Decoupling data collocation and storage performance using scientifically meaningful object identifiers

HPC I/O in the Data Center Workshop

Olivier Iffrig, James Hawkes, Simon Smart, Tiago Quintino ECMWF olivier.iffrig@ecmwf.int



ECMWF's Forecasting Systems

Established in 1975, Intergovernmental Organisation

- 23 Member States | 12 Cooperating States
- 350+ staff

24/7 operational service

- Operational NWP 4x HRES+ENS forecasts / day
- Supporting NWS (coupled models) and businesses

Research institution

- Experiments to continuously improve our models
- Reforecasts and Climate Reanalysis

Operate 2 EU Copernicus Services

ECFCMWF



- Climate Change Service (C3S)
- Atmosphere Monitoring Service (CAMS)
- Support Copernicus Emergency Management Service (CEMS) •



Starting from weather forecasting...

Handling weather data at ECMWF

The ECMWF operational workflow



Storage View of Workflow



Labelling and accessing data: the MARS language

•

•	Operational data:
	RETRIEVE,

TYPE	=	FC,
CLASS	=	OD,
EXPVER	=	0001,
STREAM	=	OPER,
DATE	=	20160517,
TIME	=	1200,
STEP	=	12/24/36,
PARAM	=	Z/T,
LEVTYPE	=	PL,
LEVELIST	=	1000/500,
TARGET	=	out.grib

Research expen	rim	ent:
ARCHIVE,		
TYPE	=	FC,
CLASS	=	RD,
EXPVER	=	ABCD,
STREAM	=	OPER,
DATE	=	20160517,
TIME	=	1200,
STEP	=	12/24/36,
PARAM	=	Z/T,
LEVTYPE	=	PL,
LEVELIST	=	1000/500,
SOURCE	=	in.grib

Product genera	atior	ר:
DISSEMINATE	,	
TYPE	=	FC,
CLASS	=	OD,
EXPVER	=	0001,
STREAM	=	OPER,
DATE	=	20160517,
TIME	=	1200,
STEP	=	12/24/36,
PARAM	=	Z/T,

T T TT / TT.T		Δ/ ± /
LEVTYPE	=	PL,
LEVELIST	=	1000/500,
GRID	=	0.5/0.5,
AREA	=	europe
TARGET	=	XXX:ZZ

Observations: ۲ RETRIEVE,

۲

TYPE = OB, OBSGROUP = CON, OBSTYPE = SSD, = 20150201,DATE = 00,TIME = 1439RANGE

Unique and semantic way to describe all ECMWF data

Curating data collections in a meaningful way

POSIX namespace

/scratch/user1/exp42/2020-03-18/output_004
/scratch/user1/exp42/20220415/out39
/scratch/user1/exp42/20220415/out48.txt
/scratch/user2/DATA/2021/FC/1208/STEP13
/shared/project4/exp_abcd/20220502/step_002
/prod/0001/20220329/step053/rank0012
/prod/0001/20220329/step053/rank0014
/prod/0001/20220329/step053/rank0015

Semantic namespace

```
source=user1, experiment=42, date=20200318, step=4
source=user1, experiment=42, date=20220415, step=39
source=user1, experiment=42, date=20220415, step=48
source=user2, experiment=fc, date=20211208, step=13
source=project4, experiment=abcd, date=20220502, step=2
source=prod, experiment=1, date=20220329, step=53, rank=12/to/15
```

FDB (version 5)



steps=0/240/by/3

date=01011999/to/31122015.

FDB 5 Semantics

- 1. ACID Transactional
- 2. Write blocks until data handed over Asynchronous
- 3. flush() blocks until data is visible Consistent
- 4. Write-once, don't overwrite *Immutable*
- 5. Data can be masked *Versioned*

- All I/O operations are asynchronous, so computation can continue
- Distributed to all servers using a *Rendezvous Hash*, so no synchronisation needed

To summarise

- Semantic data access system at ECMWF
 - Based on meteorological keywords
 - Unique description of all datasets
 - Language defined by a schema
 - Decouple data specification from the actions of the system
- FDB: semantic object store software
 - github.com/ecmwf/fdb
 - Robust and performant
 - Implementation tied to meteorology





... Towards a domain-agnostic storage system

The Data Access and Storage Interface



This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 955811. The JU receives support from the European Union's Horizon 2020 research and innovation programme and France, the Czech Republic, Germany, Ireland, Sweden, and the United Kingdom.

Enter DASI: the Data Access and Storage Interface

- Goal: Lower the entry barrier to effective semantic data management
- Pillars:

12

- Declaration of domain-specific languages via a schema
- Efficient data indexing
- Abstraction of the underlying storage system



The IO-SEA project



No-Sea

DASI: decoupling data collocation and storage performance

DASI in the context of IO-SEA



14



Interacting with the DASI

- Identifying data
 - Sets of keyword/value pairs
 - Defined by a schema
 - Key: identifies a single object **uniquely**
 - Query: multiple values allowed: lists, wildcard
- Actions
 - Put(Key, Bytes)
 - Get(Query) \rightarrow [(Key, Bytes)]
 - Evaluates the cartesian product of all combinations
 - List(Query) \rightarrow [Key]
 - Supports incomplete queries
 - Delete(Query)

Key

model: covid-spread
date: 20210112
experiment: 42
epoch: 123
variable: R0

Query

model: any
date: any
experiment: 42
epoch: [123,124,125]
variable: any



The process in a nutshell

ed

16

01/06/2022



The DASI schema

- DASI provides an API for scientifically meaningful access to data
- The schema follows a domain-specific hierarchical taxonomy
- The schema is split into three levels that can be interpreted by the backend

```
    model > experiment
        b date > epoch
        b variable
        b location > variable
        b date > location
        b variable
        variable
        b institute > location
        b date > method
        b variable
```

model: covid-spread
date: 20210112
experiment: 42
variable: R0
epoch: 123

17



Building your own schema

- Identify your data collection
 - Which data do you want to store together?
 - What are your base objects?
- Define your domain-specific keywords
 - How do you uniquely identify an object?
 - If needed, what are the different sets of keywords that you need?
- Order your taxonomy
 - Build a tree that lists all the possible sets of keywords
 - Choose the relevant order for the keywords
- Decide on where to put the boundaries for the three storage levels

model: covid-spread
date: 20210112
experiment: 42
variable: R0
epoch: 123

01/06/2022

The DASI backends

- Index backend
 - Resolves object locations
 - Resolve(SplitKey) \rightarrow Location
 - Lists available data
 - List(SplitQuery) \rightarrow [SplitKey]
 - Creates / Deletes entries
 - Create(SplitKey, Location)
 - Delete(SplitKey)

- Storage backend
 - Implements atomic operations
 - Put(Location, Bytes)
 - Get(Location) \rightarrow Bytes
 - Delete(Location)
 - Implements three storage levels
 - Allow for performance optimisation
 - Small objects \rightarrow aggregate at third level
 - Many objects → split between first and second levels



Wrapping up



- Semantic data access provides
 - Description of data in scientific terms
 - Consistency and reproducibility
 - Decoupling between performance and data collocation
- Software
 - FDB: github.com/ecmwf/fdb
 - DASI: <u>github.com/ecmwf-projects/dasi</u>



This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 955811. The JU receives support from the European Union's Horizon 2020 research and innovation programme and France, the Czech Republic, Germany, Ireland, Sweden, and the United Kingdom.