Small or medium-scale focused research project (STREP)

ICT SME-DCA Call 2013 FP7-ICT-2013-SME-DCA

### Data Publishing through the Cloud: A <u>Da</u>ta- and <u>P</u>latform-<u>a</u>s-<u>a</u>-<u>S</u>ervice Approach to Efficient Open Data Publication and Consumption

DaPaaS



### **Deliverable D2.1:**

# Open PaaS requirements, design & architecture specification

| Date:                | 30.01.2014   |
|----------------------|--|
| Author(s):           | Brian Elvesæter, Dumitru Roman, Martin Fagereng Johansen,<br>Arne J. Berre, Marin Dimitrov, and Alex Simov |
| Dissemination level: | PU   |
| WP:                  | WP2  |
| Version:             | 1.0  |



### Document metadata

#### **Quality assurors and contributors**

| Quality assuror(s) | Bill Roberts, Rick Moynihan, Amanda Smith |
|--------------------|---|
| Contributor(s)     | DaPaaS Consortium                         |

#### Version history

| Version | Date       | Description   |
|---------|------------|---|
| 0.1     | 04.12.2013 | Initial outline and Table of Contents (TOC).  |
| 0.2     | 09.12.2013 | Restructuring of deliverable with comments.   |
| 0.3     | 16.12.2013 | Draft section on requirements specification.  |
| 0.4     | 10.01.2014 | Updated the requirements specification with description of key roles.   |
| 0.5     | 16.01.2014 | Updated the technology review section and description of requirements.  |
| 0.6     | 17.01.2014 | Consistency check and update of<br>sections 1 (Introduction), 2<br>(Requirements Specification) and<br>3.1 (High-Level Architecture of the<br>DaPaaS Platform). |
| 0.7     | 21.01.2014 | Updated the architecture description of the Platform Layer.   |
| 0.8     | 23.01.2014 | Updated the review of relevant technologies.  |
| 0.9     | 28.01.2014 | Finalized technology review.<br>Deliverable ready for internal<br>technical review.   |
| 0.95    | 29.01.2014 | Addressed comments by internal technical review. Deliverable ready for quality assurors.  |
| 1.0     | 30.01.2014 | Addressed comments by quality assurors. Final formatting and layout.  |



### **Executive Summary**

The main goal of the DaPaaS project is to provide an integrated Data-as-a-Service (DaaS) and Platformas-a-Service (PaaS) environment, together with associated services, for open data, where 3rd parties can publish and host both datasets and data-driven applications that are accessed by end user data consumers in a cross-platform manner.

This document provides:

- An overview of the DaPaaS Platform and the relevant roles played in the DaPaaS context;
- The requirements for the DaPaaS Platform;
- An initial architecture design for the Platform Layer of the DaPaaS Platform;
- A state-of-the-art overview of relevant solutions and technologies for the Platform Layer and some recommendations on reuse of existing solutions to be considered in the next phase implementation of the first prototype.



### **Table of Contents**

| E2           | EXECUTIVE SUMMARY   |  |  |
|--------------|---|--|--|
| T/           | TABLE OF CONTENTS   |  |  |
| L            | LIST OF ACRONYMS  |  |  |
| L            | IST OF FIG  | SURES  | 7  |
| LI           | IST OF TA   | BLES   |  |
| 1            | INTRO   | DICTION  | 0  |
| 1            |   |  |  |
|              | 1.1 DAP<br>1.2 STRU   | AAS OVERVIEW AND KEY ROLES   | 9  |
| 2            | DAPA  | AS PLATFORM REQUIREMENTS SPECIFICATION   | 11   |
| _            | 2.1 INST  |  | 11   |
|              | 2.1 INST  | A PUBLISHER  | 12   |
|              | 2.3 Appl  | ICATION DEVELOPER  | 13   |
|              | 2.4 END   | USER DATA CONSUMER   | 15   |
| 3            | ARCH  | TECTURE OVERVIEW   | 17   |
|              | 3.1 High  | -Level Architecture of DaPaaS Platform   | 17   |
|              | 3.2 Arc   | HITECTURE OF THE PLATFORM LAYER  | 17   |
|              | 3.2.1   | User Management & Access Control   | 18   |
|              | 3.2.2   | Data Cleaning & App Development  | 19   |
|              | 3.2.3   | Notification   | 20   |
|              | 3.2.4   | App Management & Deployment  | 20   |
|              | 3.2.5   | Catalog  | 21   |
|              | 3.2.0<br>3.3 Sum  | Administration   | 21   |
| 4            | DEVIE   |  | 22   |
| 4            | KEVIE   | W UF KELEVANT TECHNOLOGIES FOK THE PLATFORM LATEK  | 24   |
|              |   |  |  |
|              | 4.1 TECH  | INOLOGY SELECTION APPROACH   | 24   |
|              | 4.1 TECH<br>4.2 PAAS  | INOLOGY SELECTION APPROACH   | 24<br>26   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>26<br>27   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3   | INOLOGY SELECTION APPROACH<br>S CAPABILITIES SOLUTIONS<br>Docker<br>Cocaine<br>Deis  | 24<br>26<br>26<br>27<br>28   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4  | INOLOGY SELECTION APPROACH<br>S CAPABILITIES SOLUTIONS<br>Docker<br>Cocaine<br>Deis<br>Juju  | 24<br>26<br>26<br>27<br>28<br>29   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.4<br>4.2.5  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>26<br>27<br>28<br>29<br>29   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.5<br>4.2.6   | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>26<br>27<br>28<br>29<br>29<br>29<br>30   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7   | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>29<br>29<br>30   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8  | INOLOGY SELECTION APPROACH   | 24<br>26<br>27<br>27<br>28<br>29<br>29<br>29<br>29<br>30<br>30<br>30<br>31   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10   | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>29<br>30<br>30<br>31<br>32   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11   | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>29<br>30<br>30<br>31<br>32<br>33<br>33   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DAT  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>30<br>30<br>31<br>32<br>33<br>34<br>35   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.10<br>4.2.11<br>4.3 DATL<br>4.3.1  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>30<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>35   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DAT<br>4.3.1<br>4.3.2  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>29<br>30<br>30<br>31<br>31<br>32<br>33<br>34<br>35<br>35<br>36   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DATL<br>4.3.1<br>4.3.2<br>4.3.3  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>36<br>37   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DAT<br>4.3.1<br>4.3.2<br>4.3.3<br>4.3.4  | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>36<br>37<br>37   |
|              | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DATL<br>4.3.1<br>4.3.2<br>4.3.3<br>4.3.4<br>4.3.5  | INOLOGY SELECTION APPROACH<br>SCAPABILITIES SOLUTIONS.<br>Docker   | 24<br>26<br>27<br>28<br>29<br>29<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>35<br>36<br>37<br>37<br>38   |
| 5            | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DAT<br>4.3.1<br>4.3.2<br>4.3.3<br>4.3.4<br>4.3.5<br>SUMM   | INOLOGY SELECTION APPROACH<br>CAPABILITIES SOLUTIONS.<br>Docker.<br>Cocaine.<br>Deis.<br>Juju<br>Cozy Cloud.<br>OpenCivic<br>OpenStack<br>Ansible<br>Puppet Open Source<br>Chef<br>Nagios Core.<br>A INTEGRATION CAPABILITIES SOLUTIONS.<br>Talend Open Studio for Data Integration<br>OpenRefine<br>Karma.<br>Cascading<br>Data Pipes.<br>ARY AND OUTLOOK.  | 24<br>26<br>26<br>26<br>26<br>26<br>27<br>28<br>30<br>30<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>36<br>37<br>37<br>37<br>38   |
| 56           | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DATL<br>4.3.1<br>4.3.2<br>4.3.3<br>4.3.4<br>4.3.5<br>SUMM<br>APPEN   | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>27<br>28<br>29<br>29<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>35<br>36<br>37<br>37<br>38<br><b>39</b>  |
| 5<br>6<br>S( | 4.1 TECH<br>4.2 PAAS<br>4.2.1<br>4.2.2<br>4.2.3<br>4.2.4<br>4.2.5<br>4.2.6<br>4.2.7<br>4.2.8<br>4.2.9<br>4.2.10<br>4.2.11<br>4.3 DAT<br>4.3.1<br>4.3.2<br>4.3.3<br>4.3.4<br>4.3.5<br>SUMM<br>APPEN<br>OLUTIONS  | INOLOGY SELECTION APPROACH<br>G CAPABILITIES SOLUTIONS.<br>Docker<br>Cocaine<br>Deis<br>Juju<br>Juju<br>Cozy Cloud.<br>OpenCivic<br>OpenCivic<br>OpenStack<br>Ansible<br>Puppet Open Source<br>Chef<br>Nagios Core<br>A INTEGRATION CAPABILITIES SOLUTIONS.<br>Talend Open Studio for Data Integration<br>OpenRefine<br>Karma.<br>Cascading<br>Data Pipes.<br>ARY AND OUTLOOK.<br>DIX A: COMMERCIAL / CLOSED SOURCE INTEGRATED DAAS & PAAS | 24<br>26<br>26<br>26<br>26<br>27<br>28<br>29<br>30<br>30<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>36<br>37<br>37<br>38<br>38<br>39   |
| 5<br>6<br>S( | <ul> <li>4.1 TECH</li> <li>4.2 PAAS</li> <li>4.2.1</li> <li>4.2.2</li> <li>4.2.3</li> <li>4.2.4</li> <li>4.2.5</li> <li>4.2.6</li> <li>4.2.7</li> <li>4.2.8</li> <li>4.2.9</li> <li>4.2.10</li> <li>4.2.11</li> <li>4.3 DAT.</li> <li>4.3.1</li> <li>4.3.2</li> <li>4.3.3</li> <li>4.3.4</li> <li>4.3.5</li> <li>SUMM</li> <li>APPEN</li> <li>DLUTIONS</li> <li>6.1 DAT.</li> </ul>                                     | INOLOGY SELECTION APPROACH<br>GCAPABILITIES SOLUTIONS.<br>Docker   |  |
| 5<br>6<br>S( | <ul> <li>4.1 TECH</li> <li>4.2 PAAS</li> <li>4.2.1</li> <li>4.2.2</li> <li>4.2.3</li> <li>4.2.4</li> <li>4.2.5</li> <li>4.2.6</li> <li>4.2.7</li> <li>4.2.8</li> <li>4.2.9</li> <li>4.2.10</li> <li>4.2.11</li> <li>4.3 DAT.</li> <li>4.3.1</li> <li>4.3.2</li> <li>4.3.3</li> <li>4.3.4</li> <li>4.3.5</li> <li>SUMM</li> <li>APPEN</li> <li>DLUTIONS</li> <li>6.1 DAT.</li> <li>6.2 SPLU</li> </ul>                   | INOLOGY SELECTION APPROACH   |  |
| 5<br>6<br>S( | <ul> <li>4.1 TECH</li> <li>4.2 PAAS</li> <li>4.2.1</li> <li>4.2.2</li> <li>4.2.3</li> <li>4.2.4</li> <li>4.2.5</li> <li>4.2.6</li> <li>4.2.7</li> <li>4.2.8</li> <li>4.2.9</li> <li>4.2.10</li> <li>4.2.11</li> <li>4.3 DATL</li> <li>4.3.1</li> <li>4.3.2</li> <li>4.3.3</li> <li>4.3.4</li> <li>4.3.5</li> <li>SUMM</li> <li>APPEN</li> <li>DLUTIONS</li> <li>6.1 DATL</li> <li>6.2 SPLU</li> <li>6.3 WINT</li> </ul> | INOLOGY SELECTION APPROACH   | 24<br>26<br>26<br>26<br>27<br>28<br>29<br>30<br>30<br>31<br>32<br>33<br>34<br>35<br>35<br>36<br>37<br>37<br>38<br>39<br>39<br>39<br>39<br>39<br>39<br>39<br>39<br>39<br>31<br>32<br>39<br>31<br>32<br>33<br>34<br>35<br>36<br>37<br>37<br>38<br>39<br>34<br>35<br>36<br>37<br>37<br>38<br>39<br>34<br>35<br>36<br>37<br>38<br>39<br>34<br>35<br>36<br>37<br>38<br>39<br>34<br>35<br>36<br>37<br>38<br>39<br>34<br>35<br>36<br>37<br>38<br>34<br>34<br>39<br>34<br>34<br>37<br>37<br>38<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34<br>34 |

DaPaaS

| 6.5 | TABLEAU SOFTWARE | 41 |
|-----|------------------|----|
| 6.6 | INFOCHIMPS       | 41 |

## List of Acronyms

| API    | Application Programming Interface   |
|--------|-------------------------------------|
| CSV    | Comma Separated Values (format)     |
| DaaS   | Data-as-a-Service                   |
| GUI    | Graphical User Interface            |
| HTTPS  | Hypertext Transfer Protocol Secure  |
| JSON   | JavaScript Object Notation (format) |
| PaaS   | Platform-as-a-Service               |
| REST   | Representational state transfer     |
| RDF    | Resource Description Framework      |
| SLA    | Service Level Agreement             |
| SOA    | Service Oriented Architecture       |
| SPARQL | SPARQL Protocol and RDF Query Lan-  |
|        | guage                               |
| SSH    | Secure Shell                        |
| UML    | Unified Modeling Language           |
| XML    | eXtensible Markup Language          |



## List of Figures

| Figure 1: DaPaaS artefacts                                   | 9   |
|--|-----|
| Figure 2: Key roles in a typical DaPaaS context              | 10  |
| Figure 3: Instance Operator (IO) requirements                | .11 |
| Figure 4: Data Publisher (DP) requirements                   | 12  |
| Figure 5: Application Developer (AD) requirements            | 14  |
| Figure 6: End User Data Consumer (EU) requirements           | 15  |
| Figure 7: High-level architecture of the DaPaaS Platform     | 17  |
| Figure 8: Architecture of the Platform Layer                 | 18  |
| Figure 9: Basic Docker functions                             | 27  |
| Figure 10: Cocaine architecture                              | 28  |
| Figure 11: Juju administration GUI                           | 29  |
| Figure 12: OpenStack conceptual architecture                 | 31  |
| Figure 13: How Puppet works                                  | 32  |
| Figure 14: How Chef works                                    | 33  |
| Figure 15: Operating principle of Nagios                     | 34  |
| Figure 16: Talend Open Studio for Data Integration           | 35  |
| Figure 17: Edit cells in OpenRefine – Common transformations | 36  |
| Figure 18: Modelling data in Karma                           | 37  |
| Figure 19: Cascading architecture                            | 38  |
|  |     |



### List of Tables

| Table 1: Description of requirements from Instance Operator (IO)          | 11 |
|---|----|
| Table 2: Description of requirements from the Data Publisher (DP)         | 13 |
| Table 3: Description of requirements from the Application Developer (AD)  | 14 |
| Table 4: Description of requirements from the End User Data Consumer (EU) | 15 |
| Table 5: User Management & Access Control                                 |    |
| Table 6: App Development  | 19 |
| Table 7: Notification   | 20 |
| Table 8: App Management & Deployment                                      | 21 |
| Table 9: Catalog  | 21 |
| Table 10: Administration  | 21 |
| Table 11: Addressed requirements by components of the Platform Layer      | 22 |
| Table 12: Overview of relevant open source technologies                   |    |
|   |    |



### 1 Introduction

This report represents Deliverable D2.1 "Open PaaS requirements, design & architecture specification" of the DaPaaS project. This deliverable is a result of Task T2.1 "Requirements, analysis & design of the Open Platform-as-a-Service infrastructure".

The aim of this deliverable is two-fold:

- 1. To introduce the DaPaaS platform, the relevant roles played in the DaPaaS context, and their requirements towards a Data- and Platform-as-a-Service infrastructure for open data;
- To provide details on the Platform Layer of the DaPaaS platform, with a focus on the architecture and evaluation of existing relevant technologies that could be reused for the implementation of the Platform Layer.

#### 1.1 DaPaaS Overview and Key Roles

The main goal of the DaPaaS project is to provide an integrated Data-as-a-Service (DaaS) and Platformas-a-Service (PaaS) environment for open data, where 3<sup>rd</sup> parties can publish and host both datasets and data-driven applications that are accessed by end user data consumers in a cross-platform manner. The DaPaaS project will deliver the software that enables platform operators to deploy such an environment in the cloud. Figure 1 below illustrates the idea that the DaPaaS software (DaaS and PaaS functionalities) can have several deployed instances.



Figure 1: DaPaaS artefacts

As the main results of the DaPaaS project two major artefacts are expected:

- 1. Software consisting of DaaS, PaaS, and associated services;
- 2. One deployed instance of the Software in an XaaS manner. In the rest of this deliverable we will refer to this deployed instance as "DaPaaS Platform".

The key roles involved in a typical DaPaaS context and their relationships to the main DaPaaS artefacts, the software and the platform, are illustrated in Figure 2 below. The roles are:

- The **DaPaaS Developer** implements DaPaaS software components and services for the integrated DaaS and PaaS environment. During the course of the project, this role is expected to be primarily played by the DaPaaS consortium.
- A deployed instance of DaPaaS software, i.e. the DaPaaS Platform, is operated and maintained by an **Instance Operator**. During the course of the project, this role is played by the DaPaaS consortium.
- The **Data Publisher** publishes data on the DaPaaS Platform which stores the data and makes it available for 3<sup>rd</sup> party application developers and end user data consumers.



- The **Application Developer** develops data-driven applications that use the data made available via the DaPaaS Platform. The applications are deployed and hosted in the DaPaaS Platform.
- End Users Data Consumers consume data resulting from the deployed applications.



Figure 2: Key roles in a typical DaPaaS context

This document outlines the requirements for the DaPaaS platform from a role-based point of view with a focus on the services and functionality required by the Instance Operator, Data Publisher, Application Developer and End Users Data Consumer.

#### **1.2 Structure of this Report**

The rest of this document is structured as follows:

- Section 2 describes the core requirements for the DaPaaS Platform from the perspective of the key roles introduced above;
- Based on the requirements identified in Section 2, Section 3 outlines the high-level architecture of the DaPaaS Platform, and details the Platform Layer in terms of core components and their relationships;
- Section 4 provides a review of relevant open source technologies for the implementation of Platform Layer;
- Section 5 summarizes this document and provides technical recommendations for the implementation phase of the DaPaaS project; and
- Appendix A provides a brief summary of selected commercial/closed source solutions that provide capabilities relevant to the DaPaaS Platform.



### 2 DaPaaS Platform Requirements Specification

In the following subsections we use a UML use case-inspired notation and technique to describe the requirements of the DaPaaS Platform, e.g. capabilities or services that should be offered for the roles introduced above.

#### 2.1 Instance Operator

The Instance Operator role is played by organizations that want to operate and maintain an instance of the DaPaaS Platform, e.g. acting as data brokers or creating data markets in various domains (e.g. environmental domain). Figure 3 below shows the requirements the Instance Operator poses on the DaPaaS Platform.



#### Figure 3: Instance Operator (IO) requirements

Descriptions of these requirements are given in Table 1 below.

#### Table 1: Description of requirements from Instance Operator (IO)

| ID    | Name   | Brief description  |
|-------|--|--|
| IO-01 | Secure access to plat-<br>form                           | The Instance Operator shall have secure access (e.g. HTTPS/SSH) to the platform.   |
| IO-02 | Platform performance monitoring                          | The Instance Operator shall be able to monitor the performance (e.g. storage and memory usage, bandwidth, CPU usage, etc.).  |
| IO-03 | Statistics monitoring<br>(users, data, apps, us-<br>age) | The Instance Operator shall be able to retrieve statistics about<br>users (e.g. number, profiles), data (e.g. number, size), apps and<br>usage (e.g. dataset access, data consumption, number of ser-<br>vice calls) as a basis for e.g. billing/invoicing for the usage of<br>the platform. |

| IO-04 | Usage accounts man-<br>agement                  | The Instance Operator shall be able to manage user accounts (e.g. add, delete, assign roles).   |
|-------|---|---|
| IO-05 | Policy/quota configura-<br>tion and enforcement | The Instance Operator shall be able to configure usage policies,<br>e.g. data/apps quotas per user. The platform shall ensure en-<br>forcement of these policies, e.g. support deployment of applica-<br>tions subject to quotas and additional restrictions. |
| IO-06 | UI for Instance Opera-<br>tor                   | The Instance Operator shall be able to access the platform ser-<br>vices through appropriate user interface (graphical and/or con-<br>sole).  |

#### 2.2 Data Publisher

The Data Publisher role is played by organizations that want to publish data via the DaPaaS Platform. Figure 4 below depicts the Data Publisher poses on the DaPaaS Platform.



Figure 4: Data Publisher (DP) requirements

Descriptions of these requirements are given in Table 2 below.



#### Table 2: Description of requirements from the Data Publisher (DP)

| ID    | Name  | Brief description  |
|-------|---|--|
| DP-01 | Dataset import  | The Data Publisher should have the ability to import open data into the DaPaaS platform. The data is <i>not</i> restricted to RDF / Linked Data and it may include other formats such as CSV, JSON, etc. |
| DP-02 | Data storage & querying   | The Data Publisher should have access to APIs and query<br>endpoints for accessing, querying and updating data stored<br>on the platform.  |
| DP-03 | Dataset search & explo-<br>ration                                   | The Data Publisher should have the possibility to explore the dataset catalog & select relevant datasets.  |
| DP-04 | Data interlinking   | The Data Publisher should have the possibility to semi-auto-<br>matically interlink data from different datasets. This applies<br>only to data which is already in RDF form.                             |
| DP-05 | Data cleaning & transfor-<br>mation                                 | The Data Publisher should have the possibility to apply sim-<br>ple data cleanup & transformation (incl. RDFization) over<br>legacy data.  |
| DP-06 | Dataset bookmarking & notifications                                 | The Data Publisher should have possibility to subscribe to datasets and receive notifications on datasets changes.   |
| DP-07 | Dataset metadata man-<br>agement, statistics & ac-<br>cess policies | The Data Publisher should have possibility to specify metadata, descriptions and access control policies for the datasets.   |
| DP-08 | Data scalability  | The platform should include mechanisms to scale to large data volumes.   |
| DP-09 | Data availability   | The platform should include mechanisms to provide high availability of data and limited downtime.  |
| DP-10 | User registration & profile management                              | The Data Publisher shall be able to register as a data pub-<br>lisher and gain access to the relevant DaaS services.   |
| DP-11 | Secure access to plat-<br>form                                      | The Data Publisher shall have secure access (e.g. HTTPS/SSH) to the platform.  |
| DP-12 | UI for Data Publisher   | The Data Publisher shall be able to access the DaaS ser-<br>vices through appropriate user interfaces (graphical and/or<br>console).   |
| DP-13 | Data publishing method-<br>ology support                            | The data publication process should be accompanied by a tool-supported methodology outlining steps containing various data operations  |

#### 2.3 Application Developer

The Application Developer role is played by Open Data application developers that for various reasons (e.g. transparency, new business models, new services) want to develop new applications and services around data and want to do so as fast as possible and as easy as possible. Figure 5 below depicts the requirements the Application Developer poses on the DaPaaS Platform.



Figure 5: Application Developer (AD) requirements

Descriptions of these requirements are given in Table 3 below.

| ID    | Name   | Brief description  |
|-------|--|--|
| AD-01 | Access to Data Pub-<br>lisher services (DP-01<br>– DP-13)          | The Application Developer shall have access to APIs and librar-<br>ies to access, import, transform, store, query, etc., datasets to<br>be used in the development of applications. Basically the Appli-<br>cation Developer has similar requirements as outlined in DP-01<br>– DP-13. This includes also requirements for secure access to<br>the platform, profile management. |
| AD-02 | Data export  | The Application Developer shall have the possibility to export data in various formats.  |
| AD-03 | Develop applications<br>in state-of-art pro-<br>gramming languages | The Application Developer shall have the possibility to develop applications in the common state-of-art programming lan-<br>guages, e.g. Java, Scala, Go, Ruby.  |
| AD-04 | Configure application deployment                                   | The Application Developer shall have the possibility to configure<br>use of common cloud resources, e.g. database/storage, possi-<br>ble also graphical widgets.   |
| AD-05 | Deploy and monitor application                                     | The Application Developer shall have access to a controlled ap-<br>plication hosting environment where data-intensive applications<br>can be easily deployed, as well as monitoring facilities for the<br>deployed applications.   |

| Table 3: Description | of requirements | from the Ap | oplication D | Developer ( | (AD) |
|----------------------|-----------------|-------------|--------------|-------------|------|
|                      | •••••           |             |              |             | ~ /  |

| AD-06 | Application metadata<br>management, statis-<br>tics & access policies | The Application Developer shall have the possibility to update metadata about applications (e.g. description) and retrieve statistics about the usage of the application. |
|-------|---|---|
| AD-07 | UI for Application De-<br>veloper                                     | The Application Developer shall have the possibility to access<br>the relevant DaaS and PaaS services through appropriate user<br>interfaces (graphical and/or console).  |
| AD-08 | Application develop-<br>ment methodology<br>support                   | Application Developers should have access to a tool-supported methodology outlining steps for developing and deploying data-<br>intensive applications.                   |

#### 2.4 End User Data Consumer

The End User Data Consumer role is played by organizations or individuals that want to consume data and applications deployed on the platform. Figure 6 below shows the requirements the End User Data Consumer poses on the DaPaaS Platform.



Figure 6: End User Data Consumer (EU) requirements

Descriptions of these requirements are given in Table 4 below.

| ID    | Name   | Brief description  |
|-------|--|--|
| EU-01 | User registration & profile management                       | End Users shall be able to register as application con-<br>sumers and manage their profiles.   |
| EU-02 | Search & explore datasets and applications                   | End Users shall be able to search and explore datasets and applications available in the platform.   |
| EU-03 | Datasets and applications bookmarking and notifica-<br>tions | End Users shall be able to bookmark and receive notifi-<br>cations (e.g. updates) of datasets and applications to<br>which they subscribe. |

| Table 4: Description of red | auiromonte from the En | d Llear Data Concumar (ELI) |
|-----------------------------|------------------------|-----------------------------|
| Table 4. Description of rec | quirements nom the En  | u usel Dala Consumer (EU)   |

| EU-04 | Mobile and desktop GUI access              | End Users shall be able to access applications on both<br>mobile and desktop devices, which requires UX compo-<br>nents to support both mobile and desktop users in an ap-<br>propriate manner. The End User Data Consumers shall<br>be able to access the relevant platform services, e.g.,<br>search for datasets, applications, run applications, visual-<br>ize datasets, etc., through appropriate graphical user in-<br>terfaces (GUIs), e.g. pie charts, time series and maps. |
|-------|--|---|
| EU-05 | Data export and download                   | End Users shall have the possibility to export data in various formats and download data from the platform.   |
| EU-06 | High availability of data and applications | High availability of data and apps  |



### **3** Architecture Overview

This section outlines the high-level architecture of the DaPaaS Platform (Section 3.1), and details the Platform Layer in terms of core components and their interactions (Section 3.2).

#### 3.1 High-Level Architecture of DaPaaS Platform

The requirements outlined in the previous section imply a layered architecture consisting of a Data-asa-Service layer (**Data Layer**) for scalable data hosting, a Platform-as-a-Service layer (**Platform Layer**) for application development and hosting, and a User Experience Layer (**UX Layer**) for user-friendly access to data and applications. These three core layers cross-cut vertical layers that are related to methodology support for data publishing and applications. Figure 7 below illustrates a simplified, high-level architecture of the DaPaaS Platform.



Figure 7: High-level architecture of the DaPaaS Platform

The rest of this deliverable focuses on the Platform Layer of the architecture. An architecture of the Platform Layer is described in the next (Section 3.2) and evaluation of exiting relevant technologies that could be reused for the implementation of the Platform Layer are discussed in Section 4.

The Data Layer and the UX Layer fall within the scope of Deliverable D1.1 and Deliverable D3.1, respectively, and further details about those layers can be found in those deliverables.

#### 3.2 Architecture of the Platform Layer

Figure 8 depicts the main software components, their relationships and associated APIs of the Platform Layer. The software components of the Platform Layer extend the capabilities offered by the Data Layer (described in Deliverable D1.1) in five main service categories plus an administration service:

- User Management & Access Control which manages user profiles and secure access control to apps and datasets.
- **Data Cleaning & App Development** which provides functionalities for applications development, and support for data cleaning & transformation and data workflows.
- **Notification** which provides functionality for subscribing to apps and datasets events and notifications.
- **App Management & Deployment** which gives developers control over the deployed applications and configuration settings for the application-hosting environment.
- Catalog for searching and exploring apps and datasets.
- Administration which allows the management and monitoring of the DaPaaS Platform, focusing on aspects related to the users, apps, datasets and services of the platform.

The functionality of these services can be used by the UX Layer through well-defined APIs. This allows for the creation of e.g. management consoles in the UX Layer. The public APIs of the DaPaaS Platform can also be used by 3<sup>rd</sup> party applications and services.

The design of the Platform Layer involves a set of software components: User Manager, Access Control Manager, Data Cleaning and Design-Time App Development Services, Run-Time App Hosting Environment, Notification Service and Apps Catalog. Each component is responsible for a functional area of the Platform Layer and communicates with each other through internal service interfaces. The components can access and use the functionality offered through the Data Layer API.



Figure 8: Architecture of the Platform Layer

#### 3.2.1 User Management & Access Control

The **User Management & Access Control** service category will be implemented by the **User Manager** and **Access Control Manager** components. The User Manager component is responsible for managing the registered users of the DaPaaS Platform and their user profiles. Users can register as Data Publishers, Application Developers and End User Data Consumers. The user profile may contain information about data quotas, apps and datasets access which are enforced by the Access Control Manager.

The User Management & Access Control services will be offered through a REST API. The API will provide create, read, update and delete (CRUD) actions for managing accounts, user profiles and access control for apps and datasets.



| Component                    | Addressed<br>Requirement | Description  |
|------------------------------|--------------------------|--|
| User Manager                 | DP-06                    | The user profile stores information about bookmarked datasets.<br>The User Manager component uses the Notification Service to<br>provide notification on dataset changes.  |
|                              | DP-10                    | User registration is implemented by the User Manager. It allows<br>users of the DaPaaS Platform to sign up as Data Publisher,<br>Application Developer and/or End User Data Consumer. User<br>account details, preferences and other relevant information are<br>stored in a user profile. The User Manager component provides<br>basic user functionality for all types of users. |
|                              | AD-01                    | Covers the functionality such as DP-06 and DP-10, but for the Application Developer.   |
|                              | EU-01                    | Same functionality as DP-10, but for the End User Data Consumer.   |
|                              | EU-03                    | The user profiles for an End User Data Consumer provides additional support for bookmarked apps and notifications.   |
| Access<br>Control<br>Manager | DP-07                    | The Access Control Manager allows access policies for the datasets to be specified.  |
|                              | AD-06                    | The Access Control Manager allows access policies for the apps to be specified.  |

#### 3.2.2 Data Cleaning & App Development

The Data Cleaning & App Development service category is a collection of Data Cleaning & Design-Time App Development Services, of which App Configuration, Data Workflows and Data Cleaning & Transformation are the main sub-components. The App Configuration is responsible for providing standardized mechanisms to configure cloud resources (e.g. data storage), DaaS services and UX components to be used by the deployed instance of the app at run-time. The Data Cleaning & Transformation component provides additional data management functionalities that complement Data Layer functionality with capabilities for data cleaning (duplicate removal), data transformation, as well as data mapping and alignment. The Data Workflows component provides the capability to define simple data-driven pipelines that use functionality of the Data Layer (e.g., import/export and publish data) and the added functionality offered by the Data Cleaning & Transformation component in order to support simple sequential data transformations.

The App Development services will be offered through a REST API. The API will provide create, read, update and delete (CRUD) actions for managing app configurations. In addition the API will provide actions for accessing and using the DaaS services, including the ability to create simple data workflows.

| Component            | Addressed<br>Requirement | Description  |
|----------------------|--------------------------|--|
| App<br>Configuration | AD-03                    | Support for state-of-the-art programming languages will be an evaluation criterion in the review and selection of the PaaS Infrastructure (Section 4).           |
|                      | AD-04                    | The App Configuration provides a standard way of configuring<br>the required cloud resources, DaaS services and UX<br>components to be used for app development. |

Table 6: App Development

| Data Workflows                 | AD-01              | The Application Developer shall have access to a range of<br>DaaS capabilities that can be used for app development. The<br>set of DaaS capabilities will be provided through the Data<br>Workflows component.   |
|--------------------------------|--------------------|--|
|                                | AD-02              | The Application Developer shall have the possibility to export data in various formats. Such DaaS functionalities will be provided through the Data workflows component.   |
|                                | AD-08 and<br>DP-13 | Partial methodology support for the development of data-<br>intensive applications can be given by pre-defined data<br>workflows that are useful for application developers and data<br>publishers, e.g., assisting data owners with the process of<br>RDFizing their data. The implementation support for<br>methodology will be developed in the context of WP4<br>(DaPaaS Methodology) further elaborated in Deliverable<br>D4.1. |
| Data Cleaning & Transformation | DP-05              | The Data Cleaning & Transformation will implement capabilities for data cleaning (e.g. duplicate removal).   |

#### 3.2.3 Notification

The **Notification** service category is implemented by a **Notification Service**. It provides functionality for subscribing to events and notifications about apps and datasets.

The Notification Service will be offered through a REST API. The API will provide actions for creating, deleting and listing topics, actions for subscribing and unsubscribing to such topics, and actions for publishing messages to subscribers of specific topics.

| Component               | Addressed<br>Requirement | Description  |
|-------------------------|--------------------------|--|
| Notification<br>Service | DP-06                    | The Notification Service must implement support for publishing and subscribing different types of events and notifications for datasets. |
|                         | EU-03                    | The Notification Service must also implement support for publishing and subscribing events and notifications for users and apps.         |

#### **Table 7: Notification**

#### 3.2.4 App Management & Deployment

The **App Management & Development** service category will be implemented by components such as **Application Container** and **App Monitoring**. The Application Container provides a cloud-provisioned environment where developers can deploy and run applications. The App Monitoring component will provide functionality for monitoring the usage of services and consumption of data on the DaPaaS Platform as well as enforcing resource quotas which guarantee the "fair use" of the Platform by 3<sup>rd</sup> parties.

The App Management & Development services will be offered through a RESTAPI. The API will provide actions for setting up the application container and deploy and undeploy apps. The API will also provide capabilities for logging relevant data about running apps for monitoring and statistics purposes.



| Table 8 | : Арр | Management | & | Deployment |
|---------|-------|------------|---|------------|
|---------|-------|------------|---|------------|

| Component                | Addressed<br>Requirement | Description   |
|--------------------------|--------------------------|---|
| Application<br>Container | AD-05                    | The Application Container provides the ability to deploy and run an application.  |
|                          | EU-06                    | The Application Container should provide mechanisms such as load balancing and app scalability to ensure high availability. |
| App Monitoring           | DP-07                    | The App Monitoring provides functionality to track usage of datasets for which statistics can be collected.                 |
|                          | AD-06                    | The App Monitoring provides functionality to track usage of the apps for which useful statistics can be collected.          |

#### 3.2.5 Catalog

The **Catalog** service category will provide capabilities for application metadata management, similar to the dataset catalog in the DaaS layer. The Catalog service will be offered through a REST API. The API will provide search actions for apps and create read, update and delete (CRUD) actions for managing metadata of apps.

Table 9: Catalog

| Component       | Addressed<br>Requirement | Description   |
|-----------------|--------------------------|---|
| Apps<br>Catalog | AD-06                    | The Apps Catalog will contain relevant metadata about the apps, e.g. description and number of users. |
|                 | EU-02                    | The Apps Catalog will provide functionality for searching and exploring apps based on metadata.       |

#### 3.2.6 Administration

The **Administration** service category will provide capabilities primarily targeting the administrator of the platform, i.e., the Instance Operator. Administration actions provide extended functionality that makes it easier for administrators to manage users, access control policies, apps and datasets. This extended functionality will be implemented by the corresponding component of the Platform Layer, but the actions will not be exposed in the public APIs available to the other users of the platform (i.e., data publishers, apps developers, or end user data consumers). Instead, separate APIs are defined for the administrator.

| Component                        | Addressed<br>Requirement | Description  |
|----------------------------------|--------------------------|--|
| User Management & Access Control | IO-01 and DP-11          | The Instance Operator and Data Publisher (and implicitly Application Developer) shall have secure access (e.g. HTTPS/SSH) to the platform. |
|                                  | IO-02                    | An Administration API will provide monitoring and logging capabilities of the complete run-time app hosting environment.                   |
|                                  | IO-03                    | An Administration API will provide extended actions for<br>retrieving statistics about users and their usage of apps<br>and datasets.      |



| IO-04 | An Administration API will provide actions for managing user accounts.  |  |  |  |  |
|-------|---|--|--|--|--|
| IO-05 | An Administration API will provide actions for configuring usage policies, e.g., predefining different policies that users can select during sign-up. |  |  |  |  |

#### 3.3 Summary of Addressed Requirements

Table 11 below depicts which DaPaaS Platform requirements (introduced in Section 2) are addressed by which components at the Platform Layer (a '+' is used to indicate this relation for each requirement).

The requirements that are not addressed by the Platform Layer are marked in grey. The requirements DP-01, DP-02, DP-03, DP-04, DP-08, DP-09 and EU-05 are addressed at the Data Layer and further elaborated in Deliverable D1.1. The requirements DP12 and EU-04 are addressed at the UX Layer and further elaborated in Deliverable D3.1.

The requirements AD-08 and DP-13 focus on methodology support. In the Platform Layer architecture partial methodology support can be provided by the Data Workflows component (see Section 3.2.2). A separate methodology component (e.g. providing online guidelines and wizards) may be introduced in the architecture. This methodological aspect is out of the scope of this document and will be addressed as part of Deliverable D4.1 in WP4 (DaPaaS Methodology).

| DaPaaS Platform<br>Requirement | User<br>Management &<br>Access Control | Data Cleaning &<br>App<br>Development | Notification | App<br>Management &<br>Deployment | Catalog | Administration |
|--------------------------------|--|---------------------------------------|--------------|-----------------------------------|---------|----------------|
| IO-01                          | +                                      |                                       |              |                                   |         | +              |
| IO-02                          |  |                                       |              | +                                 |         | +              |
| IO-03                          | +                                      |                                       |              | +                                 |         | +              |
| IO-04                          | +                                      |                                       |              |                                   |         | +              |
| IO-05                          | +                                      |                                       |              |                                   |         | +              |
| IO-06                          |  |                                       |              |                                   |         | +              |
| DP-01                          |  |                                       |              |                                   |         |                |
| DP-02                          |  |                                       |              |                                   |         |                |
| DP-03                          |  |                                       |              |                                   |         |                |
| DP-04                          |  |                                       |              |                                   |         |                |
| DP-05                          |  | +                                     |              |                                   |         |                |
| DP-06                          | +                                      |                                       | +            |                                   |         |                |
| DP-07                          | +                                      |                                       |              | +                                 |         |                |
| DP-08                          |  |                                       |              |                                   |         |                |
| DP-09                          |  |                                       |              |                                   |         |                |
| DP-10                          | +                                      |                                       |              |                                   |         |                |
| DP-11                          | +                                      |                                       |              |                                   |         |                |

Table 11: Addressed requirements by components of the Platform Layer

#### Deliverable D2.1: Open PaaS requirements, design & architecture specification Dissemination level: PU

| DP-12 |   |   |   |   |   |  |
|-------|---|---|---|---|---|--|
| DP-13 |   |   |   |   |   |  |
| AD-01 | + | + |   |   |   |  |
| AD-02 |   | + |   |   |   |  |
| AD-03 |   | + |   |   |   |  |
| AD-04 |   | + |   |   |   |  |
| AD-05 |   |   |   | + |   |  |
| AD-06 | + |   |   | + | + |  |
| AD-07 |   |   |   |   |   |  |
| AD-08 |   |   |   |   |   |  |
| EU-01 | + |   |   |   |   |  |
| EU-02 |   |   |   |   | + |  |
| EU-03 | + |   | + |   |   |  |
| EU-04 |   |   |   |   |   |  |
| EU-05 |   |   |   |   |   |  |
| EU-06 |   |   |   | + |   |  |

DaPaas

### 4 Review of Relevant Technologies for the Platform Layer

This section presents the state-of-the-art analysis of relevant technologies for the Platform Layer and includes an explanation of the technology selection (Section 4.1) and review of relevant technologies for the design and implementation of the Platform Layer of the DaPaaS Platform (Section 4.2 and 4.3).

### 4.1 Technology Selection Approach

There are many PaaS vendors in the market nowadays. The website <u>http://clouds360.com</u> lists the top 20 PaaS vendors. The three dominant players in the PaaS market are Amazon with its Amazon Web Services (AWS)<sup>1</sup> solution, Google with its Google App Engine<sup>2</sup> and Microsoft with its Windows Azure<sup>3</sup>. Another main commercial vendor in the market is Salesforce.com with its Salesforce1 Platform<sup>4</sup>. Open-source alternatives also exist, e.g. AppScale<sup>5</sup> and OpenStack<sup>6</sup>.

The implementation strategy for the DaPaaS Platform is to utilize open source technologies and be based on open standards as much as possible. This does not necessarily exclude commercial components to be used in the platform (such as the OWLIM RDF database from Ontotext for the DaaS layer) as long as they are compliant with relevant open standards (such as those produced by the W3C for RDF and SPARQL, in the case of OWLIM) and can easily be replaced by alternative open source products.

Due to many available open source solutions relevant for the Platform Layer, for the development of the Platform Layer the aim is to be based entirely on open source technologies. At the very least the core component of the Platform Layer, i.e., the Data Cleaning & App Development and the App Management & Deployment, will be based on existing open source tools and frameworks. For this reason the review of relevant technologies for the Platform Layer is focused on open source solutions for PaaS and data integration capabilities:

- **PaaS capabilities:** A Platform-as-a-Service (PaaS) provides the capability for 3<sup>rd</sup> party application developers to deploy on a cloud infrastructure 3<sup>rd</sup> party applications developed using programming languages, libraries, services and protocols provided by the PaaS provider. The PaaS provides 3<sup>rd</sup> party developers control over the deployed applications and configuration settings for the application-hosting environment. Typical PaaS service offerings include capabilities for application design, development, deployment, data access, security, scalability, storage, application instrumentation, service monitoring, workflow management, discovery, etc.
- **Data integration capabilities:** The DaaS services of the Data Layer primarily mainly provide automated functions for data management (e.g. import, linking, etc.). The data integration capabilities of the Platform Layer provide additional semi-automated functions that help the Data Publisher in publishing data and the Application Developer in developing data-driven applications. Examples of such capabilities are data cleaning, transformation, workflows and methodology that often requires a graphical user interface to interact with.

Table 12 below lists the selected solutions that have been analysed with respect to PaaS and data integration capabilities for the Platform Layer design and implementation. Further details on the evaluations are provided in the following sub-sections.

<sup>&</sup>lt;sup>1</sup> <u>http://aws.amazon.com/</u>

<sup>&</sup>lt;sup>2</sup> <u>https://developers.google.com/appengine/</u>

<sup>&</sup>lt;sup>3</sup> www.windowsazure.com/

<sup>&</sup>lt;sup>4</sup> <u>http://www.salesforce.com/eu/platform/overview/</u>

<sup>&</sup>lt;sup>5</sup> <u>http://www.appscale.com/</u>

<sup>&</sup>lt;sup>6</sup> http://www.openstack.org/



|                   |                       |   | 1   |
|-------------------|-----------------------|---|---|
| Capabilities      | Solution              | Short description   | Open source<br>license  |
| PaaS<br>Solutions | Docker                | Docker is an open source solution to easily create lightweight, portable, self-sufficient containers from any application.  | Apache License<br>Version 2.0   |
|                   | Cocaine               | Cocaine is an open-source PaaS system for creating custom cloud hosting apps that are similar to Google App Engine or Heroku.   | GNU Lesser<br>General Public<br>License (GPL)<br>Version 3                          |
|                   | Deis                  | Deis is an open source PaaS that makes it<br>easy to deploy and scale Docker containers<br>and Chef nodes used to host applications,<br>databases, middleware and other services.   | Apache License<br>Version 2.0   |
|                   | Juju                  | Juju is a tool for configuring, managing,<br>maintaining and deploying applications<br>based on Charms application configurations.  | GNU Affero General<br>Public License<br>(aGPL) Version 3                            |
|                   | Cozy Cloud            | Cozy Cloud allows users to build a Personal<br>Cloud Platform-as-a-Service (PaaS) where<br>users can deploy personal web applications<br>(official apps such as Calendar, Contacts<br>and Photos are available from a<br>marketplace) or write apps from scratch.                                       | GNU Lesser<br>General Public<br>License (LGPL)<br>Version 3                         |
|                   | OpenCivic             | OpenCivic is an open source resource<br>cataloguing, hackathon and app store<br>management platform designed to help<br>organizations better collaborate in<br>developing, sharing and maintaining<br>information and apps that solve civic<br>problems.  | Available as open<br>source on GitHub,<br>but no licensing<br>information is given. |
|                   | OpenStack             | OpenStack is a cloud operating system that<br>controls large pools of compute, storage,<br>and networking resources throughout a<br>datacenter, all managed through a<br>dashboard that gives administrators control<br>while empowering their users to provision<br>resources through a Web interface. | Apache License<br>Version 2.0   |
|                   | Ansible               | Ansible is an automation engine that aims to make systems and apps simple to deploy.  | GNU General Public<br>License (GPL)<br>Version 3                                    |
|                   | Puppet Open<br>Source | Puppet Open Source is a flexible,<br>customizable framework designed to help<br>system administrators automate the many<br>repetitive tasks they regularly perform.   | Apache License<br>Version 2.0   |
|                   | Chef                  | Chef is a systems and cloud infrastructure<br>automation framework that aims to make it<br>easy to deploy servers and applications to<br>any physical, virtual, or cloud location.  | Apache License<br>Version 2.0   |
|                   | Nagios Core           | Nagios Core is an open source IT<br>monitoring system that monitors critical IT<br>infrastructure components, including system<br>metrics, network protocols, applications,   | GNU General Public<br>License (GPL)<br>Version 2                                    |

| Table 12: Overview of relevant | t open source | technologies |
|--------------------------------|---------------|--------------|
|--------------------------------|---------------|--------------|



|                                  |  | services, servers, and network infrastructure.   |  |
|----------------------------------|--|--|--|
| Data<br>Integration<br>Solutions | Talend Open<br>Studio for<br>Data<br>Integration | Talend Open Studio for Data Integration<br>provides a set of data integration tools to<br>access, transform and integrate data from<br>any business system in real time or batch to<br>meet both operational and analytical data<br>integration needs. | GNU Lesser<br>General Public<br>License (LGPL)<br>Version 3  |
|                                  | OpenRefine                                       | OpenRefine is a tool for working with messy<br>data, cleaning it, transforming it from one<br>format into another, extending it with Web<br>Services and linking to databases.   | Open source<br>Google license with<br>dependencies to a<br>number external<br>open source<br>licenses. |
|                                  | Karma  | Karma is an information integration tool that<br>enables users to quickly and easily integrate<br>data from a variety of data sources including<br>databases, spreadsheets, delimited text<br>files, XML, JSON, KML and Web APIs.                      | Apache License<br>Version 2.0  |
|                                  | Cascading  | Cascading is an application framework for<br>Java developers to develop robust data<br>analytics and data management applications<br>on Apache Hadoop.   | Apache License<br>Version 2.0  |
|                                  | Data Pipes                                       | Data Pipes is a service to provide<br>streaming, "pipe-like" data transformations<br>on the web – things like deleting rows or<br>columns, find and replace, head, grep etc.   | MIT License  |

#### 4.2 PaaS Capabilities Solutions

#### 4.2.1 Docker

Docker<sup>7</sup> is an open source project to easily create lightweight, portable, self-sufficient containers from any application. It has a very active user community<sup>8</sup> and has more than 200 contributors.

The main feature of Docker is to provide portable deployment across machines by packaging software in a common kind of container. Docker defines a format for bundling an application and all its dependencies into a single object which can be transferred to any Docker-enabled machine, and executed there with the guarantee that the execution environment exposed to the application will be the same. This allows the packaged applications to run in different environments without being reconfigured again.

Figure 9 illustrates the basic functionality of Docker<sup>9</sup>. A container comprises both an application and all of its dependencies. Containers can either be created manually or automatically in a source code repository (requires a DockerFile). Subsequent modifications to a baseline Docker image can be committed to a new container using the Commit function and then Pushed to a Central Registry. Containers can be found in a Docker Registry using Search. Containers can be pulled from the registry using Pull and can be run, started, stopped, etc. using Run commands. The target of a run command can be self-owned servers, public instances, or a combination.

<sup>&</sup>lt;sup>7</sup> <u>https://www.docker.io/</u>

<sup>&</sup>lt;sup>8</sup> http://blog.docker.io/2013/11/docker-project-community-stats/

<sup>&</sup>lt;sup>9</sup> https://www.docker.io/the\_whole\_story/#What-are-the-Main-Features-of-Docker



Figure 9: Basic Docker functions<sup>10</sup>

In addition the Docker tool offers features for application deployment, automatic build, versioning and component re-use. Docker also provides an API<sup>11</sup> for automating and customizing the creation and deployment of containers.

Docker is released under the open source Apache License Version 2.0 license<sup>12</sup>.

#### 4.2.2 Cocaine

Cocaine<sup>13</sup> (Configurable Omnipotent Custom Applications Integrated Network Engine) is an opensource PaaS system for creating custom cloud hosting apps that are similar to Google App Engine<sup>14</sup> or Heroku<sup>15</sup>, which are not open source and locked in to the Google and Amazon cloud platforms respectively. Cocaine is developed by Yandex, the leading search engine in Russia, which organizes the community and conferences.

The Cocaine architecture (see Figure 10) simplifies the creation of cloud hosting apps by hiding the infrastructure details and the applications environment settings from the developer. The developer only needs to send the code to the Cocaine server and write a special manifest for executing the code. It is not necessary to set up anything else, such as databases, as these are handled by services in the infrastructure. Any library or service can be implemented as a service in Cocaine using a special API<sup>16</sup>. From the programmer's point of view, these services look like native modules for the programming language the code is written in.

One of the notable features of Cocaine is that apps are driven by events. There are two sources of events for every app, and there exists lots of predefined plugins providing those sources. Firstly there exists services such as publish-subscribe notification and secondly there are event drivers allowing developers to generate app events.

Cocaine uses Docker as the underlying application container. Cocaine provides a plugin which connects to Docker and controls it using a rich REST API.

<sup>&</sup>lt;sup>10</sup> Figure taken from <u>https://www.docker.io/the\_whole\_story/#What-are-the-Main-Features-of-Docker</u>

<sup>&</sup>lt;sup>11</sup> http://docs.docker.io/en/latest/api/

<sup>&</sup>lt;sup>12</sup> https://github.com/dotcloud/docker/blob/master/LICENSE

<sup>&</sup>lt;sup>13</sup> http://api.yandex.com/cocaine/

<sup>&</sup>lt;sup>14</sup> https://developers.google.com/appengine/

<sup>&</sup>lt;sup>15</sup> <u>https://www.heroku.com/</u>

<sup>&</sup>lt;sup>16</sup> https://github.com/cocaine/cocaine-docs/blob/v0.11/doc/contents.md

Cocaine is released under the open source GNU Lesser General Public License (LGPL) Version 3 license<sup>17</sup>.



Figure 10: Cocaine architecture<sup>18</sup>

#### 4.2.3 Deis

Deis<sup>19</sup> is an open source PaaS that makes it easy to deploy and scale Docker containers and Chef nodes used to host applications, databases, middleware and other services. Deis leverages Chef, Docker and Heroku Buildpacks<sup>20</sup> to combine a Heroku-inspired application platform for public and private clouds. Deis is a fairly new open source project, but has high ambitions with a regular release schedule with a production ready 1.0 release expected soon<sup>21</sup>.

The main component of Deis is the controller component<sup>22</sup>. Controllers are tied to a configuration management backend where data about users, applications and formations is stored. A formation is a set of infrastructure used to host applications.

The controller is in charge of:

- Processing client API calls
- Managing nodes that provide services to a formation
- Managing containers that perform work for applications
- Managing proxies that route traffic to containers
- Managing users, providers, flavors, keys and other base configuration

<sup>&</sup>lt;sup>17</sup> https://github.com/cocaine/cocaine-core/blob/master/LICENSE

<sup>&</sup>lt;sup>18</sup> Figure taken from <u>http://api.yandex.com/cocaine/</u>

<sup>&</sup>lt;sup>19</sup> <u>http://deis.io/overview/</u>

<sup>&</sup>lt;sup>20</sup> https://devcenter.heroku.com/articles/buildpacks

<sup>&</sup>lt;sup>21</sup> <u>http://deis.io/deis-devops-and-the-future-of-open-paas/</u>

<sup>&</sup>lt;sup>22</sup> http://docs.deis.io/en/latest/components/controller/

Deis provides a REST API<sup>23</sup> for integration with other tools and systems. Deis is released as open source under the Apache License 2.0 license<sup>24</sup>.

#### 4.2.4 Juju

Juju<sup>25</sup> is a tool for configuring, managing, maintaining and deploying Charms application architectures. Charms<sup>26</sup> encapsulate application configurations, define how services are deployed, how they connect to other services and are scaled. A Charm is essentially a structure bundled of files that contain metadata, configuration data and hooks (e.g., executable files).

Juju is a project launched by Canonical, which are the developers of the Ubuntu Linux-based operating system, and has a strong community behind it. Juju allows users to deploy, manage and scale software and interconnected services across one or more Ubuntu servers and cloud platforms.

Juju provides a GUI (see Figure 11) and command-line interface to define, configure, deploy, manage, monitor and scale services to any public or private cloud.

Juju is released under the open source GNU Affero General Public License (aGPL) Version 3 license<sup>27</sup>.



Figure 11: Juju administration GUI<sup>28</sup>

#### 4.2.5 Cozy Cloud

Cozy Cloud<sup>29</sup> allows users to build a Personal Cloud Platform-as-a-Service (PaaS) where users can deploy personal web applications (official apps such as Calendar, Contacts and Photos are available from a marketplace) or write apps from scratch. Cozy Cloud provides a centralized storage with data types and access control that can be shared/used by the apps.

Cozy Cloud is developed by a young startup company located in France with a small team of eight persons. Cozy Cloud was developed to be a private **personal** cloud solution that allows users to host their own personal applications in a single place that they control. This way, users can manage their data from anywhere while protecting their privacy. The design of the platform is very focused on the

<sup>&</sup>lt;sup>23</sup> <u>http://docs.deis.io/en/latest/server/</u>

<sup>&</sup>lt;sup>24</sup> https://github.com/opdemand/deis/blob/master/LICENSE

<sup>&</sup>lt;sup>25</sup> https://juju.ubuntu.com/

<sup>&</sup>lt;sup>26</sup> <u>https://juju.ubuntu.com/charms/</u>

<sup>&</sup>lt;sup>27</sup> https://launchpad.net/juju-core

<sup>&</sup>lt;sup>28</sup> Figure taken from <u>https://juju.ubuntu.com/features/</u>

<sup>&</sup>lt;sup>29</sup> https://www.cozycloud.cc/

personal single-user perspective. Thus capabilities for handling large data and multi-users have not been addressed in the architecture. Because of the focus on the single-user perspective, the Cozy Cloud solution does not seem suitable for the implementation of the Platform Layer in DaPaaS.

Cozy Cloud is released as open source under the GNU Lesser Generic Public License (LGPL) Version 3 license<sup>30</sup>.

#### 4.2.6 OpenCivic

OpenCivic<sup>31</sup> is an open source resource cataloguing, hackathon and app store management platform designed to help organizations better collaborate in developing, sharing and maintaining information and apps that solve civic problems. OpenCivic is based on Drupal. The main goal of the Drupal distro is to help build websites that enable people to share information about software applications.

OpenCivic provides features for app store and hackathons:

- Catalog applications and their metadata, including pictures and their descriptions
- Catalog organizations that use and contribute to apps
- Store and publish open data for developers to use in application development
- Catalog deployments of specific apps by organizations and locations
- Publish and manage hackathons events
- Collect, define, rate and refine problems for the events
- Define and manage development teams for the events

The OpenCivic distribution for Drupal is available as open source on GitHub<sup>32</sup>, but no licensing information is given.

#### 4.2.7 OpenStack

OpenStack<sup>33</sup> is a cloud operating system that controls large pools of compute, storage, and networking resources throughout a datacenter, all managed through a dashboard that gives administrators control while empowering their users to provision resources through a web interface. OpenStack has a very large community and top industry involvement behind it, organized as the OpenStack Foundation<sup>34</sup>.

Figure 12 shows a conceptual architecture of the OpenStack from the operator side of the cloud, i.e. the Instance Operator. The OpenStack provides APIs for:

- Dashboard (Web frontend)
- Store and retrieval of virtual disks/images and associated metadata, as well as the virtual disk files themselves
- Virtual network and storage volumes for computing
- Identity and service authentication

<sup>&</sup>lt;sup>30</sup> <u>https://github.com/mycozycloud/cozy-setup/blob/master/LICENSE</u>

<sup>&</sup>lt;sup>31</sup> <u>http://nucivic.com/opencivic/</u>

<sup>&</sup>lt;sup>32</sup> <u>https://github.com/civic-commons/opencivic</u>

<sup>&</sup>lt;sup>33</sup> http://openstack.org

<sup>&</sup>lt;sup>34</sup> http://www.openstack.org/community/



#### Figure 12: OpenStack conceptual architecture<sup>35</sup>

The OpenStack software supports allocating a large amount of servers to provide resources for computation. These resources can then be consumed in a uniform manner through the OpenStack abstraction layer<sup>36</sup>.

The OpenStack project is provided<sup>37</sup> under the open source Apache License 2.0 license<sup>38</sup>.

#### 4.2.8 Ansible

Ansible<sup>39</sup> is a powerful automation engine that makes systems and apps simple to deploy. Ansible is an IT automation tool. It can configure systems, deploy software, and orchestrate more advanced IT tasks such as continuous deployments or zero downtime rolling updates. There is a large community around the tool<sup>40</sup>.

Ansible provides features such as<sup>41</sup>:

- Command Line Tools
- Application Deployment
- Continuous Delivery
- Multi-Tier Orchestration
- Configuration Management
- Agentless Architecture

<sup>&</sup>lt;sup>35</sup> Figure taken from <u>http://docs.openstack.org/grizzly/openstack-compute/admin/content//conceptual-architecture.html</u>

<sup>&</sup>lt;sup>36</sup> <u>http://api.openstack.org/api-ref.html</u>

<sup>&</sup>lt;sup>37</sup> https://wiki.openstack.org/wiki/Getting\_The\_Code

<sup>&</sup>lt;sup>38</sup> <u>https://wiki.openstack.org/wiki/Open</u>

<sup>&</sup>lt;sup>39</sup> http://www.ansibleworks.com/

<sup>&</sup>lt;sup>40</sup> <u>http://www.ansibleworks.com/community/</u>

<sup>&</sup>lt;sup>41</sup> <u>http://www.ansibleworks.com/pricing/</u>



• SSH-Based Security

The tool provides a Python API<sup>42</sup>. Modules and plugins can be developed in any programming language<sup>43</sup>.

Ansible is available as an open source engine<sup>44</sup> released under the GNU General Public License (GPL) Version 3 license<sup>45</sup>.

#### 4.2.9 Puppet Open Source

Puppet Open Source<sup>46</sup> is a flexible, customizable framework available under the Apache 2.0 license designed to help system administrators automate the many repetitive tasks they regularly perform. As such it is similar to Ansible, namely an IT automation tool. It has a large and active community<sup>47</sup>.

Puppet uses a declarative, model-based approach to IT automation<sup>48</sup>.

- 1. Define the desired state of the infrastructure's configuration using Puppet's declarative configuration language.
- 2. Simulate configuration changes before enforcing them.
- 3. Enforce the deployed desired state automatically, correcting any configuration drift.
- 4. Report on the differences between actual and desired states and any changes made enforcing the desired state.



- <sup>43</sup> <u>http://docs.ansible.com/developing.html</u>
- <sup>44</sup> <u>https://github.com/ansible/ansible</u>
- <sup>45</sup> http://www.ansibleworks.com/opensource/
- <sup>46</sup> <u>http://puppetlabs.com/puppet/puppet-open-source</u>
- <sup>47</sup> http://puppetlabs.com/community/overview
- <sup>48</sup> <u>http://puppetlabs.com/puppet/what-is-puppet</u>
- <sup>49</sup> Figure taken from <u>http://puppetlabs.com/puppet/what-is-puppet</u>

<sup>42</sup> http://docs.ansible.com/developing\_api.html



#### 4.2.10 Chef

Chef<sup>50</sup> is a systems and cloud infrastructure automation framework that makes it easy to deploy servers and applications to any physical, virtual, or cloud location, no matter the size of the infrastructure. It is an IT automation and configuration management tool similar to Ansible and Puppet. It has a large and active community<sup>51</sup>.



Figure 14: How Chef works<sup>52</sup>

Chef relies on reusable definitions known as cookbooks and recipes (see Figure 14) that are written using the Ruby programming language. Cookbooks and recipes automate common infrastructure tasks. Their definitions describe what your infrastructure consists of and how each part of your infrastructure should be deployed, configured and managed. Chef applies those definitions to servers to produce an automated infrastructure. The Chef server stores your network's configuration data and recipes.

Chef is released under the Apache License Version 2.0 license<sup>53</sup>.

<sup>&</sup>lt;sup>50</sup> <u>http://docs.opscode.com/chef\_overview.html</u>

<sup>&</sup>lt;sup>51</sup> http://community.opscode.com/

<sup>&</sup>lt;sup>52</sup> Figure taken from <u>http://www.getchef.com/chef/</u>

<sup>53</sup> https://github.com/opscode/chef/blob/master/LICENSE



#### 4.2.11 Nagios Core

Nagios Core<sup>54</sup> is an open source host, service and network monitoring tool. It allows monitoring entire IT infrastructures to ensure systems, applications, services, and business processes are functioning properly. In the event of a failure, it can alert technical staff of the problem, allowing them to begin remediation processes before outages affect business processes, end users, or customers. It is capable to manage different types of services and hosts running on different operating systems such as Linux, Netware, Windows, AIX, etc. It's flexible in configuration and can be extended as much as it is necessary. It's configured within text files and managed with a Web browser.



Figure 15: Operating principle of Nagios<sup>55</sup>

Features overview:

- Monitoring of network services (SMTP, POP3, HTTP, NNTP, PING, etc.)
- Monitoring of host resources (processor load, disk usage, etc.)
- Simple plugin design that allows users to easily develop their own service checks
- Parallelized service checks
- Ability to define network host hierarchy using "parent" hosts, allowing detection of and distinction between hosts that are down and those that are unreachable
- Contact notifications when service or host problems occur and get resolved (via email, pager, or user-defined method)
- Ability to define event handlers to be run during service or host events for proactive problem resolution
- Automatic log file rotation
- Support for implementing redundant monitoring hosts
- Optional web interface for viewing current network status, notification and problem history, log file, etc.

<sup>&</sup>lt;sup>54</sup> <u>http://www.nagios.com/products/nagioscore/</u>

<sup>&</sup>lt;sup>55</sup> Figure taken from <u>http://en.wikipedia.org/wiki/Nagios</u>

Nagios provides a support forum<sup>56</sup> and an exchange site<sup>57</sup> for the community. Nagios Core is released as open source under the GNU General Public License (GPL) Version 2 license<sup>58</sup>.

#### 4.3 Data Integration Capabilities Solutions

#### 4.3.1 Talend Open Studio for Data Integration

Talend Open Studio for Data Integration<sup>59, 60</sup> provides a set of data integration tools to access, transform and integrate data from any business system in real time or batch to meet both operational and analytical data integration needs.

The product provides a graphical business modeller tool for designing business logic for data-intensive applications with data flow sequencing using components and connectors (see Figure 16). The product provides 450+ native database and storage connectivity components<sup>61</sup> that allow users to connect to almost any data source, including databases and data warehouses such as such as Amazon S3, DB2, Ingres, JDBC, Microsoft SQL Server, MySQL, Oracle and PostgreSQL.

There is an active community around the tool and the Talend teams are willing to answer any question on their forum<sup>62</sup>. The product is released as open source <sup>63</sup> under the GNU Lesser General Public License (LGPL) Version 3 license<sup>64</sup>.



Figure 16: Talend Open Studio for Data Integration

<sup>&</sup>lt;sup>56</sup> <u>http://support.nagios.com/forum/</u>

<sup>&</sup>lt;sup>57</sup> http://exchange.nagios.org/

<sup>58</sup> http://sourceforge.net/p/nagios/nagioscore/ci/master/tree/LICENSE

<sup>59</sup> http://www.talend.com/products/data-integration

<sup>60</sup> http://www.talend.com/download/data-integration

<sup>&</sup>lt;sup>61</sup> http://www.talend.com/products/specifications-data-integration

<sup>62</sup> http://www.talendforge.org/forum/

<sup>63</sup> http://www.talendforge.org/trac/tos/

<sup>&</sup>lt;sup>64</sup> http://www.gnu.org/licenses/lgpl.html



#### 4.3.2 OpenRefine

OpenRefine<sup>65</sup> is a tool for working with messy data, cleaning it, transforming it from one format into another, extending it with Web Services and linking to databases. The main features of the tool are:<sup>66</sup>

- Importing (formats TSV, CSV, Excel, XML, RDF as XML, JSON, Google Spreadsheets and RDF N3 triples)
- Filtering and faceting (exploring data by applying multiple filters)
- Editing cells, columns and rows (see Figure 17)
- Editing with Google Refine Regular Expression Language (GREL)<sup>67</sup>
- Exporting (formats TSV, CSV, Excel, HTML tables and JSON)
- History (undo and redo)
- Reconciliation against FreeBase<sup>68</sup> schemas using schema alignment dialogs
- Extending data calling Web Services

| Owner Category |                                      |                         | <b>District</b>            | Province   |                           |            |  |
|----------------|--------------------------------------|-------------------------|----------------------------|--|---------------------------|------------|--|
| Facet          | •                                    | pital                   | Chegutu District           | MASH   | ONALAND WEST PRO∨INCE     |            |  |
| Text filter    |                                      | tal                     | Chegutu District           | MASH   | ONALAND WEST PROVINCE     |            |  |
| Edit cells     | •                                    | Tra                     | Cheautu District<br>NSform | Cheanth District MASHONALAND WEST PROVI<br>Isform DNALAND WEST PROVI |                           |            |  |
| Edit colur     | mn 🕨 🕨                               | Co                      | mmon transforms            | •  | Trim leading and trailing | whitespace |  |
| Transpos       | Transpose Fill down                  |                         |                            |  | Collapse consecutive wh   | hitespace  |  |
| Sort           | Sort Bla                             |                         | ank down                   |  | Unescape HTML entities    |            |  |
| View           | View > Spl                           |                         | it multi-valued cells      |  | To titlecase              |            |  |
| Reconcil       | e 🕨                                  | Joi                     | n multi-valued cells       |  | To uppercase              |            |  |
| Govt.<br>Govt. | RHC                                  | Clu                     | ister and edit             |  | To lowercase              |            |  |
| Govt.          | Clinic                               |                         | Hururungwe District        | MASH   | To number                 |            |  |
| Govt.          | RHC                                  |                         | Hururungwe District        | MASH   | <sup>1</sup> To date      |            |  |
| Govt.          | RHC                                  | Hururungwe District MA  |                            | MASH   | <sup>1</sup> To text      |            |  |
| Govt.          | Clinic                               | Hururungwe District MAS |                            | MASH   | 10 toxe                   |            |  |
| Govt.          | Clinic                               |                         | Hururungwe District        | MASH   | SHI Blank out cells       |            |  |
| Govt           | Sovt Clinic Hururungwe District MASH |                         |                            | MASH   | ONAL AND WEST PROVINCE    |            |  |

Figure 17: Edit cells in OpenRefine – Common transformations<sup>69</sup>

OpenRefine was originally developed by Google and thus named Google Refine. In July 2010 the product was renamed OpenRefine as it transitioned to a community-supported product<sup>70</sup>. OpenRefine is licensed under an open source Google license<sup>71</sup> with dependencies to a number external open source licenses.

There are several extensions for OpenRefine. Among them, the most relevant for this project would be:

<sup>65</sup> http://openrefine.org/

<sup>&</sup>lt;sup>66</sup> https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-Users

<sup>&</sup>lt;sup>67</sup> https://github.com/OpenRefine/OpenRefine/wiki/Google-refine-expression-language

<sup>68</sup> http://www.freebase.com/

<sup>&</sup>lt;sup>69</sup> Figure taken from <u>http://schoolofdata.org/handbook/recipes/cleaning-data-with-refine/</u>

<sup>&</sup>lt;sup>70</sup> http://openrefine.org/community.html

<sup>&</sup>lt;sup>71</sup> https://github.com/OpenRefine/OpenRefine/blob/master/LICENSE.txt



- CKAN storage extension<sup>72</sup> which allows RDF data to be registered and uploaded on CKAN storage.
- RDF Refine<sup>73</sup> which is used to reconcile and link data (against SPARQL endpoints, or search for related RDF datasets) or export data as RDF format.

#### 4.3.3 Karma

Karma<sup>74</sup> is an information integration tool that enables users to quickly and easily integrate data from a variety of data sources including databases, spreadsheets, delimited text files, XML, JSON, KML and Web APIs. Users integrate information by modelling it according to an ontology of their choice using a graphical user interface (see Figure 18) that automates much of the process.

Karma learns to recognize the mapping of data to ontology classes and then uses the ontology to propose a model that ties together these classes. Once the model is complete, users can published the integrated data as RDF or store it in a database.

Karma is a research prototype developed by the Information Sciences Institute at the University of Southern California (USC)<sup>75</sup> and does not seem to be used by a large community outside of the research group at USC.

| Karma v1.90               |                     |                 |                                  |               |   |            |   |  |
|---------------------------|---------------------|-----------------|----------------------------------|---------------|---|------------|---|--|
| Import Database Table     | Import from Service | + Import File   |                                  |               |   |            | Reset   |  |
| Command                   | History             | - artworks      | -list.xm                         | I             |   |            |   |  |
| Import XML File: artworks | list.xml imported   | accessionNumber | artist                           | birthDeath    | creditLine                                | dimensions | imageURL  |  |
|                           |                     |                 |                                  |               |   |            |   |  |
|                           |                     | 26.1            | Frishmuth,<br>Harriet<br>Whitney | 1880-<br>1980 | Gift of the<br>Friends of<br>American Art | H: 61 in.  | http://www.imamuseum.org/sites/default/files/mercuryc |  |

Karma is available as open source<sup>76</sup> under the Apache License 2.0 license.

Figure 18: Modelling data in Karma<sup>77</sup>

#### 4.3.4 Cascading

Cascading<sup>78</sup> is an application framework for Java developers to simply develop robust data analytics and data management applications on Apache Hadoop<sup>79</sup>. Apache Hadoop is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage.

Cascading provides APIs (see Figure 19) for:

- Data processing, enabling complex data flows and data-oriented frameworks.
- Data integration, enabling creation and testing of integration points.
- Process scheduler, allowing scheduling units of work from 3<sup>rd</sup> party applications.

<sup>&</sup>lt;sup>72</sup> http://lab.linkeddata.deri.ie/2011/grefine-ckan/

<sup>73</sup> http://refine.deri.ie/

<sup>&</sup>lt;sup>74</sup> http://www.isi.edu/integration/karma/

<sup>&</sup>lt;sup>75</sup> http://www.isi.edu/integration/

<sup>&</sup>lt;sup>76</sup> https://github.com/InformationIntegrationGroup/Web-Karma

<sup>&</sup>lt;sup>77</sup> Figure taken from <a href="https://github.com/InformationIntegrationGroup/Web-Karma/wiki/Modeling-Data">https://github.com/InformationIntegrationGroup/Web-Karma/wiki/Modeling-Data</a>

<sup>78</sup> http://www.cascading.org/

<sup>79</sup> http://hadoop.apache.org/



Community support is provided through a mailinglist, IRC and GitHub<sup>80</sup>. Cascading is released as open source under the Apache License 2.0 license<sup>81</sup>.



Figure 19: Cascading architecture<sup>82</sup>

#### 4.3.5 Data Pipes

Data Pipes<sup>83</sup> is a service to provide streaming, "pipe-like" data transformations on the web – things like deleting rows or columns, find and replace, head, grep, etc. It provides a REST API that supports the following data transformation operations:

- none (aka raw) = no transform but file parsed
- csv = parse / render csv
- head = take only first X rows
- tail = take only last X rows
- delete = delete rows
- strip = delete all blank rows
- grep = filter rows based on pattern matching
- cut = select / delete columns
- replace = find and replace (not yet implemented)
- html = render as viewable HTML table

Data Pipes seems like a small tool project developed and used by the Open Knowledge Foundation Labs (OKFL). Data Pipes is released as open source under the MIT license<sup>84</sup>.

<sup>&</sup>lt;sup>80</sup> http://www.cascading.org/support/

<sup>&</sup>lt;sup>81</sup> <u>https://github.com/Cascading/cascading/blob/2.5/LICENSE.txt</u>

<sup>82</sup> Figure taken from http://www.cascading.org/about-cascading/

<sup>83</sup> http://datapipes.okfnlabs.org/

<sup>&</sup>lt;sup>84</sup> https://github.com/okfn/datapipes/blob/master/LICENSE.md



### 5 Summary and Outlook

This document provided an overview of the DaPaaS Platform, introduced the relevant roles played in the DaPaaS context and outlined a set of requirements for the DaPaaS Platform from the perspectives of the key roles. Furthermore, the document focused on the PaaS aspect of the platform and provided an initial architecture design for the Platform Layer of the DaPaaS Platform. A state-of-the-art overview of relevant solutions and technologies for the Platform Layer has been presented.

Following the technology evaluation performed as part of this document, the following remarks are to be considered for the implementation phase of the DaPaaS Platform, with a particular focus on the Platform Layer:

- Docker is a promising solution to be reused as an application packaging system for DaPaaS. It has a very active user community and looks like a very promising application container technology which is used by other open source PaaS solutions such as Cocaine and Deis.
- Docker together with either Cocaine or Deis are promising solutions for the implementation of the deployment and hosting environment of the Platform Layer. Cocaine and Deis provides the ability to deploy and run the applications packaged by Docker. Further investigation is needed, i.e., setting up a test infrastructure and experimenting with writing and deploying sample applications, to decide which provides the better foundation. There are differences in the capabilities offered by the two solutions, e.g., Deis provides support for user accounts and SSH public key for authentication and Cocaine provides a publish-subscribe notification service in development that one can build on.
- Talend Open Studio for Data Integration and OpenRefine for data cleaning and integration are promising frameworks for addressing the data integration aspect of the DaPaaS Platform Layer. However further investigation is needed to see if the tools should be fully integrated, whether only a useful subset of the functionality are exposed through DaPaaS-developed frontends, or whether they should be used as inspiration for relevant aspects of the Platform Layer implementation.
- Ansible, Puppet and Chef are IT automation tools that can be used to administrate and manage a set of nodes from a single place, typically the Instance Operator. All three have strong communities and the choice for one of them may depend on the technology choice for the PaaS solution. For example, the Deis PaaS solution already leverages Chef.
- Nagios Core is a monitoring engine that, with necessary adaptations, is relevant to and may be used by the Instance Operator to monitor cloud infrastructure resources where the DaPaaS platform will be deployed.
- DaPaaS-specific extensions to the user management and access control will be implemented on top of the PaaS solution that is chosen, extending it with user profile, and dataset and app access control. Application monitoring in DaPaaS will be coupled to user quotas and policies.
- The features offered by OpenCivic should be considered for the implementation of Catalog features of the Platform Layer. However, since it is based on Drupal which is an open source content management system tailored for creating websites and not PaaS, the features of OpenCivic will likely only be an inspiration for the implementation of the Catalog in Platform Layer.
- For the DaPaaS-specific data workflows services, Cascading is probably too generic for the needs of DaPaaS. It is unclear if we can easily benefit from it as compared to OpenRefine, Talend or even Karma, and further investigation is needed to decide on what solutions can be used for the data workflows.

The above remarks will be taken into consideration for the implementation of the first prototype of the DaPaaS platform, due at M12. In parallel, we will monitor the development of the integrated DaaS/PaaS commercial / closed source solutions, in particular Windows Azure Marketplace, Datameer, and Splunk (briefly introduced in Appendix A).

### 6 Appendix A: Commercial / Closed Source Integrated DaaS & PaaS Solutions

In this section we briefly introduce some relevant commercial / closed source offerings that fit within the as-a-service platform for load/store/analysis/visualization/publication of data and development/hosting of 3rd party apps. Whereas DaPaaS targets an open source environment for data publishing/hosting and apps development/hosting, the commercial / closed source offerings briefly presented here cover functionalities that are relevant to the overall requirements outlined at the beginning of this document.

#### 6.1 Datameer

Datameer (<u>http://www.datameer.com</u>) was founded in 2009 by some of the original contributors to Apache Hadoop, which is an open source software framework for storage and large-scale processing of datasets. Datameer provides analytics software (see <u>http://www.datameer.com/product/</u>) on top of the Hadoop framework that allows to integrate, analyze and visualize data. The software application comes in three versions: 1) use on your desktop computer, 2) install on your infrastructure as a server, 3) use a provided enterprise instance with your data.

Datameer is an analytics application natively built on Hadoop, focused on leveraging the linear scalability and flexibility of Hadoop for data analytics. Its product is focused on data integration (structured such as relational data, tabular, etc, as well as unstructured data such as social data, email, log files, etc.), multidevice drag-and-drop-style data visualization, and data management (data import, export, data links, storage, data partitioning, compression, etc.).

Datameer is built on open standards/technologies, with a focus on Hadoop and HTML5. It offers extensibility features in the form of plugins. Datameer includes an SDK for writing custom plugins for import, export, functions, and visualizations (documentation is available at <a href="http://documentation.datameer.com/documentation/">http://documentation.datameer.com/documentation/</a>). Datameer's components are exposed through REST APIs.

Datameer provides an applications market at <u>http://www.datameer.com/apps</u>. At time of this writing the application store consists of 48 apps.

#### 6.2 Splunk

Splunk (<u>http://www.splunk.com</u>) shipped its first software in 2006 designed to manage unstructured data generated by machines (websites, applications, servers, networks, mobile devices, sensors and RFID assets). The Splunk Enterprise software product provides features for collecting and indexing any machine data, including the capability to handle massive live datastreams, statistical analysis and real-time dashboards. The software can also be provides as a Software-as-a-Service in the Cloud. This version of the product is named Splunk Cloud. The company also provides a software product named Hunk for analyzing and visualizing data in Hadoop.

Splunk was initially developed to allow organizations to search and analyze data generated by applications, servers and network devices in IT infrastructures. It includes capabilities for developing applications and provides an applications store at <u>http://apps.splunk.com/</u>. At time of this writing the application store consists of 488 apps.

#### 6.3 Windows Azure Marketplace

Windows Azure Marketplace (<u>http://datamarket.azure.com/</u>), supported by Microsoft, is an "online market for customers and partners to share, buy, and sell finished Software-as-a-Service applications, building block components and premium datasets." The marketplace comes with a set of APIs targeted at developers to work with datasets in the same way on many different platforms, enabling the developers to develop applications for desktop, Web, mobile, and other clients. The marketplace also offers developers common security, billing, auditing, and authentication mechanisms. At time of this

writing the application store consists of 644 apps. Further info about the Windows Azure Marketplace can be found in Deliverable D1.1.

#### 6.4 GoodData

GoodData (<u>http://www.gooddata.com/</u>) offers a business intelligence platform. It provides support for extract, transform and load (ETL) processes, connectors to various data sources, scalable data storage, a data analytics engne, reporting and various forms of data visualization.

GoodData was primarily designed for the cloud and is built on technologies such as Vertica, MongoDB, Cassandra, NetApp, and Rackspace.

GoodData is an open platform in the sense that it can be embedded in existing applications, or allows applications to be built on top of the platform, through a set of RESTful APIs.

#### 6.5 Tableau Software

Tableau software (<u>http://www.tableausoftware.com</u>) provides a set of interactive data visualization products focused on business intelligence, ranging from desktop to hosted solutions. The focus of Tableau is on providing usability and ease of use for its software for business intelligence with a strong emphasis on the data visualization capabilities. Tableau offers four main products: Tableau Desktop, Tableau Server, Tableau Online, and Tableau Public.

Tableau Desktop is a drag&drop-style data analysis and visualization desktop tool. It provides data connection capabilities, data visualization, and creation of interactive dashboards.

Tableau Server is a business intelligence application that provides browser-based analytics. It comes with a Web-based dashboard where users can import and integrate data, and analyse and visualize it.

Tableau Online is a hosted version of Tableau Server. It provides a central place to manage data, data sources and metadata, and focuses on scalability.

Tableau Public is delivered as a service and aims at creating interactive visuals and publishing them without the help of programmers or IT. It targets organizations that want to enhance their websites with interactive data visualizations. It offers carious visualization types, such as maps, bar and line charts, lists, heat maps, etc.

#### 6.6 Infochimps

Infochimps (<u>http://www.infochimps.com/</u>) provides a managed cloud service called Infochimps Cloud (<u>http://www.infochimps.com/infochimps-cloud/overview/</u>) that streamlines building and managing complex Big Data environments, and aims to make it faster and less complex to develop and deploy Big Data applications. The platform provides capabilities for data streaming, storage, queries and administration.