



RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

# Nova generacija računalne obrade jezika

Damir Čavar

Odjel za lingvistiku u.o., Sveučilište u Zadru

34. skup IT profesionalaca u Splitu 2009



RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- 1 Jezik
- 2 Modeliranje lingvističke jezgre
- 3 Namjena
- 4 Modeli
- 5 Comments



# Što je jezik?

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Govor: percepcija
  - Kontinuirani nediskretni akustički događaji u vremenu
    - Spektrum varijacije energije na frekvencijama od 100–11000 Hz
    - Formanti: koncentracija energije na određenim frekvencijama
    - Prijelazi između šuma i tišine
- Govor: artikulacija
  - Kontinuirani nediskretne promjene u vokalnome traktu
    - Put zraka; položaj jezika, usana; stanje glasnice itd.
  - kao niz kompleksnih motornih instrukcija



# Što je jezik?

RaLi

D. Čavar

Outline

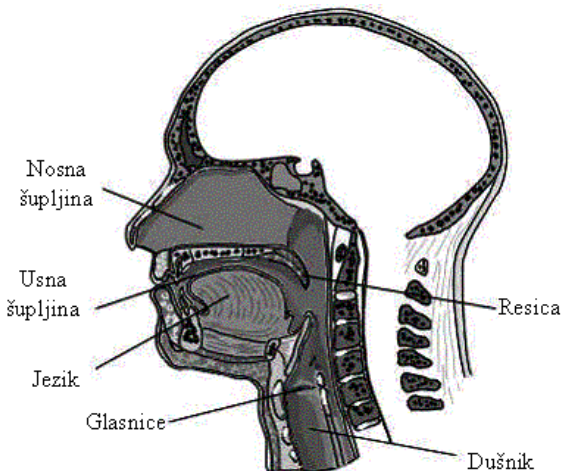
Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments



© Davor Petrinović

<http://dog.zesoi.fer.hr/predavanja/HTML/Osnoveprocesanastajanjagovora.htm>



# Što je jezik?

RaLi

D. Čavar

Outline

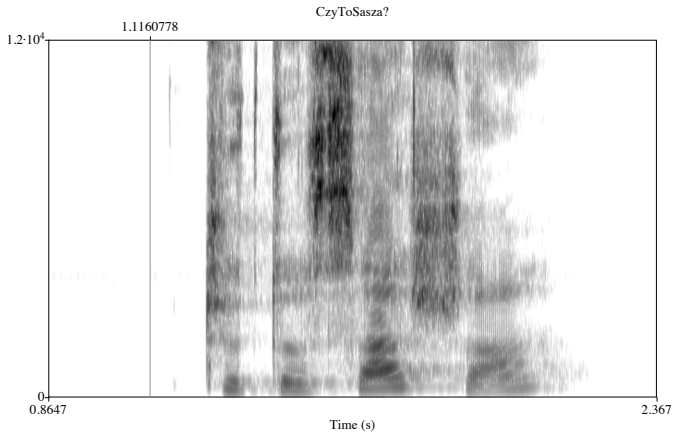
Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments





# Što je jezik?

## Lingvističke osnove

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Razine lingvističke analize (teorije i modeli):
  - Fonetika i fonologija: zvukovi i fonemi
  - Morfologija: morfemi i riječi
  - Sintaksa: rečenica (i možda kontekst)
  - Semantika: značenje rečenice (možda u kontekstu)
  - Pragmatika: govorni čina, itd.
  - itd.
- Iluzija zato što:
  - lingvističke jedinice ne koreliraju nužno s fizičkim aspektima jezika,
  - nego su kognitivne interpretacije akustičkog događaja.



# Što je jezik?

## Lingvističke osnove

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Ekstralingvistička dimenzija:
  - Kognitivni sustav i njegove osobine (npr. *Lazy evaluation*, *Least Effort*, *Last Resort*, pamćenje)
  - Govorna situacija (npr. šum, događaji, biološki uvjeti)
- Lingvistička kognitivna jezgra:
  - Neovisne formalne osobine jezika



# Što je jezik?

## Lingvističke osnove

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Fonološka razina

- Različiti (nediskretni) zvukovi klasificirani kao jedna linkvistička jedinica → fonem
- Osnova: teorija ovisnosti i interdependencije zvukova i njihova kombinatorika
- Primjer: hrvatski i španjolski “r” (*torero* – *onaj koji se bori s bikom*; *torrero* – npr. *stražar u svjetioniku*)

- Fonotaktička razina:

- Hrvatski prihvaća “*dla*” a ne “*lda*” kao slog ili početak riječi





# Što je jezik?

## Lingvističke osnove

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Morfološka razina
  - Fonemi se slažu u morfeme, tj. najmanje jedinice koje imaju značaj ili neku funkciju, i koje se spajaju u riječi
  - Osnova: teorija značenja i funkcija, ovisnosti i interdependencije morfema i njihova kombinatorika
  - Primjer: hrvatski glagol “*čitati*” se može razdvojiti u dva minimalna dijela “*čita-*” i “*-ti*” s posebnim značenjem i funkcijama
  
- Morfotaktička razina:
  - Hrvatski glagoli tipa “*čita*” se mogu kombinirati sa sufiksima kao “*-m*” i “*-š*”, ali ne s “*-om*”, iako je “*-om*” legitiman sufiks hrvatskog jezika (npr. u riječi *ruk-om*)



# Što je jezik?

## Lingvističke osnove

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Sintaktiča razina

- Riječi se slažu u rečenice
- Iako imamo dojam da su skoro sve kombinacije moguće, hrvatski je sintaktički jako ograničen

- Primjer:

- Može biti: *Ivan se penje na krov.*
- Ne može biti: *Krov Ivan se penje na.* ili *Ivan se krov penje na.* itd.

- Dodatni problemi:

- Što znači: *Ivan je nazvao nekog čovjeka iz Pariza.*
- Tko je *on* u: *Ivan **ga** je nazvao.* i *Ivan tvrdi da **ga** je Marija nazvala.*



# Što je jezik?

## Sintaktička stabla i hijerarhijska struktura

RaLi

D. Ćavar

Outline

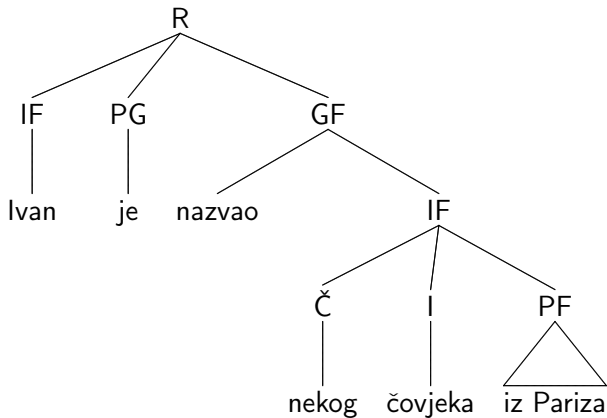
Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments





# Što je jezik?

## Sintaktička stabla i hijerarhijska struktura

RaLi

D. Ćavar

Outline

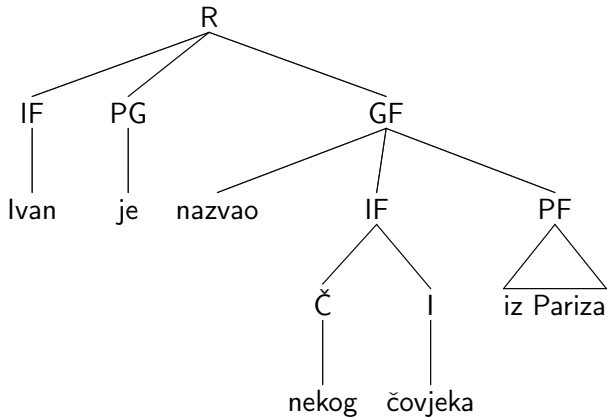
Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments





# Što je jezik?

## Sintaktička stabla i hijerarhijska struktura

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

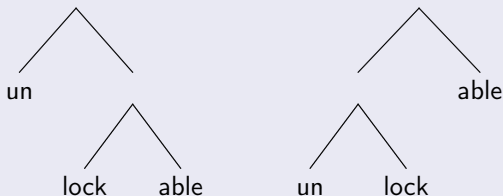
Namjena

Modeli

Comments

- **Strukturalna višeznačnost**

- Za jednu rečenicu ili riječ postoji više struktura u skladu s gramatikom, npr.



- **Leksička višeznačnost:**

- Jedna riječ ima više značenja: npr. *duga, pita, je*



# Lingvistička jezgra

## Formalni aspekti jezika

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Gramatike (elementi i pravila) opisuju moguću kombinatoriku na svim lingvističkim razinama
- Deskriptivne gramatike
  - Opis zvučnih osobina jezika
  - Riječnici
  - Preskriptivne gramatike za standardni jezik
  - Dijalektološke gramatike



# Lingvistička jezgra

## Formalni aspekti jezika

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Formalne gramatike
  - koriste eksplicitnu teoriju i formalizam i omogućavaju falsifikaciju, teoretske predikcije itd.
  - Automati: generatori i prepoznavajući jezičnih izraza (na svim lingvističkim razinama) (niski generativni kapacitet)
  - Parseri: analizatori jezičnih izraza (visoki generativni kapacitet)
- Palatalizacija (poljski): *krok* – *krocze*; *mózg* – *mózdze*; *duch* – *dusze*  
k,g,h → č,dž,š/\_\_\_ i,e [ Deminutiv | Vokativ ]
- Sintaksa:  
S → NP VP  
NP → (Adj) N (PP)  
...



RaLi

D. Čavar

Outline

Jezik

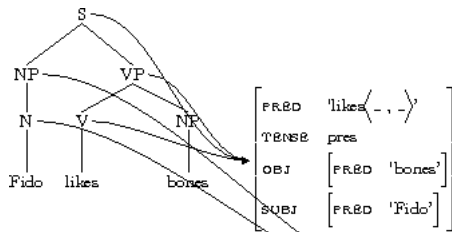
Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Dodatni formalizmi: Unifikacijske gramatike (LFG, HPSG itd.)







# Lingvistički modeli

## Razlike: formalni i prirodni jezici

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Prirodni jezici su višeznačni na svim višim lingvističkim razinama
  - Sve razine su povezane i doprinose riješenju problema višeznačnosti → paralelizam u analizi, procesiranju itd.
  - Gramatike prirodnih jezika su rekurzivne (tj. regularne, kontekstno neovisne i ovisne), što objašnjava neograničen broj izraza, rečenica itd.
- 
- Formalne osobine:
    - Regularna: fonologija i fonotaktika, morfologija
    - Kontekstno neovisna: sintaksa (možda djelomično kontekstno ovisna)
    - Semantika itd.: kontekstno ovisna



# Lingvistički modeli

## Formalne osobine

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Kontekstno neovisne i kompleksnije gramatike (i jezici) se formalno ne mogu usvojiti samo s pozitivnom evidencijom (Gold, 1967)  
iako sada postoje istraživanja koja to relativiraju
- Takve gramatike kompleksne su u procesiranju
- Ne pokrivaju nikada 100% podatke
- Ne predviđaju razlučivanje višeznačnosti



# Lingvistički modeli

## Statistička revolucija (ponovo)

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Generiranje gramatika, riječnika i lingvističkih modela iz obilježenih lingvističkih podataka (npr. korpusa)
- Kontekstno neovisne gramatike s vjerojatnosti pravila  
 $S \rightarrow NP VP$   $p:0.021$   
 $NP \rightarrow (Adj) N (PP)$   $p:0.001$   
...
- Konačni automati s vjerojatnosti na prijelazima (i/ili prijelaznim akcijama kod transduktora)
- $n$ -gram modeli
- Nesimbolički statistički modeli (npr. neuronske mreže)
- ...



# Namjena lingvističkih modela

## Osnovno procesiranje tekstualnih oblika jezika

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Transkripcija u fonetski i/ili fonemski oblik
  - IPA transkripcija
  - za npr. phonex, soundex pretraživanje, statističke fonetske modele za prepoznavanje govora i sintezu itd.
- Morfološka segmentacija i obilježje:
  - *izponapijali*: aspektualni prefiks – aspektualni prefiks – korijen i lema *napiti* – sufiks participa u množini
  - Dodatno obilježje:
    - do neke mijere – malo – “opiti se” od korijenske leme *piti* – prošlost
  - za parsiranje i semantičku analizu



# Namjena lingvističkih modela

## Osnovno procesiranje tekstualnih oblika jezika

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Sintaktičko parsiranje
  - Stabla s kategorijama i hijerarhijskom strukturom skopusa sintaktičkih fraza i riječi
  - za npr. razlučivanje višeznačnosti, semantičku analizu
- Semantičko obilježje i analiza
  - Stabla i mreže relacija i povezivanje s reprezentacijom koncepata i funkcija
  - za npr. strojno prevođenje, prepoznavanje govornog čina, analizu sadržaja itd.



# Konačni automati

RaLi

D. Čavar

Outline

Jezik

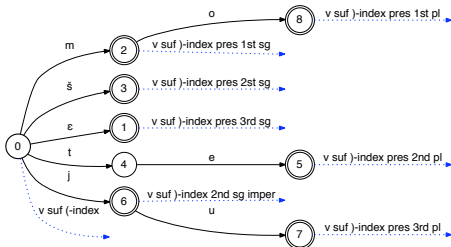
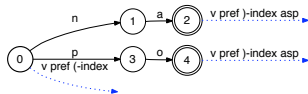
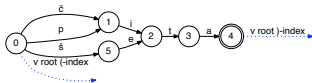
Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

Morfemi kao deterministički konačni automati (DKA) (Mealy ili Moore automati):





# Spajanje u monolitičke automate uz regularne izraze

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Višeznačnost kao emisija više oznaka: lista emisija 1 do  $n$
- Oznaka DKA s imenom
- Pravila koja koriste ta imena i modeliraju morfotaktička distribucijska pravila:

glagolAspektPref\* . glagolAtiKorijeni . glagolFleksSuf



# Generiranje monolitičkih automata

RaLi

D. Čavar

Outline

Jezik

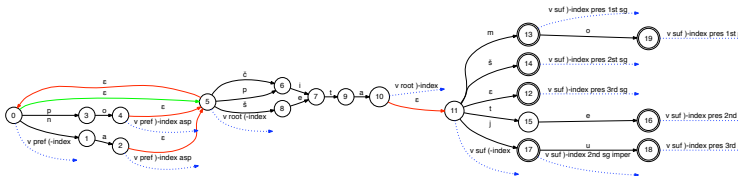
Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

Monolitički automati, mogu biti ciklički DKA:







# Sintaktičko parsiranje

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

Obično u računalnoj lingvistici Earley parser i varijacije tog algoritma (dinamično programiranje):

- s dodatnom vjerojatnosti za razlučivanje najvjerojatnije analize u slučaju višesnačnosti
- s unifikacijom obilježja za pravila kongruencije i perkolaciju oznaka
- s obilježjem semantičkih osobina i funkcija



# Statistički modeli

za npr. obilježje i prepoznavanje

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- n-gram modeli
  - lingvističko obilježje ovisno o (obično lokalnom) kontekstu
  - distribucijske osobine fonema, morfema, riječi u kontekstu otkriva osobine teksta, riječi itd.
  - za npr. klasifikaciju teksta u jezike, zadržajno; klasifikaciju nepoznatih riječi itd.



# Glavni zadatci

RaLi

D. Čavar

Outline

Jezik

Modeliranje  
lingvističke  
jezgre

Namjena

Modeli

Comments

- Stvaranje lingvističkih resursa
- Stvaranje alata za lingvističku analizu
  - lematizacija riječi u tekstu za pretraživače i daljnu analizu
  - gramatike i transfer pravila za strojno prevođenje
  - prepoznavanje jezičnih jedinica i klasifikacija u npr. ime osobe, ime tvrtke, ime produkta, datum i vrijeme, lokacija itd.
  - klasifikacija tekstova
  - analiza govora i procesiranje govornog dijaloga
  - prepoznavanje zadržaja za forenzičku analizu
  - ekstrakcija znanja i generiranje novih saznanja
  - itd.