# Enabling a Data-informed Public Sector: An Introductory Course to BDTI Essentials

## Session 1: Data Access and Exploration

**DIGITAL EUROPE PROGRAMME**

## Welcome!
## The webinar will begin shortly.

**EC-BDTI-PILOTS@ec.europa.eu**

**DG CNECT**
Directorate-General for Communications
Networks, Content and Technology

**DG DIGIT**
Directorate-General
for Digital Services

European Commission

# Some housekeeping

- The Chat and Q&A widgets are located to the right of the presenter screen
- Please use the Chat during the session, but wait until the dedicated Q&A portion to post your questions
- Microphones and video are off by default, attendees cannot turn them on
- The session will be recorded and available shortly after the end of the session on our website

**Do**
- Participate in the chat during the session
- Ask questions during the Q&A portion at the end
- Keep your questions concise
- Give feedback after the session
- Join the dedicated discussion board

**Don't**
- Use chat as a watercooler (no personal conversations)
- Self-promote
- Disrespect anyone

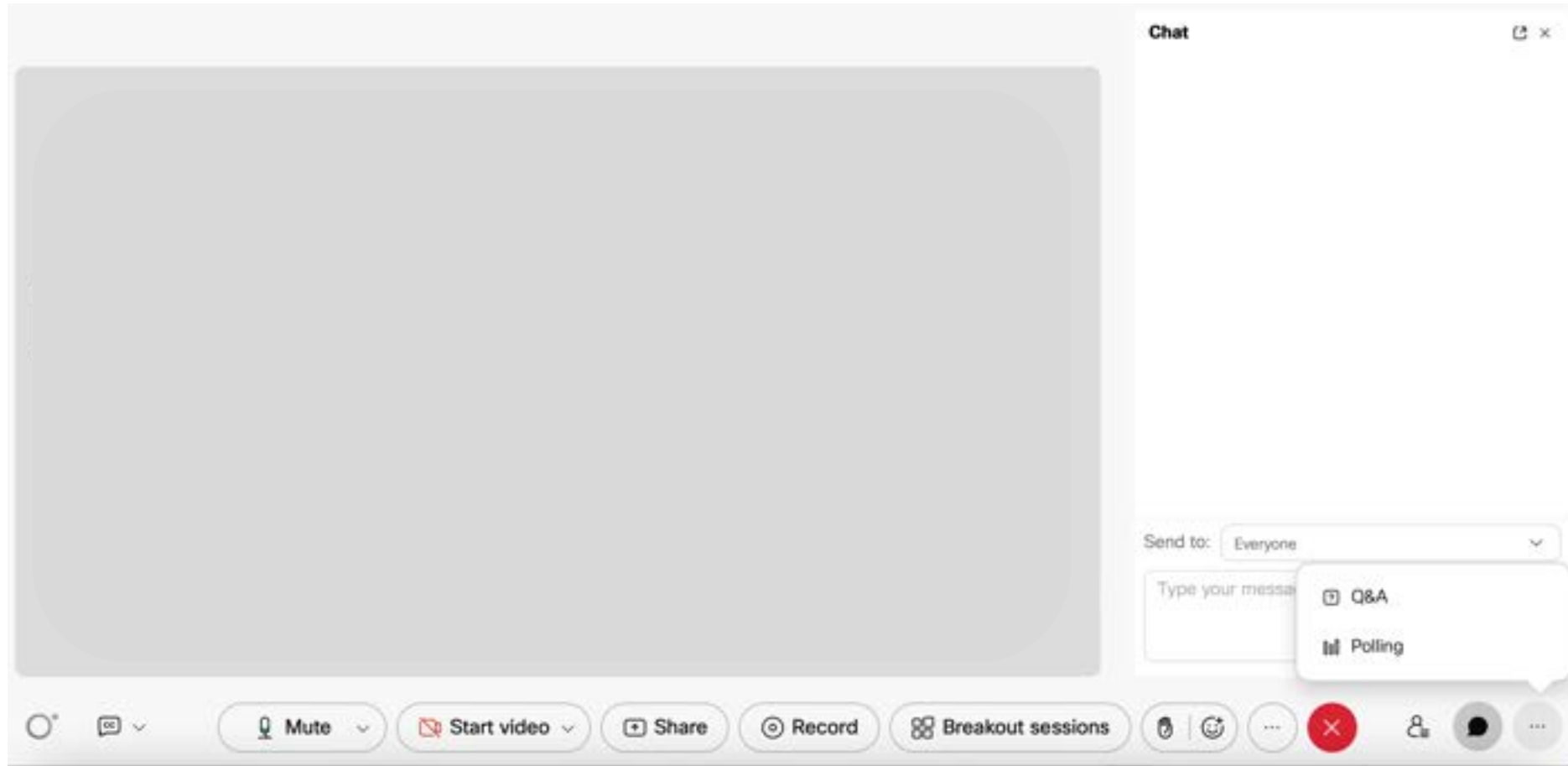Enabling a Data-informed Public Sector: An Introductory Course to BDTI Essentials

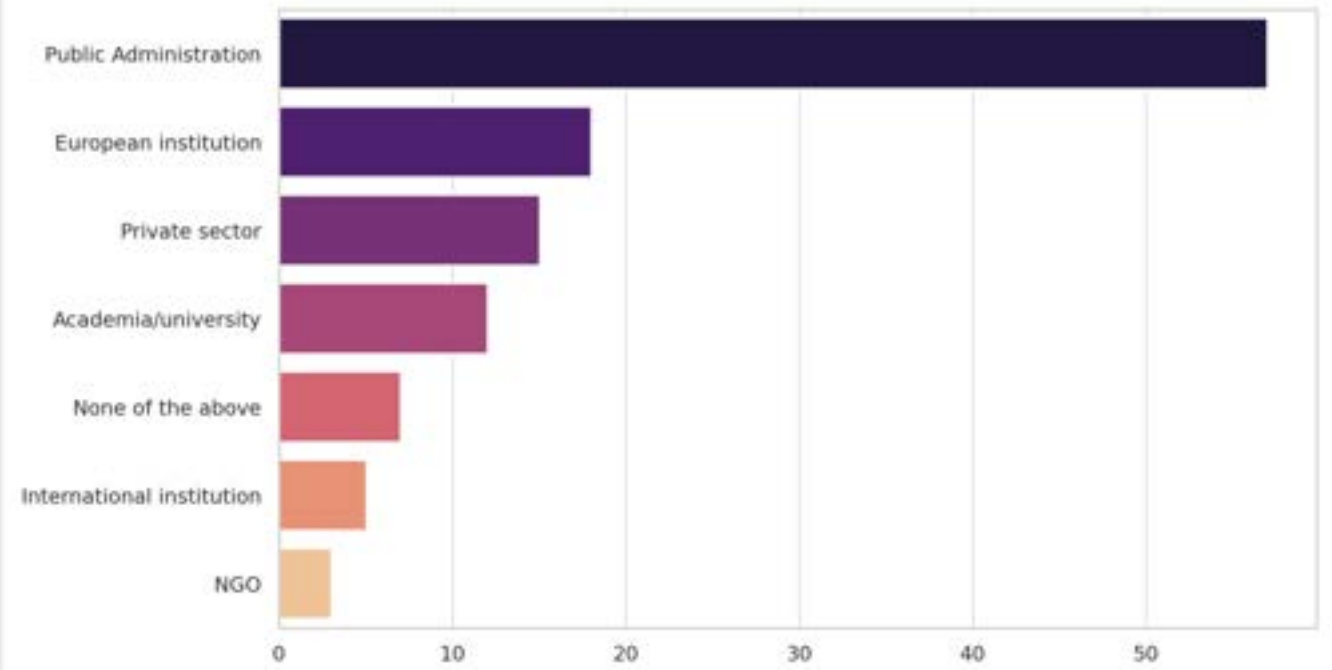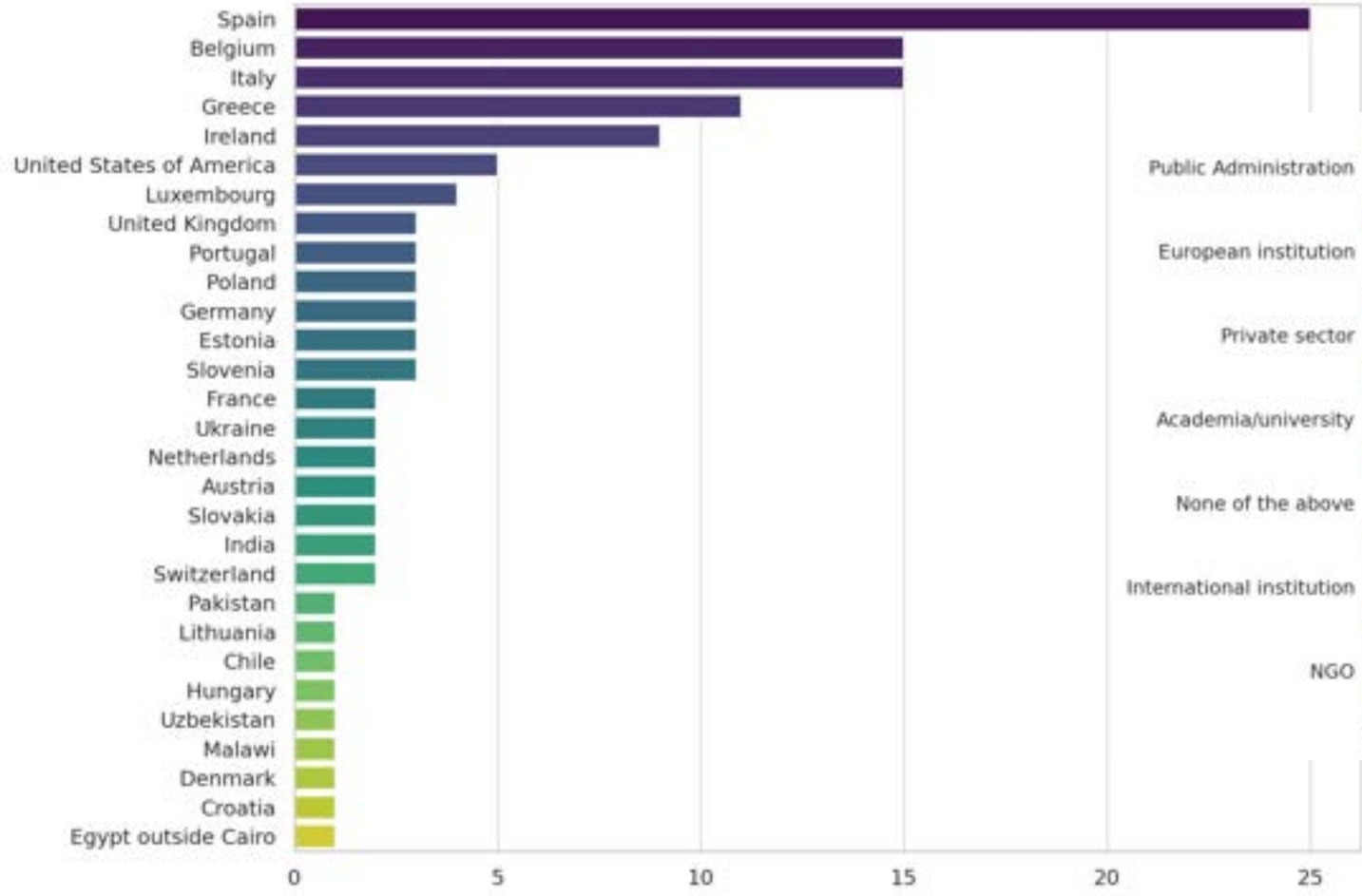Session 1: Data Access and Exploration

## Agenda

o Welcome [11:00-11:05]
o The Big Data Test Infrastructure in nutshell [11.05-11.10]
o Data.Europa.eu: where to find EU open data for your projects [11.10-11.15]
o Data access and Exploration [11:15-12:00]
o Q&A [12:00-12:15]

# Locate Chat and Q&A - Webex

# Our data enthusiasts ☺

# Intro

- The data and AI revolution: an increasing demand for professionals equipped with the skills to analyse and extract intelligence

- To draw attention to the new skills that the digital and green transition require, the European Commission has declared **2023 the European Year of Skills** (from May 2023 to May 2024)

- The Digital Europe Programme: Supports public sector in acquiring the necessary skills to extract value from data and new technologies

- The Big Data Test Infrastructure: A powerful tool offered to all EU public administrations

- This course is about **getting familiar with Data analytics**, building a concrete use case

- After this course, you will be ready to apply for BDTI and build a public sector data use case using the platform

# The Data Analysis Process is like cooking a dish

Business challenge/question/problem statement

**Ingredients + tools**

1. **Find** the data you need → *gather the right ingredients, good quality*
2. **Get, clean** and **prepare** your data → *slice and dice*
3. **Analyse** your data → *mix ingredients together and try different combinations*
4. **Present** the results and **create** knowledge → *serve and consume*

**Data→ Information → Knowledge**

# What is the Big Data Test Infrastructure (BDTI) ?

Not **only** for big data, for **public sector in general (open data)**

You have the key ingredients (datasets),
we equip you with the best **open-source tool**
to create amazing recipes for **public good.**

**Six months free of charge service**
for EU public administrations *

**Ready-to-use
data analytics stack** and
support

**Cloud platform** based on
**open-source** tools

* The cost of the pilot project must fit within the funding boundaries of the BDTI pilot budget

# Open-source tools to support your data journey

# Who is the Big Data Test Infrastructure (BDTI)  for?

**European Public Administrations**
All European Public Administrations at **local, regional and national level** can independently apply for a BDTI pilot project

**Ecosystem with academia and private sector**
Academia, spin-off, startups can apply for pilot projects once there is a **clear collaboration** with a Public Administration which will be the main point of contact for the project (**Master/PhD, GovTech startups**)

European Commission

# data.europa.eu

## Our services in a nutshell

**Data**

Providing access to free public data resources across Europe via a single platform (the portal).

**Academy**

Supporting the EU institutions, EU Member States and European countries to publish metadata of their open data in a harmonised manner.

**Community**

Organising open data competitions and conferences; communicating via social media and newsletters.

**Publications**

Assessing open data maturity in Europe; providing reports, studies and training via data.europa academy.

# Data.europa.eu as data hub

- **More than 1.6 million datasets**, grouped in 183 data catalogues

- **Navigate or search to get to the data** or catalogue you are looking for

- Benefit from **many filters**

- **Metadata translations** in all EU languages, machine translations for other text

- **Download and transform CSV files** automatically in many different formats

- Get **quick visualisations** for geo datasets

- Get **feedback for the metadata quality** on how to improve it

# Data.europa.eu: get inspired

- Consult our list of **use cases** (more than 900 examples)

- Follow re-users through the **use case observatory**

- Read our **data stories**

# Learning Overview

# Structure of sessions

| Session | Topic | Duration |
|---------|-------|----------|
| Webinar 1 | Data Access and Exploration | 45 min |
| Webinar 2 | Data Cleaning and Transformation | 75 min |
| Webinar 3 | Data Blending and Storage | 75 min |
| Webinar 4 | Analytics | 75 min |
| Webinar 5 | Advanced Module: Gathering Data from the Web and Geo Visualisation | 75 min |

- Slides, hands-on exercises, and solutions will be provided on the webinar
- Gitlab repository https://code.europa.eu/bdti/bdti-essentials-course

# Target Audience

- Assume beginner level (spreadsheet users): KNIME

- Advanced users: Jupyter lab and R-studio

# From tables and cells to process thinking

| Product ID | Store | Category | # Ordered items |
|------------|--------|-------------|-----------------|
| P 1 | Online | Clothing | 2 |
| P 2 | Onsite | Home | 3 |
| P 3 | Onsite | Clothing | 1 |
| P 4 | Online | Clothing | 5 |
| P 5 | Online | Electronics | 7 |
| P 6 | Online | Electronics | 5 |

Data Sources
Data Processing
Data Insights
Reporting

Data Professional

# The Data Process

Access the data from multiple sources

Clean and transform the data to the correct shape

Make sense of data, extract insights and visualise them in a report

| Data Access | → | Data Cleaning and Transformation | → | Data Blending and Export | → | Analytics: aggregation, visualisation, reporting | → | Advanced Analytics |

Merge the data and store it in the desired format (end of ETL process)

Enhance analytics by gathering data from the internet

# Use case: "Investigating Funding of Green Energy projects on EU Universities"

## Use Case

Zoi and her team, work in the Educational Department of a highly polluted EU region.

- Their **purpose** is to identify potential ***"energy partnerships"*** between universities on a national and international level to address pollution issues.

- To achieve this purpose, they track **two main sets data**. The first concerns the research and innovation funding of ***energy*** projects at universities across EU member states. The second dataset contains CO2 emissions of these countries. Moreover, they use a matching table to relate these datasets to each other.

- Zoi reached out the BDTI team to help learn how to **create *insights*** from the data sets and produce a ***report*** accordingly for their department, which can be **easily updated** when new data is available.

# Breakdown of objectives

- Visualise Energy projects per Member State and for **Organizations** per **Country**

- **Identify and Visualise Member States with the Most and Least **CO2 Emissions**

- Identify **Energy projects** per Country per **year**

- Explore the Relationship Between CO2 Emissions and Energy Project Funding

- Locate on a **map** the Universities that got funding for Energy projects from Horizon

# Mapping objectives to data

- First, finding the correct data to achieve the objectives.

- Then, Zoi has to **explore** the data to make decisions according to the different objectives.

- Map her objectives with the data she has.

# Searching for the (open) data

# Data Access - Sources

- Zoi found open data at https://data.europa.eu/en, she searched: "Horizon 2020 cordis" and "Horizon 2021".

# Data Access - Sources

- The team also searched for CO2 emissions and they landed on "Our World in Data" website https://ourworldindat a.org/

# Horizon data

- Organizations dataset - *Visualise Energy projects per Member State and for* ***Organizations*** *per* ***Country***



| RowID | projectID Number (Inte... | projectAc... String | organisat... Number (Inte... | vatNumber String | name String | shortName String | SME String | activityTy... String | street String | postCode String | city String | country String |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Row0 | 879926 | EEN SACHSEN | 999913831 | ⏱ | INDUSTRIE-U... | IHK C | false | PUB | STRASSE DE... | 09111 | CHEMNITZ | DE |
| Row1 | 879926 | EEN SACHSEN | 994556036 | DE176093942 | ZTS-ZENTRU... | ZTS | true | OTH | INDUSTRIEST... | 01612 | GLAUBITZ | DE |
| Row2 | 879926 | EEN SACHSEN | 999797625 | ⏱ | IND...IE- U... | IHK DRESDEN | false | PUB | LANGER WEG... | 01239 | DRESDEN | DE |
| Row3 | 879926 | EEN SACHSEN | 999940409 | ⏱ | ...E... | HWK DRESDE... | false | PUB | AM LAGERPL... | | DRESDEN | DE |
| Row4 | 879926 | EEN SACHSEN | 999877844 | DE140857609 | ...E... | ⏱ | false | HES | STRASSE DE... | | Chemnitz | DE |
| Row5 | 879926 | EEN SACHSEN | 999914219 | ⏱ | INDUSTRIE-U... | IIHK-L | false | PUB | GOERDELERR... | 04109 | LEIPZIG | DE |
| Row6 | 879926 | EEN SACHSEN | 969246408 | ⏱ | HOCHSCHUL... | Hochschule Z... | false | HES | THEODOR-KO... | 02763 | Zittau | DE |
| Row7 | 879926 | EEN SACHSEN | 999917905 | DE141484057 | AGENTUR FU... | AGIL | false | PRC | LESSINGSTR... | 04109 | LEIPZIG | DE |
| Row8 | 740689 | HEIMDALL | 997710476 | IT015032900... | CENTRO INT... | FONDAZIONE... | false | REC | VIA ARMAND... | 17100 | SAVONA | IT |
| Row9 | 740689 | HEIMDALL | 940057847 | DK37222135 | FREDERIKSB... | FBBR | false | PRC | LOGISMOSE 3 | 3600 | FREDERIKSS... | DK |
| Row.. | 740689 | HEIMDALL | 999703438 | ESS0811001G | Departament ... | INT | false | PUB | Carrer Diputa... | 08009 | Barcelona | ES |
| Row.. | 740689 | HEIMDALL | 919049102 | IT136697210... | ASSOCIAZIO... | ⏱ | false | OTH | VIA BERNAR... | 00151 | Roma | IT |
| Row.. | 740689 | HEIMDALL | 999778322 | EL094149709 | SPACE HELL... | SPACE HELL... | true | PRC | MESSOGION ... | 153 41 | Aghia Parask... | EL |
| Row.. | 740689 | HEIMDALL | 999638739 | ESG62616586 | CENTRE TEC... | CTTC | false | REC | AVINGUDA C... | 08860 | Castelldefels ... | ES |
| Row.. | 740689 | HEIMDALL | 996569950 | FR44130005... | UNIVERSITE ... | UNISTRA | false | HES | RUE BLAISE P... | 67081 | Strasbourg | FR |
| Row.. | 740689 | HEIMDALL | 935977542 | ESQ0801980D | INSTITUT CA... | ICGC | false | PUB | PARC DE MO... | 08038 | Barcelona | ES |
| Row.. | 740689 | HEIMDALL | 997822026 | ESB24352296 | TECNOSYLVA... | ⏱ | true | PRC | Parque Tecno... | 24009 | Leon | ES |
| Row.. | 740689 | HEIMDALL | 999981731 | DE121965658 | DEUTSCHES ... | DLR | false | REC | LINDER HOHE | 51147 | Koln | DE |
| Row.. | 740689 | HEIMDALL | 917614763 | GB158535583 | SCOTTISH FI... | SCOTTISH FI... | false | PUB | SCOTTISH FI... | G72 7NA | CAMBUSLANG | UK |
| Row.. | 740689 | HEIMDALL | 989367991 | GB830703457 | AVANTI COM... | ⏱ | true | PRC | COBHAM HO... | EC4V 6EB | London | UK |
| Row.. | 740689 | HEIMDALL | 958896217 | ESG55536098 | FUNDACIO D'... | PAU COSTA F... | true | REC | AV. MOSSEN ... | 08552 | Taradell | ES |

# Horizon data

- Projects dataset - *Identify Energy **Projects** per Country per **year***

| ids<br>Number (inte... | legalBasi...<br>String | masterCa...<br>String | natures<br>String | objectives<br>String | rcns<br>Number (inte... | startDates<br>String | statuss<br>String | subCalls<br>String | titl...<br>String | topicss<br>String | totalCosts ↓<br>Number (double) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 633053 | H2020-Eurato... | EURATOM-Ad... | | A Roadmap t... | 193159 | 2014-01-01 | SIGNED | EURATOM-Ad... | Imple... | EURATOM | 1,329,649,211.5 |
| 807097 | H2020-EU.3.4. | H2020-IBA-C... | | The challeng... | 218781 | 2014-01-09 | CLOSED | H2020-IBA-C... | Large ... | CS2-GAM... | 247,9...2.58 |
| 807085 | H2020-EU.3.4. | H2020-IBA-C... | | Engines ITD ... | 216640 | 2014-01-03 | CLOSED | H2020-IBA-C... | Engine... | CS2-GAM... | 228,776,991.75 |
| 945583 | H2020-EU.3.4. | H2020-IBA-C... | | Main objectiv... | 231017 | 2020-01-01 | SIGNED | H2020-IBA-C... | Large ... | CS2-GAM... | 228,776,991.75 |
| 807083 | H2020-EU.3.4. | H2020-IBA-C... | | The Airframe ... | 216645 | 2014-01-02 | CLOSED | H2020-IBA-C... | AIRFR... | CS2-GAM... | 210,183,580.85 |
| 853989 | H2020-EU.3.1. | H2020-JTI-IM... | | The European... | 226660 | 2020-01-01 | SIGNED | H2020-JTI-IM... | EURO... | IMI2-2018. | 207,963,891 |
| 101007873 | H2020-EU.3.1. | H2020-JTI-IM... | | Current anti-t... | 236634 | 2021-06-01 | SIGNED | H2020-JTI-IM... | ACAD... | IMI2-2020... | 185,000,000 |
| 777389 | H2020-EU.3.1. | H2020-JTI-IM... | crisisPrepare... | Paediatric me... | 218657 | 2018-05-01 | SIGNED | H2020-JTI-IM... | conect... | IMI2-2016... | 182,018,216 |
| 737417 | H2020-EU.2.1... | H2020-ECSEL... | | R3-POWERUP... | 210524 | 2017-11-01 | SIGNED | H2020-ECSEL... | 300m... | ECSEL-20... | 180,318,047.5 |
| 681463 | H2020-EU.4.e. | H2020-Adhoc... | | COST is an in... | 198528 | 2015-05-01 | CLOSED | H2020-Adhoc... | COST ... | COST-Net... | 178,169,171 |
| 662338 | H2020-EU.2.1... | ECSEL-2014-2 | | The SeNaTe ... | 198669 | 2015-04-01 | CLOSED | ECSEL-2014-2 | Seven ... | ECSEL-02... | 177,731,624.33 |
| 945542 | H2020-EU.3.4. | H2020-IBA-C... | | The challeng... | 229352 | 2020-01-01 | SIGNED | H2020-IBA-C... | Fast R... | CS2-GAM... | 171,739,896.46 |
| 807081 | H2020-EU.3.4. | H2020-IBA-C... | | The Systems ... | 216639 | 2014-01-01 | CLOSED | H2020-IBA-C... | Syste... | CS2-GAM... | 158,177,923.26 |
| 881603 | H2020-EU.1.2. | H2020-SGA-F... | | This proposal... | 227202 | 2020-04-01 | SIGNED | H2020-SGA-F... | Graph... | SGA-FET-... | 150,000,000 |
| 945539 | | H2020-SGA-F... | | The last of fo... | 229984 | 2020-04-01 | SIGNED | H2020-SGA-F... | Huma... | SGA-FETF... | 150,000,000 |
| 101018100 | H2020-EU.3.1. | H2020-FPA-S... | crisisResponse | Funding of re... | 232261 | 2020-07-24 | SIGNED | H2020-SGA-S... | PROP... | SGA-SC1-... | 150,000,000 |
| 692522 | H2020-EU.2.1... | H2020-ECSEL... | | The TAKE5 pr... | 203403 | 2016-04-01 | CLOSED | H2020-ECSEL... | Techn... | ECSEL-15... | 149,882,181 |
| 101036970 | H2020-EU.3.3. | H2020-LC-GD... | | REFHYNE II w... | 237006 | 2021-10-01 | SIGNED | H2020-LC-GD... | Clean ... | LC-GD-2-2... | 148,956,405 |

# Horizon data

- EuroSciVoc dataset - *Identify **Energy** Projects per Country per year*

| RowID | projectID<br>String | euroSciVocCode<br>String | euroSciVocPath<br>String | euroSciVocTitle<br>String | |
|---|---|---|---|---|---|
| Row0 | 869855 | /29/101/555/1359 | /social sciences/sociolog... | automation | |
| Row1 | 869855 | /25/63/399 | /engineering and technol... | textiles | |
| Row2 | 869855 | /25/73/453/58525161 | /engineering and technol... | sensors | |
| Row3 | 869855 | /25/73/453/459 | /engineering and technol... | robotics | |
| Row4 | 869855 | /29/93/47293307 | /social sciences/psychol... | ergonomics | |
| Row5 | 870148 | /25/63/409 | /engineering and technol... | colors | |
| Row6 | 883285 | /21/33/121/44109686/5... | /medical and health scien... | ebola | |
| Row7 | 883285 | /21/33/137/133/9678651 | /medical and health scien... | pandemics | |
| Row8 | 883285 | /21/33/121/44109686/7... | /medical and health scien... | influenza | |
| Row9 | 883285 | /29/91/523/1313 | /social sciences/econom... | employment | |
| Row... | 871385 | /23/47/297 | /natural sciences/comput... | artificial intelligence | |
| Row... | 871385 | /23/47/307 | /natural sciences/comput... | software | |
| Row... | 871385 | /29/101/555/1359 | /social sciences/sociolog... | automation | |
| Row... | 871385 | /29/97/67681549/64785... | /social sciences/political ... | revolutions | |
| Row... | 871385 | /23/47/295 | /natural sciences/comput... | computer security | |
| Row... | 682402 | /23/51/359/1073/1655 | /natural sciences/mathe... | linear algebra | |
| Row... | 682402 | /23/51/359/1067/81178... | /natural sciences/mathe... | functional analysis | |

# CO2 emissions data

- Identify and Visualise Member States with the Most and Least **CO2 Emissions**
- One single file with 79 columns, and information about all the countries in the world.
  *But Zoi only needs EU members!*

| | country | year | iso_code | population | gdp | cement_co2 | cement_co2_per_capita |
|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 1850 | AFG | 3752993.0 | NaN | NaN | NaN |
| 1 | Afghanistan | 1851 | AFG | 3767956.0 | NaN | NaN | NaN |
| 2 | Afghanistan | 1852 | AFG | 3783940.0 | NaN | NaN | NaN |
| 3 | Afghanistan | 1853 | AFG | 3800954.0 | NaN | NaN | NaN |
| 4 | Afghanistan | 1854 | AFG | 3818038.0 | NaN | NaN | NaN |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 50593 | Zimbabwe | 2017 | ZWE | 14751101.0 | 2.194784e+10 | 0.469 | 0.032 |
| 50594 | Zimbabwe | 2018 | ZWE | 15052191.0 | 2.271535e+10 | 0.558 | 0.037 |
| 50595 | Zimbabwe | 2019 | ZWE | 15354606.0 | NaN | 0.570 | 0.037 |
| 50596 | Zimbabwe | 2020 | ZWE | 15669663.0 | NaN | 0.570 | 0.036 |
| 50597 | Zimbabwe | 2021 | ZWE | 15993525.0 | NaN | 0.570 | 0.036 |

50598 rows × 79 columns

# Use case: Approach

- Use the Big Data Test Infrastructure (BDTI) initiative from the DEP program (European Commission)

- Received specialised training for effective data analysis

- Leverage open data from different sources

  o EU Funding Program Data – Horizon2020/1:
    https://data.europa.eu/data/datasets/cordish2020projects?locale=en
  o CO2 Emissions Data (Our World in Data)
    https://ourworldindata.org/co2-emissions
  o Collect external data if needed

# Outputs according to Zoi goals

**The report/dashboard needs to be informative!**

- Insightful plots

- Tables with a rank of countries according to CO2 emissions

- Bar chart with CO2 emissions

- Line plot with both Horizon 2020 and 2021

- Locate on a map the Universities that got funding for Energy projects from Horizon



## Example Dashboard Report with BDTI Framework

**Total Energy Projects per Country**

**CO2 emissions perCountry**

**Projects per country by years**

**Zoom per Country**

| countryNa... | projectID | cumulativ... | rank |
|---|---|---|---|
| Germany | 2142 | 93,290.508 | 1 |
| United Kingdom | 1372 | 78,508.719 | 2 |
| France | 1753 | 39,105.766 | 3 |
| Poland | 283 | 28,196.094 | 4 |

# But...

- The flipside of the coin is the technical side
- Raw file (CSV) vs table

# But…

- File formats
- Raw (JSON) vs table

```
, {
  "acronym" : "INVEST",
  "contentUpdateDate" : "2022-08-02 17:23:54",
  "ecMaxContribution" : 3998440,
  "ecSignatureDate" : "2022-06-01",
  "endDate" : "2025-03-31",
  "frameworkProgramme" : "HORIZON",
  "fundingScheme" : "HORIZON-CSA",
  "grantDoi" : "10.3030/101058526",
  "id" : 101058526,
  "legalBasis" : "HORIZON.2.4",
  "masterCall" : "HORIZON-CL4-2021-RESILIENCE-01",
  "nature" : "",
  "objective" : "The introduction of the Innovation Radar is a great step f
  "rcn" : 238957,
  "startDate" : "2022-04-01",
  "status" : "SIGNED",
  "subCall" : "HORIZON-CL4-2021-RESILIENCE-01",
  "title" : "Matching investors and EU funded innovations",
  "topics" : "HORIZON-CL4-2021-RESILIENCE-01-27",
  "totalCost" : 3998440
, {
  "acronym" : "ReMFra",
  "contentUpdateDate" : "2022-08-02 17:20:44",
```

| RowID | projectID Number (inte… | projectAc… String | organisat… Number (inte… | vatNumber String | name String | shortName String | SME String | activityTy… String | street String | postCode String | city String | country String |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Row0 | 879926 | EEN SACHSEN | 999913831 | | INDUSTRIE-U… | IHK C | false | PUB | STRASSE DE… | 09111 | CHEMNITZ | DE |
| Row1 | 879926 | EEN SACHSEN | 994556036 | DE176093942 | ZTS-ZENTRU… | ZTS | true | OTH | INDUSTRIEST… | 01612 | GLAUBITZ | DE |
| Row2 | 879926 | EEN SACHSEN | 999797625 | | INDUSTRIE- U… | IHK DRESDEN | false | PUB | LANGER WEG… | 01239 | DRESDEN | DE |
| Row3 | 879926 | EEN SACHSEN | 999940409 | | HANDWERKS… | HWK DRESDE… | false | PUB | AM LAGERPL… | 01099 | DRESDEN | DE |
| Row4 | 879926 | EEN SACHSEN | 999877844 | DE140857609 | TECHNISCHE… | | false | HES | STRASSE DE… | 09111 | Chemnitz | DE |
| Row5 | 879926 | EEN SACHSEN | 999914219 | | INDUSTRIE-U… | IIHK-L | false | PUB | GOERDELERR… | 04109 | LEIPZIG | DE |
| Row6 | 879926 | EEN SACHSEN | 969246408 | | HOCHSCHUL… | Hochschule Z… | false | HES | THEODOR-KO… | 02763 | Zittau | DE |
| Row… | 879926 | EEN SACHSEN | 999917905 | DE141484057 | AGENTUR FU… | AGIL | false | PRC | LESSINGSTR… | 04109 | LEIPZIG | DE |
| Row… | 40689 | HEIMDALL | 997710476 | IT015032900… | CENTRO INT… | FONDAZIONE… | false | REC | VIA ARMAND… | 17100 | SAVONA | IT |
| Row… | 40689 | HEIMDALL | 940057847 | DK37222135 | FREDERIKSB… | FBBR | false | PRC | LOGISMOSE 3 | 3600 | FREDERIKSS… | DK |
| Row… | 40689 | HEIMDALL | 999703438 | ESS0811001G | Departament … | INT | false | PUB | Carrer Diputa… | 08009 | Barcelona | ES |
| Row… | 40689 | HEIMDALL | 919049102 | IT136697210… | ASSOCIAZIO… | | false | OTH | VIA BERNAR… | 00151 | Roma | IT |
| Row… | 40689 | HEIMDALL | 999778322 | EL094149709 | SPACE HELL… | SPACE HELL… | true | PRC | MESSOGION… | 153 41 | Aghia Parask… | EL |
| Row… | 740689 | HEIMDALL | 999638739 | ESG62616586 | CENTRE TEC… | CTTC | false | REC | AVINGUDA C… | 08860 | Castelldefels … | ES |
| Row… | 740689 | HEIMDALL | 996569950 | FR44130005… | UNIVERSITE … | UNISTRA | false | HES | RUE BLAISE P… | 67081 | Strasbourg | FR |
| Row… | 740689 | HEIMDALL | 935977542 | ESQ0801980D | INSTITUT CA… | ICGC | false | PUB | PARC DE MO… | 08038 | Barcelona | ES |
| Row… | 740689 | HEIMDALL | 997822026 | ESB24352296 | TECNOSYLVA… | | true | PRC | Parque Tecno… | 24009 | Leon | ES |
| Row… | 740689 | HEIMDALL | 999981731 | DE121965658 | DEUTSCHES … | DLR | false | REC | LINDER HOHE | 51147 | Koln | DE |
| Row… | 740689 | HEIMDALL | 917614763 | GB158535583 | SCOTTISH FI… | SCOTTISH FI… | false | PUB | SCOTTISH FI… | G72 7NA | CAMBUSLANG | UK |
| Row… | 740689 | HEIMDALL | 989367991 | GB830703457 | AVANTI COM… | | true | PRC | COBHAM HO… | EC4V 6EB | London | UK |
| Row… | 740689 | HEIMDALL | 958896217 | ESG55536098 | FUNDACIO D'… | PAU COSTA F… | true | REC | AV. MOSSEN … | 08552 | Taradell | ES |

# But...

- Need to work with all files in one environment
- Enter the BDTI Framework!

**Session 1:**
**Data Access and Exploration**

# Session 1: Learning outcomes

**At the end of this session, you should be able to:**

- Access BDTI framework

- Navigate BDTI Service Catalog

- Launch various data analytic tools:
    - KNIME Analytics Platform,
    - R-studio,
    - Jupyter Notebook.

- Read different data file types into the analytic tools:
    - CSV,
    - Excel,
    - JSON.

- Complete basic data exploration in relation to the use case

# How to access the BDTI framework?

- To have access, public administrations need to **_apply on the official website_**

https://big-data-test-infrastructure.ec.europa.eu/apply-bdti_en

- Once you have access follow these steps:

    o Step 1 > Access the portal via the website

    o Step 2 > Navigate to the "Service Catalog" and launch the tools

    o Step 3 > Go to "My Services" and open the tool you created

# Tour of the BDTI framework

**Step 1:**

Once you applied succesfully,
Access the portal via the website

# Tour of the BDTI framework

**Step 2:**

Navigate to the "Service Catalog" and launch the desired tools

# Tour of the BDTI framework

Step 3:

Go to "My Services" and open the tool you created

# All these tools can be downloaded locally

As the nature of these tools is **open-source**, in case you are not a public administration interested in apply for the BDTI playground, participants can download them separately before getting access to the BDTI, that way you can follow along with the sessions and practice at home!

- KNIME Analytic Platform: https://www.knime.com/downloads

- R-studio: https://cran.r-project.org/

- Jupyter Notebook: https://jupyter.org/install

# What is KNIME Analytics Platform

- A tool to make sense of your data
  - Data analysis, Data Science, Data Engineering
  - Transformation, visualisation, reporting
- Open source
- Visual programming paradigm
  - No coding required

# Nodes and Workflows

- **NODES** perform tasks on data



- **Nodes** are combined to create WORKFLOWS

# Tour of the User Interface



Open an example workflow

Browse local space

Create new workflow

# Data Access

- The beginning of every data process
- Data can be stored in many ways
  - Locally
  - In different data format (.csv, .xls...)
  - On the cloud
  - In a database
  - ...
- We need a way to access all of them

DIGITAL EUROPE PROGRAMME

**Data Access** → Data Cleaning and Transformation → Data Blending and Export → Analytics: aggregation, visualisation, reporting → Advanced Analytics

European Commission

# Data Access in KNIME Analytics Platform



- Files
  - CSV, txt, Excel, Word, PDF
  - SAS, SPSS
  - XML, JSON, PMML
  - Images, texts, networks
- Databases
  - MySQL, PostgreSQL, Oracle
  - Theobald
  - any JDBC (DB2, MS SQL Server)
  - Amazon DynamoDB
- Other
  - Twitter, Google
  - Amazon S3, Azure Blob Store
  - Sharepoint, Salesforce
  - Kafka
  - REST, Web services

# Data Access in KNIME Analytics Platform

- Typically characterised by:
- Orange color
- By default no input ports, 1-2 output ports
- Many nodes for many data formats
- Support reading from different File Systems

**Excel Reader**

Status

Output port

Read file

Node label

# File types in Use Case

- Excel Files

**Excel Reader**



- CSV Files

**CSV Reader**



- JSON Files

**JSON Reader**

# Read Excel files

- Excel Reader - reads .xls and .xlsx file from Microsoft Excel
- Supports reading from multiple sheets

**Excel Reader**

**Read Excel Sheet Names**

# Excel Reader - Configuration



**Excel Reader**

File path

File system

Sheet specific settings

Preview

# Common Settings: File Path

- A path consists of three parts:
  - **Type**: Specifies the file system type - e.g., local, relative, mountpoint, custom URL or connected
  - **Specifier**: Optional string with additional file system specific information - e.g. relative to which location (knime.workflow, LOCAL mountpoint...)
  - **Path**: Specifies the location within the file system



- Examples:
  - (LOCAL, , C:\Users\username\Desktop)
  - (RELATIVE, knime.workflow, file1.csv)
  - (MOUNTPOINT, MOUNTPOINT_NAME, /path/to/file1.csv)
  - (CONNECTED, amazon-s3:eu-west-1, /mybucket/file1.csv)

# Read CSV files

- CSV – Comma-separated Values

# Read CSV files

- Reads either one or multiple .csv and .txt files

- Further tabs to
  - Select columns
  - Limit the rows
  - Handle quotes
  - Select encoding

**CSV Reader**



Read data.csv

Advanced settings

File system

File path

Basic settings

Preview

Help button

# Read CSV files



Filter, sort, rename and convert columns

Skip top or bottom rows. Useful if they contain unstructured data e.g. metadata

# Read JSON format

- Use the JSON Reader (or GET Request) node to get a JSON cell
- Use the JSON Path node to query the JSON file and extract parameters
  - Editor window simplifies construction of JSON queries by auto-generating them

# JSON Objects

- Stores data in a key-value pair format (e.g., "Category": "Electronics")
- Supports nested structures
  - JSON objects inside JSON objects

```
{
    "Category"  : "Electronics",
    "Product"   : "256GB SSD",
    "Price"     : "119.95",
    "Rating"    : "4.2"
}
{
    "Category"  : "Electronics",
    "Product"   : "Wireless Keyboard",
    "Price"     : "59.99",
    "Rating"    : "3.8"
}
{
    . . .
```

Key

Value

# JSON Path

- Extracts fields of interest from JSON
  - A particular entry
  - All occurrences with the same key



**JSON Path**

Add only the selected occurrence

Add all occurrences with the same key

Select field from preview

# Ungroup

- Values for a JSON Path node are returned as collections
  - Multiple values in a cell
- Ungroup node converts collection cells to rows of data

# KNIME Analytic Platform

- Sign in with your credentials stored in "My Data"

# R-studio

- Sign in with your credentials stored in "My Data"



## Use Case

Some colleagues from Zoi are **more advanced** and know how to code in R-studio, they want to know how to use this tool within the BDTI framework to integrate and collaborate with the rest of the team.

# Jupyter lab - Python

- Sign in with your credentials stored in "My Data"



## Use Case

Some colleagues from Zoi are **more advanced** and know how to code in Python, they want to know how to use this tool within the BDTI framework to integrate and collaborate with the rest of the team.

# Summary

- **Now you should be able to:**

- Access BDTI framework

- Navigate BDTI Service Catalog

- Launch various data analytic tools
  - KNIME Analytics Platform, R-studio, Jupyter Notebook.

- Read different data file types into the analytic tools:
  - CSV, Excel, JSON.

- Complete basic data exploration in relation to the use case

# Next session: 16/02/24
## 11:00 AM - 12:15 PM CEST

| Data Access | → | **Data Cleaning and Transformation** | → | Data Blending and Export | → | Analytics: aggregation, visualisation, reporting | → | Advanced Analytics |

https://big-data-test-infrastructure.ec.europa.eu/resources/courses-and-training/bdti-essentials-course_en

# Practice, practice and practice!



- Slides, hands-on exercises, and solutions will be provided on the webinar Gitlab repository

https://code.europa.eu/bdti/bdti-essentials-course

# Practice, practice and practice!

- Slides, hands-on exercises, and solutions will be provided on the webinar Gitlab repository

https://code.europa.eu/bdti/bdti-essentials-course

Please check all the material in our open source repo:
https://code.europa.eu/bdti/bdti-essentials-course

**Q&A**

# Your BDTI journey starts here

Congrats! You are on the first step to data-driven innovation.

**BDTI Essentials Course** → **Brainstorm your data project** → **Apply for BDTI Pilot-light process**

Apply for BDTI: https://big-data-test-infrastructure.ec.europa.eu/apply-bdti_en

# Course discussion board

# Get in touch and follow the BDTI activities

**Are you working for a public administration in need of infrastructure for data analytics?**

**EC-BDTI-PILOTS@ec.europa.eu**

Visit
BDTI's
website
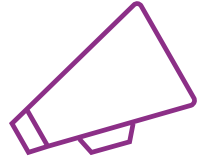
Subscribe to
BDTI's
newsletter

Subscribe to
BDTI's
Joinup

Register for BDTI
Essentials online
course

https://big-data-test-infrastructure.ec.europa.eu/

# Thanks for attending, see you in 2 weeks!

**Are you working for a public administration interested in testing a data analytics playground?**

**Schedule a 1-to-1 with the BDTI team: EC-BDTI-PILOTS@ec.europa.eu**

Don't miss the next session.
Register at the QR code below.

Questions? Comments? Ideas?
Join the dedicated discussion board.

https://big-data-test-infrastructure.ec.europa.eu/