

Különírás-egybeírás – automatikusan

Ludányi Zsófia^{1,2}, Miháltz Márton², Hussami Péter³

¹ ELTE BTK Nyelvtudományi Doktori Iskola

² MTA Nyelvtudományi Intézet, Nyelvtechnológiai és Alkalmazott Nyelvészeti Osztály

³ Alkalmazott Logikai Laboratórium

{ludanyi.zsofia, mihaltz.marton}@nytud.mta.hu, hussami@all.hu

Kivonat: Jelen tanulmány a helyesiras.hu automatikus helyesírási tanácsadó rendszer külön- vagy egybeírással foglalkozó webes alkalmazását mutatja be. A modul attribútum-érték struktúrák környezetfüggetlen nyelvtani elemzésen alapszik. Az elemzés morfológiai és szemantikai tulajdonságokra támaszkodik. A rendszer általános működésének, illetve a nyelvtani elemző felépítésének és működésének ismertetése után a rendszer egyik fontos alappilléret képező formális nyelvtan részletes bemutatása következik. Végezetül néhány bonyolultabb helyesírási probléma nyelvtechnológiai megoldását ismertetjük példákkal illusztrálva (mozgószabályok, szótagszámlálási szabály stb.).

1 Bevezetés

Az MTA Nyelvtudományi Intézete 2009 óta dolgozik egy olyan szakértői rendszeren, amely nyelvtechnológiai eszközök segítségével kísérel meg a felhasználók helyesírási kérdéseire automatikus választ adni (Miháltz et al 2012, Pintér et al 2009). A helyesiras.hu névre keresztelt készülő tanácsadó portál hét különböző helyesírási területen próbál interaktív segítséget nyújtani: külön- és egybeírás, helyesírás-ajánló, elválasztás, tulajdonnevek írása, számnevek helyesírása, keltezés, betűrendbe sorolás. A felhasználónak a nyitóoldalon felkínált menüből kell kiválasztania, hogy milyen típusú helyesírási kérdésre szeretne választ kapni (1. ábra).

Jelen tanulmány célja a helyesiras.hu projekt külön- és egybeírással foglalkozó moduljának bemutatása.

Üdvözljük portálunkon

Oldalunkon többféle – nyelvtechnológiai eszközökkel működő – tanácsadó alkalmazás található, valamint sok más információ (A magyar helyesírás szabályai (az Akadémiai Kiadó engedélyével), Gyakran feltett kérdések) a helyesírási kérdésekben való tájékozódáshoz.

ab

Külön vagy egybe?

Ellenőrizendő szó
kakukkosóra

Helyes alak
kakukkos óra

Kipróbálok

sz

Helyes-e így?

Ellenőrizendő szó
hejesírás

Helyes alak
helyesírás

Kipróbálok

AB

Névkereső

Keresett kifejezés
Széch...

Találatok
Széchenyi stb.

Kipróbálok

ab

Elválasztás

Ellenőrizendő szó
elválasztás

Elválasztási helyek
el-vá-lasz-tás

Kipróbálok

5öt

Számok

Számjegyekkel
2010

Betűkkel
kétezer-tíz

Kipróbálok

2013

Dátumok

Dátum
2012-08-30

A következő módokon látható
2012. aug. 30. stb.

Kipróbálok

abc

ABC-be rendezés

Adjon meg egy listát
tej, tojás, kenyér

A rendezés eredménye
kenyér, tej, tojás

Kipróbálok

1. ábra. A helyesiras.hu nyitóoldala

1.1 A külön- és egybeírás problémája

A magyar helyesírás egyik legproblematicusabb kérdésköre a különírás és az egybeírás. A szabályok megfelelő alkalmazása némi grammatikai alaptudást igényel, mivel a helyesírás rendszerszerűségét a magyar nyelvtan szabályai alakították ki. Különbséget kell tudni tehát tenni a szó szerkezetek, illetve szóösszetételek között. Ez sok esetben problémás, mivel a két nyelvtani kategória között sokszor nem éles a határ, gyakran fordulnak elő nem egyértelmű, többféleképpen megítható esetek (Laczkó–Mártonfi 2004).

Az összetétellé válás oka sok esetben például a jelentésváltozás. Olyan esetekben, amikor az adott kifejezés konkrét és elvont jelentésben is szerepel (s a megfelelő írásmódot éppen ez dönti el), emberi beavatkozás nélkül nem adható egyértelmű megoldási javaslat. A különírás-egybeírás sok egyéb területére is jellemző ez. Kimondható, hogy jelenleg nem lehetséges olyan helyesírás-ellenőrző, helyesírási tanácsadó alkalmazás kifejlesztése, amely teljesen önállóan, az ember anyanyelvi kompetenciáját segítségül hívó beavatkozás nélkül képes a külön- és egybeírás minden területét hatékonyan kezelni (Pintér et al. 2009).

1.2 A helyesiras.hu újszerűsége

A létező online helyesírási tanácsadók egyszerű szójegyzéken alapulnak, és csak akkor adnak kielégítő eredményt, ha a beírt szó eleve helyesen van leírva, és megtalálható a rendszer mögött álló szótárban (Pintér et al. 2009). Léteznek olyan szójegyzék alapú online tanácsadók, amelyek bizonyos hiányos bemeneteket is elfogadnak (pl. a www.magyarhelyesiras.hu online szótár elfogadja a *j-ly* cserével keletkezett hibás bemeneteket, ékezet nélküli alakokat), de a pusztán szótár alapú megközelítés nem elég hatékony.

A helyesiras.hu külön- és egybeíró modulja mögött ezzel szemben egy formális nyelvtan áll, amelyet felhasználva a kifejlesztett nyelvtani elemző létrehozza a megadott bemenetből generálható lehetséges jó megoldásokat.

2 A különírás-egybeírás webalkalmazás felépítése

2.1 Általános felépítés

A rendszer felhasználókkal történő interakcióját szemlélteti a 2. ábra, a 3. ábra pedig moduljainak kapcsolatát.

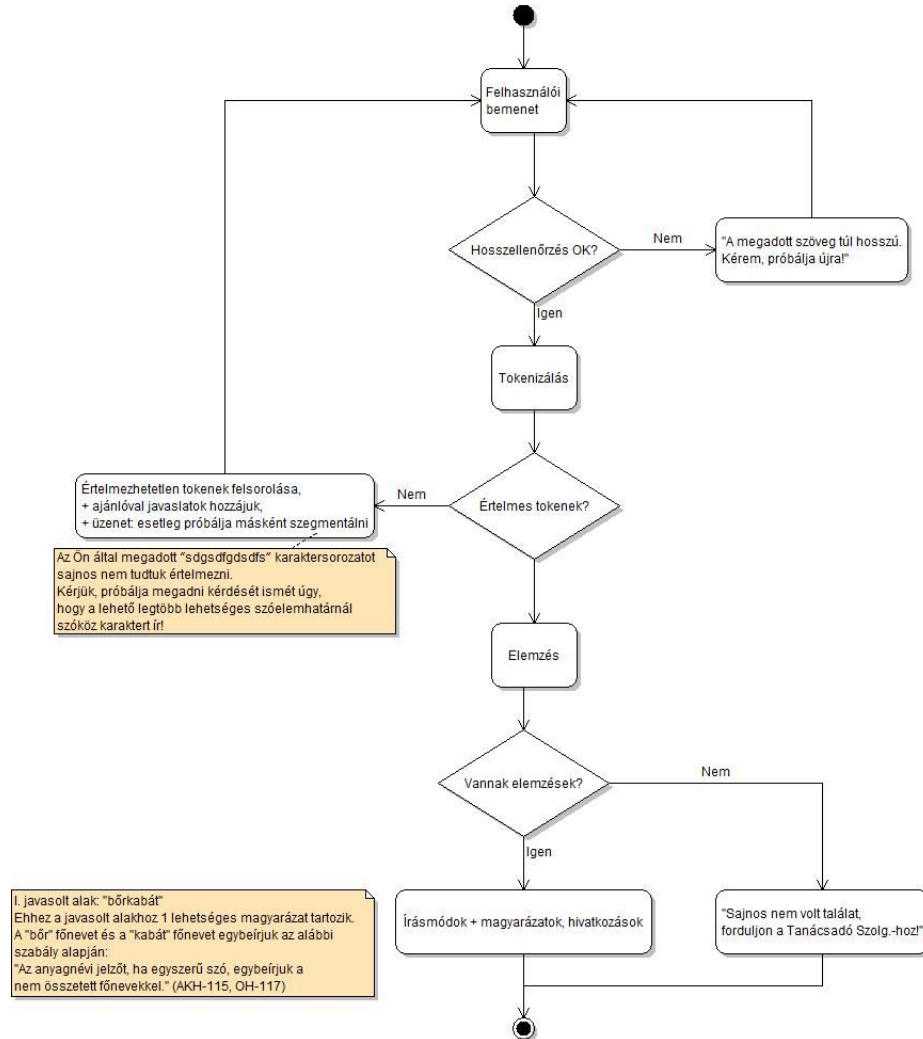
A felhasználói bemenetet néhány egyszerűbb ellenőrzésnek vetjük alá. A beírt szöveg hossza legfeljebb 70 karakter lehet; amennyiben ennél hosszabb bemenetet kapunk, a rendszer hibaüzenettel válaszol. Ha a karakterhossz megfelelő, következik a tokenizálás: a rendszer megkísérli tokenekre bontani a bemenetet.

A tokenizáló modul eltávolítja a felhasználó által megadott kötőjeleket (ha vannak), azokat szóközre cseréli, és az így kapott elemeket próbálja atomi szintű tokenekre bontani, illetve a HuMor morfológiai elemzővel (Novák–Pintér 2009) értelmezni (szófaji, morfológiai információkkal ellátni). Ha a tokenizálás sikertelen (nem sikerült minden tokent a morfológiai elemzővel azonosítani, illetve továbbbontani), a rendszer megkéri a felhasználót, hogy újból adja meg a kívánt bemenetet, az összes lehetséges helyen szóközzel elválasztva. Sikeresen értelmezés esetén következik az elemzés, ellenkező esetben újabb hibaüzenetet kapunk, illetve az oldal átirányítja a felhasználót a *Helyes-e így?* és a *Névkereső* modulokhoz.

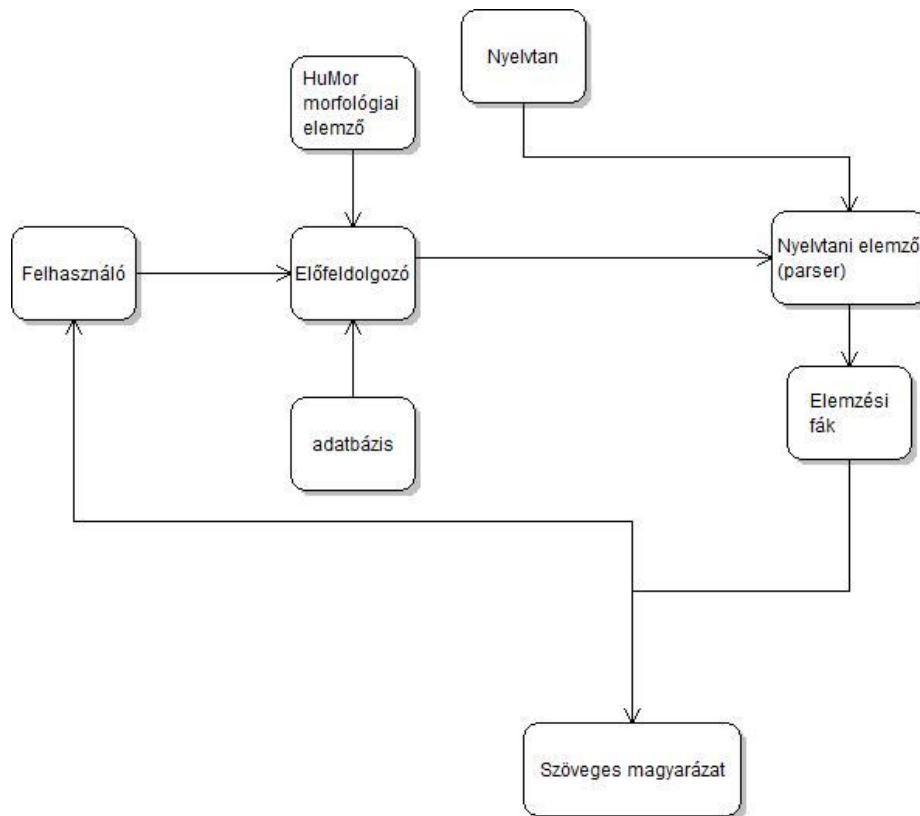
Amennyiben a tokenizálás sikeres volt, és megtörtént a bemenet morfológiai elemzése (szófaji, morfológiai, szótagszámra és összetételi tagok számára vonatkozó információkkal történő ellátása), az adatbázisban eltárolt szemantikai kategóriákat is hozzárendeljük a tokenekhez (ha vannak). Ilyen szemantikai kategóriák pl. a színnevek, foglalkozások és rangok, számnevek, földrajzi jellegű jelzők és köznevek, közterületek nevei, keresztnévek, népek és nyelvek nevei, rövidítések, közszei betűszók, önálló szóként nem használatos előtagok, a helyesírási szabályzatban az egyes szabályokban hivatkozott további kategóriák és különösen az egyes kivételek listája, melyeknek száma jelenleg 2100 körülire tehető.

A morfológiai, szemantikai tulajdonságokkal felruházott tokenek képezik az elemző modul bemenetét. Az elemző a mögöttes nyelvtanban található, speciális formális nyelven megfogalmazott helyesírási szabályokat próbálja alkalmazni a bemeneti tokenekre. Sikeres elemzés esetén megkapjuk a lehetséges megoldásokat a hozzájuk tartozó magyarázatokkal és az érvényben lévő akadémiai helyesírási szabályzat (AkH. (esetleg OH.)-beli hivatkozásokkal együtt. Ha az elemzés sikertelen,

azaz a megadott helyesírási szabályok egyike sem alkalmazható a bemenetre, a rendszer felajánlja a humán szakértői segítséget: a Nyelvtudományi Intézet Közönségszolgálatának telefonos vagy e-mailes igénybevételét (Miháltz et al. 2012).



2. ábra. A rendszer általános felépítése



3. ábra. Az elemző felépítése

2.2 Az elemző felépítése és működése

Az elemző modul bemenetét az előző részben említett szegmentált, morfológiai és szemantikai jegyekkel ellátott tokenek képezik.

A szabályok illesztésére használt algoritmus lényegében egy bottom-up modell szerint működő elemző.

A célja az, hogy előállítsa az összes olyan szintaktikai fát, amely a bizonytalanságokat is tartalmazó bemeneti adatokból a rendelkezésre álló szabályok sorozatos alkalmazásával létrejöhet. A bemeneti tokenek fölé exponenciálisan sok fát lehet építeni, az algoritmus gyakorlati megvalósítása ezt optimalizálja. Az optimalizáláshoz felhasználunk egy ún. „kiértékelési lánc” konstrukciót. Ez ilyen lánc nem más, mint egy hipotézis a szabályalkalmazási sorozatra, nevezetesen, hogy melyik szabályt hol alkalmaztuk. Egy n argumentumú szabály alkalmazása helyettesíti a lánc n csomópontját egy eredménycsomóponttal. Egy láncon általánosságban több szabály is alkalmazható, így a lánc utódjai többen is lehetnek. Ha egy láncon nem alkalmazható több szabály, de a kilépési kritériumnak még nem felel meg, akkor az a lánc lekerül a jelöltek listájáról.

Kezdetben a lánc maga a tokenlánc. Kilépési kritériumnak azt választottuk, hogy a lánc egyelemű legyen, azaz egy teljes fát reprezentáljon. A láncok evolúciójának alapművelete a helyi szabályalkalmazás: az algoritmus a lánc összes csomópontján begyűjti, milyen szabályokat tudna ott alkalmazni. Ahhoz, hogy egy lánc k -edik elemén alkalmazhassuk az n argumentumú X szabályt, arra van szükség, hogy a lánc k -edik, $k+1$ -edik, ..., $k+n-1$ -ik eleme megfeleljen X első, második, ... n -edik elemének. A megfelelés szükséges és elégséges feltétele, hogy az adott csúcsok morfológiai címkéi egybeessenek, és a bemenet attribútumai megfeleljenek a szabályban megkövetelt feltételeknek.

2.3 A nyelvtan

A modul alapját képező környezetfüggetlen, jegystruktúrárs nyelvtan formális leírása független a modul programkódjától, így könnyen karbantartható, fejleszthető.

A nyelvtan jelenleg mintegy 230 darab szabályt tartalmaz. Egy szabály felépítését az alábbi példán keresztül mutatjuk be:

```
id: M_EK_ANYAGNEV_1_2_1
rule: N(sem="Material", ncomparts>=2, type!="Qualificative") +
N(ncomparts=1) == N(sep=' ', hasnesep="1")
comment: Ha az anyagnévi jelzős kapcsolatnak valamelyik vagy mindkét
tagja össze tett szó, az anyagnevet különírjuk jelzett szavától.
refs: AKH-115, OH-117
ex: műbőr + kabát = műbőr kabát, nyersselyem + ing = nyersselyem ing
kill: M_EK_JELOLETLEN_BIRTOKOS
```

A kettőspontra végződő mezők jelentése a következő:

- Az **id** mező a szabály egyedi azonosítóját tartalmazza.
- A **rule** mezőben található maga az újraíró szabály (ennek kifejtését lásd alább).
- A **comment** mező tartalmazza a szabálynak a szöveggel történő megfogalmazását.
- A **refs** mezőben találjuk az AkH. megfelelő szabálypontjaira és/vagy az Osiris-helyesírás releváns témaköreire történő hivatkozásokat (az AkH. esetében szabálypontot, az OH. esetében oldalszámot).
- Az **ex** mezőben példákat találunk. Ezekre a szabályok automatizált tesztelésénél van szükség.
- A **kill** mezőbe (opcionális) azoknak a szabályoknak az egyedi azonosítóit tüntetjük fel, amelyek a a szabályalkalmazó algoritmus futtatásakor konkurensak lehetnek az adott szabályra nézve. Ilyenkor a kill mezőben megadott címkéjű szabályt letiltjuk, így szűkíthető a lehetséges jó megoldások halmaza (illetve az esetlegesen illeszkedő, de valójában nem jó megoldások is eltávolíthatók a kimenetről).

A **rule** mezőben található újraíró szabályok felépítése a következő:

$$(1) X(a=v, \dots) + \dots == Y(a=v, \dots),$$

ahol X a bal oldali, Y a jobb oldali szimbólumot jelenti; a az attribútum nevét, v pedig annak értékét jelöli. A bal oldali szimbólumokban az attribútumok és az értékek közötti operátorok értékvizsgálatot, a jobb oldalon értékadást jelentenek.

A leíró nyelvtan szimbólumai az angol szófaji kategóriák kezdőbetűi vagy -betűcsoportjai: N (főnév), A (melléknév), V (ige), Adv (határozószó), Num (számnév).

A szabályok bal oldalán a következő attribútumok állhatnak:

- **sem**: A token szemantikai tulajdonságait tartalmazza (egyszerre több értéke is lehet). Szemantikai tulajdonság például: színnév, anyagnév, foglalkozásnév stb. A példában szereplő *sem*="Material" attribútum-érték páros jelentése: a bemenetnek olyan tokennek kell rendelkeznie, amely rendelkezik a „Material”, azaz anyagnév szemantikai jeggyel.
- **match**: Értéke egy reguláris kifejezés, amely illeszkedik a morfológiai elemző által előállított címkesorozatra. Ha például azt szeretnénk, hogy a folyamatos melléknévi igenevekre illeszkedjen a szabály, úgy tudjuk beállítani, hogy a *match* attribútumnak megadjuk a következő értéket:

[1] `match~"IGE, _OKEP",`

amelynek jelentése: egy ige és egy -ó/-ő képző. A *match* után szereplő `~` egy speciális operátor, amely reguláris kifejezések illeszkedését vizsgálja (az operátorokról lásd később).

- **wordform**: A bemenet felszíni alakja.
- **stem**: A felszíni alak töve.
- **ncomparts**: Értéke egy egész szám. Megadja, hogy hány összetételi tagból áll az adott szimbólumnak megfelelő token(rész)sorozat. Igazodva a szótagszámlálási szabály (AkH. 138.) előírásához, a két vagy több szótagból álló igekötők összetételi tagnak számítanak (pl. *ellen-*, *elő-*), míg az egy szótagos igekötők nem. Így például az *előadás* token *ncomparts* értéke 1, míg az azonos felépítésű *beadás*-é csupán 1 (azaz egyszerű – nem összetett – szó).
- **nsylls**: Értéke egy egész szám. Megadja, hogy hány szótagból áll az adott token. Az *ncomparts* jegyhez hasonlóan a szótagszámlálási szabály előírásainak megfelelően számolódik ki az értéke: az összetett szó jel és rag nélküli alakjának szótagszámát értjük alatta. Így pl. mind a *kerékpárjavítás*, mind a *kerékpárjavításnak* tokenek *nsylls* értéke 6. A képzők viszont már beleszámítanak a szótagszámba: a *kerékpár-javítási* alak *nsylls* értéke 7 az *-i* képző miatt.
- **ntoks**: A bemenetben megadott tokenek száma. Ha egy vagy több token összetett szó, a tokenek és az összetételi tagok száma eltérő lehet (*ncomparts* és *ntoks* értéke nem mindig egyenlő).
- **join1**, **join2**, **join3**: Ezen attribútumok a kivételes (nem formalizálható) írásmódú összetételek kezelésére szolgálnak. Az előfeldolgozás során, ha a tokenek felszíni alakjai valamilyen kombinációban szerepeltek a

kivételszótárban, megkapják értékül a kivétel kategóriáját (pl. *Jelentessurito*), így az adott kivételeket kezelő szabályok érvényesek lesznek rájuk.

– **type**: Bizonyos speciális típusú főnévi csoportok megjelölésére szolgáló attribútum.

– **sep**: Ez az attribútum kötelezően szerepel egy értékadásban minden szabály jobb oldalán, ahol a bemeneti tokenek közé kerülő szeparátort (üres sztring, szóköz, kötőjel stb.) jelöli.

– **ortho**: az adott elemzési lépésben, a sep attribútum segítségével kiszámított helyesen írt alak.

– **hasnesep**: Értéke egész szám. Azt jelzi, hogy a generált alak hány darab nem egybeírást jelző szeparátort tartalmaz (azaz hány szóközt, kis- vagy nagyköjtjelet). Ha nem tartalmaz egyiket sem a felsoroltak közül, értéke 0. A szótagszámlálási szabályoknál van rá szükség: ezek a szabályok csak akkor hajtódnak végre, ha a generált alak egyáltalán nem tartalmaz kötőjelet vagy szóközt, vagyis egybeírt alakok esetén.

– **63exception**: Értéke 'YES' lehet. A 6:3-as szabály alól kivételt jelentő írásmódú szavak megjelölésére szolgál. A 6:3-as szabályok csak akkor hajthatók végre, ha a 63exception attribútum értéke nem 'YES' (azaz nem kivételes írásmódú szavakról van szó).

– **3idcons**: Ha új összetétel keletkezésekor három azonos mássalhangzó kerül egymás mellé, speciális szabályt kell alkalmazunk, ezen esetek megjelölésére szolgál ez a jegy. (Bővebben l. 3.5, 3.6 pontokban.)

A szabályok jobb oldalán az alábbi attribútumok állhatnak:

– **sep**: Az attribútum értéke a bemeneti szóelemek közé kerülő elválasztó karaktert kódolja. Értékei lehetnek:

- '' (üres sztring), amely egybeírást jelöl;
- ' ' (szóköz), amely különírást jelöl;
- '-', kötőjellel írást jelöl;
- '--' (kötőjel-kötőjel), nagyköjtjellel írást jelöl;
- '@' (kukac-szóköz), amely az anyagnévi mozgószabály speciális szeparátora; a különírt anyagnévi jelzős szókapcsolatot alkalmilag egybeírja („összerántja”), és a jelzett szót különírva kapcsolja hozzá;
- '@-' (kukac-kötőjel), amely a második mozgószabály speciális szeparátora; a különírt minőségjelzős szókapcsolatot alkalmilag egybeírja („összerántja”), és az új összetételi utótagot kötőjellel kapcsolja hozzá;
- '-@' (kötőjel-kukac), amely a második mozgószabály speciális szeparátora; a különírt minőségjelzős szókapcsolatot alkalmilag egybeírja („összerántja”), és az új összetételi előtagot kötőjellel kapcsolja hozzá;
- '-1', '-2', ..., '-n', amely egy többszörös összetétel 1., 2., ..., n. szava mögé kötőjelet szűr be (-n esetén az utolsó tag elé) – a 6:3-as szabályok használják;

- Novák, A., M. Pintér T. 2006. Milyen a még jobb Humor? In: Alexin, Z., Csendes D. (szerk.) 2006. *MSZNY 2006. Magyar Számítógépes Nyelvészeti Konferencia*. Szeged: Szegedi Tudományegyetem, 60–69.
- Pintér, T., Oravecz Cs., Mártonfi A. 2009. Online helyesírási szótár és megvalósítási nehézségei. In: Tanács, A., Szauder D., Vince V. (szerk.) 2009. *MSZNY 2009. Magyar Számítógépes Nyelvészeti Konferencia*. Szeged: JATEPress, 172–182.
- Pomázi Gyöngyi (szerk.) 2000. *A magyar helyesírás szabályai*. 11. kiadás, 12. (példaanyagában átdolgozott) lenyomat. Budapest: Akadémiai Kiadó.