

# A Software Framework for Procedural Knowledge based Collaborative Data Analytics for IoT

Snehasis Banerjee  
TCS Research & Innovation  
Tata Consultancy Services  
Kolkata, India  
snehasis.banerjee@tcs.com

Mariswamy Girish Chandra  
TCS Research & Innovation  
Tata Consultancy Services  
Bangalore, India  
m.gchandra@tcs.com

**Abstract**—The outburst of data generation by machines and humans, along with emergence of sophisticated data processing algorithms have created a demand for a wide number of data analytics based services and applications. The paper presents a collaborative framework and system to carry out a large number of data processing tasks based on semantic web technology and a combination of reasoning and data analysis approaches using software engineering guidelines. The paper serves as a first step for systematic fusion of symbolic and procedural reasoning that is programming language agnostic. This approach helps in reducing development time and increases developer's productivity. The proposed software system's logical functionality is explained with the help of a healthcare case study, and the same can be extended for other applications.

**Index Terms**—Procedural Reasoning, IoT Analytics, Software Framework, Software Orchestration

## I. INTRODUCTION

Internet-of-Things (IoT) [1] and the services and applications built around it is poised for a groundbreaking growth in near future mainly because of huge number of sensor deployments across the world and advances in networking technologies. In fact, Gartner predicts 21 Billion IoT Devices to be deployed by 2020 [2]. As data is considered the new 'oil' (fuel) for innovations in Industry 4.0 [3], it has been found that data analytics on IoT is in huge demand across all business domains and services around it [4].

Among the various sensors that are deployed (apart from camera and allied sensors) most produce one-dimensional (1-D) readings (say temperature sensor) or multiples of 1-D readings in combination (say accelerometer sensor). Doing analytics on such type of sensor data can be mapped to modules of algorithms belonging to signal processing and machine learning fields. This can be abstractly formalized in a flow diagram for IoT Analytics (Fig. 1).

A typical IoT workflow will comprise (a) Data retrieval from sensors (b) Pre-processing the data by applying various methods like formatting, noise cleaning and anomaly removal (c) Applying data transformation to another representation form like Fourier transform or Wavelet transform (d) Applying Feature Engineering principles of feature extraction from a region of interest followed by selecting features by application of a set of feature selection algorithms (e) finding a suitable Model for the data and problem post tuning modeling algo-

rithm parameters (f) optional inferencing step to derive high level deductions from discovered model (g) Visualizing the results (preferably in a way domain experts can understand).

All of the aforementioned steps need a repository of algorithms for each module and invocation as per need. The traditional approaches for IoT Analytics have confined themselves to quantitative analysis of the data based on the problem at hand. A standing issue is that different suites of algorithms are available in different programming language platforms and this leads to manual intervention or asynchronous and sequential analysis of an IoT Analytics task. The work presented here tries to address this problem. Some work [5] [6] [7] [8] has been done on realizing a data analysis framework for IoT, however they mostly lack the semantic angle and efforts to create a generic framework where problems of different domains can be plugged. In the IoT landscape, employment of semantic techniques has been surveyed in [9], however, only a few IoT applications are found to utilize semantic technologies or in general symbolic techniques of analysis. In this work, an attempt has been made to combine the symbolic and quantitative techniques suitable for a typical IoT Analytics workflow execution. To make this happen, the proposed system has both a Quantitative Analytics Engine and a Symbolic Reasoning Engine that works in synergy by an implicit central controller mechanism. Another burden in a typical IoT Analytics problem is the variety of domain dependent solutions, which is overcome by separating the logic from generic procedures and storing them in rules, ontologies and knowledge stores. This lays foundation for a domain-independent general purpose reasoning based software framework for IoT data analytics. A major issue in a Reasoning enabled Analytics system is lack of support of procedural evaluations and keeping check of facts that have turned obsolete or false. This is overcome by (a) sticking to non-monotonous reasoning following the philosophy of truth maintenance systems [10] and (b) allowing seamless integration of procedural execution into the heart of reasoning following software engineering principles.

The main contributions of this work are:

- (1) a well designed collaborative software framework and system is presented that leverages semantic techniques
- (2) enabling seamless procedural evaluation embedded in semantic rules targeted for IoT data analysis tasks, that leverages



Fig. 1. A typical IoT Analytics Workflow most suited for one-dimensional sensor data

diverse algorithm repositories and code bases (3) a case study in health-care that illustrates the utility of this approach.

Section 2 discusses about the background for IoT Analytics work-flow from a knowledge based perspective. Section 3 describes the proposed system framework and section 4 shows the solution approach taken. Section 5 illustrates procedural reasoning approach of the system and its implications with a healthcare case study. Finally section 6 concludes the paper after laying down the future scope of work.

## II. BACKGROUND

Analysis of data generated from IoT devices has gained considerable attention in the research community. Some work has happened on data analytics on IoT [11] [12] [13], however they do not provide the developer's practical perspective of code re-usability. Also semantics has been mostly ignored with full focus on traditional data science based approaches extended for IoT. This paper explores a first step of systematic application of semantic techniques for IoT Analytics. The semantic web technologies that can be used in developing IoT Analytics applications are enlisted here:

a) RDF<sup>1</sup>: A RDF triple, also called a fact, contains three components: (a) subject, in form of a RDF URI (can be a blank node) (b) predicate, which is a RDF URI representing some property (c) object, which is RDF URI reference that can be a regular reference or a blank node or literal.

An example in triple (`<subject><predicate><object>`) form: `<sensor:TemperatureSensor><rdf:property><temp:Celsius>`

b) Ontologies: used for defining class relationships and set theoretic constraints on relationships. This is typically manifested in OWL<sup>2</sup>, a knowledge representation language in machine interpretable form. OWL enables reasoning from RDF enabled sources. Some of the well known ontologies in IoT space are OGC<sup>3</sup>, SensorML of SWE<sup>4</sup> and SSN<sup>5</sup>; eg. :

```
<sml:input name="SurfaceTemperature"> <sml:
ObservableProperty definition="http://sensors.ws/ont/
NMMO/sensor/SeaSurfaceTemperature"/ > </sml:input>
```

c) SPARQL<sup>6</sup>: used to query RDF models; eg. :  
 SELECT avg(?value) WHERE {  
 <sensor:TemperatureSensor> prop:hasValueInCelsius ?value}

d) Semantic Rules: Rule based reasoning is ideally based on description logic (DL) or first-order predicate logic to infer conclusions from a sequence of statements (premises) derived by predefined rules. A reasoner deduces facts from existing semantic data and ontologies based on predefined

rules. Popular inference engines such as Pellet<sup>7</sup> and Jena<sup>8</sup> are based on different rule languages and provide standard ontology support (such as OWL). Some reasoning engines support SWRL<sup>9</sup> and RIF<sup>10</sup> rule languages, whereas others have instrumented their preferred rule syntax. For reasoning on IoT Analytics workflow, Jena was found suitable for its support of custom rules and ease of implementation of extensions. A typical Jena Rule looks like:

```
[transitiveRule: (?P rdf:type set:TransProp)
(?A ?P ?B) (?B ?P ?C) - > (?A ?P ?C) ]
```

e) Semantic Graph Databases: used when large amounts of facts are to be stored and processed either in batch or on-demand. These are also called RDF Databases and popular ones include Virtuoso<sup>11</sup>, Jena TDB<sup>12</sup> and AllegroGraph<sup>13</sup>.

### A. Reasoning

AI Reasoning depends on generating unambiguous rules from a set of well-formed statements [14]. One way to classify reasoning approaches is in terms of monotonicity. In case of Non-monotonic reasoning (such as truth maintenance system), facts can become false in future, which directly relates to practical scenarios. This requires monitoring of facts so that when a fact becomes invalid because of a context change, facts derived as a consequence of that fact is removed and so on, like a cascading effect. Reasoning techniques involve working with predicates and facts. A predicate on a domain, D, maps every element of D to either true or false T,F, which can be viewed as a relation on  $D \times \{T,F\}$ . A function on a domain, S, maps every element of S to a unique element of the co-domain (or range), R; so it can be viewed as a relation on  $S \times R$ . Any function,  $f : S \rightarrow R$ , can be converted to a predicate,  $p : S \times R$  by defining  $p(s,r) \Leftrightarrow f(s) = r$ . Similarly, any predicate  $p:S$  can be converted to a function  $f : S \rightarrow \{T,F\}$  by defining  $f(s) = T \Leftrightarrow p(s)$  and  $f(s) = F \Leftrightarrow \neg p(s)$ . Functional predicates are useful when in a reasoning enabled system, procedural evaluation is needed to evaluate inferences based on dynamic facts.

### B. Knowledge based Data Analytics

Knowledge has been found to play a crucial role in data analytics. Most of the time, the solutions are dependent on knowledge of domain, algorithms and application. The data analysts use these varied knowledge to filter out non-consistent models as well as apply correct steps in deriving inferences from data at hand. A survey [15] on Semantic Web and Data

<sup>1</sup>Resource Description Framework: <https://www.w3.org/RDF/>

<sup>2</sup>Web Ontology Language: <https://www.w3.org/OWL/>

<sup>3</sup>Open Geospatial Consortium: [www.opengeospatial.org](http://www.opengeospatial.org)

<sup>4</sup>Sensor Web Enablement: [www.sensorml.com/standards.html](http://www.sensorml.com/standards.html)

<sup>5</sup>Semantic Sensor Networks: [w3.org/2005/Incubator/ssn/ssnx/ssn](http://w3.org/2005/Incubator/ssn/ssnx/ssn)

<sup>6</sup>SPARQL Protocol and RDF Query Language: [www.w3.org/sparql](http://www.w3.org/sparql)

<sup>7</sup>Pellet Reasoner: <https://www.w3.org/2001/sw/wiki/Pellet>

<sup>8</sup>Jena under Apache License: <https://jena.apache.org>

<sup>9</sup>Semantic Web Rule Language: [www.w3.org/Submission/SWRL](http://www.w3.org/Submission/SWRL)

<sup>10</sup>Rule Interchange Format: [www.w3.org/TR/rif-overview](http://www.w3.org/TR/rif-overview)

<sup>11</sup>OpenLink Virtuoso: <https://virtuoso.openlinksw.com/>

<sup>12</sup>TDB: <https://jena.apache.org/documentation/tdb>

<sup>13</sup>Franz AllegroGraph: <https://allegrograph.com>

Analytics convey with examples how semantic web technologies can become an integral part of a typical data analysis workflow. It specifically mentions three types of ontologies namely (a) Domain: that express background knowledge about the application domain, (b) Ontologies for data analytics process: that define knowledge about the data analysis process, its steps and algorithms and their possible parameters, (c) Metadata ontologies: that describe meta knowledge about the data, such as provenance, format and source.

Some recent works on knowledge based data analytics are: (a) [16] Publish-Subscribe based Knowledge Analytics that uses automated schema mapping to overcome data heterogeneity and uses semantic inferences on ontology enhanced data, (b) [17] shows how Knowledge graphs can aid in knowledge discovery for data mining, (c) [18] shows how formal semantics in ontologies can be incorporated into data analytics as the formal structure of ontology makes it a natural way to encode domain knowledge, (d) [19] describes an ontology-based system that deals with instance-level integration and data pre-processing. It uses a pre-processing ontology to store information about the required transformations. (e) [20] uses a Knowledge base comprising domain ontology and stream data store. The main aim of the work is to handle real time streaming data using Advanced Message Queuing Protocol and online semantic annotation specifically for smart city scenario.

The value of the proposed system lies in overcoming the main drawbacks in the prior art as listed here: (a) Exploiting procedural execution and symbolic reasoning to create a synergistic system with superior logical expressiveness (b) Understanding the presence of diverse algorithms (procedures) for analytics that needs to be brought together in a single platform to finally give a seamless integrated user experience that supports collaboration (c) Posing IoT Analytics as a workflow recommendation problem by orchestrating diverse procedures based on logical rules and learning in order to develop an application executable guided by knowledge.

### III. SYSTEM DESCRIPTION

The proposed system can be represented as a tuple  $\{ D, M, W, O, R, Q, K, A, L \} \rightarrow I$ , where  $D$  = data,  $M$  = metadata,  $W$  = workflow templates,  $O$  = ontologies,  $R$  = rules,  $Q$  = queries,  $K$  = other knowledge sources,  $A$  = algorithm repository of registered procedural evaluation,  $L$  = learning algorithm recommendation knowledge and  $I$  = inferences drawn post application of reasoning and learning.

The major modules of the system (Fig. 2) are as follows:

1) Data Handler: the purpose of this module is to (a) handle bursts of data when running in live mode by using priority queues (b) using batch to process large data files when in offline mode. The module checks any inconsistency in data (such as mis-matched rows and columns, data format support) and raises a flag to the user accordingly.

2) Workflow Manager: this dedicated module comprises of a workflow engine whose main task is to compose and execute data analytics workflows. This is achieved by composing an apt workflow from templates based on given data, meta-data

and knowledge around it. The workflow manager controls the workflow instance in execution and can enable to and fro data and control exchange between Analytics Engine module and Reasoning Engine module. Further details of this module are kept out of scope to focus on the knowledge aspects.

3) Metadata store: this module contains both specific and generic metadata around algorithms used for analytics. Also meta-data related to given dataset is used as a knowledge input at different stages of analysis.

4) Analytics Engine: this module comprises of the list of analytics algorithm instances (such as those drawn from machine learning and statistical analysis) and their mappings to specific data analytics task. The workflow in execution links with this module when any data analytics task (such as classification) or sub task (such as data transformation) needs to be carried out.

5) Knowledge Store: this module consists of knowledge that is related to datasets, application, domain, workflow templates and algorithms. The knowledge can be in form of (a) ontologies such as OWL (b) facts such as RDF (c) formatted schematic files such as XML (d) flat files (such as csv) with tags or mappings for machine understanding of content (e) RDF stores (f) external databases with defined schema (for relational) or entities (for graph based).

6) Data Unifier: this module is responsible for unification of different data and knowledge formats into a consistent form, which is required for reasoning. Based on standard templates and keywords, mappings are carried out to represent data and knowledge into standard semantic form like RDF and OWL. One important thing in this respect is that the raw data is not converted into RDF form, only the meta-data around the data such as inherent columns and externally supplied metadata is. This saves overhead of re-converting the RDF facts to the initial numeric or categorical values. Each numeric value is tagged as a literal (such as string, double, etc.) in RDF, that adds individual literal tags to each instance of data value. RDF schemas are descriptive such as vocabulary definition, not prescriptive like relational schemas where data must take the prescribed form. This will create a large RDF data file - hence this approach is avoided. The pointers to data sources are kept as RDF facts in the current working memory. However, aggregation of data values that may lead to meaningful inferences can be stored as RDF facts such as mean, median, standard deviation of data fields.

7) Ontology Repository: this module maintains a list of references to all ontology type knowledge files with description. Users are free to add their own custom ontologies or link to an ontology on the web.

8) Rule Repository: this module contains all rules registered in the system via Rule Registry, including procedural rules. Rules have names and are mapped to general usage or specific applications so that selective invocation can be done on demand. System has a large repository of in-built rules which include set theoretic operations as well as commonly used logical expressions. Rules need to be binded to a Reasoner instance for execution. There is a provision to add custom

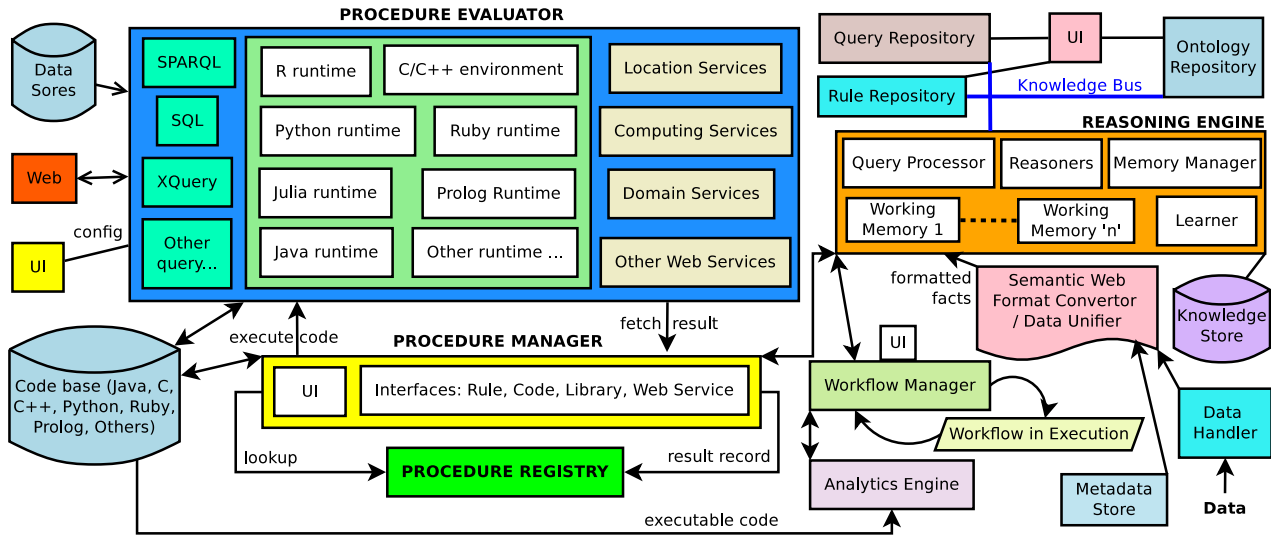


Fig. 2. Proposed Collaborative Framework for Data Analytics based on Hybrid Knowledge Processing

rules including those having procedural functions.

9) Query Repository: this module contains queries which are used in order to view the inferences drawn via rule firing. The repository is pre-filled with simple template SPARQL like semantic queries which can be extended by user as per application demand. Some logging queries are pre-registered that track execution of the Reasoning engine.

10) Reasoning Engine: this important module is concerned with three principal tasks: a) supporting procedural evaluation when firing rules b) provision of apt workflow recommendation based on inferences drawn on given problem and dataset c) algorithm selection for each of the steps of a workflow in execution. The 'Knowledge Bus' connects the Reasoning module to the different knowledge repositories discussed above. A Listener service is registered in the Knowledge Bus that checks if any notable change happens in any of the repositories. Knowledge bus can be simply thought of as a buffer where new knowledge tuples are put by a module and consumed by another module. A provision is kept to store inferred knowledge and facts in apt repositories for future use.

Reasoning engine contains the following sub-modules:

(a) Working Memories: all the relevant facts and ontologies are loaded in a dedicated memory space which is called working memory in computer reasoning terminology. It usually has a limited capacity that is responsible for temporarily holding information available for processing. With advancements in distributed computing and memory technology, it is possible to hold large amounts of facts in memory for fast processing. Facts (or tuples) lying in a working memory under the non-monotonous reasoning paradigm are considered as true, and obsolete facts are deleted from it. In other words, working memory represent the current logical state of a system. The patterns of facts loaded in a working memory is based on the patterns existing in left hand side of rules (productions) and query expressions for a specified workflow execution.

(b) Query Processor: this module supports two types of queries: 1) Ad-hoc: this can be issued by users on demand to query the working memory for situations like result visualization 2) Registered queries: this are pattern seeking queries already linked to the working memory and trigger at regular intervals or on some condition satisfaction.

(c) Reasoners: this module consists of a library of 1) pre-defined standard reasoners such as RDFS or OWL reasoners, where logical inferences are fixed 2) custom reasoners made by binding custom rules including procedural rules to embed logic as per need. A choice of hybrid, forward and backward chaining are available to the user to choose. If the number of rules for a problem is large, then forward chaining is better as rules get triggered only when matching patterns in data are observed. As forward chaining is data driven, rule processing continues until no more rule firing can occur or some limit condition in cycle is met. As this can generate a huge number of inferences (new RDF facts) that can bloat working memory, forward chaining is advised for usage in event driven scenarios. Backward chaining on the contrast is a goal driven inference technique. It looks for needed data (facts) patterns by beginning from possible goals or conclusions. This makes backward chaining suitable when number of rules are few and when working memory has limitations in capacity. For small datasets and rules, the earlier two approaches or a hybrid one usually suffice, but for large amounts of data and large number of rules, Rete is the optimal choice of reasoning. The popular and fast Rete algorithm builds upon a rule engine based on pattern-matching production system with main aim of matching facts against rules in an optimal way. Rete circumvents repetitive searching through the data elements by keeping the current contents of the conflict set in memory. Addition and deletion of conflict set items are triggered only when data elements are updated (addition/deletion) in memory.

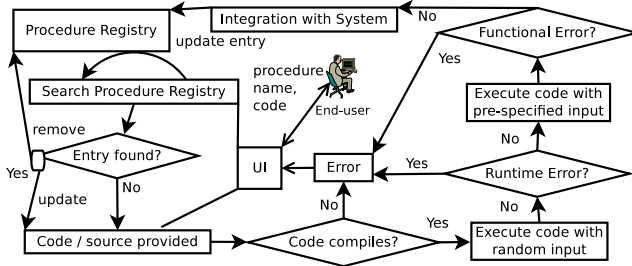


Fig. 3. Flowchart of Procedure Update Validation

Rete algorithm is designed to sacrifice memory (by creating and maintaining Rete Network of Alpha and Beta Nodes) for increased speed. Rete has three components: 1) Alpha Network - left side of the node graph (rule patterns) that forms a discrimination network and selects elements from working memory based on simple conditions 2) Beta Network - the right side of node graph which mainly performs joins between elements in working memory, 3) Agenda - typically implemented as prioritized queues, that finds the order in which a rule is fired in a match-resolve-act cycle.

(d) Memory Manager: this module retrieves resources from various repositories via Knowledge Bus. It manages the working memories and allocates and deallocates memory areas to individual reasoner instances based on need. What to load in the working memory is defined based on the dataset, application and problem description set by the user, apart from support of direct manual loading from the user's end. The system dedicates a working memory for each application instance. Due to the nature of knowledge and applications, it is often found that there is a significant overlap of the knowledge needed across use cases. Hence, keeping redundant facts across disjoint working memories adds to overhead in processing as well as space. This can be overcome by having multiple working memories with sharing support [21].

(e) Learning Module (Learner): This module's learning is based on execution logs of the reasoner. System learns which facts and rule bindings need to be kept in logical cache memory based on frequency of procedures getting called. The execution times and any error or crashes are noted so that the user who provided the procedure can rectify the problem. This also makes other users aware of probable code runtime failures and in a way influences their selection of procedures from the collaborative procedure selection list.

11) Procedure Manager: this module is the central control to manage procedural evaluations while reasoner is in execution. When a new procedural rule needs to be registered, the rule binding happens in the reasoner (implicit reasoner registry) while the procedure mapping happens in Procedure Registry. Procedure Registry maintains list of algorithm codes with pointer, names, arguments, parameters and meta-data. Users can define their own procedures to extend the built in functions of the reasoning vocabulary of keywords and pointers to executable methods. Fig 4. illustrates the user's perspective on

semantic rule formation as per demand of the solution logic. The registry contains following types of functions (methods or procedures):

a) Standard functions: popular and often used methods like mean or average, b) Defined functions: custom methods that are registered by the user in the system, c) Shared functions: functions included by other users and allow sharing for re-usability, d) Web Service functions: ready to use wrappers on existing web service APIs integrated in the system, e) External functions: wrappers on database connections and external tools bound to system registry.

12) Code Base: the executable files are kept here. Support for various languages such as C / C++, Ruby, Python, Java, Julia, Prolog, R is kept as a) different languages has different functionality that can be leveraged, b) users usually have their own preference of languages, c) different library of functions are available for each language and a specific required code library may be available only in one language.

13) Procedure Evaluator: this is composed of runtime environments of different programming languages, where the values are passed and code gets executed. In the instantiation of the framework, Java Virtual Machine (JVM) was selected as the central control due its ease of portability. JVM is connected with various runtimes with appropriate wrappers and encapsulating typed value exchange. As an illustration for seamless integration, Python runtime was connected via Jython<sup>14</sup>, Ruby was fused using JRuby<sup>15</sup>, Prolog was joined by employing JiProlog<sup>16</sup>, C/C++ was coupled through Java Native Interface (JNI<sup>17</sup>) and R was binded via Java/R Interface (JRI<sup>18</sup>). Connectivity to external databases (such as those listed in LOADB<sup>19</sup> and DBPedia<sup>20</sup>) supporting queries via SPARQL, SQL and other querying methods are maintained to link aptly at the time of procedural evaluation. Various services and their APIs are kept registered to use the functionalities offered. Services include Location (like Google Distance Matrix API), Computing (like evaluation via Wolfram Alpha) and Domain (like information around disease outbreaks). Inter connections between run-times are kept to minimize context switching.

The system has multiple front end user interfaces (UIs) binded to specific modules that serve as means to update files, configurations and defining relations and placeholders in the system. Due to limited space, discussion on UI is skipped.

#### IV. METHOD DESCRIPTION

The proposed system's underlying principle has been chosen to be non-monotonic forward chaining deductive reasoning supporting functional predicates to suit needs of IoT workflows. Logic can be expressed both in form of semantic rules and SPARQL like graph pattern queries. However, to

<sup>14</sup><https://www.jython.org>

<sup>15</sup><https://www.jruby.org/>

<sup>16</sup><http://www.jiprolog.com/>

<sup>17</sup><https://docs.oracle.com/javase/8/docs/technotes/guides/jni/>

<sup>18</sup><https://www.rforge.net/JRI/>

<sup>19</sup><http://www.loadb.org>

<sup>20</sup><https://wiki.dbpedia.org/>

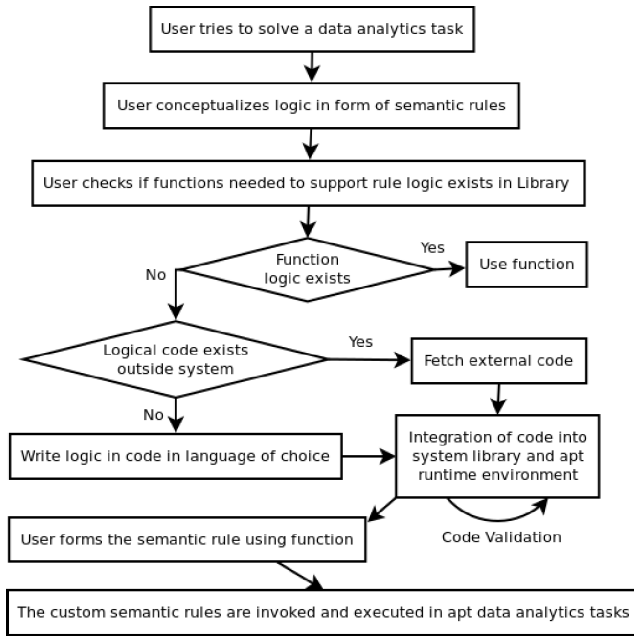


Fig. 4. Flowchart of custom semantic rules formation from user perspective

maintain the truth maintenance aspect, logic in this case is restricted to semantic rules having provision for fact removal. SPARQL is used as a way to view and retrieve answers to queries (patterns under scrutiny). Rete [22] and its improved variations [23] [24] have stood the test of time for handling fast pattern matching in a non-monotonous reasoning scenario. For IoT based workflow execution, knowledge decisions based on Jena's implementation of Rete forward reasoner was found suitable as Rete is able to handle large number of rules and fast fact deletion based on its alpha-beta networks and efficient conflict resolution strategies.

#### A. Handling External Procedures

External procedures in any supported programming language can be added to the system. Fig 3. shows how a user can register a procedure to be binded in Procedure Registry. User starts with searching for a procedure name in the existing registry and if entry is found, the user can view the code if it is public and can update (also remove) the same if access rights are granted. Else, user creates a new entry of procedure by linking or uploading code in language of choice. The system checks (in the runtime environment) if the code gets compiled and executes code with random input to see if any runtime errors or exceptions happen. Next, system asks user to provide a set of input and expected output and sees if the procedure yields functionally correct results. Next, the procedure gets integrated with the Registry, thereby enabling a user to use the function names and arguments in defining custom rules. Any errors in the process are shown to user for rectification.

#### B. Collaborative Development

To the best of our knowledge the proposed system and approach is the first step to enable a collaborative platform for IoT data analytics by harvesting the power of community shared logic in form of reusable code snippets.

The users can do the following in the UI connected to the Procedure Registry: a) Updating / adding / removing custom functions (algorithmic methods) b) Sharing functions with all users or a group or keeping it private c) Viewing source code or other implementation details of shared functions with apt permissions d) Commenting, rating and providing suggestions for different Registry functions e) Tag the self created functions with meaningful meta-data f) In the particular session, users can turn on and off certain functions. This switching comes handy in case of overriding functions like the different interpretations of a function bearing the same name.

Fig 4. illustrates the user's perspective on semantic rule formation as per demand of the solution logic. In order to solve a data analytics task, user thinks of the problem at hand and how to conceptualize a solution about it. Usually the path taken to solve that problem can be thought of semantically connected steps whose logic can be expressed in rules. User next checks in the system where building blocks of rules (combining which a complete rule expression will be formed) exist in the system, that were registered by other developers or packaged by default. If a building block (in this case, a functionality) exists in the rule processing system, then the user composes a rule based on logical requirement of the task. Else, the user is given an option to integrate required procedural logic as building blocks of rules to be added to the system. If the procedure already exists outside the system, user can just fetch the code and let the system integrate it automatically. Else the user can write the code in an editor and integrate it in the system. Code is validated before registering the function with a URI based name in the system. User forms a set of rule expressions to execute the same in the system to serve a specific task on a given data.

Functions or methods stored in the system are analyzed about their rate of usage, popularity, time and space requirements, updates, extension and related particulars. This are stored as logger table per registered function. This aids in creating effective recommendation of functions when a user searches the Registry for existing procedural functions. User created functions can be rated, shared and extended to enable a collaborative development ecosystem. Recommendation can be context-based [25] as well. Functionality of added functions can be tested and validated, with feedback comments given to the initial procedural code provider. There are well known approaches [26] [27] [28] in literature for sharing code between users and collaborative development - similar adoptions can be taken in future. In summary, the system can be thought of as an ever growing repository of algorithms that can be integrated with IoT Analytics system as per need of solution.

System creates a new URI (uniform resource identifier) for the new function being added with its programming language

as a sub part of URI. This helps in distinguishing functions with same names based on domain and nomenclature. This URI assignment is done automatically following standard pre-defined standard templates. An example is shown:- `http://com.algorepo/number/c/2/1/resample` where 'number' is the domain, 'c' represents programming language; '2' shows number of input arguments; '1' means number of declared return values; finally the name of the function as 'resample'. System creates a field and knowledge store entry to check the number of accesses to the calling function (this is updated on each call as hit) to monitor popularity in collaborative platform. Another field is created to monitor average memory and time requirement the given function consumes. A trailing number is added to the function name at the end of generated code that can be used to call the function when writing the rule so that a particular user is not required to type down the entire URI in order to refer to that specific code. This comes as a overlay in the UI when function is used in a rule. Example: *resample2*, when two functions (may be expressed in different languages) already exist. This works in the format: function name + previous counter value + unit increment. Apart from URI representation, a number block based unique ID is also kept for fast reference lookup and processing. This is inspired by memory addressing in computer architecture literature.

## V. PROCEDURAL REASONING

While using Semantic Web technology for reasoning covers the description logic aspects, enabling procedure calls from bounded variables while reasoner is executing makes a reasoner powerful in handling procedural reasoning. Some limited work [29] has happened on extending SPARQL with external execution calls. However, a collaborative extension approach for reasoning was not envisioned. Some example procedural reasoning scenarios along with corresponding logic expressions in rule syntax is shown for a healthcare use case.

The task of 2016 PhysioNet / CinC Challenge<sup>21</sup> was to classify heart sound recordings (phonocardiogram / PCG) into normal and abnormal recordings. Main steps of the task are pre-processing, segmentation and modeling (classification). The following logic snippets show the approach's usefulness for pre-processing step for illustration purpose.

A section of domain knowledge used for the problem:

1. `<problem:heartSound><sig>windowSizeInSeconds> 5`
  2. `<problem:heartSound><sig>upperBoundFreqInHz> 400`
- The above knowledge when expanded states that for a problem type related to heart (cardiac) sound recording, window size recommended is 5 seconds. Implicit knowledge not shown is that doctors have specified that 5 seconds monitoring is good enough to find a significant pattern in the heart sound. Note that window can be of many types [30], however in this case overlapping window is implicitly considered.

A section of meta-data facts supplied for the problem:

1. `<data:dataInstance><sig>samplingRateInHz> 2000`

2. `<data:dataInstance><file:format><file:wav>`

3. `<data:dataInstance><data:hasSource> 'file path'`

The above knowledge describes that the sampling rate for a specific data instance should be 2000 Hz. Also the file path of the data instance and the type ('wav') is mentioned.

A section of algorithmic knowledge stored in ontology:

1. `<sig:MedianFilter><sig:type><sig:NonLinearFilter>`
2. `<sig:LowPassFilter><sig:type><sig:LinearFilter>`
3. `<sig:LinearFilter><sig:type><sig:Filter>`
4. `<task:PreProcessing><task:subTask><sig:Filter>`

The above knowledge states that median filter, a well known signal processing filter is a type of non-linear filter. On the other hand, low pass filter is type of Linear Filter which in turn is a specialization of a Filter. Filtering is a sub task of pre-processing step of IoT Analytics. In this way other IoT Analytics algorithms can be described.

A section of procedural rules (custom defined) used:

*Rule 1.* [ `AutoResample: (?data <data:hasSource> ?url ) .`

`(?data <problem:category> ?p) .`

`(?data <sig:samplingRateInHz> ?val1) .`

`(?data <sig:upperBoundFreqInHz> ?val2) ->`

`c:resampleA(?url, ?val1, ?val2, ?urln) .`

`(?p <task:completion> <sig:AutoResample>) .`

`(?data <data:newSource> ?urln) ]`

*Rule 2.* [ `LowPassFilter: (?data <data:hasSource> ?url ) .`

`(?data <problem:category> ?p) .`

`(?p <sig:upperBoundFreqInHz> ?val) ->`

`c:lowpass(?url, ?val, ?urln) .`

`(?p <task:completion><sig:LowPassFilter>) .`

`(?data <data:newSource> ?urln) ]`

The above procedural rule no. 1 attempts to resample data using a custom algorithm. Resampling data to a lower rate without missing patterns is recommended to enable faster processing. The auto resampling algorithm takes in sampling rate and the upper frequency range with in which meaningful frequencies exist and using principles of Nyquist Theorem and sampling [31], resamples data to an optimal sampling rate. The rule takes the data instance source, checks problem tagging and sees if any facts about sampling rate and upper bound of frequency exists, next calls the autoresampling algorithm (c: is the initial prefix notation to denote custom procedural calls) coded and registered by user and outputs a fact on completion of operation and link to processed data. Rule no. 2 takes a data source, checks the problem tagging and sees if any domain knowledge of cutoff frequency exists, and next does a procedural call to low pass filter algorithm and gives out a link to processed data as well as adding a fact in working memory that lowpass filter operation is done. In this case, an available algorithm is used that is found in standard signal processing libraries of Matlab or Python. For other domains, it may be the case that required libraries for solution exists in a mixture of Java, R and C codes. Appropriate mapping via algorithm registration keeps the implementation aspects hidden from logical rule writer. Hence it is showcased how symbolic and procedural knowledge based processing can be used for data analysis tasks in general.

<sup>21</sup><https://www.physionet.org/challenge/2016/>



## VI. CONCLUSION

The paper has presented a collaborative framework and system to carry out a large number of data processing tasks based on semantic web technology with support of procedural rules that is programming language agnostic. Example procedural rules were explained with the help of a healthcare use case. For healthcare domain some automation work [32] [33] [34] for IoT Analytics has been reported. However, this work introduces the merging of symbolic knowledge and procedural analysis to create a new approach for IoT Analytics. Future work will include extensions to 2-D (image) and 3-D (video) sensor data. Fast and incremental reasoning is a need of the system. Hence, to suit event-based IoT scenarios an improved Rete [24] [23] algorithm in the reasoning module is planned for integration. As IoT scenario deals with streaming realtime data, explorations in integrating stream reasoning [35] with the framework is planned. Finally, the system will be tested for different Analytics pipelines to gather further information on usability of procedural rules. The approach reuses other developer's efforts in semantic rule formation and workflow execution, thereby reducing development time.

## REFERENCES

- [1] P. Sethi and S. R. Sarangi, "Internet of things: architectures, protocols, and applications," *Journal of Electrical and Computer Engineering*, vol. 2017, 2017.
- [2] N. Eddy, "Gartner: 21 billion iot devices to invade by 2020," *InformationWeek*, Nov. 10, 2015.
- [3] H. Lasi, P. Fettke, H.-G. Kemper, T. Feld, and M. Hoffmann, "Industry 4.0," *Business & Information Systems Engineering*, vol. 6, no. 4, pp. 239–242, 2014.
- [4] M. S. Mahdaveinejad, M. Rezvan, M. Barekatin, P. Adibi, P. Barnaghi, and A. P. Sheth, "Machine learning for internet of things data analysis: A survey," *Digital Communications and Networks*, 2017.
- [5] M. Strohbach, H. Ziekow, V. Gazis, and N. Akiva, "Towards a big data analytics framework for iot and smart city applications," in *Modeling and processing for next-generation big-data technologies*. Springer, 2015, pp. 257–282.
- [6] A. Sinharay, A. Pal, S. Banerjee, R. Banerjee, S. Bandyopadhyay, P. Deshpande, and R. Dasgupta, "A novel approach to unify robotics, sensors, and cloud computing through iot for a smarter healthcare solution for routine checks and fighting epidemics," in *International Internet of Things Summit*. Springer, 2015, pp. 536–542.
- [7] K. Hwang and M. Chen, *Big-data analytics for cloud, IoT and cognitive computing*. John Wiley & Sons, 2017.
- [8] D. Pizzolli, G. Cossu, D. Santoro, L. Capra, C. Dupont, D. Charalampos, F. De Pellegrini, F. Antonelli, and S. Cretti, "Cloud4iot: A heterogeneous, distributed and autonomic cloud platform for the iot," in *IEEE CloudCom*. IEEE, 2016, pp. 476–479.
- [9] P. Barnaghi, W. Wang, C. Henson, and K. Taylor, "Semantics for the internet of things: early progress and back to the future," *IJWSIS*, vol. 8, no. 1, pp. 1–21, 2012.
- [10] H. Beck, "Reviewing justification-based truth maintenance systems from a logic programming perspective," Tech. Rep. INFYSYS RR-1843-17-02, Institute of Information Systems, TU Vienna. July, Tech. Rep., 2017.
- [11] A. Zaslavsky, C. Perera, and D. Georgakopoulos, "Sensing as a service and big data," *International Conference on Advances in Cloud Computing (ACC)*, arXiv preprint arXiv:1301.0159, 2013.
- [12] M. Marjani, F. Nasaruddin, A. Gani, A. Karim, I. A. T. Hashem, A. Siddiqi, and I. Yaqoob, "Big iot data analytics: architecture, opportunities, and open research challenges," *IEEE Access*, vol. 5, pp. 5247–5261, 2017.
- [13] E. Ahmed, I. Yaqoob, I. A. T. Hashem, I. Khan, A. I. A. Ahmed, M. Imran, and A. V. Vasilakos, "The role of big data analytics in internet of things," *Computer Networks*, vol. 129, pp. 459–471, 2017.
- [14] R. S. Michalski, "Inferential theory of learning as a conceptual basis for multistrategy learning," *Machine Learning*, vol. 11, no. 2-3, pp. 111–151, 1993.
- [15] P. Ristoski and H. Paulheim, "Semantic web in data mining and knowledge discovery: A comprehensive survey," *Web semantics: science, services and agents on the World Wide Web*, vol. 36, pp. 1–22, 2016.
- [16] C. Esposito, M. Ficco, F. Palmieri, and A. Castiglione, "A knowledge-based platform for big data analytics based on publish/subscribe services and stream processing," *Knowledge-Based Systems*, vol. 79, pp. 3–17, 2015.
- [17] P. Ristoski, "Exploiting semantic web knowledge graphs in data mining," Ph.D. dissertation, University of Mannheim, 2018.
- [18] D. Dou, H. Wang, and H. Liu, "Semantic data mining: A survey of ontology-based approaches," in *Semantic Computing (ICSC)*, 2015. IEEE, 2015, pp. 244–251.
- [19] D. Perez-Rey, A. Anguita, and J. Crespo, "Ontodataclean: Ontology-based integration and preprocessing of distributed data," in *International Symposium on Biological and Medical Data Analysis*. Springer, 2006, pp. 262–272.
- [20] S. Kolozali, M. Bermudez-Edo, D. Puschmann, F. Ganz, and P. Barnaghi, "A knowledge-based approach for real-time iot data stream annotation and processing," in *iThings, GreenCom and CPSCOM*. IEEE, 2014, pp. 215–222.
- [21] S. Banerjee and D. Mukherjee, "Towards a universal notification system," in *Proceedings of the 2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)-Volume 03*. IEEE Computer Society, 2013, pp. 286–287.
- [22] C. L. Forgy, "Rete: A fast algorithm for the many pattern/many object pattern match problem," in *Readings in Artificial Intelligence and Databases*. Elsevier, 1988, pp. 547–559.
- [23] B. Berstel, "Extending the rete algorithm for event management," in *Temporal Representation and Reasoning, 2002. TIME 2002. Proceedings. Ninth International Symposium on*. IEEE, 2002, pp. 49–51.
- [24] T. Gao, X. Qiu, and L. He, "Improved rete algorithm in context reasoning for web of things environments," in *Green Computing and Communications (GreenCom), 2013 IEEE and Internet of Things (iThings/CPSCOM), IEEE International Conference on and IEEE Cyber, Physical and Social Computing*. IEEE, 2013, pp. 1044–1049.
- [25] D. Mukherjee, S. Banerjee, S. Bhattacharya, and P. Misra, "A context-aware recommendation system considering both user preferences and learned behavior," in *2011 7th International Conference on Information Technology in Asia*. IEEE, 2011, pp. 1–7.
- [26] H. Shen, "Interactive notebooks: Sharing the code," *Nature News*, vol. 515, no. 7525, p. 151, 2014.
- [27] S. Haeffiger, G. Von Krogh, and S. Spaeth, "Code reuse in open source software," *Management science*, vol. 54, no. 1, pp. 180–193, 2008.
- [28] E. M. Lucas, T. C. Oliveira, K. Farias, and P. S. Alencar, "Collabrdl: A language to coordinate collaborative reuse," *Journal of Systems and Software*, vol. 131, pp. 505–527, 2017.
- [29] S. Banerjee and D. Mukherjee, "On demand sparql extension: A case study of extending geo-sparql for sensor data exploration in semantic cities," in *Workshops at the Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014.
- [30] S. Banerjee et al., "Windowing mechanisms for web scale stream reasoning," in *Proceedings of the 4th international workshop on Web-scale knowledge representation retrieval and reasoning*. ACM, 2013, pp. 17–18.
- [31] E. J. Candes, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Transactions on information theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [32] S. Banerjee, T. Chattopadhyay, S. Biswas, R. Banerjee, A. D. Choudhury, A. Pal, and U. Garain, "Towards wider learning: Experiments in healthcare," *NIPS Workshop, ML4Health*, arXiv preprint arXiv:1612.05730, 2016.
- [33] S. Banerjee, T. Chattopadhyay, A. Pal, and U. Garain, "Automation of feature engineering for iot analytics," *ACM SIGBED Review*, vol. 15, no. 2, pp. 24–30, 2018.
- [34] S. Banerjee, T. Chattopadhyay, and A. Mukherjee, "Interpretable feature recommendation for signal analytics," *CIKM Workshop, IDM*, arXiv preprint arXiv:1711.01870, 2017.
- [35] D. Mukherjee, S. Banerjee, and P. Misra, "Towards efficient stream reasoning," in *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*. Springer, 2013, pp. 735–738.