# Serializing Data with Protocol Buffers

Vinicius Vielmo Cogo
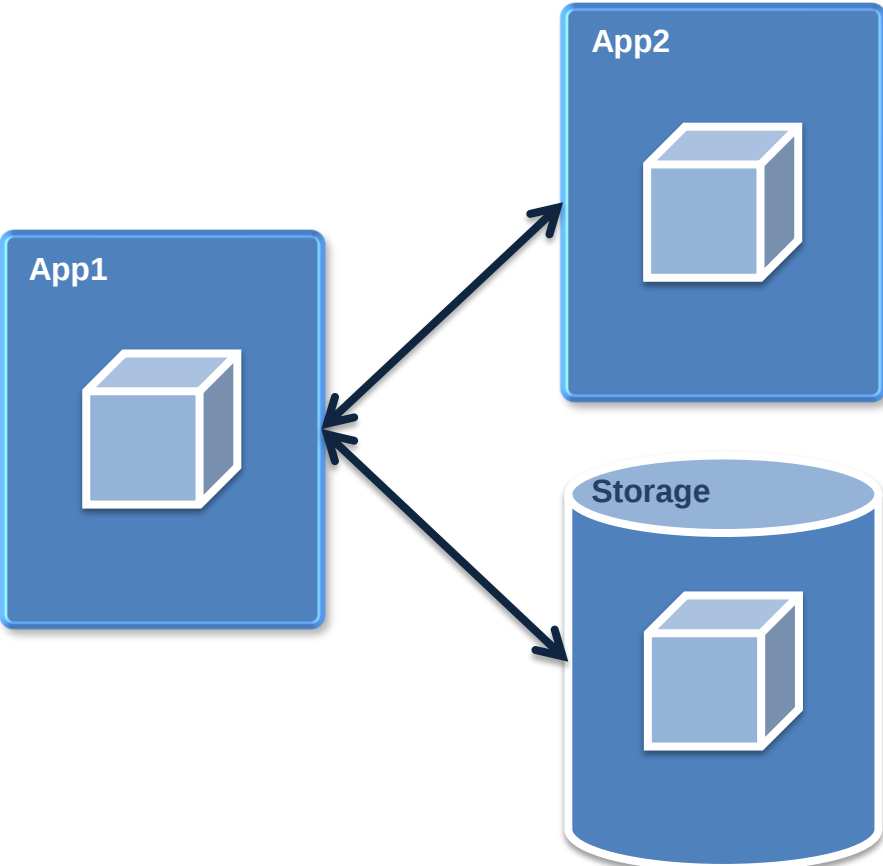Smalltalks, DI, FC/UL. February 12, 2014.
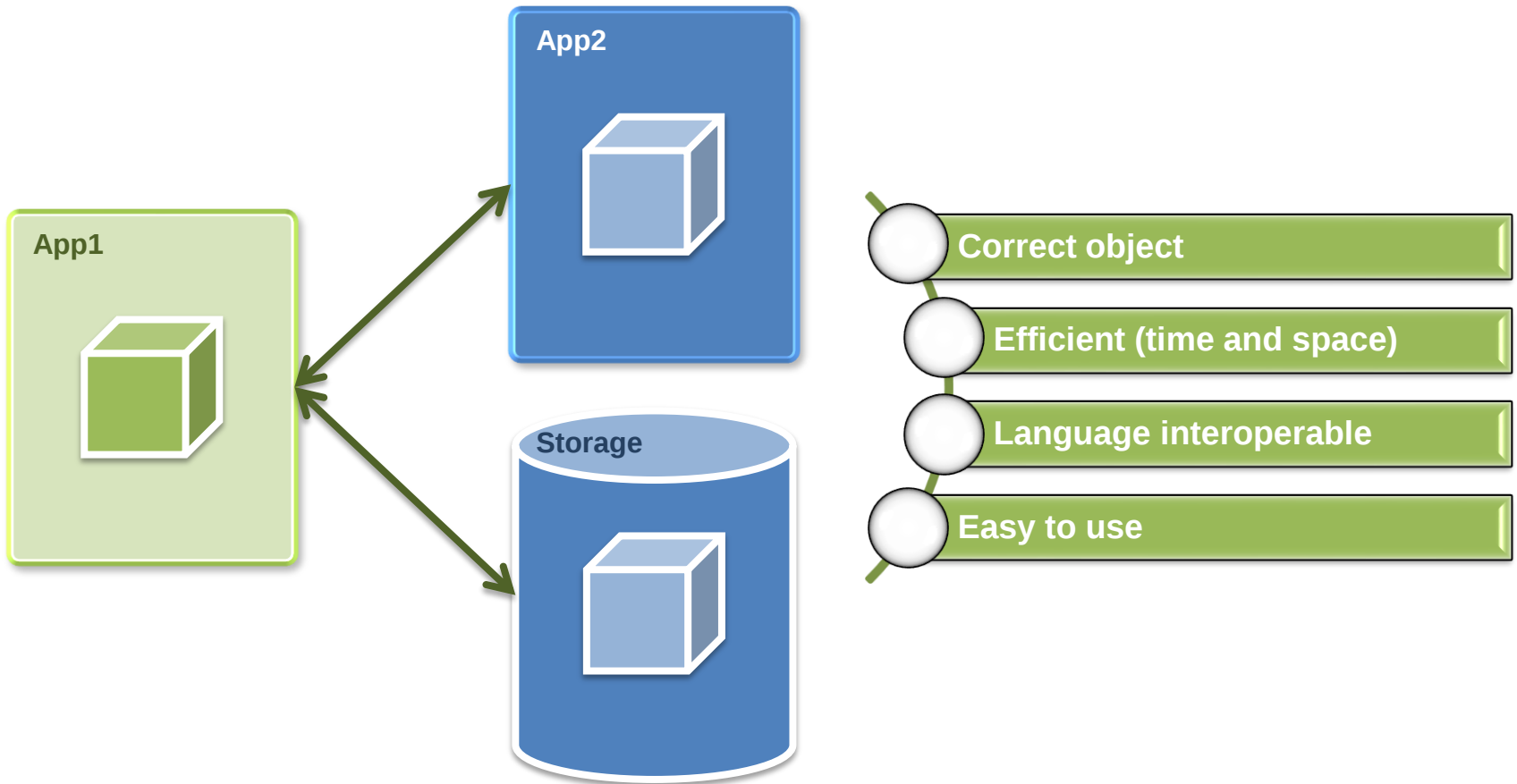
# Problem statement

App2

App1

Storage

App2

App1

Storage

Correct object

Efficient (time and space)

Language interoperable

Easy to use

# Context

- Data serialization (in Java):

1. Java built-in serialization
   - Ex.: writeObject/readObject(aObject)
   - Easy to use, but inefficient in terms of space (extra fields)
   - No language interoperability

2. Manual binary encoding
   - Ex.: 4 ints as *"12:3:-23:67"* and aObject.getBytes()
   - Space efficient, but time efficiency depends on parsing methods
   - Difficult for complex objects
   - Language interoperable

3. Human-readable formats
   - Ex.: Using XML, JSON, DOM, SAX, STAX, JAXB, JAXP, etc.
   - Inefficient (space and time w/ human readable format)
   - Language interoperable

# Protocol Buffers

- Protocol Buffers (Protobuf) is a solution to data serialization:
  - A description language
  - A compiler
  - A library

- Easy-to-use, efficient automatic binary encoding
- Created by and in production at Google Inc.
- Publicly launched in 2008.

- Language interoperable:
  - Officially: Java, C++, and Python
  - Unofficially: C, C#, Erlang, Perl, PHP, Ruby, etc.

- Download and install the Protocol Buffers

Available at:
https://developers.google.com/protocol-buffers/

# Generic workload

1. • Designing objects

2. • Describing objects

3. • Compiling the description

4. • Obtaining the generated source-code

5. • Importing objects into your project

6. • Instantiating objects

7. • Using objects

- Designing objects

## Person:
Id
Name
Age
Email
Phone(s)
...

## Protocol Message:
ClientId
Sequence
Operation
Signature
...

## Sentence:
Id
Language
Subject
Predicate
Verb(s)
Pronoun(s)
Object
...

## Protein:
Id
Organism
Function(s)
Sequence
...

- Describing objects

# Person:
required int32 id
required string name
optional string email
repeated string phone
...

# Protocol Message:
required int32 clientId
required int32 sequence
optional string operation
optional string signature
...

# Sentence:
required int32 id
optional string language
optional string subject
optional string predicate
repeated string verbs
repeated string pronouns
optional string object

# Protein:
required int32 id
optional string organism
repeated string function
optional string sequence
...

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

• Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";
```

Initial configuration

```
message Person {
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {                                              message = object
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {
        required string name = 1;
        required int32 id = 2;                                     Basic fields
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

Enumerations

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];        Internal objects
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];        Default values
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

# • Describing objects

**addressbook.proto**

```
package tutorial;
option java_package = "com.example.tutorial";
option java_outer_classname = "AddressBookProtos";

message Person {
        required string name = 1;
        required int32 id = 2;
        optional string email = 3;
        enum PhoneType {
                MOBILE = 0;
                HOME = 1;
                WORK = 2;
        }
        message PhoneNumber {
                required string number = 1;
                optional PhoneType type = 2 [default = HOME];
        }
        repeated PhoneNumber phone = 4;
}

message AddressBook {
        repeated Person person = 1;
}
```

List of objects

# • Compiling the description

Terminal

```
$ protoc --java_out=$DST_DIR addressbook.proto
```

Output folder

.proto fiile
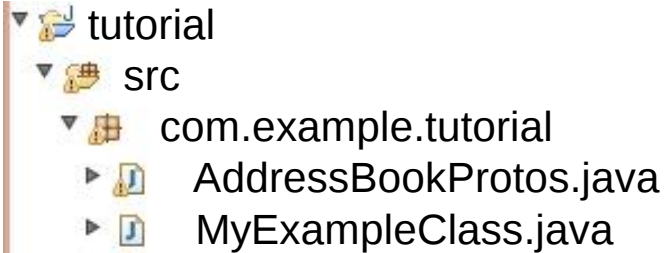
# • Obtaining the generated source-code

**Terminal**

```
$ cat $DST_DIR/com/example/tutorial/AddressBookProtos.java
2478    java.lang.String[] descriptorData = {
2479      "\n\021addressbook.proto\022\010tutorial\"\332\001\n\006Person" +
2480      "\022\014\n\004name\030\001 \002(\t\022\n\n\002id\030\002 \002(\005\022\r\n\005email\030\003 \001(" +
2481      "\t\022+\n\005phone\030\004 \003(\0132\034.tutorial.Person.Phone" +
2482      "Number\032M\n\013PhoneNumber\022\016\n\006number\030\001 \002(\t\022.\n" +
2483      "\004type\030\002 \001(\0162\032.tutorial.Person.PhoneType:" +
2484      "\004HOME\"+\n\tPhoneType\022\n\n\006MOBILE\020\000\022\010\n\004HOME\020\001" +
2485      "\022\010\n\004WORK\020\002\"/\n\013AddressBook\022 \n\006person\030\001 \003(" +
2486      "\0132\020.tutorial.PersonB)\n\024com.example.tutor" +
2487      "ialB\021AddressBookProtos"
2488    };
2489    com.google.protobuf.Descriptors.FileDescriptor.InternalDescriptorAssigner assigner =
2490      new com.google.protobuf.Descriptors.FileDescriptor.InternalDescriptorAssigner() {
2491        public com.google.protobuf.ExtensionRegistry assignDescriptors(
2492          com.google.protobuf.Descriptors.FileDescriptor root) {
2493          descriptor = root;
2494          internal_static_tutorial_Person_descriptor =
2495          getDescriptor().getMessageTypes().get(0);

2496          internal_static_tutorial_Person_fieldAccessorTable = new
2497            com.google.protobuf.GeneratedMessage.FieldAccessorTable(
2498              internal_static_tutorial_Person_descriptor,
2499              new java.lang.String[] { "Name", "Id", "Email", "Phone", });
2500          internal_static_tutorial_Person_PhoneNumber_descriptor =
2501            internal_static_tutorial_Person_descriptor.getNestedTypes().get(0);
2502          internal_static_tutorial_Person_PhoneNumber_fieldAccessorTable = new
2503            com.google.protobuf.GeneratedMessage.FieldAccessorTable(
2504              internal_static_tutorial_Person_PhoneNumber_descriptor,
2505              new java.lang.String[] { "Number", "Type", });
2506          internal_static_tutorial_AddressBook_descriptor =
2507            getDescriptor().getMessageTypes().get(1);
2508          internal_static_tutorial_AddressBook_fieldAccessorTable = new
2509            com.google.protobuf.GeneratedMessage.FieldAccessorTable(
2510              internal_static_tutorial_AddressBook_descriptor,
2511              new java.lang.String[] { "Person", });
2512          return null;
2513        }
2514      };
2515    com.google.protobuf.Descriptors.FileDescriptor
2516      .internalBuildGeneratedFileFrom(descriptorData,
2517        new com.google.protobuf.Descriptors.FileDescriptor[] {
2518        }, assigner);
2519  }
2520
2521  // @@protoc_insertion_point(outer_class_scope)
2522 }
```

# • Importing objects into your project

**Terminal**

```
$ cp $DST_DIR/com/example/tutorial/AddressBookProtos.java
~/workspace/tutorial/src/com/example/tutorial/
```

- ▼ 📁 tutorial
  - ▼ 📦 src
    - ▼ 📦 com.example.tutorial
      - ▶ 📄 AddressBookProtos.java
      - ▶ 📄 MyExampleClass.java

**MyExampleClass.java**

```
package com.example.tutorial;

import com.example.tutorial.AddressBookProtos.AddressBook;
import com.example.tutorial.AddressBookProtos.Person;

...
```

• Instantiating objects

**MyExampleClass.java**

```
...
Person  john = Person.newBuilder()
          .setId(12345)
          .setName("John Foo Bar")
          .setEmail("john@foobar.pt")
          .addPhone(Person.PhoneNumber.newBuilder()
                    .setNumber("+351 999 999 999")
                    .setType(Person.PhoneType.HOME)
                    .build())
          .build();
...
```

# • Using objects: Storage

**MyStorageExample.java**

```java
// Writing data to a file

FileOutputStream  aOutput = new FileOutputStream("theFilename");

Person  aPerson = Person.newBuilder().set… //instantiate a Person

aPerson.writeTo(aOutput);

aOutput.close();

// Reading data from a file

Person  aPerson = Person.parseFrom(new FileInputStream("theFilename"));

// Do something with the received Person
```

# • Using objects: TCP Communication

**MyTCPCommunicationExample.java**

```
// Server-side

ServerSocket  aServerSocket = new ServerSocket(10000);

Socket  aConnection = aServerSocket.accept();

Person  aPerson = Person.parseDelimitedFrom(aConnection.getInputStream());

// Do something with the received Person

aPerson.writeDelimitedTo(aConnection.getOutputStream());

// Client-side

Person  aPerson = Person.newBuilder().set... //instantiate a Person

Socket  aSocket = new Socket("127.0.0.1", 10000);

aPerson.writeDelimitedTo(aSocket .getOutputStream());

Person  aPerson = Person.parseDelimitedFrom(aConnection . getInputStream());
// Do something with the received Person
```

# • Using objects: UDP Communication

**MyUDPCommunicationExample.java**

```java
// Receive a packet
DatagramSocket   aServerSocket = new DatagramSocket(10000);
byte[]   aReceiveData = new byte[1024];
DatagramPacket   aReceivePacket = new DatagramPacket(aReceiveData, aReceiveData.length);
aServerSocket.receive(aReceivePacket);
ByteArrayInputStream   aInput = new ByteArrayInputStream(aReceiveData);
Person   aPerson = Person.parseDelimitedFrom(aInput);
// Do something with the received Person
// Send a packet
DatagramSocket   aClientSocket = new DatagramSocket(10001);
ByteArrayOutputStream   aOutput = new ByteArrayOutputStream(1024);
Person   aPerson = Person.newBuilder().set... //instantiate a Person
aPerson .writeDelimitedTo(aOutput);
byte   aSendData[] = aOutput.toByteArray();
InetAddress   aIp = InetAddress.getLocalHost(); // or aReceivePacket.getAddress();
DatagramPacket   aSendPacket = new DatagramPacket(aSendData, aSendData.length, aIp, 10001);
aClientSocket.send(aSendPacket);
```
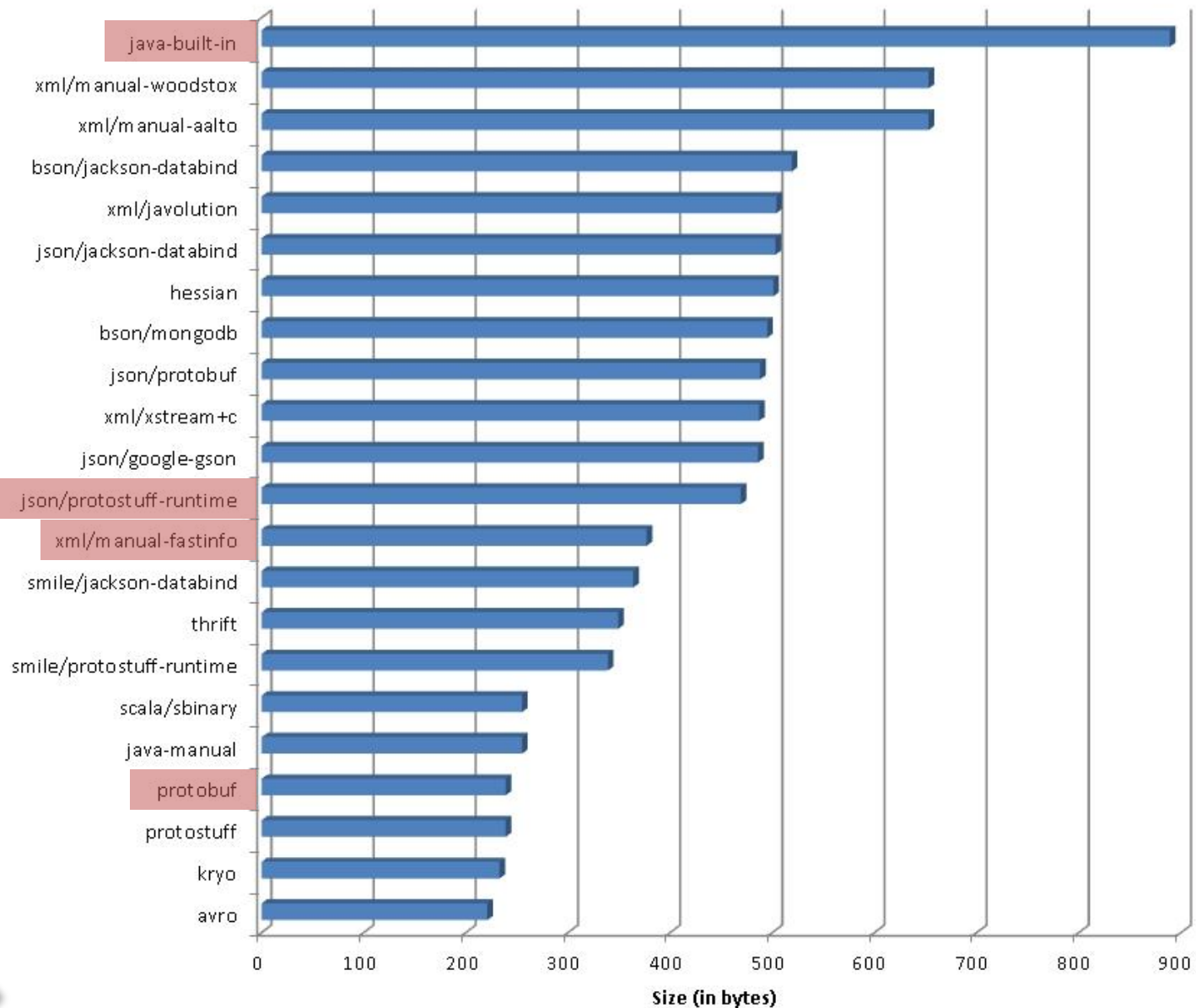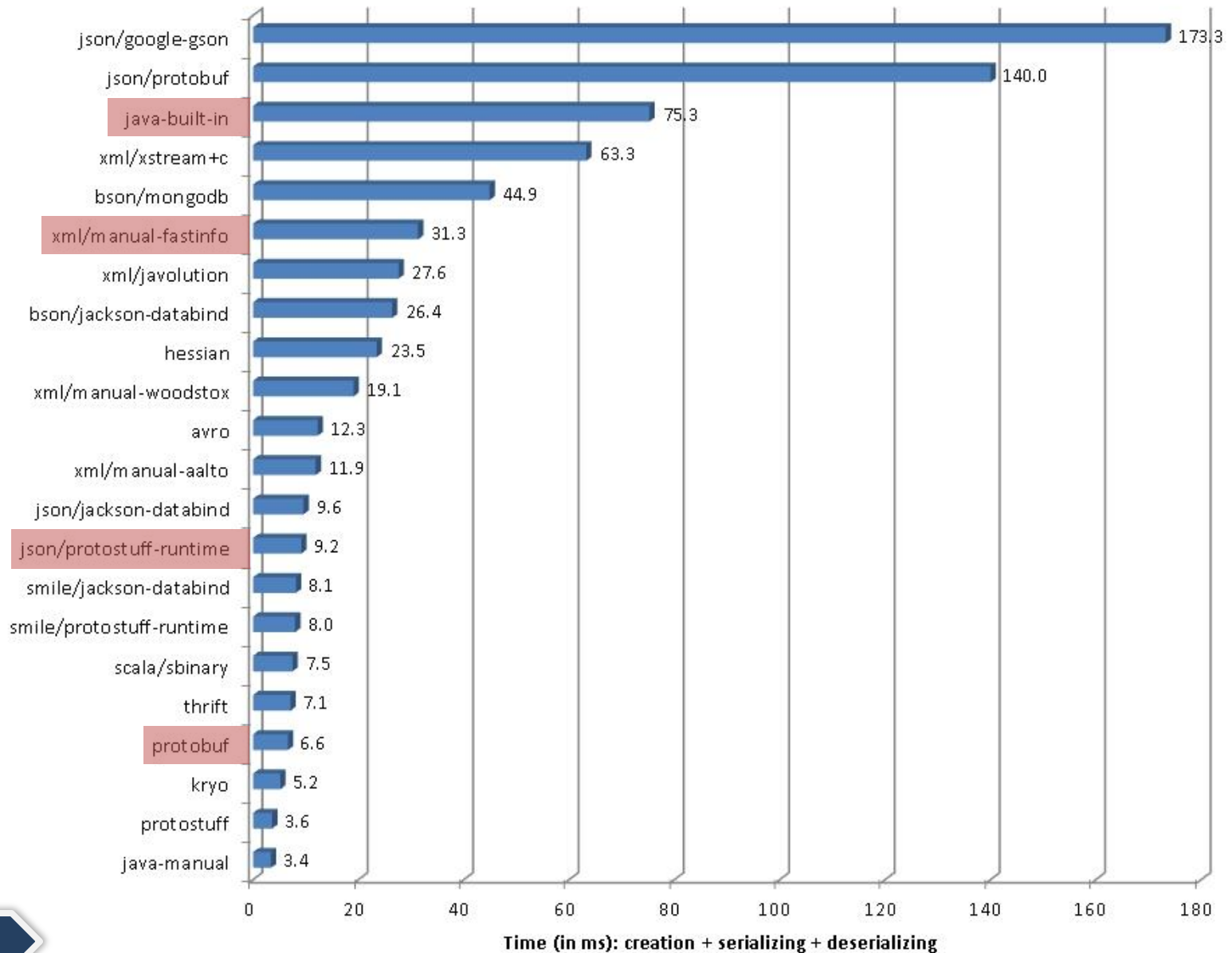
# Performance: Size

Time (in ms): creation + serializing + deserializing

| Method | Value |
|---|---|
| json/google-gson | 173.3 |
| json/protobuf | 140.0 |
| java-built-in | 75.3 |
| xml/xstream+c | 63.3 |
| bson/mongodb | 44.9 |
| xml/manual-fastinfo | 31.3 |
| xml/javolution | 27.6 |
| bson/jackson-databind | 26.4 |
| hessian | 23.5 |
| xml/manual-woodstox | 19.1 |
| avro | 12.3 |
| xml/manual-aalto | 11.9 |
| json/jackson-databind | 9.6 |
| json/protostuff-runtime | 9.2 |
| smile/jackson-databind | 8.1 |
| smile/protostuff-runtime | 8.0 |
| scala/sbinary | 7.5 |
| thrift | 7.1 |
| protobuf | 6.6 |
| kryo | 5.2 |
| protostuff | 3.6 |
| java-manual | 3.4 |

# Final remarks

- Protobuf focuses on:
    - Efficiency (space and time)
    - Language interoperability
    - Usability

- It is a good alternative for built-in serialization

- But there are other available solutions:
    - Avro (http://avro.apache.org/)
    - Thrift (http://thrift.apache.org/)
    - Kryo (https://github.com/EsotericSoftware/kryo)

# Extra 1: .proto types

| .proto Type | Notes | C++ Type | Java Type | Python Type[2] |
|---|---|---|---|---|
| double | | double | double | float |
| float | | float | float | float |
| int32 | Uses variable-length encoding. Inefficient for encoding negative numbers – if your field is likely to have negative values, use sint32 instead. | int32 | int | int |
| int64 | Uses variable-length encoding. Inefficient for encoding negative numbers – if your field is likely to have negative values, use sint64 instead. | int64 | long | int/long[3] |
| uint32 | Uses variable-length encoding. | uint32 | int[1] | int/long[3] |
| uint64 | Uses variable-length encoding. | uint64 | long[1] | int/long[3] |
| sint32 | Uses variable-length encoding. Signed int value. These more efficiently encode negative numbers than regular int32s. | int32 | int | int |
| sint64 | Uses variable-length encoding. Signed int value. These more efficiently encode negative numbers than regular int64s. | int64 | long | int/long[3] |
| fixed32 | Always four bytes. More efficient than uint32 if values are often greater than $2^{28}$. | uint32 | int[1] | int |
| fixed64 | Always eight bytes. More efficient than uint64 if values are often greater than $2^{56}$. | uint64 | long[1] | int/long[3] |
| sfixed32 | | int32 | int | int |
| sfixed64 | | | | nt/long[3] |
| bool | | | | oolean |
| string | ain UTF-8 encoded or 7-bit ASCII text. | string | String | str/unicode[4] |
| bytes | May contain any arbitrary sequence of bytes. | string | ByteString | str |

See the Protobuf language guide:
https://developers.google.com/protocol-buffers/docs/proto

# Extra 2: Polymorphism

```
message BaseType
{
    // Reserve field numbers 100 to 199 for extensions.
    extensions 100 to 199;
    // All other field numbers are available for use here.
    required string name = 1;
    optional uint32 quantity = 2;
}

extend BaseType
{
    // This extension can only use numbers 100 to 199.
    optional float price = 100;
}
```

```
message Cat
{
    optional bool declawed = 1;
}
message Dog
{
    optional uint32 bones_buried = 1;
}
message Animal
{
    required float weight = 1;
    optional Dog dog = 2;
    optional Cat cat = 3;
}
```

See the tutorial about Protocol Buffers Polymorphism:
http://www.indelible.org/ink/protobuf-polymorphism/

# Extra 3: Netty and Protocol Buffers



Transport Services

| Socket & Datagram | | |
| HTTP Tunnel | | |
| In-VM Pipe | | |

Protocol Support

| HTTP & WebSocket | SSL · StartTLS | Google Protobuf |
| zlib/gzip Compression | Large File Transfer | RTSP |
| Legacy Text · Binary Protocols with Unit Testability | | |

Core

| Extensible Event Model |
| Universal Communication API |
| Zero-Copy-Capable Rich Byte Buffer |

See the following Netty packages and examples:
 io.netty.handler.codec.protobuf
 io.netty.example.worldclock

# Extra 4: To delimit or not?

| writeTo(OutputStream o) | writeDelimitedTo(OutputStream o) |
|---|---|
| • Serializes the message<br>• Writes the message to output<br>• No flush()<br>• No close() | • Writes the size of the message<br>• Serializes the message<br>• Writes the message to output |

| parseFrom(InputStream i) | parseDelimitedFrom(InputStream i) |
|---|---|
| • Reads the InputStream until EOF<br>• No close() | • Reads the size of message<br>• Reads the message data |

```
record Image = {
  uri: String
  title: String?
  width: Int32
  height: Int32
  size: Size

  enum Size = { SMALL,
LARGE, }
}
```

```
record Media = {
  uri: String
  title: String?
  width: Int32
  height: Int32
  format: String
  duration: Int64
  size: Int64
  bitrate: Int32?
  persons: List<String>
  player: Player
  copyright: String?
  enum Player = { JAVA,
FLASH, }
}
```

```
record MediaContent = {
  images: List<Image>
  media: Media
}
```

See the following links:
http://code.google.com/p/thrift-protobuf-compare/wiki/BenchmarkingV2
https://groups.google.com/forum/#!forum/java-serialization-benchmarking

# Extra 6: General Tips

- Updating a .proto file can be compatible with previous versions

- Prefer optional fields than required (required is forever)

- **Generic messages with one "TYPE" field and multiple** optional specific messages fields is a good solution for complex protocols

- Same object with different structure are supported through optional fields