

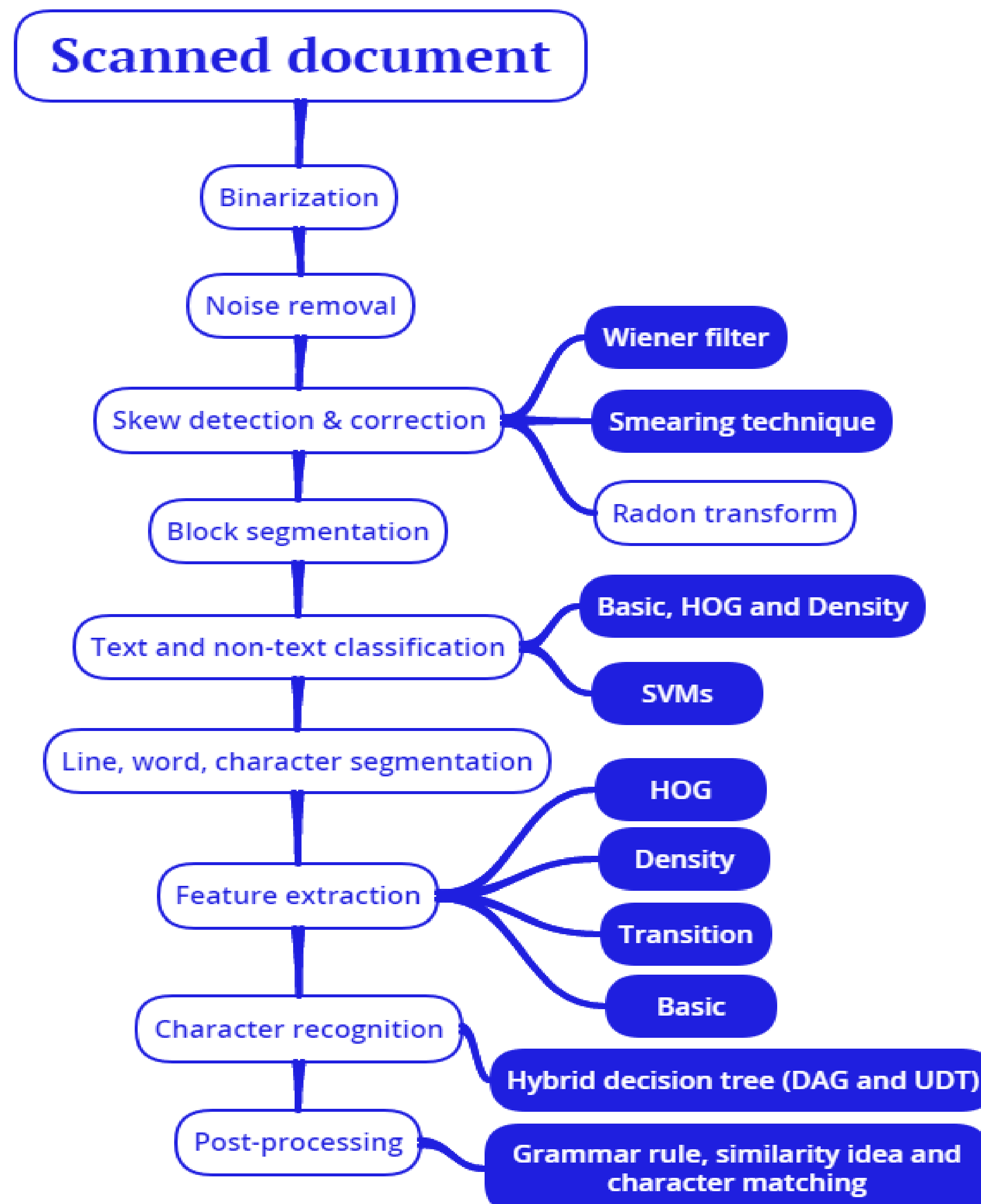
- Scanned images are needed to be in editable format in order to perform the following:
Searching, Copying, Editing, Grammar checking, Spell checking, Formatting, etc.
- Why Tamil OCR?
To preserve or reproduce ancient or decrepit Tamil books.

1. Printed Tamil documents poses challenges owing to:
One line may have different font styles, presence of pictures, multi columns, etc.
2. Existing Tamil OCRs show moderate recognition rate.
Not evaluated on a common or standard datasets due to the absence of a standard dataset.

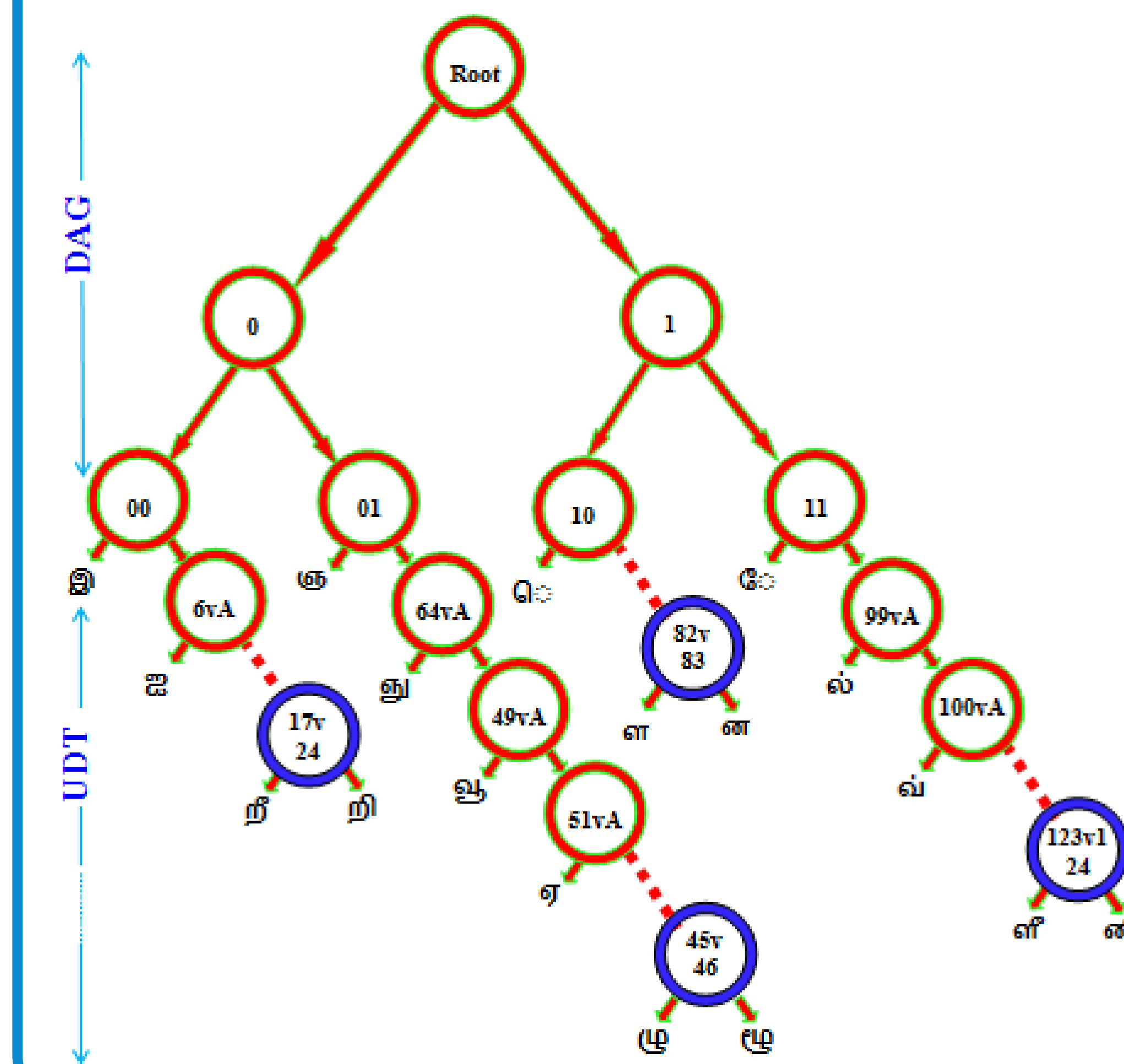
- To improve the recognition rate for printed Tamil text.
- To provide a standard and challenging dataset to the research community consisting scanned printed Tamil text.

- Created a dataset for the research community working on printed Tamil text.
- Proposed a skew detection and textual classification using Wiener filter, smearing and Radon transform methods.
- Proposed a hybrid decision tree using SVMs for character recognition of Tamil text.
- Proposed a post-processing error correction technique.

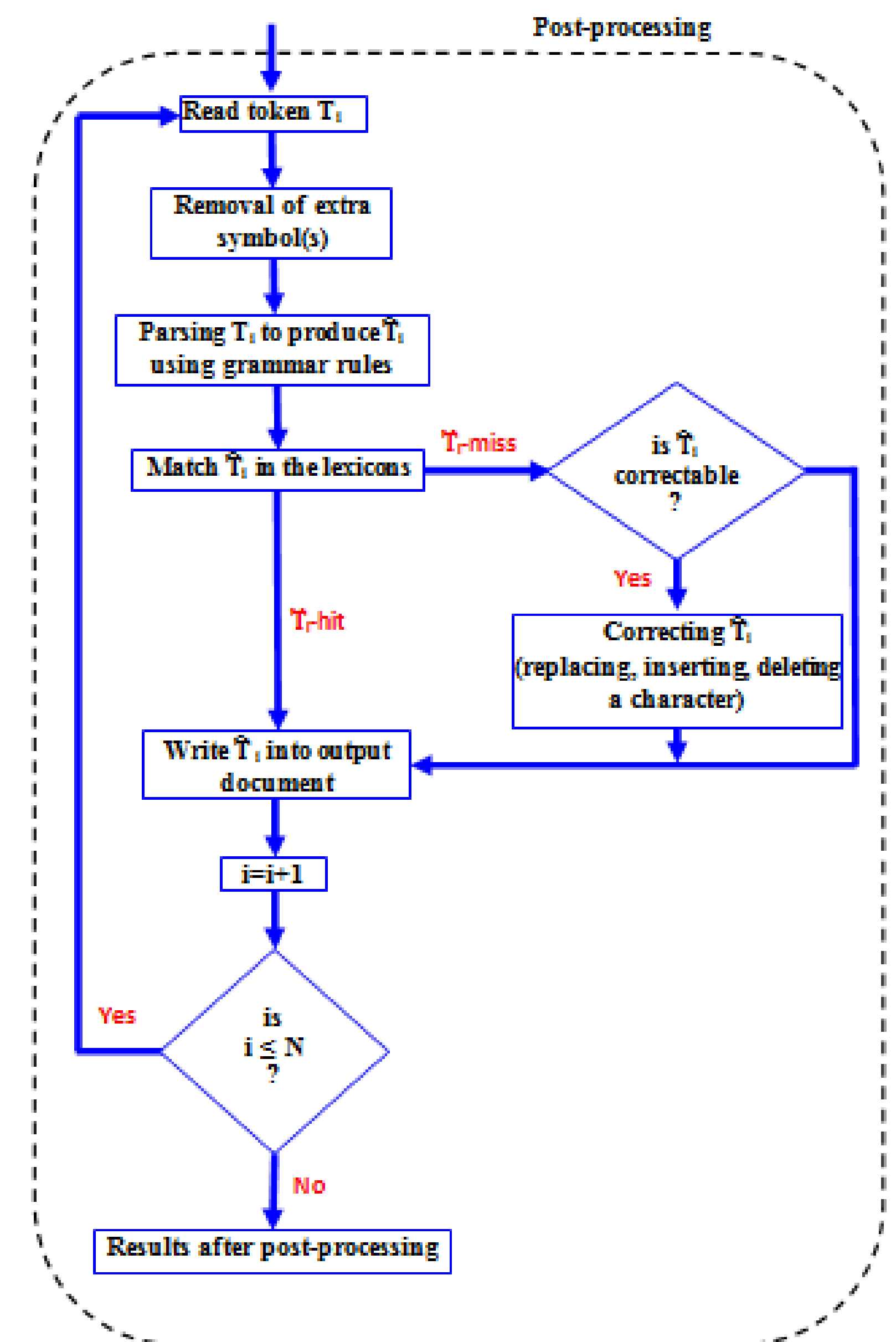
The system consists of pre-processing, feature extraction, character recognition and post-processing.



Features used in this work: Basic (6 dimension), density (89 dimension), HOG (1154 dimension) and transition features (372 dimension). The extracted feature vectors are analysed using the novel hybrid decision tree of DAG [2] and UDT-SVMs [1].



The flowchart of the proposed post-processing technique consists of various steps.



1. Binarisation:
Gray scale image \Rightarrow binary image.
2. Noise removal:
Applying median filter to the binarised image.
3. Skew correction and detection:
Using Wiener filter, smearing and Radon transform technique.
4. Block segmentation:
Using run length smearing algorithm (RLSA) and connected component analysis.
5. Text or non text classification:
Feature vectors (Basic, HOG and density) are analysed using SVMs to classify text or non-text block.
6. Character segmentation: Using projection technique and connected component analysis

- Characters were grouped by applying K-means algorithm to find the root node of the hybrid decision tree.
- The same technique was applied to the next level of the hybrid decision tree.

[illegible]

- Each decision node of the UDT is found based on the features that best separates the characters in recognition and the order of nodes is fixed in the decision tree.
- This process is followed for every decision node of UDT.

Four diverse types of printed Tamil documents: books, magazines, newspapers and pamphlets. Five different examples/document; \Rightarrow four pages / examples \Rightarrow Total 80 examples [3].

- The recognition rate for skew correction is 95.83%
- The character recognition Proposed for hybrid decision tree is 98.80%.
- The error correction technique reduces the overall average error rate by nearly 5% in the output produced by the Tamil OCR.

1. Ramanan *et al.*, "Unbalanced Decision Trees for Multi-class Classification", 2007.
2. Platt *et al.*, "Large Margin DAGs for Multiclass Classification", 2000.
3. Tamil Digitising Project, Department of Computer Science, University of Jaffna, Sri Lanka, 2014. <http://www.csc.jfn.ac.lk/tdp>