

SPEAKER RECOGNITION SYSTEM



K. ROHITH MANIKANTA

N. CHERISH

Y. HARSHA VARDHAN

22EC01010

22EC01016

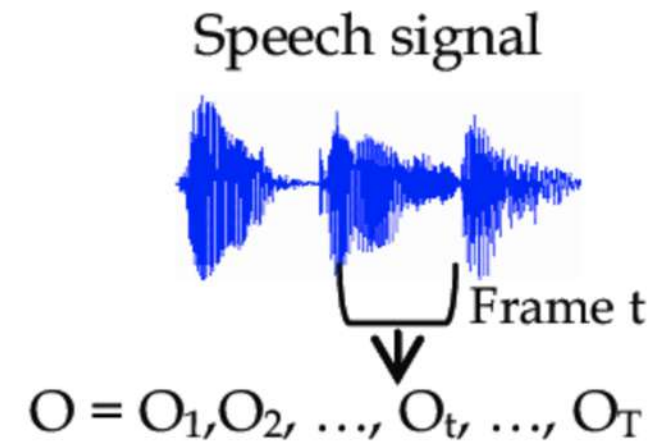
22EC01007



Speaker Recognition System

Final Implementation

- Developed a Speaker Recognition System using MFCC and Pitch features.
- Implemented using MATLAB with support for different windowing functions.
- Trained 3 models based on:
 1. Rectangular Window
 2. Hanning Window
 3. Hamming Window
- Classification performed using Support Vector Machines (SVM).
- Feature normalization applied to enhance consistency and performance.



Windowing :

$$y_t(n) = x_t(n)w(n), 0 \leq n \leq N-1$$

Fourier Transform :

$$X_n = \sum_{k=0}^{N-1} x_k e^{-2\pi jkn/N}$$

Mel Frequency Wrapping by using M filters
For each filter, compute i^{th} mel spectrum, X_i :

$$X_i = \log_{10} \left(\sum_{k=0}^{N-1} |X(k)| H_i(k) \right), i=1, 2, 3, \dots, M$$

$H_i(k)$ is i^{th} triangle filter

Compute the J cepstrum coefficients using Discrete Cosine Transform

$$C_j = \sum_{i=1}^M X_i \cos \left(j(i-1) / 2 \frac{\pi}{M} \right)$$

$j=1, 2, 3, \dots, J$; J=number of coefficients

Windowing in Feature Extraction

- Frame blocking followed by windowing improves feature stability.
- Window functions applied:

Rectangular

Simple shape, wider main lobe, higher side lobes cause more spectral leakage.

Hamming

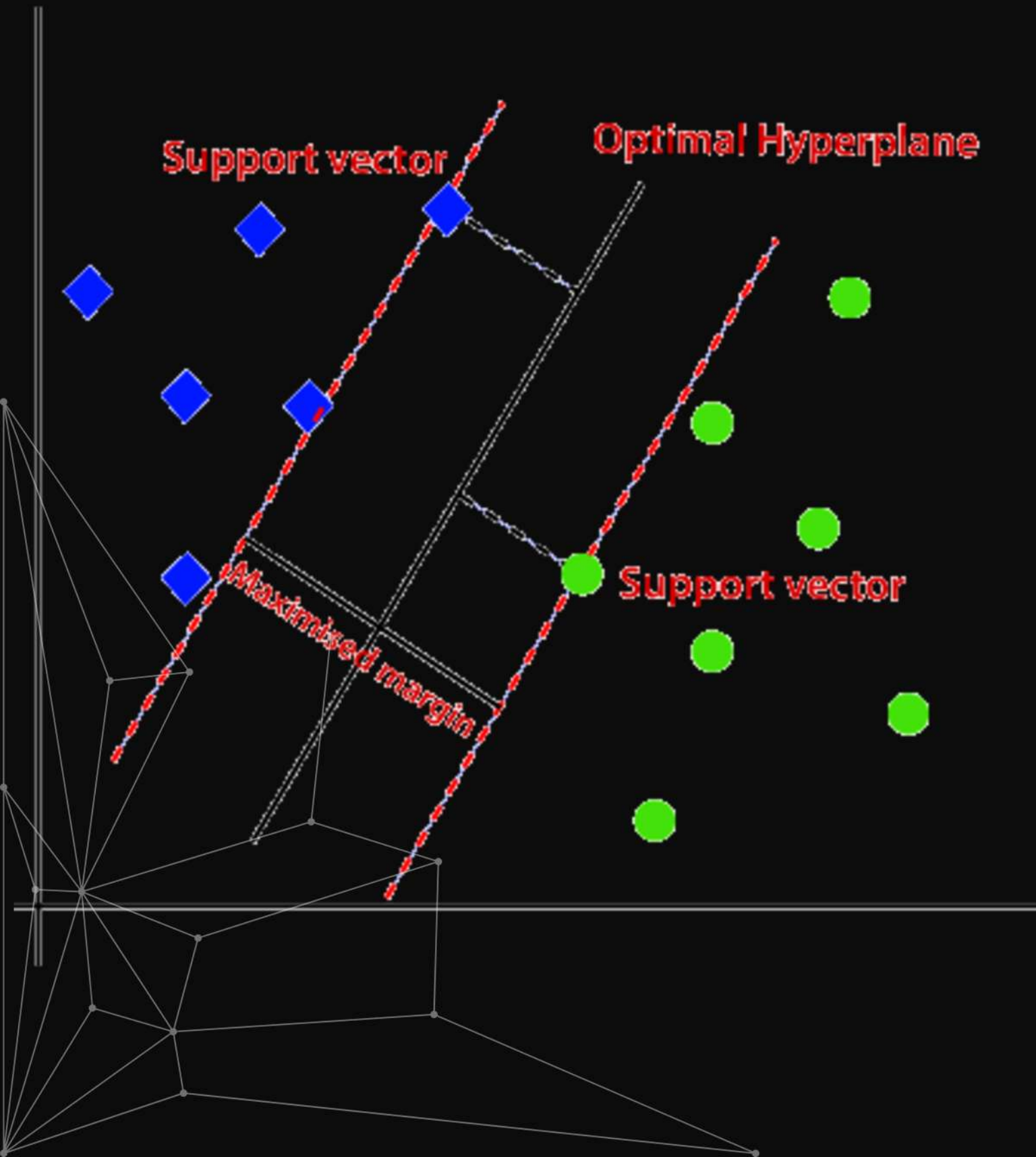
Smother taper reduces side lobes, improving signal-to-noise ratio by ~5% over rectangular.

Hanning

Even smoother taper, better side lobe attenuation, lowers spectral leakage for cleaner features.

- Speech signals are non-stationary. So instead of analyzing the whole signal at once, we break it into short segments (30 ms).
- Applying a window helps focus on that segment, reducing the effect of the rest.

Machine Learning: Classification Using SVM

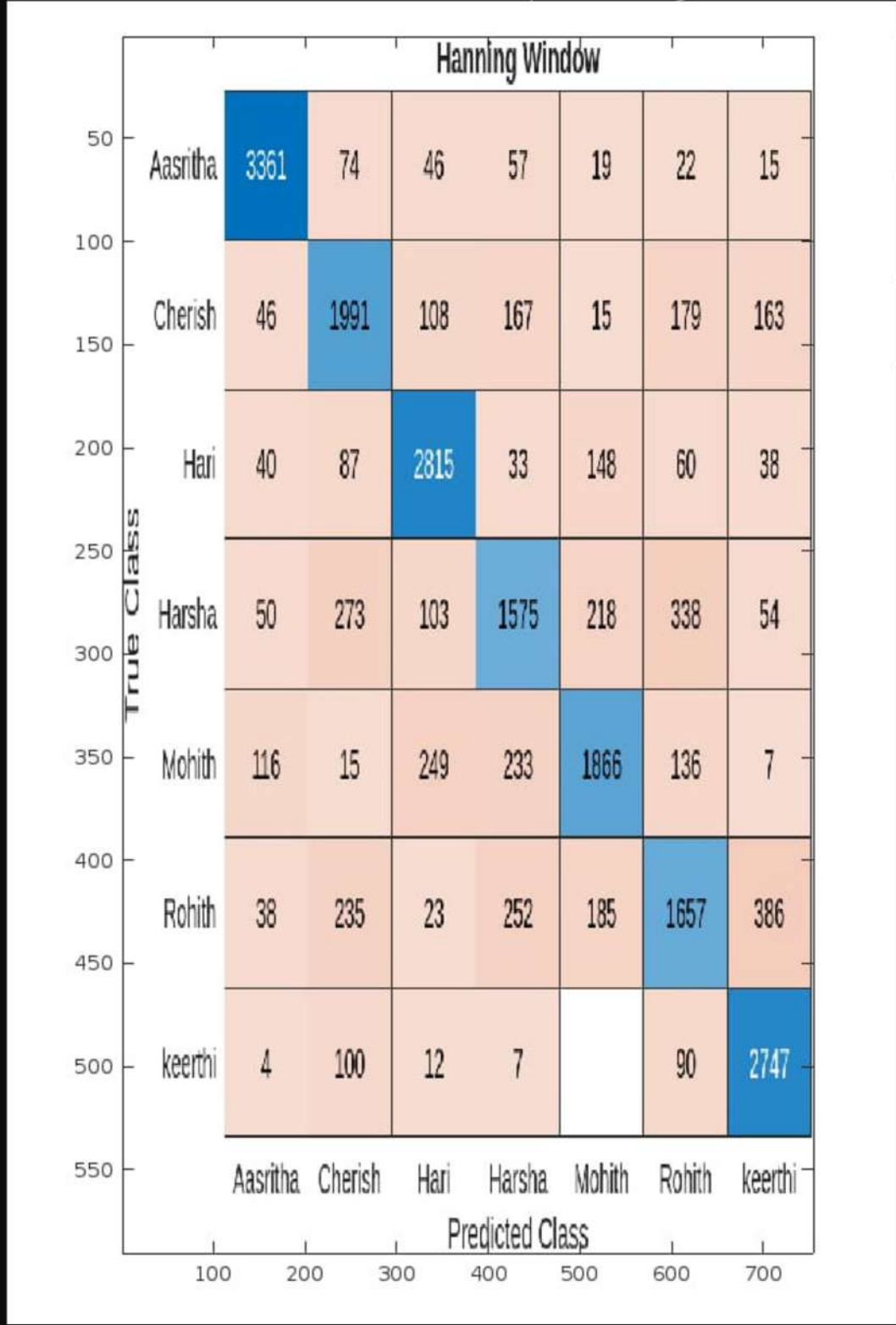
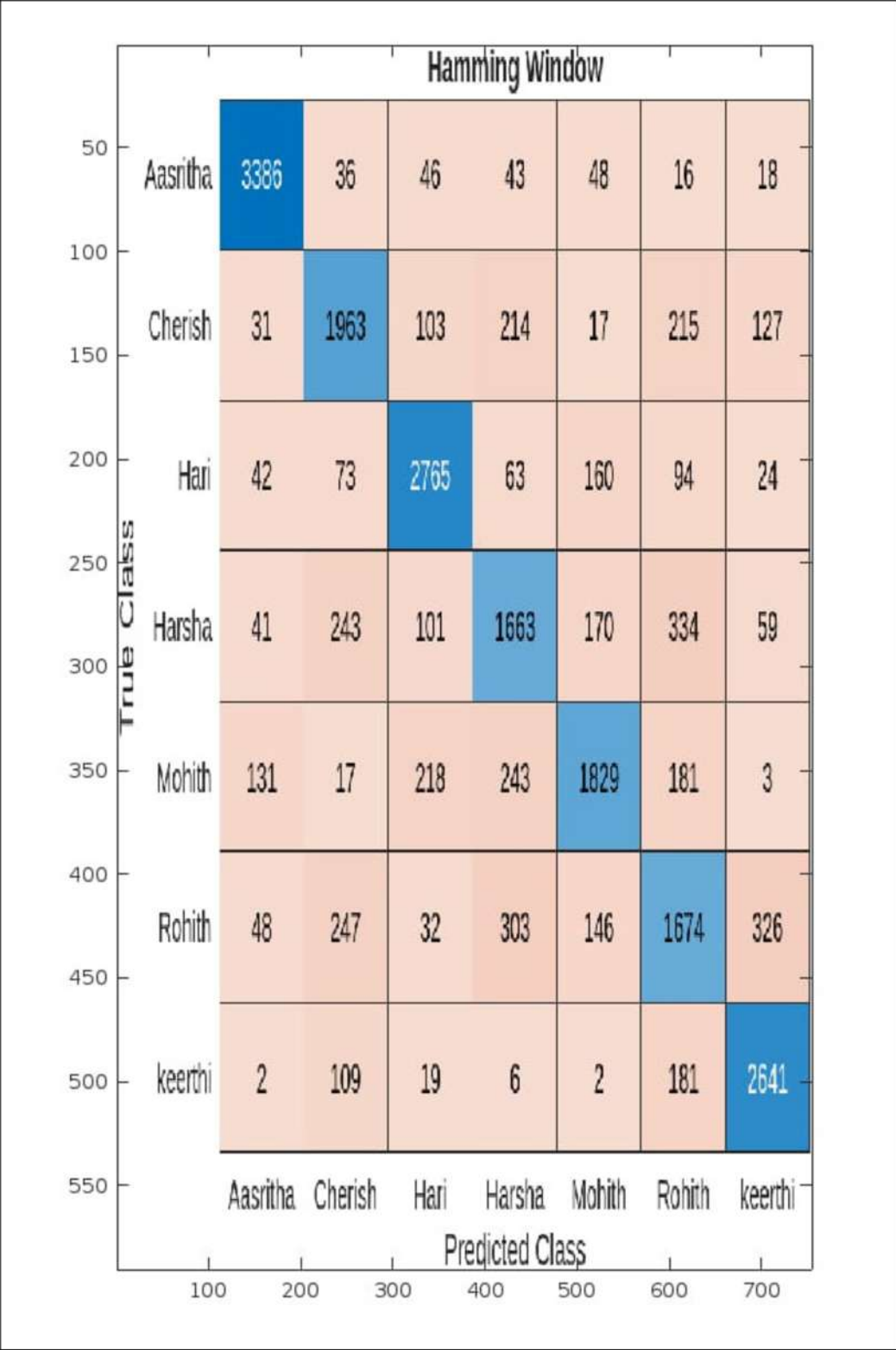
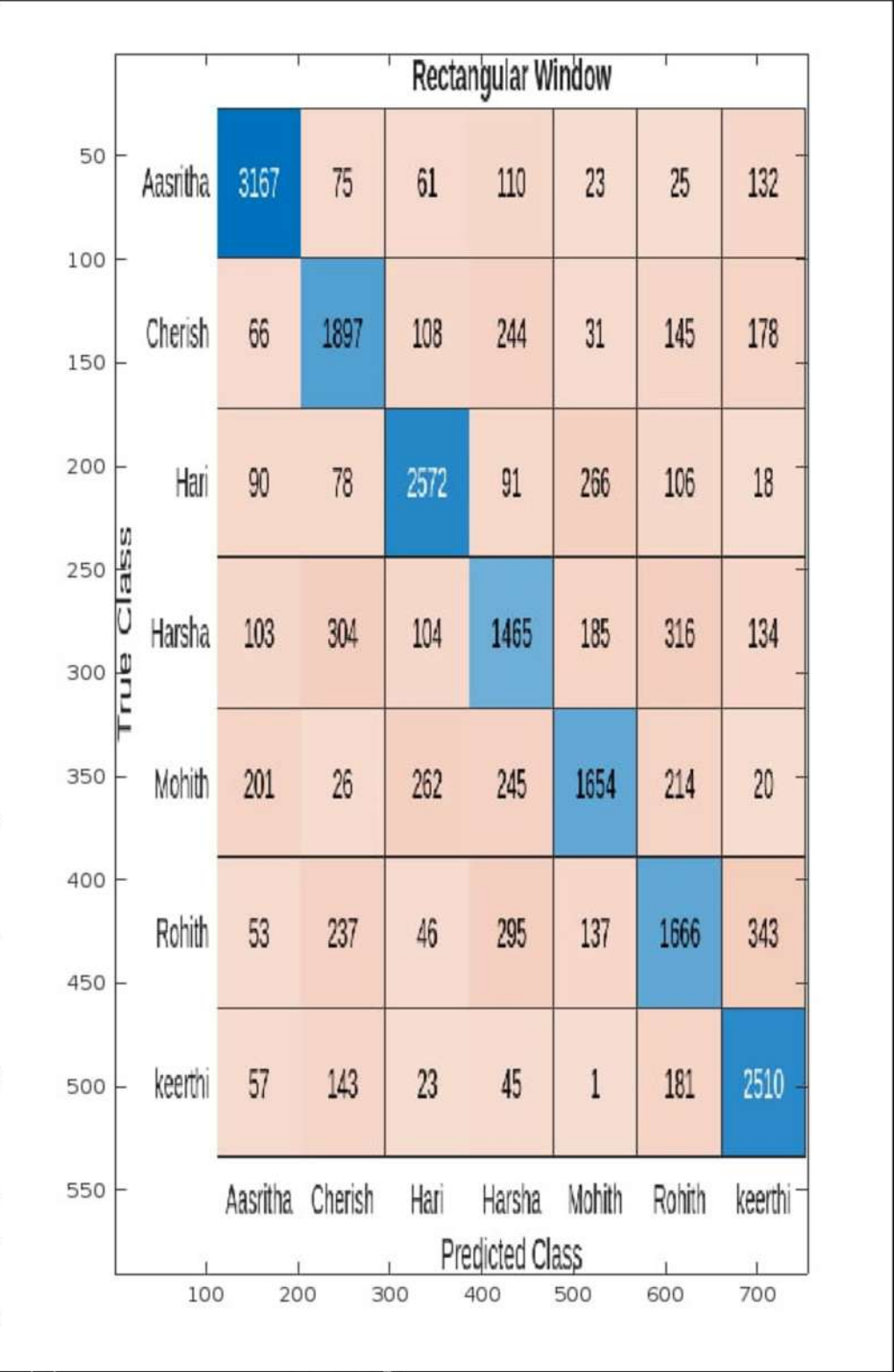


- Trained a Multi-class SVM classifier using `fitcecoc` in MATLAB.
- Used linear kernel for fast, effective classification.
- Dataset split: 80% training, 20% testing.
- Normalized feature vectors ensure balanced scaling across speakers.
- SVM distinguishes speakers by finding optimal decision boundaries.

Recognition Accuracy & Results

- Tested system with 3 different windowing functions.
- Accuracy Results:
 1. Rectangular: ~ 73.26%
 2. Hamming: ~ 78.36%
 3. Hanning: ~ 78.19%
- Hanning window gave highest accuracy
- Evaluated with confusion matrices showing correct and misclassified predictions.

CONFUSION TABLES





Limitations & Scope for Improvement

Current Limitations:

- Accuracy drops in noisy environments.
- GUI-Based Prediction is Inconsistent.
- Linear SVM may miss nonlinear patterns.
- No post-processing like speaker smoothing or Voice Activity detection (VAD).

Future Directions

- Integrate Deep Learning (CNNs, RNNs) for complex patterns.
- Add noise filtering and voice activity detection (VAD).
- Expand to multilingual or text-independent recognition.
- Port system to mobile platforms for real-world deployment.

The image features a dark background with white line art in the corners. The art consists of interconnected points and lines forming various geometric shapes, including triangles and polygons, creating a network-like structure. The central text is a large, bold, white sans-serif font.

THANK YOU