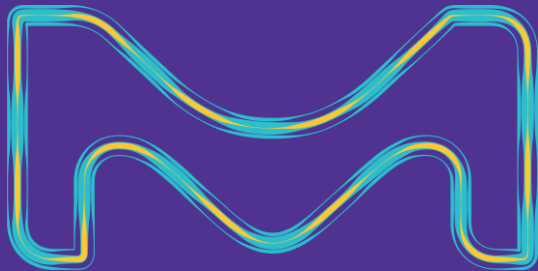


# x-OMics platform

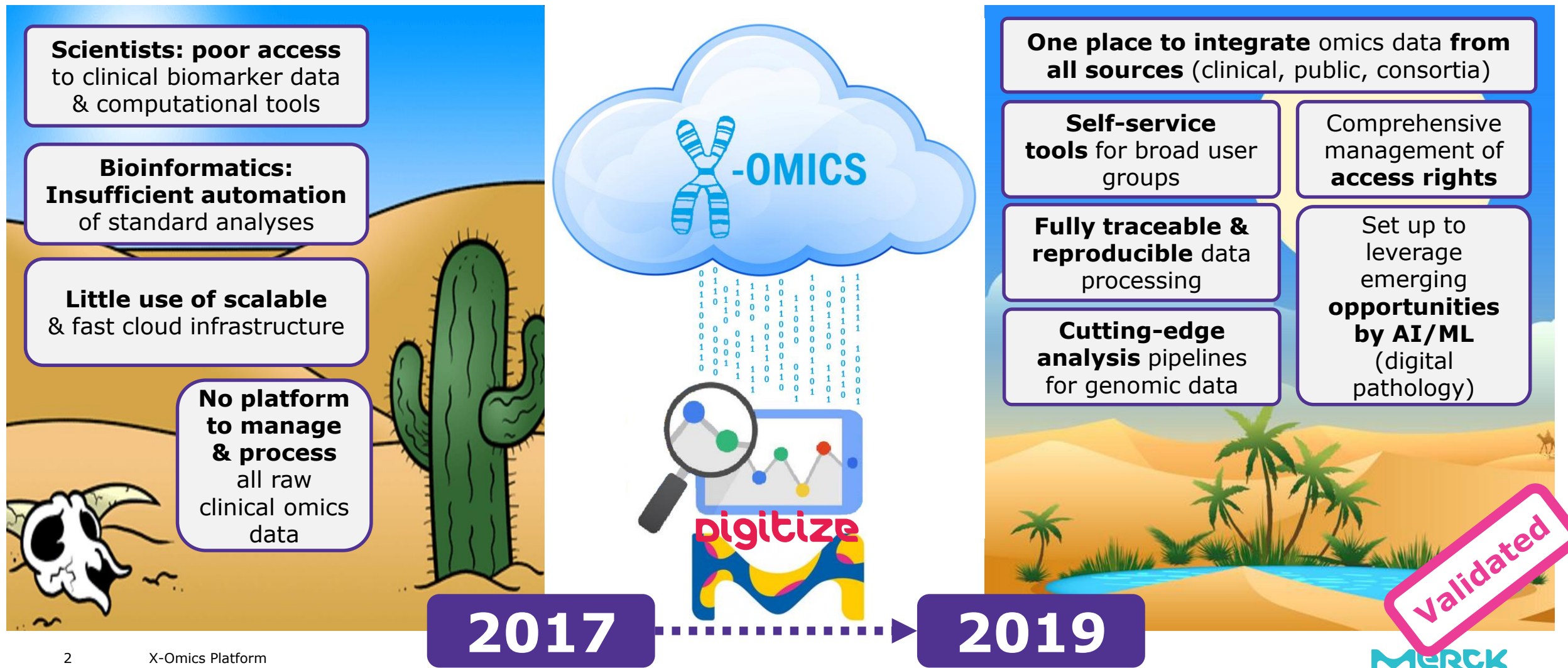
A validated platform for multi-omics data analysis

Stefan Pinkert on behalf of the XOP team  
R/Pharma 2020, 15<sup>th</sup> OCT 2020 1:40PM



MERCK

# X-OMICS Platform to Boost Computational Capabilities in Biomarker R&D



# X-OMICS Platform successfully co-developed with Genedata



Joern Peter Halle  
Head of Immuno-oncology  
and External Innovation  
Merck KGaA

Press Release



**Working closely with Genedata, we have generated a digital biomarker research platform that will enable Merck scientists to better utilize clinical and translational data to generate innovative ideas for new biomarkers and drug targets.**

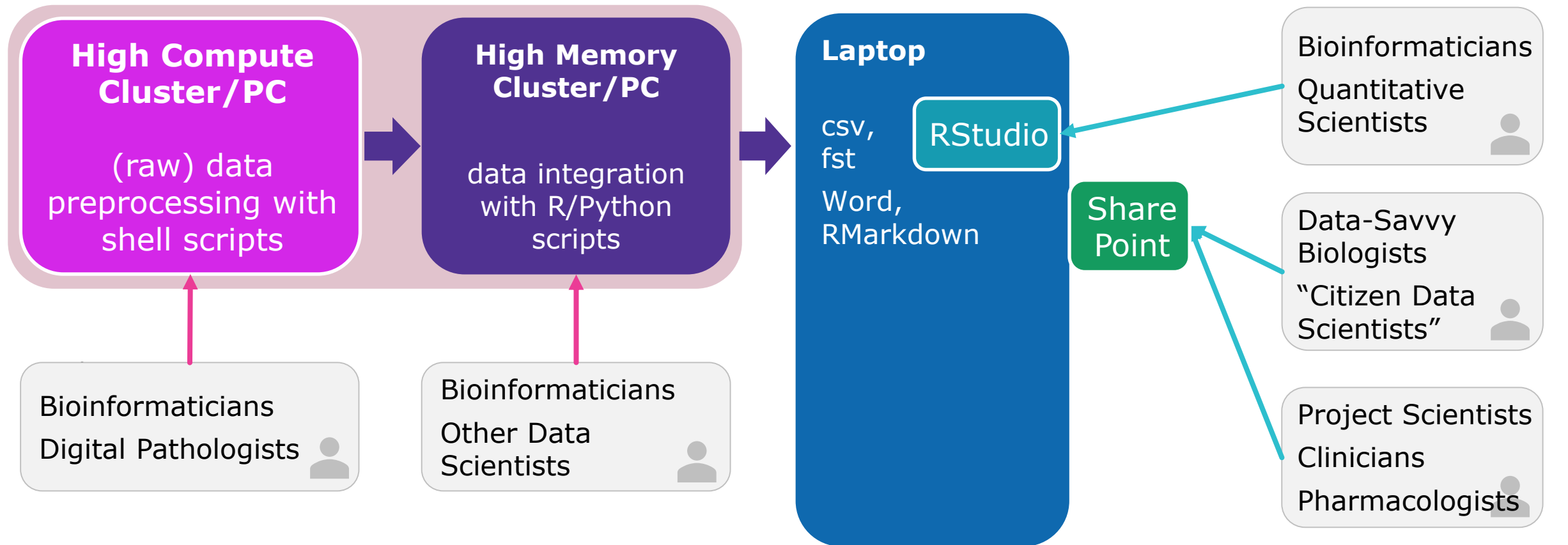


Source: <https://www.genedata.com/products/profiler>

# exploratory research with R

01

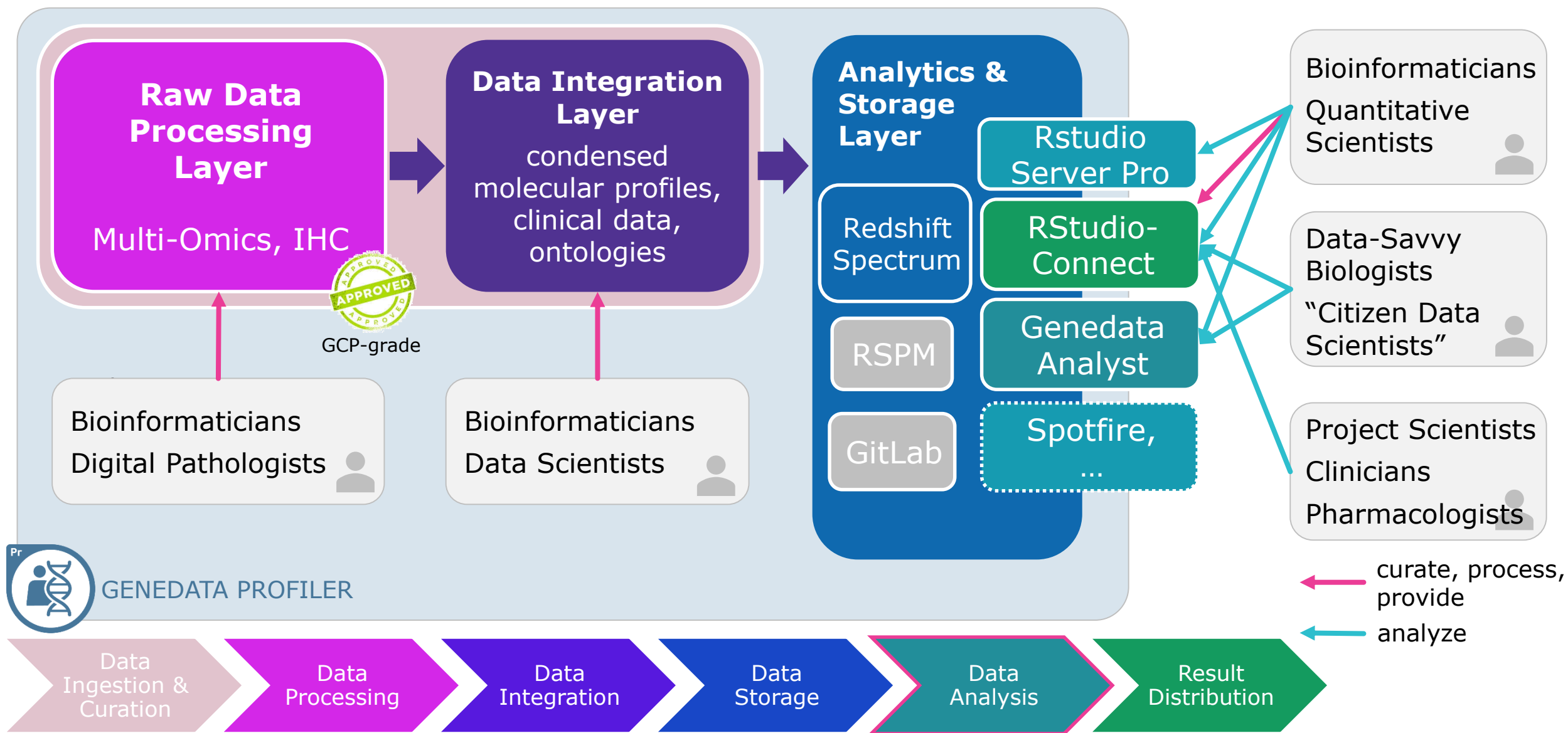
# How do exploratory researchers work?



← curate, process, provide  
← analyze



# How do researchers work with the X-OMICS Platform!

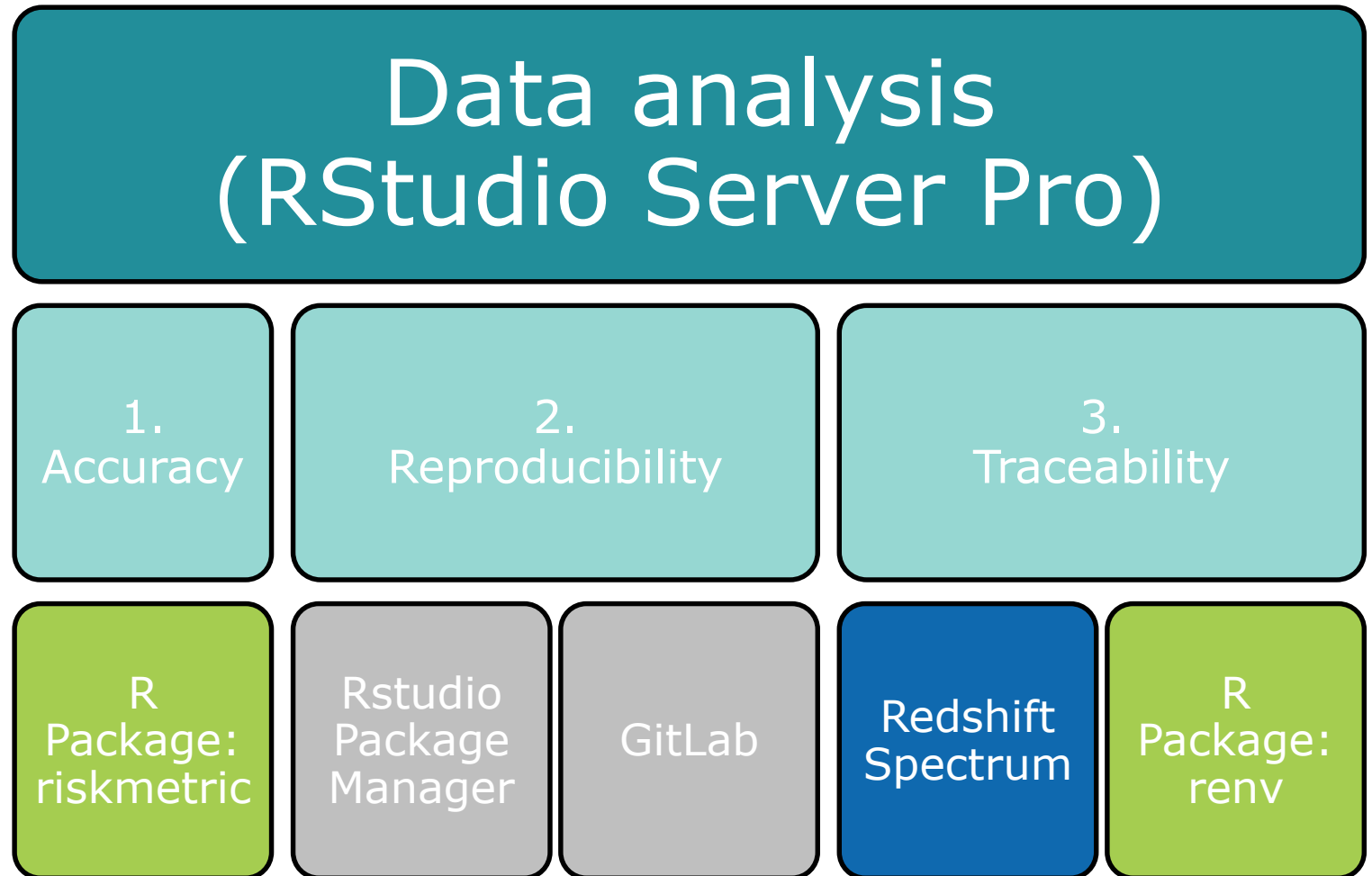


## Which elements are required for a validated system?

Since the [FDA does not require use of any specific software for statistical analyses](#), the programming language R can be used if the R installation incorporates all of the following elements:

1. Accuracy
2. Reproducibility
3. Traceability

Source:  
[https://www.pharmar.org/presentations/r\\_packages-white\\_paper.pdf](https://www.pharmar.org/presentations/r_packages-white_paper.pdf)



02

The xop supports  
reproducibility



## Reproducibility

### X-OMICS Platform uses RSPM and GitLab to support Reproducibility

It is important to acknowledge that R (like other open source languages) presents additional challenges with respect to the reproducibility of an environment. The **evolution of R packages is effectively continuous** and thus maintaining a stable R installation that allows for the addition of new and/or updated packages can be a challenge.

Source:

[https://www.pharmar.org/presentations/r\\_packages-white\\_paper.pdf](https://www.pharmar.org/presentations/r_packages-white_paper.pdf)

#### RStudio Package Manager

- Local copies of public packages
- Internal repositories / packages
- **Date based checkpoints**
- Hierarchical order of repositories

#### GitLab

- Analysis / package version control
- Role based code review & release
- Documentation
- Collaborative guiding documents

## The preconfigured X-OMICS platform supports reproducible workflows



03

The xop supports  
traceability

# Traceability

## X-OMICS Platform uses **renv** and **Redshift** to support Traceability

### Traceability of R installations

Develop system and process controls to automatically **document** the R **packages** and **installation dependencies** that are used in R analyses.

Source:

<https://www.pharmar.org/overview>

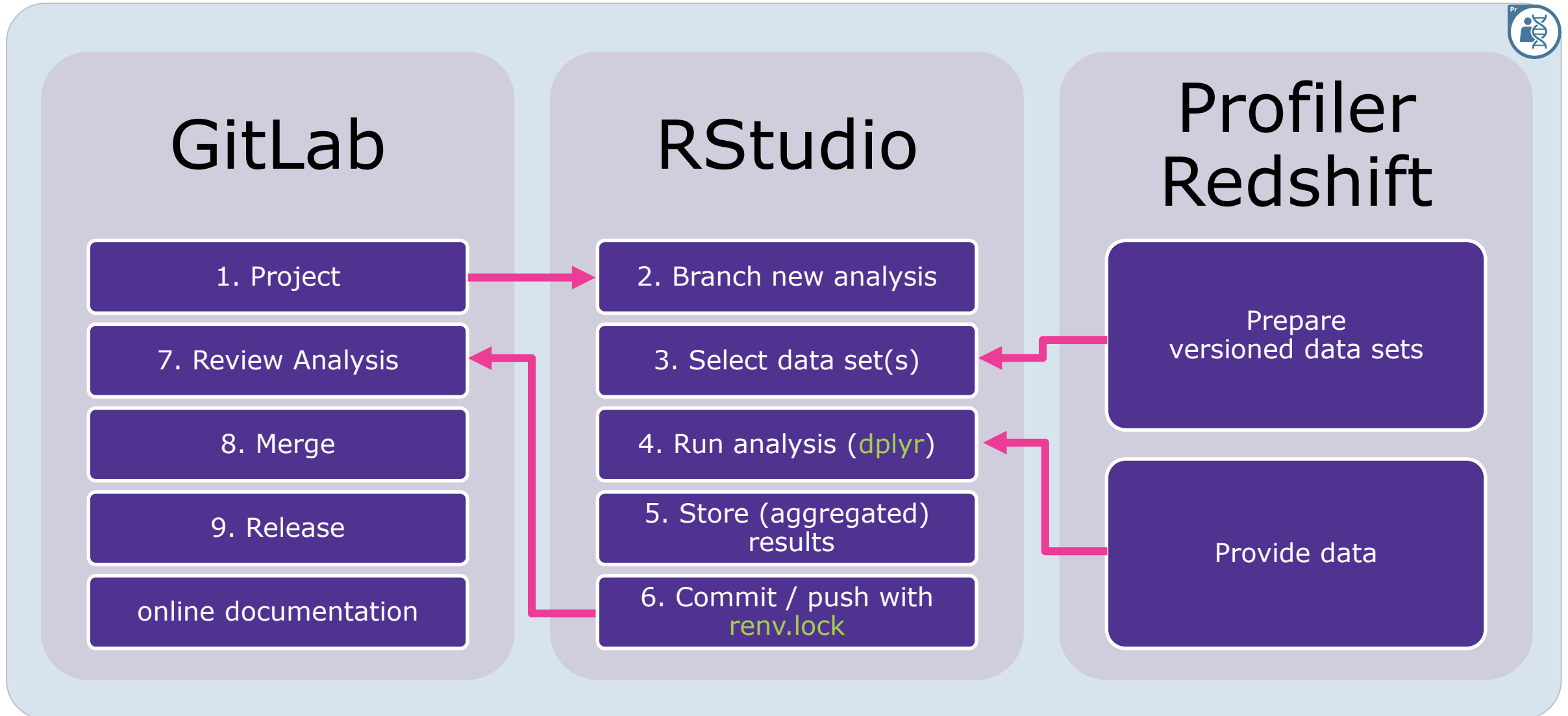
#### R package: renv

- Project specific package management tool
- References all packages & repositories with versions in a file **renv.lock**
- (soft)links all used packages

#### Redshift Spectrum (Profiler)

- S3 (parquet) file based; “unlimited” parallel versions
- Access with: SQL or **dplyr**
- Elastic scalability of storage and compute
- Caching

## The preconfigured X-OMICS platform supports traceable workflows



# Traceability

## Interactive Apps guide the traceable naming of data sets



### 1.3 Controlled vocabulary for View Names

[↑ go back to overview](#)

**Genedata Profiler (GP)** Views are the most important data entities stored and shared by XOP and the names of these views have to follow a specific standard. Each view name on XOP has to consists of a concatenation of the following nine lower-case, alphanumeric, underscore ( `_` ) separated fields of the form:

`source_dataset_analyte_assay_analysis_release_maturity_curator_version`, with a specific example being:

```
---1-----2-3-4-5-6-7-8-9
cbioportal_tcga_dna_wgs_cnv_20q1_prod_kreis_0
```

Since the last field, `version` is automatically assigned by **Genedata Profiler (GP)**, only the first eight fields have to be assigned by the user. Note that only lower-case alphanumeric characters and `_` are allowed in the final View name. Since `_` is the concatenation symbol of the view name it must not be used within any of the individual fields of the View name.

In detail, the View name fields are defined as follows:

Table 3: Field composition of the **Genedata Profiler (GP)** View name

Position	Field name	Example	Description
1	source	depmap	Lower-case shorthand name of the source of the data, most often an external organization/consortium such as <code>depmap</code> , <code>cbioportal</code> , <code>publication</code> , or the name of an internal unit or function at Merck if the data has been internally generated
2	dataset	ccle	Lower-case shorthand name of the dataset, often defined by the originating organization. Examples are <code>tcga</code> , <code>mskimpact</code> , <code>pdxdatabase</code> and so forth.
3	analyte	dna	Lower-case shorthand name for the (usually biological) analyte from which the data has been derived by measurement. Examples are <code>rna</code> / <code>dna</code> (for sequencing measurements from nucleotides), <code>img</code> (for imaging), and <code>phe</code> for phenotypic data
4	assay	wgs	Lower-case shorthand name for the experimental assay that has been applied on the sample, such as whole-genome sequencing ( <code>wgs</code> ), whole-slide imaging ( <code>ws_i</code> ), or extracts from a hospital's clinical data repository ( <code>cdr</code> )
5	analysis	snv	Lower-case shorthand name for the in-silico analysis that has been applied on the assay data, such as short nucleotide variant calling ( <code>snv</code> ), processing of patient dempgraphics data ( <code>dem</code> ), or computing cell counts based on imaging ( <code>cnt</code> )
6	release	20q1	Lower-case, lexicographically sortable release date of the data as relating to the original data source. Two formats are foreseen and supported by <code>xopdata yyqq (19q4)</code> and <code>yyyymmdd (20191231)</code> . The selected format should be kept constant, at least per dataset, but ideally per study.
7	maturity	prod	Lower-case term that captures the maturity of the data in the View; the only allowed names are <code>src</code> (for raw data originating from a data source), <code>dev</code> (for data of intermediary quality used for software development purposes), and <code>prod</code> (for data of production quality, i.e., the highest level of accuracy and robustness)
8	curator	kreis	Lower-case last name of the person within Merck who curated the data for import into XOP. Note that this role may imply initially ownership of the data within XOP; however, since ownership is transferable and View names are immutable, the owner of a View is instead captured in the View's metadata.
9	version	3	Numerical version identifier; this field is <i>automatically</i> added by <b>Genedata Profiler (GP)</b> during View generation by incrementing on existing versions of the same view. Versions start with <code>0</code> .

For the fields `analyte`, `assay`, and `analysis`, multiple examples are defined in a controlled vocabulary as follows: users are encouraged to mix and match to add new permutations while striving for re-use of existing identifiers and consistency within their own Studies:



### Naming Guide: Datasets

Please select all fields:

1. Source

gdc

2. Dataset

tcga

3. Analyte, 4. Assay, 5. Analysis

aaa

dna

rna

pro

phe

img

wt

fus

gexp

texp

gset

snv

srna

dec

6. Release

20q1

7. Maturity

dev

8. Curator

pinkert

### Recommended name:

gdc\_tcga\_rna\_wts\_gexp\_20q1\_dev\_pinkert

Copy

### Description

Position	Field.name	Example	Description
1	source	depmap	Lower-case shorthand name of the source of the data, most often as depmap, cbioportal, publication, or the name of an internal unit internally generated
2	dataset	ccle	Lower-case shorthand name of the dataset, often defined by the tcga, mskimpact, pdxdatabase and so forth.
3	analyte	dna	Lower-case shorthand name for the (usually biological) analyte from measurement. Examples are rna/dna (for sequencing measurement and phe for phenotypic data
4	assay	wgs	Lower-case shorthand name for the experimental assay that has whole-genome sequencing (wgs), whole-slide imaging (wsi), or e repository (cdr)
5	analysis	snv	Lower-case shorthand name for the in-silico analysis that has been nucleotide variant calling (snv), processing of patient dempgraphics based on imaging (cnt)
6	release	20q1	Lower-case, lexicographically sortable release date of the data as formats are foreseen and supported by xopdata yyqq (19q4) and format should be kept constant, at least per dataset, but ideally p
7	maturity	prod	Lower-case term that captures the maturity of the data in the View data originating from a data source), dev (for data of intermediary purposes), and prod (for data of production quality, i.e., the highe
8	curator	kreis	Lower-case last name of the person within Merck who curated the role may imply initially ownership of the data within XOP; however names are immutable, the owner of a View is instead captured in
9	version	3	Numerical version identifier; this field is automatically added by G generation by incrementing on existing versions of the same view



# Traceability

## X-OMICS Platform supports data traceability on multiple levels

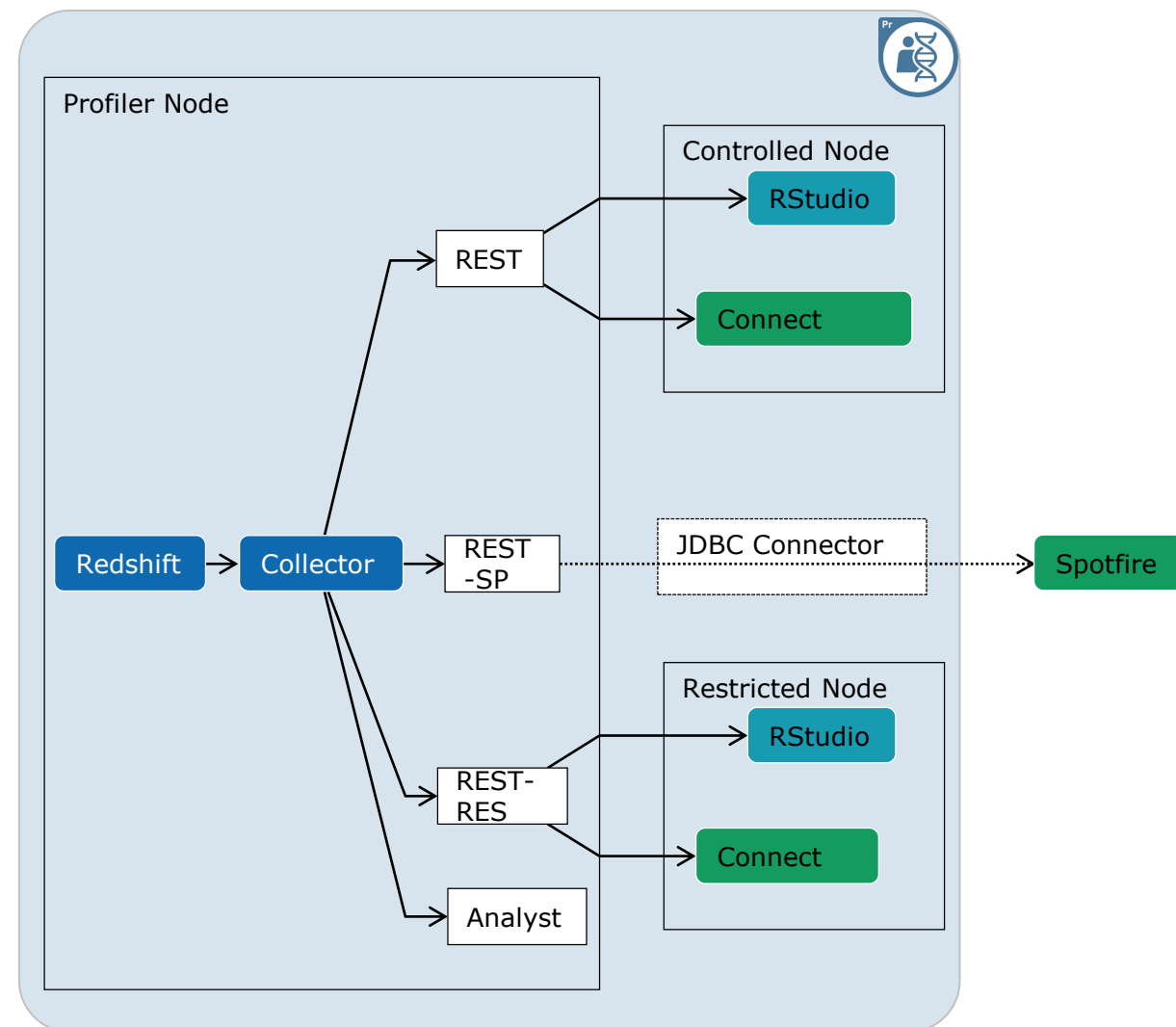
The Analytics DB (Redshift) contains the processed data.

The access from RStudio and RS-Connect is governed by **REST** interfaces.

The REST interfaces are unique for each requesting service (e.g. restricted RStudio) and controlled by **Permission Tags**.

Additionally access to data is restricted by the user having the appropriate **role** in the **study** containing the data.

Each data access is noted in the **audit trail**.



It's a teams effort!  
**Acknowledgments**

