



# ENSURING TRUSTED FAIR DATA FOR RESEARCH: INFRASTRUCTURE AND POLICY ISSUES

Carthage Smith, OECD Global Science Forum

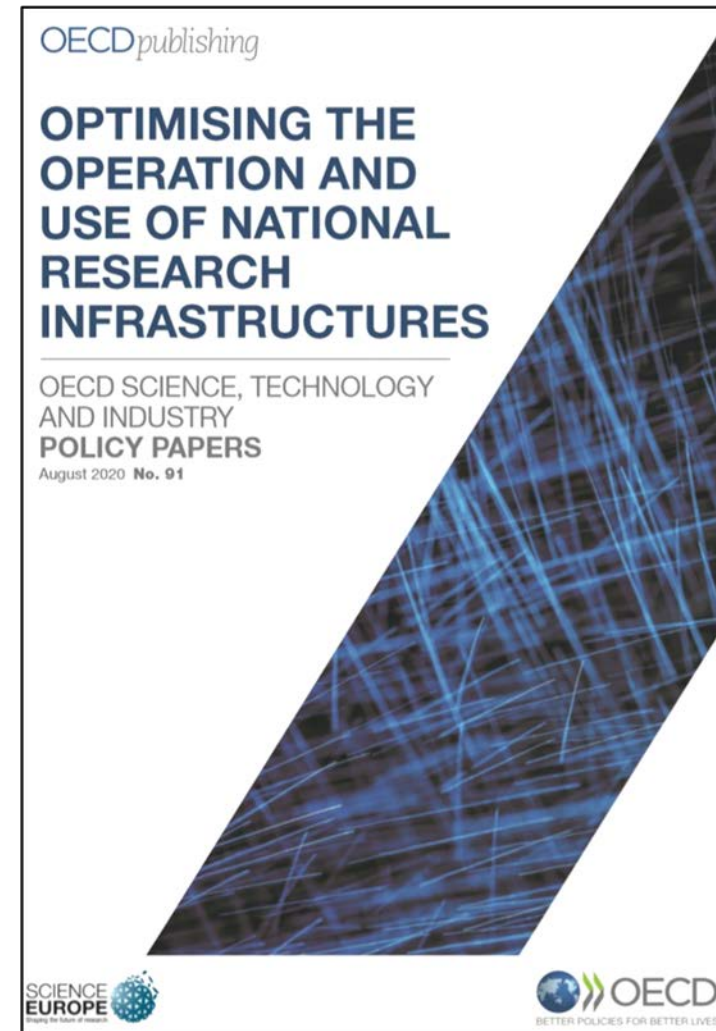


- I. Outcomes from the ICRI 2018 session on Research Infrastructures, data and trust
- II. Business models for sustainable data repositories
- III. International cooperation and data networks
- IV. Digital workforce capacity and skills for data intensive science



# I. Research Infrastructures and trusted data

e 2 0  
u 1 8  
- a t  
IC  RI





## RI ecosystem and data

---

- Trusted data and trusted cyber-infrastructure and trusted connections – need open PIDs
- Think in terms of the research data life-cycle
- Must design data services/cyber-infrastructure around specific user needs
- Build on best practices: different domains at different levels of preparedness
- Support development of open tools and software for connecting
- Develop a researcher-centric data ecosystem



# Challenges

---

- Plethora of competing standards
  - Standards for interoperability (disciplines – domains – cross-domain)
  - Gaps between HPC infrastructures and traditional repositories
  - Ownership of data
  - Data changes in real time – requires snapshots, versioning
  - Incentives and measurement - consequences of metrics
- Social, behavioural and technical barriers



# Human infrastructure

---

- Need to build a workforce of data scientists and stewards
- Need new cohorts of data specialists and stewards
- Need critical skills and education to deal with the data rich World of the future
- Technology (AI) can help but not replace human analysis.
- Generational latency because of mentorship. Need cultural change



# Google & Co.

---

- Google solns for storage and access to computing power meet a need
- Google are solving discovery challenge and driving semantic web, AI etc.

But

- Google pays more and is more attractive vs academia
- Google de facto sets the standards
- Google threaten to drive out (public) competitors
- ?long-term sustainability and customer-dependency/lock-in



## Funder (policy) perspectives

---

- On-line Poll: - ~70% of audience saw incentives and mandates as a critical function
- Data sharing policies are not the same as 'open data'.
- Balance national benefits and internat. cooperation?
- Communities need to triage and decide what must be kept and what is disposable
- Data life-cycle approach and importance of data management plans





## II. Who pays for FAIR data and how?





# Objectives

---

1. Identify and describe existing revenue sources and business models
2. Test potential business models with various stakeholders, including funders
3. **Make policy recommendations to promote sustainable business models** for research data repositories



# Elements of a successful business model



<b>Funding Source</b>	<b>Pros</b>	<b>Cons</b>
<b>Structural funding</b>	<ul style="list-style-type: none"> <li>• Compatible with open data principles.</li> <li>• Longer-term stability.</li> <li>• Larger-scale and efficiencies.</li> <li>• Flexible as to allocation.</li> </ul>	<ul style="list-style-type: none"> <li>• Fixed, multi-year may not scale easily.</li> <li>• Competes with research funding.</li> <li>• Too many eggs in few baskets.</li> </ul>
<b>Host or institutional funding</b>	<ul style="list-style-type: none"> <li>• Compatible with open data principles.</li> <li>• Longer-term stability.</li> <li>• Efficiencies through sharing services.</li> <li>• Close to researchers (customers).</li> </ul>	<ul style="list-style-type: none"> <li>• Limited purview, with focus on local community.</li> <li>• May lead to fragmentation of domain data and lower interoperability.</li> <li>• Limited incentive to add value to data and develop related services.</li> </ul>
<b>Data deposit fees</b>	<ul style="list-style-type: none"> <li>• Compatible with open data principles.</li> <li>• Demand oriented and scales with demand (data ingest).</li> <li>• Researchers price sensitivity ensures cost constraint.</li> <li>• Open data is part of research and its funding.</li> </ul>	<ul style="list-style-type: none"> <li>• Cost disincentive to depositing, so depends on strong mandates.</li> <li>• May lead to low level of curation to contain costs (price).</li> <li>• May be difficult for repository to compete for deposits with comparable repositories that do not charge.</li> </ul>
<b>Data access charges (subscriptions or use fees)</b>	<ul style="list-style-type: none"> <li>• Users pay for what they want, so funding reflects value.</li> <li>• More market-oriented approach may provide incentive for cost constraint.</li> </ul>	<ul style="list-style-type: none"> <li>• Not compatible with open data principles and many funder mandates, limiting the potential market size.</li> <li>• Charges limit use and will reduce the value of data.</li> <li>• Revenue scales with use and not ingest or curation costs.</li> <li>• Vulnerable to funding cuts.</li> </ul>
<b>Diversification of revenue sources</b>	<ul style="list-style-type: none"> <li>• No single source of failure.</li> <li>• Can maintain compatibility with open data principles.</li> <li>• Flexible and enables experimentation with new services.</li> </ul>	<ul style="list-style-type: none"> <li>• May lead to higher transaction costs (managing multiple funding sources).</li> <li>• May lead to Mission drift.</li> </ul>



# Recommendations

---

1. All stakeholders should recognise that **research data repositories are an essential part of the infrastructure for open science.**
2. All research data repositories should have a **clearly articulated business model** and value proposition(s).
3. Sponsors need to **consider the ways in which data repositories are funded**, and the pros and cons of various funding mechanisms in different circumstances.
4. Research data repository business models are **constrained by, and need to be aligned with, policy regulation** (mandates) and **incentives** (including funding).
5. In the context of financial sustainability, opportunities for **cost optimisation should be explored.**



# III. International networking





# Objectives

---

Establish principles and policy actions that can support open and sustainable international research data networks:

1. When is a data network needed?
2. How can governments use networks to maximize research data openness and reuse?
3. What is the best governance model for a particular network?
4. What interoperability arrangements are necessary for the effective operation of the network?
5. What business models can sustain a network over time?



# Findings

---

- The most successful networks have **engaged and supportive users** who clearly understand and value the services of the network.
- The top issue faced by data networks in open sharing of data is **the varying attitudes and policies across countries.**
- Different research communities require **different data networks** because **the cultures** of data sharing vary.
- The most difficult aspects of **interoperability** are rooted in **human relationships and trust.**
- Developing a **coherent and sustainable business model** is a central challenge for virtually all data networks.





# Recommendations

---

1. work toward **common definitions of, and agreements on, open data**. What is open data in different domains?
2. work toward commonly agreed and enforced **legal and ethical frameworks for the sharing of different types of public research data**.
3. Funders and host institutions should view internationally coordinated data networks as a **long-term strategic investment**
4. Networks should have clear business models, including value propositions and **measures of success** that are relevant to their different stakeholders and these measures **should be monitored**.



## IV. Human infrastructure





## Objectives

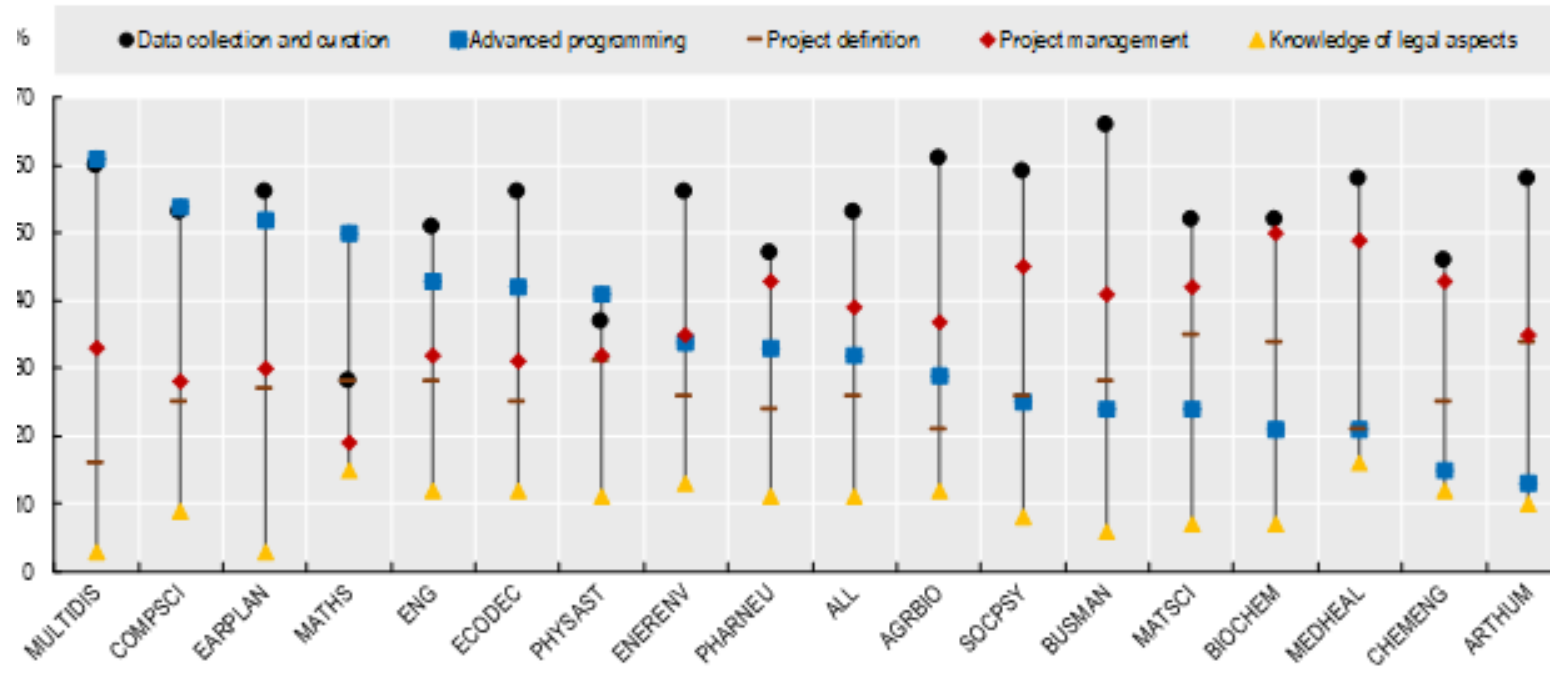
---

- Identify skill and capacity needs and gaps for data-intensive science in academia
- Identify policy actions to address these gaps
- Promote mutual learning and exchange of good practices



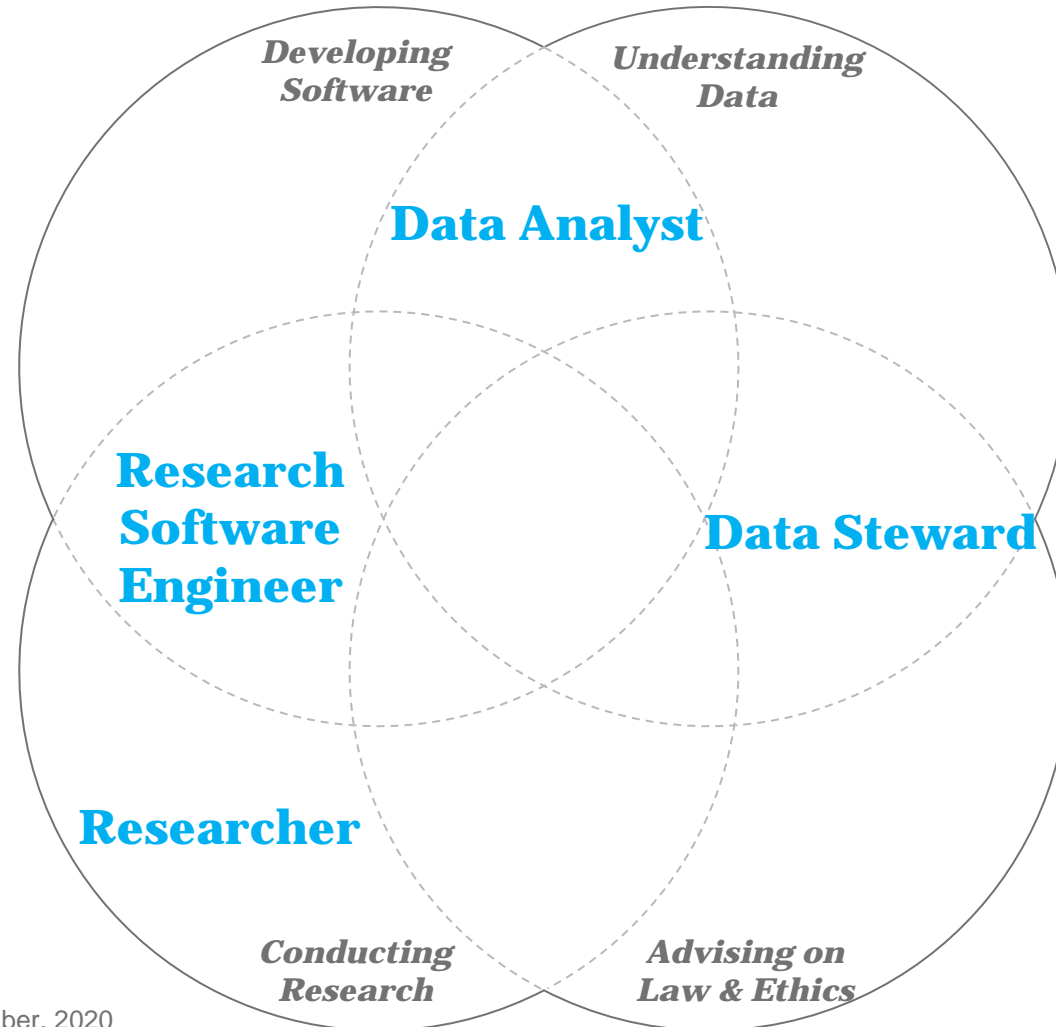
# Different needs in different research domains

**Figure 5.3. Most important skills for scientific authors' research work**  
Percentage of authors who deem each type of skill as important





# 1. Digital skills and stereotype roles





## 5 areas for (policy) action

Integrate digital workforce capacity development into broader science policy frameworks and actions, e.g. for open science and research integrity

Enablers for digital workforce capacity development

Identify the key competencies, skills and roles required for data-intensive science in different contexts.

Defining needs: digital skills, frameworks and roles

Career paths and reward structures

Implement changes in academic evaluation and reward systems in order to attract and retain diverse digitally skilled staff.

Provision of training

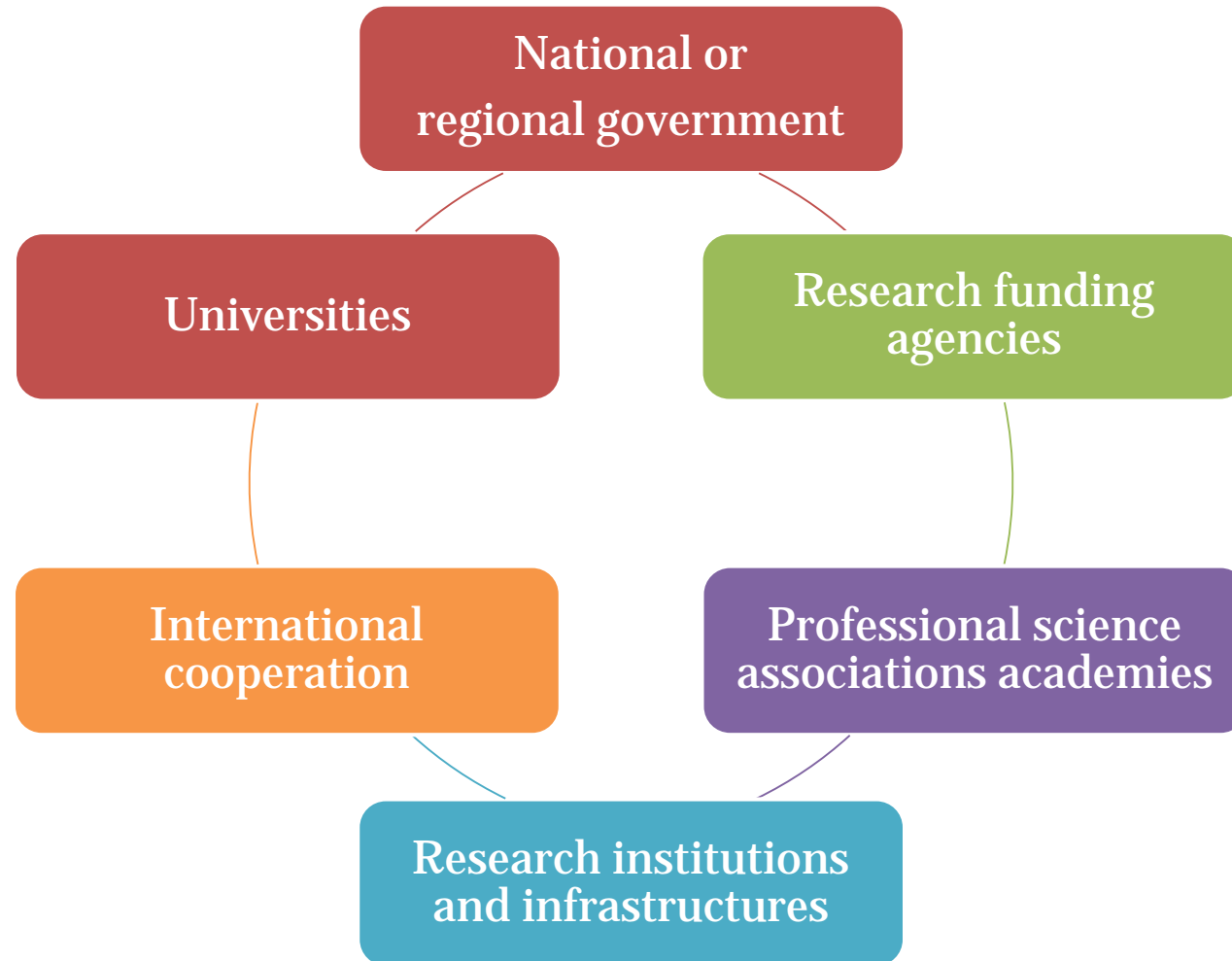
Support training in foundational digital skills and more specialized skills for scientists and research support professionals.

Community development

Support development of communities for new professional roles, learners and trainers.



# Requires complementary action at different scales and from different actors





# OECD policy reports



2015: Making Open Science a reality

2016: Research Ethics and New Forms of Data for social and economic research

2017: Business models for sustainable research data repositories

2017: Co-ordination and support of international research data networks

2020: Digital workforce capacity and skills for data intensive science

**2021: revised OECD Recommendation on Access to Research Data from Public Funding.**