# Reproducible Research in Geoinformatics: Concepts, Challenges and Benefits

## Christian Kray 🔘
Institute for Geoinformatics (ifgi), University of Münster, Germany
c.kray@uni-muenster.de

## Edzer Pebesma 🔘
Institute for Geoinformatics (ifgi), University of Münster, Germany
edzer.pebesma@uni-muenster.de

## Markus Konkol 🔘
Institute for Geoinformatics (ifgi), University of Münster, Germany
m_konk01@uni-muenster.de

## Daniel Nüst 🔘
Institute for Geoinformatics (ifgi), University of Münster, Germany
daniel.nuest@uni-muenster.de

──── **Abstract** ────

Geoinformatics deals with spatial and temporal information and its analysis. Research in this field often follows established practices of first developing computational solutions for specific spatiotemporal problems and then publishing the results and insights in a (static) paper, e.g. as a PDF. Not every detail can be included in such a paper, and particularly, the complete set of computational steps are frequently left out. While this approach conveys key knowledge to other researchers it makes it difficult to effectively re-use and reproduce the reported results. In this vision paper, we propose an alternative approach to carry out and report research in Geoinformatics. It is based on (computational) reproducibility, promises to make re-use and reproduction more effective, and creates new opportunities for further research. We report on experiences with executable research compendia (ERCs) as alternatives to classic publications in Geoinformatics, and we discuss how ERCs combined with a supporting research infrastructure can transform how we do research in Geoinformatics. We point out which challenges this idea entails and what new research opportunities emerge, in particular for the COSIT community.

## 1 Introduction

Spatial and temporal information and their analysis play a central role in many scientific disciplines such as Geography or Economics and are essential in understanding and solving many pressing societal issues. Regardless of which disciplinary perspective is taken, the way in which the corresponding research is carried out and reported is very similar. On an abstract level, a researcher will work on a specific issue or problem that has a spatiotemporal component, for example, computing a route that meets certain criteria, or checking whether

a number of spatiotemporal constraints can be fulfilled given a set of moving objects. She might develop a representation to computationally describe the problem and/or an algorithm/approach that can be applied to the problem in order to answer the research question. These might be entirely new or derived from existing ones. Or, she might carry out some studies to explore spatiotemporal aspects in the real world (e.g. how humans perform when given specific navigational instructions) and then derive some potentially generalisable insights (e.g. performance varies according to the length of instructions). She will then report her findings in a publication, i.e. a static document consisting of text, tables, and figures that concisely describe what she found out, how she went about her investigation, and (ideally) how others can independently repeat the steps she took to confirm her findings. Once peer-reviewed and published, others will use her results for further research.

This overall approach is very similar to many other scientific disciplines, and has existed in this form for a long time. While generally well established and effective, there are a number of key shortcomings inherent to this process. Firstly, it is usually not possible to actually or easily reproduce results as the required data, analysis procedures (e.g. source code), and the necessary configuration to re-run the analysis (e.g. compiler version) on the data are not available. Even when they are, the effort of reproducing the results can be prohibitive, e.g. due to having to restore the exact configuration of the computational environment that was used by the author of the paper. This configuration may come with little documentation, e.g. since it evolved over a long time period in a trial-and-error fashion. If this information is available at all, it can still involve issues such as having to retrieve and install outdated libraries, which also hinders re-use [15]. Secondly, the paper is usually provided as unstructured text for human readers that makes a deeper (computational) analysis difficult. For example, finding all papers that investigate the same spatial region or use the same spatiotemporal calculus is not trivial. In addition, this lack of semantic structure prevents presenting results in different ways that might be better suited for different (non-academic) audiences or purposes. Thirdly, with a static document it is very difficult for the reader or reviewer to explore the results more deeply. For example, a reader cannot easily check whether the reported results are robust (e.g. when parameters are slightly modified or a particular axiom is weakened), or whether they are consistent with previously reported results that might have used a different dataset. These issues apply for authors as well as readers and reviewers.

In this paper, we present several concepts and ideas for overcoming these issues for research which deals with spatiotemporal information and its computational analysis. After first reviewing the status-quo of how research in Geoinformatics is done currently, we describe the concept of executable research compendia (ERCs), which addresses some of the issues outlined above. Based on these concepts, we then propose our vision for how computational research on spatiotemporal aspects can be carried out in the future and propose an open research infrastructure for Geoinformatics (OpenRIG). This approach comes with a number of key challenges and opportunities, which we discuss subsequently. The paper concludes by outlining limitations and summarising our key contributions.

## 2    Background

Spatial and temporal information is essential when dealing with a variety of problems. In the past decades, several scientific sub-disciplines have emerged which tackle these kinds of problems across traditional disciplinary boundaries and often act as an integrator, e.g. between computer scientists and geologists. *GI Science* can be defined in different ways (cf. [8] for a short review of definitions). The research conducted under this label stretches from

semantic interoperability and ontologies, theories of spatial-temporal information systems, and spatial data modelling and services, to user-generated data [26]. Donoho critically discusses *Data Science* as a scientific field in its own right, which concerns not only the extraction of information from data, but "each and every step that the professional must take, from getting acquainted with the data all the way to delivering results based upon it, and extending even to that professional's continual review of the evidence about best practices of the whole field itself" [7, p. 794]. Spatiotemporal data is one important type of data covered by Data Science, with *Spatial Data Science* being an emerging subfield focusing on this. We define *Geoinformatics* as the discipline that generally deals with spatial and temporal information computationally (e.g. geostatistics, geosimulation, geovisualisation, interaction with spatiotemporal information). The definitions listed above overlap and describe similar disciplines that each take a different perspective. Whereas GI Science emphasises the geographic aspect and the use of geographic information systems (GIS), Data Science focuses on data and its properties, while Geoinformatics more generally considers spatial and temporal information and its analysis and use.
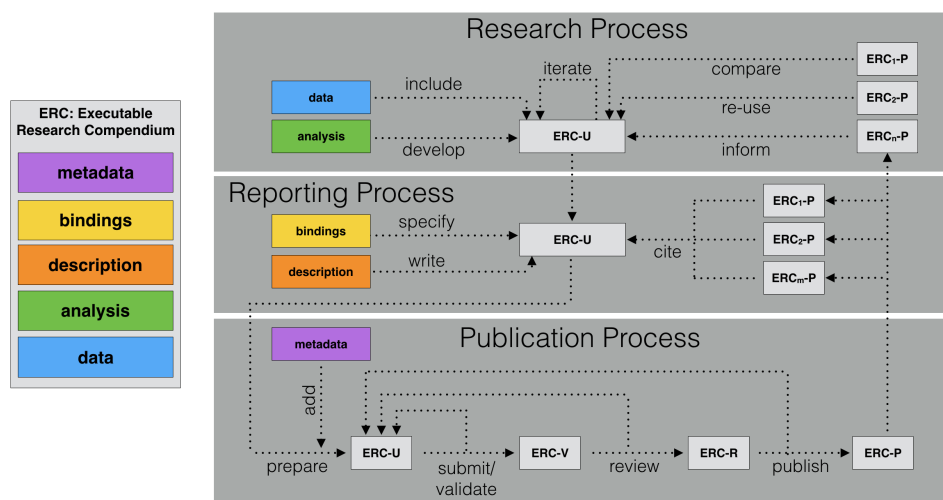
Besides investigating spatiotemporal aspects, the scientific disciplines introduced above share another commonality: the way in which a substantial part of research is carried out and reported today. After inception of an idea, scientists formulate a hypothesis or research questions, review existing work on the topic, gather data and information, and analyse them. At some point, researchers author a scholarly publication that describes their methods and results and then submit it to an academic outlet, where it undergoes peer-review. If accepted, a document is created by the publisher and made available to readers of the outlet. This process today mostly takes place in the digital realm – frequently including the research work itself [12]. Specific instances of such a publication process may vary, especially with recent initiatives for more openness and quality such as public reviews [27], preregistration (e.g. against HARKing) [4], Open Access licensing [24], preprint servers [2], independent non-profit journal publishers, open data publishing [11], FAIR data [31], or open code [9].

Although these practices have demonstrated their potential, the habits of researchers are slow to change and require action by all involved stakeholders [19]. That is why a paper published today is still very much a snapshot in time describing the results of a researcher's work. As Buckheit and Donoho put it: "An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures" [3, p. 59]. The parts that make up the publication are difficult to re-use and applying the "good enough" practices [32] is still perceived as additional work, not as a way to improve the review process and make an article's argument stronger. At the same time, the citation as a mechanism to give credit is disputed, with software and data citations [13, 22] as well as altmetrics [23] emerging as complementary means to recognise and refer to activities that demand a lot of time and effort from researchers.

In summary, we can thus conclude that though there are different disciplines dealing with spatiotemporal information, scientists in these fields do research in a similar way and face similar challenges. This way is common in other scholarly disciplines as well but suffers from a number of key issues such as low reproducibility, difficult re-use, and no deeper inspection. The following sections will propose an alternative approach that tackles these problems in Geoinformatics. We focus on Geoinformatics here due to its computational emphasis while dealing with spatiotemporal information but the ideas reported below are applicable to other fields as well. It is important to point out that our focus is on computational research here. Other types of research carried out in the fields discussed above (e.g. qualitative or explorative research) are outside of the scope of this paper.

## 3 Executable Research Compendia

In order to overcome the issues we highlighted above, we have developed an approach that combines essential elements related to a specific research activity in Geoinformatics into a coherent compendium. In the following we first describe the basic concept and its realisation before discussing how it fits into the research workflow. We also briefly contrast this approach with current practice and highlight implications of using the new approach.



**Figure 1** Executable Research Compendium (ERC) with its five key components (left) and how ERCs can be integrated into the research, reporting and publication processes (right). ERC-U stands for an unvalidated ERC, ERC-V for a validated one, ERC-R for a reviewed ERC, and ERC-P for a published one. Processes are sequentialised to make the figure more readable.

### 3.1 Concept and realisation

An *Exectuable Research Compendium (ERC)* includes all research components that are needed to reproduce the computational results in a paper [20]. It is meant to replace the "classic" static paper, e.g. a PDF or HTML file made up of textual, tabular, and graphical elements that describe the research work, its outcomes, context, and meaning. The five ERC components are depicted in Figure 1 (left) and consist of the following five items (from bottom to top):

- the *data* that was the input to a process creating the results that are being reported; this could consist, e.g., of satellite imagery, a formal specification of spatio-temporal configurations or a transcript of navigational instructions produced by study participants.

- the *analysis* or computational steps that were applied to the data in order to generate the results; examples for this component may include source code implementing a geostatistical method, a set of rules formally specifying a spatial reasoning process, or a code snippet that generates word clouds of textual navigational instructions.

- the *description* of the reported research; this part corresponds largely to a "classic" paper: it contains text providing motivation and background of the research as well as a description of the process, its outcome and the implications/meaning of the latter

- the *bindings* are an optional[1] component describing the links between the three components listed above on a fine-grained level. They can specify which part of the analysis produced which result reported in the description using which part of data. For example, a binding can encode that a specific figure in the description was produced using a particular function included in the analysis component that was applied to a specific dataset included in the data component. Bindings can also specify user interface (UI) controls that enable readers to deeply interact with results, for example by moving a slider next to a figure to change a parameter in the computations that produced this figure, which in turns results in a re-run of the computations and an updated figure.
- the *metadata* component contains meta information about the entire research component. For example, it may include author names, keywords or version information about the interpreters and libraries that were used to produce the results reported in the ERC.

By combining these five elements in a coherent compendium, all relevant information is readily available to carry out a number of essential research activities that currently are difficult to perform with "classic" publications. For example, an ERC makes it easy to re-run the computational steps behind the reported results both for human readers and systems supporting ERCs as a digital object. Section 3.3 includes a more in-depth comparison between ERCs and traditional publications.

The feasibility of the proposed concept is demonstrated by a prototypical platform that was implemented using containerisation to encapsulate the five components of an ERC [21]. The container recipe, in case of Docker[2]: the *Dockerfile*, also specifies the computational environment. The current prototype supports research that uses the R language[3] [28] to carry out spatiotemporal analysis on arbitrary data. The prototypical platform can be extended to include other analysis methods (e.g. the Python language[4] or a specific theorem prover). The prototype also provides methods for creating, running and comparing ERCs as well as further functionalities such as support for interactive figures and "one-click reproduce".

## 3.2 Integrating ERCs into the research workflow

While ERCs on their own offer some useful properties that can benefit researchers individually, it makes sense to consider their use in the larger context of a typical research workflow. Generally speaking, researchers carry out multiple, potentially overlapping activities when investigating a particular topic. In addition to performing the research itself, for example developing, testing and evaluating a model to simulate navigation, they also work on reporting and publishing the research so that other researchers can use (i.e. extend or build upon) it in their work. ERCs can be used at each step in this overall workflow as detailed in the following and depicted in Figure 1 (right).

Starting with the *research process* itself, a researcher can work on an unvalidated ERC (ERC-U) while deciding which data to include and while developing the analysis procedure. In doing so, she can use previously published ERCs (ERC-P) in several ways. She can compare her data, analysis or results to those reported in an ERC-P. She can also directly re-use parts of an ERC-P, for example, its data to test her analysis. Finally, she can use an ERC-P much like a "classic" publication to inform her work on a more general level, e.g. to

---

[1] Bindings are optional as they require additional effort from authors and basic reproducibility can already be achieved without them.

[2] `https://www.docker.com`, accessed on August 21, 2019.

[3] `https://www.r-project.org`, accessed on August 21, 2019.

[4] `https://www.python.org`, accessed on August 21, 2019.

modify the assumptions or thresholds she is using in her analysis. Once the research has progressed sufficiently, the researcher will want to *report* on it so that others can benefit from her insights. For this purpose, she will write a description much like a 'classic' publication, consisting of text, tables and figures, and add it to the ERC-U that contains the data and analysis she developed previously. During this process, she will contextualise her work by referring to related work, citing other ERCs or parts of them. In addition, she can create a bindings component for her ERC-U to make explicit how data, analysis and description are connected at a fine-grained level. For example, she might specify for a figure which function in her analysis component was applied to what part of the data to generate that figure.

Typically, the researcher will then want to publish the compendium in an academic outlet (e.g. a journal or a conference) to ensure interested readers learn about the new insights she gained through her research and to receive feedback. For this purpose, she first adds relevant metadata to ensure that others can easily use the ERC she has produced and that the ERC conforms to the requirements of the outlet. (Semi-)automatic validation mechanisms can help with this preparation and ensure that the analysis component in the ERC produces the results in the description component. Once the validation is successful, the validated ERC (ERC-V) can be submitted to peer-review. Reviewers are now able to investigate a much larger part of the scholarship underlying the "classic" paper, which is usually not available at all or not integrated into the review process. For example, if reviewers possess the skills to review the code, they can do so, but even if they do not, they can at least confirm that the computations produce the results reported in the paper and check how the results vary when parameters are changed. Depending on the outlet the review process can take different forms and may result in several iterations during which the author needs to revise her ERC. If accepted, the ERC can then be published, resulting in a published ERC (ERC-P). The final publishing process includes multiple steps, such as updating the meta-information to specify which issue/year the ERC was published in. The ERC-P then becomes available for other scientists, who can use it for their own research.

## 3.3   Benefits and challenges

ERCs come with a number of benefits that enrich a reader's workflow while studying a paper [14]. First, the source code is directly reproducible/executable in a predefined computational environment. In contrast, today's papers are, if at all, supplemented by a folder including files that contain data and code or an (incomplete) reference to the software that was used (e.g. without version, and not persistently stored). This leaves readers with the (daunting) task of figuring out how to run the included analysis. Reviewers face the same challenge with traditional papers while ERCs enable them to easily validate the analysis described in a paper by rerunning it. Second, ERCs can also enhance how researchers work with scientific publications. When searching for relevant articles, they can do so in a more fine-grained way, e.g. by also considering spatial and temporal information derived from the data or by specifying spatiotemporal constraints. While reading an article, they can in parallel investigate how the authors produced a specific figure. In case the analysis code exposes control parameters, it is possible to provide readers with UI controls that enable them to interactively manipulate the initial configuration within a range predefined by the author to see how the results change. Finally, researchers can substitute the original input dataset by data that resulted from their own experiments or another paper, or replace the included code with an alternative one. To fully realise those benefits, an easy-to-use UI is highly desirable not only for inspecting results or interacting with them but also for comparing the computational outputs to quickly identify differences. Otherwise, it might be difficult to spot the differences of an original figure and the one produced by changing a parameter.

The use of ERCs also introduces a number of challenges that need consideration. One such issue is that the concept of an ERC can only fully work if the used research software and data are open source and available with a suitable license (e.g. analysis code written in R and the data could be provided under Creative Commons Zero). Otherwise, it is not possible to include all necessary software to reproduce the results of the paper, to create bindings or to fully re-use the analysis/data. However, many computational analyses are realised using restrictively licensed and proprietary software or data formats, such as Matlab [18] or ArcGIS [10]. Though specific licenses or managed execution by the companies producing the proprietary software might facilitate their use in ERCs, the inherent lack of source code counteracts the principle of complete transparency of the research. A further challenge linked to ERCs is that researchers may be unable (e.g. due to privacy concerns) or reluctant to open their research to the degree required by ERCs despite the known benefits of open data [25] or efficiency, continuity, and reputation [17]. Researchers might fear that others "steal' their data/ideas, that the additional information exposes them to heightened scrutiny, or that the code they wrote for the analysis is not worth publishing or a publication might be harmful [1]. Finally, if an ERC reports on research that involves time-consuming computations or very large data sets, then local execution or transmission of the ERC become unfeasible.
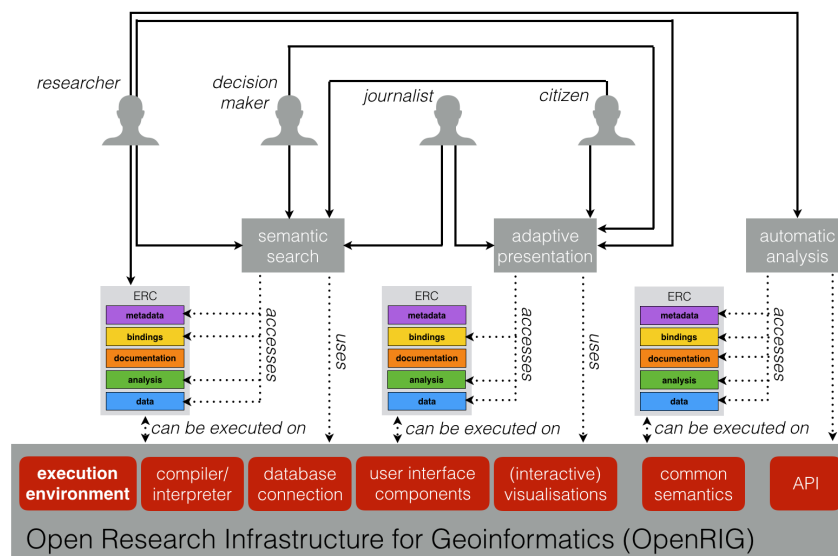
In order to establish ERCs as a desirable alternative to current practice and to realise the benefits outlined above, it is thus necessary to not only overcome some technical and legal challenges but also to change the mindset and behaviour of researchers working in Geoinformatics. Clearly communicating the potential benefits and adding further ones (such as interactive figures, one-click-reproduce or easy re-use) can help to address the latter issue. While we cannot deny that creating ERCs requires additional effort compared to submitting a PDF file, we believe that making research results easier to find, easier to understand, and reusable does not only lead to a higher impact but also strongly benefits authors and the discipline as a whole. The legal aspects are more difficult to tackle. Requiring researchers to only use open source software with permissive licenses could negatively affect their work. It also seems unlikely that vendors of proprietary software would easily agree to their software becoming part of an ERC as this would enable third parties to run their software for free. Many of the (technical) challenges outlined above (big data/long computation) can be addressed by the support infrastructure for ERCs we envision in the following section.

## 4 An Open Research Infrastructure for Geoinformatics (OpenRIG)

ERCs encapsulate individual research contributions and by combining all key elements linked to those contributions in a coherent way, they offer a number of benefits as outlined in section 3.3. In order to realise their full potential, there needs to be a supporting infrastructure that provides various functions, in particular those that affect multiple ERCs. Since previous work has focused on the realisation of ERCs individually [14, 15, 21], we will outline our vision for such an infrastructure in the following sections. We first give a rationale for why it is needed and then discuss components and services it should provide. In addition, we point out key challenges and opportunities that result from working towards and with such an infrastructure. Figure 2 provides an overview of what we envision for this infrastructure.

### 4.1 Rationale

There are several reasons why a research infrastructure for ERCs is beneficial. From a technical perspective, local execution (i.e. on the researcher's computer) is not always feasible: the computer might not be powerful enough to execute the ERC, the computation might take

**Figure 2** Open Research Infrastructure for Geoinformatics (OpenRIG): key components (red), essential functionalities enabled by it (grey boxes) and different stakeholders wanting to access them.

too long, and/or the amount of data might be to large to transmit/store. In addition, some desirable functions such as comparing multiple ERCs / their components or a *semantic search* require access to *many* ERCs at the same time, which is not feasible inside an individual ERC. For the semantic search in particular, it is necessary to analyse specific components of several ERCs and to apply constraints/filters across them.

Another opportunity that arises from the structure inherent to ERCs and the inclusion of bindings in particular is the ability to adapt the way in which results are presented. This is beneficial, for example, to harmonise diagrams across multiple ERCs so that it is easier to relate the reported results to one another. *Adaptive presentations* also allow for generating figures and diagrams that meet the needs and preferences of different stakeholders. For example, researchers might want a higher level of detail than journalists and colour-blind readers might prefer contrast-rich colour palettes.

Furthermore, ERCs also pave the way for *automatic analysis* processes. Unlike traditional papers, ERCs are highly structured, semantically annotated, and they include data and code. Automatic processes can use this structured information to perform tasks such as automatic validation of individual ERCs (are the results reported in the description component produced by the included analysis applied to the data) or cross-validation (are the results reported in an ERC consistent with those reported in ERCs with the same/similar data set and/or the same/similar analysis method). More sophisticated automated analyses are also possible, e.g. by using spatiotemporal reasoning to combine results from multiple papers.

Additionally, there are a number of practical reasons why a research infrastructure for ERCs is beneficial. It would allow for the seamless realisation of the publication process depicted in Figure 1. In addition, it could also provide a means for archival of ERCs so that long-term storage (and execution) can be guaranteed. Finally, a common infrastructure could also contribute towards ensuring that all researchers have access to the same resources thereby levelling the playing field.

## 4.2 Components

An infrastructure supporting functionalities as discussed in the previous sections fundamentally needs to be open, i.e. to facilitate the inclusion of runtime environments and to enable bindings. The term "open" in this context specifically refers to "open source" and a license that is permissive. We envision an *Open Research Infrastructure for Geoinformatics (OpenRIG)* that contains the following core components/elements.

The most obvious component of the OpenRIG is an *execution environment* that can execute ERCs "in the cloud". This is particularly important in cases where the analysis included in the ERC is computationally demanding, e.g. would take hours, days or longer on a standard PC. Such an execution environment would also benefit the realisation of automatic analysis processes for similar reasons. Automatic analysis processes could potentially be realised as ERCs as well: the analysis itself would be contained in the analysis component and the documentation could describe the purpose of the automatic analysis process. However, this might require means to persist the state of the analysis between subsequent executions of the automatic analysis. Otherwise, the periodic re-running of such a process might entail unnecessarily re-analysing all ERCs that were processed already in previous executions. Automatic execution of ERCs might also come into play in case of updates of a software library that was used in the analysis code. The resulting changes could affect the final output making it necessary to re-execute the analysis and to see if the results are still the same. A further possibility is related to metadata which is usually entered manually: an automatic extraction might relieve authors from this task, which is disliked by many researchers. An automatic analysis could also be part of an "ongoing" meta-analysis determining, for example, how many papers address research about a specific topic. Once the computational analysis is implemented, the input data could be updated regularly.

Independently of how such automatic analysis processes are realised, they have a substantial potential to improve the way in which research is done in Geoinformatics. For example, they could automatically compare ERCs that are newly published to similar existing ERCs, e.g. to check for anomalies or to infer new knowledge resulting from a newly added ERC. A new method to predict navigation errors might thus be applied to the data in previous papers that investigated the same issue. This, in turn, might help to better determine the overall performance of the new method compared to previously published ones.

Another useful element of a research infrastructure for Geoinformatics is a persistent and standardised *connection to relevant databases*. Particularly when dealing with very large datasets such as satellite data, it is not feasible to include a full copy of the data in each ERC that used it in its analysis. Instead, using a persistent, versioned interface to access this data from an ERC does not only reduce the size of an ERC but also has benefits in terms of efficiency. Given such a database connection, the infrastructure can make sure that the analysis is executable "near" the actual data, e.g. on the same server, to avoid unnecessary delays resulting from transmitting large amounts of data.

Two further beneficial and related elements of an execution environment are a set of *user interface (UI) components* and *(interactive) visualisations*. The former could provide a library of controls that enable authors to add interactivity to their ERCs, e.g. to provide readers with means to explore more deeply how results change when certain assumptions change. The latter targets visualisations particularly as they play a key role in the understanding of scientific publications [6]. In addition to standard diagrams such as bar charts or box plots, this includes in particular spatiotemporal visualisations such as maps or space-time cubes. When ERCs use these sets of standardised UI components and visualisations, this opens up the possibility to change how specific results are presented. For example, the projection,

colour palette and colour break values of a map in one ERC could be adapted to be the same as the one used in another ERC to facilitate comparison. Different stakeholders might also prefer certain visualisations (e.g. varying degree of detail, colour scheme), which constitutes another type of adaptive presentations that could be realised in this way.

In addition, we envision the OpenRIG to contain a *common semantic framework* that formally describes core spatiotemporal concepts, the components of ERCs as well as their relationships and interactions. This will not only enable reasoning about individual ERCs but also strengthen the capabilities of semantic search functions and automatic analysis processes as they can operate on a more abstract level. Semantic information about spatiotemporal aspects of ERCs can be included in their metadata component. Furthermore, semantically describing different visualisations and UI components opens the door for reasoning about how information is presented in line with the idea of interface plasticity [5]. This can also help addressing one of the key challenges in geovisualisation: adapting geovisualisations [16].

The final element that a research infrastructure for Geoinformatics should provide is a standardised *Application Programming Interface (API).* This more technical aspect is important for providing easy access to all the functions provided by the OpenRIG. For example, in order to implement semantic search, adaptive presentations or automatic analysis, it is necessary to interact with various elements of the OpenRIG. Providing a versioned and persistent API is not only important to ensure automatic analysis processes remain operational but also for the long-term archival (and execution) of ERCs.

The open research infrastructure for Geoinformatics (OpenRIG) envisioned in this section thus facilitates various desirable functions. In addition, it provides a number of interesting opportunities for future research and poses several challenges that require further investigation. The following section provides an overview over these aspects.

## 4.3 Opportunities and Challenges

Among the opportunities offered by the envisioned research infrastructure, being able to *perform reasoning* on top of it is very promising. Obviously, a well-designed semantic framework is required to realise this, which constitutes another opportunity for interesting research. The envisioned reasoning includes the automatic analysis processes mentioned above but can be extended substantially – given that ERCs are highly structured, semantically annotated entities with a strong spatiotemporal component. An example for such work could be research into which example cases the region connection calculus (RCC) [29] was applied to in different publications. The collected datasets describing the example cases could then be used as a corpus to compare RCC to a different approach in terms of whether the latter can solve all example cases as well, and whether the results are identical for both approaches. In addition to reasoning about the content of an ERC, there is an opportunity to further investigate how key results are presented. As mentioned in section 4.2, bindings semantically link data, analysis and description. Formally describing descriptions of spatiotemporal results – for example, by establishing equivalence between different visualisation types – would enable the seamless adaptation of how outcomes are presented to a human reader.

A further opportunity is the assessment of *meaningfulness* of a specific type of analysis approach to check whether the computations making up the analysis of a specific ERC are meaningful (and not just mathematically possible). Spatiotemporal properties can play a key role in determining whether a specific computation is meaningful. In the context of spatial prediction and aggregation an initial approach of this type has been proposed [30]. With the OpenRIG and a broad availability of ERCs, this idea could be extended to other (qualitative) calculi such as RCC and use the metadata of an ERC and the common semantics of the

OpenRIG as a basis. Eventually, this could even become a safeguard that is incorporated into the validation process of ERCs, e.g. upon submission to an outlet.

The solutions presented here also come with several challenges. Some papers published in Geoinformatics do not rely on computational analysis but on qualitative data and analysis, e.g. interviews. While such publications thus are not executable, they still include data. One option to integrate them with the proposed approach could be to extend the bindings concept to guide readers through the data analysis, for example, by connecting statements in the results section of the paper with (anonymised) quotes and the aggregated higher-level themes. Still, for some types of scientific output, e.g. purely conceptual papers, although relevant in general, the proposed approach is potentially less beneficial. The higher-level challenge is of course to create opportunities and benefits so good that it incentivises researchers in Geoinformatics to actively adopt open science principles. Over the last decade we have learned that only talking about open science has not led to a substantial change in research practice, despite broad agreement about research ethics. Changing practice will need a concerted effort from scholars, reviewers, publishers, and libraries as well as science-funding bodies. Creating rewards for open science activities that go beyond publishing text papers is an important component of this. All of these parties must collaborate to operate the OpenRI(G), or instances of it. It would make sense to extend and connect existing building blocks to tackle the financial and organisational issues of operating such an infrastructure. These building blocks include, for example, review and publishing services by libraries or scholarly societies, computing resources of research institutes, data hosting by domain observatories, and funding schemes for sustainable software development or open access journals. We also envision that an OpenRIG would be useful to scholars from many disciplines. It would thus make sense to extend an existing research infrastructure such as Zenodo[5] with spatiotemporal search- and link-capabilities. That would provide this functionality to scholars in any discipline where spatially and temporally referenced data are used, including hydrologists, ecologists, meteorologists, geographers, archaeologists, and so on.

## 5 Concluding Remarks

In this paper, we propose the development and the adoption of an *open research infrastructure for Geoinformatics* to help us, scientists, to transform the centuries old process of *only* sharing textual and pictorial descriptions of our research findings into one where also the data and the data analysis *procedures* are shared, comprehensibly and reproducibly. This will not only have the advantage of increased transparency and enable more trust in science in general, but also create new options for interacting with the data and procedures, searching and finding particular datasets or applications of methods, and ways to link together research components. This also directly applies for some of the research that is reported at COSIT, e.g. new calculi or algorithms to tackle specific spatiotemporal problems and apply them to example scenarios to demonstrate their usefulness. In addition, the vision outlined in this paper provides new opportunities for fundamental research in spatial information theory, e.g. in terms of designing a semantic framework capturing spatiotemporal aspects of ERCs so that deep reasoning about multiple ERCs is enabled. Previous work shows that the technical realisation of such an infrastructure is by all means possible [20, 21, 14, 15]. In order to make this rather disruptive proposal for a transformation to open science a reality, the key challenge is a social one: although scientists agree that openness is essential, most of them

---

[5] `https://zenodo.org`, accessed on August 21, 2019.

hesitate to bear the (initial) costs themselves. As a starting point, we suggest that scientists who review manuscripts reporting on computational research start to decline doing this in cases where data and reproducible procedures are not made available to the reviewers.

### References

**1** Nick Barnes. Publish your computer code: it is good enough. *Nature News*, 467(7317):753–753, October 2010. `doi:10.1038/467753a`.

**2** Philip E. Bourne, Jessica K. Polka, Ronald D. Vale, and Robert Kiley. Ten simple rules to consider regarding preprint submission. *PLOS Computational Biology*, 13(5):e1005473, May 2017. `doi:10.1371/journal.pcbi.1005473`.

**3** Jonathan B. Buckheit and David L. Donoho. WaveLab and Reproducible Research. In Anestis Antoniadis and Georges Oppenheim, editors, *Wavelets and Statistics*, number 103 in Lecture Notes in Statistics, pages 55–81. Springer New York, 1995. `doi:10.1007/978-1-4612-2544-7_5`.

**4** Andy Cockburn, Carl Gutwin, and Alan Dix. HARK No More: On the Preregistration of CHI Experiments. In *Proceedings of CHI 2018*, pages 141:1–141:12, New York, NY, USA, 2018. ACM. `doi:10.1145/3173574.3173715`.

**5** Joëlle Coutaz. User interface plasticity: Model driven engineering to the limit! In *Proceedings of ACM EICS 2018*, pages 1–8. ACM, 2010.

**6** David DiBiase, Alan M MacEachren, John B Krygier, and Catherine Reeves. Animation and the role of map design in scientific visualization. *Cartography and geographic information systems*, 19(4):201–214, 1992. `arXiv:https://doi.org/10.1559/152304092783721295`.

**7** David Donoho. 50 Years of Data Science. *Journal of Computational and Graphical Statistics*, 26(4):745–766, October 2017. `doi:10.1080/10618600.2017.1384734`.

**8** Matt Duckham. GI Expertise. *Transactions in GIS*, 19(4):499–515, 2015. `doi:10.1111/tgis.12166`.

**9** Steve M Easterbrook. Open code for open science? *Nature Geoscience*, 7(11):779, 2014. `doi:10.1038/ngeo2283`.

**10** Esri. ArcGIS Pro product page. `https://www.esri.com/en-us/arcgis/products/arcgis-pro/overview`. Accessed April 15, 2019.

**11** Virginia Gewin. Data sharing: An open mind on open data. *Nature*, 529(7584):117–119, January 2016. `doi:10.1038/nj7584-117a`.

**12** Simon Hettrick. It's impossible to conduct research without software, say 7 out of 10 UK researchers. *Software and research*, 5:1536, 2014.

**13** Daniel S. Katz and Neil P. Chue Hong. Software Citation in Theory and Practice. In James H. Davenport, Manuel Kauers, George Labahn, and Josef Urban, editors, *Mathematical Software – ICMS 2018*, Lecture Notes in Computer Science, pages 289–296. Springer International Publishing, 2018. `doi:10.1007/978-3-319-96418-8_34`.

**14** Markus Konkol and Christian Kray. In-depth examination of spatio-temporal figures in open reproducible research. *Cartography and Geographic Information Science*, 2018. `doi:10.1080/15230406.2018.1512421`.

**15** Markus Konkol, Christian Kray, and Max Pfeiffer. Computational reproducibility in geoscientific papers: Insights from a series of studies with geoscientists and a reproduction study. *International Journal of Geographical Information Science*, pages 1–22, 2018. `doi:10.1080/13658816.2018.1508687`.

**16** Alan M MacEachren and Menno-Jan Kraak. Research challenges in geovisualization. *Cartography and geographic information science*, 28(1):3–12, 2001.

**17** Florian Markowetz. Five selfish reasons to work reproducibly. *Genome Biology*, 16:274, 2015. `doi:10.1186/s13059-015-0850-7`.

**18** MathWorks. Matlab product page. `https://www.mathworks.com/products/matlab.html`. Accessed April 15, 2019.

**19**   Daniel Nüst, Carlos Granell, Barbara Hofer, Markus Konkol, Frank O Ostermann, Rusne Sileryte, and Valentina Cerutti. Reproducible research and GIScience: an evaluation using AGILE conference papers. *PeerJ Preprints*, 6:e26561v1, 2018.

**20**   Daniel Nüst, Markus Konkol, Edzer Pebesma, Christian Kray, Marc Schutzeichel, Holger Przibytzin, and Jörg Lorenz. Opening the publication process with executable research compendia. *D-Lib Magazine*, 23(1/2), 2017.

**21**   Daniel Nüst. Reproducibility Service for Executable Research Compendia: Technical Specifications and Reference Implementation, December 2018. `doi:10.5281/zenodo.2203844`.

**22**   Hyoungjoo Park and Dietmar Wolfram. An examination of research data sharing and re-use: implications for data citation practice. *Scientometrics*, 111(1):443–461, April 2017. `doi:10.1007/s11192-017-2240-2`.

**23**   Heather Piwowar. Altmetrics: Value all research products. *Nature*, 493:159, January 2013. `doi:10.1038/493159a`.

**24**   Heather Piwowar, Jason Priem, Vincent Larivière, Joan Pablo Alperin, Lisa Matthias, Bree Norlander, Ashley Farley, Jevin West, and Stefanie Haustein. The state of OA: a large-scale analysis of the prevalence and impact of Open Access articles. *PeerJ*, 6(e4375), February 2018. `doi:10.7717/peerj.4375`.

**25**   Heather A. Piwowar and Todd J. Vision. Data reuse and the open data citation advantage. *PeerJ*, 1:e175, October 2013. `doi:10.7717/peerj.175`.

**26**   Hardy Pundt and Fred Toppen. 20 Years of AGILE. In Arnold Bregt, Tapani Sarjakoski, Ron van Lammeren, and Frans Rip, editors, *Societal Geo-innovation*, pages 351–367. Springer, Cham, 2017. `doi:10.1007/978-3-319-56759-4_20`.

**27**   Ulrich Pöschl. Interactive open access publishing and public peer review: The effectiveness of transparency and self-regulation in scientific quality assurance. *IFLA Journal*, 36(1):40–46, March 2010. `doi:10.1177/0340035209359573`.

**28**   R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2018. URL: `https://www.R-project.org/`.

**29**   David A Randell, Zhan Cui, and Anthony G Cohn. A spatial logic based on regions and connection. *KR*, 92:165–176, 1992.

**30**   Christoph Stasch, Simon Scheider, Edzer Pebesma, and Werner Kuhn. Meaningful spatial prediction and aggregation. *Environmental Modelling & Software*, 51:149–165, 2014.

**31**   Mark D. Wilkinson *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, 3:160018, March 2016. `doi:10.1038/sdata.2016.18`.

**32**   Greg Wilson, Jennifer Bryan, Karen Cranston, Justin Kitzes, Lex Nederbragt, and Tracy K. Teal. Good enough practices in scientific computing. *PLOS Computational Biology*, 13(6):e1005510, June 2017. `doi:10.1371/journal.pcbi.1005510`.