

**PLN**

**Introducción.  
Técnicas clásicas**

1. CONTEXTO HISTÓRICO
2. CADENAS DE PROCESAMIENTO
3. PREPROCESAMIENTO DE TEXTOS
4. TOKENIZACIÓN
5. SEGMENTACIÓN DE FRASES
6. ANÁLISIS MORFOLÓGICO
7. ANÁLISIS SINTÁCTICO
8. ANÁLISIS SEMÁNTICO

# 1. CONTEXTO HISTÓRICO

- 50's Machine Translation en universidades de U.K. y U.S.A.
- 1963 1<sup>a</sup> reunión Association for Computational Linguistics
- Técnicas simbólicas
- 80's Aproximaciones estadísticas o empíricas: basadas en corpus
- Revolución estadística: aproximaciones simbólicas vs. estadísticas
- Técnicas híbridas: incorporación de conocimiento lingüístico en procesamientos estadísticos
- Evolución histórica: gramáticas sobreviven vs métodos estadísticos

## 2. CADENAS DE PROCESAMIENTO CLÁSICAS

### **Teoría:**

- A nivel de frase:
  - Sintáxis (estructura)
  - Semántica (significado)
- A nivel de discurso
  - pragmática (discurso, contexto)

### **Práctica:**

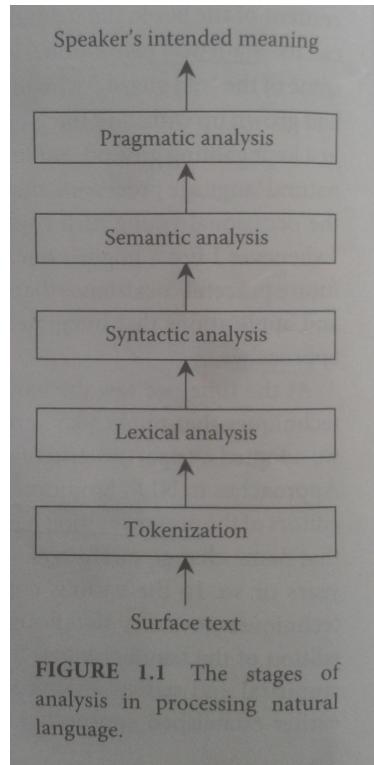
- No es tan fácil separar
- Cuanto más profundizamos en la cadena, más difícil es el procesamiento, y peores resultados. cada vez más abstracto, menos lingüístico

## 2. CADENAS DE PROCESAMIENTO CLÁSICAS

### **Análisis del lenguaje Natural vs. Generación del Lenguaje Natural:**

- Mucho más trabajo publicado en análisis que en generación. Análisis empieza por algo concreto (palabras), pero por dónde empezar con Generación? con ideas?
- NLG: para cuando enlatar respuestas no es suficiente

## 2. CADENAS DE PROCESAMIENTO CLÁSICAS



**FIGURE 1.1** The stages of analysis in processing natural language.

# 3. PREPROCESAMIENTO DE TEXTOS

## - Document triage

- Construcción de corpus limpios
- Internet Download
- Identificación automática de encoding de caracteres: 7 bits=128 (ASCII), 8 bits=256 (ISO-8859), 16bits=65.536 (chino: Big-5, GB), 1-4 bytes (Unicode 5.0, UTF-8)
- Detección de idioma:
  - Fácil: caracteres que solo hay en un cierto idioma, distancia vectorial de caracteres + frecuentes
  - Fácil? documentos multilingües, fragmentos cortos
- Filtrado de textos: suprimimos html, imágenes, publicidad y otras basuras

## - Segmentación de textos

- Tokenización: segmentación de palabras
- Análisis léxico o morfológico

## 4. TOKENIZACIÓN

- Tokenización: segmentación de palabras
  - Fácil
    - Idiomas delimitados por espacios
  - Fácil?:
    - Idiomas no segmentados: chino, japonés, tailandés..
    - textos reales: que es una palabra
    - sistemas de escritura: logográfico, silábico, alfabetico
    - palabras compuestas
    - ambigüedad
    - morfologías: aislante, aglutinante, inflexional, polisintética

# 5. SEGMENTACIÓN DE FRASES

- Detección de límites de frases (puntos...)
- Fácil:
  - Idiomas con puntuación
- Fácil?
  - Qué es una frase donde empieza y acaba
  - Textos de internet
  - Idiomas sin puntuación

# 6. ANÁLISIS LÉXICO O MORFOLÓGICO

- Descomposición de tokens en partes:
- Normalización de textos: formas canónicas (1,uno)
- Lematizar
- Asignar etiquetas POS
- Diccionarios vs. reglas de combinación
- Palabras desconocidas (word guessing)
- Identificación de abreviaturas (Sr. U.S.A.) importante para siguiente fase de segmentación de frases
- Ambigüedad
- Lematización<->generación morfológica
  - usos: MT, IR, stemming, nuevas palabras, correctores, POS tagging, tokenizacion en idiomas sin puntuación
- palabra = morfemas = lexema+morfemas derivativos (pre/in/su - fijo)

# 6. ANÁLISIS LÉXICO O MORFOLÓGICO

- **Transductores de Estados Finitos (FSTs)**

- son un tipo de Maquina/Automata de Estados Finitos (FSM) con 2 cintas: entrada y salida
- facil: raíz+afijos
- facil? idealizado

- **3 Modelos de análisis léxico:**

- I&A (Item and Arrangement) Finite State Morphology
- I&P (Item and Process) Finite State Morphonology
- W&P (Word and Paradigm)

# 6. ANÁLISIS LÉXICO O MORFOLÓGICO

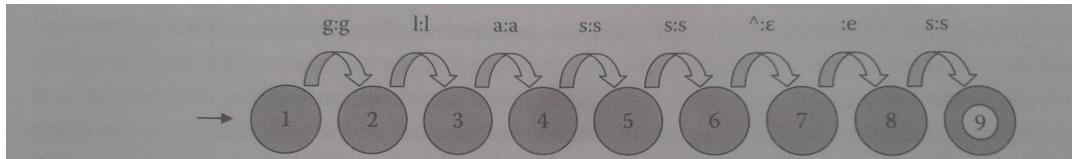


FIGURE 3.1 A spelling rule FST for *glasses*.

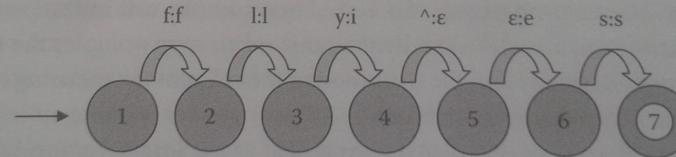


FIGURE 3.2 A spelling rule FST for *flies*.

# 6. ANÁLISIS LÉXICO O MORFOLÓGICO

```
(3.5)    Es_Spelling_Rule:  
        <$abc> == $abc <>  
        <$sx ^> == $sx e <>  
        <^> == <>  
        <> == .  
  
<g l a s s #> = glass.  
<g l a s s ^ s #> = glasses.  
<f o x #> = fox.  
<f o x ^ s> = foxes.  
<c a t #> = cat.  
<c a t ^ s #> = cats.
```

# 7. ANÁLISIS SINTÁCTICO

- Árboles sintácticos. Representación: árboles o paréntesis
- Tipos de analizadores:
  - Recognizer vs. parser vs robust parser
  - top-down vs bottom-up
  - determinístico vs no determinístico
  - de izquierda a derecha vs. otros órdenes
- Estructura de frases.
- Basado en teorías de Lingüística Generativa.
- Facilita el siguiente paso: semántica.
- Parsers NLP != Parsers de programas informáticos:
  - Distinta capacidad generativa. En NLP mejor SGC > WGC
  - ambigüedad del lenguaje natural
  - Ruido: errores y nuevas palabras

## 7. ANÁLISIS SINTÁCTICO

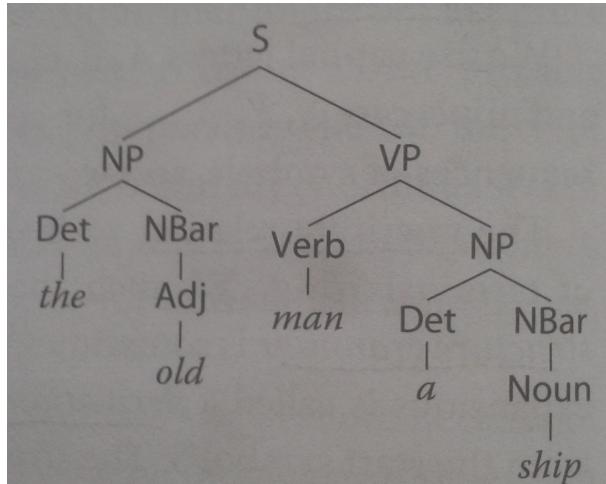


FIGURE 4.2 Syntax tree of the sentence “the old man a ship.”

```
[S [NP [Det the] [NBar [Adj old] ] ] [VP [Verb man] [NP [Det a] [NBar [Noun ship] ] ] ] ]]
```

## 7. ANÁLISIS SINTÁCTICO

Who did you sell the car to \_\_?

Who do you think that you sold the car to \_\_?

Who do you think that he suspects that you sold the car to \_\_?

Put the block in the box on the table

(4.4)

Assuming that "put" subcategorizes for two objects, there are two possible analyses of (4.4):

Put the block [in the box on the table]

(4.5)

Put [the block in the box] on the table

(4.6)

# 7. ANÁLISIS SINTÁCTICO

## - Gramáticas Libres de Contexto (CFGs)

- $G=\{T,N,S,R\}$
- Formas normales:
  - FN de Chomsky (CNF) -> algoritmo CKY
    - A->w
    - A->BC
  - CNF relajada I
    - A->w
    - A->BC
    - A->B
  - CNF relajada II
    - A->w
    - A->BC...
  - Reglas no vacías

# 7. ANÁLISIS SINTÁCTICO

$S \rightarrow NP \quad VP$	$Det \rightarrow a \mid an \mid the$
$NP \rightarrow Det \quad NBar$	$Adj \rightarrow old$
$NBar \rightarrow Adj \quad Noun$	$Noun \rightarrow man \mid men \mid ship \mid ships$
$NBar \rightarrow Noun$	$Verb \rightarrow man \mid mans$
$NBar \rightarrow Adj$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb \quad NP$	

FIGURE 4.1 Example grammar.

# 7. ANÁLISIS SINTÁCTICO

		1	2	3	4	5
		Det	NP	NP, S		S
		<i>the</i>	Adj, NBar	NBar		
0						
1						
2		<i>old</i>		Noun, Verb, NBar, VP		VP
3			<i>man</i>		Det	NP
4				<i>a</i>		Noun, NBar
						<i>ship</i>

FIGURE 4.3 CKY matrix after parsing the sentence “the old man a ship.”

# 7. ANÁLISIS SINTÁCTICO

## - Gramáticas basadas en Restricciones o de Unificación y Rasgos

$$X_0 \rightarrow X_1 X_2$$

$$\langle X_0 \text{ category} \rangle = \text{NP}$$

$$\langle X_1 \text{ category} \rangle = \text{Det}$$

$$\langle X_2 \text{ category} \rangle = \text{Noun}$$

$$\langle X_1 \text{ agreement} \rangle = \langle X_2 \text{ agreement} \rangle$$

Any such rule description can be represented as a phrase-structure rule where the symbols consist of feature terms. Below is a feature term rule corresponding to the previous rule (where  $\boxed{1}$  indicates identity between the associated elements):

$$[\text{category : NP}] \rightarrow \begin{bmatrix} \text{category : Det} \\ \text{agreement : } \boxed{1} \end{bmatrix} \begin{bmatrix} \text{category : Noun} \\ \text{agreement : } \boxed{1} \end{bmatrix}$$

The basic operation on feature terms is unification, which determines if two terms are compatible by merging them to the most general term compatible with both. As an example, the unification  $A \sqcup B$  of the terms  $A = [\text{agreement : [number : plural]}]$  and  $B = [\text{agreement : [gender : neutr]}]$  succeeds with the result:

$$A \sqcup B = \left[ \text{agreement : } \begin{bmatrix} \text{gender : neutr} \\ \text{number : plural} \end{bmatrix} \right]$$

However, neither  $A$  nor  $A \sqcup B$  can be unified with

$$C = [\text{agreement : [number : singular]}]$$

# 7. ANÁLISIS SINTÁCTICO

## - Características de un buen parser:

- Robustez
- Poder de desambiguación
  - undergeneration vs. overgeneration
  - gramáticas especializadas
  - ranking estadístico
- Eficiencia
  - Eficiencia teórica > Eficiencia práctica

# 8. ANÁLISIS SEMÁNTICO

- Menos estudiado que los anteriores pasos de la cadena
- Menos consenso
- Entender
- Metalenguaje semántico
- aplicaciones: IR, IE, Resumen textos, data mining, MT, user queries, ontologías, KR
- Semántica léxica vs. Semántica supra léxica
- 3 ambigüedades semánticas:
  - Léxica (homonimia, polisemia)
  - De alcance
  - Referencial
- Teorías semánticas:
  - Formales vs. cognitivas
  - compositivas vs léxicas (descomposicionales y relaciones)

# 8. ANÁLISIS SEMÁNTICO

## - Lógica. Lógica de Predicados.

- variables, términos, predicados, conectivas, cuantificadores existenciales

- (1) a. *Some politicians are mortal.*  
b.  $\exists x (\text{politician}(x) \wedge \text{mortal}(x))$   
[There is an x (at least one) so that x is a politician and x is mortal.]
- (2) a. *All Australian students like Kevin Rudd.*  
b.  $\forall x ((\text{student}(x) \wedge \text{Australian}(x)) \rightarrow \text{like}(x, k))$   
[For all x with x being a student and Australian, x likes Kevin Rudd.]

## 8. ANÁLISIS SEMÁNTICO

(3) a. Modus ponens:

- (i) P (premise)
- (ii)  $P \rightarrow Q$  (premise)
- (iii) Q (conclusion)

b. (i) *Conrad is tired* (P: **tired(c)**)

(ii) *Whenever Conrad is tired, he sleeps* (P: **tired(c)**, Q: **sleep(c)**,  $P \rightarrow Q$ )

(iii) *Conrad sleeps* (Q: **sleep(c)**)

# 8. ANÁLISIS SEMÁNTICO

## - Teoría de Representación del discurso (DRT)

- Estructuras de Representación del Discurso (DRS)
  - dinámicas
  - recursivas
- DRS = referentes + condiciones
- Ventajas: capaz de soportar frases nominales indefinidas, cuantificación y resolución de anáfora

## 8. ANÁLISIS SEMÁNTICO

(4) *A man sleeps. He snores.*

x, y
man (x)
sleep (x)
y = x
snore (y)

Recursiveness is an important feature. DRSSs can handle recursive structures. An example is given in (5). Notice that according to the DRSS, the subject of the sentence is a singular noun phrase, even though on the face of it *every man* is a singular noun phrase.

(5) *Every man sleeps. He snores.*

y
x
man (x)
⇒
sleep (x)
y = ?
snore (y)

# 8. ANÁLISIS SEMÁNTICO

## - Teoría del Lexicón Generativo de Pustejovsky

- Dinámica
- Basada en lexical items
- Representación en 4 niveles:
  - Argumentos
  - Eventos
  - Qualia
  - Herencia Léxica
- No orientado a cuantificación, anáfora ni presuposición sino descomposicional

# 8. ANÁLISIS SEMÁNTICO

build	
EVENTSTR =	$\left[ \begin{array}{l} E_1 = e_1 : \text{process} \\ E_2 = e_2 : \text{state} \\ \text{RESTR} = <_\alpha \\ \text{HEAD} = e_1 \end{array} \right]$
ARGSTR =	$\left[ \begin{array}{ll} \text{ARG}_1 = \boxed{1} & \left[ \begin{array}{l} \text{animate\_individ[u]al} \\ \text{FORMAL} = \text{physobj} \end{array} \right] \\ \text{ARG}_2 = \boxed{2} & \left[ \begin{array}{l} \text{artifact} \\ \text{CONST} = \boxed{3} \\ \text{FORMAL} = \text{physobj} \end{array} \right] \\ \text{D-ARG}_1 = \boxed{3} & \left[ \begin{array}{l} \text{material} \\ \text{FORMAL} = \text{mass} \end{array} \right] \end{array} \right]$
QUALIA =	$\left[ \begin{array}{l} \text{create-lcp} \\ \text{FORMAL} = \text{exist}(e_2, \boxed{2}) \\ \text{AGENTIVE} = \text{build\_act}(e_1, \boxed{1}, \boxed{3}) \end{array} \right]$

The lexical representation for the English verb *build*. (From Pustejovsky, J., *The Generative Lexicon*, MIT Press, Cambridge, MA, 1995.)

# 8. ANÁLISIS SEMÁNTICO

## - Metalenguaje Semántico Natural (NSM)

- Descompositivo, cognitivo
- Universales semánticos (semantic primes)
- Moléculas semánticas (semantic molecules)

# 8. ANÁLISIS SEMÁNTICO

TABLE 5.1 Semantic Primes, Grouped into Related Categories

I, YOU, SOMEONE, SOMETHING/THING, PEOPLE, BODY	Substantives
KIND, PART	Relational substantives
THIS, THE SAME, OTHER/ELSE	Determiners
ONE, TWO, SOME, ALL, MUCH/MANY	Quantifiers
GOOD, BAD	Evaluators
BIG, SMALL	Descriptors
KNOW, THINK, WANT, FEEL, SEE, HEAR	Mental predicates
SAY, WORDS, TRUE	Speech
DO, HAPPEN, MOVE, TOUCH	Actions, events, movement, contact
BE (SOMEWHERE), THERE IS, HAVE, BE (SOMEONE/SOMETHING)	Location, existence, possession, specification
LIVE, DIE	Life and death
WHEN/TIME, NOW, BEFORE, AFTER, A LONG TIME, A SHORT TIME, FOR SOME TIME, MOMENT	Time
WHERE/PLACE, HERE, ABOVE, BELOW, FAR, NEAR, SIDE, INSIDE	Space
NOT, MAYBE, CAN, BECAUSE, IF	Logical concepts
VERY, MORE	Intensifier, augmentor
LIKE/WAY	Similarity

Notes: Primes exist as the meanings of lexical units (not at the level of lexemes). Exponents of primes may be words, bound morphemes, or phrasemes. They can be formally complex. They can have combinatorial variants (allolexes). Each prime has well-specified syntactic (combinatorial) properties.

# 8. ANÁLISIS SEMÁNTICO

TABLE 5.2 Semantic Roles and Their Conventional Definitions

Role	Description
<i>agent</i>	a wilful, purposeful instigator of an action or event
<i>effector</i>	the doer of an action, which may or may not be wilful or purposeful
<i>experiencer</i>	a sentient being that experiences internal states, such as perceivers, cognizers, and emoters
<i>instrument</i>	a normally inanimate entity manipulated by an agent in carrying out an action
<i>force</i>	somewhat like an instrument, but it cannot be manipulated
<i>patient</i>	a thing that is in a state or condition, or undergoes a change of state or condition
<i>theme</i>	a thing which is located or is undergoing a change of location (motion)
<i>benefactive</i>	the participant for whose benefit some action is performed
<i>recipient</i>	someone who gets something (recipients are always animate or some kind of quasi-animate entity)
<i>goal</i>	destination, which is similar to recipient, except that it is often inanimate
<i>source</i>	the point of origin of a state of affairs
<i>location</i>	a place or a spatial locus of a state of affairs
<i>path</i>	a route

Source: Van Valin, R.D. and LaPolla, R.J., *Syntax: Structure, Meaning and Function*, Cambridge University Press, Cambridge, U.K., 1997, 85-89.

# 8. ANÁLISIS SEMÁNTICO

[B] Semantic explication for *Someone X felt sad*

- a. someone X felt something bad
  - like someone can feel when they think like this:
- b. “I know that something bad happened
  - I don’t want things like this to happen
  - I can’t think like this: I will do something because of it now
  - I know that I can’t do anything”

[C] Semantic explication for *Someone X felt unhappy*

- a. someone X felt something bad
  - like someone can feel when they think like this:
- b. “some bad things happened to me
  - I wanted things like this not to happen to me
  - I can’t not think about it”
- c. this someone felt something like this, because this someone thought like this

## 8. ANÁLISIS SEMÁNTICO

[D] Semantic explication for *Someone X felt bei* [Chinese]

- a. someone X felt something very bad
  - like someone can feel when they think like this:
    - "something bad happened now
    - I know that after this good things will not happen anymore
    - I don't want things like this to happen
    - I want to do something if I can
    - I know that I can't do anything
    - because I know that no one can do anything when things like this happen'
- b. "something bad happened now"
  - I know that after this good things will not happen anymore
  - I don't want things like this to happen
  - I want to do something if I can
  - I know that I can't do anything
  - because I know that no one can do anything when things like this happen'
- c. this someone felt something like this, because this someone thought like this

[E] Semantic explication for *Someone X felt chou* [Chinese]

- a. someone X felt something bad
  - like someone can feel when they think like this:
    - "something bad is happening to me
    - before this, I did not think that this would happen to me
    - I don't want things like this to happen to me
    - because of this, I want to do something if I can
    - I don't know what I can do
    - I can't not think about this all the time"
- b. "something bad is happening to me"
  - before this, I did not think that this would happen to me
  - I don't want things like this to happen to me
  - because of this, I want to do something if I can
  - I don't know what I can do
  - I can't not think about this all the time"
- c. this someone felt something like this, because this someone thought like this

# 8. ANÁLISIS SEMÁNTICO

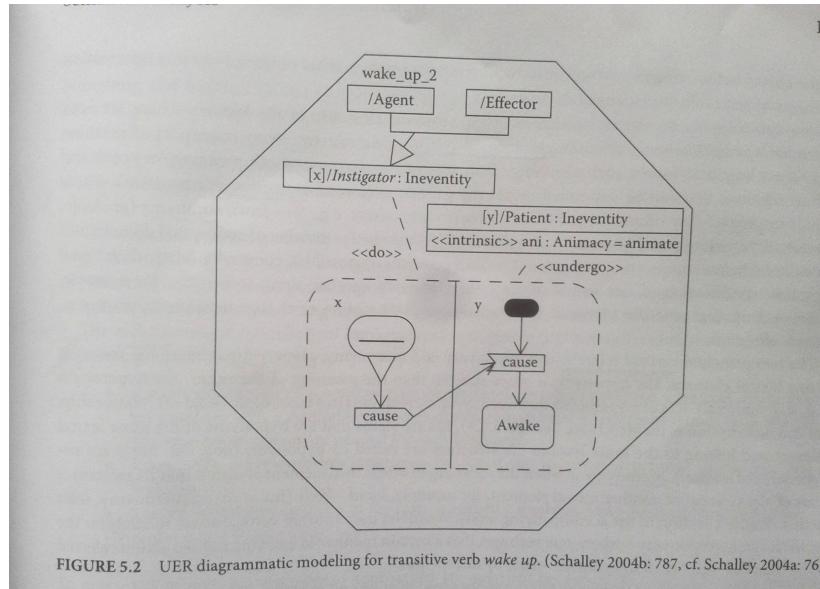
TABLE 5.3 Three Different Types of Classificatory Relationships in English

True Taxonomic Category Words	"Plural-Mostly" Functional Macro-Category Words	Singular-Only Functional-Collective Macro-Category Words
<i>birds—sparrow, wren, eagle, ...</i>	<i>vegetables—carrots, peas, celery, ...</i>	<i>furniture—table, chair, bed, ...</i>
<i>fish—trout, tuna, bream, ...</i>	<i>herbs—basil, oregano, rosemary, ...</i>	<i>cutlery—knife, fork, spoon, ...</i>
<i>animal—dog, cat, horse, ...</i>	<i>cosmetics—lipstick, powder, mascara, ...</i>	<i>jewelry—ring, earring, necklace, ...</i>

# 8. ANÁLISIS SEMÁNTICO

## - Semántica Orientada a Objetos

- Centrada en el verbo
- Representación UER (Unified Eventity Representation)



# **8. ANÁLISIS SEMÁNTICO**

## **- Relaciones semánticas y Ontologías**

### **· 1. Relaciones semánticas.**

- Wordnet

- Relaciones horizontales:

  - Sinonimia

  - Oposición

    - Incompatibilidad

    - Antonimia

    - Complementaridad

    - Conversidad

    - Reversibilidad

- Relaciones verticales

  - Hiponimia

  - Meronimia

  - Troponimia

# **8. ANÁLISIS SEMÁNTICO**

- **2. Relaciones ontológicas**
- **3. Semántica ontológica**
  - Arquitectura:
    - Fuentes de Conocimiento
    - Lenguajes de Representación
    - Módulos de procesamiento

# 8. ANÁLISIS SEMÁNTICO

## - Roles semánticos

- Centrada en el verbo
- FrameNet

TABLE 5.2 Semantic Roles and Their Conventional Definitions

Role	Description
<i>agent</i>	a wilful, purposeful instigator of an action or event
<i>effector</i>	the doer of an action, which may or may not be wilful or purposeful
<i>experiencer</i>	a sentient being that experiences internal states, such as perceivers, cognizers, and emoters
<i>instrument</i>	a normally inanimate entity manipulated by an agent in carrying out an action
<i>force</i>	somewhat like an instrument, but it cannot be manipulated
<i>patient</i>	a thing that is in a state or condition, or undergoes a change of state or condition
<i>theme</i>	a thing which is located or is undergoing a change of location (motion)
<i>benefactive</i>	the participant for whose benefit some action is performed
<i>recipient</i>	someone who gets something (recipients are always animate or some kind of quasi-animate entity)
<i>goal</i>	destination, which is similar to recipient, except that it is often inanimate
<i>source</i>	the point of origin of a state of affairs
<i>location</i>	a place or a spatial locus of a state of affairs
<i>path</i>	a route

Source: Van Valin, R.D. and LaPolla, R.J., *Syntax: Structure, Meaning and Function*, Cambridge University Press, Cambridge, U.K., 1997, 85-89.