

# Words and transducers

---

Bognár Ivett

2020. Március 6.

Bevezetés

Véges állapotú morfológiai elemzés

Véges állapotú lexikon

Véges állapotú transzducer

FST morfológiai elemzésre

Lexikon + szabályok

Porter Stemmer

Tokenizáció

Helyesírási hibák

Minimum Edit Distance

# Bevezetés

---

# Mit szeretnénk?

Meg szeretnénk találni az összeset az összes állatból:

Nem csak annyit szeretnénk tudni, hogy a *woodchuck* többesszáma *woodchucks*, hanem felismerni az olyan többesszámokat is, mint:

*fox, foxes; peccary, peccaries; goose, geese*

Ehhez szükségünk van:

- Helyesírási szabályok (*peccary, peccaries*)
- Morfológiai szabályok

Voltaképpen morfológiai elemzőt szeretnénk.

# Miért szeretnénk?

Mire jó a morfológiai elemzés?

- Kap egy felszíni/bemeneti formát *going*  
⇒ *VERB – go + GERUND – ing*
- Automatikus keresés egy adott szó minden alakjára.
- POS tagging
- Helyesírás ellenőrzés
- Gépi fordítás
- A **produktivitás** miatt muszáj, hogy legyen morfológiai elemzőnk.

# Véges állapotú morfológiai elemzés

---

# Véges állapotú morfológiai elemzés

English		Spanish		
Input	Morphological Parse	Input	Morphological Parse	Gloss
cats	cat +N +PL	pavos	pavo +N +Masc +Pl	‘ducks’
cat	cat +N +SG	pavo	pavo +N +Masc +Sg	‘duck’
cities	city +N +Pl	bebo	beber +V +PInd +1P +Sg	‘I drink’
geese	goose +N +Pl	canto	cantar +V +PInd +1P +Sg	‘I sing’
goose	goose +N +Sg	canto	canto +N +Masc +Sg	‘song’
goose	goose +V	puse	poner +V +Perf +1P +Sg	‘I was able’
gooses	goose +V +1P +Sg	vino	venir +V +Perf +3P +Sg	‘he/she came’
merging	merge +V +PresPart	vino	vino +N +Masc +Sg	‘wine’
caught	catch +V +PastPart	lugar	lugar +N +Masc +Sg	‘place’
caught	catch +V +Past			

**Figure 3.2** Output of a morphological parse for some English and Spanish words. Spanish output modified from the Xerox XRCE finite-state language tools.

Kell hozzá:

1. **lexikon**: tövek és affixumok listája, minimális információval (igei vagy főnévi tő, stb.)
2. **morfotaktika**
3. **helyesírás**

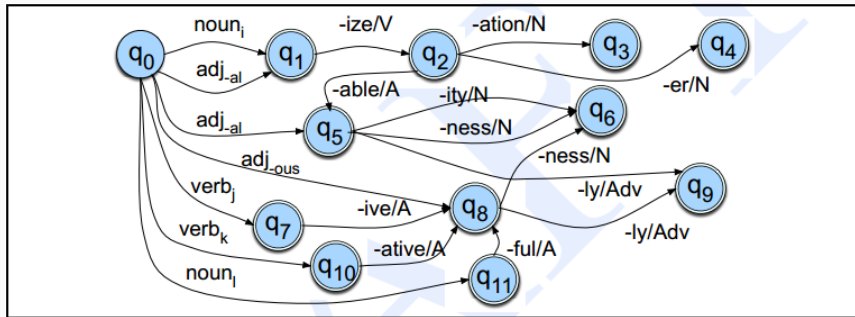


# Véges állapotú lexikon

---

- Egy lexikonba nem tudunk minden szót felvenni.
- A lexikonok általában tövek és affixumok listája + morfolaktika, ami elmondja, hogyan kombinálhatjuk őket.
- Morfolaktika modellezésére: **véges állapotú automaták.**

# Angol melléknevek



# Véges állapotú transzducer

---

# Véges állapotú transzducer

- "Finite-state transducer or FST is a type of finite automation which maps between two sets of symbols."
- Az FST relációt definiál két string-halmaz között.
- Vagyis: Olvas egy stringet és generál egy másikat.
  1. **Felismerőként:** vesz egy string párt inputként és outputként, **elfogadja**, ha a pár része a nyelvnek, és **elutasítja**, ha nem.
  2. **Generátor:** a nyelv string párait adja vissza.
  3. **Fordító:** olvas egy stringet, és visszaad egy másikat.
  4. **Set relater**

$Q$	a finite set of $N$ states $q_0, q_1, \dots, q_{N-1}$
$\Sigma$	a finite set corresponding to the input alphabet
$\Delta$	a finite set corresponding to the output alphabet
$q_0 \in Q$	the start state
$F \subseteq Q$	the set of final states
$\delta(q, w)$	the transition function or transition matrix between states; Given a state $q \in Q$ and a string $w \in \Sigma^*$ , $\delta(q, w)$ returns a set of new states $Q' \subseteq Q$ . $\delta$ is thus a function from $Q \times \Sigma^*$ to $2^Q$ (because there are $2^Q$ possible subsets of $Q$ ). $\delta$ returns a set of states rather than a single state because a given input may be ambiguous in which state it maps to.
$\sigma(q, w)$	the output function giving the set of possible output strings for each state and input. Given a state $q \in Q$ and a string $w \in \Sigma^*$ , $\sigma(q, w)$ gives a set of output strings, each a string $o \in \Delta^*$ . $\sigma$ is thus a function from $Q \times \Sigma^*$ to $2^{\Delta^*}$

Ahogy az FSA-k izomorfok a regulális nyelvekre, úgy az FST-k izomorfok a regulális relációkra.

A **reguláris relációk** olyan kifejezések, melyek két reguláris kifejezésből állnak, melyeket  $:$  választ el. Egy ilyen azt jelenti, hogy a kettőspont előtti kifejezésnek (a reláció felső tagjának) megfelelő nyelv minden eleme relációban áll a kettőspont után álló kifejezésnek (a reláció alsó tagjának) megfelelő nyelv minden elemével.

FST-k két fontos művelete:

- **inverzió**: megfordítja az input és az output címkéket.
- **kompozíció**

Míg az FSA-król elmondhattuk, hogy mindegyik átalakítható determinisztikussá, az FST-kről ez már nem mondható el.

- **Sequential transducers:** determinisztikusak az input tekintetében.
- **Subsequential transducers:** a sequential transducer általánosítása, generál további output stringet a végállapotnál, hozzákonkatenálva az eddig produkált outputhoz.



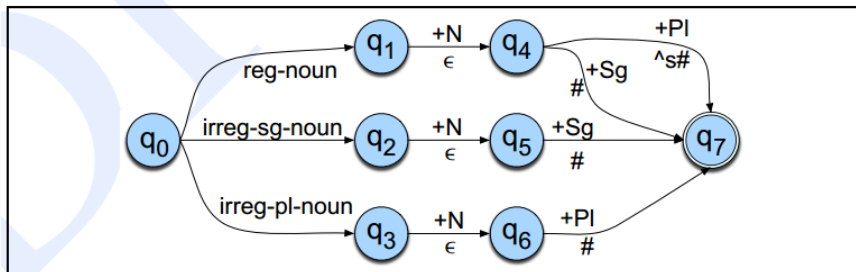
A szavakat úgy vesszük, hogy van egy lexikális szintjük (**lexical level**) és egy felszíni szintjük (**surface level**).

FST esetén ez azt jelenti, hogy van kettő tape:

- **lexical tape** from  $\Sigma$
- **surface tape** from  $\Delta$

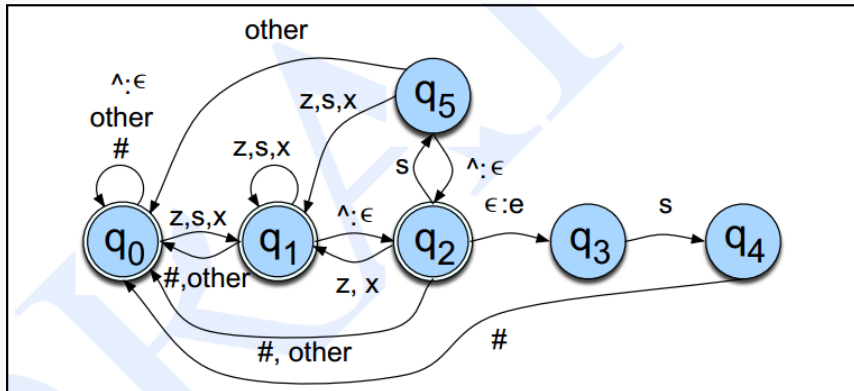
$\Sigma + \Delta = \Sigma'$  : *komplexszimbólumokvégesábécéje.*

A szimbólumpárokak ebben megvalósítható pároknak (**feasible pairs**), illetve találhatunk benne párokat, mint pl.  $a:a$ , amit **default párnak** hívunk.



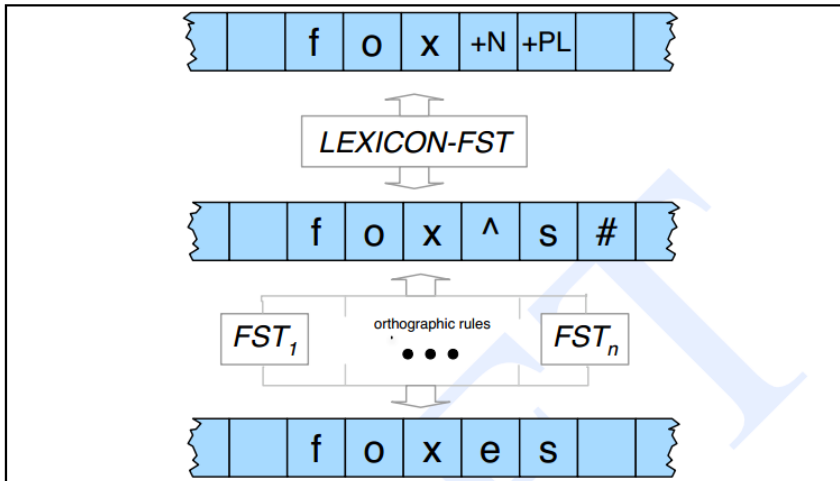
**Figure 3.13** A schematic transducer for English nominal number inflection  $T_{num}$ . The sym-

# Angol többesszám - jobban

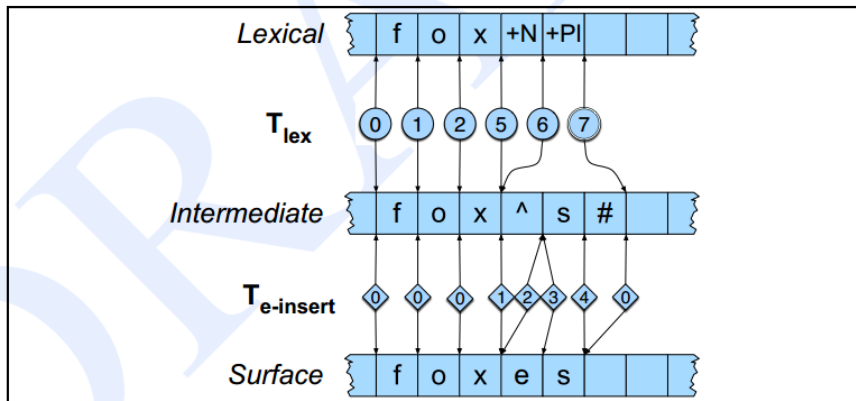


## Lexikon + szabályok

---



**Figure 3.19** Generating or parsing with FST lexicon and rules



**Figure 3.20** Accepting *foxes*: The lexicon transducer  $T_{lex}$  from Fig. 3.14 cascaded with the E-insertion transducer in Fig. 3.17.

- A *foxes* lehet ige is.
- A transducer ilyen esetekben nem képes dönteni, szükség van a szövegkörnyezet ismeretére is. (*That trickster foxes me every time!*)  
A legjobb, amit tehetünk, hogy felsoroljuk az összes lehetőséget.  
Amit viszont kezelniük kell: local ambiguity.  
"Transducers in parallel can be combined by **automaton intersection**."

# Porter Stemmer

---



- lemmatizálás  
stemming / (szó)tövesítés vs. lemmatizálás A **lemma** mindig értelmes szóalak, a **szótő** pedig olykor nem.
- tokenizáció – példányosítás

**Porter Stemmer** (1980) Egy egyszerű és viszonylag hatékony algoritmus tövesítésre.

## Szabályok:

ATIONAL → ATE (e.g., relational → relate)

ING → ε if stem contains vowel (e.g., motoring → motor)

SSES → SS (e.g., grasses → grass)

Hibák:

<b>Errors of Commission</b>		<b>Errors of Omission</b>	
organization	organ	European	Europe
doing	doe	analysis	analyzes
numerical	numerous	noise	noisy
policy	police	sparse	sparsity

# Tokenizáció

---

Azt hinnénk, hogy csak levágjuk a whitespace-t és megvagyunk, de azért ennyire nem egyszerű.

- A tokenizálók általában külön szedik a többszavas kifejezéseket.  
(*New York*)  
A kérdéskör kicsit a névfelismerésé is.
- Az írásjelek nem mindig egyértelműek: *Mr.*, *Inc.*, stb.
- A mondat szegmentállása lesz az első lépés!  
Szabály alapon, vagy machine learninggel.
- Egyszerű, de nagyszerű tokenizálót lehet csinálni reguláris kifejezésekkel.

# Helyesírási hibák

---

Miért szeretjük felismerni őket?

- Keresésekhez
- Kézírás felismerőkhöz

Milyen problémákat kell kezelni?

- **non-word error detection**

Olyan helyesírási hibák felismerése, amik nem létező szóhoz vezetnek, pl. *graffe*, *giraffe* helyett

- **isolated-word error correction**

Kijavítjuk a zsiráfot, de csak önmagában nézzük.

- **context-dependent error detection and correction**

A kontextust is használjuk, hogy felismerjünk egy hibát, még akkor is, ha az véletlenül éppen egy valódi szó. (**real-word error**) Pl. *desert*, *dessert* helyett.

# Minimum Edit Distance

---

## Két string távolsága

Két string távolsága az a szám, ami kifejezi, hogy mennyire hasonló két string egymáshoz.

		i	n	t	e	n	t	i	o	n
delete i	→									
substitute n by e	→	n	t	e	n	t	i	o	n	
substitute t by x	→	e	t	e	n	t	i	o	n	
insert u	→	e	x	e	n	t	i	o	n	
substitute n by c	→	e	x	e	n	u	t	i	o	n
		e	x	e	c	u	t	i	o	n

**Figure 3.24** Operation list transforming *intention* to *execution* (after Kruskal 1983)



$$(3.5) \quad distance[i, j] = \min \begin{cases} distance[i-1, j] + \text{ins-cost}(target_{i-1}) \\ distance[i-1, j-1] + \text{sub-cost}(source_{j-1}, target_{i-1}) \\ distance[i, j-1] + \text{del-cost}(source_{j-1}) \end{cases}$$

# Levenshtein távolság

<http://www.let.rug.nl/~kleiweg/lev/>

<b>n</b>	<i>9</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>12</i>	<i>11</i>	<i>10</i>	<i>9</i>	<i>8</i>
<b>o</b>	<i>8</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>10</i>	<i>9</i>	<i>8</i>	<i>9</i>
<b>i</b>	<i>7</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>9</i>	<i>8</i>	<i>9</i>	<i>10</i>
<b>t</b>	<i>6</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>
<b>n</b>	<i>5</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>10</i>
<b>e</b>	<i>4</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>9</i>
<b>t</b>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>8</i>
<b>n</b>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>7</i>	<i>8</i>	<i>7</i>
<b>i</b>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>6</i>	<i>7</i>	<i>8</i>
<b>#</b>	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>
	<b>#</b>	<b>e</b>	<b>x</b>	<b>e</b>	<b>c</b>	<b>u</b>	<b>t</b>	<b>i</b>	<b>o</b>	<b>n</b>

**Figure 3.26** Computation of minimum edit distance between *intention* and *execution* via algorithm of Fig. 3.25, using Levenshtein distance with cost of 1 for insertions or deletions, 2 for substitutions. In italics are the initial values representing the distance from the empty string.

*Köszönöm a figyelmet!*