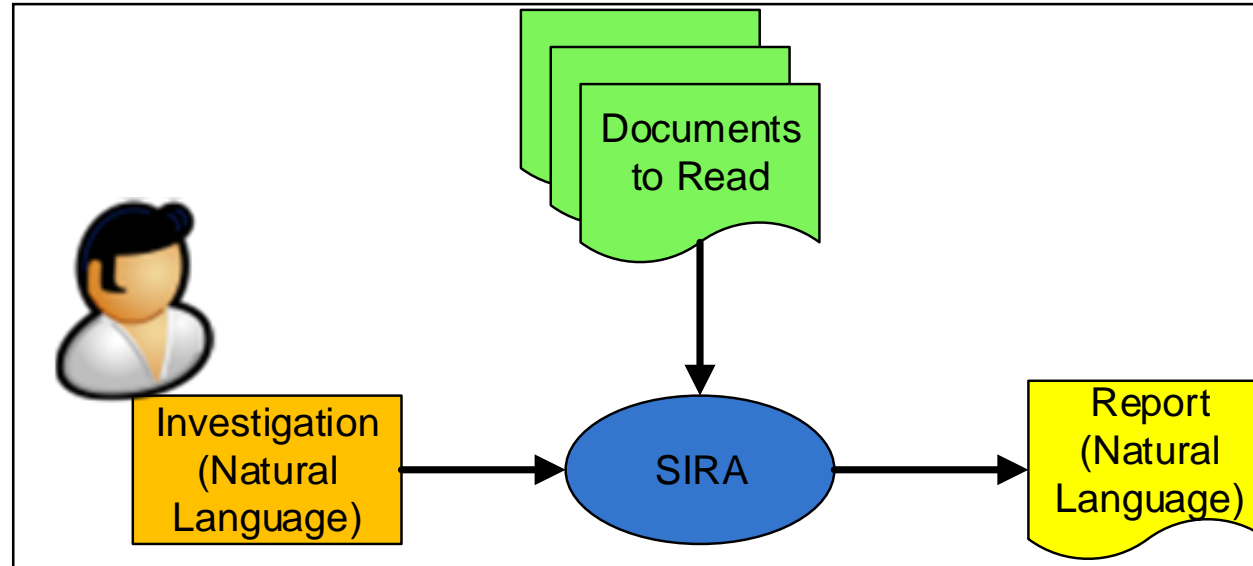# Natural Language Understanding Applied to Research
## (an introduction and demo)

Chuck Rehberg

Chief Scientist, Semantic Insights

# Semantic Insights Research Assistant (SIRA) Technology

- The Semantic Insights Research Assistant (SIRA) Technology was developed to address the need to:

  - Capture and apply enough human knowledge to automate reading Natural Language prose to:

    - Read a potentially a vast number of documents,

    - Look for specific information of interest and

    - Render the findings in useful ways

- Along the way, we found ourselves in the thick of Natural Language Understanding, Reasoning, Machine Learning, Discovery, Natural Language generation, Big Data, Semantic Web, and Graph Databases (e.g. Triplestores).

- For more information visit: http://www.semanticinsights.com

SEMANTIC INSIGHTS™

# View of SIRA by General Users



1. An *Investigation* (i.e. query) is stated in Natural Language
2. The corpus of *documents* to "machine read" are identified
3. The kind of *report* is selected and generated

# About "The Investigation"

- *The Investigation* (actually the "query") is a natural language text (English for now) consisting of one or more sentences.

- These sentences can also contain *open variables*
  - For example, the term "#?#" stands for an "unknown thing"
  - Sample usage: "#?# is traveling to Tunisia on Tuesday."

- Sentences can also contain *closed variables*
  - For example; let #my-equities# refer to the list of stocks I am interested in.
  - Sample usage: "the price of #my-equities# fell."

# About "The Source Documents"

- Must contain some form of natural language text

- Can be encoded in .html, .docx, .rtf, .txt, .pptx., .pdf (non-image), others

- Text can come from databases, email, web pages, documents, and social media feeds.

- Source Documents are machine read line-by-line/ sentence-by-sentence

- Source Documents are not preprocessed

# About "The Report"

- The Kind of Report is selected

    - The format of the report is taken from a template that has been previously created via a SIRA Report Editor.

- SIRA can currently generate reports in these renderings (.txt, .csv, .html, .pdf); others can be easily created

- The most common general report contains verbatim Natural Language text with Bibliography.

# Demos (two basic on-line tools)

- **Research Assistant™** is a Google Chrome™ plugin that
  - Starting from the current web page, and a research description you provide,
  - Research Assistant "reads" the given web page + links + links of links and
  - generates a research report with bibliography
  - Demo
    - "What causes Autism?"
    - http://en.wikipedia.org/wiki/Autism
- **Research Librarian™** is a website that
  - Starting from a selected set of information sources, and a research description you provide,
  - Research Librarian "reads" each document in the selected source + links and
  - Generates a research report with bibliography
  - Demo
    - Example: "What inhibits CYP2D6?"
    - http://52.205.34.110:8082/ResearchAssistant/loginResearchLibrarian
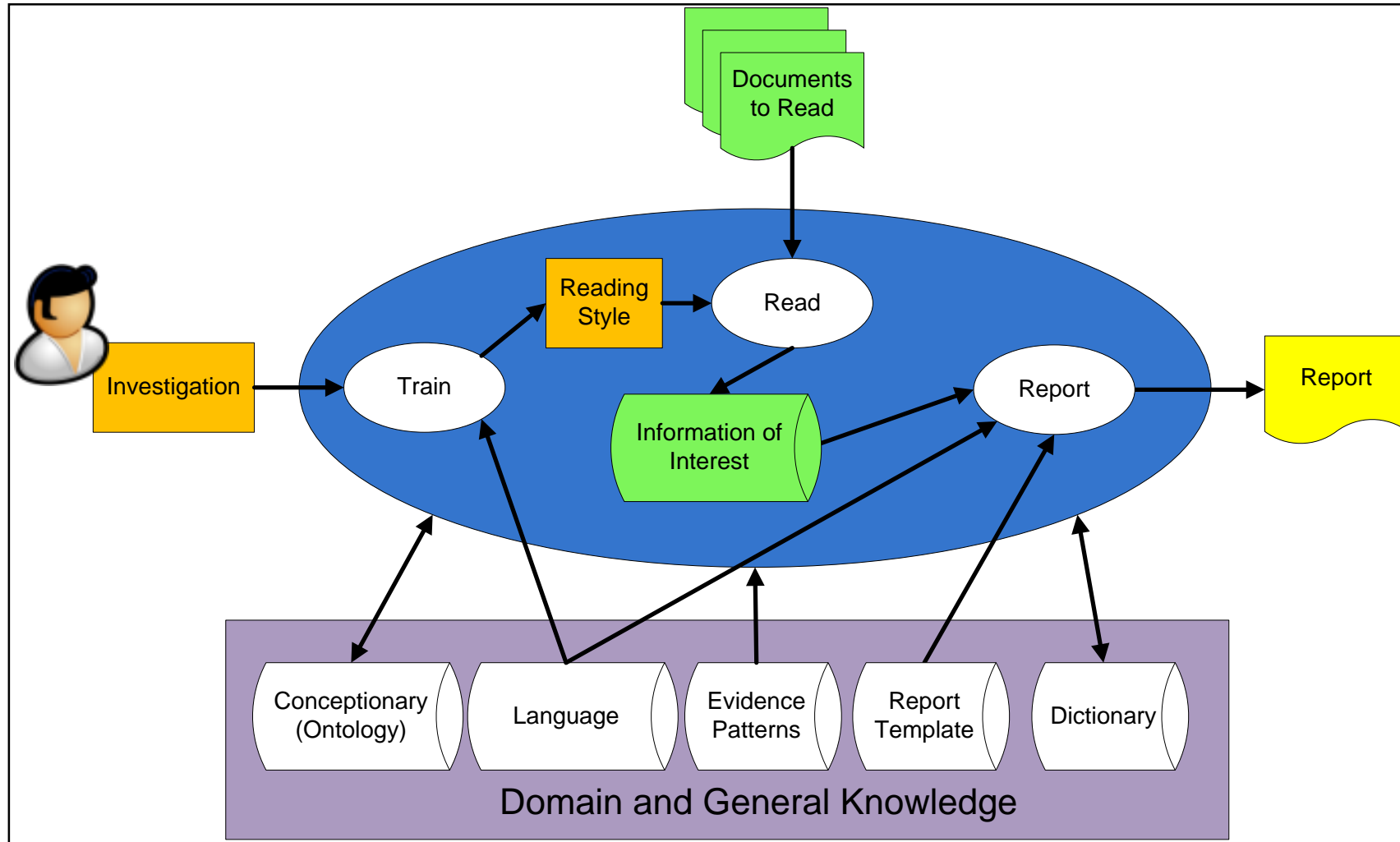
# Another on-line tool with more information

- PriArt is an on-line Research tool that
    - Gather's information based on your research statements and
    - Generates research reports with bibliography and hyperlinks containing information relevant to your research
    - PriArt provides more control of language and presentation than Research Assistant or Research Librarian

- Demo
    - http://www.autism-society.org/what-is/

# A little more detail

# Unsupervised Learning of a given Domain

➢ *Domain Knowledge Bootstrap*

1. Load pre-existing domain-dictionaries
2. Run the *Unrecognized Term Identifier* on a set of domain-specific documents to identify terms missing from the dictionary and looks them up terms in common and domain-specific web dictionaries.
3. Use the *Hypernym Recognizer* to generate hierarchy relationships in the Ontology and add the terms to the dictionary.
4. Use a *set of domain-specific documents as the investigation* to bootstrap the domain dictionary (terms) and Ontology (relationships)

➢ *To learn what else is known about a given topic*

1. Start with a set of statements of interest (i.e. the initial Investigation)
2. Read a corpus of documents
3. Add the relevant sentences in the results to the current investigation and re-read the corpus
4. Continue this process until condition is reached (e.g. no new relevant sentences)

# About reading the Internet...

- PriArt doesn't really read the whole internet. Of course that would take a long time. The main limiting factor is bandwidth.

- However, if you want PriArt to effectively read a large document corpus (like the internet), you can direct PriArt to automatically gather relevant documents by using a keyword search engine.

- PriArt will create and execute a search strategy for you. This includes executing a number of keyword queries, each designed to find the most relevant documents based on examining your statement of investigation.

- You still may get millions of documents to read. So, you can direct PriArt to read a specific number of most relevant documents found according to the search engine ranking.
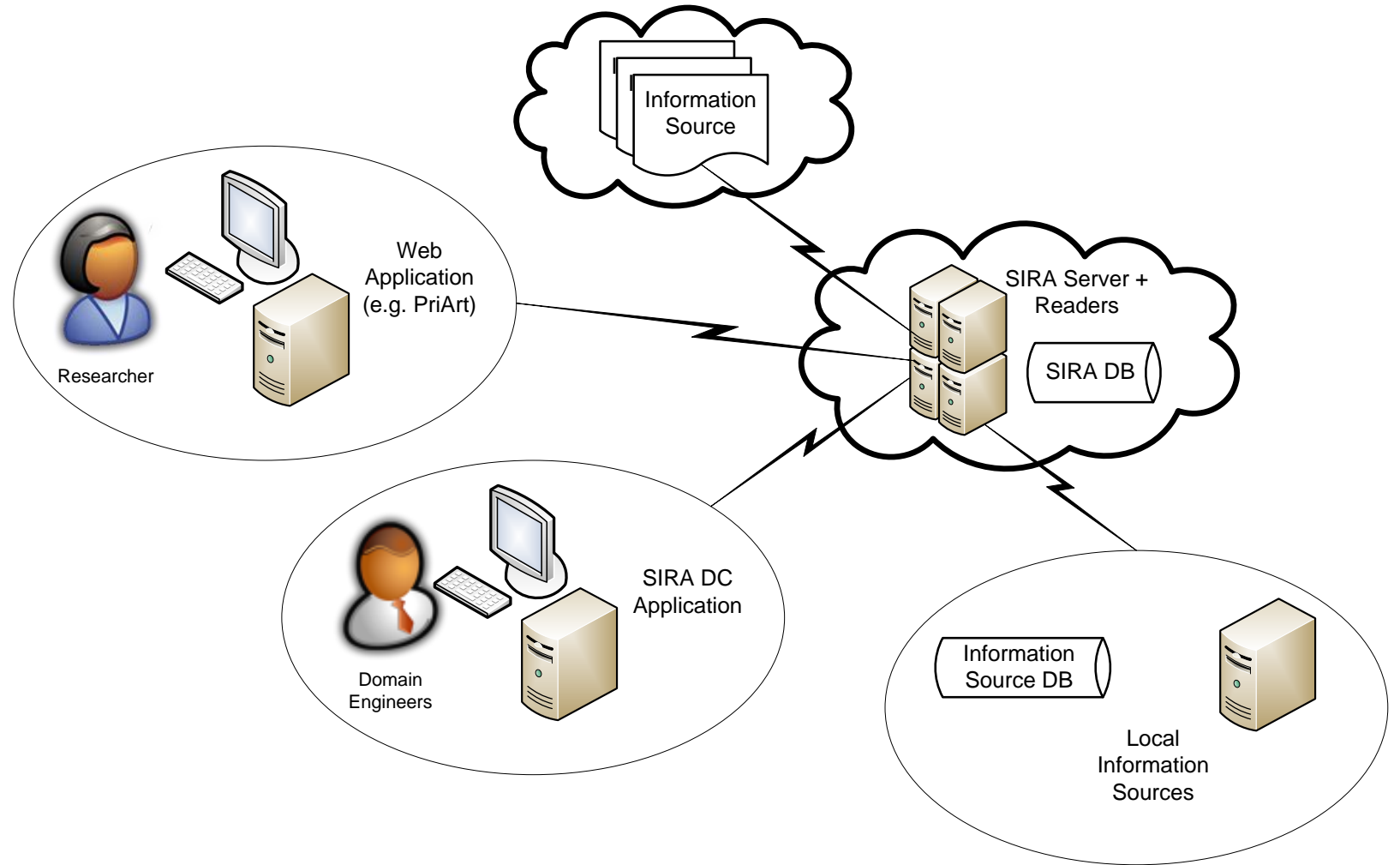
# Learning as you go

- Every time you provide SIRA with an investigation, the investigation is subjected to deep linguistic analysis

- This analysis can add to the Dictionary and create new relationships in the Ontology

- Other forms on experienced-based learnings are under development

# Physical SIRA Configuration

- ➢ SIRA Server
- ➢ SIRA Readers
- ➢ Web Applications
- ➢ SIRA Development Center (thick Client)
- ➢ Multiple Information Sources



Information Source

Researcher
Web Application (e.g. PriArt)

SIRA Server + Readers
SIRA DB

Domain Engineers
SIRA DC Application

Information Source DB
Local Information Sources

# Development Tools Provided with SIRA

1. **PriArt**: A Semantic Research Assistant
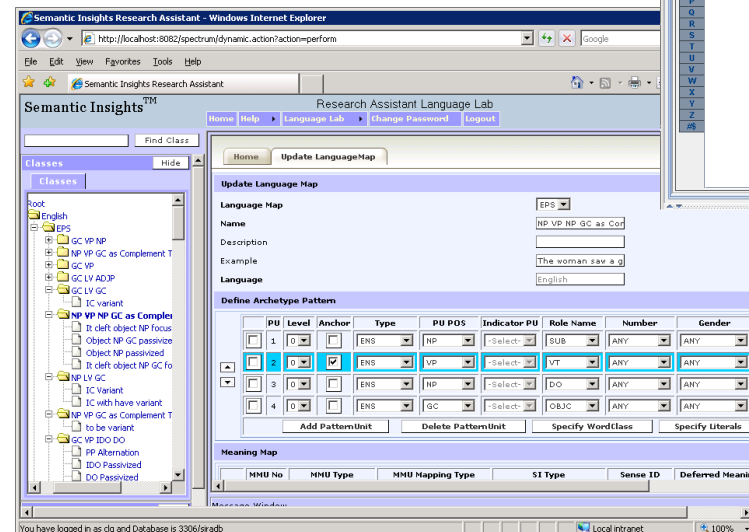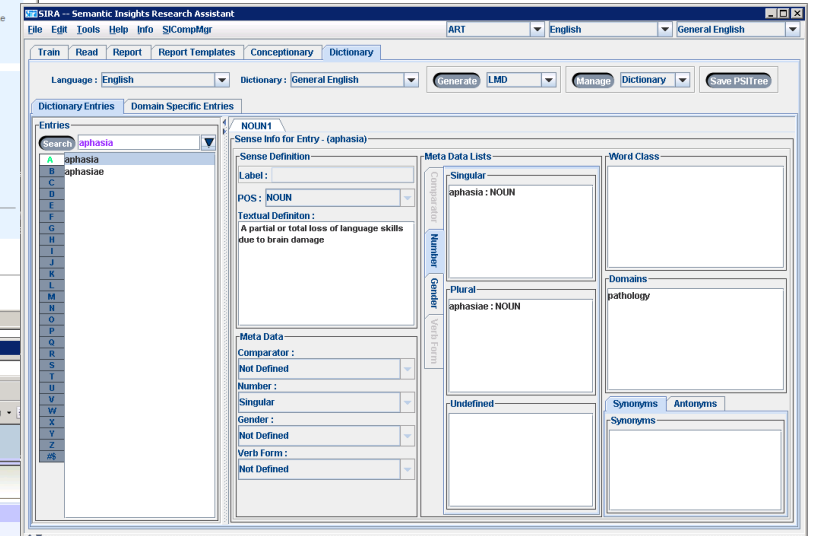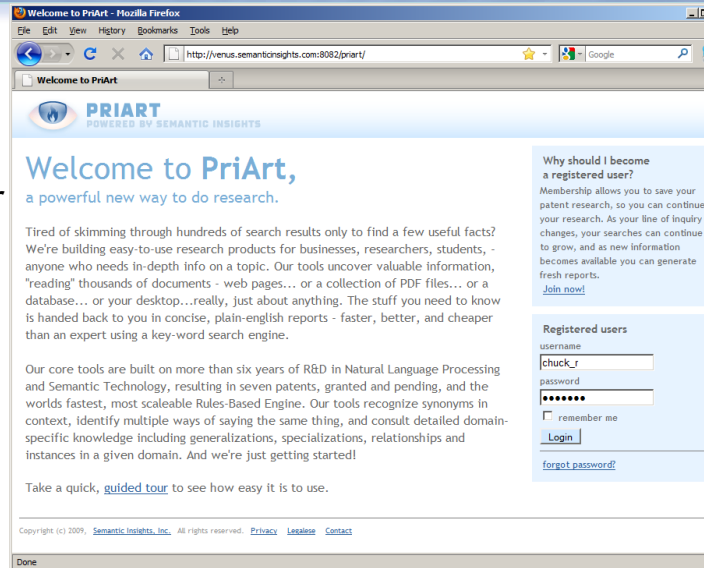   - Browser-based UI
   - Used to demonstrate and test via Browser Interface

2. **SIRA Development Center**
   - Desktop Application
   - Used to develop and manage World Knowledge
     - Ontologies, Dictionaries, Testing and Training
   - Used to develop and manager Report Templates

3. **Language Lab**
   - Browser-based UI
   - Used to define Language and Genre
   - Syntax, Grammar and Meaning Maps

# SEMANTIC INSIGHTS™

A Division of Trigent Software, Inc.