

## Érzelmes beszéd gépi előállítására érzelmelem specifikus beszédadatbázisok felhasználásával

Fék Márk, Zainkó Csaba, Németh Géza

Beszédtechnológiai Laboratórium,  
Távközlési és Médiainformaticai Tanszék,  
Budapesti Műszaki és Gazdaságtudományi Egyetem,  
1117 Budapest, Magyar tudósok körútja 2.  
{fek, zainko, nemeth}@tmit.bme.hu

**Kivonat:** Tanulmányunkban megvizsgáljuk hogyan lehet érzelmelem specifikus beszédadatbázisok felhasználásával gépileg érzelmes beszédet előállítani. Kísérletünket magyar nyelvre végeztük, de a módszer nyelvfüggetlen. Felvettünk egy szemantikailag semleges tartalmú mondatot és 26 logatomot amelyek a mondat szintetizálásához szükséges diádokat és CVC triádokat tartalmazták. A hanganyagot egy profi színész mondta fel a hat alapérzelmelemnek megfelelően, illetve semleges érzelmi változatban. A logatomok felhasználásával 7 érzelmelemfüggő beszédelem adatbázist hoztunk létre. A 7 beszédelem adatbázist összepárosítva a természetes mondatokból kinyert 7 prozódiai kontúrral 49 szintetizált mondatot állítottunk elő. A logatomokban, illetve a természetes és a szintetizált mondatokban hallható érzelmelemeket 194 tesztalany értékelt ki. A tesztelők a logatomok 99%-ban, illetve az összes természetes mondatban szignifikánsan a véletlen találgatás szintje felett ismerték fel a színész által kifejezett érzelmelemeket. Az érzelmelem azonosítási aránya egyes szintetizált mondatok esetén meghaladta a természetes mondatokét.

### 1 Bevezetés

A napjainkban elterjedten használt elemösszefűzéses beszéd szintetizátorok általában adnak valamilyen lehetőséget a prozódiai paraméterek (alapfrekvencia, hangidőtartamok, intenzitás) vezérlésére. A gépi beszéd hangszínezését azonban alapvetően a beszédadatbázis határozza meg, amelyből az összefűzendő elemeket kiválasztják. Ismert, hogy a beszédben az érzelmelemeket a prozódia és a hangszínezés együttesen hordozza [1].

Korábbi kísérletek több nyelv esetében vizsgálták [2,3] az érzelmelem specifikus beszédadatbázisok használatának lehetőségét elemösszefűzéses beszéd szintetizálás esetében. A kísérletek mindegyikében a bemondók felolvastak egy szöveget az előzetesen definiált érzelmelemnek megfelelő változatokban. A felvett beszédből minden érzelmelemhez egy-egy érzelmelem specifikus adatbázist hoztak létre. Elemkiválasztáson alapuló beszéd szintetizálás segítségével szemantikailag semleges mondatokat állítottak elő az egyes érzelmi adatbázisokból.

A [2] és [3]-ban leírt kísérletekben a szintetizált mondatok prozódiai paramétereit természetes mondatokról másolták, illetve [2] második kísérletében a beszéd-szintézis algoritmus állította elő a prozódiát. Ezután meghallgatásos tesztekkel végeztek, hogy meghatározzák a szintetikus mondatokkal kifejezni szándékolt érzelmek felismerési arányát.

Montero és szerzőtársai [2] szemantikailag semleges szöveget vettek fel spanyolul egy professzionális színész segítségével, aki négy érzellemmel és semleges stílusban mondta be a szöveget. A szerzők a természetes mondatokról átmásolt prozódia segítségével szintetizált mondatokon a véletlen szintnél magasabb azonosítási arányt értek el (semleges 76%, öröm 62%, meglepetés 91%, szomorúság 81%, harag 95%) egy 6 választási lehetőséget megengedő meghallgatásos tesztben (amely a "nem azonosítható érzelmek" opciót is tartalmazta). Hasonló eredményeket értek el (kivéve a meglepetésre 53%) az érzelmi adatbázisokon betanított, automatikusan generált prozódia használata esetében is.

Bulut és szerzőtársai [3] örömet, szomorúságot, haragot és semlegességet közvetítő rövid szövegrészeket vettek fel egy fél-profí színész segítségével. Az adatbázisokból előállított mondatokat a véletlen találgatás szintje felett azonosították (harag 86%, szomorúság 89%, öröm 44%, semleges 82%) egy 4 választási lehetőséget tartalmazó meghallgatásos tesztben.

Schröder és Grice [4] egy teljes német diádos elemkészletet vettek fel három különböző vokális erőfeszítés mellett (lágy, modális és hangos). Meglátásuk szerint a különböző hangszalag feszességeknek megfelelő vokális erőfeszítések fontos szerepet játszanak az érzelmek kifejezésében. A diádotok értelem nélküli hordozó szavakba (logatomokba) ágyazták és egy német anyanyelvű bemondó olvasta fel őket, akit megkértek, hogy állandóan tartott hangmagasságon beszéljen. Egy (elemkiválasztás nélküli) elemösszefűzés szintetizátort használtak a tesztmondatok előállításához. A meghallgatásos teszt folyamán kimutatták, hogy a szintetizált mondatokban a szándékoltnak megfelelően észlelhető a vokális erőfeszítés.

A fent felsorolt munkák ellenére még mindig kérdéses, hogyan lehet olyan érzelmes beszéd-szintetizátort készíteni, amely a megfelelő hangszínezettel állítja elő az egyes érzelmeket. Tanulmányunkban az érzelmek specifikus diád és triád elemkészlet használatát vizsgáljuk érzelmes beszéd előállítására.

Kísérletünk hasonló a korábbiakhoz a következő lényegi eltérésekkel. Az általunk használt szintézis módszer elemkiválasztás nélküli elemösszefűzésen alapul és [4]-hez hasonlóan szintén értelem nélküli logatomokat vettünk fel, de a bemondott logatomok közvetlenül hordozták az egyes érzelmeket, hasonlóan a [2] és [3]-ban felolvasott szövegekhez. A mi beszédelem-készletünk mind diádotokat, mind CVC triádotokat tartalmazott növelve az összefűzés hangzásának folytonosságát a csak diádotokból álló elemkészlethez képest. Meghallgatásos tesztel ellenőriztük az érzelmek felismerhetőségét az egyes logatomokban. A forrásanyag ilyen jellegű formális ellenőrzése a korábbi munkákban nem történt meg. Ellentétben a [4]-ben leírtakkal, a színész szabadon variálhatta a hangmagasságát, amely jelentős változatosságot mutatott a logatomokon belül és azok között, illetve a mondatokon belül is. Emiatt a szintézis folyamán alkalmanként jelentős alaphang módosításra volt szükség, amely nemkívánatos torzulást vitt a szintetizált jelbe. Az elemkiválasztáson alapuló szintézis módszerek is érzékenyek az érzelmes beszédre jellemző jelentős alaphang

változásokra, ami fokozottan jelentkezhet a [2] és [3]-ban használt kisméretű elemkészletek esetében. Az alapprofrendencia változtatás okozta torzulások feltehetően befolyásolták ezen kísérletek eredményeit. Az idézett szerzőkhöz képest több érzelmét vizsgáltunk, ami több választási lehetőséghez vezetett a meghallgatásos tesztek folyamán rontván a várható érzelemazonosítási arányt. Kísérletünket magyar nyelvre végeztük, de a korábbi kísérletekhez hasonlóan a módszer más nyelvek esetében is alkalmazható.

## 2 A hanganyag felvétele

A kísérletben a következő szemantikailag semleges mondatot használtuk: „A menüben minden szükséges információ elhangzik”. Meghatároztuk mondat fonetikus átírásának megfelelő, diádokból és CVC triádokból álló szekvenciát. Összeségében 14 CVC triádot, 3 VV és 9 CC diádot használtunk. A diádokat és triádokat jelentés nélküli szavakba (logatomokba) ágyasztuk, amelyek a CC diádok esetében kétszótagosak, a VV diádok és a CVC triádok esetében három szótag hosszúak voltak.

A semleges mondatot és a 26 logatomot 7 változatban vettük fel a hat alapérzelmét (öröm, harag, meglepetés, undor, szomorúság, félelem) és semlegességet kifejezve. A hanganyagot egy professzionális (30 éves) magyar színész (mondta be, akinek már volt tapasztalata beszédintézéshez szükséges elemkészlet felvételével, illetve korábbi kísérletekhez már mondott fel érzelmes mondatokat. Minden egyes érzelmhez iteratívan több (4-7) változatban kerültek bemondásra a logatomok, amíg a cikk szerzői megfelelőnek nem ítélték őket. A felvétel után informális meghallgatásos tesztet végeztünk, melynek során a cikk két szerzője kiválasztotta az érzelmileg legmeggyőzőbb logatomokat és mondatokat. A logatomok kb. 20%-át nem volt elég kifejező, ezért azokat egy második alkalommal újra bemondattuk. Az első felvétel folyamán készült legmeggyőzőbb mintákat lejátszottuk a színésznek az egyes érzelmek másodszori felvétele előtt. A logatomok több (4-6) változatban is bemondásra kerültek, melyek közül a cikk két szerzője kiválasztotta a legjobban sikerülteket.

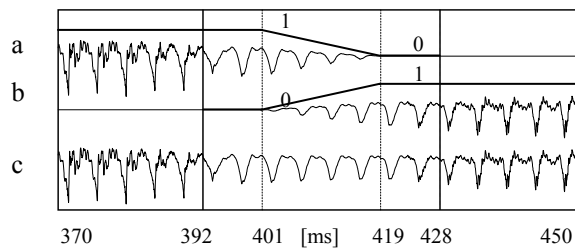
A felvételeket süket szobában készítettük AKG C-414 B-ULS kondenzátor mikrofonnal, a kardoid karakterisztika és a 75 Hz-es aluláteresztő szűrés bekapcsolásával. A mikrofon elé pop szűrőt helyeztünk. A beszédjelet 44.1 kHz-en, mintánként 16 biten digitalizáltuk.

## 3 Szintetizált minták előállítása

A logatomokból készített hét érzelmű beszédanyag adatbázist és a hét változatban bemondott mondatokat az összes lehetséges módon összepárosítottuk. Mind a 49 pár esetén a természetes mondat prozódiját az érzelmű beszédanyag adatbázisból szintetizált mondatra másoltuk.

A szintézist és a prozódia másolását a következő módon végeztük. A hanghatárokat manuálisan bejelöltük a felvett logatomokban és mondatokban. Az alap-

frekvencia menetet a Praat-ban implementált autokorreláción alapuló algoritmus [5] segítségével detektáltuk. A cél lapfrekvencia görbék periódusonként egy frekvencia értéket definiáltak. Az intenzitást 8 ms-onként számítottuk 32 ms széles ablak segítségével. A Praat-ban implementált PSOLA alapú prozódia módosító algoritmust [5] használtuk a hangidőtartamok, az alapfrekvencia és az intenzitások diádok és triádokra másolásához.



**1. ábra:** Szomszédos triádok összefűzése időtartománybeli átmosással (50%-os átfedéssel).

Fontos megjegyezni, hogy a diádok két egymást követő hangot teljesen lefednek, míg a triádok három egymás követő hang teljes időtartamán át tartanak. Mint az 1. ábrán látható, a prozódia módosítás után a diádokat és a triádokat időtartománybeli átmosással fűztük össze, oly módon, hogy a szomszédos szegmensek első, illetve utolsó hangjai között 50% átfedés legyen. Az ábrán a függőleges folytonos vonalak a hanghatárokat, a függőleges szaggatott vonalak pedig az átfedési tartományt jelölik. A felső két görbe (a és b) mutatja az összefűzött triádokat az átlapoló ablak feltüntetésével. Az alsó görbe (c) mutatja az összefűzés után kapott jelet.

## 4 Meghallgatásos teszt

Egy webes felületen végzett meghallgatásos teszt során meghatároztuk, hogy a tesztelők milyen érzelmeket azonosítanak a logatomokban, illetve a természetes és a szintetizált mondatokban. 208 magyar anyanyelvű alany vett részt a vizsgálatban. A tesztalanyok többsége motivált informatika szakos hallgató volt.

14 tesztelő eredményét nem vettük figyelembe, mert vagy nem fejezték be a tesztet, vagy véletlenszerű válaszokat adtak. A kizárt tesztelők közül néhányan hanglejátszási problémát jelzetek. A maradék 194 résztvevő közül 159 férfi és 35 nő volt, átlagos koruk 23 év. 83 tesztelő fej- vagy fülhallgatót használt, míg 111 tesztelő hangszóró segítségével hallgatta meg a felvételeket.

A teszt hat részből állt és átlagosan 18 percet vett igénybe. A résztvevők az első részben a logatomokat, a negyedikben pedig a természetes és szintetizált mondatokat értékelték ki. A második és a negyedik rész egy a jelen tanulmánytól független vizsgálathoz tartozott. A harmadik és a hatodik részt azért iktattuk be, hogy ki tudjuk szűrni a véletlenszerűen válaszoló tesztelőket. Ebben a két részben egy 5-fokozatú skálán kellett minősíteni egy mondat természetes, illetve három különböző

szintetizátorral előállított változatát. Négy tesztelőt kizártunk, mivel az általuk adott minősítések nem tükrözték a felvételek közötti különbséget.

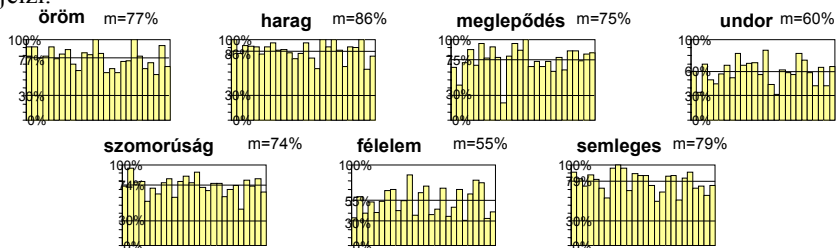
Hogy csökkentsük a tesztelők terhelését, egy résztvevőnek nem kellett az összes mondatot, illetve logatomot meghallgatnia. A 49 szintetizált és 7 természetes mondatot összekevertük és véletlen választással két 28-28 mondatot tartalmazó csoportra bontottuk. 98 tesztelő csak az első, míg 94 csak a második csoportot hallgatta meg. Hasonlóképpen a logatomokat is hét egyforma méretű csoportra osztottuk. Egy ilyen csoportot legalább 23-an hallgattak meg.

A tesztelőknek 7 lehetőség (hat alapérzelem, illetve semlegesség) közül kellett kiválasztani a logatomok, illetve a természetes és a szintetizált mondatok által kifejezett érzelmet. A tesztelőkkel kapcsolatos ismertető szöveget ugyanaz a színésznő olvasta fel, mint aki a hangját adta a kísérlethez, így a résztvevők megismerhették érzelmileg semleges beszédstílusát. Azért, hogy ne lehessen ezt a részt átugorni, a továbblépéshez meg kellett adni egy az ismertető végén elhangzó, véletlenszerűen generált kódot. Az egyes elemek között a tesztelők maguk szabályozták a továbblépést. Egy mintát akárhányszor meg lehetett hallgatni, viszont egy korábban már kiértékelte mintára nem volt megengedett a visszatérés. A minták lejátszási sorrendje tesztelőnként véletlenszerűen változott. A tesztben szereplő hangminták honlapunkon elérhetőek: <http://speechlab.tmit.bme.hu>

## 5 Eredmények

Annak eldöntésére, hogy az azonosítási arányok szignifikánsan a 14%-os véletlen szint felett vannak-e binomiális próbát ( $p < 0.05$ ) használtunk nemegyenlő (1:7) arányokkal. A logatomos teszt esetében (23 tesztelő) a 30% feletti azonosítási eredmény volt szignifikánsan a véletlen szint felett. A mondatok kiértékelésénél (94 tesztelő) a 21% feletti azonosítási arány volt szignifikánsan a véletlen szint felett.

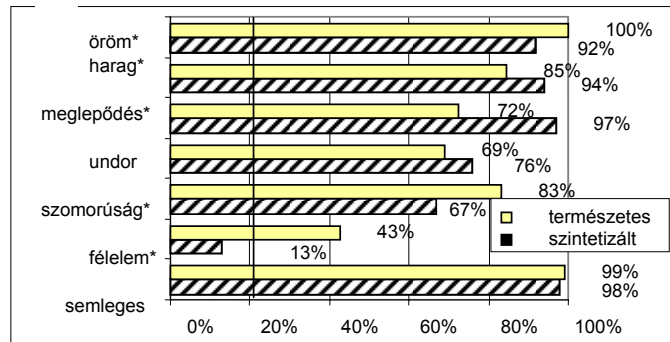
A 2. ábra minden egyes logatomra mutatja a kifejezni szándékolt érzelm felismerési arányát. Az alsó vonal mutatja a szignifikánsan a véletlen szint feletti szintet. A felső vonal az adott érzelmhez tartozó adatbázis átlagos felismerési arányát jelzi.



2. ábra: Érzelmek azonosítási aránya a 26 logatomra.

A kifejezni szándékolt érzelmet a 182 logatomból 180 esetében szignifikánsan a véletlen szint felett azonosították. A harag volt az átlagban legjobban felismert érzelem (86%), míg az undor (60%) és a félelem (55%) sokkal kevésbé volt felismer-

hető. Ugyanakkor több undort és félelmet kifejező logatom 80% felett teljesített, ami mutatja, hogy a tesztelők ezeket az érzelmeket képesek rövid mintákon is felismerni. A többi logatomra kapott alacsonyabb eredmény oka valószínűleg az, hogy a színésznőnek nehéz volt ezen érzelmelek konzisztensen az összes mintán kifejezni.



**3. ábra:** Érzelmek azonosítási aránya a természetes mondatokra, illetve az egyező érzelmi adatbázissal és prozódiaíval előállított szintetizált mondatokra. A szignifikáns különbséget \* jelöli.

A 3. ábra mutatja szándékolt érzelem felismerési arányát természetes mondatokra, illetve az egyező érzelmi adatbázissal és prozódiaíval rendelkező szintetizált mondatokra. Két-mintás t-próba ( $p < 0.05$ ) segítségével ellenőriztük, hogy a természetes és a szintetizált mondatok felismerési arányai között szignifikáns-e az eltérés. A függőleges vonal által jelölt értékénél magasabb felismerési arányok szignifikánsan a véletlen szintje feletti.

A természetes mondatok közül a legjobb eredményt az öröm (100%) érte el. A haragot (85%) semlegesnek ítélte a tesztelők 9%-a. A meglepetést (72%) örömmek ítélte a tesztelők 28%-a. Ennek az lehetett az oka, hogy a színésznő pozitív meglepetést kifejezendő mosolyogva fejezte be a mondatot. A színésznőnek nehéz volt az undor (69%) kifejezése, amit 12% örömmek, illetve 10% haragnak gondolt. A szomorúságot (83%) semlegesnek ítélte a tesztelők 11%-a. A legkevésbé felismert félelmet (43%) a tesztelők 30%-a semlegesnek, míg 22%-a szomorúnak tartotta. A színésznő a mondat második felében abbahagyta a hangmagasság remegtetését, ami magyarázatot adhat a semleges válaszok magas arányára. Ráadásul a színésznő lágy hangon mondta be a mondatot, ami megmagyarázhatja miért tartották azt sokan szomorúnak. A mondat még a félelmet kifejező logatomok átlagánál is rosszabb eredményt ért el, ami szintén mutatja a kevésbé szerencsés realizációt. A semleges mondatot (99%) könnyen azonosították a tesztelők.

A szintetizált örömet (92%) könnyen felismerték. A kapott arány jelentősen meghaladja a [2]-ben és [3]-ban közölt korábbi eredményeket. [2] szerzői feltételezték, hogy a mosolyt tartalmazó, illetve nem tartalmazó beszédekem keveredése csökkenti a szintetizált mondat örömteli hangszínezetét. Az öröm viszonylag alacsony felismerési arányát [3] szerzői részben a kevésbé sikeres színészi teljesítménnyel magyarázták. A természetes örömet kifejező mondatra, illetve logatomokra kapott magas felismerési arányok mutatják, hogy a fenti problémák egyike sem jelentkezett

az általunk végzett kísérletben. A szintetizált haragot (94%) szintén könnyen felismerték. Érdekes módon a szintetizált harag felismerési aránya szignifikánsan meghaladta a természetes harag felismerését. Megfigyeltük, hogy a prozódia módosító algoritmus érdekesebb hangzásúvá tette a beszédet, ami magyarázatot adhat a jobb eredményre. Hasonló megfigyelést [2] is közölt. A szintetizált meglepetést (97%) könnyen felismerték, és a természetes mondattal ellentétben nem tévesztették össze az örömmel. Ez megmagyarázható, ha figyelembe vesszük, hogy a szintetizált mondat hangszínezete (a természetes mondattal ellentétben) nem ment át mosolygásba a mondat végén. A szintetizált undort (76%) a tesztelők 13%-a a haraggal keverte össze. Továbbá a szintetizált felismerési arány meghaladta (de nem szignifikánsan) a természetes mondatét. Mindkét megfigyelés magyarázza a szintetizált mondat érdekes hangszínezete. A szintetizált szomorúságot (67%) a tesztelők 17%-a semlegesnek, míg 7%-uk félelemnek ítélte. Megfigyeltük, hogy a szintetizált szomorú mondat a többi érzelemhez képest torzabban hangzott, ami magyarázza a kapott viszonylag alacsony eredményt. A szintetizált mondaton leginkább eltorzult szakaszokat megvizsgálva észrevettük, hogy az azoknak megfelelő szomorú logatomok szabálytalan zöngé periódusokat (glottalizációt) tartalmaztak. A szomorú beszédelem adatbázis használatával előállított összes mondatban hasonló torzítást találtunk. Továbbá azt is észrevettük, hogy a természetes szomorú mondat is tartalmaz glottalizációt. A színesítő valószínűleg részben glottalizációval fejezte ki a szomorúságot, de a szabálytalan zöngé periódusok megzavarták az alapprozódia detektáló és módosító algoritmusok működését, ami torzulást okozott. A szintetizált félelmet (13%) a véletlen szint alatt ismerték csak fel. A mondatot a tesztelők 64%-a haragnak, 16%-a semlegesnek, 6%-a pedig undornak azonosította. Az alacsony felismerési eredmény részben a természetes félelme kapott alacsony eredménnyel magyarázható. A szintetizált semlegesség (98%) majdnem olyan jó eredményt ért el, mint természetes párja.

Az 1. táblázat mutatja az összes (49) adatbázis és prozódia párosításra kapott felismerési eredményeket. Az egy-egy adott érzelmek legjobb kifejező párosításhoz tartozó cellákat árnyékolással jelöltük meg. Eltérő adatbázis és a prozódia esetén a függőleges nyíl jelzi, ha a mondatot az adatbázis érzelmi tartalmának megfelelően azonosították. A vízszintes nyíl jelzi, ha a mondatot a célprozódiaival megegyezően azonosították. A felismert érzelmek külön kiírtuk, ha az sem az adatbázisnak, sem a prozódiajának nem felelt meg. A félkövérrel jelzett számok jelzik az adott párosításhoz tartozó legmagasabb felismerési arányt. A táblázatban csak a véletlen találgatás szintjét szignifikánsan meghaladó értékeket tüntettük fel.

Az öröm és a félelem kivételével minden érzelmek esetében az egyező adatbázis és prozódia párosítása adta a legjobb eredményt. A félelmet csak 13% azonosította egyező prozódia és adatbázis esetén. Sokkal jobb eredményt hozott, ha természetes félelem mondat prozódiaját a szomorú (62%), vagy a meglepődés (60%) adatbázisból előállított mondatra másoltuk. Az így kapott felismerési arányok még a természetes félelme (43%) kapott eredményt is meghaladták. Megfigyeltük, hogy a félelem adatbázisban szereplő logatomok „présselt” hangzással lettek bemondvá, ellentétben a szomorúságot illetve félelmet kifejező logatomokkal. Ez megmagyarázhatja, miért lett „lágyabb” az utóbbi két adatbázisból szintetizált mondatok hangzása. A 2. ábra alapján a „présselt” logatomok többé-kevésbé a szándékolt érzelmek fejezték ki

(átlag=55%), de a prozódia módosítás által bevitt érdes hangzás a harag felé (64%) tolta el a mondat felismerését.

**1. táblázat:** Érzelmek azonosítási aránya szintetikus mondatokra az összes prozódia és adatbázis párosítás esetén.

↑adatbázis ←prozódia	öröm	harag	meglepetés	undor	szomorúság	félelem	semleges
öröm	<b>92%</b>	← 2%, <b>↑74%</b>	← <b>35%</b> , ↑14% félelem 33%	← 2%, ↑ 14% <b>félelem</b> <b>37%</b>	← 3%, <b>↑ 50%</b> félelem 29%	←11%, <b>↑35%</b> szomorú 23%	←1%, <b>↑ 50%</b> szomorú 23%
harag	←1%, <b>↑ 93%</b>	<b>94%</b>	←12%, ↑12% <b>félelem</b> <b>41%</b>	← <b>63%</b> , ↑ 22%	← 3%, <b>↑39%</b> semleges 27% félelem 26%	← <b>60%</b> , ↑ 2%	←35%, <b>↑ 39%</b>
meglepetés	← <b>68%</b> , ↑ 31%	← <b>56%</b> , ↑41%	<b>97%</b>	← <b>82%</b> , ↑ 5%	← <b>75%</b> , ↑ 5%	← <b>85%</b> , ↑ 1%	← <b>87%</b> , ↑ 4%
undor	←13%, <b>↑ 82%</b>	← 48%, <b>↑49%</b>	←23%, ↑10% öröm 29%	<b>76%</b>	← 14%, <b>↑45%</b>	← <b>71%</b> , ↑ 2%	← <b>68%</b> , ↑ 9%
szomorúság	← 1%, <b>↑ 78%</b>	← 0%, ↑ 33% <b>semleges</b> 34% undor 30%	← 11%, ↑ 12% <b>semleges</b> 37%	←16%, ↑ 43%	<b>67%</b>	← 20%, ↑12% <b>undor</b> 36% semleges 28%	← 32%, ↑ 52%
félelem	← 3%, <b>↑ 74%</b>	←1%, <b>↑87%</b>	← <b>60%</b> , 11%	←6%, ↑22% <b>harag</b> <b>58%</b>	← <b>62%</b> , ↑18%	13% <b>harag</b> <b>64%</b>	← 9%, <b>↑77%</b>
semleges	←22%, <b>↑ 75%</b>	←33%, <b>↑54%</b>	← <b>79%</b> , ↑ 9%	← <b>80%</b> , ↑13%	← <b>79%</b> , ↑13%	← <b>85%</b> , ↑ 0%	<b>98%</b>

Az öröm adatbázis harag prozódiaival kombinálva (93%) némileg magasabb felismerési arányt ért el, mint öröm prozódiaival (92%). Ez azt jelentheti, hogy egy közös prozódiai modell használható ezen érzelmek szintetizálásához.

Ha az 1. táblázatban megvizsgáljuk a semleges prozódiahoz tartozó párosításokat, észrevehetjük, hogy a szintetizált harag (54%) és öröm (75%) kifejezésében inkább az adatbázis játszik szerepet. A semleges adatbázishoz tartozó párosítások alapján látható, hogy a meglepetést (87%) és az undort (68%) inkább a prozódiaival lehet kifejezni. Hasonló eredményre jutottak [2] és [3] szerzői is, az undort (ami nem szerepelt a korábbi kísérletekben) és a szomorúságot kivéve, mely inkább a prozódia-tól függpt mindkét korábbi kísérletben. Az utobbi eltérés egy lehetséges magyarázata az, hogy a mi kísérletünkben a színész nő glottalizáció segítségével fejezte ki a szomorúságot. A szintézishez használt módszer nem tudta kezelni a szabálytalan zöngé periódusokat, ami akadályozta a szomorú prozódia szintetikus semleges mondatra (13%) másolását. Ezt a feltételezést megerősíti a szomorú adatbázis párosítása öröm



(szomorúság 50%), harag (szomorúság 39%) és undor (szomorúság 45%) prozódiaival. Ezen mondatok prozodiáját jobban sikerült a szomorú adatbázisból szintetizált mondatra ültetni, és mindegyik esetben a szomorúság volt a legnagyobb arányban felismert érzelem.

## 6 Összefoglalás

A kísérlet igazolta, hogy a szintetizálással előállított harag és öröm meggyőzően kifejezhető érzelem specifikus adatbázisok segítségével. Az örömteli beszéd elem adatbázis és a haragos prozódia kombinációja meggyőzően fejezte ki az örömet. Ennek alapján elégséges lehet egy közös prozódiai modul használata ezen érzelmek esetében. A meglepetés és az undor pusztán prozódia segítségével is kifejezhető, de a megfelelő hangszínezet tovább növeli felismerhetőségüket.

Az undort és félelmet kifejező természetes mondatok kevésbé voltak meggyőzőek, mint a többi mondat. Ez részben megmagyarázza ezek szintetizált változataira kapott rosszabb felismerési arányokat is. A prozódia módosító algoritmus érdekessége a félelem adatbázisban szereplő erősen „préselt” logatomok hangzását. Az adott technológiai korlátok mellett előnyös lehet kevésbé „préselten” képzett logatomok használata.

A természetes szomorú felvételek szabálytalan zöngperiódusokat tartalmaztak, amelyeket nem tudott megfelelően kezelni a szintézis algoritmus. A szintetikus mintákban megjelenő járulékos torzítás megmagyarázza a gépi szomorúságra kapott rosszabb felismerési arányt. A szintetizált szomorú kifejező ereje javulhat mesterséges glottalizáció hozzáadásával, ugyanakkor a prozódiamódosítás technológiai korlátai miatt a beszéd elem adatbázisokban kerülni kell a glottalizációt.

Az elemösszefűzéses beszéd szintézis segítségével lehetséges meggyőzően érzelmeket kifejezni, viszont a különböző érzelmekhez tartozó beszéd elem adatbázisok felvétele igen megterhelő. Az egyes érzelmekre jellemző hangszínezet feltehetően inkább marad konzisztens egy rövid logatomon, mint egy teljes mondaton belül. Egy hangszínezet módosító algoritmust egy viszonylag kisméretű, de konzisztens logatom adatbázison betanítva akár hosszabb, elemkiválasztáshoz használt beszéd adatbázisokat is érzelmileg kifejezővé lehetne tenni.

## 7 Köszönetnyilvánítás

A munkát a Nemzeti Kutatási és Technológiai Hivatal támogatta (NKFP 2/034/2004).

## **Bibliográfia**

1. Ladd, D.R., Silverman, K., Tolkmitt, F., Bergmann, G., and Scherer, K.R., “Evidence for the independent function of intonation contour type, voice quality, and f0 range in signalling speaker affect”, *Journal of the Acoustic Society of America*, 78 (2), pp. 435–444, 1985.
2. Montero, J.M., Arriola, G.J., Colas, J., Enriquez, E., and Pardo, J.M., “Analysis and Modeling of Emotional Speech in Spanish”, *Proc. of ICPhS*, pp. 957-960, 1999.
3. Bulut, M., Narayanan, S. S., and Syrdal, A. K.: “Expressive Speech Synthesis Using a Concatenative Synthesizer”, In. *ICSLP-2002*, pp. 1265-1268, 2002.
4. Schröder, M., Grice, M., “Expressing Vocal Effort in Concatenative Synthesis”, *Proc. of ICPhS, Barcelona, Spain*, pp. 2589-2592, 2003.
5. Boersma, P., “Praat, a system for doing phonetics by computer”, *Glott International*, vol. 5:9/10, pp. 341-345, 2001.