

FAIR

FAIR-prinsippene

Fra forskningsmiljø til samfunnsinfrastruktur

Per Hovde



Agenda



1

Bakgrunnen for FAIR-prinsippene

2

Hvorfor er FAIR viktig – også utenfor forskningsverdenen?

3

FAIR-prinsippene i detalj

4

Persistente identifikatorer (PID) sentralt i FAIR

5

Behov for nye tjenester for tilordning av PIDs?

Agenda



1

Bakgrunnen for FAIR-prinsippene

2

Hvorfor er FAIR viktig – også utenfor forskningsverdenen?

3

FAIR-prinsippene i detalj

4

Persistente identifikatorer (PID) sentralt i FAIR

5

Behov for nye tjenester for tilordning av PIDs?

www.nature.com/scientificdata

SCIENTIFIC DATA

Amended: Addendum

OPEN Comment: The FAIR Guiding Principles for scientific data management and stewardship

SUBJECT CATEGORIES
» Research data
» Publication characteristics

Mark D. Wilkinson *et al.**

Received: 10 December 2015
Accepted: 12 February 2016
Published: 15 March 2016

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measurable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplar implementations in the community.

Supporting discovery through good data management
Good data management is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process. Unfortunately, the existing digital ecosystem surrounding scholarly data publication prevents us from extracting maximum benefit from our research investments (e.g., ref. 1). Partially in response to this, science funders, publishers and governmental agencies are beginning to require data management and stewardship plans for data generated in publicly funded experiments. Beyond proper collection, annotation, and archival, data stewardship includes the notion of 'long-term care' of valuable digital assets, with the goal that they should be discovered and re-used for downstream investigations, either alone, or in combination with newly generated data. The outcomes from good data management and stewardship, therefore, are high quality digital publications that facilitate and simplify this ongoing process of discovery, evaluation, and reuse in downstream studies. What constitutes 'good data management' is, however, largely undefined, and is generally left as a decision for the data or repository owner. Therefore, bringing some clarity around the goals and desiderata of good data management and stewardship, and defining simple guideposts to inform those who publish and/or preserve scholarly data, would be of great utility.

This article describes four foundational principles—Findability, Accessibility, Interoperability, and Reusability—that serve to guide data producers and publishers as they navigate around these obstacles, thereby helping to maximize the added-value gained by contemporary, formal scholarly digital publishing. Importantly, it is our intent that the principles apply not only to 'data' in the conventional sense, but also to the algorithms, tools, and workflows that led to that data. All scholarly digital research objects²—from data to analytical pipelines—benefit from application of these principles, since all components of the research process must be available to ensure transparency, reproducibility, and reusability.

There are numerous and diverse stakeholders who stand to benefit from overcoming these obstacles: researchers wanting to share, get credit, and reuse each other's data and interpretations; professional data publishers offering their services; software and tool-builders providing data analysis and processing services such as reusable workflows; funding agencies (private and public) increasingly

*Correspondence and requests for materials should be addressed to B.M. (email: baend.mors@dtf.isn.li).
#A full list of authors and their affiliations appears at the end of the paper.

SCIENTIFIC DATA | 3:160018 | DOI: 10.1038/sdata.2016.18 | 1

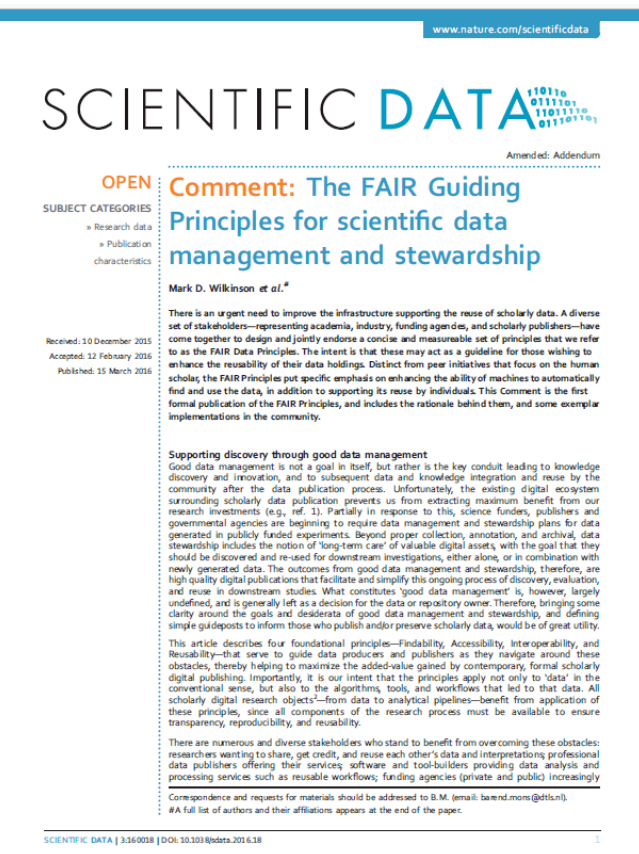
Opprinnelsen til FAIR:

Artikkel i Nature 15. mars 2016 (Wilkinson *et al.*, 2016)

Artikkelen pekte på:

- Datamengden i forskning økte eksponensielt – men gjenbruk var krevende
- Å samle data til én analyse kunne kreve månedsvis med manuelt arbeid
- En ny utfordring: maskiner og automatiserte agenter møtte de samme barrierene som mennesker – men uten evnen til å tolke kontekst og mening uten forklarende metadata
- Akademia, industri, forskningsinansierer og forlag sammen om å definere et felles minimumssett av prinsipper på tvers av systemer og fagfelt dom de kalte :

- **FAIR (Findable, Accessible, Interoperable, Reusable)**



Oppsummering av artikkelen:

- FAIR retter seg mot **både mennesker og maskiner**
- Introduserte begrepet **Machine-actionability**
- FAIR-prinsippene destillert:
 - **Findable:** Persistente identifikatorer, rike metadata og registrering i søkbare ressurser
 - **Accessible:** Standardiserte, åpne protokoller for tilgang, med støtte for autentisering der nødvendig, og varig tilgang til metadata
 - **Interoperable:** Bruk av delte formelle språk, standarder og vokabularer, samt eksplisitte referanser mellom data
 - **Reusable:** Tydelige lisenser, dokumentert proveniens, rike beskrivelser og samsvar med relevante fellesskapsstandarder
- Prinsippene er **teknologi-nøytrale og modulære**
- **FAIR en grunnleggende forutsetning for god dataforvaltning i et distribuert forskningslandskap, og en nødvendig respons på desentraliseringen av datalagring og -publisering**

Ferskt eksempel på utfordringer med data som ikke er FAIR



KUNSTIG INTELLIGENS

Hollywood-kjendiser i falskt datasett: — Komisk dårlig

Svært tvilsomme data har blitt brukt for å utvikle KI-modeller som skal oppdage hjerneslag. Det kan få alvorlige konsekvenser for pasienter, advarer forskere.



Disse bildene av kjente skuespillere lå i et billedatasett som ble brukt til å trene et KI-verktøy som skulle oppdage hjerneslag. Øverst fra venstre Sylvester Stallone, Angelina Jolie og Clint Eastwood. Nederst fra venstre: Pierce Brosnan, George Clooney og Forest Whitaker. Foto: Hentet fra datasettet Droopy på nettsiden Kaggle

Agenda



1

Bakgrunnen for FAIR-prinsippene

2

Hvorfor er FAIR viktig – også utenfor forskningsverdenen?

3

FAIR-prinsippene i detalj

4

Persistente identifikatorer (PID) sentralt i FAIR

5

Behov for nye tjenester for tilordning av PIDs?

Datavekst og desentralisering

Verdens datamengde vokser eksponentielt – og spres på stadig flere infrastrukturer (desentralisert dataforvaltning)

220 ZB

Prognostisert mengde innen utgangen av 2026

1000+

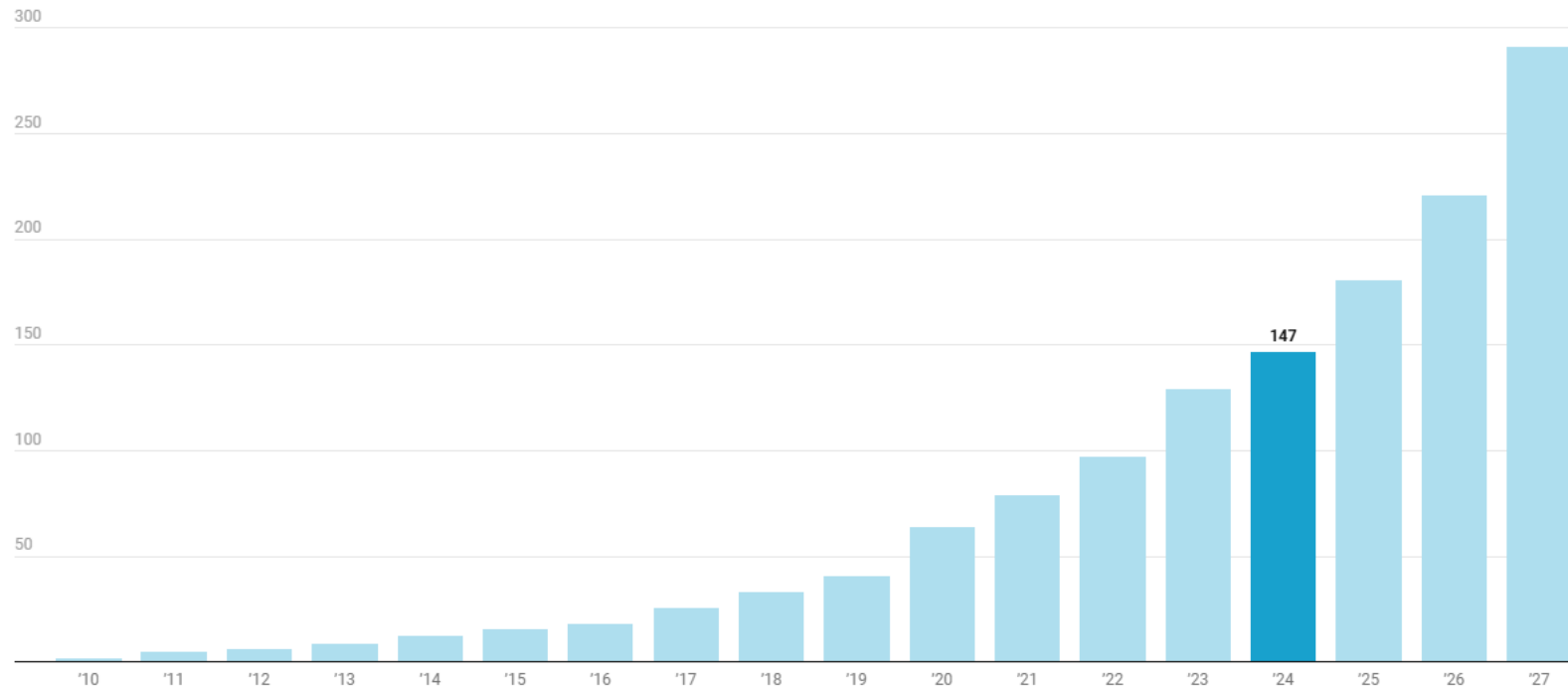
separate dataregistre og infrastrukturer i norsk offentlig sektor

<10%

av tilgjengelige data faktisk gjenbrukt i nye analyser

Worldwide data creation 2010-2027 (in zettabytes)

Volume of information created and consumed globally to reach 147 zettabytes (ZB) in 2024 forecast by IDC with continued growth through 2027.



Based on forecast and compound annual growth rate identified by source.

Chart: CBL Data Recovery/Platter Chatter • Source: IDC DataAge 2025; IDC Worldwide Global DataSphere Forecast, 2023–2027 • Created with [Datawrapper](#)

Datavekst og desentralisering

Verdens datamengde vokser eksponentielt – og spres på stadig flere infrastrukturer (desentralisert dataforvaltning)

220 ZB

Prognostisert mengde innen utgangen av 2026

1000+

separate dataregistre og infrastrukturer i norsk offentlig sektor

<10%

av tilgjengelige data faktisk gjenbrukt i nye analyser

1 Zettabyte (ZB) = 2^{70} =

1 180 591 620 717 411 303 424 bytes

(1 trilliard, 180 trillioner, 591 billiarder, 620 billioner, 717 milliarder, 411 millioner, 303 tusen og 424 bytes)

Datavekst og desentralisering

Verdens datamengde vokser eksponentielt – og spres på stadig flere infrastrukturer (desentralisert dataforvaltning)

220 ZB

Prognostisert mengde innen utgangen av 2026

1000+

separate dataregistre og infrastrukturer i norsk offentlig sektor

<10%

av tilgjengelige data faktisk gjenbrukt i nye analyser

Utfordringene i dag

- Data finnes overalt, men er svært vanskelig å oppdage
- Metadata er fragmentert, inkonsistente eller mangler helt
- Ulike formater og protokoller hindrer sammenstilling
- Tilgangsbetingelser er uklare eller udokumenterte
- Analyse krever uker/måneder av manuelt arbeid
- Maskinlesbare data er unntaket, heller enn regelen

Manglende etterlevelse av FAIR gjør at samfunnet milliarder i tapt innovasjon

FAIR er viktig langt ut over forskningsmiljøene

Prinsippene er relevante for alle sektorer i en datadrevet økonomi (samfunnsinfrastruktur)

Opprinnelse: Behov i forskning

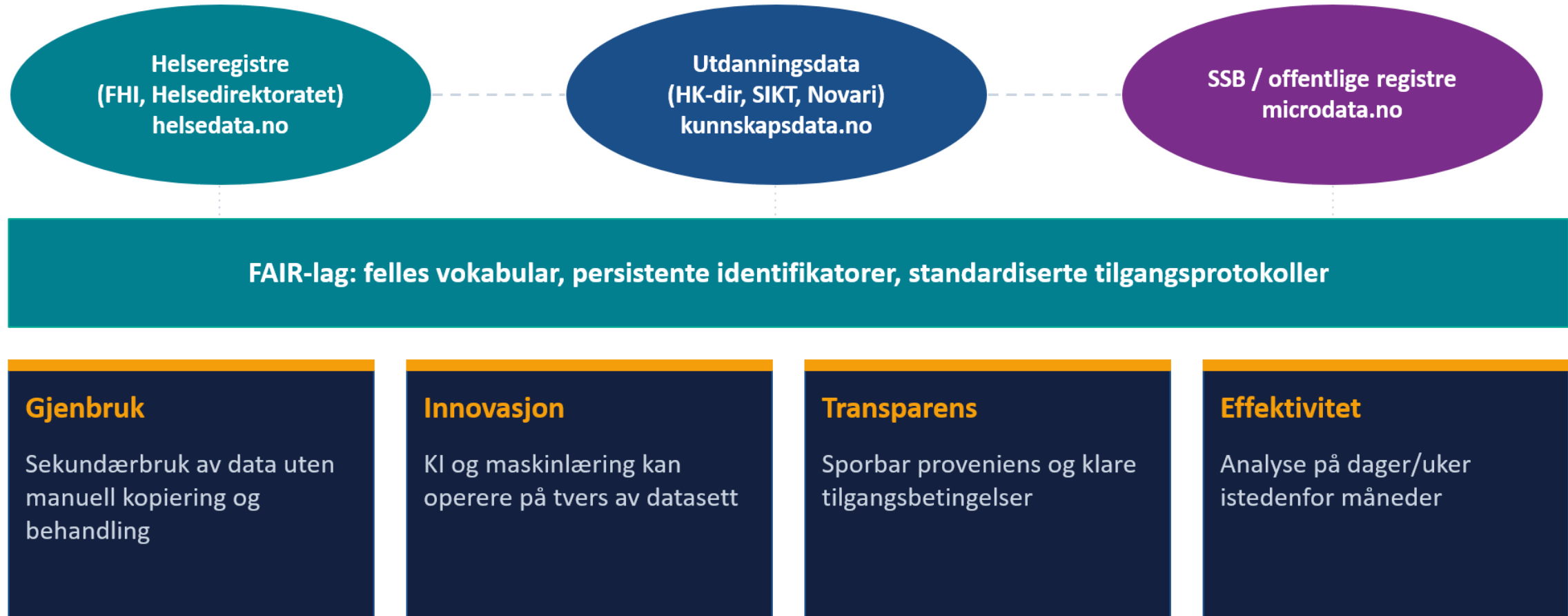
- ✓ Reproduserbarhet av vitenskapelige resultater
- ✓ Etterprøvbarehet av forskningsresultater
- ✓ Deling og gjenbruk av datasett mellom forskningsmiljøer
- ✓ Automatisert dataintegrasjon i eVitenskap
- ✓ Sitering og kreditering av datapublikasjoner



Bredt potensial: Alle datasektorer

- Offentlig sektor: deling på tvers av registre og etater
- Helse: sammenstilling av pasientdata, helseregistre og sosio-økonomiske data (SSB)
- Utdanning: kobling av utdannings- og arbeidsmarkedsdata
- Næringsliv: gjenbruk og videresalg av data som tjeneste for å skape innovasjon
- KI og automatisering: maskinlesbare data som grunnlag
- Økende grad av desentralisering av data
- Regelverk som "tvinger" oss til å dele mer data

FAIR på variabelnivå – på tvers av desentraliserte datainfrastrukturer



FAIR data er grunnmuren for kunstig intelligens

KI er bare så god som dataene den opererer på

«Machine-actionable»

FAIR-artikkelen fra 2016 hadde maskinlesbarhet og autonome agenter som et kjernemål – lenge før dagens KI-bølge. Det prinsippet er nå mer aktuelt enn noen gang.

FAIR-data gir pålitelig KI

KI som opererer på FAIR-strukturerte data gir mer pålitelige resultater fordi datagrunnlaget er entydig beskrevet, sporbart og dokumentert.

Forklarbarhet og tillit

Krav til sporbar proveniens (R1.2) gjør det mulig å forklare hvilke data en KI-modell bygger på – kritisk for regulert bruk.

Lisens og samtykke

FAIR krever tydelig lisens (R1.1) – som er avgjørende for å avklare om data lovlig kan brukes til KI-trening.

Automatisk dataoppdagelse

KI-agenter kan selv finne og vurdere relevante datasett når metadata er maskinlesbar og standardisert (F + I).

Tverrdomene-analyse

Interoperable data (I) gjør at KI kan analysere på tvers av helse, utdanning og sosio-økonomi uten manuell datavask.

Vektorer og RAG

Retrieval-Augmented Generation (RAG) er helt avhengig av at grounding data er søkbar, identifiserbar og velstrukturert (F + A).

Uten FAIR data: KI som svarer feil, ikke kan forklares, og ikke etterprøves | Med FAIR data: KI mer pålitelig, transparent og skalerbar

Agenda



1

Bakgrunnen for FAIR-prinsippene

2

Hvorfor er FAIR viktig – også utenfor forskningsverdenen?

3

FAIR-prinsippene i detalj

4

Persistente identifikatorer (PID) sentralt i FAIR

5

Behov for nye tjenester for tilordning av PIDs?

FAIR - Fire prinsipper for god dataforvaltning og -publisering

F

Findable

- F1.** (meta)data are assigned a *globally unique and persistent identifier*
- F2.** data are described with rich *metadata* (defined by R1 below)
- F3.** metadata clearly and *explicitly include the identifier of the data they describe*
- F4.** (meta)data are registered or *indexed in a searchable resource*

FAIR - Fire prinsipper for god dataforvaltning og -publisering

F

Findable

- F1.** (meta)data are assigned a *globally unique and persistent identifier*
- F2.** data are described with rich *metadata* (defined by R1 below)
- F3.** metadata clearly and *explicitly include the identifier of the data they describe*
- F4.** (meta)data are registered or *indexed in a searchable resource*

A

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 authorisation procedure where necessary*
- A2.** *Metadata are accessible*, even when the data are no longer available.

FAIR - Fire prinsipper for god dataforvaltning og -publisering

F

Findable

- F1.** (meta)data are assigned a *globally unique and persistent identifier*
- F2.** data are described with rich *metadata* (defined by R1 below)
- F3.** metadata clearly and *explicitly include the identifier of the data they describe*
- F4.** (meta)data are registered or *indexed in a searchable resource*

A

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 authorisation procedure where necessary*
- A2.** *Metadata are accessible*, even when the data are no longer available.

I

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*

FAIR - Fire prinsipper for god dataforvaltning og -publisering

F

Findable

- F1.** (meta)data are assigned a *globally unique and persistent identifier*
- F2.** data are described with rich *metadata* (defined by R1 below)
- F3.** metadata clearly and *explicitly include the identifier of the data they describe*
- F4.** (meta)data are registered or *indexed in a searchable resource*

A

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 authorisation procedure where necessary*
- A2.** *Metadata are accessible*, even when the data are no longer available.

I

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*

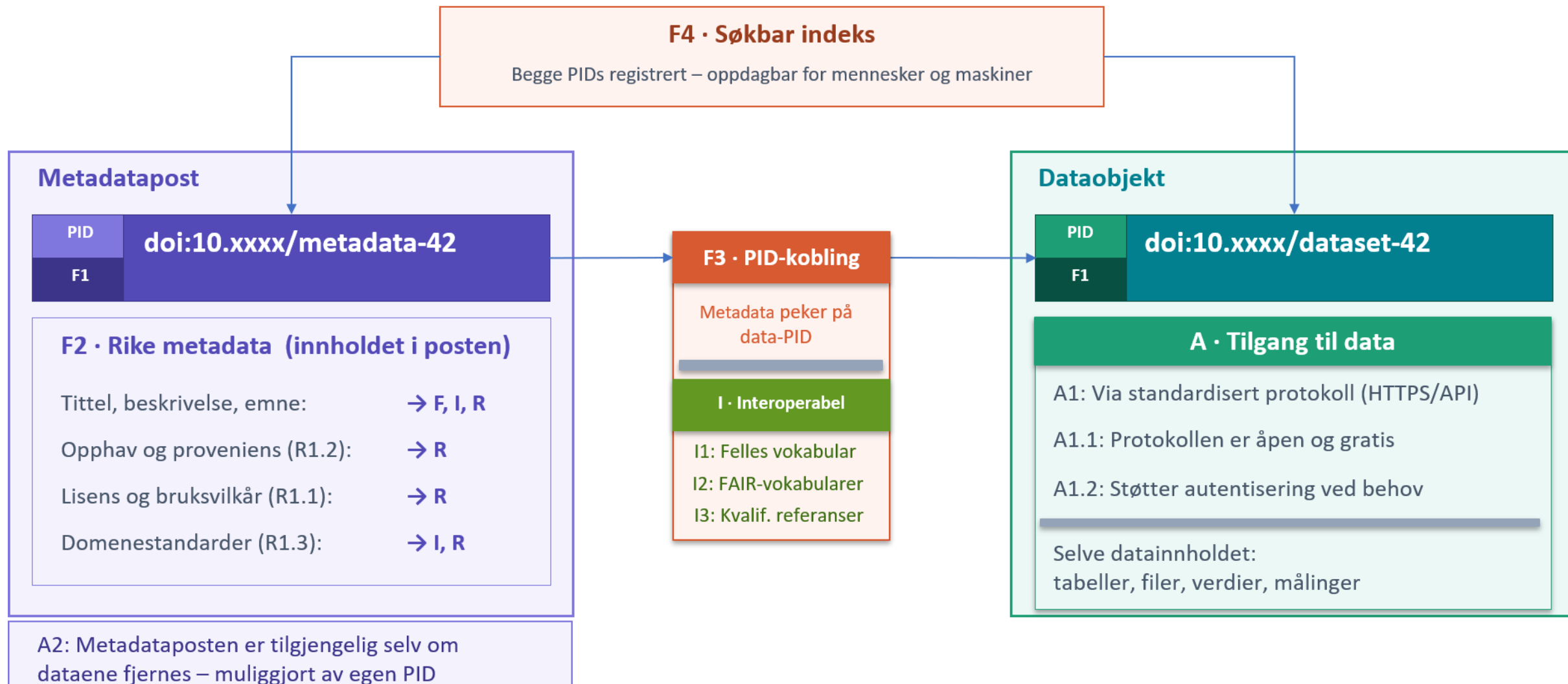
R

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*

Sammenhengen mellom FAIR-dimensjonene

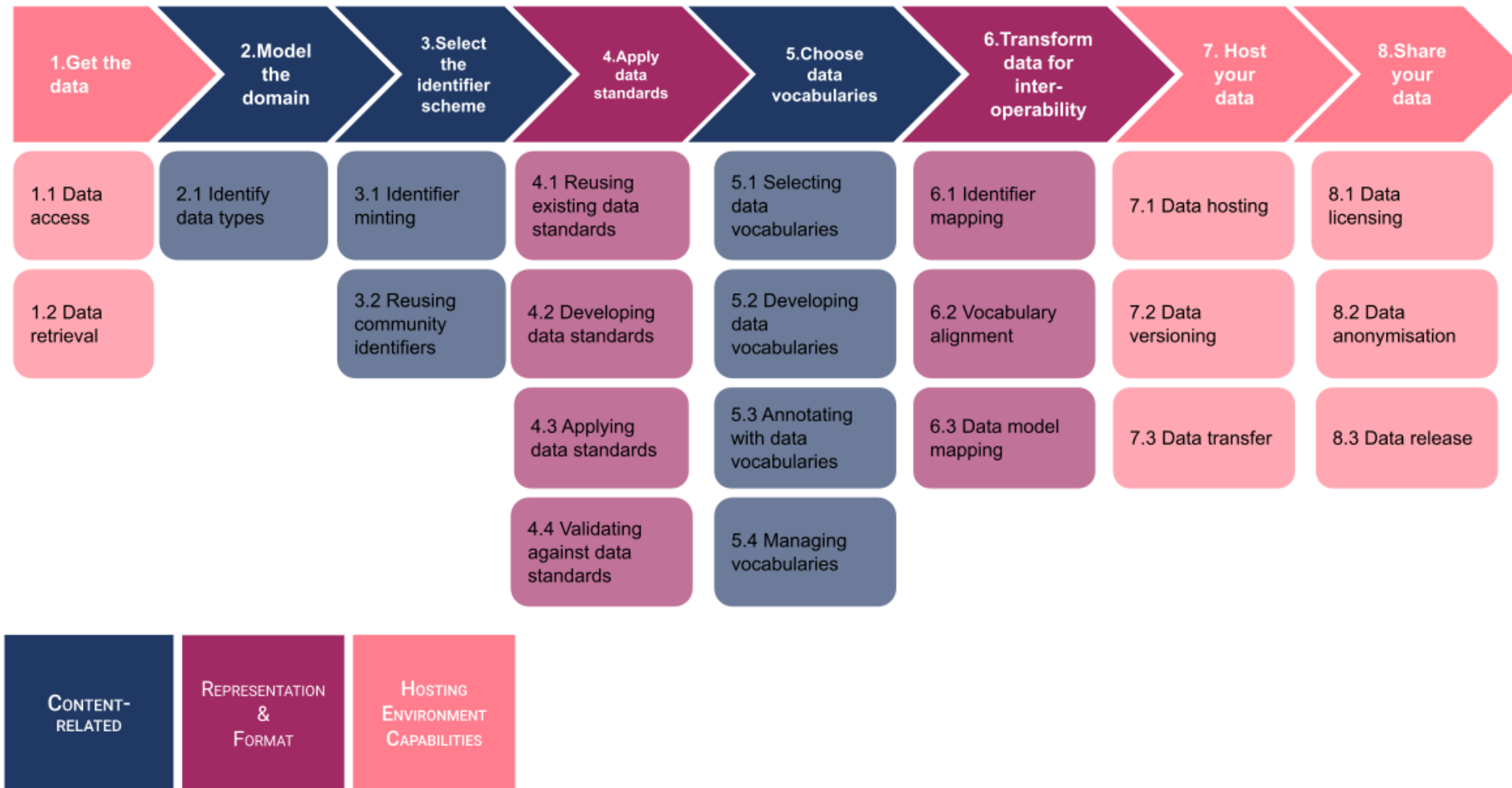
F1 · Globalt unike og persistente identifikatorer (PIDs) tildeles både metadataposten og dataobjektet



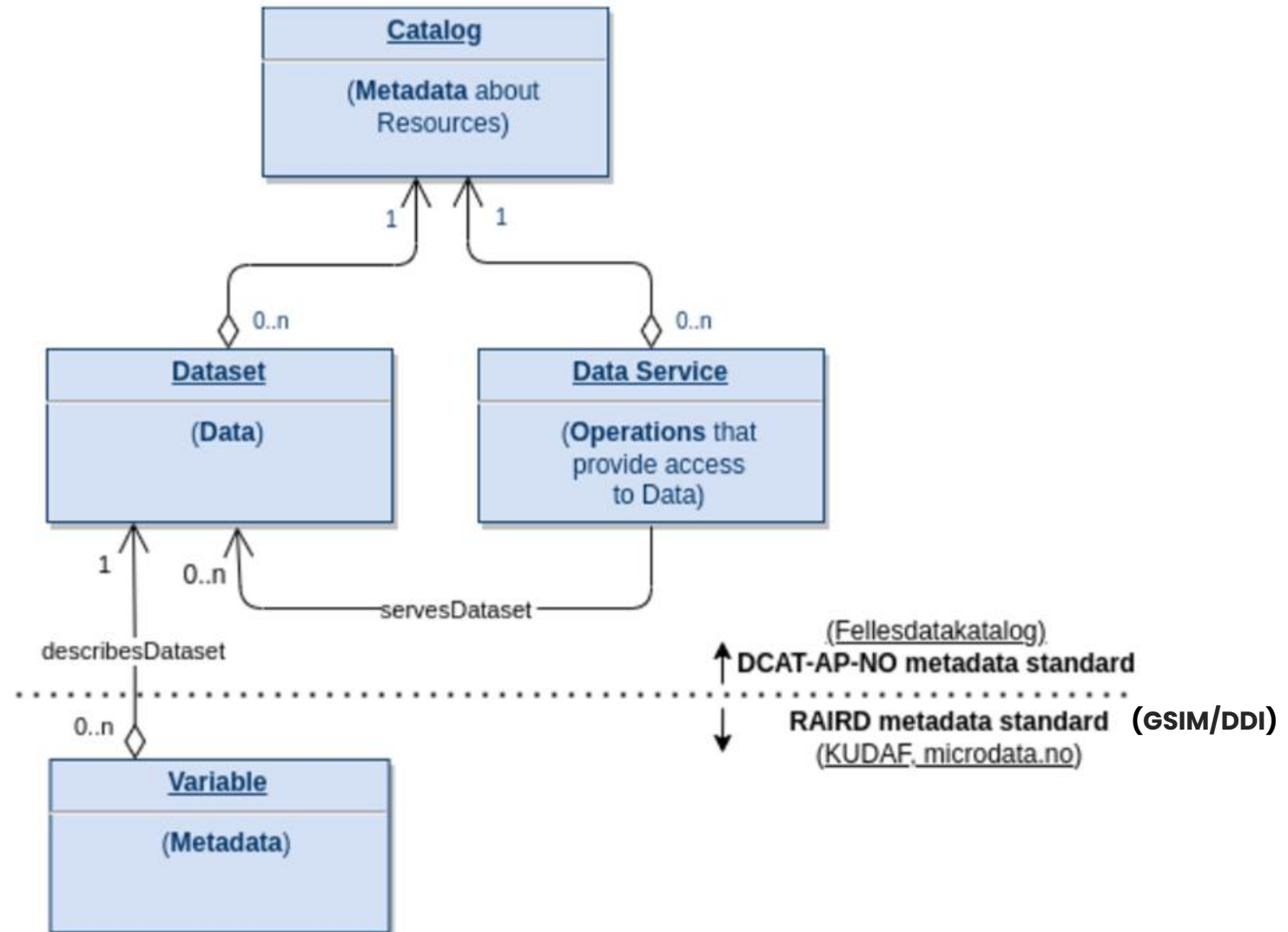
The FAIRification process



The FAIRification process

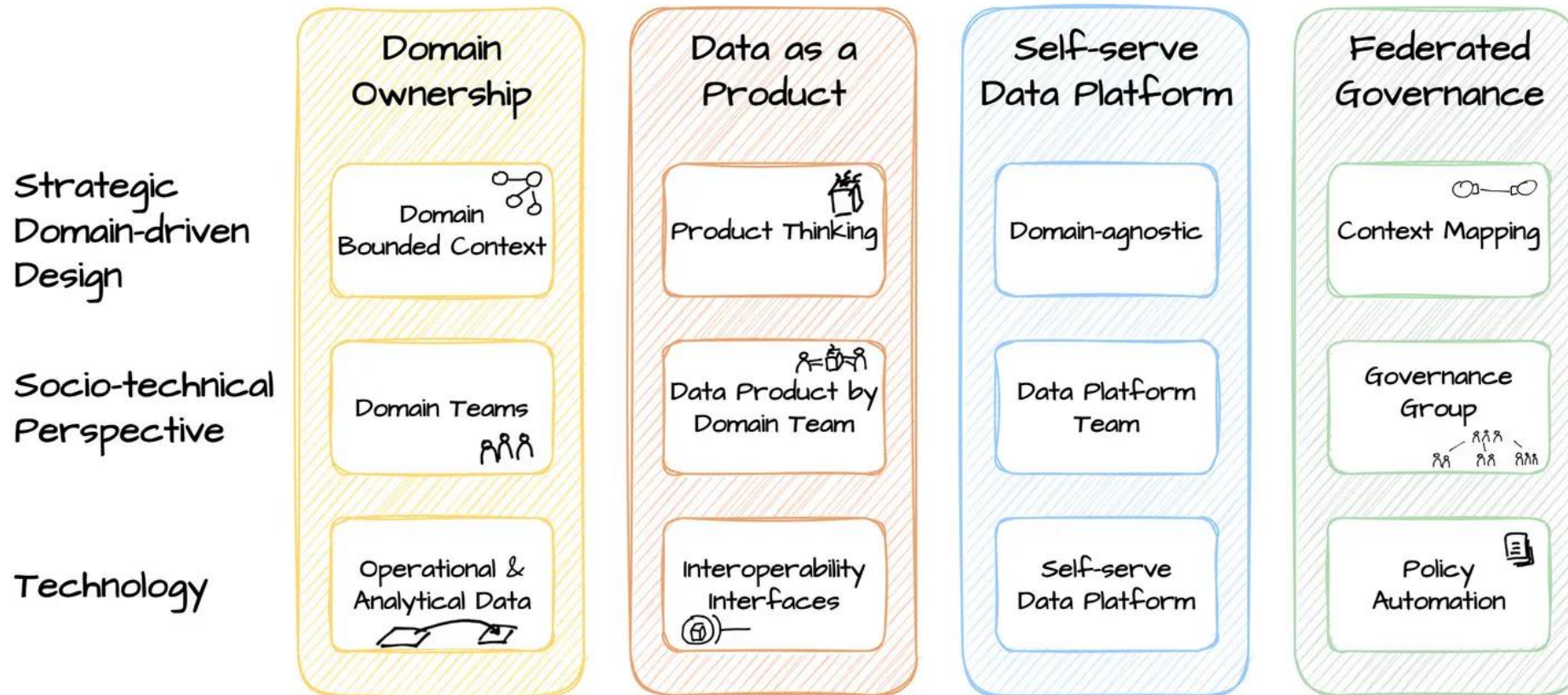


FAIR i praksis i et nasjonalt økosystem



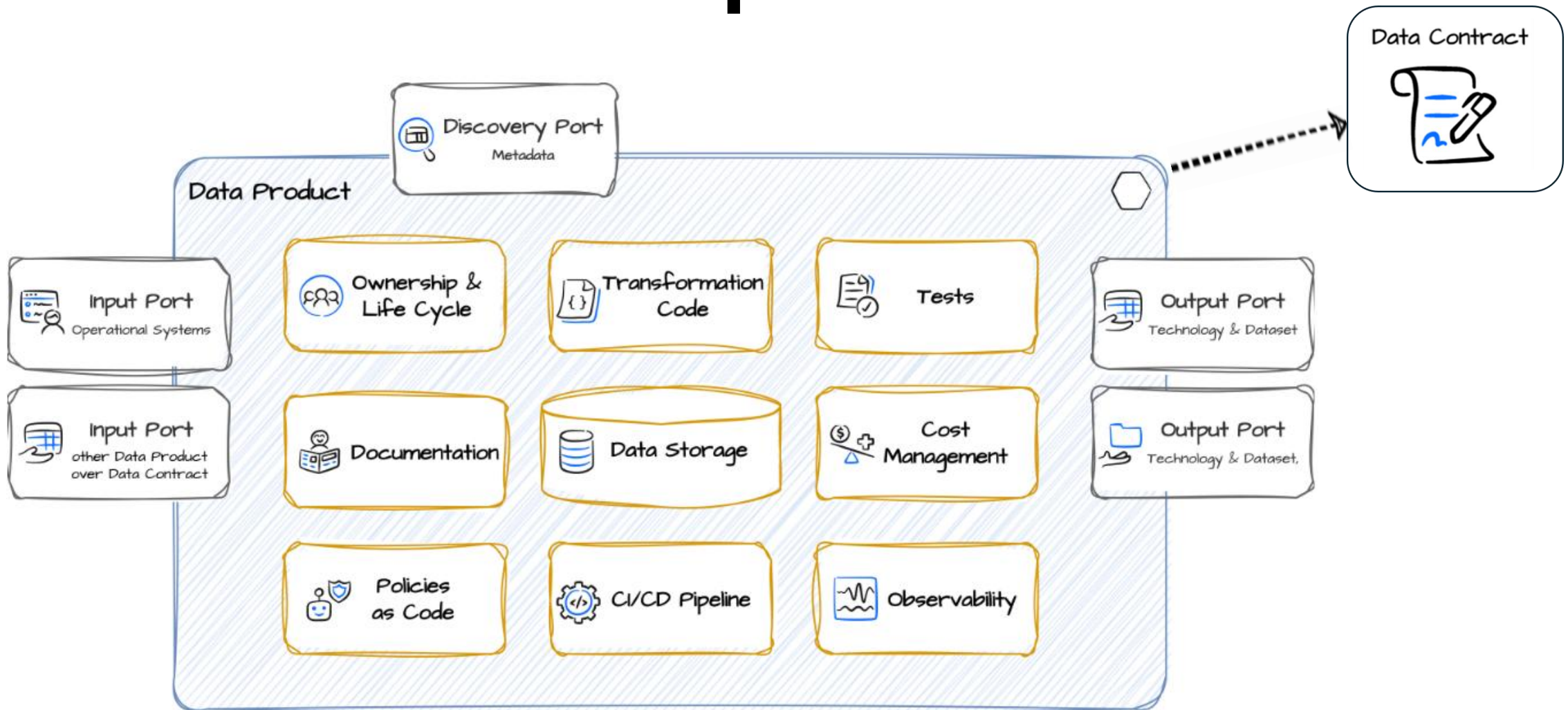
Data Mesh

What Is Data Mesh?



datamesh-architecture.com

Dataprodukt



datamesh-architecture.com

Agenda



1

Bakgrunnen for FAIR-prinsippene

2

Hvorfor er FAIR viktig – også utenfor forskningsverdenen?

3

FAIR-prinsippene i detalj

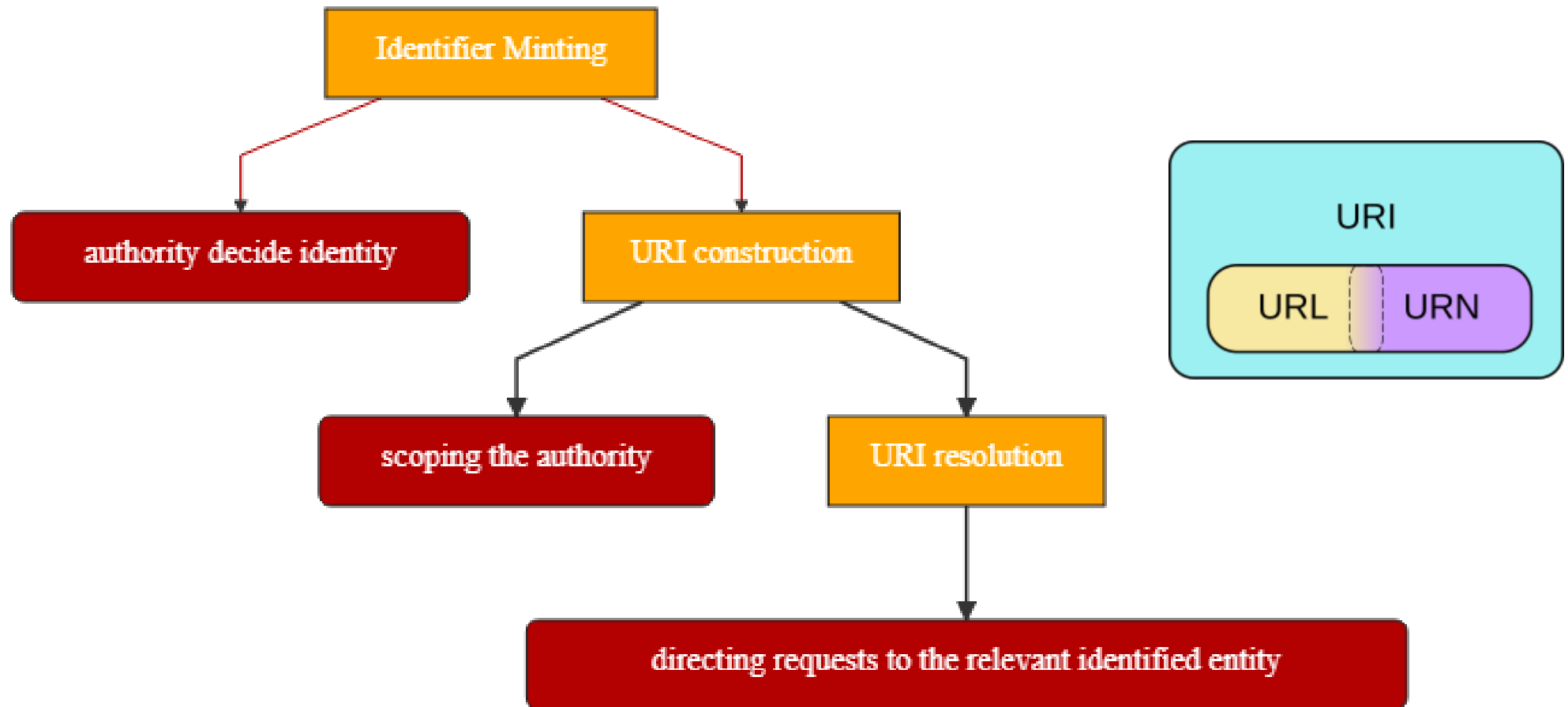
4

Persistente identifikatorer (PID) sentralt i FAIR

5

Behov for nye tjenester for tilordning av PIDs?

Hvordan sikres unike og persistente identifikatorer (PID)?



Eksempler på globalt unike og persistente identifikatorer som benyttes i dag

PID-infrastruktur og FAIR-prinsippene

FAIR: Findable · Accessible · Interoperable · Reusable

Hva krever FAIR?

FAIR-prinsippene F1 og A1 krever at data og metadata er tilordnet globalt unike, persistente identifikatorer (PIDs)

1

**Utgivere /
Organisasjoner**

ROR · ISNI · LEI

2

**Datakilder /
Datasett**

DOI · Handle · ARK

3

**Variabler /
Konsepter**

URI · DDI · GSIM

4

**Teknisk
Infrastruktur**

HTTP · Resolver · RDF

Nivå 1 – Utgivere og organisasjoner (datatilbydere)

Identifikator	Beskrivelse	Status
ROR (Research Organization Registry)	De facto standard for forskningsinstitusjoner globalt. Åpen, kuratert av Community of Scholars. Eks: ror.org/05xg72x27 (NTNU)	★ Anbefalt
ISNI (ISO 27729)	ISO-standard for personer og organisasjoner. Brukes i forlag og bibliotek.	Bibliotek
LEI (Legal Entity Identifier)	ISO 17442. Primært finanssektoren, men vokser som generell org-ID.	Finans
Org.nr. + prefix	Norske organisasjonsnumre kan brukes, men mangler global kontekst uten namespace.	Norsk bruk

ROR er gullstandarden for kunnskaps- og forskningssektoren

Nivå 2 – Datakilder og datasett

DOI

Digital Object Identifier

ISO 26324. Mest utbredt globalt.
Tildeles via DataCite / Crossref / Sikt i Norge via DataCite.

Handle

Handle System

Underliggende teknologi for DOI.
Brukes direkte i noen systemer.

ARK

Archival Resource Key

Alternativ til DOI. Brukes særlig
i bibliotek og arkiv.

W3ID

W3C Persistent URI

Stabil HTTP-URI for semantiske og lenkede data-
ressurser. Krever ikke et eget eksternt registersystem.

URN:NBN

Nasjonalt bibliografisk nr.

Brukes av Nasjonalbiblioteket for norske ressurser. I praksis
en URN som er bygd opp slik
<https://www.nb.no/items/URN:NBN:<id>>

I EU/EØS-kontekst

DCAT-AP

EU-standard for datakatalogmetadata — anbefaler
stabile URler

Europeana

Bruker primært DOI via DataCite for alle digitale
kulturarv-ressurser

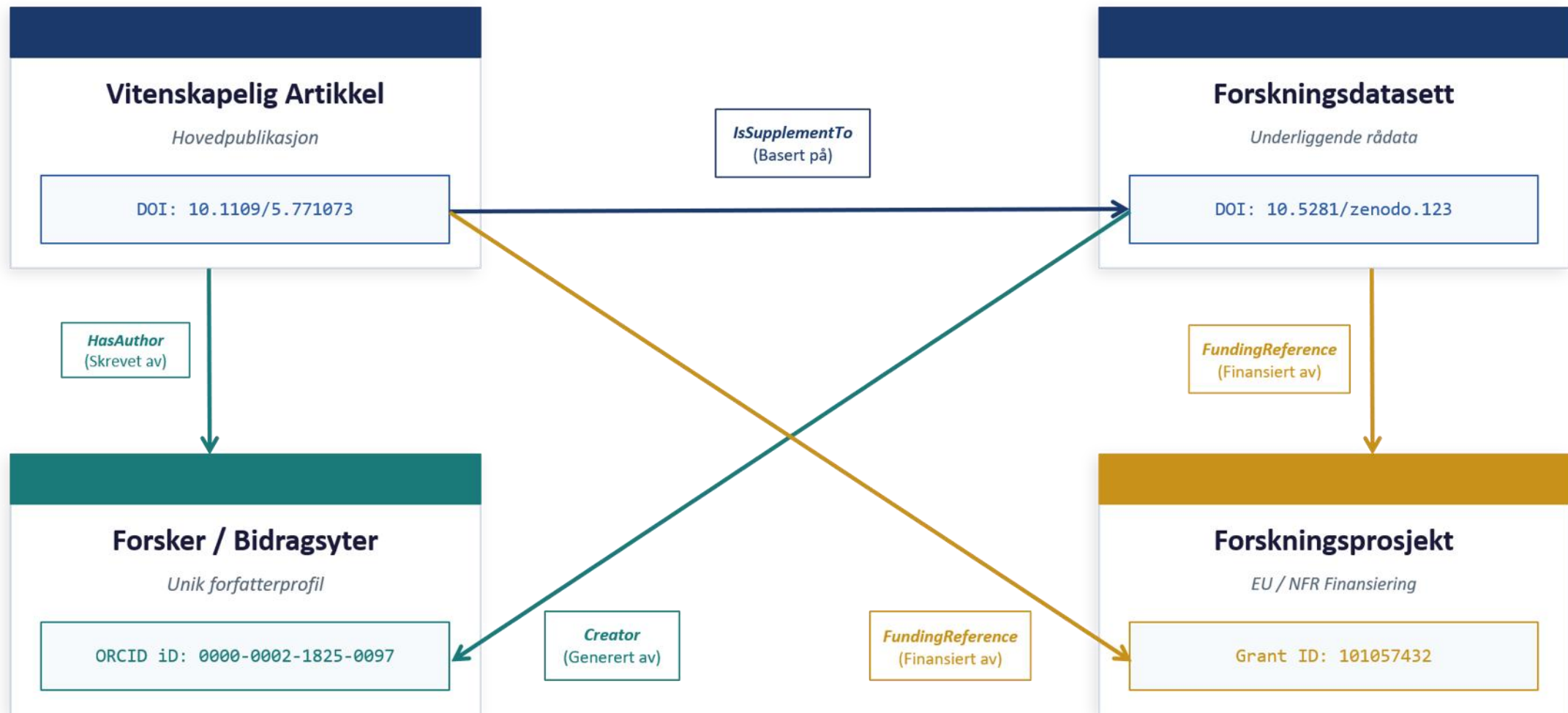
OpenAIRE

Europeisk open access-infrastruktur — DOI som
primær PID-type

EOSC

European Open Science Cloud — beveger seg mot
PID-krav for all finansiering

OpenAIRE Graph: Hvordan PID-er kobler forskning sammen



Nivå 3 – Variabler, konsepter og begreper

Det mest komplekse nivået — praksisen er ennå fragmentert

URI / IRI

W3C-standard

Grunnleggende mekanisme — alle semantiske webløsninger bygger på dette

GSIM

UNECE / Eurostat

Statistisk informasjonsmodell for variabler og datastrukturer

SKOS + URI

W3C SKOS

Standard for begrepsregistre og kodelister (EU Publications Office, SSB)

Wikidata QID

Wikimedia Foundation

Voksende praksis for å lenke konsepter til globalt delte identifikatorer

DDI

DDI-CDI

Internasjonal standard for sosialvitenskapelige og statistiske data

SNOMED / LOINC

Helsespesifikke

Domenespesifikke variabel-IDer innen helsesektoren

Nivå 4 – Den tekniske infrastrukturen

Felles krav på tvers av alle nivåer:

01 HTTP URI-er som resolves

PID-ene må være reelle HTTP-adresser som kan slås opp og returnerer innhold — ikke bare interne strenger eller lokale koder.

02 Registrar / Resolver som garanterer persistens

En betrodd tredjepart sikrer at URLen alltid peker til riktig ressurs, selv om eier skifter domene eller system (f.eks. doi.org, Handle System).

03 Metadata ved oppslag (Content Negotiation)

Samme URI returnerer HTML for mennesker og JSON-LD/RDF for maskiner — muliggjør automatisert dataintegrasjon og maskinlesbarhet.

04 Opasitetsprinsippet

PID-en bør ikke inneholde informasjon om eier, struktur eller innhold. Gjør IDen robust mot omorganisering og navneskift.

Oppsummering

F1

Globalt unik ID er en FAIR-forutsetning — ikke et teknisk tillegg

Org

ROR for utgivere, DOI (DataCite) for datasett — allerede i bruk i norsk sektor

Var

Variabler på URI/DDI/GSIM — praksisen er fragmentert, men standarder finnes.
Men hva med tjenester for tilordning?

Infra

HTTP + resolver + content negotiation + opasitet = teknisk minimumskrav

EU

EOSC og DCAT-AP beveger seg mot PID-krav som betingelse for finansiering



Follow the pattern

e.g. `http://{domain}/{type}/{concept}/{reference}`

Re-use existing identifiers

e.g. `http://education.data.gov.uk/id/school/123456`

Link multiple representations

e.g. `http://data.example.org/doc/foo/bar.html`

e.g. `http://data.example.org/doc/foo/bar.rdf`

Implement 303 redirects for real-world objects

e.g. `http://www.example.com/id/alice_brown`

Use a dedicated service

i.e. independent of the data originator

10 rules for persistent URIs



Avoid stating ownership

e.g. `http://education.data.gov.uk/ministryofeducation/id/school/123456`

Avoid version numbers

e.g. `http://education.data.gov.uk/doc/school/v1/123456`

Avoid using auto-increment

e.g. `http://education.data.gov.uk/id/school1/123456`

e.g. `http://education.data.gov.uk/id/school1/123457`

Avoid query strings

e.g. `http://education.data.gov.uk/doc/school?id=123456`

Avoid file extensions

`http://education.data.gov.uk/doc/schools/123456.cs`