

A6M33BIO - Biometrie

Biometrické metody založené na rozpoznávání hlasu I

Doc. Ing. Petr Pollák, CSc.

13. prosince 2018 - 16:24

- **Úvod**

- Aplikace a principy hlasové biometrické verifikace
- Základní popis vzniku řeči (hlasu)

- **Řečové charakteristiky a příznaky pro identifikaci**

- Základní frekvence řeči
- Spektrální charakteristiky
- Formanty
- Expertní identifikace, spektrografické metody
- Kepstrum, kepstrální vzdálenost

I. část

Identifikace řečníka
(základní principy, variabilita hlasu)

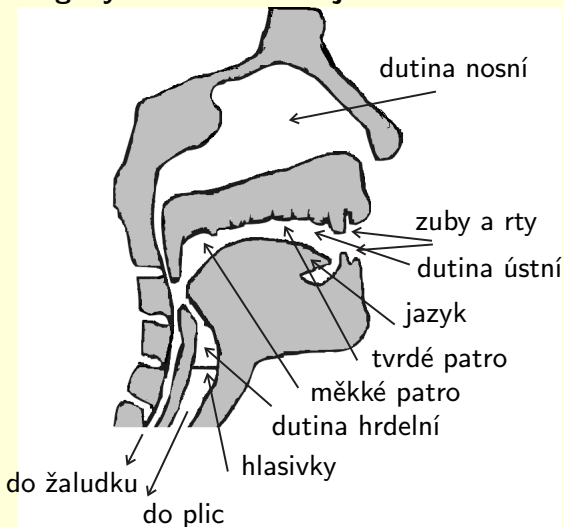
1) přesná identifikace totožnosti mluvího

- kriminalistická a soudní praxe - forenzní aplikace (dosud subjektivní fonetická a lingvistická analýza)
- identifikace pro přístup k zabezpečeným systémům (bankovní účty, přístup do chráněných objektů/systémů, bezpečnostní kontroly, apod.)

2) identifikace mluvího s největší podobností hlasu

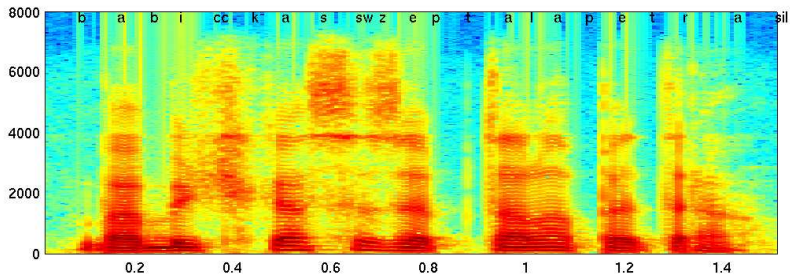
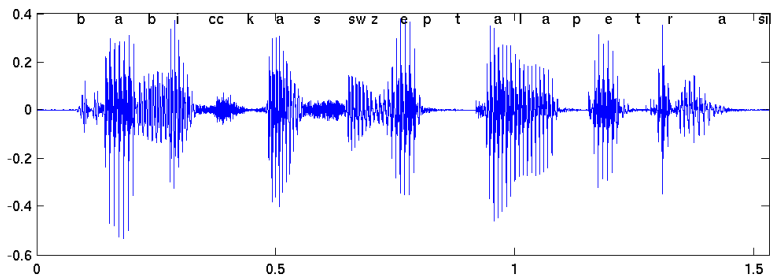
- př. - identifikace volajících v call-centrech
- komplexní rozpoznávače řeči (LVCSR - diktovací systémy, transkripční systémy pro přepis rozhlasových/TV zpravodajství)
 - modely pro konkrétního mluvího
 - skupinové modely
 - modely závislé na pohlaví mluvího

Artikulační orgány hlasového ústrojí člověka



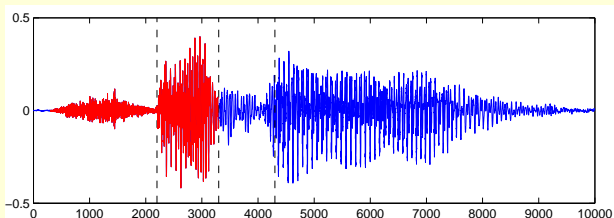
Produkce řeči : frekvenční modifikace širokopásmového buzení proudem vzduchu průchodem dutinami (rezonátory) hlasového ústrojí

Časový a spektrální obsah řeči

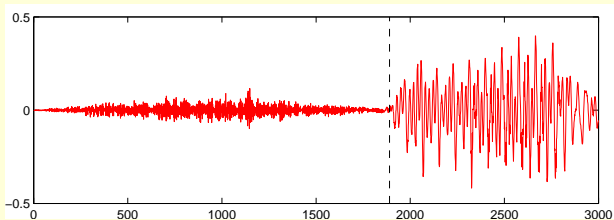


Řečové hlásky v časové oblasti

Slovo “šedý”



Slabika “še”

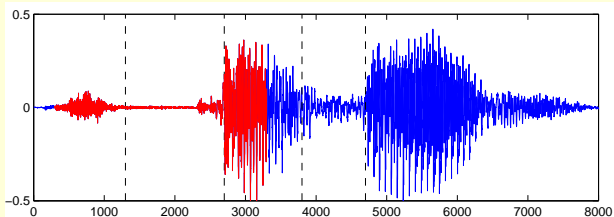


Hláška “š” ... neznělá, šumový charakter

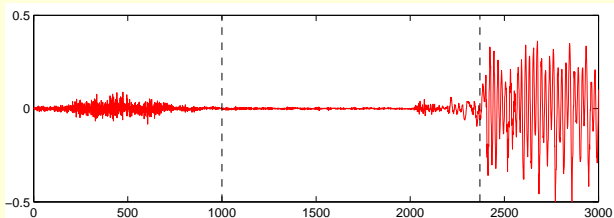
Hláška “e” ... znělá, periodický charakter (harmonická struktura)

Řečové hlásky v časové oblasti

Slovo "čtyři"



Slabika "čty"

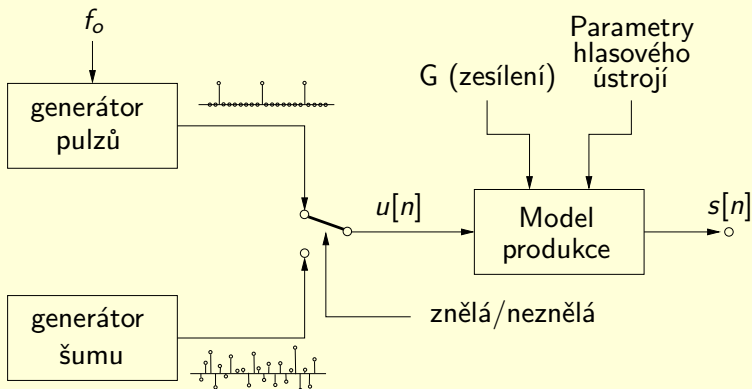


Hláška "č" ... neznělá, šumový charakter

Hláška "t" ... plozivní, okluze (závěr) + exploze

Hláška "y" ... znělá, periodický charakter (harmonická struktura)

Produkce řeči - signálový model vzniku řeči

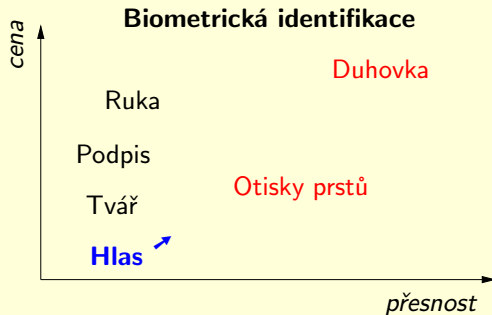


Závislost na mluvčím (rozdíly v rozměrech hlasového ústrojí):

- 1 vokální trakt (barva hlasu)**
 - souvislost s anatomii rezonátorů (dutin) hlasového ústrojí
 - AR model, snadná identifikace parametrů (LPC)
- 2 generování pulzů (výška hlasu)** - daná vlastnostmi hlasivek

Odlišnost mluvčích :

- originalita hlasu - výška a barva
- originalita stylu - různá doba trvání hlásek
- obecná variabilita jednotlivých realizací - PROBLÉM
- možnost napodobení hlasu - PROBLÉM

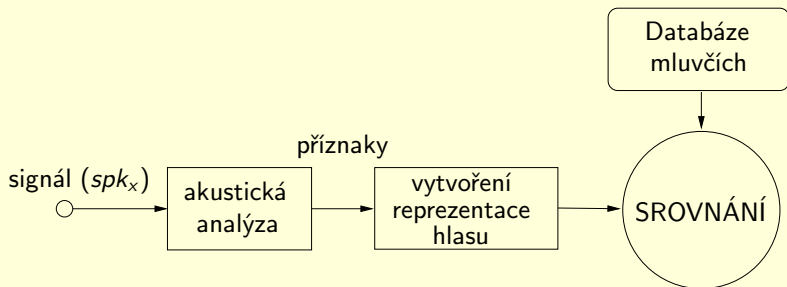


Motivace pro použití hlasové identifikace

- přirozenost komunikace, relativně jednodušší realizace
- jediná volba při dostupnosti pouze hlasového záznamu

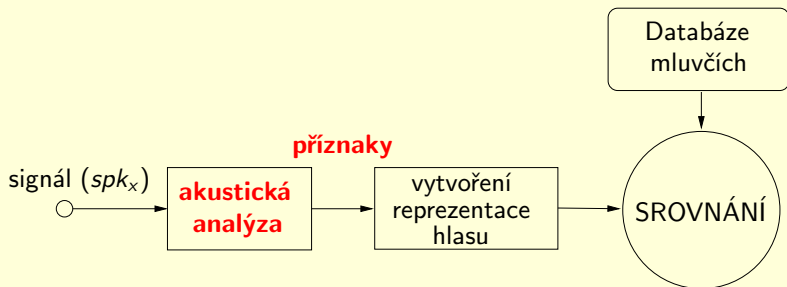
Možnosti identifikace mluvčího

- historické klasické přístupy :
 - **expertní rozhodování** (fonetici, lingvisté)
- moderní automatizované metody :
 - **rozpoznávání mluvčího**



Možnosti identifikace mluvčího

- historické klasické přístupy :
 - **expertní rozhodování** (fonetici, lingvisté)
- moderní automatizované metody :
 - **rozpoznávání mluvčího**



III. část

**Řečové charakteristiky
a možnosti využití pro identifikaci**

Obecné požadavky pro příznaky resp. systémy identifikace

- **vysoká variabilita pro různé mluvčí**
 - **nízká variabilita pro jednoho mluvčího**
(možné vlivy - aktuální stav, nálada, stres, hluk, styl promluvy)
-

- snadný a efektivní výpočet
- odolnost vůči šumu a zkreslení (výše zmiňované jevy)
- odolnost proti imitaci hlasu



Vnitřní charakteristiky - související s vytvářením řeči

Získané charakteristiky - souvisejí s dynamikou pohybů hlasového traktu (dané prostředím)

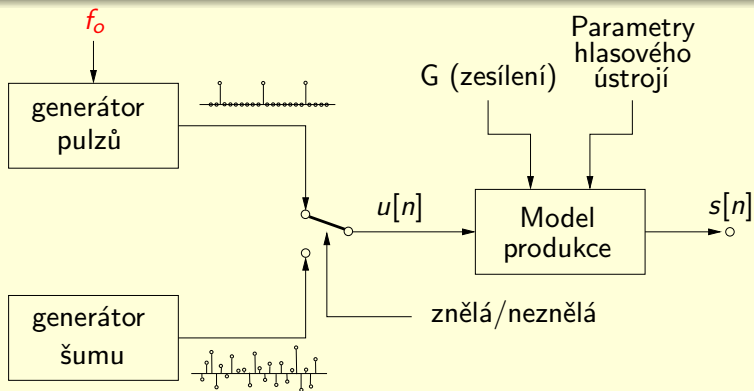
Vnitřní charakteristiky řečníka

- dané fyzikálními rozměry hlasového ústrojí
(lze obtížně napodobit)
- ovlivnitelné zdravotním stavem
(např. nosní dutina : neměnné rozměry při artikulaci,
mírné nachlazení = zásadní změna)

Získané charakteristiky řečníka

- styl mluvy (časování, intonace, hrubost, živost, síla, srozumitelnost)
→ jako celek komplexní charakteristika řečníka
(používáno člověkem při přirozené identifikaci)
- nemusí být snadno modelovatelné různými modely
- způsob řeči lze snadno napodobit

Základní tón řeči



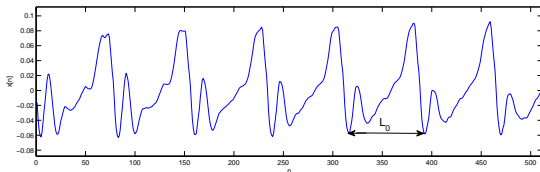
- základní frekvence $f_o = 1/T_o$
- pro znělé hlásky s harmonickou strukturou
- souvisí s kmitáním hlasivek
- hodnota f_o je ovlivněna vlastnostmi hlasivek (pružnost, hmotnost, délka)
→ hrubá charakteristika mluvčího

Odhad základního tónu řeči

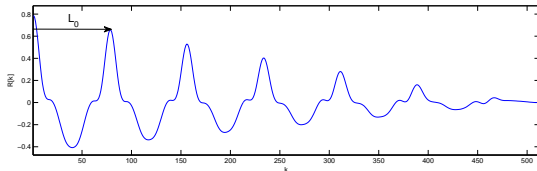
f_o základní tón (frekvence) řeči $f_o = \frac{1}{T_o}$
 T_o (L_o) základní perioda (v sekundách vs. ve vzorcích)

Nejčastější metoda odhadu - na bázi autokorelační funkce
(hledání postranního maxima autokorelační funkce)

segment signálu

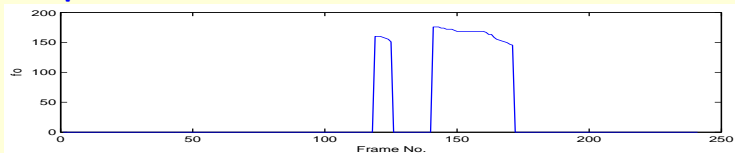


odhad autokorelační funkce

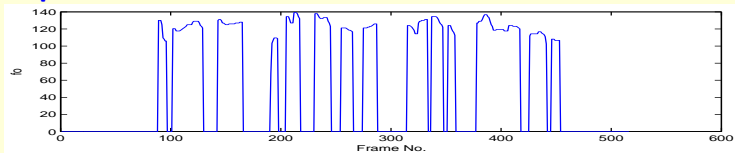


Průběh základního tónu v promluvě

Krátká promluva - slovo



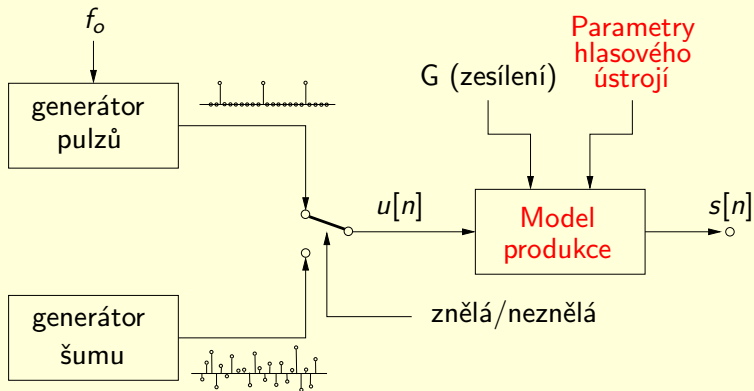
Delší promluva - věta



Průběh f_0 v promluvě → získaná (naučená) charakteristika

Průměrná hodnota f_0 → vnitřní charakteristika (výška hlasu)

Spektrální charakteristiky řeči



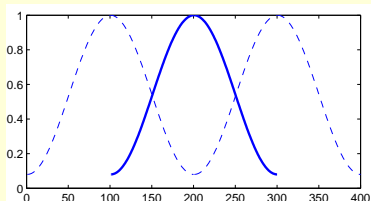
- spektrální charakteristiky souvisí s vokálním traktem
- otázka vhodné reprezentace pro identifikaci

Odhad spektra na bázi DFT:

- **řeč je obecně nestacionární signál** \Rightarrow nutná segmentace a sledování vývoje krátkodobého spektra (spektrogram)
- **řeč je kvazistacionární**
(tj. stacionární v krátkém časovém intervalu - cca 10-100 ms)
 \Rightarrow 20-30 ms - typická délka krátkodobého segmentu
- **DFT spektrum je ovlivněno proakováním**
 \Rightarrow nutné váhování vhodným oknem (Hammingovo)
 \Rightarrow nutná segmentace s překryvem (obvykle 50%)

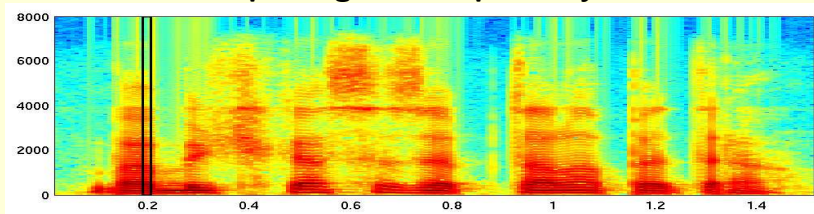
$$w[n] = 0,54 - 0,46 \cos \frac{2\pi n}{N}$$

$$\text{pro } 0 \leq n \leq N - 1.$$



Přehled možností spektrální reprezentace promluvy

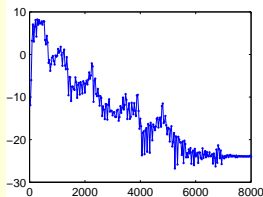
Spektrogram celé promluvy



Spektrální reprezentace vybraného segmentu

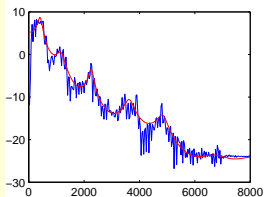
Sklon amplitudového spektra - vyšší kmitočty - nižší energie

DFT spektrum:



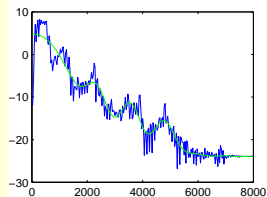
256 vzorků spektra
(amplitudové sp.)

LPC reprezentace:



16 koeficientů a_k
(autoregresní koef.)

Kepstrální koeficienty:



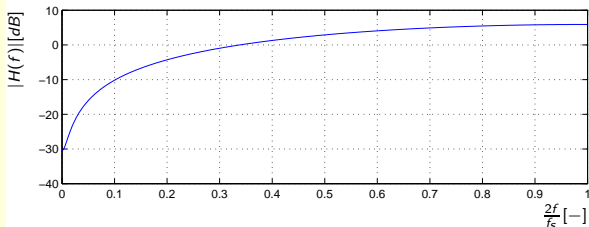
16 koeficientů c_n
(reálné keprstrum)

Preemfáze signálu

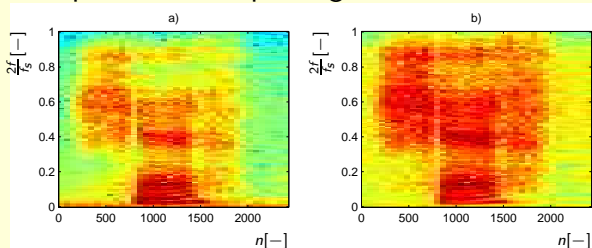
Kompensace sklonu amplitudového spektra

$$s'[n] = s[n] - m \cdot s[n - 1], \quad m = 0,97$$

Frekvenční charakteristika preemfázového filtru

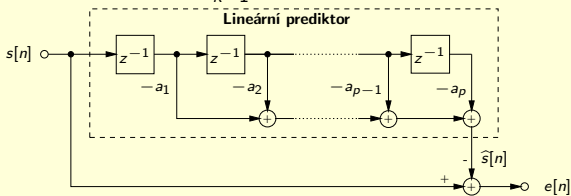


Ilustrace vlivu preemfáze ve spektrogramu



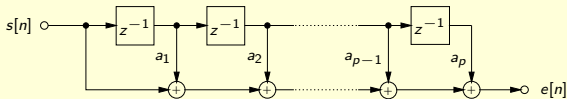
Lineární prediktivní analýza

Lineární predikce : $\hat{s}[n] = - \sum_{k=1}^p a_k s[n - k] .$



Chybový signál (míra kvality prediktoru)

$$e[n] = s[n] - \hat{s}[n] = s[n] + \sum_{k=1}^p a_k s[n - k] = \sum_{k=0}^p a_k s[n - k] .$$



IDEA: přesnější predikce \rightarrow nižší úroveň chybového signálu

Kritérium - výkon chybového signálu

$$J = E \left\{ e^2[n] \right\}$$

Hledání koeficientů $a_k \equiv$ Minimalizace chyby predikce
 \equiv hledání minima J , i.e.

$$\frac{\partial J}{\partial a_k} = 0, \quad \text{for } k = 1, 2, \dots, p \quad \Rightarrow \quad p \text{ lineárních rovnic}$$

Řešení a metody výpočtu (pro různé definice J):

- **autokorelační metoda** - nejčastěji používaný přístup
- Levinson-Durbinův algoritmus (rychlý výpočet autokor.met.)
- Burgův algoritmus - vychází z křížové struktury filtru

Autokorelační metoda, Yuleovy-Walkerovy rovnice

$$\begin{bmatrix} R[0] & R[1] & R[2] & \dots & R[p-1] \\ R[1] & R[0] & R[1] & & R[p-2] \\ R[2] & R[1] & R[0] & \ddots & R[p-3] \\ \vdots & & \ddots & \ddots & \vdots \\ R[p-1] & R[p-2] & R[p-3] & \dots & R[0] \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ \vdots \\ a_p \end{bmatrix} = - \begin{bmatrix} R[1] \\ R[2] \\ \vdots \\ \vdots \\ R[p] \end{bmatrix}$$

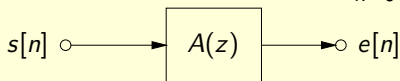
$R[k]$ autokorelační koeficienty analyzovaného signálu

VÝSLEDEK:

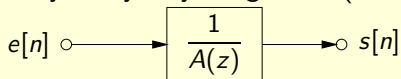
a_k autoregresní koeficienty (AR model signálu)

$P_p = R[0] + \sum_{k=1}^p a_k R[k]$ výkon chybového signálu

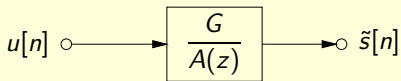
Dekorelační (analyzující) filtr : $A(z) = \sum_{k=0}^p a_k z^{-k}$



Syntéza se skutečným chybovým signálem (ideální případ)



Syntéza s umělým signálem s jednotkovým výkonem (AR model)
- G závisí na úrovni analyzovaného signálu ($G = \sqrt{P_p}$)



Spektrální vlastnosti AR modelu

Obecný popis AR modelu v Z-oblasti

$$\tilde{S}(z) = H(z) \cdot U(z)$$

Popis AR modelu ve frekvenční oblasti

$$S_{\tilde{s}}(e^{j\Theta}) = |H(e^{j\Theta})|^2 \cdot S_u(e^{j\Theta})$$

Vlastnosti a důsledky: - $S_u(e^{j\Theta})$ je ploché

→ tvar $S_{\tilde{s}}(e^{j\Theta})$ je kompletně zahrnut v AR modelu



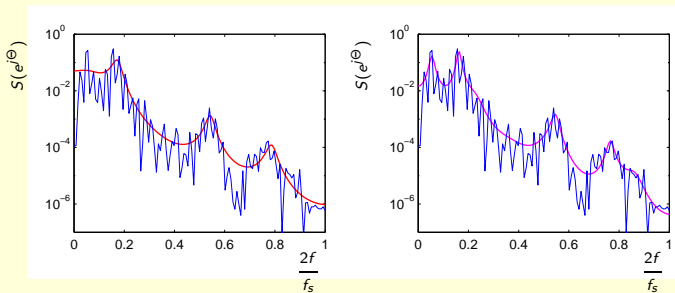
LPC spektrum (pokud $S_u(e^{j\Theta}) = 1$)

$$S_{\tilde{s}}(e^{j\Theta}) = |H(e^{j\Theta})|^2 = \frac{G^2}{|A(e^{j\Theta})|^2}$$



koeficienty a_k komprimovaná spektrální reprezentace

$$S_{\tilde{S}}(e^{j\Theta}) = |H(e^{j\Theta})|^2 \approx \frac{|S[k]|^2}{N}$$



- AR model: “all-pole” filtr, modeluje pouze špičky ve spektru (rezonátory v dutinách vokálního traktu)
- obecná špička = dvojice komplexně združených pólů
- vyšší řád AR modelu = více špiček v LPC spektru
→ typické hodnoty: $p = 10$ pro $f_s = 8$ kHz, $p = 16$ pro $f_s = 16$ kHz

- **Formant (formantové frekvence)**

→ centrální kmitočty rezonátorů vokálního traktu

- významné špičky ve **VYHLAZENÉM** krátkodobém spektru
- významné formanty F1 - F4 v pásmu do 4 kHz
- F5 - méně významný (obtížně odhadnutelný formant)
- !! Nezaměňovat se základním tónem řeči f_0
(f_0 není detekovatelné ve vyhlazeném spektru)



Souvislost s fyziologií vokálního traktu = vhodný vnitřní příznak
(formantové frekvence jsou nepřímo úměrné délce vok. traktu)

$$F_i = \frac{(2i - 1) \cdot c}{4 \cdot VTL}$$



Pro rozlišení mluvčích - vzdálenost sousedních formantů

- špičky LPC spektra - rezonátory = formanty
- F_i - formantová frekvence (centrální kmitočet rezonátoru)
- B_i - šířka pásma formantu
- špičky LPC spektra - určené **póly přenosové funkce** p_i

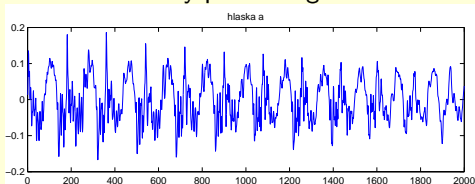
$$F_i = f_s \cdot \arg p_i / 2\pi$$

$$B_i = -f_s \cdot \ln |p_i| / 2\pi$$

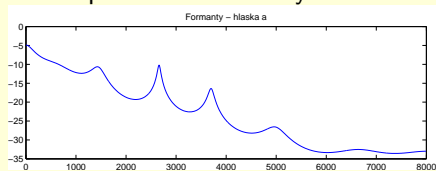
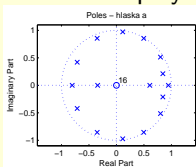
- Problémy:
 - obecně menší robustnost LPC analýzy (závislost na datech)
 - určení vhodného řádu (vliv přítomnosti šumu)
 - seřazení vypočítaných pólů (sledování stejného formantu)
 - vyřazení nadbytečného pólu (méně významné špičky)

Odhad formantů na bázi LPC - příklad

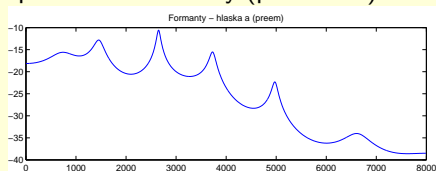
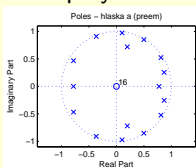
Časový průběh signálu



póly & LPC spektrum s formanty



póly & LPC spektrum s formanty (preemfáze)



Speciální příznaky pro rozpoznávání mluvčího

- F2 v “n”
- F3 v “u”
- F2 v “i”
- délka trvání “k”

- *obecnější formulace*
- hodnota formantu ve vybrané hlásce
- šířka pásma vybraného formantu ve vybrané hlásce
- směrnice poklesu formantu ve vybrané hlásce
- Průběh F0 ve vybrané větě (slově)
- průměrná hodnota F0 ve větě (slově)
- apod.

Kepstrum - definice a základní vlastnosti

Základní definice pomocí Z-transformace

$$\hat{c}[n] = \mathcal{Z}^{-1}\{\ln \mathcal{Z}\{x[n]\}\}$$

Přímý výpočet pomocí DFT

$$c_k[n] = \text{IDFT}\{\ln \text{DFT}\{x[n]\}\} \quad \dots \text{komplexní kepstrum}$$

$$c_r[n] = \text{IDFT}\{\ln |\text{DFT}\{x[n]\}|\} \quad \dots \text{reálné kepstrum}$$

$$c_v[n] = \text{IDFT}\{\ln |\frac{1}{N} \text{DFT}\{x[n]\}|^2\} \quad \dots \text{výkonové kepstrum}$$

Vlastnosti:

- $\hat{c}[n]$ nekonečně dlouhé, rychle ubývá k nule
- $c_k[n]$ konečně dlouhé, nesymetrické, informace o fázi
- $c_r[n]$ **konečně dlouhé, symetrické, inf. o ampl. spektru**
- $c_v[n]$ oproti $c_r[n]$ se liší pouze měřítkem a hodnotou $c[0]$
- ve všech případech vždy reálné hodnoty

Základní slovní přesmyčka: **spektrum** vs. **kepstrum**

Další vybrané přesmyčky:

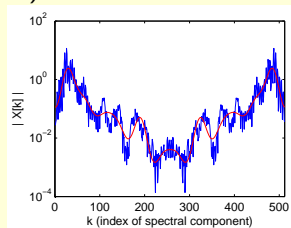
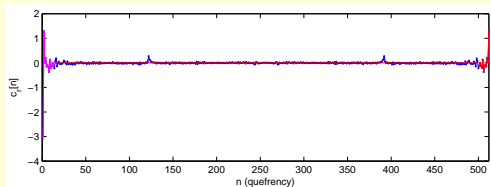
- **kvefrence** (frekvence) - základní proměnná kepra [čas]
- **liftr** (filtr)
- **liftrace** (filtrace) = **modifikace kepra**
(váhování, oříznutí, zkracování)
- krátko-dobý liftr (dolno-frekvenční filtr)
- dlouho-dobý liftr (horno-frekvenční filtr)
- **gamnitude** (magnituda, amplituda)
-

Vlastnosti reálného DFT kepra

$$c_r[n] = \text{IDFT}\{\ln |\text{DFT}\{x[n]\}|\}$$

Vlastnosti:

- DFT keprum - numerický výpočet (period. a symetr.)
- První část - hlavní informace o tvaru amplitudového spektra
tj. spektrum neperiodické složky signálu
tj. spektrální obálka (vyhlazené spektrum)



Vyhlažené spektrum: $\overline{|X[k]|} = e^{\text{DFT}\{c_n \cdot w_n\}}$

Výchozí veličiny: parametry AR modelu - a_k , $G = \sqrt{E_p}$

$$c_0 = \ln G$$

$$c_n = -a_n - \frac{1}{n} \sum_{k=1}^{n-1} (n-k)a_k c_{n-k}, \text{ pro } n = 1, 2, \dots, p,$$

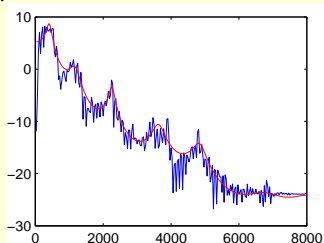
$$c_n = -\frac{1}{n} \sum_{k=1}^p (n-k)a_k c_{n-k}, \text{ pro } n = p+1, p+2, \dots$$

Vlastnosti:

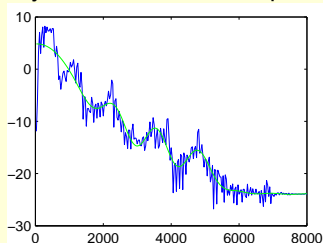
- Koeficienty Taylorova rozvoje $\ln |H(z)|$ (inverzní Z-transf.)
- Nekonečně dlouhé, první hodnoty opět nejvýznamnější
- Lze spočítat rekurentně, neobsahuje náhodnou složku
- Tvar spektra kopíruje LPC spektrum

Kepstrální analýza pro zpracování řeči

LPC spektrum:



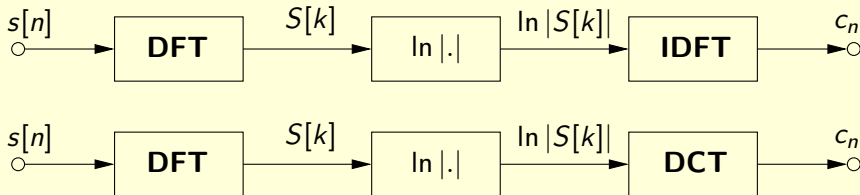
Vyhlazený odhad z reálného kepra:



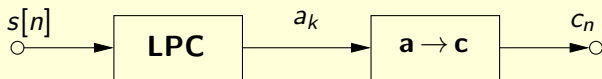
- První koeficienty nesou hlavní informaci o tvaru amplitudového spektra (12-20 keprálních koeficientů)
- **kepra podobných segmentů tvoří shluky**
⇒ použití jako příznaky pro rozpoznávání
- **keprum nese informaci o vokálním traktu**
⇒ použití pro automatickou hlasovou identifikaci

DFT a LPC keprstrum - bloková schémata výpočtu

Blokové schéma výpočtu keprstrálních koeficientů pomocí DFT



Blokové schéma výpočtu LPC keprstrálních koeficientů



Mel-kepstrum - melodická frekvenční stupnice

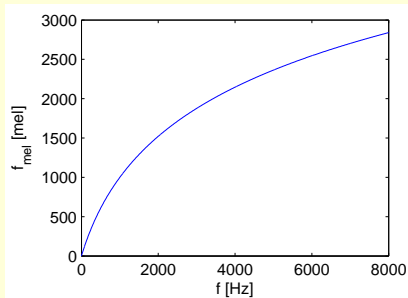
Vhodnější výpočet kepstra :

≈ **modelování nelinearity vnímání frekvence lidským sluchem**

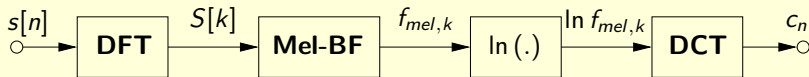
Nelineární zkreslení frekvenční osy - *melodická stupnice*

$$f_{mel} = \text{Mel}(f) = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$$

$$f = \text{InvMel}(f_{mel}) = 700 \cdot \left(10^{\frac{f_{mel}}{2595}} - 1 \right)$$

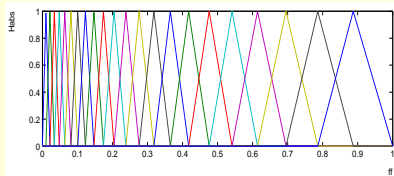


Blokové schéma výpočtu mel-kepstrálních koeficientů:



Výpočet energie v jednom pásmu

$$g_j = \ln \sum_{k=0}^{N/2} |S[k]|^2 H_{mel,j}[k].$$



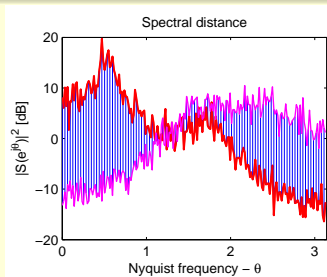
Výpočet kepstra pomocí DCT

$$c_i = \sqrt{\frac{2}{P}} \sum_{j=1}^P g_j \cos\left(\frac{\pi i}{P}(j - 0.5)\right)$$

- **nejrozšířenější příznaky používané v ASR**
- nyní hodně používané i pro **hlasovou identifikaci**
(v úlohách textově nezávislého rozpoznávání řečníka)

Spektrální vzdálenost (L_2 -norma)

$$L_2 = \int_{-\pi}^{\pi} \ln \frac{|S_1(e^{j\theta})|^2}{|S_2(e^{j\theta})|^2} d\theta$$

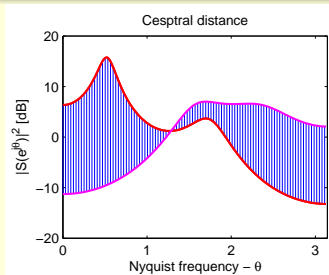
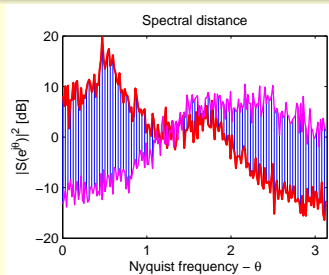


Spektrální vzdálenost na bázi L_2 -normy

→ kvantifikuje plochu ohraničnou dvěma spektry (křivkami)

Kepstrální vzdálenost

$$CD = \sqrt{(c_s[0] - c_x[0])^2 + 2 \sum_{k=1}^L (c_s[k] - c_x[k])^2}$$



- CD aproximuje spektrální vzdálenost na bázi L_2 -normy
- používá první keprální koeficienty
- vzdálenost je vypočítána ze spektrální obálky (tj. z vyhlazených spekter)

Různé definice kepstrální vzdálenosti

- Euklidovská vzdálenost:
$$CD = \sqrt{\sum_{k=0}^L (c_s[k] - c_x[k])^2}$$
- Euklidovská vzdálenost bez $c[0]$:
$$CD = \sqrt{\sum_{k=1}^L (c_s[k] - c_x[k])^2}$$
- kvadrát Euklidovské vzdálenosti bez $c[0]$:
$$CD = \sum_{k=1}^L (c_s[k] - c_x[k])^2$$
- vážená (liftrovaná) kepstrální vzdálenost:
$$CD = \sum_{k=1}^L (L_k c_s[k] - L_k c_x[k])^2$$

-
- Vždy kvantifikace rozdílů ve spektru
 - Varianty - různá citlivost a různé měřítka

● Forezní lingvistika a fonetika

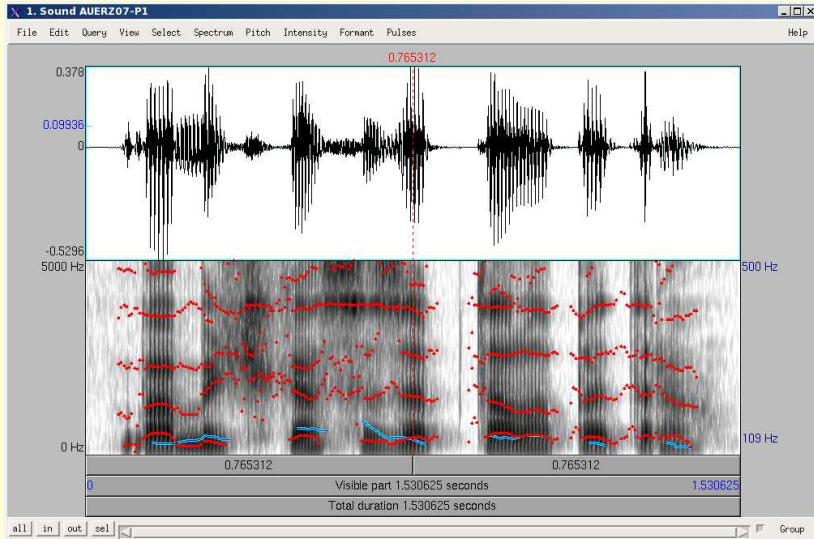
- sledování osobitých rysů projevu řečníka
- zaměření na artikulační zvláštnosti
- typické vedení melodie řeči (intonace)
- většinou na bázi poslechu

● Spektrografické metody

- Využívají možnost zobrazení diskutovaných hlasových charakteristik (spektrogramů, průběhu f_0 , trajektorií formantů, apod.)
- řešeno opět na expertní bázi
- detaily realizace vybraných hlásek

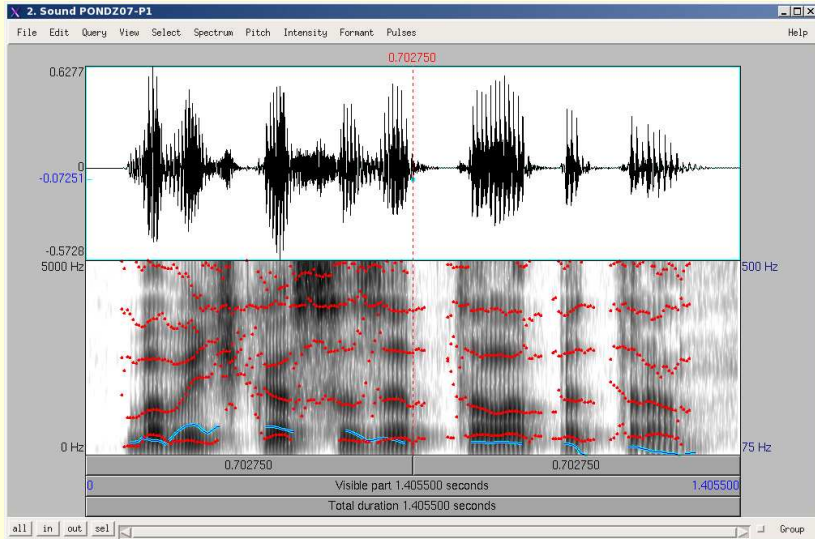
Formanty & základní tón - odhad Praat

Muž 1



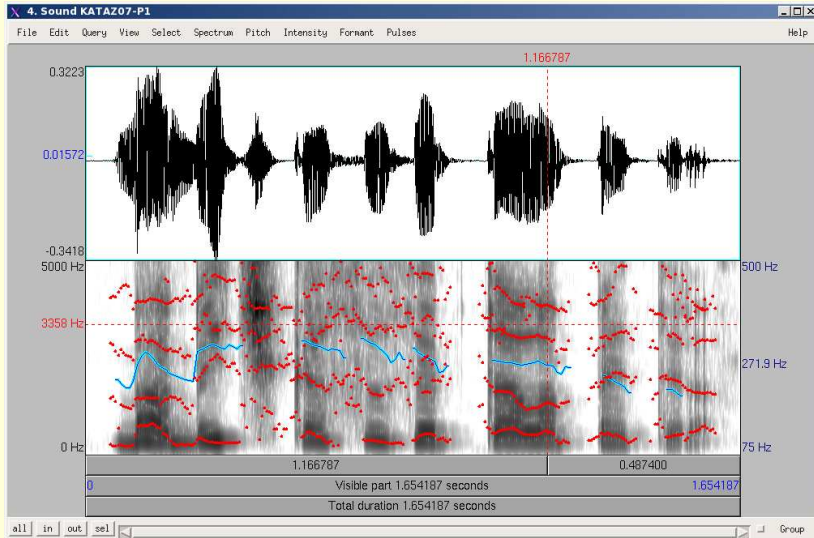
Formanty & základní tón - odhad Praat

Muž 2



Formanty & základní tón - odhad Praat

Žena 1



Děkuji vám za pozornost !