

DATA FORMATS

CS 346: Application
Development

APPLICATIONS = DATA + COMPUTATION

Meaningful applications interact with data.

- e.g., text editors write text files; graphics editors load and modify images.

Consider what data your application manipulates:

- What is the source of data: is it loaded from a file or streamed from a service?
- Does it have a format that I need to be able to manipulate? e.g., JPEG images, MP4 video.
- Do I need to cache data locally, or do I reload it as-needed? If so, when and how?
- How do I represent this data in my code?

You also need to consider how the data will be used:

- Is the data specific to a user account (e.g., password), or an application (e.g., window size).
- Do you need to export or transmit the data, or store is in a shared location?
- What are the privacy and security implications of transmitting or storing this data?

TERMINOLOGY

Data can have a simple structure (e.g. the user's name), or complex one (e.g. a customer with a name, address, job title).

We often use these terms:

- **field:** a particular piece of data, corresponding to members or variables in a program. e.g. "name" in our Customer class.
- **record:** a collection of fields that together comprise a single instance of a class or object. e.g. an instance of our Customer class.
- **collection:** a group of records that belong together. e.g. all Customer records.

We might change up the nomenclature for different forms of data, but this structure applies to all kinds of data:

- e.g., Image editor: collection of images, and an image data model which contains text (name, creation date), possibly integer data (image size, checksum) and binary data (the actual image).

DATA MODELS

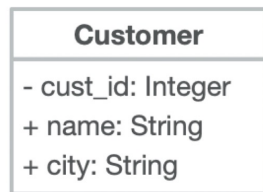
We typically have **data models** that represents our data. You may choose to express this in different ways, depending on where it is stored.

- e.g., struct in C, class in Kotlin.

We usually work with logical models, which are abstract enough to work with any concretion.

General approach

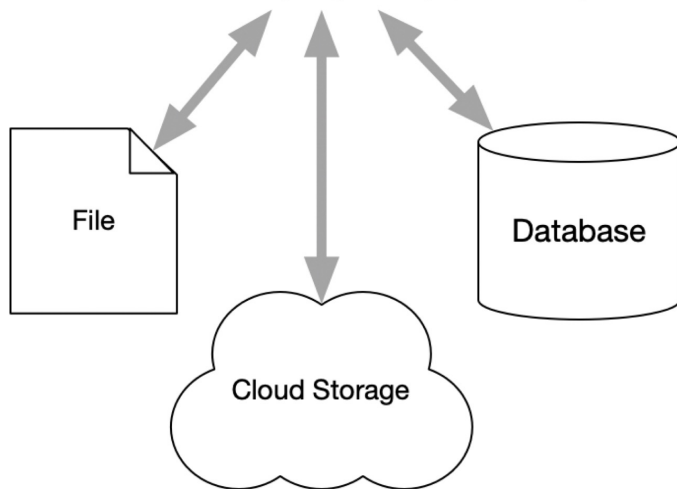
- Data models: class models that reflect the data that you need to store.
- Mechanisms to convert these data models to other representations as needed.
- e.g., subset of data to be expressed in the UI; records to be stored in a database or in a file.



```

val c1 = Customer("Jeff Avery", "Waterloo")
val c2 = Customer("Marie Curie", "Paris")
val c3 = Customer("Billy Bishop", "Ottawa")

```



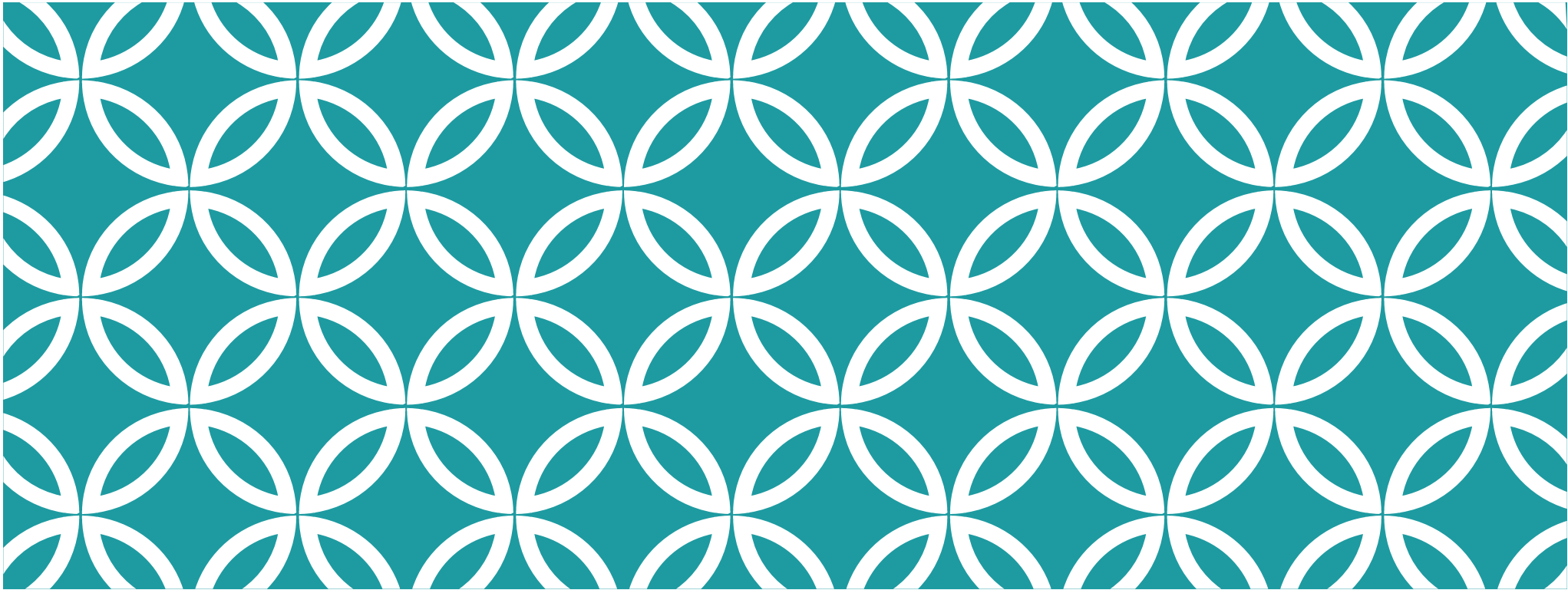
← class representation (UML diagram)

← logical records (fields+data)

We can define a Customer Class in Kotlin that reflects this data model. We want the flexibility to also:

- save to a file
- send to a remote machine
- save to a database

Our objective is to simplify moving data around our system (while ensuring data integrity).



DATA MODELS IN KOTLIN

CS 346: Application
Development

DATA CLASSES: PRIMITIVE DATA

The obvious choice for storing records would be to use a data class for the record class, and a collection class for the set of records.

```
data class Customer (  
    val cust_id: Int,  
    val name: String,  
    val city: String  
)
```

```
val customers = List<Customer>() // List of customers  
customers.add(Customer(1, "John Hall", "New York"))
```

DATA CLASSES: BINARY DATA

Binary data isn't much more challenging; you just need to use the built-in Byte types: Byte (UByte), and ByteArray (UByteArray) for signed (unsigned) values.

```
data class Customer (  
    val cust_id: Int,  
    val name: String,  
    val city: String,  
    val avatar: Array<UByte> = ubyteArrayOf() // unsigned 8-bit ints  
)  
  
val customers = List<Customer>()  
customers.add(Customer(1, "John Hall", "New York",  
    ubyteArrayOf(UByte.MIN_VALUE, 130U, 131u, UByte.MAX_VALUE)))  
) // 0, 130, 131, 255
```


PROCESSING DATA CLASSES

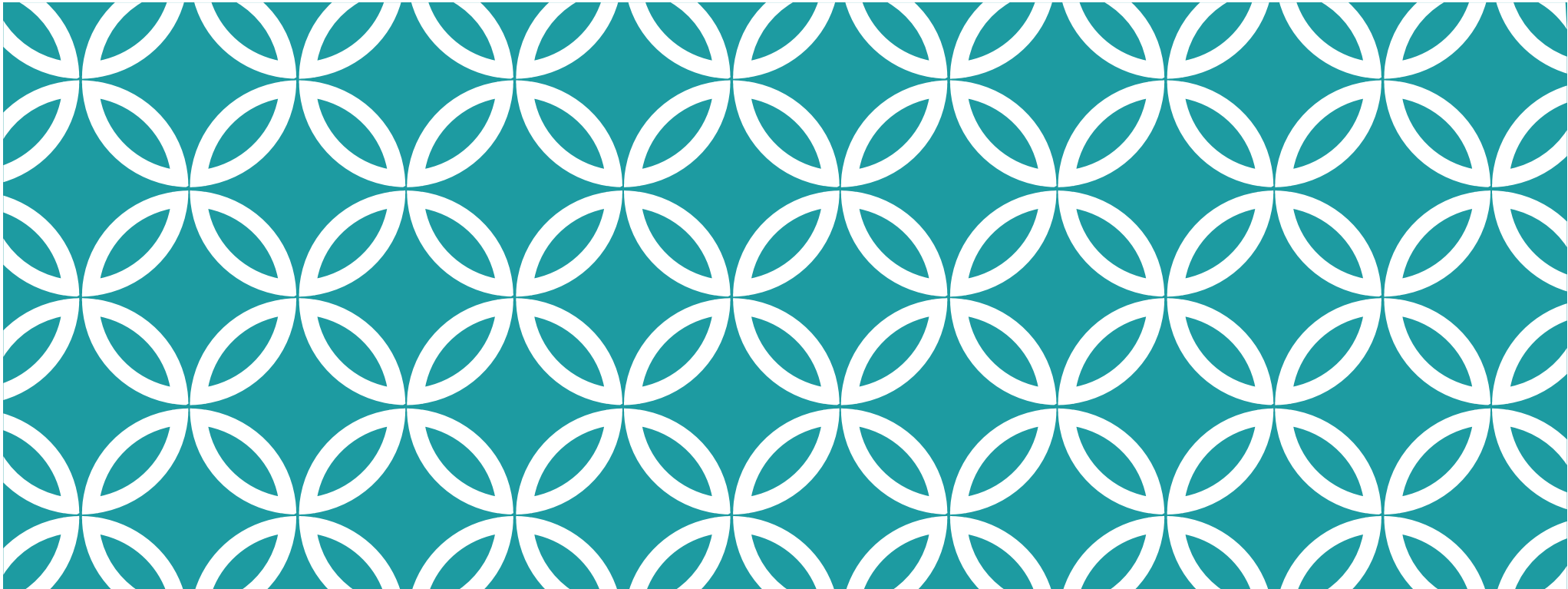
How do you convert a data class into something that you can:

- write to a file,
- transmit over a network connection,
- save to a database?

You will likely need to do all these things, even for a trivial application.

- e.g., local application settings should be saved to a file
- e.g., core data will need to be saved to a database or service.

What formats are suitable for these scenarios?



FILE STORAGE

CS 346: Application
Development

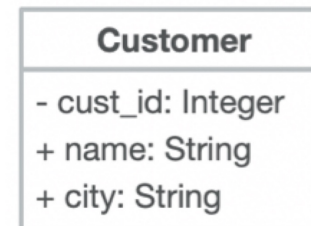
STRUCTURED CSV

The simplest way to store records might be to use a CSV (comma-separated values) file. We use this structure:

- Each row corresponds to one record (i.e. one object)
- The values in the row are the field for each record, separated by commas.

For example, transaction data file stored in a comma-delimited file:

```
1001, Jeff Avery, Cambridge  
1002, Allison Barnett, Waterloo  
1003, John McAfee, Delphi
```



CSV realization of this class data.

READING/WRITING OBJECTS TO CSV

Customer
- cust_id: Integer
+ name: String
+ city: String

```
data class Customer (val cust_id: Int, val name: String, val city: String)
```

```
val customers = List<Customer>() // list of customers
customers.add(Customer(1001, "John Hall", "New York"))
customers.add(Customer(1002, "Allison Barnett", "Waterloo"))
customers.add(Customer(1003, "John McAfee", "Delphi"))
```

```
File("output.csv").open("w").use {
    it.write("Customer ID, Name, City\n")
    for (customer in customers) {
        it.write("${customer.cust_id},${customer.name},${customer.city}\n")
    }
}
```

Reading the file will require loading a line and splitting at each delimiter (comma).

PRO/CON OF CSV FILES

CSV is literally the *simplest possible thing* that we can do.

As a file format, it has some **advantages**:

- Programming languages can easily work with CSV files (they're just text!)
- It's pretty space efficient.
- It's human-readable. Kind-of.

However, CSV comes with some big **disadvantages**:

- It doesn't work very well if your data contains the delimiter (e.g. a comma in your city field).
- It assumes a **fixed structure** and doesn't handle variable length records.
- It's hard to read! **There is no semantic information** to make sense of it. (i.e., there is no simple way to interpret the structure, no schema file format).
- It doesn't work for complex, multi-dimensional data. e.g. Customer transactions.

STRUCTURED DATA FORMATS: XML

Extensible Markup Language (XML) is a markup language that **designed** for data storage and transmission.

Defined by the World Wide Web Consortium's XML specification, it was the first major standard for markup languages. It's structurally similar to HTML, with a focus on data transmission (vs. presentation).

Structure consists of pairs of tags that enclose data elements. Attributes can be added.

```
<name>Jeff</name>
```

```
This is a caption</img>
```

You can have a schema that describes the data structure! You can validate data.

XML EXAMPLE

Example of a music collection **structured in XML**¹.

```
<catalog>
  <album>
    <title>Empire Burlesque</title>
    <artist>Bob Dylan</artist>
    <company>Columbia</company>
    <price>10.90</price>
    <year>1985</year>
  </album>
  <album>
    <title>Innervisions</title>
    <artist>Stevie Wonder</artist>
    <company>The Record Plant</company>
    <price>9.90</price>
    <year>1973</year>
  </album>
</catalog>
```

Album is a record, containing
Fields for title, artist etc.

Notice the opening and closing tags. If
XML looks like HTML, that's because
they're both descended from a
common ancestor language, SGML.

1. Check out these albums!

READING/WRITING OBJECTS TO XML

Customer
- cust_id: Integer
+ name: String
+ city: String

```
data class Customer (val cust_id: Int, val name: String, val city: String)
```

```
val customers = List<Customer>() // list of customers
customers.add(Customer(1001, "John Hall", "New York"))
customers.add(Customer(1002, "Allison Barnett", "Waterloo"))
customers.add(Customer(1003, "John McAfee", "Delphi"))
```

```
File("output.csv").open("w").use {
    it.write("<customers>")
    for (customer in customers) {
        it.write("<customer>")
        it.write("<cust_id>${customer.cust_id}</cust_id>")
        it.write("<name>${customer.name}</name>")
        it.write("<city>${customer.city}</city>")
    }
    it.write("</customers>")
}
```

Reading the file will require a very complex parser e.g., Stax for Java.

PROS/CONS OF XML

XML structures textual data. The use of tags, and the optional use of a schema file, means that we can formally define the semantic structure of our data!

This provides some major advantages compared to CSV.

- Can rely on structure to infer the meaning of data.
- You can nest elements e.g., collections of records.

However, XML is rarely used except in legacy systems. Why?

- Tags “bloat” the data, which results in excessive space requirements.
- Practically impossible to parse without a complex library.

STRUCTURED DATA FORMAT: YAML

YAML Ain't Markup Language ([YAML](#)) is a data serialization language. It's easy for humans to read, and it's commonly used for configuration.

- Three dashes: start of YAML document
- Key: value pairs
- Lists: dash for each element

Thoughts on human-readable formats

- Human readable! Used extensively for **config files**.
- Indentation used for structure; difficult to parse manually.
- Not as widely supported as other formats :/

Example from <https://www.cloudbees.com/>

```
---
doe: "a deer, a female deer"
ray: "a drop of golden sun"
pi: 3.14159
xmas: true
french-hens: 3
calling-birds:
  - huey
  - dewey
  - louie
  - fred
xmas-fifth-day:
  calling-birds: four
  french-hens: 3
  golden-rings: 5
  partridges:
    count: 1
    location: "a pear tree"
  turtle-doves: two
```

STRUCTURED DATA FORMAT: JSON

JSON (JavaScript Object Notation) is an open standard file format, and data interchange format that's commonly used on the web.

It's based on JavaScript object notation but is language independent. It was standardized in 2013 as ECMA-404.

JSON has a much simpler syntax compared to XML or YAML.

- Data elements consist of name/value pairs
- Fields are separated by commas
- Curly braces hold objects
- Square brackets hold arrays

JSON is preferred for communications, data persistence. It is **widely supported** by existing programming languages.

JSON EXAMPLE

```
{ "catalog":  
  {  
    "albums": [  
      {  
        "title": "Empire Burlesque",  
        "artist": "Bob Dylan",  
        "company": "Columbia",  
        "price": "10.90",  
        "year": "1988"  
      },  
      {  
        "title": "Innervision",  
        "artist": "Stevie Wonder",  
        "company": "The Record Plant",  
        "price": "9.90",  
        "year": "1973"  
      }  
    ]  
  }  
}
```

Album is a record, containing albums.
[] denotes an array, { } encloses an object

No tags, just keys and values! Much easier
to read, since it's all meaningful data.

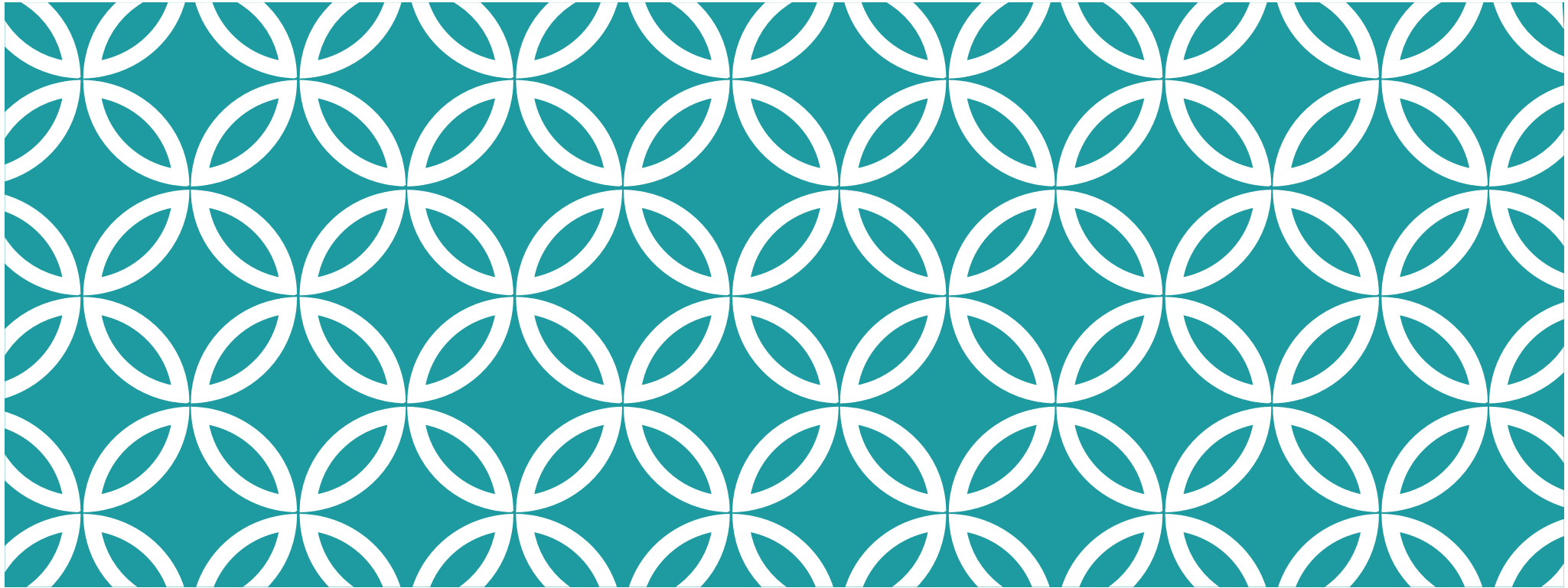
Condensing closing tags makes JSON easier to read.

```
{ "employees":[
  { "first":"John", "last":"Zhang", "dept":"Sales" },
  { "first":"Anna", "last":"Smith", "dept":"Engineering" }
]
```

Compare this to the corresponding XML:

```
<employees>
  <employee><first>John</first> <last>Zhang</last> <dept>Sales</dept></employee>
  <employee><first>Anna</first> <last>Smith</last> <dept>Eng.</dept></employee>
</employees>
```

This is a small record. I had to remove fields to fit the slide.



SERIALIZATION

CS 346: Application
Development

SERIALIZING DATA

We have objects in memory. We'd like to convert them to JSON to save them. How can we accomplish this?

Serialization is a mechanism to convert a data object to a useful format that you can save/stream or otherwise manipulate outside of your program.

- **Serialization:** save your object to a stream (file or network).
- **Deserialization:** instantiate an object from your stream (file or network).

```
class Emp(var name: String, var id:Int) : Serializable {}  
var file = FileOutputStream("datafile")  
var stream = ObjectOutputStream(file) // binary format  
  
var ann = Emp(1001, "Anne Hathaway", "New York")  
stream.writeObject(ann) // serialize to a file
```

READING/WRITING OBJECTS TO JSON

We can use serialization to convert objects directly into JSON format!

- Serialize data objects into JSON strings.
- Save those strings (aka text) to disk/stream over a network/save to a database.
- Deserialization can be used to reverse the process (convert stream → object in mem)

To add serialization support, install these plugins/dependencies (newest versions):

```
plugins {  
    id 'org.jetbrains.kotlin.plugin.serialization' version '1.9.10'  
}  
dependencies {  
    implementation "org.jetbrains.kotlinx:kotlinx-serialization-json:1.5.1"  
}
```



```

@Serializable
data class Project(val name: String, val owner: Account, val group: String)

@Serializable
data class Account(val userName: String)

val moonshot = Project("Moonshot", Account("Jane"), "R&D")
val cleanup = Project("Cleanup", Account("Mike"), "Maintenance")

val string = Json.encodeToString(listOf(moonshot, cleanup))
// [ {"name":"Moonshot","owner":{"userName":"Jane"},"group":"R&D"},
//   {"name":"Cleanup","owner":{"userName":"Mike"},"group":"Maintenance"} ]

val projectCollection = Json.decodeFromString<List<Project>>(string)
// [ Project(name=Moonshot, owner=Account(userName=Jane), group=R&D),
//   Project(name=Cleanup, owner=Account(userName=Mike), group=Maintenance) ]

```

JSON AS A STANDARD

We'll use JSON for storing and transmitting data anywhere that we need it.

Structure + data means that we can process it consistently.

- We can easily convert JSON to/from object format

Easy to work with it! It's just a string.

- Read it, print out to the console
- Save in a text file, using standard File classes
- Saved to a database — *we will discuss further*
- Send over a network — *we will discuss in web services lecture*

Don't underestimate the value in being able to read your data in a debugger, or text editor as you're working with it. JSON being human-readable text is one of its biggest advantages as a data file format.