

SPECRA: MONITOR DEGENERATIVE REPETITION IN LLM AGENTS USING RANDOMIZED FFT

Anonymous authors

Paper under double-blind review

ABSTRACT

LLM-based agents also suffer from "degenerative repetition" like chatbots, which leads to task failure and results in significant waste of computational resources and API costs until token limit is reached. Existing methods require modification of training process or customization of model deployment, and detection algorithms are brittle to approximate or structural recurrence. We therefore introduce SpecRA, a simple yet effective algorithm for detection of self-repetitions in text. Via a randomized projection from the large LLM vocabulary onto a unit-norm complex sequence, our method leverages the power of the Fast Fourier Transform (FFT) to compute the sequence's autocorrelation. Peaks in the autocorrelation function robustly reveal the underlying periodicity of the content, with tolerance to minor variations. Through an analysis of 813 repetitive samples identified from 1.13M records of anonymized agent outputs, we build a taxonomy of repetition modes in agents and show that SpecRA offers a lightweight, non-intrusive mechanism for constructing more reliable and cost-efficient LLM agents across both standard open-source model deployments and proprietary models.

1 INTRODUCTION

The promise of LLM-based agents to automate complex, multi-step tasks is transforming diverse fields, from deep research and software engineering to scientific discovery and gaming (Huang et al., 2025b; Jimenez et al., 2024; Chen et al., 2025; Hu et al., 2024). However, their practical deployment is constrained by degenerative repetition, where the model becomes trapped in a recursive loop, generating near-identical sequences repeatedly (Holtzman et al., 2019; Huang et al., 2025a).

This behavior not only leads to task failure but also incurs significant computational waste or API costs. The severity is amplified by the scale of modern agent deployments: a single long-running agent task may involve dozens to hundreds of model invocations, meaning that even a seemingly low failure rate of 1 in 10,000 can affect a substantial number of tasks at scale. Moreover, these failures can cascade, as users may repeatedly re-trigger a malfunctioning agent, inadvertently creating a denial-of-service-like scenario that degrades performance for all concurrent users.

Current approaches prove largely inadequate for addressing this challenge. Penalty-based methods can mitigate repetition but often compromise overall performance (OpenAI, 2025; Keskar et al., 2019) and frequently require custom model deployments that are impractical for many applications (Dong et al., 2025; Ginart et al., 2025). Moreover, many commercial LLM APIs either fix decoding hyperparameters such as temperature or do not expose repetition-related penalties to end users, making it difficult for most agent developers to rely solely on decoding-time mitigation.

Classical exact string matching techniques (such as n-grams and suffix trees) fail entirely when faced with minor variations, while edit-distance algorithms suffer from prohibitive polynomial runtime complexity, making them unsuitable for streaming applications where pattern lengths are unknown a priori.

We introduce **SpecRA**, a fast spectral detector for approximate repetition. Rather than analyzing the text directly, we recast the input sequence as a discrete signal of constant energy and exploit tools from signal processing. Each token is projected to a uniformly distributed random complex of unit magnitude, producing a sequence $S = \{s_1, \dots, s_N\}$. The projection makes non-repetitive text behave like white noise, so peaks provide a robust, quantitative signal for periodicity. We then

054 compute its autocorrelation of lag k via the Wiener-Khinchin theorem:

$$055 \quad R(k) = \mathcal{F}^{-1}(|\mathcal{F}(S)|^2)[k],$$

056 where \mathcal{F} denotes the discrete Fourier transform. The procedure is streaming-friendly with overall
057 $\mathcal{O}(N \log N)$ time.
058

059 Our contributions are: **(1)** Formulation of detection of degenerative repetition as approximate
060 periodicity detection over discrete streams; **(2)** An efficient approximate autorepetition detection
061 algorithm; **(3)** Theoretical bounds on false positives and detection efficacy; **(4)** A taxonomy of
062 repetition modes from 1.13M agent traces and practical guidance on setting the detection threshold.
063

064 2 RELATED WORK

065
066 Degenerative repetition is a widespread issue affecting modern LLMs, including current frontier
067 models. This phenomenon has been observed across diverse LLM-driven tasks such as code genera-
068 tion, translation, and dialogue (Dong et al., 2025; Wang et al., 2024; Xi et al., 2021). Prior research
069 has analyzed its underlying causes and mechanisms, including exposure bias and likelihood-driven
070 decoding that over-amplify frequent patterns, duplicated training data with skewed token frequencies,
071 and high-inflow dynamics that trap the generation process in self-reinforcing attractors (Holtzman
072 et al., 2019; Li et al., 2023; Fu et al., 2021; Mahaut & Franzon, 2025).

073 The most common approach involves applying penalties during the decoding process to discourage
074 repetitive behavior. Frequency and presence penalties, popularized by OpenAI-compatible APIs,
075 penalize tokens that have already appeared in the context window, while repetition penalty (Keskar
076 et al., 2019) suppresses the generation of duplicate n-grams. However, their effectiveness is highly
077 sensitive to hyperparameter tuning, and overly aggressive penalties can degrade output quality and
078 coherence.

079 More advanced decoding techniques such as contrastive search (Su et al., 2022; Sen et al., 2025),
080 information-theoretic penalties (Ginart et al., 2025), and grammar-aware penalties (Dong et al., 2025)
081 have been proposed to further reduce repetition rates. Nevertheless, these methods remain less widely
082 adopted, as they are not supported by many LLM inference providers or require custom model
083 deployments. Alternative approaches such as model editing (Wang et al., 2024) target the problem at
084 the model level but require significant computational effort and specialized expertise, making them
085 impractical for most agent developers.

086 An alternative paradigm focuses on post-hoc detection rather than prevention during generation.
087 Classical exact-match detection approaches include n-gram overlap and suffix trees, but these fail
088 entirely when faced with minor lexical variations. While edit-distance methods (Landau et al.,
089 1998) can tolerate some variations, they suffer from quadratic or higher complexity that becomes
090 prohibitive for streaming applications processing long sequences. Specialized periodicity detection
091 algorithms (Kolpakov & Kucherov, 1999; Main & Lorentz, 1984) achieve linear $\mathcal{O}(N)$ runtime but
092 are designed specifically for exact repetitions and require computationally expensive extensions to
093 handle approximate matching scenarios typical in LLM outputs.

094 Methods from bioinformatics offer alternative approaches to repetition detection (Kurtz et al., 2001).
095 K-mer based techniques, widely used in genomic sequence analysis, are also brittle to minor variations.
096 More promisingly, Fourier transforms have been successfully used in bioinformatics for detecting
097 tandem repeats (Silverman & Linsker, 1986) by mapping nucleotides to complex symbols. SpecRA
098 adapts this insight to LLM token vocabularies with randomized projection, achieving $\mathcal{O}(N \log N)$
099 complexity while maintaining robustness to lexical variation.

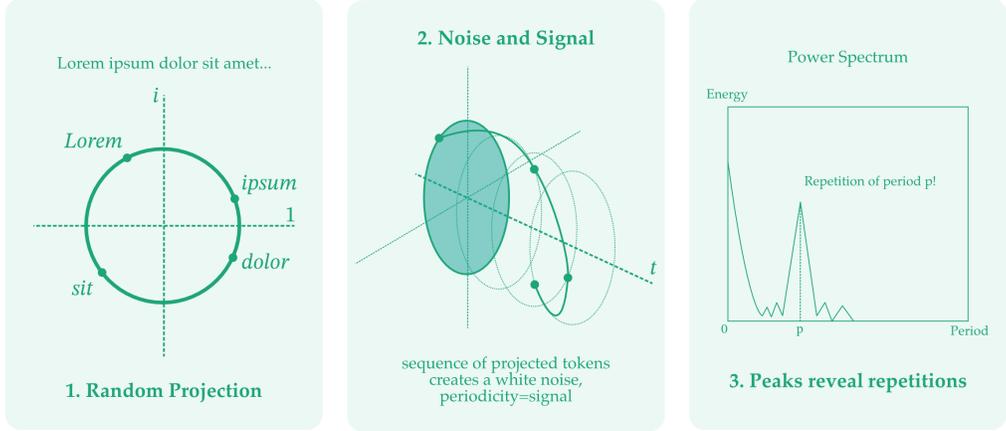
100 3 PROBLEM DEFINITION

101
102 **Task intuition.** Given a live token stream from an LLM agent, we want to raise an alarm as soon as
103 the agent falls into a "loop", namely, when its output becomes approximately periodic after allowing
104 up to εN mismatches per period.
105

106 While the excerpt below appears to show perfect repetition at first glance, the "P_P" sequence in the
107 middle disrupts the otherwise regular pattern. This exemplifies *approximate periodicity*, which we
formally define below.

162 FFT, efficiently identifying periodic patterns across multiple candidate periods; and (iii) **statistical**
 163 **detection** compares the maximum autocorrelation peak against a threshold derived from theoretical
 164 false-positive bounds.

165 The core insight is that repetitive text exhibits strong autocorrelation peaks at the repetition period,
 166 while non-repetitive text behaves like white noise with near-zero autocorrelation. The randomized
 167 projection ensures robustness to minor variations (e.g., number increments, minor spelling changes)
 168 that would confound exact string matching approaches.
 169



170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185 Figure 1: SpecRA workflow diagram showing the complete pipeline from token stream input to
 186 periodicity detection output.

187
188 **Randomized Token Projection.** Let V be the model’s vocabulary of size $|V|$. For each token
 189 $v \in V$, we draw an independent random phase $\theta_v \sim \mathcal{U}[0, 2\pi)$ and define the projection function
 190 $\phi : v \mapsto e^{i\theta_v}$. This maps each token to a point on the unit circle in the complex plane, ensuring
 191 constant signal energy $|\phi(v)| = 1$ regardless of token identity. This constraint is crucial: if tokens had
 192 unequal magnitudes, those with larger norms would disproportionately dominate the autocorrelation,
 193 biasing the repetition score based on token magnitude rather than structural periodicity. By enforcing
 194 unit magnitude, we treat all tokens symmetrically.

195 Given a token stream (x_1, x_2, \dots) , we obtain the complex sequence $S = (s_1, s_2, \dots)$ where $s_t =$
 196 $\phi(x_t)$. The key property is that identical tokens always map to the same complex number, preserving
 197 exact repetition structure, while different tokens map to (nearly) orthogonal directions. This design
 198 makes the detector robust to lexical variations: swapping "large" with "big" changes the complex
 199 representation but preserves the overall periodic structure if the substitution occurs consistently across
 200 repetitions.
 201

202 **Spectral Autocorrelation.** We apply the Wiener-Khinchin theorem to compute the autocorrelation
 203 efficiently via FFT. For a sliding window of length W , the circular autocorrelation at lag k is:

$$204 R_k = \mathcal{F}^{-1}(|\mathcal{F}(S)|^2)[k] = \sum_{t=1}^W s_t s_{t-k}^*,$$

205 where \mathcal{F} denotes the discrete Fourier transform and s_{t-k}^* is the complex conjugate with indices taken
 206 modulo W .

207 The power spectrum $|\mathcal{F}(S)|^2$ captures the frequency content of the signal, and its inverse FFT
 208 yields the autocorrelation across all lags simultaneously. For repetitive sequences with period P ,
 209 the autocorrelation R_P exhibits a large magnitude because many terms $s_t s_{t-P}^*$ align constructively.
 210 For non-repetitive sequences, these terms behave like independent random rotations, resulting in
 211 near-zero autocorrelation due to destructive interference.
 212
 213

214 **Repetition Score.** To detect repetitive patterns, we focus on the real part of the autocorrelation,
 215 which captures the alignment between tokens at different lags. Let P_{\min} and P_{\max} define the range

of plausible periods (e.g., 2 to 256 tokens). We compute the normalized repetition score:

$$S_{\text{rep}} = \frac{\max_{l=P_{\min}}^{P_{\max}} \Re(R_l)}{\Re(R_0)}$$

The denominator $\Re(R_0) = W$ normalizes by the total signal energy, ensuring the score is invariant to window size. The numerator $\Re(R_l) = \sum_{t=1}^W \Re(s_t s_{t-l}^*) = \sum_{t=1}^W \cos(\theta_t - \theta_{t-l})$ measures how well the sequence aligns with itself in the sense of cosine similarity when shifted by l positions. For perfect repetition with period P , we have $S_{\text{rep}} \approx 1$, while for random sequences, $S_{\text{rep}} \approx 0$.

We trigger a repetition alarm when $S_{\text{rep}} > S_{\text{th}}$ for some threshold $S_{\text{th}} \in (0, 1)$. The period range $[P_{\min}, P_{\max}]$ excludes trivial cases: periods smaller than $P_{\min} = 2$ are not meaningful, while periods larger than P_{\max} would require prohibitively long sequences to establish reliable patterns.

Algorithm 1: SPECRA-BATCH: Batch processing with FFT

Input: token stream (x_1, x_2, \dots) , window size W , period range $[P_{\min}, P_{\max}]$, threshold S_{th} , batch size B (e.g., $B = W$)

Output: alarm bit A_t every B timesteps

Offline: draw phases $\theta_v \sim \mathcal{U}[0, 2\pi)$ and set $\phi(v) = e^{i\theta_v}$;

Initialize: batch buffer \mathcal{B} of size B ;

Online:

for $t \leftarrow 1, 2, \dots$ **do**

$s_t \leftarrow \phi(x_t)$;

 append s_t to \mathcal{B} ;

if $|\mathcal{B}| = B$ **or end-of-stream** **then**

$F \leftarrow \mathcal{F}(\mathcal{B})$;

$R \leftarrow \mathcal{F}^{-1}(|F|^2)$;

$S_{\text{rep}} \leftarrow \max_{l=P_{\min}}^{P_{\max}} \frac{\Re(R_l)}{\Re(R_0)}$;

$A_{\text{batch}} \leftarrow \mathbf{1}\{S_{\text{rep}} > S_{\text{th}}\}$;

 emit A_{batch} ;

 clear \mathcal{B} ;

end

end

Computational Complexity. Initializing one window costs $\mathcal{O}(W \log W)$. Afterwards each batch has a complexity of $\mathcal{O}(W \log W)$ and can be amortized to $\mathcal{O}(\log W)$ per token when processed in batches, meeting the streaming constraints of Section 3.

5 EFFECTIVENESS ANALYSIS

5.1 BEHAVIOR UNDER THE NULL HYPOTHESIS

Under a non-repetitive stream, the projected tokens $s_t = e^{i\theta_t}$, with $\theta_t \sim \mathcal{U}[0, 2\pi)$, form i.i.d. isotropic noise. For any non-zero lag $l \neq 0$, the real part of the circular autocorrelation, $\Re(R_l) = \sum_{t=1}^W \Re(s_t s_{t-l}^*)$, is a sum of W i.i.d. random variables. Each term $\Re(s_t s_{t-l}^*) = \cos(\theta_t - \theta_{t-l})$ is a random variable bounded in $[-1, 1]$ with zero mean.

By applying Hoeffding’s inequality to this sum, we can bound the probability of observing a large repetition score purely by chance. Let $M = P_{\max} - P_{\min} + 1$ be the number of candidate periods. A union bound over these periods yields:

Lemma 1 (False-positive bound). *For any threshold $S_{\text{th}} > 0$, $\Pr_{\text{null}}[S_{\text{rep}} > S_{\text{th}}] \leq M \cdot \exp\left(-\frac{W S_{\text{th}}^2}{2}\right)$.*

Choosing $S_{\text{th}^*} = \sqrt{\frac{2}{W} \log(M/\delta)}$ guarantees a false-alarm rate no greater than δ .

Proof. Under the null hypothesis, the projected tokens $s_t = e^{i\theta_t}$ are i.i.d. with phases θ_t uniformly distributed in $[0, 2\pi)$. For any non-zero lag l , the term $s_t s_{t-l}^* = e^{i(\theta_t - \theta_{t-l})}$ is a random rotation. Let $Y_t = \Re(s_t s_{t-l}^*) = \cos(\theta_t - \theta_{t-l})$. Since θ_t and θ_{t-l} are independent and uniformly distributed, their difference modulo 2π is also uniform in $[0, 2\pi)$.

The variables $\{Y_t\}_{t=1}^W$ are thus i.i.d.¹ with $\mathbb{E}[Y_t] = \frac{1}{2\pi} \int_0^{2\pi} \cos(u) du = 0$ and are bounded in the interval $[-1, 1]$.

The real part of the autocorrelation is $\Re(R_l) = \sum_{t=1}^W Y_t$. The repetition score for this lag is $\frac{\Re(R_l)}{R_0} = \frac{1}{W} \sum_{t=1}^W Y_t$, since $R_0 = W$. We want to bound the probability of the event $\{\frac{\Re(R_l)}{W} > S_{th}\}$, which is equivalent to $\{\sum_{t=1}^W Y_t > WS_{th}\}$.

We apply Hoeffding's inequality. For a sum $S_W = \sum Y_t$ of W independent random variables where $Y_t \in [a_t, b_t]$, the inequality states $\Pr[S_W - \mathbb{E}[S_W] \geq \epsilon] \leq \exp(-\frac{2\epsilon^2}{\sum (b_t - a_t)^2})$. Here, $\mathbb{E}[S_W] = 0$, $\epsilon = WS_{th}$, $a_t = -1$, and $b_t = 1$, so $b_t - a_t = 2$. For a single lag, we have:

$$\Pr\left[\frac{\Re(R_l)}{W} > S_{th}\right] \leq \exp\left(-\frac{2(WS_{th})^2}{\sum_{t=1}^W (1 - (-1))^2}\right) = \exp\left(-\frac{2W^2 S_{th}^2}{4W}\right) = \exp\left(-\frac{WS_{th}^2}{2}\right).$$

The score S_{rep} is the maximum over $M = P_{max} - P_{min} + 1$ candidate lags. Applying the union bound gives the final result:

$$\Pr_{null}[S_{rep} > S_{th}] \leq \sum_{l=P_{min}}^{P_{max}} \Pr\left[\frac{\Re(R_l)}{W} > S_{th}\right] \leq M \cdot \exp\left(-\frac{WS_{th}^2}{2}\right),$$

proving the claim. \square

5.2 POWER UNDER ϵ -MISMATCH APPROXIMATE PERIODICITY

Assume a true period P . For each position t , with probability $1 - \epsilon$ we have an exact repeat $x_t = x_{t-P}$; with probability ϵ a mismatch occurs where x_t is independent of x_{t-P} (and independent across t). Under the fixed random projection ϕ above, define $X_t = \Re(s_t s_{t-P}^*) \in [-1, 1]$.

When $x_t = x_{t-P}$, $X_t = 1$; when a mismatch occurs, s_t and s_{t-P} are independent unit phases so $\mathbb{E}[X_t] = 0$. Therefore

$$\mathbb{E}[X_t] = (1 - \epsilon) \cdot 1 + \epsilon \cdot 0 = 1 - \epsilon, \quad \mathbb{E}[\Re(R_P)] = \sum_{t=1}^W \mathbb{E}[X_t] = W(1 - \epsilon).$$

Since $R_0 = W$, the normalized score for the true period is $\frac{\Re(R_P)}{R_0} = \frac{1}{W} \sum_{t=1}^W X_t$ with mean $1 - \epsilon$.

Theorem 1 (Exponential bound under ϵ -mismatch). *If $0 < S_{th} < 1 - \epsilon$ and mismatches occur independently across t , then*

$$\Pr[S_{rep} \leq S_{th}] \leq \exp\left(-\frac{W(1 - \epsilon - S_{th})^2}{2}\right).$$

Proof. Set $S_W = \sum_{t=1}^W X_t$ and $\mu' = \mathbb{E}[S_W] = W(1 - \epsilon)$. Each $X_t \in [-1, 1]$ and, by assumption, the $\{X_t\}$ are independent. The miss event $\{S_{rep} \leq S_{th}\}$ implies $\Re(R_P) \leq S_{th}W$, i.e., $S_W - \mu' \leq -(\mu' - S_{th}W)$. Hoeffding's inequality yields

$$\Pr[S_W - \mu' \leq -(\mu' - S_{th}W)] \leq \exp\left(-\frac{(\mu' - S_{th}W)^2}{2W}\right) = \exp\left(-\frac{W(1 - \epsilon - S_{th})^2}{2}\right).$$

\square

¹Although terms like Y_t and Y_{t+l} share the phase θ_t , the phase differences $(\theta_t - \theta_{t-l}) \pmod{2\pi}$ remain pairwise independent. The only statistical dependency in the full sum arises from the circular boundary condition. To ensure strict validity for Hoeffding's inequality, we may discard the l wrap-around terms, which breaks the circular dependency and makes the remaining terms strictly i.i.d., with negligible energy loss since $W \gg l$. For clarity we write W as the window length without the wrap-around terms.

6 EMPIRICAL ANALYSIS

6.1 ROBUSTNESS AGAINST SYNTHETIC NOISE

Theoretical analysis in Section 5 suggests that SpecRA can resist minor substitutions, while the effect of different noise levels and robustness against insertions and deletions remain to be investigated. We empirically validate this by generating synthetic sequences and evaluating the repetition scores.

To isolate the effect of perturbations we generate synthetic sequences in the form $T = \underbrace{(P \parallel P \parallel \dots \parallel P)}_{L/p \text{ copies}} \oplus \mathcal{N}(\varepsilon)$, where P is a base pattern of length p drawn uniformly at random

from a vocabulary of size $V = 32768$, $L = 1024$ is the total window length, and $\mathcal{N}(\varepsilon)$ applies one of substitution, deletion and insertion at rate $\varepsilon \in [0, 0.2]$. We tested $p \in \{4, 16, 64\}$ and report the median repetition score S_{rep} over 10^4 Monte-Carlo trials per setting.

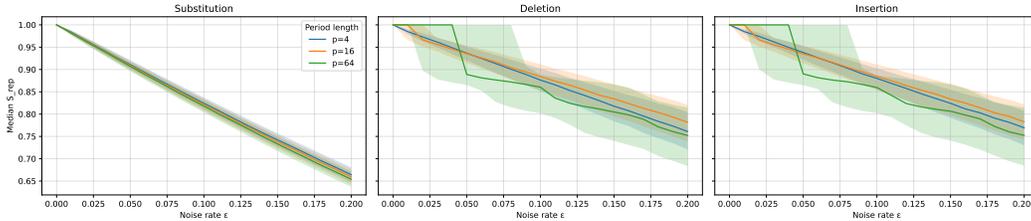


Figure 2: Median repetition score S_{rep} as a function of noise rate ε . Shaded bands denote the inter-quartile range. Insertions and deletions apply to the same randomly chosen anchor index within each period copy to simulate structural repetitions that LLMs produce.

For substitution noise, curves for different p almost overlap and S_{rep} decays almost linearly with ε for all p , confirming its insensitivity to noise level and the underlying period length under substitutions. SpecRA also tolerates indels on short patterns, as long as major structure is preserved.

To better approximate the token statistics of natural language, we further evaluate repetition scores on synthetic sequences where both the base pattern P and the corruption noise are drawn from a truncated Zipf distribution over the vocabulary.

In this setting we treat edit distance between adjacent period copies as a gold standard measure of autorepetition, and compare SpecRA against several practical heuristics. As pairwise sequence comparison methods, SimHash and MinHash are not directly applicable, so they are given access to the true period P (Oracle). NaiveFFT directly translates classical FFT-based approximate string matching by mapping n -th token to corresponding unit root $n \mapsto e^{2\pi i n/V}$. Experiment details are provided in Appendix D.

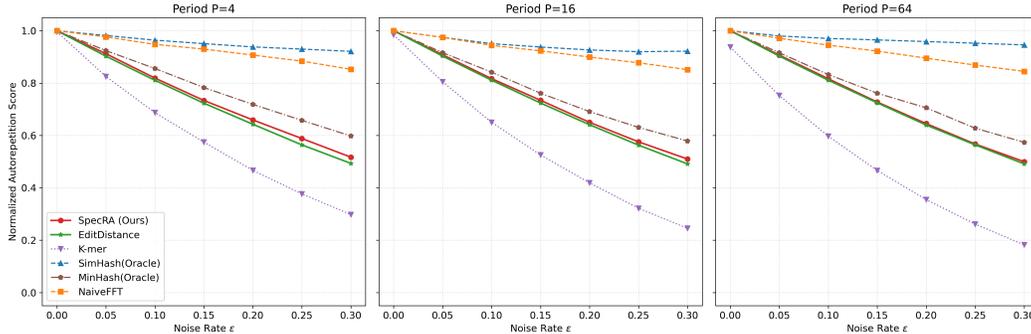


Figure 3: Normalized autorepetition score as a function of substitution noise rate ε on synthetic sequences with Zipfian token frequencies. The closer to EditDistance, the better.

378 Across all noise levels and period lengths, SpecRA almost exactly matches the EditDistance curve,
 379 indicating that our spectral self-similarity score is a highly calibrated proxy for edit distance.
 380

381 In contrast, NaiveFFT and SimHash systematically overestimate repetition and fail to decay ap-
 382 propriately with noise, maintaining high scores even at high noise levels. On the other hand, the
 383 K-mer heuristic underestimates repetition and decays too aggressively. While MinHash follows the
 384 decay trend better, it still lacks the precision of our method. SpecRA thus achieves a near-oracle
 385 approximation to edit distance while being the fastest method among all baselines (Fig. 4).
 386

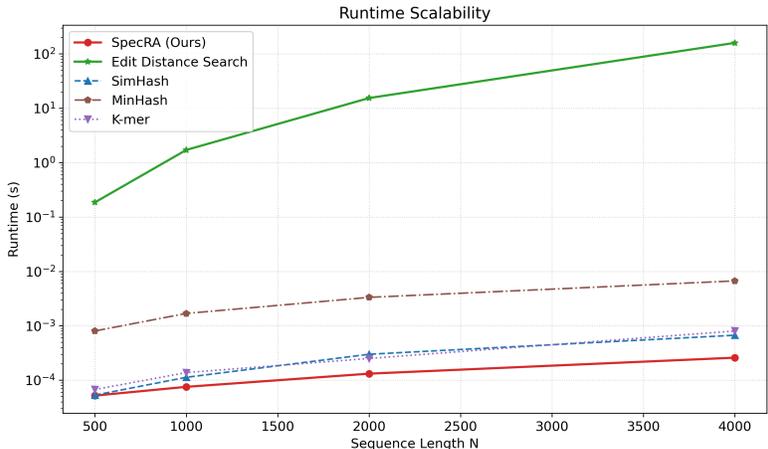


Figure 4: Runtime performance comparison between SpecRA and other baselines.

406 6.2 REPETITION FEATURES IN REAL-WORLD DATA

407 To properly set the detection threshold for SpecRA, we need to first understand the distribution of
 408 the repetition score S_{rep} in real-world data. We sampled 153,060 passages from Wikipedia (89,359
 409 passages in English and 63,701 passages in Chinese) (Foundation, 2023), 208,414 code snippets
 410 from GitHub Code (CodeParrot, 2022), and collected repetition scores from 1,133,797 rounds of
 411 LLM outputs from a general-purpose agent powered by state-of-the-art proprietary models (or their
 412 equivalents with >100B parameters) working on tasks spanning different domains involving software
 413 development, data processing, analysis/visualization, and general NLP. The S_{rep} distribution is shown
 414 in Figure 5.
 415

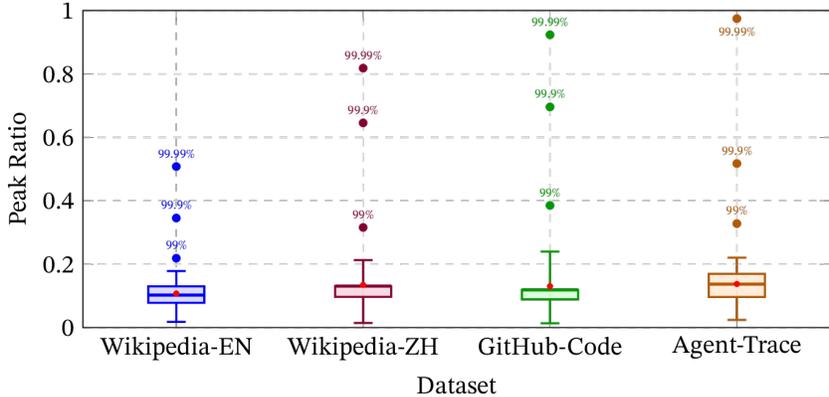


Figure 5: Repetition score distribution comparison across Wikipedia-EN, GitHub-Code, and Agent-Trace datasets. Red dots indicate mean values.

Industrial deployments that run the detector on all text streams (e.g., logging, analytics) typically demand far stricter guarantees: $FPR < 10^{-3}$ or even $< 10^{-4}$. Empirically the GitHub-Code corpus dominates the upper tail, with the 99.9th and 99.99th percentiles located at $S_{99.9} = 0.69$ and $S_{99.99} = 0.92$, respectively. We therefore recommend the following tiers:

- **Balanced:** $S_{th} = 0.69$ for general applications
- **Safe:** $S_{th} = 0.92$ for coding agents

Users with custom vocabularies or window sizes can re-estimate the quantile with a single offline scan and plug them into the same decision rule.

6.3 AGENT REPETITION TAXONOMY

SpecRA flagged 813 suspicious repetition samples out of 1,133,797 agent traces (0.071%), using a threshold of $S_{th} = 0.69$. We excluded 264 samples that were too short to classify reliably, likely due to incomplete generation from network errors or early termination. The remaining 549 samples were classified into four distinct repetition categories (see Appendix E for annotation details):

Structural repetition. Systematic iteration over semantically related content patterns, where agents generate sequences of structurally similar elements with incremental variation. Examples include enumeration of chemical elements in periodic table order, systematic generation of numbered function definitions, and iterative construction of similar data structures. This pattern reflects the model’s attempt to complete structured tasks through template-based generation.

Syntactic degradation. Purely syntactic repetition without semantic coherence, where models generate identical token sequences or character patterns with no underlying logical structure. This includes infinite repetition of single characters, alphabetical and numerical sequences (e.g., ",3,3,3..."), representing complete semantic breakdown in the generation process.

Binary data generation. A notable repetition pattern observed in agents attempting to emit binary data (105 / 549; **19.1%**), manifesting in three distinct sub-patterns: **(1) Multimedia encoding loops:** cyclical repetition of base64-encoded character sequences representing multimedia content (images, audio, video files) or structured documents; **(2) URI malformation cycles:** iterative generation of malformed data URIs or embedded content for visualization services (e.g., mermaid.ink, plantuml.com), often producing corrupted markup with repeated URI fragments; and **(3) Direct binary emission:** direct attempts to emit binary file headers and control characters (e.g., ZIP/Office file magic numbers "PK", PNG signatures), interspersed with repeated `\uFFFE` patterns.

Legitimate repetition. Cases where structurally necessary repetition is misclassified as degenerative, including large-scale data serialization (JSON arrays, CSV records), ASCII art containing repeated patterns, systematic progress tracking with templated status reports, and algorithmic output requiring repetitive formatting patterns that serve a functional purpose.

These categories constitute 46.26%, 21.68%, 19.13%, 12.93% of the flagged cases, respectively.

Table 1: Breakdown of the 549 repetitive turns by error category ($S_{th} = 0.69$).

Category	# samples	Share
Structural repetition	254	46.26%
Syntactic degradation	119	21.68%
Binary total	105	19.13%
Multimedia encoding loops	55	10.02%
URI malformation cycles	31	5.65%
Binary header emission	19	3.47%
Legitimate repetition	71	12.93%
Total	549	100%

486 7 DISCUSSION

487 7.1 USAGE NOTE ON PRACTICAL DEPLOYMENT

488 For agent developers who rely on LLM API providers without access to inference infrastructure,
489 original tokens may not be available. In such cases, SpecRA can still be effectively applied to
490 character-level streams to detect repetitive failures. However, when applied to smaller vocabularies
491 (e.g., ASCII-only streams), performance may degrade due to increased collision probability in the
492 hash space.

493 Phase values θ_v are sampled i.i.d. from a continuous distribution. A potential issue arises when two
494 distinct tokens are assigned similar phases, causing mismatches to contribute $\cos(\Delta) \approx 1$ to the
495 autocorrelation, mimicking matches.

496 The risk depends on the interplay between random phase collisions and token co-occurrence statistics.
497 Large vocabularies increase the number of potential collision pairs ($\binom{|V|}{2}$), while small vocabularies
498 concentrate statistical weight on fewer pairs. If frequently co-occurring tokens happen to receive
499 similar phases, the distortion effect is amplified.

500 Our proposed mitigation using K independent projections is highly effective. For a mismatch to
501 consistently distort the signal, it must be a near-collision across all K mappings—a vanishingly
502 unlikely event that ensures detector reliability regardless of vocabulary size or input statistics.

503 Additionally, legitimate repetition may occasionally be misclassified as degenerative. In such cases,
504 LLMs can serve as a secondary validation mechanism.

505 7.2 LIMITATIONS

506 While SpecRA demonstrates effectiveness in detecting repetitive failures, it has several inherent
507 limitations. First, although SpecRA excels at identifying simple structural repetitions, it may struggle
508 with more complex patterns that require deeper contextual understanding. For instance, it may fail
509 to detect repetitive failures in code generation tasks that produce semantically similar code snippets
510 with varying lengths, as the length variations introduce phase shifts across repetitive blocks (Dong
511 et al., 2025).

512 Within the (ε, p) -approximate periodicity framework of Section 3, we did not observe false negatives
513 on our agent-trace dataset, consistent with the exponential miss-detection bound in Theorem 1.
514 However, more semantic forms of repetition that violate the fixed-period assumption, such as
515 grammatically repetitive code blocks whose lengths grow across iterations, fall outside this formal
516 setting and remain an important direction for future extensions of SpecRA.

517 As discussed in Section 7.1, SpecRA’s performance is sensitive to vocabulary size. Smaller vocabu-
518 laries increase collision probability, potentially leading to performance degradation. While we have
519 proposed mitigation strategies using multiple independent projections, their effectiveness requires
520 further empirical evaluation.

521 8 CONCLUSION AND FUTURE WORK

522 We framed degenerative repetition in LLM agents as an approximate periodicity detection problem and
523 introduced **SpecRA**, which combines randomized phase projection with FFT-based autocorrelation
524 analysis. Our method achieves $\mathcal{O}(W \log W)$ processing complexity with $\mathcal{O}(\log W)$ amortized time
525 per token, while providing provable bounds on both false-alarm and miss-detection probabilities.
526 Extensive experiments across public corpora and real agent traces demonstrate that SpecRA offers a
527 lightweight, non-intrusive solution for building more reliable and cost-efficient LLM agents.

528 Future work can extend this research in three immediate avenues: **(1) Inference-time integration**, by
529 incorporating SpecRA scores as decoding penalties to steer models away from repetitive attractors; **(2)**
530 **Cross-modal generalization**, by adapting the spectral approach to detect cyclic artifacts in vocoder
531 waveforms, embedding streams, or tool-use trajectories; and **(3) An enhanced signal-processing**
532 **toolkit**, exploring techniques like wavelet coherence or adaptive filtering to build a comprehensive
533 suite of guards for trustworthy AI.

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

REPRODUCIBILITY STATEMENT

To ensure reproducibility, we provide a reference Python implementation in Appendix C. Our experiments use synthetic data, generated as described in Section 6.1, and public corpora (Wikipedia, GitHub Code) detailed in Section 6.2. The full experimental setup and parameters are specified in Section 6.

REFERENCES

- 594
595
596 Qiguang Chen, Mingda Yang, Libo Qin, Jinhao Liu, Zheng Yan, Jiannan Guan, Dengyun Peng, Yiyang
597 Ji, Hanjing Li, Mengkang Hu, Yimeng Zhang, Yihao Liang, Yuhang Zhou, Jiaqi Wang, Zhi Chen,
598 and Wanxiang Che. Ai4research: A survey of artificial intelligence for scientific research, 2025.
599 URL <https://arxiv.org/abs/2507.01903>.
- 600 CodeParrot. Github code dataset. <https://huggingface.co/datasets/codeparrot/github-code>, 2022.
601
- 602 Yihong Dong, Yuchen Liu, Xue Jiang, Zhi Jin, and Ge Li. Rethinking repetition problems of llms in
603 code generation, 2025. URL <https://arxiv.org/abs/2505.10402>.
- 604 Wikimedia Foundation. Wikimedia downloads. <https://dumps.wikimedia.org>, 2023.
605
- 606 Zihao Fu, Wai Lam, Anthony Man-Cho So, and Bei Shi. A theoretical analysis of the repetition
607 problem in text generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):
608 12848–12856, May 2021. doi: 10.1609/aaai.v35i14.17520. URL [https://ojs.aaai.org/
609 index.php/AAAI/article/view/17520](https://ojs.aaai.org/index.php/AAAI/article/view/17520).
- 610 Antonio A. Ginart, Naveen Kodali, Jason Lee, Caiming Xiong, Silvio Savarese, and John R. Emmons.
611 Lz penalty: An information-theoretic repetition penalty for autoregressive language models, 2025.
612 URL <https://arxiv.org/abs/2504.20131>.
- 613 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text
614 degeneration. *arXiv preprint arXiv:1904.09751*, 2019. URL [https://arxiv.org/abs/
615 1904.09751](https://arxiv.org/abs/1904.09751).
- 616 Sihao Hu, Tiansheng Huang, Gaowen Liu, Ramana Rao Kompella, Fatih Ilhan, Selim Furkan Tekin,
617 Yichang Xu, Zachary Yahn, and Ling Liu. A survey on large language model-based game agents,
618 2024. URL <https://arxiv.org/abs/2404.02039>.
- 619 Donghao Huang, Thanh-Son Nguyen, Fiona Llausvia, and Zhaoxia Wang. RAP: A metric for
620 balancing repetition and performance in open-source large language models. In Luis Chiruzzo,
621 Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas
622 Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume
623 1: Long Papers)*, pp. 1479–1496, Albuquerque, New Mexico, April 2025a. Association for
624 Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.69. URL
625 <https://aclanthology.org/2025.naacl-long.69/>.
626
- 627 Yuxuan Huang, Yihang Chen, Haozheng Zhang, Kang Li, Huichi Zhou, Meng Fang, Linyi Yang,
628 Xiaoguang Li, Lifeng Shang, Songcen Xu, Jianye Hao, Kun Shao, and Jun Wang. Deep research
629 agents: A systematic examination and roadmap, 2025b. URL [https://arxiv.org/abs/
630 2506.18096](https://arxiv.org/abs/2506.18096).
- 631 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
632 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL
633 <https://arxiv.org/abs/2310.06770>.
634
- 635 Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher.
636 Ctrl: A conditional transformer language model for controllable generation, 2019. URL [https://arxiv.org/abs/
637 1909.05858](https://arxiv.org/abs/1909.05858).
- 638 Roman Kolpakov and Gregory Kucherov. Finding maximal repetitions in a word in linear time. In
639 *Proceedings of the 40th Annual Symposium on Foundations of Computer Science, FOCS '99*, pp.
640 596, USA, 1999. IEEE Computer Society. ISBN 0769504094.
- 641 Stefan Kurtz, Jomuna V. Choudhuri, Enno Ohlebusch, Chris Schleiermacher, Jens Stoye, and Robert
642 Giegerich. Reputer: the manifold applications of repeat analysis on a genomic scale. *Nucleic Acids
643 Research*, 29(22):4633–4642, 11 2001. ISSN 0305-1048. doi: 10.1093/nar/29.22.4633. URL
644 <https://doi.org/10.1093/nar/29.22.4633>.
645
- 646 Gad M. Landau, Eugene W. Myers, and Jeanette P. Schmidt. Incremental string comparison. *SIAM
647 Journal on Computing*, 27(2):557–582, 1998. doi: 10.1137/S0097539794264810. URL <https://doi.org/10.1137/S0097539794264810>.

648 Huayang Li, Tian Lan, Zihao Fu, Deng Cai, Lema Liu, Nigel Collier, Taro Watanabe, and Yixuan
649 Su. Repetition in repetition out: Towards understanding neural text degeneration from the data
650 perspective. *ArXiv*, abs/2310.10226, 2023. URL [https://api.semanticscholar.org/
651 CorpusID:264146506](https://api.semanticscholar.org/CorpusID:264146506).

652 Matéo Mahaut and Francesca Franzon. Repetitions are not all alike: distinct mechanisms sustain
653 repetition in language models, 2025. URL <https://arxiv.org/abs/2504.01100>.

654 Michael G Main and Richard J Lorentz. An $o(n \log n)$ algorithm for finding all repetitions in
655 a string. *Journal of Algorithms*, 5(3):422–432, 1984. ISSN 0196-6774. doi: [https://doi.org/
656 10.1016/0196-6774\(84\)90021-X](https://doi.org/10.1016/0196-6774(84)90021-X). URL [https://www.sciencedirect.com/science/
657 article/pii/019667748490021X](https://www.sciencedirect.com/science/article/pii/019667748490021X).

658 OpenAI. API reference - OpenAI API. [https://platform.openai.com/docs/
659 api-reference](https://platform.openai.com/docs/api-reference), 2025. Archived at: <https://archive.li/IIXy2>.

660 Jaydip Sen, Rohit Pandey, and Hetvi Waghela. Context-enhanced contrastive search for improved
661 llm text generation, 2025. URL <https://arxiv.org/abs/2504.21020>.

662 B.D. Silverman and R. Linsker. A measure of dna periodicity. *Journal of Theo-
663 retical Biology*, 118(3):295–300, 1986. ISSN 0022-5193. doi: [https://doi.org/10.
664 1016/S0022-5193\(86\)80060-1](https://doi.org/10.1016/S0022-5193(86)80060-1). URL [https://www.sciencedirect.com/science/
665 article/pii/S0022519386800601](https://www.sciencedirect.com/science/article/pii/S0022519386800601).

666 Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. A contrastive
667 framework for neural text generation, 2022. URL <https://arxiv.org/abs/2202.06417>.

668 Weichuan Wang, Zhaoyi Li, Defu Lian, Chen Ma, Linqi Song, and Ying Wei. Mitigating the language
669 mismatch and repetition issues in llm-based machine translation via model editing, 2024. URL
670 <https://arxiv.org/abs/2410.07054>.

671 Yadong Xi, Jiashu Pu, and Xiaoxi Mao. Taming repetition in dialogue generation, 2021. URL
672 <https://arxiv.org/abs/2112.08657>.

673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A PRIVACY AND ETHICS CONSIDERATIONS

Our analysis of agent logs was conducted under strict privacy safeguards and ethical guidelines. All LLM outputs were anonymized and de-identified prior to analysis, with access restricted exclusively to patterns flagged as anomalous by our detection algorithm. No personally identifiable information, proprietary content, or sensitive user-generated data was examined during the analysis process.

SpecRA provides inherent privacy advantages: it operates on statistical properties of token sequences rather than semantic content, enabling detection of repetitive failures without requiring persistent storage or detailed inspection of user data. This design preserves the confidentiality of user-agent interactions while delivering essential protection against computational waste and system instability. Furthermore, the randomized projection mechanism ensures that even if projection parameters were compromised, recovering original token sequences would remain impractical.

B USE OF LLMs

LLMs were utilized during various stages of this paper’s development. Specifically, LLMs assisted with: (a) translating our original derivations and algorithmic ideas into initial \LaTeX formulations and suggesting more standard notation, (b) language polishing and stylistic refinement, and (c) comprehensive proofreading. The core ideas, problem formulation, and all mathematical proofs, including Lemma 1 and Theorem 1, were developed and verified by the authors, who take full responsibility for their validity. We acknowledge that [Anonymous frontier model] helped us identify a flaw in an early draft of the proof of Theorem 1, we subsequently corrected and validated the final version ourselves.

Additionally, LLMs were employed to assist with the initial classification and anonymization of potentially suspicious repetitive cases prior to manual analysis. All automated labels were subsequently verified and when necessary, corrected by the authors to ensure accuracy and consistency.

C REFERENCE IMPLEMENTATION

This section provides a minimal Python implementation of the SpecRA algorithm. The code demonstrates the core spectral analysis pipeline described in Section 4, implementing offline batch processing suitable for research and prototyping.

```

import numpy as np
from typing import List, Tuple

# Build a random unit-magnitude complex phase mapping for each token id.
# In practice, VOCAB_SIZE should be at least the vocabulary size of
# all models being used.
VOCAB_SIZE = 32768
SPECRA_MAP = np.exp(1j * 2 * np.pi * np.random.rand(VOCAB_SIZE))

# SpecRA detects approximate repetition in a token sequence
# using spectral analysis.
# Parameters:
# - token_ids: input token sequence as integer ids
# - threshold: detection threshold for normalized repetition score, range [0, 1]
# Returns:
# - peak_ratio: maximum normalized autocorrelation value
# - detected: True if repetition detected above threshold
def specra(token_ids: List[int], threshold: float) -> Tuple[float, bool]:
    # Step 1: Random projection
    seq = np.array([SPECRA_MAP[t] for t in token_ids], dtype=np.complex128)
    n = len(seq)

    # Step 2: FFT-based autocorrelation via Wiener-Khinchin theorem
    # Forward FFT: F = FFT(seq)
    coeffs = np.fft.fft(seq)
    # Power spectrum: power = |F|^2

```

```

756     power = np.abs(coeffs) ** 2
757     # Inverse FFT: autocorr = IFFT(|F|^2)
758     autocorr = np.fft.ifft(power)
759
760     # Step 3: Compute normalized repetition score
761     r0 = float(np.real(autocorr[0]))
762     if r0 <= 0.0:
763         return 0.0, False
764     p_min, p_max = 1, n // 2 + 1
765     if p_max <= p_min:
766         return 0.0, False
767
768     real_corr = np.real(autocorr[p_min:p_max])
769     peak = float(np.max(real_corr)) if real_corr.size > 0 else 0.0
770     return peak / r0, peak > threshold * r0

```

771 D SYNTHETIC NOISE EXPERIMENT DETAILS

772 **Data Generation.** We generate synthetic repetitive sequences of fixed window length $W = 1,024$ with underlying periods $p \in \{4, 16, 64\}$ and vocabulary size of $V = 32,768$.

773 For the Uniform Noise setting, base patterns are drawn uniformly from the vocabulary, and noise replaces tokens with new ones drawn uniformly at random. For the Zipfian Noise setting, to better approximate natural language statistics, tokens are drawn from a Zipf distribution with exponent $a = 1.05$. Substitutions also draw replacement tokens from the same distribution, ensuring the noise follows the underlying term frequency distribution and maintaining the "long tail" property of language. For Indel Noise, we apply insertions and deletions at $\varepsilon \cdot N$ randomly chosen positions, fixing the mutation positions across all periods.

782 **Baselines Implementation.** All baselines are implemented in Python. Scores are normalized to $[0, 1]$, where 1 corresponds to perfect repetition and 0 indicates that every token is unique. We compare SpecRA against:

- 786 • **EditDistance:** We compute the average edit distance between adjacent non-overlapping blocks of length p in the sequence. We report $1 - (\text{distance}/p)$ as the similarity score. This baseline represents an idealized detector that knows the period p but not the ground truth pattern.
- 787 • **EditDistance Search:** We search for the period p in the sequence using a brute-force approach with period range $[2, N/10]$.
- 788 • **NaiveFFT:** A direct spectral method that maps token ID k to the k -th root of unity $e^{2\pi ik/V}$ before applying FFT. This tests whether randomized projection is necessary compared to deterministic mapping.
- 789 • **K-mer:** We use $K = 4$. The score is calculated as $1 - \frac{\text{unique } K\text{-mers}}{\text{total } K\text{-mers}}$, measuring the redundancy of local sub-patterns.
- 790 • **SimHash:** We compute 64-bit SimHash fingerprints for each period-sized block. In practice, SimHash is usually used together with stop-word removal or TF-IDF, while in a streaming setting, TF-IDF is not feasible. For Zipfian data, we remove the top-10 most frequent tokens as stop-words to improve robustness, as high-frequency tokens dominate the hash signatures. We compare these fingerprints against the oracle pattern's fingerprint, giving this baseline a significant advantage.
- 791 • **MinHash:** We generate signatures using 16 permutations and compute Jaccard similarity against the oracle pattern's signature.

800 Runtime benchmarks were conducted on a standard CPU (Apple M1). We measured the average wall-clock time over 10 trials for sequence lengths $N \in \{500, 1000, 2000, 4000\}$. SpecRA's complexity is dominated by the FFT operation ($\mathcal{O}(N \log N)$), whereas edit distance search scales poorly ($\mathcal{O}(N^3)$).

E AGENT REPETITION TAXONOMY DETAILS

The 549 repetitive samples were annotated by multiple LLMs with major voting, and then carefully verified and when necessary, corrected by the authors.

The 264 examples were excluded because they were too short or incomplete to be classified reliably (e.g., due to network failure, user cancellation or safety filters). While they manifest repetitive features, it is uncertain whether the model would have entered an infinite loop (or equivalently whether the model will exit eventually after a some more periods). This may slightly shift the distribution, but likely not significantly.

During annotation, we allowed for an "Other" label for edge cases that are hard to classify, and manually labeled them after thorough inspection. The process of classification can be described as follows:

- (i) Does the pattern carry coherent meaning to a human reader?
 - Yes → (ii)
 - No → (iii)
 - Uncertain → (Other)
- (ii) Given that the content is meaningful, is the repetition legitimate for the user’s task (e.g., translating a repetitive table)?
 - Yes → (Legitimate repetition)
 - No → (Structural repetition)
 - Uncertain → (Other)
- (iii) Given that the content is not meaningful to humans, does it still correspond to machine-interpretable binary formats (e.g., base64-encoded multimedia content, binary file headers)?
 - Yes → (Binary data generation)
 - No → (Syntactic degradation)
 - Uncertain → (Other)