



LarKC – a platform for large-scale Semantic Web Reasoning

Alexey Cheptsov

High Performance Computing Center Stuttgart





Outline

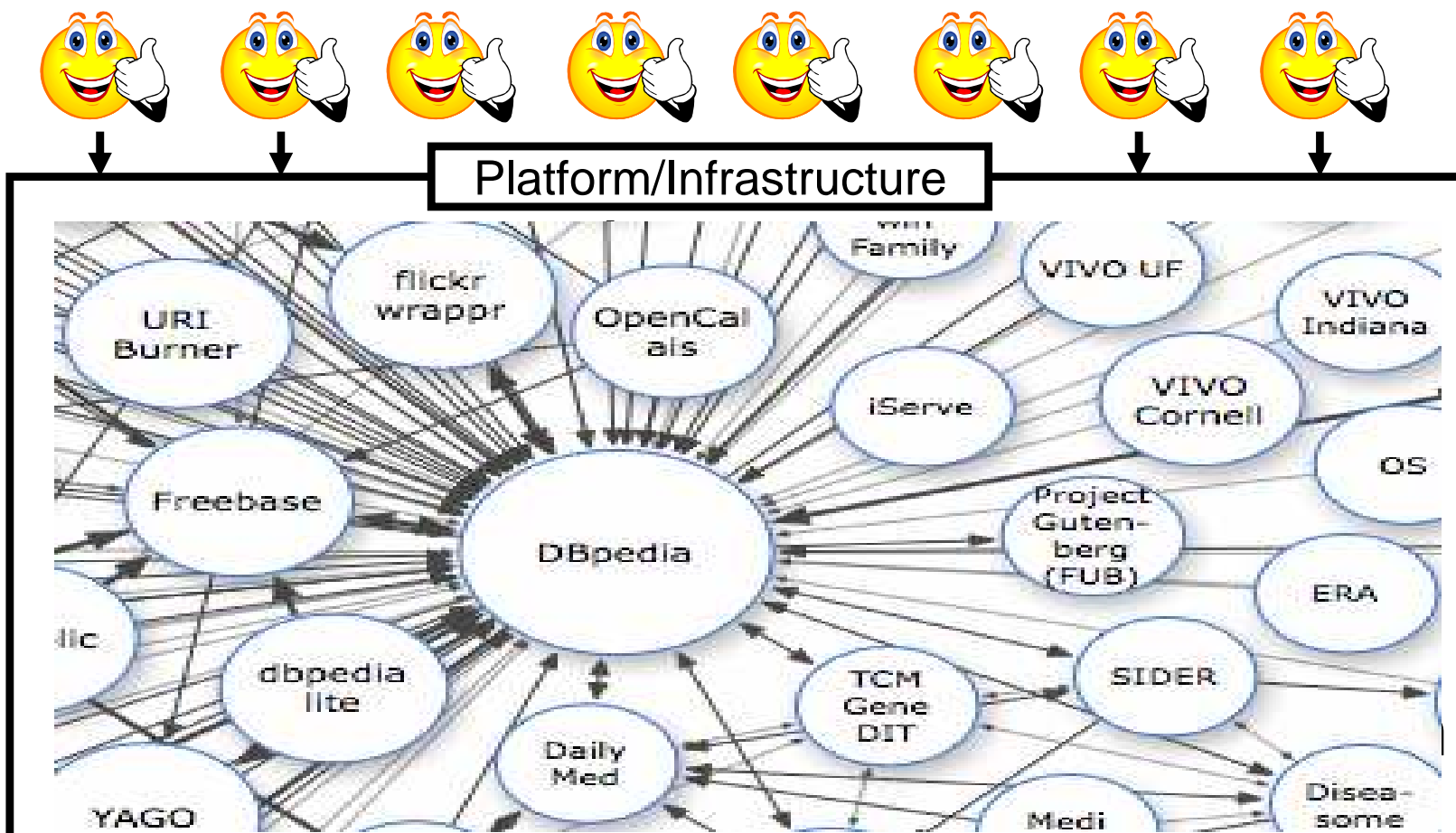
- **Semantic Web.** New challenges
- **LarKC.** A platform for web-scale Reasoning
 - Main motivation
 - Architecture overview
- **LarKC Platform.** Main subsystems
- **LarKC Software Redistribution.** [LarKC@SourceForge](#)
- **Conclusions.**

Large Knowledge Collider



① Semantic Web. New challenges

Semantic Web enables machine-supported inferencing over the WWW-distributed data



① *Semantic Web. New challenges*

Web-scale reasoning

Gartner (May 2007):

"By 2012,
70% of public Web pages will have some level of semantic markup,
20% will use more extensive Semantic Web-based ontologies"

- Semantic Technologies at Web Scale?
 - 20% of 30 billion pages @ 1000 triples per page = **6 trillion triples**
 - 30 billion and 1000 are underestimated, imagine in 6 years from now...
 - data-integration and semantic search at web-scale?

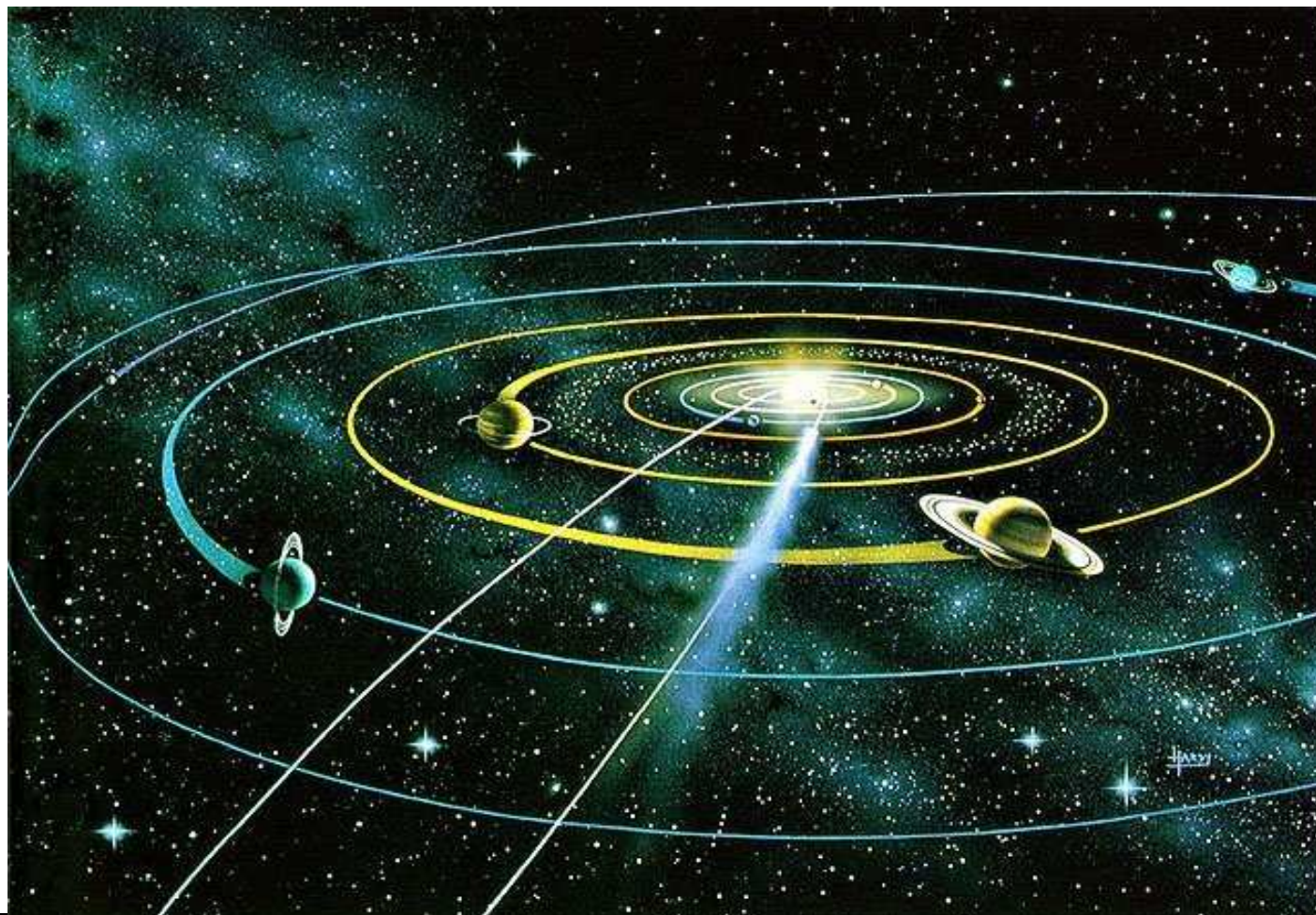
① *Semantic Web. New challenges*

Web-scale reasoning



1 Triple

$\sim 10^{14}$ Triples



① *Semantic Web. New challenges*

Data on the web scale

not only from

large numbers

- from performant data layer (OWLIM)
- from parallel deployment and execution of reasoners (IRIS)
- from load-balancing strategies
- ...

but also from

interaction of multiple components

allowing for **incompleteness** and **anytime behaviour**
in the reasoning process

② *The idea of LarkKC*

An experimental platform for large-scale reasoning

requires

not only: deductive inference over given axioms

but also:



Reasoning + Search

where do the axioms come from? (**IDENTIFY**)

which part of knowledge & data is required (**SELECT**)

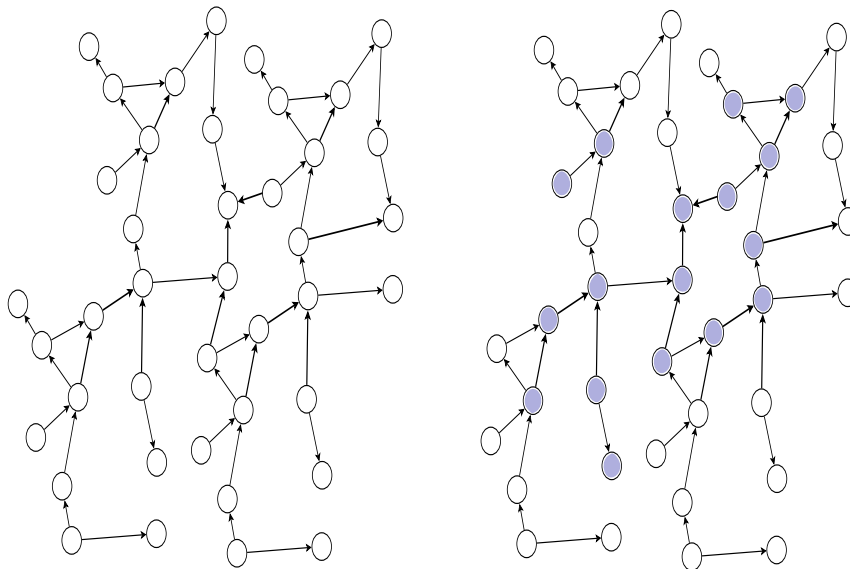
when is an answer "good enough" or "best possible" (**DECIDE**)

non-deductive inference (inductive, statistical) (**REASON**)

② The idea of LarkC

Statistical Semantics tasks

Subsetting



Query expansion

SPARQL Query

```
SELECT ?S ?P ?O  
WHERE {  
  ?S ?P ?O .  
  FILTER(?O='ultrasound')  
}
```

↑ refine

```
UNION  
{?S ?P ?O .  
  FILTER (?S="reflection")}
```

extract keywords →

ultrasound

find similar URIs/literals →

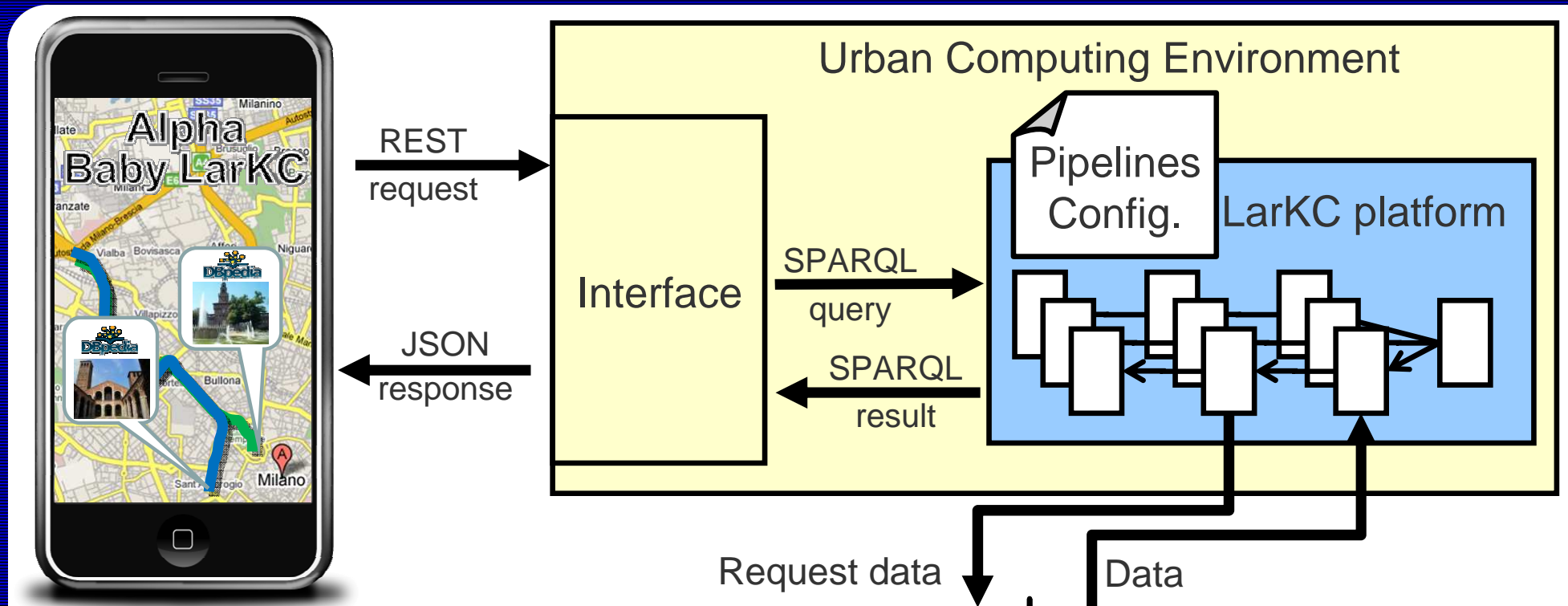
Semantic Index

Similarity	Literal/URI
1.0	ultrasound
0.96	reflection
0.94	sonography
...	...

Large Knowledge Collider



② The idea of LarKC

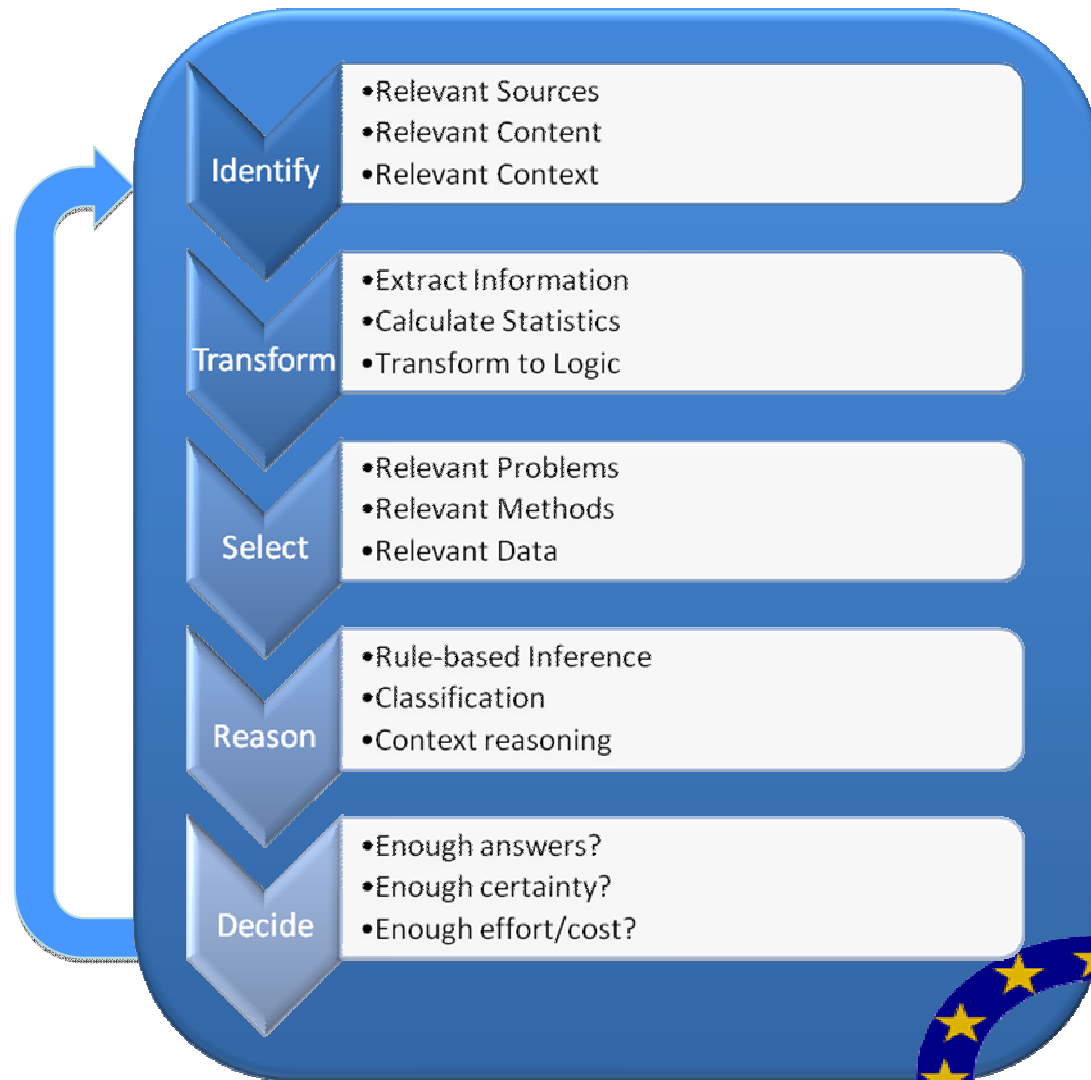


PROBLEM:
Which Milano monuments or events or friends can I quickly get to from here?



Large Knowledge Collider

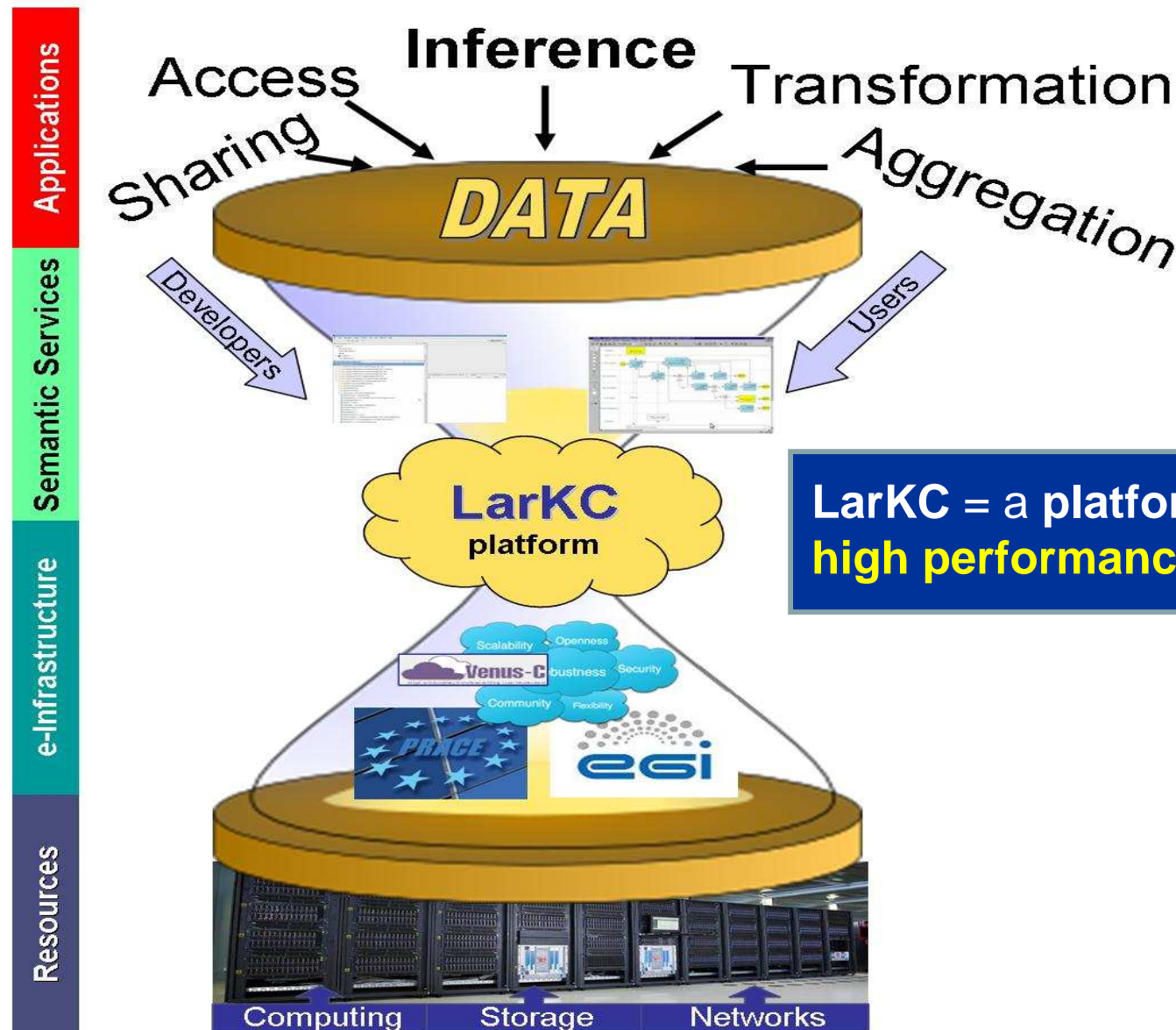
② *The idea of LarKC*



Large Knowledge Collider



② The idea of LarKC

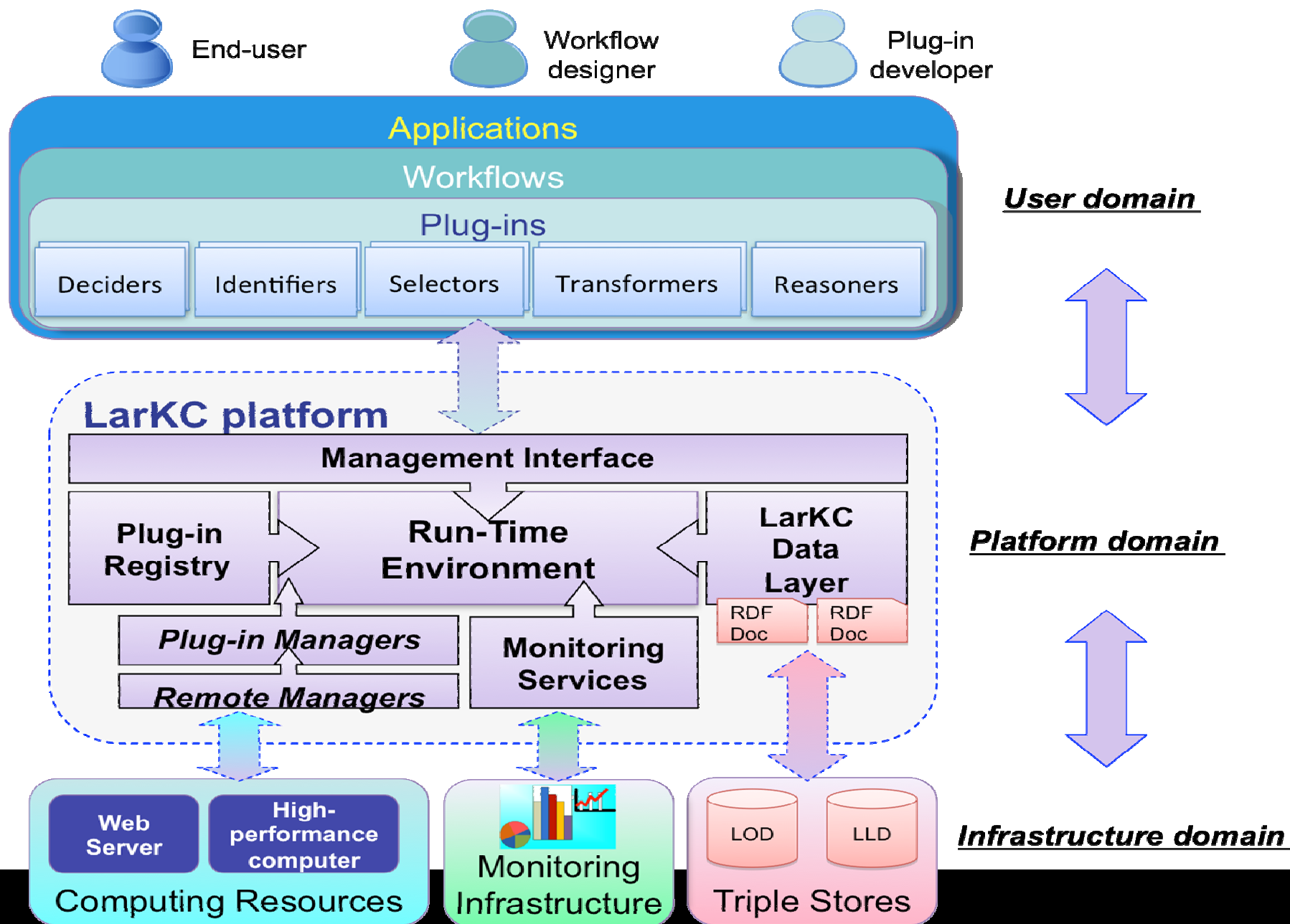


LarKC = a platform for large scale, **high performance** reasoning

Large Knowledge Collider

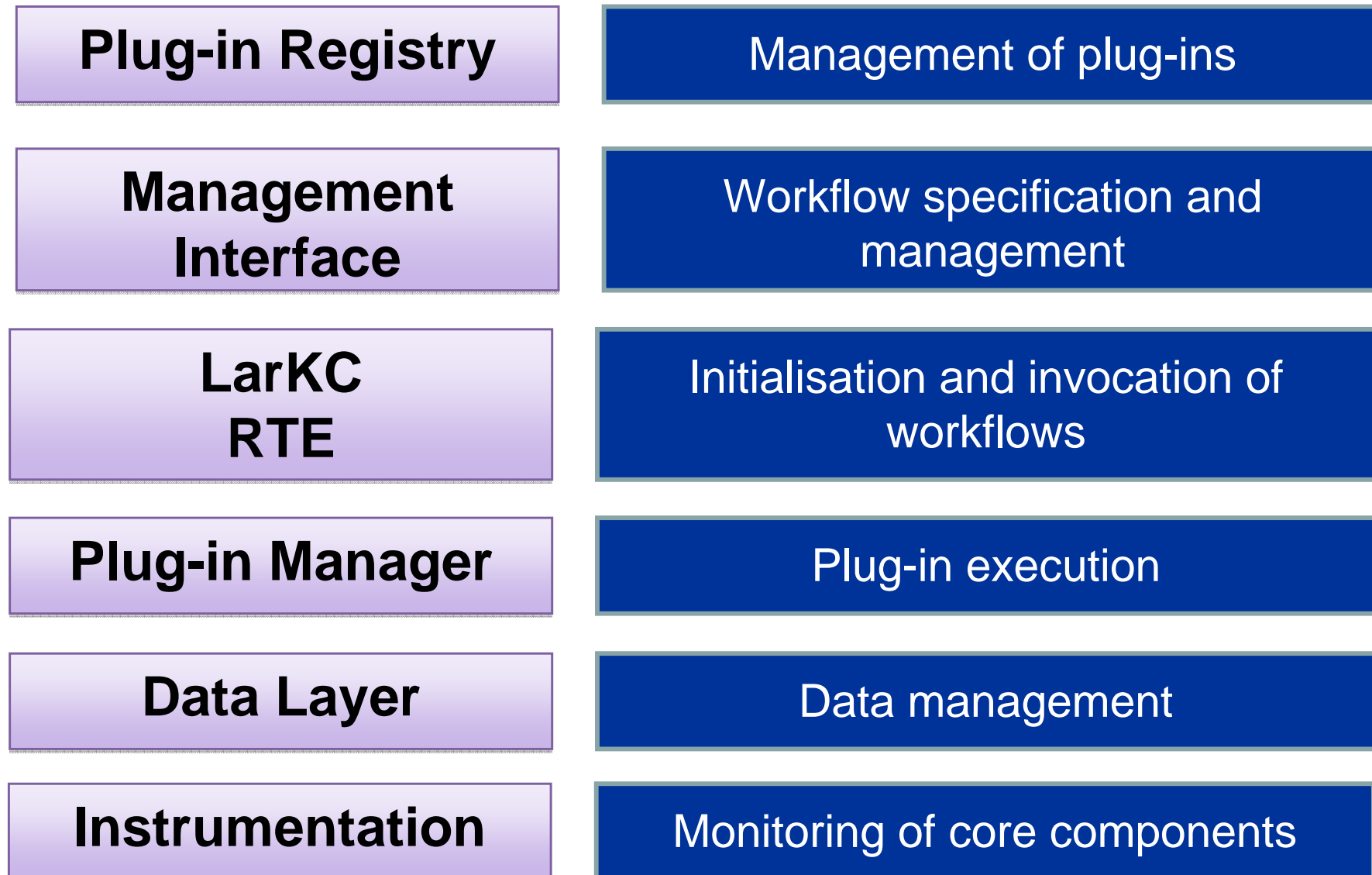


③ The architecture of LarKC



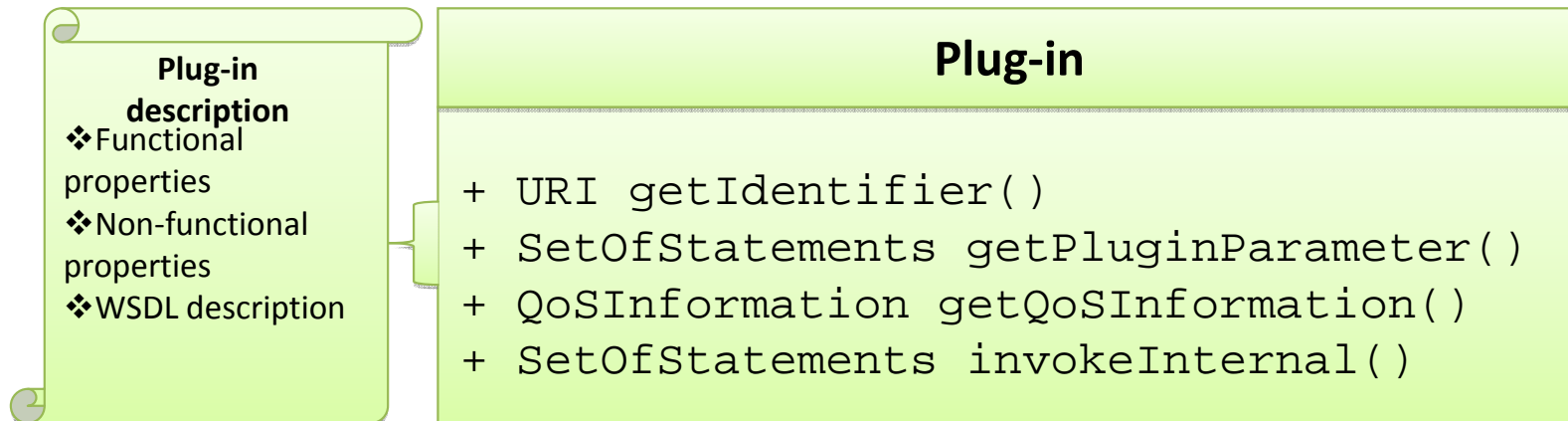


③ *The architecture of LarKC*



③ *The architecture of LarKC*

Plug-Ins



- Plug-ins are assembled into workflows to realise a LarKC experiment or application
- Plug-ins are identified by a URI (Uniform Resource Identifier)
- Plug-ins provide metadata about what they do and what they consume/produce (Functional properties): e.g. type = Selector
- Plug-ins provide information about their needs, including QoS information (Non-functional properties): e.g. Throughput, MinMemory, etc.
- Plug-ins are equipped with functionalities for data caching and messaging



③ *The architecture of LarKC*

New plug-in API – both input and output are represented in RDF

Automatic parallelization (thread-based) on the statement level

Data caching, instrumentation and event processing

Maven based build system

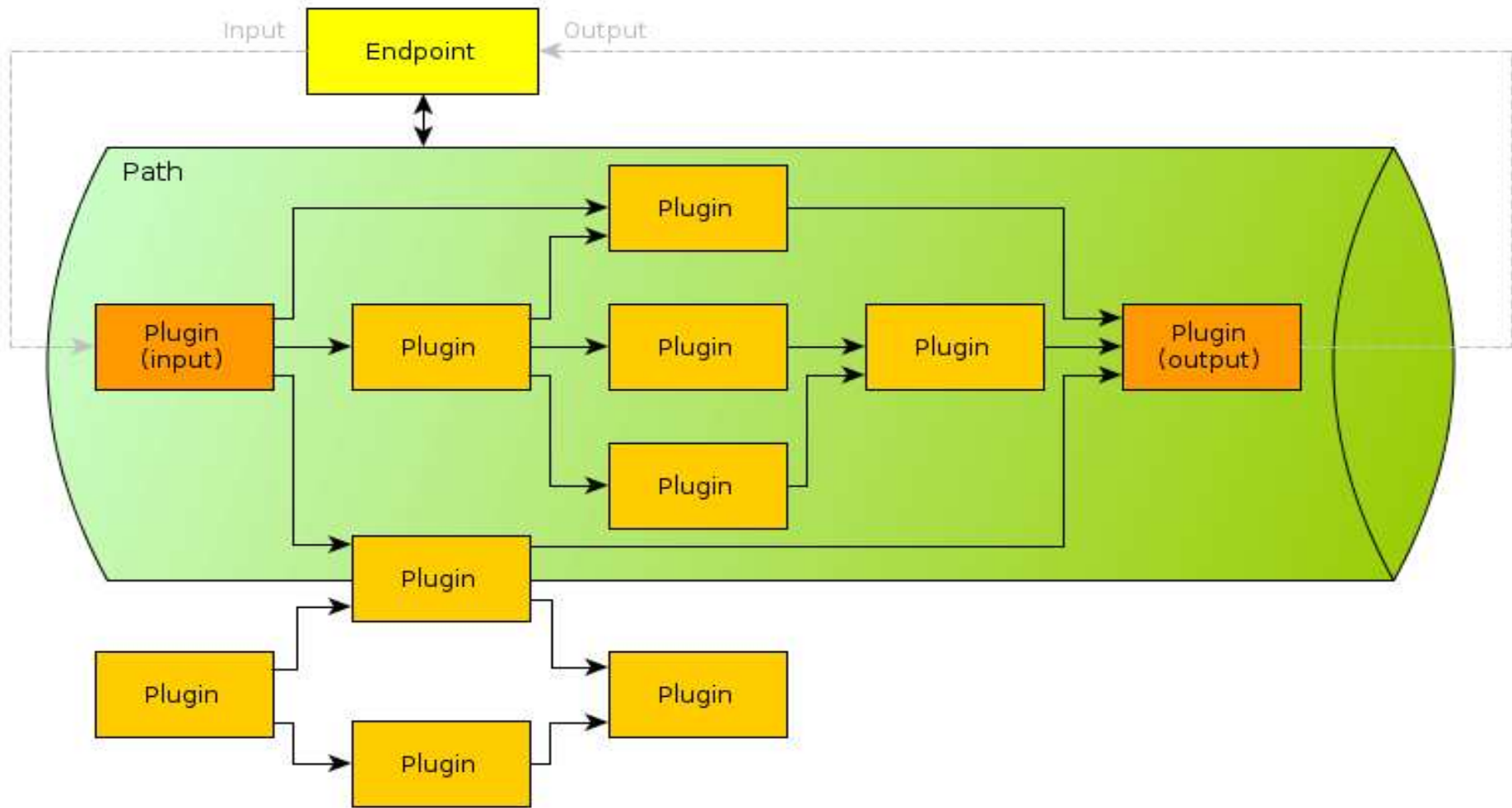
- Improved version controlling and dependency management
- Simplified procedure of new plug-in creation

Large Knowledge Collider



③ The architecture of LarKC

Workflows



③ *The architecture of LarkC*

Workflow Ontology

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix larkc: <http://larkc.eu/schema#> .

# Define two plug-ins
_:plugin1 a <urn:eu.larkc.plugin.identify.TestIdentifier> .
_:plugin2 a <urn:eu.larkc.plugin.transform.TestTransformer> .

# Connect the plug-ins
_:plugin1 larkc:connectsTo _:plugin2 .

# Define a path to set the input and output of the workflow
_:path a larkc:Path .
_:path larkc:hasInput _:plugin1 .
_:path larkc:hasOutput _:plugin2 .

# Connect an endpoint to the path
_:ep a <urn:eu.larkc.endpoint.sparql> .
_:ep larkc:links _:path .
```



③ *The architecture of LarKC*

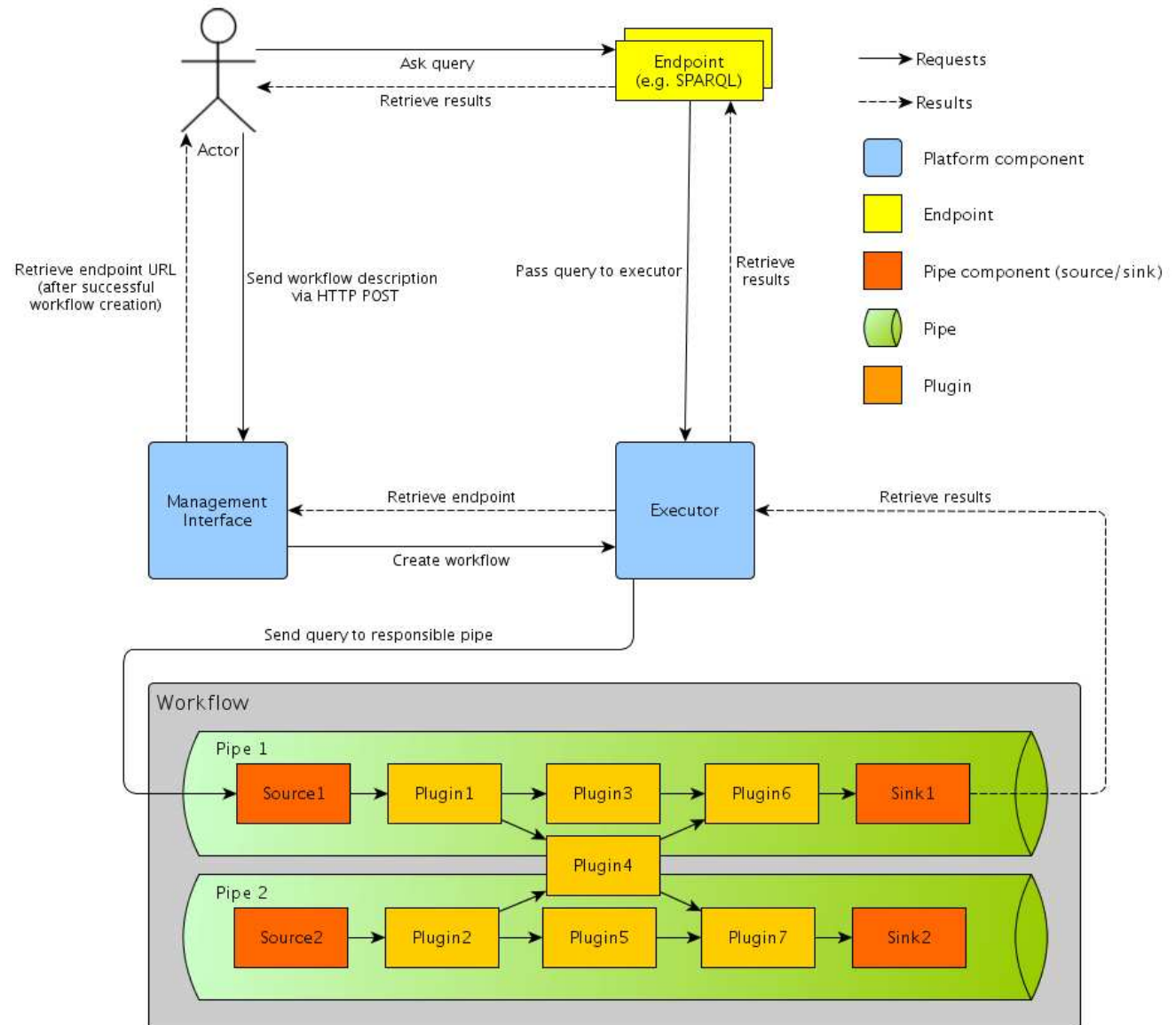
Management Interface

- implemented as a RESTful service
- allows the submission/monitoring/deletion of workflows
- supports N3 and RDF/XML workflow representations
- equipped with additional features for
 - retrieving all registered plug-ins from the platform
 - retrieving configuration templates of all supported remote hosts
- accessible via a simple HTML interface

③ The architecture of LarKC

Run-Time Environment

- comprises the components ()
- responsible for workflows
- collaborates with Managers
- takes care of loops, workflow





③ *The architecture of LarKC*

Data Layer

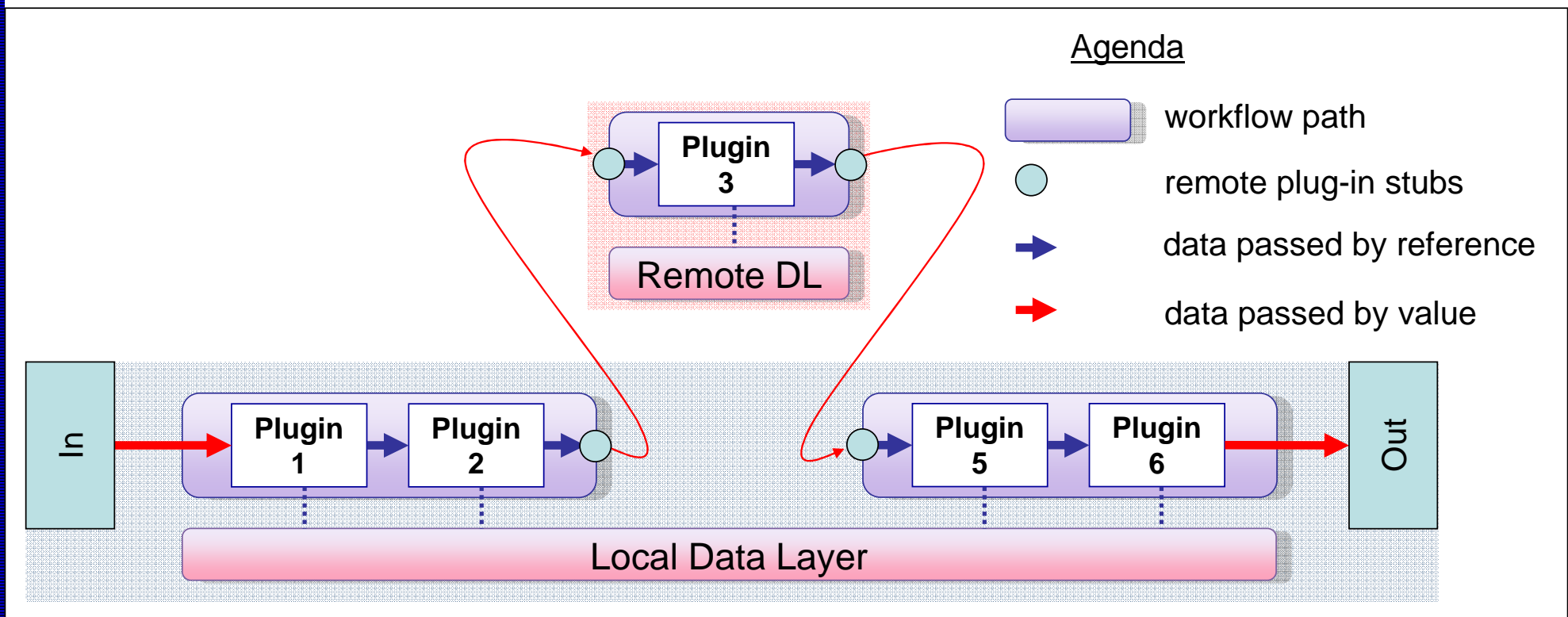
- The Data Layer supports LarKC plug-ins:
 - storage, retrieval and light-weight inference on top of large volumes of data
 - Reference implementation of ORDI data model
 - Retrieval data exposed via standard SPARQL endpoints
 - Configurable forward-chaining reasoning OWL2-RL
 - automates the exchange of RDF data by reference and by value
 - offers other utility tools to manage data (e.g. data streaming, querying remote data)
 - can run in cluster mode to improve resilience and provide scalable query answering

Large Knowledge Collider



③ The architecture of LarKC

Data Layer

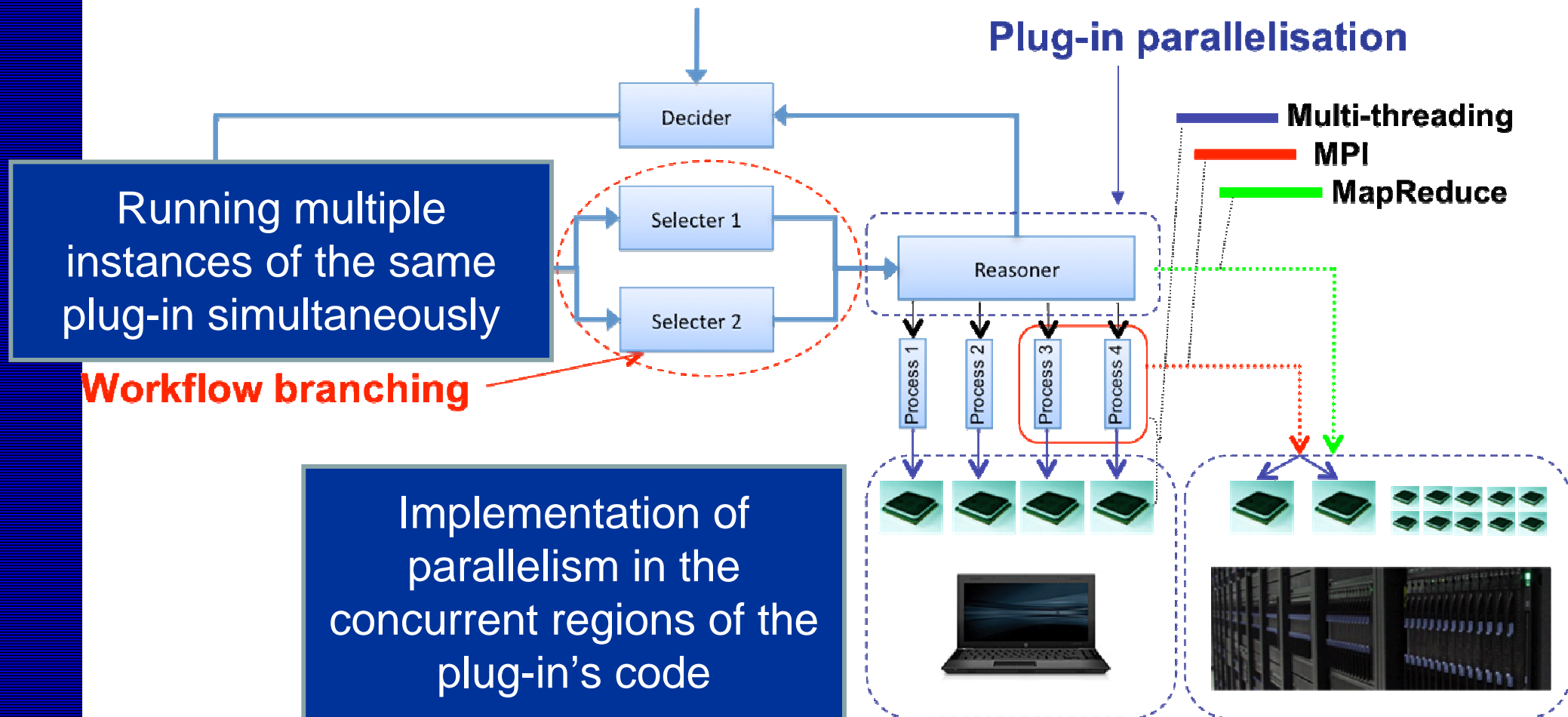


Large Knowledge Collider



③ The architecture of LarKC

Parallelisation



Large Knowledge Collider

③ The architecture of LarKC



User-level Tools

LarKC WorkFlow Designer [New] [Load] [Initialize] [Show Description] [Execute Query] [About]

LarKC Components

- Plugins
 - TestRemoteIdentifier
 - TestIdentifier
 - TestTransformer
 - TestDecider
 - RISearchPlugin
 - SpanningWorkflowPlugin
 - LLDSelector
 - QueryExpansionPlugin
 - SparqlQueryEvaluationReasoner
- Path
 - Input
 - Output
- Hosts
 - GlobusToolkit4
 - Local
 - Tomcat
 - SshWithUserName
 - SshWithPublicKey

Workflow Diagram:

- Input** (diamond) connects to **Local** (rectangle) and **RISearchPlugin** (rectangle).
- RISearchPlugin** connects to **QueryExpansionPlugin** (rectangle).
- QueryExpansionPlugin** connects to **LLDSelector** (rectangle).
- LLDSelector** connects to **SparqlQueryEvaluationReasoner** (rectangle).
- SparqlQueryEvaluationReasoner** connects to **Output** (diamond).
- Local** connects to **Tomcat** (rectangle).
- Tomcat** connects to **SparqlQueryEvaluationReasoner**.

Local Host Configuration:

- larkc:hostType
- larkc:JEE
- larkc:jeeUri
- <http://localhost:8080>

Query Editor:

Query: `SELECT ?s ?p ?o WHERE { { ?s ?p ?o . ?s ?p "asthm`

Endpoints: Endpoint Sparql

Minimap

Result: download RDF

```
<xmp><?xml version="1.0" encoding="U
<sparql xmlns="http://www.w3.org/200
<head>
  <variable xmlns:ns0="http://www.v
  <variable xmlns:ns0="http://www.v
  <variable xmlns:ns0="http://www.v
</head>
<results>
  <result>
    <binding xmlns:ns0="http://www
      <uri>http://linkedlifedata.co
    </binding>
    <binding xmlns:ns0="http://www
      <uri>http://www.w3.org/2004/0
    </binding>
    <binding xmlns:ns0="http://www
      <literal>asthma</literal>
    </binding>
  </result>
</results>
<result>
  <binding xmlns:ns0="http://www
    <uri>http://rdf.freebase.com,
  </binding>
  <binding xmlns:ns0="http://www
    <uri>http://rdf.freebase.com,
  </binding>
```



④ *The LarkC redistributions*

Release 2.5

Plug-in / workflow descriptions and plug-in parameter are in RDF

Separation of workflow specification and execution

Integration of various endpoints (e.g. SPARQL) and applications

Workflow branching, conditional loops, splits / merges of data flow

(Remote) plug-in execution, parallelisation support, anytime behaviour

Data caching, instrumentation and event processing

Data storage, data streaming, parallel request handling



④ The LarKC redistributions

LarKC@SourceForge

- Software releases
 - Platform binaries
 - User and developer guide
- Source code
 - **Maven based build system**
 - Subversion repository
- Additional Support Tools
 - Mailing lists
 - Discussion forums
 - Tracker










































































A screenshot of the SourceForge project page for "Large Knowledge Collider (LarKC) Alpha". The page includes a search bar at the top with the text "Find Open Source Software". Below the search bar, the project name "Large Knowledge Collider (LarKC) Alpha" is displayed, along with the names of the developers: "axeltenschert, bradeskojest, czink, hpcassel, hpcochep". A navigation menu contains links for "Summary", "Files", "Reviews", "Support", "Development", "Hosted Apps", "Tracker", "Mailing Lists", and "Forums". The main content area shows a description of the project as a "platform for..." and "ability barriers of currently existing...". A prominent orange starburst graphic with the text "APACHE license" is overlaid on the page. To the right of the starburst, there is a green "Download" button with the text "larkc-release-2...zip". Below the starburst, there are links for "Other Versions", "Browse all files", "License", "More Detail", and "Show". The "License" link points to "Apache License V2.0".

4 The LarkC redistributions

Plug-In Market Place

Plug-In Marketplace

Plug-ins

Name 	Platform Version	Type	Description 	www 	Download 	Contact person 
Base-Line Full-Text Search Selector	1.0	select				
CRION reasoner	Alpha release	reason				
DIGReasoner	1.0	reason				
EventIdentifier	1.0, 1.1	identify				
GWADecider	1.1	decide				
GWASIdentifier	1.1	identify				
GWASQueryTransformer	1.1	transform				
Information Retrieval Selector	1.0	select				
Interest-Based Reasoner	Alpha Release	reason				
Interest-Based Selector	1.0	select				
Keyphrase Selector	1.0	select				
OWLAPI Reasoner	1.0	reason				
PION Reasoner	1.0	reason				
RDF2MatrixTransformer	1.0	transform				
Random Indexing Reasoner	1.0	reason				
RandomIndexingDecider	1.1	decide				
RandomIndexingIdentifier	1.1	identify				

> 25 Plug-ins available



⑤ Conclusions



is ...

... A platform for rapid semantic web application prototyping

based on light-weight flexible components (plug-ins)

easy integration in a workflow

powered by platform's features (data layer, remote execution, etc.)

... An infrastructure that enables web scale and high performance

... A friendly developer team looking forward to your requests