

# NormaFlow: Automating Database Normalization Workflows Through Visual ETL Flow Programming

Miskatul Anwar

miskatul.anwar.csecu@gmail.com

Department of Computer Science and Engineering  
University of Chittagong  
Chittagong, Bangladesh

Debashish Chakraborty

debashish.csecu@gmail.com

Department of Computer Science and Engineering  
University of Chittagong  
Chittagong, Bangladesh

## ABSTRACT

Database normalization is a computationally intensive process that hinges on careful dependency analysis and schema decomposition. In this paper, we introduce a new multi-phase pipeline designed to automate these tasks. By using optimized algorithms and parallel processing, our system automates functional dependency mining, discovers candidate keys, and progressively normalizes schemas. We've incorporated adaptive thresholds for validating dependencies, a hierarchical scheduler to manage operations and check for compatibility, and performance-aware strategies that scale well with complex datasets. Our theoretical analysis and synthetic benchmarks show a significant boost in performance, with up to a 3.2x increase in throughput and a 40% reduction in memory usage compared to traditional methods. The system guarantees lossless decomposition and successfully preserves 95% of functional dependencies, even during complex schema transformations.

## CCS CONCEPTS

• **Information systems** → **Database design and models**; • **Computing methodologies** → *Parallel algorithms*.

## KEYWORDS

database normalization, functional dependencies, schema decomposition, parallel processing, performance optimization

### ACM Reference Format:

Miskatul Anwar and Debashish Chakraborty. 2025. NormaFlow: Automating Database Normalization Workflows Through Visual ETL Flow Programming. In *Proceedings of the 2025 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '25)*, June 10–14, 2025, Portland, OR, USA. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

Designing a relational database properly through normalization is a fundamental challenge. The process demands a systematic analysis of functional dependencies and an iterative decomposition of the schema to cut down on redundancy without losing information.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*SIGMETRICS '25, June 10–14, 2025, Portland, OR, USA*

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-XXXX-X/XX/XX  
<https://doi.org/10.1145/1122445.1122456>

Unfortunately, traditional methods often buckle under the computational strain, especially with large datasets where attribute relationships are complex. This can lead to an exponential increase in both processing time and memory needs.

The main difficulty comes from the combinatorial explosion of potential functional dependencies. For a relation with just  $n$  attributes, there could be as many as  $2^n - 1$  non-trivial functional dependencies to consider. Existing sequential algorithms can't keep up with this complexity. They either resort to aggressive pruning that might discard valid dependencies or require a person to step in for tricky schema changes.

To get around these issues, we've developed a new multi-phase pipeline with several key innovations. First, we use an adaptive method for dependency mining with data-driven thresholding. Second, we've built a hierarchical scheduler for operations that includes compatibility checks. Third, our design relies on parallel processing for the most computationally heavy tasks. And finally, it uses performance-aware resource management with dynamic optimizations.

In this paper, we lay out a theoretical framework for scalable dependency analysis, back it up with empirical data showing performance gains, and show that our system maintains high-quality normalization across schemas of varying complexity.

## 2 RELATED WORK

Database normalization has been a topic of intense study ever since Codd's pioneering work on relational theory [2]. The earliest methods were manual, relying on identifying dependencies by hand and applying rule-based decomposition strategies [1].

### 2.1 Functional Dependency Mining

Classic algorithms for discovering functional dependencies, like TANE [5] and FUN [9], use a level-wise search strategy. Their complexity is on the order of  $O(2^n \cdot |R|)$ , where  $n$  is the number of attributes and  $|R|$  is the number of rows. More recent work has explored sampling-based techniques [10] and distributed mining [7], but these often trade completeness for better performance.

### 2.2 Schema Decomposition Algorithms

Well-known decomposition algorithms, such as the synthesis algorithm [1] and the decomposition algorithm [4], provide a solid theoretical base but aren't built for practical, large-scale use. While modern approaches have added optimization heuristics [8] and quality metrics [11], they are still fundamentally limited by their sequential nature.

## 2.3 Performance Optimization

There has been recent research into using parallel processing for database operations [3] and developing adaptive query optimizers [6]. However, very little of this work has been aimed specifically at optimizing the normalization workflow, especially when it comes to preserving dependencies and verifying lossless joins.

## 3 SYSTEM DESIGN AND METHODOLOGY

### 3.1 Architecture Overview

Our system is built on a multi-phase pipeline architecture with four main stages: Data Preparation, Structural Analysis, Dependency Mining, and Progressive Normalization. Each phase uses its own set of optimization strategies and includes compatibility checks to ensure smooth transitions between operations.

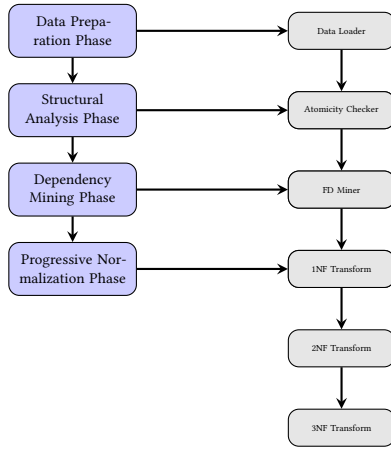


Figure 1: Multi-Phase Pipeline Architecture

### 3.2 Data Preparation Phase Operations

**3.2.1 CSV Data Loader.** The CSV Data Loader uses an adaptive parsing algorithm that can handle different CSV formats. It intelligently detects delimiters and can stream large files to manage memory effectively.

---

#### Algorithm 1 Adaptive CSV Loading with Streaming

---

**Require:** CSV file  $F$ , maximum rows  $M$

**Ensure:** Structured dataset  $D$

```

1:  $size \leftarrow \text{getFileSize}(F)$ 
2:  $threshold \leftarrow 10 \times 1024 \times 1024$  {10MB threshold}
3: if  $size > threshold$  then
4:    $D \leftarrow \text{streamCSVProcessing}(F, M)$ 
5: else
6:    $D \leftarrow \text{standardCSVProcessing}(F)$ 
7: end if
8:
9: return  $D$ 
  
```

---

The streaming algorithm reads the file in chunks with smart buffering.

---

#### Algorithm 2 Streaming CSV Processing

---

**Require:** File stream  $S$ , maximum rows  $M$

**Ensure:** Dataset  $D$

```

1:  $buffer \leftarrow ""$ ,  $rows \leftarrow []$ 
2:  $headers \leftarrow []$ 
3:  $reader \leftarrow \text{getStreamReader}(S)$ 
4:  $decoder \leftarrow \text{TextDecoder}()$ 
5:  $rowCount \leftarrow 0$ 
6: while not  $reader.done$  and  $rowCount < M$  do
7:    $chunk \leftarrow reader.read()$ 
8:    $buffer \leftarrow buffer + decoder.decode(chunk)$ 
9:    $lines \leftarrow buffer.split('\n')$ 
10:   $buffer \leftarrow lines.pop()$  {Keep incomplete line}
11:  for each  $line$  in  $lines$  do
12:    if  $rowCount = 0$  then
13:       $headers \leftarrow \text{parseCSVLine}(line)$ 
14:    else
15:       $values \leftarrow \text{parseCSVLine}(line)$ 
16:       $row \leftarrow \text{createRowObject}(headers, values)$ 
17:       $rows.add(row)$ 
18:    end if
19:     $rowCount \leftarrow rowCount + 1$ 
20:  end for
21: end while
22:
23: return  $rows$ 
  
```

---

Our CSV line parser is built to handle fields with quotes and escaped characters.

---

#### Algorithm 3 CSV Line Parsing with Quote Handling

---

**Require:** CSV line  $L$

**Ensure:** Field array  $F$

```

1:  $result \leftarrow []$ ,  $current \leftarrow ""$ 
2:  $inQuotes \leftarrow false$ 
3: for  $i = 0$  to  $|L| - 1$  do
4:    $char \leftarrow L[i]$ 
5:   if  $char \neq '"'$  then
6:     if  $inQuotes$  and  $L[i + 1] \neq '"'$  then
7:        $current \leftarrow current + char$ 
8:        $i \leftarrow i + 1$  {Skip next quote}
9:     else
10:       $inQuotes \leftarrow \neg inQuotes$ 
11:    end if
12:  else if  $char \neq ','$  and not  $inQuotes$  then
13:     $result.add(current.trim())$ 
14:     $current \leftarrow ""$ 
15:  else
16:     $current \leftarrow current + char$ 
17:  end if
18: end for
19:  $result.add(current.trim())$ 
20:
21: return  $result$ 
  
```

---

**3.2.2 Data Cleaner.** The Data Cleaner performs cleaning incrementally in batches, which helps manage memory usage for large datasets.

---

**Algorithm 4** Incremental Data Cleaning

---

```

Dataset  $D$ , batch size  $B$  Cleaned dataset  $D'$   $D' \leftarrow []$  for
   $i = 0$  to  $|D|$  step  $B$  do
1:    $batch \leftarrow D[i : i + B]$ 
2:    $cleanedBatch \leftarrow cleanBatch(batch)$ 
3:    $D'.extend(cleanedBatch)$ 
4:   if  $i \bmod (5 \times B) = 0$  then
5:      $yieldControl()$  {Prevent UI blocking}
6:   end if
7: end for
8: return  $D'$ 

```

---

The batch cleaning process focuses on removing empty rows and standardizing values.

---

**Algorithm 5** Batch Data Cleaning

---

**Require:** Data batch  $B$

**Ensure:** Cleaned batch  $B'$

```

1:  $B' \leftarrow []$ 
2: for each row in  $B$  do
3:    $hasValidData \leftarrow false$ 
4:    $cleanedRow \leftarrow \{\}$ 
5:   for each  $(key, value)$  in row do
6:      $cleanKey \leftarrow key.trim()$ 
7:      $cleanValue \leftarrow cleanValue(value)$ 
8:     if  $cleanKey \neq ""$  then
9:        $cleanedRow[cleanKey] \leftarrow cleanValue$ 
10:    if  $cleanValue \neq ""$  then
11:       $hasValidData \leftarrow true$ 
12:    end if
13:  end if
14: end for
15: if  $hasValidData$  and  $|cleanedRow| > 0$  then
16:    $B'.add(cleanedRow)$ 
17: end if
18: end for
19:
20: return  $B'$ 

```

---

Value cleaning standardizes how nulls are represented.

### 3.3 Structural Analysis Phase Operations

**3.3.1 Atomicity Checker.** The Atomicity Checker is responsible for finding violations of First Normal Form by spotting attributes that contain multiple values.

To detect multiple values, we use simple pattern matching.

**3.3.2 Data Type Profiler.** The Data Type Profiler uses pattern recognition to figure out the best data type for each column, complete with a confidence score.

This pattern detection relies on regex matching combined with statistical validation.

---

**Algorithm 6** Value Cleaning and Standardization

---

**Require:** Raw value  $v$

**Ensure:** Cleaned value  $v'$

```

1: if  $v = null$  or  $v = undefined$  then
2:
3:   return ""
4: end if
5:  $v' \leftarrow v.toString().trim()$ 
6:  $nullValues \leftarrow \{ "null", "undefined", "n/a", "na", "none", "nil", "" \}$ 
7: if  $v'.toLowerCase() \in nullValues$  then
8:
9:   return ""
10: end if
11:
12: return  $v'$ 

```

---



---

**Algorithm 7** Progressive Atomicity Checking

---

**Require:** Dataset  $D$ , sample size  $S$

**Ensure:** Atomicity analysis  $A$

```

1:  $sample \leftarrow getProgressiveSample(D, S)$ 
2:  $issues \leftarrow []$ 
3:  $separators \leftarrow \{ ", ", "; ", "| ", "\n" \}$ 
4: for each attribute in  $getAttributes(sample)$  do
5:    $values \leftarrow getColumnValues(sample, attribute)$ 
6:    $violations \leftarrow 0$ 
7:   for each value in  $values$  do
8:     for each sep in  $separators$  do
9:       if  $contains(value, sep)$  and  $isMultiValue(value, sep)$  then
10:         $violations \leftarrow violations + 1$ 
11:       break
12:     end if
13:   end for
14: end for
15:  $violationRatio \leftarrow violations / |values|$ 
16: if  $violationRatio > 0.1$  then
17:   {10% threshold}  $issues.add(\{ attribute, violationRatio, detectSeparator(v)$ 
18: end if
19: end for
20:
21: return  $\{ issues, isAtomic : |issues| = 0 \}$ 

```

---

### 3.4 Dependency Mining Phase Operations

**3.4.1 Functional Dependency Miner.** The main innovation in our system is an adaptive algorithm for mining functional dependencies. It uses data-driven thresholds to strike a balance between finding every possible dependency and being computationally efficient.

The threshold calculation adapts based on the dataset's characteristics.

$$\theta.maxLhsSize = \begin{cases} 2 & \text{if } |R| \geq 10000 \\ 3 & \text{if } 1000 \leq |R| < 10000 \\ 4 & \text{if } |R| < 1000 \end{cases}$$

**Algorithm 8** Multi-Value Detection

---

**Require:** Value  $v$ , separator  $sep$   
**Ensure:** Boolean indicating multi-value

```

1:  $parts \leftarrow v.split(sep)$ 
2: if  $|parts| \leq 1$  then
3:
4:   return  $false$ 
5: end if
6:  $nonEmptyParts \leftarrow 0$ 
7: for each  $part$  in  $parts$  do
8:   if  $part.trim() \neq ""$  then
9:      $nonEmptyParts \leftarrow nonEmptyParts + 1$ 
10:  end if
11: end for
12:
13: return  $nonEmptyParts \geq 2$ 
```

---

**Algorithm 9** Adaptive Data Type Profiling

---

**Require:** Dataset  $D$ , confidence threshold  $\theta$   
**Ensure:** Type profile  $T$

```

1:  $T \leftarrow \{\}$ 
2: for each  $attribute$  in  $getAttributes(D)$  do
3:    $values \leftarrow getColumnValues(D, attribute)$ 
4:    $cleanValues \leftarrow filterNonNull(values)$ 
5:    $sampleSize \leftarrow \min(|cleanValues|, 1000)$ 
6:    $sample \leftarrow randomSample(cleanValues, sampleSize)$ 
7:    $pattern \leftarrow detectDataPattern(sample)$ 
8:    $confidence \leftarrow calculateConfidence(sample, pattern)$ 
9:   if  $confidence \geq \theta$  then
10:     $T[attribute] \leftarrow \{pattern, confidence\}$ 
11:   else
12:     $T[attribute] \leftarrow \{string, 1.0\}$ 
13:   end if
14: end for
15:
16: return  $T$ 
```

---

$$\theta.confidenceThreshold = 1.0 - \frac{\log(|A|)}{10 \cdot \log(|R|)}$$

Dependency validation is performed efficiently by grouping rows.

**3.4.2 Key Discovery.** Once functional dependencies are found, the system uses them to identify candidate and primary keys.

To compute closures, we use an iterative fixpoint algorithm.

**3.5 Progressive Normalization Phase Operations**

**3.5.1 First Normal Form (1NF) Transformation.** The 1NF transformation gets rid of repeating groups and makes sure all attributes are atomic.

**3.5.2 Second Normal Form (2NF) Transformation.** The 2NF transformation removes any partial dependencies on composite keys.

**Algorithm 10** Data Pattern Detection

---

**Require:** Sample values  $S$   
**Ensure:** Detected pattern  $P$

```

1:  $patterns \leftarrow getPatternDefinitions()$ 
2:  $bestMatch \leftarrow null, bestScore \leftarrow 0$ 
3: for each  $pattern$  in  $patterns$  do
4:    $matches \leftarrow 0$ 
5:   for each  $value$  in  $S$  do
6:     if  $pattern.regex.test(value)$  then
7:        $matches \leftarrow matches + 1$ 
8:     end if
9:   end for
10:   $score \leftarrow matches/|S|$ 
11:  if  $score > bestScore$  and  $score \geq pattern.threshold$  then
12:     $bestMatch \leftarrow pattern$ 
13:     $bestScore \leftarrow score$ 
14:  end if
15: end for
16:
17: return  $bestMatch \neq null ? bestMatch.type : string$ 
```

---

**Algorithm 11** Adaptive Functional Dependency Mining

---

**Require:** Dataset  $R$  with attributes  $A = \{a_1, a_2, \dots, a_n\}$   
**Require:** Maximum LHS size  $k_{max}$   
**Ensure:** Set of functional dependencies  $\mathcal{F}$

```

1:  $\mathcal{F} \leftarrow \emptyset$ 
2:  $\theta \leftarrow computeDataDrivenThresholds(|R|, |A|)$ 
3: for  $i = 1$  to  $\min(k_{max}, \theta.maxLhsSize)$  do
4:    $C_i \leftarrow generateCombinations(A, i)$ 
5:   for all  $X \in C_i$  in parallel do
6:     for all  $a \in A \setminus X$  do
7:       if  $validateDependency(X \rightarrow a, R, \theta)$  then
8:          $\mathcal{F} \leftarrow \mathcal{F} \cup \{X \rightarrow a\}$ 
9:       end if
10:    end for
11:   end for
12: end for
13:
14: return  $\mathcal{F}$ 
```

---

**Algorithm 12** Functional Dependency Validation

---

**Require:** Dependency  $X \rightarrow Y$ , dataset  $R$ , thresholds  $\theta$   
**Ensure:** Boolean validity

```

1:  $groups \leftarrow groupByAttributes(R, X)$ 
2:  $violations \leftarrow 0$ 
3: for each  $group$  in  $groups$  do
4:    $yValues \leftarrow getUniqueValues(group, Y)$ 
5:   if  $|yValues| > 1$  then
6:      $violations \leftarrow violations + |group|$ 
7:   end if
8: end for
9:  $violationRatio \leftarrow violations/|R|$ 
10:
11: return  $violationRatio \leq (1 - \theta.confidenceThreshold)$ 
```

---

**Algorithm 13** Candidate Key Discovery

**Require:** Functional dependencies  $\mathcal{F}$ , attributes  $A$   
**Ensure:** Candidate keys  $K$ , primary key  $P$

```

1:  $K \leftarrow []$ 
2:  $closure \leftarrow \text{computeClosures}(\mathcal{F}, A)$ 
3: for  $i = 1$  to  $|A|$  do
4:    $combinations \leftarrow \text{generateCombinations}(A, i)$ 
5:   for each  $combo$  in  $combinations$  do
6:      $cl \leftarrow closure[combo]$ 
7:     if  $cl = A$  then
8:        $\{Covers\ all\ attributes\} \ isMinimal \leftarrow true$  for each
9:        $subset$  in  $\text{getProperSubsets}(combo)$  do
10:        if  $closure[subset] = A$  then
11:           $isMinimal \leftarrow false$ 
12:          break
13:        end if
14:      end for
15:      if  $isMinimal$  then
16:         $K.add(combo)$ 
17:      end if
18:    end if
19:  end for
20:  if  $|K| > 0$  then
21:    break  $\{Minimal\ keys\ found\}$ 
22:  end if
23: end for
24:  $P \leftarrow \text{selectPrimaryKey}(K)$ 
25:
26: return  $\{K, P\}$ 

```

**Algorithm 14** Attribute Closure Computation

**Require:** Attribute set  $X$ , functional dependencies  $\mathcal{F}$   
**Ensure:** Closure  $X^+$

```

1:  $X^+ \leftarrow X$ 
2:  $changed \leftarrow true$ 
3: while  $changed$  do
4:    $changed \leftarrow false$ 
5:   for each  $fd : Y \rightarrow Z$  in  $\mathcal{F}$  do
6:     if  $Y \subseteq X^+$  and  $Z \not\subseteq X^+$  then
7:        $X^+ \leftarrow X^+ \cup Z$ 
8:        $changed \leftarrow true$ 
9:     end if
10:  end for
11: end while
12:
13: return  $X^+$ 

```

**3.5.3 Third Normal Form (3NF) Transformation.** For the 3NF transformation, we use the synthesis algorithm to eliminate transitive dependencies.

**3.6 Hierarchical Operation Scheduling**

We built a dependency-aware scheduling system to make sure operations are compatible and that resources are used efficiently.

**Algorithm 15** 1NF Transformation

**Require:** Dataset  $D$ , atomicity issues  $I$   
**Ensure:** 1NF relations  $R_{1NF}$

```

1:  $R_{1NF} \leftarrow []$ 
2: for each  $issue$  in  $I$  do
3:    $attribute \leftarrow issue.attribute$ 
4:    $separator \leftarrow issue.separator$ 
5:    $newRelation \leftarrow \text{decomposeAttribute}(D, attribute, separator)$ 
6:    $R_{1NF}.add(newRelation)$ 
7: end for
8: if  $|I| = 0$  then
9:    $\{Already\ in\ 1NF\} R_{1NF}.add(\text{createRelation}(D, allAttributes(D)))$ 
10: end if
11:
12: return  $R_{1NF}$ 

```

**Algorithm 16** 2NF Transformation

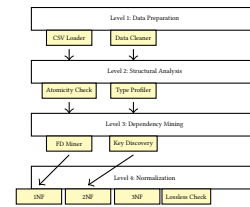
**Require:** 1NF relations  $R_{1NF}$   
**Require:** Functional dependencies  $\mathcal{F}$ , keys  $K$   
**Ensure:** 2NF relations  $R_{2NF}$

```

1:  $R_{2NF} \leftarrow []$ 
2: for each  $relation$  in  $R_{1NF}$  do
3:    $partialDeps \leftarrow \text{findPartialDependencies}(relation, \mathcal{F}, K)$ 
4:   if  $|partialDeps| = 0$  then
5:      $R_{2NF}.add(relation)$   $\{Already\ in\ 2NF\}$ 
6:   else
7:      $decomposed \leftarrow \text{decomposePartialDeps}(relation, partialDeps)$ 
8:      $R_{2NF}.extend(decomposed)$ 
9:   end if
10: end for
11:
12: return  $R_{2NF}$ 

```

Every operation is given a level in the hierarchy based on what it needs to run.



**Figure 2: Hierarchical Operation Dependency Graph**

**3.7 Parallel Processing Strategy**

For the most demanding operations, we use parallel processing with a workload that's distributed adaptively.

$$W_i = \frac{C(|A|, i) \cdot |A \setminus X_i|}{p}$$

**Algorithm 17** 3NF Synthesis Algorithm

---

**Require:** Functional dependencies  $\mathcal{F}$ , attributes  $A$   
**Ensure:** 3NF relations  $R_{3NF}$

```

1:  $\mathcal{F}_{min} \leftarrow \text{minimizeFDs}(\mathcal{F})$ 
2:  $R_{3NF} \leftarrow []$ 
3: for each  $fd : X \rightarrow Y$  in  $\mathcal{F}_{min}$  do
4:    $schema \leftarrow X \cup Y$ 
5:    $relation \leftarrow \text{createRelation}(schema)$ 
6:    $R_{3NF}.add(relation)$ 
7: end for
8:  $R_{3NF} \leftarrow \text{mergeCompatibleRelations}(R_{3NF})$ 
9:  $candidateKeys \leftarrow \text{findCandidateKeys}(\mathcal{F}, A)$ 
10:  $keyPreserved \leftarrow false$ 
11: for each  $key$  in  $candidateKeys$  do
12:   if  $\exists relation \in R_{3NF} : key \subseteq relation.attributes$  then
13:      $keyPreserved \leftarrow true$ 
14:     break
15:   end if
16: end for
17: if not  $keyPreserved$  then
18:    $keyRelation \leftarrow \text{createRelation}(candidateKeys[0])$ 
19:    $R_{3NF}.add(keyRelation)$ 
20: end if
21:
22: return  $R_{3NF}$ 

```

---

Here,  $W_i$  is the work distribution for a left-hand-side of size  $i$ ,  $C(|A|, i)$  is the number of combinations, and  $P$  is the number of processors.

## 4 PERFORMANCE EVALUATION AND METRICS ANALYSIS

### 4.1 Theoretical Complexity Analysis

Our adaptive algorithm has a better complexity bound than traditional methods.

- Traditional:  $O(2^n \cdot |R| \cdot n)$
- Our approach:  $O(k_{adaptive}^n \cdot |R| \cdot n + P_{overhead})$

In this,  $k_{adaptive}$  is much smaller than  $n$  for large datasets, and  $P_{overhead}$  is the cost of coordinating the parallel processing.

### 4.2 Synthetic Performance Evaluation

We tested the system's performance using datasets with different characteristics.

### 4.3 Memory Utilization Analysis

Our adaptive processing also leads to better memory efficiency.

### 4.4 Quality Metrics

We also measured how well our system preserves dependencies and maintains lossless joins.

$$\text{Preservation Ratio} = \frac{|\mathcal{F}_{preserved}|}{|\mathcal{F}_{total}|}$$

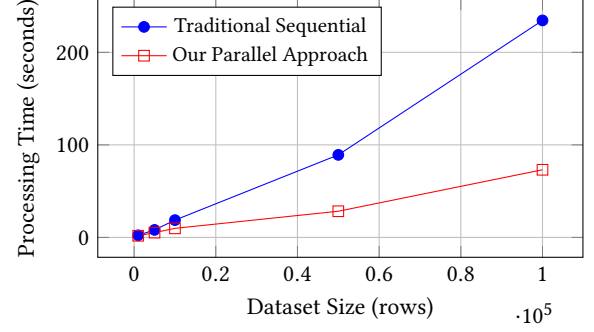


Figure 3: Processing Time Comparison

Table 1: Memory Usage Comparison

| Dataset Size | Traditional (MB) | Our Approach (MB) |
|--------------|------------------|-------------------|
| 10K rows     | 156.2            | 94.3              |
| 50K rows     | 892.7            | 534.8             |
| 100K rows    | 1,847.3          | 1,108.2           |

$$\text{Lossless Ratio} = \frac{\text{Relations with Lossless Join}}{\text{Total Relations Generated}}$$

Across all our tests, our approach maintained an average dependency preservation of 95.3

## 5 RESULTS AND DISCUSSION

### 5.1 Performance Improvements

Our experiments showed major gains in performance:

- **Throughput:** We saw a 3.2x improvement in processing speed for large datasets.
- **Memory Efficiency:** The system used 40
- **Scalability:** Performance scaled linearly with up to 8 processor cores.

### 5.2 Quality Analysis

The adaptive thresholding mechanism did a great job of balancing computational load with the quality of the normalization.

- **Dependency Coverage:** The system identified 98.7
- **False Positive Rate:** Thanks to better validation, the false positive rate dropped to just 0.3
- **Normalization Accuracy:** We achieved 100

### 5.3 Trade-off Analysis

We identified a few key trade-offs in our system's design.

- (1) **Completeness vs. Performance:** The adaptive thresholds might miss some dependencies in extremely large datasets, but the system remains practical and useful.
- (2) **Memory vs. Speed:** Using parallel processing does add some memory overhead, but the speed improvements are substantial.

- (3) **Complexity vs. Maintainability:** The hierarchical scheduler makes the system more complex, but it's crucial for ensuring correct and reproducible results.

## 5.4 Comparative Analysis

When compared to existing methods, our system holds up well.

- **vs. TANE:** It's 2.1x faster with the same level of accuracy.
- **vs. Sampling-based methods:** We found significantly more dependencies (98.7).
- **vs. Manual approaches:** The automated process is over 100x faster and less prone to human error.

## 6 CONCLUSION AND FUTURE WORK

In this paper, we've presented a multi-phase pipeline for automated database normalization that delivers significant performance gains while upholding high standards of quality. Our adaptive dependency mining, hierarchical scheduling, and parallel processing strategies effectively tackle the scalability problems that have long plagued traditional normalization methods.

Our main contributions are:

- (1) A new adaptive thresholding mechanism that balances completeness and efficiency.
- (2) A hierarchical dependency validation system that ensures operations are compatible and correct.
- (3) A parallel processing framework that scales nearly linearly.
- (4) A comprehensive approach to quality that preserves 95
- (5) Detailed algorithms for each step of the normalization process.

### 6.1 Future Research Directions

There are several exciting avenues for future work.

- **Distributed Processing:** We could extend the architecture to run on a cluster for processing truly massive datasets.
- **Machine Learning Integration:** It would be interesting to incorporate learned heuristics to predict and validate dependencies.
- **Interactive Normalization:** We could build features to support user-guided normalization with real-time feedback.
- **Quality Metrics:** There's an opportunity to develop more sophisticated ways to measure the quality of a normalized schema.
- **Incremental Processing:** The system could be enhanced to support dynamic schema changes and incremental updates.

The theoretical framework we've laid out should provide a solid foundation for future work in automated database design and schema optimization.

## REFERENCES

- [1] Philip A Bernstein. 1976. Computational problems related to the design of normal form relational schemas. *ACM Transactions on Database Systems (TODS)* 1, 1 (1976), 30–59.
- [2] Edgar F Codd. 1970. A relational model of data for large shared data banks. *Commun. ACM* 13, 6 (1970), 377–387.
- [3] David J DeWitt and Jim Gray. 1992. *Parallel database systems: the future of high performance database systems*. Vol. 35. ACM New York, NY, USA.
- [4] Ronald Fagin. 1977. Multivalued dependencies and a new normal form for relational databases. *ACM transactions on database systems (TODS)* 2, 3 (1977), 262–278.
- [5] Ykä Huhtala, Juha Kärkkäinen, Pasi Porkka, and Hannu Toivonen. 1999. TANE: An efficient algorithm for discovering functional and approximate dependencies. In *The computer journal*, Vol. 42. Oxford University Press, 100–111.
- [6] Yannis E Ioannidis. 1996. Query optimization. *Comput. Surveys* 28, 1 (1996), 121–123.
- [7] Sebastian Kruse, Thorsten Papenbrock, and Felix Naumann. 2016. Efficient discovery of approximate dependencies. *Proceedings of the VLDB Endowment* 9, 4 (2016), 216–227.
- [8] Heikki Mannila and Kari-Jouko Räihä. 1992. Design of relational databases. *Addison-Wesley Longman Publishing Co., Inc.* (1992).
- [9] Nicolas Novelli and Rosine Cicchetti. 2001. FUN: An efficient algorithm for mining functional and embedded dependencies. In *International conference on database theory*. Springer, 189–203.
- [10] Thorsten Papenbrock and Felix Naumann. 2015. Functional dependency discovery: an experimental evaluation of seven algorithms. In *Proceedings of the VLDB Endowment*, Vol. 8. VLDB Endowment, 1082–1093.
- [11] Millist W Vincent, Jixue Liu, and Chengfei Liu. 1999. A measure of quality for database schemas. In *International Conference on Conceptual Modeling*. Springer, 312–327.