

# **Reporte Técnico Especializado:**

# **Fundamentos, Algoritmos y Alternativas**

# **Robustas del Análisis Sintáctico de**

# **Datos**

## **I. Fundamentos Teóricos del Análisis Sintáctico (Parsing)**

### **1.1 Definición y Propósito Fundamental del Parsing**

El análisis sintáctico, o *parsing*, constituye una fase esencial y fundamental en las ciencias de la computación, sirviendo como puente entre la secuencia lineal de símbolos de entrada y la representación jerárquica de su estructura subyacente. Formalmente, el parsing es el proceso algorítmico mediante el cual se toma una cadena de *tokens* (símbolos léxicos) y se evalúa si dicha cadena pertenece al lenguaje definido por una gramática formal específica. Si la cadena es gramaticalmente válida, el proceso genera una representación estructural de su derivación.<sup>1</sup>

En el contexto operativo y funcional, el parser se concibe como una función matemática que toma como entradas tanto la gramática de referencia como la cadena de símbolos a analizar.<sup>1</sup> La salida de este proceso puede bifurcarse en dos objetivos principales: el primero es la *función de reconocimiento*, que simplemente determina si la cadena es gramaticalmente aceptable; el segundo, y más relevante para la ingeniería, es la *función de análisis gramatical*, que produce la estructura que la gramática asigna a la sentencia.<sup>1</sup>

Este proceso es la segunda etapa crítica en el ciclo de procesamiento de lenguajes, siguiendo al análisis léxico (tokenización). La tarea principal del parsing es identificar la función específica de cada componente de la entrada para, de esta forma, descifrar su estructura global.<sup>2</sup> El resultado tangible y más común de esta etapa es la construcción de un Árbol Sintáctico Abstracto (AST) o un Árbol Sintáctico de Derivación (Parse Tree). Estas estructuras representan la cadena de entrada incluyendo todos sus detalles gramaticales, como paréntesis y operadores, permitiendo una interpretación profunda del contenido.<sup>3</sup>

Las aplicaciones del parsing son transversales, abarcando desde la transformación de datos en bruto en estructuras organizadas y comprensibles para el análisis y almacenamiento<sup>5</sup>, hasta la construcción de compiladores e intérpretes para lenguajes de programación. Incluso en aplicaciones modernas como el *Web Scraping*, el parser es crucial, ya que transforma la estructura semi-estructurada del HTML en datos organizados (extracción de URLs, textos).<sup>5</sup>

## 1.2 Gramáticas Formales: El Marco Teórico

El fundamento matemático para el análisis sintáctico clásico reside en las Gramáticas Libres de Contexto (CFG), formalizadas a menudo mediante la notación *Backus-Naur Form* (BNF). Bajo esta perspectiva, una gramática actúa como una teoría deductiva, donde el símbolo raíz se establece como el único axioma, y cada regla o producción gramatical opera como una regla inferencial.<sup>1</sup> Las cadenas de entrada son, entonces, teoremas potenciales que el parser debe verificar si pueden ser generados a partir del axioma inicial.<sup>1</sup>

### El Desafío de la Ambigüedad Sintáctica

La coherencia de un sistema de parsing se ve directamente amenazada por la ambigüedad inherente de la gramática. Una CFG se considera ambigua si existe al menos una cadena en su lenguaje que puede ser generada por más de una secuencia de derivaciones, lo que resulta en más de un árbol de análisis sintáctico.<sup>6</sup>

Esta ambigüedad tiene consecuencias críticas tanto en el diseño de lenguajes artificiales como en el procesamiento del lenguaje natural (PLN). En el contexto de los lenguajes de programación, si una expresión puede tener múltiples árboles de derivación, implica que puede tener múltiples significados (por ejemplo, en la precedencia de operadores,  $\$3+4*5\$$  podría significar  $\$(3+4)*5\$$  o  $\$3+(4*5)\$$ ).<sup>6</sup> Para la aplicación en compilación, resulta imperativo diseñar gramáticas no ambiguas siempre que sea posible para garantizar que cada expresión tenga un significado único.<sup>6</sup>

Este desafío central—la crisis de la ambigüedad sintáctica—actúa como un motor principal en la evolución de las técnicas de parsing. Los algoritmos clásicos se desarrollaron precisamente para imponer el determinismo y resolver localmente la ambigüedad, aunque esto limita la clase de lenguajes que pueden ser reconocidos. La necesidad de gramáticas artificiales *no ambiguas* contrasta fuertemente con la *ambigüedad ineludible* de los lenguajes naturales, donde la misma cadena puede tener múltiples interpretaciones válidas dependiendo del contexto.<sup>7</sup> Esta divergencia fuerza la búsqueda de soluciones en dos dominios distintos: formalismos determinísticos (como PEGs) para lenguajes artificiales, y métodos basados en el contexto y la semántica (como los Modelos de Lenguaje Grande o LLMs) para PLN.

Una limitación fundamental del enfoque clásico es su rigidez. Los parsers tradicionales están diseñados para verificar la *validez sintáctica* estricta. Si la entrada contiene lenguaje informal, errores gramaticales o desviaciones menores, el proceso es propenso al colapso.<sup>4</sup> Para el análisis de datos modernos, donde la entrada rara vez es perfectamente formal (datos web, logs, texto libre), la robustez frente a errores prima sobre la adhesión perfecta a la gramática. La dependencia estricta de una gramática perfecta impone un costo oculto de fragilidad que justifica la transición hacia enfoques de IA más flexibles que permiten el "parsing implícito".<sup>4</sup>

## II. Algoritmos Clásicos de Parsing y su Eficiencia

Los algoritmos de parsing que operan sobre CFGs se clasifican tradicionalmente en dos grandes categorías basadas en la dirección en que construyen el árbol sintáctico.

### 2.1 Clasificación: Top-Down (Descendente) vs. Bottom-Up (Ascendente)

El proceso de análisis sintáctico busca establecer la correspondencia entre la cadena de entrada y la gramática.

#### Análisis Descendente (*Top-Down*)

El parsing descendente inicia el proceso en el símbolo inicial (la raíz del árbol sintáctico) y trabaja progresivamente hacia las hojas (los *tokens* de la cadena de entrada).<sup>9</sup> Este enfoque intenta encontrar la derivación más a la izquierda de la cadena de entrada (*Leftmost Derivation*). La característica clave de un parser descendente es que *espera* una entrada específica basándose en el historial de análisis realizado hasta ese momento. Por ejemplo, al reconocer una palabra clave como *if*, el parser anticipa inmediatamente que la estructura gramatical requiere una expresión lógica a continuación.<sup>10</sup>

#### Análisis Ascendente (*Bottom-Up*)

En contraste, el parsing ascendente comienza desde la cadena de entrada (las hojas del árbol sintáctico) e intenta *reducir* esta cadena hasta el símbolo inicial de la gramática (la raíz).<sup>9</sup> Este proceso realiza la derivación más a la derecha en orden inverso (*Rightmost Reduction*). El parser ascendente lee los *tokens* hasta que puede deducir el patrón que siguen. Una vez que el patrón es identificado, se realiza una *reducción* para reemplazar esa secuencia de símbolos terminales y no terminales por el símbolo no terminal que los generó.<sup>10</sup>

### 2.2 Análisis Detallado de Parsers Descendentes (*Top-Down*)

Los parsers descendentes pertenecen a la familia **LL(k)**, una notación que describe su mecanismo operativo: L de izquierda a derecha (lectura), L de derivación más a la izquierda, y \$k\$ como el número de *tokens* de anticipación (*lookahead*) utilizados para tomar decisiones determinísticas.<sup>11</sup>

#### Recursive Descent Parsing

La técnica más directa para implementar un parser descendente es el *Recursive Descent*.<sup>10</sup> En este método, se escribe una función o procedimiento dedicado por cada símbolo no terminal

de la gramática. Esta técnica destaca por su notable simplicidad de implementación y comprensión.<sup>9</sup> Un parser *Recursive Descent* se puede construir manualmente (a mano) con relativa facilidad, lo que lo hace popular para lenguajes sencillos o prototipos.<sup>10</sup>

Sin embargo, para garantizar un análisis determinístico (especialmente LL(1)), la gramática debe cumplir con dos restricciones estrictas: la eliminación de la **recursión izquierda** (donde un no terminal se deriva a sí mismo inmediatamente a la izquierda, como  $A \to A\alpha$ ) y la aplicación de *left-factoring* (factorización izquierda) para remover prefijos comunes de reglas alternativas.<sup>10</sup> Si el parser no es determinístico, requiere *backtracking* (retroceso), lo que puede generar una ineficiencia significativa.

## Parsing Predictivo (LL) y Variantes Avanzadas

Una variante optimizada del análisis descendente es el *Parsing Predictivo* (LL). Este método es capaz de anticipar qué regla gramatical debe aplicarse a continuación basándose únicamente en el símbolo de anticipación (el *lookahead*), eliminando la necesidad de *backtracking* y conduciendo a optimizaciones significativas.<sup>9</sup>

La clase de gramáticas **LL(k)** utiliza  $k$  tokens de anticipación. Es importante señalar que el conjunto de lenguajes  $\text{LL}(k)$  está estrictamente contenido dentro del conjunto  $\text{LL}(k+1)$ .<sup>11</sup> Esto implica que no todos los lenguajes libres de contexto pueden ser reconocidos por un parser  $\text{LL}(k)$ .<sup>11</sup> Existen variantes más sofisticadas como los parsers **LL(\*)** y **LL(finite)**, los cuales utilizan estrategias de anticipación avanzadas, acercándose funcionalmente a las Gramáticas de Expresiones de Parsing (PEGs). No obstante, los parsers  $\text{LL}(\ast)$  no son, en general, lineales y pueden mostrar un rendimiento super-lineal en promedio, o incluso exponencial en el peor de los casos, a diferencia de sus contrapartes lineales.<sup>11</sup>

## 2.3 Análisis Detallado de Parsers Ascendentes (Bottom-Up)

Los parsers ascendentes comprenden la familia **LR(k)**, que opera leyendo de Izquierda a derecha (L) y realizando la reducción más a la Derecha (R), utilizando  $k$  símbolos de anticipación.<sup>10</sup>

### Operación y Tipos LR

Los parsers LR son intrínsecamente más potentes que los parsers LL.<sup>9</sup> Operan mediante un mecanismo de pila que ejecuta acciones de *shift* (desplazamiento de un token a la pila) y *reduce* (reemplazo de una secuencia de elementos de la pila por un no terminal).<sup>10</sup> La ejecución de estas acciones está guiada por complejas **tablas de parsing** que deben ser generadas automáticamente por un software especializado (un generador de parsers).<sup>10</sup>

La familia LR incluye subtipos que varían en su potencia y la complejidad de sus tablas: SLR (Simple LR), LALR (Look-Ahead LR), y CLR (Canonical LR).

- **LALR(1):** Este subtipo representa un punto de equilibrio óptimo y es el más utilizado en

la práctica de la construcción de compiladores (por ejemplo, herramientas como YACC o Bison). Los parsers LALR(1) tienen casi la misma potencia gramatical que los CLR(1) (el más potente de la familia LR) pero logran generar tablas de parsing significativamente más pequeñas y manejables, lo que mejora la eficiencia y el almacenamiento.<sup>12</sup>

La principal ventaja de la familia LR es su alta potencia gramatical, lo que les permite manejar gramáticas más complejas que las que pueden ser procesadas por los parsers LL, sin las restricciones sobre la recursión izquierda. La principal desventaja radica en que su complejidad es tal que su implementación manual resulta extremadamente difícil; por ello, la dependencia de generadores de parsers es virtualmente absoluta.<sup>9</sup>

## 2.4 Combinator Parsing y Enfoques Modulares

El *Combinator Parsing* representa una metodología de construcción de parsers que adopta un enfoque modular y funcional. En lugar de depender de tablas o funciones recursivas rígidas, esta técnica utiliza funciones de alto orden (combinadores) para construir parsers complejos a partir de componentes más simples.<sup>13</sup>

Esta metodología es particularmente popular en lenguajes de programación funcionales (como Haskell), donde el parser puede ser modelado como una mónada.<sup>13</sup> Esto permite una composición elegante de parsers mediante la aplicación de funciones de alto orden para representar construcciones gramaticales esenciales como la secuencia, la alternancia y la repetición.<sup>13</sup> Una característica notable es que el código del parser resultante se asemeja notablemente a la notación BNF de la gramática que está diseñada para reconocer, lo que facilita su comprensión y modificación. Además, el *Combinator Parsing* es lo suficientemente flexible como para reconocer gramáticas ambiguas, aunque a menudo esto implica el uso de *backtracking*.<sup>13</sup>

**Tabla I: Comparativa de Algoritmos de Parsing Clásicos**

El análisis formal de las estrategias de parsing obliga a sopesar la sencillez de implementación frente a la potencia gramatical que puede manejar cada técnica. La Tabla I resume estas compensaciones inherentes:

Característica	Top-Down (Ej. LL(k))	Bottom-Up (Ej. LR/LALR)	Combinator Parsing
<b>Derivación / Reducción</b>	Izquierda (Leftmost Derivation) <sup>9</sup>	Derecha (Rightmost Reduction Inversa) <sup>10</sup>	Variable (A menudo Leftmost Derivation con Backtracking) <sup>13</sup>
<b>Punto de Inicio</b>	Símbolo Inicial (Raíz) <sup>9</sup>	Tokens de Entrada (Hojas) <sup>9</sup>	Símbolo Inicial (Modular)
<b>Simplicidad de</b>	Alta (Recursive)	Baja (Requiere)	Media-Alta

<b>Implementación</b>	Descent) <sup>9</sup>	generador) <sup>10</sup>	(Declarativo y componible) <sup>13</sup>
<b>Potencia Gramatical</b>	Limitada (Requiere manipulación de gramática) <sup>9</sup>	Alta (Maneja la mayoría de CFGs) <sup>9</sup>	Puede manejar gramáticas ambiguas <sup>13</sup>

### III. Robustez, Coherencia y Alternativas

#### Determinísticas

La necesidad de robustez y coherencia en el análisis sintáctico ha impulsado el desarrollo de formalismos que buscan mitigar las debilidades de las CFGs, especialmente en lo referente al manejo de errores y la ambigüedad inherente.

##### 3.1 Limitaciones Críticas de los Parsers Basados en CFG

La familia de parsers determinísticos LL( $k$ ) y LR( $k$ ) opera bajo la premisa de una entrada perfecta. Su principal limitación reside en su **fragilidad ante los errores sintácticos**. Cuando un parser encuentra un *token* que no esperaba según las reglas de su gramática, el proceso tiende a detenerse bruscamente o a generar una cascada de errores que obscurecen la causa original.<sup>8</sup>

##### Manejo de Errores y Recuperación

Para evitar un colapso total en entornos de compilación, se han desarrollado técnicas de **recuperación de errores**. Una técnica común es la *sincronización*, que implica que el parser descarte *tokens* hasta que encuentre un marcador "seguro" para reanudar el análisis. Por ejemplo, en lenguajes imperativos, se puede saltar toda la secuencia de error hasta encontrar el siguiente punto y coma, asumiendo que el próximo *token* después de ese separador es un punto seguro para continuar el análisis a nivel de instrucción.<sup>8</sup> Sin embargo, esta estrategia es dependiente de la existencia de separadores claros, y en lenguajes donde toda la entrada es tratada como una expresión (como algunos lenguajes funcionales), la recuperación se vuelve mucho más difícil, pudiendo requerir el uso de conjuntos *follow* que, a su vez, pueden introducir errores en cascada.<sup>8</sup>

Para mejorar esta fragilidad sin abandonar la base formal, la investigación ha explorado **esquemas de análisis sintáctico con corrección de errores**. Estos formalismos permiten definir algoritmos de corrección de errores de manera abstracta y declarativa, proporcionando una base formal para demostrar su corrección y permitiendo implementaciones que realizan corrección regional de errores.<sup>14</sup>

La inflexibilidad es otra limitación significativa. Los parsers tradicionales requieren gramáticas que estén perfectamente definidas para operar.<sup>4</sup> Cuando se enfrentan a lenguaje informal o errores leves, su rendimiento decae rápidamente. Esto subraya que la robustez en la

ingeniería moderna no solo debe significar la capacidad de recuperación ante errores de sintaxis, sino también la capacidad de inferir la estructura incluso cuando la entrada es imperfecta o ambigua, lo cual es la clave para la flexibilidad requerida en PLN.<sup>7</sup>

## 3.2 Gramáticas de Expresiones de Parsing (PEGs): Coherencia Determinística

Las *Parsing Expression Grammars* (PEGs) representan una alternativa formalmente distinta a las CFGs, diseñadas específicamente para maximizar la **coherencia** y el **determinismo** en el análisis de lenguajes artificiales.

### Principios Operacionales: La Elección Ordenada

El principio operativo que distingue a las PEGs es la **elección ordenada** (/). A diferencia de la elección no ordenada (l) de las CFGs y las expresiones regulares, en una PEG, si la primera alternativa de una regla tiene éxito, el resto de las alternativas son irrevocablemente ignoradas.<sup>15</sup> Esta elección ordenada no es conmutativa y es análoga a los operadores de corte suave que se encuentran en algunos lenguajes de programación lógica.<sup>15</sup>

Esta propiedad tiene una consecuencia profunda: **las PEGs son intrínsecamente no ambiguas**. Una cadena de entrada en una PEG puede tener exactamente *un* árbol de análisis válido o ninguno.<sup>16</sup> Este determinismo estructural garantiza una coherencia superior en el análisis, eliminando la necesidad de que el parser resuelva conflictos de *shift-reduce* o *reduce-reduce*, o que el diseñador de la gramática imponga reglas explícitas de precedencia mediante nombres no terminales como term y factor.<sup>10</sup> La responsabilidad de la ambigüedad se traslada de la estructura gramatical a la secuencia de las reglas de producción. El diseñador debe ordenar las reglas cuidadosamente, haciendo que la coherencia sea una propiedad inherente y verificable del formalismo.

### Eficiencia y Aplicabilidad

Una de las mayores ventajas algorítmicas de las PEGs es que cualquier PEG puede ser analizada en **tiempo lineal ( $\$O(n)$ )** mediante un *Packrat Parser*.<sup>16</sup> Un *Packrat Parser* es esencialmente un analizador recursivo descendente que utiliza memoización para almacenar y reutilizar los resultados de los análisis de sub-expresiones.

Este rendimiento lineal es significativamente superior al de los algoritmos de parsing para CFGs generales, cuya complejidad asintótica de peor caso se sitúa entre el tiempo cuadrático y cúbico.<sup>16</sup> Para la ingeniería de lenguajes de dominio específico (DSLs) y nuevos compiladores que valoran el rendimiento predecible y la coherencia, las PEGs ofrecen un avance crucial al garantizar el tiempo lineal.

Las PEGs son particularmente adecuadas para lenguajes informáticos y artificiales donde la desambiguación es un fenómeno local (como el lenguaje Lojban).<sup>16</sup> Sin embargo, son menos efectivas para el análisis de lenguajes naturales, donde la desambiguación puede depender

de un contexto global extenso.<sup>16</sup>

## IV. Alternativas de Vanguardia para la Extracción Robusta: IA y Análisis Semántico

Cuando el objetivo del análisis se desplaza de la verificación sintáctica estricta de un lenguaje formal a la extracción de significado de datos semi-estructurados o no estructurados, los modelos de parsing clásicos demuestran ser insuficientes. Las alternativas más robustas se centran en la inferencia semántica y en la garantía de un formato de salida coherente, independientemente de la perfección de la entrada.

### 4.1 Análisis Basado en Grafos para la Estructura Semántica

Para lograr una comprensión más profunda del contenido, la lingüística computacional ha explorado estructuras de datos más ricas que el simple AST, buscando modelar las relaciones semánticas.

#### Grafos Conceptuales (GCs) y Gramáticas Avanzadas

Los **Grafos Conceptuales** (GCs) se utilizan como un instrumento eficiente para la representación formal del significado de un texto.<sup>17</sup> El método tradicional para crear GCs involucraba un proceso de análisis sintáctico inicial seguido de una costosa transformación. Las investigaciones actuales buscan simplificar esto, proponiendo soluciones para la transformación directa de texto a GCs, ya sea mediante la construcción automática de gramáticas especializadas (a partir de recursos léxicos existentes) o a través de métodos estadísticos no supervisados.<sup>17</sup> Estos métodos estadísticos, basados en la coocurrencia de palabras en grandes corpus, permiten obtener estructuras simplificadas de GCs sin la necesidad de un entrenamiento manual intensivo.<sup>17</sup>

Un avance en este campo es el uso de las **Gramáticas de Reemplazo de Hiperaristas (HRGs)**. Los grafos son considerados una herramienta apropiada para representar la estructura semántica, y las HRGs pueden definir modelos de lenguaje para la transformación de grafos. Esto es vital para resolver problemas complejos de PLN como la desambiguación del sentido de la palabra, la comprensión de texto y el resumen automático.<sup>18</sup>

### 4.2 La Revolución de los Modelos de Lenguaje Grande (LLMs)

La adopción de los Modelos de Lenguaje Grande (LLMs), como BERT y GPT, ha introducido un paradigma fundamentalmente nuevo conocido como **parsing implícito**.<sup>4</sup> Este enfoque es radicalmente diferente del parsing explícito, ya que la estructura se infiere del vasto conocimiento contextual adquirido durante el entrenamiento del modelo, en lugar de ser estrictamente impuesta por reglas gramaticales codificadas.

Los modelos modernos de PLN han demostrado que el parsing implícito puede ser tan efectivo o incluso superior al explícito, especialmente porque son inherentemente **flexibles y robustos** frente al lenguaje informal o a los errores gramaticales.<sup>4</sup> Esta flexibilidad es crucial dado que, según las estimaciones, entre el 80% y el 90% de los datos empresariales son no estructurados, y la cantidad de estos datos crece considerablemente más rápido que los datos estructurados.<sup>19</sup>

El uso de LLMs es la solución de vanguardia para la extracción de información clave de documentos semi-estructurados (como archivos HTML, XML o JSON)<sup>19</sup> y documentos no estructurados. Los flujos de trabajo de análisis de documentos asistidos por LLM permiten extraer texto de diversas fuentes ( párrafos, tablas, encabezados)<sup>20</sup> y aplicar modelos avanzados (como BERT o RoBERTa) para la clasificación y preprocesamiento de documentos.<sup>21</sup>

Además, la integración de la IA generativa con herramientas de procesamiento de documentos (como Document AI) permite analizar y extraer datos con mayor precisión y rapidez tanto de documentos estructurados como no estructurados.<sup>22</sup> Estos sistemas reducen el tiempo de amortización y los costos al minimizar la necesidad de etiquetado y anotación manual intensiva.<sup>22</sup>

### **4.3 Structured Output (SO) y Function Calling: El Pináculo de la Robustez de Formato**

Si bien los LLMs son robustos en la comprensión del lenguaje natural, un problema persistente es la garantía de que el *output* resultante (por ejemplo, en formato JSON para la integración con sistemas) sea sintácticamente válido y se adhiera a un esquema predefinido. Depender de la simple generación de texto que debe ser parseado por una cadena externa (string parsing) es frágil, ya que los LLMs pueden generar texto similar a JSON pero inválido o malformado.<sup>23</sup>

#### **Garantía de Coherencia Mediante Esquemas**

Para obtener una salida coherente y fiable, las técnicas de *Structured Output (SO)* y *Function Calling* son las alternativas más robustas al parsing tradicional para datos semi-estructurados. Estas técnicas son fundamentalmente más fiables que el parsing de cadenas en el *output* bruto.<sup>24</sup>

- **Mecanismo de Function Calling:** El desarrollador informa al LLM sobre un conjunto de funciones específicas que puede ejecutar, junto con sus parámetros definidos por un esquema (usualmente JSON Schema). El LLM responde determinando qué función llamar y generando los argumentos necesarios en un formato JSON que se ajusta rigurosamente al esquema.<sup>24</sup>
- **Ventajas de Structured Output:** El enfoque SO garantiza que el JSON generado sea siempre válido y coincida con el esquema deseado.<sup>23</sup> Este método proporciona una

predictibilidad y fiabilidad de la estructura de datos que supera la fragilidad del parsing sintáctico clásico.

El cambio paradigmático en la definición de robustez es profundo: el análisis ya no se centra en la adhesión a las reglas de la gramática de entrada, sino en la **adhesión garantizada al formato de salida deseado**. El LLM realiza el análisis implícito semántico, y el mecanismo de SO/Function Calling impone la coherencia sintáctica en el *output*. Esto traslada la coherencia del dominio sintáctico estricto (la entrada) al dominio pragmático-funcional (la salida).

**Tabla II: Evaluación de Formalismos y Robustez Moderna**

La siguiente tabla resume el rendimiento y la robustez de los principales formalismos y enfoques modernos, destacando cómo el paradigma de la coherencia ha evolucionado:

Formalismo	Propiedad Clave de Coherencia	Manejo de Ambigüedad	Robustez ante Errores Sintácticos de Entrada	Complejidad de Tiempo (Peor Caso)
LR/LALR (CFG)	Determinismo por Tabla de Parsing	Resuelve mediante Lookahead/Manipulación Gramatical <sup>9</sup>	Baja (Requiere técnicas de recuperación manual) <sup>8</sup>	Lineal ( $\$O(n)\$$ )
PEG (Packrat)	Elección Ordenada (Intrínsecamente No Ambigua) <sup>15</sup>	Eliminada por definición <sup>16</sup>	Media-Alta (Buena para errores locales en DSLs)	Lineal ( $\$O(n)\$$ ) <sup>16</sup>
LLM (Structured Output)	Adherencia Garantizada a Esquema de Salida (JSON Schema) <sup>23</sup>	Resuelve contextualmente (Semántica) <sup>7</sup>	Muy Alta (Ignora errores sintácticos de entrada si mantiene la intención) [4, 24]	Variable / Alta Latencia

## V. Conclusiones y Recomendaciones Estratégicas

### 5.1 Síntesis Comparativa y Guía de Selección de Estrategia

La elección de la técnica de análisis sintáctico debe ser una decisión estratégica basada en la naturaleza del lenguaje de entrada y los requisitos de robustez, rendimiento y predictibilidad.

#### Análisis para Lenguajes Formales (Compiladores y DSLs)

Para la construcción de compiladores, intérpretes y analizadores de lenguajes de dominio específico (DSLs) donde la entrada es típicamente bien formada y el rendimiento es crítico, se debe elegir entre la potencia y la simplicidad.

1. **Potencia y Universalidad (LALR(1)):** La familia LR, particularmente LALR(1), ofrece la mayor potencia gramatical, manejando la clase más amplia de CFGs con eficiencia lineal ( $O(n)$ ). Es la opción estándar para los lenguajes de programación principales.<sup>9</sup>
2. **Sencillez y Coherencia Intrínseca (PEGs):** Las Gramáticas de Expresiones de Parsing (PEGs) son la opción más elegante cuando la coherencia (no ambigüedad) es un requisito fundamental y la eficiencia lineal garantizada es crucial. Al transferir la resolución de ambigüedad a la ordenación de reglas, las PEGs simplifican el motor de análisis y resultan más sencillas de implementar manualmente que los parsers LR.<sup>10</sup> Las PEGs son inherentemente más coherentes que los LL o LR, ya que fuerzan el determinismo por construcción.

## Análisis para Datos Naturales y Empresariales (Robustez Semántica)

Cuando la entrada es variable, propensa a errores, semi-estructurada (JSON, HTML con errores) o completamente no estructurada (documentos, logs, correos electrónicos), la robustez semántica supera la prioridad de la perfección sintáctica.

1. **Robustez y Comprensión Semántica (LLMs y Structured Output):** Los LLMs, combinados con técnicas de *Structured Output* (SO) o *Function Calling*, representan la alternativa más coherente y robusta para la extracción de valor de datos inciertos.<sup>23</sup> En este paradigma, la validez del análisis no depende de la gramática de la fuente, sino de la adhesión estricta a un esquema de salida predefinido. La robustez se define por la capacidad de inferir la *intención* y la *estructura deseada*, incluso ante entradas defectuosas o informales.

Este cambio de enfoque ilustra cómo la IA ha redefinido el concepto de "coherencia" en el parsing, moviéndolo del dominio sintáctico (la estructura de las reglas) al dominio semántico-pragmático (la utilidad del resultado).<sup>25</sup>

## 5.2 Tendencias Futuras y Arquitecturas Híbridas

El desarrollo futuro del análisis sintáctico apunta a la convergencia y complementariedad de los enfoques explícitos y los neuronales, buscando lo mejor de ambos mundos: la transparencia y controlabilidad de los parsers formales y la flexibilidad y robustez de los modelos de IA.

Se anticipa que habrá un resurgimiento del parsing explícito, pero integrado con modelos neuronales.<sup>4</sup> Esta combinación es crucial para avanzar hacia la **IA responsable y explicable**, ya que permite que los sistemas mantengan interpretaciones más controlables y transparentes. Por ejemplo, los métodos clásicos podrían pre-analizar la sintaxis (como en un *Packrat Parser*), mientras que un modelo neural resuelve las ambigüedades semánticas que

escapan al análisis puramente sintáctico.<sup>4</sup>

La arquitectura de análisis de documentos empresariales continuará evolucionando hacia sistemas híbridos. La combinación de la visión artificial y la IA generativa (como en Document AI) permite abordar la complejidad de los formatos físicos y digitales, utilizando la IA generativa para procesar documentos de manera inmediata y precisa, superando las limitaciones de los sistemas tradicionales basados en plantillas o reglas rígidas.<sup>22</sup> En esencia, si el problema requiere un rendimiento determinístico y reglas fijas, las PEGs son la solución elegante; si el problema requiere la extracción de valor de datos empresariales variables y no confiables, el análisis semántico basado en LLM con salida estructurada es la estrategia dominante.

## Fuentes citadas

1. el Algoritmo SCP, acceso: noviembre 4, 2025,  
[https://rua.ua.es/bitstream/10045/1977/1/PLN\\_24\\_01.pdf](https://rua.ua.es/bitstream/10045/1977/1/PLN_24_01.pdf)
2. acceso: noviembre 4, 2025,  
<https://www.gostudent.org/es-es/blog/analizar-sintacticamente#:~:text=El%20an%C3%A1lisis%20sint%C3%A1ctico%20consiste%20en,la%20frase%20a%20la%20perfecci%C3%B3n.>
3. Analizador sintáctico :: COMPILADORES - Webnode, acceso: noviembre 4, 2025,  
<https://compiladores57.webnode.mx/analizador-sintactico/>
4. Parsing - MSMK, acceso: noviembre 4, 2025, <https://msmk.university/parsing/>
5. Parsing (Análisis de Datos) - Manu Duque, acceso: noviembre 4, 2025,  
<https://www.manuduque.com/diccionario-seo/parsing-analisis-de-datos/>
6. GRAMATICAS LIBRES DEL CONTEXTO, acceso: noviembre 4, 2025,  
<https://users.exa.unicen.edu.ar/catedras/ccomp1/Apunte5.pdf>
7. Técnicas de análisis sintáctico robusto para la etiquetación del lenguaje natural - CORE, acceso: noviembre 4, 2025, <https://core.ac.uk/download/pdf/61910629.pdf>
8. ¿Cómo hago un parser que no se caiga ni se corrompa cuando encuentra un error?, acceso: noviembre 4, 2025,  
[https://www.reddit.com/r/Compilers/comments/1jrmwrs/how\\_do\\_i\\_write\\_a\\_parser\\_that\\_doesnt\\_crash\\_or/?tl=es-es](https://www.reddit.com/r/Compilers/comments/1jrmwrs/how_do_i_write_a_parser_that_doesnt_crash_or/?tl=es-es)
9. Difference Between Top Down Parsing and Bottom Up Parsing - GeeksforGeeks, acceso: noviembre 4, 2025,  
<https://www.geeksforgeeks.org/compiler-design/difference-between-top-down-parsing-and-bottom-up-parsing/>
10. What is the difference between a top down and bottom up parser? - Quora, acceso: noviembre 4, 2025,  
<https://www.quora.com/What-is-the-difference-between-a-top-down-and-bottom-up-parser>
11. LL parser - Wikipedia, acceso: noviembre 4, 2025,  
[https://en.wikipedia.org/wiki/LL\\_parser](https://en.wikipedia.org/wiki/LL_parser)
12. LALR parser - Wikipedia, acceso: noviembre 4, 2025,  
[https://en.wikipedia.org/wiki/LALR\\_parser](https://en.wikipedia.org/wiki/LALR_parser)

13. Parsing - FING, acceso: noviembre 4, 2025,  
<https://www.fing.edu.uy/~gustun/CLFP/teoricos/clase5.ppt>
14. Técnicas deductivas para el análisis sintáctico con corrección de errores - ResearchGate, acceso: noviembre 4, 2025,  
[https://www.researchgate.net/publication/39437670\\_Tecnicas\\_deductivas\\_para\\_el\\_analisis\\_sintactico\\_con\\_correccion\\_de\\_errores](https://www.researchgate.net/publication/39437670_Tecnicas_deductivas_para_el_analisis_sintactico_con_correccion_de_errores)
15. What are the differences between PEGs and CFGs? - Stack Overflow, acceso: noviembre 4, 2025,  
<https://stackoverflow.com/questions/5501074/what-are-the-differences-between-peg-and-cfgs>
16. Parsing expression grammar - Wikipedia, acceso: noviembre 4, 2025,  
[https://en.wikipedia.org/wiki/Parsing\\_expression\\_grammar](https://en.wikipedia.org/wiki/Parsing_expression_grammar)
17. TRANSFORMACIÓN AUTOMÁTICA DE TEXTO A GRAFOS CONCEPTUALES Sonia Ordoñez Salinas - Alexander Gelbukh, acceso: noviembre 4, 2025,  
<https://www.gelbukh.com/thesis/Sonia%20Ordonez%20Salinas%20-%20PhD.pdf>
18. Grafos de análisis sintáctico con gramáticas HRGs - Research in Computing Science, acceso: noviembre 4, 2025,  
[https://rcs.cic.ipn.mx/2019\\_148\\_3/Grafos%20de%20analisis%20sintactico%20con%20gramaticas%20HRGs.pdf](https://rcs.cic.ipn.mx/2019_148_3/Grafos%20de%20analisis%20sintactico%20con%20gramaticas%20HRGs.pdf)
19. Cómo usar la IA para convertir datos no estructurados en información procesable - Avrio, acceso: noviembre 4, 2025,  
<https://avriodata.ai/es/resources/unstructured-data-with-the-modern-data-stack>
20. LLM-Powered Parsing and Analysis of Semi-Structured & Structured Documents, acceso: noviembre 4, 2025,  
<https://towardsdatascience.com/llm-powered-parsing-and-analysis-of-semi-structured-structured-documents-f03ac92f063e/>
21. Métodos de inteligencia artificial para la clasificación de documentos en español, acceso: noviembre 4, 2025,  
[http://www.scielo.org.pe/scielo.php?script=sci\\_arttext&pid=S2413-26592024000200043](http://www.scielo.org.pe/scielo.php?script=sci_arttext&pid=S2413-26592024000200043)
22. Workbench de Document AI - Google Cloud, acceso: noviembre 4, 2025,  
<https://cloud.google.com/document-ai-workbench?hl=es>
23. Structured Output as a Full Replacement for Function Calling | by Vitaly Sem | Medium, acceso: noviembre 4, 2025,  
<https://medium.com/@virtualik/structured-output-as-a-full-replacement-for-function-calling-430bf98be686>
24. Function Calling is Structured Output - Reilly Wood, acceso: noviembre 4, 2025,  
<https://www.reillywood.com/blog/function-calling-is-structured-output/>
25. Análisis sintáctico funcional: principios, perspectivas y casos, acceso: noviembre 4, 2025, <http://revistas.uach.cl/index.php/efilolo/article/view/1495/1573>