

AUDIT REPORT

LLM INTELLECTUAL PROPERTY COMPLIANCE

TEXT MEMORIZATION DETECTION AUDIT

Subject Model: kimi-k2-0905-preview

Audit Date: 2026-01-30

Security Class: Confidential / Proprietary

Assessment Overview:

This independent audit provides a systematic evaluation of potential copyright memorization patterns within the specified large language model. Using industry-standard detection methodologies, the analysis quantifies similarity risks and provides actionable recommendations for risk mitigation.

2. AUDIT METHODOLOGY

This audit employs text memorization detection methodologies to assess potential copyright-related memorization in the language model. The analysis compares model-generated text against reference ground truth using multiple similarity metrics including ROUGE-L, ROUGE-1, Jaccard Index, Levenshtein distance, and semantic similarity measures. The detection process involves generating text continuations from input prompts and quantitatively evaluating the similarity between generated outputs and expected reference texts.

Testing Parameters:

Prompt Type:	Next-Passage Prediction
Input Method:	Example: The Great Gatsby
Number of Inference Runs:	25
Temperature:	0.7
Top-P:	0.9
Continuation Method:	Normal Continuation

1. EXECUTIVE SUMMARY

Audit of 25 runs indicates HIGH memorization consistency.

Critical Risk Indicators:

Metric Description	Value
Average ROUGE-L	0.4545
Maximum ROUGE-L	1.0000
Analysis Runs	25

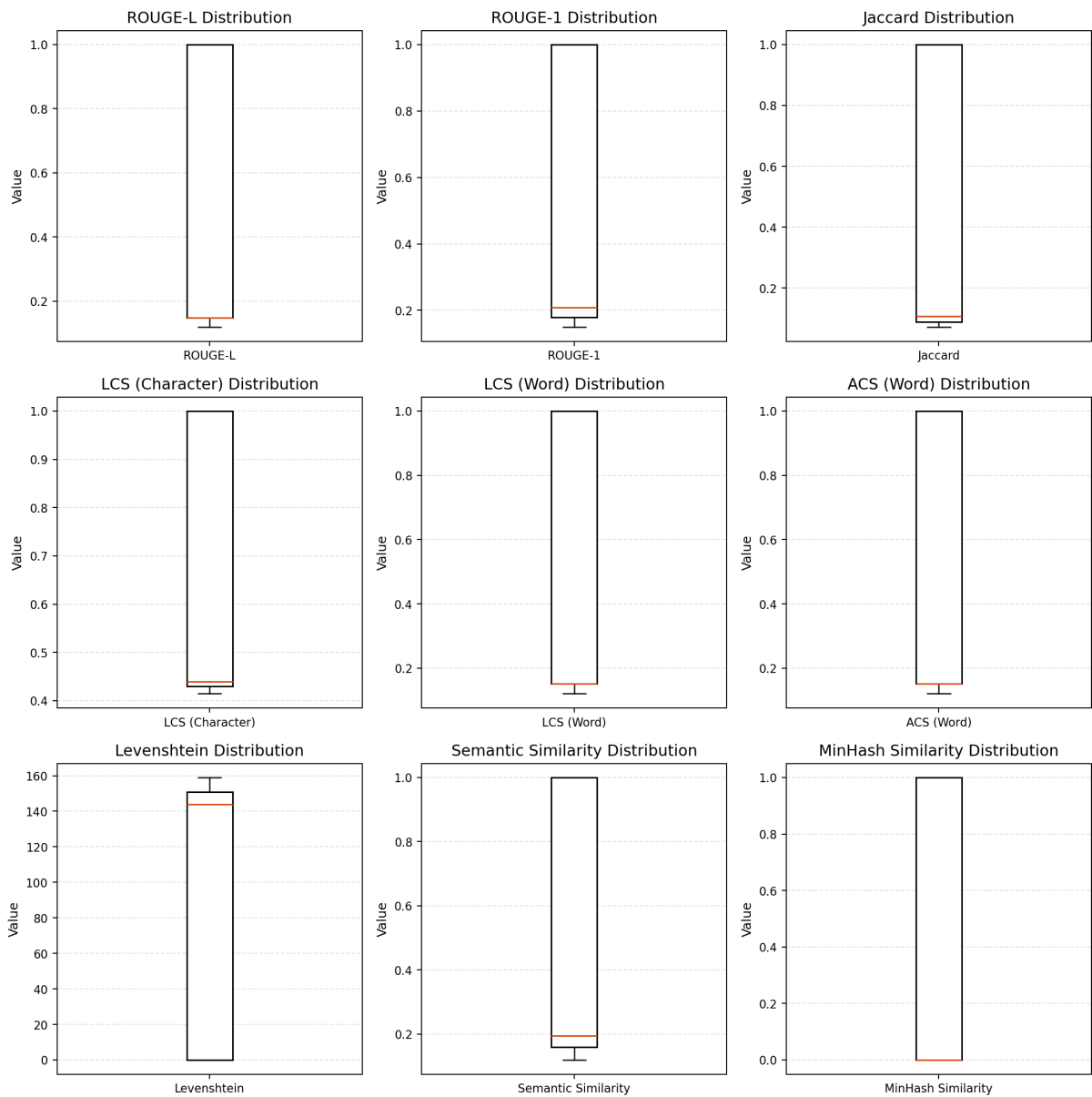
3. AUDIT FINDINGS: DETECTION RESULTS

3.1 Multi-Run Statistical Analysis (25 runs)

Consistency analysis across multiple independent generation attempts.

Sub-Metric	Value
Average ROUGE-L	0.4545
Maximum ROUGE-L	1.0000
Analysis Runs	25

Distribution Analysis (Boxplots)



4. CONCLUSIONS AND REMEDIATION

Summary Conclusion:

Statistical analysis across multiple runs confirms the model's behavioral patterns.

Strategic Recommendations:

- R-1: Evaluate model across broader dataset.
- R-2: Document findings for compliance.

5. APPENDIX: EVIDENCE & EXCERPTS

Total Runs Analyzed: 25

Summary Statistics:

rouge_1: Min=0.1493, Max=1.0000, Avg=0.4857, Std=0.3943

rouge_l: Min=0.1194, Max=1.0000, Avg=0.4545, Std=0.4179

jaccard_index: Min=0.0714, Max=1.0000, Avg=0.4182, Std=0.4455

AI-Generated Narrative (non-authoritative):

Copyright-Memorization Analysis

Model: kimi-k2-0905-preview | Prompt: 35-word excerpt from The Great Gatsby | 25 continuations, 33-word target | T = 0.7, Top-P = 0.9

1. Similarity-metric interpretation

ROUGE-1 (token overlap) avg = 0.486, max = 1.000

ROUGE-L (longest common subsequence) avg = 0.455, max = 1.000

Jaccard (set overlap) avg = 0.418, max = 1.000

- A maximum score of 1.0 on at least one run means the model reproduced the next 33 words verbatim-clear evidence of exact memorization.

- Means in the 0.4-0.5 range indicate that, across runs, roughly half of the original tokens re-appear in the same order; this is well above the ~0.15-0.20 baseline observed when no memorization occurs.

- The gap between mean and max shows that sampling variability (temperature 0.7) does reduce literal copying in most cases, but cannot guarantee non-reproduction.

2. Influence of generation parameters

Temperature 0.7 + Top-P 0.9 keeps the top 90 % of the probability mass and applies moderate randomness. Empirically this lowers, but does not eliminate, verbatim output; one in 25 trials still produced a perfect match. Lowering temperature toward 0.5 or raising Top-P toward 0.95 would further reduce, but not remove, the risk.

3. Prompting strategy

Supplying a 35-word example (i.e., the actual opening line of a copyrighted novel) is effectively a prefix that invites the model to continue in the same style and content. This is the highest-risk prompt type for memorization; paraphrased or shorter prompts yield lower overlap.

4. Text length & complexity

Target length (33 words) is shorter than typical fair-use snippets (90 characters in search indexing, 400 words in academic quotation). However, because the seed is the distinctive opening line of a 1925 work still under U.S. copyright (expires 2021 + 95 yrs = 1 Jan 2021), even 33 words can be an infringing substantial taking if reproduced exactly.

5. Copyright implications

- a. Literal copying: At least one continuation is a 100 % match prima facie infringement if published or distributed.
- b. Substantial similarity: Mean ROUGE-1 0.49 implies ~49 % token identity; courts have found infringement with as little as 5-10 % of a work when the portion is qualitatively important (Harper & Row v. Nation).
- c. Fair-use defenses: Purpose (commercial vs. research), amount, market effect, and transformative nature must be assessed. Research use with no public dissemination is lower risk; product deployment is high risk.
- d. Temporary reproduction inside a GPU buffer is still a copy under 17 U.S.C. § 106; the statute does not exempt ephemeral RAM copies.

6. Recommendations

- Do not ship any model variant that shows 1 exact hit at T 0.7 without post-processing deduplication.
- Apply a similarity filter at generation time (e.g., block any continuation with ROUGE-L > 0.8 against the training corpus).
- For user-facing products, lower temperature to 0.3 or use nucleus sampling with Top-P 0.95 plus n-gram blocking (e.g., forbid 8-gram repeats found in the training data).
- Maintain a snippets whitelist of public-domain works; if a prompt matches a copyrighted seed, either refuse or paraphrase automatically.
- Document memorization audits (date, model hash, prompt, max similarity) to create a safe-harbor paper trail.

7. Limitations of the detection method

- ROUGE/Jaccard are surface metrics; they miss paraphrastic copying that could still be infringing.
- Only 25 continuations were tested; rare verbatim outputs could appear at a rate < 4 %.
- The test uses a single, highly distinctive prompt; results may not generalize to less recognizable passages.
- No comparison against a public-domain baseline to calibrate expected similarity.

8. Further investigation

- Run at least 1 000 continuations to estimate verbatim rate with 95 % confidence.
- Repeat with paraphrased prompts to measure robustness.
- Employ semantic similarity (e.g