

Egy általános célú morfológiai annotáció

REBRUS PÉTER – KORNAI ANDRÁS – VARGA DÁNIEL

Bevezetés

Cikkünk a morfológia annotáció általános kérdéseit tárgyalja a magyar példáján keresztül. Az első részben a morfémákra közvetlenül támaszkodó „*konkrét*” annotációs sémák és a *rögzített kódhosszúságú* rendszerek problémáit írjuk le. A második részben a magyar főnévi, igei, és egyéb inflexiók paradigmák részletes kódolásáról írunk. Annotációs rendszerünk, a ***hunmorph*** az említett kódolásokkal szemben *absztrakt, változó kódhossz* használó rendszer, amelynek alapelvei teljesen általánosak és nyelvfüggetlenek. A harmadik részben pedig röviden érintjük a deriváció és a szóösszetétel kezelését. Cikkünk záró részében az annotációs rendszert használó nyílt forráskódú számítógépes nyelvészeti eszközöket ismertetjük.

Az allomorf-alapú annotáció problémái

Kiindulópontunk az, hogy a morfológiai annotáció elsődleges célja az adott szóalakban levő ***morfoszintaktikai információk*** megjelenítése. Morfoszintaktikainak tekintjük a szóalakban meglévő olyan információkat, amelyeknek közvetlen szintaktikai hatása van, azaz amelyek az adott szóalak mondatbeli formai viselkedését (disztribúcióját) befolyásolják – ilyen elsősorban az a szintaktikai ***pozíció***, ahol a szóalak a grammatikus mondatban megjelenhet, illetve az ***egyeztetés***, amikor egy szóalak morfológiai jegyei befolyásolják egy másik szóalak morfológiai jegyeit. Ennélfogva az alábbi módszertani elvet követtük: a kizárólag a ***jelentésre*** és a ***hangalakra*** (vagy az írásképre) vonatkozó információk *nem* részei a morfoszintaktikai reprezentációnak. Egyes esetekben a szemantikai és a szintaktikai információk éles elkülönítése nehézségekbe ütközik, ezért egy általános célú morfológiai annotáció tervezésekor mérlegelnünk kell, hogy a potenciális alkalmazások számára mely szemantikai információk lehetnek lényegesek. A fenti módszertani elvet azért is érdemes szem előtt tartani, mert az annotáció elveinek transzparensnek kell lenniük: egy formai tulajdonságot bárkinek könnyű betanítani (és így emberi erőforrás segítségével előállítani egy nagy pontossággal címkézett korpuszt), míg a szemantikai tulajdonságok nagy részére ez nem áll. Az egyszerre szintaktikai és szemantikai tulajdonságokon alapuló ilyen jegyek körébe tartozik többek között a főnév – melléknév megkülönböztetés vagy igeiknél a modális (ható

ige) és a múlt idő, amelyeket annotációs rendszerünk is megkülönböztet (erről ld. később).

Allomorfia

A hangalakra (fonológiai formára) vagy az írásképre vonatkozó információknak a morfológiai annotációban való megjelenítése azért sem lenne szerencsés, mert nagyon gyakran önkényes döntéseket kellene hozni arról, hogy milyen alakot adjunk meg *allomorfia* esetén (azaz akkor, ha az adott morféma több alakban jelenhet meg). Vegyünk néhány példát: a *fára* alak a következő információkat hordozza: (i) lemmája: FA, (ii) morfoszintaktikai jegyei: SZÁM: EGYES, BIRTOKOS: NINCS, BIRTOK: NINCS, ESET: SUBLATIVUS. Ha az annotációban ezeken a jegyeken túl azt is meg akarnánk jeleníteni, hogy a szóban forgó *fára* alakban a többeli magánhangzó hosszú (szemben más alakokkal, pl. ilyen a *fa* tőalak, a toldalékolt *faként* alak, vagy a *facipő* szóösszetétel), akkor az elemzésben a tövet esetleg a $f\acute{a}$ és nem a f_a alakban adhatnánk meg. A szóban forgó *fára* alakban jelenlevő toldalék azonban előlképzett magánhangzójú változatban is megjelenhet (pl. *kép-re*), így dönthetnénk úgy is, hogy az esetragnak ezt a jellegzetességét az annotációnak tükröznie kell, azaz valamilyen alulspecifikált alakban adhatnánk meg a toldalékot (pl. $-r_A$, ahol a nagy A szimbólum a középnnyílt elülső e és hátulsó a magánhangzók helyett áll). Hasonló a helyzet máskor is, ahol a szóalakban szereplő morfémák allomorfiát szenvednek el. Például a *szelek* vagy a *sarki* alakokban szintén tőallomorfiát találunk: *szél* – *szelek*, *sarok* – *sarki*, sőt az első esetben a többes szám jelölője más szóalakokban más és más alakban jelenhet meg (pl. *kár-ok*, *ház-ak*, *sün-ök*, *zokni-k*), ezért ennek a morfémának a jelölése sem nyilvánvaló (lehetne az előzőhöz hasonlóan alulspecifikált magánhangzóval $-v_k$ vagy magánhangzó nélkül csupán $-k$).

Látható tehát, hogy ha a morfológiai annotációt az allomorfokkal vagy az allomorfoknak valamilyen absztrakt alakjával adjuk meg, akkor az esetek jelentős részében legalább három megoldást követhetünk (természetesen lehetségesek kevert megoldások is): (a) a „*konkrét*” elemzésben az adott szóalakban megjelenő allomorfokat (tulajdonképpen a teljes sztringet eredeti formájában) szerepeltetjük (pl. $f\acute{a}+ra$, $szel+ek$ és $sark+i$); (b) az „*allomorfiamentes*” elemzésben az allomorfook közül a leggyakoribbat vagy az „alap”allomorfot választjuk ki (ilyen, amikor a f_a , a *szél*, és a *sarok* töveket adjuk meg a *fára*, *szelek*, illetve *sarki* alakok elemzésénél); és (c) az „*absztrakt*” elemzés, ahol allomorfia esetén az összes allomorfot lehetőség szerint szerepeltetjük: ez a $-r_A$ és a $-v_k$ toldalékok vagy a $f\acute{A}$, $sz\acute{E}L$ és a $sarOk$ tövek esete, ahol egy alulspecifikált (nagybetűs) szimbólum mutatja a váltakozás helyszíneit (ez lehet nyúlás, rövidülés, hangkivetés, magánhangzó-harmónia stb.). Ez az alulspecifikációs megoldás azonban nem mindig lehetséges: vannak az allomorfiának olyan esetei, amelyben a váltakozó szekvencia nem adható meg alulspecifikált szimbólummal: ilyen ún. nem-fonológiai allomorfiákat találunk az igei paradigmában, ha a toldalékváltozatok között nincs fonológiai kapcsolat (sőt gyakran a szekvenciák

hossza sem azonos): ilyen az. E.2 alakokban a tövégtől függő *sz~ol/el/öl* váltakozás (pl. *kap-sz ~ mos-ol*) az E.3 definit alakokban a tő hangrendjétől függő *ja~i* váltakozás (pl. *lop-ja ~ lep-i*).

Egy további probléma a tő- és toldalékallomorfok azonos alakúságával függ össze: a *szelek* tőalakja a *szél*, viszont van egy másik, nem-rövidülő magánhangzót tartalmazó azonos alakú lexéma, vö. *szél – szélek*. Hasonló igaz a *sarki* alakra: ez lehet a SAROK lexémájához tartozó (pl. *sarki bolt*), de lehet a SARK lexémához tartozó (pl. *sarki expedíció*). Tehát ha ezek az alakok önmagában utalnának az aktuális lexémára, az nem lenne elegendő (ez természetesen más, nem allomorfikus esetben is így van, ekkor a lexémákat a lexikográfiai gyakorlatban sorszámok használatával – pl. ÁR1, ÁR2, ÁR3 – különítik el egymástól). Hasonló homonímia-jelenségek léphetnek fel a toldalékokban is: a *-k* toldalék nemcsak névszók többes számára utalhat (pl. *ház-ak*), hanem igéknél az E.1 (pl. én *kap-j-ak*) és bizonyos esetekben a T.3 (pl. ők *kap-t-ak*) szám/személyre is. Ugyanígy az *-i* toldalék nemcsak melléknévképző (ld. *sarki*), hanem utalhat a birtok többes számára (pl. *hajó-i*, *Pál-é-i*). Tehát a morfoszintaktikai kódoláshoz a toldalékok alakja sem ad elégséges információt. Tanulságos összjátékot mutat a homonímia és a nem-fonológiai allomorfia az olyan alakoknál, mint amilyen az indefinit E.1 *kap-ta-m*, ahol a „szokásos” indefinit E.1 *-k* toldalék helyett *-m* toldalékot találunk. Ekkor a „konkrét” elemzésben (kap+t+am) az *-m* szerepel, ami félrevezető lehet, hiszen az *-m* a szokásos *definit* E.1 toldalékkal azonos (pl. *kap-om*, *kap-j-am* (*azt*)). Az „allomorfiamentes” elemzésben ezzel szemben a „szokásos” *-k indefinit* E.1 végződés szerepelne (kap+t+k), ami szintén félrevezető, hiszen egy E.3 alak pontosan ennek megfelelő alakú (*ők kaptak*). Az „absztrakt” elemzésben viszont szerepeltetni kellene mindkét allomorfot, hiszen a *-k* és *-m* toldalékallomorfokat nem lehet értelmes módon alulspecifikálni: kap+t+m/k. Az alábbi táblázatban összefoglaljuk az említett alakok háromféle elemzését (félkövérrel jelölve a problémás eseteket; a kérdőjel a többféle lehetséges elemzést jelöli valódi morfoszintaktikai különbség nélkül).

(1) Főbb elemzési lehetőségek allomorfia esetén :

	(a) “konkrét” elemzés:	(b) “allomorfiamentes” elemzés:	(c) “absztrakt” elemzés:
<i>fára:</i>	fá+ra	fa+ra	fÁ+rA
<i>szelek:</i>	szel+ek	szél+k (vö. <i>szélek</i>)	?szEl+Vk ?szEl+k
<i>sarki:</i>	sark+i (vö. <i>sark – sarok</i>)	sarok+i, sark+i	sarOk+i, sark+i
<i>kaptam (vmit):</i>	kap+t+am (vö. <i>kapok vmit</i>)	kap+t+k (vö. <i>ők kaptak</i>)	?kap+t+m/k, ...

Az imént bemutatott allomorfia-alapú elemzések tehát több szempontból problematikusak: **(i)** nincs módszertani eszközünk arra, hogy *eldöntsük*, hogy a három ideáltípus közül melyik elemzési módot kövessük (pl. fá+ra vagy fa+ra avagy fa+rA); **(ii)** az absztrakt (és részben az allomorfiamentes) annotáció használata mögött hallgatólagosan olyan vitatott elemzések és így nyelvészeti *elméletek*

kaphatnak szerepet, amelyekről a nyelvtudománynak nincs egységes álláspontja (pl. a kötőhangzó része-e a toldaléknak vagy sem, vagy nemfonológiai allomorfia esetén mely allomorfo(ka)t szerepeltessük); **(iii)** egyes elemzések összemossák a tő- vagy a toldalékallomorfofokban potenciálisan jelen levő **homonímiákat**, így önmagukban nem elegendőek a szóalakban levő morfológiai információk megadásához (ld. pl. a *szélek*, a *sarki* és a *kaptam vmit* alakok fenti esetét).

Szegmentálás

A fentebb bemutatott annotációs megközelítéseknek egy további súlyos következménnyel is szembesülniük kell, ez pedig a morfológiai **szegmentálás** bizonytalansága. A morf-alapú annotációnak tartalmaznia kell egy határjelölőt, amely elválasztja a morfokat egymástól (a fenti hipotetikus elemzésekben erre a célra a + szimbólumot választottuk). Ez az elválasztás azonban sok esetben önkényes és nem ritkán problémákba ütközik. Lássunk néhány példát! A problémás esetek első típusa az **írásképpel** kapcsolatos. A grafémikus alakban a kettőzött digráfok speciális írásmódja miatt nem lehetséges az eredeti szóalak karaktereit megfelelő módon elválasztani; ez történik pl. a *hússzor*, *ésszerű* stb. alakok elemzésénél: a konkrét elemzésben a kettőzött digráfot meg kell osztani a tő és a toldalék között, ami félrevezető (pl. *hús+szor*, *és+szzerű*); az absztraktabb elemzésben viszont nem pontosan a szóalak karakterei találhatóak (pl. *húsz+szor*, *ész+szzerű*). Hasonló a helyzet akkor, ha a szóalak kettős mássalhangzóra végződik, és a toldalék ugyanezzel a mássalhangzóval kezdődik (pl. *szebből*, *halottal*). A következő problematikus típus a morfhataron lezajló **hasonulás**sokkal kapcsolatos. Így például a *-val/-vel* toldaléknak vagy a felszólító mód *-j* toldalékának egyes mássalhangzó utáni változatai esetén nem világos a szegmentálás (*hát+tal*, *hátt+al* vagy *hát+val*, illetve *fus+sa* vagy *fut+ja*). Ha a tő digráfra végződik, akkor a két említett probléma együtt jelentkezik: pl. *ác+csal*, *áccs+al*, *ács+csal* vagy *ács+val*, illetve *ed+dze*, *edz+dze* vagy talán *edz+je*. A következő táblázatban ezeket az elemzési lehetőségeket foglaltuk össze.

(2) Szegmentálási lehetőségek különböző elemzések esetén

	(a) eredeti sztring	(a') "átelemezett" sztring	(b) allomorfiamentes
<i>hússzor</i> :	hús+szor (vö. <i>hús</i>)	<i>húsz+szor</i>	
<i>szebből</i> :	szeb+ból	<i>szebb+ból</i>	
<i>háttal</i> :	hát+tal hátt+al		<i>hát+val</i>
<i>áccsal</i> :	ác+csal áccs+al	<i>ács+csal</i>	<i>ács+val</i>
<i>fussa</i> :	<i>fus+s+a</i>		fut+j+a (vö.: <i>futja</i>)
<i>eddz</i> :	ed+dz	edz+dz	<i>edz+j</i>

A fentiekhez hasonló technikai problémákkal minden morfológiai elemzőprogramnak meg kell birkóznia. Az, hogy egy elemző technikailag melyik módszert követi az aktuális szóalakok (sztringek) manipulációja során, az elemzőprogram (és az erőforrások) felépítésétől, lehetőségeitől függ, és nem a végső kódolásban megjelenítendő morfoszintaktikai információ. Azaz a szegmentálás és az allomorfofok kiválasztása az elemző belügye, és *nem* lehet *része* a morfoszintaktikai *annotációnak*.

Fúzió és szuppletivizmus

Ki kell térnünk egy további problémakörre, amely azt is megmutatja, hogy egyes esetekben a sztringalapú elemzést nem is lehetséges ésszerű módon megvalósítani. Az ún. **fúziós morfémák** esetén több funkció szételemezhetetlenül társul egy morfhoz; legismertebb példa a magyarban a birtokos alakok és igék szám/személy jelölése. Az igazán problematikus esetek azonban azok, amikor a fúzió csak bizonyos esetekben áll fenn, máskor a morfémák agglutinatív módon jelennek meg. Ezt a jelenséget láthatjuk az igei definitjelölésnél: E.1 és E.2 egyértelműen fúziós (pl. *ad-om*, *ad-od*), E.3 és T.2 agglutinatív (*ad-ja*, *ad-já-tok*). Néha még az is előfordul, hogy az E.3 definitjelölés a módjelölővel fuzionál, pl. az *ad-ná* alakban a *-ná* toldalék együtt fejezi ki a feltételes módot és a definitiséget, tehát ezek az allomorfofok nem alkalmasak a morfológiai annotáció jelölésére. A fúzióhoz tartozó jelenség az, amikor a toldaléktömb formailag szételemezhető, viszont nem egyértelmű, hogy mely funkcióhoz mely szekvenciák tartoznak; ilyen a többes számú birtokosjelölős alakok esete (pl. *kalapjaim*), ahol a szételemezett *kalap+ja+i+m* alakban a *ja* szekvencia nem bír morfoszintaktikai szereppel (ennek absztraktabb nyelvészeti elemzését, ld pl. Melcsuk 1965).

A sztringalapú elemzés lehetetlenségét az. ún. **szuppletív** alakok mutatják leginkább, ahol ugyanazon lexémához tartozó alakok töve teljesen különbözik (pl. *van* vs. *lehet*, *jön* vs. *gyere*, *sok* vs. *több* stb.). Hasonló jelenség lép fel egyes kis zárt szóosztályoknál, így a személyes és birtokos névmásoknál: az *engem*, *téged* stb. accusativusi alakok nem állíthatók elő, mint $\text{én}+\text{t}$, $\text{te}+\text{t}$ stb.; teljes szuppletivizmusra példa a *benneteket*, *benniünket* alakok, melyek “tövének” alakja inessivusi, ennek ellenére ezek egyszerű accusativusi alakok: t_i+t , m_i+t . Hasonlóan az *enyém*, *tied* stb. alakok morfoszintaktikailag nem birtokosjelölős, hanem birtokjelölős alakok, tehát morfoszintaktikailag $\text{én}+\text{é}$, $\text{te}+\text{é}$ elemzést kellene, hogy kapjanak (a személyes és birtokos névmások esetéről később részletesen is írunk). A *gyere*, *gyertek* alakok ugyanígy a *JÖN* lexémához tartoznak és kötő-felszólító módúak, annak ellenére, hogy alakilag sem a *tő*, sem az idő/mód jelölő nem látszik.

(3) Fúziós és szuppletív alakok elemzési problémái

	„formai” elemzés	„morfoszintaktikai” elemzés
<i>adná</i>	?ad+na+a	ad+NÁ+JA
<i>kalapjaim</i>	?kalap+ja+i+m	kalap+K+m
<i>engem</i>	?én+m	ÉN+T
<i>enyém</i>	?én+m	ÉN+É
<i>benneteket</i>	?benn+etek+et	TI+t
<i>gyere</i>	?gyer+e	JÖN+J

A következő részben azt tekintjük át, hogy milyen alternatív annotációs megoldás lehetséges.

A kizárólag morfoszintaktikai kategóriákon alapuló annotáció

Általános annotációs elvek

Az előző részben láttuk, hogy a morfológiai annotáció problémájára a megoldás nem a szóalakok (fonológiai vagy grafémikus) formáján alapuló kódolás, hanem egy nyelvészetileg megalapozott morfoszintaktikai kategóriákra épülő formalizmus adhat választ. Egy ilyen elterjedt annotáció az ún. MSD-kódrendszer (*Morphosyntactic Description*, ld Erjavec 1997), amelyben a morfoszintaktikai **kód rögzített hosszúságú**: egy jegyértékekből álló sztring, mely sztring minden pozíciójához eleve rögzített módon vannak hozzárendelve a jegyek: azaz a pozíciók azt adják meg, hogy mely értéknek a jegyeit töltjük ki. Lássunk néhány példát: a *fiú* és a *fiatokéival* főnévi, illetve az *ad* és az *adtatok* igei alakok MSD-kódrendszer szerinti annotációi az alábbiak

(4) Két főnévi és igealak MSD-annotációja

<i>fiú</i>	Nc-sn-y---	<i>ad</i>	Vmip3s---n-----
<i>fiatokéival</i>	Nc-pi-yp2p	<i>adtatok</i>	Vmis2p---y-----

Amint a példákból is kitűnik, ennek a kódolásnak a hátránya az, hogy egyrészt **nehezen kezelhető** (rosszul olvasható a sok üresen hagyott érték és az értékek nem vagy csak kevéssé transzparens kódjai miatt). Másrészt **nem hierarchikus**, azaz az annotáció közvetlenül nem tükrözi az egyes értékek közötti összefüggéseket – például ilyen összefüggés az, hogy csak birtokos alakoknál van szükség a birtokos számának és személyének megjelölésére, vagy az, hogy a magyarban van egy speciális *-lak/-lek* toldalék, amely 2. személyű tárgyra utal, viszont az alanynak E.1-nek kell lennie: pl. *(én) lát-lak (téged/titeket)*; azaz ez a morfoszintaktikai érték függ az ige szám/személyétől. Harmadrészt nem képes a **morfológiai jelöltséget** tükrözni: azaz egy formailag és funkcionálisan komplex szóalak (pl. *fiatokéinak* vagy *adhattatok*) és egy ilyen szempontból jelöletlen szóalak (pl.

fiú vagy *ad*) annotációja ugyanolyan komplexitású. További problémája az, hogy egyelőre **csak inflexiós** kódrendszer, és nem nyilvánvaló, hogy a morfoszintaktikailag releváns képzések hogyan illeszthetők bele (különösen igaz ez a szófajváltó képzésekre).

Jegy-érték szerkezetek

A fenti problémák egy részére megoldást jelent a hierarchikus jegy-érték struktúrák (pl. az ún. AVS-ek, *Attribute-Value Structures*, ld. Trón 2002) használata. Az AVS-ek előnye a nyelvészeti és formális megalapozottság: ezt a formalizmust több szintaktikai elmélet használja. A teljesen kitöltött AVS-eknek is problémája azonban az, hogy az annotáció nem tesz különbséget morfológiailag jelölt és jelöletlen szóalakok között. Lássunk egy példát: a fenti *fiatokéinak* és a *fiú* alak a következő morfoszintaktikai információkat hordozza (itt és a későbbiekben a jegyeket és értékeiket kiskapitálissal jelöltük, ezen belül félkövérrel a jegyeket, és kurzívval az értékeket; a hierarchikus viszonyok jelölésére tabulálást alkalmaztunk).

(5) Két főnévi alak sematikus jegy-érték struktúrája

(a) *fiatokéinak*

LEMMA	<i>FIÚ</i>
KATEGÓRIA	<i>FŐNÉV</i>
SZÁM	<i>TÖBBES</i>
BIRTOKOS	<i>IGEN</i>
	SZÁMA <i>TÖBBES</i>
	SZEMÉLYE <i>2.</i>
BIRTOK	<i>IGEN</i>
	SZÁMA <i>TÖBBES</i>
ESET	<i>DATIVUS</i>

(b) *fiú*

LEMMA	<i>FIÚ</i>
KATEGÓRIA	<i>FŐNÉV</i>
SZÁM	<i>EGYES</i>
BIRTOKOS	<i>NEM</i>
	(SZÁMA <i>X</i>)
	(SZEMÉLYE <i>X</i>)
BIRTOKOS	<i>NEM</i>
	(SZÁMA <i>X</i>)
ESET	<i>NOMINATIVUS</i>

Hasonlóan az említett *adhattatok* és *ad* igék szokásos specifikációja az alábbi.

(6) Két igei alak sematikus jegy-érték struktúrája

(a) *adhattatok*

LEMMA	<i>AD</i>
KATEGÓRIA	<i>IGE</i>
MODÁLIS	<i>IGEN</i>
IDŐ	<i>MÚLT</i>
MÓD	<i>KIJELENTŐ</i>
SZÁM	<i>TÖBBES</i>
SZEMÉLY	<i>2</i>
DEFINITISÉG	<i>IGEN</i>

(b) *ad*

LEMMA	<i>AD</i>
KATEGÓRIA	<i>IGE</i>
MODÁLIS	<i>NEM</i>
IDŐ	<i>JELEN</i>
MÓD	<i>KIJELENTŐ</i>
SZÁM	<i>EGYES</i>
SZEMÉLY	<i>3</i>
DEFINITISÉG	<i>NEM</i>

A fenti (5) és (6) szerkezetekből látható, hogy nincs jelentős különbség a jelölt és jelöletlen alakok jegy-érték struktúrájának „bonyolultsága” között: azok ugyanazokat a jegyeket tartalmazzák. Ez azonban nem intuitív és nem is praktikus, hiszen a morfológiailag jelöletlen alakok általában rövidebbek (több zérusmorfológia vagy morfémat tartalmaznak) és jelentősen gyakoribbak (funkciójuk általánosabb, használatuk kiterjedtebb).

Bináris és unáris jegyek: főnevek

Ha azonban az AVS-ekben a jegyeket úgy fogalmazzuk meg, azok értéke csak az igen/nem (+/-) lehessen, és az értékek közül szisztematikusan az egyik a **jelöltet** (szokásosan a +), a másik a jelöletlent (ez általában a -) jelentse, akkor ezen a **bináris** jegyrendszeren jelentős egyszerűsítést tehetünk (ld. Kornai 1989). Ha megengedjük további jegyek és hierarchia bevezetését, akkor ezt mindig megtehetjük, hiszen többértékű jegyek esetén ezek értékeit mindig átírhatjuk bináris jeggyé (pl. ilyen a SZEMÉLY vagy az ESET jegy a főneveknél vagy az IDŐ vagy a MÓD az igéknél ld. fenti (5), ill (6)). (Vegyük észre, hogy az alakok helyett a jegyek megcímkézése jelölt és jelöletlen értékekre csak akkor tehető meg, ha egy jegyérték jelöltsége nem függ egy másik jegy értékétől, azaz ebben az értelemben környezetfüggetlen. Ez egyes nyilvánvaló esetekben nem igaz, pl. a jelöltség függhet a lexémától: az ún. relációs főneveknél (pl. *barát, anya* stb.) a birtokos alak jelöletlenebb a nem-birtokos alaknál. További ismert eset a felszólító módú igék: itt a 2. személyű alakok – univerzálisan is – jelöletlenebbek, míg más módban általában a 3. személy jelöletlen (pl. a magyarban is E.2 indefinit alak állhat zérus szám/személyjelölővel (pl. *ad-j*), de az E.3 indefinit alak viszont a többi móddal ellentétben toldalékkal áll (*ad-j-on*). Ezek azonban az egész rendszer szempontjából elhanyagolható mértékű hátrányok: elfogadjuk, hogy a jelöltség jegyértékekre való értelmezésével az alakok jelöltsége jól közelíthető).

Lássuk, hogy az (5) és (6)-beli példáink milyen bináris jegyszerkezetet kapnak (az újonnan bevezetett jegyek nyelvészeti értelmezéséről ld. szintén Kornai 1989-et). Az alábbi (7a)-ban a jelölt főnévi alakot látjuk: itt a legtöbb bináris jegyérték pozitív, míg a (7b)-beli alak esetében az összes másodlagos morfoszintaktikai kategória értéke negatív (az áttekinthetőség érdekében csak a pozitív értékkel bíró jegyek vannak félkövérrel szedve). A (7a) és (7b)-beli AVS-ek ugyanazokat az információkat tartalmazzák, mint az (5a), illetve (5b) szerkezetek. Fontos, hogy a negatív értékkel bíró jegyek alá rendelt jegyeknek semmikor nincs szerepük, ez három esetben állhat elő: (i) az alárendelt jegy negatív értékű domináns jegy esetén nem értelmezhető, vagy (ii) a domináns jegy megfogalmazásából következik, hogy az alárendelt jegy (az adott nyelvben) csak negatív értéket vehet fel, vagy (iii) az adott nyelvben az alárendelt jegynek csak a domináns jegy pozitív értéke esetén releváns morfoszintaktikailag. Az (i) esetre példa a BIRTOKOS vagy a BIRTOK jegyek, amelyek negatív értéke esetén – vagyis ha nincs a főnéven birtokos- vagy birtokjelölés – nincs értelme a

birtokos számáról vagy személyéről beszélni (ezt látjuk pl. a *fiú* alak esetén (5b)-ben és (7b)-ben). Egy másik eset a familiáris többes: a *Péterék, szomszédék* stb. alakok morfoszintaktikailag többes számúak, ez azonban egy speciális többes szám: a szóalakban jelölt alakkal familiáris viszonyban álló emberek csoportjára utal; így a FAMILIÁRIS jegyet ésszerű a TÖBBES jegy alá rendelni (erről ld. részletesen Kornai 1989). A (ii) eset akkor áll elő, ha a jegy megfogalmazásából következik, hogy az alárendelt jegy(ek) negatív értékű domináns jegy esetén egyértelműen csak negatív értékeket vehetnek fel. Ilyen jegy a NEM-3 SZEMÉLYŰ és a NEM NOMINATIVUSI ESETŰ jegyek, hiszen ha ezek értéke negatív, akkor a birtokos 3. személyű, illetve az eset nominativusi, így az alárendelt jegyeknek (melyek a további lehetőségeket adják meg) kötelezően negatív értékkel kell rendelkezniük: egy szó a nominativusszal együtt más esettel nem rendelkezhet). A (iii) lehetőségre az igei rendszer bemutatásánál térünk vissza. A (7) ábrában ezeket a „default módon” kitölthető vagy érték nélküli jegy–érték párokat zárójelbe tettük.

(7) Két főnévi alak bináris jegy–érték struktúrája

(a) *fiaitokéinak*

fiú	+
FŐNÉV	+
TÖBBES SZÁMÚ	+
FAMILIÁRIS	–
BIRTOKOS	+
TÖBBES SZÁMÚ	+
NEM-3. SZEMÉLYŰ	+
1. SZEMÉLYŰ	–
2. SZEMÉLYŰ	+
BIRTOK	+
TÖBBES SZÁMÚ	+
NEM NOM. ESETŰ	+
ACCUSATIVUS	–
DATIVUS	+
INSTRUMENTALIS	–
SUPERESSIVUS	–
...	

(b) *fiú*

fiú	+
FŐNÉV	+
TÖBBES SZÁMÚ	–
(FAMILIÁRIS	X)
BIRTOKOS	–
(TÖBBES SZÁMÚ	X)
(NEM-3. SZEMÉLYŰ	X)
(1. SZEMÉLYŰ	X)
(2. SZEMÉLYŰ	X)
BIRTOK	–
(TÖBBES SZÁMÚ	X)
NEM NOM. ESETŰ	–
(ACCUSATIVUS	–)
(DATIVUS	–)
(INSTRUMENTALIS	–)
(SUPERESSIVUS	–)
(...)	

Vegyük észre, hogy ha a morfoszintaktikai információkat tartalmazó jegyeket rögzítjük, akkor bármilyen negatív értékű jegy redundánssá válik, és elegendő csak a pozitív jegyeket megadnunk. Ezt a tulajdonságot felhasználhatjuk arra, hogy a bináris jegyrendszert egyértékűvé (*unáris*sá) tegyük. Ehhez elég a pozitív értékű jegyeket tekintetbe venni, és ha kizárólag ezen jegyek neveit soroljuk fel, akkor egy teljes értékű annotációt kapunk. Az alábbi (8a,b) ábrában ez a hierarchikus unáris jegyrendszer látható, amit úgy kaptunk, hogy a (7a,b) bináris jegyrendszerből elhagytuk a negatív értékű jegyeket és a pozitív értékeket. Ezzel az unáris jegyrendszerrel aztán közvetlenül használható annotációs rendszert jön létre: (8a',b')-ben a hierarchikus rendszert zárójel

segítségével linearizáltuk (az annotációs formalizmus a következő: a lexémát / jel választja el a morfoszintaktikai annotációtól, ez utóbbi a főkategóriával indul, és utána a további morfoszintaktikai jegyek szerepelnek a hierarchiának megfelelően zárójellel; az e mögött álló formalizmusról részletesebben a következő részben írunk).

(8) Két főnévi alak elemzése unáris jegyekkel (a redundáns információk nélkül):

(a) hierarchikus formában:

fiatokéinak

fiú
FŐNÉV (NOUN)
TÖBBES SZÁMÚ (PLUR)
BIRTOKOS (POSS)
TÖBBES SZÁMÚ (PLUR)
NEM-3. SZEMÉLYŰ (--)
2. SZEMÉLYŰ (2)
BIRTOK (ANP)
TÖBBES SZÁMÚ (PLUR)
NEM NOM. ESETŰ (CAS)
DATIVUS (DAT)

(b) hierarchikus formában:

fiú

fiú
FŐNÉV (NOUN)

(a') linearizált formában:

fiatokéinak

fiú/NOUN<PLUR><POSS<PLUR><2>><ANP<PLUR>><CAS<DAT>>

(b') linearizált formában

fiú

fiú/NOUN

A jegyeknek a legvégső formában látható megnevezései az angol nyelvészeti szakirodalomban elterjedt rövidítéseket követik: PLUR: *plural* (többes szám), POSS: *possessive* (birtokos), ANP: *anaphoric possessive* (birtok), CAS: *case* (eset) stb.). A fent vázolt jegyrendszer úgy van tervezve, hogy a lehető legegyszerűbben feldolgozható formában tükrözze a morfológiai jelöltségi viszonyokat: éppen ezért ahol nem szükséges, ott az adott jegyet elhagytuk; ilyen a NEM-3. SZEMÉLYŰ jegy, amelyet a linearizált annotáció nem is jelöl (erre nincs szükség, mert a személyre utaló jegyek amúgy is a BIRTOKOS jegy alá vannak rendelve). A birtokos alakok jelölése így egyszerűbbé válik: a POSS jegy alatti személyre utaló jegyek kétféleképpen lehetnek: <POSS<1>> vagy <POSS<2>>. A 3. személyű birtokos alakokban a POSS jegy az 1 és a 2 jegy nélkül szerepel, azaz jelölése <POSS>, ez egybevágg azzal a megfigyeléssel, hogy a három szám/személy közül a 3. a jelöletlen. A birtokosjelölővel ellátott alakok sémája tehát a következő.

(9) Birtokos alakok annotációja

<i>fiam</i>	fiú/NOUN<POSS<1>>
<i>fiad</i>	fiú/NOUN<POSS<2>>
<i>fia</i>	fiú/NOUN<POSS>
<i>fiunk</i>	fiú/NOUN<POSS<PLUR><1>>
<i>fiatok</i>	fiú/NOUN<POSS<PLUR><2>>
<i>fiuk</i>	fiú/NOUN<POSS<PLUR>>

Megemlítjük, hogy a PLUR jegyre a hierarchia három különböző helyén is szükség van: közvetlenül a főkategória-jegy (NOUN) alatt (ekkor a lemmában megadott entitás többes számát jelzi), a POSS alatt (ekkor az entitást birtokló birtokos többes számát jelzi), és az ANP alatt (ekkor az entitás által birtokolt birtok többes számát jelzi) – a hierarchikus elrendezés azonban biztosítja, hogy ugyanannak a PLUR jegynek a használata nem vezet félreértéshez, hiszen ezek más jegyek alatt helyezkednek el, amit a linearizált kódban a zárójelezés mutat. Ezt mutatják az alábbi alakok, ahol a PLUR különböző pozíciókban külön-külön és egyszerre is megjelenhet (itt megjegyzendő, hogy a birtok többes számának jelzése a beszélt köznyelvben állítmányi helyzetben nem kötelező, sőt egyes beszélőknél tiltott: pl. *A könyvek a %fiúéi / %fiúé*).

(10) A PLUR jegy különböző használatai

mi többes számú?

nincs birtokos- és birtokjelölés:

<i>fiúk</i>	fiú/NOUN<PLUR>	(entitás)
-------------	----------------	-----------

csak birtokosjelölés (itt 3. személyű):

<i>fiai</i>	fiú/NOUN<PLUR><POSS>	(entitás)
<i>fiuk</i>	fiú/NOUN<POSS<PLUR>>	(birtokos)
<i>fiaik</i>	fiú/NOUN<PLUR><POSS<PLUR>>	(entitás és birtokos)

csak birtokjelölés:

<i>fiúké</i>	fiú/NOUN<PLUR><ANP>	(entitás)
<i>fiúéi</i>	fiú/NOUN<ANP<PLUR>>	(birtok)
<i>fiúkái</i>	fiú/NOUN<PLUR><ANP<PLUR>>	(entitás és birtok)

birtokos- és birtokjelölés is (csak azok, ahol a birtok többes számú):

<i>fiáéi</i>	fiú/NOUN<POSS><ANP<PLUR>>	(birtok)
<i>fiaiéi</i>	fiú/NOUN<PLUR><POSS><ANP<PLUR>>	(entitás és birtok)
<i>fiukéi</i>	fiú/NOUN<POSS<PLUR>><ANP<PLUR>>	(birtokos és birtok)
<i>fiaikéi</i>	fiú/NOUN<PLUR><POSS<PLUR>><ANP<PLUR>>	(entitás, birtokos és birtok)

Itt kell kitérnünk a többes szám egy speciális használatára: a *familiáris többes* alak morfoszintaktikailag többes számú, de nem a lexémával kifejezett entitás többes számára, hanem az azzal valamilyen „familiáris” viszonyban levők összességére (család, ismerősök stb.) utal: pl. *sógorék, szomszédék* stb. Ez a viszony kombinálódhat a birtokosjelölős alakokkal (pl. *sógorodék*) és a

birtokjelölős alakokkal (*sógoréké*). Ezért az annotáció egy a NOUN alatti PLUR általi dominált FAM jegy segítségével történik.

(11) Familiáris többes alakok

<i>fiúk</i>	fiú/NOUN<PLUR<FAM>>	az entitás familiáris csoportja
<i>fiáék</i>	fiú/NOUN<PLUR<FAM>><POSS>	a birtokolt entitás fam. csoportja
<i>fiúéké</i>	fiú/NOUN<PLUR<FAM>><ANP>	az entitás fam. csoportjának birtoka
<i>fiáéké</i>	fiú/NOUN<PLUR<FAM>><POSS><ANP>	a birtokolt entitás fam. csoportjának birtoka

Az esetek kódolása is megfelel a morfológiai jelöltségnek: mivel a jelöletlen eset a nominativus, ezért az alanyesetű alakokat külön nem jelöljük, a többi 17 eset kódolására az elterjedt latin elnevezéseik három betűs rövidítéseit használjuk. A CAS jegy azt jelzi, hogy itt egy jelölt (azaz nem nominativusi) alakkal van dolgunk. A 18 eset annotációja és az esetek elnevezése az alábbi.

(12) Az esetek annotációja

„strukturális esetek”

<i>fiú</i>	fiú/NOUN	nominativus
<i>fiút</i>	fiú/NOUN<CAS<ACC>>	accusativus
<i>fiúnak</i>	fiú/NOUN<CAS<DAT>>	dativus

„lexikális esetek”

helyhatározói

forrás

<i>fiúról</i>	fiú/NOUN<CAS>	delativus
<i>fiúból</i>	fiú/NOUN<CAS<ELA>>	elativus
<i>fiútól</i>	fiú/NOUN<CAS<ABL>>	ablativus

hely

<i>fiún</i>	fiú/NOUN<CAS<SUE>>	superessivus
<i>fiúban</i>	fiú/NOUN<CAS<INE>>	inessivus
<i>fiúnál</i>	fiú/NOUN<CAS<ADE>>	adessivus

cél

<i>fiúra</i>	fiú/NOUN<CAS<SBL>>	sublativus
<i>fiúba</i>	fiú/NOUN<CAS<ILL>>	illativus
<i>fiúhoz</i>	fiú/NOUN<CAS<ALL>>	allativus
<i>fiúig</i>	fiú/NOUN<CAS<TER>>	terminativus

egyéb

<i>fiúval</i>	fiú/NOUN<CAS<INS>>	instrumentalis-comitativus
<i>fiúért</i>	fiú/NOUN<CAS<CAU>>	causalis-finalis
<i>fiúként</i>	fiú/NOUN<CAS<FOR>>	formativus
<i>fiúvá</i>	fiú/NOUN<CAS<TRA>>	translativus-factivus
<i>húsvétkor</i>	húsvét/NOUN<CAS<TEM>>	temporalis

Bináris és unáris jegyek: igék

Az igék specifikációjában az eddigi elveknek megfelelően a hierarchikus elrendezést a bináris jegyekkel kombináljuk. Az említett *adhattátok* és *ad* igealakok kétértékű jegyekkel való annotációja a következő.

(13) Két igealak bináris jegy-érték struktúrája

(a) *adhattátok*

ad	+
IGE	+
MODÁLIS	+
NEM JELEN.KIJ	+
MÚLT IDŐ	+
FELTÉTELES M.	-
KÖTŐ-FELSZ M.	-
INFINITÍVUSZ	-
TÖBBES SZÁM:	+
NEM 3. SZEMÉLY+	
1. SZEMÉLY	-
	(TÁRGY 2. SZEMÉLYŰ X)
2. SZEMÉLY	+
DEFINITSÉG	+

(b) *ad*

ad	+
IGE	+
MODÁLIS	-
NEM JELEN.KIJ	-
	(MÚLT IDŐ -)
	(FELTÉTELES M. -)
	(KÖTŐ-FELSZ M. -)
	(INFINITÍVUSZ -)
TÖBBES SZÁM	-
NEM 3. SZEMÉLY	-
	(1. SZEMÉLY -)
	(TÁRGY 2. SZEMÉLYŰ X)
	(2. SZEMÉLY -)
DEFINITSÉG	-

(14) Elemzés unáris jegyekkel (a redundáns információk nélkül):

(a) hierarchikus formában:

adhattátok

ad
IGE (VERB)
MODÁLIS (MODAL)
NEM JELEN.KIJ (--)
MÚLT IDŐ (PAST)
TÖBBES SZÁM (PLUR)
NEM 3. SZEMÉLY (PERS)
2. SZEMÉLY (2)
DEFINITSÉG (DEF)

(b) hierarchikus formában:

ad

ad
IGE (VERB)

(a') linearizált formában:

adhattátok

ad/VERB<MODAL><PAST><PLUR><PERS<2>><DEF>

(b') linearizált formában

ad

ad/VERB

Az igék idő/módjának annotációja úgy történik, hogy közvetlenül a VERB jegy alatt szerepel az erre vonatkozó információ (azaz a NEM JELEN.KIJELENTŐ MÓDÚ jegy a linearizált formából hiányzik). A

jegyrendszer felépítése biztosítja, hogy a zérusmorfémát tartalmazó jelöletlen jelen idő kijelentő módú alakok nem kapnak külön jelölést.

(15) Az igék négy idő/módjának annotációja

<i>ad</i>	ad/VERB	jelen idő kijelentő mód
<i>adott</i>	ad/VERB<PAST>	múlt idő kijelentő mód (<i>past</i>)
<i>adna</i>	ad/VERB<COND>	jelen idő feltételes mód (<i>conditional</i>)
<i>adjon</i>	ad/VERB<SUBJUNC-IMP>	kötő-felszólító mód (<i>subjunctive-imperative</i>)

Az igei személyjelölés annotációja a főnévi birtokos mintát követi: a jelöletlen 3. személy jegy nélkül áll, a 1. és 2. személyek jegyei a PERS jegy alatt szerepelnek. A speciális, csak E.1. személyű igéknél megfigyelhető, 2. személyű tárgyra utaló *-lak/lek* toldalékos alakok annotációja egy a <1> jegy alá bevezetett <OBJ2> jeggyel történik. A definit–indefinit (általános–határozott) igealakok megkülönböztetésére a DEF jegy szolgál, hiszen a definit alakok a morfológiailag jelöltek. A számjelölés a független <PLUR> jeggyel történik

(16) Igei indefinit és definit szám/személyjelölés annotációja

<i>adok</i>	ad/VERB<PERS<1>>	<i>adom</i>	ad/VERB<PERS<1>><DEF>
<i>adlak</i>	ad/VERB<PERS<1<OBJ2>>>		
<i>adsz</i>	ad/VERB<PERS<2>>	<i>adod</i>	ad/VERB<PERS<2>><DEF>
<i>ad</i>	ad/VERB	<i>adja</i>	ad/VERB<DEF>
<i>adunk</i>	ad/VERB<PLUR><PERS<1>>	<i>adjuk</i>	ad/VERB<PLUR><PERS<1>><DEF>
<i>adtok</i>	ad/VERB<PLUR><PERS<2>>	<i>adjátok</i>	ad/VERB<PLUR><PERS<2>><DEF>
<i>adnak</i>	ad/VERB<PLUR>	<i>adják</i>	ad/VERB<PLUR><DEF>

Az infinitívusz szám/személy jelölésének annotációja igen hasonló az igékéhez. Az infinitívusz, mint jegye az igéknek (VERB<INF>) sajátos jegykombinációkat enged csak meg: az infinitívusznak nincsen idő/módja és definitése, viszont lehet szám/személye, melyet a <PERS> jeggyel fejezünk ki (az %*adhatni* és az ?*adnalak* típusú infinitívuszi alakok is csak periferiálisan léteznek). Az egyetlen jelentős eltérés az igék annotációjához képest, hogy a <PERS> jegy hiánya ebben az esetben nem a 3. személyű alakot (pl. *adnia*), hanem a szám/személyjelölés nélküli alakot (pl. *adni*) kódolja. (Ez fontos különbség, mert az infinitívuszt vonzó igék közül azok, amelyek szám/személyjelöléssel rendelkeznek, kizárólag a szám/személyjelölés nélküli infinitívuszt engedik meg: pl. *Dolgozni(*a) akar.*)

(17) A szám/személyjelölővel nem rendelkező és az azzal rendelkező infinitívusz annotációja

<i>adni</i>	ad/VERB<INF>		
<i>adnia</i>	ad/VERB<INF><PERS>	<i>adniuk</i>	ad/VERB<INF><PLUR><PERS>
<i>adnom</i>	ad/VERB<INF><PERS<1>>	<i>adnunk</i>	ad/VERB<INF><PLUR><PERS<1>>
<i>adnod</i>	ad/VERB<INF><PERS<2>>	<i>adnotok</i>	ad/VERB<INF><PLUR><PERS<2>>

Összefoglalva, az unáris jegyekkel való hierarchikus ábrázolás lehetőséget teremt arra, hogy egyszerűen megfogalmazható és nyelvészeti alátámasztott morfoszintaktikai jegyek segítségével olyan annotációt adjunk, amely teljes és általában véve tükrözi a morfológiai jelöltségi viszonyokat. Azaz anélkül, hogy közvetlenül hivatkoznunk kellene az elemzett szóalak formai tulajdonságaira (allomorfolk, szegmentálás stb.) az annotációs kód változó hosszúságú: hossza nagyjából megfelel a szóalak morfológiai komplexitásának. Ez azt is jelenti, hogy zérusmorfémák esetén az annotáció – mivel bináris jegyeik mind negatívak – kizárólag a lexémából és a főkategória címkéjéből állnak. Minden további morféma tovább növeli az annotáció bonyolultságát. Az alábbi táblázatban néhány ilyen „monotonon bővülő” komplexitású alaksort és annotációikat adtunk meg a zérustoldaléktól a maximális alakokig (az összehasonlítás kedvéért a hozzávetőleges morfémahatárokat a szóalakokban jelöltük).

(18) Főnévi és igei szóalakok és annotációik egy-egy monoton növvő komplexitású sora

<i>fiú</i>	fiú/NOUN
<i>fiú-k</i>	fiú/NOUN<PLUR>
<i>fi-a-i</i>	fiú/NOUN<PLUR><POSS>
<i>fi-a-i-d</i>	fiú/NOUN<PLUR><POSS<2>>
<i>fi-a-i-tok</i>	fiú/NOUN<PLUR><POSS<PLUR><2>>
<i>fi-a-i-tok-é</i>	fiú/NOUN<PLUR><POSS<PLUR><2>><ANP>
<i>fi-a-i-tok-é-i</i>	fiú/NOUN<PLUR><POSS<PLUR><2>><ANP<PLUR>>
<i>fi-a-i-tok-é-i-t</i>	fiú/NOUN<PLUR><POSS<PLUR><2>><ANP<PLUR>><CAS<ACC>>
<i>ad</i>	ad/VERB
<i>ad-hat</i>	ad/VERB<MODAL>
<i>ad-hat-ott</i>	ad/VERB<MODAL><PAST>
<i>ad-hat-t-ak</i>	ad/VERB<MODAL><PAST><PLUR>
<i>ad-hat-t-atok</i>	ad/VERB<MODAL><PAST><PLUR><PERS<2>>
<i>ad-hat-t-á-tok</i>	ad/VERB<MODAL><PAST><PLUR><PERS<2>><DEF>

Fontos rámutatni, hogy az annotáció semmilyen értelemben nem használja az alulspecifikációt (unáris jegyek esetén ez nem is lehetséges), azaz nem lehetséges megadni úgy egy morfoszintaktikai leírást, hogy az valamilyen értékre ne legyen meghatározva – ez bináris vagy többértékű jegyeket alkalmazó rendszerekben egyszerűen a szóban forgó jegy értékének kitöltetlenül hagyásával történhet. Mivel minden annotáció a morfofonológiai értékekre nézve teljesen specifikált, ezért a potenciálisan alulspecifikáltként kezelhető eseteket kétértelműségként

kell kezelniük. Ilyen eset a magyarban meglehetősen ritka. például az E.1 és E.2 birtokos alakok esetjelölés nélkül jelenthetnek nominativus vagy accusativust; vagy egyes igealakok a definitiség mindkét értékét felvehetik. Néhány példa.

(19) Morfológiai kétértelműségek kezelése alulspecifikáció nélkül

<i>fiam</i>	fiú/NOUN<POSS<1>> fiú/NOUN<POSS<1>><CAS<ACC>>	nominativus (pl. <i>A fiam látott engem.</i>) accusativus (pl. <i>Láttam a fiam.</i>)
<i>adtam</i>	ad/VERB<PAST><PERS<1>> ad/VERB<PAST><PERS<1>><DEF>	indefinit (pl. <i>Egy almát adtam neki.</i>) definit (pl. <i>Az almát adtam neki.</i>)

Az inflexiós annotáció

Formalizmus:

Az itt következő részben pontosítjuk az inflexiós jegyrendszernek azt a formalizmusát, amelyet az előző részben mutattunk be. Formálisan az **inflexiós annotáció** két komponensből áll, az egyik komponens a jegy-érték struktúra, amelyben a **bináris morfoszintatikai jegyek** és ezeknek pozitív vagy negatív **értékei** szerepelnek. A másik komponens a hierarchiáért felelős, ezt a legegyszerűbb egy **irányított körmentes gráfként** (azaz irányított faként) meghatározni, melyben minden csomóponthoz egy bináris jegy-érték pár van rendelve, az irányított élek pedig megfelelnek a jegy-érték párok közötti dominancia viszonyoknak. Mivel ez a gráf egy fa, ezért összefüggő és egy csomópont (a gyökércsomópont) kivételével minden csomóponthoz van olyan csomópont, amelyik őt közvetlenül dominálja; a körmentesség pedig azt biztosítja, hogy ne lehessen egy csomópontnak több közvetlenül domináns csomópontja. A jegy-érték párokkal címkézett gráfra egy további feltételnek kell teljesülnie: csak a **pozitív** értékkel rendelkező jegy-érték párok csomópontjai **dominálnak** más csomópontokat (azaz a negatív értékkel címkézett csomópontok a fában levelek lesznek). Ez a feltétel az előző részben elmondottak alapján lehetővé teszi, hogy a bináris jegyes hierarchikus szerkezet unáris jegyessé alakítható legyen a hierarchia megtartásával, és így a(z unáris) jegyek száma tükrözze a morfológiai jelöltséget. (Valójában az annotációt közvetlenül unáris jegyekkel címkézett fagráffal is definiálhatnánk, ekkor egy annotáció ennek a jegyekkel címkézett fának az olyan részfája lenne, amelynek a gyökércsomópontja megegyezik a bővebb fáéval.)

Ahogy az előző részben láttuk, a **gyökércsomópont** tartalmazza az inflektált szóalak **kategóriáját** (szófaját, POS (*part-of-speech*)-címkéjét): a gyökércsomópont egy olyan jegy-érték párral van címkézve, ahol a jegy valamely főkategória-jegy (az előző részben ezek közül a NOUN és a VERB szerepelt). A hunmorph annotációs rendszer aktuális változata által használt főkategória-jegyek listája megtalálható a függelék (A1) ábrájában. Minden inflektálható kategóriához tartozik

egy rögzített inflexiós jegy–érték struktúra, azaz bináris jegy–értékekkel címkézett csomópontú fa-gráf. Inflektálható kategória azonban csak öt van: a három névszói és a ragozható determinánsi és az egy igei kategória, ezek jegy–érték szerkezeteiről ld. az előző, illetve a következő részt.

Az inflexiós annotáció linearizálása úgy történik, hogy a pozitív értékkel bíró jegyeket írjuk le a megfelelő zárójelezéssel. Mivel egy fában az ugyanazon csomópont által dominált csomópontok (az ún. testvércsomópontok) egymás közötti sorrendje lényegtelen, ilyen esetekben a linearizálás az összes sorrendben lehetséges. Praktikus okokból azonban a jegyek sorrendjét úgy rögzítettük, hogy a félreolvasás lehetősége a lehető legkisebb legyen (az inflektálható kategóriákhoz tartozó jegyek kimerítő listáját és sorrendjüket ld. a következő részben). Így a linearizált annotáció már egy egyértelmű, kódok és zárójelekből álló sztring lesz.

Mivel a linearizált kód – jelöletlen szóalak esetén – egyetlen főkategória-jegyből is állhat, ezért fontos megjegyeznünk, hogy elvi különbség van egy főkategória-jegy és az ilyen "rövid" inflexiós annotáció között. Például a `NOUN` és a `<NOUN>` két különböző dologra utal: az első egy **jegy** *neve*, amely a gyökércsomópontban állhat; a második egy morfoszintaktikailag **teljesen specifikált** alak, azaz jegyekkel címkézett fa-gráf, amelynek minden főnévi jegye negatív (azaz esetünkben az egyes számú nem-birtokos nem-birtok nominativusi alak, ld. (7b), illetve unáris formában (8b, b')); hasonlóan igékre, ld. (7a), illetve (8a, a')). Ezt a különbséget a végső linearizált kódban azonban nem használjuk: a főkategória (és így az egész morfoszintaktikai jegyrendszer) praktikus okokból mindig külső zárójelek nélkül szerepel – ez nem vezethet félreértéshez, hiszen az annotáció úgyszólván mindig teljes elemzést ad vissza. Az elemzés általános formája – amely már az előző részből ismerős – a következő:

(21) Az inflexiós annotáció sémája

szóalak lemma/FŐKATEGÓRIA<TOVÁBBI_INFLEXIÓS_JEGYEK>

Morfológiai elemzésnek általánosan egy olyan hozzárendelést nevezünk, amely minden egyes jólformált szóalakhhoz (sztringhez) hozzárendel egy lexéma–annotáció párt. Ez a hozzárendelés azonban nem egyértelmű (nem függvény), mivel ugyanahhoz a szóalakhhoz több különböző elemzést is rendelhet morfológiai homonímia esetén -- ld. pl. a (19)-beli eseteket. A hozzárendelés megfordítása (inverze) sem függvény, mert ugyanolyan lexémának ugyanolyan annotációval különböző szóalakok felelhetnek meg: ez a helyzet áll elő morfofonológiai ingadozás esetén (pl. *fotelban* -- *fotelben* vagy *fürdenek* -- *fürödnek*), vagy olyan alakoknál, ahol a szuppletív tő megjelenése nem kötelező (pl. *jöjj* --- *gyere* vagy *volna* -- *lenne*).

(22) Ingadozó alakok azonos annotációt kapnak

<i>fotelban</i>	fotel/NOUN<CAS<INE>>
<i>fotelben</i>	fotel/NOUN<CAS<INE>>
<i>fürdenek</i>	fürdik/VERB<PLUR>
<i>fürödnek</i>	fürdik/VERB<PLUR>
<i>gyere</i>	jön/VERB<PERS<2>><SUBJUNC-IMP>
<i>jöjj</i>	jön/VERB<PERS<2>><SUBJUNC-IMP>
<i>jöjjél</i>	jön/VERB<PERS<2>><SUBJUNC-IMP>

A következőkben a korábbi főnévi és igei elemzéseket kiegészítjük a többi inflektálható elem annotációjával.

Névszói kategóriák

Régi problémája a leíró nyelvtanoknak, hogy be lehet-e (és ha igen, hogyan) sorolni egyértelműen a névszói alakokat valamelyik névszói kategóriába (ld. többek között Moravcsik 1997). A melléknevek és a számnevek a főnevekkel átfedő osztályokat alkotnak, és nehéz egyértelmű disztribúciós tesztek adni, amelyek alapján ezek a kategóriák egyértelműen megkülönböztethetőek lennének. Ezen a helyzeten a morfológiai vizsgálatok sem segítenek, mivel mind a melléknevek, mind a számnevek felvehetik az összes főnévi inflexiót egyes „elliptikus” és „nominalizáló” kontextusokban: pl. *Nem szeretem a kíváncsiakat.*, *Ez az én nagy labdám, az meg a te kicsid.*, *Bátraké a szerencse.*, *Összeültek a nyolcak.*, *Négyet rendeltem.*, *Az ő öt könyve meg az én háromam.* Itt és a többi hasonló példában vitatható, hogy az adott melléknév vagy számnév a saját „prototipikus” mondattani *funkciójában* szerepel-e, de melléknév, illetve számnév voltak mellett számos érv szól. Nyilvánvaló, hogy a mondatokban különböző funkciókban álló ugyanazon elemek megkülönböztetése nem lehet a feladata egy csak szóalakokat vizsgáló morfológiai elemzőnek, és így az annotációnak sem. Így például a *pék barátom* és a *szomszéd Józsi* -féle szerkezetekben az első főnév módosító szerepű (ahogyan tipikusan a melléknevek), a *szépek imádata* és a *kevés is sok*-féle szerkezetekben a melléknév, illetve a számnév főnévi jellegű (birtokos szerkezeten belül, illetve alanyként áll); ezt a tényt azonban nem érdemes az adott alakok többszófajúsága mellett felhozni, mert akkor a névszók jelentős többségével ezt kellene tennünk, és így értelmetlenül sok többszörös annotációt kapnánk. (A kizárólag deadjektíválisnak tartott képzések, mint amilyen a közép- és felsőfok, sem adnak jobb fogódzót, ezek a mellékneveken kívül egyes számnevekkel is lehetségesek (pl. *több*, *kevesebb*, *legelső*), és egyes konstrukciókban főnevekhez is: pl. *székebb a széknél*.)

A hunmorph kategóriarendszerének összeállításánál arra is figyelemmel kellett lennünk, hogy az elérhető elektronikus adatbázisok (pl. szótárak) és a rendelkezésre álló elemzett korpuszok (pl. a *Szeged Korpusz*, ld. Csendes 2004) valamilyen módon mégis megkülönböztetik a három fő névszói kategóriát (ezt nagyon sokszor nem formai–disztribúciós, hanem szemantikai–funkcionális

alapokon teszik). Ezért az információvesztés elkerülése végett érdemes ezt a kategorizációt megtartani. A három névszói kategória morfoszintaktikai jegyrendszere viszont azonos lesz: bármely névszó felveheti az összes főnévi inflexiós kategóriát. Az alábbi néhány példa inflektált alakokra.

(23) Melléknévi és számnévi alakok névszói inflexiókkal

<i>kíváncsi</i>	kíváncsi/ADJ
<i>kíváncsijaitokét</i>	kíváncsi/ADJ<PLUR><POSS<PLUR><2>><ANP><CAS<ACC>>
<i>kétezer</i>	kétezer/NUM
<i>kétezeinkével</i>	kétezer/NUM<PLUR><POSS<PLUR><1>><ANP><CAS<INS>>

Determinánsok

A negyedik inflektálható kategória, a determinánsok (DET), ld. (20). Pontosabban a determinánsok egy része inflektálható, az olyan szerkezetekben, mint pl. *ezeké a lányoké, abban a házban*. Más részük viszont nem inflektálható, pl. *e lányoké, ama házban, azon gondolatoknak*. Az inflektálható determinánsok inflexiós jegyszerkezetükben megegyeznek a többi névszóval. (Meg kell jegyeznünk, hogy a szokásosan a determinánsok közé számított névelők a hunmorph-ban külön kategóriát képeznek (ART), amit rendkívül gyakori előfordulásuk és speciális funkciójuk indokol – ide csupán három lemma tartozik: *a, az, egy*). Néhány példa determinánsokra (az utolsókét felsorolt típus *ezen, azon* stb.) kétértelmű: lehet inflektálhatatlan determináns, de lehet superessívusi esetű inflektálható is: vö. *azon emberekkel* vs. *azon az embereken*.

(24) Inflektált és nem inflektált determinánsok

<i>emez</i>	emez/DET
<i>ugyanazokéval</i>	ugyanaz/DET<PLUR><ANP><CAS<INS>>
<i>e</i>	e/DET
<i>azon</i>	azon/DET
	az/DET<CAS<SUE>>

Főnévi, melléknévi, számnévi névmások

A névmások a hunmorph rendszerben nem képeznek külön kategóriát (szemben a más alapokon nyugvó annotációkkal, pl. a már említett MSD-kódrendszerrel). A disztribúciós elemzés (és funkcionális megfontolások is) azt az elképzelést támogatják, miszerint a névmások szétoszthatók a négy névszói (NOUN, ADJ, NUM, DET) és a határozószerkezet (ADV) kategóriák között. Hely hiányában a névmások elemzésére itt részletesen nem tudunk kitérni, álljon itt néhány példa a hagyományos besorolásuk szerint:

(25) főnévi, melléknévi és számnévi névmások annotációja

mutató

<i>ez</i>	ez/NOUN
<i>azokéval</i>	az/NOUN<PLUR><ANP><CAS<INS>>
<i>ilyen</i>	ilyen/ADJ
<i>olyanjainak</i>	olyan/ADJ<PLUR><POSS><CAS<DAT>>
<i>ennyi</i>	ennyi/NUM
<i>annyinkat</i>	annyi/NUM<POSS<PLUR><2>><CAS<ACC>>

kérdő

<i>micsoda</i>	micsoda/NOUN
<i>kikét</i>	ki/NOUN<PLUR><ANP><CAS<ACC>>
<i>melyik</i>	melyik/ADJ
<i>milyeneken</i>	milyen/ADJ<PLUR><CAS<SUE>>
<i>hány</i>	hány/NUM
<i>mennyivel</i>	mennyi/NUM<CAS<INS>>

egyéb (vonatkozó, általános, tagadó)

<i>amely</i>	amely/NOUN
<i>valakijeitékét</i>	valaki/NOUN<PLUR><ANP<PLUR>><CAS<ACC>>
<i>bármelyik</i>	bármelyik/ADJ
<i>semmilyenekkel</i>	semmilyen/ADJ<PLUR><CAS<INS>>
<i>mindahány</i>	mindahány/NUM
<i>akármennyiért</i>	akármennyi/NUM<CAS<CAU>>

Személyes névmások

A hagyományosan személyes és birtokos névmásoknak nevezett szóosztály annotálása érdekében az eddig bemutatott névszói annotációs jegyrendszert kismértékben ki kell bővítenünk. A személyes névmások annotációs rendszerünk szerint speciális főnevek, melyeknek névszói inflexiók jegyeik lehetnek (alakjuk nagyon gyakran szuppletív, pl. *engem, bennünket, velük, rá*). A különböző személyű személyes névmásokkal való egyeztetési jelenségek indokolják, hogy a névszói jegyrendszert kiegészítsük az – igéknél ismert – a személyre utaló PERS jeggyel. Ez a PERS jegy az infinitívuszoknál látott módon jelöli a személyt (ld. (17)): magában a <PERS> 3. személyre utal, míg az e jegy által dominált személyjegyekkel az 1., illetve 2. személyre. Ekkor a személyes névmások annotációja a következő (a formális – „önöző”, illetve „magázó” – személyes névmásokat is szerepeltettük, ezek morfoszintaktikailag 3. személyűek):

(26) A személyes névmások annotációja

<i>én</i>	én/NOUN<PERS<1>>	<i>mi</i>	mi/NOUN<PLUR><PERS<1>>
<i>te</i>	te/NOUN<PERS<2>>	<i>ti</i>	ti/NOUN<PLUR><PERS<2>>
<i>ő</i>	ő/NOUN<PERS>	<i>ők</i>	ők/NOUN<PLUR><PERS>
<i>ön</i>	ön/NOUN<PERS>	<i>önök</i>	önök/NOUN<PLUR><PERS>
<i>maga</i>	maga/NOUN<PERS>	<i>maguk</i>	maguk/NOUN<PLUR><PERS>

A személyes névmások esetekkel ellátott alakjai között több morfofonológiailag kivételes, illetve szuppletív alak van (ld. pl. (3)), ezenkívül néhány alak a legtöbb alak a nem-formális személyes névmásoknál hiányzik (TRA: **énné* FOR, **teként*, TER: **őig*, TEM: **önkor*), illetve többszörös alakváltozatok is előfordulnak; néhány példa:

(27) Inflektált személyes névmások

<i>engem engemet</i>	én/NOUN<PERS<1>><CAS<ACC>>
<i>neked néked</i>	te/NOUN<PERS<2>><CAS<DAT>>
<i>vele véle</i>	ő/NOUN<PERS><CAS<INS>>
<i>önhöz</i>	ön/NOUN<PERS><CAS<ADE>>
<i>magáig</i>	maga/NOUN<PERS><CAS<TER>>
<i>bennünket minket</i>	mi/NOUN<PLUR><PERS<1>><CAS<ACC>>
<i>bennetek</i>	ti/NOUN<PLUR><PERS<2>><CAS<INE>>
<i>rajtuk</i>	ők/NOUN<PLUR><PERS><CAS<SUE>>
<i>önökké</i>	önök/NOUN<PLUR><PERS><CAS<TRA>>
<i>magukként</i>	maguk/NOUN<PLUR><PERS><CAS<FOR>>

Birtokos névmások

Az ún. „birtokos” névmások nem birtokosjelölővel, hanem *birtok*jelölővel vannak ellátva, hiszen nem a személyes névmás által kifejezett személy birtokosát, hanem annak birtokát jelölik, és szintaktikai disztribúciójuk is ennek felel meg: *A könyv a fiúé / tied / övé*. Ezért ezek annotációja az ANP jeggyel történik (megjegyzendő, hogy a POSS jegy főnévi személyes névmásokra nem is használatos, hiszen az ezt kifejező alakok szisztematikusan hiányoznak: **éned*, **öm*, **önötök*, **magánk*):

(29) Az egyes és többes számú birtokra utaló névmások annotációja

<i>enyém</i>	én/NOUN<PERS<1>><ANP>	<i>enyémek</i>	mi/NOUN<PERS<1>><ANP<PLUR>>
<i>tied tied tiedé</i>	te/NOUN<PERS<2>><ANP>	<i>tieid tiedek</i>	ti/NOUN<PERS<2>><ANP<PLUR>>
<i>övé</i>	ő/NOUN<PERS><ANP>	<i>övéi</i>	ő/NOUN<PERS><ANP<PLUR>>
<i>öné</i>	ön/NOUN<PERS><ANP>	<i>önéi</i>	ön/NOUN<PERS><ANP<PLUR>>
<i>magáé</i>	maga/NOUN<PERS>	<i>magáéi</i>	maga/NOUN<PERS><ANP<PLUR>>
<i>mienk miénk</i>	mi/NOUN<PLUR><PERS<1>><ANP>	<i>mieink</i>	mi/NOUN<PLUR><PERS<1>><ANP<PLUR>>
<i>tietek tiétek</i>	ti/NOUN<PLUR><PERS<1>><ANP>	<i>tieitek</i>	ti/NOUN<PLUR><PERS<2>><ANP<PLUR>>
<i>övék övéké</i>	ők/NOUN<PLUR><PERS><ANP>	<i>övéik</i>	ők/NOUN<PLUR><PERS><ANP<PLUR>>
<i>önöké</i>	önök/NOUN<PLUR><PERS><ANP>	<i>önökéi</i>	önök/NOUN<PLUR><PERS><ANP<PLUR>>
<i>maguké</i>	mi/NOUN<PLUR><PERS><ANP>	<i>maguk</i>	maguk/NOUN<PLUR><PERS><ANP<PLUR>>

A birtokos névmások viselkedése jól példázza azt, hogy a birtokjelölés bizonyos erősen korlátozott esetekben és módon egy alakon belül megismételhető. A többszörös birtokviszonyok lerövidítésére szolgálnak az olyan szerkezetek, mint *Az én kutyám pórása – Kinek a pórása? – Az én kutyámé / Az enyéme / ?# Az enyém*. Ez utóbbi alak nem egyszerű birtokltságot, hanem a birtok

általi újabb birtokoltságát fejez ki, és bizonyos mértékig itt is kifejezhető mindkét birtok többes száma: *A kutyáim pórása – az enyémeké.*; *A kutyám pórásai -- ?az enyémei.*; *A kutyáim pórásai -- az ?enyémekéi.* (Megjegyezzük, hogy ugyanez a szerkezet nem-névmási alakokkal csak nagyon nehezen fogadható el: *?*Pálée*, talán kicsit grammatikusabbá válik, ha az első birtokjelölő után többes számjelölés van: *??Páléié*). Ezt a szerkezetet az ANP csomópont alatti újabb ANP csomóponttal lehet kódolni.

(30) Többszörös birtokjelölés

<i>enyémé</i>	én/NOUN<PERS<1>><ANP<ANP>>
<i>tieidé</i>	te/NOUN<PERS<2>><ANP<PLUR><ANP>>
<i>mienkéi</i>	mi/NOUN<PLUR><PERS<1>><ANP<ANP<PLUR>>>
<i>tieitekéi</i>	ti/NOUN<PLUR><PERS<2>><ANP<PLUR><ANP<PLUR>>>

A birtokos névmások szintén felvehetnek további főnévi inflexiós jegyeket (a POSS jegy kivételével). Így az előbbiekkal és az esetjelöléssel megkaphatjuk a maximálisan komplex PERS jegyet tartalmazó főnévi alakot (az alábbi listában az összehasonlítás végett a szóalakokban bejelöltük a morféma-k hozzávetőleges határait):

(31) Inflektált birtokos névmások növekvő komplexitású sora

(<i>ő</i>	ó/NOUN<PERS>)
<i>öv-é</i>	ó/NOUN<PERS><ANP>
<i>eny-é-m</i>	én/NOUN<PERS<1>><ANP>
<i>ti-e-i-d</i>	te/NOUN<PERS<2>><ANP<PLUR>>
<i>mi-e-i-nk</i>	mi/NOUN<PLUR><PERS<1>><ANP<PLUR>>
<i>ti-e-i-tek-et</i>	ti/NOUN<PLUR><PERS<2>><ANP<PLUR>><CAS<ACC>>

Névutók

A névutók külön főkategória (POSTP), de bizonyos névutós alakok érintik a főnévi annotációt is. Az ún. „személyragozott névutók” (*utánam, eléd, nélküle* stb.) olyan elemek, amelyek formája névutó + szám/személyjelölő, de szintaktikai viselkedésük a főnév+névutó sémát követi (pl. *A fiú nélkül / Nélkülem jött el.*). Ez problematikus, mivel a morfológiai elemző csak egyes szavakat képes annotálni, a szavak között fennálló konstrukciók elemzése közvetlenül nem feladata. Vegyük azonban észre, hogy a főnév+névutó szerkezetek szintaktikailag az esetjelölős főnévi szerkezetekkel rokoníthatók (vö. az előző mondatokat a következőkkel: *A fiúval / Velem jött el.*). Így a személyragozott névutókat tekinthetjük speciális névutós személyes névmásoknak, és mivel a névmások PERS jeggyel ellátott főnevek, ezért érdemes bevezetni a főnévi inflexiós rendszerbe a POSTP jegyet (amely a POSTP főkategória jegytől különbözik). A főnév alatti POSTP jegynek aljegye lesz az összes névutós jegy, mely jegyek nevei a névutó lemmájával azonosak. Ekkor a

személyragozott névutók annotációja az esetjelölős személyes névmásokéval rokonítható. Figyeljük meg, hogy a személyrag nélküli névutók és a személyraggal ellátott névutók különböző főkategóriához tartoznak: az első névutó (POSTP), a második főnév (NOUN), ami természetes, hiszen szintaktikai disztribúciójuk nagyban eltér: a POSTP kategóriát közvetlenül megelőzi egy főnév (pl. *fiú mögött*), míg a névutós NOUN kategória határozóként vagy vonzatként áll a mondatban (pl. *mögöttünk*). Megfigyelhetjük azt is, hogy a 3. személyű személyragozott névutók lemmája sohasem lehet a formális önözés/magázás ön, illetve maga lemmája: *előtte* értelmezése 'öelőtte', és nem 'ön/maga előtt', és hasonlóan a T.3 *közéjük* 'öközējük' és nem 'önök/maguk közé'.

(32) Személyragozott és sima névutók annotációja

<i>nélkül</i>	nélkül/POSTP
<i>nélküled</i>	te/NOUN<POSTP<NÉLKÜL>><PERS<2>>
<i>közül</i>	közül/POSTP
<i>közülük</i>	ők/NOUN<POSTP<KÖZÜL>><PLUR><PERS>
<i>mellé</i>	mellé/POSTP
	ő/NOUN<POSTP<MELLÉ>><PERS>
<i>melléje</i>	ő/NOUN<POSTP<MELLÉ>><PERS>
<i>számára</i>	számára/POSTP
	ő/NOUN<POSTP<SZÁMÁRA>><PERS>
<i>számunkra</i>	mi/NOUN<POSTP<SZÁMÁRA>><PLUR><PERS<1>>

Figyeljük meg, hogy egyes névutók kétértelműek: egyrészt vannak olyanok, amelyek E.3 alakjuk megegyezhet az inflektálatlan alakkal (pl. *mellé*, *mögé*); másrészt azok, amelyek birtokjelölős alakúak, az E.3 személyben kétértelműek, pl. a *számára* funkciója kettős: egyrészt névutó (pl. *Pál számára*), másrészt az E.3 névutós névmási alak: 'öszámára' (pl. *Száma nincsen kegyelem.*). Bizonyos esetekben a mutató névmás és a névutó egybeírható (pl. *anélkül*, *ezelőtt*, *amiatt*). Ekkor az annotáció szintén a főnévi kategóriát használja a POSTP jeggyel.

(33) Személyragozott és sima névutók annotációja

<i>anélkül</i>	az/NOUN<POSTP<NÉLKÜL>>
<i>ezelőtt</i>	ez/NOUN<POSTP<ELŐTT>>

A névszói alakok teljes inflexiós specifikációja a függelék (A2) táblázatában található.

Ige

Az igei alakok jelentős rész morfológiailag defektív, azaz egyes lemmák bizonyos – egyébként megengedett – jegykombinációkkal nem állnak. Ilyenek többek között a személytelen igeik,

amelyeknek minden idő/módjuk megvan, de csak E.3 indefinit alakban állhatnak: ilyen a *hajnalodik* ige vagy a *kell*, *lehet* segédigék. A *rejlík*, *történik* stb. igeiknek vannak többes számú alakjaik is, de csak 3. személyben (és indefinitként) állnak. Az intranszitiv igeiknek nem vagy csak periferikusan van definit alakjuk (pl. *kimosakodja magát*, *lejárja a távot* stb.), de van néhány ige, amelyeknek egyáltalán nincs (pl. *jön*, *megy*, *van*). Van olyan ige, amelynek formailag csak múlt ideje van, más idő/módban, infinitívusban nem állhat, ilyen a *szokott* segédige, vagy pedig csak kijelentő mód 3. személyben áll: pl. *nincs*, *nincsenek*. A *fog* segédigének viszont nincs sem kifejezett idő/módja, sem infinitívusza. A legtöbb ige azonban minden idő/módban, szám/személyben és definit, illetve indefinit alakban is állhat. Az igeik teljes annotációs specifikációja a függelék (A3) ábrájában látható.

Képzés és szóösszetétel

A képzés annotációja

Bizonyos feladatok (tartalomelemzés, gépi fordítás) megoldásához szükség lehet az inflexiós réteg kielemezésén túli morfológiai elemzésre is. Tipikus példa erre a szintaktikai elemzés (parsing), hiszen a melléknévi igeves szerkezetek a *Koreában terjedő bankók* főnévi csoport helyes elemzése megköveteli, hogy a *terjedő* ne csak mint melléknév kerüljön elemzésre (hisz ezt az elemző az őt követő főnévhez csatolná) hanem ki tudjuk elemezni a *terjed* igeot is, hogy a *Koreában* alakot mint ennek bővítményét tudjuk elemezni. Képzőnek azt a toldalékot tekintjük, amely főkategóriához egy másik főkategóriát rendel. Ennek az elvnek a következménye, hogy a képzés bemenete nem éríthet inflektált alakot, morfológiai megfogalmazásban: inflektált alakok nem képezhetők tovább. Képzett alakoknak viszont kötelező inflexiót felvenniük (ha a kimeneti kategória inflektálható). Ezért a képzés annotációját formálisan olyan irányított gráfként tudjuk megadni, melynek csomópontjai a főkategóriák: egy képző két (nem feltétlenül különböző) **kategóriacímke közötti irányított él**. Az elemzésben a képzés alapjául szolgáló lemma és a lemma kategóriája mellett meg kell adnunk a képző elnevezését és a kimeneti kategóriát, amely a következő módon történik.

(38) Az elemzés formalizmusa képzett alak esetén

szóalak lemma/LEMMA_KATEGÓRIA[KÉPZŐ]/VÉGSŐ_KATEGÓRIA<INFL_ANNOTÁCIÓ>

Az alábbiakban néhány példát adunk meg különböző kategóriákból történő képzések elemzésére.

(39) Példák inflektált képzésekre

<i>terjedő</i>	terjed/VERB[IMPERF_PART]/ADJ	(foly. mn-i igenév)
<i>csináltathatnék</i>	csinál/VERB[CAUS]/VERB<MODAL><COND><PERS<1>>>	(kauzatív)
<i>székestül</i>	szék/NOUN[COM]/ADV	(társ)
<i>lábuákat</i>	láb/NOUN[INAL_ATTRIB]/ADJ<PLUR><CAS<ACC>>	(elidegeníthetetlen tul.)
<i>oroszul</i>	orosz/ADJ[MANNER]/ADV	(mód)
<i>okosabbjaival</i>	okos/ADJ[COMPAR]/ADJ<PLUR><PERS><CAS<INS>>	(középfok)
<i>sokszor</i>	sok/NUM[MULTIPL-ITER]/ADV	(multiplikatív–iteratív)
<i>nyolcadikak</i>	nyolc/NUM[ORD]/NUM<PLUR>	(sorszám)
<i>szembeni</i>	szemben/POSTP[ATTRIB]/ADJ	(tulajdonság)
<i>utániakról</i>	után/POSTP[ATTRIB]/ADJ<PLUR><CAS>	

Többszörös képzés esetén a teljes képzési gráf kerül linearizálásra az előbbihez hasonló módon, az alábbiakban az általános sémát és példákat adunk meg.

(40) Többszörös képzést annotáló elemzés

szóalak

lemma/KATEGÓRIA_1[KÉPZŐ_1]/...KATEGÓRIA_n[KÉPZŐ_n]/VÉGSŐ_KATEGÓRIA<INFL_ANNOTÁCIÓ>

(41) Példák többszörös képzésekre

<i>láthatósági</i>	lát/VERB[MODAL_PART]/ADJ[ABSTRACT]/NOUN[MET_ATTRIB]/ADJ
<i>faxolgatás</i>	fax/NOUN[ACT]/VERB[FREQ]/VERB[GERUND]/NOUN
<i>lányosabban</i>	lány/NOUN[ATTRIB]/ADJ[COMPAR]/ADJ[MANNER]/ADV
<i>tárgyasít</i>	tárgy/NOUN[MET_ATTRIB]/ADJ[ATTRIB]/ADJ[TRANS_RESULT]/VERB

A VERB, NOUN, ADJ, NUM kategóriák között minden irányban lehetséges képzés (a denumerális ige képzés kivételével). Az ADV kategória nem igazán képezhető, a POSTP kategóriából egyetlen melléknévképzés lehetséges. A függelék (A4) táblázatában megtalálható a jelenleg használt képzők listája.

Szóösszetételek formalizmusa

A szóösszetételek annotációjához szükség van az összetételi tagok elemzésére. Szerencsére inflektált összetételi tag tipikusan csak az összetétel utolsó eleménél fordul elő, ezért az összetételi tagokat elegendő csak az (esetlegesen előforduló képzési annotációval együtt) kategóriájukkal elemezni. A formalizmust és néhány példát az alábbiakban találunk.

(42) Az elemzés formalizmusa szóösszetétel esetén

Szóalak lemma_1/KÉPZÉS_ELEMZÉSE1+...+lemma_n/KÉPZÉS_ELEMZÉSE<INFL_ANNOTÁCIÓ>

(43) Példák összetételek elemzésére

inflektált összetételek:

<i>eladja</i>	e1/PREV+ad/VERB<PERS>
<i>vérfarkasok</i>	vér/NOUN+farkas/NOUN<PLUR>
<i>sötétzöldje</i>	sötét/ADJ+zöld/ADJ<POSS>
<i>kanárisárgáé</i>	kanári/NOUN+sárga/ADJ<ANP>

képzett összetételi tagok:

<i>zúzottkő</i>	zúz/VERB[PERF_PART]/ADJ+kő/NOUN
<i>háromszínű</i>	három/NUM+szín/NOUN[INAL_ATTRIB]/ADJ

többszörös összetétel:

<i>birsalmasajt</i>	birs/NOUN+alma/NOUN+sajt/NOUN
---------------------	-------------------------------

Alkalmazások

A BME MOKK kutatócsoportjának irányításával számos olyan nyelvtechnológiai eszköz fejlesztése történt, amelyek a cikkünkben bemutatott morfológiai reprezentációt felhasználva végeznek automatikus szövegfeldolgozási műveleteket. Az alábbiakban bemutatjuk ezek közül a cikkünk szempontjából a legfontosabbakat.

Morfológiai leírás

A *hunlex* formalizmus egy morfológiai leírások reprezentálására alkalmas formális nyelv. A nyelvész szakértő ebben írhatja le egy adott nyelv morfológiai szabályait, illetve az ezekhez kapcsolódó morfológiai szótárakat. A *hunlex* reprezentáció specifikálja egy morfológiai elemző (vagy generáló) program viselkedését az adott nyelvre, azaz szóalakokhoz morfológiai annotációkat rendel. Munkatársaink számos nyelvre dolgoztak ki morfológiai leírást ebben a formalizmusban, ezek közül a *morphdb.hu* erőforrás a jelen cikkben bemutatott magyar nyelvű annotációs rendszer megvalósítása (Trón 2006). A szintén *hunlex* névre hallgató szoftver a nyelvészek által könnyen olvasható és módosítható magas szintű morfológiai erőforrásból egy úgynevezett aff/dic (affixumlista/szótár) formátumú, "gépközeli" erőforrást állít elő. Ez az egyszerű formalizmus szabványnak tekinthető, a Firefox és OpenOffice szoftverek komponenseként több száz millió számítógépre telepített hunspell helyesírásellenőrző rendszer is ezt használja a morfológiai információ leírására.

Morfológiai elemzés

Egy morfológiai elemző-szoftver feladata egy szóalakhoz megadni az összes legális morfológiai elemzést. Több olyan morfológiai elemző szoftverimplementáció is létezik (*rfst*, *ocamorph*, *jmorph*), amely az aff/dic formátumú morfológiai erőforrásra támaszkodva elvégzi ezt a feladatot (Trón 2005).

Morfológiai egyértelműsítés

A *hunpos* morfológiai egyértelműsítő azt a problémát orvosolja, hogy egy adott szóalakhoz a morfológiai elemző gyakran több legális elemzést is talál (Halácsy 2007). Az egyértelműsítő ezek közül választja ki a mondatbeli szövegekörnyezet alapján legvalószínűbbet. Ennek a feladatnak a hibátlan megoldásához a mondat teljes megértése lenne szükséges, amely a technológia mai állása mellett természetesen nem lehetséges. Ugyanakkor magyar nyelvre a morfológiai egyértelműsítés feladata – automatikusan feltárt egyszerű statisztikus szabályszerűségek kiaknázásával – elfogadható (tokenenként 95% feletti) pontossággal megoldható.

Tulajdonnév-felismerés (NER) és a főnévi csoportok felismerése (NP-chunking)

A morfológiai egyértelműsítés hasznos előfeldolgozási lépésként szolgálhat magasabb szintű nyelvfeldolgozási műveletek elvégzése előtt. Ilyen rendszerre példa a *hunner* tulajdonnév-felismerő, amely a szövegben automatikusan azonosítja és klasszifikálja a személyneveket, helyneveket, szervezetneveket és egyéb tulajdonneveket (Varga 2007). A jelen kötetben bemutatott *hunchunk* NP-daraboló (főnévi csoport felismerő) is támaszkodik működése során a morfológiailag feldolgozott szövegre (ld. Recski–Varga tanulmányát a jelen kötetben).

Gyakorisági korpusz

Korpusznyelvészeti és pszicholingvisztikai kutatásokhoz hasznos segédeszköz a *Szószablya* webes ***gyakorisági szótár*** (Szalai 2008). Ennek alapja a magyar nyelvű webről reprezentatív módon kigyűjtött több, mint fél milliárd tokennyi szöveg automatikus morfológiai egyértelműsítése nyomán épített gyakorisági táblázat. A keresések eredményei a legkülönbözőbb morfológiai és fonetikai szempontok alapján szűrhetők és rendezhetők. A morfológiai reprezentáció hierarchikus mivolta megkönnyíti, hogy a gyakorisági szótár nyelvész felhasználója kereséseihez a legkülönbözőbb módon specifikálhasson morfoszintaktikai relációkat.

Bibliográfia

Csendes, D, J Csirik, T Gyimóthy (2004) The Szeged Corpus: A POS tagged and syntactically annotated Hungarian natural language corpus, Lecture notes in computer science, ISSU 3206, pages 41-48, Springer

Erjavec, T. and Monachini, M. (szerk.) (1997). Specifications and Notation for Lexicon Encoding. Deliverable D1.1 F. Multext-East Project COP-106. <http://nl.ijs.si/ME/CD/docs/mte-d11f/>.

Szalai András, Halácsy Péter (2008). Szószablya Gyakorisági Szótár. <http://szotar.mokk.bme.hu/szoszablya/>

Halácsy, P, A Kornai, Cs Oravecz (2007) HunPos – an open source trigram tagger In.: Proceedings of the ACL 2007 Demo and Poster Sessions, pages 209-212, Prague 2007. Association for Computational Linguistics

Kornai András (1989) A főnévi csoport egyeztetése. In Telegdi and Kiefer, editors, Általános Nyelvészeti Tanulmányok, XVII. Akadémiai Kiadó, Budapest.

Kornai András, Rebrus Péter, Vajda Péter, Halácsy Péter, Rung András, Trón Viktor (2004). Általános célú morfológiai elemző kimeneti formalizmusa. In: Alexin Zoltán, Csendes Dóra (szerk.). II. Magyar Számítógépes Nyelvészeti Konferencia. SZTE, Szeged, p. 172-176.

Melcsuk, Igor (1965) A magyar főnév birtokos személyragjainak morfológiai felépítéséről. *Magyar Nyelv* **61**: 264–275.

Moravcsik, Edith (1997) Hungarian adjectives from a typological point of view. Kézirat. University of Wisconsin-Milwaukee, Milwaukee.

Recski Gábor, Varga Dániel (ebben a kötetben) Magyar főnévi csoportok azonosítása. Általános Nyelvészeti Tanulmányok.

Trón Viktor (2002) Attribútum–érték struktúrák. In Kálmán László, Trón Viktor és Varasdi Károly (szerk.) *Lexikalista elméletek a nyelvészetben*. Tinta Könyvkiadó, Budapest.

Trón, V, Németh, L, Halácsy, P, Kornai, A, Gyepesi, G, and Varga, D (2005). Hunmorph: open source word analysis, In: Proceedings of the ACL. ACL 2005.

Trón, V, P Halácsy, P Rebrus, A Rung, P Vajda és E Simon (2006) Morphdb.hu: Hungarian lexical database and morphological grammar, In: Proceedings of 5th International Conference on Language Resources and Evaluation. ELRA, pages 1670-1673.

Varga Dániel, Simon Eszter. Hungarian named entity recognition with a maximum entropy approach. *Acta Cybernetica* Volume 18 Issue 2. 2007.

Függelék

(A1) A hunmorph főkategória-jegyei

(i) inflektálható (és továbbképezhető)

névszóként:

NOUN	főnév (<i>noun</i>)
ADJ	melléknév (<i>adjective</i>)
NUM	számnév (<i>numeral</i>)
DET	determináns (<i>determiner</i>) [8]

igeként:

VERB	ige (<i>verb</i>)
------	---------------------

(ii) nem inflektálható (és nem továbbképezhető)

ADV	határozószó (<i>adverb</i>)
POSTP	névutó (<i>postposition</i>) – marginálisan esetjelölés és mn-képzés lehetséges
ART	névelő (<i>article</i>) [3]
CONJ	kötőszó (<i>conjunction</i>) [198]
PREP	prepozíció (<i>preposition</i>) [5]
PREV	igekötő (<i>preverb</i>) [105]
UTT-INT	mondatszó/indulatszó (<i>utterance/interjection</i>) [273]
ONO	hangutánzó (<i>onomatopoeic</i>)

(iii) egyéb, írott nyelvi kategória

PUNCT	központozási jel (<i>punctuation</i>)
-------	---

(A2) A névszói inflexiók jegyhierarchia

{NOUN ADJ NUM DET}

[<POSTP	(POSTP jegy csak NOUN esetén)
<ELÓTT>	
...	(az összes névutó)
>]
<PLUR	
<FAM>	
>	
<PERS	
<1>	
<2>	
>	
<POSS	
<PLUR>	
<1>	
<2>	
>	
<ANP	
<PLUR>	
<ANP	
<PLUR>	
>	
>	
<CAS	
<ACC>	
<DAT>	
...	(összesen 17 jelölt esetrag)

(A2b) Nemlétező névszói jegykombinációk

egymást kizáró testvérjegyek:

- *<PERS<<1><2>>> személyes névmás nem lehet egyszerre 1. és 2. személyű
- *<POSS<<1><2>>> birtokos nem lehet egyszerre 1. és 2. személyű
- *<CAS<<A>>> egyszerre két esetjelölés nem lehetséges
- *<POSTP<<A>>> egyszerre két névutó nem állhat egy „főnévi” alakon belül

kötelezően folytatandó jegyek:

- *<CAS> jelölt (nem-nominativusi) esetről meg kell jelölni a konkrét esetet
- *<POSTP> a névutónál meg kell jelölni a konkrét névutót

(A3) Az igei inflexiók jegyhierarchia

```
VERB
  <MODAL>
  <PAST>
  <COND>
  <SUBJUNC-IMP>
  <INF>
  <PLUR>
  <PERS <1
    <OBJ2>
  >
  <2>
  >
  <DEF>
```

(A3b) Nemlétező igei jegykombinációk

egymást kizáró testvérjegyek:

- <PAST> | <COND> | <SUBJUNC-IMP1> | <INF> a jelölt idő/módok és az inf. páronként kölcsönösen kizárók
- *<PERS<<1><2>>> ige nem lehet egyszerre 1. és 2. személyű

kötelezően folytatandó jegyek:

- *<-INF><PERS> véges ige esetén a nem-3. személy esetén a személyt meg kell jelölni

(A4) A hunmorph képzői (a más elemzőkben gyakran inflexióként elemzettek kiemelve)

	VERB	NOUN	ADJ	ADV
VERB	-gat [FREQ] -tat [CAUS] -ódik [MEDIAL] 	-ás [GERUND]	-ó [IMPERF_PART] -ott [PERF_PART] -andó [FUT_PART] -ható [MODAL_PART] -hatatlan [NEG_MODAL_PART] -atlan [NEG_PERF_PART]	-va [PART] -ván [PERF_PART]
NOUN	-ozik [ACT] -ol [ACT2] -kodik [REG_ACT] 	-né [MRS] -cska [DIMIN]	-os [ATTRIB] -i [MET_ATTRIB] -jú [INAL_ATTRIB] -tlan [NEG_ATTRIB] -mentes [NEG_ATTRIB2] -nyi [QUANTITY] -szerű [TYPE1] -féle [TYPE2] -beli [LOC_INE]	-stul [COM] -képpen [ESS-FOR] -nként [PERIOD]
ADJ	-ít [TRANS_RESULT] -ul/-odik [INTRANS_RESULT] 	-ság [ABSTRACT]	-bb [COMPAR] -bbik [COMPAR_DESIGN] leg- -bb [SUPERLAT] leg- -bbik [SUPERLAT_DESIGN] -cska [DIMIN] -as [ATTRIB]	-an/-ul [MANNER]
NUM	-- 	-dika [DATE]	-as [ATTRIB] -szori [ITER_ATTRIB] -szoros [MULTIPL_ATTRIB]	-an [AGGREG] -szor [MULTIPL-ITER] -adszor [ORD-ITER]

egyéb:

NUM – NUM (számnévből számnevet képzők):	-ad [FRACT]	-adik [ORD]
POSTP – ADJ (névutóból melléknevet képző):	-i [ATTRIB]	