# A Language Data Space for Europe

**Philippe Gelin**

**Head of Sector Multilingualism**
Data Directorate
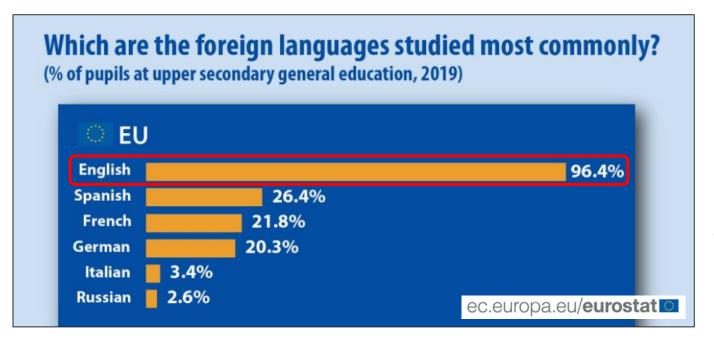Directorate-General for Communications Networks, Content and Technology
European Commission

*My Europe, My Language*

Virtual Berlin, 15th November 2021 – 9:15 – 9:40
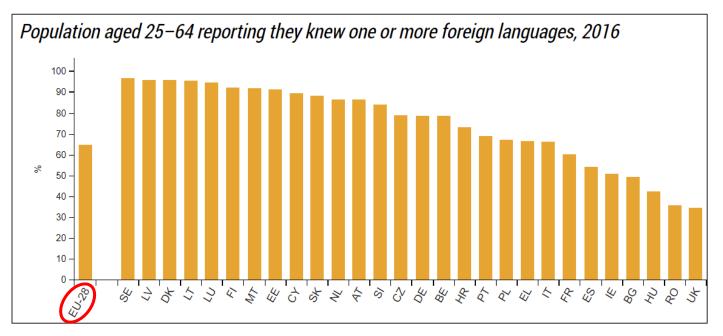
# DIGITAL Background Information:
# Language Learning

## Which are the foreign languages studied most commonly?
(% of pupils at upper secondary general education, 2019)

**EU**

| Language | Percentage |
|----------|-----------|
| English | 96.4% |
| Spanish | 26.4% |
| French | 21.8% |
| German | 20.3% |
| Italian | 3.4% |
| Russian | 2.6% |

ec.europa.eu/**eurostat**

In 2019, **96 %** of pupils in upper secondary general education in the EU were **learning English** as a **foreign language**.
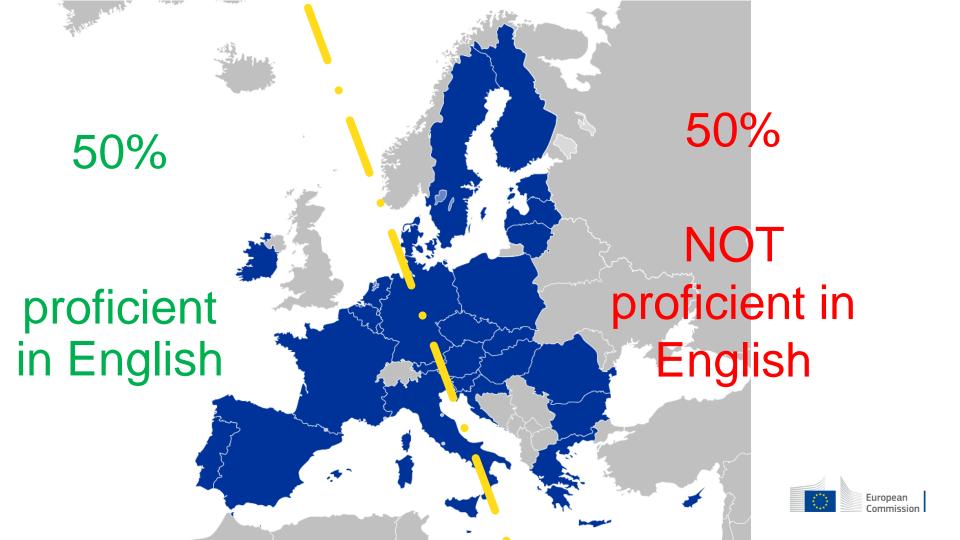
Education and training in the EU - facts and figures

Which are the foreign languages studied most commonly
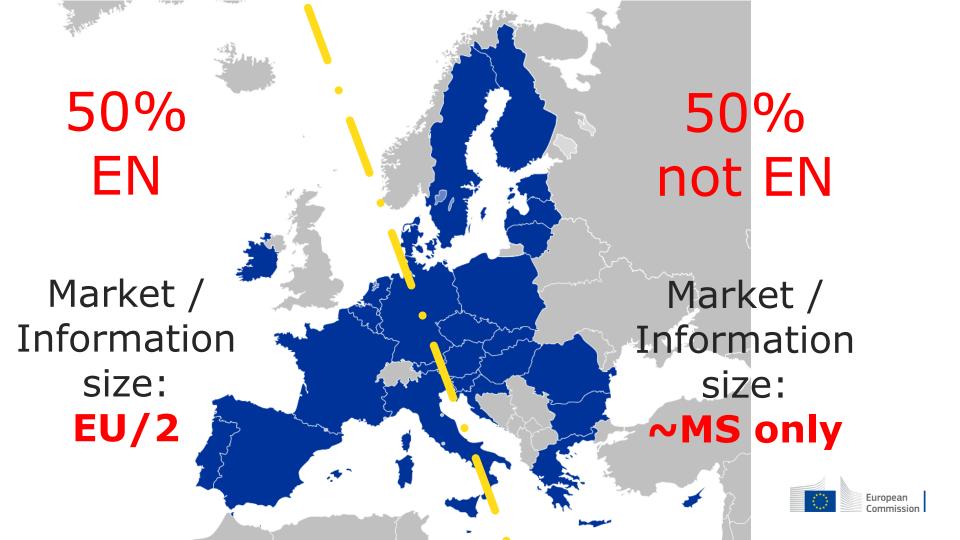
European Commission

# DIGITAL Background Information:
# Foreign Language



Population aged 25–64 reporting they knew one or more foreign languages, 2016

50%

proficient in English

50%

NOT proficient in English

European Commission

50%
EN

Market /
Information
size:
**EU/2**

50%
not EN

Market /
Information
size:
**~MS only**

50%
EN

50%
not EN

Market /
Information
size:
**EU/2**

Market /
Information
size:
**~MS only**

**US, CN**

European
Commission

50%
25+ Lang

Culture

US, CN

Market /
Information
size:
EU+

50%
25+ Lang

Market /
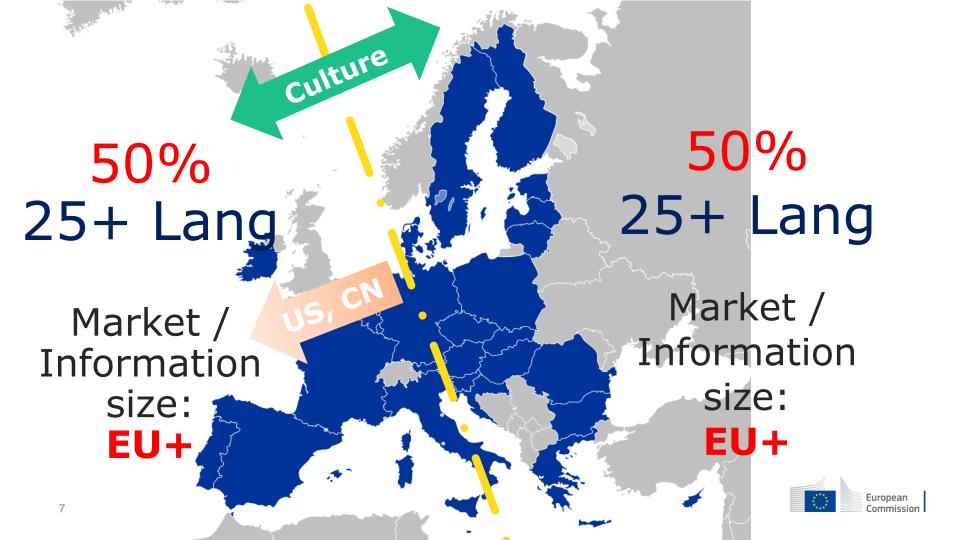Information
size:
EU+

European Commission

# DIGITAL in a Nutshell:
## The Programme *at a Glance*



Digital Europe Programme

### VISION
Digital Decade and digital targets for 2030

### FOCUS
Building the **strategic digital capacities** of the Union and facilitating the wide **deployment** of digital technologies

### GOAL
**Bridging the gap** between digital technology research and market deployment

European Commission

# DIGITAL in a Nutshell:
## The *Synergies*

# DIGITAL in a Nutshell:
# The Big Tickets
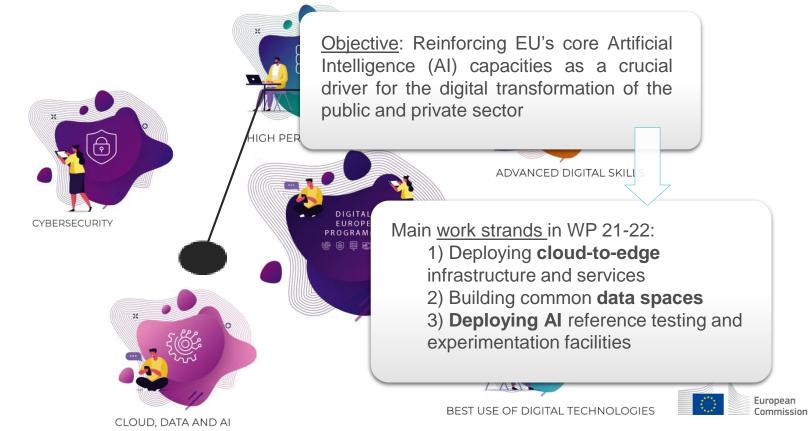
HIGH PERFORMANCE COMPUTING

ADVANCED DIGITAL SKILLS

CYBERSECURITY

DIGITAL EUROPE PROGRAMME

CLOUD, DATA AND AI

BEST USE OF DIGITAL TECHNOLOGIES

European Commission

# DIGITAL in a Nutshell:
# The Cloud, Data and AI

CYBERSECURITY

HIGH PER...

ADVANCED DIGITAL SKILLS

DIGITAL EUROPE PROGRAM...

CLOUD, DATA AND AI

BEST USE OF DIGITAL TECHNOLOGIES

Objective: Reinforcing EU's core Artificial Intelligence (AI) capacities as a crucial driver for the digital transformation of the public and private sector

Main work strands in WP 21-22:
1) Deploying **cloud-to-edge** infrastructure and services
2) Building common **data spaces**
3) **Deploying AI** reference testing and experimentation facilities

European Commission

# DATA SPACES Background Information:
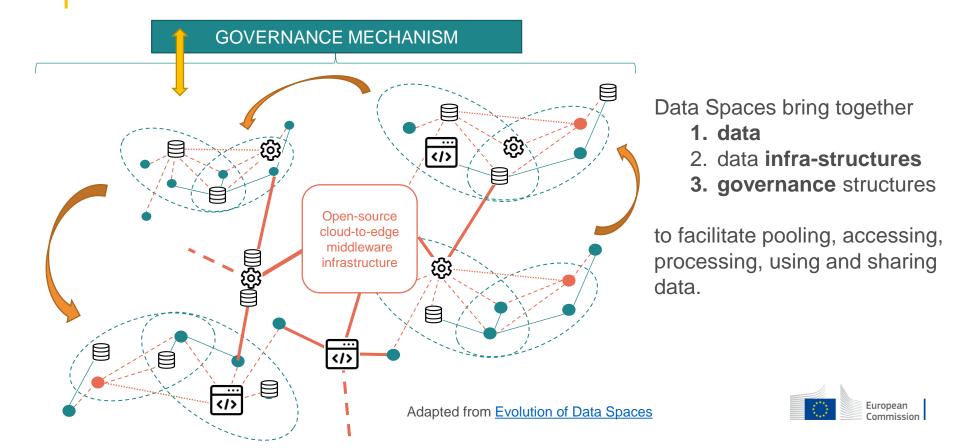# European Data Strategy

*High Impact Project*: Data spaces should foster an ecosystem (of companies, civil society and individuals) creating new products and services based on more accessible data.
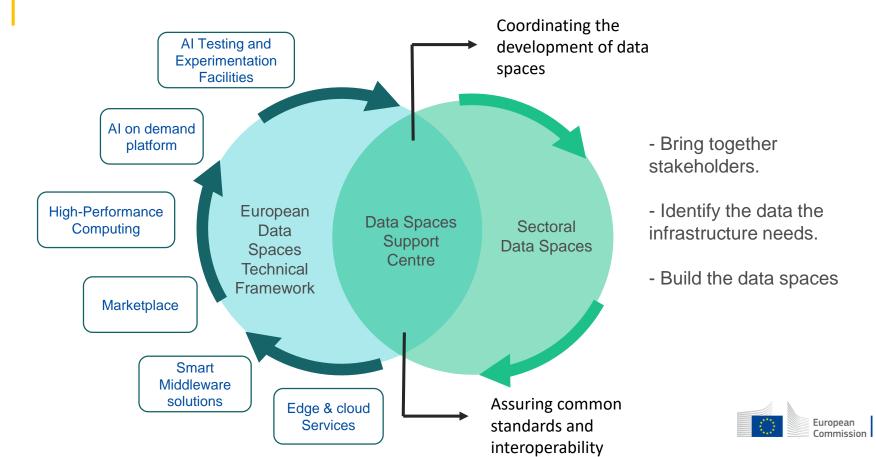
DATA SPACES

- EU-wide, common, interoperable;

- Including:
  - the deployment of data-sharing tools and platforms;
  - the creation of data governance frameworks;

- Improving the availability, quality and re-usabiliy of data– both in domain-specific settings and across sectors.

A European strategy for data

# DATA SPACES in a Nutshell (1):



GOVERNANCE MECHANISM

Open-source cloud-to-edge middleware infrastructure

Data Spaces bring together
1. **data**
2. data **infra-structures**
3. **governance** structures

to facilitate pooling, accessing, processing, using and sharing data.

Adapted from Evolution of Data Spaces

European Commission

# DATA SPACES in a Nutshell (2):



AI Testing and Experimentation Facilities

AI on demand platform

High-Performance Computing

Marketplace

Smart Middleware solutions

Edge & cloud Services

European Data Spaces Technical Framework

Data Spaces Support Centre

Sectoral Data Spaces

Coordinating the development of data spaces

- Bring together stakeholders.

- Identify the data the infrastructure needs.

- Build the data spaces

Assuring common standards and interoperability

European Commission

# Data Spaces in DIGITAL:
## An Overview of all DSs in WP 21-22

| | | |
|---|---|---|
| Green Deal | Smart communities | Mobility |
| Manufacturing | Agriculture | Cultural Heritage |
| Health - Genomics | Health Cancer Images | Media |
| Financial | Skills | Language |
| Public Procurement | Security and Law Enforcement | Tourism |

European Commission

A Language
Data
Space for Europe

# Worth Noting

# Reasoning

**Multimodal Data**

- Text, audio and video are a large part of the data produced every day. *However, …*

**Data Ecosystem**

- Such data need to be aggregated and organized for AI-based language services to process and extract relevant information. *Hence, …*
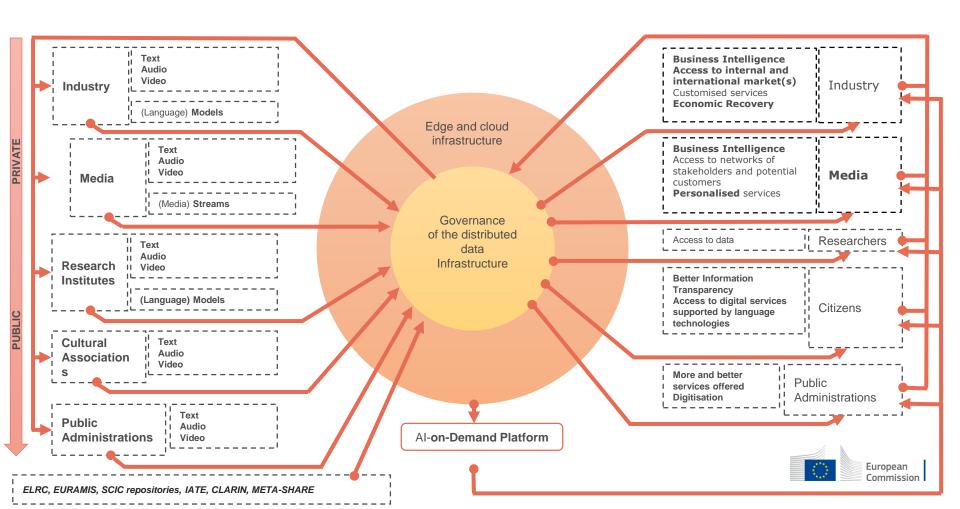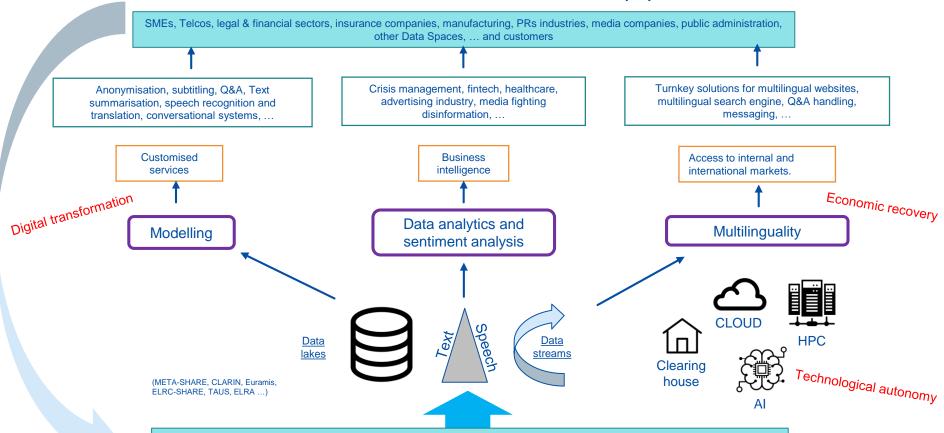
**Language Data Space**

- for the collection, creation, sharing and reuse of multimodal language data &

- as a support to the deployment of large multimodal language models + AI language technologies services through the AI platform.

# LANGUAGE DATA SPACE in a Nutshell



PRIVATE

PUBLIC

**Industry**
- Text
- Audio
- Video

(Language) **Models**

**Media**
- Text
- Audio
- Video

(Media) **Streams**

**Research Institutes**
- Text
- Audio
- Video

(Language) Models

**Cultural Associations**
- Text
- Audio
- Video

**Public Administrations**
- Text
- Audio
- Video

*ELRC, EURAMIS, SCIC repositories, IATE, CLARIN, META-SHARE*

Edge and cloud infrastructure

Governance of the distributed data Infrastructure

AI-**on-Demand Platform**

**Business Intelligence
Access to internal and international market(s)**
Customised services
**Economic Recovery**

Industry

**Business Intelligence**
Access to networks of stakeholders and potential customers
**Personalised** services

**Media**

Access to data

Researchers

**Better Information
Transparency
Access to digital services** supported by language technologies

Citizens

**More and better services offered
Digitisation**

Public Administrations

European Commission

# LANGUAGE DATA SPACE in a Nutshell (2)

SMEs, Telcos, legal & financial sectors, insurance companies, manufacturing, PRs industries, media companies, public administration, other Data Spaces, … and customers

Anonymisation, subtitling, Q&A, Text summarisation, speech recognition and translation, conversational systems, …

Crisis management, fintech, healthcare, advertising industry, media fighting disinformation, …

Turnkey solutions for multilingual websites, multilingual search engine, Q&A handling, messaging, …

Customised services

Business intelligence

Access to internal and international markets.

Modelling

Data analytics and sentiment analysis

Multilinguality

*Digital transformation*

*Economic recovery*

*Technological autonomy*

Data lakes

Text   Speech

Data streams

CLOUD

HPC

Clearing house

AI

(META-SHARE, CLARIN, Euramis, ELRC-SHARE, TAUS, ELRA …)

Existing repositories and aggregators, EU and MS administrations, media companies, Telcos, publishers, other Data Spaces

# Language Data Space

- Collection, creation, sharing and reuse of **multimodal language data**.

- This will support the deployment of **large multimodal language models** and a wide range of **AI language technologies services** to be offered through the AI platform.

Two Strands:

- Establish an institutional Centre of Excellence for Language Technologies (**CELT**)

- **Deployment** of the language data space

# 1ˢᵗ Work Strand

- Establish an institutional Centre of Excellence for Language Technologies (CELT) **to coordinate across the Member States** the creation and collection of multimodal language data and models.

*This will make use of existing EU initiatives of language data collections such as ELRC, EURAMIS, SCIC repositories, IATE, CLARIN, META-SHARE.*

# 1ˢᵗ Work Strand - Detail

- Develop in close collaboration with the Member States a multi-stakeholder **data and services governance scheme**, bringing together large industrial entities, public stakeholders and small-and-mid-size enterprise stakeholders.

- **Identify and bring together** existing European stakeholders and initiatives.

- Elaborate a blueprint for the **language data ecosystem** based on existing EU legislation and language data policies to build innovative business models across all stakeholders.

- Identify the **large multimodal language models** to be deployed.

- **Identify datasets and data streams** (public, private, citizen-collected….) relevant to the creation of large language models to be brought into conformity with the new blueprint standards and principles

- Establish a **detailed roadmap** on how to deploy the language data space.

European Commission

# 2nd work strand - Deployment

- deploying the necessary **infrastructure** for the collection and sharing of multimodal language data and models;

- implementing the **data & services governance**, business models and strategies as defined by CELT;

- supporting the increased uptake and usage of the language data space and its derived language technologies services among European **private and public sectors**, including various DIGITAL platforms and endeavours/initiatives, also through **piloting and deployment projects**;

- establishing Member States and industrial collaboration to bring the set of identified language datasets and data streams into **conformity** with the new blueprint standards and principles;

- establishing Member States and industrial collaboration to **create and deploy the identified large multimodal language models**;

- deploying advanced **AI-based language processing services and tools** to be made available through the AI-on-demand platform.

# Links & contacts

**PHILIPPE GELIN**
European Commission

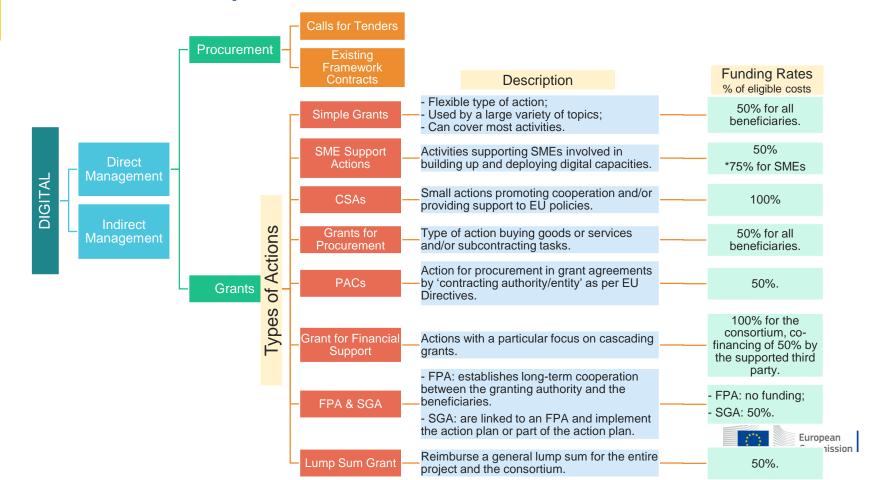DG CONNECT - Communications Network, Content and Technology

Dir G: Data

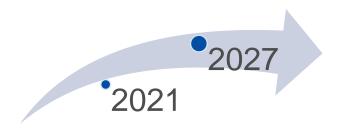Unit G3: Accessibility, **Multilingualism** and Safer Internet

Philippe.Gelin@ec.europa.eu

# DIGITAL Implementation in WP 21-22

# DIGITAL in a Nutshell:
## The Programme *Budget*



DIGITAL in MFF 2021-2027

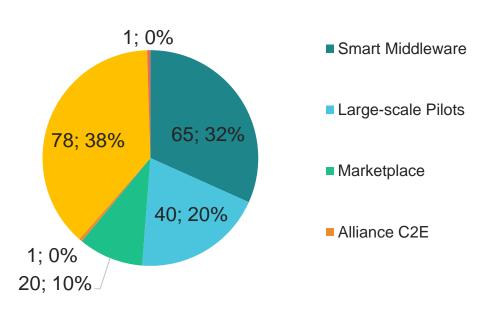| Total financing €7.6 Bio | Total financing for *Cloud, data and Artificial Intelligence* €2.1 Bio |
|---|---|
| Funding & tenders (europa.eu) | |

DIGITAL WP 21-22

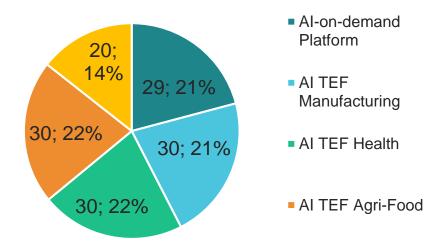| Total financing €1.98 Bio | Cloud to Edge: €204 Mio Data Spaces: €206 Mio TEF & AI-on-demand Platform: €139 Mio Total financing for *Cloud, data and AI*: €549 Mio |
|---|---|
| DIGITAL Work Programme | |

■ Cloud-to-Edge

■ Data Spaces

■ TEF & AI-on-Deman Platform

# DATA SPACES Implementation in WP 21-22

## 2. Cloud to Edge: **€204 Mio**



- Smart Middleware (65; 32%)
- Large-scale Pilots (40; 20%)
- Marketplace (20; 10%)
- Alliance C2E (1; 0%)
- (78; 38%)
- (1; 0%)

## 3. TEF & AI-on-demand Platform: **€139 Mio**



- AI-on-demand Platform (29; 21%)
- AI TEF Manufacturing (30; 21%)
- AI TEF Health (30; 22%)
- AI TEF Agri-Food (30; 22%)
- (20; 14%)

European Commission

# Data Spaces in DIGITAL:
## Indicative Budget according to WP 21-22 (1)

DATA SPACES Indicative Budget (Mio)
over three possible calls
TOT €205 Mio



European Commission