# Needs and Challenges for Putting FAIR into Practice
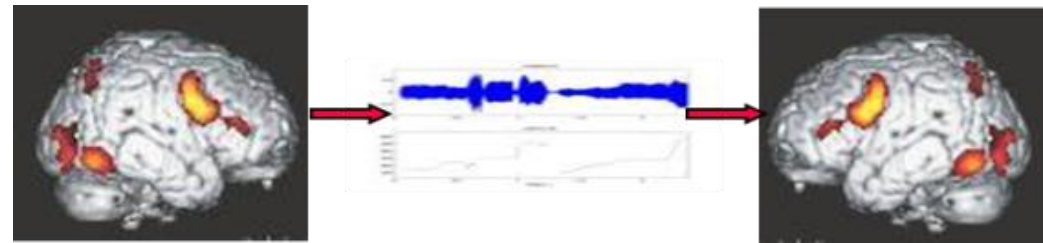
Peter Wittenburg

Max Planck Compute & Data Facility

Research Data Alliance GEDE

MPCDF

GEDE

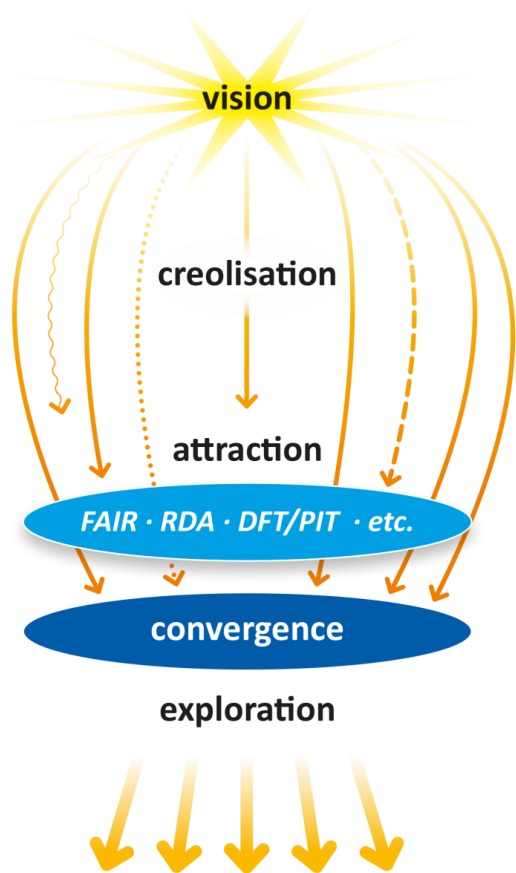Group of European Data Experts

# Who am I – short Intro

- Responsibility for Methodology and Technology at Max Planck Institute for Psycholinguistics
  (what happens in the brain while listening, speaking, acquainting language)

- Responsibility for some large Research Infrastructures (DOBES, CLARIN, EUDAT)

- Co-Founder of Research Data Alliance and co-chairing groups
  (Data Foundation&Technology, Data Fabric, Group of European Data Experts)

- Pushing the Concept of (FAIR) Digital Objects

- Co-Editor of some "relevant" Papers
  (Riding the Wave, FAIR Principles, PID Usage, Turning FAIR into Practice, Revolutionary Infrastructures, etc.)
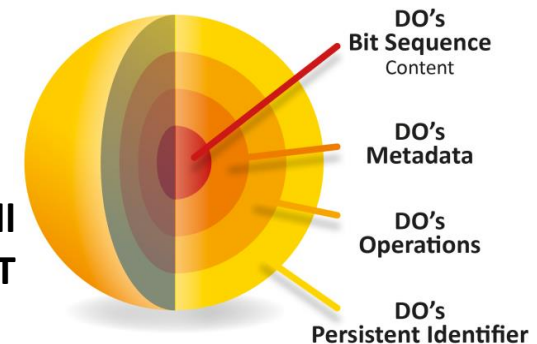
# Dreams



- TCP/IP brought us the world wide unified Computer Network (Internet)
- HTTP brought us the world wide unified Information Network (Web)
- ??? brought us the world wide unified Data Network (???)

- FAIR Principles – great summary of discussions (but paper work)
- RDA with about 10.000 experts working on data issues – grass-roots initiative and yet no systemic approach
- concept of FAIR Digital Objects implements FAIR principles but still no agreement about its usefulness

- but revolutionary inventions take much time (Internet ~ 30 years)

**Taken from Wittenburg & Strawn**
Common Patterns in Revolutionary Infrastructures and Data

**Digital Object Modell
from RDA DFT**



DO's
**Bit Sequence**
Content

DO's
**Metadata**

DO's
**Operations**
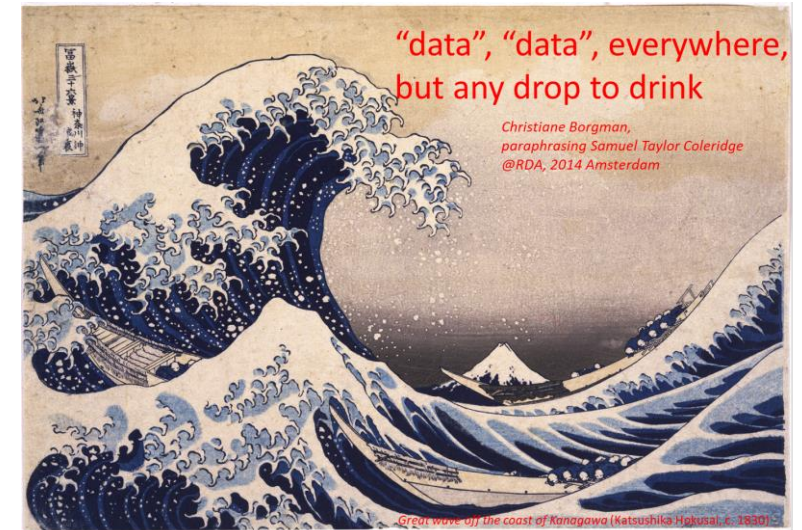
DO's
**Persistent Identifier**

# Reality I

- 80% of time & effort in data projects is spent on wrangling (science, industry, etc.)
- 60% of data projects in industry fail
- many researchers are excluded from data driven science (cross-silo/disciplinary)

- just studied ~60 RI reports deeply – some paradoxes

  - "Standards" are good for science, but researchers don't want to change if no clear benefit.

  - Great FAIR Principles, but researchers shift changes to the end stage of a project.

  - Have huge number of tools, but they don't help to create the unified FAIR domain.

Standards

Researcher

MPCDF
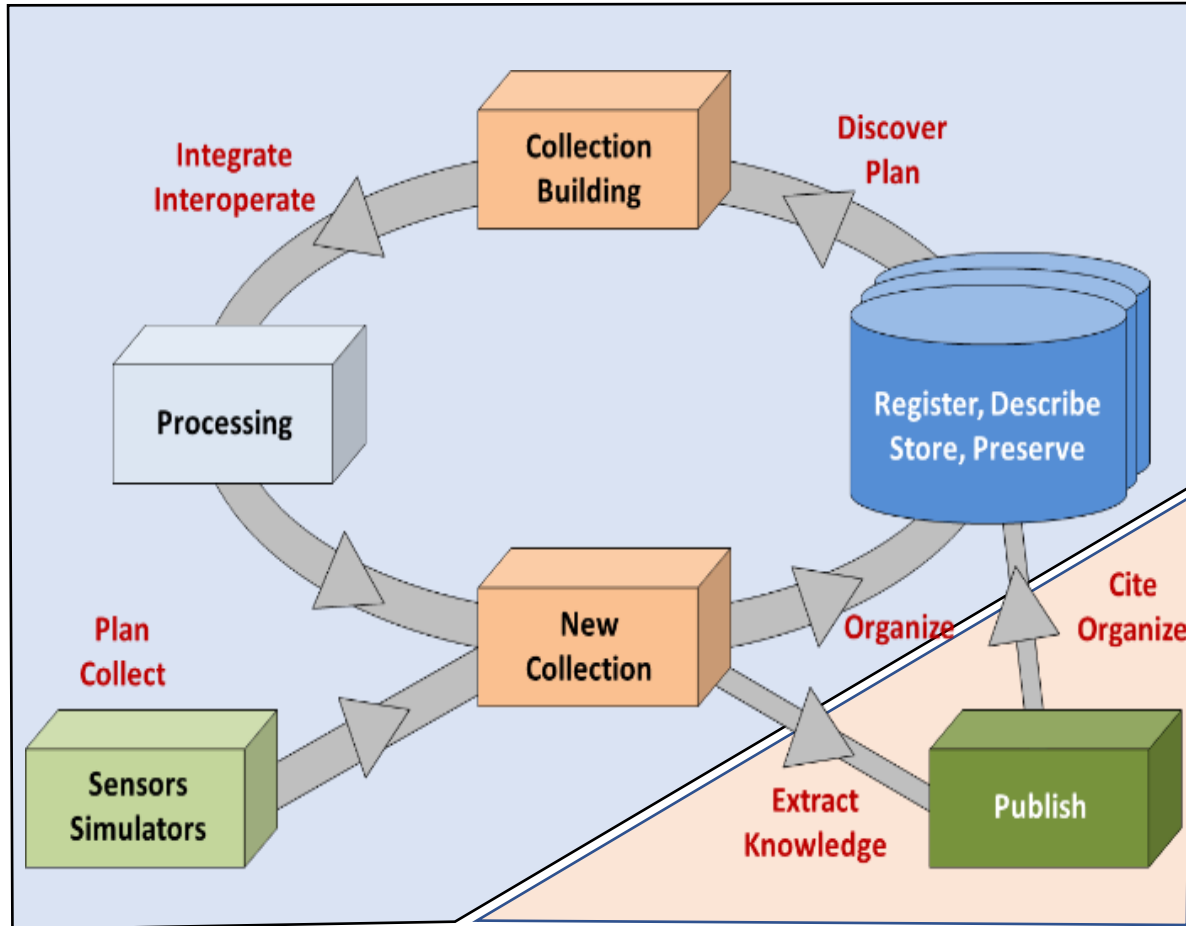
GEDE
Group of European Data Experts

# Reality II

- just studied ~60 RI reports deeply – some paradoxes (ctnd)

  - Having increasing number of regulations (legal, ethical, formal, DMPs), but researchers shift to the start/end stage and hope on copy&paste

  - >90% of data is in the processes and little data will be published, but researchers shift actions to the last step, i.e. Open Science remains a myth – data sharing without metadata?

  - Discipline experts believe that their practices are unique, however, there are re-occurring patterns in data creation, management and processing



"data", "data", everywhere, but any drop to drink

*Christiane Borgman, paraphrasing Samuel Taylor Coleridge @RDA, 2014 Amsterdam*

*Great wave off the coast of Kanagawa (Katsushika Hokusai, c. 1830)*

# Data Cycle Studied in RDA DF

**Data Lab Fabrics**



**Data Publishing**

The results confirmed RDA DF studies in 2014 that led to founding RDA Data Fabric:

- Much has been done to improve the last step: publication (Librarians & Publishers are very active)

- Practices in the Labs did not really change, but there is the mass of data to be re-used

- FAIR Digital Objects as a WayOut to improve practices in the labs !?

# Digital Objects: Model Development I

**some applications**

**FTP**

**SMTP**

**GOPHER**

**etc.**

**early 80s**

Processing / Exchanging
Meaningful Data Entities

**Data Centres**
Management/Curation/
Processing

p

**Data Centres**
Management/Curation/
Processing

**Internet Device**
TCP/IP

p

**Internet Device**
TCP/IP

Message Exchange
without "Meaning"

**complex,
many different types,
scientifically driven**

**„simple", few types,
technologically driven**

MPCDF

GEDE
Group of European Data Experts

# DO: Model Development II

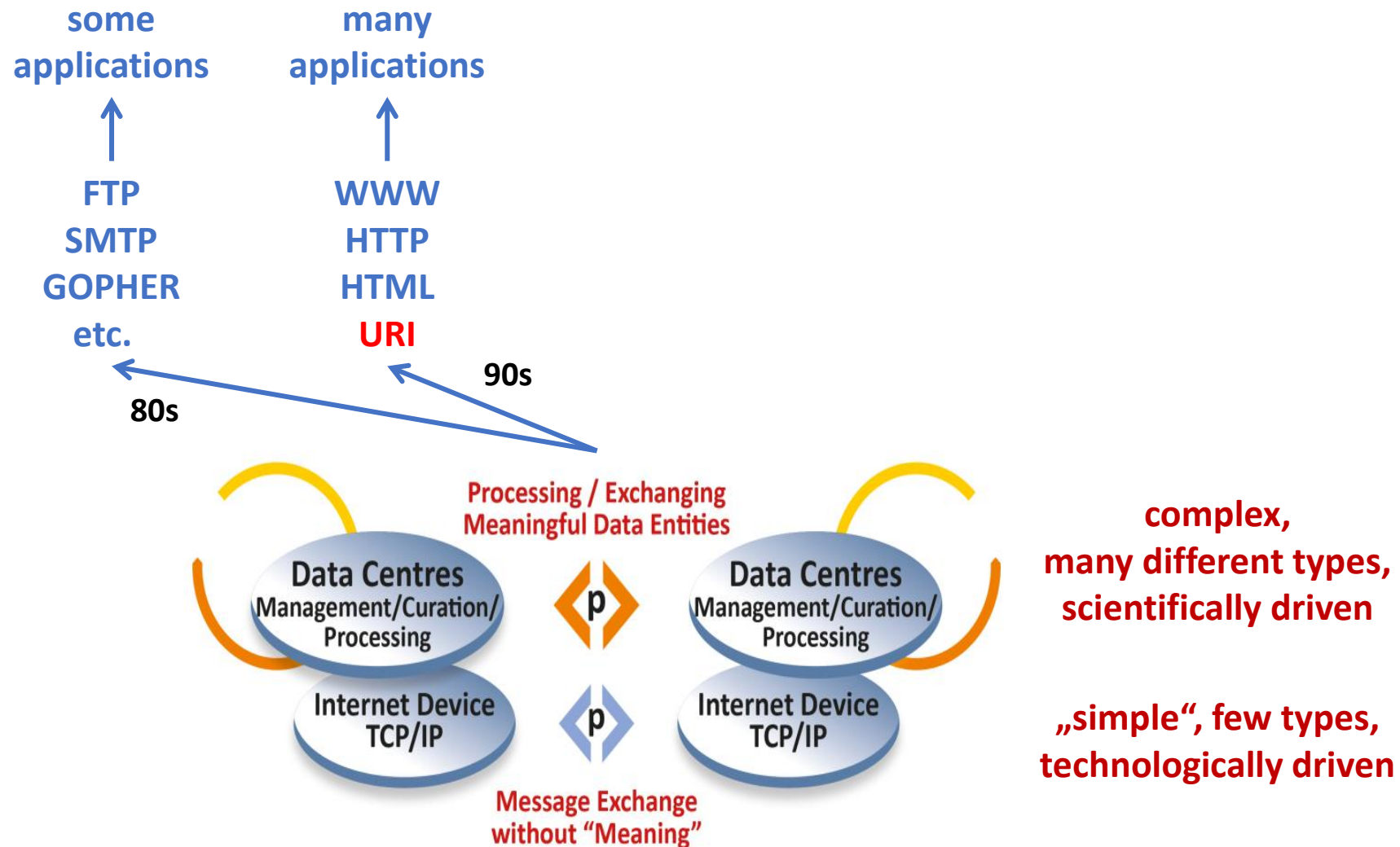some
applications

many
applications

↑

↑

**FTP**

**WWW**

**SMTP**

**HTTP**

**GOPHER**

**HTML**

**etc.**

**URI**

**90s**

**80s**

Processing / Exchanging
Meaningful Data Entities

**complex,
many different types,
scientifically driven**

Data Centres
Management/Curation/
Processing

Data Centres
Management/Curation/
Processing

Internet Device
TCP/IP

Internet Device
TCP/IP

**„simple", few types,
technologically driven**

Message Exchange
without "Meaning"

MPCDF

GEDE
Group of European Data Experts

# DO: Model Development III

some
applications

many
applications

Handle
System

**90s**

FTP
SMTP
GOPHER
etc.

WWW
HTTP
HTML
URI

**90s**

**80s**



Processing / Exchanging
Meaningful Data Entities

Data Centres
Management/Curation/
Processing

Data Centres
Management/Curation/
Processing

Internet Device
TCP/IP

Internet Device
TCP/IP

Message Exchange
without "Meaning"

complex,
many different types,
scientifically driven

„simple", few types,
technologically driven

MPCDF

GEDE
Group of European Data Experts

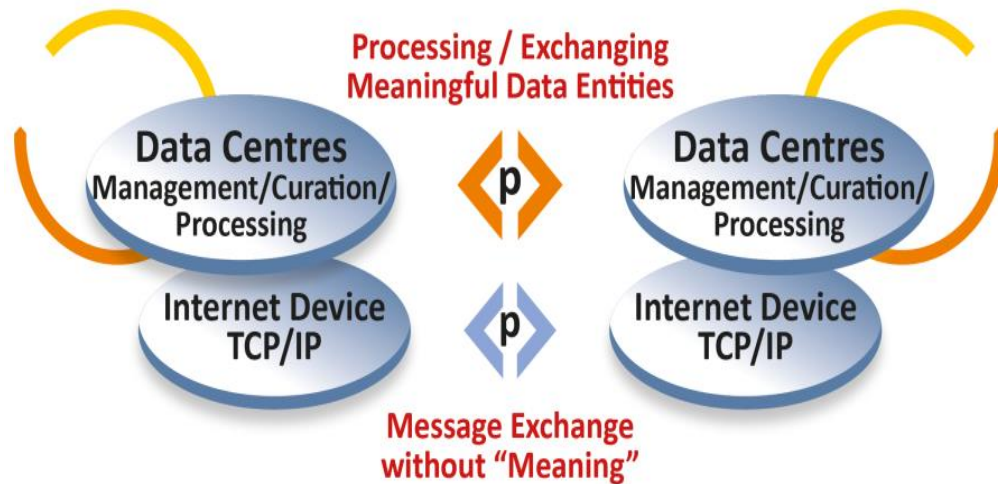# DO: Model Development IV

some
applications

many
applications

**00s**   repositories

**FTP**
**SMTP**
**GOPHER**
etc.

**WWW**
**HTTP**
**HTML**
**URI**

**Handle
System**

**90s**

**00s**   **Publishers
DOI**

**90s**

**80s**



Processing / Exchanging
Meaningful Data Entities

Data Centres
Management/Curation/
Processing

Data Centres
Management/Curation/
Processing

**complex,
many different types,
scientifically driven**

Internet Device
TCP/IP

Internet Device
TCP/IP

**„simple", few types,
technologically driven**

Message Exchange
without "Meaning"

# DO: Model Development V

**00s**

DO Architecture

Handle System

→ repositories

→ Publishers DOI

some applications

many applications

**FTP**
**SMTP**
**GOPHER**
**etc.**

**WWW**
**HTTP**
**HTML**
**URI**

**A Framework for Distributed Digital Object Services (Kahn & Wilensky)**

**80s**

**90s**

**95/06**

Processing / Exchanging Meaningful Data Entities

Data Centres
Management/Curation/Processing

Data Centres
Management/Curation/Processing

Internet Device TCP/IP

Internet Device TCP/IP

Message Exchange without "Meaning"
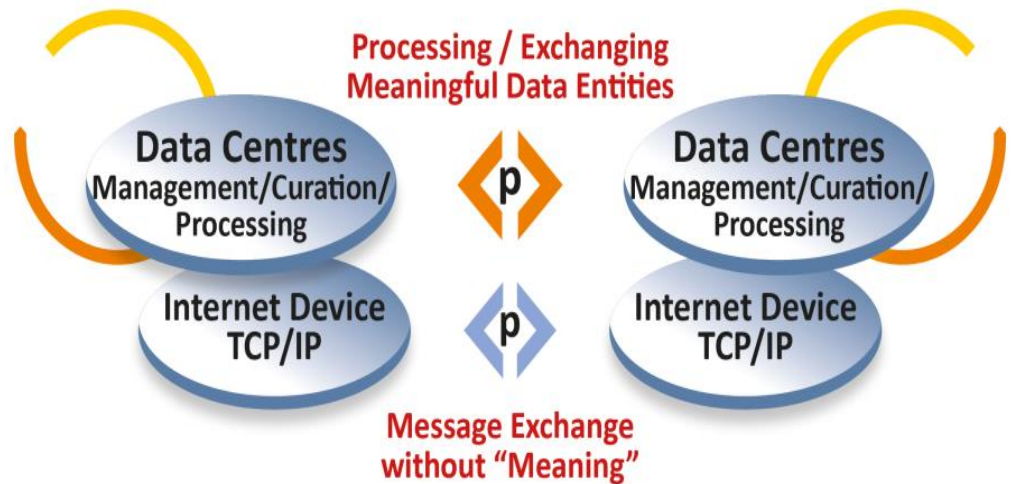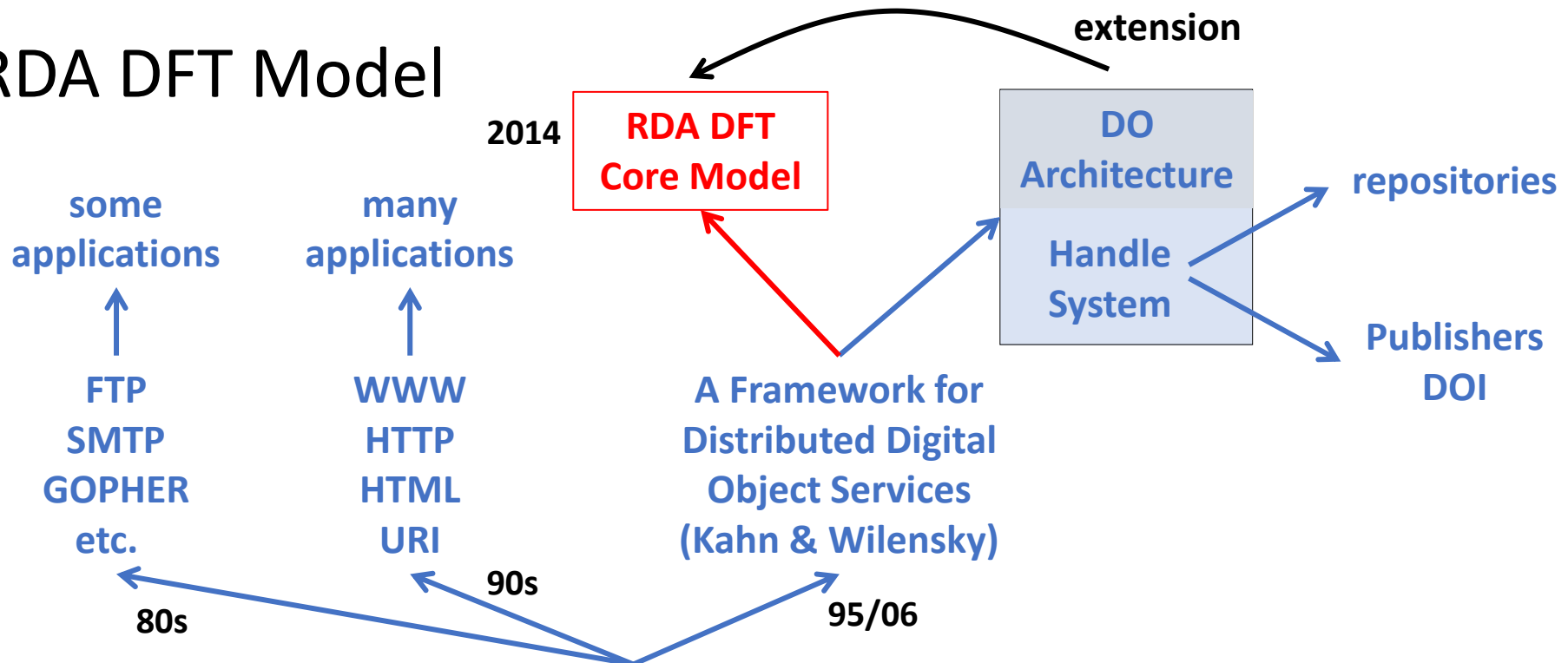
complex,
many different types,
scientifically driven

„simple", few types,
technologically driven

# DO: RDA DFT Model

**extension**

**2014**

**RDA DFT Core Model**

**DO Architecture**

**Handle System**

**repositories**

**Publishers DOI**

**some applications**

**many applications**

**FTP SMTP GOPHER etc.**

**WWW HTTP HTML URI**

**A Framework for Distributed Digital Object Services (Kahn & Wilensky)**

**80s**

**90s**

**95/06**

Processing / Exchanging Meaningful Data Entities

Data Centres
Management/Curation/Processing

Data Centres
Management/Curation/Processing

Internet Device TCP/IP
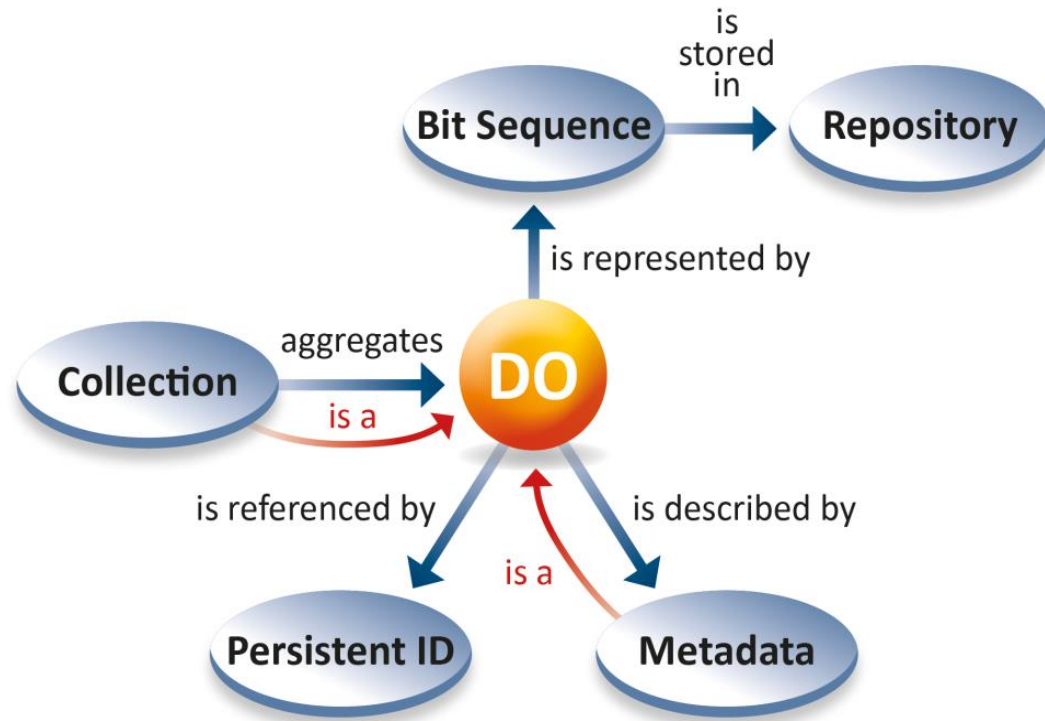
Internet Device TCP/IP

Message Exchange without "Meaning"

**complex, many different types, scientifically driven**

**„simple", few types, technologically driven**

MPCDF

# DO: RDA Data Foundation & Terminology (2014)



**RDA Specs**
**-PIT**
**-Kernel**
**-DTR**

**RDA DFT**: a DO has a structured bit sequence stored in some repositories, is assigned a PID and is described by metadata.
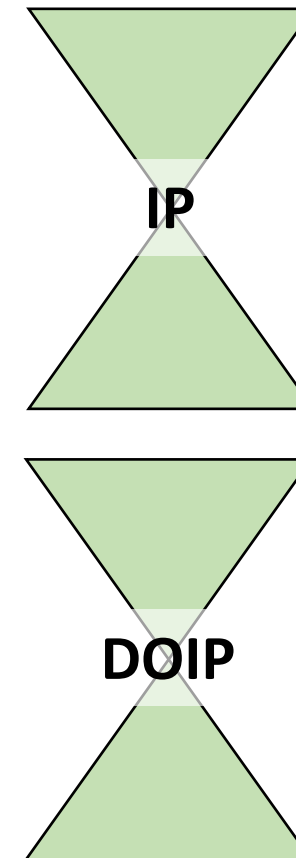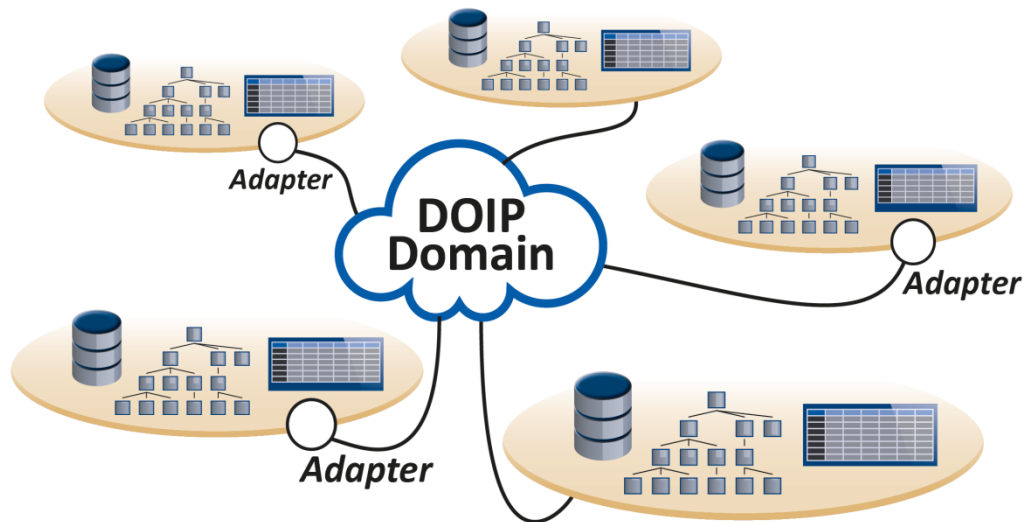
DOs can be aggregated to collections which are also DO. Metadata descriptions are DOs.

DO's PID Record is resolved to machine-actionable attributes enabling human/machine actions.

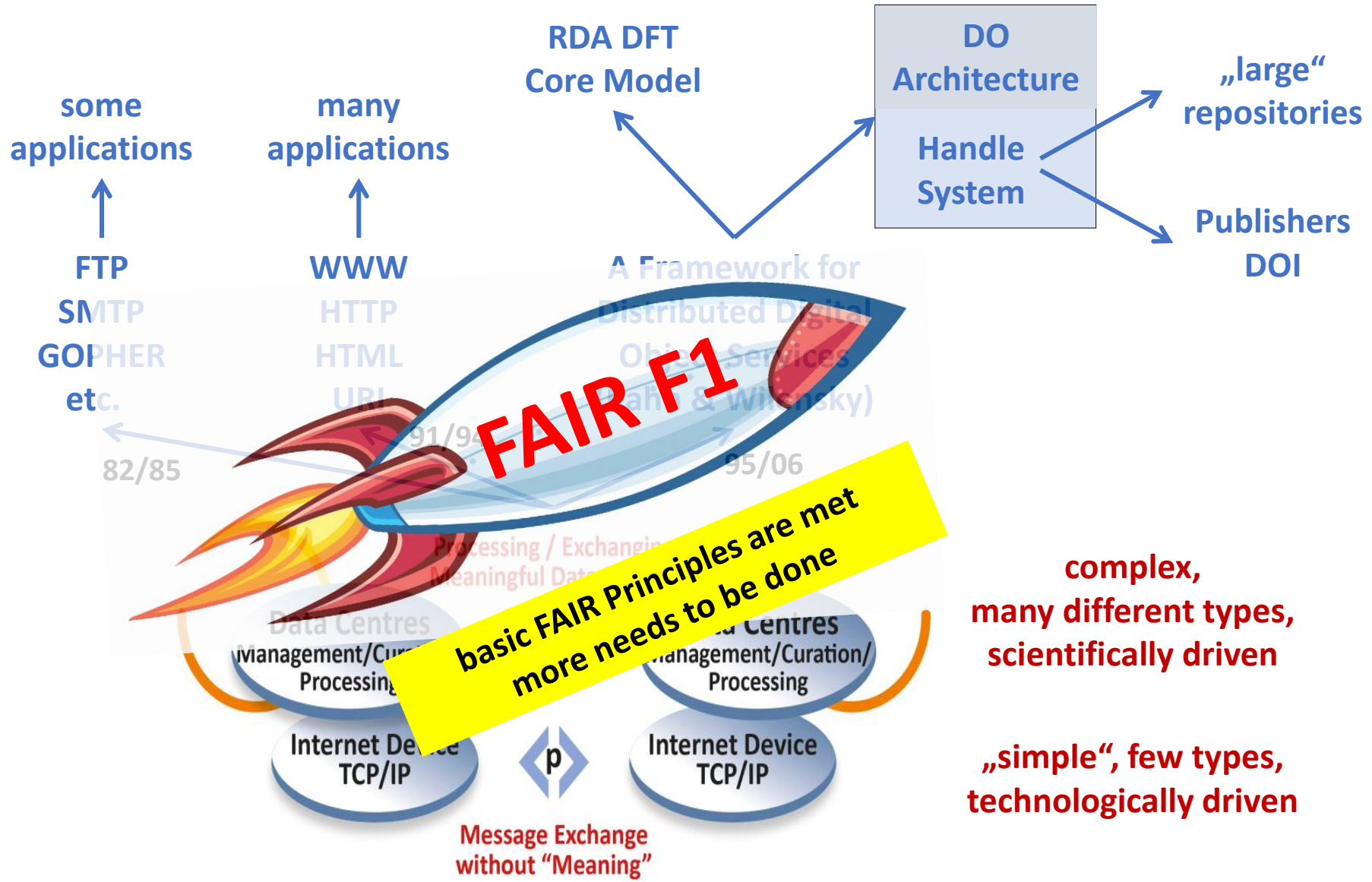PID = globally unique persistent resolvable identifier (Handle, DOI)

# also Software available: DOIP V2.0 (DONA)

- improved specification and implementation of DO Architecture
- DOIP V2.0 specifying unified client – DO Server interaction
  - CORDRA reference implementation ready
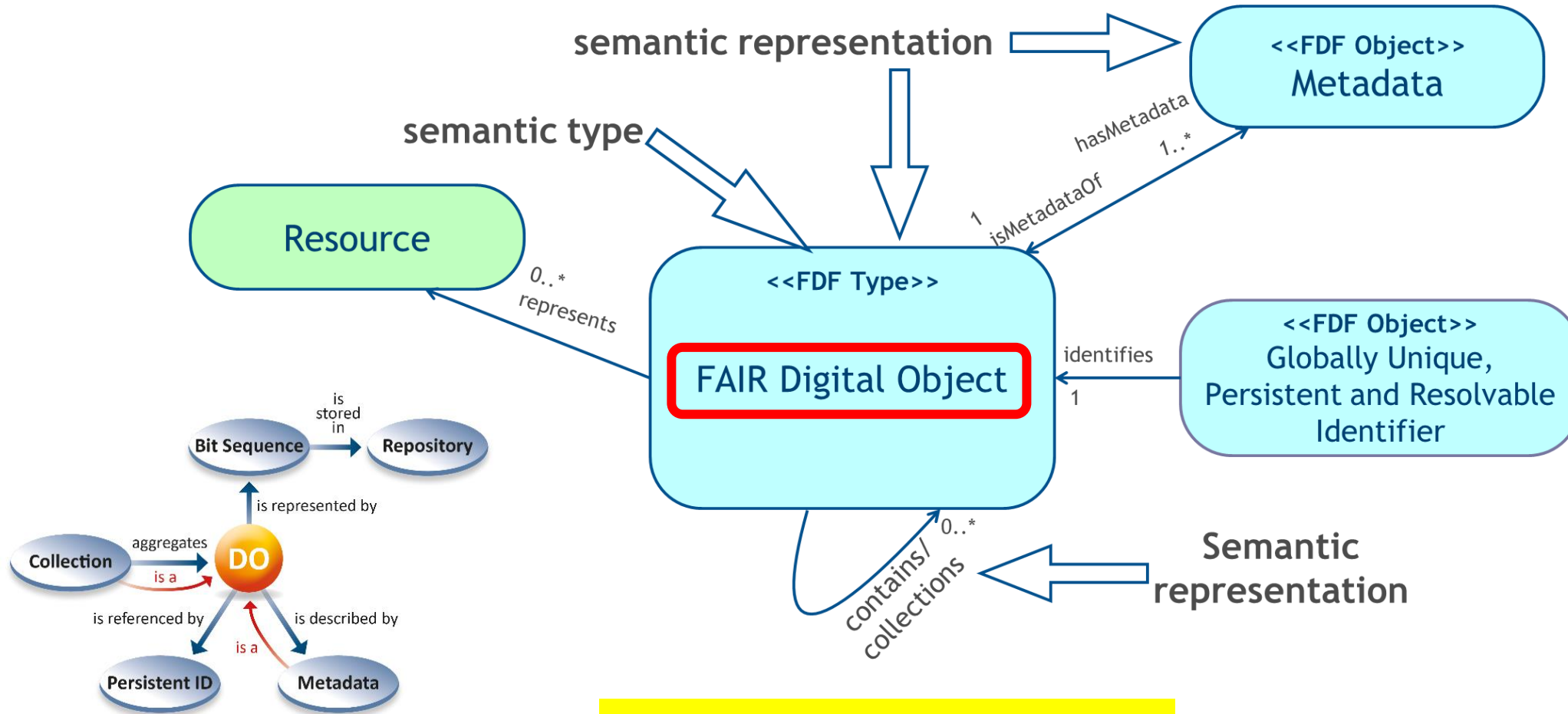  - DOIPV2.0 SDK ready
  - all based on PIDs



**IP**

**DOIP**

**Cannot expect people to start from Scratch**

# Do DOs support FAIR?

RDA DFT
Core Model

**DO
Architecture**

**Handle
System**

„large"
repositories

some
applications

many
applications

Publishers
DOI

FTP
SMTP
GOPHER
etc.

WWW
HTTP
HTML
URL

A Framework for
Distributed Digital
Object Services
(Kahn & Wilensky)

82/85

91/94

**FAIR F1**

95/06

Processing / Exchanging
Meaningful Data

Data Centres
Management/Curation/
Processing

Data Centres
Management/Curation/
Processing

Internet Device
TCP/IP

Internet Device
TCP/IP

**basic FAIR Principles are met
more needs to be done**

**complex,
many different types,
scientifically driven**

**„simple", few types,
technologically driven**

Message Exchange
without "Meaning"

# FAIR requires Semantic Explicitness

(in close collaboration with Luiz Bonino, applying mechanisms from LD)

# Long Term Vision & Identification (FAIR F1)

- V. Cerf: warning for a dark digital age
- why?
  - it's about persistence of relevant bit-sequences, describing metadata AND **relations** for **100+** years
  - and relations will express much of our cumulative scientific knowledge



ISBN → Book MD → Catalogue

Handle/DOIs

DO's Bit Sequence Content
DO's Metadata
DO's Operations
DO's Persistent Identifier

URLs — **location**

DOs identified by Handles abstract from details, bind relevant information and encapsulate!

# All Ready for a Big Change?



- **NOoooo₀₀₀₀**

- FDO not yet accepted broadly – many different voices how to build a global unified data infrastructure (yet no help from EOSC)

- Researchers are right to be careful:
  - no stability yet – still much dynamics in convictions, trends
  - miss supporting software to reduce the load for researchers

- Thus, if we want to change practices
  - need to take the researchers with us who are not interested in technicalities
    - offer the obvious (Zenodo, B2Share, Handle/DOI, etc.)
    - address data sovereignty
    - need to be patient, nevertheless work hard on DO SW components
  - Interested? – Join the GEDE DO and CWFS discussions (Canonical Workflow Frameworks for Science)

# Thanks for the attention.

all can be found under GEDE – Github: https://github.com/GEDE-RDA-Europe/GEDE
just search for „Github GEDE"