

## IDŹ DO

PRZYKŁADOWY ROZDZIAŁ



SPIS TREŚCI

## KATALOG KSIĄŻEK

KATALOG ONLINE

ZAMÓW DRUKOWANY KATALOG

## TWÓJ KOSZYK

DODAJ DO KOSZYKA

## CENNIK I INFORMACJE

ZAMÓW INFORMACJE  
O NOWOŚCIACH

ZAMÓW CENNIK

## CZYTELNIA

FRAGMENTY KSIĄŻEK ONLINE

## XML na poważnie

Autorzy: Przemysław Kazienko, Krzysztof Gwiazda  
 ISBN: 83-7197-765-4  
 Format: B5, stron: 444  
 Zawiera CD-ROM



Książka ta to przegląd głównych standardów związanych z językiem XML:

- XML 1.0.
- DTD – opis struktury dokumentu.
- XSLT – transformacje dokumentów.
- XPath – nawigacja i wyszukiwanie.
- XPointer – wskazania wewnątrz dokumentów.
- DOM (Level 1, 2, 3) – dostęp do dokumentu jako drzewa węzłów.
- SAX – zdarzeniowe przetwarzanie dokumentów.
- Namespaces – przestrzenie nazw.

a także najnowszych standardów zatwierdzonych w 2001 r.:

- XML Schema – schematy zawartości.
- XLink – odsyłacze hipertekstowe.
- XSL FO – prezentacja danych.
- XML Base – adresy odniesienia.
- SVG – grafika wektorowa.

Interesująco zapowiada się część poświęcona zastosowaniom języka XML:

- serwisy internetowe,
- elektroniczna wymiana danych (EDI),
- bazy danych XML – XDBMS,
- pseudostrukturalne dane tekstowe,
- systemy prezentacji i systemy interaktywne,
- dane specjalistyczne (np. MathML, CML).

Autorzy umieścili informacje o ograniczeniach poszczególnych standardów, jak również różnorodne zalecenia będące owocem ich przemyśleń a dotyczących np. sposobów tworzenia (modelowania) dokumentów, strategii budowy schematów XML Schema czy stosowania polskich liter.

Książka może być z powodzeniem wykorzystywana jako podręcznik języka XML dla studentów kierunków informatycznych. Przeznaczona jest także dla programistów, projektantów oraz analityków, jak również zwykłych użytkowników komputerów, posiadających jednak pewne obycie informatyczne.

„XML na poważnie” zawiera wiele przykładów, a większość z nich jest dostępna na załączonej płycie CD-ROM w postaci prawie 300 plików.



# Spis treści

<b>Wprowadzenie .....</b>	<b>11</b>
<b>Rozdział 1. Geneza i stan obecny języka XML.....</b>	<b>19</b>
1.1. Początki języka znaczników. Historia XML.....	19
1.2. Ograniczenia języka HTML.....	21
1.3. XML — wzbogacony HTML czy zubożony SGML? .....	24
1.4. Języki znaczników.....	24
1.5. Standardy języka XML .....	25
1.6. Przyszłość języków systemu WWW .....	29
<b>Rozdział 2. XML — pojęcia podstawowe .....</b>	<b>31</b>
2.1. Dokument XML .....	31
2.2. Podstawowe składniki dokumentu XML .....	32
2.2.1. Elementy .....	32
2.2.2. Atrybuty .....	35
2.2.3. Elementy puste.....	36
2.3. Podelementy. Hierarchia elementów.....	38
2.4. Rodzaje elementów .....	40
2.5. Postać elementów .....	41
2.6. Deklaracja XML.....	43
2.7. Poprawny dokument XML.....	44
2.8. Przetwarzanie dokumentów XML .....	45
2.9. Inne składniki języka XML.....	46
2.9.1. Deklaracje .....	46
2.9.2. Instrukcje przetwarzania .....	49
2.9.3. Entity — jednostki .....	50
2.10. Nośniki informacji w dokumencie XML .....	52
2.11. Konwersja dokumentów HTML do XML .....	53
<b>Rozdział 3. DTD — opis struktury dokumentu.....</b>	<b>55</b>
3.1. Poprawność strukturalna dokumentu .....	56
3.2. Struktura DTD.....	58
3.2.1. DTD wewnętrzne .....	58
3.2.2. DTD zewnętrzne .....	58
3.3. Deklaracja elementu.....	61
3.3.1. Deklaracja elementu tekstowego .....	61
3.3.2. Deklaracja elementu pustego .....	62
3.3.3. Deklaracja elementu dowolnego.....	62
3.4. Deklarowanie podelementów — sekwencje i wybory.....	63
3.4.1. Sekwencja podelementów.....	63
3.4.2. Wybór podelementów .....	64
3.4.3. Połączenie sekwencji i wyboru.....	65

3.5. Wskaźniki liczby wystąpień .....	67
3.6. Modele zawartości elementu .....	71
3.7. Budowanie hierarchii elementów .....	73
3.8. Problemy zawartości elementowej .....	76
3.9. Deklaracje atrybutów .....	78
3.9.1. Typy atrybutów .....	78
3.9.2. Wartości domyślne atrybutów .....	81
3.10. Dylemat: podelementy czy atrybuty .....	83
3.11. Namespaces — przestrzenie nazw .....	85
3.11.1. Wykorzystanie wielu nazw z przestrzeni .....	87
3.11.2. Wiele przestrzeni. Przestrzeń domyślne .....	88
3.11.3. DTD a przestrzenie nazw .....	89
3.11.4. Przestrzenie standardowe .....	89
3.12. Atrybuty predefiniowane .....	91
3.12.1. Atrybuty xml:lang oraz xml:space .....	91
3.12.2. Atrybut xml:base .....	92
3.13. Jednostki (encje) .....	94
3.13.1. Jednostki parametryczne .....	95
3.13.2. Jednostki ogólne .....	98
3.14. Deklaracje notacji .....	102
3.15. Sekcje warunkowe .....	102
3.16. Problemy i ograniczenia DTD .....	104
3.17. Podsumowanie .....	105
<b>Rozdział 4. XML Schema — schematy dokumentów .....</b>	<b>107</b>
4.1. Budowa schematu .....	108
4.1.1. Przestrzeń nazw XML Schema. Korzeń schematu .....	108
4.1.2. Zasadnicze składowe schematu .....	109
4.2. Łączenie schematu z dokumentem XML .....	110
4.2.1. Przydzielanie schematów do przestrzeni nazw .....	110
4.2.2. Odwołania do schematów bez przestrzeni nazw .....	112
4.2.3. Odwołania do schematów z przestrzenią nazw .....	113
4.3. Sprawdzanie poprawności dokumentu XML .....	115
4.4. Typy proste a typy złożone .....	116
4.5. Wbudowane typy danych .....	119
4.6. Aspekty .....	125
4.6.1. Aspekty określające długość .....	127
4.6.2. Aspekty ograniczające wartości liczbowe .....	128
4.6.3. Wzorce wartości .....	129
4.6.4. Wyliczenia .....	134
4.6.5. Łączenie aspektów .....	135
4.6.6. Blokowanie wartości aspektów .....	135
4.7. Wymuszenia występowania .....	135
4.7.1. Liczby wystąpień elementu .....	135
4.7.2. Liczby wystąpień w deklaracjach DTD a schematy XML Schema .....	136
4.7.3. Występowanie atrybutu .....	137
4.8. Wyprowadzanie typów .....	137
4.9. Typy złożone zawierające podelementy. Składacze .....	140
4.10. Typy złożone z atrybutami — bez podelementów. Elementy puste .....	146
4.11. Deklaracje globalne i lokalne. Odesłania do elementów i atrybutów .....	148
4.12. Grupy elementów i atrybutów .....	151
4.13. Domyślne oraz stałe wartości elementów i atrybutów .....	154
4.14. Wartości niepowtarzalne .....	155

4.15. Listy .....	160
4.16. Kombinacje .....	161
4.17. Elementy zastępcze .....	162
4.18. Strategie budowy schematów .....	163
4.18.1. Metoda zagnieżdżenia .....	163
4.18.2. Metoda płaskiego katalogu .....	167
4.18.3. Metoda definiowania typów .....	171
4.18.4. Wnioski .....	175
4.19. XML Schema a DTD .....	175
4.20. Podsumowanie .....	179
<b>Rozdział 5. XPath — nawigacja i wyszukiwanie .....</b>	<b>181</b>
5.1. Budowa wyrażeń XPath .....	181
5.2. Wędrówka po drzewie dokumentu XML .....	182
5.2.1. Drzewo węzłów .....	183
5.2.2. Dostęp do węzłów drzewa .....	184
5.3. Wyszukiwanie węzłów .....	188
5.3.1. Operatory .....	189
5.3.2. Numer kolejnego węzła .....	190
5.3.3. Wyszukiwanie poprzez zawartość elementu .....	192
5.3.4. Wyszukiwanie poprzez wartość atrybutu .....	192
5.4. Funkcje .....	193
5.4.1. Operacje dotyczące tekstu .....	193
5.4.2. Operacje dotyczące liczb .....	196
5.4.3. Operacje dotyczące wartości logicznych .....	197
5.4.4. Pozostałe operacje .....	198
5.5. Ograniczenia języka XPath .....	200
<b>Rozdział 6. XLink — odsyłacze hipertekstowe .....</b>	<b>201</b>
6.1. Podstawy języka XLink .....	202
6.1.1. Deklaracja przestrzeni nazw .....	202
6.1.2. Rodzaje odsyłaczy .....	202
6.2. Odsyłacze proste .....	203
6.3. Podelementy odsyłaczy rozszerzonych .....	204
6.4. Atrybuty odsyłaczy XLink .....	205
6.5. Dwa sposoby wskazywania struktur XML .....	208
6.6. Łuki. Reguły przechodzenia .....	209
6.7. Odsyłacze rozszerzone. Bazy odsyłaczy .....	211
6.7.1. Prezentacja łuków .....	213
6.7.2. Wiązanie wielu zasobów za pomocą jednego łuku .....	214
6.7.3. Prezentacja rozbudowanych łuków .....	218
6.7.4. Zastosowanie baz odsyłaczy .....	219
6.8. Zachowanie odsyłaczy .....	220
6.9. Definiowanie struktur zawierających odsyłacze .....	221
6.10. Podsumowanie .....	224
<b>Rozdział 7. XPointer — wskazania wewnątrz dokumentów .....</b>	<b>225</b>
7.1. Rodzaje wskazań .....	225
7.2. Sekwencje dzieci .....	228
7.3. Funkcje XPointer .....	228
7.4. Punkty .....	229
7.4.1. Punkt początkowy .....	229
7.4.2. Punkt końcowy .....	230

7.5. Zakresy .....	231
7.5.1. Wskazanie zakresu poprzez wyszukanie ciągu znaków .....	231
7.5.2. Wskazanie zakresu poprzez podanie początku i końca .....	233
7.5.3. Zamiana elementu na zakres .....	234
7.6. Wskazania względne .....	234
7.7. Kodowanie znaków specjalnych .....	235
7.8. Wyrażenia XPath zawierające wyrażenia XPath .....	236
7.9. Podsumowanie .....	237
<b>Rozdział 8. XSLT — transformacje dokumentów .....</b>	<b>239</b>
8.1. Transformacje .....	241
8.2. Budowa dokumentu XSLT .....	242
8.3. Wzorce — definiowanie, wywoływanie .....	244
8.3.1. Definiowanie szablonów .....	244
8.3.2. Wywoływanie szablonów .....	247
8.4. Wydobywanie informacji z dokumentu wejściowego .....	250
8.5. Kontrola przebiegu transformacji .....	251
8.6. Zmienne, zbiory atrybutów .....	255
8.7. Elementy sterujące .....	258
8.8. Elementy kopiujące .....	259
8.9. Klucze .....	261
8.10. Przetwarzanie białych znaków .....	263
8.11. Numerowanie .....	263
8.12. Funkcje .....	270
8.13. Dołączanie arkuszy .....	273
8.14. Tworzenie nowego arkusza stylów .....	275
8.15. Rodzaje dokumentu wynikowego .....	277
8.16. Wykorzystanie kaskadowych arkuszy stylów CSS .....	279
8.17. Transformacje do XSL FO .....	283
8.18. Transformacje do dokumentu XML o innej strukturze .....	285
8.19. Przeglądarki internetowe a aplikacje serwerowe .....	287
<b>Rozdział 9. XSL FO — prezentacja danych .....</b>	<b>289</b>
9.1. XSLT a XSL FO .....	289
9.2. Budowa arkusza .....	290
9.2.1. Deklaracja przestrzeni nazw. Element główny .....	290
9.2.2. Struktura dokumentu XSL .....	291
9.3. Struktura strony .....	293
9.4. Szablony .....	296
9.5. Zawartość strony .....	301
9.6. Atrybuty prezentacji .....	303
9.7. Obszar bloku .....	310
9.7.1. Bloki .....	310
9.7.2. Tabele .....	310
9.7.3. Listy .....	314
9.8. Obiekty graficzne .....	316
9.9. Linie rozdzielające .....	318
9.10. FOP — transformacja do formatu PDF .....	319
9.11. Podsumowanie .....	321
<b>Rozdział 10. DOM — budowanie i dostęp do drzewa dokumentu XML .....</b>	<b>323</b>
10.1. Core — zasadnicza część DOM .....	324
10.1.1. DOMImplementation .....	324
10.1.2. Document .....	326

10.1.3. Node oraz NodeList .....	329
10.1.4. Element .....	332
10.1.5. Attr .....	333
10.1.6. CharacterData, Comment, Text .....	334
10.2. DOM StyleSheets .....	335
10.3. Events, Views, Traversal, Range .....	336
10.4. Podsumowanie .....	337
<b>Rozdział 11. SAX — przetwarzanie zdarzeniowe .....</b>	<b>339</b>
11.1. Zasada działania .....	340
11.2. Klasy i interfejsy SAX .....	341
11.2.1. XMLReader, XMLFilter .....	343
11.2.2. ContentHandler .....	346
11.2.3. Error handler, Locator .....	348
11.2.4. DTD Handler .....	349
11.2.5. Entity resolver, InputSource .....	350
11.2.6. Attributes .....	351
11.3. Klasy pomocnicze .....	352
11.4. Podsumowanie .....	353
<b>Rozdział 12. Zastosowanie języka XML .....</b>	<b>355</b>
12.1. Serwisy WWW .....	356
12.1.1. Warstwa treści i wyglądu — rozdzielenie informacji zawartych na stronach HTML .....	356
12.1.2. Transformacje dokonywane na serwerze .....	358
12.1.3. Transformacje dokonywane po stronie klienta .....	358
12.1.4. Strony WWW o złożonej strukturze .....	360
12.1.5. Wielokrotne wykorzystanie struktur XML. Różne sposoby prezentacji .....	361
12.2. Wymiana danych EDI .....	361
12.2.1. Koncepcja elektronicznej wymiany danych EDI .....	361
12.2.2. Wymiana danych z zastosowaniem języka XML a tradycyjna wymiana EDI .....	362
12.2.3. Obszary zastosowania języka XML w wymianie EDI .....	364
12.2.4. Bezpośrednia wymiana pomiędzy dwoma partnerami .....	365
12.2.5. Zastosowanie formatu wymiennego .....	367
12.2.6. Inicjatywy EDI XML .....	368
12.2.7. ebXML .....	368
12.2.8. BizTalk .....	370
12.2.9. RosettaNet .....	373
12.2.10. BASDA .....	374
12.2.11. Inne inicjatywy .....	375
12.3. Bazy danych .....	376
12.3.1. XDBMS — bazy danych w formacie XML .....	377
12.3.2. Tamino XML Server .....	379
12.3.3. XML jako interfejs pomiędzy bazą danych a serwisem WWW .....	381
12.3.4. XML a Microsoft SQL Server .....	382
12.3.5. XML a bazy Oracle .....	385
12.4. Pseudostrukturalne dane tekstowe .....	386
12.5. Systemy prezentacji oraz systemy interaktywne .....	390
12.5.1. SVG — grafika wektorowa .....	390
12.5.2. SMIL — integracja i synchronizacja mediów .....	392
12.5.3. VoiceXML — dźwiękowa prezentacja informacji .....	395
12.5.4. WML — prezentacje internetowe w telefonie komórkowym .....	396

12.6. Przechowywanie danych specjalistycznych.....	398
12.6.1. CML — opis struktur chemicznych.....	398
12.6.2. MathML — wzory matematyczne .....	399
12.6.3. MML — zapis muzyczny .....	401
12.6.4. NVML — nawigacja w terenie.....	401
12.6.5. Inne języki specjalistyczne .....	402
<b>Rozdział 13. System informacyjny „Plan zajęć uczelni” .....</b>	<b>405</b>
13.1. Struktura danych systemu (DTD) .....	405
13.2. Architektura systemu.....	406
13.3. SAX.....	407
13.3.1. Część prezentacyjna systemu.....	407
13.3.2. Opis techniczny.....	408
13.4. DOM.....	411
13.4.1. Część administracyjna systemu .....	411
13.4.2. Opis techniczny.....	412
13.5. Podsumowanie .....	414
<b>Dodatek A Procesory msxml .....</b>	<b>415</b>
<b>Dodatek B Atrybuty XSL FO .....</b>	<b>419</b>
<b>Dodatek C Słownik wyrażzeń angielskich .....</b>	<b>423</b>
<b>Dodatek D Słownik skrótów .....</b>	<b>429</b>
<b>Literatura .....</b>	<b>435</b>
<b>Skorowidz.....</b>	<b>445</b>

## Rozdział 2.

# XML — pojęcia podstawowe

XML, ponieważ jest językiem znaczników, służy do opisywania zawartości dokumentów elektronicznych w sposób zrozumiały zarówno dla komputerów, jak i dla ludzi. Jego zaletą jest to, że nawet osoba nie mająca z nim nic wspólnego będzie w stanie w znacznym stopniu rozszyfrować informację zakodowaną w dokumencie XML. Wynika to z prostoty i harmonii tego języka.

## 2.1. Dokument XML

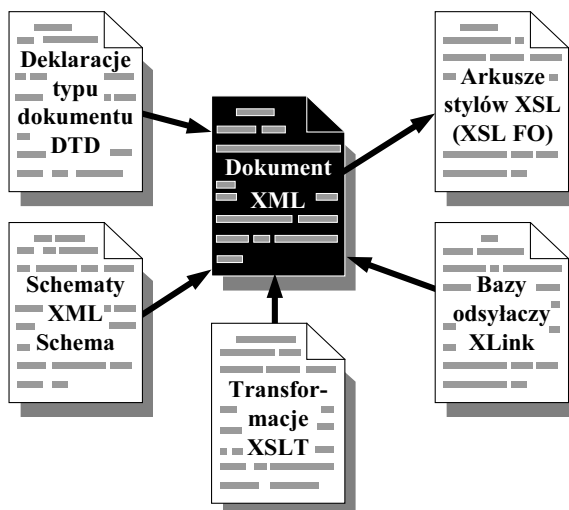
Aby rozgraniczyć niektóre pojęcia, ustalmy — dla potrzeb niniejszej książki — że **dokument XML** to dokument tekstowy (przechowywany zwykle w postaci pojedynczego pliku), w którym za pomocą elementów języka XML autor umieścił jakąś informację. Pozostałe dokumenty, np. DTD, schematy, arkusze stylów — jakkolwiek także wyrażone za pomocą języka XML — będą określane innymi terminami, np. „definicje typu dokumentu DTD”, „schematy XML Schema”, „arkusze stylów XSL” itd.

Definicje DTD, schematy, arkusze stylów, transformacje, a nawet bazy odsyłaczy są dodatkowymi specyfikacjami języka XML, odnoszącymi się do zasadniczej informacji zawartej w dokumencie XML — rysunek 2.1. Oczywiście często są to specyfikacje bardzo istotne dla konkretnych zastosowań. W dokumencie XML można więc wykorzystać definicje z DTD i schematów. Jednocześnie do dokumentu XML można zastosować transformacje, uzyskując np. postać gotową do wydruku — arkusze stylów. Specyfikacje XML przechowywane są zwykle w osobnych plikach, aczkolwiek mogą być także umieszczane wewnątrz dokumentu XML.

Ograniczenie znaczenia terminu „dokument XML” służy przede wszystkim łatwiejszemu zrozumieniu opisów, w których występuje więcej niż jedno pojęcie związane z językiem XML, np. „dokument XML zgodny z deklaracjami DTD” oznacza zbiór elementów XML (czyli dokument XML), który jest strukturalnie poprawny względem pewnego DTD.



**Rysunek 2.1.**  
Dokument XML



## 2.2. Podstawowe składniki dokumentu XML

Podstawowymi składnikami dokumentu XML są:

- ◆ **elementy**,
- ◆ **atrybuty**, które są umieszczane w elementach, jako dodatkowe informacje.

Elementy mogą być przy tym:

- ◆ **nie puste** — posiadające treść,
- ◆ **puste** — bez treści.

Opócz tego w dokumencie XML można umieszczać:

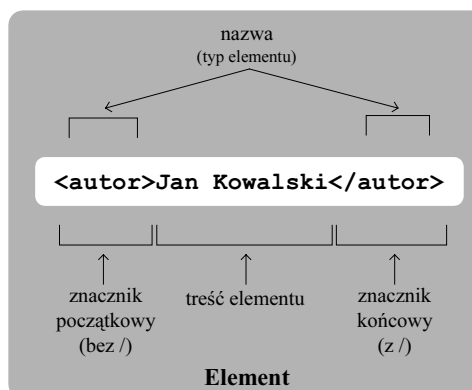
- ◆ **deklaracje**,
- ◆ **instrukcje przetwarzania**,
- ◆ **jednostki**.

Przyjrzyjmy się bliżej ww. składnikom dokumentów XML.

### 2.2.1. Elementy

Podstawowym budulcem dokumentu XML jest **element** (*element*). Posiada on zwykle tekstową **treść** poprzedzoną **znacznikiem początkowym** (*start-tag*) i zakończoną **znacznikiem końcowym** (*end-tag*) — rysunek 2.2. Oba znaczniki zawierają **nazwę elementu** objętą znakami „<” i „>”, przy czym w znaczniku końcowym nazwa jest zawsze poprzedzona znakiem „/”.

**Rysunek 2.2.**  
Składowe elementy



Jeżeli element nie jest pusty (posiada jakąś treść), to w języku XML znacznik początkowy i końcowy zawsze muszą wystąpić. Brak któregoś ze znaczników oznacza, że dany element — a co za tym idzie cały dokument — są niepoprawne.

W nazwach elementów małe i duże litery są rozróżniane. Tworząc dokument XML, należy zwrócić uwagę na to, by nazwa występująca w znaczniku początkowym była taka sama jak nazwa umieszczona w znaczniku końcowym.

Element o nazwie autor jest inny niż element o nazwie Autor, AUTOR lub aUtor. W praktyce w nazwach używa się zwykle małych liter. Nie jest to jednak regułą i są tacy, którzy preferują duże litery, np. [125].

Poniższy element nie jest poprawny, gdyż nazwa umieszczona w jego znaczniku początkowym jest inna niż w znaczniku końcowym:

```
<autor> Jan Kowalski </AUTOR>
```

Natomiast elementy

```
<autor> Anna Nowicka </autor>
<autor> Jan Kowalski </autor>
<AUTOR> trzeci </AUTOR>
```

są poprawne. Pierwsze dwa są tego samego rodzaju, zaś trzeci jest inny — ma inną nazwę, a co za tym idzie zwykle także inne znaczenie.



Nie poleca się jednak wykorzystywać małych i dużych liter w celu tworzenia elementów o nazwach tak samo brzmiących (jak wyżej). Powoduje to zwykle więcej zamieszania, niż daje pożytku.

W nazwach elementów mogą występować polskie znaki diakrytyczne (ą, Ą, ć, Ć, ...) <sup>1</sup>, jednak wymaga to odpowiedniego zadeklarowania i ścisłego przestrzegania kodowania znaków (patrz: punkt 2.6) we wszystkich dokumentach XML, a także dokumentach związanych, np. arkuszach stylów.

Niektóre programy przetwarzające mogą niepoprawnie interpretować polskie znaki w nazwach.

<sup>1</sup> Podobnie jak i inne znaki narodowe. Lista dozwolonych znaków w standardzie Unicode jest umieszczona w specyfikacji XML 1.0 pod adresem: <http://www.w3.org/TR/2000/REC-xml-20001006#NT-Letter>.

Dlatego też poprawnym elementem jest:

```
<ąĄćĆęĘłńŃóóśŚżŻź> Anna Nowicka </ąĄćĆęĘłńŃóóśŚżŻź>
```

Problemy z polskimi literami w nazwach elementów mogą wystąpić także w przypadku zastosowania dokumentów XML w kontaktach międzynarodowych. Kłopoty mogą również pojawić się przy wykorzystaniu dokumentów XML w połączeniu z bazami danych, które miewają różne standardy kodowania<sup>2</sup>. Z tego powodu w dużych projektach należałoby zalecić powstrzymanie się od polskich liter w nazwach elementów. W przypadku małych systemów, zwłaszcza wykorzystywanych przez zamknięte grono odbiorców, stosowanie polskich znaków może jednak zwiększyć czytelność dokumentów.



Nazwy elementów mogą zawierać cyfry, myślniki<sup>3</sup> i kropki, ale nie mogą się od nich zaczynać. Natomiast nazwy mogą się zaczynać od podkreśleń i dwukropków.

Poniższe elementy są poprawne<sup>4</sup>:

```
<Nazwa.1-2_3:ą> Treść </Nazwa.1-2_3:ą>
<_Nazwa.1-2_3:ą> Treść </_Nazwa.1-2_3:ą>
<:Nazwa.1-2_3:ą> Treść </:Nazwa.1-2_3:ą>
```

Natomiast:

```
<1Nazwa.1-2_3> Treść </1Nazwa.1-2_3> <!-- cyfra na początku -->
<.Nazwa.1-2_3> Treść </.Nazwa.1-2_3> <!-- kropka na początku -->
<_Nazwa.1-2_3> Treść </_Nazwa.1-2_3> <!-- myślnik na początku -->
```

są niezgodne ze standardem XML 1.0 [22].



Nazwa elementu niesie ze sobą bardzo ważną informację. Dzięki niej można określić znaczenie treści zawartej między znacznikami.

Nazwę można więc traktować jako określenie **typu elementu**. Z tego względu należy nadawać takie nazwy, aby były one zrozumiałe dla człowieka. Nie zaleca się korzystać z jakiegokolwiek sposobu ich kodowania, tak powszechnego w przypadku nazewnictwa funkcji i zmiennych w językach programowania, np. zamiast `fktr` lepiej użyć po prostu `faktura`. Oczywiście wielu programistów preferuje zapis zwięzły, co zwykle wynika z ograniczenia długości nazwy. W języku XML jednak nie ma żadnych ograniczeń długości nazwy, więc nie ma konieczności używania skrótów. Istotne jest tutaj także to, że dokumenty XML są najczęściej przeznaczone nie dla programistów, ale dla zwykłych użytkowników komputerów, którzy wolą czytelniejsze wersje nazw.



Autorzy nie powinni się bać nadawania elementom długich, opisowych nazw.

Treścią elementu jest tekst, tzn. nawet jeżeli będzie on miał postać liczby, to i tak można ją traktować jako znakowy ciąg cyfr.

<sup>2</sup> Problem polskich liter w nazwach jest także omawiany nieco dalej, w punkcie 2.6.

<sup>3</sup> Dotyczy to także podobnych znaków: półpauzy, pauzy oraz innych łączników (*Combining Characters*).

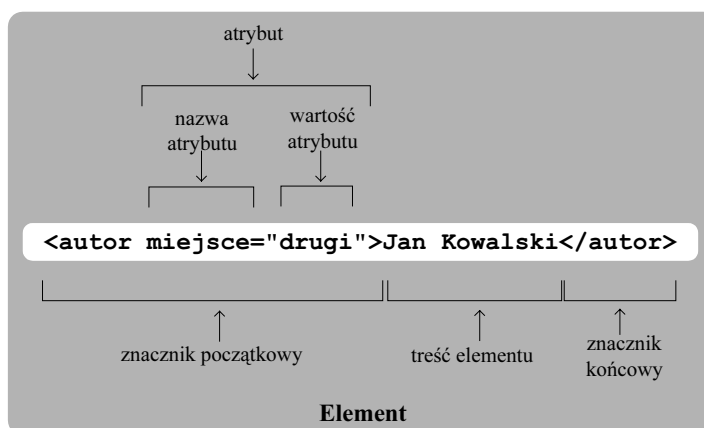
<sup>4</sup> Znak dwukropka rozdziela prefiks przestrzeni nazw od nazwy właściwej (patrz: punkt 3.11) i w związku z tym nie zaleca się używania go w innym celu. Przykładowo, parser XML zawarty w Microsoft Internet Explorerze nie dopuszcza stosowania znaku dwukropka na początku nazwy.

Ponieważ znaki „<” oraz „>” oznaczają początek i koniec znacznika, to nie mogą one występować w treści ujętej pomiędzy znacznikiem początkowym a końcowym (w treści elementu). Dlatego należy używać specjalnych symboli, tzw. jednostek wbudowanych (o których będzie mowa dalej), odpowiadających powyższemu znakom. I tak: znak „<” jest zastępowany przez „&lt;”, natomiast „>” przez „&gt;”.

## 2.2.2. Atrybuty

Autor dokumentu XML może przekazywać informacje nie tylko przez treść i nazwę elementu. Dodatkowe informacje można również umieszczać, w postaci **atrybutów**, w znaczniku początkowym elementów — rysunek 2.3. Atrybut zawsze posiada **nazwę** i poprzedzoną znakiem równości **wartość**.

**Rysunek 2.3.**  
*Element*  
zawierający atrybut



Atrybuty mogą występować jedynie w znaczniku początkowym lub w znaczniku elementu pustego<sup>5</sup>.

Podobnie jak w przypadku nazw elementów, małe i duże litery w nazwach atrybutów są rozróżniane.

W jednym elemencie nie mogą istnieć dwa atrybuty o takiej samej nazwie.

Dzięki umieszczeniu wartości atrybutu w cudzysłowie możliwe jest przechowywanie w nich dowolnych struktur tekstowych. Oczywiście — tak jak w przypadku treści elementów — ewentualne liczby traktowane są jako ciągi cyfr.

Zgodnie z powyższymi zasadami poniższy element jest poprawny:

```
<autor miejsce="pierwszy" Miejsce='drugi' MIEJSCE='trzeci'> Jan Kowalski </autor>
```

natomiast element

```
<autor miejsce="drugi" miejsce='trzeci'> Jan Kowalski </autor>
```

jest niepoprawny, gdyż posiada dwa atrybuty o tej samej nazwie.

<sup>5</sup> Elementy puste będą omówione w następnym podrozdziale.



W nazwach atrybutów również mogą występować polskie znaki diakrytyczne. Jednakże niektóre programy przetwarzające mogą je niepoprawnie interpretować<sup>6</sup>.

Dlatego też poprawnym atrybutem jest:

```
<żżżżąĆćĘęłńńÓóŚśŻżż="Dziwny atrybut"> Anna Nowicka </żżżż>
```



Wartości atrybutów muszą być zawsze umieszczone w cudzysłowach (") lub apostrofach (').

Element:

```
<autorzy liczba_autorow='2'> Anna Nowicka, Jan Kowalski </autorzy>
```

jest poprawny, natomiast elementy:

```
<autorzy liczba_autorow=2> Anna Nowicka, Jan Kowalski </autorzy>  
<autorzy liczba_autorow="2"> Anna Nowicka, Jan Kowalski </autorzy>
```

są niepoprawne, gdyż brak jest w nich takich samych znaków cudzysłowu lub apostrofów okalających wartość atrybutu.



Według standardu XML 1.0 kolejność atrybutów jest dowolna.

Może się jednak czasami zdarzyć, że program przetwarzający konkretny rodzaj dokumentów XML ma wbudowaną kolejność atrybutów lub ściśle przestrzega kolejności podanej w definicjach typu dokumentu DTD. Dobrym nawykiem jest więc zachowywanie tej samej kolejności atrybutów.

Podobnie jak znaki: „<” oraz „>” nie mogły występować w treści elementu, tak i tutaj zarówno apostrof ('), jak i cudzysłów (") nie powinny się pojawiać w polu wartości atrybutu. Dlatego należy zastąpić je odpowiednio przez „&quot;” oraz „&apos;”.



W praktyce wystarczy zwykle przestrzegać, aby wartość atrybutu zaczynała się i kończyła tym samym znakiem (apostrofu lub cudzysłowu), zaś w środku tej wartości może wystąpić drugi ze znaków (odpowiednio — cudzysłów lub apostrof).

Prawidłowy jest więc element:

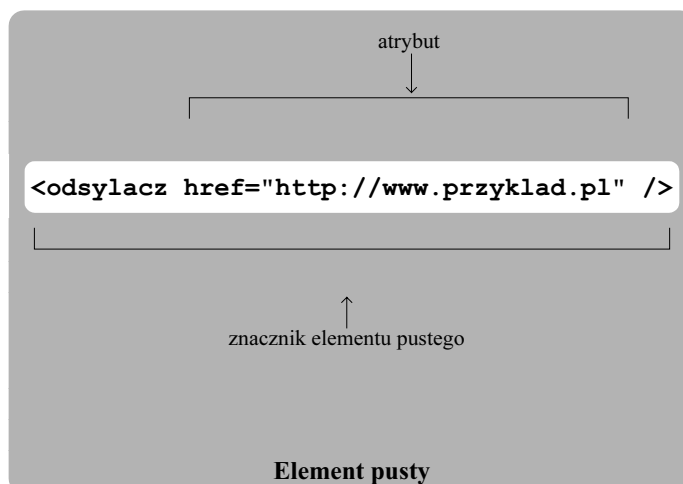
```
<autor pierwszy=' "Dziwny" "atrybut" ' drugi=' 'apostrof w środku' '> Anna Nowicka  
</autor>
```

### 2.2.3. Elementy puste

W języku XML elementy mogą nie posiadać treści. W takim przypadku można zamiast dwóch znaczników — początkowego i końcowego — użyć pojedynczego, innego znacznika: **znacznika elementu pustego** (*empty-element tag*). Składa się on z nazwy elementu, po której następuje znak „/” — rysunek 2.4.

<sup>6</sup>Ograniczenia nałożone na nazwy atrybutów są takie same jak przy nazwach elementów, np. nazwa nie może zaczynać się od cyfry ani od kropki.

**Rysunek 2.4.**  
*Element pusty  
 zawierający atrybut*



Znacznik elementu pustego zawsze kończy się znakiem ukośnika „/”, po którym dopiero następuje symbol „>”.

Oczywiście można także wykorzystać oba znaczniki: początkowy i końcowy i nie umieszczać między nimi żadnej treści.

Dlatego też oba poniższe przykłady są poprawne:

```
<faktura_usunieta/>
<faktura_usunieta></faktura_usunieta>
```

Naturalnie zalecaną formą jest ta pierwsza — krótsza.

Elementy puste zwykle (ale nie zawsze) zawierają atrybuty, np.:

```
<odsylnacz href='http://www.onet.pl' nazwa='Portal Onet' />
```

Zgodnie z powyższą uwagą można oczywiście także użyć:

```
<odsylnacz href='http://www.onet.pl' nazwa='Portal Onet'></odsylnacz>
```

Taka postać jest jednak rzadziej stosowana.



Elementy puste wykorzystuje się najczęściej wtedy, gdy zawierają one wyłącznie atrybuty.

Element pusty zazwyczaj składa się tylko ze znacznika elementu pustego, w którym można umieszczać atrybuty.

Elementy puste można także stosować wtedy, gdy potrzebne jest przekazanie dodatkowych informacji do programów przetwarzających dokument XML, np.:

```
<uruchom_mikrofon/>
```

## 2.3. Podelementy. Hierarchia elementów

Elementy mogą przechowywać tekst (pomiędzy znacznikiem początkowym i końcowym), mogą być puste, ale mogą także zawierać inne elementy — **podelementy** — dla których są wtedy **elementami nadrzędnymi**, np.:

```
<autorzy>
  <wspolautor> Anna Nowicka </wspolautor>
  <wspolautor> Jan Kowalski </wspolautor>
</autorzy>
```

Element `autorzy` zawiera w sobie dwa inne elementy (podelementy) — `wspolautor` — przy czym oba są tego samego rodzaju.

Oczywiście zagnieżdżanie może być bardziej złożone:

```
<sprawozdanie>
  <autorzy liczba_autorow="2">
    <wspolautor>
      <imie> Anna </imie>
      <nazwisko> Nowicka </nazwisko>
    </wspolautor>
    <wspolautor>
      <imie> Jan </imie>
      <nazwisko> Kowalski </nazwisko>
    </wspolautor>
  </autorzy>

  <tytul> Zestawienie zysków / obrotów </tytul>
</sprawozdanie>
```

Element `sprawozdanie` składa się tutaj z podelementu `autorzy` (także posiadającego podelementy), po którym następuje podelement `tytul`. Kolejność podelementów jest istotna i zamiana podelementu `autorzy` na podelement `tytul` da zupełnie inny dokument.



Od strony standardu języka XML nie ma żadnych ograniczeń na liczbę poziomów zagnieżdżeń.

Zagnieżdżając jedne elementy w drugich, umieszczamy w dokumencie XML informację o zależnościach występujących pomiędzy tym, co zostało zawarte w podelemencie a samym elementem, a także o zależnościach pomiędzy podelementami. Relacje te nazywane są **hierarchią elementów** — rysunek 2.5.



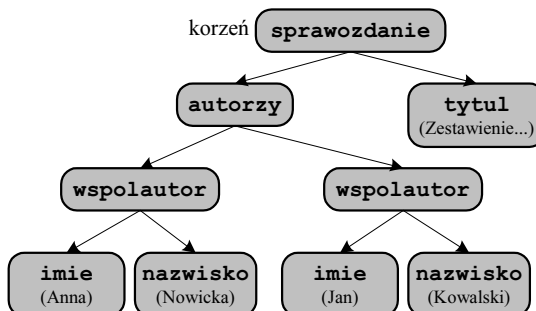
Podelement zawiera zwykle informację bardziej szczegółową niż element nadrzędny.

W związku z tym język XML nadaje się szczególnie dobrze do przechowywania informacji, w których występują relacje podległości (typu: nadrzędny — podrzędny), np. ogół — szczegół, całość — część, lista — element listy.

Nieco gorzej jest zbudować hierarchię dla relacji między obiektami o podobnym poziomie ogólności, np. klient — towar<sup>7</sup>, uczeń — przedmiot<sup>8</sup>, mąż — żona<sup>9</sup>.

### Rysunek 2.5.

Hierarchia elementów



Zbudowanie odpowiedniej hierarchii elementów jest często najtrudniejszym problemem występującym przy tworzeniu dokumentów XML.

Zagadnienie budowy właściwej hierarchii jest zwykle rozważane podczas definiowania typu dokumentu DTD (rozdział 3.) lub schematu dokumentu (rozdział 4.).

W powyższym przykładzie widać, że nazwisko umieszczone jest na tym samym poziomie co imię, natomiast elementem najwyższym, głównym, czyli **korzeniem**, jest sprawozdanie.



W dokumentach XML może istnieć tylko jeden korzeń — element główny.

Cały dokument XML jest więc pojedynczą hierarchią z jednym korzeniem.

Elementy o takiej samej nazwie — nie będące korzeniem — mogą występować w różnych miejscach dokumentu, na różnych poziomach i mogą zawierać różnego rodzaju treść:

```

<dane_firmowe>
  <nazwa> Producent Jogurtów „Super-Zdrowie” S.A. </nazwa>

  <dostawa>
    <odbiorca>
      <nazwa> Hurtownia Spożywcza „Milk-XML” Sp. z o.o. </nazwa>
      <adres> ul. Niecodzienna 17, Zielona Góra </adres>
    </odbiorca>
  </dostawa>
</dane_firmowe>
  
```

<sup>7</sup> Relacja pomiędzy klientem a towarem wyraża się zwykle poprzez zakup lub sprzedaż. Jak jednak zbudować relację hierarchii dla dokumentu XML przechowującego informacje o miesięcznych zakupach jakiejś firmy (dane zakupu, dane klienta i informacje o towarze)? Dane o dostawcach lub towarach powielająby się. Rozwiązaniem tutaj mogłyby się okazać odsyłacze (patrz rozdział 6.). Tradycyjne relacyjne bazy danych lepiej się sprawdzają w podobnych przypadkach.

<sup>8</sup> Uczeń jest w relacji z przedmiotem poprzez klasę, do której uczęszcza. Rozwiązanie podobnego problemu za pomocą odsyłaczy standardu XLink umieszczono w rozdziale 6.

<sup>9</sup> W relacji mąż – żona zwykle bezpieczniej nie określać stopnia podległości.



```

<towar>
  <nazwa> Jogurt naturalny </nazwa>
  <magazyn>
    <nazwa> 5 </nazwa>
  </magazyn>
  <opakowanie>
    <nazwa> paleta </nazwa>
    <jednostka_detal> szt. </jednostka_detal>
  </opakowanie>
</towar>
</dostawa>
</dane_firmowe>

```

Element nazwa został umieszczony w kilku różnych kontekstach (nazwa firmy, nazwa dostawcy, nazwa towaru, nazwa magazynu, nazwa jednostki miary) i posiada różne rodzaje danych (tekstowe — nazwy własne firm, numeryczne — nazwa magazynu, pojedyncze słowa — nazwa jednostki miary opakowania).



Elementy mogą być zagnieżdżone, ale nie mogą na siebie zachodzić, tzn. znacznik końcowy podelementu musi wystąpić przez znacznikiem końcowym elementu nadrzędnego. Innymi słowy, każdy element (oprócz korzenia) musi być w całości otoczony przez inny element.

W związku z tym, poniższa postać jest nieprawidłowa:

```
<imie> Jan <nazwisko> Kowalski </imie></nazwisko>
```

Zamiast tego należy użyć:

```
<imie> Jan <nazwisko> Kowalski </nazwisko></imie>
```

lub:

```
<nazwisko> Kowalski <imie> Jan </imie></nazwisko>
```

Ze względu na równorzędność imienia i nazwiska najlepszym rozwiązaniem jest tutaj:

```
<imie> Jan </imie>
<nazwisko> Kowalski </nazwisko>
```

## 2.4. Rodzaje elementów

W przypadku, gdy element zawiera w sobie jakieś dane (tekst lub podelementy), jest nazywany pojemnikiem (*container element*).

Element — pojemnik może mieć następujące rodzaje zawartości:

- 1. Zawartość elementową** (*element content*) — gdy pojemnik zawiera tylko i wyłącznie inne elementy, np.:

```

<autor>
  <opis> To jest współautor sprawozdania </opis>
  <imie> Jan </imie>
  <nazwisko> Kowalski </nazwisko>
</autor>

```

**2. Zawartość mieszana** (*mixed content*) — gdy pojemnik zawiera zarówno tekst, jak i inne elementy, np.:

```
<autor>
  To jest współautor sprawozdania
  <imie> Jan </imie>
  <nazwisko> Kowalski </nazwisko>
</autor>
```

**3. Zawartość tekstową** (*data content*) — gdy pojemnik zawiera tylko i wyłącznie tekst, np.:

```
<autor> To jest współautor sprawozdania: Jan Kowalski </autor>
```

W przypadku zawartości mieszanej elementy mogą być nawet wplecione w tekst:

```
<sprawozdanie>
  Autorami sprawozdania są: <autor> Anna Nowicka </autor> oraz <autor> Jan Kowalski
  </autor>. Innych autorów nie ma.
</sprawozdanie>
```

Zawartość mieszana nie jest jednak godna polecenia, gdyż może prowadzić do pewnego zamętu. Można ją stosować dla zwykłych tekstów, w których chcemy jedynie wydzielić pewne szczególne informacje. Jest to przydatne na etapie przejściowym, w którym tekst napisany w języku naturalnym przekształcamy w formy bardziej ustrukturalizowanej.

Dla ww. przypadku lepszym rozwiązaniem może być zawartość elementowa z dodatkowymi elementami, czyli z redundancją (powieleniem) informacji:

```
<sprawozdanie>
  <autor> Anna Nowicka </autor>
  <autor> Jan Kowalski </autor>
  <opis> Autorami sprawozdania są: Anna Nowicka oraz Jan Kowalski.
  Innych autorów nie ma. </opis>
</sprawozdanie>
```

lub (rozwiązanie gorsze):

```
<sprawozdanie>
  <opis>Autorami sprawozdania są: </opis>
  <autor> Anna Nowicka </autor>
  <opis> oraz </opis>
  <autor> Jan Kowalski </autor>
  <opis>. Innych autorów nie ma. </opis>
</sprawozdanie>
```

## 2.5. Postać elementów

Tworząc dokumenty XML, należy zwrócić uwagę na sposób zapisu dokumentu, czyli na postać dokumentu. Zapis:

```
<sprawozdanie><autorzy liczba_autorow="2"><autor><imie> Anna </imie> <nazwisko>
Nowicka </nazwisko></autor><autor><imie> Jan </imie> <nazwisko> Kowalski
</nazwisko></autor></autorzy><okres> 7 miesięcy </okres></sprawozdanie>
```

jest zupełnie nieczytelny! Poniższy zapis jest lepszy:

```
<sprawozdanie>
<autorzy liczba_autorow="2">
<autor>
<imie> Anna </imie>
<nazwisko> Nowicka </nazwisko>
</autor>
<autor>
<imie> Jan </imie>
<nazwisko> Kowalski </nazwisko>
</autor>
</autorzy>
<okres> 7 miesięcy </okres>
</sprawozdanie>
```

Brakuje mu jednak odpowiednich wcięć uwydatniających hierarchiczne zależności pomiędzy elementami. Jeszcze lepszą postacią zapisu jest więc:

```
<sprawozdanie>
  <autorzy liczba_autorow="2">
    <autor>
      <imie> Anna </imie>
      <nazwisko> Nowicka </nazwisko>
    </autor>
    <autor>
      <imie> Jan </imie>
      <nazwisko> Kowalski </nazwisko>
    </autor>
  </autorzy>
  <okres> 7 miesięcy </okres>
</sprawozdanie>
```

Zauważmy, że niektóre elementy zawarte są w jednej linii, zaś w innych znacznik początkowy i końcowy umieszczono w osobnych liniach. Można oczywiście znaczniki wszystkich elementów umieścić w osobnych liniach:

```
<sprawozdanie>
  <autorzy liczba_autorow="2">
    <autor>
      <imie>
        Anna
      </imie>
      <nazwisko>
        Nowicka
      </nazwisko>
    </autor>
    <autor>
      <imie>
        Jan
      </imie>
      <nazwisko>
        Kowalski
      </nazwisko>
    </autor>
  </autorzy>
  <okres>
    7 miesięcy
  </okres>
</sprawozdanie>
```

Wydaje się jednak, że to ostatnie rozwiązanie nie tyle poprawia czytelność dokumentu, ile zwiększa jego objętość, a co za tym idzie, raczej zmniejsza jego przejrzystość.

## 2.6. Deklaracja XML

Jak już wspomniano, podstawowym składnikiem dokumentu XML są elementy. To one zwykle zawierają zasadniczą informację, którą chcemy umieścić w dokumencie. Aby jednak dany dokument mógł być identyfikowany jako dokument XML, należy na jego początku umieścić **deklarację XML**. W deklaracji tej należy:

- ♦ określić wersję języka XML użytego w dokumencie — parametr `version`,
- ♦ określić rodzaj kodowania znaków — parametr `encoding`,
- ♦ podać informację o pominięciu (lub nie) przetwarzania zewnętrznych definicji typu dokumentu (DTD) oraz zewnętrznych jednostek — parametr `standalone` (informacja ta mówi o tym, czy dokument jest „samodzielny”).

Deklaracja XML zaczyna się znakami `<?xml` i kończy sekwencją `?>`. Oto przykład deklaracji XML:

```
<?xml version="1.0" encoding="UTF-8" standalone="yes" ?>
```

Parametr `version` jest obowiązkowy i zawsze powinien mieć wartość 1.0 (na razie istnieje tylko jedna wersja standardu XML). Pozostałe parametry (`encoding` oraz `standalone`) deklaracji są opcjonalne (mogą nie wystąpić).



Zalecanym rodzajem kodowania znaków w języku XML jest standard Unicode.

Standard Unicode jest obsługiwany przez większość systemów operacyjnych<sup>10</sup>. Każdy program przetwarzający dokumenty XML powinien potrafić odpowiednio interpretować kody Unicode w standardzie UTF-8 i UTF-16<sup>11</sup>. W UTF-8 wszystkie znaki zakodowane są w zmiennej liczbie bajtów (od 1 do nawet 6 bajtów), przy czym pierwsze 127 kodów ASCII (m.in. alfabet angielski, cyfry) zajmują tylko jeden bajt [71]. Polskie znaki są mniej uprzywilejowane (zajmują 2 bajty), zaś ich kody można znaleźć w [101].

Deklaracje XML wskazujące na kodowanie wg standardu Unicode wyglądają następująco:

```
<?xml version="1.0" encoding="UTF-8" ?>  
<?xml version="1.0" encoding="UTF-16" ?>
```

Znaki mogą być także zakodowane za pomocą jednego ze standardów ISO, np. Latin 1 (bez polskich znaków diakrytycznych):

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

<sup>10</sup> Obsługa ta nie zawsze jest pełna. Systemy często mają zaimplementowany ograniczony zestaw znaków Unicode, np. domyślne fonty Unicode w Windows 98 zawierają tylko 1 znak azjatycki.

<sup>11</sup> Wymóg ten został zapisany w standardzie XML 1.0 [22].

lub — jeżeli chcemy używać polskich znaków — Latin 2 (ISO 8859-2):

```
<?xml version="1.0" encoding="ISO-8859-2"?>
```

Dostępne są także standardy proponowane przez firmę Microsoft, znane z systemu Windows, np. dla polskich znaków może to być Windows-1250<sup>12</sup>:

```
<?xml version="1.0" encoding="windows-1250" ?>
```

Rodzaj kodowania zależy od tego, czy programy przetwarzające dokument będą rozumiały dany rodzaj kodowania.



Nie można stosować sposobu kodowania innego dla nazw elementów i atrybutów, a innego dla ich treści i wartości.

Możliwość stosowania tylko jednego kodowania może być pewnym ograniczeniem w przypadku zastosowania języka XML w połączeniu z bazami danych, w których często stosuje się różne standardy kodowania<sup>13</sup>. Dodatkowo nazwy elementów są często wykorzystywane w dokumentach związanych, np. w definicjach typu dokumentu DTD, schematach XML Schema, transformacjach XSLT czy wskazaniach XPath. Ewentualna zmiana kodowania polskich znaków (lub niezgodność z kodowaniem w dokumentach XML) spowoduje konieczność modyfikacji tych powiązanych dokumentów. Z tych powodów zwykle lepiej nie stosować polskich znaków diakrytycznych w nazwach elementów i atrybutów.

Przykładowy poprawny dokument XML, zawierający deklarację XML, może wyglądać następująco (plik *Sprawozdanie1.xml*):



```
<?xml version="1.0" encoding="ISO-8859-2"?>
<sprawozdanie>
  <autorzy liczba_autorow="2">
    <wspolautor> Anna Nowicka </wspolautor>
    <wspolautor> Jan Kowalski </wspolautor>
  </autorzy>
</sprawozdanie>
```

## 2.7. Poprawny dokument XML

Standard języka XML, jak już można było zauważyć, wyznacza pewne ograniczenia dotyczące zapisu elementu oraz relacji pomiędzy elementami. Przypomnijmy te ograniczenia:

1. Dokument musi się zaczynać od deklaracji XML, aby mógł być traktowany jako dokument XML.

<sup>12</sup> Takie kodowanie może jednak nie być akceptowane przez parsery inne niż firmy Microsoft.

<sup>13</sup> Wyobraźmy sobie sytuację, w której nazwy elementów w standardowej fakturze zawierają polskie znaki w standardzie Unicode, natomiast dane w bazach danych, niezbędne do wygenerowania ostatecznej wersji dokumentu przechowywane są w standardzie ISO Latin-2. Wówczas, podczas procesu wymiany faktury pomiędzy jedną a drugą bazą, niepotrzebnie należałoby dwukrotnie przekodowywać dane: najpierw ze standardu ISO Latin-2 do Unicode, a następnie odwrotnie.

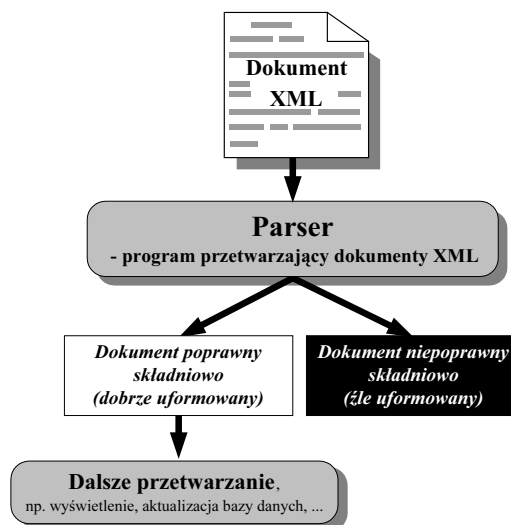
2. W dokumencie może istnieć tylko jeden unikatowy element główny — korzeń. Element o nazwie takiej jak korzeń nie może wystąpić jako podelement w tym dokumencie. W związku z tym wszystkie elementy występujące w dokumencie muszą się w całości zawierać wewnątrz korzenia.
3. Wszystkie niepuste elementy muszą posiadać znacznik początkowy i końcowy.
4. Element pusty może składać się z pojedynczego znacznika elementu pustego.
5. Elementy mogą być zagnieżdżone, jednak nie mogą na siebie zachodzić.
6. Wartości atrybutów muszą być umieszczone w cudzysłowach (lub apostrofach).

Przestrzeganie wymienionych wyżej reguł prowadzi do **dobrze uformowanego** (*well-formed*), czyli **poprawnego składniowo**, dokumentu XML. Większość programów przetwarzających dokumenty XML umożliwia sprawdzenie, czy dany dokument XML jest dobrze uformowany.

## 2.8. Przetwarzanie dokumentów XML

Dokumenty XML mogą być czytane przez użytkowników w zwykłych edytorach tekstowych, jednak najczęściej są one przetwarzane przez programy, zwane **parserami**. Programy te powinny w pierwszej kolejności sprawdzić poprawność składniową dokumentu, tzn. stwierdzić czy jest on dobrze uformowany — rysunek 2.6.

**Rysunek 2.6.**  
Przetwarzanie dokumentów XML



Następnie mogą one wykorzystać informację, zawartą w dokumencie, w jakimś specyficznym celu, np.:

- ♦ wyświetlić ją użytkownikowi w określonym formacie;
- ♦ zmienić postać informacji, np. przekształcając dokument XML na stronę w języku HTML;

- ♦ wygenerować nowy dokument XML i przekazać go innemu programowi;
- ♦ przetworzyć informację do postaci zgodnej z określonym formatem bazy danych i zaktualizować bazę danych.

Do dalszego przetwarzania dokumentu XML zwykle niezbędne są dodatkowe informacje, np. arkusze stylów, transformacje, definicje typu dokumentu DTD itd.

## 2.9. Inne składniki języka XML

Oprócz zwykłych elementów, zawierających informacje przydatne wprost dla użytkownika, dokument XML może zawierać także instrukcje dla programu przetwarzającego — parsera. Są to:

- ♦ deklaracje (w tym m.in. komentarze i bloki tekstu),
- ♦ instrukcje przetwarzania,
- ♦ jednostki.

### 2.9.1. Deklaracje

Deklaracje umieszczane są w tzw. **znacznikach deklaracji** (*markup declaration*). Różnią się nieznacznie od zwykłych znaczników. Otaczane są znakami „<!” oraz „>”.

Na początku deklaracji musi być umieszczona nazwa jednoznacznie identyfikująca daną deklarację:

```
<!NAZWA_DEKLARACJI ... >
```

Znaczniki deklaracji wykorzystywane są także do grupowania kilku innych deklaracji zamkniętych w zbiór przy użyciu znaków „[” oraz „]”, przy czym niektórych rodzajów deklaracji nie można grupować:

```
<!zbiór_deklaracji [  
  <!deklaracja1>  
  <!deklaracja2>  
>
```

Standard XML 1.0 określa następujące deklaracje:

- ♦ deklaracje podstawowe (umieszczane w dokumentach XML),
- ♦ deklaracje wykorzystywane w definicjach typu dokumentu DTD.

Deklaracje podstawowe to:

- ♦ komentarze,
- ♦ sekcje CDATA.

### 2.9.1.1. Komentarze

Istnieje możliwość umieszczania w dokumencie XML **komentarzy** (*comments*). Zawierają one zwykle informację w języku naturalnym zrozumiałą jedynie dla człowieka i dlatego są przez programy przetwarzające dokument (parsery) pomijane. Komentarz zaczyna się czterema znakami „<!--”, a kończy trzema „-->”.

```
<!-- To są elementy opisujące autorów sprawozdania -->
<autor>Anna Nowicka</autor> <!-- Autor pierwszy, co nie oznacza, że najważniejszy -->
<autor>Jan Kowalski</autor> <!-- Autor drugi, wcale nie gorszy od pierwszego -->
```



Komentarze można umieszczać w dowolnym miejscu dokumentu.

Poprawny komentarz zawiera tylko dwie pary myślników występujących w części początkowej i końcowej. Nie mogą więc one wystąpić w treści zasadniczej komentarza:

```
<!-- To nie jest poprawny -- komentarz, bo ma w treści 2 myślniki -->
<!-- To jest poprawny komentarz - nie ma w treści dwóch myślników -->
```

Komentarze najczęściej występują w różnych specyfikacjach związanych z dokumentem XML (patrz rysunek 2.1), np. w definicjach typu dokumentu DTD czy schematach, gdyż tam istnieje zwykle największa potrzeba bliższego wyjaśnienia znaczenia elementów i relacji między nimi.

### 2.9.1.2. Sekcje CDATA — bloki tekstu

Jak już wcześniej zostało wspomniane, nie wszystkie znaki (symbole) mogą być zawarte w treści elementu, gdyż mają szczególne znaczenie dla przetwarzania dokumentu. W przypadku, gdy wymagane jest użycie znaków „<”, „>” czy „&”, zamiast używania kodów odpowiadających tymże znakom (&lt;, &gt;, &amp;) można wykorzystać tzw. **blok tekstu**, czyli **sekcję CDATA**. Pozwala ona na podanie dowolnego tekstu, zaś wszystkie znaki ujęte wewnątrz deklaracji CDATA będą potraktowane jako zwykły tekst — dosłownie. Sekcja CDATA ma postać:

```
<![CDATA[tekst]]>
```

Załóżmy, że dokument XML zawiera w treściach elementów komunikaty wyświetlane przez pewien system SuperXML, w tym m.in. komunikat:

Naciśnij <<< Ctrl-Enter >>> aby pójść dalej

Odpowiedni dokument XML może mieć postać (plik *System1.xml*):



```
<?xml version="1.0" ?>
<system>
  <nazwa_systemu> SuperXML wersja 2.0 </nazwa_systemu>
  <komunikat nr_komunikatu="102">
    Naciśnij &lt;&lt;&lt; Ctrl - Enter &gt;&gt;&gt; aby przejść dalej
  </komunikat>
</system>
```

Zamiast symbolu < oraz > konieczne było zastosowanie odpowiednio &lt; i &gt;. Można jednak zastosować tutaj sekcję CDATA (plik *System2.xml*):





```
<?xml version="1.0" ?>
<system>
  <nazwa_systemu> SuperXML wersja 2.0 </nazwa_systemu>
  <komunikat nr_komunikatu="102">
    <![CDATA[Naciśnij <<< Ctrl-Enter >>> aby przejść dalej]]>
  </komunikat>
</system>
```

Bloki tekstu są przydatne także wtedy, gdy w treści elementu (lub jako wartość atrybutu) chcemy umieścić fragment dokumentu XML w postaci zwykłego tekstu. Przykładowo, aby umieścić tekst:

```
Złe zagnieżdżenie: <nazwisko> <imie></nazwisko></imie>
```

w treści elementu przykład, należy napisać:

```
<przyklad>
  <![CDATA[Złe zagnieżdżenie: <nazwisko> <imie></nazwisko></imie>]] >
</przyklad>
```

Tekst zawarty w sekcji CDATA będzie traktowany przez parser jako ciąg znaków, a nie jako dwa niepoprawne elementy.

### 2.9.1.3. Deklaracje wykorzystywane w DTD

Istnieje kilka deklaracji wykorzystywanych w definicjach typu dokumentu DTD. Będą one szczegółowo omówione w rozdziale 3. W tym miejscu zostaną tylko wymienione — tabela 2.1.

**Tabela 2.1.** Deklaracje wykorzystywane w DTD

Deklaracja	Znaczenie deklaracji
<!DOCTYPE ..... >	Określa <b>typ dokumentu</b> i zawiera lub wskazuje definicję typu dokumentu DTD
<!ENTITY ..... >	Definiuje <b>jednostki</b> , dzięki którym można jeden raz zdefiniować pewien fragment dokumentu, a potem wielokrotnie się do niego odwoływać
<!NOTATION ..... >	Definiuje <b>notacje</b> , poprzez które można określać typy danych zewnętrznych
<!ELEMENT ..... >	Deklaruje element (opisuje zawartość elementu)
<!ATTLIST ..... >	Deklaruje listę atrybutów elementu
<![IGNORE [ ..... ] >	Umożliwia odpowiednio wyłączenie i włączenie do przetwarzania pewnych fragmentów dokumentu
<![INCLUDE [ ..... ] >	

Deklaracja typu dokumentu DOCTYPE nie zawsze musi być związana z DTD. Może ona być także wykorzystywana jedynie do nazwania dokumentu XML, np. deklaracja:

```
<?xml version="1.0" ?>
<!DOCTYPE sprawozdanie>
<sprawozdanie>
  ...
</sprawozdanie>
```

umieszczona na początku dokumentu (po deklaracji XML) nadaje nazwę sprawozdanie całemu dokumentowi. Jest to jednocześnie nazwa elementu głównego — korzenia.

Takie zastosowanie deklaracji DOCTYPE nie ma jednak większego praktycznego sensu i dlatego też nie przez każdy parser deklaracja taka jest obsługiwana.

## 2.9.2. Instrukcje przetwarzania

Instrukcje przetwarzania (*processing instruction* — *PI*) zawierają informacje potrzebne dla programów przetwarzających dokument XML. Dla odróżnienia od zwykłego elementu instrukcje przetwarzania są umieszczane w bloku ograniczonym znakami „<?” i „?>”. Instrukcje te zawierają **cel** (*PITarget*) oraz **instrukcję do wykonania**. Celem jest słowo kluczowe identyfikujące aplikację, dla której przeznaczona jest instrukcja do wykonania. Jeżeli program przetwarzający dokument poprawnie rozpozna cel, zadana instrukcja zostanie przez niego wykonana:

```
<!-- instrukcje dla robotów internetowych: zaindeksować i nie   przechodzić dalej
      po odsyłaczach -->
<?robots index="yes" follow="no" ?>

<!-- Instrukcje dla edytorów tekstowych -->
<?DTPSystem DO:page-break ?> <!--Dla programu DTP: wstawić stronę-->

<?EdytorTXT {New Page} ?>    <!-- Edytor tekstowy ma dodać stronę-->

<!-- Jeżeli programem przetwarzającym dokument XML będzie serwer php, to powinien
      on odpowiednio zinterpretować poniższy fragment programu -->
<?php
mysql_connect("baza.pwr.wroc.pl", "kowalski", "password");
$result = mysql("STUDENCI", "SELECT Nazwisko, Imie FROM Studenci
ORDER BY Nazwisko, Imie");
$i = 0;
while ($i < mysql_numrows ($result)) {
    $fields = mysql_fetch_row($result);
    echo "<student>$fields[1] $fields[0] </student>\r\n";
    $i++;
}
mysql_close();
?>
```

Jeżeli cel nie jest znany aplikacji przetwarzającej dokument, wtedy cała instrukcja zostanie przez nią pominięta.

Pewnym ograniczeniem jest to, że dwuznak „?>” nie może być zawarty w instrukcji do wykonania.



Instrukcje przetwarzania umożliwiają umieszczanie w dokumentach XML fragmentów innych języków, zwłaszcza języków programowania.

Gdy jako cel podana zostanie nazwa `<xml`, wtedy otrzymujemy deklarację XML opisaną w punkcie 2.6. Nie jest to typowa instrukcja przetwarzania, zaś ciąg `<xml` (w dowolnej kombinacji małych i dużych liter) jest zarezerwowany i nie należy używać go jako celu we własnych instrukcjach przetwarzania.

Instrukcje przetwarzania mogą być:

- ♦ **standardowe**, zdefiniowane w którymś ze standardów związanych z językiem XML, np. `<?xml ...?>`, `<?xml-stylesheet ... ?>` — wykorzystywane w języku XSLT;
- ♦ **niestandardowe** — niezwiązane ze standardami języka XML, czyli utworzone dla specjalistycznych potrzeb, np. dla skryptów PHP, konkretnego edytora tekstów czy dla potrzeb przetwarzania w procesie wymiany danych między firmą X a firmą Y.

### 2.9.3. Entity — jednostki

Język XML umożliwia fizyczną separację części przechowywanego dokumentu. Przykładowo, gdyby dokument XML opisywał książkę, to jej rozdziały można by umieścić w osobnych plikach. Umożliwia to sprawniejszą kontrolę, prostszą edycję czy nawet krótszy czas ładowania przez Internet. Każda taka osobna część jest nazywana **jednostką** lub **encją** (*entity*). Każda jednostka ma swoją nazwę, przez którą jest identyfikowana. Jednostka jest definiowana słowem kluczowym ENTITY, które wchodzi w skład **deklaracji jednostki** (*entity declaration*), a ta z kolei może być wstawiona jedynie w opisanej w następnym rozdziale deklaracji typów dokumentów (DTD). Wykorzystanie jednostki w dokumencie polega na umieszczeniu w żądanym miejscu odesłania do danej jednostki. Odesłanie składa się z nazwy jednostki poprzedzonej znakiem „&” oraz zakończonej znakiem średnika. W dokumencie może wystąpić wiele odesłań do tej samej jednostki.

Najczęściej jednostki stosuje się, gdy:

- ♦ ta sama informacja jest używana wielokrotnie;
- ♦ informacja jest częścią dużego dokumentu i z powodów praktycznych jest podzielona na łatwe w zarządzaniu części;
- ♦ informacja zawiera dane nie przetwarzane przez XML, lecz przez inną aplikację.

Najprostszą formą jednostki jest **jednostka wewnętrzna tekstowa** (*internal text entity*), zwana także **ogólną jednostką wewnętrzną** (*general internal entity*). Jej deklaracja znajduje się w przetwarzanym dokumencie XML; zawiera tylko i wyłącznie tekst. Stosuje się ją w przypadku częstego powtarzania fragmentu tekstu. Następująca deklaracja:

```
<!ENTITY XML "eXtensible Markup Language" >
```

oznacza, że nazwie XML przyporządkowano tekst „*eXtensible Markup Language*”, zaś każde odwołanie do tej nazwy (czyli `&XML;`) będzie zastępowane poprzez ten tekst.

Przykładowo odwołanie:

```
<opis> To jest tekst o XML czyli o &XML;. Angielski termin "&XML;" oznacza
    rozszerzalny język znaczników</opis>
```

jest równoznaczne z:

```
<opis> To jest tekst o XML czyli o eXtensible Markup Language. Angielski termin
    "eXtensible Markup Language" oznacza rozszerzalny język znaczników</opis>
```

Deklaracje jednostek powinny być umieszczane w deklaracji typu dokumentu DOCTYPE. Cały dokument zawierający powyższe opisy ma więc postać (plik *Jednostka.xml*):



```
<?xml version="1.0"?>
<!DOCTYPE opisy [
  <!ENTITY XML "eXtensible Markup Language">
]>

<opisy>

  <!-- oba opisy mają taka samą treść -->

  <opis> To jest tekst o XML czyli o &XML;. Angielski termin "&XML;"
  oznacza rozszerzalny język znaczników </opis>

  <opis> To jest tekst o XML czyli o eXtensible Markup Language.
  Angielski termin "eXtensible Markup Language" oznacza rozszerzalny
  język znaczników </opis>

</opisy>
```

Kolejnym typem jednostek jest **jednostka zewnętrzna tekstowa** (*external text entity*) zwana także **ogólną jednostką zewnętrzną** (*general external entity*). Zdefiniowanie jednostki zewnętrznej różni się od definicji jednostki wewnętrznej. W tym przypadku konieczne jest podanie, z jakiego pliku należy zacytować informację o treści danej jednostki. Deklaracja:

```
<!ENTITY XML SYSTEM "http://www.ksiazka.pl/jednostki/XML_skrot.xml">
```

oznacza, że w pliku [http://www.ksiazka.pl/jednostki/XML\\_skrot.xml](http://www.ksiazka.pl/jednostki/XML_skrot.xml) znajduje się angielska nazwa języka XML.

Odwołanie do ogólnej jednostki zewnętrznej jest takie samo jak do wewnętrznej:

```
<opis> To jest tekst o XML czyli o &XML;. Angielski termin "&XML;" oznacza
rozszerzalny język znaczników</opis>
```

Różnica między wymienionymi jednostkami jest jedynie w umiejscowieniu rozwinięcia nazwy jednostki: w jednostce wewnętrznej tekst rozwinięcia znajduje się w deklaracji, zaś w jednostce zewnętrznej — w osobnym, zewnętrznym pliku. W obu przypadkach każde odwołanie do jednostki jest zastępowane przez program przetwarzający tekst rozwinięcia.

Kolejnymi, aczkolwiek różnymi od powyżej wymienionych, są **predefiniowane jednostki ogólne** (*predefined entities*). Umożliwiają one umieszczanie w treści elementów, w wartościach atrybutów i innych miejscach dokumentu XML symboli zarezerwowanych. I tak:

- ♦ &lt;: zastępuje znak „<” (*less than*),
- ♦ &gt;: zastępuje znak „>” (*greater than*),
- ♦ &amp;: zastępuje symbol „&” (*ampersand*),
- ♦ &apos;: zastępuje znak apostrofu ' (*apostrophe*),
- ♦ &quot;: zastępuje znak cudzysłowu " (*double quote*).



Jednostki predefiniowane nie wymagają deklarowania za pomocą ENTITY.

Jednostki predefiniowane nie podlegają żadnym zmianom. Nie istnieje także możliwość dodania własnych jednostek, tak by stały się jednostkami predefiniowanymi.

WWW Consortium dla potrzeb współpracy pomiędzy systemami wykorzystującymi język SGML oraz HTML w [84] zamieściło także szerszy zestaw jednostek predefiniowanych. Nie są to jednak jednostki standardowe dla języka XML i niekoniecznie są obsługiwane przez poszczególne parsery.

Wszystkie zaprezentowane dotąd jednostki są **jednostkami ogólnymi** (*general entity*). Można się do nich odwoływać w każdej części dokumentu, a także w deklaracjach DTD.

Istnieją także **jednostki parametryczne** (*parameter entity*), do których odwołania są dostępne jedynie w deklaracjach typu dokumentu DTD.

Jednostki zostały szerzej omówione w rozdziale 3. (patrz: punkt 3.13).

## 2.10. Nośniki informacji w dokumencie XML

Zasadniczym zadaniem dokumentów XML jest przechowywanie informacji. Autorzy tworzący dokumenty mogą w nich umieszczać informacje w różnych miejscach i na różny sposób (patrz rysunek 2.7), a mianowicie:

- ◆ w treściach — tekstach elementów,
- ◆ w nazwach elementów, które określają typ elementu i informują o tym, co zawiera cały element,
- ◆ w hierarchii (strukturze, układzie) elementów, która zawiera informacje o zależnościach między elementami,
- ◆ w nazwach atrybutów,
- ◆ w wartościach atrybutów.

Informacje zawierające się w nazwach elementów i atrybutów są opisem zrozumiałym dla człowieka. W systemach nie wykorzystuje się ich zwykle w sposób bezpośredni, ale dopiero po przetworzeniu, zastępując je najczęściej etykietami, np. element:

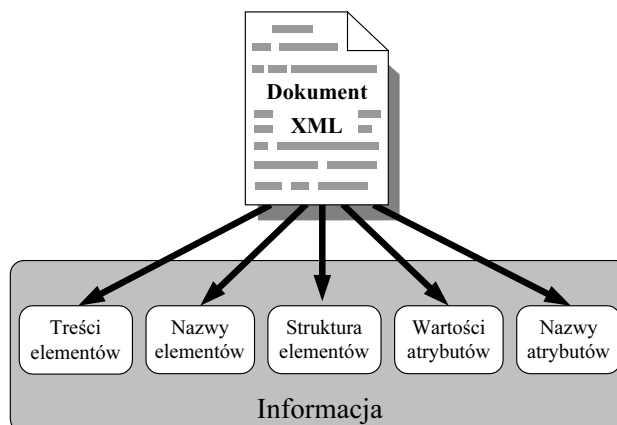
```
<autor> Jan Kowalski </autor>
```

po odpowiednim przekształceniu (transformacji za pomocą XSLT) może mieć końcową postać:

```
Autor sprawozdania: Jan Kowalski
```

Znacznikowi autor odpowiada tutaj tekst Autor sprawozdania:, zaś treść elementu została podana w sposób dosłowny.

**Rysunek 2.7.**  
*Nośniki informacji  
 w dokumencie XML*



Nazwy atrybutów często nie są wykorzystywane w sposób widoczny dla użytkownika końcowego, a jedynie stosowane są w procesie przetwarzania (do identyfikacji), np. znacznik:

```
<faktura rodzaj="korygująca" nr_faktury="12/2001" />
```

może być przekształcony do postaci:

```
Faktura korygująca numer 12/2001
```

Pominięto tutaj nazwę atrybutu `rodzaj`, a użyto jedynie jego wartości (`korygująca`).

Autor dokumentu XML musi określić, w którym miejscu umieszczać odpowiednie informacje, jak nazywać elementy i atrybuty, jaką strukturę nadać dokumentowi. Jest to bardzo ważny etap stanowiący **modelowanie rzeczywistości**. Dokonuje się tego zwykle podczas określania typu dokumentu (przy okazji tworzenia DTD) lub przy tworzeniu schematów XML Schema. W następnych rozdziałach przedstawione zostaną niektóre problemy związane z modelowaniem.

## 2.11. Konwersja dokumentów HTML do XML

Obecnie większość serwisów informacyjnych dostępnych w systemie WWW jest przechowywana w postaci dokumentów utworzonych w języku HTML. Wydaje się, że w przyszłości wiele z nich będzie przetwarzanych do postaci dokumentów XML. Proces takiego przekształcania można opisać za pomocą następujących kroków:

1. Zadbanie o to, by dokument był dobrze uformowany (poprawny składniowo), co nie zawsze ma miejsce w przypadku stron HTML.
2. Przekonwertowanie do języka XHTML. W tym kroku zwykle wystarcza, że dokument jest dobrze uformowany, aczkolwiek język XHTML daje więcej możliwości (np. poprzez modularyzację).

3. Zaprojektowanie struktury informacji zawartej w dokumencie, np. poprzez zdefiniowanie DTD, nawet jeżeli nie powstaną formalne deklaracje DTD. Ustalenie nazw i zawartości elementów, ich atrybutów, jak również struktury elementów (określenie hierarchii elementów).
4. Przeniesienie tekstów występujących w dokumencie HTML do zaprojektowanych elementów XML (i atrybutów). Uzyskujemy w ten sposób dokument XML odpowiadający początkowemu dokumentowi HTML.
5. Wydzielenie elementów związanych z wyświetlaniem dokumentu. Utworzenie transformacji XSLT umożliwiających przejście do odpowiedniego formatu prezentacyjnego dokumentu, np. HTML. Ewentualne utworzenie dokumentu arkuszy stylów XSL zawierających przetworzony dokument źródłowy.



Autorzy stron w języku HTML powinni już obecnie starać się o to, by ich dokumenty były dobrze uformowane (patrz: punkt 2.7).

W poprzednim rozdziale na rysunku 1.1 przedstawiono dokument HTML zawierający przykładowe sprawozdanie, umieszczone w wewnętrznym intranecie pewnej korporacji. Stosując powyższe kroki (na razie bez określania definicji DTD oraz arkuszy stylów) można ten dokument przekształcić w dokument XML. Mógłby on wyglądać następująco (plik *Sprawozdanie.xml*):



```
<?xml version='1.0' ?>
<sprawozdanie rodzaj='roczne'>
  <autorzy liczba_autorow='2'>
    <autor>
      <imie> Anna </imie>
      <nazwisko> Nowicka </nazwisko>
    </autor>
    <autor miejsce='drugi'>
      <imie> Jan </imie>
      <nazwisko> Kowalski </nazwisko>
    </autor>
  </autorzy>
  <tytul> Zestawienie zysków / obrotów </tytul>
  <dane>
    <pozycja waluta='PLN' okres='kwartal' typ_danych='zysk'>
      <nr_okresu> I </nr_okresu>
      <wartosc> 20000.00 </wartosc>
    </pozycja>
    <pozycja okres='kwartal' typ_danych='zysk'>
      <nr_okresu> II </nr_okresu>
      <wartosc> 22000.00 </wartosc>
    </pozycja>
    <pozycja okres='miesiac' typ_danych='obrot' waluta='PLN'>
      <nr_okresu> 7 </nr_okresu>
      <wartosc> 5000.00 </wartosc>
    </pozycja>
  </dane>
</sprawozdanie>
```

Deklaracje typu dokumentu DTD oraz schemat XML Schema dla powyższego dokumentu zamieszczono w dalszej części książki, w punkcie 4.19.