

1 We wish to thank all of the reviewers for their time and thorough reading of our paper! Specific concerns are addressed:

2 **Reviewer #1** We appreciate the reviewer’s suggestions regarding clarity. To improve this, we have: (1) Added a  
3 preliminaries section which introduces our mathematical notation; (2) Highlighted key results and observations that  
4 were previously buried inside paragraphs in sections 4.2, 4.4, and 4.5; and (3) Added more detail to the captions  
5 in Figures 1, 3, and 4. We chose not to restructure the mathematical results as theorem/lemma/result (as suggested by  
6 the reviewer) as we felt that most of our mathematical statements were largely definitions (such as eq. (3)).

7 **Reviewer #2** Addressing the suggested improvements: (1) We have added the suggested summary sentence “the key  
8 activity performed by the RNN for sentiment analysis is simply counting the number of positive and negative words  
9 used” to the discussion. (2) We started with binary sentiment classification, but are actively working on more tasks. For  
10 multi-level sentiment classification (e.g. 5-way), our hypothesis is that the networks will still use a 1D line attractor,  
11 but that this attractor will be curved such that different readouts will partition different sections of the line attractor  
12 (corresponding to different levels of evidence), but still yielding low-dimensional dynamics. We are currently running  
13 this experiment and will include its results in the final version of the paper (likely in the supplement, due to space  
14 constraints). (3) The classification accuracy of the Jacobian linearized model is much worse than the LSTM, due to  
15 small errors in the linear approximation that accrue as the network processes a document. Note that if we directly train  
16 a linear model, the performance is quite high (only around 3% worse than the LSTM), which suggests that the error  
17 of the linearized model has to do with errors in the approximation, not from having less expressive power. We have  
18 included a few sentences about this in the discussion. (4) We have added a derivation of the expression relating the  
19 eigenvalue ( $\lambda$ ) to the time constant ( $\tau$ ) in the supplement, along with a corresponding reference in the main text.

20 **Reviewer #3 (Major point 1.)** We agree with  
21 the reviewer that a systematic study of the vari-  
22 ability we see in the dynamical structures in our  
23 analysis is warranted. Assessing if and how this  
24 variability is related to performance differences is  
25 something we wish to pursue in future work. We  
26 have begun some of these investigations, and have  
27 found that the small differences in drift and Q val-  
28 ues do not seem to affect the performance (their  
29 values are too small to have an effect over typical  
30 document lengths in these datasets). **(Major point**  
31 **2.)** As mentioned in the discussion, we have yet  
32 to systematically analyze negation bigrams. We  
33 have done some preliminary analysis (see Figure  
34 1, at right) which suggests that RNNs are capable  
35 of correctly accounting for ‘not’ tokens. We have  
36 a few ideas for how to uncover these mechanisms  
37 (e.g. using switched linear approximations), how-  
38 ever, this remains as future work. **(Minor point**  
39 **4.)** Added more detail to the Figure 1 caption, ex-  
40 plaining that it is many neurons for one document.  
41 **(Minor points 5. and 3.)** Added a reference to the  
42 accuracy table in the appendix in Section 3.1. The  
43 bag of words does have performance close to that  
44 of RNNs, especially for smaller datasets (providing  
45 further support for using linear approximations of  
46 the RNNs). **(Minor point 6.)** Regarding the input  
47 point around which we linearize: in the paper, we  
48 linearized around zero input. We also tried lineariz-  
49 ing around the average embedding of all words, this  
50 does not change the results (indeed, the average embedding of all words is very close to the zeros vector—the norm  
51 of the average embedding is  $7.6 \times 10^{-3}$ ). We have added a footnote noting this in the main text. **(Minor point 7.)**  
52 Removed the incorrect reference to Fig. 1D. **(Minor point 8.)** Fixed typo. **(Minor point 9.)** We have not correlated  
53 the performance with things like the input projections. The projection histograms are over the top positive and negative  
54 words, whereas the performance (over test examples) depends on the particular words that show up in those examples;  
55 as such, it may not make sense to correlate them.

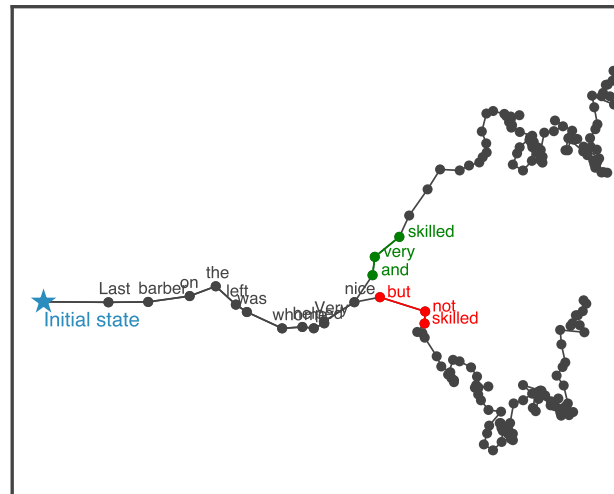


Figure 1: Probing the RNN with negation bigrams. Projection of RNN hidden states onto the top two PCs for two different input sequences that differ only by two tokens (replacing ‘and very’ with ‘but not’ in the middle of the sequence). The trajectories start out the same as the initial tokens are identical. They then diverge at the critical tokens, moving in opposite directions along the readout (the readout is aligned with the y-axis; not shown). After these two tokens, the rest of the sequence is also identical (tokens not shown to remove clutter). Note how the presence of the negation bigram changes the effect of future tokens on the hidden state.

50 does not change the results (indeed, the average embedding of all words is very close to the zeros vector—the norm  
51 of the average embedding is  $7.6 \times 10^{-3}$ ). We have added a footnote noting this in the main text. **(Minor point 7.)**  
52 Removed the incorrect reference to Fig. 1D. **(Minor point 8.)** Fixed typo. **(Minor point 9.)** We have not correlated  
53 the performance with things like the input projections. The projection histograms are over the top positive and negative  
54 words, whereas the performance (over test examples) depends on the particular words that show up in those examples;  
55 as such, it may not make sense to correlate them.