LECTURER: NGHIA DUONG-TRUNG
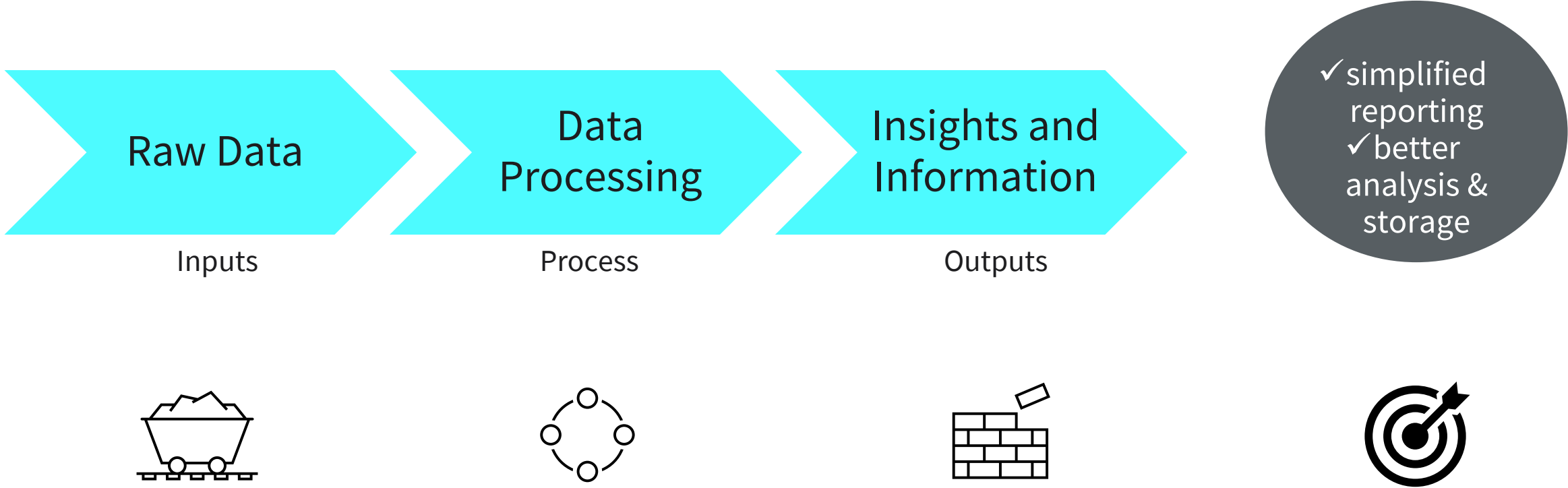
# DATA SCIENCE

**TOPIC OUTLINE**

# PROCESSING OF DATA

On completion of this unit, you will have learned …

— the concepts of data, information, and data processing.

— the stages and cycles of data processing.

— the different methods and types of data processing.

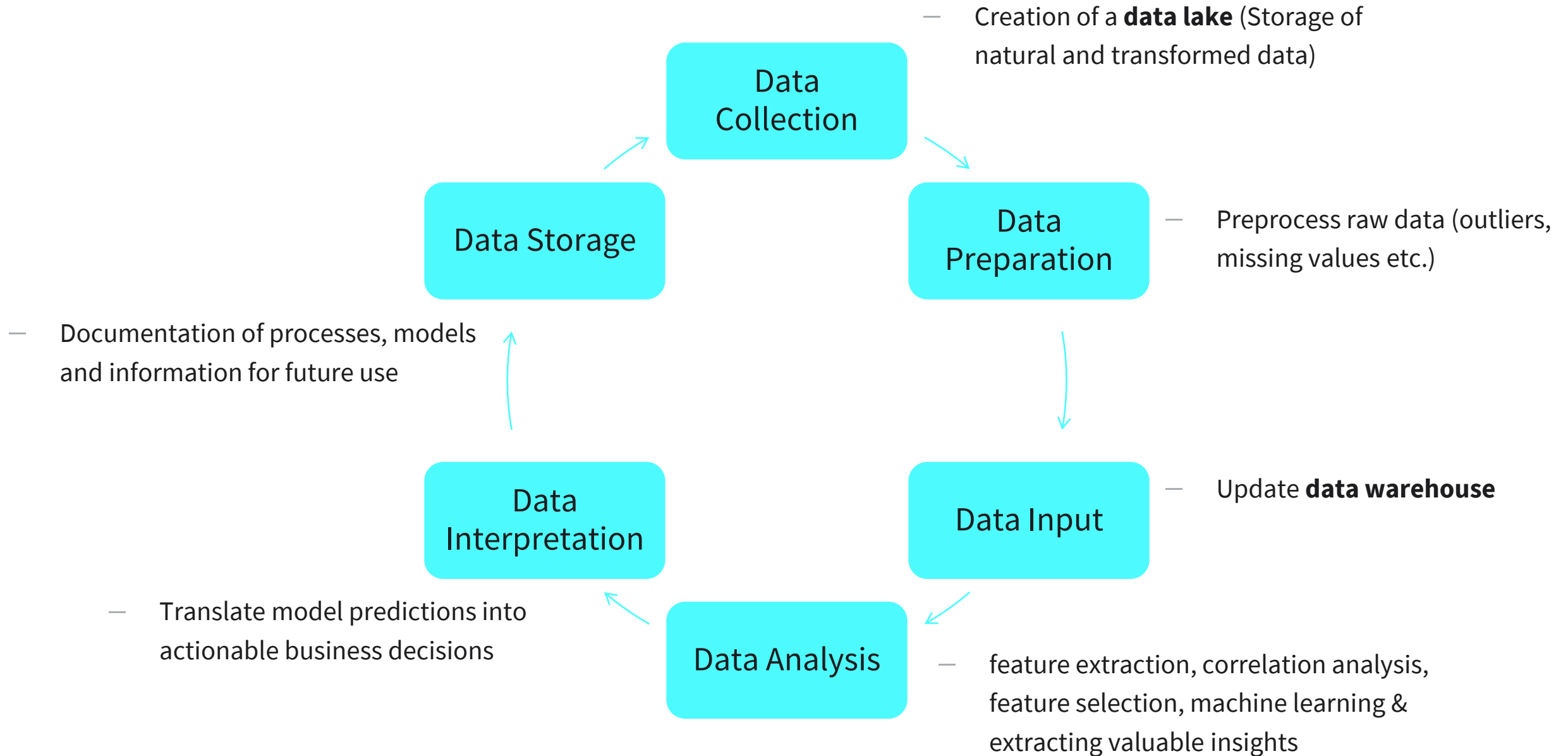— the output forms and file formats for processed data.

1. Explain why data processing is important.
2. In what way benefit data science projects from data processing.
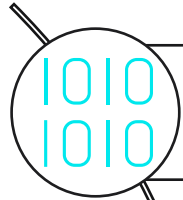3. Describe the five types of electronic data processing.

**DATA PROCESSING INTRODUCTION**

| Raw Data | Data Processing | Insights and Information | ✓ simplified reporting ✓ better analysis & storage |
|---|---|---|---|
| Inputs | Process | Outputs | |

**DATA PROCESSING CYCLE**



— Creation of a **data lake** (Storage of natural and transformed data)

— Preprocess raw data (outliers, missing values etc.)

— Update **data warehouse**

— feature extraction, correlation analysis, feature selection, machine learning & extracting valuable insights

— Translate model predictions into actionable business decisions

— Documentation of processes, models and information for future use

Data Collection

Data Storage

Data Preparation

Data Interpretation

Data Input

Data Analysis

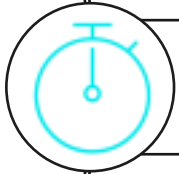# ELECTRONICAL DATA PROCESSING

**Batch**
split data into batches to permit sequential processing (mostly offline)
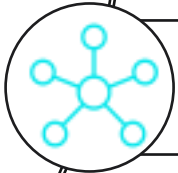
**Online**
make use of internet connections

**Real-time**
immediate response to requests

**Distributed**
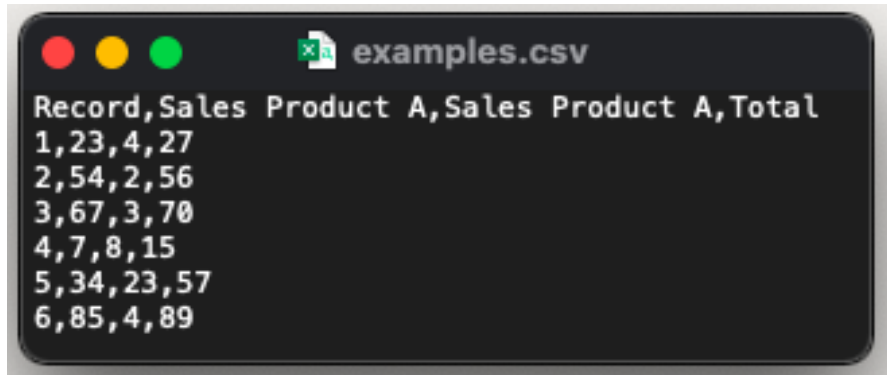multiple remote workstations connected to a large server

**Time-sharing**
computing unit is utilized by multiple users

## CSV (comma-separated-value)

— row-based: every line represents one record

— features are separated by comma

## XLS (Excel spreadsheet)
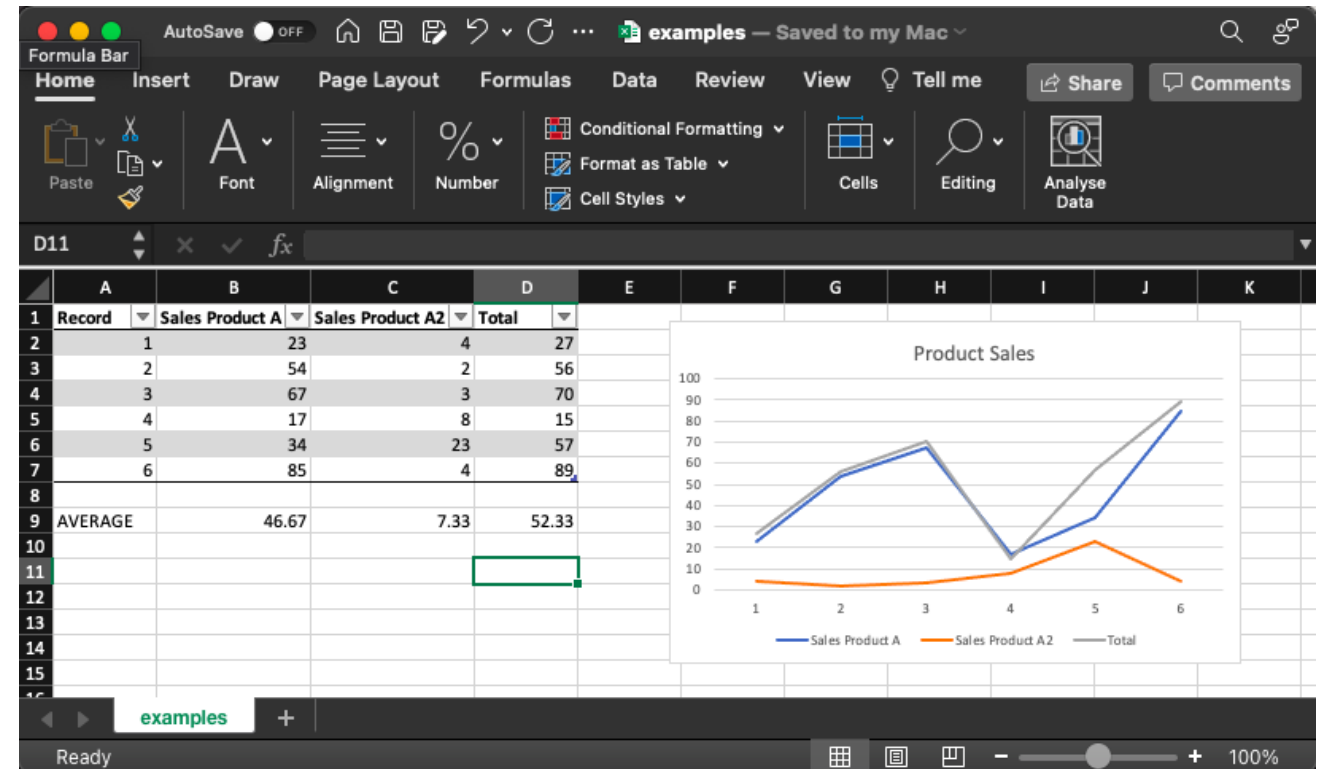
— tabular format of records and features

— possibility to add graphs, computations



Source of the graphics: Antonia Schulze, 2021.

**OUTPUT FORMATS OF PROCESSED DATA**

# XML (extensible markup language)

— structured, non-tabular data written as text with annotations
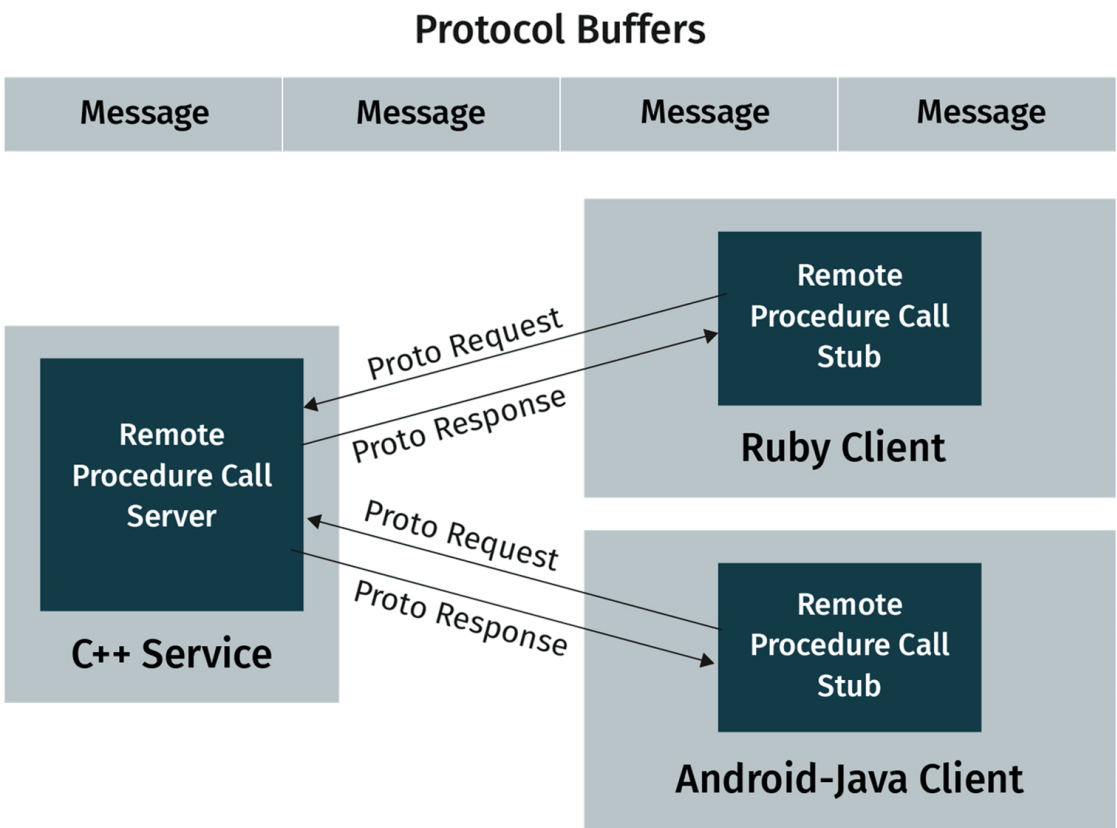
```
<student>
 <name>
    <firstname> John </firstname>
    <lastname> Miller</lastname>
 </name>
 <birthdate> 1999-12-12</birthdate>
</student>
<student>
 <name>
    <firstname> Alice </firstname>
    <lastname> Doe</lastname>
 </name>
 <birthdate> 1987-01-06</birthdate>
</student>
```
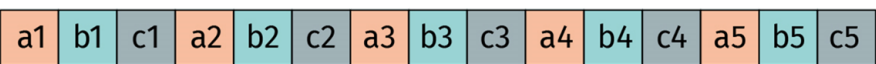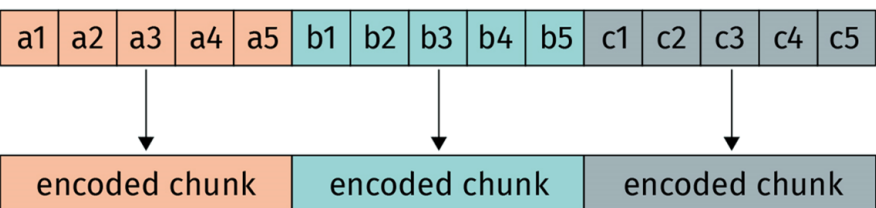
# Protobuf (protocol buffers)

— reduced XML version



**Protocol Buffers**

| Message | Message | Message | Message |

Remote Procedure Call Stub
**Ruby Client**

Remote Procedure Call Server
**C++ Service**

Proto Request
Proto Response

Remote Procedure Call Stub
**Android-Java Client**

Proto Request
Proto Response

# Apache Parquet
— column-based file format

# JSON (Java script object notation)
— a list of key-value pairs

You have learned …

— the concepts of data, information, and data processing.
— the stages and cycles of data processing.
— the different methods and types of data processing.
— the output forms and file formats for processed data.

# TRANSFER TASK

Create a framework that helps data practitioners to choose the best sub-type of electronic data processing.
Which questions should they ask themselves?

Please present your results.

The results will be discussed in plenary.

1. In which step is data with missing values handled?

a) feature selection

b) machine learning

c) correlation analysis

d) data pre-processing

2. The data provided in this format, <img fig="Alice.jpg" tag="Alice" /> , represents the…

    a) SQL data format.
    b) XLS data format.
    c) XML data format.
    d) CSV data format.

3. The patterns and relationships among data elements are defined as ...

a) data.
b) properties.
c) information.
d) features.