FAIR4HE – FAIR for Higher Education

# Data Stewardship Professional Competence Framework (CF-DSP)
# and Body of Knowledge (DSP-BoK)

Yuri Demchenko

University of Amsterdam

FAIRsFAIR Final WP7 Workshop

15 February 2022

**FAIRSFAIR**
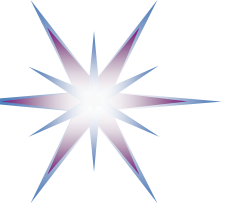Fostering Fair Data Practices in Europe
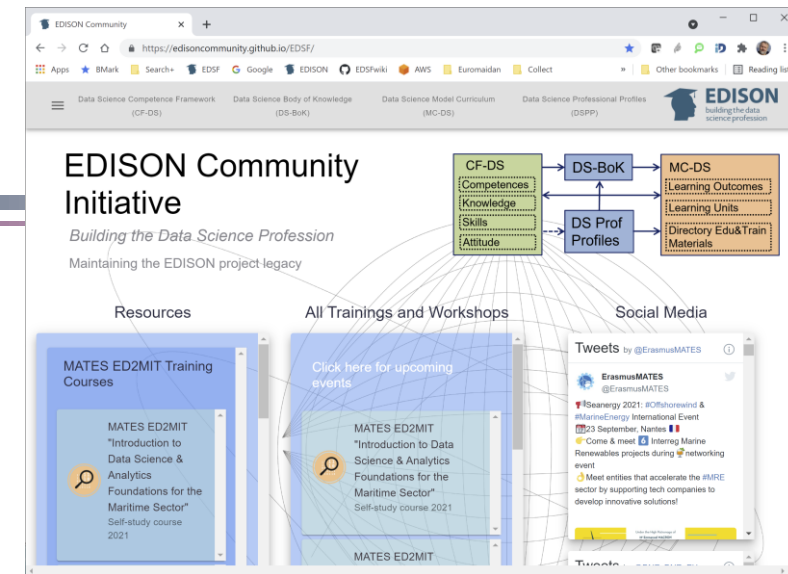
**EDISON**
building the data
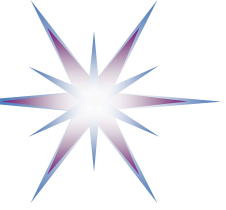science profession

# Outline

- FAIRsFAIR WP7 and Data Stewardship Professional Competence Framework (Deliverable D7.3)
- Job market analysis for Data Stewardship and related professions
  – Data collection and analysis outcome
- Proposed Data Stewardship Professional Competence Framework (CF-DSP)
  – Essential competence groups
- DSP Body of Knowledge: Advice for curriculum design and practical teaching
- Reference dataset for Data Stewardship competences assessment
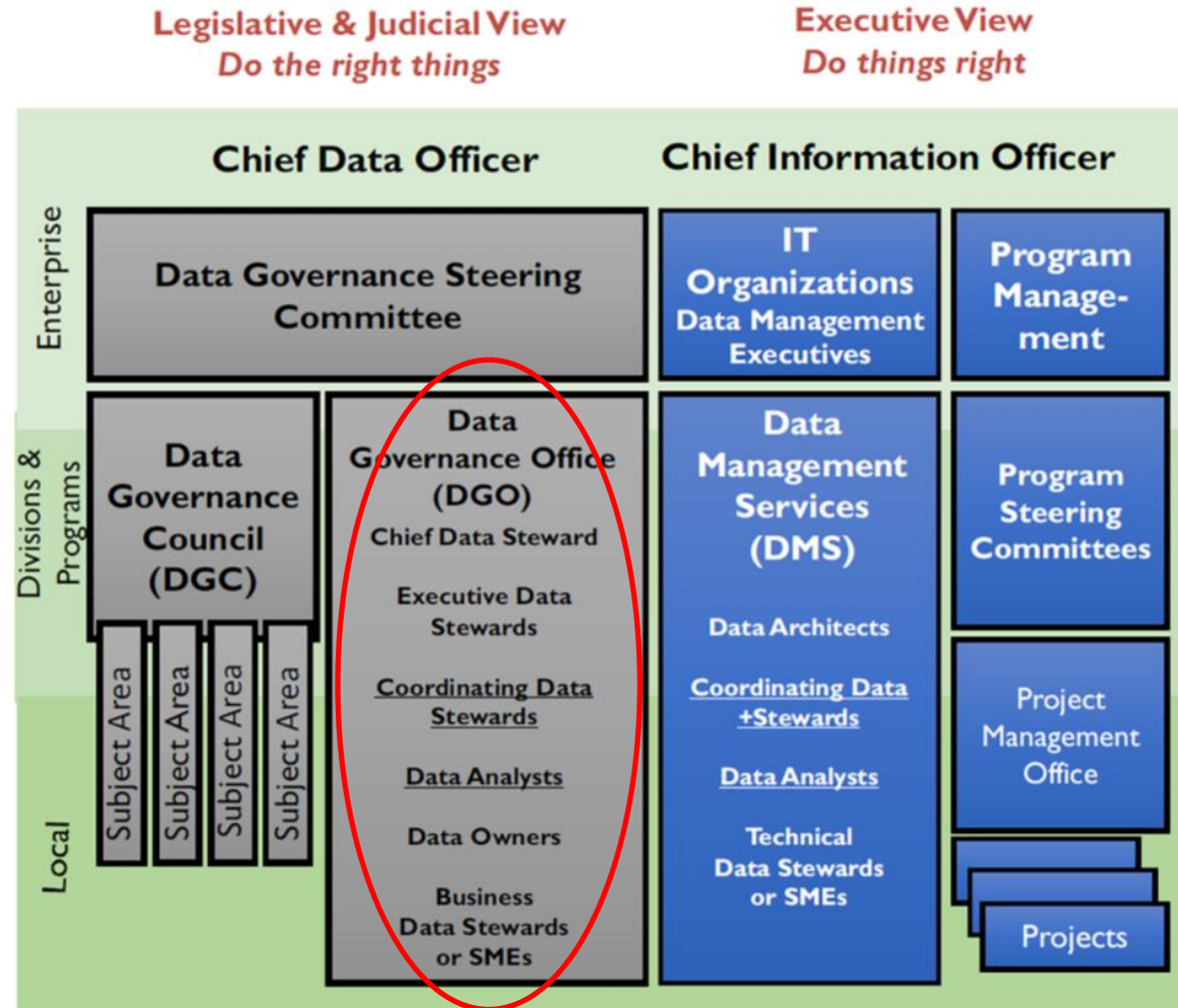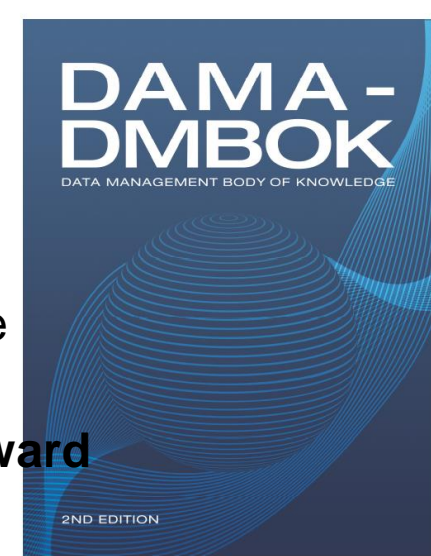- Discussion

# Methodology: How to put all together?



- EDISON Data Science Framework (EDSF) as a base methodology https://edisoncommunity.github.io/EDSF/
  - CF-DS: Data Science Competence groups:
    DSDA, **DSENG, DSDM,** DSRMP, DSDK structure and mapping
  - EDSF Data Management Body of Knowledge
    (DS-BoK, includes KAG-DSDM)
- DAMA Body of Knowledge (DAMA BoK)
- Job market Analysis: Evidence based and market driven
- Existing Frameworks for Data Stewardship and Research Data Management (RDM)

# DMBOK Framework: Data Governance Organisation Parts



- Separation of governance responsibilities
- Multi-layer
- CDO
- CIO
- Councils

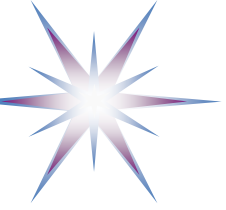Data Governance Office (DGO)
- Chief **Data Steward**
- Executive **Data Steward**
- Business **Data Steward** or SME
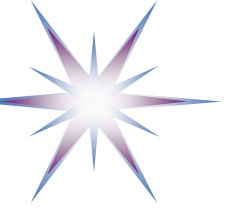
## Data Steward functions

- Creating and managing core Metadata
- Documenting rules and standards
- Managing data quality issues
- Executing operational data governance activities

*"Best Data Steward is not made but found" DMBOK1 (2009)*

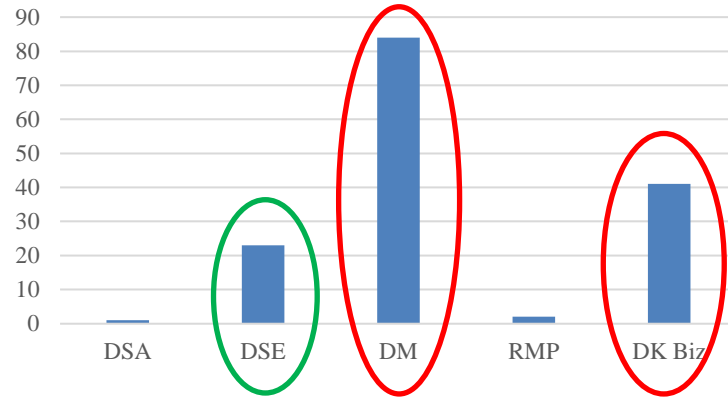# Data Stewardship Job Market Analysis (Snapshot Sept 2020)

- First stage: Exploratory – manual collection

- Period data collected: 30 August – 1 September 2020
- Sites: Indeed.com – NL, UK, DE, USA (large number of vacancies); monsterboard.nl, IEEE Jobs – NL (single vacancies)
- Days vacancy open: >50% more than 30 days
- Data Steward and related vacancies discovered: NL – 51, UK – 30+, DE ~20, US – 300+
- Information collected/downloaded:
- Key skills snapshot: – for all or first 200 for USA
- Full vacancy texts analysed: – approx. 40 in total
- Detailed analysis of sample vacancies
- Number of companies and organisations posted Data Steward related jobs: – more than 50

- New job market analysis (massive data collection): Preliminary May 2021, Advanced – Sept 2021
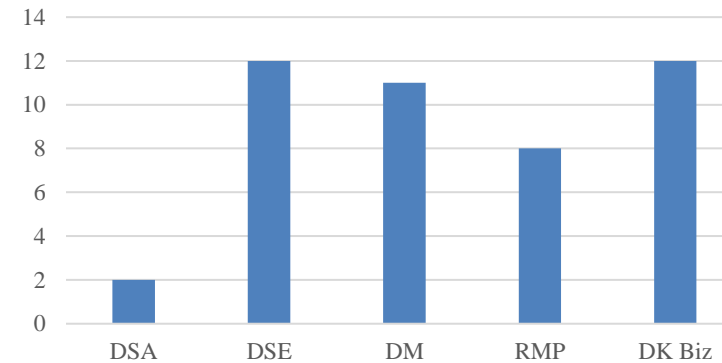  - Datasets and trained data analytics model has been published

Wide range of Competences: Responsibility, Functions, activities, however focus om
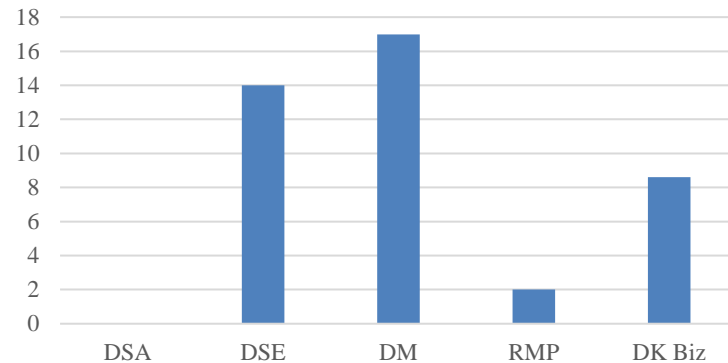Data Management and Business/Domain specific competences



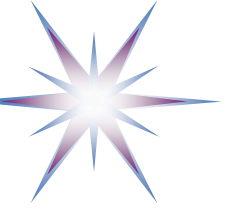Functions/Abilities - Competences

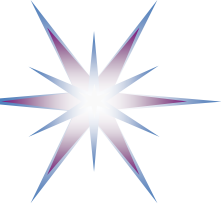Knowledge topics

Required Experience/skills

DSA – Data Science and Analytics

DSE – Data Science Engineering

DM – Data Management and Governance

RMP – Research Methods and Project Management

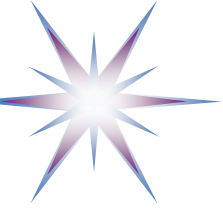DK Biz – Domain Knowledge, particular Business domain

# Important Knowledge Items extracted from Job vacancies (indeed.com – NL, DE, UK, US, Sept 2020)

- Data Management techniques
- FAIR data principles
- Data Management and Data Governance principles
- Data integrity
- Metadata, PID and linked data
- Ontology and Semantics
- FAIR metrics and Maturity framework, FAIR certification
- Data compliance regulations and standards
- Data privacy law
- GDPR
- Ethics
- Business process management
- Marketing
- Banking financial services and data management
- Multilevel Bill of Materials
- Data Warehouses

- Version control system
- Master Data Management (MDM) and Reference Data
- Research methods
- Project management
- Data analysis and visualisation tools
- Data lifecycle, lineage, provenance
- Visual Basic for Applications (VBA) and interface design
- WebAPI use for data access, collection and publishing
- DevOps, Agile, Scrum methods and technologies
- Data formats, standards
- Data modeling (SQL and EDBMS, NoSQL)
- Modern data infrastructure: Data registries, Data Factories, Semantic storage, SQL/NoSQL

# Proposed CF-DSP Competences DSDM01-DSDM04 as extension to CF-DS general Data Management Competences (1)

| EDSF Competences | Proposed DSP Competence |
|---|---|
| **Data Management (DSDM)** | **Relevance and proposed changes and extensions** (posted as revised text and bulleted extensions) |
| DSDM<br>Develop and implement data management strategy for data collection, storage, preservation, and availability for further processing. | DSDM<br>Develop and implement data management strategy for data collection, storage, preservation, and availability for further processing,<br>• ensure compliance with FAIR data principles. |
| DSDM01<br>Develop and implement data strategy, in particular, in a form of data management policy and Data Management Plan (DMP) | DSDM01<br>Develop and implement data management and governance strategy, in particular, in a form of Data Governance Policy and Data Management Plan (DMP)<br>• Ensure **compliance with standards and best practices in Data Governance and Data Management** |
| DSDM02<br>Develop and implement relevant data models, define metadata using common standards and practices, for different data sources in variety of scientific and industry domains | DSDM02<br>Develop and implement relevant data models, define metadata using common standards and practices, for different data sources in variety of scientific and industry domains.<br>• **Ensure metadata compliance with FAIR requirements**<br>• Be familiar with the metadata management tools |
| DSDM03<br>Integrate heterogeneous data from multiple sources and provide them for further analysis and use | DSDM03<br>Integrate heterogeneous data from multiple sources and provide them for further analysis and use<br>• Perform data preparation and cleaning<br>• Match/transfer data model |
| DSDM04<br>Maintain historical information on data handling, including reference to published data and corresponding data sources (data provenance) | DSDM04<br>Maintain historical information on data handling, including reference to published data and corresponding data sources<br>• **Publish data, metadata and related metrics**<br>• Perform and maintain data archiving<br>• **Develop necessary archiving policy, comply with Open Science and Open Access policies** if applicable<br>• Ensure data provenance continuity through the whole data lifecycle |

Data Stewardship Competence Framework

# Proposed CF-DSP Competences DSDM01-DSDM04 as extension to CF-DS general Data Management Competences (2)

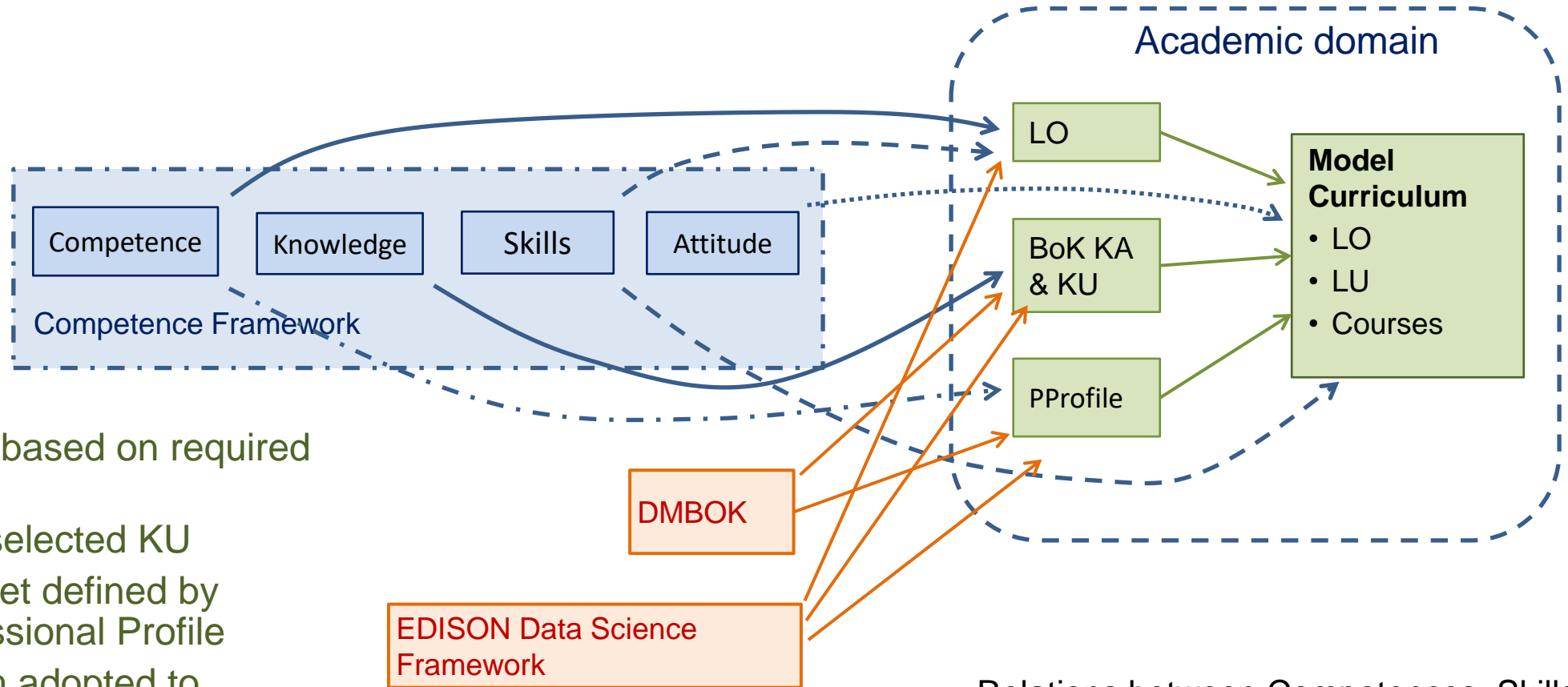| EDSF Competences | Proposed DSP Competence |
|---|---|
| DSDM05<br>Ensure data quality, accessibility, interoperability, compliance to standards, and publication (data curation) | DSDM05<br>Develop policy and metrics for data quality management (e.g. Altmetrix), maintain data quality and compliance to standards, perform data curation<br>• Interact/Collaborate with data providers and data owners to ensure data quality |
| DSDM06<br>Develop and manage/supervise policies on data protection, privacy, IPR and ethical issues in data management | DSDM06<br>Develop and manage/supervise policies on data protection, privacy, IPR and ethical issues in data management, address legal issues if necessary.<br>• Ensure GDPR compliance in data management and access<br>• Develop data access policies and coordinate their implementation and monitoring, including security breaches handling |
| None | DSDM07* (new)<br>**Manage Data Management/Data Stewards team**, coordinate related activity between organisational departments, external stakeholder to fulfil Data Governance policy requirements, provide advice and training to staff. Define domain/organisation specific data management requirements, communicate to all departments and supervise/coordinate their implementation. Coordinate/supervise data acquisition. |
| None | DSDM08* (new)<br>**Develop organisational policy and coordinate activities for sustainable implementation of the FAIR data principles** and Open Science, define corresponding requirements to data infrastructure and tools, ensure organisational awareness. |
| None | DSDM09* (new)<br>**Specify requirements to and supervise the organisational infrastructure for data management** and (and archiving), maintain the park for data management tools, provide support to staff (researchers or business developers), coordinate solving problems. |

# Necessary Data Stewardship Competences in Data Engineering and Data Infrastructure (DSENG)

Data Steward must work with the IT services and coordinate/supervise implementation of necessary organization specific tools and services

- Coordinate development of new tools and applications for data management, ensure support of the data FAIRness requirements by existing and new tools and applications
- Implement requirements for data storage facilities to comply with the data management policies and FAIR data principles in particular.
- Define and implement (coordinate) data access policies for different stakeholders and organisational roles
- Define, implement and maintain data model, reference data, master data definitions, implement consistent metadata
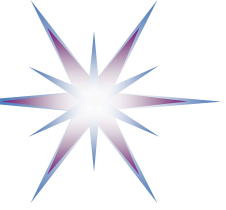
# How to use CF-DSP and DSP-BoK for curriculum Design



- LO are defined based on required competences
- LU defined by selected KU
- Curriculum target defined by intended Professional Profile
- Final curriculum adopted to available resources
- Look *Adoption Handbook* for suggested courses

Academic domain

Competence | Knowledge | Skills | Attitude

Competence Framework

LO

BoK KA & KU

PProfile

Model Curriculum
- LO
- LU
- Courses

DMBOK

EDISON Data Science Framework

- Relations between Competences, Skills, Knowledge/BoK, Professional Profiles
- Mapping between Competence elements and Academic domain elements
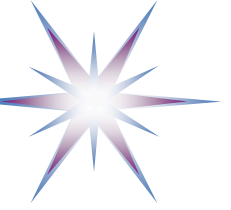
# Data Stewardship Professional Body of Knowledge

https://docs.google.com/spreadsheets/d/1TImJ4ujo79KvCpBS-5sHLFWvT1bwdCG2dzd_y5dTsnk/online/edit#gid=0
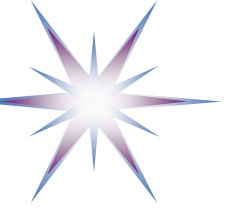
- Defined as a profile/subset and extension to the Data Science Body of Knowledge
  - Working document after FAIR Adoption Handbook sprint
  - Defined set of Knowledge Units (KU) Data Management Knowledge Area Group (KAG-DSDM)
- Competence level defined for Bachelor, Master, PhD, Postdoc, Professional
  - Mapped to EQF

- Valuable contribution from community as part of the FAIR Adoption Handbook Sprint in summer 2021
- To be published as a separate document: Zenodo, FAIRsFAIR, part of EDSF at the EDSF website

# Reference dataset: Data Stewardship vacancies collection (2021)

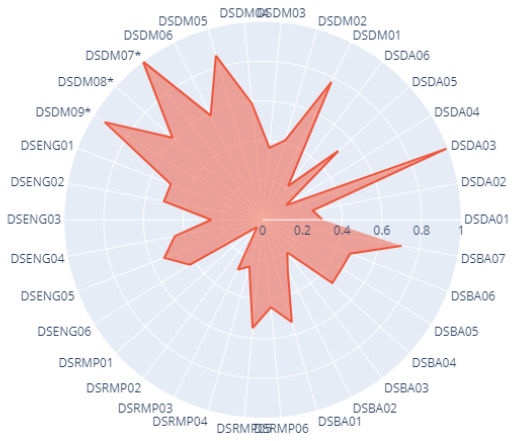- Published o Zenodo [https://zenodo.org/record/6008122](https://zenodo.org/record/6008122)
  - Data Steward Professional: Reference dataset of Data Steward related job vacancies for competences assessment
- Structure:
  - README
  - 32 vacancies from indeed.com (2021): Primarily business oriented (well described)
  - Metadata
  - Application code
- Data Steward related roles
  - Data curators, Data archivists
  - Data Architect, Data Managers
  - Data Quality Managers
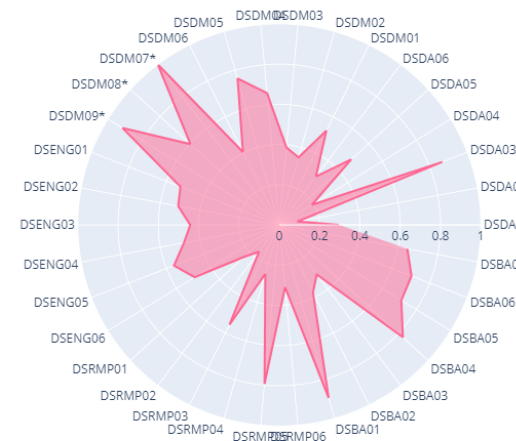- TODO: Extended dataset with research oriented vacancies worldwide

# Example Data Steward Competences structure
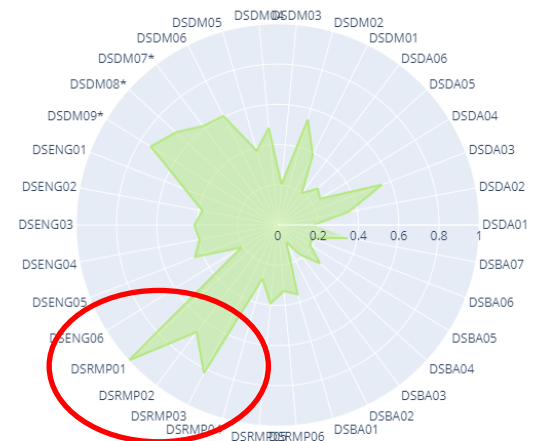


Business Consultant - Organisation Structure Data - ABN AMRO Bank

Business Developer II - ABN AMRO Bank

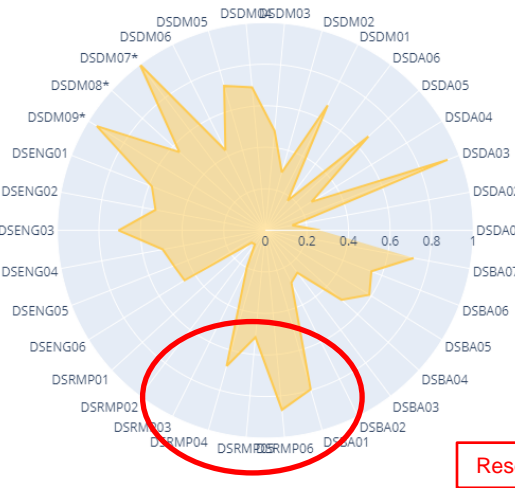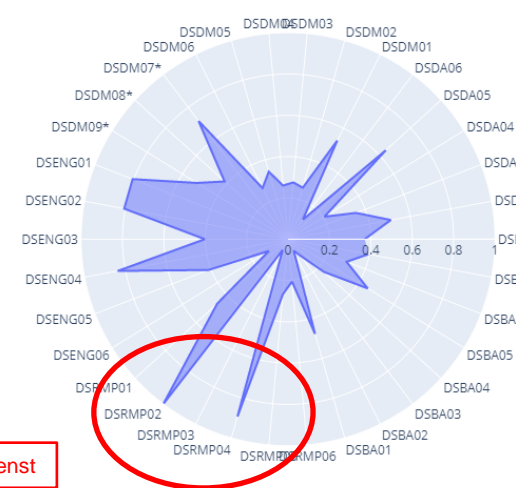Regulatory Reporting & Data Steward - Robert Walters
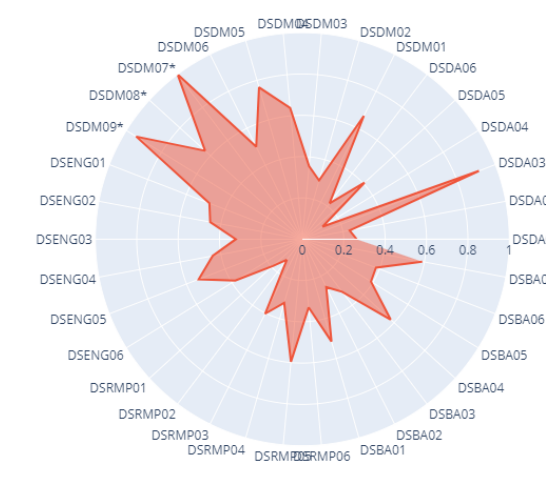
Research Data Officer - Vrije Universiteit Amsterdam

Principal Data Architect - Arcadis

Informatiemanager Groen & Water - BlueTrail

Finance Traineeship - Calco

Specialist Data Management - DAS

Research focus/elemenst

# How to use CF-DSP and DSP-BoK

- Customised Curriculum Design
- Competences – Vacancies assessment
- Data Science/Data Steward Team building

# Outcome Based Educations and Training Model: Addressing target competences for the profession



**Data Science Competence Framework (CF-DS)** → **DS Professional Profiles (DSP-P)** ⟵ **Target Competences**

**Learning Outcomes (LO) (OBE Learning Model)**

**Knowledge Units (KU) (DS-BoK)** → **Learning Units (LU) (MC-DS)**

**Tracks/Specialisations (based on DSP-P) (LO, LU, KU, courses/modules)**

**DS Academic Programmes** — Courses → LO/KU/LU/DSP-P

**DS Professional Training** — Modules → LO/KU/LU/DSP-P

From Competences and DSP Profiles

to Learning Outcomes (LO) and

to Knowledge Unites (KU) and Learning Units (LU)

- EDSF allow for customized educational courses and training modules design

RESEARCH DATA MANAGEMENT AND ETICAL ISSUES 2
RESEARCH METHOD 5
MODELLING AND EXPERIMENT PLANNING 4
DATA SELECTION AND QUALITY EVALUATION 4

DATA MINING 3
DATA PREPARATION AND PREPROCESSING 3
PERFORMANCE ANALYSIS 3
MACHINE LEARNING 6

1° YEAR I SEMESTER — 1° YEAR II SEMESTER — 2° YEAR I SEMESTER — 2° YEAR II SEMESTER

USE CASES ANALYSIS 3
DATA MINING 2
DATA PREPARATION AND PREPROCESSING 4
PERFORMANCE ANALYSIS 4

PREDICTIVE ANALYSIS 6
REINFORCED LEARNING 5
GENERALIZED LINEAR MODELS 6
MODELLING AND SIMULATION THEORY 2

- Planning on sequence and duration of courses for maximum learning outcome
- Importance of splitting core courses over 2 semesters
- Theory and practice oriented courses
- Importance of pre-requisite knowledge:
  - Statistics for Data Scientists
  - Organisational management for Data Stewards

MATCHING – COMPETENCE PROFILES

## Individual Education/Training Path based on Competence benchmarking

- Red polygon indicates the chosen professional profile: Data Scientist (general)
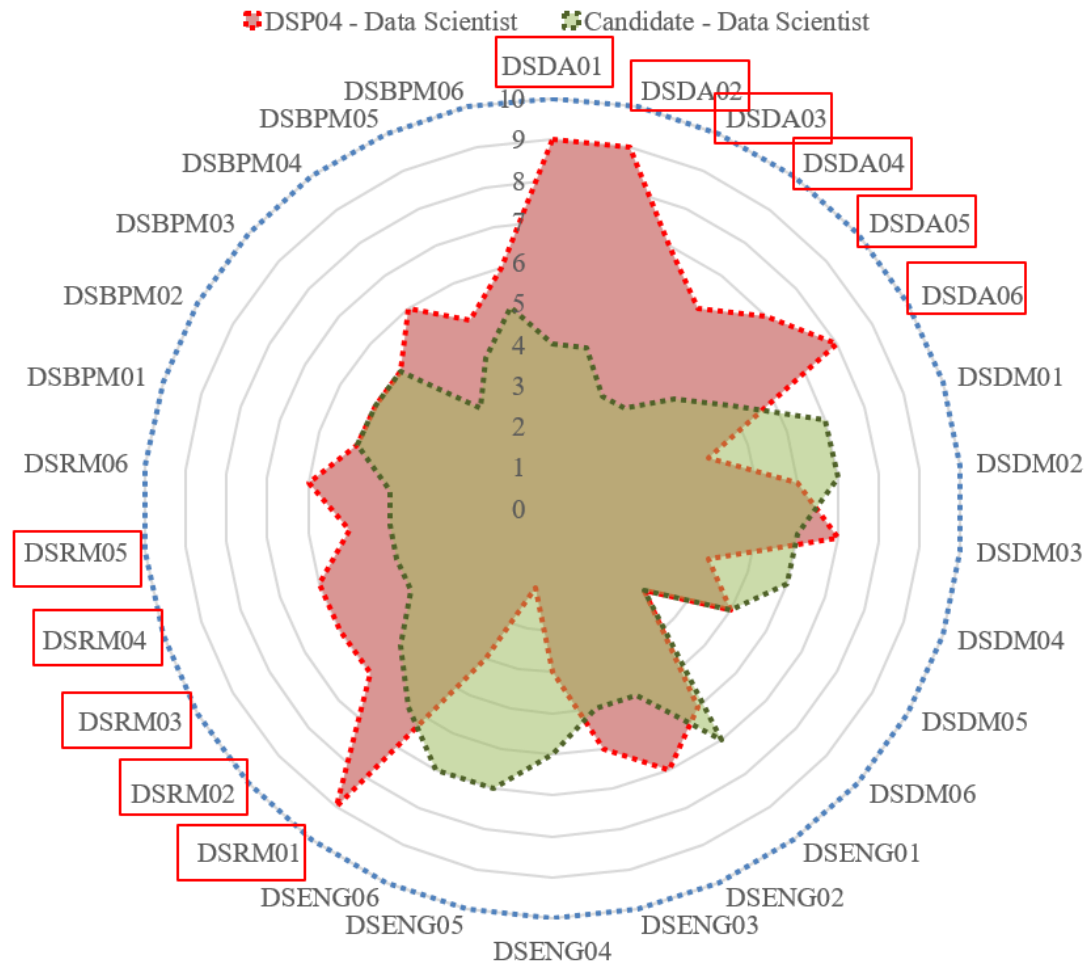- Green polygon indicates the candidate or practitioner competences/skills profile
- Insufficient competences (gaps) are highlighted in *red*
  - *DSDA01 – DSDA06 Data Science Analytics*
  - *DSRM01 – DSRM05 Data Science Research Methods*
- Can be use for team skills match marking and organisational skills management

[ref] For DSP Profiles definition and for enumerated competences refer to EDSF documents CF-DS and DSP Profiles.

# Building a Data Science Team

# Discussion

- How to use the Data Stewardship Professional Competence Framework (CF-DSP) and Body of Knowledge (DSP-BoK) for teaching and organizational staff/HR management

- How to implement FAIR and Data Stewardship in university curriculum and teaching: Challenges, opportunities and experience

- How to contribute to CF-DSP and DSP-BoK update and maintenance

# References

- FAIR Competence Framework for Higher Education (Data Stewardship Professional Competence Framework), FAIRsFAIR Project Deliverable D7.3
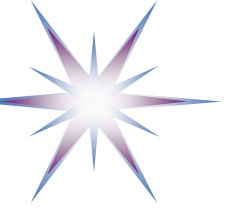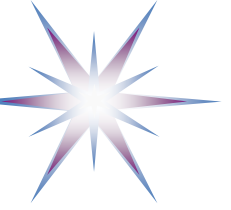    - https://zenodo.org/record/4562089#.YIBZeegzZPZ
- D7.2 Briefing on FAIR Competences and Synergies
    - https://zenodo.org/record/4009007#.YL86SfkzaF4
- Data Stewardship Curricula in Denmark
    - https://www.deic.dk/sites/default/files/Data%20Steward%20Education%20in%20Denmark_0.pdf
- ZonMw & ELIXIR-NL funded project "Towards FAIR Data Steward as profession for the Life Sciences"
    - Final report (Oct 3, 2019): https://doi.org/10.5281/zenodo.3471707
- Yuri Demchenko, Lennart Stoy, Research Data Management and Data Stewardship Competences in University Curriculum, In Proc. Data Science Education (DSE), Special Session, EDUCON2021 – IEEE Global Engineering Education Conference, 21-23 April 2021, Vienna, Austria
    - http://www.uazone.org/demch/papers/educon2021-data-stewardship-competence-fw-v02.pdf
- EDISON Data Science Framework (EDSF)
    - https://edisoncommunity.github.io/EDSF/
- FAIRsFAIR FAIR Adoption Handbook: D7.4 How to be FAIR with your data. A teaching and training handbook for higher education institutions:
    - https://zenodo.org/record/5905866
- Data Steward Professional: Reference dataset of Data Steward related job vacancies for competences assessment
    - https://zenodo.org/record/6008122

# Additional information

- ESDF development and EDSF Release 4
- FAIR data principles, technical context and organisational roles

# EDISON Project (2015-2017) and EDISON Data Science Framework (EDSF)

- EDISON project website - http://edison-project.net/

- EDISON Data Science Framework (EDSF) – main outcome of the project
- Currently maintained by EDISON Community Initiative, coordinated by UvA

- EDSF Release 3 published in 2018 – Currently active
- EDSF Release 4 (EDSF2021) to be published by the end of 2021 (initially planned end 2020)
  - Reviewed at EDSF Release 4 Design Workshop – 20 Nov 2019, UvA

# Competences Map to Knowledge and Skills

- **Competence** is a demonstrated ability to apply knowledge, skills and attitudes for achieving observable results

Model Curriculum and Learning Units, Courses

Learning Outcomes

Knowledge Topics

Knowledge Area, Knowledge Units
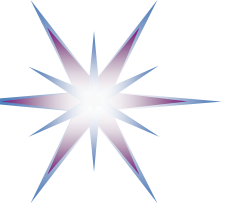
Competence vs Competency

Learning Outcomes vs Learning Objectives

**Competences**

**Knowledge**

**Ability to perform Org functions**

**Skills**

Education (BoK)

Experience & Workplace Train

**Professional Profiles**

# FAIR Data Principles: Metadata Management (GO FAIR recommendations)

## Findable:

- F1 (meta)data are assigned a globally unique and persistent identifier;

- F2 data are described with rich metadata;

- F3 metadata clearly and explicitly include the identifier of the data it describes;

- F4 (meta)data are registered or indexed in a searchable resource;
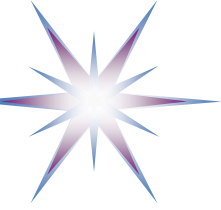
## Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.

- I2. (meta)data use vocabularies that follow FAIR principles;

- I3. (meta)data include qualified references to other (meta)data;

## Accessible:

- A1 (meta)data are retrievable by their identifier using a standardized communications protocol;
  - A1.1 the protocol is open, free, and universally implementable;
  - A1.2 the protocol allows for an authentication and authorization procedure, where necessary;

- A2 metadata are accessible, even when the data are no longer available;
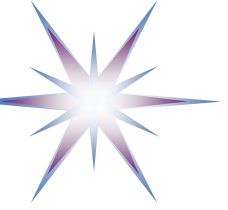
## Reusable:

- R1 meta(data) are richly described with a plurality of accurate and relevant attributes;

- R1.1 (meta)data are released with a clear and accessible data usage license;

- R1.2 (meta)data are associated with detailed provenance;

- R1.3 (meta)data meet domain-relevant community standards;

# FAIR adoption and Ecosystem Sustainability Elements

- FAIR must be accepted by all roles in organisational data management and governance process
  - FAIR must be endorsed by top management C-level
  - Roles and responsibilities to be defined and staffed
  - Inter-role functions as factor for modern agile organisations
- FAIR must be adopted for the whole Research/industrial Data lifecycle
- FAIR must be practiced by all participants along data lifecycle and specifically started from the data producers i.e. researchers or facility operators or sale agents
- FAIR must be supported by infrastructure and tools
- FAIR must be embedded into applications development
- Organisational capability and capacity management
- Education and training – To enable them all
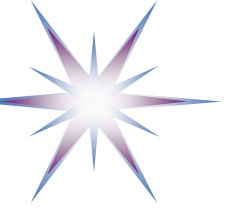  - Basic academic and professional education + continuous education

# FAIR from the technical point of view

- Findable
  - Metadata and PDI – infrastructure and tools
  - Registries and handles resolution, API
  - Policies and SLA
- Accessible
  - Repositories and data storage: infrastructure and management
  - Policy and access control: infrastructure and API management
  - Data access protocols
  - Usage Policy and Sovereignty
  - Data protection, compliance, privacy and GDPR
- Interoperable
  - Standard data formats
  - Metadata and API
  - FAIR maturity level and certification
- Reusable
  - Data provenance and lineage
  - Preservation
  - Metadata, PID and API – linked or embedded into datasets

This motivates Data Stewards' interaction with both **Data Analytics and Applications developers** roles and **Data Infrastructure** roles
- Consequently related competences from Data Stewards are needed

# FAIR Data Management and Organisational Roles

FAIR data principles to be adopted cross organisation for the whole data lifecycle

- Data collection
  - Researchers, Data Engineers, data entry workers
- Data preservation and curation
  - Data curators, Data Custodians/Archivists
- Data Analysis
  - Data Scientists, Data Architects, Application developers
- Data publication, sharing access
  - Data Stewards, Data Curators
- Data Governance and Data management
  - Data Stewards and CDO
    - Data policy and data delivery agreements
- Data Infrastructure and tools for data storage and handling
  - Storage, database engineers/managers
    - Metadata and PID services, Master data and Reference data