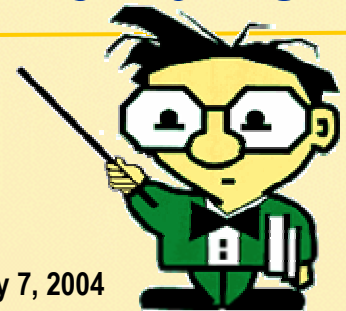# Toward a Knowledge-Based Solution for Information Discovery in Complex and Dynamic Domains

**Eloise Currie and Mary Parmelee**
**SAS Institute, Cary NC**

**7th International Protégé Conference: July 7, 2004**

# About *SAS*: *The Power to Know*®

- *SAS*: The Market Leader in Business Intelligence Software
- Founded: 1976
- World's Largest Privately Held Software Company
- Worldwide Offices: 269
- Worldwide Employees: 9,238
- Worldwide Revenue 2003: $1.34 billion
- Reinvestment in R&D 2003: 26%
- *SAS* Solutions:
  - Used at more than 40,000 sites
  - Used by 96 of the top 100 of the 2003 Fortune Global 500

# The Problem with Information

- From a Global Perspective

  - As the volume of online information grows, information retrieval (IR) has become a major challenge.

- How much is it changing?

  - In 1995, over 90% of corporate documents were in paper form. By 2005, less than 30% will remain in paper form.

  - Within the next 3 years, the world will produce as much data as has been produced since **THE DAWN OF TIME!**

# The Problem with Information

From a SAS Perspective

1. Address Customer Pains:

   – **SAS** customers are satisfied with the quality of our documentation, but they have difficulty locating information.

2. Manage Growth and Complexity

   – How much is it changing?

   – SAS product growth: 1 > 170+

   – SAS user growth: 4 million users worldwide

**7th International Protégé Conference: July 7, 2004**

# Toward a Solution: Preliminary Research

- **Observation**

  - Online information and queries are communicated via natural language, which has two main properties:

    - semantics gives meaning in **context**

    - syntactics give structure and order

    - Yet, most IR systems match only syntactics

- **Objective**

  - Create an IR system that leverages the semantics of natural language.

- **Investigation**

  - Emerging technologies, initiatives and standards: Semantic Web, Ontologies, RDF

  - Consulted IR Experts (UNC-Chapel Hill)

  - Tools: Protégé, Jena Toolkit

**7th International Protégé Conference: July 7, 2004**

BACK GROUND

# Toward a Solution: Progress to Date

2001: Proof of concept project

- **Tiny domain** (subject area): two pages of documentation

- **Rudimentary UI**

- **Deliverable:** development methodology and repeatable process

**Browsable directory tree**
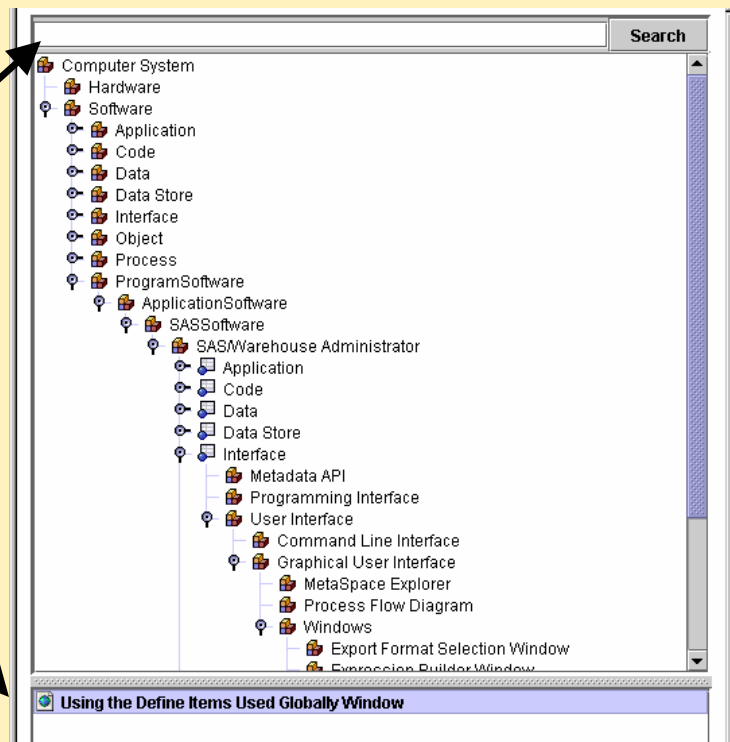
**Resources associated with**

**a node in the tree**



Output Delivery System
- Data Component
- Destination
  - HTML Destination
  - Listing Destination
  - Output Destination
  - Printer Destination
- Output Object

Creating Data Sets from Sir
Identifying Output Objects
ODS SELECT Statement
Selecting Output Objects to
Using Selection and Exclus

**7th International Protégé Conference: July 7, 2004**

BACK
GROUND

# Toward a Solution: Progress to Date

2002: Pilot Project

- **Larger domain**: a complete product user's guide
- **More robust UI**
- **Deliverable**: A functioning mini system

**Added search functionality**

**Search results pane**

BACK GROUND

# Toward a Solution: Progress to Date

## February 2004: Prototype Project

- **Large domain**: Includes several products
- **Intelligence layer:** Enables advanced search and reasoning capability
- **Advanced UI**: Delivers information in context while minimizing the complexity surfaced to the user
- **Deliverable:** Fully Functional Prototype

BACK GROUND

# Knowledge Base System Architecture

Content Development ➡ Intelligence Layer ➡ Contextual Information Delivery

## Create Resource Set (Domain)

Authored Content

Resource Repository

Static Resource Assembly

Hardcopy

Compact Disc

HTML

## Classify Resources

Knowledge Model (Ontologies)

## Process Knowledge Model

Knowledge Repository

Retrieves Resources

Unstructured Content

Process Controller

Relay Queries and Retrieves Results

Full Text Engine

Inference Engine

Queries KR and Retrieves Metadata

## Deliver

Delivers Resources

Client Queries

KB Portal

Client Queries

Happier Client !

7th International Protégé Conference: July 7, 2004

# System Development Process

**We begin with a document collection ( a "resource set")..**

1. Use SAS® Text Miner to create a hierarchy of resource clusters

2. Use a custom Protégé plugin to generate a **Domain ontology** that categorizes resources based on content

3. Use a custom Protégé plugin to extract resource information and generate a **Resource Manager ontology**

4. Merge Domain and Resource ontologies into a **Master ontology**

5. Use a custom Protégé plugin to map Resource instances to Domain instance slots

6. Use Protégé to develop the merged ontology into a production Master ontology

7. Use a custom Protégé plugin to reverse map Domain Instances to Resource Instance Slots

**7th International Protégé Conference: July 7, 2004**

# Step 1: Use SAS® Text Miner to Create a Hierarchy of Resource Clusters

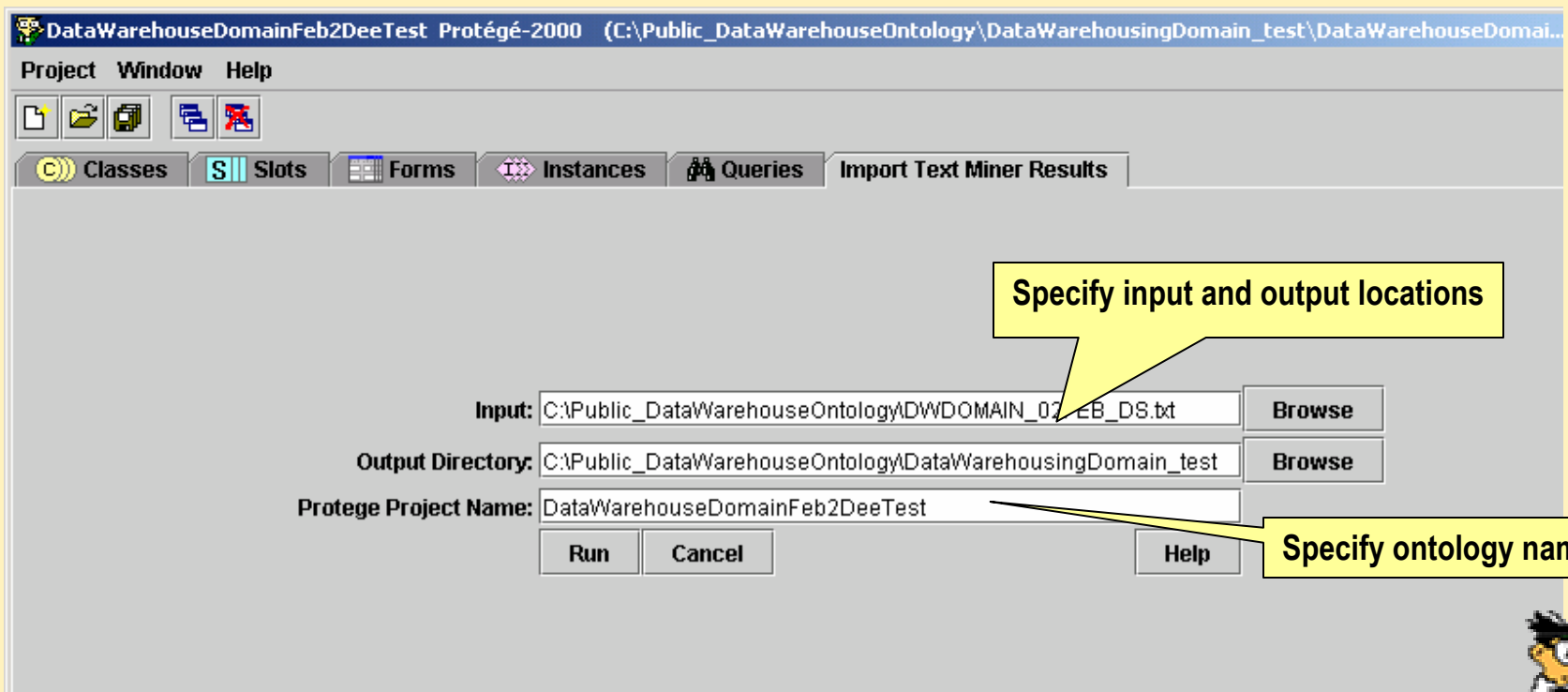| Steps | Result |
|---|---|
| **Preprocessing** | Create a SAS data set from the document collection. |
| **Text parsing** | Generate quantitative representation of the content |
| **Transformation** | Consolidate quantitative representation |
| **Document analysis** | Cluster documents by concept |

# Step 1:Hierarchical Output

Resource Set

"window" 1402 docs

36

"metadata" 1402 docs

29

"data" 326 docs

3

"display" 1076 docs

15

"property" 159 docs

7

"define" 821 docs

22

"table" 599 docs

1

"option" 477 docs

2

"method" 292 docs

6

"interface" 529 docs

8

"permission" 329 docs

4

5

"Hybrid" cluster, built up from root clusters

"Root" cluster

"SAS 9.1" 200 docs

# Step 2: Use a Custom Protégé Plugin to Generate a Domain Ontology in Protégé

- Saves Text Miner hierarchy as a Protégé ontology
- Preserves URI's of documents in a Resource ID slot

# Step 2: Domain Ontology in Protégé

- Each direct instance is a Text Miner document cluster
- Generic instance and class names will be edited



Direct Instances have generic names generated by Protégé

Classes have generic names generated by Text Miner

Significant cluster terms identified by Text Miner

Resource IDs of those resources assigned to this direct instance of the Domain ontology

**7th International Protégé Conference: July 7, 2004**

# Step 3: Use a Custom Protégé Plugin to Generate a Resource Manager Ontology in Protégé

- Extracts resource information
- Preserves URI's of documents in a Resource ID slot
- Saves resource information in a Protégé ontology

**Resource Manager**

Specify path and name of Resource Manager ontology

Select a Resource Manager Protege Project

File Name: C:\Public_DataWarehouseOntology_PRODUCTION\DW.pprj

Browse    New

Register    Retrieve    Help

# Step 3: Resource Manager Ontology

Each Direct Instance is a document in the Resource class



Resource Manager Resource class

Instances of the Resource class

**7th International Protégé Conference: July 7, 2004**

# Step 4: Merge the Domain and Resource Ontologies into a Master Ontology

# Step 5: Use a Custom Protégé Plugin to Map Resource Instances to Domain Instance Slots

- Maps by Resource ID
- Populates resource instance slot

**Link Resources with Categories**

**Link Resources with Categories**

Merged Ontology: elopmentDocs\MAproduction\MarketingAutomation_1.0Master_Protege1.8.ppn]  Browse

□ Map Resources to Categories

Map Categories to Resources

Run     Help

**Specify path of Master Ontology**

**Select "Map Resources to Categories"**

# Step 5: Resource Instances
# Mapped to Domain Instance Slots



Domain instance

Domain class

Resource Manager Description slot facilitates ontology refinement

Resource Manager instances are mapped to Domain instance slots
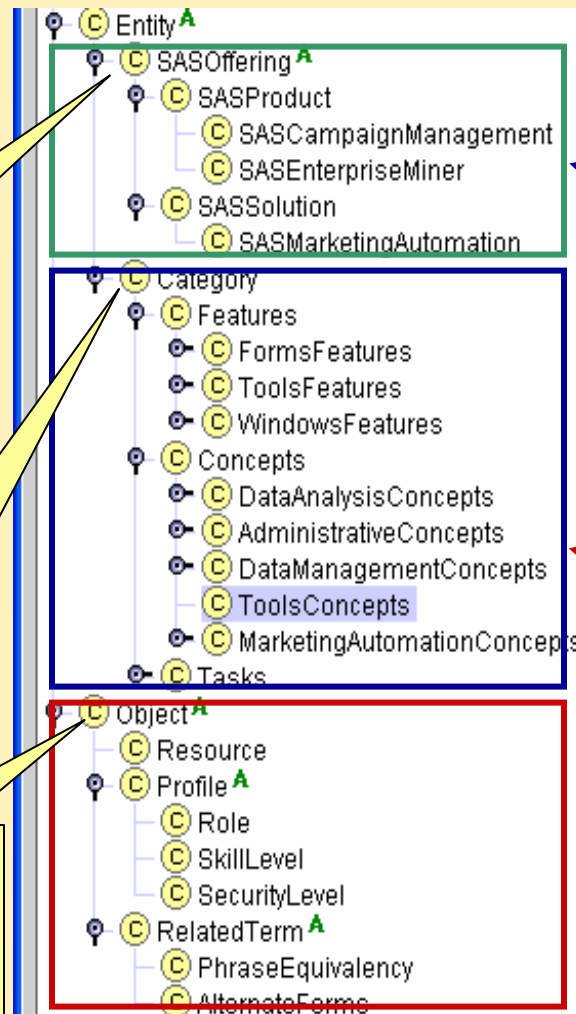
7th International Protégé Conference: July 7, 2004

# Step 6: Refine and Expand the Merged Ontology into a Production Master Ontology

- Add Related Terms to enable fuzzy matching of misspellings, synonymous phrases and alternative word forms

- Refine the Domain Hierarchy
  - Add the SAS Offerings Model
    - SAS products and solutions
  - Develop full ontology from Text Miner hierarchy
  - Add metaclasses to assign URIs at the class level

# Step 6: Production Master Ontology



**SAS Offerings Model**

**Domain Ontology Classes with edited, meaningful names**

**Resource Manager ontology classes**

# Step 6: Add Metaclasses



Add Metaclasses to assign URIs at the class level

**7th International Protégé Conference: July 7, 2004**

# Step 7: Use a Custom Protégé Plugin to Map Domain Instances to Resource Instance Slots

**Reverse mapping of Resource to Domain instance (Step 5)**



Specify path of existing Master Ontology

Select "Map Categories to Resources"

# Step 7: Domain Instances are Mapped to Resource Instance Slots

# Ontologies Define the Intelligence Layer



Knowledge Base Resources optimized for:

Browse          Search          Preferences

Intelligence Layer

Teradata

MySQL    DB2

XML    PDF
Word    HTML

Databases          Multimedia          Distributed Systems          Unstructured Resources

# Knowledge Base Prototype

## Delivers information in context using

- Browsable categories

- Categorized search results

- Hover text descriptions

- Category bread crumb trails

- Category and full text search

- Fuzzy matching

**7th International Protégé Conference: July 7, 2004**

# Browse View: Browsable Directory and Hover Text Contextual Cues

# Search Results View: Search Expansion Fuzzy Match Synonymous Phrase

enter phrase "grouping reports" and push the search button

**SAS Marketing Automation 3.1 Knowledge Base**

Browse   Search

Search: grouping reports          all words ▼    Search
         Show only results that match my preferences       Search tip

Reports > Portfolios

**Running Portfolios**
All the reports stored in a report portfolio ca...        For example, you might have a portfolio of...
SAS Marketing Automation 3.1 > SAS Campaign Management 7.4 > Tasks > Data Management > Reports > Portfolios

Search for more results in *Reports* >>>

Contact
Cop

system matches to phrase equivalency "portfolios"

breadcrumb trail to the root category provides context

Optional full text search filtered by the current category

# Search Results View: Contextual Cues
# Grouped Results of Full Text Search

**Reports > Portfolios**

**Running Portfolios**
All the reports stored in a report portfolio can be run together. You might want to do this to run a number of reports several times. For example, you might have a p
SAS Marketing Automation 3.1 > SAS Campaign Management 7.4 > Tasks > Data Management > Reports > Portfolios

**Word Processing > Save**

**Saving Reports**
A report can be saved on the system for later use. Reports are usually saved in portfolios. These are used to group together several reports in a logical way. For ex
SAS Marketing Automation 3.1 > SAS Campaign Management 7.4 > Tasks > Data Management > Word Processing > Save

**Report > Campaign Report**

**Saving Reports**
A report can be saved on the system for later use. Reports are usually saved in portfolios. These are used to group together several reports in a logical way. For ex
SAS Marketing Automation 3.1 > SAS Campaign Management 7.4 > Features > Windows > Dialog Boxes > Data Management > Report > Campaign Report

**Renaming and Deleting Reports**
Portfolios and reports which you own can both be renamed and deleted. Follow the steps below: Click the Open icon. The Report and Portfolio Management dialog
SAS Marketing Automation 3.1 > SAS Campaign Management 7.4 > Features > Windows > Dialog Boxes > Data Management > Report > Campaign Report

**Retrieving and Editing Reports**
Saved reports can be opened for viewing, copying to the clipboard, and printing. You might also want to retrieve saved reports to run them against an updated data
SAS Marketing Automation 3.1 > SAS Campaign Management 7.4 > Features > Windows > Dialog Boxes > Data Management > Report > Campaign Report

**Reports > Edit Reports**

**system returns all search results "grouped by" category**

**7th International Protégé Conference: July 7, 2004**

# Our Vision for an Integrated Solution:
# The Value of XML-Based Modular Content

- **Why XML?**
  - Accepted standard-already used by a critical mass
  - Facilitates platform independent interoperability
  - Provides a content development framework that supports modular writing

- **Why Modular Writing?**
  - Reusability
    - Controls work redundancy
    - Reduces semantic heterogeneity
      - using the same terminology to mean different things
      - using different terminology to mean the same thing
  - Facilitates content classification: "about" one thing
  - Enables advanced information retrieval and delivery techniques
    - dynamic assembly of complex resources that are relevant to a user's current context

**7th International Protégé Conference: July 7, 2004**

# Knowledge Base System Architecture

**Content Development** ➡ **Intelligence Layer** ➡ **Dynamic Information Delivery**
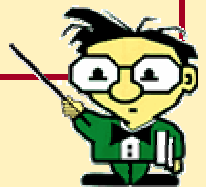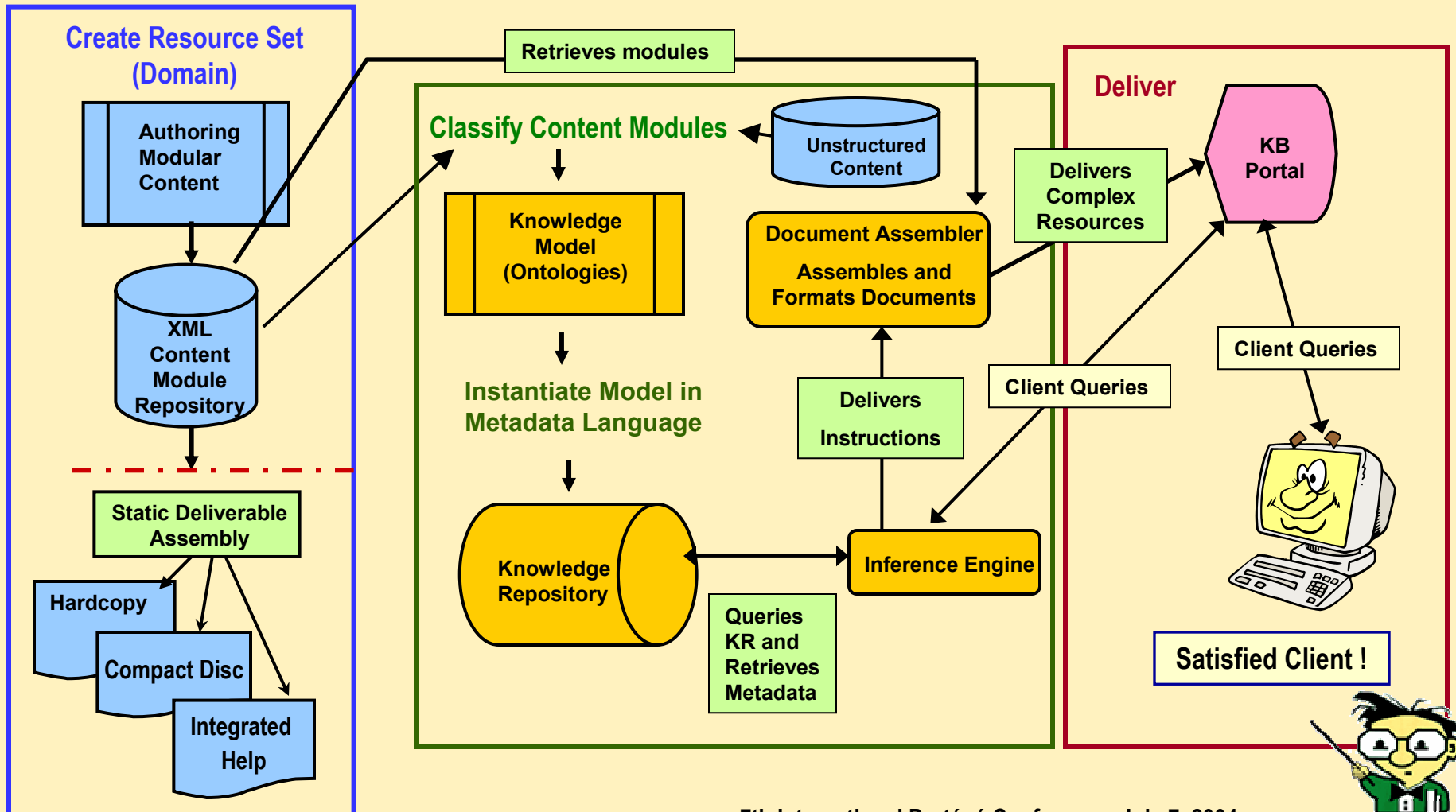
**Create Resource Set (Domain)**

Authored Content

Resource Repository

Static Resource Assembly

Hardcopy

Compact Disc

HTML

**Retrieves Resources**

**Classify Content Modules**

Unstructured Content

Knowledge Model (Ontologies)

**Instantiate Model in Metadata Language**

Knowledge Repository

Queries KR and Retrieves Metadata

Relay Queries and Retrieves Results

Process Controller

Full Text Engine

Inference Engine

**Deliver**

Delivers Resources

Client Queries

KB Portal

Client Queries

Happier Client !

# One Vision for an Integrated Solution

7th International Protégé Conference: July 7, 2004

# Questions/Comments?

**Thanks to Contributors: Dee Stribling and Chris Goolsby**

**7th International Protégé Conference: July 7, 2004**