

# Szótári névelemek felismerése és morfológiai annotálása

TIKK DOMONKOS, KARDKOVÁCS ZSOLT TIVADAR, MAGYAR GÁBOR  
BME Távközlési és Médiainformaticai Tanszék, {tikk,kardkovacs,magyar}@tmit.bme.hu

SZIDAROVSKY FERENC P.  
Szidarovszky Kft., ferenc.szidarovszky@szidarovszky.com

**Kulcsszavak:** internetes keresők, mélyháló, névelemek felismerése, morfológiai annotálás

„A szavak hálójában” projekt keretében készülő internetes keresőszolgáltatásnak egyik célja az, hogy lehetőséget nyújtson természetes nyelvű magyar kérdésekkel internetes adatbázisok tartalmában – az úgynevezett mélyhálóban – való keresésre. Az adatbázisokból ki lehet nyerni azokat az egyedi azonosítókat, amelyek együttese lehetővé teszi, hogy a felhasználói keresések információigénye és a mélyhálós tartalmak között kapcsolatot teremtsünk. Az egyedi azonosítókat névelemnek nevezzük. A természetes nyelvű kérdések feldolgozásának kiemelt fontosságú része a bennük szereplő ismert névelemek felismerése, valamint a kérdésben betöltött szerepük meghatározásához a felismert névelemek morfológiai jegyeinek meghatározása. Cikkünkben bemutatjuk a probléma megoldására javasolt és megvalósított algoritmusunkat, amely számítási igényt tekintve is hatékonyan oldja meg a felvázolt feladatokat.

## 1. Bevezetés

Cikkünk felépítése a következő: először meghatározuk az általunk feldolgozott névelemek körét, és ismeretjük, hogy milyen problémákat kell megoldania a névelem felismerő algoritmusnak, majd részletesen ismeretjük az általunk javasolt névelemfelismerő algoritmust. Ezt követően a működését is bemutatjuk példákon keresztül. Végül az utolsó szakaszban röviden összegezzük a cikk lényeges eredményeit.

A mélyháló jellegzetességei és keresésének jelentősége [1,6], valamint a projekt<sup>1</sup> keretében kidolgozott mélyhálós internettartalmak keresését végző rendszerünk [4,5,7] felől érdeklődő Olvasók számára a megadott irodalmi forrásokat ajánljuk.

## 2. Névelemek és felismerésük problematikája

Az egyedi azonosítókat *szótári*, vagy *ismert névelemnek* nevezzük, amelyeket a *névelemtárban* tárolunk. A szótári jelzőt a *minták alapján felismert névelemektől* (például dátumok, postai és internetes címek stb.) való megkülönböztetésre használjuk, hangsúlyozandó azt, hogy a névelemtárban szereplő névelem bejegyzéseket szótári (kanonikus) alaknak tekintjük.

A szótári névelemek nagy részét a fenti meghatározás miatt a tulajdonnevek teszik ki, azonban alkalmazásunkban a fogalomba beleértjük az olyan rögzített alakú közneveket is, amelyeknek kiemelt szerepe van bizonyos minták alapján felismert névelem típusok (menyisígek, címek stb.) és egyéb, az elemzett kérdés további feldolgozása szempontjából fontos fogalmak azonosítása során. Eszerint névelemnek tekintjük például az alábbi csoportokba tartozó közneveket: a pénz-

nemek jelölései (forint, euró stb.), nemzetiségnevek (magyar, angol, szlovák stb.), közterület típus (út, utca, tér) stb.

A névelemtárnak az adatbázisból történő feltöltése során szemantikai információkat rendelünk az egyes elemekhez, amelyeket az adat adatbázisbeli séma- és attribútum-információiból nyerünk ki. A névelemtárban lehetőség van a kanonikus alak lehetséges szinonimáinak<sup>2</sup> megadására is (például ‘Petőfi Sándor’ bejegyzéshez a ‘Petőfi’ szinonima, vagy a ‘forint’ bejegyzéshez a ‘HUF’ szinonima).

A névelemtár elemei meghatározzák azt az információs teret, amelyben a felhasználó kérdésére választ tudunk adni. Ez azt jelenti, hogy csak azokat a kérdéseket tudjuk megválaszolni a mélyhálós tartalmak segítségével, amelyekben ezen tartalmakból kinyert névelemek szerepelnek. Összességében az alábbi megszorításokat tesszük a felhasználó kérdéseire vonatkozóan, a listában szerepelnek a tartalmi vonatkozású megkötések is:

- csak egyszerű, azaz nem összetett mondatokat fogadunk el;
- csak helyesen írt, és nyelvtanilag helyes mondatokat fogadunk el;
- csak kérdőszóval kezdődő, nem eldöntendő kérdést fogadunk el; a lehetséges kérdőszavakat is korlátozzuk;
- szubjektív (‘Hány éves a kapitány?’), ok-okozati viszonyra irányuló (‘Miért tört ki a II. világháború?’), vagy egyéb nem tényszerű, illetve nem a fenti információs térben található mondatok helyes megválaszolását nem garantáljuk.

A természetes nyelvű kérdések feldolgozásának tehát kiemelt fontosságú része a bennük szereplő ismert névelemek felismerése, valamint a kérdésben betöltött

<sup>1</sup> NKFP-0019/2002 projekt

<sup>2</sup> Nem toldalékolt alakok, csak különböző lehetséges előfordulásai a kanonikus alaknak.

szerepük meghatározásához a felismert névelemek morfológiai jegyeinek meghatározása. Ez a toldalékoló magyar nyelv esetén korántsem egyszerű feladat, mivel a névelemek nem feltétlenül rögzített alakjukban (beleértve a szinonimákat) fordulnak elő, hanem többnyire toldalékolt alakban. A toldalék megváltoztathatja a névelem szótóvét, illetve ha a szótári alak már eleve toldalékolt, akkor ezt is módosíthatja<sup>3</sup>.

További gondot jelenthet az egymásba ágyazott névelemnél a névelem határainak meghatározása<sup>4</sup> [3]. Ha ez utóbbi esetben több értelmezés lehetséges, akkor alternatívákat állítunk elő. A morfológiai jegyek meghatározásánál a nem alanyesetű kanonikus alakok és a nem magyar (azaz morfológiai elemző által fel nem ismert) névelemek *speciális* esetei kívánnak külön megfontolást.

Cikkünkben bemutatjuk a probléma megoldására javasolt és megvalósított algoritmusunkat, amely azon kívül, hogy a fenti feladatokat megoldja, mindezt a számítási igényt tekintve hatékonyan valósítja meg. Az ismertetett módszer a HunMorph [2] szabad forráskódú statisztikai alapú morfológiai elemzőt használja, ennek megfelelően a példákban található morfológiai elemző eredmények is a HunMorph kódolása szerint vannak megadva.

Fontosnak tartjuk kiemelni, hogy a módszer *nem felügyelt tanuláson alapul*, mivel célja nem ismeretlen névelemek felismerése, hanem az ismertek pontos azonosítása.

### 3. Szótári névelemek felismerése

A szótári névelem (ezentúl itt csak *névelem*) felismerőnek két fő célja van:

- keresés: a mondatban szereplő névelemek megtalálása;
- annotálás (vagy címkézés): névelemek morfológiai jegyeinek meghatározása.

A keresés és annotálás folyamata általában összekapcsolódik, így önmagukban nem hajthatók végre.

Mivel egy névelem több szóból is állhat, a kérdőmondat tetszőleges szegmense (szavak rögzített sorrendű sorozata) lehet névelem. Egy  $n$  szavas kérdőmondat szegmenseinek száma  $n(n+1)/2$ . Egy átlagos kérdőmondat 7-10 szóból áll, míg a névelemtár mérete  $10^6$  nagyságrendű is lehet. Így sokkal hatékonyabb a mondatsegmensekből kiindulva keresni, mint a névelemtárból kiindulva. Egy kifejezés keresése a névelemtárban gyorsítható a névelemtár elemeinek hash-elérésével. A mondatsegmensek összevetése a névelemtárral a szegmensek hossza szerint csökkenő sorrendben történik.

A névelem felismerés egy másik problémája, hogy egy névelem tartalmazhat egy másikat (például: 'a The New York Times egy napilap'). Míg a Blitz NL feldolgozó [3] a felismert névelemek közül csak egyet választ ki

konfidencia értékek alapján, mi fel kívánjuk ismerni az összes névelemet, különböző mondat alternatívákat létrehozva. Ebből kifolyólag az összevetés a keresés eredményétől függetlenül tovább folytatódik a rövidebb szegmensekkel.

A szegmensek összevetése az alábbi sorrendben történik:

1. A teljes mondattal kezdjük:  $[1, \dots, n]$ , és vesszük az első szóval kezdődő egyre rövidebb szegmenseket:  $[1, \dots, j]$ , ahol  $j = n-1, \dots, 1$ .
2. Vesszük a második szóval kezdődő egyre rövidebb szegmenseket:  $[2, \dots, j]$ , ahol  $j = n, \dots, 2$ .
3. Általánosan, az összes szegmenst megvizsgáljuk a kezdőszó mondatbeli pozíciója szerint növekvő, majd azon belül a szegmens hossza szerint csökkenő sorrendben:  $[i, \dots, j]$ , ahol  $i = 3, \dots, n$ ,  $j = n, \dots, i$ .

**Megjegyzés:** Nyilván nem mindegyik mondatsegmens lehet valóban névelem. Ha figyelembe vesszük, hogy a mondat első szavának a megszorítások miatt feltétlenül kérdőszónak kell lennie, akkor kezdetünk a 2. lépéssel ( $[2, \dots, n]$  szegmenstől), a vizsgálandó részeket  $n(n-1)/2$ -re csökkentve.

#### 3.1. A névelem felismerő algoritmus

A továbbiakban a névelem felismerést egy konkrét mondatsegmens (ezentúl *jelölt*) kapcsolatban ismeretjük. A magyar nyelvben a szavak töve változhat toldalékolásnál. Az esetek nagy részében a szótónek csak az utolsó két betűje változhat (tűz→tűzet; álom→álmot) például rövidülés, hangkivetés miatt. Hasonlóan, egy toldalék megváltozhat egy következő toldaléktól (ez csak akkor fordul elő, ha a névelem magában is toldalékolt, és azt a mondatban tovább toldalékoljuk, lásd 3. lábjegyzet), ekkor azonban csak az utolsó betű változhat. Mindezeket figyelembe kell vennünk a névelem felismerési keresés fázisában.

A névelemek jelentős része nem magyar nyelvű, így a morfológiai elemző nem képes azokat elemezni. Ennek ellenére a névelem felismerő ezen névelemeket is el kell lássa morfológiai jegyekkel. Erre a feladatra úgynevezett *helyettesítő szavakat* használunk (a helyettesítő szavak előállításáról a 3.2. szakaszban részletesebben is értekezünk).

A helyettesítő szónak a névelemek toldalékainak meghatározásánál van szerepe. Feltételezzük, hogy minden névelemhez rendelkezünk egy helyettesítő szóval, mely morfológiailag elemezhető és pontosan ugyanúgy ragozódik (kiejtés szerint azonos hangrendű, főnév), mint a névelem utolsó szava. A helyettesítő szónak mindig főnévnek kell lennie, mivel az ismert névelemek előfordulásai egyedi entitásokat jelölnek, tehát a mondatban főnévi szerepben állnak és eszerint kapnak toldalékokat. Kivételt képeznek a 2. szakaszban ismertetett egyéb névelemtípusok egyes esetei, de ezek a morfológiai elemző által ismert magyar szavak, ahol tehát a morfológiai jegyek megállapítására nincs szükség helyettesítő szóra.

<sup>3</sup> lásd: *Vissza a jövőbe és Hol adják a Vissza a jövőbét?*

<sup>4</sup> *New York Times* sport rovata tartalmazza a *New York, York, Times*, és *New York Times*-t.

Az alábbi jelöléseket használjuk:

- $last(x)$  jelöli az  $x$  kifejezés utolsó szavát,
- $length(x)$  jelöli az  $x$  szó betűinek számát,
- $trunc(x,i)$  jelöli az  $x$  szót az utolsó  $i$  betűje nélkül,
- $lchar(x)$  jelöli az  $x$  szó utolsó betűjét.

Továbbá jelölje  $C$  a jelöltet,  $S$  a helyettesítő szót és  $E$  a névelemet.

Az algoritmus folyamatábráját az 1. ábra szemlélteti.

1. Ha  $last(E)$  toldalékolható, alanyesetű, magyar szó (azaz a morfológiai elemző felismeri).

**1.1. keresés:**

- 1.1.a ha  $length(last(E)) \geq 3$ , ellenőrizzük, hogy  $C$   $trunc(E,2)$ -vel kezdődik-e.
- 1.1.b ha  $length(last(E)) < 3$ , ellenőrizzük, hogy  $C$   $E$ -vel kezdődik-e.

**1.2. szótó ellenőrzés:**

Ha 1.1.a igaz, azaz  $C$   $trunc(E,2)$ -vel kezdődik, akkor meg kell határozni, hogy  $last(C)$  és  $last(E)$  szótöve megegyezik-e. Erre azért van szükség, mert a betűelhagyás miatt a csonkolt szó több értelmes szónak is a prefixe lehet.

Ez a lépés kihagyható, ha 1.1.b igaz.

**1.3. annotáció:**

Ha 1.2.-ben a szótövek megegyeznek, akkor  $C$  az  $E$  névelem, melynek morfológiai jegyei a  $last(C)$  jegyei. Ha  $E$  és  $C$  egyaránt rendelkezik nem záró morfémával, azt kihagyjuk az annotációból (lásd 4. példa).

2. Ha  $last(E)$  nem felel meg az 1. feltételeinek, azaz a morfológiai elemző nem ismeri fel, vagy nem toldalékolható, vagy nem alanyesetű.

**2.1. keresés:**

- 2.1.a Ha  $lchar(last(E))=a$  vagy  $=e$ , ellenőrizzük, hogy  $C$   $trunc(E,1)$ -vel kezdődik-e.
- 2.1.b Ha  $lchar(last(E)) \neq a$  és  $\neq e$ , ellenőrizzük, hogy  $C$   $E$ -vel kezdődik-e.

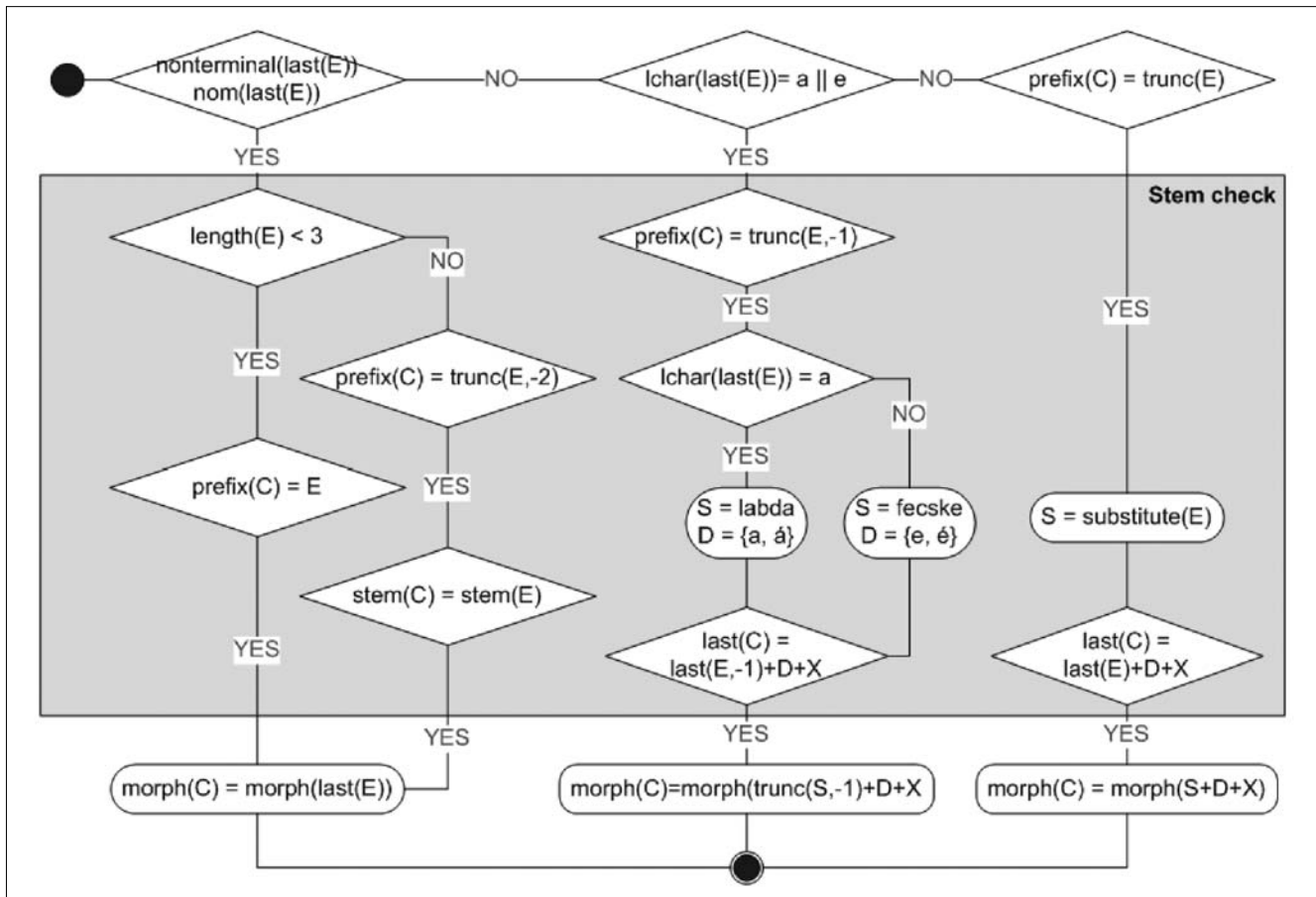
**2.2. helyettesítő szó megállapítása:**

- 2.2.a Ha 2.1.a igaz és  $lchar(last(E))=a$ , akkor  $S=labda$ , ha  $lchar(last(E))=e$ , akkor  $S=fecske$ .
- 2.2.b Ha 2.1.b igaz, akkor vesszük a névelem-tárban  $E$ -hez megadott  $S$ -t.

**2.3. annotáció:**

- 2.3.a  $C$  utolsó szavának alakja a következő:  $[trunc(last(E),1)\{a,e\}marad]$ , ahol *marad* a ( $C$ ) végén lévő maradék betűkből áll (ha vannak). A következő szövegeket elemeztetjük a morfológiai elemzővel:  $[trunc(last(S),1)\{a\}marad]$ , illetve  $[trunc(last(S),1)\{e\}marad]$ , ha  $lchar(E) = a$ , illetve  $lchar(E) = e$ , azaz a szóvégi magánhangzót hosszúra cseréljük. Csak az egyik szöveg lesz helyes szó, és ismeri fel a morfológiai elemző. A  $C$  morfológiai jegyei a helyes szó jegyei lesznek.

1. ábra Az algoritmus folyamatábrája



2.3.b C utolsó szavának alakja a következő: [last(E) *marad*]. A következő szöveget elemeztetjük a morfológiai elemzővel: [S *marad*]. A C morfológiai jegyei az [S *marad*] szó jegyei lesznek.

**1. megjegyzés:** Látható, hogy az első esetben a keresés bonyolultabb, mert a toldalékolható szavak esetén a helyes szót azonosítása nehezebb. A második esetben viszont az annotálás a bonyolultabb, mert a toldalékok meghatározása csak egy megfelelő helyettesítő szóval lehetséges.

**2. megjegyzés:** A névelemek keresett alakja a névelemtár feltöltésekor számítható és tárolható, így jelentős időt nyerünk a keresésnél.

**3. megjegyzés:** A 2.3.-nál ha  $\text{length}(\textit{marad})=0$ , akkor kihagyható a morfológiai elemző használata, mert ez azt jelenti, hogy a névelemen nincsenek toldalékok és az egy alanyesetű főnévnek tekinthető.

**4. megjegyzés:** A 2.2.b-ben használt, a névelemhez rendelt helyettesítő szó meghatározásánál egy félheurisztikus algoritmust használunk. A helyettesítő szavakat már a névelemtár feltöltésekor offline, a névelem utolsó mássalhangzója és az utolsó szavának magánhangzói alapján határozzuk meg. Míg ez (kiejtett) magánhangzóra végződő szavak esetén triviális, mássalhangzóra végződő szavak esetén több körülményt igényel. Ez az eljárás például az idegen szavak kiejtés követő toldalékolása miatt nem 100%-osan tökéletes, de az esetek túlnyomó többségében (több mint 98%-ban) jó helyettesítő szavakat eredményez.

### 3.2. Helyettesítő szavak automatikus előállítás

A helyettesítő szavakat a névelem utolsó szava alapján állítjuk elő. A helyettesítő szónak mindig főnévnek kell lennie, mivel a szótári névelem egyedi entitásokat jelölnek, tehát a mondatban főnévi szerepben állnak és eszerint kapnak toldalékokat.

1. Egyszerű esetben, amennyiben a névelem utolsó szava a morfológiai elemző által felismert alanyesetű főnév vagy melléknév<sup>5</sup>, akkor a helyettesítő szó azonos a névelem utolsó szavával.
2. Ha a névelem utolsó szava a morfológiai elemző által felismert szó, de más szófajú, illetve nem alanyesetű főnév, vagy melléknév, akkor ezek más para-

digma szerint kapnak a mondatban toldalékokat, mint ha nem névelem részét képeznék, hiszen ekkor a például nem alanyesetű főnevek újabb zárómorfémát kaphatnak. (Lásd a 3. lábjegyzetet és 3. példát). Ebben az esetben a következőképpen járunk el:

**2.1.** Meghatározzuk a vizsgált szó hangrendjét a benne szereplő magánhangzók számának és pozíciójának alapján.

Van néhány háromalakú rag (például hoz/hez/höz) is a magyarban, ekkor a magas hangrendű síkban pedig labiális és illabiális változatok vannak. A hasonló morfémák helyes illesztéséhez – például *instrumentalis* (-vAl esetrag) – a helyettesítő szó a vizsgált szó utolsó betűjé-  
től is függ.

**2.2.** Egy táblázatból kikeressük a hangrendnek és az utolsó betűnek megfelelő helyettesítő szót. Az alkalmazott táblázat egy részletét, illetve a példákat lásd az 1. táblázatban.

**3.** Abban az esetben, ha a szót nem ismeri fel a morfológiai elemző, akkor általában idegen nyelvű a névelem utolsó szava<sup>6</sup>, ami elég gyakori eset a névelemként előforduló idegen nyelvű tulajdonnevek nagy száma miatt (személynevek, földrajzi nevek stb.). A magyar ebben az esetben többnyire a kiejtés alapján közvetlenül, azaz nem kötőjellel kapcsolja a toldalékokat az idegen szóhoz [4]. Ez alól kivételt jelent, ha a tulajdonnév végén néma betű van, vagy ha a kiejtett hangot bonyolult, magyartól idegen betűkombináció jelöli (például 'Diderot-nak', 'Renault-t', 'Rousseau-val'). Ez utóbbi esetben a toldalék meghatározása és a morfológiai annotálás triviális, tehát csak a kötőjel nélküli esettel foglalkozunk.

**3.1.** Mivel a hangrend, illetve az utolsó betű kiejtése idegen szavaknál nem egyértelmű, ezért esetenként több kísérletet teszünk a helyettesítő szó meghatározására. A hangrendet a magyar szabályok szerint határozzuk meg. A helyettesítő szót a 2. esetben is használt táblázattal adjuk meg.

**3.2.** Ha helytelen a kiválasztott helyettesítő szó, akkor annak a toldalékkal bővített alakját a morfológiai elemző nem fogja felismerni, ekkor új helyettesítő szót keresünk.

1. táblázat A helyettesítő szavak meghatározása hangrend és utolsó betű alapján (részlet)

szóvég	mély hangrend	magas illabiális	magas labiális	példák
<i>magánhangzó</i>	labda	fecske	tömlő	Híd a <i>túlvilágra</i> ; Vissza a <i>jövőbe</i>
<i>b</i>	comb	seb	göb	Coulomb;
<i>c</i>	konc	férc	gönc	
<i>cs</i>	gáncs	tincs	göcs	Antics; Gerevich; Göröcs
<i>d</i>	kaland	pléd	körönd	Lund; Sutherland;

<sup>5</sup> Melléknév ragozási szempontból azonosan viselkedik a főnévvel.

<sup>6</sup> A Hunmorph számos gyakran használt idegen tulajdonnevet ismer, ezekre természetesen az előző két eset valamelyikét kell alkalmazni.

**3.3.** Először a hangrendi módosulato-  
kat vizsgáljuk, tehát például 'Beck-  
hamtól' a mély hangrendű 'karám'  
helyett, a magas hangrendű 'szem'-  
et alkalmazzuk.

**3.4.** Az utolsó szó kiejtés szerint tolda-  
lékolása esetén egy segéd szabályt  
használnak, mely a kiejtési válto-  
zatokat adja meg. Ennek alapján  
a táblázatban a kiejtett hang sze-  
rinti sorokat vizsgáljuk meg. Például  
ch-végződés esetén az alapértelme-  
zett a *h* (Bachhal), de lehetséges  
még a *cs* is (Gerevichcse), illetve a  
*k* is (Murdochkal) stb.

**1. megjegyzés:** A névelem felismerő al-  
goritmusban azért használunk helyettesítő  
szavakat, ahelyett hogy az illesztés után  
megmaradó karakterláncot próbálnánk meg  
toldalékokként felismerni, mert ez utóbbiak  
rendkívül sokfélék lennének, azaz szinte  
egy (valamelyest korlátozott) morfológiai  
elemzőt kellene írni a megvalósításához.

**2. megjegyzés:** Előfordulhat időnként  
az, hogy a 2. csoportbeli szavaknak hely-  
telenül határozzuk meg a hangrendjét. Ezt  
ugyanúgy detektáljuk, és oldjuk meg, mint  
a 3.3. esetben, azaz a más hangrendű he-  
lyettesítő szót alkalmazunk helyette.

#### 4. Példák

A továbbiakban néhány példán keresztül  
bemutatjuk az algoritmus működését.

##### 1. példa

Lásd a 2. ábrát.

*Milyen költők vannak Arany Jánostól József Attiláig?*  
*E=József Attila,*

last(*E*)-t felismeri a morfológiai elemző mint

**Attila[noun\_prs]+[NOM]**

így ez az 1-es eset. A keresés *József Atti* kifejezés-  
sel végezzük, ami alapján a *C = József Attiláig* szeg-  
menst találjuk (mivel ezekben a példákban a *C* válas-  
tása triviális, a következőkben külön nem térünk ki rá).

A last(*C*) morfológiai elemzése

**Attila[noun\_prs]+[TERM]**

Így az *E* névelemet felismertük *C*-ben és a morfoló-  
giai jegyei [TERM].

##### 2. példa

Lásd a 2. ábrát.

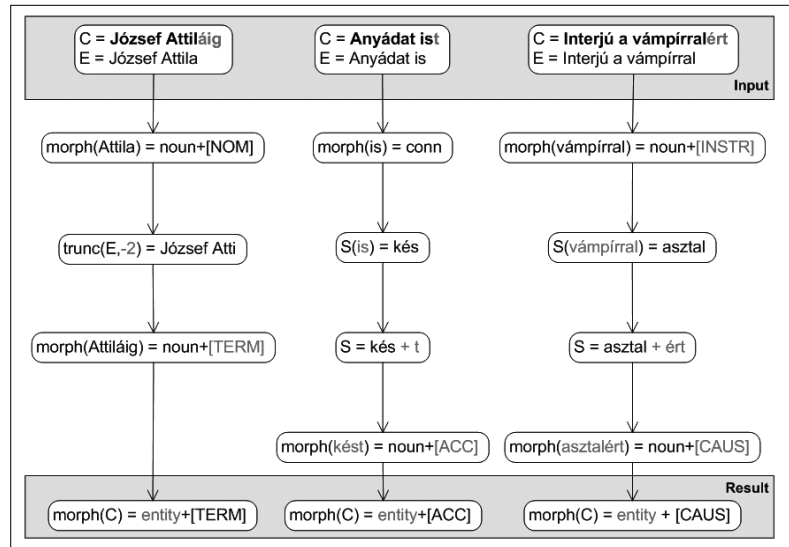
*Ki rendezte az Anyádat ist?*

*E=Anyádat is,* ez a 2 (b) eset, mert az *is* kötőszó,  
mely nem toldalékolható. Legyen *S* a *kés*, így a morfo-  
lógiai elemzővel a *kést* szöveget elemeztetjük.

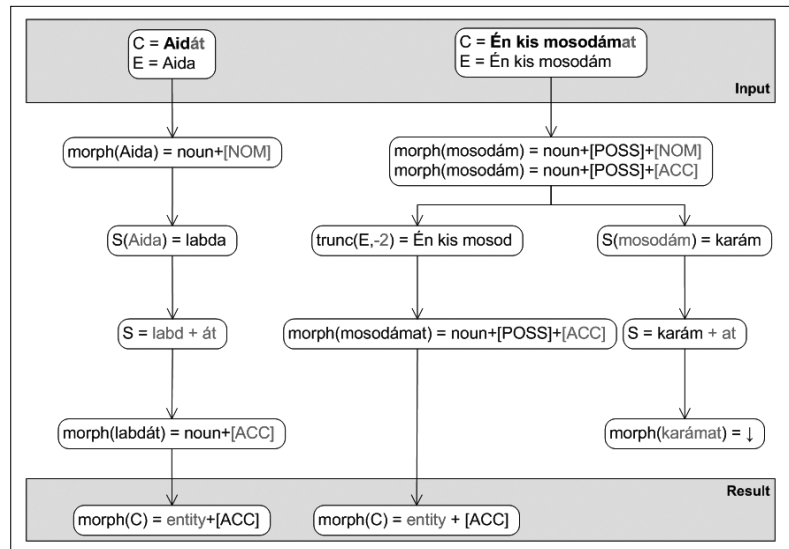
Az eredmény

**kés[noun]+[ACC]**

így a felismert névelem: *Anyádat is<sub>névelem</sub>+ [ACC]*.



2. ábra Illusztráció az 1-3. példához



3. ábra Illusztráció a 4-5. példához

##### 3. példa

Lásd a 2. ábrát.

*Mennyit kell fizetnem az Interjú a vámpírralért?*

*E=Interjú a vámpírral,* ez is a 2 eset, mert last(*E*) már  
toldalékol:

**vámpír[noun]+[INSTR]**

Legyen *S* az *asztal*, így a morfológiai elemzővel az  
*asztalért* szöveget elemeztetjük.

Az eredmény

**asztal[noun]+[CAUS/FIN]**

így a felismert névelem:

*Interjú a vámpírral<sub>névelem</sub>+ [CAUS/FIN]*.

##### 4. példa

Lásd a 3. ábrát.

*Ki rendezte Az én kis mosodámat?*

*E=Az én kis mosodám.* A névelem utolsó szava bir-  
tokos toldalékú, amit a névelem egészére mint entitás-  
ra vonatkozóan tárgyrag követ.

Ebből következően a névelemet csak a tárgyraggal kell felcímkézni. Az utolsó szó morfológiai elemzése a névelem az algoritmus mindkét fő ágát aktiválja, hiszen

**mosoda[noun]+[POSS\_SG\_1]+[ACC]**

**mosoda[noun]+[POSS\_SG\_1]+[NOM]**

Az első sor a 2-es esetet aktiválja. Legyen S a *karám*, így a morfológiai elemzővel a *karámat* szöveget elemeztetjük. Mivel ezt a szót a morfológiai elemző nem ismeri fel, ez az ág nem talál névelemet.

A második sor az 1-es esetet aktiválja. A  $last(E) = mosodám$  és  $last(C) = mosodámat$  szótöve egyezik, és C E-vel kezdődik.

Végül a morfológiai jegyeket a  $last(C)$  és  $last(E)$  morfológiai jegyeinek különbözetéből kapjuk:

*Az én kis mosodám*<sub>névelem</sub>+**[ACC]**.

### 5. példa

Lásd a 3. ábrát.

*Hol játsszák az Aidát?*

$E=Aida$ . Ez a 2 (a) eset, mert  $last(E)$ -t nem ismeri fel a morfológiai elemző. Legyen S a *labda*, így a morfológiai elemzővel a *labdát* szöveget elemeztetjük, melynek eredménye

**labda[noun]+[ACC]**

Így a névelem felismerés eredménye:

*Aida*<sub>névelem</sub>+**[ACC]**.

## 5. Összefoglalás

A fentiekben ismertettük annak a feladatnak a jelentőségét és nehézségeit, mely egy természetes magyar nyelvű kérdőmondatban a szótári névelemek összes előfordulásának megkeresése és morfológiai jegyekkel való ellátása.

Ismertettünk egy algoritmust, mely megoldás erre a feladatra, és hatékonyan végrehajtható.

## Köszönetnyilvánítás

A cikk a Nemzeti Kutatási és Fejlesztési Pályázatok NKFP-0019/2002 jelű projektjének és a Mobil Innovációs Központ támogatásával készült.

## Irodalom

- [1] M. K. Bergman:  
The deep web: surfacing hidden value.  
Journal of Electronic Publishing, 7/2001.  
[www.press.umich.edu/jep/07-01/bergman.html](http://www.press.umich.edu/jep/07-01/bergman.html)
- [2] Hunmorph,  
<http://mokk.bme.hu/resources/hunmorph/>
- [3] B. Katz, D. Yuret, J. Lin,  
S. Felshin, R. Schulman, A. Ilik:  
Blitz: A preprocessor for detecting context-independent linguistic structures. In Proc. of the 5th Pacific Rim Conference on Artificial Intelligence (PRICAI '98), Singapore, 1998.
- [4] Magyar Tudományos Akadémia:  
A magyar helyesírás szabályai (11 kiadás),  
Az idegen közzavak és tulajdonnevek írása –  
216-217.; pp.87–88., Akadémiai Kiadó, 1984.
- [5] D. Tikk, Zs. T. Kardkovács, Z. Andriská, G. Magyar,  
A. Babarczy, I. Szakadát:  
Natural language question processing for  
hungarian deep web searcher. In Proc. of IEEE  
Int. Conf. on Computational Cybernetics (ICCC04),  
pp.303–309, Wien, Austria, 2004.
- [6] D. Tikk, Zs. T. Kardkovács, G. Magyar:  
A szavak hálójában:  
szabadszavas mélyháló-kereső program,  
Híradástechnika, 60(5): pp.2–8, 2005.
- [7] H. Winkler:  
Suchmaschinen – metamedien im internet?  
In B. Becker, M. Paetau, editors,  
Virtualisierung des Sozialen, Frankfurt/NY  
pp.185–202., 1997; német nyelven,  
angol fordítás:  
[www.uni-paderborn.de/~timwinkler/suchm\\_e.html](http://www.uni-paderborn.de/~timwinkler/suchm_e.html)

