

**IIC3432 - Tópicos Avanzados en Bases de Datos**  
**Extracción de Información en XML**

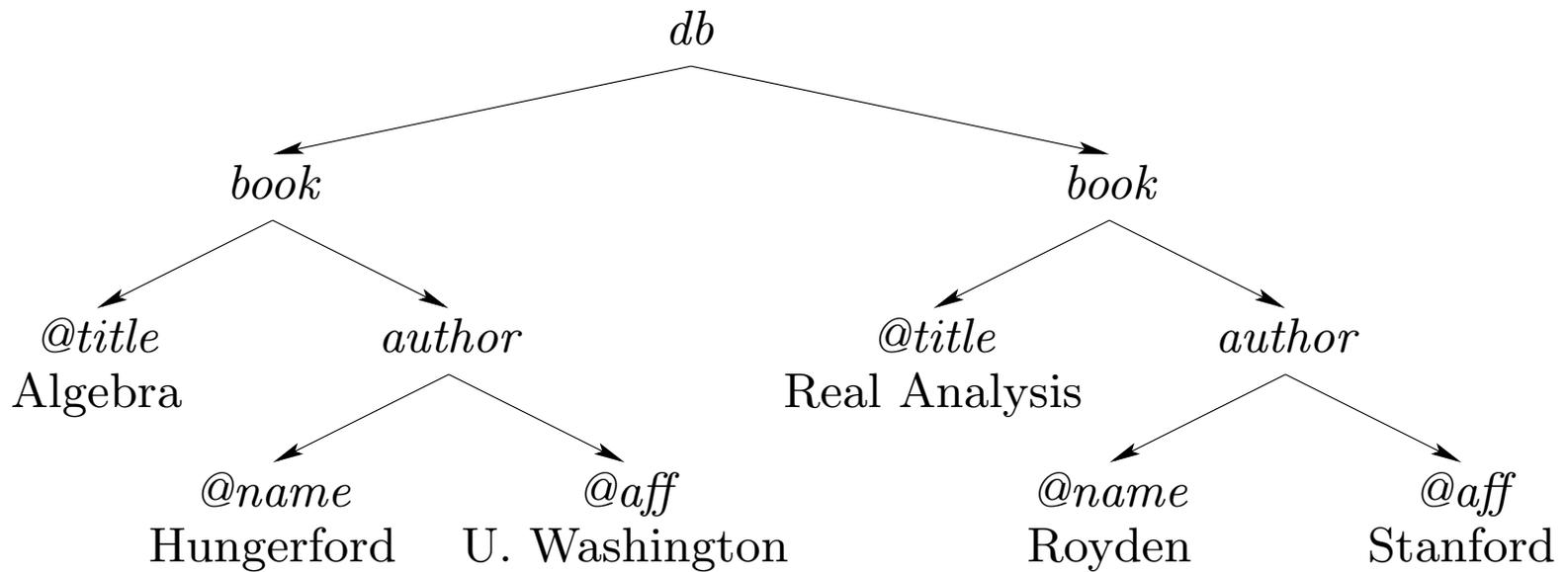
## Nuestro objetivo: Documentos XML

---

```
<db>
  <book title="Algebra">
    <author name="Hungerford" aff="U. Washington">
      </author>
    </book>
  <book title="Real Analysis">
    <author name="Royden" aff="Stanford">
      </author>
    </book>
</db>
```

## Documentos XML como árboles

---



## Lenguajes de consulta para XML: Operaciones básicas

---

**Filtrar:** Seleccionar valores desde un documento XML.

- Navegación, selección, extracción.

**Mezclar:** Integrar valores desde múltiples fuentes.

- Join, agregación.

**Transformar** valores desde un esquema a otro.

- Construcción de documentos XML.

# Lenguajes de consulta para XML: Algunos ejemplos

---

## **XPath:**

- Lenguaje más popular para navegar, seleccionar y extraer valores desde documentos XML.
- Parte de lenguajes más complejos como XQuery y XSLT.

## **XQuery:**

- Genera documentos XML como respuesta.
- Incluye join y agregación.

## **XSLT:**

- Lenguaje de patrones.
- Puede generar como respuesta documentos XML, HTML, texto u otros formatos.

# Lenguajes de consulta para XML: Algunos ejemplos

---

## **XPath:**

- Lenguaje más popular para navegar, seleccionar y extraer valores desde documentos XML.
- Parte de lenguajes más complejos como XQuery y XSLT.

## **XQuery:**

- Genera documentos XML como respuesta.
- Incluye join y agregación.

## **XSLT:**

- Lenguaje de patrones.
- Puede generar como respuesta documentos XML, HTML, texto u otros formatos.

## XML Path Language (XPath)

---

Estándar de la W3C: <http://www.w3.org/TR/xpath>.

Lenguaje para navegar, seleccionar nodos y extraer valores.

Algunas implementaciones:

XALAN : Apache Foundation (XSLT)

XT : James Clark (XSLT)

SAXON : Michael Kay (XSLT y XQuery)

## XML Path Language (XPath)

---

Estándar de la W3C: <http://www.w3.org/TR/xpath>.

Lenguaje para **navegar, seleccionar nodos** y extraer valores.

Algunas implementaciones:

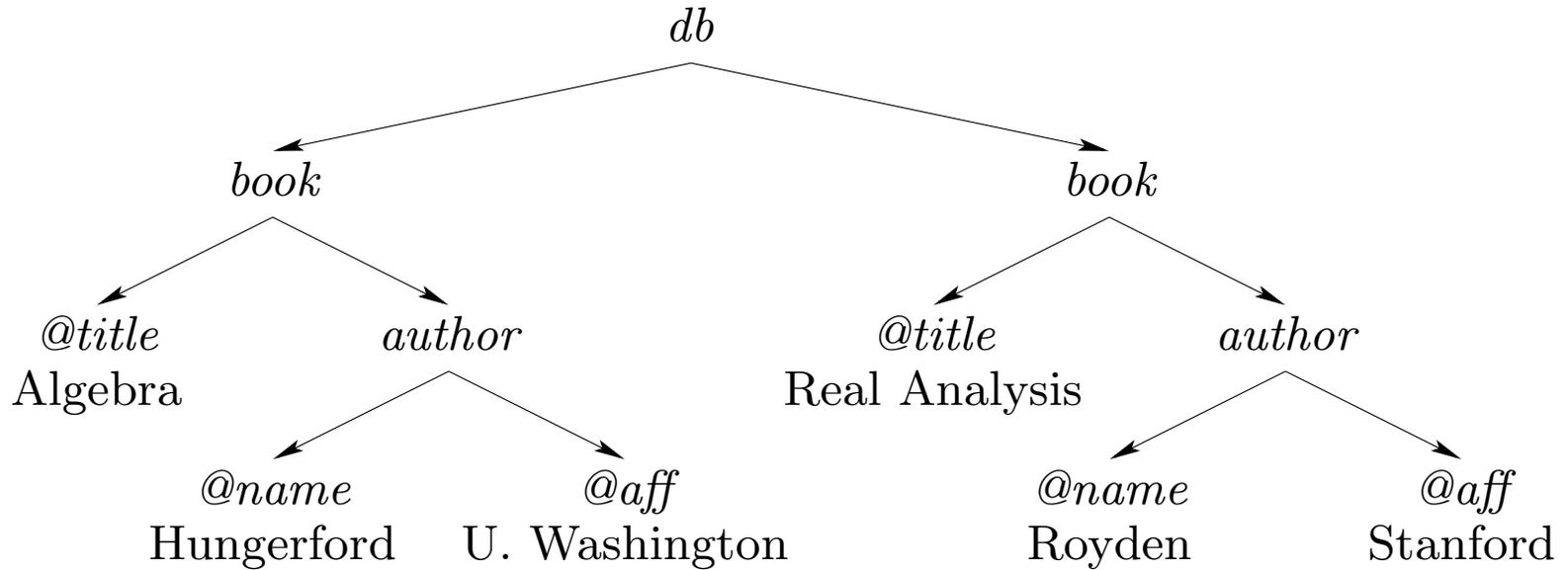
XALAN : Apache Foundation (XSLT)

XT : James Clark (XSLT)

SAXON : Michael Kay (XSLT y XQuery)

## Core XPath: Primer ejemplo

---

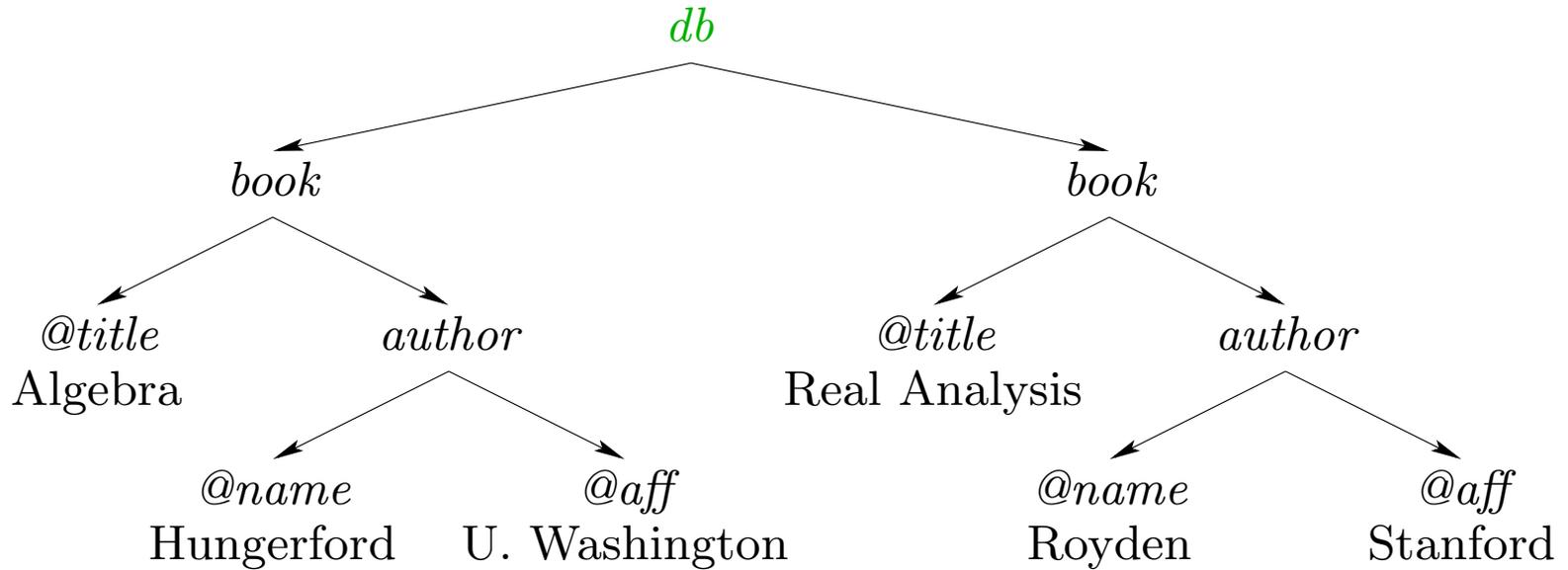


Consulta XPath: `child/?book`

Respuesta:

# Core XPath: Primer ejemplo

---

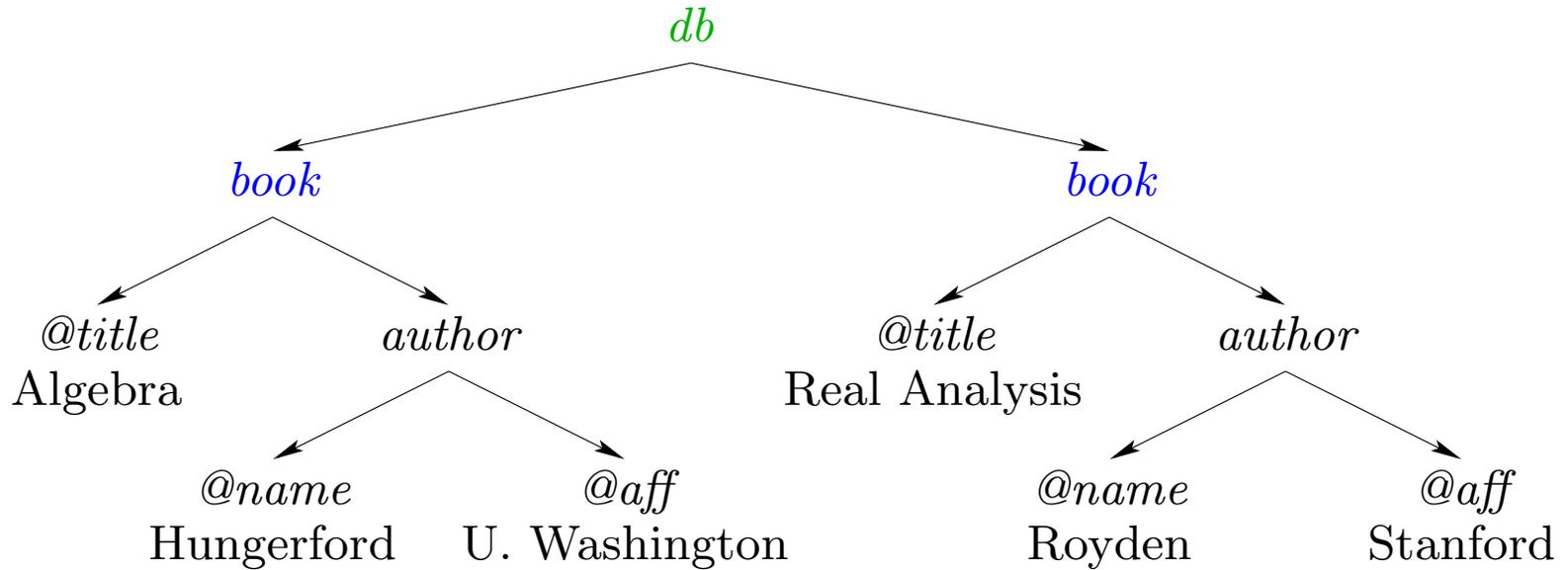


Consulta XPath: *child/?book*

Respuesta:

# Core XPath: Primer ejemplo

---

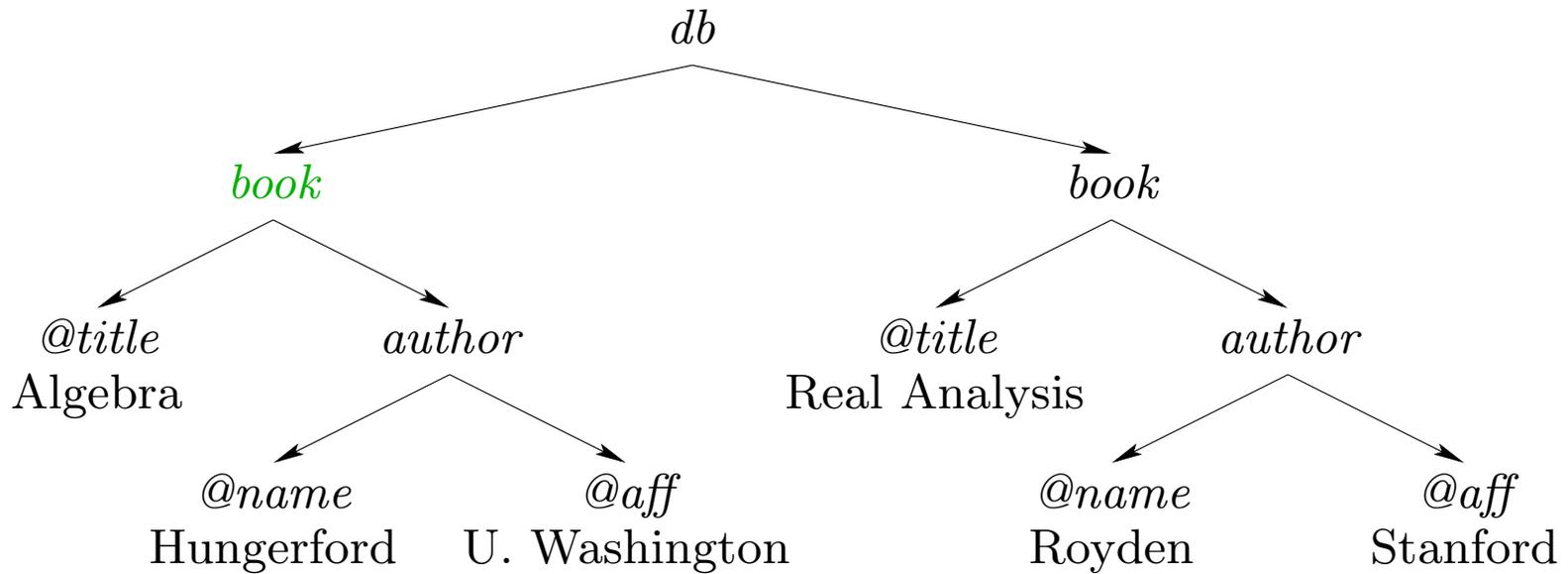


Consulta XPath: `child/?book`

Respuesta: `nodos azules`

## Core XPath: Segundo ejemplo

---

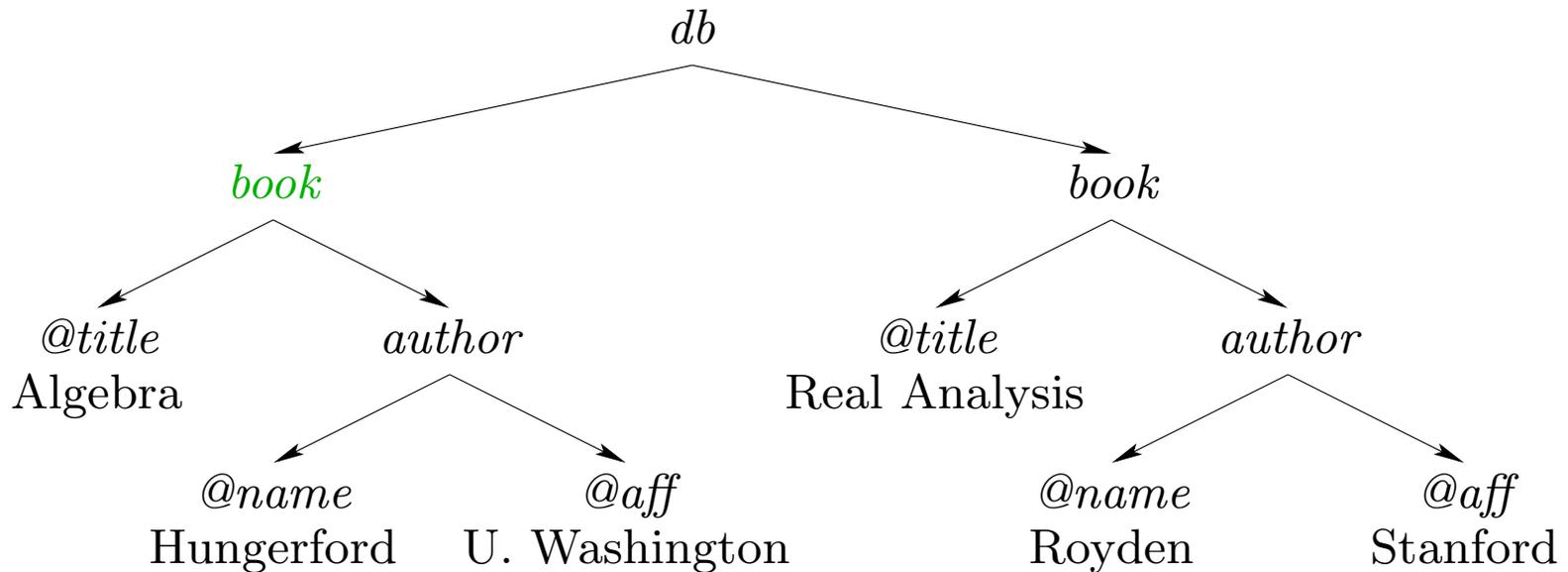


Consulta XPath: `child/?book`

Respuesta:

## Core XPath: Segundo ejemplo

---

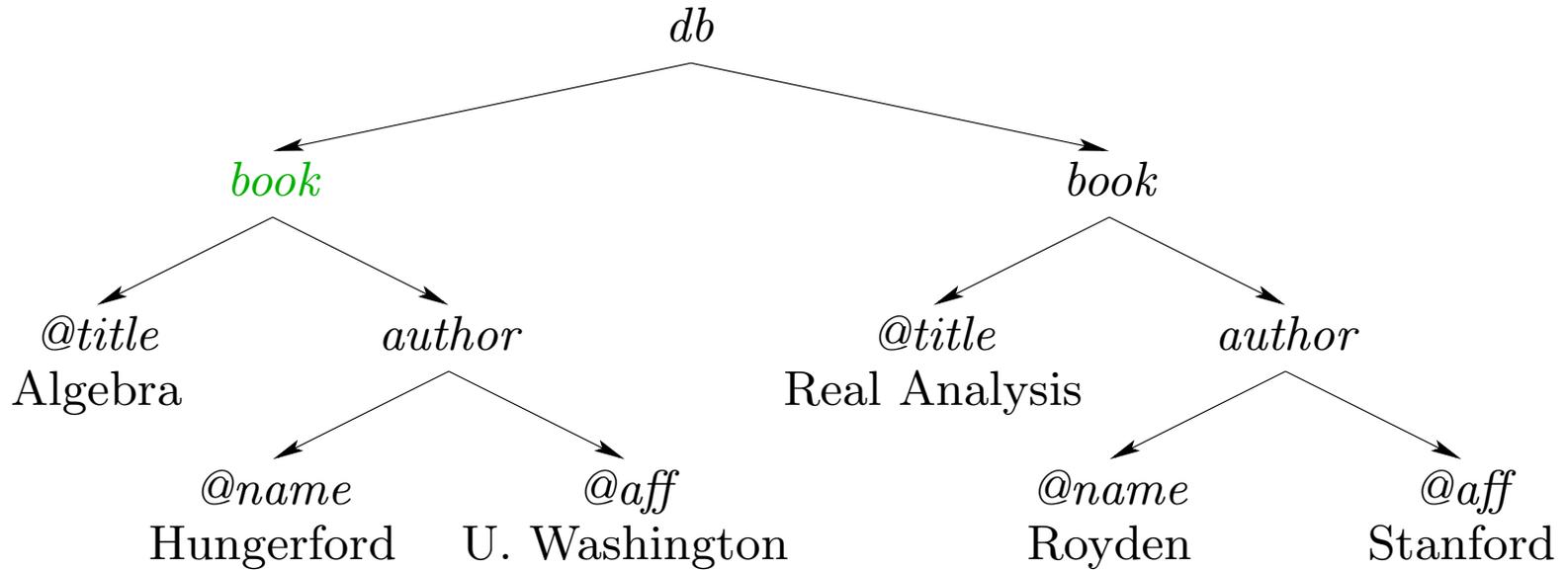


Consulta XPath: *child/?book*

Respuesta:  $\emptyset$

## Core XPath: Tercer ejemplo

---

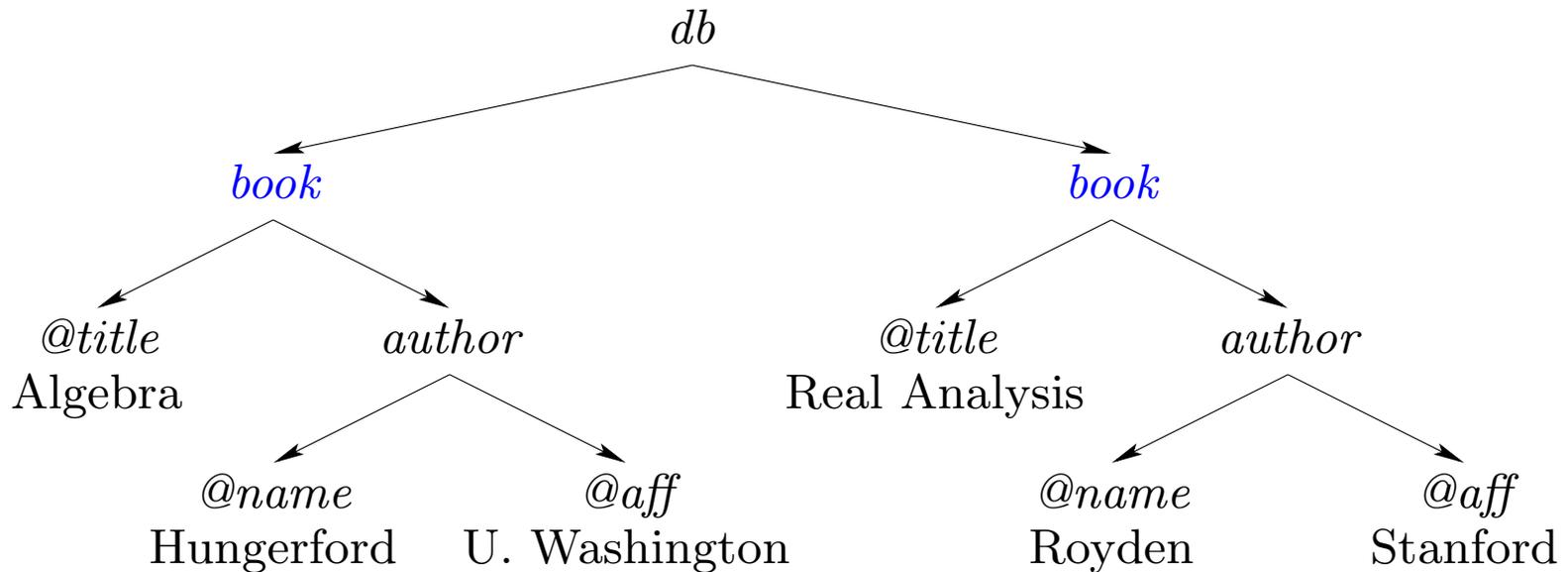


Consulta XPath: `parent/child/?book`

Respuesta:

## Core XPath: Tercer ejemplo

---

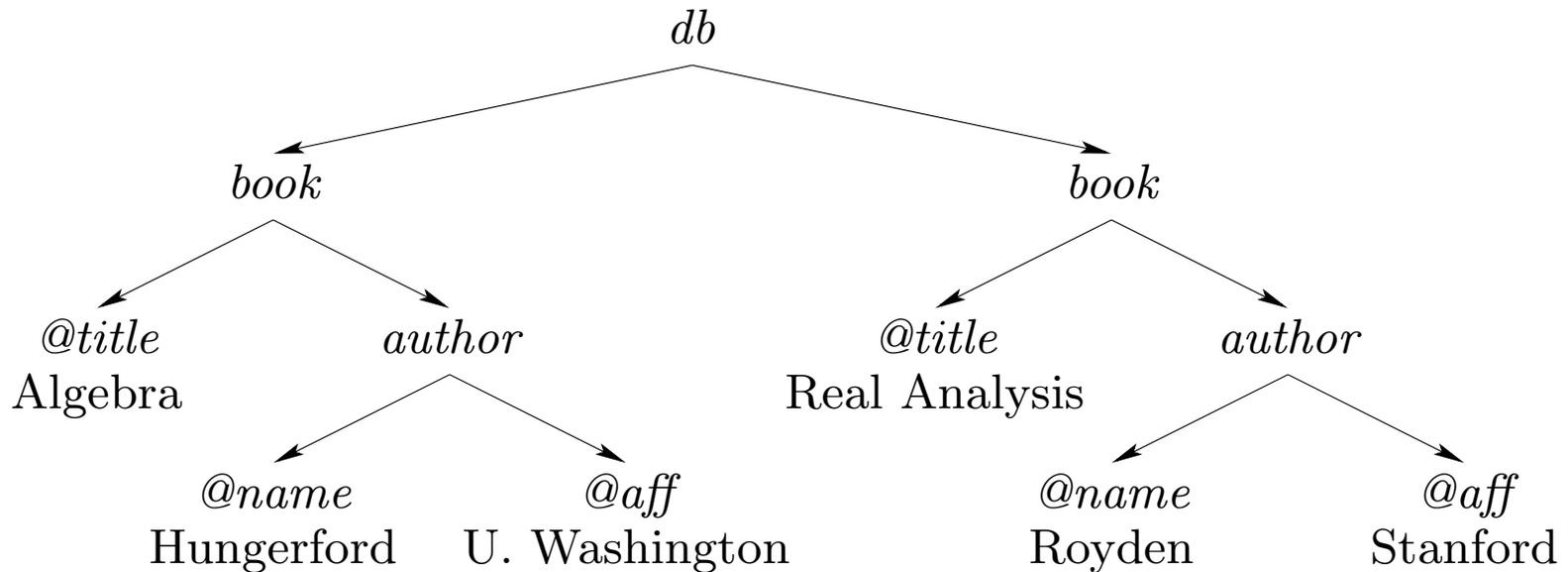


Consulta XPath: `parent/child/?book`

Respuesta: `nodos azules`

## Core XPath: Ultimo ejemplo

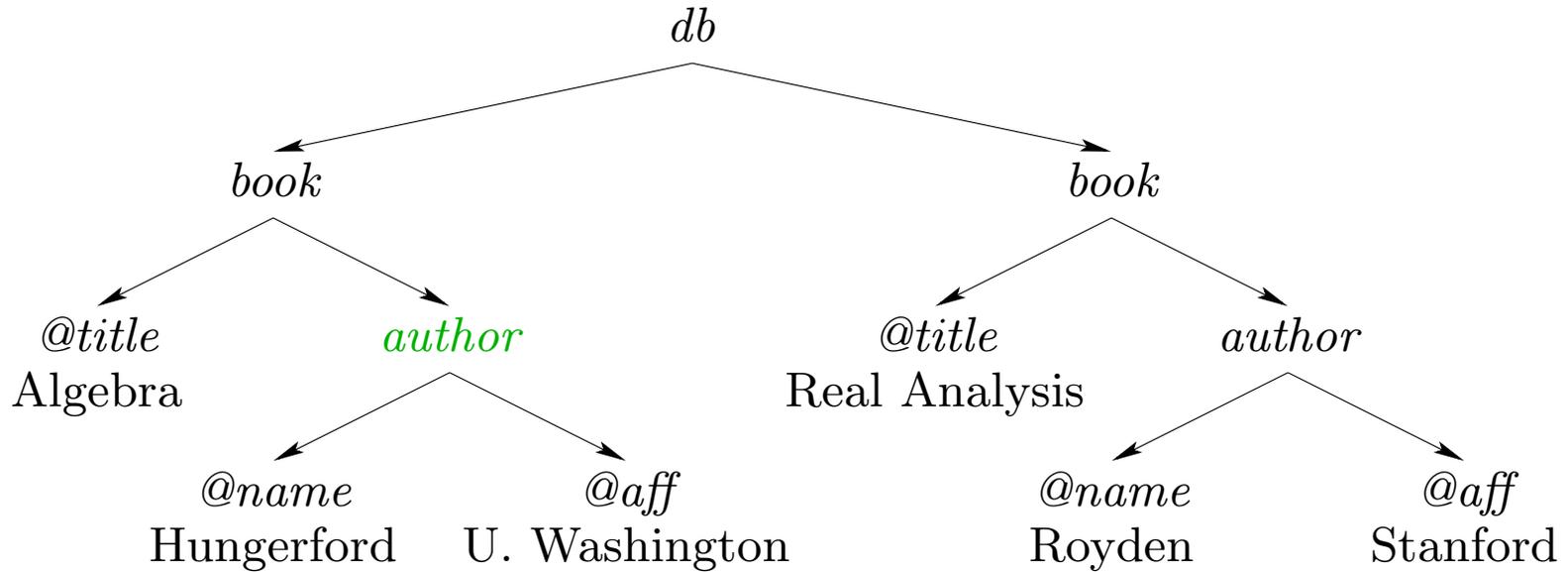
---



Consulta XPath: **parent\*/child\*/?book**

## Core XPath: Ultimo ejemplo

---

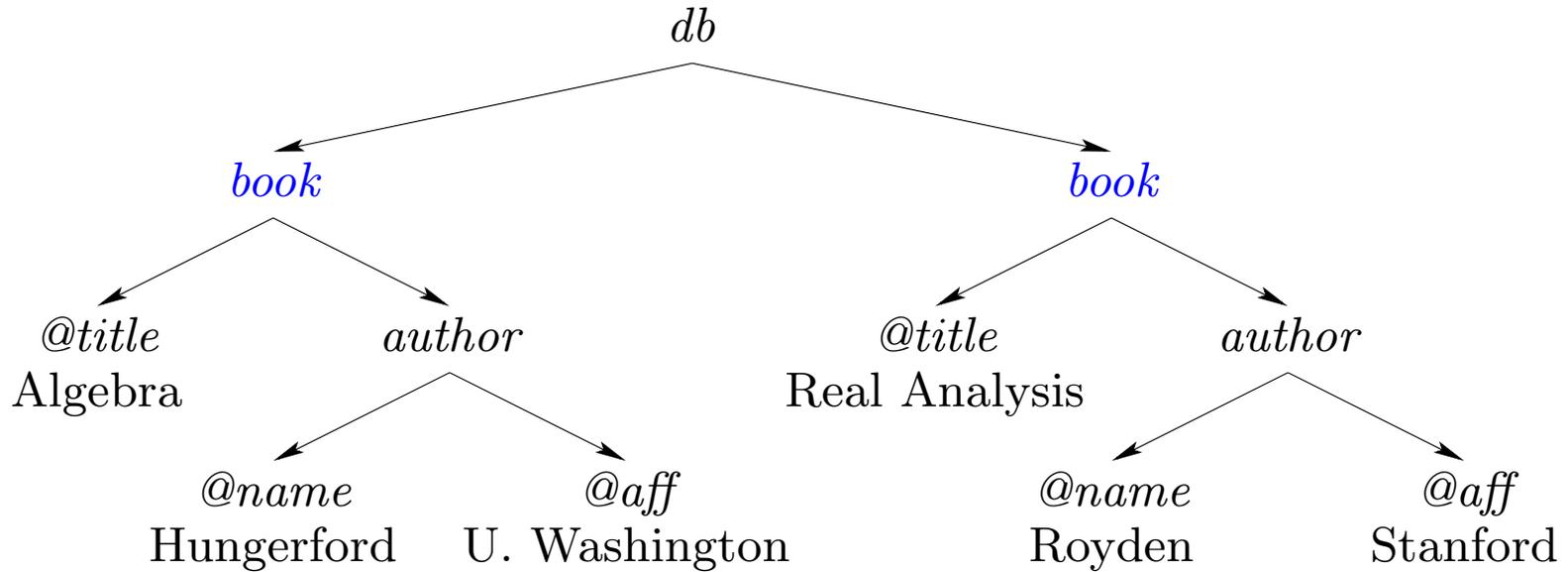


Consulta XPath: **parent\*/child\*/?book**

Respuesta:

## Core XPath: Ultimo ejemplo

---



Consulta XPath: **parent\*/child\*/?book**

Respuesta: **nodos azules**

## Core XPath: Sintaxis

---

Camino básicos:

```
paso ::= child | parent | right | left
```

Expresiones para caminos:

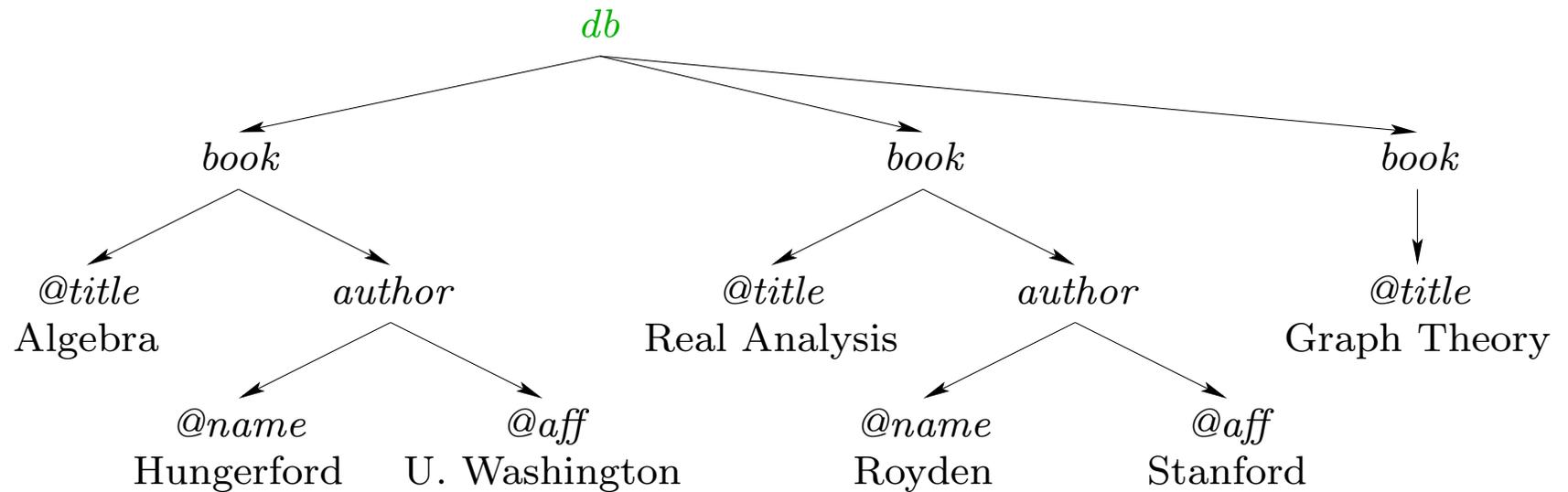
```
camino ::= paso | paso* | camino/camino |  
         camino ∪ camino | ?test
```

Filtros:

```
test ::= nombre | <camino> | ¬test | test ∧ test
```

## Core XPath: Otro ejemplo

---

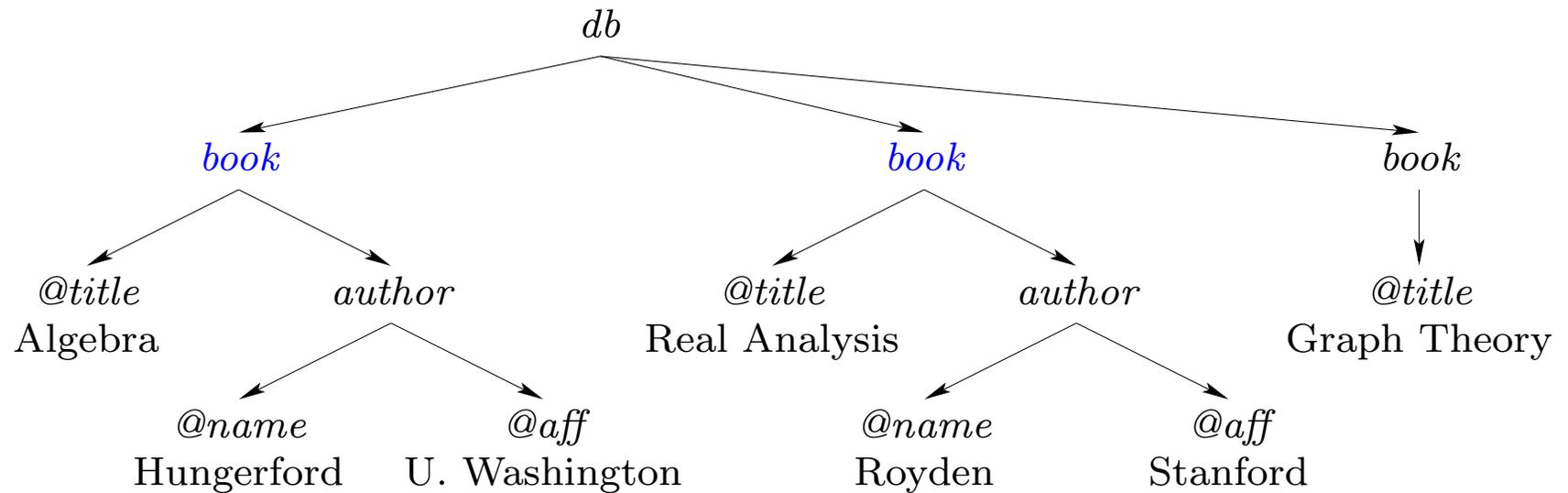


Consulta XPath: `child/?book/?<child/?author>`

Respuesta:

## Core XPath: Otro ejemplo

---



Consulta XPath: `child/?book/?<child/?author>`

Respuesta: **nodos azules**

## Core XPath: Semántica para caminos básicos

---

Dado: Árbol  $T = (D, \lambda)$ .

$$\llbracket \text{child} \rrbracket_T = \{(n_1, n_2) \mid n_2 \text{ es hijo de } n_1 \text{ en } T\}$$

$$\llbracket \text{parent} \rrbracket_T = \{(n_1, n_2) \mid n_2 \text{ es el padre de } n_1 \text{ en } T\}$$

$$\llbracket \text{right} \rrbracket_T = \{(n_1, n_2) \mid \text{existe nodo } n \text{ en } T \text{ tal que } n_1 \text{ es el } i\text{-ésimo hijo de } n \text{ y } n_2 \text{ es el } (i + 1)\text{-ésimo hijo de } n \text{ en } T\}$$

$$\llbracket \text{left} \rrbracket_T = \{(n_1, n_2) \mid \text{existe nodo } n \text{ en } T \text{ tal que } n_1 \text{ es el } (i + 1)\text{-ésimo hijo de } n \text{ y } n_2 \text{ es el } i\text{-ésimo hijo de } n \text{ en } T\}$$

## Core XPath: Semántica para caminos complejos

---

$$\begin{aligned} \llbracket \text{camino}_1 / \text{camino}_2 \rrbracket_T &= \{(n_1, n_2) \mid \text{existe } n_3 \text{ tal que} \\ &\quad (n_1, n_3) \in \llbracket \text{camino}_1 \rrbracket_T \text{ y } (n_3, n_2) \in \llbracket \text{camino}_2 \rrbracket_T\} \\ \llbracket \text{camino}_1 \cup \text{camino}_2 \rrbracket_T &= \llbracket \text{camino}_1 \rrbracket_T \cup \llbracket \text{camino}_2 \rrbracket_T \\ \llbracket \text{paso}^* \rrbracket_T &= \{(n, n) \mid n \in D\} \cup \llbracket \text{paso} \rrbracket_T \cup \llbracket \text{paso/paso} \rrbracket_T \cup \dots \\ &\quad \llbracket \text{paso/paso/paso} \rrbracket_T \cup \dots \\ \llbracket ?\text{test} \rrbracket_T &= \{(n, n) \mid n \in \llbracket \text{test} \rrbracket_T\} \end{aligned}$$

## Core XPath: Semántica para filtros

---

$$\llbracket \textit{nombre} \rrbracket_T = \{n \text{ en } T \mid \lambda(n) = \textit{nombre}\}$$

$$\llbracket \langle \textit{camino} \rangle \rrbracket_T = \{n \text{ en } T \mid \text{existe } n' \text{ en } T \text{ tal que } (n, n') \in \llbracket \textit{camino} \rrbracket_T\}$$

$$\llbracket \neg \textit{test} \rrbracket_T = \{n \text{ en } T \mid n \notin \llbracket \textit{test} \rrbracket_T\}$$

$$\llbracket \textit{test}_1 \wedge \textit{test}_2 \rrbracket_T = \llbracket \textit{test}_1 \rrbracket_T \cap \llbracket \textit{test}_2 \rrbracket_T$$

## ¿Cuándo es *bueno* un lenguaje de consulta?

---

Criterios esenciales:

- Expresividad.
- Complejidad.

¡Estos objetivos se contraponen!

- Vamos a ver que pasa en el caso de Core XPath.

## Core XPath: Complejidad

---

Enfoque ingenuo para evaluar una consulta: Usar un algoritmo recursivo que procesa secuencialmente la consulta.

- Usado en la mayoría de las implementación de XPath.

En XALAN y XT: Si  $Q = p_1/p_2/\dots/p_k$  ( $k \geq 1$ ), donde cada  $p_i$  es un paso o un test, entonces  $Q$  es procesada de la siguiente forma.

**procesar**( $Q$ : consulta,  $T$ : árbol,  $n$ : nodo)

$N := \{n' \mid (n, n') \in \llbracket p_1 \rrbracket_T\}$

**if**  $k = 1$  **then return**  $N$

**else**

$R := \emptyset$

**for each**  $n' \in N$  **do**  $R := R \cup$  **procesar**( $p_2/\dots/p_k, T, n'$ )

**return**  $R$

# Core XPath: Complejidad

---

¡Enfoque ingenuo es exponencial!

- Incluso si consideramos documentos de tamaño fijo.

Ejemplo [GKP05]: Considere el documento usado antes y la siguiente secuencia de consultas.

$Q_1 = ?db/child/?book$

$Q_2 = ?db/child/?book/parent/?db/child/?book$

$Q_3 = ?db/child/?book/parent/?db/child/?book/parent/?db/child/?book$

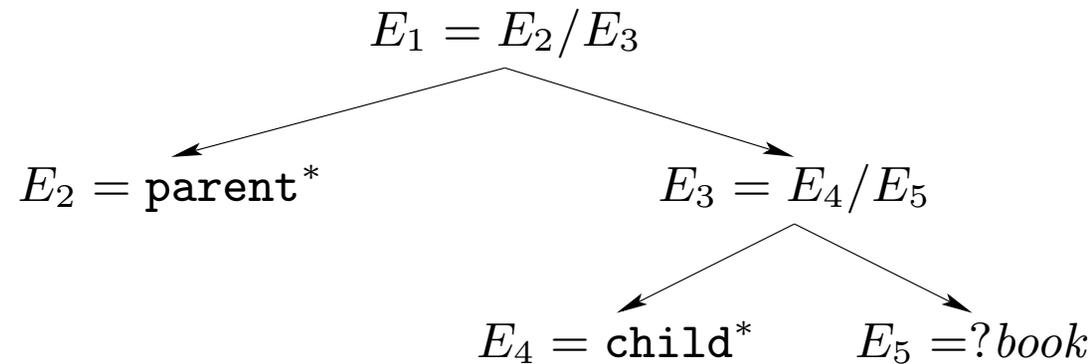
...

## Core XPath: Evaluación bottom-up

---

Algoritmo Bottom-up [GKP05]: *parent\*/child\*/?book*

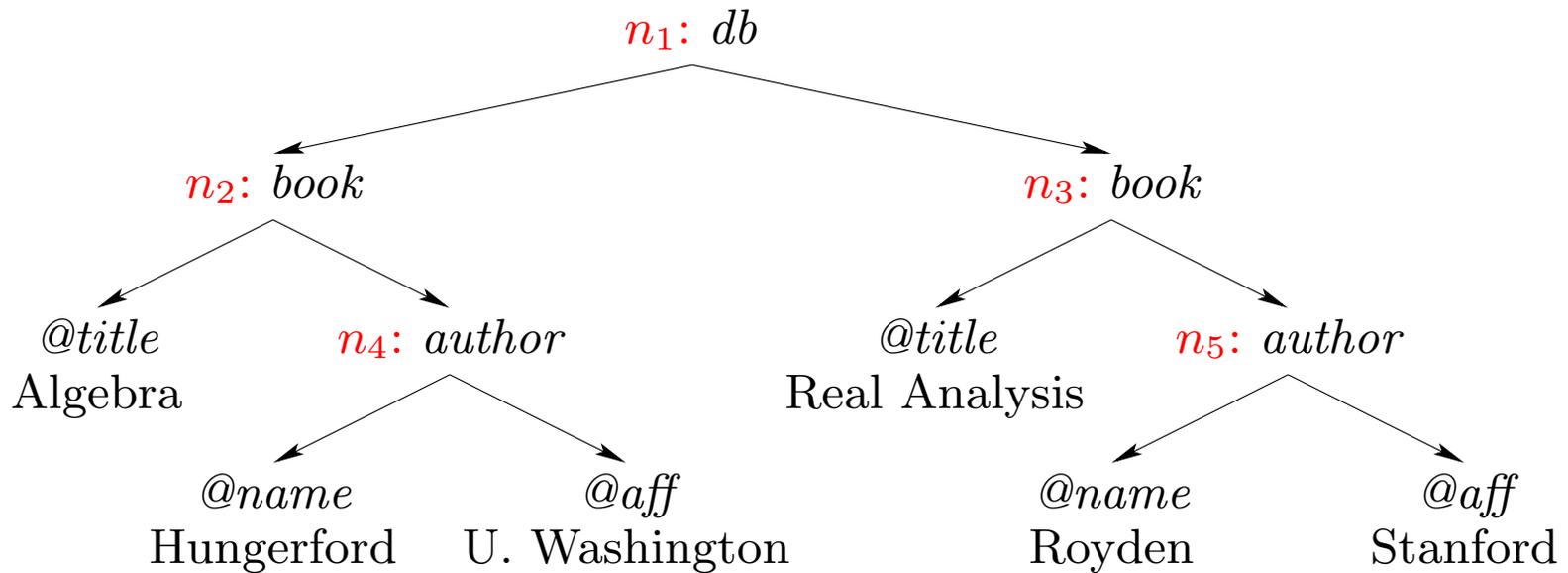
Primer paso: Construir el árbol de parsing de la consulta.



## Core XPath: Evaluación bottom-up

---

Segundo paso: Evaluar las sub-consultas de manera bottom-up.



# Core XPath: Evaluación bottom-up

---

$[[E_4]]_T$	
$n_1$	$n_1$
$n_2$	$n_2$
$n_3$	$n_3$
$n_4$	$n_4$
$n_5$	$n_5$
$n_1$	$n_2$
$n_1$	$n_3$
$n_1$	$n_4$
$n_1$	$n_5$
$n_2$	$n_4$
$n_3$	$n_5$

# Core XPath: Evaluación bottom-up

---

$\llbracket E_4 \rrbracket_T$		$\llbracket E_5 \rrbracket_T$	
$n_1$	$n_1$	$n_2$	$n_2$
$n_2$	$n_2$	$n_3$	$n_3$
$n_3$	$n_3$		
$n_4$	$n_4$		
$n_5$	$n_5$		
$n_1$	$n_2$		
$n_1$	$n_3$		
$n_1$	$n_4$		
$n_1$	$n_5$		
$n_2$	$n_4$		
$n_3$	$n_5$		

# Core XPath: Evaluación bottom-up

---

$[[E_4]]_T$		$[[E_5]]_T$		$[[E_3]]_T$	
$n_1$	$n_1$	$n_2$	$n_2$	$n_2$	$n_2$
$n_2$	$n_2$	$n_3$	$n_3$	$n_3$	$n_3$
$n_3$	$n_3$			$n_1$	$n_2$
$n_4$	$n_4$			$n_1$	$n_3$
$n_5$	$n_5$				
$n_1$	$n_2$				
$n_1$	$n_3$				
$n_1$	$n_4$				
$n_1$	$n_5$				
$n_2$	$n_4$				
$n_3$	$n_5$				

# Core XPath: Evaluación bottom-up

---

$[[E_2]]_T$	
$n_1$	$n_1$
$n_2$	$n_2$
$n_3$	$n_3$
$n_4$	$n_4$
$n_5$	$n_5$
$n_2$	$n_1$
$n_4$	$n_2$
$n_4$	$n_1$
$n_3$	$n_1$
$n_5$	$n_3$
$n_5$	$n_1$

# Core XPath: Evaluación bottom-up

---

$[[E_2]]_T$		$[[E_3]]_T$	
$n_1$	$n_1$	$n_2$	$n_2$
$n_2$	$n_2$	$n_3$	$n_3$
$n_3$	$n_3$	$n_1$	$n_2$
$n_4$	$n_4$	$n_1$	$n_3$
$n_5$	$n_5$		
$n_2$	$n_1$		
$n_4$	$n_2$		
$n_4$	$n_1$		
$n_3$	$n_1$		
$n_5$	$n_3$		
$n_5$	$n_1$		

# Core XPath: Evaluación bottom-up

---

$[[E_2]]_T$		$[[E_3]]_T$		$[[E_1]]_T$	
$n_1$	$n_1$	$n_2$	$n_2$	$n_1$	$n_2$
$n_2$	$n_2$	$n_3$	$n_3$	$n_1$	$n_3$
$n_3$	$n_3$	$n_1$	$n_2$	$n_2$	$n_2$
$n_4$	$n_4$	$n_1$	$n_3$	$n_2$	$n_3$
$n_5$	$n_5$			$n_3$	$n_2$
$n_2$	$n_1$			$n_3$	$n_3$
$n_4$	$n_2$			$n_4$	$n_2$
$n_4$	$n_1$			$n_4$	$n_3$
$n_3$	$n_1$			$n_5$	$n_2$
$n_5$	$n_3$			$n_5$	$n_3$
$n_5$	$n_1$				

## Core XPath: Evaluación bottom-up

---

**Teorema [GKP05]:** Una Consulta  $Q$  en Core XPath puede ser evaluada en tiempo  $O(|T| \cdot |Q|)$ .

Nota: Dado un nodo  $n$ , es posible calcular  $\{n' \mid (n, n') \in \llbracket Q \rrbracket_T\}$  en tiempo  $O(|T| \cdot |Q|)$ .

Sabemos que en términos de complejidad Core XPath es un buen lenguaje. ¿Es también bueno en términos de expresividad?

## ¿Qué tan expresivo es Core XPath?

---

La respuesta es relativa.

- Es *más, tan o menos* expresivo que ...

## ¿Con qué lenguaje podemos comparar?

- Un poco de historia: Bases de datos relacionales.

## Bases de datos relacionales: Lógica de primer orden

---

Base de datos relacional: Información es almacenada en relaciones (tablas).

Lenguaje de consulta natural: **Lógica de primer orden (FO)**.

- Ha sido estudiada por más de 100 años.
- Sintaxis y semántica bien definida.
- Expresividad bien entendida: que se puede y que no se puede decir.
- Complejidad bien entendida.

Pero existe un problema: Difícil de optimizar.

# Bases de datos relacionales: Algebra relacional

---

Un segundo lenguaje de consulta: **Algebra Relacional**.

- Combinación de operaciones algebraicas: selección, proyección, join, unión, diferencia, ...

Ventaja: Fácil de implementar y optimizar.

- Una de las razones para el éxito de las bases de datos relacionales.

Pero ¿cuál es la expresividad del álgebra relacional?

**Teorema [Cod72]:** Algebra relacional = FO.

## FO y XML

---

¿Podemos utilizar lógica de primer orden como un lenguaje de consulta para XML?

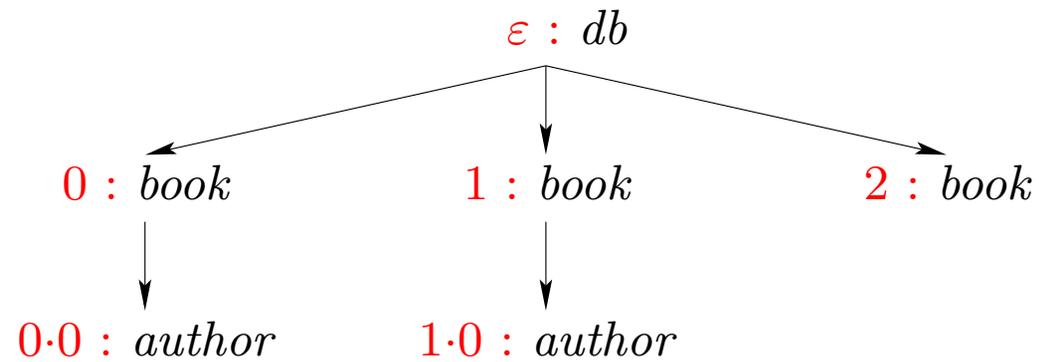
- ¿Es fácil expresar una consulta?

¿Cuál es la relación entre Core XPath y FO?

## FO sobre árboles

---

Dado árbol  $T$ :



Representamos  $T$  como una estructura relacional:

$$I_T \quad : \quad \{\text{child, desc, ns, sb, } P_{db}, P_{book}, P_{author}\}$$

## FO sobre árboles

---

Donde:

$$\begin{aligned} D &= \{\varepsilon, 0, 1, 2, 0\cdot 0, 1\cdot 0\} \\ \text{child} &= \{(\varepsilon, 0), (\varepsilon, 1), (\varepsilon, 2), (0, 0\cdot 0), (1, 1\cdot 0)\} \\ \text{desc} &= \{(\varepsilon, \varepsilon), (\varepsilon, 0), (\varepsilon, 1), (\varepsilon, 2), (\varepsilon, 0\cdot 0), (\varepsilon, 1\cdot 0), (0, 0), \\ &\quad (0, 0\cdot 0), (1, 1), (1, 1\cdot 0), (2, 2), (0\cdot 0, 0\cdot 0), (1\cdot 0, 1\cdot 0)\} \\ \text{ns} &= \{(0, 1), (1, 2)\} \\ \text{sb} &= \{(\varepsilon, \varepsilon), (0, 0), (0, 1), (0, 2), (1, 1), (1, 2), (2, 2), \\ &\quad (0\cdot 0, 0\cdot 0), (1\cdot 0, 1\cdot 0)\} \\ P_{db} &= \{\varepsilon\} \\ P_{book} &= \{0, 1, 2\} \\ P_{author} &= \{0\cdot 0, 1\cdot 0\} \end{aligned}$$

## FO sobre árboles: Ejemplos

---

Veamos como escribir las consultas que formulamos en Core XPath:

`child/?book:`

$$Q_1(x, y) = \text{child}(x, y) \wedge P_{\text{book}}(y).$$

`parent/child/?book:`

$$Q_2(x, y) = \exists z (\text{child}(z, x) \wedge \text{child}(z, y) \wedge P_{\text{book}}(y)).$$

`parent*/child*/?book:`

$$Q_3(x, y) = \exists z (\text{desc}(z, x) \wedge \text{desc}(z, y) \wedge P_{\text{book}}(y)).$$

`child/?book/?<child/?author>:`

$$Q_4(x, y) = \text{child}(x, y) \wedge P_{\text{book}}(y) \wedge \exists z (\text{child}(y, z) \wedge P_{\text{author}}(z)).$$

## Core XPath y FO

---

**Core XPath  $\subseteq$  FO:** Para cada consulta  $Q$  en Core XPath existe una fórmula  $\varphi(x, y)$  en FO tal que

$$(n, n') \in \llbracket Q \rrbracket_T \quad \text{si y sólo si} \quad I_T \models \varphi(n, n').$$

**Ejercicio:** Haga la demostración.

Pero: **FO  $\not\subseteq$  Core XPath.**

- No podemos expresar algunas consultas naturales.

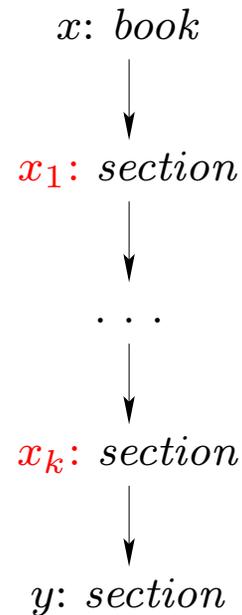
# FO no está contenido en Core XPath

---

Un ejemplo:

$$Q(x, y) = P_{book}(x) \wedge \mathbf{desc}(x, y) \wedge P_{section}(y) \wedge \\ \forall z (\mathbf{desc}(x, z) \wedge x \neq z \wedge \mathbf{desc}(z, y) \wedge z \neq y \rightarrow P_{section}(z)).$$

Buscamos:



## Core XPath

---

¿Qué tenemos que agregar a Core XPath para alcanzar la expresividad de FO?

¿Es una buena idea agregar nuevos elementos a Core XPath?

- Si lo que vamos a agregar es *natural, útil e implementable eficientemente*, entonces conviene agregarlo.

## Una pregunta natural: ¿Por qué no usar FO directamente?

---

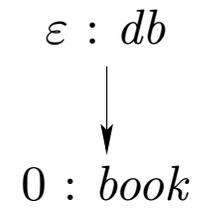
Tres desventajas:

- A los usuarios no les gusta trabajar con variables.
- Difícil de optimizar.
- Costoso evaluar una consulta.

# Complejidad de FO sobre árboles

---

Sea  $T$ :



# Complejidad de FO sobre árboles

---

Sea  $T$ :



Podemos usar  $\varepsilon$  como falso y  $0$  como verdadero, y así reducir desde SAT.

# Complejidad de FO sobre árboles

---

Sea  $T$ :



Podemos usar  $\varepsilon$  como falso y  $0$  como verdadero, y así reducir desde SAT.

Ejemplo:  $(x \vee \neg y \vee \neg z) \wedge (\neg x \vee y \vee \neg z)$  es representado como

$$\exists x \exists y \exists z ((P_{book}(x) \vee P_{db}(y) \vee P_{db}(z)) \wedge (P_{db}(x) \vee P_{book}(y) \vee P_{db}(z))).$$

# Complejidad de FO sobre árboles

---

Sea  $T$ :



Podemos usar  $\varepsilon$  como falso y  $0$  como verdadero, y así reducir desde SAT.

Ejemplo:  $(x \vee \neg y \vee \neg z) \wedge (\neg x \vee y \vee \neg z)$  es representado como

$$\exists x \forall y \exists z ((P_{book}(x) \vee P_{db}(y) \vee P_{db}(z)) \wedge (P_{db}(x) \vee P_{book}(y) \vee P_{db}(z))).$$

Conclusión: FO sobre árboles es NP-hard.

# Conditional XPath

---

Buena noticia: Hay que agregar poco a Core XPath para obtener FO y la complejidad no cambia.

## Conditional XPath:

- Caminos básicos:

`paso ::= child | parent | right | left`

- Expresiones para caminos:

`camino ::= paso | (paso/?test)* | camino/camino |  
camino  $\cup$  camino | ?test`

- Filtros:

`test ::= nombre | <camino> |  $\neg$ test | test  $\wedge$  test`

## Conditional XPath

---

Dado un árbol  $T = (D, \lambda)$ :

$$\begin{aligned} \llbracket (\text{paso}/?\text{test})^* \rrbracket_T &= \{(n, n) \mid n \in D\} \cup \llbracket \text{paso}/?\text{test} \rrbracket_T \cup \\ &\quad \llbracket \text{paso}/?\text{test}/\text{paso}/?\text{test} \rrbracket_T \cup \dots \end{aligned}$$

Nótese que Core XPath  $\subseteq$  Conditional XPath:

$$\text{paso}^* \text{ es equivalente a } (\text{paso}/?\neg(a \wedge \neg a))^*.$$

## Conditional XPath y FO

---

**Teorema [Mar04a, Mar05]:** Sobre árboles se tiene que  
Conditional XPath = FO.

**Teorema [Mar04b]:** Una Consulta  $Q$  en Conditional XPath  
puede ser evaluada en tiempo  $O(|T| \cdot |Q|)$ .

¡Conditional XPath es un buen lenguaje!

## Otra pregunta natural

---

Si la complejidad de Conditional XPath es  $O(|T| \cdot |Q|)$  y de FO es NP-hard, ¿Cómo puede ser que Conditional XPath = FO?

Del teorema anterior se puede concluir que **no hay una traducción eficiente de FO a Conditional XPath.**

Del teorema anterior **no** se puede inferir qué lenguaje es mejor para el usuario.

## ¿Necesitamos más lenguajes de consulta?

---

Conditional XPath es un buen lenguaje, pero ¿es suficiente?

Para navegar XML es necesario utilizar **expresiones regulares arbitrarias**.

- La mayoría de los lenguajes propuestos las incluyen.

Veamos un lenguaje con expresiones regulares.

## Regular XPath: Sintaxis

---

Camino básicos:

`paso ::= child | parent | right | left`

Expresiones para caminos:

`camino ::= paso | camino* | camino/camino |  
camino  $\cup$  camino | ?test`

Filtros:

`test ::= nombre | <camino> |  $\neg$ test | test  $\wedge$  test`

## Regular XPath: Semántica

---

Dado un árbol  $T = (D, \lambda)$ :

$$\begin{aligned} \llbracket \text{camino}^* \rrbracket_T &= \{(n, n) \mid n \in D\} \cup \llbracket \text{camino} \rrbracket_T \cup \\ &\quad \llbracket \text{camino/camino} \rrbracket_T \cup \llbracket \text{camino/camino/camino} \rrbracket_T \cup \dots \end{aligned}$$

Ejemplo:  $(\text{child}/?section/\text{child}/?section)^*$

## La expresividad de Regular XPath

---

Nótese que  $\text{Conditional XPath} \subseteq \text{Regular XPath}$ :

$(\text{paso}/?\text{test})^*$  es una fórmula en Regular XPath.

¿Qué tan expresivo es Regular XPath?

- ¿Es  $\text{Conditional XPath} = \text{Regular XPath}$ ?
- ¿Con qué otro lenguaje lo podemos comparar?

Otro poco de historia: **Autómata y lógica.**