

Automatyczne rozpoznawanie mowy

Autor: mgr inż. Piotr Bratoszewski

Rys historyczny

- 1930-1950 – pierwsze systemy Automatycznego rozpoznawania mowy (ang. *Automatic Speech Recognition* – ASR), metody holistyczne; „ad-hoc”; izolowane słowa; małe słowniki; Bell Laboratories
- 1950-1960 – pierwsze systemy ASR oparte na zależnościach fonetycznych; małe słowniki;

Rys historyczny

- 1960-1980 – systemy oparte o rozpoznawanie wzorca (ang. *pattern recognition*); wykorzystanie parametrów kodowania predykcyjnego (LPC); sekwencje izolowanych lub połączonych słów; małe i średnie słowniki
- 1980-2000 – wprowadzenie statystycznego modelowania zależności dynamicznych i statycznych w **mowie ciągłej**; zastosowanie ukrytych modeli Markowa (ang. *Hidden Markov Models* - HMM)

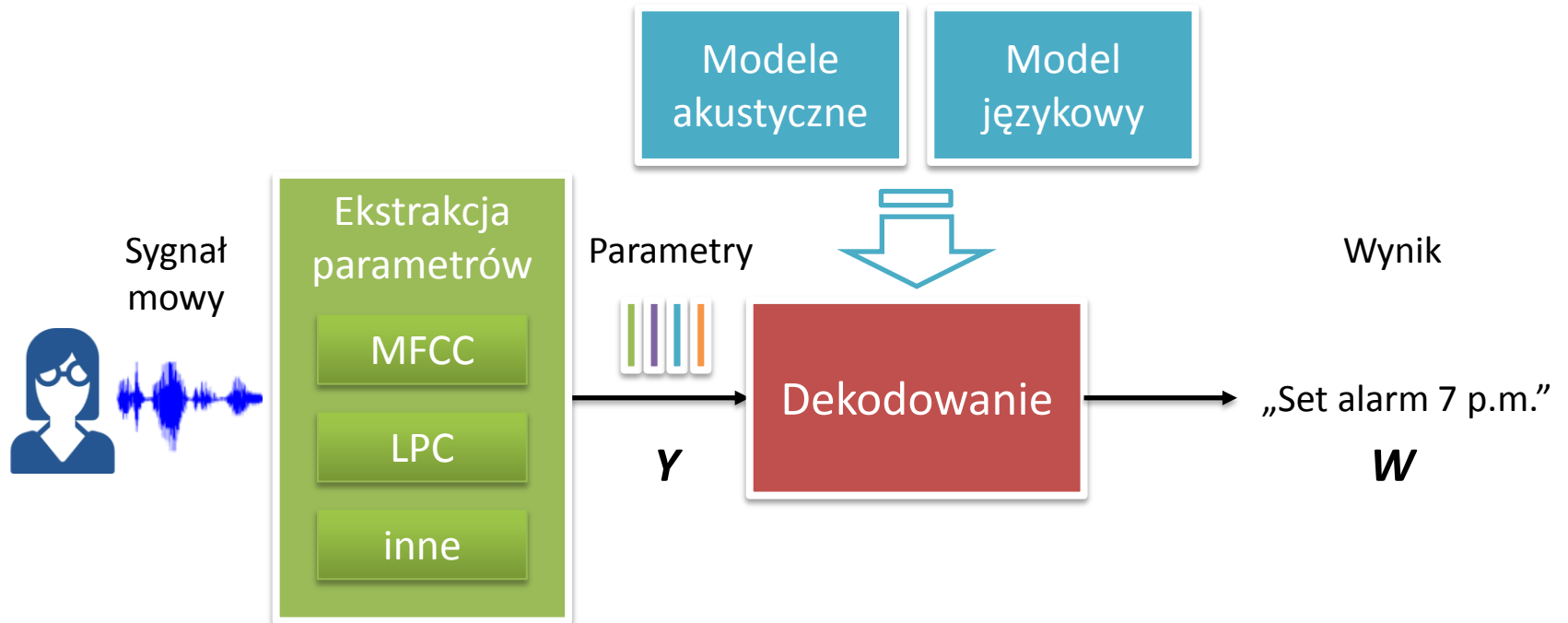
Rys historyczny

- 2000-teraz – kombinacje modeli HMM oraz zależności akustyczno fonetycznych w celu znajdowania i korekcji nieregularności językowych, deep learning, systemy pracujące w chmurze; zwiększanie odporności systemu na pracę w środowisku szumowym; rozpoznawanie wielomodalne

Istotne terminy

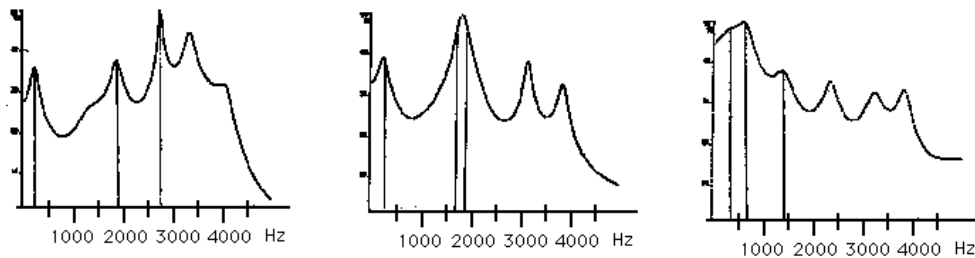
- Słownik – ilość wyrazów które system jest w stanie rozpoznać:
 - Mały słownik: 2 – 100 wyrazów
 - Średni słownik: 100 – 1000 wyrazów
 - Duży słownik: ponad 1000 wyrazów (w tej chwili 50 tys. słów)
- System „zależny/niezależny od mówcy”
- Rozpoznawanie mowy ciągłej/izolowanej
- Składnia (ang. *syntax*) – mowa naturalna/wydawanie poleceń/rozpoznawanie cyfr

Schemat systemu ASR

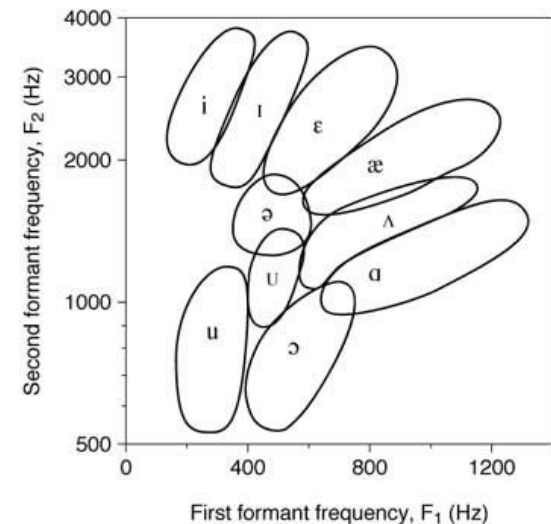


Parametry widmowe

- Podstawowymi parametrami mowy są parametry widmowe uzyskiwane poprzez analizę transformaty Fouriera sygnału mowy
- Analiza rozkładu formantów pozwala na rozpoznawanie samogłosek



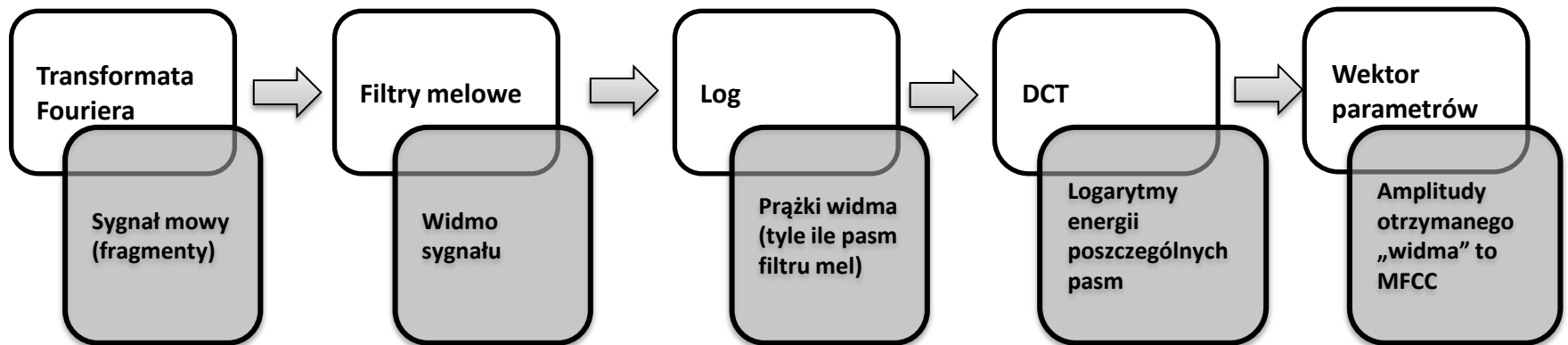
Formanty



Metody parametryzacji mowy

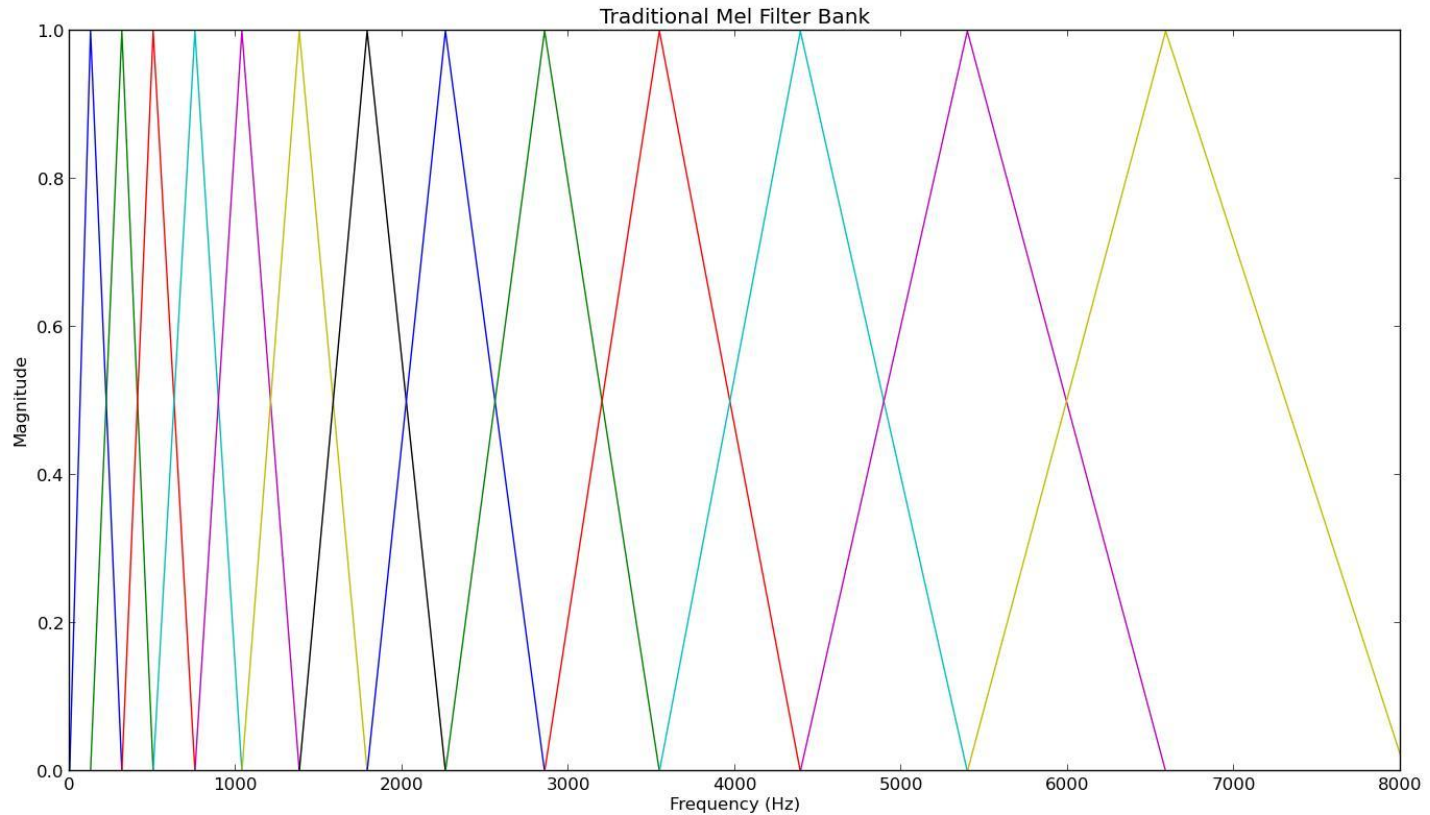
Cepstralne (np. MFCC)

– Metody efektywne i łatwe w implementacji



Metody parametryzacji mowy

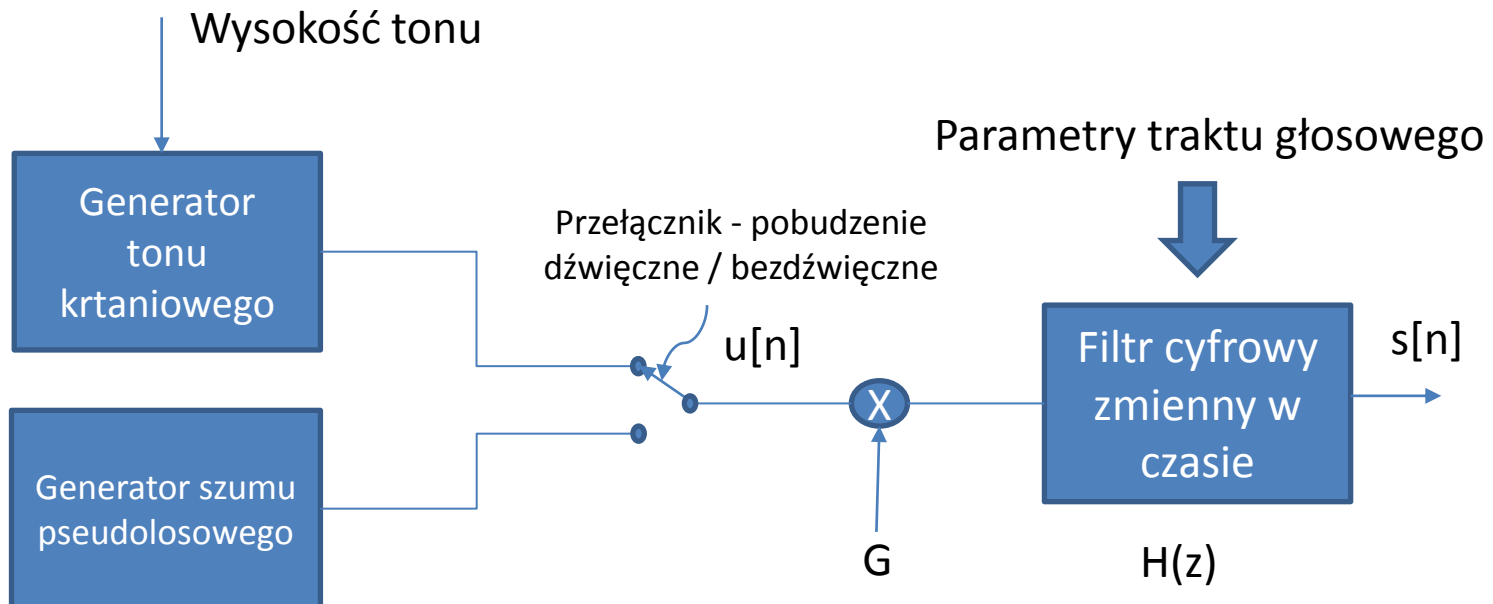
- Filtry w skali melowej



Metody parametryzacji mowy

Predykcyjne (LPC, PLP)

- Skuteczność porównywalna do MFCC w warunkach optymalnych
- Większa skuteczność w przypadku rozbieżności danych testowych z treningowymi



Model akustyczny

- Wymawiane słowa dzieli się na sekwencję podstawowych jednostek mowy – fonemy
- Model akustyczny reprezentuje zależności pomiędzy sygnałem akustycznym a fonemami

BROWSER b r a w z a x

BROWSER b r a w z a x r

CALCULATOR k a e l k y u h l e y t a x

CALCULATOR k a e l k y u h l e y t a x r

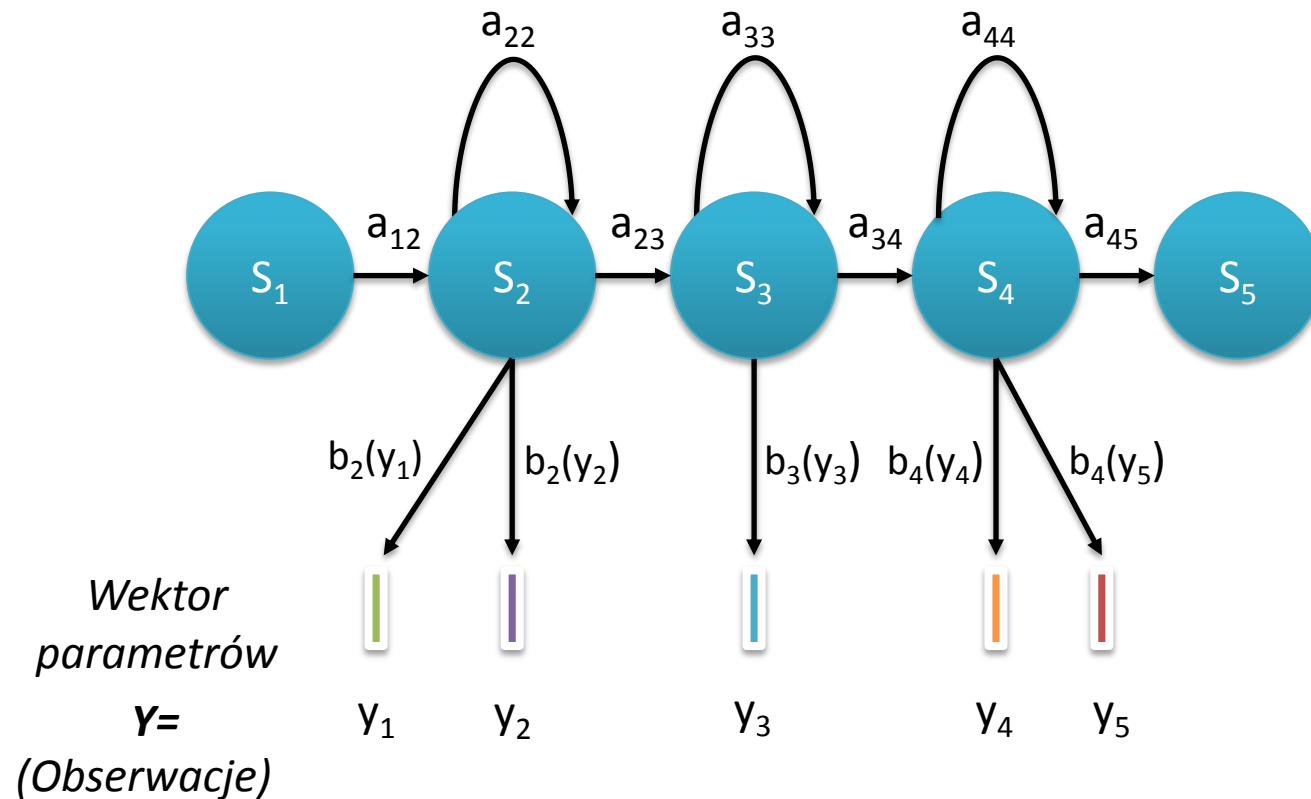
COMPUTER k a x m p y u w t a x

COMPUTER k a x m p y u w t a x r

Model akustyczny

- W celu zamodelowania najkrócej trwających fonemów (głoski wybuchowe – ang. *plosive phonemes*: -p; -t; -k) należy dobrać odpowiednie okno analizy – typowo o długości 10 ms
- Przy tworzeniu modelu akustycznego korzysta się najczęściej z parametrów mel-cepstralnych (MFCC) lub parametrów LPC
- Typowo stosuje się modele trifonowe
- Istotne – osobny model ciszy

Ukryte Modele Markowa



Model *Bakisa*, (ang. *left-right HMM*)

Ukryte Modele Markowa

- Modelowanie procesu na podstawie skończonej liczby stanów S
- Opisywane są przez:
 - N , liczba stanów
 - M , liczba obserwacji
 - Prawdopodobieństwo przejścia pomiędzy stanami $A=\{a_{ij}\}$
 - Prawdopodobieństwo wygenerowania danej obserwacji w stanie j : $B=\{b_j(Y)\}$
 - Rozkład początkowy prawdopodobieństwa $\pi=\{\pi_i\}$

Dekodowanie

- Zadaniem dekodera jest rozpoznanie wymawianego słowa
- Posiadając na wejściu wektory $Y=y_1, \dots, y_T$ dekodery ma za zadanie rozpoznać sekwencję odpowiadających im wyrazów $W= w_1, \dots, w_K$ zgodnie z:

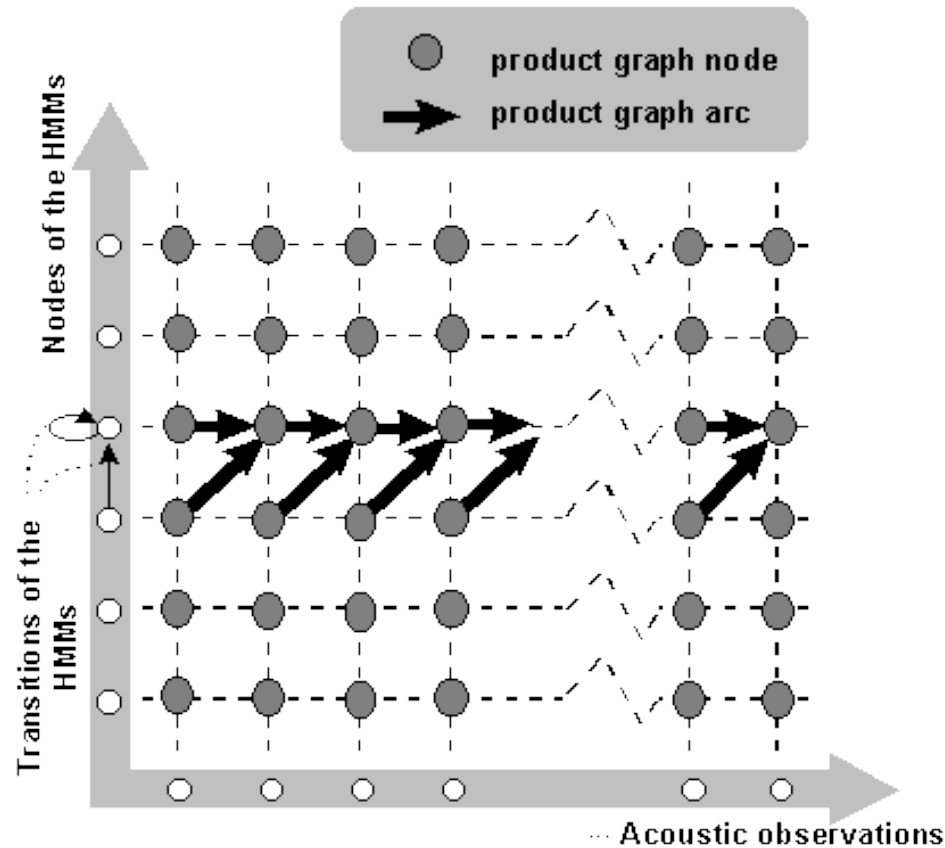
$$\hat{W} = \underset{W}{\operatorname{arg\,max}} [p(W | Y)]$$

Korzystając z twierdzenia Bayesa:

$$\hat{W} = \underset{W}{\operatorname{arg\,max}} [p(Y | W) \cdot p(W)]$$



Dekodowanie



W fazie dekodowania (odkrywania ukrytej sekwencji stanów modelu HMM) często stosuje się algorytm *Viterbiego*

Model językowy

$$p(W) = \prod_{k=1}^K p(w_k | w_{k-1}, \dots, w_{k-N+1})$$

Model N-gramowy – biorący pod uwagę N wyrazów „wstecz”

Bigram:

$$P(\text{ala ma małego kota}) = P(\text{ala} | \langle \text{start} \rangle) \cdot P(\text{ma} | \text{ala}) \\ \cdot P(\text{małego} | \text{ma}) \cdot P(\text{kota} | \text{małego}) \cdot P(\langle \text{end} \rangle | \text{kota})$$

Unigram:

$$P(\text{ala ma małego kota}) = P(\text{ala}) \cdot P(\text{ma}) \cdot P(\text{małego}) \cdot P(\text{kota})$$

Skuteczność rozpoznawania

- Do oceny skuteczności systemu ASR stosowana jest miara „wyrazowej stopy błędu” (ang. *Word Error Rate* – WER)

$$WER = \frac{D+S+I}{H+D+S} \cdot 100\%$$

H – liczba poprawnie rozpoznanych słów

D – liczba nie rozpoznanych słów (ang. *deletions*)

S – liczba błędnie rozpoznanych słów (ang. *substitutions*)

I – liczba wstawionych słów (ang. *insertions*)

Skuteczność rozpoznawania

- Wynik rzeczywistego rozpoznawania

""*/SPEAKER03_C1_AUD1_13.rec"

999999 5889999 TUESDAY

7889999 14239999 WEDNESDAY

16239999 22729999 THURSDAY

24729999 31399999 FRIDAY

33399999 41129999 SATURDAY

43129999 49709999 SUNDAY

51709999 57859999 MOUSE

59859999 64749999 MOVE

66749999 71429999 MUSIC

73429999 78809999 MUTE

""*/SPEAKER03_C1_AUD1_13.rec"

600000 7600000 TUESDAY

7600000 16000000 WEDNESDAY

16000000 24700000 THURSDAY

24700000 33400000 FRIDAY

33400000 43200000 SATURDAY

43200000 51900000 SUNDAY

51900000 55200000 MY

55200000 60200000 AS

60200000 67200000 MOVE

67200000 74000000 MUSIC

74000000 80500000 MUTE

Dynamic Time Warping

- Algorytm DTW – dynamiczne marszczenie czasu

Dwa przebiegi czasowe:

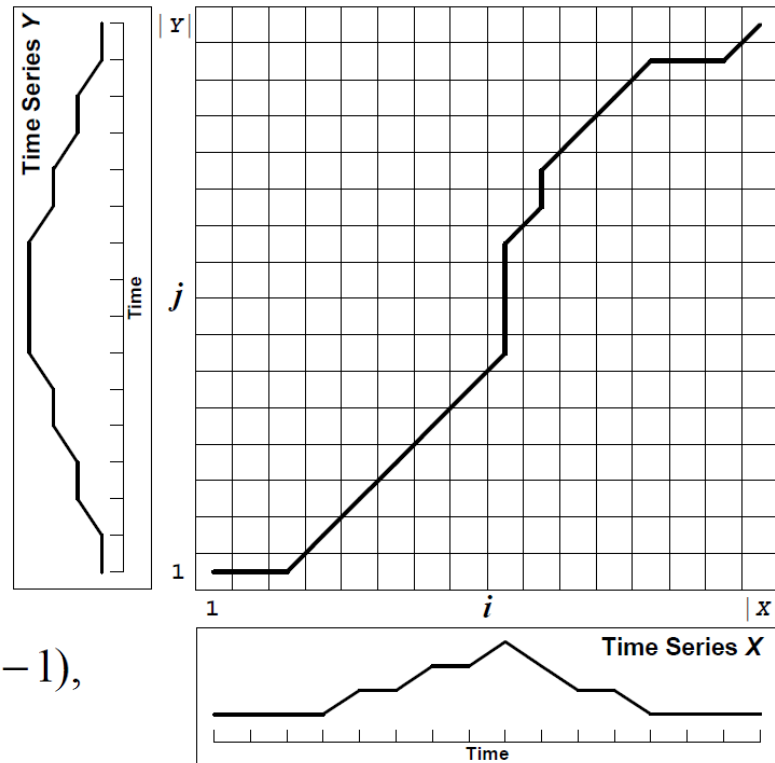
$$X = x_1, x_2, \dots, x_i, \dots, x_{|X|}$$

$$Y = y_1, y_2, \dots, y_j, \dots, y_{|Y|}$$

Tworzenie macierzy kosztów

Obliczanie ścieżki o
najmniejszym koszcie

$$D(i, j) = \text{Dist}(i, j) + \min[D(i-1, j), D(i, j-1), D(i-1, j-1)]$$



Poprawa rozpoznawania w szumie

- Stosowanie filtracji widmowej (ang. *Spectral Subtraction*)
- Stosowanie filtracji adaptacyjnej
- Stosowanie filtracji przestrzennej – *Beamforming*
- Stosowanie różnych modeli akustycznych

Bimodalne rozpoznawanie mowy

- Dołączenie do wektora parametrów akustycznego parametrów wizyjnych
- Fonemy = wizemy
- Widowiskowe podejście – czytanie z ruchu warg
- Teoretycznie wzrost skuteczności w warunkach szumowych
- Wiele problemów do rozwiązania (detekcja ust, framerate, cechy osobnicze)

Przykładowe systemy ASR

Apple Siri



Google Now



Przykładowe systemy ASR

Intel RealSense



Command Mode

Command mode is for issuing commands tied to discrete actions (e.g., saying "fire" to trigger cannon fire in a game).

((("e-mail photo" "fire!" "open" "search")
"mute" "next" "save"))

In command mode, the SDK module recognizes only from a predefined list of context phrases that you have set. The developer can use multiple command lists, which we will call grammars. Good application design would create multiple grammars and activate the one that is relevant to the current application state (this limits what the user can do at any given point in time based on the command grammar used). You can get recognition confidence scores for command and control grammars. To invoke the command mode, provide a grammar.



Dragon Anywhere
Complete work on the go with professional-grade mobile dictation



Dragon Professional Individual
Put your voice to work for greater efficiency on the job



Dragon for Mac
Get the all-in-one Mac solution that let's you dictate, edit and transcribe by voice



Dragon Professional Group
Improve documentation productivity across the enterprise

Przykładowe otwarte systemy ASR



HTK Toolkit



Bibliografia

- HTK Book:
speech.ee.ntu.edu.tw/homework/DSP_HW2-1/htkbook.pdf
- Rabiner L., A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition
- Benesty, Springer Handbook of Speech Processing