

# HETEROGENEOUS PATTERN SAMPLING ACCORDING TO FREQUENCY

Rayane Lachache<sup>1</sup>, Djawad Bekkoucha<sup>1,3</sup>, Abdelkader Ouali<sup>1</sup>, Bruno Crémilleux<sup>1</sup>, Thi-Bich-Hanh Dao<sup>2</sup>, and Christel Vrain<sup>2</sup>

<sup>1</sup>University of Caen Normandy, CNRS, UMR 6072 GREYC, 14032 Caen, France

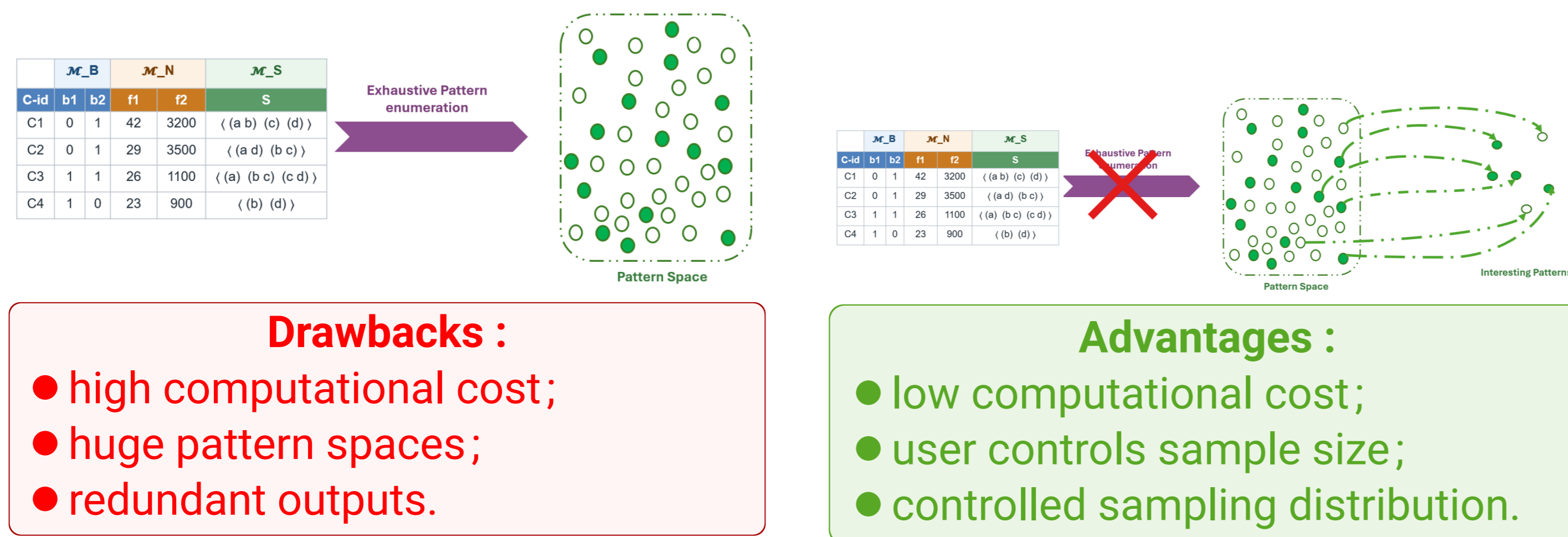
<sup>2</sup>University of Orléans, INSA Centre Val de Loire, LIFO EA 4022, Orléans, France

<sup>3</sup>LISN, CNRS (UMR 9015), Paris Saclay University, Gif-sur-Yvette, France

rayane.lachache@unicaen.fr



## 1. Motivation



### Drawbacks :

- high computational cost;
- huge pattern spaces;
- redundant outputs.

### Advantages :

- low computational cost;
- user controls sample size;
- controlled sampling distribution.

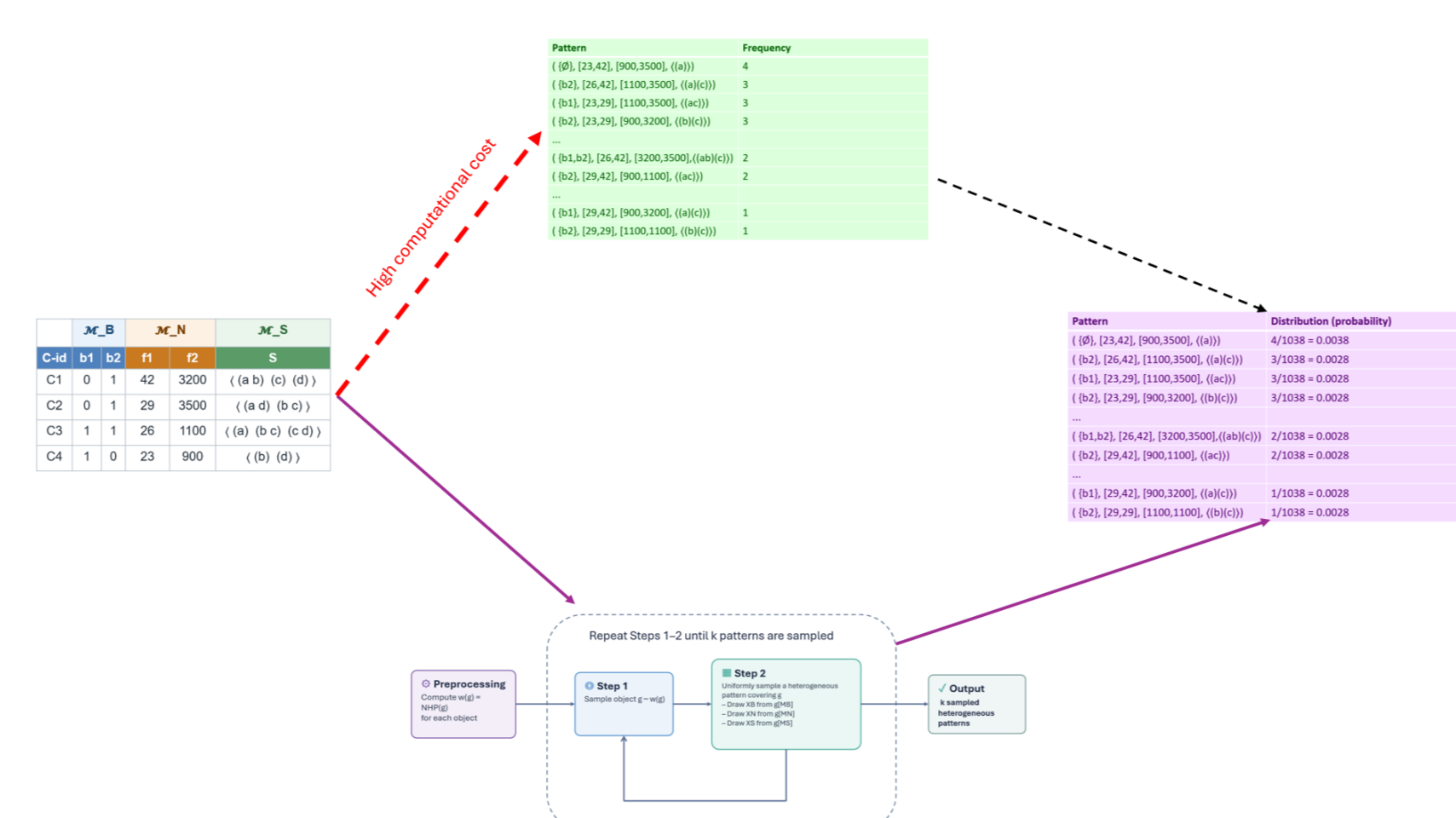
## 2. Problem Statement

### Sampling goal :

draw a heterogeneous pattern proportionally to its frequency.

$$P(X_H) = \frac{\text{freq}(X_H, H)}{\sum_{X'_H \in \mathcal{L}_H} \text{freq}(X'_H, H)}$$

$$X_H = \langle X_B, X_N, X_S \rangle$$



## 3. State of the Art

### Stochastic methods :

- **Graphs** : [El Hasan et al., VLDB 2009], [Bendimerad et al., ICDM 2016]
- **Formal concepts** : [Boley et al., DSM 2010]
- **Tiles** : [Bendimerad et al., IDA 2020]

### Declarative method :

- **Itemsets** : [Dzyuba et al., DAMI 2017]

### Multi-step methods :

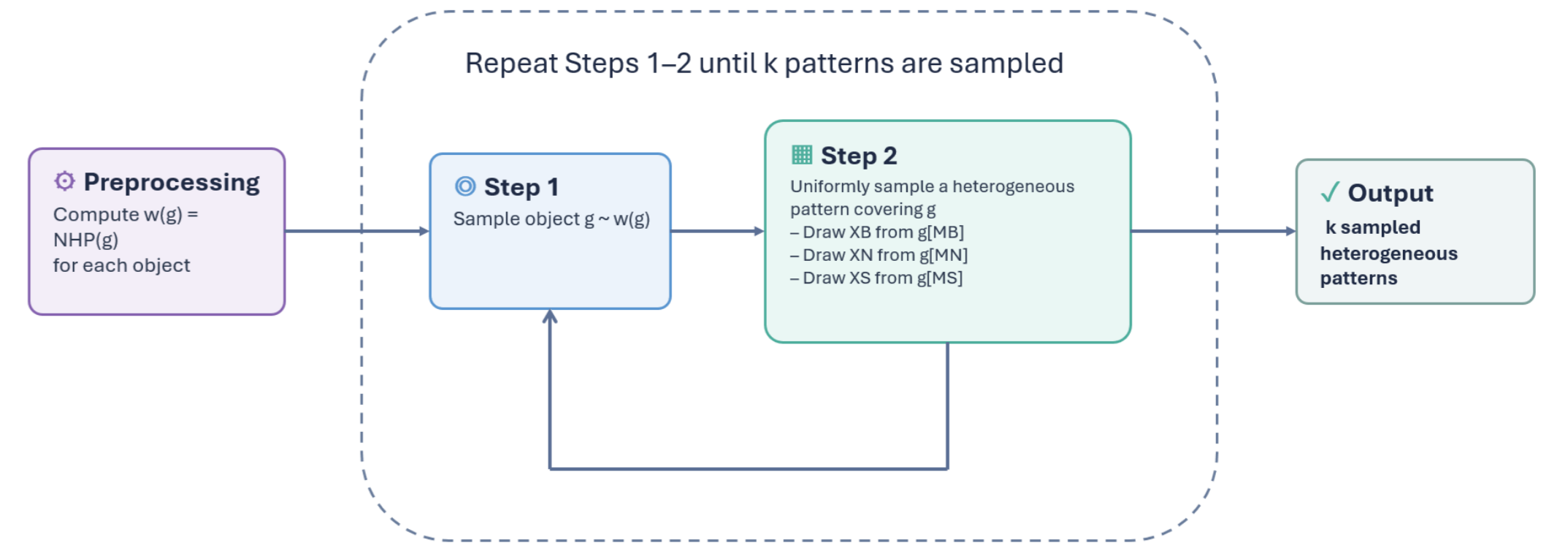
- **Itemsets** : [Boley et al., KDD 2011]
- **Sequences** : [Diop et al., ICDM 2018]
- **Numerical data** : [Giacometti et al., SDM 2018], [Bekkoucha et al., DKE 2026]

**Our contribution** : SEHP is the first frequency-based sampling method for heterogeneous tabular data.

**Published work** : [Lachache et al., IDA 2026] *Heterogeneous Pattern Sampling according to Frequency*.

## 4. SEHP Method

### Sampling Heterogeneous Patterns Proportionally to Frequency



### Preprocessing : compute one weight per object

$$w(g) = \text{NHP}(g) \quad \text{NHP}(g) = \text{NI}(g) \times \text{NIP}(g) \times \text{NSP}(g)$$

C-id	M_B		M_N		M_S	W(g)
	b1	b2	f1	f2	S	
C1	0	1	42	3200	\langle (a b) (c) (d) \rangle	12
C2	0	1	29	3500	\langle (a d) (b c) \rangle	18
C3	1	1	26	1100	\langle (a) (b c) (c d) \rangle	24
C4	1	0	23	900	\langle (b) (d) \rangle	46

### Step 1

Draw an object  $g$  :

$$P(g) = \frac{w(g)}{\sum_{g' \in G} w(g')}$$

C-id	M_B		M_N		M_S	W(g)
	b1	b2	f1	f2	S	
C1	0	1	42	3200	\langle (a b) (c) (d) \rangle	12
C2	0	1	29	3500	\langle (a d) (b c) \rangle	18
C3	1	1	26	1100	\langle (a) (b c) (c d) \rangle	24
C4	1	0	23	900	\langle (b) (d) \rangle	46

Step 1 draw one object proportionally to its weight

### Step 2

Uniformly draw a pattern covering  $g$  :

$$X_H = \langle X_B, X_N, X_S \rangle$$

C-id	M_B		M_N		M_S	W(g)
	b1	b2	f1	f2	S	
C1	0	1	42	3200	\langle (a b) (c) (d) \rangle	12
C2	0	1	29	3500	\langle (a d) (b c) \rangle	18
C3	1	1	26	1100	\langle (a) (b c) (c d) \rangle	24
C4	1	0	23	900	\langle (b) (d) \rangle	46

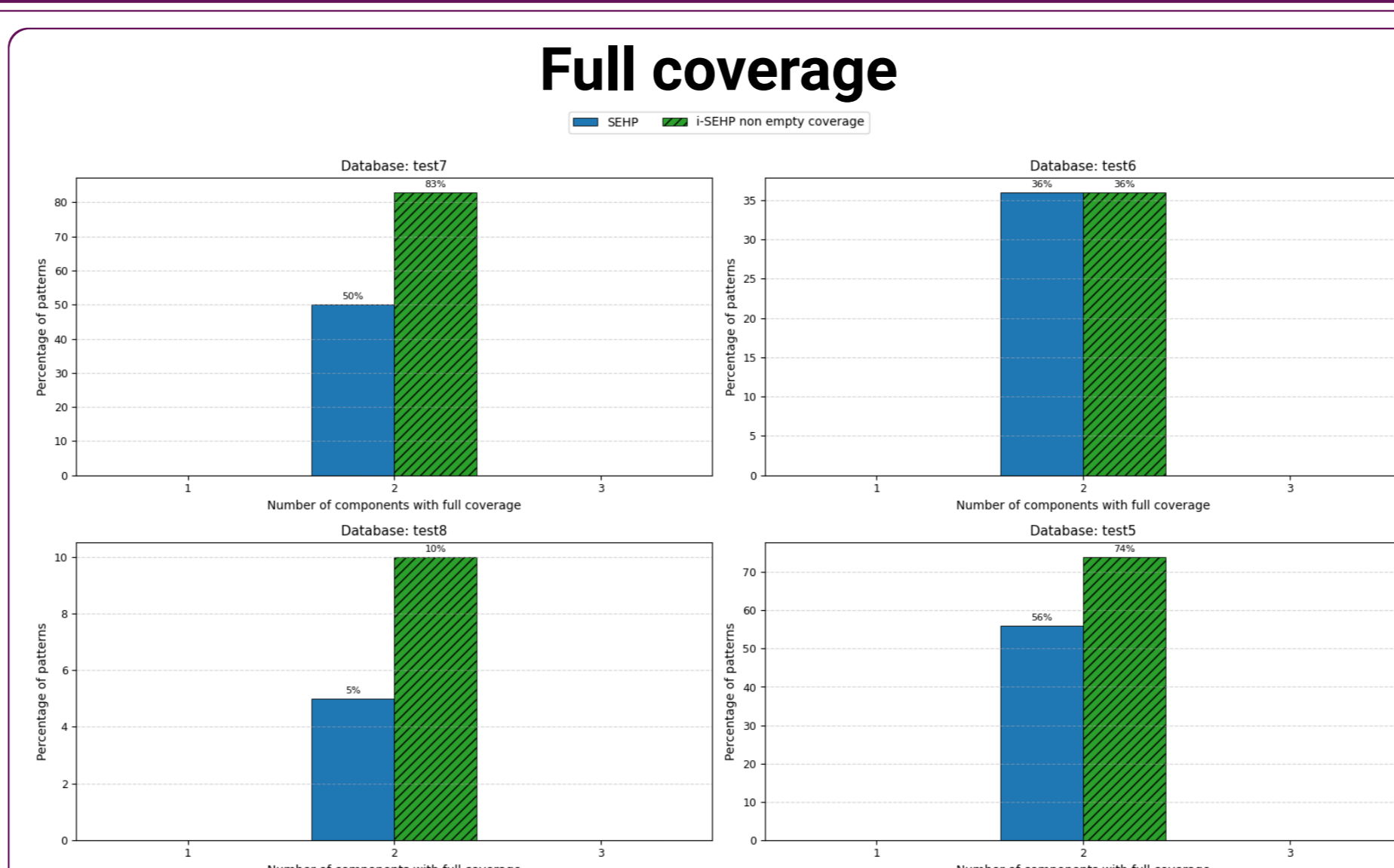
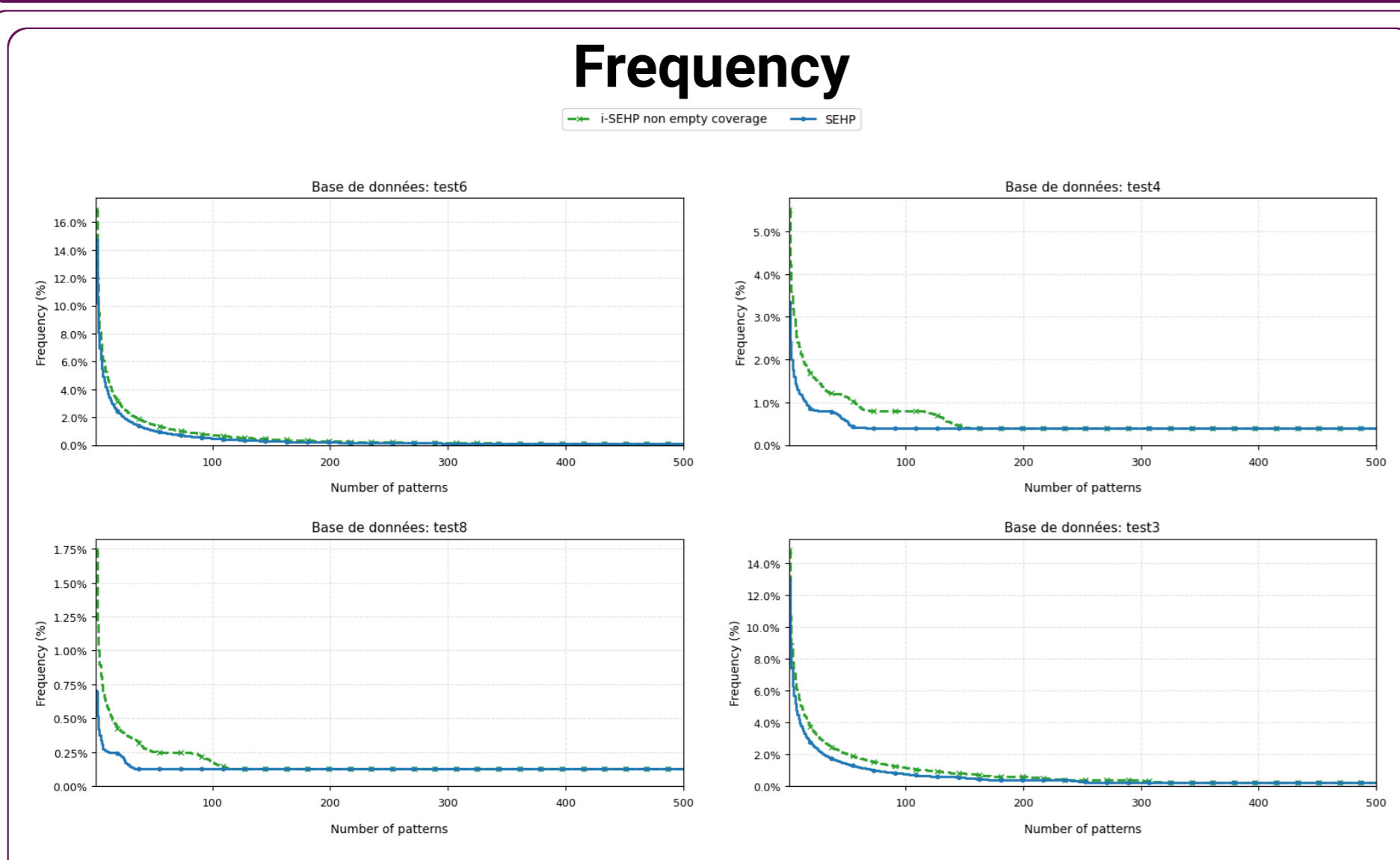
Step 2: uniform draw of a pattern covering  $g$   
 $X_H = \langle (b), (23, 29), (900, 1100) \rangle, \langle (b) (d) \rangle$

## 5. Theoretical Guarantee

$$Z = \sum_{g \in G} \text{NHP}(g) \implies P(X_H) = \sum_{g \in \text{cover}(X_H, H)} \frac{\text{NHP}(g)}{Z} \frac{1}{\text{NHP}(g)} = \frac{\text{freq}(X_H, H)}{Z}$$

✓ SEHP samples heterogeneous patterns proportionally to their frequency.

## 6. Experimental Results



### Baseline : I-SEHP

I-SEHP independently samples :

$$X_B \sim \text{freq}(X_B), \quad X_N \sim \text{freq}(X_N), \quad X_S \sim \text{freq}(X_S)$$

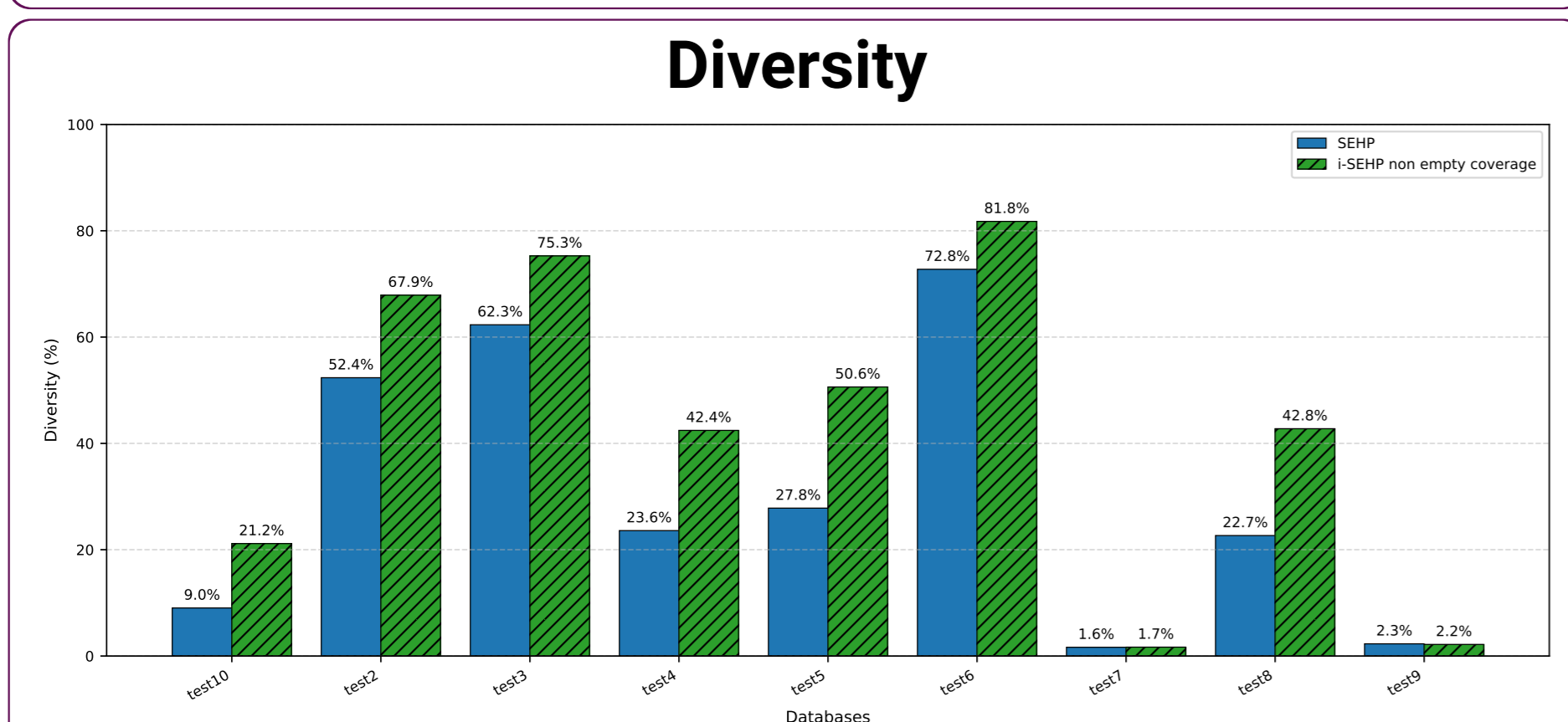
Then it aggregates them :

$$X_H = \langle X_B, X_N, X_S \rangle$$

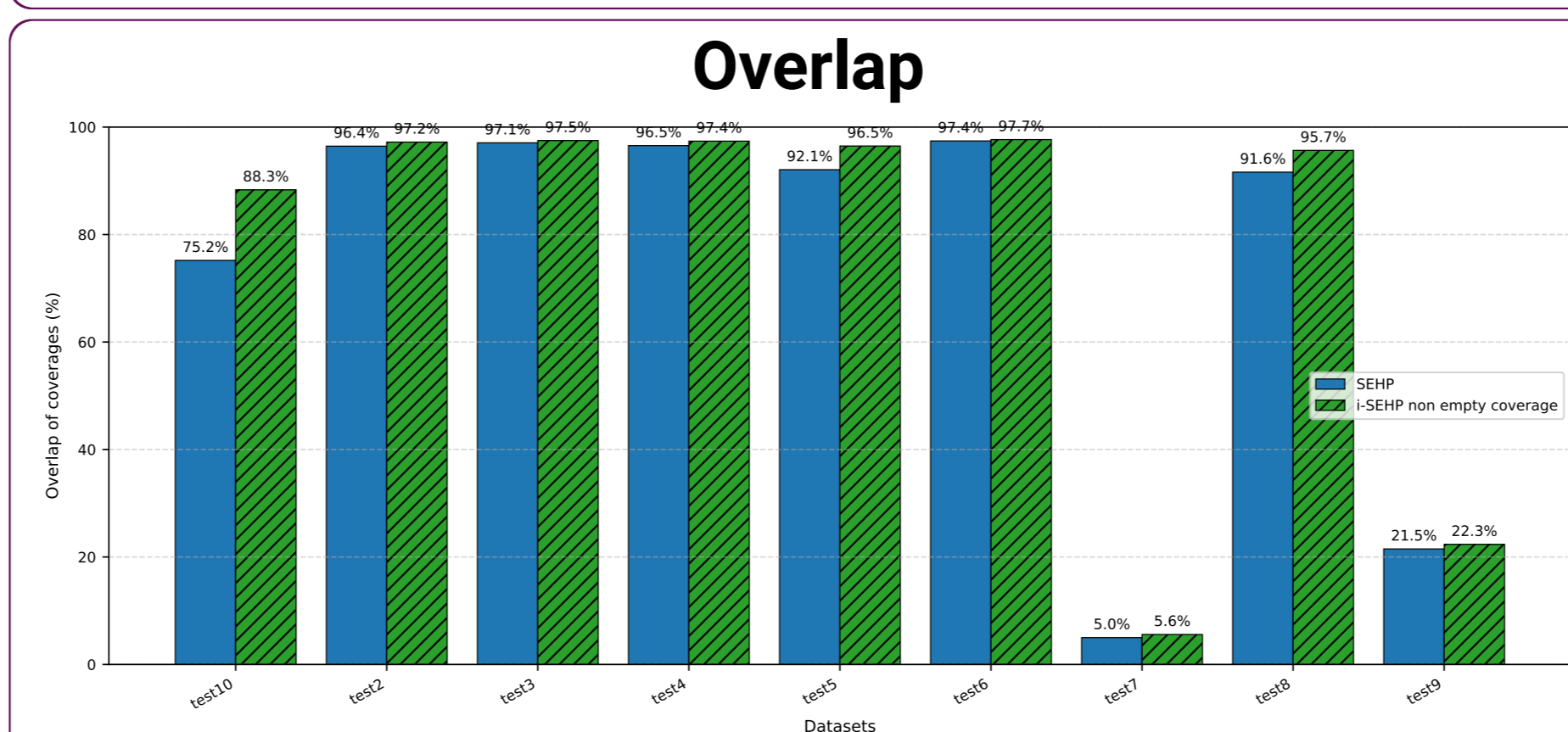
with coverage

$$\text{cover}(X_H, H) = \text{cov}_B \cap \text{cov}_N \cap \text{cov}_S$$

Independent draws may produce **empty-cover** patterns, requiring rejection.



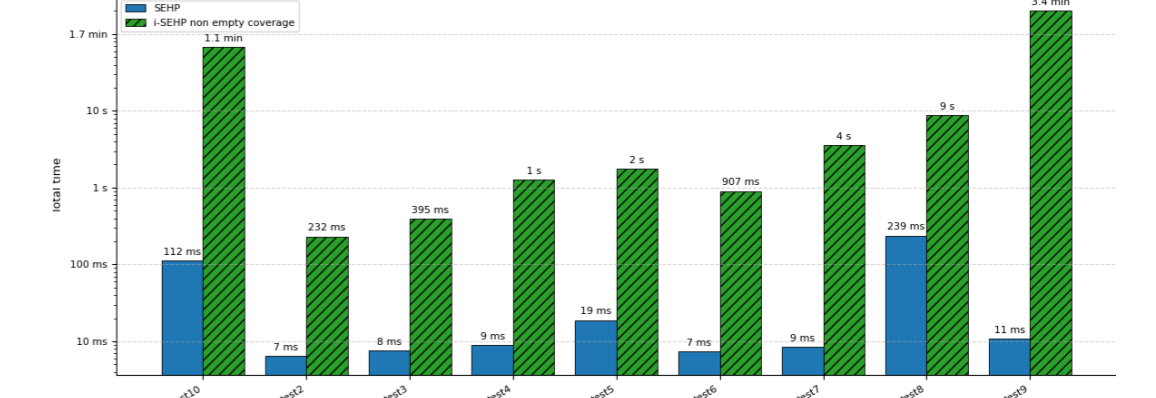
$$\text{Diversity}(K, H) = \frac{|\{\text{cover}(X_H^1, H), \dots, \text{cover}(X_H^K, H)\}|}{|K|}$$



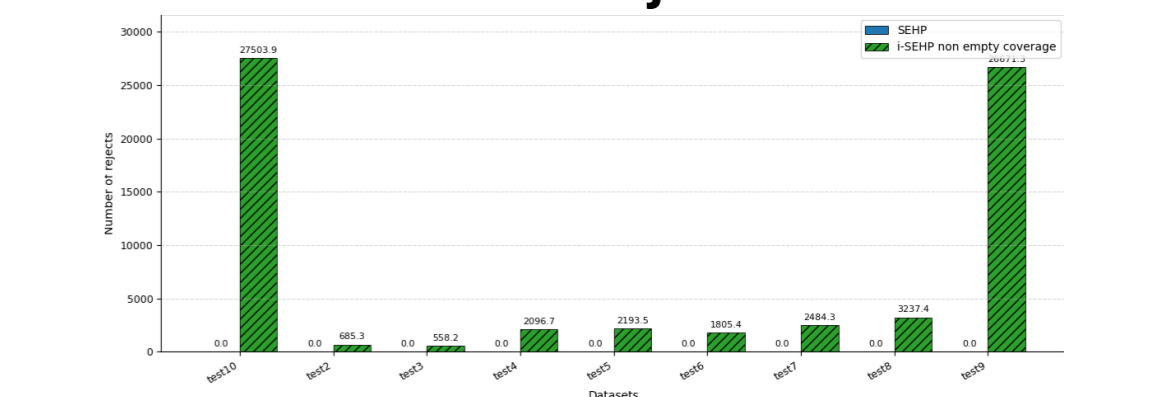
$$\text{Overlap}(K, H) = \frac{2}{|K|(|K|-1)} \sum_{1 \leq i < j \leq |K|} \left( 1 - \frac{|\text{cover}(X_H^i, H) \cap \text{cover}(X_H^j, H)|}{|\text{cover}(X_H^i, H) \cup \text{cover}(X_H^j, H)|} \right)$$

### Runtime analysis

#### CPU time



#### Number of rejections



### Main messages

- ✓ SEHP avoids empty-cover patterns by construction.
- ✓ SEHP requires no rejection.
- ✓ SEHP is faster on most datasets.
- ✗ I-SEHP may be slowed down by repeated rejection.
- ✗ Both methods are affected by the long-tail phenomenon.
- ✗ I-SEHP produces more full-coverage sub-patterns than SEHP.
- ✓ SEHP provides a formal frequency-based guarantee.
- ✓ Diversity and overlap characterize exploration of the pattern space.

## 7. Future Directions

- Extending SEHP to richer interestingness measures beyond frequency;
- Integrating user-defined constraints into the heterogeneous sampling process;
- Leveraging sampled heterogeneous patterns as features for supervised and unsupervised learning;
- Embedding SEHP within interactive pattern mining workflows.