

# MPEG-4 modell alkalmazása szájmozgás megjelenítésére

Takács György, Tihanyi Attila, Bárdi Tamás, Feldhoffer Gergely, Srancsik Bálint

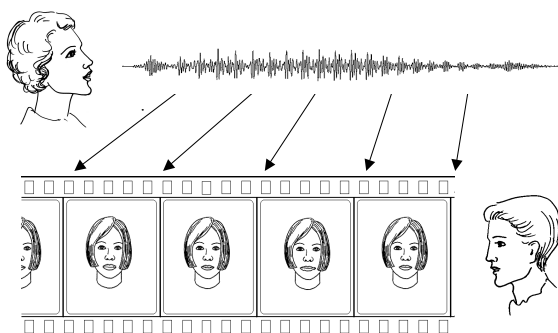
Pázmány Péter Katolikus Egyetem Információs Technológia Kar  
1083 Budapest Práter utca 50/a  
Telefon: +(36) 18864763, Fax: +(36) 18864724,  
email {takacs.gyorgy, tihanya, bardit, flugi, sraba}@itk.ppke.hu

## 1. ÖSSZEFOGLALÓ

Az MPEG-4 kódolást multimédia alkalmazások, mozgó fejek élethű megjelenítése figyelembe vételével fejlesztették. Az általános célú, nyílt forráskódú LUCIA modellt alkalmassá tettük mozgó száj képének megjelenítésére. A fejlesztésben kiemelten kezeltük a siketek speciális igényeit. Törekedtünk számítási erőforrások minimális igénybe vételére, hogy az alkalmazás mobil eszközökön is megvalósítható legyen. Fontos eredménynek tartjuk, hogy az MPEG-4 animáció működik akkor is, ha nem képpontok mintavételezése alapján származtattuk a tartópont paramétereket, hanem beszédjelből számoltuk azokat.

## 2. ELŐZMÉNYEK

Egy teljes rendszert dolgoztunk ki, amely alkalmas arra, hogy beszédjelből mozgó száj képét állítsa elő. A mozgó szájról siketek képesek a beszédet leolvasni. A rendszer ismertetése ugyanezen folyóirat számban megtalálható [3]. Itt azokat a részleteket és általános megfontolásokat taglaljuk, amelyek kifejezetten a megjelenítő egységre vonatkoznak.



1. ábra Mozgó száj előállítási vázlata

Folyamatos beszédjelből mozgóképfolyamat hozunk létre. Ez egy olyan transzformáció amelynek lényegi részét egy neurális hálózat hajtja végre. A neurális hálózat komplexitását korlátok között kellett tartani, ezért elengedhetetlen volt az emberi beszéd folyamat lényegét jól megragadó, tömör és hatékony leírása a vizuális beszédnek.

A neurális hálót előfeldolgozott hangadatokkal tanítottunk és képi koordinátákon vártunk a kimeneteken. Főkomponens analízist alkalmaztunk a képi koordináták tömör reprezentálására. Így mindössze 6 kimeneti jellemző kisebb, mint 2% hibával leírta a szükséges képi koordinátákat.

A rendszer fejlesztésében külön kezelt probléma volt a mozgóképfolyamat megjelenítés modellje.

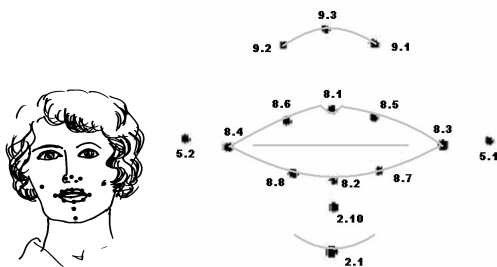
## 3. AZ MPEG-4 SZABVÁNY FEJMOZGÁSOK TÖMÖRÍTETT KÓDOLÁSÁRA

Az MPEG (Moving Picture Expert Group) szabványok fő célja a hang és videó jelek tömörítése. A tömörítés alapvető követelményei a hatékonyság és élethűség. A multimédia-alkalmazásokban elterjedt az MPEG-2 kódolás. Az ezt meghaladó MPEG-4 kódolás is ígéretes jövő előtt áll, ugyanakkor céljainkat közvetlenül támogatja. Az MPEG-4 nem csak nagy tömörítésre alakították ki, hanem figyelembe vettek olyan multimédia alkalmazásokat is, mint 3D-s jelenetek, animációk, szintetizált hangok, képek, szövegek, grafikák külön vagy akár együttes kezelése és élethű megjelenítése.

Az MPEG-4 szabvány egyik legösszetettebb része a fej és az emberi test megjelenítése és mozgása (FBA - Face and Body Animation). Az FBA-ra vonatkozó szabványrész leírja az arc és a test alakjának és mozgásának kódolási alapelveit. Az FBA egyik legfontosabb tulajdonsága tehát, hogy nem adja meg pontosan a kódolási és a dekódolási eljárást, csak a küldött adat formáját és értelmezését.

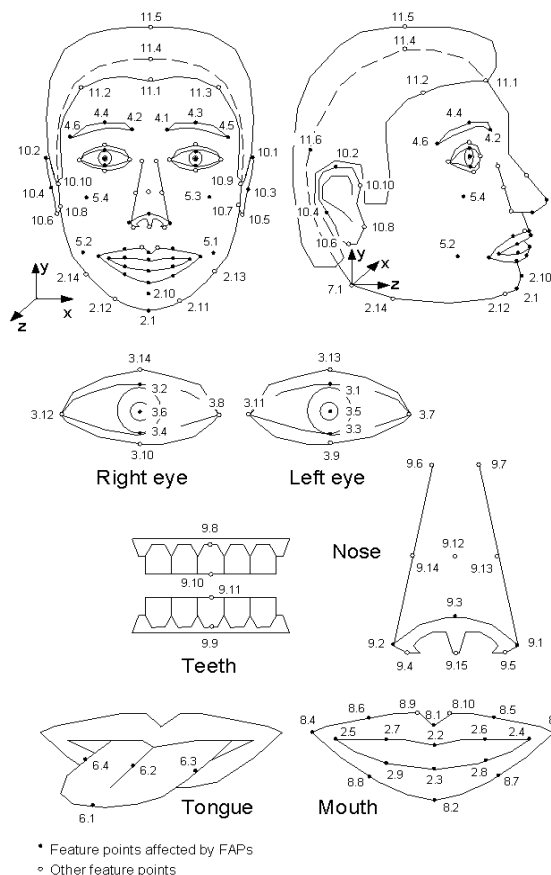
Az MPEG-4 szabvány az arc modelljét az arc normál állapotával írja le, megad a több tartópontot (Feature Point - FP) és az arc mozgását leíró paramétereket (Facial Animation Parameter - FAP), mely lényegében a normál archoz képesti elmozdulást jellemzi. Az elmozdulások méretét és arányát a szabvány szerint mindig az emberi arcra jellemző alapvető méretek alapján fejezi ki. A szakirodalomban ennek elterjedt rövidítése FAPU (Face Animation Parameter Unit). A 6. ábra mutat magyarázatot erre. A FAPU-kat az arc olyan jellegzetes távolságaiból kell számolni, mint például a szemgolyók távolsága vagy a száj szélessége.

A szabványban 84 tartóponttal írják le az arcot. (Az adatbázisunk összeállítása során mi 15 FP-t használtunk a száj és környékének leírására).



2. ábra Felhasznált tartópontok

A tartópontok fő feladata, hogy referenciaként szolgáljanak a FAP-ok számára. A FAP-ok által leírt összetett mozgások mindig a normál tartópontok által leírt fejre vonatkoznak. A normál fej csükköt száját és semleges arckifejezést jelent. Vannak olyan FP-k is, melyekre egy FAP sincs közvetlen hatással (pl.: az orr szélei). Ezeket mindössze az arc alakjának meghatározására használják. Az FP-ket minden MPEG-4 kompatibilis modellen a 3. ábra alapján kell elhelyezni.



3. ábra A tartópontok szabványos elhelyezkedése a fejen

FAP-ból a szabvány 68-at különböztet meg, melyet 10 csoportba sorol az alapján, hogy az arc mely részét mozgatja.

Az első két FAP magas szintű paraméter. Ez azt jelenti, hogy ezekkel előre beállított komplexebb

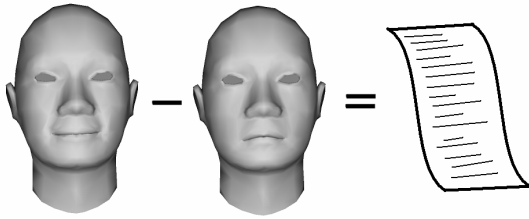
mozgást lehet kivitelezni. Az első FAP egy megadott vízéma szerinti megjelenést határoz meg. A vízéma a fonéma képi megfelelője. A második FAP a hat alap érzelm megjelenítésére szolgál, úgy mint öröm, bánat, harag, félelem, undor és meglepetés. Tovább érzelmkifejezéseket az alap érzelmek keveréséből lehet megjeleníteni.

A többi FAP alacsony szintű. Ezek abban különböznek a magas szintű FAP-októl, hogy itt a mozgás irányát és amplitúdóját kell megadni, nem pedig egy összetett feladatra előre összeszerkesztett mozgásvezérlést kell kezdeményezni. Az alacsony szintű FAP-ok általában egy-két tartópontot mozgatnak. Előfordul olyan FAP is, amely az összes FP-t mozgatja, ilyen például a fej forgatása. Az alacsony szintű FAP-oknál a szabvány meghatározza, hogy a mi a hozzá illő FAPU, amiből a mozgás mérték alapja. A FAP előjele a tartópont mozgásirányára vonatkozó információt hordoz, például a száj nyitására vonatkozó paraméterek pozitív, a zárásra vonatkozók negatív előjelűek. Ez független attól is, hogy a tartópont a száj alsó vagy felső részéhez tartozik. A mozgatás lehet eltolás, forgatás vagy skálázás.

#### 4. LUCIA MODELL

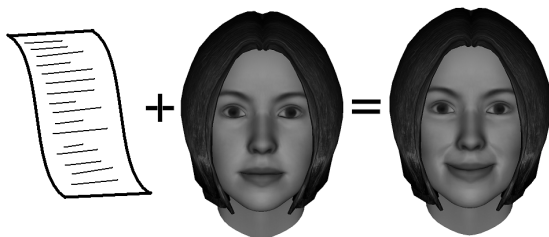
A legtöbb modell, legyen az két- vagy háromdimenziós, hálóból áll. A háló (mesh) több egymáshoz illeszkedő nem feltétlenül egy síkban levő sokszöget tartalmazó felület. A hálóban a csúcspontok koordinátáin kívül a lapok, az élek és a csúcsok illeszkedési viszonyait is nyilván kell tartani [12]. A modell felületi jellemzői, textúrája erre a rácsra van ráhúzva. Ahogy mozgatjuk a háló csúcspontjait, úgy mozog vele a textúra is. Ám az MPEG-4 szabványban csak az FP-k mozgatására van mód, az egyes hálókérra közvetlenül nincs. Egy modell tetszőleges számú és finomságú hálóból állhat, a szabvány erre nem terjed ki. Minden MPEG-4 kompatibilis fejmodell azonban azonos tartópont rendszerre épül. A háló mozgatása a tartópontok mozgatásával történik.

A LUCIA modellt Cosi vezetésével olasz kutatók fejlesztették ki [1]. Ez egy nyílt forráskódú mozgó fejmodell. A LUCIA egy MPEG-4 megvalósítás, ami alkalmas vízémák és érzelmi állapotok FAP paraméter alapú közvetlen megjelenítésére. Az MPEG-4 modell tömörítést kifejtő (decompress) része egy grafikus modell mozgatási feladat, alapvetően az 5. ábra szemléltetése szerinti információk felhasználásával képes átvinni a mozgás jellegzetességeit. A szabványosított eljárás során az alaphelyzetű fej teljes képének meghatározása és vevő oldalra történő átvitele valósul meg, és a továbbiakban csak az alaphelyzettől történő eltérések átvitelére van szükség a tömörített adatközlés során.



**4. ábra Az MPEG-4 rendszerű tömörítés koncepciója**

Az MPEG-4 tömörítési folyamat (4. ábra) azon az elven működik, hogy a tömörítendő mosolygós fej lényeges paramétereinek valamint az alaphelyzetű fej paramétereinek különbségéből meghatározza a tömörített jellemzőket. Az MPEG-4 koncepció szerint ez a jellemzősor fej alakjától és környezetétől független adatokat tartalmaz.



**5. ábra Az MPEG-4 rendszerű visszaállítás koncepciója**

A visszaállítási folyamat (5. ábra) során a tömörített jellemzőkhöz, amely jelenleg a mosolygós adatait tartalmazza „hozzáadva” egy tetszőleges alaphelyzetű fej paramétereit egy mosolygós fej képét kapjuk. Az alaphelyzetű fej meghatározó adatai között kell elhelyezni a felületi jellemzőket valamint az esetleges további adatokat mint például a modell haja szeme stb. A visszaállítás során kell létrehozni a felületeket azok megvilágítástól függő színezésével együtt [2].



**6. ábra Az emberi arcra jellemző méretek**

- ESO=szemgolyók távolsága;**
- IRISD0=Az írisz étmérő;**
- ENS0=Az orr hossza;**
- MNS0=Az orr és azáj távolsága;**
- MW0=A száj szélessége**

Az MPEG-4-ban a tömörítés során meghatározott és felhasznált távolság mértékrendszer (6. ábra) lehetőséget biztosít arra, hogy a tömörített információ felhasználásával tetszőleges más alaphelyzetű fejre lehessen alkalmazni a visszaállítást,

és így lehessen változtatni a visszaállítás folyamatát.

Az ESO; IRISD0; ENS0; MNS0; MW0; távolságok határozzák meg az adott arcberendezésen alkalmazandó távolságegységek halmazát. A távolságmérésnek ez a módszere biztosítja azt a lehetőséget, hogy a visszaállítás során az eredetitől jelentősen eltérő felépítésű alaphelyzetben álló fejre is visszaállíthatók legyenek a tömörített információk.

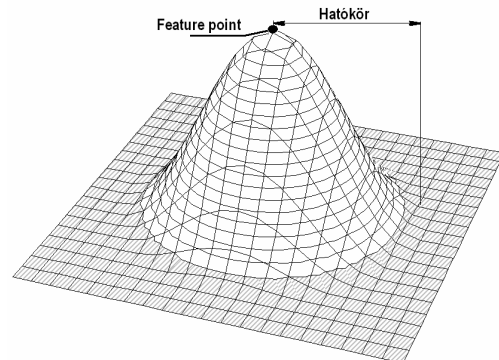
Az MPEG-4 szabványnak köszönhetően az arc mozgásához nem kell képkockáról képkockára megadni a videó minden egyes pixelét, mindössze a mozgató FP-khez tartozó FAP-okat kell továbbítani. Ennek köszönhetően igen alacsony sávszélességen keresztül is elérhető a real-time arcanimáció.

Az MPEG-4 szabvány előnyeit leginkább Internetes alkalmazásokban használják. Találkozhatunk olyan rendszerrel, mely az E-mail-eket alakítja át olyan videóvá, ahol az általunk kiválasztott személy mondja el az üzenetet. Léteznek olyan alkalmazások, melyek Internetes áruházakban „eladókat” alkalmaznak, vagyis egy MPEG-4 szabványú modell ad segítséget az árakról, a minőségről vagy éppen a készletről.

A szintetikusan létrehozott szájmozgás megjelenítésére felhasznált LUCIA modell egy szokásos 3D grafikus modell, amely animálható és így a céljaink megvalósítására alkalmas.

Az adatbázis felvételnél az arcra felfestett pontok kis mértékben eltérhetnek a szabványban előírt tartópontok helyétől. Ezt a hibát úgy korrigáltuk, hogy a tartópontokat ráillesztettük a felfestett pontokra úgy, hogy szintetizáláskor egybeessenek.

Az 7. ábra bemutatja egy eredetileg vízszintesen elhelyezkedő négyzögháló felhasználásával készített animáló eljárás hatását abban az egyszerű esetben, ha az eredeti helyzetből függőleges irányban felfelé kívánjuk elmozdítani a síknak egyetlen pontját.



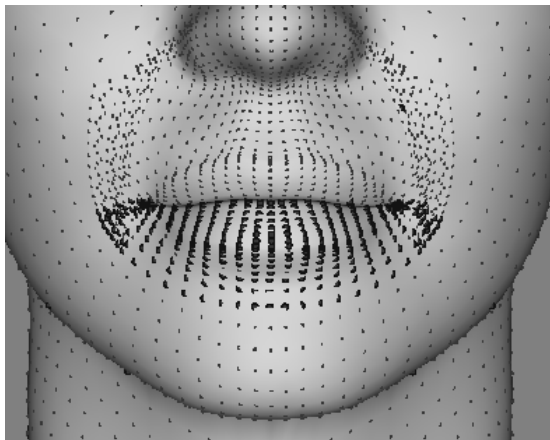
**7. ábra Az FP függőleges elmozdításának hatása a vízszintes felületre**

A hatókörükben összeérő, egymásmellé eső pontok egymásra hatását megfelelő súlyozással kell kiküszöbölni. Elképzelhető, hogy egy hálórész több tartópont is mozgatni akar. Ilyenkor természetesen súlyozottan összegződnek az elmozdulások. A súlyozás meghatározásánál az elmozdítást eredmé-

nyező pont hatását annak távolságával fordított arányban határozzuk meg, ez a módszer azt eredményezi, hogy a modell rácspontjainak elmozdulását a FP-hez közeli rácspontok esetén nagymértékben az FP helyzete határozza meg. A vázolt eljárással lehetséges kijelölt pontok és hozzájuk tartozó területek rögzítése. Ilyen technikával oldottuk meg a 3D-s LUCIA fej állának mozgását.

Annak érdekében, hogy az állcsont a megfelelő forgáspont körül elforduljon az állcsúcsot (2.1-es FP) mozgattuk. Az állcsont miatt nagy hatókörrel kell a 2.1-es FP-t mozgatni aminek az a hatása, hogy szemből nézve úgy tűnik mintha az egész áll leesne. A jelenséget meg lehet szüntetni oly módon, hogy az arc körvonalához tartozó 2.13 és 2.14-es FP-t minden irányban 0-val mozdítjuk el, ennek hatására a 2.13 és 2.14-ös FP-k nagy súllyal helyben tartják az arc körvonalát és csak elenyésző mértékben mozdul a környezetük a 2.1 és 2.10 pontok mozgásának hatására. Az alkalmazott technika teljesen kiküszöböli az áll leesésének a jelenségét.

A LUCIA modell tartalmazza az alsó és felső fogsort valamint a nyelvet is. Az alsó fogsor mozgását kizárólagosan az állcsúcs mozgása határozza meg, a felső fogsor mozgását az orr megfelelő pontjaihoz kötöttük, így annak elmozdulása minimális, hiszen az orr középpontját tekintettük a munka során referenciának. A nyelv mozgásával a projekt nem foglalkozott.



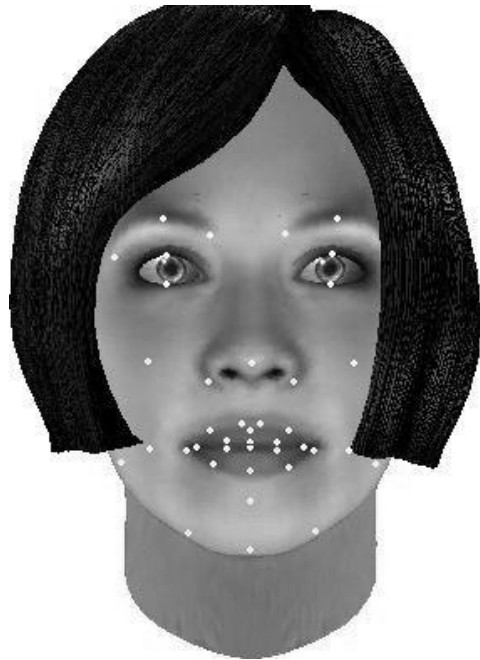
**8. ábra** Az alsó szájszélet meghatározó hálópontok

A mozgatható felületen a háló törése, szakadása (például szem, száj) azt a problémát jelenti, hogy a szakadási vonalnál tovább azon átnyúlva nem alkalmazhatjuk az előzőekben vázolt módszert. Például az alsó ajak mozgása nem húzza magával a felső ajkak hálórészét, pedig azok a hatókörön belül esnek. Ezzel a módszerrel kezelhető a száj, a szemek természetes nyitása. Azt a módszert választottuk, hogy minden mozgatható FP-hez meghatároztuk a modellünk egy-egy háló csúcspontokkal leírt egybefüggő részét. Ez jelentősen gyorsítja a mozgató algoritmusokat, mivel nem kell a teljes fej összes rácspontjának távolságát meghatározni min-

den egyes FP helyzetétől, hanem elegendő a kijelölt részhalmoz pontjainak a figyelembe vétele a számítások folyamán. Az alsó és felső ajakrész szétválasztását szemlélteti a 8. ábra.

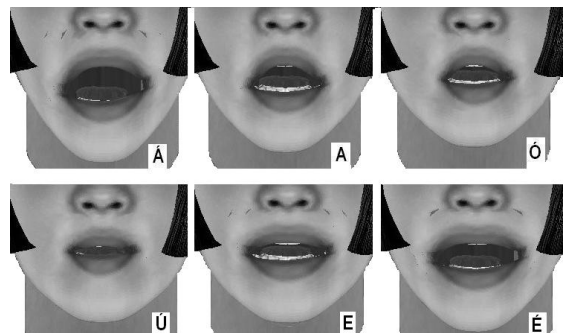
Az ábrán sötétebb pontok jelölik a száj alsó szélét. Ezekre a pontokra hatnak ezeket mozgatják a 8.2; 8.7; 8.8 tartópontok (3. ábra).

Minden FP-hoz tartozik egy mozgatható hatókör - egy gömb alakú térrész- és azon hatókörön belül levő rácspontok elmozdulását határozza meg az adott FP elmozdulása az MPEG-4 rendszerben meghatározott skálázás szerint.



**9. ábra** LUCIA modellen alkalmazott FP-k

A 3D grafikus modellt az MPEG-4 rendszernek megfelelően ki kell egészíteni a 3 dimenzióban értelmezett FP-vel, és azok hatókörének meghatározásával valamint az egyes FP-k által mozgatható rácspontok halmazával (9. ábra).



**10. ábra** Példák magyar nyelvű jellegzetes magánhangzó szájaállásokra (vizémákra)

A beszédjelből szájmozgást előállító project során az előzőekben részletezett módon kialakított LUCIA modellt alkalmaztuk. A project eredeti elképzelései szerint a megvalósításkor a beszédjelből közvetlenül az FP mozgatható paramétereket

állítottunk elő, tehát nem volt szükség arra, hogy az egyes vizémákat külön-külön meghatározzuk és előállítsuk, de a hosszan kitartott magánhangzók tiszta fázisainál jól megkülönböztethető szájállásokat hozott létre a fejmodell (10. ábra).

## 5. MÉRÉSI EREDMÉNYEK ÉS KÖVETKEZTETÉSEK

Az animációs rendszerünk komponenseinek ellenőrzésére szájrol olvasási kísérleteket végeztünk siket tesztalanyokkal. A szájrol olvasási feladatok nehézségét úgy állítottuk be, hogy körülbelül 95 és 100% közötti felismerési arányt kapjunk a vetített eredeti videó felvételekre, hogy referenciaként szolgálhasson. Ilyen jó arányt az előzetes kísérletek leírásánál [3] már ismertetett módon a felismerendő szövegben használt szókinccs és nyelvtan erős szűkítésével, valamint egy jól artikuláló jeltolmács szerepeltetésével értük el. Ezután mértük a felismerési arányt, úgy, hogy a videó felvétel helyett az animált beszélőfej-modell volt látható, ugyanakkor minden más kísérleti körülményt változatlanul hagytunk.

A fejmodellre való áttérés két lépcsőben történt. Az elsőben a felvételeken festékpöttyel megjelölt MPEG-4 pontok koordinátáit igyekeztünk átvinni a modellre: vagyis a fejmodell vázát képező háló megfelelő csomópontjait minden képkockán a felvételen mért koordinátájú pozíciókba mozgattuk. Ezzel azt kívántuk elérni, hogy a modell közvetlenül utánozza a jeltolmács artikulációját, ebben a lépésben a hang még nem játszott szerepet. A második lépcsőben a beszédhang alapján számított koordináták szerint vezéreltük a fejmodellt. Itt már csak a hangbemenetre volt szükség az animáció előállításához [3].

A kísérlet során a felismerési arányok a következők szerint alakultak:

eredeti felvételek (referencia) – 97,1%;

animáció a jeltolmácsra festett tartópontok alapján vezérelt modell (1. lépcső) – 54,9%;

animáció a hang alapján (2. lépcső) – 47,9%.

Jelen cikk szempontjából a felvételekről a LUCIA modellre való áttérés, vagyis az első lépcső érdekes. Itt elég jelentős romlás tapasztalható a felismerési arányban, ennek lehetséges (valószínű) okaira térünk ki röviden.

Megállapítható, hogy az általunk kiválasztott és a felvételeken megjelölt MPEG-4 pontok helyzete hiányosan (információ veszteséggel) reprezentálja azokat a látható beszédképzési jellemzőket, melyek a szájrol olvasásban szerepet játszanak. A kísérletek után minden alkalommal kikértük a résztvevő siketek véleményét, hogy mely tényezők gátolták őket leginkább a szájrol olvasásban. A felvételek és az animációk között talán a legfontosabb különbség, hogy a fejmodellnek nincs nyelve. De ha a LUCIA modell lehetővé tenné a nyelv animálását, akkor is problémát jelentene, hogy nincsenek referencia adataink a nyelv pillanatnyi helyzetéről, nem

tudjuk, hogyan is kéne mozogni. A nyelvre a felvételeken nem festhettünk pontot. A nyelv hiányában pl. *kilenc* vagy a *nulla* szavak felismerése gondot okozott az animáció esetében, míg a felvételeken jól látható volt a nyelv főtről-le csapódása az *l* hang után, így valamennyi tesztalanyunk könnyedén felismerte azokat.

A másik problémánk volt, hogy a felvételekhez csak az ajkak külső kontúrján tudunk pontokat megjelölni, beljebb nem. Ezek viszont az ajkerekítésről kevés információt tartalmaznak. Az animációkon elsősorban ajkerekítéses magánhangzók (pl. *u*, *ü*) voltak kifogásolhatók. Szintén a pontok elhelyezésére vezethető vissza, hogy nincs elegendő információnk a fogak láthatóságáról. Pedig elsősorban ettől függ az ajkakon belüli terület világossága, ami egy igen karakteres és könnyen észlelhető vizuális jellemző [4].

Az MPEG-4 szabvány eredeti célja egy olyan modell megalkotása, aminek segítségével tömöríteni, majd rekonstruálni lehet mozgó fej adatokat. Munkánk során megoldottuk, hogy a szabványra építve olyan minőségben mozgatható a száj és környezete, hogy ennek alapján siketek a beszédet képesek szájrol leolvasni.

További fontos eredménynek tartjuk, hogy az animáció működik akkor is, ha nem képpontok mintavételezése alapján származtattuk a tartópont paramétereket, hanem beszédjelből számoltuk. Az eredményeink azt mutatják, hogy igen kis különbség van a mintavételezéssel vezérelt arc, és a beszédjel alapján vezérelt arcmodell felismerhetősége között.

További fejlesztést igényel a fejmodell finomítása. A száj külső körvonalán túl a belső kontúr, fogak vagy nyelv láthatósága tűnik a továbblépés első lehetőségeinek.

## 6. KÖSZÖNETNYILVÁNÍTÁS

A szerzők ezúton is kifejezik köszönetüket a Nemzeti Kutatási és Technológiai Hivatalnak a 472/04 szerződés keretében nyújtott támogatásáért.

## 7. IRODALOM

- [1] **Cosi P., Fusaro A., Tisato G.**, "LUCIA a New Italian Talking-Head Based on a Modified Cohen-Massaro's Labial Coarticulation Model", **Proceedings of Eurospeech 2003, Geneva, Switzerland, September 1, 2003, Vol. III, pp. 2269-2272**
- [2] **Szirmai-Kolos László, Antal György, Csonka Ferenc**, "Háromdimenziós grafika animáció és játékfejlesztés", **ComputerBook Kiadó Kft, Budapest 2003**
- [3] **Takács György, Tihanyi Attila, Bárdi Tamás, Feldhoffer Gergely, Srancsik Bálint**: „Beszédjel átalakítása mozgó száj képévé

*siketek kommunikációjának segítésére”*  
Híradástechnika 2006 xxxx

- [4] **László Czap, János Mátyás** „*Virtual Speaker*” Híradástechnika, selected papers, 2005 Június. Vol LX. pp 2-5.
- [4] **I. Pandzic and R. Forchheimer**, „*MPEG-4 Facial Animation: The Standard, Implementation and Applications*”, Wiley, 2002.