

IEEE – Entity Event Knowledge Graph for Powerful Health Informatics

As part of Franz's participation in the IEEE – ICHI conference, our paper has been published and is available from the IEEE Website.



ICH2022 is a premier community forum concerned with the application of computer science, information science, data

science, and informatics principles, as well as information technology, and communication science and technology to address problems and support research in healthcare, medicine, life science, public health, and everyday wellness.

Franz Inc. presented on June 14th – **Entity Event Knowledge Graph for Powerful Health Informatics**

Download Franz's IEEE Publication – Entity Event Knowledge Graph for Powerful Health Informatics.

Conference Website

No-Code Queries Can Accelerate AI and Data Analytics

By Dr. Jans Aasman, CEO

The low-code, no-code methodology is becoming highly sought-after throughout the modern IT ecosystem—and with good reason. Options that minimize manually writing code capitalize on the self-service, automation idiom that's imperative in a world in which working remotely and doing more with less keeps organizations in business.

Most codeless or low-code approaches avoid the need for writing language-specific code and replace it with a visual approach in which users simply manipulate on-screen objects via a drag-and-drop, point-and-click interface to automate code generation. The intuitive ease of this approach – which is responsible for new standards of efficiency and democratization of no-code development – has now extended to no-code query writing.

No-code querying provides two unassailable advantages to the enterprise. First, it considerably expedites what is otherwise a time-consuming ordeal, thereby accelerating data analytics and AI-driven applications and second, it can help organizations overcome the talent shortage of developers and knowledge engineers. Moreover, it does so by furnishing all the above benefits that make codeless and low-code options mandatory for success.

Read the full article at [DZone](#).

Data-Centric Architecture Forum – DCAF 2021

Data and the subsequent knowledge derived from information are the most valuable strategic asset an organization possesses. Despite the abundance of sophisticated technology developments, most organizations don't have disciplines or a plan to enable data-centric principles.

DCAF 2021 will help provide clarity.

Our overarching theme for this conference is to **make it REAL**. Real in the sense that others are becoming data-centric, it is achievable, and you are not alone in your efforts.

Join us in understanding how data as an open, centralized resource outlives any application. Once globally integrated by sharing a common meaning, internal and external data can be readily integrated, unlike the traditional "application-centric" mindset predominantly used in systems development.

The compounding problem is these application systems each have their own completely idiosyncratic data models. The net result is that after a few decades, hundreds or thousands of applications implemented have given origin to a segregated family of disparate data silos. Integration debt rises and unsustainable architectural complexity abounds with every application bought, developed, or rented (SaaS).

Becoming data-centric will improve data characteristics of findability, accessibility, interoperability, and re-usability (FAIR principles), thereby allowing data to be exported into any needed format with virtually free integration.\



**Dr. Jans Aasman to present –
Franz's approach to Entity Event
Data Modeling for Enterprise
Knowledge Fabrics**

Text Analytics Forum 2020 – KMWorld Connect

Join us November 17, 2020 – Text Analytics has the ability to add depth, meaning, and intelligence to any organization's most under-utilized resource – text. Through text analytics, enterprises can unlock a wealth of information that would not otherwise be available. Join us as we explore the power of text analytics to provide relevant, valuable, and actionable data for enterprises of all kinds.

Jans Aasman to present – Analyzing Spoken Conversations for Real-Time Decision Support in Mission-Critical Applications

November 17, 2020 at 2PM Eastern

Sharing Ontologies Globally To Speed Science And Healthcare Solutions – OntoPortal

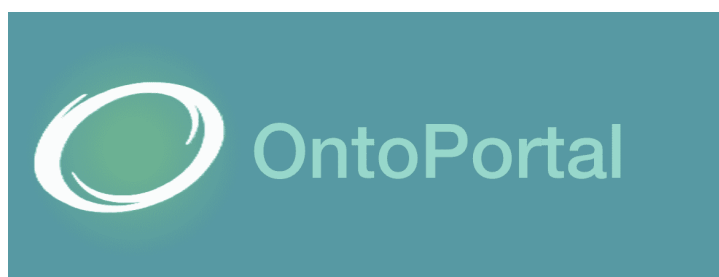
International Ontology Sharing Is Becoming A Reality

A consortium of researchers recently formed an organization dedicated to standardizing how scientists define their ontologies, which are essential for retrieving datasets as well as understanding and reproducing research. The group called OntoPortal Alliance is creating a public repository of internationally shared domain-specific ontologies. All the repositories will be managed with a common OntoPortal appliance that has been tested with AllegroGraph Semantic Knowledge Graph software. This enables any OntoPortal adopter to get all the power, features, maintainability, and support benefits that come from using a widely adopted, state-of-the-art semantic knowledge graph database.

Read the full article at [HealthIT Outcomes](#) –

As Dr. Jans Aasman, CEO of Franz Inc. explains, “When building a Knowledge Graph as your enterprise’s single source of truth, it’s critical to include ontologies and taxonomies. AI applications and complex reasoning analytics require information from both databases and knowledge bases that contain domain information, taxonomies, and ontologies to solve complex questions. To make this possible, we developed a novel hybrid sharding technology called FedShard, which facilitates the combination of data and knowledge required by applications like Montefiore’s PALM. But this approach is not unique or specific to Healthcare, it is applicable in many

other industries, which is why we are excited about OntoPortal's plans to bring sharing of domain ontologies to a broad audience."



**Knowledge Graphs: A Single
Source of Truth for the
Enterprise**



The notion of a “single source of truth” for the enterprise has been the proverbial moving goalpost for generations of CIOs, only to be waylaid by brittle technology and unending legacy systems. Truth-seeking visions rebuffed by technological trends have continuously confounded business

units trying to achieve growth and market penetration. But technology innovation has finally led us to a point where CIOs can now deliver that truth.

Graphing the Truth

Knowledge graphs possess the power to deliver a single source of truth by linking together any assortment of data sources required, standardizing their diversity of data elements, and eliminating silos. They support the most advanced analytics options and decentralized transactions, which is why they’re now deployed as systems of records for some of the most significant, mission-critical use cases affecting our population.

Because they scale to include almost any number of applications – and link to other knowledge graphs as well – these repositories are the ideal solution for real-time information necessary to inform business users’ performances with concrete, data-supported facts. Most importantly, users can get an exhaustive array of touchpoints pertaining to any customer, product, or interaction with an organization from the knowledge graph, making it a single source of truth.

Read the full article at [Dataversity](#).

Using Microsoft Power BI with AllegroGraph

There are multiple methods to integrate AllegroGraph SPARQL results into Microsoft Power BI. In this document we describe two best practices to automate queries and refresh results if you have a production AllegroGraph database with new streaming data:

The first method uses Python scripts to feed Power BI. The second method issues SPARQL queries directly from Power BI using POST requests.

Method 1: Python Script:

Assuming you know Python and have it installed locally, this is definitely the easiest way to incorporate SPARQL results into Power BI. The basic idea of the method is as follows: First, the Python script enables a connection to your desired AllegroGraph repository. Then we utilize AllegroGraph's Python API within our script to run a SPARQL query and return it as a Pandas dataframe. When running this script within Power BI Desktop, the Python scripting service recognizes all unique dataframes created, and allows you to import the dataframe into Power BI as a table, which can then be used to create visualizations.

Requirements:

1. You must have the AllegroGraph Python API installed. If you do not, installation instructions are here: <https://franz.com/agraph/support/documentation/current/python/install.html>
2. Python scripting must be enabled in Power BI Desktop. Instructions to do so are here: <https://docs.microsoft.com/en-us/power-bi/connect-data/desktop-python-scripts>

a) As mentioned in the article, pandas and matplotlib must be installed. This can be done with 'pip install pandas' and 'pip install matplotlib' in your terminal.

The Process:

Once these requirements have been met, create a Python file with whatever script editor you usually use. The following code will create a connection to your desired repository. For this example, we will be using the Kennedy dataset that is available with the AllegroGraph distribution (See the 'Tutorial' directory). Load the Kennedy.ntriples file into your running AllegroGraph. (Replace the '****' in the code with your corresponding username and password.)

#the necessary imports

```
import os

from franz.openrdf.connect import ag_connect

from franz.openrdf.query.query import QueryLanguage

import pandas as pd
```

#connect to your agraph repository

```
def setup_env_var(var_name, value, description):

    os.environ[var_name] = value

    print("{}: {}".format(description, value))

setup_env_var('AGRAPH_HOST', 'localhost', 'Hostname')

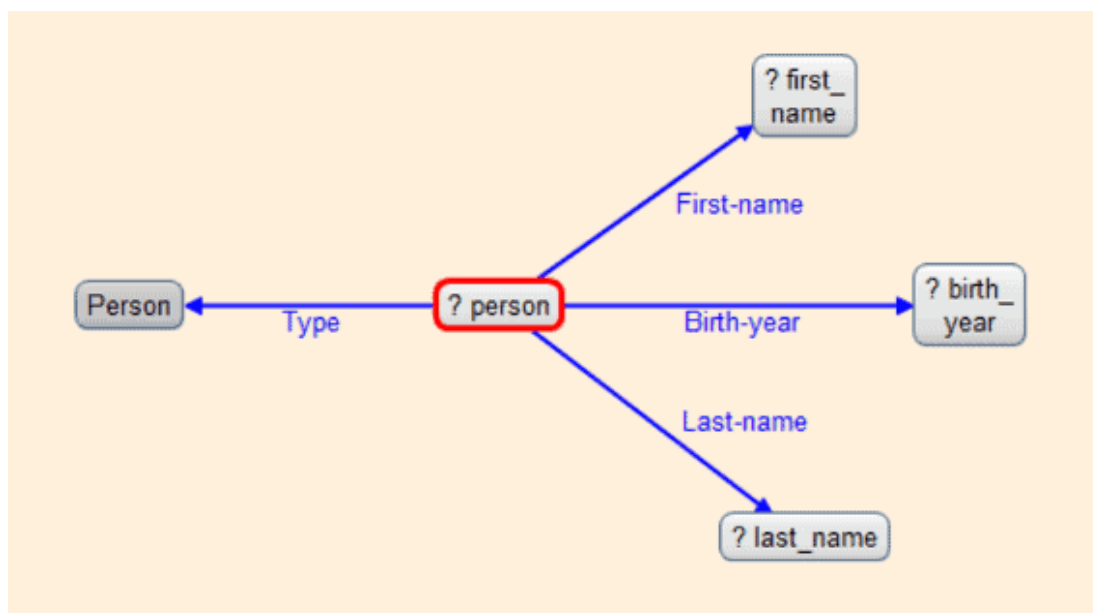
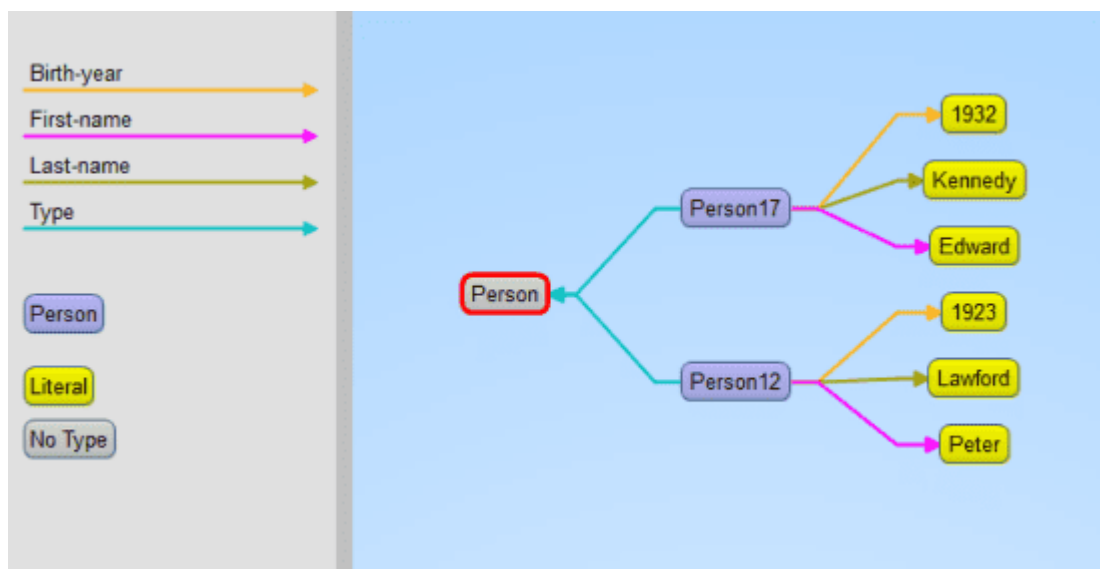
setup_env_var('AGRAPH_PORT', '10035', 'Port')

setup_env_var('AGRAPH_USER', '****', 'Username')
```

```
setup_env_var('AGRAPH_PASSWORD', '****', 'Password')
```

```
conn = ag_connect('kennedy', create=False, clear=False)
```

2. We then want to create a query. For this example, we will first show what our data looks like, what the visual query of the information is, and what the written query looks like. With the following query we want every person's first and last names, as well as their birth years. Here is a small portion of the data visualized in Gruff, and then the visualization of the query:



3. Then add the written query to the python script as a variable string (we added an additional line to the query to sort on birth year). Next use the API functionality to simply execute the query and turn the results into a pandas dataframe.

```
query = """select ?person ?first_name ?last_name ?birth_year
where
{ ?person <http://www.franz.com/simple#first-name> ?first_name
;
      <http://www.franz.com/simple#birth-year> ?birth_year
;
      rdf:type <http://www.franz.com/simple#person> ;
      <http://www.franz.com/simple#last-name> ?last_name .
}
order by desc(?birth_year)"""

with conn.executeTupleQuery(query) as result:
    df = result.toPandas()
```

When looking at the result, we see that we have a DataFrame!

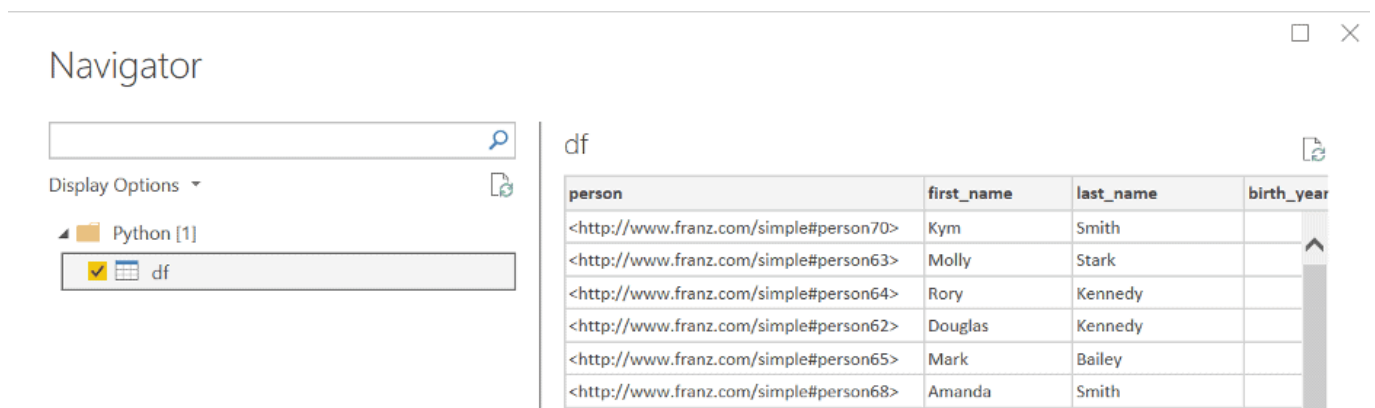


```
[5] df.head()
```

	person	first_name	last_name	birth_year
0	<http://www.franz.com/simple#person70>	Kym	Smith	1972
1	<http://www.franz.com/simple#person63>	Molly	Stark	1968
2	<http://www.franz.com/simple#person64>	Rory	Kennedy	1968
3	<http://www.franz.com/simple#person62>	Douglas	Kennedy	1967
4	<http://www.franz.com/simple#person65>	Mark	Bailey	1967

4. Now we will use this script in Power BI. When in Power BI Desktop, go to 'Get Data' and look for the python script option. Then simply copy and paste your entire script into the

text box, and run the script. In this case, our output looks like this:



5. Next simply 'Load' the data, and then you can use the Power BI Desktop interface to create whatever visualizations you want! If you do have a lot of additional operations to perform on your dataframe, we recommend doing these in your python script.

Method 2: POST Request:

For the SPARQL query via POST requests to work you need to url-encode the query. Every modern programming language will support that, but in our example we will be using Python again. This method is better for when you do not have python locally installed or prefer a different programming language.

It is possible to send a GET request from Power BI, but once the results from the query reach a certain size, a POST request is required, which is confusing to do within the Power BI Desktop interface. The following steps will show you how to do SPARQL Queries using POST requests. It looks a bit odd but it works well.

The Process:

1. In your AG WebView create an 'anonymous' user. (Go to admin -> Users -> [add a user] -> and add 'anonymous' as

username without adding a password). You can use these settings:

Users

anonymous [\[remove\]](#)

Roles: None

[\[suspend\]](#) [\[disable\]](#) [\[expire password\]](#)

☐ Superuser ☐ Start sessions ☐ Evaluate arbitrary code ☐ Control replication ☐ Two-phase commit

☒ Allow user attributes via HTTP header `x-user-attributes`

☐ Allow user attributes via SPARQL PREFIX `franzOption_userAttributes`

◦ `read/write` on all [\[remove\]](#)

Grant on catalog repository [\[ok\]](#)

Security Filters: `None` [\[add\]](#)

2. Go to your desired repository in WebView and Click on 'Queries' -> 'New'
3. Write a simple SPARQL query, and run it to make sure you get the correct response back.
4. In python create the following script: (Assuming your AllegroGraph is on your localhost port 10035 and your repo is called 'kennedy')

```
import urllib
```

```
def CreatePOSTquery(query):
```

```
    start = "http://anonymous:@localhost:10035/repositories/kennedy?queryL\nn=SPARQL&limit=1000&infer=false&returnQueryMetadata=false&chec\nkVariables=false&query="\n    response = start + urllib.parse.quote(query)\n    return response
```

This function url-encodes the query and attaches it to the POST request. Replace the 'localhost:10035' and 'kennedy' strings in the start variable with your corresponding data. Then, using the same query as our previous example, we create

our url-encoded POST query:

```
query = """select ?person ?first_name ?last_name ?birth_year
where
{ ?person <http://www.franz.com/simple#first-name> ?first_name
;
    <http://www.franz.com/simple#birth-year> ?birth_year
;
    rdf:type <http://www.franz.com/simple#person> ;
    <http://www.franz.com/simple#last-name> ?last_name .
}
order by desc(?birth_year)"""
```

```
result = CreatePOSTquery(query)
print(result)
```

This gives us the following result:

[illegible]

5. Within Power BI Desktop we go to 'Get data' and create a 'Blank query' and go into the 'Advanced Editor' window. Using the following format we will get our desired results (please note that due to the length of the url-encoded request, it did not all fit in the image. Copy and pasting into the url field works fine. The 'url' variable needs to be in quotes and have a comma at the end):

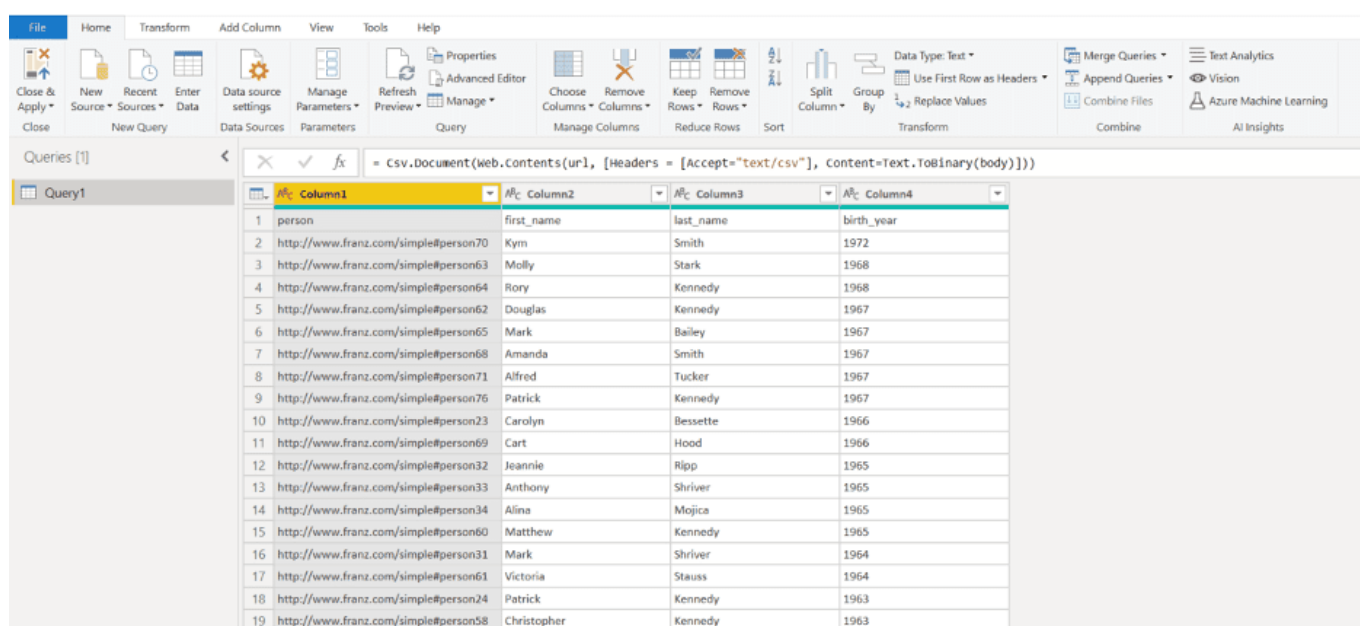
Query1

```

let
    url = "http://anonymous:@localhost:10035/repositories/kennedy?queryLn=SPARQL&limit=1000&infer=false&returnQuery",
    body = "",
    Source = Csv.Document(Web.Contents(url), [Headers = [Accept="text/csv"], Content=Text.ToBinary(body)]])
in
    Source

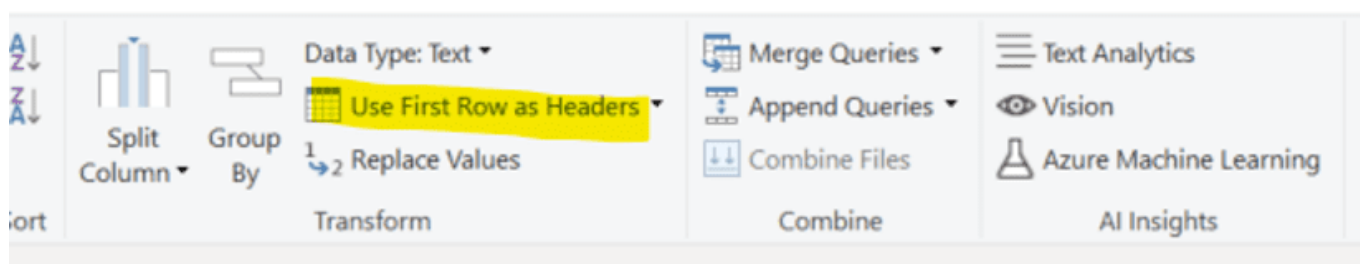
```

We see the following results:



Column1	Column2	Column3	Column4
person	first_name	last_name	birth_year
http://www.franz.com/simple#person70	Kym	Smith	1972
http://www.franz.com/simple#person63	Molly	Stark	1968
http://www.franz.com/simple#person64	Rory	Kennedy	1968
http://www.franz.com/simple#person62	Douglas	Kennedy	1967
http://www.franz.com/simple#person65	Mark	Bailey	1967
http://www.franz.com/simple#person68	Amanda	Smith	1967
http://www.franz.com/simple#person71	Alfred	Tucker	1967
http://www.franz.com/simple#person76	Patrick	Kennedy	1967
http://www.franz.com/simple#person23	Carolyn	Bessette	1966
http://www.franz.com/simple#person69	Cart	Hood	1966
http://www.franz.com/simple#person32	Jeannie	Ripp	1965
http://www.franz.com/simple#person33	Anthony	Shriver	1965
http://www.franz.com/simple#person34	Alina	Mojica	1965
http://www.franz.com/simple#person60	Matthew	Kennedy	1965
http://www.franz.com/simple#person31	Mark	Shriver	1964
http://www.franz.com/simple#person61	Victoria	Stauss	1964
http://www.franz.com/simple#person24	Patrick	Kennedy	1963
http://www.franz.com/simple#person58	Christopher	Kennedy	1963

6. One last step is to turn the top row into the column names, which can be achieved by pressing the 'Use first row as headers':



The best part about both of these methods is that once the query has been created, Power BI can refresh the visuals using the same queries if your data changed. This can be achieved by

scheduling refreshes within the Power BI Desktop interface (<https://docs.microsoft.com/en-us/power-bi/connect-data/refresh-data#configure-scheduled-refresh>)

Please send any questions or issues to: support@franz.com

AllegroGraph Named to 100 Companies That Matter Most in Data

Franz Inc. Acknowledged as a Leader for Knowledge Graph Solutions

Lafayette, Calif., June 23, 2020 – Franz Inc., an early innovator in Artificial Intelligence (AI) and leading supplier of Semantic Graph Database technology for Knowledge Graph Solutions, today announced that it has been named to The 100 Companies That Matter in Data by Database Trends and Applications. The annual list reflects the urgency felt among many organizations to provide a timely flow of targeted information. Among the more prominent initiatives is the use of AI and cognitive computing, as well as related capabilities such as machine learning, natural language processing, and text analytics. This list recognizes companies based on their presence, execution, vision and innovation in delivering products and services to the marketplace.

“We’re excited to announce our eighth annual list, as the industry continues to grow and evolve,” remarked Thomas Hogan, Group Publisher at Database Trends and Applications. “Now,

more than ever, businesses are looking for ways transform how they operate and deliver value to customers with greater agility, efficiency and innovation. This list seeks to highlight those companies that have been successful in establishing themselves as unique resources for data professionals and stakeholders.”

“We are honored to receive this acknowledgement for our efforts in delivering Enterprise Knowledge Graph Solutions,” said Dr. Jans Aasman, CEO, Franz Inc. “In the past year, we have seen demand for Enterprise Knowledge Graphs take off across industries along with recognition from top technology analyst firms that Knowledge Graphs provide the critical foundation for artificial intelligence applications and predictive analytics.

Our recent launch of AllegroGraph 7 with FedShard, a breakthrough that allows infinite data integration to unify all data and siloed knowledge into an Entity-Event Knowledge Graph solution will catalyze Knowledge Graph deployments across the Enterprise.”

Gartner recently released a report “How to Build Knowledge Graphs That Enable AI-Driven Enterprise Applications” and have previously stated, “The application of graph processing and graph databases will grow at 100 percent annually through 2022 to continuously accelerate data preparation and enable more complex and adaptive data science.” To that end, Gartner named graph analytics as a “Top 10 Data and Analytics Trend” to solve critical business priorities. (*Source: Gartner, Top 10 Data and Analytics Trends, November 5, 2019*).

“Graph databases and knowledge graphs are now viewed as a must-have by enterprises serious about leveraging AI and predictive analytics within their organization,” said Dr. Aasman “We are working with organizations across a broad range of industries to deploy large-scale, high-performance Entity-Event Knowledge Graphs that serve as the foundation for AI-

driven applications for personalized medicine, predictive call centers, digital twins for IoT, predictive supply chain management and domain-specific Q&A applications – just to name a few.”

Forrester Shortlists AllegroGraph

AllegroGraph was shortlisted in the February 3, 2020 Forrester Now Tech: Graph Data Platforms, Q1 2020 report, which recommends that organizations “Use graph data platforms to accelerate connected-data initiatives.” Forrester states, “You can use graph data platforms to become significantly more productive, deliver accurate customer recommendations, and quickly make connections to related data.”

Bloor Research covers AllegroGraph with FedShard

Bloor Research Analyst, Daniel Howard noted “With the 7.0 release of AllegroGraph, arguably the most compelling new capability is its ability to create what Franz refers to as “Entity-Event Knowledge Graphs” (or EEKGs) via its patented FedShard technology.” Mr. Howard goes on to state “Franz clearly considers this a major release for AllegroGraph. Certainly, the introduction of an explicit entity-event graph is not something I’ve seen before. The newly introduced text to speech capabilities also seem highly promising.”

AllegroGraph Named to KMWorld’s 100 Companies That Matter in Knowledge Management

AllegroGraph was also recently named to KMWorld’s 100 Companies That Matter in Knowledge Management. The KMWorld 100 showcases organizations that are advancing their products and capabilities to meet changing requirements in Knowledge Management.

Franz Knowledge Graph Technology and Services

Franz’s Knowledge Graph Solution includes both technology and

services for building industrial strength Entity-Event Knowledge Graphs based on best-of-class tools, products, knowledge, skills and experience. At the core of the solution is Franz's graph database technology, AllegroGraph with FedShard, which is utilized by dozens of the top F500 companies worldwide and enables businesses to extract sophisticated decision insights and predictive analytics from highly complex, distributed data that cannot be uncovered with conventional databases.

Franz delivers the expertise for designing ontology and taxonomy-based solutions by utilizing standards-based development processes and tools. Franz also offers data integration services from siloed data using W3C industry standard semantics, which can then be continually integrated with information that comes from other data sources. In addition, the Franz data science team provides expertise in custom algorithms to maximize data analytics and uncover hidden knowledge.

Ubiquitous AI Demands A New Type Of Database Sharding

Forbes published the following article by Dr. Jans Aasman, Franz Inc.'s CEO.



The notion of sharding has become increasingly crucial for selecting and optimizing database architectures. In many cases, sharding is a means of horizontally distributing data; if properly implemented, it results in near-infinite scalability. This option enables database availability for business continuity, allowing organizations to replicate databases among geographic locations. It's equally useful for load balancing, in which computational necessities (like processing) shift between machines to improve IT resource allocation.

However, these use cases fail to actualize sharding's full potential to maximize database performance in today's post-big data landscape. There's an even more powerful form of sharding, called "hybrid sharding," that drastically improves the speed of query results and duly expands the complexity of the questions that can be asked and answered. Hybrid sharding is the ability to combine data that can be partitioned into shards with data that represents knowledge that is usually unshardable.

This hybrid sharding works particularly well with the knowledge graph phenomenon leveraged by the world's top data-driven companies. Hybrid sharding also creates the enterprise scalability to query scores of internal and external sources for nuanced, detailed results, with responsiveness commensurate to that of the contemporary AI age.



[Read the full article at Forbes.](#)

Natural Language Processing and Machine Learning in AllegroGraph

The majority of our customers build Knowledge Graphs with Natural Language and Machine learning components. Because of this trend AllegroGraph now offers strong support for the use of Natural Language Processing and Machine learning.

Franz Inc has a team of NLP engineers and Taxonomy experts that can help with building turn-key solutions. In general however, our customers already have some expertise in house. In those cases we train customers in how to take the output of NLP and ML processing and turn that into an efficient Knowledge Graph based on best practices in the industry.

This document primarily describes the NLP and ML plug-in AllegroGraph.

Note that many enterprises already have a data science team with NLP experts that use modern open source NLP tools like Spacy, Gensim or Polyglot, or Machine Learning based NLP tools like BERT and Scikit-Learn. In another blog about Document Handling we describe a pipeline of how to deal with NLP in Document Knowledge Graphs by using our NLP and ML plugin and mix that with open source tools.

PlugIn features for Natural Language Processing and Machine Learning in AllegroGraph.

Here is the outline of the plugin features that we are going to describe in more detail.

Machine learning

- data acquisition
- classifier training
- feature extraction support
- performance analysis
- model persistence

NLP

- handling languages
- handling dictionaries
- tokenization
- entity extraction
- Sentiment analysis
- basic pattern matching

SPARQL Access

- Future development

Machine Learning

ML: Data Acquisition

Given that the NLP and ML functions operate within AllegroGraph, after loading the plugins, data acquisition can be performed directly from the triple-store, which drastically simplifies the data scientist workflow. However, if the data is not in AllegroGraph yet we can also import it directly from ten formats of triples or we can use our additional capabilities to import from CSV/JSON/JSON-LD.

Part of the Data Acquisition is also that we need to pre-process the data for training so we provide these three functions:

- prepare-training-data
- split-dev-test

- equalize (for resampling)

Machine Learning: Classifiers

- Currently we provide simple linear classifiers. In case there's a need for neural net or other advanced classifiers, those can be integrated on-demand.
- We also provide support for online learning (online machine learning is an ML method in which data becomes available in a sequential order and is used to update the best predictor for future data at each step, as opposed to batch learning techniques which generate the best predictor by learning on the entire training data set at once). This feature is useful for many real-world data sets that are constantly updated.
- The default classifiers available are Averaged Perceptron and AROW

Machine Learning: Feature Extraction

Each classifier is expecting a vector of features: either feature indices (indicative features) or pairs of numbers (index – value). These are obtained in a two-step process:

1. A classifier-specific extract-features method should be defined that will return raw feature vector with features identified by strings of the following form: prefix|feature.

The prefix should be provided as a keyword argument to the collect-features method call, and it is used to distinguish similar features from different sources (for instance, for distinct predicates).

2. Those features will be automatically transformed to unique integer ids. The resulting feature vector of indicator features may look like the following: #(1 123 2999 ...)

Note that these features may be persisted to AllegroGraph for repeated re-use (e.g. for experimenting with classifier hyperparameter tuning or different classification models).

Many possible features may be extracted from data, but there is a set of common ones, such as:

1. individual tokens of the text field
2. ngrams (of a specified order) of the text field
3. presence of a token in a specific dictionary (like, the dictionary of slang words)
4. presence/value of a certain predicate for the subject of the current triple
5. length of the text

And in case the user has a need for special types of tokens we can write specific token methods, here is an example (in Lisp) that produces an indicator feature of a presence of emojis in the text:

```
(defmethod collect-features ((method (eql :emoji)) toks &key
pred)
(dolist (tok toks)
(when (some #'(lambda (code)
(or (<= #x1F600 code #x1F64F)
(<= #x1F650 code #x1F67F)
(<= #x1F680 code #x1F6FF))))
(map 'vector #'char-code tok))
(return (list "emoji")))))
```

Machine Learning: Integration with Spacy

The NLP and ML community invents new features and capabilities at an incredible speed. Way faster than any database company can keep up with. So why not embrace that? Whenever we need something that we don't have in AllegroGraph yet we can call out to Spacy or any other external NLP tool. Here is an example of using feature extraction from Spacy to collect

indicator features of the text dependency parse relations:

```
(defmethod collect-features ((method (eql :dep)) deps &key
pred dep-type dep-labels)
  (loop :for ds :in deps :nconc
    (loop :for dep :in ds
      :when (and (member (dep-tag dep) dep-labels)
        (dep-head dep)
        (dep-tok dep))
      :collect (format nil "dep|~a|~a_~a"
        dep-type
        (tok-word (dep-head dep)
          (tok-word (dep-tok dep)))))))
```

The demonstrated integration uses Spacy Docker instance and its HTTP API.

Machine Learning: Classifier Analysis

We provide all the basic tools and metrics for classifier quality analysis:

- accuracy
- f1, precision, recall
- confusion matrix
- and an aggregated classification report

Machine Learning: Model Persistence

The idea behind model persistence is that all the data can be stored in AllegroGraph, including features and classifier models. AllegroGraph stores classifiers directly as triples. This is a far more robust and language-independent approach than currently popular among data scientists reliance on Python pickle files. For the storage we provide a basic triple-based format, so it is also possible to interchange the models using standard RDF data formats.

The biggest advantage of this approach is that when adding

text to AllegroGraph we don't have to move the data externally to perform the classification but can keep the whole pipeline entirely internal.

Natural Language Proccession (NLP)

NLP: Language Packs

Most of the NLP tools are language-dependent: i.e. there's a general function that uses language-specific model/rules/etc. In AllegroGraph, support for particular languages is provided on-demand and all the language-specific is grouped in the so called "language pack" or langpack, for short – a directory with a number of text and binary files with predefined names.

Currently, the langpack for English is provided at `nlp/langs/en.zip`, with the following files:

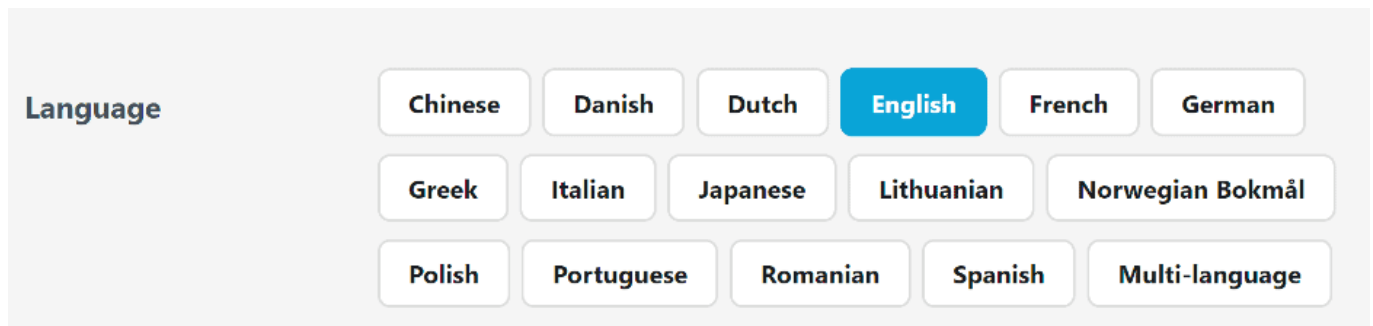
- `contractions.txt` – a dictionary of contractions
- `abbrs.txt` – a dictionary of abbreviations
- `stopwords.txt` – a dictionary of stopwords
- `pos-dict.txt` – positive sentiment words
- `neg-dict.txt` – negative sentiment words
- `word-tok.txt` – a list of word tokenization rules

Additionally, we use a general dictionary, a word-form dictionary (obtained from Wiktionary), and custom lexicons.

Loading a langpack for a particular language is performed using `load-langpack`.

Creating a langpack is just a matter of adding the properly named files to the directory and can be done manually. The names of the files should correspond to the names of the dictionary variables that will be filled by the pack. The dictionaries that don't have a corresponding file will be just skipped. We have just finished creating a langpack for Spanish and it will be published soon. In case you need other

dictionaries we use our AG/Spacy infrastructure. Spacy recently added a comprehensive list of new languages:



NLP: Dictionaries

Dictionaries are read from the language packs or other sources and are kept in memory as language-specific hash-tables. Alongside support for storing the dictionaries as text files, there are also utilities for working with them as triples and putting them into the triple store.

Note that we at Franz Inc specialize in Taxonomy Building using various commercial taxonomy building tools. All these tools can now export these taxonomies as a mix of SKOS taxonomies and OWL. We have several functions to read directly from these SKOS taxonomies and turn them into dictionaries that support efficient phrase-level lookup.

NLP: Tokenization

Tokenization is performed using a time-proven rule-based approach. There are 3 levels of tokenization that have both a corresponding specific utility function and an :output format of the tokenize function:

- :parags – splits the text into a list of lists of tokens for paragraphs and sentences in each paragraph
- :sents – splits the text into a list of tokens for each sentence
- :words – splits the text into a plain list of tokens

Paragraph-level tokenization considers newlines as paragraph delimiters. Sentence-level tokenization is geared towards western-style writing that uses dot and other punctuation marks to delimit sentences. It is, currently, hard-coded, but if the need arises, additional handling may be added for other writing systems. Word-level tokenization is performed using a language-specific set of rules.

NLP: Entity Extraction

Entity extraction is performed by efficient matching (exactly or fuzzy) of the token sequences to the existing dictionary structure.

It is expected that the entities come from the triple store and there's a special utility function that builds lookup dictionaries from all the triples of the repository identified by certain graphs that have a `skos:prefLabel` or `skos:altLabel` property. The lookup may be case-insensitive with the exception of abbreviations (default) or case-sensitive.

Similar to entity extraction, there's also support for spotting sentiment words. It is performed using the positive/negative words dictionaries from the langpack.

One feature that we needed to develop for our customers is 'heuristic entity extraction'. In case you want to extract complicated product names from text or call-center conversations between customers and agents you run into the problem that it becomes very expensive to develop altLabels in a taxonomy tool. We created special software to facilitate the automatic creation of altlabels.

NLP: Basic Pattern Matching for relationship and event detection

Getting entities out of text is now well understood and supported by the software community. However, to find complex concepts or relationships between entities or even events is

way harder and requires a flexible rule-based pattern matcher. Given our long time background in Lisp and Prolog one can imagine we created a very powerful pattern matcher.

SPARQL Access

Currently all the features above can be controlled as stored procedures or using Lisp as the command language. We have a new (beta) version that uses SPARQL for most of the control. Here are some examples. Note that `fai` is a magic-property namespace for “AI”-related stuff and `inc` is a custom namespace of an imaginary client:

1. Entity extraction

```
select ?ent {  
  ?subj fai:entityTaxonomy inc:products .  
  ?subj fai:entityTaxonomy inc:salesTerms .  
  ?subj fai:textPredicate inc:text .  
    ?subj fai:entity(fai:language "en", fai:taxonomy  
inc:products) ?ent .  
}
```

The expressions `?subj fai:entityTaxonomy inc:products` and `?subj fai:entityTaxonomy inc:salesTerms` specify which taxonomies to use (the appropriate matchers are cached).

The expression `?subj fai:entity ?ent` will either return the already extracted entities with the specified predicate (`fai:entity`) or extract the new entities according to the taxonomies in the texts accessible by `fai:textPredicate`.

2. `fai:sentiment` will return a single triple with sentiment score:

```
select ?sentiment {  
  ?subj fai:textPredicate inc:text .  
  ?subj fai:sentiment ?sentiment .  
  ?subj fai:language "en" .  
  ?subj fai:sentimentTaxonomy franz:sentiwords .  
}
```

3. Text classification:

Provided `inc:customClassifier` was already trained previously, this query will return labels for all texts as a result of classification.

```
select ?label {  
  ?subj fai:textPredicate inc:text .  
  ?subj fai:classifier inc:customClassifier .  
  ?subj fai:classify ?label .  
  ?label fai:storeResultPredicate inc:label .  
}
```

Further Development

Our team is currently working on these new features:

- A more accessible UI (python client & web) to facilitate NLP and ML pipelines
- Addition of various classifier models
- Sequence classification support (already implemented for a customer project)
- Pre-trained models shipped with AllegroGraph (e.g. English NER)
- Graph ML algorithms (deepwalk, Google Expander)
- Clustering algorithms (k-means, OPTICS)