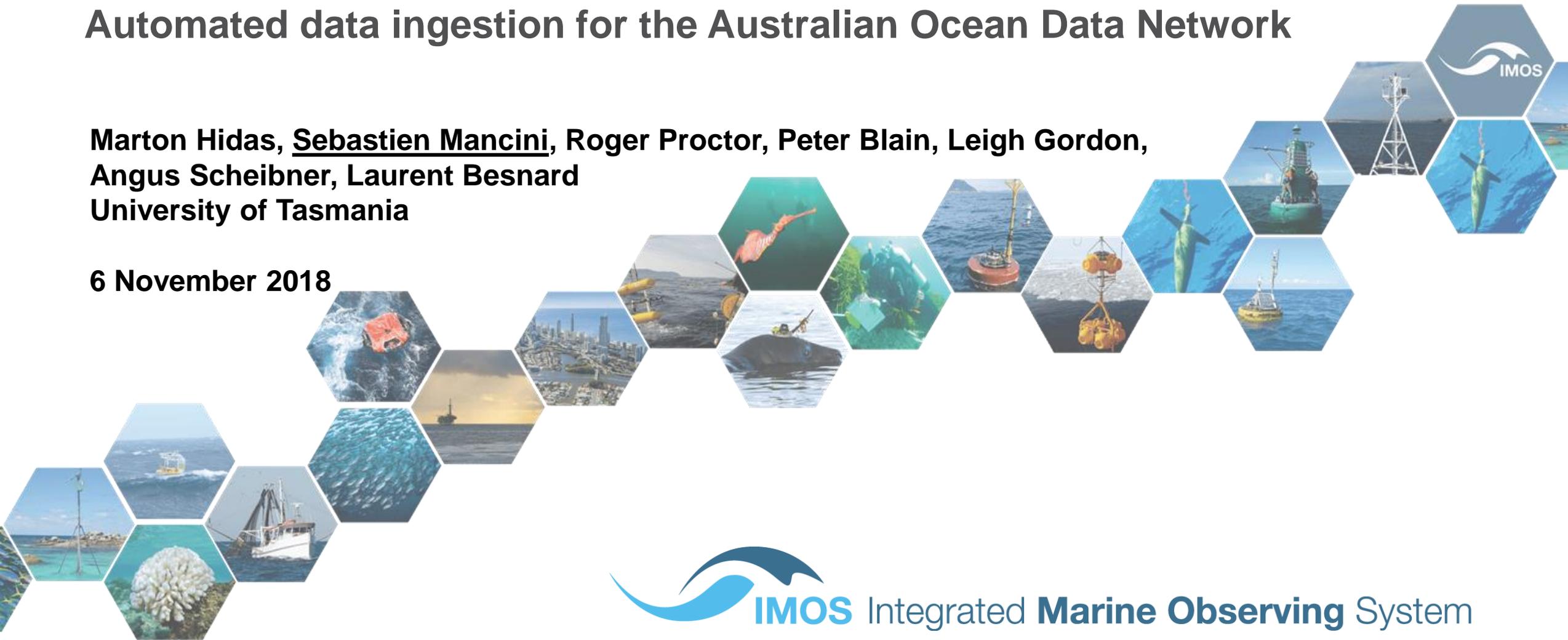


Australia's Integrated Marine Observing System (IMOS)

Automated data ingestion for the Australian Ocean Data Network

**Marton Hidas, Sebastien Mancini, Roger Proctor, Peter Blain, Leigh Gordon,
Angus Scheibner, Laurent Besnard**
University of Tasmania

6 November 2018



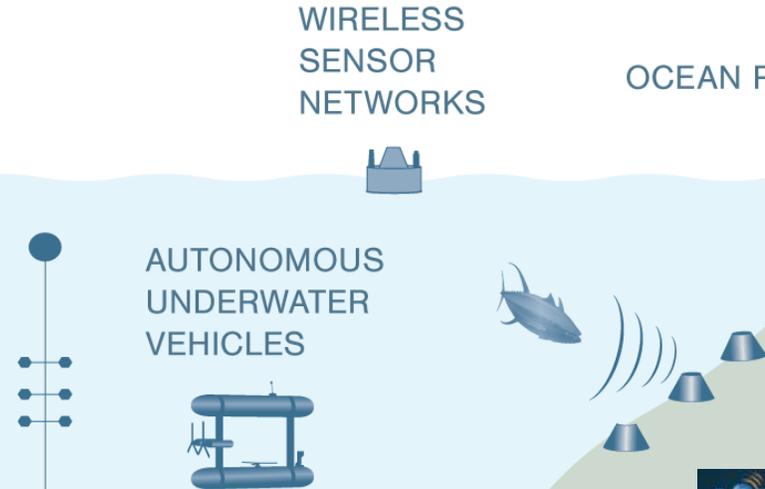
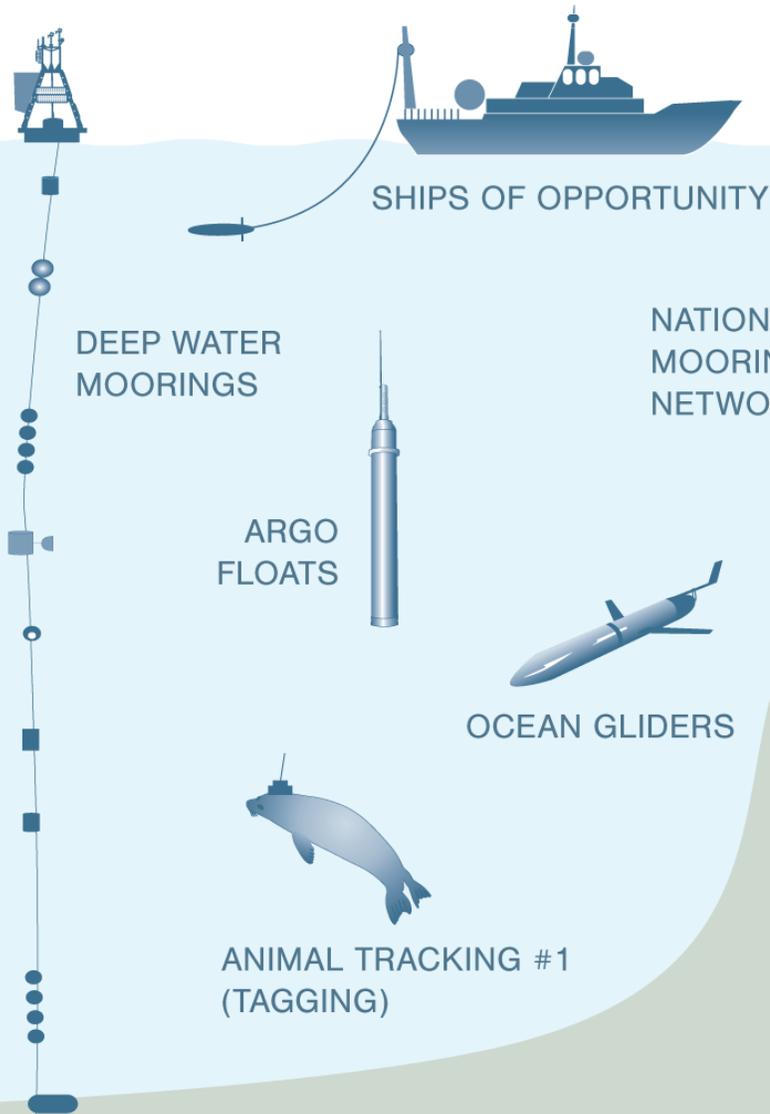
Outline of the talk

12 minutes, plus 3 minutes for questions

1. About IMOS
2. Data providers versus data users
3. Automated pipeline
4. Future improvements



IMOS Facilities



All data discoverable, accessible, usable and reusable



<https://portal.aodn.org.au>

Why AODN?

Data providers



Diversity in

- Technologies
- Instruments
- Platforms
- Organisations
- People

Diversity in

- Data products required
- Data volume required
- Method of access
- Preferred formats
- Tools used

Data users/use cases

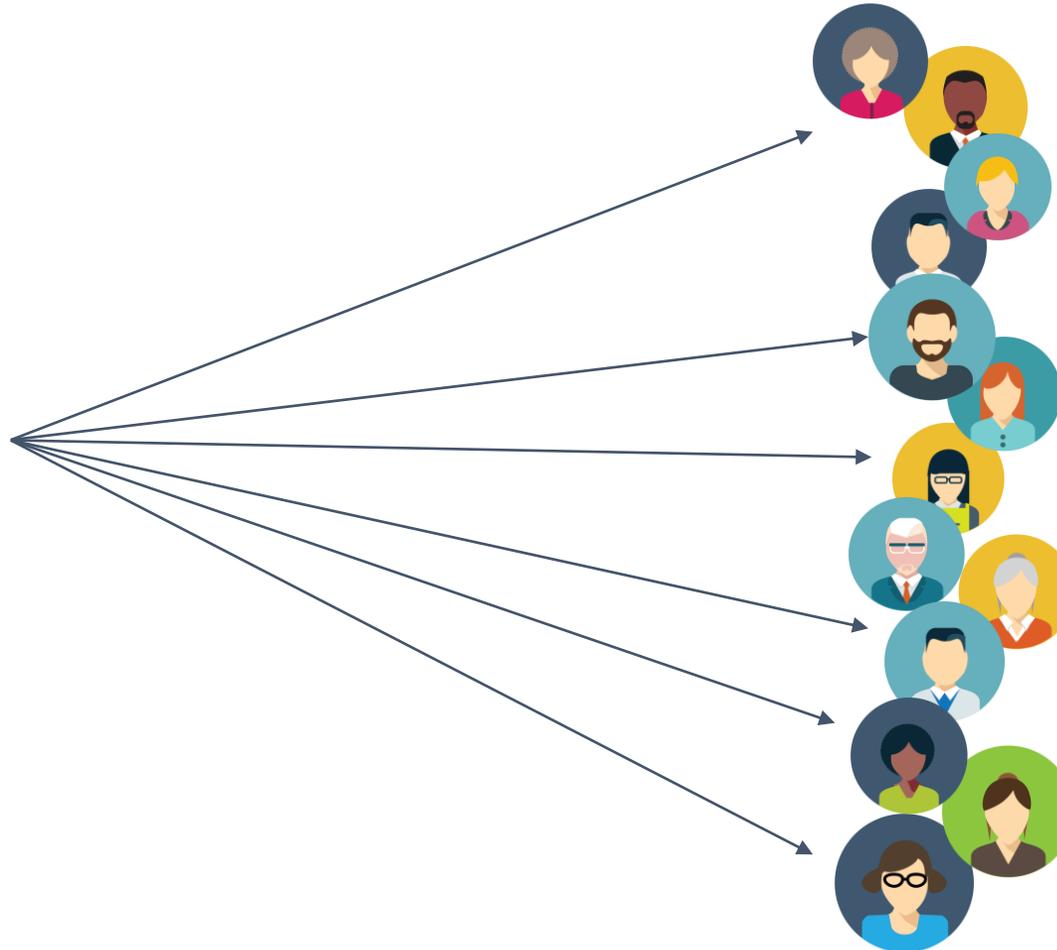


- Research scientists
- Government departments
- Private industry
- Fisheries
- Managers
- Policy makers
- Shipping
- Recreational fishing/boating
- Etc...

Why AODN?

Data providers

Data users/use cases



Each provider has to cater for many users

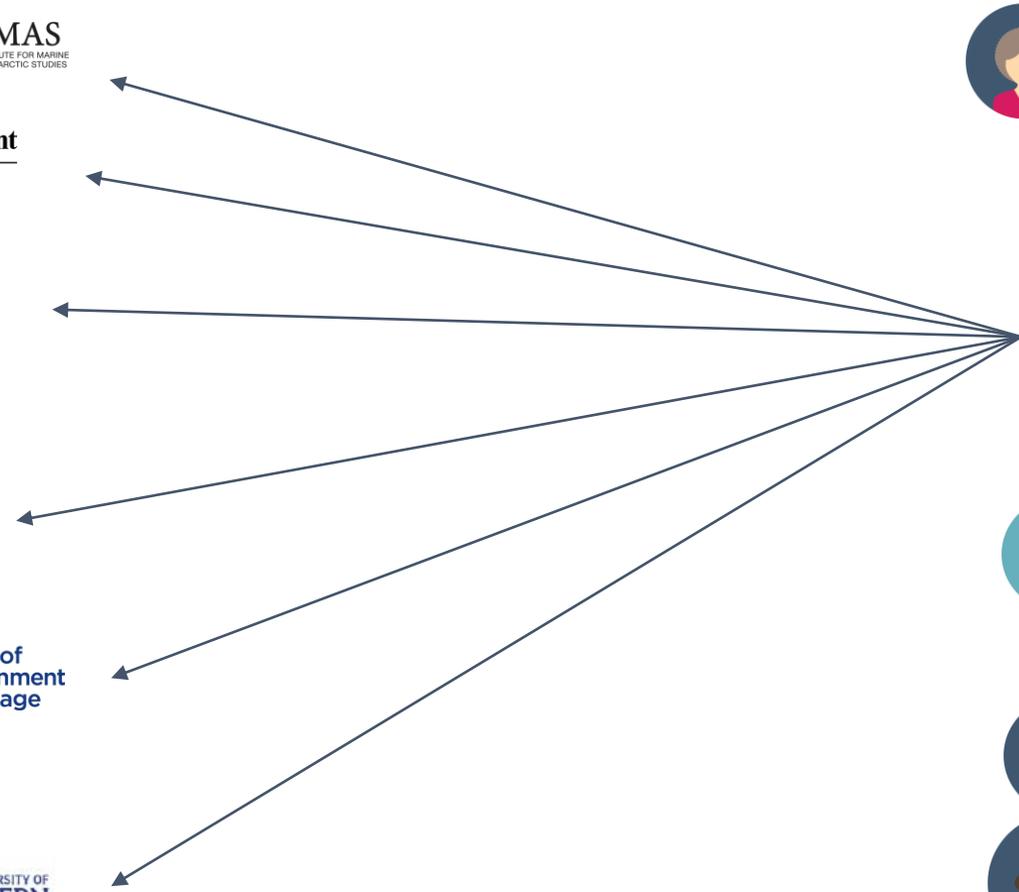
Why AODN?

Data providers

Data users/use cases



Each user may need to obtain data from multiple sources



Why AODN?

Data providers

Data users/use cases



Each user may need to obtain data from multiple sources

Each provider has to cater for many users

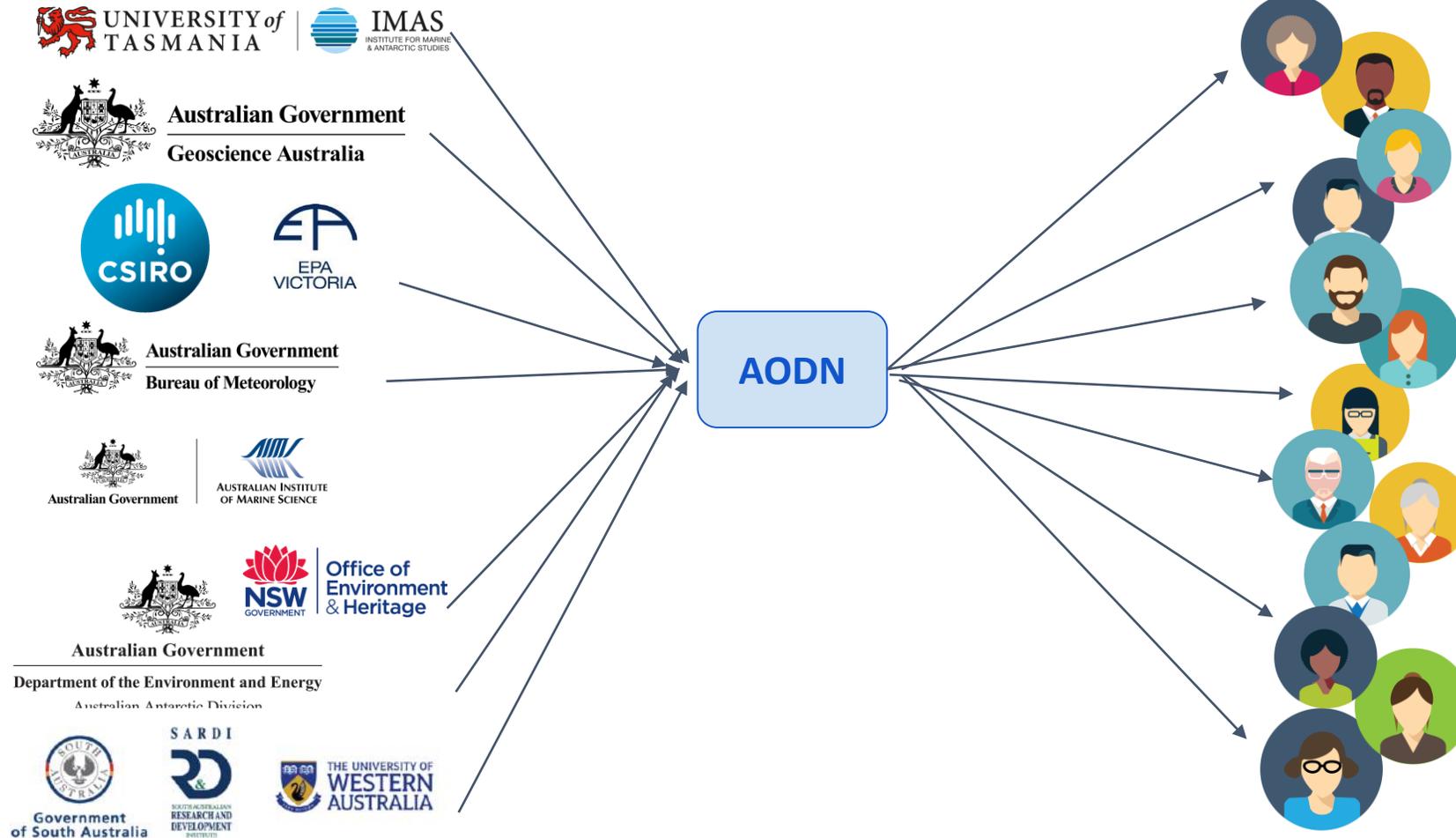
Why AODN?

Data providers

Data users/use cases

Standard ingestion

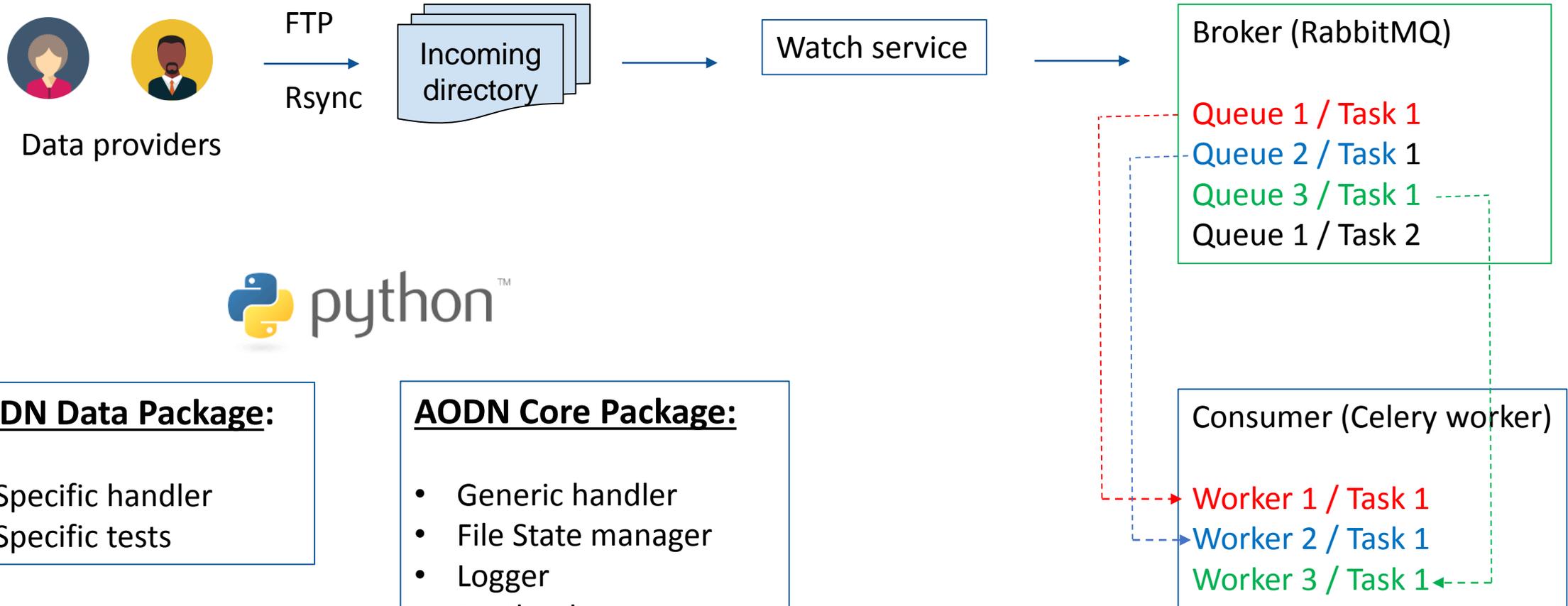
Standard access



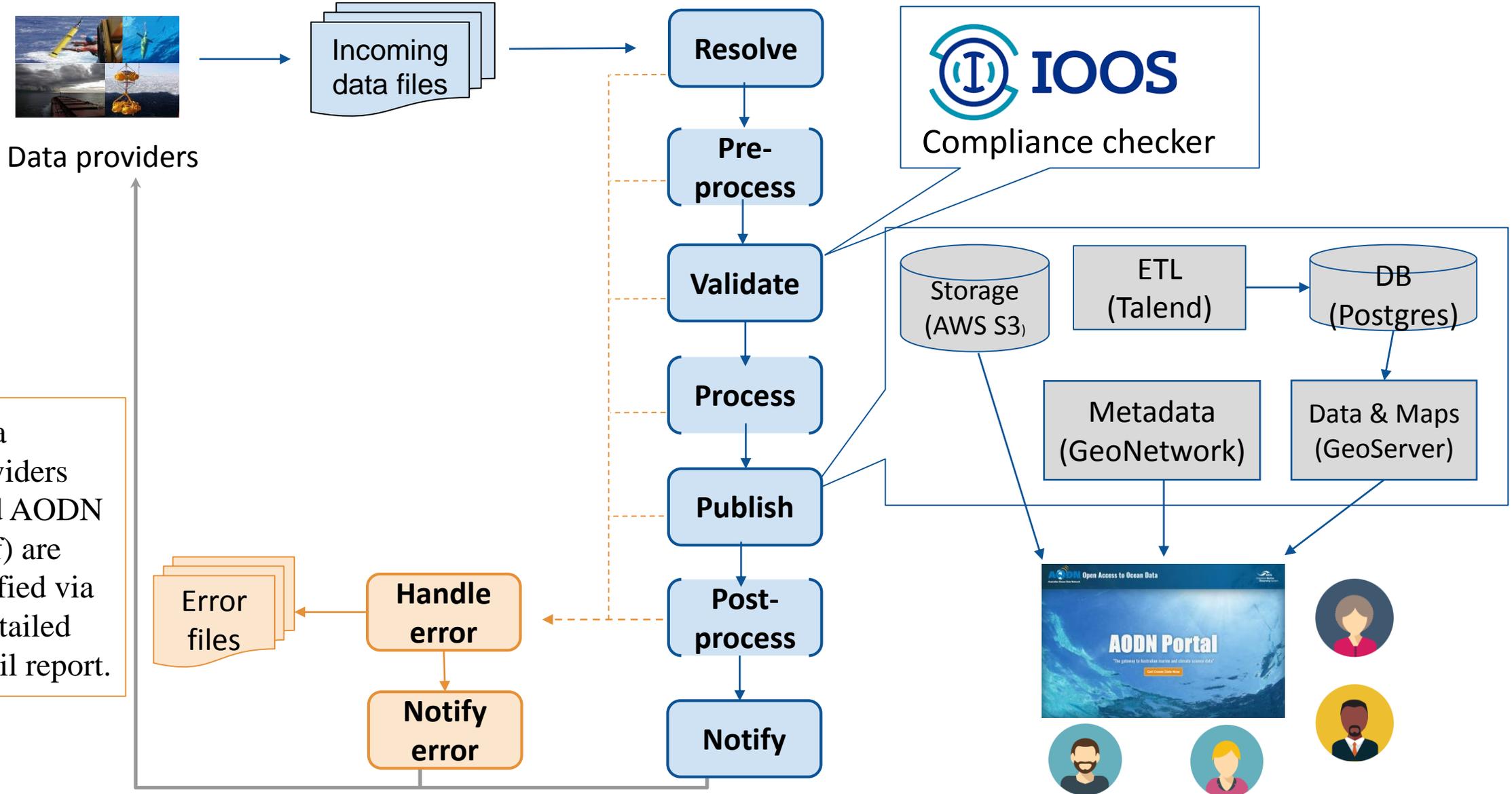
IMOS standard ingestion: Key design criteria

- Keep data safe
- Make only high quality data for users
- Do not increase load on front facing systems
- Make data available as quickly as possible
- Robustness
- Transparency

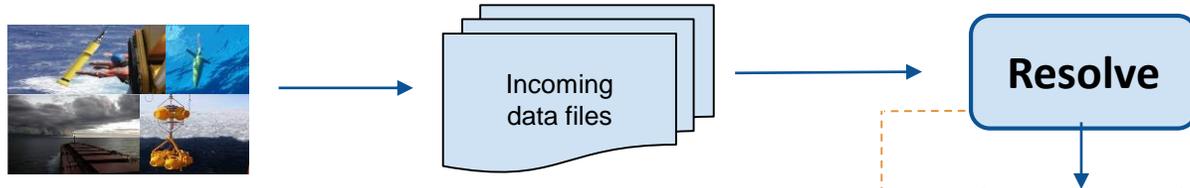
Architecture overview of the pipeline



Pipeline workflow: the state machine



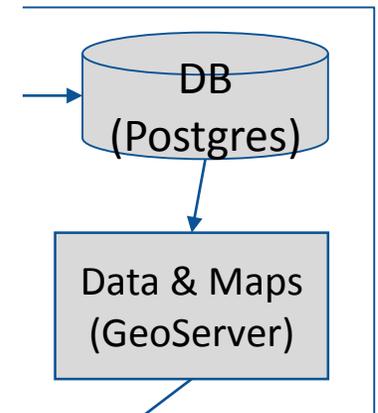
Pipeline workflow: the state machine



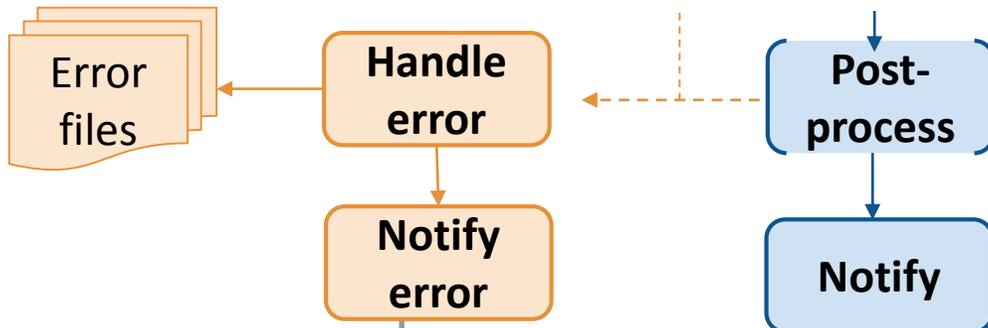
Data providers

Some Numbers:

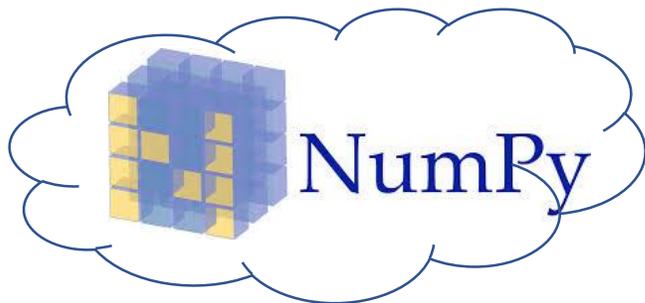
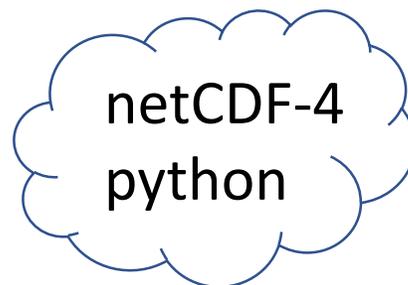
- 32 Organisations providing data
- 60 Data uploaders (+ a handful that we pull from)
- 24 Handler classes
- 43 Ingestion pipelines
- 86 Data/metadata “Harvesters” (95 Harvest jobs)
- 204 Data collections in the AODN Portal



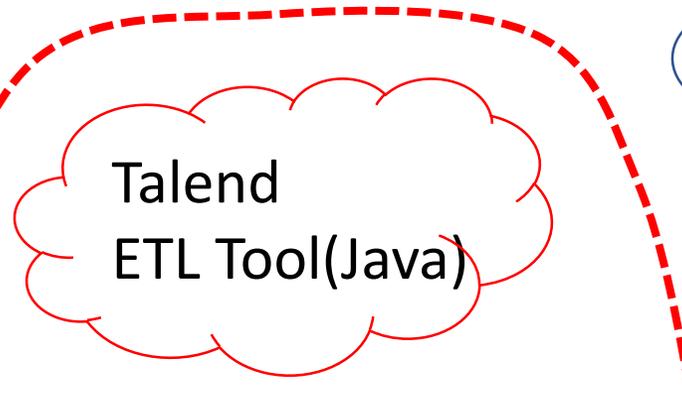
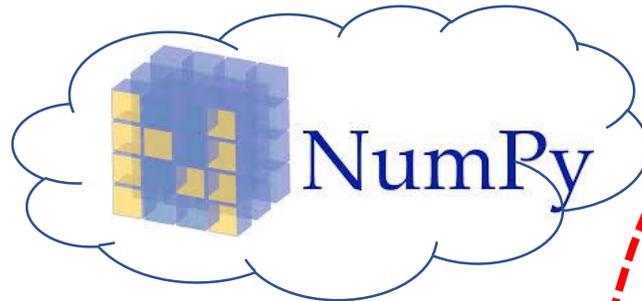
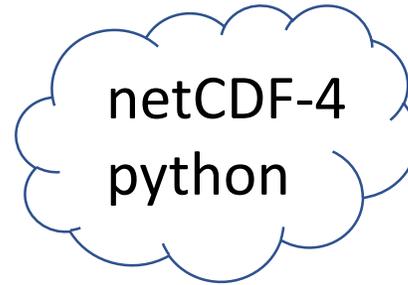
Data providers (and AODN staff) are notified via a detailed email report.



Benefits: Python ecosystem python™



Benefits: Python ecosystem – One outlier



Benefits: Test and deploy code in confidence

The screenshot shows the GitHub repository page for `aodn/python-aodncore`. The browser address bar displays the URL `https://github.com/aodn/python-aodncore`. The repository name is `aodn / python-aodncore`, with 15 watchers, 0 stars, and 0 forks. The repository is categorized as `Code` and has 13 issues, 5 pull requests, insights, and settings. The repository description is "AODN pipeline core library" with a link to `https://aodn.github.io/python-aodncore/`. The repository statistics include 396 commits, 8 branches, 41 releases, 1 environment, 9 contributors, and the GPL-3.0 license. The repository is currently on the `master` branch. The commit history shows the latest commit `a64d1ad` from `aodn-ci` with the message "Jenkins version bump (0.21.2)" 6 days ago. The commit history also includes updates to `.idea`, `sphinx`, `test_aodncore`, `.coveragerc`, `.gitignore`, and `travis.yml`.

GitHub, Inc. (US) | <https://github.com/aodn/python-aodncore>

Search or jump to... Pull requests Issues Marketplace Explore

aodn / python-aodncore Unwatch 15 Star 0 Fork 0

Code Issues 13 Pull requests 5 Insights Settings

AODN pipeline core library <https://aodn.github.io/python-aodncore/> Edit

Manage topics

396 commits 8 branches 41 releases 1 environment 9 contributors GPL-3.0

Branch: master New pull request Create new file Upload files Find file Clone or download

aodn-ci Jenkins version bump (0.21.2) Latest commit a64d1ad 6 days ago

.idea	Update HTML docs	7 months ago
aodncore	Jenkins version bump (0.21.2)	6 days ago
sphinx	Update README.md	7 days ago
test_aodncore	Remove cc-plugin-imos dependency	a month ago
.coveragerc	Update .coveragerc	2 months ago
.gitignore	Update doc generation to push to 'gh-pages' branch, to keep autogener...	5 months ago
travis.yml	Update travis.yml	6 days ago

Benefits: Test and deploy code in confidence

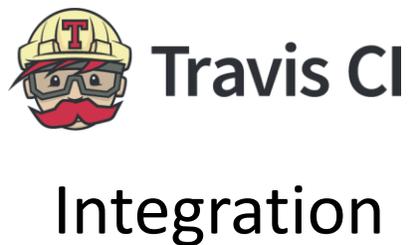
 COPYING	Add GPLv3 licence	8 months ago
 README.md	Add travis codecov integration	3 months ago
 _config.yml	Set theme jekyll-theme-minimal	9 months ago
 aodncore_demo.ipynb	Update example notebook testlib path	a year ago
 requirements.txt	Update setup.py/requirements.txt relationship as described at: https:...	3 months ago
 setup.py	update required compliance-checker version to 4.1.1	6 days ago
 state_machine.png	Initial commit of aodncore package	a year ago

 [README.md](#) 

python-aodncore

build passing | codecov 89%

Project documentation is hosted at: <https://aodn.github.io/python-aodncore/index.html>



Benefits: Test and deploy code in confidence



Travis CI

627 Tests
passed !!!

The screenshot shows the Travis CI web interface for a repository named 'aodn / python-aodncore'. The page displays a 'build passing' status. A pull request #135 is highlighted, with a green checkmark indicating it passed. The pull request details include the commit hash '434c330', the branch 'master', and the author 'Leigh Gordon'. The build job #627.1 is shown to have passed, running for 2 minutes and 46 seconds, and completed about 11 hours ago. The environment is specified as 'Python: 2.7'. At the bottom, a 'Job log' section is visible, showing the start of the build process with 'Worker information' and 'Build system information'.

Benefits:

- Generic code
- Increased confidence:
 - All features can be tested from one version to the next
- Easier to debug
- Better control of the different version
 - Workflow to build, package and deploy
 - What version is deployed on which environment (Release candidate or Production)
- Faster publication of new data
- Improved consistency within data collections

Future improvements

- Reporting, creation of dashboard in Sumologic
- Improve code:
 - Replace Talend (Java) by a similar tool in Python for a consistent environment
- Use of AWS Batch to improve scalability:
 - No conflict between pipelines
 - Multiple queues for different data streams
 - Only running when needed

GitHub

<https://github.com/aodn/python-aodncore>

<https://github.com/aodn/python-aodndata>



IMOS is a national collaborative research infrastructure, supported by Australian Government. It is operated by a consortium of institutions as an unincorporated joint venture, with the University of Tasmania as Lead Agent. www.imos.org.au

PRINCIPAL PARTICIPANTS



SIMS is a partnership involving four Universities.

ASSOCIATE PARTICIPANTS

