

# EDISON Project Overview:

Building the Data Science Profession  
for Research and Industry

Yuri Demchenko, University of Amsterdam



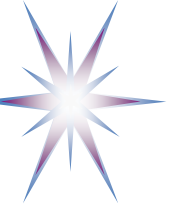
**EDISON**  
building the data  
science profession

EDISON – Education for **D**ata Intensive  
Science to **O**pen **N**ew science frontiers

2nd EDISON Data Science Champions  
Conference

15-16 March 2017

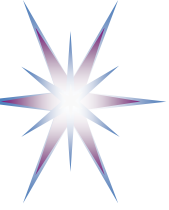
Universidad Carlos III de Madrid, Madrid



# Outline

- Background and motivation
  - European initiatives related to Digital Single Market and Digital Skills Agenda
  - EDISON building network to promote establishing Data Science profession
- EDISON Data Science Framework (EDSF)
  - From Data Science Competences to Body of Knowledge and Model Curriculum
- Data Science Competence Framework: Essential competences and skills
- Education and training focus
  - Data Science Body of Knowledge (DS-BoK)
  - Data Science Model Curriculum (MC-DS)
  - Example Research Data Management Literacy curriculum
- Further developments to formalise EDSF and Data Science profession establishment





# EDISON Objectives, Impact and Actions

IMPACT

**Increase the number of Data Scientists and Market for establishing Data Science Profession**

Objectives and Actions

**Create a Data Science profession**

Data Science professional profiles

Interact with demand and supply sides

Career path building and skills transferability

**Services to education and training**

Define Model Curriculum and design tools

Support for accreditation and certification

Collaborating and sharing expertise and materials

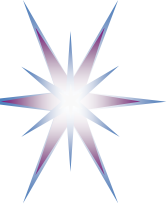
**Engage stakeholder communities**

Sustain platforms of communities of practice

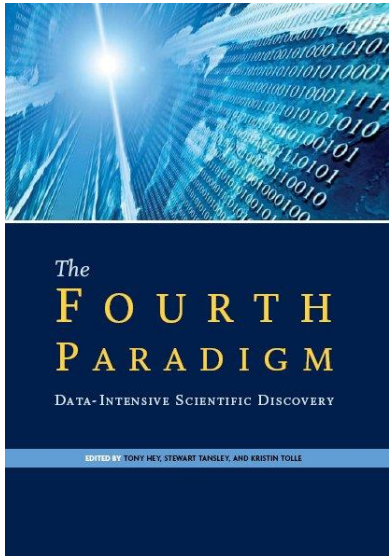
Create community of "champion" universities

Interact with Expert Liaison Groups

**Data Science Competence Framework and Body of Knowledge**



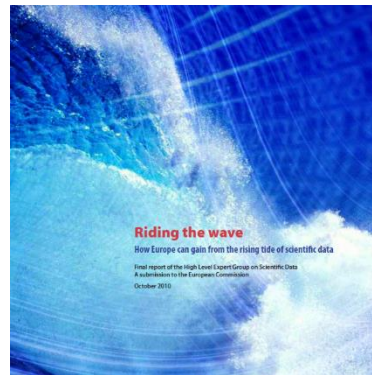
# Visionaries and Drivers: Seminal works, High level reports, Activities



## The Fourth Paradigm: Data-Intensive Scientific Discovery.

By Jim Gray, Microsoft, 2009. Edited by Tony Hey, Kristin Tolle, et al.

<http://research.microsoft.com/en-us/collaboration/fourthparadigm/>



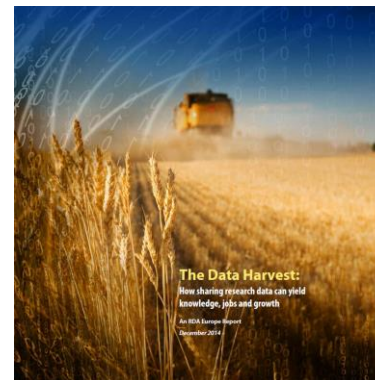
## Riding the wave: How Europe can gain from the rising tide of scientific data.

Final report of the High Level Expert Group on Scientific Data. October 2010.

<http://cordis.europa.eu/fp7/ict/e-infrastructure/docs/hlg-sdi-report.pdf>



## HLEG report on European Open Science Cloud (October 2016)

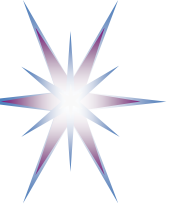


## The Data Harvest: How sharing research data can yield knowledge, jobs and growth.

An RDA Europe Report. December 2014

<https://rd-alliance.org/data-harvest-report-sharing-data-knowledge-jobs-and-growth.html>

**Emergence of Cognitive Technologies**  
(IBM Watson and others)



# Recent European Commission Initiatives 2016

**Digitising European Industry: Reaping the full benefits of a **Digital Single Market**.**  
COM(2016) 180 final, Brussels, 19.4.2016

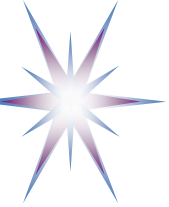
- The need for **new multidisciplinary and digital skills in particular Data Scientist**
  - Expected rapidly growing demand will lead to more than 800 000 unfilled vacancies by 2020

**European Cloud Initiative** - Building a competitive data and knowledge economy in Europe, COM(2016) 178 final, Brussels, 19.4.2016

- **European Open Science Cloud (EOSC)** and European digital research and data infrastructure
  - To offer 1.7 million European researchers and 70 million professionals in science and technology open and seamless services for **storage, management, analysis and re-use** of research data
- Address growing demand and shortage of data-related skills

**A New Skills Agenda for Europe**, COM(2016) 381 final Brussels, 10.6.2016

- Addresses the need for digital and complementary skills, ensure young talents flow into data driven research and industry
- Launch **Digital Skills and Jobs Coalition (1st December 2016, Brussels)** to develop comprehensive national digital skills strategies by mid-2017



# HLEG report on European Open Science Cloud (October 2016) – Demand for Data Scientists/Stewards

**Realising the European Open Science Cloud.** First report and recommendations of the Commission High Level Expert Group on the European Open Science Cloud, October 2016

[https://ec.europa.eu/research/openscience/pdf/realising\\_the\\_european\\_open\\_science\\_cloud\\_2016.pdf](https://ec.europa.eu/research/openscience/pdf/realising_the_european_open_science_cloud_2016.pdf)

- Definition of the **Data Steward** as a distinctive role and profession
  - Core Data Experts need to be trained and their career perspective improved
- **Estimation: More than 80,000 data stewards to serve 1.7 mln scientists in Europe (1 per every 20 scientists)**
  - Based on 5% grant funding for Data management and preservation
- **Clash of cultures** between domain specialists and e-Infrastructure specialists (i.e. IT/Computer Science)



# OECD and UN on Digital Economy and Data Literacy

## OECD

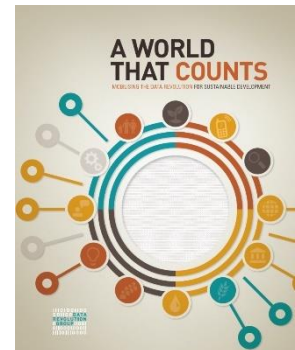
- Demand for new type of *“dynamic self-re-skilling workforce”*
- Continuous learning and professional development to become a shared responsibility of workers and organisations

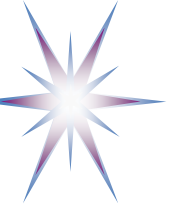
[ref] SKILLS FOR A DIGITAL WORLD, OECD, 25-May-2016

[http://www.oecd.org/officialdocuments/publicdisplaydocumentpdf/?cote=DSTI/ICCP/IIS\(2015\)10/FINAL&docLanguage=En](http://www.oecd.org/officialdocuments/publicdisplaydocumentpdf/?cote=DSTI/ICCP/IIS(2015)10/FINAL&docLanguage=En)

## UN

- Data Revolution Report "A WORLD THAT COUNTS" Presented to Secretary-General (2014)  
<http://www.undatarevolution.org/report/>
- Data Literacy is defined as key for digital revolution
- **Data literacy** = critically analyse data collected and data visualised





# Approach

- Task is not for one project – **Need collaboration**
- Task is not for **science** only in isolation from **industry**
- Needs strong **conceptual approach**
  - Use science to solve the problems of science
- **Standardisation** is an important factor of sustainability and development





# EDISON Network and Engagement Activity (1)

- Cooperative relations and exchange of developments with RI projects
  - ELIXIR, CORBEL/RIttrain, EUDAT, ENVRI
  - FOSTER2, EOSCpilot
- Cooperation with Big Data and Data Science projects
  - EDSA, BDVA
- Active contribution to the Research Data Alliance (RDA) activities
  - RDA IG on Education and Training on Handling Research Data (IG-ETHRD)
  - BoFs and proposed WG on Certification and accreditation,
  - Proposed WG on Text Data Mining
  - Proposed WG on Research Data Management Curriculum
  - BoF on Data Champions



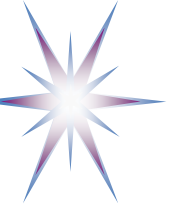
# EDISON Network and Engagement Activity (2)

- Workshops to promote a common approach towards addressing growing demand for Data Science and critical data competences and skills as required by **European Research Infrastructures (RI)**, future **European Open Science Cloud (EOSC)** and generally European **Digital Single Market (DSM)**.
  - **Joint EDISON and EC workshop** “Data Infrastructure Competences and Skills Framework: a European and Global Challenge” (Brussels, 9th February, 2016)
  - Joint IEEE, STC CC and RDA Workshop on Curricula and Teaching Methods (DTW2015 and DTW2016 collocated with IEEE CloudCom)
- EDISON initiated a set of **national action meetings** to address Data Science and digital skills by bringing together key stakeholders from universities, employer associations, and government
  - A first workshop jointly organised by the EDISON project and Dutch Ministry of Education, Culture and Science in June 2016 (during Netherlands Presidency in EU)



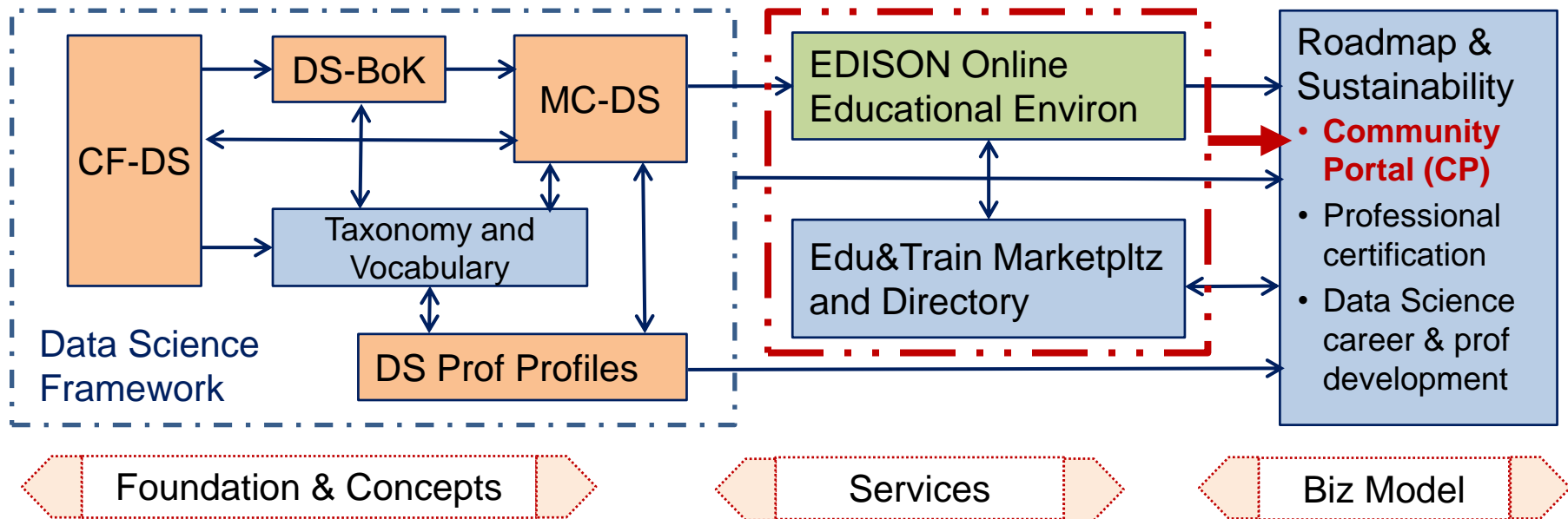
# EDISON Network and Engagement Activity (3)

- European and international standardisation bodies and professional organisations
  - CEN TC426 Committee (former e-Competence Framework e-CFv3.0 workshop)
  - ESCO (European Skills, Competence, Occupations)
  - CEPIS and association **ICT Professionalism Europe** (co-signed 21 Nov 2016, Amsterdam)
- EDISON Booth at the Launch event of the **Digital Skills and Jobs Coalition: Boosting Europe's Digital skills**, 1 December 2016, Brussels
  - Part of actions toward European Digital Single Market (DSM) and
- Contribution to International standardization bodies, professional organisations and initiatives
  - **Business Higher Education Forum (BHEF) in USA**
  - **DARE project for APEC** (Asia Pacific Economic Cooperation) to develop a Data Analytics checklist for APEC countries
  - **Data Science Curriculum Meeting of Professional and Academic Societies in USA** (4 March 2017, Alexandria) including ACM



# EDISON Data Science Framework (EDSF)

Release 1, October 2016

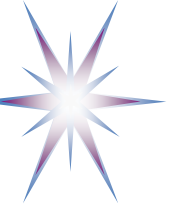


## EDISON Framework components

- CF-DS – Data Science Competence Framework
- DS-BoK – Data Science Body of Knowledge
- MC-DS – Data Science Model Curriculum
- DSP – Data Science Professional profiles
- Data Science Taxonomies and Scientific Disciplines Classification
- EOEE - EDISON Online Education Environment

## Methodology

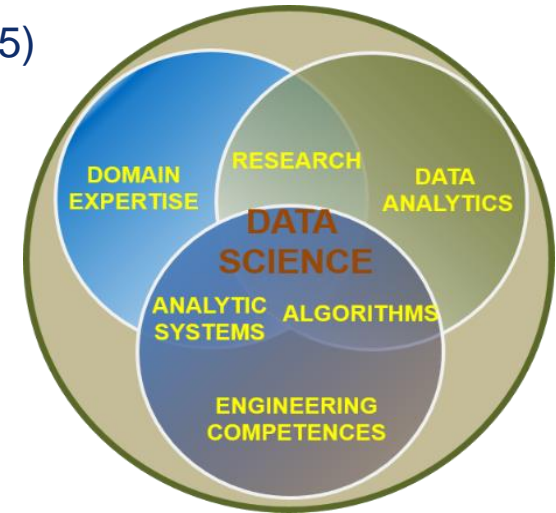
- Job market study, existing practices in academic, research and industry.
- Compliance with related standards
- Review and feedback from the ELG, expert community, domain experts.
- Input from the champion universities and community of practice.



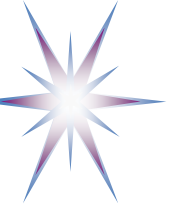
# Data Scientist definition

Based on the definitions by NIST Big Data WG (NIST SP1500 - 2015)

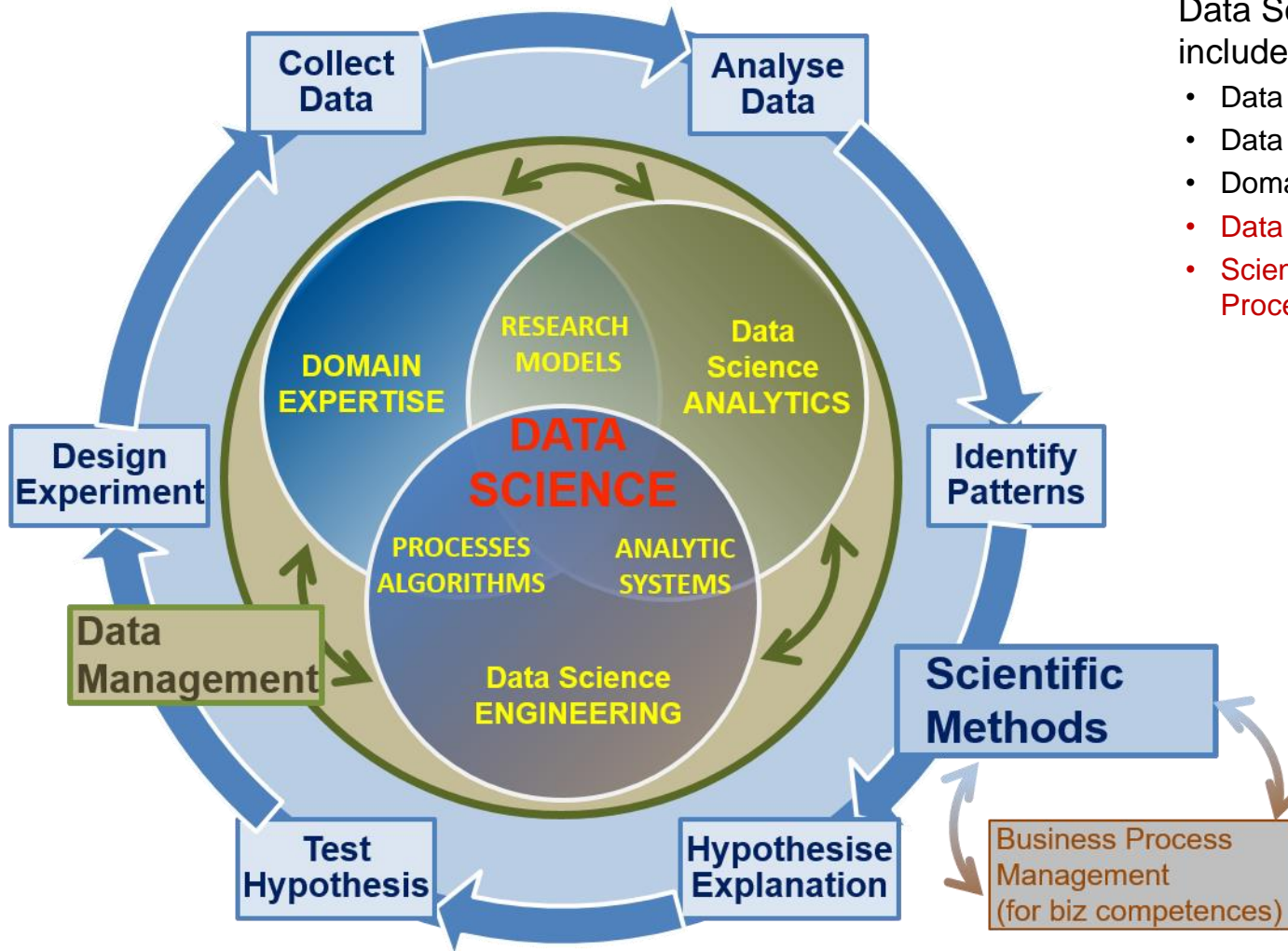
- **A Data Scientist is a practitioner who has sufficient knowledge in the overlapping regimes of expertise in business needs, domain knowledge, analytical skills, and programming and systems engineering expertise to manage the end-to-end scientific method process through each stage in the **big data lifecycle****
  - ... Till the delivery of expected scientific and business value to science or industry
- **Other definitions to admit such features as**
  - Ability to solve variety of business problems
  - Optimize performance and suggest new services for the organisation
  - Develop a special mindset and be statistically minded, **understand raw data and “appreciate data as a first class product”**
- **Data science is the empirical synthesis of actionable knowledge and technologies required to handle data from raw data through the complete data lifecycle process.**
- **Big Data is the technology to build system and infrastructures to process large volume of structurally complex data in a time effective way**



[ref] Legacy: NIST BDWG definition of Data Science



# Data Science Competence Groups - Research



Data Science Competence includes 5 areas/groups

- Data Analytics
- Data Science Engineering
- Domain Expertise
- **Data Management**
- **Scientific Methods (or Business Process Management)**

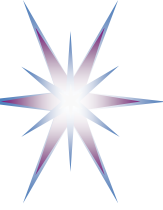
## Scientific Methods

- Design Experiment
- Collect Data
- Analyse Data
- Identify Patterns
- Hypothesise Explanation
- Test Hypothesis

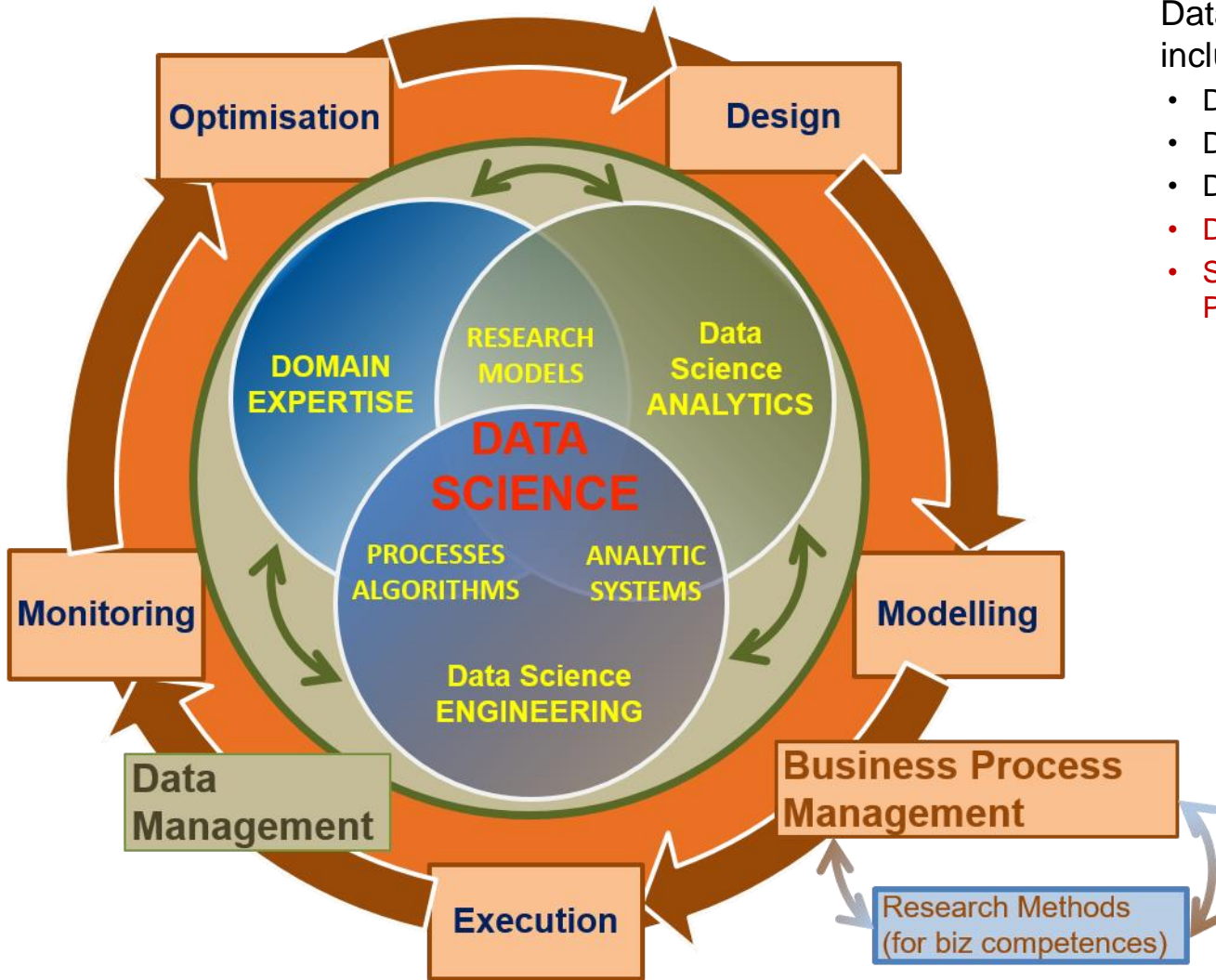
## Business Operations

- Operations Strategy
- Plan
- Design & Deploy
- Monitor & Control
- Improve & Re-design





# Data Science Competences Groups – Business



Data Science Competence includes 5 areas/groups

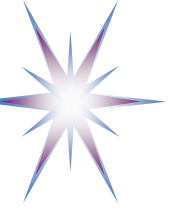
- Data Analytics
- Data Science Engineering
- Domain Expertise
- **Data Management**
- **Scientific Methods (or Business Process Management)**

## Scientific Methods

- Design Experiment
- Collect Data
- Analyse Data
- Identify Patterns
- Hypothesise Explanation
- Test Hypothesis

## Business Process Operations/Stages

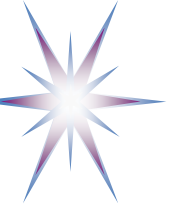
- Design
- Model/Plan
- Deploy & Execute
- Monitor & Control
- Optimise & Re-design



# Identified Data Science Competence Groups

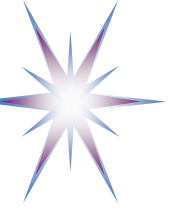
	Data Science Analytics (DSDA)	Data Management (DSDM)	Data Science Engineering (DSENG)	Research/Scientific Methods (DSRM)	Data Science Domain Knowledge, e.g. Business Processes (DSDK/DSBPM)
0	Use appropriate statistical techniques and predictive analytics on available data to deliver insights and discover new relations	Develop and implement data management strategy for data collection, storage, preservation, and availability for further processing.	Use engineering principles to research, design, develop and implement new instruments and applications for data collection, analysis and management	Create new understandings and capabilities by using the scientific method (hypothesis, test/artefact, evaluation) or similar engineering methods to discover new approaches to create new knowledge and achieve research or organisational goals	Use domain knowledge (scientific or business) to develop relevant data analytics applications, and adopt general Data Science methods to domain specific data types and presentations, data and process models, organisational roles and relations
1	<b>DSDA01</b> Use predictive analytics to analyse big data and discover new relations	<b>DSDM01</b> Develop and implement data strategy, in particular, Data Management Plan (DMP)	<b>DSENG01</b> Use engineering principles to design, prototype data analytics applications, or develop instruments, systems	<b>DSRM01</b> Create new understandings and capabilities by using scientific/ research methods or similar domain related development methods	<b>DSBPM01</b> Understand business and provide insight, translate unstructured business problems into an abstract mathematical framework
2	<b>DSDA02</b> Use statistical techniq to deliver insights	<b>DSDM02</b> Develop data models including metadata	<b>DSENG02</b> Develop and apply computational solutions	<b>DSRM02</b> Direct systematic study toward a fuller knowledge or understanding of the observable facts	<b>DSBPM02</b> Participate strategically and tactically in financial decisions
3	<b>DSDA03</b> Develop specialized ...	<b>DSDM03</b> Collect integrate data	<b>DSENG03</b> Develops specialized tools	<b>DSRM03</b> Undertakes creative work	<b>DSBPM03</b> Provides support services to other
4	<b>DSDA04</b> Analyze complex data	<b>DSDM04</b> Maintain repository	<b>DSENG04</b> Design, build, operate	<b>DSRM04</b> Translate strategies into actions	<b>DSBPM04</b> Analyse data for marketing
5	<b>DSDA05</b> Use different analytics	<b>DSDM05</b> Visualise cmplx data	<b>DSENG05</b> Secure and reliable data	<b>DSRM05</b> Contribute to organizational goals	<b>DSBPM05</b> Analyse optimise customer relatio





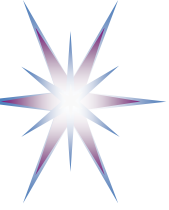
# Identified Data Science *Skills/Experience* Groups

- **Group 1: Skills/experience related to competences**
  - Data Analytics and Machine Learning
  - Data Management/Curation (including both general data management and scientific data management)
  - Data Science Engineering (hardware and software) skills
  - Scientific/Research Methods or Business Process Management
  - Application/subject domain related (research or business)
  - Mathematics and Statistics
- **Group 2: Big Data (Data Science) tools and platforms**
  - Big Data Analytics platforms
  - Mathematics & Statistics applications & tools
  - Databases (SQL and NoSQL)
  - Data Management and Curation platform
  - Data and applications visualisation
  - *Cloud based platforms and tools*
- **Group 3: Programming and programming languages and IDE**
  - General and specialized development platforms for data analysis and statistics
- **Group 4: Soft skills or Social Intelligence**
  - Personal, inter-personal communication, team work, professional network



# Practical Application of the CF-DS

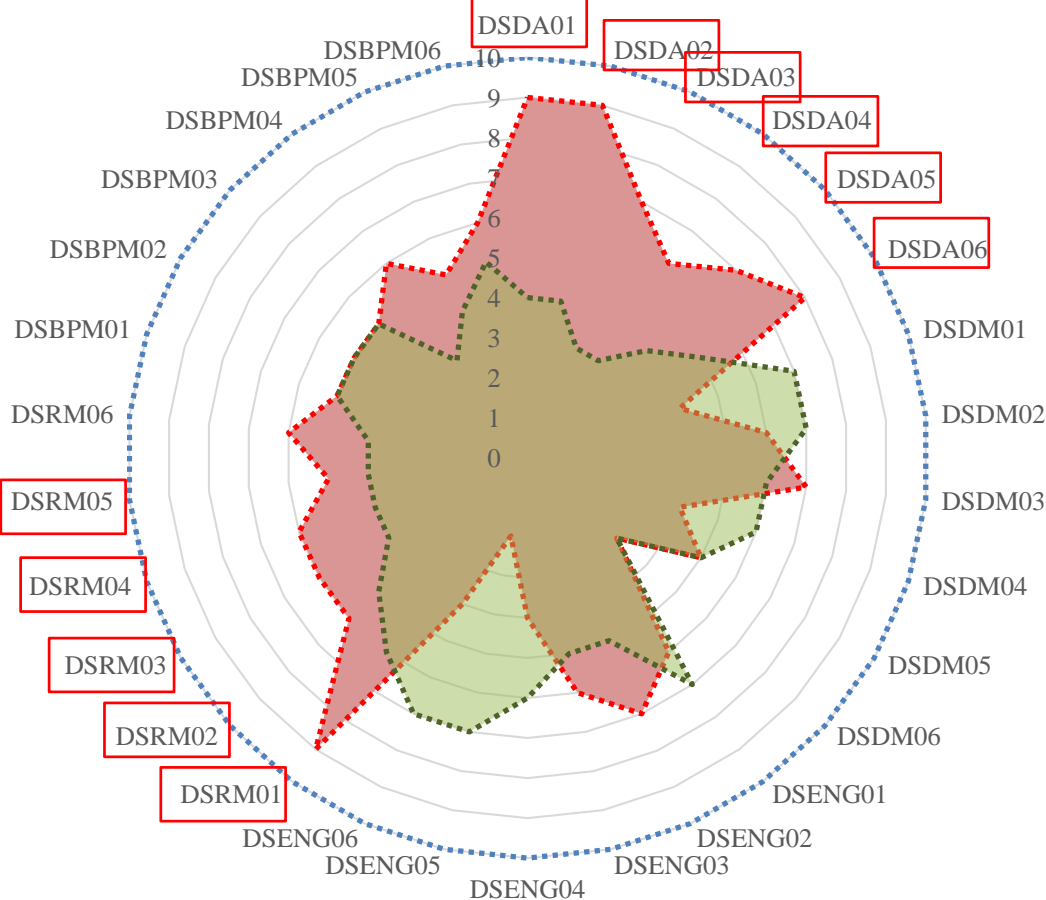
- Basis for the definition of the Data Science Body of Knowledge (DS-BoK) and Data Science Model Curriculum (MC-DS)
  - CF-DS => Learning Outcomes (MC-DS) => Knowledge Areas (DS-BoK)
  - CF-DS => Data Science taxonomy of scientific subjects and vocabulary
- Data Science professional profiles definition
  - Extend existing EU standards and occupations taxonomies: e-CFv3.0, ESCO, others
- Professional competence benchmarking
  - For customizable training and career development
  - Including CV or organisational profiles matching
- Professional certification
  - In combination with DS-BoK professional competences benchmarking
- Vacancy construction tool for job advertisement (for HR)
  - Using controlled vocabulary and Data Science Taxonomy



# Individual Competences Benchmarking

## MATCHING – COMPETENCE PROFILES

■ DSP04 - Data Scientist   ■ Candidate - Data Scientist

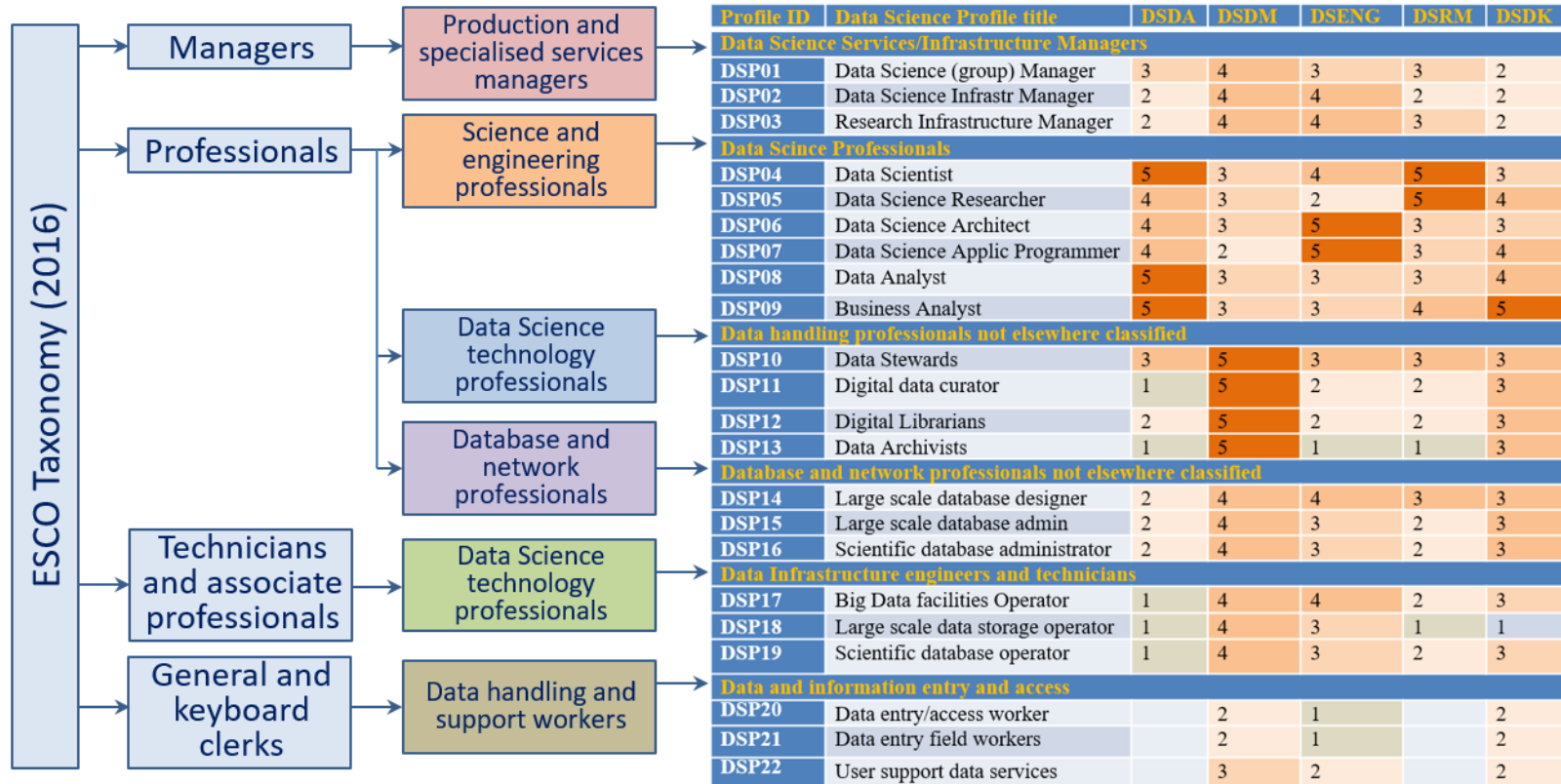


## Individual Education/Training Path based on Competence benchmarking

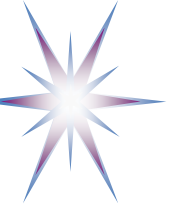
- Red polygon indicates the chosen professional profile: Data Scientist (general)
- Green polygon indicates the candidate or practitioner competences/skills profile
- Insufficient competences (gaps) are highlighted in *red*
  - DSDA01 – DSDA06 Data Science Analytics
  - DSRM01 – DSRM05 Data Science Research Methods
- Can be used 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.

# DSP Profiles mapping to ESCO Taxonomy High Level Groups



- DSP Profiles mapping to corresponding CF-DS Competence Groups
  - Relevance level from 5 – maximum to 1 – minimum



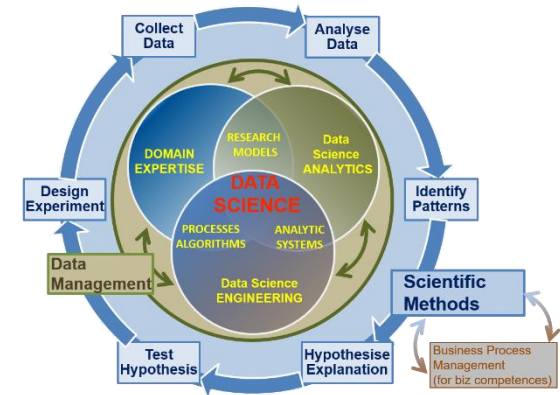
# Education and Training


- Foundation and methodological base
  - Data Science Body of Knowledge (DS-BoK)
    - Taxonomy and classification of Data Science related scientific subjects
  - Data Science Model Curriculum (MC-DS)
    - Set Learning Units mapped to CF-DS Learning and DS-BoK Knowledge Areas/Units
  - Instructional methodologies and teaching models
- Platforms and environment
  - Virtual labs, datasets, developments platforms
  - Online education environment and courses management
- Services
  - Individual benchmarking and profiling tools (competence assessment)
  - Knowledge evaluation tools
  - Certifications and training for self-made Data Scientists practitioners
  - Education and training marketplace: Courses catalog and repository

# Data Science Body of Knowledge (DS-BoK)

## DS-BoK Knowledge Area Groups (KAG)

- KAG1-DNA: Data Analytics group including Machine Learning, statistical methods, and Business Analytics
- KAG2-DSE: Data Science Engineering group including Software and infrastructure engineering
- **KAG3-DSDM: *Data Management group including data curation, preservation and data infrastructure***
- **KAG4-DSRM: *Scientific/Research Methods group***
- KAG5-DSBP: Business process management group
  
- Data Science domain knowledge to be defined by related expert groups

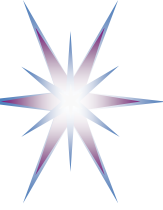




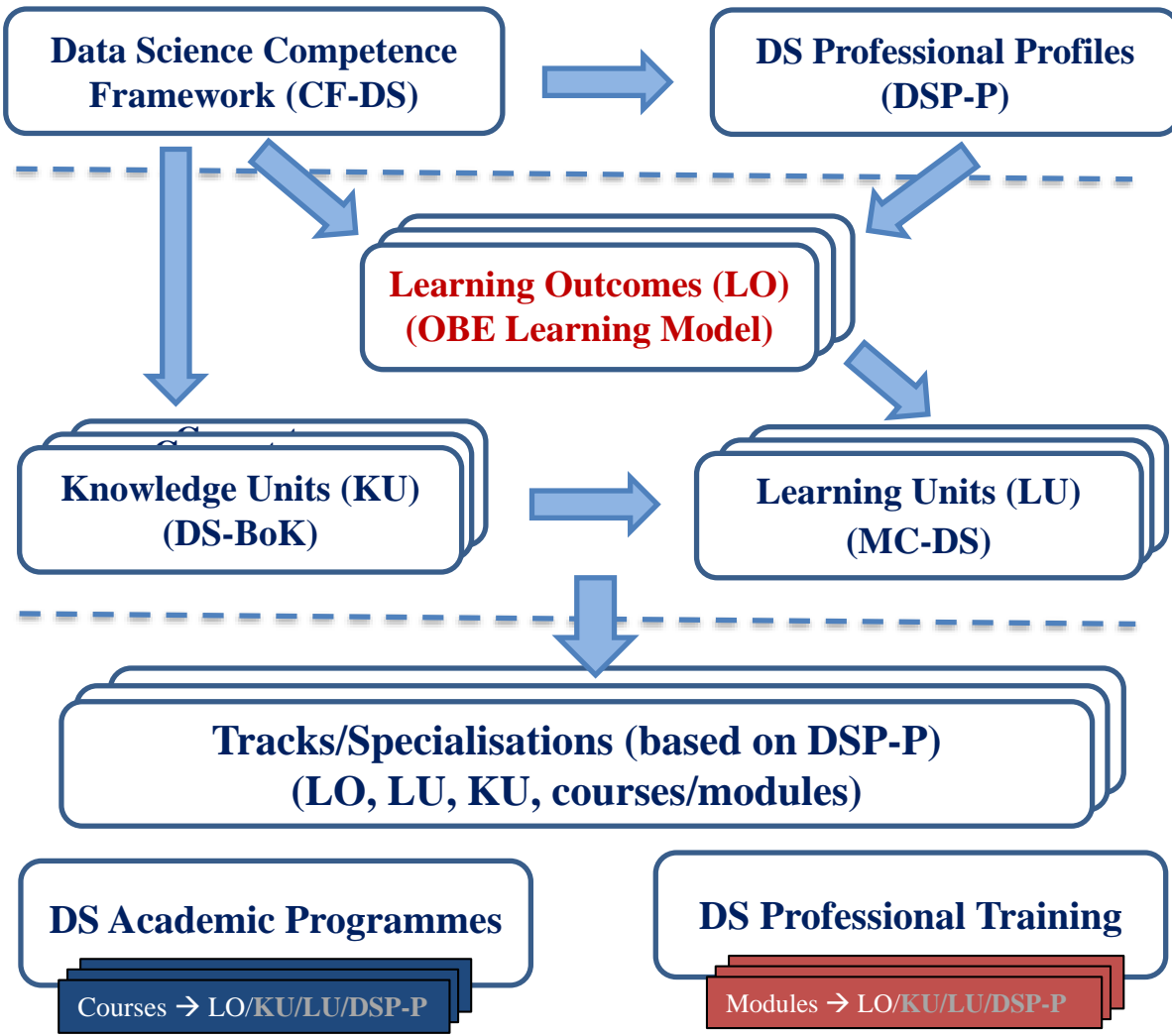
# Data Science Model Curriculum (MC-DS)

## Data Science Model Curriculum includes

- Learning Outcomes (LO) definition based on CF-DS
  - LOs are defined for CF-DS competence groups and for all enumerated competences
- LOs mapping to Learning Units (LU)
  - LUs are based on CCS(2012) and universities best practices
  - Data Science university programmes and courses inventory (interactive)  
<http://edison-project.eu/university-programs-list>
- LU/course relevance: Mandatory Tier 1, Tier 2, Elective, Prerequisite
- Learning methods and learning models (in progress)



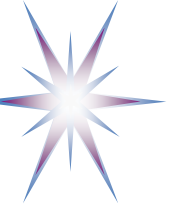
# Outcome Based Educations and Training Model



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

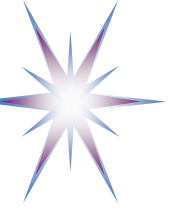




# Example DS-BoK Knowledge Areas definition and mapping to existing BoKs and CCS (2012)

Knowledge Area Groups (KAG)	Knowledge Areas (KA)	Suggested Knowledge Units (KU)	Mapping to CCS2012 (including suggested Data Science extensions) and existing BoKs
KAG1-DSDA: Data Analytics group (including Machine Learning, statistical methods)	Theory of computation	Design and Analysis of Algorithms	CCS2012: Theory of computation Design and analysis of algorithms Data structures design and
		Machine Learning Theory	
Knowledge Area Groups (KAG)	Knowledge Areas (KA)	Suggested Knowledge Units (KU)	Mapping to CCS2012 (including suggested Data Science extensions) and existing BoKs
KAG2-DSENG: Data Science Engineering group including Software an infrastru engineering	Computer systems organisation for Big Data	Parallel and Distributed Computer Architecture	CCS2012: Computer systems organization Architectures Parallel architectures
		Computer networks architectures	
Knowledge Area Groups (KAG)	Knowledge Areas (KA)	Suggested Knowledge Units (KU)	Mapping to CCS2012 (including suggested Data Science extensions) and existing BoKs
	Data Management and Enterprise data infrastructure	Data management, including Reference and Master Data	DM-BoK selected KAs (1) Data Governance, (2) Data Architecture, (3) Data Modelling and Design, (4) Data Storage and Operations, (5) Data Security, (6) Data Integration and Interoperability, (7) Documents and Content, (8) Reference and Master Data, (9) Data Warehousing and Business Intelligence, (10) Metadata, and (11) Data Quality.
		Data Warehousing and Business Intelligence	
		Data storage and operations	
		Data archives/storage compliance and certification	
		Metadata, linked data, provenance	
		Data infrastructure, data registries and data factories	
		Data security and protection	
		Data governance, data quality, data Integration and Interoperability	

- Mapping suggested to CCS2012 and existing BoKs



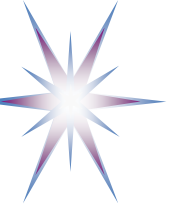
# Example MC-DS Mapping Learning Units to DS-BoK and CCS (2012)

KAG/ LU# (*)	Learning Unit (course name) <sup>2</sup>	Type/relevance <sup>3</sup>				Map to DS-BoK, CCS2012 and known BoKs	
		Tier 1	Tier 2	Elective	Pre requisite	CCS2012 based academic subjects	DS-BoK and other BoKs
	Software requirements and design					Extensions are suggested from SWEBOK	SWEBOK selected KAs <ul style="list-style-type: none"> <li>Software requirements</li> </ul>

KAG/ LU# (*)	Learning Unit (course name) <sup>2</sup>	Type/relevance <sup>3</sup>				Map to DS-BoK, CCS2012 and known BoKs		
		Tier 1	Tier 2	Elective	Pre requisite	CCS2012 based academic subjects	DS-BoK and other BoKs	
	Information theory					Mathematical analysis		production engineering configuration engineering
	Mathematical analysis							
	<i>Extensibility point for adding new courses</i>							
	Artificial Intelligence					Computing methodologies	No specific BoK are defined	engineering process engineering models and
	Natural Language Processing					Artificial intelligence		

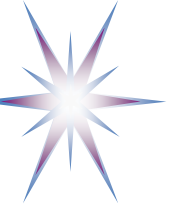
KAG/ LU# (*)	Learning Unit (course name) <sup>2</sup>	Type/relevance <sup>3</sup>				Map to DS-BoK, CCS2012 and known BoKs		
		Tier 1	Tier 2	Elective	Pre requisite	CCS2012 based academic subjects	DS-BoK and other BoKs	
	Knowledge Representation and Reasoning					CCS2012 based academic subjects	DS-BoK and other BoKs	
	Data mining and knowledge discovery					Extended with the general Data Management Knowledge Areas and related academic subjects.	General Data Management KA's Data Lifecycle Management Data archives/storage compliance and certification New KAs to support RDA recommendations and community data management models (Open Access, Open Data, etc) Data type registries, PIDs Data infrastructure and Data Factories TBD – To follow RDA and ERA community developments	
	Text analysis, Data mining							
	Text analytics including linguistic, and structural techniques to analyze and unstructured data							
	Machine Learning theoretical algorithms							
	<i>Extensibility point for adding new courses</i>							
	Research methodology, research cycle					Extended with the general Scientific/Research Methods subjects and related academic subjects.	Suggested KAs to develop DSRM related competences: Research methodology, research cycle (e.g. 4 step model Hypothesis – Research Methods – Artefact –	
	Modelling and experiment planning							

- Mapping suggested to ACM CCS2012, DS-BoK and other related BoKs



# New courses currently missing

- Data Management / Research Data Management
  - Data Curation, Data Stewardsip
- Professional issues in Data Science
  - + Ethics and responsible use of Data Science



# KAG3-DSDM: *Data Management group: data curation, preservation and data infrastructure*

DM-BoK version 2 “Guide for performing data management”

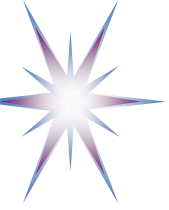
– 11 Knowledge Areas

- (1) Data Governance
- (2) Data Architecture
- (3) Data Modelling and Design
- (4) Data Storage and Operations
- (5) *Data Security***
- (6) Data Integration and Interoperability
- (7) *Documents and Content***
- (8) Reference and Master Data
- (9) Data Warehousing and Business Intelligence
- (10) *Metadata***
- (11) Data Quality

Other Knowledge Areas motivated by RDA, European Open Data initiatives, European Open Data Cloud

- (12) *PID, metadata, data registries***
- (13) *Data Management Plan***
- (14) *Open Science, Open Data, Open Access, ORCID***
- (15) *Responsible data use***

- Highlighted in red: Considered (Research) Data Management literacy (minimum required knowledge)



# Research Data Management Model Curriculum – Part of the EDISON Data Literacy Training

## **A. Use cases for data management and stewardship**

- Preserving the Scientific Record

## **B. Data Management elements (organisational and individual)**

- Goals and motivation for managing your data
- Data formats
- Creating documentation and metadata, metadata for discovery
- Using data portals and metadata registries
- Tracking Data Usage
- Handling sensitive data
- Backing up your data
- Data Management Plan (DMP) - to be a part of hands on session

## **C. Responsible Data Use Section (Citation, Copyright, Data Restrictions)**

## **D. Open Science and Open Data (Definition, Standards, Open Data use and reuse, open government data)**

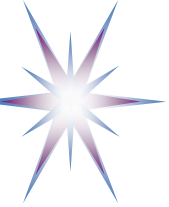
- Research data and open access
- Repository and self- archiving services
- ORCID identifier for data
- Stakeholders and roles: engineer, librarian, researcher
- Open Data services: ORCID.org, Altmetric Doughnut, Zenodo

## **E. Hands on:**

- a) Data Management Plan design
- b) Metadata and tools
- c) Selection of licenses for open data and contents (e.g. Creative Common and Open Database)

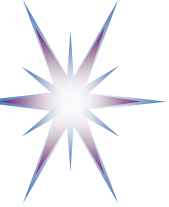
Collaboration with the Research Data Alliance (RDA) on developing model curriculum on Research Data Literacy:

- Modular, Customisable, Localised, Open Access
- Supported by the network of trainers via resource swap board



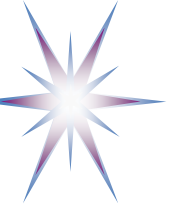
# Further Developments and Actions

- Involve academic and industry experts and professional organisations to the definition of DS-BoK following from CF-DS
- **Work with champion universities to practically validate the proposed EDSF**
- Formally provide suggestions to ESCO for the definition of the Data Science professional profiles (occupations) family
- **Formally provide suggestions for e-CF3.0 extensions for Data Science to CEN/PC 428**
  - **Involve national e-CF bodies and adopters where available**
- Suggest required ACM CCS(2012) Classification extensions and proposal for Data Science curriculum definition



# Discussion

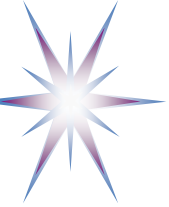
- Questions
- Comments
- **Invitation to contribution and cooperation:**
  - Forum, EDISON Liaisons Groups, Champions Conference (Spring & Summer 2017)
- EDISON project website <http://edison-project.eu/>
- EDISON Data Science Framework Release 1 (EDSF)  
<http://edison-project.eu/edison-data-science-framework-edsf>
  - Data Science Competence Framework  
<http://edison-project.eu/data-science-competence-framework-cf-ds>
  - Data Science Body of Knowledge  
<http://edison-project.eu/data-science-body-knowledge-ds-bok>
  - Data Science Model Curriculum  
<http://edison-project.eu/data-science-model-curriculum-mc-ds>
  - Data Science Professional Profiles  
<http://edison-project.eu/data-science-professional-profiles-definition-dsp>
- Survey Data Science Competences: Invitation to participate  
[https://www.surveymonkey.com/r/EDISON\\_project\\_-\\_Defining\\_Data\\_science\\_profession](https://www.surveymonkey.com/r/EDISON_project_-_Defining_Data_science_profession)



# Definitions (according to e-CFv3.0)

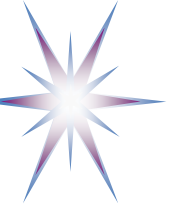
- **Competence** is a demonstrated ability to apply knowledge, skills and attitudes for achieving observable results.
  - Competence vs Competency (e-CF vs ACM)
    - Competence is ability acquired by training or education (linked to learning outcome)
    - Competency is similar to skills or experience (acquired feature of a person)
- Competence is not to be confused with process or technology concepts such as, 'Cloud Computing' or 'Big Data'. These descriptions represent evolving technologies and in the context of the e-CF, they may be integrated as elements within knowledge and skill examples.
- **Knowledge** in the context of competence definition is treated as something to know, to be aware of, familiar with, and obtained as a part of education.
- **Skills** is treated as provable ability to do something and relies on the person's experience.



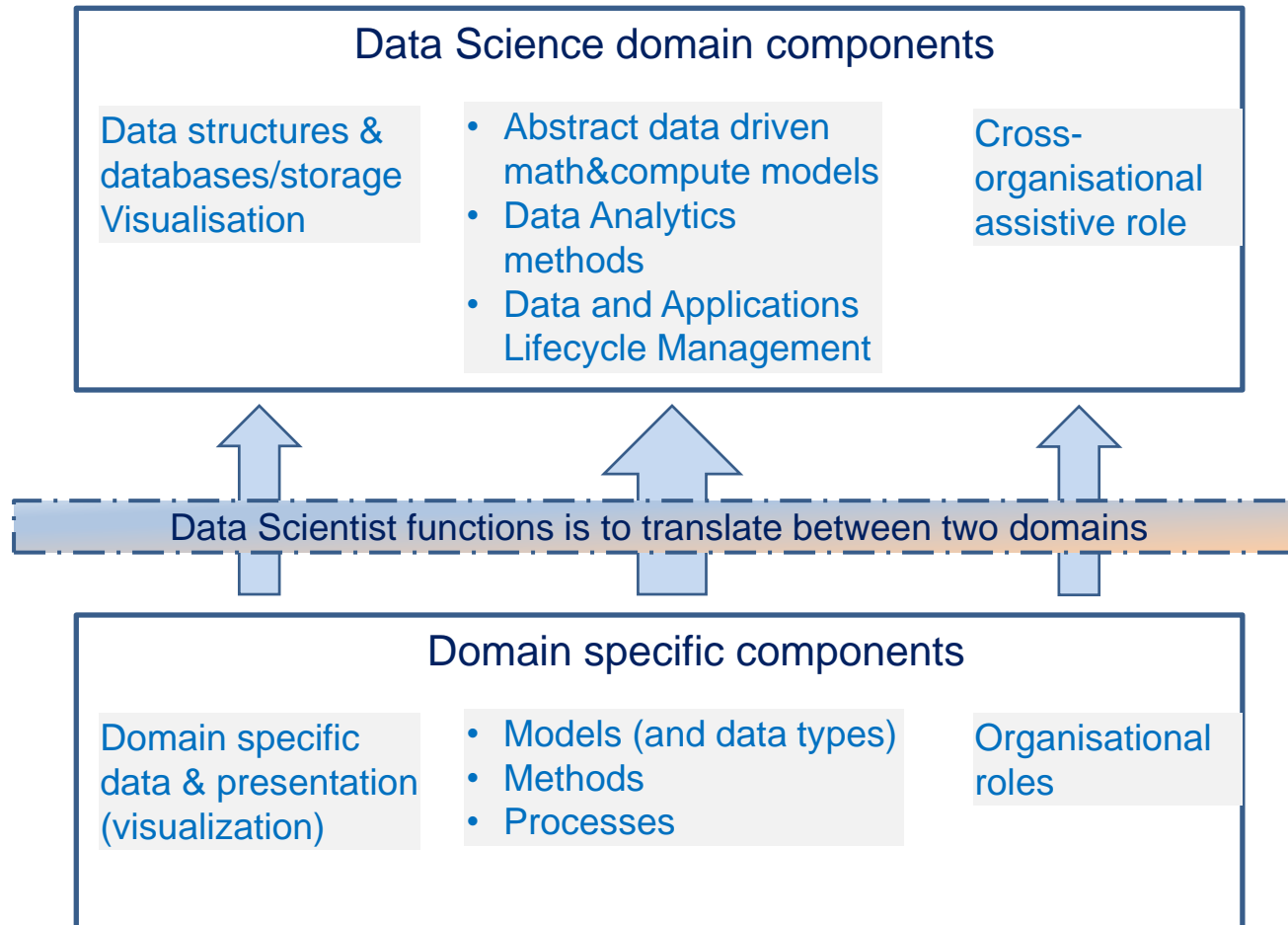


# Data Scientist and Subject Domain Specialist

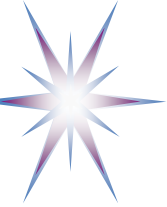
- **Subject domain components**
  - Model (and data types)
  - Methods
  - Processes
  - Domain specific data and presentation/visualization methods
  - Organisational roles and relations
- **Data Scientist is an assistant to Subject Domain Specialists**
  - Translate subject domain Model, Methods, Processes into abstract data driven form
  - Implement computational models in software, build required infrastructure and tools
  - Do (computational) analytic work and present it in a form understandable to subject domain
  - Discover new relations originated from data analysis and advice subject domain specialist
  - Present/visualise information in domain related actionable way
  - Interact and cooperate with different organizational roles to obtain data and deliver results and/or actionable data



# Data Science and Subject Domains

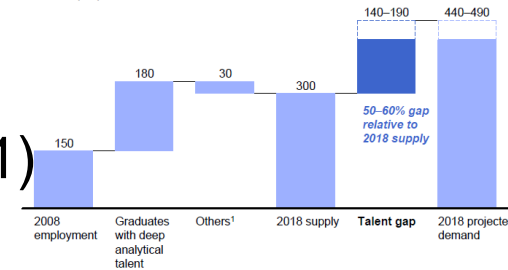


- Data Scientist role is to maintain the Data Value Chain (domain specific):**
- Data Integration => Organisation/Process/Business Optimisation => Innovation



# Demand for Data Science and data related professions

- McKinsey Global Institute on Big Data Jobs (2011)  
[http://www.mckinsey.com/mgi/publications/big\\_data/index.asp](http://www.mckinsey.com/mgi/publications/big_data/index.asp)
  - Estimated gap of 140,000 - 190,000 data analytics skills by 2018
- IDC Report on European Data Market (2015)
  - Number of data workers 6.1 mln (2014) - increase 5.7% from 2013
  - Average number of data workers per company 9.5 - increase 4.4%
  - Gap between demand and supply 509,000 (2014) or 7.5%
- HLEG report on European Open Science Cloud (October 2016) identified need for data experts and data stewards
  - **Recommendation: Allocate 5% grant funding for Data management and preservation**
  - **Estimation: More than 80,000 data stewards (1 per every 20 scientists)**
  - Core Data Experts (as defined) need to be trained and their career perspective improved



- **Clash of cultures**

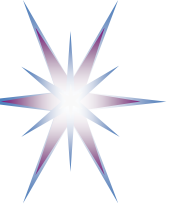
- between domain specialists and e-Infrastructure specialists (i.e. IT/Computer Science)

- New data experts come from **scientific** and **engineering** cultures

- with very different reward systems and incentives,
- different jargons and very different skill sets.

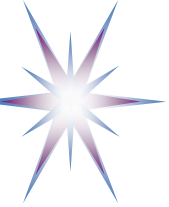
- **Consequences**

- Evident a divide between researchers and those that support research with data processing and software
  - Two communities that are both essential to Open Science have not closely co-evolved and do not converge
- Lack of data scientists that venture out from classical computer or data science departments into other scientific fields.



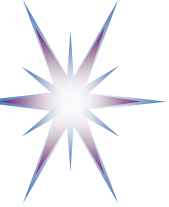
# HLEG EOSC Report Essentials – Core Data Experts

- **Core Data Experts** is a new class of colleagues with core scientific professional competencies and the communication skills to fill the gap between the two cultures.
  - **Core data experts** are neither computer savvy research scientists nor are they hard-core data or computer scientists or software engineers.
  - They should be technical data experts, though proficient enough in the content domain where they work routinely from the very beginning (experimental design, proposal writing) until the very end of the data discovery cycle
  - Converge two communities:
    - Scientists need to be educated to the point where they hire, support and respect Core Data Experts
    - Data Scientists (Core Data Experts) need to bring the value to scientific research and organisations
- Implementation of the EOSC needs to include instruments to help train, retain and recognise this expertise,
  - In order to support the 1.7 million scientists and over 70 million people working in innovation.



# GO FAIR and IFDS

- Global Open FAIR
  - Findable – Accessible – Interoperable - Reusable
- IFDS – Internet of FAIR Data and Services = EOSC
- GO FAIR implementation approach
  - GO-BUILD
  - GO-CHANGE
  - GO-TRAIN: Training of data stewards capable of providing FAIR data services
- A critical success factor is availability of expertise in data stewardship
  - Training of a new generation of FAIR data experts is urgently needed to provide the necessary capacity.



# EOSC Report Recommendations – Implementation on training and skills

- **I2.1: Set initial guiding principles to kick-start the initiative as quickly as possible. -> Bridge two cultures/communities**
  - A first cohort of core data experts should be trained to translate the needs for data driven science into technical specifications to be discussed with **hard-core data scientists and engineers**.
  - This new class of core data experts will also help translate back to the **hard-core scientists** the technical opportunities and limitations
- **I3: Fund a concerted effort to develop core data expertise in Europe.**
  - Substantial training initiative in Europe to locate, create, maintain and sustain the required core data expertise.
  - **By 2022, to train (hundreds of thousands of) certified core data experts** with a demonstrable effect on ESFRI/e-INFRA activities and prospects for long-term sustainability of this critical human resource
    - Consolidate and further develop assisting material and tools for Data Management Plans and Data Stewardship plans (including long-term preservation in FAIR status)
- **I7: Provide a clear operational timeline to deal with the early preparatory phase of the EOSC.**
  - **Define training needs for the necessary data expertise and draw models for the necessary training infrastructure**