



Technology Labs

Research Directions in Enterprise Knowledge Management

Rayid Ghani

Enterprise Search Today

[Advanced Search](#)

Accenture Technology Labs

SIT - **Technology** > **Accenture Technology Labs** **Accenture Technology Labs** Asset Monitoring & Optimization Practice **Accenture Technology Labs** Asset Monitoring & Optimization Practice

<https://kx.accenture.com/Organizations/Pages/AccentureTechnologyLabs.aspx>

Content Current Date - 6/14/2005

Last Index Date - 2/13/2007

[\[View duplicates\]](#)

Accenture Technology Labs Patent Profiles

Accenture Technology Labs Patent Profiles **Accenture Technology Labs** Patent Profiles **Accenture Technology Labs**

<https://kx.accenture.com/TechMeth/Pages/LabsPatent-AllProfiles.aspx>

Content Current Date - 4/7/2005

Last Index Date - 12/29/2006

[\[View duplicates\]](#)

Accenture Technology Labs: Contacts - Meet the **Labs**

Accenture Technology Labs: **Accenture Technology Labs** falls under Systems Integration & **Technology** under the direction of the director of the **Accenture Technology Labs**, responsible for setting strategic direction and managing

<https://kx.accenture.com/Organizations/Pages/AccentureTechnologyLabs-Contacts.aspx>

Content Current Date - 5/24/2005

Last Index Date - 7/6/2007

Accenture Technology Labs - Sell & Deliver

paper describing **Accenture Technology Labs**' "Command and Control (C2) Interactive Wall" and how it Short Video Clips of **Technology** Prototypes and Solutions from **Accenture Technology Labs** showcasing the **Accenture Technology Labs** latest prototypes and solutions.

<https://kx.accenture.com/Organizations/Pages/AccentureTechnologyLabs-SellDeliver.aspx>

Content Current Date - 6/13/2005

Last Index Date - 12/29/2006

Accenture Technology Labs India

SIT - **Technology** > **Accenture Technology Labs** **Accenture Technology Labs** India in Collaboration - **Accenture Technology Labs** creates a Virtual Corridor between **Accenture** offices i

<https://kx.accenture.com/Organizations/Pages/AccentureTechnologyLabs-Bangalore.aspx>

Content Current Date - 3/20/2006

Last Index Date - 6/28/2007

Accenture Technology Labs - About the Organization

SIT - **Technology** > **Accenture Technology Labs** **Accenture Technology Labs**, the **technology** research and development (R&D) organization within Acce **Accenture Technology Labs** falls under Systems Integration & **Technology** under the direction of the

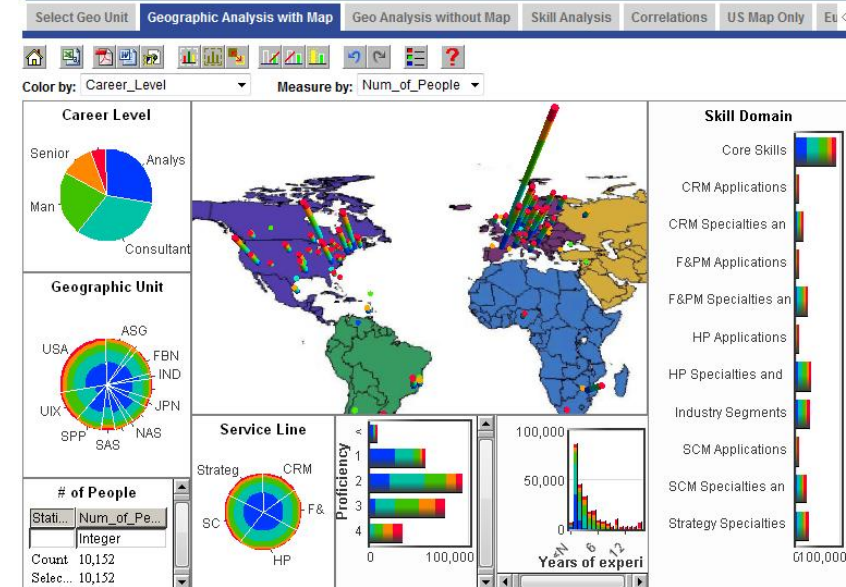
<https://kx.accenture.com/Organizations/Pages/AccentureTechnologyLabs-About.aspx>

Enterprise IR is different from Web IR

- Scale
- Hyperlinks
- Static ranking
- Structured data
- Version Control
- Expectations (THE document vs a document)
- Redundancy
- Security/access control
- Users (roles/needs/context)
- Tasks (Business processes)

Enterprise Tasks involve specialized business processes and roles

- Proposal Writing
- Marketing
- Selling/Estimation
- Risk management
- Requirements Analysis
- Software Testing
- Training/Learning
- Outsourcing
- Vendor evaluation
- Procurement
- Business Intelligence
- Project Staffing
- Recruiting



Collaborative Document Development Process

- Identify requirements
- Identify collaborators/team
- Develop high level themes
- Create outline with sections
- Assign sections to individuals
- Support individuals finding content to complete sections
- Support checkpoints and alerts
- Consolidate drafts of sections and provide support to ensure consistency, compliance, theme integration, tracking changes
- Consolidate to produce final document
- Review for consistency, compliance, integration
- finalize

http://betterthinkfast.techlabs.accenture.com/saable/mantlepreview.html?search=ria - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites

Links my del.icio.us post to del.icio.us NYT Meetingplace

Address http://betterthinkfast.techlabs.accenture.com/saable/mantlepreview.html?search=ria

techlabs Highlight Alliance Partners Highlight Acronyms Highlight Clients Client Directors Similar in KX Select Search SABLE

Go

Search sap utilities Go

accenture Technology Labs saable

Text Results Rich View Histogram View

Solutions for the Natural Gas Industry Based

SAP Gas Presentation

Proposal Excerpt: CMS Energy CEA SAP

Global CCS Meeting - Prague - June 7-8, 2004

Waldorf Connection 03/2004

SAP Utility Customer E- Services

Global CCS Meeting - Amsterdam,

Proposal Excerpt: Enbridge SAP

SAP Utilities Recommendation Letter

Accenture Spain SAP Utilities Solution

Alliance Tea Overview

Download

► All Documents (11018)

► customer e-services, sap ag, sap is

► alliance team, sap crm, sap utilities i

► call center, mysap crm, sap crm (5

► accenture, , (5)

► client profile, crm, project status (5

► accenture sap, bc hydro (4)

► accenture lp, , sap ccs (4)

► cms, organizational change manag

► hydro one's, supply chain (4)

► key messages, sap global, utilities s

► sap utilities discussion group (3)

► energy data management, sap is-u

► internal accenture, sales effectiveness

► award, excellence delivering value (

► master data, sap ag, speaker name

Direct Access from Call Center to Back end Business Objects and Processes

Direct Access to:
✓ Bills/Invoices
✓ Budget Billing Plans
✓ Security Deposits
✓ Devices
✓ Dunning
✓ Open Items
...

SAP CRM System

Business Objects, such as Business Partners, Activities, Contracts, Products, Opportunities, etc.

Electronic Business (Internet)

Tele-Business

Mobile Components

SAP IS-U/ CCS

Back End

Front End

THE BEST-RUN E-BUSINESSES RUN SAP

SAP

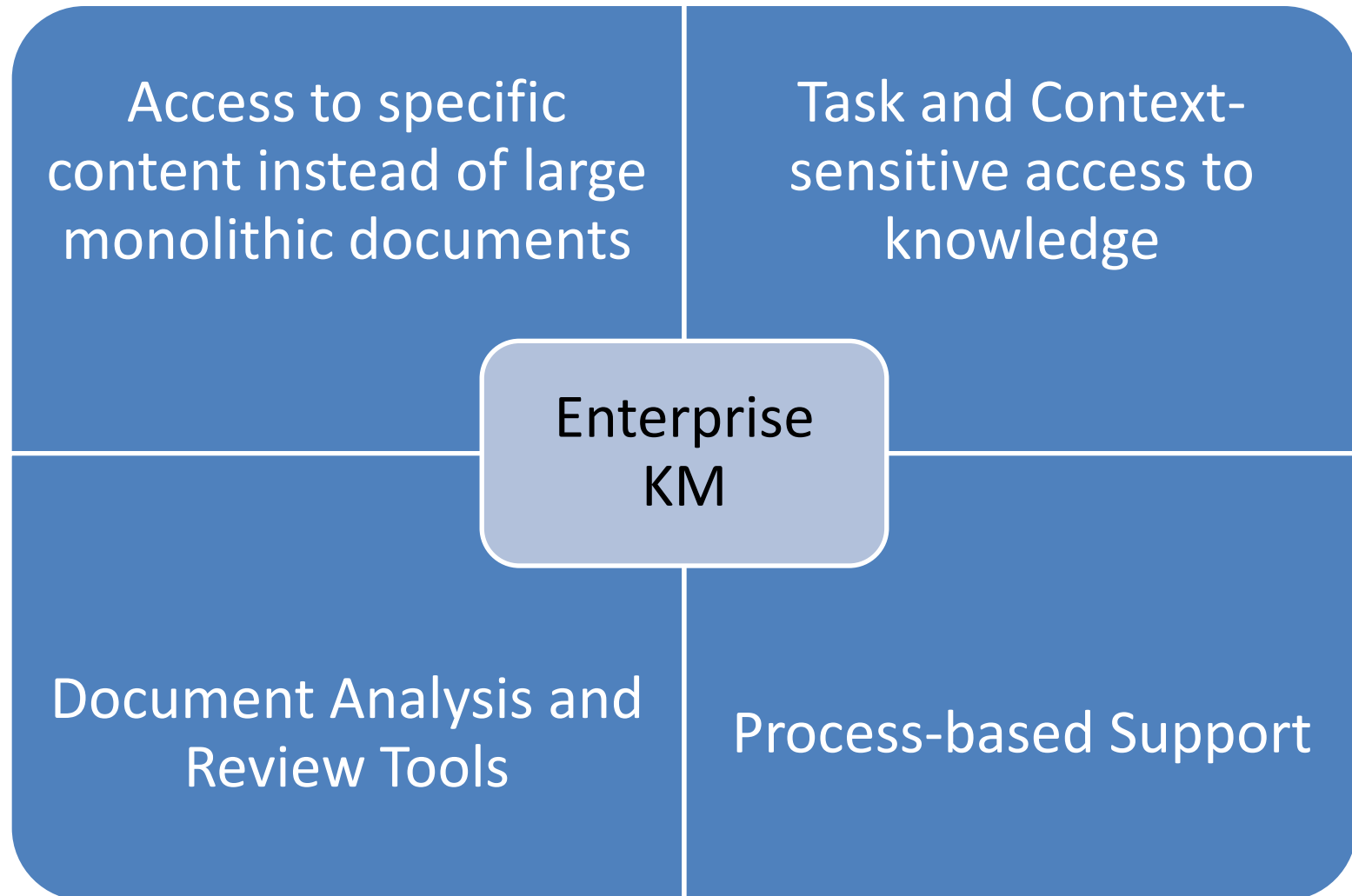
Global CCS Meeting - Amsterdam, June 24 & 25, 2002

Slide Page: 1

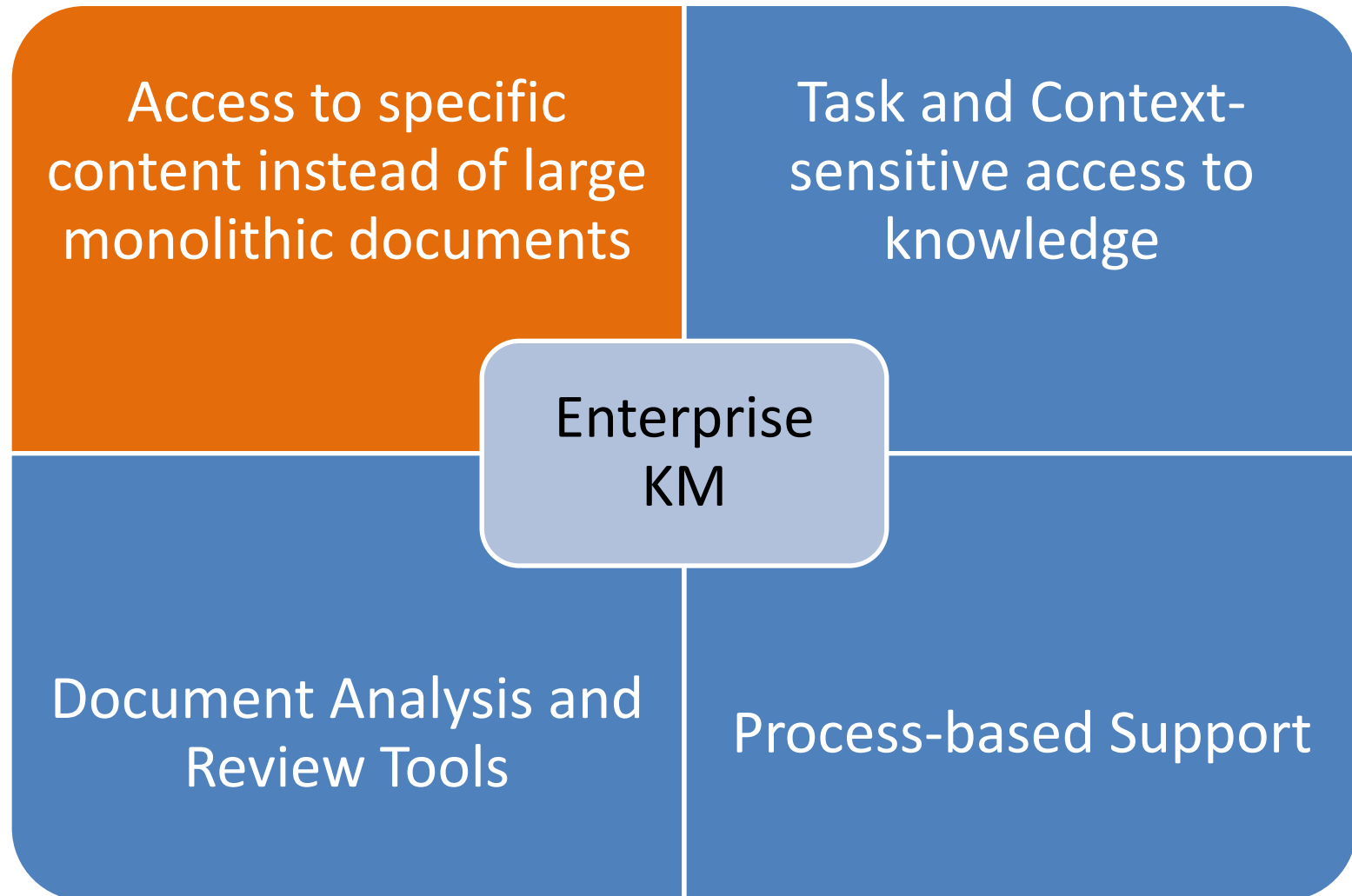
Done

Internet

Areas of Improvement



Areas of Improvement



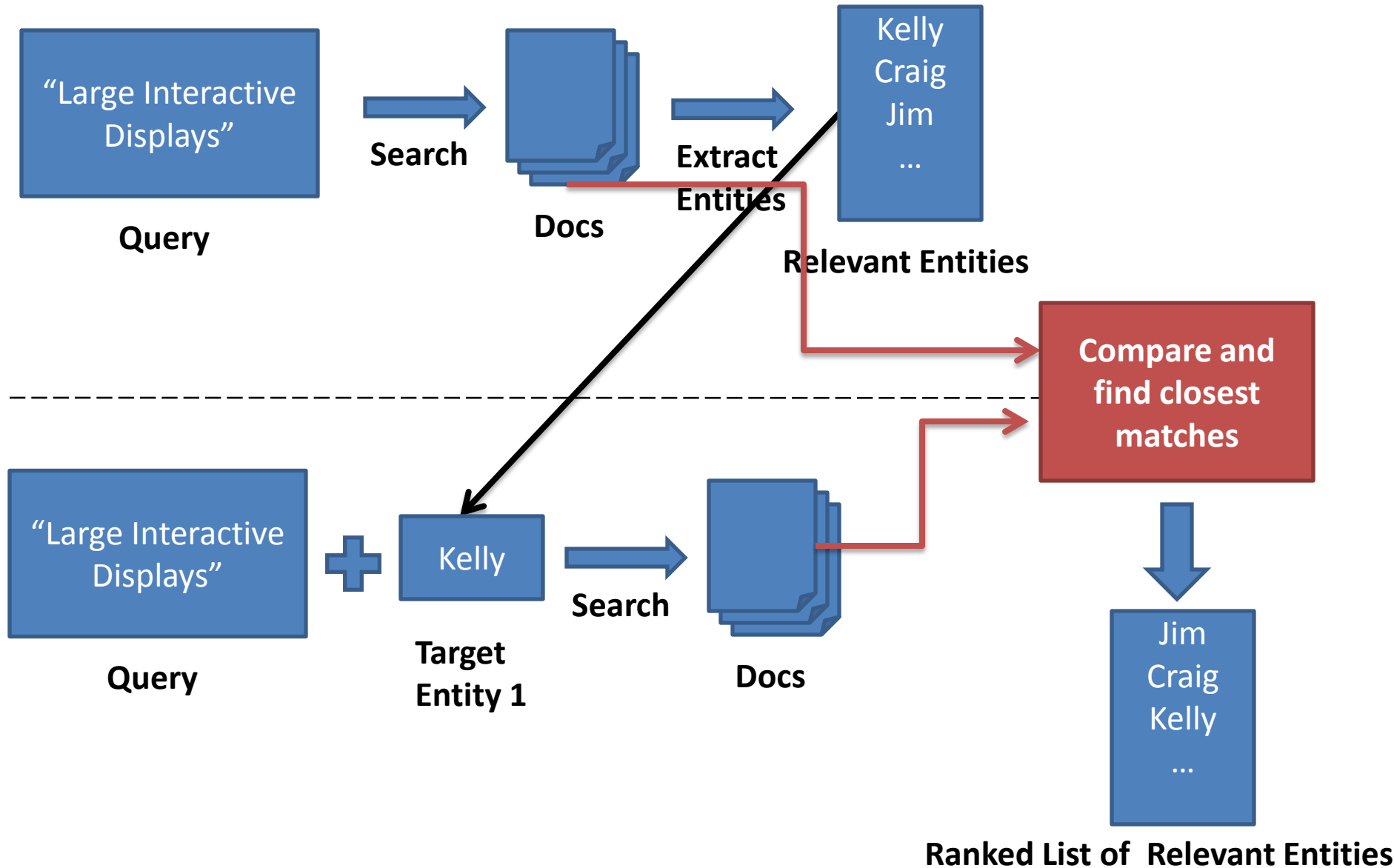
From monolithic documents to reusable pieces of information

- Entities (experts for example - [demo](#))
- Document Chunking ([demo](#))
- Images ([demo](#))

Motivating business task

- Task: Developing a consulting project proposal
- Information Needs
 - *People* that have worked on similar proposals
 - *Clients* for which Accenture has done similar work
 - *Vendors* that Accenture has used for similar work
 - *Alliances* Accenture has with companies that we can partner with
 - *People* available to work on the project
 - ...

Entity Retrieval & Ranking



Online experiments

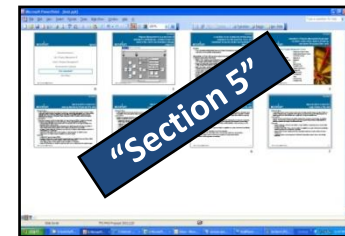
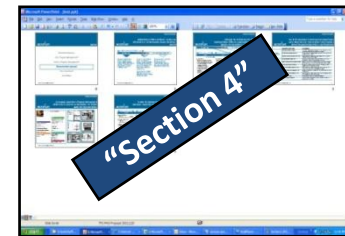
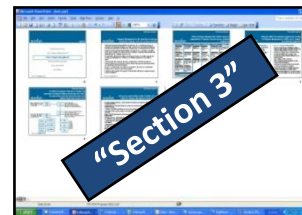
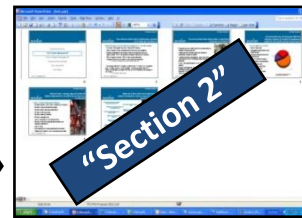
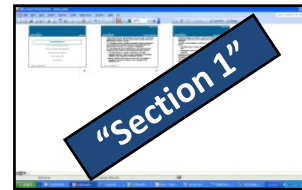
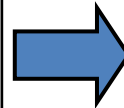
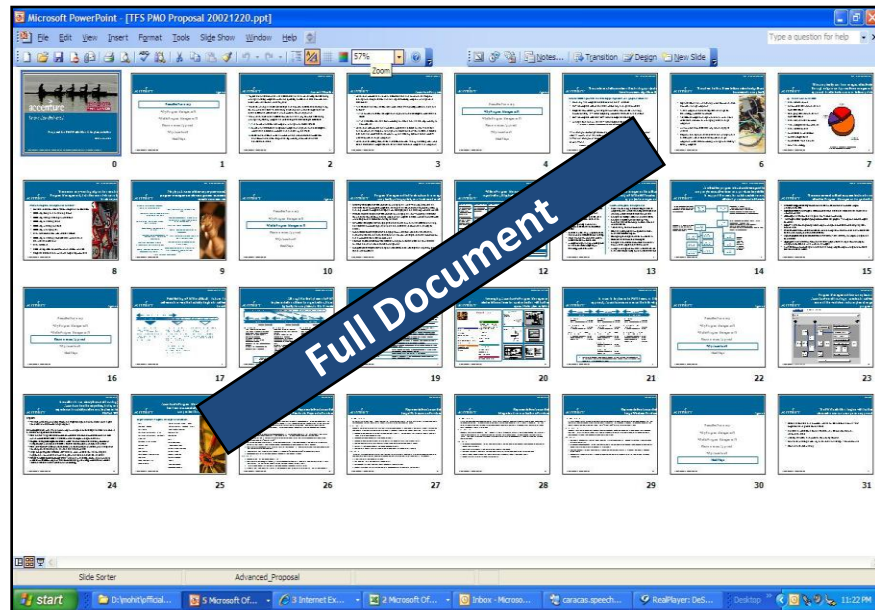
- Expert Search as an extension to enterprise search



- IM bot to find and interact with experts



Document Chunking



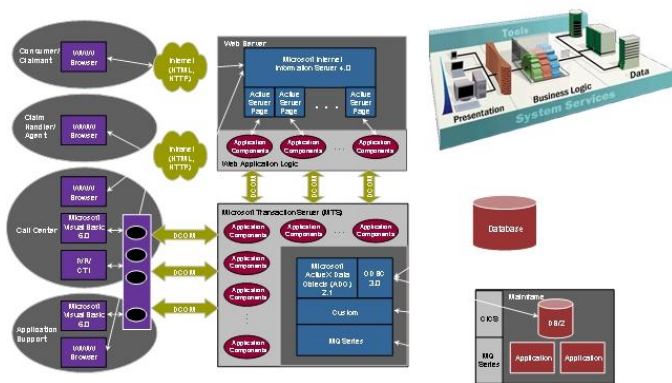
- Approach:

- Determine boundary of sections by segmenting documents
 - Use clustering to find similar sections
 - Classify sections into predetermined section classes

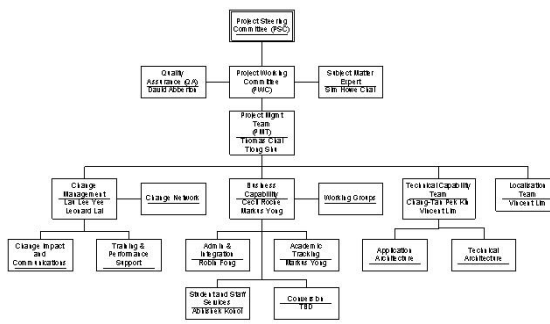
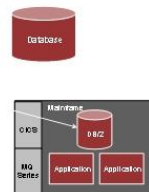
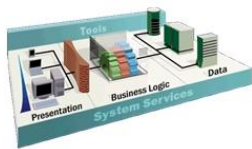
- Output:

- List of labeled individual sections from the document

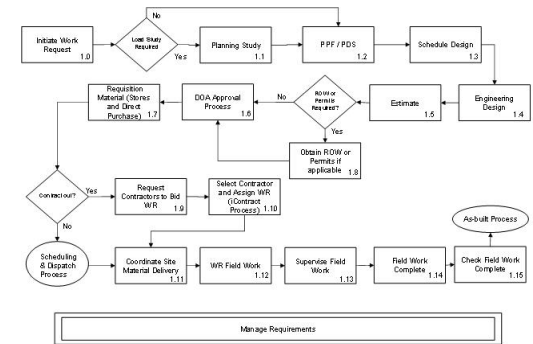
Enterprise Image Classification and Retrieval



Architecture Diagram



Org Chart

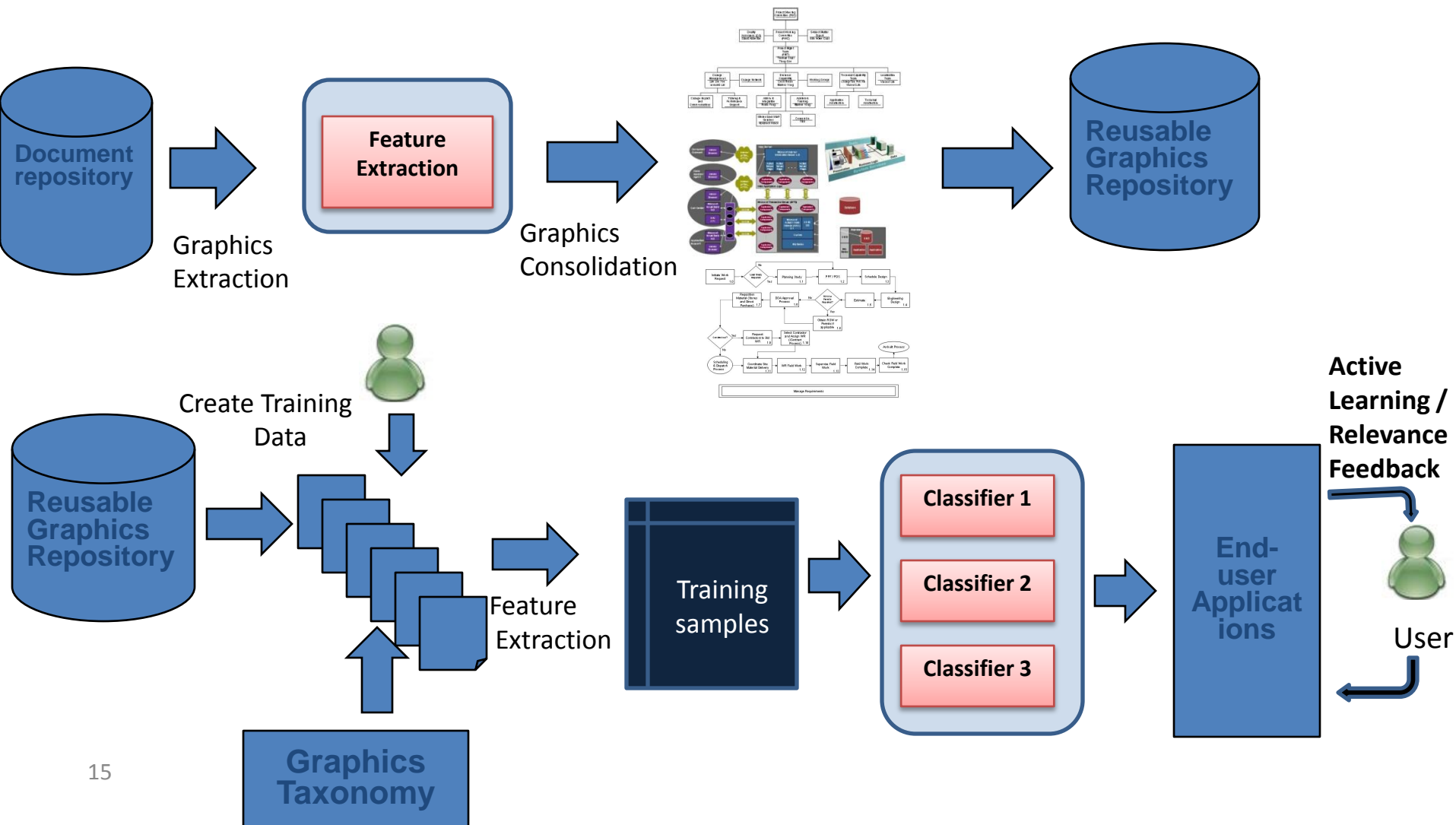


Process Flow

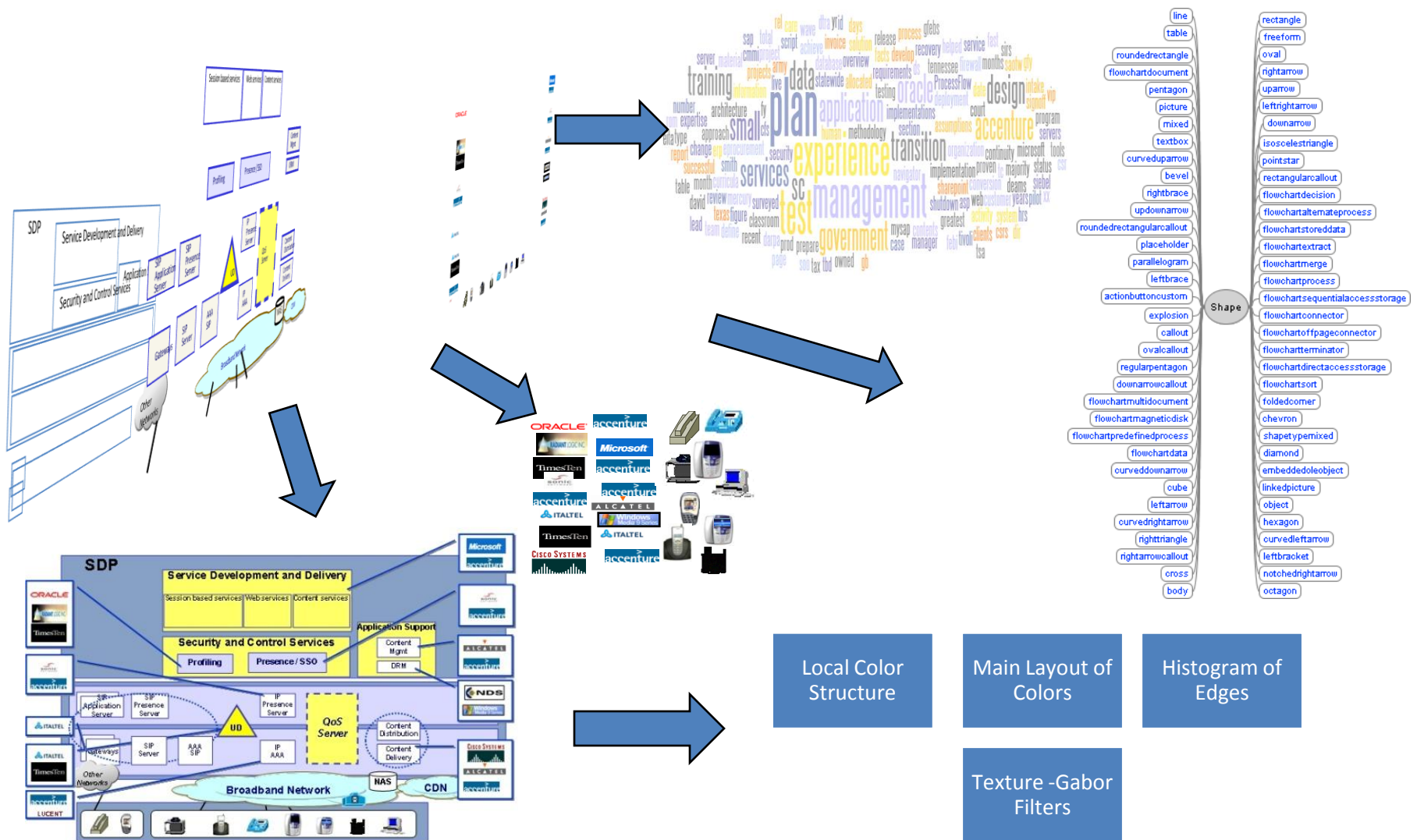
[Demo](#)

- Approach
 - Image extraction from office Documents
 - Supervised learning (6-10 categories) with structural, visual, and text features of the image
 - Index image using words and image category

Overview of Approach



Feature Extraction

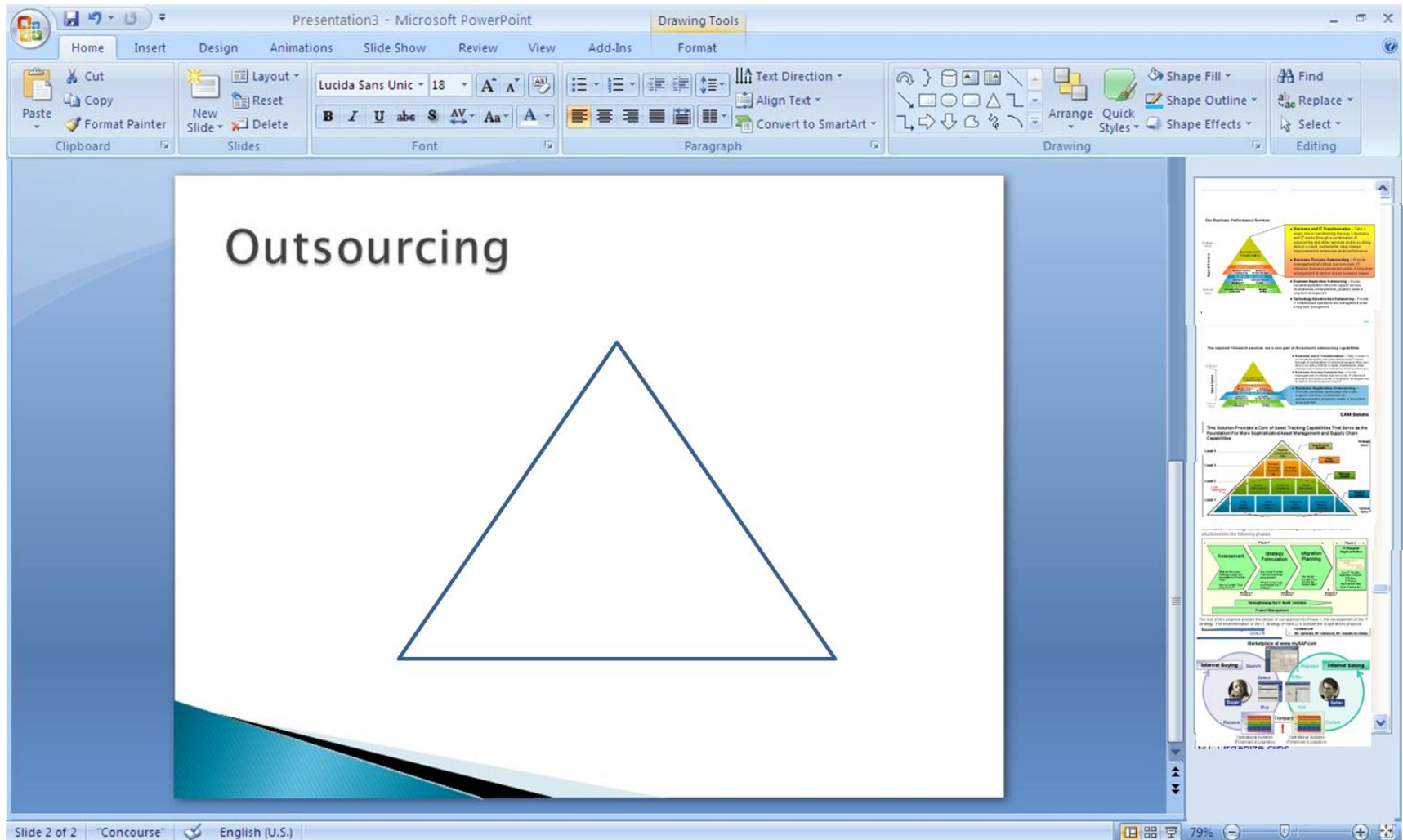


Task specific application

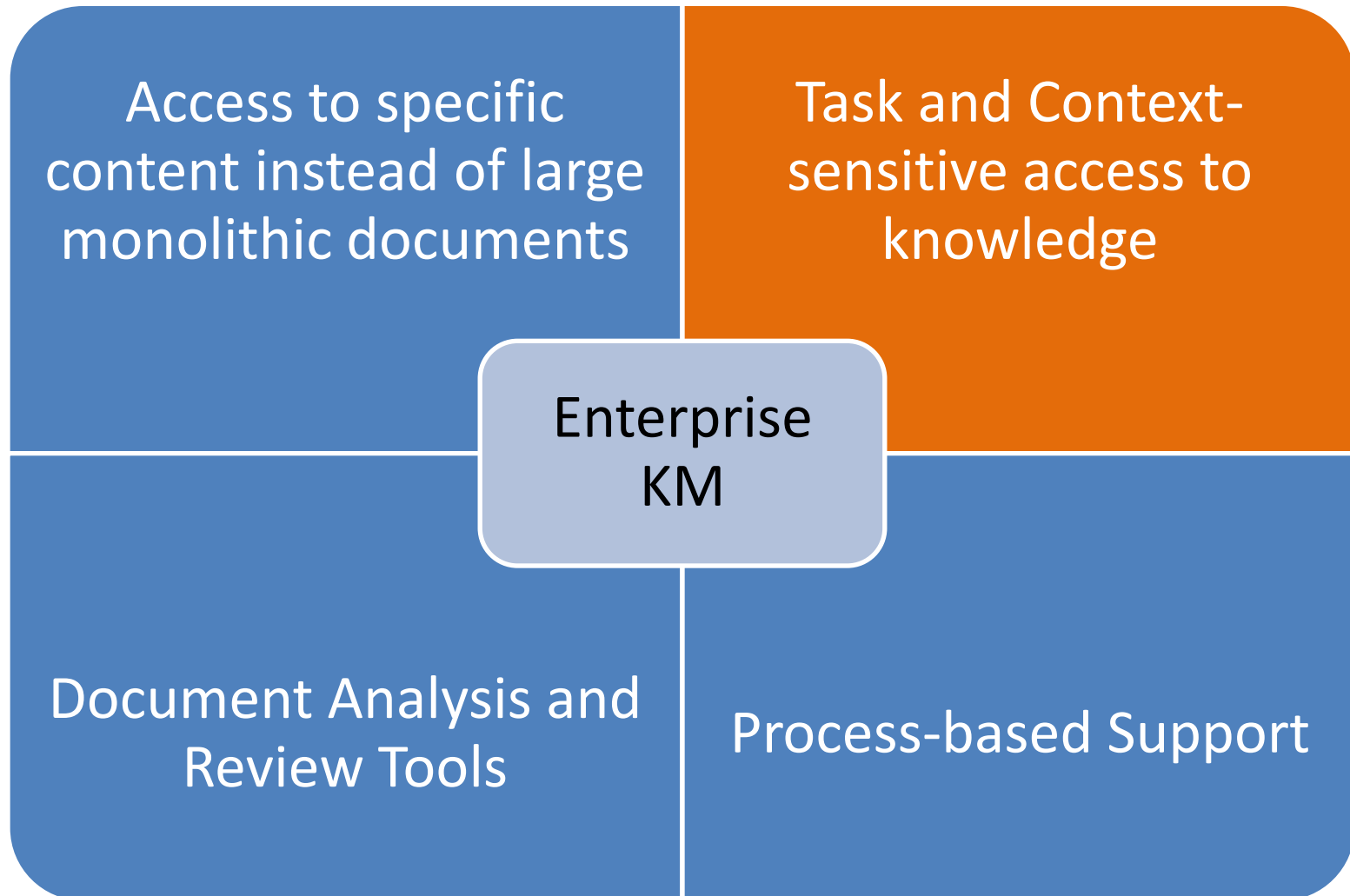
The screenshot displays the Microsoft PowerPoint 2010 interface. The title bar reads 'Presentation3 - Microsoft PowerPoint'. The ribbon includes 'Home', 'Insert', 'Design', 'Animations', 'Slide Show', 'Review', 'View', 'Add-Ins', and 'Format'. The 'Format' ribbon is active, showing 'Font' (Lucida Sans Unic, 18), 'Paragraph' (Text Direction, Align Text, Convert to SmartArt), and 'Drawing' (Shape Fill, Shape Outline, Shape Effects, Arrange, Quick Styles, Find, Replace, Select). The main slide area shows the title 'Outsourcing' and a diagram consisting of two horizontal lines with square nodes at the ends and a central square node connected by vertical lines. The 'Diagrams' task pane on the right has a search bar with 'Solution overview' and a 'Go' button. It displays a search result for 'Solution overview' with a diagram showing a process flow from 'Assessment' to 'Strategy Formulation' to 'Migration Planning'. Below this is a circular diagram with 'Internet Buying' and 'Internet Selling' and a pyramid diagram at the bottom.

Slide 2 of 2 "Concourse" English (U.S.) 79%

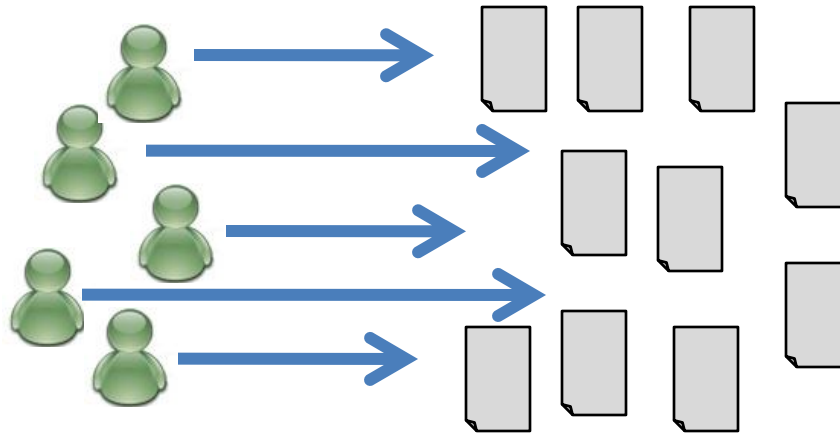
3. Task specific application



Areas of Improvement

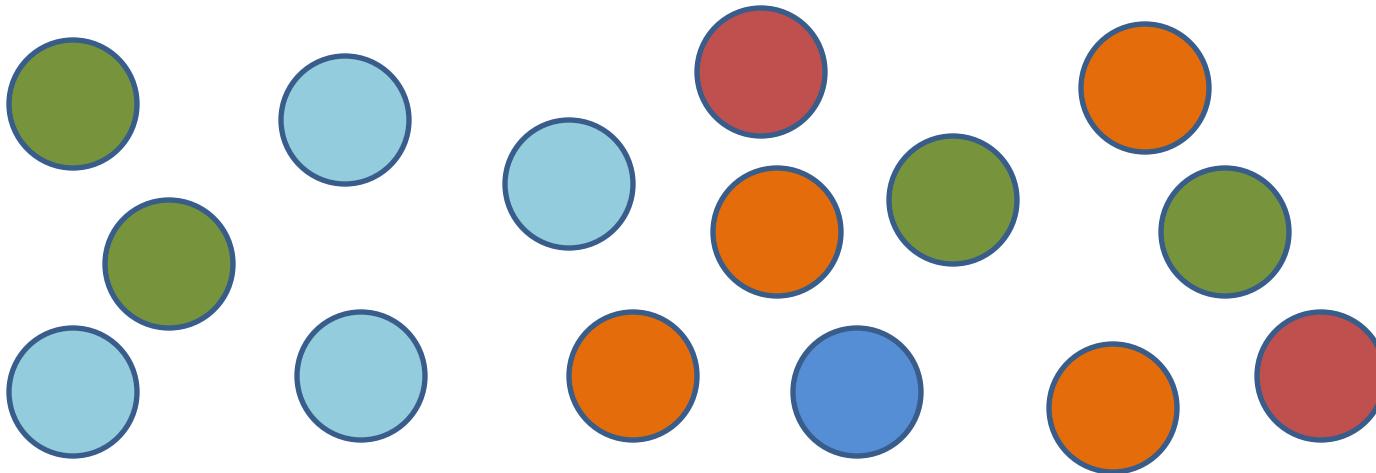


Context

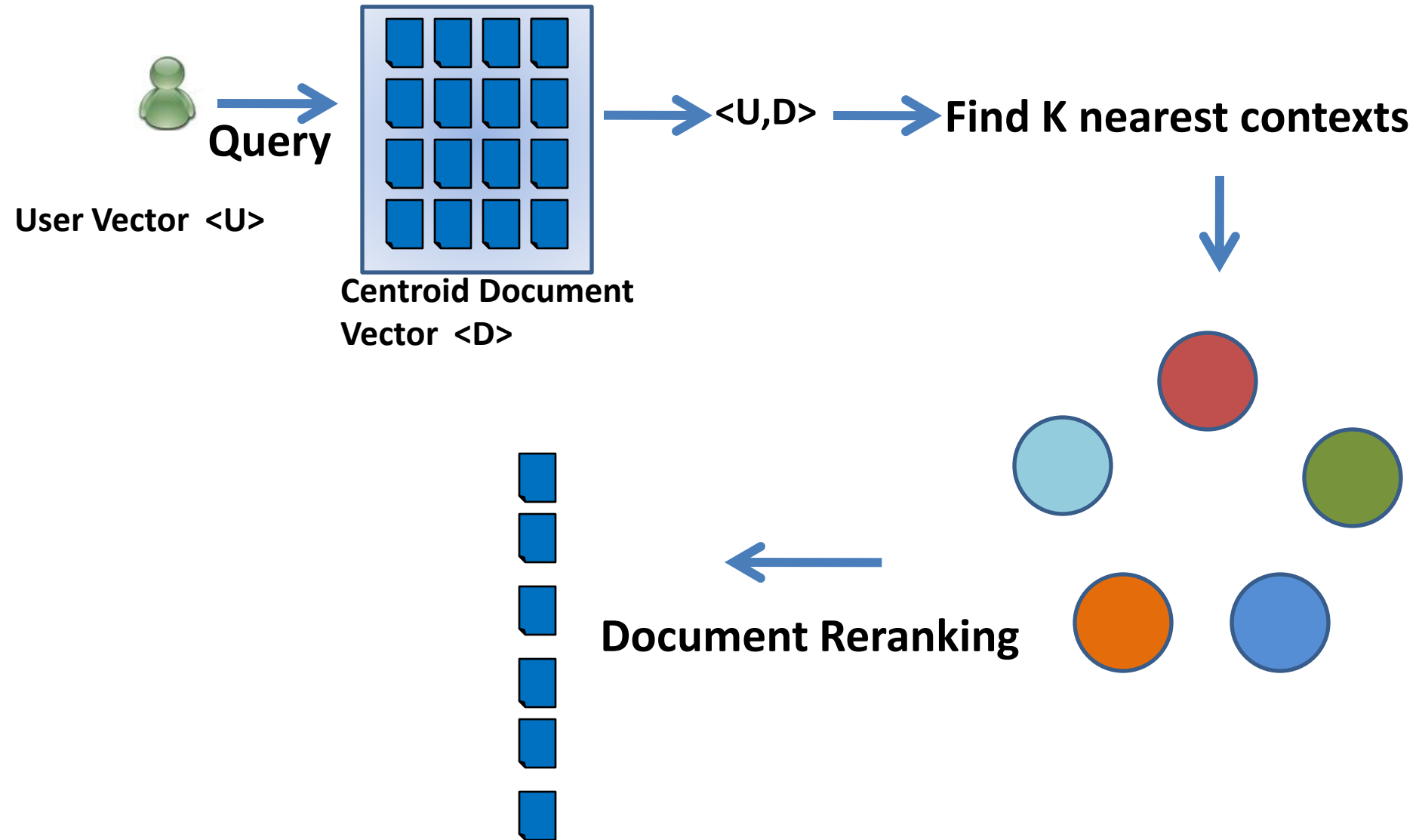


Context = **User Role** + **Document Attributes**

Super Contexts = Cluster All contexts over time



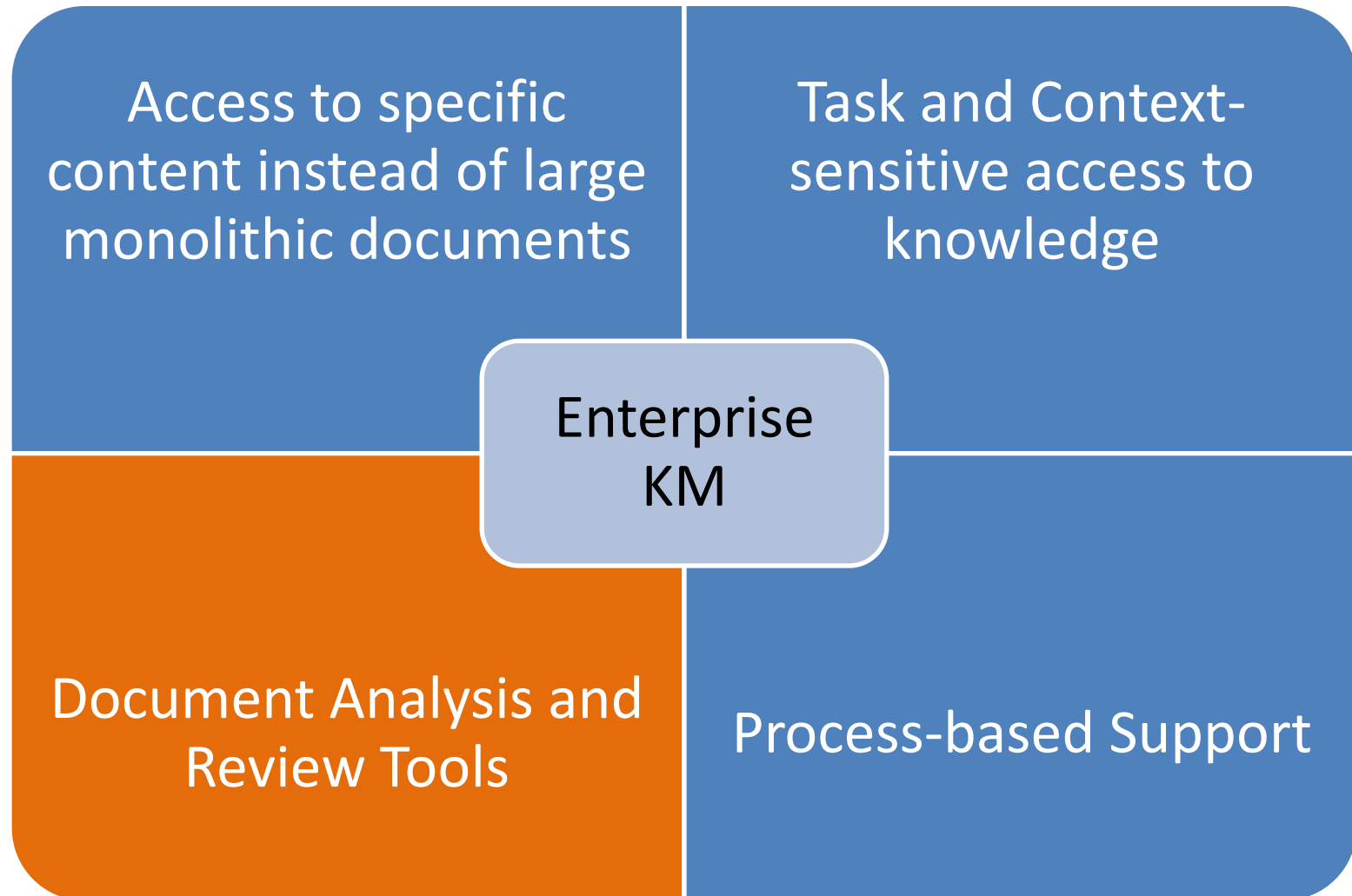
Context-sensitive ranking

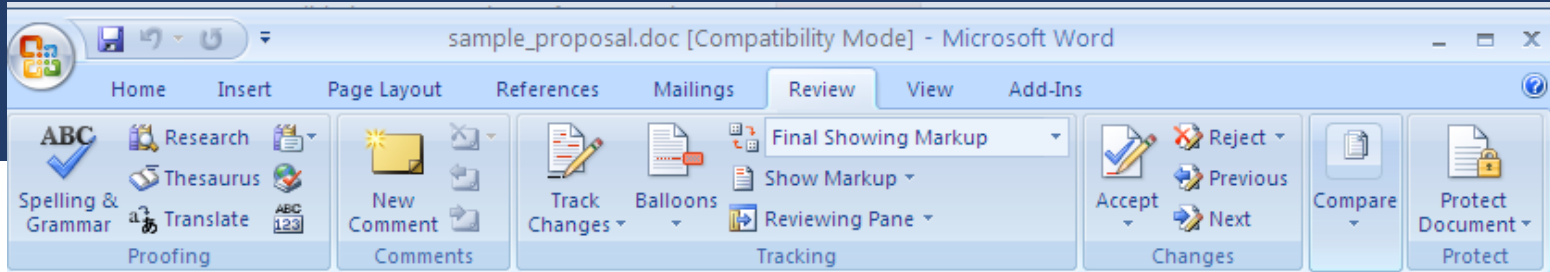


Contextual Access to Information

- Current search engines are based on keywords and do not work well with long (sentence-like) queries
- Need to be used differently (example)
 - We are writing a response to a proposal for an online retailer that is exploring new marketing capabilities to improve their customer and marketing results through direct mail, email, web and outbound telemarketing programs. The focus is on flexible delivery of an end-to-end CRM solution, as well as alternate delivery methods that could appeal to the client. The proposed end-to-end option will leverage our alliances with Acxiom, Unica and Teradata
- Long Query ([prototype demo](#))

Areas of Improvement





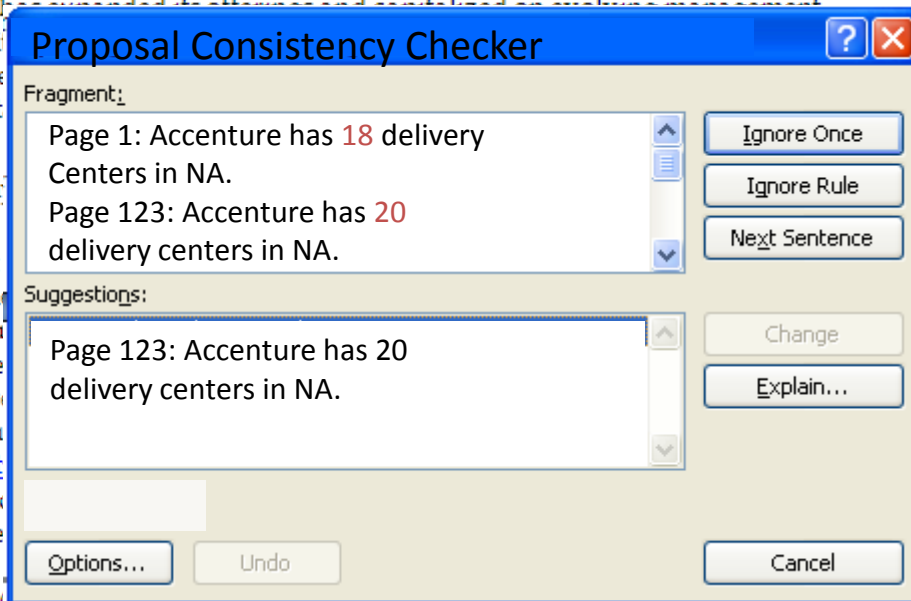
Established in 1989 primarily as a technology consultant and systems integrator, Accenture soon began offering a new breed of business integration solutions to clients – solutions that aligned organizations' technologies, processes and people with their strategies.

Throughout its history, Accenture has succeeded in identifying and capitalizing on evolving management trends and technologies to benefit its clients. This approach has led to the deployment of e-business integration; led the deployment of e-commerce and electronic services; and has established a strong presence in the technology consulting and systems integration markets.

2.1.2 Overview of the Company

Accenture is a global management services company. Committed to deliver solutions that help them become high-performance businesses, Accenture leverages its broad global resources and a proven technology to help clients improve their performance. The company generated net revenue of \$10.1 billion in 2006. The company's home page address is www.accenture.com. Accenture provides a Solution to all Novelis locations, in addition to the corporate executive offices.

Accenture is a "Featured Partner" (the highest level of partnership) with Hyperion. The Accenture and Hyperion alliance brings critical business performance management capabilities to clients by combining Hyperion's leading Business Performance Management (BPM) software with Accenture's vast business process and global industry knowledge. Together, we help organizations realize the benefits of



sample_proposal.doc [Compatibility Mode] - Microsoft Word

HomeInsertPage LayoutReferencesMailingsReviewViewAdd-Ins

Font

Paragraph

Styles

1.1. AaAaBbCcL1. AaE

DV Head...EmphasisHeading 1

Change Styles

Editing

ACTIVEKnowledge-Powered Enterprise

Search

Win Theme Analyzer

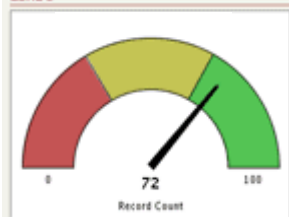
1 Cost Efficiency

2 Team Experience

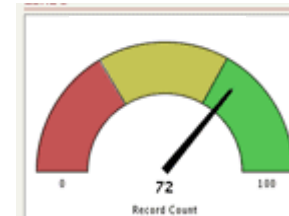
3 Brand Expansion

4 Long-term relationship

Analyze Win Themes



Cost Efficiency



Team Experience

All Word

Connected to Office Online

2. Vendor Profile

2.1. Company Overview and Ownership

2.1.1 Brief Company History

Continuous innovation and rapid transformation have been themes throughout Accenture's history. Established in 1989 primarily as a technology consultant and systems integrator, Accenture soon began offering a new breed of business integration solutions to clients – solutions that aligned organizations' technologies, processes and people with their strategies.

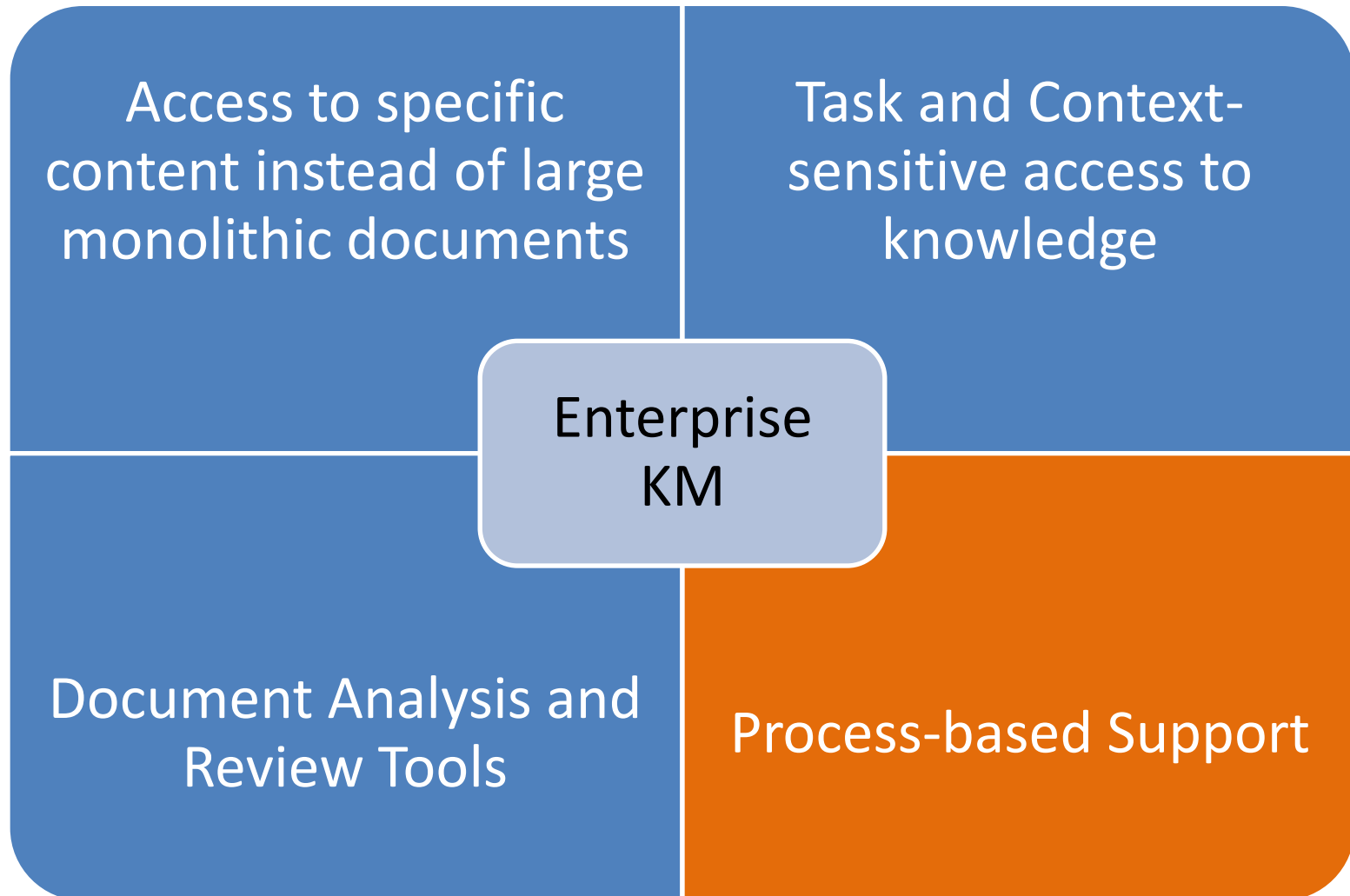
Throughout its history, Accenture has expanded its offerings and capitalized on evolving management trends and technologies to benefit its clients. The company pioneered systems integration and business integration; led the deployment of enterprise resource planning, customer relationship management and electronic services; and has established itself as a leader in today's global marketplace.

2.1.2 Overview of the Company and Its Subsidiaries

Accenture is a global management consulting, systems integration, technology and outsourcing services company. Committed to delivering innovation, Accenture collaborates with its clients to help them become high-performance businesses. With deep industry and business process experience, broad global resources and a proven track record, Accenture can mobilize the right people, skills and technologies to help clients improve their performance. With more than 140,000 people in 48 countries, the company generated net revenues of US\$16.7 billion for the fiscal year ended Aug. 31, 2006. Our home page address is www.accenture.com. This global reach will be an asset in delivering the HFM Solution to all Novelis locations, including North America, South America, Asia Pacific, and Europe in addition to the corporate executive offices.

Accenture is a "Featured Partner" (the highest level) with Hyperion. The Accenture and Hyperion alliance brings critical business performance management capabilities to clients by combining

Areas of Improvement



Publishing Enterprise Content Today

Gather

- KM team gathers project materials
- Filters out “irrelevant” stuff

Upload

- They bundle up docs & tag them
- Bundles are uploaded

Tag

- Create indices and tag with enterprise taxonomies to organize bundles

Encouraging Content Submission

- Constantly monitor laptop hard drives and index new documents
- Classification system to tag new documents

Some Metadata fields

Industry	<ul style="list-style-type: none">Automotive, Financial, Electronics, Government, Retail, Forestry, Products,...
Offerings (OGs)	<ul style="list-style-type: none">Products, Strategy, CRM, Supply Chain, Resources,...
Doc type	<ul style="list-style-type: none">Offering, Proposal, Credential,...
Business Function	<ul style="list-style-type: none">Accounting, Learning, HR, Outsourcing, Program Management,...

Level 2 subcategories branch out x10

Ordering Sequential Prediction Tasks

- Determine order of tasks to maximize overall performance
- Intractable
- Approximation (Lad et al. SDM 2009):
 - Find Pairwise Preferences
 - Combine to form an optimal ordering

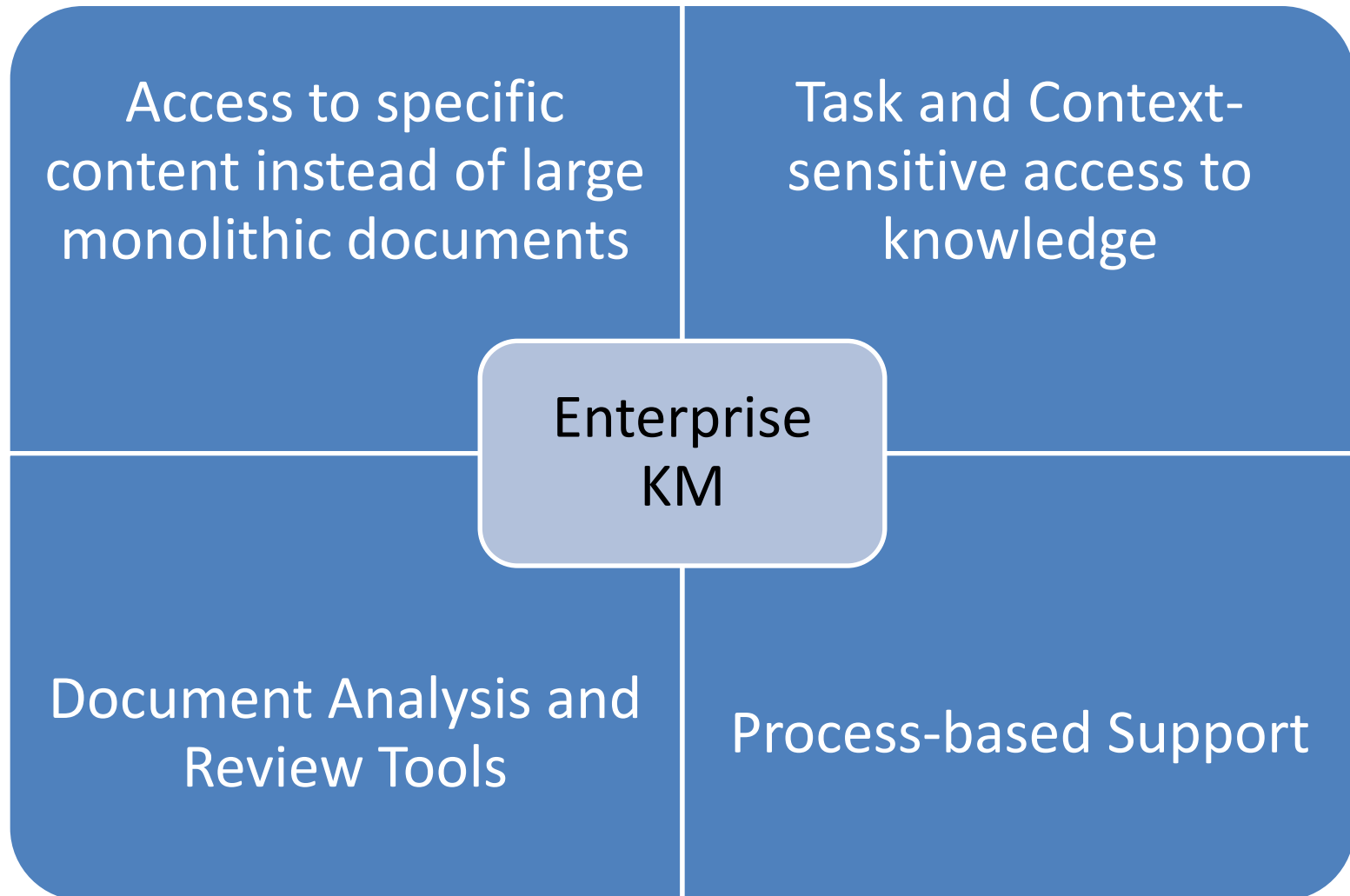
Distributed Active Learning

- Querying the user offers many interesting problems for active learning
 - Online only? No pool of instances available.
 - Query for all contributions by a user – Structured example?
 - Constraints/Dependencies among labels for contributions
- Querying KM team members might even have structure

Confidentiality Issues

- Structured Data – Guessing Anonymity (Rachlin et al. PinKDD 2008)
- Challenge: Unlike structured data, text does not contain identified sensitive attributes
- Tension: Redact documents while still preserving some aspect of the utility (Cumby, ECML 2009 Demo)
- K-confusability

Summary



Enterprise KM: A different perspective

The Consumer's Life



The Consumer's Life



Te Knowledge:
Na Pervasive
Bo Dense
Co Rich
Co Accessible

Crowd:

Loose Organization
Scalable Efforts
Transparency



Work Life



Work Life



Staffing....Communities

accnture

Advanced Search

on... Accenture Portal

First Name Direct Dial

Knowledge:

Te

Sparse
A Siloed
R Stale
Internal

Crowd:

A burden, rather than
an opportunity.



Personalized Knowledge Models



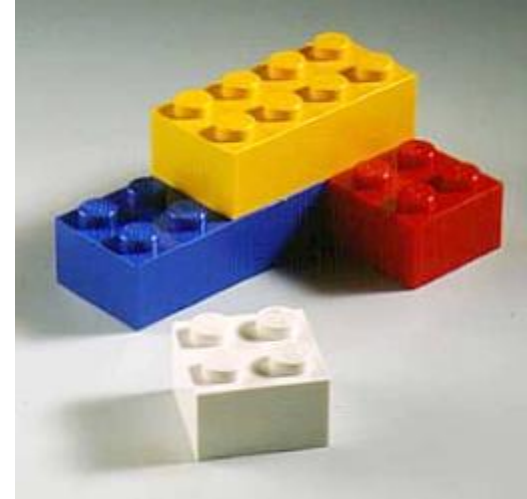
Vs.



In the enterprise: Top Down:
"This is what you will use"

In the outside world: Bottom Up:
"I choose or make what I'll use"

"Lego blocks" that will let our community make the tools they can use to provide and consume enterprise information.



Clothing Retailer

Product Database : Table					
	SKU	UPC	Item	Quantity	Price
	005654	1212234455	DKNYJTee	1	\$49.00



Lauren by Ralph Lauren Catlin Twill Pant \$89.50

A sophisticated pant in substantial cotton twill is contoured for a sleek silhouette.

details

shipping & returns

- Cotton
- Machine washable
- Imported
- Web ID: 273846

Overall Rating ★★★★★ 4.5 out of 5

6 of 6 (100%) customers would recommend this product to a friend.

[Read all reviews](#) [Write a review](#)

Share this Product: [f](#) [g+](#) [p](#) [t](#) [v](#) [e](#)

size

color

qty

Select Size ▼

Select Color ▼

1 ▼

[size chart](#)


[reset selection](#)

Sporting Goods

Shop by Sport / Department

SEARCH By Keyword or [View Item](#) #

GO



LARGER IMAGE

Asics ZD361 Unrestrained Wrestling Knee Pad

Item no: 1306778

Our Price:
\$12.99

SAVED ALERT
Eligible for FREE Ground Shipping on orders of \$99 or more
(See Promotion Details)

The Asics® ZD361 Unrestrained wrestling knee pad utilizes dual-density padding with Kinetoflex for impact protection and an innovative pad design for a better fit. The terry loop inner fabric and mesh back panel helps to enhance moisture management to keep you comfortable.

Color : Size
Please select a color : size

Qty:

[ADD TO CART](#) [ADD TO WISHLIST](#)

Bill Me Later BUY NOW! NO PAYMENTS FOR 90 DAYS
on purchases over \$250

AVAILABILITY: In stock, leaves warehouse in 1 - 2 full bus. days. - [Details](#)

Standard Ground Delivery available - [Details](#)

Browse:

[Knee Pads](#)


Product Tools:

[Email A Friend](#)
Product Reviews
[ADD YOUR REVIEW](#)

Related Categories:

[Ear Guards](#)
[Knee Straps](#)
[Singlets](#)
[See all Asics](#)


Related Items:



[Asics ZW602 Junior GEL Wrestling Ear Guard](#)
Our Price: \$21.99

Black


Add to Cart ☐



[Asics ZW352 Unrestrained Wrestling Ear Guard](#)
Our Price: \$24.99

Red

Add to Cart ☐



[Bute High Cut Lycra Jam Wrestling Singlet](#)
Our Price: \$29.97

Navy : Small

Add to Cart ☐

[ADD TO CART](#)

[back to top](#)

Features

- Innovative pad design
- Dual-density padding with Kinetoflex
- Terry loop inner fabric and perforated sleeve

[back to top](#)

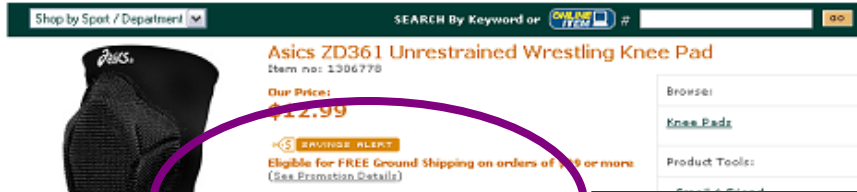
*List price is for reference only. No sales may have occurred at this price.

*List price is for reference only. No sales may have occurred at this price.

**See product page for details.

[secure shopping](#) | [privacy policy](#) | [terms and conditions](#) | [investor relations](#) | [contact us](#) | [product availability and price](#) © 2001-2006.

Sporting Goods



The Asics® ZD361 Unrestrained wrestling knee pad utilizes dual-density padding with Kinetof foam for impact protection and an innovative pad design for a better fit. The terry loop inner fabric and mesh back panel helps to enhance moisture management to keep you comfortable.

Attribute

Value

padding

dual-density

#material# inner fabric

terry loop

sleeve

perforated

pad design

innovative

Features

- Innovative pad design
- Dual-density padding with Kinetof foam
- Terry loop inner fabric and perforated sleeve

- Dual-density padding with Kinetof foam
- Terry loop inner fabric and perforated sleeve
- Innovative pad design

[back to top](#)

*List price is for reference only. No sales may have occurred at this price.

*List price is for reference only. No sales may have occurred at this price.

**See product page for details.

[secure shopping](#) | [privacy policy](#) | [terms and conditions](#) | [investor relations](#) | [contact us](#) | [product availability and price](#) © 2001-2006.

Challenge: How do you describe these products?



Extreme Case: eBay

All Categories

☐ Search titles & descriptions

Browse Categories

Category	Format	Listings	Location	
All Categories	All Items	All Active	Available on: eBay.com	Show

☒ Show number of items in category ☐ Show category numbers

[Antiques](#) (213451)
[Antiquities](#) (5061)
[Architectural & Garden](#) (13710)
[Asian Antiques](#) (25012)
[Books & Manuscripts](#) (6562)
[Decorative Arts](#) (33460)
[Ethnographic](#) (4220)
[Furniture](#) (16049)
[Home & Hearth](#) (913)
[Linens & Textiles \(Pre-1930\)](#) (9232)
[Maps, Atlases & Globes](#) (9777)
[Maritime](#) (2634)
[Mercantile, Trades & Factories](#) (744)
[Musical Instruments \(Pre-1930\)](#) (341)
[Periods & Styles](#) (5977)
[Primitives](#) (13542)
[Restoration & Care](#) (46)
[Rugs & Carpets](#) (19060)
[Science & Medicine \(Pre-1930\)](#) (1805)
[Sewing \(Pre-1930\)](#) (2260)
[Silver](#) (36947)
[Reproduction Antiques](#) (1369)
[Other](#) (4730)
[See all Antiques categories...](#)

[Collectibles](#) (1974383)
[Advertising](#) (132014)
[Animals](#) (93120)
[Animation Art & Characters](#) (94804)
[Arcade, Jukeboxes & Pinball](#) (12238)
[Autographs](#) (8956)
[Banks, Registers & Vending](#) (5340)
[Barware](#) (10342)
[Beads](#) (1514)
[Bottles & Insulators](#) (13618)
[Breweriana, Beer](#) (31822)
[Casino](#) (15407)
[Clocks](#) (11449)
[Comics](#) (147817)
[Cultures & Ethnicities](#) (59523)
[Decorative Collectibles](#) (205113)
[Disneyana](#) (67015)
[Fantasy, Mythical & Magic](#) (22536)
[Historical Memorabilia](#) (92234)
[Holiday & Seasonal](#) (35719)
[Kitchen & Home](#) (57285)
[Knives, Swords & Blades](#) (65112)
[Lamps, Lighting](#) (17714)
[Linens, Fabric & Textiles](#) (22115)
[Metalware](#) (9781)

[Jewelry & Watches](#) (2105770)
[Children's Jewelry](#) (13765)
[Designer Brands](#) (38434)
[Engagement & Wedding](#) (105485)
[Ethnic, Regional & Tribal](#) (49543)
[Fashion Jewelry](#) (831412)
[Fine Jewelry](#) (121625)
[Handcrafted, Artisan Jewelry](#) (67967)
[Jewelry Boxes, Cases & Display](#) (25253)
[Jewelry Design & Repair](#) (57839)
[Loose Beads](#) (211374)
[Loose Diamonds & Gemstones](#) (100458)
[Men's Jewelry](#) (70898)
[Vintage & Antique Jewelry](#) (112976)
[Watches](#) (255139)
[Other](#) (3821)
[Wholesale Lots](#) (39781)
[See all Jewelry & Watches categories...](#)

[Music](#) (1852118)
[Accessories](#) (3390)
[Cassettes](#) (17288)
[CDs](#) (1534322)
[DVD Audio](#) (1403)
[Records](#) (256823)

Having a product attribute database allows....

- Transfer learning from campaigns to related products and categories and **scale** to a large number of products efficiently
- Dealing with the Cold-start problem especially for fast-changing categories
- Ad networks and aggregators to learn across client

craigslist

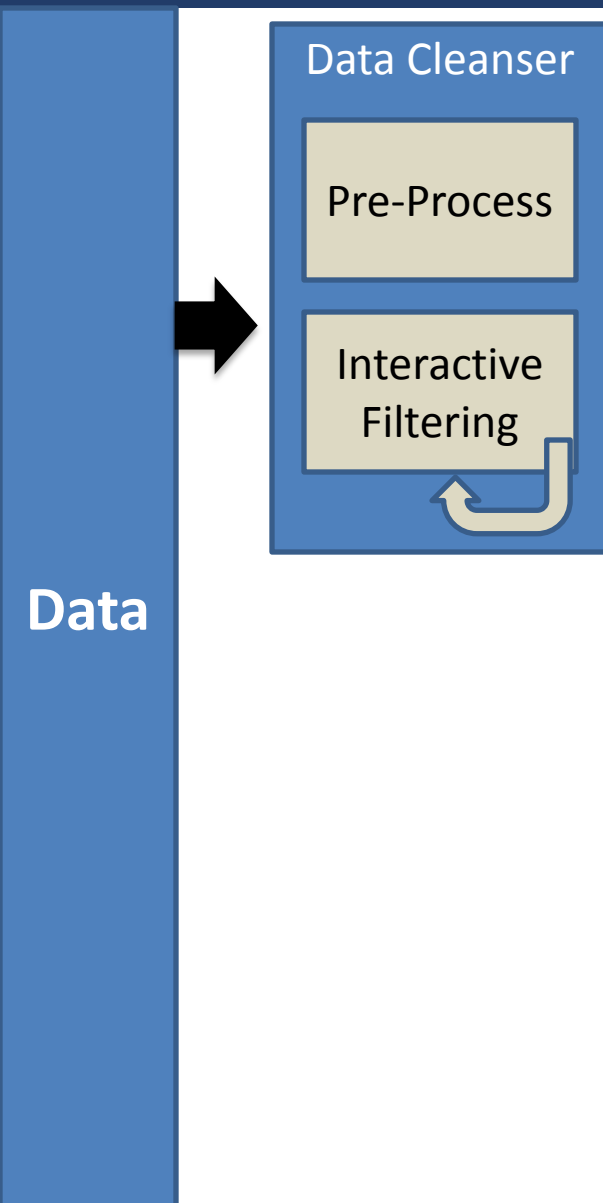


What can businesses do with such enriched databases?

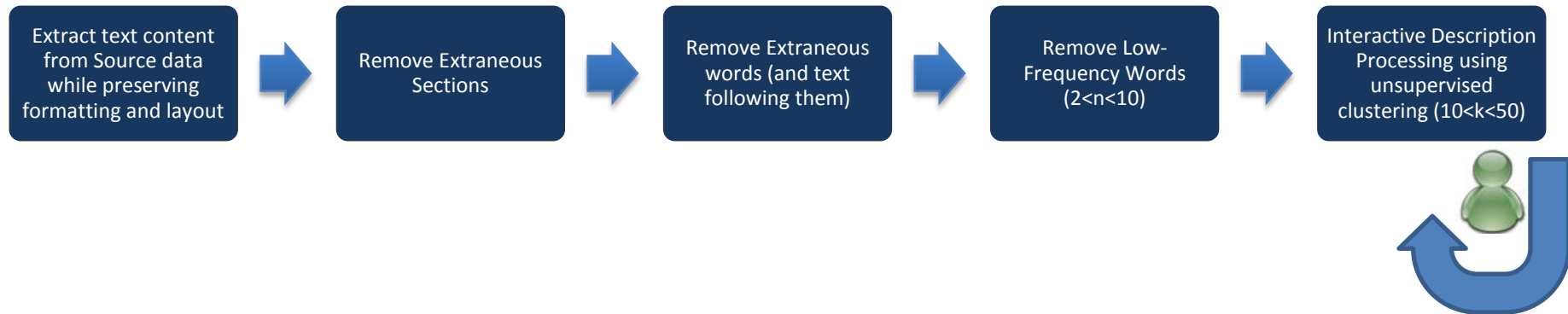
- **Assortment optimization**
 - Sample Questions: Do people buy Bic Gel Grip Pen because of the brand being Bic or the Gel grip?
- **Supply chain management**
 - Sample Question: Are these two products comparable?
- **Procurement**
 - Sample Question: What should I buy and from whom?
- **Marketing**
 - Sample Question: How should I talk about my products?
- **Competitive Intelligence:**
 - Sample Question: How many high-end TVs does my competitor sell as compared to me?
- **Product lifecycle planning**
 - Sample Question: How do products change over time?

Challenges

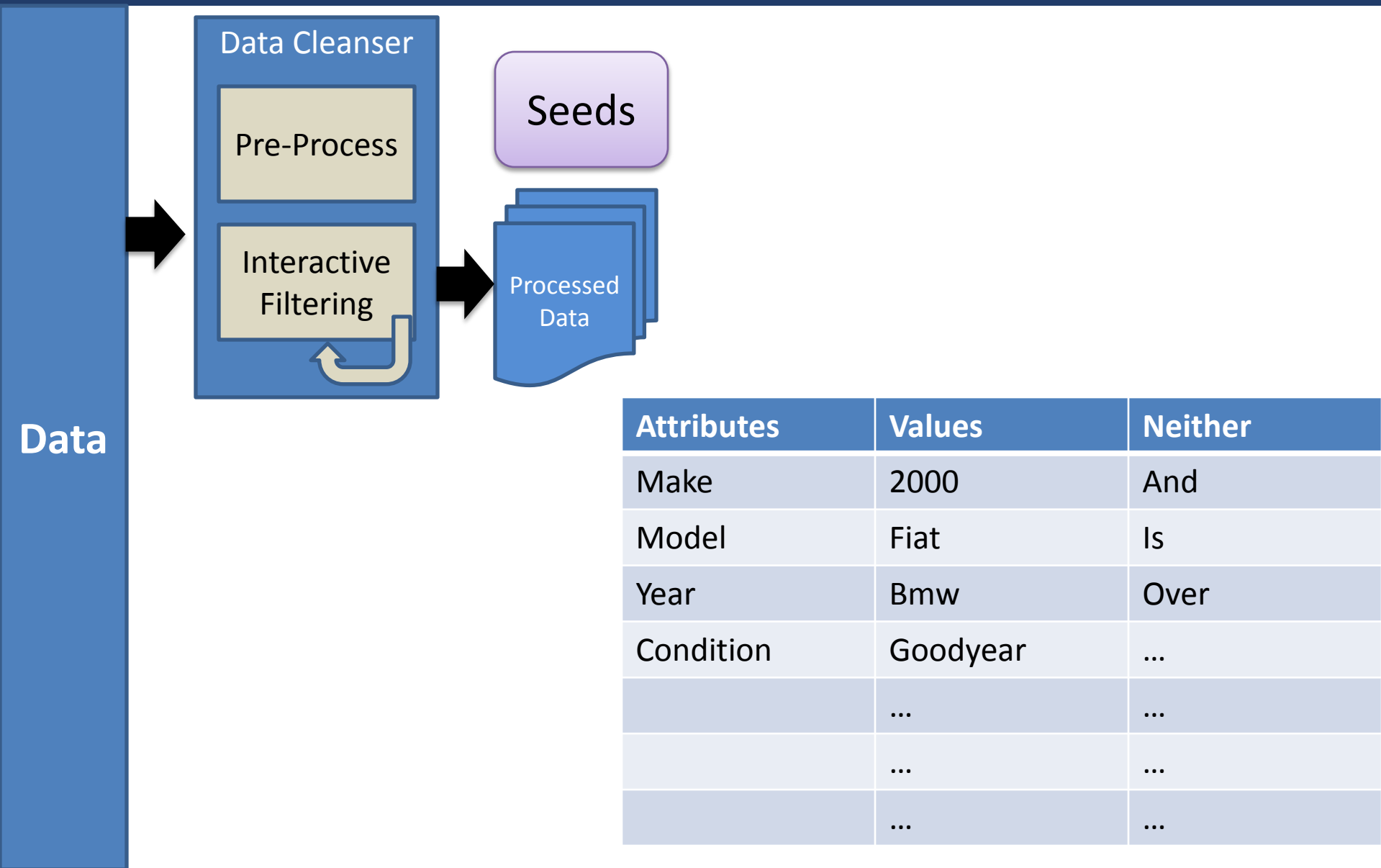
- Signal/noise ratio
 - Extra information
- Labeled data is noisy/misleading
- Heterogeneous sources (title, desc, tables)
- Adhoc sellers
- Typos
- Looooong sentences



Pre-Processing



Keep	Cluster ID	Descriptive Words	Size
<input type="checkbox"/>	1	guarantee, free, worldwide, ship, anywhere	35432
<input checked="" type="checkbox"/>	2	Chevy, grand, front, wheel, silver	21211
<input checked="" type="checkbox"/>	3	Tire, rubber, tread, michelin, goodyear	5000
<input checked="" type="checkbox"/>	4	#number#, model, make, year, all	4323
<input type="checkbox"/>	5	Return, carefully, package, accepted	3244



Unsupervised Seed Augmentation

- Goal: High Precision and Low Recall
- Intuition: Attributes and Values often occur as pairs of consecutive words where the 1st word is the value and 2nd is the attribute
 - Back pocket, side pocket, front pocket
- 2nd word shouldn't occur with the same word every time (phrase) and also shouldn't occur with “too many” words
 - Pittsburgh Steelers, Fifth Avenue, running across

Unsupervised Seed Generation

- All bigrams $w_j w$ are considered candidates for attribute-value pairs:
 - w_j is a potential value and w is a potential attribute
- Let w_j (with $0 < j < k$) be the set of unique words that occur just before w
- Sort all w_j from highest to lowest $P(w_j | w)$
- Retain all w_j such that the sum of the highest-ranking $p(w_j / w) = z$ (we use $z=0.5$)

- Compute **cumulative mutual information**:

- Let $p(w, w_{1...k}) = \sum_{j=1}^k p(w, w_j)$. Then,

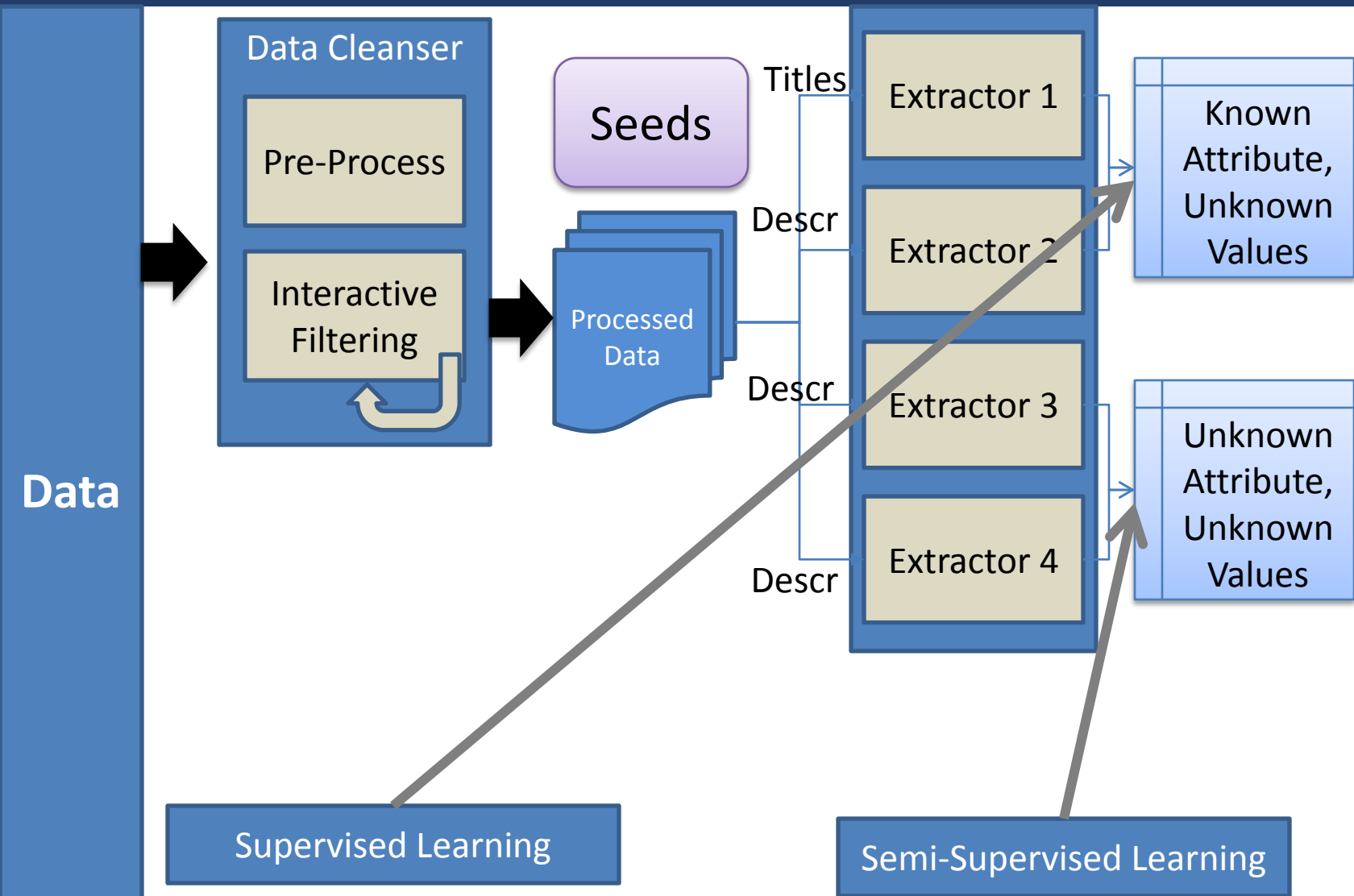
$$cmi(w_{1...k}; w) = \log \frac{p(w, w_{1...k})}{(\lambda * \sum_{j=1}^k p(w_j)) * ((\lambda - 1) * p(w))}$$

where $0 < \lambda < 1$.

- Retain as pairs those whose c.m.i. exceeds a threshold.

Examples of extracted attribute-value pairs

Attribute	Values
case	carrying, storage
compartment	main, racquet
pocket	ball, welt, side-seam, key



Titles

- Supervised Learning
- Sequential Classification

Descriptions

- Supervised Learning from Title: Apply to Description
- Supervised Learning from Title: Apply to Description
- Semi-Supervised Learning from Title: Apply to Description
- Learning Patterns from Tables

Classification: Co-EM with Naïve

- Split data into two views:
 - View1: stemmed word itself and POS tag
 - View2: 8-word context and their POS tags
- Training data for first classification iteration:
 - **Seeds** extracted in previous step
 - Lists of **generic attributes**: countries, colors, materials

Classification: Co-EM with Naïve Bayes

- General Co-EM procedure:
 1. Initialize based on labeled data
 2. Train *view1* classifier, label *view1* of unlabeled data
 3. Use *view1* labels to train *view2* classifier
 4. Label *view2* of unlabeled data using *view2* classifier
 5. Repeat 2, 3, and 4
 6. Use both classifiers to get final probabilities on unlabeled data

Classification: Co-EM with Naïve Bayes

- When labeling *view2*, estimate word and class probabilities using:
 - Current assignments of classes to *view1* words (use probability distribution over all classes)
 - Co-occurrence counts of *view1* words with each *view2* word (how often did a *view1* word appear in *view1*'s context)
- Reverse procedure when labeling *view1*

Training View2 from View1

- Estimate new *view2* class prior probabilities:

$$P(c_k) = \frac{1 + \sum_i^{n_1} \text{cnt}(\text{view1}_i) * P(c_k | \text{view1}_i)}{\text{numclasses} + \sum_i^{n_1} \text{cnt}(\text{view1}_i)}$$

- Estimate new *view2* word probabilities:

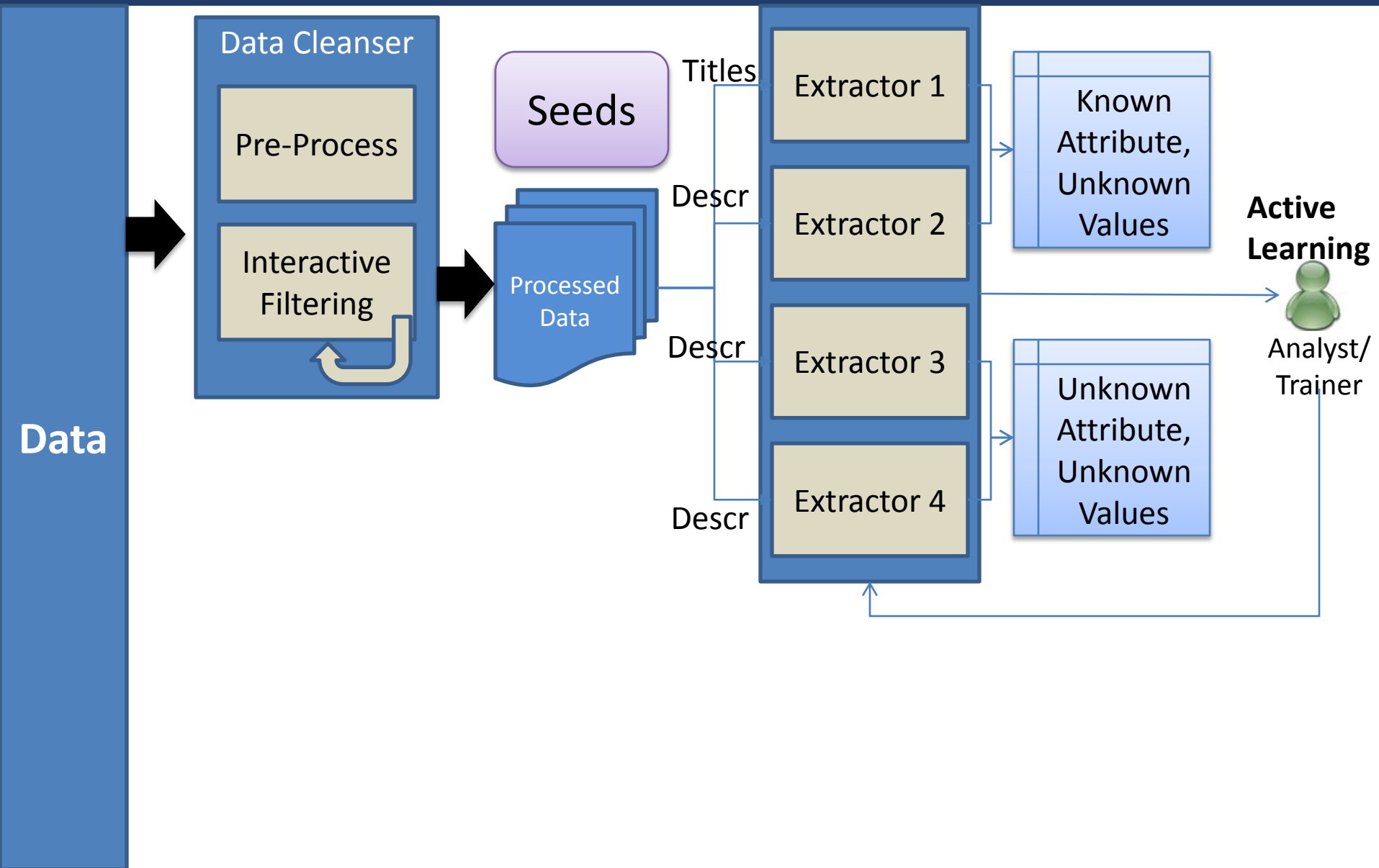
$$P(\text{view2}_j | c_k) = \frac{1 + \sum_{i=1}^{n_1} \text{cooc}(\text{view1}_i, \text{view2}_j) * P(c_k | \text{view1}_i)}{n_2 + \sum_{i=1}^{n_1} \text{cooc}(\text{view1}_i, \text{view2}_j)}$$

- Label *view2* words:

$$P(c_k | \text{view2}_i) \propto P(c_k) * P(\text{view2}_i | c_k)$$

- Use both view1 and view2 classifiers to label unlabeled data

$$P(c_k | \langle view1_i, view2_j \rangle) = \frac{P(c_k | view1_i) + P(c_k | view2_j)}{2}$$



Active Learning

- Multi-view Semi-supervised + Active Learning
- Variety of metrics
 - Frequency
 - KL
 - Frequency weighted KL
- Interactive version of co-EM + NB
 - Orders of magnitude faster (from 15 minutes to 10 seconds)

Current extraction results

Attribute	Value	Sentence	Score
polycotton	1 1/2-inch	1 1/2-inch polycotton blend tape	0.89
rolls	4	4 rolls white athletic tape	0.89
#color# rolls	white	4 rolls white athletic tape	0.89
#material# ghillies	Metal	Metal ghillies	0.87
Short	sleeves	Short sleeves	0.87
Moisture-wicking	design	Moisture-wicking design	0.87
Design	Slim-fitting	Slim-fitting design	0.87
#material# loop inner	Terry	Terry loop inner fabric and perf...	0.78
loop inner	fabric	Terry loop inner fabric and perf...	0.78
#material# Seamless	knit	Seamless knit construction	0.76
perforated	sleeve	Terry loop inner fabric and perf...	0.76
#material# collar	Ribbed	Ribbed collar and cuffs	0.74
logo	1	Embroidered And 1 logo	0.74
#material# piece	Metal	Metal piece on the underside	0.73

Training Progress

Attribute	Value	Sentence	Score
tips	Actual	Actual tips demonstrated on court	1.00
tape	blend	1 1/2-inch polycotton blend tape	1.00
tape	athletic	4 rolls white athletic tape	1.00
#material# upper	leather	Synthetic leather upper	1.00
support	Performance	Performance support striping	1.00
#material#	94% cotton 4%	94% cotton, 4% spandex	1.00
#material#	spandex	94% cotton, 4% spandex	1.00
#material#	Cotton terry	Cotton terry to absorb perspiration	1.00
seam	neck	Twill tape on neck seam and side sli...	1.00
#material# tape	Twill	Twill tape on neck seam and side sli...	1.00
#country#	USA	Made in USA	1.00
seam	Flat	Flat seam stitching	1.00
#material#	90% cotton 10% ...	90% cotton, 10% Lycra	1.00
support	Knitted-in	Knitted-in support features	1.00

Interactive Training

Please specify whether the following attribute-value pair is a correct pair, or what needs to be corrected. To finish the process, select "Finish correcting".

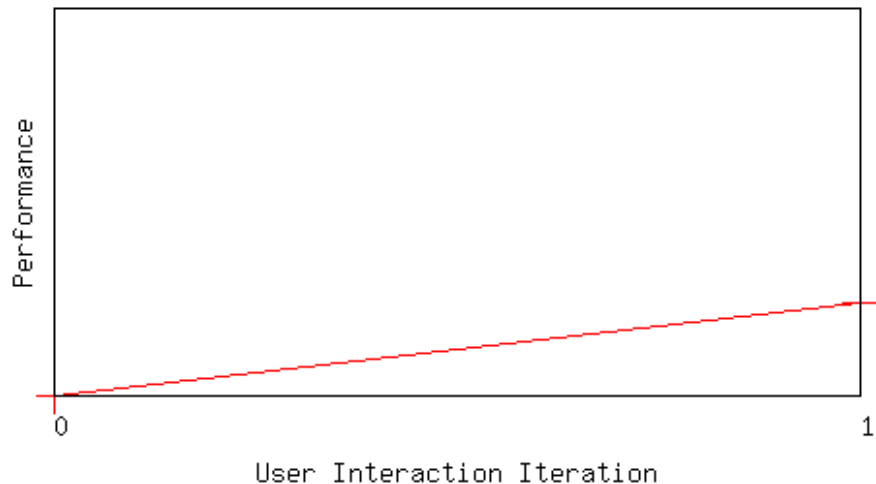
Attribute: 4-way stretch

Value: fabric

Sample context: 4-way stretch fabric

Show web page

- ☐ Correct pair
- ☒ Flip attribute and value
- ☐ Correct attribute, value should be:
- ☐ Correct value, attribute should be:
- ☐ Corrected attribute/value pair:



Number of User Interactions	% Correct	% Partially Correct
1	0.19	0.95

Current extraction results

Attribute	Value	Sentence	Score
tips	Actual	Actual tips demonstrated on co...	1.00
tape	blend	1 1/2-inch polycotton blend tape	1.00
tape	athletic	4 rolls white athletic tape	1.00
#material# upper	leather	Synthetic leather upper	1.00
support	Performance	Performance support striping	1.00
#material#	94 % cotton 4 %	94 % cotton, 4 % spandex	1.00
#material#	spandex	94 % cotton, 4 % spandex	1.00
#material#	Cotton terry	Cotton terry to absorb perspirat...	1.00
seam	neck	Twill tape on neck seam and si...	1.00
#material# tape	Twill	Twill tape on neck seam and si...	1.00
#country#	USA	Made in USA	1.00
fabric	4-way stretch	4-way stretch fabric	1.00
seam	Flat	Flat seam stitching	1.00
#material#	90 % cotton 10 % ...	90 % cotton, 10 % Lycra	1.00

Training Progress

Attribute	Value	Sentence	Score
fabric	4-way stretch	4-way stretch fabric	1.00
fabric	Moisture-wicki...	Moisture-wicking fabric	0.90
#material# inner fabric	Terry loop	Terry loop inner fabric and p...	0.73
weave	4-way stretch	4-way stretch weave	0.42

Interactive Training

Please specify whether the following attribute-value pair is a correct pair, or what needs to be corrected. To finish the process, select "Finish correcting".

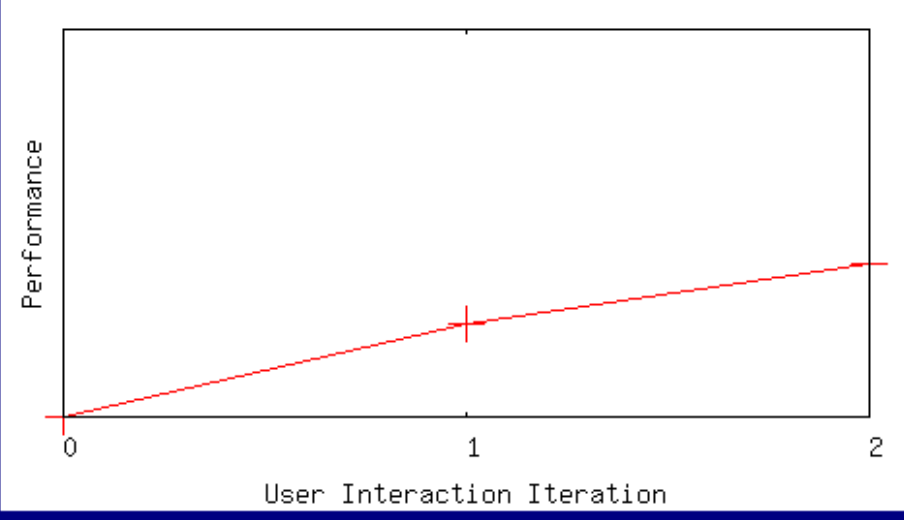
Attribute:

Value:

Sample context:

Show web page

- ☒ Correct pair
- ☐ Correct attribute, value should be:
- ☐ Flip attribute and value
- ☐ Correct value, attribute should be:
- ☐ Corrected attribute/value pair:



Number of User Interactions	% Correct	% Partially Correct
1	0.19	0.95
2	0.32	1.00

Setup User Training

Current extraction results

Attribute	Value	Sentence	Score
#material#	90% cotton 10% ...	90% cotton, 10% Lycra	1.00
support	Knitted-in	Knitted-in support features	1.00
#material#	90% nylon 10% L...	90% nylon, 10% Lycra	1.00
#material#	90% nylon 10% L...	90% nylon; 10% Lycra spandex	1.00
support	Knitted-in added	Knitted-in feature ensure comf...	1.00
design	pad	Innovative pad design	1.00
padding	padding	Dual-density padding with Kin...	1.00
video	72 minute	72 minute video	0.98
video	40 minute	40 minute video	0.98
video	36 minute	36 minute video	0.98
video	35 minute	35 minute video	0.98
#material# waffle	100% cotton	100% cotton waffle knit	0.97
#material# waffle	knit	100% cotton waffle knit	0.97
#material# ghillies	Metal	Metal ghillies	0.90

Interactive Training

Please specify whether the following attribute-value pair is a correct pair, or what needs to be corrected. To finish the process, select "Finish correcting".

Attribute: Dual-density

Value: padding

Sample context: Dual-density padding with Kinetof foam

Show web page

☐ Correct pair

☐ Correct attribute, value should be:

☒ Flip attribute and value

☐ Correct value, attribute should be:

☐ Corrected attribute/value pair:

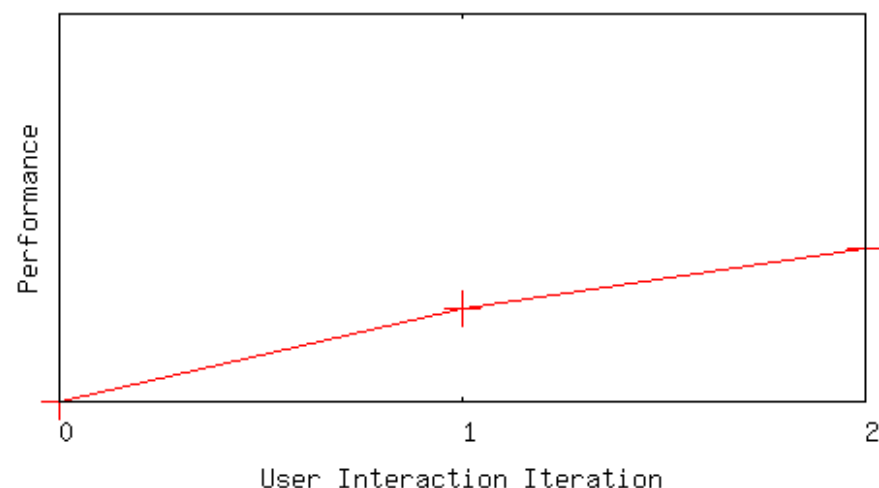
Retrain / Next

Finish training

☐ In all or most contexts

Training Progress

Attribute	Value	Sentence	Score
fabric	4-way stretch	4-way stretch fabric	1.00
fabric	Moisture-wicki...	Moisture-wicking fabric	0.90
#material# inner fabric	Terry loop	Terry loop inner fabric and p...	0.73
weave	4-way stretch	4-way stretch weave	0.42



Number of User Interactions	% Correct	% Partially Correct
1	0.19	0.95
2	0.32	1.00

Current extraction results

Attribute	Value	Sentence	Score
tips	Actual	Actual tips demonstrated on co...	1.00
tape	blend	1 1/2-inch polycotton blend tape	1.00
tape	athletic	4 rolls white athletic tape	1.00
#material# upper	leather	Synthetic leather upper	1.00
support	Performance	Performance support striping	1.00
#material#	94% cotton 4%	94% cotton, 4% spandex	1.00
#material#	spandex	94% cotton, 4% spandex	1.00
#material#	Cotton terry	Cotton terry to absorb perspirat...	1.00
seam	neck	Twill tape on neck seam and si...	1.00
#material# tape	Twill	Twill tape on neck seam and si...	1.00
#country#	USA	Made in USA	1.00
fabric	4-way stretch	4-way stretch fabric	1.00
seam	Flat	Flat seam stitching	1.00
#material#	90% cotton 10% ...	90% cotton, 10% Lycra	1.00

Training Progress

Attribute	Value	Sentence	Score
padding	Dual-density	Dual-density padding with Ki...	1.00
pad design	Innovative	Innovative pad design	0.96
#material# ghillies	Metal	Metal ghillies	0.92
Short	sleeves	Short sleeves	0.92
design	Slim-fitting	Slim-fitting design	0.92

Interactive Training

Please specify whether the following attribute-value pair is a correct pair, or what needs to be corrected. To finish the process, select "Finish correcting".

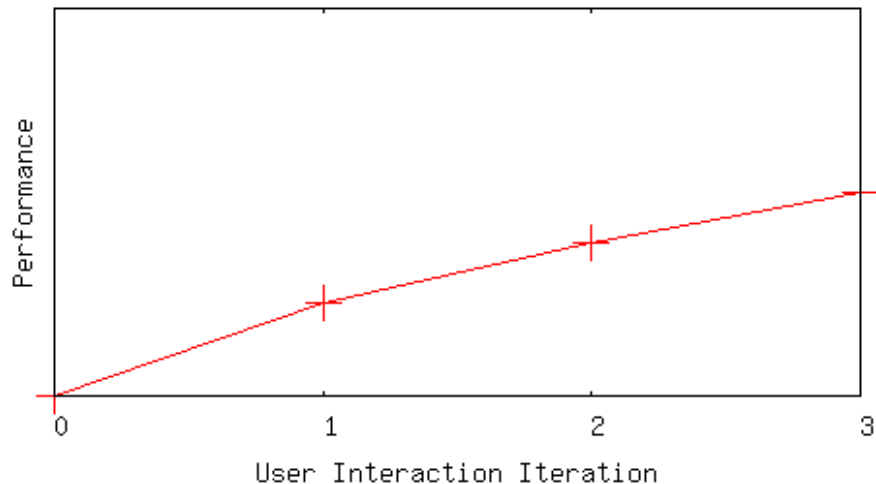
Attribute: Short

Value: sleeves

Sample context: Short sleeves

Show web page

- ☐ Correct pair
- ☐ Correct attribute, value should be:
- ☐ Flip attribute and value
- ☐ Correct value, attribute should be:
- ☐ Corrected attribute/value pair:



Number of User Interactions	% Correct	% Partially Correct
1	0.19	0.95
2	0.32	1.00
3	0.42	1.00

Experimental Results

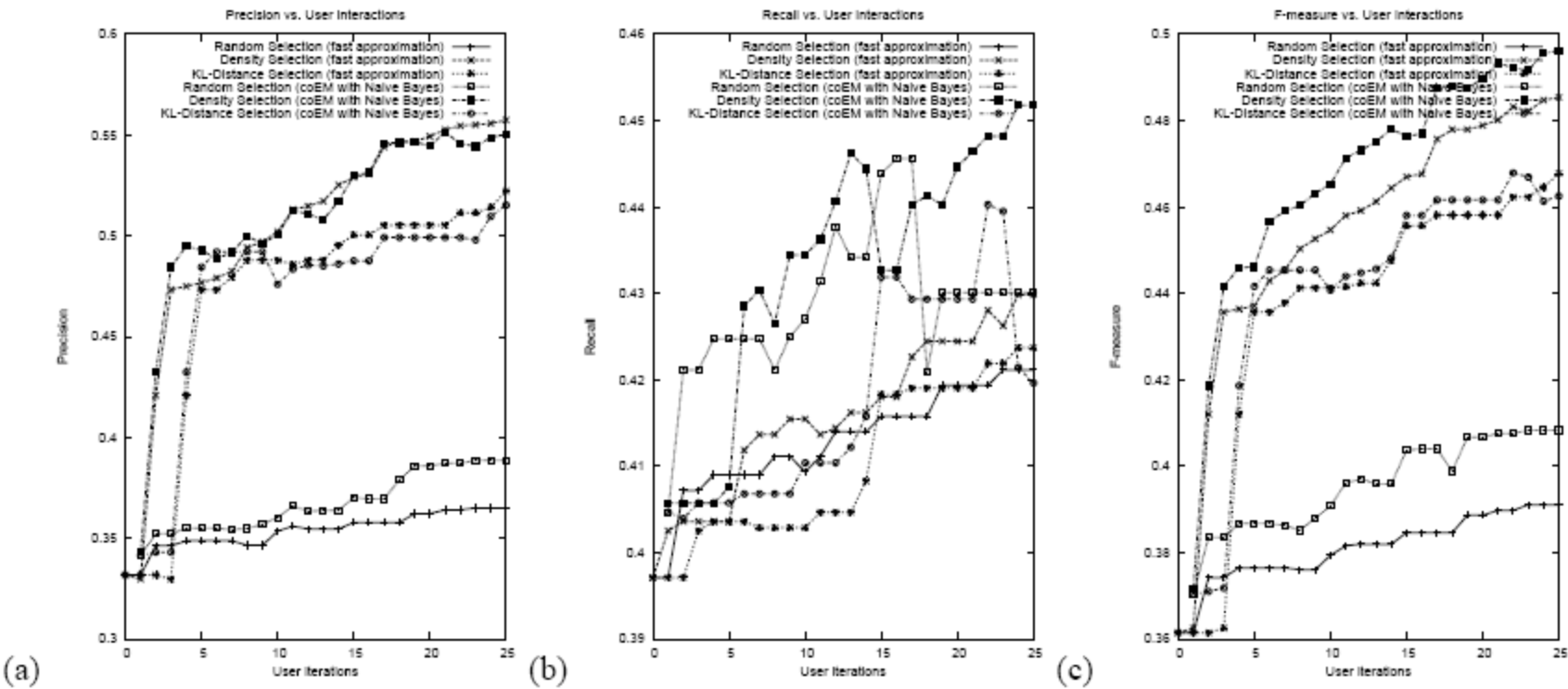


Fig. 1. Precision, Recall, and F-measure for fast algorithm compared to coEM with naive Bayes. The y-value for k indicates the recall, precision, and F-measure after the k^{th} user interaction.

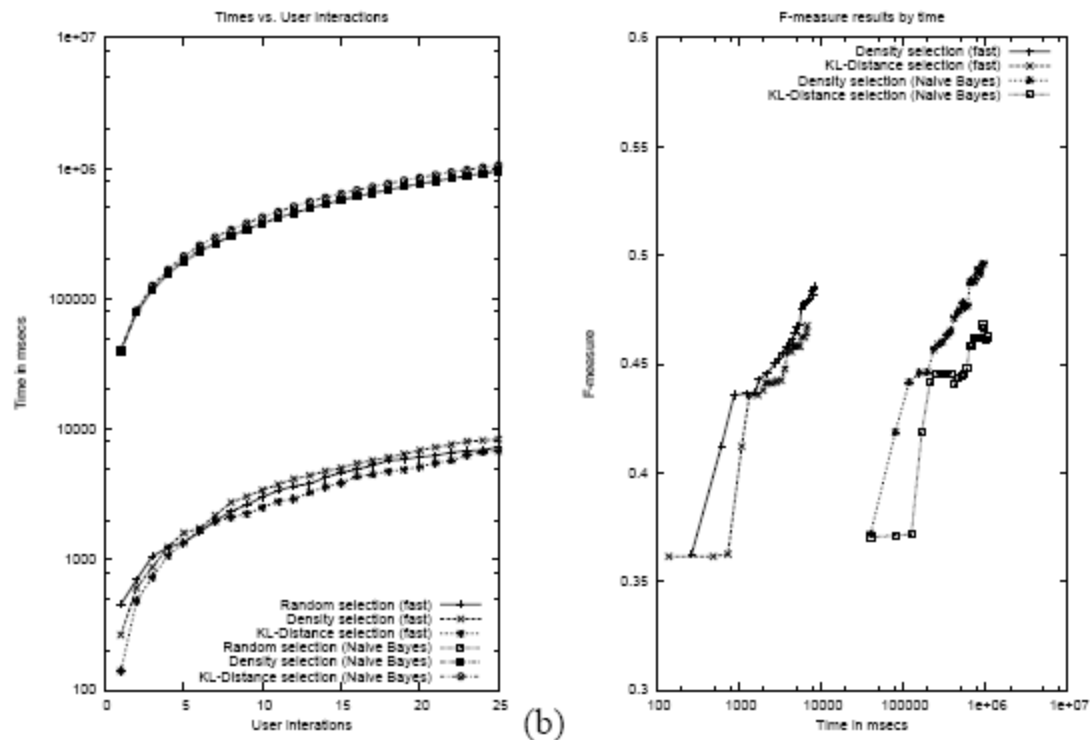
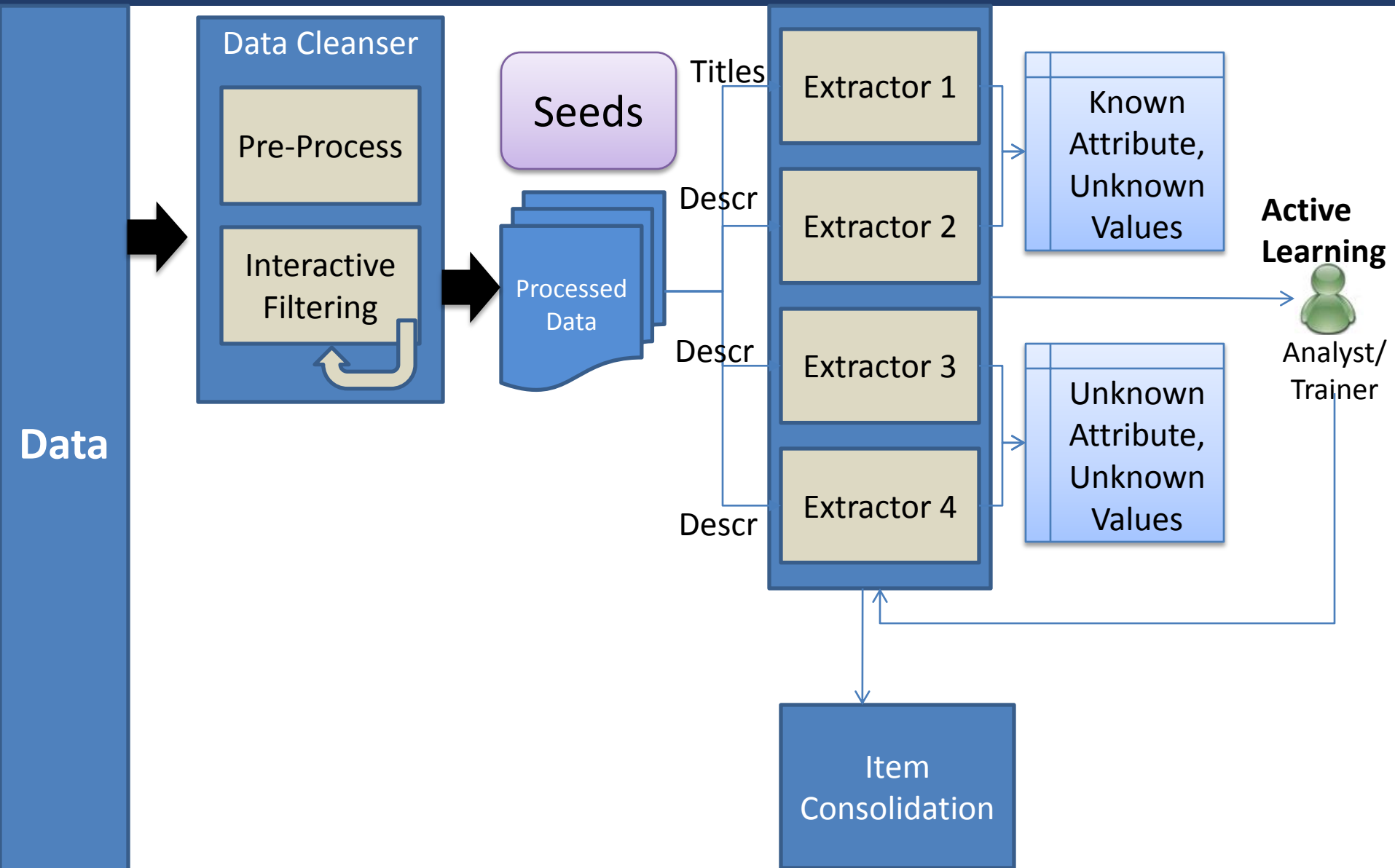


Fig. 2. Time comparison between Fast algorithm and coEM with naive Bayes. In (a), the y-value for k indicates the time the user needed to wait until the k^{th} user interaction.



Forming attribute-value pairs

- Classification results in single words labeled as attributes, values, or neither, but many attributes are phrases
 - Merge adjacent words with the same label and with high correlation scores (χ^2 , yule's Q, m.i.)
- Link attribute phrases and value phrases to form pairs
 - Syntactic dependencies (Minipar)
 - As fallback, also link phrases that have high correlation score or are adjacent
 - Add known attributes that are implicit but present in our KB
 - Assess whether unlinked attributes are “binary”

Examples

Example	Attribute	Value
1 1/2-inch polycotton blend tape	polycotton blend tape	1 1/2-inch
1 roll underwrap	underwrap	1 roll
1 tape cutter	tape cutter	1
Extended Torsion bar	bar	Torsion
Synthetic leather upper	#material# upper	leather
Metal ghillies	#material# ghillies	Metal
adiWear tough rubber outsole	rubber outsole	adiWear tough
Imported	Imported	#true#
Dual-density padding with Kinetof foam	padding	Dual-density
Contains 2 BIOflex concentric circle magnet	BIOflex concentric circle magnet	2
93% nylon, 7% spandex	#material#	93% nylon 7% spandex
10-second start-up time delay	start-up time delay	10-second

Examples of extracted pairs for system run with co-EM

Interesting Results

- New Domain Attributes
 - Cyl (Cylinders)
 - Logo
 - Hub Bore
 - Lip
 - Warranty
 - Included
 - Donor Vehicle
 - Shown

Interesting Results

- New values discovered for Tire Type:
 - **Suv**
 - **Minivan**
 - **Snow**
 - **trailer**
 - **Offroad**
 - **boat.**

Context-dependent extraction

- Misleading terms
 - ‘Other:’ , skype:’
- Typos:
 - ‘finnish’ as an attribute
- POS-dependent
 - ‘Bolt’ (NP) = attribute (Bolt Pattern)
 - ‘Bolt’ (VP) = neither (bolts on to)
- ‘Style’
 - *4 new 17x8 Morxchn Mustang Cobra R style deep lip wheels*
 - *We only guarantee it to fit the same Make Model Year and Body Style as its donor vehicle*

KL metric is useful in detecting confusable words

- interga
- quntity

Summary

- Businesses are not very good at capturing all the structure in their data
- Machine learning can help make that process faster, “better”, and cheaper
- The same techniques are applicable in the semantic web to help create, maintain , and update ontologies

Summary

- Traditional Enterprise Knowledge Management has interesting research challenges
- Next generation KM systems are likely to be personalized mash-ups with personalized knowledge models
- Opportunity for semi-automated semantic learning approaches