

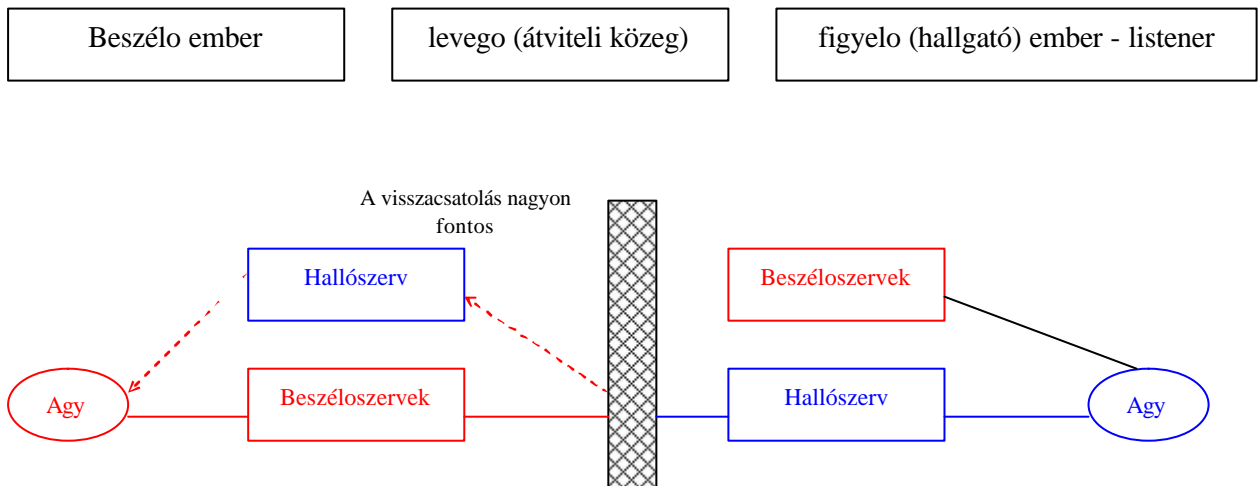
1. BEVEZETÉS

1.1. Alapfogalmak

1.1.1. Nyelv

- Minimum az emberi kommunikáció illetve az emberi gondolkodás legfontosabb eszköze.
- A nyelv elsődleges megnyilvánulási formája a beszéd (a beszéd az emberi kommunikációnak nem az egyetlen formája – nem verbális kommunikáció)
- természetes körülmények között az emberi kommunikáció alapvető jellemzője a multimodalitás, de a beszéd az egyetlen olyan kommunikációs eszköz, mely önmagában is érthető, ezért a beszédnek az emberi kommunikációban kimelt szerepe van

1.1.2. A természetes beszédlánc



1.1.3. Beszédfeldolgozás (beszédtechnológia)

A *beszédfeldolgozás* a természetes beszédlánc egy vagy több elemének mesterséges (gépi) feldolgozásával foglalkozik

Primer alkalmazások

- *beszédatvitel*: időbeni vagy térbeli távolságon keresztül és mindezt hatékonyan (sávszélesség éhség). A 30-as években került napirendre az az igény, hogy a beszéd sávszélességigényét úgy csökkentsék, hogy a felismerhetősége ne vagy csak alig romoljon. Később ennek módja a digitalizálás és tömörítés (MPEG)
- *beszédszintézis*: a beszéd mesterséges előállítása. Célja az informatikai folyamatok segítése.
- *beszédfelismerés*
- *beszélő azonosítás*: beléptető rendszereknél alkalmazzák, egy előre bementett mondat alapján azonosítják a beszélőt
beszélő felismerés: nagy adattárakban beszédmintákat tárolnak, a rendszernek ez alapján kell felismernie, hogy ki a beszélő, vagy esetleg nincs is rá vonatkozó információ az adatbázisban
- *beszédkorrektor*: például beszéd vizualizálása, hogy a süket ember is megtanulhasson beszélni

- beszédmanipuláció: pl megváltoztatni a beszéd sebességét úgy, hogy a hangmagasság (és hangsín) változatlan maradjon
- ember-gép kapcsolatok megváltozása: a kezelés és szemlélésen keresztül megvalósuló hagyományos ember-gép kapcsolatot felváltja a verbális ember-gép kapcsolat
- beszédinformációs rendszerek: a beszédfeldolgozás az információs rendszerek belső, inherens részévé válik (pl. bementett telefonszám alapján működő tudakozó)

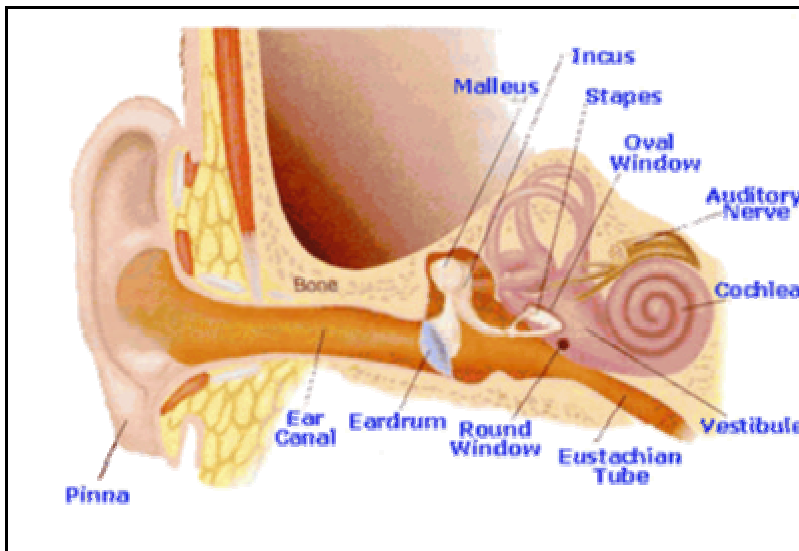
1.2. A hang fizikai leírása

- a hangot le lehet írni nyomással (p) illetve térfogatsebességgel (v)
- P_0 szinten lévő nyomás (1 atm) és ezen picurka longitudinális hullámok
- $P = P_0 + p(t)$
- $P_{eff} = \sqrt{\int (P_0 + p(t))^2 dt}$
- a hang a levegő nyomásváltozása, a levegőben longitudinális hullámként terjed (a térfogatrészek mozognak)
- normális viszonyok között: $\frac{p}{v} = 410 \frac{kg}{m^2 s}$, ha síkhullámnak tekinthető
- $c = f\lambda = 340m/s$
- a hang energiát visz magával (intenzitása van), $[I] = \frac{W}{m^2}$ területegységen időegység alatt áthaladó energia
- a hangnyomásszint (akusztikai decibel, intenzitásszint),

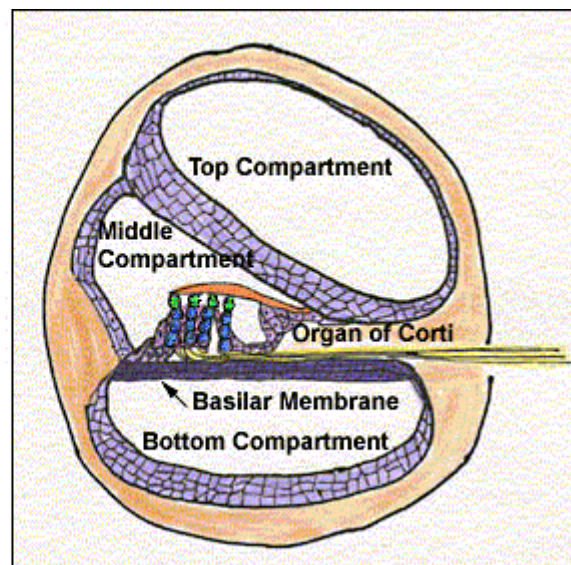
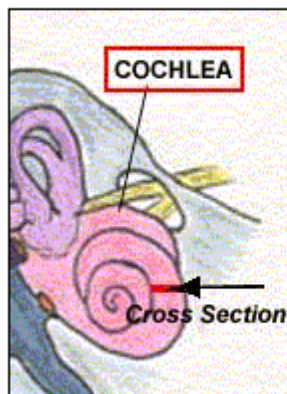
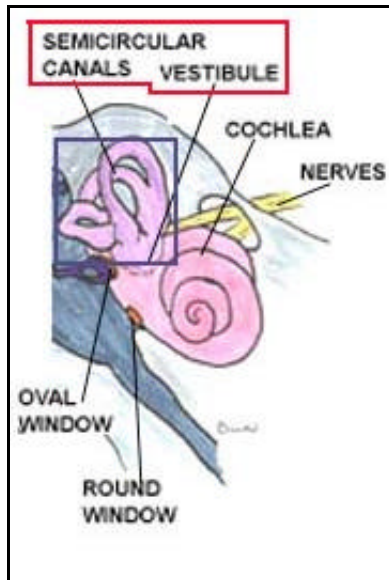
$$L = 20 \cdot \lg \frac{P}{20 \cdot 10^{-6} Pa} dB = 10 \cdot \lg \frac{I}{10^{-12} W} dB$$

1.3. Hallás

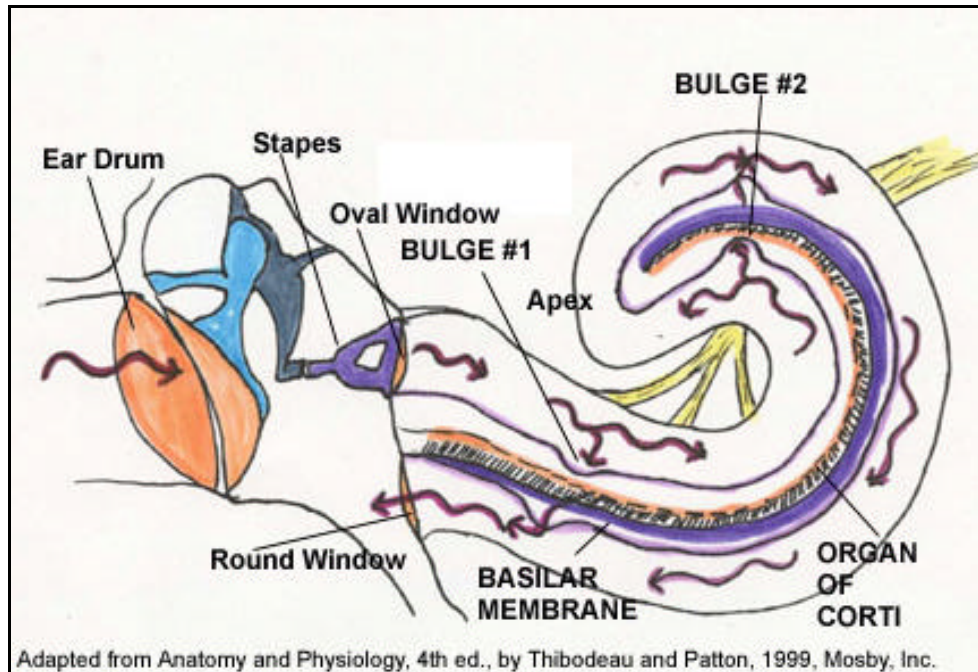
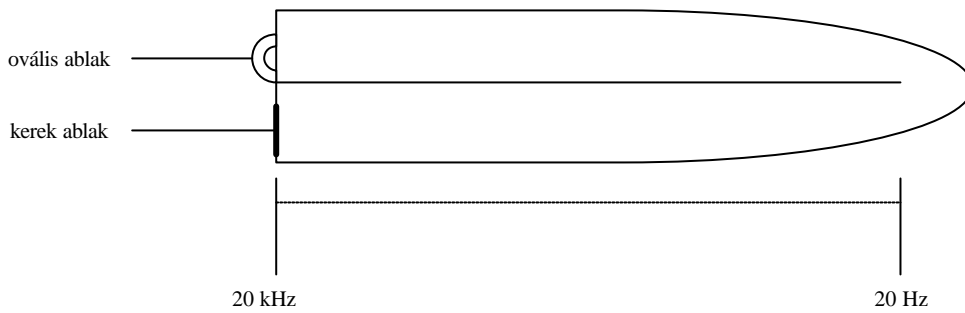
Hallószervek: fülkagyló, külső fülcsatorna (3000Hz rezonanciafrekvencia), dobhártya, halócsontok (kalapács, üllő, kengyel), belső fül



A belső fülhöz a hallócsontokhoz csatlakozik (kengyel, így veszi át a mechanikai hullámokat) egy tömlöcske (auditoria tube), amely folyadékkal van kitöltve. Ennek közepén van az alapmembrán, mely hang hatására rezgésbe jön. A membránon 3 sor ún. szorsejt (kb. 30000db) helyezkedik el, ezen sejtek végéhez idegek csatlakoznak, melyek közvetlenül az agyba mennek. Egy-egy ideg több ilyen szorsejttel is kapcsolatban lehet. A szorsejtek hozzáérnek a corti-szerv hártájához, így a mozgás hatására elektromos kisülések keletkeznek, amely az idegpályákon továbbterjed.



Az elektromos kisülés helye arra jellemző, hogy milyen frekvenciájú hangot hallottunk. Tiszta szinuszos hangoknál az ovális ablaktól való távolság számítható. A hang intenzitását az elektromos impulzusok sűrűsége (frekvenciája) mutatja. A hallás során az agyból is jönnek jelek – gátló jelek. Ezek teszik az ember frekvenciameghatározó képességét ilyen pontossá (a legjobban gerjesztett sejtek környezete blokkolódik).



A hallószervből jövő idegek nem közvetlenül mennek az agykéreg azon területére, ahol a hallásérzet keletkezik (mint pl. a szemidegeknél), hanem 5 központon keresztül. A hallásmechanizmus a legbonyolultabb az emberi érzékelések között, utnzására egyelőre semmi esély.

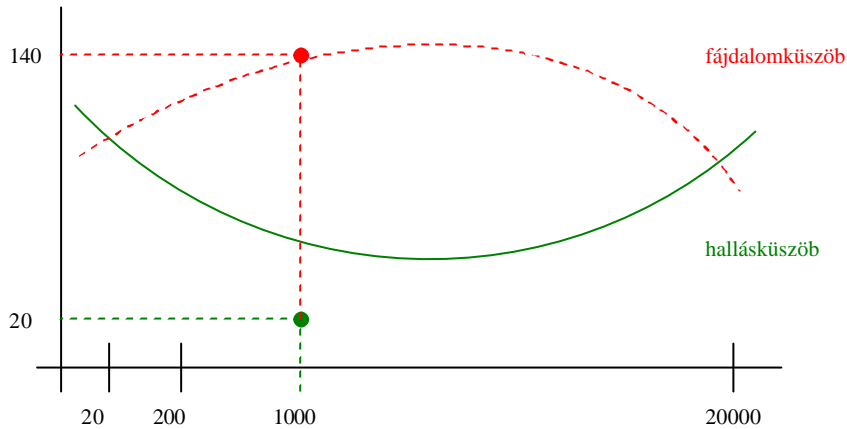
1.4. A hallás pszichofizikai (pszoakkusztikai) jellemzői

A hangjelenséggel kapcsolatos mérhető, fizikai mértékek és a hangérzet között nincs 1-1 értelmű megfeleltetés.

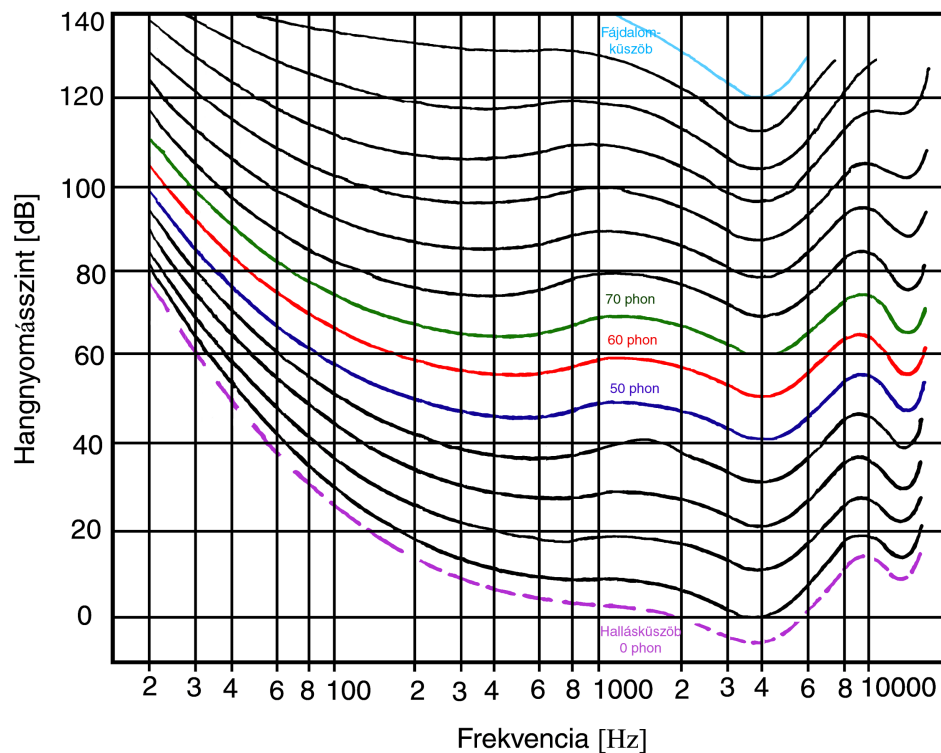
fizikai		hangérzet
intenzitás		hangosságérzet
spektrum		hangszín
frekvencia		hangmagasság

1.4.1. Azonos hangosságú (phon) görbék

Ezeket tiszta, szinuszos nagokra vizsgálják. Az n -phon az a görbe, amelyet az emberek statisztikailag azonos erőddégűnek hallanak és 1kHz-en n dB az erőssége.



- mértékegység a [phon]=hangosság szint
- a phon elég jól megadja a hangosságérzetet, de nem teljesen (pl. 30 phon mennyivel kevesebb, mint 40 phon?)



1.4.2. Hangosságérzet

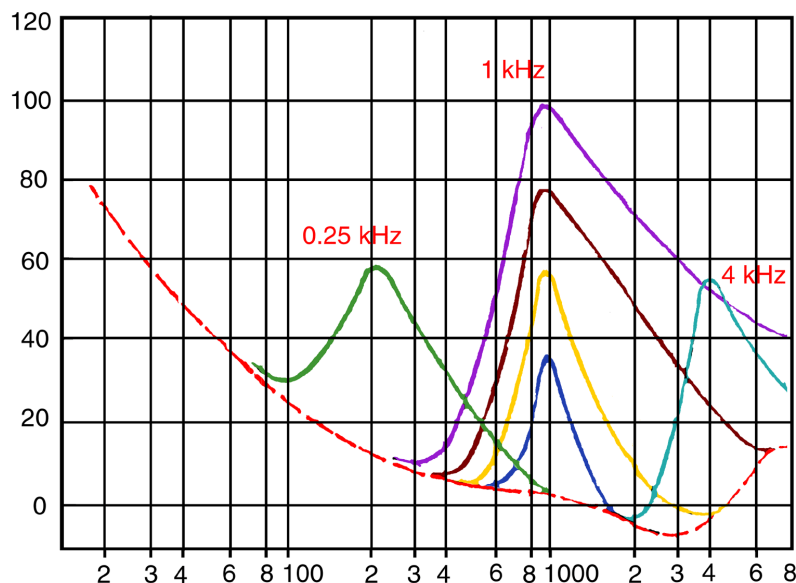
- általában son-ban adják meg
- megmutatja, hogy két phonban mért hangosságérzet aránya mekkor
- egységnek a 40 phon-t tekintik, az 50 phon 2 son, és így tovább

1.4.3. Kritikus sávok

- alkalmazzunk egy keskenysávú fehérzaj-gerjesztést (ennek intenzitása a görbe alatt lévő terület)
- a következő vizsgálo jel legyen szélesebb spektrumú, de ugyanolyan intenzitású
- bizonyos spektrumszélesség fölött a hangosságérzet no annak ellenére, hogy a kibocsátott zaj intenzitása nem változik
- **kritikus sávnak**¹ azt nevezzük, amelyen belül a hangosságérzet nem érzékeny a sáv szélességre.
- a kritikus sávokat kísérlettel szokták meghatározni
- bark (zwicker) skála: két frekvencia különbségét pszichoakusztikai szempontból megmutató skála; a különbség jellemzője, hogy hányszor lehet felmérni a kritikus sáv szélességet

1.4.4. Elfedés a frekvenciatartományban

- elfedő hangjelenség: 1000 Hz-en megszólaltatunk egy nagy intenzitású hangot, és a kritikus sáv szélességen belül szóló kisebb intenzitású hangokat a fülünk nem érzékeli
- pl. $f = 1000 \text{ Hz}$, $\Delta f = 160 \text{ Hz}$, $L = 80 \text{ dB}$ esetén a 1010 Hz-en 40 dB intenzitással szóló hangot nem érzékeljük
- tehát ha van elfedő hang, akkor a többi hang számára a hallásküszöb megemelkedik
- ezek az elfedési görbék alacsony frekvenciákon keskenyebbek, magasabb frekvenciákon pedig szélesebben elterülnek
- ezeket a tulajdonságokat hangtömörítésnél alkalmazzák elsősorban



^{1 1} Kitéző: A spektrálsűrűség

Tekintsük véletlenszerű függvényosság egy elemét (ilyen módon sokmindent érdemes modellezni, pl. emberi beszédet, zajt, stb.). Ezt a sokaságok sztochasztikus folyamatnak nevezzük. Jelöljük a sztochasztikus folyamatot \mathbf{X}_t -vel, ennek egy elemét $X_a(t)$ -vel, $S_x(f)$ -fel pedig ennek a sztochasztikus folyamatnak a Fourier transzformáltját. A $X_a(t)$ -t rávezetve egy keskenysávú szűrőre (f_a és f_b között) mérjük a folyamat teljesítményét. A $\mathbf{X}(t)$ -sztochasztikus folyamat spektrálsűrűsége megmutatja, hogy $f_a - f_b$ tartományban ennek a sztochasztikus folyamatnak mennyi a teljesítménye.

$$P(\{f_a, f_b\}) = \int_{f_a}^{f_b} S_x(f) df$$

1.4.5. Elfedés az idotartományban

- ha van egy nagyobb intenzitású hangjelenség, akkor ha ennél kicsit kisebb intenzitású megszólal, nem vesszük észre
- ha befejeződik a hangjelenség, a következő megjelenését nem azonnal vesszük észre, a fülnek van egy kis tehetetlensége – ez akár 150ms is lehet
- a fenti jelenség „visszafelé” is működik, csak sokkal kisebb időértékkel – 20ms

1.4.6. Írányérzékenység

- kismegjelenések az irányérzékenység a két fülbe érkező jel közötti fáziskülönbségből adódik
- nagyobb frekvenciákon az irányérzékenység az intenzitáskülönbségen alapul

1.4.7. Frekvencia – idő felbontóképesség

- Kérdés: milyen hibával (Δf) találjuk el a Δt ideig tartó hang frekvenciáját? $\Delta f := s$
- Harkevic és Gábor Dénes bebizonyították, hogy lineáris rendszerekben $\Delta f \cdot \Delta t \approx 0,01$.
- Ugyanezt a fül kb. 1000 Hz-ig jobban csinálja, noha különböző maszkolási jelenségekkel becsapható, azért nem annyira, mint a szem (pl. a szem számára RGB-ból gyakorlatilag minden szín kikeverhető)

1.5. A beszéd nyelvi szerkezete

- A hangot kétféle szinten vizsgáljuk:
 - ☞ akkusztikai szinten: valamilyen hanghullám
 - ☞ agyi szint: képesek vagyunk ezeket valamilyen diszkrét elemek (hangok) sorozatára bontani (ezen diszkrét elemek a beszédhangok, melyeket a hallás során az ember érzékelni képes)
- Lehetséges hierarchiaszintek: beszédhang – szótag – szó – mondatrész – ... – mondat – ... ezek nem mindegyike precízen meghatározott
- egy diszkrét beszédhang megszámlálhatatlanul végtelen sok alakú időfüggvényből absztrahálódik
 - ☞ intraindividuális: ugyanazon ember ugyanazta a hangot kétszer egymás után nem ugyanúgy mondja
 - ☞ interindividuális: két különböző ember ugyanazta a hangot nem ugyanúgy ejti ki
- A beszéd (nyelvek) úgy alakult(ak) ki, hogy ritkák az átlapolódások. Nem használja ki az összes lehetséges hullámformát (redundáns)
- **Artikulációs (akkusztikai) bázisnak** nevezzük a beszédalkeltés folyamán használt elemi folyamatokat, és **percepció bázisnak** hívjuk a beszéd megértése folyamán használt elemi folyamatokat.

1.5.1. Fonetikai megfontolások

- fonéma készlet – elemeknek olyan minimális számosságú halmaza, amelyből minden közlemény jelentéshelyesen, de csak egyféleképpen állítható elő agyi szinten, vagyis ha egy közleményben egy fonémát kicserélnénk, akkor megváltozna a közlemény értelme vagy elveszíti értelmét.
- allofonok – egyazon fonéma különböző akkusztikai megjelenései.
 - ☞ például a magyarban a nyílt és a zárt *e* megkülönböztethető (akkusztikailag), de egy fonéma
 - ☞ harang szóban az *ng*-t egy hangnak ejtjük, nem külön *n* és *g* egymásutánjaként, de ez a hang nem külön fonéma
 - ☞ összeállították a magyar nyelv fonémakészletét, ebben leggyakoribb az *e* (*eke*) és legritkább a *h ahhoz*).

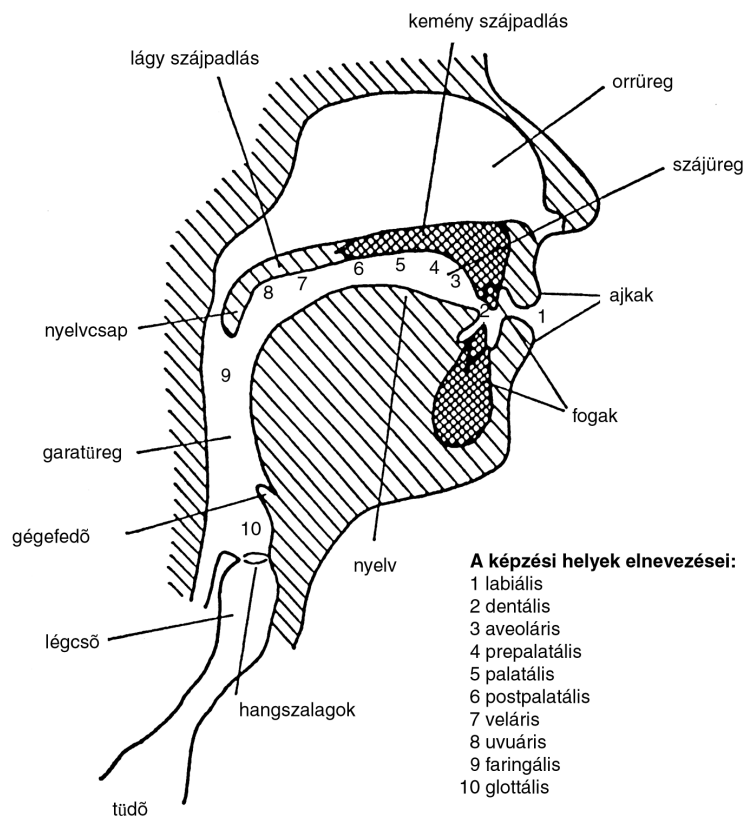
1.5.2. Az írás és a beszéd kapcsolata

- sok nyelvet karakterekkel írnak le, de léteznek olyan karakterek, amelyek szimbólumok (pl. 1, #)
 - ☞ ortografikus karakterek: azoka karakterek, amelyeket ki tudunk ejteni
 - ☞ graféma: pl. egészség szóban az *sz* hang egy graféma
- szöveg leírásakor ortografikus karaktereket használunk, de ha a jelentésig el akarunk jutni, akkor a graféma-konverziót meg kell tenni. Ezek már fonémákkal írhatók le.
- fonetikus leírások: a hangzást is megpróbáljuk leírni
 - ☞ IPA átírás – ASCII karakterekkel nem írhatók le
 - ☞ SAMPA – a 7 bites ASCII karakterekkel leírhatók
- a hangok átmenetekkel kapcsolódnak egymáshoz, természetes átkötés (ha nem stimmel, érzékeljük)

1.6. A hangképzés

A *tüdő* energiát szolgáltat, levegot pumpál a *légcsobe*. A gégeben elhelyezkedő *hangszalagokat* a hangképzés során vagy használjuk, vagy nem. A levego a gégeből a *garatüregen* keresztül a *száj-* illetve az *orrüregbe* jut. A *szájüregben* a *nyelv*, *fogak*, *ajkak* is részt vehetnek a hangképzésben. A hangszalagoktól felelő sor rész, mely részt vesz a hangképzésben az ún. *vokális traktus* (üregrendszer, toldalékcső).

A hangképző szervek vázlatos képe



Alapvetően három emberi hangkeltési mechanizmus van.¹

¹¹ Léteznek más hangképzési mechanizmusok is (pl. csettintés Afrikai törzseknél), de mi nem használjuk őket. Valószínűleg a hangképzési mechanizmusok és maguknak a hangoknak is a nyelvben lévő előfordulása attól függ, hogy az adott nyelv milyen körülmények között alakult ki (zajviszonyok, stb.)

1.6.1. Zöngé

- a hangszalagok tulajdonképpen két rostos, izmos hártya, melyek összeérnek, az izmok pedig a gégeben lévő porcokhoz tapadnak
- ha az izmok megfeszülnek, akkor a levegot nem engedik ki (féldobok)
- a levego kiáralik a tüdobol, és ha a hangszalagokat az izmokkal működésbe hozzuk, akkor azok előtt kis túlnyomás alakul ki (kb. 3-4cm vízoszlop nyomásának megfelelo), melynek hatására a hártya szétfeszül. Emiatt azonban a túlnyomás lecsökken és a hártya ismét visszazár. Így egy kvázi-periodikus jelenség alakul ki. Ennek a kvázi-periodikus jelenségnek a periódusideje férfiak esetében 8-12 ms, nőkél 4-7 ms.

1.6.2. Turbulens áramlás

- szignifikáns szükületeket képezünk a vokális traktusban (pl. az *f*, *s* hangok képzésénél)
- a levegorészecskék a szükületet elhagyva véletlenszeruen leszakadoznak
- az, hogy a hang hogyan hangzik attól is függ, hogy hol a résképzés helye és mi van után

1.6.3. Lökéshullám

- miközben a levegot préseljük ki, a vokális traktusban zárat képezünk (nem a hangszalagokkal) – ezért egy ilyen lökéshullámot 3 szakaszra tudunk bontani
- néma szakasz (néma fázis) – ez a zár képzésének ideje, ilyenkor a levego nem tud továbbhaladni
 - ☞ zárfelpattanás – a zár felpattan, és ennek hatására valamilyen hang jön létre (pl. *p* hangnál)
 - ☞ elhalás (átmenet)
 - ☞ látszólag véletlenszeruen jön, amplitúdóban nagyobb, mint a turbulens áramlás

1.7. A beszédhangok osztályozása

magánhangzók: a, á, e, é, i, o, ö, u, ü

mássalhangzók

nazálisok: m,n,ng, ny

likvidák: l, j

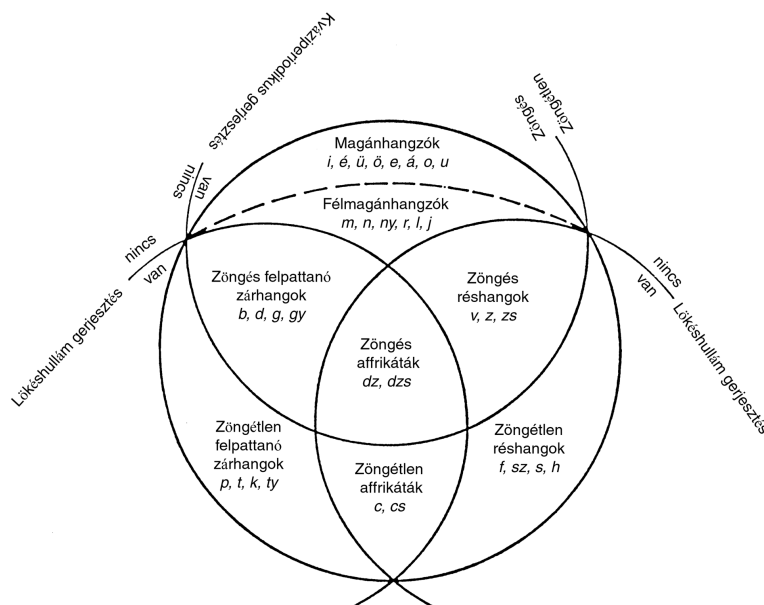
pergo: r

felpattanó zárhangok (explozívák, plozívák, stops): t, d, p, b, k, g, ty, gy

részhangok (frikatívák): f, v, s, zs, sz, z, h

zár-részhangok (affrikáták): c, dz, cs, dzs

A magyar beszédhangok osztályozása a fonetikai és a gépi feldolgozás alapján



1.8. Spektrális módszerek

1.8.1. Periodikus jelek – Fourier sor

$$f(t) = f(t + kT_0), \quad k \in \mathbb{Z}$$

$$f(t) = c_0 + \sum_{n=1}^{\infty} c_n \cos(n\Omega_0 t + \varphi_n), \quad \Omega_0 = \frac{2\pi}{T_0}$$

- egy jel spektrumán azt értjük, hogy adott frekvencián milyen amplitúdójú és fázisú az adott jel. Tehát a spektrumot a $\{n\omega_0, c_n, \varphi_n\}$ hármas határozza meg.
- az időben periodikus jelek vonalas spektrumúak
- a Fourier sor komplex alakja:

$$\cos(x) = \frac{e^{jx} + e^{-jx}}{2} \Rightarrow f(t) = c_0 + \sum_{n=1}^{\infty} \left(\frac{c_n}{2} e^{j\varphi_n} e^{jn\Omega_0 t} + \frac{c_n}{2} e^{-j\varphi_n} e^{-jn\Omega_0 t} \right)$$

bevezetve az alábbi jelöléseket:

$$C_0 = c_0, \quad C_n = \frac{c_n}{2} e^{j\varphi_n} \quad \text{és} \quad C_n^* = C_{-n} = \frac{c_n}{2} e^{-j\varphi_n}$$

$$f(t) = \sum_{n=-\infty}^{\infty} C_n e^{jn\Omega_0 t}, \quad \text{ahol} \quad C_n = \int_{t_1}^{t_1+T_0} f(t) \cdot e^{-jn\Omega_0 t} dt$$

- egy általános vonalas spektrumú jel Fourier-sora

$$f(t) = \sum_{n=-\infty}^{\infty} D_n e^{j\Omega_n t}, \quad \text{ahol} \quad \frac{\Omega_n}{\Omega_m} \text{ irracionális is lehet.}$$

1.8.2. Egyszeri folyamat – Fourier-integrál

$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt = F\{f(t)\}$$

- Jelentése: kontinuum sok szinuszos hullám összege. Ez egy komplex függvény abszolútértékkel és arkusszal.
- az időfüggvényt az inverz-Fourier transzformációval állíthatjuk elő

$$f(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} F(\omega) e^{j\omega t} d\omega$$

- fontos paraméter még a jel fajlagos energiája és teljesítménye:

$$E = \int_{-\infty}^{\infty} f^2(t) dt \quad \text{illetve} \quad P = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-T/2}^{T/2} f^2(t) dt$$

- azért fajlagos, mert megfelelő konstanssal szorozva energiát illetve teljesítményt kapunk
- periodikus jelek esetében $E = \infty$ és P véges, egyszeri folyamatoknál E véges és $P = 0$.

1.8.3. Lineáris rendszerek hatása az átvitt jelre

- a mínusz végtelentől végtelenig való integrálás probléma, hiszen akkor meg kell várni, amíg a teljes folyamat lejátszódik, erre viszont nincs mindig lehetőségünk
- a teljes folyamatot kiablakozzuk: az időfüggvényt megszorozzuk a t_1 időpillanatokra eltolt ablakfüggvénnyel

$F_w(\omega, t_1) = F\{f(t) \cdot a(t, t_1)\}$, kérdés, hogy mennyire rontja el az ablakolás a spektrumot

$$F_w(\omega, t) = F(\omega) \cdot A(\omega, t_1)$$

- legjobban az ún. szeretjük *Hamming-ablakot*, mert ennek spektrumában a főmaximum és a második maximum között 50 dB erősítéskülönbség van, tehát a távoli frekvenciákat a konvolúcióban ez az ablakolás gyakorlatilag nem veszi figyelembe.
- Hamming-ablak; $a(t) = 0.54 - 0.46 \cdot \cos 2\pi \frac{t}{T_a}$
- digitális jelfeldolgozás során DFT-ket (Discrete Fourier Transformat alkalmazunk, általában a DFT-nek egy gyorsan elvégezhető módszerét alkalmazzuk, az ún. FFT-t (Fast Fourier Transform)
- spektrogram: gördülő spektrumot tekintve (idő – frekvencia sík) ahol a spektrum értéke nagy, ott erosen befeketítjük,
 - ☞ a spektrális viszonyok az idő függvényében változnak
 - ☞ az 1900-as évek közepén az ún. *szonográfot* alkalmazták, amely ezt a gördülő spektrumot közelítette (a közelített ábra a *szonogram*, *spektrogram*)

1.9. Beszédhangok finom szerkezete

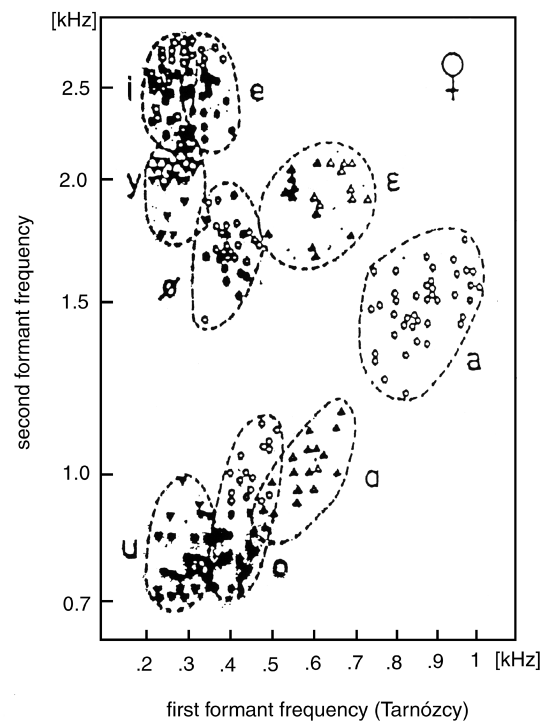
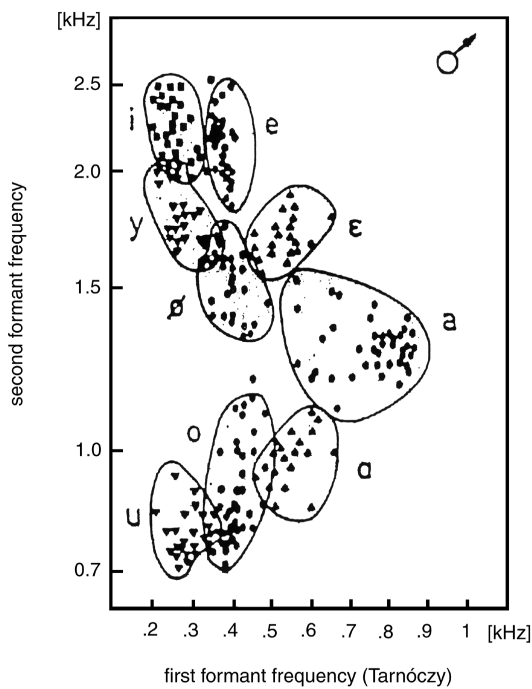
1.9.1. Hosszan tartható zöngés hangok

- olyan hangok, amelyekben a zöngén kívül nincs más hangkeltési mechanizmus
- időtartományban közel háromszög alakú térfogatsebességgel ábrázoljuk
- a hangképzés során keletkező hang leíró függvénye kvázi-periodikus, ezért Fourier-sorba fejthető (gazdag felharmonikus tárral rendelkeznek)
- pl. a beszédet a telefon 300-3400 Hz között viszi át
- ha az alapharmonikus (f_0) nincs benne az átvitt jelben, a fül akkor is képes azt kiejteni, mert a felharmonikus tár igen gazdag
- a vokális traktuson keresztül a hang a száj illetve orrüregben keresztül távozik. Az orrüreg minden hangra hatást gyakorol direkt illetve indirekt módon. Ezt úgy modellezhetjük, hogy
 - ☞ a vokális traktusnak van egy átviteli függvénye: $V(\omega)$
 - ☞ létezik egy ún. sugárzás (radiációs) ellenállás: $R(\omega)$, ahogy a szájból eltávozik a hang, a magas frekvenciákat „lennyomja”
 - ☞ $R(\omega) \cdot V(\omega)$ a zöngé által keltett hangot *formálja*
 - ☞ $P(\omega)$ – ezen hang spektrális megjelenése
- a vonalas spektrumra illeszthető burkológörbe maximumhelyeit *formánsoknak* nevezzük (F_1, F_2, \dots, F_n formáns frekvenciák)
- a hosszan tartható zöngés hangok formáns struktúrával rendelkeznek
- hogyan értelmezzük a burkológörbét: minden spektrumvonalra ültessünk egy $\frac{\sin x}{x}$ függvényt.

Az adott helyen a függvény értéke legyen a spektrumvonal magassága és a hullámátmenetek F_0 távolságban legyenek (ahol F_0 az alapharmonikus frekvenciája).

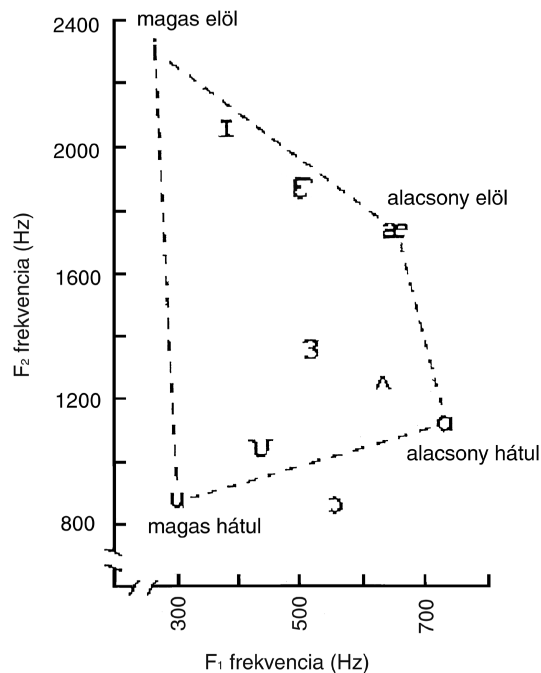
1.9.1.1. Magánhangzók csoportja

- a magánhangzókat két formáns közel, három formáns szinte teljes biztonsággal megkülönbözteti egymástól, érzeti szempontból leglényegesebbek a formánsfrekvenciák
- az $A_1 \dots A_n$ formánsok „amplitúdóit” az elsore (A_1) szokták normálni
- a helyi maximumok alatt 3 dB-lel meghúzott vonal és a burkológörbe metszéspontja jelöli ki $B_1, B_2, \dots B_n$ formánsok sávszélességét.
- a magánhangzókat F_1 - F_2 síkon szoktuk ábrázolni
- a nők F_1 -ben és F_2 -ben is magasabb területre kerülnek, mint a férfiak, további jellegzetes különbsége nők és a férfiak beszéde között, hogy a vonalas spektrum léceit a férfiaknál 100 Hz, a nőknél 200 Hz választja el egymástól (ezért pl. a nők magas hangon nem tudják az *u* hangot kiénekelni)
- további érdekesség, hogy a magánhangzók magasságát F_2 és nem pedig F_1 határozza meg¹
- a magánhangzók hossza kb. 30-60 ms



¹ Tehát pl. az *uhang* F_1 frekvenciája nagyobb, mint ugyanazon hangmagasságban lévo *i-é*, de az F_2 már az *i-nél* magasabb.

Angol magánhangzók F_1 , F_2 grafikonja



1.9.1.2. Nazálisok

- 250-300 Hz környékén van az F_1 , de magasabb frekvenciákon nincsenek formánsok
- sokkal kisebb az energiájuk, mint a magánhangzóknak

1.9.1.3. Likvidák

- van formáns struktúrájuk, de nem jellegzetes
- nagy szabadságfokú, hogy hol képződik
- sokkal kisebb az energiájuk, mint a magánhangzóknak

1.9.2. Felpattanó zárhangok

- általános szerkezetük: néma fázis (60-120 ms), folytott zöng/zöngétlen (30 ms), zárfelpattanás (40 ms) hallható aspiráció/átvezetés a következő hangra
- a hang hosszítása a néma fázis hosszabbításával történik

1.9.3. Réshangok

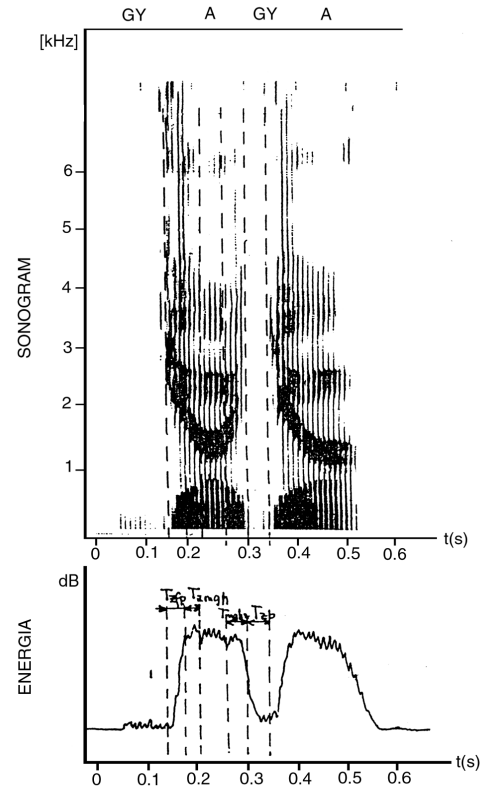
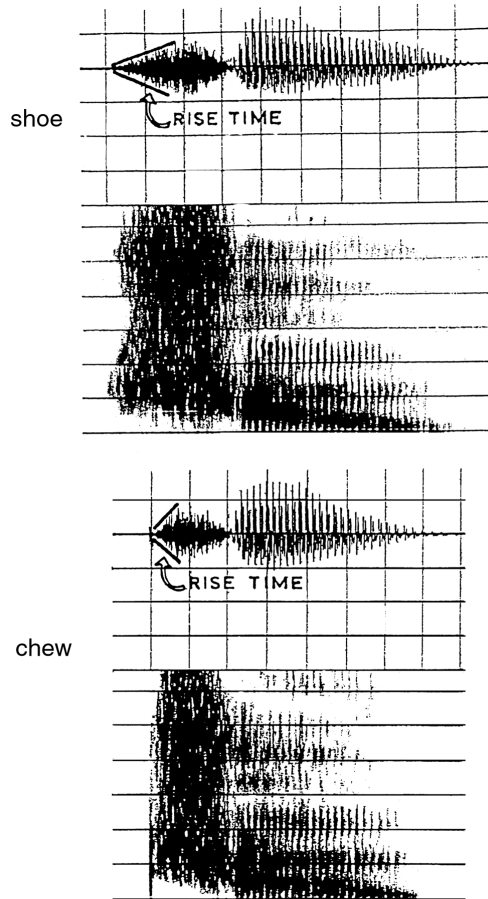
- lényeges, hogy frekvenciában 3,5 kHz felett van aspektrumuk (ezért pl. telefonban ezeket a hangokat nem tudjuk megkülönböztetni)
- ha a hang zöngés, akkor erre a spektrumra szuperponálódik rá egy formáns szerkezetű vonalas spektrum

1.9.4. Zár-rés hangok (affrikáták)

- általános szerkezetük: néma fázis, folytott zöng/zöngétlen, zárfelpattanás, réshang
- fontos az idoszerkeze: a néma fázis az zár-rés hangok előtt rövidebb, mint a zárhangoknál és a réshang kialakulása sokkal gyorsabb (10-15 ms), mint a tiszta réshangé, továbbá a tiszta réshang tartási ideje sokkal kisebb (rövidebb)

1.9.5. Hangámenetek

- az ember a beszéloszerveit nem tudja ugrásszerűen változtatni, ezért beszéd közben átmeneteket képez a hangok között (végtelen sok félet)
- az átkötéseket alapvetően a szomszédok határozzák meg, bizonyos esetekben az 1, 2 –vel után elhelyezkedő hangok is hatással vannak a kiejtett hangra
- tipikus vizsgált kombinációk: cv, vc, cvc, vcv¹



1.9.5.1. Hosszan tartható hangok környezete

- (szomszéd) – (átmenet – tiszta fázis – átmenet) – szomszéd
 - ☞ átmenet: látszólag összevissza függvény, egyre jobban hasonlít a tiszta fázis jelalakjához
 - ☞ tiszt fázis: kvázi-periodikus jel
 - ☞ átmenet: többnyire amplitúdó csökkenés történik, de a legvégén itt is vannak torzulások, mint a kezdő átmenetnél

1.9.5.2. Locus

- a cv átmenet jellegzetessége a locus: megfigyelték, hogy pl. a d után ejtett magánhangzók felfutó szakaszait, ha visszafelé meghosszabbítjuk, ezek egy pontban metszik egymást – a legtöbb mássalhangzó az ot követő magánhangzó vagy ot megelőző magánhangzó második formánsát a szóban forgó mássalhangzót jellemző frekvenciára kényszeríti, ezek a locusok. Legjellemzőbb locusok a felpattanó zárhangoknál vannak, ezek elég jól jellemezhetők a locusaikkal.

¹ c: consonant (mássalhangzó), v: vowel (magánhangzó)

1.10. A folyamatos beszéd akkusztikai szerkezete

1.10.1. Hangsúly

- a *hangsúly* az, hogy a mondanivaló függvényében ugyanazt a beszédhangot különbözőképpen ejtjük
- *intonáció*: $F_0(t)$, a zöngé alapfrekvenciájának változása az időben
 - ☞ ezt szekundumos nagyságrendben érzékeljük¹

1.10.2. Intenzitás

- rövid idejű energia (a jel által hordozott energia egy ún. kiablakolt energia)
- ugyanazt a magánhangzót egy kérdés mondatban más intenziással ejtjük, mint egy kijelentő mondatban

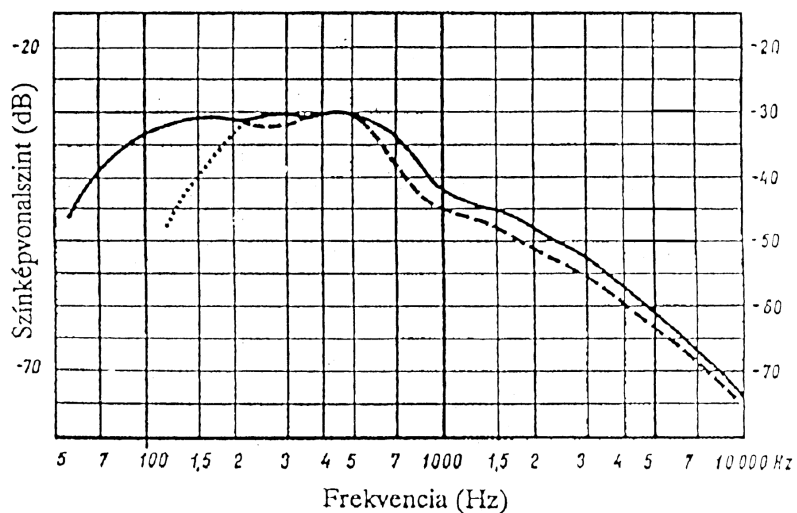
1.10.3. Ritmus

- a különböző hangokat különböző hosszúsággal ejtjük
- a magánhangzók tiszta fázisainak hosszával szabályozzuk

1.10.4. Statisztikai jellemzők

- stacionárius sztochasztikus folyamatként tekintve a beszédet
- tetszőleges helyen kiválasztott valószínűségi változó sűrűségfüggvényével tökéletesen jellemezhető egy stacionárius sztochasztikus folyamat, ez az ún. *amplitúdóeloszlás*
 - ☞ normáljuk ezt a sűrűségfüggvényt a saját szórására
 - ☞ a beszéd sokkal jobban feszíti az eszközök amplitúdótartományát
 - ☞ a beszédben kb. 30 dB dinamika-különbség és még ennyi hangerő
 - ☞ σ szórás a hangerőre jellemző
- ugyanezt a folyamatot a frekvenciatartományban vizsgálva a spektrális sűrűség írja le jól
 - ☞ $10 \lg \left(\frac{s(f)}{s(f_r)} \right)$, ahol $s(f)$ a spektrális sűrűség a kérdéses frekvencián, $s(f_r)$ pedig a spektrális sűrűség a referenciakézfrequencián

¹ Létezik mikrointonáció is (ez egy hangon belüli alapfrekvencia-változás), ettől lesz személyre jellemző, természetes hangzás



62 magyar beszélő átlagos színeképe. A folytonos vonal a férfiak és nők burkológörbéjét jelzi. A női beszédészíneképe ez alatt marad, pontozott vonal jelzi az eltérést. Ahol pedig a férfi beszédészíneképe marad alatta, szaggatott vonal jelzi az eltérést.

2. BESZÉDKÓDOLÁS ÉS TÖMÖRÍTÉS

2.1. Bevezetés

- a beszéd folyamata analóg jel: időben folytonos, értékészlete folytonos
- a beszédfeldolgozásban digitális eszközökkel dolgozunk, tehát a beszédet bitek sorozatává kell alakítanunk, ennek módja: mintavételezés → kvantálás → kódolás¹
- analóg jelsorozat bitsorozattá való étalakításakor az alábbi jellemzőkre kell figyelni
 - ☞ huség: valamelyem matematikai mértékkel mérhető (pl. négyzetes hiba) vagy pszichoakusztai mérések eredménye alapján
 - ☞ sebesség
 - ☞ komplexitás: a megvalósíthatóság szempontjából érdekes
- ezek a paraméterek egymás ellen dolgoznak, u, a tömörítés ugyan a sebességre irányítja a figyelmet, de mindhárom egyaránt fontos.

2.2. Mintavételezés

Shannon-Kotvelnyikov: Ha egy jel B sávra korlátozott, azaz a jelben adott B frekvencia fölött nem fordul elő komponens, akkor $f_0 \geq 2B$ suruséggel vett mintáiból a jel egyértelműen visszaállítható.

Visszaállításhoz $\frac{\sin x}{x}$ interpoláló függvényt alkalmazunk minden mintavételi pontban az adott minta értéke és a többi mintavételi pontban pedig 0. Probléma: ez egy matematikai függvény, nem megvalósítható.

2.2.1. PAM² típusú simító visszaállítás

- $x(t)$ jelet $f_0 = 1/T_0$ frekvenciával mintavételezzük, és így előáll $\{x_i\}$ jel. Legyen az elemi jelkelto $m(t)$ súlyfüggvényu, kimenete pedig változzék az x_i -vel arányosan: $x_i m(t - iT_0)$. Vegyünk továbbá egy $g(t)$ súlyfüggvényu szurot, amelyre ráeresztve az elemi jelkelto által generált jelet, egy $\tilde{x}(t)$ függvényt kapunk. Az átvitel akkor huséges, ha $x(t) = \tilde{x}(t)$.
- Legyen $M(f)$ az $m(t)$ és $G(f)$ a $g(t)$ Fourier transzformáltja.

Tétel: Egy B frekvenciasávra korlátozott jel $T_0 = 1/f_0$ idoközu mintáiból akkor állítható PAM típusú simító visszaállítással helyre, ha $f_0 > 2B$ és $H(f) = M(f)G(f)$ $-B$ és B között konstans³, a mintavételi frekvenciák B sugarú környezetében 0, egyébként pedig tetszoleges.

a Shannon-tétel ennek egy speciális esete, amikor (B és f_0 között kijelölve egy pontot: x) $-x$ és x között $H(f)$ értéke 1, azon kívül pedig 0. Az ilyen szögletes spektrum idotartománybeli megfeleloje:

$$h(t) = \frac{\sin x}{x} = m(t) * g(t) = \int_{-\infty}^{\infty} m(\mathbf{t})g(t - \mathbf{t})d\mathbf{t}$$

- gyakorlati szempontból: megvalósíthatósági megfontolások – a gyakorlatban impulzusokat állítunk elő, melyeknek spektrum egy lefelé görbülo $\frac{\sin x}{x}$ -hez hasonló jelalak. Ez az $M(f)$. Ahhoz, hogy $H(f)$ a tételnek megfelelo legyen, $-B$ és B között $G(f)$ -nek felfelé kell görbülnie, majd $f_0 - B$ -nél már közel 0⁴ értéket kell képviselnie. (Általában $f_0 - B$ -ben a szuronek pólusa van)

¹ A kódolás nem feltétlenül egy mintát vesz alapul, lehetséges kódolás minták egy véges halmazára is.

² Pulse Amplitude Modulation

³ Ha az elemi jelkelto elé egy T_0 -as szorzót teszünk, akkor ez a konstans éppen 1 lesz.

⁴ Frekvenciában nem tudunk tartósan zérus átvitelt biztosítani: 40 – 60 – 80 dB-nyit.

- további szempont, hogy a fáziskarakterisztika lineáris legyen, ami azt is jelenti, hogy a csoportfutási idő konstans. Meredek lefutások előtt viszont a csoportfutási időnek mindig csúcsa van
- Manapság már inkább $G(f)$ -et tervezik és ehhez alakítják $m(t)$ -t.

2.2.2. Azonos mintájú jelek

- Vegyünk azonos mintájú, de különböző időfüggvényű jeleket: $x_i(kT_0) = x_j(kT_0)$, $\forall k$. Van-e spektrális rokonság az azonos minták között?

Tétel: az azonos mintájú jelek halmozott spektrumai megegyeznek, azaz

$$\sum_{l=-\infty}^{\infty} x_i(f - lf_0) = \sum_{k=-\infty}^{\infty} x_j(f - kf_0)$$

ahol $\sum_{l=-\infty}^{\infty} x_i(f - lf_0)$ az $x_i(t)$ halmozott spektruma.

- A halmozott spektrum a képzés módja miatt periodikus. Az $-\frac{f_0}{2}, \frac{f_0}{2}$ közötti szakaszt nevezzük Nyquist-ekvivalensnek. A Nyquist-ekvivalens az azonos mintájú jelek közül a legkisebb sávszélességű. Ha a jel eredetileg teljesíti a mintavételi tétel frekvencia-feltételét, akkor a Nyquist-ekvivalense önmaga.
- A minták egyértelműen meghatározzák a halmozott spektrumot és a mintákra halmozott spektrumot határozzák meg egyértelműen.
- A halmozott spektrum tulajdonképpen egy Fourier sor a frekvenciatartományon (innen is látszik, hogy ez periodikus jel):

$$\sum_{l=-\infty}^{\infty} x_i(f - lf_0) = T_0 \sum_{n=-\infty}^{\infty} x(nT_0) e^{-j2\pi f n T_0}$$

2.2.3. PAM típusú simító visszaállító kimeneti jele

$$X(f) = T_0 \left[\sum_{n=-\infty}^{\infty} x(nT_0) e^{-j2\pi f n T_0} \right] M(f)G(f) = \left[\sum_{l=-\infty}^{\infty} X(f - lf_0) \right] M(f)G(f)$$

és $\sum X(f - kf_0)$ azon jelnek halmozott spektruma, amelyből a minták származnak.

- A gyakorlati mintavételezés során $x(t)$ halmozott spektruma $H(f)$ -fel lép interakcióba.
- Helyes spektrumfeltétel esetén a halmozott spektrumban nincsenek átlapolódások, a mintavett jel spektruma megegyezik a Nyquist-ekvivalensével, a visszaállítás során így éppen az eredeti jelet kapjuk.
- Ha a spektrumfeltételt nem tartjuk be, akkor a mintavett spektrumok a halmozott spektrum képzése során összegződnek, ebből az összegzett spektrumból az eredeti nem állítható elő hűségesen. Ilyen esetben törekedhetünk arra, hogy levágjuk azokat az széleket, ahol a kicsúcsosodások vannak, illetve a jelet elotte is sávkorlátozhatjuk.
- Idegen átlapolódás jelensége: ha a jelünk ugyan sávkorlátozott és betartja a frekvencia-feltételt, de a nagyfrekvenciás tartományokban van valami idegen jel, akkor ez a halmozás folyamán megjelenik a hasznos sávban és a jelből többé nem távolítható el. Ezen hibák kiküszöbölésére bemeneti aluláteresztő szűrőt alkalmazunk, amely ugyan levágja a jel egy részét, viszont így a jel oly módon lesz sávkorlátozott, hogy betartja a frekvencia-feltételt, és ha a bemeneti szűrő nem ideális voltából adódó torzításokat $H(f)$ -fel kikompenzáljuk, akkor a nem levágott részeket torzításmentesen tudjuk átvinni. Tehát torzítás csak az átlapolódás megszüntetése miatti sávkorlátozásból adódik.

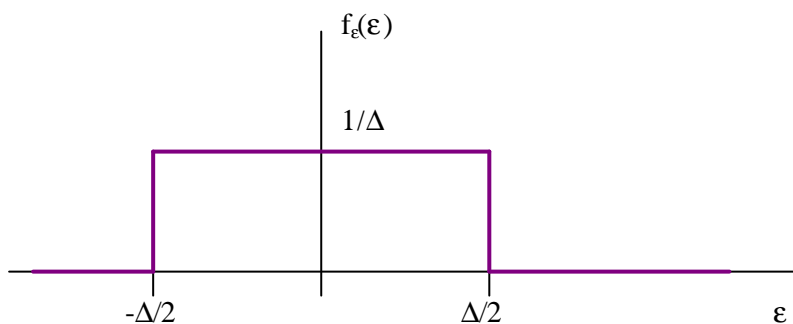
2.3. (Lineáris) kvantálás

a kvantálás általában a jelet nem időben, hanem amplitúdó értékészletben teszi diszkrétte

- többnyire mintavett jelre alkalmazzuk
- véges amplitúdó értékészlet:
 - ☞ kvantálási szintek: a megengedett szintek
 - ☞ kvantálási lépcső: a megengedett szintek közötti lépcsők
- lineáris kvantálás esetén $\Delta_i = \Delta \forall i$ -re, és a kvantálási szabály az, hogy minden értéket a hozzá legközelebb eső kvantálási szinttel helyettesítjük.¹

2.3.1. A kvantálás ára: a kvantálási zaj

Legyen \mathbf{e} a kvantálási hiba értéke, \hat{x} a kvantált és x a kvantálandó érték. Ekkor $\hat{x} = x + \mathbf{e}$. Ésszerű modellt választva ε -ra, a egy valószínűségi változó $f_\varepsilon(\varepsilon)$ sűrűségfüggvénnyel.



A kvantálási hibát, mint számsorozatot tekintve: ε_i és ε_j időben egymás után megjelenő valószínűségi változók. Ezek az ε értékek egymástól függetlenek, köztük korreláció nincs: $M(\varepsilon_i \varepsilon_j) = M(\varepsilon_i) \cdot M(\varepsilon_j)$ $\forall i \neq j$ esetén. Mivel az eloszlás szimmetrikus a 0-ra, ezért a várható érték 0.

A mintasorozat visszaállítása – $\hat{x} \rightarrow \tilde{x}(t)$ – során kizárólag lineáris műveleteket feltételezve: $x + \mathbf{e} \rightarrow x(t) + \mathbf{e}(t)$. Ez az $\varepsilon(t)$ additív jelenség, zaj. A függetlenségből következően a számunkra fontos $(-f_0/2, f_0/2)$ frekvenciasávban jó közelítéssel állandónak, tehát fehérzajnak tekinthető.

¹ A kvantálást 1938-ban Reeves „találta fel”. Észrevette, hogy a nem kvantált jelet átvéve a jelhez zaj adódik, és így a jel már nem állítható vissza. Ha azonban a zaj egy kvantált jelben okoz torzulást, és a hozzáadódó zaj kisebb, mint a kvantálási lépcső fele, akkor észrevehető, hogy hiba történt, és ennek megfelelően újra kvantálhatunk. Ezt a folyamatot nevezik regenerálásnak. Manapság a kvantálás a bináris ábrázoláshoz és a kódolhatósághoz szükséges. A kvantálással való átalakítást az ADC (analóg/digitál konverter) végzi.

2.3.2. A kvantálási zaj teljesítménye

Mivel e egy sztochasztikus folyamatnak tekinthető, ezért

$$P_e = M(e^2) = \int_{-\infty}^{\infty} e^2 f_e de = \int_{-\frac{\Delta}{2}}^{\frac{\Delta}{2}} e^2 \frac{1}{\Delta} de = \left| \frac{1}{\Delta} \cdot \frac{e^3}{3} \right|_{-\frac{\Delta}{2}}^{\frac{\Delta}{2}} = \frac{\Delta^2}{12}$$

Ezzel azonban nem kaptuk meg a hasznos jel és a zaj viszonyát. Tekintsünk egy referencia jel/zaj viszonyt: vegyünk egy, a teljes kvantálási tartományt kitöltő C amplitúdójú szinuszjelet. N szinten kvantálva $N=2C$. Ennek a szinuszjelnek a teljesítménye $P_s = \frac{C^2}{2}$. A kvantálás jel/zaj viszonya SNR (Signal to Noise Ratio).

$$SNR = \frac{\frac{C^2}{2}}{\frac{\Delta^2}{12}} = 6 \frac{C^2}{\Delta^2}. \text{ Tehát } N \text{ kvantálási szint esetén } \frac{C}{\Delta} = \frac{N}{2} \text{ és } SNR = \frac{3}{2} N^2 = \frac{3}{2} 2^{2n}, \text{ ha az } N$$

kvantálási szintet n biten kódoljuk. Ebből már adódik, hogy $SNR^{[dB]} = 10 \cdot \lg(SNR) = 1.74 + n \cdot 6.02 \text{ dB}$

A kvantálás jel/zaj viszonya tehát kifejezhető a kvantálás során felhasznált kódszavak hosszával: 8 bites kódok esetén kb. 49.7 dB, 16 bites kódszavak esetén pedig kb. 97.74 dB.¹

2.4. Logaritmikus (PCM) kvantálás

Megfigyelték, hogy az analóg telefonvonalon a távoli elofizeto hangja nagyon kis teljesítménnyel, míg a közeli elofizeto hangos (esetleg üvöltő) hangja nagy teljesítménnyel szerepel. Ez a teljesítményviszony elérheti akár a 60-70 dB-t is. A tapasztalat azt mutatta, hogy $n=12$ biten kellene kvantálni egyenletes kvantálással ahhoz, hogy a távoli elofizeto hangja is hallható legyen.

Stevens azonban észrevette, hogy az emberi fül a nagy amplitúdók esetén kevésbé érzékeny a hibákra, és kimondta az Stevens törvényt, miszerint $\frac{dx}{x} = c$, azaz ha a hiba és amplitúdó aránya konstans, akkor az érzeti világunk egyensúlyban van. Tehát a kvantálási lépcsők a 0-tól elfele nonek. Az ilyen kvantálást *nem lineáris kvantálásnak* nevezzük.

Kvantálási karakterisztikának nevezzük azt a karakterisztikát, amely a nem lineárist lineáris transzformálja (ez a lineáris kvantálás esetén egy egyenes).

$\Delta x/x = c_1$ és $\Delta y = C_2$. $f(x) = ?$ Elosztva egymással a két egyenletet, majd átrendezve az oldalakat, a következő összefüggést kapjuk:

$$f'(x) = \frac{c_3}{x}, \text{ ahol } c_3 = \frac{c_2}{c_1} \text{ és így } f(x) = \ln(x) + c_4$$

Ezt tovább finomítva két szabvány alakult ki, az európai PCM karakterisztika (A-law) és az amerikai PCM karakterisztika (μ -law).

A 12 bites lineáris kvantálással 96 kbit/s, míg a logaritmikus kvantálással 64 kbit/s átviteli sebességre van szükség. A fenti tömörítési módszert hívják érzeti tömörítésnek.

¹ Az itt felsorolt értékekből is látszik, hogy $SNR^{[dB]} \approx 6n$ és $SNR \approx 2^{2n}$.

2.5. Lineáris predikció

Motiváció: a beszéd 8kHz-es mintavételezése esetén sokszor fordul elő, hogy két egymást követő minta nem nagyon különbözik egymástól. Tehát ha nem a mintát, hanem a minták különbségét kvantáljuk, akkor kevesebb biten tudjuk átvenni ugyanazt az információt (tömörítés). Legyen ξ_k a k-adik időpillanatban lévő minta. Ekkor

$$\tilde{\mathbf{x}}_k = \sum_{j=1}^M w_j \mathbf{x}_{k-j}, \text{ ahol } w_j \text{ az ún. súlytényező (weight factor)}$$

Gyakorlatban ezeket a súlytényezőket és az így számított értékeket akkor használják, amikor a beszéd stacionáriusnak tekinthető (hasonló, és ezért ezek a súlytényezők néhány mintán keresztül érvényesek). A w súlytényezőket keretről keretre a mintákból határozzuk meg. Egy keret 10-20 ms hosszú (ez 80-200 mintát jelent), pl. egy magánhangzó tiszta fázisa lehet egy keret.

A w -k előállításánál arra törekszünk, hogy az így kapott $\tilde{\mathbf{x}}_k$ minél jobban megközelítse ξ_k -t. A predikció hibája v_k , de mivel ezek valószínűségi változók, össze-vissza ugrálnak, ezért ezeket jól jellemezhetjük a négyzetes várható értékükkel. Tehát w -k akkor optimálisak, ha

$$E = M(\tilde{\mathbf{x}}_k - \mathbf{x}_k)^2 \text{ minimális.}$$

$$E = M\left(\mathbf{x}_k - \sum_{j=1}^M w_j \mathbf{x}_{k-j}\right)^2 = M\left(\sum_{j=0}^M w_j \mathbf{x}_{k-j}\right)^2.$$

$$E = M\left(\sum_{j=0}^M (w_j \mathbf{x}_{k-j}) \cdot (w_j \mathbf{x}_{k-j})\right) = M\left(\sum_{i=0}^M \sum_{j=0}^M w_i \cdot \mathbf{x}_{k-i} \cdot \mathbf{x}_{k-j} \cdot w_j\right) = \sum_{i=0}^M \sum_{j=0}^M w_i \cdot M(\mathbf{x}_{k-i} \cdot \mathbf{x}_{k-j}) \cdot w_j$$

A szorzatok eredményét össze kell adni és el kell osztani. De ha a folyamat ergodikus, akkor 1 folyamat is magán hordozza a sokaság tulajdonságait, tehát a k -val is végigfuttatható. Így az összeg az i - j távolságra lévő minták átlagával közelíthető.

Véve tehát $R_{ij} = R_{ji} = M(\mathbf{x}_{k-i} \cdot \mathbf{x}_{k-j})$ autokorrelációs függvényt, a fenti összefüggés így alakítható át:

$$\sum_{i=0}^M \sum_{j=0}^M w_i R_{ij} w_j + w_m \sum_{j=0}^M R_{ij} w_j + \left(\sum_{i=0}^M w_i R_{im}\right) w_m + w_m R_{nm} w_m$$

$$\frac{\partial E}{\partial w_m} = \sum_{j=0}^M R_{mj} w_j + \sum_{i=0}^M w_i R_{im} + 2w_m R_{nm}$$

$$\frac{\partial E}{\partial w_m} = 2 \sum_{j=0}^M w_j R_{mj} = 0$$

$$\sum_{j=1}^M w_j R_{mj} = R_{m0} = R_{0m}$$

A harmadik sorában a deriváltat azért tettük 0-vá, mivel a minimalizáláskor a j -nek nem volt igazi változó tartalma, azt mi definiáltuk -1 -nek. Végül valójában M db. egyenletünk van. Ezeket felírva $w_1, w_2 \dots w_M$ -re:

$$\begin{bmatrix} R_{11}w_1 + R_{12}w_2 + \dots + R_{01} \\ \dots \\ R_{M1}w_1 + R_{M2}w_2 + \dots + R_{0M} \end{bmatrix}$$

M egyenlet és M ismeretlenes egyenletrendszer. Ez a lineáris predikció alapegyenlete. R-et a korrelációs mátrixnak nevezzük. A korrelációs mátrix további tulajdonságai:

- $R_{ii} = k$ függetlenül i -től, tehát $R_{11} = R_{22} = \dots = R_{MM}$; s ez legyen R_{∞} .
- Ez a tulajdonság továbbra is igaz a főátlóval párhuzamos átlókra: $R_{i,j} = R_{k-1}$, ha $i-j = k-1$. Ez az ún. Toeplitz mátrix.

$$\underline{\underline{R}} = \begin{bmatrix} R_{11} & R_{12} & \dots & R_{1M} \\ R_{21} & R_{22} & \dots & R_{2M} \\ \dots & & & \dots \\ R_{M1} & R_{M2} & \dots & R_{MM} \end{bmatrix} = \begin{bmatrix} R_{\infty} & R_1 & & R_{M-1} \\ R_1 & R_{\infty} & & \\ & & & R_1 \\ R_{M-1} & & R_1 & R_{\infty} \end{bmatrix}$$

$$\dots \dots \underline{\underline{b}} = \begin{bmatrix} R_{01} \\ R_{02} \\ \dots \\ R_{0M} \end{bmatrix}$$

Vagyis a lineáris predikció alapegyenlete tömör formában a következőképpen írható fel:

$$\underline{\underline{R}} \cdot \underline{\underline{w}} = \underline{\underline{b}} \dots$$

és ennek $|\underline{\underline{R}}| \neq 0$ esetén \exists megoldása $\underline{\underline{w}}$ -re,

$$\text{hogy } E = M(\mathbf{x}_k - \tilde{\mathbf{x}}_k)^2 \text{ min imális.}$$

Az egyenletrendszer megoldása nem egyszerű, de léteznek rá algoritmusok:

- Durbin
- Levinson
- rekurzív algoritmus
- adaptív algoritmusok

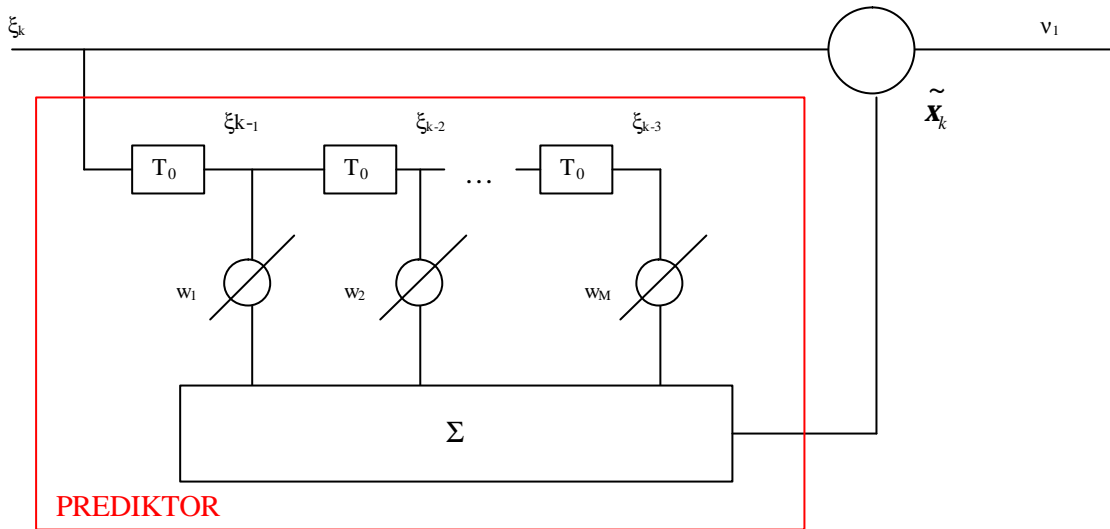
$$\mathbf{n}_k = \mathbf{x}_k - \tilde{\mathbf{x}}_k$$

Ha a közelítés jó, akkor \mathbf{n}_k kicsi és ekkor 1-1 \mathbf{n} -beli kódolási minta kvantálásához kevesebb kvantálási szint szükséges.

- \mathbf{x} kvantálása: n_2 bit
- \mathbf{n} kvantálása: n_1 bit
- $n_1 < n_2$

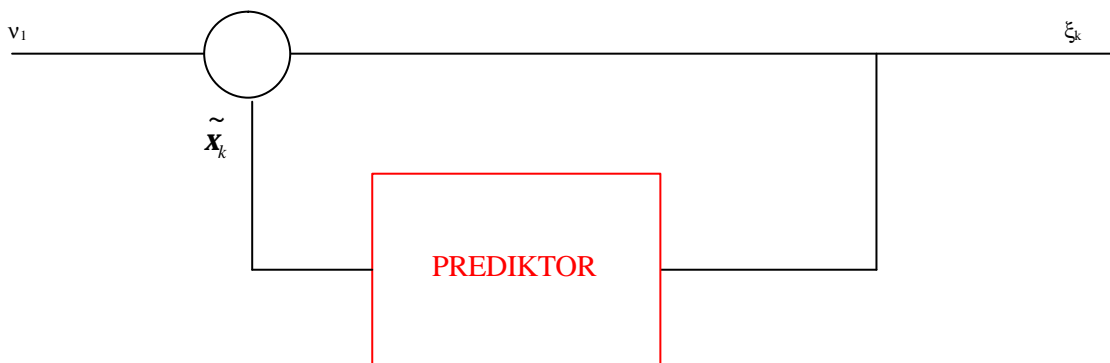
Kérdés: ez a tömörítés alkalmas-e arra, hogy az eredeti jel pontosan vagy elfogadható hibávan visszaállítható legyen?

2.5.1. Tömörítés és visszaállítás



- a tömörítő: transzverzális (digitális) szuro
- impulzusválasz függvénye: véges, azaz FIR (Finite Impulse Response)
- ez a lineáris predikció beszéd ananlízis modellje

A visszaállítás:

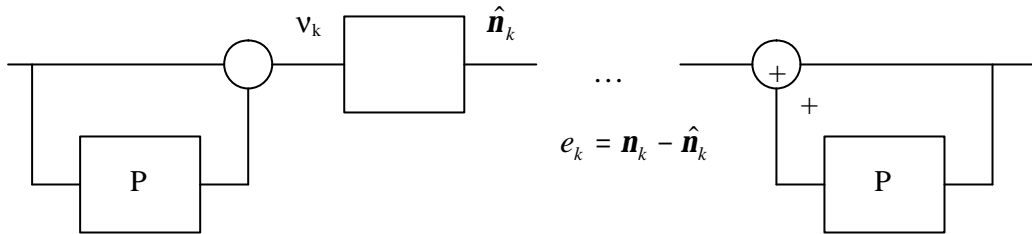


- A lineáris preditív beszékkódolás szintézi modellje
- A bemenet rögtön megjelenik a kimeneten + visszacsatolás
- 1 bemenő impulzusnak elvileg végtelen válasza van
- IIR (Infinite Impulse Response)

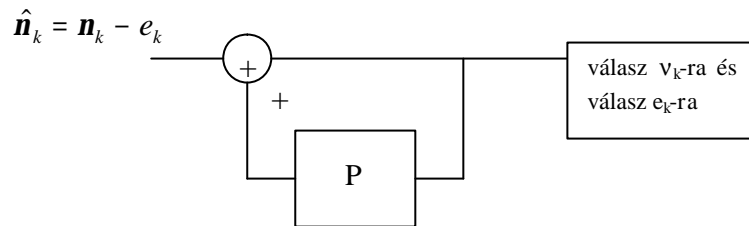
Fizikai értelem:

- másik, jól szegmentált magánhangzó esetén hasonló v-ket kapunk, de a w_i -k nagyon különböznek, jellegre azonban hasonlóak
- az egészet helyettesíthetjük azzal, mintha a rendszert a zöngé ütemében gerjesztenénk
- a szintézis modell ehhez a szemlélethez hasonlóan működik
- más, pl. hosszan tartható zöngétlen hangok esetében (f, s) ez a gerjesztő hang fehérzajszerű, véletlenszerű, hasonló, mint a vokális traktusban (itt is véletlenszajgenerátor, turbulencia, leszakadó levegőrészecskék, stb.)

2.5.2. A kódoló gyakorlati megvalósítása

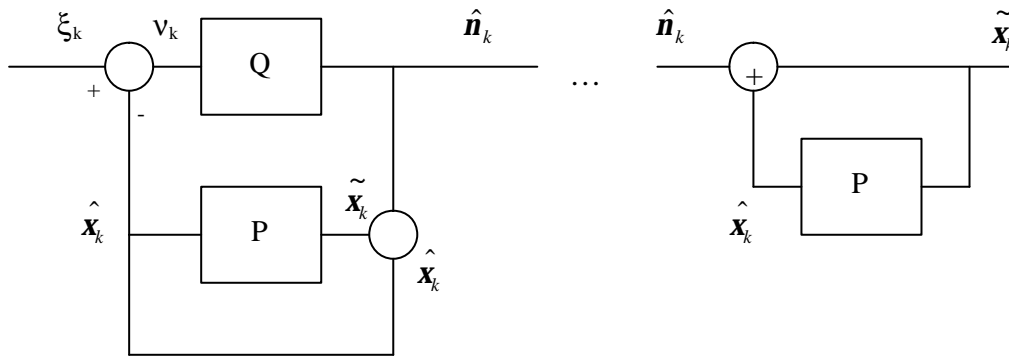


Mi lesz a kvantálási hibával a visszaállítás során?



- e_k : a tényleges és a kvantált érték közötti különbség, a kvantálási hiba
- ha csak egy ilyen is van, akkor a vevőben lévő IIF szűrés miatt végtelen választ ad
- állandó hibát okoz, halmozódik a kvantálási hiba, egy idő után kitér a dekódoló tartományból
- minden hibát, amit v_k elszűnved, a kimeneten halmozottan kapunk vissza

A kvantálási hibát azonban kézben lehet tartani: valójában a kódolóba is beépítjük a dekódolót.



$$\mathbf{x}_k - \tilde{\mathbf{x}}_k = \hat{\mathbf{x}}_k + \mathbf{n}_k - (\hat{\mathbf{x}}_k + \hat{\mathbf{n}}_k) = \mathbf{n}_k - \hat{\mathbf{n}}_k$$

- vagyis a megmaradó hiba csak a tényleges kvantálási zaj
- a bemenet és a kimenet csak a kvantálási hibában különbözik
- a k -adik időpillanatban lévő hiba *kizárólag* a k -adik időpillanatban elkövetett hibától függ (emlékezet nélküli csatorna)

2.5.3. Lineáris predikció a gyakorlatban

- a jelet 10-15 ms-os darabokra bontjuk (ezek az ún. keretek)
- meghatározzuk v -t
- meghatározzuk azt a 10 LPC együtthatót a vevonek (Ezek a w -k)
- a túloldalon visszaállított jelből valamilyen trükkal megpróbáljuk meghatározni a predikciós együtthatókat (vagy eltároljuk, ha arról van szó)

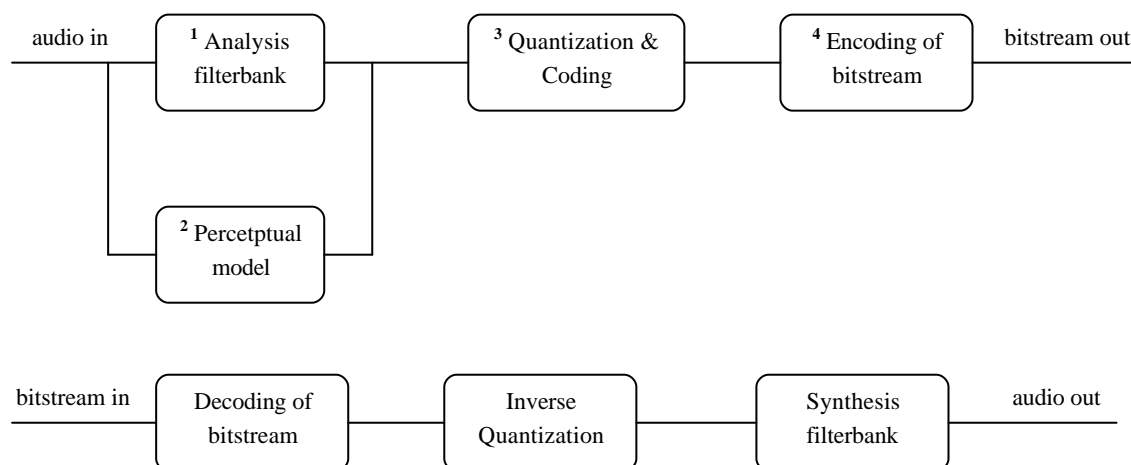
Trükkös esetek

- v -k közül a két legnagyobbat választjuk ki vagy azt a hármát, amely egymás mellett a legnagyobb
- elképzelünk 1024 hibasorozatot, a tényleges v -k helyett ebből a készletből visszük át azt, amelyik valamilyen értelemben a legjobban hasonlít a tényleges hibasorozatra (CELP)

2.6. Érzeti (részszávos kódolás) – perceptual (subband) coding

2.6.1. Frekvenciamaszkolási jelenség

- minden sávra megállapítjuk, hogy milyen energiájú összetevők vannak a jelben
- az elfedett összetevők kihagyása
- a kvantálási zaj növelésének lehetősége: úgy kvantálunk, hogy a kvantálási zaj ne legyen nagyobb, mint az elfedési szint (itt a tömörítési lehetőség)



¹ A bemenetre kerülő jelet összetevőkre bontja; elony, hogy a kisebb frekvenciájú jelet kisebb frekvencián kell mintavételezni, lejátszani (szűrosor)

² Érzeti modell: elfedési görbe meghatározása

³ Kvantáló és kódoló: az összetevők kvantálása több lépcsőben

⁴ bitfolyam kódolása szabványos formátumra (mintavételi frekvencia, szűrokomponensek)

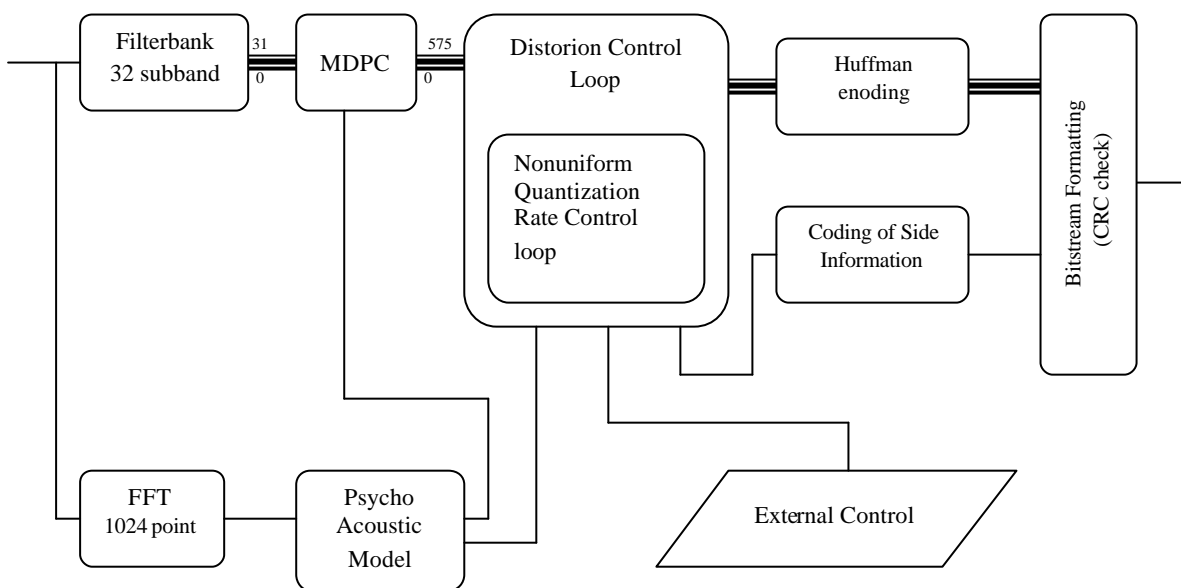
Elony az LPC-vel szemben:

- az LPC-ben át kell vinni a súlyokat, a hibát és a szegmentálást, és egy szegmensben stacionáriusnak tekintjük a beszédet
- az új módszerrel nemcsak beszédet, hanem más audio jelet is át lehet vinni

2.6.2. Motion Pictures Expert Group

- hivatalos neve: ISO/IEC JTC/SC2 9/WG11
- Feladatok:
 - ☞ digital audio bradcasting (DAB)
 - ☞ ISDN
 - ☞ tárolás
 - ☞ DVB, HDTV

- ☞ Internet streaming
- ☞ hordozható audio mp3-lejátszó
- ☞ audio filecsere
- a fentiekbol csak a hangkódolással foglalkozunk
- különböző minőségi szintek ugyanazon elv alapján (layernek hívják őket)
- ☞ MPEG1 – 1992 (192 ... 32 kb/s)
- ☞ MPEG2 – 1994 (újabb frekvenciák: 16, 22,05, 24 kHz)
- ☞ MPEG3 – HDTV-hez készült volna, de visszatértek az MPEG2-höz
- ☞ MPEG4 – 1998
 - újdonság: nem a tömörítés hatékonyságát javítják, hanem új szolgáltatásokat hoznak be (interaktív TV)
 - hangkódolási tartomány: 2kb/s ----
- ☞ MPEG7 – 2001? (kidolgozás alatt, tartalom-reprezentációs szabvány)



- bemenet: Digital Audion Signal (PCM 768 kb/s)
- kimenet: Coded Signal (2...192 kb/s)
- Filterbank: nagyjából a kritikus sávok szerinte felbontja a jelet, minden sávot további 18 részsávra lehet bontani, így lesz 0...575 sáv
- MPC: Modified Discrete Cosine Transform
- segédinformáció: pl. milyen Huffman-táblát használjunk

A szabvány a dekódert írja le

- nyílttá válik a lehetőség a kódoló implementálására (ugyanazon dekódoló jobb kódolóval jobb minőséget produkál)
- változó bitsebesség (lehetséges az is, hogy bizonyos részeket más bitsebességgel viszünk át)
- többféle forrást kell tudnia: mono, sztereo hagyományos és kombinált, kétcsatornás

Problémák

- visszaállított jel: csendes szakaszban is megjelenik valamekkora jel, ez a zaj a pre-echo jelenség (nem kauzális a kódoló)
- ha úgy tudunk szegmentálni, hogy az időelfedési jelenség miatt nem halljuk, azzal csökkenthetjük a hibát

3. BESZÉDVÁLASZÚ RENDSZEREK

3.1. Gépi beszédkeltés alapfogalmai: három kategóriát különböztetünk meg

3.1.1. Kötött szókészlet

- tudjuk, hogy a rendszernek mit kell majd mondania
- állandó üzenet („a hívott szám nem elérhető”, kiterjesztett magnetofon)
- változó elemek
 - ☞ primitív: „Önnek üzenete érkezett 2000 május”
 - ☞ bonyolultabb: „A hívott szám megváltozott, az új szám: 325-29-48”

3.1.2. Kötetlen szókészlet (text to speech, szövegfeldolgozó)

- gyakorlatilag ilyen nincs
- széles szókincssel kell rendelkeznie kiinduló állapotban, és tetszőlegesen bővíthető
- ha tudjuk a tematikát, kifejezéseket, akkor meg lehet tanítani

3.1.3. Vegyes rendszerek

- vannak állandó üzenetek (ezeket nem célszerű TTS-sel megoldani, mert fárasztó)
- vannak változó üzenetek

szókészlet * minőség = konstans

3.2. Kötött szókészletű rendszerek tervezési szempontjai

3.2.1. Tematika felderítése

- az adott rendszeren mik azok az információk, melyeket el kell juttatni a felhasználóhoz
- mik ezeknek a módjai
- a felhasználók figyelembevétele (kezdő + profi különböző)

3.2.2. bemondandó szöveg tervezése

3.2.3. szótárkészlet kialakítása

- az elozóval szinkronban
- kompromisszum a minőség és a bonyolultság között
- szótárelemek számára algoritmus

3.2.4. bemondó választása

- akusztikai arculat

3.2.5. akusztikai adatbázis elkészítése

- felvétel készítése
- elemek kivágása és feldolgozása

3.2.6. rendszerbeillesztés

3.3. **Konkrét példa:**

- többnyelvu számbemondó tervezése és megvalósítása
- Hagyományos megoldás: írásnak megfelelő számelemek összefuzése szünetekkel
- Pl: 125000 :

(English)	one	hundred	and	twenty	five	thousand
(German)	ein	hundert	fünf	und	zwanzig	tausend
(Hungarian)	száz	huszon	öt	ezer		
(Portuguese)	cento e	vinte	cinco	mil		

- Magyar: 25 db elem, portugál: 53 db elem

	Basic element	English	German	Hungarian	Portuguese
1.	1	one [wɔ̃v]	ein [ai v]	egy [eθ]	um [ũ]
2.	1--	--	eins [ai v σ]	--	--
3.	1--	--	eine [ai v ↔]	--	--
4.	2	two [tv:]	zwei [τσ ω αι]	kettő [κετ:O:]	dois [doj]
5.	3	three [Tpi:]	drei [δ ρ αι]	három [ηα:ρομ]	três [tre]
6.	4	four [φ :]	vier [φ i:]	négy [ve:θ]	quatro [kwatru]
7.	5	five [φαιw]	fünf [φ Φ v φ]	öt [Oτ]	cinco [s ã ku]
8.	6	six [σικσ]	sechs [ζ E κ σ]	hat [η τ]	seis [sEj]
9.	7	seven [sevn]	sieben [ζ i: β v]	hét [ηε:τ]	sete [σEτ↔]
10.	8	eight [eit]	acht [α ξ τ]	nyolc [j ol τσ]	oito [ojtu]
11.	9	nine [v ai v]	neun [v φ v]	kilenc [kilεv τσ]	nove [v (w ↔)]
12.	10	ten [ten]	zehn [τσ ε: v]	tíz [ti:ζ]	dez [δE]
13.	10x			tizen... [tizen]	
14.	11	eleven [i λ ε w v]	elf [E λ φ]	--	onze [j] ↔]
15.	12	twelve [τ ω ↔ λ w]	zwölf [τσ w ↵ λ φ]	--	doze [δ j] ↔]
16.	13	thirteen [T ↔ : ti: v]	dreizehn	--	treze [tpe] ↔]

17.	14	fourteen	vierzehn [φ ι ρ τ σ ε: v]	--	catorze [κ α τ ο ρ] ↔
18.	15	fifteen	fünfzehn	--	quinze [κ ῖ] ↔
19.	16	sixteen	sechszehn	--	dezasseis [δ ↔] α σ Ε φ]
20.	17	seventeen	siebzehn	--	dezassete [δ ↔] α σ Ε τ ↔
21.	18	eighteen	achtzehn	--	dezoito [δ ↔] ο φ τ υ]
22.	19	nineteen	neunzehn	--	dezanove [δ ↔] α ν [ω ↔]
23.	20	twenty [twenti]	zwanzig [τ σ ω α ν τ σ ι X]	húsz [η υ: σ]	vinte [ω ῖ τ ↔]
	Basic element	English	German	Hungarian	Portuguese
24.	2x			huszon.. [η υ σ ο ν]	vinte e [ω ῖ τ φ]
25.	30	thirty	dreizig	harminc [η ρ μ ι ν τ σ]	trinta [τ ρ ῖ τ α]
26.	3x				trinta e [τ ρ ῖ τ α φ]
27.	40	forty	vierzig	negyven [νε θ ω ε ν]	quarenta [κ ω α ρ ῆ τ α]
28.	4x				quarenta e
29.	50	fifty	fünfzig	ötven	cinquenta [σ ῖ κ ω ῆ τ α]
30.	5x				cinquenta e
31.	60	sixty	sechzig	hatvan [η τ ω ν]	sessenta [σ ↔ σ ῆ τ α]
32.	6x				sessenta e
33.	70	seventy	siebzig	hetven	setenta [σ ↔ τ ῆ τ α]
34.	7x				setenta e
35.	80	eighty	achtzig	nyolcvan	oitenta [ο φ τ ῆ τ α]
36.	8x				oitenta e
37.	90	ninety	neunzig	kilencven	noventa [ν υ ῶ ῆ τ α]
38.	9x				noventa e
39.	100	hundred [η ρ ν δ ρ ↔ δ]	hundert [η Υ ν δ τ]	száz [σ α: ζ]	cem [σ ῆ ῖ]
40.	1xx				cento e [σ ῆ τ υ φ]
41.	200				duzentos [δ υ] ῆ τ υ]
42.	300				trezentos [τ ρ ε] ῆ τ υ]

43.	400				quatrocentos [κωατρυσῆ τυ]
44.	500				quinhentos [κί]ῆ τυ]
45.	600				seiscentos [σΕφ[σῆ τυ]
46.	700				setecentos [σΕτ<→σῆ τυ]
47.	800				oitocentos [οιτυσῆ τυ]
48.	900				novacentos [ν (ῶ<→σῆ τυ]
49.	1000	thousand [Ταυζνδ]	tausend [τ αυ ζ ν τ]	ezer [ΕζΕρ]	mil [mil]
	Basic element	English	German	Hungarian	Portuguese
50.	1000x				mil e [milj]
51.	1000000	million [μιλφ<→ν]	million [μ ι λ φ ο : ν]	millió [μιλιο:]	milhão [μιλ×© ù]
52.					milhão e [μιλ×© ù φ]
53.					milhões [μιλ×ο ï]]
54.					milhões e [μιλ×ο ï] φ]
55.		billion [βιλφ<→ν]	milliarde [μιλφαρδ<→]	milliárd [μιλια:ρδ]	bilhão [βιλι© ù]
55.					bilhões [βιλι ï]]
56.	0	0 [ou]	--	--	--
57.	0-	zero [ζι<→ρου]	null [νυλ:]	nulla [νυλ:]	zero30 [ζερυ]
58.		and [Θνδ]	und [υντ]	--	e [φ]

Természetes kiejtéshez biztosítani kell:

- folyamatos kiejtés, helyes pozíciójú és hosszúságú szünetekkel
- a számelemek kiejtési helytől függő idoszerkezete
- spektrális és intenzitás folytonosság (koartikuláció figyelembe vétele) az elemhatárokon
- szóhangsúlyok és alaphangfrekvencia változások helyessége

3.3.1. Folyamatos kiejtés

A megfelelő helyeken, megfelelő hosszúsággal beiktatott szünetekkel, a 2, 3, 4, szempontok szerint kiválasztott elemek folytonos összefuzése (vágás nullátmenetnél negatívból pozitívba)

3.3.2. A számelemek kiejtési helytől függő idoszerkezete

Kezdo (B, beginning, pl. 1234567), középso (M, middle, 1231567), záró (L, last, 1234561) elem szükséges a többi szempont szerint kiválasztott minden elemből (elvileg).nagyszámú (közel ezer) kimondott szám vizsgálata alapján

3.3.3. Spektrális és intenzitás folytonosság (koartikuláció figyelembe vétele) az elemhatárokon

Minden elemre hat az előző és a következő elem

Lehetséges pozíciók:

- Egyedül áll (6)
- Felsorolás (12, 2 56.)
- Elso (elemXXX)
- Belso (XXXelemXXX)
- Záró (XXXelem)

3.3.3.1. Az 1 példája**3.3.3.1.1. Angolul one.**

oneXXX után *hundred, thousand, million, billion*, (pl., 1100)

XXXone előtt *thousand, million, billion, and, twenty, thirty.... ninety* (pl., 1100, 101, 21).

one elemkészlet:

(1) szabály: *one* felsorolásban vagy egyedül 1, 2, 3

oneXXX esetekben

(2) szabály: *one (one hundred)*, a (h) módosítja az (n)-et,

(3) szabály: *one (one thousand)* az (n) és a (t) azonos artikulációs bázisú, ezért az (n) rövidebb lesz,

(4) szabály: *one (one million)*, az (n) (m)-be megy át,

XXXone esetekben

(5) szabály *one (thousand one, hundred and one, etc.)* zárhang és (v) találkozása,

(6) szabály: *one (twenty one, etc.)* i és (v) találkozása ,

XXXoneXXX

(5) + (2), (5) + (3), (5) + (4), (6) + (3), (6) + (4)

Összesen: 11 (1+3+2 +5) elméleti lehetőség.

A hagyományos módszer minden elemére elvégezve a fenti elemzést, a spektrális és intenzitás folytonosság biztosítható.

3.3.3.1.2. Németül:

(1) szabály: *Ein* ha a szám 1-nél nagyobb, *eins* ha egyedül áll (pl., 1, 2, 3), *eine* pl. *eine million* és *eine DM*.

(2), (3), és (4) változatlan, mert az *einXXX* és a *oneXXX* kategóriái megegyeznek
XXXein különbözik

(5) szabály: *hundert ein, tausend ein, milliard ein*, zárhang és magánhangzó találkozása,

(6) szabály: *million ein*, nazális és magánhangzó találkozása.

XXXeinXXX

(5) + (2), (5) + (3), (5) + (4), (6) + (2), (6) + (3)

Összesen: 13 (3+3+2 +5) elméleti lehetőség.

3.3.3.1.3. Magyarul egy:

(1) szabály: *egy* egyedül áll (1, 2, 3 stb.),

egyXXX

(2) szabály: *egy millió* és *egy milliárd*

(3) szabály: *egy ezer*, pl. 31000), zöngés alveolo-palatális zárhang és magánhangzó találkozása,

(4) szabály: *egy száz* pl. 3125000, zöngés alveolo-palatális zárhangot zöngétleníti a *száz sz* hangja,

XXXegy

(5) szabály: *...n egy*, pl. 51, 61, 71, etc.) a nazális hang módosítja az *e-t*,

(6) szabály: *millió egy*, pl. 5000001) magánhangzó-magánhangzó kapcsolat.

XXXegyXXX

(5) + (2), (5) + (3), (6) + (3), (6) + (4)

Összesen: 10 (1+3+2 +4) elméleti lehetőség.

3.3.3.2. A legfontosabb regresszív koartikulációs szabályok

az elozo elem utolsó hangja	az alábbira változik, ha	a következő elem első hangja
b, d, g, v, z, ʒ	p, t, k, f, s, S	zöngétlen
ts	ts felpattanás (burst) nélkül	s
t	t felpattanás (burst) nélkül	n
n	n(k)	k
n	n(h)	h
n	mn	n
n	ŋ	ŋ
n	m	m, b, p
magánhangzó	átmeneti szakasz	magánhangzó
magánhangzó	palatalizált átmeneti szakasz	palatális

A legfontosabb progresszív koartikulációs szabályok

ha az elozo elem utolsó hangja	és a következő elem első hangja	akkor a következő elem első hangja az alábbira változik
nazális	magánhangzó	nazalizált átmeneti szakasz
palatális	magánhangzó	palatalizált átmeneti szakasz
magánhangzó	magánhangzó	átmeneti szakasz

previous element in concatenation	element of the inventory	next concatenated element	example number
--	1(εθ)	--	1
--	1(εθ)(m)	millió, milliárd	1564322
millió	1(o)(εθ)	--	3000001
1000 (ΕζΕρ)	1(εX)	100 (σ:ζ)	1100
..(an), ..(en), ..(on)	1(n)(eJ)	1000 (J)(ΕζΕρ)	51000
..(an), ..(en), ..(on)	1(n)(eJ)	millió, milliárd	61000000
any element	2(κετ:O:)	100 (σ:ζ), millió	200, 312
any element	2(κετ:O:)	1000 (O:)(ΕζΕρ)	2000
any element	3(ηα:ρομ)	100, 1000, millió	300, 3000
any element	4(νε:θ)	1000 (J)(ΕζΕρ)	4555
any element	4(νε:θ)	millió (J)(μιλιο:)	4000000
any element	4(νεX)	100 (σ:ζ)	400
--, 30, 100, 1000	5(Oτ),	100, 1000, millió	535, 5000
..(v), ..(en), ..(on)	5(n)(Oτ)	100, 1000, millió	65, 75, 25
any element	6 (η τ),7(ηε:τ)	100, 1000, millió	600, 700
1000, millió	8(βλ τσ) without burst in [τσ]	100	812
10--90, 100	8(βλ τσ)	1000, millió	8000, 8000000
1000, millió	9(κιλΕν τσ) without burst in [ts]	100	900
10-90, millió	9(κιλΕν τσ)	1000, millió	59000
any element	10(τι:ζ)	(ΕζΕρ),millió,	510000
--, 100, 1000, millió	variants for ending -Ev in numbers 1x, 4x, 5x, 7x, 9x	1, 5,	11, 115
" "	-Ev)(k)	2, 9	12, 142, 79
" "	-Ev)(h)	3, 6, 7	13, 53
" "	-Ev)(n)	4	14, 94
" "	-Ev)(β)	8	18, 98
" "	20 (ηυ:σ)	(ΕζΕρ), millió,	20000
" "	variants for ending -on in numbers 2x, (-on)	1, 5	21, 125
" "	-on)(k)	2, 9	22, 122
" "	-on)(h)	3, 6, 7	23, 1223
" "	-on)(n)	4	24, 224
" "	-on)(θ)	8	28
" "	30 (η ρμιν τσ)	1,2,3,4,5,6,7,8,9, (ΕζΕρ), millió	30256
" "	variants for ending - v in numbers 6x, 8x, - v)	1, 5	61, 185
" "	- v)(k)	2, 9	62, 289
" "	- v)(h)	3, 6, 7	63, 187, 666
" "	- v)(n)	4	64, 164
" "	- v)(β)	8	168, 968
1,2,3,4,5,6,7,8,9, 1000	100 (σ:ζ)	1,4,5,8,40,50,80, 1000,millió	1000010í

" "	100 (σ α :σ)	2,3,6,7,9,10,1x,20, 2x,,30,60,70,90	102
5,6,7,8,9,100	1000 (EζEρ)	any number element	5001
2	(O:)(EζEρ) 1000	" "	2000
1, 4	(∅)(EζEρ)1000	" "	4000
40,50,60,70, 80,90	(n)(EζEρ) 1000	" "	50000
3	(m)(EζEρ)	" "	3000

3.3.3.2.1. Szóhangsúlyok és alapfrekvencia változások helyessége

sample number	pronounced style	comment
121	o n e h u n d r e d a n d t w e n t y o n e. AB N N AM AL	.=full stop AL= accent and falling intonation in the last item
2151	t w o t h o u s a n d o n e h u n d r e d a n d f i f t y o n e. AB N AM N N AM AL	AB, AM= accents in the number

A számok kimondásakor több hangsúly is megjelenik.

- AB: kezdő hangsúly
- AM: közbeni hangsúly
- AL: záró hangsúly, eső intonáció
- N: semleges, hangsúlytalan elemek

Szerencsére a helyes időtartamot biztosító elemek (B, M, L) megfelelő tervezés esetén magukban hordozzák a helyes hangsúlyt is.

- Ha a számelem a mondat végén áll, (pl. Az ön számlájának egyenlege: **53424** forint) eső jellegű intonációja lesz.
- Ha a mondat közepén helyezkedik el, (pl. Az ön számláján **53424 forint** összegű tranzakció valósul meg.) a számelem intonációja laposabb, lebegőbb).

3.3.4. A számkimondó megvalósítása

Elozmény:

- az elemi (hagyományos) építokockák, számelemek meghatározása
- a kimondási szabályrendszerek (idotartam, koartikulációs, hangsúly és intonáció) meghatározása

3.3.4.1. *A felolvasandó szöveglista meghatározása*

- Vivoszóveg kialakítása az építokockák és a szabályrendszer alapján
- Example of determining the list of number elements and the source from where they will be cut out (for English)

number element	position	rule type	example of the recorded sample number which contains the element
one	B, L	(1)	1, 2, 1. (with pauses)
one(h)	B, M	(2)	121, 2151
one(t)	B, M	(3)	1121, 2001121
one(m)	B, M	(4)	1121151
(d)one(h)	M	(5), (2)	1122
(d)one	L	(5)	121101.
(ty)one(t)	M	(6), (3)	531231
(ty)one	L	(6)	541.

3.3.4.2. *A felolvasandó szöveg felvétele*

- Minden elemet a megfelelő vivoszóvegben kell felolvasni. A vivoszóveget célszerű redundánsra tervezni (minden elem legalább kétszer forduljon elő).
- Az egyes elemek között kb. 2 sec szünetet célszerű tartani.
- Nagyobb egységenként (pl. oldalanként) érdemes hosszabb szünetet tartani.
- Az oldal megkezdése előtt az előző oldal végének meghallgatása.
- Összpontosítás az egyenletes hangmagasság, hangero és beszédsebesség biztosításához.

3.3.4.3. *A hangelemek kivágása a felolvasott vivoszóvegből*

- Kivágás előtt a felolvasás helyességét ellenőrizni, hiba esetén a redundáns elem elvétele.
- Idobeli (esetleg spektrális) vizsgálat alapján határok megállapítása.
- Elemek elmentése az építőelem lista és a szabályrendszernek megfelelő logikus rendben (adatbázis, könyvtárstruktúra, stb.)

3.4. Szövegfelolvasó rendszerek (TTS)

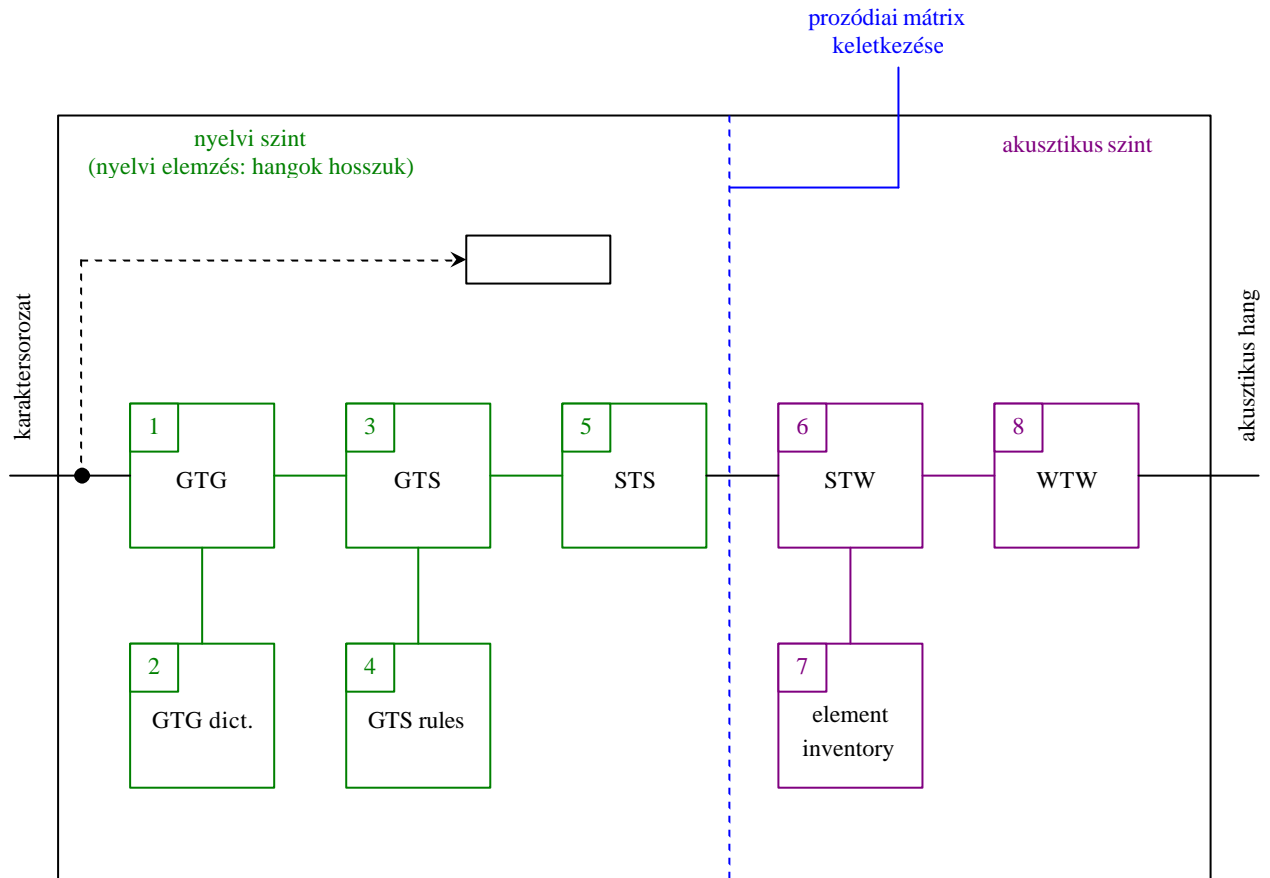
- **Szövegfelolvasó (text to speech):** adott nyelv köznapi szókincsében előforduló szövegek felolvasása (kb. egy 8 éves gyerek szókincsének megfelelő)
- **Üzenet felolvasó (concept text to speech):** a kifejezni kívánt üzenetre vonatkozó jelekkel ellátott szöveg felolvasása
 - ☞ pl. [Conf_Req] A gépkocsi típusa [Car_Type] Volkswagen Golf
- **Többnyelvu TTS (multilingual):** azonos építőelemek minél nagyobb halmazának egységes keretben történő felhasználása TTS rendszer megvalósításához több nyelven. Ideális esetben (ami cél és nem pedig a valóság) azonos program kód (és hardware), a nyelvfüggo adatok egységes szerkezetu, külső adatbázisban helyezkednek el.
- **Poligott TTS:** azonos hangon szóló TTS
 - ☞ szeniális paraméteres leírás, mely nincs emberi hanghoz kötve
 - ☞ egy bemondó sok nyelven mondja el a szöveget
- **Kötött tematikájú (domain specific) TTS:** csak egy adott témakör (pl. menetrend, időjárás, szálloda foglалás) szöveg felolvasására alkalmas rendszer. Átmenet egy hagyományos kötött szókészletu és egy TTS rendszer között.
- **képernyő felolvasó (screen reader):** számítógép monditor tartalmát értelmező vakok és gyengénlátók számára. Nem tartalmaz TTS-t, csak illesztést képes alkalmazás és TTS között.

3.4.1. Osztályozási szempontok

- milyen nyelveken szeretnék felolvasatni
- milyen szövegeket – egy teljes rendszert általában csak a TTS kimenete alapján ítélnék meg, a bemenetet nem látják.
 - ☞ szövegtípus: általános, szakszöveg, e-mail, SMS, stb.
 - ☞ mondatípus: kijelentő, kérdő, felkiáltó, egyéb érzelm kifejezése, CTS
- milyen minőségben
 - ☞ érthetőség : intelligibility
 - ☞ természetesség: naturalness
 és ezek nem is feltétlenül korrelálnak egymással
- milyen hangokon – egy illetve több hangon, amit kiemelünk, más hangon szóljon
- milyen paraméterek állíthatók
 - ☞ sebesség
 - ☞ hangmagasság
 - ☞ suttogás
 - ☞ rekedtség
 - ☞ szünetek hossza
 - ☞ betűzés
- milyen platformokon fusson
 - ☞ hardware
 - ☞ operációs rendszer (Windows, Unix, OS/2)
 - ☞ erőforrásigény, csatorna – nem mindegy, hogy mobiltelefonban vagy távközlési központban
- milyen vezérlési felületek, API-k
- bővítési, továbbfejlesztési lehetőségek – mit ad hozzá a felhasználó és mit a fejlesztő, pl. rövidítésfeloldó
- milyen speciális igények merülnek fel – pl. IT, callback egy adott szó kimondásának elején/végén, kimondás állapotának lekérdezhetősége (menetrendnél)
- milyen támogatást ad a TTS fejlesztő az alkalmazásfejlesztőnek

3.4.2. Felépítés: néhány alaprobléma

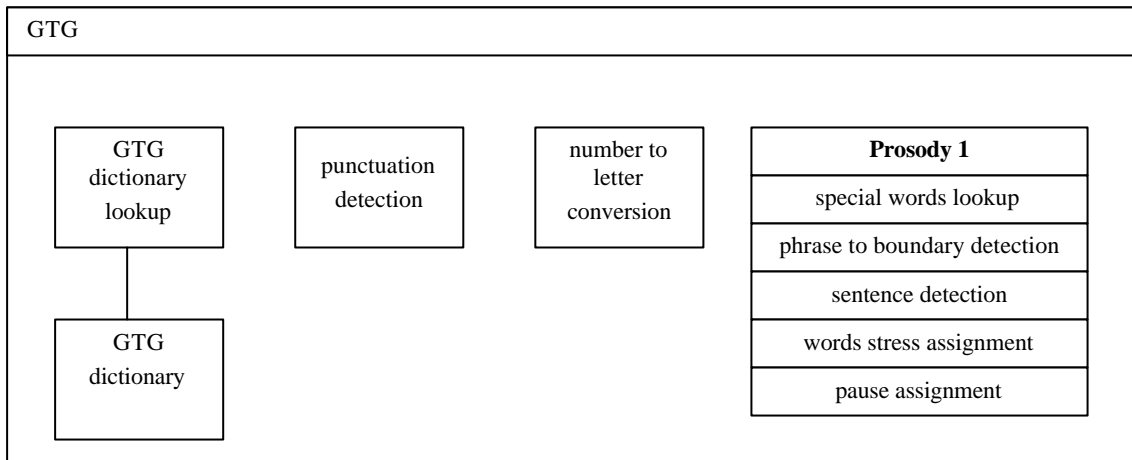
- Az írás diszkrét, a szavakat szünetek választják el. A beszédben a szavak folyamatosan következnek egymás után, csak nagyobb egységeket (prozódiai egység) választ el szünet. A beszédben a folyamatosság megértése teszi nehézé a megértést.
- Az írott hibákat másképp kezeljük: az akusztikus formára sokkal érzékenyebbek vagyunk.
- Fontos: a TTS bemenetére minél helyesebb és minél részletesebb szókimondást segítő információt tartalmazó jelsorozat érkezen.



- | | | |
|----------------------|---|---------------------------------------|
| 1. GTG | : | Grapheme to grapheme (írásjel→betű) |
| 2. GTG dict. | : | GTG dictionary (szótár) |
| 3. GTS | : | Grapheme to Sound (betű→hang) |
| 4. GTS rules | : | szabály és szótár |
| 5. STS | : | Sound to Sound (hang→hang) |
| 6. STW | : | Sound to Wave (hang→hanghullám) |
| 7. element inventory | : | hangelem-tár, akusztikai adatbázis |
| 8. WTW | : | Wave to Wave (hanghullám-feldolgozás) |

- 1-5 elvileg lehet nyelvfüggetlen, viszont 6-8 mindenképpen nyelvfüggo
- ha igazán általánossá akarjuk tenni, akkor nagyon bonyolult és nagy leíró nyelvre van szükség

3.4.2.1. GTG



- punctuation detection: azért fontos, mert pl. egy mondatban pont sok helyen lehet (rövidítések, ..., mondat vége, stb.)
- special words lookup: pontosvesszo, pont, vesszo, zárójel, csillag, bizonyos dolgokat nem mindig akarunk hallani
- phrase boundary detection: egységként kimondható szavak, frázisok között mikor tartunk szünetet
- sentence detection: intonáció egy mondatra

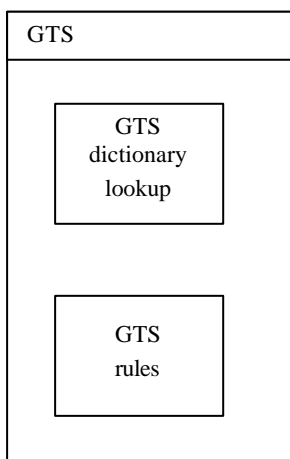
Példa: (6 soros idézet a BME szabályozási rendeletéből)

- szövegfelolvasó mit dolgozzon fel egy egységként (az egész egy mondat)
- többszintű nyelvi elemzés szükséges

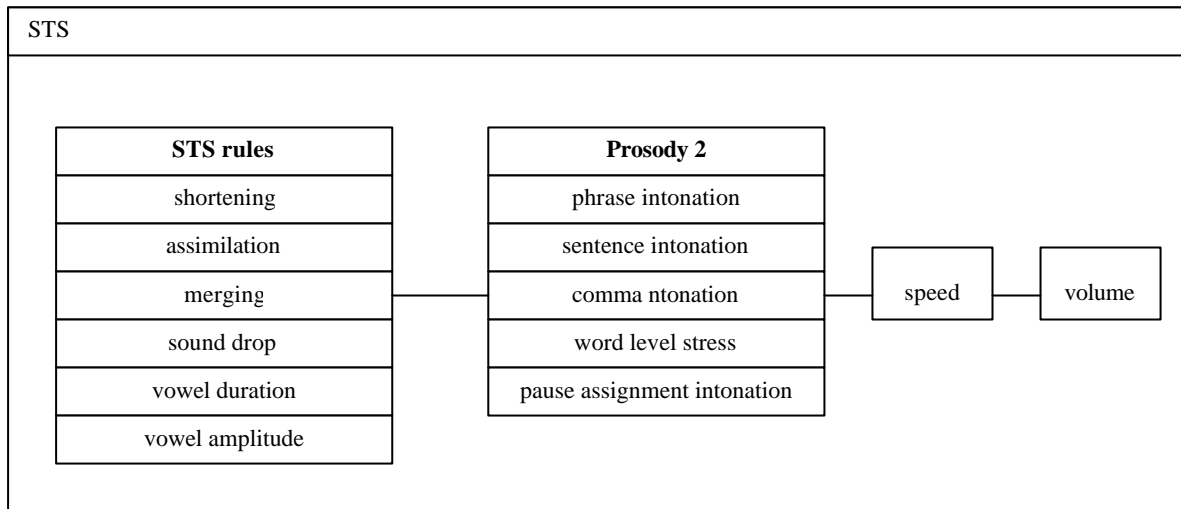
Írott szövegnek nem egy az egyben felelnek meg a kimondott hangok

- hasonulások
- röviden írjuk, hosszan ejtjük és viszont
- mássalhangzó torlódások
- betukép helyes értelmezése szó illetve morféma határon (malacság, egészség)

3.4.2.2. GTS

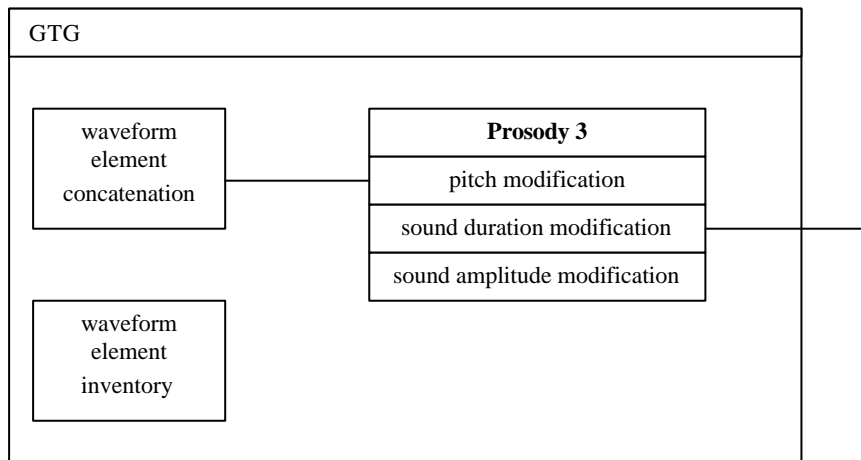


3.4.2.3. STS



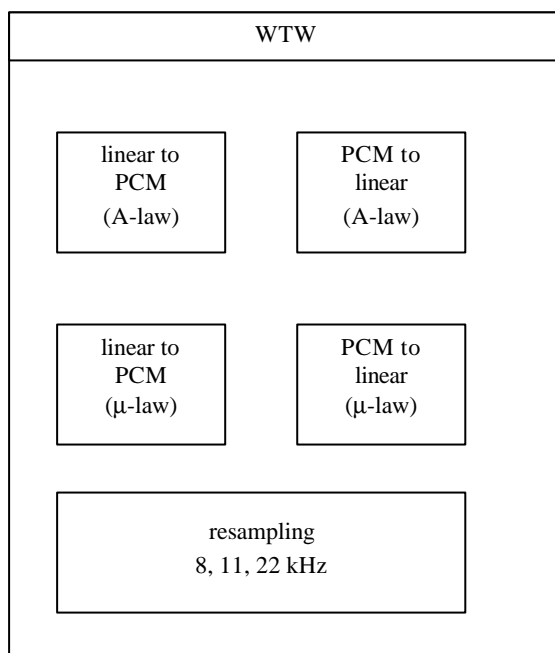
- speed: szünetek kivágása (vigyazni kell vele, mert ha figyelmetlenül vagdossuk ki a szüneteket, nem ugyanazt a hangot kapjuk)
- prosody 2: magasszintu leírás

3.4.2.4. STW



- waveform element inventory: valamilyen akusztikai adatbázis
 - ☞ paraméterek (pl. LPC) forráskódolása
 - ☞ hullámforma-kódoló
- prosody 3
 - ☞ a jó minőségű szövegfelolvasók esetén kulcsfontosságú
 - ☞ a prozódiai mátrixban előírtakat el kell végezni (módosítások és folytonosság biztosítása)

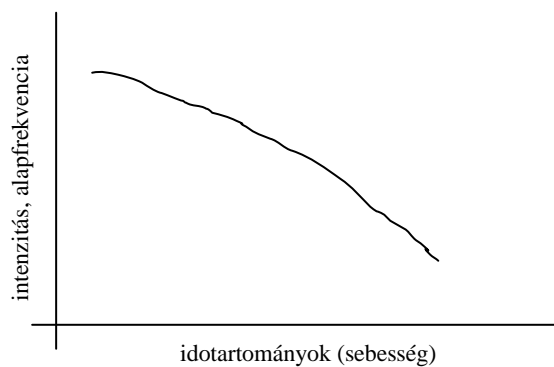
3.4.2.5. WTW



3.4.3. Néhány elvi probléma

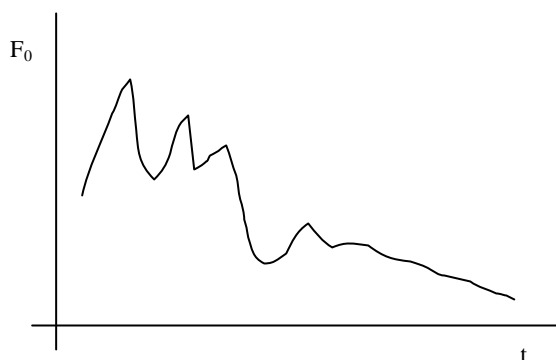
- alapvető feldolgozási egység:
- ember esetében a mondatnál nagyobb
- mondat egy gépi megoldásban: a prozódiai algoritmusok legmagasabb szintje is a mondat

Pl. Kijelentő mondat eső jellegű, ez azonban humán megközelítés, mérnöki módon hogyan fejezzük ezt ki



A prozódia három fő tulajdonsága

- intenzitás
- alaphékvencia
- idotartománybeli jellemzők



nem folytonos: F0-t csak a zöngés hangoknál tudjuk értelmezni, felpattanó hangoknál pl. nem

- absztrakciós szintek:
- hang szint (legalacsonyabb)
- szótag szint
- szó szint
- tagmondat szint (prozódiai fázis)
- mondat
- ezekre mérhető, fizikai paramétereket kellene találni
- beszédfelismerés kulcsterülete: a sokréteű szinteket megkülönböztetni, elválasztani (lehetőleg minél jobban) – és a beszédben ezek folytonosak.

3.4.4. Megoldási stratégiák

3.4.4.1. Szabályalapú

- lebontás: címkézés
 - ☞ mondat: kijelentő, kérő, felkiáltó
 - ☞ szó: alany, állítmány, tárgy, határozó, jelző; hangsúlyos/hangsúlytalan
 - ☞ szótag: hangsúlyos/hangsúlytalan
 - ☞ hang: szó eleje, szó közepe, szó vége; alacsony/mély hangrendű; magánhangzó/mássalhangzó; stb.
- szabályok megalkotása nagyon lassú – hiba esetén a szabályok kijavítása nehézkes
- a nyelv nem reguláris szerkezetű, vannak kivételek
- a nyelv változik, nem statikus

3.4.4.2. Gépi tanulás (machine learning)

- a gyakorlatból, mint nagy adatbázisból kinyerjük a szabályokat
- vegyünk sok, egymással összefüggésbe hozott, címkézett adatot – ez elég nagy adatbázis, címkézett szöveggel: neurális hálóval megvalósítva a rendszer következtetni tud
- a hosszú, absztrakciós szinten történő munkát kiváltjuk: sok adatban korrelációk, összefüggések keresése
- problémák
 - ☞ adatbázis létrehozása: több millió adatot kézzel kell felcímkézni (☹ hangsztig le kell menni)
 - ☞ a rendszer jól működik arra az adatbázisra, amelyre be lett tanítva, de a többire nincs garancia
 - ☞ milyen alapon ítéljük meg egy rendszer jóságát (Ezt a felhasználó dönti el!)

Az ember többnyire rosszul turi a minőség változását. Ha megszoktunk egy adott minőséget (még ha rossz is), nehezen viseljük, ha az megváltozik. Ebben van a szabályalapú rendszerek egyik nagy előnye: kiszámíthatóság.

A valóságban a két módszer ötvözetét használják. Egy adott, zárt problémakört fed le a gépi tanulás módszere, a kimaradó halmazra valamilyen szabályalapú megoldást alkalmaznak.

3.4.5. TTS tervezése: a hanganyag (akusztikai elmebázis)

- az akusztikai elmebázis nagyjából megfelel egy adott nyelv betukészletének (hangjainak)
- fonéma – graféma konverzió problémái:
- a fonéma minimálhalmaz, diszkrét elemek definíciója, a természetes beszéd azonban folytonos
- hangelemek (hangkód) – a fonémák kibovítése, pl. tájshólások e-jei, ng, rövid/hosszú magánhangzók
- ezek átmeneti tulajdonságait is figyelembe kell venni: megoldási ötlet

3.4.5.1. Diád (diphone)

- pl. legyenek az adott hangsorozat hangkódjai: 1,5,11,13. Ekkor vesszük a 1-t és az 5-t a közepéig, aztán az 5-t a közepétől a 11 közepéig, majd a 11-t a közepétől és a 13-mat.
- Minden elemet az alábbi módon definiálunk:

elemhatár – hanghatár – elemhatár

	1	2	3	4	5		11	12	13	...	50
1					x						
2											
3											
4											
5							x				
⋮											
11									x		
⋮											
50											

- Így az elem tartalmazza mindegyik hang bizonyos stabil szakaszát és az átmenetet.
- ez az ún. diád (diphone), ezzel csak az a probléma, hogy négyzetes a fonémák számával

Diád peremfeltétel: diádatáron a hangokban legyen folytonos átmenet (bizonyos esetekben megoldható: szünet + s; bizonyos esetben nem oldható meg: rövid magánhangzó, több hatás, elöl is és hátul is hat)

3.4.5.2. Triád (triphone)

- vegyünk hármass egységeket, a problémás elem legyen középen, az elemhatárt tegyük oda, ahol a vágás gond nélkül megtehető

elemhatár + 2 hanghatár

- a fonémák köbével arányos

3.4.5.3. További kiterjesztés

- azon elemeket, amelyek gyakran előfordulnak, tároljuk diádos/triádos elemekből
- változó méretű adatbázis (non-uniform database)

3.4.5.4. Adatbázis-elemek ábrázolása

- alapvetően aszerint, hogy a prozódiai módosításokat el akarjuk-e rajtuk végezni vagy sem, és ha igen, akkor milyen technikával
 - ☞ hullámforma: módosítás nélküliek
 - ☞ parametrikus forráskódolt: formáns, LPC

3.5. PSOLA (Pitch Synchronous OverLap Add) algoritmus

(zöngeszinkron átlapoló-összeadó)

3.5.1. Alapötlet

- tekintsük $s(n)$ -t, mint egy FIR szurot egy $i(n)$ impulzussorozatra
- ekkor szét kell szedni a jelet valamilyen pozícióban megjelenő impulzussorozatra és az azokra adott elemi válaszok összegére

$$i(n) = \sum_{k=-\infty}^{\infty} d(n - P_a(k)) \text{ és a válaszfüggvény: } h^{P_a(k)}(n) = s(n) \cdot w^{P_a(k)}(n)$$

- w : valamilyen ablakot elhelyezünk a hullámformán Legyen ez az ablakfüggvény:

$$w^{P_a(k)}(n) = \begin{cases} 0, & \text{ha } n < P_a(k-1) \text{ vagy } n > P_a(k+1) \\ 0.5 - 0.5 \cos\left(\frac{p(n - P_a(k-1))}{P_a(k) - P_a(k-1)}\right) & \text{ha } P_a(k-1) < n < P_a(k) \\ 0.5 - 0.5 \cos\left(\frac{p + p(n - P_a(k))}{P_a(k+1) - P_a(k)}\right) & \text{ha } P_a(k) < n < P_a(k+1) \end{cases}$$

- a függvény érdekessége, hogy nem kauzális

3.5.2. Szintézis:

$$\hat{s}(n) = \sum_{m=-\infty}^{\infty} i(m) \cdot h^m(n) \quad \text{és} \quad i(m) = \sum_{k=-\infty}^{\infty} d(m - P_a(k))$$

$$\hat{s}(n) = \sum_{m=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} d(m - P_a(k)) h^m(n) = \sum_{k=-\infty}^{\infty} h^{P_a(k)}(n) = s(n) \cdot \left(\sum_{k=-\infty}^{\infty} w^{P_a(k)}(n) = s(n) \right)$$

- Ha egyetlen ilyen létrejött az adatbázisban (megvannak a válaszfüggvények)
- alapfrekvencia változtatás: egy pitchmark sorozat változtatás: ha növelni szeretnénk az alapfrekvenciát, akkor közelebb hozzuk őket, így a válaszokat összeadogatva nagyobb lesz a frekvencia
 - ⊗ nagyon közel nem hozhatjuk őket, mert akkor egymásra csúsznak
 - ⊗ nagyon távol sem vihetjük őket, mert akkor a jel elhal két impulzus között, és ott csend lesz
- Ezzel az algoritmussal kb. 20%-ot lehet változtatni az alapfrekvencián.

4. BESZÉDFELISMERŐK

4.1. Bevezetés

4.1.1. Felismerési feladatok

- a gépi rendszer ismerje fel a beszédet
- beszéddetekció: annak felismerése, hogy beszéd van vagy nincs (sokszor része a beszédfelismerőknek, de önmagában is hasznos lehet¹)
- zöngés/zöngétlen meghatározás – leginkább csak támogatja a beszédfelismerőket, de néha önálló feladatnak is tekintik
- beszélő felismerés
- speciális esete a beszélő azonosítás
 - ☞ hagyományos módszer mintaszöveg felvétele, amit belépéskor el kell mondani
 - ☞ véletlenszerűen kisorsolt minta: sok mintát vesznek fel, és ezek közül egyet véletlenszerűen sorsol a rendszer a belépéskor

4.1.2. Beszédfelismerés osztályozása

- kis (kötött) szótáras, kb. 100 szó ↔ nagy szótáras (kötetlen szótáras), 20-80000 szó²
- személyfüggő ↔ személyfüggetlen
 - ☞ személyfüggő: egy személy beszédét ismeri fel, általában adaptív rendszer, egy adott személyre rátanul
 - ☞ személyfüggetlen: nagyon sok mintával dolgozik, a személyfüggőséget megpróbálja kiátlagolni
- izolált szavas ↔ kapcsolt szavas ↔ folyamatos beszéd
 - ☞ izolált szavas: egymástól hosszú idővel elválasztott szavak (pl. utasítások)
 - ☞ kapcsolt szavas: a szavak közti szünetek minimálisak
 - ☞ folyamatos beszéd: diktáló rendszerek
- jó minőségű beszédből felismerők ↔ robusztus rendszerek
- jó minőségű beszédből felismerők: mindig innen indul a felismerés, és valamilyen trükkel sikerül robusztussá tenni
- robusztus rendszerek: elég nagy zaj mellett is felismerik a beszédet³

4.1.3. Filozófiai probléma

- 1965-ig az a tendencia volt, hogy próbáljuk meg leutánozni, hogy mi történik az agyban a beszédfelismerés során
- 1965-ben új elmélet: az agyban született mondanivaló kódolása a beszéd, tehát a szájon keresztül kijövő hangból kell tudnunk következtetni a mondanivalóra
- nehézség: egyazon mondanivalónak végtelen sok reprezentációja van akusztikai szinten
 - ☞ nyelvi szinten ugyanazon dolgok akusztikai megjelenése különböző
 - ☞ a végtelen sokféle reprezentációk közül melyek azok, melyek 1 fonémának/átmenetnek tekinthetők és melyek nem

¹ Pl. zajos környezetben segélykiáltás detektálása

² Bizonyos nyelveken ez már gyakorlatilag diktáló rendszernek (STT: Speech to Text) tekinthető. Magyarban ez nem így van, mivel a magyar toldalékozó nyelv. Így egy adott szó felvétele nem jelenti a szó összes megjelenési formájának felvételét. Erre illusztratív példa, hogy míg az angolban (amely nem egy toldalékozó nyelv) egy 20000 szavas felismerő jó minőségű, addig egy ugyanennyi szót tartalmaz magyar nyelvű felismerő kb. az időjárás jelentést tudná viszonylag nagy biztonsággal felismerni.

³ Pl. harci helikopterekben használnak beszédfelismerőket bizonyos parancsokra (izolált szavas, de nagyon robusztus)

4.1.4. A beszédfelismerés három komponense

- lényegkiemelés (feature abstraction)
 - ☞ a hanghullám változásából olyan elemeket próbálunk kiemelni, melyeknek kicsi az intrinviduális és az interindividuális jellemzője (függetlenül attól, hogy ki mondta, milyen érzelmi állapotban mondta)
 - ☞ redukáljuk az adatmennyiséget (kb. 1/10 részére csökkenthető)
- mintaillesztés
 - ☞ ugyanazt a szót nem lehet kétszer ugyanabban a ritmusban kimondani, ezért a mintaillesztés legfontosabb feladata a különböző ritmikájú kiejtések közötti ritmuskülönbségek kiküszöbölése
- utó/előfeldolgozás
 - ☞ utófeldolgozás: biztosabbá teszi a felismerést
 - ☞ előfeldolgozás: könnyebbé teszi a felismerést

4.2. Lényegkiemelés

- az időfüggvényt keretekre (ablakokra) bontjuk, ezek 10-30 ms hosszú ablakok
- az ablakokat 50%-os fedésben helyezük egymásra

4.2.1. Az ablak alakja

4.2.1.1. *Rectangular (négyzetlejtés)*

4.2.2. Hamming ablak

Ha a beszéd időfüggvénye $f(t)$ és az ablak időfüggvénye $w(t)$, akkor a kiablakolt függvény $a(t)=f(t) \cdot w(t)$. A spektrális jellemzésre igen jó a Fourier transzformáció:

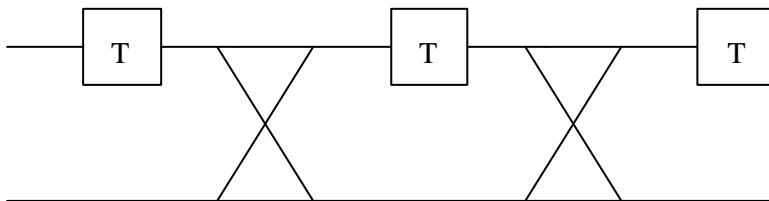
$$A(\omega) = F(\omega) * W(\omega) = \int_{-\infty}^{\infty} F(a)W(\omega - a) da$$

Olyan ablakfüggvényt kell választanunk, amelynél a kiablakolt beszéd spektrálisan legjobban hasonlít az eredetire, azaz az ablakfüggvény legjobban hasonlít a Dirac függvényre: $\delta(t)$. A Hamming ablak bizonyos értelemben jobban hasonlít a Dirac függvényre

- ☺ ω_0 -tól távol nagyon gyorsan lecseng, kevés összetevőt vesz figyelembe
- ☹ a közvetlen ω_0 melletti összetevők jobban befolyásolják a kiablakolt jelet¹

4.2.3. Lényegi jellemzők

- szegmens-energia
- F_0 , amennyiben a hang zöngés
- nullátmenetek száma a szegmensen belül (beszéddetekciónál fontos)
- LPC paraméterek (w_i)
- Gráfstruktúrában (rácsszerkezetű gráffal) meg lehet határozni a kódolót és a dekódolót is



¹ Az a sáv amelyen belül simít, pszichoakusztikailag nem zavaró, mert a fül is simít, és ez a kritikus sávon nem nyúlik túl

- $k_i \leq 1$ – ezek az ún. parcor együtthatók, melyek egyre több korrelációt vonnak kis a jelből, végül teljesen korrelálatlan lesz
- Más megoldás, ha valaki leírja a vokális traktust, mint egy változó keresztmetszetű csövet, ekkor kapjuk az ún. área együtthatókat.
- $r_i = \frac{A_{i+1} - A_i}{A_{i+1} + A_i}$, és érdekes módon $r_i = -k_i$
- A beszédhullámokból a beszédkeltőre lehet következtetni, amikor a parcor együtthatókat határozzuk meg.

4.2.3.1. Spektrális jellemzők

- általában minták állnak rendelkezésre, ezért DFT-t (diszkrét Fourier transzformációt) alkalmazunk.
- a DFT a halmozott spektrum mintáit adja (számsorozat DFT-je a számsorozathoz tartozó spektrum elégséges mintáit adja)
- 1965-ben Cooley és Tuckey felfedezték az FFT-t (fast Fourier transform)¹

4.2.3.2. Kepsztrális jellemzők

1962-ben Bogert, Healy és Tuckey észrevette, hogy a beszédet nagyon jól jellemzi a kepsztruma, amelyet az alábbi módon definiáltak:

$$c(q) = \left| \mathcal{F}^{-1} \left\{ \log(|\mathcal{F}\{a(t)\}|^2) \right\} \right|^l \quad \text{ahol } l=1,2$$

- $c(q)$ a kepsztrum, idő jellegű mennyiség (de nem idő)
- konvolváltság időfüggvények kepsztruma
 - ☞ természetes beszéd: gerjesztés * vokális traktus súlyfüggvénye
 - ☞ rögzített beszéd: beszéd * mikrofon súlyfüggvénye
- ha sikerül különválasztani a kepsztrum segítségével a vokális traktus és a beszéd súlyfüggvényét, akkor a beszéd főbb jellemzőit különválaszthatjuk.

4.2.3.3. Kepsztrum analízis

$$c(q) = \left| \mathcal{F}^{-1} \left\{ \log(|\mathcal{F}\{a(t)\}|^2) \right\} \right|^l \quad \text{ahol } l=1,2$$

- $a(t)$ pedig a kiablakolt időfüggvény: $a(t)=f(t)w(t)$
- a logaritmalás lényege, hogy ha $a(t)$ egy eleme nem szorzat, hanem konvolúció, akkor össze lehet őket adni.
- az inverz Fourier transzformáció hatására idő jellegű térbe transzformáljuk a jelet, melynek változója a frekvencia.

Legyen $a(t) = x(t) * y(t)$, ahol $x(t)$ a zöngé, $y(t)$ pedig a vokális traktus súlyfüggvénye

$$\mathcal{F}\{a(t)\} = A(\omega) = X(\omega)Y(\omega) \quad /^2, \log$$

$$\log |A(\omega)|^2 = \log |X(\omega)|^2 + \log |Y(\omega)|^2$$

$$\mathbf{g}_a(q) = \mathbf{g}_x(q) + \mathbf{g}_y(q)$$

$$C_f(q) = (\mathbf{g}_x(q) + \mathbf{g}_y(q)) \cdot (\bar{\mathbf{g}}_x(q) + \bar{\mathbf{g}}_y(q)) = |\mathbf{g}_x(q)|^2 + |\mathbf{g}_y(q)|^2 + \mathbf{g}_x(q)\bar{\mathbf{g}}_y(q) + \mathbf{g}_y(q)\bar{\mathbf{g}}_x(q)$$

- az első két tag pedig $c_x(q)$ és $c_y(q)$, azaz a zöngé és a vokális traktus kepsztruma

¹ Az FFT DFT esetén egy n elemű számsorozat, és míg DFT-hez n^2 szorzás és n^2 komplex összeadást kell végezni, míg az FFT-hez $\frac{n \cdot \log_2 n}{2}$ szorzást és $\log_2 n$ összeadást szükséges.

Ha x kepsztruma valamilyen kefrencia tartományban domináns, és y egy másikban, akkor a kettejük szorzata éppen 0-t ad. Ha a tartók elkülönülök, akkor a kepsztrum a két jel kepsztrumainak összege. Ilyen esetben a kepsztrum analízis dekompozícióra alkalmazható.

- Lényegkiemelést alkalmazva keretenként egy paramétervektort kapunk
- 10/20-adfokú LPC analízist használva 10/20-adfokú kepsztrumot kapunk

A paramétervektort az alábbi módon származtatjuk

$$p_t = \begin{bmatrix} p_i^t & i = 1..l & \text{primer paraméterek} \\ p_i^t - p_i^{t-1} & & \text{delta paraméterek} \\ \sum_{j=2}^2 a_j \cdot p_i^{t+j} & & \text{delta paraméterek lineáris kmbinációj a} \end{bmatrix}$$

- sokáig azzal próbálkoztak, hogy a paraméterek számát növelték, később azonban kiderült, hogy bizonyos paraméterszám után ez nem ad többjellemzot
- egy szegmensbe ne csak a saját paramétereit vegyük bele, hanem az elozo szegmensbol is néhány paramétert
- a beszédet az a változás is jellemzi, hogy mi kerül a következő szegmensbe
- a p_i^t , p_i^{t+1} és p_i^{t-1} segítségével lineáris kombinációval megkapható
- a mai kísérletek nagy rész arra irányul, hogy a vektorparamétereket milyen súllyal vegyük figyelembe ahhoz, hogy optimális legyen

4.3. Mintaillesztés alapjai

Adottak pl. izolált szavak (lényegkiemelt vektorsorozatokkal), kefrenciák, prototípusok (r_k vektorsorozatok) illetve fonémák (ezekhez is vektorsorozatok tartoznak), vagyis a felismerés alapjául szolgáló nyelvi egységet reprezentáló vektorsorozatok. Feladat, hogy az ún. tesztkiejtésbol (felismerhető kiejtés) meg tudjuk állapítani, hogy melyik referenciához hasonlít a legjobban.

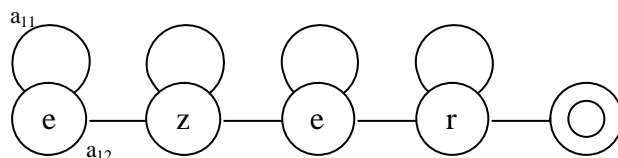
Legnagyobb probléma, hogy ugyanazt a szót az emberek különböző ritmusban képesek kiejteni, de ugyanez igaz egy embernél ugyanazon szó kétszeri kiejtésénél. Meg kell tehát találnunk az a technikát, amivel a megfelelő dolgok lesznek összeillesztve.

Erre három módszer létezik, ebbol kettonek statisztikai megfigyelés az alapja, a harmadik sablon (template) alapú.

- HMM – Hidden Markov Model (statisztikus)
- ANN – Artificial Neural Network (statisztikus)
- DTW – Dynamic Time Warping (sablon alapú)

4.3.1. Mintaillesztés HMM-mel

Ritmikai változások figyelembevétele pl. izolált szavas beszéd felismerésnél:



A modell: lépni kényszerül minden 10. ms-ban, de nem kényszerül ellépni onnan. Így ezzel a technikával alapvetően ki lehet küszöbölni az alapvető ritmusbeli különbségeket. Lehetséges ugró él is, ha valamelyik hangot nem ejtjük ki.

Miért hívják ezt rejtett Markov modellnek? Azért, mert a véges automatáknál megszokott módtól eltérően itt nem tudjuk, hogy a folyamat milyen állapotban van. Erre a megfigyelésből kell következtetnünk.

A modell kiad P_T vektort, miközben az állapotokban eljut az N . állapotig. A megfigyelési sorozatot O -val jelöljük (mint observation). Aközben az emissziókból (azon vektorok, amelyeket a Markov folyamat emittál) nem tudjuk megállapítani, hogy melyik állapotban vagyunk.

Egy állapothoz sokféle vektor tartozhat, ezért inkább valószínűségekkel számolunk.

$P(O_t/q_j)$ egy valószínűsége jellemző érték, ahol q_j a j -edik állapot O_t pedig folytonos értékészlet. Minden egyes ponthoz sűrűségfüggvény-értéket rendelünk. Ezt a sűrűségfüggvényt adatbázisokból kell meghatározni. Az egyszerűség kedvéért csak a Gauss-eloszlások lineáris kombinációit tekintjük.

$$b_j(O_t) = \sum_{m=1}^M C_{jm} \cdot G(O_t, \mathbf{m}_{jm}, \mathbf{s}_{jm}), \text{ ahol } G(x, \mu, \sigma) \text{ egy } \mu \text{ várható értékű } \sigma \text{ szórású Gauss eloszlás.}$$

Szemléletesen ennek az a jelentése, hogy a j . állapotban mennyi a valószínűsége, hogy O_t megfigyelés adódik. Az adatbázisok alapján meg kell határozni az átmeneti valószínűségeket. A tanítás során:

- e_1 hozzárendeljük az első Markov folyamatot minden felismerhető elemhez
- e_j $b_j(O_t)$ $j=1,2,\dots,N, \{a_{ij}\}$

Végül minden e_i -hez tartozik egy $b_j(O_t)$ készlet.

Felismerés: érkezik egy tényleges megfigyelést tartalmazó vektor, ekkor ki kell keresni egy olyan Markov láncot, amely a legnagyobb valószínűséggel tudja ezt követni. Meg kell nézni, hogy az egyes Markov modellek milyen valószínűséggel képesek ezt az O_t -t produkálni, és a legnagyobb valószínűségű modell által leírt nyelvi elemet tekintjük felismernek.

A Baum-Welch algoritmus egyértelműen meghatározza a c , μ , σ -t, de adott N esetében. A Viterbi algoritmus különböző Markov modellekre is alkalmazható.

- feltételezve O_t -t létezik egy legnagyobb valószínűségű (optimális) út a kiindulástól i -be
- $\delta_t(i)$ legyen ezen út valószínűsége
- így minden pontra meghatározható az odavezető optimális út valószínűsége
 1. inicializálás: $\mathbf{d}_0(i) = 1$
 2. indukciós lépés: $\mathbf{d}_{t+1}(j) = \max_i \{ \mathbf{d}_t(i) \cdot a_{ij} \} \cdot b_j(O_{t+1})$, és $i = 1 \dots N, j = 1 \dots N$
 3. leállás feltétele: $t = T$
 4. kimenet: $\max_i \{ \mathbf{d}_T(i) \}$ és $i = 1 \dots N$
- végül megkapjuk, hogy melyik az a Markov modell, amelyik a legvalószínűbb.
- izolált szavak esetén a modellhez egyértelműen tartozik a felismert szó
- Vannak olyan esetek, amikor az is fontos, hogy mi volt az optimális út, nem elég a célállapot. Ilyenkor minden pontba azt az értéket írjuk, amelyik állapotból jöttünk, azaz $\arg \max_i \{ \mathbf{d}_t(i) \cdot a_{ij} \}$ -t. Így T időpontból visszafelé az optimális út meghatározható.

4.3.2. Dinamikus idovetemítés (Dynamic Time Warping)

- pl. izolált szavas felismerésénél, itt ugyanis nincs lehetőség statisztikus felismerésre
- felvesszük a felismerendő szöveget és a referenciamintára próbálunk illeszteni
- a legjobban hasonlító minták a legkisebb távolságra vannak a referenciától – definiálni kell tehát egy d távolságot
- Legyen
 - ☞ \vec{t} tesztvektor
 - ☞ \vec{r}_k referenciavektor, $k= 1..M$
 - ☞ és így $\underset{\vec{r}_k}{arg \min}\{d(\vec{r}_k, \vec{t})\}$

Az elobb említett távolságot az alábbi módokon definiálhatjuk (általában az első kettőt szoktuk használni):

$$d(\vec{r}_k, \vec{t}) = \left\{ \begin{array}{l} \sum_{i=1}^M (r_k(i) - t(i))^2 \\ \sum_{i=1}^M |r_k(i) - t(i)| \\ \dots \\ \max_i |r_k(i) - t(i)| \end{array} \right\}$$

Probléma akkor van, ha \vec{t} és \vec{r}_k vektorok különböző hosszúak. Ez abból adódhat, hogy egyrészt a különböző bemondások különböző sebességgel történnek, másrészt ugyanazon bemondáson belül sebesség-ingadozás is felléphet.

Az egyszerű lineáris vetemítés akkor lenne jó, ha kizárólag sebességkülönbségről van szó (pl. az ember is nyávogna, ha gyorsabban beszél, mint a magnóra felvett szöveg, ha gyorsabban játsszuk le).

Ezért nemlineáris vetemítésre van szükség.

- legyen $F_i(t)$ egy vetemítögörbe
- ezek kijelölnek bizonyos tartományt a $\vec{t} - \vec{r}$ síkon (bizonyos irányokra majd nem is lesz szükség, úgyhogy ezek majd bizonyos korlátokat szabnak)
- a $d(\vec{r}, \vec{t}) = \min_{F_i} \{d(r(t), t(F_i))\}$ távolságot úgy kell meghatározni, hogy az adott vetemítögörbére vett távolság a legkisebb legyen
- a vetemítögörbe tulajdonságai
 - ☞ monoton növekszik
 - ☞ lokális korlátok (extrém körülmények között túlléphető)
 - ☞ teljes optimum: lokális optimumokon keresztül valósul meg¹

¹ Bárhol megállva a vetemítögörbén, addig a pontig vezető út optimális kell, hogy legyen

Példa

$$\vec{t} = \{2,6,8,9,8,3\}$$

$$\vec{r} = \{1,6,9,6,5\}$$

$$d := \sum_i |r(i) - t(i)|$$

1. lokális tulajdonságok kiszámítása

r ₆	2	3	6	3	2
r ₅	7	2	1	2	3
r ₄	8	3	0	3	4
r ₃	7	2	1	2	3
r ₂	5	0	3	0	1
r ₁	1	4	7	4	3
	t ₁	t ₂	t ₃	t ₄	t ₅

2. akkumulált tulajdonságok kiszámítása

r ₆	30	11	9	6	6
r ₅	28	8	3	4	7
r ₄	21	6	2	5	8
r ₃	13	3	2	4	7
r ₂	6	1	4	4	5
r ₁	1	5	12	16	19
	t ₁	t ₂	t ₃	t ₄	t ₅

Az eljárás elonyei

- ☺ gyorsan tanítható, gyakorlatilag 1 mintával is valamilyen szinten működhet
- ☺ a minták a vetemítőfüggvény mentén átlagolhatók

4.3.3. Nyelvi modell integrálása HMM illesztéséhez

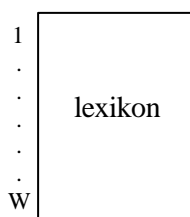
ha nem csak izolált szavakat akarunk használni

Alapötlet: az izolált szavas HMM Markov modelljeit összevonjuk, így s rendszerbe nyelvtani információt viszünk a modell topológián keresztül

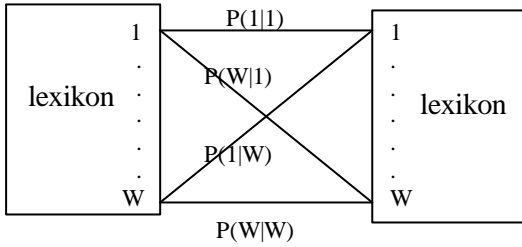
- ☹ a természetes nyelvek nem írhatók le determinisztikus nyelvtannal
- ☹ méret is probléma: sztochasztikus n-gram modellek (n=1,2,3): milyen valószínűséggel fordul elő n adott szó egymás után
- ☹ a szöveg nem ugyanaz, mint a hangsor, tehát kiejtésmodellezés is szükséges

Legyen W a szótárméret

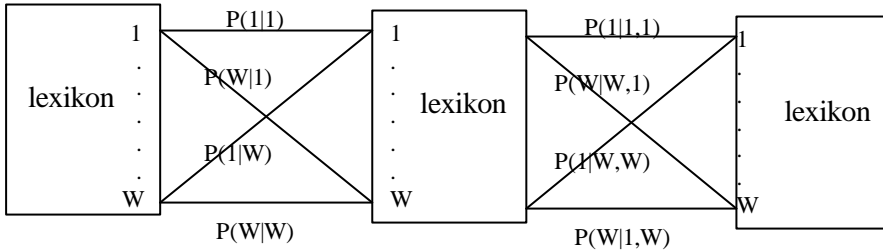
1-gram modell



bigram modell



trigram modell



Felmerül az optimalizálási igény: $N = 20000$, 20 modellállapot (szó), akkor a HMM állapotok szám: $20000 \cdot 3 \cdot 20 = 1.2$ millió, súlyok: W^3 -bel arányosak.

Keresés optimalizálása:

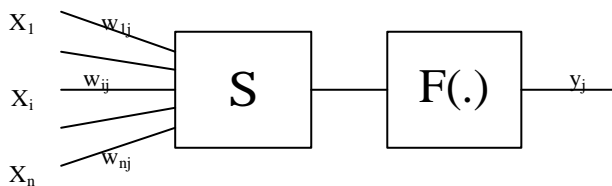
- keresési igény megadása (beam search)
- token passing (útvonal + útvonal valószínűsége), teljesen értelmetlen irányok elhagyása

4.3.4. Mesterséges neurális háló (ANN)

izolált szavas beszéd felismerésre alkalmas, arra bevált

4.3.4.1. Mesterséges neuron

A j-edik neuron

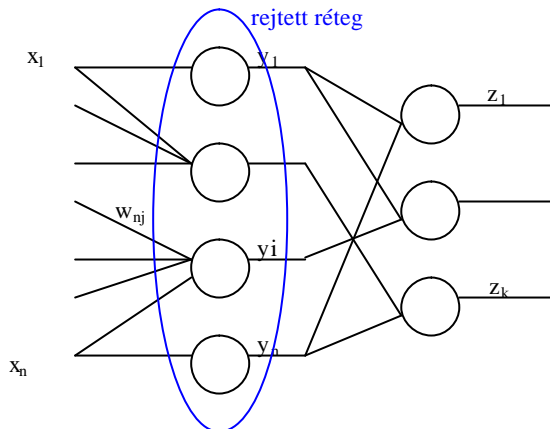


w_{ij} : a j-edik neuron i-edik bementének súlyvektora (weighting factor)

$$y_j = F\left(\sum_{i=1}^n w_{ij} \cdot x_i + \Theta_j\right) = F(\mathbf{a}), \text{ ahol } \Theta_j \text{ a j. neuronra vonatkozó küszöbérték.}$$

4.3.4.2. Egyszerűsített Multi-Layer Perceptron, 1 rejtett réteg

Minden bemenet hathat akármelyik azt követő neuronsorozat bármely elemére, és minden kimenet hathat az azt követő neuronsorozat bármely elemére.



Az általános neuronhálótól abban különbözik, hogy a neuronok rétegekbe rendezettek, nem tetszőlegesen összekötöttek, nincs visszaút, egy réteg mindig csak egy magasabb rétegbe mehet, tehát a neuronháló, mint irányított gráfot tekintve, ez a gráf DAG. Ezenkívül minden rétegben ugyanazt az F függvényt alkalmazzák (beszédfeldolgozásban rétegenként nem alkalmaznak külön függvényt, hanem az egész hálóban 1 F függvényt használnak)

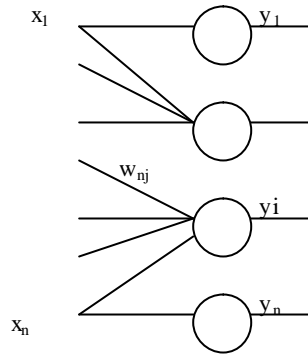
4.3.4.3. Beszéd felismerés ANN-nel: példák

1. $x = (x_1, x_2 \dots x_N)$ legyen az egy kerethez tartozó tulajdonságvektor: $z_k \approx 1$, ha x abból a hangból származik, amit mi a k -adik sorszámmal látunk el, minden más esetben egy 0-hoz közeli érték. Így keretenként fonémákat akarunk felismerni
2. $x = (x_1, x_2 \dots x_l \mid x_{l+1}, x_{l+2} \dots x_{2l} \mid x_{2l+1}, x_{2l+2} \dots x_{3l} \mid \dots \mid x_{4l+1}, x_{4l+2} \dots x_{5l})$, $N = 5l$, ezek pedig a $t-2, t-1, t, t+1, t+2$ keretből származó tulajdonságvektorok: $z_k \approx 1$, ha ez olyan szótagból származik, amelynek közepén t fonéma van, minden más esetben 0. Az ilyen perceptronokat időkéleltetéses perceptronnak nevezzük.
3. sok bemenetet megengedve (pl. izolált szavas számfelismerő): 30-40 keretnyi bemenetet engedünk meg pl. 1000 neuron bementre, ez 10 kimeneten jön ki. A rejtett rétegben kb. 60-70 neuron található. Ez igen gyakori alkalmazása a MLP-nak.
4. előfeldolgozásra is használják: a fonémákat osztályokba sorolják (magánhangzó, félmagánhangzó, afrikáta, felpattanó, zárrés), és ezzel a szegmentáláson kívül a szegmensekhez a fonetikai osztályt is hozzárendeli. Ha az elotte/utána lévő is megmondja, akkor címkéz. A címkézés finomságát lehet változtatni. Tehát izolált szavakat, szótagokat NN-nel fel lehet ismerni, de magasabb szintű elemzéshez előkészítésnek is alkalmazható.

4.3.4.4. A működés lényege: w_{ij} -k és F helyes megválasztása

- ezekre nem ismerünk analitikus módszert
- manapság w_{ij} -ket tanítással határozzuk meg
- nyelvi elemekkel bombázzuk a neurális hálót, és w_{ij} -ket úgy változtatgatjuk, hogy jól működjön
- matematikailag bizonyítható, hogy ha az N dimenziós vektor olyan osztályokból származik, mely osztályok az N dimenziós térben – akár nem lineáris területtel is – elválnak, akkor a kimenet meghatározható.

4.3.4.5. Tanítás szemléltetése egyrétegu, nemlinearitás-mentes hálózatokra



x^P – a P-edik tanításkor bemenő sorozat
 y^P – a P-edik tanításkor megjelenő tényleges kimenet
 t^P – a tanítás során elvárt kimenet (target)

$y_j^P = \sum_i w_{ij} x_i^P$, ehhez kell definiálnunk egy hibát és egy súlytényezőt, amivel korrigálni fogunk.

A hiba: $E = \sum_P \left[\frac{1}{2} \cdot \sum_j (t_j^P - y_j^P)^2 \right]$ $E^P = \frac{1}{2} \cdot \sum_j (t_j^P - y_j^P)^2$, E^P pedig a P-edik bemenéssel kapcsolatos hiba.

Optimalizálandó faktorok az y -okban lévő w_{ij} -k. Gradiens módszert alkalmazva a súlytényezők:

$$\Delta w_{ij} = \mathbf{h} \left(- \frac{\partial E}{\partial w_{ij}} \right), \text{ ahol } \eta \text{ egy alkalmasan választott konstans } 0.02 \text{ és } 0.1 \text{ között.}^1$$

$$\Delta w_{ij} = \mathbf{h} \left(- \frac{\partial E}{\partial w_{ij}} \right) = \mathbf{h} \left(- \frac{\partial}{\partial w_{ij}} \sum_P E^P \right) = \sum_P \mathbf{h} \left(- \frac{\partial}{\partial w_{ij}} E^P \right) = \sum \Delta w_{ij}^P$$

miel

$$\frac{\partial E^P}{\partial w_{ij}} = \frac{\partial E^P}{\partial y_j^P} \cdot \frac{\partial y_j^P}{\partial w_{ij}},$$

hiszen y_j^P az w_{ij} függvénye. Alkalmazva a parciális deriválás tulajdonságait: E^P parciális deriváltja csak az adott j esetében nem lesz 0.

Tehát a $t_j^P - y_j^P = -\mathbf{d}_j^P$, y_j^P parciális deriváltja pedig éppen $i=j$ esetben nem 0. Vagyis

$$\Delta w_{ij}^P = \mathbf{h} \cdot \mathbf{d}_j^P x_i^P$$

azaz

$$\Delta w_{ij} = \sum_P w_{ij}^P = \mathbf{h} \sum_P \mathbf{d}_j^P \cdot x_i^P.$$

¹ Ha η túl nagy, akkor oszcillál a kívánt érték körül, ha túl kicsi, akkor pedig nagyon lassan konvergál a kívánt értékhez.

4.4. **Beszélofelismerés**

4.4.1. **Bevezetés**

4.4.1.1. ***Ki volt a beszélő?***

- megállapítható-e az elhangzó beszéd alapján a beszélő személye ha ismerjük az illetőt, illetve ha nem

4.4.1.2. ***Feltételezés***

- az agyunkban létrejövő neurális spektrogram tartalmazza a beszélő ismérveit
- vajon ez az információ, ezek a paraméterek olyan mértékben jellemzőek-e a beszélőre, mint pl. az ujjlenyomat?

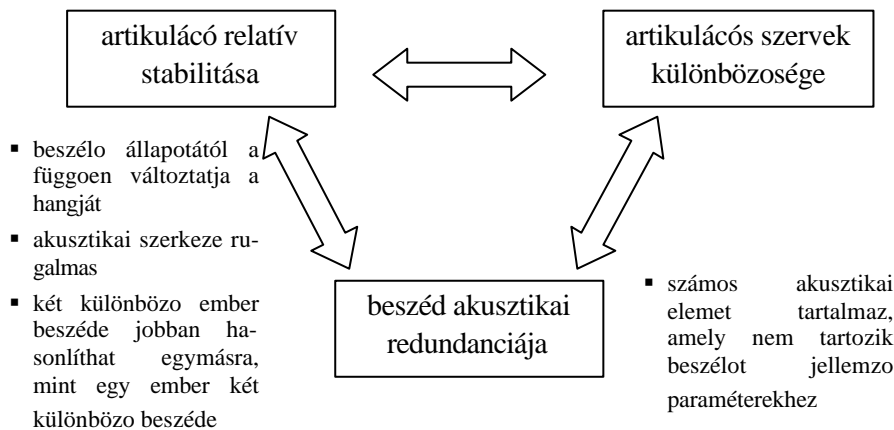
4.4.1.3. ***Alkalmazás***

- a beszélő személy vizuális megjelenése nem fontos (beléptető rendszerek, információs rendszerek, telefonos lekérdező rendszerek)
- kötött szótáras esetben létezik 100%-os felismerés, itt a beszélőnek érdeke a felismerés
- kriminalisztikai megközelítés
 - ⊗ itt a beszélő érdekeivel ellentétes az, hogy őt felismerjék

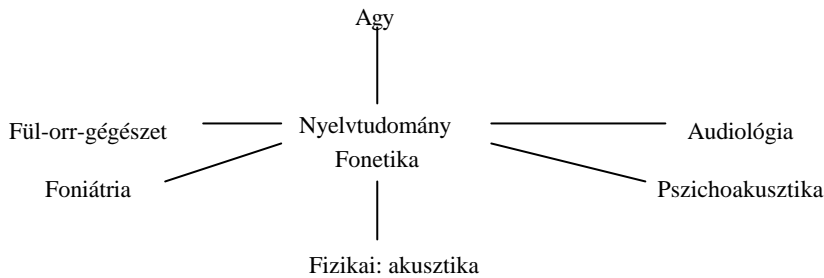
4.4.1.4. ***A beszélőfelismerésnek alapvetően két iránya van***

- az n lehetőségből kizárható-e nx személy
- az n lehetőségből melyik az nx esemény
- a két kombinációja: benne van-e, és ha igen, ki lehet?

4.4.1.5. ***A beszélőfelismerés paradoxona***



4.4.1.6. A témához kapcsolódó tudományágak:



4.4.2. A fonetikán belül: hangszínezet (ez alapján azonosítja a beszélot)

- milyen mértékben jellemző az emberre a hangja, beszéde (feltételezés: teljes mértékben)
- miképpen határozható meg az egyéni hangszínezet
- mely beszédképzési konfigurációval mutatja a legszorosabb kapcsolatot
- a zöngé, a toldalécszó, artikulációs mozgások vagy az összes együtt
- Miként fejezhető ki: akusztikai-fonetikai, percepció-fonetikai vagy mind együtt
- mindennapi életben: szubjektív benyomások a hangról (minden embernek vannak bizonyos elvárásai, amik csak akkor tudatosulnak, ha valami nem stimmel, nem illik a környezetbe), a hallgatóknak bizonyos feltételezései vannak a beszélőről

4.4.3. Beszédprodukció

- makrotervezés
- gondolati folyamat ↔ ismeretek, fogalmak
- pragmatika (hol mondom)
- szemantikai szerkezet, jelentés (milyen kifejezéseket válogatok ki) – stratégiák
- mikrotervezés
- szintaktikai szerkezet (milyen sorrendben, hogyan mondom őket) – transzformációs szabályok
- felszíni szerkezet – fonológiai szabályok
- fonetikai szerkezet
- artikulációs működések
- akusztikai hullámforma

4.4.4. További problémák

- egyetlen beszélő egyetlen hangja laboratóriumi körülmények között és egyéb tényezők közül többet is elhanyagolva is nagyon sokféle lehet
- felpattanó zárhangok esetén a zöngésedési idő (VOT [ms]) nagyon jellemző az adott mássalhangzóra (p, t, k), azonban a határértékek nagyon széles tartományokba esnek
- a spontán beszédben eltűnik a szó, a magánhangzók nagy rész 'svá', azaz 'ö' lesz a kiejtés pongyolasága miatt (nem tökéletes hangképzés)
- az alaphang magasságának változása: a koraival változik
 - ☞ a nőké felnttkorig kicsit mélyül, de alapvetően nem változik, idosebb kora jobban mélyül
 - ☞ a férfiaké kamaszkorban nagyon mélyül, felnttkorban egész mély, majd idosebb kora újra magasodik
- a beszéd akusztikájának relatív állandósága
 - ☞ egy nap elteltével jobb tényezők tekintetében a két kiejtés akusztikai jegyei között nincs különbség, de a részletekben már nagy eltérés mutatkozhat
 - ☞ gyermek és felnttkorban láthatóan eltérő
 - ☞ 20 év alatt az invariáns elemek nem változnak, de a hangszínkép között jelentős különbségek vannak
 - ☞ fiatal és idős között még nagyobb különbségek

- ☞ az érzelmek is befolyásolják (öröm, bánat, stb., de ez a kettő azonosítható a legjobban)
- ☞ a prozódia jellemzőbb, mint a szegmentális szerkezet: alaphang-magasság, tempó, intenzitás, szünetstruktúra, ritmus, artikulációs változás
- ☞ alkoholos állapot: bizonyos szintig nincs hatása a beszédre, a mennyiség egyénfüggő
 - beszédtempó lassul
 - alaphang magasabb lesz
 - a szünetek előfordulása és időtartama nő
 - a beszéd intenzitása nő
 - az artikulációs mozgások elnagyolódnak
 - a beszéd tempójának változása
- ☞ a beszéd tempójának változása: meghatározó a város, környezet.
 - 1869 óta vannak erre adatok: 26-ról 1995-ig 65 szó / percre növekedett a beszéd sebessége
 - ami a beszédgyorsulást meggátolja: a hallgató megértése
 - ha automatizálni szeretnénk a beszélőfelismerést, az időviszonyok változása gondot jelenthet: két beszédminta tempója nem azonos, akkor most ugyanaz a beszélő volt-e vagy sem, továbbá befolyásolhatja a beszéd sebességét az érzelmi állapot és legfőképpen a zajviszonyok.
- környezeti zaj hatása a beszédre
 - ☞ beszédtempó gyorsulása
 - ☞ alaphangmagasság nő
 - ☞ intenzitás nő
 - ☞ monoton jellegű lesz a beszéd (moduláció csökken)
 - ☞ hangsúlyhibák gyakoribb előfordulása
 - ☞ ejtéshibák (szegmentális, szupraszegmentális), pongyola kiejtés, hanghibák
- tempó
- artikulációs tempó: mennyi hasznos beszédjel esik egy adott időtartamra (beszédképzésre fordított idő)
- beszédtempó: a nem hasznos beszédjeleket is beleértve (megakadás, szünet, 'ö', ismétlés, levegővétel)

Olyan esetek is elképzelhetők, amikor a teljes beszédből csak 300 – 3500 Hz-re sávkorlátozott spektrum áll rendelkezésre (tipikusan telefon)

- redundancia csökken
- hiányoznak bizonyos invariánsok
- sokkal zajosabb a jel

4.4.5. Használt akusztikai vizsgálatok

- az [e] formánsszerkezeteinek különbözősége: alsóbb indexszámú formánsok között nincs különbség, de a magasabbakban igen, továbbá a formánsok sáv szélessége is mutathat változásokat
- LPC analízis jó segítség lehet
- felhang szerkezete (telefon és studiofelvétel esetén, de a teljes spektrum alapján nem lehet eldönteni)

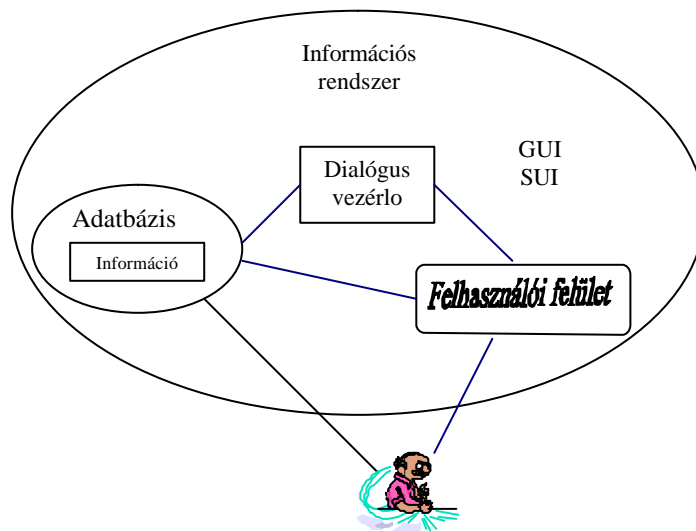
4.4.5.1. Akusztikai- fonetikai paraméterek a beszélő felismerésben

- központi formánsfrekvenciák és átmenetek
- maximumpontok
- formán sáv szélessége
- r és zárhangok zörejfrekvenciái
- sajátos spektrografikus alakzat
- felhangszerkezet
- magánhangzók időtartama
- beszéd és artikulációs sebesség
- csendes szünetek száma
- dallammenet
- hangsúlyozás
- egyéni ejtési sajátosság akusztikai tükröződése
- áthangolt spektrum (LPC analízis)
- Mi a helyzet a torzítással (amikor szándékosan el akarja változtatni a hangját)
 - ☞ az átlagember erre képtelen
- életkor, természet, súly meghatározása hang alapján (életkor elég jól meghatározható, természet kevésbé, de a súly nem jellemző)
- olyan beszédtorzítás, ami nem ismerhető fel (jelenleg): suttozás

5. BESZÉDINFORMÁCIÓS RENDSZEREK

- alapvető kérdés: mi lesz a rendszer célja, kinek készül a rendszer, ki fogja használni
 - ☞ a rendszert alapvetően a felhasználónak készítjük
- a beszédinformációs rendszer elemi építőkockákból épül fel
 - ☞ mit rakunk össze: milyen elemeket használunk, jó elemek-e ezek
 - ☞ hogyan rakjuk össze: az építőelemeket megfelelően rakjuk-e össze

5.1. Beszédinformációs rendszer felépítése

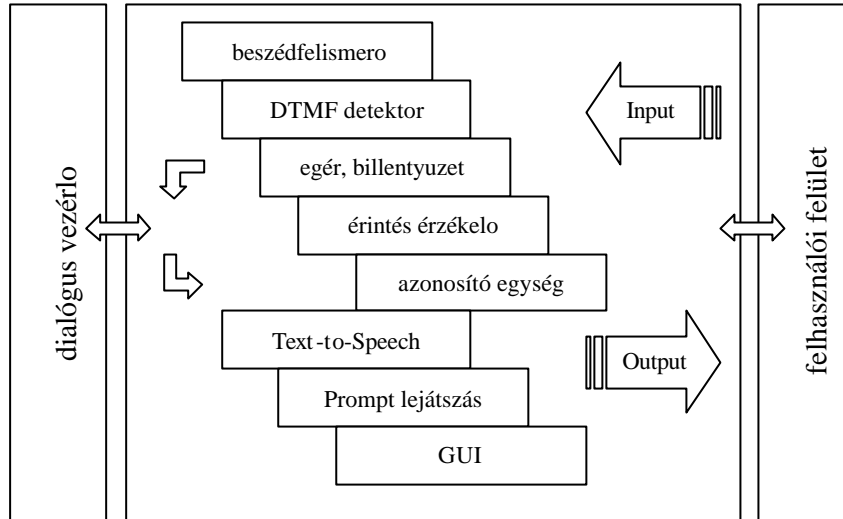


GUI – graphical user interface
 SUI – speech user interface

- valamilyen információt valahogyan el kell juttatnunk a felhasználóhoz

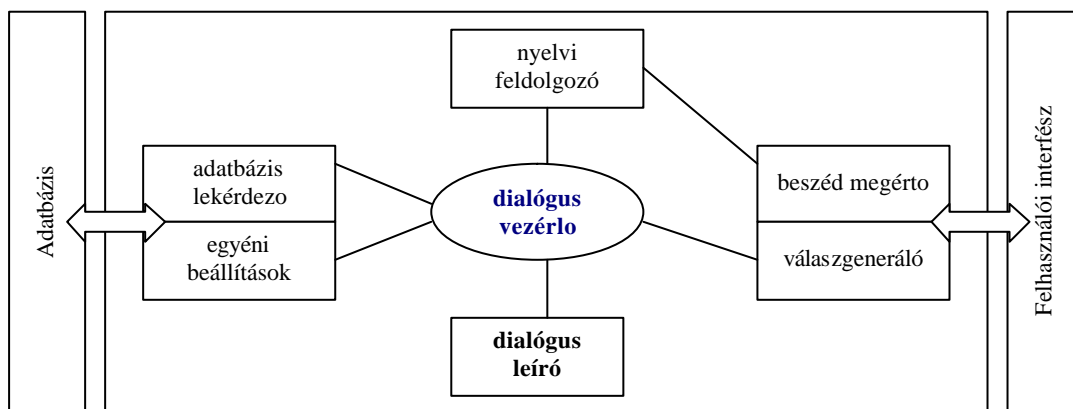
- a dialógusvezérlo arra szolgál, hogy a felhasználót rávegye, hogy hogyan érje el ezt az információt
- a felhasználónak hozzá kell férnie az adatbázishoz, viszont a közvetlen adatbázis-hozzáférés leszukíti a felhasználók körét

5.1.1. A felhasználói felület



- a kimenetek beszédinformációs rendszer függöek: a felhasználó felöl többféle kommunikáció lehetséges attól függően, hogy milyen rendszert szeretnének megvalósítani (beszéd, grafikai, egér, billentyuzet, érintésvezérlo)
- az interaktív működés érdekében: visszahallgatás
- prompt lejátszás: előre rögzített felvételek lejátszása, pl. „Üzenethallgatáshoz nyomja meg a 1. gombot!”
- az output lehet pl. bankjegykiadó, számlanyomtató

5.1.2. Dialógus vezérlo



- vezérlo
 - ☞ irányítja a feldolgozást (processzor, több gép, esetleg AI)
 - ☞ a beérkezo információ alapján a DB lekérdezozt hívja meg
- válaszgeneráló
 - ☞ biztosítja a kapcsolatot: a kimeneti egységek közül melyikkel és hogyan, mit kell kommunikálni
- felhasználói interfész felöl érkezik valamilyen metanyelven az input (DTMF karakterek, beszédfelismerotol valamilyen adat)
 - ☞ a beszédmegérto: nem érvényes input kiszurése

- ☞ nyelvi feldolgozó: további egyszerűsítések
- egyéni beállítások arra szolgálnak, hogy az egész dialógus egyéni lehessen a felhasználónak (rendszerfüggő illetve rendszerfüggetlen)
- dialógus leíró:
 - ☞ az inputok közül melyik micsoda, mikor mit várunk, mire mit kell csinálni (a protokoll)

5.1.2.1. *A rendszer modalitása*

- unimodális: 1 érzékre hat
- multimodális: több érzékre hat
- és ezeket input és output oldalról is vizsgáljuk
- PC
 - ☞ egér
 - ☞ billentyűzet
 - ☞ GUI
 - ☞ beszéd
- Információs pult
 - ☞ GUI
 - ☞ érintő képernyő
 - ☞ mutató detekció
 - ☞ gesztus

5.1.2.2. *Telefonos alkalmazás*

- vezeték (klasszikus: kizárólag hang továbbítása)
- mobil (rég, analóg)
- mobil (GSM, nemcsak hangátvitel, SMS, WAP)
- IP telefon
- Videotelefon

5.2. Dialógus rendszerek osztályozása

- A vezérlés jellege szerint (ki kezdeményez, ki irányítja a rendszer működését)
 - ☞ rendszer vezérelt
 - ☞ felhasználó-vezérelt
 - ☞ vegyes kezdeményezésű
- a vezérlés módja szerint (elsősorban telefonos alkalmazásoknál)
 - ☞ DTMF
 - ☞ beszédfelismerő
 - ☞ egyéb

5.2.1. Rendszer vs felhasználó vezérelt

- a rendszer határozza meg a navigációt
- menürendszer, felajánlott választási lehetőségekkel (pl. hangposta: mindig végig kell menni az összes menüponton)
- a felhasználó határozza meg a navigációt (pl. úgy döntök, hogy az utolsó két üzenetet törölöm)
- vegyes kezdeményezésű (lehetőség a navigáció módosítására, az előző két kombinációja)

5.2.1.1. *Menürendszer tervezési szempontok*

- építőkockák felhasználása (ha van egy működő menü, abból már lehet építkezni)
- 4 választási lehetőség (ennél több ne legyen, mert a felhasználó úgysem tudja megjegyezni), visszalépés biztosítása
- maximális mélység 4-5 szint
- felhasználófüggetlen menürendszer (egyértelmű funkciók, következetes rendszer)
- csak a témakörhöz tartozó információk közlése
- megfelelő részletesség

- újdonságok, fontos elemek kiemelése
- gyakran használt gombok: kényelemhez legyenek igazítva
- egyes gombokat úgy helyezünk el, hogy többletjelentést is lehessen nekik tulajdonítani (pl. nyilak)

5.2.1.2. Rendszerfüggetlen egyéni opciók

- felhasználói szint
 - ☞ kérdések, magyarázatok hossza, részletessége
 - ☞ választási lehetőségek száma
 - ☞ felajánlott választási lehetőségek száma
- felolvasás paraméterei
 - ☞ beszélő kiválasztása (férfi, nő)
 - ☞ beszédtempó beállítása
 - ☞ szünetek időtartama (pl. mondatok között)
- adaptív változtatás / felhasználó vezérelt

5.2.1.3. DTMF vezérlés

- Dual Tone Multi-Frequency (4*4 frekvencia)
- adatbevitel a telefon billentyűzetén
- ☺ nagyon megbízható
- ☺ kipróbált technológia
- ☺ olcsó
- ☹ a kialakítható menü nem felhasználóbarát
- ☹ nehézkes a használat, ha a billentyűzet nem elérhető
- ☹ humán operátor szükséges lehet

5.2.1.4. Vezérlés hanggal

- ☺ a telefonba beszélni természetes dolog
- ☺ szélesebb témakörben alkalmazható (nemcsak számok)
- ☹ megbízhatatlanabb
- ☹ bizonyos esetekben jóval lassabb, mint a DTMF (főleg adatbevitel illetve gyakorlok felhasználók esetén)
- ☹ kötött szókészlet
- Yes/No elvű rendszerek
 - ☞ lassú és természetellenes
 - ☞ legmegbízhatóbb a beszédfelismerős rendszerekben
 - ☞ jól kell megválasztani a yes/no magyar megfelelőjét
 - ☞ 2 szavas szótár nem elég (igen, jó, OK, mehet, rendben, ja, aha ...)
 - ☞ adatbevitel fa struktúrában
- Kötött szótáras
 - ☞ meghatározott (10-20) szó felismerése
 - ☞ kényelmesebb menürendszer jellegű
 - ☞ beszélofüggetlen / adaptív jellegű (beszélofüggetlen mag és a felhasználóhoz igazítják)
 - ☞ gyors elérés mély struktúrák esetén is
 - ☞ keverhető a DTMF vezérléssel „fall back” (visszalépünk DTMF rendszerre)
 - ☞ bizonyos esetekben az adatbevitel nehézkes (számlaszám, visszaellenorzés, javítás)
- adott témakörben bármilyen információ mondható. lekérdezheto(a magyar sokkal szabadabb, jól meg kell gondolni, hogy mit próbálunk felismerni)
- felhasználó-vezérelt
- emberközeli használat
- pl. menetjegy-árusító rendszer
- diktáló rendszerek (tematika nélkül)
- mesterséges intelligencia szükséges (komoly nyelvtani szabályok, komplex dolgok)
- beszédfelismerés + mozgás és kézmozdulatok felismerése

5.2.1.5. Megerosítés (verifikáció)

- szükséges (mivel nem 100% -os a megbízhatóság, hibás bevitel)
- elvart (a felhasználó biztonságérzete miatt): lehet explicit pl. közvetlen visszakerdezéssel és implicit pl. elrejtve a következő kérdésben
- explicit megerosítés minden egyes adatra: rákerdezés minden egyes adatra
 - ☞ eldöntendo kérdés
 - ☞ egyszeru struktúra
 - ☞ kényelmetlen a dialógus a felhasználó számára
- explicit megerosítés javítással: rákerdezés minden egyes adatra
 - ☞ ugyanaz, mint az elobb
 - ☞ igen/nem válasz mellett a javított adat is megadható
 - ☞ gyorsabb dialógusmenet
 - ☞ kevésbé akadozó

- explicit megerosítés több adatra: rákérdezés minden adatra egyszerre
 - ☞ kevesebb kérdés
 - ☞ igen/nem válasz mellett a javított adat is megadható / nem adható természetesebb
 - ☞ csak az adatbevitel végén van ellenorzés
 - ☞ gondot jelenthet, hogy jól értelmezi-e a rendszer a mondanivalót
 - ☞ megnóhet a kérdés összetételének bonyolultsága
- implicit megerosítés
 - ☞ követo adatbekérésbe ágyazott ellenorzés – a kérdésben benne van a régi információra vonatkozó visszakérés + az új információ
 - ☞ közelebb áll a természetes párbeszédhez
 - ☞ a kérdés hossza megnó
 - ☞ javítás nehzkesebb
 - ☞ a rendszer bonyolultabb
- hibás megerosítés
 - ☞ felismerési probléma
 - ☞ rosszul ismeri fel
 - ☞ a felhasználó nem igen/nem-mel válaszolt
 - ☞ nem ismeri el a rendszer
 - ☞ többszörös megerosítés kritikus adatoknál: nem rögtön utána, inkább a végére érdemes berakni még egy megerosítést

5.2.2. Dialógus leíró eszközök

- SAPI
- VoCAPi
 - ☞ kisméretu eszközök: telefon, mosógép, fénymásoló
- ECTF – megpróbálja összefogni a különböző technológiákat
- VoiceXML