

VOICE TO TEXT SYSTEM FOR DISABLED (V2TS)

SIVATHIBAN A/L KRISHNAMURTTU

**A thesis submitted in fulfillment of the
requirement for the award of the degree of
Bachelor of Computer Science**

**Faculty of Computer Systems & Software Engineering
Universiti Malaysia Pahang**

MAY 2010

ABSTRACT

Voice To Text System for Disabled (V2TS) is a speech recognition system that is aimed for Disabled people that can't type but have the speaking ability. V2TS deals with a number of useful functions where the disabled user can give command to the application for general operations such as reading text and closing the application itself. It has the Text To Speech (TTS) functionality where the text that has been input in the text area will be synthesized and read out by the computer. This report will discuss on the preparation, analysis, development and result that have been collected throughout the development cycle of this system. A complete reference and research details have been inserted in this document. This document will be a quick reference to refer on the specification and requirements of the system.

ABSTRAK

Sistem Suara ke Teks untuk Orang Kurang Upaya (V2TS) merupakan sebuah sistem yang disertakan dengan teknologi pengesanan pertuturan. Sistem ini adalah khas untuk Orang Kurang Upaya (OKU) yang tidak berkebolehan untuk menaip tetapi boleh bertutur. V2TS menyediakan pelbagai kemudahan dimana OKU boleh memberi arahan bersuara kepada aplikasi untuk operasi umum seperti membaca teks dan menutup aplikasi. Sistem ini juga mempunyai kebolehan untuk sintesis teks yang dimasukkan oleh pengguna dan membacanya. Laporan ini akan membincangkan perseediaan, analisis, pembangunan, dan juga keputusan yang telah dikumpul selama proses pembangunan system ini. Kesemua rujukan dan butiran mengenai sumber telah disertakan di dalam dokumen ini. Dokumen ini boleh digunakan sebagai rujukan untuk spesifikasi dan keperluan system.

TABLE OF CONTENTS

| CHAPTER | TITLE | PAGE |
|----------|---|----------|
| | SUPERVISOR'S DECLARATION | ii |
| | TITLE PAGE | iii |
| | DECLARATION | iv |
| | DEDICATION | v |
| | ACKNOWLEDGEMENT | vi |
| | ABSTRACT | vii |
| | ABSTRAK | viii |
| | TABLE OF CONTENT | ix |
| | LIST OF TABLES | xiii |
| | LIST OF FIGURES | xiv |
| | LIST OF APPENDICES | xvi |
| 1 | INTRODUCTION | 1 |
| | 1.1 Introduction | 2 |
| | 1.1.1 Research on Speech | 3 |
| | 1.1.1.1 Speech research in Microsoft Corporation | 3 |
| | 1.1.1.2 Speech research in Carnegie Mellon University | 3 |
| | 1.2 Problem Statement | 4 |
| | 1.3 Objective | 5 |
| | 1.4 Scope | 5 |
| | 1.5 Thesis Organization | 6 |
| 2 | LITERATURE REVIEW | 8 |

| | | |
|----------|--|-----------|
| 2.1 | Introduction to Speech Recognition | 9 |
| 2.2 | Existing System | 11 |
| 2.2.1 | Thai Automatic Speech Recognition | 11 |
| 2.2.2 | HMM-Based Speech Synthesis Applied to English | 12 |
| 2.2.3 | E-speaking | 14 |
| 2.2.4 | Alphabet Generator for Kids Using Speech Recognition | 15 |
| 2.3 | Techniques | 18 |
| 2.3.1 | Hidden Markov Model | 18 |
| 2.3.2 | Neural Network | 19 |
| 2.3.3 | Dynamic Time Warping | 21 |
| 2.3.4 | Linear Predictive Coding | 22 |
| 2.4 | Speech API | 24 |
| 2.4.1 | Microsoft Speech API | 25 |
| 2.4.2 | Java Speech API | 25 |
| 2.5 | Comparison of reviewed System | 27 |
| 2.6 | Recognizer | 28 |
| 2.7 | Conclusion | 29 |
| 3 | METHODOLOGY | 30 |
| 3.1 | Introduction | 31 |
| 3.2 | Analysis | 32 |
| 3.2.1 | User Requirement | 33 |
| 3.2.2 | System Requirement | 34 |
| 3.2.2.1 | Hardware Requirement (Device) | 34 |
| 3.2.2.2 | Software Requirement | 36 |
| 3.3 | Design | 37 |
| 3.3.1 | System Flow | 38 |
| 3.3.2 | Data Library (Speech Engine) | 40 |

| | | |
|----------|---|-----------|
| | 3.3.2.1 Data Library Process | 40 |
| | 3.3.3 Prototype Graphical User Interface | 41 |
| | 3.4 Construction | 42 |
| | 3.5 Testing | 43 |
| | 3.5.1 Black Box Testing | 43 |
| | 3.6 Deployment | 44 |
| | 3.7 General Requirement | 44 |
| | 3.7.1 Hardware Requirement | 44 |
| | 3.7.2 Software Requirement | 45 |
| 4 | IMPLEMENTATION | 47 |
| | 4.1 Introduction | 48 |
| | 4.2 V2TS Architecture | 48 |
| | 4.3 V2TS Interface Sketch (Adobe Photoshop) | 49 |
| | 4.3.1 Document Library of Sketch | 50 |
| | 4.3.2 V2TS Interface Design (Microsoft Expression Blend 3) | 51 |
| | 4.4 V2TS Development Environment | 60 |
| | 4.4.1 Initialization of Development | 60 |
| | 4.4.2 GUI Components | 63 |
| | 4.4.3 Speech Commands | 64 |
| | 4.5 References used in Development | 65 |
| | 4.6 Import Statements | 66 |
| | 4.6.1 Speech Library | 66 |
| | 4.6.2 Multipoint® Library | 70 |
| | 4.7 Source Codes for Controls | 72 |
| | 4.8 Source codes for functions and operations | 77 |
| 5 | RESULTS AND DISCUSSION | 84 |
| | 5.1 Introduction | 85 |

| | | |
|----------|--|-----------|
| | | 85 |
| | 5.2.1 Develop a prototype of Speech Recognition application that reacts to speech command and output words. | 85 |
| | 5.2.1 Design a simple and easy Speech Recognition application to enable Disabled people to get accessed to computer. | 86 |
| | 5.3 Output and Result | 87 |
| | 5.4 System Constraints | 88 |
| | 5.5 Suggestion | 89 |
| | 5.6 Future System | 90 |
| 6 | CONCLUSION | 91 |
| | 6.1 Conclusion | 92 |
| | REFERENCES | 94 |
| | APPENDICES | 97 |
| | Appendix A (Gantt Chart) | 98 |
| | Appendix B (User Manual) | 100 |
| | Appendix C (Software Test Report Document) | 106 |

LIST OF TABLES

| TABLE NO. | TITLE | PAGE |
|------------------|--|-------------|
| 2.1 | Parameters of speech recognition systems | 10 |
| 2.2 | Comparison of reviewed system | 27 |
| 2.3 | 6 common steps in recognizers | 28 |
| 3.1 | Hardware Requirements | 45 |
| 3.2 | Software Requirements | 46 |
| 4.1 | GUI Components | 63 |
| 4.2 | List of Speech Commands and functions | 64 |
| 5.1 | Output and Result | 87 |

LIST OF FIGURES

| FIGURE NO. | TITLE | PAGE |
|------------|--|------|
| 2.1 | Block Diagram of Speech Recognition | 10 |
| 2.2 | HMM-Based Speech Synthesis System | 13 |
| 2.3 | Main Menu Interface of E-speaking system | 14 |
| 2.4 | Main Interface of Alphabet Generator for Kids Using Speech Recognition system | 16 |
| 2.5 | Result generated for 'A' utterance | 16 |
| 2.6 | Result generated for 'B' utterance | 17 |
| 2.7 | Three state Markov model with transition probabilities | 19 |
| 2.8 | Neural Networks Model | 20 |
| 2.9 | Example of DTW implementation | 22 |
| 2.10 | Depiction of how LPC functions | 23 |
| 2.11 | Simplified model of speech production | 24 |
| 3.1 | Rapid Application Development | 32 |
| 3.2 | Flow Chart of the system | 38 |
| 3.3 | Flow of Data Library | 40 |
| 3.4 | Prototype Graphical User Interface (Suggestion only) | 41 |
| 4.1 | Sketching Graphical User Interface in Adobe Photoshop CS3 | 49 |
| 4.2 | Document Library of files used in GUI Sketch | 50 |
| 4.3 | Designing Interface in Microsoft Expression Blend 3 | 51 |
| 4.4 | Solution Explorer of V2TS in Visual Studio 2008 | 60 |
| 4.5 | MainWindow.xaml | 61 |
| 4.6 | About.xaml | 62 |
| 4.7 | Help.xaml | 62 |
| 4.8 | References that were used in V2TS development | 65 |
| 4.9 | Imported Libraries | 66 |
| 4.10 | System.Speech Library | 66 |

| | | |
|------|--|----|
| 4.11 | Classes in System.Speech.AudioFormat | 67 |
| 4.12 | Classes in System.Speech.Recognition | 67 |
| 4.13 | Classes in System.Speech.Recognition.SrgsGrammer | 68 |
| 4.14 | Classes in System.Speech.Synthesis | 68 |
| 4.15 | Classes in System.Speech.Synthesis.TtsEngine | 69 |
| 4.16 | Microsoft.Multipoint.SDK Library | 70 |
| 4.17 | Classes in Microsoft.Multipoint.SDK | 70 |
| 4.18 | Classes in Microsoft.Multipoint.SDK.Interop | 71 |
| 4.19 | Classes in Microsoft.Multipoint.SDK.Controls | 71 |
| 4.20 | Source code of btnHome | 72 |
| 4.21 | Source code of btnHelp | 72 |
| 4.22 | Source code of btnAbout | 73 |
| 4.23 | Source code of btnExit2 | 73 |
| 4.24 | Source code of btnOpen | 73 |
| 4.25 | Source code of btnSave | 74 |
| 4.26 | Source code of btnWriteStart | 74 |
| 4.27 | Source code of btnExit | 75 |
| 4.28 | Source code of btnWrite | 75 |
| 4.29 | Source code of btnRead | 75 |
| 4.30 | Source code of btnProcess | 76 |
| 4.31 | Source code of btnStop | 76 |
| 4.32 | Source code of btnClear | 76 |
| 4.33 | Source code of MainWindow | 77 |
| 4.34 | Source code of MouseConnected event | 78 |
| 4.35 | Source code of SpeechDetected event | 78 |
| 4.36 | Source code of CreateBitmapImage function | 78 |
| 4.37 | Source code of SpeechRejected event | 79 |
| 4.38 | Source code of Speech Recognition Initialization | 79 |
| 4.39 | Source code of MainWindow_KeyDown event | 79 |

LIST OF APPENDICES

| APPENDIX | TITLE | PAGE |
|-----------------|-------------------------------------|-------------|
| A | Gantt Chart | 98 |
| B | User Manual | 100 |
| C | Software Test Report Document (STR) | 106 |

CHAPTER 1

INTRODUCTION

This chapter presents an outline of the entire project and the introduction into problem statements, objectives, scopes and thesis organization.

1.1 Introduction

Voice to Text System for disabled (V2TS) is a speech recognition embedded system. Speech recognition converts spoken words to machine-readable input. It simply means that this program enables the computer to generate correct output through our voice by utilizing some related algorithm. Speech recognition can be used for isolated words or continuous speech. Isolated word means that the system will only take one word at a time while continuous speech has continuous speech characteristic and the system need to recognize and convert the utterance at once. Speech Application Programming Interface (SAPI) will be a very important part of this system because it needs to recognize wide range of voice and tones. The output for any voice input depends on the matching between available Grammar and the utterance speed.

The term "voice recognition" is sometimes used to refer to speech recognition where the recognition system is trained to a particular speaker as is the case for most desktop recognition software, hence there is an aspect of speaker recognition, which attempts to identify the person speaking, to better recognize what is being said. Speech recognition is a broad term which means it can recognize almost anybody's speech. V2TS deals with a number of useful functions where the disabled user can give command to the application for general operations such as reading text and closing the application itself. It has the Text To Speech (TTS) functionality where the text that has been input in the text area will be synthesized and read out by the computer. Furthermore to facilitate multiuser as the Disabled might need help from secondary user, the Multipoint® Technology have been implemented where more than one mouse can be utilized in the system. For processing purpose, the text can be saved as text file or send to word processor for editing and other related works.

If a user has lost the use of his hands, or for visually impaired users when it is not possible or convenient to use a Braille keyboard, the systems allow personal expression through dictation as well as control of many computer tasks. Some programs save users'

speech data after every session, allowing people with progressive speech deterioration to continue to dictate to their computers. ^[1]

1.1.1 Research on Speech

1.1.1.1 Speech research in Microsoft Corporation

Microsoft Research has a group in Redmond and another in Beijing working together to improve spoken language technologies. Their main goal is to build applications that make computers available everywhere, and work with its Speech Products Group to make this vision a reality. The research are interested not only in creating state-of-the-art spoken language components, but also in how these disparate components can come together with other modes of human-computer interaction to form a unified, consistent computing environment. Microsoft is pursuing several projects to help reach its vision of a fully speech-enabled computer ^[2]

1.1.1.2 Speech research in Carnegie Mellon University (CMU)

The Sphinx Group at Carnegie Mellon University is committed to releasing the long-time, DARPA-funded Sphinx projects widely, in order to stimulate the creation of speech-using tools and applications, and to advance the state of the art both directly in speech recognition, as well as in related areas including dialog systems and speech synthesis.

The Sphinx Group has been supported for many years by funding from the Defense Advanced Research Projects Agency, and the recognition engines to be released are those that the group used for the various DARPA projects and their respective evaluations.

The packages that the CMU Sphinx Group is releasing are a set of reasonably mature, world-class speech components that provide a basic level of technology to anyone interested in creating speech-using applications without the once-prohibitive initial investment cost in research and development; the same components are open to peer review by all researchers in the field, and are used for linguistic research as well. ^[3]

1.2 Problem Statement

This system is aimed at Disabled people therefore the problems are Disabled people that can't write or type encounters difficulties in expressing their spoken words in text or documents form. Disabled people have difficulties to communicate well in this fast paced community where they are left behind in Information Technology. Most of them are not computer literate therefore they need an application that simplifies their task by just using speech command.

1.3 Objective

- i. Develop a prototype of Speech Recognition application that reacts to speech command and output words.
- ii. Design a simple and easy Speech Recognition application to enable Disabled people to get accessed to computer.

1.4 Scope

The scopes of the project are:

Project:

- i. Recognize spoken words in English either command or grammar.

User:

- i. Disabled people
- ii. Normal user (Tutor/Assistant)

1.5 Thesis Organization

Chapter 1: Introduction

The purpose of this chapter is to introduce to the readers about the project that will be developed later. This chapter contains introduction, problem statement, objective, and scope and thesis organization.

Chapter 2: Literature review

This chapter explains about the reviews for the chosen project. This chapter is divided into two sub reviews that require students to study to get complete information about the project.

Chapter 3: Methodology

This chapter discusses the approach and framework for the project. Method, technique or approach that will be and will be used while designing and implementing the project will be included in the content. Justification and of method on approach used and hardware and software necessary is stated here.

Chapter 4: Implementation

This chapter acts to document all processes that involve in the development of the project. Designed project development is explained here. The content of this project depends on the system. It contains information of database and tools used. Data in database is shown in this chapter.

Chapter 5: Results and Discussion

The purpose of this system is to explain about the results and data analysis that had been acquired. Result analysis, project limitation and suggestion and project enhancement are contents for the chapter.

Chapter 6: Conclusion

This chapter explains briefly and summarizes the developed project.

CHAPTER 2

LITERATURE REVIEW

This chapter explains about the reviews for the chosen project. This chapter is divided into two sub reviews that require students to study to get complete information about the project.

2.1 Introduction to Speech Recognition

Speech recognition is the process of converting an acoustic signal, captured by a microphone or a telephone, to a set of words. The recognized words can be the final results, as for applications such as commands & control, data entry, and document preparation. They can also serve as the input to further linguistic processing in order to achieve speech understanding. Speech recognition systems can be characterized by many parameters. An isolated-word speech recognition system requires that the speaker pause briefly between words, whereas a continuous speech recognition system does not. Spontaneous, or extemporaneously generated, speech contains disfluencies, and is much more difficult to recognize than speech read from script. Some systems require speaker enrollment where a user must provide samples of his or her speech before using them, whereas other systems are said to be speaker-independent, in that no enrollment is necessary. Some of the other parameters depend on the specific task. Recognition is generally more difficult when vocabularies are large or have many similar-sounding words. When speech is produced in a sequence of words, language models or artificial grammars are used to restrict the combination of words. The simplest language model can be specified as a finite-state network, where the permissible words following each word are given explicitly. One popular measure of the difficulty of the task, combining the vocabulary size and the language model, is *perplexity*, loosely defined as the geometric mean of the number of words that can follow a word after the language model has been applied. Finally, there are some external parameters that can affect speech recognition system performance, including the characteristics of the environmental noise and the type and the placement of the microphone. ^[4]

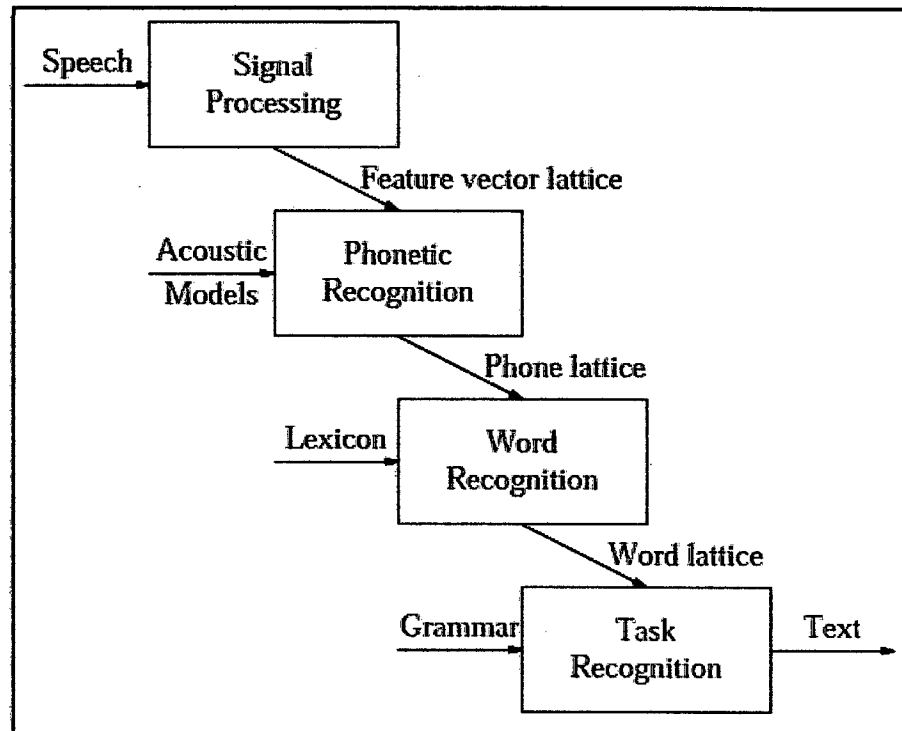


Figure 2.1: Block Diagram of Speech Recognition

Table 2.1: Parameters of speech recognition systems

| Parameters | Range |
|----------------|--|
| Speaking Mode | Isolated words to continuous speech |
| Speaking Style | Read speech to spontaneous speech |
| Enrollment | Speaker dependent to Speaker independent |
| Vocabulary | Small (less than 20 words) to large (more than 20,000 words) |
| Language Model | Finite state to context sensitive |
| Perplexity | Small (<10) to large (>100) |
| SNR | High (>30dB) to low (<10dB) |

2.2 Existing system

2.2.1 Thai Automatic Speech recognition

2.2.1.1 Introduction

This research was performed as part of the DARPA-Babylon program aimed at rapidly developing multilingual speech-to-speech translation capability in several languages. It is built on extensive background in ASR, language portability, and speech translation, the group has built Arabic-English and Thai-English Speech-to-Speech translation systems in less than 9 months per language. This system has been used in an external DARPA evaluation involving medical scenarios between an American Doctor and a naive monolingual Thai patient. ^[5]

2.2.1.2 Technique

Hidden Markov Model has been used in this system.

2.2.1.3 Objective

To develop a robust and flexible Thai Speech Recognizer that can be integrated to Thai-English speech translation.

2.2.1.4 Feature

- I. Automatic pronunciation generation.
- II. Rapid bootstrapping.
- III. Phone set and pronunciation variation.
- IV. Real Time recognizer for medical dialogs

2.2.2 HMM-Based Speech Synthesis System Applied to English

2.2.1.2 Introduction

Although many speech synthesis systems can synthesize high quality speech, they still cannot synthesize speech with various voice characteristics such as speaker individualities, speaking styles, emotions, etc. To obtain various voice characteristics in speech synthesis systems based on the selection and concatenation of acoustical units, a large amount of speech data is necessary. However, it is difficult to collect store such speech data. In order to construct speech synthesis systems which can generate various voice characteristics, the HMM-based speech synthesis system (HTS) was proposed. In the training part, spectrum and excitation parameters are extracted from speech database and modeled by context dependent HMMs. In the synthesis part, context dependent HMMs are concatenated according to the text to be synthesized. Then spectrum and excitation parameters are generated from the HMM by using a speech parameter generation algorithm. Finally, the excitation generation module and synthesis filter module synthesize speech waveform using the generated excitation and spectrum parameters. The attraction of this approach is in that voice characteristics of synthesized speech can easily be changed by transforming HMM parameters. In fact, it is shown that we can change voice characteristics of synthesized speech by applying a speaker adaptation technique, a speaker interpolation technique, or an Eigen voice technique. ^[6]

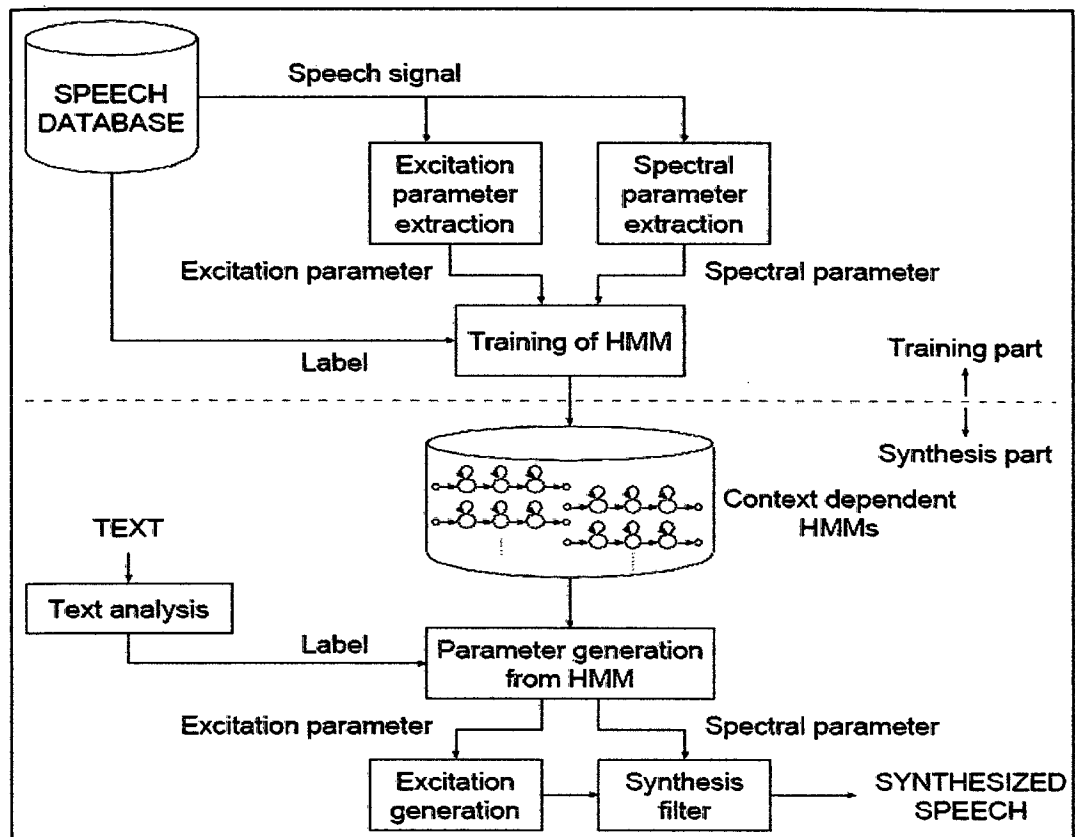


Figure 2.2: HMM-Based Speech Synthesis System

2.2.2.2 Technique

Hidden Markov Model has been used in this system.

2.2.2.3 Objective

To utilize HMM-based speech synthesis system (HTS) to English speech synthesis.

2.2.2.4 Feature

- I. Spectrum Modeling.