate activity spikes around the subject and stays approximately zero for subsequent time steps until the verb. One difference with respect to the desired property is the non-zero activity of the input gate at the time step immediately following the subject. This may be due to various reasons and requires further research. One possible explanation for this is that the network has developed this behavior as a heuristic to deal with compound nouns, given that for compound nouns the rel-

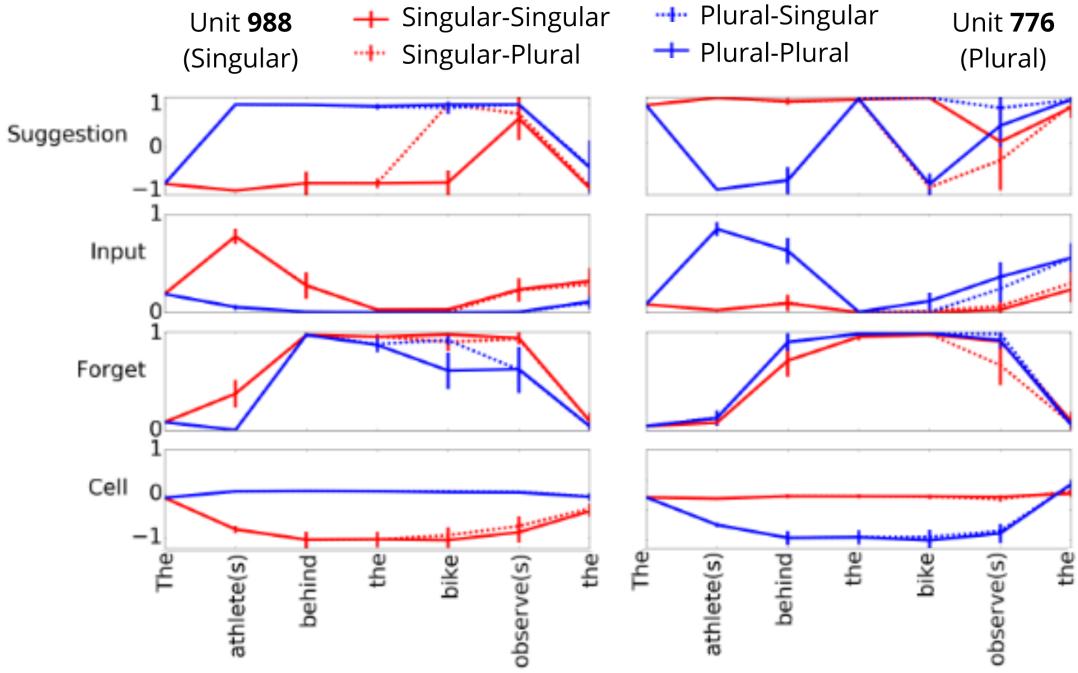
event number information resides at the second noun, whereas in the case of simple nouns there's anyway no risk of encountering an interfering noun immediately after the subject (TODO: discuss this part in the meeting to see if it makes sense to all. If yes, perhaps we could easily check this in an experiment.). Finally, note that for unit 988, the input gate is only open when the subject is singular, whereas for unit 776 it is only open when the subject is plural. This too is consistent with the labeling of these units as singular and plural.

Third, forget-gate dynamics of both number units also seem to correspond to the above properties. In both units, forget-gate activity starts at value around zero  $f_{t_{subject}} = 0$  and then goes abruptly towards its maximal value at the next step  $t_{subject} + 1$ , then stably staying at this level until after the verb  $t_{verb} + 1$ . Note that for all four conditions (SS, SP, PS and PP), the forget-gate exhibits similar dynamics, being indifferent to the grammatical number of the subject. This seems appropriate for whether the second noun is singular or plural given that the network cannot know in advance whether an interfering noun will appear, and it should anyway store number information for long-range dependencies also in the absence of any upcoming noun (TODO: explain or leave as an open question the reason for which we observe the same dynamics whether the first noun is singular or plural). Last, we note that in all cases the forget-gate activity resets at  $t_{verb} + 1$ . This seems appropriate, given that at this point the subject's number is no more useful, and the cell would be better free up to encode new number information.

Finally, cell activity should reflect the dynamics of the suggestion, input and forget gates. Indeed, the cell value becomes non-zero at  $t_{subject}$  and preserves this value until  $t_{verb} - 1$  when verb-form prediction occurs (Panels G-H). Note that this is the case only for the relevant conditions: in conditions SS and SP, unit 987 encodes singular as  $C_t = -1$  and is approximately zero during sentence processing in the other two conditions (PP and PS). Similarly, unit 776 encodes plural with a non-zero, negative, value only in the relevant conditions (PP and PS) but not in the

480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527

528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575

576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671

**Figure 2:** Cell and gate activations during processing of a sentence with a prepositional phrase between subject and verb. (A) Cell activity  $C_t$  for the two number units 775 and 987 and output activity  $h_t$  for the syntax unit 1149, for all four combinations of grammatical numbers of the two nouns. Note that the cell activity of units 775/987 is non-zero only when the first noun is plural/singular, respectively. (B) Corresponding forget-gate activity for the same number units. Note that gate activity is indifferent of the grammatical number of both nouns and that its value is close to one during the PP until after the verb. (C) Input-gate activity of the same units. Note that the gate value of unit 775/987 spikes around the first noun only when it is plural/singular. **Change bike to bike(s) in caption?** Also, I think the NLP audience will understand  $\tilde{C}$  and  $C$  better than Suggestion and Cell. Also, do we really need  $\tilde{C}$ ? Isn't the important information already there in the  $C$  plot?

irrelevant ones (SS and SP). Note that for the irrelevant conditions, cell activity is kept approximately zero thanks to the clear-up of the cell:  $f_{t_{subject}} = 0$  and  $i_{t_{subject}} = 0$ , and the following input- and forget-gate dynamics.

Taken together, these results describe the intricate mechanism underlying subject-verb agreement in LSTM number units. They also clarify why ablating either one of these two units may bring the network close to chance level on the NA-task. Without the stored information in the cell of a number unit the network hopelessly tries to solve the task.

### 5.1.3 Predicting the verb form

For long-range dependencies, storing the number of the subject in the cell is necessary for correct prediction in the NA-task (except for cases of mere chance) but is not sufficient. Two more conditions need to be met: (1) the stored number should be output from the cell at the right moment  $t_{verb} - 1$ ; and (2) The output from the number unit should increase activations only in units that represent the matching verb form in the output layer. For example, the output from unit 988 should be projected differently to singular and plural verb forms in the output layer, such that it will increase activity only in units representing the singular form. Given

672 that the encoding of singular by unit 988 is with  
 673 a negative value  $C_{t_{verb}-1} < -1$ , and therefore  
 674  $h_{t_{verb}-1} < -1$  (equation 1.x), this means that  
 675 weights from unit 988 to units in the output  
 676 layer that correspond to singular forms should  
 677 be negative too, but not those projecting to plu-  
 678 ral forms.  
 679

680 To see whether number units meet the first  
 681 condition, figure 3A shows the output-gate and  
 682 hidden state dynamics of units 988 and 776. In-  
 683 deed, the output gate opens at  $t_{verb}-1$ , reaching  
 684 its almost maximal value (Panels A-B). This en-  
 685 sures that the stored number information is out-  
 686 put from the cell and propagates to the output  
 687 layer via  $h_t$  (equation 1.X) at the right moment  
 688 (Panels C-D). Note also that for both units,  
 689 output-gate dynamics are quite similar across  
 690 all four conditions. This may seem sub-optimal,  
 691 since one may expect that the network would  
 692 learn to close the output gate for the irrele-  
 693 vant conditions (PS and PP for unit 987 and  
 694 SS and SP for unit 776). However, as we saw,  
 695 the cell value for these condition is anyway ap-  
 696 proximately zero and therefore an open output  
 697 gate will have the same effect as closed one.  
 698

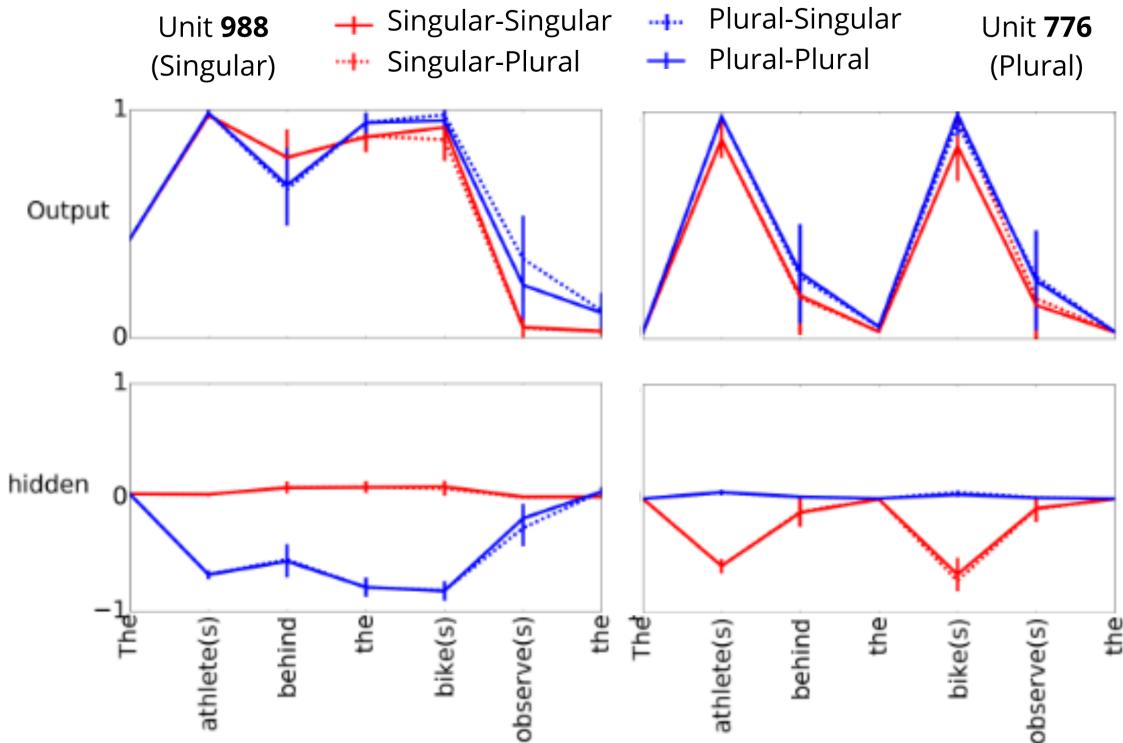
699 Next, to see whether number units meet the  
 700 second condition, figure 3B presents the distri-  
 701 bution of weight values from the two number  
 702 units and from several other units to 36 units  
 703 at the output layer - 18 corresponds to verbs in  
 704 the singular form and 18 to verbs in the plural  
 705 forms (TODO: extend to more/all verbs in the  
 706 vocab). Clearly, for number units, weights to  
 707 singular and plural forms have different values,  
 708 but for other non-number units, there's no clear  
 709 structure. Moreover, the weight values corre-  
 710 spond to the encoding of singular/plural in the  
 711 number unit. For example, weight values from  
 712 unit 988 to singular forms are indeed negative,  
 713 and those to plural forms are not. For all other  
 714 units as well, the sign of  $h_{t_{verb}-1}$  corresponds to  
 715 the sign of the relevant weight values (panels A-  
 716 B), such that their product is always positive.  
 717 This ensures that number units increase activa-  
 718 tions only in the matching units in the output  
 719 layer.

### 5.1.4 Short-range number units

I got lost in this discussion. Why does  
 the continuous code requirement have to  
 be met? The fundamental thing to be  
 shown here should be that ablating *all* the  
 candidate SR units (plus the LR units)  
 and show they cause a large decrease in  
 performance on the SR NA-tasks (con-  
 trolling with random ablations of simi-  
 lar numbers of units). What's here now  
 seems more like a *post-hoc* analysis of  
 these units. We saw in section 5.3.1 that per-  
 formance on several NA-tasks was not impaired  
 after the ablation of any unit, nor by ablating  
 the LR-number units 776 and 988, which sug-  
 gested that number information may be encoded  
 also elsewhere in the network, in other units.  
**Cite Dieuwke's paper on this.** We there-  
 fore start by defining requirements that a unit  
 should meet to count as a (grammatical) *num-  
 ber unit*. For convenience, we focus on sentences  
 that start by a subject and define two such re-  
 quirements:

1. *Continuous code*: The grammatical-  
 number of the subject should be decodable  
 from unit's cell activity in a continuous  
 manner. Continuity means that decoding  
 performance is high for *all* subsequent time  
 points and that the code for grammatical  
 number is the same code as that used by  
 the unit during the presentation of the  
 subject. This excludes units for which  
 grammatical code can be decoded from  
 their activity at subsequent time points  
 but that there is no consistent code, e.g.,  
 cell values for singular and plural flips  
 from one time point to another. Units that  
 meet this requirement for continuous code  
 could be therefore interpreted as *storing*  
 grammatical-number information.
2. *Number-segregated efferent weights*: The ef-  
 feren weights of a number unit should be  
 segregated for singular and plural forms of  
 the verb. This is similarly to what was de-  
 scribed above for the LR number units.

To quantify the first requirement, we use

768  
769  
770  
771  
772  
773  
774816  
817  
818  
819  
820  
821  
822

794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815

**Figure 3:** Hidden and output-gate activations during processing of a sentence with a prepositional phrase between subject and verb. Do we need both? Isn't it a bit redundant? I mean, I see you can make the point of the output gate behaving in the same way in both cases, but that's perhaps a bit minor?

842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

a generalization across time (GAT) approach (cite). Specifically, for each unit, we trained a model to predict the grammatical number of the subject based on cell activity during the presentation of the subject but tested its prediction based on activity *from all time points*. We evaluate model performance using the Area under of Curve (AUC), and for each time point, we define 'high' performance as AUC values greater than 0.95. To quantify the second requirement, given two sets of weights to singular and plural verb forms  $W_S, W_P$  of unit  $u$ , with means  $\bar{W}_S, \bar{W}_P$  and standard deviations  $\sigma_S, \sigma_P$ , we define a statistical distance:  $d_{W_P, W_S}^u = \frac{|\bar{W}_S - \bar{W}_P|}{\sigma_S + \sigma_P}$ . A unit is then said to have highly segregated weights if  $d_{W_P, W_S}^u$  is at the top five percentile across all units. Units that meet Finally, as discussed in

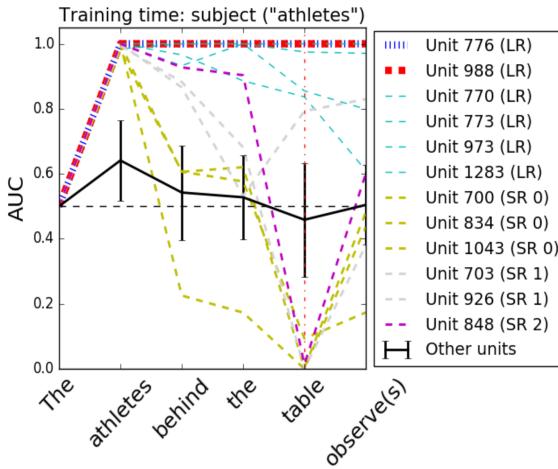
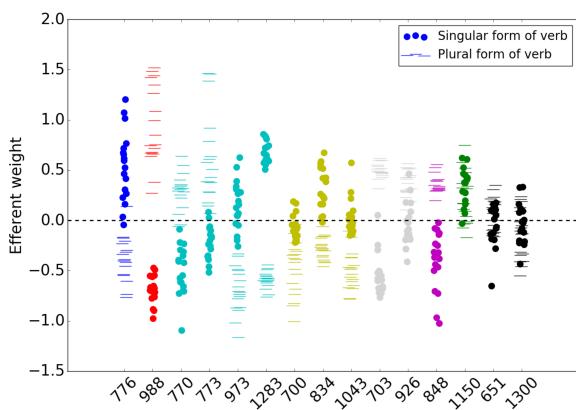
section 5.3.1, we distinguish between two types of units: (1) SR-units - units

Figure 4B presents shows the Area Under of Curve (AUC) for all units in the network that meet these two requirements.

## 5.2 Syntax units

In the previous section, a small number of units in the network was identified, which was found to store grammatical number for long-range dependencies also in the presence of interfering nouns. Remarkably, information about the two possible values of grammatical number, singular and plural, is steadily carried by two specialized units in the network: unit 776 and 998. It remains unclear, however, how is the timing of storing and updating of number information in

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911



**Figure 4:** Connectivity structure to output layer. (A) Output activity  $h_t$  of all number units during the processing of a sentence with a PP between subject and verb. (B) Weight values from various units to output layer. Note that only for number units the output weights are clearly separated between singular and plural form of the verb, either positive or negative, compare to the syntax unit (1149) and two non-number units in the second layer. (C) Visualization of 18 verbs in their plural and singular forms (36 words in total) on the plane spanned by the two first principal components of their embeddings by the output weight matrix. A clear separation is observed between the singular and plural form along the first PC.

these units controlled by the network. In particular, given that the storage of information is controlled by the input and forget gates, how is the timing of these gates is controlled by the network. We therefore hypothesized that other units in the network may encode information about the syntactic structure of the sentence, which can be propagated to LR number units at moments required to accomplish the NA-task.

To identify such units, we tested whether there are units in the network from which syntactic information can be decoded. Given that ablation of single units in the network revealed number units only, we concluded that syntactic information may be encoded by the network in a more distributed way. We therefore tested decoding from the joint activity of a group of units in the network, using a regression model to predict syntactic information from the pattern of activity in the entire network. Specifically, the dependent variable of the model was syntactic complexity, which was quantified by counting the number of open nodes in the syntactic-tree, inspired by [cite Nelso et. al 2017](#). In what

follows, we refer to the number of open nodes in the syntactic tree as syntactic-tree depth. We then used the activity of all units in the network as predictors in the model, together with word frequency as a covariate. In addition, to address possible co-linearities among unit activities, and to examine from which units syntactic-tree depth can be best predicted, we L1-regularized the regression model (section 4.2 for details [section 4.2 should explain that the model was trained on the data from 3.2, word frequency, standarization of the values, etc.](#)).

Figure 5A shows the resulting weights of all predictors in the model. We define as *outlier weights*, weights that are more than three standard-deviations away from the mean weight size [add mean and std values](#). 17 outlier weights were found in the model [?add resulting unit numbers?](#), whose corresponding units may thus be carrying syntactic information, given their high weight in predicting syntactic complexity.

If these units indeed carry information about the syntactic structure of the sentence, ablating these units would lead to a reduction in network

912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959

960 performance on the NA-task We used Linzen for  
 961 this. Do we want to also test of nounPP, etc?.  
 962 To test this, we conducted an ablation study,  
 963 in which groups of 17 random units are ablated,  
 964 and the resulting network is then tested on the  
 965 NA-task. We performed 1000 such random-  
 966 ablation experiments, and another one in which  
 967 we ablated the 17 units resulted from the re-  
 968 gression model. Figure 5B presents the null-  
 969 distribution of network performance on the NA-  
 970 task, generated from the 1000 random-ablation  
 971 experiments. Reduction in network performance  
 972 when ablating the 17 units from the regres-  
 973 sion model (black arrow) was found significant  
 974 ( $p - value = 0.024$ ). Taken together, these re-  
 975 sults suggest that the 17 units identified by the  
 976 regression model carry syntactic information re-  
 977 quired for successfully performing the NA-task.  
 978 In what follows, we therefore refer to these units  
 979 as *syntax units*.

980 Next, we looked into the dynamics of these  
 981 units by visualizing their gate and cell dynam-  
 982 ics during the processing of sentences with vari-  
 983 ous syntactic structures. In particular, we found  
 984 that unit 1150, which resulted with the highest  
 985 weight in the regression study, seems to have  
 986 an interpretable dynamics during the process-  
 987 ing of various types of sentences. Figure 2 & 3  
 988 shows the hidden and cell activity of this unit in  
 989 the nounPP, subjrel, objrel and double-subjrel  
 990 tasks. In all these cases, the value of the hidden  
 991 variable of this unit follows the structure of the  
 992 embedded phrase separating the main subject-  
 993 verb dependency. Its activity is positive dur-  
 994 ing the processing of the entire embedded phrase  
 995 and transitions to a negative value at the time  
 996 when the main verb is presented, but not be-  
 997 fore, even if other verbs preceded it (as in the  
 998 examples with relative clauses). This suggests  
 999 that unit 1150 encodes the existence of an em-  
 1000 bedded phrase separating the main subject-verb  
 1001 dependency in the sentence. Interestingly, cell  
 1002 activity of this unit consistently increased during  
 1003 the processing of the embedded phrase, simi-  
 1004 larly to what was observed in several electrodes  
 1005 in the language network in intracranial studies  
 1006 on sentence processing cite nelson 2017 - refer to  
 1007 figure. TODO: for this, also add to our figures

1008 cell activity of unit 1150. In the next section, we  
 1009 show how these dynamics of unit 1150 are im-  
 1010 portant in controlling the timing of gate activity  
 1011 in the LR number units - 776 and 988.

1012 ?Do we also want to note that unit 1150  
 1013 had an effect also on the NA-task when ablated  
 1014 alone...?

### 5.3 Processing of relative clauses

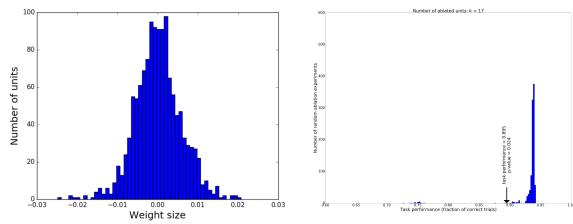
[1]  
 [1]

## 6 Discussion

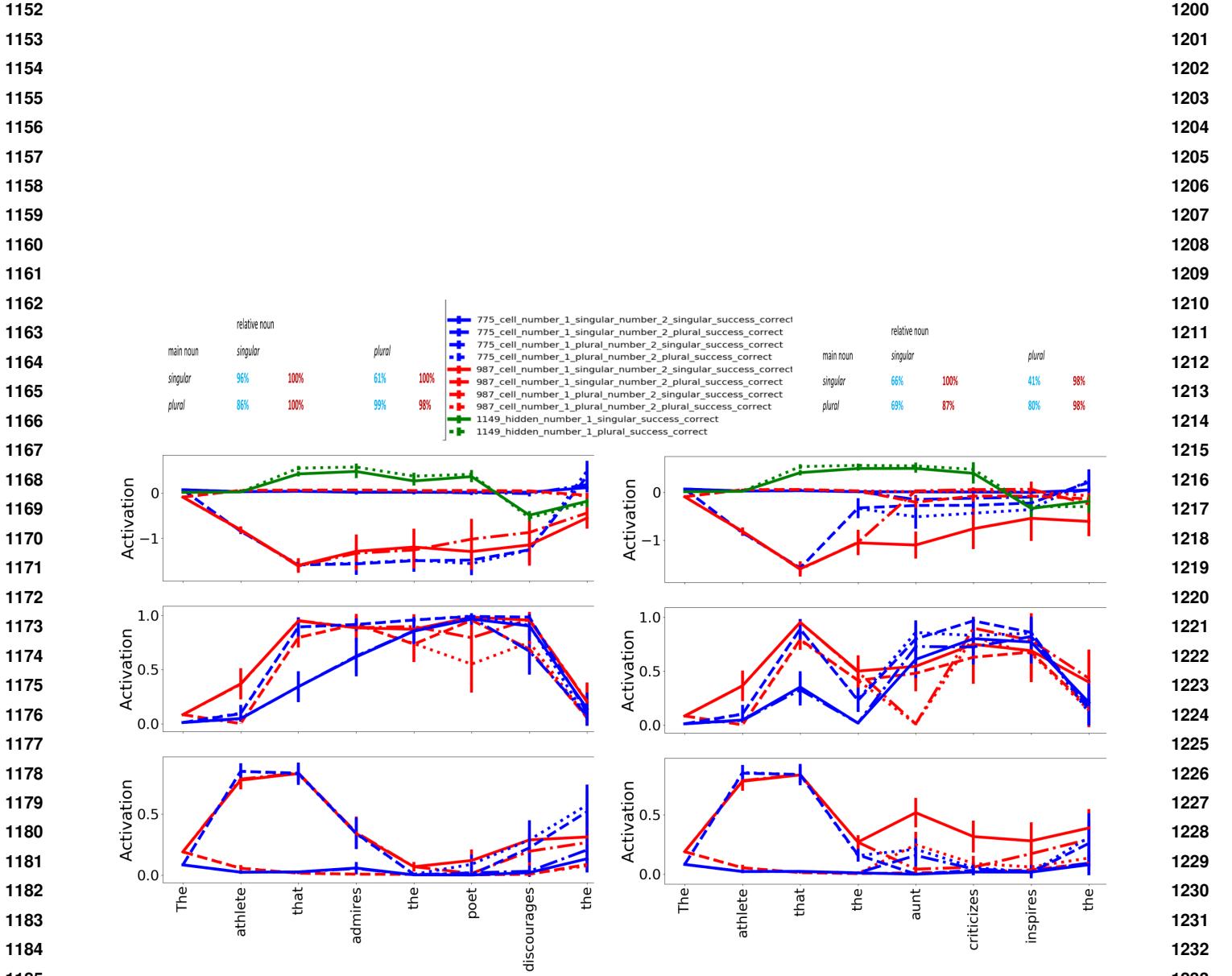
[1] Bock + asymmetry plural/singular (non-  
 phonological explanation)

1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025  
 1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055

1056	1104
1057	1105
1058	1106
1059	1107
1060	1108
1061	1109
1062	1110
1063	1111
1064	1112
1065	1113
1066	1114
1067	1115
1068	1116
1069	1117
1070	1118
1071	1119
1072	1120
1073	1121
1074	1122
1075	1123
1076	1124
1077	1125
1078	1126
1079	1127
1080	1128
1081	1129
1082	1130
1083	1131
1084	1132
1085	1133
1086	1134
1087	1135
1088	1136
1089	1137
1090	1138
1091	1139
1092	1140
1093	1141
1094	1142
1095	1143
1096	1144
1097	1145
1098	1146
1099	1147
1100	1148
1101	1149
1102	1150
1103	1151



**Figure 5:** (A) Distribution of the resulting weight values from the tree-depth regression model. Outlier weights were defined as having a value that is distant from the mean by more than three standard deviations (17 outlier weights in total - marked in red). (B) Task performance of 1000 models after ablating 17 random units (in blue) and based on the 17 outlier weights from the tree-depth regression model (black arrow). The reduction in performance due to outlier-weights ablation is statistically significant ( $p - value < 0.05$ ) when compared to the null distribution generated by the random ablations.



**Figure 6:** Subject-verb agreement in relative clauses: agreement-task accuracy for (A) subject relatives and (B) object relatives. (C & D) The corresponding cell activations for the number units (775 and 987) and the syntax unit 1149. (E & F) The corresponding forget-gate activity and (G & H) input-gate activity of the number units.

