

Nyelvtechnológia

4

BME, 2007. november 27.



a
b
c
d
f

g

h
i
j
k
l



Mire elég a szintaxis?

INPUT: Jane sold a book to Bill

MEANING:

```
(S (NP (NOUN Jane))
  (VP (VERB sold)
    (NP (DET a)
      (NOUN book))
    (PP (PREP to)
      (NP
        (NOUN Bill))))))
```

Lehet, hogy többet érne a „jelentés”?

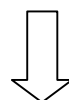
INPUT: Jane sold a book to Bill

MEANING:

(sell	(agent	Jane)
	(object	a book)
	(counter-agent	Bill)
	(tense	past))

Sok mondat - egy jelentés

Jane sold a book to Bill.
and
 A book was sold by Jane to Bill.
and
 Bill was sold a book by Jane.
and
 A book was sold to Bill by Jane.



(sell	(agent	Jane)
	(object	a book)
	(counter-agent	Bill)
	(tense	past))

Hasonló mondat - különböző jelentés

Mom baked for 3 hours

The pie baked for 3 hours

(baked	(agent	Mom)
	(duration	3 hours))

(baked	(object	pie)
	(duration	3 hours))

„Mély” esetek

agent
beneficiary
instrument
source
goal
time
duration

Híres fogalmi hálók

AZ MI kezdetén:

Quillian, Minsky, Charniak, ...

Fogalmi függőség:

Schank

Logikák:

Hendrix, Sowa (fogalmi gráfok), ...

Ontológiák:

CyC, MindNet, FrameNet, ...

WordNet (pszichológusok indították):

WordNet, EuroWordNet,
eXtendedWordNet, ...

Szemantikus web:

(?)

A fogalmi függőség igeosztályai

ATRANS	abstract transfer
PTRANS	physical transfer
MTRANS	mental transfer
MBUILD	build new information
MOVE	move a body part
PROPEL	apply force to an object
GRASP	grasp an object
INGEST	take something into body
EXPEL	expel something from body
SPEAK	make a voiced sound
ATTEND	focus a sense organ
DO	unspecified act

Eseményábrázolás a FF elméletében

Mary took a book from John.

((actor	Mary)
(action	MTRANS)
(object	book)
(direction	(to Mary)
	(fromJohn)))

A fogalmi függőség állapotosztályai

<i>state</i>	<i>scale</i>	<i>examples</i>
Health	-10 to 10	-10 dead -3 sick
Fear	-10 to 0	-9 terrified -2 anxious
Mental state	-10 to 10	-7 depressed -2 sad 3 happy 7 euphoric 9 ecstatic
Hunger	-10 to 10	-8 starving -6 ravenous 5 full 8 stuffed

Schank (1)

1. $PP \longleftrightarrow ACT$

2. $PP \longleftrightarrow PA$

3. $PP \longleftrightarrow PP$

4. PP
 \uparrow
 PA

5. PP
 $\uparrow \uparrow$
 PP

6. $ACT \xleftarrow{o} PP$

7. $ACT \xleftarrow{R} \begin{cases} PP \\ PP \end{cases}$

John \xleftrightarrow{P} PTRANS

John \longleftrightarrow height (>average)

John \longleftrightarrow doctor

boy
 \uparrow
 nice

dog
 $\uparrow \uparrow$ POSS-BY
 John

John \xleftrightarrow{P} PROPEL \xleftarrow{o} cart

John \xleftrightarrow{P} ATRANS $\xleftarrow{R} \begin{cases} John \\ Mary \end{cases}$
 $\uparrow o$
 book

John ran.

John is tall.

John is a doctor.

A nice boy

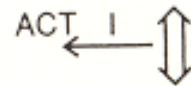
John's dog

John pushed
 the cart.

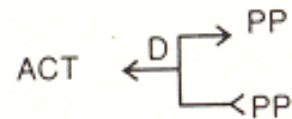
John took the
 book from Mary.

Schank (2)

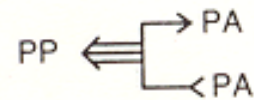
8.



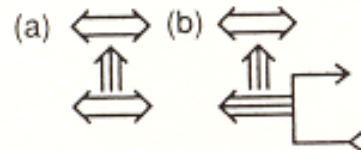
9.



10.



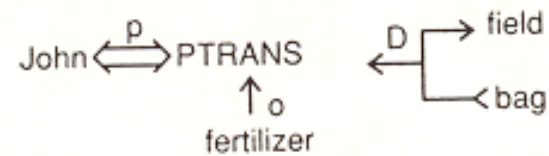
11.



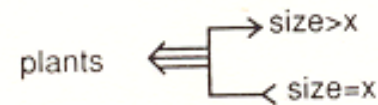
12.



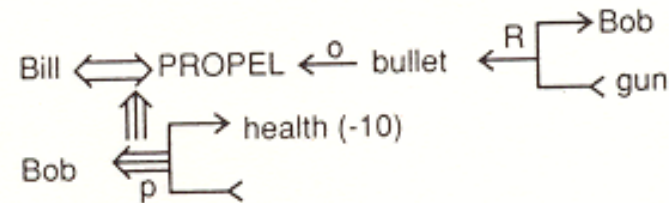
John ate ice cream.



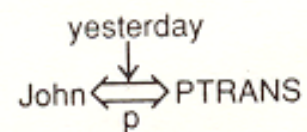
John fertilized the field.



The plants grew.



Bill shot Bob.



John ran yesterday.

a
b
c
d
f
h
i
j
k
l

Forgatókönyvek

Jane was hungry. She decided to go to a restaurant. She ordered spaghetti and a Pepsi. The waitress brought it quickly so when she left she left her a large tip.

QUESTION:

Did Jane eat anything?

Az „étterem” forgatókönyve (a tipikus eseménysor)

entering

seating

ordering

serving

eating

paying

leaving



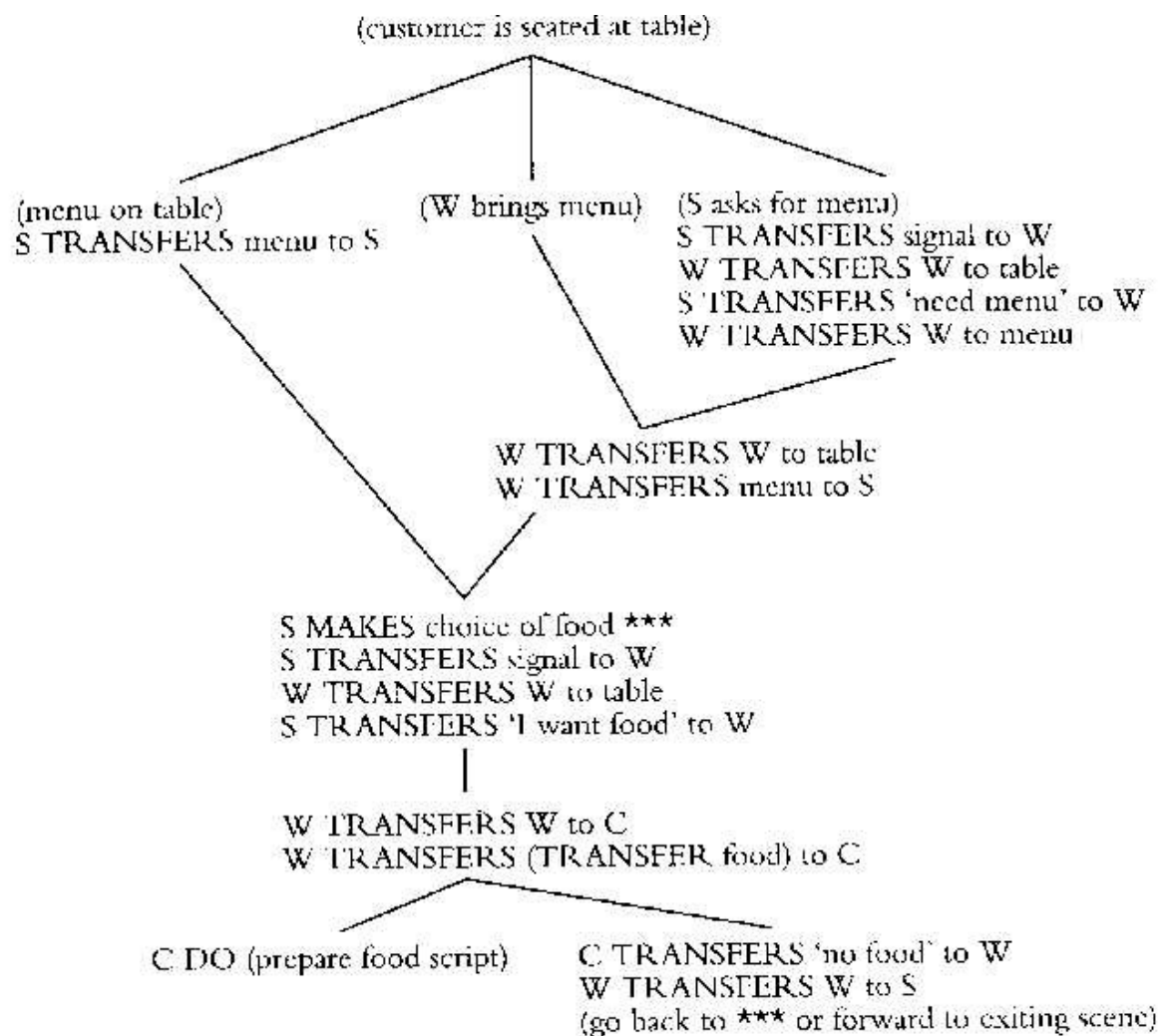
Az „étterem” forgatókönyve (alapismeretek)

ROLES	<i>constraints</i>	<i>defaults</i>
PATRON	*HUMAN*	*ADULT*
WAITER	*HUMAN*	*FEMALE* *ADULT*
COOK	...	
MANAGER	...	

HEADER	
PLANNER	PATRON
GOAL	SATISFY(HUNGER) SATISFY(SOC-INTERACTION)

BODY
EVENT_CHAIN

Az „étterem” teljes forgatókönyve



Szótárak és terminológiakezelés

- ❑ nyomtatott szótárak és elektronikus szótárak
- ❑ terminológiai adatbázisok
- ❑ közvetlen és közvetett elektronikus szótárak
- ❑ egynyelvű, kétnyelvű és többnyelvű szótárak
- ❑ a forrásnyelv és a célnyelvek aszimmetriája

Szerkesztési elvek

- ❑ Az (önálló ill. utaló) szócikkek és felépítésük
- ❑ A szócikkfej: címszó, homonimák és álhomonimák, alak- és írásváltozatok, kiejtés, elválasztás, szófaj, főbb toldalékos alakok, nyelvtani megjegyzés, stílusminősítés
- ❑ Jelentéscsoportok (alapjelentés és jelentésárnyalatok): értelmezések (ekvivalensek) és példák
- ❑ Szóláshasonlatok, közmondások, más szavakkal alkotott összetételek, származékszók

Keresés a szótár(ak)ban

- ☐ betű szerint
- ☐ csonkolt keresés
- ☐ hasonlósági keresés (fuzzy, spell)
- ☐ nyelvi alapú keresés a bemeneti oldalon
- ☐ nyelvi alapú keresés a találati oldalon
- ☐ a kifejezések kezelésének problémái:
alcímszók, kulcsszó-választás, indexek,
egyazon kifejezés több címszó alatt
- ☐ „könyvespolc”: egységes felület
- ☐ egyidejű használat: párhuzamos(nak tűnő)
keresés

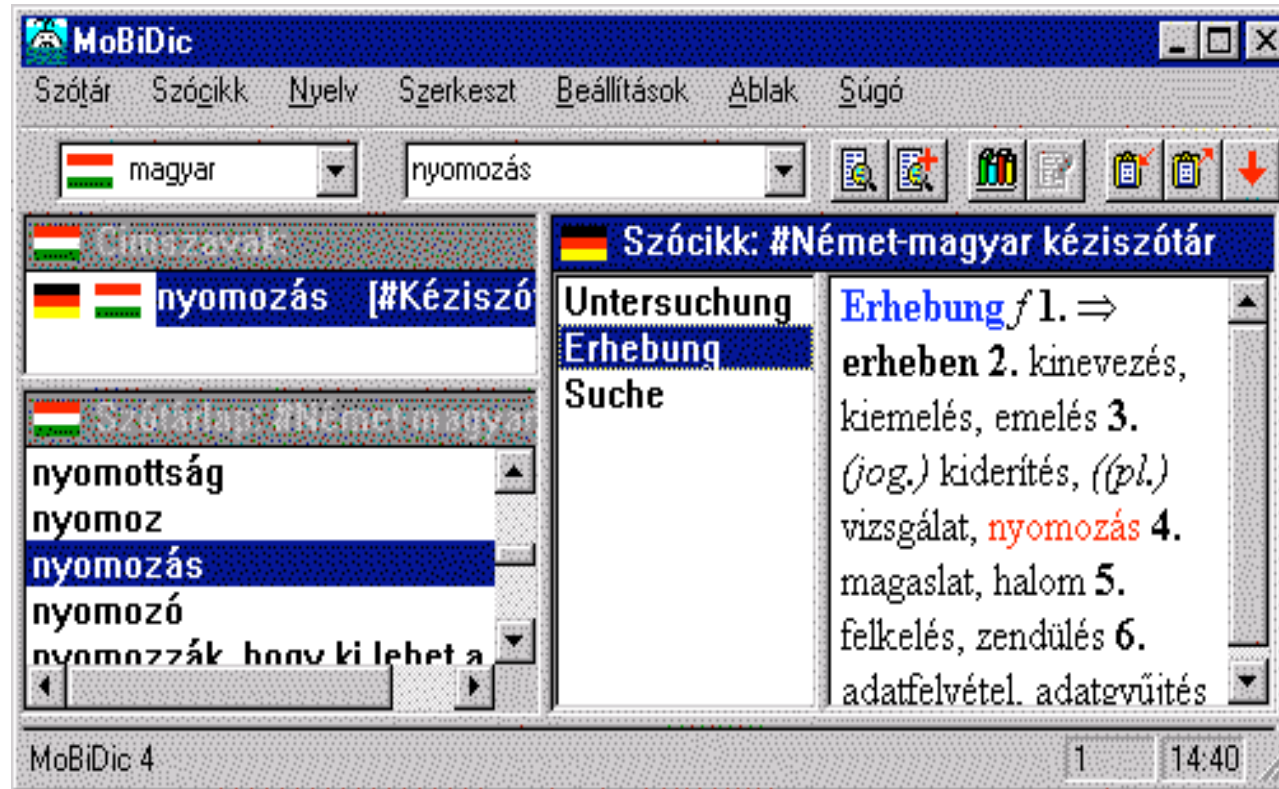
Többszavas kifejezések keresése

- ☐ csak címszóként
- ☐ betű szerint
- ☐ teljes szövegű kereséssel
- ☐ reguláris kifejezésként
- ☐ tőindexekkel: készítéskor vagy elemzési időben (is)

A szótári jobboldal szerepe

- ☐ papírszótárak esetében: csak tipográfiai
- ☐ elektronikusan: új lehetőség
- ☐ ábécé-környezet helyett szinonimák
- ☐ többféle jelentés kezelése a baloldali címszavak segítségével
- ☐ új találati ablak
- ☐ elektronikusan érdemes „kifordítani” a szótárakat

Az elektronikus szótárak megfordíthatók



Gyorsfordítók

- ❑ amikor információ kell, pl. szótári, akkor: csak amit kérek, nem többet, de azt gyorsan, kevés aktív művelettel és a lehető legautomatikusabban!
- ❑ kialakul a „pop-up” viselkedés
- ❑ a kijelölhetőség, ill. az automatikus indíthatóság szerepe

A „rávetítő” megoldás lépései

- ❑ szöveg(rész)-felismerés
- ❑ nyelvi elemzés: morfológia, lemmák, szókapcsolatok (esetleg környezetelemzés)
- ❑ szótári keresés: tövesítve vagy csak literálisan
- ❑ megjelenítés: buborékban vagy fix ablakban
- ❑ log: automatikus információgyűjtés lehetősége

A fordítómemória gondolata

A lefordítandó mondat:

After a few seconds, a window will appear in which you are expected to enter a valid User ID and (if *necessary*) a password.

Korábban már fordítottuk ezt:

After 5 seconds, a window will appear on the screen in which you are expected to enter a User ID and (if *required*) a password.

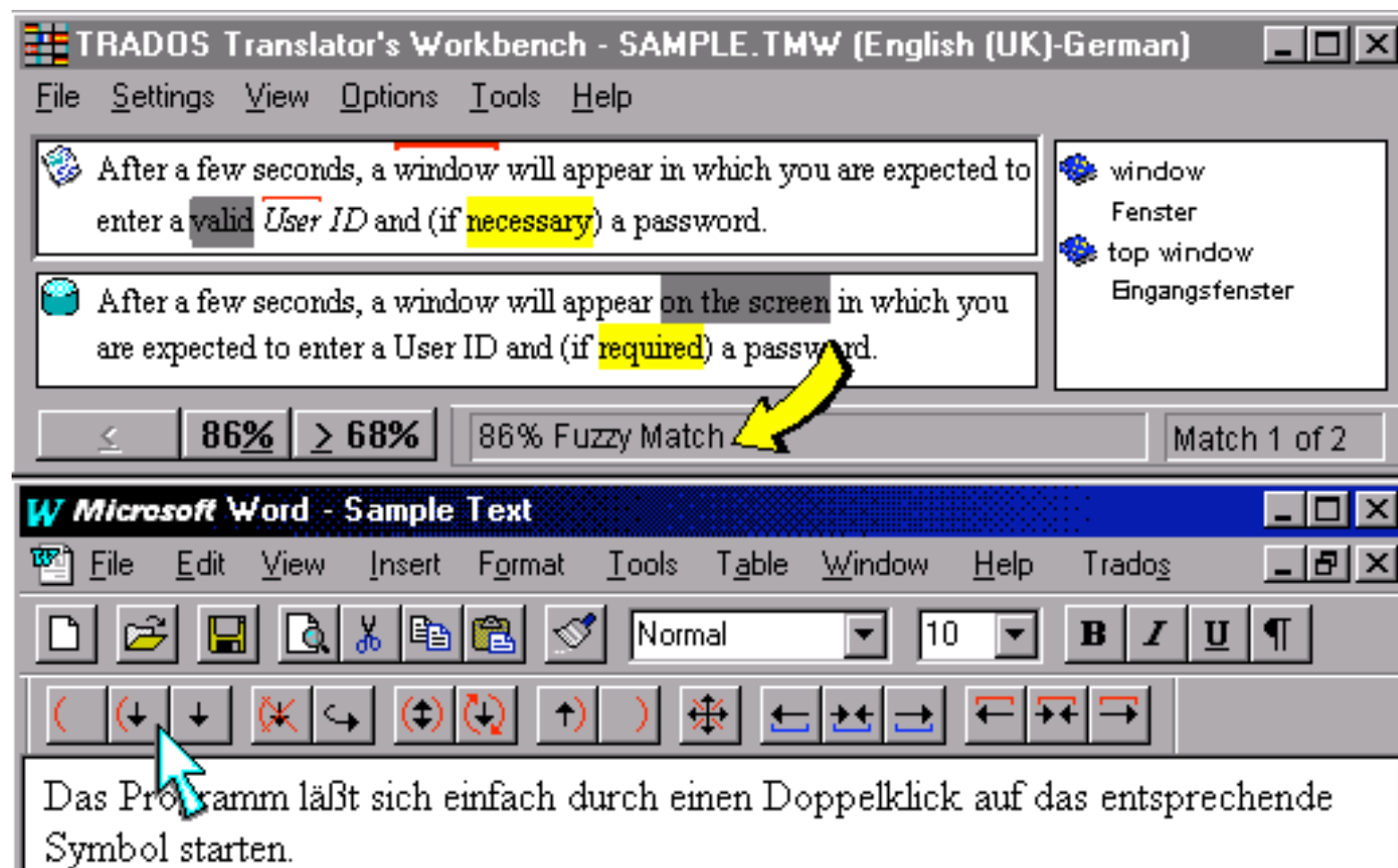
Méghez a így:

Öt másodperc múlva egy ablak jelenik meg a képernyőn, amelybe be kell gépelni egy felhasználó-azonosítót és (ha szükséges) egy jelszót.

Ebből a következő fordítás könnyen előállhat:

Néhány másodperc múlva egy ablak jelenik meg, amelybe be kell gépelni egy érvényes felhasználó-azonosítót és (ha szükséges) egy jelszót.

A fordítómemória mint eszköz



Szövegszinkronizálás

- ❑ bi-text
- ❑ párhuzamos korpuszok
- ❑ szinkronizálás: valós időben és utólag
- ❑ pl. a Biblia

„You will not surely die,” the
serpent said to the woman.
(Genesis 3:4)

A kígyó erre azt mondta az
asszonynak: „Dehogyan haltok meg!”
(Ter 3,4)

Szövegszinkronizálási szintek

- ☐ bekezdésszint
- ☐ mondat szint
- ☐ frázis-szint (?)
- ☐ szószint (??)
- ☐ mondathatár-problémák
- ☐ horgonyok
- ☐ statisztikai módszerek

Nem feltétlenül 1-1 értelmű

(1 = 1,2) O stylographe à la plume de platine, que ta course rapide et sans heurt trace sur le papier au dos satiné les glyphes alphabétiques qui trans-mettront aux hommes aux lunettes éti-ce-lantes le récit narcissique d'une double ren-contre à la cause autobusilistique.

(1 = 1) Ó, platinahegyű töltőtoll!
(2 = 1) Vajha tajtékos-gyors futásod a szaténhátú papirosra róná amaz alfabéta-cikornyákat, melyek a csillogó okulárés emberek tudomására hozzák az autóbuszilisztikus-okú találkozás önbálványozó krónikáját!

A nyelvi szerkezetek hasonlóságáról

zöld kutya
zöld macska
sárga kutya
sárga macska
piros egér
kis asztal
hét kis ágy
a tegnapi buliról
elmentem a tegnapi buliról
beléptünk az EU-ba
jó napot!

A gépi fordítás alapszervei

- szabály-alapú:
közvetlen fordítás
közvetítőnyelves fordítás
transzfer rendszerek



- statisztikai

Becslések az európai nyelven írt internetes szövegek lehetséges méretéről

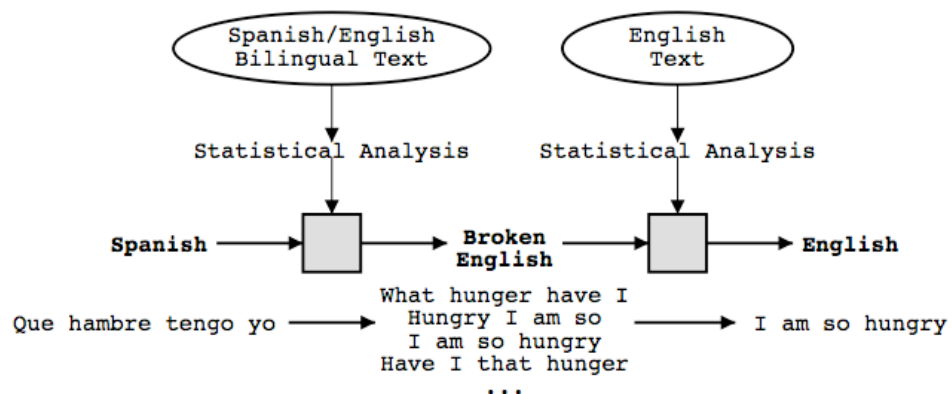
	Nyelv	Szó a weben
1	angol	76 598 718 000
2	német	7 035 850 000
3	francia	3 836 874 000
4	spanyol	2 658 631 000
5	olasz	1 845 026 000
6	portugál	1 333 664 000
7	holland	1 063 012 000
8	svéd	1 003 075 000
8	norvég	609 934 000
10	cseh	520 181 000
11	magyar	457 522 000
12	dán	346 945 000
13	finn	326 379 000
14	lengyel	322 283 000
15	szlovák	216 595 000
16	katalán	203 592 000

	Nyelv	Szó a weben
17	török	187 367 000
18	maláj	157 241 000
19	horvát	136 073 000
20	szlovén	119 153 000
21	észt	98 066 000
22	ír	88 283 000
23	román	86 392 000
24	eszperantó	57 154 000
25	latin	55 943 000
26	baszk	55 340 000
27	izlandi	53 941 000
28	lett	39 679 000
29	litván	35 426 000
30	velszi	14 993 000
31	breton	12 705 000
32	albán	10 332 000

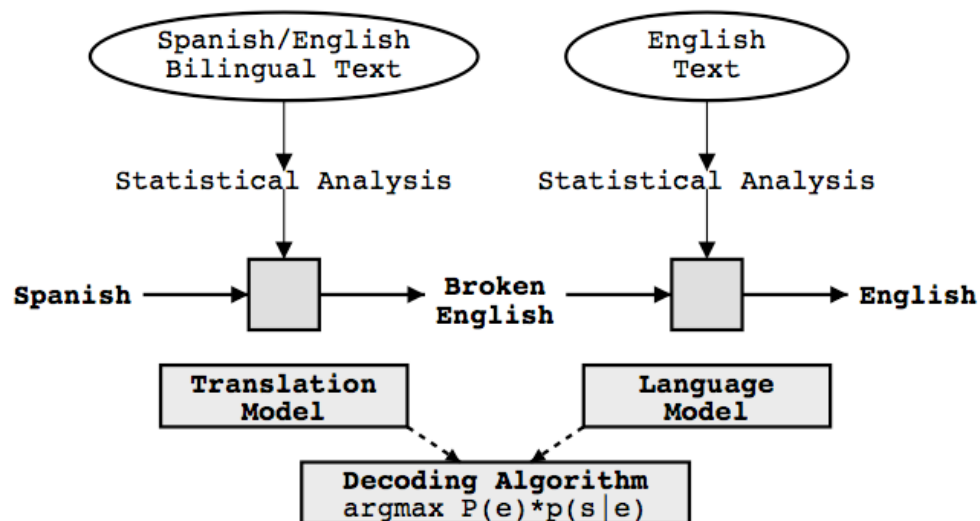
... milyen „minőségű” szövegek vannak a weben?

	Full
internet	2 460 000 000
internte	67 400
interent	681 000
intenret	116 000
intrenet	193 000
inetrnet	128 000
itnernet	66 400
ninternet	47 700
interne.	19 200 000
intern.t	1 940 000
inter.et	19 400 000
inte.net	2 480 000
int.rnet	436 000
in.ernet	522 000
i.ternet	441 000
.ninternet	1 150 000

Statisztikai gépi fordítás

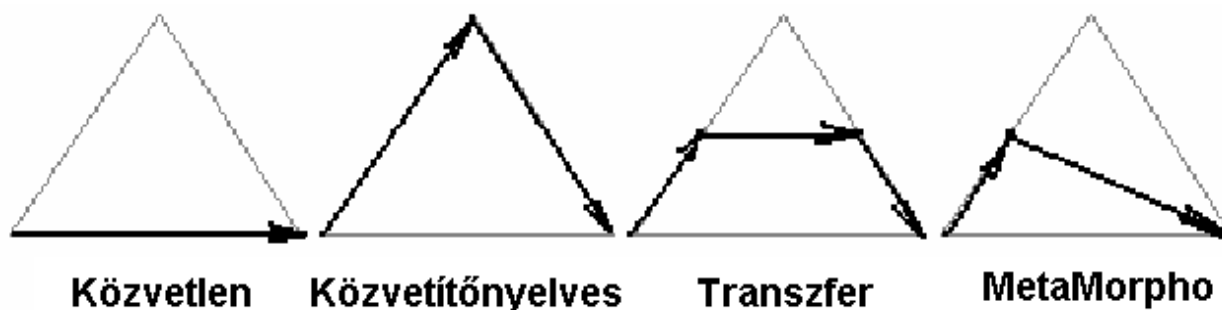


- ☐ Fordítási modell
- ☐ Nyelvmodell
- ☐ Bayes



MetaMorpho-elvek

- **Nincs** külön szótár és külön nyelvtan
- **Csak minta-párok vannak:** bemenet/interpretáció szerkezet-párok
- **Egyetlen elemzési menet:** nincs rákövetkező művelet (pl. transzfer)
- Célszerkezet-generálás:
az elemzés „melléktermékeként”
- Új:



Minták: általánosított nyelvészeti információk

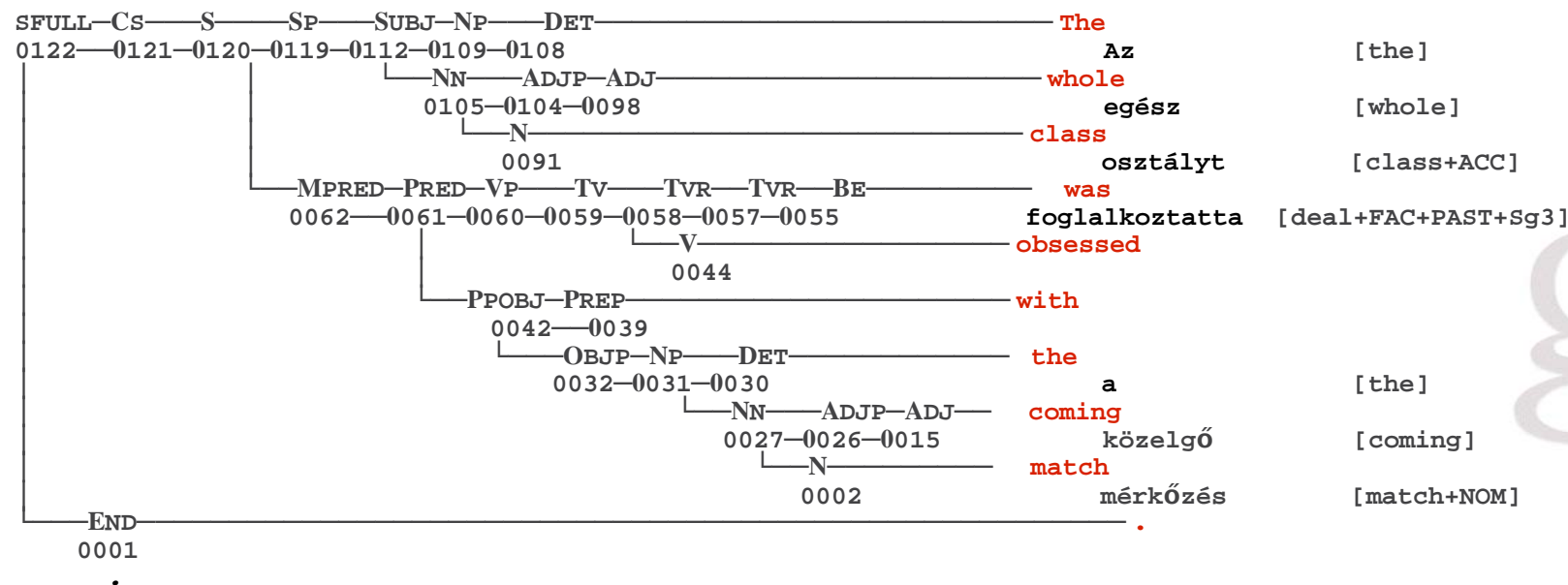
- Rövid, specifikus minták:
szótári címszavak
- Hosszabb, specifikus minták:
többtagú kifejezések
- Részlegesen alulspecifikált minták:
kollokációk, idiómák
- Teljesen alulspecifikált minták:
nyelvészeti szabályok
- Fordítástámogató nyelv:
minta-interpretáció párok

A MetaMorpho projekt

- **A projekt:** 1991-től folyamatosan készített moduljaink felhasználásával (kb. 100 emberév) 2000-ben indult, belső projektként (semmilyen külső támogatása nincs)
- **Cél:** mondatszintű fordítás – új elven: a szavak környezetének felhasználásával (egy n elemű mondatban éppen n darab $(n-1)$ elemből álló környezet van)
- **Forrásnyelv:** angol, magyar
- **Célnyelv(ek):** magyar, angol, ...
- **Szakterület:** nincs – de dinamikusan bővíthető
- **Minta-alapú:** példák (TM) és szabályok (MT) egységesen
- **Minták száma:** kb. 200.000
- **Lexikon:** kb. 100.000 alapszó
- **Elvárt sebesség:** 50 karakter/s
- **Felhasználói felület:** MoBiCAT, MoBiWAP, MMO-Office, MorphoWord, MoBiWeb, webforditas.hu

A MetaMorpho „belülről”

EN: The whole class was obsessed with the coming match.




HU: Az egész osztályt foglalkoztatta a közelgő mérkőzés.


Angol-magyar gyorsfordító szolgáltatás

MoBiCAT: teljes mondatok fordítása
(MoBiCAT-szerver akár intraneten vagy interneten)

20th cent
with show
ceiling m
consists of bathroom, kitchen and a nice pantry. The original
beams of the upper part are unique in Europe.

MoBiCAT

 It is not too far from the building where Franz Kafka lived.

 Ez nincs túl messze az épülettől, ahol Franz Kafka élt.

Visszajelzés: F2

Angol-magyar weblap-fordítás

(MorphoWeb, webforditas.hu)

ARCHIVE | CLASSIFIED | SHOPPING | PROMOTIONS | GAMES | FAST TIMES | MY T

Search

GO

CAR LOCATOR

DATING

CREME

September 13 2005

MAKE TIMES ONLINE
YOUR HOMEPAGE / BOOKMARK

TIMES ONLINE

Home

Britain

World



Petrol panic begins to spread

Emergency powers to reserve fuel for essential users will be reviewed by ministers and oil companies

[Drivers race to fill up despite call for calm](#)
[France puts the pressure on Brown](#)



MorphoWeb

Fordít

☒ URL

Nyelv: angol-magyar

Linkek: Mindent követ

Visszajel

ARCHÍVUM

APRÓHIRDETÉSEK

VÁSÁRLÁS

ELŐLÉPTETÉSEK

JÁTEKOK

BŐJTÖK

AZ

Keresés

GO

CAR LOCATOR

DATING

CREME

September 13 2005

MAKE TIMES ONLINE
YOUR HOMEPAGE / BOOKMARK

IDŐK ONLINÉK

Otthon

Britannia

Világ

Üzlet

Pénz



A benzinpánik elkezd terjedni

Hatalmakat, hogy lényeges felhasználóknak foglaljanak le üzemanyagot, fognak áttekinteni miniszterek és olajvállalatok

[A vezetők versenyeznek hogy nyugalomnak szóló hívás ellenére teljenek meg](#)

[Franciaország a nyomást gyakorolja Brownra](#)



a
b
c
d
f
g
h
i
k
l

