

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: (?).

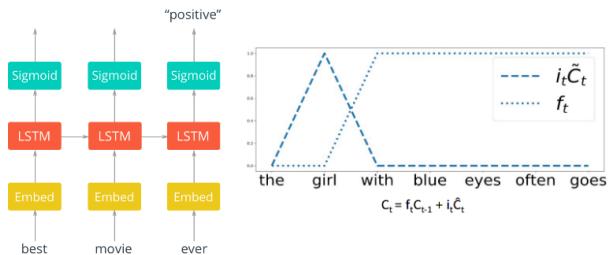


Figure 1: Caption.

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

2.2 Subject-verb agreement in English

- Subject-verb agreement from psycholinguistics (e.g., Miller and Bock, Franck and Rizzi) - SV agreement in LSTMs (Linzen 2016, Gulordava 2018, etc.). - Relate to Nelson et. al 2017 PNAS, an intracranial study that identifies electrodes whose high-gamma activity correlates with syntactic tree-depth (number of open nodes)

3 The data

[1]

3.1 Synthetic data

3.1.1 Stimuli for the number-agreement task (NA-task)

This section describes the generation process of synthetic sentence stimuli for the number-

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	agreement task (NA-task). These sets of stimuli	144
097	were used to evaluate the performance of both	145
098	full- and ablated models on the NA-task.	146
099	Each NA-task contained sentences with a	147
100	fixed syntactic structure, such as “Det Noun	148
101	Adv Verb” or “Det Noun P Det Noun Verb”, and	149
102	each task was composed of several <i>conditions</i> de-	150
103	pending on the possible assignments of gram-	151
104	matical number to the noun(s) in the sentence.	152
105	For example, one NA-task contained sentences	153
106	of the form: “Det Noun Verb”, and had two con-	154
107	ditions corresponding to the two possible num-	155
108	ber values of its noun. So the first condition	156
109	such as “The boy runs”, and the ablated net-	157
110	work was tested on predicting the correct verb	158
111	form. This task had two conditions, correspond-	159
112	ing to the two possible assignments of gram-	160
113	matical number (singular or plural) to the main	161
114	noun. Another NA-task contained sentences of	162
115	the form: “Det Noun-1 P Det Noun-1 Verb”,	163
116	such as “The boy behind the girls jumps”. This	164
117	task had four conditions, corresponding to the	165
118	four possible assignments of grammatical num-	166
119	ber to noun-1 and noun-2.	167
120	The network was evaluated on predicting the	168
121	correct verb form (singular or plural).	169
122		170
123	3.2 Corpus data	171
124	[1]	172
125		173
126	4 The models	174
127		175
128	I would describe in this section the regression	176
129	and language models.	177
130		178
131	4.1 LSTM language model	179
132		180
133	4.1.1 Architecture and dynamics	181
134	TODOs: - Number and type (embedding,	182
135	LSTM, softmax) of layers, including dimensions.	183
136	- Describe the dynamics of the network (list of	184
137	equations of the LSTM)	185
138	4.1.2 Model training and evaluation	186
139		187
140	- Task (refer to Gulordava et. al) - Describe	188
141	and give (nick) names for the other models, e.g.,	189
142	LM-high-dropout, LM-SEED1 - The model was	190
143	evaluated on a left-out test set...perplexity..	191

of a verb immediately following a subject. Finally, subsection 5.1.4 characterizes the structure of the efferent weights of number units, which propagate number information to the output layer.

5.1.1 Local vs. distributed code - an ablation study

Generally, number information may be stored in the network in either a local, sparse, or a distributed way, depending on the fraction of active units that carry number information. We hypothesized that if the network uses a local or sparse coding, meaning that there's a small set of units that encode number information, 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.

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 pre-

NA task	C	770	776	988	1283	Full	240
Simple	S	-	-	-	-	100	241
Adv	S	-	-	-	-	100	242
2Adv	S	-	-	-	-	99.8	243
Co-Adv	S	-	-	84.0	84.0	98.8	244
namePP	S	-	-	-	-	98.9	245
nounPP	SS	-	-	-	-	97.5	246
nounPP	SP	-	-	58.8	-	88.5	247
subjrel	SS	-	-	88.0	-	97.0	248
subjrel	SP	-	-	-	-	58.8	249
objrel	SS	-	-	-	-	64.7	250
objrel	SP	-	-	-	-	45.7	251
Simple	P	-	-	-	-	100	252
Adv	P	-	-	-	-	99.6	253
2Adv	P	-	-	-	-	99	254
Co-Adv	P	-	78.9	-	-	99.7	255
namePP	P	-	57.6	-	-	66.8	256
nounPP	PS	85.2	49.7	-	-	93.2	257
nounPP	PP	-	81.7	-	-	98.3	258
subjrel	PS	85.8	58.6	-	-	87.8	259
subjrel	PP	-	88.1	-	-	99.3	260
objrel	PS	-	-	-	-	69.0	261
objrel	PP	-	-	-	-	81.0	262
Linzen	-	?	?	?	?	?	263

Table 1: 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.

viously 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 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.

288	We next describe several important aspects	336
289	of the ablation-experiment results. First,	337
290	in all NA-tasks, only four units from the	338
291	entire network (1300 LSTM units in total)	339
292	had a significant effect on task performance.	340
293	This result suggests a local coding scheme	341
294	for long-range grammatical-number information	342
295	(TODO: quantify a 'significant' reduction, and	343
296	perhaps link to the emergence of local coding	344
297	in neural-network simulations (Bowers) and to	345
298	findings about grandmother neurons in humans	346
299	(e.g., Fried)). Second, we note that all number	347
300	units emerged at the second layer of the net-	348
301	work. This seems appropriate if number informa-	349
302	tion needs to be directly projected to the out-	350
303	put layer for correct verb-form prediction. In	351
304	section 5.1.4 we further explore the projection	352
305	weights from number units. Third, for simple,	353
306	1Adv and 2Adv NA-tasks, none of the units had	354
307	a significant effect on task performance. This	355
308	suggests that for short-range dependencies num-	356
309	ber information may be also encoded elsewhere	357
310	in the network, perhaps via a more distributed	358
311	code. We therefore make a distinction between	359
312	long-range (LR) and short-range (SR) number	360
313	units in what follows. We return to this point	361
314	in section 5.1.3 (TODO: complete the identifica-	362
315	tion of short-range number units from the result-	363
316	ing weights of the classifier in the generalization-	364
317	across-time experiment). Fourth, LR-number	365
318	units can be further divided into two types, de-	366
319	pending on the grammatical number of the sub-	367
320	ject. Units 770 and 776 had a significant effect	368
321	only when the first noun was plural, but not sin-	369
322	gular, and vice versa for units 988 and 1283 (blue	370
323	and red in table 1, respectively). We therefore	371
324	refer to the former as <i>plural units</i> and to the lat-	372
325	ter as <i>singular units</i> . Finally, we note that two of	373
326	the number units (776 & 988) had an exceptional	374
327	effect on network performance in both nounPP-	375
328	SP&PS conditions. These two conditions are	376
329	in particular revealing since they involve both	377
330	a long-range dependency (over a prepositional	378
331	phrase) and an interfering noun before the verb,	379
332	while performance of the non-ablated network	380
333	is still relatively high (88.5%&93.2%, respec-	381
334	tively) in contrast to these conditions in sub-	382
335	jrel and objrel. Ablating one of these two units	383
	brought the network from high performance on	
	the NA-task to around chance-level performance	
	(58.8%&49.7%, respectively). In the next sec-	
	tion, 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 par-	
	ticular, we find that in some conditions two units	
	can bring the network from relatively high per-	
	formance to around chance-level performance on	
	the NA-task (section 5.1.1). However, it remains	
	unclear what is the exact mechanism underly-	
	ing successful trials in the NA-task, and what	
	goes wrong in unsuccessful ones. To better un-	
	derstand this, we now looks into gate and state	
	dynamics of these units during the processing of	
	sentences from the nounPP NA-task.	
	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 activi-	
	ty. 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 pro-	
	duce correct number agreement, it seems that	
	number units should at least have the follow-	
	ing three properties: (1) The grammatical num-	
	ber of the subject should first be encoded by	
	$\tilde{C}_{t_{subject}}$, encoding singular and plural with <i>dif-</i>	
	<i>ferent</i> values. (2) To grant the encoded gram-	
	matical 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 informa-	
	tion 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 relevant 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

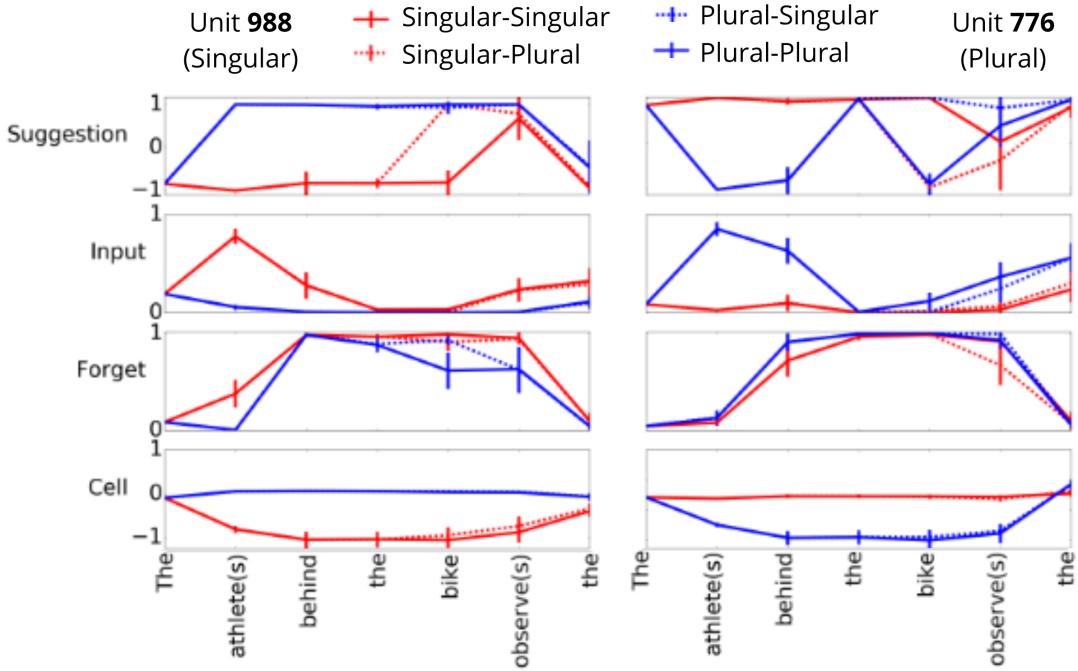
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575

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.

relevant conditions (PP and PS) but not in the 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 that the encoding of singular by unit 988 is with a negative value $C_{t_{verb}-1} < -1$, and therefore $h_{t_{verb}-1} < -1$ (equation 1.x), this means that weights from unit 988 to units in the output layer that correspond to singular forms should be negative too, but not those projecting to plu-

576	ral forms.	624
577	To see whether number units meet the first	625
578	condition, figure 3A shows the output-gate and	626
579	hidden state dynamics of units 988 and 776. In-	627
580	deed, the output gate opens at $t_{verb}-1$, reaching	628
581	its almost maximal value (Panels A-B). This en-	629
582	sures that the stored number information is out-	630
583	put from the cell and propagates to the output	631
584	layer via h_t (equation 1.X) at the right moment	632
585	(Panels C-D). Note also that for both units,	633
586	output-gate dynamics are quite similar across	634
587	all four conditions. This may seem sub-optimal,	635
588	since one may expect that the network would	636
589	learn to close the output gate for the irrele-	637
590	vant conditions (PS and PP for unit 987 and	638
591	SS and SP for unit 776). However, as we saw,	639
592	the cell value for these condition is anyway ap-	640
593	proximately zero and therefore an open output	641
594	gate will have the same effect as closed one.	642
595	Next, to see whether number units meet the	643
596	second condition, figure 3B presents the dis-	644
597	tribution of weight values from the two number	645
598	units and from several other units to 36 units	646
599	at the output layer - 18 corresponds to verbs in	647
600	the singular form and 18 to verbs in the plural	648
601	forms (TODO: extend to more/all verbs in the	649
602	vocab). Clearly, for number units, weights to	650
603	singular and plural forms have different values,	651
604	but for other non-number units, there's no clear	652
605	structure. Moreover, the weight values corre-	653
606	pond to the encoding of singular/plural in the	654
607	number unit. For example, weight values from	655
608	unit 988 to singular forms are indeed negative,	656
609	and those to plural forms are not. For all other	657
610	units as well, the sign of $h_{t_{verb}-1}$ corresponds to	658
611	the sign of the relevant weight values (panels A-	659
612	B), such that their product is always positive.	660
613	This ensures that number units increase activa-	661
614	tions only in the matching units in the output	662
615	layer.	663
616	5.1.4 Short-range number units	664
617	We saw in section 5.3.1 that performance on	665
618	several NA-tasks was not impaired after the	666
619	ablation of any unit, nor by ablating the LR-	667
620	number units 776 and 988, which suggested that	668
621	number information may be encoded also else-	669
622	where in the network and thus available for	670
623		671
	short-range dependecies without an interfering	
	noun. To explore this, we tested whether there	
	are units in the network from which grammatical	
	number of the subject can be decoded at time	
	points <i>following the presentation of the subject</i> .	
	High decoding performance would be	
	5.2 Syntax units	
	[1]	
	5.2.1 Predicting syntactic-tree depth	
	from network activity	
	[1]	
	5.2.2 Ablation study	
	[1]	
	5.3 Syntax-number units interactions	
	[t] [1]	
	5.4 Processing of relative clauses	
	[1]	
	[1]	
	6 Discussion	
	[1] Bock + asymmetry plural/singular (non-	
	phonological explanation)	
	Acknowledgments	
	[1]	

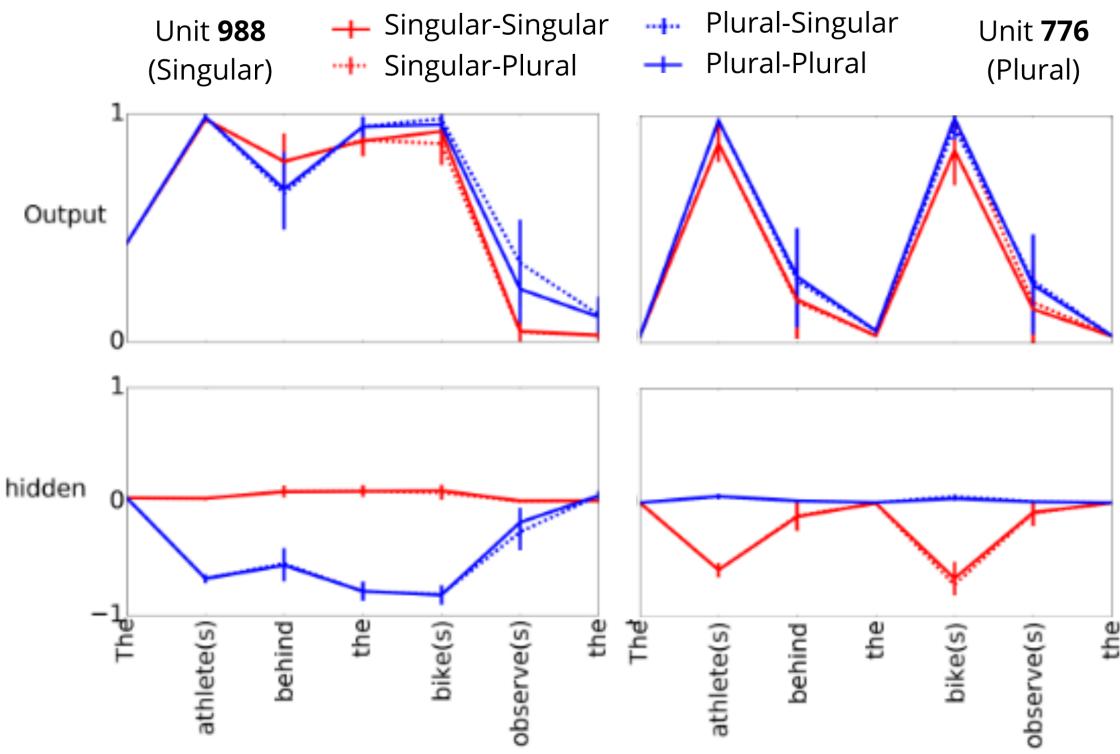


Figure 3: Hidden and output-gate activations during processing of a sentence with a prepositional phrase between subject and verb.

768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815

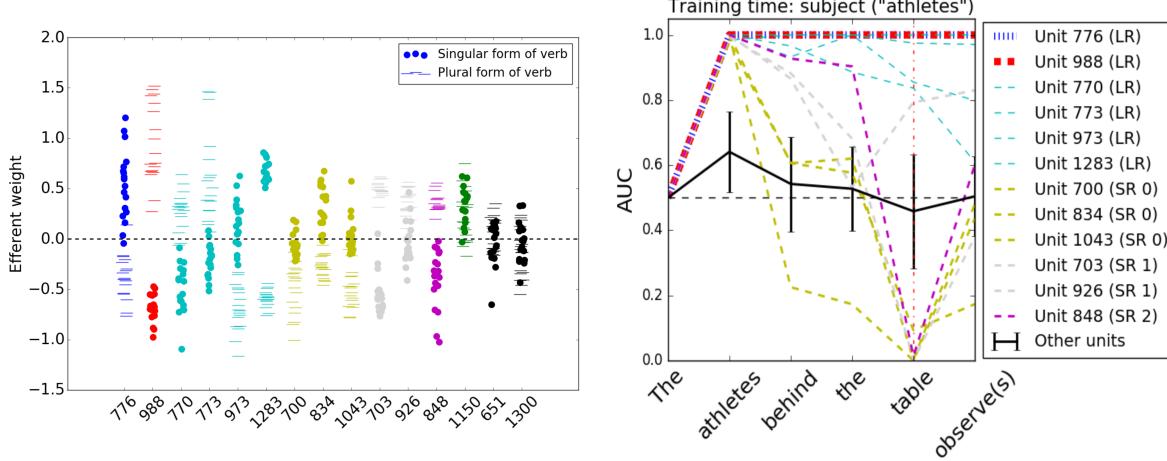


Figure 4: Connectivity structure to output layer. (A) Output activity h_t of all number units during the processing of a sentence with a PP between subject and verb. (B) Weight values from various units to output layer. Note that only for number units the output weights are clearly separated between 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.

864	912
865	913
866	914
867	915
868	916
869	917
870	918
871	919
872	920
873	921
874	922
875	923
876	924
877	925
878	926
879	927
880	928
881	929
882	930
883	931
884	932
885	933
886	934
887	935
888	936
889	937
890	938
891	939
892	940
893	941
894	942
895	943
896	944
897	945
898	946
899	947
900	948
901	949
902	950
903	951
904	952
905	953
906	954
907	955
908	956
909	957
910	958
911	959

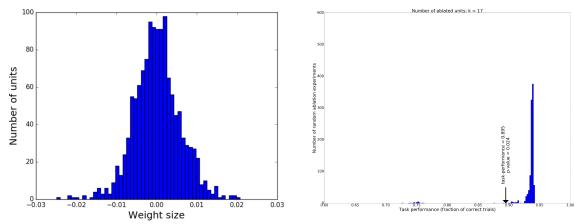
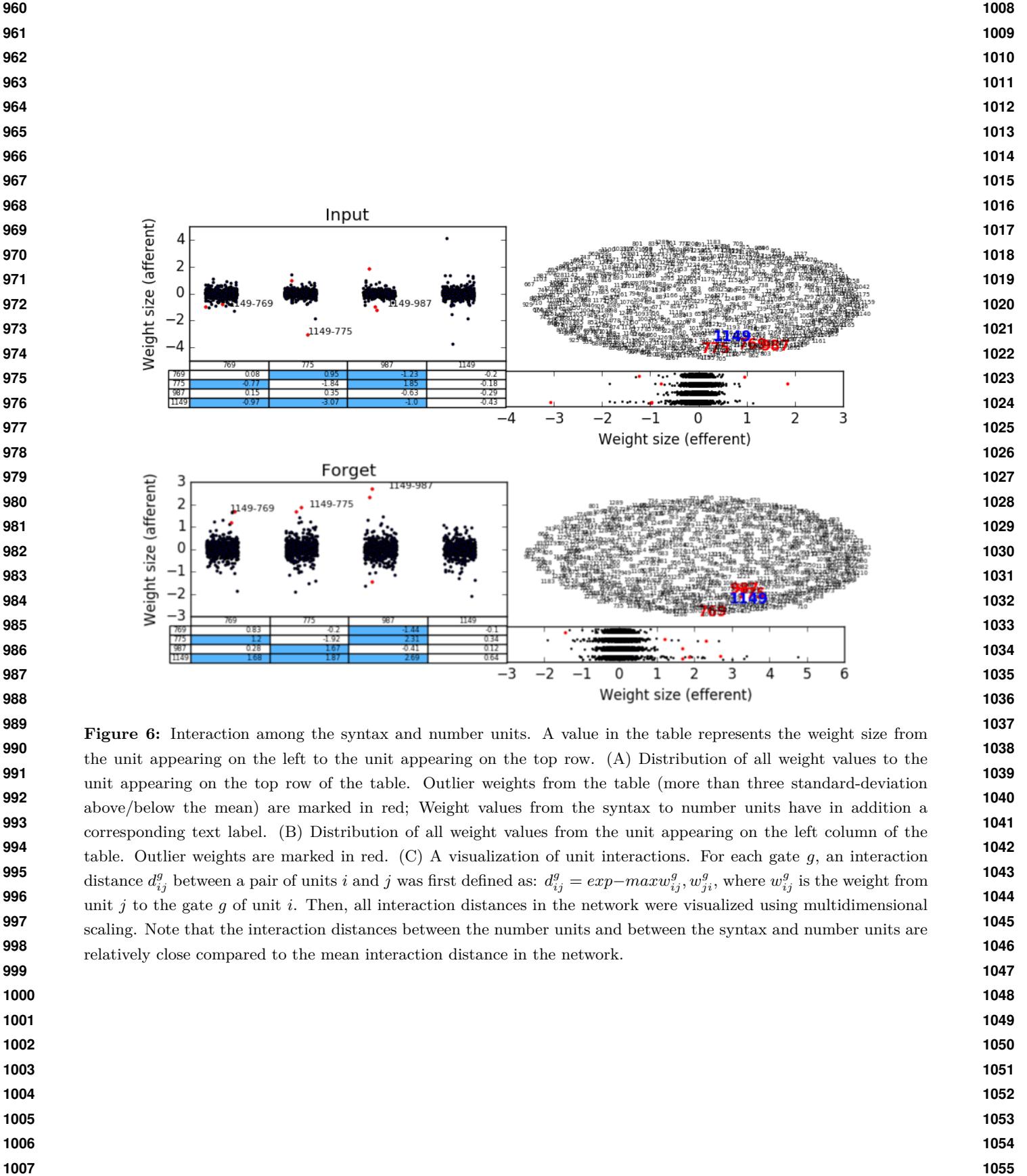


Figure 5: (A) Distribution of the resulting weight values from the tree-depth regression model. Outlier weights were defined as having a value that is distant from the mean by more than three standard deviations (17 outlier weights in total - marked in red). (B) Task performance of 1000 models after 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.



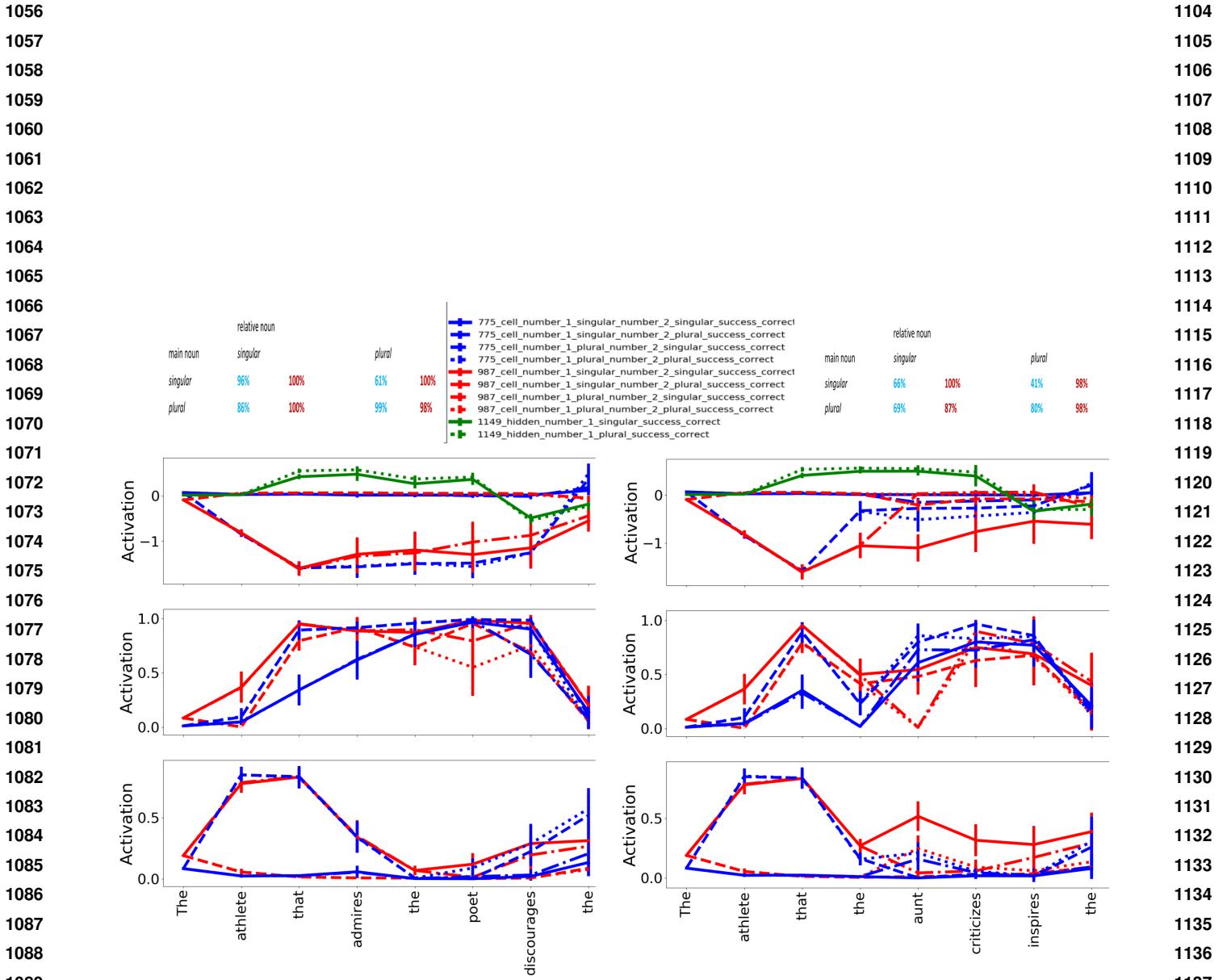


Figure 7: Subject-verb agreement in relative clauses: agreement-task accuracy for (A) subject relatives and (B) object relatives. (C & D) The corresponding cell activations for the number units (775 and 987) and the syntax unit 1149. (E & F) The corresponding forget-gate activity and (G & H) input-gate activity of the number units.

