

# Applying High-speed Pattern Recognition to Generate Queryable Semantics from Big Data

Chuck Rehberg

CTO Trigent Software

Chief Scientist at Semantic Insights

October 2014





# A Semantic Translation Solution

- The raw data represents the lowest level of abstraction available. Often it is the only level of abstraction available.
- A high-speed complex pattern recognizer can be employed to sift through and correlate, recognize and generate aggregate data representing a higher level of abstraction. Multiple possible interpretations can be maintained.
- Additional high-speed complex pattern recognizers can operate on this new higher level data to generate new aggregate data representing an even higher level of abstraction. This can be continued to allow flexible changes to each data transformation level. Again, multiple possible interpretations can be maintained.
- *The goal is to produce data representing higher levels of abstractions which constitute semantic representations required by the domain-knowledgeable investigators.*

The diagram illustrates the process of semantic translation. It starts with a block of raw binary data (0s and 1s). A curved arrow points from this data to a list of key and response events. A second curved arrow points from the list of events to a final line of text that identifies the events as potential security threats.

```
...11 1001010101000100110100100010001000100010001000100001001000101001001001000  
10010001000100010001010100001001001010010101001001001010001010001010010111101001001  
0010010010100100101010100101000101001001001001001001001001001001001001001010100  
100101010010101010100101010010101010010010010010010010010010010010010010010010  
0101001001111001010101000100110100100010001000100010001000100010001000100010001  
001001000100100010001000101010000100101001010010100100100100100100100100101011  
101001001001001001010010010101010010100010100100100100100100100100100100100100  
10101010010010101001010100101010010101010101010101010100100100100100100100100  
1010010100101001001001...
```

... Key 12, Key 37 + Resp 13, Key 33, Key 14 + Resp 11, Key 12, Key 37 + Resp 15, Key 33,  
Key 14 + Resp 19, Key 12, Key 37 + Resp 15, Key 33, Key 14 + Resp 11, Key 12, Key 39 +  
Resp 13, Key 33, Key 14 + Resp 11, Key 12, Key 37 + Resp 13, Key 33, Key 14 + Resp 11, Key  
12, Key 37 + Resp 15, Key 33, Key 14 + Resp 19, Key 12, Key 37 + Resp 15, Key 33, Key 14 +  
Resp 11, Key 12, Key 39 + Resp 13, Key 33, Key 14 + Resp 11...

... Login – Check Balance – Transfer \*failed – Check Balance – Transfer \*failed\* - Read Notice –  
Idle time out – Login \*failed – Login – View Stock - ...

... [possible hacking attempt] or [possible phishing attempt]...

# Two Key Translation Challenges must be met

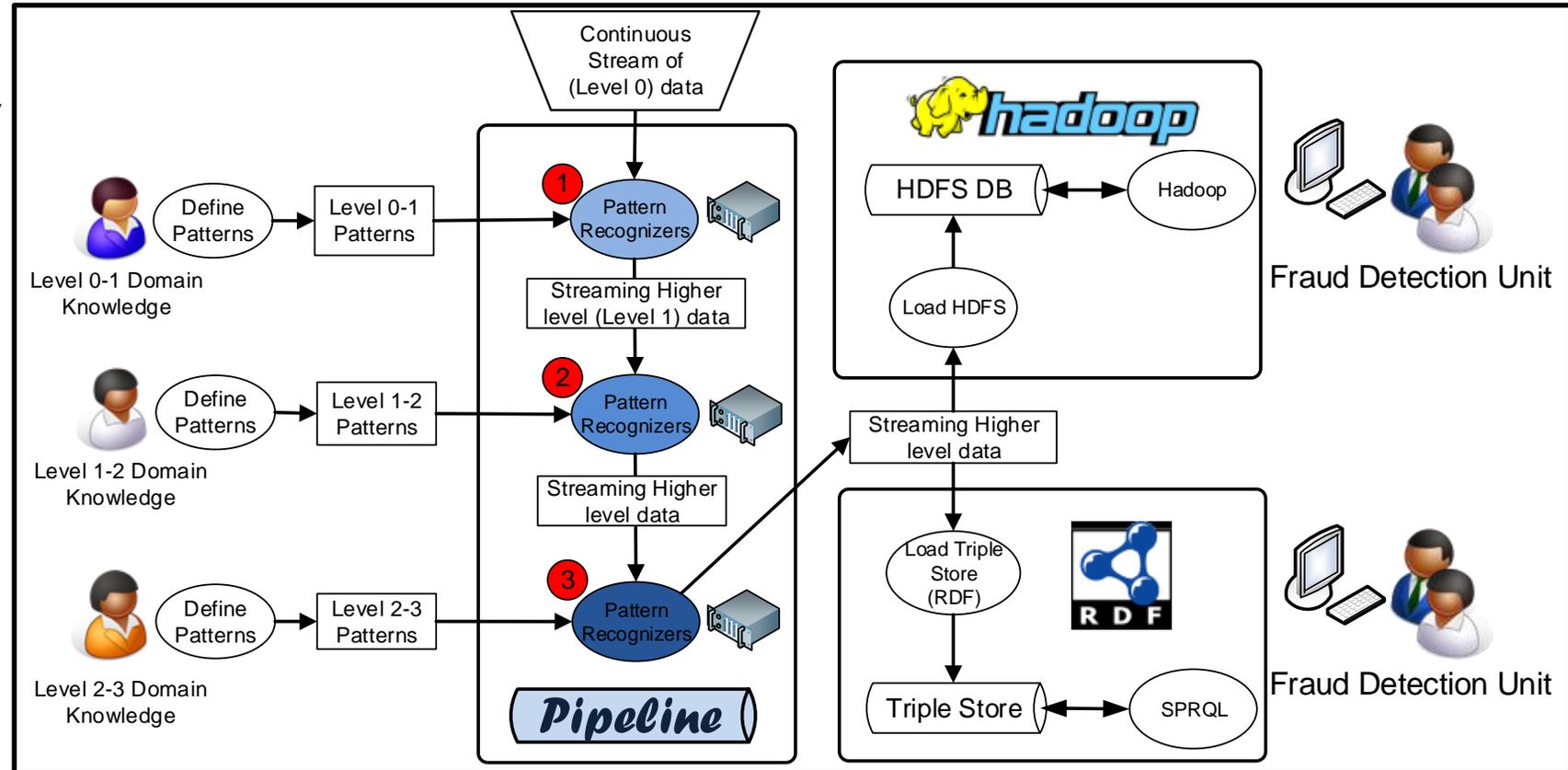
1. One key challenge is to flexibly transform the raw data into data that domain-knowledgeable investigators can directly query.
  - The transformation may require considering at a vast amount of data in sequence before a proper transformation can be made.
  - Multiple possible interpretations of the data may need to be maintained.
  - Individual transformations should be allowed to evolve over time to become more accurate
  - The data may contain information that does not contribute (or even detracts from) the semantic interpretation of the data. The data of interest may be surrounded by semantic “noise”.
2. For time-sensitive data, you will need to transform the data in real-time as it arrives.

# Real-time Semantic Transformation

- Our High-speed complex pattern recognizer is based on the Fast Rule Selection Engine (FRSE). US Patent #7,433,858
  - We have made additional modifications to the FRSE Engine to efficiently handle time-sequenced pattern matching.
  - Patterns can match adjacent tokens or span intervening non-matching tokens (noise).
  - The number of tokens to be scanned at a time can be pre-specified or adjust dynamically as needed.
- Multiple simultaneous levels of transformation in real-time is accomplished by pipelining.
  - Where each level hands off to the next level for further processing while continuing to process the new input.

# Transforming Streamed Data in Real-Time: A Multi-layered Pipelined Approach

1. The continuous stream of data is transformed **1** into objects recognized in another transformation **2** process which in turn produces objects that are sent to a third transformation **3** process, and so forth.
2. This can continue an arbitrary number of times with each transformation running independently in parallel.
3. This is well suited to problems requiring a cascade of information toward increasing levels of abstraction.



# Specifying Transformation Patterns Flexibly

Transformation Patterns processed by the FRSE Engine are powerful and simple to define:

1. Transformation Patterns are sequences of Token Expressions (i.e. expressions that match Match Tokens)
2. Transformation Patterns also produce a Pattern Result
3. Pattern Result creates a new sequence of one or more Result Tokens in the output data
4. Match Tokens are any data item that can be matched in the input data
5. Result Tokens are created by Transformation Patterns
6. Token Expressions state what must be true for a Match Token to be correctly identified
7. Result Tokens can act like Match Tokens in subsequent Transformation Patterns

Specifically:

```
<Transformation Pattern> ::= <Token Expression>(1,n) --> [<directive>] <Result Token>(1,n)
<Token Expression> ::= [<Match Operator>] <Match Token> /* default Match Operator is '=' */
<Match Operator> ::= { = | < | > | < | <= | => }
<Directive> ::= { ADD | REPLACE | RESULT } /* ADD/REPLACE modifies input; RESULT sends to output */
<Match Token> ::= <Token Set ID> | token(1,n)
<Token Set ID> ::= token(1,n) /* Token Set ID acts like a synonym */
<Result Token> ::= token
```

# Summary

- Big Data, representing real world events, often needs to be filtered and transformed into semantically meaningful data before domain-knowledgeable analysts can perform queries.
- Today, interactive data visualization techniques are often used to display *data aggregations and relationships* to help domain-knowledgeable analysts perform generalized “visual queries”.
- However, data visualization can be particularly unrevealing when data aggregation requires recognizing content-dependent patterns in a field of physical and semantic noise.
- High-speed Pattern Recognizers are capable of filtering and transforming real-time Big Data for immediate high-level evaluation.
- For more information, please visit us at [www.semanticinsights.com](http://www.semanticinsights.com).

# Chuck Rehberg



As CTO at Trigent Software and Chief Scientist at Semantic Insights, Chuck Rehberg has developed patented high performance rules engine technology and advanced natural language understanding technologies that empower a new generation of semantic research solutions.

Chuck has more than thirty years in the high-tech industry, developing leading-edge solutions in the areas of Artificial Intelligence, Semantic Technologies, analysis and large -scale configuration software.



**SEMANTIC INSIGHTS™**

A Division of Trigent Software, Inc.