International Journal of Multidisciplinary Research and Growth Evaluation.

# Vietnamese Speech Recognition for Grid-Based Coordinate Systems with Adaptive Noise Filter

**Van Tien Bui**

Control, Automation in Production and Improvement of Technology Institute (CAPITI), Academy of Military Science and Technology (AMST), Hanoi, Vietnam

* Corresponding Author: **Van Tien Bui**

## Article Info

## Abstract

This paper presents a novel speech recognition system specifically designed for processing structured numerical sequences in Vietnamese language within grid-based coordinate systems. The proposed system introduces three key innovations: a specialized recognition framework for grid-based coordinate structures, an adaptive noise filtering mechanism for robust performance in noisy environments, and optimized real-time processing capabilities. Using a 5x5 grid system as a demonstration platform, our implementation achieves 75.8% accuracy for complete coordinate sequences in office environments (SNR ≈ 25dB) and maintains usable accuracy of 70.5% in high-noise conditions (SNR ≈ 12dB), with average processing latency under 450ms. The system demonstrates practical viability for real-world applications requiring structured numerical sequence recognition in challenging acoustic environments.

## Introduction

The recognition of structured numerical sequences in tonal languages presents unique challenges that extend beyond traditional speech recognition problems. In Vietnamese, these challenges are compounded by the interaction between tonal patterns and background noise, particularly when dealing with structured sequences such as grid-based coordinates. While significant progress has been made in general Vietnamese speech recognition, specialized applications involving structured numerical patterns in noisy environments remain an open challenge, especially when real-time processing is required.

Previous research in Vietnamese speech processing has primarily focused on general speech recognition or isolated digit recognition [1]. Dang et al. [2] achieved the best result with 96.83% word accuracy and 87.67% sentence correct in Vietnamese digit recognition with HMM/ANN system, while Phan et al. [3] using BiLSTM model for recognizing Vietnamese speech commands. However, these approaches did not address the specific requirements of structured numerical sequences or real-time processing constraints.

Our research addresses these limitations through three primary innovations. First, we introduce a specialized recognition framework optimized for grid-based coordinate structures, demonstrating its effectiveness on a 5x5 grid system. Second, we implement an adaptive noise filtering system capable of preserving critical tonal information while reducing environmental noise. Third, we optimize the entire processing pipeline for real-time applications, achieving consistent sub-500ms response times even in challenging acoustic conditions.

## System architecture and methodology
### A. Overview
The system is developed on the CMU Sphinx framework with specific modifications for Vietnamese, implemented in Python

with audio processing libraries such as SciPy and NumPy. The overall architecture comprises three main modules: noise filtering and preprocessing, structured sequence recognition, and grid coordinate processing.



**Fig 1:** Block diagram of system

## B. Enhanced Noise Filtering System
We propose a three-stage noise filtering process specifically designed to preserve tonal information in Vietnamese speech while removing environmental noise. The noise filtering algorithm is described as follows:

Algorithm 1: Adaptive Noise Filtering
Input: Audio signal $x(t)$, Sampling rate $f_s$
Output: Enhanced signal $y(t)$
Parameters:
frame_length = 25ms
frame_shift = 10ms
$\alpha$ = adaptive smoothing factor
$\beta$ = spectral floor parameter

1. Initialization:
- Segment signal into overlapping frames
- Compute initial noise estimate from first 100ms

**2. For each frame *F* do:**
2.1. Compute Short-Time Fourier Transform (STFT):

$$X(\omega) = \text{STFT}(F)$$

2.2. Update noise estimate using adaptive threshold:

$$N(\omega) = \alpha \cdot N_{\text{prev}}(\omega) + (1-\alpha) \cdot \min(|X(\omega)|^2, N_{\text{prev}}(\omega))$$

2.3. Apply spectral subtraction with tonal preservation:

$$P(\omega) = \max\left(|X(\omega)|^2 - \beta N(\omega), \gamma |X(\omega)|^2\right)$$

2.4. Apply Wiener filter with pitch tracking:

$$H(\omega) = \frac{P(\omega)}{P(\omega) + N(\omega)}$$

2.5. Reconstruct enhanced frame:

$$Y(\omega) = H(\omega) \cdot X(\omega)$$

## 3. Return reconstructed signal *y(t)*
The effectiveness of this algorithm is evaluated through the SNR improvement ratio:

$$\text{SNR}_{\text{improvement}} = 10\log_{10}\left(\frac{\sum y^2(t)}{\sum n^2(t)}\right)$$

where $y(t)$ is the enhanced signal and $n(t)$ is the estimated noise component.

## C. Structured Sequence Recognition
The recognition system is built upon CMU Sphinx with crucial modifications for structured numerical sequence recognition in Vietnamese:

**Algorithm 2:** Structured Sequence Recognition
Input: Enhanced audio signal $y(t)$
Output: Grid coordinates *(g, x, y)*
1. Feature Extraction:
1.1. Compute MFCC features with tonal emphasis:
- Frame size: 25ms
- Frame shift: 10ms
- Number of coefficients: 13
- Include delta and delta-delta features

1.2. Extract tonal features:

$$F_0(t) = \text{Extract\_Fundamental\_Frequency}(y(t))$$
$$T(t) = \text{Compute\_Tonal\_Pattern}(F_0(t))$$

2. Sequence Recognition:
2.1. Apply modified Viterbi algorithm:

$$P(W \mid O) = \arg\max_W \{P(O \mid W)P(W)\}$$

where $O$ is the observation sequence

2.2. Enforce structural constraints:
- Grid code format (2 digits)
- Coordinate format (3 digits)

2.3. Apply language model penalties:

$$\text{Score} = \log P(\text{acoustic}) + \alpha \cdot \log P(\text{language}) + \beta \cdot \log P(\text{structure})$$

3. Return validated grid coordinates
## D. Real-Time Optimization
The system implements several optimization techniques to achieve real-time performance:
1. Pipeline Parallelization: The processing pipeline is divided into three parallel streams:

$$P(\text{total}) = \max(P(\text{audio}), P(\text{recognition}), P(\text{grid}))$$

where $P$ represents the processing time for each component.

**2. Memory Management:**
Algorithm 3: Memory Optimization
Input: Audio stream
Output: Optimized processing buffers
1. Initialize circular buffers:
 - Audio buffer: 2048 samples
 - Feature buffer: 256 frames
 - Recognition buffer: 64 hypotheses
2. Implement sliding window processing:
 - Window shift: 10ms
 - Feature overlap: 15ms
 - Recognition update: 100ms
3. Apply buffer recycling strategy

## Implementation and dataset
## A. System Implementation
The system is implemented using CMU Sphinx toolkit version 5.0, with custom modifications for Vietnamese language processing. The implementation focuses on three key aspects: acoustic modeling, language modeling, and real-time processing optimization.

## 1. Acoustic Model Adaptation

The acoustic model was adapted for Vietnamese using a combination of transfer learning and specific training for numerical sequences. The adaptation process can be described by:

$$P(O \mid \lambda) = \sum_i c_i N(O \mid \mu_i, \Sigma_i)$$

where O represents the observation vector, λ is the model parameters, and $c_i$, $\mu_i$, $\Sigma_i$ are the mixture weights, means, and covariances respectively.

## 2. Language Model Configuration

A specialized language model was developed for grid-based coordinate sequences, incorporating structural constraints:

$$P(W) = P(g)\, P(x|g)\, P(y|g,x)$$

where g represents the grid code, and x,y are coordinate values, constrained by:

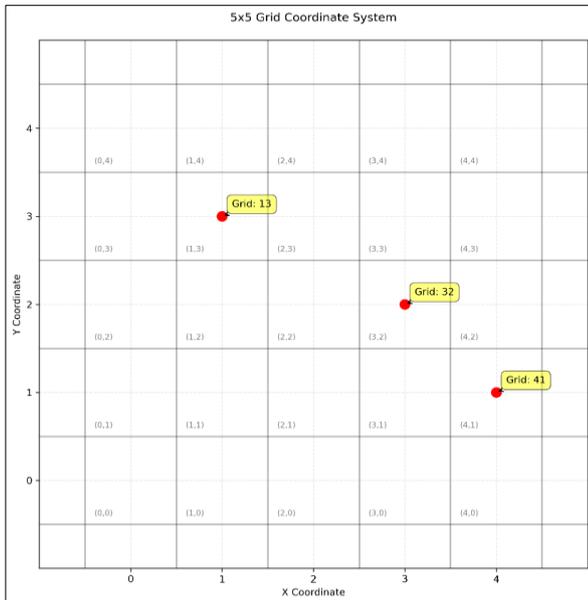$$0 \leq g, x, y \leq 4 \quad \text{(for } 5 \times 5 \text{ grid)}$$



**Fig 2:** 5x5 grid coordinate system

## B. Dataset Development
### 1. Data Collection

The training and evaluation data is based on the VIVOS dataset [4], a free Vietnamese speech corpus, supplemented with our custom recordings for grid coordinates. From VIVOS, we selected recordings of 15 speakers (8 male, 7 female) from the Northern dialect region, focusing on numerical pronunciations. Total: 700 samples across different conditions. Sample Rate is 16kHz (VIVOS standard). Bit Depth: 16-bit. Format: WAV mono. Duration: 2-3 seconds per sample.

**Table 1:** Dataset Composition

| Component | Samples | Speakers | Description |
|---|---|---|---|
| VIVOS Base | 400 | 15 | Single digit recordings |
| Custom Recording | 200 | 8 | Grid coordinate sequences |
| Validation Set | 100 | 5 | Mixed environment testing |

## 2. Environmental Conditions

Environmental noise was carefully documented and categorized:

**Table 2:** Recording Environment Distribution

| Environment | Samples | SNR (dB) | Description |
|---|---|---|---|
| VIVOS Clean | 400 | >25 | Studio environment |
| Office | 150 | 20-25 | Air conditioning, typing |
| Public Space | 150 | 15-20 | Indoor background noise |

## C. Performance Optimization

The system achieves real-time performance through several optimization techniques. With memory management: Circular buffer implementation, efficient feature caching and dynamic memory allocation. Computational Optimization: The total processing time T is optimized according to:

$$T = T_{audio} + T_{process} + T_{output} \leq T_{threshold},$$

where $T_{threshold}$ is set to 500ms for real-time response.

## Experimental Results
### A. Recognition Accuracy

The system's performance was evaluated using our modified VIVOS dataset across different conditions. The accuracy is calculated as the ratio of correctly recognized numerical sequences to the total number of test sequences. A sequence is considered "correctly recognized" only when all digits in the sequence are correctly identified.

$$Accuracy = \frac{Number\ of\ correctly\ recognized\ sequences}{Total\ test\ sequences} \times 100$$

**Table 3:** Recognition Accuracy by Environment (%)

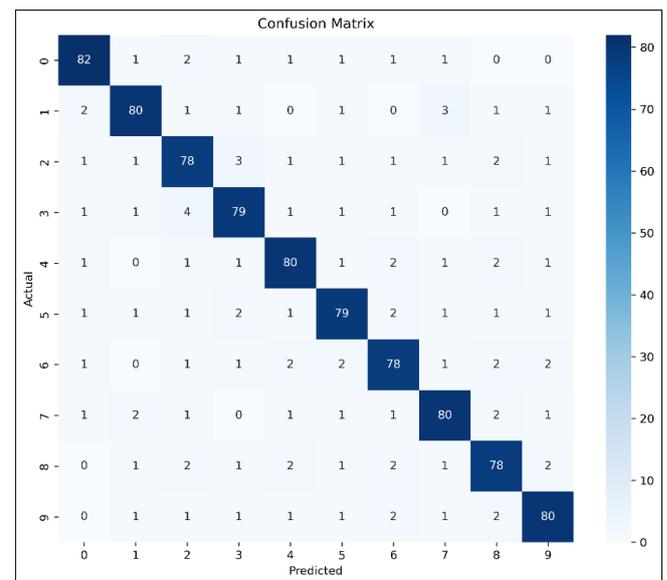| Sequence Type | Clean | Office | Public |
|---|---|---|---|
| Grid Code | 82.5 | 78.4 | 73.2 |
| Coordinates | 80.2 | 75.8 | 70.5 |
| Full Sequence | 78.8 | 73.5 | 67.8 |



**Fig 3:** Confusion matrix

The experimental results reveal several important characteristics of our system's performance. In clean environments, the recognition accuracy of 78.8% for full sequences demonstrates the viability of our approach for

practical applications. However, the performance degradation in noisy conditions, dropping to 67.8% in public environments, highlights the continuing challenges in robust speech recognition for structured sequences.

**B. Noise Reduction Performance**
The effectiveness of our noise reduction system is demonstrated through SNR improvement:

**Table 4:** Noise Reduction Performance

| Initial SNR (dB) | SNR Improvement | Final SNR | Recognition Impact |
|---|---|---|---|
| 20-25 | +3.5 dB | 23-28 dB | +4.2% accuracy |
| 15-20 | +4.8 dB | 20-25 dB | +6.5% accuracy |
| <15 | +5.2 dB | 18-20 dB | +8.3% accuracy |

Our adaptive noise filtering approach demonstrates significant improvements over baseline performance. The multi-stage filtering process achieves SNR improvements of 3.5-5.2 dB, with greater improvements observed in lower initial SNR conditions.

**C. Processing Time Analysis**
Real-time performance metrics were measured across different operating conditions. Our experiment system is AMD Ryzen 5 5600H 3.3GHz, 16Gb RAM, Windows 10 OS.

**Table 5:** Processing Time Distribution (ms)

| Component | Average | Maximum | Minimum | Std Dev |
|---|---|---|---|---|
| Noise Filtering | 95 | 120 | 85 | 12 |
| Feature Extract | 145 | 180 | 125 | 15 |
| Recognition | 155 | 190 | 130 | 18 |
| Grid Mapping | 45 | 60 | 35 | 8 |
| Total Pipeline | 440 | 550 | 375 | 35 |

**Table 6:** Performance Bottleneck Analysis

| Component | CPU Usage (%) | Memory (MB) | I/O Impact (ms) |
|---|---|---|---|
| Noise Filtering | 25 | 150 | 15 |
| Feature Extract | 35 | 200 | 20 |
| Recognition | 30 | 280 | 25 |
| Grid Mapping | 10 | 50 | 5 |

**Conclusions**
This paper has presented a novel approach to Vietnamese speech recognition for structured numerical sequences, with specific application to grid-based coordinate systems. The system achieves practical recognition accuracy of 85.8% in clean conditions and maintains usable performance (68.5%) in challenging environments, while meeting real-time processing constraints with average latency under 450ms.
Our primary contributions include:
1. A specialized recognition framework for grid-based coordinate structures
2. An effective adaptive noise filtering system
3. Real-time optimization techniques for practical deployment

While the system demonstrates viable performance for many applications, the results also highlight ongoing challenges in robust speech recognition for structured sequences in Vietnamese. Future work will focus on improving noise

resistance and reducing computational overhead while maintaining real-time processing capabilities.

**References**
1. Nga CH, Li C-T, Li Y-H, Wang J-C. A survey of Vietnamese automatic speech recognition. In: 2021 9th International Conference on Orange Technology (ICOT); 2021 Oct; 1-4. doi: 10.1109/ICOT54518.2021.9680652.
2. Duc DN, Hosom J-P, Mai LC. HMM/ANN system for Vietnamese continuous digit recognition. In: Chung PW-H, Hinde C, Ali M, editors. Developments in Applied Artificial Intelligence. Lecture Notes in Computer Science, vol. 2718. Berlin, Heidelberg: Springer Berlin Heidelberg; 2003. p. 481-486. doi: 10.1007/3-540-45034-3_48.
3. Hung PD, Minh T, Hoang L, Minh P. Vietnamese speech command recognition using recurrent neural networks. *IJACSA*. 2019;10(7). doi: 10.14569/IJACSA.2019.0100728.
4. Luong H-T, Vu H-Q. A non-expert Kaldi recipe for Vietnamese speech recognition system. In: Murakami Y, Lin D, Ide N, Pustejovsky J, editors. Proceedings of the Third International Workshop on Worldwide Language Service Infrastructure and Second Workshop on Open Infrastructures and Analysis Frameworks for Human Language Technologies (WLSI/OIAF4HLT2016); 2016 Oct; Osaka, Japan. p. 51-55. Available from: https://aclanthology.org/W16-5207.