

# Extracting Novel Ideas From Crowdsourcing Requirements

R·I·T

B. THOMAS GOLISANO  
College of COMPUTING AND  
INFORMATION SCIENCES

Diem Dao

Advisor: Dr. Pradeep Murukannaiah

## Background

- When researchers and businesses crowdsource data from the public, a large database of information can be overwhelming and redundant.
- Finding novelty ideas from a large data set can be challenging.

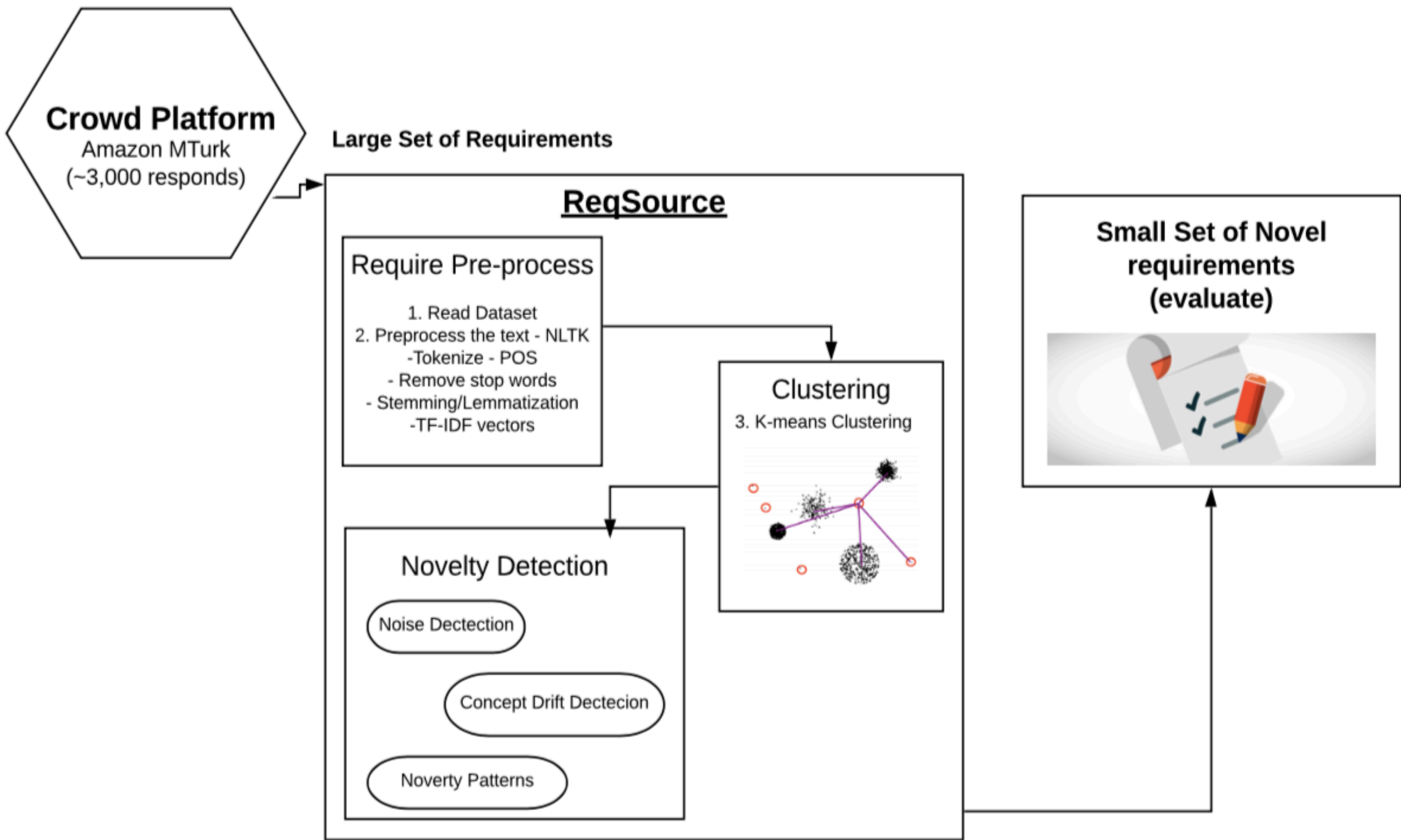
Hence, we developed a study that:

- Can automatically determine novel ideas from crowdsourcing requirements.
- Collected approximately 3,000 requirements from a MTurk survey on what people wanted in their smart home applications.
- Filtered the mundane ideas from novel suggestions (requirements).

## Current Progress

- Carried out data pre-processing techniques using NLTK (Tokenization, Stop word removal, POS Tagging, Stemming from textual requirements obtained through crowdsourcing).
- Used data mining techniques such as TF-IDF Vectorization and Count Vectorization to understand and analyze the data.
- Implemented K means algorithm and used the Elbow method with Euclidean distance to determine the optimal number of clusters.
- Find the requirements in each cluster.

Goal: Develop tools and techniques to extract novel ideas and requirements from crowdsourced data.



## Direction

- Implement MINAS: multiclass learning algorithm for novelty detection in data streams.
- Create a decision model by training and testing the data set to detect novelty patterns.
- Remove outliers or noise and find concept drifting.

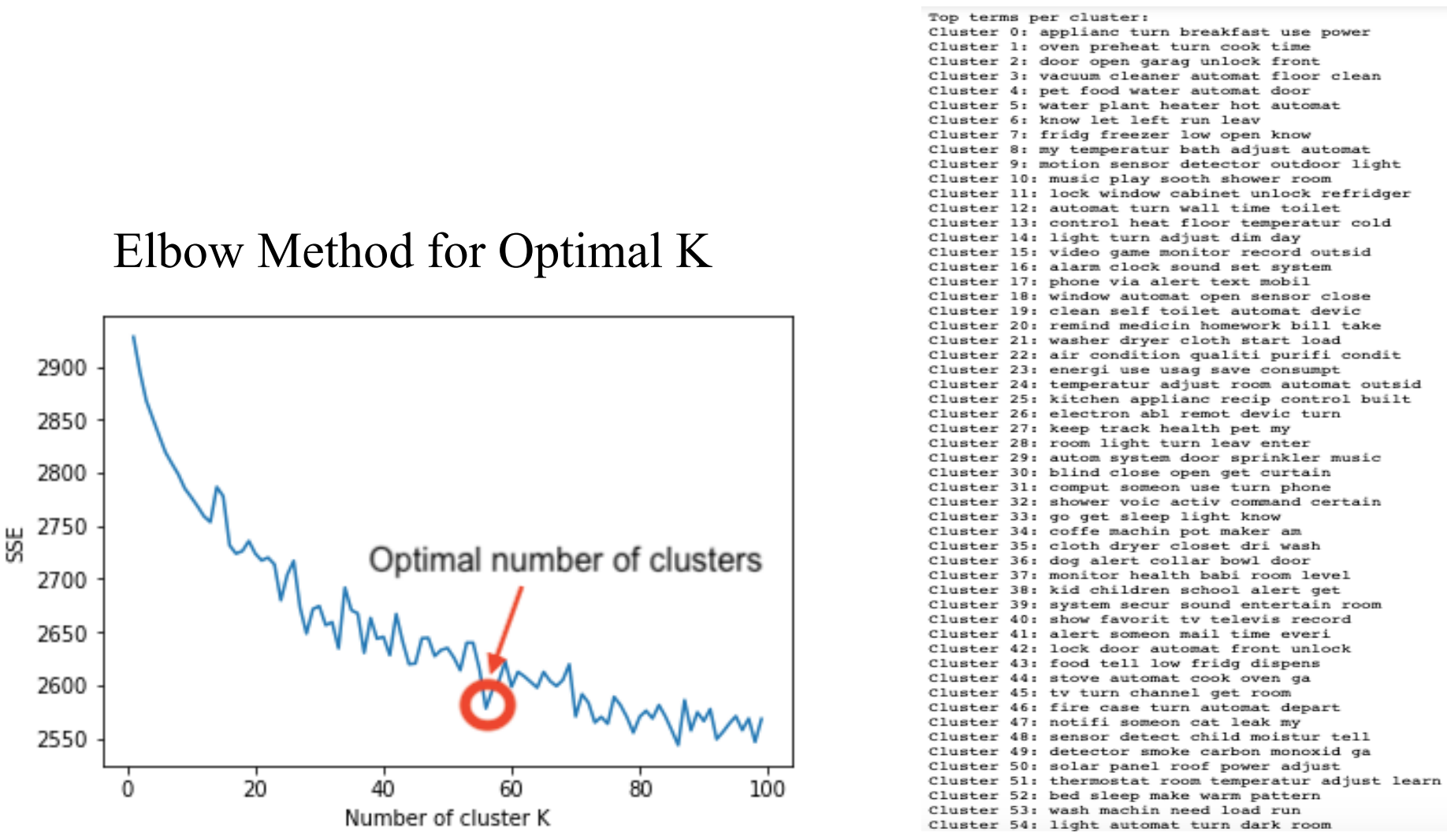


Fig. 1 shows the optimal number of K clusters



Fig. 2 shows the 55 most common ideas and concepts from all the requirements

cluster	data_index	text
4	1	4 class, remind
152	1	152 remind pre, manag, list, day, night
516	1	516 remind activ, kid, need, go
773	1	773 set, remind, start, laundri, alert, laundri, done
799	1	799 make, recommend, date, new, movi, may, enjoy, ...
848	1	848 remind day, everi, morn, get
879	1	879 remind leav, left, food, water, cat
987	1	987 remind exercis, track, progress
1010	1	1010 A, bill, remind
1011	1	1011 auto, featur, remind chang, water, softer
1099	1	1099 remind homework

Fig. 3 shows what requirements are in each cluster (subset of a group). This cluster shows concepts relevant to the topic of reminders.

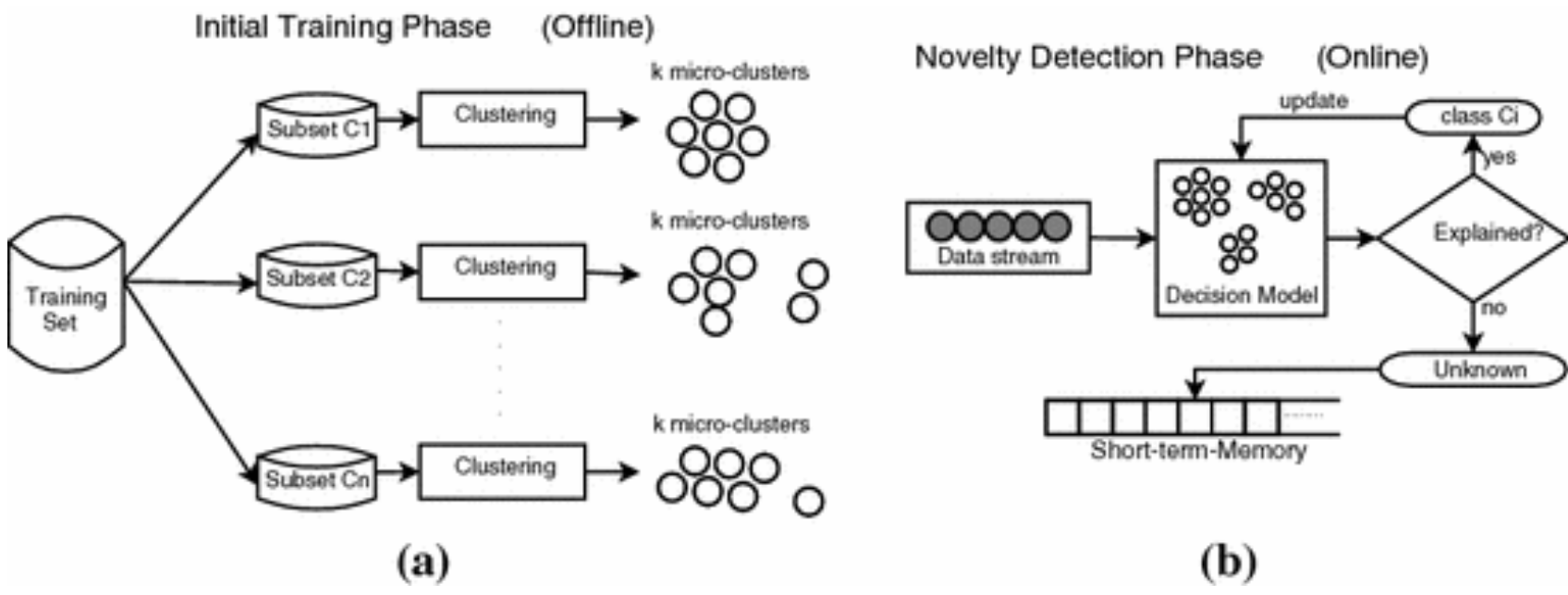


Fig. 4 MINAS algorithm in 2 phases: Initial Training Phase (Offline) and Novelty Detection Phase (Online)

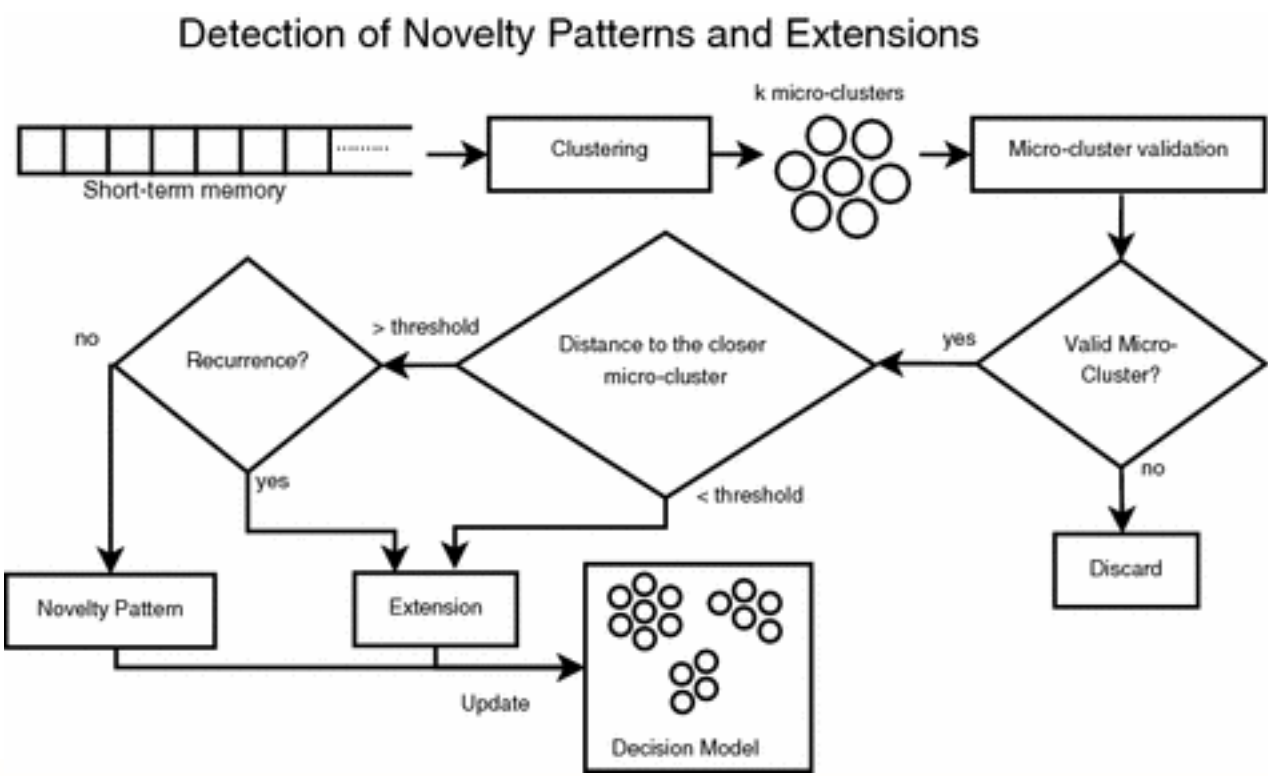


Fig. 5 Novelty Detection Process

## References:

de Faria, E.R., Ponce de Leon Ferreira Carvalho, A.C. & Gama, J. Data Min Knowl Disc (2016) 30: 640. <https://doi.org/10.1007/s10618-015-0433-y>

## Acknowledgement

This work is supported by NSF CNS-1757680.

Pradeep K. Murukannaiah, Ph.D.  
Contact: Assistant Professor  
Department of Software Engineering  
Rochester Institute of Technology  
pkmvse@rit.edu

