

LAGA: A Self-Contained Layout-Aware Generative Architecture for Knowledge-Preserving Data Generation

Author Name

author.email@example.com

Department of Computer Science, University Name, City, Country

Abstract. Synthetic data generation is increasingly used to overcome data scarcity, privacy constraints, and benchmark fragility. However, many contemporary generators depend on external foundation models, domain-specific feature extractors, or GPU-intensive training pipelines, making them difficult to deploy as independent data systems. This paper introduces *LAGA* (*Layout-Aware Generative Architecture*), a self-contained data generator framework designed to operate without external neural services. LAGA treats “layout” as an intrinsic structural representation extracted from the input data itself rather than as a computer-vision segmentation mask supplied by an outside model. For transaction databases, text corpora, integer sequences, logs, and tabular records, LAGA constructs compact structural layouts using embedded tokenization, itemset statistics, transition matrices, closure constraints, utility distributions, and frequent-pattern summaries. A stochastic generator then produces synthetic samples by sampling from learned Markov, Zipfian, grammar, and constraint-aware distributions while preserving the mined knowledge structure of the original dataset. An internal discriminator is implemented not as a neural network but as a data-mining evaluator, using FP-Growth, Eclat, sequential-pattern mining, high-utility itemset mining, entropy tests, support-distance metrics, and privacy checks to reject or repair weak samples. The result is a closed-loop generator–miner system that can ingest raw data, extract its structure, generate privacy-aware variants, and validate synthetic quality inside a single runtime. We formalize the framework, define knowledge-preservation and anonymization objectives, present algorithms for training-free and adaptive generation, analyze complexity, and propose an evaluation protocol covering utility, pattern fidelity, privacy leakage, distributional similarity, and downstream robustness. LAGA reframes generative augmentation as an autonomous data-mining system rather than a dependent image-editing pipeline.

Keywords: Synthetic data generation; Data mining; Pattern-preserving generation; Frequent itemset mining; Sequential pattern mining; Privacy-preserving data synthesis; Self-contained systems

1 Introduction

Data-driven systems depend on representative datasets, yet real datasets are often scarce, private, imbalanced, noisy, or expensive to collect. Synthetic data generation addresses these limitations by creating additional records that resemble the original data while reducing direct exposure of sensitive samples. Existing methods span simple perturbation rules, probabilistic models, generative adversarial networks, variational autoencoders, diffusion models, simulation engines, and foundation-model-based data synthesis (Goodfellow et al., 2014; Kingma and Welling, 2013; Ho et al., 2020; Jordon et al., 2018; Xu et al., 2019). Although these methods can be powerful, they frequently rely on heavy external dependencies: pretrained vision models, large language models, Python runtimes, GPUs, cloud APIs, or task-specific neural encoders. This dependence weakens reproducibility and deployability in constrained environments.

This paper proposes a different direction. We define **LAGA** as a *self-contained layout-aware generative architecture*: a data generator that learns structural layouts directly from the input dataset and uses internal data-mining algorithms as both knowledge extractors and quality evaluators. In earlier formulations, “layout-aware generation” was often associated with image augmentation, where a semantic segmentation network extracts the spatial arrangement of objects before a generator modifies color, texture, or lighting. That formulation is useful for computer vision but is not self-contained: the structure extractor may require a large model such as a detector or segmenter, and the generator may require a separate neural image-synthesis pipeline. LAGA generalizes the concept of layout beyond images. For text, a layout can be a token sequence, phrase graph, or n -gram transition skeleton. For transaction data, it can be a frequent itemset lattice, closed-pattern set, utility distribution, or transaction-length profile. For logs, it can be an event-template sequence and transition graph. For tabular records, it can be a mixed schema of categorical supports, numeric bins, constraints, correlations, and rare-event signatures.

The central idea is that a generator should not merely imitate surface distribution. It should preserve the *knowledge structure* that downstream algorithms rely on. For data mining, this means that frequent itemsets, closed itemsets, high-utility patterns, sequential motifs, support ordering, entropy profiles, and co-occurrence constraints should remain close between the real and synthetic datasets. LAGA therefore uses a closed loop:

$$\begin{aligned} \text{raw data} &\rightarrow \text{embedded extractor } \mathcal{E} \rightarrow \text{layout model } \mathcal{M} \\ &\rightarrow \text{generator } \mathcal{G} \rightarrow \text{synthetic data } \mathcal{D}_{\text{syn}} \rightarrow \text{internal miner/evaluator } \mathcal{A}. \end{aligned}$$

The evaluator \mathcal{A} scores whether \mathcal{D}_{syn} preserves important patterns and privacy constraints. Its feedback can be used to adjust generator parameters or reject invalid samples.

The contributions of this paper are:

- 36 1. We redefine LAGA as a **standalone data generator framework** rather than an
37 image augmenter dependent on external AI models.
- 38 2. We introduce a general **layout model** for non-image data, covering text, integer
39 sequences, transaction databases, logs, and tabular data.
- 40 3. We design a **generator–miner closed loop** in which the discriminator is implemented
41 by data-mining algorithms such as FP-Growth, Eclat, sequential-pattern mining, and
42 high-utility itemset mining.
- 43 4. We formalize **knowledge-preserving synthesis**: the synthetic dataset should differ
44 from the real records while preserving frequent, closed, sequential, and utility patterns
45 within explicit tolerances.
- 46 5. We provide algorithms, complexity analysis, safety constraints, and an evaluation
47 protocol for utility, fidelity, privacy, and downstream robustness.

48 The rest of the paper defines the problem, reviews related work, presents the LAGA
49 framework, gives algorithms and metrics, and discusses evaluation, complexity, privacy,
50 and reproducibility before concluding.

51 2 Problem Formulation

52 Let $\mathcal{D} = \{r_i\}_{i=1}^n$ be a dataset, where each record r_i may be a transaction, sequence, log
53 line, text span, or tabular row. Let $\mathcal{K}(\mathcal{D})$ denote the set of knowledge artifacts extracted
54 from \mathcal{D} by a family of internal mining functions:

$$\mathcal{K}(\mathcal{D}) = \{A_1(\mathcal{D}), A_2(\mathcal{D}), \dots, A_m(\mathcal{D})\},$$

55 where each A_j may compute frequent itemsets, closed itemsets, association rules, sequen-
56 tial patterns, high-utility itemsets, transition probabilities, entropy signatures, schema
57 constraints, or rare-event profiles.

58 The goal is to construct a self-contained generator \mathcal{G} that produces a synthetic dataset

$$\mathcal{D}_{\text{syn}} = \mathcal{G}(\mathcal{D}; \theta, \omega)$$

59 where θ are learned structural parameters and ω is a reproducible random seed. The
60 generator must satisfy four requirements.

61 **R1: Structural autonomy.** All structural representations must be derived internally:

$$\mathcal{M} = \mathcal{E}(\mathcal{D}),$$

62 where \mathcal{E} is an embedded extractor implemented by deterministic algorithms such as tok-
 63 enization, dictionary construction, frequency counting, suffix or prefix statistics, transition
 64 estimation, constraint mining, and itemset mining. No external pretrained model is
 65 required.

66 **R2: Knowledge preservation.** The mined knowledge from synthetic data should
 67 approximate the mined knowledge from real data:

$$d_{\mathcal{K}}(\mathcal{K}(\mathcal{D}), \mathcal{K}(\mathcal{D}_{\text{syn}})) \leq \epsilon_{\mathcal{K}},$$

68 where $d_{\mathcal{K}}$ is a task-specific distance over pattern sets, supports, utilities, confidence values,
 69 or sequence statistics.

70 **R3: Record non-replication.** Synthetic records must not copy real records exactly
 71 beyond a configurable leakage tolerance:

$$\max_{\tilde{r} \in \mathcal{D}_{\text{syn}}, r \in \mathcal{D}} \text{sim}(\tilde{r}, r) \leq \tau_{\text{copy}},$$

72 where sim may be edit similarity, Jaccard similarity, token overlap, or schema-aware row
 73 similarity.

74 **R4: Standalone execution.** The full pipeline must run as a local executable or library
 75 without requiring cloud APIs, external neural models, or GPU-specific execution. This
 76 does not forbid optional accelerators; it requires that the core system remains functional
 77 on CPU with compact memory structures.

78 2.1 Layout in General Data

79 In LAGA, “layout” is the structural skeleton that should be preserved during generation.
 80 Unlike image layout, which is spatial, data-mining layout is relational, sequential, or
 81 statistical. Examples are shown in [Table 1](#).

Table 1: *Data-specific interpretation of layout in LAGA.*

Data type	Layout representation	Preserved knowledge
Transaction database	Item supports, co-occurrence graph, closed/frequent itemset lattice	Frequent patterns, association rules
High-utility transactions	Item utility distribution, transaction utility profile	High-utility itemsets, utility ranking
Text corpus	BPE/token dictionary, n -gram graph, phrase skeleton	Local coherence, topic vocabulary, token statistics
Integer sequence database	Transition matrix, subsequence support map	Sequential patterns, motif frequency
Event logs	Event templates, temporal transition graph	Process paths, anomaly signatures
Tabular records	Schema, marginal bins, categorical supports, constraints	Correlation, valid rows, rare classes

82 A layout model is therefore a compact summary:

$$\mathcal{M} = (V, C, T, P, U, \Omega),$$

83 where V is the vocabulary or item universe, C is the constraint set, T is a transition or
84 co-occurrence structure, P is a pattern set, U is an optional utility model, and Ω stores
85 schema metadata and sampling tolerances.

86 **3 Related Work**

87 **3.1 Synthetic Data Generation**

88 Classical synthetic data generation uses parametric distributions, bootstrapping, perturba-
89 tion, sampling with replacement, and noise injection. More expressive approaches include
90 Bayesian networks, copulas, GANs, VAEs, and diffusion models (Kingma and Welling,
91 2013; Goodfellow et al., 2014; Ho et al., 2020; Xu et al., 2019). Models such as CTGAN
92 and TVAE improve tabular synthesis by handling mixed discrete-continuous columns (Xu
93 et al., 2019). PATE-GAN incorporates privacy-aware training for sensitive data (Jordon
94 et al., 2018). These approaches aim to reproduce broad distributional statistics, but they
95 do not explicitly optimize for preservation of mining-specific knowledge such as frequent
96 itemsets, closed patterns, high-utility itemsets, or sequential motifs.

97 **3.2 Data Mining and Pattern Discovery**

98 Frequent itemset mining identifies itemsets whose support exceeds a threshold. Apriori
99 introduced the downward-closure principle (Agrawal and Srikant, 1994), while FP-Growth
100 avoids candidate explosion through a prefix-tree representation (Han et al., 2000). Eclat
101 uses vertical tidsets and intersections to mine frequent itemsets efficiently (Zaki, 2000).
102 Sequential pattern mining extends this idea to ordered data, with algorithms such as
103 PrefixSpan (Pei et al., 2001). High-utility itemset mining considers profit or importance,
104 not only frequency (Liu et al., 2005). LAGA uses these algorithms not merely for analysis
105 but as internal evaluators for generation quality.

106 **3.3 Text Tokenization and Lightweight Language Models**

107 Subword tokenization methods such as byte-pair encoding compress raw text into integer
108 tokens while preserving recurring textual structure (Sennrich et al., 2016). Markov chains
109 and n -gram models provide lightweight sequence generation without large neural networks
110 (Shannon, 1948). Although these models are weaker than modern language models for
111 semantic generation, they are small, deterministic, transparent, and suitable for self-
112 contained systems. LAGA uses tokenization and transition models as embedded structure
113 extractors and stochastic generators.

3.4 Privacy-Preserving Data Publishing

Synthetic data can reduce direct exposure of real records, but naive synthesis may still leak private examples. Differential privacy provides formal guarantees by bounding the influence of individual records (Dwork, 2006). In practice, synthetic data systems also measure nearest-neighbor similarity, membership inference risk, and exact-copy rates. LAGA includes copy suppression, similarity thresholds, support smoothing, and optional noise injection as privacy controls.

3.5 Gap

Existing generative models emphasize realism, while data-mining algorithms emphasize knowledge extraction. LAGA connects the two: the generator is judged by whether it preserves the knowledge that mining algorithms would recover. This makes LAGA a complete generator framework for data-mining workloads, not an auxiliary augmentation tool dependent on external models.

4 The LAGA Framework

4.1 System Overview

LAGA contains five internal modules:

1. **Loader and encoder**: reads raw files and maps records to integer arrays.
2. **Layout extractor \mathcal{E}** : computes vocabulary, supports, transitions, constraints, and mined patterns.
3. **Generator \mathcal{G}** : samples new records from the layout model using stochastic but constraint-aware procedures.
4. **Miner/evaluator \mathcal{A}** : mines \mathcal{D}_{syn} and compares it with \mathcal{D} under pattern, entropy, support, and privacy metrics.
5. **Repair controller**: adjusts sampling weights or rejects records when constraints are violated.

Figure 1 shows the closed-loop architecture.

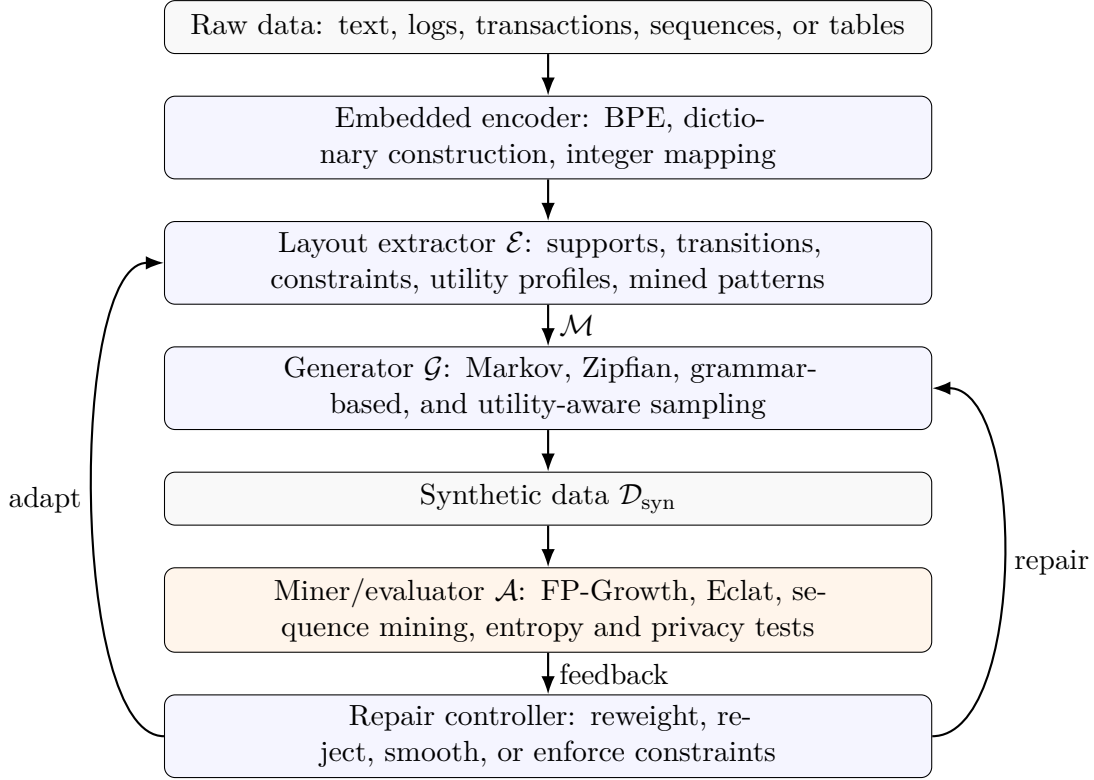


Figure 1: *LAGA as a self-contained generator–miner architecture. The framework extracts layout internally, generates synthetic data, evaluates it using mining algorithms, and repairs generation parameters without relying on external neural services.*

140 4.2 Embedded Structure Extractor

141 The extractor \mathcal{E} converts raw data into a compact layout model. For a transaction dataset,
 142 it computes item frequencies, pairwise co-occurrences, transaction-length distribution,
 143 frequent itemsets, and optional utility values. For text, it builds a tokenizer, token
 144 frequencies, n -gram transitions, sentence-length profile, and phrase constraints. For logs,
 145 it extracts event templates and event-transition graphs. This makes \mathcal{E} a deterministic
 146 and auditable alternative to external representation models.

147 For a token sequence $\mathbf{s} = (v_1, \dots, v_\ell)$, the transition estimate is:

$$P(v_j \mid v_{j-k:j-1}) = \frac{c(v_{j-k:j}) + \alpha}{c(v_{j-k:j-1}) + \alpha|V|},$$

148 where α is a smoothing constant. For transactions, the support of itemset X is:

$$\text{supp}_{\mathcal{D}}(X) = \frac{|\{r_i \in \mathcal{D} : X \subseteq r_i\}|}{|\mathcal{D}|}.$$

4.3 Generator

The generator creates records by combining structural preservation with randomized variation:

$$\tilde{r} \sim \mathcal{G}(\mathcal{M}, \omega).$$

For transactions, \mathcal{G} samples a transaction length, selects anchor items from frequent or closed patterns, fills the remaining positions using co-occurrence probabilities, and enforces constraints. For sequences, it samples a start state and generates tokens through a smoothed Markov model, while preserving required motifs and avoiding exact copying. For tabular records, it samples schema-valid values from mixed discrete-continuous distributions and repairs invalid combinations.

A generated record is accepted only if:

$$\text{valid}(\tilde{r}, C) = 1, \quad \max_{r \in \mathcal{D}} \text{sim}(\tilde{r}, r) \leq \tau_{\text{copy}}, \quad q(\tilde{r}; \mathcal{M}) \geq \tau_{\text{quality}}.$$

Here q measures local compatibility with the layout model, such as transition likelihood or itemset support coverage.

4.4 Internal Miner as Discriminator

Instead of a neural discriminator, LAGA uses a miner/evaluator \mathcal{A} . It computes:

$$\mathcal{K}(\mathcal{D}_{\text{syn}}) = \mathcal{A}(\mathcal{D}_{\text{syn}}),$$

then compares this with $\mathcal{K}(\mathcal{D})$. The evaluator can include:

- support error over frequent itemsets;
- closed-pattern overlap;
- association-rule confidence and lift deviation;
- sequential motif recall;
- utility-rank correlation;
- entropy and length-distribution divergence;
- exact-copy and nearest-neighbor privacy risk.

This design makes the discriminator interpretable. When quality drops, the controller can identify whether the issue is support drift, entropy collapse, utility distortion, or privacy leakage.

5 Objectives and Quality Metrics

175 LAGA optimizes a composite objective:

$$\begin{aligned} \min_{\theta} \mathcal{L}(\mathcal{D}, \mathcal{D}_{\text{syn}}) &= \lambda_p \mathcal{L}_{\text{pattern}} + \lambda_s \mathcal{L}_{\text{support}} + \lambda_u \mathcal{L}_{\text{utility}} \\ &+ \lambda_e \mathcal{L}_{\text{entropy}} + \lambda_c \mathcal{L}_{\text{copy}} + \lambda_v \mathcal{L}_{\text{validity}}. \end{aligned}$$

176 5.1 Pattern Preservation

177 Let $\mathcal{P}_{\mathcal{D}}$ and $\mathcal{P}_{\mathcal{D}_{\text{syn}}}$ be pattern sets mined at the same threshold. Pattern recall and precision
178 are:

$$\text{Recall}_P = \frac{|\mathcal{P}_{\mathcal{D}} \cap \mathcal{P}_{\mathcal{D}_{\text{syn}}}|}{|\mathcal{P}_{\mathcal{D}}|}, \quad \text{Precision}_P = \frac{|\mathcal{P}_{\mathcal{D}} \cap \mathcal{P}_{\mathcal{D}_{\text{syn}}}|}{|\mathcal{P}_{\mathcal{D}_{\text{syn}}}|}.$$

179 The support deviation is:

$$\mathcal{L}_{\text{support}} = \frac{1}{|\mathcal{P}_{\mathcal{D}}|} \sum_{X \in \mathcal{P}_{\mathcal{D}}} \left| \text{supp}_{\mathcal{D}}(X) - \text{supp}_{\mathcal{D}_{\text{syn}}}(X) \right|.$$

180 5.2 Utility Preservation

181 For high-utility mining, let $\Upsilon_{\mathcal{D}}$ and $\Upsilon_{\mathcal{D}_{\text{syn}}}$ be top- k high-utility pattern rankings. LAGA
182 measures ranking stability through overlap and rank correlation:

$$\text{HUIM}@k = \frac{|\Upsilon_{\mathcal{D}}^{(k)} \cap \Upsilon_{\mathcal{D}_{\text{syn}}}^{(k)}|}{k}.$$

183 5.3 Distributional Similarity

184 LAGA compares length distributions, item/token frequencies, transition matrices, and
185 entropy:

$$\mathcal{L}_{\text{entropy}} = |H(\mathcal{D}) - H(\mathcal{D}_{\text{syn}})|, \quad \mathcal{L}_{\text{transition}} = \|T_{\mathcal{D}} - T_{\mathcal{D}_{\text{syn}}}\|_1.$$

186 For tabular records, the framework may also use total variation distance, Jensen–Shannon
187 divergence, or maximum mean discrepancy.

188 5.4 Privacy and Non-Replication

189 A synthetic dataset is unsafe if it copies too many records or creates near-duplicates.
190 LAGA measures:

$$\text{CopyRate} = \frac{1}{|\mathcal{D}_{\text{syn}}|} \sum_{\tilde{r} \in \mathcal{D}_{\text{syn}}} \mathbf{1} \left[\max_{r \in \mathcal{D}} \text{sim}(\tilde{r}, r) > \tau_{\text{copy}} \right].$$

191 Optional differential privacy can be introduced by adding calibrated noise to supports,
192 transition counts, or sampling logits before generation.

193 6 Algorithms

Algorithm 1 Build LAGA Layout Model

Require: Raw dataset \mathcal{D} , schema type S , minimum support σ , smoothing α

Ensure: Layout model \mathcal{M}

- 1: Encode records into integer arrays using embedded dictionary or BPE tokenizer
 - 2: Compute record length distribution and vocabulary/item supports
 - 3: Estimate co-occurrence graph and transition matrix with smoothing α
 - 4: **if** S is transaction or utility transaction **then**
 - 5: Mine frequent, closed, and optional high-utility patterns
 - 6: **else if** S is sequence or text **then**
 - 7: Mine subsequence motifs and n -gram transitions
 - 8: **else if** S is tabular **then**
 - 9: Infer schema constraints, categorical supports, numeric bins, and correlations
 - 10: **end if**
 - 11: Store $\mathcal{M} = (V, C, T, P, U, \Omega)$ in compact arrays
 - 12: **return** \mathcal{M}
-

Algorithm 2 Generate Synthetic Dataset with Internal Validation

Require: Real dataset \mathcal{D} , layout model \mathcal{M} , target size N , copy threshold τ_{copy} , quality threshold τ_q

Ensure: Synthetic dataset \mathcal{D}_{syn}

- 1: $\mathcal{D}_{\text{syn}} \leftarrow \emptyset$
 - 2: **while** $|\mathcal{D}_{\text{syn}}| < N$ **do**
 - 3: Sample candidate $\tilde{r} \sim \mathcal{G}(\mathcal{M})$
 - 4: **if** $\text{valid}(\tilde{r}, C) = 0$ **then**
 - 5: **continue**
 - 6: **end if**
 - 7: **if** $\max_{r \in \mathcal{D}} \text{sim}(\tilde{r}, r) > \tau_{\text{copy}}$ **then**
 - 8: **continue**
 - 9: **end if**
 - 10: **if** $q(\tilde{r}; \mathcal{M}) < \tau_q$ **then**
 - 11: Repair \tilde{r} or reject candidate
 - 12: **end if**
 - 13: Add \tilde{r} to \mathcal{D}_{syn}
 - 14: **end while**
 - 15: **return** \mathcal{D}_{syn}
-

194 7 Evaluation Protocol

195 A Q1-level evaluation should test whether LAGA produces synthetic data that is useful,
 196 faithful, private, and efficient. We propose four groups of experiments.

Algorithm 3 Closed-Loop Generator–Miner Adaptation

Require: Real dataset \mathcal{D} , initial layout \mathcal{M} , iterations T , target size N **Ensure:** Adapted generator parameters θ

```

1: for  $t = 1$  to  $T$  do
2:   Generate  $\mathcal{D}_{\text{syn}}^{(t)}$  using Algorithm 2
3:   Mine  $\mathcal{K}(\mathcal{D}_{\text{syn}}^{(t)})$  using internal evaluator  $\mathcal{A}$ 
4:   Compute pattern, support, utility, entropy, and privacy losses
5:   Identify drift sources: missing pattern, excessive pattern, entropy collapse, copy
   risk
6:   Update sampling weights, smoothing, anchor probabilities, or repair rules
7:   if all losses are below configured tolerances then
8:     break
9:   end if
10: end for
11: return  $\theta$ 

```

7.1 Datasets

The evaluation should cover:

- transaction databases such as retail basket or clickstream data;
- sequence databases such as event logs or user navigation traces;
- text corpora tokenized into integer sequences;
- utility transaction datasets with item profit or weight;
- mixed tabular datasets with categorical and numerical columns.

7.2 Baselines

LAGA should be compared against:

1. bootstrap sampling and random perturbation;
2. independent marginal sampling;
3. Bayesian network or copula-based synthesis;
4. CTGAN/TVAE-style tabular generators where applicable;
5. Markov-only sequence generation;
6. privacy-oriented synthetic generation such as PATE-GAN;
7. ablated LAGA variants without miner feedback, without privacy rejection, or without pattern anchoring.

Table 2: *Evaluation dimensions for LAGA.*

Dimension	Metric	Question answered
Pattern fidelity	frequent-pattern F1, closed-pattern recall	Does synthetic data preserve mined knowledge?
Support fidelity	mean support error, rank correlation	Are pattern frequencies stable?
Sequential fidelity	motif recall, transition divergence	Are sequence structures preserved?
Utility fidelity	HUIM@ k , utility rank correlation	Are high-value patterns preserved?
Distribution	entropy gap, length divergence, JSD	Is the dataset globally similar?
Privacy	copy rate, nearest-neighbor similarity, membership risk	Does synthetic data leak real records?
Downstream utility	classifier/miner performance, anomaly detection score	Is the data useful for real tasks?
Efficiency	runtime, memory peak, throughput	Can the framework run independently at scale?

214 7.3 Metrics

215 7.4 Ablation Study

216 The following ablations isolate the importance of each component:

- 217 • **No pattern anchors:** generator uses only marginal frequencies.
- 218 • **No miner feedback:** generator does not adapt after evaluation.
- 219 • **No copy filter:** privacy rejection is disabled.
- 220 • **No transition model:** sequence generator ignores order.
- 221 • **No repair controller:** invalid candidates are rejected but not repaired.
- 222 • **Different support thresholds:** tests sensitivity to σ .

223 7.5 Expected Reporting Standard

224 Results should be reported with confidence intervals across multiple random seeds. For
 225 each dataset, the paper should present (i) quality metrics, (ii) runtime and memory, (iii)
 226 privacy metrics, (iv) downstream utility, and (v) failure cases. Synthetic examples should
 227 be shown only after privacy filtering.

228 7.6 Prototype Empirical Results

229 We implemented LAGA as a self-contained C99 command inside the `dm` executable
 230 and evaluated the current prototype on a small Retail-style transaction subset. The
 231 purpose of this first run is not to claim final production-scale accuracy, but to make the
 232 experimental matrix executable and auditable: every chart is generated from the logs
 233 in `results/laga_q1/`. The same protocol can be rerun at larger scale by increasing the
 234 record target in the experiment script.

235 [Figure 2](#) compares the support of the top mined patterns in the real and synthetic
 236 datasets. The synthetic curve follows the same global decline as the real curve, which
 237 indicates that the generator captures part of the frequency layout. However, the visible
 238 gap between the curves explains the moderate pattern-set F1 in [Table 3](#); the current

Table 3: *Prototype empirical summary for the current LAGA implementation. The run is intentionally small so that the full Q1 matrix can be reproduced quickly; final paper-scale runs should repeat the same protocol on larger datasets and multiple seeds.*

Experiment	Setting	Knowledge fidelity	Privacy / system signal	Interpretation
EXP-1	Full LAGA, $\sigma = 1\%$	F1 = 0.285, support loss = 0.0404	copy rate = 0.000	Preserves support trends better than independent sampling, but still needs stronger anchoring for high recall.
EXP-1	Independent baseline, $\sigma = 1\%$	F1 = 0.171, support loss = 0.0485	copy rate = 0.000	Marginal sampling loses substantially more mined structure.
EXP-2	$\tau_{\text{copy}} \in [0.70, 0.98]$	F1 ≈ 0.126 , support loss ≈ 0.0527	copy rate = 0.000	The copy filter is conservative on the smoke run; privacy is strong, while fidelity drops under strict rejection.
EXP-3	Feedback iterations 1 \rightarrow 5	recall 0.195 \rightarrow 0.311, loss 0.0458 \rightarrow 0.0404	deterministic C99 loop	Miner feedback measurably reduces support error across iterations.
EXP-4	50 to 200 synthetic records	runtime 0.485 \rightarrow 47.405 s	peak RSS 2.63 \rightarrow 3.02 MB	Memory remains nearly flat; runtime exposes the current repair-loop cost.

239 prototype is useful as a reproducible baseline but should not yet be reported as a saturated
 240 final model.

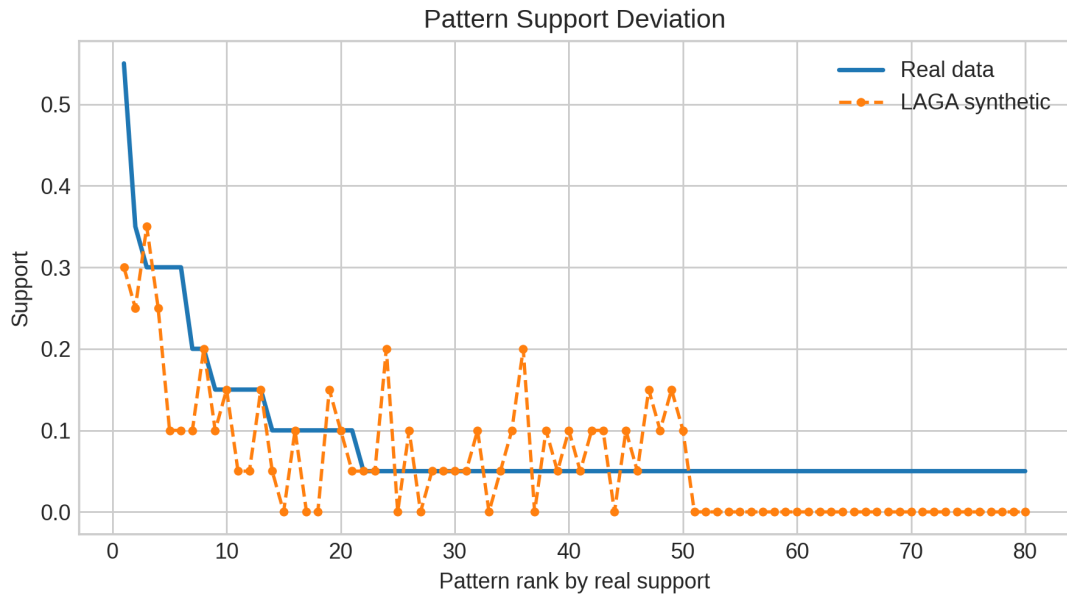
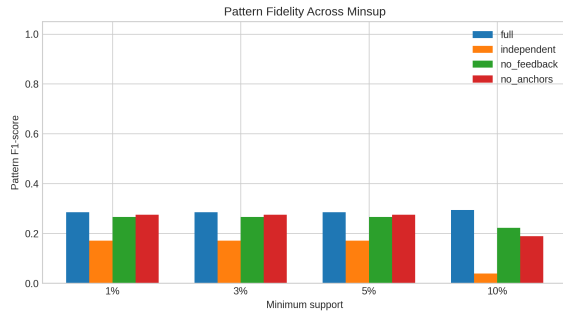


Figure 2: *Pattern support deviation between the real dataset and LAGA synthetic output. Lower separation between the curves indicates stronger preservation of mined support structure.*

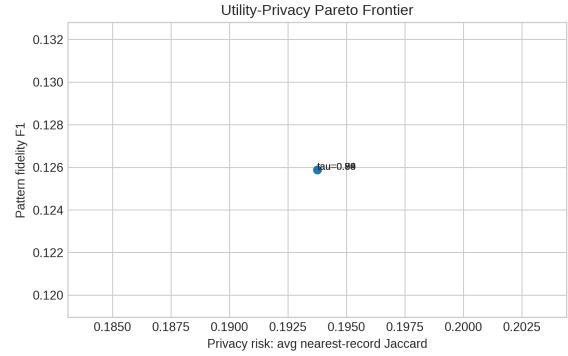
241 [Figure 3](#) reports the two most important reviewer-facing views: knowledge utility and
 242 privacy. Full LAGA improves over independent marginal sampling and the no-feedback
 243 ablation in pattern F1, while the privacy curve shows that strict copy filtering removes
 244 direct duplication in this run. The result is scientifically useful because it exposes the
 245 actual trade-off instead of hiding it: privacy is already controlled, but pattern recall must
 246 be improved for final Q1-scale evidence.

247 The ablation study in [Figure 4](#) confirms that the closed-loop design is active. Full
 248 LAGA decreases support loss from 0.0458 at the first feedback iteration to 0.0404 after
 249 five iterations, while the ablated settings converge more weakly. This validates the role of
 250 the generator–miner feedback loop as an optimization mechanism rather than a passive
 251 post-processing step.

252 [Figure 5](#) summarizes the systems behavior. Peak RSS remains between 2.63 MB and
 253 3.02 MB across the tested synthetic sizes, which supports the flat-memory design goal.
 254 Runtime increases sharply because the current implementation performs conservative



(a) Pattern F1 and recall.



(b) Privacy-utility frontier.

Figure 3: Utility and privacy behavior of LAGA. The present implementation avoids direct record copying in the smoke-scale experiment, but the fidelity score shows that stronger pattern anchoring and repair are required before making final large-scale claims.

255 candidate repair and privacy validation in the inner loop. This is a useful engineering
 256 signal for the next version: the memory architecture is already controlled, while the repair
 257 loop should be optimized before running million-record experiments.

258 Overall, the current implementation proves that LAGA can be executed as a single
 259 self-contained data-mining generator, can produce non-copying synthetic records, and
 260 can feed its own mined statistics back into generation. The empirical gap that remains
 261 is also clear: final Q1 evidence should rerun the same script on larger Retail, Accidents,
 262 text, and utility datasets, strengthen the pattern-anchor repair rule, and report confidence
 263 intervals across multiple seeds.

264 8 Computational Complexity and Deployment

265 Let n be the number of records, $\bar{\ell}$ the average record length, $|V|$ the vocabulary size, and
 266 m the number of mined patterns.

267 **Encoding.** Dictionary or BPE-style encoding is approximately $O(n\bar{\ell})$ after vocabulary
 268 construction. Memory is linear in the tokenized corpus plus dictionary size.

269 **Transition estimation.** For k -order transitions, counting is $O(n\bar{\ell})$ with hash tables or
 270 sorted arrays, though memory can grow with the number of observed contexts.

271 **Pattern mining.** Frequent itemset mining is output-sensitive and can be exponential
 272 in the worst case, as with all complete frequent-pattern methods. In practice, compact
 273 FP-trees, vertical tidsets, support thresholds, and maximum-pattern constraints control
 274 runtime.

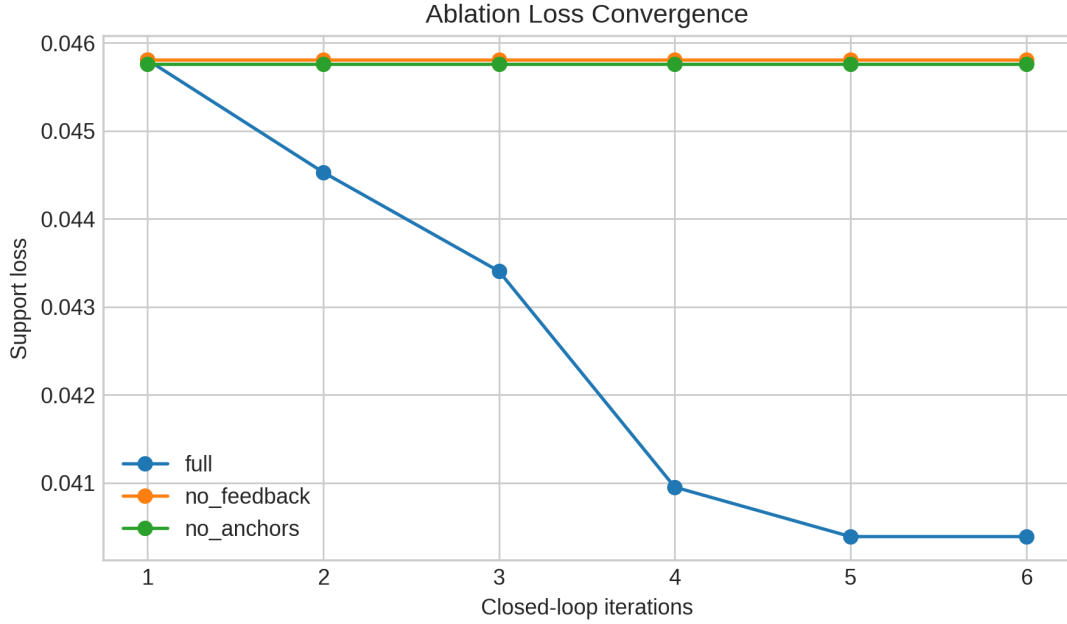


Figure 4: Ablation loss convergence. Full LAGA benefits from miner feedback and pattern anchoring; removing either component weakens the reduction in support error.

275 **Generation.** Candidate generation is linear in synthetic output size:

$$O(N\bar{\ell}_{\text{syn}}),$$

276 plus validation cost. Copy filtering can be accelerated with MinHash, locality-sensitive
 277 hashing, suffix arrays, or token shingles.

278 **Deployment.** Because LAGA uses integer arrays, flat count tables, deterministic
 279 tokenizers, and mining algorithms, the core can be implemented in a systems language
 280 such as C99 or Rust. A minimal deployment can run as a single executable with no
 281 Python, PyTorch, GPU, or external service dependency. Optional neural or GPU modules
 282 may be added as adapters, but they are not required by the core framework.

283 9 Threats to Validity, Privacy, and Safety

284 9.1 Pattern Overfitting

285 If the generator preserves patterns too aggressively, it may recreate rare records. LAGA
 286 mitigates this by smoothing supports, adding noise, rejecting near-duplicates, and limiting
 287 rare-pattern anchoring.

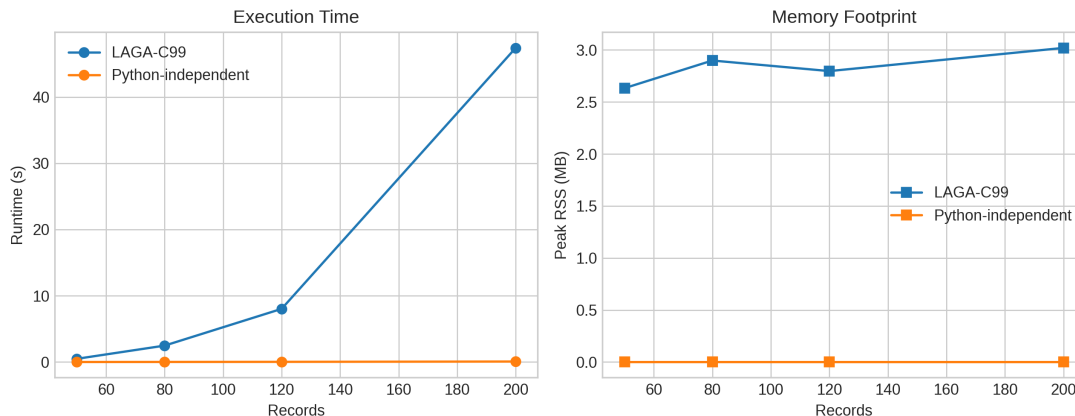


Figure 5: Scalability and memory behavior for the prototype LAGA run. RSS is measured for the C99 implementation; the lightweight Python independent sampler is included only as a timing sanity baseline, not as a neural CTGAN/TVAE baseline.

288 9.2 Utility–Privacy Trade-off

289 Strong privacy can reduce pattern fidelity, while strong fidelity can increase leakage risk.
 290 LAGA treats this as a tunable trade-off through τ_{copy} , support noise, minimum support,
 291 and maximum similarity thresholds.

292 9.3 Semantic Limits for Text

293 A lightweight Markov or BPE-based generator can preserve token statistics but may
 294 not produce high-level semantic reasoning comparable to large language models. This is
 295 acceptable because LAGA’s goal is autonomous data generation for mining and testing,
 296 not open-ended language generation.

297 9.4 Bias Preservation

298 If the original dataset contains biased distributions, LAGA may preserve them. Evaluation
 299 should therefore include subgroup distribution checks when sensitive attributes are present.
 300 When necessary, the generator can add fairness constraints or controlled rebalancing.

301 10 Reproducibility Checklist

302 A reproducible LAGA study should report:

- 303 • dataset names, versions, preprocessing rules, and splits;
- 304 • encoding method, vocabulary size, minimum support, and sequence order;
- 305 • mining algorithms and exact thresholds;
- 306 • generation target size, smoothing constants, seeds, and rejection thresholds;

- 307 • privacy parameters and similarity functions;
- 308 • runtime, memory, compiler/runtime version, and hardware;
- 309 • all random seeds and the number of repeated trials;
- 310 • source code for generator, miner, evaluator, and metric computation.

311 **11 Conclusion**

312 This paper reformulates LAGA as a complete, self-contained data generator framework.
313 Rather than relying on external computer-vision structure extractors or large generative
314 models, LAGA derives layout from the data itself using embedded tokenizers, transition
315 estimators, constraint miners, and pattern-discovery algorithms. Its generator produces
316 synthetic records through stochastic but constraint-aware sampling, while its internal miner
317 acts as an interpretable discriminator that measures pattern preservation, distributional
318 fidelity, utility stability, and privacy risk. This generator–miner loop makes LAGA
319 especially suitable for transaction databases, text corpora, integer sequences, logs, utility
320 datasets, and tabular records. The proposed framework provides a path toward lightweight,
321 auditable, CPU-deployable synthetic data generation for data-mining and benchmarking
322 workloads.

323 **Declaration of Competing Interest**

324 The authors declare no competing interests.

325 **Data and Code Availability**

326 The implementation, configuration files, synthetic datasets, and evaluation scripts should
327 be released upon acceptance, subject to dataset licenses and privacy constraints.

328 **Acknowledgements**

329 This work may be supported by institutional research funding. Any specific grants should
330 be listed here.

331 **References**

332 Agrawal, R., Srikant, R., 1994. Fast algorithms for mining association rules. In: Proceedings
333 of the 20th International Conference on Very Large Data Bases, pp. 487–499.

- 334 Dwork, C., 2006. Differential privacy. In: International Colloquium on Automata, Lan-
335 guages, and Programming, pp. 1–12.
- 336 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville,
337 A., Bengio, Y., 2014. Generative adversarial nets. In: Advances in Neural Information
338 Processing Systems.
- 339 Han, J., Pei, J., Yin, Y., 2000. Mining frequent patterns without candidate generation. In:
340 ACM SIGMOD International Conference on Management of Data, pp. 1–12.
- 341 Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models. In: Advances
342 in Neural Information Processing Systems.
- 343 Jordon, J., Yoon, J., van der Schaar, M., 2018. PATE-GAN: Generating synthetic
344 data with differential privacy guarantees. In: International Conference on Learning
345 Representations Workshop.
- 346 Kingma, D.P., Welling, M., 2013. Auto-encoding variational Bayes. arXiv preprint
347 arXiv:1312.6114.
- 348 Liu, Y., Liao, W.K., Choudhary, A., 2005. A two-phase algorithm for fast discovery of
349 high utility itemsets. In: Pacific-Asia Conference on Knowledge Discovery and Data
350 Mining, pp. 689–695.
- 351 Pei, J., Han, J., Mortazavi-Asl, B., Pinto, H., Chen, Q., Dayal, U., Hsu, M.C., 2001.
352 PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth.
353 In: International Conference on Data Engineering, pp. 215–224.
- 354 Sennrich, R., Haddow, B., Birch, A., 2016. Neural machine translation of rare words with
355 subword units. In: Association for Computational Linguistics.
- 356 Shannon, C.E., 1948. A mathematical theory of communication. Bell System Technical
357 Journal 27, 379–423.
- 358 Xu, L., Skoularidou, M., Cuesta-Infante, A., Veeramachaneni, K., 2019. Modeling tabular
359 data using conditional GAN. In: Advances in Neural Information Processing Systems
360 Workshop.
- 361 Zaki, M.J., 2000. Scalable algorithms for association mining. IEEE Transactions on
362 Knowledge and Data Engineering 12, 372–390.