# Speech-Enabled Interactive Enquiry System in Tamil

**Principal Investigator**

- Dr. T. Nagarajan, Prof. & Head/IT

**Co-Investigators**

- Dr. P. Vijayalakshmi, Prof./ECE

- Dr. B. Bharathi, Asso. Prof./CSE

- Ms. S. Sasirekha, Asst. Prof./IT

SSN College of Engineering,
Kalavakkam

# Contents

# Introduction

A speech-enabled interactive enquiry system is one in which a user interacts with a system to obtain the required information. It consists primarily of a speech recognition system, a database, and a text-to-speech synthesis system, as shown in Fig. 1.1. In the current project, such a system is being developed to provide information regarding agriculture, specifically regarding paddy, sugarcane, and ragi.
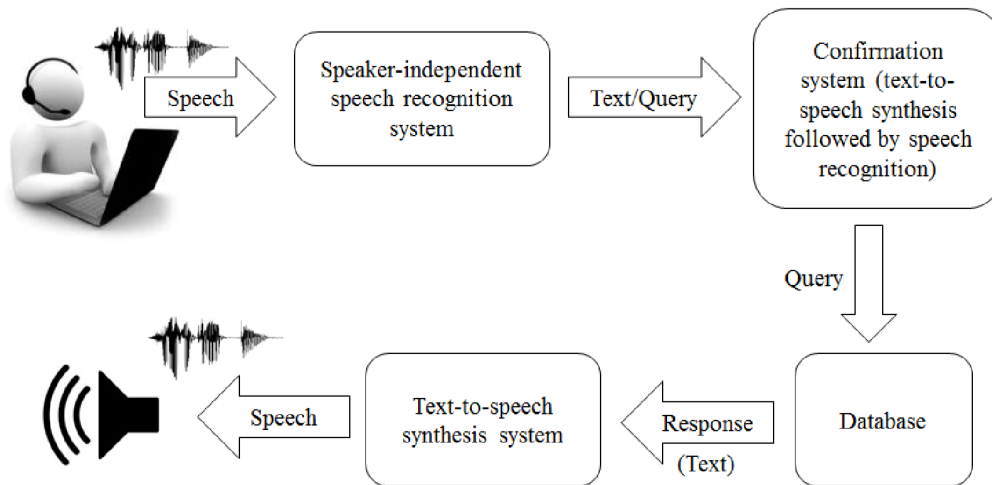


Figure 1.1: Speech-Enabled Interactive Enquiry System

The following steps have been taken in developing the enquiry system:

1. Analysis of paddy, sugarcane, and ragi related data

2. Collection of text and speech data

3. Development of speech recognition systems

4. Development of an application that combines the speech recognition systems, a data retrieval logic, and a speech synthesizer

The report elaborates these steps and is organized as follows: Chapter 2 describes the preliminary work that was carried out, specifically the development of the first prototype of the enquiry system and interactions with TNAU. Chapter 3 provides details on the text and speech data collected for use in the database and to train the recognition system. Chapter 4 describes the recognition systems that were developed and a likelihood-based approach to confirm the recognition result with the user in the event of a doubt. Chapter 5 describes the final enquiry system capable of providing information on paddy, sugarcane and ragi. Finally, Chapter 6 provides the financial details of the project and Chapter 7 summarizes the report.

# Preliminary Work

## 2.1 First Prototype

A simple prototype of the enquiry system was first created and presented to
the Tamil Virtual Academy (TVA), Chennai on 15/04/2016. This prototype
could mimic the target application but it worked only for a limited number
of use cases (ten questions related to agriculture). When the user asked the
system these ten questions the application could answer them with 100%
accuracy. The prototype also had a web user interface that displayed the
text of the question asked and the answer to the question, apart from the
speech output provided, for reference. Fig. 2.1 shows a screen-shot of this
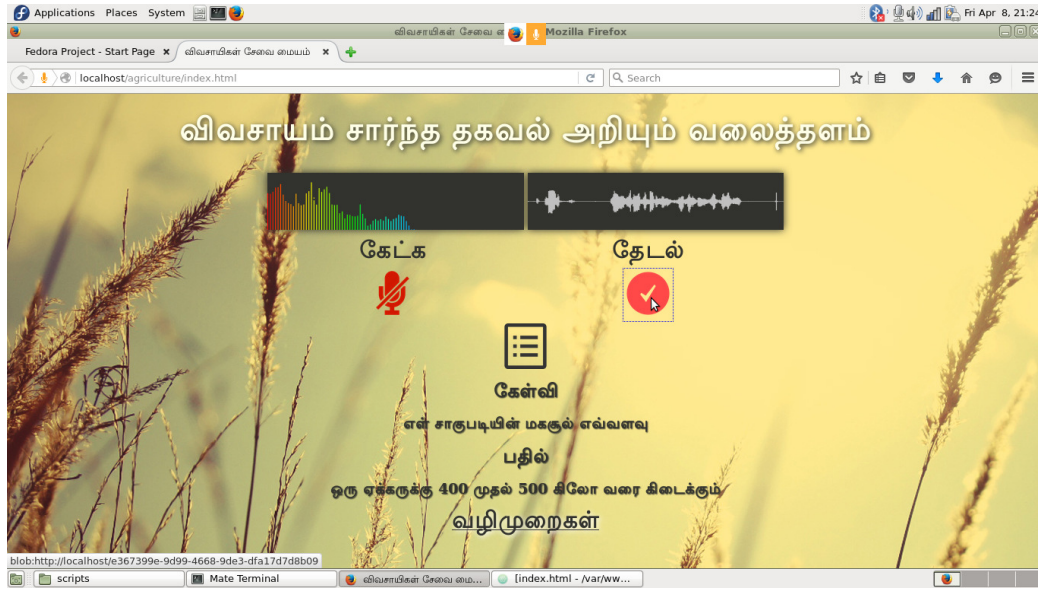prototype.



Figure 2.1: Prototype Application

## 2.2   Interactions with TNAU

Two project officers from SSN College of Engineering (SSNCE), visited Tamil Nadu Agricultural University (TNAU), Coimbatore, twice, for the purpose of knowledge transitions and data collection.

The first visit on 02/05/2016 served as an orientation to the agricultural domain. From this visit, the task at hand was understood from a whole new perspective. During this visit, the project staff met with the (i) Director, e-Extension Center, (ii) Professor and members of the e-Extension Center, (iii) Director, Center for Agricultural and Rural Development Studies (CARDS), and (iv) Technical Personal Officer to Vice Chancellor. The visit helped in understanding the breadth and depth of the agricultural domain and the approach that should be taken for building the enquiry system.

A lot of time was spent in the e-Extension center as this is the hub where all the data about agriculture are transformed into information for the common man and the farmer. During this time, the organization of the Agritech portal and the depth of information that could be retrieved from it were understood. This is currently serving as our primary reference for all agricultural information.

A visit to the Kisan Call Center provided more practical information about the kind of questions a farmer usually asks and the kind of replies they expect. This is crucial information for a speech enquiry system such as ours, since our application has to be fine tuned to the concerns of the farmer, as opposed to a book of reference.

Separate discussions were held with the directors mentioned above as well. During these discussions it was advised that it would be a daunting task to cover all aspects of agriculture in a short period of time, as it would amount to tedious research and data collection activities, and hence to concentrate on a specific aspect.

A second visit to TNAU was made on 06/06/2015 along with Dr. Dhanalakshmi, Assistant Director, TVA. During this trip, the project officers had the opportunity to meet the Vice Chancellor, TNAU and present a demo of the Prototype. They also had a discussion on the involvement of TNAU in providing intelligence concerning agriculture for building the application efficiently. The Vice Chancellor extended and assured his full support in providing any intelligence regarding the agriculture domain for building the application. He was also gracious enough to provide all the information used for developing the Expert System developed by TNAU.

A meeting with the Director and Head of the Department of the Paddy Research Center was also held and the project staff gathered information

from the Paddy Research Center concerning the current advancements in paddy and the kind of questions and clarifications that farmers usually ask them concerning paddy production and protection. Discussions were also held with an entomologist at the paddy research center to specially discuss about the crop protection aspects.

## 2.3  Analysis on Agritech Portal

A detailed analysis was performed on the information available in the Agritech portal and it was initially decided that the enquiry system would focus on the following aspects of agriculture:

1. Crop production - Deals with the production of crops, technologies associated with it, latest breed varieties, latest cultivation methods, etc.

2. Crop protection - Deals with insects, pests and diseases, their symptoms, identification, and management

The number of questions possible from a farmer, across all crops for crop production and crop protection, were then estimated from the Agritech portal. These are portrayed in Table 2.1.

| Category | Sub-Category | Number of Questions |
|---|---|---|
| Crop Production | | 5250 |
| Crop Protection | Agricultural Crop - Insects/Pests | 7170 |
| | Horticultural Crop - Insects/Pests | 9990 |
| | Grain Storage | 900 |
| | Agricultural Crop - Diseases | 1431 |
| | Horticultural Crop - Insects/Pests | 2268 |
| | Post Harvest  Diseases | 1116 |
| **Totals** | | **28125** |

Table 2.1: Number of possible questions across all crops for crop production and protection

As evident from the figures given in Table 2.1, the number of possible questions that a farmer might ask the system is very huge. Also the number of questions mentioned here is only limited by human imagination, since in the real world scenario we may encounter several other questions than what is accounted for here. Therefore, based on the enormity of the task at hand and the advice from TNAU experts, as mentioned previously, it was

decided that the enquiry system would initially focus on one crop. Further, since paddy is the most commonly cultivated crop and since information on the production and protection of paddy amounts to about 40% of the total information available in the Agritech portal, the first two phases of the project focused on developing an enquiry system that provides information on paddy. In the final phase of the project, two other commonly cultivated crops in Tamil Nadu, namely, sugarcane and ragi, were also added to the enquiry system.

# Data Collection

## 3.1 Analysis of Paddy, Sugarcane, and Ragi Related Data

Although restricting to three crops would limit the number of questions, each of these questions may be posed in several ways. As a result, it would be more suitable to guide the user towards the question at hand, through a series of interactive questions, thus reducing the complexity of the system by a huge margin. In this regard, the information available on paddy, sugarcane, and ragi in the agritech portal and the expert system acquired from TNAU were analyzed and questions to be posed by the enquiry system to the user and the anticipated word or phrase responses were formulated. The possible responses were then recorded from multiple speakers to use in the development of speech recognition systems. These are elaborated below.

### 3.1.1 Paddy

Information on paddy is broadly classified into two categories, namely, protection and production. The protection category provides information on measures to be taken in preventing or treating symptoms caused by pests and diseases. The user may specify the name of the disease or pest to retrieve relevant information or specify the symptoms observed to allow the system to identify the disease/pest and provide the appropriate information. A total of 8 questions, depicted in Fig. 3.1, that would elicit 1 to 3-word responses, have been finalized. All possible answers to these questions were then derived. These responses must also take into account possible variations of the same word. In this regard, 181 responses related to paddy protection have been derived.

Of the information on paddy production three categories, namely, seasons and varieties, nursery management, and paddy ecosystem were considered. An outline of the information, categorized in this form, is depicted in Fig. 3.2.
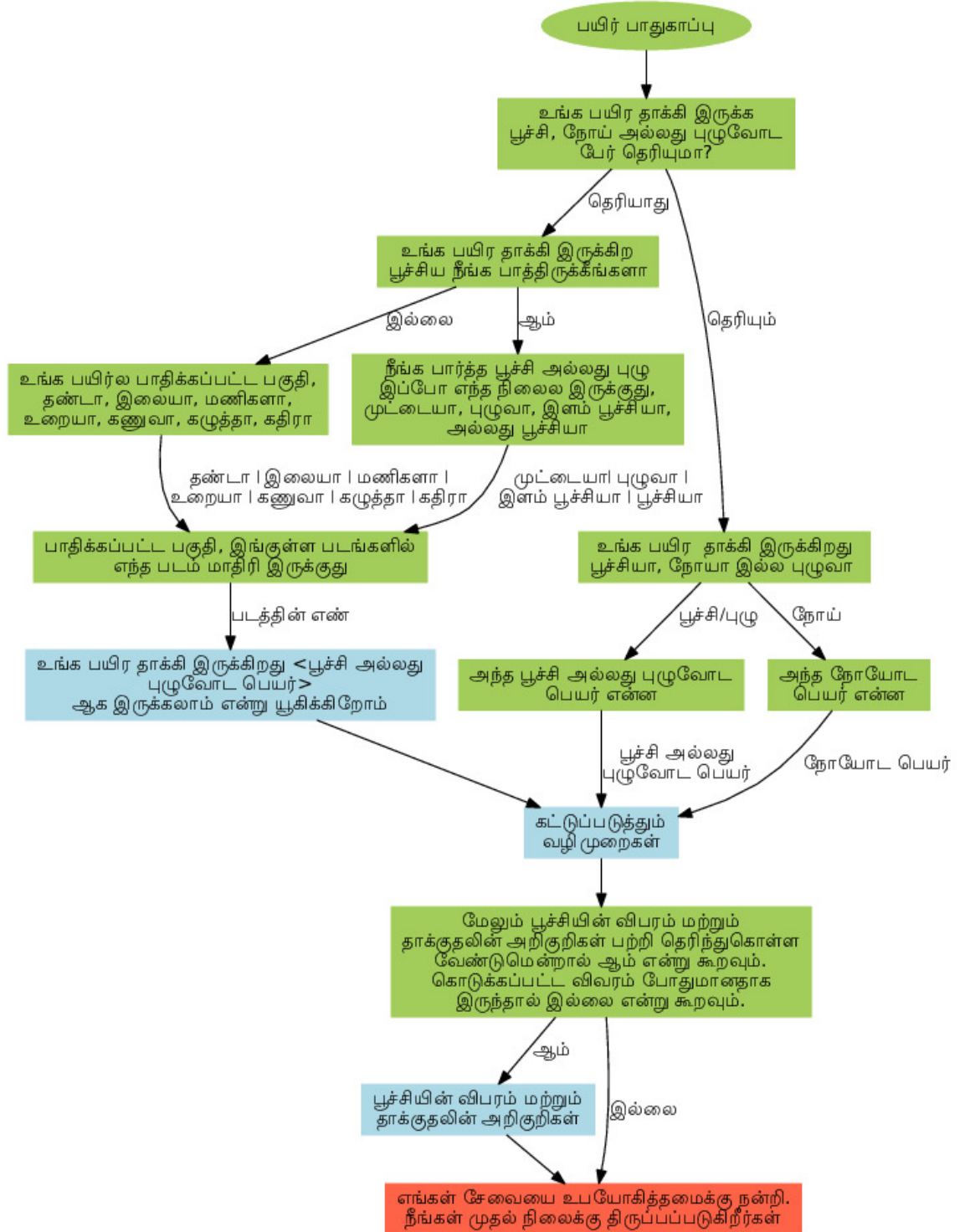
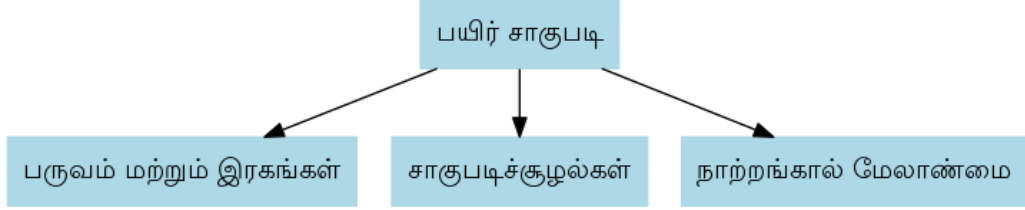Figure 3.1: Crop protection flow for paddy

Figure 3.2: Production overview

Discussions at TNAU revealed that farmers most frequently requested for information regarding the variety of paddy that provides the best yield in a particular season. In this regard, details on the varieties of paddy that may be planted in each season, in each district of Tamil Nadu, have been collected. The user can navigate to the seasons section and provide the name of the season and district and obtain the names of the varieties that can be planted for the specified season and district. The user can also specify directly the name of the variety to obtain information about it. This is shown in Fig. 3.3. The details collected include information on 28 districts, 14 seasons, and 45 varieties of paddy.
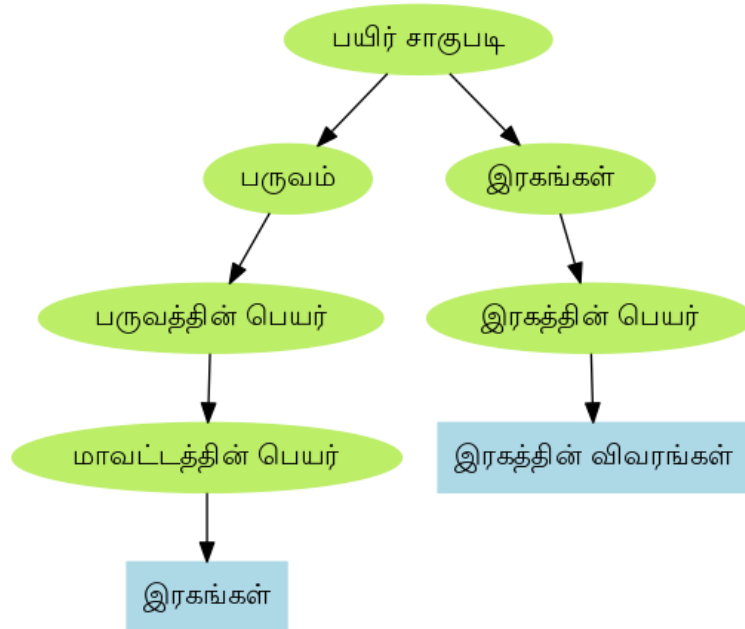
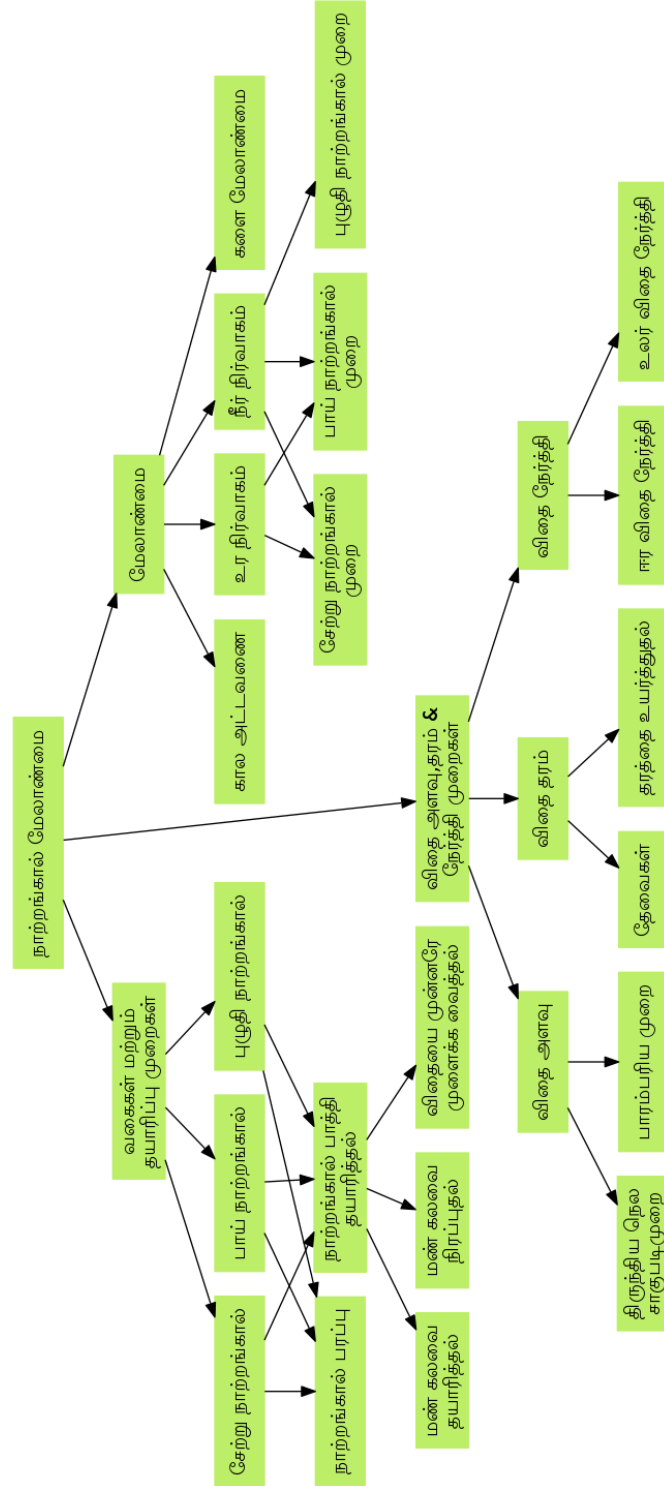

Figure 3.3: Seasons and Varieties
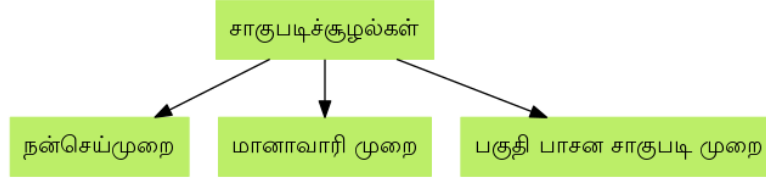
Figure 3.4: Nursery management
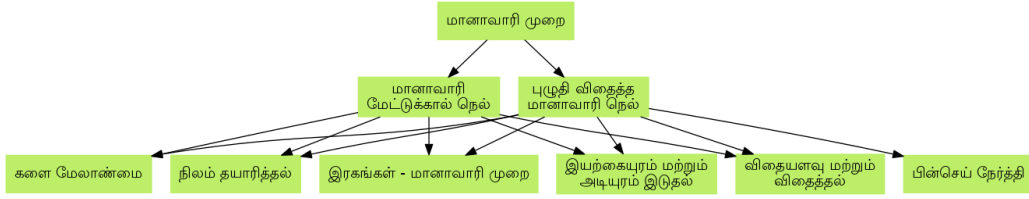
Figure 3.5: Paddy ecosystem - Overview



Figure 3.6: Paddy ecosystem - Dry system



Figure 3.7: Paddy ecosystem - Semidry system



Figure 3.8: Paddy ecosystem - Wet system - Overview

The section on nursery management contains information on the types of nurseries, seed management, seed treatment, and some management strategies for water, nutrients, and weeds, as depicted in Fig 3.4. The section on paddy ecosystem deals with the kind of ecosystem/environment/growing conditions that are required by the paddy plants to give a good yield. Fig. 3.5 depicts an overview of paddy ecosystem. Figures 3.6 to 3.10 expand the nodes shown in Fig. 3.5.

For paddy production, a total of 47 questions have been formulated and 174 responses have been derived.

Figure 3.9: Paddy Ecosystem - Wet System - Part 1

Figure 3.10: Paddy ecosystem - Wet System - Part 2

### 3.1.2   Sugarcane

Information on sugarcane is categorized into 12 primary sections, namely, main field preparation, fertilizer, management of main field operations, weed manage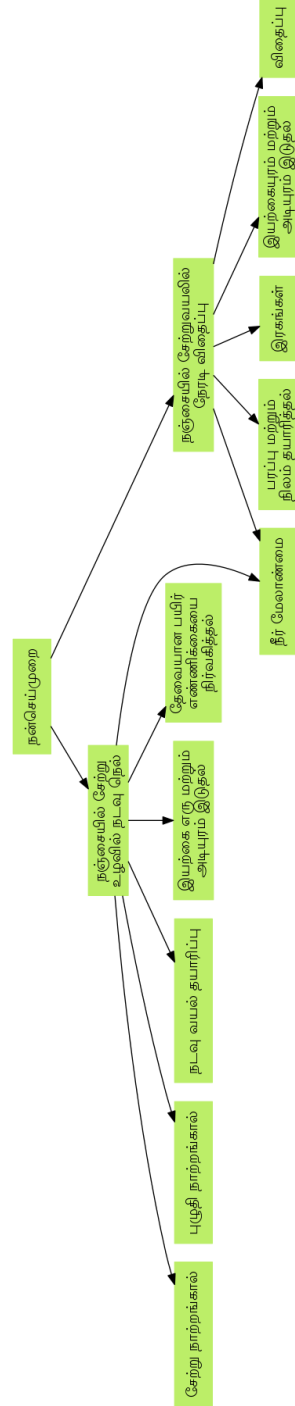ment, nitrogen saving, nutrition, water management, contingent plan, harvesting, raising inter crops, propping, and filling gaps. An outline of the information, categorized in this form, is depicted in Fig. 3.11. A total of 10 questions and 40 responses have been formulated.

### 3.1.3   Ragi

The information available on the production of ragi has been categorized into 12 sections as well, namely, land preparation, forming raised beds, crop management, water management, main field management, harvesting, manuring and fertilization, cropping system, and seeds. This is portrayed in Fig. 3.12. Based on this information, 8 questions and 40 responses have formulated on ragi.

## 3.2   Speech Data Collection

The responses derived for all three crops were recorded from multiple speakers over a head-mounted microphone, that is a part of the Philips SHG7980 USB Headset, in a laboratory environment. The data was recorded at a sampling rate of 16 kHz, over a single channel, with 16 bit precision.

The 181 responses derived for paddy protection were recorded from 20 speakers (8 male and 12 female). Four speakers recorded two versions of the responses, while the rest recorded each response once. The 169 responses derived for paddy production were recorded by 11 speakers, 5 of whom were male and 6 female. One male speaker recorded three variations of each response, three recorded two variations, and one speaker recorded each response once. Among the female speakers, three recorded each response once, while three recorded them twice. Therefore, there are now 10 sets of data collected from male speakers and nine from female speakers. For sugarcane and ragi, 77 responses were recorded from 12 speakers, 6 male and 6 female. One male speaker and all the female speakers recorded each response twice, while the others recorded them once.

Figure 3.11: Sugarcane production flow

Figure 3.12: Ragi production flow

## 3.3 Segmentation

First it was ensured that the amount of silence at the start and end of each response was between 70 and 120 ms. In order to derive time-aligned phonetic transcriptions, around 15 minutes of data were manually segmented. Monophone models with 5 states and 3 mixture components per state were then trained from these manually segmented files. The monophone models were then used to segment the entire database using the forced Viterbi algorithm iteratively.

# Speech Recognition System

HMM-based speech recognition systems were developed with the data described in the previous chapter. The details of the systems developed are as follows:

## 4.1 Analysis on Recognition Systems

In order to determine the configuration of the recognition system to be used, initially a number of recognition systems were trained on the paddy protection data described in Section 3.2. The system that yielded the best result was identified and this system was trained on the entire speech data collected and used in the enquiry system.

### 4.1.1 Isolated Word Recognizer

Initially, an isolated word recognizer was developed, consisting of 181 word models, with the number of states equal to thrice the number of phonemes in the word and 5 mixture components per state. The performance of this system however, was not satisfactory and so it was decided that triphone-based recognizers would be more suitable.

### 4.1.2 Triphone Recognizer

This system consists of 546 word-internal triphone models with three states. The number of mixture components per state has been varied from 1 to 16, in powers of 2. The triphone system when tested on isolated words, performed much better than the isolated word recognizer and using 8 mixture components per state yielded the best result.

### 4.1.3   Triphone Recognizer with Garbage Models

Although the questions in the enquiry system are framed such that they would elicit one of the responses derived in Chapter 3, it is possible that the user might occasionally answer in phrases or sentences consisting of out-of-vocabulary words. This would necessitate the identification of the keyword in the phrase/sentence. In this regard, garbage word models are included to accommodate the out-of-vocabulary words that might appear. For this, the responses recorded from each speaker are split into 15 groups based on the number of the phonemes in the word. A garbage model is trained on each of these groups of words, with the number of states equal to thrice the number of phonemes in the smallest word in the group and 8 mixture components per state. This system was observed to perform better than the rest.

Analysis was then carried out on the use of a single recognition system/common dictionary against the use of a different recognition system for each question, that is, the use of multiple dictionaries. Further, the effect of vocal tract length normalization on the performance of the system was also analyzed. In this regard, test data consisting of 11 sentences each from 12 speakers were recorded, resulting in a total of 132 test utterances. The results of these analyses are described below:

#### Using a Common Dictionary

The results obtained when a common dictionary was used, for different number of mixture components per state, are shown in Table 4.1. As mentioned earlier, using 8 mixture components per state yielded the best result, where 84.1% of responses were correctly identified.

| No. of mixture components | Performance% |
|:---:|:---:|
| 2 | 84.1 |
| 4 | 82.6 |
| **8** | **84.1** |
| 16 | 82.6 |

Table 4.1: Performance with a common recognizer. Performance% indicates the percentage of responses that were recognized correctly

#### Using Multiple Dictionaries

Using multiple dictionaries yielded better results compared to using a common one. Since only a particular set of responses can be given to a particular

question, it makes sense to have multiple recognizers in the system, one for each set of responses. The results are depicted in Table 4.2 and it is observed that 98.4% of the responses were recognized correctly, when 8 mixture components were used per state.

| No. of mixture components | Performance% |
|:---:|:---:|
| 2 | 96.2 |
| 4 | 97.7 |
| **8** | **98.4** |
| 16 | 96.9 |

Table 4.2: Performance with multiple recognizers. Performance% indicates the percentage of responses that were recognized correctly.

**Vocal Tract Length Normalization**

Vocal Tract Length Normalization (VTLN) is a means by which the possible reduction in the performance of a recognition system, due to variations in vocal tract length among the users, can be avoided. The difference in vocal tract length can be compensated by changing a factor called the warping factor ($\alpha$). The warping factor changes the way the filter banks are positioned across the frequency spectrum when features are derived. Its value ranges from 0.8 to 1.2.

VTLN can be performed, (i) during training alone, (ii) during testing alone, or (iii) during training as well as testing. In this analysis, VTLN has been performed during the training phase alone. Various $\alpha$ values were set for each speaker and the performance of the recognition system was evaluated. The $\alpha$ value that resulted in the best performance was chosen as the optimal value for that speaker. This process was repeated for all speakers. The results of this experiment are as follows:

| No. of mixture components | Performance% |
|:---:|:---:|
| 2 | 97.7 |
| 4 | 98.4 |
| 8 | 98.4 |
| 16 | 96.9 |

Table 4.3: Performance with multiple recognizers and VTLN. Performance% indicates the percentage of responses that were recognized correctly.

## 4.2 Final Recognition Systems

Based on the analysis in the previous section, triphone-based recognition systems were trained for all questions posed by the enquiry system. The triphone models were trained with 3 states and 8 mixture components per state. Since the use VTLN did not significantly alter the performance of the recognition system, it was not used in the final systems. The recognition systems were trained on the entire data collected and consisted of a total of 1509 triphone models.

## 4.3 Likelihood analysis

To ensure that response recognized by the speech recognition system is accurate, the enquiry system could confirm the recognized response with the user every time. However, this might be annoying to the user. Therefore, to restrict the system to confirming a response only in the event of a doubt, a likelihood ratio-based analysis was performed. The goal was to confirm a response with the user only when the likelihood ratio exceeded a threshold. This likelihood ratio was calculated by deriving the two best possible recognition results for each utterance and calculating the ratio between the log-likelihoods with which each word/result was recognized. Based on the value of this ratio and a threshold, a decision was made as to whether a confirmation question is posed to the user or not. This is elaborated below:

### 4.3.1 Training phase

In the training phase, the threshold to be used is computed for each word. In this regard, all examples of each word in the training data are recognized by the appropriate recognition systems, and the two best recognition results and the ratio between the corresponding log-likelihoods ($R_{w_i}$) are obtained. The threshold, $T_{w_i}$, to be used for each word could be set to the minimum, maximum, or mean value of the ratios obtained for each word, or a scaled version of either of the three values, as shown below:

$$T_{w_i} = cf(R_{w_i}) \tag{4.1}$$

where $c$ is a scaling factor and $f(R_{w_i})$ is either the minimum, maximum, or mean of all ratios corresponding to the word $w_i$.

### 4.3.2 Testing phase

During the testing phase, a two step process is used to make a decision. First the boundary for the target word is estimated. Secondly, using the boundaries detected in the first step, 2-best recognition is performed on the target word, to derive the log-likelihood ratio. This ratio is then compared with the threshold calculated in the training phase. When the ratio exceeds the threshold, the confirmation question is asked, else the first best recognized word is used to proceed to the next stage.

### 4.3.3 Setting an appropriate threshold

The following four conditions may occur when using a likelihood ratio-based approach:

1. First best recognition is wrong and confirmation is posed to the user

2. First best recognition is wrong but confirmation is not posed to the user

3. First best recognition is right but confirmation is posed to the user

4. First best recognition is right and confirmation is not posed to the user

Here cases 1 and 4 are true cases and do not cause a problem. But cases 2 and 3 are cases that need attention. Case 3 is not problematic, but if the user is continually asked for confirmations, then the use of this method will not be maximized. Case 2 is problematic and should be avoided at all costs because the user will be navigated to another flow if the wrong word is considered to be right, without a confirmation. Hence the key rule for deciding the threshold is that, the case 2 should be avoided at all costs while reducing case 3 as much as possible. In this regard, selecting the mean of $R_{w_i}$ as the threshold, appeared to be the best choice.

# Enquiry System

The enquiry system, as mentioned before, is the combination of three components, namely, (i) a speech recognition system that recognizes the user's speech and converts it to text, (ii) a database that contains relevant information, and (iii) a speech synthesis system that converts text to speech and plays it to the user. In this regard, the final recognition systems developed in Chapter 4 and the data collected in Chapter 3 have been used along with an HMM-based text-to-speech synthesis system. This synthesizer was trained on five hours of data, recorded from a male, native-Tamil, professional speaker, in a studio environment. The synthesizer was developed as a part of a MeitY (Ministry of Electronics and Information Technology)-funded project titled, "Development of Text-to-Speech System for Indian Languages - High Quality TTS and Small Footprint TTS Integrated with Disability Aids", which was a joint venture with 11 other organizations, with IIT Madras at the head.
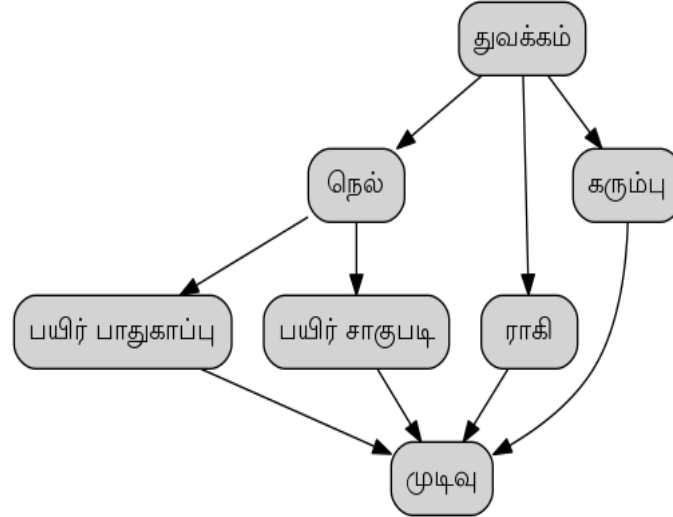
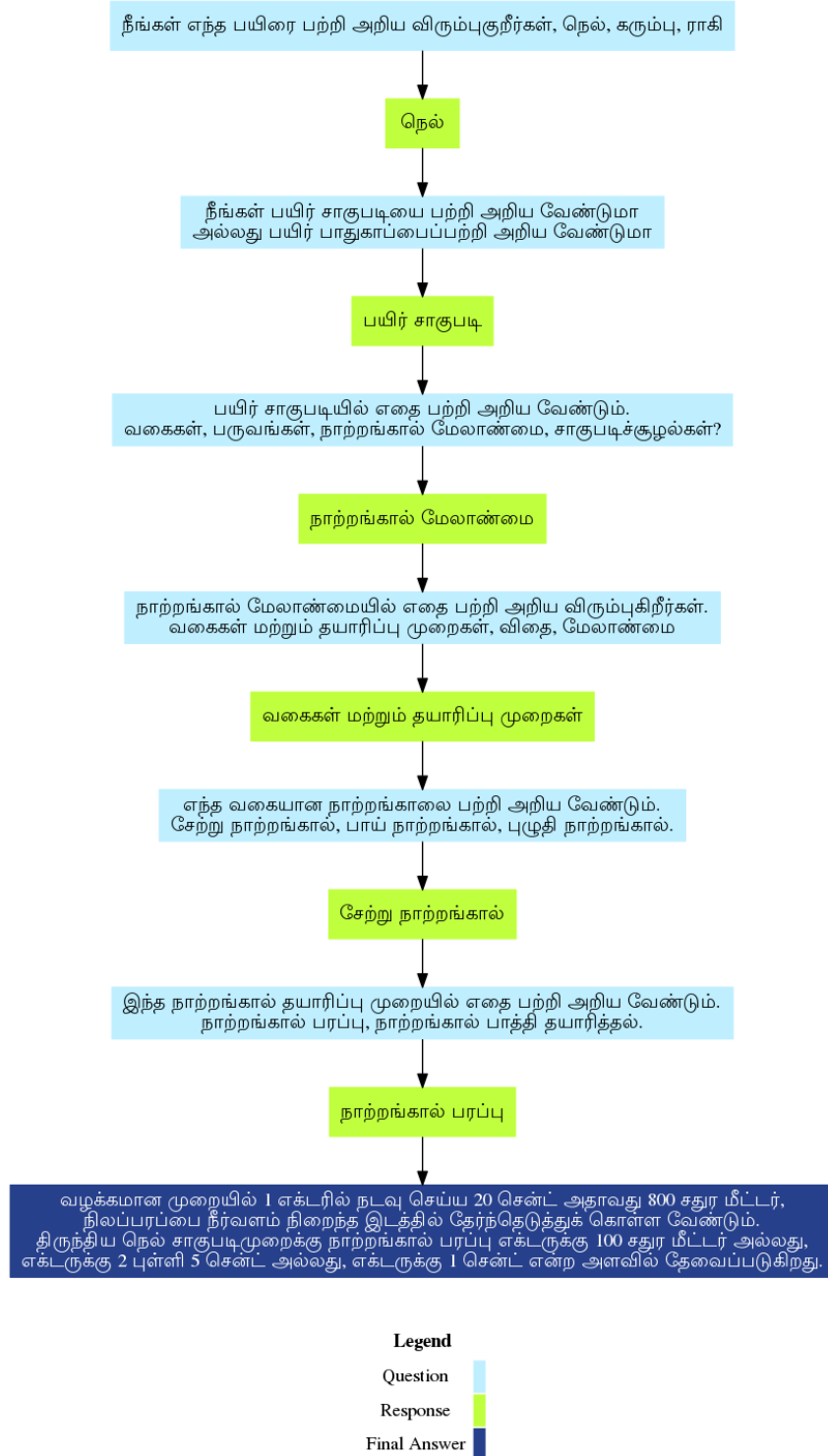Figure 5.1: Enquiry System - Overview

Figure 5.2: Sample question-answer flow

An abstract flowchart for the enquiry system capable of providing information on paddy production and protection and sugarcane and ragi production is shown in Fig. 5.1. The flow has been abstracted for simplicity. Each leaf node in the flow expands and corresponds to its own flow, described in Chapter 3. A sample question-answer flow in the enquiry system is depicted in Fig. 5.2. It shows the question that is asked by the system to the user along with corresponding options, the response given by the user to the system, and the final answer, which is the information provided by the system based on the series of questions answered. The questions and answers in the enquiry system have been framed in literary Tamil.

The components of the enquiry system are linked to each other by shell scripts. They are presented to the user via a web interface that uses Javascript, Ajax, jQuery, and Php to fetch/process data from the server and work in sync with the shell scripts to dynamically update the app without any page reloads. Fig. 5.3 - 5.6 show some screen-shots of the web app for reference. CD attached contains a copy of the app with all the source code. A guide explains how to set it up for use.
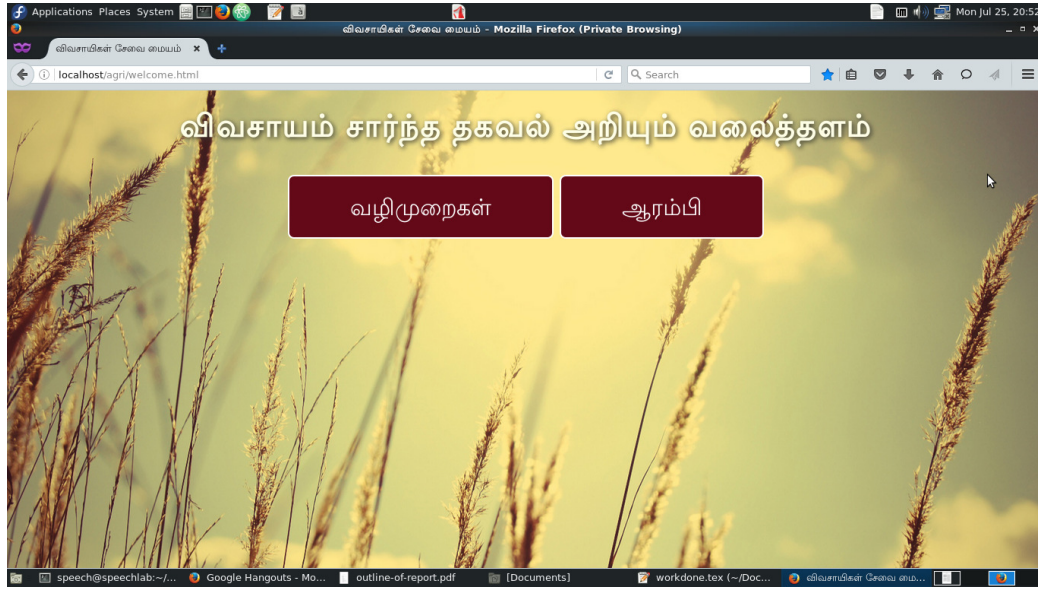


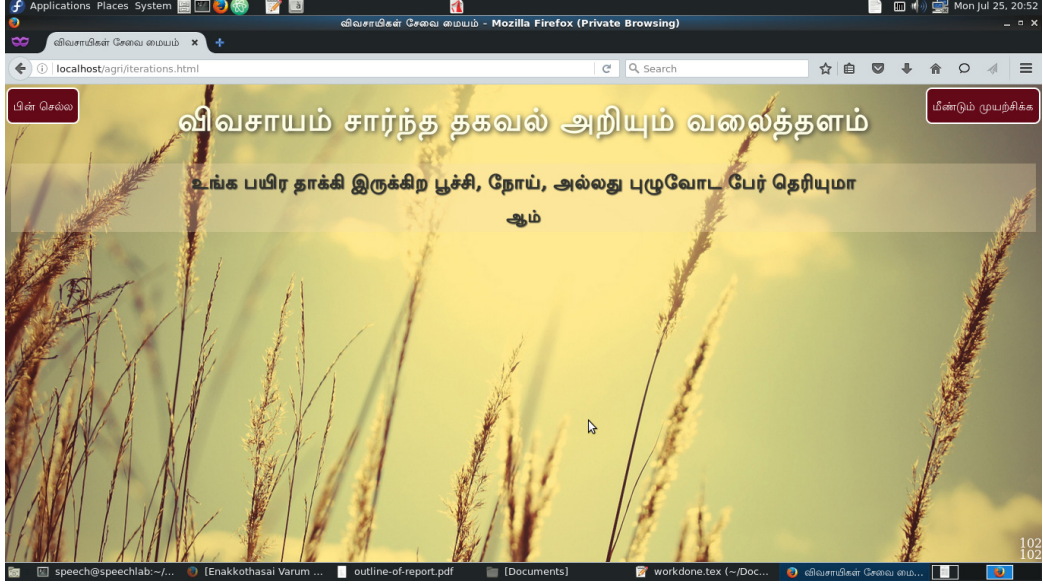Figure 5.3: Welcome Screen of the Web App

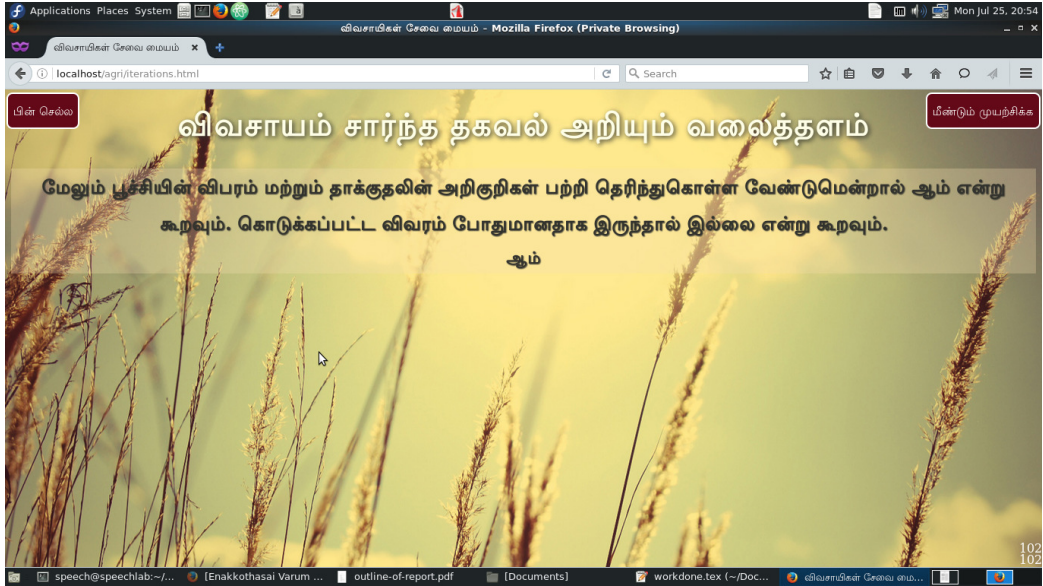Figure 5.4: Example of a Question-Answer iteration being displayed in the app



Figure 5.5: Another example of a Question-Answer iteration being displayed in the app
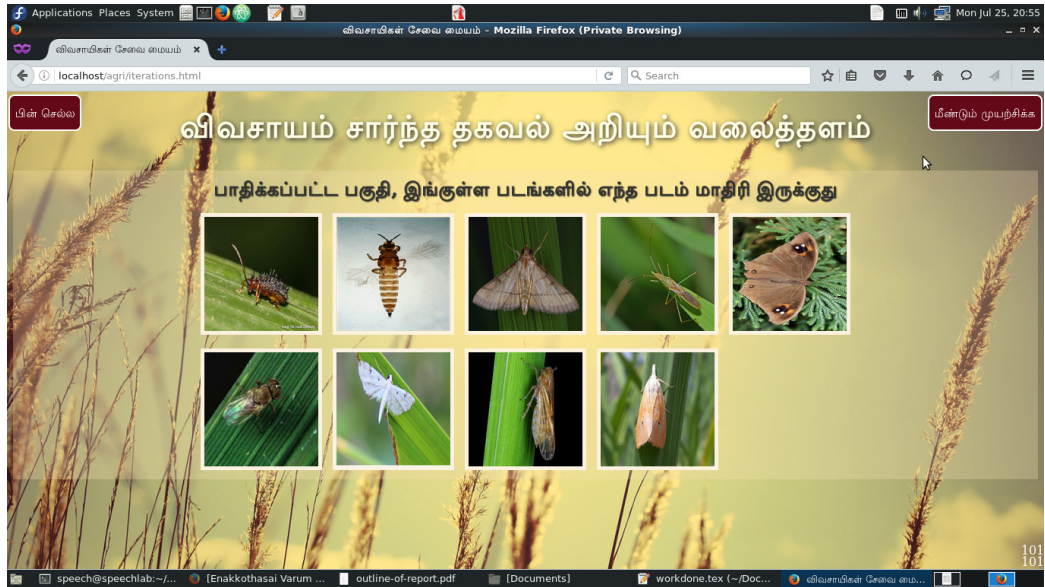
Figure 5.6: Example of images being displayed in the app

# Financial Details

An amount of Rs. 9.52 lakh was sanctioned for the project. Till date, an amount of Rs. 7.71 lakh (after 10% TDS deduction) has been received. Distribution of this fund, as mentioned in the project proposal, is portrayed in Table 6.1. Accordingly, an amount of Rs. 86,500 was directed towards institutional overhead. Till date, an amount of Rs. 12,143 has been spent towards travel expenses. Since the project was initially sanctioned for 6 months but carried out over the course of a year, the remaining amount has been used for manpower.

Table 6.1: Distribution of Funds

|  | Item | BUDGET (in Rupees) | | | |
|---|---|---|---|---|---|
|  |  | 1st phase (2 months) | 2nd phase (3 months) | 3rd phase (1 month) | Total |
| A. | Recurring |  |  |  |  |
|  | 1. Salaries/wages | 1,60,000 | 2,40,000 | 80,000 | 4,80,000 |
|  | 2. Consumables | 20,000 | 20,000 | 20,000 | 60,000 |
|  | 3. Travel | 15,000 | 15,000 | 10,000 | 40,000 |
|  | 4. Miscellaneous | 20,000 | 30,000 | 20,000 | 70,000 |
|  | 5. Contingency | 10,000 | 15,000 | 5,000 | 30,000 |
| B. | Non-recurring |  |  |  |  |
|  | 1.High-end Server | 1,50,000 | - | - | 1,50,000 |
|  | 2. Laser Printer | 20,000 | - | - | 20,000 |
| C. | Consultancy | 5,000 | 7,500 | 2,500 | 15,000 |
| D. | Overhead (10% of total) | 40,000 | 32,750 | 13,750 | 86,500 |
|  | Grand total | **4,40,000** | **3,60,250** | **1,51,250** | **9,51,500** |

# Summary

In summary, a speech-enabled interactive enquiry system has been developed to provide information on the production and protection of paddy and on the production of sugarcane and ragi. The enquiry system consists of a speech recognition system, a database containing relevant information, and a text-to-speech synthesis system. Initially, the agritech portal and TNAU's expert system were analyzed and possible questions and responses were formulated. Information to be stored in the database was also collected from these sources. Speech data to be used to train the recognition system was then collected from multiple speakers. A triphone-based recognition system was then developed with this data. To ensure that the system did not proceed with an incorrect recognition result, a likelihood ratio-based confirmation system was employed, which confirmed the recognition result with the user, in the event of a doubt. For synthesis, an HMM-based text-to-speech synthesis system, developed as a part of a MeitY-funded project was used. Finally, the components of the enquiry system were integrated into a web-based application.