# Speeding-up data analysis: DIVAnd interpolation tool in the Virtual Research Environment

C. Troupin, A. Barth, M. Buurman, S. Mieruch,
L. Bruvry Lagadec, T. Zamani & P. Thijsse

# A few definitions before we start

DIVA: software tool designed for spatial interpolation
 https://github.com/gher-ulg/DIVA

DIVAnd: n-dimensional version of DIVA
 https://github.com/gher-ulg/divand.jl

Julia: high-level, high-performance dynamic programming language for numerical computing
 https://julialang.org/

Notebook: documents that contain live code, equations, visualizations and narrative text

Jupyter: open-source web application to create and share notebook
 http://jupyter.org/

Jupyterhub: multi-user server for Jupyter notebooks
 https://jupyterhub.readthedocs.io/en/latest/

# Virtual Research Environment (VRE)

Definition: a VRE is made up of software tools, datasets and other resources that several users can access and use in a cooperative way.

Advantages:

- Availability of tools – no need to install anything
- Availability of data – no download needed
- Availability of computing power

# 1
# Context for DIVAnd in the VRE

Back in 2007

or in 2013

## What have we improved since then?

1. New mathematical formulation  📄 Barth et al. 2014
2. Julia language  instead of Fortran
3. Only 2 (!!!) input files  data & bathymetry
4. Applications as Jupyter notebooks  all in one

# Jupyter notebooks as guidelines

# Jupyter notebooks as guidelines

# Jupyter notebooks as guidelines

# Jupyter notebooks as guidelines

# Jupyter notebooks as guidelines

# What should we improve?

+ Access to computing power
+ Data availability
+ Documentation

# What should we improve?

+ Access to computing power
+ Data availability
+ Documentation

_____

= Virtual Research Environment!

# DIVAnd
# in a few
# clicks

# DIVAnd in a few clicks

# DIVAnd in a few clicks

## Work directory

# DIVAnd in a few clicks

nextcloud_sync allows users to access their private workspace

# DIVAnd in a few clicks

directory where webODV gave me the data files

# DIVAnd in a few clicks

sub-directory where webODV gave me the data files

# DIVAnd in a few clicks

data file to use in the notebooks

# DIVAnd in a few clicks

DIVAnd-Workshop =  repository cloned into the VRE

# DIVAnd in a few clicks

Notebooks divided into subdirectories

# DIVAnd in a few clicks

Full-analysis using the extracted data file

# Jupyter notebooks to produce climatologies

# DIVAnd outputs

1. A netCDF file storing the analysis fields and the observation ID's
2. A XML file that can be used for the Sextant catalog
3. Figures (here Autumn salinity at 5 m depth)

# Improving DIVAnd in the VRE

# DIVAnd in 2020

1. Updated mathematical formulation      📄 Barth et al. 2014
2. Julia language
3. Only 2 input files      data & bathymetry
4. User guides = Jupyter notebooks      all in one

# Providing users access to notebooks



2017 | 1st deployment at **CINECA**

2019 | Deployment transferred to **DKRZ**
for the first training workshop

2020 | Multiple machines at
**DKRZ**
**GRNET**
**STFC**

Multiple copies of a Docker container are run
(https://hub.docker.com/r/abarth/divand-jupyterhub)

# What's shipped inside the container?

☑ Libraries: netCDF, unzip, git, ...

🔗 https://github.com/gher-ulg/DIVAnd-jupyterhub

☑ Libraries: netCDF, unzip, git, …
☑ Julia language          (V1.4.1, released on April 14, 2020))

# What's shipped inside the container?

☑ Libraries: netCDF, unzip, git, …

☑ Julia language         (V1.4.1, released on April 14, 2020))

☑ Julia packages: PyPlot, NCDatasets, DataStructures, …

# What's shipped inside the container?



🔗 https://github.com/gher-ulg/DIVAnd-jupyterhub

☑ Libraries: netCDF, unzip, git, …
☑ Julia language          (V1.4.1, released on April 14, 2020))
☑ Julia packages: PyPlot, NCDatasets, DataStructures, …
☑ DIVAnd.jl                          (#master)

# What's shipped inside the container?



🔗 https://github.com/gher-ulg/DIVAnd-jupyterhub

- ☑ Libraries: netCDF, unzip, git, ...
- ☑ Julia language            (V1.4.1, released on April 14, 2020))
- ☑ Julia packages: PyPlot, NCDatasets, DataStructures, ...
- ☑ DIVAnd.jl            (#master)
- ☑ DIVAnd *workshop* notebooks       (#master)

1. Julia language

1. Julia language
2. DIVAnd package

1. Julia language
2. DIVAnd package
3. Training notebooks (Diva-Workshop)

1. Julia language
2. DIVAnd package
3. Training notebooks (Diva-Workshop)
4. Docker container in the VRE

v1.4.1 (April 14, 2020)

New features: changes in language features, multi-threading, build systems, library functions, ...

PackageCompiler.jl : allows to pre-compile the modules
Reduce the famous "*time-to-first-plot*":
when you first load and run a package in a session,
Julia needs to compile it first



**Roadmap for a faster time-to-first-plot?**
■ Internals & Design

dlfivefifty                                                    Apr '19

Having unintentionally "thrown shade" and performed a "kvetchfest" in another thread, let me first apologise. I'm very happy with the improvements made in Julia, especially the interpreter and debugger, and am very much excited for future developments.

That said, perhaps it would help ease frustration to have a "roadmap" for when the time-to-first-plot issues are planned to be addressed. It actually seems like it's already a solved problem with PackageCompiler.jl, if that were cleaned up in an easy-to-use framework.

4 ♡  &

# DIVAnd.jl package

## v2.6.0 (April 27, 2020)

New feature: heatmap generation

Closed issues:
- DIVAnd_qc doesn't work with advection constrain
- Missing surfextend documentation
- Wrong long name for lat
- netCDF scalefactore is NaN resulting in an empty output variable

# Example of heatmaps

Data points are locations acquired by GPS device
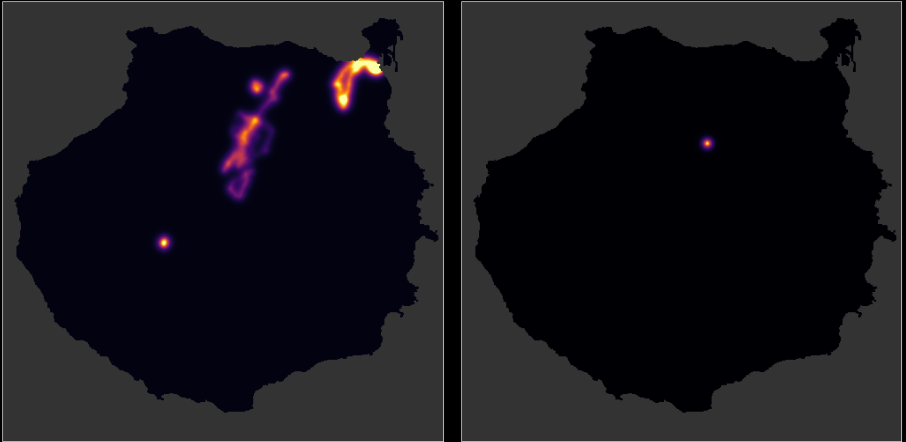during running activities



Figure 1: Left: heatmap before the confinement; Right: after the confinement

# Training notebooks (Diva-Workshop)

### v1.2 (April 27, 2020)

- Compatibility with Julia 1.4
- CI (Travis)
- Cleaned notebooks
- Sub-directories with notebooks arranged by topics.
- Minor bugs

# Docker container (DIVAnd-jupyterhub)

## Precompiled DIVAnd with PackageCompiler

1. DIVAnd is deployed in the SeaDataCloud VRE
2. Continuous improvements in:
   Julia, DIVAnd, Training notebooks and in the Docker
3. Making available data, tools and computing power is essential to speed up research

# Thanks
# for your attention

🐦 CharlesTroupin

#shareEGU20