

# Análisis Exhaustivo de Polars: Arquitectura, Métodos y Ventajas Comparativas

## I. Introducción: Polars como un Motor Moderno de DataFrames

En el ecosistema de la ciencia de datos y la ingeniería de datos, la manipulación eficiente de datos estructurados es una piedra angular. Durante años, la librería pandas ha sido la herramienta predominante en el mundo de Python para estas tareas. Sin embargo, el crecimiento exponencial en el volumen de datos y la creciente demanda de un procesamiento más rápido han revelado las limitaciones inherentes a su arquitectura. En este contexto, emerge Polars, una librería de DataFrames que no se presenta como una simple mejora incremental, sino como una reimaginación fundamental de cómo se debe abordar el procesamiento de datos en hardware moderno.

### 1.1. Definiendo Polars: Más Allá de un "Pandas más Rápido"

Polars es una librería de código abierto diseñada desde cero para la manipulación de datos estructurados a alta velocidad.<sup>1</sup> Aunque su interfaz principal es para Python, su núcleo está escrito en Rust y también ofrece APIs para Rust y JavaScript/NodeJS, lo que demuestra su versatilidad y su creciente adopción en diferentes ecosistemas de desarrollo.<sup>1</sup>

La descripción de Polars como una simple alternativa a pandas subestima su verdadera naturaleza. Es más preciso definirlo como un **motor de consultas analíticas** para DataFrames.<sup>2</sup> Esta terminología, comúnmente asociada con bases de datos, es deliberada y fundamental para comprender su diseño. Polars internaliza conceptos avanzados como la optimización de consultas, la ejecución en paralelo y el procesamiento fuera del núcleo (out-of-core), funcionalidades que en el ecosistema de pandas a menudo requieren la integración de herramientas externas como Dask o un sistema de base de datos completo.<sup>4</sup>

El núcleo de Polars, al estar escrito en Rust, le confiere un rendimiento comparable al de lenguajes de bajo nivel como C o C++, junto con garantías de seguridad de memoria y un control granular sobre el entorno de ejecución.<sup>1</sup> Esta es una decisión arquitectónica clave que le permite superar las limitaciones de rendimiento inherentes a lenguajes interpretados de más alto nivel como Python, especialmente en lo que respecta a la concurrencia.<sup>6</sup>

## 1.2. La Filosofía Central: Una Tríada de Rendimiento, Predictibilidad y Rigurosidad

La filosofía de diseño de Polars se puede resumir en tres principios fundamentales que guían su desarrollo y lo diferencian de sus predecesores.

1. **Rendimiento:** El objetivo principal de Polars es ser una librería de DataFrames "ultrarrápida" (*blazingly fast*). Para lograrlo, está diseñada para utilizar todos los núcleos de CPU disponibles en una máquina, optimizar las consultas para reducir el trabajo redundante y las asignaciones de memoria innecesarias, y procesar datos en fragmentos para minimizar la contención en el paralelismo.<sup>2</sup>
2. **Predictibilidad:** La API de Polars busca ser consistente y predecible. Esto se manifiesta en un diseño que a menudo ofrece una única forma idiomática y clara de realizar una operación, en contraste con la API de pandas, que a veces es percibida como sobrecargada con múltiples métodos para lograr el mismo resultado.<sup>8</sup> Un aspecto central de esta filosofía es la eliminación deliberada del índice de filas, un concepto fundamental en pandas. En Polars, las filas se identifican por su posición entera, lo que simplifica la semántica de las consultas y evita que el resultado de una operación dependa del estado de un índice.<sup>4</sup>
3. **Esquema Estricto:** Polars opera bajo un paradigma de esquema estricto, lo que significa que los tipos de datos de cada columna deben ser conocidos antes de que se ejecute una consulta.<sup>2</sup> Esta rigurosidad, aunque puede requerir una mayor explicitud por parte del usuario, es un prerequisite para el funcionamiento del optimizador de consultas. Garantiza la integridad de los datos, minimiza los errores en tiempo de ejecución debidos a la coerción de tipos inesperada y permite al motor planificar la ejecución de la manera más eficiente posible.<sup>2</sup>

## 1.3. Visión General de las Estructuras de Datos Clave

La arquitectura de Polars se basa en un conjunto de estructuras de datos jerárquicas y eficientes, diseñadas para el procesamiento columnar.

- **DataFrame:** Es la estructura de datos principal, bidimensional y tabular, análoga a una tabla en una base de datos SQL. Internamente, se puede conceptualizar como una abstracción sobre un `Vec<Series>`, es decir, un vector o lista de objetos Series, lo que subraya su naturaleza fundamentalmente columnar.<sup>11</sup> La salida impresa de un DataFrame de Polars incluye información útil como las dimensiones (filas, columnas) y el tipo de dato (`dtype`) debajo de cada nombre de columna, proporcionando un contexto inmediato sobre la estructura de los datos.<sup>7</sup>

- **Series:** Representa una única columna de datos. Es la estructura de datos unidimensional fundamental de Polars. Una Series es agnóstica al tipo en su interfaz, actuando como un contenedor que permite aplicar operaciones de manera uniforme sobre diferentes tipos de datos subyacentes. La mayoría de las operaciones en Polars se definen como una función que transforma una o más Series de entrada en una Series de salida.<sup>11</sup>
- **ChunkedArray<T>:** Es la estructura de datos raíz y la implementación concreta que respalda a una Series. Una Series es esencialmente un contenedor sin tipo genérico alrededor de un ChunkedArray<T>, donde T es el tipo de dato específico (por ejemplo, i64, f64, str). A su vez, un ChunkedArray<T> es un contenedor que agrupa uno o más arrays de Apache Arrow (Vec<dyn ArrowArray>).<sup>11</sup> Este diseño "fragmentado" o "en trozos" (*chunked*) es crucial para la eficiencia: permite que los datos de una columna se almacenen en múltiples bloques de memoria no contiguos. Esto facilita operaciones como la adición de datos sin necesidad de reasignar toda la columna en memoria, y es fundamental para el motor de streaming que procesa los datos en lotes.<sup>11</sup>

La relación entre estas estructuras revela una estrategia de diseño deliberada. Al construir el DataFrame a partir de Series independientes y columnarmente orientadas, y al basar las Series en la especificación de memoria de Apache Arrow a través de los ChunkedArray, Polars establece una base arquitectónica que es inherentemente propicia para el paralelismo y la eficiencia de la memoria caché, diferenciándose fundamentalmente de los enfoques basados en arrays de NumPy orientados a filas.

## II. Los Fundamentos Arquitectónicos del Rendimiento de Polars

El rendimiento excepcional de Polars no es el resultado de una única optimización, sino de una pila de tecnologías y decisiones de diseño sinérgicas que trabajan en conjunto. Cada capa de su arquitectura, desde el lenguaje de programación del núcleo hasta el modelo de memoria, está elegida para maximizar la velocidad y la eficiencia en el hardware moderno.

### 2.1. El Núcleo en Rust: Aprovechando la Velocidad y la Concurrencia a Nivel de Sistema

La decisión más fundamental en la arquitectura de Polars es el uso de Rust para su núcleo computacional.<sup>1</sup> Rust es un lenguaje de programación de sistemas que ofrece un rendimiento comparable al de C y C++, pero con garantías de seguridad de memoria que eliminan clases enteras de errores comunes en esos lenguajes.<sup>5</sup> Esto permite a Polars operar "cerca del metal" (

*close to the metal*), con un control total sobre la gestión de la memoria y la ejecución de las operaciones, y sin dependencias externas para sus cálculos principales.<sup>1</sup>

Una de las ventajas más significativas de Rust sobre Python para tareas de computación intensiva es la ausencia de un **Bloqueo Global del Intérprete (Global Interpreter Lock, GIL)**.<sup>6</sup> El GIL en CPython (la implementación estándar de Python) es un mecanismo que impide que múltiples hilos ejecuten código de Python en paralelo dentro de un mismo proceso. Aunque librerías como pandas delegan algunas operaciones a código C que puede liberar el GIL, su núcleo sigue estando fundamentalmente limitado a un solo hilo para muchas tareas.<sup>4</sup> Polars, al estar escrito en Rust, no tiene esta limitación. Puede lograr un paralelismo real y nativo, permitiendo que múltiples hilos trabajen simultáneamente en diferentes partes de una consulta, compartiendo memoria de forma segura y escalando el rendimiento de manera efectiva con el número de núcleos de CPU disponibles.<sup>5</sup>

## 2.2. El Modelo de Memoria Apache Arrow: La Ventaja Columnar

Polars se basa en el modelo de memoria de Apache Arrow, una especificación estandarizada e independiente del lenguaje para representar datos tabulares en memoria en formato columnar.<sup>1</sup> Esta elección tiene profundas implicaciones para el rendimiento.

- **Disposición Columnar:** A diferencia de los arrays de NumPy utilizados por defecto en pandas, que suelen tener una disposición orientada a filas, Arrow almacena los datos columna por columna en bloques de memoria contiguos.<sup>5</sup> Las consultas analíticas (por ejemplo, calcular la media de una columna, filtrar por el valor de otra) operan típicamente sobre un subconjunto de columnas. La disposición columnar es extremadamente eficiente para la caché de la CPU en estos casos, ya que los datos necesarios para una operación sobre una columna están juntos en la memoria, minimizando los saltos de memoria (*cache misses*).<sup>11</sup>
- **Operaciones de Copia Cero (Zero-Copy):** El formato estandarizado de Arrow permite un intercambio de datos de "copia cero" entre Polars y otras herramientas compatibles con Arrow, como DuckDB, PySpark o librerías de GPU como CuDF.<sup>1</sup> Esto significa que los datos pueden ser compartidos entre diferentes procesos o librerías sin la necesidad de copiarlos o de pasar por costosos procesos de serialización y deserialización. Dado que se estima que la serialización/deserialización puede representar entre el 80% y el 90% de los costos computacionales en los flujos de trabajo de datos, evitar este paso proporciona una ganancia de rendimiento masiva y una mayor eficiencia de memoria.<sup>5</sup>

Es crucial destacar una distinción importante: aunque Polars utiliza la *especificación* de memoria de Arrow, ha desarrollado su propia implementación personalizada en Rust de los núcleos computacionales y los búferes de memoria.<sup>2</sup> No está construido sobre una librería como PyArrow, que es el backend que utiliza pandas 2.0. Esta decisión estratégica, aunque implica una mayor carga de mantenimiento, otorga al equipo de Polars un control absoluto sobre las rutas de código críticas para el rendimiento, permitiéndoles

optimizar, corregir errores e implementar nuevas características sin depender del ciclo de lanzamiento o las decisiones de diseño de otro proyecto. Esto es un componente clave de su filosofía de estar "cerca del metal".<sup>1</sup>

### 2.3. "Vergonzosamente Paralelo": El Motor de Consultas Multihilo

Polars está diseñado desde su concepción para ser "vergonzosamente paralelo" (*embarrassingly parallel*), lo que significa que paraleliza automáticamente las cargas de trabajo en todos los núcleos de CPU disponibles sin necesidad de configuración adicional por parte del usuario.<sup>1</sup>

Este paralelismo no se limita a operaciones simples. Gracias a la API de Expresiones (que se detallará más adelante), el motor de Polars puede analizar el plan de consulta completo, identificar cómputos independientes y distribuirlos entre los hilos disponibles.<sup>5</sup> Por ejemplo, si un usuario crea tres nuevas columnas en una sola operación, y ninguna de ellas depende de las otras, Polars puede calcular las tres en paralelo. Esta capacidad de paralelización por defecto es una de las mayores ventajas de rendimiento sobre el núcleo monohilo de pandas.<sup>4</sup>

### 2.4. Procesamiento Vectorizado con SIMD

A un nivel aún más bajo, el motor de consultas de Polars aprovecha las instrucciones **SIMD (Single Instruction, Multiple Data)**, una característica presente en las CPUs modernas.<sup>1</sup> Las instrucciones SIMD permiten que una única instrucción de la CPU (como una suma o una multiplicación) se aplique simultáneamente a múltiples puntos de datos (un vector de números).

Esta forma de paralelismo a nivel de instrucción es perfectamente adecuada para la disposición de datos columnar de Apache Arrow. Como los valores de una columna numérica están almacenados en un bloque de memoria contiguo, pueden ser cargados en los registros SIMD de la CPU y procesados en lotes, lo que resulta en una aceleración significativa de los cálculos numéricos en comparación con la ejecución de la misma operación de forma secuencial, elemento por elemento.<sup>1</sup>

En conjunto, esta pila tecnológica —un núcleo en Rust sin GIL, un modelo de memoria columnar eficiente para la caché, un motor multihilo por defecto y el uso de vectorización SIMD— crea un efecto de rendimiento compuesto. Cada capa refuerza a las demás, dando como resultado un sistema que está intrínsecamente diseñado para extraer el máximo rendimiento del hardware contemporáneo.

## III. La API de Polars: Un Cambio de Paradigma en la Manipulación de Datos

Más allá de su arquitectura interna, lo que define la experiencia de usuario en Polars es su API. Esta no es simplemente una imitación de la API de pandas, sino que introduce un cambio de paradigma centrado en dos conceptos clave: la distinción entre ejecución *lazy* y *eager*, y el poder de una API de Expresiones declarativa. Juntos, estos conceptos permiten a los usuarios escribir código que no solo es legible y expresivo, sino también altamente optimizable y paralelizable.

### 3.1. Ejecución Lazy vs. Eager: Los Dos Modos de Operación

Polars ofrece dos modos de ejecución distintos, cada uno adecuado para diferentes casos de uso.<sup>10</sup>

- **Modo Eager (Ansioso o Inmediato):** Este es el modo de ejecución por defecto y es conceptualmente similar a cómo funciona pandas. Cada operación se ejecuta inmediatamente en el momento en que se invoca, y el resultado está disponible al instante.<sup>15</sup> Este modo es ideal para la exploración de datos interactiva, por ejemplo, en un Jupyter Notebook. Permite a los usuarios inspeccionar los resultados intermedios de cada paso de su análisis, lo cual es valioso para comprender y depurar los datos.<sup>17</sup>
- **Modo Lazy (Perezoso o Diferido):** Este es el modo preferido para scripts, pipelines de producción y cualquier carga de trabajo donde el rendimiento es crítico.<sup>15</sup> En el modo lazy, las operaciones no se ejecutan de inmediato. En su lugar, cada transformación se añade a un **plan de consulta lógico**, que es una representación abstracta de la secuencia de operaciones deseadas. La computación real se difiere hasta que se invoca explícitamente el método `.collect()`.<sup>13</sup>

El flujo de trabajo lazy típicamente comienza con los métodos `pl.scan_*` (por ejemplo, `pl.scan_csv`, `pl.scan_parquet`). A diferencia de los métodos `pl.read_*` del modo eager, `scan` no lee el archivo de datos completo en memoria. Simplemente escanea los metadatos del archivo (como el esquema de columnas) y devuelve un objeto `LazyFrame`, que es un marcador de posición para el conjunto de datos.<sup>2</sup> A partir de este

`LazyFrame`, se pueden encadenar una serie de transformaciones (`.filter()`, `.with_columns()`, `.group_by()`, etc.). Ninguna de estas operaciones consume CPU o memoria significativa. Solo cuando se llama a `.collect()`, Polars toma el plan de consulta lógico acumulado, lo pasa a su **optimizador de consultas**, genera un plan de ejecución físico eficiente y finalmente ejecuta la consulta para producir un `DataFrame` real.<sup>16</sup>

El impacto en el rendimiento de este enfoque es drástico. Al diferir la ejecución, Polars obtiene una visión global de toda la cadena de operaciones. Esto le permite aplicar optimizaciones potentes que son imposibles en un modelo de ejecución paso a paso. Los benchmarks demuestran que las consultas lazy pueden ser entre 5 y 12 veces más rápidas que sus equivalentes en modo eager, y la diferencia puede ser aún mayor para operaciones de I/O complejas.<sup>15</sup>

## 3.2. El Poder de las Expresiones

El segundo pilar de la API de Polars es su sistema de Expresiones. Las expresiones son el lenguaje que se utiliza para construir los planes de consulta y comunicar la intención del usuario al motor de Polars.

- **Definición:** Formalmente, una expresión es un mapeo funcional que toma una Series como entrada y produce una Series como salida (es decir, una función con la firma `Fn(Series) -> Series`).<sup>11</sup> Son los bloques de construcción fundamentales de casi todas las operaciones de manipulación de datos en Polars.
- **Sintaxis Central:** Las expresiones se crean típicamente usando `pl.col("nombre_columna")` para referirse a una o más columnas, o `pl.lit(valor)` para referirse a un valor literal (constante). Estas expresiones se pueden combinar usando operadores estándar de Python (`+`, `*`, `==`, `>`) o encadenarse con una vasta biblioteca de métodos específicos de Polars (`.sum()`, `.alias("nuevo_nombre")`, `.cast(pl.Int64)`, etc.).<sup>7</sup>
- **Contextos de Ejecución:** El verdadero poder de las expresiones se manifiesta cuando se aplican dentro de "contextos" o "verbos" específicos del DataFrame. Los principales contextos son:
  - `select()`: Para seleccionar, renombrar o crear nuevas columnas a partir de expresiones.
  - `filter()`: Para seleccionar filas basándose en una expresión que evalúa a un booleano.
  - `with_columns()`: Para añadir o sobrescribir columnas existentes con nuevas expresiones.
  - `group_by().agg()`: Para aplicar expresiones de agregación dentro de cada grupo de datos.<sup>7</sup>

La combinación de la API de Expresiones y el modo Lazy es lo que desbloquea el máximo potencial de Polars. La API de Expresiones proporciona una forma estructurada y declarativa de describir una consulta, mientras que la API Lazy proporciona el tiempo de ejecución que difiere la ejecución para que esta descripción pueda ser analizada, optimizada y luego ejecutada de la manera más eficiente posible. Cuando un usuario escribe `df.filter(pl.col("a") > 5)`, no está ejecutando un filtro inmediatamente; está construyendo un nodo en un grafo de computación que representa la operación de filtrado. Este grafo es el plan de consulta lógico que el optimizador de Polars reescribirá para obtener el máximo rendimiento.

Los beneficios de esta API de Expresiones son múltiples:

1. **Paralelismo:** Cuando se aplican múltiples expresiones en el mismo contexto (por ejemplo, `df.with_columns([ (pl.col("a") * 2).alias("a2"), (pl.col("b") + 1).alias("b1") ])`), Polars puede ejecutar estas expresiones en paralelo, ya que entiende que son computacionalmente independientes.<sup>12</sup>
2. **Optimización:** La API de Expresiones es el lenguaje que el optimizador de consultas entiende. Permite a Polars realizar optimizaciones como el *predicate pushdown* y el

*projection pushdown* (detallados más adelante), que son la clave de su eficiencia en operaciones de I/O.<sup>12</sup>

3. **Legibilidad:** El encadenamiento de expresiones a menudo conduce a un código más declarativo ("qué quiero lograr") en lugar de imperativo ("cómo hacerlo paso a paso"). Esto puede hacer que las transformaciones de datos complejas sean más fáciles de leer y mantener en comparación con el código de pandas equivalente, que podría requerir la creación de variables intermedias o el uso de una sintaxis de indexación más compleja.<sup>1</sup>

## IV. Un Catálogo Completo de Métodos y Operaciones de Polars

La interfaz de programación de aplicaciones (API) de Polars es exhaustiva y abarca un conjunto integral de herramientas para el ciclo de vida del análisis de datos, desde la adquisición hasta la exportación. Esta sección presenta las principales categorías de métodos y operaciones, enfatizando la sintaxis idiomática basada en expresiones.

### 4.1. Ingesta y Exportación de Datos: Un Subsistema de I/O Versátil

Polars ofrece soporte de primera clase para una amplia gama de formatos de datos, lo que facilita su integración en pilas de datos existentes.<sup>1</sup> Para cada formato, generalmente existen dos variantes de lectura: una *eager* (`read_*`) que carga los datos inmediatamente en un `DataFrame`, y una *lazy* (`scan_*`) que devuelve un `LazyFrame` para su posterior optimización.

- **Formatos Soportados:**

- **Texto:** CSV (`read_csv`, `scan_csv`) y JSON (`read_json`).<sup>1</sup>
- **Binario Columnar:** Parquet (`read_parquet`, `scan_parquet`), Delta Lake (`scan_delta`), y Feather/Arrow IPC (`read_ipc`, `scan_ipc`). Estos formatos son altamente recomendados para el rendimiento debido a su naturaleza columnar.<sup>1</sup>
- **Otros Formatos Binarios:** AVRO (`read_avro`) y Excel (`read_excel`).<sup>1</sup>
- **Bases de Datos:** Polars puede leer datos directamente desde bases de datos SQL a través de conectores para sistemas como PostgreSQL, MySQL, SQL Server, SQLite, Redshift y Oracle.<sup>1</sup>
- **Almacenamiento en la Nube:** Soporte nativo para leer y escribir en sistemas de almacenamiento en la nube como Amazon S3, Azure Blob Storage y Azure File Storage.<sup>1</sup>

### 4.2. Verbos de Manipulación Central



Estos son los métodos principales utilizados para dar forma y transformar un DataFrame, contruidos sobre la API de Expresiones.

- `select()`: Se utiliza para seleccionar un subconjunto de columnas, reordenarlas o crear nuevas columnas a partir de expresiones. Por ejemplo: `df.select([pl.col("a"), pl.col("b").alias("b_renamed")])`.<sup>7</sup>
- `filter()`: Selecciona filas que cumplen una o más condiciones booleanas. Las condiciones se expresan como predicados. Ejemplo: `df.filter(pl.col("quantity") > 50)`.<sup>7</sup>
- `with_columns()`: Es la forma idiomática de añadir nuevas columnas o sobrescribir las existentes. Es más eficiente que crear columnas una por una. Ejemplo: `df.with_columns((pl.col("price") * pl.col("quantity")).alias("total_cost"))`.<sup>7</sup>
- `drop()`: Elimina una o más columnas del DataFrame. Ejemplo: `df.drop("columna_innecesaria")`.<sup>19</sup>
- `sort()`: Ordena el DataFrame según una o más columnas. Ejemplo: `df.sort("fecha", descending=True)`.<sup>11</sup>

### 4.3. Agregación y Agrupamiento

Polars proporciona una sintaxis potente y expresiva para operaciones de agregación, similar en concepto a GROUP BY en SQL.

- `group_by()`: Es el punto de entrada para las operaciones de agrupamiento. Por sí solo, crea un objeto `GroupBy`. Se combina con el método `.agg()` para realizar los cálculos. Ejemplo: `df.group_by("categoria")`.<sup>7</sup>
- `agg()`: Se aplica a un objeto `GroupBy` y toma una lista de expresiones de agregación. Cada expresión define un cálculo a realizar en cada grupo. Ejemplo: `df.group_by("categoria").agg([pl.col("ventas").sum().alias("ventas_totales"), pl.col("producto").n_unique().alias("productos_unicos"), pl.len().alias("conteo_filas")])`.<sup>7</sup>
- `pivot()`: Crea una tabla dinámica estilo hoja de cálculo a partir de los datos, transformando valores de una columna en nuevas columnas. Ejemplo: `df.pivot(values="ventas", index="fecha", columns="tienda", aggregate_function="sum")`.<sup>19</sup>
- `group_by_dynamic()` y `rolling()`: Funciones especializadas para agregaciones en ventanas de tiempo o secuenciales, extremadamente potentes para el análisis de series temporales.<sup>7</sup>

### 4.4. Unión y Fusión de DataFrames

Polars implementa un conjunto completo de operaciones de unión (`join`) de estilo SQL.

- `join()`: Realiza uniones entre dos DataFrames. Soporta diferentes estrategias como `inner` (por defecto), `left`, `outer`, `cross`, `semi` y `anti`. Ejemplo: `df1.join(df2, on="id_cliente",`

how="left").<sup>19</sup>

- **join\_asof():** Realiza una unión "asof", que es una operación especializada para series temporales. En lugar de buscar una coincidencia exacta en las claves, busca el valor más cercano en el tiempo sin pasarse. Esto es muy útil para alinear observaciones que no ocurren exactamente al mismo tiempo.<sup>11</sup>
- **concat():** Concatena DataFrames, ya sea verticalmente (apilando filas) u horizontalmente (añadiendo columnas).

## 4.5. Espacios de Nombres de Expresiones Especializadas

Para manejar tipos de datos complejos, Polars ofrece "espacios de nombres" (namespaces) dentro de su API de Expresiones. Estos agrupan funcionalidades específicas para cada tipo de dato, lo que resulta en una API limpia y organizada.

- **str (Cadenas de texto):** Se accede a través de `pl.col("col_texto").str`. Proporciona una rica funcionalidad para la manipulación de cadenas, incluyendo métodos como `.contains()`, `.extract()`, `.replace()`, `.to_lowercase()`, `.len_chars()`, `.json_path_match()` y muchos más.<sup>20</sup>
- **dt (Fecha y hora):** Se accede a través de `pl.col("col_fecha").dt`. Ofrece métodos para extraer componentes de fechas y horas (`.year()`, `.month()`, `.weekday()`), formatear (`.strftime()`) y realizar aritmética de fechas.<sup>7</sup>
- **list y arr (Listas):** Se accede a través de `pl.col("col_lista").list` o `.arr`. Permite realizar operaciones sobre columnas que contienen listas, como obtener la longitud (`.len()`), sumar los elementos (`.sum()`), acceder a un elemento por índice (`.get()`) o expandir la lista en filas separadas (`.explode()`).<sup>20</sup>
- **struct (Estructuras):** Se accede a través de `pl.col("col_struct").struct`. Permite trabajar con tipos de datos anidados (structs), accediendo a sus campos internos (`.field("nombre_campo")`).

La siguiente tabla resume algunas de las expresiones más comunes, organizadas por categoría, para servir como una guía de referencia rápida.

Categoría	Expresiones / Métodos Comunes	Ejemplo de Uso
<b>Selección y Renombrado</b>	<code>pl.col()</code> , <code>pl.all()</code> , <code>.alias()</code> , <code>.exclude()</code>	<code>df.select([pl.all().exclude("id"), pl.col("name").alias("nombre")])</code>
<b>Transformación y Casting</b>	<code>+</code> , <code>-</code> , <code>*</code> , <code>/</code> , <code>.cast()</code> , <code>.round()</code>	<code>df.with_columns((pl.col("price") * 1.21).cast(pl.Float32).alias("precio_con_iva"))</code>
<b>Agregación</b>	<code>.sum()</code> , <code>.mean()</code> , <code>.median()</code> , <code>.max()</code> , <code>.min()</code> , <code>.std()</code> , <code>.count()</code> , <code>.n_unique()</code> , <code>pl.len()</code>	<code>df.group_by("category").agg(pl.col("value").mean())</code>

<b>Condicionales</b>	<code>pl.when().then().otherwise()</code>	<code>df.select(pl.when(pl.col("score") &gt; 5).then("Aprobado").otherwise("Reprobado"))</code>
<b>Manipulación de Cadenas</b>	<code>.str.contains(), .str.extract(), .str.replace_all(), .str.to_uppercase()</code>	<code>df.filter(pl.col("product_name").str.contains("(?i)special"))</code>
<b>Operaciones de Fecha/Hora</b>	<code>.dt.year(), .dt.month(), .dt.weekday(), .dt.strftime(), .dt.total_days()</code>	<code>df.with_columns(pl.col("timestamp").dt.strftime("%Y-%m").alias("año_mes"))</code>
<b>Operaciones de Lista</b>	<code>.list.len(), .list.sum(), .list.get(), .list.explode(), .list.join()</code>	<code>df.select(pl.col("tags").list.len().alias("num_tags"))</code>
<b>Funciones de Ventana</b>	<code>.over()</code>	<code>df.with_columns(pl.col("sales").sum().over("department").alias("ventas_departamento"))</code>

## V. Análisis Comparativo: Polars vs. Pandas

La pregunta más frecuente al considerar Polars es cómo se compara con pandas, el estándar de facto en el ecosistema de Python. La comparación va más allá de la simple velocidad; abarca diferencias fundamentales en arquitectura, filosofía de API y gestión de recursos.

### 5.1. Benchmarks de Rendimiento: Una Síntesis Cuantitativa

Múltiples benchmarks independientes y de los propios desarrolladores de Polars demuestran una ventaja de rendimiento significativa y consistente en una amplia gama de operaciones, especialmente a medida que aumenta el tamaño de los datos.

- **Lectura y Escritura de Datos:** En la lectura de archivos, Polars muestra una superioridad notable. Para un archivo CSV grande (casi 3 GB), Polars fue aproximadamente **11 veces más rápido** que pandas (1.25 segundos vs. 14.14 segundos). Incluso para archivos pequeños, Polars fue 2.5 veces más rápido. En la escritura de datos, la ventaja se mantiene, con Polars siendo aproximadamente **2 veces más rápido** que pandas para archivos grandes.<sup>22</sup> Esta diferencia se debe a la implementación optimizada en Rust y a la capacidad de paralelizar la lectura y el parseo.
- **Operaciones de GroupBy y Agregación:** Esta es una de las áreas donde Polars brilla con más intensidad. En conjuntos de datos de tamaño mediano a grande, Polars supera a pandas por un factor de **4.7x a 8.6x**.<sup>23</sup> La razón es su motor multihilo, que puede

procesar las agregaciones de diferentes grupos en paralelo, y sus algoritmos de hashing eficientes.

- **Uniones (Joins):** Las uniones de datos también son significativamente más rápidas en Polars. Los benchmarks muestran que para unir DataFrames de 10,000 filas, Polars puede ser hasta **13.75 veces más rápido** que pandas.<sup>23</sup> De nuevo, los algoritmos de join optimizados y la paralelización son los factores clave.
- **Selección de Columnas y Filtros:** La selección de columnas es una operación donde la arquitectura columnar de Polars ofrece una ventaja masiva. Es entre **6 y 11 veces más rápido** que pandas, ya que solo necesita acceder a los bloques de memoria de las columnas deseadas.<sup>23</sup> En las operaciones de filtrado, la ventaja es de **3 a 4 veces** a favor de Polars en datasets grandes.<sup>23</sup>

El benchmark PDS-H, una derivación del estándar industrial TPC-H, sitúa a Polars en la cima de las librerías de DataFrames de una sola máquina. En la escala SF-10, el motor de streaming de Polars completó el benchmark en 3.89 segundos, mientras que pandas tardó 365.71 segundos, lo que representa una diferencia de aproximadamente **94 veces**.<sup>24</sup> Es importante notar que, para conjuntos de datos muy pequeños (por ejemplo, 10,000 filas), pandas puede ser ocasionalmente más rápido en operaciones de filtrado muy simples, pero Polars domina en cualquier operación más compleja o a cualquier escala mayor.<sup>23</sup>

La siguiente tabla resume los resultados de los benchmarks para una comparación rápida.

Operación	Conjunto de Datos	Rendimiento Polars	Rendimiento Pandas	Factor de Aceleración (Polars vs. Pandas)
<b>Lectura Parquet</b>	Grande (~905 MB)	1.25 s	14.14 s	<b>11.3x</b>
<b>Escritura Parquet</b>	Grande (~905 MB)	10.39 s	20.55 s	<b>2.0x</b>
<b>GroupBy &amp; Agg</b>	Grande (1,000,000 filas)	~0.07 s	~0.60 s	<b>8.6x</b>
<b>Join</b>	Mediano (100,000 filas)	~0.02 s	~0.27 s	<b>13.8x</b>
<b>Selección de Columna</b>	Grande (1,000,000 filas)	~0.001 s	~0.01 s	<b>10.9x</b>
<b>Benchmark PDS-H (SF-10)</b>	Complejo (10 GB)	3.89 s (streaming)	365.71 s	<b>94.0x</b>

## 5.2. Divergencia Arquitectónica

Las diferencias de rendimiento son una consecuencia directa de elecciones arquitectónicas fundamentalmente distintas.

- **Núcleo y Modelo de Memoria:** Polars tiene un núcleo en Rust y utiliza el modelo de memoria columnar de Apache Arrow. Pandas tiene un núcleo en Python/C y utiliza por defecto arrays de NumPy, que tienen una disposición orientada a filas.<sup>4</sup> Aunque pandas 2.0 puede usar PyArrow como backend, la integración no es tan profunda ni el motor está diseñado desde cero en torno a este paradigma como en Polars.
- **Paralelismo:** Polars es multihilo por defecto y de forma nativa. Pandas es fundamentalmente monohilo en su núcleo y requiere librerías externas como Dask para escalar a múltiples núcleos.<sup>4</sup>
- **Gestión de Memoria:** Polars es significativamente más eficiente en el uso de la memoria. Generalmente, requiere entre 2 y 4 veces la cantidad de RAM del tamaño del conjunto de datos para realizar operaciones, en comparación con las 5 a 10 veces que puede necesitar pandas.<sup>5</sup> Esto se debe a algoritmos que minimizan las copias de datos, la eficiencia del formato Arrow y la capacidad de su motor de streaming para procesar datos en lotes.

### 5.3. Diferencias Conceptuales y de API

Más allá del rendimiento, existen diferencias filosóficas en el diseño de la API que afectan la forma en que los usuarios interactúan con la librería.

- **La Ausencia del Índice:** Esta es quizás la diferencia conceptual más importante. Pandas se centra en el concepto de un Index (y MultiIndex), que etiqueta cada fila. Polars elimina este concepto por completo. Las filas se identifican únicamente por su posición entera (de 0 a n-1). La justificación de Polars es que la ausencia de un índice conduce a una API más simple, predecible y legible. Las operaciones no cambian su comportamiento dependiendo del estado del índice (por ejemplo, si está ordenado o si es único), y se elimina la necesidad de llamadas frecuentes a `.reset_index()` que son comunes en el código de pandas.<sup>4</sup>
- **Consistencia de la API:** La API de Polars es a menudo descrita como más consistente, lógica y menos "sobrecargada".<sup>8</sup> Generalmente hay una forma idiomática y clara de realizar una tarea, basada en el sistema de expresiones. En contraste, pandas a menudo ofrece múltiples formas de lograr el mismo resultado (por ejemplo, para la selección de datos se puede usar ```, `.loc`, `.iloc`, `.at`, `.iat`), lo que puede llevar a inconsistencias y a un código más difícil de razonar.
- **Rigurosidad vs. Flexibilidad:** Polars es más estricto con los tipos de datos. Esto ayuda a prevenir errores sutiles que pueden surgir en pandas debido a su coerción de tipos más flexible y a veces inesperada.<sup>4</sup>

La siguiente tabla ofrece una comparación directa de las características clave.

Característica	Polars	Pandas	Implicaciones
<b>Lenguaje del Núcleo</b>	Rust	Python / C	Rendimiento, seguridad de memoria,

			paralelismo nativo.
<b>Modelo de Memoria</b>	Apache Arrow (Columnar)	NumPy (Orientado a filas por defecto)	Eficiencia de caché, velocidad en operaciones analíticas, uso de memoria.
<b>Indexación de Filas</b>	No tiene (posición entera)	Índice / Multi-índice	Simplicidad de la API, predictibilidad, menos estado que gestionar.
<b>Paralelismo</b>	Multihilo por defecto	Monohilo en el núcleo	Utilización del hardware moderno, escalabilidad en una sola máquina.
<b>Paradigma de API</b>	Lazy / Eager (basado en Expresiones)	Eager (Imperativo)	Optimización de consultas, legibilidad del código, rendimiento.
<b>Optimización de Consultas</b>	Optimizador incorporado (automático)	No tiene (requiere optimización manual)	Rendimiento en I/O, eficiencia en pipelines complejos.
<b>Procesamiento Fuera del Núcleo</b>	Motor de Streaming incorporado	Requiere librerías externas (ej. Dask)	Capacidad para manejar datasets más grandes que la RAM.

## VI. Capacidades Avanzadas: Optimización de Consultas y Procesamiento Fuera del Núcleo

Las características más avanzadas de Polars son las que verdaderamente lo elevan de una simple librería de DataFrames a un motor de consultas completo. Su optimizador de consultas y su motor de streaming permiten un nivel de rendimiento y escalabilidad en una sola máquina que antes estaba reservado para sistemas de computación distribuida.

### 6.1. Dentro del Optimizador de Consultas

El optimizador de consultas es el componente central del motor lazy de Polars. Su función es tomar el plan de consulta lógico (la secuencia de operaciones definida por el usuario) y reescribirlo en un plan de ejecución físico más eficiente, sin cambiar el resultado final.<sup>1</sup> Este proceso es automático y transparente para el usuario. Dos de las optimizaciones más

importantes que realiza son el *predicate pushdown* y el *projection pushdown*.

- **Predicate Pushdown (Empuje de Predicados):** Esta técnica consiste en mover las operaciones de filtrado (los "predicados", como `col("a") > 10`) tan cerca de la fuente de datos como sea posible.<sup>16</sup> En lugar de leer un archivo completo en memoria y luego filtrarlo, el optimizador "empuja" el filtro hacia la capa de lectura.
  - **Para archivos CSV o JSON:** El filtro se aplica a cada fila *mientras se está leyendo y parseando el archivo*. Las filas que no cumplen la condición se descartan inmediatamente, sin llegar a ser asignadas en memoria. Esto reduce drásticamente el uso máximo de memoria y el trabajo computacional posterior.<sup>25</sup>
  - **Para archivos Parquet:** Esta optimización es aún más potente. El formato Parquet almacena metadatos para cada grupo de filas (row group), incluyendo estadísticas como los valores mínimo y máximo de cada columna. Antes de leer los datos, el optimizador de Polars lee estos metadatos. Si un predicado como `col("ventas") > 1000` se aplica, y los metadatos de un grupo de filas indican que el valor máximo de la columna "ventas" en ese grupo es 950, el optimizador sabe que ninguna fila en ese grupo puede cumplir la condición y, por lo tanto, **omite la lectura de ese bloque de datos por completo**.<sup>25</sup>
- **Projection Pushdown (Empuje de Proyecciones):** De manera similar, esta optimización se asegura de que solo las columnas que son estrictamente necesarias para la consulta se lean desde la fuente de datos.<sup>16</sup> Si un usuario escanea un archivo Parquet de 200 columnas pero la consulta final solo utiliza las columnas "usuario" y "fecha", el optimizador instruirá al lector de Parquet para que solo descomprima y cargue en memoria esas dos columnas, ignorando las 198 restantes. Esto ahorra una cantidad inmensa de tiempo de I/O y memoria.

Estas técnicas, que son estándar en bases de datos y sistemas distribuidos como Spark, son democratizadas por Polars y puestas a disposición del científico de datos en su máquina local. La implicación es profunda: permite a los usuarios escribir código de manera intuitiva, y el motor se encarga de reordenarlo para una ejecución óptima.

El método `.explain()` se puede utilizar en un `LazyFrame` para visualizar el plan de consulta optimizado. Esto permite a los usuarios verificar que optimizaciones como el *predicate pushdown* se están aplicando correctamente, proporcionando una ventana transparente al funcionamiento interno del motor.<sup>26</sup>

## 6.2. Motor de Streaming: Procesamiento de Conjuntos de Datos Más Grandes que la RAM

El motor de streaming es una capacidad del motor lazy que permite procesar conjuntos de datos que no caben en la memoria RAM disponible (procesamiento "fuera del núcleo" u *out-of-core*).<sup>1</sup>

En lugar de cargar todos los datos de una vez, el motor de streaming procesa los datos en

lotes o fragmentos de un tamaño manejable.<sup>29</sup> El plan de consulta optimizado se aplica a cada lote de forma secuencial. Los resultados intermedios se van combinando hasta producir el resultado final. Para activar el motor de streaming, el usuario simplemente añade el argumento

`streaming=True` a la llamada `.collect()` (por ejemplo, `query.collect(streaming=True)`).<sup>3</sup>

Una revelación clave de los benchmarks es que el motor de streaming no es solo una solución para problemas de memoria. A menudo, es **más rápido que el motor en memoria** incluso para conjuntos de datos que sí caben en la RAM. Los benchmarks PDS-H muestran que el motor de streaming puede ser de **3 a 7 veces más rápido** que el motor en memoria de Polars.<sup>24</sup> Esto se debe a que el procesamiento en lotes puede hacer un uso más eficiente de la caché de la CPU y permite una mejor utilización de los núcleos, ya que diferentes etapas del pipeline de procesamiento pueden operar en diferentes lotes de datos simultáneamente.

Sin embargo, el motor de streaming tiene limitaciones. No todas las operaciones se pueden realizar de forma puramente en streaming. Por ejemplo, una ordenación global (sort) requiere ver todos los datos a la vez para determinar el orden correcto. En estos casos, Polars ejecutará en streaming todas las partes de la consulta que pueda, y luego materializará los datos necesarios en memoria para realizar la operación de bloqueo (como el sort), antes de continuar, si es posible, en modo streaming.<sup>29</sup>

La combinación del optimizador de consultas y el motor de streaming difumina la línea entre las librerías de DataFrames para análisis en memoria y las herramientas de ETL (Extracción, Transformación y Carga) para el procesamiento de grandes volúmenes de datos. Polars permite realizar complejas tareas de ingeniería de características y agregación (típicas de un DataFrame) directamente sobre flujos de datos que exceden la memoria del sistema (típico de una herramienta ETL), todo dentro de una única API consistente y de alto rendimiento.

## VII. Conclusión: Adopción Estratégica de Polars en el Ecosistema de la Ciencia de Datos

Polars representa una evolución significativa en el campo de las herramientas de procesamiento de datos en una sola máquina. Su diseño, fundamentado en Rust, Apache Arrow y un enfoque de paralelismo por defecto, aborda directamente los cuellos de botella de rendimiento y memoria que han desafiado a las herramientas de la generación anterior. La introducción de un potente optimizador de consultas y un motor de streaming dentro de una API de DataFrame cohesiva lo posiciona como una solución robusta y escalable para los desafíos de datos modernos.

### 7.1. Resumen de Ventajas Clave

El análisis exhaustivo de Polars revela un conjunto de ventajas claras y convincentes:



- **Rendimiento Extremo:** Gracias a su núcleo en Rust, su arquitectura multihilo y la vectorización SIMD, Polars ofrece una velocidad de ejecución órdenes de magnitud superior a la de pandas en la mayoría de las operaciones analíticas comunes.
- **Eficiencia de Memoria:** El uso del formato columnar de Apache Arrow y algoritmos que minimizan las copias de datos resulta en un consumo de RAM significativamente menor, permitiendo trabajar con conjuntos de datos más grandes en el mismo hardware.
- **Escalabilidad Fuera del Núcleo:** Su motor de streaming integrado permite procesar de forma nativa conjuntos de datos que exceden la memoria RAM disponible, una capacidad que tradicionalmente requería el uso de sistemas distribuidos más complejos.
- **API Expresiva y Predictible:** La API basada en expresiones, junto con la eliminación del índice, fomenta un estilo de codificación declarativo, legible y menos propenso a errores, con un comportamiento consistente y predecible.
- **Optimización Automática:** El optimizador de consultas, que aplica de forma transparente técnicas como el *predicate* y *projection pushdown*, libera al usuario de la carga de la optimización manual y garantiza que incluso las consultas complejas se ejecuten de la manera más eficiente posible.

## 7.2. Guía para la Adopción: Cuándo Elegir Polars

La decisión de adoptar Polars debe basarse en las necesidades específicas del proyecto y del equipo.

**Se recomienda elegir Polars cuando:**

- **El rendimiento es crítico:** Para pipelines de producción, aplicaciones en tiempo real o cualquier escenario donde la velocidad de procesamiento de datos impacta directamente en el negocio o los costos de infraestructura en la nube.<sup>5</sup>
- **Los conjuntos de datos son de tamaño mediano a grande:** Cuando se trabaja con conjuntos de datos que superan los cientos de miles o millones de filas, los beneficios de rendimiento y eficiencia de memoria de Polars se vuelven innegables.<sup>23</sup>
- **Se inician nuevos proyectos:** Comenzar un nuevo proyecto con Polars permite aprovechar desde el principio su API limpia y sus paradigmas de rendimiento, sin la carga de migrar un código base existente.
- **Las transformaciones son complejas:** Los pipelines de datos con múltiples pasos de filtrado, agregación y unión se benefician enormemente del optimizador de consultas del modo lazy.<sup>16</sup>

**Se podría considerar mantener pandas cuando:**

- **Existe un código base extenso y maduro:** La migración de un ecosistema grande y estable profundamente integrado con la API de pandas y sus librerías satélite puede suponer un costo de desarrollo significativo.<sup>23</sup>
- **Los conjuntos de datos son consistentemente pequeños:** Si los datos rara vez superan unas pocas decenas de miles de filas y el rendimiento no es una preocupación

principal, la familiaridad con pandas puede ser un factor decisivo.<sup>23</sup>

- **La dependencia de librerías de terceros es crítica:** Si el flujo de trabajo depende en gran medida de librerías que tienen una integración muy estrecha y específica con el objeto DataFrame de pandas y su índice, y no tienen alternativas compatibles con Polars.

### 7.3. La Trayectoria Futura de los DataFrames de Alto Rendimiento

Polars, junto con otras herramientas modernas como DuckDB, forma parte de una nueva generación de herramientas de análisis de datos construidas sobre las lecciones aprendidas de las limitaciones de sus predecesores.<sup>5</sup> Esta nueva ola se caracteriza por un enfoque implacable en el rendimiento aprovechando el hardware moderno, la adopción de estándares abiertos como Apache Arrow para la interoperabilidad, y la integración de técnicas de optimización de bases de datos en APIs más amigables para el usuario.

El futuro del ecosistema de datos parece dirigirse hacia un modelo en el que los datos pueden fluir entre las mejores herramientas de su clase (para consultas, manipulación, machine learning, visualización) con una fricción mínima y un rendimiento máximo. En este futuro interoperable y de alto rendimiento, Polars está firmemente posicionado no solo como una alternativa, sino como un nuevo estándar para el procesamiento de DataFrames a escala en una sola máquina.

#### Fuentes citadas

1. Polars — DataFrames for the new era, acceso: septiembre 15, 2025, <https://pola.rs/>
2. Polars user guide: Index, acceso: septiembre 15, 2025, <https://docs.pola.rs/>
3. pola-rs/polars: Extremely fast Query Engine for DataFrames, written in Rust - GitHub, acceso: septiembre 15, 2025, <https://github.com/pola-rs/polars>
4. Coming from Pandas - Polars user guide, acceso: septiembre 15, 2025, <https://docs.pola.rs/user-guide/migration/pandas/>
5. Polars vs. pandas: What's the Difference? | The PyCharm Blog, acceso: septiembre 15, 2025, <https://blog.jetbrains.com/pycharm/2024/07/polars-vs-pandas/>
6. Polars: Pandas DataFrame but Much Faster | by Travis Tang | TDS Archive | Medium, acceso: septiembre 15, 2025, <https://medium.com/data-science/pandas-dataframe-but-much-faster-f475d6be4cd4>
7. Introduction to Polars - Practical Business Python -, acceso: septiembre 15, 2025, <https://pbpython.com/polars-intro.html>
8. Polars vs Pandas : r/Python - Reddit, acceso: septiembre 15, 2025, [https://www.reddit.com/r/Python/comments/1jg402b/polars\\_vs\\_pandas/](https://www.reddit.com/r/Python/comments/1jg402b/polars_vs_pandas/)
9. Modern Polars: A comparison of the Polars and Pandas dataframe libraries | Hacker News, acceso: septiembre 15, 2025,

- <https://news.ycombinator.com/item?id=34275818>
10. Python Polars: High-Performance DataFrame Library in Rust - Analytics Vidhya, acceso: septiembre 15, 2025, <https://www.analyticsvidhya.com/blog/2024/06/python-polars/>
  11. polars - Rust - Docs.rs, acceso: septiembre 15, 2025, <https://docs.rs/polars/latest/polars/>
  12. What is a Polars expression? | Rho Signal, acceso: septiembre 15, 2025, <https://www.rhosignal.com/posts/what-is-an-expression/>
  13. How Polars Can Help You Build Fast Dash Apps for Large Datasets - Plotly, acceso: septiembre 15, 2025, <https://plotly.com/blog/polars-to-build-fast-dash-apps-for-large-datasets/>
  14. Bristech: Polars DataFrame library built on Apache Arrow, acceso: septiembre 15, 2025, <https://pola.rs/posts/talk-bristech-bytesize-conference/>
  15. Optimizations - Hugging Face, acceso: septiembre 15, 2025, <https://huggingface.co/docs/hub/datasets-polars-optimizations>
  16. Cookbook Polars for R - 4 Lazy execution - GitHub Pages, acceso: septiembre 15, 2025, [https://ddotta.github.io/cookbook-rpolars/lazy\\_execution.html](https://ddotta.github.io/cookbook-rpolars/lazy_execution.html)
  17. When To Use Polars Eager API Instead Of Lazy API? - YouTube, acceso: septiembre 15, 2025, <https://www.youtube.com/shorts/o5sL6x5qbrY>
  18. Polars: Speed Up Data Processing 12x with Lazy Execution - CodeCut, acceso: septiembre 15, 2025, <https://codecut.ai/polars-speed-up-data-processing-12x-with-lazy-execution/>
  19. DataFrame — Polars documentation, acceso: septiembre 15, 2025, <https://docs.pola.rs/py-polars/html/reference/dataframe/index.html>
  20. Python API reference — Polars documentation, acceso: septiembre 15, 2025, <https://docs.pola.rs/api/python/stable/reference/>
  21. Expressions — Polars documentation, acceso: septiembre 15, 2025, <https://docs.pola.rs/api/python/version/0.18/reference/expressions/index.html>
  22. Polars vs. Pandas — An Independent Speed Comparison | Towards Data Science, acceso: septiembre 15, 2025, <https://towardsdatascience.com/polars-vs-pandas-an-independent-speed-comparison/>
  23. Pandas vs Polars: A Comprehensive Performance Comparison | by Sina Mirshahi | Medium, acceso: septiembre 15, 2025, <https://medium.com/@neurogenou/pandas-vs-polars-a-comprehensive-performance-comparison-31a9296e4cd4>
  24. Updated PDS-H benchmark results (May 2025) - Polars, acceso: septiembre 15, 2025, <https://pola.rs/posts/benchmarks/>
  25. The power of predicate pushdown - Polars, acceso: septiembre 15, 2025, <https://pola.rs/posts/predicate-pushdown-query-optimizer/>
  26. Lazy API - Polars user guide, acceso: septiembre 15, 2025, <https://docs.pola.rs/user-guide/concepts/lazy-api/>
  27. polars.LazyFrame.explain — Polars documentation, acceso: septiembre 15, 2025, <https://docs.pola.rs/api/python/dev/reference/lazyframe/api/polars.LazyFrame.explain.html>

28. Out-of-core Polars - Hussain Sultan, acceso: septiembre 15, 2025,  
<https://hussainsultan.com/posts/out-of-core-polars/>
29. Streaming - Polars user guide, acceso: septiembre 15, 2025,  
<https://docs.pola.rs/user-guide/concepts/streaming/>