

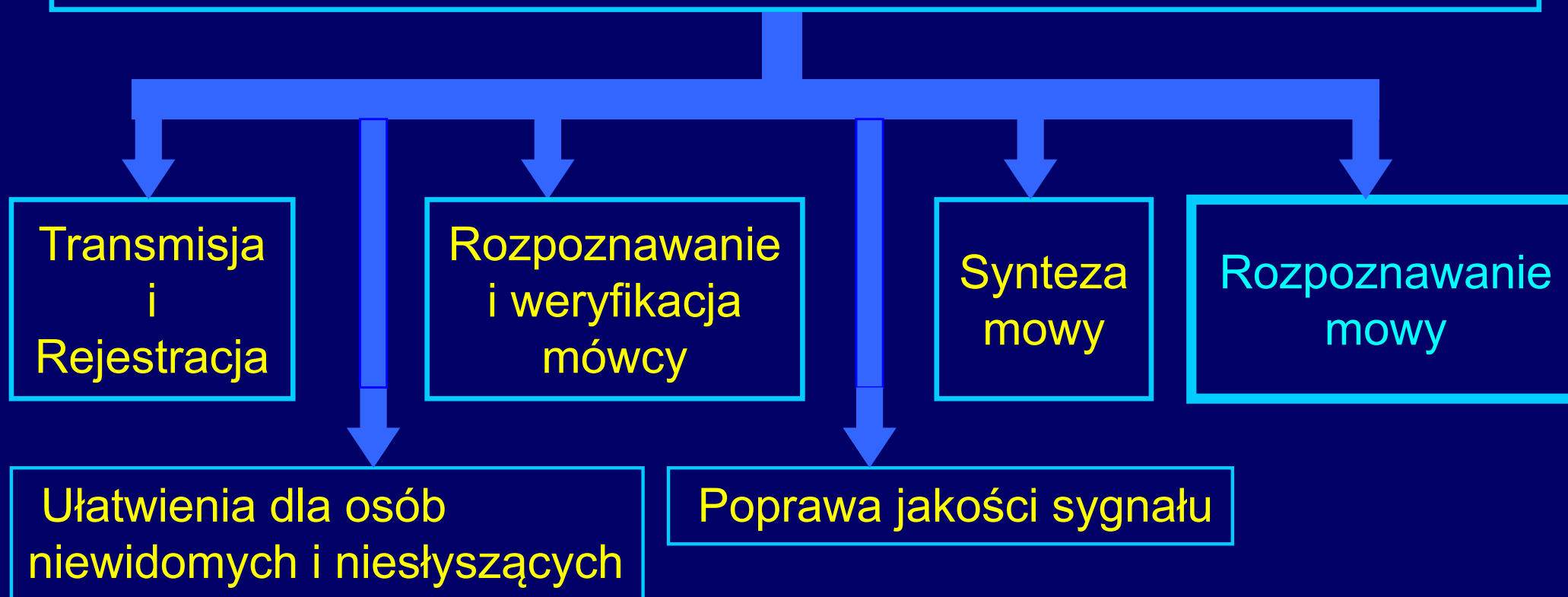
Komputerowe przetwarzanie mowy

Plan wykładu

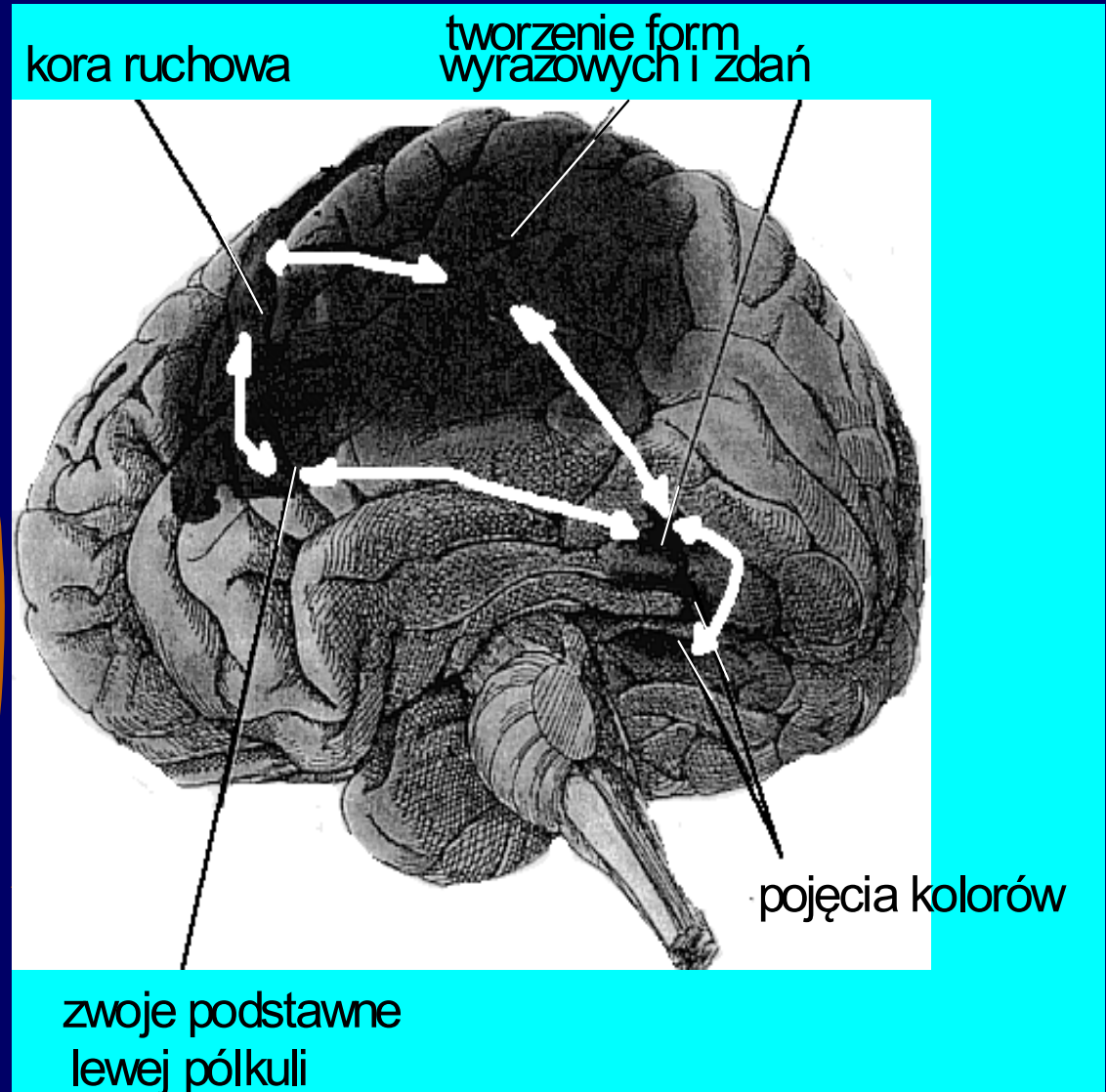
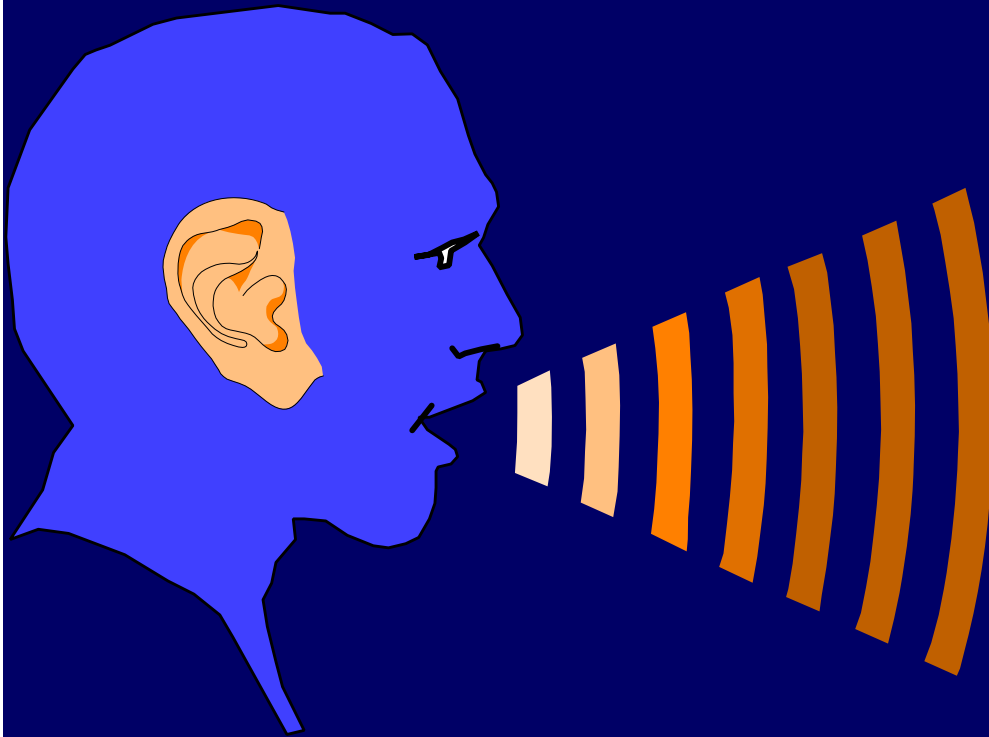
- 1. Wprowadzenie – zagadnienia podstawowe**
- 2. Podział systemów rozpoznawania mowy**
- 3. Charakterystyka metod rozpoznawania mowy**
- 4. Model fizyczny traktu głosowego**
- 5. Ekstrakcja parametrów sygnału mowy**
- 6. Przykładowe algorytmy rozpoznawania mowy**
- 7. Podsumowanie**

Komputerowe przetwarzanie mowy

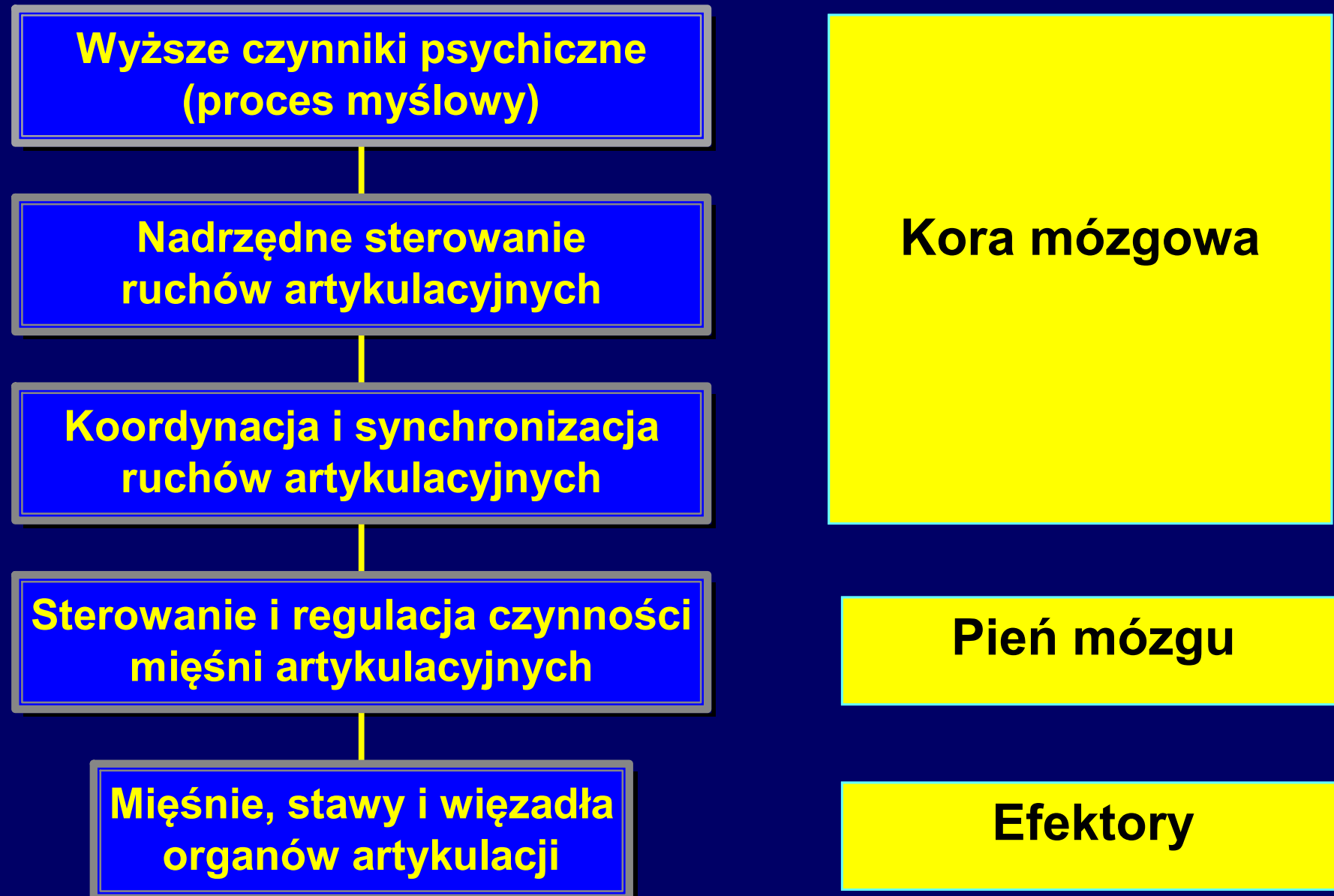
Cyfrowe techniki przetwarzania sygnału mowy



Wprowadzenie – komunikacja werbalna



Struktura systemu wytwarzania mowy



Wytwarzanie mowy

APARAT ARTYKULACYJNY

Składa się z narządów, które modyfikują strumień powietrza. Na styku jamy gardłowej, ustnej i nosowej powstają głoski ustne i nosowe.

Położenie języka w jamie ustnej decyduje o wytwarzaniu głosek twardych i miękkich.

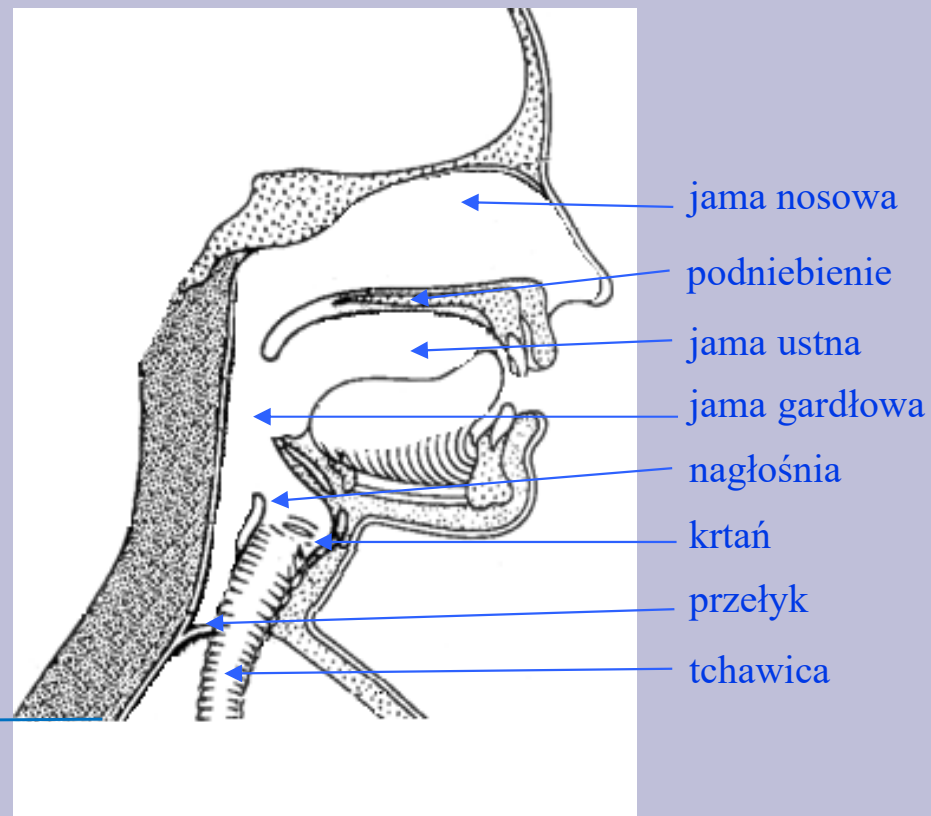
W artykulacji biorą też udział wargi, zęby, dziąsła, podniebienie twarde.

APARAT FONACYJNY

Przy udziale krtani powstają głoski dźwięczne i bezdźwięczne, a położenie więzadeł głosowych decyduje o ich dźwięczności

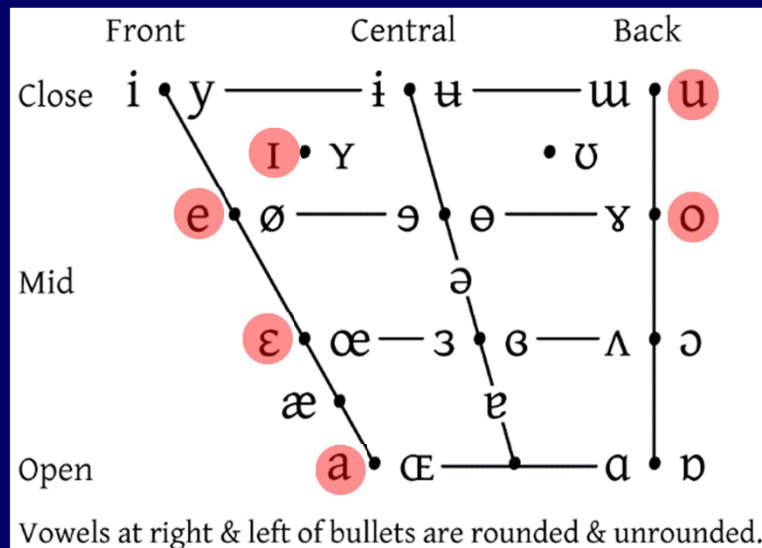
APARAT ODDECHOWY

Dostarczają energię, generującą falę dźwiękową



Międzynarodowy alfabet fonetyczny (angl. International Phonetic Alphabet)

Samogłoski



Jest to standardowy zapis fonetyczny dla wszystkich języków świata.

Alfabet IPA przyjęty w 1886 roku przez Międzynarodowe Towarzystwo Fonetyczne.

Najnowsza wersja alfabetu opublikowana w roku 2005 roku.

Międzynarodowy alfabet fonetyczny

Spółgłoski

	Dwuwargowe (Bilabial)	Wargowo- zębowe (Labiodental)	Zębowe (Dental)	Dziąsłowe (Alveolar)	Zadziąsłowe (Postalveolar)	Retrofleksyjne (Retroflex)	Podniebienne (Palatal)	Miękko- podniebienne (Velar)	Języczkowe (Uvular)	Gardłowe (Pharyngeal)	Krtaniowe (Glottal)
Zwarty-wybuchowe (Plosive)						t d	c ɟ		q ɢ		ʔ
Nosowe (Nasal)	m	ɱ		n		ɳ	ɲ	ŋ	ɴ		
Drżące (Trill)	ʙ			ɾ					ʀ		
Uderzeniowe (Tap or Flap)				ɾ		ɽ					
Szczelinowe (Fricative)	ɸ β	v	θ ð			ʂ ʐ	ç ʝ		χ ʁ	ħ ʕ	h ɦ
Boczne szczelinowe (Lateral fricative)				ɬ ɮ							
Aproksymanty (Approximant)		ʋ		ɹ		ɻ	j	ɰ			
Boczne aproksymanty (Lateral approximant)				l		ɭ	ʎ	ʟ			

Podział według miejsca artykulacji

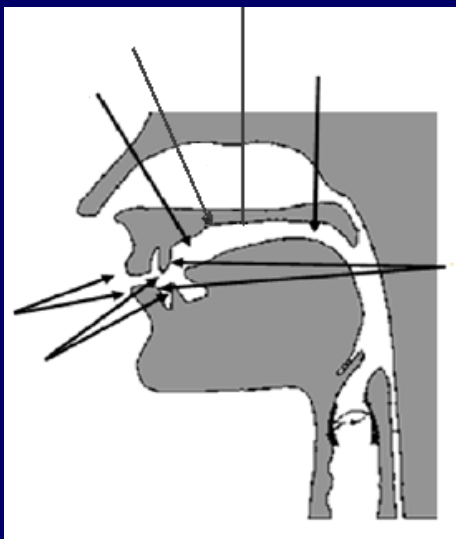
Podniebienne

Zadziąsłowe

Dziąsłowe

Dwu
wargowe

Zębowe



Miękko-
podniebienne

Wargowo
-zębowe

Spółgłoski	IPA	Przykłady
Dwuwargowe	p, p' b, b' m, m'	póvas, peteliškė brolis, labiáu ãmatas, smėgenys
Wargowo-zębowe	f, f'	fãbricas, figūrã
Zębowe	t, t' d, d'	tãkas, šaltėkšnis dãrbas, liũdesys
Dziąsłowe	s, s' z, z, n, n' l, l'	sãulė, vaĩsius zylė, zirzėti nãmas, nėšti vãlsas, valiã
Zadziąsłowe	ʃ, ʃ' ʒ, ʒ, r, r'	šaka, šiãudas žvãkė, žiógas rãtas, kriãušė
Podniebienne	j	ãidas
Miękko-podniebienne	k, k' g, g' x, x' ɣ, ɣ'	kãtinas, kiaũlė gañdras, gėrvė chòras, chėmija harmònija, hiacintas

Podział według sposobu artykulacji

- **Spółgłoski zwarto-wybuchowe**
Zwarcie w jamie ustnej zakańcza się wybuchem
- **Spółgłoski nosowe:**
W jamie ustnej powstaje zwarcie, natomiast w jamie nosowej następuje przepływ powietrza.
- **Spółgłoski drżące:**
między językiem a dziąslami powstaje zwarcie, przez które w przechodzi powietrze
- **Spółgłoski boczne aproksymanty:**
język zwiera się z zębami. Powietrze przechodzi przez boczną powierzchnią języka a zębami.
- **Spółgłoski szczelinowe**
Powstaje nieduża szczelina, przez którą dostarczane jest powietrze.

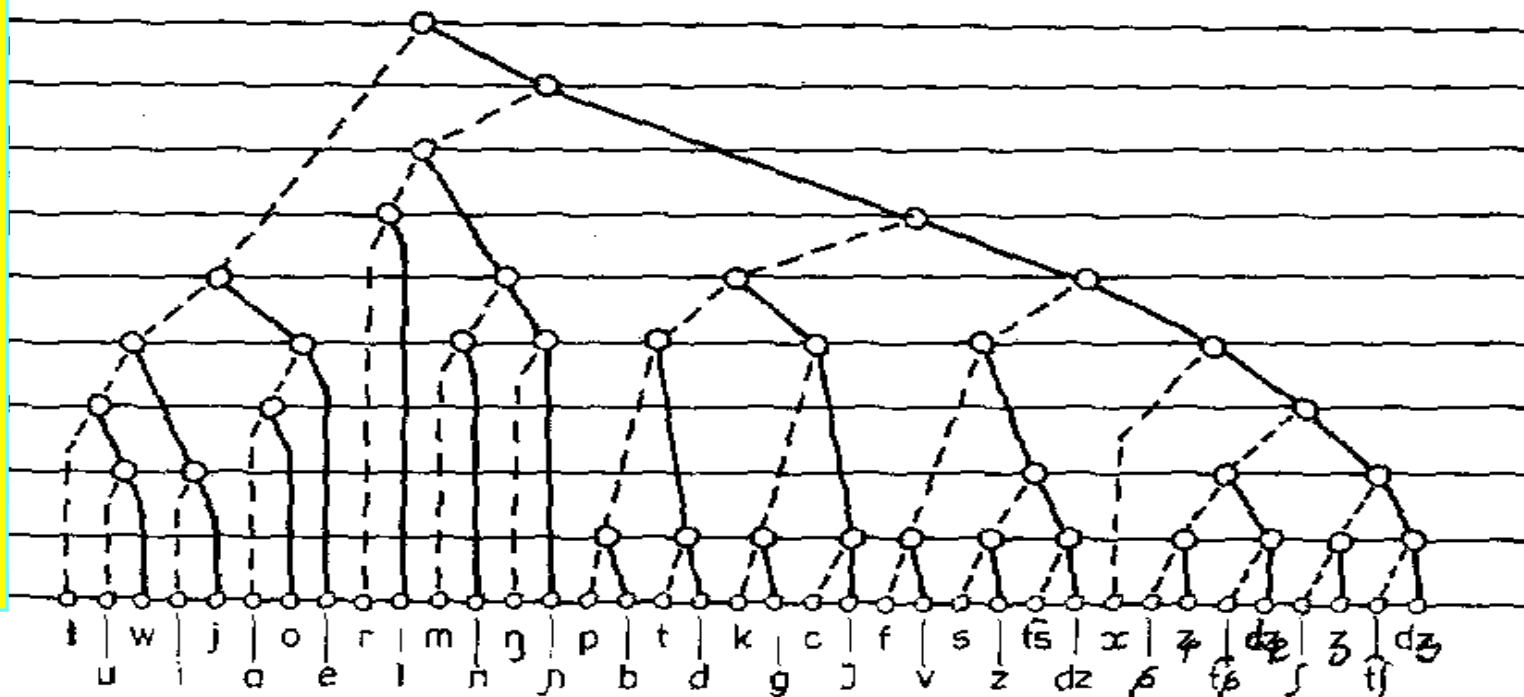
Spółgłoski	IPA	Przykłady
Zwarto-wybuchowe	p, p' b, b' t, t' d, d' k, k' g, g'	póvas, peteliškė brolis, labiáu tākas, šaltėkšnis dārbas, liūdesys kātinas, kiaulė gañdras, gėrvė
nosowe	m, m' n, n'	matas, smėgenys nāmas, nėšti
drżące	r, r'	rātas, kriāušė
boczne aproksymanty	l, l'	vālsas, valià
Szczelinowe	f, f', s, s' z, z' ʃ, ʃ' ʒ, ʒ' x, x' ɣ, ɣ'	fābrikas, figūrā sāulė, vaĩsius zylė, zirzėti šakà, šiaudas žvākė, žiogas chòras, chėmija harmònija, hiacintas

Cechy widma mowy - przykład

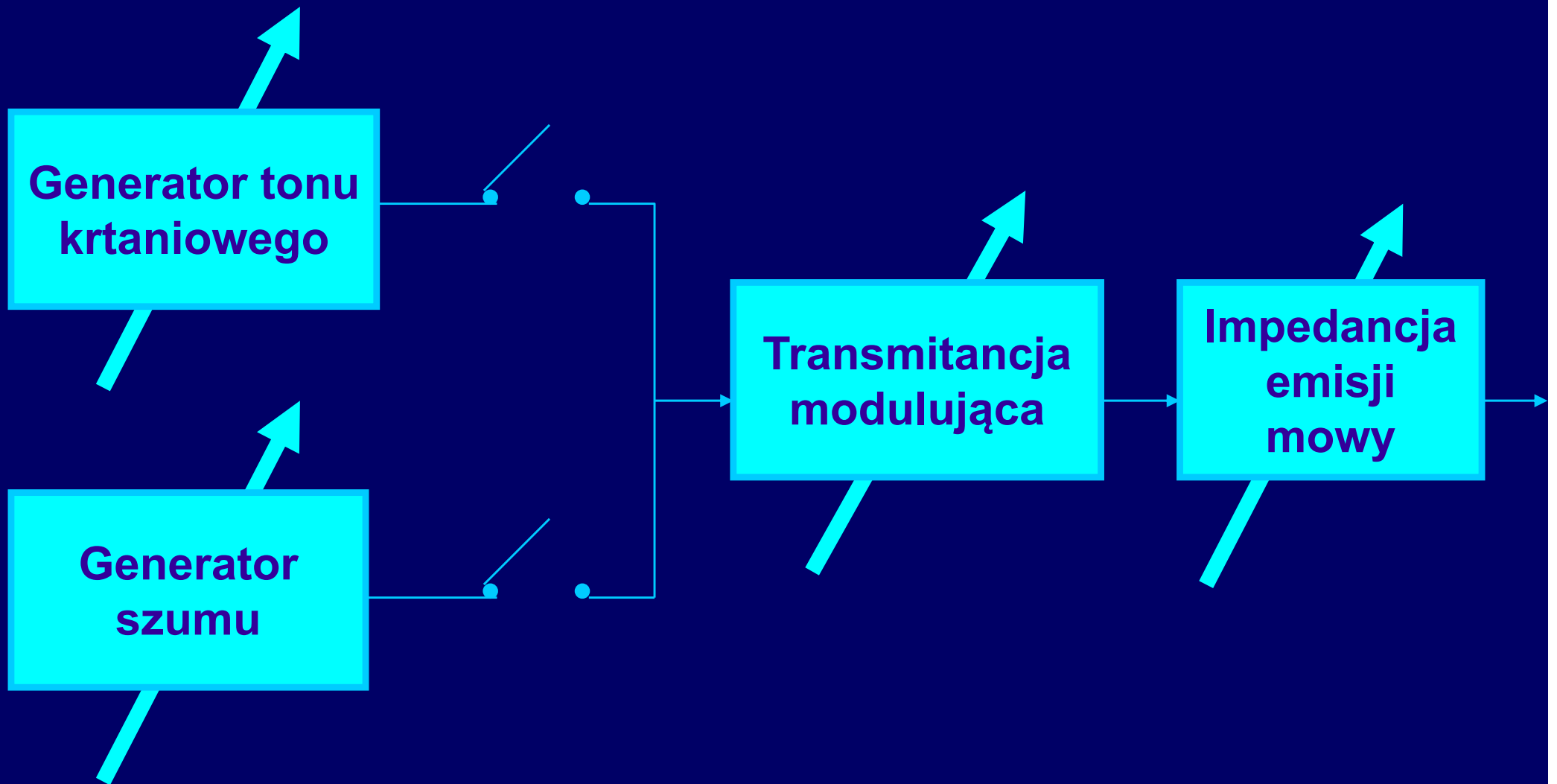
- **nosowy-ustny** - jeżeli widmo mowy wykazuje więcej niż dwa formanty poniżej 2 kHz, to jest to fonem nosowy. W przeciwnym przypadku fonem jest ustny
- **dźwięczny-bezdźwięczny** – fonemy dźwięczne charakteryzuje obecność składowej periodycznej, której z kolei brak w fonemach bezdźwięcznych

Najprostszy system rozpoznawania fonemów

Spółgłoskowe
Ponadkrtaniowe
Nosowe
Łagodne
Skupione
Jasne
Niskotonowe
Krótkie
Dźwięczne



Wytwarzanie mowy

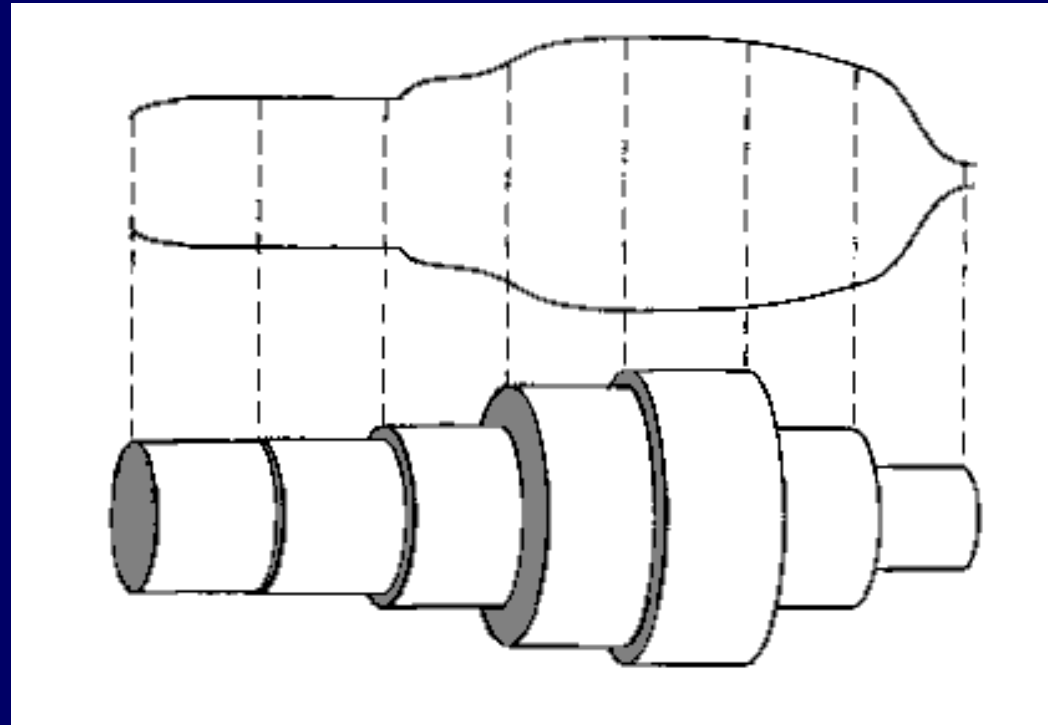


Schemat zastępczy systemu artykulacyjnego

Modelowanie fizyczne - model falowodowy

System złożony z N cylindrów o długości L_i i powierzchni A_i
($i = 1, 2, \dots, N$)

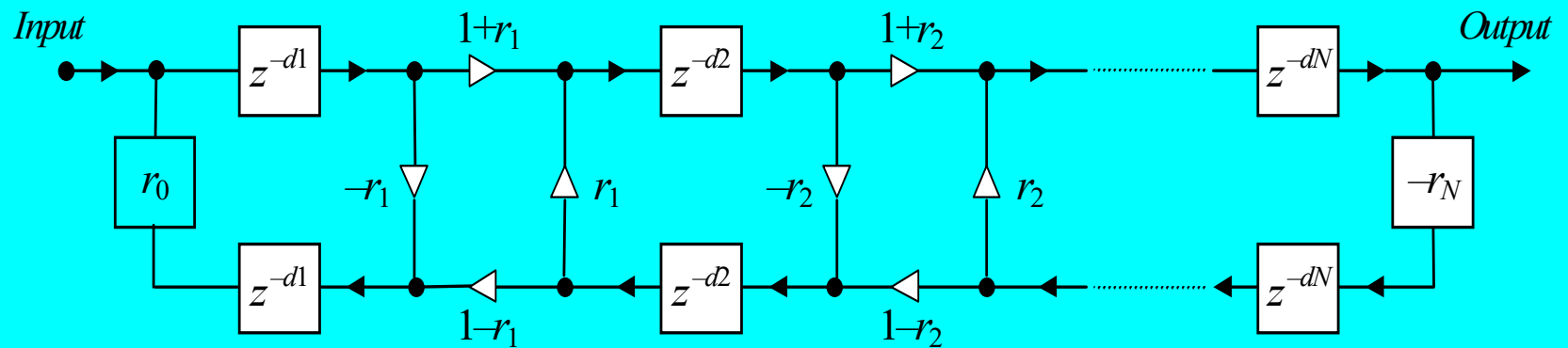
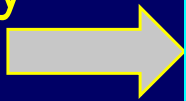
Model fizyczny



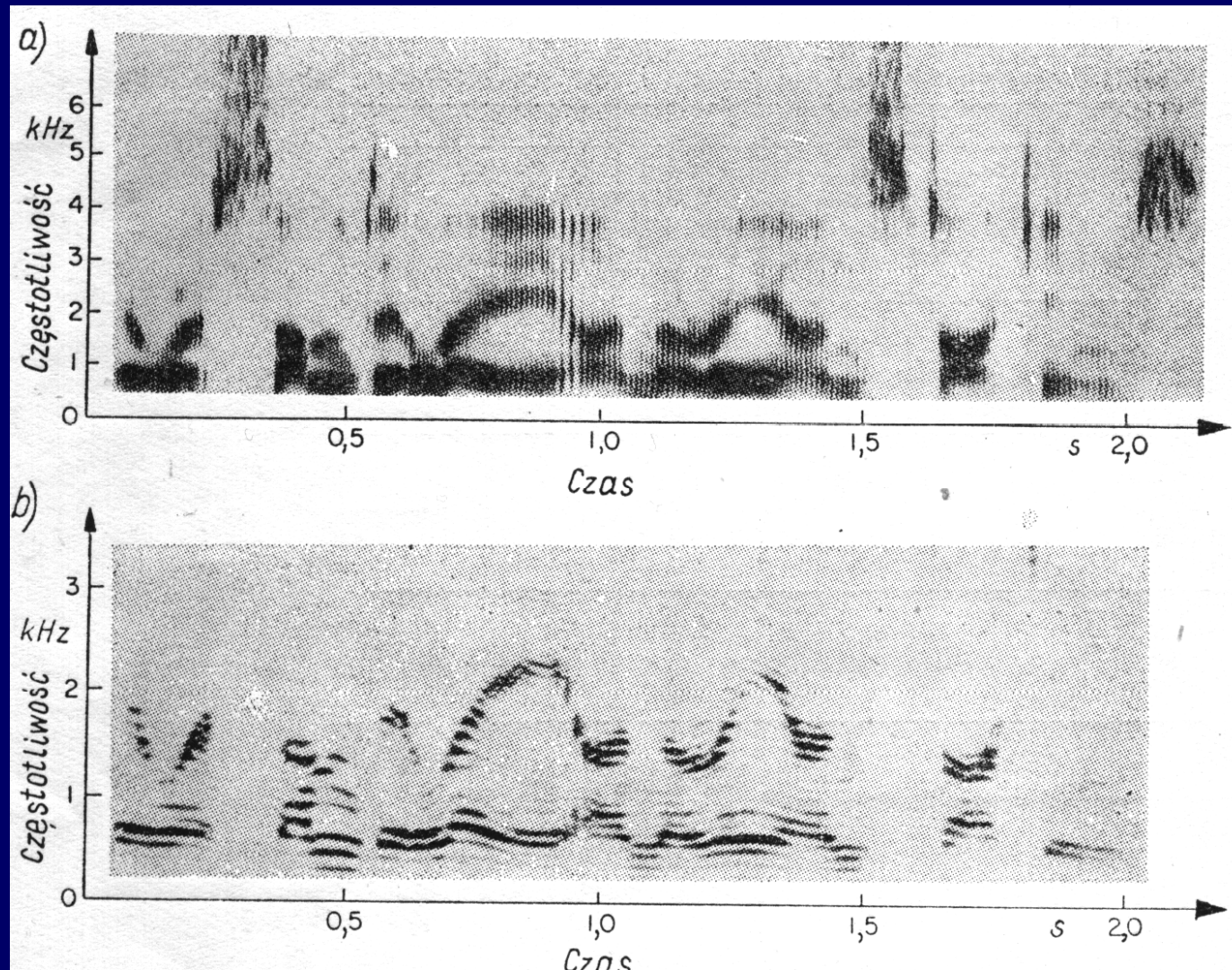
Zbiór rezonatorów
cylindrycznych



„Cyfrowy”
model
falowodowy



Cechy widma mowy



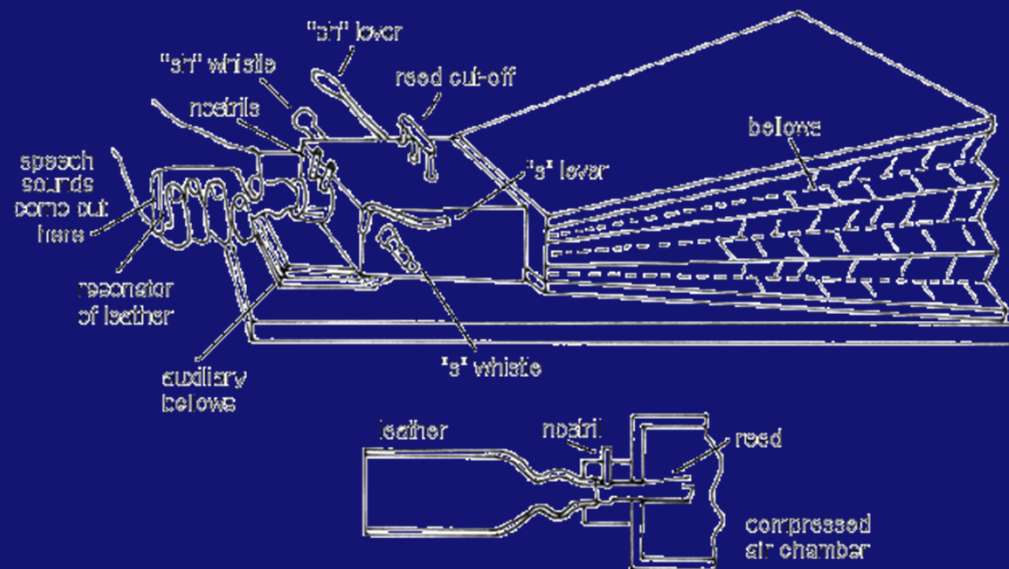
Przykład analizy sonograficznej

Historia syntezy mowy

1773 r. pierwsze badania nad syntezą mowy (profesor Ch.G. Kratzenstein, Kopenhaga)

1846 r. Joseph Faber zaprezentował urządzenie nazwane jako "Euphonia", które generowało nie tylko mowę ludzką, ale także śpiew.

1939 r. pierwszy elektryczny syntezytor mowy wykonany przez Homera Dudley'a ("VODER,,)



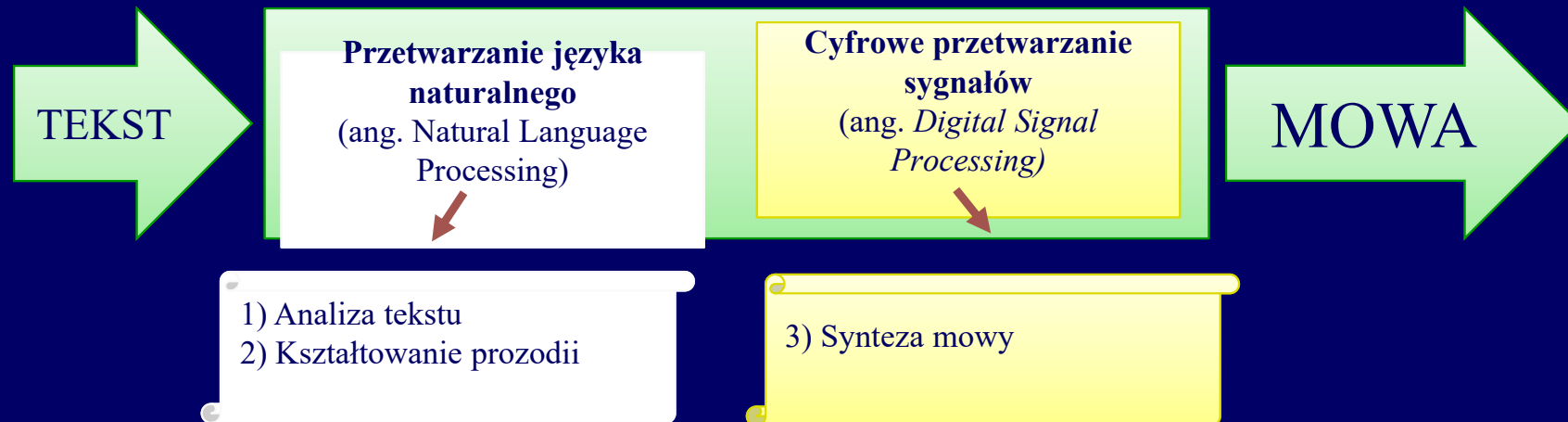
Pierwszy mechaniczny syntezytor (von Kempelen, 1791)

Synteza mowy (ang. Text-To-Speech)

Zmiana tekstu na sygnał akustyczny

Podstawowe cele:

- Zrozumiałość wypowiedzi
- Naturalny dźwięk



Przetwarzanie języka naturalnego

Analiza tekstu



Cel analizy:

Przekształcenie tekstu na zapis fonetyczny

Analiza morfologiczna tekstu

Przydzielenie formy podstawowej i wartości cech gramatycznych dla każdego ze słów.

Przykłady:

szafy	<i>szafa, l. poj., dopełniacz l. mnoga, mianownik</i>
domem	<i>dom, l. poj., narzędnik</i>
mówiła	<i>mówić, czas przeszły, 3osoba l. poj., rodzaj żeński</i>

Analiza kontekstowa

Zadaniem analizatora kontekstowego jest ograniczenie znaczenia poszczególnych słów. W tym celu badane są części mowy słów znajdujących się w sąsiedztwie.

Analiza kontekstowa obejmuje

- Analizę syntaktyczną (rozpoznanie fraz i ich powiązań składniowych)
- Analizę semantyczną (rozpoznanie obiektów, relacji między nimi)
- Analizę pragmatyczną (interpretacja wypowiedzi w konkretnym kontekście, związki logiczne)

Na danym etapie analizy stosowane są

- Metody n-gramów
- Modele Markowa
- Sieci neuronowe

Analiza prozodyczna

Analizowane są brzmieniowe właściwości mowy nakładające się na głoskowy, sylabiczny i wyrazowy ciąg wypowiedzi.

Prozodie odzwierciedlają:

- Osobiste cechy mówcy
- Stan emocjonalny mówcy
- Cechy wypowiedzi (ironiczny lub sarkastyczny)
- Nacisk, kontrast i ostrość

Kształtowanie prozodii

Kształtowanie prozodii jest niezbędnym procesem dla każdego systemu mowy. Bez zaprogramowania cech emocjonalnych synteza brzmi sztucznie (jak „głos robota”)



- Podwyższenie lub obniżenie tonu
- Zwiększenie lub zmniejszenie intensywności amplitudy
- Wydłużenie lub skrócenie czasu trwania głoski/wyrazu

Jak stany emocjonalne znajdują
swoje odbicie w mowie ?

Ryszard Gubrynowicz

Interpretacja aktorska

happy 

sad 

angry 

interested 

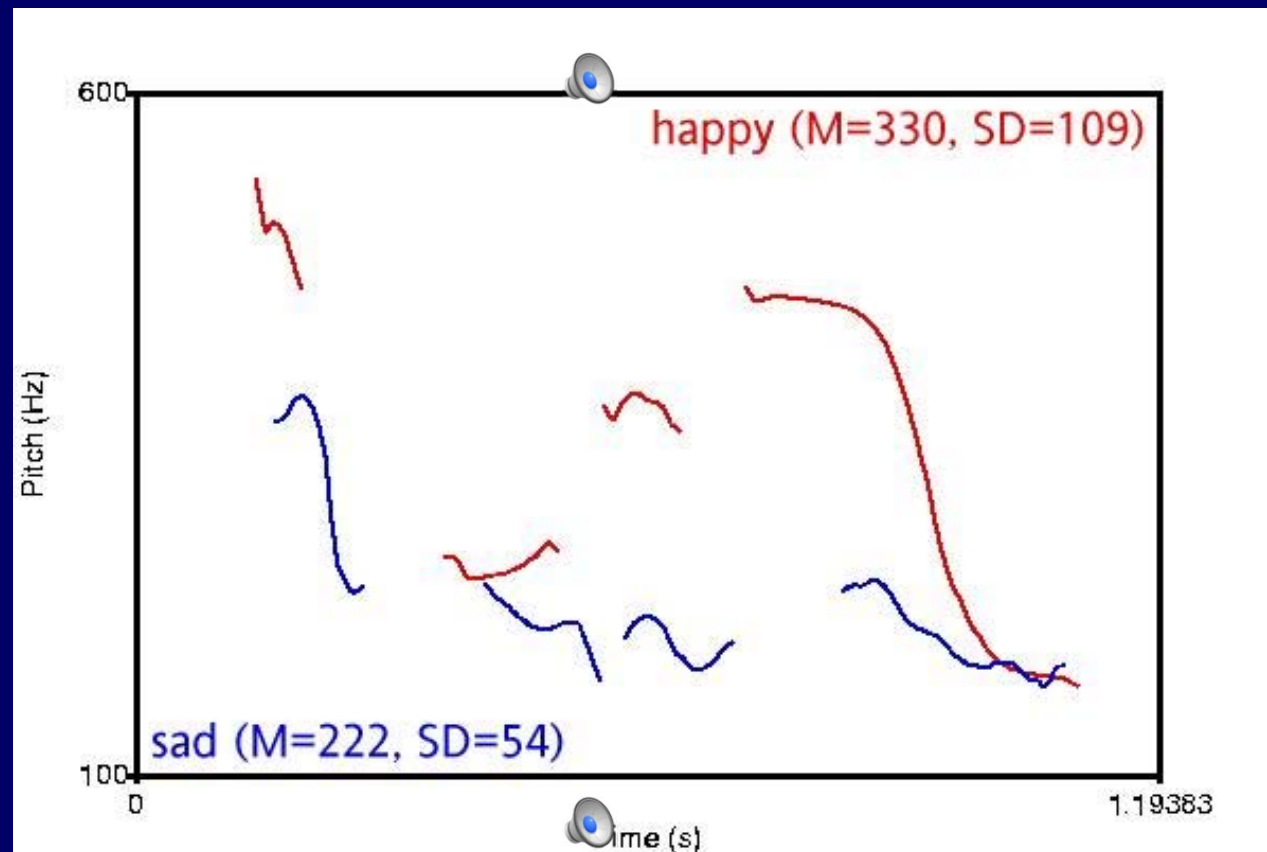
bored 



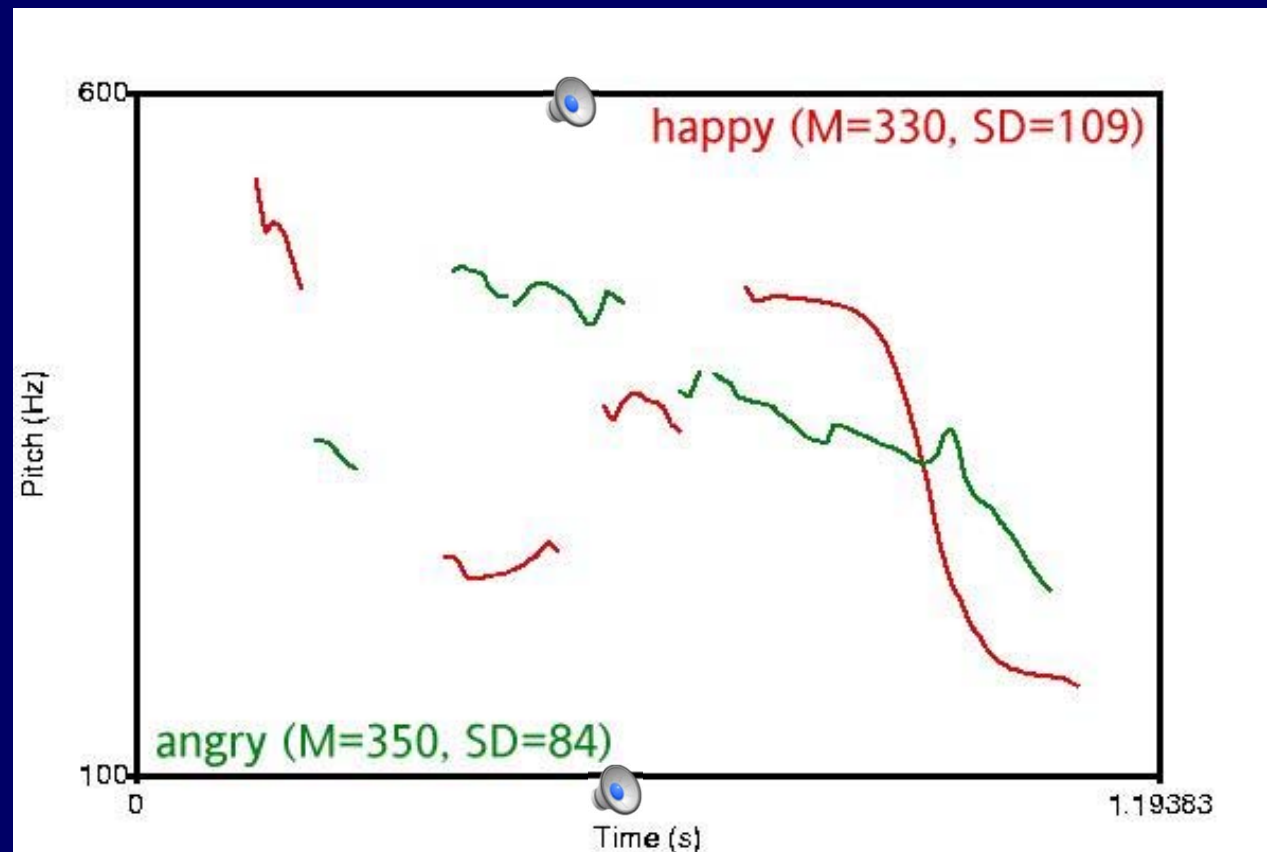
Funkcje emocjonalne cech prozodycznych

Słuchacz na ogół kontroluje w wypowiedzi swój stan emocjonalny. W jego wyrażeniu posługuje się przede wszystkim tempem mówienia, głośnością, wprowadzaniem dodatkowych pauz, przedłużaniem niektórych dźwięków, a także modulowaniem melodii. W wypowiedziach nacechowanych emocjonalnie wahania melodii są znacznie większe, niż w wypowiedziach o charakterze neutralnym. Neutralne – 3-4 tony, z dużym ładunkiem emocjonalnym - > 1 oktawy.

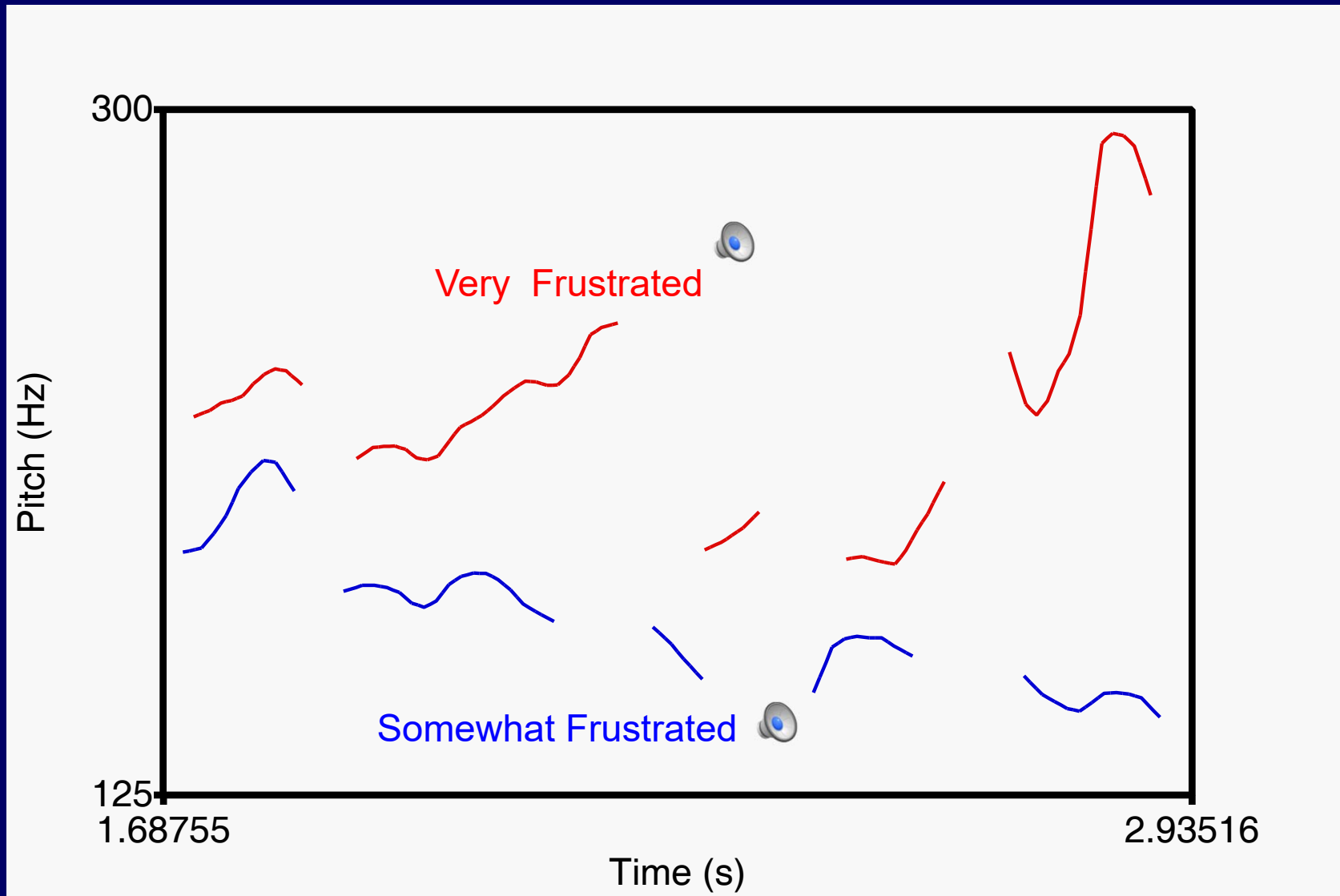
Przykład opozycji szczęśliwy – smutny w konturze melodycznym



Przykład opozycji szczęśliwy – gniewny w konturze melodycznym



Przykład z dialogu typu HMIHY – How may I help you ?

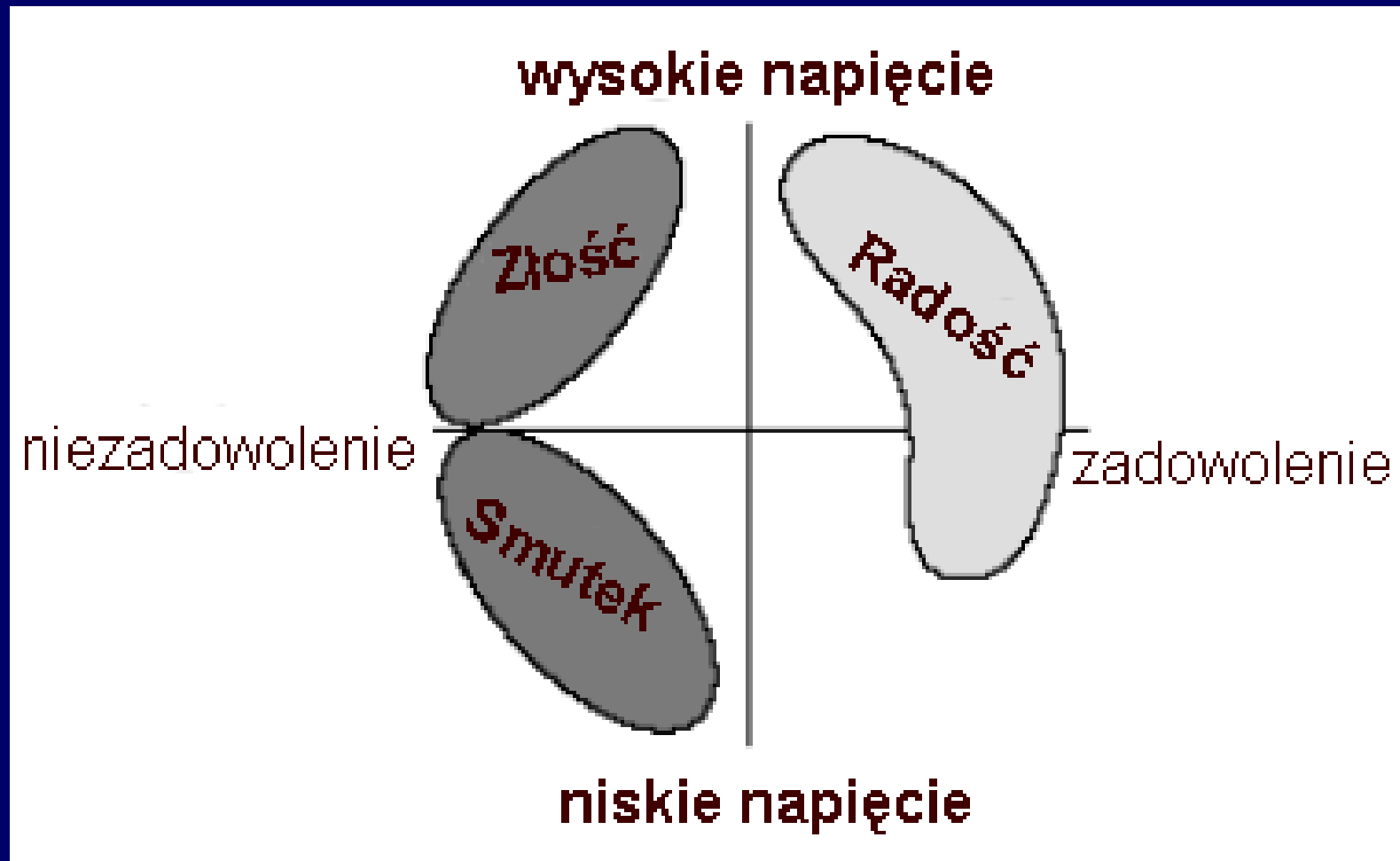


Trudności w określaniu emocji

Nadanie wypowiedzi określonego typu emocji jest zadaniem bardzo złożonym. Osoby określające typ wypowiedzi pod względem emocji rzadko są zgodne w swych ocenach, z wyjątkiem krańcowych, lub wyraźnie kontrastowych typów emocji

Słuchacze w swojej ocenie głównie opierają się na cechach prozodycznych, zwłaszcza na iloczynach i stylizowanym przebiegu F0.

Emocje w płaszczyźnie subiektywnej



Emocje kontrastowe w płaszczyźnie akustycznej

Strach/złość

- zwiększona prędkość i głośność wypowiedzi
- podwyższone F0
- zwiększony zakres F0
- zaburzony rytm mowy
- dokładniejsza artykulacja
- zwiększona energia w zakresie wyższych częstotliwości

Smutek/odprężenie

- zmniejszona prędkość i głośność wypowiedzi
- obniżone F0
- zmniejszony zakres F0
- wyrównany rytm mowy, płynna mowa
- niedokładna artykulacja
- obniżona energia w zakresie wyższych częstotliwości

Article

Recognition of Emotions in Speech Using Convolutional Neural Networks on Different Datasets

Marta Zielonka [†], Artur Piastowski [†], Andrzej Czyżewski ^{*†}, Paweł Nadachowski [†], Maksymilian Operlejn [†] and Kamil Kaczor

Faculty of Electronics, Telecommunications and Informatics, Gdańsk University of Technology, Gabriela Narutowicza 11/12, 80-233 Gdańsk, Poland

^{*} Correspondence: ac@pg.edu.pl[†] These authors contributed equally to this work.

Citation: Zielonka, M.; Piastowski, A.; Czyżewski, A.; Nadachowski, P.; Operlejn, M.; Kaczor, K. Recognition of Emotions in Speech Using Convolutional Neural Networks on Different Datasets. *Electronics* 2022, 11, 3831. <https://doi.org/10.3390/electronics11223831>

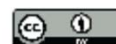
Academic Editors: Daniel Hładek, Matěj Pleva, Piote Szczerko and Andrzej Zgank

Received: 9 August 2022

Accepted: 15 November 2022

Published: 21 November 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Abstract: Artificial Neural Network (ANN) models, specifically Convolutional Neural Networks (CNN), were applied to extract emotions based on spectrograms and mel-spectrograms. This study uses spectrograms and mel-spectrograms to investigate which feature extraction method better represents emotions and how big the differences in efficiency are in this context. The conducted studies demonstrated that mel-spectrograms are a better-suited data type for training CNN-based speech emotion recognition (SER). The research experiments employed five popular datasets: Crowd-sourced Emotional Multimodal Actors Dataset (CREMA-D), Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS), Surrey Audio-Visual Expressed Emotion (SAVEE), Toronto Emotional Speech Set (TESS), and The Interactive Emotional Dyadic Motion Capture (IEMOCAP). Six different classes of emotions were used: happiness, anger, sadness, fear, disgust, and neutral. However, some experiments were prepared to recognize just four emotions due to the characteristics of the IEMOCAP dataset. A comparison of classification efficiency on different datasets and an attempt to develop a universal model trained using all datasets were also performed. This approach brought an accuracy of 55.89% when recognizing four emotions. The most accurate model for six emotion recognition was trained and achieved 57.42% accuracy on a combination of four datasets (CREMA-D, RAVDESS, SAVEE, TESS). What is more, another study was developed that demonstrated that improper data division for training and test sets significantly influences the test accuracy of CNNs. Therefore, the problem of inappropriate data division between the training and test sets, which affected the results of studies known from the literature, was addressed extensively. The performed experiments employed the popular ResNet18 architecture to demonstrate the reliability of the research results and to show that these problems are not unique to the custom CNN architecture proposed in experiments. Subsequently, the label correctness of the CREMA-D dataset was studied through the employment of a prepared questionnaire.

Keywords: speech emotion recognition; SER; machine learning; artificial intelligence; classification; convolutional neural networks

1. Introduction

The recognition of emotions is a relatively difficult and complex task [1], even for humans. Many people could say that they can perform this task efficiently; however, they often have the opportunity to recognize emotions based on a few different aspects, such as body language, facial expression, and voice timbre or prosody. Meanwhile, speech emotion recognition (SER) is a potentially significant step toward the future as it presents a huge variety of use cases.

SER considers recognizing emotions using only one modality, voice recordings, which makes it more complex. Thus, it uses one additional medium—a microphone that may also capture some noise [2]. Achieving decent results on this type of problem could lead to the

Rodzaje syntezy mowy

- **Metoda formantowa**

Odwzorowanie widma sygnału mowy

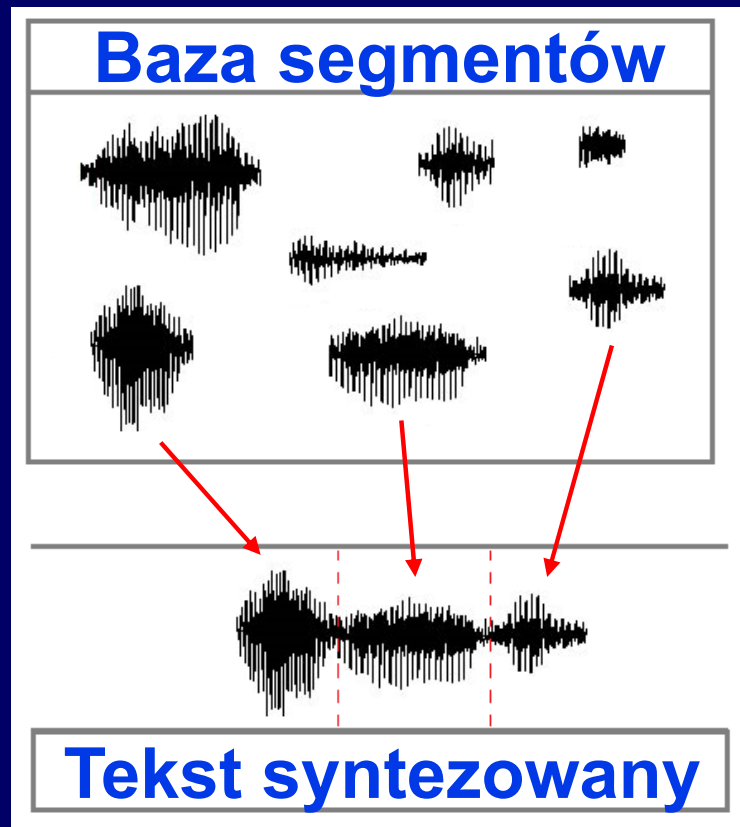
- **Metoda artykulacyjna**

Fizyczne odwzorowanie mechanizmów
wytwarzania mowy

Metoda konkatenacyjna

Wykorzystanie nagranych próbek sygnału mowy

Konkatenacyjna synteza mowy



Łączenie wypowiedzi z mniejszych jednostek nagranych przez lektora

Wykorzystywane jednostki:

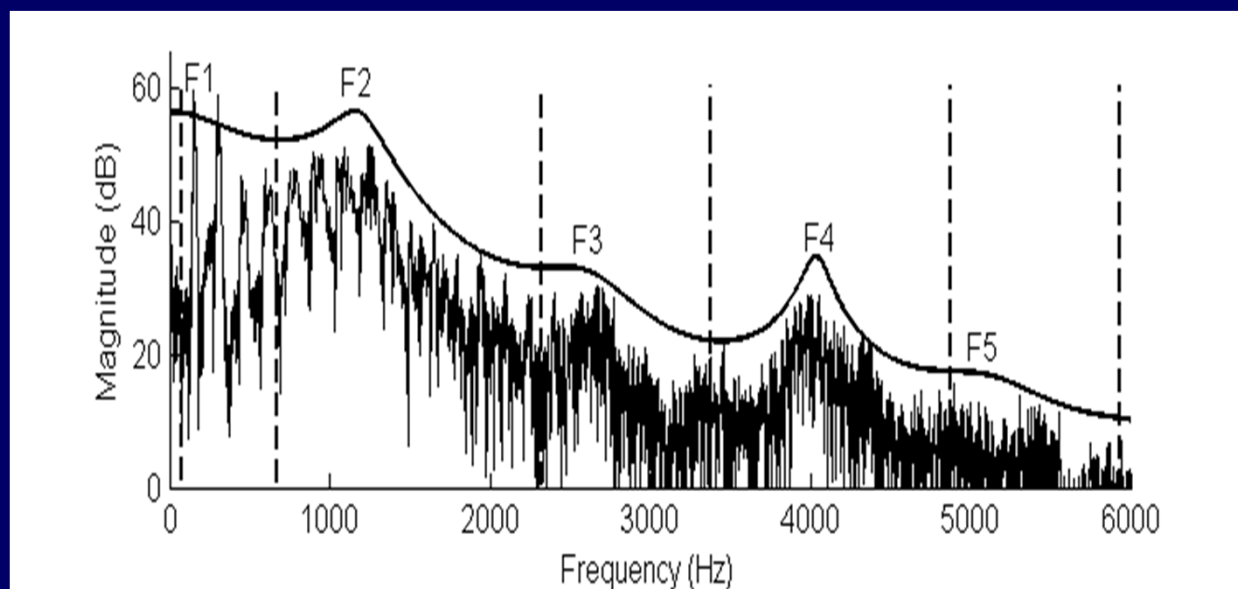
- fonem (głoska)
- difony
- trifony
- sylaby
- całe wyrazy

Jest to najczęściej spotykana metoda syntezy.

Formantowa synteza mowy

Modelowanie traktu głosowego jako połączenie rezonatorów – filtrów elektrycznych lub cyfrowych.

Podejście to ma w założeniu odwzorować formantowy charakter sygnału mowy.



Formant - skupisko energii w widmie sygnału mowy.

Od rozmieszczenia formantów zależy zrozumiałość mowy.

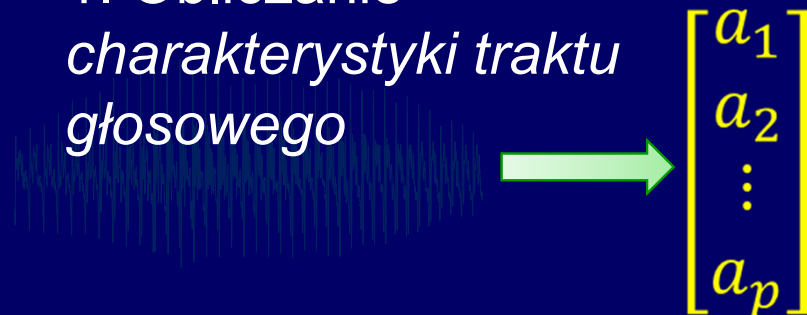
Artykulacyjna synteza mowy

Zakłada się, że głos powstaje w trakcie głosowym (układ filtrów - rezonatorów o zmiennych parametrach) za pomocą sygnału pobudzającego

Sygnał pobudzający - struny głosowe (oddziaływanie strumienia powietrza i fałd głosowych lub szumu białego)

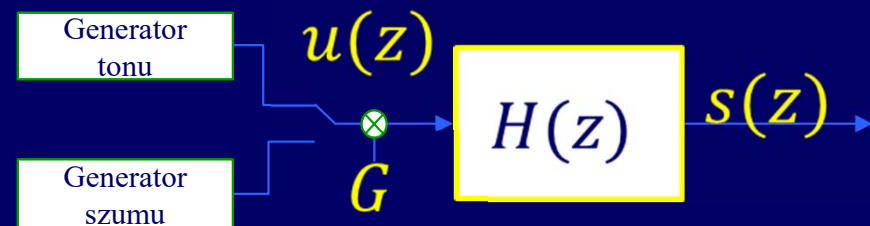
Najczęściej używa się kodowania predykcyjnego (*Linear Predictive Coding*).

1. Obliczanie
*charakterystyki traktu
głosowego*



Analiza LPC

2. Odwzorowanie
*charakterystyki traktu
głosowego* za pomocą modelu
matematycznego.



$$\frac{1}{-\sum_{k=1}^p a_k}$$

Zastosowania syntezy mowy

- Urządzenia dla osób niewidomych
- Mówiące telefony, komputery, planszety
- Słowniki językowe
- Udźwiękowanie stron internetowych, aplikacji, gier edukacyjnych
- Odczyt poczty elektronicznej

Historia rozpoznawania mowy



1937 r. Stevens i Newman zdefiniowali melową skalę częstotliwości

1952 r. Naukowcy z Bell Labs wynaleźli system rozpoznawania cyfr izolowanych.

1965 r. Cooley i Tukey opracowali algorytm szybkiej transformacji Fouriera.

Zabawka Radio Rex powstała w 1920 roku

Rozpoznawanie mowy



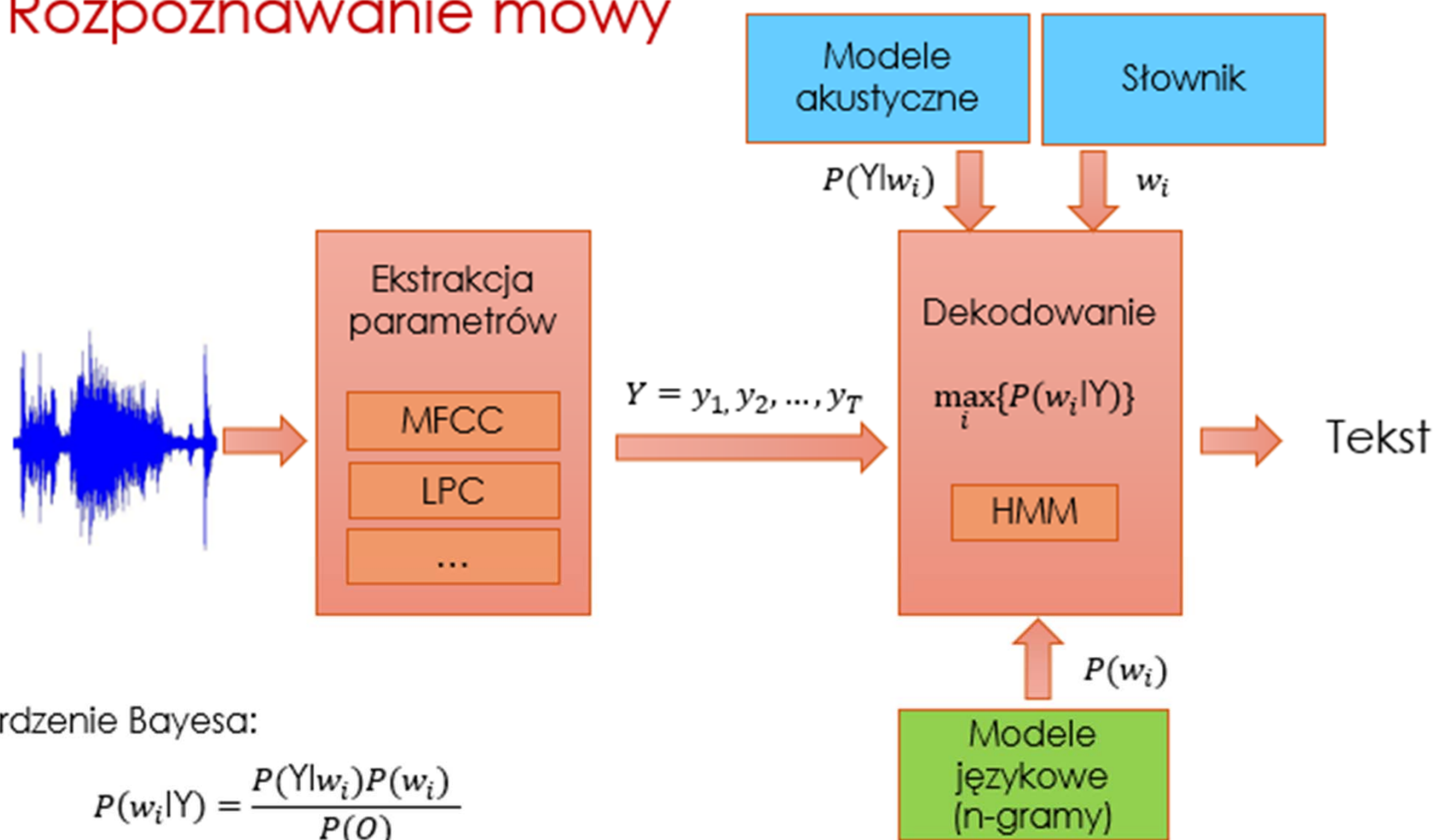
System może być zależny i niezależny od mówcy

Wielkość słownika

Słownik	Ilość wyrazów
Mały	2 – 100 wyrazów
Średni	100 – 1000 wyrazów
Duży	ponad 1000 wyrazów

Obecnie systemy są w stanie rozpoznać 50 tysięcy słów

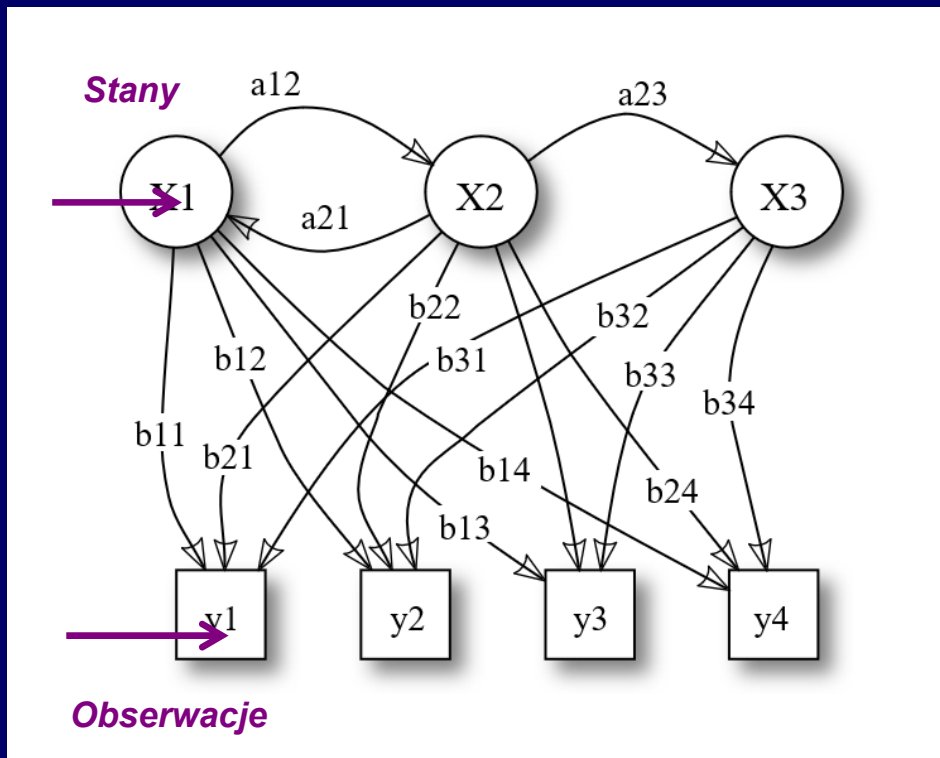
Rozpoznawanie mowy



Twierdzenie Bayesa:

$$P(w_i|Y) = \frac{P(Y|w_i)P(w_i)}{P(O)}$$

Dekodowanie sygnału za pomocą ukrytych modeli Markowa



ang. Hidden Markov Models (HMM)

Obliczenie prawdopodobieństwa $P(Y|w_i)$ sprowadza się do obliczenia sumarycznego prawdopodobieństwa (zdarzeń i przejść).

W ukrytym modelu Markowa stan nie jest widoczny, jednak wyjście zależne od niego jest znane.

Do odkrywania ukrytej sekwencji stanów modelu HMM stosuje się algorytm *Viterbiego*

<https://upload.wikimedia.org/wikipedia/commons/8/8a/HiddenMarkovModel.svg>

Ekstrakcja parametrów - metody cepstralne

ang. *Mel Frequency Cepstral Coefficient (MFCC)*

- Cepstrum - to transformata Fouriera logarytmu widma $\hat{X}(T) = F[\ln(X(f))]$.
- Skala cepstrum odpowiada dziedzinie czasu
- Współczynniki cepstralne niosą informacje o trakcie głosowym i o tonie krtaniowym
- Skala melowa, określająca subiektywny odbiór wysokości dźwięku przez ludzkie ucho względem skali w hercach $F_{mel} = 1127 \log_e 1 + f / 700$

Podział
sygnału
na ramki

Zastosowanie
okna
Hamminga na
każdej z ramek

Transformata
Fouriera na
każdej z ramek

Filtracja danych
bankiem filtrów i
obliczenie
logarytmu energii

Transformata
kosinusowa,
której wynikiem
są współczynniki
cepstralne

Analiza mowy – parametryzacja

- współczynniki cepstralne (MFCC) w skali nieliniowej (melowej)

$$M_i = \sum_{k=1}^{20} X_k \cos[i(k - 0.5)\pi / 20]$$

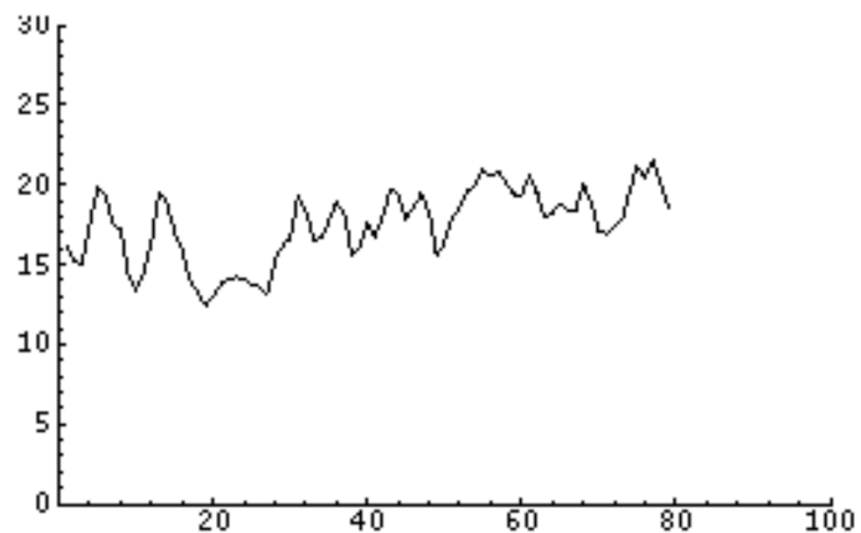
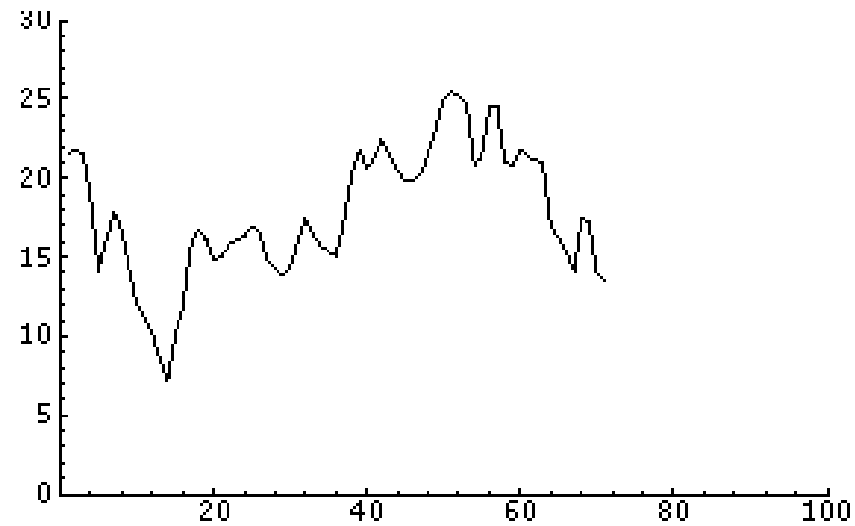
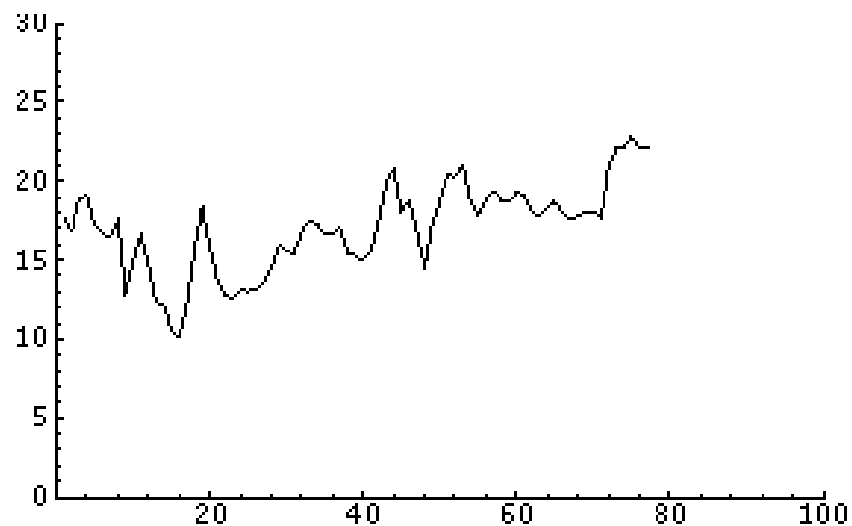
gdzie: i - numer współczynnika cepstralnego;

k - liczba pasm częstotliwości

X_k - logarytm energii w danym paśmie częstotliwości k

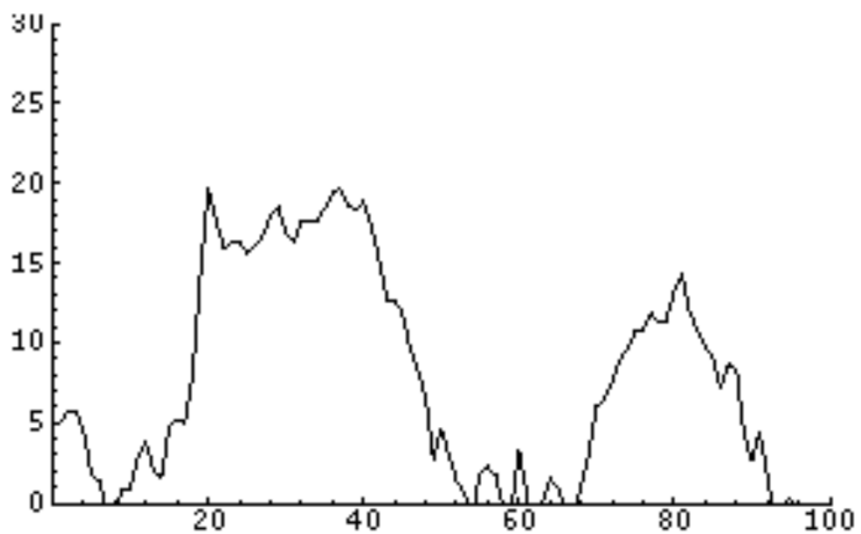
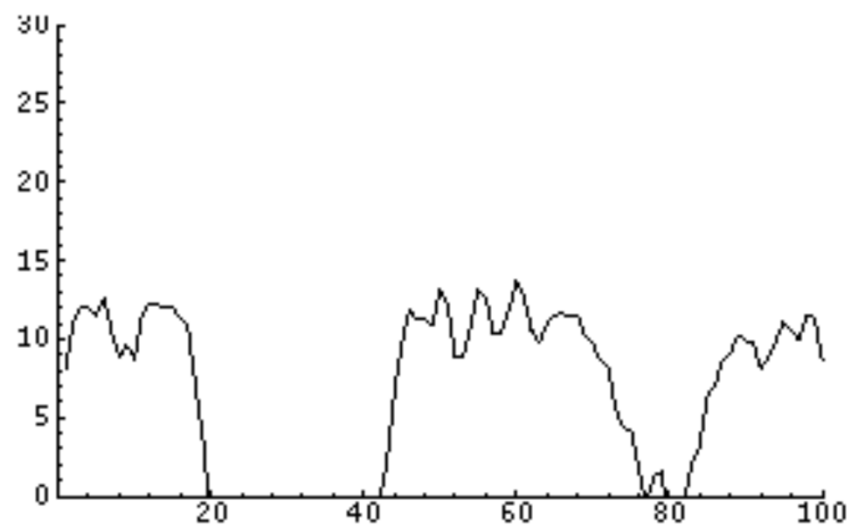
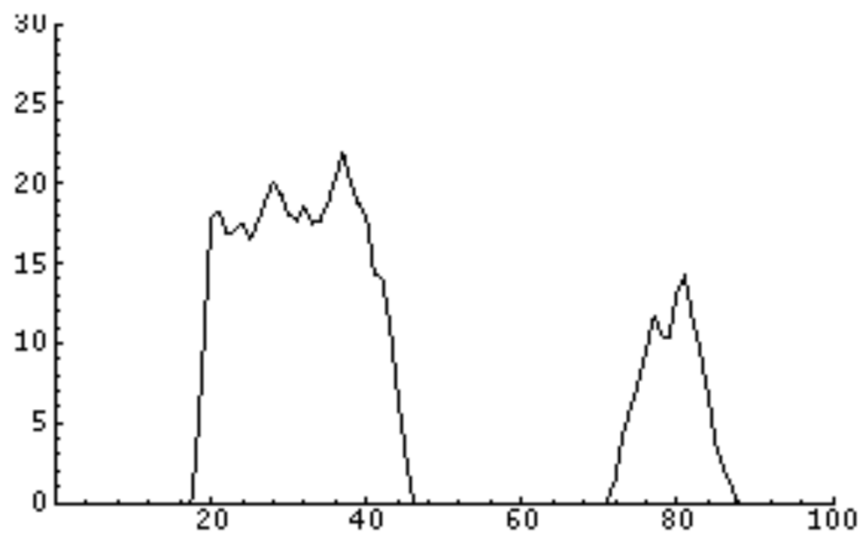
- proste parametry, np. liczba przejść przez zero lub przez inną wartość (w celu ograniczenia wpływu składowej stałej)
- analiza LPC – współczynniki LPC

Współczynniki mel-cepstralne



**Słowo „zero” -
trzech mówców**

Współczynniki mel-cepstralne

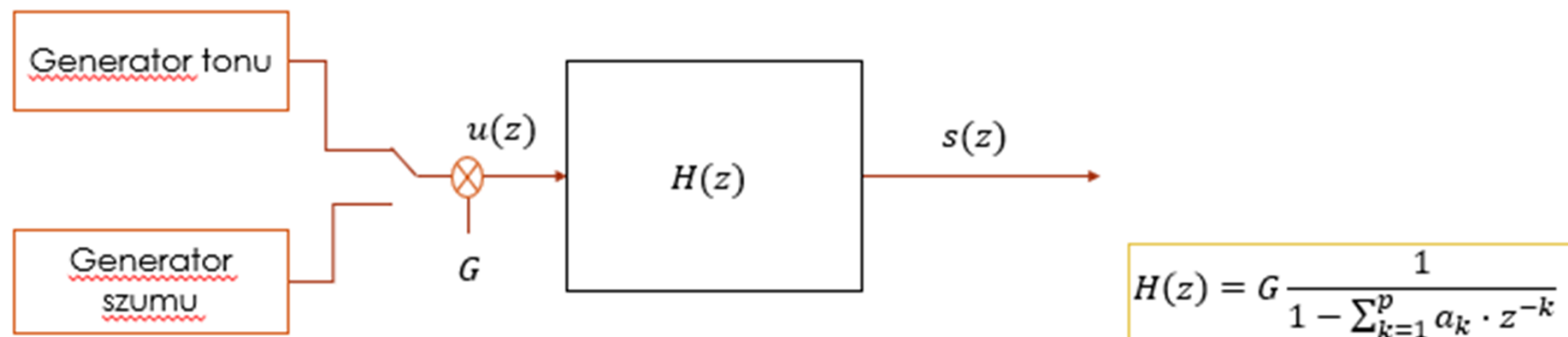


**Słowo „sześć” -
trzech mówców**

Ekstrakcja parametrów - metody predykcyjne

ang. Linear Predictive Coding (LPC)

Odwzorowuje rezonansową strukturę traktu głosowego



Sygnal mowy - odpowiedź filtra na pobudzenie

Filtr - rezonansowa charakterystyka traktu głosowego

Pobudzenie - sygnał tonu krtaniowego

Podział systemów ARM



Proces rozpoznawania sygnału mowy

Analiza i przetwarzanie wstępne sygnału

Ekstrakcja parametrów

Identyfikacja elementów fonetycznych

Analiza leksykalna, gramatyczna, semantyczna

"Rozumienie"



Analiza mowy – przetwarzanie wstępne

- Normalizacja energetyczna
- Segmentacja sygnału (detekcja granic wyrazów)

Przykładowo:

- Segmentacja poprzez analizę obwiedni amplitudowej

$$p_i - p_{i-1} > k \vee p_i - p_{i+1} > k$$

gdzie:

p_i - i -ta próbka sygnału

k - arbitralnie przyjęta wartość progowa

$$c = \frac{\int_{t_1}^{t_2} ts(t)dt}{\int_{t_1}^{t_2} s(t)dt}$$

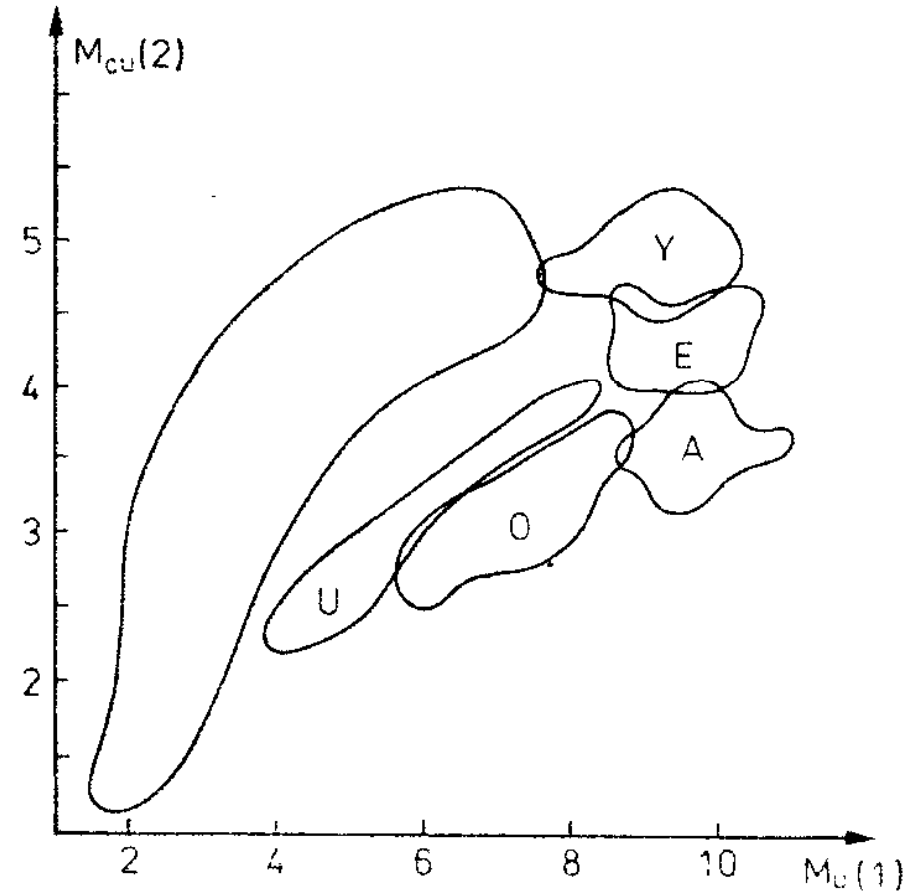
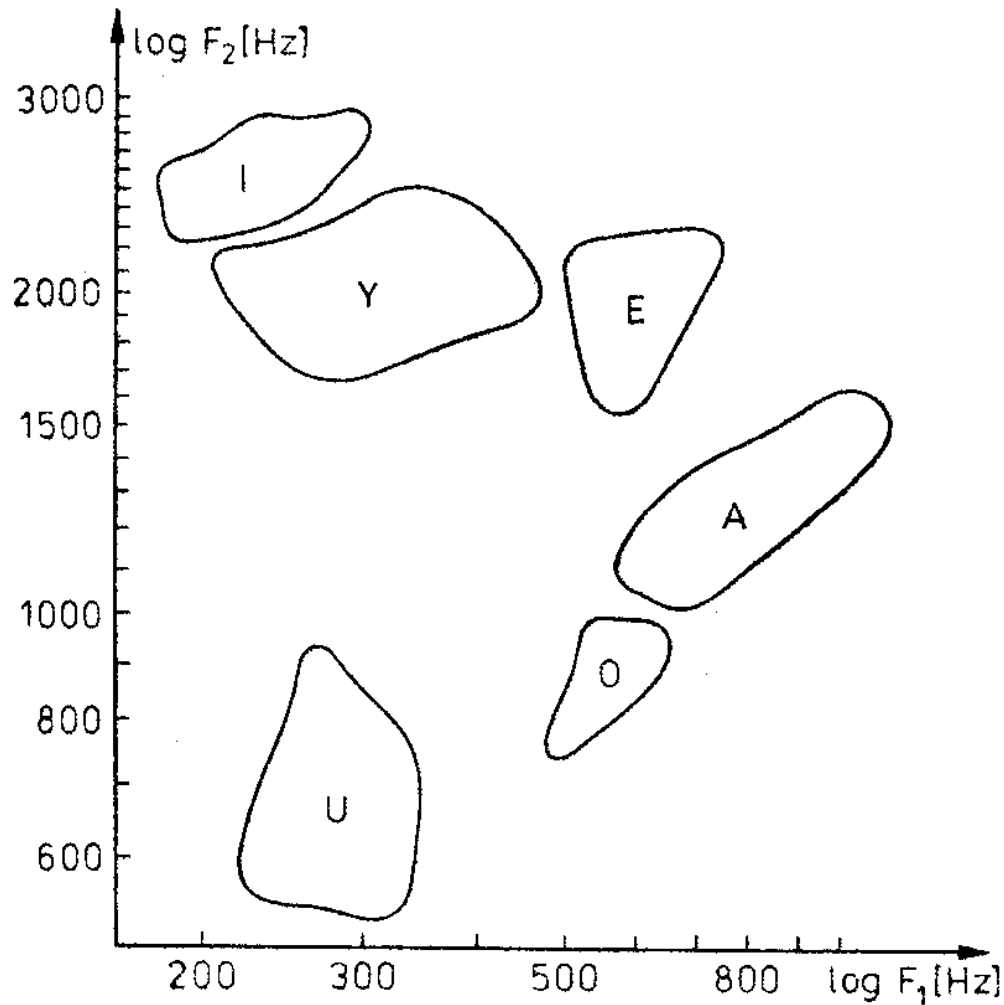
$$d = \sqrt{\frac{\int_{t_1}^{t_2} (t-c)^2 s(t)dt}{\int_{t_1}^{t_2} s(t)dt}}$$

$$k_1 \cong c - wd$$

$$k_2 \cong c + wd$$

gdzie: c - środek ciężkości, d – dyspersja, t_1, t_2 – dowolna próbka „przed” i „za” wyrazem, $s(t)$ – rozkład czasowy funkcji gęstości p , k_1, k_2 – granice wyrazu (numer próbki),

Ekstrakcja parametrów - fonemy samogłoskowe



Formanty F_1 i F_2

Momenty centralne M_{c1} i M_{cu2}

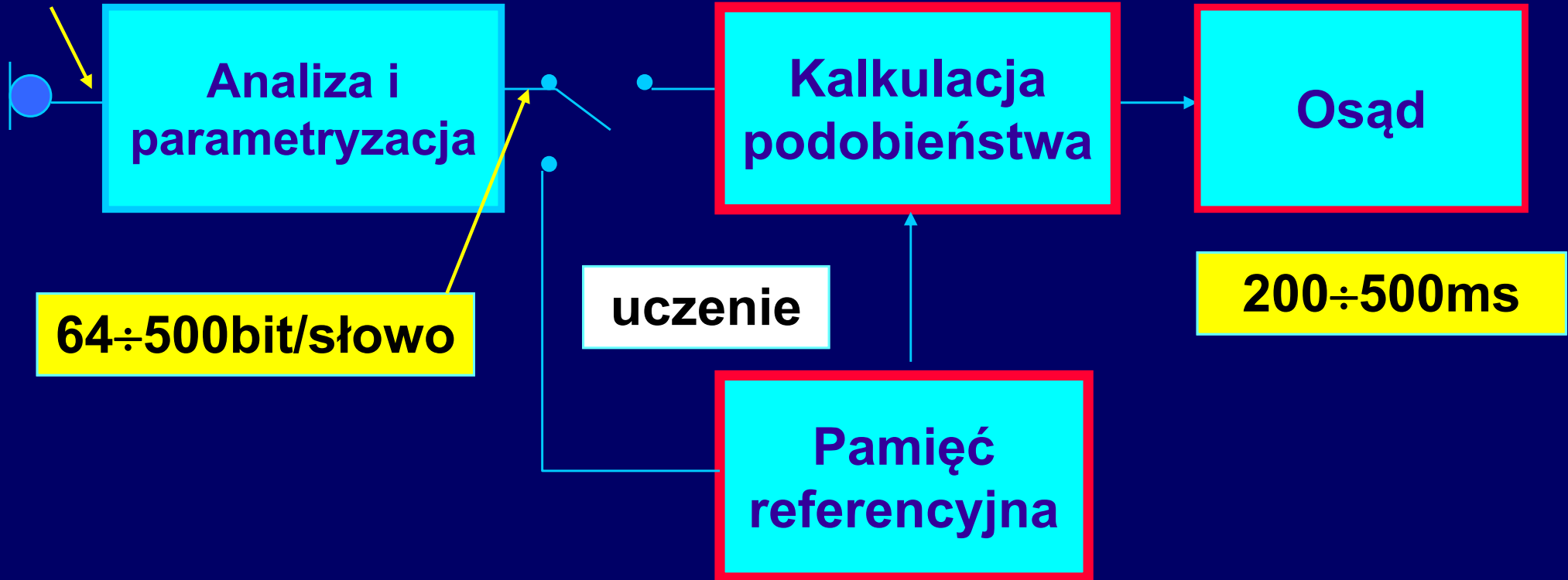
Automatyczna klasyfikacja

Segmentacja
redukcja
danych

Badanie odległości
ciągów binarnych

64kbit/s

rozpoznawanie



64÷500bit/słowo

uczenie

200÷500ms

Pamięć
referencyjna

Osąd

Kalkulacja
podobieństwa

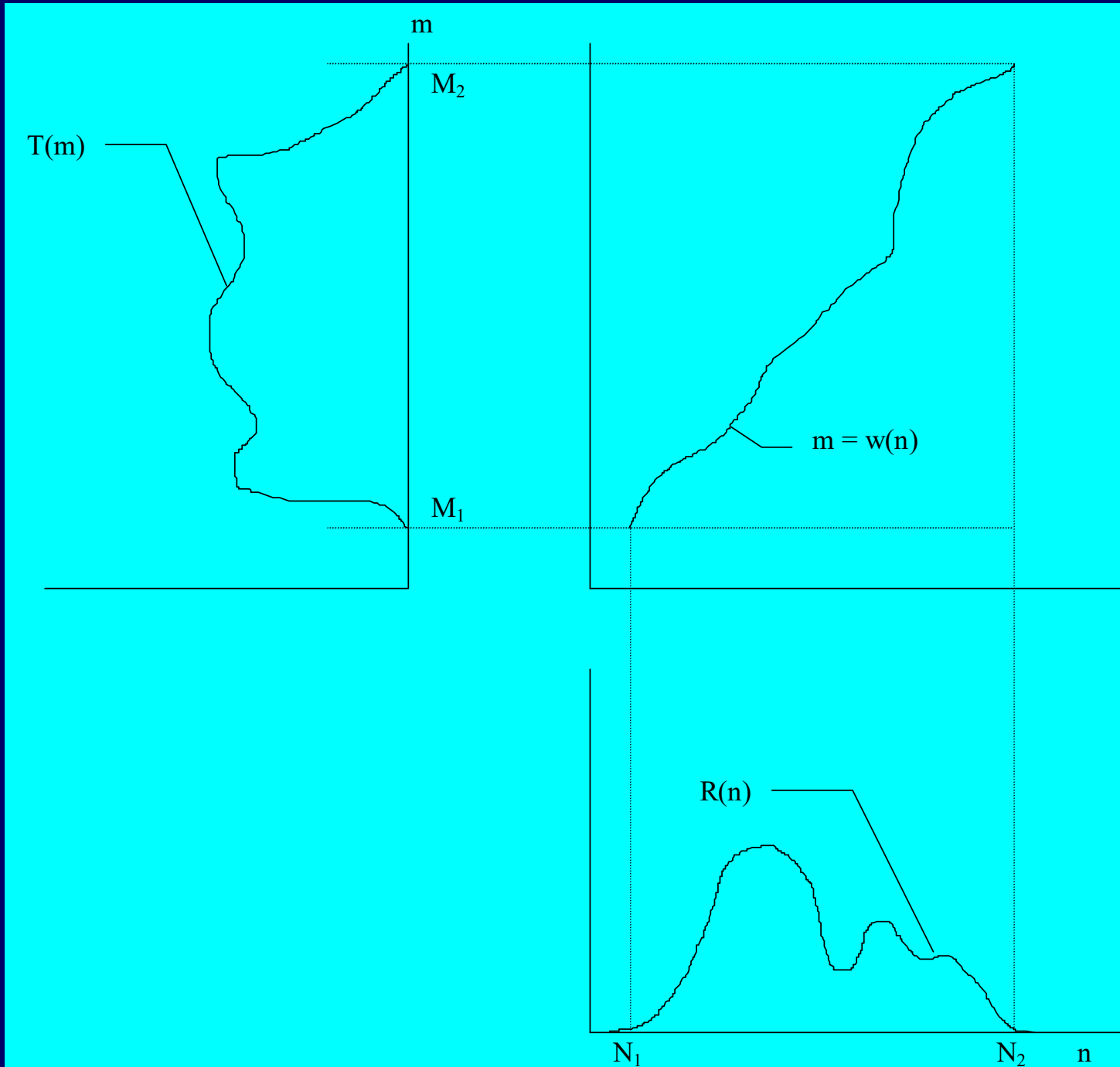
Analiza i
parametryzacja

ARM – systemy decyzyjne

Metody rozpoznawania izolowanych wyrazów:

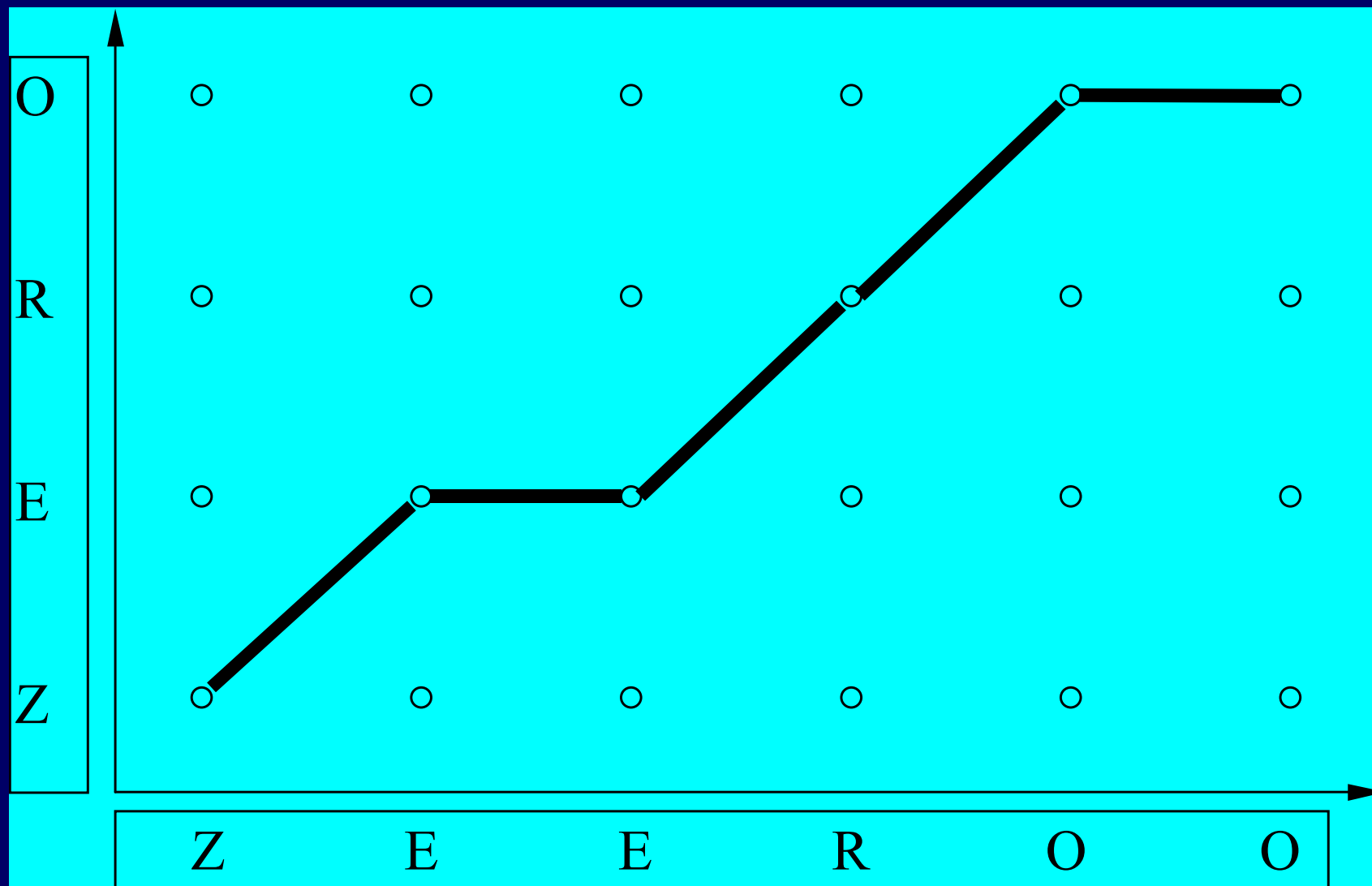
- nieliniowa normalizacja czasowa
- niejawne modele Markowa (HMM)
- sztuczne sieci neuronowe
- metoda zbiorów przybliżonych

Algorytm nieliniowego dopasowania czasowego



**Proces
nieliniowego
dopasowania
czasowego**

Algorytm nieliniowego dopasowania czasowego



Ilustracja procesu nieliniowego dopasowania czasowego w przypadku izolowanych wyrazów

Algorytm nieliniowego dopasowania czasowego

- Dopasowanie można przedstawić jako funkcję $m = w(n)$

przy spełnionych warunkach brzegowych:

$$M_1 = w(N_1) \quad M_2 = w(N_2)$$

oraz warunków ciągłości (następstwo segmentów)

$$w(n+1) - w(n) = 0,1,2 \quad (w(n) \neq w(n-1))$$

$$w(n+1) - w(n) = 1,2 \quad (w(n) = w(n-1))$$

- Dystans skumulowany jest miarą wskazującą na podobieństwo obiektu do wzorca:

$$D_T = \min_{\{w(n)\}} \sum_{n=1}^N D(R(n), T(w(n)))$$

$$D_A(n, m) = D(n, m) + \min_{q \leq m} D_A(n-1, q)$$

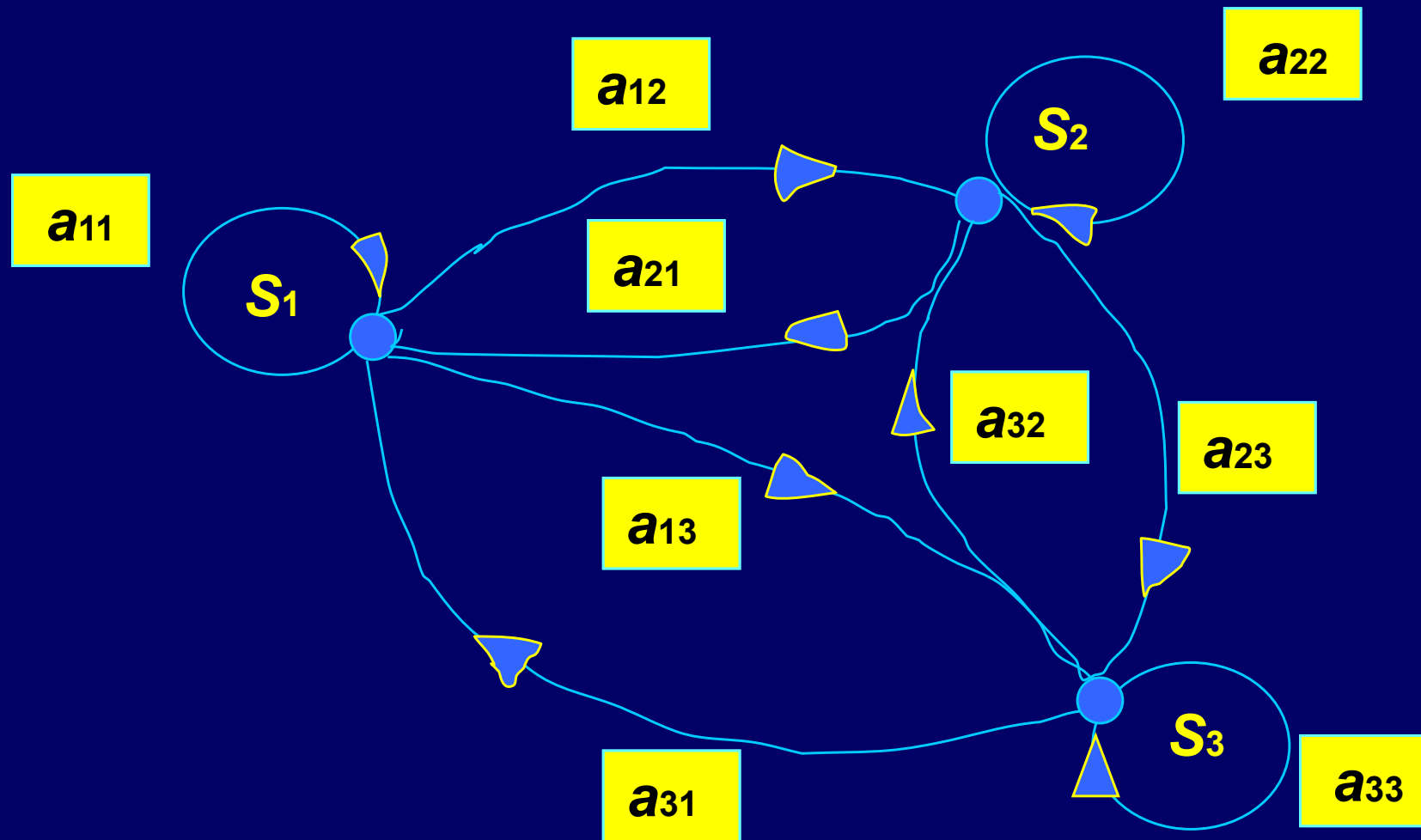
gdzie: $D_A(n, m)$ - minimalny dystans skumulowany względem punktu (n, m) siatki

HMM

- Dane słowo S_m w słowniku M możliwych słów jest reprezentowane ciągiem m zdarzeń O
- Każde słowo w słowniku jest opisane łańcuchem Markowa (HMM), dla M słów $\Rightarrow M \cdot \text{HMM} \{L_1, L_2, \dots, L_M\}$
- procedura dopasowania polega na obliczeniu sumarycznego prawdopodobieństwa (zdarzeń i przejść), że dany ciąg zdarzeń O został wygenerowany przez dany model L
- Prawdopodobieństwo to dane jest wzorem:

$$P_m = \Pr(O|L_m)$$

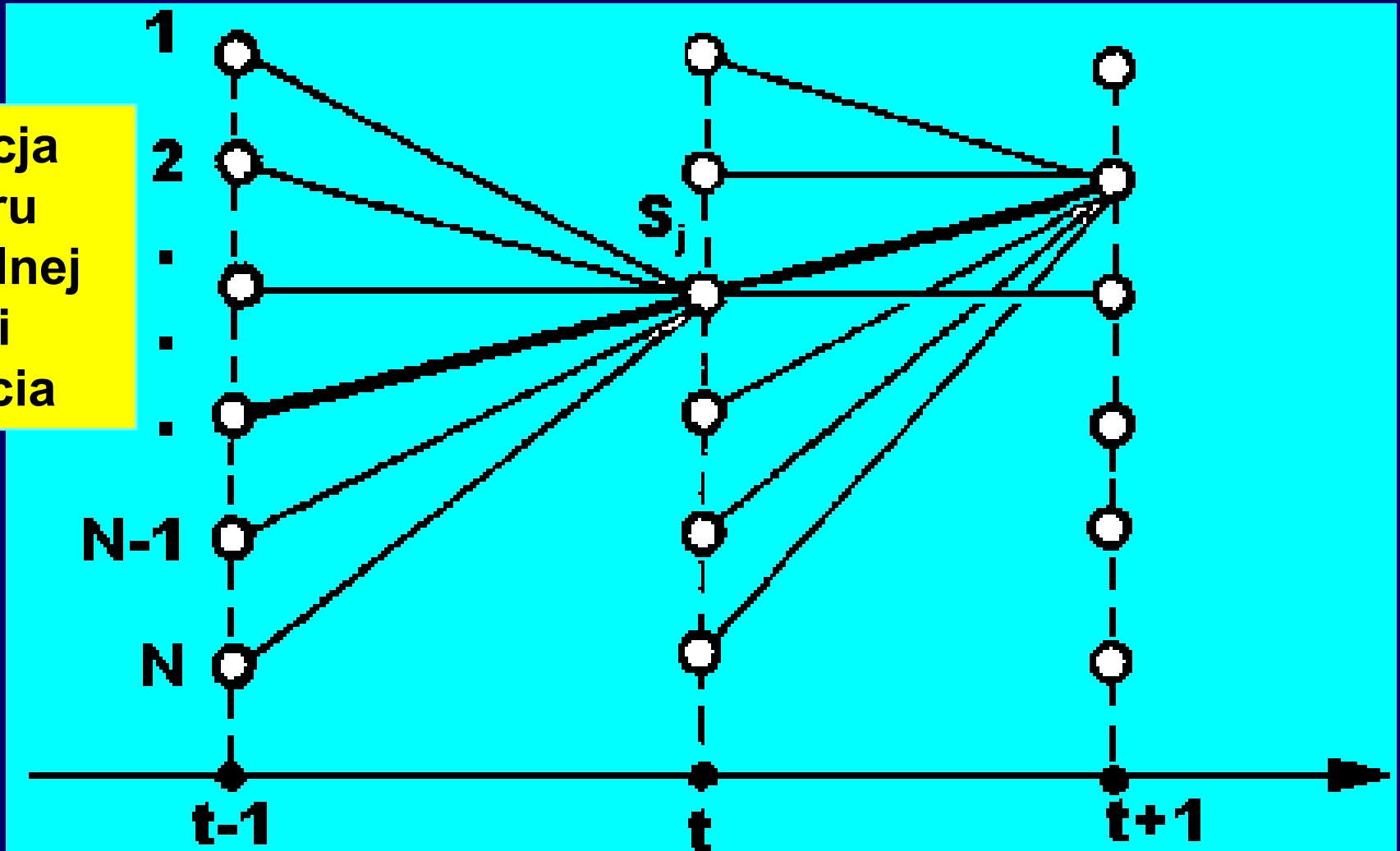
HMM



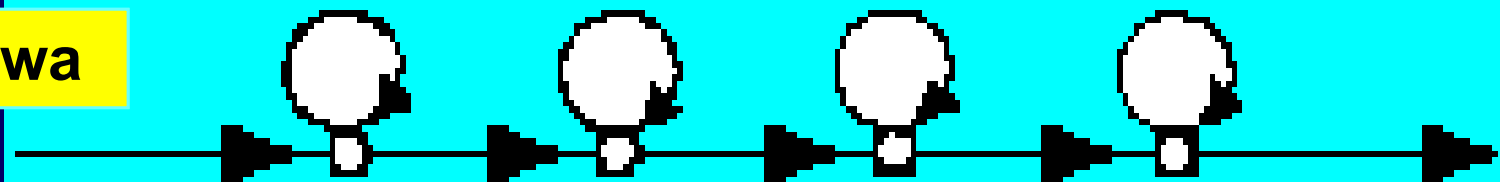
Ilustracja stanów i prawdopodobieństwa zdarzeń procesu Markowa

HMM

Ilustracja
wyboru
optymalnej
drogi
przejścia



Model słowa



- Sing
- SPASM Information
- Vocal Tract Shape v
- Tract Radii Display t
- Glottis Editor g
- Performance Editor p
- Phoneme Synthesis s
- Diphone Synthesis d
- Shape Interpolator i
- Formant Editor f
- Record Voice r
- Record Consonant c
- Demo Soundfiles x
- Set Sys. Defaults a
- Nasal Tract Shape n
- Help ?
- Hide h
- Quit q

Demo Window

Demo
Play Soundfiles
Generated by the
SPASM System

- sounds/aheeoo
- sounds/machines
- sounds/arpeg
- sounds/nasals
- sounds/fricatives
- diphones/diph
- diphones/nasals
- diphones/plosives
- diphones/cresc
- nphones/shiela
- nphones/vocaliz

Glottal Excitation Editor

Time Waveform

Log Magnitude Spectrum

Mix: Last Glottis This Glottis

Edge Beg.

Edge End

Harm.

Gain

Noise Gain

Load Save

Vocal Tract Shape

Vocal Tract Editor

Glott. Refl. Gain

Lip Refl. Gain

Mono
 Lp/Nk
 Lp/Ns
 Nk/Ns

Neck

Lip

Nose

Velum Opening

Phoneme Synthesis and Library

Time Waveform

Log Magnitude Spectrum

Perform Transform and Display Log Magnitude

Noise Impulse From:

Sing Synth

Load Save

Performance Feature Editor

Pitch Hz.

Vib. Freq. Hz.

Vib. Amt. %

Vib. Del. sec.

Rnd. Per. sec.

Rnd. Amt. %

Gliss. Amt. %

Gliss. Time sec.

Load Save

Tract Section Radius

Noise (turbulence) generator
Position of Noise Injection

Diphone Synthesis and Library

Init. Shape

Init. Glottis

Fin. Shape

Initial | Transition | Final

Curve
 Linear
 HypTn
 Expon.
 LowFrm

Shape Space Interpolator

File Nam

- shapes/ah
- shapes/ee
- shapes/oo
- shapes/rrr
- shapes/mn
- shapes/nn

Mouse controls shape
 Slow
 Med.
 Fast

Tract Noise Editor

Log Magnitude Spectrum

Cutoff Hz.

Reson. (< 1.0)

Cutoff Hz.

Reson. (< 1.0)

Gain

Formant Editor

Transfer Function from Glottis to Mixed Output

Number of Formants to Adapt: 2 3 4 5



NET

5:05 PM

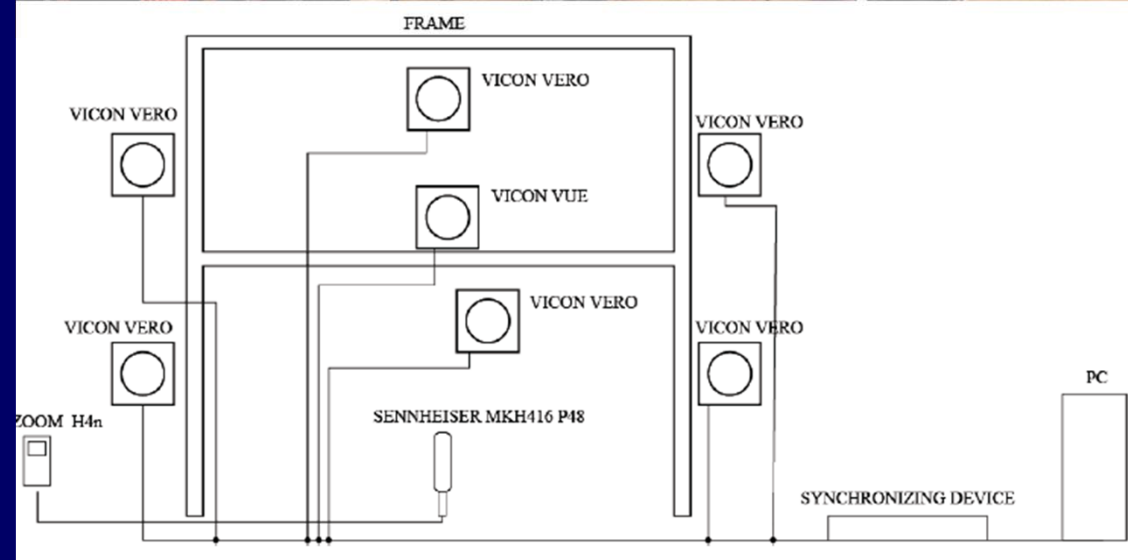
WED 25 SEP

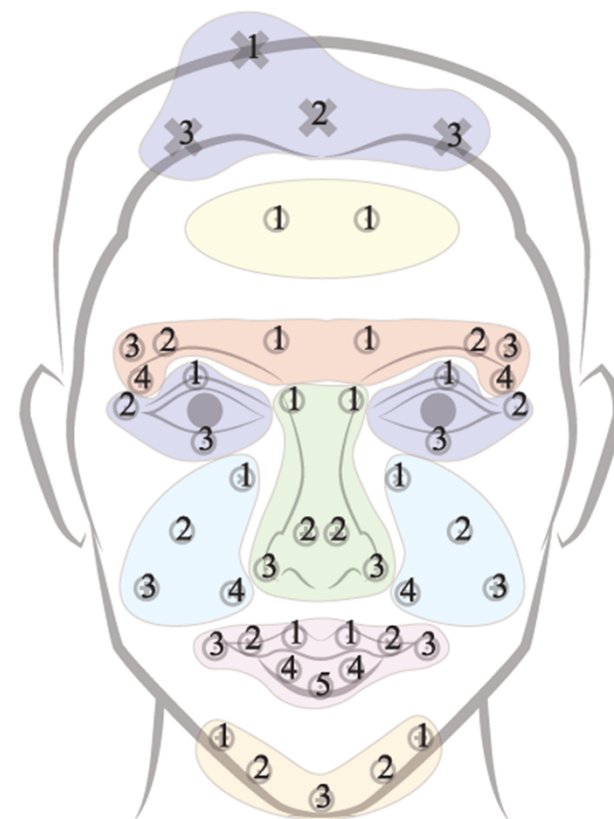
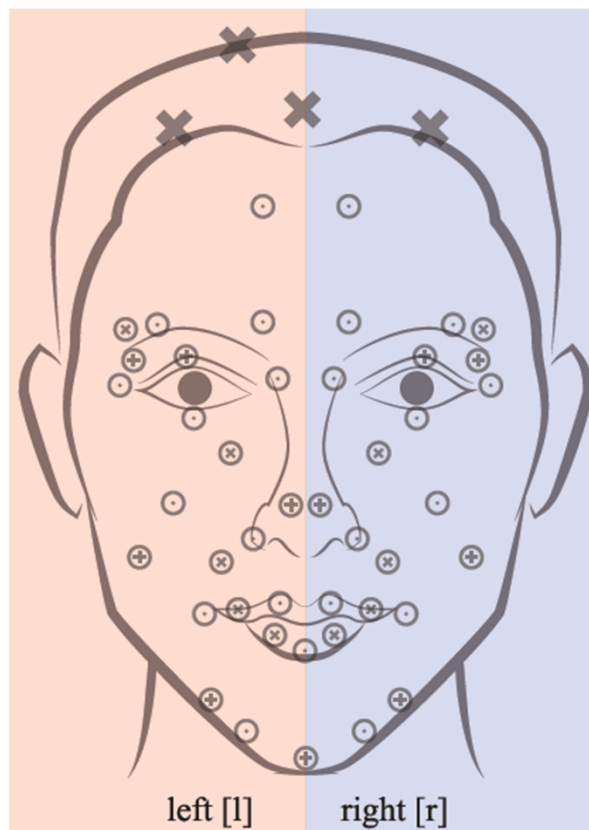
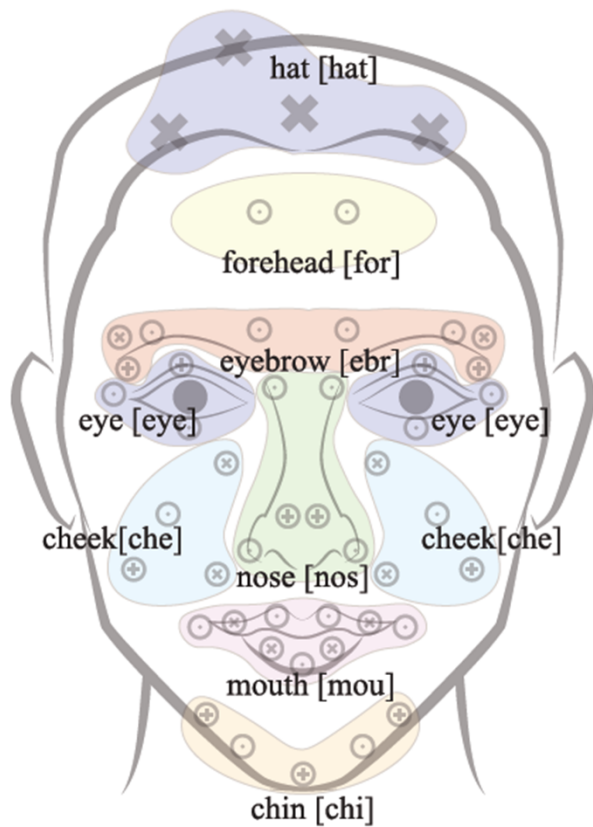
Audiowizualne rozpoznawanie mowy

Autor:

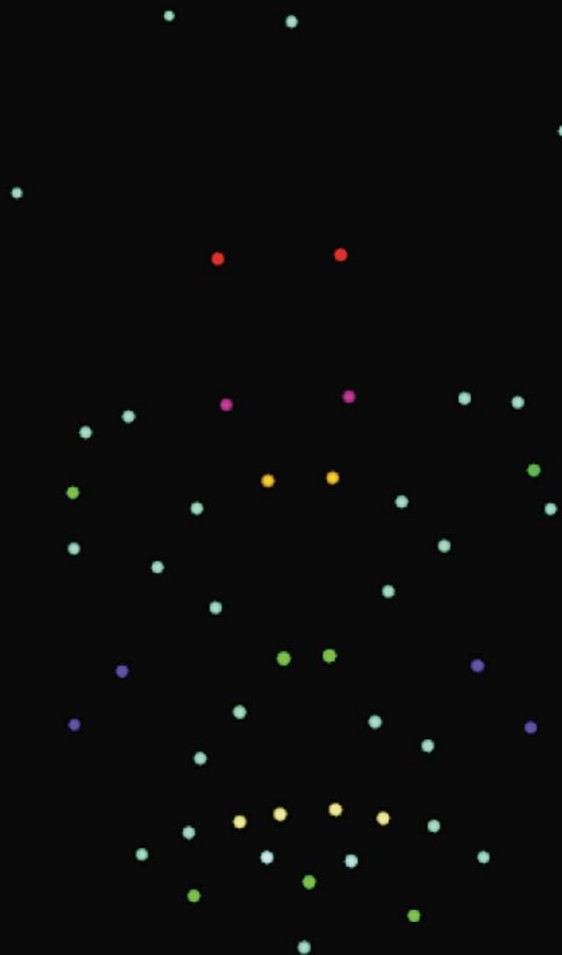
Piotr Bratoszewski

WYKORZYSTANE URZĄDZENIA

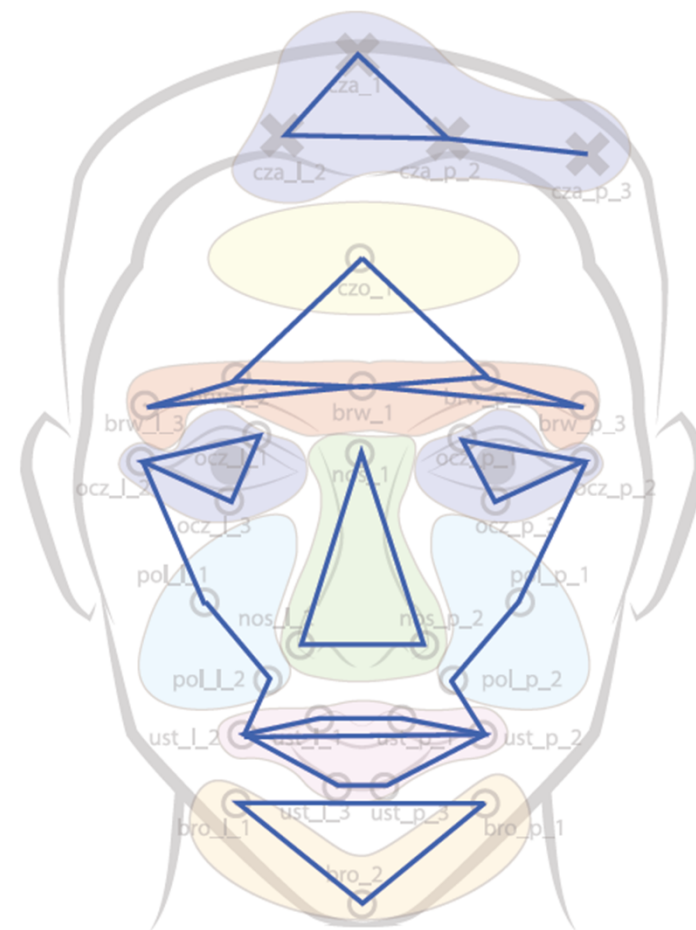
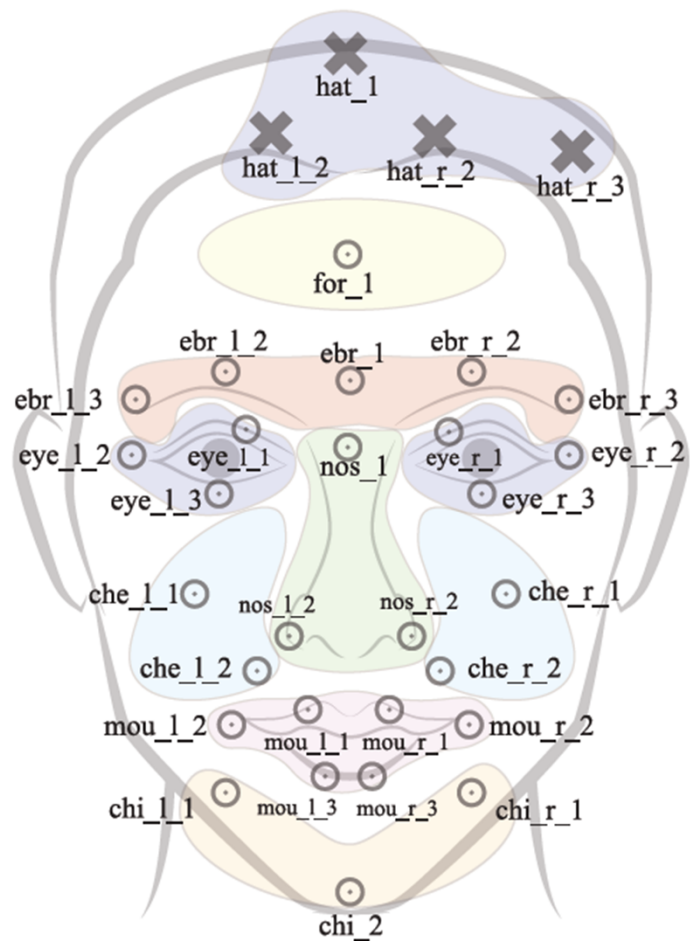




SPOSÓB NAZEWNICTWA MARKERÓW



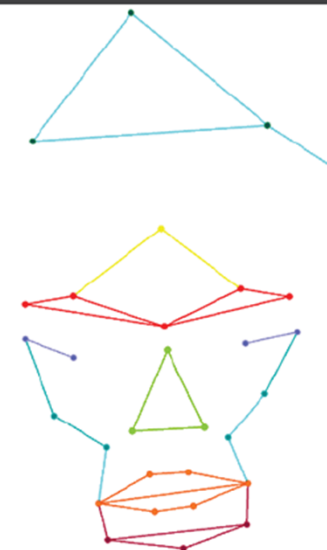
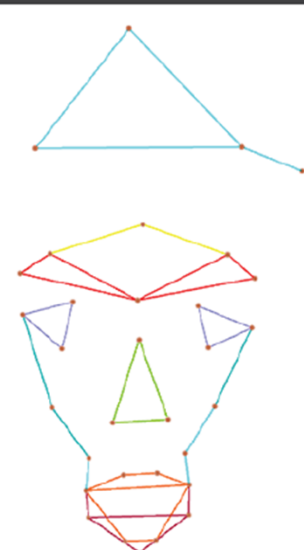
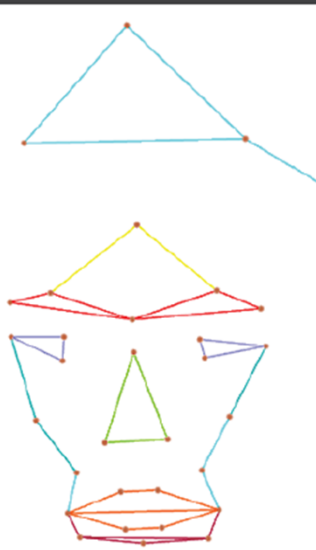
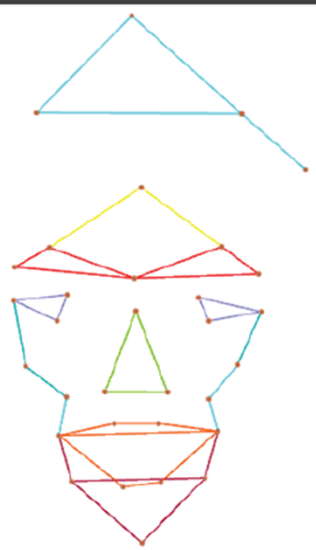
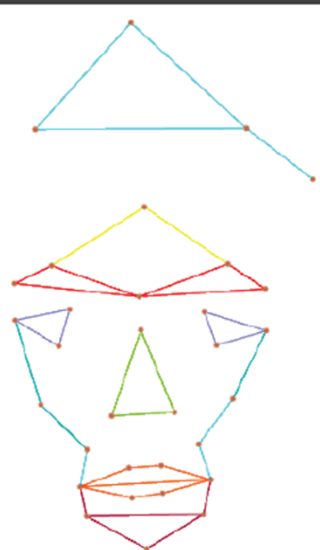
WIZUALIZACJA PUNKTÓW W PROGRAMIE BLADE



FINALNE UŁOŻENIE I NAZEWNICTWO



**OSOBY BIORĄCE UDZIAŁ
W NAGRANIACH**



**ZESTAW PRZYKŁADOWYCH
EKSPRESJI OSOBY NR 3**

What is Modality Corpus?

The MODALITY corpus consists of over **30 hours** of multimodal recordings. The database contains high-resolution, high-framerate stereoscopic video streams and audio signals obtained from a microphone array and a laptop microphone. The corpus can be employed to develop an AVSR system, as every utterance was labelled. Recordings in noisy conditions can be used to test the robustness of speech recognition systems.

[READ MORE >](#)

License



Explore Alofon

Corpus

Select speaker

Please select the type data you want to play and download



Words	Allophones
spi	spe sp
nac	ck spa
sp	sp sp
spar	ot ort
king	



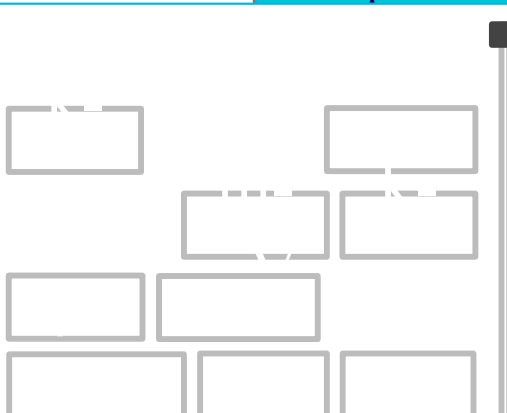
Explore Alofon

Corpus

Select speaker

Please select the type data you want to play and download

Words **Allophones**



Speaker selection row: 15 buttons, the 3rd button contains the number '3'.





Grupy wizemów

W1	W2	W3	W4	W5	W6
WARGOWE	DZIAŚŁOWE	MIĘKKO- PODNIEBIENIOWE	ZĘBOWO- WARGOWE	PODNIEBIENIOW O- DZIAŚŁOWE	ZĘBOWE
					



Grupy wizemów c.d.

W7	W8	W9	W10	W11
ROZSZERZONE	OTWARTE-ROZSZERZONE	OTWARTE-NEUTRALNE	OTWARTE-ZAOKRĄGLONE	WYSUNIĘTO-ZAOKRĄGLONE
				

Proces postępowania z danymi

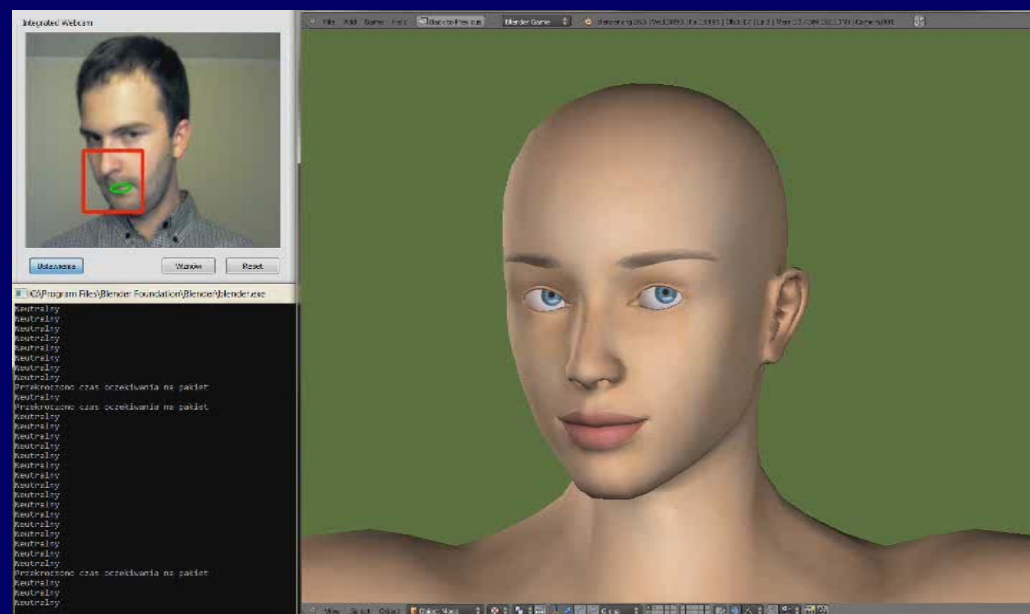


Parametry

Geometryczne	Teksturalne
<ul style="list-style-type: none">• Odległościowe – 39 parametrów• Kątowe – 20 parametrów• Powierzchniowe – 8 parametrów	<ul style="list-style-type: none">• Histogram obrazu ust w skali szarości – 32 parametry• Histogram obrazu ust HSV – 32 parametry• Histogram obrazu ust w skali szarości po equalizacji – 32 parametry• Histogram obrazu ust w skali szarości po filtrowaniu Clahe – 32 parametry• Wartości DCT dla obszaru ROI – 32 parametry

Wirtualny terapeuta – zastosowania w logopedii

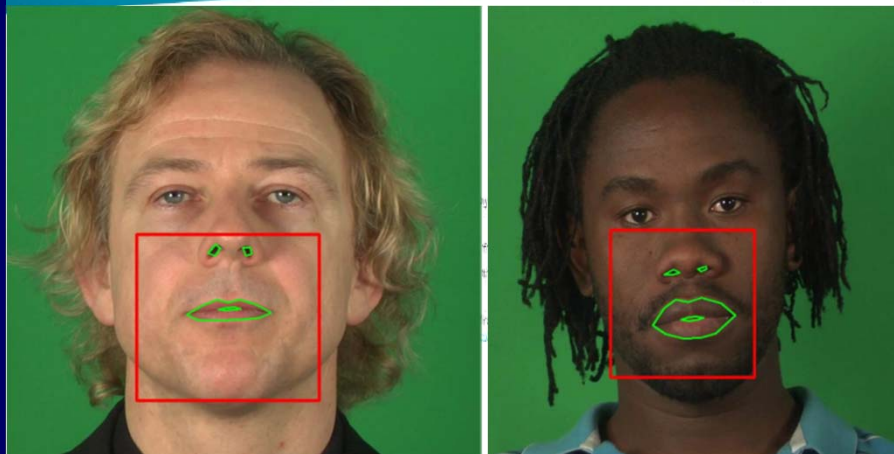
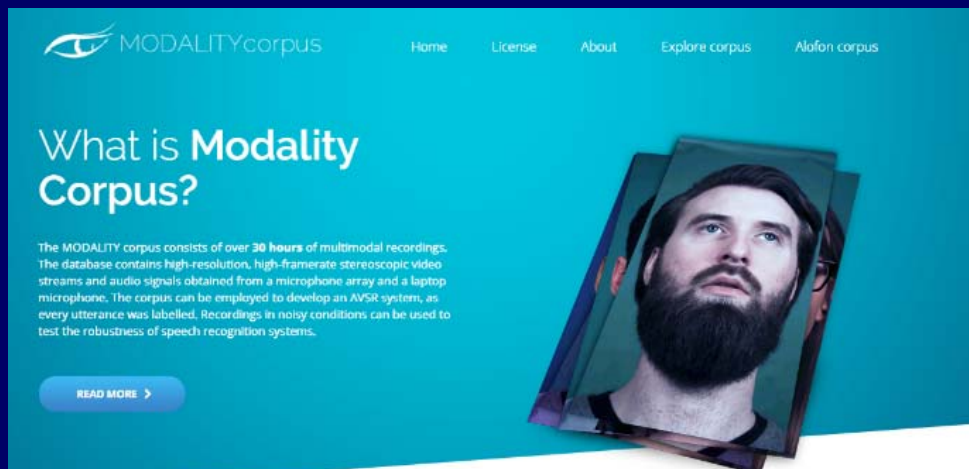
- Osadzona sztuczna sieć neuronowa (perceptron wielowarstwowy) wytrenowana na obrazach



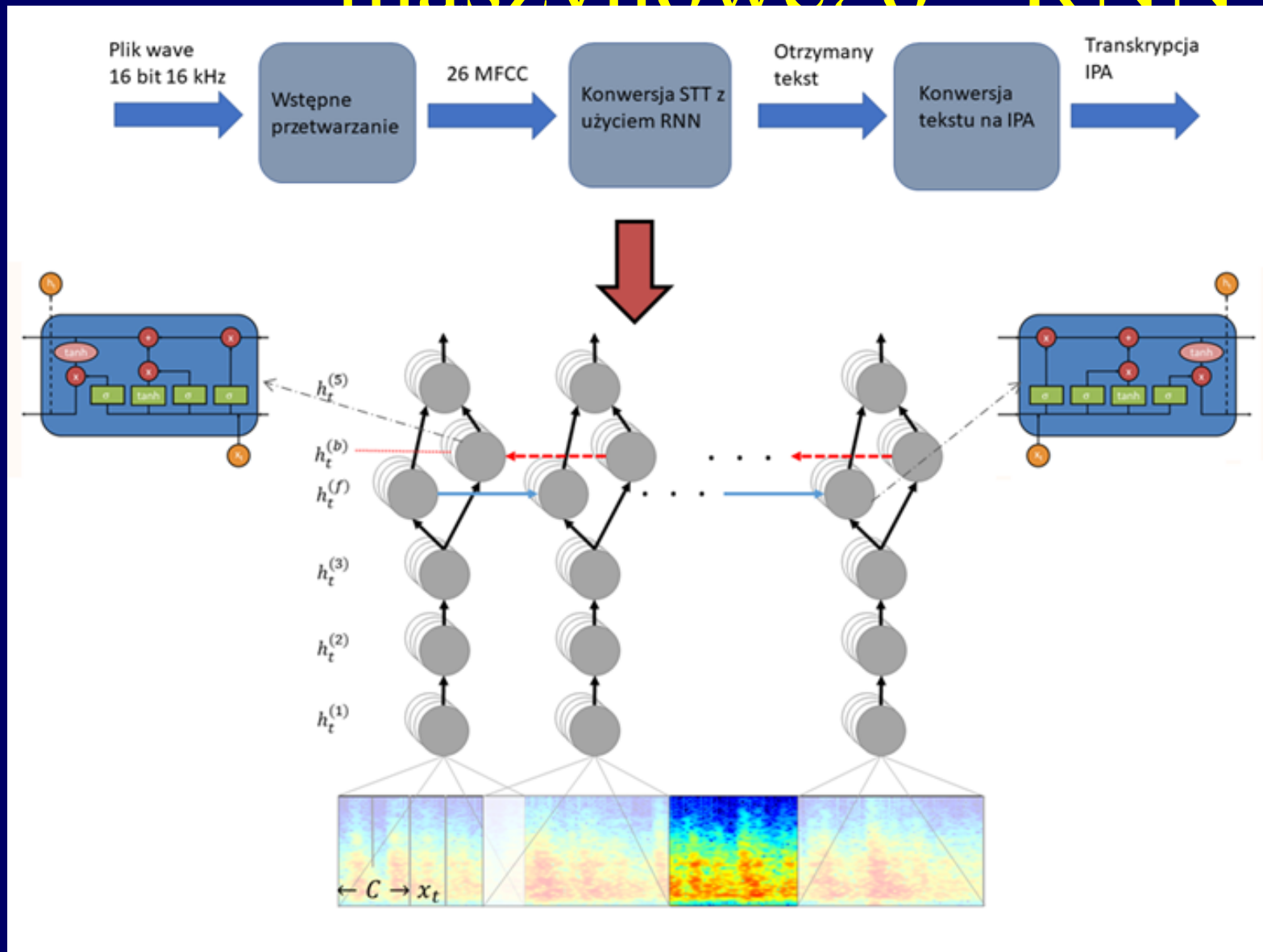
Agata.
wmv

Audiowizualna analiza mowy

- **Korpus MODALITY** (31 godzin nagrań audio-video mówców) www.modality-corpus.org
- **Korpus ALOFON**
- (1,5h godziny nagrań audio, video i Facial Motion Capture)

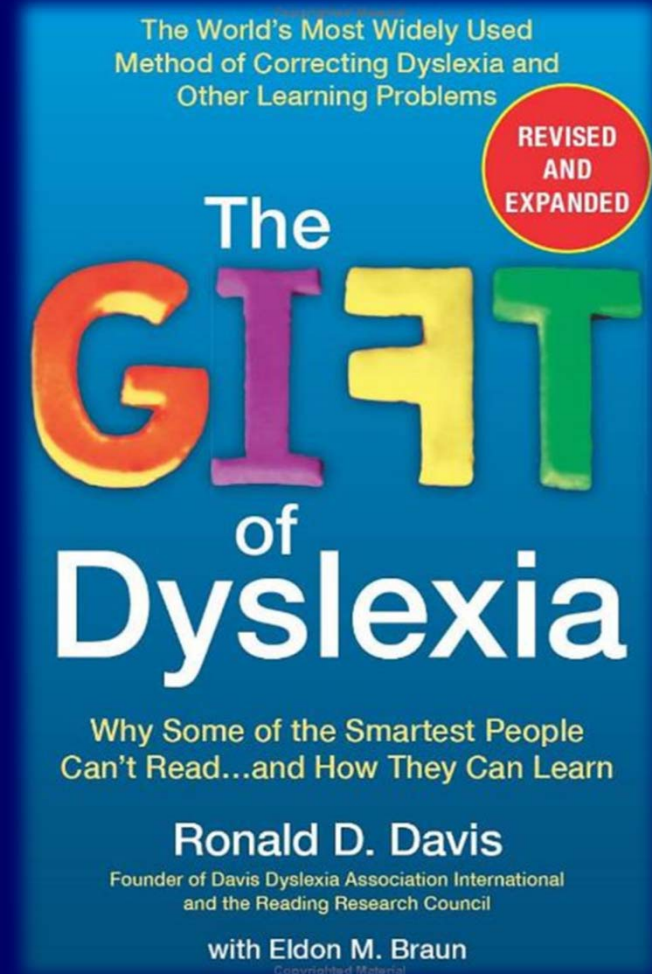
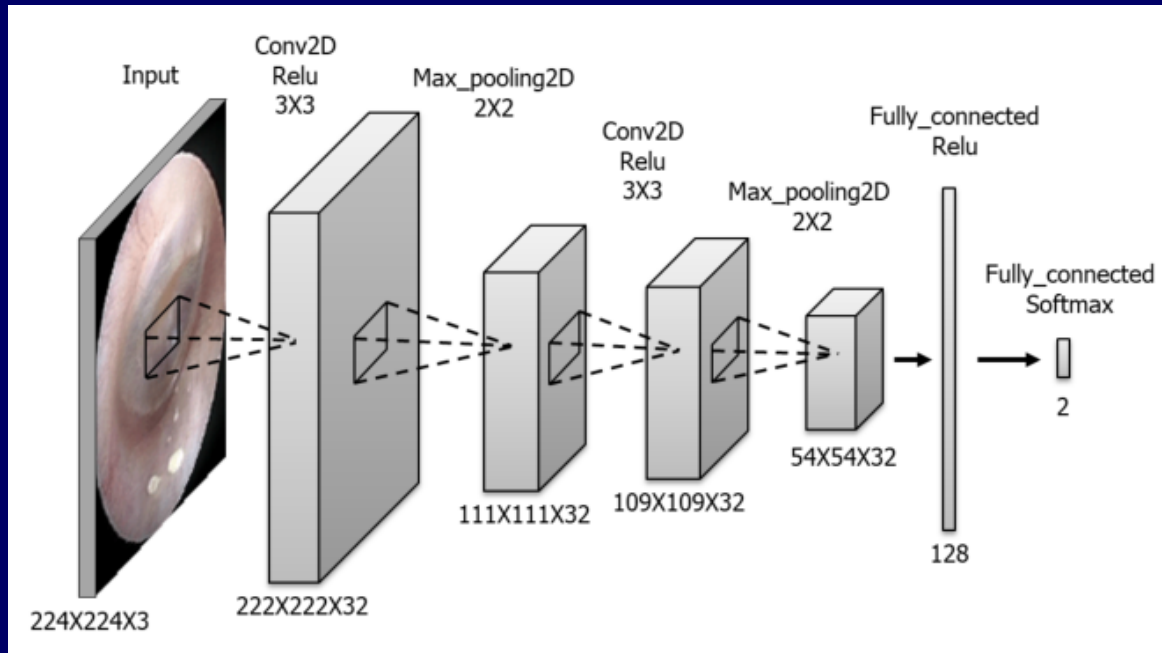


Wykorzystanie uczenia maszynowego - RNN



„Dar dysleksji” – mózg i głębokie uczenie

Davis D. Ronald Eldon M. Braun



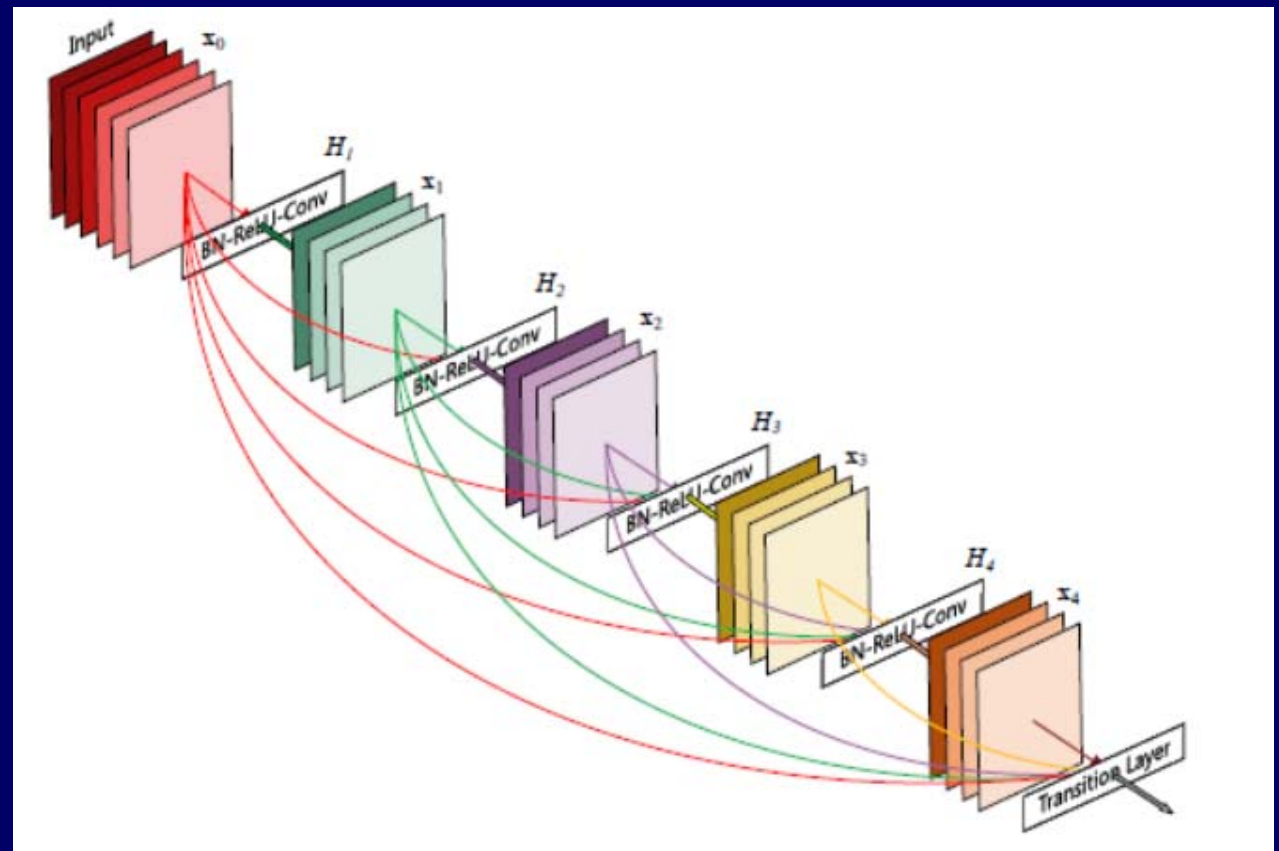
Sieci neuronowe DenseNet

Prace badawcze (publikacje Cornell University) wykazały, że sieci splotowe mogą być znacznie bardziej efektywne w treningu, jeśli zawierają krótsze połączenia między warstwami w pobliżu wejścia i tymi zbliżonymi do wyjścia. Dense Convolutional Network (DenseNet) łączy każdą warstwę do każdej innej warstwy z wyprzedzeniem.

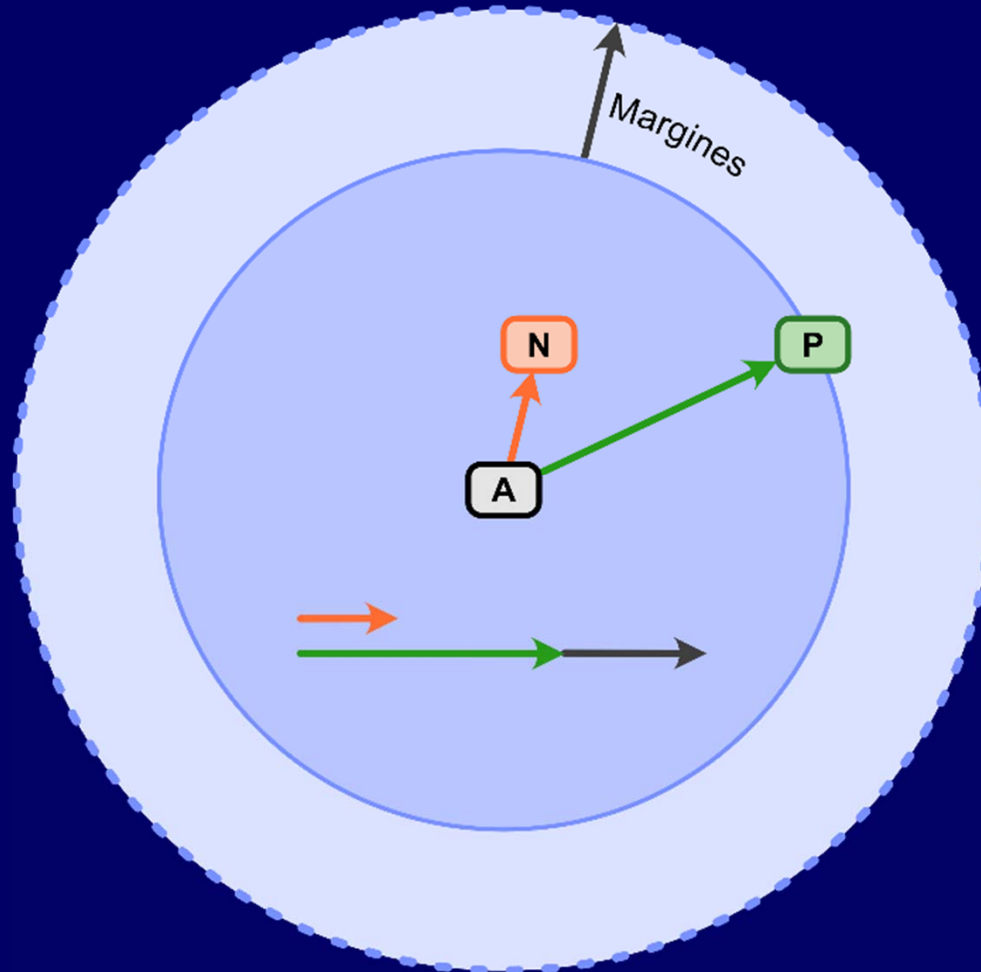
9	5	7	6	3	1	2	8	4
4 ₁	2 ₂	3 ₁	8	7	2	6	5	9
2 ₂	4 ₂	3 ₂	4	5	9	7	3	1
5 ₁	3 ₂	9 ₁	7	8	6	4	1	2
8	4	2	9	1	5	3	6	7
7	6	1	2	4	3	5	9	8
2	8	5	3	9	7	1	4	6
1	7	4	5	6	8	9	2	3
3	9	6	1	2	4	8	7	5

analogia do podstawowych operacji na obrazie

Kod i wstępnie wytrenowane modele są dostępne na stronie <https://github.com/liuzhuang13/DenseNet>.



Sieć neuronowa trenowana metodą Triplet Loss



Przykład trójki niespełniającej warunku odległości dla algorytmu triplet loss. Pomarańczowy wektor (AN) powinien być dłuższy od zielonego wektora (AP) o wartość marginesu.

Algorytm uczenia maszynowego, w którym bazowe (kotwiczące) dane wejściowe są porównywane z pozytywnymi (prawdziwymi) i negatywnymi (fałszywymi) danymi wejściowymi.

Zastosowania rozpoznawania mowy

- Programy i urządzenia przeznaczone dla osób niepełnosprawnych
- Sterowanie urządzeniami za pomocą głosu, np. telefonu komórkowego, komputera, inteligentnego domu, urządzeń samochodowych
- Nawigacja stroną internetową
- Gry edukacyjne
- Rozpoznawanie osób
- Pisanie tekstu
- Aplikacje multimedialne
- Zabawki dla dzieci
- Robotyka

Dziękuję za uwagę