



1 **eddy4R: A community-extensible processing, analysis and**
2 **modeling framework for eddy-covariance data based on R,**
3 **Git, Docker and HDF5**

4
5 **Stefan Metzger¹, David Durden¹, Cove Sturtevant¹, Hongyan Luo¹, Natchaya**
6 **Pingintha-Durden¹, Torsten Sachs², Andrei Serafimovich², Jörg Hartmann³,**
7 **Jiahong Li⁴, Ke Xu⁵, Ankur R. Desai⁵**

8 ¹Battelle Ecology, 1685 38th Street, Boulder, CO 80301, USA

9 ²GFZ German Research Centre for Geosciences, Telegrafenberg, 14473 Potsdam, Germany

10 ³Alfred Wegener Institute – Helmholtz Centre for Polar and Marine Research, Am
11 Handelshafen 12, 27570 Bremerhaven, Germany

12 ⁴LI-COR Biosciences, 4647 Superior Street, Lincoln, NE 68504, USA

13 ⁵University of Wisconsin-Madison, Dept. of Atmospheric and Oceanic Sciences, 1225 West
14 Dayton Street, Madison, WI 53706, USA

15 Correspondence to: Stefan Metzger (smetzger@battelleecology.org)

16

17 **Keywords:** computing, container, continuous development, continuous integration, devOps,
18 eddy4R, eddy-covariance, image, NEON, reproducibility, science code

19

20 **Abstract**

21 This study presents the systematic development of an open-source, flexible and modular eddy-
22 covariance (EC) data processing framework. This is achieved through adopting a Development
23 and Systems Operation (DevOps) philosophy, building on the eddy4R family of EC code
24 packages in the R Language for Statistical Computing as foundation. These packages are
25 community-developed via the GitHub distributed version control system and wrapped into a
26 portable and reproducible Docker filesystem that is independent of the underlying host
27 operating system. The HDF5 hierarchical data format then provides a streamlined
28 mechanism for highly compressed and fully self-documented data ingest and output.

29 This framework is applicable beyond EC, and more generally builds the capacity to deploy
30 complex algorithms developed by scientists in an efficient and scalable manner. In addition,
31 modularity permits meeting project milestones while retaining extensibility with time.

32 The efficiency and consistency of this framework is demonstrated in the form of three
33 application examples. These include tower EC data from first instruments installed at a
34 National Ecological Observatory (NEON) field site, aircraft flux measurements in combination
35 with remote sensing data, as well as a software intercomparison. In conjunction with this
36 study, the first two eddy4R packages and simple NEON EC data products are released publicly.
37 While this proof-of-concept represents a significant advance, substantial work remains to arrive
38 at the automated framework needed for the streaming generation of science-grade EC fluxes.



39 1 Introduction

40 Answering grand challenges in earth system science and ecology requires combining
41 information from hierarchies of environmental observations (tower, aircraft, satellite; Raupach
42 et al., 2005; Running et al., 1999; Turner et al., 2004). Eddy-covariance (EC) measurements
43 serve as crucial observations in this hierarchy to study landscape-scale surface-atmosphere
44 exchange processes that both inform and anchor earth system models. Networks of EC towers
45 such as FLUXNET (Baldocchi et al., 2001), AmeriFlux (Law, 2007), ICOS (Sulkava et al.,
46 2011), and others are vital for providing the necessary distributed observations covering the
47 climate space, with the longest running towers now reaching two decades of observations.

48 A current challenge for EC tower networks in informing regional and continental scale
49 processes is instrument and computational compatibility. Much progress has been made in
50 developing community standards for processing algorithms and workflows (Aubinet et al., 2012;
51 Papale et al., 2006). However, the computations involved in EC processing are complex and
52 developmentally dynamic, making code portability, extensibility, and documentation
53 paramount. Many authors have included code in publication, or have developed sharable tools
54 (e.g. EddyPro and TK3 by Fratini and Mauder (2014), EddyUH by Mammarella et al. (2016),
55 EdiRe by Clement et al. (2009). Still, large differences in instrumentation, site setup, data
56 format, and operating system stymie the adoption of a universal EC processing environment,
57 exacerbated by the significant and often unfunded effort required to adequately document and
58 generalize code. In 50% of published scientific code, one cannot even replicate the necessary
59 software dependencies (Collberg et al., 2014) let alone develop tailored workflows to
60 incorporate additional data streams, automate and scale processing across large compute
61 facilities, or inject additional algorithms to address specific needs or synergistic research
62 questions.

63 The National Ecological Observatory Network (NEON), once fully operational, will represent
64 the largest single-provider EC tower network globally, with a standardized measurement suite
65 designed explicitly for cross-site comparability and analysis of continental-scale ecological
66 change (Schimel et al., 2007). This capability is accompanied by a strong need for a flexible
67 and scalable processing framework that can incorporate specific data streams, take advantage
68 of tight hardware-software integration for problem tracking and resolution, provide traceability
69 and reproducibility of outputs, and seamlessly integrate distributed and dynamic community-
70 developed code within existing cyberinfrastructure.

71 Here, we describe and demonstrate a framework that enables these capabilities by embracing a
72 Development and Systems Operation (DevOps) approach. DevOps is a philosophy arising from
73 within the software development community that emphasizes collaboration among developers
74 and operators to continuously iterate the development, building, testing, packaging, and release
75 of software (Erich et al., 2014; Loukides, 2012). Tools are adopted that control and automate
76 these processes, allowing distributed development and rapid iteration. A key aspect of DevOps
77 that can aid improving accessibility, extensibility, and reproducibility of scientific software is
78 through recipe- or script-based generation and packaging of computation environments rather
79 than abstracted documentation (Boettiger, 2015; Clark et al., 2014). The recipe automates the



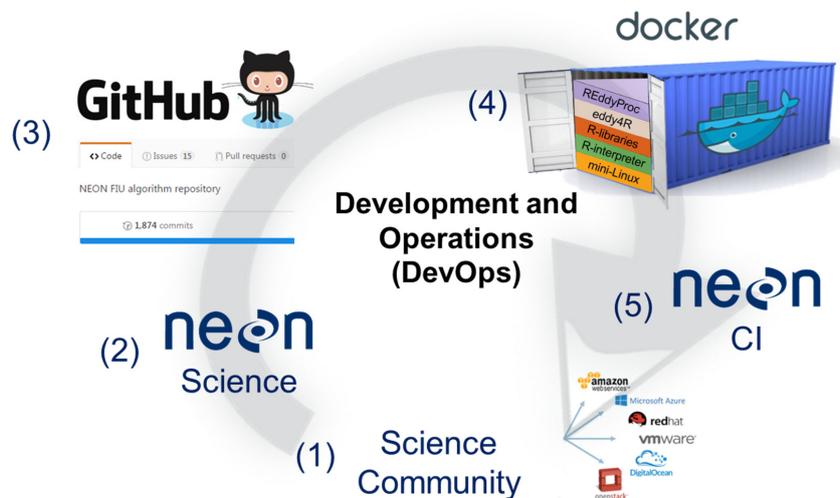
80 loading of the software including all dependencies so that the most significant hurdle of
81 reproducing the computational environment is overcome. At the same time, the recipe serves
82 as explicit documentation, and can be easily extended (added to or changed), shared, and
83 versioned. The entire computational environment including any necessary data are packaged
84 into Docker images that work identically across different computers and operating systems, can
85 be deployed at scale, and archived for ultimate reproducibility.

86 In the following we present this framework and demonstrate its success in producing EC data
87 products via a family of modular, open-source R packages wrapped in Docker images. We
88 emphasize that this paper is not a presentation of EC processing software (although this is the
89 ultimate application). Rather, it is a presentation of the developmental framework that facilitates
90 scalability, portability, and extensibility of EC processing software. Section 2 describes the
91 DevOps framework, Sect. 3 provides tower and aircraft example applications including
92 NEON's first set of EC data, as well as a software cross-validation. Section 4 summarizes the
93 work remaining to operationally produce EC fluxes from 47 NEON sites, and provides an
94 outlook on future capabilities. Code and data availability information is provided in Sect. 5.

95 2 The development and operations framework

96 NEON's DevOps framework consists of a periodic sequence (Figure 1): The science
97 community contributes algorithms and best practices (1) which together with NEON Science
98 (2) are compiled into eddy4R packages via the GitHub distributed version control system (3).
99 NEON Science releases an eddy4R version from GitHub, which automatically builds an
100 eddy4R-Docker image on DockerHub as specified in a "Dockerfile" (4). The eddy4R-Docker

101



102

103 Figure 1. DevOps workflow of the eddy4R-Docker image. Please see text in Sect. 2 for detailed
104 explanation.



105 image is immediately available for deployment by NEON Cyberinfrastructure (CI; 5), the
106 Science Community (1) and NEON Science (2) alike. This DevOps cycle can be repeated for
107 continuous development and integration of requests and future methodological improvements,
108 resulting in the next release. Two principal types of releases are provided: stable versions are
109 tagged with “v1.0.0”, “v1.0.1” etc., and the most recent development built is tagged with
110 “latest”.

111 In the following we describe the key infrastructure components of this DevOps framework,
112 namely the eddy4R family of code packages (Sect. 2.1), Git-based distributed code
113 development (Sect. 2.2), packaging of the computational environment in Docker images
114 (Sect. 2.3), hierarchical data formats (Sect. 2.4), integration with NEON’s cyberinfrastructure
115 (Sect. 2.5), and installation and deployment (Sect. 2.6).

116 2.1 The eddy4R family of R-packages

117 eddy4R is a family of open-source packages for EC raw data processing, analyses and modeling
118 in the R Language for Statistical Computing (R Core Team, 2016). It is being developed by
119 NEON scientists with wide input from the micrometeorological community (e.g., De Roo et al.,
120 2014; Kohnert et al., 2015; Lee et al., 2015; Metzger et al., 2012; Metzger et al., 2013; Metzger
121 et al., 2016; Sachs et al., 2014; Salmon et al., 2015; Serafimovich et al., 2013; Starkenburg et
122 al., 2016; Vaughan et al., 2015; Xu et al., 2017). eddy4R currently consists of four packages
123 eddy4R.base, eddy4R.qaqc, eddy4R.turb, and eddy4R.erf. Of these, eddy4R.base and
124 eddy4R.qaqc are published here in conjunction with NEON’s release of EC Level 1 data
125 products: descriptive statistics of calibrated instrument output. In addition, previews of
126 eddy4R.turb and eddy4R.erf are provided, which will be published along NEON’s release of
127 EC Level 4 data products (derived quality-controlled fluxes and related variables).
128 Development of two additional R-packages eddy4R.stor and eddy4R.ucrt has started, which
129 provide functionalities for storage flux computation and uncertainty quantification, respectively.
130 These packages are not covered here, and will be published once available.

131 Each eddy4R package consists of a hierarchical set of reusable definition functions, wrapper
132 functions and workflow templates. Following best practices, eddy4R is written in controlled
133 and strictly hierarchical terminology consisting of base names and modifiers, which ensures
134 modular extensibility over time. Interactive documentation is provided through the use of
135 Roxygen (<http://roxygen.org/>) tags during development, and follows the Comprehensive R
136 Archive Network (CRAN; <https://cran.r-project.org/>) guidelines for package dissemination. In
137 addition, expanded documentation is available in the form of Algorithm Theoretical Basis
138 Documents from the NEON data portal (data.neonscience.org/home).

139 eddy4R.base provides natural constants and basic functions for usability, regularization,
140 transformation, lag-correction, aggregation and unit conversion ensuring consistency of internal
141 units at any point in the workflow. Next, eddy4R.qaqc provides the general quality assurance
142 and quality control (QA/QC) tests of Taylor and Loeschner (2013), along the Smith et al. (2014)
143 framework for tracking quality information in large datasets, and functions for de-spiking
144 (Brock, 1986; Fratini and Mauder, 2014; Mauder and Foken, 2011; Mauder et al., 2013;



145 Metzger et al., 2012; Vickers and Mahrt, 1997). eddy4R.turb provides standard, Reynolds-
 146 decomposed turbulent flux calculation (Foken, 2008), accompanied by facilities for planar fit
 147 transformation (Wilczak et al., 2001) and spectral correction (Nordbo and Katul, 2012).
 148 Additional functionalities include Fourier transform, the determination of detection limit
 149 (Billesbach, 2011), integral length scales and statistical sampling errors (Lenschow et al., 1994),
 150 and flux-specific QA/QC (Foken and Wichura, 1996; Vickers and Mahrt, 1997). Also, basic
 151 scaling variables, atmospheric stability and roughness length (Stull, 1988), as well as the flux
 152 footprint (Kljun et al., 2015; Kormann and Meixner, 2001; Metzger et al., 2012) can be
 153 determined. Lastly, edd4R.erf provides time-frequency de-composed flux processing and
 154 artificially intelligent functionality to determine environmental response functions and project
 155 the flux fields underlying the EC observations (Metzger et al., 2013; Xu et al., 2017).

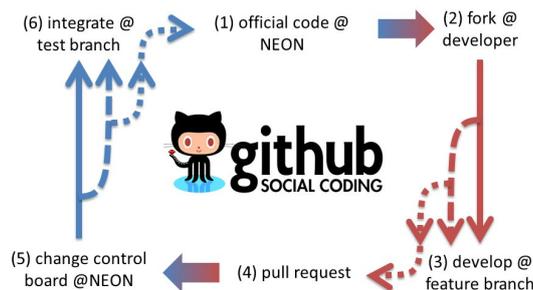
156 eddy4R can be used with a fully adaptive single-pass workflow template (Sect. 3.1), which
 157 makes it computationally efficient compared to the multiple passes required by other flux
 158 processing schemes. In addition, eddy4R is fully parallelized and memory efficient leveraging
 159 R's snowfall (<https://cran.r-project.org/package=snowfall>) and ff (<https://cran.r-project.org/package=ff>)
 160 facilities, respectively. This makes eddy4R seamlessly scalable from
 161 local laptop development to deployment across massively parallel computing facilities. Lastly,
 162 its unique modularity permits straightforward adjustments and versioning as science and/or
 163 hardware progresses.

164 2.2 Git distributed version control

165 The eddy4R source code resides on a version-controlled Git repository on the hosting service
 166 GitHub (<https://github.com/>). In general, a developer community uses a version control system
 167 to manage and track different states of their works over time. GitHub provides distributed
 168 version control and has become widely used by scientific research groups because it is free,
 169 open-source, and provides several features that make it useful for managing artifacts of
 170 scientific research (Ram, 2013).

171 Git allows multiple users and developers to simultaneously access and collaborate on a remote
 172 repository by means of independent 'forks' or replicas of the entire repository

173



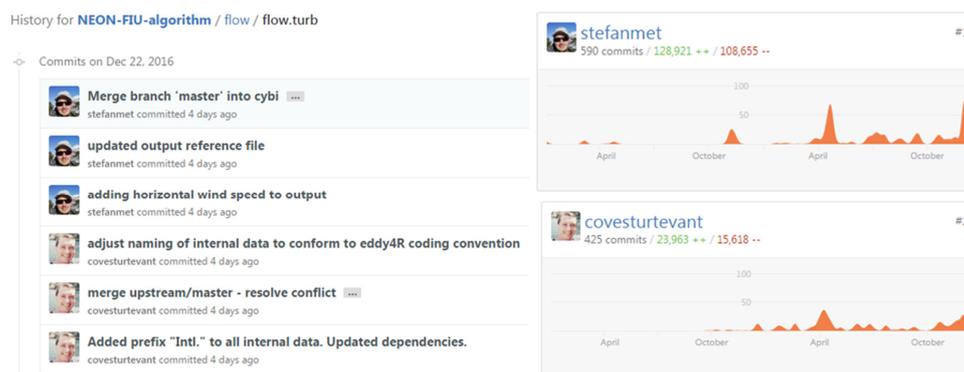
174

175 Figure 2. NEON's Git workflow. Please see text in Sect. 2.2 for detailed explanation.



176 (Paarsch and Golyaev, 2016). Figure 2 shows NEON's Git workflow: At any given time (1) the
177 official, stable eddy4R source code resides on NEON's GitHub repository. A user can install
178 the eddy4R packages directly from there, and (2) a developer can 'fork' or copy the repository
179 and create 'branches' for modification. After (3) 'committing' or creating a new feature, the
180 developer (4) can propose the feature for inclusion in the official eddy4R source code by issuing
181 a pull request to (5) NEON's change control board. After (6) thorough review and testing, the
182 feature can be 'merged' or integrated into the next release of (1) the official, stable eddy4R
183 source code. This cycle can be repeated to accommodate requests and future developments,
184 resulting in subsequent releases. The developers can periodically update their 'forks' from the
185 remote repository, ensuring that they always work on basis of the most current eddy4R source
186 code.

187



188

189 Figure 3. Example of GitHub facilities for exploring code authorship, development history.

190

191 The ultimate advantages of Git are provenance, reproducibility and extensibility (Figure 3):
192 Every copy of the code repository includes the complete history of all changes, and authorship
193 that can be viewed and searched by anyone (Ram, 2013). This allows developers to build from
194 any stage of the versioned project and makes it easy to collaborate as an integrated the scientific
195 community.

196 2.3 Docker image build and deployment

197 Docker images (<https://www.docker.com/what-docker>) wrap a piece of software in a complete
198 filesystem that contains only the minimal context an application needs to run: code, runtime,
199 system tools, system libraries. This guarantees that it always performs the same, regardless of
200 the compute environment it is deployed in. By running as native processes, Docker bypasses
201 the overhead encountered in the similar but more cumbersome virtual machine approach.
202 Docker is used by many organizations (e.g., National Center for Atmospheric Research,
203 National Snow and Ice Data Center, NSF Agave API), and widely supported across large-scale
204 cloud compute environments (e.g., Amazon EC2 Container Service, Google Container



205 Engine, NSF Xsede). It is particularly well suited to NEON’s DevOps strategy: combining
206 development, operation and quality assurance to enable creating, testing, deploying and
207 updating scientific software rapidly and reliably (Figure 1).

208 Docker can build images automatically by reading the instructions from a Dockerfile. A
209 Dockerfile is a text document that contains all the commands a user would call on the command
210 line to assemble an image. Using e.g. a cloud hosting platform like DockerHub
211 (<https://hub.docker.com/>), the image build and distribution can be automated. This is realized
212 through executing the series of command-line instructions defined in the Dockerfile whenever
213 a new eddy4R source code version is available on GitHub. A key feature of eddy4R-Docker is
214 that it builds upon “Rocker” pre-built Docker images, maintained by the rOpenSci group
215 (<https://ropensci.org/>). This ensures access to stable, up-to-date base images containing R and
216 a variety of packages commonly used. The eddy4R-Docker image (0.1.0) released in this study
217 was built based on the rocker/ropensci/latest image containing R (3.3.2;
218 <https://hub.docker.com/t/rocker/ropensci/builds/>). As specified in the eddy4R Dockerfile, our
219 R packages eddy4R.base (0.1.0) and eddy4R.qaqc (0.1.0) and their dependencies were
220 automatically built on top of this base image. To complete the eddy4R-Docker processing,
221 analysis and modeling environment, the NEON data portal API Client nneo (0.0.3.9100) as well
222 as the REddyProc (0.8-2) high-level utilities for aggregated EC data were also included. In
223 addition, the user can install any desired R packages to customize the environment.

224 Docker’s benefits to scientific software development are described in detail in Boettiger (2015).
225 For NEON’s purposes, several Docker properties are particularly important:

- 226 • **Portability:** Docker images are portable and independent of the underlying operating
227 system. This enables scientists to develop code on local computers or virtual machines
228 without worrying about the deployment architecture.
- 229 • **Reproducibility:** The DevOps principles are ingrained into the Docker build process,
230 thus ensuring a fully traceable and documented Docker image.
- 231 • **Streamlined interface between Science and CI:** Defined inputs, outputs and
232 instructions provide an ideal framework to isolate and package algorithmic services for
233 operational deployment.
- 234 • **Continuous development and integration:** Docker provides a modular and extensible
235 framework, permitting NEON’s data processing to remain up-to-date with the latest
236 algorithmic developments. As shown by the nneo and REddyProc examples, it enables
237 directly leveraging community-developed code. In this way eddy4R-Docker is
238 functionally extensible, while making it easy for the community to incorporate NEON-
239 developed code into their own data processing.

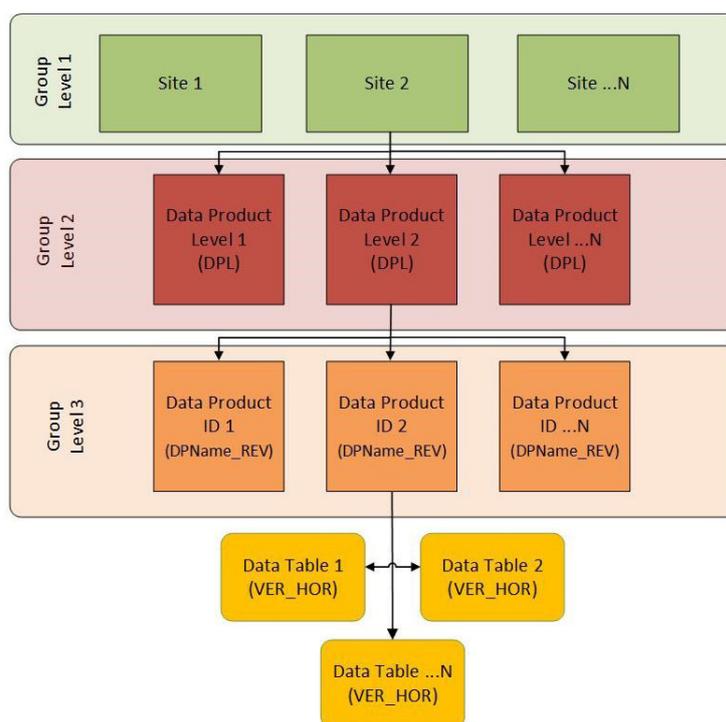
240 **2.4 Hierarchical Data Format version 5**

241 The capability to process large data sets is reliant upon efficient input and output of data, data
242 compressibility to reduce compute resource loads, and the ability to easily package and access
243 metadata. The Hierarchical Data Format (HDF5) is a file format that can meet these needs. A
244 NEON standard HDF5 file structure and metadata attributes allow users to explore larger data



245 sets in an intuitive “directory-like” structure that is based upon the NEON data product naming
 246 convention (see Figure 4). This provides a streamlined data-delivery mechanism for the
 247 eddy4R-Docker processing framework. For the tower datasets analyzed in this study, including
 248 sonic anemometer, infrared gas analyzer and mass flow controller data, file sizes ranged from
 249 1 GB for the uncompressed data to 0.1 – 0.2 GB in HDF5 format, depending on the amount of
 250 missing data. Another important function of the HDF5 file format is the ability to attach
 251 metadata as attributes. The data in this study has the units and variable names as metadata
 252 attached to the data tables in the HDF5 file. As a result, HDF5 and similar self-documenting
 253 hierarchical data formats are gaining traction in a community that has traditionally relied on
 254 ASCII text column or comma-delimited files, especially as tools for viewing, manipulating, and
 255 extracting data from HDF5 become more commonplace.

256



257

258 Figure 4 The NEON HDF5 file structure based on the NEON data product naming convention.

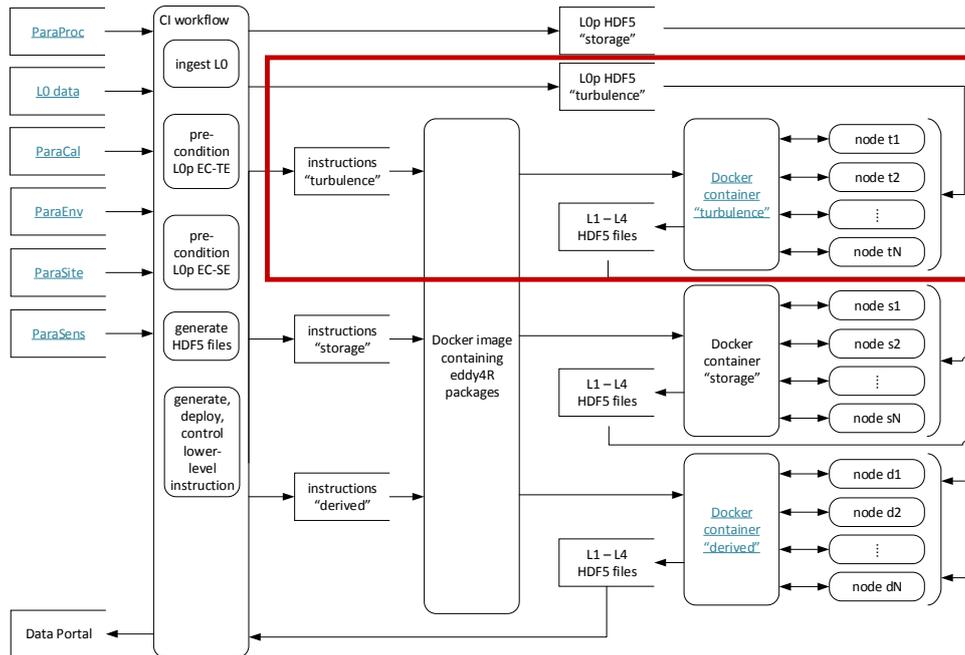
259

260 2.5 Modular compatibility with existing compute infrastructure

261 To perform a defined series of processing steps, a Docker image is called with an instruction
 262 file, resulting in a running instance called Docker container (Figure 5). Through this mechanism,
 263 an arbitrary number of Docker containers can be run simultaneously performing identical or



264 different services depending on the instruction file. This provides an ideal framework for scaled
 265 deployment using e.g. high-throughput compute architectures, cloud-based services etc.
 266



267
 268 Figure 5. NEON's EC workflow. The red box visualizes the scope of the present study, and
 269 individual workflow components are described in the text.

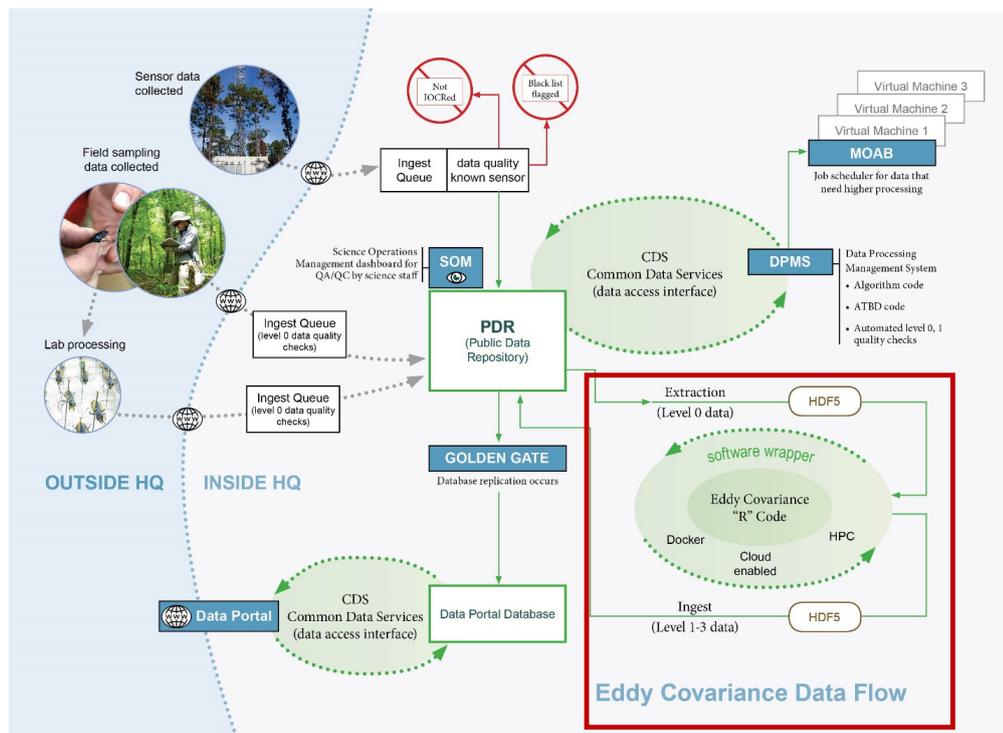
270
 271 NEON's eddy4R-Docker workflow begins with ingesting information from various data
 272 sources on a site-by-site basis (Figure 5 top left). This includes EC raw data (Level 0, or L0
 273 data) alongside contextual information on measurement site (ParaSite), environment (ParaEnv),
 274 sensor (ParaSens), calibration (ParaCal), as well as processing parameters (ParaProc). Next,
 275 the raw data is preconditioned and all information is hierarchically combined into a compact and
 276 easily transferable HDF5 file (Figure 5 panel "CI workflow"). Each file contains the calibrated
 277 raw data (L0 prime, or L0p) for one site and one day, either for EC turbulent exchange or
 278 storage exchange. In this manuscript, we focus on demonstrating the turbulence data process
 279 and analysis in the red box of Figure 5. Together with the "turbulence" instruction file the HDF5
 280 L0p data file is passed to the eddy4R-Docker image, where a running Docker container is
 281 created that scales the computation over a specified number of compute nodes (Figure 5 top
 282 right). The resulting higher-level data products (Level 1 – Level 4, or L1-L4)
 283 are collected from the compute nodes and, together with all contextual information, are combined
 284 into a daily L1-L4 HDF5 data file that is served on the data portal (Figure 5 bottom left). This
 285 sequence is performed analogously for different combinations of instructions and data, and it is
 286 possible for the instruction sets to interact with each other. For example, the "turbulence" and "storage"



287 containers are processing in parallel, and starting the “derived” container once all intermediary
 288 results are available (Figure 5 bottom right).

289 This eddy4R-Docker workflow modularly integrates into pre-existing data processing pipelines,
 290 such as the one of NEON (Figure 6): in NEON’s pre-existing framework the CI group encoded
 291 simple algorithms (e.g. temporal means) in Java, based on algorithm documentation provided
 292 by Science staff. The key difference of the eddy4R-Docker workflow is that instead of
 293 algorithm documentation, NEON Science staff now provides documented algorithms that
 294 perform a complex series of processing steps, which can be directly deployed by CI. Not only
 295 does this adoption of the DevOps workflow (Figure 1) streamline end-to-end operational
 296 implementation and efficiency, it empowers the Science community at large by putting the key
 297 to the scientific algorithms into the hand of scientists.

298



299

300 Figure 6. NEON’s streaming processing framework for EC data. The red box visualizes the
 301 Docker deployment within the overall CI framework.

302

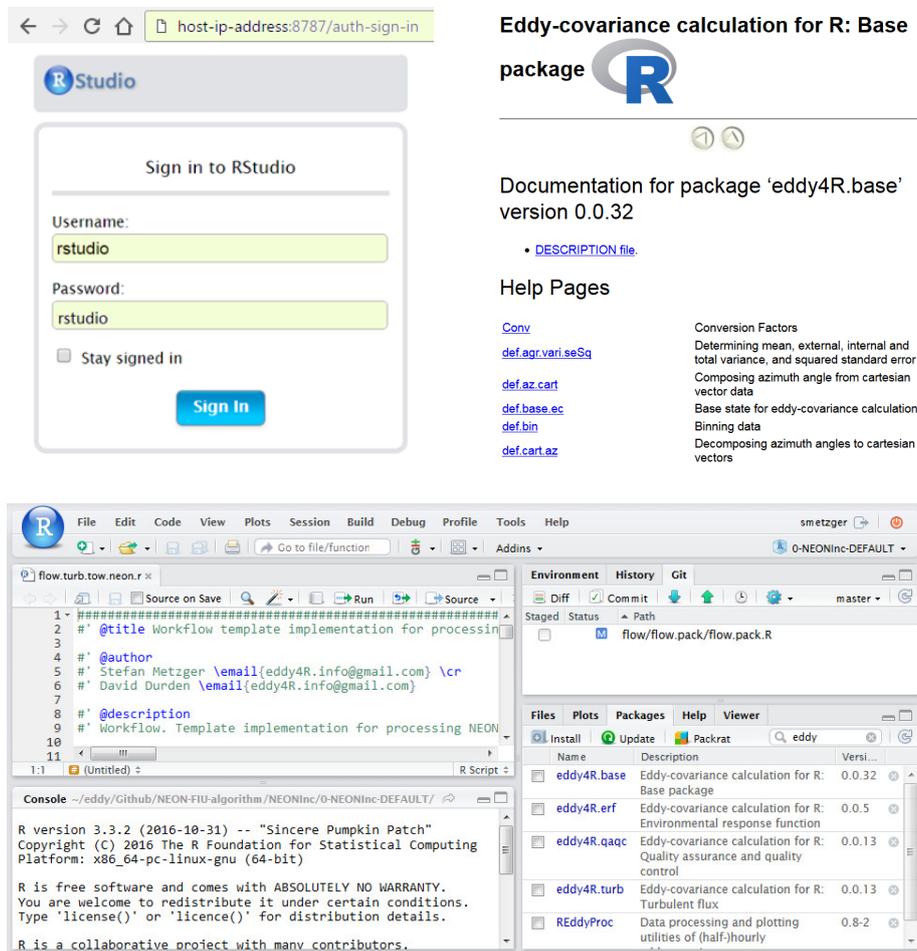
303 2.6 Installation and operation

304 One source of resistance to reproducible research is the initial burden of learning a new
 305 workflow. The eddy4R-Docker image aims to reduce the initial setup effort and learning



306 requirements. This is achieved by providing a computational environment that contains all the
307 necessary software dependencies, the Rstudio graphical development environment
308 (<https://www.rstudio.com/>), and a code base consisting of workflow templates and easily
309 accessible functions. Combined with simple and thoroughly documented installation procedure
310 it provides a similar feel to working locally.

311



312 Figure 7 Docker-based Rstudio server session login via web-browser. Top left panel: Sign-in
313 screen with highlighted areas showing information to input by the user. Top right panel:
314 Interactive help for the eddy4R.base package. Bottom panel: Integrated development
315 environment with workflow template, R console, Git staging area and eddy4R packages.
316

317 To work with the eddy4R-Docker image, one first needs to sign up at DockerHub
318 (<https://hub.docker.com/>) and install the Docker host software following the Docker installation



319 instructions (https://docs.docker.com/engine/getstarted/step_one/). Next, the download of the
320 eddy4R-Docker image and subsequent creation of a container can be performed by two simple
321 commands in an open shell (Linux/Mac) or the Docker Quickstart Terminal (Windows):

```
322     docker login  
323     docker run -d -p 8787:8787 stefanmet/eddy4r:v0.1.0
```

324 The first command will prompt for the user's DockerHub ID and password. The second
325 command will download the latest eddy4R-Docker image and start a Docker container that
326 utilizes port 8787 for establishing a graphical interface via web-browser. The release version of
327 the Docker image can be specified, or alternatively the specifier `latest` provides the most up-
328 to-date development image.

329 The interactive Rstudio Server session running inside the Docker container can then be accessed
330 via a web browser at `http://host-ip-address:8787`, using the IP address of the Docker host
331 machine. The IP address of the Docker host can be determined by typing `localhost` in a shell
332 session (Linux/Mac) or by typing `docker-machine ip default` in `cmd.exe` (Windows).
333 Lastly, in the web browser the user can log into the RStudio session with username and
334 password `rstudio` (see Figure 7 top left panel). Figure 7 also shows the Rstudio integrated
335 development environment and interactive help for the `eddy4R.base` package in the bottom and
336 top left panels, respectively. Additional information about the use of Rstudio and eddy4R
337 packages in Docker containers can be found on the `rocker-org/rocker` website
338 (<https://github.com/rocker-org/rocker/wiki/Using-the-RStudio-image>) and the eddy4R Wiki
339 pages, such as scaled deployment from the command line without graphical user interface.

340 **3 Example applications**

341 In the following we present three example applications of eddy4R-Docker. The calculations
342 were performed on 12 Intel Xeon X5550 2.67GHz CPUs, 32 GB memory with 10 Mbit
343 interconnects and 10 Mbit access to 8 TB storage on an Oracle Zettabyte File System. The
344 software specifications were CentOS 7 (3.10.0-327.el7.x86_64) with docker-engine (1.11.0).
345 In Sect. 3.1, results of 12 days of EC data from a fixed tower at a NEON field site are shown.
346 Next, in Sect. 3.2, we present EC fluxes from a 1-hour recording of a moving platform: airborne
347 observations in a convectively mixed boundary layer. Lastly, a validation via software
348 intercomparison is provided in Sect. 3.3.

349 **3.1 Tower eddy-covariance measurements**

350 Here, we use tower EC measurements to illustrate a typical implementation of the eddy4R
351 processing framework. The Smithsonian Environmental Research Center (SERC) in
352 Edgewater, MD, USA is located on the Rhode and West Rivers, and hosts the NEON SERC
353 tower (38°53'24.29" N, 76°33'36.04" W; 30 m a.s.l.). The ecosystem at SERC is a closed-
354 canopy hardwood deciduous forest dominated by tulip popular, oak and ash, with a mean
355 canopy height of approximately 38 m (Figure 8). EC turbulent flux sensors are mounted at the
356 tower top at 62 m above ground or 24 m above the forest canopy.



357 An enclosed infrared gas analyzer (IRGA, LI-COR Biosciences, Lincoln, NE, USA, model: LI-
358 7200, firmware v7.3.1. was used to measure the turbulent fluctuations of H₂O and CO₂. A mass
359 flow controller (Alicat Scientific, Burlington, VT, USA, model: MCRW-20 SLPM-DS-NEON)
360 was used to maintain a constant flow rate of 12 SLPM through the IRGA cell. A sonic
361 anemometer (Campbell Scientific, Logan, UT, USA, model: CSAT3, firmware v3) was used to
362 measure the 3-dimensional turbulent wind components. Data from the IRGA and the sonic
363 anemometer was synchronized using triggering and network timing protocol, and collected
364 simultaneously at 20 Hz sampling rate.

365 Here, data from April 22 to May 3, 2016 were used. The mean temperature during this time
366 period was 15°C, with a maximum temperature of 29°C and a minimum of 8°C. A total of
367 15 mm of precipitation was observed at nearby Annapolis Naval Academy.

368



369 Figure 8. Left panel: Ecosystem at the NEON SERC tower (credit: Stephen Voss Photography;
370 <http://www.stephenvoss.com/blog/neon-tower-smithsonian>). Right panels: EC instrumentation
371 on top of the NEON SERC tower. Right top panel: Campbell Scientific CSAT-3 three-
372 dimensional sonic anemometer (front) and LI-COR Biosciences LI-7200 infrared gas analyser
373 (back) on the retracted tower-top boom. Right bottom panel: Same instrumentation but with the
374 tower-top boom extended at 230° from true north.



375 **3.1.1 Algorithm settings and profiling**

376 The eddy4R workflow file was configured to ingest on the order of 50 data streams at 20 Hz,
377 including 3-D wind components, sonic temperature, and H₂O and CO₂ concentrations. The data
378 were processed to half-hourly L1 data products and turbulent fluxes. The L1 data products are
379 essentially state variables (wind, temperature, concentrations) with basic statistical products
380 derived, i.e. mean, minimum, maximum, standard error of the mean and variance. The
381 algorithmic processing for the L4 flux calculations requires additional scientific and procedural
382 complexity to test theoretical assumptions of the EC theory. The resultant fluxes represent half-
383 hourly vertical turbulent exchanges between the earth's surface and the atmosphere
384 corresponding to these state variables.

385 For the datasets analyzed in this study, the file sizes ranged from 0.1 – 0.2 GB in HDF5 format
386 depending on the amount of missing data, with metadata attached as attributes. We used the
387 simple data format for our HDF5 files, as opposed to compound data type, this resulted in
388 reduced read in time from 60 seconds to 3 seconds for 20 Hz IRGA data. Elementary testing
389 indicates that in this framework 6 CPU-minutes were required to process 1 day of 20 Hz L0
390 data, and 1.2 CPU-minutes per 1 day of L0p data (100,000,000 observations). No reduction in
391 efficiency was observed between direct software deployment and its Docker implementation.
392 Once scientific QA/QC and uncertainty budget is implemented, the computational expense will
393 likely increase by a factor of two to three. This suggests that eddy4R performs comparably to
394 other flux processors. Memory usage is kept below 2 GB through the use of fast access file-
395 backed objects, enabling more sophisticated scientific analyses through access to multiple days
396 of data without overloading random access memory (RAM) resources. Additionally, the
397 snowfall R package allows for logical parallelization frameworks to be implemented in the
398 processing framework, even at low-level analysis steps.

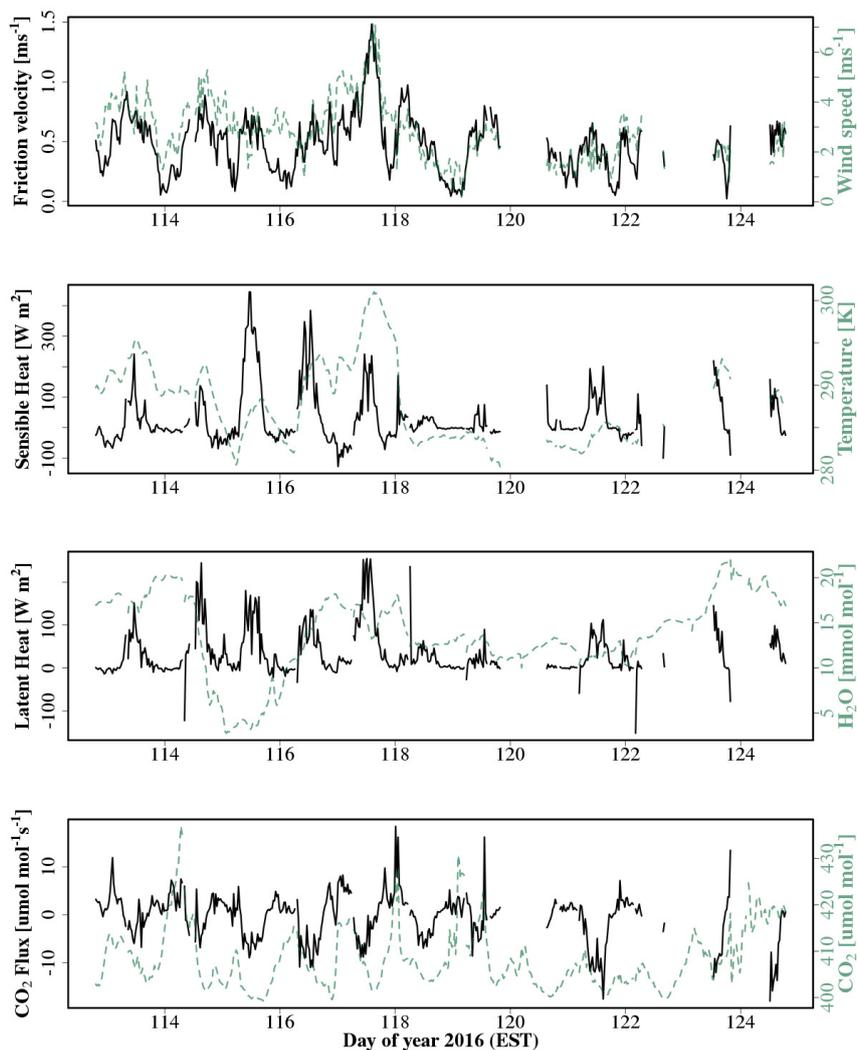
399 **3.1.2 Results and discussion**

400 The time series ranging from April 22 to May 3, 2016 was processed to deliver both state (L1)
401 and flux (L4) quantities; however, the initial eddy4R package release will only contain
402 functions necessary to report state variables or L1 data products in the NEON data product
403 description. The QA/QC and uncertainty frameworks were not fully implemented during the
404 processing of the proof-of-concept results, but averaging periods with >10% missing data (incl.
405 bad sensor diagnostic flags) were removed.

406 Figure 9 shows the resultant time series of shear stress (friction velocity), sensible heat, latent
407 heat and CO₂ flux. The derived values fall into typical ranges for mid-latitude hardwood forests
408 in spring. As expected, fluxes follow the general trends in the scalar quantities. Good data
409 coverage can be seen for the LI-7200 measurements even during the rainy period at the end of
410 the analysis. A footprint analysis revealed that 90% of the flux measurement signals were
411 sourced within 800 m from the tower, and 80% were within 500 m from the tower at our site.
412 Data coverage was reduced after day of year (DOY) 120 due to inclement weather conditions.
413 The spiky results preceding and following periods with >10% invalid data highlight the need
414 for scientific QA/QC and uncertainty budget to provide science-grade fluxes. Nonetheless, this



415 implementation of eddy4R in a Docker image, as it will interact with NEON CI, clearly
416 demonstrates its core capability to generate L1-L4 data products.
417



418
419 Figure 9. Time-series of turbulent fluxes derived from EC measurements atop the NEON SERC
420 tower. Top to bottom: Vertical turbulent exchange of shear (friction velocity) and wind speed,
421 sensible heat and temperature, latent heat and H₂O dry mole fraction and CO₂ flux and CO₂ dry
422 mole fraction.



423 **3.2 Aircraft eddy-covariance measurements**

424 Here, we use aircraft EC measurements to illustrate more advanced scientific capabilities of the
425 eddy4R processing framework. Airborne turbulent flux observations were performed along
426 more than 3100 km of low level (i.e. 50 m above ground level) flights across the North Slope
427 of Alaska in July 2012, using the research aircraft Polar 5 (Tetzlaff et al., 2015). The example
428 data used in this manuscript were recorded during a SSW-NNE flight line near the village of
429 Atkasuk, Alaska, above tundra dominated by sedges and emerging herbaceous wetland
430 vegetation. Large, often oriented, lakes and the meandering Meade River characterize the
431 surrounding landscape.

432 The aircraft was equipped with a 3 m nose boom holding a 5-hole probe for wind measurements,
433 an open wire Pt100 in an unheated Rosemount housing for air temperature measurements, and
434 an HMT-330 (Vaisala, Helsinki, Finland) in a Rosemount housing for relative humidity.
435 Sample air was drawn from an inlet above the cabin at about 9.7 l s^{-1} , analysed in an RMT-200
436 (Los Gatos Research Inc., Mountain View, California, USA) and recorded at 20 Hz. Aircraft
437 position and attitude was provided by several Global Positioning Systems (NovAtel Inc.,
438 Calgary, Alberta, USA) and an Inertial Navigation System (Laseref V, Honeywell International
439 Inc., Morristown, New Jersey, USA), altitude was determined by a radar altimeter (KRA 405B/
440 Honeywell International Inc., Morristown, New Jersey, USA) and a laser altimeter (LD90/
441 RIEGL Laser Measurements Systems GmbH, Horn, Austria). After spike removal the sampling
442 frequency of the original data was reduced from 100 Hz to 20 Hz resolution using block
443 averaging.

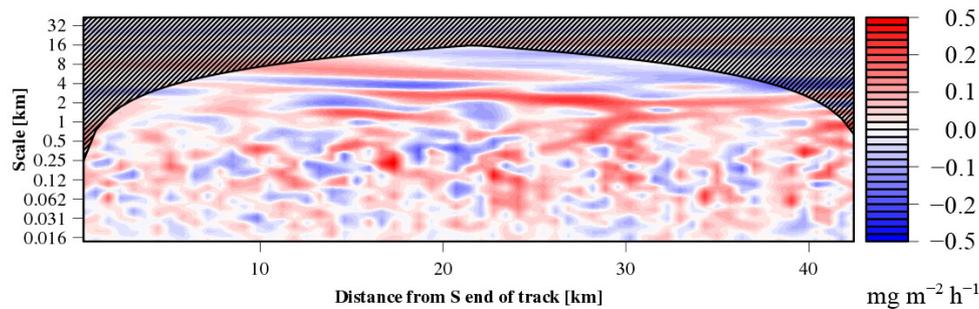
444 **3.2.1 Algorithm settings and profiling**

445 Here, aircraft-measured vertical wind speed and CH_4 dry mole fraction were analysed to
446 determine CH_4 emissions by means of a time-frequency-resolved version of the EC method
447 (Metzger et al., 2013). For this purpose a combination of settings were chosen in the eddy4R
448 workflow file that differ from Sect. 3.1: Initially the small (<1 MB) EC raw data file consisting
449 of 17 variables and 12,800 data points (or 42 km flight data) was read in ASCII Gzip format –
450 standard R capabilities for data ingest can be used to read data in various formats, frequencies
451 and units. Aircraft-measured vertical wind speed and CH_4 dry mole fraction were then
452 correlated using a Wavelet transform (Metzger et al., 2013). This process considers ranging and
453 de-spiking of unphysical raw data values (Mauder et al., 2013; Metzger et al., 2012), fast dry
454 mole fraction derivation (e.g., Burba et al., 2012) and spectroscopic correction (Tuzson et al.,
455 2010) of cavity-ringdown CH_4 trace gas observations, and high-frequency spectral correction
456 (Ammann et al., 2006) by means of applying a sigmoidal transfer function (Eugster and Senn,
457 1995) directly in Wavelet space. This permits estimating turbulent fluxes with improved spatial
458 discretization and determining ~100 biophysically relevant surface properties in the flux
459 footprint. The analysis took 56 minutes with 8-fold parallelization and consumed <3 GB RAM
460 thanks to the use of fast access file-backed objects.



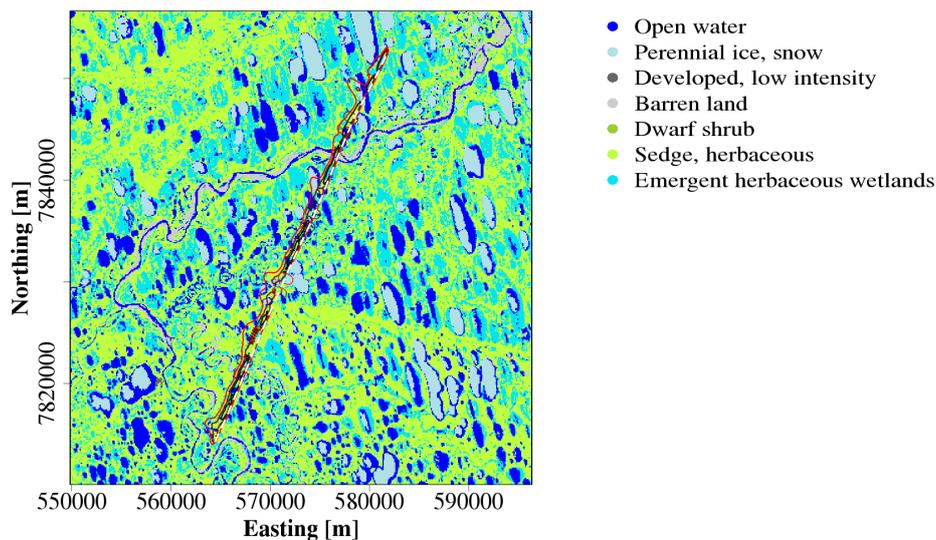
461 **3.2.2 Results and discussion**

462 The resulting Wavelet cross-scalogram (Figure 10) is integrated in frequency over transport
463 scales up to 20 km, and along the flight path over a 1000 m moving window with 100 m step
464 size, similar to the resolution of the land surface data. The result is an in-situ observed space-
465 series of the CH₄ surface-atmosphere exchange at 100 m spatial resolution. Analogously,
466 turbulence statistics characterizing shear stress and buoyancy are determined for characterizing
467 the atmospheric transport between the emitting land surface and the aircraft position.
468



469
470 Figure 10. Wavelet cross-scalogram of the CH₄ flux equivalent to a time (x-axis) frequency (y-
471 axis) resolved version of EC.

472



473 Figure 11. The composite flux footprint along the flight line (30 %, 60 %, 90% contour lines)
474 superimposed over the National Land Cover Database. The white dashed line represents the
475 aircraft flight track.

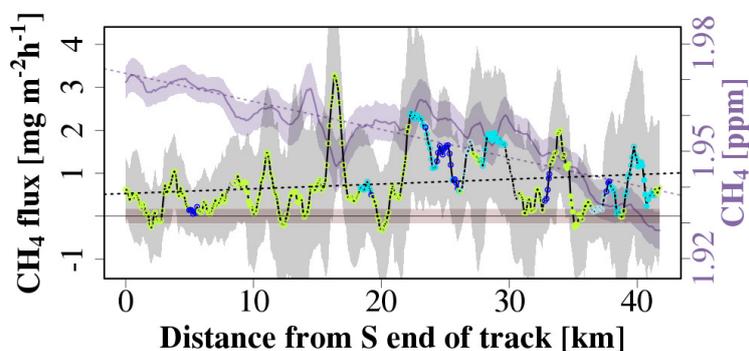


476 Corresponding systematic and random statistical errors are calculated following Lenschow and
477 Stankov (1986) and Lenschow *et al.* (1994), and the flux detection limit is calculated after
478 Billesbach (2011).

479 The relationship between the aircraft-observed CH₄ surface-atmosphere exchange and land
480 surface properties is established through an atmospheric transport operator, the so-called flux
481 footprint function (e.g., Schmid, 1994). Here we use a computationally efficient one-
482 dimensional parameterization of a Lagrangian particle model for the along-wind footprint
483 extent (Kljun *et al.*, 2002; Kljun *et al.*, 2004), combined with an analytical approach to
484 determine cross-wind surface contributions to each 100 m aircraft measurement, depending on
485 aircraft position (Figure 10; Metzger *et al.*, 2012).

486 For each 100 m observation of the CH₄ surface-atmosphere exchange an individual footprint
487 weight matrix derived from the footprint parameterization is convolved with the land surface
488 drivers. The results are space-series of land surface contributions accompanying the CH₄
489 measured surface-atmosphere exchange (Figure 12).

490



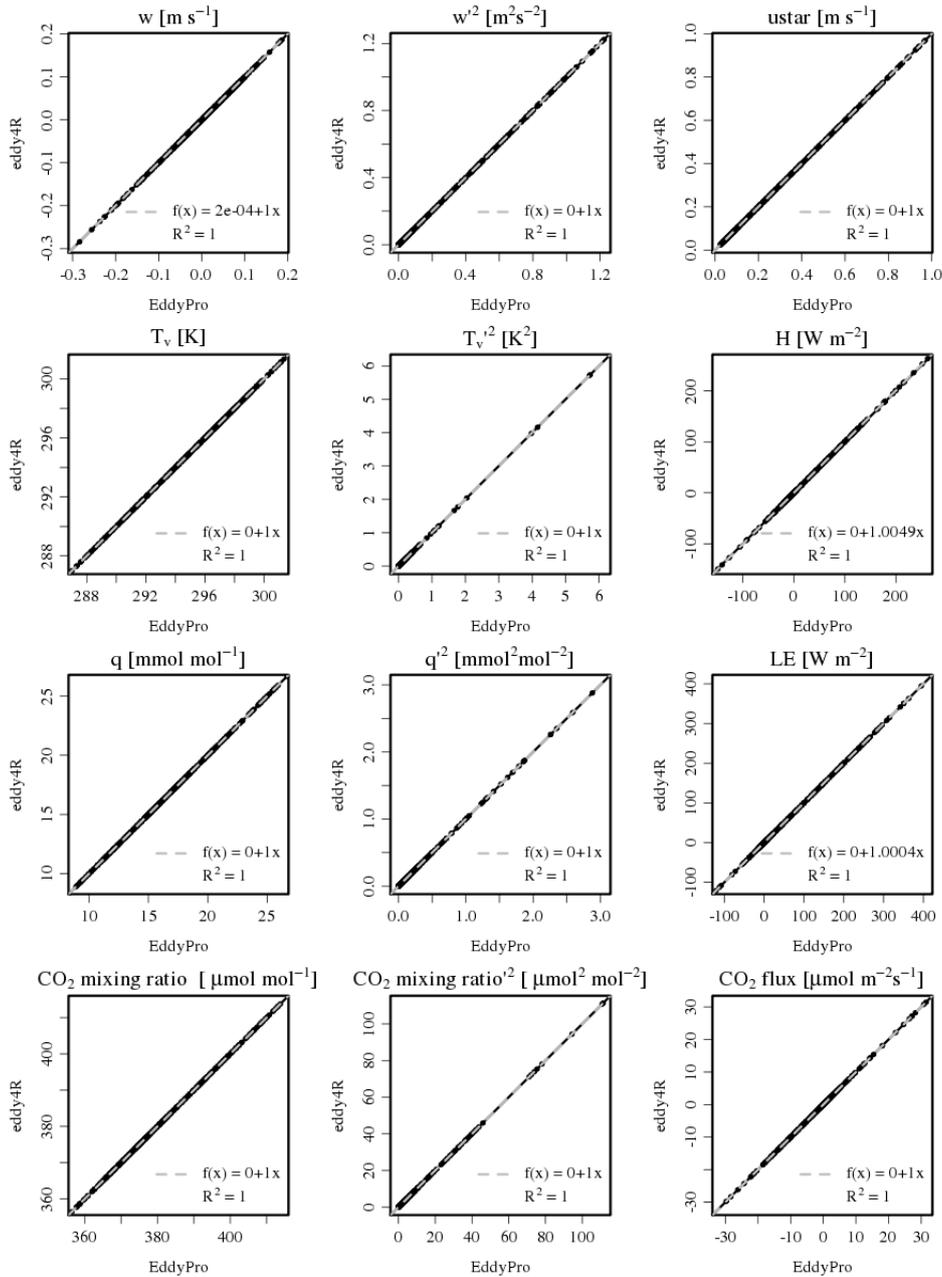
491

492 Figure 12. Space-series of 399 CH₄ concentration (purple line) and flux (black line)
493 observations each 100 m, averaged over 1000 m windows. The random sampling errors are
494 indicated by the shaded areas enveloping each line, and the flux detection limit is shown as
495 salmon envelope around the abscissa. Circles indicate the dominating land cover in the footprint
496 of each observation (Figure 11) with full circles corresponding to ‘pure’ fluxes (>80% surface
497 contribution).

498 3.3 Validation and verification

499 eddy4R includes a verification script which automatically processes subsets of the tower and
500 aircraft data introduced in Sect. 3.1 and Sect. 3.2, and verifies the results against a reference,
501 e.g. generated with a different software.

502 Here, we demonstrate such an approach at the Park Falls, Wisconsin WLEF very tall tower
503 Ameriflux site (US-PFa). The 447 m tall WLEF television tower (45.946°N, 90.272°W) has
504 been instrumented for EC measurements in 1996, and is part of the AmeriFlux network. Flux
505 measurements at 30 m, 122 m and 396 m sample a mixed landscape of forests and wetlands



506

507 Figure 13. Scatterplot of means (vertical wind speed, w ; sonic temperature, T_v ; H_2O dry mole
 508 fraction (q); and CO_2 dry mole fraction), variances and fluxes (friction velocity, u_{star} ; sensible
 509 heat flux, H ; latent heat, LE ; CO_2 flux). Data are generated from 2011 July to Aug WLEF data
 510 in EddyPro and eddy4R. Each point represents a one-hour averaging period. Black lines are 1:1
 511 lines, and dashed lines are robust regressions (Salibian-Barrera and Yohai, 2006).



512 (Desai et al., 2015). The surrounding forest canopy has approximately 70% deciduous and 30%
513 coniferous trees, and a mean canopy height of 20 m. The site has an interior continental climate.
514 Instrumentation at each level consists of fast response wind speed and temperature from a sonic
515 anemometer (Applied Technologies., Inc., Seattle, USA, ATI Type K). 10 Hz dry mole fraction
516 of CO₂ and H₂O at the 122 m level used here were measured by a closed-path infrared gas
517 analyzer (LI-COR, Inc., Lincoln, USA, LI-6262) located on the tower.

518 A data set from July 27 to August 19, 2011 was used in the intercomparison between eddy4R
519 and the reference software EddyPro (LI-COR, Inc., Lincoln, USA, v6.2.0). EddyPro was
520 released in April 2011 and is being widely used in the EC community.

521 3.3.1 Algorithm settings

522 Several preprocessing steps were applied, and the resulting data and settings were used in both,
523 eddy4R and EddyPro: (i) The raw data was pre-cleaned in eddy4R using the Brock (1986) de-
524 spiking algorithm with a filter width of 9 data points for all variables. (ii) EddyPro was used to
525 calculate the planar-fit rotation parameters (Wilczak et al., 2001) over the entire dataset (offset
526 = -0.06 ms^{-1} , pitch = -5.27° , roll = -1.81°). (iii) Time lags for dry mole fractions of CO₂ (0.8 s
527 behind vertical wind) and H₂O (0.1 s behind vertical wind) were calculated in eddy4R using
528 maximum correlation (median lag time over entire dataset).

529 Because CO₂ and H₂O fluxes were calculated from dry mole fractions, the Webb et al. (1980)
530 density correction was not necessary and therefore not applied (Burba et al., 2012). Frequency
531 correction was not considered in this validation and therefore not applied. Means, variances and
532 fluxes were calculated on the basis of one-hour block averages. Based on Schotanus et al.
533 (1983), sensible heat flux was calculated from point-by-point conversion of sonic temperature
534 in eddy4R, and with the half-hourly statistical correction in EddyPro.

535 3.3.2 Results and discussion

536 eddy4R and EddyPro produce nearly identical results (Figure 13), and the gain error is within
537 0.04% for most outputs. Sensible heat flux values produced by eddy4R have slightly larger
538 magnitude compared to EddyPro, by 0.49%. This is likely a result of the different methods
539 applied when converting sonic temperature to air temperature. A detailed end-to-end
540 intercomparison considering additional processing steps and EC software is planned for a
541 separate manuscript accompanying NEON's release of flux data products.

542 4 Summary and conclusions

543 Adopting a DevOps philosophy has facilitated the creation of a flexible, scalable, and extensible
544 processing environment for producing NEON's EC data products. Git-distributed version
545 control facilitates simultaneous internal-external collaboration on scientific algorithms, the
546 outcome being a modular family of open-source R packages. The use of Hierarchical Data
547 Format allows for efficient, self-describing data input and output. Docker images package the
548 entire processing environment for robust, scalable, and portable deployment. The capability of
549 this framework was demonstrated with cross-validated tower and aircraft fluxes.



550 The results presented here are from a file-based implementation of the eddy4R Docker
551 workflow, with EC instrument data accessed directly e.g. from the NEON site and manually
552 processed into the HDF5 ingest format (Sect. 2.5). Focus now shifts to the operational
553 implementation of the eddy4R-Docker workflow for reporting means and variances. This
554 includes: (i) Automated ingest of streaming raw data into the NEON database; (ii) Processing
555 of raw data into the standard, defined inputs required by the eddy4R-Docker in HDF5 format,
556 and (iii) Developing the software and hardware infrastructure to pass data and instructions back
557 and forth to the eddy4R-Docker workflow, and control program execution in a distributed
558 computing framework. Lessons learned here will profit the community at large, e.g. through
559 enabling streaming processing directly at an EC site or over cellular modems with the same
560 eddy4R-Docker open-source software as used for sophisticated analyses (Sect. 3.2).

561 Thereafter, remaining scientific algorithms will be integrated in eddy4R-Docker for producing
562 defensible turbulent exchange data products. These algorithms include on-the-fly de-spiking,
563 lag correction, planar-fit and spectral correction, scientific QA/QC, and uncertainty budget
564 estimation. Finally, the eddy4R-Docker will be expanded to include “storage” and “derived”
565 workflows (Figure 6) for producing defensible net ecosystem exchange data products in 2018.

566 While our sole focus in developing this framework has been to facilitate generating EC data
567 products with the unique capabilities and constraints of NEON, it has become clear that the
568 framework has the potential for enabling the implementation of a suite of complex processing
569 algorithms, such as temporal gap filling of sensor time series data or modeling re-aeration rates.
570 There exist many potential synergies between NEON, other tower networks, and the user
571 community for producing high level EC data products. We hope this framework can serve as a
572 model for implementing community-sourced, distributed-development scientific code while
573 combatting the deficiencies of current computational frameworks that limit accessibility,
574 reproducibility, and extensibility.

575 **5 Code and data availability**

576 The source code packages eddy4R.base (0.1.0) and eddy4R.qaqc (0.1.0) used in this study are
577 archived at https://w3id.org/smetzger/Metzger-et-al_2017_eddy4R-Docker/code, under the
578 GNU Affero General Public License (GNU AGPLv3). Similarly, at
579 https://w3id.org/smetzger/Metzger-et-al_2017_eddy4R-Docker/docker the corresponding
580 eddy4R-Docker image (0.1.0) is available. Lastly, a data supplement is provided at
581 https://w3id.org/smetzger/Metzger-et-al_2017_eddy4R-Docker/data, including an extended
582 abstract and all NEON SERC raw data used in this study, accompanied by variable
583 documentation.

584 **Acknowledgements**

585 Many colleagues at Battelle Ecology supported this study. In particular, Santiago Bonarrigo
586 provided pre-parsed high-frequency data from the SERC site, Andrew Fox (now: National
587 Center for Atmospheric Research), Mike SanClements and David Hulslander commented on
588 an earlier version of the manuscript, and Andrea Thorpe, Thomas Gulbransen and Michael
589 Kuhlman helped shepherding this study and its publication through required administrative
590 procedures. The National Ecological Observatory Network is a project sponsored by the



591 National Science Foundation and managed under cooperative agreement by Battelle Ecology,
592 Inc. This material is based upon work supported by the National Science Foundation under the
593 grant DBI-0752017. Any opinions, findings, and conclusions or recommendations expressed in
594 this material are those of the author(s) and do not necessarily reflect the views of the National
595 Science Foundation. Ankur Desai acknowledges support from NSF DBI-1457897 and DOE
596 Office of Science Ameriflux Network Management Project core site support to the ChEAS
597 cluster. Torsten Sachs and Andrei Serafimovich are supported by the Helmholtz Association of
598 German Research Centres through a Helmholtz Young Investigators Group grant to Torsten
599 Sachs (grant VH-NG-821).

600 References

601 Ammann, C., Brunner, A., Spirig, C., and Neftel, A.: Technical note: Water vapour
602 concentration and flux measurements with PTR-MS, *Atmos. Chem. Phys.*, 6, 4643-4651,
603 doi:10.5194/acp-6-4643-2006, 2006.

604 Aubinet, M., Vesala, T., and Papale, D., (Eds.): *Eddy covariance: A practical guide to
605 measurement and data analysis*, Springer, Dordrecht, Heidelberg, London, New York, 438 pp.,
606 2012.

607 Baldocchi, D., Falge, E., Gu, L., Olson, R., Hollinger, D., Running, S., Anthoni, P., Bernhofer,
608 C., Davis, K., Evans, R., Fuentes, J., Goldstein, A., Katul, G., Law, B., Lee, X., Malhi, Y.,
609 Meyers, T., Munger, W., Oechel, W., U, K., Pilegaard, K., Schmid, H., Valentini, R., Verma,
610 S., Vesala, T., Wilson, K., and Wofsy, S.: FLUXNET: A new tool to study the temporal and
611 spatial variability of ecosystem-scale carbon dioxide, water vapor, and energy flux densities,
612 *Bull. Am. Meteorol. Soc.*, 82, 2415-2434, doi:10.1175/1520-
613 0477(2001)082<2415:FANTTS>2.3.CO;2, 2001.

614 Billesbach, D. P.: Estimating uncertainties in individual eddy covariance flux measurements:
615 A comparison of methods and a proposed new method, *Agric. For. Meteorol.*, 151, 394-405,
616 doi:10.1016/j.agrformet.2010.12.001, 2011.

617 Boettiger, C.: An introduction to Docker for reproducible research, with examples from the R
618 environment, *Operating Systems Review*, 49, 71-79, doi:10.1145/2723872.2723882, 2015.

619 Brock, F. V.: A nonlinear filter to remove impulse noise from meteorological data, *J. Atmos.
620 Oceanic Technol.*, 3, 51-58, doi:10.1175/1520-0426(1986)003<0051:anftri>2.0.co;2, 1986.

621 Burba, G., Schmidt, A., Scott, R. L., Nakai, T., Kathilankal, J., Fratini, G., Hanson, C., Law,
622 B., McDermitt, D. K., Eckles, R., Furtaw, M., and Velgersdyk, M.: Calculating CO₂ and H₂O
623 eddy covariance fluxes from an enclosed gas analyzer using an instantaneous mixing ratio,
624 *Global Change Biol.*, 18, 385-399, doi:10.1111/j.1365-2486.2011.02536.x, 2012.

625 Clark, D., Culich, A., Hamlin, B., and Lovett, R.: BCE: Berkeley's common scientific compute
626 environment for research and education, *Proceedings of the 13th Python in Science Conference
627 (SCIPY 2014) Austin, USA, 2014.*



- 628 Clement, R. J., Burba, G. G., Grelle, A., Anderson, D. J., and Moncrieff, J. B.: Improved trace
629 gas flux estimation through IRGA sampling optimization, *Agric. For. Meteorol.*, 149, 623-638,
630 doi:10.1016/j.agrformet.2008.10.008, 2009.
- 631 Collberg, C., Proebsting, T., Moraila, G., Shankaran, A., Shi, Z., and Warren, A. M.: Measuring
632 reproducibility in computer systems research, University of Arizona, Department of Computer
633 Science, Tucson, USA, 37, 2014.
- 634 De Roo, F., Abdul Huq, S. U., Metzger, S., Desai, A. R., Xu, K., and Mauder, M.: On the benefit
635 of driving large-eddy simulation with spatially resolved surface fluxes derived from
636 environmental response functions, TERENO International Conference, Bonn, Germany, 29
637 September - 2 October, 2014.
- 638 Desai, A. R., Xu, K., Tian, H., Weishampel, P., Thom, J., Baumann, D., Andrews, A. E., Cook,
639 B. D., King, J. Y., and Kolka, R.: Landscape-level terrestrial methane flux observed from a very
640 tall tower, *Agric. For. Meteorol.*, 201, 61-75, doi:10.1016/j.agrformet.2014.10.017, 2015.
- 641 Erich, F., Amrit, C., and Daneva, M.: A mapping study on cooperation between information
642 system development and operations, 15th International Conference on Product-Focused
643 Software Process Improvement, PROFES 2014, Helsinki, Finland, 2014.
- 644 Eugster, W., and Senn, W.: A cospectral correction model for measurement of turbulent NO₂
645 flux, *Boundary Layer Meteorol.*, 74, 321-340, doi:10.1007/bf00712375, 1995.
- 646 Foken, T., and Wichura, B.: Tools for quality assessment of surface-based flux measurements,
647 *Agric. For. Meteorol.*, 78, 83-105, doi:10.1016/0168-1923(95)02248-1, 1996.
- 648 Foken, T.: *Micrometeorology*, Springer, Berlin, Heidelberg, 306 pp., 2008.
- 649 Fratini, G., and Mauder, M.: Towards a consistent eddy-covariance processing: An
650 intercomparison of EddyPro and TK3, *Atmos. Meas. Tech.*, 7, 2273-2281, doi:10.5194/amt-7-
651 2273-2014, 2014.
- 652 Kljun, N., Rotach, M. W., and Schmid, H. P.: A three-dimensional backward lagrangian
653 footprint model for a wide range of boundary-layer stratifications, *Boundary Layer Meteorol.*,
654 103, 205-226, doi:10.1023/A:1014556300021, 2002.
- 655 Kljun, N., Calanca, P., Rotach, M. W., and Schmid, H. P.: A simple parameterisation for flux
656 footprint predictions, *Boundary Layer Meteorol.*, 112, 503-523,
657 doi:10.1023/B:BOUN.0000030653.71031.96, 2004.
- 658 Kljun, N., Calanca, P., Rotach, M. W., and Schmid, H. P.: A simple two-dimensional
659 parameterisation for Flux Footprint Prediction (FFP), *Geosci. Model Dev.*, 8, 3695-3713,
660 doi:10.5194/gmd-8-3695-2015, 2015.
- 661 Kohnert, K., Serafimovich, A., Metzger, S., Hartman, J., and Sachs, T.: Geogenic sources
662 strongly contribute to the Mackenzie River Delta's methane emissions derived from airborne
663 flux data, 48th AGU annual Fall Meeting, San Francisco, U.S.A., 14 - 18 December, 2015.
- 664 Kormann, R., and Meixner, F. X.: An analytical footprint model for non-neutral stratification,
665 *Boundary Layer Meteorol.*, 99, 207-224, doi:10.1023/A:1018991015119, 2001.



- 666 Law, B.: AmeriFlux network aids global synthesis, *Eos, Transactions American Geophysical*
667 *Union*, 88, 286-286, doi:10.1029/2007eo280003, 2007.
- 668 Lee, J., Vaughan, A., Lewis, A., Shaw, M., Purvis, R., Carlslaw, D., Hewitt, C., Misztal, P.,
669 Metzger, S., Beevers, S., Goldstein, A., Karl, T., and Davison, D.: Spatially resolved emissions
670 of NO_x and VOCs and comparison to inventories, 48th AGU annual Fall Meeting, San Francisco,
671 U.S.A., 14 - 18 December, 2015.
- 672 Lenschow, D. H., Mann, J., and Kristensen, L.: How long is long enough when measuring
673 fluxes and other turbulence statistics?, *J. Atmos. Oceanic Technol.*, 11, 661-673,
674 doi:10.1175/1520-0426(1994)011<0661:HLILEW>2.0.CO;2, 1994.
- 675 Loukides, M.: *What is DevOps? Infrastructure as Code*, O'Reilly Media, Ebook, Safari Books
676 Online, 15 pp., 2012.
- 677 Mammarella, I., Peltola, O., Nordbo, A., Järvi, L., and Rannik, Ü.: Quantifying the uncertainty
678 of eddy covariance fluxes due to the use of different software packages and combinations of
679 processing steps in two contrasting ecosystems, *Atmos. Meas. Tech.*, 9, 4915-4933,
680 doi:10.5194/amt-9-4915-2016, 2016.
- 681 Mauder, M., and Foken, T.: *Documentation and instruction manual of the eddy-covariance*
682 *software package TK3*, Universität Bayreuth, Arbeitsergebnisse Abteilung Mikrometeorologie,
683 46, Bayreuth, Germany, 60 pp., ISSN 1614-8924, 2011.
- 684 Mauder, M., Cuntz, M., Drüe, C., Graf, A., Rebmann, C., Schmid, H. P., Schmidt, M., and
685 Steinbrecher, R.: A strategy for quality and uncertainty assessment of long-term eddy-
686 covariance measurements, *Agric. For. Meteorol.*, 169, 122-135,
687 doi:10.1016/j.agrformet.2012.09.006, 2013.
- 688 Metzger, S., Junkermann, W., Mauder, M., Beyrich, F., Butterbach-Bahl, K., Schmid, H. P.,
689 and Foken, T.: Eddy-covariance flux measurements with a weight-shift microlight aircraft,
690 *Atmos. Meas. Tech.*, 5, 1699-1717, doi:10.5194/amt-5-1699-2012, 2012.
- 691 Metzger, S., Junkermann, W., Mauder, M., Butterbach-Bahl, K., Trancón y Widemann, B.,
692 Neidl, F., Schäfer, K., Wieneke, S., Zheng, X. H., Schmid, H. P., and Foken, T.: Spatially
693 explicit regionalization of airborne flux measurements using environmental response functions,
694 *Biogeosciences*, 10, 2193-2217, doi:10.5194/bg-10-2193-2013, 2013.
- 695 Metzger, S., Burba, G., Burns, S. P., Blanken, P. D., Li, J., Luo, H., and Zulueta, R. C.:
696 Optimization of an enclosed gas analyzer sampling system for measuring eddy covariance
697 fluxes of H₂O and CO₂, *Atmos. Meas. Tech.*, 9, 1341-1359, doi:10.5194/amt-9-1341-2016,
698 2016.
- 699 Nordbo, A., and Katul, G.: A wavelet-based correction method for eddy-covariance high-
700 frequency losses in scalar concentration measurements, *Boundary Layer Meteorol.*, 146, 81-
701 102, doi:10.1007/s10546-012-9759-9, 2012.
- 702 Paarsch, H. J., and Golyaev, K.: *A gentle introduction to effective computing in quantitative*
703 *research: What every research assistant should know*, MIT Press, Cambridge, USA, 776 pp.,
704 2016.



- 705 Papale, D., Reichstein, M., Aubinet, M., Canfora, E., Bernhofer, C., Kutsch, W., Longdoz, B.,
706 Rambal, S., Valentini, R., Vesala, T., and Yakir, D.: Towards a standardized processing of Net
707 Ecosystem Exchange measured with eddy covariance technique: algorithms and uncertainty
708 estimation, *Biogeosciences*, 3, 571-583, doi:10.5194/bg-3-571-2006, 2006.
- 709 R Core Team: R: A language and environment for statistical computing, R Foundation for
710 Statistical Computing, Vienna, Austria, 2016.
- 711 Ram, K.: Git can facilitate greater reproducibility and increased transparency in science, *Source
712 Code Biol. Med.*, 8, 7, doi:10.1186/1751-0473-8-7, 2013.
- 713 Raupach, M. R., Rayner, P. J., Barrett, D. J., DeFries, R. S., Heimann, M., Ojima, D. S., Quegan,
714 S., and Schimmlus, C. C.: Model-data synthesis in terrestrial carbon observation: Methods, data
715 requirements and data uncertainty specifications, *Global Change Biol.*, 11, 378-397,
716 doi:10.1111/j.1365-2486.2005.00917.x, 2005.
- 717 Running, S. W., Baldocchi, D. D., Turner, D. P., Gower, S. T., Bakwin, P. S., and Hibbard, K.
718 A.: A global terrestrial monitoring network integrating tower fluxes, flask sampling, ecosystem
719 modeling and EOS satellite data, *Remote Sens. Environ.*, 70, 108-127, doi:10.1016/S0034-
720 4257(99)00061-9, 1999.
- 721 Sachs, T., Serafimovich, A., Metzger, S., Kohnert, K., and Hartmann, J.: Low permafrost
722 methane emissions from arctic airborne flux measurements, 47th AGU annual Fall Meeting, San
723 Francisco, U.S.A., 15 - 19 December, 2014.
- 724 Salibian-Barrera, M., and Yohai, V. J.: A fast algorithm for S-regression estimates, *Journal Of
725 Computational And Graphical Statistics*, 15, 414-427, 2006.
- 726 Salmon, O., Caulton, D., Shepson, P., Brian, S., Metzger, S., and Musinsky, J.: Attributing
727 airborne measurements of forest CO₂ exchange to finer spatial scales, 5th NACP Principal
728 Investigators Meeting, Washington D.C., U.S.A., 26 - 29 January, 2015.
- 729 Schimel, D., Hargrove, W., Hoffman, F., and MacMahon, J.: NEON: a hierarchically designed
730 national ecological network, *Frontiers in Ecology and the Environment*, 5, 59-59,
731 doi:10.1890/1540-9295(2007)5[59:nahdne]2.0.co;2, 2007.
- 732 Schmid, H. P.: Source areas for scalars and scalar fluxes, *Boundary Layer Meteorol.*, 67, 293-
733 318, doi:10.1007/bf00713146, 1994.
- 734 Schotanus, P., Nieuwstadt, F. T. M., and Bruin, H. A. R.: Temperature measurement with a
735 sonic anemometer and its application to heat and moisture fluxes, *Boundary Layer Meteorol.*,
736 26, 81-93, doi:10.1007/BF00164332, 1983.
- 737 Serafimovich, A., Metzger, S., Kohnert, K., Hartmann, J., and Sachs, T.: The airborne
738 measurements of methane fluxes (AIRMETH) arctic campaign, 46th AGU annual Fall Meeting,
739 San Francisco, U.S.A., 9 - 13 December, 2013.
- 740 Smith, D. E., Metzger, S., and Taylor, J. R.: A transparent and transferable framework for
741 tracking quality information in large datasets, *PLoS One*, 9, e112249,
742 doi:10.1371/journal.pone.0112249, 2014.



- 743 Starkenburg, D., Metzger, S., Fochesatto, G. J., Alfieri, J. G., Gens, R., Prakash, A., and
744 Cristóbal, J.: Assessment of de-spiking methods for turbulence data in micrometeorology, *J.*
745 *Atmos. Oceanic Technol.*, doi:10.1175/jtech-d-15-0154.1, 2016.
- 746 Stull, R. B.: *An Introduction to Boundary Layer Meteorology*, Kluwer Academic Publishers,
747 Dordrecht, The Netherlands, 670 pp., 1988.
- 748 Sulkava, M., Luysaert, S., Zaehle, S., and Papale, D.: Assessing and improving the
749 representativeness of monitoring networks: The European flux tower network example, *J.*
750 *Geophys. Res.*, 116, G00J04, doi:10.1029/2010jg001562, 2011.
- 751 Taylor, J. R., and Loescher, H. L.: Automated quality control methods for sensor data: A novel
752 observatory approach, *Biogeosciences*, 10, 4957-4971, doi:10.5194/bg-10-4957-2013, 2013.
- 753 Tetzlaff, A., Lüpkes, C., and Hartmann, J.: Aircraft-based observations of atmospheric
754 boundary-layer modification over Arctic leads, *Q. J. R. Meteorol. Soc.*, 141, 2839-2856,
755 doi:10.1002/qj.2568, 2015.
- 756 Turner, D. P., Ollinger, S. V., and Kimball, J. S.: Integrating remote sensing and ecosystem
757 process models for landscape- to regional-scale analysis of the carbon cycle, *BioScience*, 54,
758 573-584, doi:10.1641/0006-3568(2004)054[0573:irsaep]2.0.co;2, 2004.
- 759 Tuzson, B., Hiller, R. V., Zeyer, K., Eugster, W., Neftel, A., Ammann, C., and Emmenegger,
760 L.: Field intercomparison of two optical analyzers for CH₄ eddy covariance flux measurements,
761 *Atmos. Meas. Tech.*, 3, 1519-1531, doi:10.5194/amt-3-1519-2010, 2010.
- 762 Vaughan, A. R., Lee, J., Misztal, P., Metzger, S., Shaw, M. D., Lewis, A. C., Purvis, R., Carslaw,
763 D., Goldstein, A., Hewitt, C. N., Davison, B., Beevers, S. D., and Karl, T.: Spatially resolved
764 flux measurements of NO_x from London suggest significantly higher emissions than predicted
765 by inventories, *Faraday Discuss.*, doi:10.1039/c5fd00170f, 2015.
- 766 Vickers, D., and Mahrt, L.: Quality control and flux sampling problems for tower and aircraft
767 data, *J. Atmos. Oceanic Technol.*, 14, 512-526, doi:10.1175/1520-
768 0426(1997)014<0512:QCAFSP>2.0.CO;2, 1997.
- 769 Webb, E. K., Pearman, G. I., and Leuning, R.: Correction of flux measurements for density
770 effects due to heat and water vapour transfer, *Q. J. R. Meteorol. Soc.*, 106, 85-100,
771 doi:10.1002/qj.49710644707, 1980.
- 772 Wilczak, J. M., Oncley, S. P., and Stage, S. A.: Sonic anemometer tilt correction algorithms,
773 *Boundary Layer Meteorol.*, 99, 127-150, doi:10.1023/A:1018966204465, 2001.
- 774 Xu, K., Metzger, S., and Desai, A. R.: Upscaling tower-observed turbulent exchange at fine
775 spatio-temporal resolution using environmental response functions, *Agric. For. Meteorol.*, 232,
776 10-22, doi:10.1016/j.agrformet.2016.07.019, 2017.
- 777
- 778