1 We thank the reviewers for their helpful comments. Overall, reviewers said the work was **timely** (R2, R3) and **important**
2 (R1, R3), with R3 adding: "the impact of the released model and released dataset is likely to be huge in the community."

3 We are excited that reviewers rated our Grover generator as being a **simple** (R3) yet **novel** (R1, R2) approach for
4 multi-field document generation. Reviewers noted that Grover generates "extremely credible" articles (R2) and that due
5 to its inductive bias as a generator (R3) it is also a state-of-the-art detector of neural fake news.

6 Last, reviewers said that our experiments on neural fake news detection in a semi-supervised setting were extensive
7 (R1) and "very well done and exciting" (R3), particularly analysis about when the adversary has more resources.

─────────────── **Reviewer 1** ───────────────

9 (Contribution 3, Comment 1) / (Clarity Clarification 1): **machine authorship vs. fake?** Good point, we will clarify our
10 terminology in revision. Just to clarify, Grover is pretrained only on human-written truthful news; propaganda websites
11 are excluded. We did additional experimentation and found that Grover discriminates between human-written real news
12 vs. human-written propaganda with **98% accuracy**, and will add this (and discussion) to the revision.

13 (Contribution 3, Comment 2): **evaluating consistency with metadata?** Good idea: we did additional human evaluation
14 on the consistency of the article body with the headline, date, and author. We found that **generations are largely**
15 **consistent** overall. For instance, human propaganda articles are consistent with the headline with an average score of
16 2.85/3 (higher is better) but machine-written propaganda gets 2.64/3; the news scores are similar.

17 (Contribution 3, Comment 3): **spotting generations of other models?** We ran additional experiments here: Grover
18 detects GPT2-generated news (released by OpenAI) with **96% accuracy**, even without finetuning on GPT2 generations.

19 (Philosophical comment): **fake news term?** We appreciate this point and will revisit the word choice. In our draft,
20 we followed the lead of the following political science paper, which has an extensive discussion about recommended
21 terminology: *Lazer, David MJ, et al. "The science of fake news." Science 359.6380 (2018): 1094-1096.*

22 (Contribution 1, Q): **All fields share the same vocabulary** - we'll clarify this in revision. We haven't seen the model
23 generate field-inappropriate tokens (like an invalid date), perhaps due to extensive pretraining and nucleus sampling.

─────────────── **Reviewer 2** ───────────────

25 **Novelty of Grover over GPT2?** We believe that our "novel way to guide generation" makes Grover novel, not just an
26 'adaptation'. Indeed, GPT(2), BERT, XLnet, and Grover share the same backbone but learn from different objectives.

27 **What is given to the turkers?** We will provide the full prompt in revision along with other details (we used 3
28 annotators) and discussion. For overall trustworthiness for instance, we asked "Does the article read like it comes
29 from a trustworthy source?" Thus it emphasizes style, while content sensibility measures whether the semantic
30 content is believable; thereby **disassociating style vs. content**. The results in the paper show that news written in a
31 propaganda style appears less trustworthy, but the *content* of human and machine propaganda is equally sensible. The
32 real news→Grover news results likely generalize to propaganda→Grover propaganda, were this desired.

33 **"It takes a thief to catch a thief"?** We thank the reviewer for this comment and will revise the paper to be more precise.
34 However, we *did* test this hypothesis beyond Grover. We trained our BERT model on RealNews in a multi-field setting.
35 Nevertheless, BERT is worse at neural fake news discrimination compared with Grover. We found this surprising
36 because NLP leaderboards for discriminative tasks, like GLUE, show the dominance of deep bidirectional models like
37 BERT over unidirectional ones like GPT. Even though Grover cannot handle right-to-left dependencies, it still is a
38 state-of-the-art discriminator because its inductive bias matches that of a generator's (R3).

39 **Choice of sources?** We used only those propaganda sites whose strong political affiliations make it so they spread
40 disinformation, and as rated by the Media Bias Chart (an updated data-driven list of news sources in terms of bias and
41 truthfulness). This includes extreme-left websites like naturalnews.com. We will clarify this in revision.

─────────────── **Reviewer 3** ───────────────

43 **Examples that support the analyses?** We will investigate ways to visualize the discriminator in action, including
44 color-coding predictability like (Strobelt and Gehrmann, 2019).

45 **Experiments on human propaganda?** see response to R1 (Contribution 3, Comment 1).

46 **'counterintuitive'?** Thanks for the comment. We will change the word choice accordingly. See our response to R2 "it
47 takes a thief to catch a thief" for why we found it interesting.

48 **Minimal overlap between models?** Good point - the discriminator is initialized from generator parameters at an
49 earlier stage of training (L182) but measuring overlap between models is challenging. Still, in our paper we consider
50 discriminating from a larger/smaller generator (so no parameter overlap) and even in this case Grover is the best detector.

51 **Clarity of section 6?** Thanks for the comment; we will better disentangle these two key ideas in revision.