

# Jezikovne tehnologije

**Hiter razvoj računalništva in informatike je posegel na številna področja človekovega udejstvovanja. Tudi obdelava naravnih jezikov se temu ni mogla izogniti. Tako smo dobili številna orodja, brez katerih si danes težko zamisljamo tudi povsem običajna opravila, kot je npr. pisanje.**

**Miro Romih**

Preprosto povedano so jezikovne tehnologije vse tiste informacijske tehnologije, ki se ukvarjajo z naravnimi jeziki, med katere sodi tudi slovenščina. Jezikove tehnologije postajajo vse pomembnejši sestavni del informacijskih tehnologij, to pa je v osnovi omogočilo šele dovolj velika procesna moč sodobnih računalnikov in drugih naprav, ki so hkrati sposobne hraniti velikansko število informacij. Naravni jeziki so namreč izjemno velika zakladnica besed, zapletenih pravil in znanja, ki ga je človeštvo ustvarjalo tisoče let. Zaradi zahtevnosti bi zato jezikovne tehnologije zlahka uvrstili ob bok drugim, bolj znanim visokim tehnologijam, kot so npr. nanotehnologija, telekomunikacije, umetna inteligenca in robotika.

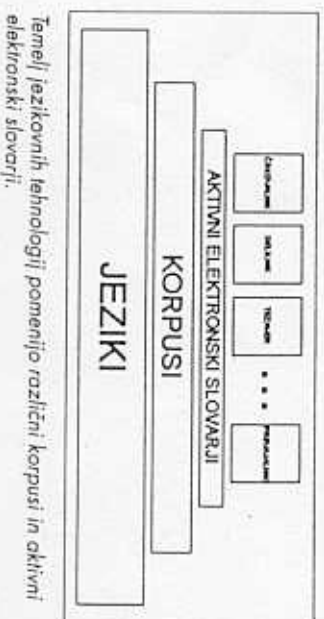
Nistejmo nekaj značilnih primerov, kako lahko jezikovne tehnologije vplivajo na naše vsakdanje življenje. Crkvalnik ali celo slovnčni pregledovalnik sta že danes marsikomu nepogrešljiv pripomoček, brez katerega si je težko zamišljati pisanje besedil. Prevajanje besedil iz tujih jezikov in obratno v veliki meri pokaže ob pomoči katerega od elektronskih slovarjev, programa za pomoč pri prevajanju ali kar neposredno s pomočjo strojnega prevajanja. V zadnjem času je precej govora o pametnih hiši prihodnosti, kjer je ena izmed pomembnejših funkcionalnosti tudi govorna komunikacija z napravami. Murskije so vloga človeškega operaterja prevzeli govorni telefonski odzivniki. Pri iskanju koristnih informacij si že danes lahko pomagamo s katerim od programov, sposobnih komuniciranja v naravnem jeziku; nekateri izmed teh sistemov posajajo celo sesavni del najnovejših modelov avtomobilov. Slepi in slabovidni

lažje uporabljajo marsikatero napravo ali računalniški program, ki je posebej zanje prilagojen z dodatnim govornim vmesnikom (ŽIT 2002/7-8, str. 44). Vsega tega in še marsičesa drugega brez razvoja jezikovnih tehnologij ne bi bilo. Na eni strani nam torej jezikovne tehnologije pomagajo pri naših običajnih opravilih v zvezi z jezikom (npr. pravilno pisanje, hitrejša in lažja prevajanje), po drugi strani pa nam lahko približajo uporabo nekaterih naprav, še zlasti določenih funkcij računalnika, na nam priljubljenejši način (npr. uporaba govora).

Na področju jezikovnih tehnologij deluje veliko število podjetij, organizacij in ustanov, razvijalci in proizvajalci se srečujejo na mnogih mednarodnih konferencah, kakovost številnih rešitev ter sistemov s tega področja je vsaj za večino glavnih svetovnih jezikov na precej visoki ravni. Slovenija oz. slovensčina sicer nekoliko zaostaja za najrazvitejšimi, vendar je hkrati v prednosti pred številnimi narodi oz. jeziki, ki jo po številu »uporabnikov« krepko presenjajo. Razvojne skupine v Sloveniji so združene v Slovenskem društvu za jezikovne tehnologije (SDJT), ki vsa-ki dve leti organizira posebno konferenco na to temo. Člani društva sodelujejo v številnih mednarodnih organizacijah in projektih, tako da je zagotovljen stalni stik z razvojnimi dosežki v tujini.

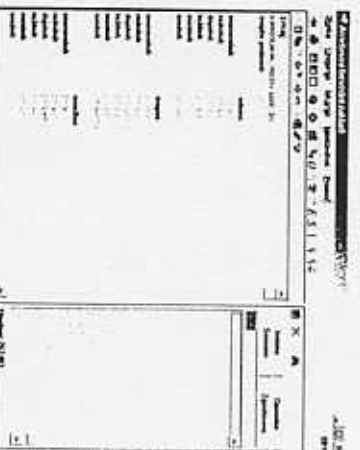
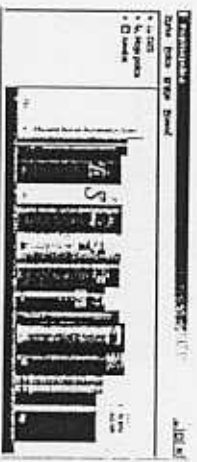
## STRUKTURA JEZIKOVNIH TEHNOLOGIJ

Jezikove tehnologije obsegajo vrsto področij in značilnih (računalniških) aplikacij, njihovo zgradbo, medsebojno povezanost in

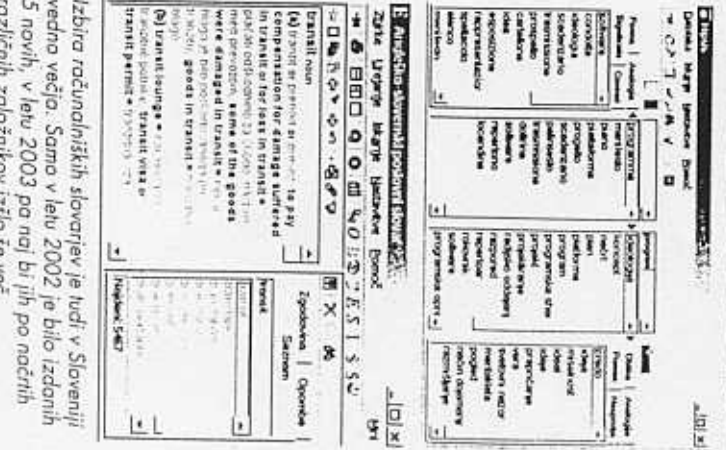
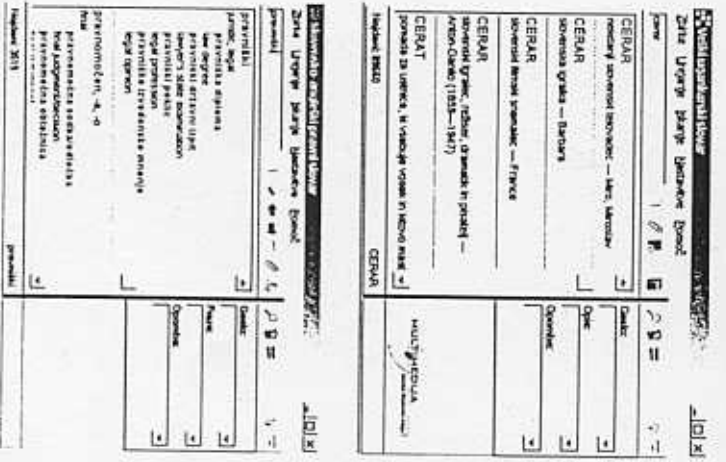


Temelj jezikovnih tehnologij pomenijo različni korpusi in aktivni elektronski slovarji.

odvisnost pa si lahko ponazorimo na različne načine. Vsekakor je poleg razvoja informacijske tehnologije osnovni razlog za obstoj jezikovnih tehnologij obstoj naravnih jezikov samih, sistema pisave ter množica pisa-



nih in govornih besedil, dostopnih v raznih oblikah, predvsem elektronski. Nadgradnja tej neurejeni množici besedil in hkrati temelj jezikovnih tehnologij pomenijo t. i. korpusi, ki prinašajo urejenost ter s tem možnost nadaljnjih avtomatiziranih obdelav. Ker pa samo urejena besedila ne morejo biti sama po sebi osnova za



Izbira računalniških slovarjev je tudi v Sloveniji vedno večja. Samo v letu 2002 je bilo izdanih 5 novih, v letu 2003 pa naj bi jih po načrtih različnih založnikov izšlo še več.

Ueto različnih jezikovnih orodij, je treba korpus nadgraditi s t. i. aktivnimi elektronskimi slovarji, ki vsebujejo dodatne jezikovne informacije, pridobljene bodisi na osnovi analize korpusa, bodisi z ročnim vnosom na osnovi pridobljenega jezikovnega znanja. Z zadostno količino podatkov in znanja o jeziku, vsebovanem v aktivnih elektronskih slovarjih, lahko ti postajo osnova za različna jezikovna področja in njihove aplikacije.

## ELEKTRONSKI SLOVAR

**Število dostopnih domačih slovarjev v elektronski obliki se iz leta v leto povečuje. Njihova kakovost pa izboljšuje. Razvoj slovarskih orodij gre predvsem v smeri enostavnejše in hitrejšje uporabe ter integracije slovarjev z drugimi orodji, kot so npr. prevajalniki. Močan razvoj pa poteka tudi v smeri gradnje t. i. aktivnih elektronskih slovarjev. Ti niso preveč namenjeni ljudem, temveč predvsem programom in aplikacijam, zato je pri njih osnovna zahteva, da so strojno berljivi. Tak slovar je npr. sestavni del prevajalnega sistema in je v osnovi zgrajen drugače kot običajni računalniški slovar.**

Uspešna infrastruktura je ključnega pomena za razvoj jezikovnih tehnologij. To lahko sestavljajo bolj ali manj urejene oz. popolne zbirke besed, besedil in izgovarjaja

ter različni slovarji, leksikoni, tezavri itd. Brez zadostnega števila ustrezno urejenih jezikovnih podatkov si razvoja in izdelave splošno uporabnih programov ter sistemov s področja jezikovnih tehnologij res ne moremo zamisljati. Za gradnjo aktivnih (strojno berljivih) elektronskih slovarjev, ki vsebujejo različne informacije (znanje) o jeziku, so bistvenega pomena besedilni korpusi. To so zbirke besedil, urejene po določenih pravilih in pretvorjene v enotno obliko zapisa, ter usreznih programov, s katerimi preiskujemo besedila in nad njimi izvajamo določene statistične operacije. V zadnjih letih se je težišče jezikoslovnega raziskovalnega dela zelo prevesilo v smeri besedilnih korpusov, saj so ti – še zlasti, če so zelo obsežni in pravilno uravnoteženi –, edini pravi pokazatelj dejanskega stanja jezika. Poleg besedilnih obstajajo tudi korpusi govornih besedil, ki se uporabljajo kot infrastruktura za razvoj razpoznavne in sinteze govora. Ti so manjši od večine besedilnih korpusov, saj je na voljo manj posnetih zvočnih materialov kot samih besedil, njihova obdelava pa je precej zahtevnejša in dolgotrajnejša.

no prilagoditi tudi zbiranje materialov in načrtin gradnje. Za potrebe referenčnega korpusa, ki naj bi predstavljal čim bolj verodostojno sliko stanja jezika, je treba zbrane materiale na podlagi določenih pravil uravnotežiti. Za uravnoteženje se uporabljajo določena razmerja glede lektorske obdelave (lektorično, nelektorično), izvora (dnevnik, tehniški revije, knjige, ...) in vrst besedil (drame, leposlovje, proza, ...) tako, da ti predstavljajo pravo sliko jezika. Če korpus ni referenčni

Besedilni korpus lahko uporabimo v različne namene, med drugim kot koristno orodje pri razvoju slovarjev. Kot primer je prikazan postopek iskanja besed v korpusu FIDA, ki se najpogosteje neposredno vežejo s samostalnikom *znanje*. Te podatke lahko uporabimo pri izdelavi slovarskega gesla kot primera najpogostejše uporabljanih zvez.

### 3. Izpolni statično - ekster - Microsoft Internet Explorer

Št. št.	Indikator	1997
1	prejemanje	354
2	prejemanje	242
3	prejemanje	184
4	prejemanje	101
5	prejemanje	353
6	prejemanje	789
7	prejemanje	155
8	prejemanje	104
9	prejemanje	125
10	prejemanje	121

### 3. korak: Ko dobimo vsa poljavljana besede znanje, uporabimo funkcijo iskanja prave besede na levi.

### 2. korak: Dobljene rezultate uredimo po ustreznosti statistični funkciji. Rezultat je po pogostosti urejen seznam besed, ki nastopajo v slovarščini neposredno pred besedo znanje.

### 3. korak: Ko dobimo vsa poljavljana besede znanje, uporabimo funkcijo iskanja prave besede na levi.

### 2. korak: Dobljene rezultate uredimo po ustreznosti statistični funkciji. Rezultat je po pogostosti urejen seznam besed, ki nastopajo v slovarščini neposredno pred besedo znanje.

### 3. korak: Ko dobimo vsa poljavljana besede znanje, uporabimo funkcijo iskanja prave besede na levi.

### 2. korak: Dobljene rezultate uredimo po ustreznosti statistični funkciji. Rezultat je po pogostosti urejen seznam besed, ki nastopajo v slovarščini neposredno pred besedo znanje.



uravnoveženje ni potrebno. Naslednji korak pri gradnji korpusa je običajno še oblikoslovo označevanje besed, nazadnje pa je treba korpus spraviti v obliko, ki je primerna za hitro iskanje in različne statistične operacije nad njim. Zaradi velikosti korpusov so ti v glavnem nameščeni na strežnikih, uporabljeni pa do njih dostopajo prek spleta.

Večina jezikov že ima zgrajene različne vrste korpusov, ki pa se seveda nenehno dograjujejo. Tudi slovenski jezik ima vrsto različnih korpusov. Ta čas največji uravnovežen referenčni besedilni korpus slovenskega jezika je korpus FTDA, ki vsebuje dobrih 100 milijonov besed, kar je zelo dobra osnova za različne jezikovne raziskave. V sodelovanju številnih domačih partnerjev pa se pripravljata tudi nov projekt za izgradnjo še večjega jezikovnega korpusa, ki bo omogočil še večji napredek na področju jezikoslovja in jezikovnih tehnologij.

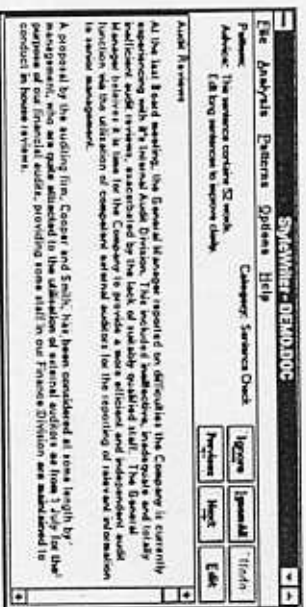
## ČRKOVALNIK

Črkovalnik je najnovejši jezikovni modul in je običajno sestavni del urejevalnikov besedil. Njegova naloga je odkrivanje tipkarskih napak v besedilu, njegovo delovanje pa je v osnovi zelo preprosto. S pomočjo slovarja pravih besednih oblik preveri, ali določena beseda obstaja. Če besede v slovarju ne najde, jo proglasi za napako. Seveda je lahko beseda povsem pravilno zapisana, vendar slovar te oblike ne vsebuje. Zato je pri črkovanju zelo pomembno zadostno število pravih oblik, ki jih vsebuje slovar. Navadno je to nekaj sto tisoč besednih oblik, pri preglebnih jezikih, med katere sodi tudi slovensčina, pa se to število lahko povzpne tudi krepko čez milijon. Ker v njem seveda nikdar ni vseh besed, saj te neprestano nastajajo in se spreminjajo, črkovalnik običajno omogoča tudi dodajanje novih besednih oblik v dodatni slovar. Poleg tega ima lahko še nekaj drugih funkcij, s katerimi pomaga uporabniku pri delu, kot so npr. nasveti pri napakah besedah. Poleg urejevalnikov besedil črkova-

valnik pri svojem delu s pridom uporablja tudi drugi programi, kot so npr. programi za optično prepoznavanje znakov (OCR), ki z njegovo pomočjo skušajo samodejno izločiti čim več napak pri razpoznavanju znakov, zato so lahko tudi precej bolj zanesljivi.

## SLOVNIČNI PREGLEDOVALNIK

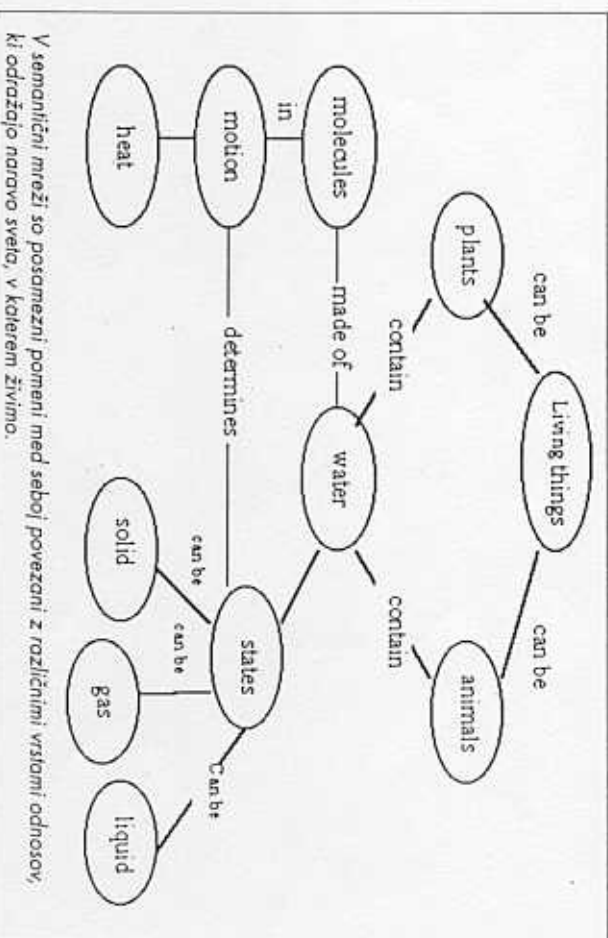
Slovniki pregledovalnik je naslednji jezikovni modul, ki pa v urejevalnikih besedil ni tako pogosto zastopan. V nasprotju s črkovalnikom, ki preverja le očine tipkarske napake, slovniki pregledovalnik preverja sintakso celotnega stavka, zato odkrije tudi slovnične napake, ki jih črkovalnik ne zna. Seveda je izdelava tega modula neprimerno zahtevnejša, predvsem pa jezikovno zelo odvisna. Že samo izdelava ustreznega slovarja je zahtevno in dolgotrajno delo. V vsaki besedi mora namreč vsebovati številne informacije, ki so odvisne od njene besedne vrste. Pri samostalniku npr. potrebujemo informacijo o spolu, vse pripadajoče oblike za različna števila in skloni, svojilni pridevnik (če obstaja), podatek o živosti – in še kaj bi se našlo. Poleg slovarja je treba za skoraj vsak jezik posebej zgraditi tudi sintaktični analizator, ki uporablja prej omenjeni slovar in ki je jedro tega modula. Običajno je naloga slovnicega pregledovalnika tudi svetovanje pri odpravi napake.



Slovniki pregledovalniki so sposobni poleg tipkarskih odkrivati tudi zapletenejša slovnične napake.

## TEZAYER

Trejni jezikovni modul, ki je tudi pogosto del urejevalnikov, je tezaver. Ta poseben slovar



V semantični mreži so posamezni pomeni med seboj povezani z različnimi vrstami odnosov, ki odražajo naravo sveta, v katerem živimo.

med seboj povezuje različne besede oz. pomeni, ki jih te besede predstavljajo. Koristi nam v primerih, ko pri pisanju potrebujemo soroden izraz, sinonim, protipomenko ali referenčno besedo (npr. pridevnik samostalnika, žensko obliko, nadpomen). Tezaver pomeni neke vrste semantično mrežo, v kateri povezuje med pomeni kot osnovnimi elementi teza- vira odseva naravo sveta, y katerem živimo. Informacije, ki so zbrane v tezaveru, lahko pri svojem delu koristno uporabljajo tudi nekateri drugi programi s področja jezikovnih tehnologij. Kot primer lahko vzamemo sistem za strojno prevajanje, ki v primeru, da za kak pomen nima prevoda, lahko uporabi prevod njegovega nadpomena, če ta seveda obstaja.

## DEILNIK

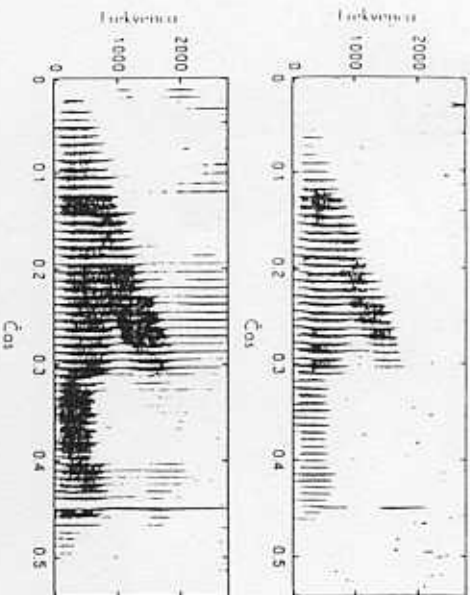
Deilnik za večino jezikov sodi med preprostejšje jezikovne module. V urejevalnikih v primeru, ko je treba čim bolj zapolniti vrstico besedila, poskrbi, da se ustrezna beseda deli na praviem mestu. Običajno deilnik pri svojem delu uporablja poseben jezikovno odvisen algoritem, pri izjemah pa si pomaga z dodatnim slovarjem, ki ga uporabniki lahko tudi sami dopolnjujejo. Deilnik oz. modul za zlogovanje poleg urejevalnikov koristno upo-

rabljajo še drugi programi. Npr. sintetizator govora ga s pridom uporablja pri tvorjenju govora, saj je hitrost izgovarjave med drugim odvisna tudi od števila zlogov v besedi.

## RAZPOZNAVNA GOVORA

Govorne tehnologije so zelo poznano področje jezikovnih tehnologij. V grobem jih razdelimo na razpoznavo in sintezo govora. Vsaka ima svoje razvojne probleme in možnosti, obema pa je cilj vzpostaviti govorno komunikacijo med ljudmi in različnimi napravami – seveda je tam, kjer je to smiselno in uporabno. S pomočjo razvitih modulov že delujejo nekatere splošno znane rešitve, kot so spletni agenti in informacijski portali.

Pretrgovna govora v besedilo je bila na začetkih računalništva le tih želja razvijalcev, danes pa je že povsem običajna funkcija marsikaterega sistema ali naprave. Nekateri boljši sistemi precej dobro razpoznavajo celo neprekinjeni govor. Glavni razlog za hiter razvoj tega področja v zadnjih letih je v zadostni procesorski in pomnilniški moči sodobnih naprav, kar je omogočilo splošno uporabo razpoznavalnikov govora na povsem običajnih računalnikih in s tem njihovo popularizacijo.



Osnova za predstavljanje zvočnega zapisa besede v računalniku so spektrogrami, ki dajejo časovno in frekvenčno sliko govornega signala. Zgornji spektrogram kaže zapis izgovorjene angleške besede *word*, spodnji pa računalniško generirano enako besedo.

Postopek razpoznavanja govora v osnovi temelji na primerjanju kratkih delov govora z vzorci, ki so v posebni podatkovni bazi. Seveda se pred tem po posebnem postopku obdelajo (pretvorijo iz časovnega v frekvenčni prostor), potem pa se s pomočjo različnih modelov izvršita primerjava in sestavljanje manjših enot v večje. Končni cilj je seveda pretvorba govornega signala v besedilo s čim manj napakami. Pri tem so odločilni nekateri dejavniki, ki pomembno vplivajo na kakovost razpoznavanja. Med njimi je npr. spreminjanje šum ali zvočni okolice, ki ga človek namam neznani način uspešno »izloča«, medtem ko imajo umetni sistemi z njim precej težav. Podobno je z menjavo govornca, ki se mu človek v trenutku prilagodi, pa čeprav ga sliši prvič, sistemi za razpoznavo govora pa pri tem postanejo precej manj zanesljivi. Kljub številnim težavam je razvoj omenjenih sistemov v precejšnjem razmahu, zato lahko na tem področju še veliko pričakujemo.

## SINTEZA GOVORA

Sinteza govora je jezikovno precej bolj odvisna od izbranega jezika kot razpoznavanje. Zato je prvi pogoj za uspešno pretvorbo be-

sedila v govor izgradnja ustrezne baze izgovorjav in algoritmov. Pod sintezo govora uvrščamo predvsem tiste rešitve, kjer ni vnapij posnetih celih besed ali celo stavkov, ampak se govor generira splošno, z lepljenjem osnovnih glasov. Nekatere rešitve, ki ta čas prevladujejo, govor tvorijo z lepljenjem sestavljenih glasov, ker so rezultati glede na vložek najbolj optimalni. Seveda obstajajo tudi druge rešitve, ki npr. simulirajo govorni trakt in dajejo zelo dobre končne rezultate. Žal je časovno krmiljenje takih sistemov izjemno zahtevno, razvoj pa izredno dolgotruden ter drag. Vsekakor se sinteza govora tako kot razpoznavanje iz laboratorijev počasi vse bolj seli k uporabnikom, z njeno pomočjo pa številne naprave postajajo vedno bolj prijazne za uporabo, kar je glavni cilj jezikovnih tehnologij.

## STROJNO (PODPRTO) PREVAJANJE

Strojno (podprto) prevajanje postaja eno izmed pomembnejših področij jezikovnih tehnologij, zato mu na tem mestu posvetimo malo več pozornosti. V času vsesplošne globalizacije je za ohranjanje dobljenega jezika obstoj takih sistemov posebej pomemben. To še zlasti velja za tako majhen jezik, kot je slovenščina. Velika večina prevajalnih sistemov je namenjena predvsem podpori prevajanja, kar pomeni, da so le v pomoč prevajalcu pri njegovem delu, še vedno pa je prevajalec tisti, ki besedilo dokončno prevede. Seveda obstoječe hitre napredke strojne in programirane opreme nekateri prevajalni sistemi postajajo tako dobri, še posebej na ozko specializiranih področjih, da lahko govorimo že kar o izključno strojnem prevajanju, saj poseg prevajalca pogosto sploh ni več potreben.

Pod izrazom »strojno podprto prevajanje« razumemo prevajanje, ki ga v prvi vrsti izvaja človek prevajalec ob podpori različnih ra-

čunalniških orodij, med katerimi igrajo posebno vlogo t. i. pomnilniki prevodov. Ti so jezikovno neodvisni, delujejo pa tako, da prevajalcu poskušajo pomagati na podlagi prevodov celih stavkov, ki jih je za kolikor toliko uspešno prevajanje treba prej vnesti v prevajalno bazo. In tih nikakor ni malo. Zato se za splošno prevajanje taka vrsta programa zdi precej neuporabna, se pa toliko bolj uporabljajo na omejenih področjih, kjer se določeni stavki precej ponavljajo (npr. navodila za uporabo, zakonodaja). Pomnilniki prevodov imajo vgrajene posebne mehanizme približnega iskanja, brez katerih bi bili v večini primerov povsem neuporabni. Tako pa lahko prevajalcu tudi v primeru, da konkretnega stavka ni v prevajalni bazi, predlagajo stavek in njegov prevod, ki je po obliki oz. uporabljenih besedah najboljši. Prevajalec mora seveda prevod ustrezno spremeniti, kar pa je običajno še vedno precej hitreje, kot če bi moral stavek prevedeti v celoti.

Strojno prevajanje poteka brez pomoči človeka prevajalca. Taki sistemi so jezikovno odvisni, njihov razvoj pa zahteva izredno veliko znanja in časa. Brez poprejšnjega učenja so sposobni prevesti poljubno besedilo, kar je velika prednost pri njihovi uporabi. V splošnem njihov prevod ni vedno najboljši, zato so pregled in popravki s strani uporabnika še vedno zaželeni.

Prevajalne sisteme s področja strojnega prevajanja bi lahko razvrstili po različnih merilih. Če kot merilo vzamemo »globoino« oz. način dela, na kakršnega vsebinsko prevajajo, bi jih lahko v grobem razdelili v štiri skupine:

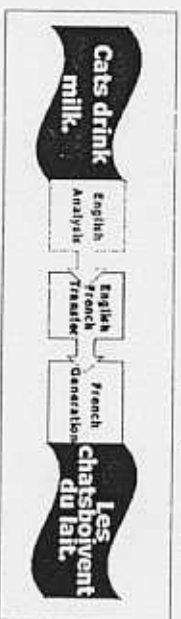
1. Prevajalni sistemi, ki le menjajo besedno obliko za besedno obliko. Taki sistemi so za pregledne jezike, kakšen je tudi slovenščina, v splošnem neuporabni. V poštev bi lahko prišli le pri zelo omejenem področju prevajanja.

2. Prevajalni sistemi, ki znajo vhodno besedilo lematizirati (poiskati osnovno obliko besede). Tak sistem brez kakršne koli dodatne logike bi npr. za slovenščino že dal koli-

kor toliko ubrbe rezultate. Seveda pa v splošnem stavki ne bi bili posebno berljivi, ker prevajanje še vedno poteka le na podlagi zamenjave besede za besedo, pa tudi sintaksa izvornega in ciljnega jezika se ne upošteva.

3. Prevajalni sistemi, ki poleg morfološke upoštevanje tudi sintakso obeh jezikov. Prevod takega sistema je v splošnem že kar lepo berljiv, ker pa sistem nima pomenske analize, prevod vsebinsko ne ustreza vedno najbolj.

4. Prevajalne sisteme, ki prevajajo na podlagi morfološke, sintakčne in pomenske analize. Ti dajejo pri prevajanju najboljše rezultate, a seveda niso neizmisljivi. Da bi pravi »razumeli« stavek, ki ga prevajajo, zaradi boljše analize del informacij prenašajo iz prejšnjih stavkov, saj se pogosto zgodi, da se kak stavek nanaša na prejšnjega ali prejšnje stavke. Osnovna enota prevajanja so tako odstavki. Tudi rezultati prevajanja so običajno bolj pri daljših stavkih, saj tam prevajalni sistem lahko natančneje določi pomen posameznih besed kot pri krajših stavkih.



Strojno prevajalnik najprej analizira stavek v izvornem jeziku, v drugem koraku prevede posamezne dele stavka, na koncu pa oblikuje stavek po pravilih ciljnega jezika.

V splošnem se prevajanje izvaja v treh osnovnih korakih.

Naprej se z analizo poskuša čim bolj natančno ugotoviti, kakšen je pomen posamezne besede v izvornem stavku ter kako so te med seboj pomensko povezane. Včasih – še zlasti pri na videz preprostitih stavkih – se zdi, da analiza ni posebej zahtevna stvar. Vendar ni tako. Analizator mora najti in upoštevati vse mogoče pomenske in sintakčne možnosti, saj nikoli ne moremo natančno vedeti, katere od njih je prava: ali tista, ki je sicer najboljše, ali katere druga. Kot primer vzemimo stavek *To je Miro Cerar*. Na prvi pogled ni nobenega dvoma. Pomen stavka si lahko razlagamo, kot da je to oseba moškega



spozna imenom in primiklo  
Vendar pa si isti stavek lahko  
dveh drugih pomenih, ki sicer  
sintaktično in pomensko pa su-  
na. Stavek lahko razumemo ti-  
to (*nekaj*) je žensko z imeno-  
to (*nekaj*) je to (*nekaj*). Če bi  
Miro Cerar je to (*nekaj*). Če bi  
tični analizator, potem so vsi  
možnosti sintaktično pravilne  
ne. Šele z dodatno pomensk-  
pomaga z določenimi statistič-  
jami, pridobljenimi s pomoč-  
meznega jezika (korpusa), ter  
informacijami, pridobljenimi  
na podlagi prejšnjih stavkov,  
lahko te možnosti učežimo po-  
verjetnosti ter s tem določimo  
najverjetnejši prevod.

Ko tako izluščimo določen pomen posameznih besed ali fraz v izvornem stavku, moramo osnovne oblike teh besed ustrezno prevoriti v besede in fraze ciljnega jezika. Pri tem si prevajalnik običajno pomaga z vgrajenim splošnim slovarjem. Ker ima lahko določena beseda ali fraza več možnih prevodov, je treba vse prenesti na prejšjo tretjega modula, ki se potem odloči za »pravi« prevod. Poleg osnovnega splošnega slovarja ima prevajalnik lahko še dodatne terminološke slovarje, osebne slovarje uporabnika ali vzorce prevodov celih stavkov, katerih naloga je čim bolj in s čim več možnimi besedami prevesti v ciljni jezik gledat na področje, s katerega je prevajano besedilo. Vsako področje ima namreč svoje zakonitosti tvorjenja stavkov in specifične prevode določajo besedi.

Ko je prevajalnik oprem-  
ljen s podatki o pomenih po-  
sameznih besed in njihovih

možnih prevodih, mora opraviti še treto tazo prevajanja, tj. sintezo. Ta ima nalogo, da oblikuje stavke po slovničnih pravilih ciljnega jezika, ki so lahko bistveno drugačna od pravih izvirnega jezika. Poleg tega se mora odločiti za najboljši prevod ustreznega pojma. Pri tem lahko preprosto vrne tisti prevod, ki je v splošnem (na podlagi statistične analize) najpogostejši, lahko pa na podlagi drugih besed v stavku (na podlagi statistične analize kolokacij) določi prevod, ki je v primeru konkretnega stavka verjetno boljši. Vsekakor ima statistična analiza obeh jezikov, tako iz

## Prevalajni sistemi

Od tujih prevajalnih sistemov so pri nas najbolj poznani in uporabljeni izdelki podjetij Tredos in Systron, ki po so za slovensko praktično uporabo in kot pomnilniki prevodov. Najnovejši izdelek na področju strojnega prevajanja slovenske je prevajalni sistem *Presis* podjetja Amehis iz Kamnika. V osnovi je to »pravilni prevajalnik, ki pa ima vgrajen pomnilnik prevodov, tako da je uporaben za različne vrste uporabnikov. Za zdaj je na voljo le slovensko-angliški prevajalni modul, če bo šlo vse po načrtih, pa bo *Presis* že čez nekaj let prevajal tudi v obrnjeni smeri, torej iz angleščine v slovensko. Da bi čim bolj zadoštil potrebam uporabnikov, *Presis* v odvisnosti od kupljenega paketa vsebuje več programov z različnimi uporabniškimi vmesniki, ki pa uporabljajo isto prevajalno jedro. Ker zasnova celoga sistema oz. njegovih programov temelji na strukturi odjema teč-sisteznik, se lahko prevajalno jedro po dogovoru integrirajo v različne aplikacije (spletna, mobilna, mrežna itd.).

1. **Memorandum**  
 2. **Formal Letter**  
 3. **Informal Letter**  
 4. **Notice**  
 5. **Advertisement**  
 6. **Report**  
 7. **Journal Entry**  
 8. **Accountant's Statement**  
 9. **Business Plan**  
 10. **Business Proposal**  
 11. **Business Letter**  
 12. **Business Card**  
 13. **Business Contract**  
 14. **Business Agreement**  
 15. **Business Invoice**  
 16. **Business Receipt**  
 17. **Business Check**  
 18. **Business Stamp**  
 19. **Business Seal**  
 20. **Business Logo**  
 21. **Business Brand**  
 22. **Business Identity**  
 23. **Business Strategy**  
 24. **Business Plan**  
 25. **Business Proposal**  
 26. **Business Letter**  
 27. **Business Card**  
 28. **Business Contract**  
 29. **Business Agreement**  
 30. **Business Invoice**  
 31. **Business Receipt**  
 32. **Business Check**  
 33. **Business Stamp**  
 34. **Business Seal**  
 35. **Business Logo**  
 36. **Business Brand**  
 37. **Business Identity**  
 38. **Business Strategy**  
 39. **Business Plan**  
 40. **Business Proposal**  
 41. **Business Letter**  
 42. **Business Card**  
 43. **Business Contract**  
 44. **Business Agreement**  
 45. **Business Invoice**  
 46. **Business Receipt**  
 47. **Business Check**  
 48. **Business Stamp**  
 49. **Business Seal**  
 50. **Business Logo**  
 51. **Business Brand**  
 52. **Business Identity**  
 53. **Business Strategy**  
 54. **Business Plan**  
 55. **Business Proposal**  
 56. **Business Letter**  
 57. **Business Card**  
 58. **Business Contract**  
 59. **Business Agreement**  
 60. **Business Invoice**  
 61. **Business Receipt**  
 62. **Business Check**  
 63. **Business Stamp**  
 64. **Business Seal**  
 65. **Business Logo**  
 66. **Business Brand**  
 67. **Business Identity**  
 68. **Business Strategy**  
 69. **Business Plan**  
 70. **Business Proposal**  
 71. **Business Letter**  
 72. **Business Card**  
 73. **Business Contract**  
 74. **Business Agreement**  
 75. **Business Invoice**  
 76. **Business Receipt**  
 77. **Business Check**  
 78. **Business Stamp**  
 79. **Business Seal**  
 80. **Business Logo**  
 81. **Business Brand**  
 82. **Business Identity**  
 83. **Business Strategy**  
 84. **Business Plan**  
 85. **Business Proposal**  
 86. **Business Letter**  
 87. **Business Card**  
 88. **Business Contract**  
 89. **Business Agreement**  
 90. **Business Invoice**  
 91. **Business Receipt**  
 92. **Business Check**  
 93. **Business Stamp**  
 94. **Business Seal**  
 95. **Business Logo**  
 96. **Business Brand**  
 97. **Business Identity**  
 98. **Business Strategy**  
 99. **Business Plan**  
 100. **Business Proposal**  
 101. **Business Letter**  
 102. **Business Card**  
 103. **Business Contract**  
 104. **Business Agreement**  
 105. **Business Invoice**  
 106. **Business Receipt**  
 107. **Business Check**  
 108. **Business Stamp**  
 109. **Business Seal**  
 110. **Business Logo**  
 111. **Business Brand**  
 112. **Business Identity**  
 113. **Business Strategy**  
 114. **Business Plan**  
 115. **Business Proposal**  
 116. **Business Letter**  
 117. **Business Card**  
 118. **Business Contract**  
 119. **Business Agreement**  
 120. **Business Invoice**  
 121. **Business Receipt**  
 122. **Business Check**  
 123. **Business Stamp**  
 124. **Business Seal**  
 125. **Business Logo**  
 126. **Business Brand**  
 127. **Business Identity**  
 128. **Business Strategy**  
 129. **Business Plan**  
 130. **Business Proposal**  
 131. **Business Letter**  
 132. **Business Card**  
 133. **Business Contract**  
 134. **Business Agreement**  
 135. **Business Invoice**  
 136. **Business Receipt**  
 137. **Business Check**  
 138. **Business Stamp**  
 139. **Business Seal**  
 140. **Business Logo**  
 141. **Business Brand**  
 142. **Business Identity**  
 143. **Business Strategy**  
 144. **Business Plan**  
 145. **Business Proposal**  
 146. **Business Letter**  
 147. **Business Card**  
 148. **Business Contract**  
 149. **Business Agreement**  
 150. **Business Invoice**  
 151. **Business Receipt**  
 152. **Business Check**  
 153. **Business Stamp**  
 154. **Business Seal**  
 155. **Business Logo**  
 156. **Business Brand**  
 157. **Business Identity**  
 158. **Business Strategy**  
 159. **Business Plan**  
 160. **Business Proposal**  
 161. **Business Letter**  
 162. **Business Card**  
 163. **Business Contract**  
 164. **Business Agreement**  
 165. **Business Invoice**  
 166. **Business Receipt**  
 167. **Business Check**  
 168. **Business Stamp**  
 169. **Business Seal**  
 170. **Business Logo**  
 171. **Business Brand**  
 172. **Business Identity**  
 173. **Business Strategy**  
 174. **Business Plan**  
 175. **Business Proposal**  
 176. **Business Letter**  
 177. **Business Card**  
 178. **Business Contract**  
 179. **Business Agreement**  
 180. **Business Invoice**  
 181. **Business Receipt**  
 182. **Business Check**  
 183. **Business Stamp**  
 184. **Business Seal**  
 185. **Business Logo**  
 186. **Business Brand**  
 187. **Business Identity**  
 188. **Business Strategy**  
 189. **Business Plan**  
 190. **Business Proposal**  
 191. **Business Letter**  
 192. **Business Card**  
 193. **Business Contract**  
 194. **Business Agreement**  
 195. **Business Invoice**  
 196. **Business Receipt**  
 197. **Business Check**  
 198. **Business Stamp**  
 199. **Business Seal**  
 200. **Business Logo**  
 201. **Business Brand**  
 202. **Business Identity**  
 203. **Business Strategy**  
 204. **Business Plan**  
 205. **Business Proposal**  
 206. **Business Letter**  
 207. **Business Card**  
 208. **Business Contract**  
 209. **Business Agreement**  
 210. **Business Invoice**  
 211. **Business Receipt**  
 212. **Business Check**  
 213. **Business Stamp**  
 214. **Business Seal**  
 215. **Business Logo**  
 216. **Business Brand**  
 217. **Business Identity**  
 218. **Business Strategy**  
 219. **Business Plan**  
 220. **Business Proposal**  
 221. **Business Letter**  
 222. **Business Card**  
 223. **Business Contract**  
 224. **Business Agreement**  
 225. **Business Invoice**  
 226. **Business Receipt**  
 227. **Business Check**  
 228. **Business Stamp**  
 229. **Business Seal**  
 230. **Business Logo**  
 231. **Business Brand**  
 232. **Business Identity**  
 233. **Business Strategy**  
 234. **Business Plan**  
 235. **Business Proposal**  
 236. **Business Letter**  
 237. **Business Card**  
 238. **Business Contract**  
 239. **Business Agreement**  
 240. **Business Invoice**  
 241. **Business Receipt**  
 242. **Business Check**  
 243. **Business Stamp**  
 244. **Business Seal**  
 245. **Business Logo**  
 246. **Business Brand**  
 247. **Business Identity**

Presis PRO je najzmogljivejši program prevajalnega sistema Presis.

vojnega kol čimprej, bistveno vloga pri kakovostnem prevajanju. To je v bistvu tisto znanje, ki se pri človeku (prevajalcu) odraža v obliki izkušenj, in na podlagi katerih se odloča tudi sam. Nisorej dovolj le veliki in dobri slovarji prevodov ter algoritmi za stavčno analizo in sintezo; prevajalnik mora vsebovati čim bolj popolno sliko o obeh jezikih, da se lahko kar najbolj približa dobremu prevodu.

Najnovější prevajalni sistemi poskušajo kombinirati prednosti tako »čistih« prevajalnikov kot pomnilnikov prevajalnikov. Tako dajejo uporabniku možnost, da po želji uporablja tiste mehanizme, ki so za njih najugodnejši glede hitrosti in kakovosti prevajanja. Hitrosti prevajanja, ki jih zmorejo sodobni prevajalni sistemi, so zelo odvisne od strojne in programske opreme, na kateri tečejo, navadno pa so sposobni prevesti nekaj deset strani (formata A 4) besedila na minuto. To pa zlasti pri prevajanju večjih količin besedila seveda lahko prinese velikeanske prihranke časa.

Težja, ki jih je treba rešiti pri tako zahtevnem opravilu, kot je prevajanje, je precej. Posopno izboljševanje različnih delov prevajalnega sistema, npr. slovarskega dela, ki že zaradi spreminjanja jezikov verjetno ne bo nikoli končano, pa bo prejel ali slej prinese rezultate, s katerimi bodo zadovoljni tudi najzahtevnejši uporabniki.

## DIALOG V NARAVNEM JEZIKU

Naravna komunikacija človeka z drugimi ljudmi poteka v obliki govora, za pismene tudi prek pisnih sporočil. V prvem in drugem primeru spoznavanje poteka v obliki stavkov, oblikovanih po pravilih določenga jezika. Ta način komunikacije nam omogoča postopno doseganje vsega, kar želimo nekemu sporočiti, enako sprejemati in se na osnovi tako sprejetih informacij učiti. Seveda pa za uspešnim dialogom stoji nekakšna »notranja logika«, ki jo ljudje brez vsakih težav obvladujejo.

Ask Jeeves sodi med najbolj znane informacijske portale na spletu, ki komunicirajo v naravnem jeziku.



mo, zračunalnik pa pomeni zelo veliko oviro. Z njeinim odstranjevanjem se ukvarjajo posebna veja jezikovnih tehnologij – razvoj sistemov dialoga. Napredek v tej smeri bi ljudem precej olajšal delo z računalniki oz. njim podobnimi napravami, saj bi bila z dodatno uporabo govornih tehnologij takšna komunikacija med človekom in računalnikom praktično enaka komunikaciji med ljudmi.

## SKLEP

V tem prispevku načrta področja in rešitve seveda nikakor niso vse, kar bi lahko uvrstili k jezikovnim tehnologijam. Obstaja namreč še vrsta drugih, manj znanih področij, z razvojem katerih skušajo strokovnjaki uporabo sodobnih informacijskih naprav približati ljudem. S hitrim razvojem informacijske tehnologije pa se pred jezikovne tehnologije postavljajo tudi vedno novi in novi izzivi, na katere bo treba najti odgovor. Počakajmo in se pustimo presenetiti. Morda se večina tistega, o čemer danes le sanjamo, kmalu tudi uresniči.

**http://...**  
nljz-sf.sfj (Slovenka družba za izdajanje tehnologije - SOTI)  
www.fica.net (Slovenski referenčni korpus FICA)  
www.amebis.si (posredni kmeti)  
prod.amebis.si (prejeto sistem frakt)  
www.osk.com (informacijski portal del javne)