

SOFTWARE

Building applications for interactive data exploration in systems biology

Bjørn Fjukstad¹, Vanessa Dumeaux³, Karina Standahl Olsen², Eiliv Lund², Michael Hallett³ and Lars Ailo Bongo^{1*}

Abstract

Background: With the advent of high-throughput genomics, there is a growing need for interactive data exploration tools as the systems biology community collects data at an unprecedented rate. These tools need to combine advanced statistical analyses, relevant prior knowledge, and interactive visualizations in an application with simple user interfaces. Although data analyses, user interfaces and visualizations are specialized, application developers can share and reuse the underlying components of the application. To answer specific research questions tools must provide specialized user interfaces and visualizations. While these are application-specific, the underlying components of a data analysis tool can be shared and reused later. Application developers can therefore compose applications of reusable services rather than implementing a single monolithic application.

Results: Our approach for developing data exploration applications in systems biology builds on the microservice architecture. The resulting applications integrate advanced statistical software, up-to-date information from biological databases and modern data visualization libraries. We demonstrate the viability through the MlxT Blood-Tumor web application that explore and compare transcriptional profiles from blood and tumor samples in breast cancer patients. This approach was reused in two other web applications and several command-line tools. Thus our microservice approach building on software container technology enables re-use and sharing of key components between application reducing development, deployment and maintenance time.

Conclusions: Our approach and reference implementation Kvik is open-sourced at github.com/fjukstad/kvik. The web application for exploring transcriptional profiles, MlxT, is available at mixt-blood-tumor.bci.mcgill.ca and its source code at github.com/fjukstad/mixt.

Keywords:

Background

In recent years the biological community has generated and collected an unprecedented amount of data. While the cost of data collection has drastically decreased, data analysis continues to be a larger fraction of the total experiment cost.[1] This calls for novel methods in data analysis and exploration.

Several tools are now available in various programming languages for data analysis. These include novel statistical/bioinformatics methodologies and graphical analysis tools. In the R statistical programming language developers share software through repositories

such as CRAN^[1] or Bioconductor^[2]. In other languages, libraries for biological computation are often available like BioPython^[?] and biogo^[?] written in Python or Go, respectively.

Although these frameworks are of tremendous help for scientists with some computing know-how, even computationally savvy researchers can get tangled up when wrestling with software and big data. Novel applications are needed to better integrate statistical packages, biological databases, and interactive visualizations. Different programming languages solve different tasks. For example, new biological data analysis techniques are quickly released in R and its package repositories, high performance computer vision tasks

* Correspondence: larsab@cs.uit.no

¹Department of Computer Science, UiT – The Arctic University of Tromsø, 9037 Tromsø, NO

Full list of author information is available at the end of the article

^[1] cran.r-project.org

^[2] bioconductor.org

are performed optimally in C++ and OpenCV, and portable user interfaces more easily built in HTML, CSS and JavaScript. Therefore applications that integrate novel statistical analysis tools, interactive visualizations, and biological databases likely need to include several components written in different languages.

A microservice architecture structures an application into small reusable, loosely coupled parts. These communicate via lightweight programming language-agnostic protocols such as HTTP, making it possible to write single applications in multiple different programming languages. This way the most suitable programming language is used for each specific part. To build a microservice application, developers bundle each service in a container. Containers are built from configuration files which describe the operating system, software packages and versions of these. This makes reproducing an application a trivial task. The most popular implementation of a software container is Docker^[3], but others such as Rkt^[4] exist. Initiatives such as BioContainers^[5] now provide containers pre-installed with different bioinformatics tools.

In this paper, we describe a novel approach for building data exploration applications in systems biology. We show that by building applications as a set of services it is possible to reuse and share its components between applications. In addition, by packaging the services using container technology we promote reproducible research and simplify application deployment. We have used our approach to build a number of applications, both command-line and web-based. In this paper we describe how we used our approach to develop MIXT, a web application for exploring and comparing transcriptional profiles from blood and tumor samples. The MIXT web application integrates statistical analysis together with biological databases and interactive visualizations.

Requirements

First, we identified a set of reusable services that application developers can use to build a wide range of applications. The key services of a biological data exploration application are i) a compute service for executing statistical analyses in languages such as R, and ii) a database query service for retrieving information from biological databases.

To build these services we need a framework that fulfills the following requirements:

- 1 It provides a language-independent approach for integrating, or embedding, statistical software, such as R, directly in interactive data exploration applications.

- 2 It has an interface to online reference databases to provide meta-data to biological entities.
- 3 Its components should be easy to develop, maintain, deploy and share between projects.
- 4 It scales to meet the performance requirements of interactive data exploration applications.

Related Work

In this section we aim to cover some of the existing systems for building interactive data exploration applications in systems biology.

Integrate Statistical Analyses

OpenCPU is a system for embedded scientific computing and reproducible research.^[2] It offers an HTTP API to the R programming language to provide an interface with statistical methods. It enables users to make function calls to any R package and retrieve the results in a wide variety of formats such as json or pdf. Users can choose to host their own R server or use public servers, and OpenCPU works in a single-user setting within an R session, or a multi-user setting facilitating multiple parallel requests. This makes OpenCPU suitable for building a service that can run statistical analyses. OpenCPU provides a Javascript library for interfacing with R, as well as Docker containers for easy installation and OpenCPU has been used to build multiple applications.^[6]

Renjin is a JVM-based interpreter for the R programming language.^[3] It targets developers who want to integrate the R interpreter in web applications. Since it is built on top of the JVM it allows developers to write data exploration applications that interact directly with R code, both running on top of the JVM. Although Renjin supports a large number of CRAN packages it cannot access any R package (i.e. any package from BioConductor or CRAN) without modification. This makes the programming effort to use Renjin as an interface to R higher.

Shiny is a web application framework for R.^[7] It allows developers to build web applications in R without having to have any knowledge about HTML, CSS or Javascript. Its widget library provides more advanced Javascript visualizations such as Leaflet^[8] for maps or three.js^[9] for WebGL-accelerated graphics. Developers can choose to host their own web server with the user-built Shiny Apps, or host them on public servers. Shiny forces users to implement data exploration applications in R, limiting the functionality to the widgets and libraries in Shiny.

^[3] docker.com

^[4] coreos.com/rkt

^[5] biocontainers.pro

^[6] opencpu.org/apps.html

^[7] shiny.rstudio.com

^[8] leafletjs.com

^[9] threejs.org

Biogo

this one is a bit out of place.

Biogo is a bioinformatics library written in Go. It provides functionality to analyze genomic and metagenomic datasets in the go programming language.[4] Using the go programming language the developers are able to provide high-performance parallel processing in a clean and simple programming language.

Visualization frameworks

Cytoscape is an open source software platform for visualizing complex networks and integrating these with any type of attribute data[5]. It allows for analysis and visualization in the same platform. Users can add additional features, such as databases connections or new layouts, through Apps. One such app is cyREST which allows external network creation and analysis through a REST API[6]. To bring the visualization and analysis capabilities to the web the creators of Cytoscape have developed Cytoscape.js^[10], a Javascript library to create interactive graph visualizations.

Caleydo is a framework for building applications for visualizing and exploring biomolecular data[?]. Until 2014 it was a standalone tool that needed to be downloaded, but the Caleydo team are now making the tools web-based. There have been several applications built using Caleydo: StratomeX for exploring stratified heterogeneous datasets for disease subtype analysis[7]; Pathfinder for exploring paths in large multivariate graphs[8]; UpSet to visualize and analyse sets, their intersections and aggregates[9]; Entourage and enRoute to explore and visualize biological pathways [10][11]; LineUp to explore rankings of items based on a set of attributes[12]; and Domino for exploring subsets across multiple tabular datasets[13].

BioJS is an open-source JavaScript framework for biological data visualization.[14] It provides a community-driven online repository with a wide range components for visualizing biological data contributed by the bioinformatics community. BioJS builds on node.js^[11] providing both server-side and client-side libraries.

WIP: Biological Databases

Maybe some words here on how to get data out of the different biological databases?

WIP: Microservices, Docker etc.

...

^[10] js.cytoscapejs.org

^[11] nodejs.org

Methods

In this section we first motivate our microservice approach based on our experiences developing the MIXT web application. We describe the process from initial data analysis to the final application, highlighting the importance of language-agnostic services to facilitate the use of different tools in different parts of the application. We then generalize the ideas to a set of principles and services that can be reused and shared between applications.

Matched Interactions Across Tissues (MIxT)

The aim of the MIXT study was to identify genes and pathways in the primary breast tumor that are tightly linked to genes and pathways in the patient blood cells. We generated and analyzed expression profiles from blood and matched tumor cells in 173 breast cancer patients included in the Norwegian Women and Cancer (NOWAC) study.[?] The MIXT analysis starts by identifying sets of genes tightly co-expressed across all patients in each tissue. Each group of genes or modules were annotated based on known a priori biological knowledge about gene functionality. Focus was placed on the relationships between tissues by asking if specific biologies in one tissue are linked with (possibly distinct) biologies in the second tissue, and this within different subgroup of patients (i.e. subtypes of breast cancer).

An R package^[12] was built to contain the statistical methods and static visualizations for identifying associations between modules across tissues. The exploration of the results encompass the examination of ~ 20 modules and their functional enrichments. That is $\sim 23 \times 19 = 437$ associations computed for each of the 22 patient subgroups. An application integrating many types of information was therefore needed to explore this large amount of data and results. The system needed to directly interface with the R package should be accessible through standard protocols, such as HTTP, not enforcing any programming language or platform on the application developer.

A large part of biological data research is to link the results to known biology from literature or reference databases. In the MIXT project we needed an application that could interface with a set of different databases, keeping the information up-to-date. We interface with MSigDB to get gene set meta-data and the Entrez Programming Utilities (E-utils) to get gene meta-data. As with the R interface the interface to the databases should be accessible from any programming language.

A key feature that motivated the design of MIXT was that we want to have the flexibility to run the application on any platform. Bundling each component in a

^[12] github.com/vdumeaux/mixt

software container such as Docker allows us to deploy the application on a wide range of hardware, from local installations to deployments to cloud providers such as Amazon Web Services^[13].

Kvik

Our experience in developing Mlxt can be generalized into the following design principles and microservices encapsulated in the Kvik framework.

Principle 1: Build applications as collections of language-agnostic microservices. This makes it possible to re-use key components and build specialized data exploration applications in the most suitable programming language.

Principle 2: Deploy each service using container technology such as Docker. This has a number of benefits. It simplifies deployment itself, it makes it trivial to share services between projects and research groups, and it ensures reproducible services.

Principle 3: Package statistical methods and data as software packages that can be used by power-users and the data exploration tools themselves. An example is to build an application using R packages and OpenCPU or Kvik. This makes it possible to either explore the data and methods through the data exploration application or an R session.

From these three main principles we built a set of software packages to build microservices used in data exploration application in systems biology. We built a compute service for executing statistical analyses and a database service for retrieving relevant information on genes and biological processes. Using these it is possible to develop specialized data exploration application in any modern programming language.

Building applications

Kvik provides an interface to perform database lookup and execute statistical analyses. Application developers focus on implementing the application-specific user interfaces and statistical analyses, using the microservices from Kvik to execute the analyses and get meta-data for the results.

Figure 1 illustrates the relationship between the Mlxt application and Kvik. In Mlxt we built a specialized web application that interfaces with Kvik to get data from biological databases and to run statistical analyses from the mixtR R package.

It is important to note that since the end applications interface directly with R, developers can leverage this to produce dynamic visualizations. For example, if an application uses a clustering method to color nodes

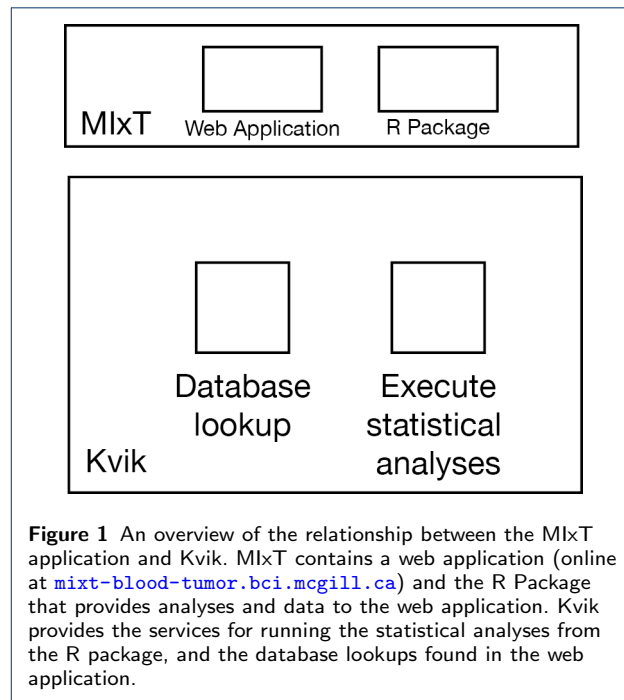


Figure 1 An overview of the relationship between the Mlxt application and Kvik. Mlxt contains a web application (online at mixt-blood-tumor.bci.mcgill.ca) and the R Package that provides analyses and data to the web application. Kvik provides the services for running the statistical analyses from the R package, and the database lookups found in the web application.

in a graph, end-users can tweak parameters that interactively change the node coloring in graph visualization.

In Kvik we do not specify what programming language or set of tools to use to build an application. We believe that by building an application as a set of services that communicate over standard protocols, application developers can choose the most suitable language, or framework, to build specialized visualizations and user interfaces.

Compute Service

The compute service in Kvik is built using a hybrid state pattern^[2]. We provide three main operations for interfacing with R: Call, Get, and RPC. The Call operation is used to execute and run a function from an R package. It takes as input an R package name, a function name, and optional arguments. It returns a unique identifier that later can be used by the Get operation to retrieve results. The Get operation is used to get results in different output formats, e.g. JSON, CSV, PDF, or PNG. The RPC is just a combination of a Call and a subsequent Get.

Database Service

The database service provides an interface to biological databases for retrieving meta-data on genes and processes. In Kvik we have built packages for interfacing with E-utilities^[14], MSigDB^[15], Hugo Gene Nomencla-

^[13] aws.amazon.com

^[14] eutils.ncbi.nlm.nih.gov

^[15] software.broadinstitute.org/gsea/msigdb

ture Committe (HGNC)^[16], and Kyoto Encyclopedia of Genes and Genomes (KEGG)^[17].

TODO: Talk about licensing issues.

The database service uses a caching mechanism to reduce the load on the online databases. It will also speed up subsequent queries for a cached object, since the query can be served out of cache and not having to be fetched from a remote database. We allow application developers to specify the cache eviction policy, but on default the database service does not evict anything from its cache before the service is restarted. This can be modified by the application developer.

Implementation

In this section we describe the implementation details in Kvik and the microservices required to build the MlXt web application.

Kvik is implemented as a collection of Go packages required to build services that can integrate statistical software in a data exploration and provide an interface to up-to-date biological databases. We chose the Go programming language because of its performance, ease of development, and simple deployment. To integrate R we provide two packages *gopencpu* and *r*, that interface with OpenCPU and Kvik R servers respectively. To interface with biological databases we provide the packages *eutils*, *gsea*, *genenames*, and *kegg* that interface with E-utils, MsigDB, HGNC and KEGG respectively. In addition to these packages we provide Docker images that implement the two required microservices.

Both the compute and the databases service in Kvik builds on the standard *http* library in Go. On start the compute service launches a user-defined number of R sessions that execute analyses on demand. This allows for parallel execution of analyses. We provide a simple FIFO queue for queuing of requests. The compute service also provides the opportunity for users to cache analysis results to speed up subsequent calls. The database service use the *gocache*^[18] package to cache any query to an online database.

Results

We show the viability of the microservices approach in Kvik by describing the MlXt web application for exploring and comparing transcriptional profiles from blood and tumor samples.

We define six data exploration tasks that the web application should help users perform:

Explore co-expression gene sets in tumor and blood tissue. We want to simplify the process of

exploring the computed co-expression gene sets, or modules, through the web-application. The application should visualize gene expression patterns together with clinicopathological variables for each module. In addition we want to enable users to study the underlying biological functions of each module by including gene set analyses between the module genes and known gene sets.

Explore co-expression relationships between genes. Users should be able to explore the co-expression relationship as a graph visualization. The network should visualize each gene as a node and a significant co-expression relationship as an edge.

Explore relationships between modules from each tissue. Users should be able to explore the relationship between modules from different tissues. We provide two different metrics to compare modules, and the web application should enable users to interactively browse these relationships. In addition to providing visualizations the compare modules from each tissue, users should also be able to explore the relationships, but for different breast cancer patient groups.

Explore relationships between clinical variables and modules. In addition to comparing the association between modules from both tissues, users also have the possibility of exploring the association with a module and a specific clinical variable. This should also be possible for the different breast cancer patient group.

Explore association between user-submitted gene lists and computed modules. We want to enable users to explore their own gene lists to explore them in context of the co-expression gene sets. The web application must handle uploads of gene lists and compute association between the genelist and the MlXt modules on demand.

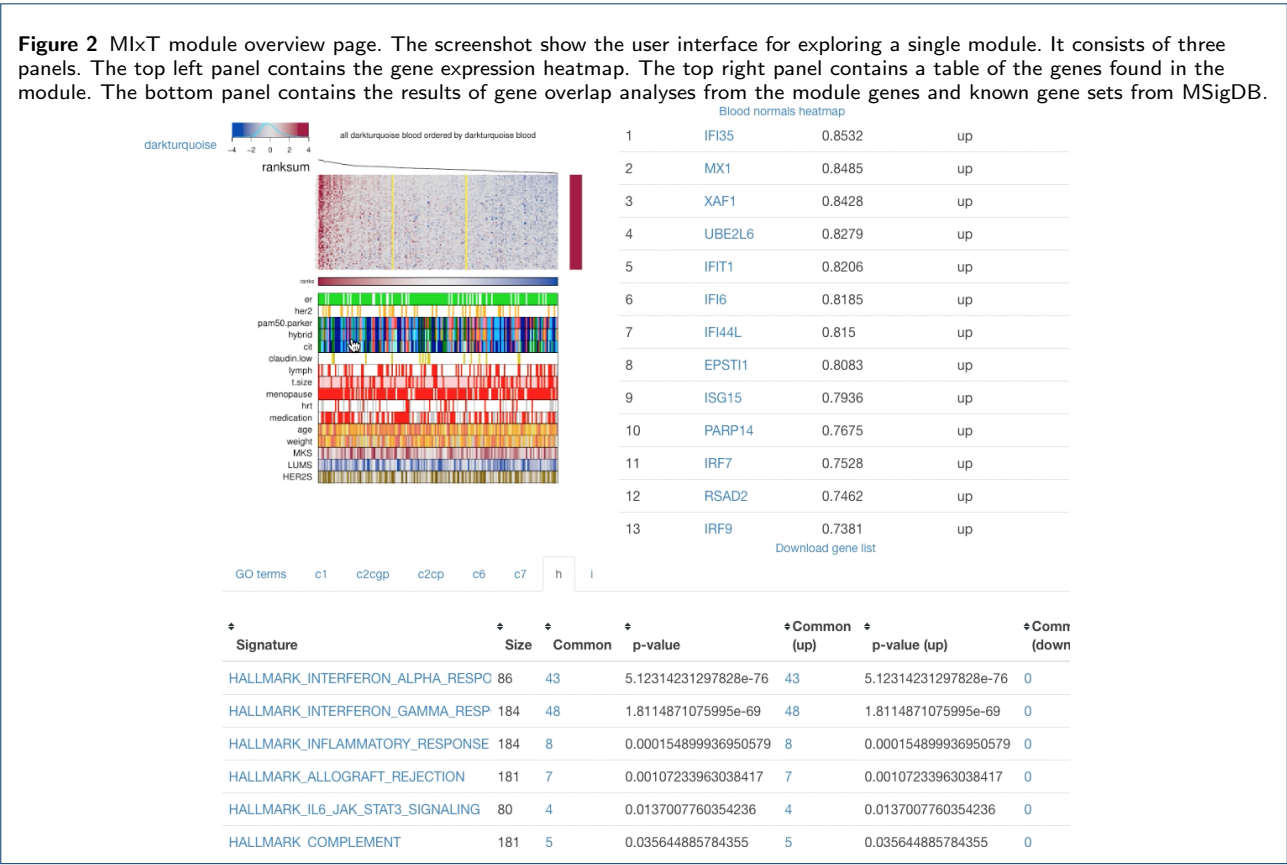
Search for genes or gene lists of interest. To facilitate faster lookup of genes and biological processes, the web application should provide a search functionality that lets users locate genes or gene lists and show association to the co-expression gene sets.

From these six analysis tasks we designed and implemented MlXt as a web application that integrates statistical analyses and information from biological databases together with interactive visualizations. The MlXt web application consists of three services: i) the web application itself containing the user-interface and visualizations; ii) the compute service performing the MlXt analyses delivering data to the web application; and iii) the database service providing up-to-date information from biological databases. Each of these services run within Docker containers making the process of deploying the application simple. By composing the application of a set of services we can substitute parts

^[16]genenames.org

^[17]kegg.jp

^[18]github.com/fjukstad/gocache



of the application without re-writing the entire application. For example if we want to use OpenCPU to interface with data analysis we can do so by simply exchanging the Kvik compute service with OpenCPU. Both services communicate over HTTP and their interfaces are the same.

We structured the MlXt application with a separate view for each analysis task.

To explore the co-expression gene sets we built a view that combines both static visualizations from R together with interactive tables with gene overlap analyses. Figure 2 shows the web page presented to users when they access the co-expression gene set 'darkturquoise' from blood. Using the Kvik compute service we can generate plots on demand and provide users with high-resolution PDFs or PNG files.

To explore the co-expression relationship between genes we use an interactive graph visualization build with Sigmajs^[19]. We have built visualization for both tissues, with graph sizes of 2705 nodes and 90 348 edges for the blood network, and 2066 nodes and 50 563 edges for the biopsy network. The sigmajs visualization library has functionality for generating a layout

for large networks, but we generate this layout server-side to reduce the computational load on the client. To generate this layout we use the GGally package^[20]. By generating the network layout using the compute service we relieve the clients.

To visualize relationships between modules from different tissues, or their relationship to clinical variables we built a heatmap visualization using the d3^[21] library. Figure 3 shows an example of this heatmap visualization, showing the association between the clinical variables and the modules from biopsy for all samples.

WIP: Discussion

The compute service in Kvik follows many of the design patterns in OpenCPU. Both systems interface with R packages using a hybrid state pattern over HTTP. Both systems provide the same interface to execute analyses and retrieve results. While OpenCPU is implemented on top of R and Apache, Kvik is implemented from the ground up in Go. Because of the similarities in the interface to R in Kvik we provide packages for interfacing with our own R server or OpenCPU R servers through the *gopencpu* package.^[22]

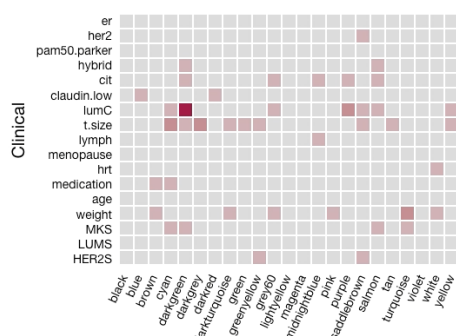
^[19] sigmajs.org

^[20] cran.r-project.org/web/packages/GGally

^[21] d3js.org

^[22] github.com/fjukstad/kvik/tree/master/gopencpu

Figure 3 Heatmap visualization of the association between clinical variables and the modules in biosy. The visualization is built using the d3 JavaScript library.



WIP: Future work

While we have built

We have built a database service that provides a sufficient interface for the MIXT web application. While we have developed the software packages for interfacing with more databases, these haven't been included in the database service yet.

One large concern that we haven't addressed in this paper is security. In particular one security concern that we plan to address in Kvik is the restrictions on

Docker compose

This includes both access restrictions to parts of an application, e.g. access to a database, and security concerns regarding remote code execution within the compute service.

to the biological databases that are

Security. Documentation. Scale.

WIP: Conclusions

We have designed an approach for building data exploration applications in systems biology that builds on a microservice architecture. Using this approach we have built a web application that leverages this architecture to integrate statistical analyses, interactive visualizations and data from biological databases.

List of abbreviations used

Competing interests

The authors declare that they have no competing interests.

Author details

¹Department of Computer Science, UiT – The Arctic University of Tromsø, 9037 Tromsø, NO. ²Department of Community Medicine, UiT – The Arctic University of Tromsø, 9037 Tromsø, NO. ³Department of Biology, Concordia University, H4B 1R6 Montréal, CA.

References

- Sboner, A., Mu, X.J., Greenbaum, D., Auerbach, R.K., Gerstein, M.B.: The real cost of sequencing: higher than you think! *Genome biology* **12**(8), 125 (2011)
- Ooms, J.: The opencpu system: Towards a universal interface for scientific computing through separation of concerns. *arXiv preprint arXiv:1406.4806* (2014)
- Bertram, A.: Renjin: The new r interpreter built on the jvm. In: *The R User Conference, useR! 2013 July 10–12 2013 University of Castilla-La Mancha, Albacete, Spain*, vol. 10, p. 105 (2013)
- Kortschak, R.D., Adelson, D.L.: biogo: a simple high-performance bioinformatics toolkit for the go language. *bioRxiv* (2014). doi:[10.1101/005033](https://doi.org/10.1101/005033). <http://biorxiv.org/content/early/2014/05/12/005033.full.pdf>
- Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., Ideker, T.: Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome research* **13**(11), 2498–2504 (2003)
- Ono, K., Muetze, T., Kolishovski, G., Shannon, P., Demchak, B.: Cyrest: Turbocharging cytoscape access for external tools via a restful api. *F1000Research* **4** (2015)
- Lex, A., Streit, M., Schulz, H.-J., Partl, C., Schmalstieg, D., Park, P.J., Gehlenborg, N.: Stratomex: Visual analysis of large-scale heterogeneous genomics data for cancer subtype characterization. In: *Computer Graphics Forum*, vol. 31, pp. 1175–1184 (2012). Wiley Online Library
- Partl, C., Gratzl, S., Streit, M., Wassermann, A.M., Pfister, H., Schmalstieg, D., Lex, A.: Pathfinder: Visual analysis of paths in graphs. In: *Computer Graphics Forum*, vol. 35, pp. 71–80 (2016). Wiley Online Library
- Lex, A., Gehlenborg, N., Strobel, H., Vuilleumot, R., Pfister, H.: Upset: visualization of intersecting sets. *IEEE transactions on visualization and computer graphics* **20**(12), 1983–1992 (2014)
- Lex, A., Partl, C., Kalkofen, D., Streit, M., Gratzl, S., Wassermann, A.M., Schmalstieg, D., Pfister, H.: Entourage: Visualizing relationships between biological pathways using contextual subsets. *IEEE transactions on visualization and computer graphics* **19**(12), 2536–2545 (2013)
- Partl, C., Lex, A., Streit, M., Kalkofen, D., Kashofer, K., Schmalstieg, D.: enroute: Dynamic path extraction from biological pathway maps for in-depth experimental data analysis. In: *Biological Data Visualization (BioVis)*, 2012 IEEE Symposium On, pp. 107–114 (2012). IEEE
- Gratzl, S., Lex, A., Gehlenborg, N., Pfister, H., Streit, M.: Lineup: Visual analysis of multi-attribute rankings. *IEEE transactions on visualization and computer graphics* **19**(12), 2277–2286 (2013)
- Gratzl, S., Gehlenborg, N., Lex, A., Pfister, H., Streit, M.: Domino: Extracting, comparing, and manipulating subsets across multiple tabular datasets. *IEEE transactions on visualization and computer graphics* **20**(12), 2023–2032 (2014)
- Gómez, J., García, L.J., Salazar, G.A., Villaveces, J., Gore, S., García, A., Martín, M.J., Launay, G., Alcántara, R., Ayllón, N.D.T., et al.: Biojs: an open source javascript framework for biological data visualization. *Bioinformatics*, 100 (2013)