

Capítulo 3

Algoritmos de análisis sintáctico para TAG

Se describen varios de los algoritmos de análisis sintáctico que han sido propuestos para las gramáticas de adjunción de árboles desde que éstas fueron por descritas por primera vez en [93]. La mayor parte de estos algoritmos están basados en algoritmos de análisis sintáctico independientes del contexto extendidos para el tratamiento de la adjunción. Asimismo es de destacar que prácticamente todos estos algoritmos utilizan técnicas de programación dinámica, obteniendo una complejidad polinómica de orden $\mathcal{O}(n^6)$, donde n es la longitud de la cadena de entrada. La aportación de este capítulo es múltiple: se proporciona un camino evolutivo que relaciona los algoritmos de análisis de TAG más populares; se describen por vez primera algunos algoritmos de análisis de TAG, como por ejemplo, las diferentes versiones de los algoritmos de tipo Earley ascendente y las versiones propuestas de los algoritmos de tipo Earley sin la propiedad del prefijo válido; por último, se realiza una descripción conjunta de la mayor parte de los algoritmos de análisis existentes para TAG. Este capítulo está basado en [8, 9].

3.1 Introducción

A la hora de describir los algoritmos de análisis para TAG, tenemos que elegir una representación adecuada para indicar el reconocimiento parcial de los árboles elementales. En las gramáticas independientes del contexto se utilizan habitualmente producciones con punto para separar la parte de la producción que ya ha sido procesada de la que no lo ha sido aún. En los algoritmos que proceden unidireccionalmente, un solo punto es suficiente, mientras que en aquellos que proceden bidireccionalmente se necesitan dos puntos [189]. En el caso de TAG, al no ser las estructuras elementales producciones sino árboles, deberemos representar árboles con punto. Existen varias alternativas:

1. Al igual que en las gramáticas independientes del contexto se escriben producciones completas con punto, en TAG se podría escribir cada árbol elemental completo con su punto correspondiente. Un ejemplo de utilización de este tipo de notación puede encontrarse en [176].
2. Si identificamos unívocamente todo elemento en una producción independiente del contexto, para lo cual podemos utilizar subíndices correspondientes al número de la producción y a la posición que cada elemento ocupa en la producción de tal modo que la producción $k : N \rightarrow MPQ$ se representaría como $N_{k,0} \rightarrow N_{k,1}N_{k,2}N_{k,3}$, podemos representar una producción con punto simplemente indicando un elemento de la producción contiguo al

punto e indicando si este se encuentra situado a derecha o izquierda de dicho elemento. Del mismo modo, podemos indicar la posición del punto en un árbol elemental de una TAG indicando simplemente un nodo del árbol adyacente al punto y la posición relativa del punto respecto a dicho nodo: arriba-izquierda, abajo-izquierda, abajo-derecha o arriba-derecha. Este es el tipo de notación utilizada en [168].

3. Como un árbol elemental γ puede considerarse constituido por un conjunto de producciones independientes del contexto $\mathcal{P}(\gamma)$, podemos indicar la posición del punto en el árbol simplemente indicando la posición en la producción correspondiente. Este tipo de representación recibe el nombre de *multicapa* en [64].
4. Aplicar un aplanamiento a los árboles, esto es, una conversión de los mismos a cadenas de caracteres parentizadas en las cuales cada nuevo nivel de paréntesis indica un nivel adicional de profundidad en el árbol elemental. Este tipo de representación, que recibe el nombre de *plana* en [64], implica que cada árbol inicial se representa únicamente por una cadena de caracteres mientras que cada árbol auxiliar se representa por dos cadenas correspondientes a las partes situadas a la izquierda y a la derecha de la espina. Dichas cadenas pueden ser utilizadas como partes derechas de las producciones de una gramática independiente del contexto que permitiría representar de modo conciso los árboles elementales de una gramática de adjunción [64].

La primera alternativa presenta el inconveniente de que las representaciones lineales de árboles suelen dificultar su comprensión. Alternativamente podrían utilizarse representaciones pictóricas de los árboles, que si bien son incómodas a lo hora de describir los pasos deductivos del algoritmo pueden ser útiles para representar gráficamente el comportamiento del algoritmo. La segunda alternativa, aunque muy manejable y concisa tiene el inconveniente de que es necesario recurrir a las descripciones de los árboles elementales para poder ver el contexto en el que se está aplicando cada paso del algoritmo. La tercera alternativa resuelve este problema al proporcionar un contexto formado por los hermanos del nodo considerado, que si bien es limitado, proporciona al lector la información suficiente la mayor parte de las veces. La desventaja es que normalmente deberemos añadir pasos deductivos extra que realizan el recorrido del árbol en un orden determinado. La cuarta alternativa, aunque en principio supone una reducción del número de producciones independientes al contexto, en la práctica su utilización no implica generalmente una reducción del número de pasos deductivos utilizados para describir el procesamiento realizado por los algoritmos de análisis sintáctico¹ y contribuye a oscurecer la notación.

En nuestro caso hemos preferido utilizar la representación multicapa, por lo que indicaremos la posición del punto en un árbol mediante una producción $N \rightarrow \delta \bullet \nu$, donde $\delta \nu$ son los hijos de N . En el caso de los elementos del lado derecho de las producciones cuyas etiquetas pertenecen a $V_T \cup \epsilon$, puesto que no pueden tener descendientes ni son susceptibles de ser nodos de adjunción, consideraremos directamente su etiqueta a la hora de describir las producciones, en lugar del nodo propiamente dicho². Adicionalmente y con el fin de simplificar los esquemas de análisis de los algoritmos más complejos, seguiremos el enfoque de [125] de considerar la producción adicional $\top \rightarrow \mathbf{R}^\alpha$ para cada árbol inicial α y las dos producciones adicionales siguientes para cada árbol auxiliar β : $\top \rightarrow \mathbf{R}^\beta$ y $\mathbf{F}^\beta \rightarrow \perp$, donde \mathbf{R}^β y \mathbf{F}^β se refieren a los nodos raíz y pie de β , respectivamente. Con el fin de no modificar la capacidad generativa de las gramáticas,

¹ Como ejemplo, podemos comparar el algoritmo descrito en [64] con el representado por el esquema 3.7.

² Con esta convención, que se adopta porque simplifica considerablemente la notación, la única forma de distinguir diferentes hojas de un árbol elemental con la misma etiqueta es referenciando los respectivos padres o la producción independiente del contexto de la que forman parte.

los nuevos nodos \top y \perp no pueden ser nodos de adjunción, por lo cual la nueva gramática se asemeja a una TAG *limpia*³.

Con respecto a las restricciones de adjunción, $\beta \in \text{adj}(N^\gamma)$ si $\beta \in \mathbf{A}$ puede ser adjuntado en N^γ . Si $\beta = \mathbf{nil}$ entonces la adjunción en el nodo N^γ es opcional. Si \mathbf{nil} es el único valor que retorna la función para un nodo N^γ , entonces no es posible realizar ninguna adjunción en dicho nodo.

Ejemplo 3.1 El árbol inicial α de la figura 3.1 se representa en notación multicapa mediante el siguiente conjunto de producciones independientes del contexto:

$$\begin{aligned}\top &\rightarrow \langle \alpha, 0 \rangle \\ \langle \alpha, 0 \rangle &\rightarrow a \langle \alpha, 2 \rangle \\ \langle \alpha, 2 \rangle &\rightarrow c\end{aligned}$$

El árbol auxiliar β de la misma figura se representa mediante el siguiente conjunto de producciones:

$$\begin{aligned}\top &\rightarrow \langle \beta, 0 \rangle \\ \langle \beta, 0 \rangle &\rightarrow a \langle \beta, 2 \rangle \\ \langle \beta, 2 \rangle &\rightarrow \langle \beta, 2.1 \rangle c \\ \langle \beta, 2.1 \rangle &\rightarrow \perp\end{aligned}$$

En ambos casos $\langle \alpha, g \rangle$ denota el nodo de α que ocupa la posición g según el direccionamiento de Gorn⁴. ¶

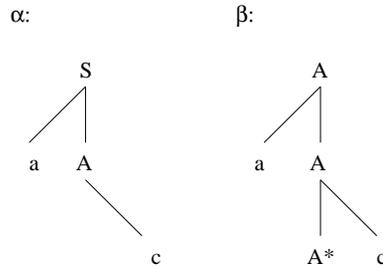


Figura 3.1: Ejemplos de árbol inicial y auxiliar

La relación \Rightarrow de derivación sobre $\mathcal{P}(\gamma)$ se define como $\delta \Rightarrow \nu$ si existen $\delta', \delta'', M^\gamma, v$ tal que $\delta = \delta' M^\gamma \delta''$, $\nu = \delta' v \delta''$ y existe una producción $M^\gamma \rightarrow v \in \mathcal{P}(\gamma)$. Mediante $\overset{*}{\Rightarrow}$ denotaremos el cierre reflexivo y transitivo de \Rightarrow .

Sea $\mathcal{P}(\mathcal{T}) = \bigcup_{\gamma \in \mathbf{I} \cup \mathbf{A}} \mathcal{P}(\gamma)$ el conjunto de todas las producciones independientes del contexto presentes en los árboles de una gramática de adjunción de árboles \mathcal{T} . Con el fin de ser capaces de representar derivaciones que incluyan adjunciones extendemos la relación \Rightarrow a $\mathcal{P}(\mathcal{T})$, de tal modo que $\delta \overset{*}{\Rightarrow} \nu$ si existen $\delta', \delta'', M^\gamma, v$ tal que $\delta = \delta' M^\gamma \delta''$, $\mathbf{R}^\beta \overset{*}{\Rightarrow} v_1 \mathbf{F}^\beta v_3$, $\beta \in \text{adj}(M^\gamma)$, $M^\gamma \rightarrow v_2$ y $\nu = \delta' v_1 v_2 v_3 \delta''$.

³Poller y Becker [149] denominan TAG *limpias* a las gramáticas de adjunción en las que los nodos raíz de los árboles elementales y los nodos pie de los árboles auxiliares tienen restricciones de adjunción nulas.

⁴En el direccionamiento de Gorn se utiliza 0 para referirse al nodo raíz, k para referirse al k -ésimo hijo del nodo raíz y $p.q$ para referirse al q -ésimo hijo del nodo con dirección p .

La mayoría de los algoritmos de análisis sintáctico diseñados para TAG se corresponden con adaptaciones de algoritmos para el análisis de gramáticas independientes del contexto y por lo tanto recibirán el mismo nombre en ambos casos. Solamente en aquellos casos en los que pueda existir confusión diferenciaremos explícitamente entre unos y otros⁵. Para describir los diferentes algoritmos de análisis sintáctico haremos uso de *esquemas de análisis*, una estructura para realizar descripciones de alto nivel de algoritmos de análisis desarrollada por Sikkel en [189] y que se describe en el apéndice A.

3.2 Algoritmo de tipo CYK

A continuación mostramos una extensión para TAG del algoritmo CYK de análisis sintáctico (ver sección B.1) basado en el algoritmo descrito en [209, 206] pero con algunas correcciones. Asumiremos que todo nodo de un árbol elemental de la gramática tiene a lo sumo dos descendientes. Este condicionante puede verse como una trasposición de la forma normal de Chomsky [85] al caso de las gramáticas de adjunción.

El reconocimiento de los nodos de un árbol elemental se realiza aplicando casi literalmente el algoritmo para el caso independiente del contexto. Sin embargo ahora es preciso definir cómo reconocer la operación de adjunción de forma totalmente ascendente. Para ello se definen ítems de la forma

$$\left\{ \begin{array}{l} [N^\gamma, i, j \mid p, q \mid adj] \mid \begin{array}{l} N^\gamma \xrightarrow{*} a_{i+1} \dots a_p \mathbf{F}^\gamma a_{q+1} \dots a_j \xrightarrow{*} a_{i+1} \dots a_j \quad \text{sii } (p, q) \neq (-, -) \\ N^\gamma \xrightarrow{*} a_{i+1} \dots a_j \quad \text{sii } (p, q) = (-, -) \end{array} \end{array} \right\}$$

que contienen un nuevo componente *adj* con respecto a los ítems definidos en [209, 206], que toma su valor del conjunto {true, false} y que indica

- si su valor es *true* significa que el ítem es el resultado de una operación de adjunción ya totalmente realizada;
- si valor es *false* significa que el ítem no es el resultado de una adjunción.

Con ello podemos limitar a una sola la cantidad de adjunciones que se pueden realizar sobre cada nodo de un árbol elemental y podemos también asegurar el cumplimiento de las restricciones de adjunción obligatoria, ajustándonos de este modo a lo establecido en el formalismo de las gramáticas de adjunción [175].

Esquema de análisis sintáctico 3.1 El sistema de análisis \mathbb{P}_{CYK} que se corresponde con el algoritmo CYK para una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\text{CYK}} = \left\{ [N^\gamma, i, j \mid p, q \mid adj] \mid \begin{array}{l} \text{etiqueta}(N^\gamma) \in V_N, \gamma \in \mathbf{I} \cup \mathbf{A}, \quad 0 \leq i \leq j, \\ (p, q) \leq (i, j), \quad adj \in \{\text{true}, \text{false}\} \end{array} \right\}$$

$$\mathcal{H}_{\text{CYK}} = \{ [a, i-1, i] \mid a = a_i, 1 \leq i \leq n \}$$

$$\mathcal{D}_{\text{CYK}}^{\text{Scan}} = \frac{[a, i, i+1]}{[N^\gamma, i, i+1 \mid -, - \mid \text{false}]} \quad N^\gamma \rightarrow a \in \mathcal{P}(\gamma)$$

⁵En el apéndice B se describen sucintamente los algoritmos de análisis sintáctico para gramáticas independientes del contexto que son relevantes en este capítulo

$$\mathcal{D}_{\text{CYK}}^\epsilon = \overline{[N^\gamma, i, i \mid -, - \mid \text{false}]} \quad N^\gamma \rightarrow \epsilon \in \mathcal{P}(\gamma)$$

$$\mathcal{D}_{\text{CYK}}^{\text{Foot}} = \overline{[\mathbf{F}^\gamma, i, j \mid i, j \mid \text{false}]}$$

$$\mathcal{D}_{\text{CYK}}^{\text{LeftDom}} = \frac{[M^\gamma, i, k \mid p, q \mid \text{adj1}], [P^\gamma, k, j \mid -, - \mid \text{adj2}]}{[N^\gamma, i, j \mid p, q \mid \text{false}]}$$

$$\begin{aligned} N^\gamma &\rightarrow M^\gamma P^\gamma \in \mathcal{P}(\gamma), \\ M^\gamma &\in \text{espina}(\gamma), \\ \text{adj1} &= \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma), \\ \text{adj2} &= \text{false sii } \mathbf{nil} \in \text{adj}(P^\gamma) \end{aligned}$$

$$\mathcal{D}_{\text{CYK}}^{\text{RightDom}} = \frac{[M^\gamma, i, k \mid -, - \mid \text{adj1}], [P^\gamma, k, j \mid p, q \mid \text{adj2}]}{[N^\gamma, i, j \mid p, q \mid \text{false}]}$$

$$\begin{aligned} N^\gamma &\rightarrow M^\gamma P^\gamma \in \mathcal{P}(\gamma), \\ P^\gamma &\in \text{espina}(\gamma), \\ \text{adj1} &= \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma), \\ \text{adj2} &= \text{false sii } \mathbf{nil} \in \text{adj}(P^\gamma) \end{aligned}$$

$$\mathcal{D}_{\text{CYK}}^{\text{NoDom}} = \frac{[M^\gamma, i, k \mid -, - \mid \text{adj1}], [P^\gamma, k, j \mid -, - \mid \text{adj2}]}{[N^\gamma, i, j \mid -, - \mid \text{false}]}$$

$$\begin{aligned} N^\gamma &\rightarrow M^\gamma P^\gamma \in \mathcal{P}(\gamma), \\ N^\gamma &\notin \text{espina}(\gamma), \\ \text{adj1} &= \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma), \\ \text{adj2} &= \text{false sii } \mathbf{nil} \in \text{adj}(P^\gamma) \end{aligned}$$

$$\mathcal{D}_{\text{CYK}}^{\text{Unary}} = \frac{[M^\gamma, i, j \mid p, q \mid \text{adj}]}{[N^\gamma, i, j \mid p, q \mid \text{false}]} \quad \begin{aligned} N^\gamma &\rightarrow M^\gamma \in \mathcal{P}(\gamma), \\ \text{adj} &= \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma) \end{aligned}$$

$$\mathcal{D}_{\text{CYK}}^{\text{Adj}} = \frac{[\mathbf{R}^\beta, i', j' \mid i, j \mid \text{adj}], [N^\gamma, i, j \mid p, q \mid \text{false}]}{[N^\gamma, i', j' \mid p, q \mid \text{true}]} \quad \beta \in \mathbf{A}, \beta \in \text{adj}(N^\gamma)$$

$$\mathcal{D}_{\text{CYK}} = \mathcal{D}_{\text{CYK}}^{\text{Scan}} \cup \mathcal{D}_{\text{CYK}}^\epsilon \cup \mathcal{D}_{\text{CYK}}^{\text{Foot}} \cup \mathcal{D}_{\text{CYK}}^{\text{LeftDom}} \cup \mathcal{D}_{\text{CYK}}^{\text{RightDom}} \cup \mathcal{D}_{\text{CYK}}^{\text{NoDom}} \cup \mathcal{D}_{\text{CYK}}^{\text{Unary}} \cup \mathcal{D}_{\text{CYK}}^{\text{Adj}}$$

$$\mathcal{F}_{\text{CYK}} = \{ [\mathbf{R}^\alpha, 0, n \mid -, - \mid \text{adj}] \mid \alpha \in \mathbf{I}, \text{adj} = \text{false sii } \mathbf{nil} \in \text{adj}(\mathbf{R}^\alpha) \}$$

donde $(p, q) \leq (i, j)$ se cumple si $i \leq p \leq q \leq j$ en el caso de que $(p, q) \neq (-, -)$ y si $i \leq j$ en el caso de que $(p, q) = (-, -)$. §

La definición de las hipótesis realizada en este sistema de análisis sintáctico se corresponde con la estándar y es la misma que se utilizará en los restantes sistemas de análisis del capítulo. Por consiguiente, no nos volveremos a referir explícitamente a ellas.

Los pasos clave en el sistema de análisis \mathbb{P}_{CYK} son $\mathcal{D}_{\text{CYK}}^{\text{Foot}}$ y $\mathcal{D}_{\text{CYK}}^{\text{Adj}}$, puesto que son los que definen el mecanismo de adjunción. Los demás pasos se encargan de recorrer de forma ascendente los árboles elementales, propagando la información relativa a la porción de la cadena de entrada cubierta por el pie en el caso de los árboles auxiliares. El paso deductivo $\mathcal{D}_{\text{CYK}}^{\text{Scan}}$ permite resolver la limitación planteada en [209, 206] con respecto a nodos frontera etiquetados con ϵ .

El paso deductivo $\mathcal{D}_{\text{CYK}}^{\text{Foot}}$ hace posible el análisis ascendente de los árboles auxiliares, puesto que predice todas las posibles porciones de la cadena de entrada que puede cubrir el pie. Si

no fuese así, sería necesario esperar al reconocimiento total del pie antes de poder continuar el reconocimiento del árbol auxiliar, con lo cual el algoritmo de análisis no sería totalmente ascendente.

Una de las consecuencias de la utilización del paso $\mathcal{D}_{\text{CYK}}^{\text{Foot}}$ es que existirán múltiples análisis diferentes para cada árbol auxiliar, diferentes únicamente en las posición de la cadena de entrada que se *supone* en cada caso que cubre el pie. Pero finalmente sólo uno o unos pocos de esos árboles podrán formar parte del árbol de derivación, aquellos que hayan predicho correctamente el pie. La realización de esta comprobación le corresponde al paso $\mathcal{D}_{\text{CYK}}^{\text{Adj}}$. Efectivamente, cuando se ha alcanzado la raíz de un árbol auxiliar, se comprueba si existe un subárbol de un árbol elemental cuya raíz pueda ser un nodo de adjunción para dicho árbol auxiliar y que además cubra la porción de la cadena de entrada que espera ser cubierta por el pie. En caso de que así sea, se eliminará la posibilidad de realizar otras adjunciones en ese nodo y el análisis continuará ascendiendo por el nuevo árbol elemental. En la figura 3.2 se muestra gráficamente la aplicación de este paso, donde los dos árboles de la izquierda se corresponden con las descripciones de los ítems antecedentes y el árbol de la derecha con la descripción del ítem consecuente. Se utilizan líneas punteadas para indicar partes de árboles elementales que no han sido aún analizadas, mientras que las espinas se resaltan mediante líneas gruesas.

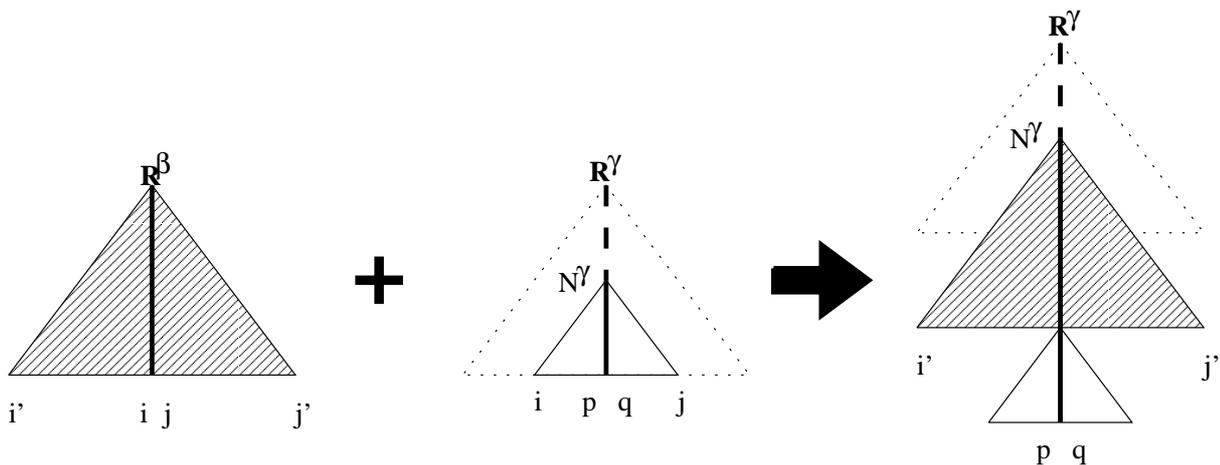


Figura 3.2: Descripción gráfica de un paso $\mathcal{D}_{\text{CYK}}^{\text{Adj}}$

La complejidad temporal del algoritmo con respecto a la longitud n de la cadena de entrada es $\mathcal{O}(n^6)$ y viene dada por el paso deductivo encargado de la adjunción, que debe combinar seis posiciones de la cadena de entrada. La complejidad espacial con respecto a la cadena de entrada es $\mathcal{O}(n^4)$, puesto que cada ítem almacena cuatro posiciones de la cadena de entrada.

3.3 Algoritmos de tipo Earley ascendente

El algoritmo CYK presenta una limitación muy importante: sólo es aplicable a gramáticas en las cuales un nodo no tiene más de dos descendentes. Nuestro objetivo ahora es extender el esquema **CYK** a la clase general de TAG, obteniendo lo que podríamos llamar un esquema Earley ascendente (ver sección B.2) extendido a gramáticas de adjunción de árboles. Debemos reseñar que no conocemos ninguna adaptación anterior para TAG de este algoritmo.

Como primer paso para la realización de un algoritmo de tipo Earley ascendente para gramáticas de adjunción de árboles precisamos introducir puntos en las producciones que repre-

sentan los árboles elementales, por lo que los ítems que utilizaremos tendrán la forma

$$\left\{ \begin{array}{l} [N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid adj] \mid \delta \xrightarrow{*} a_{i+1} \dots a_p \mathbf{F}^\gamma a_{q+1} \dots a_j \xrightarrow{*} a_{i+1} \dots a_j \quad \text{sii } (p, q) \neq (-, -) \\ \delta \xrightarrow{*} a_{i+1} \dots a_j \quad \text{sii } (p, q) = (-, -) \end{array} \right\}$$

Los ítems del nuevo esquema de análisis sintáctico, que denominaremos \mathbf{buE}_1 , son por tanto un refinamiento de los ítems de \mathbf{CYK} . Sobre los pasos deductivos aplicaremos también un *refinamiento* puesto que los pasos de tipo LeftDom, RightDom y NoDom serán separados en diferentes pasos de tipo Init y Comp, mientras que los pasos de tipo Unary y ϵ ya no serán necesarios puesto que todas las producciones serán tratadas uniformemente con independencia de su longitud. Finalmente, se realizará una extensión del dominio de las producciones permitiendo árboles con un grado arbitrario de ramificación.

Esquema de análisis sintáctico 3.2 El sistema de análisis $\mathbb{P}_{\mathbf{buE}_1}$ que se corresponde con una primera extensión del algoritmo de Earley ascendente para TAG, dada una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\mathbf{buE}_1} = \left\{ [N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid adj] \mid \begin{array}{l} N^\gamma \rightarrow \delta \nu \in \mathcal{P}(\gamma), \gamma \in \mathbf{I} \cup \mathbf{A}, \\ 0 \leq i \leq j, (p, q) \leq (i, j), adj \in \{\text{true}, \text{false}\} \end{array} \right\}$$

$$\mathcal{D}_{\mathbf{buE}_1}^{\text{Init}} = \overline{[N^\gamma \rightarrow \bullet \delta, i, i \mid -, - \mid \text{false}]}$$

$$\mathcal{D}_{\mathbf{buE}_1}^{\text{Foot}} = \overline{[\mathbf{F}^\beta \rightarrow \perp \bullet, i, j \mid i, j \mid \text{false}]}$$

$$\mathcal{D}_{\mathbf{buE}_1}^{\text{Scan}} = \frac{[N^\gamma \rightarrow \delta \bullet a\nu, i, j \mid p, q \mid \text{false}], [a, j, j+1]}{[N^\gamma \rightarrow \delta a \bullet \nu, i, j+1 \mid p, q \mid \text{false}]}$$

$$\mathcal{D}_{\mathbf{buE}_1}^{\text{Comp}} = \frac{\begin{array}{l} [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}], \\ [M^\gamma \rightarrow \nu \bullet, i, j \mid p, q \mid adj] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]} \quad adj = \text{false} \text{ sii } \mathbf{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\mathbf{buE}_1}^{\text{Adj}} = \frac{[\top \rightarrow \mathbf{R}^\beta \bullet, k, j \mid k', j' \mid \text{false}], [M^\gamma \rightarrow \nu \bullet, k', j' \mid p, q \mid \text{false}]}{[M^\gamma \rightarrow \nu \bullet, k, j \mid p, q \mid \text{true}]} \quad \beta \in \mathbf{A}, \beta \in \text{adj}(\gamma)$$

$$\mathcal{D}_{\mathbf{buE}_1} = \mathcal{D}_{\mathbf{buE}_1}^{\text{Init}} \cup \mathcal{D}_{\mathbf{buE}_1}^{\text{Foot}} \cup \mathcal{D}_{\mathbf{buE}_1}^{\text{Scan}} \cup \mathcal{D}_{\mathbf{buE}_1}^{\text{Comp}} \cup \mathcal{D}_{\mathbf{buE}_1}^{\text{Adj}}$$

$$\mathcal{F}_{\mathbf{buE}_1} = \left\{ [\top \rightarrow \mathbf{R}^\alpha \bullet, 0, n \mid -, - \mid \text{false}] \mid \alpha \in \mathbf{I} \right\}$$

y donde $p \cup q$ se refiere a la operación de unión de índices, función parcial de enteros a enteros definida como

$$p \cup q = \begin{cases} p & \text{si } q = - \\ q & \text{si } p = - \\ - & \text{si } p = q = - \end{cases}$$

donde $-$ se utiliza para indicar que el valor de uno de los parámetros está indefinido. §

Los pasos $\mathcal{D}_{\text{buE}_1}^{\text{Init}}$ son los encargados de lanzar el análisis ascendente. Los ítems generados por estos pasos son siempre válidos, puesto que no expanden ninguna porción la cadena de entrada. Tan sólo expresan la expectativa de que una determinada producción pueda ser aplicada para reconocer una porción de la cadena que comienza en una determinada posición.

Los pasos $\mathcal{D}_{\text{buE}_1}^{\text{Foot}}$ son utilizados, al igual que en el caso CYK, para predecir la porción de la cadena de entrada cubierta por el pie de un árbol auxiliar.

El algoritmo procede a reconocer ascendentemente los árboles auxiliares mediante la aplicación de pasos $\mathcal{D}_{\text{buE}_1}^{\text{Comp}}$, propagando también la información correspondiente a la porción de la cadena de entrada cubierta por el pie en el caso de los árboles auxiliares.

El paso deductivo $\mathcal{D}_{\text{buE}_1}^{\text{Adj}}$ se comporta de modo idéntico al paso homónimo del algoritmo CYK, comprobando que las predicciones respecto al pie se corresponden con las de una adjunción que se ha realizado realmente.

Proposición 3.1 $\text{CYK} \xrightarrow{\text{ir}} \text{CYK}' \xrightarrow{\text{sr}} \text{ECYK} \xrightarrow{\text{ext}} \text{buE}_1.$

Demostración:

Como primer paso definiremos el sistema de análisis $\mathbb{P}_{\text{CYK}'}$ para una gramática de adjunción \mathcal{T} y una cadena de entrada $a_1 \dots a_n$.

$$\mathcal{I}_{\text{CYK}'} = \left\{ [N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid \text{adj}] \mid \begin{array}{l} N^\gamma \rightarrow \delta \nu \in \mathcal{P}(\gamma), \gamma \in \mathbf{I} \cup \mathbf{A} \\ 0 \leq i \leq j, (p, q) \leq (k, j), \text{adj} \in \{\text{true}, \text{false}\} \end{array} \right\}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{Scan}} = \frac{[a, i, i+1]}{[N^\gamma \rightarrow a \bullet, i, i+1 \mid -, - \mid \text{false}]}$$

$$\mathcal{D}_{\text{CYK}'}^\epsilon = \frac{N^\gamma \rightarrow \epsilon}{[N^\gamma \rightarrow \bullet, i, i \mid -, - \mid \text{false}]}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{Foot}} = \frac{[\mathbf{F}^\gamma \rightarrow \perp \bullet, i, j \mid i, j \mid \text{false}]}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{LeftDom}} = \frac{\begin{array}{l} [M^\gamma \rightarrow \delta \bullet, i, k \mid p, q \mid \text{adj1}], \\ [P^\gamma \rightarrow \nu \bullet, k, j \mid -, - \mid \text{adj2}] \end{array}}{[N^\gamma \rightarrow M^\gamma P^\gamma \bullet, i, j \mid p, q \mid \text{false}]} \quad \begin{array}{l} N^\gamma \rightarrow M^\gamma P^\gamma \in \mathcal{P}(\gamma), \\ M^\gamma \in \text{espina}(\gamma), \\ \text{adj1} = \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma), \\ \text{adj2} = \text{false sii } \mathbf{nil} \in \text{adj}(P^\gamma) \end{array}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{RightDom}} = \frac{\begin{array}{l} [M^\gamma \rightarrow \delta \bullet, i, k \mid -, - \mid \text{adj1}], \\ [P^\gamma \rightarrow \nu \bullet, k, j \mid p, q \mid \text{adj2}] \end{array}}{[N^\gamma \rightarrow M^\gamma P^\gamma \bullet, i, j \mid p, q \mid \text{false}]} \quad \begin{array}{l} N^\gamma \rightarrow M^\gamma P^\gamma \in \mathcal{P}(\gamma), \\ P^\gamma \in \text{espina}(\gamma), \\ \text{adj1} = \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma), \\ \text{adj2} = \text{false sii } \mathbf{nil} \in \text{adj}(P^\gamma) \end{array}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{NoDom}} = \frac{\begin{array}{l} [M^\gamma \rightarrow \delta \bullet, i, k \mid -, - \mid \text{adj1}], \\ [P^\gamma \rightarrow \nu \bullet, k, j \mid -, - \mid \text{adj2}] \end{array}}{[N^\gamma \rightarrow M^\gamma P^\gamma \bullet, i, j \mid -, - \mid \text{false}]} \quad \begin{array}{l} N^\gamma \rightarrow M^\gamma P^\gamma \in \mathcal{P}(\gamma), \\ N^\gamma \notin \text{espina}(\gamma), \\ \text{adj1} = \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma), \\ \text{adj2} = \text{false sii } \mathbf{nil} \in \text{adj}(P^\gamma) \end{array}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{Unary}} = \frac{[M^\gamma \rightarrow \delta \bullet, i, j \mid p, q \mid \text{adj}]}{[N^\gamma \rightarrow M^\gamma \bullet, i, j \mid p, q \mid \text{false}]} \quad \begin{array}{l} N^\gamma \rightarrow M^\gamma \in \mathcal{P}(\gamma), \\ \text{adj} = \text{false sii } \mathbf{nil} \in \text{adj}(M^\gamma) \end{array}$$

$$\mathcal{D}_{\text{CYK}'}^{\text{Adj}} = \frac{[\top \rightarrow \mathbf{R}^{\beta \bullet}, i', j' \mid i, j \mid \text{adj}], [N^\gamma \rightarrow \delta \bullet, i, j \mid p, q \mid \text{false}]}{[N^\gamma \rightarrow \delta \bullet, i', j' \mid p, q \mid \text{true}]} \quad \beta \in \mathbf{A}, \beta \in \text{adj}(N^\gamma)$$

$$\mathcal{D}_{\text{CYK}'} = \mathcal{D}_{\text{CYK}'}^{\text{Scan}} \cup \mathcal{D}_{\text{CYK}'}^\epsilon \cup \mathcal{D}_{\text{CYK}'}^{\text{Foot}} \cup \mathcal{D}_{\text{CYK}'}^{\text{LeftDom}} \cup \mathcal{D}_{\text{CYK}'}^{\text{RightDom}} \cup \mathcal{D}_{\text{CYK}'}^{\text{NoDom}} \cup \mathcal{D}_{\text{CYK}'}^{\text{Unary}} \cup \mathcal{D}_{\text{CYK}'}^{\text{Adj}}$$

$$\mathcal{F}_{\text{CYK}'} = \{ [\top \rightarrow \mathbf{R}^\alpha \bullet, 0, n \mid -, - \mid \text{adj}] \mid \alpha \in \mathbf{I}, \text{adj} = \text{false sii } \mathbf{nil} \in \text{adj}(\mathbf{R}^\alpha) \}$$

Para demostrar que $\text{CYK} \xrightarrow{\text{ir}} \text{CYK}'$, definiremos la siguiente función

$$f([N^\gamma \rightarrow \delta \bullet, i, j \mid p, q \mid \text{adj}]) = [N^\gamma, i, j \mid p, q \mid \text{adj}]$$

de la cual se obtiene directamente que $\mathcal{I}_{\text{CYK}} = f(\mathcal{I}_{\text{CYK}'})$ y que $\Delta_{\text{CYK}} = f(\Delta_{\text{CYK}'})$ por inducción en la longitud de las secuencias de derivación. En consecuencia, $\mathbb{P}_{\text{CYK}} \xrightarrow{\text{ir}} \mathbb{P}_{\text{CYK}'}$, con lo que hemos probado lo que pretendíamos.

Definiremos ahora el sistema de análisis sintáctico \mathbb{P}_{ECYK} para una gramática de adjunción de árboles \mathcal{T} en la que ningún nodo puede tener más de dos descendientes y una cadena de entrada $a_1 \dots a_n$ dada:

$$\mathcal{I}_{\text{ECYK}} = \mathcal{I}_{\text{CYK}'} = \mathcal{I}_{\text{buE}_1}$$

$$\mathcal{D}_{\text{ECYK}}^{\text{Init}} = \mathcal{D}_{\text{buE}_1}^{\text{Init}}$$

$$\mathcal{D}_{\text{ECYK}}^{\text{Foot}} = \mathcal{D}_{\text{buE}_1}^{\text{Foot}}$$

$$\mathcal{D}_{\text{ECYK}}^{\text{Scan}} = \mathcal{D}_{\text{buE}_1}^{\text{Scan}}$$

$$\mathcal{D}_{\text{ECYK}}^{\text{Comp}} = \mathcal{D}_{\text{buE}_1}^{\text{Comp}}$$

$$\mathcal{D}_{\text{ECYK}}^{\text{Adj}} = \mathcal{D}_{\text{buE}_1}^{\text{Adj}}$$

$$\mathcal{D}_{\text{ECYK}} = \mathcal{D}_{\text{ECYK}}^{\text{Init}} \cup \mathcal{D}_{\text{ECYK}}^{\text{Foot}} \cup \mathcal{D}_{\text{ECYK}}^{\text{Scan}} \cup \mathcal{D}_{\text{ECYK}}^{\text{Comp}} \cup \mathcal{D}_{\text{ECYK}}^{\text{Adj}}$$

$$\mathcal{F}_{\text{ECYK}} = \mathcal{F}_{\text{buE}_1}$$

Para demostrar que $\text{CYK}' \xrightarrow{\text{sr}} \text{ECYK}$, deberemos demostrar que para todo sistema de análisis $\mathbb{P}_{\text{CYK}'}$ y \mathbb{P}_{ECYK} se cumple que $\mathcal{I}_{\text{CYK}'} \subseteq \mathcal{I}_{\text{ECYK}}$ y que $\vdash_{\text{CYK}'}^* \subseteq \vdash_{\text{ECYK}}^*$. Lo primero es cierto por definición, puesto que $\mathcal{I}_{\text{CYK}'} = \mathcal{I}_{\text{ECYK}}$. Para lo segundo debemos mostrar que $\mathcal{D}_{\text{CYK}'} \subseteq \vdash_{\text{ECYK}}^*$. Consideremos caso por caso:

- un paso deductivo $\mathcal{D}_{\text{CYK}'}^{\text{Scan}}$ es equivalente a la secuencia de pasos deductivos constituida por la aplicación de un paso $\mathcal{D}_{\text{ECYK}}^{\text{Init}}$ y un paso $\mathcal{D}_{\text{ECYK}}^{\text{Scan}}$:

$$\overline{[N^\gamma \rightarrow \bullet a, i, i, \mid -, - \mid \text{false}]}$$

$$\frac{[N^\gamma \rightarrow \bullet a, i, i, \mid -, - \mid \text{false}], [a, i, i + 1]}{[N^\gamma \rightarrow a \bullet, i, i + 1, \mid -, - \mid \text{false}]}$$

- un paso deductivo $\mathcal{D}_{\text{CYK}'}^\epsilon$ es equivalente a un paso $\mathcal{D}_{\text{ECYK}}^{\text{Init}}$:

$$\overline{[N^\gamma \rightarrow \bullet, i, i, \mid -, - \mid \text{false}]}$$

- $\mathcal{D}_{\text{ECYK}}^{\text{Foot}} = \mathcal{D}_{\text{CYK}'}^{\text{Foot}}$.

- un paso deductivo $\mathcal{D}_{\text{CYK}'}^{\text{LeftDom}}$ es equivalente a la secuencia de pasos deductivos constituida por un paso $\mathcal{D}_{\text{ECYK}}^{\text{Init}}$ y dos pasos $\mathcal{D}_{\text{ECYK}}^{\text{Comp}}$:

$$\frac{\frac{\frac{[N^\gamma \rightarrow \bullet M^\gamma P^\gamma, i, i, | -, - | \text{false}]}{[N^\gamma \rightarrow \bullet M^\gamma P^\gamma, i, i, | -, - | \text{false}], [M^\gamma \rightarrow \delta \bullet, i, k | p, q | \text{adj}]}{[N^\gamma \rightarrow M^\gamma \bullet P^\gamma, i, k, | p, q | \text{false}]}}{[N^\gamma \rightarrow M^\gamma \bullet P^\gamma, i, k, | p, q | \text{false}], [P^\gamma \rightarrow \nu \bullet, k, j | -, - | \text{adj}]}{[N^\gamma \rightarrow M^\gamma P^\gamma \bullet, i, j, | p, q | \text{false}]}}$$

- un paso $\mathcal{D}_{\text{CYK}'}^{\text{RightDom}}$ es equivalente a la secuencia formada por un paso $\mathcal{D}_{\text{ECYK}}^{\text{Init}}$ y dos pasos $\mathcal{D}_{\text{ECYK}}^{\text{Comp}}$:

$$\frac{\frac{\frac{[N^\gamma \rightarrow \bullet M^\gamma P^\gamma, i, i, | -, - | \text{false}]}{[N^\gamma \rightarrow \bullet M^\gamma P^\gamma, i, i, | -, - | \text{false}], [M^\gamma \rightarrow \delta \bullet, i, k | -, - | \text{adj}]}{[N^\gamma \rightarrow M^\gamma \bullet P^\gamma, i, k, | -, - | \text{false}]}}{[N^\gamma \rightarrow M^\gamma \bullet P^\gamma, i, k, | -, - | \text{false}], [P^\gamma \rightarrow \nu \bullet, k, j | p, q | \text{adj}]}{[N^\gamma \rightarrow M^\gamma P^\gamma \bullet, i, j, | p, q | \text{false}]}}$$

- un paso $\mathcal{D}_{\text{CYK}'}^{\text{NoDom}}$ es equivalente a la secuencia formada por un paso $\mathcal{D}_{\text{ECYK}}^{\text{Init}}$ y dos pasos $\mathcal{D}_{\text{ECYK}}^{\text{Comp}}$:

$$\frac{\frac{\frac{[N^\gamma \rightarrow \bullet M^\gamma P^\gamma, i, i, | -, - | \text{false}]}{[N^\gamma \rightarrow \bullet M^\gamma P^\gamma, i, i, | -, - | \text{false}], [M^\gamma \rightarrow \delta \bullet, i, k | -, - | \text{adj}]}{[N^\gamma \rightarrow M^\gamma \bullet P^\gamma, i, k, | -, - | \text{false}]}}{[N^\gamma \rightarrow M^\gamma \bullet P^\gamma, i, k, | -, - | \text{false}], [P^\gamma \rightarrow \nu \bullet, k, j | -, - | \text{adj}]}{[N^\gamma \rightarrow M^\gamma P^\gamma \bullet, i, j, | -, - | \text{false}]}}$$

- un paso $\mathcal{D}_{\text{CYK}'}^{\text{Unary}}$ es equivalente a la secuencia formada por un paso $\mathcal{D}_{\text{ECYK}}^{\text{Init}}$ y un paso $\mathcal{D}_{\text{ECYK}}^{\text{Comp}}$:

$$\frac{\frac{[N^\gamma \rightarrow \bullet M^\gamma, i, i, | -, - | \text{false}]}{[N^\gamma \rightarrow \bullet M^\gamma, i, i, | -, - | \text{false}], [M^\gamma \rightarrow \delta \bullet, i, j | -, - | \text{adj}]}{[N^\gamma \rightarrow M^\gamma \bullet, i, j, | -, - | \text{false}]}}$$

- $\mathcal{D}_{\text{ECYK}}^{\text{Adj}} = \mathcal{D}_{\text{CYK}'}^{\text{adj}}$.

El esquema de análisis sintáctico **ECYK** está definido para gramáticas de adjunción de árboles en las cuales ningún nodo puede tener más de dos descendientes mientras que el esquema de análisis **buE₁** está definido para cualquier TAG. Es fácil mostrar que **ECYK** $\xrightarrow{\text{ext}}$ **buE₁** puesto que **ECYK**(\mathcal{T}) = **buE₁**(\mathcal{T}) es cierto para toda gramática de adjunción de árboles, ya que por definición $\mathbb{P}_{\text{ECYK}} = \mathbb{P}_{\text{buE}_1}$. \square

Se puede eliminar el elemento de los ítems que indica si se ha realizado una adjunción en el nodo situado en el lado izquierdo de la producción contenida en dicho ítem, quedando definido el conjunto de ítems

$$\left\{ \begin{array}{ll} [N^\gamma \rightarrow \delta \bullet \nu, i, j | p, q] & \delta \xrightarrow{*} a_{i+1} \dots a_p \mathbf{F}^\gamma a_{q+1} \dots a_j \xrightarrow{*} a_{i+1} \dots a_j \quad \text{sii } (p, q) \neq (-, -) \\ & \delta \xrightarrow{*} a_{i+1} \dots a_j \quad \text{sii } (p, q) = (-, -) \end{array} \right\}$$

Para ello es necesario aplicar un filtro al esquema de análisis sintáctico **buE₁**, consistente en la contracción de pasos deductivos: puesto que el ítem generado por un paso deductivo de tipo

Adj sólo podrá ser utilizado en un paso de tipo Comp para avanzar el punto en la producción que predijo el no terminal del lado izquierdo de su producción, podemos crear un nuevo tipo AdjComp de paso deductivo y eliminar los pasos de tipo Adj. Obtendremos así el esquema de análisis **buE** cuyo sistema de análisis \mathbb{P}_{buE} mostramos a continuación.

Esquema de análisis sintáctico 3.3 El sistema de análisis \mathbb{P}_{buE} que se corresponde con una nueva versión de la extensión del algoritmo de Earley ascendente para TAG, dada una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\text{buE}} = \left\{ [N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q] \mid \begin{array}{l} N^\gamma \rightarrow \delta \nu \in \mathcal{P}(\gamma), \gamma \in \mathbf{I} \cup \mathbf{A}, \\ 0 \leq i \leq j, (p, q) \leq (i, j) \end{array} \right\}$$

$$\mathcal{D}_{\text{buE}}^{\text{Init}} = \overline{[N^\gamma \rightarrow \bullet \delta, i, i \mid -, -]}$$

$$\mathcal{D}_{\text{buE}}^{\text{Foot}} = \overline{[\mathbf{F}^\beta \rightarrow \perp \bullet, i, j \mid i, j]}$$

$$\mathcal{D}_{\text{buE}}^{\text{Scan}} = \frac{[N^\gamma \rightarrow \delta \bullet a\nu, i, j \mid p, q], [a, j, j + 1]}{[N^\gamma \rightarrow \delta a \bullet \nu, i, j + 1 \mid p, q]}$$

$$\mathcal{D}_{\text{buE}}^{\text{Comp}} = \frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p, q], [M^\gamma \rightarrow \nu \bullet, k, j \mid p', q']}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p \cup p', q \cup q']} \quad \text{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{buE}}^{\text{AdjComp}} = \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, k, j \mid l, m], \\ [M^\gamma \rightarrow \nu \bullet, l, m \mid p', q'], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p, q], \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p \cup p', q \cup q']} \quad \beta \in \mathbf{A}, \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{buE}} = \mathcal{D}_{\text{buE}}^{\text{Init}} \cup \mathcal{D}_{\text{buE}}^{\text{Foot}} \cup \mathcal{D}_{\text{buE}}^{\text{Scan}} \cup \mathcal{D}_{\text{buE}}^{\text{Comp}} \cup \mathcal{D}_{\text{buE}}^{\text{AdjComp}}$$

$$\mathcal{F}_{\text{buE}} = \{ [\top \rightarrow \mathbf{R}^\alpha \bullet, 0, n \mid -, -] \mid \alpha \in \mathbf{I} \}$$

§

Como se puede observar, el paso $\mathcal{D}_{\text{buE}}^{\text{AdjComp}}$ solamente permite que un árbol auxiliar sea adjuntado en un nodo de un árbol elemental, pues una vez realizada la adjunción, el punto de la producción correspondiente al padre del nodo de adjunción se mueve a la derecha mientras que la producción del nodo de adjunción permanece sin cambios: si posteriormente otro árbol auxiliar es adjuntado en dicho nodo, representará una ambigüedad en el análisis sintáctico de la cadena de entrada, no la adjunción simultánea de dos árboles auxiliares en un mismo nodo [175]. En la figura 3.3 se muestra una representación gráfica de la aplicación de este paso deductivo para su caso más complejo, aquel en el que γ es un árbol auxiliar y el nodo de adjunción pertenece a su espina.

La complejidad temporal del esquema de análisis **buE** con respecto a la cadena de entrada es aparentemente $\mathcal{O}(n^7)$ puesto que son 7 los índices involucrados en el paso deductivo $\mathcal{D}_{\text{buE}}^{\text{AdjComp}}$: i, j, k, l, m y bien p y q o bien p' y q' . Sin embargo, podemos observar que los índices l y m sólo son necesarios para relacionar los dos primeros antecedentes y son irrelevantes para el resto de los ítems involucrados en el paso deductivo. Por tanto, mediante la aplicación parcial o *currificación* del paso deductivo $\mathcal{D}_{\text{buE}}^{\text{AdjComp}}$ la complejidad del mismo se reduce a $\mathcal{O}(n^6)$. De forma equivalente, también se podría incluir dicha aplicación parcial en el propio esquema de análisis, sustituyendo el paso $\mathcal{D}_{\text{buE}}^{\text{AdjComp}}$ por otros dos, uno que combine los dos primeros ítems antecedentes y genere un ítem intermedio y otro que combine dicho ítem con el tercer antecedente del paso original para obtener el ítem resultado. Sin embargo, con ello oscureceríamos la comprensión del algoritmo descrito en el esquema y estaríamos vulnerando la filosofía de los esquemas de análisis, ya que la aplicación parcial es un detalle de implementación del que nos debemos abstraer.

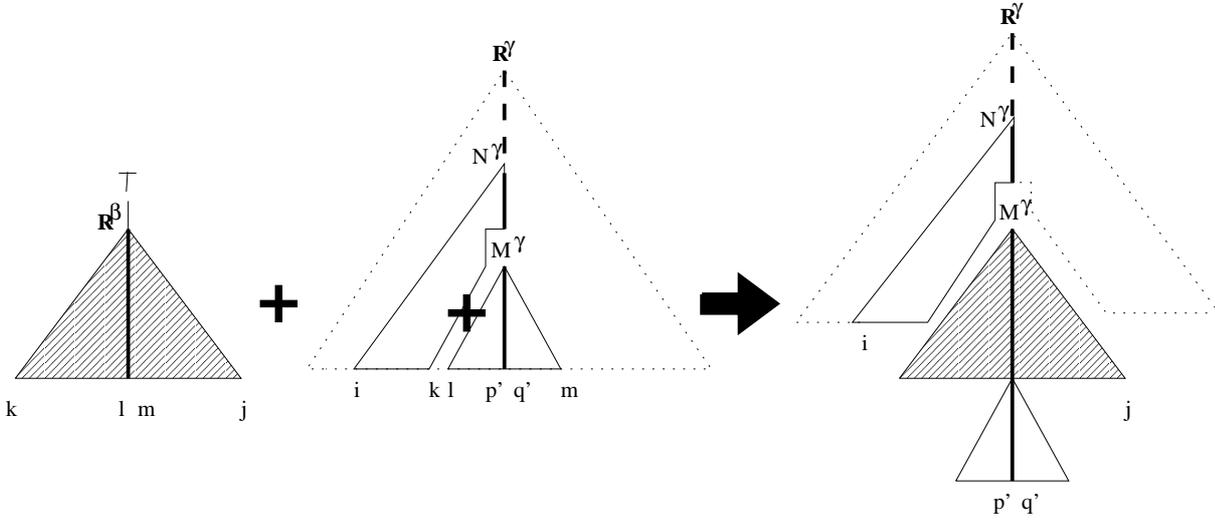


Figura 3.3: Descripción gráfica de un paso $\mathcal{D}_{\text{buE}}^{\text{AdjComp}}$

Proposición 3.2 $\text{buE}_1 \xrightarrow{\text{df}} \text{buE}'_1 \xrightarrow{\text{sc}} \text{buE}_2 \xrightarrow{\text{ic}} \text{buE}$.

Demostración:

Como primer paso definiremos el sistema de análisis $\mathbb{P}_{\text{buE}'_1}$ para una gramática de adjunción \mathcal{T} y una cadena de entrada $a_1 \dots a_n$, en el cual el paso Comp ha sido desdoblado en Comp^1 y Comp^2 , el primero de los cuales se encarga de avanzar el punto en aquellos casos en que el nodo sobre el que se avanza no ha sido utilizado como nodo de adjunción, mientras que el segundo avanza el punto sobre un nodo de adjunción.

$$\mathcal{I}_{\text{buE}'_1} = \mathcal{I}_{\text{buE}_1}$$

$$\mathcal{D}_{\text{buE}'_1}^{\text{Init}} = \mathcal{D}_{\text{buE}_1}^{\text{Init}}$$

$$\mathcal{D}_{\text{buE}'_1}^{\text{Foot}} = \mathcal{D}_{\text{buE}_1}^{\text{Foot}}$$

$$\mathcal{D}_{\text{buE}'_1}^{\text{Scan}} = \mathcal{D}_{\text{buE}_1}^{\text{Scan}}$$

$$\mathcal{D}_{\text{buE}'_1}^{\text{Comp}^1} = \frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{false}]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]} \quad \text{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{buE}'_1}^{\text{Comp}^2} = \frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{true}]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]} \quad \exists \beta \in \mathbf{A}, \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{bu}E_1}^{\text{Adj}} = \frac{[\top \rightarrow \mathbf{R}^\beta \bullet, k, j \mid k', j' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k', j' \mid p, q \mid \text{false}]}{[M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{true}]} \quad \beta \in \mathbf{A}, \beta \in \text{adj}(\gamma)$$

$$\begin{aligned} \mathcal{D}_{\text{bu}E'_1} &= \mathcal{D}_{\text{bu}E'_1}^{\text{Init}} \cup \mathcal{D}_{\text{bu}E'_1}^{\text{Foot}} \cup \mathcal{D}_{\text{bu}E'_1}^{\text{Scan}} \cup \mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^1} \cup \mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^2} \cup \mathcal{D}_{\text{bu}E'_1}^{\text{Adj}} \\ \mathcal{F}_{\text{bu}E'_1} &= \mathcal{F}_{\text{bu}E_1} \end{aligned}$$

Los pasos deductivos $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^1}$ y $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^2}$ son el resultado de aplicar un filtro dinámico al paso $\mathcal{D}_{\text{bu}E_1}^{\text{Comp}}$. En el caso de $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^1}$ dicho filtro consiste en comprobar si el último elemento de los ítems contiene el valor false, mientras que en el caso de $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^2}$ consiste en comprobar que su valor es true.

Para demostrar que $\text{bu}E_1 \xrightarrow{\text{df}} \text{bu}E'_1$, deberemos demostrar que para todo esquema de análisis $\mathbb{P}_{\text{bu}E_1}$ y $\mathbb{P}_{\text{bu}E'_1}$ se cumple que $\mathcal{I}_{\text{bu}E_1} \supseteq \mathcal{I}_{\text{bu}E'_1}$ y $\vdash_{\text{bu}E_1} \supseteq \vdash_{\text{bu}E'_1}$. Lo primero se cumple puesto que por definición $\mathcal{I}_{\text{bu}E_1} = \mathcal{I}_{\text{bu}E'_1}$. Para lo segundo debemos mostrar que $\vdash_{\text{bu}E_1} \supseteq \mathcal{D}_{\text{bu}E'_1}$. Para ello sólo necesitamos considerar aquellos pasos de $\mathbb{P}_{\text{bu}E'_1}$ a los que se les ha aplicado el filtro dinámico, pues los restantes pasos permanecen inalterados:

- un paso $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^1}$ es equivalente a la aplicación del paso $\mathcal{D}_{\text{bu}E_1}^{\text{Comp}}$:

$$\frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{false}]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]}$$

- un paso $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^2}$ es equivalente a la aplicación del paso $\mathcal{D}_{\text{bu}E_1}^{\text{Comp}}$:

$$\frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{true}]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]}$$

Sólo se pueden generar ítems que tengan el último componente con valor true como consecuencia de la aplicación de un paso $\mathcal{D}_{\text{bu}E_1}^{\text{Adj}}$ y a su vez estos ítems sólo puede ser utilizados como antecedentes en un paso $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^2}$. Por tanto podemos juntar ambos pasos en uno sólo. A continuación definimos el sistema de análisis $\mathbb{P}_{\text{bu}E_2}$ para una gramática de adjunción \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ que incorpora esta transformación.

$$\begin{aligned} \mathcal{I}_{\text{bu}E_2} &= \mathcal{I}_{\text{bu}E'_1} \\ \mathcal{D}_{\text{bu}E_2}^{\text{Init}} &= \mathcal{D}_{\text{bu}E'_1}^{\text{Init}} \\ \mathcal{D}_{\text{bu}E_2}^{\text{Foot}} &= \mathcal{D}_{\text{bu}E'_1}^{\text{Foot}} \\ \mathcal{D}_{\text{bu}E_2}^{\text{Scan}} &= \mathcal{D}_{\text{bu}E'_1}^{\text{Scan}} \\ \mathcal{D}_{\text{bu}E_2}^{\text{Comp}} &= \mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^1} \\ \mathcal{D}_{\text{bu}E_2}^{\text{Adj}} &= \frac{[\top \rightarrow \mathbf{R}^\beta \bullet, k, j \mid k', j' \mid \text{false}], \\ & [M^\gamma \rightarrow v \bullet, k', j' \mid p, q \mid \text{false}], \\ & [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]} \quad \beta \in \mathbf{A}, \beta \in \text{adj}(\gamma) \\ \mathcal{D}_{\text{bu}E_2} &= \mathcal{D}_{\text{bu}E_2}^{\text{Init}} \cup \mathcal{D}_{\text{bu}E_2}^{\text{Foot}} \cup \mathcal{D}_{\text{bu}E_2}^{\text{Scan}} \cup \mathcal{D}_{\text{bu}E_2}^{\text{Comp}} \cup \mathcal{D}_{\text{bu}E_2}^{\text{Adj}} \\ \mathcal{F}_{\text{bu}E_2} &= \mathcal{F}_{\text{bu}E_1} \end{aligned}$$

Para demostrar que el sistema de análisis sintáctico $\mathbb{P}_{\text{bu}E_2}$ es el resultado de aplicar una contracción de pasos al sistema de análisis $\mathbb{P}_{\text{bu}E'_1}$ tenemos que demostrar que $\mathcal{I}_{\text{bu}E'_1} \supseteq \mathcal{I}_{\text{bu}E_2}$ y que $\vdash_{\text{bu}E'_1}^* \supseteq \vdash_{\text{bu}E_2}^*$. Lo primero es cierto por definición, puesto que $\mathcal{I}_{\text{bu}E'_1} = \mathcal{I}_{\text{bu}E_2}$. Para lo

segundo es suficiente con demostrar que $\vdash_{\text{bu}E'_1}^* \supseteq \mathcal{D}_{\text{bu}E_2}$. Puesto que los demás pasos deductivos son idénticos en ambos esquemas, debemos mostrar únicamente que $\vdash_{\text{bu}E'_1}^* \supseteq \mathcal{D}_{\text{bu}E_2}^{\text{Adj}}$, tarea que no presenta complicación alguna puesto que como se ha mencionado anteriormente, un paso $\mathcal{D}_{\text{bu}E_2}^{\text{Adj}}$ es equivalente a la aplicación consecutiva de un paso $\mathcal{D}_{\text{bu}E'_1}^{\text{Adj}}$ y un paso $\mathcal{D}_{\text{bu}E'_1}^{\text{Comp}^2}$:

$$\frac{[\top \rightarrow \mathbf{R}^\beta \bullet, k, j \mid k', j' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k', j' \mid p, q \mid \text{false}]}{[M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{true}]}$$

$$\frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p', q' \mid \text{false}], [M^\gamma \rightarrow v \bullet, k, j \mid p, q \mid \text{true}]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p' \cup p, q' \cup q \mid \text{false}]}$$

Finalmente observamos que todos los ítems del sistema de análisis $\mathbb{P}_{\text{bu}E_2}$ tienen la forma $[N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid \text{false}]$. Puesto que el último elemento tiene un valor constante, su eliminación no producirá ningún efecto con respecto al conjunto de ítems que puedan ser generados. El sistema de análisis $\mathbb{P}_{\text{bu}E}$ se obtiene precisamente al aplicar esta transformación al sistema de análisis $\mathbb{P}_{\text{bu}E_2}$. Dicha transformación constituye un caso peculiar de contracción de ítems puesto que no se rompe un ítem en varios, sino que se establece una relación biyectiva entre los conjuntos $\mathcal{I}_{\text{bu}E_2}$ y $\mathcal{I}_{\text{bu}E}$. Vemos que la relación de contracción de ítems se mantiene entre los dos esquemas de análisis puesto que definiendo la función

$$f([N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid \text{false}]) = [N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q]$$

obtenemos que $f(\mathcal{I}_{\text{bu}E_2}) = \mathcal{I}_{\text{bu}E}$ y $f(\Delta_{\text{bu}E_2}) = \Delta_{\text{bu}E}$ por inducción en la longitud de las secuencias de derivación. \square

3.4 La propiedad del prefijo válido en los algoritmos de análisis

Los analizadores sintácticos que satisfacen la *propiedad del prefijo válido* (VPP) garantizan que, en tanto que leen la cadena de entrada de izquierda a derecha, las subcadenas leídas son prefijos válidos del lenguaje definido por la gramática. A la propiedad del prefijo válido también se la denomina a veces propiedad de detección de errores porque implica que los errores son detectados tan pronto como es posible en una lectura de la cadena de entrada de izquierda a derecha. La carencia de la propiedad del prefijo válido no significa que los errores no sean detectados, sino simplemente que estos lo serán más tarde. En contra de la idea mantenida en ciertos momentos [168], la propiedad del prefijo válido es independiente de la propiedad de análisis *en línea* [169].

Más formalmente, un analizador sintáctico satisface la propiedad del prefijo válido si al leer la subcadena $a_1 \dots a_k$ de la cadena de entrada $a_1 \dots a_k a_{k+1} \dots a_n$ garantiza que hay una cadena $b_1 \dots b_m$, donde b_i no tiene porque formar parte de la cadena de entrada, tal que $a_1 \dots a_k b_1 \dots b_m$ es una cadena válida del lenguaje.

El mantenimiento de la propiedad del prefijo válido exige ir reconociendo los posibles árboles derivados de forma prefija. Este recorrido consta de dos fases, una descendente que dado un nodo expande sus nodos hijos y una ascendente que agrupa los nodos hijos para indicar el reconocimiento del nodo padre. Cuando se desea mantener la propiedad del prefijo válido estas dos fases deben actuar coordinadamente. La fase descendente debe ser además restringida, para evitar expansiones que lleven al reconocimiento de prefijos no válidos [169].

En el caso de gramáticas independientes del contexto, existen numerosos algoritmos de análisis que preservan la propiedad del prefijo válido (por ejemplo, Earley) y que muestran una complejidad en el peor caso igual a aquellos algoritmos que no la preservan (por ejemplo, CYK). Esto se debe a que la operación de sustitución puede ser aplicada de forma totalmente

independiente del contexto, lo que conlleva que el conjunto de caminos de una gramática independiente del contexto sea un lenguaje regular. Como consecuencia, el mantenimiento de la propiedad del prefijo válido se puede asegurar sin tener que aplicar restricciones complejas sobre la fase descendente de los algoritmos.

En el caso de las gramáticas de adjunción, la operación de adjunción no es completamente independiente del contexto en el cual se aplica. Durante el reconocimiento de un árbol derivado en forma prefija, la expansión de un nodo puede depender de operaciones de adjunción previamente realizadas en la parte recorrida del árbol. Esta dependencia del contexto conlleva que el conjunto de caminos ya no sea un lenguaje regular sino un lenguaje independiente del contexto [230, 217]. Un algoritmo básicamente ascendente (p.ej.: de tipo CYK, algunas variantes de tipo Earley) puede simplemente hacer uso de una pila para ir almacenando las dependencias indicadas por el lenguaje que define el conjunto de caminos. Con ello se conseguiría un algoritmo de análisis correcto pero sin la propiedad del prefijo válido. Para preservar esta propiedad es necesario disponer de una fase descendente, que también tendría que disponer de una pila para satisfacer las restricciones impuestas por el lenguaje que define el conjunto de caminos. Schabes [169] argumentaba que entonces, al tener que coordinar las pilas de la fase ascendente y descendente, la complejidad del algoritmo de análisis sintáctico resultante aumentaría. Sin embargo, el algoritmo descrito por Nederhof en [125] mantiene la propiedad del prefijo válido con una complejidad $\mathcal{O}(n^6)$, igual a la de aquellos algoritmos que no la mantienen.

3.5 Algoritmos de tipo Earley sin la propiedad del prefijo válido

Mediante la aplicación de un filtrado dinámico a los esquemas de análisis sintácticos anteriores es posible obtener un esquema de análisis sintáctico de un algoritmo al estilo del de Earley pero extendido al caso de gramáticas de adjunción de árboles. Concretamente, el filtrado que se realiza es el siguiente:

- El paso deductivo $\mathcal{D}_{\text{buE}}^{\text{Init}}$ sólo contendrá producciones cuyo lado izquierdo corresponda con la raíz de un árbol inicial.
- Un conjunto de pasos deductivos predictivos controlan la generación de nuevos ítems tratando de limitarla únicamente a aquellos que puedan resultar útiles en el proceso de análisis.

Los algoritmos descritos en esta sección no preservan la propiedad del prefijo válido puesto que la fase predictiva no es lo suficientemente restrictiva como para evitar que las predicciones realizadas durante el análisis del pie de un árbol no conlleven el análisis de subárboles no válidos.

Una primera aproximación a un esquema de análisis sintáctico para un algoritmo de tipo Earley para gramáticas de adjunción consiste en transformar el esquema de análisis sintáctico **buE**. Al nuevo esquema de análisis, que presenta ciertas semejanzas con los algoritmos descritos por Schabes en [168, 169, 172], lo denominaremos **E** y su correspondiente sistema de análisis \mathbb{P}_{E} se define a continuación.

Esquema de análisis sintáctico 3.4 El sistema de análisis \mathbb{P}_{E} que se corresponde con una versión del algoritmo de Earley para TAG sin la propiedad del prefijo válido, dada una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\text{E}} = \mathcal{I}_{\text{buE}}$$

$$\mathcal{D}_{\text{E}}^{\text{Init}} = \frac{}{[\top \rightarrow \bullet \mathbf{R}^\alpha, 0, 0 \mid -, -]} \quad \alpha \in \mathbf{I}$$

$$\mathcal{D}_E^{\text{Scan}} = \mathcal{D}_{\text{buE}}^{\text{Scan}}$$

$$\mathcal{D}_E^{\text{Pred}} = \frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[M^\gamma \rightarrow \bullet v, j, j \mid -, -]} \quad \text{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_E^{\text{Comp}} = \mathcal{D}_{\text{buE}}^{\text{Comp}}$$

$$\mathcal{D}_E^{\text{AdjPred}} = \frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[\top \rightarrow \bullet \mathbf{R}^\beta, j, j \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_E^{\text{FootPred}} = \frac{[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -]}{[M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_E^{\text{FootComp}} = \frac{[M^\gamma \rightarrow \delta \bullet, k, l \mid p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -]}{[\mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_E^{\text{AdjComp}} = \mathcal{D}_{\text{buE}}^{\text{AdjComp}} = \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [M^\gamma \rightarrow v \bullet, k, l \mid p, q], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p \cup p', q \cup q']} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_E = \mathcal{D}_E^{\text{Init}} \cup \mathcal{D}_E^{\text{Scan}} \cup \mathcal{D}_E^{\text{Pred}} \cup \mathcal{D}_E^{\text{Comp}} \cup \mathcal{D}_E^{\text{AdjPred}} \cup \mathcal{D}_E^{\text{FootPred}} \cup \mathcal{D}_E^{\text{FootComp}} \cup \mathcal{D}_E^{\text{AdjComp}}$$

$$\mathcal{F}_E = \mathcal{F}_{\text{buE}}$$

§

El paso $\mathcal{D}_E^{\text{AdjComp}}$, aunque es idéntico al del sistema de análisis $\mathbb{P}_{\text{buE}_1}$ se repite aquí para facilitar la comprensión del algoritmo. Como se ha dicho anteriormente, el esquema de análisis **E** presenta ciertas similitudes con los algoritmos descritos por Schabes en [168, 169, 172]. Podemos establecer la siguiente relación entre los pasos deductivos de **E** y los procesos definidos en dichos algoritmos, tal y como se muestra en la tabla 3.1. La razón del cambio de nombre en las pasos encargados de la adjunción se debe a que creemos que los nombres utilizados por Schabes pueden llevar a confusión, pues se tiende a pensar que cada *Predictor* está asociado con el *Completor* correspondiente a su mismo lado, cuando no es así. Es por ello que hemos decidido emparejar los pasos deductivos por las funciones que realizan: predicción-compleción de adjunción y predicción-compleción de pie.

Respecto al comportamiento del algoritmo, diremos que el análisis comienza creando un ítem correspondiente a una producción del nodo raíz de un árbol inicial con el punto situado en el extremo izquierdo. Posteriormente, los pasos $\mathcal{D}_E^{\text{Pred}}$ y $\mathcal{D}_E^{\text{Comp}}$ se encargan de ir recorriendo el árbol de modo descendente y ascendente, respectivamente, de modo similar a como actúa el algoritmo de Earley para gramáticas independientes del contexto. Se puede predecir la adjunción de un árbol β en un nodo de un árbol elemental γ mediante la aplicación de un paso $\mathcal{D}_E^{\text{AdjPred}}$,

Pasos deductivos de \mathbf{E}	Algoritmo de Schabes
$\mathcal{D}_E^{\text{Init}}$	<i>Initial item</i>
$\mathcal{D}_E^{\text{Scan}}$	Scanner
$\mathcal{D}_E^{\text{Pred}}$	Move Dot Down
$\mathcal{D}_E^{\text{Comp}}$	Move Dot Up
$\mathcal{D}_E^{\text{AdjPred}}$	Left Predictor
$\mathcal{D}_E^{\text{FootPred}}$	Left Completor
$\mathcal{D}_E^{\text{FootComp}}$	Right Predictor
$\mathcal{D}_E^{\text{AdjComp}}$	Right Completor

Tabla 3.1: Relación entre \mathcal{D}_E y el algoritmo tipo Earley sin VPP de Schabes

con lo que se comienza el análisis del árbol β . Una vez alcanzado el pie de dicho árbol auxiliar, deberemos retomar el análisis de γ , concretamente del subárbol que pende del nodo de adjunción. El problema es que al no conocer en qué nodo de qué árbol elemental se ha producido la adjunción, deberemos predecir todos los posibles nodos donde esté permitida la adjunción de β , predicción realizada por un paso deductivo $\mathcal{D}_E^{\text{FootPred}}$. Es la predicción que se realiza en los pasos $\mathcal{D}_E^{\text{FootPred}}$ lo que provoca que el algoritmo no posea la propiedad del prefijo válido, puesto que se puede comenzar a analizar una parte de la cadena de entrada que no es gramatical, al predecir un subárbol que no se corresponde con el árbol desde el que se realizó la adjunción.

Una vez terminado de analizar todo el subárbol predicho por un paso $\mathcal{D}_E^{\text{FootPred}}$, deberemos retomar el análisis del árbol auxiliar β a partir del pie, tarea encomendada a los pasos deductivos $\mathcal{D}_E^{\text{FootComp}}$. Una vez terminado de analizar completamente el árbol auxiliar β , deberemos concluir la operación de adjunción aplicando un paso $\mathcal{D}_E^{\text{AdjComp}}$. Es en estos pasos en los que se verifica que el subárbol escindido del nodo de adjunción y el árbol auxiliar han sido correctamente reconstruidos. Las adjunciones simultáneas sobre un mismo nodo [175] son evitadas por los pasos $\mathcal{D}_E^{\text{AdjComp}}$ puesto que cuando se ha terminado de recorrer completamente el árbol auxiliar, se verifica que se ha analizado la parte correspondiente al subárbol del nodo de adjunción y se avanza el punto de la producción que contiene a este, sin cambiar el ítem correspondiente al nodo de adjunción. Si posteriormente otro árbol auxiliar es adjuntado en dicho nodo, representará una ambigüedad en el análisis sintáctico de la cadena de entrada pero no la adjunción simultánea de dos árboles auxiliares en un mismo nodo.

Proposición 3.3 $\text{buE} \stackrel{\text{df}}{\implies} \mathbf{E}$.

Demostración:

Para demostrar que el esquema de análisis sintáctico \mathbf{E} es el resultado de aplicar un filtrado dinámico al esquema de análisis buE , debemos demostrar que $\mathcal{I}_{\text{buE}} \supseteq \mathcal{I}_E$ y que $\vdash_{\text{buE}} \supseteq \vdash_E$ para los sistemas de análisis sintáctico \mathbb{P}_{buE} y \mathbb{P}_E . Lo primero es cierto por definición puesto que $\mathcal{I}_{\text{buE}} = \mathcal{I}_E$. Respecto a lo segundo, es suficiente con mostrar que $\vdash_{\text{buE}} \supseteq \mathcal{D}_E$.

Los pasos $\mathcal{D}_E^{\text{Scan}}$, $\mathcal{D}_E^{\text{Comp}}$ y $\mathcal{D}_E^{\text{AdjComp}}$ son idénticos a sus homónimos del sistema \mathbb{P}_{buE} y los pasos $\mathcal{D}_E^{\text{Init}}$ generan un subconjunto de los ítems generados por $\mathcal{D}_{\text{buE}}^{\text{Init}}$. Respecto a los otros pasos:

- Dado un paso $\frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j | p, q]}{[M^\gamma \rightarrow \bullet \nu, j, j | -, -]} \in \mathcal{D}_E^{\text{Pred}}$ existe un paso $\overline{[M^\gamma \rightarrow \bullet \nu, j, j | -, -]} \in \mathcal{D}_{\text{buE}}^{\text{Init}}$ y por tanto existe la inferencia

$$[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j | p, q] \vdash_{\text{buE}} [M^\gamma \rightarrow \bullet \nu, j, j | -, -]$$

- Dado un paso $\frac{[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j | p, q]}{[\top \rightarrow \bullet \mathbf{R}^\beta, j, j | -, -]} \in \mathcal{D}_E^{\text{AdjPred}}$ existe un paso $\overline{[\top \rightarrow \bullet \mathbf{R}^\beta, j, j | -, -]} \in \mathcal{D}_{\text{buE}}^{\text{Init}}$ y por tanto existe la inferencia

$$[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j | p, q] \vdash_{\text{buE}} [\top \rightarrow \bullet \mathbf{R}^\beta, j, j | -, -]$$

- Dado un paso $\frac{[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k | -, -]}{[M^\gamma \rightarrow \bullet \delta, k, k | -, -]} \in \mathcal{D}_E^{\text{FootPred}}$ existe un paso $\overline{[M^\gamma \rightarrow \bullet \delta, k, k | -, -]} \in \mathcal{D}_{\text{buE}}^{\text{Init}}$ y por tanto existe la inferencia

$$[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k | -, -] \vdash_{\text{buE}} [M^\gamma \rightarrow \bullet \delta, k, k | -, -]$$

- Dado un paso $\frac{[M^\gamma \rightarrow \delta \bullet, k, l | p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k | -, -]}{[\mathbf{F}^\beta \rightarrow \perp \bullet, k, l | k, l]} \in \mathcal{D}_E^{\text{FootComp}}$ existe un paso $\overline{[\mathbf{F}^\beta \rightarrow \perp \bullet, k, l | k, l]} \in \mathcal{D}_{\text{buE}}^{\text{Foot}}$ y por tanto existe la inferencia

$$[M^\gamma \rightarrow \delta \bullet, k, l | p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k | -, -] \vdash_{\text{buE}} [\mathbf{F}^\beta \rightarrow \perp \bullet, k, l | k, l]$$

□

La complejidad temporal del esquema de análisis sintáctico **E** con respecto a la cadena de entrada es $\mathcal{O}(n^6)$ puesto que la aparente complejidad $\mathcal{O}(n^7)$ del paso deductivo $\mathcal{D}_E^{\text{AdjComp}}$ puede reducirse a $\mathcal{O}(n^6)$ mediante la aplicación parcial o *currificación* de dicho paso, ya que los índices l y m sólo involucran a los dos primeros ítems antecedentes.

Con el fin de definir un esquema de análisis sintáctico que se corresponda con un algoritmo más cercano al espíritu del algoritmo de Earley, debemos fortalecer la fase predictiva de los esquemas de análisis anteriores, puesto que estos no utilizan toda la información que tienen a su disposición. En concreto:

- Los pasos deductivos $\mathcal{D}_E^{\text{FootPred}}$ en los que se realiza la predicción del pie no comprueban que previamente se haya iniciado la adjunción del árbol β en el nodo M^γ .
- Idem para los pasos deductivos $\mathcal{D}_E^{\text{FootComp}}$ que finalizan el reconocimiento del pie.
- Los pasos deductivos $\mathcal{D}_E^{\text{FootComp}}$ no comprueban que el árbol auxiliar predicho para adjunción en el nodo M^γ sea el mismo que el que ha provocado el reconocimiento del subárbol enraizado en dicho nodo mediante la aplicación de los pasos FootPred.

A continuación definimos un nuevo esquema de análisis **Ear** derivado a partir de **E**, sobre el que hemos aplicado las siguientes modificaciones:

- La aplicación de un filtro dinámico a los pasos deductivos FootPred y FootComp consistente en la verificación de la existencia de los ítems que representan el comienzo de la operación de adjunción en el nodo M^γ .
- La aplicación de un refinamiento al paso AdjComp, que se divide en dos pasos AdjComp¹ y AdjComp², el primero realizando las comprobaciones pertinentes en el caso de que se haya producido la adjunción de un árbol auxiliar en un nodo de la espina de otro árbol auxiliar. En la figura 3.4 se muestra una representación gráfica de la aplicación del paso deductivo AdjComp¹, el más complejo de los dos ya que el árbol γ en el que se realiza la adjunción es un árbol auxiliar y el nodo de adjunción pertenece a su espina. Las partes de

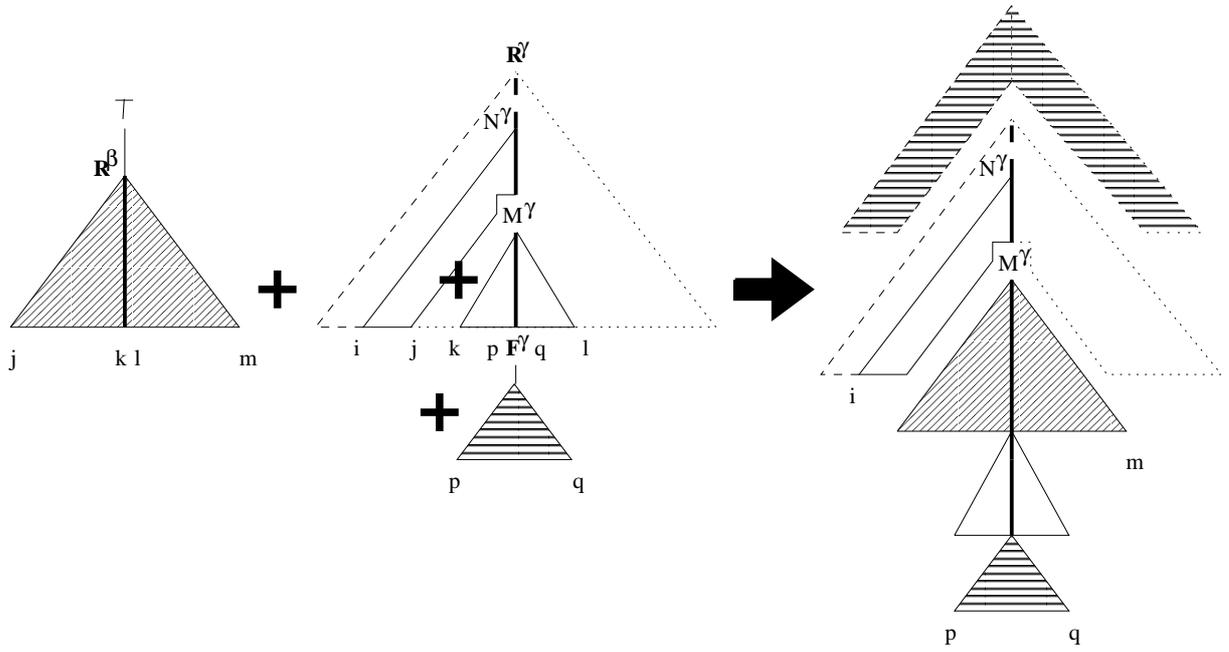


Figura 3.4: Descripción gráfica de un paso $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1}$

los árboles involucrados que no se encuentran representados por los ítems que intervienen en el paso deductivo se muestran en línea discontinua si el algoritmo de análisis tiene que haber pasado obligatoriamente por dicha parte del árbol al menos en la fase predictiva y en línea punteada si se trata de partes que serán analizadas posteriormente.

Esquema de análisis sintáctico 3.5 El sistema de análisis \mathbb{P}_{Ear} que se corresponde con una versión del algoritmo de Earley sin la propiedad del prefijo válido con predicción fuerte, dada una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\text{Ear}} = \mathcal{I}_{\text{buE}}$$

$$\mathcal{D}_{\text{Ear}}^{\text{Init}} = \mathcal{D}_{\text{E}}^{\text{Init}}$$

$$\mathcal{D}_{\text{Ear}}^{\text{Scan}} = \mathcal{D}_{\text{buE}}^{\text{Scan}}$$

$$\mathcal{D}_{\text{Ear}}^{\text{Pred}} = \mathcal{D}_{\text{E}}^{\text{Pred}}$$

$$\mathcal{D}_{\text{Ear}}^{\text{Comp}} = \mathcal{D}_{\text{buE}}^{\text{Comp}}$$

$$\mathcal{D}_{\text{Ear}}^{\text{AdjPred}} = \mathcal{D}_{\text{E}}^{\text{AdjPred}}$$

$$\mathcal{D}_{\text{Ear}}^{\text{FootPred}} = \frac{[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Ear}}^{\text{FootComp}} = \frac{[M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q']}{[\mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]} \quad \begin{array}{l} \beta \in \text{adj}(M^\gamma), \\ p \cup p' \text{ y } q \cup q' \text{ está definido} \end{array}$$

$$\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1} = \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [M^\gamma \rightarrow v \bullet, k, l \mid p, q], \\ [\mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^2} = \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [M^\gamma \rightarrow v \bullet, k, l \mid -, -], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Ear}} = \mathcal{D}_{\text{Ear}}^{\text{Init}} \cup \mathcal{D}_{\text{Ear}}^{\text{Scan}} \cup \mathcal{D}_{\text{Ear}}^{\text{Pred}} \cup \mathcal{D}_{\text{Ear}}^{\text{Comp}} \cup \mathcal{D}_{\text{Ear}}^{\text{AdjPred}} \cup \mathcal{D}_{\text{Ear}}^{\text{FootPred}} \cup \mathcal{D}_{\text{Ear}}^{\text{FootComp}} \cup \mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1} \cup \mathcal{D}_{\text{Ear}}^{\text{AdjComp}^2}$$

$$\mathcal{F}_{\text{Ear}} = \mathcal{F}_{\text{buE}}$$

§

Proposición 3.4 $\mathbf{E} \xrightarrow{\text{sr}} \mathbf{E}' \xrightarrow{\text{df}} \mathbf{E}_{\text{ar}}$.

Demostración:

Como primer paso definiremos el esquema de análisis \mathbf{E}' que se obtiene a partir de \mathbf{E} simplemente rompiendo el conjunto de pasos deductivos $\mathcal{D}_{\text{E}}^{\text{AdjComp}}$ en dos conjuntos $\mathcal{D}_{\text{E}'}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{E}'}^{\text{AdjComp}^2}$. El sistema de análisis $\mathbb{P}_{\text{E}'}$ sería por tanto el siguiente:

$$\begin{aligned} \mathcal{I}_{\text{E}'} &= \mathcal{I}_{\text{buE}} \\ \mathcal{D}_{\text{E}'}^{\text{Init}} &= \mathcal{D}_{\text{E}}^{\text{Init}} \\ \mathcal{D}_{\text{E}'}^{\text{Scan}} &= \mathcal{D}_{\text{buE}}^{\text{Scan}} \\ \mathcal{D}_{\text{E}'}^{\text{Pred}} &= \mathcal{D}_{\text{E}}^{\text{Pred}} \\ \mathcal{D}_{\text{E}'}^{\text{Comp}} &= \mathcal{D}_{\text{buE}}^{\text{Comp}} \\ \mathcal{D}_{\text{E}'}^{\text{AdjPred}} &= \mathcal{D}_{\text{E}}^{\text{AdjPred}} \\ \mathcal{D}_{\text{E}'}^{\text{FootPred}} &= \mathcal{D}_{\text{E}}^{\text{FootPred}} \\ \mathcal{D}_{\text{E}'}^{\text{FootComp}} &= \mathcal{D}_{\text{E}}^{\text{FootComp}} \\ \mathcal{D}_{\text{E}'}^{\text{AdjComp}^1} &= \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [M^\gamma \rightarrow v \bullet, k, l \mid p, q], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]} \quad \beta \in \text{adj}(M^\gamma) \\ \mathcal{D}_{\text{E}'}^{\text{AdjComp}^2} &= \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [M^\gamma \rightarrow v \bullet, k, l \mid -, -], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']} \quad \beta \in \text{adj}(M^\gamma) \end{aligned}$$

$$\mathcal{D}_{\text{E}'} = \mathcal{D}_{\text{E}'}^{\text{Init}} \cup \mathcal{D}_{\text{E}'}^{\text{Scan}} \cup \mathcal{D}_{\text{E}'}^{\text{Pred}} \cup \mathcal{D}_{\text{E}'}^{\text{Comp}} \cup \mathcal{D}_{\text{E}'}^{\text{AdjPred}} \cup \mathcal{D}_{\text{E}'}^{\text{FootPred}} \cup \mathcal{D}_{\text{E}'}^{\text{FootComp}} \cup \mathcal{D}_{\text{E}'}^{\text{AdjComp}^1} \cup \mathcal{D}_{\text{E}'}^{\text{AdjComp}^2}$$

$$\mathcal{F}_{E'} = \mathcal{F}_{\text{bu}E}$$

Las condiciones a verificar son que $\mathcal{I}_E \subseteq \mathcal{I}_{E'}$ y que $\vdash_E \subseteq^* \vdash_{E'}$. La primera condición se verifica por la propia definición de los ítems mientras que la segunda se obtiene directamente considerando que el único cambio que se ha producido en $\mathbb{P}_{E'}$ es hacer explícita la incompatibilidad del par de índices (p, q) con el par (p', q') de los pasos $\mathcal{D}_E^{\text{AdjComp}}$: si son p' y q' los índices que están definidos, entonces el paso $\mathcal{D}_E^{\text{AdjComp}}$ se convierte en $\mathcal{D}_{E'}^{\text{AdjComp}^1}$, mientras que si son p y q los índices que están definidos, entonces el paso $\mathcal{D}_E^{\text{AdjComp}}$ se convierte en $\mathcal{D}_{E'}^{\text{AdjComp}^2}$.

Para demostrar que el esquema de análisis sintáctico **Ear** es el resultado de aplicar un filtrado dinámico al esquema de análisis **E'**, debemos demostrar que $\mathcal{I}_{E'} \supseteq \mathcal{I}_{\text{Ear}}$ y que $\vdash_{E'} \supseteq \vdash_{\text{Ear}}$ para los sistemas de análisis sintáctico $\mathbb{P}_{E'}$ y \mathbb{P}_{Ear} . Lo primero es cierto por definición puesto que $\mathcal{I}_{\text{bu}E} = \mathcal{I}_{E'} = \mathcal{I}_{\text{Ear}}$. Respecto a lo segundo es suficiente con mostrar que $\vdash_{E'} \supseteq \vdash_{\text{Ear}}$.

Los pasos $\mathcal{D}_{\text{Ear}}^{\text{Init}}$, $\mathcal{D}_{\text{Ear}}^{\text{Scan}}$, $\mathcal{D}_{\text{Ear}}^{\text{pred}}$, $\mathcal{D}_{\text{Ear}}^{\text{Comp}}$ y $\mathcal{D}_{\text{Ear}}^{\text{AdjPred}}$ son idénticos a sus homónimos del sistema $\mathbb{P}_{E'}$. Respecto a los otros pasos:

$\mathcal{D}_{\text{Ear}}^{\text{FootPred}}$: Dado un paso $\frac{[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]}$ existe un paso $\frac{[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -]}{[M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]} \in \mathcal{D}_{E'}^{\text{FootPred}}$ y por tanto existe la inferencia

$$[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q] \vdash_{E'} [M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]$$

$\mathcal{D}_{\text{Ear}}^{\text{FootComp}}$: Dado un paso $\frac{[M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q']}{[\mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]}$ existe un paso $\frac{[M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -]}{[\mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]} \in \mathcal{D}_{E'}^{\text{FootComp}}$ y por tanto existe la inferencia

$$[M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \\ \vdash_{E'} [\mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]$$

$\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1}$: Dado un paso $\frac{[\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [\mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]}$ existe un paso $\frac{[\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]} \in \mathcal{D}_{E'}^{\text{AdjComp}}$ y por tanto existe la inferencia

$$[\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], [\mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -] \\ \vdash_{E'} [N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]$$

$\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^2}$: Dado un paso $\frac{[\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [M^\gamma \rightarrow \nu \bullet, k, l \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']}$ existe un paso $\frac{[\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [M^\gamma \rightarrow \nu \bullet, k, l \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q']}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']} \in \mathcal{D}_{E'}^{\text{AdjComp}}$ y por tanto existe la inferencia

$$[\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [M^\gamma \rightarrow \nu \bullet, k, l \mid -, -], [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \\ \vdash_{E'} [N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']$$

□

La complejidad temporal del esquema de análisis sintáctico **Ear** con respecto a la longitud n de la cadena de entrada es $\mathcal{O}(n^6)$ puesto que la aparente complejidad $\mathcal{O}(n^7)$ de los pasos deductivos $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^2}$ puede reducirse a $\mathcal{O}(n^6)$ mediante la aplicación parcial o *currificación* de dichos pasos, puesto que el índice l sólo involucra a los dos primeros ítems antecedentes en cada uno de ellos.

3.6 Algoritmos de tipo Earley con la propiedad del prefijo válido

El primer algoritmo de análisis sintáctico para TAG que satisfacía la propiedad del prefijo válido fue el descrito por Schabes y Joshi en [173] y por Schabes en [168]. La principal particularidad de dicho algoritmo es que su complejidad temporal con respecto a la cadena de entrada es $\mathcal{O}(n^7)$, tal como muestran Díaz Madrigal et al. en [65], mientras que los algoritmos sin la propiedad del prefijo válido presentan una complejidad $\mathcal{O}(n^6)$. Durante mucho tiempo cobró fuerza la opinión de que aquellos algoritmos que cumplieren la propiedad del prefijo válido deberían tener una complejidad más alta que aquellos que no la cumplieren. Sin embargo, Nederhof presentó en [125] un algoritmo para el análisis de TAG que cumple la propiedad del prefijo válido⁶ y que presenta una complejidad temporal $\mathcal{O}(n^6)$. Veremos que dicho algoritmo es fácilmente derivable a partir del esquema de análisis **Ear** presentado en la sección anterior, correspondiente al algoritmo de tipo Earley sin la propiedad del prefijo válido.

El esquema de análisis sintáctico **Ear** describe un algoritmo que no cumple la propiedad del prefijo válido porque los pasos deductivos que se encargan de reconocer el nodo correspondiente al pie de un árbol auxiliar no pueden verificar la contigüidad de las fronteras del árbol al que pertenece el nodo de adjunción y del árbol auxiliar. Analicemos detalladamente dichos pasos:

- El paso deductivo $\mathcal{D}_{\text{Ear}}^{\text{FootPred}}$ puede verificar, mediante el ítem $[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]$, que existe un nodo M^γ en el que el árbol auxiliar β puede ser adjuntado para reconocer la cadena de entrada a partir de la posición j . También puede verificar, mediante el ítem $[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -]$, que se ha alcanzado el nodo pie del árbol auxiliar β . Lo que no puede verificar este paso deductivo es que el árbol al que pertenece dicho nodo pie se corresponda con la instancia de β que ha sido utilizada en la operación de adjunción que nos ocupa, pues para ello tendría que verificar que el extremo izquierdo de la frontera de la instancia de β es j , información que no es posible obtener a partir de los ítems definidos para el esquema de análisis **Ear**.
- El paso deductivo $\mathcal{D}_{\text{Ear}}^{\text{FootComp}}$ puede verificar, mediante el ítem $[N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]$, que existe un nodo M^γ en el que el árbol auxiliar β puede ser adjuntado para reconocer la cadena de entrada a partir de la posición j . También puede verificar, mediante el ítem $[\mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -]$, que se ha alcanzado el nodo pie del árbol auxiliar β en la posición k de la cadena de entrada. Por último, el ítem $[M^\gamma \rightarrow \delta \bullet, k, l \mid p, q]$ permite verificar que la frontera del subárbol enraizado en M^γ comienza en la posición k de la cadena de entrada. Pero al igual que en el caso anterior y por las mismas razones, no puede verificar que el árbol al que pertenece el nodo pie se corresponde con la instancia de β que ha sido utilizada en la operación de adjunción que nos ocupa.

En consecuencia, para obtener un esquema de análisis que se corresponda con un algoritmo del tipo Earley para TAG que posea la propiedad del prefijo válido es necesario modificar la forma de los ítems para incluir un nuevo elemento, un índice que indique la posición del extremo izquierdo de la frontera del árbol al que se refieren los nodos de cada ítem que se genere. Esta operación se corresponde con la aplicación de un refinamiento de los ítems utilizados hasta el momento. Los nuevos ítems son de la forma

$$\left\{ \begin{array}{l} [h, N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q] \mid \exists \alpha \in \mathbf{I}, \mathbf{R}^\alpha \xrightarrow{*} a_1 \dots a_h \mathbf{R}^\gamma \mu, \mathbf{R}^\gamma \xrightarrow{*} a_{h+1} \dots a_i \delta \nu \text{ y además:} \\ \delta \xrightarrow{*} a_i \dots a_p \mathbf{F}^\gamma a_{q+1} \dots a_j \xrightarrow{*} a_i \dots a_j \text{ sii } (p, q) \neq (-, -) \\ \delta \xrightarrow{*} a_i \dots a_j \text{ sii } (p, q) = (-, -) \end{array} \right\}$$

con lo cual un ítem $[N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q]$ de **Ear** se corresponde ahora con el conjunto de ítems $[h, N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q] \forall h \in [0, n]$.

⁶Para una comparación experimental entre los algoritmos de Schabes y Nederhof, consultar [63].

Una vez definidos los nuevos ítems podemos pasar a describir el esquema de análisis **Earley** que corresponde a la primera versión de un algoritmo de tipo Earley para TAG que cumple la propiedad del prefijo válido. El correspondiente sistema de análisis sintáctico se define a continuación.

Esquema de análisis sintáctico 3.6 El sistema de análisis $\mathbb{P}_{\text{Earley}}$ que se corresponde con la el algoritmo de análisis sintáctico de tipo Earley para TAG que cumple la propiedad del prefijo válido, dada una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\text{Earley}} = \left\{ [h, N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q] \mid \begin{array}{l} N^\gamma \rightarrow \delta \nu \in \mathcal{P}(\gamma), \gamma \in \mathbf{I} \cup \mathbf{A}, \\ 0 \leq h \leq i \leq j, (p, q) \leq (i, j) \end{array} \right\}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Init}} = \frac{}{\vdash [0, \top \rightarrow \bullet \mathbf{R}^\alpha, 0, 0 \mid -, -]} \quad \alpha \in \mathbf{I}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Scan}} = \frac{[h, N^\gamma \rightarrow \delta \bullet a\nu, i, j \mid p, q], [a, j, j+1]}{[h, N^\gamma \rightarrow \delta a \bullet \nu, i, j+1 \mid p, q]}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Pred}} = \frac{[h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[h, M^\gamma \rightarrow \bullet \nu, j, j \mid -, -]} \quad \mathbf{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Earley}}^{\text{Comp}} = \frac{[h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p, q], [h, M^\gamma \rightarrow \nu \bullet, k, j \mid p', q']}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p \cup p', q \cup q']} \quad \mathbf{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Earley}}^{\text{AdjPred}} = \frac{[h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[j, \top \rightarrow \bullet \mathbf{R}^\beta, j, j \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Earley}}^{\text{FootPred}} = \frac{[j, \mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[h, M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Earley}}^{\text{FootComp}} = \frac{\begin{array}{l} [h, M^\gamma \rightarrow \delta \bullet, k, l \mid p, q], \\ [j, \mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[j, \mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]} \quad \begin{array}{l} \beta \in \text{adj}(M^\gamma), \\ p \cup p' \text{ y } q \cup q' \text{ está definido} \end{array}$$

$$\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1} = \frac{\begin{array}{l} [j, \top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [h, M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], \\ [h, \mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -] \end{array}}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2} = \frac{\begin{array}{l} [j, \top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [h, M^\gamma \rightarrow \nu \bullet, k, l \mid -, -], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']} \quad \beta \in \text{adj}(M^\gamma)$$

$$\begin{aligned} \mathcal{D}_{\text{Earley}} = & \mathcal{D}_{\text{Earley}}^{\text{Init}} \cup \mathcal{D}_{\text{Earley}}^{\text{Scan}} \cup \mathcal{D}_{\text{Earley}}^{\text{Pred}} \cup \mathcal{D}_{\text{Earley}}^{\text{Comp}} \cup \mathcal{D}_{\text{Earley}}^{\text{AdjPred}} \\ & \cup \mathcal{D}_{\text{Earley}}^{\text{FootPred}} \cup \mathcal{D}_{\text{Earley}}^{\text{FootComp}} \cup \mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1} \cup \mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2} \\ \mathcal{F}_{\text{Earley}} = & \{ [0, \top \rightarrow \mathbf{R}^{\alpha \bullet}, 0, n \mid -, -] \mid \alpha \in \mathbf{I} \} \end{aligned}$$

§

En la figura 3.5 se muestra una representación gráfica de la aplicación del paso deductivo $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$. Es interesante comparar esta nueva figura con la 3.4 correspondiente al mismo paso deductivo del algoritmo de tipo Earley sin la propiedad del prefijo válido con el fin de observar cómo se restringen los árboles candidatos en la aplicación del paso deductivo. Se puede observar que la única diferencia entre ambas radica en que el extremo izquierdo del árbol γ está explícitamente indicado por el índice h en la figura 3.5, mientras que en la figura 3.4 se consideraba universalmente cuantificado.

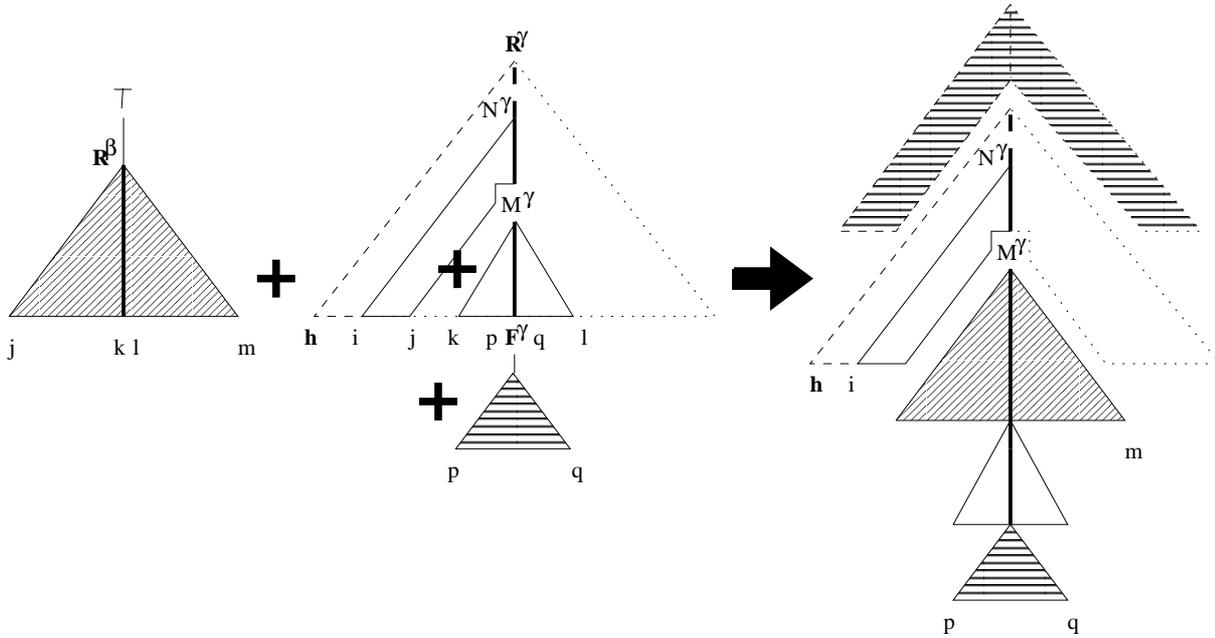


Figura 3.5: Descripción gráfica de un paso $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$

Proposición 3.5 $\text{Ear} \xrightarrow{\text{ir}} \text{Earley}$.

Demostración:

Para demostrar que el esquema de análisis **Earley** es derivable del esquema de análisis **Ear** mediante refinamiento de los ítems definiremos la siguiente función:

$$f([h, N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid adj]) = [N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q \mid adj]$$

de la cual se obtiene directamente que $\mathcal{I}_{\text{Ear}} = f(\mathcal{I}_{\text{Earley}})$ y que $\Delta_{\text{Ear}} = f(\Delta_{\text{Earley}})$ por inducción en la longitud de las secuencias de derivación. En consecuencia, $\mathbb{P}_{\text{Ear}} \xrightarrow{\text{ir}} \mathbb{P}_{\text{Earley}}$, con lo que hemos probado lo que pretendíamos. \square

Un aspecto interesante a tener en cuenta es que el ítem $[h, \mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q]$ es redundante en el paso $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$, puesto que su existencia viene implícitamente determinada por la existencia del ítem $[h, M^\gamma \rightarrow \delta \bullet, k, l \mid p, q]$, ya que en otro caso este último sería inconsistente y por consiguiente algoritmo sería incorrecto. La finalidad de su presencia en dicho conjunto de pasos deductivos, así como en $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1}$, es facilitar la transición hacia el esquema de análisis **Nederhof**. Si prescindimos de dicho ítem los pasos deductivos $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$ se podrían fundir en uno solo:

$$\mathcal{D}_{\text{Earley}}^{\text{AdjComp}} = \frac{\begin{array}{l} [j, \top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [h, M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p \cup p', q \cup q']} \quad \beta \in \text{adj}(M^\gamma)$$

Por idénticas razones, los pasos $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}^1}$ del esquema **Ear** podrían fundirse en un nuevo paso $\mathcal{D}_{\text{Ear}}^{\text{AdjComp}}$:

$$\mathcal{D}_{\text{Ear}}^{\text{AdjComp}} = \frac{\begin{array}{l} [\top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], \\ [N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p \cup p', q \cup q']} \quad \beta \in \text{adj}(M^\gamma)$$

En consecuencia, en lugar de la evolución $\mathbf{E} \xrightarrow{\text{sr}} \xrightarrow{\text{df}} \mathbf{Ear} \xrightarrow{\text{ir}} \mathbf{Earley}$ podríamos haber definido una línea evolutiva $\mathbf{E} \xrightarrow{\text{df}} \mathbf{Ear}' \xrightarrow{\text{ir}} \mathbf{Earley}' \xrightarrow{\text{sr}} \xrightarrow{\text{df}} \mathbf{Earley}$, donde \mathbf{Ear}' y \mathbf{Earley}' son como **Ear** y **Earley**, respectivamente, excepto por la sustitución de los pasos AdjComp^1 y AdjComp^2 por AdjComp . Este resultado viene a mostrar una vez más que existen varios caminos para transformar un esquema de análisis sintáctico en otro, tal como establece Sikkel en [189] para el caso de los algoritmos de análisis de gramáticas independientes del contexto y que nosotros mostramos aquí para el caso de las gramáticas de adjunción de árboles.

El algoritmo descrito por el esquema **Earley** presenta una complejidad temporal de $\mathcal{O}(n^7)$. Aunque aparentemente la utilización de 8 índices con respecto a la cadena de entrada en los pasos deductivos $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$ hace pensar en una complejidad $\mathcal{O}(n^8)$, la utilización de aplicación parcial o *currificación* en dichos pasos reduce la complejidad hasta $\mathcal{O}(n^7)$. En el caso de $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ la aplicación parcial sobre los dos primeros ítems involucra combinar 7 índices, de donde resulta la complejidad $\mathcal{O}(n^7)$ con respecto a la cadena de entrada, pero únicamente los 5 índices j, m, h, p y q forman parte del resultado intermedio puesto que son los únicos que se necesitarán en posteriores aplicaciones parciales. La siguiente aplicación parcial combina el ítem intermedio con el tercer ítem del paso deductivo, operación que involucra únicamente a los 5 ítems mencionados anteriormente, que son conservados en el resultado intermedio producido. Por último, la aplicación parcial con el cuarto ítem involucra la combinación de los 6 índices h, i, j, m, p, q . El caso de $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$ es análogo al de $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$.

El aumento de la complejidad de $\mathcal{O}(n^6)$ a $\mathcal{O}(n^7)$ se debe al índice adicional incorporado en los ítems y que indica la posición del extremo izquierdo de la frontera del árbol que se está considerando. La inclusión en los ítems de este índice había surgido por la necesidad de controlar que se están utilizando los árboles correctos para reconocer el pie de un árbol auxiliar. En consecuencia, dicho índice sólo es de real utilidad en los pasos **FootPred** y **FootComp**. El resto de los pasos deductivos únicamente propagan el valor de dicho índice. Por consiguiente, en caso de que sea necesario estos últimos pasos deductivos pueden ser refinados, dividiéndolos en varios pasos con el fin de generar ítems intermedios carentes de dicho índice. Esta técnica, si

se aplica convenientemente, puede llegar a reducir la complejidad de los algoritmos de análisis. Evidentemente, hay que verificar que los ítems intermedios portan la información necesaria para que el resultado de la composición de los pasos deductivos obtenidos mediante el refinamiento de uno dado sea equivalente al resultado obtenido por aplicar directamente el paso deductivo sin refinar.

En el caso completo del esquema de análisis **Earley**, para reducir la complejidad a $\mathcal{O}(n^6)$ es suficiente con hacer uso de la *propiedad de independencia del contexto de TAG* [214]. Básicamente, lo que dicha propiedad establece es que cada operación de adjunción es independiente de la previa o posterior aplicación de cualquier otra operación de adjunción. Una consecuencia de esta propiedad es que si en un nodo M^γ de un árbol γ está permitida la adjunción de un árbol auxiliar β y se cumplen las tres condiciones siguientes:

1. la parte izquierda de la frontera de $\gamma - M^\gamma$ se extiende desde la posición h hasta la posición j de la cadena de entrada, donde $\gamma - M^\gamma$ denota el resultado de escindir el subárbol enraizado por M^γ de γ ;
2. la frontera del árbol β se expande desde la posición j hasta la posición m de la cadena de entrada con una discontinuidad en el pie desde la posición k hasta la l ;
3. la frontera del subárbol enraizado por M^γ abarca precisamente desde la posición k hasta la posición l ;

como resultado de la adjunción de β en M^γ la frontera del subárbol enraizado por este último nodo se expande desde la posición j hasta la m sin discontinuidad y dicho subárbol puede ser insertado en todo árbol $\beta - M^\gamma$ cuya frontera izquierda finalice en la posición j independiente de la posición h en la que se sitúe el extremo izquierdo de dicha frontera. Precisamente, los algoritmos de tipo Earley para TAG que no cumplen la propiedad del prefijo válido hacen uso de esta propiedad para verificar la corrección de la adjunción realizada en los pasos AdjComp. Como se recordará de la sección precedente, los ítems de los esquemas de análisis de dichos algoritmos no incorporan el índice h del extremo izquierdo.

En consecuencia, los ítems que utilizaremos en el esquema de análisis sintáctico **Nederhof** correspondiente al algoritmo de tipo Earley para TAG presentado por Nederhof en [125] que preserva la propiedad del prefijo válido con una complejidad $\mathcal{O}(n^6)$, son de dos tipos: los definidos para el esquema **Earley** y los pseudo-ítem que definimos a continuación

$$\left\{ \begin{array}{l} [[N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q]] \mid \begin{array}{l} \delta \xrightarrow{*} a_i \dots a_p \mathbf{F}^\gamma a_{q+1} \dots a_j \xrightarrow{*} a_i \dots a_j \quad \text{sii } (p, q) \neq (-, -) \\ \delta \xrightarrow{*} a_i \dots a_j \quad \text{sii } (p, q) = (-, -) \end{array} \end{array} \right\}$$

Ahora podemos definir el esquema de análisis **Nederhof** cuyo sistema de análisis mostramos a continuación.

Esquema de análisis sintáctico 3.7 El sistema de análisis $\mathbb{P}_{\text{Nederhof}}$ que se corresponde con la el algoritmo de análisis sintáctico de tipo Earley para TAG que cumple la propiedad del prefijo válido y posee una complejidad $\mathcal{O}(n^6)$, dada una gramática de adjunción de árboles \mathcal{T} y una cadena de entrada $a_1 \dots a_n$ se define como sigue:

$$\mathcal{I}_{\text{Nederhof}}^{(1)} = \mathcal{I}_{\text{Earley}} = \left\{ [h, N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q] \mid \begin{array}{l} N^\gamma \rightarrow \delta \nu \in \mathcal{P}(\gamma), \gamma \in \mathbf{I} \cup \mathbf{A}, \\ 0 \leq h \leq i \leq j, (p, q) \leq (i, j) \end{array} \right\}$$

$$\mathcal{I}_{\text{Nederhof}}^{(2)} = \left\{ [[N^\gamma \rightarrow \delta \bullet \nu, i, j \mid p, q]] \mid \begin{array}{l} N^\gamma \rightarrow \delta \nu \in \mathcal{P}(\gamma), \gamma \in \mathbf{I} \cup \mathbf{A}, \\ 0 \leq h \leq i \leq j, (p, q) \leq (i, j) \end{array} \right\}$$

$$\mathcal{I}_{\text{Nederhof}} = \mathcal{I}_{\text{Nederhof}}^{(1)} \cup \mathcal{I}_{\text{Nederhof}}^{(2)}$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{Init}} = \mathcal{D}_{\text{Earley}}^{\text{Init}} = \frac{}{\vdash [0, \top \rightarrow \bullet \mathbf{R}^\alpha, 0, 0 \mid -, -]} \quad \alpha \in \mathbf{I}$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{Scan}} = \mathcal{D}_{\text{Earley}}^{\text{Scan}} = \frac{[h, N^\gamma \rightarrow \delta \bullet a\nu, i, j \mid p, q], [a, j, j+1]}{[h, N^\gamma \rightarrow \delta a \bullet \nu, i, j+1 \mid p, q]}$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{Pred}} = \mathcal{D}_{\text{Earley}}^{\text{Pred}} = \frac{[h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[h, M^\gamma \rightarrow \bullet \nu, j, j \mid -, -]} \quad \mathbf{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{Comp}} = \mathcal{D}_{\text{Earley}}^{\text{Comp}} = \frac{[h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, k \mid p, q], [h, M^\gamma \rightarrow \nu \bullet, k, j \mid p', q']}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, j \mid p \cup p', q \cup q']} \quad \mathbf{nil} \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{AdjPred}} = \mathcal{D}_{\text{Earley}}^{\text{AdjPred}} = \frac{[h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[j, \top \rightarrow \bullet \mathbf{R}^\beta, j, j \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{FootPred}} = \mathcal{D}_{\text{Earley}}^{\text{FootPred}} = \frac{[j, \mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p, q]}{[h, M^\gamma \rightarrow \bullet \delta, k, k \mid -, -]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{FootComp}} = \mathcal{D}_{\text{Earley}}^{\text{FootComp}} = \frac{\begin{array}{l} [h, M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], \\ [j, \mathbf{F}^\beta \rightarrow \bullet \perp, k, k \mid -, -], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[j, \mathbf{F}^\beta \rightarrow \perp \bullet, k, l \mid k, l]} \quad \begin{array}{l} \beta \in \text{adj}(M^\gamma), \\ p \cup p' \neq y \cup q' \text{ está definido} \end{array}$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0} = \frac{\begin{array}{l} [j, \top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], \\ [h, M^\gamma \rightarrow \nu \bullet, k, l \mid p, q], \end{array}}{[[M^\gamma \rightarrow \nu \bullet, j, m \mid p, q]]} \quad \beta \in \text{adj}(M^\gamma)$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1} = \frac{\begin{array}{l} [[M^\gamma \rightarrow \nu \bullet, j, m \mid p, q]], \\ [h, \mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid -, -] \end{array}}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p, q]}$$

$$\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2} = \frac{\begin{array}{l} [[M^\gamma \rightarrow \nu \bullet, j, m \mid -, -]], \\ [h, N^\gamma \rightarrow \delta \bullet M^\gamma \nu, i, j \mid p', q'] \end{array}}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet \nu, i, m \mid p', q']}$$

$$\mathcal{D}_{\text{Nederhof}} = \mathcal{D}_{\text{Nederhof}}^{\text{Init}} \cup \mathcal{D}_{\text{Nederhof}}^{\text{Scan}} \cup \mathcal{D}_{\text{Nederhof}}^{\text{Pred}} \cup \mathcal{D}_{\text{Nederhof}}^{\text{Comp}} \cup \mathcal{D}_{\text{Nederhof}}^{\text{AdjPred}} \cup \mathcal{D}_{\text{Nederhof}}^{\text{FootPred}}$$

$$\cup \mathcal{D}_{\text{Nederhof}}^{\text{FootComp}} \cup \mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0} \cup \mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1} \cup \mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$$

$$\mathcal{F}_{\text{Nederhof}} = \mathcal{F}_{\text{Earley}} = \{ [0, \top \rightarrow \mathbf{R}^\alpha \bullet, 0, n \mid -, -] \mid \alpha \in \mathbf{I} \}$$

Obsérvese que se ha aplicado un refinamiento al paso deductivo $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ para obtener los pasos $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$ y $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1}$. Análogamente, el paso deductivo $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$ ha sido refinado en los pasos $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$ y $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$. Para garantizar la corrección se verifica que:

- La aplicación consecutiva de $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$ y $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ (resp. $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$ y $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$) es equivalente a la aplicación de $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ (resp. $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$). Es fácil comprobar que ambos utilizan la misma información: todos los antecedentes de $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1}$ (resp. $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$) son utilizados por pasos del esquema **Nederhof** y toda la información presente en los antecedentes de los pasos del esquema **Nederhof** es utilizada por el paso $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1}$ (resp. $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$), puesto que el ítem intermedio no crea nueva información, sino que simplemente es un “resumen” de información contenida en los otros antecedentes. Es también fácil comprobar que en ambos casos se genera la misma información, puesto que los ítems generados (excluyendo pseudo-ítems) son idénticos en ambos casos.
- El paso deductivo $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1}$ (resp. $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$) solo puede ser aplicado si previamente se ha aplicado el paso $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$. Se puede verificar fácilmente puesto que el paso $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$ genera un ítem intermedio y los únicos pasos que toman un ítem intermedio como antecedente son $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^2}$.

En la figura 3.5 se muestra una representación gráfica de la aplicación de los pasos deductivos $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^0}$ y $\mathcal{D}_{\text{Nederhof}}^{\text{AdjComp}^1}$.

La complejidad del algoritmo descrito por el esquema de análisis **Nederhof** es $\mathcal{O}(n^6)$, puesto que en la combinación de los ítems de cualquier paso intervienen activamente a los sumo 6 índices con respecto a la cadena de entrada.

En [64] se describe una versión de este algoritmo en la que se utiliza una representación plana de los árboles elementales en lugar de la representación multicapa que se ha utilizado aquí.

Proposición 3.6 $\text{Earley} \xrightarrow{\text{sr}} \text{Nederhof}$.

Demostración:

Para demostrar que el esquema de análisis **Nederhof** puede ser obtenido mediante un refinamiento de los pasos deductivos del esquema **Earley** debemos probar que para todo sistema de análisis $\mathbb{P}_{\text{Earley}}$ y $\mathbb{P}_{\text{Nederhof}}$ se cumple $\mathbb{P}_{\text{Earley}} \xrightarrow{\text{sr}} \mathbb{P}_{\text{Nederhof}}$. Ello conlleva demostrar que $\mathcal{I}_{\text{Earley}} \subseteq \mathcal{I}_{\text{Nederhof}}$ y que $\vdash_{\text{Earley}} \subseteq \vdash_{\text{Nederhof}}$. Lo primero se obtiene directamente puesto que $\mathcal{I}_{\text{Earley}} \subseteq \mathcal{I}_{\text{Nederhof}}$ por definición de los sistemas de análisis. Lo segundo se obtiene demostrando que $\mathcal{D}_{\text{Earley}} \subseteq \mathcal{D}_{\text{Nederhof}}$. Los únicos pasos deductivos de $\mathbb{P}_{\text{Earley}}$ que no se han incorporado directamente en $\mathbb{P}_{\text{Nederhof}}$ son $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ y $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$ y para ellos se cumple que:

- Un paso $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^1}$ es equivalente a la aplicación de un paso $\mathcal{D}_{\text{Nederhof}}^{\text{Comp}^0}$ seguido de la aplicación de un paso $\mathcal{D}_{\text{Nederhof}}^{\text{Comp}^1}$:

$$\frac{[j, \top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [h, M^\gamma \rightarrow v \bullet, k, l \mid p, q],}{[[M^\gamma \rightarrow v \bullet, j, m \mid p, q]]}$$

$$\frac{[[M^\gamma \rightarrow v \bullet, j, m \mid p, q]], [h, \mathbf{F}^\gamma \rightarrow \perp \bullet, p, q \mid p, q], [h, N^\gamma \rightarrow \delta \bullet M^\gamma v, i, j \mid -, -]}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet v, i, m \mid p, q]}$$

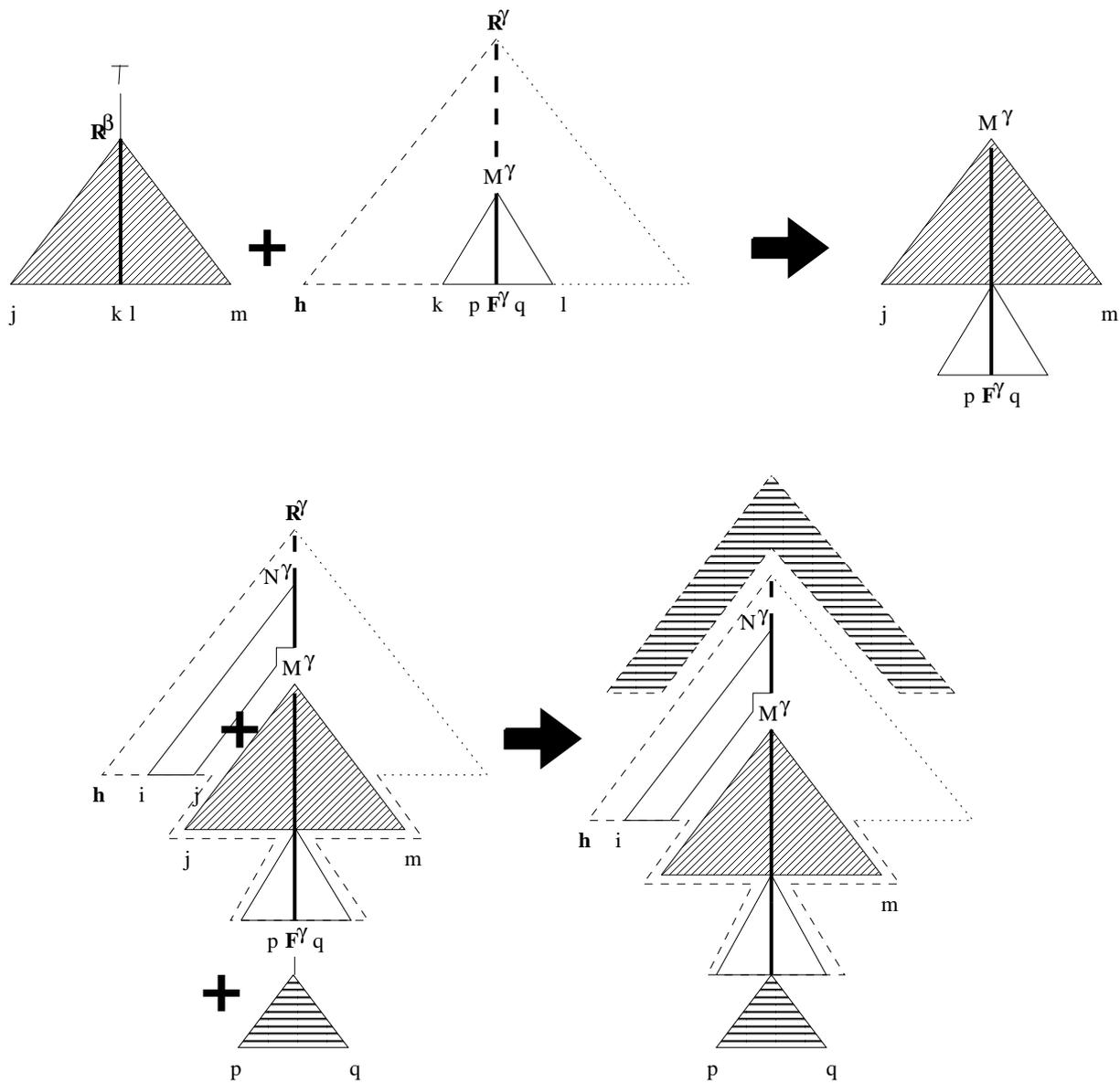


Figura 3.6: Descripción gráfica de la aplicación consecutiva de los pasos $\mathcal{D}_{Nederhof}^{AdjComp^0}$ y $\mathcal{D}_{Nederhof}^{AdjComp^1}$

- Un paso $\mathcal{D}_{\text{Earley}}^{\text{AdjComp}^2}$ es equivalente a la aplicación de un paso $\mathcal{D}_{\text{Nederhof}}^{\text{Comp}^0}$ seguido de la aplicación de un paso $\mathcal{D}_{\text{Nederhof}}^{\text{Comp}^2}$:

$$\frac{[j, \top \rightarrow \mathbf{R}^\beta \bullet, j, m \mid k, l], [h, M^\gamma \rightarrow v \bullet, k, l \mid p, q],}{[[M^\gamma \rightarrow v \bullet, j, m \mid p, q]]}$$

$$\frac{[[M^\gamma \rightarrow v \bullet, j, m \mid p, q]], [h, N^\gamma \rightarrow \delta \bullet M^\gamma v, i, j \mid p', q']}{[h, N^\gamma \rightarrow \delta M^\gamma \bullet v, i, m \mid p', q']}$$

□

3.7 Análisis sintáctico de TAG lexicalizadas

Las gramáticas lexicalizadas poseen una propiedad muy interesante desde el punto de vista del análisis sintáctico: son finitamente ambiguas. Puesto que cada componente de la gramática (en el caso de TAG, cada árbol elemental) está asociado con un componente léxico, solamente un conjunto finito de tales estructuras pueden ser utilizadas para el análisis de una cadena de entrada dada y además solamente existe un número finito de combinaciones de dichas estructuras. En resumen, las gramáticas lexicalizadas impiden la aparición de análisis cíclicos.

El análisis de gramáticas lexicalizadas puede realizarse en dos fases, una primera en la cual se seleccionan todas las estructuras relevantes para la cadena de entrada que se pretende analizar y una segunda en la cual se aplica un algoritmo de análisis sintáctico que combine dichas estructuras. Este tipo de procesamiento se corresponde con un análisis fuera de línea. Las gramáticas lexicalizadas también pueden ser analizadas en línea, de tal modo que según se va avanzado en la lectura de la cadena de entrada se proporcionen las estructuras elementales correspondientes. Algunos autores [174, 168] sugieren que en análisis fuera de línea está mejor adaptado a este tipo de gramáticas puesto que las estructuras seleccionadas en la primera fase posibilitan al analizador sintáctico la utilización de información ascendente no local⁷, restringiendo de este modo las posibilidades de combinación de las estructuras e incluso el número de estructuras a considerar. En efecto, al actuar de este modo, el analizador sintáctico solamente considerará aquellas estructuras relevantes para la cadena a analizar, por lo que se podría decir que trabaja sobre una subgramática relevante para la cadena de entrada.

Los beneficios que se obtiene de la lexicalización dependen del algoritmo de análisis que se vaya a aplicar. Los algoritmos puramente ascendentes, del tipo CYK, únicamente se benefician de la reducción del número de estructuras a considerar durante el proceso de análisis. Un algoritmo puramente descendente, basado en una exploración en profundidad con retroceso, conseguiría mayores beneficios, puesto que al ser las gramáticas lexicalizadas finitamente ambiguas el espacio de búsqueda es finito y por lo tanto el análisis terminará en todos los casos⁸. Los algoritmos mixtos que utilizan información ascendente y descendente, como por ejemplo los algoritmos de tipo Earley, se ven también beneficiados por la lexicalización. Una primera ventaja surge del hecho de que ningún árbol elemental tiene la cadena vacía por frontera, lo cual significa que una adjunción no puede ser predicha y completada sin avanzar en el reconocimiento de la cadena de entrada. Por tanto, al terminación está asegurada. Adicionalmente, la

⁷ Esta información puede incluso no estar acotada con respecto a la distancia [168], de tal modo que no se puede imitar su efecto mediante la utilización de un número limitado de símbolos de preanálisis. Esta característica se puede aplicar por ejemplo al reconocimiento de frases hechas con expresiones arbitrarias intercaladas.

⁸ Un analizador sintáctico puramente descendente puede no terminar para una gramática no lexicalizada puesto que puede intentar repetir indefinidamente el análisis del mismo conjunto de estructuras sin avanzar en el reconocimiento de la cadena de entrada.

utilización de una estrategia de dos fases permite que la selección de estructuras de acuerdo a los componentes léxicos presentes en la cadena de entrada ayude al analizador sintáctico en la tarea de filtrar las predicciones y/o compleciones para la adjunción y la sustitución. Resultados experimentales realizados Schabes y Joshi [174, 168] muestran que la estrategia de dos fases aumenta considerablemente la eficiencia de un algoritmo de análisis sintáctico de tipo Earley para TAG.

Todos los algoritmos mostrados en este capítulo pueden ser fácilmente adaptados a gramáticas de adjunción de árboles lexicalizadas. Para ello sólo es preciso incluir un paso deductivo para tratar la sustitución de un árbol en un nodo de sustitución. Puesto que dicha operación es independiente del contexto, no afecta a la complejidad espacial ni temporal de los algoritmos.

3.8 El bosque de análisis

Los algoritmos mostrados hasta el momento, tal y como han sido descritos, son realmente reconocedores y no analizadores sintácticos, puesto que no construyen árboles de derivación. Sin embargo, cada uno de los pasos deductivos contiene la información necesaria para generar la parte correspondiente de un árbol de derivación y, puesto que todos los algoritmos recorren todas las posibles derivaciones, se pueden reconstruir todos los posibles árboles de derivación.

Puesto que estamos tratando con analizadores no deterministas se trata de construir una estructura, denominada *bosque de análisis* [30, 215] que permita representar todas las derivaciones de un forma compacta, compartiendo subderivaciones comunes, y que permita extraer cada una de las derivaciones en tiempo lineal con respecto al tamaño del bosque de análisis. El problema de la construcción del bosque de análisis para TAG ha sido estudiado con anterioridad por Vijay-Shanker y Weir en [215], que han propuesto dos posibles soluciones: la utilización de gramáticas independientes del contexto y la utilización de gramáticas lineales de índices. Cualquiera de las soluciones es aplicable a los algoritmos de análisis sintáctico mostrados en este capítulo.

3.8.1 Gramáticas independientes del contexto como bosque de análisis

Es posible representar el bosque compartido mediante una gramática independiente del contexto que capture la independencia al contexto de la operación de adjunción. Los no-terminales de la gramática serán de la forma $\langle tb, N^\gamma, i, j, p, q \rangle$ donde $tb \in \{\top, \perp\}$ se utiliza para indicar si el no-terminal representa al nodo N^γ antes (\perp) o después (\top) de una adjunción. Es interesante observar que los no-terminales son casi idénticos a los ítems utilizados en el esquema de análisis **CYK**. Mediante una pequeña modificación en los pasos deductivos⁹ sería posible hacer $adj = \text{true}$ siempre que $tb = \top$ y que $adj = \text{false}$ siempre que $tb = \perp$. Puesto que los ítems de los restantes esquemas son un refinamiento de los ítems de **CYK** la información necesaria para los no-terminales se puede obtener directamente a partir de los ítems.

Respecto a la forma de las producciones, a modo de ejemplo, mostramos la producción correspondiente a la adjunción del árbol auxiliar β en el nodo N^γ :

$$\langle \top, N^\gamma, i, j, r, s \rangle \rightarrow \langle \top, \mathbf{R}^\beta, i, j, p, q \rangle \langle \perp, N^\gamma, p, q, r, s \rangle$$

la cual se corresponde con el paso deductivo de adjunción del esquema de análisis **CYK**. Al igual que ocurría con los no-terminales, las producciones del bosque de análisis se pueden obtener directamente a partir de los pasos deductivos en los diferentes esquemas de análisis.

⁹Esencialmente la adición de un paso deductivo $\mathcal{D}_{\text{CYK}}^{\text{NoAdj}} = \frac{[N^\gamma, i, j]_{p, q}[\text{false}]}{[N^\gamma, i, j]_{p, q}[\text{true}]}$ es aplicable siempre que la realización de una adjunción sobre N^γ sea opcional.

El número de producciones es $\mathcal{O}(n^6)$ y la construcción de la gramática tienen una complejidad temporal $\mathcal{O}(n^6)$, por lo que la complejidad temporal de los algoritmos permanece inalterable, aunque la complejidad espacial aumenta de $\mathcal{O}(n^4)$ ó $\mathcal{O}(n^5)$ a $\mathcal{O}(n^6)$.

Un aspecto interesante a destacar es que aunque el bosque de análisis construido de esta forma codifica las derivaciones para una cadena de entrada dada, el lenguaje derivado por la gramática independiente del contexto no es importante. Lo que importa es que el lenguaje generado es no vacío si la cadena pertenece a la TAG original y en tal caso las derivaciones para la TAG original puede ser obtenidas en tiempo lineal a partir de las derivaciones de la gramática independiente del contexto que codifica el bosque compartido, siempre que esta haya sido podada para eliminar los símbolos inútiles.

3.8.2 Gramáticas lineales de índices como bosque de análisis

Se puede representar el bosque compartido mediante una gramática lineal de índices utilizando la transformación de TAG en LIG definida en [214]. A modo de ejemplo, las siguiente producciones representan la adjunción del árbol auxiliar β en el nodo N^γ :

$$\begin{aligned} \langle \top, i, j \rangle [\circ \circ N^\gamma] &\rightarrow \langle \top, i, j \rangle [\circ \circ N^\gamma \mathbf{R}^\beta] \\ \langle \perp, p, q \rangle [\circ \circ N^\gamma \mathbf{F}^\beta] &\rightarrow \langle \perp, p, q \rangle [\circ \circ N^\gamma] \end{aligned}$$

La primera producción representa el final de la adjunción mientras que la segunda representa el reconocimiento del nodo pie del árbol auxiliar.

La información contenida en los no-terminales y producciones de la gramática lineal de índices puede obtenerse directamente a partir de los ítems y pasos deductivos de los esquemas de análisis sintáctico.

El tamaño de la gramática es $\mathcal{O}(n^3)$ y el tiempo necesario para construirla es $\mathcal{O}(n^6)$, por lo que las complejidades temporal y espacial de los algoritmos no se ven afectadas.

El inconveniente de esta representación estriba en que no se pueden extraer directamente los árboles de derivación individuales, sino que es preciso construir una estructura auxiliar que toma la forma de un autómata finito que reconoce las pilas de índices asociadas a cada terminal. Una vez construido dicho autómata finito, cada uno de los árboles de derivación puede ser extraído con una complejidad temporal que tiene por cota inferior $\mathcal{O}(n^4)$ y que en el caso de las gramáticas de adjunción de árboles lexicalizadas tiene como cota superior $\mathcal{O}(n^5)$. Puesto que el tamaño del bosque compartido es $\mathcal{O}(n^3)$, no se asegura que en todos los casos la recuperación de los árboles de derivación individuales se pueda realizar en tiempo lineal con respecto al tamaño del bosque compartido.

3.9 Otros algoritmos de análisis sintáctico para TAG

Presentamos a continuación un conjunto de algoritmos de análisis sintáctico de TAG que si bien no están en el camino principal de la evolución de los algoritmos de análisis mostrado anteriormente, presentan interés bien por constituir caminos laterales del camino principal de evolución, bien por haber constituido hitos importantes aunque posteriormente hayan quedado relegados, o bien por constituir ejemplos singulares de aplicación a TAG de ciertas áreas del análisis sintáctico, como la incrementalidad o el paralelismo.

3.9.1 El algoritmo de Lang

Lang describe en [106] un algoritmo tabular de análisis de TAG, con el objetivo principal de mostrar que técnicas de análisis muy generales pueden ser utilizadas para desarrollar un analizador sintáctico de tipo Earley para TAG. Concretamente, Lang utiliza las siguientes técnicas: