

Janiform Intra-Document Analytics for Repeatable Research

Jens Dittrich

Patrick Bender

Saarland University
infosys.cs.uni-saarland.de

ABSTRACT

Peer-reviewed publication of research papers is a corner stone of science. However, one of the many issues of our publication culture is that our publications only publish a snapshot of the final result of a long project. This means, we put well-polished graphs describing (some) of our experimental results into our publications. However, the algorithms, input datasets, benchmarks, raw result datasets, as well as scripts that were used to produce the graphs in the first place are rarely published and typically not available to other researchers. Often they are only available when personally asking the authors. In many cases, however, they are not available at all. This means from a long workflow that led to producing a graph for a research paper, we only publish the final result rather than the entire workflow. This is unfortunate and has been lamented upon in various scientific communities. In this demo we argue that one part of the problem is our dated view on what a “document” and hence “a publication” *is*, *should*, and *can be*. As a remedy, we introduce portable database files (PDbf). These files are janiform, i.e. they are at the same time a standard static pdf as well as a highly dynamic (offline) html-document. PDbfs allow you to access the raw data behind a file, perform portable OLAP-style analysis, and reproduce your own graphs from the raw data — all of this *within* a portable document. We demo a tool allowing you to create PDbfs smoothly from within \LaTeX . This tool allows you to preserve the connection of raw data to its final graphical output through all stages of the workflow. Notice that this pdf already showcases our technology: rename this file to “.html” and see what happens (currently Firefox, Chrome, or Safari on a Desktop machine).

1. INTRODUCTION

Irreproducibility is a problem frequently lamented upon in various scientific communities [?, ?]. In the context of computer science it has recently been coined “The Real Software Crisis” [?]. The database community has identified it more than ten years ago and is attacking it through repeatability committees, e.g. [?]. In these committees a separate committee reruns the experiments of accepted papers using the datasets and code provided by the authors. Obviously, given the sheer size and complexity of some projects, in

many cases these boils down to blackbox testing, i.e. it can neither be tested if the tested code actually implements the algorithms presented in the paper nor whether the code reports results in a proper way. Another problem of repeatability committees is that they cannot remedy inherent publication bias: “reviewers don’t like negative results”. Hence, for an experimental evaluation you need the “right” queries, the “right” dataset and the “right” baselines. This naturally leads to a flood of papers with positive results. And to papers where the “improvements don’t add up” [?]. Publication bias was attacked by the inauguration of Experiments&Analysis papers at (P)VLDB. These kind of papers reevaluate existing work in a uniform setting and may also publish negative results. These experimental evaluations may then serve as landmarks in the flood of papers with (overly) positive results giving clear advice on the strength and weaknesses of a particular method.

This small demo is neither the place to even summarize nor defend the different arguments in the debate on repeatability and our experimental culture. It is an emotional topic where the esteemed reader of these lines probably has strong opinions in one way or the other. This is just fine. In the following, we will simply accept that there is a problem [?, ?]¹. And that this problem is calling “for a new model for the way how we publish our results” [?]. Handling this problem can be regarded an instance of “small data” [?, ?]².

2. SIGNIFICANCE OF THE CONTRIBUTION

We believe that our contribution is significant for the following reasons:

1. **Portable DataBase Files.** We provide Portable DataBase Files a general model to overlay a static pdf document with additional highly dynamic content. In order to specify an overlay, we simply require access to the static document S , the dynamic content D , and a PDbf-configuration file C defining where to place the dynamic content.
2. **PDbf-Compiler.** We provide a compiler taking as its input the triple (S, D, C) . Our compiler outputs a janiform document. That document is **at the same time** a valid pdf **and** a valid html-document. Thus, if you open the file with a pdf-viewer, you will see the static content, i.e. only the S -part. However, if you rename the file to “.html” and open it with a Web browser, you will be able to inspect the dynamic part, i.e. D and S .

¹Just read the “ten simple rules for reproducible results” [?] and then ask yourself how little of this we follow for our papers.

²Also notice an upcoming Dagstuhl perspectives workshop on Artifact Evaluation for Publications [?] where the first author participates.

3. **Full L^AT_EX integration.** We instrument L^AT_EX to output not only the static pdf-file S , but also a valid PDBF-configuration file C . In addition, we provide an extension to L^AT_EX allowing you to create graphs directly from within L^AT_EX — without requiring additional tools. Like that we are able to preserve the connectivity of raw data and graphs *through the L^AT_EX compiler*. In addition, this process creates dynamic variants of the graphs as a side-effect, i.e. the D -part. This means, the user can create her graphs seamlessly without worrying about different representations of graphs in static and dynamic representations. The result of this instrumentation is again a triple (S, D, C) which can be fed into our PDBF-Compiler (see Contribution 2) to create a janiform PDBF-document.
4. **Preservation of Raw Data and Graph-Connectivity.** Our compilers preserve the connection between raw data and the graphs and/or tables produced from that raw data. In addition, we are able to ship that part of the workflow, i.e. data, graphs, and the code producing the graphs within a “document”.
5. **Alternative Dynamic Views on the Data.** Our technology allows you to perform **offline** OLAP-style analytics on the **data shipped within the document**. All you need is a Web Browser (currently Firefox, Chrome, or Safari on a Desktop machine). We support a rich feature list. See Section 3.1 for our currently supported features (as of March 31, 2015).
6. **Longterm Preservation of Raw Data.** As our tool embeds the raw data within the publication, PDBF-documents naturally archive the raw data with the document. Therefore, the raw data may be “downloaded” directly from within the PDBF-document.
7. **Impact On Research in General.** We believe that our technology may not only be interesting to the database community. Our tool may be interesting for all research communities working with experimental data. Therefore, we are planning to open source our tool upon publication of this demo paper.

3. THE PDBF FRAMEWORK

PATRICK

core compilation steps (with graphs could use the ones I did for my BTW talk)

different output format options

pros and cons of those output formats

3.1 Currently Supported Features

1. support for DESKTOP browsers

3.2 Future Work

- 1.

4. THE DEMO

We invite the reader to change the suffix “.pdf” of this submission to “.html”. Your browser will open and you will be able to see the dynamic features. Notice that the dynamic content is completely offline and runs entirely in your Web browser’s sandbox.

PATRICK

4.1 Feature 1...

Switch between data representations for better visualization
Chart, Pivot Table, Table

4.2 Feature 2...

Get raw data, source code
raw data can be exported as csv with header
maybe other file attachments?

4.3 Feature 3...

Analyze data in browser Remove filtering functions from graphs
Look which method performed best under certain condition

4.4 Else

JD: bitte latex code-Beispiel zur Erzeugung eines Graphen zeigen

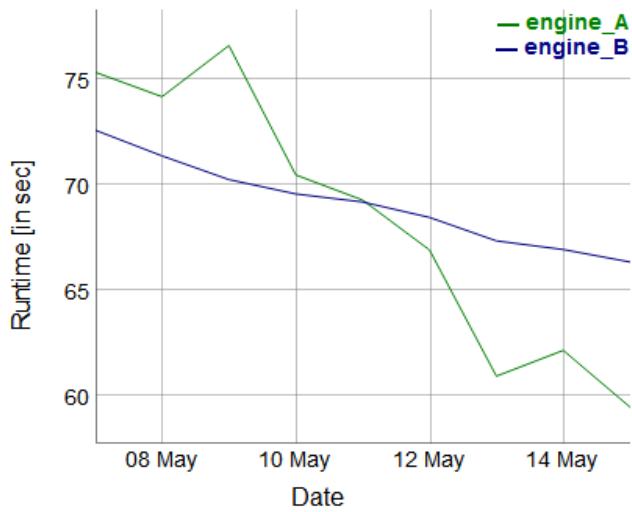


Figure 1: Multi-column Line Chart

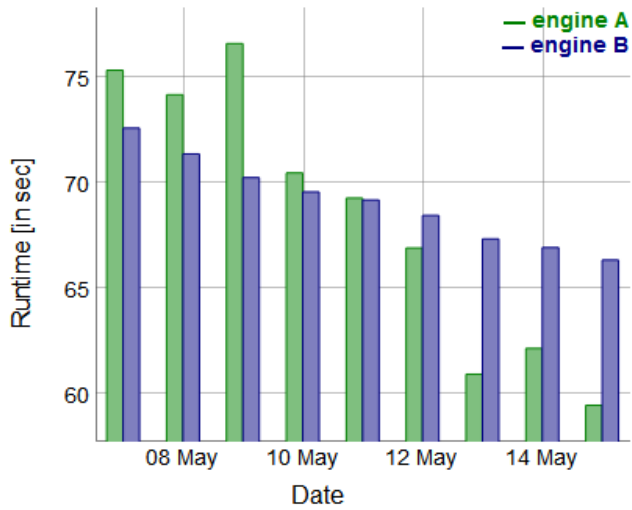


Figure 2: Multi-column Bar Chart

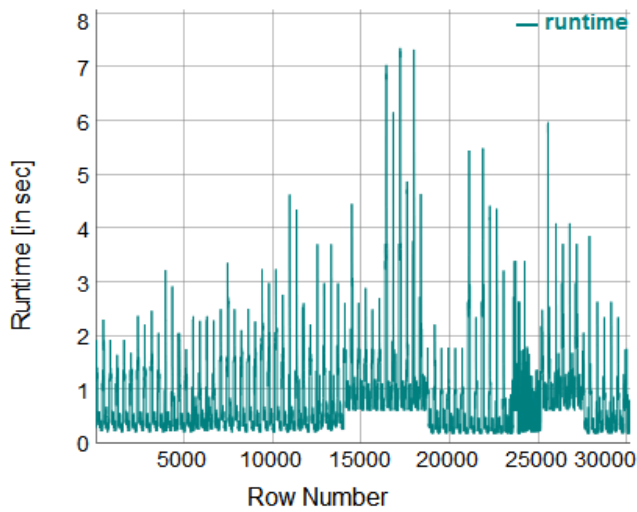


Figure 3: Lot of points

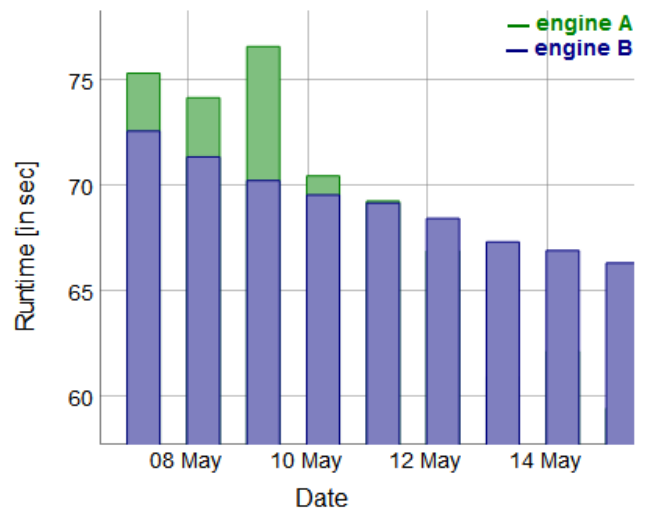


Figure 4: Multi-column overlapping Bar Chart

Figure 5: Pivot table

5. RELATED WORK

hmm, how much should we really say here?
SIGMOD reproducibility initiatives and papers (two, right?)
upcoming Dagstuhl perspectives workshop
experiments and analysis track and its problems: access to data,
access to code
Pweave <http://mpastell.com/pweave/index.html>

6. CONCLUSIONS

This demo opened the book for portable database files (PDbF). PDbfs allow you to embed raw data into a document while preserving the entire data-pipeline from raw data to graphs. This allows readers to perform in-document OLAP-style analytics on the published document. PDbFs are janiform documents that are at the same time a valid pdf and a html document. Thus they can be viewed in either way. This pdf is an example of such a janiform PDbF. In addition, the conference demo show-cases two compilers allowing you to seamlessly create PDbFs from within \LaTeX .

As part of future work we want to research how to include other parts of the research pipeline as part of the pdf. For instance, in future, one might consider to embed a virtual machine emulator, an operating system image, and all software that is required to run the original code *within the pdf* (~3GB of data). This may sound infeasible in 2015. But 4K video streaming over the Internet sounded equally silly in 1995.

7. REFERENCES