

# A Declarative Approach to Information Extraction using Web Services API

Slides  
before 1st  
Section  
Divider

Studied issue

Locks and keys

Declarative data  
extraction

Data integration with  
WS

Main  
scenarios

Tutorial ICWE 2016  
June, 7<sup>th</sup>

Results and  
discussion

Additional  
slides

John SAMUEL  
Christophe REY

# A Declarative Approach to Information Extraction using Web Services API

ICWE 2016 Tutorial  
June, 7<sup>th</sup>

John SAMUEL

Université de Lyon - FRANCE

Christophe REY

LIMOS, Université Clermont Auvergne - FRANCE

# Outline

- Studied issue → Integrating with WS
- Locks and Keys → SME, few developers, numerous WS  
→ be declarative !
- Declarative data extraction → A proposal
- Scientific basis → LAV mediation with access patterns
- Main scenarios
- Results and discussion

# [Studied issue]

## Studied issue

Locks and Keys

Declarative data extraction

Data integration with WS

Main scenarios

Results and discussion

# Small and Medium Scale Enterprises dependent on Numerous Web Services



# Small and Medium Scale Enterprises (SME)

- ❑ Less human resources
- ❑ Dependent on numerous web services for their day to day activities
- ❑ Data spread across different data centers
- ❑ No direct control over their own data
- ❑ Need to access historical data
- ❑ Easier mechanism to compute performance indicators

# Web Services (WS)

- Services accessible on internet
- Primary audience: users with web browsers
- Cater to needs of numerous clients
- Manage user resources
- Provide ability to access and manipulate user resources
- Heterogeneous and autonomous
- For machine to machine access, service providers provide API  
(Application Programming Interface)

# Integration with WS

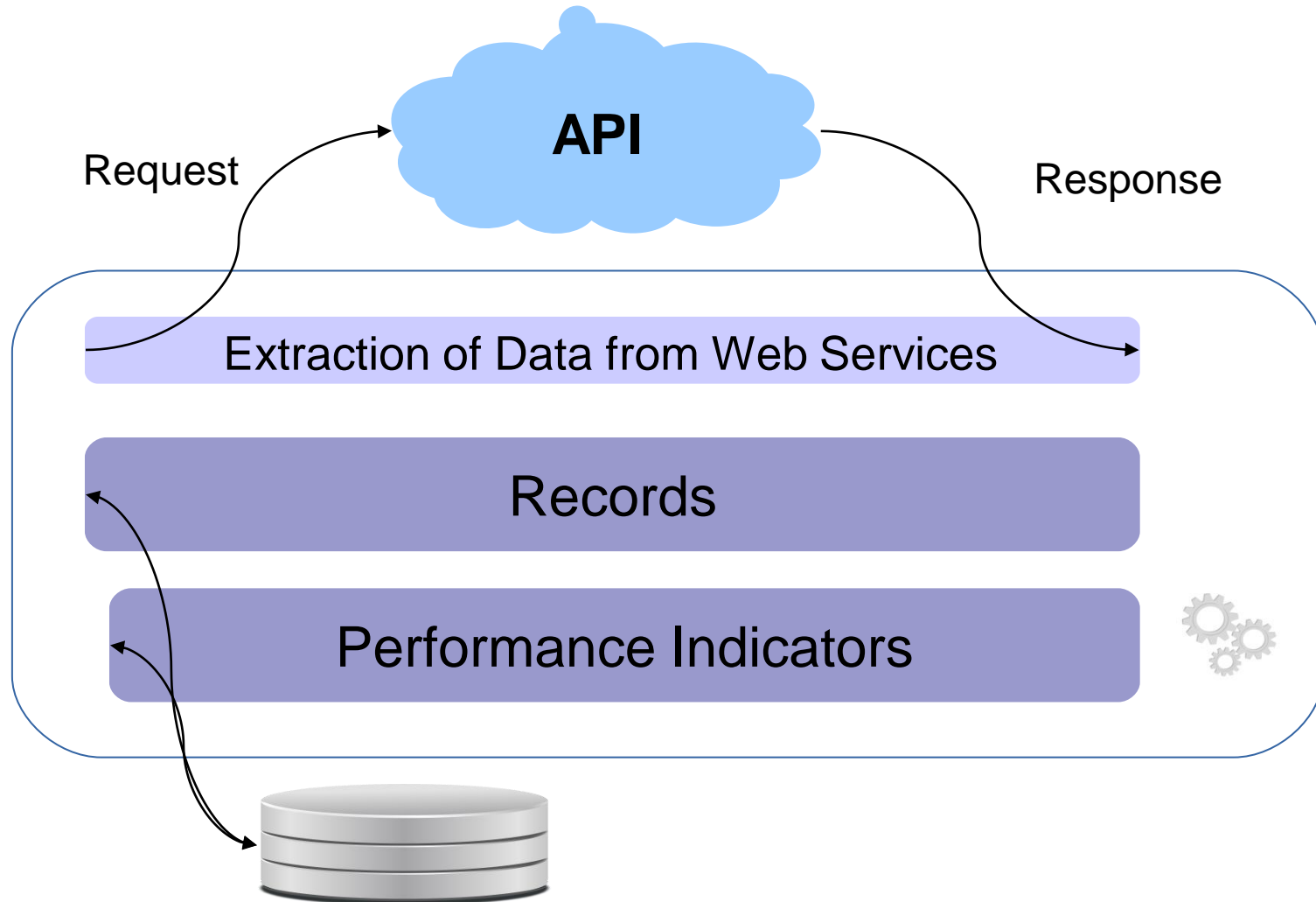
- Problem:
  - Integrate with **many** web services using their APIs
  - More precisely:  
“integration” = “data extraction”
- Constraint
  - **Very few** developers
- Initially proposed by Rootsystem, Vichy, France



# Integration with WS

- Goal 1  
Extract and transform data from WS  
using APIs  
Historical data
- Goal 2  
Compute performance indicators  
using historical data

# Integration with WS



# Real Life Example

- Social networking service: Twitter API
  - Resources: user, status, followers, lists,
  - Analytics: favorites, retweets
- REST API
  - Operations: Read, Write, Update and Delete
  - Simple and widely used
- Message formats
- Authentication mechanism

# Performance Indicators

- Performance Indicators like profit, sales etc.
- Historical data for comparing the evolution of enterprise performance
- Use of various web services like social networking sites. E.g. Twitter
- Example: Number of followers on a social media site

# Twitter Documentation

## Documentation

[Fabric](#)

[Twitter for Websites](#)

[Cards](#)

[OAuth](#)

[REST APIs](#)

[Streaming APIs](#)

[Ads APIs](#)

[MoPub](#)

[Best Practices](#)

[API Overview](#)

[API Status](#)

[Playbooks](#)

[Events](#)


[Case Studies](#)

[Manage My Apps](#)

[Terms of Use](#)

## Documentation

The Twitter Platform connects your website or application with the worldwide conversation happening on Twitter.

 **API Changelog** [Follow](#)  
Get auto-notified when **Twitter API** changes  
[View Details](#)

## Fabric

[Fabric](#) is a flexible, modular set of mobile development tools called “Kits” that help you make your app more stable, add social features like sharing and login, and turn your app into a business with easy monetization.

## Twitter for Websites

[Twitter for Websites](#) is a suite of embeddable widgets, buttons, and client-side scripting tools to integrate Twitter and display Tweets on your website or JavaScript application, including a [single Tweet](#), [multiple Tweets](#), [Twitter Moments](#), [Tweet Button](#), and the [Follow Button](#).

## Cards

# Twitter REST API

[API Console Tool](#)

[Public API](#)

[The Search API](#)

[The Search API: Tweets by Place](#)

[Working with Timelines](#)

[API Rate Limits](#)

[API Rate Limits: Chart](#)

[GET statuses/mentions\\_timeline](#)

[GET statuses/user\\_timeline](#)

[GET statuses/home\\_timeline](#)

[GET statuses/retweets\\_of\\_me](#)

[GET statuses/retweets/:id](#)

[GET statuses/show/:id](#)

[POST statuses/destroy/:id](#)

[POST statuses/update](#)

[POST statuses/retweet/:id](#)

[POST statuses/unretweet/:id](#)

[POST statuses/update\\_with\\_media](#)

[GET statuses/oembed](#)

[GET statuses/retweeters/ids](#)

[GET statuses/lookup](#)

## REST APIs

The [REST APIs](#) provide programmatic access to read and write Twitter data. Author a new Tweet, read author profile and follower data, and more. The REST API identifies Twitter applications and users using [OAuth](#); responses are available in JSON.

If your intention is to monitor or process Tweets in real-time, consider using the [Streaming API](#) instead.

## Overview

Below are the documents that will help you get going with the REST APIs as quickly as possible

- [API Rate Limiting](#)
- [API Rate Limits](#)
- [Working with Timelines](#)
- [Using the Twitter Search API](#)
- [Uploading Media](#)
- [Multiple Media Entities in Statuses](#)
- [Finding Tweets about Places](#)

## Latest Updates

As of version 1.1 of the Twitter API, the more recent updates to our API are highlighted below. We're excited about what it means for developers and we've captured all the meaningful changes here so you don't miss a thing.

# Twitter Message Format: JSON

API Console Tool

Public API

The Search API

The Search API: Tweets by Place

Working with Timelines

API Rate Limits

API Rate Limits: Chart

GET statuses/mentions\_timeline

GET statuses/user\_timeline

GET statuses/home\_timeline

GET statuses/retweets\_of\_me

GET statuses/retweets/:id

GET statuses/show/:id

POST statuses/destroy/:id

POST statuses/update

POST statuses/retweet/:id

POST statuses/unretweet/:id

POST statuses/update\_with\_media

GET statuses/oembed

## GET followers/list

Returns a cursored collection of user objects for users following the specified user.

At this time, results are ordered with the most recent following first — however, this ordering is subject to unannounced change and eventual consistency issues. Results are given in groups of 20 users and multiple “pages” of results can be navigated through using the `next_cursor` value in subsequent requests. See [Using cursors to navigate collections](#) for more information.

## Resource URL

<https://api.twitter.com/1.1/followers/list.json>

## Resource Information

Response formats	JSON
------------------	------

Requires authentication?	Yes
--------------------------	-----

Rate limited?	Yes
---------------	-----

Requests / 15-min window (user auth)	15
--------------------------------------	----

Requests / 15-min window (app auth)	30
-------------------------------------	----

# Twitter Status

POST mutes/users/create  
POST mutes/users/destroy  
GET mutes/users/ids  
GET mutes/users/list  
GET users/suggestions/:slug  
GET users/suggestions  
GET users/suggestions/:slug/  
members  
GET favorites/list  
POST favorites/destroy  
POST favorites/create  
GET lists/list  
GET lists/statuses  
POST lists/members/destroy  
GET lists/memberships  
GET lists/subscribers  
POST lists/subscribers/create  
GET lists/subscribers/show  
POST lists/subscribers/destroy  
POST lists/members/create\_all  
GET lists/members/show  
GET lists/members  
POST lists/members/create  
POST lists/destroy

## Example Request

GET

[https://api.twitter.com/1.1/statuses/mentions\\_timeline.json?count=2&since\\_id=14927799](https://api.twitter.com/1.1/statuses/mentions_timeline.json?count=2&since_id=14927799)

## Example Result

```
[
  {
    "coordinates": null,
    "favorited": false,
    "truncated": false,
    "created_at": "Mon Sep 03 13:24:14 +0000 2012",
    "id_str": "242613977966850048",
    "entities": {
      "urls": [

    ],
      "hashtags": [

    ],
      "user_mentions": [
        {
          "name": "Jason Costa",
          "id_str": "14927800",
          "id": 14927800,
          "indices": [
            0,
            11
```



# Twitter Authentication: OAuth

Overview

Application-Only  
Authentication

3-Legged OAuth

PIN-Based OAuth

xAuth

OAuth Echo

## OAuth

Send secure authorized requests to the  
Twitter API

Twitter uses [OAuth](#) to provide authorized access to its API.



# Running (toy) example

- Two social networking services:
  - $SN_1$
  - $SN_2$
- $SN_1$  uses XML data format
- $SN_2$  uses JSON data format
- Both provide API operations to access
  - user status
  - view counts

# Example: $SN_1$ operations

- $sn1myid^o(userid)$ 
  - Output: the user identifier of the respective owner of the social networking service
- $sn1user^i(userid, name)$ 
  - Input: the user identifier
  - Output: the name of the user
- $sn1status^{io}(userid, statusid, creationdate)$ 
  - Input: the identifier of a user
  - Output: the details of statuses like the status identifier and the creation date
- $sn1viewcount^o(statusid, count)$ 
  - Input: the status identifier
  - Output: returns the view count

**sn1myid**

userid
1
...

**sn1user**

userid	name
1	John
...	...

**sn1status**

userid	statusid	creationdate
1	12	05-12-02
...	...	...

**sn1viewcount**

statusid	count
12	5
...	...

# Example: SN<sub>2</sub> operations

- $\text{sn2myid}^{\circ\circ}(\text{userid}, \text{name})$ 
  - Output: the user identifier and name of the respective owner of the social networking service.
- $\text{sn2update}^{\text{i}\circ\circ\circ}(\text{userid}, \text{statusid}, \text{creationdate}, \text{count})$ 
  - Input: the identifier of a user
  - Output: the details of statuses like the status identifier, creation date and view count.

**sn2myid**

userid	name
54	Chris
...	...

**sn2update**

userid	statusid	creationdate	count
54	12	05-12-02	4
...	...	...	...

# [Locks and keys]

Studied issue

**Locks and Keys**

Declarative data extraction

Data integration with WS

Main scenarios

Results and discussion

# Commonly Used Solution

1. Read the documentation
2. Coding a specific wrapper  
(with a programming language)
3. Coding business specific code  
(with a programming language)

Example of programming languages

JAVA, C++, C#, ...

➔ procedural languages

# Commonly Used Solution

1. Read the web service API documentation
    - Determine interesting operations
    - Determine message formats (XML, JSON etc.)
    - Determine API authentication mechanisms
    - For every operation, determine the input and output parameters, required HTTP URL, header, method and body contents
- ➔ In real life, only example inputs and outputs are given.

# Commonly Used Solution

2. Coding a specific wrapper with a (procedural) progr. language
  - Translate API input/output parameters to programming language related data structures (int, string, struct, class...)
  - Code to make web service API operation call (framing HTTP method, URL, body etc. on a per-operation basis)
  - Code to receive API operation responses and extract and transform this response to chosen internal format (data structures)



# Commonly Used Solution

3. Coding business specific code with a (procedural) progr. language
  - Determining what information is really interesting for the enterprises
  - Determining the sequence of operations to obtain this information from the web service API operations using the code previously written

# Estimated cost ? (in time)

- Relative to developer experience
- Relative to complexity of web service API operations
- Approximation:  
it may take hours/days to develop a program to integrate with one web service API.

# Twitter Libraries: Java

Documentation

Best Practices

API Overview

Object: Users

Object: Tweets

Object: Entities

Object: Entities in Objects

Object: Places

Twitter IDs

Connecting to Twitter API using  
TLS

Using cursors to navigate  
collections

Error Codes & Responses

Twitter Libraries

API Status

Playbooks

## Twitter Libraries

These libraries, while not necessarily tested by Twitter, should support Twitter API v1.1.

Want your library to be included in this index or need to update the details we have? [Submit your library for inclusion!](#)

## Libraries built and maintained by Twitter

### Java

- [hbc](#) — A Java HTTP client for consuming Twitter's Streaming API
- [twitter-kit-android](#) — Twitter Kit is a multi-module gradle project containing several Twitter SDKs including TweetComposer, TwitterCore, and TweetUi. It is built on the Fabric platform and uses many shared components.

## Libraries built for the Twitter Platform

### Multi-platform

- [Temboo](#) — *by @temboo* — Framework for working with Twitter via many platforms including iOS, Android, Java, PHP, Python, Ruby, and Node.js

# Twitter Library: HBC

```
120 public void run(String consumerKey, String consumerSecret, String token, String tokenSecret)
121     throws InterruptedException, ControlStreamException, IOException {
122     // Create an appropriately sized blocking queue
123     BlockingQueue<String> queue = new LinkedBlockingQueue<String>(10000);
124
125     // Define our endpoint: By default, delimited=length is set (we need this for our processor)
126     // and stall warnings are on.
127     List<Long> followings = new ArrayList<Long>();
128     followings.add(111111111L);
129     followings.add(222222222L);
130
131     SitestreamEndpoint endpoint = new SitestreamEndpoint(followings);
132     Authentication auth = new OAuth1(consumerKey, consumerSecret, token, tokenSecret);
133
134     // Create a new BasicClient. By default gzip is enabled.
135     BasicClient client = new ClientBuilder()
136         .hosts(Constants.SITESTREAM_HOST)
137         .endpoint(endpoint)
138         .authentication(auth)
139         .processor(new StringDelimitedProcessor(queue))
140         .build();
141
142     // Create an executor service which will spawn threads to do the actual work of parsing the incoming messages and
143     // calling the listeners on each message
144     int numProcessingThreads = 4;
145     ExecutorService service = Executors.newFixedThreadPool(numProcessingThreads);
146
147     // Wrap our BasicClient with the twitter4j client
148     Twitter4jSitestreamClient t4jClient = new Twitter4jSitestreamClient(
149         client, queue, Lists.newArrayList(listener), service);
150
151     // Establish a connection
152     t4jClient.connect();
153     for (int threads = 0; threads < numProcessingThreads; threads++) {
154         // This must be called once per processing thread
155         t4jClient.process();
156     }
157
```

# Problem 1: API modifications

- Change in message formats  
(XML ↔ JSON...)
- Change in authentication mechanisms
- Change in input/output parameters
- Change in sequence of operations  
required to extract desired information

## Problem 2: WS lifetime

- What if the user decides to change WS or if WS shuts down ?
- Loss of historical data (unless some data export option is provided)
- Rewrite previous code or make use of some adapter to integrate with new web service

# Proposal

- Constraints (recalls)
  - Many WS APIs to integrate
  - Very few developers
- Idea
  - Replacing procedural programming languages by declarative ones
  - View integration as an administration task (not a coding one anymore)

# Procedural vs declarative

- Procedural program: « HOW »
  - Developers describe the steps (computations) the system has to follow to reach the intended result.
- Declarative program: « WHAT »
  - Developers: describe the intended result.
  - System: finds the steps to compute it.
  - Same outputs for same inputs
  - No side-effect



# Procedural vs declarative

- Procedural → declarative
  - Loops → recursive calls
  - Mutable variables → immutable variables
- Examples of more or less declarative languages
  - SQL (database querying)
  - Prolog (logic programming)
  - Lisp (functional programming)
  - CHR (constraint programming)
  - XML, JSON (content and presentation description)

# Recalls about declarative lang.

- Any programming language can be used declaratively when
  - It has a sufficient number of libraries, or
  - It is based on a generic reasoning
    - Prolog and SLD-resolution
    - SQL and associated processing engines (query optimizer and memory management handler).
- Declarative programming
  - an additional abstraction layer
  - often a new syntaxe

# Advantages of declarativeness

- Smaller code → safer code
  - less syntax errors
  - conception errors: easier to locate and understand
- Closer to natural language
  - easier to understand
- → improved development time (e.g. prototyping)
- → improved scalability

# Drawbacks of declarativeness

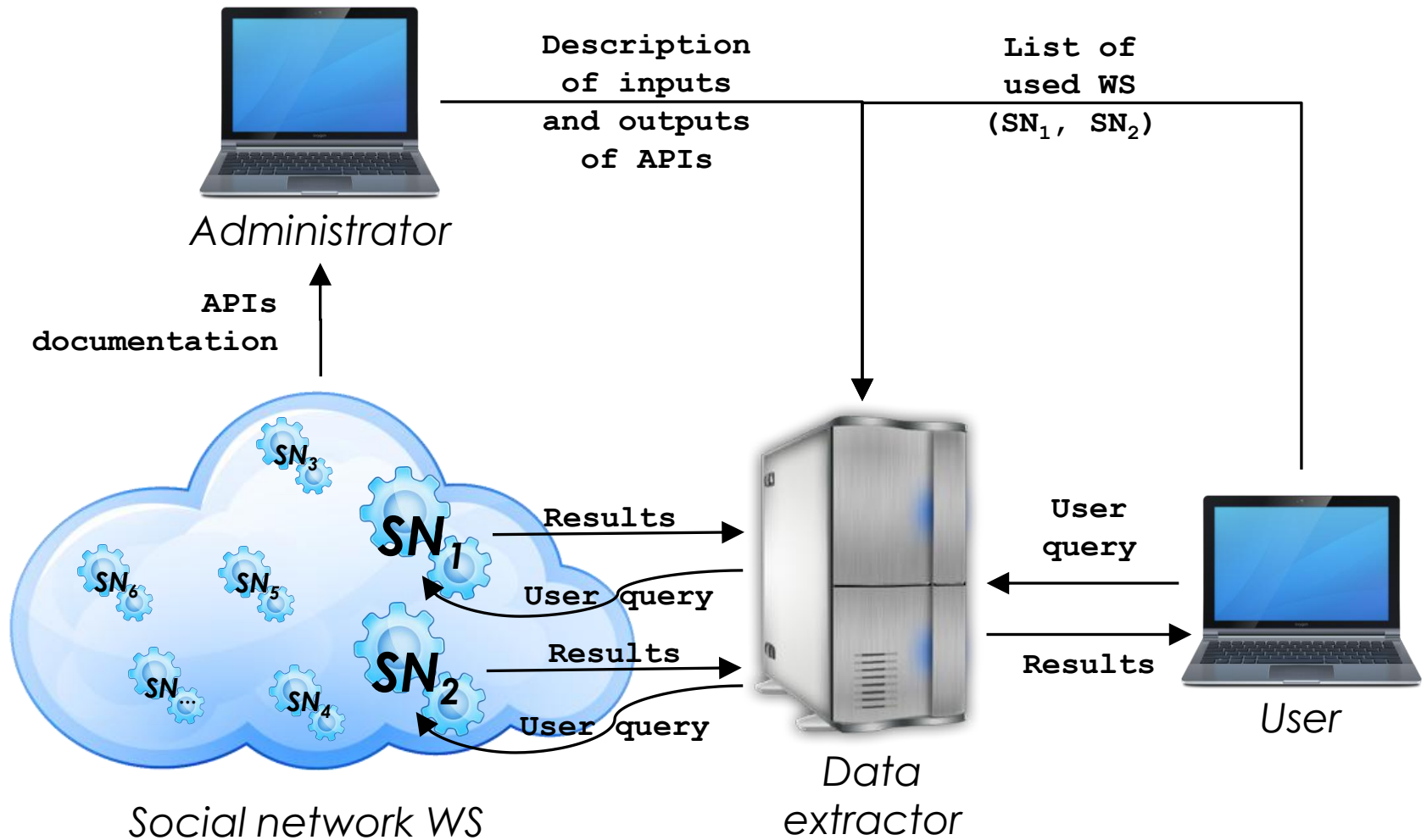
- Dependence on the used libraries, reasoning or engine  
→ declarative languages often are specialized (associated to a specific usage)
- The recursive way of thinking has to be practised.
- Handling complex data structures may be tricky.

# Two main usages

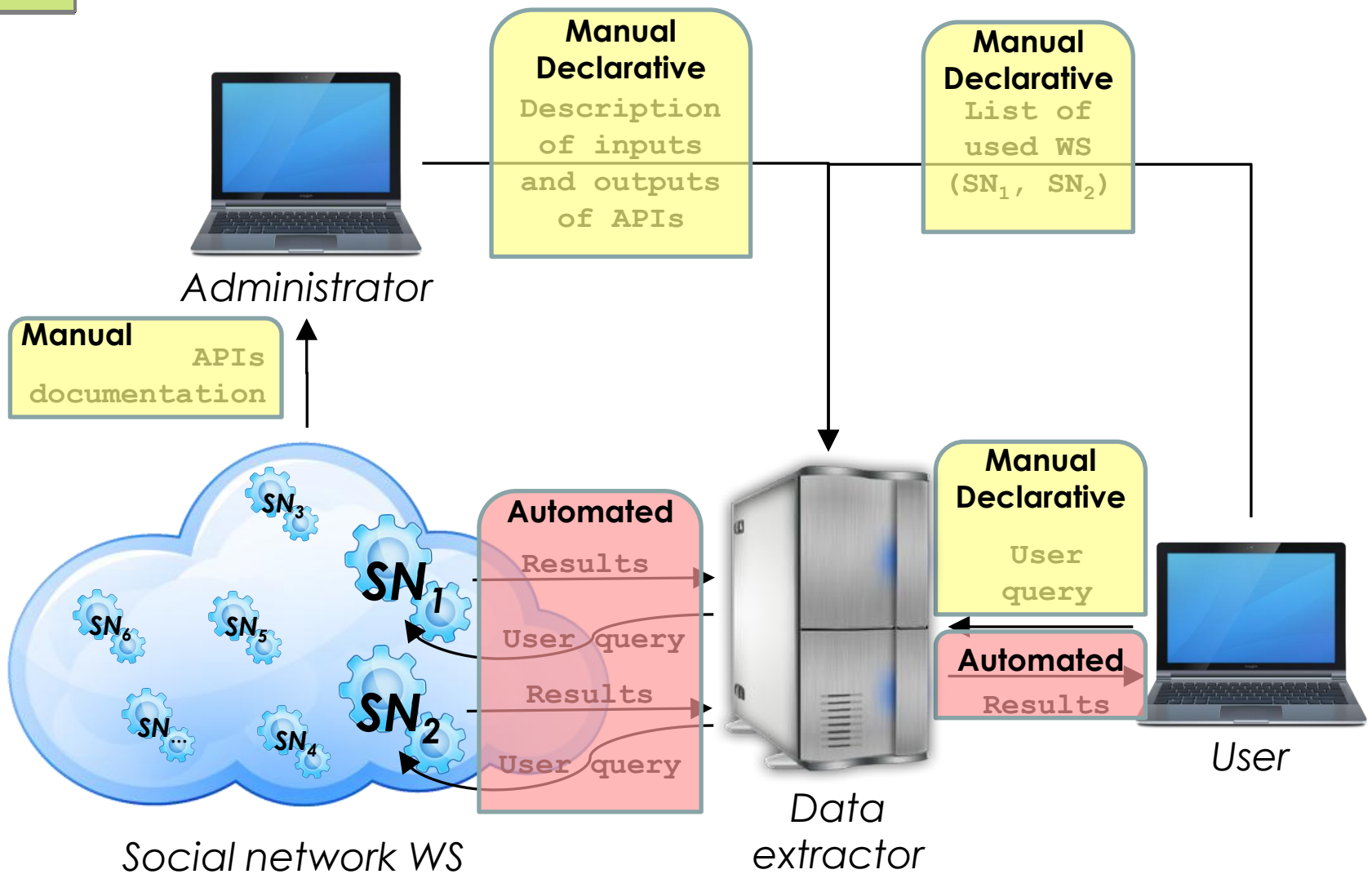
1. Building big programs
  - declarative software development
  - needs developers
2. Configuring an existing program or reasoning
  - Modifying parameters
  - Adding components, data
  - Expressing queries
  - Specifying presentation rules
  - declarative software administration
  - needs administrators

→ We choose the second option.

# Declarative data extractor



# Declarative data extractor



# Proposed workflow

1. Read the web service documentation
2. Declaratively set up the data extractor
3. Consuming data by querying
  1. Asking for data by querying
  2. Dealing with results via a database (DB)



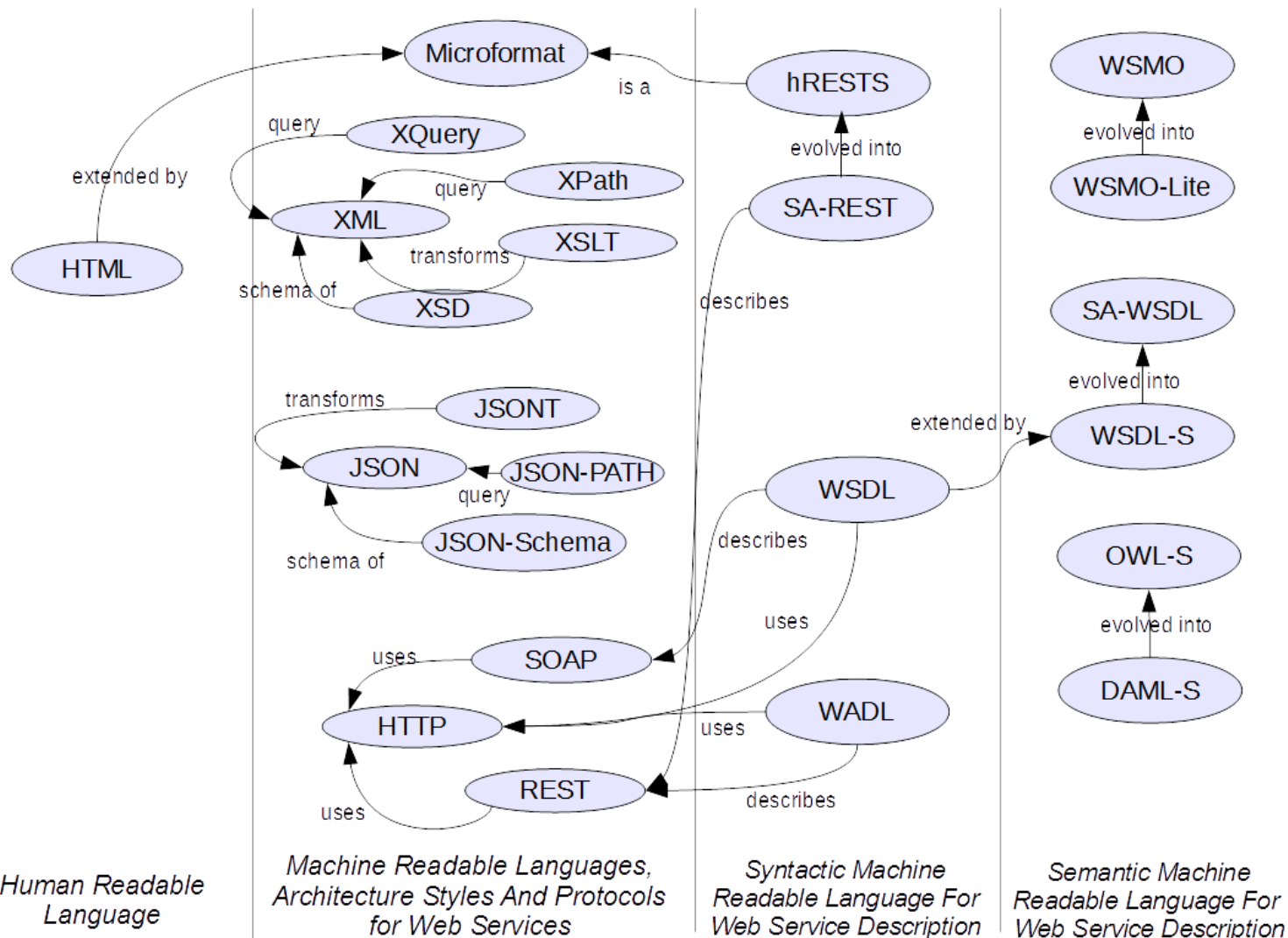
# Envisioned benefits

- Reduced administration time
  - A new WS = a new DB entry
  - A new WS version = a new DB entry
  - Api modifications = DB updates  
(e.g. change of input/output parameters)  
→ no need to code any more
- Estimated cost in time ?
  - Shifting to minutes/hours administration tasks (and not hours/days coding tasks)

# Which declarative languages ?

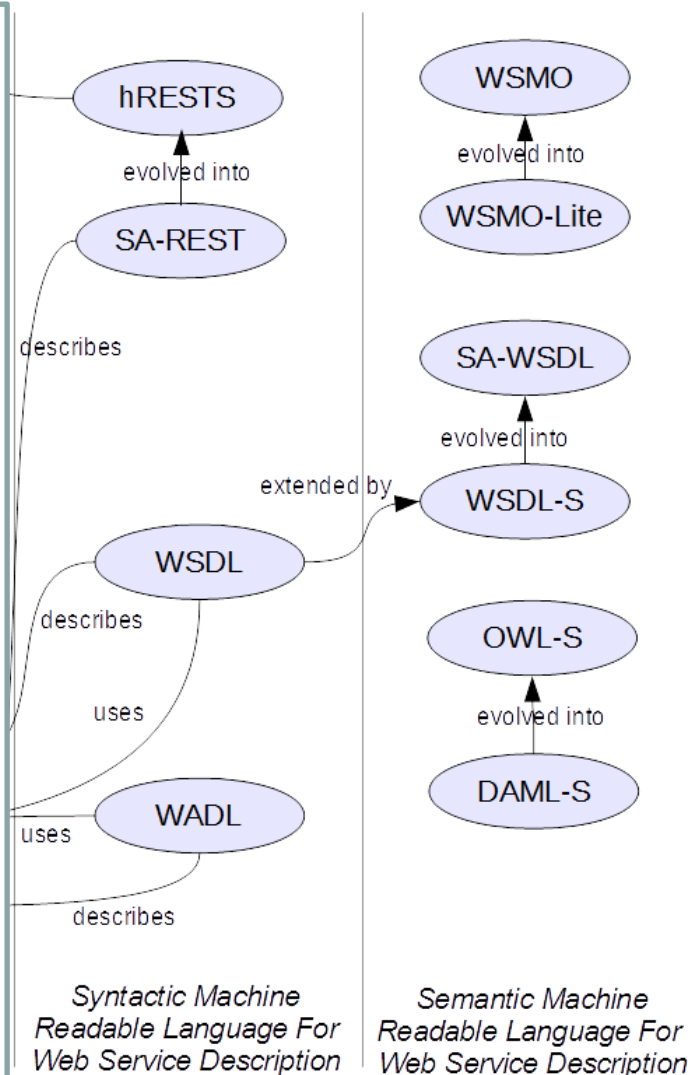
- Assume the idea is feasible
- Which declarative language ?
- What to do exactly ?
- Chosen declarative languages must be as much as possible
  - Simple
  - Well-known
  - Widespread  
(already used by many people)

# WS languages overview



# WS languages overview

Languages →  
are not  
widely used  
(yet).



# About WSDL

- Definition
  - Web Service Description Language
  - → possible to automatically generate WS wrappers
- Problems
  - WSDL-generated wrappers are pieces of (procedural) code to be inserted in a larger program → need a developer
  - Chaining API calls
    - One output becomes an input → API calls plans
    - Not easy to handle such plans
    - Solution: manage as query plans → declarative approach
  - WSDL is not widely used
    - → not declarative enough ?

# Which declarative language ?

- Previous candidates
  - not widely used
- What a pity !
  - ~~□ WSDL → automatically generated wrapper~~
  - ~~□ Semantic web languages → automatic discovery of web services~~
  - ~~□ ...~~

# In the real world

- Documentation
    - Natural language
    - Examples with code snippets
  - Message format
    - JSON
    - XML
- ➔ now, we must deal with this !

# Challenge

- Limit the manual work to
    - Reading APIs documentation
    - Declaratively describing APIs to the system (system set up)
    - Declaratively express user queries
  - Automate
    - Calls to WS APIs
    - User query processing
    - Results management
  - Allow a quite expressive user query language
- ➔ possible with widely used declarative languages ?



# [Declarative data extraction]

Getting into the details...

Studied issue

Locks and Keys

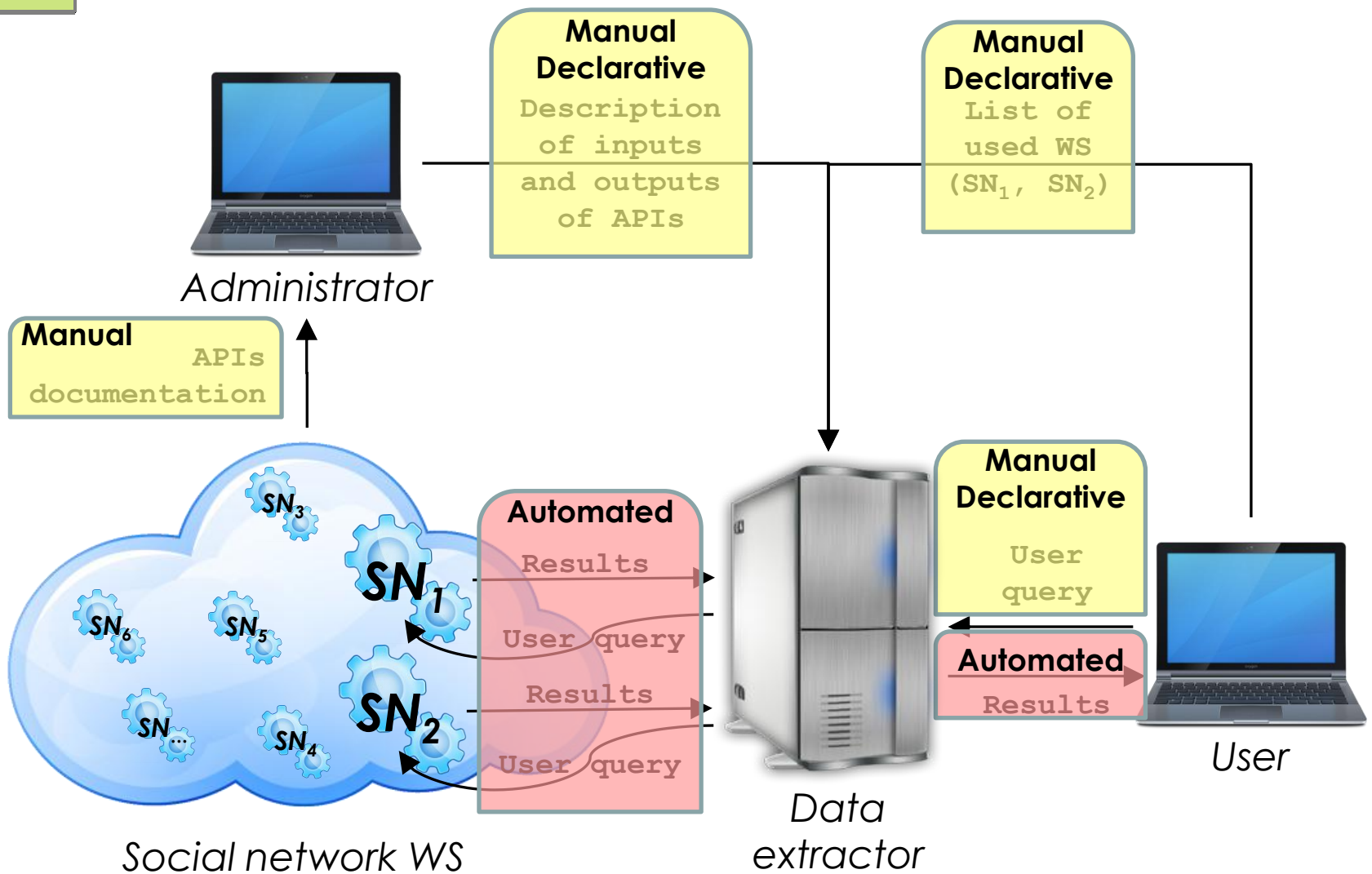
**Declarative data extraction**

Data integration with WS

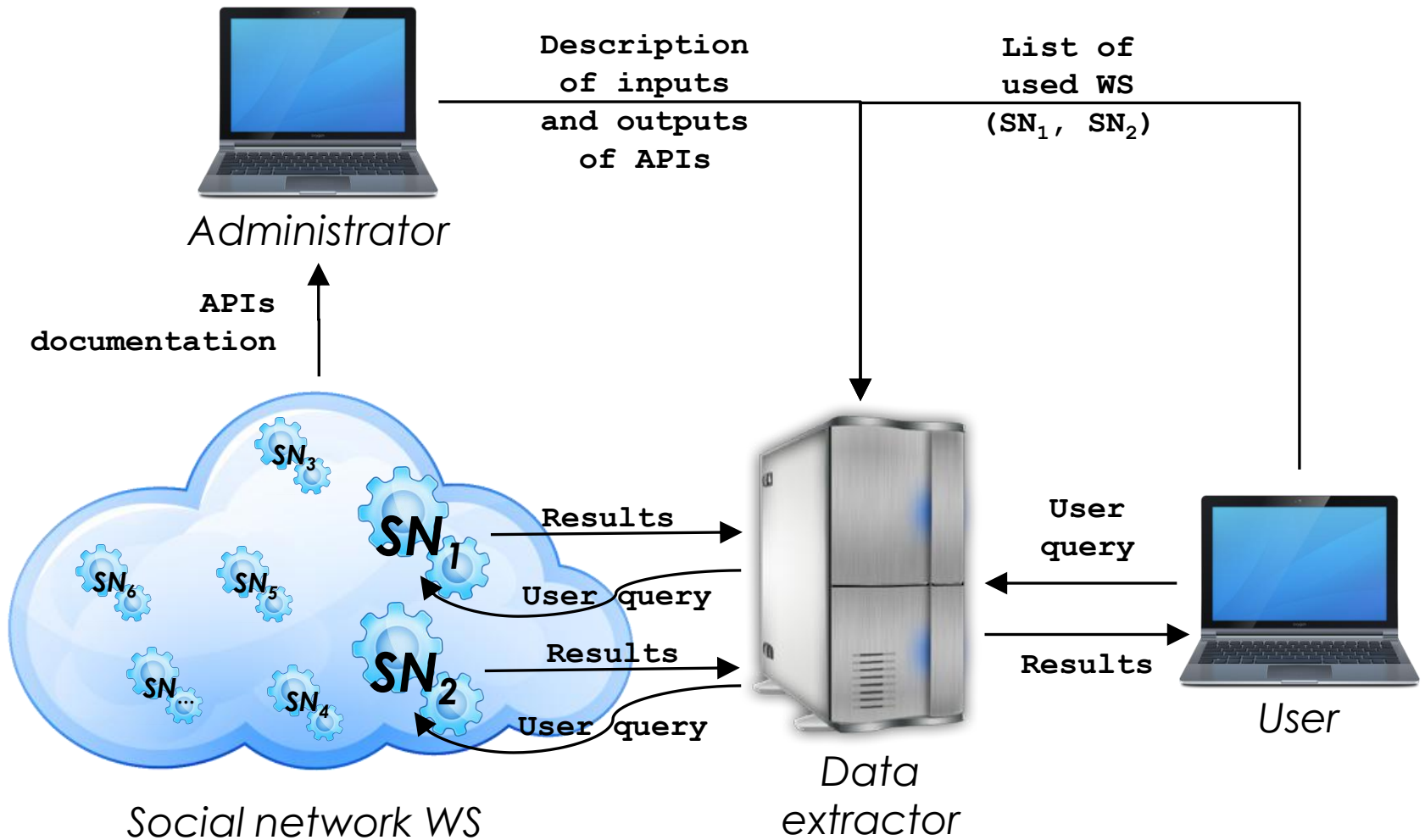
Main scenarios

Results and discussion

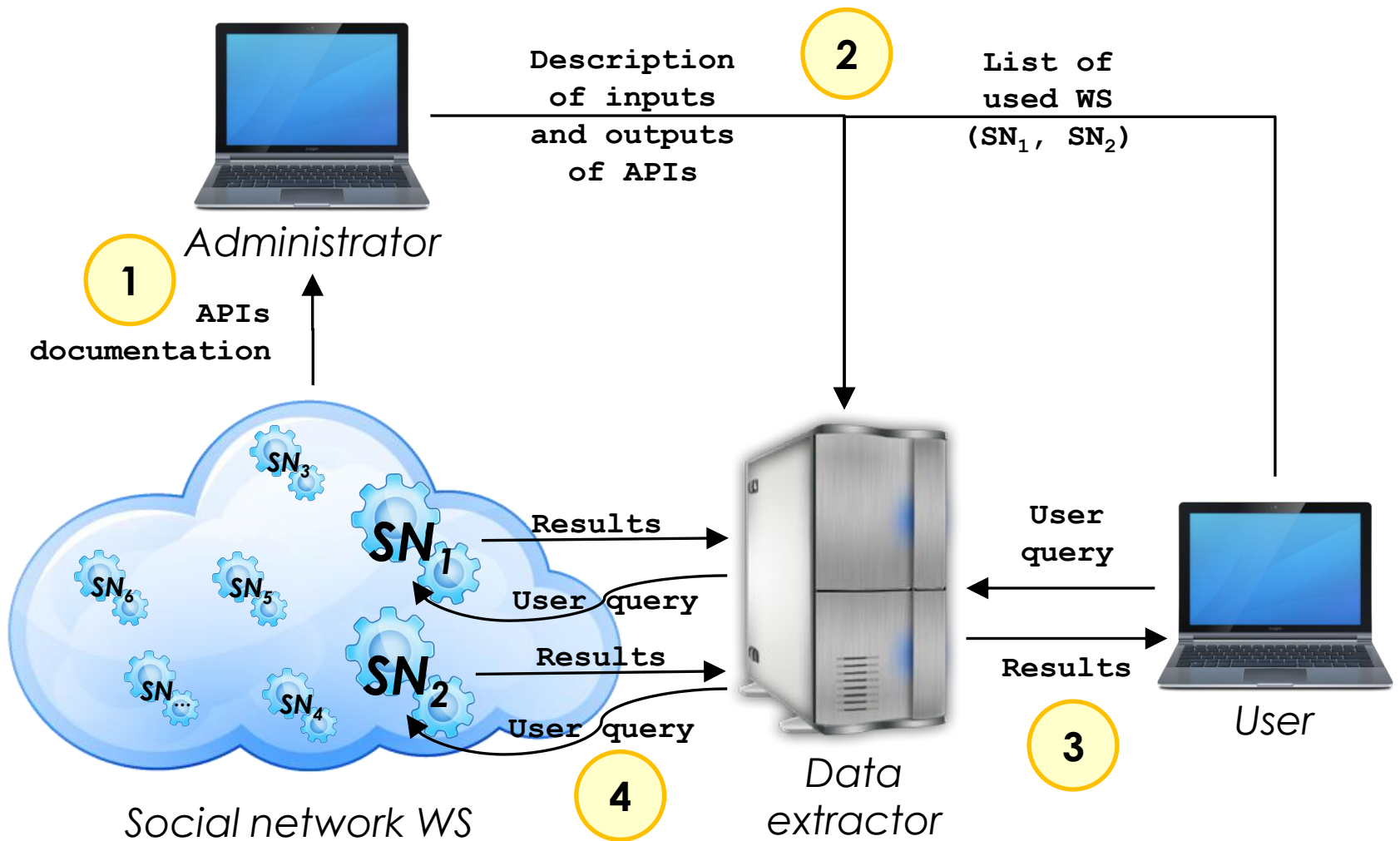
# Declarative data extractor



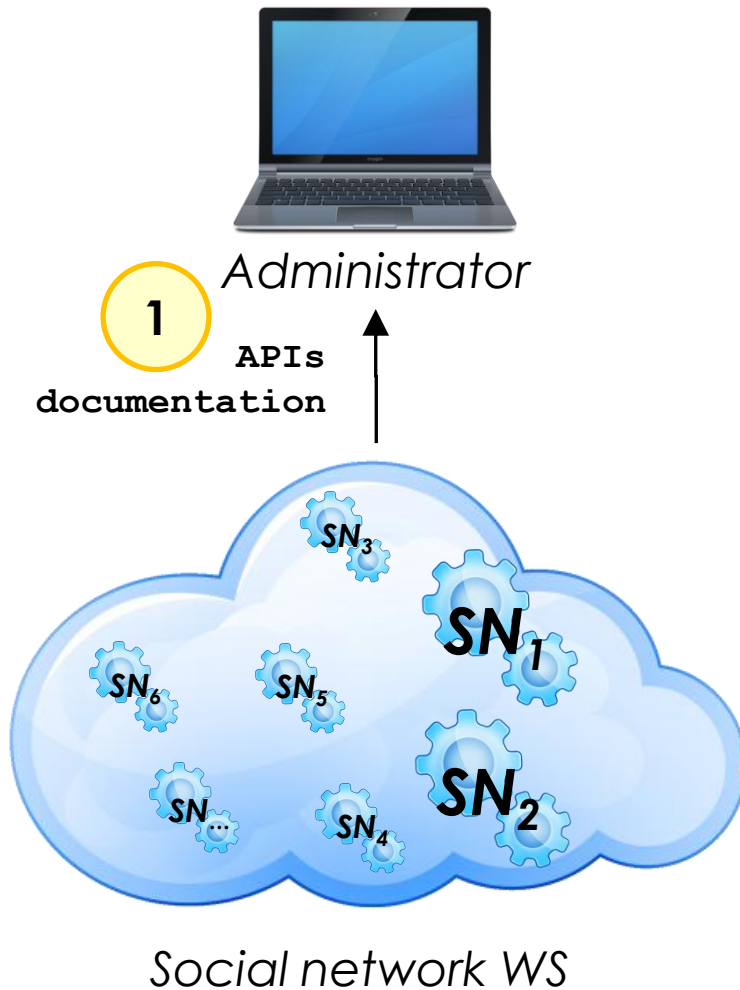
# Declarative data extractor



# 4 use cases



# Read the WS API doc.



- Determine operations
- Determine message and response formats
- Determine API authentication mechanisms
- For every operation, determine
  - the input and output parameters, and
  - required HTTP URL, header, method and body contents

# Inside the data extractor



*Administrator*

Description  
of inputs  
and outputs  
of APIs

2

List of  
used WS  
( $SN_1$ ,  $SN_2$ )

- Query and answer processor
  - Get queries
  - Process queries
  - Build answers
- Generic WS wrapper
  - Build API calls
  - Response and failure management

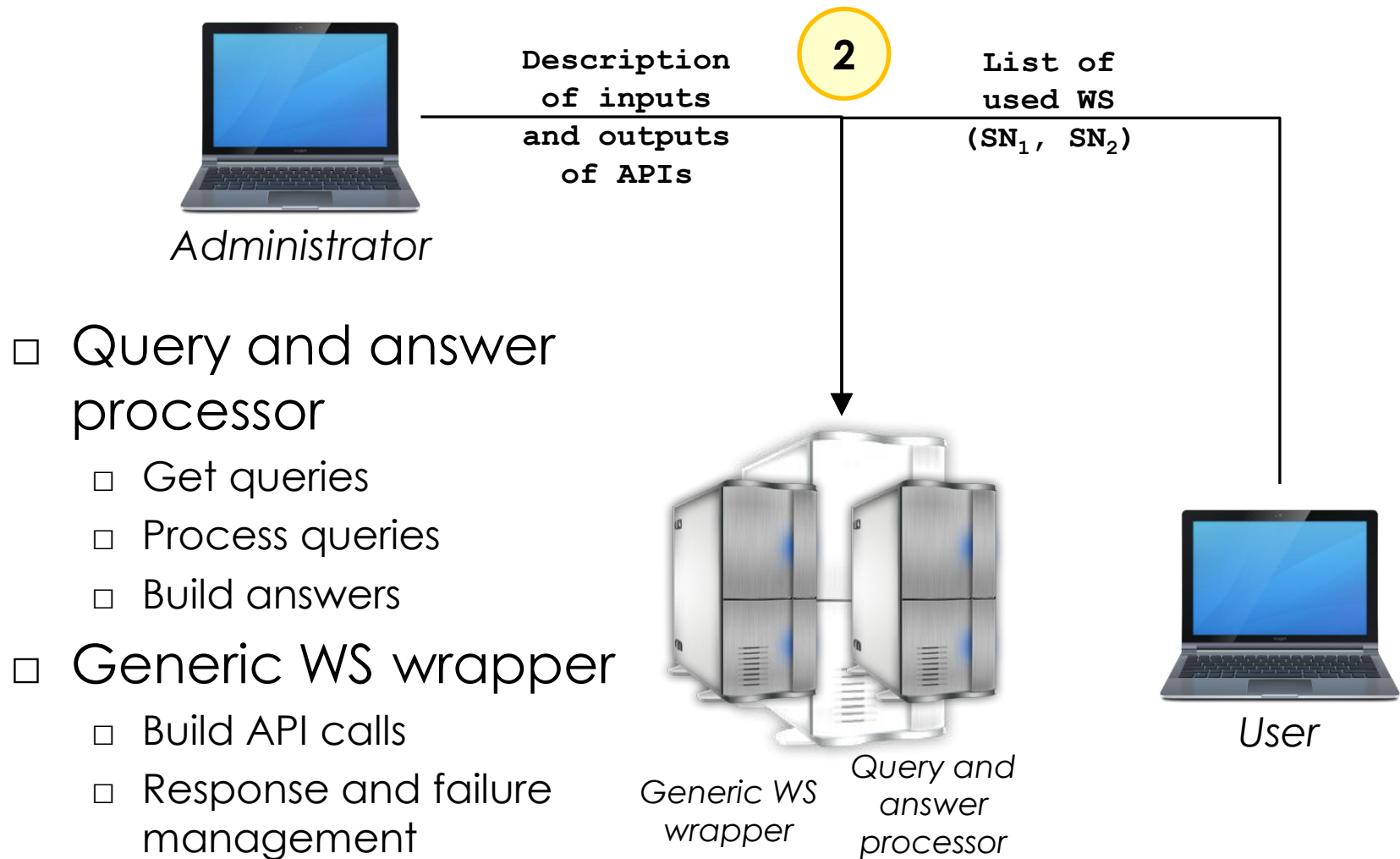


*Data  
extractor*



*User*

# Inside the data extractor



# Declarative set up



*Administrator*

- Describe WS API operation call (HTTP method, URL template, ...) → language ?
- Set up the Generic Wrapper
  - Building of operation calls
  - Processing of responses into chosen internal formats
- Set up HTTP body template
- Set up validation, extraction and transformation details required to build responses

Description  
of inputs  
and outputs  
of APIs

2

List of  
used WS  
( $SN_1$ ,  $SN_2$ )



*Generic WS  
wrapper*

*Query and  
answer  
processor*

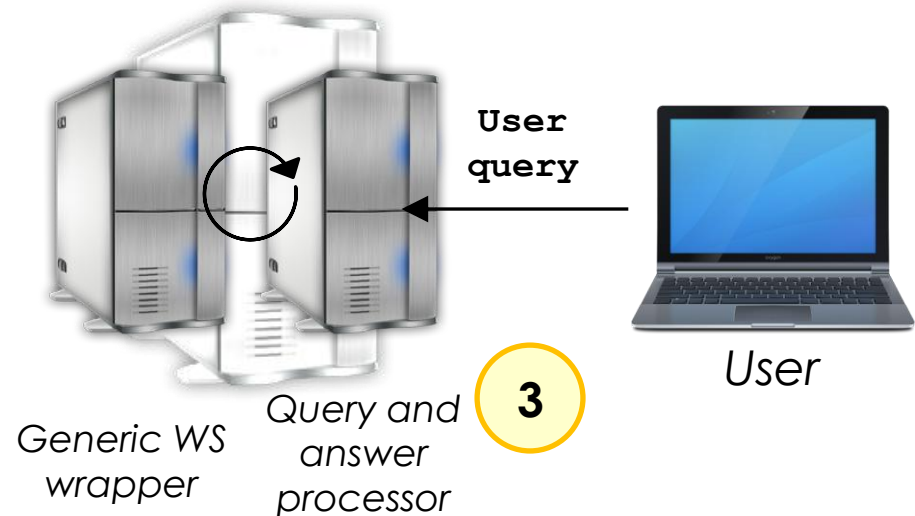


*User*



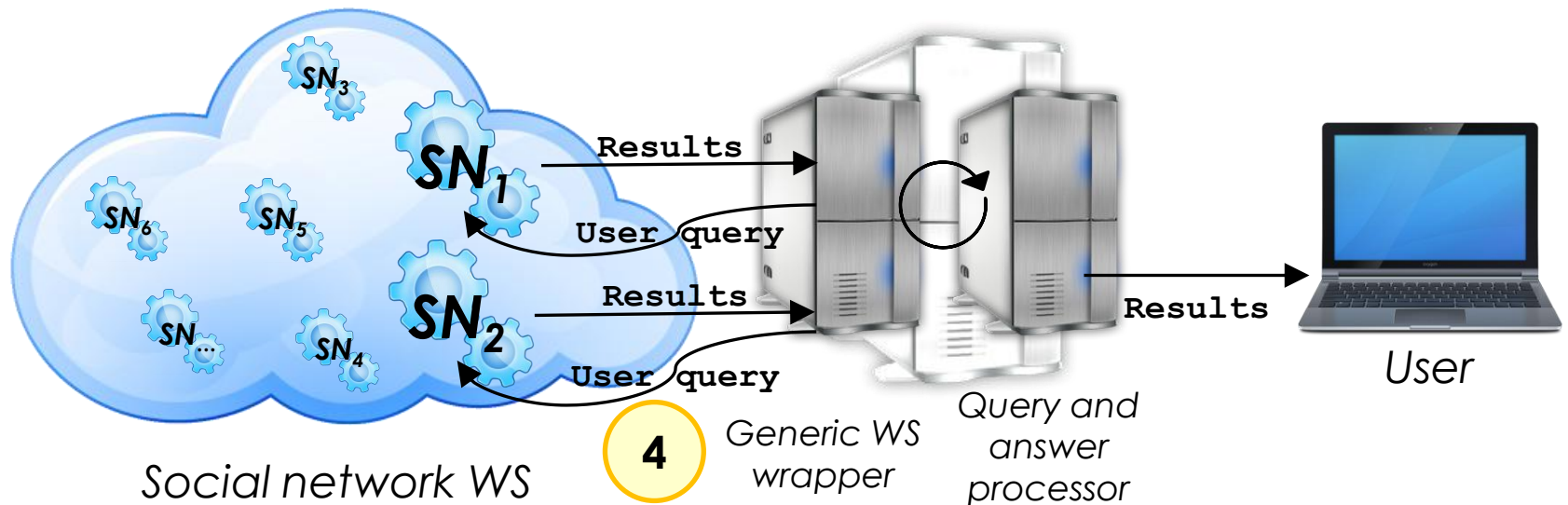
# Consuming data by querying

- ❑ Querying is declarative !
- ❑ Query language : SQL ? Other ?
- ❑ Internal formats : ?

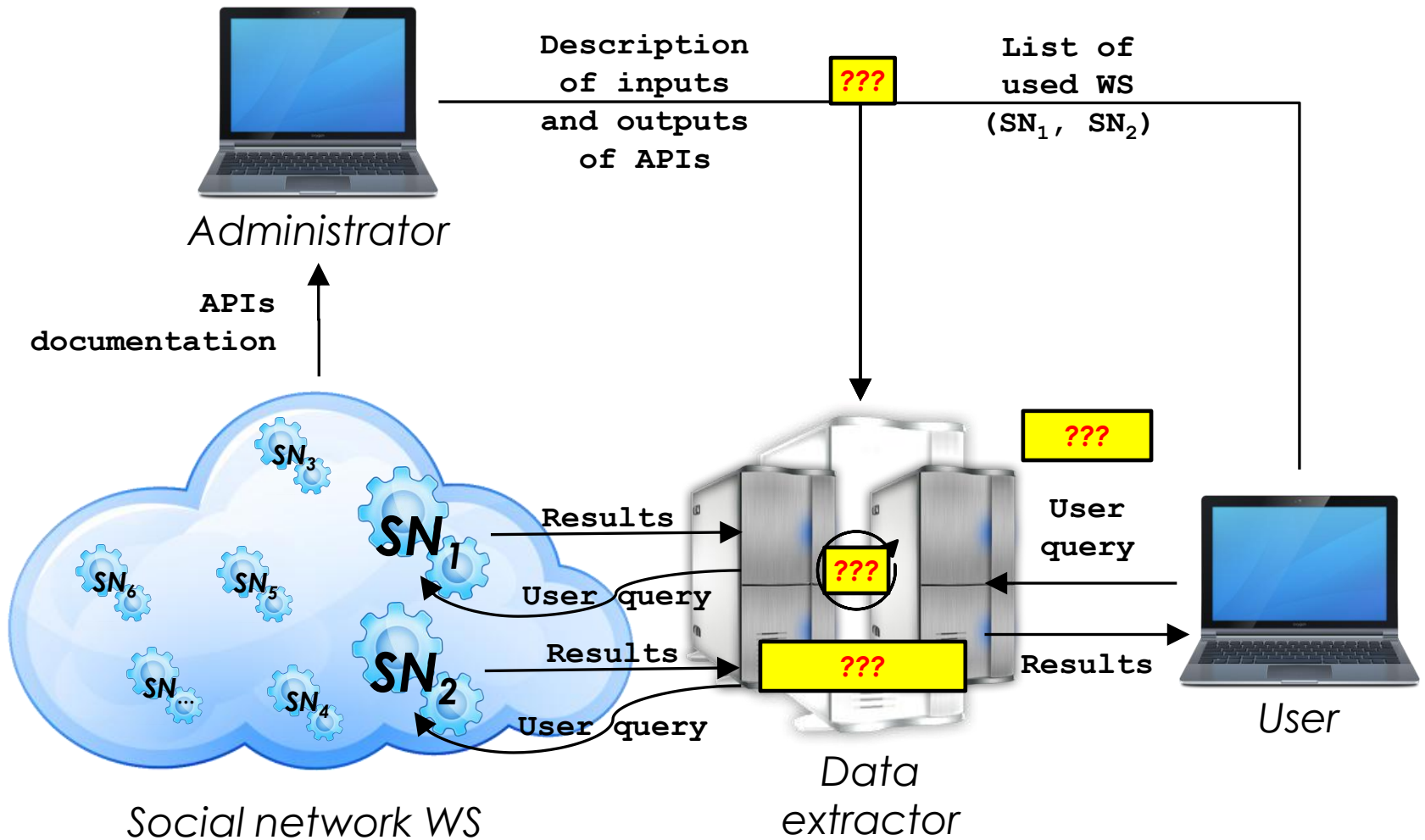


# Automated querying processing

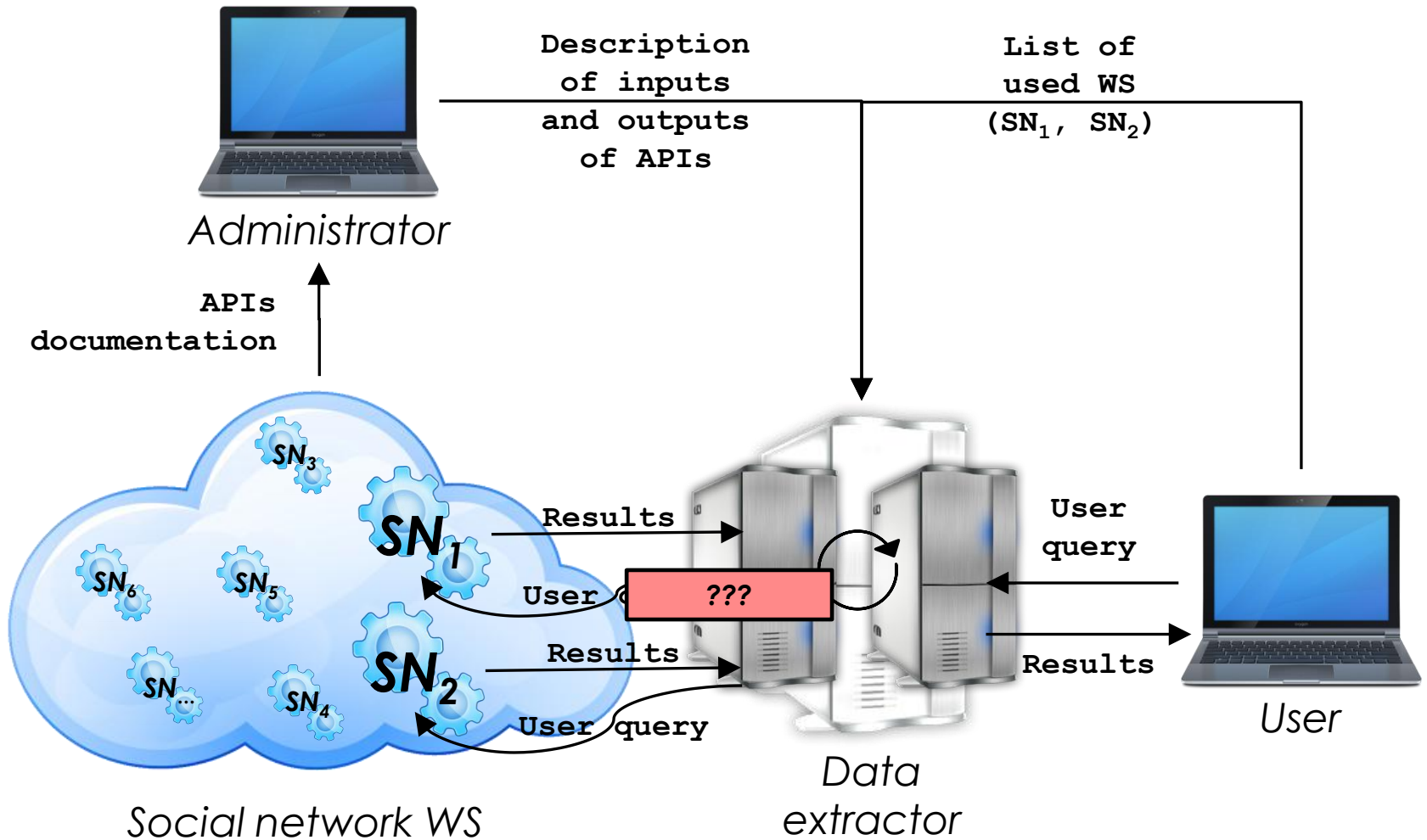
- How to automatically build the API calls ?
- How to deal with responses and potential failures ?
- Response processing formats ?



# Languages ?



# Automated processing ?



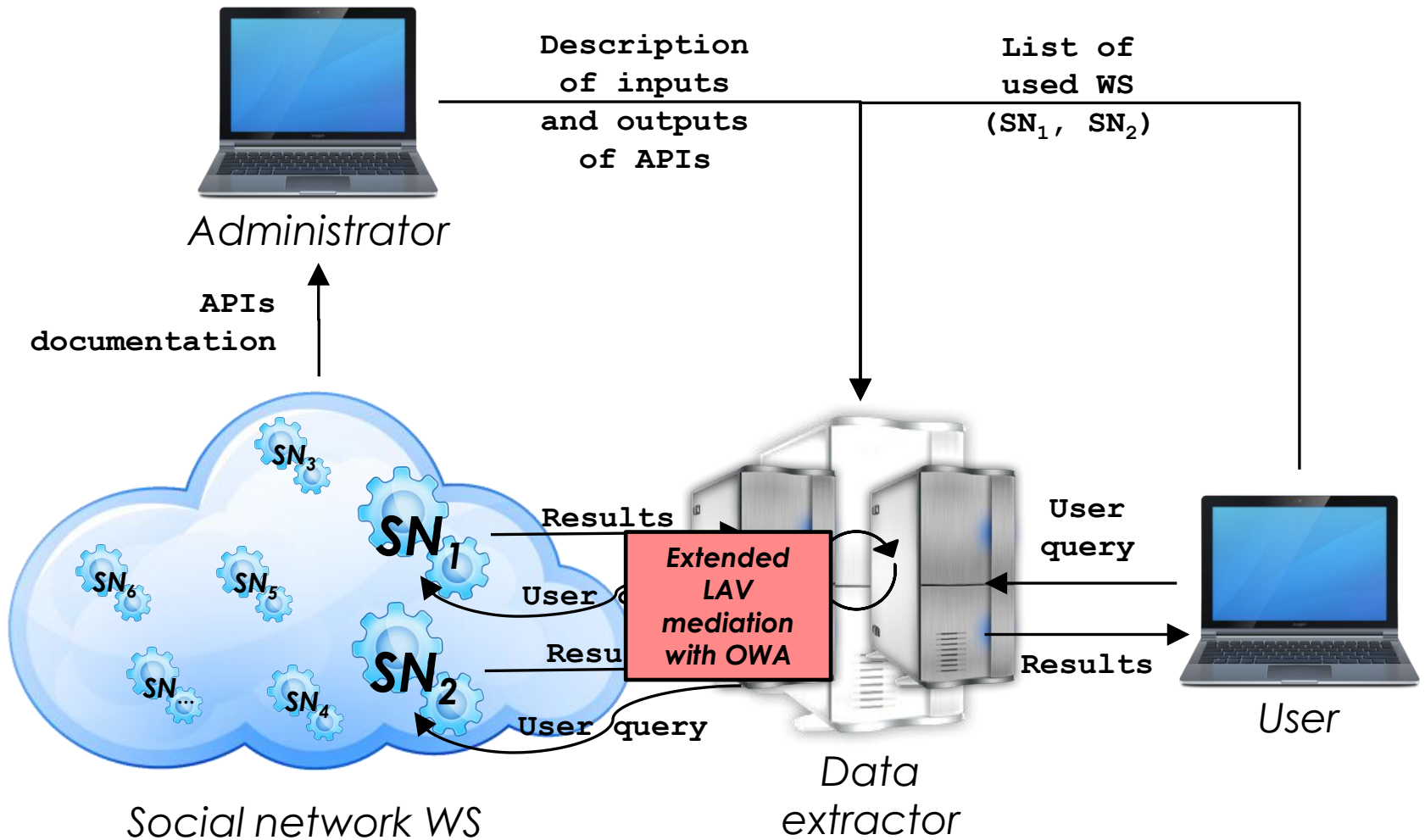
# A data integration problem

- WS = data sources
- WS operations = relations with access patterns  
(some attributes are inputs, some are outputs)
- One user query to be posed to many WS operations via WS APIs.
- Constraints:
  - Declarative links between WS and the data extractor
  - Use of widely spread declarative languages

# Followed approach

- Basis : mediation with LAV mappings and OWA.
- Use of the well-known Inverse Rules algo
  - Access patterns
  - Full and functional dependencies
- Mappings are CQ with access patterns
- User queries are datalog queries.
- Extension
  - To deal with incomplete information (heuristics)
  - To handle pagination

# Followed approach

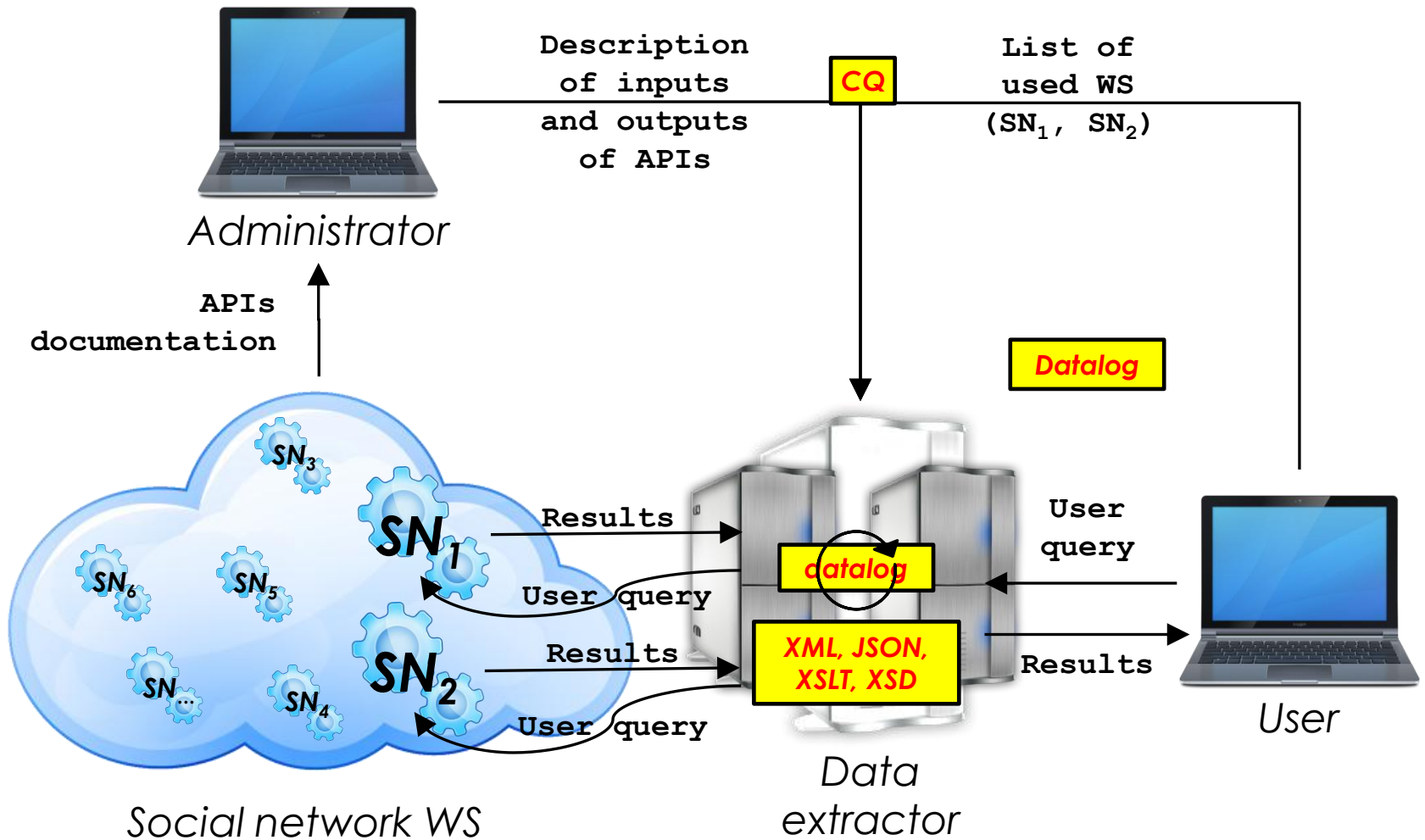


# Declarative languages used in this work

- Presentation languages/ message formats → those which are used
  - JSON / XML
  - XSD
  - XSLT
- Query languages
  - Conjunctive queries
  - Datalog



# Back to the extractor



# Conjunctive queries CQ

- Equivalent to
  - select-project-join queries
  - i.e. SELECT... FROM ... WHERE ... ; in SQL
- Example
  - Relational schema:  $\{r1(A,B), r2(C,D), r3(E,F,G)\}$
  - $\text{ans}(X) \leftarrow r1(X,Y), r2(Y,Z), r3(Z,X,W).$
  - SELECT r1.A  
FROM r1, r2, r3  
WHERE (r1.B = r2.C) and (r2.D = r3.E) and (r1.A=r3.F);

# Datalog queries

- Equivalent to
  - Recursive unions of conjunctive queries
- Example
  - Transitive closure of a graph
  - Relational schema:  $\{\text{edge}(X,Y)\}$
  - $\text{path}(X,Y) \leftarrow \text{edge}(X,Y).$
  - $\text{path}(X,Y) \leftarrow \text{edge}(X,Z), \text{path}(Z,Y).$

# JSON

```
1 {  
2   "updates": [{  
3     "identifier": 1245,  
4     "update": "I hope the weather is good",  
5     "date": "2016-03-04T10:20:18+02:00",  
6     "count": 14  
7   }, {  
8     "identifier": 1246,  
9     "update": "It's sunny here",  
10    "date": "2016-03-04T14:25:16+02:00",  
11    "count": 14  
12  }]  
13 }
```

# XML

```
<?xml version="1.0" encoding="UTF-8"?>  
<xml>  
  <viewcount>  
    <count>123</count>  
  </viewcount>  
</xml>
```

# XSD

```
<xs:schema attributeFormDefault="unqualified"  
  elementFormDefault="qualified" xmlns:xs="http://www.w3.org/2001/XMLSchema">  
  <xs:element name="xml">  
    <xs:complexType>  
      <xs:sequence>  
        <xs:element name="viewcount">  
          <xs:complexType>  
            <xs:sequence>  
              <xs:element type="xs:byte" name="count"/>  
            </xs:sequence>  
          </xs:complexType>  
        </xs:element>  
      </xs:sequence>  
    </xs:complexType>  
  </xs:element>  
</xs:schema>
```

# XSLT

```
<?xml version="1.0" encoding="UTF-8"?>
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform" version="1.0">
  <xsl:template match="/">
    <xsl:for-each select="statuses/status">
      <xsl:value-of select="identifier" />
      ,
      <xsl:value-of select="text" />
      <xsl:text />
    </xsl:for-each>
  </xsl:template>
</xsl:stylesheet>
```

# [Data integration with WS]

Studied issue

Locks and Keys

Declarative data extraction

**Data integration with WS**

Main scenarios

Results and discussion



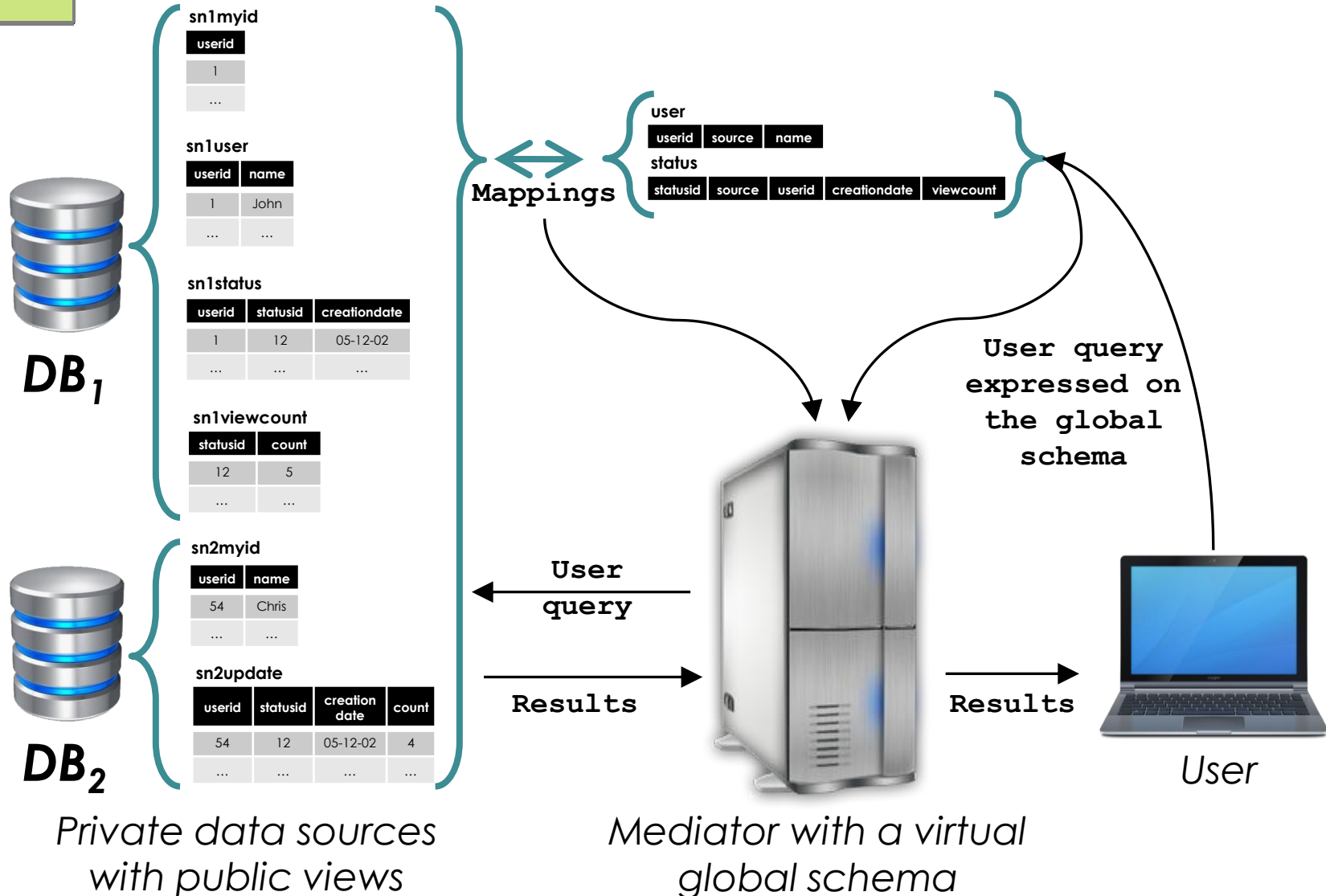
# Data integration approaches

- Materialized → data warehousing
  - Source data are stored locally and maintained over time
  - Main purposes: statistics, decision support
  - Analytical queries
- Virtual → mediation
  - Source descriptions only (no data) are stored locally
  - Main purpose: « lightweight » heterogeneous sources querying
  - Transactional queries (no modification)

# Dealing with WS

- Access patterns
  - ➔ complete data unavailable
- Great number of WS
- High versatility
  - Data keep on changing
  - WS keep on evolving
- Networking costs
- Transactional « get » queries only
  - ➔ The mediation data integration approach fits well.

# Mediation in a nutshell



# Mediation in a nutshell

- Main problem:  
(User) query answering using (data from the) views
- Certain answers semantics
- Depend on the mappings:
  - GAV mappings: « Global As View »
  - LAV mappings: « Local As view »
  - GLAV mappings: « Global and Local As View »

# Mediation in a nutshell

- GAV mappings:
  - Each global schema predicate = a query on the schema made up with all views (i.e. the source schema).
  - Query answering by unfolding the user query:  
replace global schema predicates by their definition
- LAV mappings:
  - Each view = a query on the global schema.
  - Query answering by query rewriting:
    - The user query is rewritten (if possible) into one or more query rewritings.
    - Each rewriting is a query plan that is answered by calling one or many views.
- GLAV mappings:
  - Sort of generalized LAV case
  - Link (inclusion) between a query on the global schema and a query on the views schema
  - Query answering by query rewriting

# Examples of mappings



**DB<sub>1</sub>**

sn1myid

userid
1
...

sn1user

userid	name
1	John
...	...

sn1status

userid	statusid	creationdate
1	12	05-12-02
...	...	...

sn1viewcount

statusid	count
12	5
...	...



**DB<sub>2</sub>**

sn2myid

userid	name
54	Chris
...	...

sn2update

userid	statusid	creation date	count
54	12	05-12-02	4
...	...	...	...

Mappings

user

userid source name

status

statusid source userid creationdate viewcount

- LAV mapping

**sn1user(userid,name) ← user(userid,'SN1',name).**

- GAV mapping

**user(userid,'SN1',name) ← sn1user(userid,name),  
sn1myid(userid).**

**user(userid,'SN2',name) ← sn2myid(userid,name).**

- GLAV mapping

**user(userid,source,name), ∃ status(statusid,source,userid,creationdate,viewcount) ← sn1user(userid,name), sn1status(userid,statusid,creationdate), sn1viewcount(statusid,count).**

*Private data sources  
with public views*

# Mediation in a nutshell

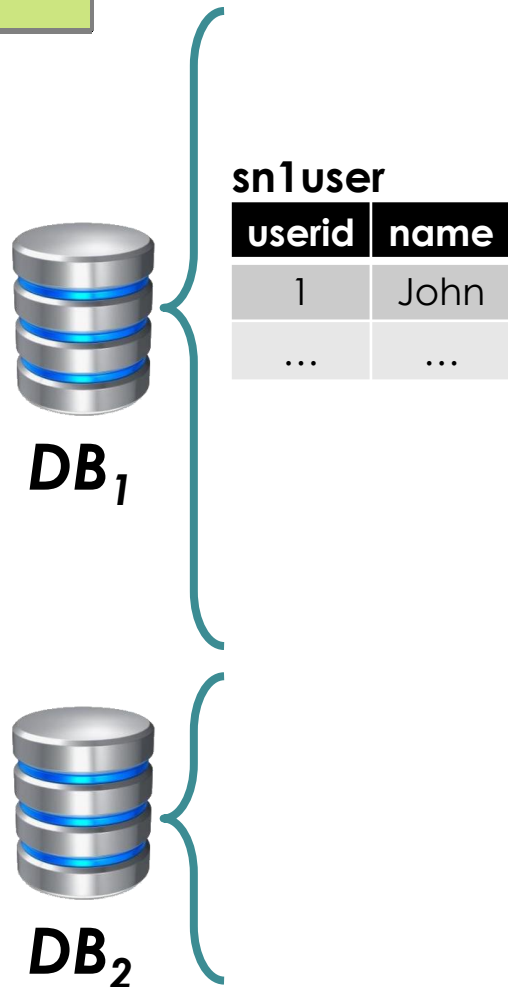
- Pros and cons GAV
  - (+) Query processing expected to be easier
  - (-) Defining global schema not always an easy task
  - (-) Adding new sources impacts the global schema
- Pros and cons LAV
  - (+) Modularity/Scalability (e.g., adding new sources do not impact the global schema)
  - (-) Query processing is more complex

# Mediation in a nutshell

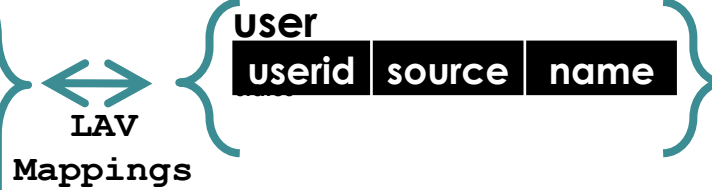
- Pros and cons GLAV
  - (+) Modularity/Scalability (e.g., adding new sources do not impact the global schema)
  - (+) Highly expressive mappings
  - (-) Query processing is more complex
  - (-) Handle recursive queries ?
  
- Conclusion
  - ➔ LAV best suited for WS integration



# LAV mappings



Private data sources  
with public views



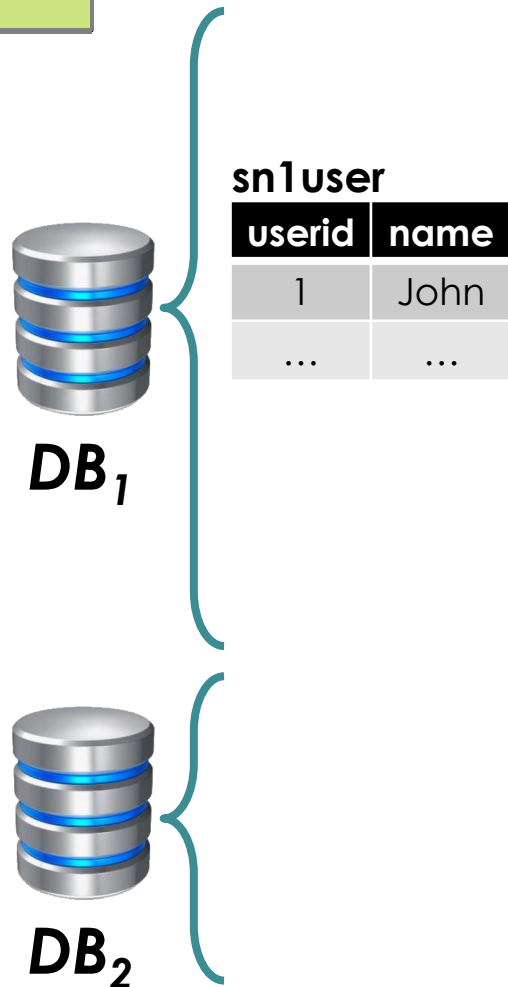
Mediator's  
virtual global  
schema GS

- Example

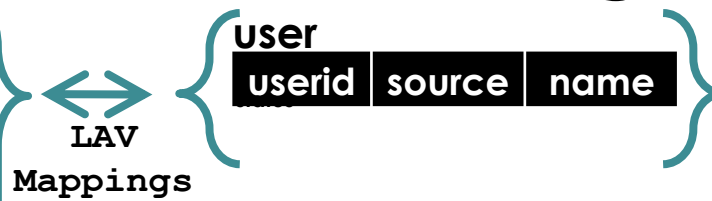
```
sn1user(userid,name) ← user(userid,'SN1',name).
```

- If « GS was materialized », then:  
couples (userid,name) of sn1user would be in triples (userid,'SN1',name) from user
- A conjunctive query
- Declarative !

# LAV mappings



Private data sources  
with public views



Mediator's  
virtual global  
schema *GS*

## □ Example

```
sn1user(userid,name) ← user(userid,'SN1',name).
```

## □ Assumptions

- Sound views  
« all couples from sn1user are included in triples of user »
- Incomplete views  
« a couple that is not in sn1user may be true or false (we don't know) »  
→ a.k.a. the open world assumption (OWA)

→ LAV with OWA

# LAV mappings

Mediator's virtual global schema GS



**DB<sub>1</sub>**

sn1myid

userid
1
...

sn1user

userid	name
1	John
...	...

sn1status

userid	statusid	creationdate
1	12	05-12-02
...	...	...

sn1viewcount

statusid	count
12	5
...	...



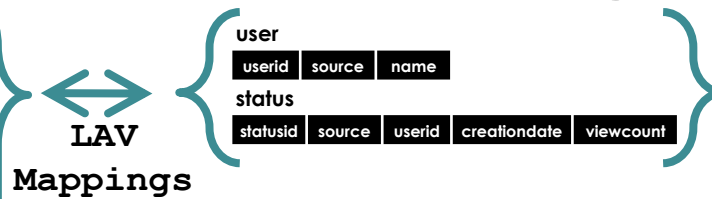
**DB<sub>2</sub>**

sn2myid

userid	name
54	Chris
...	...

sn2update

userid	statusid	creation date	count
54	12	05-12-02	4
...	...	...	...

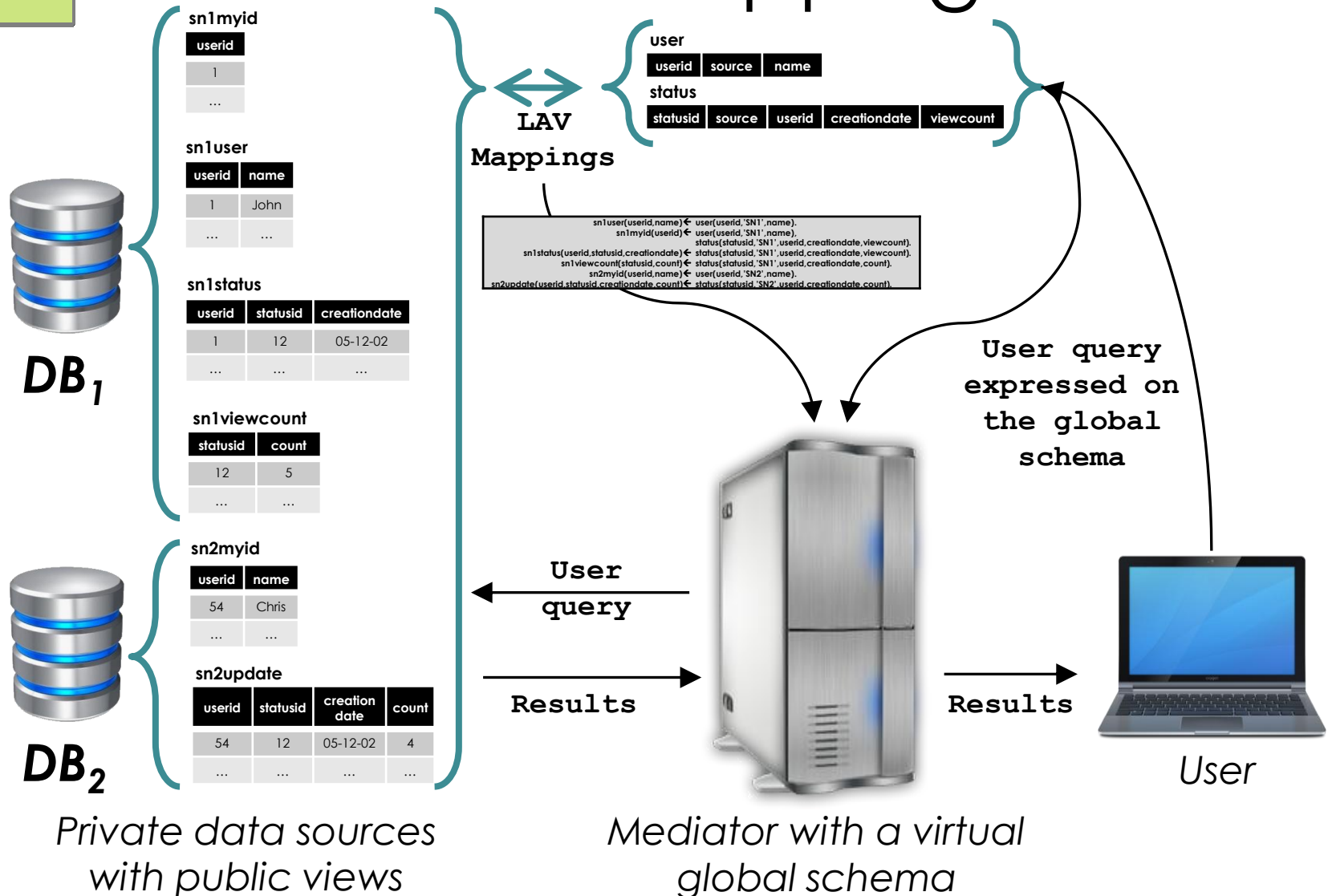


```

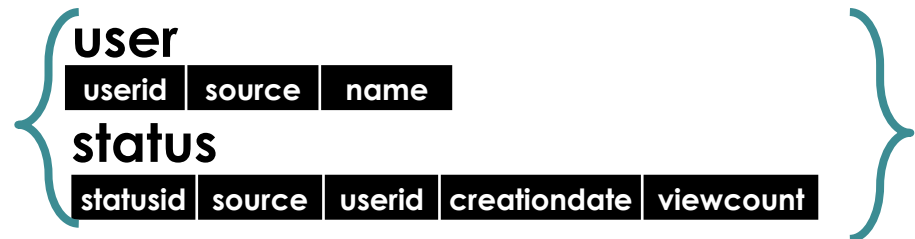
sn1user(userid,name) ← user(userid,'SN1',name).
sn1myid(userid) ← user(userid,'SN1',name),
                 status(statusid,'SN1',userid,creationdate,viewcount).
sn1status(userid,statusid,creationdate) ← status(statusid,'SN1',userid,creationdate,viewcount).
sn1viewcount(statusid,count) ← status(statusid,'SN1',userid,creationdate,count).
sn2myid(userid,name) ← user(userid,'SN2',name).
sn2update(userid,statusid,creationdate,count) ← status(statusid,'SN2',userid,creationdate,count).
    
```

Private data sources with public views

# LAV mappings



# User query



- Example 1

```
q(userid, statusid, viewcount) ← user(userid, 'SN1', name),  
                                status(statusid, 'SN1', userid, '2016-03-03', viewcount)
```

- Give me all triples  
(userid,statusid,viewcount)  
which creation date is 2016-03-03  
and source is 'SN1'.
- ➔ Conjunctive query  
(subset of datalog)

# User query

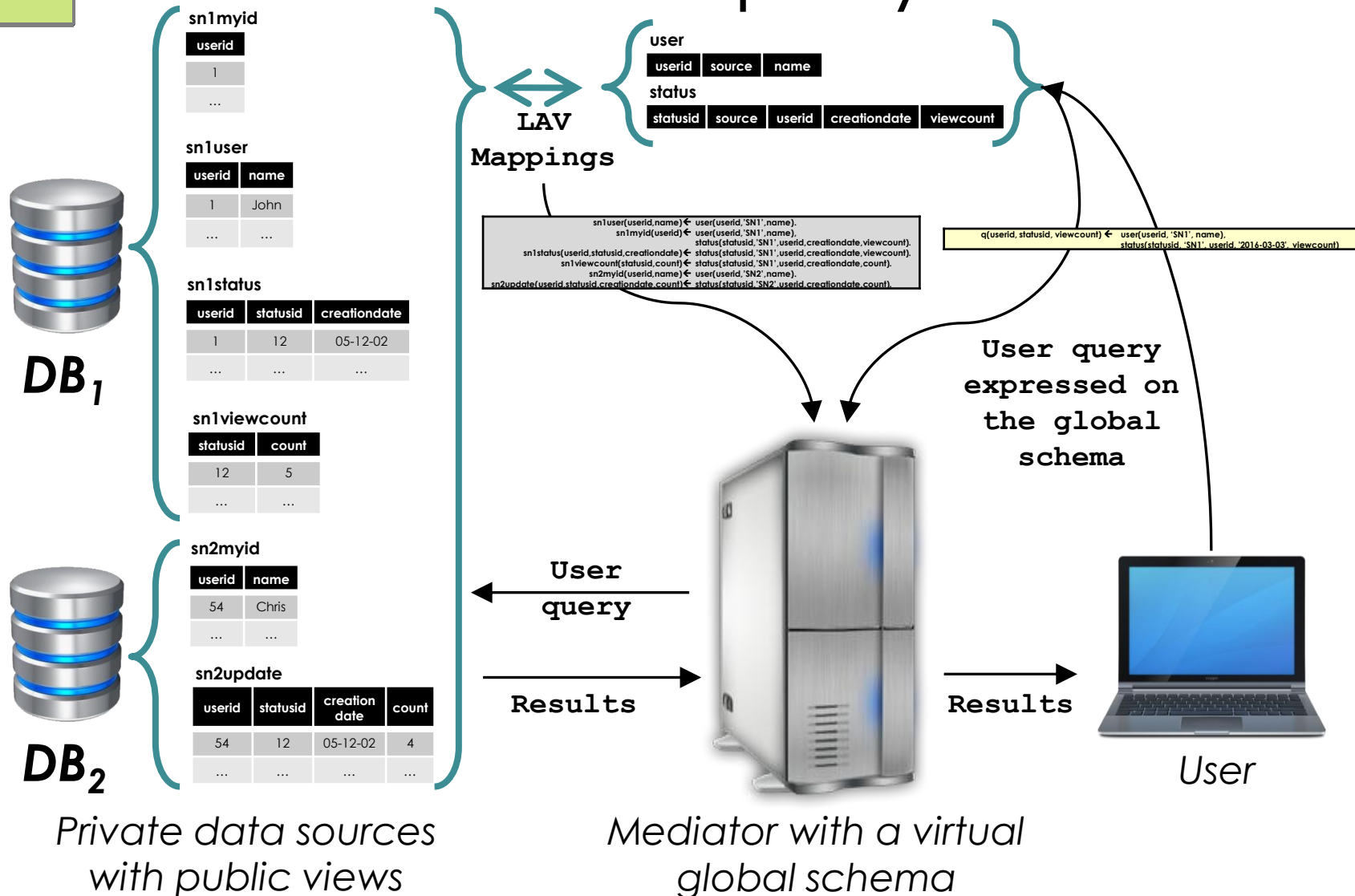


## □ Example 2

$q(\text{statusid}, \text{source}, \text{viewcount}) \leftarrow \text{status}(\text{statusid}, \text{source}, \text{userid}, '2016-03-03', \text{viewcount})$   
 $q(\text{statusid}, \text{source}, \text{viewcount}) \leftarrow \text{status}(\text{statusid}, \text{source}, \text{userid}, '2016-03-04', \text{viewcount})$   
 $q(\text{statusid}, \text{source}, \text{viewcount}) \leftarrow \text{status}(\text{statusid}, \text{source}, \text{userid}, '2016-03-05', \text{viewcount})$

- Give me all triples  
(statusid,source,viewcount)  
which creation date is 2016-03-03, or 2016-03-04 or 2016-03-05.
- ➔ Union of conjunctive queries  
(subset of datalog)

# User query

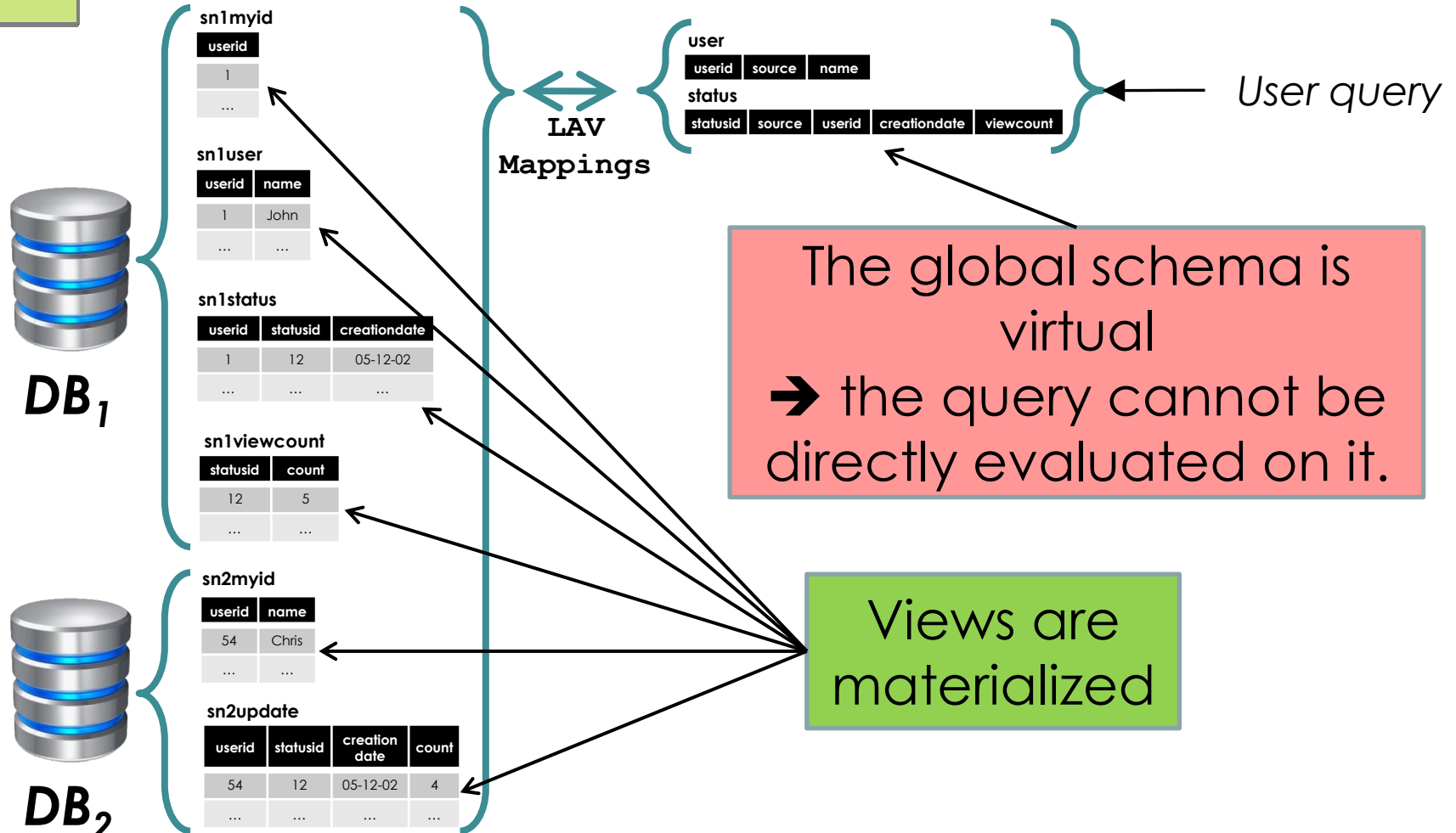


# Query answering using views

- Assume
  - LAV mappings with OWA
  - A user query posed to the global schema
- Certain answers semantics [AD98]
  - Finding « certain » answers
  - i.e. answers that would be obtained, whatever materialization of the global schema consistent with the LAV mappings with OWA, if the user query was directly evaluated on the global schema

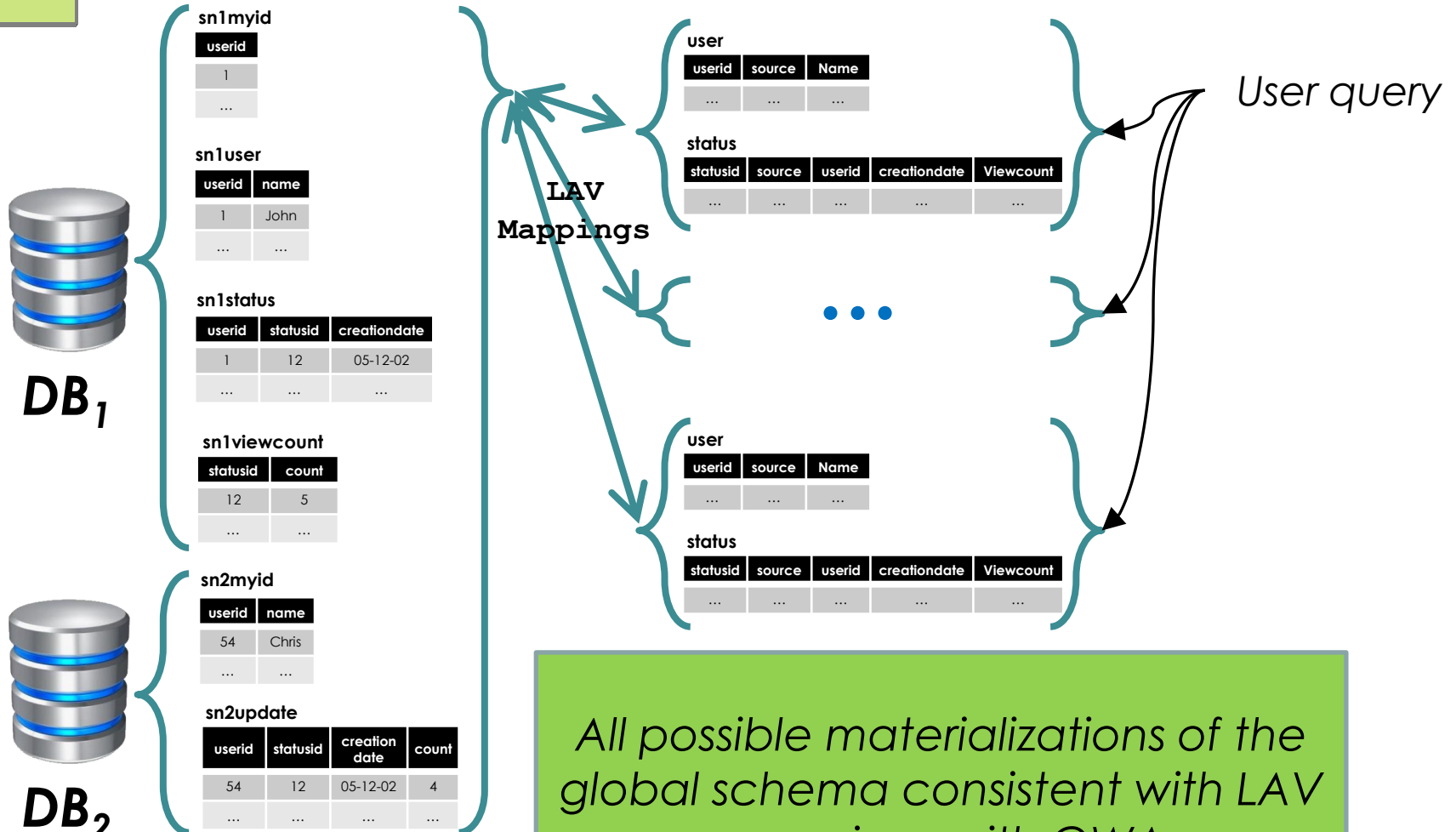


# Certain answers semantics



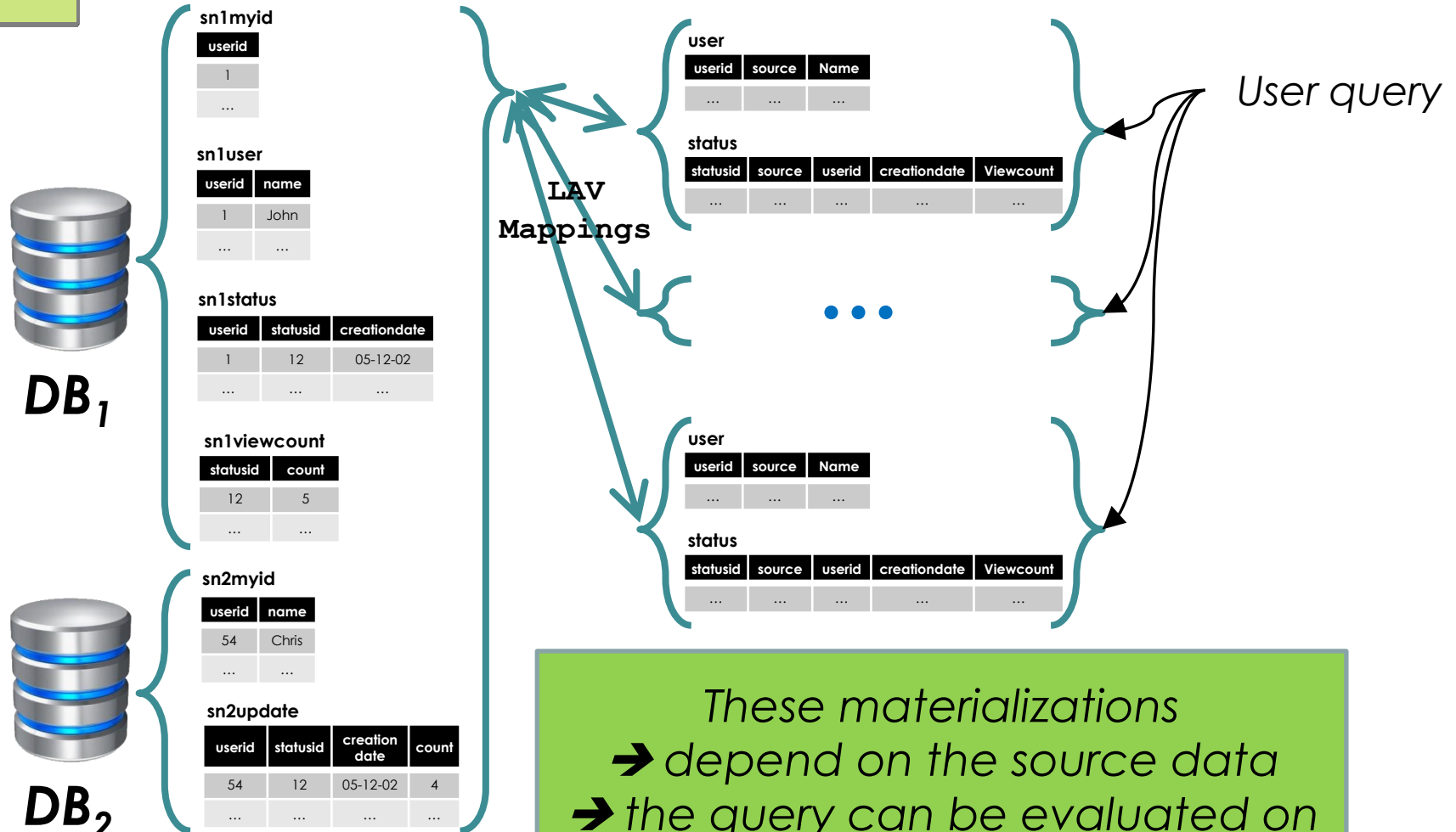
Private data sources  
with public views

# Certain answers semantics



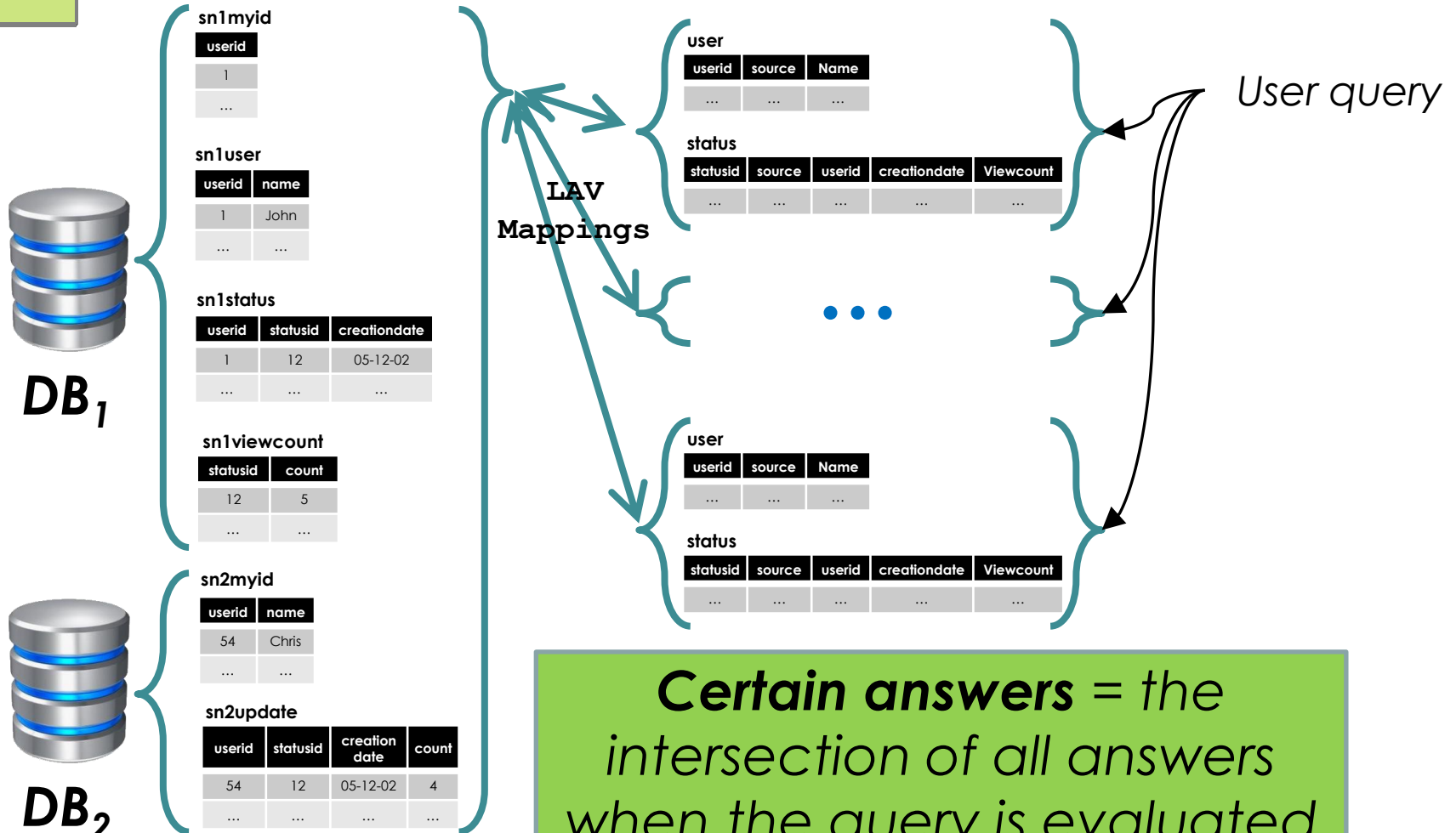
*All possible materializations of the global schema consistent with LAV mappings with OWA*

# Certain answers semantics



*These materializations  
 → depend on the source data  
 → the query can be evaluated on them.*

# Certain answers semantics



**Certain answers = the intersection of all answers when the query is evaluated on all materializations.**

# Computing certain answers

- Assume
  - A set of views with their extensions  $I_s$  (i.e. their data)
  - A set  $S$  of LAV mappings as conjunctive queries, with OWA
  - A datalog user query  $Q$  posed on the global schema
- Then [AD98]
  - Certain answers of  $Q$  wrt  $S$  and  $I_s$  are computed by any query plan  $Q'$  that is maximally contained in  $Q$  wrt  $S$ .
  - A query plan = a datalog query with only view predicates

# Max. contained query plans

- A query plan  $P$ 
  - a datalog query with only view predicates
  - Can be « expanded » or « unfolded » predicates are replaced by their definitions (from the LAV mappings)
  - $P^{\text{exp}}$  is the expanded query plan
- $P$  is a maximally contained query plan in a datalog query  $Q$  iff
  - $P^{\text{exp}}$  contained in  $Q$
  - For each query plan  $P'$  s.t.  $(P')^{\text{exp}}$  contained in  $Q$  there is  $(P')^{\text{exp}}$  contained in  $P^{\text{exp}}$

# Query containment

- Classical result
  - For CQ [CM77]
  - For datalog [S93]
- CQ containment is NP-complete.
  - $P \subseteq Q$  iff there exists a substitution  $\theta$  on variables(Q) s.t.
    - $\text{atoms}(\theta(Q)) \subseteq \text{atoms}(P)$ , and
    - $\text{var}(\text{head}(\theta(Q))) = \text{var}(\text{head}(P))$
- Datalog containment is undecidable.

# Complexity results [AD98]

- Data complexity of the problem of computing certain answers under OWA

views $s$	query $Q$				
	$CQ$	$CQ \neq$	$PQ$	<i>datalog</i>	$FO$
$CQ$	$PTIME$	$Co-NP$	$PTIME$	$PTIME$	und
$CQ \neq$	$PTIME$	$Co-NP$	$PTIME$	$PTIME$	und
$PQ$	$Co-NP$	$Co-NP$	$Co-NP$	$Co-NP$	und
<i>datalog</i>	$Co-NP$	und	$Co-NP$	und	und
$FO$	und	und	und	und	und

"und" means undecidable

- Depends on query containment complexity



# Running example

- LAV mappings

```
sn1user(userid,name) ← user(userid,'SN1',name).
sn1myid(userid) ← user(userid,'SN1',name),
                 status(statusid,'SN1',userid,creationdate,viewcount).
sn1status(userid,statusid,creationdate) ← status(statusid,'SN1',userid,creationdate,viewcount).
sn1viewcount(statusid,count) ← status(statusid,'SN1',userid,creationdate,count).
sn2myid(userid,name) ← user(userid,'SN2',name).
sn2update(userid,statusid,creationdate,count) ← status(statusid,'SN2',userid,creationdate,count).
```

- User query Q

```
q(userid, statusid, viewcount) ← user(userid, 'SN1', name),
                                status(statusid, 'SN1', userid, '2016-03-03', viewcount)
```

- Example of query plans

```
q(userid, statusid, viewcount) ← sn1myid(userid)
```

```
q(userid, statusid, viewcount) ← sn1user(userid, name),
                                sn1status(userid,statusid, '2016-03-03'),
                                sn1viewcount(statusid, viewcount).
```

```
q(userid, statusid, viewcount) ← sn2update(userid,statusid, '2016-03-03', viewcount).
```

# Running example

- LAV mappings

```
sn1user(userid,name) ← user(userid,'SN1',name).
sn1myid(userid) ← user(userid,'SN1',name),
                 status(statusid,'SN1',userid,creationdate,viewcount).
sn1status(userid,statusid,creationdate) ← status(statusid,'SN1',userid,creationdate,viewcount).
sn1viewcount(statusid,count) ← status(statusid,'SN1',userid,creationdate,count).
sn2myid(userid,name) ← user(userid,'SN2',name).
sn2update(userid,statusid,creationdate,count) ← status(statusid,'SN2',userid,creationdate,count).
```

- User query Q

```
q(userid, statusid, viewcount) ← user(userid, 'SN1', name),
                                status(statusid, 'SN1', userid, '2016-03-03', viewcount)
```

- Example of query plans

- P1 = `q(userid, statusid, viewcount) ← sn1myid(userid)`

- P1<sup>exp</sup> = `q(userid, statusid, viewcount) ← user(userid,'SN1',_1),
 status(_2,'SN1',userid,_3,_4).`

# Query plan P1

- User query Q

$q(\text{userid}, \text{statusid}, \text{viewcount}) \leftarrow \text{user}(\text{userid}, \text{'SN1'}, \text{name}),$   
 $\text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{'2016-03-03'}, \text{viewcount})$

- Query plan P1

- P1 =  $q(\text{userid}, \text{statusid}, \text{viewcount}) \leftarrow \text{sn1myid}(\text{userid})$

- P1<sup>exp</sup> =  $q(\text{userid}, \text{statusid}, \text{viewcount}) \leftarrow \text{user}(\text{userid}, \text{'SN1'}, \text{\_1}),$   
 $\text{status}(\text{\_2}, \text{'SN1'}, \text{userid}, \text{\_3}, \text{\_4}).$

- Problem 1

- Head variables statusid and viewcount are not in the body.
- When calling the view, we get only userid !

- Problem 2 : P1<sup>exp</sup> is not contained in Q  
there cannot be any substitution  $\theta$  s.t.

$$\theta(\text{'2016-03-03'}) = \text{\_3}$$

# Query plan P2

- User query Q

```
q(userid, statusid, ← user(userid, 'SN1', name),  
viewcount)      status(statusid, 'SN1', userid, '2016-03-03', viewcount)
```

- Query plan P2

- P2 = 

```
q(userid, statusid, viewcount) ← sn1user(userid, name),  
sn1status(userid, statusid, '2016-03-03'),  
sn1viewcount(statusid, viewcount).
```

- P2<sup>exp</sup> = 

```
q(userid, statusid, viewcount) ← user(userid, 'SN1', name),  
status(statusid, 'SN1', userid, '2016-03-03', _1),  
status(statusid, 'SN1', _3, _2, viewcount),
```

- P2<sup>exp</sup> is NOT contained in Q

$\theta$ :  
userid → userid            'SN1' → 'SN1'  
name → name                '2016-03-03' → '2016-03-03'  
statusid → statusid

viewcount → **\_1**

**We should have  
viewcount → viewcount !**

# Query plan P3

- User query Q

$q(\text{userid}, \text{statusid}, \text{viewcount}) \leftarrow \text{user}(\text{userid}, \text{'SN1'}, \text{name}), \text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{'2016-03-03'}, \text{viewcount})$

- Query plan P3

- P3 =

$q(\text{userid}, \text{statusid}, \text{viewcount}) \leftarrow \text{sn2update}(\text{userid}, \text{statusid}, \text{'2016-03-03'}, \text{viewcount}),$

- P3<sup>exp</sup> =

$q(\text{userid}, \text{statusid}, \text{viewcount}) \leftarrow \text{status}(\text{statusid}, \text{'SN2'}, \text{userid}, \text{'2016-03-03'}, \text{viewcount}),$

- Problem: P3<sup>exp</sup> is not contained in Q  
there cannot be any substitution  $\theta$  s.t.

$$\theta(\text{'SN1'}) = \text{'SN2'}$$

moreover

there is no atom  $\text{user}(\dots)$  in P3<sup>exp</sup>.

# Running example

- LAV mappings

```
sn1user(userid,name) ← user(userid,'SN1',name).
sn1myid(userid) ← user(userid,'SN1',name),
                 status(statusid,'SN1',userid,creationdate,viewcount).
sn1status(userid,statusid,creationdate) ← status(statusid,'SN1',userid,creationdate,viewcount).
sn1viewcount(statusid,count) ← status(statusid,'SN1',userid,creationdate,count).
sn2myid(userid,name) ← user(userid,'SN2',name).
sn2update(userid,statusid,creationdate,count) ← status(statusid,'SN2',userid,creationdate,count).
```

- User query Q

```
q(userid, statusid, ← user(userid, 'SN1', name),
   viewcount)      status(statusid, 'SN1', userid, '2016-03-03', viewcount)
```

**There is no (max.)  
contained rewriting  
for Q.**

# Updated running example

- LAV mappings

```
sn1user(userid,name) ← user(userid,'SN1',name).
sn1myid(userid) ← user(userid,'SN1',name),
                 status(statusid,'SN1',userid,creationdate,viewcount).
sn1status(userid,statusid,creationdate) ← status(statusid,'SN1',userid,creationdate,viewcount).
sn1viewcount(statusid,userid,creationdate, ← status(statusid,'SN1',userid,creationdate,count).
              count)
sn2myid(userid,name) ← user(userid,'SN2',name).
sn2update(userid,statusid,creationdate,count) ← status(statusid,'SN2',userid,creationdate,count).
```

- User query Q

```
q(userid, statusid, viewcount) ← user(userid, 'SN1', name),
                                status(statusid, 'SN1', userid, '2016-03-03', viewcount)
```

- Maximally contained rewriting (only one)

```
q(userid, statusid, viewcount) ← sn1user(userid,name),
                                sn1viewcount(statusid,userid,'2016-03-03',viewcount)
```

# Max contained query plans computation

- Naive algorithm (for CQ, not datalog) [LMSS95,RSU95]
- Bucket Algorithm [LRO96]
- Unification-join algorithm [Q96]
- Inverse rules [DG97,DGL00]
- Minicon [PH01]
- Inverse Minicon [K04]



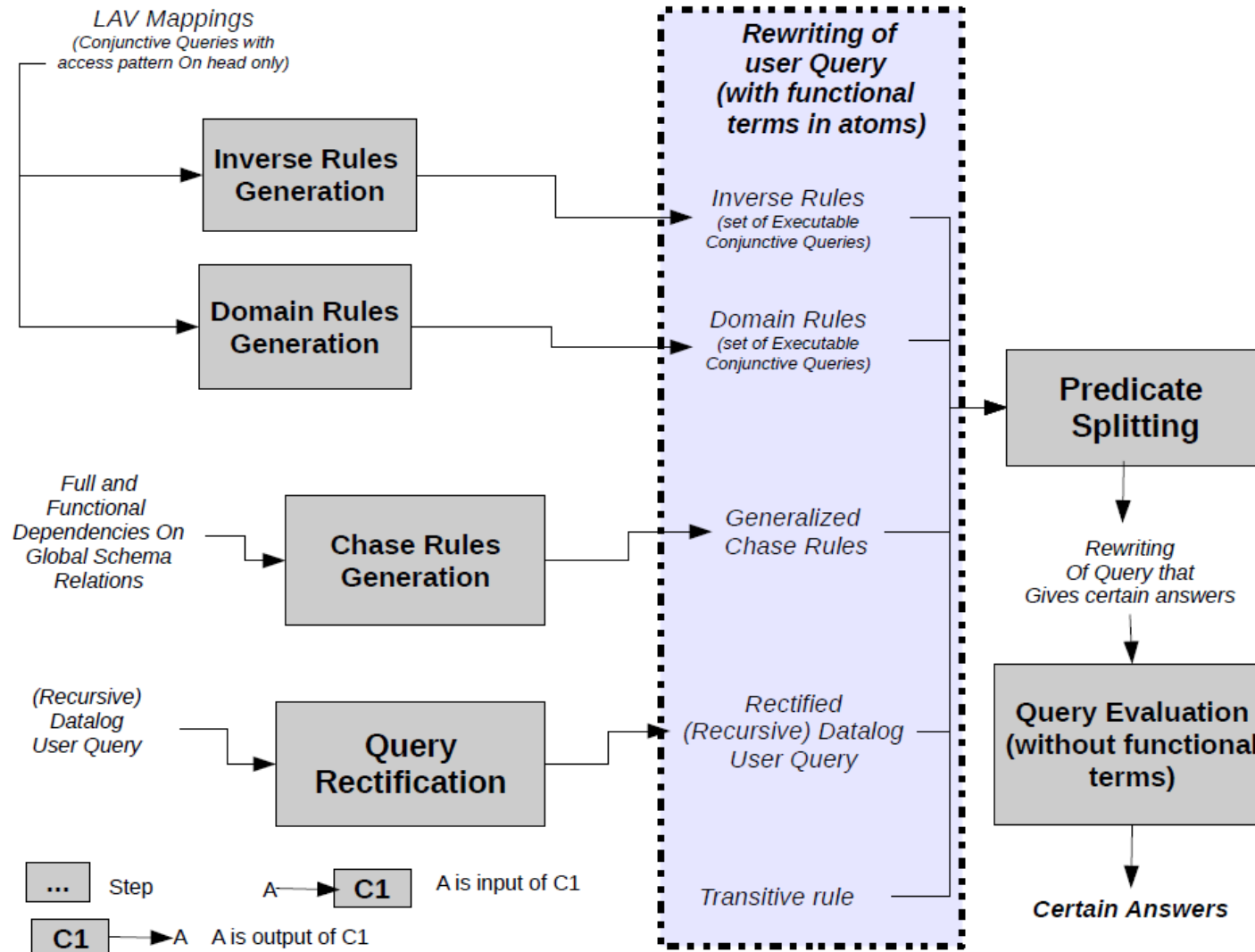
# Max contained query plans computation

- Naive algorithm (for CQ, not datalog) [LMSS95,RSU95]
  - Bucket Algorithm [LRO96]
  - Unification-join algorithm [Q96]
  - Inverse rules [DG97,DGL00]
  - Minicon [PH01]
  - Inverse Minicon [K04]
- The only algorithm that handles at the same time**
- recursive queries
  - access patterns (on head of mappings)
  - full and functional dependencies
  - PTIME complexity

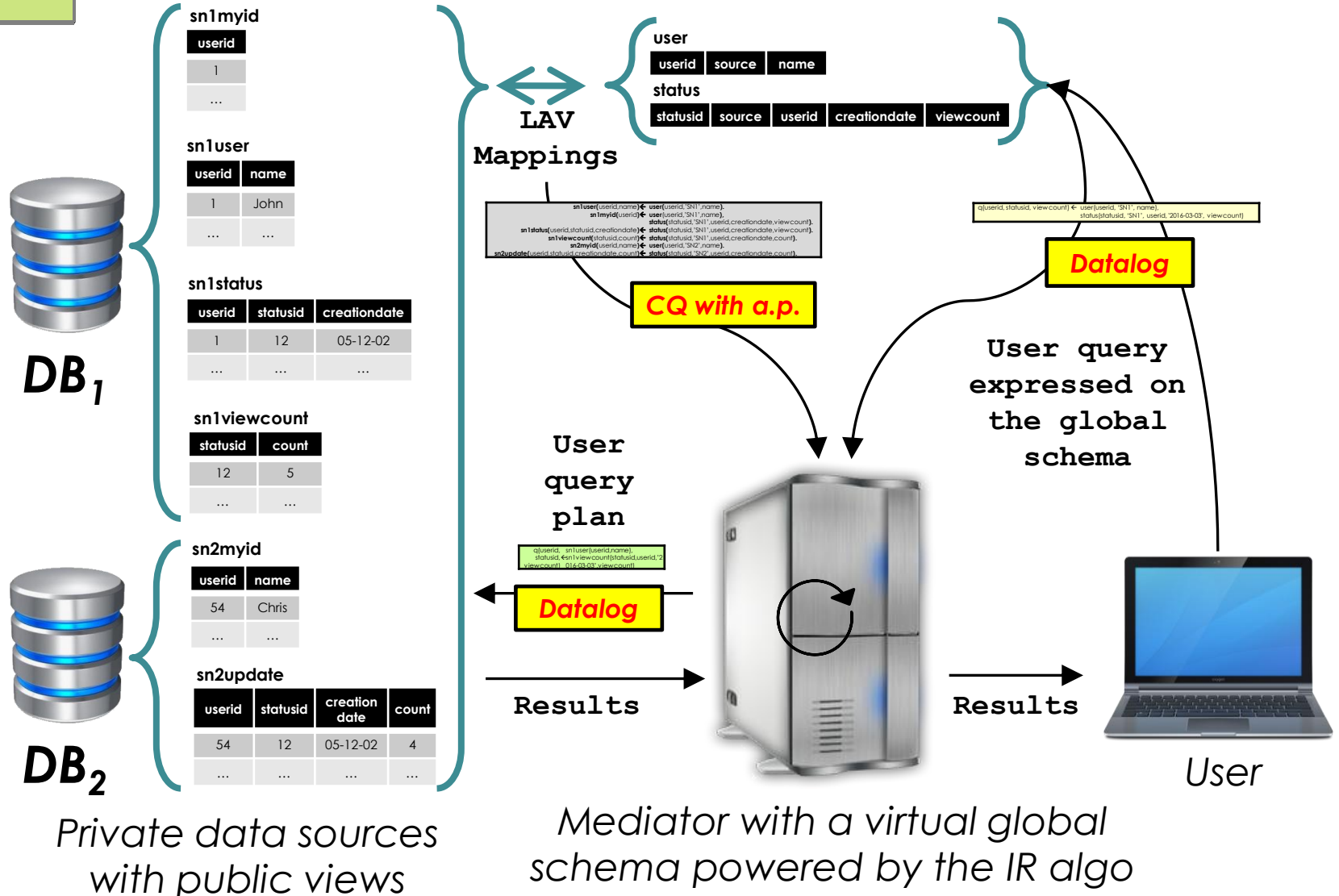
# Inverse rules [DG97,DGL00]

- LAV mappings: CQ with access patterns on heads
- User query: datalog
- Maximally contained rewritings: datalog queries
- Generate maximally contained rewritings by adding to (rectified) query rules
  - Inverse rules: handle LAV mapping
  - Domain rules: handle access patterns
  - Generalized chase rules: handle dependencies
  - Transitive rules: handle the equality between variables
- Get rid of functional terms by predicate splitting
- Compute certain answers by query evaluation

# Inverse rules [DG97,DGL00]



# From DI to WS data extraction



# From DI to WS data extraction

sn1myid

userid
1
...

sn1user

userid	name
1	John
...	...

sn1status

userid	statusid	creationdate
1	12	05-12-02
...	...	...

sn1viewcount

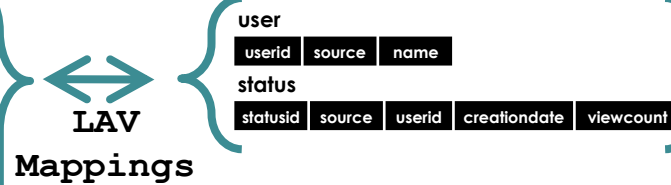
statusid	count
12	5
...	...

sn2myid

userid	name
54	Chris
...	...

sn2update

userid	statusid	creation date	count
54	12	05-12-02	4
...	...	...	...



```

sn1user(userid, name) ← user(userid, SN1, name),
sn1myid(userid) ← user(userid, SN1, name),
sn1status(userid, statusid, creationdate) ← status(statusid, SN1, userid, creationdate, viewcount),
sn1viewcount(statusid, count) ← status(statusid, SN1, userid, creationdate, viewcount),
sn2myid(userid, name) ← user(userid, SN2, name),
sn2update(userid, statusid, creationdate, count) ← status(statusid, SN2, userid, creationdate, viewcount)
    
```

CQ with a.p.

```

q(userid, statusid, viewcount) ← user(userid, SN1, name),
status(statusid, SN1, userid, 2016-03-03, viewcount)
    
```

Datalog

User query expressed on the global schema

User query plan

```

q(userid, sn1user(userid, name),
statusid, sn1viewcount(statusid, userid, 2, viewcount))
    
```

Datalog

Results

Results

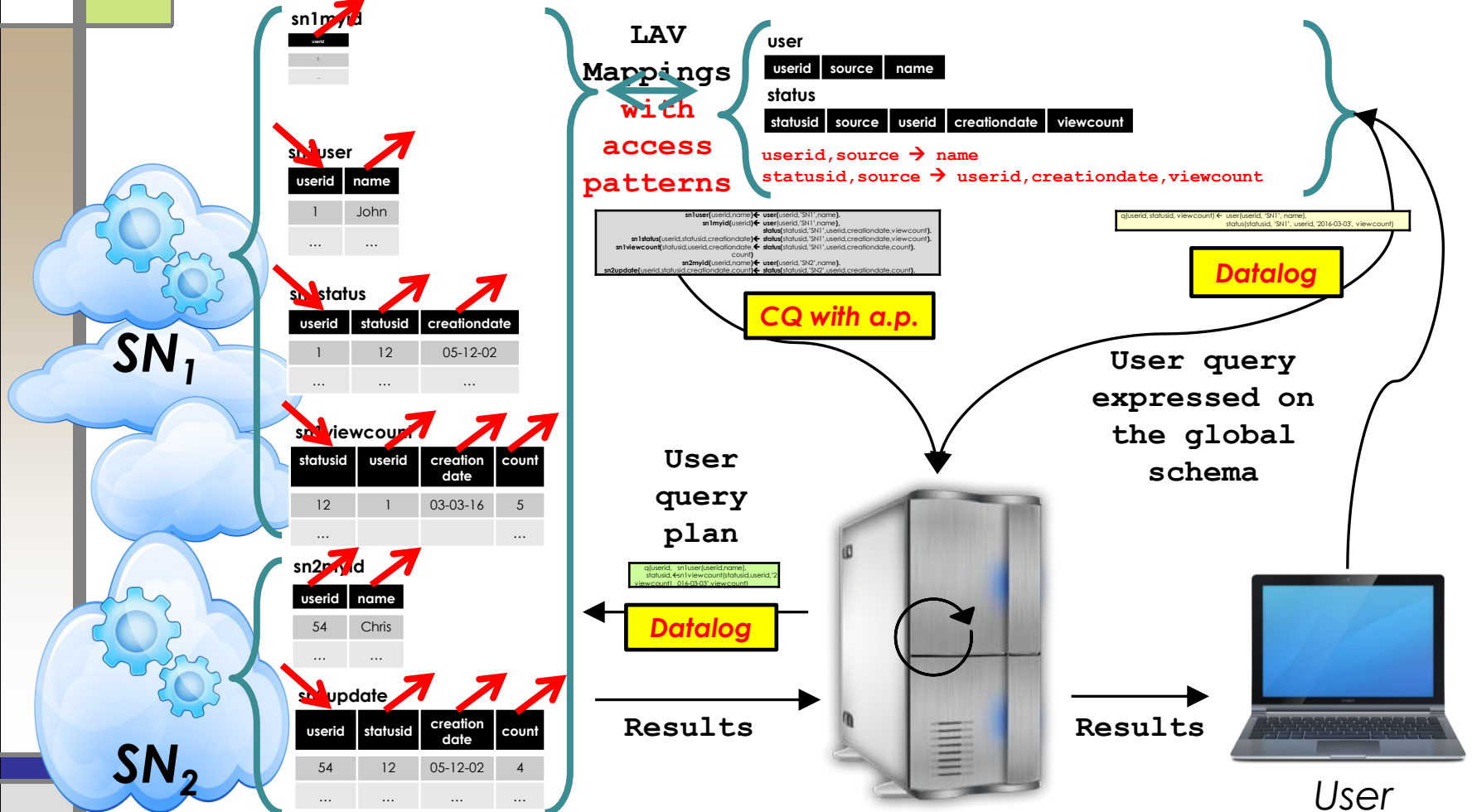


User

Private data sources with public views

Mediator with a virtual global schema powered by the IR algo

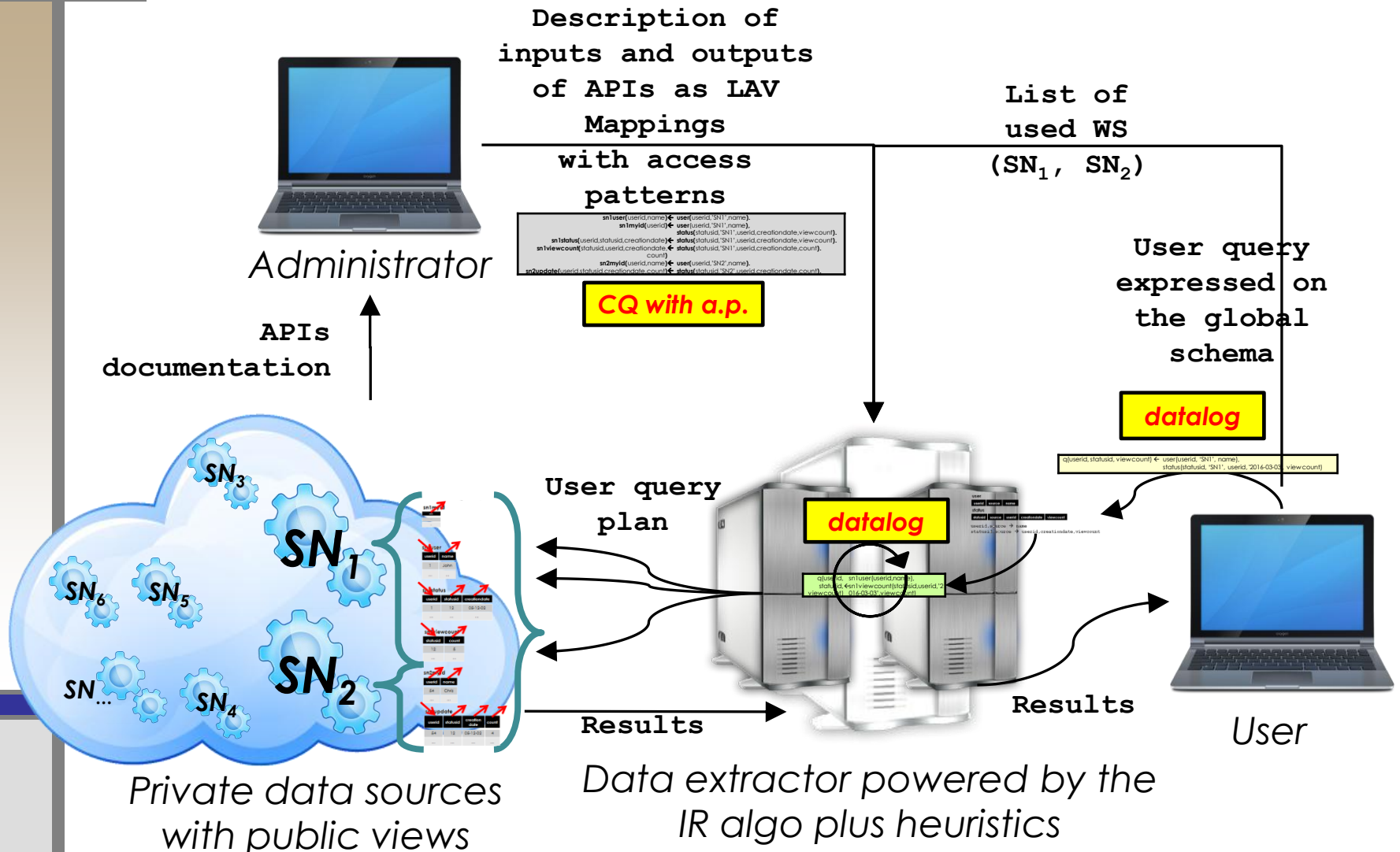
# From DI to WS data extraction



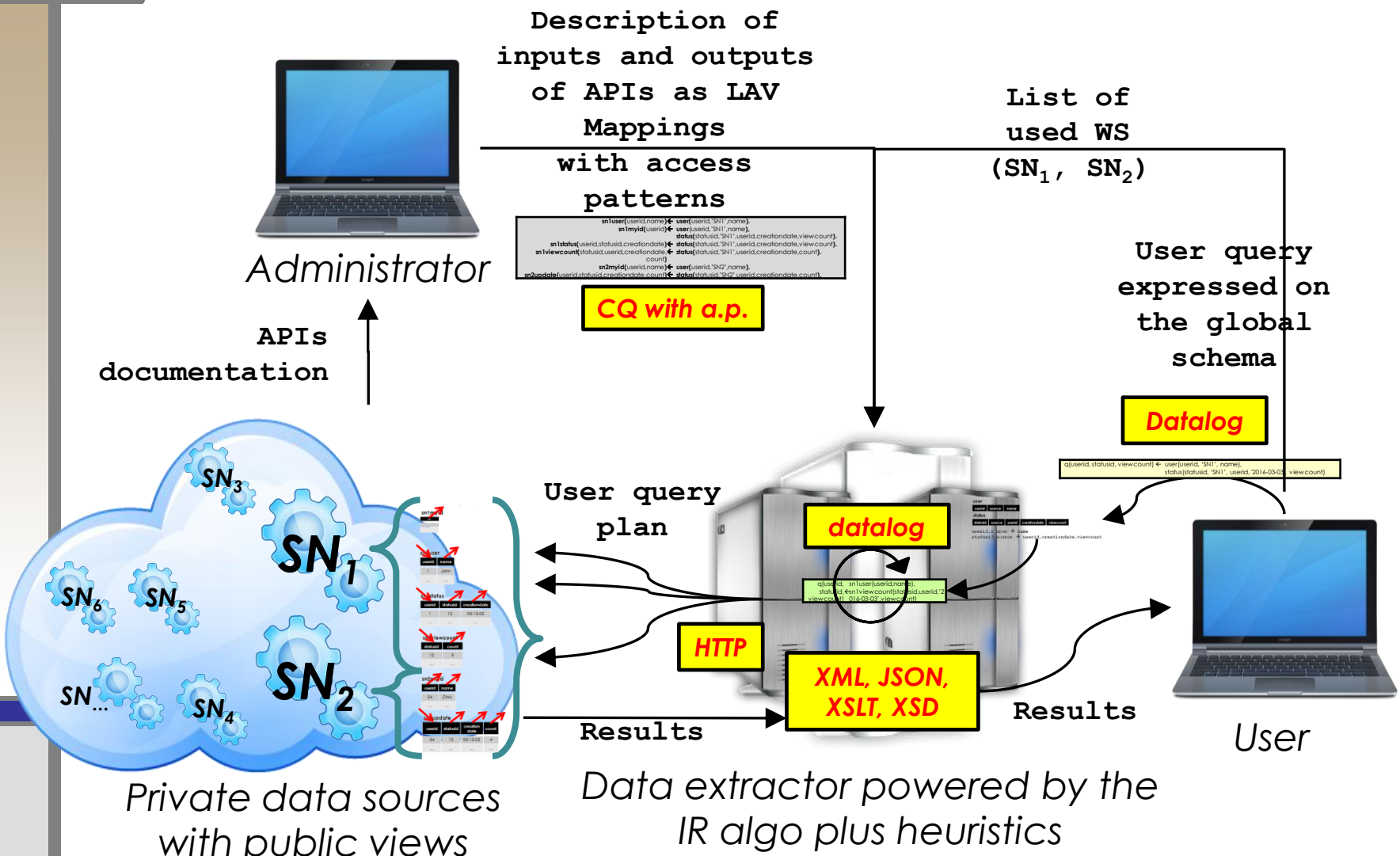
Private data sources with public views

Mediator with a virtual global schema powered by the IR algo

# From DI to WS data extraction



# From DI to WS data extraction





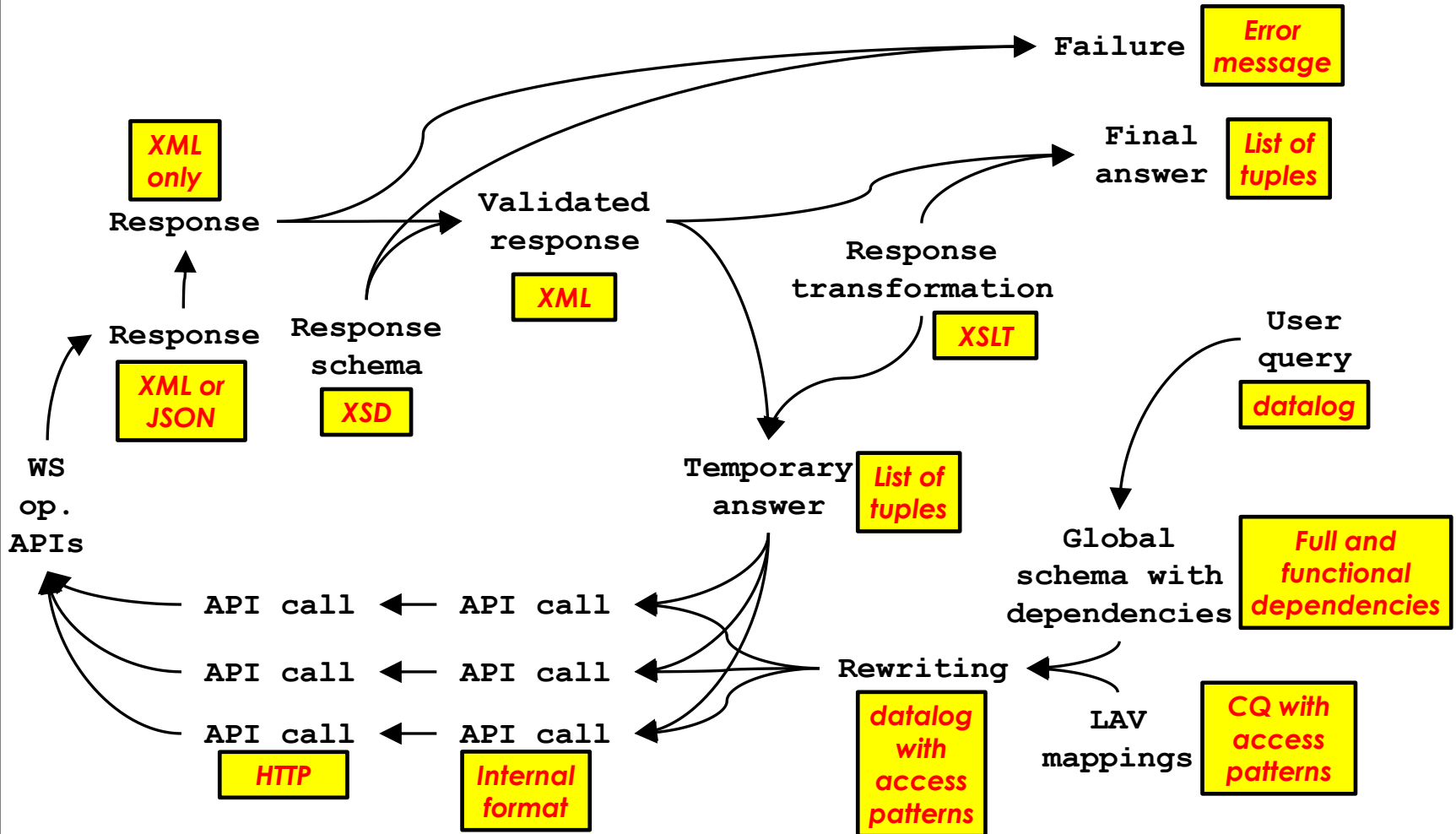
# Declarative languages

- **CQ with access patterns:** express the LAV mappings that define WS operations in terms of the global schema
- **Datalog:**
  - express the user queries posed to the global schema
  - express the max. cont. rewritings, i.e. the way the user query is distributed over WS operations
- **XML / JSON:**
  - Send requests to WS operations
  - Get back the response of WS operations
- **XSD:** validate operation responses and detecting API changes
- **XSLT:** extract and transform desired information from operation response
- **HTTP:** call WS operations APIs

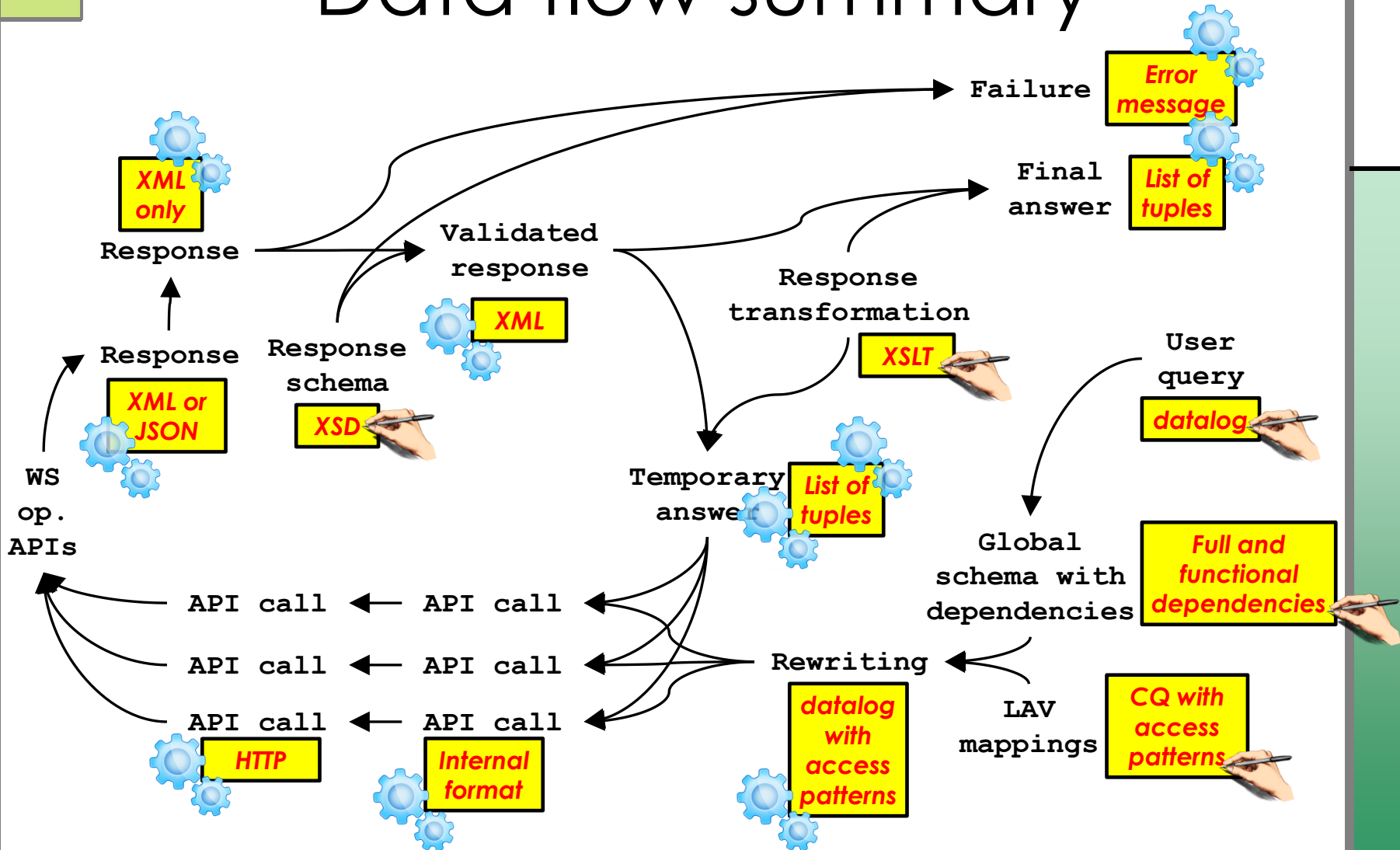
# Handling incomplete information

- Certain answers semantics:
  - A bit restricted (cf examples) for real use
  - Incomplete query answers may still be interesting.
- Relaxation:
  - Handle functional terms in the IR algo
  - Allows to return incomplete but additional answers
- Proposal of a heuristics
- Proposal of a semantics for datalog with access patterns
- Work in progress

# Data flow summary



# Data flow summary



# [Main scenarios]

Studied issue

Locks and Keys

Declarative data extraction

Data integration with WS

**Main scenarios**

Results and discussion

# Scenarios outline

- 1. Introducing the WS and WADL (sn1, sn2)
- 2. Defining the Global Schema Relation (user, status)
- 3. Defining the WS (Category, Provider, API, Operations)
- 4. Defining the Local Schema Relations (LAV Mapping)
- 5. Defining Records
- 6. Defining Performance Indicators
- 7. Defining the user (User's choice of WS, Records, Performance Indicators)
- 8. Obtaining the user record values, performance indicators

# [Results and discussion]

Studied issue

Locks and Keys

Declarative data extraction

Data integration with WS

Main scenarios

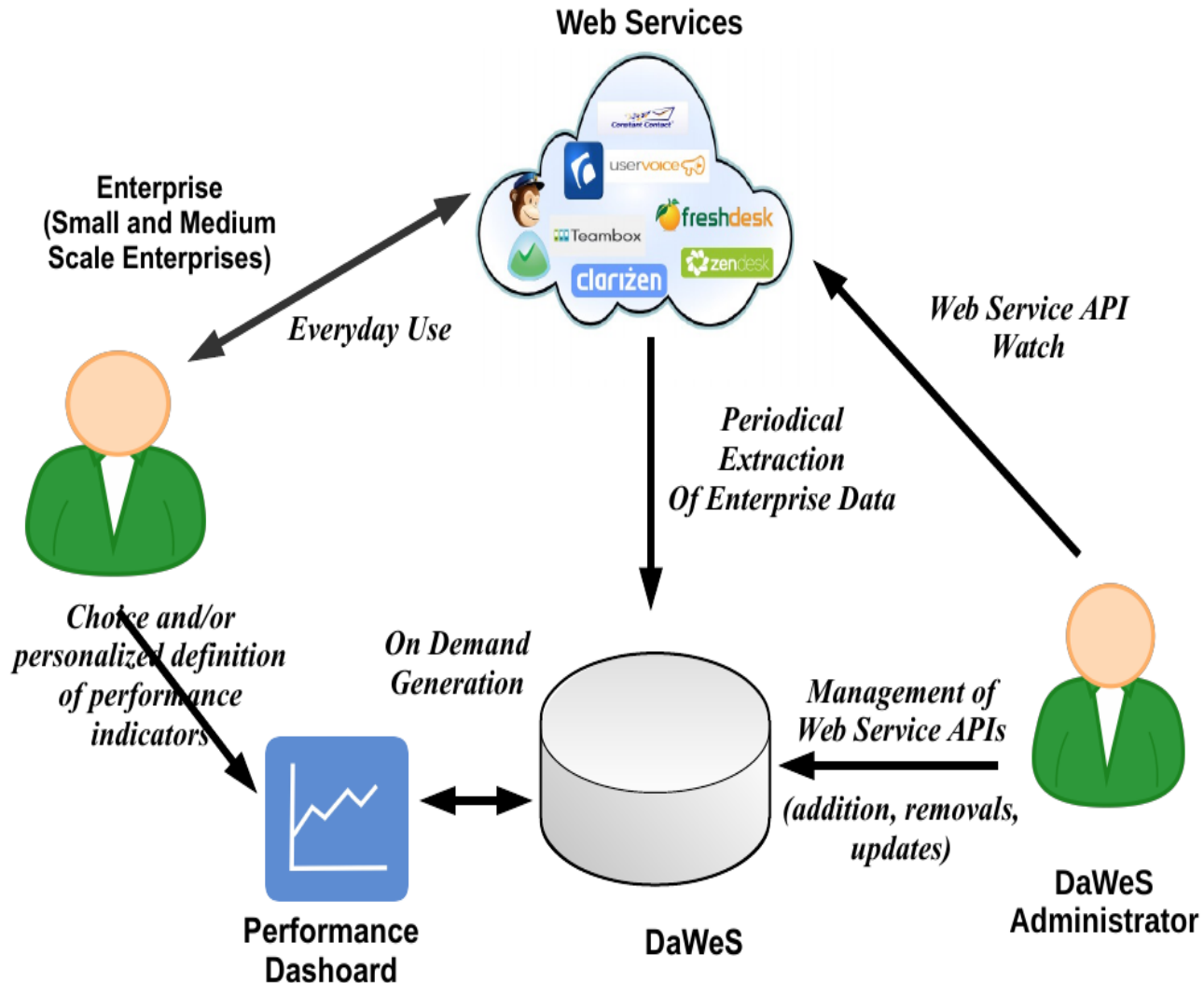
**Results and discussion**

# Application to DW

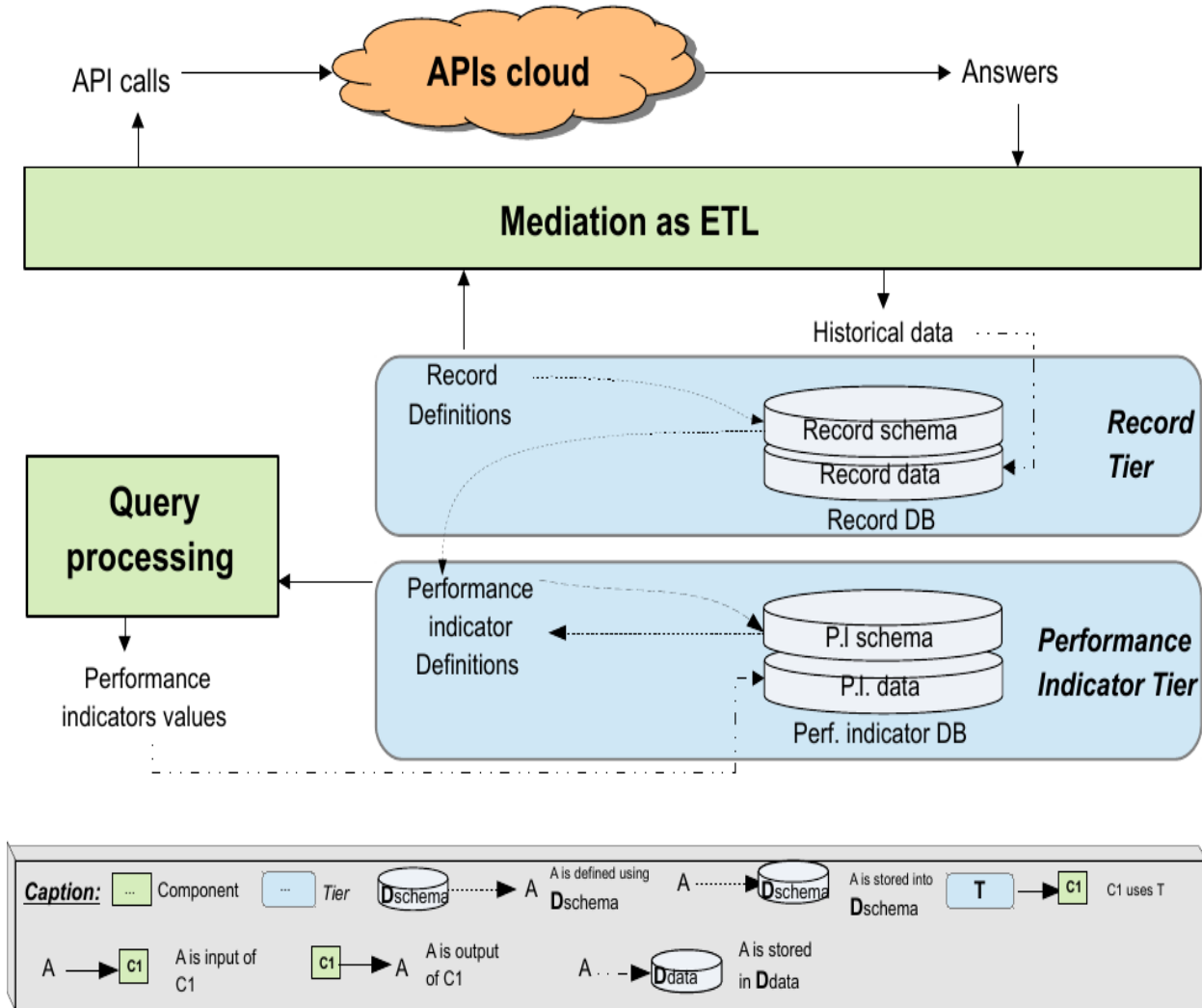
- Use of this declarative approach to WS data extraction to feed a datawarehouse (DW)
- « Mediation as ETL »
- Prototype : DaWeS  
« Data Warehouse fed with Web Services »



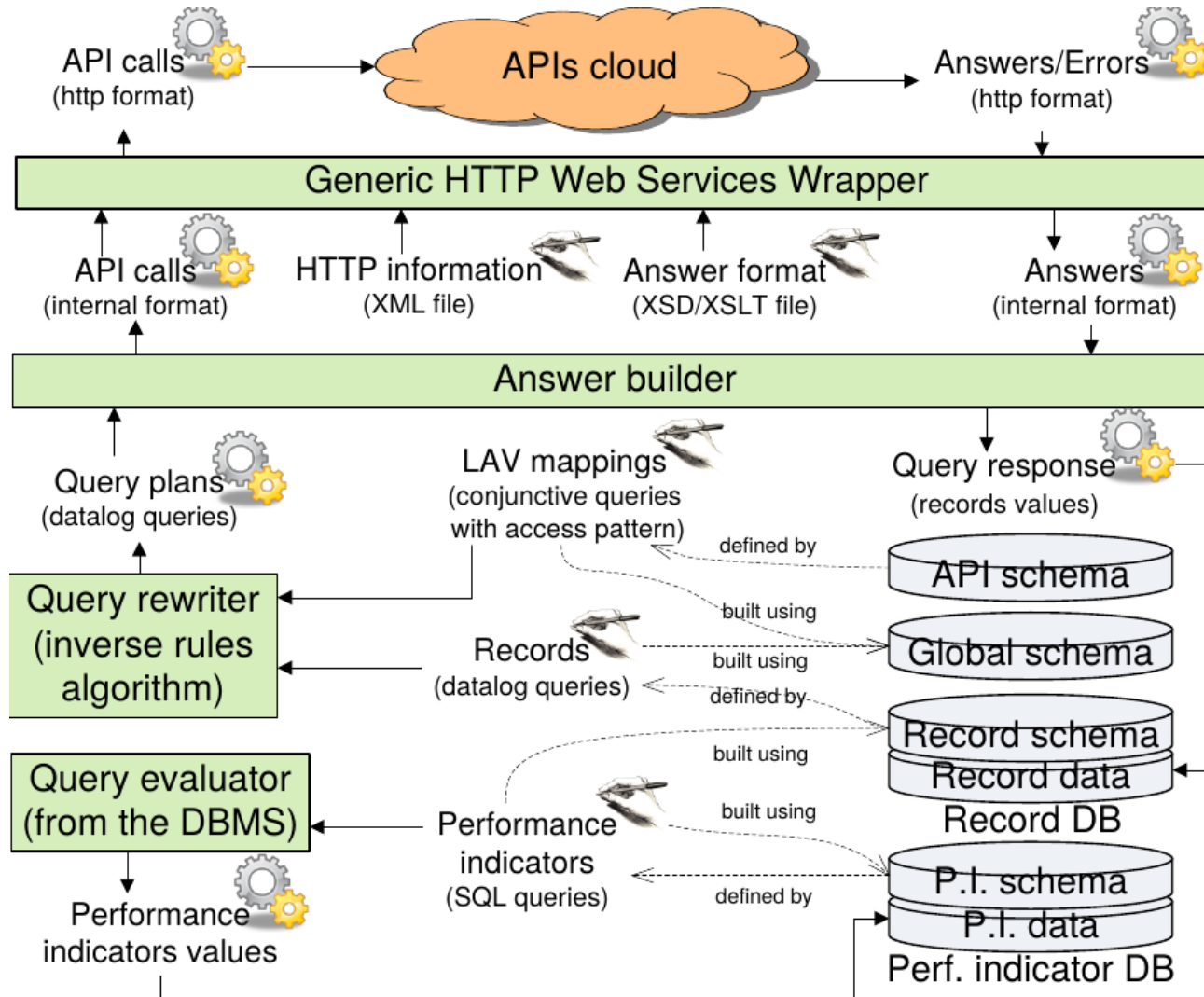
# DaWeS overview



# DaWeS architecture



# DaWeS declarativeness



**Caption:** ... Component B → C → A B is an input of C A is an output of C A (hand icon) A is manually obtained A (gear icon) A is automatically computed

# DaWeS experiments

- Number of domain of Web Services: 3
- Number of Web services considered : 12
- Number of API operations considered: 35
- Number of Global Schema Relations: 12
- Number of Test Organizations: 100
- Message Formats: 2  
XML, JSON
- Authentication Mechanisms: 2  
HTTP basic authentication, OAuth
- Number of Record Definitions: 17
- Number of Performance Indicator Queries: 20
- Operation details:
  - None, one or more input parameters
  - Pagination

# Discussion

- Mediation: not recent domain  
About 20 years old
- Industrial development of mediators ?  
No one, up to our knowledge
- Main reasons (we suppose)
  - Computationally complex algorithms
  - Not easy to implement
  - Need to add extra features to be really useful
    - Access patterns
    - Incomplete information
    - Dependencies
    - Aggregation functions
    - More expressive view language
    - GAV, GLAV mappings

# Perspectives

- Extra features management
  - Adding features ?
  - Already done:
    - One heuristics to handle incomplete info
    - Definition of the semantics of datalog with access patterns
  - To do:
    - Going on with the state-of-the-art [DLN07, YKC06, CM08, CCM09, ...]
    - Complexity study of this heuristics (bound on the number of accesses)
    - Enhancing our heuristics to algorithm
    - Complexity study of this algorithm
- Concerning DaWeS
  - Optimizing the number of operation calls
  - Access patterns and errors
  - Automated error handling

# Perspectives

- Industrial scale experimentations
  - Already done
    - 3 domains: project management, email marketing, support/helpdesk
    - 12 WS
    - 35 operations
  - To do: other domains, far more operations
- Broaden the approach
  - To ontological query answering of WS APIs
  - To other applications (than feeding a DW)

# DaWeS future

- Looking for industrial partners
  - Other domains
  - Other applications
- Looking for academic partners
  - Towards a test platform for rewriting algorithms
    - Data integration extensions
    - Ontological query answering algorithms
    - ...



Thank you.

**[Additional slides]**

# Bibliography

- [AD98] Abiteboul S. and O. M. Duschka. Complexity of Answering Queries Using Materialized Views. In *Proc. 17th Annual ACM Symp. Principles of Databases (PODS '98)*, Seattle, Washington 1998, pp. 254–263.
- [S93] Oded Shmueli. Equivalence of DATALOG Queries is Undecidable. [J. Log. Program. 15\(3\)](#) 231-241 (1993)
- [CM77] A. K. Chandra, P. M. Merlin. Optimal implementation of conjunctive queries. In *Proc. ACM SIGACT Symp. on the Theory of Computing (STOC '77)*, pp. 77–90, 1977.
- [FGMP05] Fagin, R., Kolaitis, P. G., Miller, R. J., and Popa, L. Data exchange: semantics and query answering. *Theor. Comput. Sci.*, 336(1):89–124.
- [ABFL13] Arenas, M., Barceló, P., Fagin, R., and Libkin, L. (2013). Solutions and query rewriting in data exchange. *Inf. Comput.*, 228:28–61.
- [DLN07] Deutsch, A., Ludäscher, B., and Nash, A. (2007). Rewriting queries using views with access patterns under integrity constraints. *Theor. Comput. Sci.*, 371(3):200–226.
- [YKC06] Yang, G., Kifer, M., and Chaudhri, V. K. (2006). Efficiently ordering subgoals with access constraints. In *PODS*, pages 183–192.
- [CM08] Cali, A. and Martinenghi, D. (2008). Querying data under access limitations. In *ICDE*, pages 50–59.
- [CCM09] Cali, A., Calvanese, D., and Martinenghi, D. (2009a). Dynamic query optimization under access limitations and dependencies. *J. UCS*, 15(1):33–62.

# Bibliography

- [LMSS95] Alon Y. Levy, Alberto O. Mendelzon, Yehoshua Sagiv, Divesh Srivastava: Answering Queries Using Views. PODS 1995: 95-104
- [RSU95] Anand Rajaraman, Yehoshua Sagiv, Jerrey D. Ullman: Answering Queries Using Templates with Binding Patterns. PODS 1995: 105-112
- [LRO96] Alon Y. Levy, Anand Rajaraman, Joann J. Ordille: Querying Heterogeneous Information Sources Using Source Descriptions VLDB 1996: 251-262
- [LRO96] Alon Y. Levy, Anand Rajaraman, Joann J. Ordille: Query-Answering Algorithms for Information Agents. AAI/IAAI, Vol. 1 1996: 40-47
- [Q96] Xiaolei Qian: Query Folding. ICDE 1996:48-55
- [DR97] Oliver M. Duschka, Michael R. Genesereth: Answering Recursive Queries Using Views. PODS 1997:109-116
- [DRL00] Oliver M. Duschka, Michael R. Genesereth, Alon Y. Levy: Recursive query plans for data integration. The Journal of Logic Programming 43 (2000):49-73
- [PH01] Rachel Pottinger, Alon Y. Halevy: MiniCon: A scalable algorithm for answering queries using views. VLDB J. (VLDB) 10(2-3):182-198 (2001)
- [K04] Christoph Koch: Query rewriting with symmetric constraints. AI Commun. (AICOM) 17(2):41-55 (2004)

# Characteristics of some real WS

# A study of some real WS in the field of project management

Web Service	Basecamp	Liquid Planner	Teamwork	Zoho Projects	Wrike	Teambox
1. API Description	HTML page	HTML page	HTML page	HTML page	HTML page	HTML page
2. Conformance to REST	REST like	REST like	REST like	Not REST	Not REST	REST like
3. Version	v1	3.0.0	N.A.	N.A.	v2	1
4. Authentication	Basic HTTP, OAuth 2	Basic HTTP	Basic HTTP	Basic HTTP	OAuth 1.0	OAuth 2.0
5. Resources Involved	Project, Todo List, Todo	Project, Task	Project, Task List, Task	Project, Task List, Task	Task	Project, Task
6. Message Formats	JSON	JSON	XML, JSON	XML, JSON	XML, JSON	JSON
7. Service Level Agreement	Max 500 requests /10s from same IP address for same account	Max 30 requests /15s for same account	Max 120 requests /1min	Error code:6403 on exceeding the limit	N.A.	N.A.
8. HTTP Resource Access	GET	GET	GET	POST	POST	GET
9. Data Types (dt)	Enumerated dt (Project and Todo Status), Date	Enumerated dt (Project and Task Status), Date	Enumerated dt (Project and Task Status), Date	Enumerated dt (Project and Task Status), Date	Enumerated dt (Task Status), Date	Enumerated dt (Project and Task Status), Date
10. Dynamic nature of the resources	Yes (Project and Todo Status)	Yes (Project and Task Status)	Yes (Project and Task Status)	Yes (Project and Task Status)	Yes (Task Status)	Yes (Project and Task Status)
11. Operation Invocation Sequence Required	Yes	No	Yes	Yes	Yes	Yes
12. Pagination	No	No	No	Yes	Yes	No

# A study of some real WS in the field of email marketing

Web Service	Mailchimp	Campaign Monitor	iContact	Constant Contact	AWeber
1. API Description	HTML page	HTML page	HTML page	HTML page	HTML page
2. Conformance to REST	Not REST	REST like	REST like	REST	REST
3. Version	1.3	v3	2.2	N.A.	1.0
4. Authentication	Basic HTTP	Basic HTTP, OAuth 2	Basic HTTP (with Sandbox)	OAuth 2.0	OAuth 1.0
5. Resources Involved	Campaign, Campaign Statistics	Campaign, Campaign Statistics	Campaign, Campaign Statistics	Campaign, Campaign Statistics	Campaign, Campaign Statistics
6. Message Formats	XML, JSON, PHP, Lolcode	XML, JSON	XML, JSON	XML	JSON
7. Service Level Agreement	N.A.	N.A.	75,000 requests /24h, with a max of 10,000 requests /1h	N.A	60 requests per minute
8. HTTP Resource Access	GET	GET	GET	GET	GET
9. Data Types (dt)	Enumerated Data types (Campaign Status), Date	Enumerated Data types (Campaign Status), Date	Enumerated Data types (Campaign Status), Date	Enumerated Data types (Campaign Status), Date	Enumerated Data types (Campaign Status), Date
10. Dynamic nature of the resources	Yes (Campaign Status)	Yes (Campaign Status)	Yes (Campaign Status)	Yes (Campaign Status)	Yes (Campaign Status)
11. Operation Sequence Required	Yes	Yes	No	Yes	Yes
12. Pagination	Yes	No	No	Yes	Yes

# A study of some real WS in the field of support/ helpdesk

Web Service	Zendesk	Desk	Zoho Support	Uservoice	Freshdesk	Kayako
1. API Description	HTML page	HTML page	HTML page	HTML page	HTML page	HTML page
2. Conformance to REST	REST like	REST	Not REST	REST like	REST like	REST
3. Version	v1	v2	N.A.	v1	N.A.	N.A.
4. Authentication	Basic HTTP	Basic HTTP, OAuth 1.0a	Basic HTTP	OAuth 1.0	Basic HTTP	Basic HTTP
5. Resources Involved	Forum, Topic, Ticket	Case	Task	Forum, Topic, Ticket	Forum, Topic, Ticket	Ticket, Topic
6. Message Formats	XML, JSON	JSON	XML, JSON	XML, JSON	JSON	XML
7. Service Level Agreement	Limit exists (but unknown)	60 requests per minute	250 calls /day /org (Free)	N.A.	N.A.	N.A.
8. HTTP Resource Access	GET	GET	GET	GET	GET	GET
9. Data Types (dt)	Enumerated dt (Ticket Status), Date	Enumerated dt (Case Status), Date	Enumerated dt (Task Status), Date	Enumerated dt (Ticket Status), Date	Enumerated dt (Ticket Status), Date	Enumerated dt (Ticket Status), Date
10. Dynamic resources	Yes (Ticket Status)	Yes (Case Status)	Yes (Task Status)	Yes (Ticket Status)	Yes (Ticket Status)	Yes (Ticket Status)
11. Operation Sequence Required	Yes	Yes	Yes	Yes	Yes	Yes
12. Pagination	Yes	Yes	Yes	Yes	No	Yes



# JSON schema

```
1 {
2   "$schema": "http://json-schema.org/draft-04/schema#",
3   "id": "http://jsonschema.net",
4   "type": "object",
5   "properties": {
6     "updates": {
7       "id": "http://jsonschema.net/updates",
8       "type": "array",
9       "items": {
10        "id": "http://jsonschema.net/updates/1",
11        "type": "object",
12        "properties": {
13          "identifier": {
14            "id": "http://jsonschema.net/updates/1/identifier",
15            "type": "integer"
16          },
17          "update": {
18            "id": "http://jsonschema.net/updates/1/update",
19            "type": "string"
20          },
21          "date": {
22            "id": "http://jsonschema.net/updates/1/date",
```

# Computing certain answers complexity

# Complexity [AD98]

Let  $\mathbf{IS} = ((G, I_G), (S, I_S), M = (s : I_G \rightarrow s(I_G)))$  be a LAV integration system under OWA. Let  $Q$  be a query on  $G$ . Let's call "views" the predicates defined in  $S$  by  $s$ .

- ▶ Languages for queries and views may be:
  - ▶ CQ and CQ with inequality ( $CQ^{\neq}$ )
  - ▶ Positive query ( $PQ$ ): nr-datalog program with a particular query predicate
  - ▶ Positive query + inequality ( $PQ^{\neq}$ )
  - ▶ Datalog queries (*datalog*)
  - ▶ First order queries ( $FO$ )
- ▶ Complexity may be:
  - ▶ Data complexity: function of the size of  $I_S$
  - ▶ Query complexity: function of the sizes of the view definition  $s$  and the query  $Q$
  - ▶ Combined complexity: function of all (both the sizes of the view definition  $s$  and the query  $Q$ , and the size of  $I_S$ )

# Complexity [AD98]

- ▶ For  $s \in PQ^{\neq}$  and  $Q \in datalog^{\neq}$ , the problem of determining, given  $I_s$ , whether a tuple is a certain answer under OWA or CWA is in co-NP (combined complexity).
- ▶ Data complexity of the problem of computing certain answers under OWA.

views $s$	query $Q$				
	$CQ$	$CQ^{\neq}$	$PQ$	$datalog$	$FO$
$CQ$	$PTIME$	$Co-NP$	$PTIME$	$PTIME$	und
$CQ^{\neq}$	$PTIME$	$Co-NP$	$PTIME$	$PTIME$	und
$PQ$	$Co-NP$	$Co-NP$	$Co-NP$	$Co-NP$	und
$datalog$	$Co-NP$	und	$Co-NP$	und	und
$FO$	und	und	und	und	und

"und" means undecidable

# Complexity [AD98]

Let

- ▶  $\mathbf{IS} = ((G, I_G), (S, I_S), M = (s : I_G \rightarrow s(I_G)))$  be a LAV integration system under OWA
- ▶  $\mathcal{L}_1, \mathcal{L}_2 \in \{CQ, CQ^\neq, PQ, Datalog, FO\}$
- ▶  $s \in \mathcal{L}_1$  and  $Q \in \mathcal{L}_2$

Then

1. 

(View pb)	<b>iff</b>	(Containment pb)
Computing		the containment problem of a query
$cert(Q, s, I_S)$		in $\mathcal{L}_1$ in a query in $\mathcal{L}_2$ is decidable
is decidable		
2. If the two problems are decidable: the combined complexity of the view pb = the query complexity of the containment pb  
( $\Rightarrow$  The data complexity of the problem of computing certain answers under OWA is **at most** the query complexity of the query containment problem)

# The bucket algo example 1

**sn1user**(userid,name) ← **user**(userid,'SN1',name).  
**sn1myid**(userid) ← **user**(userid,'SN1',name),  
**status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1status**(userid,statusid,creationdate) ← **status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1viewcount**(statusid,count) ← **status**(statusid,'SN1',userid,creationdate,count).  
**sn2myid**(userid,name) ← **user**(userid,'SN2',name).  
**sn2update**(userid,statusid,creationdate,count) ← **status**(statusid,'SN2',userid,creationdate,count).

**q**(userid, statusid, viewcount) ← **user**(userid, 'SN1', name),  
**status**(statusid, 'SN1', userid, '2016-03-03', viewcount)

<b>user</b> (userid, 'SN1', name)	<b>status</b> (statusid, 'SN1', userid, '2016-03-03', viewcount)
<b>sn1user</b> (userid,name)	<b>sn1myid</b> (userid)
<b>sn1myid</b> (userid)	<b>sn1status</b> (userid,statusid,'2016-03-03')
	<b>sn1viewcount</b> (statusid,viewcount)

### Candidate query plans

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1myid**(userid)

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1status**(userid,statusid, '2016-03-03')

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1viewcount**(statusid,viewcount)

**q**(userid, statusid, viewcount) ← **sn1myid**(userid), ~~**sn1myid**(userid)~~

**q**(userid, statusid, viewcount) ← **sn1myid**(userid), **sn1status**(userid,statusid, '2016-03-03')

**q**(userid, statusid, viewcount) ← **sn1myid**(userid), **sn1viewcount**(statusid,viewcount)

**sn1user**(userid,name) ← **user**(userid,'SN1',name).  
**sn1myid**(userid) ← **user**(userid,'SN1',name),  
**status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1status**(userid,statusid,creationdate) ← **status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1viewcount**(statusid,count) ← **status**(statusid,'SN1',userid,creationdate,count).  
**sn2myid**(userid,name) ← **user**(userid,'SN2',name).  
**sn2update**(userid,statusid,creationdate,count) ← **status**(statusid,'SN2',userid,creationdate,count).

**q**(userid, statusid, viewcount) ← **user**(userid, 'SN1', name),  
**status**(statusid, 'SN1', userid, '2016-03-03', viewcount)

### Candidate query plans

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1myid**(userid)

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1status**(userid,statusid, '2016-03-03')

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1viewcount**(statusid,viewcount)

**q**(userid, statusid, viewcount) ← **sn1myid**(userid), ~~**sn1myid**(userid)~~

**q**(userid, statusid, viewcount) ← **sn1myid**(userid), **sn1status**(userid,statusid, '2016-03-03')

**q**(userid, statusid, viewcount) ← **sn1myid**(userid), **sn1viewcount**(statusid,viewcount)

### Candidate expanded query plans (bodys only)

**user**(userid,'SN1',name), **user**(userid,'SN1',name2), **status**(statusid2,'SN1',userid,creationdate2,viewcount2)

**user**(userid,'SN1',name), **status**(statusid,'SN1',userid, '2016-03-03',viewcount2)

**user**(userid,'SN1',name), **status**(statusid,'SN1',userid2,creationdate2,viewcount)

**user**(userid,'SN1',name2), **status**(statusid2,'SN1',userid,creationdate2,viewcount2)

**user**(userid,'SN1',name2), **status**(statusid2,'SN1',userid,creationdate2,viewcount2), **status**(statusid,'SN1',userid, '2016-03-03',viewcount2)

**user**(userid,'SN1',name2),  
**status**(statusid2,'SN1',userid,creationdate2,viewcount2), **status**(statusid,'SN1',userid2,creationdate3,viewcount)



**sn1user**(userid,name) ← **user**(userid,'SN1',name).  
**sn1myid**(userid) ← **user**(userid,'SN1',name),  
**status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1status**(userid,statusid,creationdate) ← **status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1viewcount**(statusid,count) ← **status**(statusid,'SN1',userid,creationdate,count).  
**sn2myid**(userid,name) ← **user**(userid,'SN2',name).  
**sn2update**(userid,statusid,creationdate,count) ← **status**(statusid,'SN2',userid,creationdate,count).

q(userid, statusid, viewcount) ← user(userid, 'SN1', name),  
 status(statusid, 'SN1', userid, '2016-03-03', viewcount)

### Candidate query plans

q(userid, statusid, viewcount) ← sn1user(userid,name), sn1myid(userid)

Not safe

q(userid, statusid, viewcount) ← sn1user(userid,name), sn1status(userid,statusid, '2016-03-03')

Not safe

q(userid, statusid, viewcount) ← sn1user(userid,name), sn1viewcount(statusid,viewcount)

q(userid, statusid, viewcount) ← sn1myid(userid), sn1myid(userid)

Not safe

q(userid, statusid, viewcount) ← sn1myid(userid), sn1status(userid,statusid, '2016-03-03')

Not safe

q(userid, statusid, viewcount) ← sn1myid(userid), sn1viewcount(statusid,viewcount)

### Candidate expanded query plans (bodys only)

user(userid,'SN1',name), user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2)

user(userid,'SN1',name), status(statusid,'SN1',userid, '2016-03-03',viewcount2)

user(userid,'SN1',name), status(statusid,'SN1',userid2,creationdate2,viewcount)

Not contained in Q

user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2)

user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2), status(statusid,'SN1',userid, '2016-03-03',viewcount2)

user(userid,'SN1',name2),

status(statusid2,'SN1',userid,creationdate2,viewcount2), status(statusid,'SN1',userid,creationdate2,viewcount)

Not contained in Q

$sn1user(userid,name) \leftarrow user(userid,'SN1',name).$   
 $sn1myid(userid) \leftarrow user(userid,'SN1',name),$   
 $status(statusid,'SN1',userid,creationdate,viewcount).$   
 $sn1status(userid,statusid,creationdate) \leftarrow status(statusid,'SN1',userid,creationdate,viewcount).$   
 $sn1viewcount(statusid,count) \leftarrow status(statusid,'SN1',userid,creationdate,count).$   
 $sn2myid(userid,name) \leftarrow user(userid,'SN2',name).$   
 $sn2update(userid,statusid,creationdate,count) \leftarrow status(statusid,'SN2',userid,creationdate,count).$

$q(userid, statusid, \leftarrow user(userid, 'SN1', name),$   
 $viewcount) \quad status(statusid, 'SN1', userid, '2016-03-03', viewcount)$

### Candidate query plans

~~$q(userid, statusid, viewcount) \leftarrow sn1user(userid,name), sn1myid(userid)$~~

Not safe

~~$q(userid, statusid, viewcount) \leftarrow sn1user(userid,name), sn1status(userid,statusid,'2016-03-03')$~~

Not safe

~~$q(userid, statusid, viewcount) \leftarrow sn1status(userid,statusid,'2016-03-03'), sn1myid(userid)$~~

Not safe

~~$q(userid, statusid, viewcount) \leftarrow sn1myid(userid), sn1status(userid,statusid,'2016-03-03')$~~

Not safe

~~$q(userid, statusid, viewcount) \leftarrow sn1status(userid,statusid,'2016-03-03'), sn1myid(userid)$~~

~~$q(userid, statusid, viewcount) \leftarrow sn1myid(userid), sn1status(userid,statusid,'2016-03-03')$~~

There is no contained rewriting for this query.

### Candidate expressions

~~$user(userid,'SN1',name), status(statusid,'SN1',userid,creationdate2,viewcount2)$~~

~~$user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2)$~~

~~$user(userid,'SN1',name), status(statusid,'SN1',userid2,creationdate2,viewcount)$~~

Not contained in Q

~~$user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2)$~~

~~$user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2), status(statusid,'SN1',userid,'2016-03-03',viewcount2)$~~

~~$user(userid,'SN1',name2),$~~

~~$status(statusid2,'SN1',userid,creationdate2,viewcount2), status(statusid,'SN1',userid2,creationdate2,viewcount2)$~~

Not contained in Q

**sn1user**(userid,name) ← **user**(userid,'SN1',name).  
**sn1myid**(userid) ← **user**(userid,'SN1',name),  
**status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1status**(userid,statusid,creationdate) ← **status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1viewcount**(statusid,count) ← **status**(statusid,'SN1',userid,creationdate,count).  
**sn2myid**(userid,name) ← **user**(userid,'SN2',name).  
**sn2update**(userid,statusid,creationdate,count) ← **status**(statusid,'SN2',userid,creationdate,count).

**q**(userid, statusid, viewcount) ← **user**(userid, 'SN1', name),  
**status**(statusid, 'SN1', userid, '2016-03-03', viewcount)

**user**(userid, 'SN1', name)      **status**(statusid, 'SN1', userid, '2016-03-03', viewcount)

**sn1user**(userid,name)      **sn1myid**(userid)

**sn1myid**(userid)      **sn1status**(userid,statusid,'2016-03-03')

**sn1viewcount**(statusid,viewcount)

### Extra candidate query plan

**q**(userid, statusid, viewcount) ← **sn1user**(userid,name),  
**sn1status**(userid,statusid, '2016-03-03'),  
**sn1viewcount**(statusid,viewcount)

### Expanded extra candidate query plan

**q**(userid, statusid, viewcount) ← **user**(userid,'SN1',name),  
**status**(statusid,'SN1',userid, '2016-03-03',viewcount2)  
**status**(statusid,'SN1',userid2,creationdate2,viewcount)

Even if we had the FD  
 statusid → userid, the  
 extra candidate is not  
 contained in Q.

# The bucket algo example 2

Same mappings as example 1 but simplified query.

**sn1user**(userid,name) ← **user**(userid,'SN1',name).  
**sn1myid**(userid) ← **user**(userid,'SN1',name),  
**status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1status**(userid,statusid,creationdate) ← **status**(statusid,'SN1',userid,creationdate,viewcount).  
**sn1viewcount**(statusid,count) ← **status**(statusid,'SN1',userid,creationdate,count).  
**sn2myid**(userid,name) ← **user**(userid,'SN2',name).  
**sn2update**(userid,statusid,creationdate,count) ← **status**(statusid,'SN2',userid,creationdate,count).

**q**(statusid, viewcount) ← **user**(userid, 'SN1', name),  
**status**(statusid, 'SN1', userid, creationdate, viewcount)

<b>user</b> (userid, 'SN1', name)	<b>status</b> (statusid, 'SN1', userid, creationdate, viewcount)
<b>sn1user</b> (userid,name)	<b>sn1myid</b> (userid)
<b>sn1myid</b> (userid)	<b>sn1status</b> (userid,statusid, creationdate)
	<b>sn1viewcount</b> (statusid,viewcount)

### Candidate query plans

**q**(statusid, viewcount) ← **sn1user**(userid,name), **sn1myid**(userid)

**q**(statusid, viewcount) ← **sn1user**(userid,name), **sn1status**(userid,statusid, creationdate)

**q**(statusid, viewcount) ← **sn1user**(userid,name), **sn1viewcount**(statusid,viewcount)

**q**(statusid, viewcount) ← **sn1myid**(userid), ~~**sn1myid**(userid)~~

**q**(statusid, viewcount) ← **sn1myid**(userid), **sn1status**(userid,statusid, creationdate)

**q**(statusid, viewcount) ← **sn1myid**(userid), **sn1viewcount**(statusid,viewcount)

$\text{sn1user}(\text{userid}, \text{name}) \leftarrow \text{user}(\text{userid}, \text{'SN1'}, \text{name}).$   
 $\text{sn1myid}(\text{userid}) \leftarrow \text{user}(\text{userid}, \text{'SN1'}, \text{name}),$   
 $\text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{creationdate}, \text{viewcount}).$   
 $\text{sn1status}(\text{userid}, \text{statusid}, \text{creationdate}) \leftarrow \text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{creationdate}, \text{viewcount}).$   
 $\text{sn1viewcount}(\text{statusid}, \text{count}) \leftarrow \text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{creationdate}, \text{count}).$   
 $\text{sn2myid}(\text{userid}, \text{name}) \leftarrow \text{user}(\text{userid}, \text{'SN2'}, \text{name}).$   
 $\text{sn2update}(\text{userid}, \text{statusid}, \text{creationdate}, \text{count}) \leftarrow \text{status}(\text{statusid}, \text{'SN2'}, \text{userid}, \text{creationdate}, \text{count}).$

$q(\text{statusid}, \text{viewcount}) \leftarrow \text{user}(\text{userid}, \text{'SN1'}, \text{name}),$   
 $\text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{creationdate}, \text{viewcount})$

### Candidate query plans

~~$q(\text{statusid}, \text{viewcount}) \leftarrow \text{sn1user}(\text{userid}, \text{name}), \text{sn1myid}(\text{userid})$~~

Not safe

~~$q(\text{statusid}, \text{viewcount}) \leftarrow \text{sn1user}(\text{userid}, \text{name}), \text{sn1status}(\text{userid}, \text{statusid}, \text{creationdate})$~~

Not safe

$q(\text{statusid}, \text{viewcount}) \leftarrow \text{sn1user}(\text{userid}, \text{name}), \text{sn1viewcount}(\text{statusid}, \text{viewcount})$

~~$q(\text{statusid}, \text{viewcount}) \leftarrow \text{sn1myid}(\text{userid}), \text{sn1myid}(\text{userid})$~~

Not safe

~~$q(\text{statusid}, \text{viewcount}) \leftarrow \text{sn1myid}(\text{userid}), \text{sn1status}(\text{userid}, \text{statusid}, \text{creationdate})$~~

Not safe

$q(\text{statusid}, \text{viewcount}) \leftarrow \text{sn1myid}(\text{userid}), \text{sn1viewcount}(\text{statusid}, \text{viewcount})$

### Candidate expanded query plans (bodys only)

~~$\text{user}(\text{userid}, \text{'SN1'}, \text{name}), \text{user}(\text{userid}, \text{'SN1'}, \text{name2}), \text{status}(\text{statusid2}, \text{'SN1'}, \text{userid}, \text{creationdate2}, \text{viewcount2})$~~

~~$\text{user}(\text{userid}, \text{'SN1'}, \text{name}), \text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{creationdate}, \text{viewcount2})$~~

$\text{user}(\text{userid}, \text{'SN1'}, \text{name}), \text{status}(\text{statusid}, \text{'SN1'}, \text{userid2}, \text{creationdate2}, \text{viewcount})$

Not contained in Q

~~$\text{user}(\text{userid}, \text{'SN1'}, \text{name2}), \text{status}(\text{statusid2}, \text{'SN1'}, \text{userid}, \text{creationdate2}, \text{viewcount2})$~~

~~$\text{user}(\text{userid}, \text{'SN1'}, \text{name2}), \text{status}(\text{statusid2}, \text{'SN1'}, \text{userid}, \text{creationdate2}, \text{viewcount2}), \text{status}(\text{statusid}, \text{'SN1'}, \text{userid}, \text{creationdate}, \text{viewcount2})$~~

$\text{user}(\text{userid}, \text{'SN1'}, \text{name2}),$

Not contained in Q

$\text{status}(\text{statusid2}, \text{'SN1'}, \text{userid}, \text{creationdate2}, \text{viewcount2}), \text{status}(\text{statusid}, \text{'SN1'}, \text{userid2}, \text{creationdate3}, \text{viewcount})$

# The bucket algo example 3

One different mapping from example 1 but same query (as in example 1).

<b>sn1user</b> (userid,name)←	<b>user</b> (userid,'SN1',name).
<b>sn1myid</b> (userid)←	<b>user</b> (userid,'SN1',name),
<b>sn1status</b> (userid,statusid,creationdate)←	<b>status</b> (statusid,'SN1',userid,creationdate,viewcount).
<b>sn1viewcount</b> (statusid,userid,creationdate,count)←	<b>status</b> (statusid,'SN1',userid,creationdate,viewcount).
<b>sn2myid</b> (userid,name)←	<b>user</b> (userid,'SN2',name).
<b>sn2update</b> (userid,statusid,creationdate,count)←	<b>status</b> (statusid,'SN2',userid,creationdate,count).

q(userid, statusid, viewcount) ← user(userid, 'SN1', name),  
 status(statusid, 'SN1', userid, '2016-03-03', viewcount)

<b>user</b> (userid, 'SN1', name)	<b>status</b> (statusid, 'SN1', userid, '2016-03-03', viewcount)
<b>sn1user</b> (userid,name)	<b>sn1myid</b> (userid)
<b>sn1myid</b> (userid)	<b>sn1status</b> (userid,statusid,'2016-03-03')
	<b>sn1viewcount</b> (statusid,userid, '2016-03-03',viewcount)

### Candidate query plans

q(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1myid**(userid)

q(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1status**(userid,statusid,'2016-03-03')

q(userid, statusid, viewcount) ← **sn1user**(userid,name), **sn1viewcount**(statusid,userid,'2016-03-03',viewcount)

q(userid, statusid, viewcount) ← **sn1myid**(userid), ~~sn1myid~~(userid)

q(userid, statusid, viewcount) ← **sn1myid**(userid), **sn1status**(userid,statusid,'2016-03-03')

q(userid, statusid, viewcount) ← **sn1myid**(userid), **sn1viewcount**(statusid,userid,'2016-03-03',viewcount)



<b>sn1user</b> (userid,name)←	<b>user</b> (userid,'SN1',name).
<b>sn1myid</b> (userid)←	<b>user</b> (userid,'SN1',name),
	<b>status</b> (statusid,'SN1',userid,creationdate,viewcount).
<b>sn1status</b> (userid,statusid,creationdate)←	<b>status</b> (statusid,'SN1',userid,creationdate,viewcount).
<b>sn1viewcount</b> (statusid,userid,creationdate,count)←	<b>status</b> (statusid,'SN1',userid,creationdate,count).
	<b>sn2myid</b> (userid,name)←
	<b>user</b> (userid,'SN2',name).
<b>sn2update</b> (userid,statusid,creationdate,count)←	<b>status</b> (statusid,'SN2',userid,creationdate,count).

q(userid, statusid, viewcount) ← user(userid, 'SN1', name),  
status(statusid, 'SN1', userid, '2016-03-03', viewcount)

### Candidate query plans

q(userid, statusid, viewcount) ← sn1user(userid,name), sn1myid(userid)	Not safe
q(userid, statusid, viewcount) ← sn1user(userid,name), sn1status(userid,statusid,'2016-03-03')	Not safe
q(userid, statusid, viewcount) ← sn1user(userid,name), sn1viewcount(statusid,userid,'2016-03-03',viewcount)	
q(userid, statusid, viewcount) ← sn1myid(userid)	Not safe
q(userid, statusid, viewcount) ← sn1myid(userid), sn1status(userid,statusid,'2016-03-03')	Not safe
q(userid, statusid, viewcount) ← sn1myid(userid), sn1viewcount(statusid,userid,'2016-03-03',viewcount)	

### Candidate expanded query plans (bodys only)

<del>user(userid,'SN1',name), user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2)</del>	
<del>user(userid,'SN1',name), status(statusid,'SN1',userid,'2016-03-03',viewcount2)</del>	
user(userid,'SN1',name), status(statusid,'SN1',userid,'2016-03-03',viewcount)	Contained in Q
<del>user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2)</del>	
<del>user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2), status(statusid,'SN1',userid,'2016-03-03',viewcount2)</del>	
user(userid,'SN1',name2), status(statusid2,'SN1',userid,creationdate2,viewcount2), status(statusid,'SN1',userid,'2016-03-03',viewcount)	Contained in Q

<b>sn1user</b> (userid,name)←	<b>user</b> (userid,'SN1',name).
<b>sn1myid</b> (userid)←	<b>user</b> (userid,'SN1',name),
	<b>status</b> (statusid,'SN1',userid,creationdate,viewcount).
<b>sn1status</b> (userid,statusid,creationdate)←	<b>status</b> (statusid,'SN1',userid,creationdate,viewcount).
<b>sn1viewcount</b> (statusid,userid,creationdate,count)←	<b>status</b> (statusid,'SN1',userid,creationdate,count).
	<b>sn2myid</b> (userid,name)←
	<b>user</b> (userid,'SN2',name).
<b>sn2update</b> (userid,statusid,creationdate,count)←	<b>status</b> (statusid,'SN2',userid,creationdate,count).

q(userid, statusid, viewcount) ← user(userid, 'SN1', name),  
 status(statusid, 'SN1', userid, '2016-03-03', viewcount)

### Candidate query plans

q(userid, statusid, viewcount) ← sn1user(userid,name), sn1myid(userid)	Not safe
q(userid, statusid, viewcount) ← sn1user(userid,name), sn1status(userid,statusid,'2016-03-03')	Not safe
q(userid, statusid, viewcount) ← sn1user(userid,name), sn1viewcount(statusid,userid,'2016-03-03',viewcount)	
q(userid, statusid, viewcount) ← sn1myid(userid)	Not safe
q(userid, statusid, viewcount) ← sn1myid(userid), sn1status(userid,statusid,'2016-03-03')	Not safe
q(userid, statusid, viewcount) ← sn1myid(userid), sn1viewcount(statusid,userid,'2016-03-03',viewcount)	

### Candidate expanded query plans (bodys only)

<b>user</b> (userid,'SN1',name), <b>user</b> (userid,'SN1',name2), <b>status</b> (statusid2,'SN1',userid,creationdate2,viewcount2)	
<b>user</b> (userid,'SN1',name), <b>status</b> (statusid,'SN1',userid,'2016-03-03',viewcount2)	
<b>user</b> (userid,'SN1',name), <b>status</b> (statusid,'SN1',userid,'2016-03-03',viewcount)	Max. contained in Q
<b>user</b> (userid,'SN1',name2), <b>status</b> (statusid2,'SN1',userid,creationdate2,viewcount2)	
<b>user</b> (userid,'SN1',name2), <b>status</b> (statusid2,'SN1',userid,creationdate2,viewcount2), <b>status</b> (statusid,'SN1',userid,'2016-03-03',viewcount2)	
<b>user</b> (userid,'SN1',name2),	
<b>status</b> (statusid2,'SN1',userid,creationdate2,viewcount2), <b>status</b> (statusid,'SN1',userid,'2016-03-03',viewcount)	Not maximally contained

# Inverse rules algo.

# Inverse rules algo. overview

Algorithm to build maximally contained rewritings (i.e. maximally contained datalog query plans)

Inputs:  $Q \in \text{datalog}$  and  $s \in \mathcal{C}Q$

Output: the maximally contained query plan  $\mathcal{P}$  (i.e.  $\mathcal{P} \in \text{datalog}$ ) in  $Q$  wrt  $s$

1. Building a logic program (datalog program with functions) with views as only EDB predicates
  - a. remove from  $Q$  rules containing EDB predicates that do not appear in any view definition
  - b. add inverse view definitions
2. Transforming the logic program into a datalog program (without function)
3. Optimizing
  - a. remove unnecessary parts
  - b. reduce the size of the obtained program

# Inverse rules algo. Step 1a

1a. Building a logic program (datalog program with functions) with views as only EDB predicates

Running example

- ▶ Query to rewrite  $\mathcal{P}$ :

$q(X, Y) : \neg edge(X, Y)$

$q(X, Y) : \neg edge(X, Z), q(Z, Y)$

- ▶ Views  $s$ :

$v(X, Y) : \neg edge(X, Z), edge(Z, Y)$

- ▶ Logic program  $(\mathcal{P}^-, s^{-1})$ :

$q(X, Y) : \neg edge(X, Y)$

$q(X, Y) : \neg edge(X, Z), q(Z, Y)$

$edge(f(X, Y), Y) : \neg v(X, Y)$

$edge(X, f(X, Y)) : \neg v(X, Y)$

# Inverse rules algo. Step 1a

The obtained logic program is  $(\mathcal{P}^-, s^{-1})$ :

$q(X, Y) : \text{--edge}(X, Y)$

$q(X, Y) : \text{--edge}(X, Z), q(Z, Y)$

$\text{edge}(f(X, Y), Y) : \text{--v}(X, Y)$

$\text{edge}(X, f(X, Y)) : \text{--v}(X, Y)$

$(\mathcal{P}^-, s^{-1})$  produces answers with and without function symbols ( $f$  for the running example).

Let's suppose that  $(\mathcal{P}^-, s^{-1}) \downarrow$  is a datalog program that computes the answers of  $(\mathcal{P}^-, s^{-1})$  with no function symbol.

# Inverse rules algo. Step 1b

1b. Inversing view definition  $v(X_1, \dots, X_m)$ :

- ▶ distinguished variables  $X_1, \dots, X_m$  remain unchanged
- ▶ each existential variable must be replaced by a term  $f(X_1, \dots, X_m)$  where  $f$  is a new function symbol for each existential var.

Another example

- ▶ Views  $s$ :

$$v_1(X, Y) : \neg \text{edge}(X, Z), \text{edge}(Z, W), \text{edge}(W, Y)$$

$$v_2(X) : \neg \text{edge}(X, Z)$$

- ▶ Inversed views  $s^{-1}$

$$\text{edge}(X, f_1(X, Y)) : \neg v_1(X, Y)$$

$$\text{edge}(f_1(X, Y), f_2(X, Y)) : \neg v_1(X, Y)$$

$$\text{edge}(f_2(X, Y), Y) : \neg v_1(X, Y)$$

$$\text{edge}(X, f_3(X)) : \neg v_2(X)$$

# Inverse rules algo. Step 2

## 2. Eliminating function symbols

In a bottom-up evaluation like process:

**2.a** functions terms in  $s^{-1}$  are replaced by the list of their variables and the associated predicate is replaced by an annotated one

Running example:

$edge(X, f(X, Y)) : -v(X, Y)$

is replaced by

$edge^{<*, f(*,*)>}(X, X, Y) : -v(X, Y)$

**2.b** new annotated predicates replace original ones only in the rules of  $(\mathcal{P}^-, s^{-1})$



# Inverse rules algo. Example

Beginning  $(\mathcal{P}^-, s^{-1})$ :

$$(1) q(X, Y) : -edge(X, Y)$$

$$(2) q(X, Y) : -edge(X, Z), q(Z, Y)$$

$$(3) edge(f(X, Y), Y) : -v(X, Y)$$

$$(4) edge(X, f(X, Y)) : -v(X, Y)$$

Result of 2.a 2 rules have been generated:

$$(3') edge^{<f(*,*),*>}(X, Y, Y) : -v(X, Y)$$

$$(4') edge^{<*,f(*,*)>}(X, X, Y) : -v(X, Y)$$

after elimination of function symbol  $f$ .

# Inverse rules algo. Example

- $(\mathcal{P}^-, s^{-1})$
- (1)  $q(X, Y) : \neg edge(X, Y)$
  - (2)  $q(X, Y) : \neg edge(X, Z), q(Z, Y)$
  - (3)  $edge(f(X, Y), Y) : \neg v(X, Y)$
  - (4)  $edge(X, f(X, Y)) : \neg v(X, Y)$

- Result of 2.a
- (3')  $edge^{<f(*,*),*>}(X, Y, Y) : \neg v(X, Y)$
  - (4')  $edge^{<*,f(*,*)>}(X, X, Y) : \neg v(X, Y)$

## Result of 2.b Step 1

- (5)  $q^{<*,f(*,*)>}(X, Y_1, Y_2) : \neg edge^{<*,f(*,*)>}(X, Y_1, Y_2)$

- ▶ by replacing  $edge(X, Y)$  in (1) by  $edge^{<*,f(*,*)>}(\dots, \dots, \dots)$  from (4')

- (6)  $q^{<f(*,*),*>}(X_1, X_2, Y) : \neg edge^{<f(*,*),*>}(X_1, X_2, Y)$

- ▶ by replacing  $edge(X, Y)$  in (1) by  $edge^{<f(*,*),*>}(\dots, \dots, \dots)$  from (3')

# Inverse rules algo. Example

$(\mathcal{P}^-, s^{-1})$  (1)  $q(X, Y) : \neg \text{edge}(X, Y)$

(2)  $q(X, Y) : \neg \text{edge}(X, Z), q(Z, Y)$

(3)  $\text{edge}(f(X, Y), Y) : \neg v(X, Y)$

(4)  $\text{edge}(X, f(X, Y)) : \neg v(X, Y)$

2.a (3')  $\text{edge}^{\langle f(*, *), * \rangle}(X, Y, Y) : \neg v(X, Y)$

(4')  $\text{edge}^{\langle *, f(*, *) \rangle}(X, X, Y) : \neg v(X, Y)$

2.b Step 1 (5)  $q^{\langle *, f(*, *) \rangle}(X, Y_1, Y_2) : \neg \text{edge}^{\langle *, f(*, *) \rangle}(X, Y_1, Y_2)$

(6)  $q^{\langle f(*, *), * \rangle}(X_1, X_2, Y) : \neg \text{edge}^{\langle f(*, *), * \rangle}(X_1, X_2, Y)$

2.b Step 2 (7)  $q^{\langle *, * \rangle}(X, Y) : \neg$

$\text{edge}^{\langle *, f(*, *) \rangle}(X, Z_1, Z_2), q^{\langle f(*, *), * \rangle}(Z_1, Z_2, Y)$

▶ by replacing  $\text{edge}(X, Z)$  in (2) by  $\text{edge}^{\langle *, f(*, *) \rangle}(\dots, \dots, \dots)$  from (4')

▶ and by replacing  $q(Z, Y)$  in (2) by  $q^{\langle f(*, *), * \rangle}(\dots, \dots, \dots)$  from (6)

(8)  $q^{\langle f(*, *), f(*, *) \rangle}(X_1, X_2, Y_1, Y_2) : \neg$

$\text{edge}^{\langle f(*, *), * \rangle}(X_1, X_2, Z), q^{\langle *, f(*, *) \rangle}(Z, Y_1, Y_2)$

▶ by replacing  $\text{edge}(X, Z)$  in (2) by  $\text{edge}^{\langle f(*, *), * \rangle}(\dots, \dots, \dots)$  from (3')

▶ and by replacing  $q(Z, Y)$  in (2) by  $q^{\langle *, f(*, *) \rangle}(\dots, \dots, \dots)$  from (5)

# Inverse rules algo. Example

- ( $\mathcal{P}^-, s^{-1}$ )
- (1)  $q(X, Y) : \neg edge(X, Y)$
  - (2)  $q(X, Y) : \neg edge(X, Z), q(Z, Y)$
  - (3)  $edge(f(X, Y), Y) : \neg v(X, Y)$
  - (4)  $edge(X, f(X, Y)) : \neg v(X, Y)$
- 2.a
- (3')  $edge^{<f(*,*),*>}(X, Y, Y) : \neg v(X, Y)$
  - (4')  $edge^{<*,f(*,*)>}(X, X, Y) : \neg v(X, Y)$
- 2.b Step 1
- (5)  $q^{<*,f(*,*)>}(X, Y_1, Y_2) : \neg edge^{<*,f(*,*)>}(X, Y_1, Y_2)$
  - (6)  $q^{<f(*,*),*>}(X_1, X_2, Y) : \neg edge^{<f(*,*),*>}(X_1, X_2, Y)$
- 2.b Step 2
- (7)  $q^{<*,*>}(X, Y) : \neg edge^{<*,f(*,*)>}(X, Z_1, Z_2), q^{<f(*,*),*>}(Z_1, Z_2, Y)$
  - (8)  $q^{<f(*,*)f(*,*)>}(X_1, X_2, Y_1, Y_2) : \neg edge^{<f(*,*),*>}(X_1, X_2, Z), q^{<*,f(*,*)>}(Z, Y_1, Y_2)$
- 2.b Step 3
- (9)  $q^{<f(*,*),*>}(X_1, X_2, Y) : \neg edge^{<f(*,*),*>}(X_1, X_2, Z), q^{<*,*>}(Z, Y)$ 
    - ▶ by replacing  $edge(X, Z)$  in (2) by  $edge^{<f(*,*),*>}(\dots, \dots, \dots)$  from (3')
    - ▶ and by replacing  $q(Z, Y)$  in (2) by  $q^{<*,*>}(\dots, \dots)$  from (7)
  - (10)  $q^{<*,f(*,*)>}(X, Y_1, Y_2) : \neg edge^{<*,f(*,*)>}(X, Z_1, Z_2), q^{<f(*,*)f(*,*)>}(Z_1, Z_2, Y_1, Y_2)$ 
    - ▶ by replacing  $edge(X, Z)$  in (2) by  $edge^{<*,f(*,*)>}(\dots, \dots, \dots)$  from (4')
    - ▶ and by replacing  $q(Z, Y)$  in (2) by  $q^{<f(*,*)f(*,*)>}(\dots, \dots, \dots, \dots)$  from (8)

2.b End

# Inverse rules algo. Example

The result of part 2 of the algorithm is the following datalog program, called  $(\mathcal{P}^-, s^{-1})^{datalog}$  or  $(\mathcal{P}^-, s^{-1})^{split}$ :

- (3')  $edge^{<f(*,*) , * >} (X, Y, Y) : -v(X, Y)$
- (4')  $edge^{<*, f(*,*) >} (X, X, Y) : -v(X, Y)$
- (5)  $q^{<*, f(*,*) >} (X, Y_1, Y_2) : -edge^{<*, f(*,*) >} (X, Y_1, Y_2)$
- (6)  $q^{<f(*,*) , * >} (X_1, X_2, Y) : -edge^{<f(*,*) , * >} (X_1, X_2, Y)$
- (7)  $q^{<*, * >} (X, Y) : - edge^{<*, f(*,*) >} (X, Z_1, Z_2), q^{<f(*,*) , * >} (Z_1, Z_2, Y)$
- (8)  $q^{<f(*,*) , f(*,*) >} (X_1, X_2, Y_1, Y_2) : -$   
 $edge^{<f(*,*) , * >} (X_1, X_2, Z), q^{<*, f(*,*) >} (Z, Y_1, Y_2)$
- (9)  $q^{<f(*,*) , * >} (X_1, X_2, Y) : -$   
 $edge^{<f(*,*) , * >} (X_1, X_2, Z), q^{<*, * >} (Z, Y)$
- (10)  $q^{<*, f(*,*) >} (X, Y_1, Y_2) : -$   
 $edge^{<*, f(*,*) >} (X, Z_1, Z_2), q^{<f(*,*) , f(*,*) >} (Z_1, Z_2, Y_1, Y_2)$

Part 2 always terminates since function symbols are introduced into inverse rules that are not recursive.

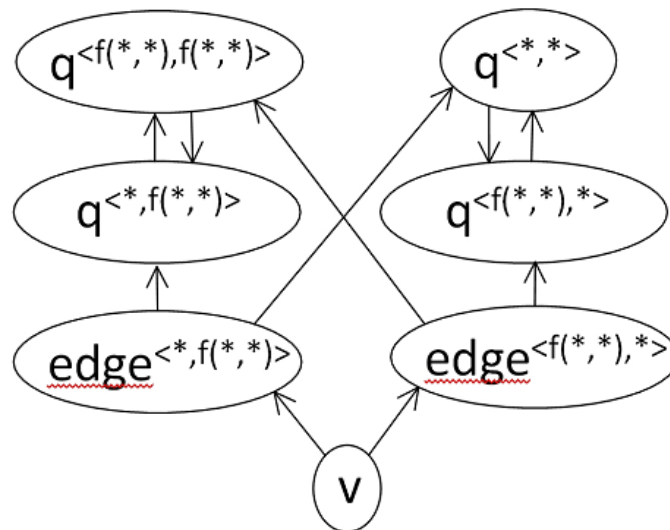
# Inverse rules algo. Optim.

Optimizations can be applied to  $(\mathcal{P}^-, s^{-1})^{datalog}$  :

- 3.a If a predicate  $p$  cannot contribute answers to the query (because there is no path from  $p$  to the query predicate in the dependency graph) then rules defining  $p$  can be deleted
- 3.b IDB predicates defined by one rule with one atom body can be replaced by their body atom (be careful of variables).

# Inverse rules algo. Example

Here is the dependency graph of the program:



So rules defining predicates  $q^{<*,f(*,*)>}$  and  $q^{<f(*,*),f(*,*)>}$  can be deleted. We obtain then:

$$(3') \text{ edge}^{<f(*,*),*>}(X, Y, Y) : \neg v(X, Y)$$

$$(4') \text{ edge}^{<*,f(*,*)>}(X, X, Y) : \neg v(X, Y)$$

$$(6) q^{<f(*,*),*>}(X_1, X_2, Y) : \neg \text{edge}^{<f(*,*),*>}(X_1, X_2, Y)$$

$$(7) q^{<*,*>}(X, Y) : \neg \text{edge}^{<*,f(*,*)>}(X, Z_1, Z_2), q^{<f(*,*),*>}(Z_1, Z_2, Y)$$

$$(9) q^{<f(*,*),*>}(X_1, X_2, Y) : \neg \text{edge}^{<f(*,*),*>}(X_1, X_2, Z), q^{<*,*>}(Z, Y)$$

# Inverse rules algo. Example

In the current program, rule 3.b cannot be applied because of variables. So the result is not changed wrt 3.a.

$(\mathcal{P}^-, s^{-1})^{datalog}$  :

(3')  $edge^{<f(*,*),*>}(X, Y, Y) : -v(X, Y)$

(4')  $edge^{<*,f(*,*)>}(X, X, Y) : -v(X, Y)$

(6)  $q^{<f(*,*),*>}(X_1, X_2, Y) : -edge^{<f(*,*),*>}(X_1, X_2, Y)$

(7)  $q^{<*,*>}(X, Y) : -edge^{<*,f(*,*)>}(X, Z_1, Z_2), q^{<f(*,*),*>}(Z_1, Z_2, Y)$

(9)  $q^{<f(*,*),*>}(X_1, X_2, Y) : -edge^{<f(*,*),*>}(X_1, X_2, Z), q^{<*,*>}(Z, Y)$

is the maximally contained query plan for our problem.



# Inverse rules algo. Results

For a set of source description  $\mathcal{D}$ , let  $\mathcal{P}$  be a query plan.  
Then we note  $\mathcal{P} \downarrow$  the plan that computes the same answers than  $\mathcal{P}$  on  $\mathcal{D}$  without tuples that use a function symbol.

About  $(\mathcal{P}^-, s^{-1})$  and  $(\mathcal{P}^-, s^{-1}) \downarrow$  ▶  $\forall \mathcal{P}, \forall s, \forall I_s(\text{finite})$ ,  
 $(\mathcal{P}^-, s^{-1})$  has a unique minimal fixpoint (cf semantics of a logic program) and naive and semi naive evaluation always terminate and produce this fixpoint.

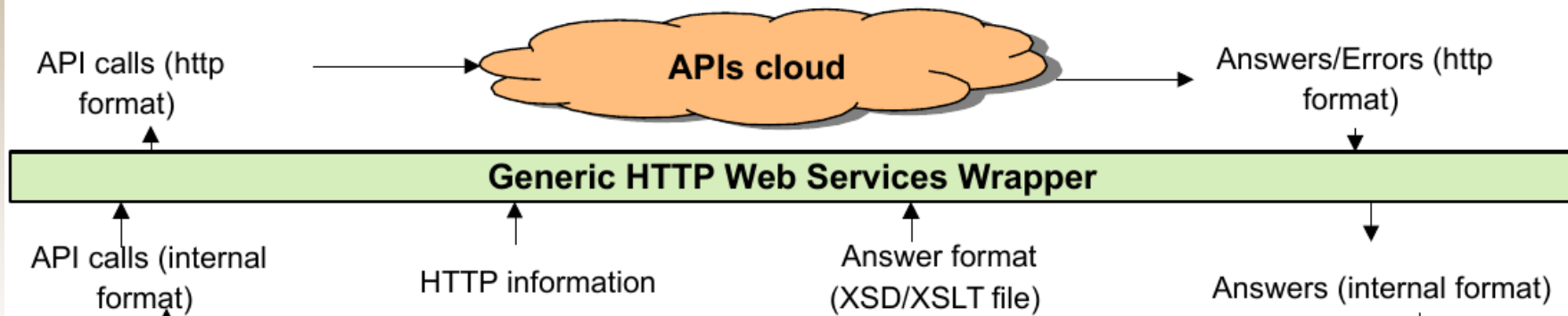
- ▶  $(\mathcal{P}^-, s^{-1})$  can be built in time polynomial in the sizes of  $\mathcal{P}$  and  $s$ .
- ▶  $\forall \mathcal{P}, \forall s, (\mathcal{P}^-, s^{-1}) \downarrow$  is maximally contained in  $\mathcal{P}$

About  $(\mathcal{P}^-, s^{-1})^{\text{datalog}}$  ▶  $\forall \mathcal{P}, \forall s, (\mathcal{P}^-, s^{-1}) \downarrow$  is equivalent to  $(\mathcal{P}^-, s^{-1})^{\text{datalog}}$

- ▶  $\forall \mathcal{P}, \forall s, (\mathcal{P}^-, s^{-1})^{\text{datalog}}$  is maximally contained in  $\mathcal{P}$ .
- ▶ If there exists a query plan equivalent to  $\mathcal{P}$  then  $(\mathcal{P}^-, s^{-1})^{\text{datalog}}$  is equivalent to  $\mathcal{P}$ .

# Architecture of the generic wrapper

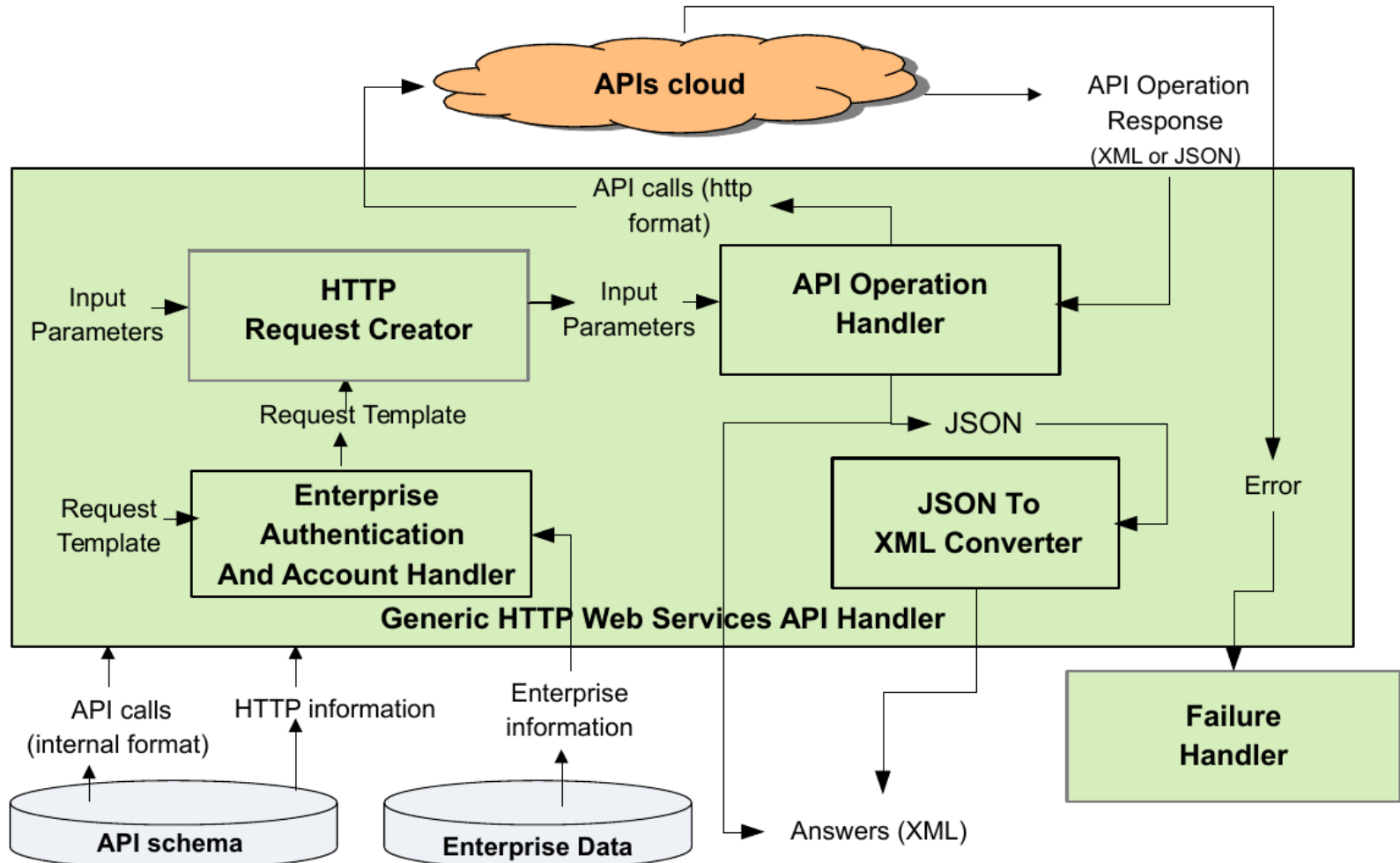
# Wrapper overview



# Inputs and output

- Inputs
  - Enterprise Information (Authentication parameters, subdomains, profile related information)
  - HTTP URL template, method
  - HTTP body template (XML)
  - XML Schema (XSD)
  - XML Transformation (XSLT)
- Output
  - Transformed (desired) information from WS
  - Or Error

# More detailed architecture



**Caption:**  Component  
 A → C1 A is input of C1  
 C1 → A A is output of C1  
 D (cylinder) -.-> C1 C1 uses D