1 **Toward cyberinfrastructure to facilitate collaboration and reproducibility for marine**

2 **Integrated Ecosystem Assessments**

3

4 Stace E. Beaulieu[1]*, Peter A. Fox[2], Massimo Di Stefano[1,2,3], Andrew Maffei[1], Patrick West[2],

5 Jonathan A. Hare[4], Michael Fogarty[4]

6

7 [1] Woods Hole Oceanographic Institution, Woods Hole, MA 02543 USA

8 [2] Tetherless World Constellation, Rensselaer Polytechnic Institute, Troy, NY 12180 USA

9 [3] Center for Coastal and Ocean Mapping, University of New Hampshire, Durham, NH 03824

10 USA

11 [4] Northeast Fisheries Science Center, National Oceanic and Atmospheric Administration, Woods

12 Hole, Massachusetts 02543 USA

13 * stace@whoi.edu, +1-508-289-3536

14

15 20 October 2016

17  **Abstract**

18  There is a growing need for cyberinfrastructure to support science-based decision making in

19  management of natural resources. In particular, our motivation was to aid the development of

20  cyberinfrastructure for Integrated Ecosystem Assessments (IEAs) for marine ecosystems. The

21  IEA process involves analysis of natural and socio-economic information based on diverse and

22  disparate sources of data, requiring collaboration among scientists of many disciplines and

23  communication with other stakeholders. Here we describe our bottom-up approach to developing

24  cyberinfrastructure through a collaborative process engaging a small group of domain and

25  computer scientists and software engineers. We report on a use case evaluated for an Ecosystem

26  Status Report, a multi-disciplinary report inclusive of Earth, life, and social sciences, for the

27  Northeast U.S. Continental Shelf Large Marine Ecosystem. Ultimately, we focused on sharing

28  workflows as a component of the cyberinfrastructure to facilitate collaboration and

29  reproducibility. We developed and deployed a software environment to generate a portion of the

30  Report, retaining traceability of derived datasets including indicators of climate forcing, physical

31  pressures, and ecosystem states. Our solution for sharing workflows and delivering reproducible

32  documents includes IPython (now Jupyter) Notebooks. We describe technical and social

33  challenges that we encountered in the use case and the importance of training to aid the adoption

34  of best practices and new technologies by domain scientists. We consider the larger challenges

35  for developing end-to-end cyberinfrastructure that engages other participants and stakeholders in

36  the IEA process.

37

38  **Keywords**

39    e-Science, executable workflow, indicator, IPython Notebook, open science, use case

40    methodology

41

50

51

52

## Introduction

There is a growing need for cyberinfrastructure to support science-based decision making in management of natural resources (e.g., Acreman 2005; Reichman et al. 2011; Palmer 2012; Muste et al. 2013; Horsburgh 2015). Over the past decade the U.S. has moved toward an ecosystem-based management approach for marine ecosystems, and there is a need for development of cyberinfrastructure to support the science teams who are reporting on these ecosystems and provisioning services such as fisheries. We were motivated to develop cyberinfrastructure to provide a transparent pathway from data to knowledge to action, responding to the U.S. National Ocean Policy Implementation Plan, in particular "improving science-based products and services for informed decision-making" (National Ocean Council 2013). Here, we define cyberinfrastructure as infrastructure that comprises "*both technology and human expertise necessary to support scientific research processes and collaboration*" (Jirotka et al. 2013). Levin et al. (2009, 2014) and Samhouri et al. (2014) describe a formal process for an Integrated Ecosystem Assessment (IEA), involving natural and social scientists working together to assess a marine ecosystem with respect to management objectives (Fig. 1). Data collected, integrated, and interpreted in a marine IEA may be as diverse as climate indices, satellite-derived sea surface temperature, counts of phyto- and zooplankton from net tows, and landings data from commercial fisheries.

For any coupled natural and human system it is challenging to develop cyberinfrastructure to enable multi- and inter-disciplinary research to understand, model, and make predictions for the system as a whole. Technical challenges include handling, integrating, analyzing, and tracking provenance of very heterogeneous data (e.g., Reichman et al. 2011). In an IEA to make sense of a plethora of data, it is common practice to focus on a select subset of indicators of natural or

76   anthropogenic drivers or ecosystem states that can be monitored for changes over time and space

77   (Samhouri et al. 2012). Indicators tend to be derived datasets and are often "synthesized

78   products" (term used in NOAA 2014), resulting from complex data processing workflows that

79   integrate not only data and models but also subjective choices made by scientists based on

80   knowledge in their domain. Social challenges include scientists of different domains using

81   different terms to describe their data and different software and tools to work with data (e.g.,

82   Pennington 2011; Cooke and Hilton 2015). E-Science teams inclusive of scientists and

83   information technology (IT) experts face the additional challenge that "IT experts cannot

84   understand the needs of the scientists – and scientists cannot understand what is even possible –

85   without conceptual integration between the scientists and IT experts" (Pennington 2011).

86   Here we report on the ECO-OP (an abbreviation joining ECOsystem and interOPerability)

87   project involving fisheries scientists, oceanographers, computer scientists, information modelers,

88   and software developers. As part of this project, we identified and conducted a use case to

89   support the bi-annual generation of an Ecosystem Status Report (hereinafter the Report) as part

90   of an IEA for the Northeast U.S. Continental Shelf Large Marine Ecosystem. The Report is

91   composed of chapters, each of which is prepared by different specialists for climate forcing,

92   physical pressures, primary and secondary production, benthic invertebrates, fish communities,

93   protected species, anthropogenic factors, and integrated ecosystem measures (Ecosystem

94   Assessment Program 2012). The software framework to be developed needed to enable these

95   different specialists to process heterogeneous data and provide products for the Report. The

96   framework would be flexible to allow for addition and subtraction of indicators from the Report

97   and portable to accommodate assessment of marine ecosystems in other managed regions of the

98   ocean.

99   The ECO-OP project addressed challenges in developing cyberinfrastructure for e-Science teams

100  participating in marine IEAs. Following our definition of cyberinfrastructure above, our use case

101  for the Report involved integrating *technologies* ranging from data sharing (including access and

102  re-usability) to executable workflows and *human expertise* including knowledge and practices in

103  multiple natural and social science domains. In the spirit of open science (Reichman et al. 2011;

104  Nosek et al. 2015), we aimed beyond transparency toward the reproducibility standard in the

105  U.S. NOAA Information Quality Guidelines (NOAA 2014) for indicators and other data

106  products in the Report. Below, we describe the software prototype that we developed and how

107  we aided its adoption by the scientists producing the Report. We discuss how to scale the

108  prototype and other considerations for the larger cyberinfrastructure to be developed for the IEA

109  process.

110

111  **Methods**

112  *Methodology to develop cyberinfrastructure and evaluate the use case*

113  We employed a bottom-up approach in which a small team with diverse skills worked closely to

114  evaluate use cases with very specific goals as representative of a larger set of goals. This

115  approach engages domain scientists directly in the collaborative development of a software

116  solution. The use cases were iteratively developed to articulate specific goals of fisheries

117  scientists delivering indicators and data products, capture detail on what went into reaching those

118  goals, and the outcomes they needed to evaluate success. Computer scientists and software

119  developers provided options for technologies which were then evaluated to determine how they

120  could be adopted and then how they could be incorporated into a larger framework of

121  cyberinfrastructure. In addition to engaging with fisheries scientists in the use case evaluation,

122    informatics and software experts in the small team also regularly attended science meetings to

123    learn more about the science, understand concepts, share ideas, and build trust. This

124    methodology is in contrast to top-down approaches that prescribe technologies for domain

125    scientists as end users.

126    The use case for the Report explored options for the portion of the IEA process including

127    "Develop Indicators," "Monitoring of Ecosystem Indicators," and "Assess Ecosystem" (Fig. 1).

128    We provide a diagram as an overview of the data-level and application-level mediation

129    requirements to compile the Report (Fig. 2). We also show representative temporal and spatial

130    indicators as derived data products in the Report (Fig. 3). We evaluated the use case through the

131    Tetherless World Constellation (TWC) Semantic Web Methodology (hereafter, TWC

132    Methodology), a collaborative process of rapid prototyping based on a small team including

133    domain scientists (Fox and McGuinness 2008). Essentially, the small team was a subset of a

134    larger e-Science team collaborating on a prototype Report. The TWC Methodology is a cycle

135    involving ten stages (Fig. 4):

136    (1) The use case defines the interactions between people, hardware, software, and desired

137    products and can be adjusted or refined after each iteration of the cycle. The initial goal of the

138    use case for the Report was to efficiently generate figures representing ecosystem data and

139    information products; this goal was expanded to be inclusive of generating the Report documents

140    [portable document format (PDF) and associated webpages].

141    (2) The small team with mixed skills met initially to define the use case and then subsequently

142    (in stage 10 described below) to evaluate each prototype to complete an iteration of the cycle.

143    The authors of this paper comprise the team for the use case: facilitator (Fox, Maffei), domain

144    experts [Hare, Fogarty, and other scientists in the Ecosystem Assessment Program at NOAA's

145   Northeast Fisheries Science Center], knowledge representation and information modeling

146   (West), software engineering (Di Stefano), and scribe (Beaulieu). The larger group of fisheries

147   scientists contributing to the Report comprises ~40 individuals working at ~10 different NOAA

148   offices and academic institutions.

149   (3) Analysis of the use case included identifying the actors and source data, writing a narrative

150   description, outlining a flow, and drawing an activity diagram (Fig. 5). Expectations ultimately

151   were refined to the following: The framework should retrieve data, report quality

152   assurance/quality control, conduct standard analyses, provide iterative and interactive

153   visualization, allow for interpretation, and generate final graphics to embed into webpages and

154   PDF. In addition, the data represented in each figure should be available. The framework should

155   also document the specific process for each data and information product, including source data,

156   code, and related contextual information suitable for traceability, repeatability, explanation,

157   verification, and validation. The framework should use the same components/structure for each

158   data and information product, thereby allowing the addition and subtraction of data and

159   information products in future Reports.

160   (4) Neither an information model nor ontology was formally developed in the Report use case.

161   However, we explored and mapped concepts that were important to document as metadata, due

162   to different terms being used by different actors in the use case. In this project our use of

163   "semantics" in the TWC Methodology involved "developing shared conceptualizations across

164   disciplinary boundaries" *sensu* Pennington (2011).

165   (5) The TWC Methodology advocates finding and using relevant tools; thus, we tested a number

166   of existing open source tools as we iterated the prototype including Drupal, Wt (the C++ Web

167   Toolkit), and the IPython (now Jupyter) Notebook (Pérez and Granger 2007; Ragan-Kelley et al.

168    2014; Shen 2014). In particular, the IPython Notebook is an "interactive computational

169    environment" with a web application and "notebooks, for recording and distributing the results of

170    the rich computations" (https://github.com/ipython/ipython-

171    website/blob/b578013e545d18deafa0f9e1567e3db5368f0cf6/notebook.rst l, accessed 17 October

172    2016).

173    (6) Science/expert reviews occurred within each iteration of the cycle as the prototype was being

174    developed for the next major group evaluation.

175    (7 & 8) We adopted technologies that were available as open source and leveraged the

176    technology infrastructure (hardware and software) that the fisheries scientists were already using

177    to generate indicators. Cooke and Hilton (2015) provide a comprehensive list of factors to

178    consider when selecting technologies for e-Science teams (e.g., ease of use, accessibility,

179    security, compatibility).

180    (9) The initial rapid prototype acted "to glue the components together and connect them to

181    interfaces and visualization tools. ...latter stages of the prototype must pay increasing attention to

182    non-functional aspects of the use case, such as scalability, reliability, etc." (Fox and McGuinness

183    2008).

184    (10) The final stage is evaluation of the prototype to determine whether/how it should be

185    redesigned and redeployed. In practice this stage involves demonstration of the software

186    prototype to the larger e-Science team and then an evaluation by the small team to complete the

187    iteration of the cycle.

188    We developed prototypes for the Report use case during three complete iterations of the TWC

189    Methodology. Each iteration of the cycle took a few to several months, accounting for the time to

190     develop and test software, and demonstrate and evaluate each prototype. The fisheries scientists

191     requested transfer of the technologies after demonstration of the third iteration prototype, which

192     focused on the "Climate Forcing" and "Physical Pressures" chapters in the Report (Ecosystem

193     Assessment Program, 2009). Prior to the delivery to fisheries scientists, the small team

194     conducted three small "spin-off" use cases to further test the software prototype. These small use

195     cases were intended to examine whether the prototype that was successful for one portion of the

196     Report could also be adapted for indicators and data products from other chapters in the Report

197     (Ecosystem Assessment Program, 2012). We delivered the prototype software environment to

198     the fisheries scientists in two ways: in a virtual machine (VM) provided to individuals, and by

199     installation on a server at the Narragansett facility with the aid of NOAA's IT staff.

200

201     *Training to aid adoption of the technologies*

202     During each iteration of the cycle described above, the e-Science team gains some exposure to

203     the cyberinfrastructure inclusive of technologies and others' expertise, but it is mainly the small

204     team that gains hands-on experience with the software prototype. Additional training and hands-

205     on experience is desired to aid adoption of the technologies by the larger team. We provided

206     training opportunities and technical support in groups and for individuals, as recommended by

207     Cooke and Hilton (2015). In the first iteration prototype, fisheries scientists were introduced to

208     several applications that were new to them: interactive programming software (IPython

209     Notebook), version control software (Subversion), and content management systems (including

210     Trac and Drupal). Ultimately we focused the training on IPython Notebook and changed to

211     version control with GitHub. We offered three group training workshops, two of which were

212     specific to ECO-OP cyberinfrastructure. The first workshop, which involved the second iteration

213     prototype, was essentially an introduction to IPython Notebooks utilizing a shared online server

214     that the e-Science team logged into as users. During the one-day workshop and for a few months

215     afterward (as we were conducting the third iteration of the use case), users were provided folders

216     on the shared server to store their notebooks and data products. The second workshop was

217     provided after we completed the final prototype and was aimed towards learning Python

218     programming and best practices for version control. This training involved a two-day Software

219     Carpentry Bootcamp (Wilson 2014) held at Northeast Fisheries Science Center and was also

220     open to other fisheries scientists. The third workshop was to assist the e-Science team in using

221     the final prototype - i.e., ECO-OP pyecoop software library distributed within a VM - to generate

222     data products specific to their chapters of the Report. The purpose of this final training over 2.5

223     days was to assist with user-specific, individual needs (we asked participants to come with their

224     own data and code).

225

226     **Results**

227     *Initial prototypes*

228     As a first step towards developing the prototype Report, the small team sketched an activity

229     diagram which identified the primary actors in the collaboration, including many people (e.g.,

230     data preparation reviewer, Report compiler/editor) and a software agent  (Fig. 5). Pre-conditions

231     for the use case included that source data are accessible. The basic flow for the use case may be

232     described as: Source data are retrieved > Source data are processed into preliminary data

233     products (which are stored) > Intermediate and final data products including indicators are

234     calculated, analyzed, and plotted in an iterative and interactive process (and stored) > Indicators

235     are interpreted > Text is written for context, interpretation, and synthesis > Report is compiled

236   (and stored). Post-conditions for the use case, not explicitly addressed in the prototype, included

237   storage and archiving of the preliminary, intermediate, and final data and visualization products

238   and the Report itself.

239   During the first two iterations of the TWC Methodology, we were developing multiple software

240   prototypes corresponding to different components of the desired cyberinfrastructure. The first

241   iteration prototype targeted software tools for data access, data processing, metadata acquisition,

242   and data visualization. We focused on the first two chapters in the Report, "Climate Forcing"

243   which included climate indices [e.g., North Atlantic Oscillation; Fig. 2.1 in the 2009 Report

244   (Ecosystem Assessment Program, 2009)] and "Physical Pressures" which included sea surface

245   temperature anomalies [e.g., Fig. 3.5 in the 2009 Report (Ecosystem Assessment Program,

246   2009)]. The first iteration prototype separately considered a tool for data access and processing

247   (IPython Notebook), tools for manual contribution of metadata in controlled vocabularies (Trac

248   and Drupal), and other web applications for interactive display of final datasets. In practice, we

249   utilized IPython Notebooks to output comma-separated value files for time-series indicators, we

250   manually input metadata for these indicators to other file formats, we stored the data and

251   metadata files at specific addresses, and the web applications called to these addresses to display

252   one or more indicators. As a result of the evaluation of the first iteration prototype, the fisheries

253   scientists were intrigued but not comfortable with IPython Notebook, mainly because this first

254   demo involved converting code from one programming language (MATLAB) to another

255   (Python) [not necessary in further iterations due to the availability of a Python-MATLAB bridge

256   (and, now, also a Matlab kernel for Jupyter; Jupyter Team 2015)]. The fisheries scientists were

257   not keen to learn tools to manually contribute metadata and requested that we focus on

258   automated acquisition of metadata. They also requested that we further customize a web

259 application for interactive display of the indicators. In response the small team sketched a

260 Graphical User Interface (GUI) with a drop-down list to select indicators, more options for

261 plotting, and buttons for exporting data and visualization products, viewing metadata, and saving

262 a session.

263 For the second iteration prototype we built a web-app GUI using Wt that could be displayed on

264 its own or within an IPython Notebook. We recorded a demo to show the larger e-Science team

265 how to use the web-app GUI for interactive display of the indicators and how to log in and use

266 both the IPython Notebook and the web-app GUI to re-calculate an indicator with the latest

267 version of code, then store and display the final data file. To support this human-oriented process

268 we implemented a shared server to contain the development environment and allow for easy

269 sharing of notebook files and the output data files, images, and PDFs. Converting notebooks into

270 PDFs was a key new development made possible with the nbconvert tool, which also handles

271 other formats including HTML and LaTeX (Frederic 2013). We continued to focus on indicators

272 in the "Climate Forcing" and "Physical Pressures" chapters of the Report but also performed

273 workflows using IPython Notebooks for ecosystem indicators, including a phytoplankton

274 abundance anomaly (Di Stefano et al., 2012) and time series of copepod abundance [Fig. 4.10 in

275 the 2009 Report (Ecosystem Assessment Program, 2009)].

276 To evaluate the second iteration prototype, we distinguished three levels of users: users of an

277 interactive PDF for the Report with hyperlinks to data and metadata (Level 1), users of the web-

278 app GUI to access final data products (Level 2), and users interacting with IPython Notebooks

279 (Level 3). A major result of the evaluation was that the fisheries scientists aspired to become

280 Level 3 users and asked to have an IPython Notebook tutorial as soon as possible. The overall

281 assessment was that the IPython Notebook technology offered the most flexibility for

282    calculating, analyzing, and plotting indicators for the Report and would also enable the

283    production of an interactive PDF. The fisheries scientists requested that we explore further the

284    conversion of notebooks to HTML, as the group was considering providing the Report directly

285    online as a website. Essentially, the IPython Notebook appeared to be a single tool that could

286    accommodate components considered separately in the first iteration prototype.

287

288    *Final prototype*

289    The third prototype focused on the IPython Notebook tool and ultimately was refined to the final

290    prototype delivered to fisheries scientists. Much of the development in the third iteration of the

291    use case involved building a software library for processing, analyzing, and visualizing

292    indicators in IPython Notebooks and an environment to accommodate all the dependencies. Our

293    first "spin-off" use case was to test the conversion of an IPython Notebook to an Ecosystem

294    Advisory webpage. We used a notebook created in the first iteration prototype for the "Physical

295    Pressures" chapter to successfully reproduce a webpage in HTML format for long-term

296    temperature trends in the Northeast U.S. Shelf ecosystem (Di Stefano et al., 2013). The

297    demonstration of the third iteration prototype included this simulated Ecosystem Advisory

298    webpage and a notebook (Fig. 6) that retrieved and processed data for two climate indicators and

299    output an interactive PDF (Fig. 7) formatted to look exactly like a portion of the "Climate

300    Forcing" chapter in the Report (Ecosystem Assessment Program, 2009). This notebook (Fig. 6),

301    which requires the installation of TeX Live [TeX distribution for several Linux distributions

302    (https://www.tug.org/texlive/)] into the environment, utilizes the pdflatex command to compile

303    text files with image files created on-the-fly as a result of data visualization in the notebook. The

304    interactive PDF (Fig. 7) included embedded links to data files plotted in the figures.

305    As a result of the evaluation of the third prototype, the fisheries scientists determined that the

306    expectations for the use case were met. However, prior to the transfer of technologies, they

307    requested that we address some of the challenges in reproducing other chapters of the Report.

308    Our second and third "spin-off" use cases examined challenges in reproducing the workflows for

309    a fisheries indicator (Fig. 3a) and a map of primary production (Fig. 3b) from other chapters in

310    the Report (Ecosystem Assessment Program, 2012). For both of these use cases, our goal was to

311    determine whether a complex workflow utilizing many data sources, multiple tools, and multiple

312    programming languages could be accommodated with an executable workflow in an IPython

313    Notebook. We worked directly with the fisheries scientists responsible for these data products in

314    the Report to determine the earliest point at which the prototype developed for the Report use

315    case (dashed box in Fig. 5) could apply to their respective workflows. The fisheries indicator is

316    constructed by a natural scientist and a social scientist working together. Their workflow had a

317    number of manual steps in accessing multiple data sources and preparing preliminary data,

318    including the use of a manual data query extraction tool. However, the remainder of the

319    workflow involving these preliminary data products could be conducted within an IPython

320    Notebook with an extension for the R programming language (now, an R kernel for Jupyter;

321    Jupyter Team 2015). The map of primary production is constructed by one scientist and involves

322    an even more complex workflow that starts with accessing thousands of source data files. The

323    scientist utilizes SeaDAS (http://seadas.gsfc.nasa.gov) tools and Interactive Data Language

324    (IDL) to process data and construct the map image. At the time although SeaDAS tools could be

325    implemented in a Python environment, there was no extension for IDL in IPython Notebook.

326    Today, Jupyter has an IDL kernel (Jupyter Team 2015), and the scientist should be able to create

327    a notebook to execute the complete workflow from source data retrieval to outputting a figure for

328    the Report, without having to convert code into Python.

329    The final prototype was a software environment for Linux operating systems inclusive of a

330    software library with general utility to enable the reproducibility of scientific workflows that

331    acquire data online, process and plot data, and package text and figures into a document.

332    Workflows are conducted within IPython Notebooks. The ECO-OP pyecoop software library is

333    available at a GitHub repository with GNU Lesser General Public License, accessible via

334    https://data.rpi.edu/xmlui/handle/10833/1756. The pyecoop software library, written in Python

335    (>=2.7 , >=3.3), has several modules including a module with utility functions (ecoop.ecooputil)

336    and a module that defines methods for data in the "Climate Forcing" chapter of the Report

337    (ecoop.cf). Dependencies for the pyecoop code include the installation of TeX Live and

338    RubyGems (https://rubygems.org/). Other Python libraries are required, including matplotlib

339    (Hunter 2007), pandas (McKinney 2010), and scipy (Jones et al. 2001). The software

340    environment includes IPython Notebook and other open source applications used in generating

341    indicators and documents, such as Geographic Resources Analysis Support System (GRASS

342    Development Team 2015), Octave (Eaton et al. 2014), and R (R Core Team 2013). The software

343    environment was distributed within a VM (important for when users are not online) and by

344    installing a single-port instance on a server at NOAA's Narragansett facility. Ultimately the

345    components of the delivered cyberinfrastructure included software and human resources

346    (including training described below) but excluded hardware resources. We did not prescribe data

347    storage or archiving, and the Report use case did not require support for high performance

348    computing (this may be required for other use cases involving ecosystem modeling).

349

*Results of training to aid adoption of the technologies*

We provide some results for our first and third group training opportunities which were specific to ECO-OP cyberinfrastructure; however, we did not conduct surveys or interviews for a more rigorous evaluation of the training. Thirteen fisheries scientists participated at the first workshop. The most positive result was that one month after the training, one of the fisheries scientists was using IPython Notebook to develop and document new indicators, utilizing extensions to enable functionality for other programming languages. Upon seeing these new notebooks, another fisheries scientist joined the shared server (available in the second prototype) as a new user and aided the development of the notebook for the Ecosystem Advisory webpage that was part of our third prototype demonstration. Eight fisheries scientists participated at the third workshop; six did not attend the first training which placed them at a disadvantage since we assumed some familiarity with IPython Notebooks. At least one attendee was able to generate a PDF with their own data and code. All attendees left the workshop with the software requirements installed and configured in a VM on their own laptops. The environment provided to each attendee with the VM was fully compatible with the software infrastructure installed on the server at NOAA's Narragansett facility. Comparing these two training opportunities, the first appeared to be more successful with the single shared software environment; we think that we lost users when each distribution was installed separately as a VM, not only due to challenges in the installation but also in terms of having to use email or other shared storage services to share notebooks. Importantly, the training was of benefit not just to the users, but also to the small team developing the software environment, to observe the challenges expressed by domain scientists with a range of skills. The first training session aided development during the third iteration of

372    the use case. The third training session was conducted after deciding upon the final prototype and

373    helped us with documentation prior to delivery.

374

375    **Discussion**

376    *Solution for sharing workflows and delivering reproducible documents*

377    Our solution for the fisheries scientists to reproduce a portion of their Report was a software

378    environment in which IPython Notebook acted as a lightweight, flexible, re-usable, scientific

379    workflow technology to document data processing, analyses, visualization, and reporting. The

380    solution is in the spirit of open science in which the sharing of workflows engenders trust in the

381    derived data products (Reichman et al. 2011; Nosek et al. 2015; Wright 2016). We recognize that

382    the delivered prototype, which reproduced a portion of the "Climate Forcing" chapter in the

383    Report (Fig. 7) and accommodated workflows for a variety of other ecosystem indicators, only

384    addressed a limited set of technical and social challenges involved in preparing and compiling

385    the Report. We addressed many challenges in terms of software required to execute the

386    workflows (e.g., use of different programming languages, integrating with open source software

387    libraries); however, we were not able to fully address challenges in the sharing of these

388    workflows. We did not go so far as to enable a repository, management system, or social

389    network for the sharing of workflows (e.g., Goble et al. 2010; Liu et al. 2015). Ultimately we

390    were limited in implementing a shared file system in the final prototype, although this may be

391    more straightforward to develop today due to recent developments for multi-user servers for

392    notebooks (e.g., Wakari, JupyterHub).

393     We successfully reproduced a portion of one chapter and additional indicators, but an ultimate

394     goal would be to enable a Report "on-demand" (at the time of this project, production of the

395     Report was manually intensive and limited to every two years). Many technical and social

396     challenges arise when considering the compilation of the entire Report as a reproducible

397     document, a reason why we drew this step outside of the dashed box in the activity diagram (Fig.

398     5). A major challenge at this time would be the accessibility of source data for the many data

399     processing workflows. For reproducibility in the future, the cyberinfrastructure would also need

400     to account for versioning of IPython Notebooks for each data visualization product. The main

401     technical challenge that we highlight here is sustaining a computational infrastructure for all of

402     the e-Science team members' software environments and dependencies inclusive of

403     repository(ies) with version control. This assemblage of very dynamic and distributed software

404     environments is analogous to a "scientific software ecosystem" in recent publications (e.g.,

405     Howison et al. 2015). In addition, to reproduce all of the chapters, all of the fisheries scientists

406     would need to adopt new technologies, which we address below.

407

408     *Training to aid adoption of the technologies*

409     Our experience with fisheries scientists provides a specific example of the general importance of

410     training and professional development when selecting technologies to support multi-disciplinary

411     e-Science teams (e.g., Cooke and Hilton 2015). We recognized with the initial prototypes that

412     training would be central to our success in transferring the software environment to fisheries

413     scientists. One measure of success for our delivered prototype is how the fisheries scientists used

414     the technologies for their subsequent Report and other work conducted for the IEA process. We

415     expected our bottom-up/user-driven approach to promote adoption of technologies based on

416 research "finding that technical systems that were well aligned with and ready to accomplish the

417 task scientists intended were more likely to be successfully adopted by the community" (Olson et

418 al. 2008). Ultimately, only a few fisheries scientists utilized the prototype to produce portions of

419 the subsequent Report. This may in part be due to technology readiness for the scientists (e.g.,

420 many had never interacted with a Linux operating system, and/or had no experience with the

421 Python programming language). As noted by the iMarine project described in the next section,

422 "in the domain of fisheries, marine biology and environmental sciences... users and researchers

423 generally lack advanced IT skills" and "it is important to bear in mind the time to learn to use

424 new tools" (iMarine 2014). Additional consultation and/or continued training was needed for

425 fisheries scientists to build on and extend our prototype to produce chapters for the next Report.

426 Pennington (2011) describes additional factors that influence technology adoption that may have

427 been factors in our project, e.g., extrinsic motivation (which would be more applicable in a top-

428 down approach).

429 In the long-term, perhaps more important than training to adopt specific technologies, our

430 training encompassed best practices that were new to many of the scientists. Because

431 technologies change frequently it is important for training to "generalise to broader classes of

432 technologies and the socio-technical arrangements to which they point" (Jirotka et al. 2013).

433 Including the Software Carpentry Bootcamp our training opportunities may be considered an

434 attempt to grow the culture of best practices for data and software management in the community

435 in which fisheries scientists work. Our training led to the broader use of open source tools and

436 version control by scientists at the Northeast Fisheries Science Center. However, to build e-

437 Science teams for new applications, there needs to be continued interaction with computer

438 scientists, software engineers, and other IT experts.

439

*Comparing our approach to other efforts to develop cyberinfrastructure for e-Science teams in IEAs*

Our project involved a bottom-up approach in which a small team addressed very specific use cases as representative of a larger body of collaborative work for marine IEAs. The approach also involved the informatics and software experts engaging with domain scientists at their regular meetings to improve understanding of concepts and to develop relationships and trust in addition to the targeted use cases. At the end of each cycle of the TWC Methodology the small team shared the latest prototype with the larger e-Science team, thus directly involving end users in the evaluation. We aspired to prototype a software environment that would enable the flexibility for these end users to also become developers, re-shaping and expanding the software environment as needed to accommodate more data and information products in the Report. This lack of "clear delineations between users and developer" has been recognized in general for the development of technologies and infrastructure for e-Science teams (Jirotka et al. 2013). Our bottom-up approach is aligned with the Computer Supported Cooperative Work "focus on the scientists' everyday work practices, with a view to enabling new collaborations" (Jirotka et al. 2013), very much focused on the individual scientist and how s/he collaborates with other scientists contributing to an IEA.

Our approach is much smaller in scale than efforts that we highlight below from the European Union and Australia that also are directed toward cyberinfrastructure for IEAs. The European iMarine project is described as "an open and collaborative initiative aimed at supporting the implementation of the Ecosystem Approach to fisheries management" (http://www.i-marine.eu/Pages/Home.aspx, accessed 31 December 2015). Many of the goals of iMarine are

462  similar to the ECO-OP project, including "facilitated retrieval, access, collaborative production

463  and sharing of information and tools" (http://www.i-marine.eu/Pages/Home.aspx, accessed 31

464  December 2015). To achieve these goals iMarine provides web-based virtual research

465  environments (VREs) through domain-specific infrastructure built onto D4Science e-

466  infrastructure, "a virtual aggregator of resources available in interoperable e-infrastructures"

467  (Taconet et al. 2014). Our interpretation is that scientists are users of the platform although they

468  may be developers of workflows incorporated into the platform. As a future research effort we

469  recommend exploring how to incorporate the ECO-OP prototype inclusive of executable

470  workflows in IPython Notebooks into the iMarine platform.

471  For Australia we highlight the eReefs project, built upon "an innovative central information

472  infrastructure reflecting best practice in environmental information management"

473  (http://ereefs.org.au/ereefs/platform, accessed 31 December 2015). We draw an analogy between

474  our Report use case and the "Report Card" of the eReefs Platform

475  (http://ereefs.org.au/ereefs/platform, accessed 15 April 2016). In our use case we explored the

476  use of a scientific workflow tool to account for processing source observational and model data

477  into data visualization products, similar to the eReefs pilot (however, they used a proprietary

478  tool; Chen et al. 2011). The ECO-OP project accounted for additional heterogeneity and issues of

479  interoperability by addressing additional "spin-off" use cases and through a provenance use case

480  described elsewhere (Ma et al. 2017). The current eReefs project (2012 - 2017) is intended to

481  develop an information architecture to "allow for the next generation of data interoperability by

482  augmenting established, standardised, services and allowing for the integration of multi-service

483  use" (Car 2013). As a future research effort we also recommend exploring how to incorporate the

484  ECO-OP prototype into the eReefs Platform.

485    We recognize that some of the challenges in scaling up and out when developing

486    cyberinfrastructure with a bottom-up approach, differ from top-down development efforts. Top-

487    down efforts may enforce policies or encourage the removal of technical or social barriers that

488    inhibit broad usage of collaborative tools. However, although the ECO-OP project only

489    addressed a small portion of the overall cyberinfrastructure that would be implemented within a

490    VRE, we see most if not all of the socio-technical issues we considered critical to the success of

491    our use case also applying to VREs (i.e., Jirotka et al. 2013, their sxn. 4.2). Our bottom-up

492    approach in which the scientists (as end users of the infrastructure) are participating directly in

493    the development of the infrastructure, was a nimble and rapid means to achieve the prototype

494    Report. Our approach aligns with the concepts of "vertical user stories" in agile software

495    development (e.g., Pulsifer et al. 2011) and participatory design (or co-design) in socio-technical

496    systems (Muller and Kuhn 1993). Moreover, the adaptation of a more agile and iterative, i.e.,

497    quicker, sequence of try, evaluate, and revise indicates that future efforts to develop

498    cyberinfrastructure for e-Science teams in IEAs (but also more generally) consider incorporating

499    an agile approach or the small team/TWC Methodology as a means to supplement the larger

500    development process.

501

502    *Toward end-to-end cyberinfrastructure for the IEA process*

503    The work conducted by scientists in the IEA process is embedded within a larger process

504    involving other stakeholders in ecosystem-based management (Fig. 1). An ultimate goal is to

505    extend the cyberinfrastructure developed for e-Science teams to address challenges at the

506    science-policy interface including "... communication and debate about assumptions, choices and

507    uncertainties, and about the limits of scientific knowledge" (van den Hove 2007). Essentially,

508  cyberinfrastructure for the IEA process should encompass a virtual organization (*sensu* Ahuja

509  and Carley 1998) of diverse stakeholders including scientists, decision makers, and the public.

510  Our work in this project is just one example of the growing need for cyberinfrastructure to

511  support science-based decision making in management of natural resources (e.g., Acreman 2005;

512  Reichman et al. 2011; Palmer 2012; Muste et al. 2013; Horsburgh 2015). Our vision was to

513  facilitate the engagement of natural and social scientists in routine ecosystem assessments, yet

514  we aspire to involve other stakeholders through presenting robust science data in forms that

515  various end users can consume and verify. This vision is shared by others developing

516  cyberinfrastructure for IEAs including iMarine (Taconet et al. 2014) and eReefs (Car 2013).

517  The ECO-OP project provided a pilot toward end-to-end transparency starting from a scientist's

518  desktop and being shared with collaborators, to a report provided to managers, policy makers,

519  and the public. IPython Notebooks can be used as electronic lab notebooks, whereby scientists

520  digitally record the steps involved in their computations and ultimate data products (Shen 2014).

521  These notebooks essentially document a provenance chain, especially useful for indicators that

522  summarize large collections of underlying heterogeneous data. Our solution included interactive

523  and transparent workflows of data analysis and delivery of a reproducible document, but did not

524  represent provenance in a machine-readable standard. After completing the use case with

525  fisheries scientists described in this paper and to respond to the Executive Order for open,

526  accessible, and machine-readable data (Obama 2013), the ECO-OP project explored a

527  provenance use case to adopt the W3C PROV-O standard (Ma et al. 2017). As an example of a

528  report using the PROV-O standard, the U.S. National Climate Assessment is incorporated into

529  the Global Change Information System (GCIS) with a knowledge base that links data products,

530  key messages, and certainty (Tilmes et al. 2013). Future efforts could bridge the ECO-OP

531    prototype with GCIS or other information systems to represent provenance chains from

532    acquisition of source data to inclusion of derived data products in interpreted figures in a report.

533    As an example of analogous efforts, we note that the eReefs project includes integration with

534    provenance and vocabulary services (Car 2013). We also note that semantic mediation may

535    facilitate discovery, access, and understanding of data products by diverse stakeholders and

536    recommend further development of a knowledge network to accommodate concepts in the IEA

537    process (Fig. 2; Fox et al. 2012).

538

539    **Conclusions**

540    Our motivation was to develop cyberinfrastructure, including technology and human expertise,

541    to enable routine, well-documented, integrated assessments of a marine ecosystem. The small

542    team approach with computer scientists and IT specialists working directly with fisheries

543    scientists and oceanographers led to rapid results, with a limiting factor being sufficient training

544    for adoption of the technologies by the larger group of domain scientists. The prototype that we

545    delivered for the Ecosystem Status Report for the Northeast U.S. Continental Shelf Large Marine

546    Ecosystem enabled the reproducibility of a portion of a collaborative, multi-disciplinary report

547    with very heterogeneous data types. However, we only addressed a limited subset of the many

548    technical and social challenges in facilitating collaboration and reproducibility for the Report as

549    a whole. This project provided a pilot toward end-to-end transparency from scientists' desks to a

550    report provided to policy makers and the public, important for science-based decision-making in

551    the U.S. National Ocean Policy Implementation Plan.

552

**List of abbreviations**

ECO-OP, abbreviation joining ECOsystem and interOPerability;

GCIS, Global Change Information System;

GUI, Graphical User Interface;

IDL, Interactive Data Language;

IEA, Integrated Ecosystem Assessment;

IT, information technology;

NOAA, National Oceanic and Atmospheric Administration;

PDF, portable document formats;

TWC, Tetherless World Constellation;

VM, virtual machine;

VRE, virtual research environment;

**References**

Acreman M (2005) Linking science and decision-making: features and experience from environmental river flow setting. Environmental Modelling & Software 20:99-109. doi: 10.1016/j.envsoft.2003.08.019

Ahuja MK, Carley KM (1998) Network structure in virtual organizations. Journal of Computer-Mediated Communication 3:0. doi: 10.1111/j.1083-6101.1998.tb00079.x

572 Car NJ (2013) The eReefs information architecture. 20th International Congress on Modelling

573 and Simulation, Adelaide, Australia, 1-6 December 2013, p. 831-837.

574 http://www.mssanz.org.au.previewdns.com/modsim2013/C7/car2.pdf. Accessed 23

575 December 2015

576 Chen Y, Minchin SA, Seaton S, Joehnk KD, Robson BJ, Bai Q (2011) eReefs – a new

577 perspective on the Great Barrier Reef. 19th International Congress on Modelling and

578 Simulation, Perth, Australia, 12-16 December 2011, p. 1195-1201.

579 http://www.mssanz.org.au/modsim2011/C4/chen.pdf. Accessed 23 December 2015

580 Cooke NJ, Hilton ML (eds) (2015) Enhancing the Effectiveness of Team Science. The National

581 Academies Press, Washington, D.C. doi: 10.17226/19007

582 Di Stefano M, Fox P, Beaulieu S, Maffei A (2012) The integrated ecosystems assessment

583 initiative - enabling the assessment of impacts on large marine ecosystems: informatics to the

584 forefront of science-based decision support. 2012 ICES Annual Science Conference, Bergen,

585 Norway. http://tw.rpi.edu/media/2012/10/08/cd52/ICES_2012.pdf. Accessed 23 December

586 2015

587 Di Stefano M, Fox P, Maffei A, West P, Hare J (2013) An open source approach to enable the

588 reproducibility of scientific workflows in the ocean sciences. American Geophysical Union

589 Fall Meeting, San Francisco, CA. http://tw.rpi.edu/media/2014/02/23/b139/AGU2013-

590 IN51A-15330-MDS.pdf. Accessed 23 December 2015

591 Eaton JW, Bateman D, Hauberg S, Wehbring R (2014) GNU Octave version 3.8.1 manual: a

592 high-level interactive language for numerical computations. CreateSpace Independent

593     Publishing Platform.  ISBN 1441413006.

594     http://www.gnu.org/software/octave/doc/interpreter/

595  Ecosystem Assessment Program (2009) Ecosystem Assessment Report for the Northeast U.S.

596     Continental Shelf Large Marine Ecosystem. U.S. Department of Commerce, Northeast

597     Fisheries Science Center Reference Document 09-11, 61 pp.

598     http://www.nefsc.noaa.gov/publications/crd/crd0911/. Accessed 23 December 2015

599  Ecosystem Assessment Program (2012) Ecosystem Status Report for the Northeast Shelf Large

600     Marine Ecosystem - 2011. U.S. Department of Commerce, Northeast Fisheries Science

601     Center Reference Document 12-07, 32 pp. http://nefsc.noaa.gov/publications/crd/crd1207/.

602     Accessed 23 December 2015

603  Fox P, Batchelder H, Lawrence S, Maffei A, Young O (2012) Information models for

604     development and evolution of complex multi-scale knowledge networks for marine

605     ecosystems. Ocean Sciences Meeting, Salt Lake City, UT.

606     https://tw.rpi.edu//web/doc/OSC2012_139_ecoop_poster. Accessed 23 December 2015

607  Fox P, McGuinness DL (2008) TWC Semantic Web Methodology.

608     http://tw.rpi.edu/web/doc/TWC_SemanticWebMethodology. Accessed 23 December 2015

609  Frederic J (2013) Nbconvert refactor. Final 1.0. http://digitalcommons.calpoly.edu/physsp/85.

610     Accessed 17 October 2016.

611  Goble CA, Bhagat J, Aleksejevs S, et al. (2010) myExperiment: a repository and social network

612     for the sharing of bioinformatics workflows. Nucleic Acids Research 38:W677-W682. doi:

613     10.1093/nar/gkq429

614     GRASS Development Team (2015) Geographic Resources Analysis Support System (GRASS)

615         Software, Version 7.0. Open Source Geospatial Foundation. http://grass.osgeo.org

616     Horsburgh JS (2015) Hydrology domain cyberinfrastructures: Successes, challenges, and

617         opportunities. American Geophysical Union Fall Meeting, abstract #H42A-07.

618         https://agu.confex.com/agu/fm15/meetingapp.cgi/Paper/66729. Accessed 15 April 2016

619     Howison J, Deelman E, McLennan MJ, Ferreira da Silva R, Herbsleb JD (2015) Understanding

620         the scientific software ecosystem and its impact: Current and future measures. Research

621         Evaluation 24:454-470. doi: 10.1093/reseval/rvv014

622     Hunter, JD (2007) Matplotlib: A 2D graphics environment. Computing in Science & Engineering

623         9:90-95. doi: 10.1109/MCSE.2007.55

624     iMarine (2014) Executive Summary: "iMarine data platform for collaborations" workshop, 7

625         March 2014, FAO, Rome, Italy. http://uripreview.i-marine.eu/be0c89a7-6eca-4ae1-ac87-

626         9a52d8800641.pdf. Accessed 31 December 2015

627     Jirotka M, Lee CP, Olson GM (2013) Supporting scientific collaboration: Methods, tools, and

628         concepts. Computer Supported Cooperative Work 22:667-715. doi: 10.1007/s10606-012-

629         9184-0

630     Jones E, Oliphant E, Peterson P, et al. (2001) SciPy: Open Source Scientific Tools for Python.

631         http://www.scipy.org/. Accessed 18 October 2016.

632     Jupyter Team (2015) Jupyter Documentation. Kernels (Programming Languages).

633         http://jupyter.readthedocs.io/en/latest/projects/kernels.html. Accessed 18 October 2016.

634    Levin PS, Fogarty MJ, Murawski SA, Fluharty D (2009) Integrated Ecosystem Assessments:

635        Developing the scientific basis for ecosystem-based management of the ocean. PLoS Biology

636        7:e1000014. doi: 10.1371/journal.pbio.1000014

637    Levin PS, Kelble CR, Shuford RL, Ainsworth C, deReynier Y, Dunsmore R, Fogarty MJ,

638        Holsman K, Howell EA, Monaco ME, Oakes SA, Werner F (2014) Guidance for

639        implementation of Integrated Ecosystem Assessments: a US perspective. ICES J. Mar. Sci.

640        71:1198-1204. doi: 10.1093/icesjms/fst112

641    Liu J, Pacitti E, Valduriez P, Mattoso M (2015) A survey of data-intensive scientific workflow

642        management. Journal of Grid Computing 13:457-493. doi: 10.1007/s10723-015-9329-8

643    Ma X, Beaulieu SE, Fu L, Fox P, Di Stefano M, West P (2017) Documenting provenance for

644        reproducible marine ecosystem assessment in open science. In: Diviacco P, Leadbetter A,

645        Glaves HM (eds) Oceanographic and marine cross-domain data management for sustainable

646        development. IGI Global, Hershey, PA, pp. 100-126, doi: 10.4018/978-1-5225-0700-0.ch005

647    McKinney W (2010) Data Structures for Statistical Computing in Python. Proceedings of the 9th

648        Python in Science Conference:51-56.

649    Muller MJ, Kuhn S (1993) Participatory design. Communications of the ACM 36:24-28. doi:

650        10.1145/153571.255960

651    Muste, M, Bennett, D, Secchi, S, Schnoor, J, Kusiak, A, Arnold, N, Mishra, S, Ding, D, Rapolu,

652        U (2013) End-to-end cyberinfrastructure for decision-making support in watershed

653        management. Journal of Water Resources Planning and Management 139:565-573. doi:

654        10.1061/(ASCE)WR.1943-5452.0000289

655     National Ocean Council (2013) National ocean policy implementation plan.

656       https://www.whitehouse.gov//sites/default/files/national_ocean_policy_implementation_plan.

657       pdf. Accessed 23 December 2015

658     NOAA (2014) National Oceanic and Atmospheric Administration Information Quality

659       Guidelines. Issue date of this revision: 30 October 2014.

660       http://www.cio.noaa.gov/services_programs/IQ_Guidelines_103014.html. Accessed 15 April

661       2016

662     Nosek BA, Alter G, Banks GC, Borsboom D, Bowman SD, Breckler SJ, Buck S, Chambers CD,

663       Chin G, Christensen G, Contestabile M, Dafoe A, Eich E, Freese J, Glennerster R, Goroff D,

664       Green DP, Hesse B, Humphreys M, Ishiyama J, Karlan D, Kraut A, Lupia A, Mabry P,

665       Madon T, Malhotra N, Mayo-Wilson E, McNutt M, Miguel E, Paluck EL, Simonsohn U,

666       Soderberg C, Spellman BA, Turitto J, VandenBos G, Vazire S, Wagenmakers EJ, Wilson R,

667       Yarkoni T (2015) Promoting an open research culture. Science 348:1422-1425. doi:

668       10.1126/science.aab2374

669     Obama B (2013) Executive order -- Making open and machine readable the new default for

670       government information. The White House, Office of the Press Secretary, May 09, 2013.

671       https://www.whitehouse.gov/the-press-office/2013/05/09/executive-order-making-open-and-

672       machine-readable-new-default-government-. Accessed 23 December 2015

673     Olson JS, Hofer EC, Bos N, Zimmerman A, Olson GM, Cooney D, Faniel I (2008) A theory of

674       remote scientific collaboration. In: Olson GM, Zimmerman A, Bos N (eds) Scientific

675       collaboration on the internet. MIT Press, Cambridge, MA, pp 73-99

676    Palmer MA (2012) Socioenvironmental sustainability and actionable science. BioScience 62:5-6.

677         doi: 10.1525/bio.2012.62.1.2

678    Pennington D (2011) Collaborative, cross-disciplinary learning and co-emergent innovation in

679         eScience teams. Earth Science Informatics 4:55-68. doi: 10.1007/s12145-011-0077-4

680    Pérez F, Granger BE (2007) IPython: A system for interactive scientific computing. Computing

681         in Science and Engineering 9:21-29. doi: 10.1109/MCSE.2007.53

682    Pulsifer PL, Collins JA, Kaufman M, Eicken H, Parsons MA, Gearheard S (2011) Applying agile

683         methods to the development of a community-based sea ice observations database. American

684         Geophysical Union Fall Meeting, abstract #IN54A-08.

685         http://adsabs.harvard.edu/abs/2011AGUFMIN54A..08P. Accessed 15 April 2016

686    R Core Team (2013) R: A language and environment for statistical computing. R Foundation for

687         Statistical Computing, Vienna, Austria. http://www.R-project.org/Ragan-Kelley M, Pérez F,

688         Granger B, Kluyver T, Ivanov P, Frederic J, Bussonier M (2014) The Jupyter/IPython

689         architecture: a unified view of computational research, from interactive exploration to

690         communication and publication. American Geophysical Union Fall Meeting, abstract

691         #H44D-07. http://adsabs.harvard.edu/abs/2014AGUFM.H44D..07R. Accessed 15 April 2016

692    Reichman OJ, Jones MB, Schildhauer MP (2011) Science challenges and opportunities of open

693         data in ecology. Science 331:703-705. doi: 10.1126/science.1197962

694    Samhouri JF, Haupt AJ, Levin PS, Link JS, Shuford R (2014) Lessons learned from developing

695         integrated ecosystem assessments to inform marine ecosystem-based management in the

696         USA. ICES Journal of Marine Science 71:1205-1215. doi: 10.1093/icesjms/fst141

697    Samhouri JF, Lester SE, Selig ER, Halpern BS, Fogarty MJ, Longo C, McLeod KL (2012) Sea

698          sick? Setting targets to assess ocean health and ecosystem services. Ecosphere 3:art41. doi:

699          10.1890/ES11-00366.1

700    Shen H (2014) Interactive notebooks: Sharing the code. The free IPython notebook makes data

701          analysis easier to record, understand and reproduce. Nature 515:151-152. doi:

702          10.1038/515151a

703    Taconet M, Ellebroek A, Castelli D, Pagano P, Caumont H, Garavelli S, Parker S (2014)

704          Sustaining iMarine: a public partnership led business model. The iMarine Sustainability

705          White Paper, final release November 2014, 65 pp.

706          ftp://ftp.fao.org/FI/DOCUMENT/FIGIS_FIRMS/2015/Inf11e.pdf. Accessed 23 December

707          2015

708    Tilmes C, Fox P, Ma X, McGuinness DL, Privette AP, Smith A, Waple A, Zednik S, Zheng JG

709          (2013) Provenance representation for the national climate assessment in the global change

710          information system. IEEE Transactions on Geoscience and Remote Sensing 51:5160-5168.

711          doi: 10.1109/TGRS.2013.2262179

712    van den Hove S (2007) A rationale for science-policy interfaces. Futures 39:807-826. doi:

713          10.1016/j.futures.2006.12.004

714    Wilson G (2014) Software Carpentry: lessons learned. F1000Research 3:62, Version 1, 19 Feb

715          2014. doi: 10.12688/f1000research.3-62.v1

716    Wright DJ (2016) Toward a digital resilience. Elementa: Science of the Anthropocene 4:000082.

717          doi: 10.12952/journal.elementa.000082

718

**Figure captions**

720 **Fig. 1.** Diagram of the Integrated Ecosystem Assessment (IEA) process, driven by the goals and

721 targets of Ecosystem-Based Management (EBM; image available online at:

722 http://www.noaa.gov/iea/loop.html).

723 **Fig. 2.** Schematic for data interoperability in the Ecosystem Status Report for the Northeast U.S.

724 Shelf Large Marine Ecosystem. The data sources (lower layer), applications (middle layer,

725 including a blank field for new tools), and the resulting integrated data products and indicators

726 for the Report (upper layer) reflect the key elements in the use case. The two gray layers indicate

727 mediation and the potential for semantic interoperability.

728 **Fig. 3.** Representative data products and indicators in the Ecosystem Status Report for the

729 Northeast U.S. Shelf Large Marine Ecosystem. (a) Time-series indicator: Mean trophic level of

730 landings by commercial fisheries [from Fig. 8.2 in Ecosystem Assessment Program (2012)]. (b)

731 Spatial data product: Mean (1998-2010) daily primary production [from Fig. 4.2 in Ecosystem

732 Assessment Program (2012)].

733 **Fig. 4.** Diagram of TWC Methodology, an iterative use case development methodology

734 [modified from Fox and McGuinness (2008)].

735 **Fig. 5.** Activity diagram for the Ecosystem Status Report use case, indicating actors, entities (i.e.,

736 data files, image products, and the Report), and activities (arrows). Note the data retriever and

737 processor is represented as a software agent (square head). The dashed box contains the activities

738 for which we built the prototype.

739 **Fig. 6.** Screen grab of a portion of the executed Climate Forcing Notebook, showing: opening a

740 document, importing text files, accessing a source data file, processing data, and plotting and

741     saving derived data products (to view details, please refer to the notebook at the GitHub

742     repository accessible via https://data.rpi.edu/xmlui/handle/10833/1756).
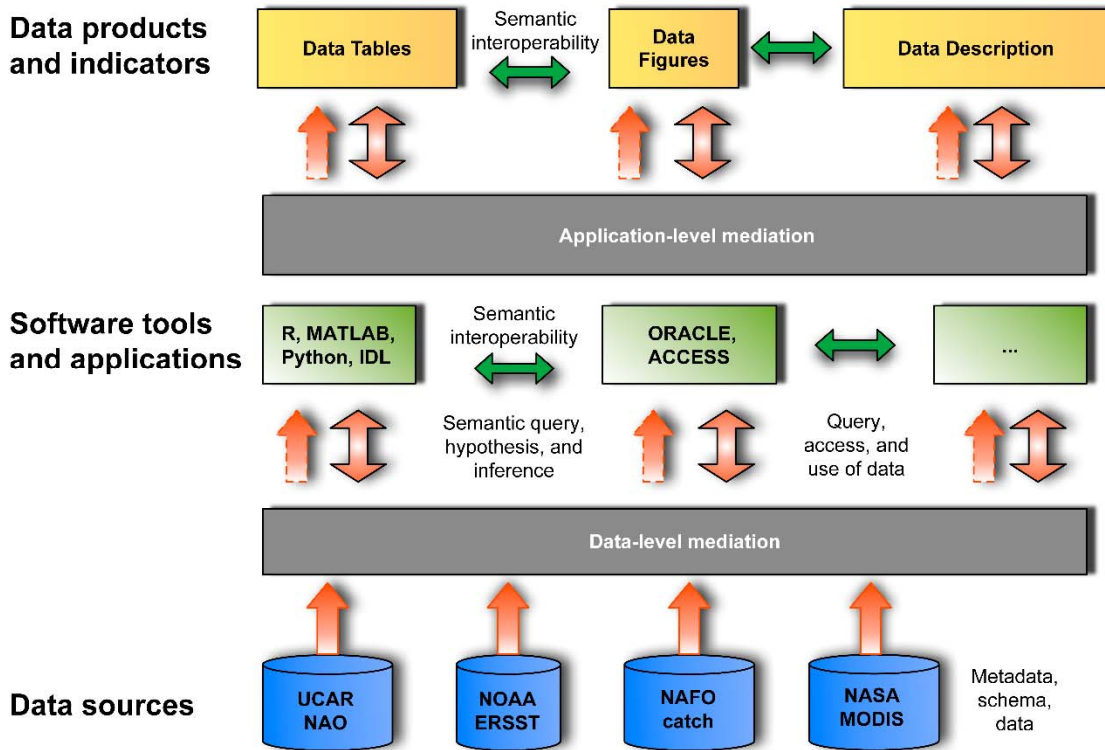
743     **Fig. 7.** Screen grab of the PDF document that results from the executed Climate Forcing

744     Notebook (to view details, please refer to the PDF at the GitHub repository accessible via

745     https://data.rpi.edu/xmlui/handle/10833/1756).

746

747     **(Figures submitted separately)**

748     **Fig. 1.**



749

750

751     **Fig. 2.**

752

753

754

755    **Fig. 3. (a) (b)**

756

757    **Fig. 4.**



758

759

760    **Fig. 5.**

761

762

763 **Fig. 6.**

**Document**

```
In [8]: ID = util.get_id('test/Climate-forcing_pdf')
        document = openDocument()
```

session data directory : test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM

**Section 1**

```
In [9]: %%writefileref (ID)/climate_forcing.txt (ecoop_username)
        Climate patterns over the North Atlantic are important drivers of oceanographic conditions and ecosystem states.
        Steadily increasing atmospheric carbon dioxide levels can not only affect climate on global and regional scales
        but alter critical aspects of ocean chemistry. Here, we describe the atmospheric forcing mechanisms related
        to climate in this region including large-scale atmospheric pressure systems, natural ocean temperature cycles in the North Atlantic,
        components of the large-scale circulation of the Atlantic Ocean, and issues related to ocean acidification.
```

Writing test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM/climate_forcing.txt

u'added references for user anonymous'

```
In [10]: section = addSection(name='Climate Forcing', data=os.path.join(ID,'climate_forcing.txt'))
```

**Sub Section 1**

```
In [11]: %%writefileref (ID)/nao.txt (ecoop_username)
        Climate and weather over the North Atlantic are strongly influenced by the relative strengths
        of two large-scale atmospheric pressure cells -- the Icelandic Low and the Azores High [4].
        As the relative strengths of these two pressure systems vary, characteristic patterns of temperature, precipitation, and wind fields are
        Iceland in the winter (December-February; see Glossary for a description of methods used to create standardised indicators).
        An index of this dipole pattern has been developed based on the standardised difference in sea level pressure between Lisbon, Portugal and
        This North Atlantic Oscillation (NAO) index has been related to key oceanographic and ecological processes in the North Atlantic basin [5]
        When the NAO index is high (positive NAO state), the westerly winds shift northward and increase in strength.
        Additionally, there is an increase in precipitation over southeastern Canada, the eastern seaboard of the United States,
        and northwestern Europe. Water temperatures are cool off Labrador and northern Newfoundland, influencing the formation of Deep Labrador S:
        but warm off the United States.
        Conversely, when the NAO index is low (negative NAO state), there is a southward shift and decrease in westerly winds, decreased stormine:
        and drier conditions over southeastern Canada, the eastern United States, and northwestern Europe.
        Water temperatures are warmer off Labrador and Newfoundland, but cooler off the eastern United States.
        Since 1972, the NAO has primarily been in a positive state (Figure 1), although notable short-term reversals to a negative state have bee
        Changes in the NAO have been linked to changes in plankton community composition in the North Atlantic, reflecting changes in both the dis
        and abundance of warm and cold-temperate species.
```

Writing test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM/nao.txt

u'added references for user anonymous'

```
In [12]: naodata = cfd.nao_get(save=ID, csvout="nao.csv", prov=True)
```

dataset used: https://climatedataguide.ucar.edu/sites/default/files/climate_index_files/nao_station_djfm.txt
nao data saved in : test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM/nao.csv

'cell-output metadata saved'

```
In [13]: # NAO
        naodata = cfd.nao_get(save=ID, csvout="nao.csv")
        cfp.plot_index(name='NAO_lowess', xticks=10, xticks_fontsize=10,
                    data=naodata, nb='y', scategory='lowess', frac=1./6, it=6,
                    output=ID, dateformat=True, figsave="nao.png", prov=True)
```

dataset used: https://climatedataguide.ucar.edu/sites/default/files/climate_index_files/nao_station_djfm.txt
nao data saved in : test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM/nao.csv
graph saved in : test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM/nao.png
NAO_lowess smoothed data saved in : test/Climate-forcing_pdf_Wednesday_25_June_2014_10_49_29_PM/NAO_lowess_lowess.csv

'cell-output metadata saved'

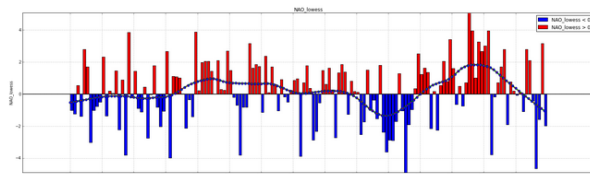Session output file 'subplots.html' already exists, will be overwritten.

764

765

766    **Fig. 7.**

# 1 Climate Forcing

Climate patterns over the North Atlantic are important drivers of oceanographic conditions and ecosystem states. Steadily increasing atmospheric carbon dioxide levels can not only affect climate on global and regional scales but alter critical aspects of ocean chemistry. Here, we describe the atmospheric forcing mechanisms related to climate in this region including large-scale atmospheric pressure systems, natural ocean temperature cycles in the North Atlantic, components of the large-scale circulation of the Atlantic Ocean, and issues related to ocean acidification.

## 1.1 North Atlantic Oscillation Index

Climate and weather over the North Atlantic are strongly influenced by the relative strengths of two large-scale atmospheric pressure cells – the Icelandic Low and the Azores High [4]. As the relative strengths of these two pressure systems vary, characteristic patterns of temperature, precipitation, and wind fields are observed. An index of this dipole pattern has been developed based on the standardized difference in sea level pressure between Lisbon, Portugal and Reykjavík, Iceland in the winter (December-February; see Glossary for a description of methods used to create standardized indicators). This North Atlantic Oscillation (NAO) index has been related to key oceanographic and ecological processes in the North Atlantic basin [5]. When the NAO index is high (positive NAO state), the westerly winds shift northward and increase in strength. Additionally, there is an increase in precipitation over southeastern Canada, the eastern seaboard of the United States, and northwestern Europe. Water temperatures are cool off Labrador and northern Newfoundland, influencing the formation of Deep Labrador Slope water, but warm off the United States. Conversely, when the NAO index is low (negative NAO state), there is a southward shift and decrease in westerly winds, decreased storminess, and drier conditions over southeastern Canada, the eastern United States, and northwestern Europe. Water temperatures are warmer off Labrador and Newfoundland, but cooler off the eastern United States. Since 1972, the NAO has primarily been in a positive state (Figure 1), although notable short-term reversals to a negative state have been observed during this period. Changes in the NAO have been linked to changes in plankton community composition in the North Atlantic, reflecting changes in both the distribution and abundance of warm and cold-temperate species.
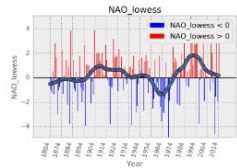


Figure 1: North Atlantic Oscillation - metadata.

## 1.2 Atlantic Multidecadal Oscillation

Multidecadal patterns in sea surface temperature (SST) in the North Atlantic are represented by the Atlantic Multidecadal Oscillation (AMO) index. The