

# Einführung in XPath

Dr. Björn Rudzewitz

pagina GmbH

18.01.2023





Publishing-Softwarelösungen

Digitale Transformation

Digital Humanities

Typografie und Templating

Digitales Lernen

Besuchen Sie unsere Webseite: <https://pagina.gmbh/>

## Lernziele:

- ✓ einsteigerfreundliche Einführung in Grundlagen
- ✓ Anwendung und Zweck von XPath
- ✓ XPath und XML-Baumstruktur
- ✓ Struktur von XPath-Ausdrücken verstehen
- ✓ Ausdrücke selber schreiben

## Nicht Teil der Lernziele:

- X komplexe, fortgeschrittene Ausdrücke mit Funktionen
- X Unterschiede zwischen Versionen von XPath
- X praktische Integration von XPath in Programme und Programmiersprachen
- X vollständige Beschreibung des Featureumfangs

# Was ist XPath?

- Auswahlsprache für XML-Dokumente
- XML Inhalte Anwendungen Menschen und Anwendungen zugänglich machen
- XPath als Basis vieler fortgeschrittener XML-Technologien (z.B. XML Schema, XSLT, XQuery, ...)
- viele Programmiersprachen unterstützen XPath (Java, Python, XSLT, ...)
  - z.B. Selektion und Verarbeitung von XML-Dokumenten nach Einlesen in Java

- XML-Dokumente haben eine Baumstruktur ( "*document tree*" )
- XPath sucht Elemente im Baum relativ zu anderen Elementen
- z.T. Ähnlichkeiten zu UNIX-Dateisystembaum

- XPath arbeitet mit Knoten (“*nodes*”) und nicht unmittelbar mit XML-Elementen, Attributen, ...
- XPath sieht jede XML-Komponente als Knoten: Elemente, Attribute, Text, Dokument, ...
- “*node set*”: Menge an Knoten

Knotenarten:

- **Wurzelknoten:** beinhaltet *alles*, inkl. *Wurzelement*,  
Verarbeitungsanweisungen, Kommentare, ...
- **Elementknoten:** entspricht einem XML-Element
- **Attributknoten:** hängt an Elementknoten und repräsentiert  
XML-Attribut
- **Textknoten:** stetiger (nicht unterbrochener) textueller Inhalt  
(#PCDATA)

Knotenarten:

- **PI-Knoten:** Verarbeitungsanweisungsknoten
- **Kommentarknoten:** für XML-Kommentare
- **Namensraumknoten**

Knotenarten:

- **PI-Knoten:** Verarbeitungsanweisungsknoten
- **Kommentarknoten:** für XML-Kommentare
- **Namensraumknoten**

⇒ XPath sieht Dokumente als *Knoten*

# Beziehungen zwischen Knoten I

- aktueller Knoten, **Kontextknoten** ("*self*")
- **Eltern, Kind, Geschwister**
- **übergeordnet**: Elternknoten, Elternknoten des Elternknotens, ... ("*ancestor*")

# Beziehungen zwischen Knoten II

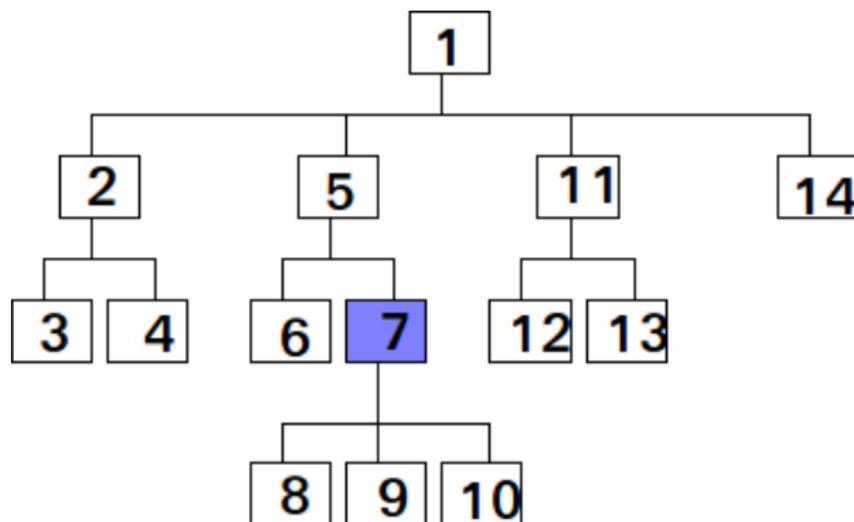
- **untergeordnet**: Kind, Kindeskind, ... (“*descendant*”)
- **Geschwister**: gleicher Elternknoten (“*sibling*”)

## Beziehungen zwischen Knoten III

- **vorangehend**: Knoten, die vor dem aktuell ausgewählten Knoten durchschritten werden ohne Vorfahren (“*preceding*”)
- **nachfolgend**: Knoten durchschritten nach dem aktuellen Knoten ohne Nachfahren ... (“*following*”)
- Kombinationen: untergeordnet-oder-selbst, ... (“*descendant-or-self*”)

guter Überblick: <https://wiki.selfhtml.org/wiki/XML/XSL/XPath>

# Beziehungen zwischen Knoten



Welche Knoten sind Eltern, übergeordnete, untergeordnete, vorangehende Knoten, nachfolgende Geschwisterknoten des aktuellen (blauen) Kontextknoten ?

Quelle: [https://wiki.selfhtml.org/wiki/Datei:Xml\\_parent.gif](https://wiki.selfhtml.org/wiki/Datei:Xml_parent.gif), zuletzt aufgerufen am 17.01.2023

- XPath basiert auf dem Konzept (virtueller) **Achsen**
- Knoten liegen auf den virtuellen Achsen, mit unterschiedlichen Knotentypen auf unterschiedlichen Achsen
- Beispiel: ein Kindelementknoten liegt auf der Elementachse unter dem Elternknotenelement, ein Attribut liegt auf einer dazu geneigten Achse
- [Simpson, 2002, page 16]:  
*an axis “turns the view” from a given starting point in the document*

Allgemeines Muster von XPath-Ausdrücken/Pfaden:

**Achse**::**Knotentest**[**Prädikat**]

- **Achse**: in welcher Richtung/Dimension
- **Knotentest**: Eigenschaften eines Knotens (Name, Typ, ...)
- **Prädikat**: Verfeinerung der Suchergebnisse

Allgemeines Muster von XPath-Ausdrücken/Pfaden:

**Achse**::**Knotentest**[**Prädikat**]

- **Achse**: in welcher Richtung/Dimension
  - **Knotentest**: Eigenschaften eines Knotens (Name, Typ, ...)
  - **Prädikat**: Verfeinerung der Suchergebnisse
- 
- meist werden mehrere solcher Ausdrücke aneinandergereiht
  - nur der Knotentest ist explizit anzugeben

## Knotentest:

- Auswahl Elemente nach Name: `/course/description`
- Auswahl aller Kindelemente: `/courses/*`
- Auswahl aller Textknoten: `/courses/description/text()`  
analog dazu: `comment()`, `processing-instruction()`
- Auswahl aller Knoten aller Art: `/courses/node()`

## Übung

Aufgabe: Schreiben Sie einen XPath-Ausdruck, um die **Titel-Elemente aller Bücher** zu bekommen.

- Beispieldaten: <https://www.w3schools.com/xml/books.xml>.
- Online-Editor zum Ausprobieren: <http://xpather.com/>

### Achse:

- **gegeben einen Kontextknoten durch den Knotentest**
- gibt eine Achse an, in welche Richtung/Dimension gesucht wird
- 13 Achsen
- spezielle Achsen nur von bestimmten Knoten erreichbar (z.B. Attributachsen nur bei Elementachsen)

```
/campus/person/attribute::node()
```

# XPath-Selektoren

## Achse

Achse	wählt aus (gegeben Kontextknoten)
child::	Kindelementknoten
parent::	Elternelementknoten
descendant::	Nachfahrenelementknoten
ancestor::	Vorfahrenelementknoten ...

# XPath-Selektoren

## Achse

Achse	wählt aus (gegeben Kontextknoten)
self::	aktuellen Kontextknoten
attribute::	Attributknoten
preceding::	vorangehende Elemente
preceding-sibling::	vorangehende Schwesterelemente

following::, following-sibling::, namespace::

# Konventionen bei Achsen

## Konventionen/Abkürzungen:

- Kindachse ist Standard:  
`child::course` äquivalent zu `course`
- Elternknoten: `parent::node()` äquivalent zu `..`
- aktueller Kontextknoten: `self::node()` äquivalent zu `.`
- `/descendant-or-self::node( )` äquivalent zu `//`
- `/attribute::` äquivalent zu `@`

## Übung

### Aufgaben:

- 1 Schreiben Sie einen XPath-Ausdruck, um **alle Attribut-Wertpaare aller Buchelemente** zu extrahieren.
- 2 Schreiben Sie einen XPath-Ausdruck, um alle **Autoren** eines Buches zu extrahieren, die **nicht an erster Stelle** stehen (d.h. Autoren, die anderen Autoren folgen):
  - Beispieldaten: <https://www.w3schools.com/xml/books.xml>.
  - Online-Editor zum Ausprobieren: <http://xpather.com/>

# Prädikate

- Prädikate spezifizieren zusätzliche Bedingungen für das Matching von Knoten
- verfeinern die Suchergebnisse mithilfe zusätzlicher XPath-Ausdrücke
- spezifiziert in eckigen Klammern [ ]
- Muster:  
`expr1 operator expr2`
- `expr1`, `expr2` sind XPath-Ausdrücke
- `operator`: =, !=, >, <, >=, <=
- `operator expr2` sind optional

- Prädikate ergeben Boolesche Ausdrücke: wahr oder falsch
- es genügt, wenn der Wert einmal wahr ist
  - z.B. wenn eine Bedingung Kindelemente berücksichtigt, reicht es, wenn die Bedingung für eines der Kindelemente wahr ist, muss nicht für alle gelten

Beispiel aus [Simpson, 2002, Ch. 3.3]:

```
<person name="John">  
  <child name="John"/>  
  <child name="Connie"/>  
  <child name="Cindy"/>  
  <child name="Mike"/>  
</person>
```

Beispiel aus [Simpson, 2002, Ch. 3.3]:

```
<person name="John">  
  <child name="John"/>  
  <child name="Connie"/>  
  <child name="Cindy"/>  
  <child name="Mike"/>  
</person>
```

Der XPath-Ausdruck `/person[child/@name='Cindy']` gibt den Wurzelementknoten zurück, da das Prädikat für eines der Kindelemente zutrifft

in der Praxis werden Ausdrücke mit Prädikaten oft verkettet:

```
/root/courses[@faculty='sfs']//information[@participants >10]
```

wählt alle “information” Elemente aus, bei denen einerseits das Attribut “participants” den Wert größer als 10 hat und die andererseits Nachfahren von “course” Elementen sind mit dem Attribut “faculty” mit dem Wert “sfs”, welche wiederum Kindelemente sind des Wurzelements mit dem Elementnamen “root”

# Compound Predicates

```
<participants>  
  <person name="John" age="20" />  
  <person name="John" age="23"/>  
  <person name="Paul" age="20"/>  
  <person name="Emma" age="20"/>  
</participants>
```

# Compound Predicates

```
<participants>  
  <person name="John" age="20" />  
  <person name="John" age="23"/>  
  <person name="Paul" age="20"/>  
  <person name="Emma" age="20"/>  
</participants>
```

- zusammengesetzte Prädikate mit **and** oder **or**
- runde Klammern für Gruppierung

```
/participants/person[@name="John"and @age="20"]
```

```
//person[(@name="John"or @name="Paul") and @age="20"]
```

# Reduced Predicates

- Prädikate können Bedingungen unabhängig vom eigentlichen Wert festlegen
- Beispiel: `//metadata[descendant::license]`  
wählt diejenigen *metadata* Elemente aus, die ein *license* Element als Nachfahre haben
- abgekürzte Positionsangabe (nur wenn nicht in zusammengesetztem Ausdruck):  
`//person[4]` äquivalent zu `//person[position() = 4]`  
wählt das vierte *person* Element aus
- weitere Positionsangaben: `first()`, `last()`

## Übung

### Aufgaben:

- 1 Schreiben Sie einen XPath-Ausdruck, um die **Titelemente** aller Bücher der **Kategorie web** zu extrahieren.
- Beispieldaten: <https://www.w3schools.com/xml/books.xml>.
- Online-Editor zum Ausprobieren: <http://xpather.com/>

- XPath erlaubt die Benutzung von Funktionen
- Beispiele:
  - arithmetische Operationen, z.B. `count()`, `sum()`
  - Textfunktionen, z.B. `text()`, `contains()`, `matches()`
- XPath-Funktionen haben immer einen Rückgabewert (nicht void)

## XPath Datentypen:

- **string**: Text
- **nodeset**: XPath Ausdruck/Verweis auf Knoten, kann leer sein
- **boolean**: wahr/falsch
- **number**: Zahl
- **anytype**: beliebig/unbestimmt

siehe [Simpson, 2002, Kapitel 4.2]

Einige wichtige XPath-Funktionen auf Knotensets:

Funktion	Rückgabewert	Beispiel
position()	Position in Knotenset (1-basiert)	<code>/participants/person</code> <code>[position() = 2]</code>
last()	Anzahl Knoten	<code>/participants/person[last()]</code>
count(nodeset)	Anzahl Knoten	<code>count(//person)</code>

Einige wichtige XPath-Funktionen für Text:

Funktion	Rückgabewert	Beispiel
<code>string(anytype)</code>	Textrepräsentation (nur 1. Knoten für Sets!)	<code>string(count(//person))</code>
<code>concat(a,b,...)</code>	Textverkettung	<code>concat(count(//person), "people")</code>
<code>contains(a,b)</code>	beinhaltet a b	<code>//person[contains(@name,"j")]</code>

Einige wichtige XPath-Funktionen für Text:

Funktion	Rückgabewert	Beispiel
<code>substring(a,1,2)</code>	Teilzeichenkette von a von 1 bis 2	<code>substring(//person[2]/@name, 1, 3)</code>

indexes start at 1 for XPath

Einige wichtige XPath-Funktionen für Zahlen:

Funktion	Rückgabewert	Beispiel
<code>sum(nodeset)</code>	Summe	<code>sum(//person/@age)</code>
<code>floor(n)</code>	Abrundungsfunktion	<code>floor(1.3 * 4)</code>
<code>ceiling(n)</code>	Aufrundungsfunktion	<code>ceiling(1.3 * 4)</code>
<code>round(n)</code>	Rundungsfunktion	<code>round(1.4*3)</code>

Hinweis: `div` zum Teilen anstelle von `/`,  
`+`, `-`, `*` werden direkt verwendet

## Übung

### Aufgaben:

- 1 Schreiben Sie einen XPath-Ausdruck, um **alle Buchtitel** zu extrahieren von Büchern, die mehr als 35 kosten.
- 2 Schreiben Sie einen XPath-Ausdruck, um die **Namen aller Autoren** zu extrahieren, die den **Großbuchstaben "K"** enthalten.
- 3 Schreiben Sie einen XPath-Ausdruck, der den **Durchschnittspreis** über alle Bücher berechnet.
  - Beispieldaten: <https://www.w3schools.com/xml/books.xml>.
  - Online-Editor zum Ausprobieren: <http://xpather.com/>

# Quellen und weiterführende Materialien

*Simpson, John. XPath and XPointer: Locating Content in XML Documents. O'Reilly Media, Inc., 2002.*

Die Struktur der Folien folgt lose diesem Buch und integriert Teile der Inhalte der Folien aus dem Text Technology Kurs von Prof. Chris Culy aus dem Jahr 2013. Die vorliegenden Folien wurden von Dr. Björn Rudzewitz erstellt und in mehreren Iterationen im SS 2017/18/19/20/21 an der Universität Tübingen zur Lehre verwendet.

Begriffe für deutsche Übersetzung:

<https://www.homepage-webhilfe.de/XML/XPath/> (zuletzt aufgerufen am 14.01.2023)

Christopher Culy. Proseminar Text Technologies SS 13, 2013.

Elliote Rusty Harold, W Scott Means, and Katharina Udemadu. *XML in a Nutshell*, volume 3. O'reilly Sebastopol, CA, 2004.

John Simpson. *XPath and XPointer: Locating Content in XML Documents*. Ö'Reilly Media, Inc.", 2002.