

RESEARCH DATA MANAGEMENT

ΔΕΔΟΜΕΝΑ, ΣΚΥΛΛΑ, ΧΑΡΥΒΔΙΣ



*All cartoons courtesy of Jørgen Stamp,
Digitalbevaring.dk. CC BY 2.5.*

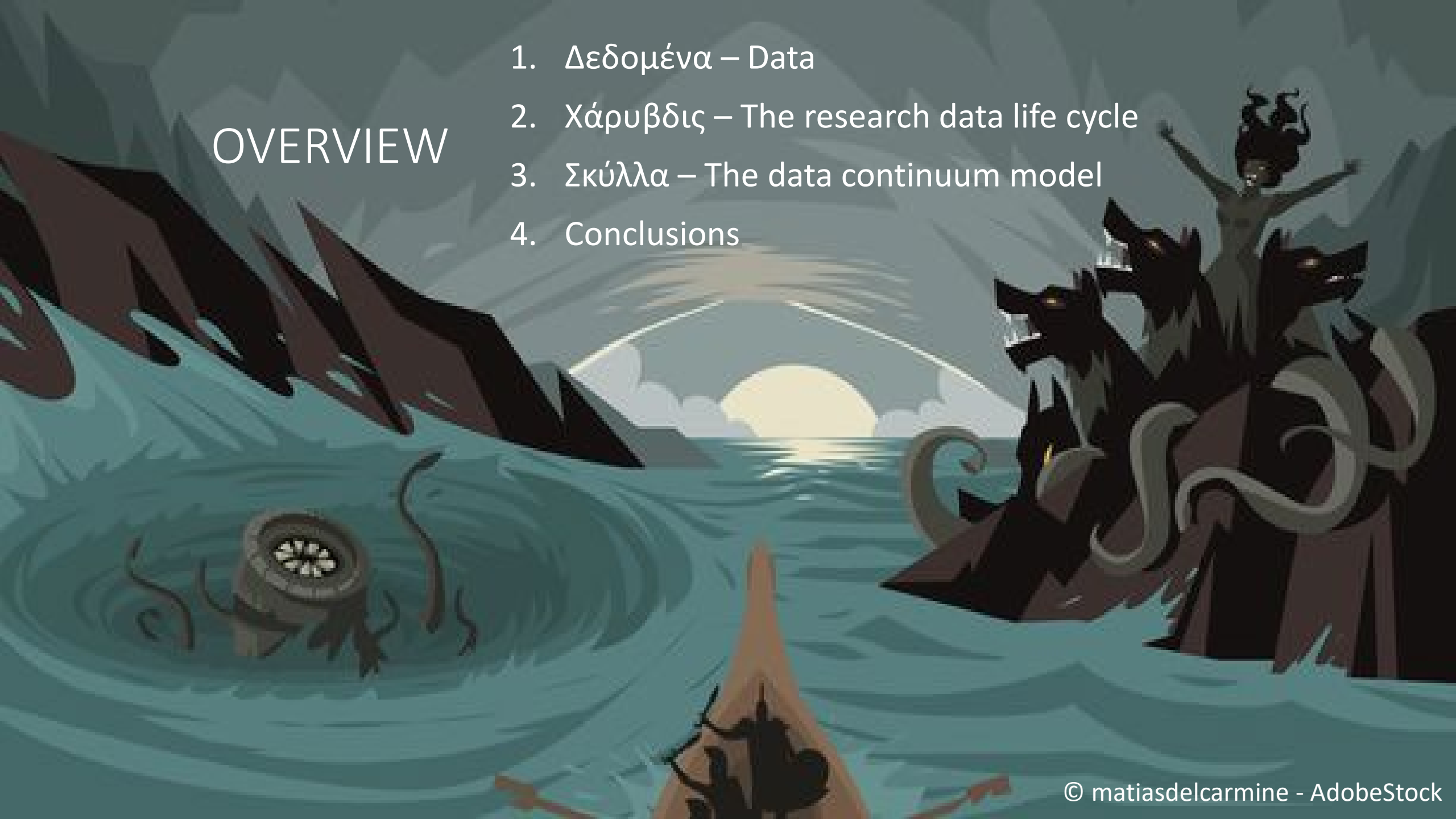
RENÉ SCHNEIDER

HAUTE ECOLE DE GESTION, GENEVA

LICENCE CC BY 4.0

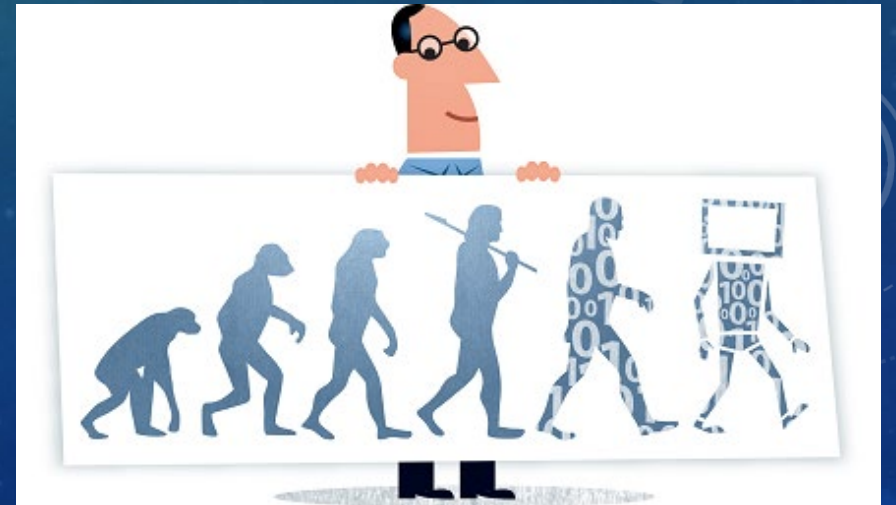
OVERVIEW

1. Δεδομένα – Data
2. Χάρυβδις – The research data life cycle
3. Σκύλλα – The data continuum model
4. Conclusions



1

Δεδομένα Data



EUKLID, 3RD CENTURY BC

Euklid wrote first the 15 books of the «Elements», still considered to be the origin of exact science.

Then he wrote «Δεδομένα», *Dedomena*, the Given, as a supplement to the «Elements».

Dedomena will then be translated to Latin as «Data» (unsurprisingly).

<https://archive.org/stream/lesuvresdeuclide03eucl#page/300/mode/2up>

EUCLIDIS

DATA.

ΟΡΟΙ.

DEFINITIONES.

α'. Δεδομένα τῶ μεγέθει λέγεται, χωρία τε, καὶ γραμμαὶ, καὶ γωνίαι, οἷς δυνάμει ἴσα περίσασθαι.

β'. Λόγος δεδύσθαι λέγεται, ᾧ δυνάμει τὸν αὐτὸν περίσασθαι.

γ'. Εὐθύγραμμα σχήματα τῶ εἶδει δεδύσθαι λέγεται, ὡν αἱ τε γωνίαι δεδομέναι εἰσὶ κατὰ μίαν, καὶ οἱ λόγοι τῶν πλευρῶν πρὸς ἀλλήλας¹ δεδομένοι.

δ'. Τῇ θέσει δεδύσθαι λέγονται², σημειῖά τε, καὶ γραμμαὶ, καὶ γωνίαι, ἃ τὸν αὐτὸν αἰὶ τὸπον ἐπέχει³.

1. Data magnitudine dicuntur, et spatia, et lineæ, et anguli, quibus possumus æqualia invenire.

2. Ratio dari dicitur, cui possumus eandem invenire.

3. Rectilineæ figuræ specie dari dicuntur, quarum et anguli dati sunt ad unum, et rationes laterum inter se datæ.

4. Positione dari dicuntur, et puncta, et lineæ, et anguli, quæ eundem semper situm obtinent.

LES DONNÉES

D'EUCLIDE.

1. Des espaces, des lignes, et des angles, auxquels nous pouvons trouver des grandeurs égales, sont dits donnés de grandeur.

2. Une raison est dite donnée, quand nous pouvons lui en trouver une qui soit la même.

3. Des figures rectilignes, dont chacun des angles est donné, et dont les raisons de leurs côtés entre eux sont données, sont dites données d'espèce.

4. Des points, des lignes, et des angles qui conservent toujours la même situation, sont dits donnés de position.

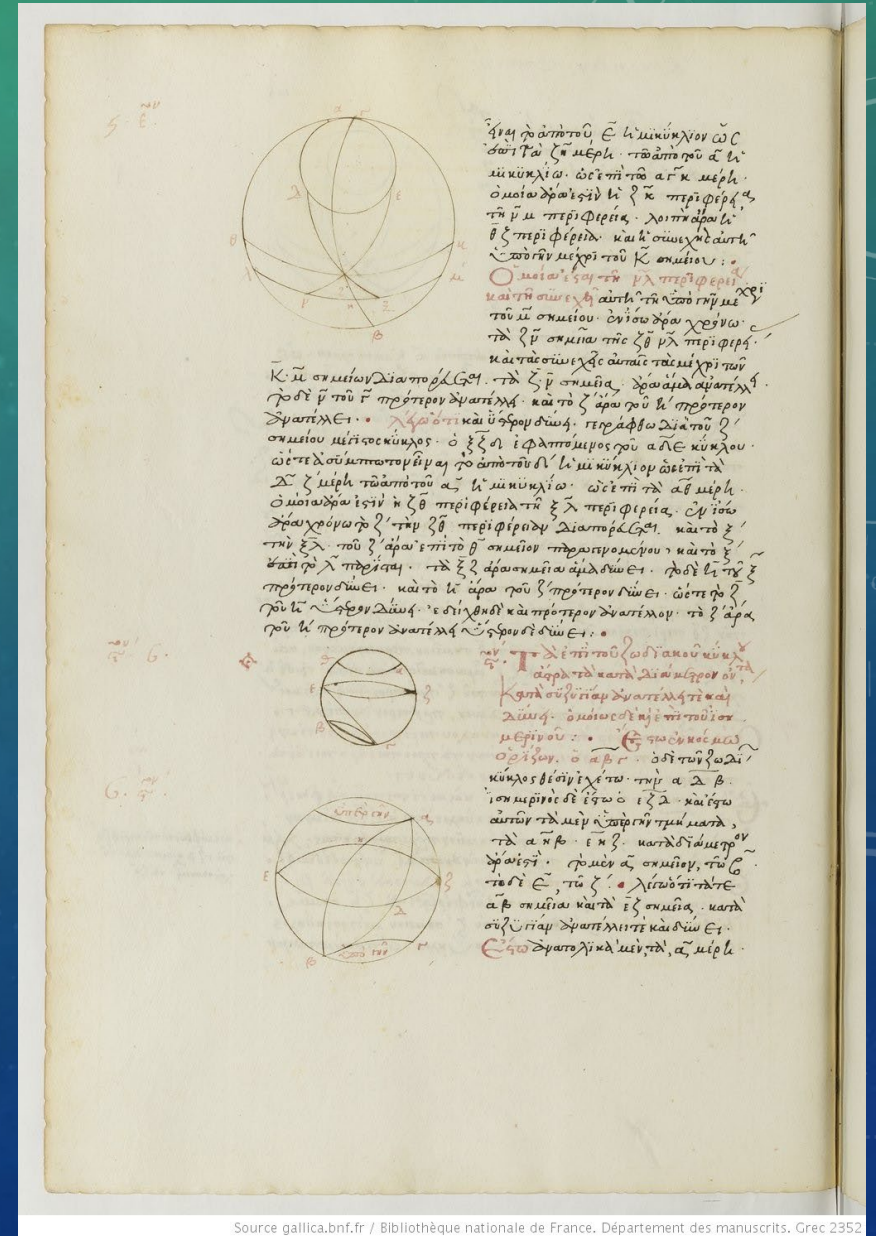
DEDOMENA

In the end, Dedomena is a collection of geometric axioms, 94 in numbers and following a strict pattern:

Given A, than B is also given.

A kind of description or deduction of «known unknowns».

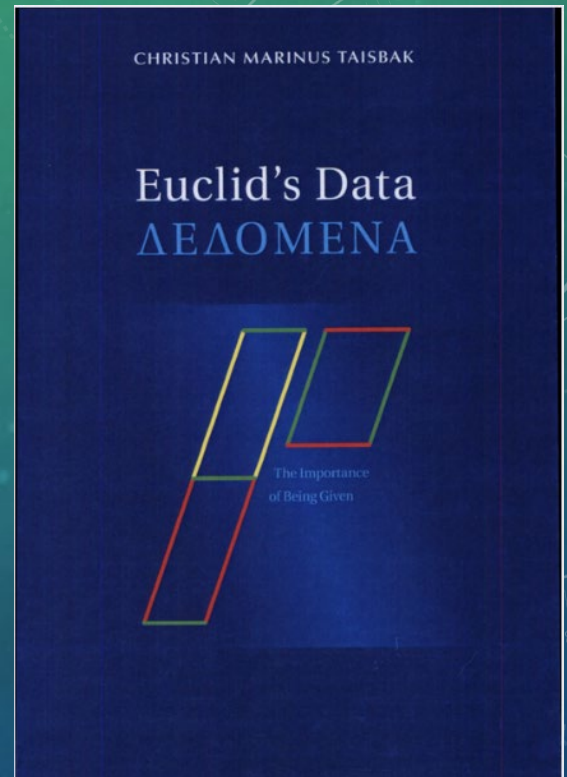
<https://gallica.bnf.fr/ark:/12148/btv1b52508678v/f224.item>



CHRISTIAN MARINUS TAISBAK

Euclid's *De Domoena*. The Importance of Being Given.
2003. p. 14:

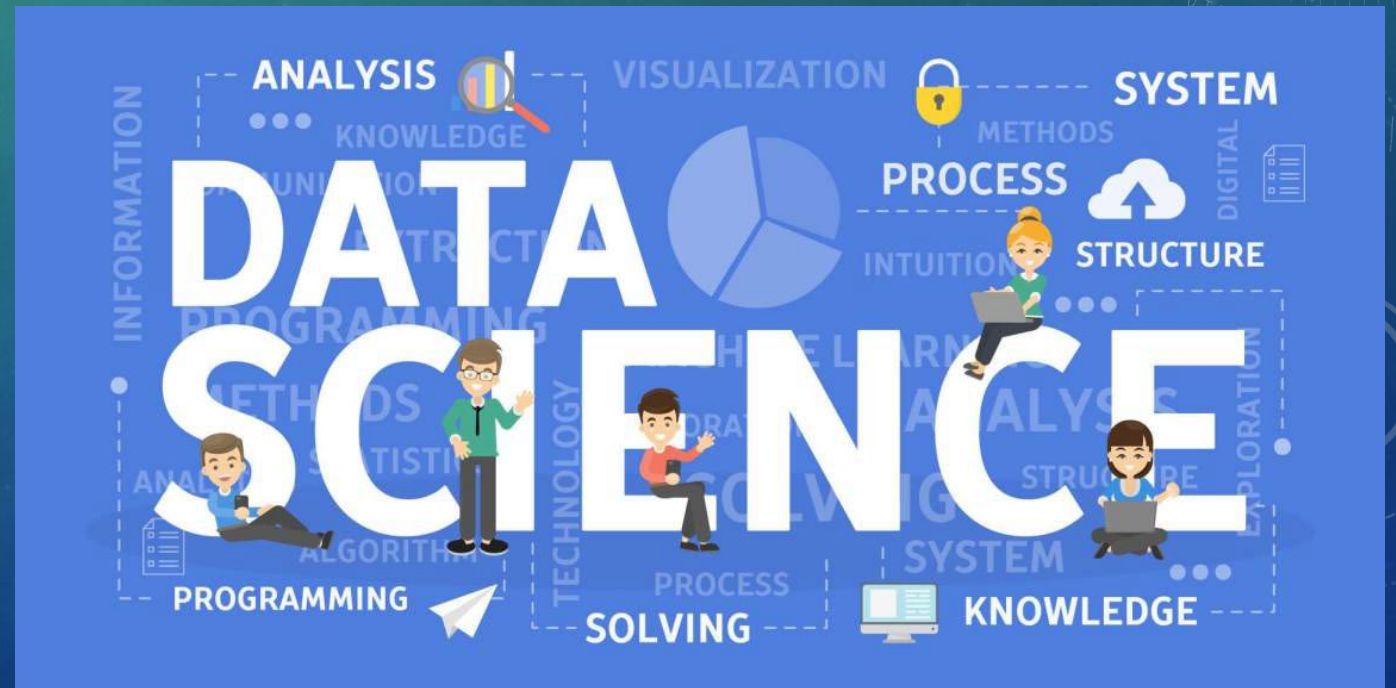
When I started to translate the Data, I found it very longwinded that a certain phrase kept popping up time and again, several times in every proposition: if this item is given, that item is also given. I decided to cancel all those alsos... But then I discovered that I was leaving out an essential feature of the "Data": the Givens hang together in chains, the purpose of any proposition being to produce more links to them.



DATA SCIENCE

Today, Data Scientists try to do the same thing
(maybe on a more complex level!)

Some would say, that
they try to find the
«unknown unknowns».

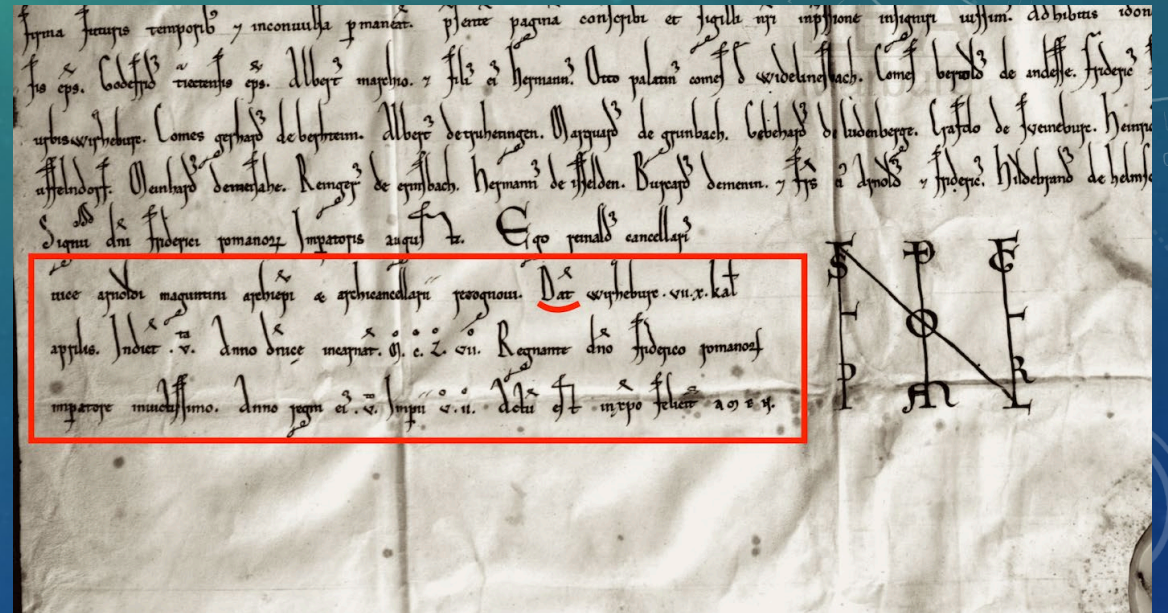


MIDDLE AGES

During the middle ages, data functions as a time stamp in documents and acts:

Pattern:

This document was created ...



17TH CENTURY

The word data reappears in the middle of the 17th century, interestingly in literature from Mathematics or Divinity.

Data is understood as «anything that is given, either as mathematical assertion or by the word of God».

Soon the term is used in other disciplines.

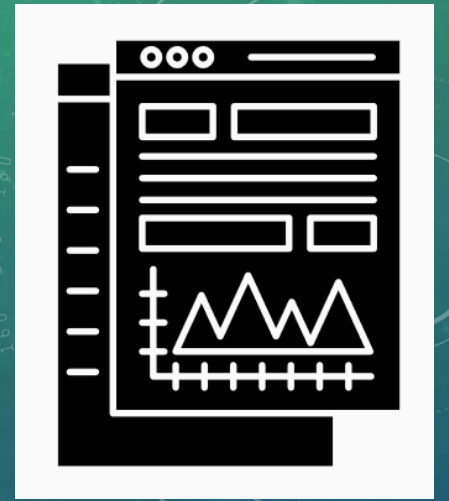
object of the first importance to an historian. On this account it is of consequence that every reader of history have it in his power to form a just idea of them from the *data* he finds in historians, and that he be guarded against the mistakes which, without some previous instruction, he would unavoidably fall into with respect to them.

R

I shall

Joseph Priestley: Lectures on History, and General Policy, 1788, S. 121

18TH CENTURY



At the end of the 18th century, the use of data becomes limited and is only found in relation with experiments, experience or the collection of facts.

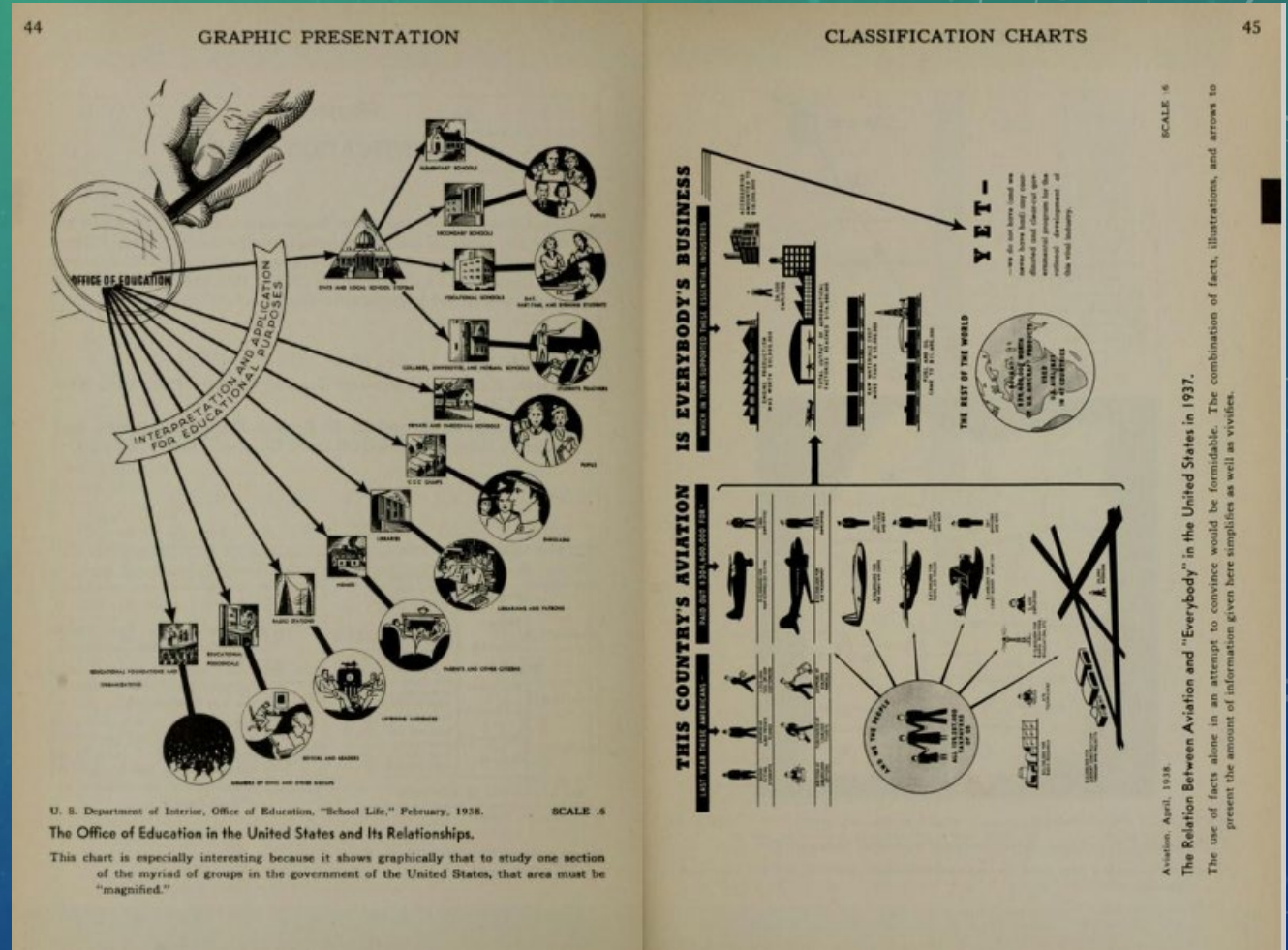
«It had become usual to think of data as the result of an investigation rather than its premise.»

Daniel Rosenberg, [“Data before the Fact”](#), in: Gitelman, Lisa (ed.): “Raw Data” is an Oxymoron. Cambridge/Mass.: MIT Press, 2013, p. 33.

MODERN TIMES

Data in modern times is firstly to be seen as related to descriptive statistics.

And its respective visualization.

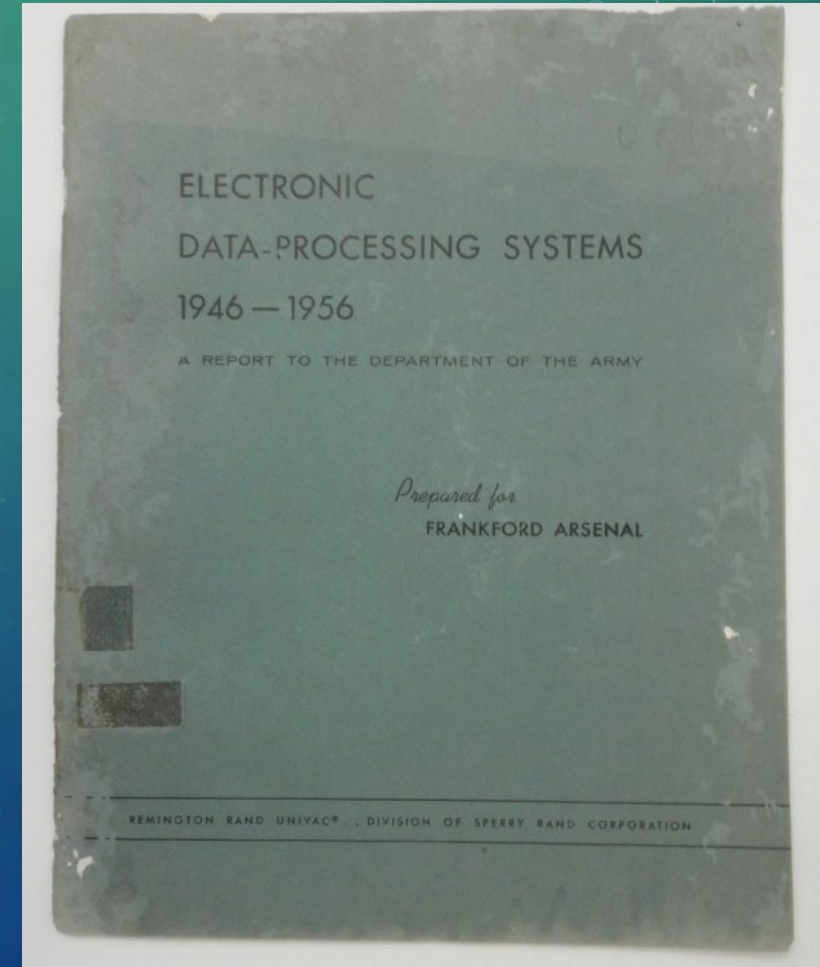


Willard Cope Brinton: [Graphic Presentation](#). 1939.

MODERN TIMES (C'TUED)

3994

Lp	A	B	C	A	B	C	Lp	Ch	n	Gn	Az	Ci	Cl	SM	If	HM	WI	A	C	E	F	G	d
Cn	D	E	F	D	E	F	Lp	Ch	n	Gn	Az	Ci	Cl	SM	If	HM	WI	A	C	E	F	G	d
Lo	G	H	I	G	H	I	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Cn	K	L	M	K	L	M	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
CS	N	O	P	N	O	P	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
LS	Q	R	S	Q	R	S	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
Kn	4	5	6	4	5	6	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
AN	7	8	9	7	8	9	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
QC	0	1	2	0	1	2	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
AV	3	4	5	3	4	5	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
So	6	7	8	6	7	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8



Secondly, with the invention of computers, data gains even more importance in the context of what is called Electronic DATA Processing.

Electronic Data-Processing Systems 1946-1956 : A report to the Department of the Army Prepared for Frankford Arsenal. no author.

DATA SCIENCE

The current state of statistical work can be described by a **Statistical Trilogy**:

1. Data Collection (experimental design, sample surveys)
2. Data Modeling and Analysis
3. Problem Understanding/Solving, Decision Making

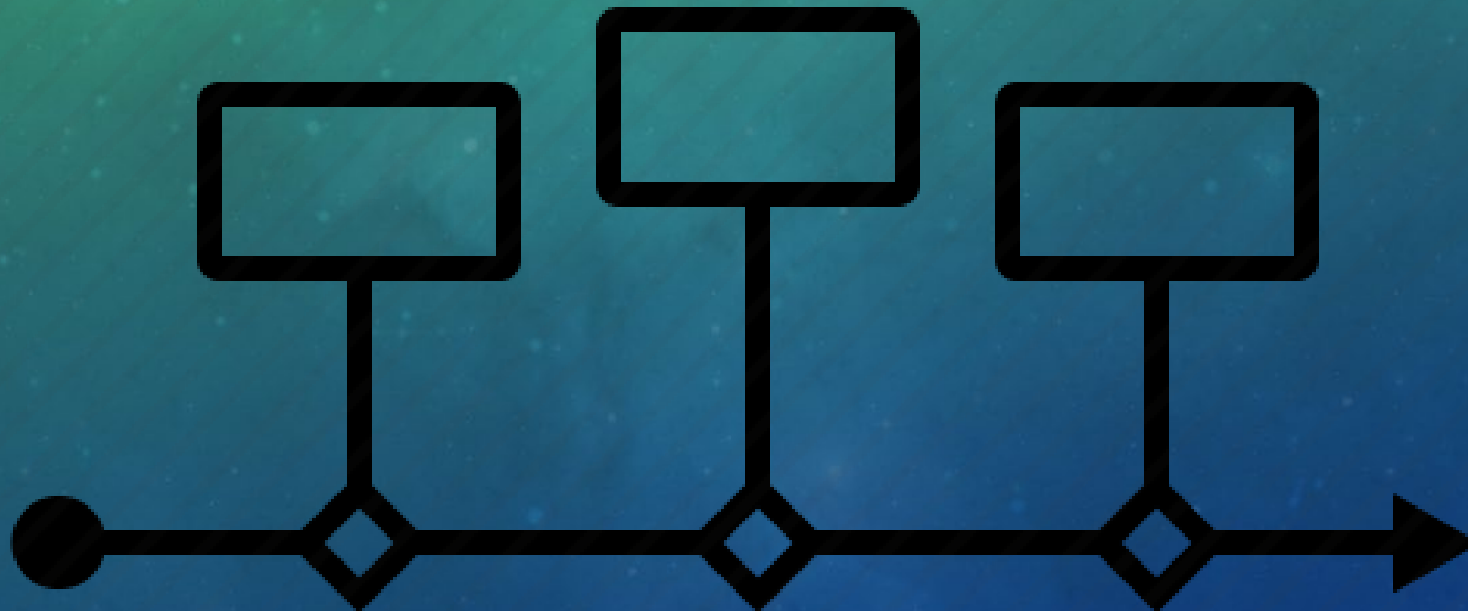
Wu, C. F. J. (1997). ["Statistics = Data Science?"](#) (PDF).

DATA OVER TIME

DESCRIPTIVE

PRESCRIPTIVE

PREDICTIVE



WHAT WE SEE SO FAR... IS THAT

The term data has a long history. As long as science.

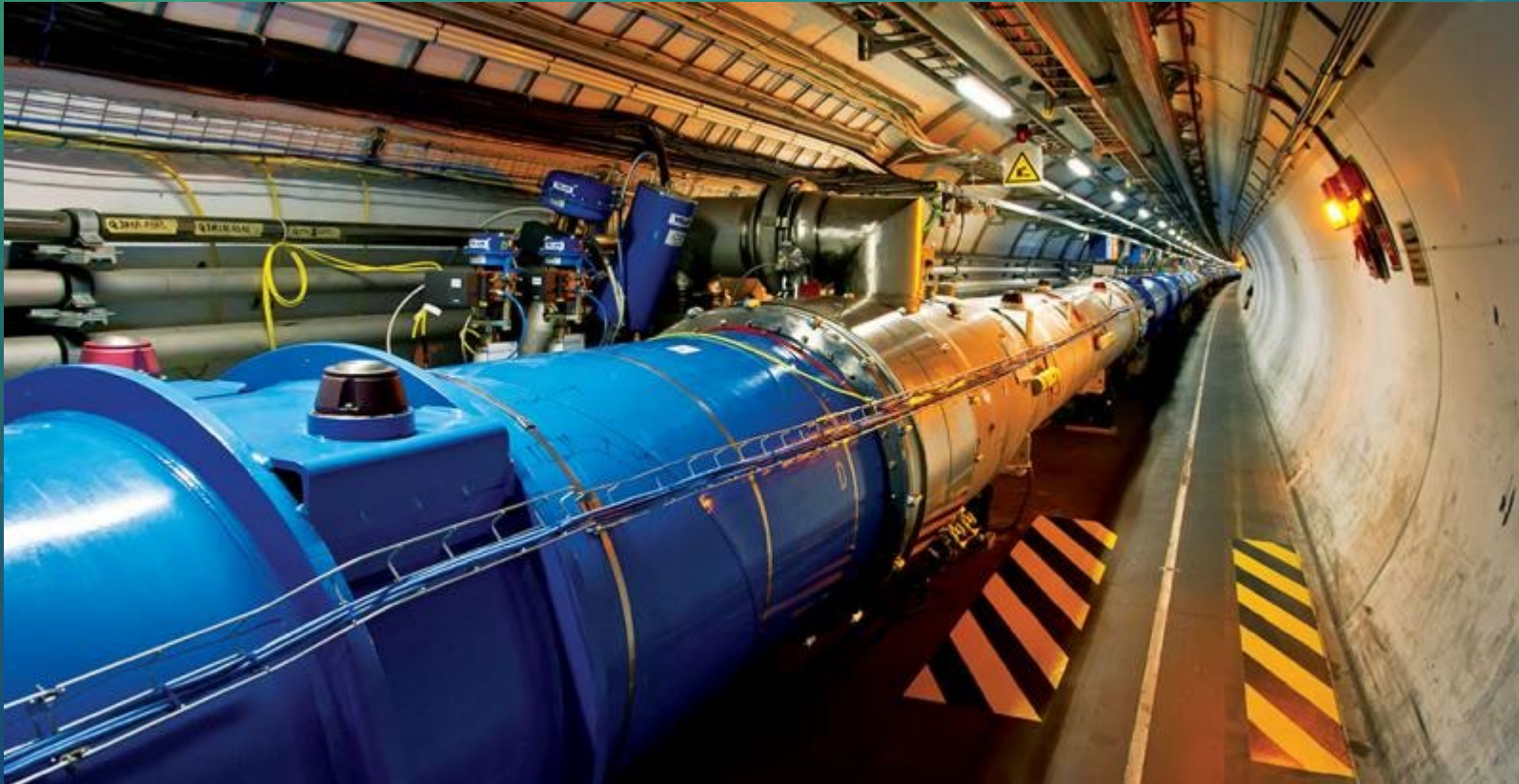
Data do have many facets.

One facet is research data. And its management.

RESEARCH DATA - DEFINITION

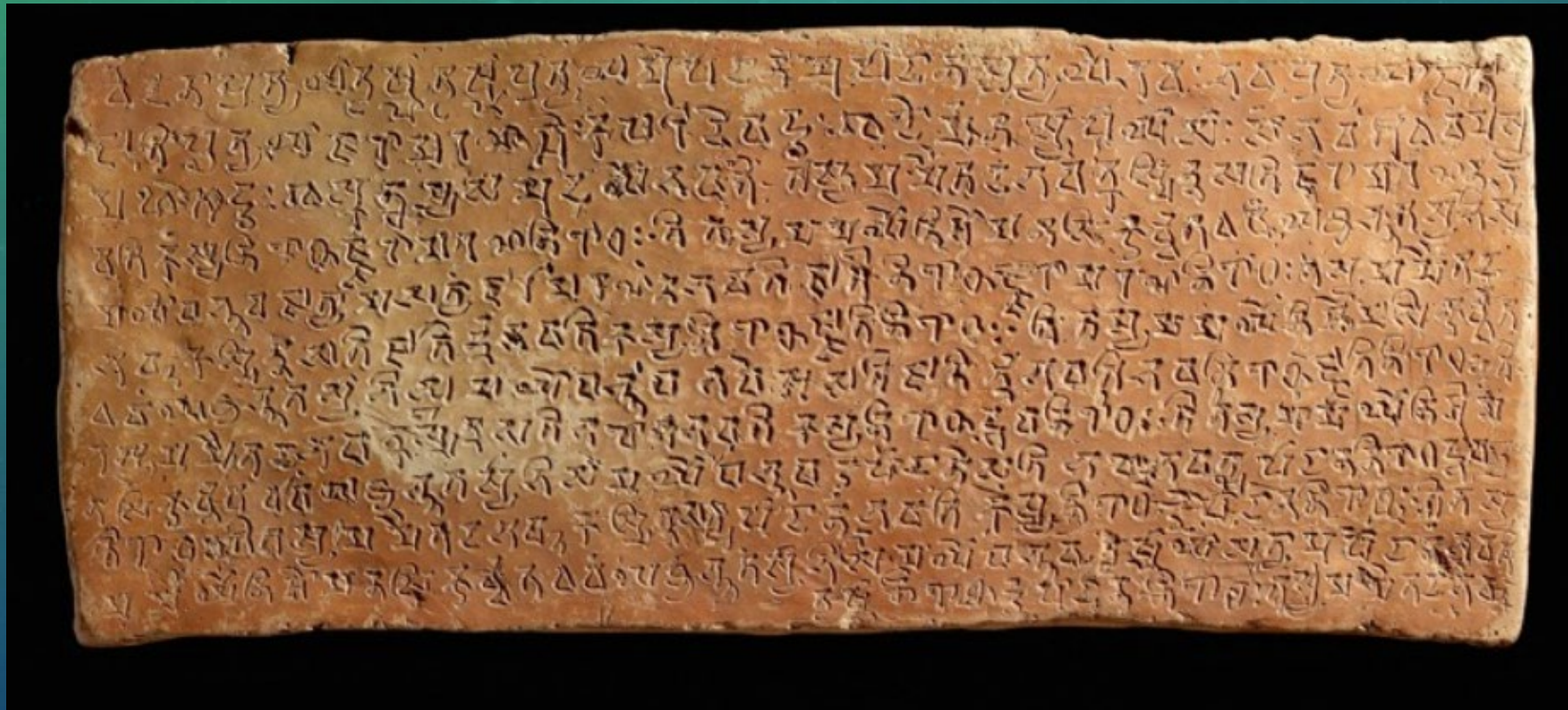
- « any information in binary digital form». Includes digital objects and data bases. (Harvey Ross. Digital Curation. 2010)
- Research data are collections, observations, models, measurements, references, digitizations etc.
- The definition of research data strongly depends on the scientific or disciplinary context of their creation and the scientific perspectives taken or methods used.

RESEARCH DATA EXAMPLE 1



Large Hadron Collider, CERN 1) Raw 2) Reconstructed 3) Reduced 4) Published

RESEARCH DATA EXAMPLE 2



Brick inscribed with the sutra of dependent origination
Gorakhpur, 5-6^{ème} siècle

RESEARCH DATA: EXAMPLE 3

STANFORD MARSHMALLOW PROJECT - 1960



STANFORD MARSHMALLOW PROJECT - REVISION

1988

1990

2006

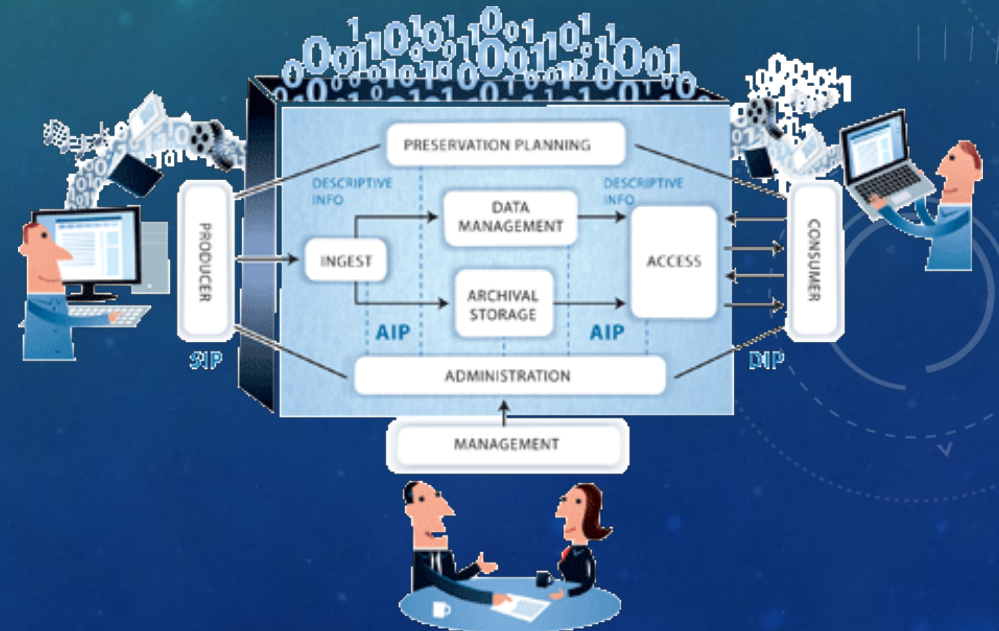
2011 ...

-> The data collected in 1960 must have been stored.

But how? On which medium? In which format?

2

Χάρυβδς The data life cycle



SIMPLE CYCLE

<http://datasupport.researchdata.nl>



DATA CURATION



“active and on-going management of data through its lifecycle of interest and usefulness to scholarship, science, and education” through activities that “enable data discovery and retrieval, maintain quality, add value, and provide for re-use over time.”

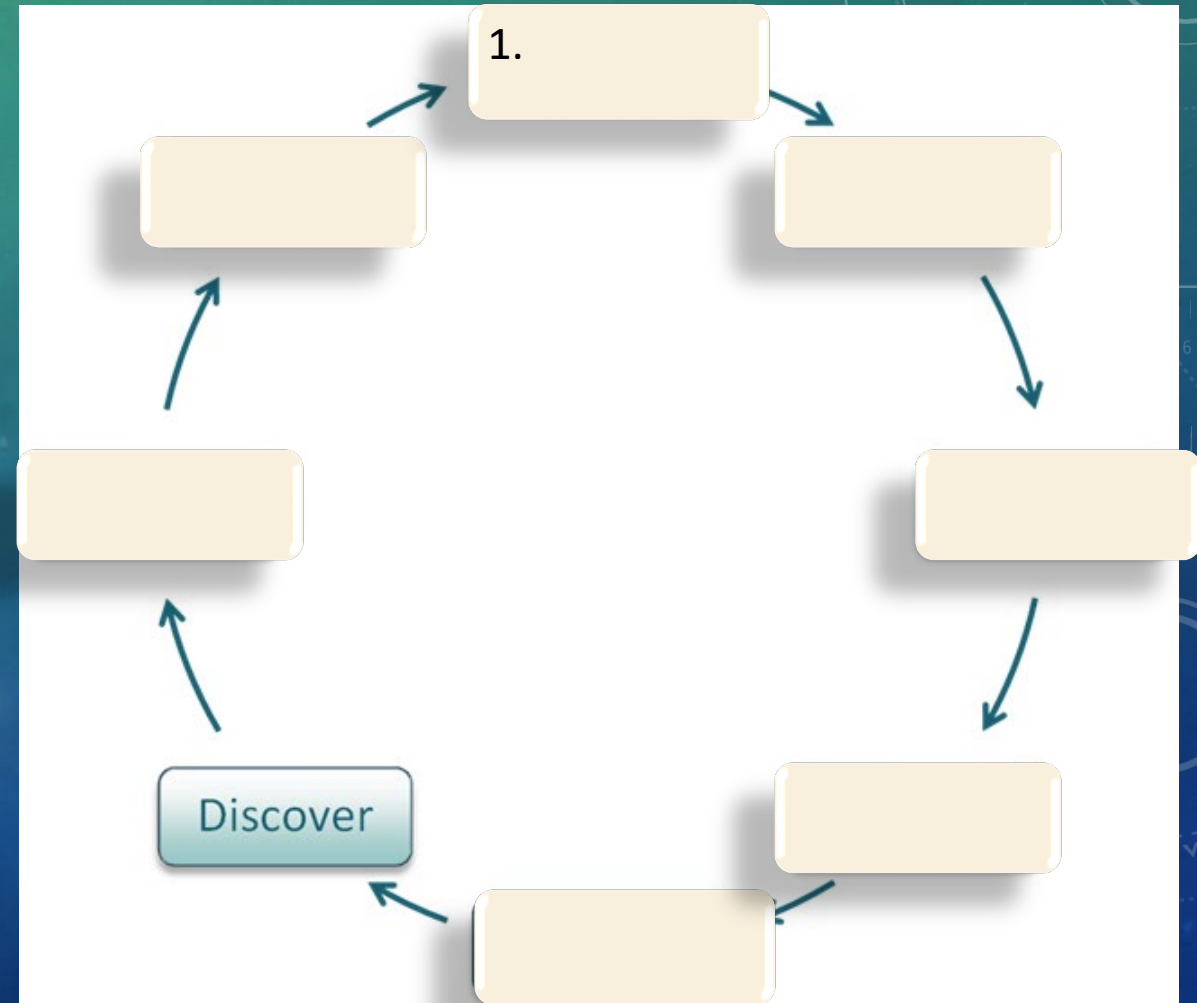
<http://hdl.handle.net/2142/3493>

EXERCICE LIFE CYCLE

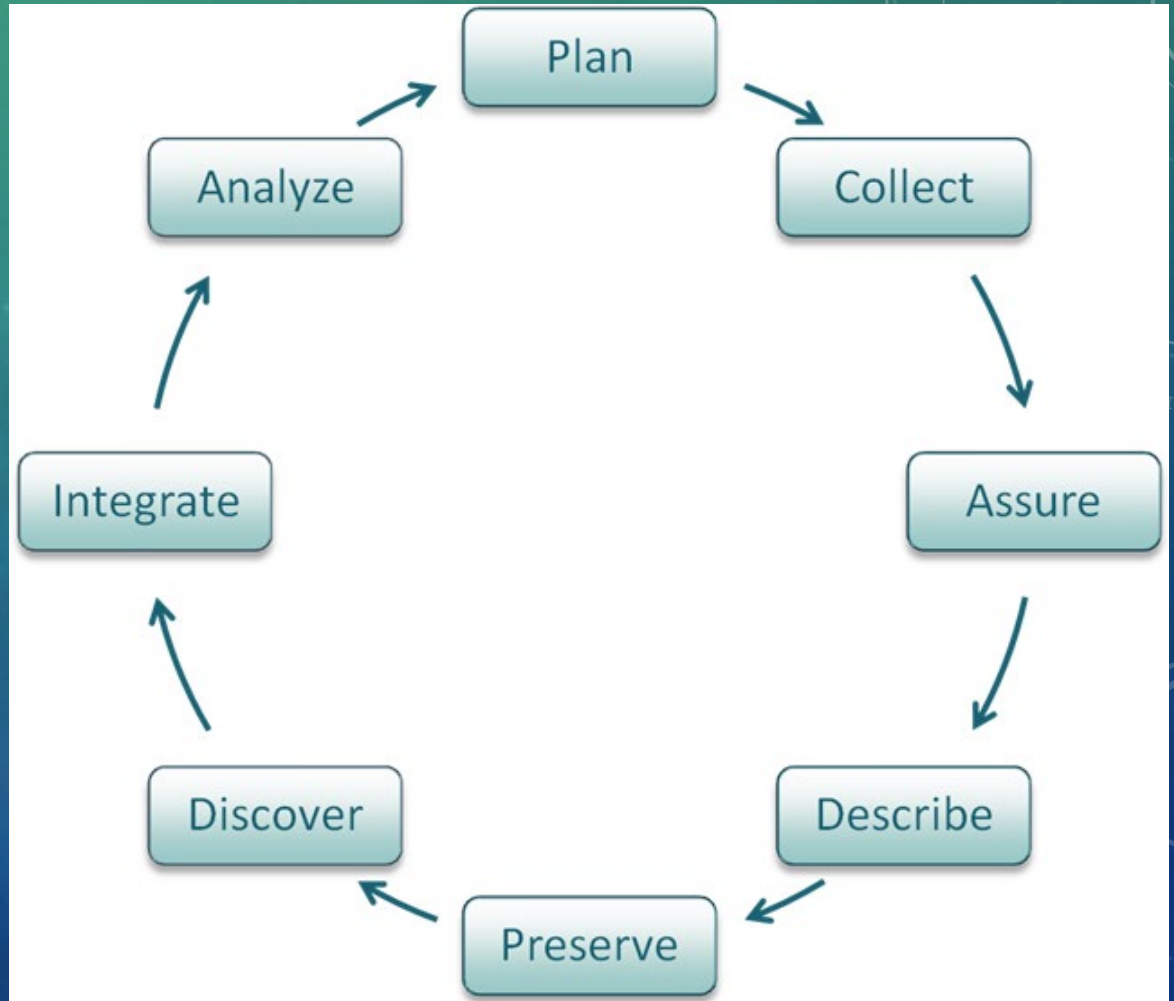


Try to put the actions of the data curation life cycle in the right order!

RESEARCH DATA LIFE CYCLE

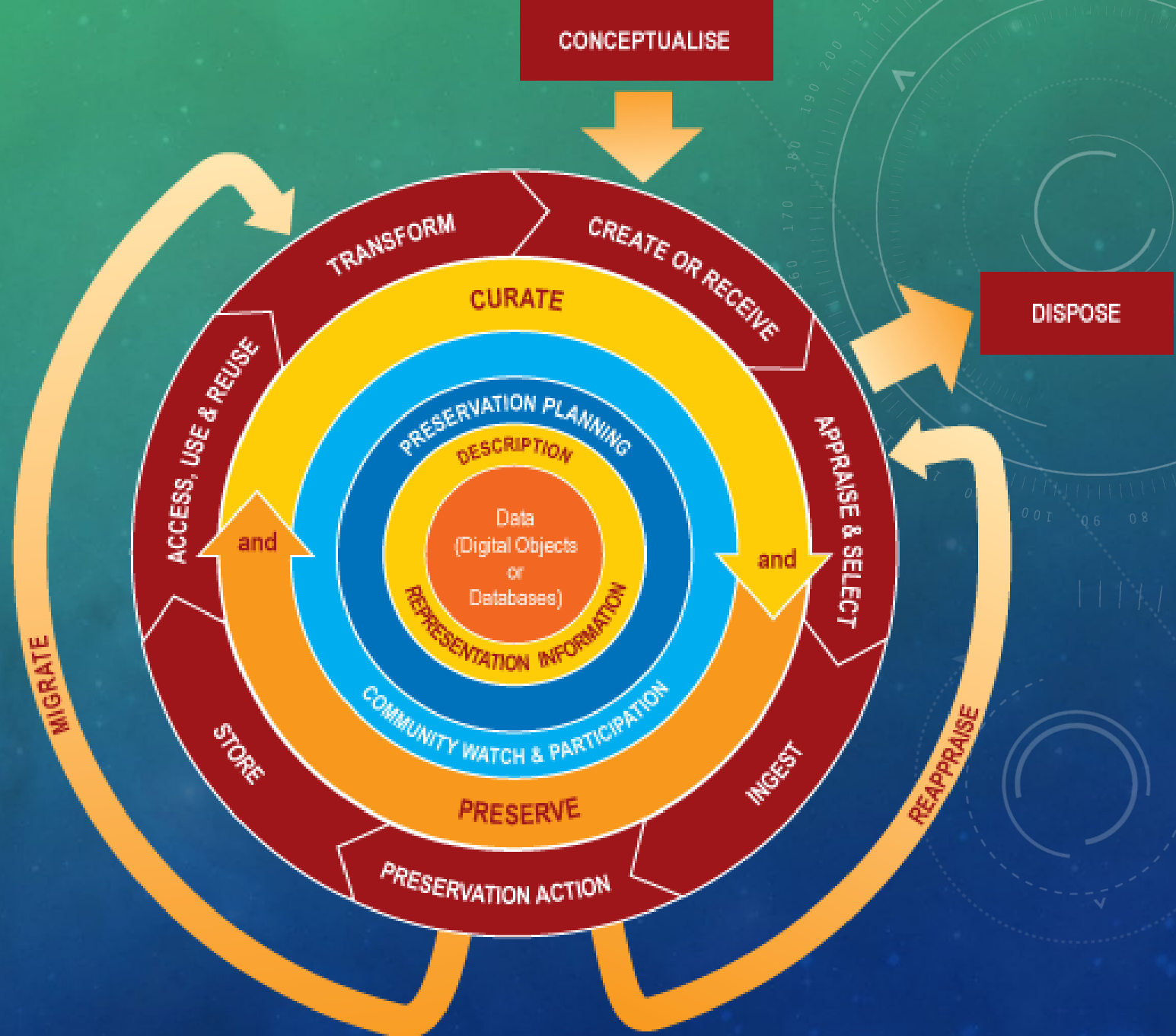


RESEARCH DATA LIFE CYCLE



Reference model Data ONE Project
<https://www.dataone.org/best-practices>

THE DCC CURATION LIFECYCLE MODEL



3

Σκύλλα

The data continuum model

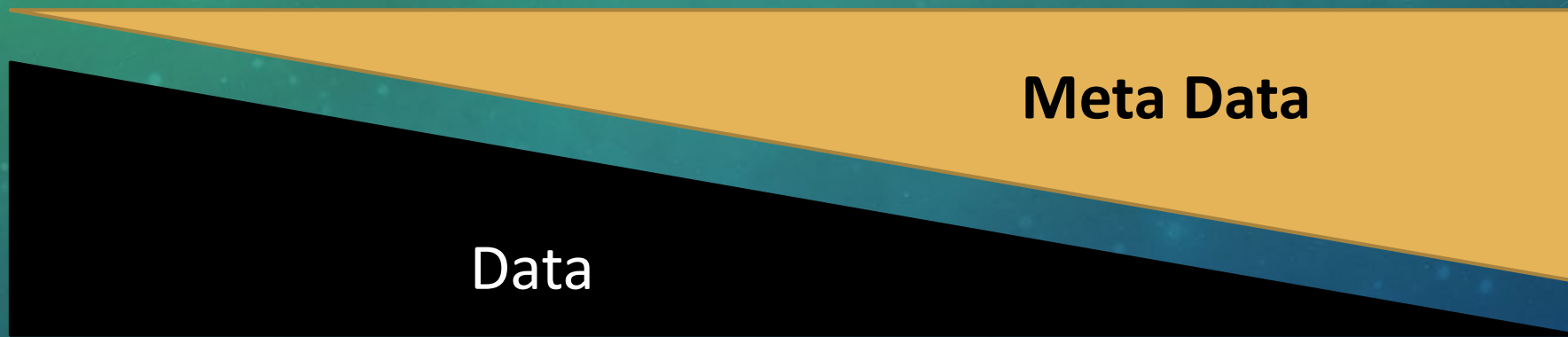


DATA CONTINUA

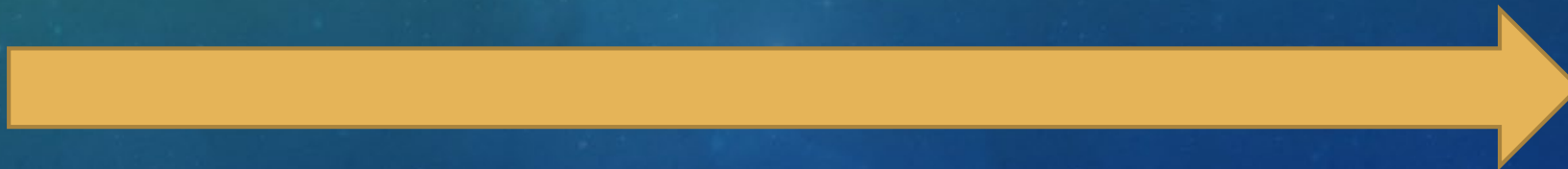
Object:	Less Metadata	←	→	More Metadata
	More Items	←	→	Fewer Items
	Larger Objects	←	→	Smaller Objects
	Objects continually updated	←	→	Objects static/derived snapshots
Management:	Researcher Manages	←	→	Organisation Manages
	Less Preservation	←	→	More Preservation
Access:	Mostly Closed Access	←	→	Mostly Open Access
	Less Exposure	←	→	More Exposure

<http://www.dlib.org/dlib/september07/treloar/treloar-table1.png>

DATA AND METADATA

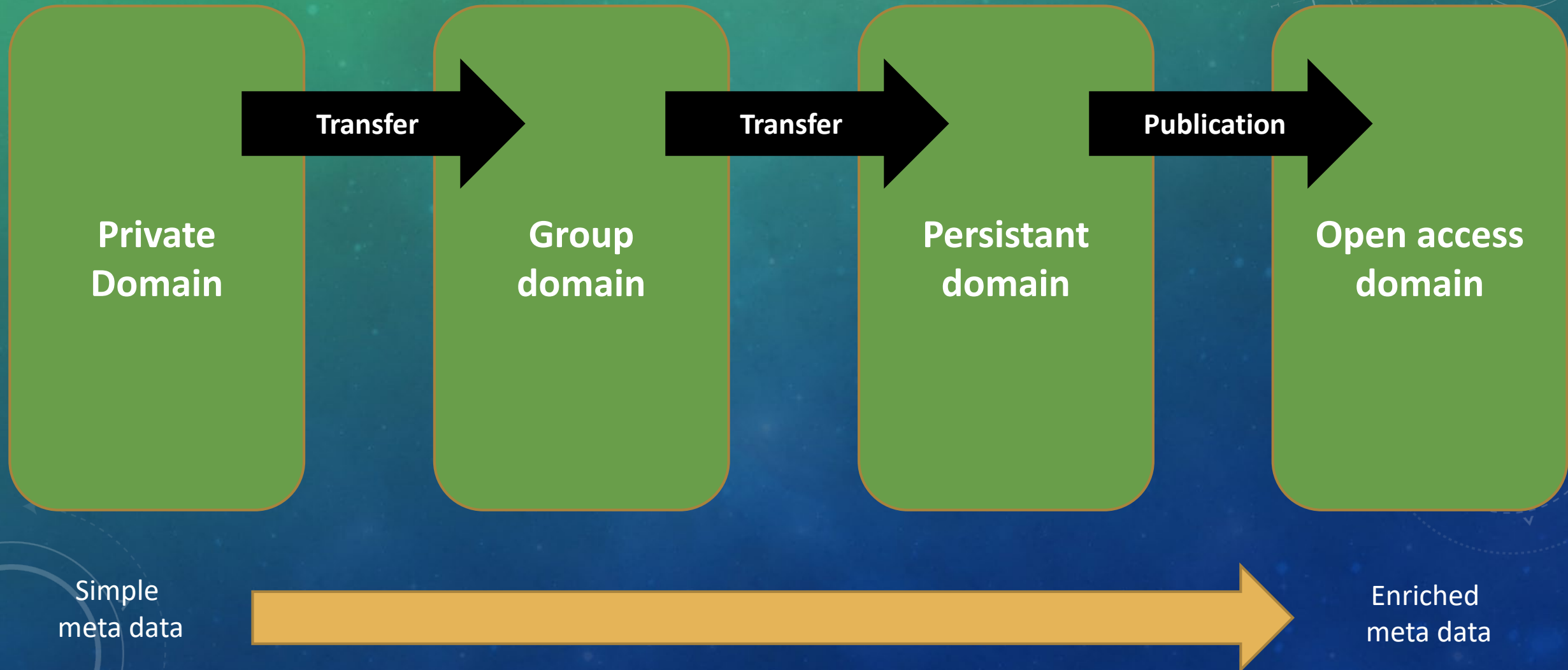


Simple
Metadata

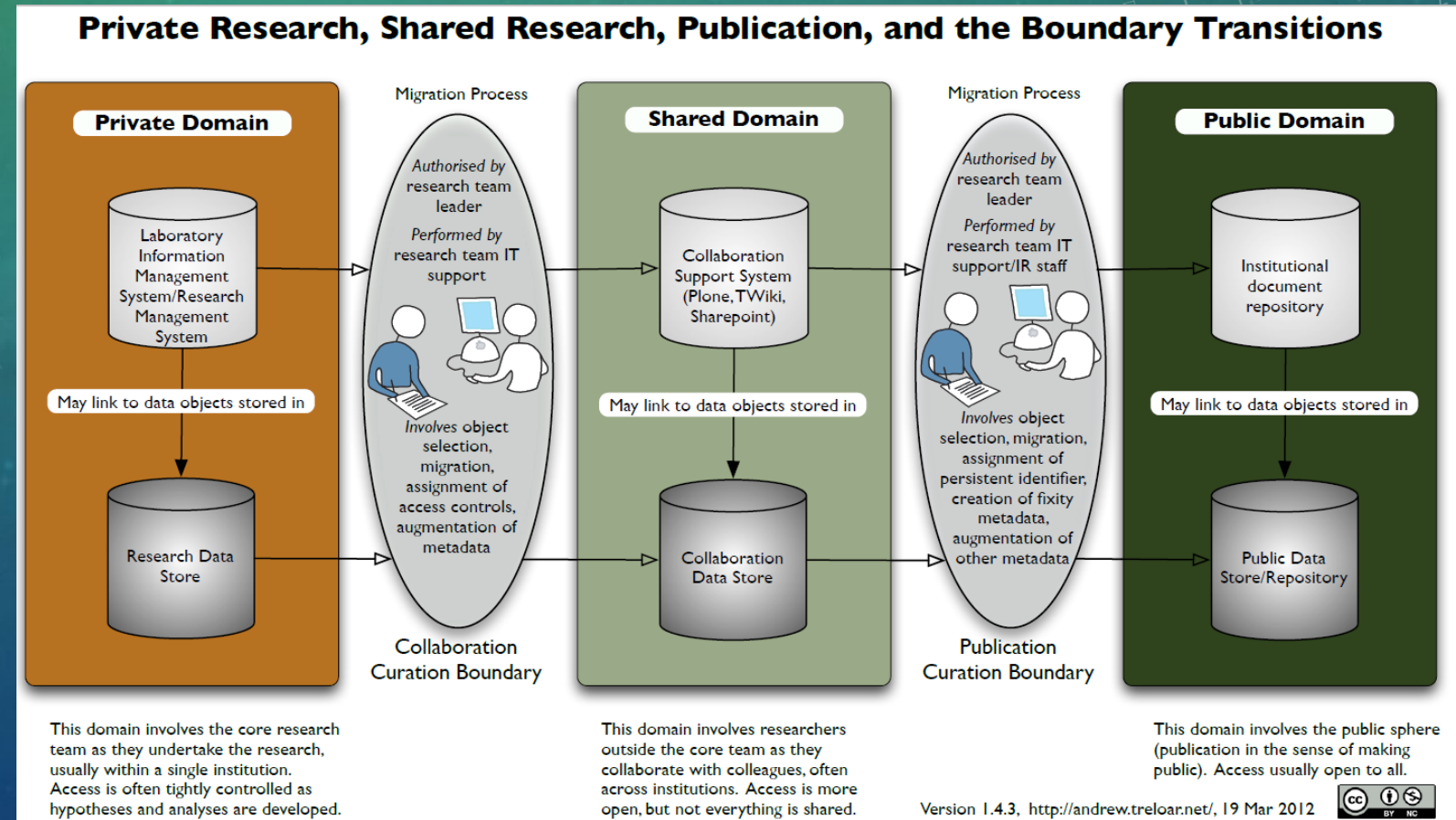


Enriched
Metadata

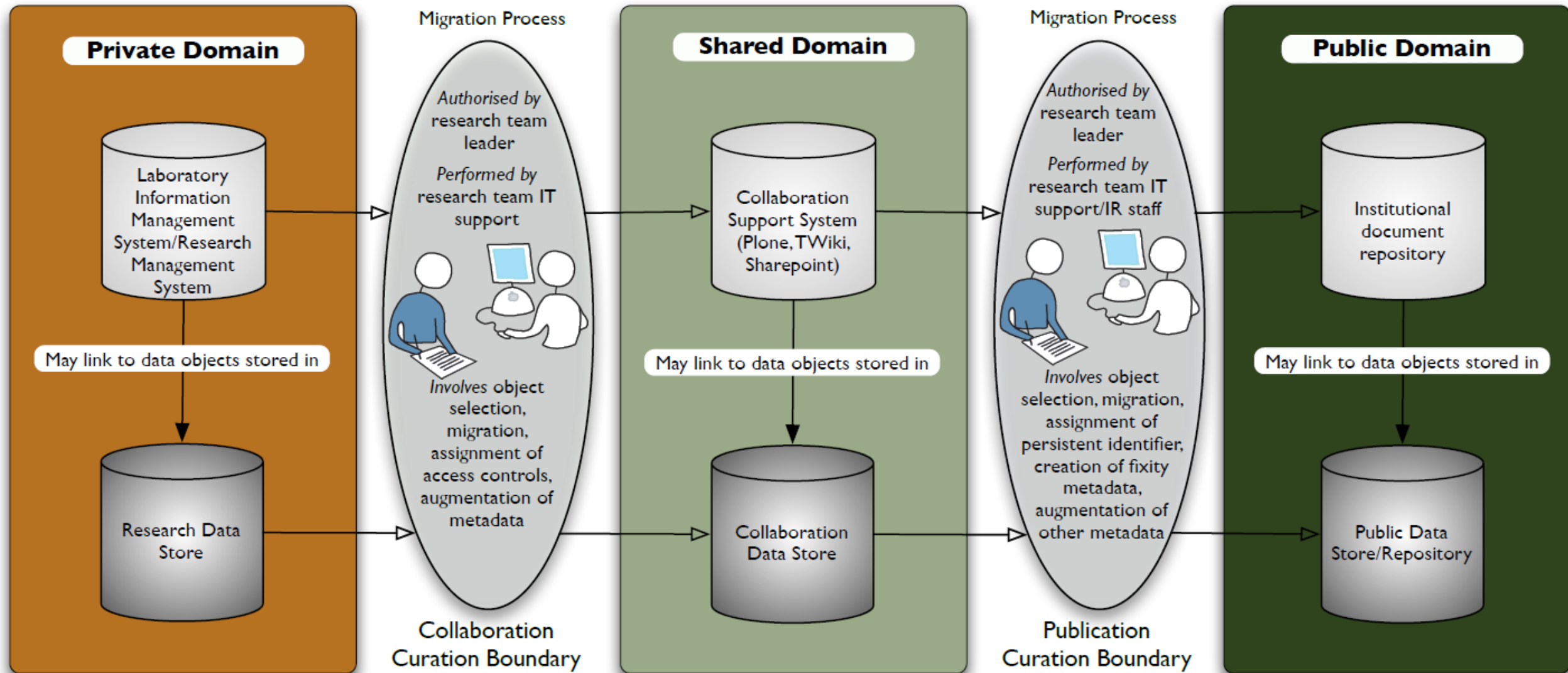
DATA CONTINUUM MODEL



DATA CONTINUUM MODEL



Private Research, Shared Research, Publication, and the Boundary Transitions

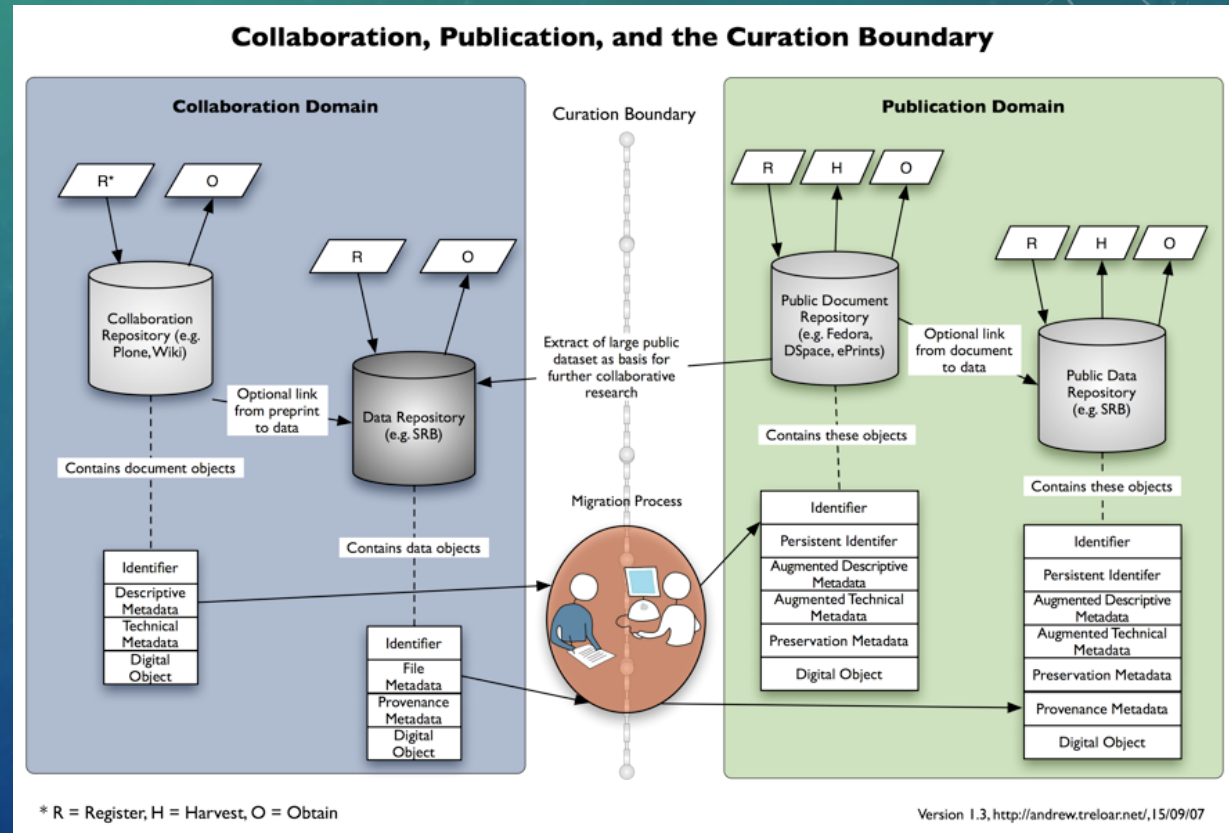


This domain involves the core research team as they undertake the research, usually within a single institution. Access is often tightly controlled as hypotheses and analyses are developed.

This domain involves researchers outside the core team as they collaborate with colleagues, often across institutions. Access is more open, but not everything is shared.

This domain involves the public sphere (publication in the sense of making public). Access usually open to all.

COLLABORATION, PUBLICATION AND THE CURATION BOUNDARY

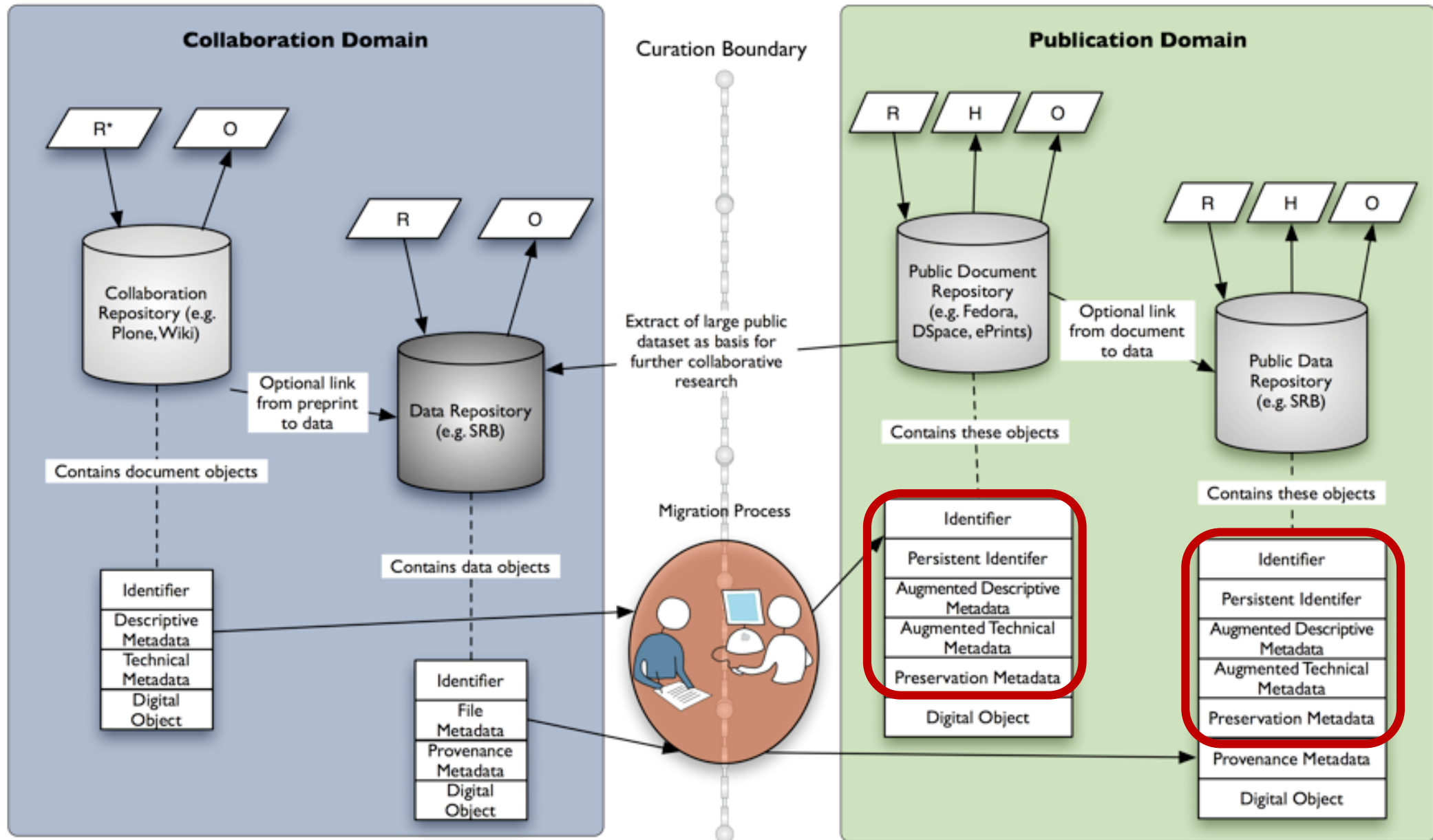


COLLABORATION, PUBLICATION AND THE CURATION BOUNDARY

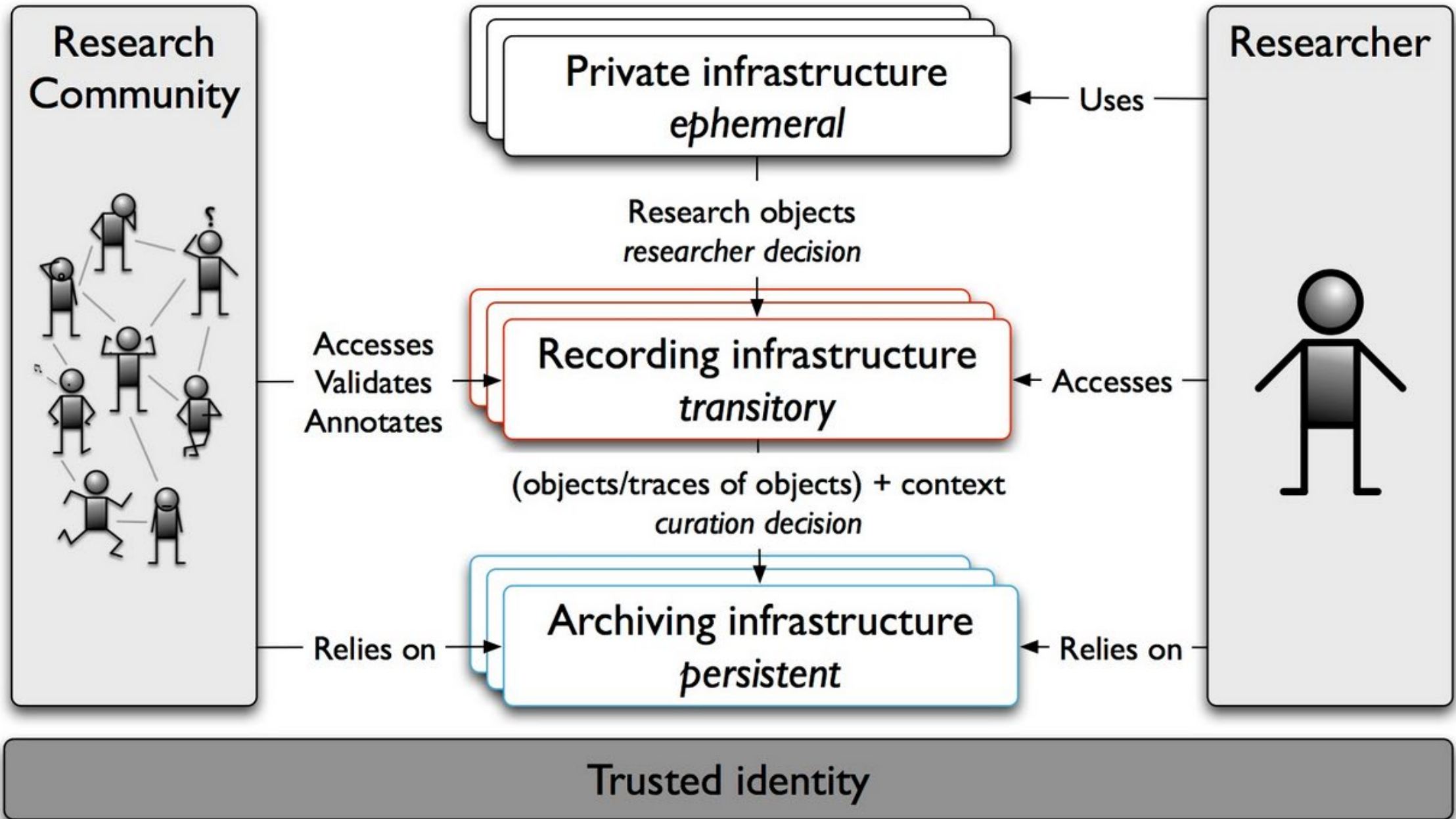
“These domains are separated by “Curation boundaries”, which are virtual decision points at which the creators of data decide what they will share, with whom, with what metadata and under what conditions.”

https://figshare.com/articles/figshare_and_Monash_University_combining_cloud_management_and_discoverability_with_institutional_storage/1224755

Collaboration, Publication, and the Curation Boundary



* R = Register, H = Harvest, O = Obtain

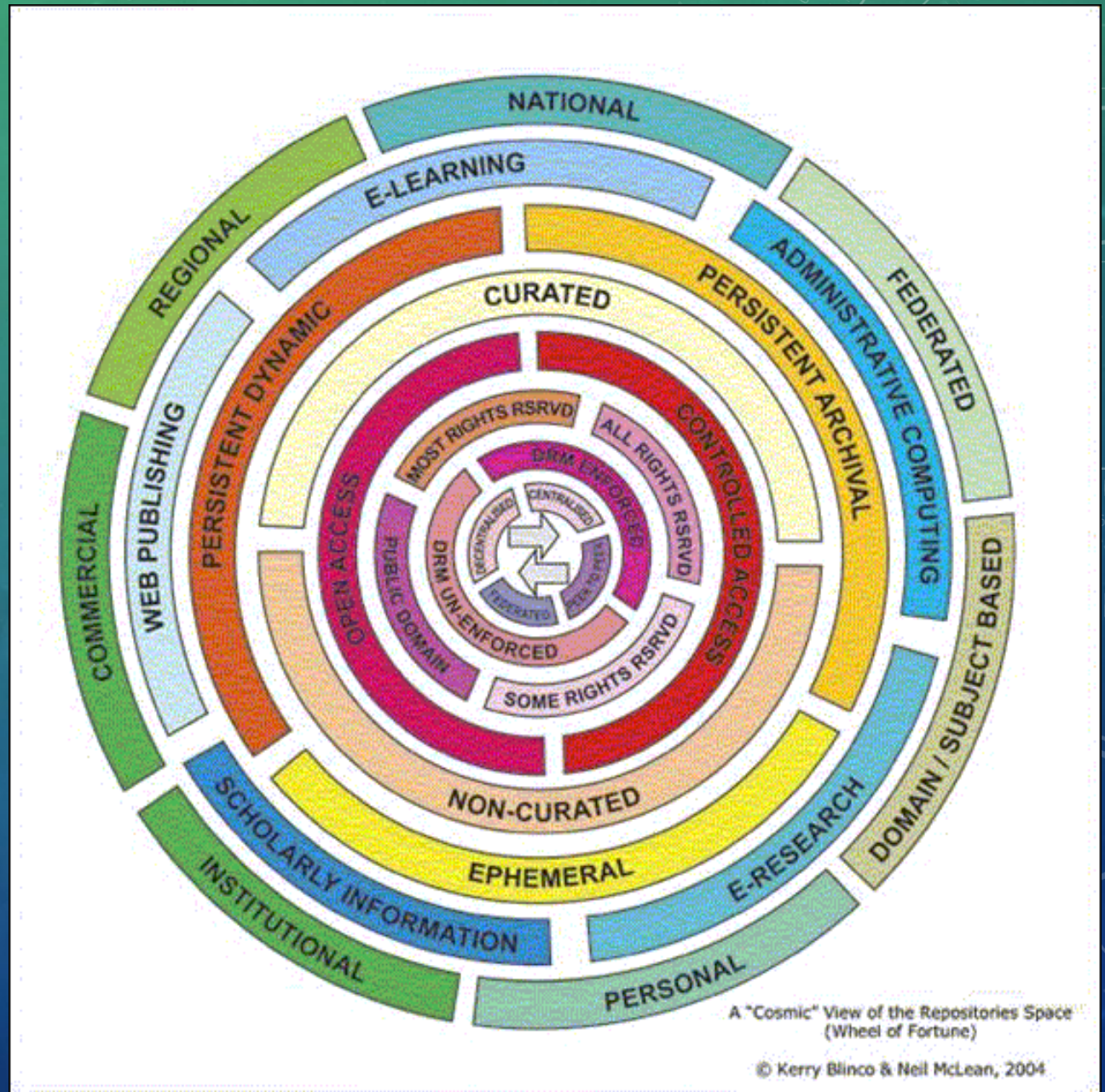


EXERCICE REPOSITORIES



Take the Wheel of Fortune in your hands!

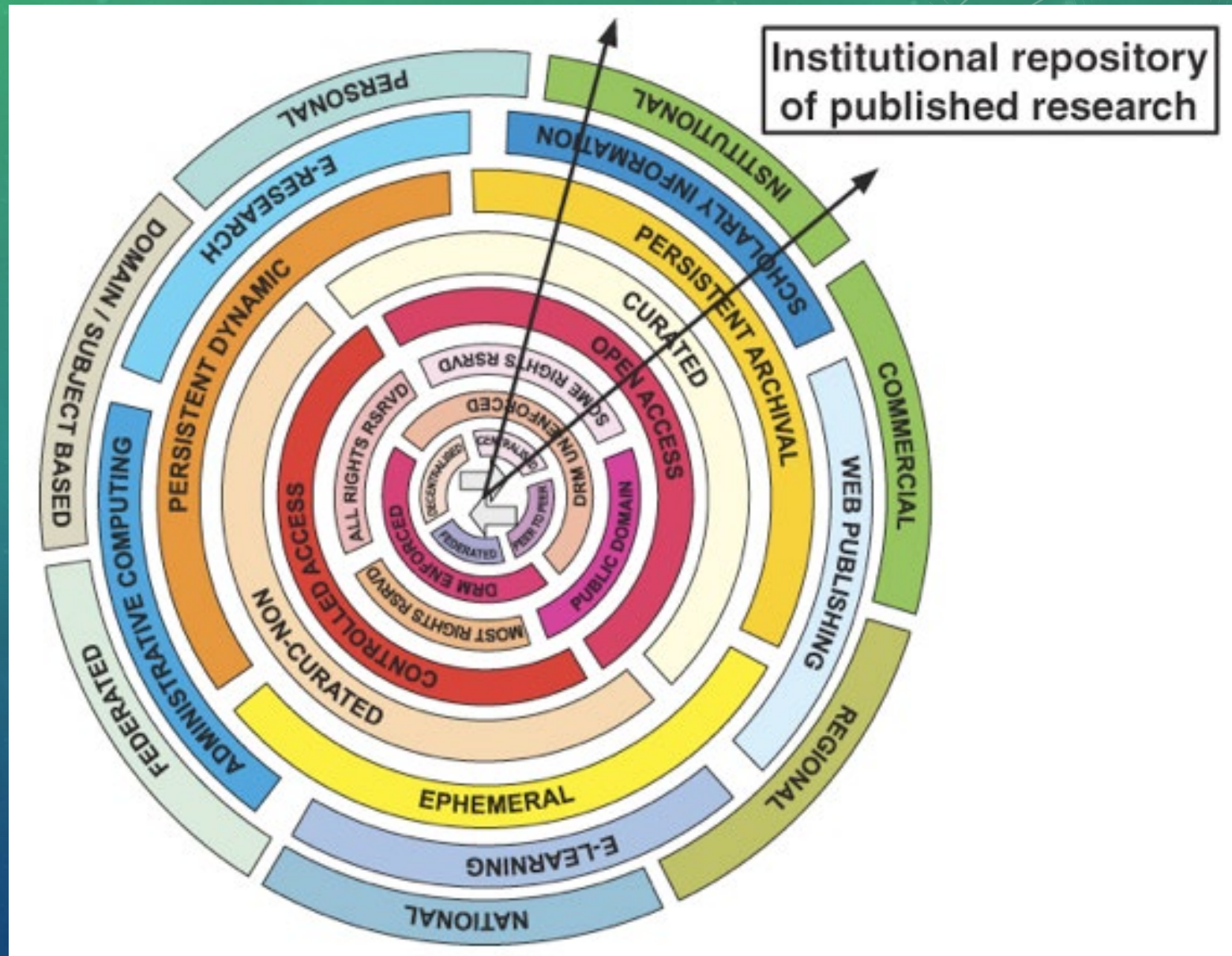
What is the position of each wheel to create an Institutional Open Access Repository for research data?



A "Cosmic" View of the Repositories Space
(Wheel of Fortune)

© Kerry Blinco & Neil McLean, 2004

SOLUTION



CONCLUSIONS



© matiasdelcarmine - AdobeStock

Σκύλλα

Χάρυβδις

life cycle

Δεδομένα

continuum



Σκύλλα

Χάρυβδις

« as long as possible »

Δεδομένα

« as soon as possible »



Σκύλλα

Χάρυβδις

Long Term Archive

Δεδομένα

Open Access Repository



Being between Scylla and Charybdis generally means two things:

1. To choose the lesser between two evils
2. Seeking to choose between equally dangerous extremes is seen as inevitably leading to a disaster.

I.e. a dilemma from which one cannot get out without damage.

An illustration depicting the mythological dilemma of Scylla and Charybdis. In the center, a small wooden boat with two rowers is navigating a narrow channel. To the left, a large, dark, jagged rock formation with a circular opening in the water represents Charybdis, a whirlpool. To the right, a large, dark, multi-headed creature with glowing eyes and long, flowing hair represents Scylla. The background shows a sunset or sunrise over the sea, with a large yellow sun partially obscured by clouds. The overall color palette is dark and moody, with shades of blue, green, and brown.

Χάρυβδις

Σκύλλα

What does this mean for us?

If you go for re-use you will either

- Be attracted by a more open access (repository) approach
- or
- be opting for a long term archival approach.

Find that out first!
Γνῶθι σεαυτόν



None of the approaches is perfect!
Be aware that each one,
the « as long as possible» approach
as well as
« as soon as possible» approach
has advantages and disadvantages.



The only question might be: Is one of them closer to reproducibility or replication?



replication

reproducibility

re-use

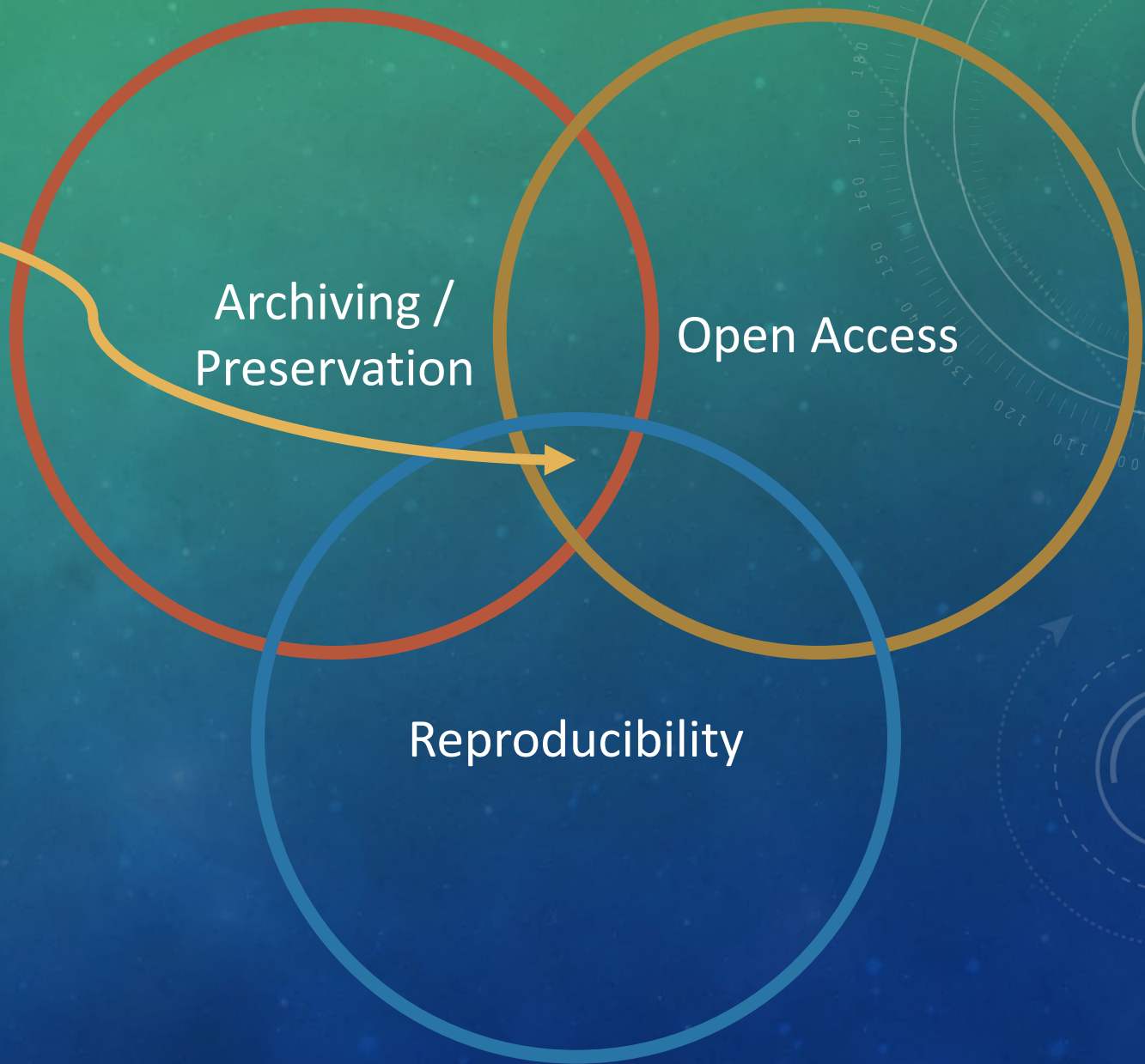
Χάρυβδις

Δεδομένα

Σκύλλα

SWEET SPOT

OF OPEN SCIENCE



Archiving /
Preservation

Open Access

Reproducibility

LITERATURE 1ST PART

Sandra Rendgen: What do we mean by «data»? Idalab. Blog.

<https://idalab.de/blog/data-science/what-do-we-mean-by-data>

Daniel Rosenberg, “Data before the Fact”, in: Gitelman, Lisa (ed.): “Raw Data” is an Oxymoron. Cambridge/Mass.: MIT Press, 2013, p. 33.

LITERATURE 2ND AND 3RD PART

Harvey Ross. Digital Curation. A How-To-Do-It-Manual. Neal Schuman Publishers, 2010.

Treloar, Andrew, David Groenewegen, and Cathrine Harboe-Ree.

"The data curation continuum: Managing data objects in institutional repositories." D-Lib magazine 13.9 (2007): 4.