

1 We thank the reviewers for their insightful feedback.

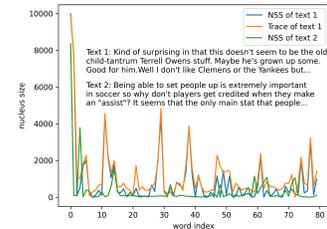
2 **Reviewer 1 Knowledge of the exact model and weights used by the victim:** Model reuse is indeed commonplace.
3 If the victim is using an off-the-shelf autocompletion app (e.g., as part of a commercial software package), the same app
4 is likely available to many other users, including potential attackers. We will clarify this in Section 4.1. Transferability
5 of fingerprints across similar (but not identical) models is an interesting topic for future research.

6 **Fully generated texts vs. intermittently accepting a completion:** For a given text, this does not affect the trace
7 measured by the attacker because in both cases the nucleus is sampled after every word (we will clarify this in the text).
8 Human editing such as deleting or rewriting might "pollute" the attacker's trace with nucleus sizes corresponding to
9 deleted subsequences. The attacker can deal with this by guessing which trace chunks originate from edits, removing
10 them, and trying to match the result. When there are many edits, this has a nontrivial computational cost, which may be
11 offset by using auxiliary information from the side channel to guide the guesses (e.g., timing, nucleus sizes, and control
12 flow of the code that operates the language model). Such methods are hard to evaluate in our lab setting because human
13 editing is indeed hard to simulate. We will expand the discussion of this limitation in Section 5.1.

14 **Conditioning NSS on longer texts:** In our experiments, the language model resets on every user post (sometimes
15 referred to as "sentence" in the submission; we will fix the terminology). We will expand on the implications of longer
16 posts in Section 3.2, providing the average post lengths for the 5 subreddits in our experiments (they range from 36 to
17 61 tokens). Average post length is inversely correlated with the fraction of "variable" sequences (variability ensures
18 that the NSS is a unique fingerprint), but even for the real-world text domains with relatively long posts (e.g., sports),
19 the fraction of variable sequences is significant (>40%). We acknowledge that for NSS generated from much longer
20 individual texts (e.g., entire documents with thousands of tokens), the fraction of variable sequences may be lower.

21 **Additional experiments:** We thank the reviewer for their suggestions. (1) We note that our mitigation makes loop
22 iteration counts independent of the nucleus sizes, therefore the attack fails completely. We will explain that it would
23 even fail to distinguish 2 known texts, let alone recognize an open-world fingerprint. (2) We will add another case study,
24 similar to the one in Section 5.3. For example, we already collected traces of 50 users' aggregated posts in the Ubuntu
25 Chat Corpus. The texts are all variable, their NSS are unique (distance $> U(N)$), traces are consistent with fingerprints
26 (noise $< d(N)$), and our attack identifies them.

27 **Reviewer 2.** As an additional illustration, we plan to add a plot similar to the one on the right, showing the NSS of two text sequences X and Y and the trace t of Y . This illustrates how the trace "matches" the fingerprint of the correct text and does not match those of other texts.



28 **Reviewer 3.** Side channels due to input-dependent branching have been studied in the computer security literature,
29 yet none of the previous work (including none of the work cited by the reviewer) identified *inputs of ML models* as
30 being vulnerable to these attacks. This observation and our analysis are nontrivial, because computation in ML models
31 involves mostly tensor arithmetics whose control flow does *not* depend on the input. Yet we demonstrate that even a
32 few lines of code in a widely used ML building block can unintentionally leak the inputs.

33 Our main contributions are to show that (1) text generation using language models creates fingerprints for many input
34 texts, and (2) these fingerprints are measurable via side channels. Uniqueness of fingerprints increases with the input
35 text's length *regardless of the specific side channel used*.

36 Our work neither requires, nor claims any innovation in side-channel techniques or defenses. We view the fact
37 that nucleus size sequences can be fingerprinted via known techniques as a strength of our paper. It means that the
38 vulnerability we identify is dangerous because it can be exploited using well-known, easily available tools and does not
39 require esoteric or new ones.

40 **Reviewer 4.** The reviewer correctly notes that loading a local shared object would thwart a Flush+Reload attack. That
41 said, other side channels such as Prime+Probe have been repeatedly shown effective for inferring loop arguments
42 (Section 2.2).

43 Users can indeed try to avoid using autocompletion features altogether, which is not always the case for cryptographic
44 primitives. Our work helps users make these decisions by identifying nucleus sampling (and ML building blocks in
45 general) as security hotspots, analogous to cryptographic primitives in that they are small functional components that
46 are (1) ubiquitously deployed and integrated into many software apps, and (2) operate on sensitive data. This motivates
47 adoption of secure coding practices for ML code.