
Linguistica Computazionale – Lezione 2

Strumenti linguistico-formali

Mercoledì 7 Marzo 2007
Cristiano Chesi, chesi@media.unisi.it

Strumenti linguistico-formali

- Indice
 - Grammatiche formali
 - Nozioni di base (grammatiche a struttura sintagmatica, PSG)
 - Gerarchia di Chomsky
 - Descrizioni Strutturali e derivazioni
 - Formalismi applicabili alla MT
 - Regole di Transfer
 - Interlingua e Ontologie
 - Grammatiche ad unificazione
 - Principi e Parametri

Lecture, approfondimenti

- **Bibliografia essenziale**
 - Hutchins & Somers (1992) Cap. 2
 - Jurafsky & Martin (2000) Cap. 1, 2
- **Approfondimenti**
 - Turcato D. (1993) *Grammatiche formali e linguaggio naturale*. Calderini Bologna
 - Allegranza, V., Mazzini G. (2000) *Linguistica generativa e grammatiche a unificazione*. Paravia scriptorium.
 - Baker M. (2001) *The Atoms of Language*. Basic Books

Competence = grammatica

Grammatiche formali

- **Competence** (data-structure, natura del problema)
 - di che tipo di struttura dati ha bisogno la conoscenza linguistica?
 - una parola può iniziare per *ma...* (*mare*) ma non per *mr...*
 - la *e* di *case* ha un valore diverso da quella di *mare*
 - "le case sono sulla collina" Vs. *"case le collina sono sulla"
 - il gatto morde il cane > sogg: gatto(agente); verbo: morde(azione); ogg: cane(oggetto)
 - ?il tostapane morde il gatto
 - l'espressione "le case" si riferisce ad un gruppo di case evidente dal contesto (Vs. "delle case")
 - ad ogni livello si devono specificare delle primitive elementari:
 - **fonemi** - tratti segmentali e sopra-segmentali
 - **morfemi** - identificazione delle regole combinatorie
 - **parole** - gruppi di morfemi significativi
 - **sintagmi** - gruppi tipizzati di parole che esprimono relazioni
 - **elementi tematici** - paziente, agente...
 - **elementi discorsivi** - convenzioni, relazioni pragmatiche pertinenti...

Livelli di adeguatezza di una grammatica

Grammatiche formali

- **Adeguatezza:** una grammatica deve fornire una descrizione adeguata rispetto alla realtà empirica a cui si riferisce. In particolare si può parlare di adeguatezza a tre livelli:
 - **osservativa:** la lingua definita dalla grammatica coincide con quella che si intende descrivere
 - **descrittiva:** l'analisi grammaticale proposta è in linea con le intuizioni linguistiche dei parlanti fornendo descrizioni strutturali adeguate delle frasi accettabili
 - **esplicativa:** i dispositivi generativi utilizzati soddisfano criteri di plausibilità psicolinguistica e riproducono realmente i meccanismi operanti nell'attività linguistica del parlante. Una grammatica si dice esplicativa quando rende conto anche dell'apprendibilità della lingua.

5

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Come si formalizza una grammatica

(inizio ...)

Grammatiche formali

- **A = Alfabeto**
insieme finito di caratteri (A^* = l'insieme di tutte le stringhe possibili costruite concatenando elementi di A ; ε è l'elemento nullo)
- **V = Vocabolario**
insieme (potenzialmente in)finito di parole, costruite concatenando elementi di A ($V \subseteq A^*$)
- **L = Linguaggio**
insieme (potenzialmente in)finito di frasi, costruite concatenando elementi di V ($L \subseteq V^*$)

6

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Come si formalizza una grammatica

(... continua ...)

Grammatiche formali

- Una **grammatica formale** per il linguaggio L è un insieme di regole che permettono di **generare/riconoscere** tutte e sole le frasi appartenenti a L e (d eventualmente) di assegnare a queste frasi un'adeguata descrizione strutturale.

Una grammatica formale G deve essere:

- **esplicita** (il giudizio di grammaticalità deve essere frutto solo dell'applicazione meccanica delle regole scelte)
- **consistente** (una stessa frase non può risultare allo stesso tempo grammaticale e non grammaticale)

7

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Come si formalizza una grammatica

(... continua ...)

Grammatiche formali

- Una grammatica formale G può essere formalizzata (**grammatica a struttura sintagmatica** o **Phrase Structure Grammar, PSG** Chomsky 1965), come una quadrupla ordinata $\langle V, V_T, \rightarrow, \{S\} \rangle$ dove:
 - V è il **vocabolario** della lingua
 - V_T è un sottoinsieme di V che racchiude tutti e soli gli **elementi terminali** (il complemento di V_T rispetto a V sarà l'insieme di tutti i vocaboli non terminali e sarà definito come V_N)
 - \rightarrow è una relazione binaria, asimmetrica e transitiva definita su V^* , detta **relazione di riscrittura**. Ogni coppia ordinata appartenente alla relazione è chiamata **regola di riscrittura**. Per ogni simbolo $A \in V_N$ $\phi A \psi \rightarrow \phi \tau \psi$ per qualche $\phi, \tau, \psi \in V^*$
 - $\{S\}$ è un sottoinsieme di V_N definito come l'insieme degli assiomi che convenzionalmente contiene il solo simbolo S .

8

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Come si formalizza una grammatica

(... fine)

Grammatiche formali

- Date due stringhe φ e $\psi \in V^*$ si dice che esiste una **φ -derivazione di ψ** se $\varphi \rightarrow^* \psi$.
- Se esiste una φ -derivazione di ψ allora si può anche dire che **φ domina ψ** . Tale relazione è riflessiva e transitiva.
- Una φ -derivazione di ψ si dice **terminata** se:
 - $\psi \in V_T^*$
 - per nessun χ esiste una ψ -derivazione di χ
- Data una grammatica G , una **lingua generata** da G , detta **$L(G)$** , è l'insieme di tutte le stringhe φ per cui esiste una S -derivazione terminata di φ .

9

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Descrizioni strutturali (cioè alberi sintattici)

Grammatiche formali

- Una **descrizione strutturale** è una quintupla **$\langle V, I, D, P, A \rangle$** tale che:
 - **V** è un insieme finito dei **vertici** (es. $v_1, v_2, v_3...$)
 - **I** è l'insieme finito degli **identificatori** (es. $S, DP, VP, la, casa...$)
 - **D** è la relazione di **dominanza**. È un ordine debole (cioè una relazione binaria, riflessiva, antisimmetrica e transitiva) definita su V
 - **P** è la relazione di **precedenza**. È un ordine stretto (cioè una relazione binaria, irreflessiva, antisimmetrica e transitiva) definita su V
 - **A** è la **funzione di assegnazione**; una funzione non suriettiva da V a I

10

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Capacità generativa e relazioni di equivalenza

Grammatiche formali

- La **capacità generativa** denota l'insieme di oggetti generati dalla grammatica; tale capacità è:
 - **debole** se riferita al solo semplice insieme di frasi generabili
 - **forte** se associa a tali frasi l'appropriata descrizione strutturale
- Due grammatiche si dicono **equivalenti** se sono in grado di generare lo stesso insieme di oggetti. Anche qua si può parlare di **equivalenza debole** o **equivalenza forte**

11

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Decidibilità

Grammatiche formali

Un insieme Σ si dice

- **decidibile** (o **ricorsivo**) se per ogni elemento e appartenente all'insieme universo esiste un **procedimento meccanico** M che permette di stabilire in un **numero finito di passi** se $e \in \Sigma$ oppure $e \notin \Sigma$ (la non appartenenza a Σ determina l'appartenenza al complemento di Σ definito come $\bar{\Sigma}$)
- **ricorsivamente enumerabile** quando esiste un procedura che enumera tutti e soli gli elementi di Σ

12

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Inclusioni tra classi di grammatiche

(inizio ...)

Grammatiche formali

- La **gerarchia di Chomsky** (1956, 59) pone in relazione grammatiche di potenza diversa ponendo restrizioni sulla struttura delle regole:
- **Tipo 0**: grammatiche non ristrette (**Turing Equivalent**):
 $\alpha \rightarrow \beta$ ($\alpha \neq \epsilon$) [Augmented Transition Networks]
- **Tipo 1**: grammatiche contestuali (**Context Sensitive**):
 $\alpha A \beta \rightarrow \alpha \gamma \beta$ ($\gamma \neq \epsilon$) [Tree Adjoining Grammars]
- **Tipo 2**: grammatiche non-contestuali (**Context Free**):
 $A \rightarrow \gamma$ [Phrase Structure Grammars]
- **Tipo 3**: grammatiche **regolari**:
 $A \rightarrow xB$ [Finite State Automata]

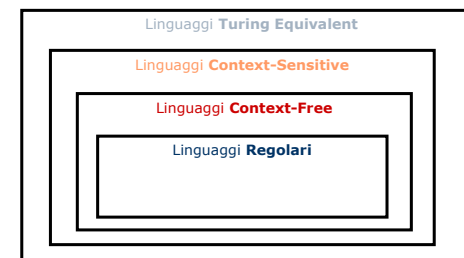
13

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Inclusioni tra classi di grammatiche

(... fine)

Grammatiche formali



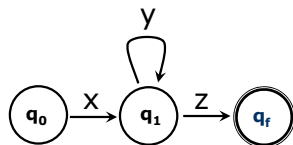
14

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Come si stabilisce l'appartenenza ad una grammatica

Grammatiche formali

- **Pumping lemmas**
Servono per verificare se una proprietà linguistica può essere catturata da una grammatica oppure no
- **Pumping lemma per le grammatiche regolari**
 $a^n b^n$ non è una stringa generabile da nessuna grammatica regolare (poiché nessuna sottostringa può essere "pompata" indefinitivamente garantendo lo stesso numero di a e di b)



15

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Dove stanno le lingue naturali?

(inizio ...)

Grammatiche formali

- Le lingue naturali **non sono generabili da grammatiche regolari** (Chomsky 1956):

If A then B (con A e B potenzialmente anch'esse nella forma "if X then Y"... quindi linguaggi di tipo $a^n b^n$)
- Le lingue naturali **non sono generabili da grammatiche context-free** (Shieber 1985):

Jan säit das mer em Hans es huus hälfed aastriche
(“famoso” dialetto svizzero tedesco)
J. dice che noi a H. la casa abbiamo aiutato a dipingere

Gianni, Luisa e Mario sono rispettivamente sposato, divorziata e scapolo
(“ABC...ABC”... quindi linguaggi di tipo XX)

16

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Dove stanno le lingue naturali?

(... continua ...)

Grammatiche formali

Ricorsività nelle lingue naturali, ovvero come fare un uso infinito di mezzi finiti:

- **Incassamento a destra** (ab^n : iterazione):
[il cane morse [il gatto [che rincorse [il topo [che scappò]]]]]
- **Incassamento centrale** ($a^n b^n$: counting recursion):
[il cane [che il gatto [che il topo che scappò], rincorse], morse]
- **Dipendenze cross-seriali** (xx, identity recursion)
Gianni, Maria e Marco sono rispettivamente sposato, nubile e divorziato

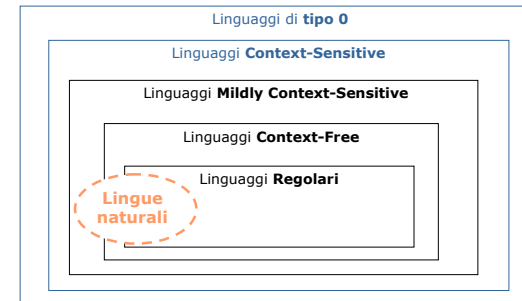
17

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Dove stanno le lingue naturali?

(... fine)

Grammatiche formali



18

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Altri fenomeni linguistici interessanti "catturabili" con CFGs

Grammatiche formali

- **accordo**
per cogliere fenomeni di accordo si deve ricorrere alla duplicazione delle regole di riscrittura. Ad esempio:
 $D_{pl} \rightarrow i$, $N_{pl} \rightarrow \text{cani}$, $D_{sg} \rightarrow \text{il}$, $N_{sg} \rightarrow \text{cane}$, $DP \rightarrow (D_{sg} N_{sg} \mid D_{pl} N_{pl})$
- **sottocategorizzazione**
ci si riferisce allo schema di sottocategorizzazione come alla possibilità di distinguere ulteriormente, all'interno di categorie maggiori (ad esempio la categoria verbale), categorie più precise che rendano conto in modo più adeguato del comportamento dei diversi elementi lessicali: transitivo, intransitivo, inaccusativo o ergativo sono solo tre delle sotto-classi verbali che servono a predire un determinato comportamento verbale (certi approcci, Levin 83, distinguono fino a 183 classi verbali distinte).
esempi di regole:
 $VP \rightarrow (V_{transitivo} DP \mid V_{intransitivo} \mid V_{frasale} S \mid V_{modale} V_{inf...})$

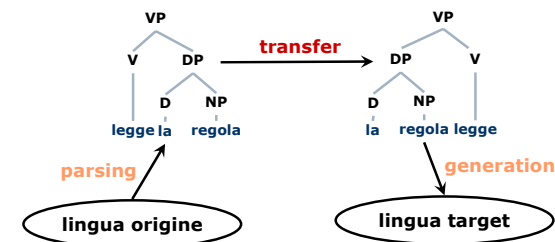
19

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Modello del Transfer

Formalismi applicabili alla MT

- **Conoscenza contrastiva**
esplicitare le differenze tra le due lingue è il primo passo verso la traduzione.
Da questo punto di vista occorre una ristrutturazione linguistica per conformarsi alle regole della lingua target



20

Linguistica Computazionale A.A. 2006-07 - C. Chesì

Esempi di Transfer con Context-Free Grammars

Formalismi applicabili alla MT

- Inglese ⇒ Italiano
DP → D₁ Agg₂ Nome₃ ⇒ DP → D₁ Nome₃ Agg₂

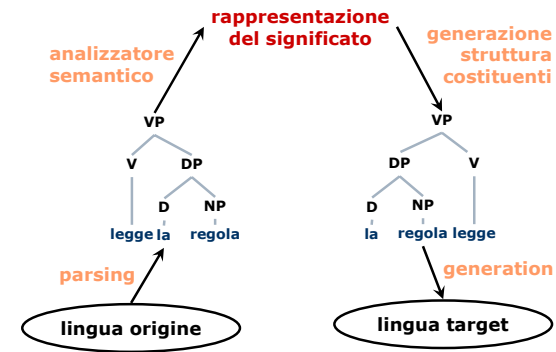
- Giapponese ⇒ Inglese
DP → relativa₁ DP₂ ⇒ DP → DP₂ relativa₁

21

Linguistica Computazionale A.A. 2006-07 – C. Chesì

Modello dell'interlingua

Formalismi applicabili alla MT



22

Linguistica Computazionale A.A. 2006-07 – C. Chesì

Ontologie

(inizio ...)

Formalismi applicabili alla MT

- **Concettualizzazione**
astrazione (e semplificazione) delle relazioni tra oggetti e concetti in un dominio di conoscenza

 - **Ontologia**
specificazione dettagliata di queste entità e relazioni
- Ogni ontologia definisce un insieme di **classi**, **relazioni**, **funzioni** e **oggetti** costanti all'interno di un dominio discorsivo, esplicitando un'**assiomatizzazione** in modo da **vincolare interpretazioni** ed **inferenze**

23

Linguistica Computazionale A.A. 2006-07 – C. Chesì

Ontologie

(... fine)

Formalismi applicabili alla MT

- **Knowledge Interchange Format**
(**KIF**, Genesereth & Fikes, 1992)
- ```
(defrelation PHYSICAL-QUANTITY
 (<=> (PHYSICAL-QUANTITY ?q)
 (and (defined (quantity.magnitude ?q))
 (double-float (quantity.magnitude ?q))
 (defined (quantity.unit ?q))
 (member (quantity.unit ?q)
 (setof meter second kilogram ampere kelvin
 mole candela))))
```

24

Linguistica Computazionale A.A. 2006-07 – C. Chesì

## Rappresentazione del significato

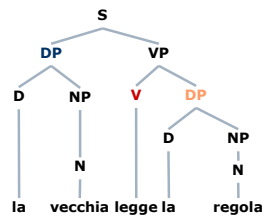
### Formalismi applicabili alla MT

#### ■ Aggiunta semantica alle regole Context Free:

$A \rightarrow \alpha_1 \alpha_2 \dots \alpha_n \quad \{f(\alpha_j, \text{sem} \dots \alpha_k, \text{sem})\}$

vecchia {vecchia}

legge  $\{\exists e, x, y \text{ Isa}(e, leggere) \wedge legge(e, x) \wedge letta(e, y)\}$



$\exists e \text{ Isa}(e, leggere) \wedge legge(e, \text{la vecchia}) \wedge letta(e, \text{la regola})$

25

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(inizio ...)

### Formalismi applicabili alla MT

#### ■ Grammatiche basate su restrizioni

rappresentazione più efficiente e significativa dell'informazione linguistica

#### ■ formalismi leggermente più "potenti" delle CFG, con cui render conto in modo

□ **compatto** (quindi più elegante) ed

□ **efficiente**

delle restrizioni linguistiche imposte da fenomeni produttivi quali quelli precedentemente elencati

#### ■ gerarchie di tratti

proprietà aggiuntive rispetto alle regole di riscrittura

26

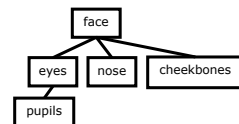
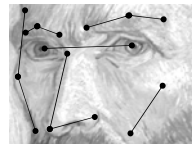
Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(... continua ...)

### Formalismi applicabili alla MT

#### ■ (ma... cosa sono i tratti?)



27

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(... continua ...)

### Formalismi applicabili alla MT

#### ■ (ma... cosa sono i tratti?)

a.



la casa  
è bella

b.



la casa  
è brutta

c.\*



\*casa la  
bella è

d.\*



\*lo casa  
è belli

e.\*



\*casa è

28

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(... continua ...)

### Formalismi applicabili alla MT

- **Struttura di tratti (FS, Feature Structures)**  
è un insieme di coppie del tipo **tratto > valore**  
(es. numero > singolare)

- due tipologie di formalizzazione (equivalenti) delle coppie **tratto > valore**:

### Matrice di Attributi e Valori (AVM, Attribute-Value Matrix)

```
Num = Sing
Gen = Femm
...
Tratton = Valoren
```

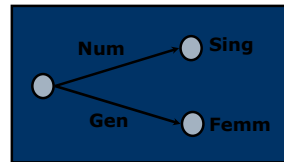


diagramma ad archi orientati ed etichettati (DAG, Direct Acyclic Graph)

29

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(... continua ...)

### Formalismi applicabili alla MT

Alcune proprietà interessanti delle strutture di tratti:

- **parzialità, maggiore o minore specificità**, ovvero alcuni elementi possono restare non specificati, ad esempio il genere:

```
N [Num = Sing]
 [Gen = []]
```

- la struttura delle AVM può essere **rientrante**, cioè un tratto che ha una qualche significatività dal punto di vista empirico, può essere definito da più sottotratti, come nel caso dell'accordo:

```
[Cat = N]
[Accordo [Num = Sing]
 [Gen = Femm]]
```

30

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(... continua ...)

### Formalismi applicabili alla MT

Altre proprietà interessanti delle strutture di tratti:

- **percorsi**, il valore di un tratto è definito in base ad un percorso univoco che lo identifica, cioè una lista di tratti lungo la struttura del tipo: **accordo>num>sing**
- **condivisione di tipo (type sharing)**, una struttura può essere condivisa tra più elementi anche se i valori non lo sono
- **condivisione di occorrenza (token sharing)**, l'occorrenza di un determinato valore può essere condivisa in tal caso può essere indicato con l'uso di un indice, es [1]:

```
[Cat = S]
[testa [Accordo [1] [Num = Sing]
 [Soggetto [Accordo [1]]]]]
```

31

Linguistica Computazionale A.A. 2006-07 - C. Chesì

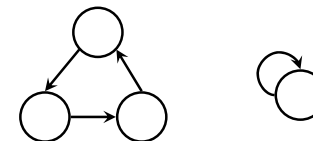
## Grammatiche ad unificazione

(... continua ...)

### Formalismi applicabili alla MT

Alcune proprietà interessanti delle strutture di tratti:

- **significatività empirica**, i tratti sembrano catturare adeguatamente, almeno a livello descrittivo, fenomeni linguisticamente produttivi
- **aciclicità**, i grafi non possono essere ricorsivi



32

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Grammatiche ad unificazione

(... fine)

### Grammatiche ad Unificazione

#### Sussunzione

stabilisce una relazione ordinata tra due strutture di tratti FS; la FS più **generica** sussume quella **specificata**. Si può quindi dire che:

$$FS_a \sqsubseteq FS_b$$

sse  $FS_b$  ha tutti i tratti di  $FS_a$  nella stessa configurazione strutturale e con uguali assegnazioni di valore

#### Unificazione

permette di combinare le informazioni per rappresentarle in formato più compatto e significativo:

$FS_a \sqcup FS_b = FS_x$  (se esiste) tale che  $FS_x$  è la più generale delle FS sussunte da  $FS_a$  e  $FS_b$

33

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Da regole a principi e parametri

(inizio ...)

### Formalismi applicabili alla MT

- **regole**  
specifiche e valide  
per una sola lingua
- **principi & parametri**  
universali linguistici +  
settaggio parametri di variazione
- Ricerca di una migliore **adeguatezza esplicativa** oltre che **descrittiva**
- **Obiettivo**: cogliere gli universali linguistici descrivendo precisamente la limitata variabilità sintattica
- I **principle-based parsers** (Barton 1984, Berwick e Fong 90, Stabler 92) si ispirano a queste idee:
  - i principi della grammatica sono **assiomi** per il parser
  - il parser è un **sistema deduttivo** che inferisce le espressioni grammaticali e la loro struttura partendo da tali assiomi

34

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Da regole a principi e parametri

(... fine)

### Formalismi applicabili alla MT

#### regole

regola passivo → frase passiva  
regola dativo → frase dativa  
regola di focalizzazione → frase focalizzata  
...

#### principi & parametri

P1 → frase passiva  
P2 → frase dativa  
P3 → frase focalizzata

- potenzialmente una decina di principi + pochi parametri possono generare migliaia di regole

35

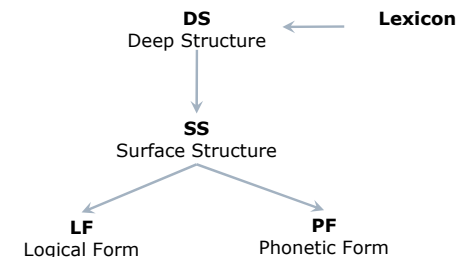
Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Principi e parametri

(inizio ...)

### Formalismi applicabili alla MT

#### Modello a "T"



36

Linguistica Computazionale A.A. 2006-07 - C. Chesì

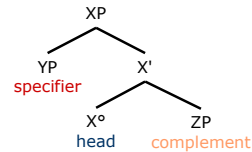
## Principi e parametri

(... continua ...)

### Formalismi applicabili alla MT

#### ■ Alcuni principi

##### □ X' theory



##### □ $\theta$ - criterion

ogni argomento deve ricevere uno ed un solo ruolo tematico (e ogni ruolo tematico è assegnato ad uno ed un solo argomento)

##### □ Case filter

ogni NP lessicale deve ricevere un caso (P e  $V_{\text{finito}}$  assegnano caso)

37

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Principi e parametri

(... continua ...)

### Formalismi applicabili alla MT

#### ■ Altri principi

##### □ Move $\alpha$

una categoria può muoversi in qualsiasi momento, ovunque

##### □ Free indexation

indici sono liberamente assegnati alle categorie in posizione A(rgomentale)

##### □ Binding theory

**condizione A** - Un'anafora (es. *se stessa*) è legata nel suo dominio di legamento (binding domain)

**condizione B** - Un pronome (es. *lei*) è libero nel suo dominio di legamento

**condizione C** - Un'espressione referenziale (es. *Maria*) è sempre libera

38

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Principi e parametri

(... fine)

### Formalismi applicabili alla MT

#### ■ Generatori

principi che producono più strutture di quante non ne ricevano in input:

##### □ Move $\alpha$

##### □ Free indexation

##### □ ...

#### ■ Filtri

principi che selezionano solo parte delle strutture che ricevono in input:

##### □ X' theory

##### □ $\theta$ - criterion

##### □ Case filter

39

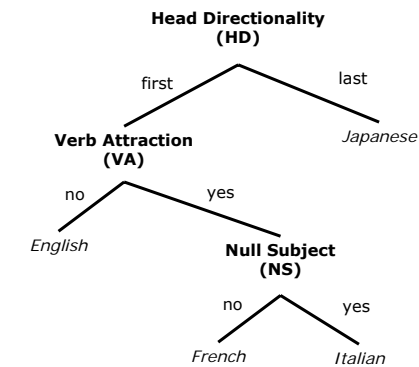
Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Principi e parametri

(... fine)

### Formalismi applicabili alla MT

#### ■ Parametri (Baker 2002)



40

Linguistica Computazionale A.A. 2006-07 - C. Chesì

## Prossima lezione

---

(Domani, Giovedì 8 Marzo, ore 16-18, Aula 456, Palazzo S. Niccolò)

### Indice

#### ■ Fondamenti

- Macchine di Turing (universali)
- Concetto di computazione e computabilità
- Dati, programmi, input e output

#### ■ Basi Dati

- Corpora
- database
- strumenti di interrogazione basi dati (esp. regolari, SQL)

#### ■ Algoritmi

- Cicli ed oggetti
- Ideazione, descrizione, formalizzazione ed implementazione di un algoritmo