

1 **MetaR: simple, high-level languages for** 2 **data analysis with the R ecosystem**

3 **Fabien Campagne^{1,2,3,*}, William ER Digan^{1,2}, and Manuele Simi^{1,2}**

4 ¹**The HRH Prince Alwaleed Bin Talal Bin Abdulaziz Alsaud Institute for Computational**
5 **Biomedicine, Weill Cornell Medical College, New York, NY, United States of America**

6 ²**Clinical Translational Science Center, Weill Cornell Medical College, New York, NY,**
7 **United States of America**

8 ³**Department of Physiology and Biophysics, Weill Cornell Medical College, New York,**
9 **NY, United States of America**

10 ***To whom correspondence should be addressed: fac2003@campagnelab.org**

11 **ABSTRACT**

12 Data analysis tools have become essential to the study of biology. Here, we applied language workbench technology (LWT) to create data analysis languages tailored for biologists with a diverse range of experience: from beginners with no programming experience to expert bioinformaticians and statisticians. A key novelty of our approach is its ability to blend user interface with scripting in a single platform. This feature helps beginners and experts alike analyze data more productively. This new approach has several advantages over state of the art approaches currently popular for data analysis: experts can design simplified data analysis languages that require no programming experience, and behave like graphical user interfaces, yet have the advantages of scripting. We report on such a simple language, called MetaR, which we have used to teach complete beginners how to call differentially expressed genes and build heatmaps. We found that beginners can complete this task in less than 2 hours with MetaR, when more traditional teaching with R and its packages would require several training sessions (6-24hrs). Furthermore, MetaR seamlessly integrates with docker to enable reproducibility of analyses and simplified R package installations during training sessions. We used the same approach to develop the first composable R language. A composable language is a language that can be extended with micro-languages. We illustrate this capability with a Biomart micro-language designed to compose with R and help R programmers query Biomart interactively to assemble specific queries to retrieve data, (The same micro-language also composes with MetaR to help beginners query Biomart.) Our teaching experience suggests that language design with LWT can be a compelling approach for developing intelligent data analysis tools and can accelerate training for common data analysis task. LWT offers an interactive environment with the potential to promote exchanges between beginner and expert data analysts.

13 **Keywords:** Data analysis, Bioinformatics, Language Workbench Technology, JetBrains MPS, Composable R, R language, Bioinformatics Education

14 **INTRODUCTION**

15 Present day biology often requires that biologists rely on software tools for data analysis. For instance,
16 software tools are required for analysis of high-throughput data, for the study of genome-wide gene
17 expression, genetic or epigenetic. Similarly, most fields of biology require specialized software tools
18 for analysis of microscopy, crystallography or other data. Most analysis software is constructed in a
19 very similar manner: writing a program as a collection of text source code compiled into one or more
20 executable analysis tools. Despite the evolution of programming languages, encoding programs as text
21 has been a constant since the invention of the first high-level programming language (FORTRAN Backus
22 [1958, 1978]).

23 In this manuscript, we discuss several drawbacks of encoding programs as text that we believe
24 make teaching of data analysis more difficult than necessary and contribute to some frequent challenges
25 encountered by data analysts. Language Workbenches (LWs) with projectional editors offer an alternative
26 platform to develop data analysis tools. LWs were conceived in the 90s Simonyi [1995] and have since led

27 to the development of robust software development environments Dmitriev [2004], Erdweg et al. [2013]. In
28 this study, we discuss an application of a LW platform to facilitate the teaching of data analysis to biologists.
29 To this end, we used the Meta-Programming System (MPS, <http://jetbrains.com/mps>), a robust
30 and open-source LW.

31 One question we were particularly interested in answering was whether we could create an analysis
32 tool that blends the boundary between programming/scripting languages and graphical user interfaces, and
33 therefore facilitate teaching complete beginners. Programming languages such as the R language Ihaka
34 and Gentleman [1996] are frequently preferred for data analysis by experts. They have so far been
35 the most flexible and powerful tools for data analysis, but require a steep learning curve. In contrast,
36 beginners tend to prefer data analysis software with a graphical user interface, which are easier to learn,
37 but eventually are found to lack flexibility and extensibility. We reasoned that blending these two types of
38 interfaces into one tool could provide a simpler way for beginners to learn elements of scripting, improve
39 repeatability and reproducibility of their analyses, and increase their productivity.

40 We found that LWT made it possible to quickly develop such a tool. We called this tool MetaR
41 because it leverages the R ecosystem, but supports meta-programming. We designed this novel analysis
42 tool using an iterative process that benefited from frequent feedback from users of the tool. In this
43 manuscript, we describe the goals of the MetaR languages, explain how the tool can be used, and highlight
44 the most innovative aspects of the languages compared to other tools used for data analysis, such as the R
45 language Ihaka and Gentleman [1996] or electronic notebooks.

46 The initial focus of MetaR was on analysis of RNA-Seq data and the creation of heatmaps, but the
47 tool is general and can be readily extended to support a broad range of data analyses. For instance, we
48 have used MetaR to analyze data in a study of association between the allogeneomics score and kidney
49 graft function Mesnard et al. [2015]. We chose to focus on the construction of heatmaps as a use case and
50 illustration for this study because this activity is of interest to many biologists who obtain high-throughput
51 data.

52 We report on our experience teaching MetaR to complete beginners and compare the duration of such
53 training sessions to that of similar training conducted with traditional approaches and tools. Importantly,
54 we found that both beginners and experts can benefit from blending user interfaces and scripting. Beginners
55 benefit because the MetaR user interface is much simpler to learn than the full R programming language.
56 Expert users benefit because they can quickly prototype and develop high-level language elements to
57 simplify repetitive aspects of data analysis.

58 RESULTS

59 Design of a High-level Data Analysis Language

60 Several decisions must be made when designing a new computational language. Most decisions are driven
61 by design goals. We have designed the MetaR language to address the following goals:

- 62 1. The language should help users who have no knowledge of programming. The goal is to offer a
63 smooth learning curve for beginners used to GUIs. We favor declarative language constructs over
64 flexibility in parts of the language aimed at beginners.
- 65 2. Since a table of data is a frequent input when working with high-throughput data, make Table a first
66 class element of the design. Leverage this element to simplify the annotation of the columns of a
67 table. We rely on the idea that a little formalism (e.g., annotation of table columns) goes a long way
68 to simplify analysis scripts.
- 69 3. Eliminate the need to know the language syntax to help beginners get started quickly. We leverage
70 the MPS LW and its projectional editor to this end (Voelter and Solomatov [2010]). The MPS
71 projectional editor provides interactive features, such as auto-completion, that provide guidance to
72 beginners and experts alike when using the language to develop analyses.
- 73 4. Provide the ability to blend a scripting language with a graphical user interface. We use language
74 composition and the ability of the MPS LW to render nodes with a mix of text and graphical user
75 interface components.
- 76 5. Offer essential data transformations (e.g., joining two tables, taking subsets of rows of a table) via
77 simple, yet composable language elements.

- 78 6. Provide means for experts to use their knowledge of the R language to work-around cases when
79 the MetaR language is not sufficiently expressive to perform a specific analysis. We offer the
80 ability to embed R code inside a MetaR analysis, as well as the ability to write scripts in the R
81 language. In both instances, this variant of the R language supports language composition and
82 enables embedding graphical user interfaces inside script fragments.

83 Overview of MetaR and Composable R

84 Figure 1 presents an overview of the features offered by MetaR and Composable R, for the full range of users that the platform supports, from beginner to expert.

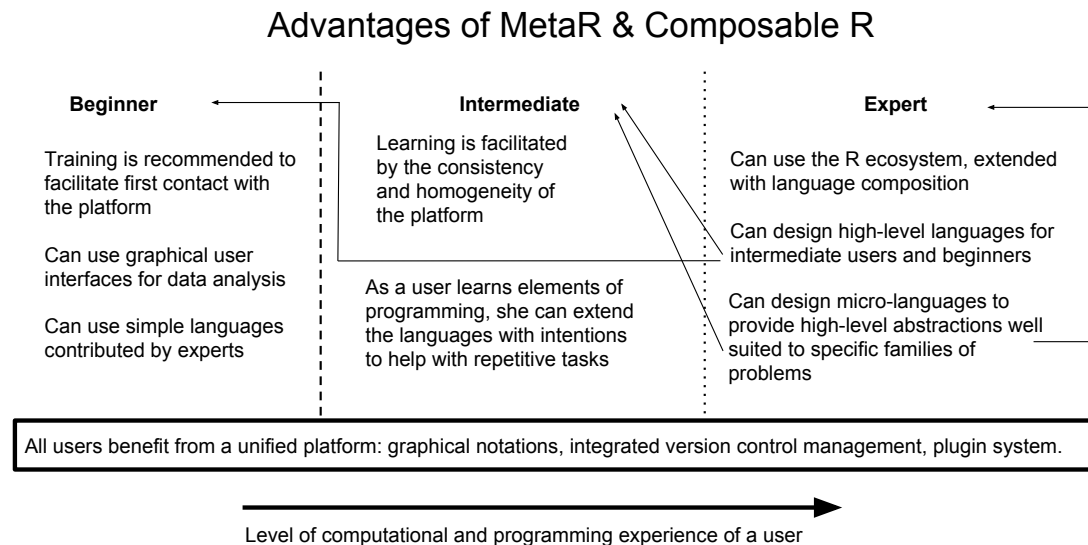


Figure 1. Overview of features provided by MetaR and Composable R in the LWT platform. We present the capabilities of the MetaR platform organized by the level of experience of a user. Beginners mostly benefit from the ability to blend graphical interfaces with scripting, and from high-level languages developed by experts on the platform. Intermediate users, who have basic programming skills, are able to customize the languages in simple ways, such as by creating intentions to help with repetitive steps of analyses. Intentions are context-dependent actions that can be added to a language at runtime (see Simi and Campagne [2014] for illustrations). Experts are users with strong programming skills who have become familiar with LWT. They can create micro-languages to extend Composable R, or design entirely new data analysis languages to help beginners with analysis for new domains. Users at all levels benefit from LWT platform features, including seamless integration of the languages with version control (see Benson and Campagne [2015] for a discussion of the integration with version control).

85

86 Teaching the MetaR data analysis language

87 Teaching a data analysis tool can smooth the learning curve and prevent un-necessary frustration that
88 students could experience if they tried working with the tool on their own. Since MetaR is a new platform
89 for data analysis, training is also important to help users get started with the software.

90 For this reason, we developed detailed training material for the MetaR language and have offered
91 training sessions at the Weill Cornell Medical College since January 2015. These monthly training
92 sessions were offered to technicians, students, post-doctoral fellows and faculty across the institutions of
93 our Clinical Translational Science Center (including three research institutions and one undergraduate
94 city college), but included participants from other institutions in NYC.

95 When advertising the training sessions, we explicitly indicated that participants required minimal
96 prior computer experience (“Analyses are executed in the R language, but no knowledge of R is needed
97 to use MetaR”, “No programming or UNIX skills are required.”). This drew a large participation from
98 attendants who had never used the R language or the command line.

99 Despite their limited computational experience, most participants were able to follow the 2h training
100 session¹ and construct a heatmap on their laptop by the end of the training. Participants who could not
101 complete the assignment were either unable to install the software on an outdated laptop operating system
102 (for instance, some problems included outdated security that prevented connection to the Wifi network
103 at our institution), or encountered installation problems for R packages that would not install on their
104 laptops on the day of training. As we recognized these problems, we developed simpler ways to install
105 software dependencies on user laptops and encouraged users to download the required software before
106 the session. We have so far trained approximately 150 participants in using MetaR to call RNA-Seq
107 differential expression and create a heatmap.

108 The duration of MetaR training (2h) compares favorably with training sessions offered with R and
109 bioConductor packages. An R/bioConductor workshop offered at the Weill Cornell Medical College
110 requires about 6 hours of training (three two hour sessions) in order to help beginners conduct similar
111 analyses to those performed in the MetaR training session. We requested help from the community to
112 try and to quantify the amount of time typical R/bioConductor training sessions require. We created
113 an online survey that trainers and trainees could fill out anonymously (<http://goo.gl/forms/3ZWESUgtmd>). The response rate of this survey was low, but indicated that between 6 and 24 hours was
114 considered by most teachers a typical amount of time needed to teach how to call differentially expressed
115 genes and creating a heatmap with R/bioConductor (two answers listed 6 hours, one answer 8 hours, two
116 listed 24 hours). Interestingly, one trainee who answered the survey noted that 40 hours were required
117 to learn the same skills, suggesting that trainers may under-estimate the amount of time needed when
118 complete beginners try to learn these skills with R/bioConductor. The responses to this community survey
119 indicate that traditional approaches require several 2 hr training sessions for a total of 6-24 hr training.

120 Assuming the responses to this survey are representative, MetaR training sessions are 3 to 6 times
121 shorter than traditional training sessions. These data strongly support the notion that simple languages
122 like MetaR can facilitate the teaching of data analysis for specialized analysis tasks. In the remainder of the
123 Results Section, we explain the design of MetaR in more detail and present the features of the platform
124 useful to experts.
125

126 High-level Design Choices

127 In addition to the design goals presented previously, the design of MetaR included several strategic
128 choices. We now present these choices and their rationales:

129 **Choice of a Target Language and Runtime System** A language needs a runtime system to execute the
130 code of programs written in the language. A possible choice for a runtime is to target another high-level
131 language (such as Java, or C) but this would require implementing all aspects of data manipulation in the
132 target language. Since the R language (Ihaka and Gentleman [1996]) is widely used for data analysis in
133 biology, we considered using it as a runtime system. Experts biostatisticians and bioinformaticians have
134 developed many R packages that implement advanced analysis for biological high-throughput data. These
135 packages can be used to simplify the implementation of a runtime system for a new data analysis language.
136 We therefore decided that the MetaR language would generate R code in order to take advantage of the
137 packages developed in this language. This decision greatly simplified the implementation of the MetaR
138 language because it removed the need to develop a custom language runtime system.

139 **Data Object Surrogates** MetaR makes extensive use of Data Object Surrogates (DOS, our terminology).
140 A data object surrogate is an object that represents other data (the source data). The surrogate often
141 contains only limited information from the original data source. The DOS contains just enough to facilitate
142 referring to the source data in another context for the purpose of data analysis, but not as much as to
143 represent the entire content of the data source in memory. A good example of DOS is the Table object,
144 which stores information about the columns of a data file. The Table DOS describes the columns of the
145 table, but does not store the data contained in the table. A DOS typically has a name which can be used
146 to refer to the DOS and its source data inside a MetaR model. References to table DOS help users refer
147 to the table as they develop an analysis. Our use of the MPS LW facilitates the creation of DOS. In the
148 MPS LW, we model DOS as concepts of the language. For instance, the Table DOS is represented by a

¹Sessions are scheduled for 2hrs, but often complete half an hour early when participants do not require software installation troubleshooting at the start of the session.

149 Table concept, whose instances can be created in a model as root nodes. DOS are also used in MetaR to
150 represent plots.

151 **Immutable Data Objects** Many programming languages (of which C, C++, Java, Perl, Python and R
152 are members) make it possible to define variables or objects whose values can be changed (so called
153 mutable variables). While this provides flexibility, it is a frequent source of confusion for beginners until
154 they have developed their own mental model of how program steps modify variable values. During the
155 design of MetaR, we chose to offer immutable objects rather than mutable variables when possible. This
156 makes MetaR analyses easier to reason about because the value of objects cannot be changed after the
157 object is created. This design decision does not prevent adding mutable variables to the MetaR language,
158 but simplifies initial learning of the language by complete beginners.

159 **Organization into Languages**

160 We designed MetaR as a collection of MPS languages. The main language, *org.campagnelab.metar* is
161 aimed at beginners with limited computational experience.

162 In the next section, we explain how the MetaR language can be used from the point of view of an
163 end-user. This section also includes highlights of features that differ from the state of the art in data
164 analysis. Please note that exhaustive reference documentation is available elsewhere (see Campagne
165 and Simi [2015]) and the goal of the following paragraphs is to provide a sufficient introduction to data
166 analysis with MetaR that readers can understand the impact of the innovations we tested in developing
167 this tool.

168 **The MetaR Language**

169 **Tables**

170 An example of an immutable DOS is the MetaR Table object. In MetaR, objects of type Table represent
171 tabular data with columns and rows of data. An example of a MetaR Table is shown in Figure 2. A
172 MetaR Table is associated to a data file that contains the actual data of the table in a Tab-Separated Value
173 (TSV). The location of the data file can be specified using Variables (i.e., $\{\text{project}\}$), which offer
174 independence from the local file system structure, and are particularly useful when keeping analyses
175 under source control). A table has columns. Columns have names and types, which determine how data
176 in each column is used. Types of data include string, numeric, boolean and enumeration (a small number
177 of pre-defined categories, such as Male and Female). Figure 2 presents a table of RNA-Seq read counts
178 which was obtained from the Gene Expression Omnibus Seguin-Estevez et al. [2014] and annotated to
179 enable analysis with MetaR.

180 Annotating a Table consists of two steps: (1) browsing to the file that contains the data. This can be
181 accomplished by clicking on the file dialog button (the little square with . . .) to locate the file. Upon
182 selection of a valid file, the MetaR table node inspects the file and determines column names and types.
183 Names and types are then shown in the Table node (under the Columns heading). (2) Specific columns
184 can be annotated with one or more Column Groups.

185 Users can define arbitrary Column Groups in a different node called “Column Groups and Usages”
186 (shown on the right of Figure 2). If two columns are related, user can define a Group Usage to explicitly
187 document the relation. For instance, in Figure 2, the usage *LPS_Treatment* is defined to indicate that
188 the Column Groups *LPS=no* and *LPS=yes* are two kinds of LPS treatments.

189 Tables and their annotations help users formalize information about data in a table. We find that
190 asking the user to provide such information early on is beneficial because the structure of annotations can
191 be leveraged in other parts of the language to provide intelligent auto-completion, customized for each
192 table of data (for instance, to provide auto-completion for column names when writing expressions, or to
193 select columns to use when joining two tables, examples of intelligent auto-completion is provided in the
194 following sections, see Figure 3).

195 For instance, in the dataset of Seguin-Estevez et al. [2014], users can indicate which columns contain
196 data for samples that were treated (*LPS=yes*) with lipopolysaccharide (LPS) or not (*LPS=no*). MetaR
197 facilitates the data curation steps of a data analysis project by offering an interactive user interface to
198 help users keep track of annotations. The interface is interactive in several ways: group names can be
199 auto-completed to the groups defined in the “Column Group and Usage” object. Menus are available to
200 add column group annotations to a set of columns that the user has selected. In addition to LPS treatment,
201 Figure 2 shows the *count* annotation, used in an RNA-Seq differential expression analysis to identify

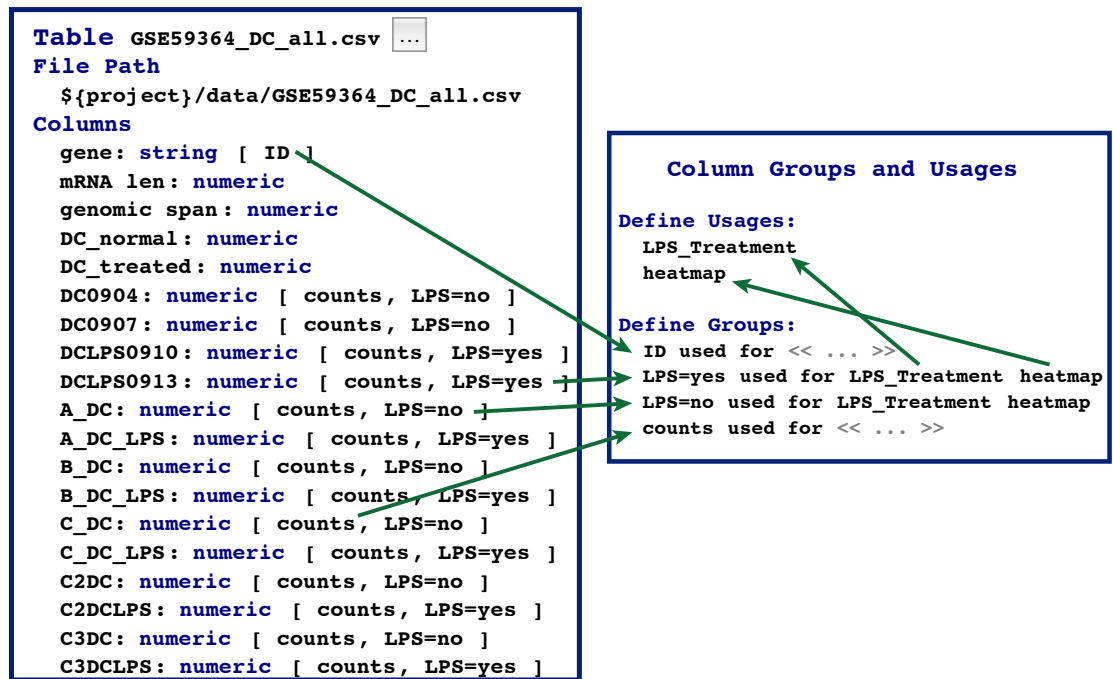


Figure 2. Table and Column Group objects. This figure presents the Table and Column Group objects. Green arrows show some cross-references among nodes of Tables and Column Groups. For instance, the ID group used to annotate the gene column is a reference to the ID group defined under the Column Group and Usage Container.

202 which columns contain read counts, the ID column group, which uniquely identifies specific rows of the
203 data table and the heatmap column group, used to choose which columns groups should be heatmap.
204 This illustrates that the table annotation mechanism is flexible and can be leveraged by specific statements
205 of the language, in order to indicate that the statement needs data annotated in a certain way.

206 **Analyses**

207 Analyses make it possible for users to express how data is to be analyzed. Figure 3 presents a MetaR
208 Analysis node. This analysis is the one we use as a worked example during training sessions we offer at
209 our institution. The editor of an analysis node offers an interface similar to that of a script in a traditional
210 editor, but provides a more interactive and intelligent user interface. For instance, auto-completion
211 is available at every point inside an analysis and suggests possible elements of the language that are
212 compatible with the context at the cursor position.

213 The user may accept a suggestion and this results in the insertion of the language element at the
214 position of the cursor. When the context calls for referencing a column of a table, for instance, only
215 columns of Tables available at this point of the analysis are shown. While it is still possible to make
216 mistakes when using this interface, mistakes created as a result of typos are less common than in programs
217 encoded as text, for two reasons:

- 218 • Auto-completion offers a convenient way to set references between objects. Accepting an auto-
219 completion suggestion helps users avoid typos.
- 220 • Some users choose not to use auto-completion to set references and instead type a referenced node
221 name. In this case, mis-typed names that cannot be resolved to a valid node are highlighted in red
222 and in the right margin of the editor (this feature of the MPS LW is available for all languages
223 developed with the MPS platform). This highlighting draws the attention of the user to the error or
224 typo. This feature is also important when merging different versions of an analysis placed under
225 source control or when combining analyses from parts of other analyses (e.g., errors will be clearly
226 marked after a code fragment is pasted into a new analysis).

Analysis Limma analysis

```

{
  import table GSE59364_DC_all.csv

  subset rows GSE59364_DC_all.csv when true: $(gene) != "Total" -> filtered
  limma voom counts= filtered model: ~ 0 + LPS
  comparing LPS=YES - LPS=NO -> Results

  join ( filtered, Results ) by group ID -> MergedResults
  subset rows MergedResults when true: $(adj.P.Val) < 0.0001 & $(logFC) > 2 | $(logFC) < -2 -> 1% FDR
  heatmap with 1% FDR select data by one or more group LPS=YES, group LPS=NO -> plot HeatmapStyle [
  annotate with 1% FDR ^myOwnTable (manuscript.Limma analysis)
  scale values ^GSE59364_DC_all.csv ^Table (manuscript)
  cluster columns ^MergedResults ^myOwnTable (manuscript.Limma analysis)
  ^Results ^myOwnTable (manuscript.Limma analysis)
  ^filtered ^myOwnTable (manuscript.Limma analysis)
]

multiplot -> PreviewHeatmap [ 1 cols x 1 rows ] 

[ plot ]

  render plot as PDF named "heatmap.pdf"  72dpi
  write Results to "results.tsv" 
}

```

Figure 3. MetaR Analysis. The Analysis node is composed of a list of statements. This analysis works with the table of data presented in Figure 2, removes the row of data where the value *Total* appears in the gene column, performs statistical modeling with Limma Voom to identify genes differentially expressed between LPS treated and control samples, constructs a heatmap and displays the plot as a preview. Finally, the analysis converts the plot to PDF format and writes the joined table (statistics and counts) in the *results.tsv* file.

227 Auto-completion help is available for the various types of references supported by the MetaR language.
 228 Examples of these can be seen on Figure 3 for tables (whose names are in green), plots (whose names are
 229 in blue), styles (names shown with a green background and white foreground, such as *HeatmapStyle*), or
 230 Column Group names (shown with a blue grey background and black foreground). Pressing control-B
 231 (or command-B on Mac) with the cursor on these nodes navigates to the destination of the reference (a
 232 menu is also available to help novice users discover this navigation mechanism). References may point to
 233 children nodes defined inside an analysis (e.g., plots), or nodes defined outside the analysis (e.g., tables
 234 and column groups).

235 Importantly, the MetaR user interface can also display buttons and images directly as part of the
 236 language. This feature takes advantage of the ability of the MPS LW to embed arbitrary graphical elements
 237 in the projectional editor. This capability is illustrated in Figure 3 by the “Hide Preview” button and by
 238 the heatmap image shown immediately below the *multiplot* keyword (pressing this button hides the
 239 plot preview).

240 The level of interactivity provided by the MetaR user interface is best conveyed by watching video
 241 recordings of its use. We provide training videos at <http://metaR.campagnelab.org> to illustrate
 242 how much more interactive the MetaR language is compared to other languages commonly used for data
 243 analysis.

244 Language Composition and Micro-Languages

245 Since MetaR is implemented as a set of MPS languages, it fully supports language composition (Voelter
 246 and Solomatov [2010]). Language composition has no equivalent in text-based programming languages

247 and many readers may be therefore unfamiliar with this technique. We will use an example to explain the
248 advantage of this technique for data analysis.

249 Consider the table of results produced by the analysis shown in Figure 3. Users are likely to need
250 to annotate the subset of genes found differentially expressed with gene names and gene descriptions.
251 Information such as this is available in the Biomart system Haider et al. [2009].

252 To illustrate language composition, we created a new kind of MetaR statement called `query`
253 `biomart`, which we defined in a micro-language. A micro-language is a language which provides
254 only a few concepts meant to extend a host language. In this case, the MetaR language is the host
255 language and `query biomart` is a concept contributed by the the micro-language. The purpose of this
256 concept is to connect to Biomart and retrieve data. In the R language, this functionality is provided as a
257 BioConductor package (called “`biomaRt`”, Durinck et al. [2005])

Analysis Micro Language Example

```
{  
  import table results.tsv  
  query biomart database ENSEMBL GENES 81 (SANGER UK) and dataset Homo sapiens genes (GRCh38.p3)  
    get attributes HGNC symbol from feature of types string with column group annotation select a group  
      Description from feature of types string with column group annotation select a group  
      Ensembl Gene ID from feature of types string with column group annotation ID  
  filters HGNC symbol(s) [e.g. NTN3] from results.tsv when true: ${adj.P.Val} < 0.01  
  -> resultFromBioMart  
  join ( resultFromBioMart , results.tsv ) by group ID -> Annotated Results  
}
```

Figure 4. Example of Micro-language Composition. The `query biomart` statement is defined in a
micro-language called `org.campagnelab.metar.biomart`, which extends the host language
`org.campagnelab.metar.tables`. The biomart language provides one statement that offers an interactive
user interface to help users retrieve data from biomart. This language reuses expressions and tables from
the host language. Micro-languages can be enabled or disabled dynamically by the end-user at the level
of a model. This example retrieves Human ENSEMBL identifiers and gene descriptions using the HGNC
gene symbols used as identifiers in the Results table (see Figure 3 for the analysis that produced Results).

258 Querying Biomart in R consists in calling one of the functions defined in the package with specific
259 parameters. The statement is very specialized, and for this reason would not typically be part of the core
260 statements of a text-based programming language. Leveraging language composition, we can offer a
261 dedicated statement that supports auto-completion in a remote Biomart instance. The statement acts as a
262 specialized user interface designed to help users retrieve data from Biomart (in very much the same way
263 that the web-based interface to Biomart helps users query this resource, but here completely integrated
264 with the MetaR host language).

265 Figure 4 illustrates how the `query biomart` statement can be used to obtain gene annotations. In or-
266 der to use these statements, end-users of MetaR would declare using both the `org.campagnelab.metar.tables`
267 (the host language) and `org.campagnelab.metar.biomart` (the micro-language). In this specific case, the
268 micro-language is provided with the MetaR distribution, but end-users can also implement other micro-
269 languages to seamlessly combine them with the host language (the process for doing so is described in
270 the MetaR documentation booklet Campagne and Simi [2015], Chapter 10). This capability makes it
271 possible to customize the data analysis process for specific problems in much more flexible ways than
272 would be possible with text-based programming languages: with the `query biomart` statement, we
273 demonstrated that it is possible to remotely query databases to support auto-completion directly in the
274 language. In contrast, text-based languages can only be extended in ways compatible with the syntax of
275 the programming language, and are not able to support such levels of interactivity.

Composable R language

276 In addition to the MetaR language illustrated in Figure 2-4, we have developed a composable R language.
277 This language models the traditional R language Ihaka and Gentleman [1996], but supports language
278 composition. Composable R is implemented in the language `org.campagnelab.metar.R` distributed with
279 MetaR. R programs can be pasted in text form into an RScript root node and the text is parsed and
280 converted to nodes of the composable R language. In Figure 5, we show the R code equivalent to the
281 analysis shown in Figure 4. This R script was pasted from the text generated automatically from the
282 MetaR analysis shown in Figure 4. Executing this script is supported in the MPS LW and yields the same
283 result that of the simpler MetaR script shown in Figure 4.
284


```
R Example.R
libDir <- "/Users/fac2003/.metaRlibs "
dir.create(file.path(libDir), showWarnings = FALSE, recursive = TRUE) .libPaths(c(libDir))
dir.create(file.path("/Users/fac2003/R_RESULTS/manuscript "), showWarnings = FALSE, recursive = TRUE)
if ( ! ( require("biomaRt") ) ){
  if ( ! require("BIOCInstaller") ){
    source("http://bioconductor.org/biocLite.R ",
           local = TRUE)
  }
  biocLite(ask = FALSE, c("biomaRt")) library("biomaRt")
}
if ( ! require("plyr") ) {...} if ( ! require("data.table") ) {...}
results.tsv <- fread("/Users/fac2003/MPSProjects/git/metar/data/manuscript/results.tsv ",
                    colClasses = c("character", "numeric", "numeric", "numeric", "numeric", "numeric", "numeric"))
cat("STATEMENT_EXECUTED/1382062817028347486/\n ")
queryBiomart_1382062817028347636 <- function ( <no parameters> ){
  output <- c()
  thisDataset <- "hsapiens_gene_ensembl "
  thisMart <- useMart("ensembl", dataset =
    thisDataset) attributes <- c("hgnc_symbol", "
    description", "ensembl_gene_id")
  filtersVector = c() valuesList = c()
  filtersVector <- c(filtersVector, "
  hgnc_symbol")
  data <- results.tsv[
    ( results.tsv$ "adj.P.Val" <0.01 )
  ]
  valuesList <- c(valuesList, list(tableIds =
    as.vector(data$ genes))) output <- getBM(
    attributes = attributes, mart = thisMart,
    filters = filtersVector, values = valuesList)
  colnames(output) <- c("
  HGNC_symbol_from_feature ", "
  Description_from_feature ", "
  Ensembl_Gene_ID_from_feature ") return(
  data.table(output, key = colnames(output)))
}

queryBiomart_1382062817028347636 ( ) -> resultFromBioMart
write.table(resultFromBioMart, "/Users/fac2003/R_RESULTS/manuscript/table_resultFromBioMart_0.tsv ",
            row.names = FALSE, sep = "\t")
cat("STATEMENT_EXECUTED/1382062817028347636/\n ")
setkey(resultFromBioMart, "Ensembl_Gene_ID_from_feature") setkey(results.tsv, "genes")
results.tsv <- rename(results.tsv, c(genes = "Ensembl_Gene_ID_from_feature "))
tableSuffixes = c(" ", " ")
joiningColumns = c("Ensembl_Gene_ID_from_feature ")
nextTableToMergeInto = resultFromBioMart nextTableToMergeFrom = results.tsv
mergedresults.tsv <- merge(nextTableToMergeInto, nextTableToMergeFrom, by = joiningColumns,
                          suffixes = tableSuffixes) nextTableToMergeInto = mergedresults.tsv
Annotated_Results <- mergedresults.tsv rm(mergedresults.tsv)
Annotated_Results <- Annotated_Results[ , "genes" := Annotated_Results$ "Ensembl_Gene_ID_from_feature " ]
results.tsv <- rename(results.tsv, c(Ensembl_Gene_ID_from_feature = "genes"))
write.table(Annotated_Results, "/Users/fac2003/R_RESULTS/manuscript/table_Annotated_Results_0.tsv ",
            row.names = FALSE, sep = "\t") cat("STATEMENT_EXECUTED/1382062817033011970/ ")
```

Figure 5. R language equivalent of the Analysis shown in Figure 4. To produce this figure, the analysis shown in Figure 4 was generated to the R language and the text was pasted in a RScript node of the composable R language. Automatic parsing of the R code into composable R objects yields a composable R version of the biomart example. Notice that boiler plate code needed to import R packages is shown only for the biomaRt package. Subsequent package import statements have been folded { . . . } to save space in the Figure. Folding is directly supported by the MPS LW. Function calls are highlighted in green and are linked to the function declaration in the package stub (end-user can navigate to each function to review its list of arguments, for instance). While it is likely that expert R programmers could produce somewhat more compact R code than this automatically generated code, comparison with Figure 4 indicates that a micro-languages can offer a concise alternative to a series of function call. The figure also illustrates the breadth of support for the language implemented in Composable R.

285 **Micro-Languages for the R Language**

286 A composable R language makes it possible to create micro-languages that compose directly with R as
287 the host language. We demonstrate this capability by adapting the `query biomart` statement shown in
288 Figure 4 to the R language. Adaptation is simple because both MetaR and R generate to the same target
289 language (R). In this case, we create a sub-concept of `Expr` (this type represents any R expression), and
290 define a field of type `Biomart` (the concept that implements `query biomart`). This simple adapter is
291 sufficient to make it possible to use the `query biomart` user interface inside an R script and is defined
292 in the language *org.campagnelab.metar.biomartToR*. The result of composing the adapter language with
293 composable R is shown in Figure 6. We also provide a short video to illustrate the interactive capabilities
294 of a micro-language combined with composable R (see <https://youtu.be/ZwGj1RPOODQ>).

295 This example illustrates that a composable R language makes it possible to mix regular R code with
296 new types of language constructs that can include user interfaces elements. This opens up new possibilities
297 to facilitate repetitive analyses in R, for instance for specific data science domains (e.g., the Biomart
298 example is useful for bioinformatic data analyses), but also for more general activities where simpler ways
299 to perform a task would be beneficial. An example of this would be a micro-language to facilitate the use
300 of packages to replace the boiler-plate package import code found at the beginning of most R scripts.

```
QueryBiomartInR.R
  if ( ! require("data.table") ) {
    install.packages("data.table", repos = "http://cran.us.r-project.org ")
    library("data.table")
  }
  if ( ! require("biomaRt") ) {...}
  if ( ! require("graphics") ) {...}

  query biomart database ENSEMBL FUNGI 29 (EBI UK) and dataset Aspergillus terreus genes (Broad (CADRE))
  get attributes % identity from aflavus homologs of types string with column group annotation select a group
  filters << ... >>
  -> resultFromBioMart
  [BioMart]
  pdf("histogram.pdf")
  hist(resultFromBioMart$percent_identity_from_aflavus_homologs )
  dev.off()
```

Figure 6. Composing Query Biomart with the composable R language. We developed an adapter that makes it possible to use the MetaR `query biomart` statement directly inside a composable R Script. This figure shows how the `query biomart` Expression adapter appears when used inside an R script. Notice how the table and column adapters are used inside a regular `hist()` function call `resultFromBioMart$percent_identity_from_aflavus_homologs`. These adapters make it possible to refer to the table produced by the statement as an R expression and provide auto-completion for column names in the table (determined dynamically based on the query expressed in the `query biomart` statement).

301 **Using R Expressions in the MetaR Language**

302 Figure 7 illustrates that language composition can also be used to embed R expressions inside a MetaR
303 analysis. This extension is possible because both analyses and R expressions generate code compatible
304 with the syntax of the R programming language. Providing a way to embed the full language in a simpler
305 analysis language offers a guarantee that the end-user will not be overly limited by restrictions of the
306 simpler language.

307 **SOFTWARE**

308 MetaR is distributed as a plugin of the MPS LW. Instructions for installing the software are available
309 online at <http://metar.campagnelab.org>. Briefly, after installing MPS, users can download
310 and activate plugins with the Preferences/Plugins (Mac) or Settings/Plugins (Windows/Linux) menu.
311 Plugins are stored as Zip files on the JetBrains Plugin repository https://plugins.jetbrains.com/category/index?pr=mps&category_id=92 and can also be downloaded and installed
312 manually from the zip file. Source code (technically, MPS languages serialized to files) are distributed
313 on GitHub at <https://github.com/campagnelaboratory/MetaR> Campagne et al. [2015].
314 MetaR (and the MPS LW) are distributed under the open-source Apache 2.0 license.
315

The screenshot shows a software interface with two main panels. The top panel, titled "Analysis MetaR with R", contains R code for simulating a dataset and performing a limma voom analysis. The code is as follows:

```
{
  simulate dataset with [
    num of samples: 100
    num of genes: 500
    mean when all factors are false: 10
    discrete factors: treatment
    effect size: 3
    continuous covariate: age , range: [ 18 - 100 ] , slope: 3
  ] -> simulatedTable

  limma voom counts= simulatedTable model: ~ 0 + treatment + age
  comparing treatment=No - treatment=Yes -> Results adjusted counts: Adjusted
  // show the top 10 hits using an R expression (limma voom binds the fit3 name to a value)
  — R
  topTable(fit3, coef = 1, number = 10)
  R —
}
```

The bottom panel, titled "Run R Script MetaR with R", shows the execution output. It includes the R script path, package loading messages for limma, methods, edgeR, Cairo, and data.table, and three hyperlinks for STATEMENT_EXECUTED. Below this is a table of results:

	genes	logFC	AveExpr	t	P.Value	adj.P.Val
323	gene_323_treatment	-0.4182245	11.03107	-5.674651	2.001635e-08	1.000817e-05
54	gene_54_treatment	-0.3870329	11.13942	-5.277938	1.722625e-07	4.020041e-05
126	gene_126_treatment	-0.3866908	11.09819	-5.213549	2.412025e-07	4.020041e-05
112	gene_112_treatment	-0.3808808	11.11778	-5.118621	3.936154e-07	4.832646e-05
166	gene_166_treatment	-0.3792082	11.12251	-5.078391	4.832646e-07	4.832646e-05
232	gene_232_treatment	-0.3729548	11.14150	-5.026791	6.274733e-07	5.228944e-05
445	gene_445_treatment	-0.3626018	11.10790	-4.859438	1.440418e-06	1.028870e-04
360	gene_360_treatment	-0.3594876	11.08581	-4.793461	1.985425e-06	1.240891e-04
328	gene_328_treatment	-0.3437431	11.10246	-4.551622	6.230697e-06	3.254523e-04
351	gene_351_treatment	-0.3307637	11.12417	-4.542160	6.509046e-06	3.254523e-04

Figure 7. Composing R Expressions with the MetaR Language. Top panel: this example illustrates that it is possible to use R code inside a MetaR analysis. In this snapshot, R code is delimited by the — R and R — markers and shown with a blue background. Embedding R code in MetaR provides flexibility to perform operations for which MetaR statements have not yet been developed. The analysis shown simulates a dataset using simple parameters and tests the ability of Limma voom, as integrated with MetaR, to call differentially expressed genes. Bottom panel: shows the result of executing the analysis inside the MPS LW. As part of execution, the analysis is converted to R code, this code is run and standard output displayed inside the LW. The STATEMENT_EXECUTED// lines hyperlink the progress of the execution with each specific analysis statement that has been executed.

316 DISCUSSION

317 Data Object Surrogates and Relation to Meta Data

318 DOS are related, but different from metadata. For instance, the Table DOS provides metadata about the
319 file that contains the tabular data represented by Table nodes. It lists columns, associates columns to
320 groups and defines group usages. This type of information can be thought of as metadata about the file
321 that contains the tabular data. However, there is an important difference between DOS and metadata. For
322 instance, a MetaR Table only provides metadata relevant to the analysis that the user needs to perform. It
323 makes no effort to provide information that would have a meaning outside of the user's analysis. This
324 simplification maximizes the benefit of annotation while keeping the effort needed to produce it minimal
325 and local to the user who actually needs the annotation.

326 Graphical User Interfaces for Data Analysis

327 Programs with graphical user interfaces (GUIs) (also called direct manipulation interfaces Galitz [2007])
328 are often popular among beginners who are starting with data analysis and have no programming or

329 scripting experience. GUIs are popular in part because they facilitate discovery of software functionality
330 directly when using the software. They do not require prior-knowledge of syntax.

331 Data Analysis software with GUIs constrains how analysis is to be performed and provides clear
332 menus and buttons that make it obvious what the program can do. A user can often discover new ways to
333 perform analysis with these tools simply by browsing the user interface and looking at choices offered
334 in menus and dialogs of the program. While such programs are favored by beginners (because they
335 are relatively easy to learn), more advanced users who need to perform similar analyses across several
336 datasets tend to strongly prefer analysis software that does not require repeating interactions with a GUI
337 for every new dataset that must be studied. The novel approaches we have used to develop MetaR share
338 these advantages with GUIs.

339 A minority of analysis software with GUIs also supports writing and running scripts in their user
340 interface. For instance, JMP from SAS Inc. is an example of a statistical analysis software with GUI that
341 also offers a scripting language. However, when scripting is offered, it is often only loosely integrated
342 with the rest of the interface. Furthermore, users who are familiar with the GUI often need to learn
343 scripting from scratch and do not benefit much from their prior experience using the GUI.

344 **Scripting and Programming Languages for Data Analysis**

345 Scripting and programming languages are popular options for data analysis because analyses encoded
346 in scripts or programs can be reused with different datasets. This makes these options popular among
347 researchers who have programming skills and engage frequently in data analysis. The popularity and
348 power of scripting for data analysis is epitomized by the development of the R language Ihaka and
349 Gentleman [1996], which has become a defacto workhorse of open data science in biology. The versatility
350 of the R language is its strength, but mastering the language requires elements of programming. Learning
351 the R programming language is not as simple as learning how to use a GUI analysis tool and many users
352 who would benefit from data analysis experience difficulties with the steep learning curve involved in
353 learning programming and the R language.

354 In contrast to R, the MetaR language offers a much simpler alternative for users who have no prior
355 programming background. At the same time, the Composable R language offers the means for expert R
356 users to extend the R language with micro-languages in order to provide custom user interfaces. Such
357 interfaces could be used to flatten the learning curve for novice data analysts or to empower expert data
358 analysts with expressive means to encode solution to specific problems. Since both these options are
359 available in the same platform (the MPS LW), users who become skilled with one language acquire
360 transferable skills that help them learn other languages available on the platform.

361 **Impact on Development of User Proficiency**

362 The MetaR high-level language shown in Figures 2 to Figure4 is aimed at novice data analysts. An
363 interesting question is whether such a language can help novice data analysts learn skills that are useful
364 when working with a variety of data analysis tasks.

365 If the language is sufficiently general, then novice users may learn skills that they can reuse when
366 learning other general data analysis languages. If the language is too limited, then novice users would
367 only learn a specialized analysis tool similar to existing GUI analysis tools. Rigorously determining to
368 which category the MetaR language belongs would require following users for several months or years
369 while they use the tool and we have not done such a study. However, we think that MetaR can help users
370 transition to more general languages for the following reasons.

371 First, users who learn the high-level MetaR language acquire basic skills that are similar to those
372 needed when working with other languages, including composable R. For instance, users learn to formalize
373 their analysis intent using the constructs offered by the language. This is a very important first step that
374 users with a strong programming background may take for granted, but that is difficult for novice users to
375 acquire when they are distracted with problems of syntax. MetaR avoids syntax distractions and helps
376 novice users focus on the logic of an analysis (e.g., how to combine language elements to achieve the
377 desired analysis).

378 Second, the high-level MetaR language does not offer loops and conditionals. Since these language
379 features are often needed for advanced analysis, many users who reach the point where they will need
380 these language features will need to learn a language like R. MetaR offers composable R for this purpose.
381 Novice users who have first learned the MetaR high-level language will be familiar with the MPS LW
382 platform where composable R is also available. Some skills that users have acquired working with the

383 high-level language will be directly transferable, including: how to run a script, how to navigate references
384 to look at definitions, how to use auto-completion or use intentions to transform the program automatically,
385 how to use source control (seamlessly integrated with the MPS LW). Subsets of the R language will still
386 need to be learned to perform more advanced analysis, but learning can occur in an environment where
387 the user is already comfortable. We believe that such an integrated environment where both high-level and
388 low-level languages of the R ecosystem are offered will facilitate teaching of the many skills needed for
389 data analysis. Formally testing whether this intuition is correct will require comparing cohorts of subjects
390 learning data analysis. Alternatively, the answer may become apparent if a large number of data analysts
391 were to transition to using composable R after initially learning the MetaR high-level language.

392 **Relation to Electronic Notebooks**

393 MetaR shares some similarities to electronic notebooks such as IPython Pérez and Granger [2007], Jupyter
394 (<https://jupyter.org/>) and Beaker (<http://beakernotebook.com/>) notebooks, but also
395 has some important differences.

396 Regarding analogies, both MetaR and notebooks can be used to present analysis results alongside the
397 code necessary to reproduce the results. For instance, the MetaR multi-view plot can be used to show a
398 plot at the location where the statement is introduced in an analysis.

399 MetaR was developed approximately over the course of one year (2015). As such the software cannot
400 be expected to be as feature-rich as software developed for many years. Beside this obvious difference,
401 MetaR has the advantage to support language composition. In contrast, current data analysis notebooks
402 support conventional programming languages constructed using text-based technology. Therefore, the
403 closest that notebooks can approach language composition is to support multiple languages in one
404 notebook, a so-called polyglot feature, available for instance in the Beaker notebook. Polyglot notebooks
405 are useful, but cannot be extended by data analysts to customize languages for the requirements of a
406 specific analysis project or domain. For instance, supporting a simple analysis language like MetaR would
407 not be possible without developing a MetaR compiler and an associated execution kernel for the notebook.
408 Developing and using micro-languages together with the traditional languages supported by the notebooks
409 is also not possible.

410 Hence, the approach taken with MetaR is different from notebooks in two major ways. First, MetaR
411 provides flexibility in designing new languages or micro-languages. It is not constrained by the syntax
412 of a full programming language. Extending MetaR often consists in adding just one statement to an
413 existing language. This promotes collaborative language design and development since many users can
414 acquire sufficient skills to create one or two statements, reusing the building blocks provided by the
415 host language (the steps needed to extend MetaR with a new language statement are described in the
416 user manual Campagne and Simi [2015]). As long as a new statement generates valid R code, a MetaR
417 Analysis that contains this statement will be executable.

418 Second, the syntax of the MetaR languages is not limited to text scripts or programs. Language
419 Workbench technology used to implement MetaR supports graphical notations and diagrams as well
420 as text. These differences combine to make it easier to design and implement custom data analysis
421 abstractions with the LWT approach than it is possible with current electronic notebooks. Interestingly,
422 the R IPython kernel could be used to execute scripts generated from MetaR analyses, which would
423 provide an interactive console similar to that offered in the IPython notebook inside the MPS LW.

424 **Reproducible Research and Education**

425 MetaR analysis and Composable R scripts can be executed seamlessly with an R environment installed
426 inside a docker image (see Methods). Users can enable this feature by providing a few details about the
427 installation of docker on their computer and checking the “Run with Docker” option in the MPS LW. This
428 feature is particularly useful to facilitate reproducible research because docker images can be tagged with
429 version numbers and always result in the same execution environment at the start of an analysis. This
430 makes it possible to pre-install specific versions of R, CRAN and Bioconductor packages in a container
431 and distribute this image with the MetaR analyses or R scripts that implement the analysis inside the
432 container. While this is possible also with R, using docker on the command line, the customization of
433 the MPS LW makes it seamless to run analyses with docker. We are not aware of a similar feature being
434 supported by current R IDEs.

435 We found this feature also particularly useful for training sessions where installation of a working
436 R environment can be challenging on trainees’ laptops. Using docker, we simply request that trainees

437 pre-install Kitematic (available on Mac and Windows), or run docker natively on Linux and download
438 the image we prepared with the packages used in the MetaR training sessions. The ability to run MetaR
439 analysis in docker container results in a predictable installation of dependencies for training session and
440 frees more of the instructor's time to present data analysis techniques.

441 **METHODS**

442 **Language Workbench Technology Primer**

443 Since many readers may not be familiar with LWT, this section briefly describes how this technology
444 differs from traditional text-based technology.

445 Text-based programming languages are implemented with compilers that internally convert the text
446 representation of the source code into an *abstract syntax tree* (AST), a data structure used when analyzing
447 and transforming programming languages into machine code.

448 In the MPS LW, the AST is also a central data structure, but the parsing elements of the compilers are
449 replaced with a graphical user interface (called a *projectional editor*) that enables users to directly edit the
450 data structure. Where text-based languages are restricted to programs written as text, a projectional editor
451 can support both textual and graphical user interfaces (such as images, buttons, tables or diagrams) Voelter
452 and Solomatov [2010]. Projectional editors can also offer distinct views of the same AST, implemented as
453 alternative editors. Projectional editors keep an AST in memory until the user saves the program. Saving
454 an AST to disk is done using serialization (loading is conversely done via deserialization to memory AST
455 data structures).

456 The choice of serialization rather than encoding with text has a profound consequence. Serialization
457 uniquely identifies the concept for each node in an AST. This method makes it possible to combine AST
458 fragments expressed with different languages, when the concept hierarchy of the languages supports
459 composition. We have presented examples of language composition in Simi and Campagne [2014],
460 Benson and Campagne [2015]. In this manuscript, we extensively use language composition to extend the
461 R language and provide the ability to embed user interfaces into R programs.

462 **Abstract Syntax Tree (AST)**

463 An AST is a data structure traditionally used by compilers as a step towards generating machine code.
464 In the MPS Language Workbench, an AST is a tree data structure, where nodes of the tree are instances
465 of concepts (in the object-oriented sense). Figure 8 illustrates the notion of AST nodes, concepts and
466 projectional editor.

467 AST concepts may have properties (values of primitive types), children (lists of other nodes they
468 contain), references (links to other nodes defined elsewhere in the AST). An AST has always a root node,
469 which is used to start traversing the tree. In the MPS LW, AST root nodes are stored in models.

470 **Languages**

471 In the MPS LW, languages are defined as collection of concepts, concept editors (which together implement
472 the user interface for the language), and other language aspects Campagne [2014]. Each language has a
473 name which is used to import, or activate, the language inside a model. After importing a language into a
474 model, it becomes possible to create ASTs with this language in the model. Creating an AST starts with
475 the creation of a root node. Children of the root node are added using the projectional editor. Children of
476 root nodes, properties and references can be edited interactively in the editor.

477 We have used the MPS Language Workbench (<http://jetbrains.com/mps>), as also described
478 in Campagne [2014] and Campagne [2015]. For an introduction to Language Workbench Technology
479 (LWT) in the context of bioinformatics see Simi and Campagne [2014] and Benson and Campagne [2015]
480 in the context of predictive biomarker model development.

481 **Language Design**

482 We designed the MetaR MPS languages through an iterative process, releasing the languages at least
483 weekly to end-users at the beginning of the project and adjusting designs and implementations according
484 to user feedback. Full language developments logs are available on the GitHub code repository (<https://github.com/CampagneLaboratory/MetaR>) Campagne et al. [2015]. Briefly, we designed
485 abstractions to facilitate specific analyses and implemented these abstractions with the structure, editor,
486 constraints and typesystem aspects of MPS languages. Generated R code is produced from nodes of the
487

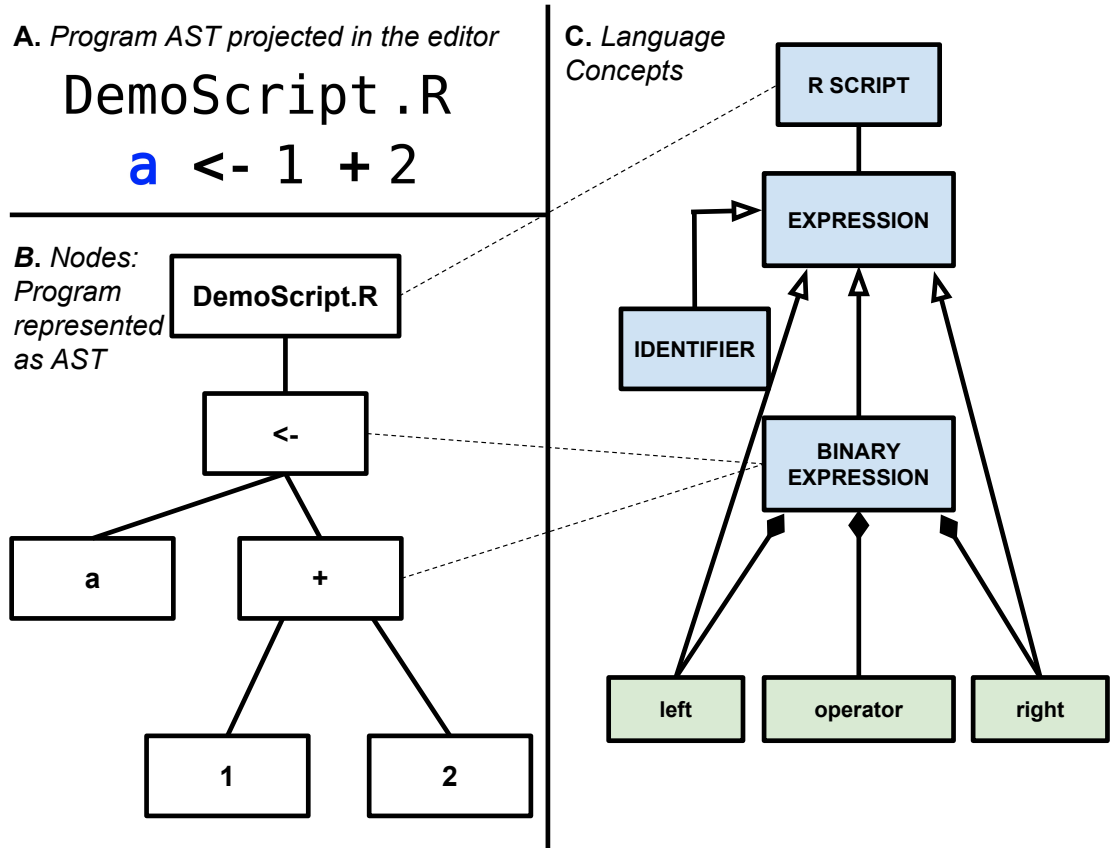


Figure 8. Concepts, Nodes and Projectional Editor. Panel A: Projectional editor showing a simple R script with one assignment expression. Panel B: An abstract syntax tree is shown with nodes that correspond to the program in panel A. Panel C: Language Concepts for the nodes in Panel (B) (shown as blue boxes). Each concept is connected to other concepts with an open-ended arrow to indicate inheritance (e.g., A <- B indicates that B is a sub-concept of A). Green boxes indicate fields of a concept and are connected to the concept that has these fields by a line with a black diamond on the concept that owns the field. This shows that BinaryExpression is a concept that is an Expression and has three fields: left, operator and right. Dotted lines connect nodes to their concept. For instance, the <- and + nodes are instances of BinaryExpression.

488 languages using the *org.campagnelab.TextOutput* plugin. An illustration of the steps required to develop
489 one language statement is available in Chapter 10 of the MetaR documentation booklet (see Campagne
490 and Simi [2015]).

491 **Table Viewer**

492 We implemented a Table viewer as an MPS Tabbed Tool, using the MPS LW mechanisms for user
493 interface extension (see Campagne [2015]). The table viewer provides the ability to inspect the data
494 content of any table produced during an analysis, or any input table. When the cursor is positioned over a
495 node that represent a FutureTable (created when running the R script generated from the MetaR Analysis),
496 and the viewer is opened, it tries to load the data file that the analysis would create for this table. If the file
497 is found, the content is displayed using a Java Swing Component in the MPS user interface of the Table
498 Viewer tool.

499 **Language Execution**

500 MetaR analyses can be executed directly from within the MPS LW. This capability was implemented with
501 Run Configurations (see Campagne [2015], Chapter 5).

502 **Execution in a Docker Container**

503 In order to facilitate reproducible execution, we implemented optional execution within a Docker container.
504 A docker image was created to contain a Linux operating system and a recent distribution of the R language
505 (provided in the rocker-base image), as well as several R packages needed when executing the MetaR
506 statements. The Run Configuration was modified to enable execution inside a docker container when the
507 user selects a checkbox "execute inside docker container". Information necessary to run with docker (i.e.,
508 location of the docker executable, docker server connection settings and image name and tag) is collected
509 under a tab in the MPS Preferences (Other Settings/Docker).

510 **ACKNOWLEDGMENTS**

511 We thank the authors of the *rocker-base* image, used in the MetaR project to facilitate the creation of
512 docker images for training sessions and reproducible research. This investigation was supported by the
513 National Institutes of Health NIAID award 5R01AI107762-02 to Fabien Campagne and by grant UL1
514 RR024996 (National Institutes of Health (NIH)/National Center for Research Resources) of the Clinical
515 and Translation Science Center at Weill Cornell Medical College. Declaration of competing interests:
516 Fabien Campagne is the author of the books "The MPS Language Workbench, Volume I and II" and
517 receives royalties from the sale of these books.

518 **REFERENCES**

- 519 J. Backus. The history of fortran i, ii, and iii. In *History of programming languages I*, pages 25–74. ACM,
520 1978.
- 521 J. W. Backus. Automatic programming: properties and performance of fortran systems i and ii. In
522 *Proceedings of the Symposium on the Mechanisation of Thought Processes*, pages 165–180, 1958.
- 523 V. M. Benson and F. Campagne. Language workbench user interfaces for data analysis. *PeerJ*, 3:e800,
524 2015.
- 525 F. Campagne. *The MPS Language Workbench*, volume I. Fabien Campagne, 2014.
- 526 F. Campagne. *The MPS Language Workbench*, volume II. Fabien Campagne, 2015.
- 527 F. Campagne and M. Simi. *MetaR Documentation Booklet*. Fabien Campagne, 2015.
- 528 F. Campagne, W. Digan, and M. Simi. MetaR: Software release described in the MetaR manuscript, Nov.
529 2015. URL <http://dx.doi.org/10.5281/zenodo.33230>.
- 530 S. Dmitriev. Language oriented programming: The next programming paradigm, 2004. URL <http://www.onboard.jetbrains.com/is1/articles/04/10/lop/>.
- 531 S. Durinck, Y. Moreau, A. Kasprzyk, S. Davis, B. D. Moor, A. Brazma, and W. Huber. BioMart and Bio-
532 conductor: a powerful link between biological databases and microarray data analysis. *Bioinformatics*,
533 21:3439–3440, 2005.
- 534 S. Erdweg, T. van der Storm, M. Völter, M. Boersma, R. Bosman, W. R. Cook, A. Gerritsen, A. Hulshout,
535 S. Kelly, A. Loh, et al. The state of the art in language workbenches. In *Software language engineering*,
536 pages 197–217. Springer, 2013.
- 537

- 538 W. O. Galitz. *The Essential Guide to User Interface Design: An Introduction to GUI Design Principles*
539 *and Techniques*. John Wiley & Sons, 2007. ISBN 0470146222. URL https://books.google.com/books?hl=en&lr=&id=Q3Xp_Auw49sC&pgis=1.
- 541 S. Haider, B. Ballester, D. Smedley, J. Zhang, P. Rice, and A. Kasprzyk. BioMart
542 Central Portal—unified access to biological data. *Nucleic Acids Res*, 2009. URL
543 http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve{\&}db=pubmed{\&}dopt=Abstract{\&}list{_}uids=19420058{\&}query{_}hl=16.
- 545 R. Ihaka and R. Gentleman. R: a language for data analysis and graphics. *Journal of computational and*
546 *graphical statistics*, 5(3):299–314, 1996.
- 547 L. Mesnard, T. Muthukumar, M. Burbach, C. Li, H. Shang, D. Dadhania, J. R. Lee, V. K. Sharma,
548 J. Xiang, C. Suberbielle, M. Carmagnat, N. Ouali, E. Rondeau, J. J. Friedewald, M. M. Abecassis,
549 M. Suthanthiran, and F. Campagne. Exome sequencing and prediction of long-term kidney allograft
550 function. Sept. 2015. ISSN 2167-9843. doi: 10.7287/peerj.preprints.854v2. URL <https://peerj.com/preprints/854>.
- 552 F. Pérez and B. E. Granger. IPython: a system for interactive scientific computing. *Computing in Science*
553 *and Engineering*, 9(3):21–29, May 2007. ISSN 1521-9615. doi: 10.1109/MCSE.2007.53. URL
554 <http://ipython.org>.
- 555 Q. Seguin-Estevez, I. Dunand-Sauthier, S. Lemeille, C. Iseli, M. Ibberson, V. Ioannidis, C. D.
556 Schmid, P. Rousseau, E. Barras, A. Geinoz, I. Xenarios, H. Acha-Orbea, and W. Reith. Ex-
557 tensive remodeling of DC function by rapid maturation-induced transcriptional silencing. *Nu-*
558 *cleic Acids Research*, 42(15):9641–9655, Aug. 2014. ISSN 0305-1048. doi: 10.1093/nar/
559 gku674. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4150779&tool=pmcentrez&rendertype=abstract>.
- 561 M. Simi and F. Campagne. Composable languages for bioinformatics: the nyosh experiment. *PeerJ*, 2:
562 e241, 2014.
- 563 C. Simonyi. The death of computer languages, the birth of intentional programming. Technical re-
564 port, 1995. URL <http://research.microsoft.com/apps/pubs/default.aspx?id=69540>.
- 566 M. Voelter and K. Solomatov. Language modularization and composition with projectional language
567 workbenches illustrated with MPS. *Software Language Engineering, SLE*, 2010.