


An Algorithm for Determining the Endpoints for Isolated Utterances


L.R. Rabiner and M.R. Sambur

The Bell System Technical Journal, Vol. 54,
No. 2, Feb. 1975, pp. 297-315




Outline

- Intro to problem
- Solution
- Algorithm
- Summary

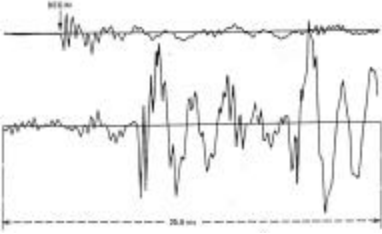


Motivation

- Word recognition needs to detect word boundaries in speech
- Recognizing silence can reduce:
 - Processing load
 - (Network not identified as savings source)
- Easy in sound proof room, with digitized tape




Visual Recognition

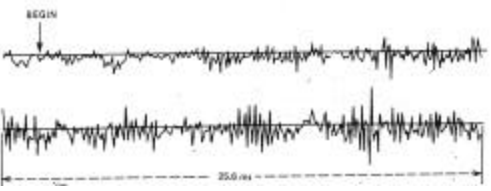


"Eight"

- Easy
- Note how quiet beginning is (tape)




Slightly Tougher Visual Recognition

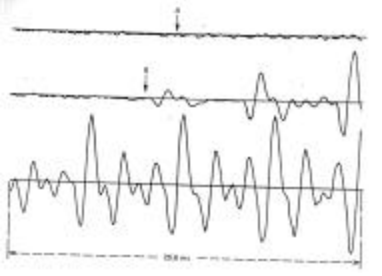


"Six"

- "sss" starts crossing the 'zero' line, so can still detect




Tough Visual Recognition



"Four"

- Eye picks 'B', but 'A' is real start
 - /f/ is a *weak fricative*



Tough Visual Recognition

"Five"

- Eye picks 'A', but 'B' is real endpoint
 - V becomes *devoiced*

Tough Visual Recognition

"Nine"

- Difficult to say where final trailing off ends

The Problem

- Noisy computer room with background noise
 - Weak fricatives: /f, th, h/
 - Weak plosive bursts: /p, t, k/
 - Final nasals
 - Voiced fricatives becoming devoiced
 - Trailing off of sounds (ex: binary, three)
- Simple, efficient processing
 - Avoid hardware costs

The Solution

- Two measurements:
 - Energy
 - Zero crossing rate
- Simple, fast, accurate

Energy

- Sum of magnitudes of 10 ms of sound, centered on interval:
 - $E(n) = \sum_{i=-50 \text{ to } 50} |s(n+i)|$

Zero (Level) Crossing Rate

- Number of zero crossings per 10 ms
 - Normal number of cross-overs during silence
 - Increase in cross-overs during speech

The Algorithm: Startup

- At initialization, record sound for 100ms
 - Assume 'silence'
 - Measure background noise
- Compute average (IZC') and std dev (σ) of zero crossing rate
- Choose IZCT (zero-crossing threshold)
 - Threshold for unvoiced speech
- $IZCT = \min(25 / 10\text{ms}, IZC' * 2 \sigma)$



The Algorithm: Thresholds

- Compute energy, $E(n)$, for interval
 - Get max, IMX
 - Have silence, IMN
 - $I1 = 0.03 * (IMX - IMN) + IMN$
(3% of peak energy)
 - $I2 = 4 * IMN$
(4x silent energy)
- Get energy thresholds
 - $ITL = \min(I1, I2)$
 - $ITU = 5 * ITL$



The Algorithm: Energy Computation

- Search sample for energy greater than ITL
 - Save as start, say s
- Search for energy greater than ITU
 - s becomes start
 - If energy falls below ITL, restart
- Results in conservative estimates
 - Endpoints may be outside

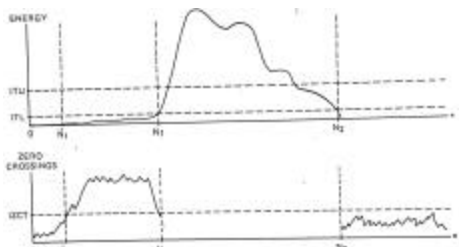


The Algorithm: Zero Crossing Computation

- Search back 250 ms
 - Count number of intervals where rate exceeds IZCT
 - If 3+, set starting point, s, to first time
 - Else s remains the same
- Do similar search after end



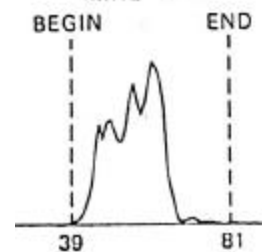
The Algorithm: Example



(Word begins with strong fricative)



Algorithm: Examples



"Half"

- Caught trailing /f/



Algorithm:
Examples

"Four"

Notice how different each "four" is

Evaluation: Part 1

- 54-word vocabulary
- Read by 2 males, 2 females
- No gross errors (off by more than 50ms)
- Some small errors
 - Losing weak fricatives
 - None affected recognition

Evaluation: Part 2

- 10 speakers
- Count 0 to 9
- No errors at all

Evaluation 3: Your Project 1

Future Work

- Three classes of speech:
 - Silence
 - Unvoiced speech
 - Voiced speech
- May be more computationally intensive solutions that are more effective