

# PRIVAaaS: Privacy Approach for a Distributed Cloud-based Data Analytics Platforms

**Tania Basso**  
**Regina Moraes**  
**Nuno Antunes**  
**Marco Vieira**  
**Walter Santos**  
**Wagner Meira Jr.**



UNICAMP



UFMG  
UNIVERSIDADE FEDERAL  
DE MINAS GERAIS



# Context

- **Big Data** → datasets with large amount of data
- **Big Data Analytics** → use advanced analytic techniques and algorithms against these datasets
- Individuals provide *Personally Identifiable Information (PII)* → **privacy concerns**

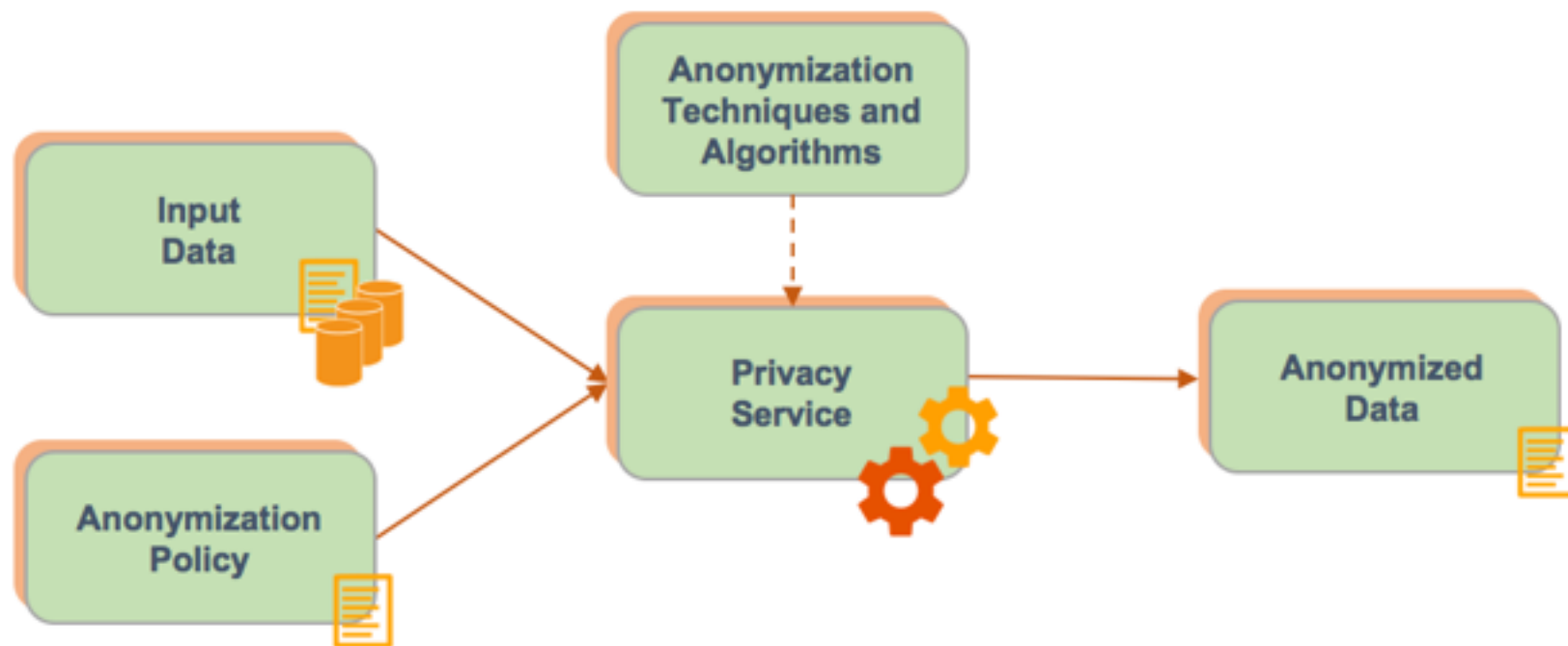
# Context

- Approach for protecting data privacy based on anonymization policies – **PRIVAAAAS**
- Discusses PRIVAAAAS integration in the LEMONADE
- LEMONADE (**L**ive **E**nvironment for **M**ining of **N**on-trivial **A**mount of **D**ata from **E**verywhere) - web-based platform for data analytics

# PRIVAaaS

- Set of tools to protect sensible information processed by data analytics
- Based on anonymization techniques and policies
- Free, open source, developed in Java
- Implements generalization, suppression, masking (replacement) and encryption.

# PRIVAaaS



# PRIVAaaS



## Anonymization Policy

```
1 [
2   [
3     {
4       "FIELD_NAME": "first_name",
5       "TYPE": "MAS",
6       "DETAIL": "FIRST"
7     },
8     {
9       "FIELD_NAME": "last_name",
10      "TYPE": "SUP",
11      "DETAIL": ""
12    }
13 ]
14 ]
```

## Database (JSON) file

```
1 @[{ "id": 1, "first_name": "Deborah", "last_name": "Carroll",
2     "email": "mcarroll0@parallels.com", "gender": "Female",
3     "ip_address": "71.188.126.112", "job": "Help Desk Technician",
4     "salary": "$4.53", "education": "Maharshi Dayanand Sarswati
5     University Ajmer"},
6     { "id": 2, "first_name": "Dorothy", "last_name": "Alvarez", "email":
7     "dalvarez1@cloudflare.com", "gender": "Female", "ip_address":
8     "75.12.204.166", "job": "Media Manager IV", "salary": "$5.51",
9     "education": "National Taiwan College of Physical Education
10    and Sports"}]
```

Privacy Service

## Anonymized Database (JSON) file

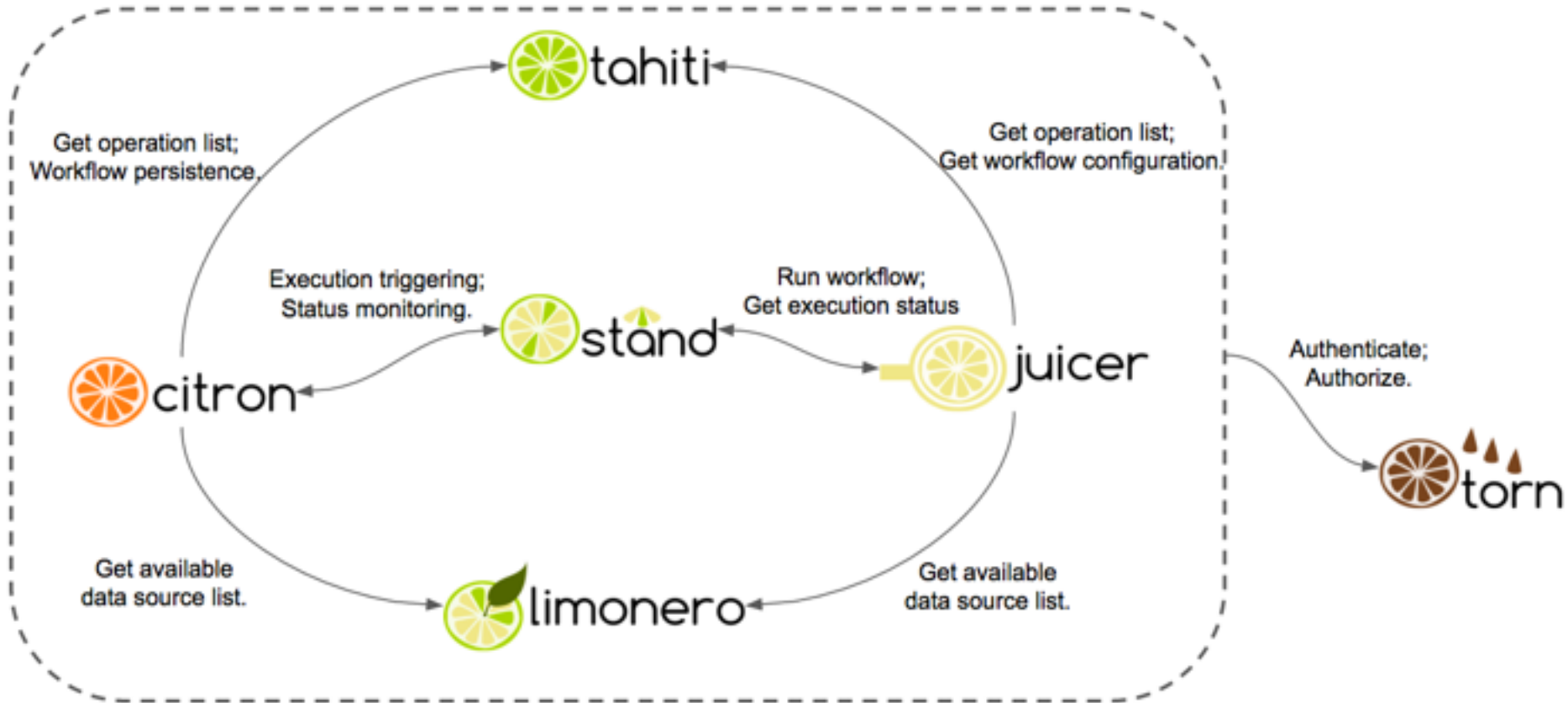
```
1 @[{ "id": 1, "first_name": "Michelle", "last_name": "", "email":
2     "mcarroll0@parallels.com", "gender": "Female", "ip_address":
3     "71.188.126.112", "job": "Help Desk Technician", "salary":
4     "$4.53", "education": "Maharshi Dayanand Sarswati University
5     Ajmer"},
6     { "id": 2, "first_name": "Susan", "last_name": "", "email":
7     "dalvarez1@cloudflare.com", "gender": "Female", "ip_address":
8     "75.12.204.166", "job": "Media Manager IV", "salary": "$5.51",
9     "education": "National Taiwan College of Physical Education
10    and Sports"}]
```

# LEMONADE

Live Environment for Mining of Non-trivial Amount of Data from Everywhere

- Interoperable, scalable and open platform for data analytics
- Users can drag and drop operations and data sources to compose different ETL processes and Machine Learning Workflows
- Runs on Spark, while providing resource allocation and QoS in a cloud-based system

# LEMONADE - Architecture

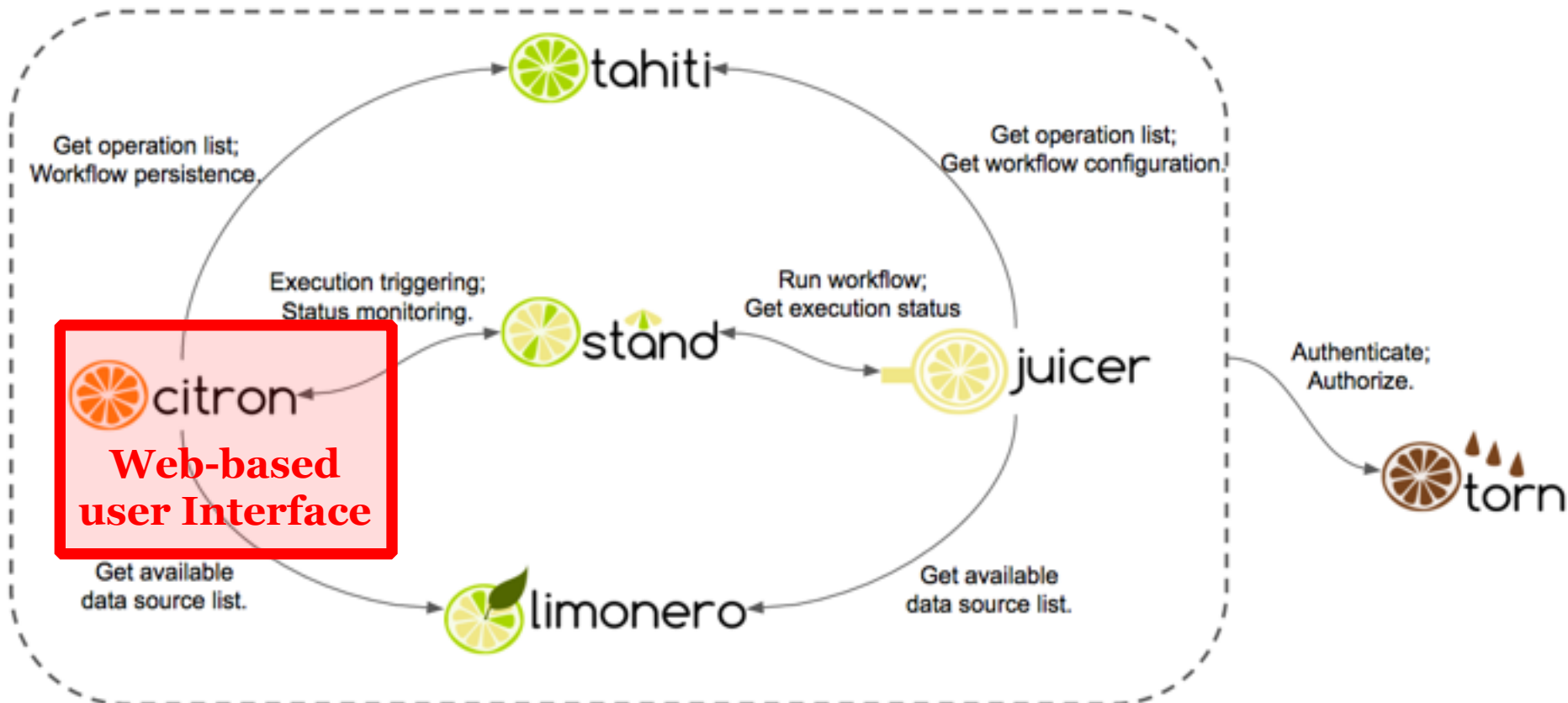




# LEMONADE



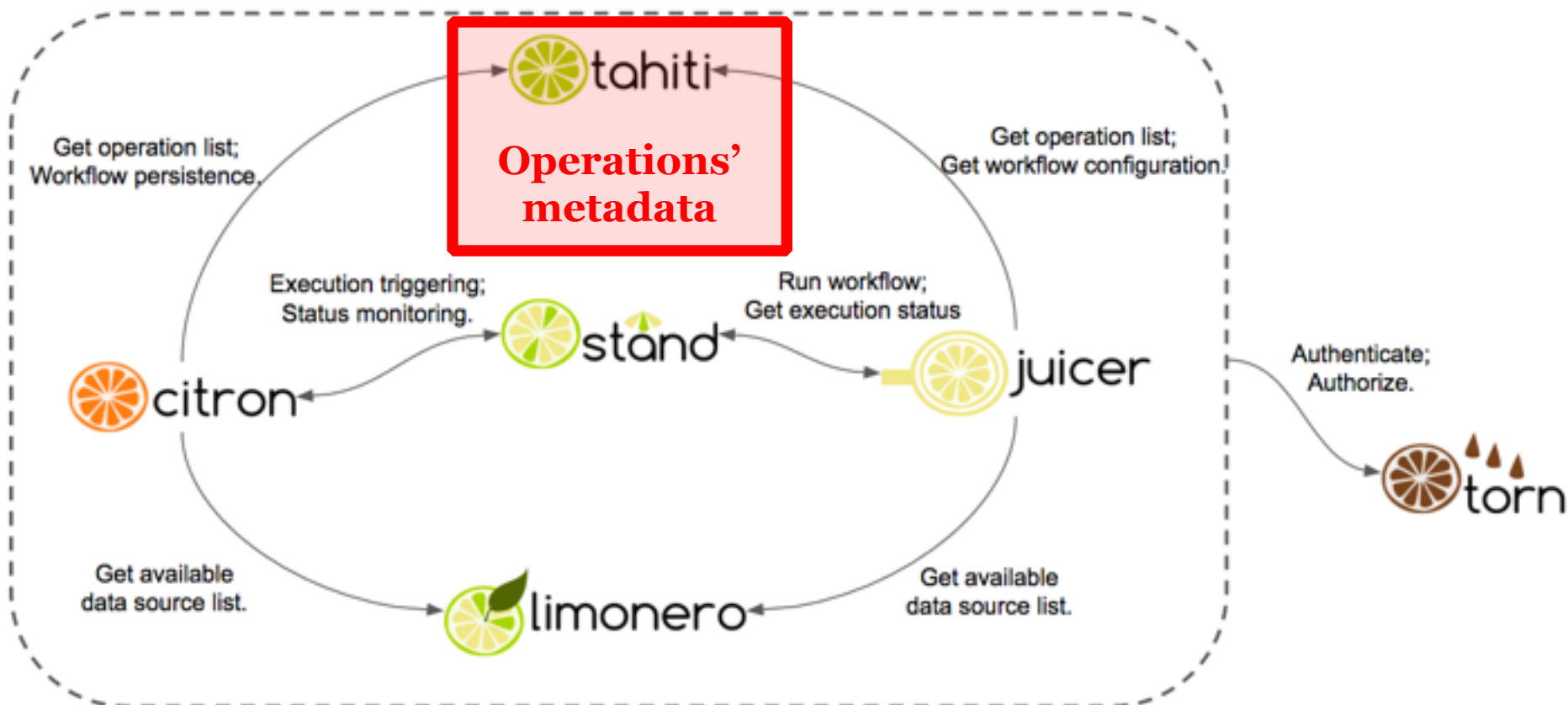
UNICAMP



# LEMONADE



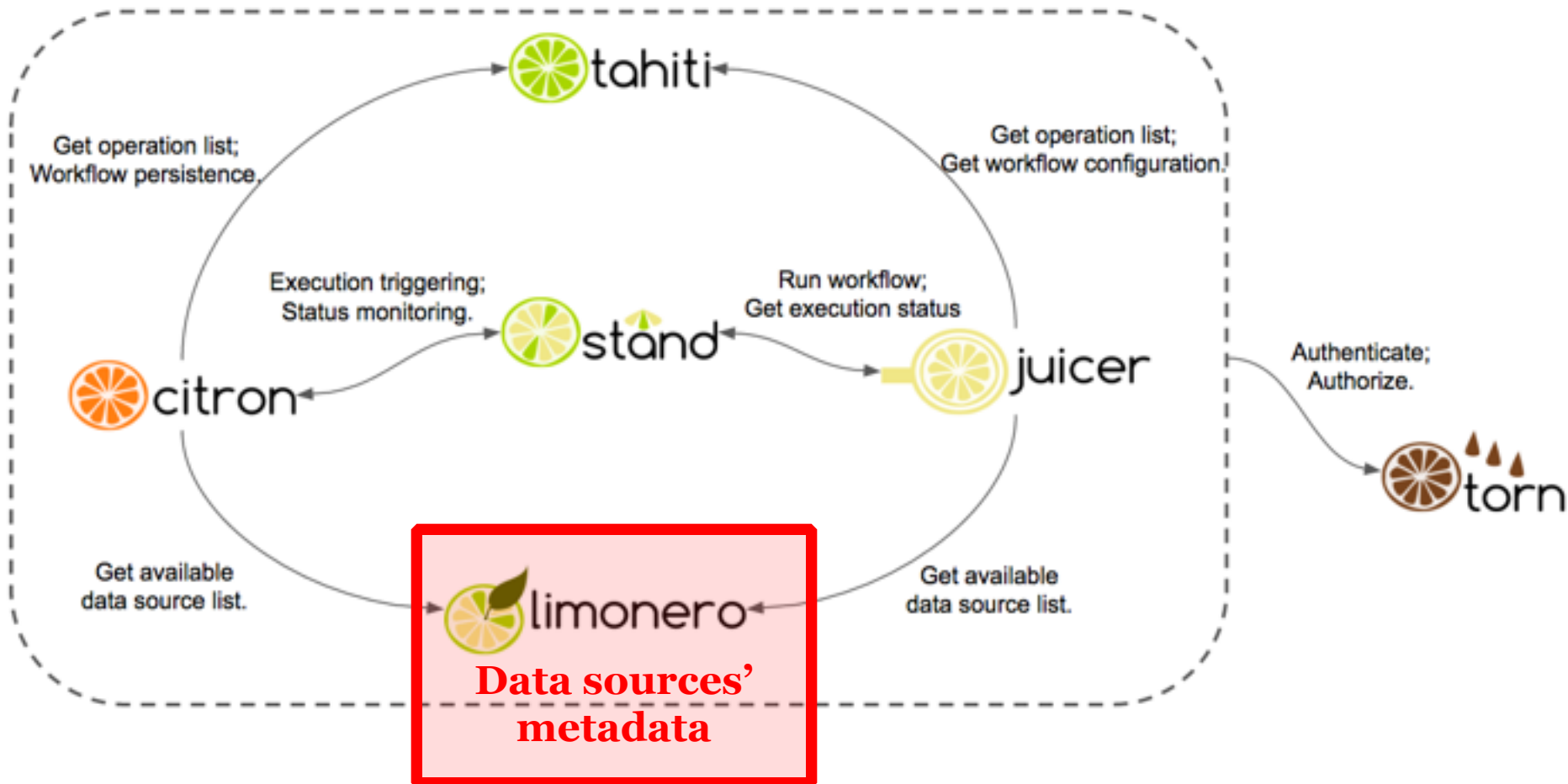
UNICAMP



# LEMONADE



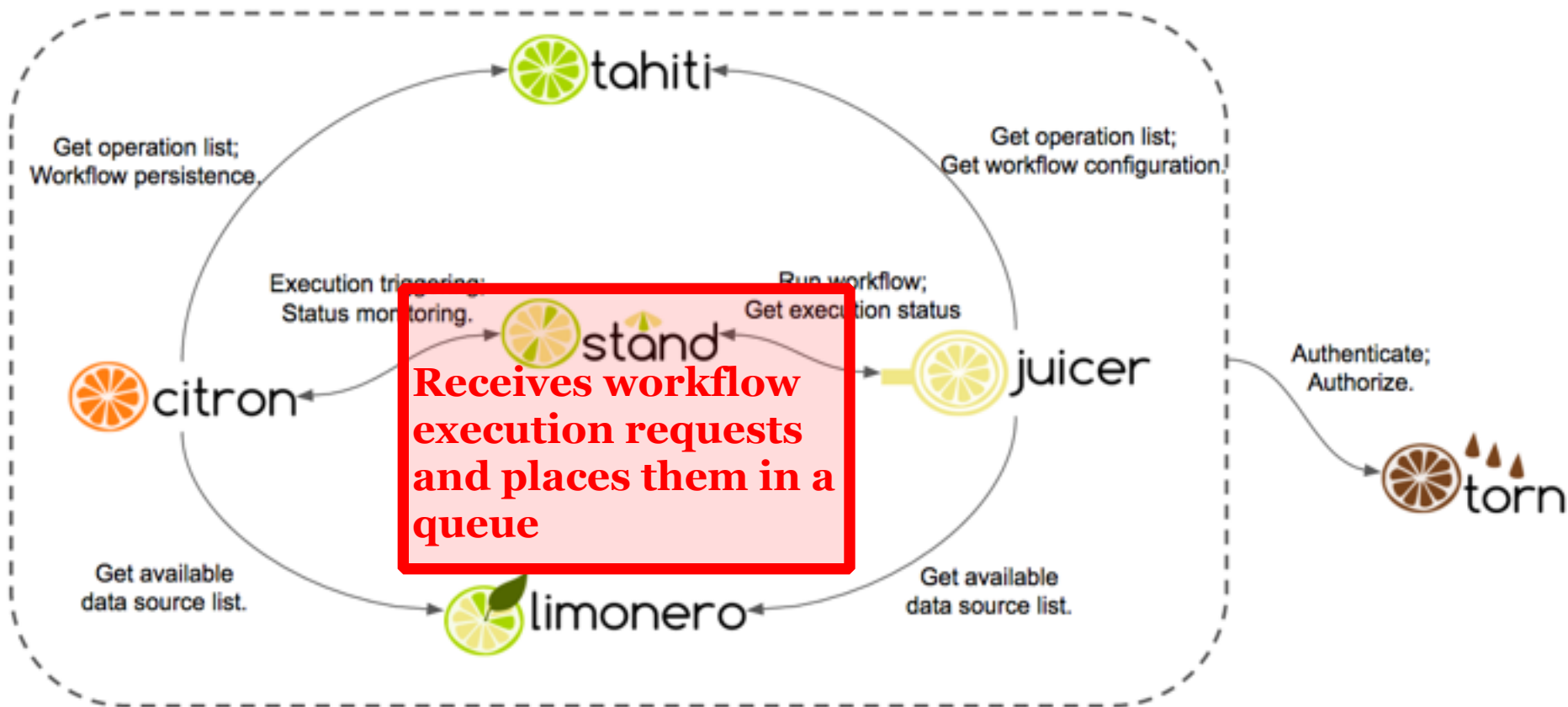
UNICAMP



# LEMONADE



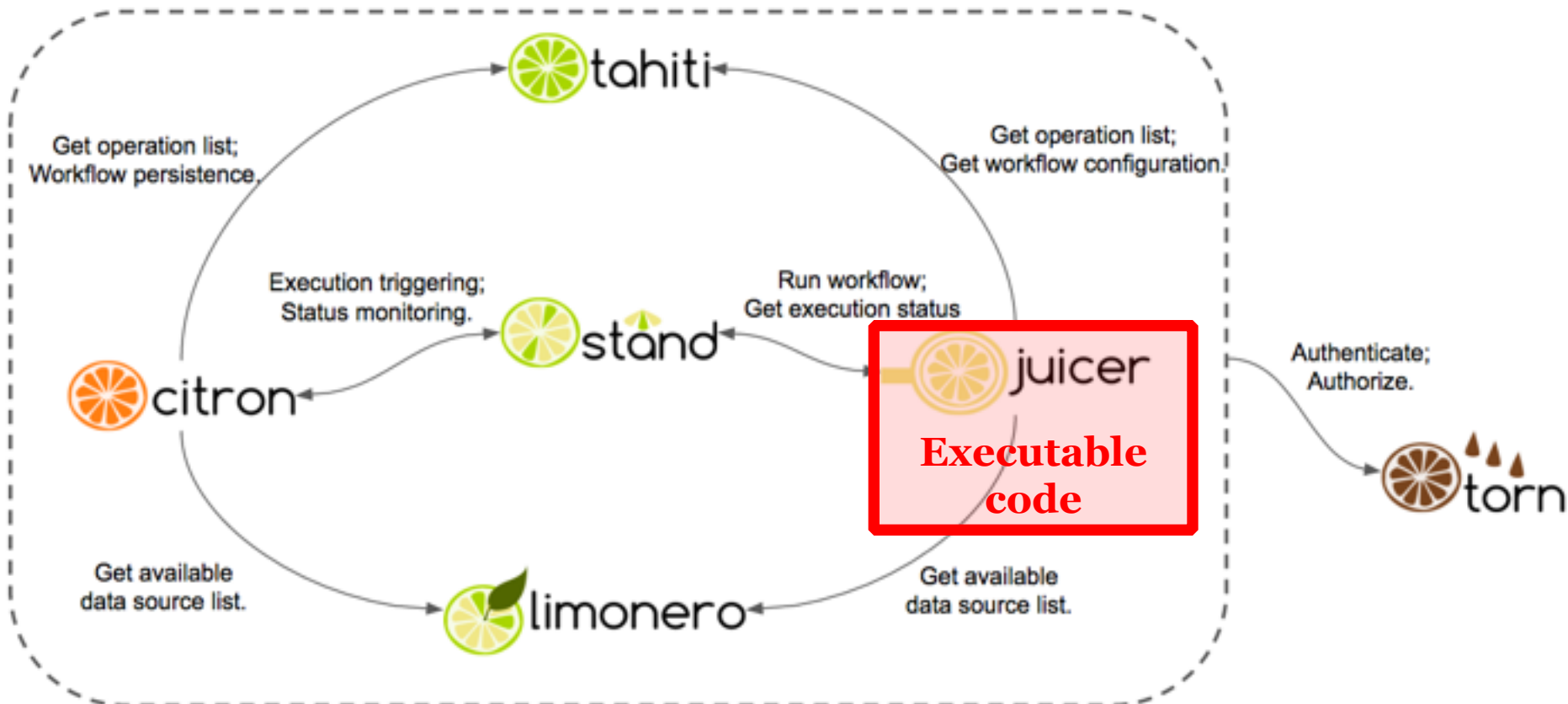
UNICAMP



# LEMONADE



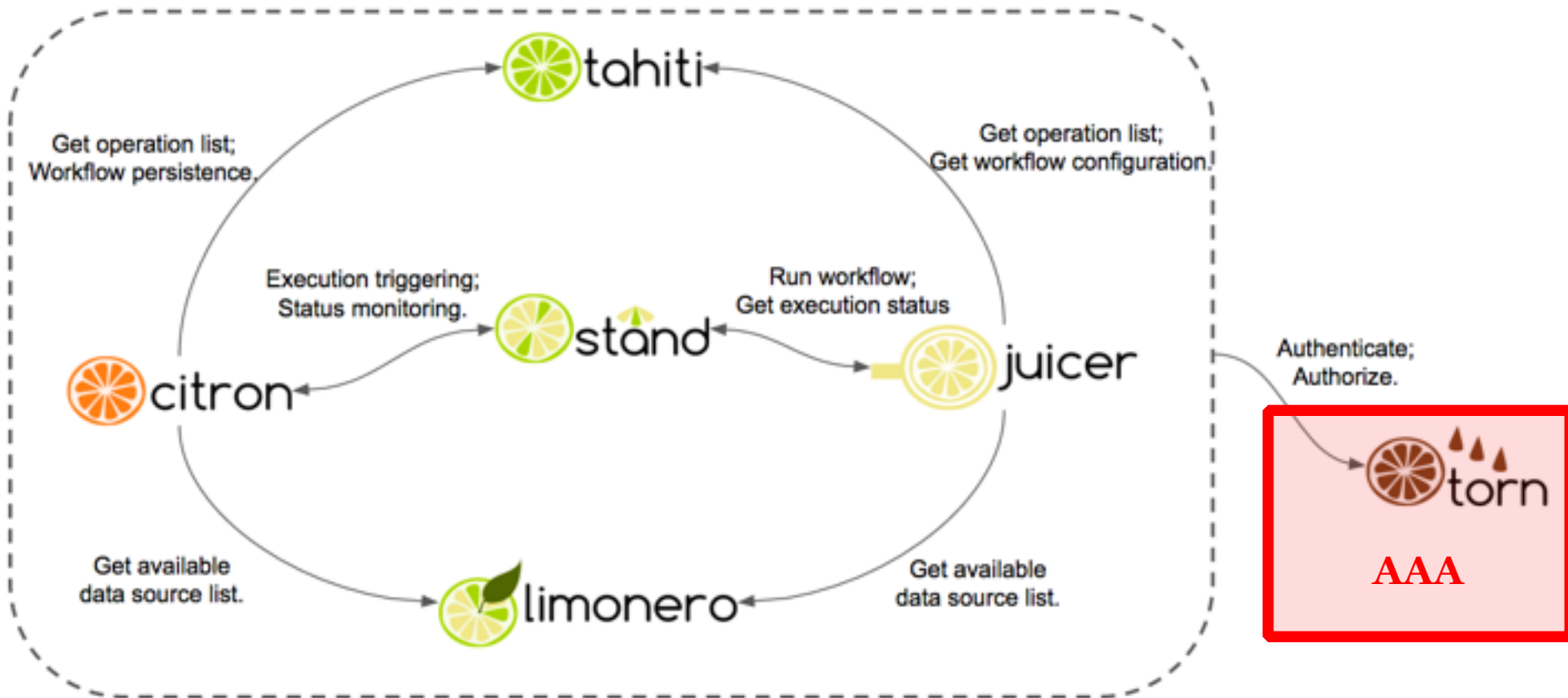
UNICAMP



# LEMONADE



UNICAMP

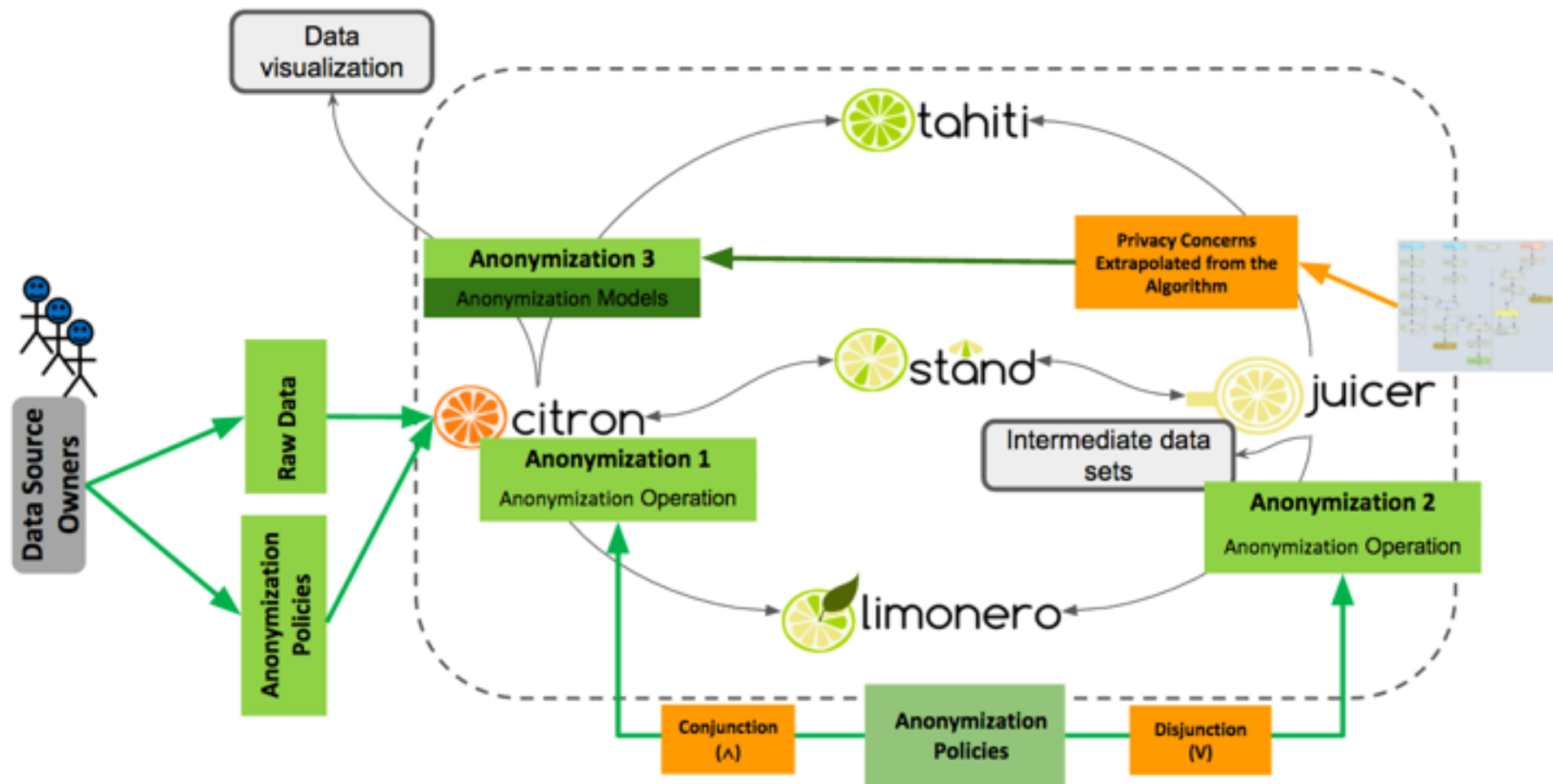


# PRIVAaaS integration on LEMONADE



- Enforcement of anonymization policies → data visualization is performed with privacy protection.
- Protect sensible information while deals with data utility, keeping it at the highest possible level
- Divided into three levels, which model different aspects of the anonymization problem

# PRIVAaaS integration on LEMONADE

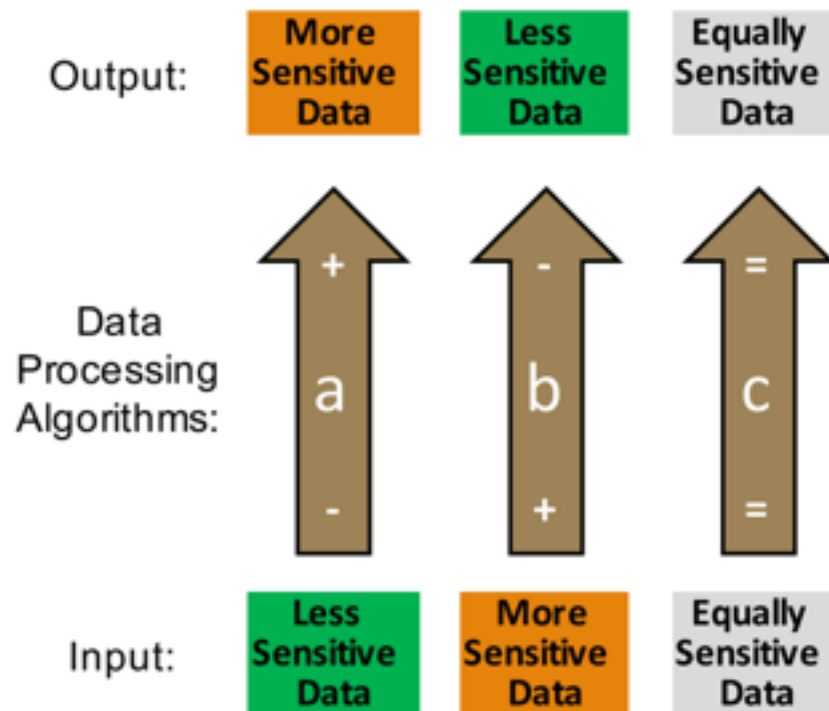




# PRIVAaaS integration on LEMONADE



- Control statistical disclosure of data analytics algorithms
- Three scenarios
- Reduce re-identification



# Conclusions

- PRIVAaas is a 3-level integration approach
  - supports not only the privacy of native input data
  - supports also the privacy of output data generated by the operations
  - uses anonymization models prior to data visualization

# Conclusions

- Use of anonymization polices allows:
  - data source owners to establish their own guidelines
  - applying conjunctive/disjunctive processes, preserving the data utility
  - enforcement guarantees that the rules established by data source owners will be respected

# Conclusions

- Provide guarantees not only to data sources, but also to the operations and their outputs
- Go beyond static elements, but also on their dynamic behavior
- Implementation may raise issues that have not been addressed in the literature
  - disruptive platform
  - data analytics applications to an unprecedented level of service

# Acknowledgments

- This work has been partially supported by the projects:
  - **EUBra-BIGSEA** ([www.eubra-bigsea.eu](http://www.eubra-bigsea.eu)), funded by the Brazilian Ministry of Science, Technology and Innovation (Project 23614 - MCTI/RNP 3rd Coordinated Call) and by the European Commission under the Cooperation Programme, Horizon 2020 grant agreement no 690116
  - **DEVASSES** ([www.devasses.eu](http://www.devasses.eu)), funded by the European Union's FP7 under grant agreement no PIRSES-GA-2013-612569
  - UFMG team is also partially supported by **CNPq**, **FAPEMIG**, as well as projects **InWeb** and **MASWeb**.

# Questions



**regina@ft.unicamp.br**

**taniabasso@ft.unicamp.br**

**nmsa@dei.uc.pt**

**mvieira@dei.uc.pt**

**walter@dcc.ufmg.br**

**meira@dcc.ufmg.br**