

Massachusetts Institute of Technology
Department of Electrical Engineering & Computer Science

6.345 Automatic Speech Recognition
Spring, 2003

Issued: 02/14/03
Due: 02/26/03

Assignment 2 Syllable Structure

A language is not only limited by the inventory of basic sound units, but also by the allowable combinations of these sounds. This assignment is intended to give you some feeling about such constraints.

To do this, we will use an interactive software facility named *Crystal*, which runs on the Linux workstations. *Crystal* is an interactive system which provides many functions for studying and displaying the distributional constraints of a lexicon. For the purpose of this lab, we will use the Merriam Pocket dictionary, which contains about 20,000 entries, as the working lexicon. To start the lab simply enter the command:

```
% start_lab2.cmd
```

Distributional Properties

We will begin our investigation by examining some of the distributional properties of this lexicon of English words.

T1: In this exercise, we will study the properties of the most common words in the English language. Click on **Sort by Brown Corpus Frequency (BCF)**¹ in the **Search Results** sub-window, which will sort the words in the dictionary according to their number of occurrences in the Brown Corpus. Study the counts and properties of the top 15 words in the list.

Q1: What are the common characteristics of the 15 most frequent words (e.g., number of syllables, part of speech, etc.?)

¹The Brown Corpus is a corpus of over one million words gathered at Brown University. These words were taken from various sources such as books, papers and magazines, and their frequencies of occurrence were recorded.

T2: In this exercise, we will study the properties of the most frequent two and three syllable words in the English language. Set the **Search Type** to **stress** and type in . . (. ?) in **Search String**. Note that all characters in the search string are separated by spaces. The first two dots match two syllables, while the third dot a question mark in parentheses matches an optional third syllable.

Q2: What are the most frequent two and three syllable words, and how highly are they ranked in the lexicon? When looking at only two syllable words by using . . as the search string, which syllable is more likely to be stressed? For the second part, use **S** to match a stressed syllable.

T3: In this exercise, we will study the distribution properties of syllable patterns for English. Restore the original lexicon by clicking on it in the history sub-window. Click on **Syllables per Word** in the **Statistics** sub-window.

The distribution of the syllable patterns in the Brown Corpus is different from that in the dictionary, because some words in the dictionary occur more often than others. To weight words by their Brown Corpus frequencies click on **Weight by BCF** in the **Statistics** sub-window. The **Syllables per Word** graph should now be weighted by Brown Corpus frequencies.

Q3: It turns out that all of the words in the lexicon contain eight or fewer syllables. What is the most frequent number of syllables per word? Describe the probability distribution for **Number of Syllables per Word**. How would your answer differ when the words are weighted by their Brown Corpus frequencies?

T4: In this exercise, we will study the distribution of stress patterns for English. Click on **Stress Pattern Occurrences** in the **Statistics** sub-window. Also, view the distribution as weighted by Brown Corpus frequencies.

Q4: What is the most frequent polysyllabic stress pattern? How would your answer differ when the words are weighted by their Brown Corpus frequencies?

T5: In this exercise, we will study the distribution properties of phonemes for English. Click on **Phoneme Occurrences** in the **Statistics** sub-window. Also, view the distribution as weighted by Brown Corpus frequencies.

Q5: Of the ten most frequently occurring phonemes in the lexicon, what are the most common manner of production and place of articulation? How would your answer differ when the words are weighted by their Brown Corpus frequencies?

Phonotactical Rules

The study of allowable sound sequences of a language are called **phonotactics**. This part of the assignment exposes you to some of the common phonotactical rules of English.

Syllable Template

The understanding of phonotactical rules can be furthered by knowledge of syllable structure. Figure 1 shows a diagram of an accepted syllable template, and Figure 2 some examples.

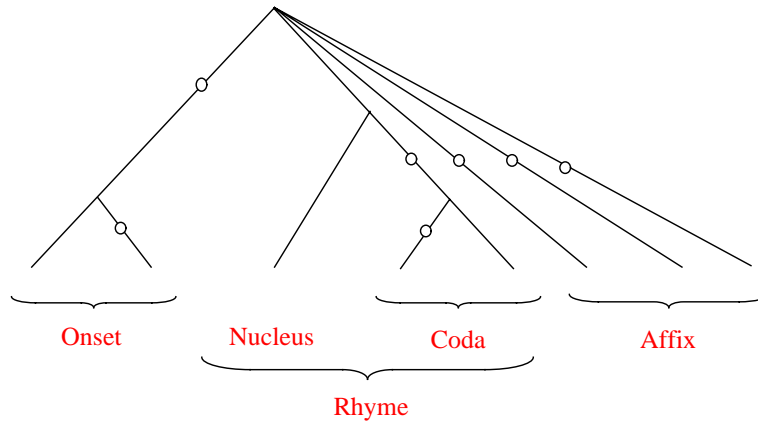


Figure 1: Syllable template: Fudge, "Syllables", J. Linguistics, 1969.

- Branches marked by \circ are optional
- Nucleus must contain a non-obstruent
- Sonority decreases away from nucleus
- Affix contains only coronals: /s, z, t, d, θ, ð, ç, j/
- Only the last syllable in a word can have an affix
- /sp/, /st/, and /sk/ are treated as single obstruents

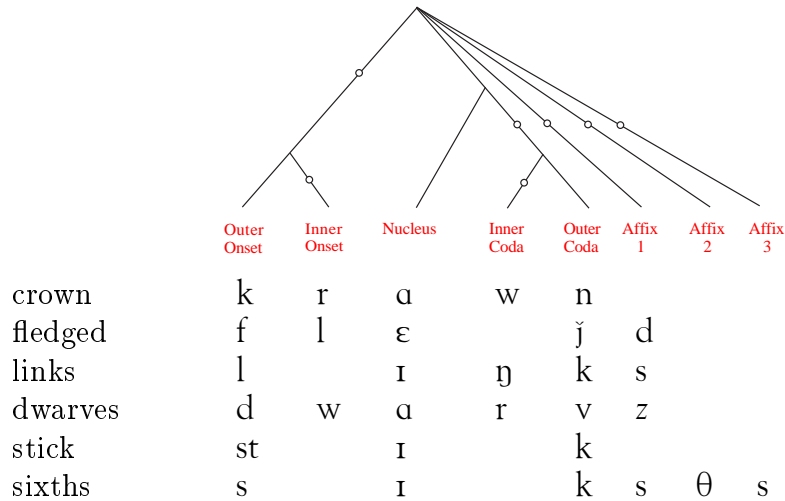


Figure 2: Some syllable examples.

Consonant Clusters

There are only a limited number of distinct word-initial and word-final consonant clusters in the English language. We will study their properties in this part of the lab.

T6: First, restore **Search Type** to **phonemic**. Search for word-initial consonant clusters in the original lexicon containing at least two consonants by typing **C C (C *) V . *** in **Search String**. The **C C (C *)** portion matches two or more consonants, while the **V** portion matches exactly one vowel. Finally, the **. *** portion matches the remaining zero or more phonemes of an arbitrary word. Pay special attention to the existence of /tk/ and /kt/ clusters.

Next, restore the original lexicon by clicking on it in the history sub-window. Search for all possible word-final consonant clusters in the lexicon by typing **. * V C *** in **Search String**. Pay special attention to the existence of /tk/ and /kt/ clusters.

Q6: We know that no word in the dictionary contains consonant cluster /tkt/ or /ktk/ (you can verify this by searching the lexicon with **. * t k t . *** or **. * k t k . ***). Are the following two phonemic transcriptions possible?

(a) /... t k t .../

(b) /... k t k .../

What is the maximum length of a word-initial consonant cluster? At this length, how many consonant clusters are there and what are they?

Vowel Clusters

T7: Search for words with two adjacent vowels by typing **. * V V . *** in **Search String**. Be sure to restore the original lexicon and ignore syllable boundaries by enabling **Ignore Syllable Boundaries**.

Q7: How many words have two vowels in a row? How many of them have a schwa as the second vowel? How many have a schwa as the first vowel? Use **(ax | ix)** to match both plain or front schwas. What do two adjacent vowels imply about the syllable structure of the two syllables to which they belong?

Homorganic Rules

T8: The homorganic nasal-stop rule states that nasal-stop clusters must agree on the place of articulation. Verify this by examining all the occurrences of nasal-stop clusters in the lexicon. You can search for all words containing nasal-stop sequences by typing **. * NASAL STOP . *** in **Search String**. You can also search for more specific examples in the resulting sub-lexicon. For example, to search for words containing /nd/, type **. * n d . *** in **Search String**; to search for words containing either /nd/ or /nt/, type **. * n (d | t) . *** in **Search String**. You will want to experiment how ignoring or accounting for syllable boundaries affects the results.

Q8: How often is nasal-stop homorganic rule violated? Can you try to generalize a rule to summarize when it is broken.

Lexical Constraints

In this part of the lab you will investigate the extent to which a given word can be disambiguated from competitors based on partial phonetic information.

T9: You have done some spectrogram reading practice in class. In this exercise, we will show that the use of lexical access can greatly assist the task. In the Figures 3, 4, and 5, you will find three spectrograms of isolated words. Start with a very coarse transcription of the spectrogram by hand. If you can not determine the phones, try to come up with phone classes such as vowel, nasal, strong fricatives, voiced stop, etc. Perform a search on the lexicon based on your partial hypothesis. If you can not determine the words, try to refine your hypothesis and search again. The search pattern should be expressed as *regular expressions*, many examples of which have already been given in the previous tasks. The following classes have been defined along with abbreviations, or you can use the OR operator, |, to create custom classes. Enable Ignore Syllable Boundaries so that you will not have to explicitly specify syllable boundaries.

CLASS	ABBREVIATION	MEMBERS
VOWEL	V	all vowels
RETROFLEXED FRICATIVE	R	r axr er
STRONG-FRICATIVE	F	s sh z zh f th v dh
WEAK-FRICATIVE	SF	s sh z zh
NASAL	WF	f th v dh
GLIDE	N	m n ng
LIQUID	G	w y
SEMIVOWEL	L	l r
ASPIRANT	SV	l r w y
STOP		hh
VOICED-STOP	S	b d g p t k
UNVOICED-STOP	VS	b d g
AFFRICATE	US	p t k
SYLLABIC-CONSONANT	A	ch jh
	SC	el em en

Q9: What are the words in each spectrogram? What is the partial phonetic hypothesis you have that leads you to the answer with the help of lexical search?

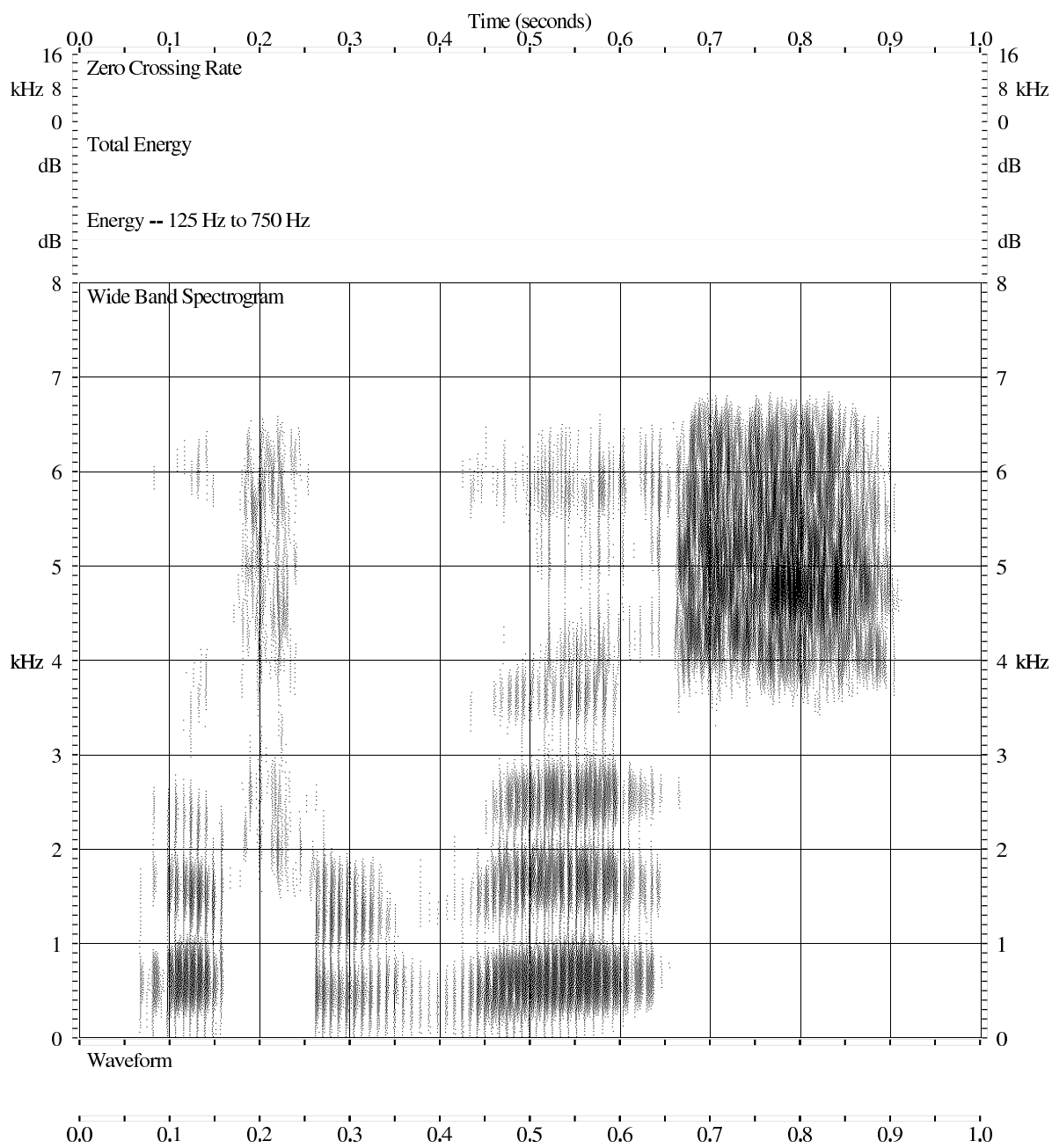


Figure 3: Mystery word #1.

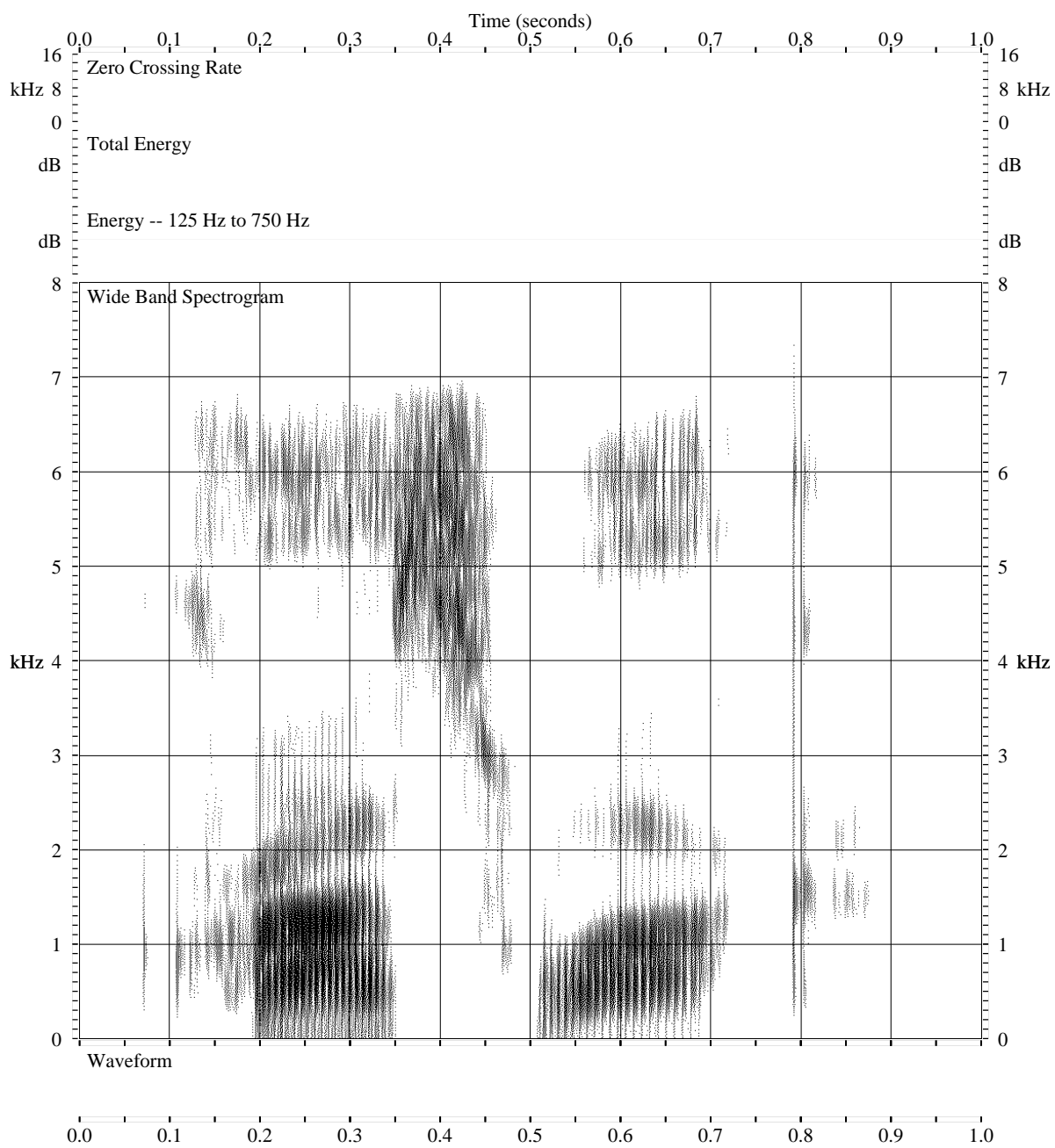


Figure 4: Mystery word #2.

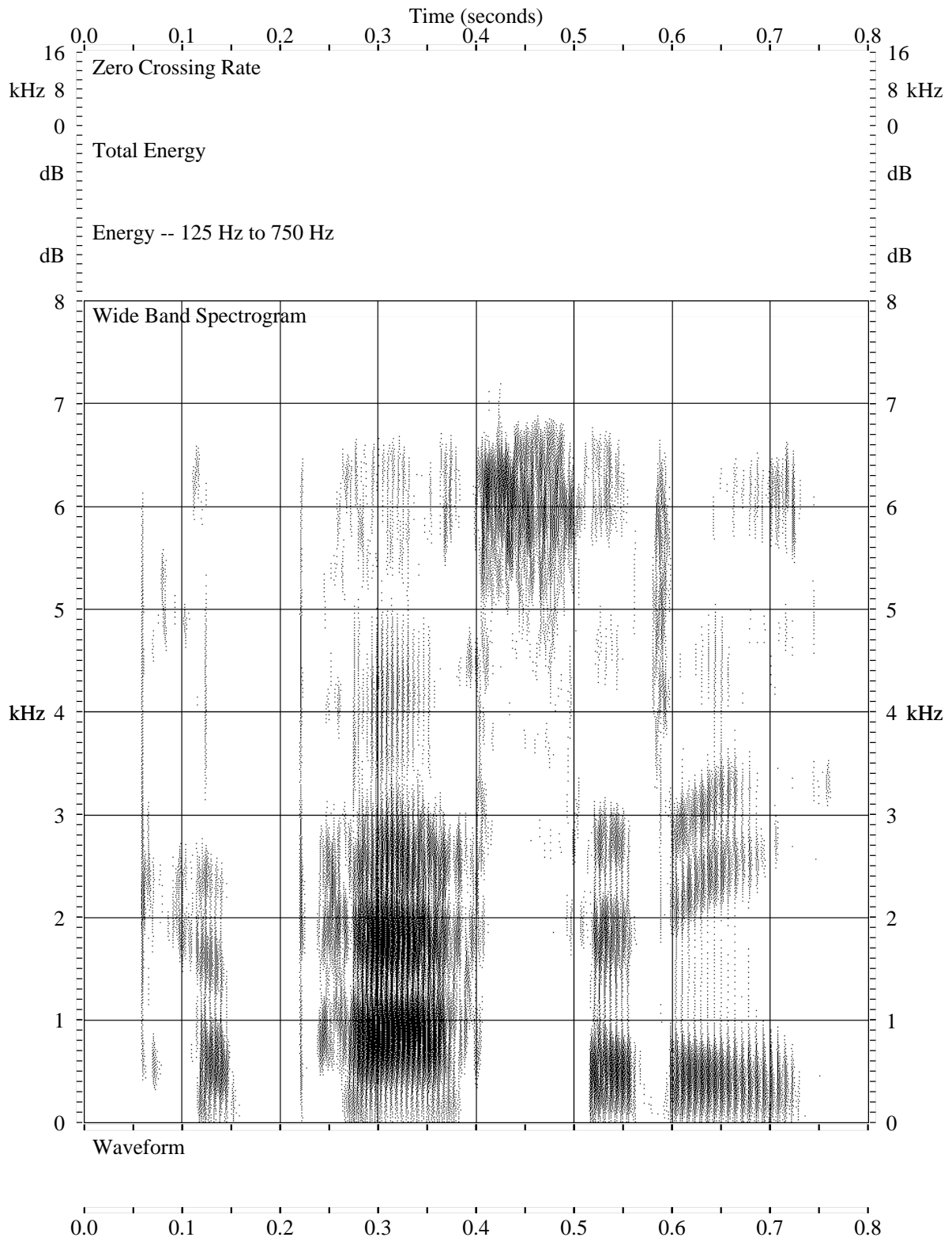


Figure 5: Mystery word #3.