

UNIVERSITY OF UDINE

DEPARTMENT OF MATHEMATICS, COMPUTER AND PHYSICAL SCIENCE

READERSOURCING 2.0: DOCUMENTATION

MICHAEL SOPRANO AND STEFANO MIZZARO

v1.0.6-alpha

Contents

List of Figures	3
List of Tables	4
1 Introduction	5
2 General Architecture	5
3 RS_Server	5
3.1 Implementation and Technology	6
3.2 Communication Paradigm	8
3.3 Database	9
3.4 Class Diagram	11
3.5 Deploy	13
3.5.1 1: Manual Way	14
3.5.1.1 Requirements	14
3.5.1.2 How To	14
3.5.1.3 Quick Cheatsheet	15
3.5.2 2: Manual Way (But Faster)	15
3.5.2.1 Requirements	15
3.5.2.2 How To	15
3.5.3 3: Heroku Deploy	17
3.5.3.1 Requirements	17
3.5.3.2 How To	18
3.5.3.3 Quick Cheatsheet	18
3.5.4 Environment Variables	18
3.5.4.1 .env file	20
3.5.4.2 Heroku App	20
4 RS_PDF	20
4.1 Implementation and Technology	20
4.2 Package Diagram	21
4.3 Class Diagram	22
4.4 Installation	22
4.4.1 Requirements	22
4.4.2 Commmand Line Interface	22
5 RS_Rate	25
5.1 Implementation and Technology	25
5.2 Installation	25

6	RS_Py	25
6.1	Implementation and Technology	26
6.2	Installation	26
6.3	Usage	26
	References	28

List of Figures

1	Architecture of Readersourcing 2.0 (NOT UML).	6
2	Intuitive scheme of the MVC pattern (NOT UML).	7
3	Entity-Relationship schema of the database of RS_Server (NOT UML).	11
4	Class diagram of RS_Server.	12
5	Representation of the token-based authentication process (NOT UML).	13
6	Package diagram of RS_PDF.	21
7	Class diagram of RS_PDF.	23

List of Tables

1	Subset of the RESTful interface of RS_Server.	10
2	Environment variables of RS_Server.	19
3	Command line options of RS_PDF.	24
4	Parameters available for the Seeder jupyter notebook.	27
5	Parameter available for the Readersourcing jupyter notebook.	27

1 Introduction

The Readersourcing 2.0 ecosystem has been built within a research project co-funded by SISSA Medialab¹ and University of Udine². It has been presented by Soprano et al. [2] during the IRCDL 2019 Conference³. The original paper can be found on Zenodo⁴, while the code and the related documentation is available on GitHub⁵.

Initially, a recap of its general architecture is presented and subsequently a brief description of the role and purpose of each of its components is presented along with some specific aspects such as the technology used, the internal architecture, the structure of the database and more. This is done by using different types of diagrams belonging to the UML standard (unless otherwise specified) which are drawn according to the set of style rules for that standard proposed by Fowler [1].

2 General Architecture

Readersourcing 2.0 is an ecosystem composed of more than one application. Indeed, there must be one (RS_Server [6], presented in Section 3) which acts as a server to gather all the ratings given by readers and one that acts as a client (RS_Rate [5], presented in Section 5) to allow the readers themselves to effectively rate publications, although it is possible to carry out every operation also directly on the web interface provided by RS_Server. There is also a component (RS_PDF [3], presented in Section 4) which has the task to annotate files in PDF format by taking advantage of an ad hoc software library; this component is exploited by the server-side application. Lastly, there is also an additional component (RS_Py, presented in Section 6) which provides a fully working implementation of the RSM and TRM models described by Soprano et al. [2]. An overview of the architecture of Readersourcing 2.0 is shown in Figure 1 and in the following we briefly describe these three components.

3 RS_Server

RS_Server [6] is the server-side application which has the task to collect and aggregate the ratings given by readers and to use RSM and TRM models described by Soprano et al. [2] in order to compute quality scores for readers and publications. An instance of RS_Server is deployed along one of RS_PDF. Then, there are up to n different browsers along with their end-users, which communicate with the server: each of them is characterized by an instance of *RS_Rate*. Both RS_PDF and RS_Rate are described in the following.

This setup means that every interaction between readers and server can be carried out

¹<https://medialab.sissa.it/>

²<https://www.uniud.it/>

³<https://ircdl2019.isti.cnr.it/>

⁴<https://zenodo.org/record/1446468>

⁵<https://github.com/Miccighel/Readersourcing-2.0>

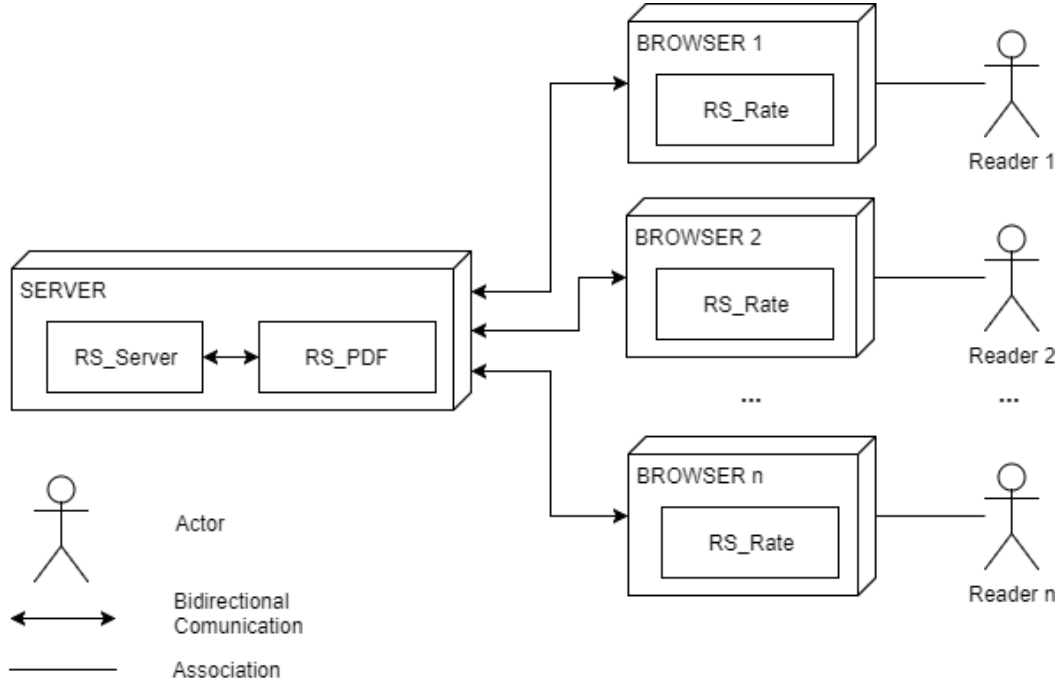


Figure 1: Architecture of Readersourcing 2.0 (NOT UML).

through clients installed on readers' browsers or by using the stand-alone web interface provided by RS_Server and these clients have to handle the registration and authentication of readers, the rating action and the download action of link-annotated publications.

During the design phase of RS_Server some strategies have been adopted to ensure its extensibility and generality. This means that: (i) it is straightforward to add new models, (ii) each model shares the same input data format, and (iii) if a model needs to save values locally to the RS_Server (i.e., in its database), there is a standard procedure to allow that.

3.1 Implementation and Technology

RS_Server is developed in Ruby on Rails,⁶ which is a framework that allows to build applications strongly based on the Model-View-Controller (MVC) architectural pattern. The technology used to develop RS_Server is an open-source web application framework called *Ruby on Rails* (it is also called *RoR* or *Rails* only; more specifically, *Rails* is the framework built above *Ruby*, the actual programming language). It allows to build applications strongly based on the Model-View-Controller (MVC) architectural pattern.

The MVC pattern allows to separate the control logic of the program from data presentation and business logic. Therefore, it allows to obtain an effective architecture since the first moments of its design phase. An intuitive representation of the structure that this pattern

⁶<https://rubyonrails.org/>

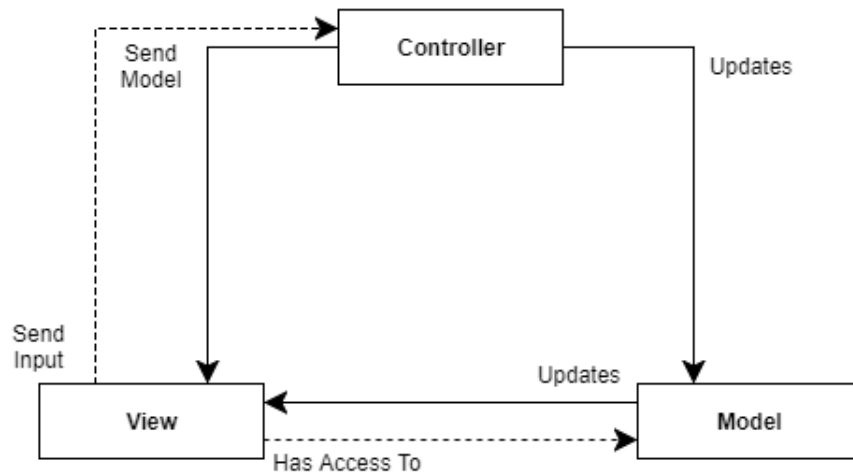


Figure 2: Intuitive scheme of the MVC pattern (NOT UML).

allows to obtain is shown in Figure 2. This structure consists of three distinct entities, called *Controller*, *Model* and *View*. These entities have the task of, respectively, managing control logic, encapsulating business logic and implementing data presentation.

The Controller has direct access to the Model and to the View; from the latter, generally, it receives the user input and on its basis the Controller itself updates the internal state of the Model using its methods. Finally, the Controller sends the updated Model to the View, which is then exploited by the View itself to obtain and display the results of the processing. A generic software can have more than one Controller, where each of them can manage more than one Model instances. In MVC frameworks dedicated to the development of web applications such as Rails, in fact, it is common practice to have a number of Controllers equal to the number of entities modeled within the application domain. Furthermore, there may be more than one View implementation to present the internal state of a specific type of Model.

The use of MVC pattern is not the only founding principle of Rails. One of the most important principles on which Rails itself is based for the development of quality applications is “Convention Over Configuration”. In other words, the framework tries to minimize the decisions that the developer must take during the construction of its application by adopting standard conventions that he can modify if he needs more flexibility. The founding principles of Rails can be deepened by reading the *Rails Doctrine*⁷.

As a last note, Rails is a continuously developing framework and is used industrially by several well-known industry players such as *GitHub*, *SoundCloud*, *Airbnb* and others. It is, therefore, a widespread and appreciated technology, for which there is an active community and a lot of learning material.

⁷<https://rubyonrails.org/doctrine/>

3.2 Communication Paradigm

A modern MVC framework such as Rails allows to develop various kind of web applications. One of the possibilities is to create a *Web Service*, which is a software component capable of carrying out various operations made remotely available through the exchange of messages encoded in a standard interchange format such as *JSON*, all thanks to a transport layer built above the basic Internet protocols like *HTTP*. All this, however, must be carried out according to a paradigm that defines precisely what are the functionalities (resources and operations) actually available and which messages must be received in order to access them.

One of the possible communication paradigms for Web Services is *RESTful* (*REpresentational State Transfer*). Within this paradigm, the functionalities of a Web Service are represented by resources identified by different URIs and the type of HTTP message sent establishes the operation to be performed. The result of the operation initiated by the message received from the Web Service is a new message encoded according to same interchange format of the one which has been sent and it is the client's responsibility to correctly interpret and use the response of the Web Service itself.

RS_Server is a Web Service (Server API-Only, according to Rails terminology) based on a communication paradigm composed of RESTful (REpresentational State Transfer) interfaces and on the exchange of messages encoded in JSON format through the transport layer provided by the HTTP protocol.

The communication interface of RS_Server is constantly evolving and, for this reason, it makes no sense to fully include it in this document. However, it is possible to consult it freely and to see examples of requests that can be made by visiting the URL below.

<https://documenter.getpostman.com/view/4632696/RWTiwfV4?version=latest>

To provide an example, a subset of the RESTful interface of RS_Server is shown in Table 1. These operations are all those available to handle one of the entities of the application domain, namely the publications. Let's then suppose that a user triggers a show operation for a publications characterized by an identifier equal to 1 by visiting the corresponding endpoint. The JSON-encoded response of RS_Server would be something like the one below.

```
1 {  
2   "id": 1,  
3   "doi": "10.1140/epjc/s10052-018-6047-y",  
4   "title": "Uncertainties in WIMP dark matter scattering revisited",  
5   "author": "John Ellis",  
6   "creator": "Springer",  
7   "producer": null,  
8   "...": ...,
```

```

9  "created_at": "2018-08-02T13:27:46.988Z",
10 "updated_at": "2018-08-02T13:27:49.135Z",
11 "...": ...,
12 }

```

3.3 Database

To implement the storage of link-annotated publications, user authentication and the other functionalities it is necessary to define a database schema, which is shown in Figure 3. There are three main entities modeled within the application domain of Readersourcing 2.0:

- **Users:** models the users of the system itself, which are characterized by their personal data, an optional *ORCID* and a boolean used to check if the user wants to receive emails or not. There are also some attributes used to store various kind of tokens to allow operations like password reset;
- **Ratings:** models the ratings given by readers of publications which are characterized by a numerical score. Also, there is a boolean to check if a rating is anonymous or not and another one to check if it has been edited at a later time;
- **Publications:** models the publications rated by their readers which are characterized by an optional *DOI*, by various metadata and by a whole series of attributes used to manage the paths on the server filesystem in order to guarantee the storage of the original and link-annotated files encoded in PDF format.

Moreover, each of these entities is characterized by further attributes (*steadiness*, *informativeness*, ...) which represent the scores/parameters computed by the Readersourcing models.

Figure 3 shows also two relationships (*gives* and *related_to*) which exist between these three entities. These relationships allow to “tie together” the entities to which they refer and to ensure compliance with the *referential integrity* constraint.

In particular, the *gives* relation establishes that a user can give $[0, \dots, n]$ different ratings, while a single rating can be expressed at most by a user. At first glance, the multiplicity equal to 0 shown in Figure 3 regarding users may look strange. The meaning of this constraint is to allow the expression of anonymous ratings. Likewise, the relation *related_to* establishes that a rating is relative to a certain publication, while a publication can be characterized by $[1, \dots, n]$ different ratings. Moreover, this structure allows to comply in a “natural” way with other constraints, such as the fact that if at least at least one publication does not exist, no ratings have to exist.

Endpoint	HTTP Message	Operation	Description
/publications.json	GET	Index	Fetches the entire collection of Publications.
/publications/list	GET	List	Fetches the entire collection of Publications (server-side view).
/publications/1.json	GET	Show	Returns the Publication with identifier equal to 1.
/publications/lookup.json	POST	Lookup	Searches for a Publication; if it doesn't exist, it is fetched from the given URL.
/publications/random.json	GET	Random	Returns a random Publication.
/publications/1/isRated.json	GET	Is Rated	Checks if the Publication with identifier equal to 1 has been rated by at least one reader.
/publications/1/isSavedForLater.json	GET	Is Saved For Later	Checks if the Publication with identifier equal to 1 has been saved for later by the current user.
/publications.json	POST	Create	Creates a new Publication.
/publications/isFetchable.json	POST	Is Fetchable	Checks if the provided URL contains a fetchable Publication.
/publications/extract.json	POST	Extract	Extract the rating url from a link-annotated Publication.
/publications/fetch.json	POST	Fetch	Fetches a Publication from the given URL.
/publications/1/refresh.json	GET	Refresh	Fetches again the Publication with identifier equal to 1.
/publications/1.json	PUT	Update	Updates the Publication with identifier equal to 1.
/publications/1.json	DELETE	Delete	Deletes the Publication with identifier equal to 1.
...

Table 1: Subset of the RESTful interface of RS_Server.

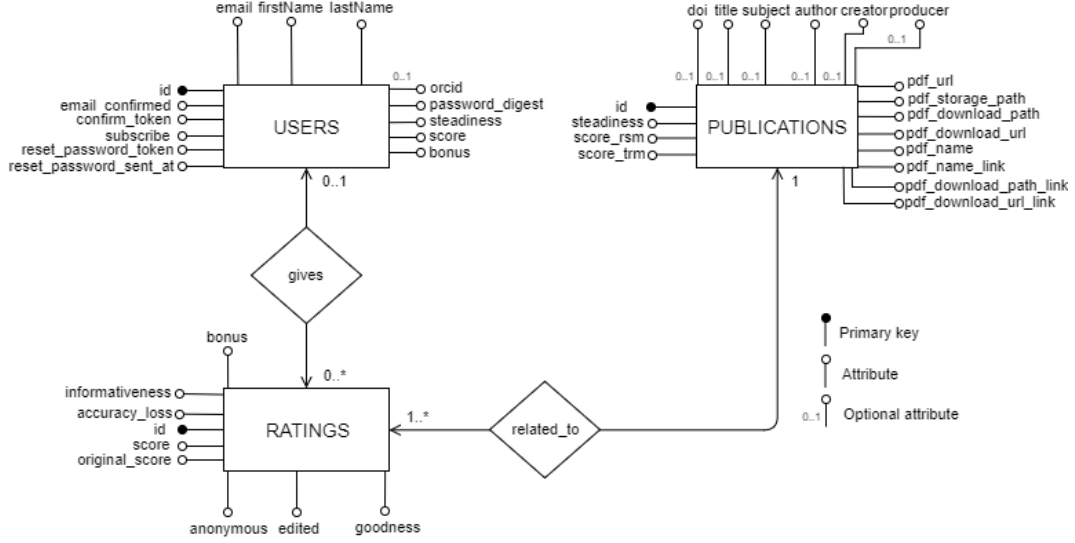


Figure 3: Entity-Relationship schema of the database of RS_Server (NOT UML).

3.4 Class Diagram

Figure 4 shows a diagram of the main classes of RS_Server. As one can see, the convention for which there is an MVC triple for each of the entities modeled in the application domain is followed, although Views are not shown in the diagram because in this case they are just simple methods. The Controller methods represent actions that a user can perform on individual entities or on collections of them, thus mapping the endpoints of the communication protocol used in order to allow the communication between RS_Server and the instances of RS_Rate. As for the Models, their attributes represent the characteristics of the reference entity, while their methods encapsulate the business logic.

Furthermore, there are two additional Controllers⁸ responsible for managing user authentication. RS_Server, as specified previously, is a Web Service; this means that the user interface is presented directly on the instances of RS_Rate and, therefore, those instances send messages to which RS_Server responds once the necessary processing has been completed, according to the RESTful communication paradigm. Because of this design choice, it is not possible to use the “classic” server-side approach to user authentication according to which some information relative to the logged user are saved in the session data, since RESTful paradigm is *stateless*. To be able to authenticate himself, therefore, the user client must attach to each request a *token* that identifies its user as valid within the system. Therefore, a *token-based* authentication approach has been implemented.

When a user performs the first request to RS_Server since some time, he must fill in the login form. If these inserted credentials exist in the database they are encrypted (a *payload* is obtained) and used together with a unique *signature* to create an alphanumeric JSON

⁸ *Application Controller* and *AuthenticationController*.

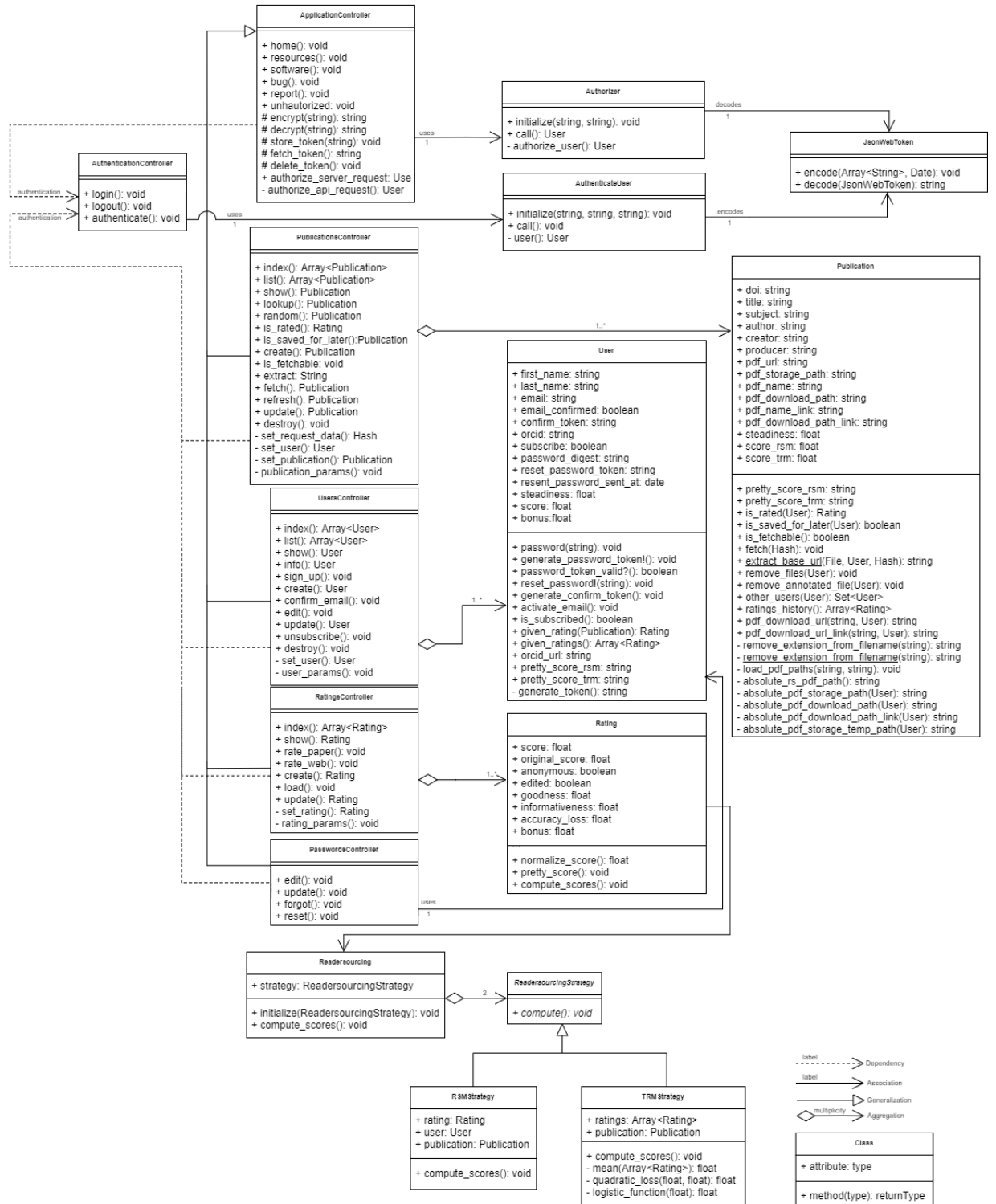


Figure 4: Class diagram of RS.Server.

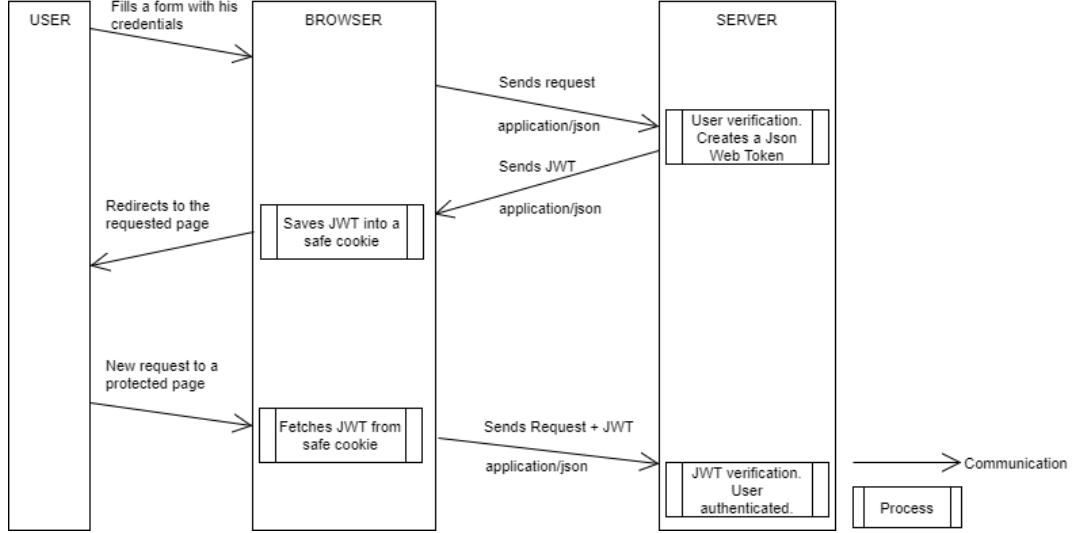


Figure 5: Representation of the token-based authentication process (NOT UML).

string, i.e. the actual token. This token thus generated is sent to the `RS_Rate` instance of the user itself which stores it in a secure cookie characterized by an expiration date after which the procedure must be repeated. At each subsequent request to `RS_Server`, the instance of `RS_Rate` attaches⁹ the previously obtained token in order to demonstrate that its user has successfully completed the authentication procedure. As for `RS_Server`, if a token is present it is extracted and decoded and if it corresponds to one of those saved in the database, then the user identified by the payload is authorized to proceed. An intuitive scheme of the process procedure is shown in Figure 5.

Finally, there is an additional set of classes which are used for different purposes than representing MVC triples for the entities of the application domain. In particular, they are exploited internally by the model containing the business logic to manage the given ratings and have the task of implementing the Readersourcing models. The structure of these classes follows a design pattern called¹⁰ *Strategy* because this pattern allows to integrate new models at a later time without having to make radical changes in the structure of `RS_Server`.

3.5 Deploy

There are three main modalities that can be exploited to deploy a working instance of `RS_Server` in *development* or *production* environment. The former environment must be used if there is the need to add custom Readersourcing model, to extend/modify the current implementation of `RS_Server` or simply to test it in a safe way and it is allowed only by deploy modalities 1 and 2, while the latter must be used if `RS_Server` is about to be made available

⁹In the *Authorization* header of the HTTP payload.

¹⁰*Readersourcing*, *ReadersourcingStrategy*, *RSMStrategy* e *TRMStrategy*.

for everyone as it is and it is allowed by every deploy modality. In the following these three deploy modalities are described, along with their requirements. It is strongly suggested to read also the sections dedicated to the *environment variables* (3.5.4.1 and 3.5.4.2), since if they are not set RS_Server will not work properly.

3.5.1 1: Manual Way

This deploy modality allows to manually download and start RS_Server locally to a machine chosen as a server. This is the most demanding modality regarding its requirements since it assumes that you have a full and working installation of *Ruby*, *JRE* (Java Runtime Environment) and *PostgreSQL* and everything must be set up manually, but it allows more flexibility if a particular setup is required for any reason.

3.5.1.1 Requirements

- Ruby \geq 2.5.3;
- JRE (Java Runtime Environment) \geq 1.8.0;
- PostgreSQL \geq 11.2.

3.5.1.2 How To

Clone RS_Server repository ¹¹ and move inside its main directory using a command line prompt (with an `ls` or `dir` command you should see `app`, `bin`, `config`, etc. folders) and type `gem install bundler`. This gem (dependency) will provide a consistent environment for Ruby projects (like RS_Server) by tracking and installing the exact gems (dependencies) and versions that are needed.

To fetch all those required by RS_Server type `bundle install` and wait for the process to complete. The next two commands are required *only before the first startup of RS_Server* because they will create and set up the database, so please be sure that the PostgreSQL service is started up and ready to accept connections on port 5432. Type `rake db:create` to create the database and `rake db:migrate` to create the required tables inside it. Now, create a `.env` file as explained in section 3.5.4.1 and set the required environment variables. After that, everything is ready to launch RS_Server in development or production mode. To do that, just type `cd bin` to move inside `bin` directory and then `rails server -b 127.0.0.1 -p 3000 -e development` with the proper values for `-b`, `-p` and `-e` options. If the previous values are used, RS_Server will be started and bound on 127.0.0.1 ip address with port 3000 and `development` environment. Every request, therefore, must be sent to `https://127.0.0.1:3000` address.

¹¹https://github.com/Miccighel/Readersourcing-2.0-RS_Server

3.5.1.3 Quick Cheatsheet

1. `cd` to main directory;
2. `gem install bundler`;
3. `bundle install`;
4. `rake db:create` (only before first startup);
5. `rake db:migrate` (only before first startup);
6. create `.env` file and set environment variables;
7. `cd bin`;
8. `rails server -b x.x.x.x -p x -e development` or
`rails server -b x.x.x.x -p x -e production`.

3.5.2 2: Manual Way (But Faster)

This deploy modality allows to download and start RS_Server locally to a machine chosen as a server similarly to the modality described in section 3.5.1, but in a way which is faster and less frustrating, despite being less flexible. Moreover, this deploy modality has less demanding requirements, since only a working installation of *Docker Desktop CE (Community Edition)* is required.

Docker is a project which allows to automate the deployment phase by distributing an *image* of an application inside a *container*. A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another. A Docker container image is a lightweight, standalone, executable package of software that includes everything needed to run an application: code, runtime, system tools, system libraries and settings. This means that there is no need to manually install the runtimes/libraries/dependencies needed to run an application since the Docker Engine will automatically fetch, install and setup them.

3.5.2.1 Requirements

- Docker Desktop CE (Community Edition);

3.5.2.2 How To

Clone RS_Server repository ¹² and move inside its main directory using a command line prompt. Now, type `ls` or `dir`; you should see a `docker-compose.yml` file and a `Dockerfile`. If you do not see them, please be sure to be in the main directory of the cloned repository.

¹²https://github.com/Miccighel/Readersourcing-2.0-RS_Server

Before proceeding, be sure that your Docker Engine has been started up, otherwise the following commands will not work. At this point two different scenarios could happen, which are outlined in the following.

Scenario 1: Deploy With Remote Images If there is no need to edit the source code of RS_Server the Docker Engine can simply fetch the dependencies required in the `docker-compose.yml` file and set up the application. To do this, open the `docker-compose.yml` file and uncomment the section between the

```
----- SCENARIO 1: DEPLOY WITH REMOTE IMAGES -----
```

and

```
----- END OF SCENARIO 1: DEPLOY WITH REMOTE IMAGES -----
```

comments and comment back the remaining lines of code. Now, create a `.env` file as explained in section 3.5.4.1 and set the required environment variables. After that, from the command line prompt type `docker-compose up` and wait for the processing to finish. Note that it may take different minutes. Once the Docker Engine completes the process, the container with a working instance of RS_Server will be started up. *If the first startup of the application is being done* type also `docker-compose run rake db:create` to create the database and `docker-compose run rake db:migrate` to create the required tables inside it. RS_Server will be started and bound on 127.0.0.1 ip address with port 3000 and `production` environment. Every request, therefore, must be sent to `https://127.0.0.1:3000` address. As it can be seen, there is no need to start the server by specifying its ip address, port and environment, since the Docker Engine will take care of that. If you want to set a custom ip address or port or switch to the `development` environment, edit the `command` key inside `docker-compose.yml` file. To shutdown and undeploy the container, simply type `docker-compose down`.

Scenario 2: Deploy With Local Build If the source code of RS_Server has been edited the application must be built locally by the Docker Engine according to the structure specified in the `Dockerfile`. After this build phase the Docker Engine itself can simply fetch the required dependencies outlined in the `docker-compose.yml` file and set RS_Server up. To do this, open the `docker-compose.yml` file and uncomment the section between the

```
----- SCENARIO 2: DEPLOY WITH LOCAL BUILD -----
```

and

```
----- END OF SCENARIO 2: DEPLOY WITH LOCAL BUILD -----
```

comments and comment back the remaining lines of code. Now, create a `.env` file as explained in section 3.5.4.1 and set the required environment variables. After that, from the command line prompt type `docker-compose up` and wait for the process to finish. Note that

it may take different minutes. Once the Docker Engine completes the process, the container with a working instance of RS.Server will be ready. *If the first startup of the application is being done* type also `docker-compose run rake db:create` to create the database and `docker-compose run rake db:migrate` to create the required tables inside it. RS.Server will be started and bound on 127.0.0.1 ip address with port 3000 and `production` environment. Therefore, every request must be sent to `https://127.0.0.1:3000` address. As it can be seen, there is no need to start the server by specifying its ip address, port and environment, since the Docker Engine will take care of that. If you want to set a custom ip address or port or switch to the `development` environment, edit the `command` key inside `docker-compose.yml` file. To shutdown and undeploy the container, simply type `docker-compose down`.

Quick Cheatsheet

- `cd` to main directory;
- create `.env` file and set environment variables;
- `docker-compose up`;
- `docker-compose run rake db:create` (only at first startup);
- `docker-compose run rake db:migrate` (only at first startup);
- `docker-compose down` (to shutdown and undeploy).

3.5.3 3: Heroku Deploy

This deploy modality allows to exploit the container registry of *Heroku* to perform a docker-based production-ready deploy of RS.Server through a working installation of the *Heroku Command Line Interface (CLI)*. Note that this modality can be used only if you choose to use RS.Server in *production* environment.

Heroku is Platform-as-a-Service (PaaS) that enables developers to build, run, and operate applications entirely in the cloud. Regarding the requirements of this modality, an *app* on Heroku must be created and *provisioned* with two addons, namely *PostgreSQL* for the database and *SendGrid* for the mailing functionalities. Follow Heroku tutorials if you do not know it and its concepts. Also, a working installation of *Docker Desktop CE (Community Edition)* on the machine used to perform the deploy is required.

3.5.3.1 Requirements

- Heroku account;
- Heroku application (PostgreSQL + SendGrid Addons);
- Heroku CLI;
- Docker Desktop CE (Community Edition);

3.5.3.2 How To

Clone RS_Server repository ¹³ and move inside the main directory using a command line prompt. Now, type `ls` or `dir`; you should see a `Dockerfile`. If you do not see it, please be sure to be in the main directory of the cloned repository. Before proceeding, *be sure that your Docker Engine has been started up, otherwise the following commands will not work.* Log in to your Heroku account by typing `heroku login` and insert your credentials. Next, log in to Heroku container registry by typing `heroku container:login`. To build and upload your instance of RS_Server type `heroku container:push web --app your-app-name` and when the process terminates type `heroku container:release web` to make it publicly accessible. Optionally, you can type `heroku open` to open the browser and be redirected on the homepage of `your_app_name` application. To create and set up the database type `heroku run rake db:create` and `heroku run rake db:migrate`. As it can be seen, there is no need to start the server by specifying its ip address, port and environment, since Heroku (through the Docker Engine) will take care of that. Now, remember to set the required environment variables on your Heroku app as explained in section 3.5.4.2.

3.5.3.3 Quick Cheatsheet

- `cd` to main directory;
- `heroku login`;
- `heroku container:login`;
- `heroku container:push web --app your-app-name`;
- `heroku container:release web --app your-app-name`;
- `heroku open --app your-app-name` (optional);
- `heroku run rake db:migrate -e production --app your-app-name` (optional);
- `heroku run rake db:seed -e production --app your-app-name` (optional);
- set environment variables on your Heroku app.

3.5.4 Environment Variables

Regardless of the chosen deploy modality, there is the need to set some environment variables which cannot be checked into a repository as a safety measure. In Table 2 each of these environment variables is described along with an explanation of where to set them on the basis of the chosen deploy modality/environment.

¹³https://github.com/Miccighel/Readersourcing-2.0-RS_Server

Env. Variable	Description	Deploy Modality	Env.	Where
SECRET_DEV_KEY	Private key used to encrypt some strings	1 - 2 (Scenario 1, Scenario 2)	development	.env file
SECRET_PROD_KEY	Private key used to encrypt some strings	1 - 2 (Scenario 1, Scenario 2) - 3	production	.env file, Heroku App
SENDGRID_USERNAME	Username of your SendGrid account	1 - 2 (Scenario 1, Scenario 2) - 3	development, production	.env file, Heroku App
SENDGRID_PASSWORD	Password of your SendGrid account	1 - 2 (Scenario 1, Scenario 2) - 3	development, production	.env file, Heroku App
SENDGRID_API_KEY	API key of your SendGrid account	1 - 2 (Scenario 1, Scenario 2) - 3	development, production	.env file, Heroku App
SENDGRID_DOMAIN	A domain registered within your SendGrid account	1 - 2 (Scenario 1, Scenario 2) - 3	development, production	.env file, Heroku App
BUG_REPORT_MAIL	An email address to receive bug reports	1 - 2 (Scenario 1, Scenario 2) - 3	development, production	.env file, Heroku App
CONTACT_MAIL	An email address to receive general questions	1 - 2 (Scenario 1, Scenario 2) - 3	development, production	.env file, Heroku App
RAILS_LOG_TO_STD	If set to true , Rails writes its logs to the standard output. Useful for debugging purposes.	3	production	.env file, Heroku App

Table 2: Environment variables of RS_Server.

3.5.4.1 .env file

To set an environment variable in a local `.env` file, create it inside the main directory of `RS_Server` and populate it in a `key=value` fashion; Listing 1 shows the content of a valid `.env` file.

Listing 1 A valid `.env` file.

```
1: SECRET_DEV_KEY=your_secret_dev_key
2: SENDGRID_USERNAME=your_sendgrid_username
3: SENDGRID_PASSWORD=your_sendgrid_password
4: SENDGRID_DOMAIN=your_sendgrid_domain
5: SENDGRID_API_KEY=your_sendgrid_secret_api_key
6: BUG_REPORT_MAIL=your_bug_report_mail
7: CONTACT_MAIL=your_contact_mail
```

3.5.4.2 Heroku App

To set an environment variable in an Heroku app, simply follow the guide provided by the platform.¹⁴ In Heroku terminology environment variables are called *config vars*.

4 RS_PDF

`RS_PDF` [3] is the software library which is exploited by `RS_Server` to actually edit the PDF files to add the URL required when a reader requests to save for later the publication that he is reading. It is a software characterized by a command line interface and this means that `RS_Server` can use it directly since they are deployed one along the other, without using complex communication channels and paradigms.

4.1 Implementation and Technology

The technology used to develop `RS_PDF` is the Kotlin object-oriented programming language, whose main feature is to be fully compatible with the Java Virtual Machine. This feature is of great importance because it allows a developer to exploit code contained in any other software published in jar format and, more generally, to import any Java class, interacting with them through the syntax of Kotlin itself.

This programming language has been chosen because it has many modern features (it has been created just three years ago) and it is supported rather intensively; furthermore, there are openings to other platforms that have greatly expanded its use possibilities. The most important reason, however, is that the underlying tool used to actually edit files encoded in PDF format is `PDFBox`,¹⁵ which is a software library developed with Java and proposed as

¹⁴<https://devcenter.heroku.com/articles/config-vars>

¹⁵<https://pdfbox.apache.org/>

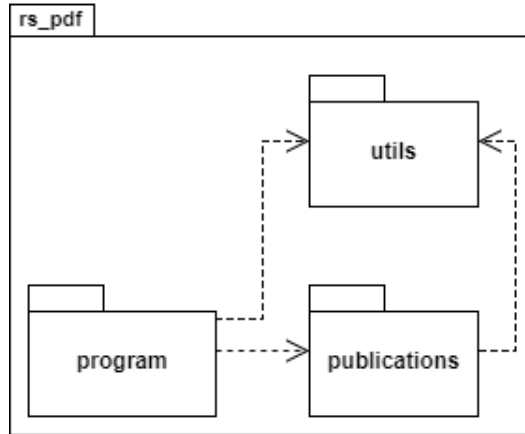


Figure 6: Package diagram of RS_PDF.

a complete toolkit to edit files in that specific format. So, RS_PDF is a wrapper for PDFBox that adds the needed links inside the PDFs requested by readers.

Kotlin has been created by JetBrains¹⁶ which, in the first half of 2017, signed an agreement with Google to let Kotlin become a first-class language for development on the Android platform¹⁷. In the same year, moreover, JetBrains announced the possibility to compile programs written in Kotlin directly into machine language, thus avoiding the use of the JVM. On the web is possible to find different pages with comparisons between Kotlin and other languages, including the official one¹⁸ made by JetBrains with Java, and several articles¹⁹ of developers enthusiastic about this programming language.

4.2 Package Diagram

Figure 6 shows a diagram of the packages in which RS_PDF is divided. This is a useful diagram since it provides a high-level overview of the internal architecture of a software.

In particular, the interaction with RS_Server takes place within the package **program**. The server-side component itself can use the functionalities of RS_PDF by executing it on the JVM, with a special set of command line options. Within this package, therefore, the parsing of the values received for each of these options and the management of the execution flow on the basis of these values take place.

The package **utils** has the task of providing useful tools to the remaining components of RS_PDF. Inside it there are shared constants and methods that allow to access to the logging functionality. As it can be seen by looking at the diagram shown in Figure 6, the other packages depend on it, in particular for some of the values of its constants.

¹⁶<https://www.jetbrains.com/>

¹⁷<https://blog.jetbrains.com/kotlin/2017/05/kotlin-on-android-now-official/>

¹⁸<https://kotlinlang.org/docs/reference/comparison-to-java.html>

¹⁹<https://medium.com/@octskyward/why-kotlin-is-my-next-programming-language-c25c001e26e3>

The package **publications** contains the business logic to handle files encoded in PDF format that must be edited. Its classes follow the logic of the MVC pattern, although its exploiting is not bound by the used technology as in the case of an application developed with Rails. There is, therefore, a Controller which takes into account the execution parameters analyzed in the package **program** and updates the internal state of one or more instances of the Model which will be as many as the files encoded in PDF format that must be annotated. This operation involves loading the input files and adding a link to RS_Server on a new page, taking advantage of the functionalities of *PDFBox*. As a last note, a View is not necessary because RS_PDF simply stores the changes in a new PDF file and, then, ends its execution.

4.3 Class Diagram

Figure 7 shows a diagram of the main classes of RS_PDF which details the internal structure of the architectural elements outlined into the diagram shown in Figure 6. The classes contained within the package **publications** are structured in a way which is similar to what Rails forces in RS_Server and most of the processing carried out by RS_Rate takes place within them. The Model contains the connections with PDFBox and its methods exploit these connections to actually edit files encoded in PDF format.

A single exception to this structure is the use of the *Parameters* class; in particular, it is only a *data class*, i.e. a class whose sole purpose is to store data of various kinds. This instance, once created, is sent to the Model by the Controller through the interfaces of the Model itself. If it is necessary to send further data, the only thing to do consists in adding them to the data class, thus avoiding modifying the signatures of the methods of the Model. Regarding the contents of the *program* and *utils* packages, there is not much else to add with respect to what was said during the description of the diagram shown in figure 6.

4.4 Installation

RS_PDF comes bundled with RS_Server, so when the latter is deployed there is no need to manually install the former. Nevertheless, it is possible to use it independently; it is sufficient to download the published *.jar* files from the release section of its repository²⁰ and place it somewhere on the filesystem.

4.4.1 Requirements

- JRE (Java Runtime Environment) $\geq 1.8.0$;

4.4.2 Command Line Interface

The behavior of RS_PDF is configured during its startup phase by RS_Server through a set of special command-line options. For this reason, it is useful to provide a list of all the options that can be used if it is necessary to use RS_PDF in other contexts, modify

²⁰https://github.com/Miccighel/Readersourcing-2.0-RS_PDF

Short	Long	Description	Values	Req.	Deps.
--pIn	--pathIn	Path on the filesystem from which to load the PDF files to be edited. It can be a file or a folder.	String representing a relative path.	No	--pOut
--pOut	--pathOut	Path on the filesystem in which to save the edited PDF files. It must be a folder.	String representing a relative path.	No	--pIn
--c	--caption	Caption of the link to add.	Any string.	Yes	No
--u	--url	Url to add.	A valid URL.	Yes	No
--a	--authToken	Authentication token received from RS_Server.	A valid authentication token received from RS_Server.	No	--pOut --pIn --pId
--pId	--publicationId	Identifier for a publication present on RS_Server.	A valid publication identifier received from RS_Server.	No	--pOut --pIn --a

Table 3: Command line options of RS_PDF.

its implementation or for any other reason. However, it is designed to work with a default configuration if no options are provided. This list of command line options is shown in Table 3. To provide an execution example, let's assume a scenario in which there is the need of edit some files encoded in PDF format with the following prerequisites:

- there is a folder containing n files to edit at path `C:\data`;
- the edited files must be saved inside a folder at path `C:\out`;
- the file in JAR format containing the library is called `RS_PDF-v1.0-alpha.jar`;
- the JAR file containing RS_PDF is located inside the folder at path `C:\lib`;
- the authentication token received from RS_Server is
`eyJhbGciOiJIUzI1NiJ9...XpC9PMX0jtjRd4NBCTB1a4SfBEi6ndgqsE3k_cEI6Wo`;
- the publication identifier received from RS_Server is 1.

The execution of RS_PDF is started with the following command:

```
java -jar C:\lib\RS_PDF-v1.0-alpha.jar -pIn C:\data -pOut C:\out -a
eyJhbGciOiJIUzI1NiJ9...XpC9PMX0jtjRd4NBCTB1a4SfBEi6ndgqsE3k_cEI6Wo -pId 1
```

5 RS_Rate

RS_Rate [5] is an extension for *Google Chrome*²¹ available on its store. It is a client for Readersourcing 2.0 that can be used by readers of publications to rate them directly from the browser without the need of using this web interface. We intend to generalize RS_Rate by providing an implementation for each of the major browsers (i.e., Firefox, Safari, ...).

5.1 Implementation and Technology

RS_Rate is an extension for Google Chrome; those extensions are developed using standard web technologies such as HTML, CSS and Javascript. Therefore, they are simple “collections” of files packaged in a CRX archive. This particular format is nothing more than a modified version of a ZIP archive with the addition of some special headers exploited by Google Chrome.

As for the Javascript component, RS_Rate does not actually uses the “pure” language but instead uses jQuery, a library developed with the aim of simplifying the selection, manipulation, management of events and the animation of DOM elements in HTML pages, as well as implementing AJAX features. These AJAX features are widely used by RS_Rate to improve the user experience during its use.

5.2 Installation

RS_Rate is freely available on the main browsers webstores. To use it, simply take advantage of the following links and install the version for your favourite browser.

- **Google Chrome** version available at: <https://chrome.google.com/webstore/detail/readersourcing-20-rsrate/hlkdlngpijhdkbdhlmgeemffaocjagg?hl=it>;
- **Firefox** version available at: *currently disabled, will be available at a later time.*

6 RS_Py

RS_Py [4] is an additional component of the Readersourcing 2.0 ecosystem which provides a fully-working implementation of the models presented by Soprano et al. [2] which are encapsulated by the server-side application of Readersourcing 2.0 described in Section 3. Developers with a background in Python programming language can take advantage of RS_Py to generate and test new simulations of ratings given by readers to a set of publications and they are allowed to alter the internal logic of the models to test new approaches without the need to fork and edit the full implementation of Readersourcing 2.0.

²¹<https://www.google.com/chrome/>

6.1 Implementation and Technology

RS.Py is a collection of *Jupyter Notebooks*. Jupyter²² is a Python-powered²³ open-source web application that allows to create and share documents that contain live code, equations, visualizations and narrative text. Notebooks can be shared with others and are made of cells that can be run independently of each other which provide a step-by-step overview to the implemented computations.

6.2 Installation

To use RS.Py notebooks it is sufficient to clone its repository²⁴ and place it somewhere on the filesystem. Be sure to install the required Python packages by taking advantage of a distribution like Anaconda²⁵. If a lightweight installation is preferred, an instance of Python 3.7.3 is needed to install the required packages like Jupyter, Pandas and others.

6.3 Usage

RS.Py is organized into five main folders on the filesystem:

- the **data** folder is used to store the dataset exploited to test the models presented by Soprano et al. [2];
- the **models** folder is used to store the output of such models;
- the **notebooks** folder contains the jupyter notebooks used to generate new datasets and to implement the models presented by Soprano et al. [2];
- the **scripts** folder contains implementations of the jupyter notebooks as pure Python scripts;
- the **src** folder contains a Python scripts which converts jupyter notebooks into pure Python scripts.

Within the **notebooks** folder there are three jupyter notebooks. **Readersourcing.ipynb** provides an implementation of the RSM model presented by Soprano et al. [2] while **TrueReview.ipynb** provides an implementation of the TRM one. Lastly, the **Seeder.ipynb** notebook allows to generate new datasets, which will be stored inside the **data** folder. Inside the **src** folder the **Convert.py** script allows to convert the notebooks in Python scripts which are then stored inside the **scripts** folder.

The behavior of **Seeder.ipynb** and **Readersourcing.ipynb** notebooks can be edited by taking advantage of the parameter setting section which is located in the first rows of both

²²<https://jupyter.org/>

²³<https://www.python.org/>

²⁴https://github.com/Miccighel/Readersourcing-2.0-RS_Py

²⁵<https://www.anaconda.com/distribution/>

Parameter	Description	Values
<code>dataset_name</code>	Name of the dataset to simulate	String
<code>papers_number</code>	Number of papers to simulate	Positive integer
<code>readers_number</code>	Number of readers to simulate	Positive integer
<code>authors_number</code>	Number of authors to simulate	Positive integer
<code>months_number</code>	Number of months of activity to simulate	Positive integer
<code>paper_frequencies</code>	Amount of papers rated by each reader group	Array of positive integers
<code>readers_percent</code>	Percentual of readers to assign to a single group	Positive integer

Table 4: Parameters available for the Seeder jupyter notebook.

Parameter	Description	Values
<code>dataset_name</code>	Name of the dataset to simulate	String
<code>day_serialization</code>	Activate serialization of data on a per day basis	True, False
<code>day_serialization_threshold</code>	Serialize data every X days	Positive integer
<code>days_number</code>	Amount of days simulated in the input dataset	Positive integer

Table 5: Parameter available for the Readersourcing jupyter notebook.

notebooks. Table 4 shows the parameters available for the former, while Table 5 shows the parameters available for the latter.

To run and use the jupyter notebooks move inside the main directory of RS_Py using a command line prompt (with an `ls` or `dir` command you should see `data`, `models`, `notebooks`, etc. folders) and type `jupyter notebook`²⁶. This command will start the Jupyter server and you will be able to take advantage of the *Notebook Dashboard* within your browser at the url of the web application (by default, `http://localhost:8888`).

²⁶<https://jupyter.readthedocs.io/en/latest/running.html#running>

References

- [1] Martin Fowler. *UML Distilled: A Brief Guide to the Standard Object Modeling Language*. 3rd ed. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2003. ISBN: 0321193687.
- [2] M. Soprano and S. Mizzaro. Crowdsourcing Peer Review: As We May Do. In: *Communications in Computer and Information Science* 988 (2019), pp. 259–273.
- [3] Michael Soprano and Stefano Mizzaro. *Readersourcing 2.0: RS_PDF*. DOI: 10.5281/zenodo.1442598. URL: <https://doi.org/10.5281/zenodo.1442597>.
- [4] Michael Soprano and Stefano Mizzaro. *Readersourcing 2.0: RS_Py*. DOI: 10.5281/zenodo.3245208. URL: <https://doi.org/10.5281/zenodo.3245208>.
- [5] Michael Soprano and Stefano Mizzaro. *Readersourcing 2.0: RS_Rate*. DOI: 10.5281/zenodo.1442599. URL: <https://doi.org/10.5281/zenodo.1442599>.
- [6] Michael Soprano and Stefano Mizzaro. *Readersourcing 2.0: RS_Server*. DOI: 10.5281/zenodo.1442630. URL: <https://doi.org/10.5281/zenodo.1442630>.