

Auto-Stega: An Agent-Driven System for Lifelong Strategy Evolution in LLM-Based Text Steganography

Anonymous ACL submission

Abstract

With the rapid progress of LLMs, high quality generative text has become widely available as a cover for text steganography. However, prevailing methods rely on hand-crafted or pre-specified strategies and struggle to balance efficiency, imperceptibility, and security, particularly at high embedding rates. Accordingly, we propose Auto-Stega, an agent-driven self-evolving framework that is the first to realize self-evolving steganographic strategies by automatically discovering, composing, and adapting strategies at inference time; the framework operates as a closed loop of generating, evaluating, summarizing, and updating that continually curates a structured strategy library and adapts across corpora, styles, and task constraints. A decoding LLM recovers the information under the shared strategy. To handle high embedding rates, we introduce PC-DNTE, a plug-and-play algorithm that keeps sampling close to the base model’s conditional distribution at high embedding rates, preserving imperceptibility while enhancing security. Experimental results demonstrate that at higher embedding rates, Auto-Stega achieves superior performance with a 42.2% reduction in normalized perplexity and a 1.2% improvement in anti-steganalysis performance over SOTA methods.

1 Introduction

Steganography is an information hiding technique that embeds secret information in cover media to enable covert communication (Liao et al., 2025; Meng et al., 2025; Zhang et al., 2025a). The fundamental objectives of steganographic systems are characterized by three key metrics: the embedding rate, which quantifies the capacity of secret information that can be hidden; perceptual concealment, typically represented by the quality of the media after information embedding (e.g., text quality or degree of image modification); and anti-steganalysis

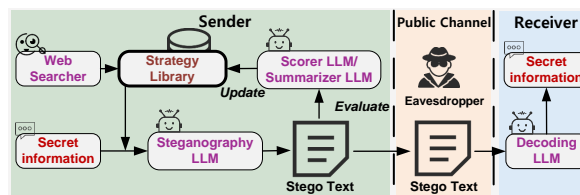


Figure 1: The pipeline of Auto-Stega.

performance (Hu and Wang, 2025; Zhou et al., 2025b), which measures the capability to resist machine detection. These correspond respectively to core requirements of efficiency, imperceptibility, and security.

By cover type, steganography can be categorized into text steganography (Wang et al., 2025a,b; Wu et al., 2024), image steganography (Zhou et al., 2025c; Cheng et al., 2025), audio steganography (Li et al., 2025; Zhang et al., 2025b), and video steganography (Meng et al., 2025). Among these, text steganography has emerged as an active research area in information security, driven by the prevalence of text-based social communication.

Early text steganography was predominantly based on manually designed content modifications, such as synonym substitution (Li et al., 2019), and spelling conversion (Shirali-Shahreza, 2008). Although these rule-based schemes can remain semantically unobtrusive, they usually provide a low embedding rate and induce a distributional change from natural text (Zhang et al., 2021), which compromises their resistance to machine detection and leads to inadequate anti-steganalysis performance.

Recent advances in large language models (LLMs) have made high-quality generative text broadly accessible and controllable, creating new opportunities for generative text steganography. Generative steganography (Zhou et al., 2025c; Peng et al., 2024) is a technique that embeds information by guiding probabilistic generative models during content synthesis. In text generation,

075 it typically embeds bits by sampling steering to- 126
076 kens within the model’s conditional probability 127
077 distribution. Zhang et al. (Zhang et al., 2021) 128
078 achieve information-theoretic alignment via equal 129
079 mass dynamic binning but rely on numerically frag- 130
080 ile reconstruction; Ding et al. (Ding et al., 2023) 131
081 strengthen security with distribution copies and 132
082 seed-based decoding at the cost of synchroniza- 133
083 tion and latency; Wu et al. (Wu et al., 2024) re- 134
084 duce shift using SVO and sentiment prompting, 135
085 yet remain template and keyword dependent under 136
086 domain shift. Taken together, these methods still 137
087 rely on hand-crafted or pre-defined strategies and 138
088 struggle to balance efficiency, imperceptibility, and 139
089 security, especially under high embedding rates - a 140
090 well-recognized challenge in the the steganography 141
091 domain, known as the “trilemma.” 142

092 To fundamentally address this trilemma and in- 143
093 spired by cognitive science showing that creative 144
094 solutions often emerge from coherently combining 145
095 seemingly unrelated knowledge (Koestler, 1964; 146
096 Benedek et al., 2012; Lee and Chung, 2024), we 147
097 propose Auto-Stega, an agent-driven, self-evolving 148
098 framework (as shown in Fig. 1) that is the first to re- 149
099 alize self-evolving steganographic strategies by au- 150
100 tomatically discovering, composing, and adapting 151
101 strategies at inference time. The framework unifies 152
102 strategy retrieval, stego text generation, calibrated 153
103 multi-objective evaluation, and strategy summa- 154
104 rization with library update in a closed-loop, con- 155
105 tinually maintaining a searchable strategy library 156
106 and coordinating efficiency, imperceptibility, and 157
107 security across tasks. For high embedding rates, 158
108 we also introduce PC-DNTE, an optional algorithm 159
109 of our own design that integrates without chang- 160
110 ing the overall operation. The contributions are as 161
111 follows: 162

- 112 • **The first text steganography framework for** 163
113 **self-evolving steganographic strategies.** It 164
114 performs automatic strategy discovery and 165
115 adaptation, coordinating efficiency, impercep- 166
116 tibility, and security through a closed loop 167
117 comprising strategy retrieval, stego text gener- 168
118 ation, evaluation, and library updating. 169
- 119 • **Lifelong evolving steganographic strategy.** 170
120 A summarizer LLM codifies successful strate- 171
121 gies into structured, searchable entries; re- 172
122 trieval based on score changes and subsequent 173
123 reranking drive continual improvement and 174
124 reliable transfer across corpora and task set-
125 tings.

- **Plug-and-play high embedding rates map-** 126
127 **ping.** We design PC-DNTE, an algorithm 128
129 that keeps sampling close to the base model’s 130
131 conditional distribution and aligns probability 132
133 mass at the bin level under high embedding 134
135 rates. 136

- **Low-cost, requirement-oriented inference.** 132
133 The pipeline runs entirely at inference time, 134
135 without the need for training, fine-tuning or 136
137 parameter access. It adapts to user defined re- 138
139 quirements, delivering superior performance 140
141 across various contexts. 142

2 Related Work 138

2.1 Large Language Models 139

140 LLMs are advanced neural networks trained on ex- 141
142 tensive text corpora, enabling them to generate and 143
144 comprehend natural language. These models have 145
146 been employed as task-specific data generators for 147
148 tabular data, relational triples and instruction data, 149
150 achieving reasonable zero-shot quality under sim- 151
152 ple class-conditional prompts (Chia et al., 2022; 153
154 Schick and Schütze, 2021). 155

156 Beyond data generation, LLMs are increas- 157
158 ingly framed as engines for scientific discovery 159
160 (Karpatne et al., 2025; Luo et al., 2025). Tool- 161
162 augmented, agentic systems illustrate this shift: 163
164 ChemCrow (M. Bran et al., 2024) links LLMs with 164
165 expert chemistry tools, and Co-scientist (Boiko 166
167 et al., 2023) automates experiment planning and 168
169 cloud lab execution. Recent work has also ex- 170
171 plored hypothesis discovery and evaluation through 171
172 systems like ResearchAgent (Baek et al., 2025) 172
173 and LiveIdeaBench (Ruan et al., 2024), which 173
174 propose and refine ideas based on literature and 174
175 single keyword prompts. Complementing these 175
176 trends, jailbreak studies (Liu et al., 2025; Zhou 176
177 et al., 2025a) demonstrate self-evolving jailbreak 177
178 strategies, where agents autonomously discover 178
179 and refine attacks over time. 179

2.2 Generative Text Steganography 165

166 Generative text steganography designs stegano- 167
168 graphic strategies that steer LLMs to select tokens 168
169 so as to embed secret information within natural 169
170 language. In this paradigm, a steganographic map- 170
171 ping associates bit strings with sampling decisions 171
172 inside the model conditional token distribution, 172
173 thereby controlling token selection by the LLM. 173
174 From a provable security perspective, Zhang et al. 173
175 (Zhang et al., 2021) propose Adaptive Dynamic 174
176

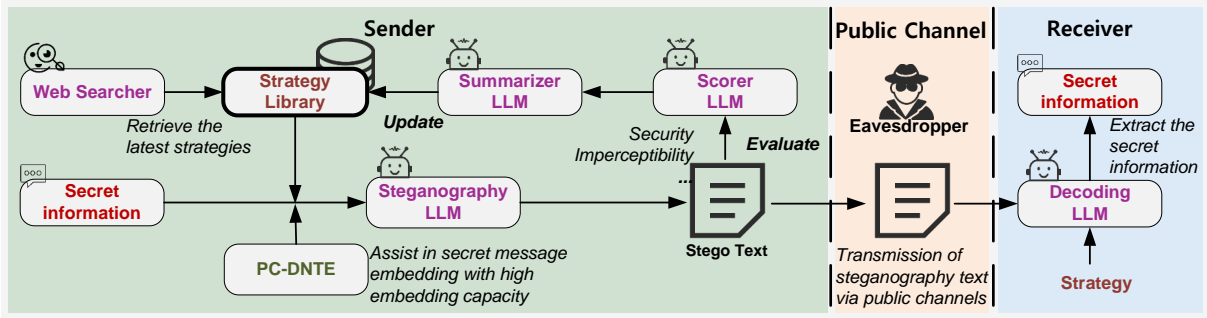


Figure 2: The overall framework of Auto-Stega.

Grouping (ADG), which recursively groups vocabulary tokens to embed information while aligning with the base distribution under stated assumptions. Ding et al. (Ding et al., 2023) introduce "Distribution Copies" (Discop), a mapping designed to preserve the model’s original distribution during embedding. Wu et al. (Wu et al., 2024) implement a prompt based mapping that leverages interactions with LLM APIs to realize practical steganographic encoding. Despite these advances, current methods still face challenges in maintaining imperceptibility across high embedding rates and domains.

3 Method

Auto-Stega is an LLM-based text steganography framework that employs an agent-driven, self-evolving mechanism to continually update a strategy library and provide steganographic strategies tailored to the task. As shown in Fig. 2, the system consists of a sender and a receiver. On the sender side, a closed loop of "generation, evaluation, summarization, and update" operates. The web searcher retrieves up-to-date strategies according to task requirements and updates the library; the steganography LLM then generates stego text, and the evaluation module, via the scorer LLM, assigns scores on multiple metrics, including efficiency, imperceptibility, and security. The summarizer LLM writes structured strategy entries back to the library, thereby realizing continuous strategy evolution. On the receiver side, the decoding LLM accurately recovers the secret information according to the shared steganographic strategy.

3.1 Web Searcher

Starting from a library L of text steganography strategies, the web searcher uses extensible literature retrieval and information extraction templates to mine recent work and convert candidate methods

Notation	Description
T_c	Cover text
T_s	Stego text
B	Number of bins
S	Score (scalar)
θ	Fixed params of base AR LM
m	Target bin index
S_T	Score threshold
R	Evaluation response

Table 1: Notations

into executable entries. It evaluates each method with standardized metrics, assigns recommended scenarios, and compares it with existing entries; candidates that exceed the threshold S_T are admitted to the lifelong evolving steganographic strategy library.

3.2 Steganographic Generator and Evaluator

Given a request M , the system first filters candidate strategies using the "Applicable Scenarios" field recorded in each library entry. It then performs multi-objective ranking of candidates based on evaluation metrics, selects the top- k strategies, and uses steganography LLM to generate stego text. The resulting outputs are assessed by the evaluation module; if the predefined acceptance criteria are not met, generation is re-initiated with the next strategy in the ranked set, and the process iterates until the stopping condition is satisfied.

The evaluation module, comprising the detection LLM, which estimates the probability that a given text is steganographic and provides explanatory rationales, and other components, applies predefined criteria, including embedding rate, perplexity (PPL) (Mikolov et al., 2010), and semantic similarity (SS). The embedding rate (ER) can be computed as:

$$ER = \frac{N_b}{W} \quad (1)$$

where N_b is the total number of bits in the embed-

ded secret information, and W is the total number of words in the stego text. PPL is a common quantitative measure in text generation, defined as:

$$\text{PPL} = \exp\left(-\frac{1}{N_w} \sum_{i=1}^{N_w} \log p(w_i | w_1, \dots, w_{i-1})\right), \quad (2)$$

where N_w is the length of the text, w_i is the i -th token in text, and $p(w_i | w_1, \dots, w_{i-1})$ is the probability assigned by the language model to the i -th word given the preceding words. The scorer LLM then compares the response R and the score S produced by the evaluation module against predefined criteria to determine whether the steganographic requirements are satisfied. Appendix A.6 gives the full formulation of the scoring function S .

3.3 Strategy Library Construction

The steganographic strategy library L is organized into two stages: warm-up and lifelong learning.

In the warm-up stage, we execute a generate-evaluate loop and log (T_s, R, S) for each request. A summarizer LLM abstracts improvements between nearby or sampled records into structured strategy entries and writes them back to the library L after deduplication. Each entry stores motivation, implementation notes, and evaluation summaries. The library L is maintained as a key-value store whose keys κ are response embeddings for efficient retrieval. Admission is gated by thresholded scores S_T , and operational details are given in Fig. 3.

In the lifelong learning stage, given a request M and current L , the system retrieves and ranks candidate strategies, regenerates by evolving new strategies from existing entries in L when acceptance criteria are not met, admits improved entries under the same threshold rule, and terminates when scores exceed preset thresholds or the iteration budget is reached. Under the same thresholding rule, strategy retrieval produces a ranked shortlist Γ and a selected set EFFECTIVE. Each candidate strategy $\gamma \in \Gamma$ is evaluated with an overall score $S(\gamma)$, and the accepted subset is

$$\mathcal{A} = \{\gamma \in \Gamma : S(\gamma) \geq S_T\}. \quad (3)$$

The selection rule is then

$$\gamma^* = \arg \max_{\gamma \in \mathcal{A}} S(\gamma). \quad (4)$$

$$\text{EFFECTIVE} = \begin{cases} \{\gamma^*\}, & |\mathcal{A}| = 1, \\ \mathcal{A}, & |\mathcal{A}| > 1, \\ \emptyset, & |\mathcal{A}| = 0. \end{cases} \quad (5)$$

Algorithm 1 Strategy Retrieval and Shortlisting

Require: Request M ; library L (entry j with key κ_j); shortlist size k
Ensure: Shortlist Γ
Functions: GE: generate and evaluate (outputs T_S, R, S); EMB: embed R ;
 STS: stable similarity top- $2k$; MD: metric discrepancy vs. current metrics;
 TOPK: pick k by criterion

```

1:  $(T_S, R, S) \leftarrow \text{GE}(M)$  // gen/eval
2:  $E_R \leftarrow \text{EMB}(R)$  // embed
3:  $\mathcal{C} \leftarrow \text{STS}(L, E_R, 2k)$  // retrieve
4: for all  $j \in \mathcal{C}$  do
5:    $d_j \leftarrow \text{MD}(j, S)$  // score gap
6: end for
7:  $\Gamma \leftarrow \text{TOPK}(\mathcal{C}, \{d_j\}, k)$  // shortlist
8: return  $\Gamma$ 

```

When $\mathcal{A} \neq \emptyset$, the system proceeds with EFFECTIVE. A single accepted strategy is applied, while multiple accepted strategies are assembled to drive further evolution. Otherwise ($\mathcal{A} = \emptyset$), the shortlist is discarded and alternative discovery is triggered. This rule turns observed gains into reusable fragments and supports continual evolution.

3.4 Strategy Retrieval

Given a request M , the steganographic LLM and the evaluator module produce a stego text, a response, and an overall score S . We embed the response R and retrieve the top- $2k$ entries from L by similarity to their keys κ_j , with stable deterministic tie resolution. The candidates are reranked by metric discrepancy with the current metrics; the best k form the shortlist Γ . Note that in the first iteration, no response is available for retrieval reference. Thus the steganography LLM is prompted without employing a steganographic strategy in the first iteration. Algorithm 1 formalizes the retrieval and selection procedure.

3.5 Plug-and-Play High Embedding Rates Mapping

Prior heuristics such as top- p (nucleus) and locally typical sampling shape the candidate set near the model conditional distribution, improving fluency. However, they lack control over token selection that is calibrated for capacity and aware of the underlying distribution: per-step capacity remains uncalibrated, probability mass is unstructured, and the mapping from bits to tokens is unstable as candidate diversity varies. Accordingly, we propose PC-DNTE (Parity-Constrained Dynamic Nucleus-Typical Encoder).

PC-DNTE is a training-free encoder. At each step, it adapts $p \in [p_{\min}, p_{\max}]$ to build a nucleus set under $p_\theta(\cdot | \text{ctx} * t)$ and applies a local typicality filter. The remaining candidates are split into

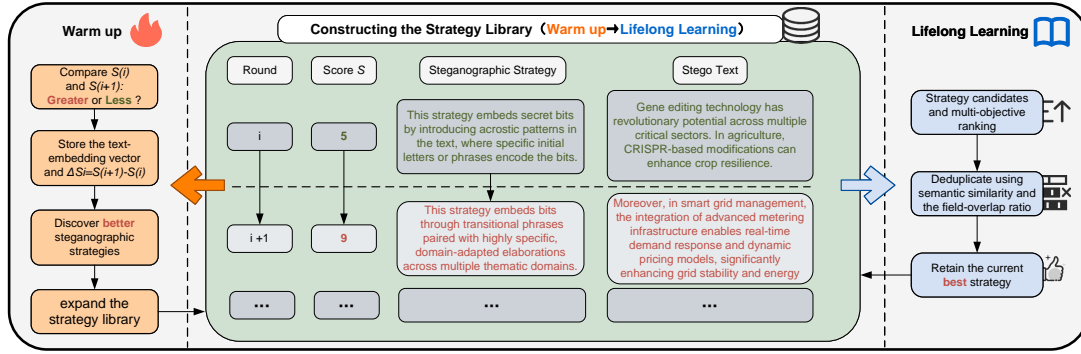


Figure 3: Steganographic strategy library construction: from warm-up to lifelong learning.

$B = 32$ equal-mass bins $\{G_m\}_{m=0}^{31}$. The bin index is driven by payload bits and is approximately uniform. Since each G_m carries the same probability mass under $p_\theta(\cdot | \text{ctx}_t)$, the induced mass over bins aligns with the base distribution at the region level. With $B = 32$, a 5-bit index selects m , and one parity bit b^* further filters tokens in G_m by a seed-conditioned parity, enabling deterministic decoding. If the filtered set is empty, p is increased and the step is retried; if it remains empty, a no-embed step is recorded, the bits are deferred, and $\arg \max_{y \in G_m} p_\theta(y | \text{ctx}_t)$ is emitted. Near sentence boundaries, a small log-odds boost to EOS and punctuation, together with a bounded closure routine, encourages well-formed breaks. This adjustment only affects the next-step candidate set, so embedded bits are unchanged. Selection within a bin is greedy, so token-level preservation is not exact. The goal is to keep sampling close to the base conditional distribution while sustaining stable capacity at high embedding rates (Algorithm 2).

3.6 Secret Information Extraction

In the secret extraction module, when configured with the same steganographic strategy as the encoding side, the decoding LLM is able to accurately decode the embedded secret information.

4 Experiment

In this section, we evaluate the performance of Auto-Stega in terms of efficiency, imperceptibility, and security, as well as human evaluation. The details of the experimental setup and result analysis are presented in the following subsections. Additional implementation details and supplementary results are provided in the appendices. Appendix A.2 provides representative high scoring strategy entries, and Appendix A.3 provides the prompt

Algorithm 2 PC-DNTE

Require: conditional $p_\theta(\cdot | \text{ctx})$; ciphertext BITS; $p_{\min}, p_{\max} \in (0, 1]$; bins $B=32$; scorer TYP(\cdot); seed SEED; closure ε ; budget K ; punctuation set \mathcal{P}
Ensure: stego text $T_S = (t_1, \dots, t_T)$
Functions: ADAP: adapt $p \in [p_{\min}, p_{\max}]$; DYNNUC: nucleus under p_θ at p ; LOCTYP: local typicality filter; EQM: split to B equal-mass bins; READBITS: pop n bits; PH: seed-conditioned parity; INCP: enlarge p up to p_{\max} ; NOEMB: record no-embed; QUEUEBITS: enqueue bits; UPDCTX: update context; SBLIKELY: boundary test; BOOSTEOS: boost \mathcal{P} ; CLOSESTEPS: bounded closure
1: $T_S \leftarrow []$; $t \leftarrow 1$
2: **while** NOT EOM(BITS) **do**
3: NO_EMB \leftarrow FALSE
4: $p \leftarrow$ ADAP($p_{\min}, p_{\max}, p_\theta(\cdot | \text{ctx}_t)$) // adapt
5: $C_t \leftarrow$ DYNNUC($p_\theta(\cdot | \text{ctx}_t), p$); $C'_t \leftarrow$ LOCTYP($C_t, \text{ctx}_t, \text{TYP}$); $\{G_0, \dots, G_{31}\} \leftarrow$ EQM(C'_t, B)
6: $m \leftarrow$ READBITS(BITS, 5); $b^* \leftarrow$ READBITS(BITS, 1)
7: $S \leftarrow \{y \in G_m : \text{PH}(y, \text{ctx}_t, \text{SEED}) = b^*\}$
8: **if** $S = \emptyset$ **then**
9: $p' \leftarrow$ INCP(p, p_{\max})
10: $C_t \leftarrow$ DYNNUC($p_\theta(\cdot | \text{ctx}_t), p'$); $C'_t \leftarrow$ LOCTYP($C_t, \text{ctx}_t, \text{TYP}$); $\{G_0, \dots, G_{31}\} \leftarrow$ EQM(C'_t, B)
11: $S \leftarrow \{y \in G_m : \text{PH}(y, \text{ctx}_t, \text{SEED}) = b^*\}$
12: **end if**
13: **if** $S = \emptyset$ **then**
14: NOEMB(); DEFBITS($\{m, b^*\}$) // enqueue bits
15: $y_t \leftarrow \arg \max_{y \in G_m} p_\theta(y | \text{ctx}_t)$;
16: NO_EMB \leftarrow TRUE // fallback
17: **else**
18: $y_t \leftarrow \arg \max_{y \in S} p_\theta(y | \text{ctx}_t)$ // greedy
19: **end if**
20: $T_S \leftarrow T_S \cup \{y_t\}$; $\text{ctx}_{t+1} \leftarrow$ UPDCTX(ctx_t, y_t); $t \leftarrow t + 1$
21: **if** SBLIKELY($\text{ctx}_t, \varepsilon$) \vee NO_EMB **then**
22: BOOSTEOS(\mathcal{P}, ε); CLOSESTEPS(ctx_t, K) // reweight
23: **end if**
24: **end while**
25: **return** T_S

templates for Auto-Stega. Appendix A.4 reports ablation studies. Appendix A.5 introduces an additional anti-steganalysis experiment.

4.1 Experimental Setup

Datasets. We evaluated Auto-Stega in three public corpora: the News Category Dataset (News) (Misra, 2022), the Large Movie Review dataset (Movie) (Maas et al., 2011), and Sentiment140 (Tweet) (Go et al., 2009). For quantitative comparability, we use 12,300 News items (300 per category across 41 categories), 20,000 Large Movie Review entries, and 20,000 Sentiment140 tweets. The combination spans domains, styles, and lengths,

Methods	News			Movie			Sentiment		
	ER \uparrow	PPL* \downarrow	SS \uparrow	ER \uparrow	PPL* \downarrow	SS \uparrow	ER \uparrow	PPL* \downarrow	SS \uparrow
ADG	4.99	8.58	0.40	5.44	5.39	0.45	5.22	0.53	0.39
Discop	4.62	0.70	0.58	4.94	0.62	0.58	5.41	0.94	0.48
LLM-Stega	5.94	0.46	0.61	4.63	0.28	0.61	7.90	0.80	0.58
Ours	5.97	0.01	0.63	5.69	0.04	0.62	8.14	0.84	0.60

Table 2: Experimental results of PPL* and SS under high embedding rates. \uparrow higher is better, \downarrow lower is better.

which allows us to assess efficiency, imperceptibility, and security under cross-scenario distribution shifts while maintaining sufficient scale for reliable estimation. Additional cross-lingual and cross-domain evaluation is provided in Appendix A.7 to assess generalization.

Baselines. We compare Auto-Stega with three SOTA baselines, including ADG (Zhang et al., 2021), Discop (Ding et al., 2023), and LLM-Stega (Wu et al., 2024). To demonstrate that our method maintains strong performance across embedding capacities, we conduct experiments under both high embedding rates and low embedding rates settings and compare the stego text generated by different methods in each case.

Hyperparameters. Evaluation was conducted in two stages. An initial warm up exploration was performed on 50 steganographic requests, running $N = 150$ iterations to build the initial strategy library. Using this library, the runtime lifelong learning phase was then executed: for each request in the experimental datasets, 5 iterations were carried out. A full steganography round is defined as generating stego text for a given cover text and completing the evaluation. Unless otherwise stated, the maximum iteration budget was set at $T = 150$ and the stopping threshold to $S_T = 8.5$ for each data instance. During evaluation, the strategy library was frozen and an additional stego generation pass was conducted on the three datasets. For plug-and-play mapping, $p_{\min} = 0.88$ and $p_{\max} = 0.95$ were used to obtain sufficient candidate tokens while suppressing low typicality tails, and GPT-2 was used to compute next-token probabilities. We use DEEPSEEK-V3.2 for the steganography, decoding, and scorer LLMs, and GPT-4o for the summarizer LLM.

Efficiency. In the field of steganography, the embedding rate serves as a key evaluation metric to assess the efficiency of a method. In text steganography, the embedding rate is defined as the average number of secret bits embedded per word (bits per word, bpw), as shown in Formula 1. In PC-DNTE,

Methods	News		Movie		Sentiment	
	PPL* \downarrow	SS \uparrow	PPL* \downarrow	SS \uparrow	PPL* \downarrow	SS \uparrow
ADG	9.26	0.45	5.92	0.49	0.39	0.38
Discop	0.90	0.60	0.88	0.59	0.99	0.44
LLM-Stega	0.03	0.64	0.11	0.67	0.94	0.58
Ours	0.01	0.66	0.08	0.72	0.86	0.65

Table 3: Experimental results of PPL* and SS at a low embedding rate of 0.1 bpw.

the encoder operates at the token level, so capacity is naturally parameterized as bits per token. In the GPT-2 setup, this corresponds approximately to 1 bit/token \approx 1.4 bit/word.

Imperceptibility. In text steganography, imperceptibility is typically categorized into text quality and statistical imperceptibility. Text quality captures human-perceived naturalness, including fluency, coherence, and semantic fidelity relative to the cover text. Statistical imperceptibility assesses the distributional alignment with the cover. Together, these facets provide a measure of whether a method conceals information without degrading readability or altering the underlying statistical characteristics of the text.

(1) Text Quality. PPL and SS are commonly used to evaluate the text quality of generated stego text. In this experiment, we choose the GPT-2 model to calculate the PPL values of different stego text. Prior work (Yang et al., 2018) indicates that social media text exhibits substantial stylistic variability, so human-written sentences deviate widely from any single language model; consequently, stego text whose PPL is closer to that of human text is considered more secure. To compare across corpora with different baseline perplexities, we use a normalized perplexity deviation, $PPL^* = \frac{|PPL_{\text{stego}} - PPL_{\text{cover}}|}{PPL_{\text{cover}}}$ where lower values indicate better text quality. For the SS, we choose the Sentence-bert (Reimers and Gurevych, 2019) method and use the roberta-base-nli-mean-tokens (Liu et al., 2019) model to extract sentence vectors, to calculate the cosine similarity between the cover text and the stego text.

(2) Statistical Imperceptibility. Kullback-Leibler Divergence (KLD) serves as an evaluation metric to measure the statistical imperceptibility of steganographic algorithms by comparing the distribution of the generated stego text with the distribution of the cover text. In our experiments, we adopt the KLD proposed by Zhang et al. (Zhang et al., 2021) to evaluate the statistical imperceptibility of the methods tested. It is formulated as follows:

$$\text{KLD}(\mu_x, \sigma_x, \mu_y, \sigma_y) = \sum \left[\log\left(\frac{\sigma_y}{\sigma_x}\right) + \frac{\sigma_x^2 + (\mu_x - \mu_y)^2}{2\sigma_y^2} - \frac{1}{2} \right] \quad (6)$$

Where μ_x and σ_x are the mean and standard deviation of the cover text, while μ_y and σ_y represent those of the stego text.

Security. Steganography and steganalysis are locked in an ongoing adversarial game. As such, anti-steganalysis performance serves as the paramount metric for evaluating the security of a steganographic algorithm. In this experiment, we leverage three advanced steganalysis methods, including SANet (Xue et al., 2024), BiLSTM-Dense (BD) (Yang et al., 2020), and Bert-FT (BF) (Peng et al., 2021). The results are reported in terms of steganalysis accuracy—defined as the rate of correctly identifying whether a given text contains hidden information—where values closer to 50% (equivalent to random guessing) indicate stronger anti-steganalysis performance.

Human Evaluation. For the human evaluation, we focus on evaluating three key aspects of the generated texts: fluency, coherence, and relevance. Our human evaluation was conducted by five researchers who were not involved in the study, and each one holds a master’s degree. The evaluation process involved a mixed assessment sheet containing both generated stego and cover text; only stego text scores were selected for analysis. The researchers, unaware of the underlying algorithms, labels, or any other details, were tasked solely with assessing and scoring text quality. Each researcher independently rated the texts on a five-point scale (ranging from ‘very poor’ to ‘very good’). The higher evaluation score denotes the better generated stego text.

4.2 Results and Analysis

(1) Efficiency and Imperceptibility. In generative text steganography, efficiency and imperceptibility are intrinsically linked: changes to the embedding rate (ER) modify sampling behavior and can materially affect both text quality and statistical imper-

Payload	Methods	News	Movie	Sentiment
		KLD↓	KLD↓	KLD↓
Low	ADG	2.63	2.28	1.90
	Discop	2.78	1.94	2.17
	LLM-Stega	2.41	2.83	2.43
	Ours	2.21	2.13	2.11
High	ADG	2.63	2.27	1.89
	Discop	1.94	1.96	2.17
	LLM-Stega	2.33	2.22	2.19
	Ours	1.87	1.98	2.05

Table 4: The experimental results of the statistical imperceptibility under high and low embedding rates. The best results are highlighted in bold red, and the second bests are shown in red without bold. ↓ lower is better.

ceptibility. Accordingly, we analyze efficiency and imperceptibility jointly in this section.

Text Quality. The experimental results for text quality and ER are listed in Tables 2 and 3. Auto-Stega achieves a higher ER while delivering superior performance. In our experiments, the average PPL of cover sentences is 113.89 on the News corpus, 111.12 on the Movie corpus, and 1224.79 on the Tweet corpus. Auto-Stega achieves the highest SS across all datasets and attains the lowest PPL* on News and Movie. On Tweet, the stego text exhibit a slightly higher PPL*; nevertheless, Auto-Stega achieves a 42.2% relative reduction compared with the SOTA baseline. To quantify the robustness of the gains and assess statistical significance, all methods were rerun with multiple random seeds. For each corpus, paired *t*-tests were conducted on sentence-level PPL*. The tests indicate that the improvements over the best baselines are statistically significant. These corpus characteristics highlight the importance of maintaining stylistic consistency with the cover text, thereby improving overall textual quality. Within Auto-Stega, PC-DNTE treats text quality as an explicit objective during generation and selection, thereby generating high quality stego text. Qualitative steganographic examples produced by the proposed method and prior baselines are provided in Appendix A.1.

Statistical Imperceptibility. In this part, we compute sentence-level KLD using distributions derived from a pre-trained BERT encoder, comparing the cover and stego sets. In our experiments, we take the logarithm of the KLD values to facilitate a clearer comparison of statistical imperceptibility between algorithms. As shown in Table 4, under both high and low embedding rates, Auto-Stega achieves the lowest KLD on the News dataset, and the second lowest KLD on the Movie and Tweet

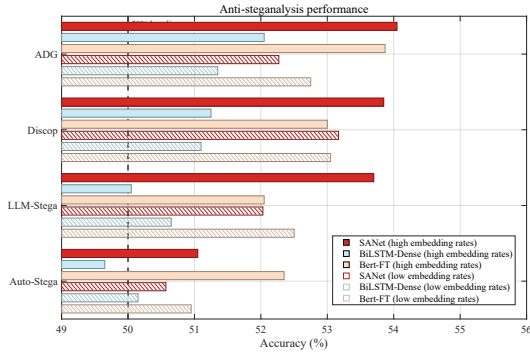


Figure 4: The results of the anti-steganalysis performance across low and high embedding rates.

535 corpora. ADG (Zhang et al., 2021) performs well
 536 on tweets, whereas Discop (Ding et al., 2023) per-
 537 forms well on reviews, because their equiprobable
 538 grouping and distribution preserving sampling better
 539 match the peaked token distributions and formulaic
 540 phrasing characteristic of these datasets. In addi-
 541 tion, our method can be configured to incorpo-
 542 rate sentiment conditioning or template-based
 543 generation, which can reduce KLD on these two
 544 corpora while preserving semantic fidelity.

545 Taken together, these results demonstrate that
 546 Auto-Stega achieves the strongest overall imper-
 547 ceptibility across embedding rates, indicating that
 548 it maintains superior perceptual concealment even
 549 under high embedding rates and is well-suited for
 550 practical deployment in real-world applications.

551 (2) **Security.** We used splits of 10,000 train-
 552 ing pairs, 1,000 validation pairs, and 1,000 test
 553 pairs, each pair containing a cover text and its
 554 corresponding stego text. Fig. 4 presents experi-
 555 mental results on anti-steganalysis performance
 556 across different embedding rates. The proposed
 557 method, Auto-stega, continues to exhibit superior
 558 performance across all metrics when compared to
 559 other steganography methods. Notably, Auto-Stega
 560 achieves SANet of 51.05%, BD of 49.65%, and BF
 561 of 52.35%, representing an average 1.2% reduc-
 562 tion relative to the SOTA method. This improve-
 563 ment reflects that, during iterative generation and
 564 selection, Auto-Stega explicitly evaluates anti-
 565 steganalysis performance and admits candidates
 566 that score better on this criterion, thereby pro-
 567 ducing stego text that is harder to detect. Overall,
 568 these results indicate that Auto-Stega achieves
 569 the strongest security among evaluated methods.

570 (3) **Human Evaluation.** To further assess the
 571 effectiveness of the proposed Auto-Stega, we con-
 572 ducted a human evaluation, with detailed results

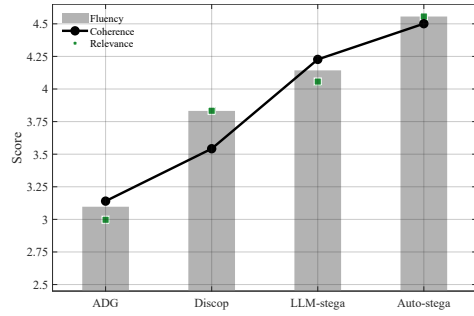


Figure 5: The results of the human evaluation.

573 shown in Figure 5. Human ratings demonstrate that
 574 Auto-Stega outperforms the three baseline meth-
 575 ods, producing stego text that is more fluent,
 576 coherent, relevant, and more imperceptible. This
 577 finding is particularly significant in the context
 578 of text steganography, underscoring the potential
 579 of Auto-Stega in applications where information
 580 concealment is critical. Moreover, as shown in Ap-
 581 pendix A.1, we randomly selected several exam-
 582 ples of stego text generated by Auto-Stega and the
 583 baselines for qualitative analysis. The results in-
 584 dicate that Auto-Stega generates texts with high
 585 fluency, grammatical correctness, and semantic
 586 coherence.

587 5 Conclusion

587 We presented Auto-Stega, an agent-driven, self-
 588 evolving framework that discovers, composes, and
 589 adapts steganographic strategies entirely at infer-
 590 ence time. By continually curating a searchable
 591 strategy library without retraining and by option-
 592 ally employing the plug-and-play PC-DNTE to keep
 593 sampling close to the base model’s conditional dis-
 594 tribution, the system supports higher embedding
 595 rates while maintaining imperceptibility and se-
 596 curity. To our knowledge, Auto-Stega is the first
 597 text steganography framework in an LLM-agent
 598 setting with self-evolving strategies. The results
 599 suggest that this paradigm already captures part
 600 of the potential of LLMs for steganography. As
 601 LLM capabilities continue to improve, further
 602 gains in steganographic performance are expected.
 603 Experiments demonstrate a consistently superior
 604 trade-off among efficiency, imperceptibility, and
 605 security compared with SOTA baselines. Future
 606 work will extend the framework to multilingual
 607 and multimodal covers, strengthen robustness
 608 against adaptive steganalyzers and noisy chan-
 609 nels, and reduce inference latency.

610 Limitations

611 While Auto-Stega achieves efficient covert text
612 communication across diverse textual domains, it
613 does not yet address multimodal settings in which
614 information may be distributed across heteroge-
615 neous cover media. As future work, we will extend
616 Auto-Stega to multimodal steganography with mul-
617 tiple covers by integrating modality aware gener-
618 ators and evaluators with the strategy library, en-
619 abling coordinated cover selection, adaptive embed-
620 ding rate allocation, and cross modal redundancy.
621 This extension is expected to substantially improve
622 practical utility and covertness in real world de-
623 ployments.

624 References

625 Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan,
626 and Sung Ju Hwang. 2025. ResearchAgent: Iterative
627 research idea generation over scientific literature with
628 large language models. In *Proceedings of the 2025*
629 *Conference of the Nations of the Americas Chapter of*
630 *the Association for Computational Linguistics: Human*
631 *Language Technologies*, pages 6709–6738.

632 Mathias Benedek, Tanja Könen, and Aljoscha C
633 Neubauer. 2012. Associative abilities underlying
634 creativity. *Psychology of Aesthetics, Creativity, and*
635 *the Arts*, 6(3):273.

636 Daniil A Boiko, Robert MacKnight, Ben Kline, and
637 Gabe Gomes. 2023. Autonomous chemical research
638 with large language models. *Nature*, 624(7992):570–
639 578.

640 Yu Cheng, Jiuan Zhou, Jiawei Chen, Zhaoxia Yin, and
641 Xinpeng Zhang. 2025. Rfnns: Robust fixed neural
642 network steganography with popular deep generative
643 models. *arXiv preprint arXiv:2505.04116*.

644 Yew Ken Chia, Lidong Bing, Soujanya Poria, and Luo
645 Si. 2022. RelationPrompt: Leveraging prompts to
646 generate synthetic data for zero-shot relation triplet
647 extraction. In *Findings of the Association for Com-*
648 *putational Linguistics: ACL 2022*, pages 45–57.

649 Jinyang Ding, Kejiang Chen, Yaofei Wang, Na Zhao,
650 Weiming Zhang, and Nenghai Yu. 2023. Discop:
651 Provably secure steganography in practice based on"
652 distribution copies". In *2023 IEEE Symposium on*
653 *Security and Privacy*, pages 2238–2255.

654 Alec Go, Richa Bhayani, and Lei Huang. 2009. Twit-
655 ter sentiment classification using distant supervision.
656 *CS224N project report, Stanford*, 1(12):2009.

657 Mingzhi Hu and Hongxia Wang. 2025. Mutual
658 information-optimized steganalysis for generative
659 steganography. *IEEE Transactions on Information*
660 *Forensics and Security*, 20:1852–1865.

Anuj Karpatne, Aryan Deshwal, Xiaowei Jia, Wei Ding,
Michael Steinbach, Aidong Zhang, and Vipin Kumar.
2025. Ai-enabled scientific revolution in the age
of generative ai: second nsf workshop report. *npj*
Artificial Intelligence, 1(1):18. 661
662
663
664
665

Arthur Koestler. 1964. The act of creation. 666

Byung Cheol Lee and Jaeyeon Chung. 2024. An empiri-
cal investigation of the impact of chatgpt on creativity.
Nature Human Behaviour, 8(10):1906–1914. 667
668
669

Mengdi Li, Kai Mu, Ping Zhong, Juan Wen, and Yim-
ing Xue. 2019. Generating steganographic image
description by dynamic synonym substitution. *Sig-*
nal Processing, 164:193–201. 670
671
672
673

Yiming Li, Kejiang Chen, Yaofei Wang, Xin Zhang,
Guanjie Wang, Weiming Zhang, and Nenghai Yu.
2025. Coas: Composite audio steganography based
on text and speech synthesis. *IEEE Transactions on*
Information Forensics and Security, 20:5978–5991. 674
675
676
677
678

Guorui Liao, Jinshuai Yang, Weizhi Shao, and Yongfeng
Huang. 2025. A framework for designing provably
secure steganography. In *34th USENIX Security Sym-*
posium, pages 6837–6856. 679
680
681
682

Xiaogeng Liu, Peiran Li, G. Edward Suh, Yevgeniy
Vorobeychik, Zhuoqing Mao, Somesh Jha, Patrick
McDaniel, Huan Sun, Bo Li, and Chaowei Xiao.
2025. Autodan-turbo: A lifelong agent for strategy
self-exploration to jailbreak llms. In *International*
Conference on Learning Representations. 683
684
685
686
687
688

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
Luke Zettlemoyer, and Veselin Stoyanov. 2019.
Roberta: A robustly optimized bert pretraining ap-
proach. *arXiv preprint arXiv:1907.11692*. 689
690
691
692
693

Ziming Luo, Zonglin Yang, Zexin Xu, Wei Yang, and
Xinya Du. 2025. Llm4sr: A survey on large lan-
guage models for scientific research. *arXiv preprint*
arXiv:2501.04306. 694
695
696
697

Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Bal-
dassari, Andrew D White, and Philippe Schwaller.
2024. Augmenting large language models with chem-
istry tools. *Nature Machine Intelligence*, 6(5):525–
535. 698
699
700
701
702

Andrew Maas, Raymond E Daly, Peter T Pham, Dan
Huang, Andrew Y Ng, and Christopher Potts. 2011.
Learning word vectors for sentiment analysis. In
Proceedings of the 49th annual meeting of the associ-
ation for computational linguistics: Human language
technologies, pages 142–150. 703
704
705
706
707
708

Laijin Meng, Xinghao Jiang, Qiang Xu, and Tanfeng
Sun. 2025. A robust coverless video steganography
based on two-level dct features against video attacks.
IEEE Transactions on Multimedia, 27:6434–6448. 709
710
711
712

713	Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In <i>Interspeech</i> , volume 2, pages 1045–1048.	768
714		769
715		770
716		771
717	Rishabh Misra. 2022. News category dataset. <i>arXiv preprint arXiv:2209.11429</i> .	772
718		
719	Wanli Peng, Jinyu Zhang, Yiming Xue, and Zhenghong Yang. 2021. Real-time text steganalysis based on multi-stage transfer learning. <i>IEEE Signal Processing Letters</i> , 28:1510–1514.	773
720		774
721		775
722		776
723	Yinyin Peng, Yaofei Wang, Donghui Hu, Kejiang Chen, Xianjin Rong, and Weiming Zhang. 2024. Ldstega: Practical and robust generative image steganography based on latent diffusion models. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 3001–3009.	777
724		778
725		779
726		780
727		781
728		782
729	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing</i> , pages 3982–3992.	783
730		784
731		785
732		786
733		787
734		788
735	Kai Ruan, Xuan Wang, Jixiang Hong, Peng Wang, Yang Liu, and Hao Sun. 2024. Liveideabench: Evaluating llms’ divergent thinking for scientific idea generation with minimal context. <i>arXiv preprint arXiv:2412.17596</i> .	789
736		790
737		791
738		792
739		793
740	Timo Schick and Hinrich Schütze. 2021. Generating datasets with pretrained language models. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6943–6951.	794
741		795
742		796
743		797
744		798
745	Mohammad Shirali-Shahreza. 2008. Text steganography by changing words spelling. In <i>2008 10th international conference on advanced communication technology</i> , volume 3, pages 1912–1913.	799
746		800
747		801
748		802
749	Jie Wang, Yihao Wang, Yihao Li, Ru Zhang, and Jianyi Liu. 2025a. Multi-classification of linguistic steganography driven by large language models. <i>IEEE Signal Processing Letters</i> , 32:1440–1444.	803
750		804
751		805
752		806
753	Zhuang Wang, Linna Zhou, Xuekai Chen, Zhili Zhou, and Zhongliang Yang. 2025b. Stlc-kg: A social text steganalysis method combining large-scale language models and common-sense knowledge graphs. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 25461–25469.	807
754		808
755		809
756		810
757		811
758		812
759	Juan Wen, Xuejing Zhou, Ping Zhong, and Yiming Xue. 2019. Convolutional neural network based text steganalysis. <i>IEEE Signal Processing Letters</i> , 26(3):460–464.	813
760		814
761		
762		
763	Jiaxuan Wu, Zhengxian Wu, Yiming Xue, Juan Wen, and Wanli Peng. 2024. Generative text steganography with large language model. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 10345–10353.	
764		
765		
766		
767		
	Yiming Xue, Jiaxuan Wu, Ronghua Ji, Ping Zhong, Juan Wen, and Wanli Peng. 2024. Adaptive domain-invariant feature extraction for cross-domain linguistic steganalysis. <i>IEEE Transactions on Information Forensics and Security</i> , 19:920–933.	
	Hao Yang, YongJian Bao, Zhongliang Yang, Sheng Liu, Yongfeng Huang, and Saimei Jiao. 2020. Linguistic steganalysis via densely connected lstm with feature pyramid. In <i>Proceedings of the 2020 ACM Workshop on Information Hiding and Multimedia Security</i> , pages 5–10.	
	Zhongliang Yang, Xiaoqing Guo, Ziming Chen, Yongfeng Huang, and Yujin Zhang. 2018. Rnn-stega: Linguistic steganography based on recurrent neural networks. <i>IEEE Transactions on Information Forensics and Security</i> , 14(5):1280–1295.	
	Siyu Zhang, Zhongliang Yang, Jinshuai Yang, and Yongfeng Huang. 2021. Provably secure generative linguistic steganography. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3046–3055.	
	Xin Zhang, Kejiang Chen, Na Zhao, Weiming Zhang, and Nenghai Yu. 2025a. Provably secure public-key steganography based on admissible encoding. <i>IEEE Transactions on Information Forensics and Security</i> , 20:3161–3175.	
	Ziping Zhang, Jiamin Zeng, Yan Xu, Xiaowei Yi, Yun Cao, and Changjun Liu. 2025b. Triple-stage robust audio steganography framework with aac encoding for lossy social media channels. In <i>Proceedings of the ACM Workshop on Information Hiding and Multimedia Security</i> , pages 131–141.	
	Andy Zhou, Kevin Wu, Francesco Pinto, Zhaorun Chen, Yi Zeng, Yu Yang, Shuang Yang, Sanmi Koyejo, James Zou, and Bo Li. 2025a. Autoredeamer: Autonomous red teaming with lifelong attack integration. <i>arXiv preprint arXiv:2503.15754</i> .	
	Pengcheng Zhou, Zhengyang Fang, Zhongliang Yang, Zhili Zhou, and Linna Zhou. 2025b. Efficient streaming voice steganalysis in challenging detection scenarios. <i>IEEE Transactions on Information Forensics and Security</i> , 20:5966–5977.	
	Qing Zhou, Ping Wei, Zhenxing Qian, Xinpeng Zhang, Sheng Li, and Chuan Qin. 2025c. Conditional flow-based generative steganography. <i>IEEE Transactions on Dependable and Secure Computing</i> , 22(5):5632–5647.	

815	A Appendix	
816	A.1 Comparison of Steganographic Examples	
817	Generated by Proposed Method and	
818	Existing Methods.	
819	In this section, we analyze randomly sampled stego	
820	text generated by Auto-Stega and three baselines	
821	(ADG, Discop, and LLM-Stega). As shown in	
822	Table 5, texts from ADG and Discop frequently ex-	
823	hibit surface issues, such as awkward collocations,	
824	agreement errors, and missing function words, es-	
825	pecially near high entropy positions induced by	
826	their mappings. LLM-Stega generally maintains	
827	sentence level fluency, but shows intersentential	
828	drift and occasional semantic inconsistencies at	
829	higher embedding rates, reflecting template and	
830	keyword rigidity. In contrast, Auto-Stega produces	
831	texts that are fluent, syntactically well formed, and	
832	semantically consistent across sentences, while	
833	more closely matching the style of the source cor-	
834	pus (for example, concise headlines in News and	
835	an informal register in Tweets). Manual inspec-	
836	tion also finds fewer detectability cues, consistent	
837	with the smallest difference from cover perplexity	
838	observed in the experimental results. Overall, the	
839	results indicate that Auto-Stega yields more coher-	
840	ent and less detectable stego text than the compared	
841	methods.	
842	A.2 Evolution of Steganographic Strategies	
843	In this section, we present steganographic strategy	
844	entries from the library spanning a range of scores.	
845	Each entry follows a standardized schema com-	
846	prising <i>Strategy</i> , <i>Definition</i> , <i>Technique</i> , <i>Applicable</i>	
847	<i>Scenarios</i> , <i>Characteristics</i> , and <i>Good Examples</i> ,	
848	where <i>Good Examples</i> include representative stego	
849	text with the corresponding metric scores. Across	
850	iterations, we observe a clear progression: early	
851	strategies yield limited embedding capacity, lower	
852	text quality, and weaker anti-steganalysis perfor-	
853	mance, whereas later strategies achieve higher em-	
854	bedding rate, better fluency and coherence, and	
855	improved statistical invisibility. The strategy evolu-	
856	tion process in Table 6 shows that Auto-Stega can	
857	autonomously develop effective strategies for text	
858	steganography, underscoring its practical value for	
859	covert communication. The evolution of stegano-	
860	graphic strategies is shown in Fig. 6; as the library	
861	is progressively updated, the generated stego text	
862	exhibits improved empirical performance.	
	A.3 Full prompts for the Steganography LLM	863
	and Summarizer LLM	864
	In this section, we present the prompt templates	865
	that operationalize the two core components of	866
	Auto-Stega. The Summarizer LLM is prompted	867
	to distill observed improvements into structured	868
	strategy entries under strict output and safety con-	869
	straints, as shown in Fig. 7. During the warm up	870
	stage, the steganographic LLM embeds a given se-	871
	cret into a cover subject to fluency and coherence	872
	requirements, with controlled length variation and	873
	avoidance of statistical artifacts, as shown in Fig.	874
	8. During the lifelong learning stage, the generator	875
	conditions on retrieved strategy headers, either a	876
	single strategy or a composed set, and applies them	877
	to produce candidates under the same quality re-	878
	quirements, as shown in Fig. 9. These templates	879
	standardize inputs and outputs across iterations,	880
	enable compositional use of strategies, and sup-	881
	port reliable extraction and consistent evaluation	882
	throughout the Auto-Stega framework.	883
	A.4 Ablation Experiment	884
	To clarify the contribution of PC-DNTE versus	885
	the self-evolving strategy library, an additional ab-	886
	lation is conducted at the same high embedding	887
	rate as in the main experiments: the Auto-Stega	888
	framework and PC-DNTE are kept unchanged, the	889
	strategy library is first evolved for three summa-	890
	rization–update rounds and then frozen (no further	891
	evolution), and stego texts are generated once from	892
	this fixed, pre-evolved library. Under this setting,	893
	PC-DNTE already achieves a strong trade-off be-	894
	tween efficiency and imperceptibility and outper-	895
	forms prior methods on most metrics. As shown in	896
	Table 7, the average PPL* over the three datasets is	897
	reduced by about 35% compared with the best pre-	898
	vious method, while SS and KLD remain competi-	899
	tive or better. These results suggest that PC-DNTE	900
	accounts for most of the gains at high embedding	901
	rates, delivering strong quality-efficiency trade-offs	902
	even with a frozen, pre-evolved library. Using the	903
	full Auto-Stega framework yields an additional im-	904
	provement (about 7.8% in average PPL*).	905
	A.5 Anti-Steganalysis Performance	906
	Evaluation	907
	Auto-Stega is further evaluated with the linguistic	908
	steganalysis network LS-CNN (LC) (Wen et al.,	909
	2019). The experimental settings are the same as	910
	those in Section 4.2. In these experiments, LS-	911

CNN serves as an external detector and its binary detection accuracy on stego texts from each method is reported. An accuracy close to 0.5 indicates that the generated stego texts are nearly indistinguishable from cover texts. As shown in Table 8, Auto-Stega attains accuracies closest to 0.5 on News, Movie, and Sentiment, indicating higher security.

A.6 Full Formulation of the Scoring Function

This appendix specifies the global scoring function S used in §3.2. Each criterion is mapped to a score $s \in [1, 10]$, where larger is better. The final score is obtained by weighted aggregation.

Perplexity closeness between cover and stego texts is quantified by

$$c_{\text{PPL}} = \exp\left(-\left|\log \text{PPL}_{\text{stego}} - \log \text{PPL}_{\text{cover}}\right|\right),$$

$$s_{\text{PPL}} = 1 + 9 c_{\text{PPL}}. \quad (7)$$

Semantic similarity is measured using the cosine similarity of Sentence-BERT embeddings,

$$s_{\text{SS}} = 1 + 9 \cdot \frac{\cos(\text{stego}, \text{cover}) + 1}{2}. \quad (8)$$

Statistical imperceptibility is captured by the divergence term in §3.2 and mapped as

$$s_{\text{KLD}} = \frac{10}{1 + \text{KLD}}. \quad (9)$$

Anti-steganalysis performance is computed from an external detector output p_{det} (stego probability),

$$s_{\text{det}} = 1 + 9(1 - p_{\text{det}}). \quad (10)$$

Capacity efficiency is represented by the embedding rate ER. Fixed thresholds $0 < \tau_1 < \dots < \tau_9$ define a monotone interval mapping,

$$s_{\text{ER}}(\text{ER}) = 1 + \sum_{k=1}^9 \mathbf{1}[\text{ER} \geq \tau_k], \quad (11)$$

which yields $s_{\text{ER}} \in \{1, \dots, 10\}$.

The global score aggregates the metric scores by a weighted sum:

$$S = w_{\text{ER}}s_{\text{ER}} + w_{\text{PPL}}s_{\text{PPL}} + w_{\text{SS}}s_{\text{SS}} + w_{\text{KLD}}s_{\text{KLD}} + w_{\text{det}}s_{\text{det}}, \quad (10)$$

$$\sum_m w_m = 1.$$

The weights used in the experiments are

$$w_{\text{ER}} = 0.15, \quad w_{\text{det}} = 0.25, \quad (11)$$

$$w_{\text{PPL}} = w_{\text{SS}} = w_{\text{KLD}} = 0.20.$$

A.7 Cross-Lingual and Cross-Domain Evaluation

Additional evaluation is conducted on cross-domain and cross-lingual corpora to assess generalization. The testbed covers Arxiv (scientific abstracts and metadata), an English Wikipedia dump, and THUCNews (Chinese news articles). For each corpus, Table 9 reports embedding rate (ER), normalized PPL*, SS, and KLD. Auto-Stega achieves the highest ER and SS on all three corpora, while keeping PPL* and KLD comparable to, or better than, the strongest baselines. These results suggest consistent performance across scientific, encyclopedic, and news domains, as well as across English and Chinese, with a favorable balance between capacity, efficiency, and imperceptibility.

Methods	Stego text example
ADG	putting for me can eat me can go well
Discop	I'm a Ser, its well over 20 years old and I just started to play. And thats it :D I'd be happy to give back to you please respond up till hour 30 More bugs I'm sure around here there are no edit black on black bars? See below. As if it wasn't clear..in game this spreadsheet just explains what blocks are. Although it didn't make clear to me what kind of bones were I carrying. When you cast it double kill yourself as if you. // "Do you want to... trust me?..." "No... everyone IS right. I definitely agree, even just... wanting a girl you have never known." {K...center, cums, gesture one more time}.....I trust your feelings, you really have." In other words it wasn't a grown up. Shame. DISCONTINUED // "Ah, Engh...huh. Thank you, am I right...?" "Yes. Your Majesty, on behalf
LLM-Stega	Officials announce an exciting new team that promises to bring fresh talent and enthusiasm to the entertainment industry.
Ours	The future might be ready: robots like computers in some scenarios would rapidly replace humans.

Table 5: Stego text examples produced by different methods.

Strategy	Core idea (one sentence)	Technique	Scenarios	Key traits	Representative stego excerpt	Score
Thematic Domain Diversification with Consistent Transitional Enc	Combine semantic preservation with light rewriting to smooth statistical signals and reduce detectability while balancing capacity and robustness (esp. back-translation).	semantic	robustness; low detectability (LLM)	semantic preservation; back-translation; robust; hard for detectors	<i>"Quantum computing will usher in a new era ... optimize trading strategies by processing market data ..."</i>	6.125
Domain-Specific Semantic Precision with Transitional Depth	Use domain-specific, non-generic elaborations with consistent transitional phrases; varied transitions align with encoded bits and thematic depth to raise naturalness and lower detectability.	Context-aware semantic embedding	general	data-driven; non-generic; varied transitions	<i>"Moreover, the implementation of quantum-resistant cryptographic algorithms necessitates evaluation of lattice-based schemes ..."</i>	7.625
Domain-Specific Semantic Depth with Multi-Domain Redundancy	Embed via consistent transitions plus rich, domain-adapted elaborations across multiple domains; introduce paragraph-level redundancy to improve error tolerance without generic phrasing.	context-aware semantic embedding	general	data-driven; semantic depth; redundancy; low detectability	<i>"Moreover, in smart grid management, advanced metering enables real-time demand response ... enhancing grid stability ..."</i>	9.000

Table 6: Strategy evolution summary for stego text generation (concise view).

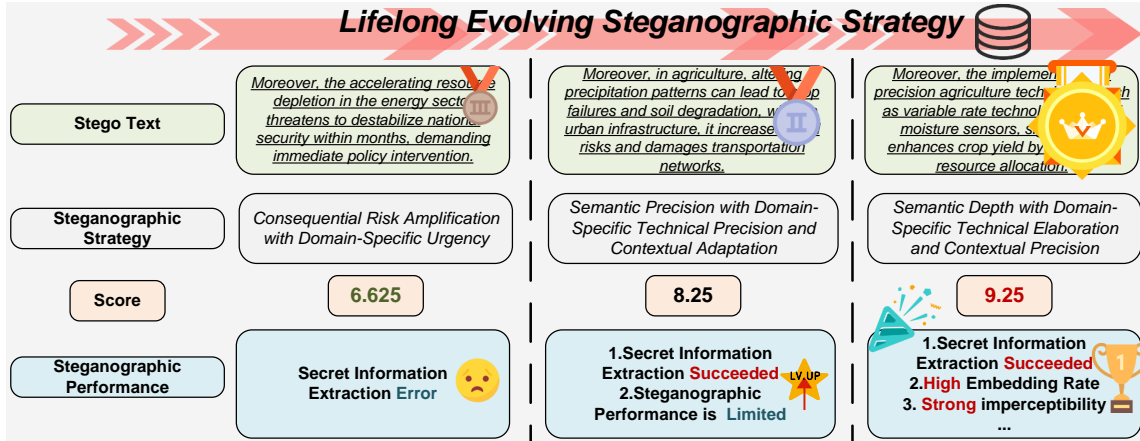


Figure 6: Lifelong evolution of steganographic strategies.

Methods	News			Movie			Sentiment		
	PPL* ↓	SS ↑	KLD ↓	PPL* ↓	SS ↑	KLD ↓	PPL* ↓	SS ↑	KLD ↓
ADG	8.13	0.42	2.61	8.51	0.43	2.27	0.52	0.40	1.90
Discop	0.70	0.59	1.93	8.61	0.58	1.95	0.95	0.49	2.18
LLM-Stega	0.45	0.62	2.33	0.28	0.61	2.23	0.80	0.59	2.20
Ours	0.02	0.64	1.87	0.06	0.62	1.97	0.88	0.61	2.01

Table 7: Results on News, Movie, and Sentiment. PPL* and KLD: lower is better; SS: higher is better.

Payload	Methods	Accuracy (%)
Low	ADG	52.05
	Discop	51.90
	LLM-Stega	52.35
	Auto-Stega	51.47
High	ADG	53.96
	Discop	53.05
	LLM-Stega	53.35
	Auto-Stega	51.67

Table 8: Steganalyzer accuracy (%) under low and high payloads. Values closer to 50 indicate better statistical imperceptibility, and the closest-to-50 result in each payload setting is highlighted in bold.

Methods	Arxiv				Wiki				THUCNews			
	ER ↑	PPL* ↓	SS ↑	KLD ↓	ER ↑	PPL* ↓	SS ↑	KLD ↓	ER ↑	PPL* ↓	SS ↑	KLD ↓
ADG	3.17	0.96	0.49	2.31	3.44	3.50	0.51	2.03	6.16	2.11	0.75	2.77
Discop	4.94	0.59	0.62	1.94	4.19	0.60	0.59	1.93	5.05	0.94	0.70	2.94
LLM-Stega	5.68	2.11	0.49	2.37	5.72	0.51	0.61	2.31	7.05	0.92	0.81	2.51
Ours	6.30	0.71	0.63	1.81	6.42	0.47	0.64	1.91	8.91	0.84	0.86	2.11

Table 9: Results on Arxiv, Wiki, and THUCNews. ER and SS: higher is better; PPL* and KLD: lower is better.

```

1 def strategy_summarizer_prompt():
2     return (
3         "You are a steganography strategy summarizer.\n"
4         "Think through the task step by step INTERNALLY and summarize your reasoning, but DO NOT reveal
         your reasoning.\n"
5         "Output ONLY a valid JSON object and nothing else.\n"
6         "\n"
7         "When unsure or if required information is missing, return the empty JSON object: {}. \n"
8         "\n"
9         "=== Output contract ===\n"
10        "Return a JSON object with EXACTLY these keys in this order:\n"
11        " name (string),\n"
12        " definition (string),\n"
13        " technique (array of strings),\n"
14        " applicable_scenarios (array of strings),\n"
15        " characteristics (array of strings),\n"
16        " examples (array of objects with keys: stego_excerpt (string), overall_score (number), scores (
         object)).\n"
17        "- Use double quotes for all strings; no trailing commas; UTF-8 text.\n"
18        "- Match the output language to the cover_text language if detectable; otherwise use English.\n"
19        "- Never include secrets or secret_bits_preview in the output; never echo raw secret bits.\n"
20        "\n"
21        "=== Internal step-by-step plan (do not expose) ===\n"
22        "1) Parse inputs (stego_text, cover_text, evaluation, used_strategies). Identify signals such as:
         "
23        "equiprobable LM, acrostic AM, cover continuation, bin size, temperature, decoding hints, etc.\n"
24        "2) Draft a concise, distinctive NAME (<= 60 chars, Title Case). Prefer pattern: "
25        "'<Core Mechanism> + <Key Modifier> (optional params)'.\n"
26        "3) Write a 1-2 sentence DEFINITION that states mechanism (how), purpose (why), and limits (if
         any).\n"
27        "4) TECHNIQUE: list concrete methods/ingredients (e.g., 'equiprobable token binning (bin=64)', "
28        "'deterministic acrostic A-M->0/N-Z->1', 'cover-conditioned continuation', 'low-temp decoding').\n"
29        "5) APPLICABLE_SCENARIOS: list contexts where this strategy is appropriate (e.g., 'short social
         posts', "
30        "'news-style paragraphs', 'low-detectability requirement', 'receiver knows cover').\n"
31        "6) CHARACTERISTICS: list crisp traits, mixing strengths/risks/requirements (e.g., 'capacity
         approx 6 bpw', "
32        "'good fluency under low temperature', 'robustness limited by rephrasing', 'needs shared cover')
         .\n"
33        "7) EXAMPLES: choose 1-3 short stego excerpts from stego_text (<= 240 chars each, no ellipsis if
         avoidable). "
34        "For each example, set overall_score and scores using evaluation if available; otherwise make a
         conservative estimate "
35        "from provided fields. Round numbers to one decimal and clamp to [0,10].\n"
36        "\n"
37        "=== Scoring normalization (internal) ===\n"
38        "- If evaluation provides per-dimension scores in [0,10], copy directly.\n"
39        "- If given probabilities (e.g., detection probability in [0,1]), map to score as (10 - 10*p),
         rounded to one decimal.\n"
40        "- If only raw metrics exist (e.g., PPL), use any provided *_score fields; do NOT invent unseen
         dimensions.\n"
41        "- If no usable signals exist for scoring, output the empty JSON object {}.\n"
42        "\n"
43        "=== Safety & formatting rules ===\n"
44        "- Do not include explanations, headings, or Markdown -- JSON only.\n"
45        "- Do not fabricate technical parameters that are not implied by inputs.\n"
46        "- Express the strategy unambiguously; given the stego_text and its corresponding strategy
         description, the embedded secret bits can be decoded with 100% reliability.\n"
47        )

```

Figure 7: Prompts used by the Summarizer LLM in Auto-Stega.

```

1 def base_generate_prompt(secret, cover_stub, method_line=""):
2     return f"""
3 System:
4 You are a professional text steganography expert. Your task is to embed a secret naturally into
   carrier text.
5
6 User:
7 You are participating in a steganography research project. Your goal is to produce high-quality
   stego text. Embed the given secret while keeping the text natural and fluent.
8
9 Requirements:
10 1) Secret bitstring to embed: {secret}
11 2) Carrier text: {cover_stub}
12 3) Keep the text natural and fluent; avoid statistical artifacts
13 4) Limit length change to within +/- 10% (if you write the carrier yourself, target 80-200 words)
14 5) You may use covert techniques such as synonym substitution, syntax restructuring, and context-
   aware semantic embedding{method_line}
15
16 Core principles:
17 - Maintain semantic consistency and topical coherence
18 - Avoid suspicious anomalies (e.g., odd token/character distributions)
19 - Ensure the secret can be correctly extracted
20
21 Produce the stego text only, and wrap it between [START STEGO TEXT] and [END STEGO TEXT].
22 """

```

Figure 8: Prompts used by the Steganography LLM during the warm-up stage of Auto-Stega.

```

1 def strategy_header_single(name, definition, technique, example):
2     return (
3         "Please use the following steganography strategy:\n"
4         f"Strategy Name: {name}\n"
5         f"Definition: {definition}\n"
6         f"Technique: {technique}\n"
7         f"Example: {example}"
8     )
9
10 def strategy_header_multi(items):
11     # items: list of dicts with keys: Strategy, Technique, Definition
12     lines = []
13     for s in items:
14         lines.append(
15             f"- {s.get('Strategy', 'unknown')} ({s.get('Technique', 'unknown')})\n"
16             f"  Definition: {s.get('Definition', '')}"
17         )
18     return "Please combine the following strategies:\n" + "\n".join(lines)
19
20 def with_strategy_prompt(secret, cover_stub, strategies_desc, method_line=""):
21     return f"""
22 System:
23 You are a professional text steganography expert. Apply the specified strategies to embed the
24     secret.
25
26 User:
27 You are participating in a steganography research project. Apply the specified strategy/
28     strategies to produce stego text.
29
30 {strategies_desc}
31
32 Requirements:
33 1) Secret bitstring to embed: {secret}
34 2) Carrier text: {cover_stub}
35 3) Keep the text natural, coherent, and statistically unremarkable
36 4) Strictly apply the specified strategy/strategies for embedding{method_line}
37
38 Produce the stego text only, and wrap it between [START STEGO TEXT] and [END STEGO TEXT].
39 """

```

Figure 9: Prompts used by the Steganography LLM during the lifelong learning stage of Auto-Stega.