

---

# Research Data Management for Computational Science

---

Christian T. Jacobs<sup>1</sup>

c.jacobs10@imperial.ac.uk

www.christianjacobs.uk

@ctjacobs\_uk

&

Alexandros Avdis<sup>1</sup>, Simon L. Mouradian<sup>1</sup>,  
Gerard J. Gorman<sup>1</sup>, Matthew D. Piggott<sup>1</sup>

<sup>1</sup>Department of Earth Science and Engineering, Imperial College London

The Data Hide, ODSI, University of Sheffield

20 October 2015

# Ocean Simulations

- ▶ Simulations of ocean dynamics are important in many applications.
  - ▶ Prediction of tsunami impacts

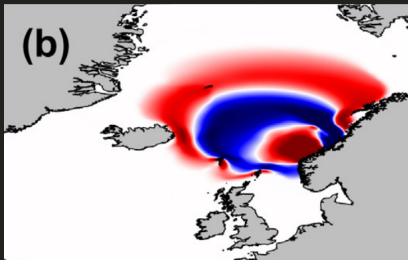


Image by Hill et al. (2014), used under CC-BY, doi:10.1016/j.ocemod.2014.08.007

- ▶ Optimisation of marine renewable energy turbines
- ▶ Estimating the range of nuclear contaminants

# Software and Data Requirements

- ▶ Simulations should be **recomputable** and **reproducible**.
- ▶ This requires:
  - ▶ the **software** itself (with info about the specific version used)
  - ▶ **raw data** (input and output files)
  - ▶ **provenance metadata**

## Problem

Unfortunately, most simulation-based publications are **not accompanied** by the data and the software (and exact version info) needed to recreate it.

# What Can Be Done?

- ▶ The level of motivation amongst researchers to share their data and software is generally quite low.
  - ▶ Extra **effort** and **time** required to gather and publish it.
  - ▶ Typically **gain little** from the process.
  - ▶ See LeVeque et al. (2012)<sup>1</sup>

## What we need

- ▶ We need a way of **publishing** data and software that is **quick and easy...**
- ▶ ...and a way of **referencing** it correctly in papers.

---

<sup>1</sup> LeVeque, R.J., Mitchell, I.M., Stodden, V. (2012). Reproducible Research for Scientific Computing: Tools and Strategies for Changing the Culture. Computing in Science & Engineering 14(4), 13--17.

# "Green Shoots Project": PyRDM

- ▶ **PyRDM**: **R**esearch **D**ata **M**anagement with **Py**thon
- ▶ Open-source, GNU GPL. [github.com/pyrdm/pyrdm](https://github.com/pyrdm/pyrdm)

- ▶ Facilitates the **automated publication** of source code and data to:

- ▶ Figshare ([figshare.com](https://figshare.com))
- ▶ Zenodo ([zenodo.org](https://zenodo.org))
- ▶ DSpace-based repositories ([dspace.org](https://dspace.org))



Jacobs et al. (2014),  
DOI: [10.5334/jors.bj](https://doi.org/10.5334/jors.bj)

- ▶ Online, citable and persistent repositories. Each code/dataset is given its own **DOI**.

# Publishing Process: Software Source Code

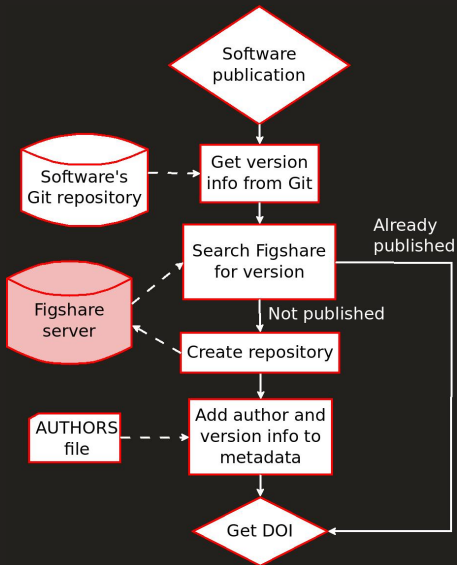


Image adapted from Jacobs et al. (2015).

# Application to Ocean Simulations

---

- ▶ A prerequisite to a reproducible simulation is the availability and reproducibility of the mesh.
- ▶ Applied PyRDM to QMesh, a tool for generating meshes from GIS data (Avdis et al., in preparation).
  - ▶ See Jacobs et al. (2015) for details about RDM implementation.

# Ocean simulations: The Mesh

- ▶ A key simulation input is the **mesh**.
  - ▶ Area of interest represented by discrete points/cells.

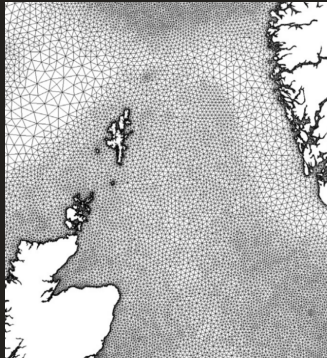


Image by Hill et al. (2014), used under CC-BY, doi:10.1016/j.ocemod.2014.08.007

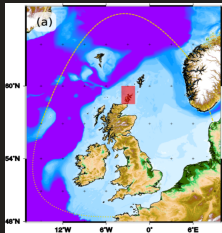
- ▶ ...but creating a realistic, high-resolution mesh by hand is **infeasible**.



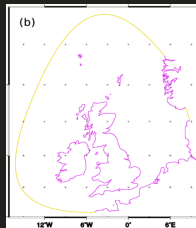
# Geographical Information Systems

- ▶ Geographical Information Systems are good at processing bathymetry and coastline data to create a realistic geometry.
  - ▶ e.g. QGIS, ArcGIS, ...

Bathymetry data



Geometry

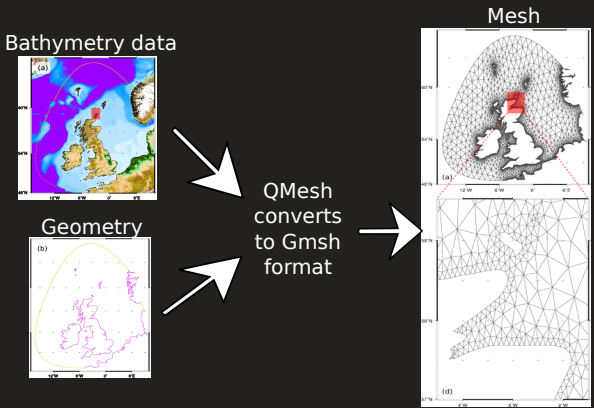


Images by Avdis et al. (2015).

- ▶ How do we create a mesh based on this input data?

# QMesh: Mesh Production using GIS Data

- ▶ **QMesh** is a software package which:
  - ▶ Takes the geometry defined in **QGIS**...
  - ▶ ...and **converts** the geometry into an appropriate format for...
  - ▶ ...**Gmsh**, a tool which generates the mesh for the domain.



## Example Workflow: Orkney and Shetland Isles

---

- ▶ Consider the area around the Orkney and Shetland Isles.
- ▶ Involves a number of GIS input data files:
  - ▶ The QGIS project file itself, comprising:
  - ▶ Geometrical layer files defining the coastlines
  - ▶ Bathymetry data in a NetCDF file

# Example Workflow: Geometry in QGIS

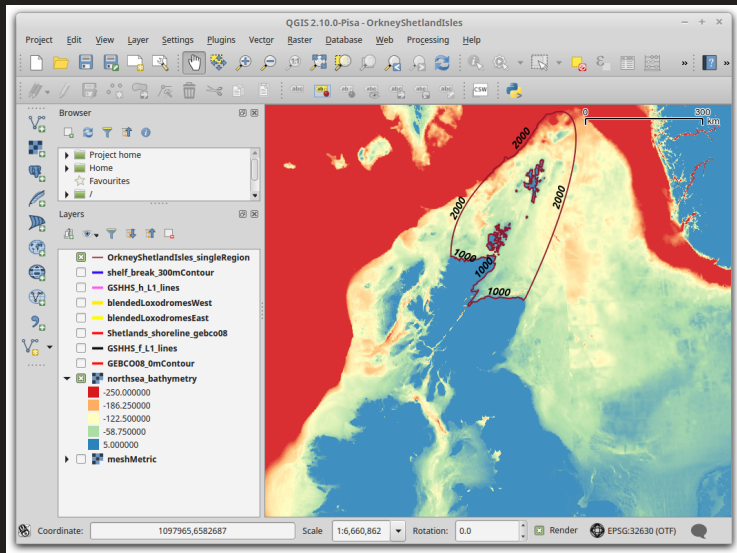


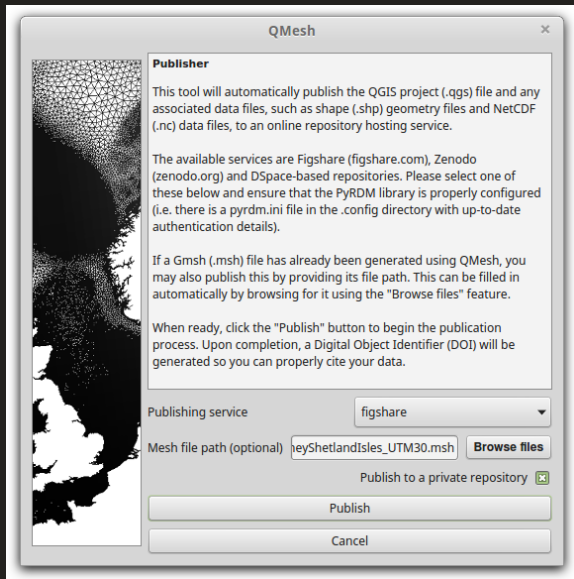
Image by Jacobs et al. (2015).

## Example Workflow: Mesh from QMesh

---

- ▶ The input data in the QGIS project is used to produce a mesh using QMesh.
- ▶ User runs their ocean simulation using this mesh.
- ▶ When results are satisfactory, user publishes the data and software using the QMesh publishing tool.

# Example Workflow: QMesh Publishing Tool



# Publishing Process: Data

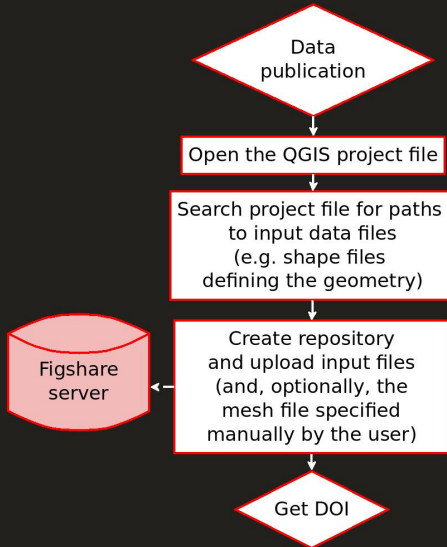


Image adapted from Jacobs et al. (2015).

# Example Workflow: QGIS project file

- ▶ Publishing tool **parses** the XML-based QGIS project file to determine **location** of all **data files** that the project comprises...

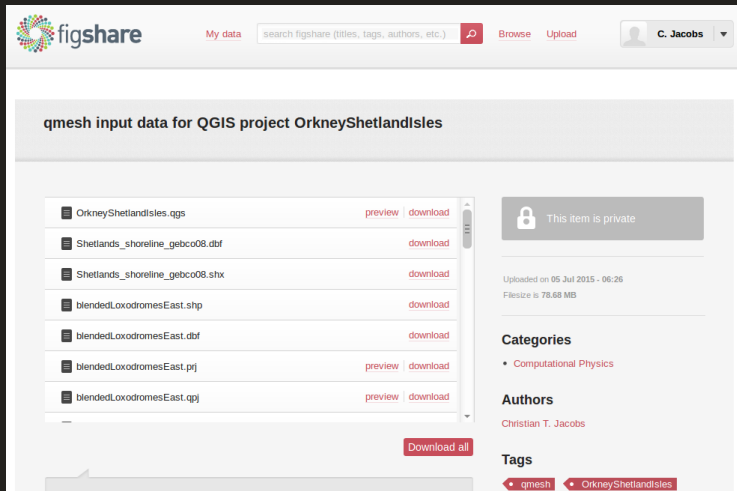
```
</edittypes>
</maplayer>
<maplayer minimumScale="0" maximumScale="1e+08" simplifyDrawingHints="1" minLabelScale="0" maxLabelScale="1" scaleBasedLabelVisibilityFlag="0">
  <id>Shetlands_shoreline_gebco0820140707115902174</id>
  <datasource>./Shetlands_shoreline_gebco08.shp</datasource>
  <title></title>
  <abstract></abstract>
  <keywordList>
    <value></value>
  </keywordList>
  <layername>Shetlands_shoreline_gebco08</layername>
  <srs>
    <spatialrefsys>
      <proj4>+proj=longlat +datum=WGS84 +no_defs</proj4>
      <srsid>3452</srsid>
      <srid>4326</srid>
      <authid>EPSG:4326</authid>
      <description>WGS 84</description>
      <projectionacronym>longlat</projectionacronym>
      <ellipsoidacronym>WGS84</ellipsoidacronym>
      <geographicflag>true</geographicflag>
    </spatialrefsys>
  </srs>

```



# Example Workflow: Files on Figshare

- ...and uploads these files to the repository hosting service via its API.



The screenshot shows the Figshare website interface. At the top, there is a navigation bar with the Figshare logo, a search bar containing the text "search figshare (titles, tags, authors, etc.)", and buttons for "My data", "Browse", and "Upload". A user profile for "C. Jacobs" is visible in the top right corner.

The main content area displays the title "qmesh input data for QGIS project OrkneyShetlandIsles". Below the title is a list of files with their respective download and preview options:

| File Name                       | Actions            |
|---------------------------------|--------------------|
| OrkneyShetlandIsles.qgs         | preview   download |
| Shetlands_shoreline_gebco08.dbf | download           |
| Shetlands_shoreline_gebco08.shx | download           |
| blendedLoxodromesEast.shp       | download           |
| blendedLoxodromesEast.dbf       | download           |
| blendedLoxodromesEast.prj       | preview   download |
| blendedLoxodromesEast.qpj       | preview   download |

Below the file list is a "Download all" button. To the right of the file list, a grey box with a lock icon indicates "This item is private". Below this, it states "Uploaded on 05 Jul 2015 - 06:26" and "Filesize is 78.68 MB".

The right sidebar contains sections for "Categories" (with a sub-item "Computational Physics"), "Authors" (listing "Christian T. Jacobs"), and "Tags" (with sub-items "qmesh" and "OrkneyShetlandIsles").

# Example Workflow: DOI

Publication ID and DOI are assigned, and presented to user once publication process is complete:

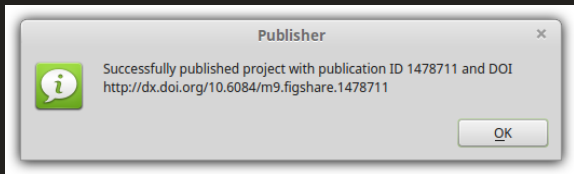


Image by Jacobs et al. (2015).

## Issues/Limitations Encountered

---

- ▶ Lack of **standardisation**. Need a better way of affiliating authors.
- ▶ Lack of **API support**. No searching in Zenodo, no server-side MD5 checksums in Figshare, ...
- ▶ Restriction on **private storage** space.
- ▶ Restriction on **number of collaborators**.
- ▶ Figshare for Institutions / cloud storage to address these restrictions?
- ▶ Publishing QMesh source code may not be enough to reproduce the exact same mesh without **knowledge of its dependencies**.

# References and Acknowledgements

---

- ▶ [Jacobs et al. \(2014\)](#). PyRDM: A Python-based library for automating the management and online publication of scientific software and data. *Journal of Open Research Software*, 2(1):e28. DOI: 10.5334/jors.bj
- ▶ [Avdis et al. \(2015\)](#). Shoreline and Bathymetry Approximation in Mesh Generation for Tidal Renewable Simulations. In *Proceedings of the European Wave and Tidal Energy Conference (EWTEC) Series*. Pre-print: <http://arxiv.org/abs/1510.01560>
- ▶ [Avdis et al. \(In Preparation\)](#). Efficient unstructured mesh generation for renewable tidal energy using Geographical Information Systems.
- ▶ [Jacobs et al. \(2015\)](#). Integrating Research Data Management into Geographical Information Systems. In *Proceedings of the 5th International Workshop on Semantic Digital Archives*. Pre-print: <http://arxiv.org/abs/1509.04729>
- ▶ Thanks to the Research Office at Imperial College London for funding.
- ▶ Slides produced using  $\text{\LaTeX}$ , with a modified version of the Wronki Beamer theme ([kaszkowiak.eu](http://kaszkowiak.eu)).

---

# Research Data Management for Computational Science

---

Christian T. Jacobs<sup>1</sup>

c.jacobs10@imperial.ac.uk

www.christianjacobs.uk

@ctjacobs\_uk

&

Alexandros Avdis<sup>1</sup>, Simon L. Mouradian<sup>1</sup>,  
Gerard J. Gorman<sup>1</sup>, Matthew D. Piggott<sup>1</sup>

<sup>1</sup>Department of Earth Science and Engineering, Imperial College London

The Data Hide, ODSI, University of Sheffield

20 October 2015