# Data Quality

## First Step to Verify Table

### CUSTOMER

| First | Maiden | SSN | Last Order |
|-------|--------|-----|------------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

# Data Quality

## Previously Are Statistical; DataPrep For Example

# Data Quality

## Previously Are Statistical; Need Semantics

✔ Maiden for married only

❌ SSN can't start with 999

✔ Last Order near current

**Semantic Detection**

False Positive

❌ Maiden has 86% NULL

False Negative

✔ SSN follows 999-XX-XXXX

False Positive

❌ Last Order Peaks in May

**Statistical Detection (Numeric/Pattern Outlier)**

CUSTOMER

| First | Maiden | SSN | Last Order |
|-------|--------|-----|------------|
| John | NULL | 999-45-6789 | 9/17/2024 |
| Emily | Davis | 999-65-4321 | 9/22/2024 |
| Mike | NULL | 999-78-9012 | 9/18/2024 |
| ... | ... | ... | ... |

# Data Quality

## Semantics Achieved by LLMs

❌ SSN can't start with 999

For SSN Column, it follows the pattern 999-XX-XXXX. Is it normal?

No, that pattern for a Social Security Number (SSN) is not normal in the United States because 999 is not a valid area number.

# How semantics is applied

**Previous Profiler**

Statistical Profile



❌ Maiden has <u>86%</u> NULL

Table

| First | Maiden | SSN | Last Order |
|-------|--------|-----|-----------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

Document

Table stores customer...

# How semantics is applied



**Previous Profiler**

Semantic Review

✓ Maiden NULL is normal

Statistical Profile

✗ Maiden has <u>86%</u> NULL

Table

| First | Maiden | SSN | Last Order |
|-------|--------|-----|------------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

Semantic Profile

"Maiden" is only not NULL for married

Semantic Context

"Maiden" is maiden name

Document

Table stores customer...

**Cocoon Profiler**

# Cocoon: Semantic Table Profiling Using LLMs

Zachary Huang, Eugene Wu
HILDA @ SIGMOD 2024

# Table profiling

First step to understand table



CUSTOMER

| First | Maiden | SSN | Last Order |
|-------|--------|-----|------------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

# Table profiling

First step to understand table



Interpretable
representation

CUSTOMER

| First | Maiden | SSN | Last Order |
|-------|--------|-----|------------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| … | … | … | … |

# Table profiling

First step to understand table



❌ Maiden has 86% NULL

✓ SSN follows 999-XX-XXXX

❌ Last Order Peaks in May

Statistical Profile
(+ Outlier Alert)

## CUSTOMER

| First | Maiden | SSN | Last Order |
|-------|--------|-----|-----------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

# Table profiling

First step to understand table

| Semantic Review | | Statistical Profile (+ Outlier Alert) |
|---|---|---|

**False Positive**

✔ Maiden for married only

❌ Maiden has <u>86% NULL</u>

**False Negative**

❌ SSN can't start with 999

✔ SSN follows <u>999-XX-XXXX</u>

**False Positive**

✔ Last Order near current

❌ Last Order Peaks in May

Semantic Review

Statistical Profile
(+ Outlier Alert)

## CUSTOMER

| First | Maiden | SSN | Last Order |
|---|---|---|---|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| … | … | … | … |

# How semantics is applied

## Previous Profiler

Statistical Profile

❌ Maiden has <u>86%</u> NULL

Table

| First | Maiden | SSN | Last Order |
|-------|--------|-----|------------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

Document

Table stores customer...

# How semantics is applied



**Previous Profiler**

Semantic Review

✔ Maiden
NULL is normal

Statistical Profile

❌ Maiden has
<u>86%</u> NULL

Table

| First | Maiden | SSN | Last Order |
|-------|--------|------|-----------|
| John | NULL | 999-45-6789 | 4/17/2024 |
| Emily | Davis | 999-65-4321 | 5/2/2024 |
| Mike | NULL | 999-78-9012 | 5/18/2024 |
| ... | ... | ... | ... |

Semantic Profile

"Maiden" is
only not NULL
for married

Semantic Context

"Maiden" is
maiden name

Document

Table stores
customer...

**Cocoon Profiler**

# LLMs for Semantic Table Profiling

Challenges

- **Many data quality issues to address**

  Duplication, Missing Values, Numeric Outliers, String Outliers...

- **Each issue needs Semantic Context, Profile and Review**

# LLMs for Semantic Table Profiling

Challenges

- **Many data quality issues to address**

  Duplication, Missing Values, Numeric Outliers, String Outliers...

- **Each issue needs Semantic Context, Profile and Review**

**Task Decomposition, but how?**

# Cocoon System Design

Interactively Verify



**1 Semantic Context**

| Table Summary |
| --- |

↓

| Column Grouping |
| --- |

for each column group

| Column Summary |
| --- |

**2 Semantic Profile & Review**

| Duplication |
| --- |

for each column group

| Column Type |
| --- |

for each column

| ... |
| --- |

for each error type

# Cocoon System Design

## Semantic Context

**1** **Semantic Context**

| Table Summary |
| :---: |
| ↓ |
| Column Grouping |

for each column group

| Column Summary |
| :---: |

**Table Summary:** Given table samples (first 5 rows) and document, describe the table in NL

**Column Grouping:** Given the Table Summary, cluster columns into groups. This is to
1. Find columns together express single concept (E.g., longitude, latitude)
2. Help human understand

**Column Summary:** Given samples for each Column Group and Table Summary, describe the group in NL

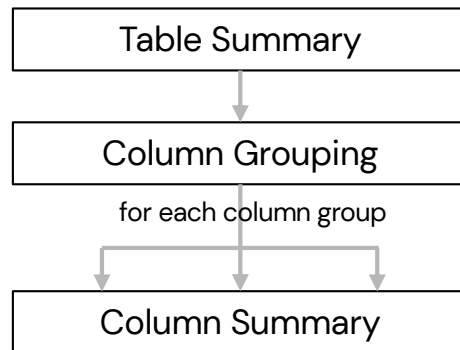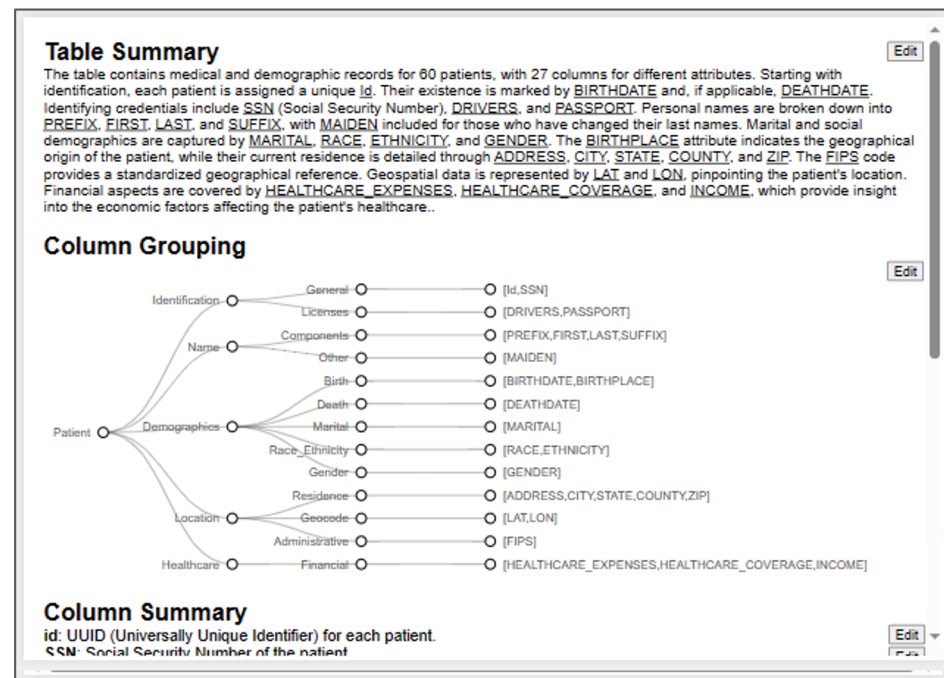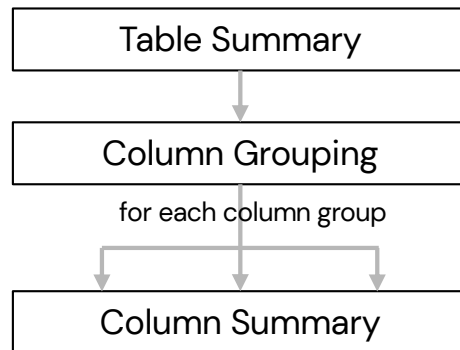# Cocoon System Design
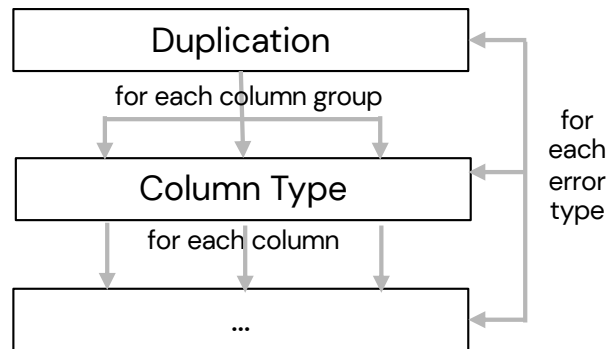
## Semantic Context



**① Semantic Context**

```
┌─────────────────────────┐
│     Table Summary       │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│    Column Grouping      │
└─────────────────────────┘
     for each column group
             │
   ┌─────────┼─────────┐
   ▼         ▼         ▼
┌─────────────────────────┐
│    Column Summary       │
└─────────────────────────┘
```

**Table Summary** [Edit]

The table contains medical and demographic records for 60 patients, with 27 columns for different attributes. Starting with identification, each patient is assigned a unique Id. Their existence is marked by BIRTHDATE and, if applicable, DEATHDATE. Identifying credentials include SSN (Social Security Number), DRIVERS, and PASSPORT. Personal names are broken down into PREFIX, FIRST, LAST, and SUFFIX, with MAIDEN included for those who have changed their last names. Marital and social demographics are captured by MARITAL, RACE, ETHNICITY, and GENDER. The BIRTHPLACE attribute indicates the geographical origin of the patient, while their current residence is detailed through ADDRESS, CITY, STATE, COUNTY, and ZIP. The FIPS code provides a standardized geographical reference. Geospatial data is represented by LAT and LON, pinpointing the patient's location. Financial aspects are covered by HEALTHCARE_EXPENSES, HEALTHCARE_COVERAGE, and INCOME, which provide insight into the economic factors affecting the patient's healthcare..

**Column Grouping** [Edit]

| | | |
|---|---|---|
| Identification | General ○ | ○ [Id,SSN] |
| | Licenses ○ | ○ [DRIVERS,PASSPORT] |
| Name | Components ○ | ○ [PREFIX,FIRST,LAST,SUFFIX] |
| | Other ○ | ○ [MAIDEN] |
| Demographics | Birth ○ | ○ [BIRTHDATE,BIRTHPLACE] |
| | Death ○ | ○ [DEATHDATE] |
| | Marital ○ | ○ [MARITAL] |
| | Race_Ethnicity ○ | ○ [RACE,ETHNICITY] |
| | Gender ○ | ○ [GENDER] |
| Location | Residence ○ | ○ [ADDRESS,CITY,STATE,COUNTY,ZIP] |
| | Geocode ○ | ○ [LAT,LON] |
| | Administrative ○ | ○ [FIPS] |
| Healthcare | Financial ○ | ○ [HEALTHCARE_EXPENSES,HEALTHCARE_COVERAGE,INCOME] |

**Column Summary**

**id:** UUID (Universally Unique Identifier) for each patient. [Edit]
**SSN:** Social Security Number of the patient

**Interface for human feedback**

# Cocoon System Design

## Semantic Profile & Review

**② Semantic Profile & Review**

| Duplication |
|---|
*for each column group*

| Column Type |
|---|
*for each column*
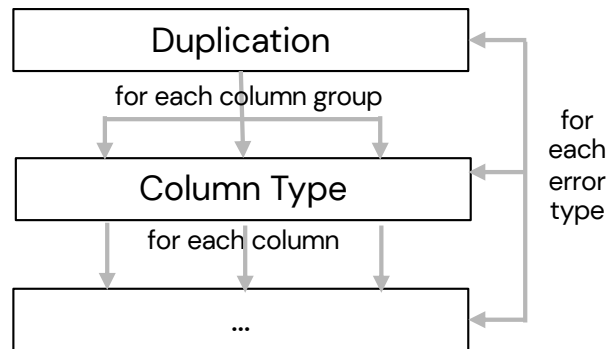
| ... |
|---|

*for each error type*

We consider errors:
- Duplication
- Column Type
- Uniqueness
- Disguised Missing Value
- Missing Value
- Numeric Outliers
- String Outliers
- Missing Record

# Cocoon System Design

## Semantic Profile & Review

**②  Semantic Profile & Review**

```
┌─────────────────────────────┐ ◄──────┐
│        Duplication          │        │
└─────────────────────────────┘        │
      for each column group            │
   ↓         ↓          ↓               │
┌─────────────────────────────┐ ◄──────┤  for
│        Column Type          │        │  each
└─────────────────────────────┘        │  error
      for each column                  │  type
   ↓         ↓          ↓               │
┌─────────────────────────────┐ ◄──────┘
│             ...             │
└─────────────────────────────┘
```
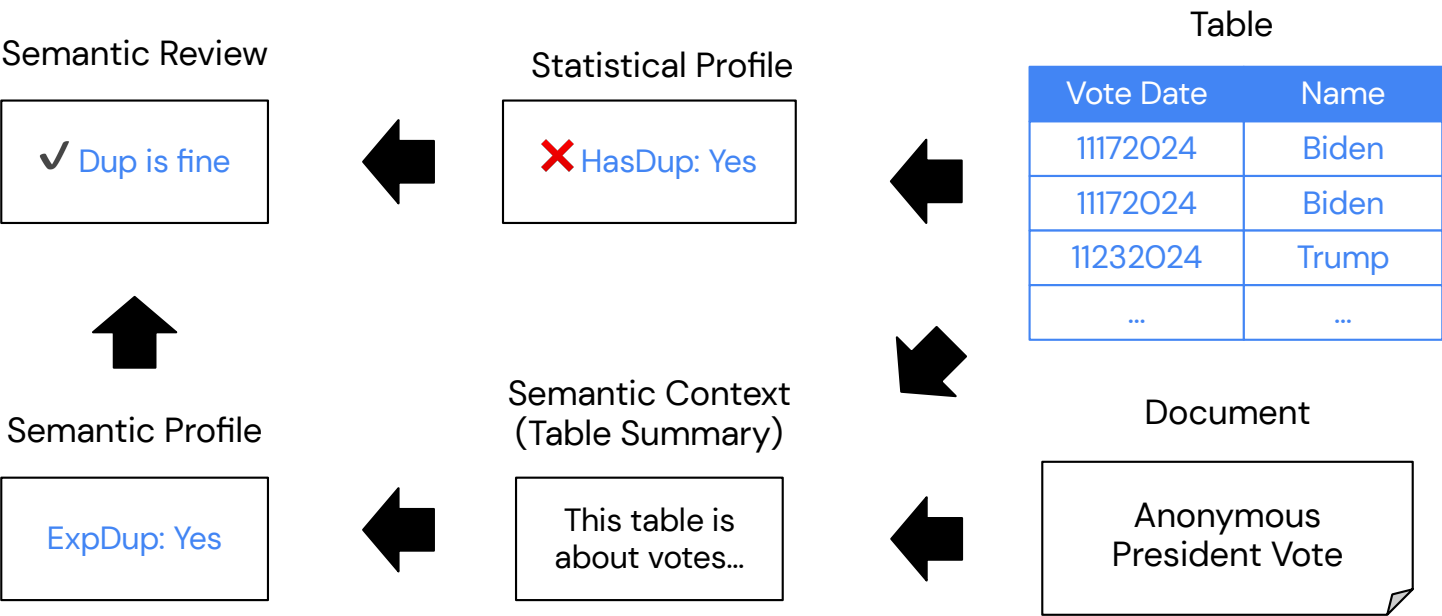
We consider errors:
- Duplication    ⬅
- Column Type    ⬅
- Uniqueness
- Disguised Missing Value
- Missing Value
- Numeric Outliers
- String Outliers
- Missing Record

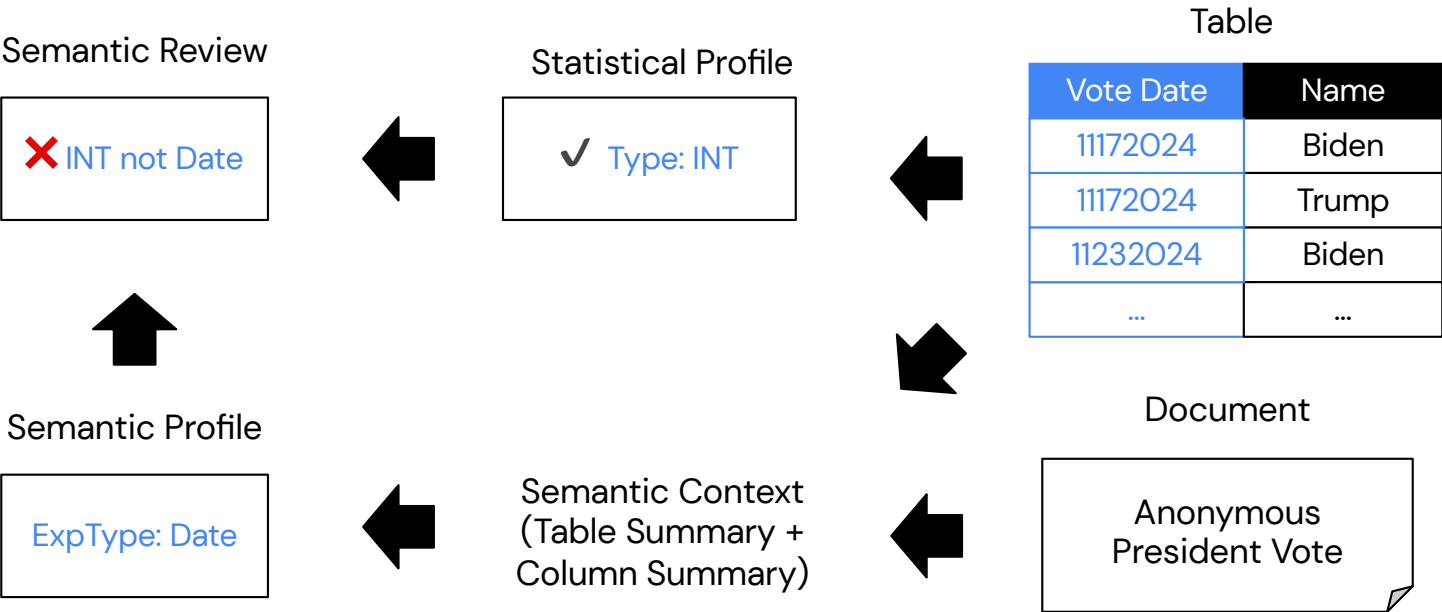**Design:** What are the Semantic Profile & Review

# Cocoon System Design

## Semantic Profile & Review for Duplication

**Semantic Review**

☑ Dup is fine

**Statistical Profile**

✖ HasDup: Yes

**Table**

| Vote Date | Name |
|-----------|-------|
| 11172024 | Biden |
| 11172024 | Biden |
| 11232024 | Trump |
| ... | ... |

**Semantic Profile**

ExpDup: Yes

**Semantic Context (Table Summary)**

This table is about votes...

**Document**

Anonymous President Vote

# Cocoon System Design

## Semantic Profile & Review for Data Type

**Semantic Review**

| ✘ INT not Date |
|---|

**Statistical Profile**

| ✔ Type: INT |
|---|

**Table**

| Vote Date | Name |
|---|---|
| 11172024 | Biden |
| 11172024 | Trump |
| 11232024 | Biden |
| ... | ... |

**Semantic Profile**

| ExpType: Date |
|---|

Semantic Context
(Table Summary +
Column Summary)

**Document**

Anonymous
President Vote

# Cocoon System Design

## Final Output

# Conclusion

- Statistical profiling has false positives/negatives

- LLMs help interactively improve profiling with semantics

- Task decomposition improves quality