

NIE-PDB: **Advanced Database Systems**

<http://www.ksi.mff.cuni.cz/~svoboda/courses/231-NIE-PDB/>

Lecture 2

Data Formats

Martin Svoboda

martin.svoboda@fit.cvut.cz

10. 10. 2023

Charles University, Faculty of Mathematics and Physics

Czech Technical University in Prague, Faculty of Information Technology

Lecture Outline

Data formats

- XML – Extensible Markup Language
- JSON – JavaScript Object Notation
- BSON – Binary JSON
- RDF – Resource Description Framework

XML

Extensible Markup Language

Introduction

XML = *Extensible Markup Language*

- **Representation and interchange of semi-structured data**
 - + a family of related technologies, languages, specifications, ...
- Derived from **SGML**, developed by **W3C**, started in 1996
- Design goals
 - Simplicity, generality and usability across the Internet
- File extension: ***.xml**, content type: `text/xml`
- Versions: 1.0 and 1.1
- W3C recommendation
 - <http://www.w3.org/TR/xml11/>

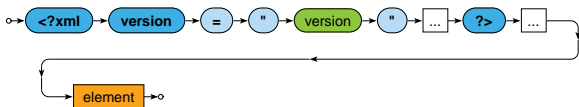
Example

```
<?xml version="1.1" encoding="UTF-8"?>
<movie year="2007">
  <title>Medvídek</title>
  <actors>
    <actor>
      <firstname>Jiří</firstname>
      <lastname>Macháček</lastname>
    </actor>
    <actor>
      <firstname>Ivan</firstname>
      <lastname>Trojan</lastname>
    </actor>
  </actors>
  <director>
    <firstname>Jan</firstname>
    <lastname>Hřebejk</lastname>
  </director>
</movie>
```

Document Structure

Document

- **Prolog**: XML version + some other stuff
- Exactly one **root element**
 - Contains other nested elements and/or other content



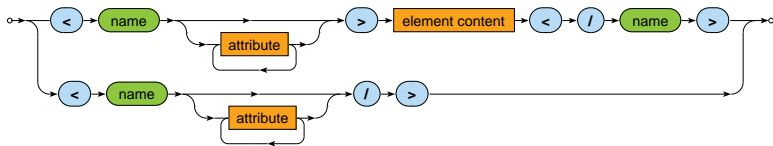
Example

```
<?xml version="1.1" encoding="UTF-8"?>
<movie>
  ...
</movie>
```

Constructs

Element

- Marked using **opening and closing tags**
 - ... or just an abbreviated tag in case of empty elements
- Each element can be associated with a **set of attributes**



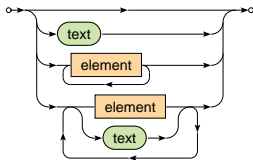
Examples

```
<title>...</title>  
<actors/>
```

Constructs

Types of element content

- **Empty** content
- **Text** content
- **Element** content
 - Sequence of nested elements
- **Mixed** content
 - Elements arbitrarily interleaved with text values



Constructs

Attribute

- Name-value pair



Escaping sequences (predefined entities)

- Used within values of attributes or text content of elements
- E.g.:
 - `<` for `<`
 - `>` for `>`
 - `"` for `"`
 - ...

XML Conclusion

XML constructs

- Basic: **element**, **attribute**, **text**
- Additional: **comment**, **processing instruction**, ...

Schema languages

- DTD, XSD (*XML Schema*), RELAX NG, Schematron

Query languages

- XPath, XQuery, XSLT

XML formats = particular languages

- XSD, XSLT, XHTML, DocBook, ePUB, SVG, RSS, SOAP, ...

JSON

JavaScript Object Notation

Introduction

JSON = *JavaScript Object Notation*

- **Open standard for data interchange**
- Design goals
 - **Simplicity:** text-based, easy to read and write
 - **Universality:** **object** and **array** data structures
 - Supported by majority of modern programming languages
 - Based on conventions of the C-family of languages (C, C++, C#, Java, JavaScript, Perl, Python, ...)
- Derived from JavaScript (but language independent)
- Started in 2002
- File extension: ***.json**
- Content type: **application/json**
- <http://www.json.org/>

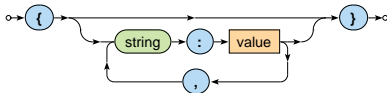
Example

```
{
  "title" : "Medvídek",
  "year" : 2007,
  "actors" : [
    {
      "firstname" : "Jiří",
      "lastname" : "Macháček"
    },
    {
      "firstname" : "Ivan",
      "lastname" : "Trojan"
    }
  ],
  "director" : {
    "firstname" : "Jan",
    "lastname" : "Hřebejk"
  }
}
```

Data Structure

Object

- **Unordered** collection of name-value pairs (**properties**)
 - Correspond to structures such as objects, records, structs, dictionaries, hash tables, keyed lists, associative arrays, ...
- Values can be of different types, **names should be unique**



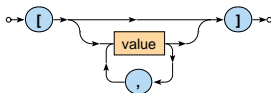
Examples

- { "name" : "Ivan Trojan", "year" : 1964 }
- { }

Data Structure

Array

- **Ordered collection of values**
 - Correspond to structures such as arrays, vectors, lists, sequences, ...
- Values can be of different types, duplicate values are allowed



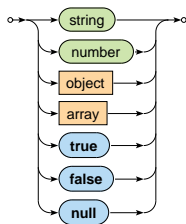
Examples

- [2, 7, 7, 5]
- ["Ivan Trojan", 1964, -5.6]
- []

Data Structure

Value

- Unicode **string**
 - Enclosed with double quotes
 - Backslash escaping sequences
 - Example: `"a \n b \" c \\ d"`
- **Number**
 - Decimal integers or floats
 - Examples: `1`, `-0.5`, `1.5e3`
- **Nested object**
- **Nested array**
- **Boolean** value: `true`, `false`
- Missing information: `null`



JSON Conclusion

JSON constructs

- Collections: **object**, **array**
- Scalar values: **string**, **number**, **boolean**, **null**

Schema languages

- JSON Schema

Query languages

- JSONiq, JMESPath, JAQL, ...

BSON

Binary JSON

Introduction

BSON = *Binary JSON*

- **Binary-encoded serialization of JSON documents**
 - Extends the set of basic data types of values offered by JSON (such as a string, ...) with a few new specific ones
- Design characteristics: **lightweight, traversable, efficient**
- Used by **MongoDB**
 - Document NoSQL database for JSON documents
 - Data storage and network transfer format
- File extension: ***.bson**
- <http://bsonspec.org/>

Example

JSON

```
{ "title" : "Medvídek", "year" : 2007 }
```

BSON

```
24 00 00 00 02 74 69 74 6C 65 00 0A 00 00 00 4D 65 64 76 C3 AD 64 65 6B  
00 10 79 65 61 72 00 D7 07 00 00 00
```

Example

JSON

```
{  
  "title" : "Medvídek",  
  "year" : 2007  
}
```

BSON

```
24 00 00 00  
02 74 69 74 6C 65 00 0A 00 00 00 4D 65 64 76 C3 AD 64 65 6B 00  
10 79 65 61 72 00 D7 07 00 00  
00
```

Document Structure

Document = serialization of **one JSON object or array**

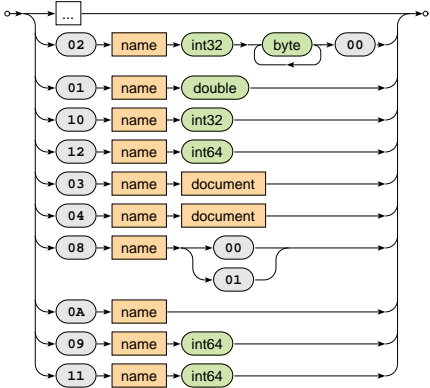
- JSON object is serialized directly
- JSON array is first transformed to a JSON object
 - Property names derived from numbers of positions
 - E.g.:

```
[ "Trojan", "Svěrák" ] →  
{ "0" : "Trojan", "1" : "Svěrák" }
```
- Structure
 - **Document size** (total number of bytes)
 - Sequence of **elements** (encoded JSON properties)
 - Terminating hexadecimal 00 byte



Document Structure

Element = serialization of **one JSON property**



Document Structure

Element = serialization of **one JSON property**

- Structure
 - **Type selector**
 - 02 (**string**)
 - 01 (double), 10 (32-bit **integer**), 12 (64-bit integer)
 - 03 (**object**), 04 (**array**)
 - 08 (**boolean**)
 - 0A (**null**)
 - 09 (datetime), 11 (timestamp)
 - ...
 - **Property name**
 - Unicode string terminated by 00



- **Property value**

RDF

Resource Description Framework

Introduction

RDF = *Resource Description Framework*

- Language for **representing information about resources** in the World Wide Web
 - + a family of related technologies, languages, specifications, ...
 - Used in the context of the **Semantic Web, Linked Data, ...**
- Developed by **W3C**
- Started in 1997
- Versions: 1.0 and 1.1
- W3C recommendations
 - <https://www.w3.org/TR/rdf11-concepts/>
 - Concepts and Abstract Syntax
 - <https://www.w3.org/TR/rdf11-mt/>
 - Semantics

Statements

Resource = any real-world entity

- **Referents** = **resources identified by IRIs**
 - Physical things, documents, abstract concepts, ...
 - E.g.: <http://db.cz/movies/medvidek>
- **Values** = **resources for literals**
 - Numbers, strings, ...
 - E.g.: [Medvídek](#)

Statement about resources = one RDF **triple**

- Three components: **subject**, **predicate**, and **object**

Examples

```
<http://db.cz/movies/medvidek>  
<http://db.cz/terms#title>  
"Medvídek" .
```

```
<http://db.cz/movies/medvidek>  
<http://db.cz/terms#actor>  
<http://db.cz/actors/trojan> .
```

Statements

Triple components

- **Subject**
 - Describes a resource the given statement is about
 - IRI or blank node identifier
- **Predicate**
 - Describes the property or characteristic of the subject
 - IRI
- **Object**
 - Describes the value of that property
 - IRI or blank node identifier or literal

Although triples are inspired by natural languages, they have nothing to do with NLP (Natural Language Processing)

Example

```
<http://db.cz/movies/medvidek>
  <http://db.cz/terms#title> "Medvídek" .

<http://db.cz/movies/medvidek>
  <http://db.cz/terms#actor> <http://db.cz/actors/machacek> .
<http://db.cz/movies/medvidek>
  <http://db.cz/terms#actor> <http://db.cz/actors/trojan> .

<http://db.cz/movies/medvidek>
  <http://db.cz/terms#year> "2007" .

<http://db.cz/movies/medvidek>
  <http://db.cz/terms#director> _:n18 .
_:n18
  <http://db.cz/terms#firstname> "Jan" .
_:n18
  <http://db.cz/terms#lastname> "Hřebejk" .
```

Identifiers and Literals

IRI = *Internationalized Resource Identifier*

- Absolute (not relative) IRIs with optional fragment identifiers
- RFC 3987
- Unicode characters
- Examples
 - `http://db.cz/movies/medvidek`
 - `http://db.cz/terms#actor`
 - `mailto:svoboda@ksi.mff.cuni.cz`
 - `urn:issn:0167-6423`
- **URLs are often used** in practice → information about given resources are then intended to be published / retrieved via standard HTTP

Identifiers and Literals

Literals

- Plain values
 - "Medvídek", "2007"
- Typed values
 - "2007"^^<http://www.w3.org/2001/XMLSchema#int>
 - XML Schema simple data types are usually used
 - `xsd:string`, `xsd:integer`, `xsd:date`, ...
- Strings with language tags
 - "Medvídek"@cs
- *Types and language tags cannot be mutually combined*

Identifiers and Literals

Blank node identifiers

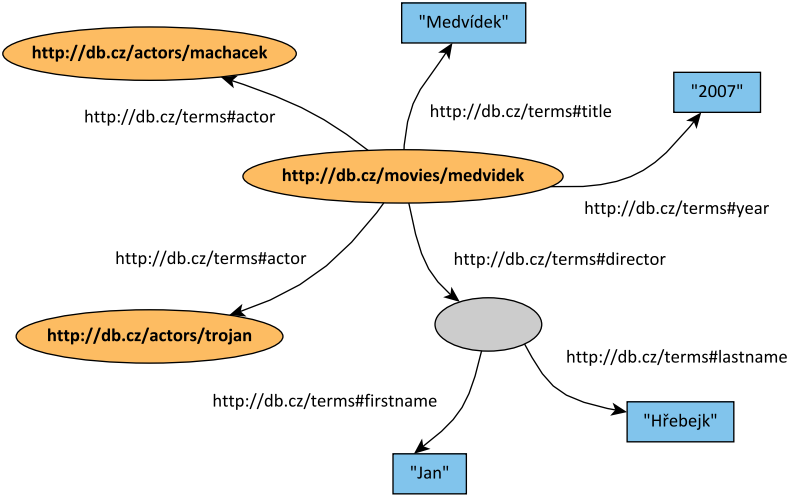
- **Blank nodes** (anonymous resources)
 - Allow to express statements about resources without explicitly naming (identifying) them
- **Blank node identifiers** only have local scope of validity
 - E.g. within a given file, query expression, ...
- Particular syntax depends on a serialization format
 - E.g.: `_:node18`

Data Model

Directed labeled multigraph

- Vertices
 - **One vertex** for each distinct **subject and object**
- Edges
 - **One edge** for each individual **triple**
 - Edges are directed $subject \xrightarrow{predicate} object$
 - Property names (predicate IRIs) are used as edge labels

Example



Serialization

Available approaches

- **N-Triples** notation
 - <https://www.w3.org/TR/n-triples/>
- **Turtle** notation (*Terse RDF Triple Language*)
 - <https://www.w3.org/TR/turtle/>
- **RDF/XML** notation
 - XML syntax for RDF
 - <https://www.w3.org/TR/rdf-syntax-grammar/>
- **JSON-LD** notation
 - JSON-based serialization for Linked Data
 - <https://www.w3.org/TR/json-ld/>
- ...

N-Triples Notation

RDF **N-Triples** notation = *A line-based syntax for an RDF graph*

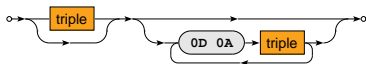
- Simple, line-based, plain text format
- File extension: ***.rdf**
- <https://www.w3.org/TR/n-triples/>

Example

- *Already presented...*

Document

- Statements are terminated by dots, delimited by EOL



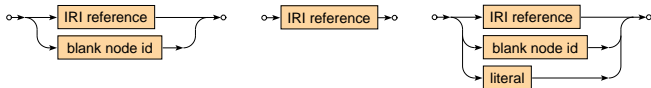
N-Triples Notation

Statement

- Individual triple components are delimited by spaces



Triple components: subject, predicate, object



N-Triples Notation

IRI reference

- IRIs are enclosed in angle brackets

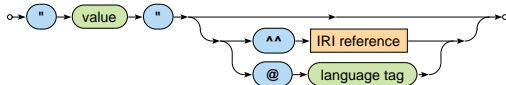


Blank node identifier



Literal

- Literals are enclosed in double quotes



Turtle Notation

Turtle = *Terse RDF Triple Language*

- Compacted format with various abbreviations
- File extension: ***.ttl**
- Content type: text/turtle
- <https://www.w3.org/TR/turtle/>

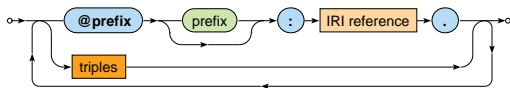
Example

```
@prefix i: <http://db.cz/terms#> .
@prefix m: <http://db.cz/movies/> .
@prefix a: <http://db.cz/actors/> .
m:medvidek
  i:title "Medvídek" ;
  i:actor a:machacek , a:trojan ;
  i:year 2007 ;
  i:director [ i:firstname "Jan" ; i:lastname "Hřebejk" ] .
```

Turtle Notation

Document

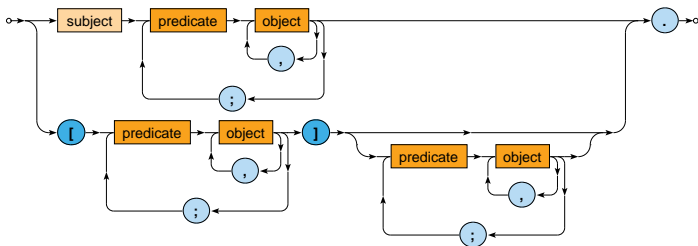
- Contains a **sequence of triples and/or declarations**
- Prefix declarations
 - **Prefixed names** can then be used instead of full IRI references
- Groups of triples
 - Individual groups are terminated by dots



Turtle Notation

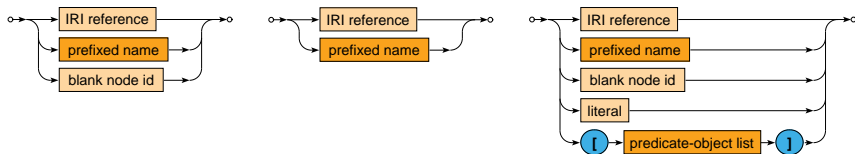
Triples

- Triples sharing the same subject and object or at least the same subject can be *grouped* together
 - **object** list for a **shared subject and predicate**
 - **predicate-object** list for a **shared subject**
- Brackets can be used to define blank nodes

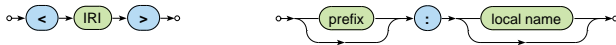


Turtle Notation

Triple components: subject, predicate, object



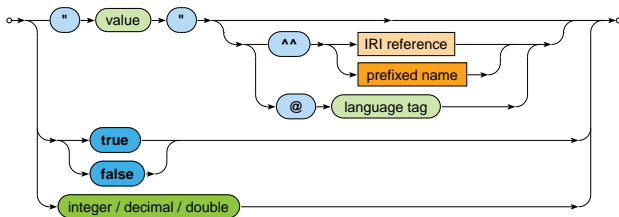
IRI reference / prefixed name



Turtle Notation

Literal

- Traditional literals
+ new abbreviated forms of numeric and boolean literals



Example

Example revisited

```
@prefix i: <http://db.cz/terms#> .  
@prefix m: <http://db.cz/movies/> .  
@prefix a: <http://db.cz/actors/> .  
m:medvidek  
  i:title "Medvídek" ;  
  i:actor a:machacek , a:trojan ;  
  i:year 2007 ;  
  i:director [ i:firstname "Jan" ; i:lastname "Hřebejk" ] .
```

RDF Conclusion

RDF statements

- Subject, predicate, and object components

Schema languages

- RDFS (*RDF Schema*)
- OWL (*Web Ontology Language*)

Query languages

- SPARQL (*SPARQL Protocol and RDF Query Language*)

Lecture Conclusion

Data formats

- Tree: **XML, JSON**
- Graph: **RDF**

Binary serializations

- **BSON**