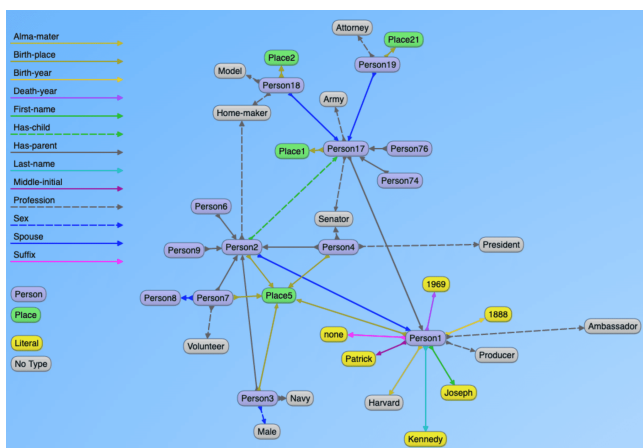


Graph Analytics with AllegroGraph and Apache Spark

AllegroGraph enables users to export data out of their Knowledge Graph and perform graph analytics with Apache Spark, one of the most popular platforms for large-scale data processing. Users immediately gain machine learning and SQL database solutions as well as GraphX and GraphFrames, two frameworks for running graph compute operations on data.

A key benefit of using Apache Spark for graph analytics within AllegroGraph is that it is built on top of Hadoop MapReduce and extends the MapReduce model to efficiently use more types of computations. Users can access interfaces (including interactive shells) for programming entire clusters with implicit data parallelism and fault-tolerance.

Apache Spark is one of the most popular platforms for large-scale data processing. In addition to machine learning, SQL database solutions, Spark also comes with GraphX and GraphFrames two frameworks for running graph compute operations on your data. In the AllegroGraph-Spark notebook, we demonstrate how to read data from AllegroGraph and then perform graph analytics with Spark.



Apache Spark was built on top of Hadoop MapReduce and it extends the MapReduce model to efficiently use more types of computations. It provides interfaces (including interactive shells) for programming entire clusters with implicit data parallelism

and fault-tolerance. For a quick start of more Spark APIs, please go to [here](#).

Visit our Github example page for more details.

The Next Step in Machine Learning's Evolution: Graph Neural Networks



Dr. Aasman recently publish this article on Graph Neural Networks at Toward Data Science.

The capacity to consistently attain enterprise value from mission critical machine learning deployments hinges on at least one of the following three applications: classifying entities, predicting events, and understanding why events happened.

No matter which technique is used, whether it includes supervised, unsupervised, or reinforcement learning, or if the scale and compute of deep learning is involved, conventional machine learning has limitations for solving these business problems.

It works well for many types of data, but incurs difficulty and outright failure when applied to high dimensionality networked datasets. These limits demand a new approach for social network research, recommendation engines, biology, chemistry, computer vision, and Natural Language Processing deployments in which context is pivotal.

Graph Neural Networks excel at predicting events, explaining them, and classifying entities at scale to deliver striking accuracy for these and other pragmatic deployments. Pairing their reasoning with semantic inferences creates additional knowledge for predicting events based on the multifaceted, contextualized relationships in high dimensionality data.

Read the full Article.

AiThority Interview with Dr. Jans Aasman



Jans Aasman, please tell us about your current role and the team / technology you handle at Franz.

As CEO of Franz Inc., I drive the overall technology vision for our Enterprise Knowledge Graph solutions and ensure our customer projects deliver the ROI results expected with graph based architectures.

Franz Inc. is composed of an expert team with skills in Graph Databases, Semantic technologies, Graph Visualization, AI, NLP and Machine Learning. Our domain knowledge encompasses large enterprises in Healthcare, Pharma, Customer Support, and Intelligence Agencies.

Our main business today revolves around AllegroGraph, a Semantic Graph platform that allows infinite data integration through a patented approach unifying all data and siloed knowledge into an Entity-Event Knowledge Graph solution that can support massive big data analytics. AllegroGraph's FedShard feature utilizes patented federated sharding capabilities that drive 360-degree insights and enable complex reasoning across a distributed Knowledge Graph. AllegroGraph is utilized by dozens of the top Fortune 500 companies worldwide.

We also offer a popular data visualization and no-code query builder called Gruff – the most advanced Knowledge Graph visualization application on the market, which we recently integrated into Franz AllegroGraph. Gruff enables users to create visual Knowledge Graphs that display data relationships in views that are driven by the user. Ad hoc and exploratory analysis can be performed by simply clicking on different graph nodes to answer questions. Gruff's unique 'Time Machine' feature provides the capability to explore temporal context and connections within data. The visual query builder within Gruff empowers both novice and expert users to create simple to highly complex queries without writing any code.

Read the full interview at AIThORITY.

NLP: Unlock the Hidden Business Value in Voice Communications



By Dr. Jans Aasman, CEO, Franz Inc.

Today organizations capture an enormous amount of information in spoken conversations, from routine customer service calls to sophisticated claims processing interactions in finance and healthcare. But most of this information remains hidden and unused due to the difficulty of turning these conversations into meaningful data that can be effectively analyzed through Natural Language Processing (NLP).

Simply applying speech recognition software to voice conversations often results in unreliable data. State-of-the-art speech recognition systems still have trouble distinguishing between homophones (words with the same pronunciation, but different meanings), as well as the difference between proper names (i.e. people, products) and separate words. In addition, there is also the challenge of identifying domain-specific words accurately. Thus, in most cases, using speech recognition software alone doesn't produce accurate enough data for reliable NLP.

Domain-specific taxonomies are key to understanding conversations via speech recognition systems. With them, we can feed conversations to knowledge graphs that understand the conversation and make connections in the data. Knowledge graphs provide the ability to extract the correct meaning of text from conversations and connect concepts in order to add business value.

Knowledge graphs fed with NLP provide two prime opportunities for monetization. First, organizations can better understand their customers to improve products and services more to their liking, which in turn boosts marketing, sales and customer

retention rates. Secondly, this analysis gives contact center agents real-time support for optimizing customer interactions to produce faster resolutions, better conversion rates, and cross-selling and up-selling opportunities. These approaches enable companies to capitalize on speech recognition knowledge graphs, accelerate their ROI, and expand their bottom lines.

Taxonomy Driven Speech Recognition

The story of taxonomy-driven speech recognition closely relates to knowledge graphs. The first wave of knowledge graphs was built from taking structured data and turning it into semantic graphs that support the linked open data movement. The next wave is all about unstructured data. People started doing Natural Language Processing on documents and textual conversations like emails and chats. Doing so accurately for a given domain requires a taxonomy to understand the words and concepts. Otherwise, downstream processes like entity extraction and event detection won't work.

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

The Future of AI: Machine Learning and Knowledge Graphs

Bringing knowledge graph and machine learning technology together can improve the accuracy of the outcomes and augment the potential of machine learning approaches. With knowledge graphs, AI language models are able to represent the relationships and accurate meaning of data instead of simply generating words based on patterns.

Read this special report to dive into key uses cases, best practices for getting started, and technology solutions every organization should know about.

The Future of AI: Machine Learning and Knowledge Graphs

Understanding What Matters With Text Analytics and NLP

Dr. Jans Aasman was quoted extensively in this KMWorld Article:



Whether employing traditional rulesbased approaches to text analytics or leveraging more modern machine learning strategies, users must initially train the systems on

relevant business domains. One way to do so is with comprehensive taxonomies of terms, their synonyms, and their meanings—which are traditionally associated with rules-based models. According to Franz CEO Jans Aasman, “There’s a part of NLP where people create taxonomies and ontologies. That is just a very acceptable way of doing NLP.” Historically, such defined hierarchies of vocabularies were paired with rules to find patterns in text and create actions such as classifications or entity extraction.

The trade-off between this approach and the taxonomic one is clear: Organizations can forsake the extensive time required to build taxonomies by simply using annotated training data.

The objective is to “just throw statistics and machine learning at the problem so it will all automatically work,” Aasman said. Although reduced time-to-value is an advantage of the deep learning approach, there are issues to consider, including the following:

- ◆ **Training data:** Machine learning models require immense amounts of training data, which organizations might not have for their domains. Transfer learning solves this problem by enabling subject matter experts to upload a couple of hundred examples (instead of thousands), highlight them, and teach dynamic models “the representative entities, key-value pairs, and classes they’re trying to derive from these documents,” Wilde noted.

- ◆ **Controlled vocabularies:** Transformers and techniques such as Bidirectional Encoder Representations from Transformers (BERT) reduce the training data quantities for machine learning models, broaden the array of training data that’s relevant, and implement a controlled vocabulary that otherwise isn’t as defined as taxonomic ones. Thus, organizations can take a phrase and “generate a similar phrase that means the same, but can be used in multiple reports in a controlled way,” Mishra said. Additionally, it’s possible to simply purchase libraries of terms and definitions. “Many companies end up buying those things to be able to incorporate those capabilities,” Shankar added.

- ◆ **Practical knowledge:** Exclusively using machine learning models to train text analytics decreases the real-world understanding and applicability of text. “People that do machine learning don’t want to spend the effort to create a vocabulary or the pragmatics or the semantics,” Aasman noted. “Machine learning has a place in all of this, but it misses part of the whole future solution where we have systems that understand what people are talking about.”

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

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:

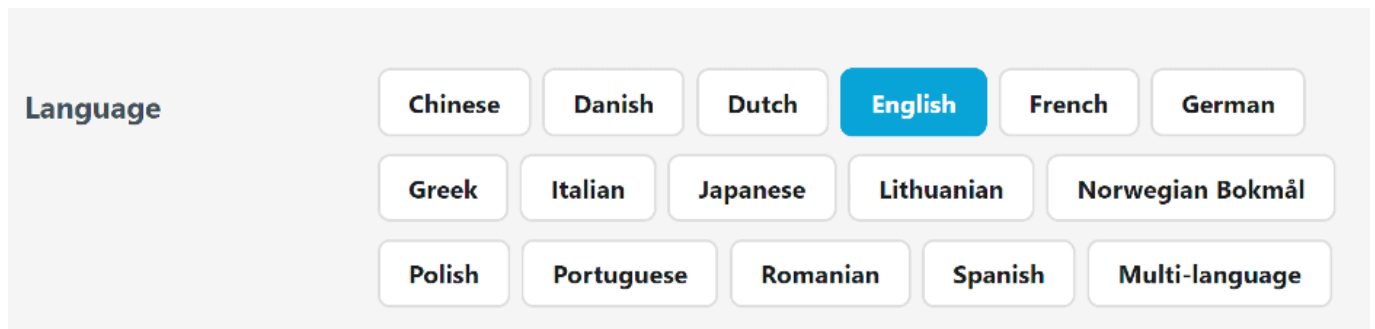
- `contractions.txt` – a dictionary of contractions
- `abbrs.txt` – a dictionary of abbreviations
- `stopwords.txt` – a dictionary of stopwords
- `pos-dict.txt` – positive sentiment words
- `neg-dict.txt` – negative sentiment words
- `word-tok.txt` – a list of word tokenization rules

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)
-

Document Knowledge Graphs with NLP and ML

A core competency for Franz Inc is turning text and documents into Knowledge Graphs (KG) using Natural Language Processing (NLP) and Machine Learning (ML) techniques in combination with AllegroGraph. In this document we discuss how the techniques described in [NLP and ML components of AllegroGraph] can be combined with popular software tools to create a robust Document Knowledge Graph pipeline.

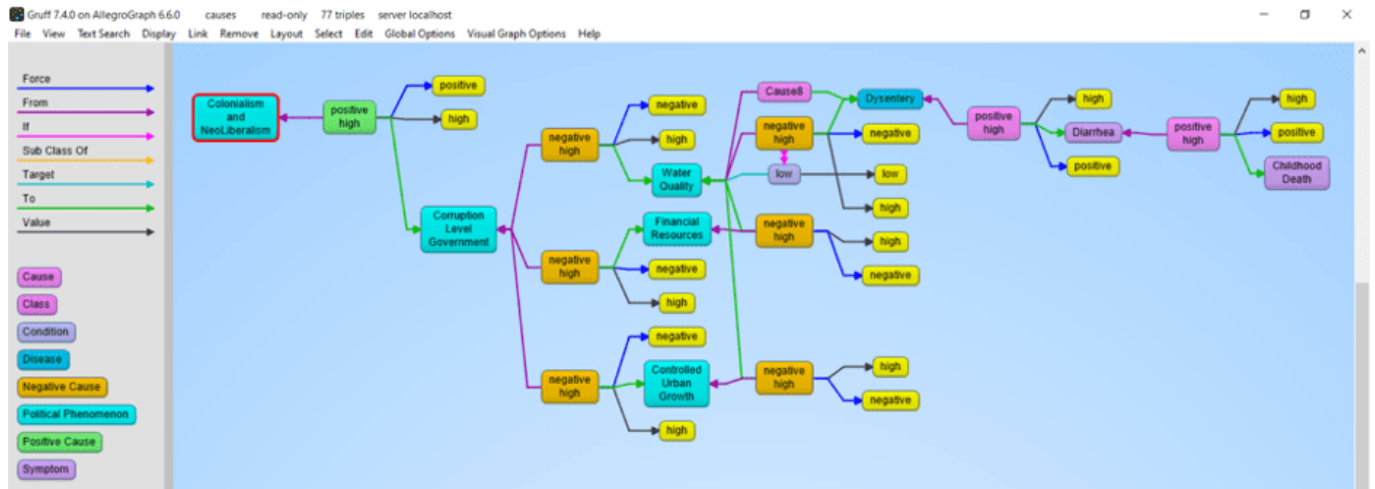
We have applied these techniques for several Knowledge Graphs but in this document we will primarily focus on three completely different examples that we summarize below. First is the Chomsky Legacy Project where we have a large set of very dense documents and very different knowledge sources, Second is a knowledge graph for an intelligent call center where we have to deal with high volume dynamic data and real-time decision support and finally, a large government organization where it is very important that people can do a semantic search against documents and policies that steadily change over time and where it is important that you can see the history of documents and policies.

Example [1] Chomsky Knowledge Graph

The Chomsky Legacy Project is a project run by a group of admirers of Noam Chomsky with the primary goal to preserve all his written work, including all his books, papers and interviews but also everything written about him. Ultimately students, researchers, journalists, lobbyists, people from the AI community, and linguists can all use this knowledge graph for their particular goals and questions.

The biggest challenges for this project are finding causal relationships in his work using event and relationship extraction. A simple example we extracted from an author

quoting Chomsky is that neoliberalism ultimately causes childhood death.



Example 2: N3 Results and the Intelligent Call Center

This is a completely different use case (See a recent KMWorld Article <https://allegrograph.com/knowledge-graphs-enhance-customer-experience-through-speed-and-accuracy/>). Whereas the previous use case was very static, this one is highly dynamic. We analyze in real-time the text chats and spoken conversations between call center agents and customers. Our knowledge graph software provides real-time decision support to make the call center agents more efficient. N3 Results helps big tech companies to sell their high tech solutions, mostly cloud-based products and services but also helps their clients sell many other technologies and services.

The main challenge we tackle is to really deeply understand what the customer and agent are talking about. None of this can be solved by only simple entity extraction but requires elaborate rule-based and machine learning techniques. Just to give a few examples. We want to know if the agent talked about their most important talking points: that is, did the agent ask if the customer has a budget, or the authority to make a decision or a timeline about when they need the new technology or whether they actually have expressed their need. But also whether the agent reached the right person, and whether the agent talked about the follow-up. In addition, if the customer

talks about competing technology we need to recognize that and provide the agent in real-time with a battle card specific to the competing technology. And in order to be able to do the latter, we also analyzed the complicated marketing materials of the clients of N3.

Example 3: Complex Government Documents

Imagine a regulatory body with tens of thousands of documents. Where nearly every paragraph has reference to other paragraphs in the same document or other documents and the documents change over time. The goal here is to provide the end-users in the government with the right document given their current task at hand. The second goal is to keep track of all the changes in the documents (and the relationship between documents) over time.

The Document to Knowledge Graph Pipeline

Process Name	Input	Output
1. Custom Taxonomy Creation	Corpus Analytics, Taxonomy tool	A SKOS taxonomy containing concepts, concept hierarchy, prefLabels, altLabels.
2. Document Preparation	Documents (pdf, word, ppt, xlsx), Apache Tika, Spacy for XML cleanup	An XML version of each document
3. Extract Document Meta Data	Document + Apache Tika	JSON dictionary of the Document MetaData
4. XML-to-Triples	XML+JSON dictionary, XMLToTriples.py	Graph-based document tree with chapters, sections, and paragraphs as triples. Also includes meta data as triples
5. Entity-Extraction	Paragraphs + taxonomies + AllegroGraph Entity extract or external extractors	Concepts, persons, places, currencies. Connected to paragraphs
6. LOD Enrichment	Paragraphs + IBM Natural Language Understanding.	Concept categories and links to DBpedia and GeoNames, etc.
7. Complex Relationship and Event extraction.	Paragraphs + Taxonomy + Rules in Spacy or AllegroGraph	Complex events and relationships, References to other document sections.
8. NLP and ML	Chapters and paragraphs + all the tools described [here], but also using Spacy, Gensim, BERT, SciKit Learn.	Similarities, sentiment, query answering, smart search, text classification, word embeddings, abstracts
9. Versioning and Document tracking	Old + New document, compare.py	Old document in historic repository, new document in current, changed graph.
10. Statistical Relationships	Concepts + OddRatio.py or OddsRatio.cl	Statistical relationships between concepts.

Let us first give a quick summary in words of how we turn documents into a Knowledge Graph.

[1] Taxonomy Creation

Taxonomy of all the concepts important to the business using open source or commercial taxonomy builders. An available industry taxonomy is a good starting point for additional customizations.

[2] Document Preparation

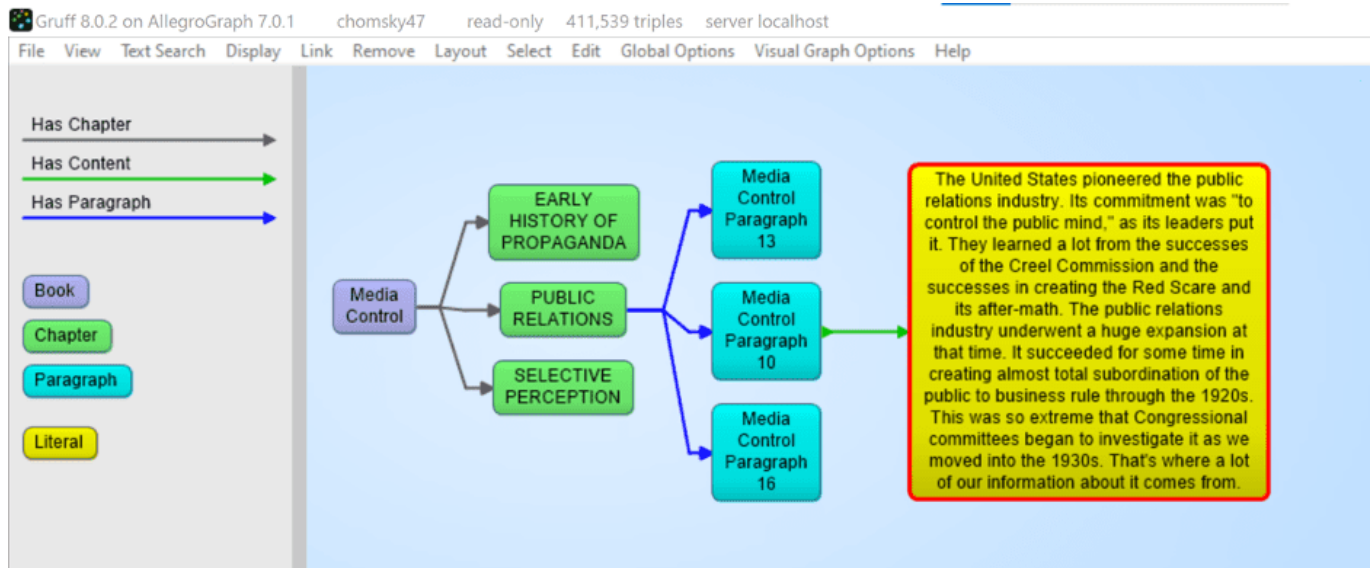
We then take a document and turn it into an intermediate XML using Apache Tika. Apache Tika supports more than 1000 document types and although Apache Tika is a fantastic tool, the output is still usually not clean enough to create a graph from, so we use Spacy rules to clean up the XML to make it as uniform as possible.

[3] Extract Document MetaData

Most documents also contain document metadata (author, date, version, title, etc) and Apache Tika will also deliver the metadata for a document as a JSON object.

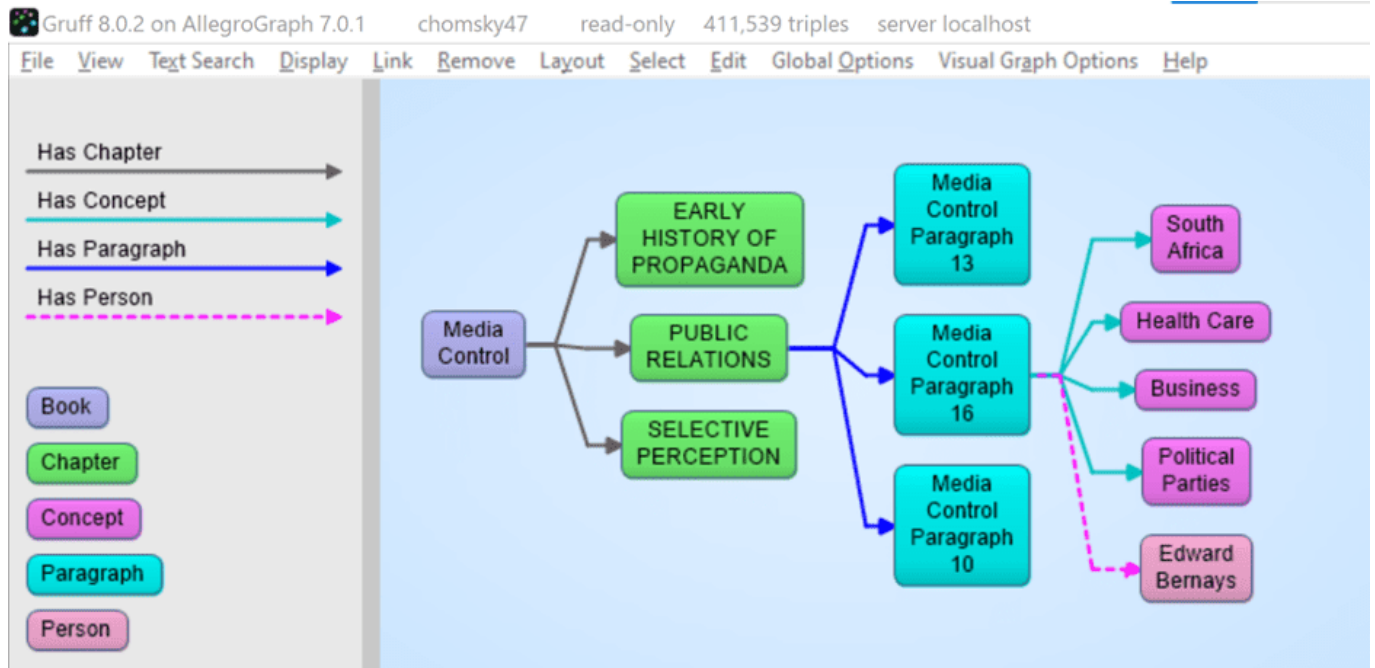
[4] XML to Triples

Our tools ingest the XML and metadata and transform that into a graph-based document tree. The document is the root and from that, it branches out into chapters, optionally sections, all the way down to paragraphs. The ultimate text content is in the paragraphs. In the following example we took the XML version of Noam Chomsky's book Media Control and turned that into a tree. The following shows a tiny part of that tree. We start with the Media Control node, then we show three (of the 11) chapters, for one chapter we show three (of the 6) paragraphs, and then we show the actual text in that paragraph. We sometimes can go even deeper to the level of sentences and tokens but for most projects that is overkill.



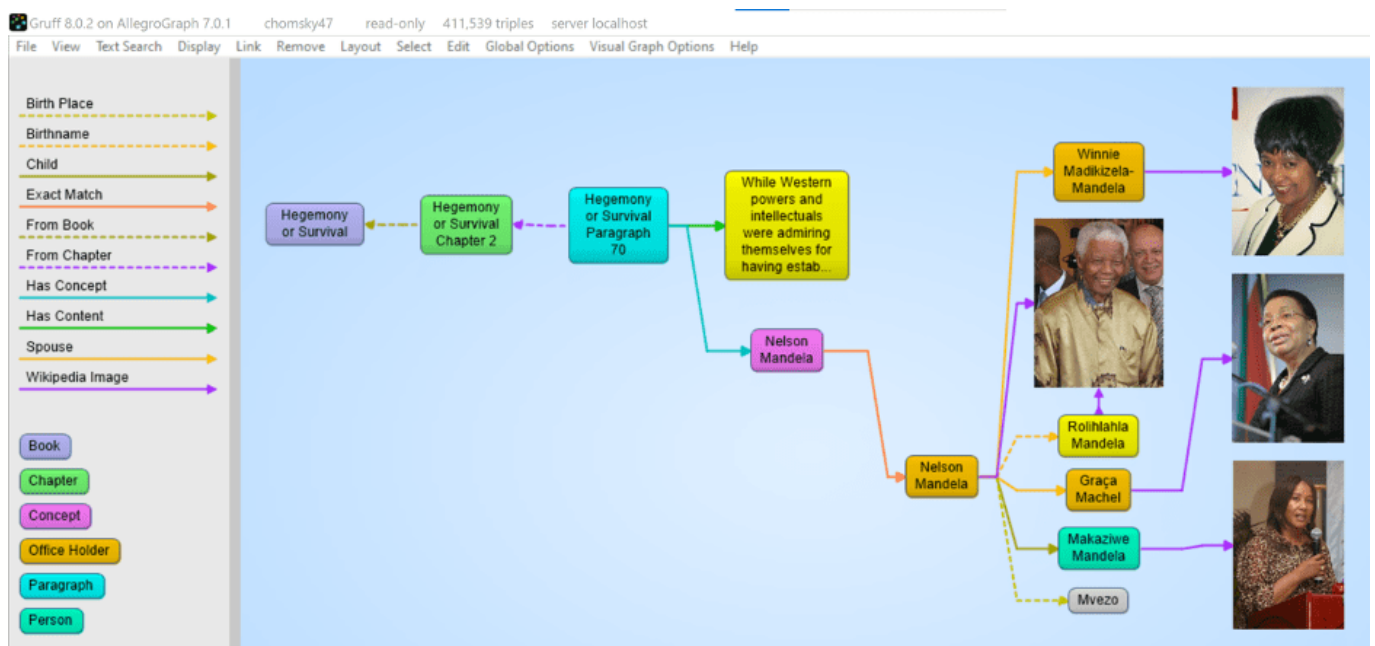
[5] Entity Extractor

AllegroGraph's entity extractor takes as input the text of each paragraph in the document tree and one or more of the taxonomies and returns recognized SKOS concepts based on prefLabels and altLabels. AllegroGraph's entity extractor is state of the art and especially powerful when it comes to complex terms like product names. We find that in our call center a technical product name can sometimes have up to six synonyms or very specific jargon. For example the Cisco product Catalyst 9000 will also be abbreviated as the cat 9k. Instead of developing altLabels for every possible permutation that human beings *will* use, we have specialized heuristics to optimize the yield from the entity extractor. The following picture shows 4 (of the 14) concepts discovered in paragraph 16. Plus one person that was extracted by IBM's NLU.



[6] Linked Data Enrichment

In many use cases, AllegroGraph can link extracted entities to concepts in the linked data cloud. The most prominent being DBpedia, wikidata, the census database, GeoNames, but also many Linked Open Data repositories. One tool that is very useful for this is IBM's Natural Language Understanding program but there are others available. In the following image we see that the Nelson Mandela entity (Red) is linked to the dbpedia entity for Nelson Mandela and that then links to the DBpedia itself. We extracted some of his spouses and a child with their pictures.



[7] Complex Relationship and Event Extraction

Entity extraction is a first good step to 'see' what is in your documents but it is just the first step. For example: how do you find in a text whether company C1 merged with company C2. There are many different ways to express the fact that a company fired a CEO. For example: Uber got rid of Kalanick, Uber and Kalanick parted ways, the board of Uber kicked out the CEO, etc. We need to write explicit symbolic rules for this or we need a lot of training data to feed a machine learning algorithm.

[8] NLP and Machine Learning

There are many many AI algorithms that can be applied in Document Knowledge Graphs. We provide best practices for topics like:

- [a] Sentiment Analysis, using good/bad word lists or training data.

- [b] Paragraph or Chapter similarity using statistical techniques like Gensim similarity or symbolic techniques where we just the overlap of recognized entities as a function of the size of a text.

- [c] Query answering using word2vec or more advanced techniques like BERT

- [d] Semantic search using the hierarchy in SKOS taxonomies.

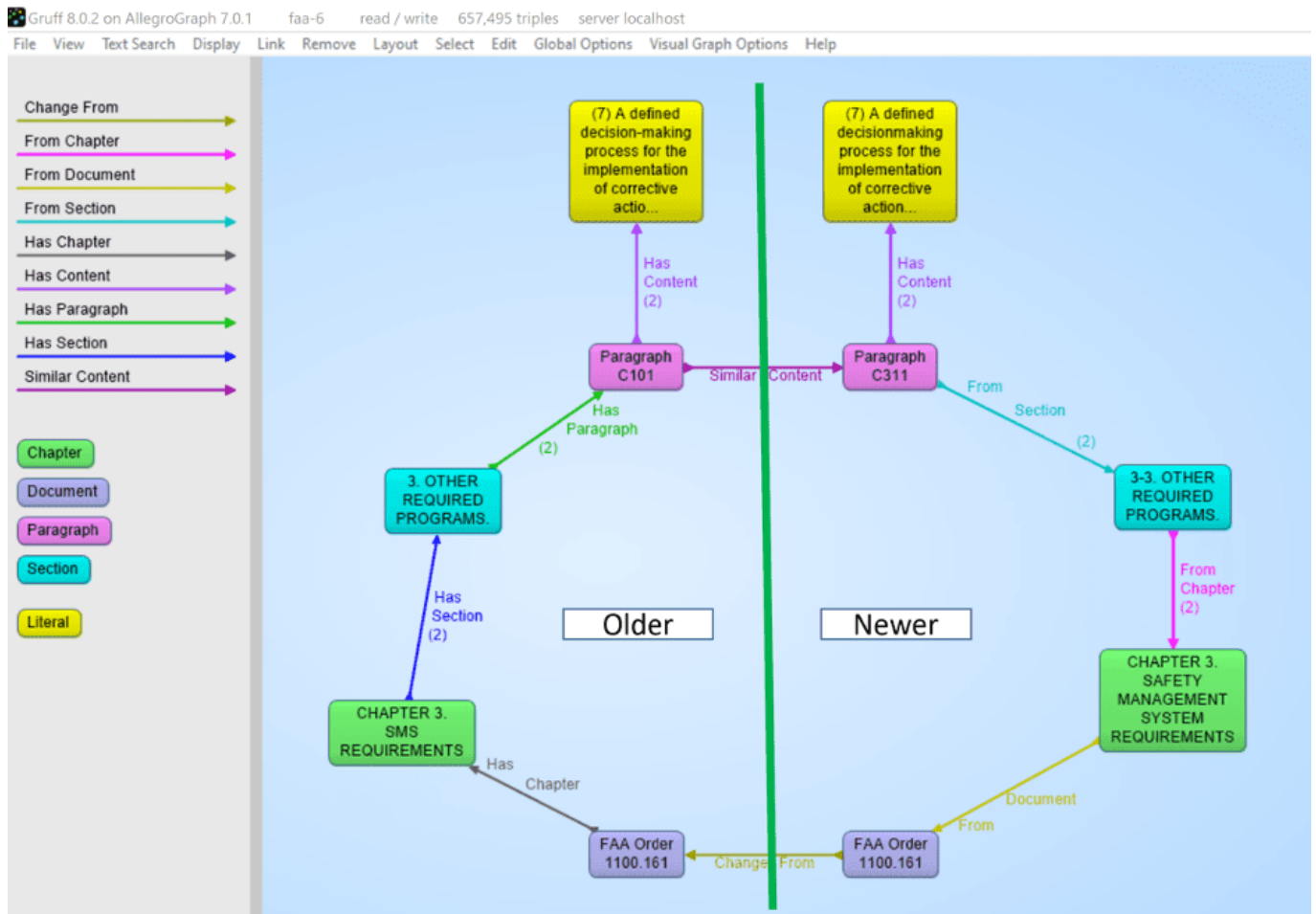
- [e] Summarization techniques for Abstractive or Extractive abstracts using Gensim or Spacy.

[9] Versioning and Document tracking

Several of our customers with Document Knowledge Graphs have noted the one constant in all of these KGs is that documents change over time. As part of our solution, we have created best practices where we deal with these changes. A crucial first step is to put each document in its own graph (i.e. the fourth element of every triple in the document tree is the document id itself). When we get a new version of a document the document ID changes but the new document will point back to the old version. We then compute which paragraphs stayed the same within a certain margin (there are always changes in whitespace) and we materialize what paragraphs disappeared in the new version and what new paragraphs appeared compared to the previous version. Part of the best practice is to put the old version of a document in a historical database that at all times can be federated with the 'current' set of documents.

Note that in the following picture we see the progression of a document. On the right hand side we have a newer version of a document 1100.161 with a chapter -> section -> paragraph -> contents where the content is almost the same as the one in

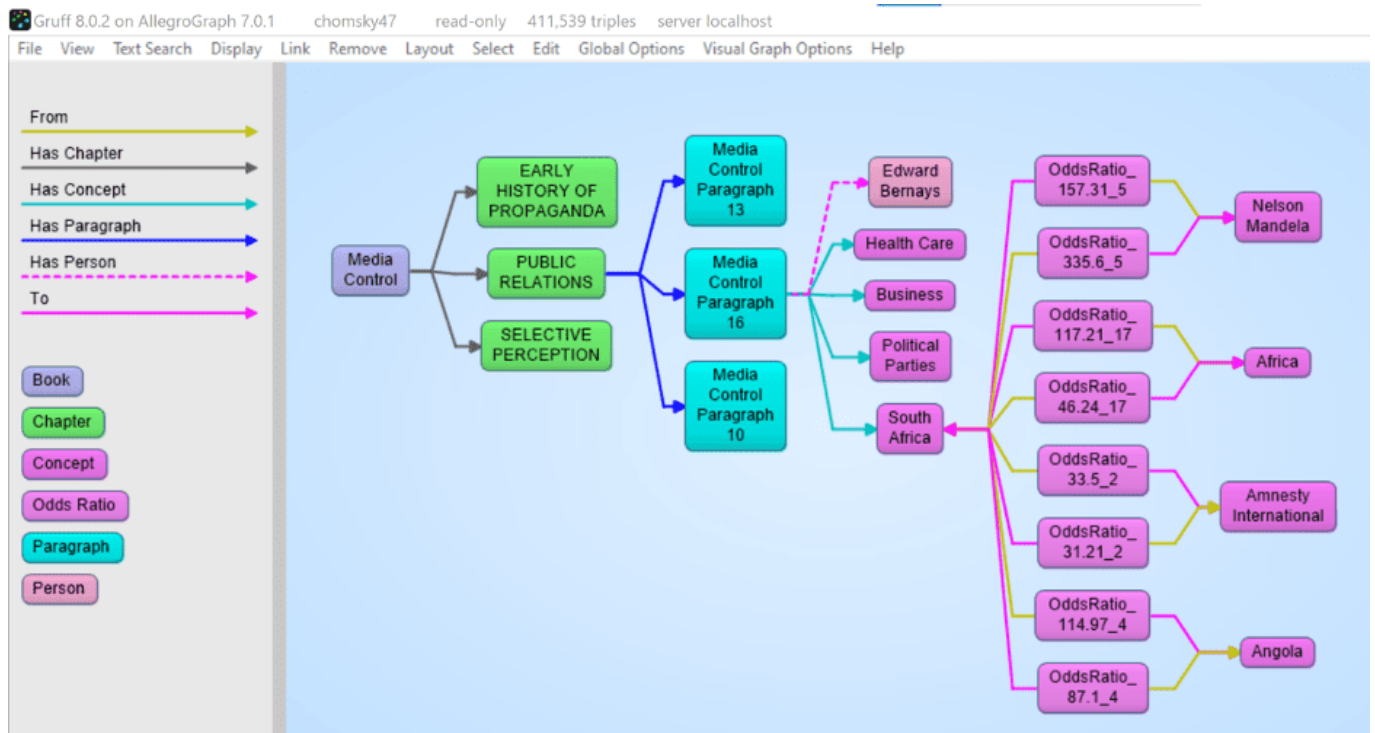
the older version. But note that the newer one spells 'decision making' as one word whereas the older version said 'decision-making'. Note that also the chapter titles and the section titles are almost the same but not entirely. Also, note that the new version has a back-pointer (changed-from) to the older version.



[10] Statistical Relationships

One important analytic one can do on documents is to look at the co-occurrence of terms. Although, given that certain words might occur more frequently in text, we have to correct the co-occurrence between words for the frequency of the two terms in a co-occurrence to get a better idea of the 'surprisingness' of a co-occurrence. The platform offers several techniques in Python and Lisp to compute these co-occurrences. Note that in the following picture we computed the odds ratios between recognized entities and so we see in

the following gruff picture that if Noam Chomsky talks about South Africa then the chances are very high he will also talk about Nelson Mandela.



Montefiore Semantic Data Lake Tackles Predictive Analytics

Montefiore Medical Center is preparing to launch a sophisticated predictive analytics program for crisis patients, which is rooted in its real-time semantic data lake technology.

Semantic computing is becoming a hot topic in the healthcare industry as the first wave of big data analytics leaders looks to move beyond the basics of population health management, predictive analytics, and risk stratification.

This new approach to analytics eschews the rigid, limited

capabilities of the traditional relational database and instead focuses on creating a fluid pool of standardized data elements that can be mixed and matched on the fly to answer a large number of unique queries.

Montefiore Medical Center, in partnership with Franz Inc., is among the first healthcare organizations to invest in a robust semantic data lake as the foundation for advanced clinical decision support and predictive analytics capabilities.

Read the full article at [Health IT Analytics](#)