



Az új magyar Braille-rövidírás kialakítása

Sass Bálint

sass.balint@nytud.mta.hu

MTA Nyelvtudományi Intézet

Nyelvtechnológiai és Alkalmazott Nyelvészeti Osztály

2013. november 11.

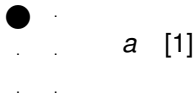


Braille-írás

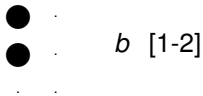
- A vakok által világszerte használt tapintáson alapuló írásmód.
- 6 pont különböző elrendezéseiből → 64 féle karakter

1	4
2	5
3	6

- A vakok által világszerte használt tapintáson alapuló írásmód.
- 6 pont különböző elrendezéseiből → 64 féle karakter



- A vakok által világszerte használt tapintáson alapuló írásmód.
- 6 pont különböző elrendezéseiből → 64 féle karakter

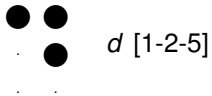


- A vakok által világszerte használt tapintáson alapuló írásmód.
- 6 pont különböző elrendezéseiből → 64 féle karakter

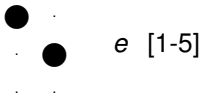


c [1-4]

- A vakok által világszerte használt tapintáson alapuló írásmód.
- 6 pont különböző elrendezéseiből → 64 féle karakter



- A vakok által világszerte használt tapintáson alapuló írásmód.
- 6 pont különböző elrendezéseiből → 64 féle karakter



A magyar Braille-írás

- önálló jelek az *ékezetes* magánhangzóknak

⠁ ⠠⠁ ⠠⠑ ⠠⠑̄ ⠠⠢ ⠠⠢̄ ⠠⠔ ⠠⠔̄ ⠠⠕ ⠠⠕̄ ⠠⠥ ⠠⠥̄ ⠠⠦ ⠠⠦̄
a á e é i í o ó ö ő u ú ü ű

- önálló jelek a *kettős* mássalhangzóknak

⠉ ⠠⠉⠰ ⠠⠒ ⠠⠒̄ ⠠⠓ ⠠⠓̄ ⠠⠔ ⠠⠔̄ ⠠⠕ ⠠⠕̄ ⠠⠖ ⠠⠖̄ ⠠⠗ ⠠⠗̄
c cs g gy l ly n ny s sz t ty z zs

- „tükrözés” – kivéve: *e, i, c, g*

- ha Braille-írást sima szöveggént ábrázolunk, a kettős mássalhangzókat egy karakterrel – az első karakter nagybetűs változatával – jelöljük: C, G, L, N, S, T, Z – pl. *maGar*

A különféle nyelvekre (pl. angol, német, magyar) kidolgozott Braille-rövidírások az általános Braille-írást rövidítési, tömörítési szabályokkal egészítik ki.

Rövidírás használatával:

- a nyomtatott forma kevesebb helyet igényel;
- az írás, jegyzetelés gyorsul;
- az olvasás, felismerés gyorsul.

A magyar Braille-rövidírás: az ún. „kis” rövidírás

Magyarországon használatos sztenderd rövidírás.

Elemei:

- 1 nagybetűjel $\cdot\cdot$ törlése $\sim rk = 2\%$
- 2 vessző utáni szóköz törlése $\sim rk = 2\%$
- 3 határozott névelő utáni szóköz törlése, az rövidítése $\sim rk = 2\%$
- 4 44 szabály: $\sim rk = 4\%$
 - 7 db szóvégi rövidítés:
ban/ben \rightarrow b, ból/ból \rightarrow bl, hoz/hez/höz \rightarrow hz ...
 - 16 db egyjelű szórövidítés:
Cak \rightarrow C, de \rightarrow d, és \rightarrow é ...
 - 21 db kétjelű szórövidítés:
aNNi \rightarrow ai, boldog \rightarrow bg, gond \rightarrow gd ...

\rightarrow a kis rövidírás mért rövidítési képessége: $rk = 9,4\%$



Kezdeményezés + követelmények

Kezdeményezés az MVGYOSZ Braille-bizottsága részéről:
bővítsük a kis rövidírást új szabályokkal annak érdekében,
hogy a rövidítési képessége jelentősen növekedjen.

Használhatósági követelmények:

- az ismert kis rövidírást egészítse ki;
- *jó olvashatóság*: a rövidítések emlékeztessenek az eredetire;
- *jó felismerhetőség*: tapintás útján könnyen felismerhető jelek alkalmazása;
- *könnyű megtanulhatóság*: kevés, egyszerű szabály.

Kompromisszumos javaslatot készül, mely . . .

- törekszik a maximális rövidítési képesség elérésére,
- miközben megfelel a használhatósági feltételeknek is.

Az új rövidítés létrehozása

Mit rövidítsünk?

A magyar nyelv gyakorisági adatai alapján a maximális rövidítési képességgel bíró ideális rövidítés automatikusan kialakítható.

Elv: a lehető leggyakoribb betűkapcsolatokat kell a lehető legrövidebbre rövidíteni.

Pl.: $ty \leftrightarrow et \sim 20\times$ -os különbség a gyakoriságban

→ automatikus algoritmus

Mire rövidítsünk?

Szakértői döntések a gyakorisági adatok és a használhatósági megfontolások alapján. Akár az automatikus eredményt felülbírálvá.

→ manuális munka

Az algoritmus működése

- 1 Számba vesszük az összes rövidíthető nyelvi elemet: előállítjuk a karaktersorozatok gyakorisági listáját.
- 2 Kiszámoljuk az elemek rövidítési képességét a következő módon:
$$rk(w, r(w)) = [l(w) - l(r(w))] \cdot fq(w)$$

w az eredeti rövidítendő karaktersorozat, $r(w)$ a rövidítés, $l()$ a hossz (karakterszám), $fq()$ a gyakoriság.
Pl.: a *sem* és a *maGar* gyakorisága nagyjából azonos
→ érdekesebb az utóbbit két karakterre rövidíteni.
- 3 Rendezzük a gyakorisági listát rövidítési képesség szerint.
- 4 Az *első* elemhez (maximális rövidítési képesség) hozzárendelünk egy megfelelő rövidítést.
- 5 Ez után *alkalmazzuk* a most kialakított rövidítő szabályt az egész listára.
- 6 → 2

Működési példa

2 karakterre rövidítünk

A gyakorisági lista eleje:

w	$f_q(w)$	$l(w) - l(r(w))$	$rk(w)$
<i>Ser</i>	10901	$\times 1 =$	10901
<i>Serint</i>	2515	$\times 4 =$	10060
<i>maGar</i>	3326	$\times 3 =$	9978

Az első elemet rövidítjük: *Ser* \rightarrow *Sr*

Működési példa

2 karakterre rövidítünk

A gyakorisági lista eleje:

w	$f_q(w)$	$l(w) - l(r(w))$	$rk(w)$
<i>Ser</i>	10901	$\times 1 =$	10901
<i>Serint</i>	2515	$\times 4 =$	10060
<i>maGar</i>	3326	$\times 3 =$	9978

Az első elemet rövidítjük: $Ser \rightarrow Sr$

w	$f_q(w)$	$l(w) - l(r(w))$	$rk(w)$
<i>maGar</i>	3326	$\times 3 =$	9978
...			
<i>Srint</i>	2515	$\times 3 =$	7545
...			
<i>Sr</i>	10901	$\times 0 =$	0

Működési példa

2 karakterre rövidítünk

A gyakorisági lista eleje:

w	$f_q(w)$	$l(w) - l(r(w))$	$rk(w)$
<i>Ser</i>	10901	$\times 1 =$	10901
<i>Serint</i>	2515	$\times 4 =$	10060
<i>maGar</i>	3326	$\times 3 =$	9978

Az első elemet rövidítjük: *Ser* \rightarrow *Sr*

w	$f_q(w)$	$l(w) - l(r(w))$	$rk(w)$
<i>maGar</i>	3326	$\times 3 =$	9978
...			
<i>Srint</i>	2515	$\times 3 =$	7545
...			
<i>Sr</i>	10901	$\times 0 =$	0

A következő rövidítendő elem a *maGar* lesz.

A *Serint* az új szabály miatt sokkal lejjebb csúszott. Ha mégis sorra kerül, tovább rövidíthetjük (így lesz).

A *Sr* kétkarakteres rövidítéssel nyilván nem rövidíthető tovább.

Több nagyon hatékony rövidítési ötletet
a használhatósági követelmények miatt végül elvetettünk:

- szóközt elnyelő szabályok (pedig: $rk(t_{_}, T) \approx 1,2\%!);$
- különböző helyzetekben különböző jelentéssel bíró rövidítések (német példa: $r(\text{immer})=r(\text{mm})=x);$
- „értelmetlen” elemek (pl.: *et, en, er, ele, ala, tás, ter*) rövidítése (pedig: $rk(\text{et}, T) \approx 0,8\%!).$

Elv: „értelmest értelmesre”


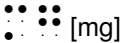

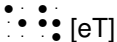
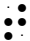




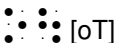
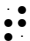






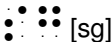


= csak értelmezhető szóelemet rövidítünk
és olvasható rövidítésjelet rendelünk hozzá

- A javaslatban csak *kétkarakteres* rövidítéseket használunk. A legtöbb esetben könnyű kiválasztani egy olyan ritka kétkarakteres rövidítést, ami megfelelően illeszkedik a rövidítendőhöz, pl.: $r(\begin{smallmatrix} \bullet\bullet & \bullet\bullet & \bullet\bullet & \bullet\bullet & \bullet\bullet \\ \bullet\bullet & \bullet\bullet & \bullet\bullet & \bullet\bullet & \bullet\bullet \end{smallmatrix} [\text{maGar}]) = \begin{smallmatrix} \bullet\bullet & \bullet\bullet \\ \bullet\bullet & \bullet\bullet \end{smallmatrix} [\text{mG}]$.
- A leghatékonyabb szabályok a nagyon gyakori betűkapcsolatok (pl.: 'et', 'el') *egy karakterre* való rövidítései lennének. Van néhány nagyon ritka jel, ami itt rövidítésként alkalmazható lenne:
pl.: $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix}$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix}$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix} [\text{@}]$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix} [\text{q}]$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix} [\text{=}]$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix} [\text{*}]$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix} [\text{ty}]$, $\begin{smallmatrix} \bullet\bullet \\ \bullet\bullet \end{smallmatrix} [\text{w}]$,
de használatuk, főleg sok ilyen rövidítésé, nagyban rontja az olvashatóságot. Bizonyos elemek egykarakteres rövidítése még megfontolás alatt van: *el*, *tt*, *meg*, *Ser*, *ás/és*, *eG*.









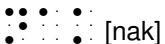




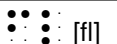

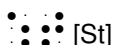

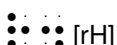

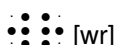
A módszer összefoglalása

- Szigorúan automatikus úton történik a bemutatott algoritmussal az épp aktuális (következő) legjobban rövidíthető elem meghatározása. Egy ajánlott, elég jó olvashatóságú rövidítést is automatikusan megad hozzá a rendszer.
- A használhatósági feltételeknek nehéz teljesen automatizált módon megfelelni, ezért szükséges az automatikusan kialakított szabályrendszer interaktív módosítása, manuális véglegesítése:
 - meghatározzuk a konkrét lehető legjobb rövidítésjelet, adott esetben megfelelő *kinézetű* jel létrehozásával;
 - szükség esetén azt is kiköthetjük, hogy az adott elemet – tudva, hogy ezzel veszítünk a rövidítési képességből – mégsem rövidítjük;
 - a gyakorisági adatok ismeretében megszoríthatjuk a rövidítés pozícióját (pl.: csak szó végén).

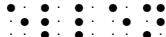


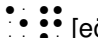

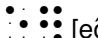

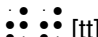
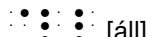
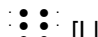








A javaslat 1/3

	rövidítendő	rövidítés	megjegyzés
1.	 [meg]	 [mg]	szó elején (98%)
2.	 [ett]	 [eT]	= ott  ~ 
3.	 [Ser]	 [Sr]	
4.	 [ott]	 [oT]	= ett  ~ 
5.	 [maGar]	 [mG]	
6.	 [jelen]	 [jn]	
7.	 [ség]	 [sg]	= ság
8.	 [lehet]	 [lT]	<i>lh</i> túl gyakori

A javaslat 2/3

	rövidítendő	rövidítés	megjegyzés
9.	 [vezet]	 [vz]	
10.	 [nek]	 [nx]	= nak  ~  szó végén
11.	 [köz]	 [kz]	
12.	 [nak]	 [nx]	= nek  ~  szó végén
13.	 [fel]	 [fl]	szó elején (86%)
14.	 [Serint]	 [St]	szó elején (95%)
15.	 [rend]	 [rH]	csúsztatott <i>d</i> (<i>rd</i> , <i>rn</i> nem jó)
16.	 [ember]	 [wr]	jel kinézete

A javaslat 3/3

	rövidítendő	rövidítés	megjegyzés
17.	 [ellen]	 [oó]	jel kinézete
18.	 [elnök]	 [eö]	
19.	 [elő]	 [eő]	
20.	 [tart]	 [tt]	szó elején
21.	 [áll]	 [LI]	jel kinézete
22.	 [mond]	 [mH]	csúsztatott <i>d</i>
23.	 [támogat]	 [tg]	
24.	 [ért]	 [éT]	
25.	 [ság]	 [sg]	= ség

A javaslat kiértékelése

23 (25) szabály.

Rövidítendő lista:

áll ellen elnök elő ember ért ett/ott fel- jelen köz lehet maGar megmond -nak/-nek rend ság/ség Ser Serint(-) támogat tart(-) vezet

Kis rövidítés rövidítési képessége: $rk = 9,4\%$

A fenti szabályokkal kiegészített. . .

új rövidítés mért rövidítési képessége: $rk = 12,3\%$

A rövidítési képesség növekedése: +2,9 százalékpont (+30%)

Kevés szabály, jelentős növekedés.

Szükség esetén egyszerűen tovább bővíthető a jövőben.

Példa

eredeti: *Bill Gates szerint az internet nem menti meg a világot*



rövidítve: *bill gates szt ,internet n menti mg ,vgot*



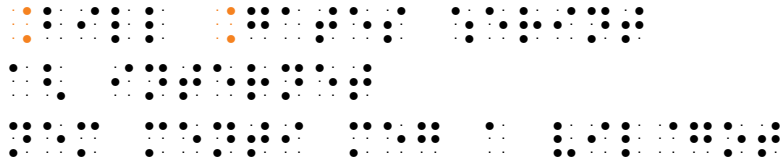
Kis rövidítés: nagybetűjel elhagyása; névelők összevonása; nem, világ

Új rövidítés: szerint, meg

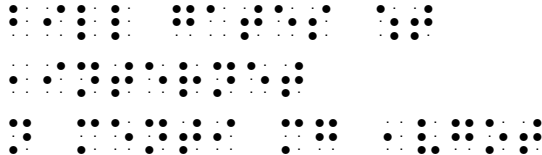
Hossz: 55 karakter

Példa

eredeti: *Bill Gates szerint az internet nem menti meg a világot*



rövidítve: *bill gates szt ,internet n menti mg ,vgot*



Kis rövidírás: **nagybetűjel elhagyása**; névelők összevonása; nem, világ

Új rövidírás: szerint, meg

Hossz: 55 → 53 karakter: $rk() = 3,6\%$

Példa

eredeti: *Bill Gates szerint az internet nem menti meg a világot*

••••• ••••• ••••• •••••• •••••• •••••• ••••••

••••• •••••• •••••• •••••• ••••••

••••• •••••• •••••• •••••• •••••• •••••• •••••• ••••••

rövidítve: *bill gates szt ,internet n menti mg ,vgot*

••••• •••••• •••••• ••••••

••••• •••••• •••••• ••••••

••••• •••••• •••••• •••••• ••••••

Kis rövidírás: nagybetűjel elhagyása; névelők összevonása; nem, világ

Új rövidírás: szerint, meg

Hossz: 55 → 50 karakter: rk() = 9, 1%



••
•• = Ö
••