# A data storage for generic and heterogenous scientific data

*Ketil Malde, Tomasz Furmanek, and Esmael Hassen*
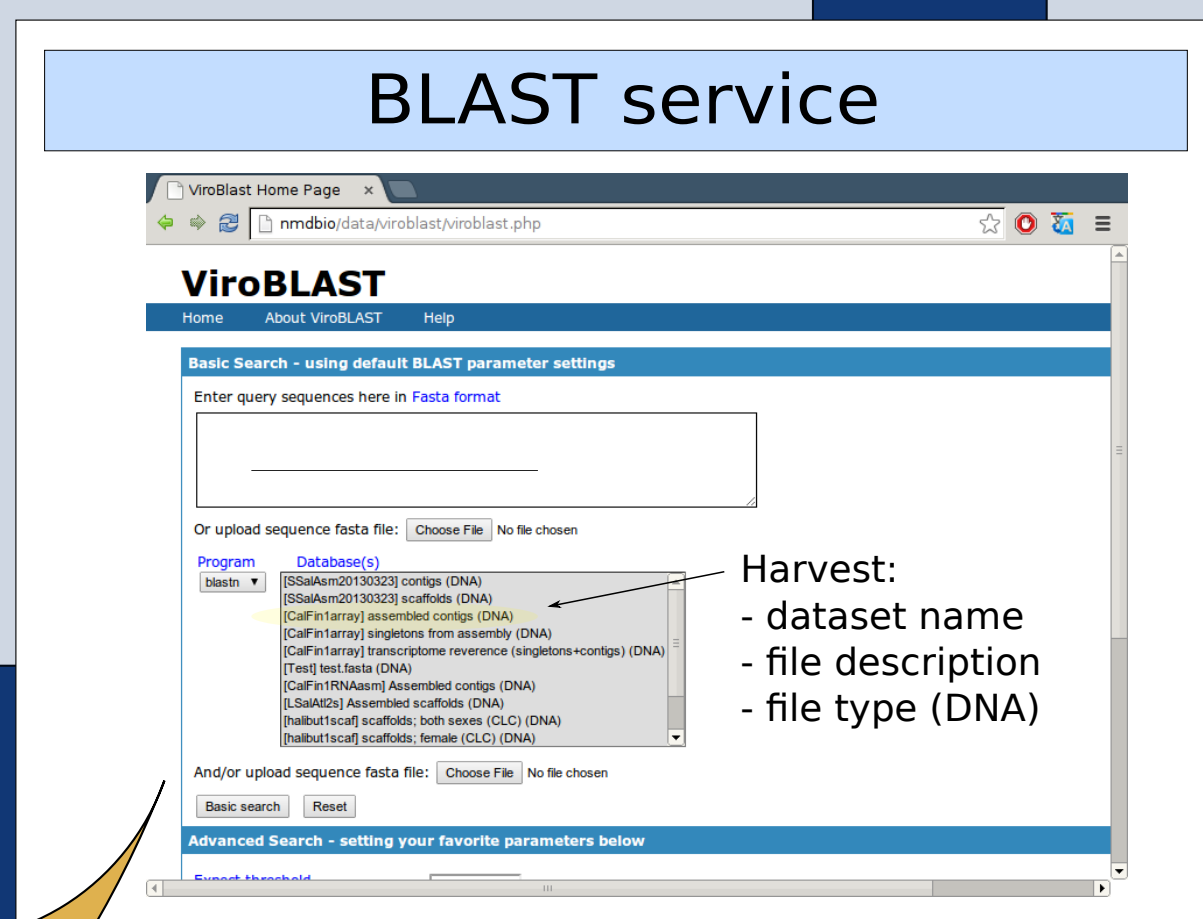
## Data Submission

**Is simple and easy**

Uses domain specific file formats
Auto-generates most metadata
Free-text descriptions/
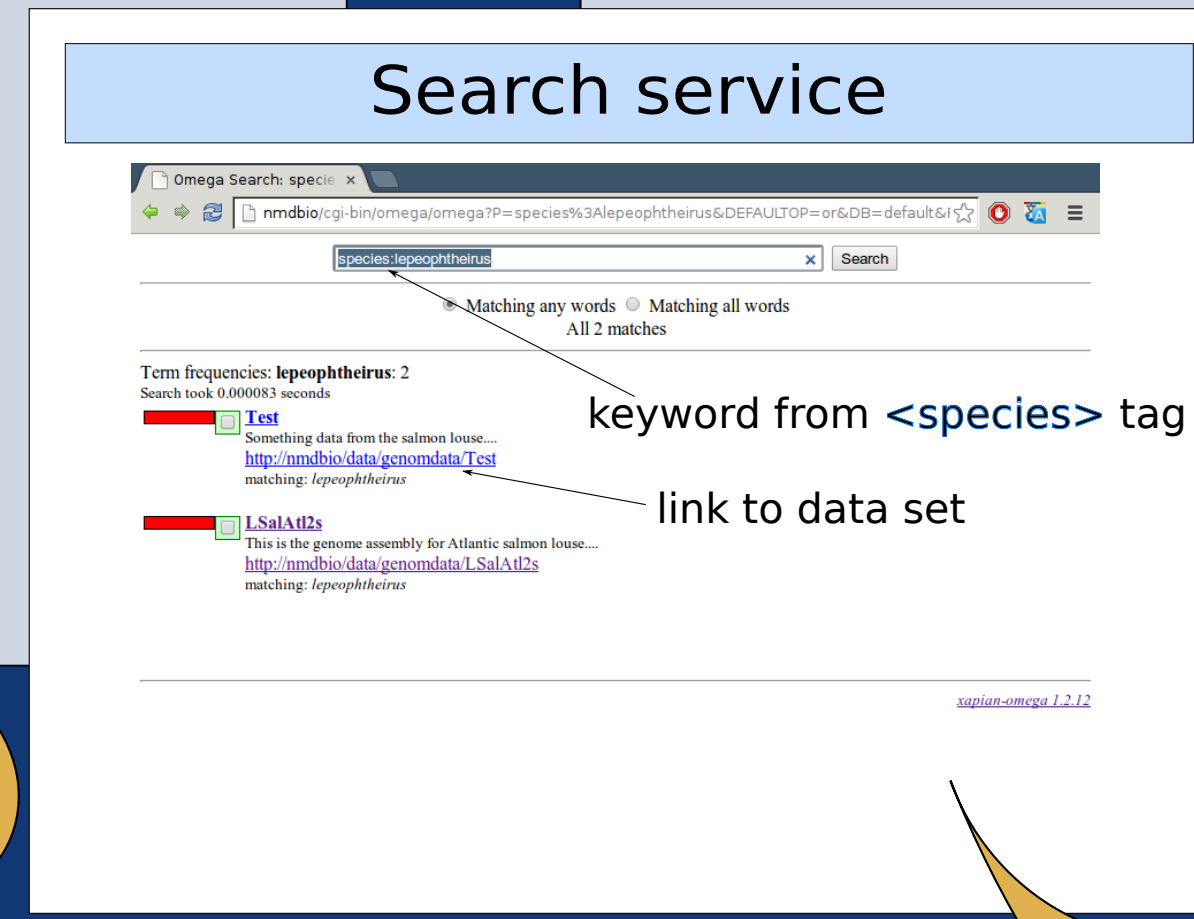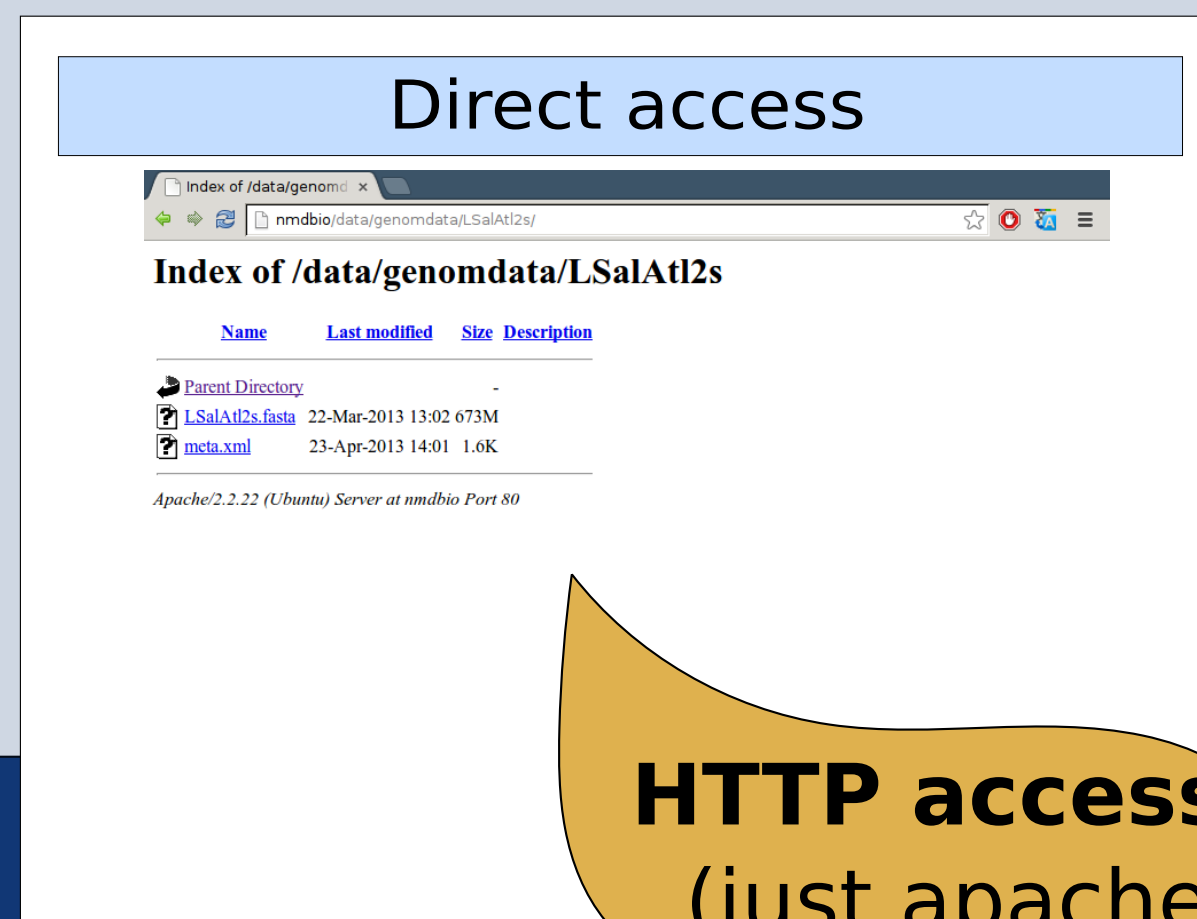(almost) no mandatory fields

**Provenance and links**

Unique, persistent
IDs (citation)
Link datasets

## Search/services

**Self-contained services**

Technology-agnostic (i.e. use
any relational database system)
Extract and index only relevant
data from data store
Independent and modular

## Data Access

**File-based storage**

Access through HTTP,
FTP, rsync, bittorrent....

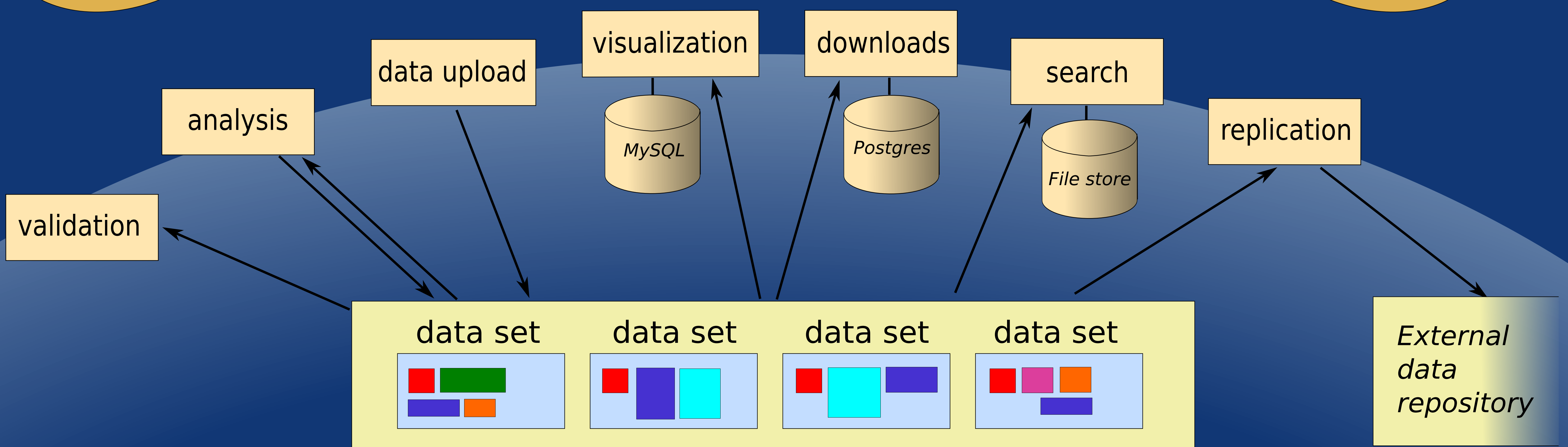Easy replication

**Domain-specific formats**

Easy for domain
experts
No data conversion

**Provenance**

Identify origins of data

---

BLAST service

ViroBLAST

Harvest:
- dataset name
- file description
- file type (DNA)

Direct access

Index of /data/genomdata/LSalAtl2s

Search service

Omega Search spec...

blocking any words ☐ Matching any words
All 2 matches

Term Sequences: hypophlebrins 2
Search took 0.000002 seconds

keyword from <species> tag

link to data set

**HTTP access**
(just apache)

**Specialized**
search service -
scans *relevant* data,
ignores rest

**Generic**
search service
scans and indexes
*metadata* only

## specialized services

visualization    downloads    search

analysis    data upload    *MySQL*    *Postgres*    *File store*    replication

validation

| data set | data set | data set | data set |
| --- | --- | --- | --- |

*External data repository*

## generic file storage

---

## Metadata

**XML format**

automatic validation
tagging (TSN, geoloc, etc)
free text descriptions

**Directory listings**

Checksums for integrity
File type tagging

**Data set relationships**

Obsolescense and replacement
Dependencies
Aggregation and extraction

## Extensibility
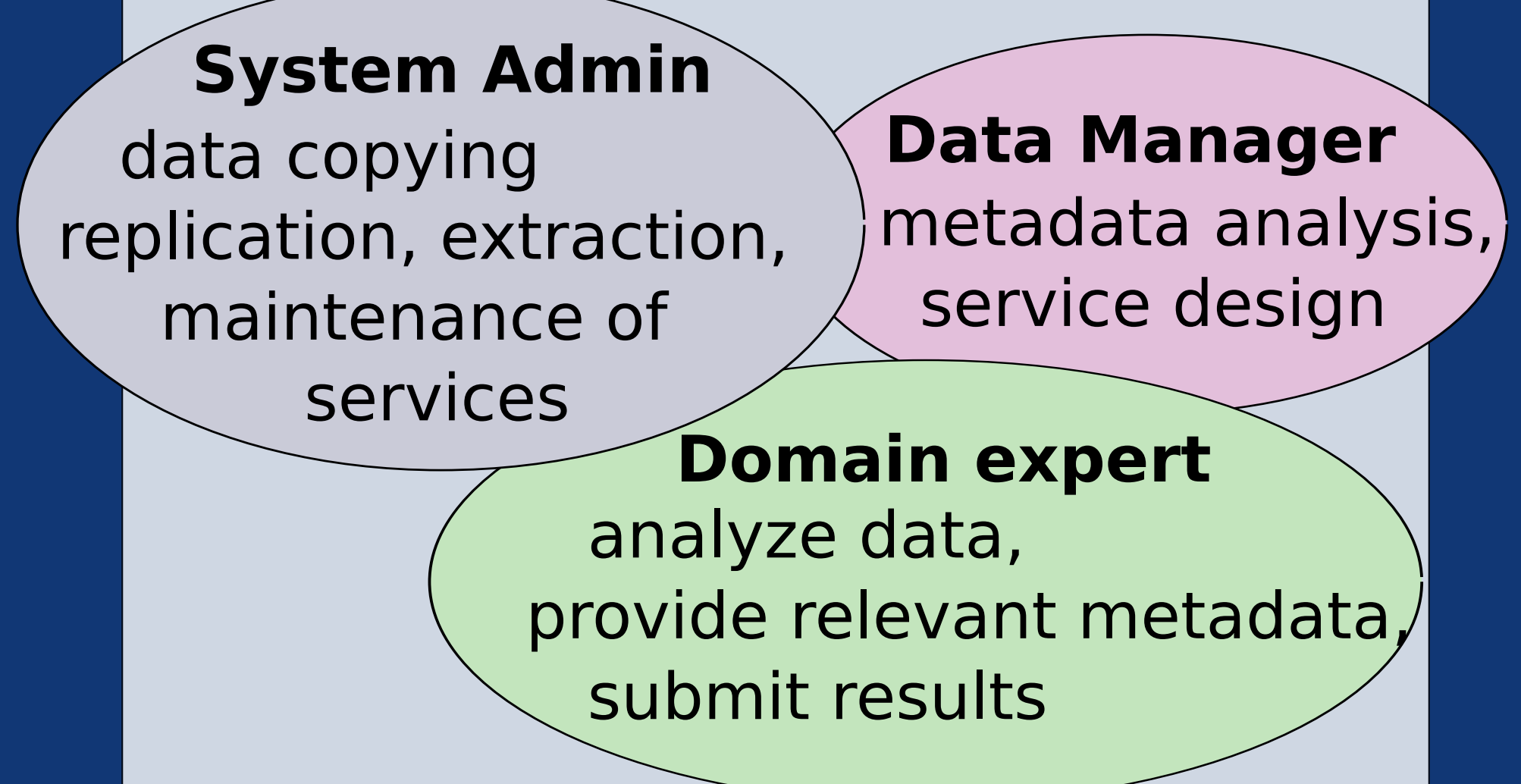
**New data types**

Technology advances means
frequent new data types
Adding a new data type is a
two minute operation

**New technologies**

Services are independent,
can use different technologies

Easy integration of off-the-shelf
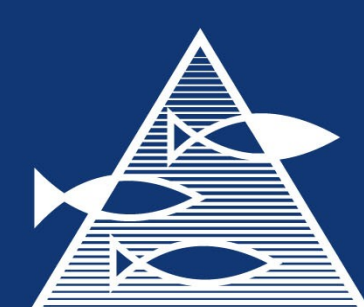products

## Separation of concerns

Separate roles with separate
skill sets

**System Admin**
data copying
replication, extraction,
maintenance of
services

**Data Manager**
metadata analysis,
service design

**Domain expert**
analyze data,
provide relevant metadata,
submit results

---