



mitp

Wilfried
Grube

XML

Grundlagen | Technologien
Validierung | Auswertung



Hinweis des Verlages zum Urheberrecht und Digitalen Rechtemanagement (DRM)

Der Verlag räumt Ihnen mit dem Kauf des ebooks das Recht ein, die Inhalte im Rahmen des geltenden Urheberrechts zu nutzen. Dieses Werk, einschließlich aller seiner Teile, ist urheberrechtlich geschützt. Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlages unzulässig und strafbar. Dies gilt insbesondere für Vervielfältigungen, Übersetzungen, Mikroverfilmungen und Einspeicherung und Verarbeitung in elektronischen Systemen.

Der Verlag schützt seine ebooks vor Missbrauch des Urheberrechts durch ein digitales Rechtemanagement. Bei Kauf im Webshop des Verlages werden die ebooks mit einem nicht sichtbaren digitalen Wasserzeichen individuell pro Nutzer signiert.

Bei Kauf in anderen ebook-Webshops erfolgt die Signatur durch die Shopbetreiber. Angaben zu diesem DRM finden Sie auf den Seiten der jeweiligen Anbieter.

Bibliografische Information der Deutschen Nationalbibliothek

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <<http://dnb.d-nb.de>> abrufbar.

Bei der Herstellung des Werkes haben wir uns zukunftsbewusst für umweltverträgliche und wiederverwertbare Materialien entschieden.

Der Inhalt ist auf elementar chlorfreiem Papier gedruckt.

ISBN 978-3-95845-755-3

1. Auflage 2018

www.mitp.de

E-Mail: mitp-verlag@sigloch.de

Telefon: +49 7953 / 7189 - 079

Telefax: +49 7953 / 7189 - 082

© 2018 mitp Verlags GmbH & Co. KG, Frechen

Dieses Werk, einschließlich aller seiner Teile, ist urheberrechtlich geschützt. Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlages unzulässig und strafbar. Dies gilt insbesondere für Vervielfältigungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Verarbeitung in elektronischen Systemen.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Lektorat: Sabine Janatschek

Sprachkorrektur: Petra Heubach-Erdmann

Covergestaltung: Christian Kalkert, www.kalkert.de

Satz: Wilfried Grupe

Druck: Plump Medienhaus GmbH, Rheinbreitbach

Bildnachweis Cover: [istock.com/lena_serditova](https://www.istock.com/lena_serditova)

Inhalt

Kapitel 1: Einleitung.....	11
Kapitel 2: XML.....	13
XML-Basics.....	19
XML: Wohlgeformte Dokumente.....	23
XML-Version.....	23
XML Encoding.....	24
XML-Entitäten.....	26
XML-Kommentare.....	29
XML: Processing-Instruction.....	36
XML-Datenstrukturen.....	40
XML: Die Sache mit den Namespaces.....	48
Namespaces in XML-Dokumenten.....	53
Die XML-Namespace-Flut.....	56
Versionierung.....	59
XML: Automatischer Namespace-Report.....	61
Wie kommt XML überhaupt zustande?.....	64
XML auswerten mit VisualBasic.NET.....	67
Kapitel 3: XML-Validierung.....	70
XML-Validierung: Wozu?.....	73
XML: Klare Strukturen.....	79
Hohe Fehlertoleranz und die Folgen.....	80
So etwas brauche ich nicht	81
Interessenkonflikte.....	84
RelaxNG compact - Beispiel.....	89
RelaxNG - Beispiel.....	89
DTD - Beispiel.....	90
XML-Schema - Beispiel.....	91
DTD.....	92
XML-Schema.....	94
XML-Schema 1.0.....	97
XML-Schema 1.1.....	118
XML-Schema Validierung in Java.....	130
XML-Schema: Datenvalidierung mit VisualBasic.NET.....	132
XML-Schema-Validierung mit ANT.....	135
XML-Schema-Datenvalidierung mit XProc.....	136
NVDL.....	137

Kapitel 4: XPath..... 139

XPath 3.0, XPath 2.0, XPath 1.0.....	140
XPath-Achsen.....	141
ancestor::*.....	143
ancestor-or-self::*.....	144
attribute::*.....	145
child::*.....	147
descendant::*.....	149
descendant-or-self::*.....	150
Verschachtelungstiefe.....	153
following::*.....	154
following-sibling::*.....	157
Positionsbestimmung bei following-sibling.....	157
namespace::*.....	160
parent::*.....	162
preceding::*.....	162
preceding-sibling::*.....	164
self::*.....	166
Automatische Generierung des XPath-Statements.....	166
XPath: Pfade, Prädikate.....	168
XPath-Operatoren.....	171
XPath-Funktionen.....	178
Zahlenfunktionen.....	200
Zeit ist Geld.....	212
Stringfunktionen.....	220
XPath: Sequenz-Funktionen.....	266
XPath 3.1: Map, xsl:map.....	330
XPath: transform.....	342
XPath 3.1: Array.....	345
available-environment-variables.....	357
system-properties.....	358
Der Namespace System.Xml.XPath.....	359
XPath in C#.NET.....	360

Kapitel 5: XSL..... 363

XSL-Übersicht.....	366
Funktionale Programmierung.....	368
XSL-Prozessoren.....	369
XSLT 3.0, XPath 3.0.....	372
xsl:accumulator.....	373
xsl:analyze-string.....	376
xsl:assert.....	377
xsl:attribute.....	378
xsl:attribute-set.....	379
xsl:apply-templates, xsl:next-match.....	381
xsl:apply-templates: Teilkonvertierung.....	386
xsl:for-each vs. xsl:apply-templates.....	388

xsl:call-template.....	389
xsl:character-map.....	390
Liste der Sonderzeichen selbst erstellen.....	391
Zeichensätze generieren mit C#.NET.....	395
xsl:choose.....	396
XSL-Analyse mit Collections.....	397
xsl:copy, xsl:copy-of.....	403
xsl:decimal-format.....	406
xsl:element.....	410
xsl:evaluate.....	412
xsl:fallback.....	414
xsl:fork.....	415
xsl:for-each select.....	416
xsl:for-each-group.....	419
xsl:function.....	430
xsl:if.....	431
xsl:include, xsl:import, xsl:apply-imports.....	432
xsl:import-schema.....	434
xsl:iterate, xsl:break.....	437
xsl:key.....	438
xsl:merge.....	440
xsl:message.....	443
xsl:namespace.....	444
xsl:number.....	446
Arbeiten mit optionalen Elementen.....	451
xsl:output.....	457
xsl:param.....	461
xsl:preserve-space, xsl:strip-space.....	467
xsl:result-document.....	469
sitemap.xml mit XSLT 3.0 generieren.....	470
xsl:sort, xsl:perform-sort, fn:sort.....	473
xsl:template.....	478
xsl:text.....	479
xsl:try/xsl:catch.....	480
xsl:value-of.....	482
xsl:variable.....	483
Schattenkabinett.....	486
XSLT 2.0: Erweiterte Syntax.....	488
XSLT-Konvertierung von XML nach HTML.....	490
Arbeiten mit xsl:for-each.....	490
Einbindung externer XML-Dokumente.....	493
Arbeiten mit xsl:apply-templates.....	496
Arbeiten mit xsl:template name/xsl:call-template.....	498
Spaltenweises Programmieren einer Tabelle.....	500
Spaltenweises Programmieren: pro Ort.....	505
Konvertierung von XML nach XML.....	507
Konvertierung von Elementen in Attribute.....	512
Arbeit mit temporären Bäumen.....	513
Erzeugen von skalierbaren Vektor-Grafiken (SVG).....	516
C#.NET in XSLT aufrufen.....	522

Konvertierung von XML nach Text.....	525
XSL-Transformationsaufrufe.....	529

Kapitel 6: XQuery..... 532

Was ist XQuery?.....	534
Arbeit mit Sequenzen.....	534
XSD-Type-Cast.....	535
Sortierung einer Sequenz.....	536
Arbeiten mit Variablen.....	539
XQuery: Arbeiten mit XML-Input.....	540
WHERE.....	541
XQuery: WHERE und Nummerierung.....	542
Geschachtelte Schleifen.....	545
FLOWR.....	546
XQuery: Element-Konstruktor.....	547
Vereinigte Sequenzen.....	548
XQuery: concat, union, intersect, except.....	549
XQuery: Generierung von 3erGruppen.....	551
XQuery: Arbeiten mit Namespaces und Funktionen.....	552
XPath 3.1: Arrays in XQuery.....	555
XQuery 3.0: switch/case.....	559
XQuery 3.0: try/catch.....	560
XQuery 3.0: Gruppierungen mit group by.....	560

Kapitel 7: XML-Datenbanken.....564

XML und Datenbanken.....	566
Der relationale Ansatz.....	566
XML-Dokumente in ORACLE 11g verwalten.....	575
XML-Datenbank: BaseX.....	577
Datenbank: INSERT und UPDATE.....	581

Kapitel 8: XProc..... 583

Kapitel 9: XML testen.....588

Geistreich, aber falsch gerechnet?.....	590
Selenium.....	598
Detailtests mit Schematron.....	600
XSLT Unit Tests mit XSpec.....	605

Kapitel 10: XML-Datenaustausch.....607

XML als Datenaustauschformat.....	607
Datenübertragung.....	610
XÖV: XML in der öffentlichen Verwaltung.....	613
Internet der Dinge (IoT).....	614

Objekt-Serialisierung mit C#.NET.....	616
Objekt-Serialisierung mit VisualBasic.NET.....	623
Objektserialisierung mit Java.....	627
JAXB.....	630
JAXB - XSLT - JAXB.....	638
XML auswerten mit Java-SAX.....	640
Java: DOM-Programmierung.....	642
JDOM-Programmierung.....	644
StAX.....	646
Maintenance.....	648
Best Practices.....	653

Kapitel 11: Formatting Objects (FO)..... 655

Die Struktur von Formatting Objects (FO).....	658
XSL-FO.....	662
Arbeiten mit XSL 3.0 und FOP.....	664
FOP mit ANT.....	667

Kapitel 12: Ratschläge für einen schlechten Programmierer.....670

Kapitel 1

Einleitung

XML hat Konjunktur. XML wird in zahlreichen Unternehmen, Behörden, Verwaltungen, Verlagen, Gerichten, Universitäten, Fachhochschulen, im weltweiten Web täglich millionenfach eingesetzt.

Wenn Sie eine Banküberweisung tätigen, Ihre Steuerklärung machen oder Bescheide erhalten, wenn Sie einen Zeitungsartikel oder ein Buch wie dieses lesen, so hat XML mit einiger Wahrscheinlichkeit einen Anteil an dessen technischem Zustandekommen. Es mag sein, dass Sie davon nichts sehen, denn dieser Anteil wird wieder ausgeblendet. Aber XML war vermutlich beteiligt, als unentbehrlicher Helfer.

Mit XML geht einfach alles. Die Verwendbarkeit ist scheinbar unbegrenzt. XML ist äußerst flexibel, bietet alle Voraussetzungen zu höchster Präzision in jeder Branche, über die verschiedensten IT-Systeme hinweg, ist unabhängig von internationalen Zeichensätzen und daher wichtig für den globalisierten Datenaustausch.

XML-Technologien gehören zu den grundlegenden Qualifikationen in der IT. Die fachlichen Kenntnisse werden immer häufiger als selbstverständlich vorausgesetzt. Daher wendet sich das Buch an alle Leser, die sich mit XML-Technologien befassen.

Das Hauptanliegen der Arbeit ist ein Überblick über die aktuellen XML-Technologien XML, XML-Schema, XPath (1.0, 2.0, 3.0, 3.1), XSLT (1.0, 1.1, 2.0, 3.0), XSL-FO, XQuery, XProc, Schematron und XSpec sowie die Unterstützung, die XML-Technologien in Java, C#.NET oder Datenbanken finden.

Zu Beginn eines jeden größeren Kapitels finden Sie einen Abschnitt "Grundlagen". Auf den jeweils folgenden Seiten gibt es vertiefende Details dazu, die vor allem bei den Abschnitten über XML-Schema, XPath und XSLT den Charakter einer Kurzreferenz haben.

Sie können das Buch so lesen, dass Sie zunächst die "Grundlagen" nacheinander durcharbeiten. Die Themen sind so aufgebaut, dass Sie sie aufgrund der bisherigen Lektüre verstehen können. Dann sind Sie je nach Lesetempo in wenig mehr als einer Stunde durch und haben einen generellen Überblick.

- Was¹⁾ ist ein XML-Dokument?
- Wie²⁾ können Sie eine einheitliche Datenstruktur absichern?

1) Seite: 13

2) Seite: 70

- Wie³⁾ können Sie XML-Dokumente via XPath auswerten?
- Wie⁴⁾ können Sie XML-Dokumente via XSL in andere Datenstrukturen transformieren?
- Inwiefern stellt XQuery⁵⁾ eine gute Alternative zu XSL dar?
- Wie⁶⁾ können Sie zahlreiche XML-Dokumente in XML-Datenbanken sichern und via XQuery auswerten?
- Wie⁷⁾ können Sie mit Schematron detailgenau testen, ob die erwarteten Ergebnisse stimmen?
- Inwiefern⁸⁾ ist XML ein ideales Datenaustauschformat?

Im Anschluß an die "Grundlagen" finden Sie reichlich Gelegenheit zum Stöbern in den fachlichen Details, die konkrete Hilfe für häufige Programmierprobleme bieten. Da jene Details fachlich ineinandergreifen, ist es sinnvoll, diesen Teil als Nachschlagewerk zu betrachten.

3) Seite: 139

4) Seite: 363

5) Seite: 532

6) Seite: 564

7) Seite: 588

8) Seite: 607

Kapitel 2

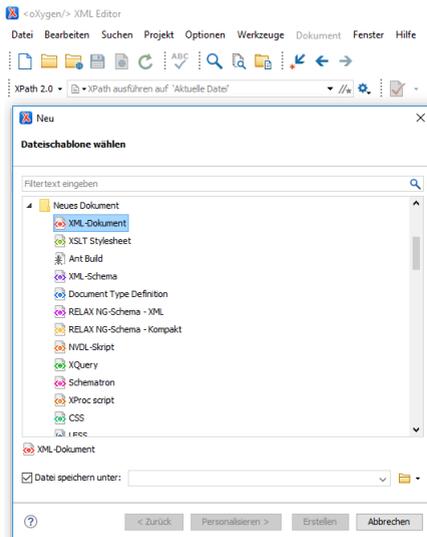
XML

Elemente, Attribute, Kommentare, Entitäten, Prolog, Processing Instruction, Namespaces sind zentrale Grundbegriffe in XML-Technologien.

XML-Grundlagen

XML-Dateien sind simple Textdokumente, die einige strukturelle Voraussetzungen erfüllen müssen. Zwar ist es möglich, XML-Dateien mit sehr einfachen Texteditoren zu schreiben. Diese bieten jedoch häufig nur eine geringe fachgerechte Unterstützung.

Weit effizienter ist daher die Arbeit mit professionellen Editoren, die Sie in jeder Hinsicht unterstützen. Einer dieser Editoren ist der Oxygen XML Editor 19.1. Über DATEI|NEU erhalten Sie dieses Fenster:



Mit der Auswahl *XML-Dokument* generiert der Editor eine Datei, in der der sogenannte XML-Prolog bereits enthalten ist.

```

Unbenannt1.xml* x
1  <?xml version="1.0" encoding="UTF-8"?>
2  |

```

Diesen Prolog können Sie zunächst löschen, um sich mit den Grundlagen zu befassen.

Elemente

Grundsätzlich hat jedes XML-Dokument genau ein Root-Element. Dieses Element kann so aussehen:

```
<Personen></Personen>
```

Das Element besteht aus einem Start-Tag `<Personen>` und einem Ende-Tag `</Personen>`.

Wichtig ist dabei die korrekte Schreibweise. Haben Sie das Start-Tag `<Personen>` genannt, dann muss das Ende-Tag genau so geschrieben werden, lediglich mit dem `"/`-Slash-Zeichen versehen. Das Ende-Tag anders zu schreiben, etwa `</personen>` oder `</PERSONEN>`, wäre von vornherein falsch: Das XML-Dokument wäre nicht wohlgeformt.

Für die Benennung der Elemente gibt es einige Einschränkungen, auf die ich später⁹⁾ noch einmal eingehe.

Nun hindert Sie niemand daran, zwischen Start- und Ende-Tag eines Elements weitere Elemente einzufügen. Zum Beispiel so:

```
<Personen>
  <Person></Person>
  <Person></Person>
</Personen>
```

Das Verfahren können Sie beliebig ausbauen, zudem können Sie zwischen den jeweiligen Start- und Ende-Tags auch normalen Text einfügen.

```
<Personen>
  <Person>
    <Vorname>Susi</Vorname>
    <Nachname>Sinnlos</Nachname>
  </Person>
  <Person>
    <Vorname>Alfons</Vorname>
    <Nachname>Achtlos</Nachname>
  </Person>
</Personen>
```

Wenn Sie dieses XML-Dokument in eine Datei "Personen.xml" abspeichern, dann können Sie diese Datei mit einem normalen Internet-Browser öffnen. Der Firefox zeigt die Datei beispielsweise so an:

9) Seite: 23

```

-<Personen>
  -<Person>
    <Vorname>Susi</Vorname>
    <Nachname>Sinnlos</Nachname>
  </Person>
  -<Person>
    <Vorname>Alfons</Vorname>
    <Nachname>Achtlos</Nachname>
  </Person>
</Personen>

```

Andere Programme können das genauso gut. Gute Power-Tools zeigen darüber hinaus noch eine andere Darstellung.

Personen	Person	Vorname	Nachname
	(2 rows)	1 Susi	Sinnlos
		2 Alfons	Achtlos

Angenommen, Sie kennen von einer Person weder den Vornamen noch den Nachnamen. Wenn zwischen dem Start- und Ende-Tag "Person" kein Inhalt abgebildet werden kann, wenn das Element also leer ist, dann handelt es sich sinnigerweise um ein "leeres Element", bei dem Sie sich das Ende-Tag sparen können, sofern der "/" am Ende des Start-Tags erscheint. Beispiel:

```

<Personen>
  <Person/>
</Personen>

```

Attribute

Alternativ zur eben dargestellten Elementschreibweise können Sie Inhalte auch als Attribute einfügen. Ein Attribut wird grundsätzlich in das Start-Tag eines Elements geschrieben. Da die ursprünglichen Kind-Elemente entfallen können, schreiben Sie die Information kurzerhand als leeres Element, aber mit Attributen.

```

<Personen>
  <Person Vorname="Susi"
    Nachname="Sinnlos" />
  <Person Vorname="Alfons"
    Nachname="Achtlos" />
</Personen>

```

Speichern Sie das in einer Datei "Personen2.xml" ab, so zeigt der Internetbrowser das so an:

```

-<Personen>
  <Person Vorname="Susi" Nachname="Sinnlos"/>
  <Person Vorname="Alfons" Nachname="Achtlos"/>
</Personen>

```

Auch das XML-Power-Tool macht mit und zeigt diese Darstellung, die sich nur in einer scheinbar winzigen Kleinigkeit von der vorherigen unterscheidet: das "@" in *@Vorname* und *@Nachname*.

Personen	Person	@Vorname	@Nachname
	(2 rows)		
	1	Susi	Sinnlos
	2	Alfons	Achtlos

Die Attributschreibweise hat gegenüber der Elementschreibweise nur einen Nachteil: Jedes Attribut darf nur ein einziges Mal vorkommen. Hat eine Person also mehrere Vornamen, so bietet es sich an, den Nachnamen als Attribut und die Vornamen in Elementschreibweise darzustellen.

```
<Personen>
  <Person Nachname="Holzflos">
    <Vorname>Hugo</Vorname>
    <Vorname>Helmut</Vorname>
    <Vorname>Horst</Vorname>
  </Person>
  <Person Nachname="Nixlos">
    <Vorname>Tanja</Vorname>
    <Vorname>Theodora</Vorname>
  </Person>
</Personen>
```

Kommentare

Ein großer Vorteil von XML-Dokumenten ist, dass Sie auch Kommentare einfügen können. Das erleichtert die Lesbarkeit sehr und ist auch in der praktischen Arbeit von erheblicher Bedeutung, nicht zuletzt bei der Fehlersuche. XML-Kommentare¹⁰⁾ beginnen mit "`<!--`" und enden mit "`-->`".

```
<Personen>
  <!-- Attributschreibweise -->
  <Person Vorname="Resi"
    Nachname="Denzschlos"/>
  <!-- Elementschreibweise -->
  <Person>
    <Vorname>Lotte</Vorname>
    <Nachname>Rielos</Nachname>
  </Person>
</Personen>
```

Entitäten

Sie haben schon bemerkt, dass jedes Tag mit einem "`<`" beginnt und mit einem "`>`" endet. Diese beiden Zeichen haben eine zentrale Bedeutung in XML. Das kann Sie jedoch in

10) Seite: 29

einige Verlegenheit bringen, wenn Sie die hochwichtige Information "3 < 4" in XML abbilden möchten.

```
<Info>3 < 4</Info>
```

Dieser Versuch geht schief. Jedes XML-Tool, das etwas auf sich hält und nicht außergewöhnlich leidensfähig ist, meckert Sie an:

```
XML-Verarbeitungsfehler: nicht wohlgeformt
```

oder im schönsten IT-Deutsch:

```
System-Fehlerlevel: error
The content of elements must consist
of well-formed character data or markup.
```

Hier bleibt Ihnen nur übrig, das "<" als Entität zu deklarieren und entsprechend zu kennzeichnen.

```
<Info>3 &lt; 4</Info>
```

Es gibt mehrere Standard-Entitäten¹¹⁾ und auch die Möglichkeit, eigene Entitäten zu definieren. Dazu später mehr.

Prolog

Häufig finden Sie am Anfang eines XML-Dokuments einige Zusatzinformationen, die fast durchweg so aussehen:

```
<?xml version="1.0"?>
<Personen>
  <Person/>
</Personen>
```

<?xml version="1.0"?> ist die Minimalinformation. Sie besagt, dass es sich um ein XML-Dokument in der Version 1.0¹²⁾ handelt.

Hin und wieder finden Sie dort auch Informationen zum verwendeten Encoding¹³⁾. Darin ist ein Hinweis auf den im Dokument verwendeten Zeichensatz enthalten. Fehlt diese Angabe, dann handelt es sich per Default um den Zeichensatz "UTF-8".

11) Seite: 26

12) Seite: 23

13) Seite: 24

```
<?xml version="1.0"
      encoding="ISO-8859-1"?>
<Personen>
  <Person/>
</Personen>
```

Processing Instruction

Neben dem Prolog können Sie noch Verarbeitungsanweisungen¹⁴⁾ mitgeben für das Programm, das das XML-Dokument auswertet. Beispiel:

```
<?versandperMail dort@woderpfefferwaechst.de"?>
```

Das verarbeitende Programm wertet diese Information aus und nimmt (hoffentlich) wohlwollend zur Kenntnis, was hier angemerkt wird.

```
<?xml version="1.0"
      encoding="ISO-8859-1"?>
<?sorry ich@habnichtigemacht.gov"?>
<Personen>
  <Person/>
</Personen>
```

Namespace

Namespaces¹⁵⁾ sind ebenfalls ein zentrales Thema in XML. Sie eröffnen weitreichende Möglichkeiten, ein XML-Dokument einem Namensraum zuzuordnen und dabei auch versionsbedingte Informationen einzubinden.

```
<Personen xmlns="www.Kundenliste.de/2018">
  <k:Person xmlns:k="www.besondererKunde.de">
    <k:Vorname>Wanja</k:Vorname>
    <k:Nachname>Wunschlos</k:Nachname>
  </k:Person>
</Personen>
```

- Zu unterscheiden sind hier Namespace-Präfixe, die ein Kürzel für einen Namensraum definieren, sowie Default-Namespaces.
- Im obigen Beispiel ist `xmlns="www.Kundenliste.de/2018"` der Default-namespace.
- `xmlns:k="www.besondererKunde.de"` definiert einen Namensraum mit einem Kürzel `k`, dem sogenannten Namespace-Präfix, das bei den Elementen `k:Person`, `k:Vorname` und `k:Nachname` zum Einsatz kommt.

Stören Sie sich nicht daran, wenn einzelne Internet-Browser sich weigern, die Namespaces anzuzeigen, nachdem Sie das Dokument in eine XML-Datei gespeichert

14) Seite: 36

15) Seite: 48

haben und den Browser auffordern, diese anzuzeigen. Die Tools arbeiten da unterschiedlich, und mit einem guten XML-Editor finden Sie ohnehin alles wieder.

XML-Datenstrukturen

Bereits die wenigen Anmerkungen verdeutlichen die enorme Flexibilität, die XML zu bieten hat. Es ist praktisch alles darstellbar, was in das menschliche Hirn hineinpasst, egal wie komplex ein Sachverhalt auch sein mag. Vorausgesetzt, Sie halten sich an einige Grundregeln der Wohlgeformtheit.

Das begrenzt sich durchaus nicht auf die klar strukturierten Daten, die ich bisher beschrieben habe. Auch Folgendes¹⁶⁾ ist nicht nur möglich, sondern kommt recht häufig vor:

```
<?xml version="1.0"
    encoding="ISO-8859-1"?>
<Meldung>Mit XML durch das Weltall,
zum <fett>Mars</fett>,
<kursiv>Jupiter</kursiv>
und in die Milchstrasse.</Meldung>
```

Die Kernfrage bei dieser überwältigenden Flexibilität lautet: Wie schaffen Sie es, Programme zu schreiben, die damit klarkommen? Darum geht es in diesem Buch.

XML-Basics

XML ist eine erweiterbare, flexible, stukturierte Markup-Sprache, die in unterschiedlichen Bereichen zum Einsatz kommt, etwa bei Transformation zu HTML, XML, Text, SVG, RTF, PNG, TIFF, PDF.

```
<?xml version="1.0" encoding="UTF-8"?>
<?lernen was="XML XSLT XSD XQuery"?>
<r:root xmlns:r="Namespaces">
    <child attribut="JA"/>
    <![CDATA[ <Spezial/> ]]>
</r:root>
```

Was ist das für eine Sprache,

- die so einfach strukturiert ist, dass man ihre Grundlagen leicht in einer halben Stunde lernen kann?

16) Seite: 40

- die nicht auf einem eigenen Compiler, Interpreter, anderen Übersetzer basiert, der sprachspezifische Kommandos in Maschinensprache umwandelt?
- die (abgesehen von einer äußerst knappen Formalstruktur) über fast keine Schlüsselworte oder Vokabular verfügt, aus denen eine Sprache normalerweise mindestens besteht?
- die von sich aus praktisch gar nichts mitbringt, aber dennoch eine äußerst präzise, versionsbezogene Datenstrukturdefinition für die unterschiedlichsten Branchen ermöglicht?
- die zudem alle Möglichkeiten für eine leichthändige Datenkonvertierung zwischen den unterschiedlichsten Datenformaten und Zeichensätzen bietet?
- die selbst de facto überhaupt nichts "tut" oder "kann", sondern sämtliche Möglichkeiten aus der breiten Unterstützung anderer Technologien bzw. Sprachen bezieht?
- die von sich aus weder Technologien zu Systemkonfiguration, Datenaustausch, GUI, Datenbankzugriffen, Prozessdefinition, Automatisierung, präzisen Definition hochkomplexer Datenstrukturen, Transformation beliebiger Datenstrukturen in beliebige Zielformate, leistungsfähiger Publishing-Standards und vielem Anderen mehr bereitstellt, aber in all diesen Bereichen hervorragend zum Einsatz kommt?

Auswertung, Transformation	XPath, XSL, XQuery
Webseiten	HTML
Verknüpfung	XInclude, XLink, XPointer
Validierung, Testen	DTD, XML-Schema (XSD), Relax NG, Schematron, XSpec, NVDL
Signatur, Verschlüsselung	XML Signature, XML Encryption
Arbeit in verteilten Systemen	XML-RPC
EDI	XML, daneben Übertragungen von EDIFACT, ANSI X.12, SAP IDoc, CSV, CargoImp u.a.m.
Finanzberichte	XBRL
Formatierung	FO, XSL-FO, MathML, SVG, EPUB, WordML
Publishing-Standards	DITA, DocBook, TEI
Grafiken	SVG, X3D
Formatting Object	TIFF, PNG, PDF, PS, PCL
Geodaten	CityGML (City Geography Markup Language), GML (Geography Markup Language), GPX (GPS Exchange Format),

	KML (Google Earth: Keyhole Markup Language)
Landwirtschaft	AgroXML
Prozessdefinition	ANT, XProc
Webservices	SOAP, WSDL, REST
Office-Anwendungen	OASIS Open Document Format for Office Applications, RTF

XML: erweiterbar, flexibel

Zunächst ist XML eXtensible, also eine erweiterbare Markup Language. Sie besteht ausdrücklich nicht aus einer endlichen Menge von einigen Dutzend Schlüsselwörtern, die man (in diversen Programmiersprachen) kennen muss, um Code schreiben zu können, den der PC dann ausführt. Sondern XML ist erweiterbar. Potenziell können unendlich viele Begriffe definiert werden, in den unterschiedlichsten Sprachen, Zeichensätzen und Schreibweisen. Die einzige Bedingung ist, einige wenige grundlegende Anforderungen einzuhalten.

Zweitens handelt es sich um eine äußerst flexible Markierungssprache, die wohlgeformte, strukturierte Daten definiert und obendrein eine präzise Zuordnung zu einem fachlichen Kontext, auch mit Versionsunterschieden erlaubt. XML zieht seine Effizienz jedoch aus einer breiten Unterstützung durch andere Technologien.

Dabei kann XML sehr schwach strukturiert (dokumentzentriert: Die XML-Elemente dienen zur semantischen Strukturierung des Textes, was eine maschinelle Verarbeitung erschwert), sehr stark strukturiert (datenzentriert: XML-Elemente, Attribute etc. folgen einer klaren Strukturdefinition zur effizienten maschinellen Verarbeitung), aber auch *mixed content* haben (semistrukturiert: XML-Dokumente als Mischung aus starker und schwacher Strukturierung; die Ansätze zur effizienten Auswertung sind hier andere als bei Datenzentrierung).

Die vorliegende Arbeit hat ihren Schwerpunkt auf datenzentrierten Dokumenten. XML ist hier nie Selbst- oder Endzweck, sondern Teil einer Verarbeitungskette: Immer ist ein Folgeprogramm nötig, das mit der jeweiligen XML-Datenstruktur umzugehen weiß. Das mag ein Webbrowser sein, der skalierbare Vektorgrafiken (SVG, eine Spezialform von XML) anzeigt. Ebenso kann es ein Systemprogramm sein, das eine in XML definierte System- oder Serverkonfiguration auswertet. Auch kann es sich um ein in Java¹⁷⁾, C#.NET¹⁸⁾, VisualBasic.NET¹⁹⁾, C++ oder in einer anderen Sprache geschriebenes Programm handeln, das in XML definierte Prozesse schrittweise abarbeitet.

17) Seite: 627

18) Seite: 616

19) Seite: 623

Die Struktur der in XML vorliegenden Daten und die Programme, die sie auswerten, müssen also Hand in Hand gehen. XML-Datenstrukturen, die ein auswertendes Programm nicht verarbeiten kann, sind wirkungslos.

Die Programme folgen unterschiedlichen Verarbeitungsmodellen. SAX verarbeitet XML als sequenziellen Datenstrom und hält für bestimmte Ereignisse spezielle *callback functions* bereit. Das sehr speicherintensive DOM betrachtet XML dagegen auf der Baumstruktur und gewährt wahlfreien Zugriff mit Manipulationsmöglichkeiten. Daneben stehen noch die Pull-API (sequenzielle Verarbeitung mit Iterator) oder die Verarbeitung auf Byte-Ebene bereit. Häufig werden Objekte in XML-Dokumente umgewandelt (Serialisierung) oder umgekehrt (marshalling); siehe JAXB oder XML-Schema-Definition-Toolkit in .NET.

Trotzdem ist es nicht notwendig, die interne Logik jener Programme, die XML-Dokumente auswerten, zu kennen, um XML-Dokumente schreiben zu können, die das Programm auswerten kann. Es reicht aus, klare Vorgaben hinsichtlich der Struktur und Detailtypen der zu übermittelnden Daten verfügbar zu haben. Hier helfen diverse Standards weiter, die eine Strukturdefinition der XML-Dokumente erlauben, unter anderem DTD und XML-Schema (XSD).

XML-Schema erlaubt, die grundsätzlich extrem gestaltungsflexiblen Möglichkeiten von XML-Dokumenten auf eine endliche Anzahl zulässiger Element- und Attributnamen einzugrenzen, verbunden mit einer begrifflichen Zuordnung zu bestimmten Namespaces. Auf diese Weise wird eine Datenstruktur definierbar, die von Folgeprogrammen ausgewertet werden kann. Es besteht die Möglichkeit einer Vorprüfung (Validierung), ob das jeweilige XML-Dokument diesen Vorgaben entspricht; falls nicht, kann die Weiterverarbeitung gestoppt werden.

Freilich kommt es recht häufig vor, dass die verfügbaren Daten in einer Struktur vorliegen, die die Folgeprogramme nicht auswerten können. Dann wird eine Konvertierung erforderlich. Hier hilft XSLT, ab XSLT 2.0 auch dann, wenn die ursprünglichen Daten gar nicht in XML vorliegen, sondern beispielsweise in Textformaten wie CSV.

XSL ist selbst auch ein XML-Dokument, jedoch mit einem exakt definierten Aufbau sowie einer Reihe klar definierter Schlüsselbegriffe. XSL bietet (in Kombination mit XPath) sehr effiziente Ansätze zur Transformation vorliegender Daten in andere Strukturen.

XSL und XPath bieten effiziente Möglichkeiten zur Transformation von (strukturierten) XML-Dokumenten in diverse Zielformate, z.B. HTML, XML, Text, SVG, RTF, PNG, TIFF, PDF und andere mehr. Dabei können mehrere XML-Quelldokumente ebenso berücksichtigt werden wie mehrere Ausgabedokumente.

Grundvoraussetzung für effizientes Programmieren mit XSL ist, dass sowohl die Struktur des XML-Quelldokuments als auch die Struktur des Zielformats zweifelsfrei klar ist. Während das XML-Inputdokument die (hoffentlich) klar strukturierten Daten liefert, stehen in XSL/XPath jene Programmieranweisungen, die die Struktur des Quelldokuments in die gewünschte Zielstruktur konvertiert. Das Kind der Ehe von XML und XSL ist das gewünschte Dokument.

XML: Wohlgeformte Dokumente

XML-Dokumente sind wohlgeformt, sofern sie nicht gegen eine XML-Regel verstoßen, z.B.:

- Jedes XML-Dokument hat genau ein Root-Element.
- Elemente dürfen nicht mehrere Attribute mit demselben Namen haben.
- Alle Elemente mit Child-Elementen (auch Textknoten) haben ein Start- und ein Ende-Tag, z.B. `<content>Inhalt</content>`, die innerhalb ihres Parent-Knotens abgeschlossen sein müssen.
- Elemente ohne Childnodes sind leer, sie benötigen kein Ende-Tag, sondern können einfach geschlossen werden mit `>` (z.B. `<content/>`). Eventuell weisen sie Attribute auf, wie im folgenden Beispiel:

```
<content mycontent="XML is my favourite"/>
```

Für die Benennung von Elementen und Attributen gibt es einige Einschränkungen. So darf das erste Zeichen keine Zahl sein. Ebenso ist eine ganze Reihe von Sonderzeichen ganz oder eingeschränkt verboten, die als logische oder Rechenoperatoren oder Bestandteile von XPath-Statements zum Einsatz kommen können:

```
+ - * / \ & " ' % # ! = ; ( ) [ ] { } @ $ § ? .
```

: kommt im Zusammenhang mit Namespaces, Blanks kommen bei Attributdeklaration zum Einsatz.

`<fe-ld/>` oder `<_fe.ld/>` ist erlaubt, als erstes Zeichen `<-.feld/>` oder `<_feld/>` aber nicht. Problemlos: `<_1feld/>`.

XML-Version

Im Normalfall gehört zu jedem XML-Dokument ein Prolog, der mindestens über die XML-Version informiert.

In den allermeisten Fällen hat die Version den Wert "1.0".

```
<?xml version="1.0"?>
```

Nur in sehr seltenen Ausnahmefällen (die Daten enthalten Steuerzeichen wie den vertikalen Tabulator, Zeilenvorschub oder Inhalte in selten verwendeten Sprachen), und wenn die verwendeten Parser damit umgehen können, kann die Versionsbenennung von 1.1 vorteilhaft sein.

XML Encoding

Speziell im internationalen Datenaustausch werden unterschiedliche Codierungen verwendet. Sofern die in XML verwendeten Zeichen nicht aus dem UTF-8-Encoding stammen, ist im XML-Prolog das verwendete Encoding anzugeben.

Wer konsequent mit UTF-8 arbeitet, erspart sich den Umgang mit Sonderzeichen. Das ist aber nicht immer möglich: Auch in ISO-8859-1 gibt es eine längere Reihe von Sonderzeichen, die in XML nicht durch HTML-Entitätsreferenzen abgedeckt sind, etwa für das EURO-Zeichen.

Ein Aufruf der HTML-Referenz `€` führt in XML zu einem Fehler, hier muss mit `€` bzw. deren Hexwert `€` gearbeitet werden. In

- http://www.w3schools.com/charsets/ref_html_8859.asp
- <https://wiki.selfhtml.org/wiki/Referenz:HTML/Zeichenreferenz>
- http://docstore.mik.ua/orelly/xml/xmlnut/ch26_01.htm

finden Sie sehr brauchbare Übersichten.

Das Encoding definiert die Zeichencodierung, die im Dokument verwendet werden soll. UTF-8 ist dabei die Standard-Codierung. Bitte beachten Sie, dass nicht alle Parser auch sämtliche Encodings unterstützen. Einige Tools ignorieren das angegebene Encoding und arbeiten grundsätzlich mit UTF-8.

UTF-8	Standard-Encoding in XML-Dokumenten. UTF-8 ist so designt, dass alle ASCII Dokumente legale UTF-8-Dokumente darstellen, was bei UTF-16 und Latin1 nicht der Fall ist.
UTF-16	Ein Zwei-Byte-Encoding von Unicode, das alle Zeichen von Unicode 3.0 und früher umfasst.
ISO-10646-UCS-2	Die Basis-Multilingual-Version von Unicode; der Zeichensatz ist weitgehend identisch mit UTF-16. Der Unterschied betrifft lediglich Unicode 3.1 und höher.
ISO-10646-UCS-4	Ein Vier-Byte-Encoding von Unicode.
ISO-8859-1	Latin-1, ASCII plus jene Zeichen, die für die meisten westeuropäischen Sprachen verwendet werden, inkl. Dänisch, Deutsch, Holländisch, Englisch, Finnisch, Flämisches, Irisch, Isländisch, Italienisch, Norwegisch, Portugiesisch, Spanisch und Schwedisch.

ISO-8859-2	Latin-2, ASCII plus jene Zeichen, die für die meisten zentraleuropäischen Sprachen verwendet werden, inkl. Kroatisch, Tschechisch, Ungarisch, Polnisch, Slowakisch, Slowenisch.
ISO-8859-3	Latin-3, ASCII plus jene Zeichen für Esperanto, Galizisch, Maltesisch, Türkisch.
ISO-8859-4	Latin-4, ASCII plus jene Zeichen für Lappländisch, Lettisch, Litauisch, Grönländisch. Wurde weitgehend ersetzt durch ISO-8859-10, Latin-6.
ISO-8859-5	ASCII plus die kyrillischen Zeichen für Belorussisch, Bulgarisch, Mazedonisch, Russisch, Serbisch, Ukrainisch.
ISO-8859-6	ASCII plus Arabisch.
ISO-8859-7	ASCII plus Griechisch.
ISO-8859-8	ASCII plus Hebräisch.
ISO-8859-9	Latin-5, weitgehend identisch mit Latin-1 (ASCII plus westeuropäisch), aber ohne bestimmte türkische und isländische Zeichen.
ISO-8859-10	Latin-6: Zeichen für nordeuropäische Sprachen wie Grönländisch, Isländisch, Lappländisch, Litauisch. Ähnlich wie Latin-4, ergänzt in ISO-8859-13.
ISO-8859-11	ASCII plus Thai. Die Unterstützung von XML-Prozessoren ist nicht optimal.
ISO-8859-12	Nicht benötigt.
ISO-8859-13	Alternativer Zeichensatz für baltische Sprachen. Vgl. Latin 6.
ISO-8859-14	Latin-8; eine Variante für Latin-1 mit Zusatzzeichen für Gälisch und Welsch.
ISO-8859-15	Latin-9; Revision von Latin-1. Weitestgehend identisch mit ISO-8859-1.
ISO-8859-16	Latin-10; für Rumänisch.

ISO-2022-JP	Sieben-Bit-Encoding mit japanischem Zeichensatz JIS X-0208-1997, in E-Mails und im Web verwendet, siehe RFC 1468.
Shift_JIS	Japanischer Zeichensatz JIS X-0208-1997, in Microsoft Windows verwendet.
EUC-JP	Japanischer Zeichensatz JIS X-0208-1997, in den meisten UNIX-Varianten verwendet.

XML-Entitäten

XML-Entitäten finden Verwendung in Form von Standard-, selbst definierten Entitäten, Einbindung separater XML-Dokumente sowie zur Entitätsdeklaration in XSLT.

Folgende Standard-Entitäten sollten geläufig sein:

```
&lt; für <  
&gt; für >  
&quot; für "  
&amp; für &
```

Selbst definierte Entitäten

Ergänzend können selbst definierte Entitäten zum Einsatz kommen:

```
<?xml version="1.0" encoding="ISO-8859-1"?>  
<!DOCTYPE mytext [  
<!ENTITY LvH "Leute von heute" >  
>  
<mytext>  
Hallo, &LvH;  
</mytext>
```

Das Ergebnis sieht im Browser dann so aus:

```
<?xml version="1.0" encoding="ISO-8859-1"?>  
<!DOCTYPE mytext>  
<mytext> Hallo, Leute von heute </mytext>
```

Entitäten als externe Dokumente

Darüber hinaus besteht die Möglichkeit, mittels selbst definierter Entitäten externe XML-Dokumente in ein größeres Dokument einzubinden. Auf diese Weise können mehrere

Autoren parallel an unterschiedlichen Dokumenten arbeiten; nach abgeschlossener Arbeit werden die XML-Dokumente in ein umfassenderes Dokument eingefügt.

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<!DOCTYPE article [
<!ENTITY KAP1 SYSTEM "Kapitell.xml">
]>
<article>
  <chapter>
    <title>XML-Basics</title>
    &KAP1;
  </chapter>
</article>
```

Wobei das externe Dokument etwa so aussehen könnte:

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<section>
  <title>Attribute</title>
  <para>Attribute sind ebenso sinnvoll wie Elemente.</para>
</section>
```

Im Arbeitsspeicher werden die Dokumente dann zusammengefügt zu:

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<article>
  <chapter>
    <title>XML-Basics</title>
    <section>
      <title>Attribute</title>
      <para>Attribute sind ebenso sinnvoll wie Elemente.</para>
    </section>
  </chapter>
</article>
```

Entitätsdeklaration in XSLT

Je nach XSLT-Prozessor²⁰⁾ können im Vorspann zu XSLT auch Entitäten deklariert werden (bei diversen Prozessorvarianten habe ich hier temporäre Probleme gefunden, die teilweise bereits beseitigt sind).

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE stylesheet [
<!ENTITY euro    "€" >
<!ENTITY auml   "ä" >
<!ENTITY Auuml  "Ä" >
<!ENTITY ouuml  "ö" >
<!ENTITY Ouuml  "Ö" >
<!ENTITY uuml   "ü" >
<!ENTITY Uuuml  "Ü" >
<!ENTITY szlig  "ß" >
```

20) Seite: 369

```
]>
<xsl:stylesheet version="2.0"
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
  <xsl:output method="html" />
    <xsl:decimal-format
      name="df"
      decimal-separator=","
      grouping-separator="."
      minus-sign="-"
      digit="#" />
  <xsl:template match="/">
    <html>
      <head>
        <meta
          http-equiv="Content-Type"
          content="text/html; charset=UTF-8" />
      </head>
      <body>
        <xsl:value-of
          select="format-number(
            sum(//Gehalt),
            '#.##0,00 &euro;',
            'df')" />
        <br />
        <table>
          <tr>
            <td>ae</td>
            <td>&auml;</td>
            <td>
              <xsl:text>&auml;</xsl:text>
            </td>
          </tr>
          <tr>
            <td>Ae</td>
            <td>&Auml;</td>
            <td>
              <xsl:text>&auml;</xsl:text>
            </td>
          </tr>
          <tr>
            <td>oe</td>
            <td>&ouml;</td>
            <td>
              <xsl:text>&ouml;</xsl:text>
            </td>
          </tr>
          <tr>
            <td>Oe</td>
            <td>&Ouml;</td>
            <td>
              <xsl:text>&Ouml;</xsl:text>
            </td>
          </tr>
          <tr>
            <td>ue</td>
            <td>&uuml;</td>
            <td>
              <xsl:text>&uuml;</xsl:text>
            </td>
          </tr>
          <tr>
            <td>Ue</td>
            <td>&Uuml;</td>
            <td>
              <xsl:text>&Uuml;</xsl:text>
            </td>
          </tr>
        </table>
      </body>
    </html>
  </template>
</xsl:stylesheet>
```

```

        <xsl:text>&Uuml;</xsl:text>
    </td>
</tr>
<tr>
    <td>szlig</td>
    <td>&szlig;</td>
    <td>
        <xsl:text>&szlig;</xsl:text>
    </td>
</tr>
</table>
</body>
</html>
</xsl:template>
</xsl:stylesheet>

```

Das Ergebnis sieht im Browser dann so aus:

23.816,77 €

ae ä ä

Ae Ä ä

oe ö ö

Oe Ö Ö

ue ü ü

Ue Ü Ü

szlig ß ß

XML-Kommentare

In XML, XSLT, XML-Schema und anderen Standards sind ergänzende Kommentare hilfreich, um die Wartung komplexer Anwendungen zu erleichtern. Neben Standard-Kommentaren stehen auch *CDATA*-Kommentare bereit.

Ergänzende Kommentare können weitere Hilfestellung geben. Ein Standard-XML-Kommentar beginnt mit `<!--` und endet mit `-->`.

```

<root>
  <!-- this is a comment, you are free to write down your CV -->
</root>

```

XML-Kommentare in XSL erzeugen

In XSL ist es angebracht, zum besseren Verständnis der Programmlogik lokale Kommentare einzubauen, die nicht im Ergebnisdokument erscheinen. Sinnvolle

Kommentare im Quelltext können die Wartung der Programme sehr erleichtern, daher sind sie unbedingt zu empfehlen.

```
<erg>
  <!-- dieser Kommentar gilt nur lokal in XSL,
        wird nicht im Ergebnisdokument erscheinen -->
</erg>
```

Die Ausgabe im Ergebnisdokument lautet wie beabsichtigt ohne Kommentar:

```
<erg/>
```

Um im Ergebnisdokument einen Kommentar sichtbar zu machen, können Sie mit `<xsl:comment>` arbeiten.

```
<erg>
  <xsl:comment>Dieser Kommentar wird im Ergebnis erscheinen</xsl:comment>
</erg>
```

Die Ausgabe im Ergebnisdokument lautet:

```
<erg>
  <!--Dieser Kommentar wird im Ergebnis erscheinen-->
</erg>
```

XML-Kommentare in XSL auswerten

Umgekehrt ist es auch möglich, mit XSLT die Kommentare in den Input-Dokumenten auszuwerten. Betrachten Sie folgendes XML-Input-Dokument, das (abgesehen von einem "root"-Node) ausschließlich XML-Kommentare aufweist.

```
<root>
  <!--Kommentar 1-->
  <!--Kommentar 2-->
  <!--Kommentar 3-->
  <!--Kommentar 4-->
</root>
```

Nun soll versucht werden, diese Kommentare in XSLT auszuwerten. Das funktioniert recht einfach:

```
<xsl:stylesheet version="1.0"
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
  <xsl:output method="xml"
    version="1.0"
    encoding="UTF-8"
    indent="yes"/>
  <xsl:template match="/">
```