# Enterprise and Desktop Search

# Lecture 5:  Desktop Search and Personal Information Management
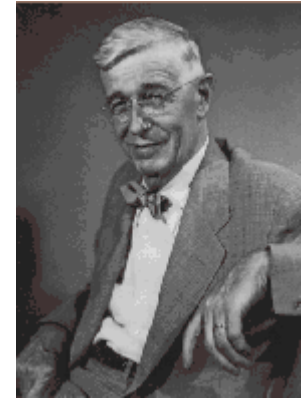
Pavel Dmitriev

Yahoo! Labs

Sunnyvale, CA

USA

Pavel Serdyukov

Delft University of

Technology

Netherlands

**Sergey Chernov**

**L3S Research Center**

**Hannover**

**Germany**

# Searching Personal Collections with Memex

*Posited by Vannevar Bush in "As We May Think"*
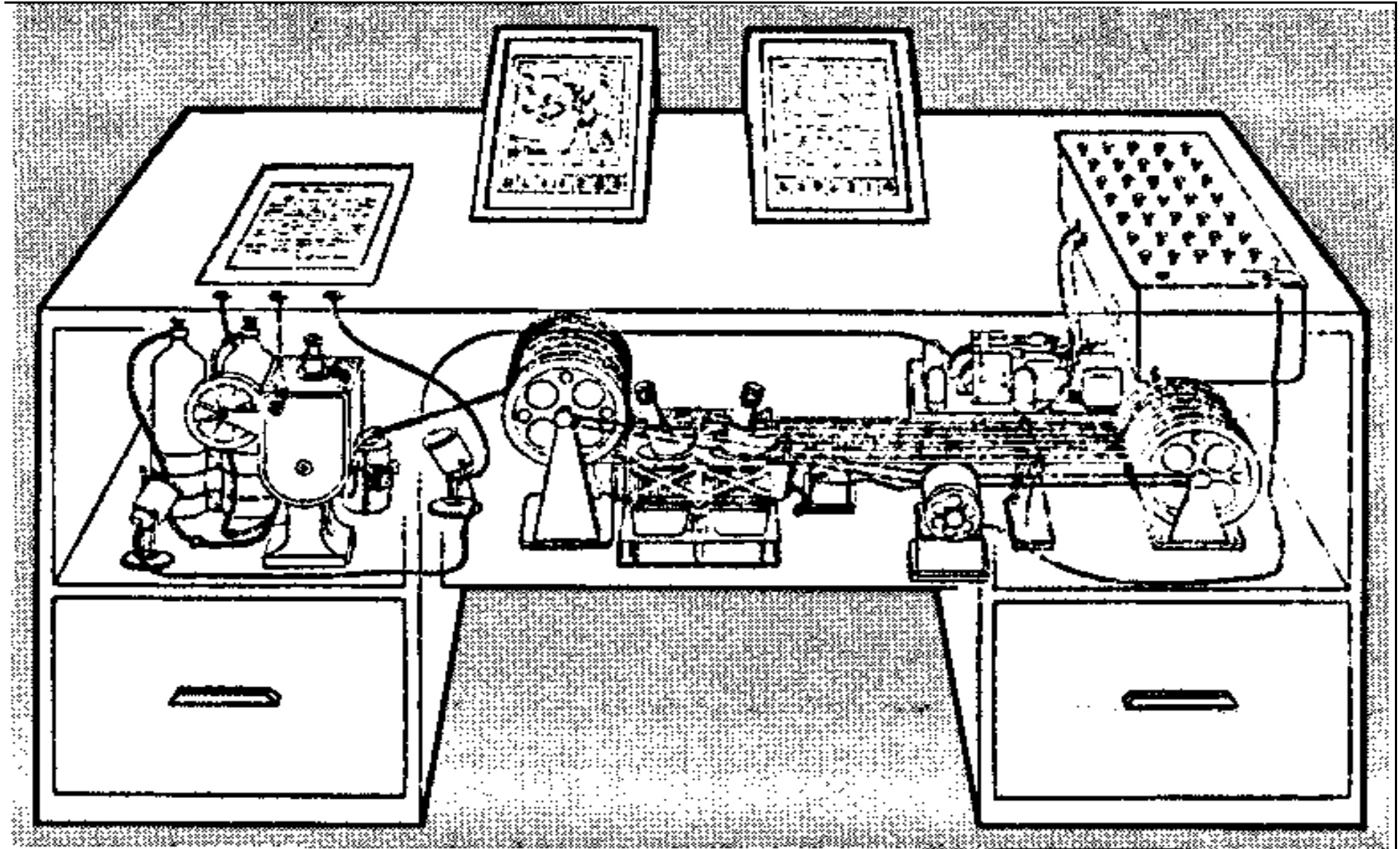*The Atlantic Monthly, July 1945*

"A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility"

Supports: Annotations, links between documents, and "trails" through the documents

"yet if the user inserted 5000 pages of material a day it would take him hundreds of years to fill the repository, so that he can be profligate and enter material freely"
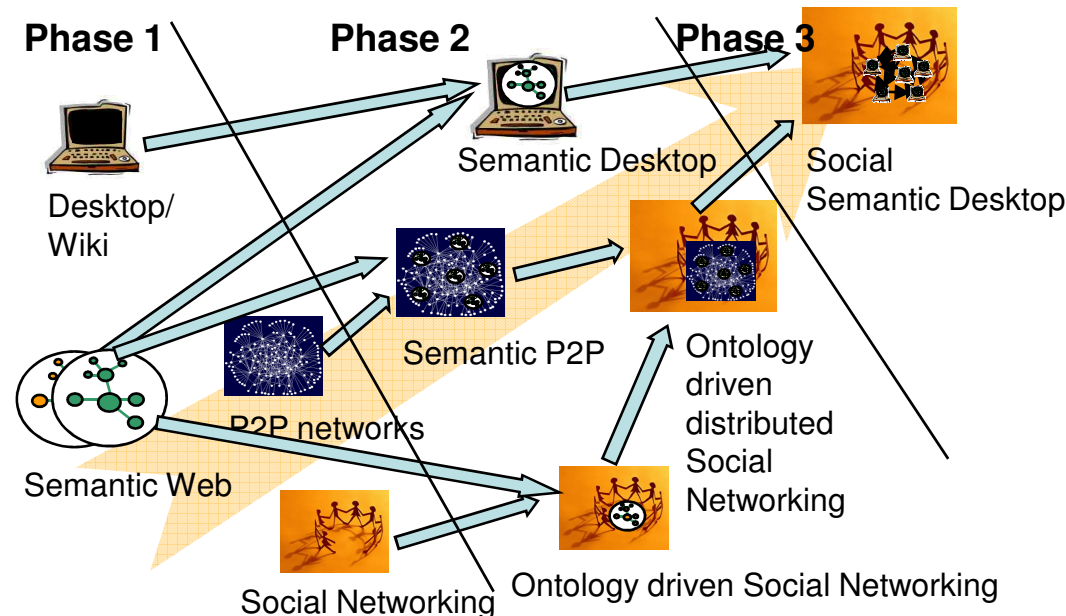
# Sketch of Memex

# Desktop Search and Personal Information Management

- **Desktop search** is the name for the field of search tools which search the contents of a user's own computer files, rather than searching the Internet. These tools are designed to find information on the user's PC, including web browser histories, e-mail archives, text documents, sound files, images and video.

- **Desktop Search** is a part of a more general field of **Personal Information Management** (**PIM**).

- **Personal Information Management** (**PIM**) refers to both the practice and the study of the activities people perform in order to acquire, organize, maintain, retrieve and use information items such as documents (paper-based and digital), web pages and email messages for everyday use to complete tasks (work-related or not) and fulfill a person's various roles (as parent, employee, friend, member of community, etc.)

Source: Wikipedia

# Desktop Search: Motivation

- Why *desktop search*?
  - Size of data on the desktop is big (50k – 500k items) and continously growing
  - Moving towards Social Semantic Desktop
  - Social – communication in a social network
  - Semantic – metadata descriptions and relations



Phase 1     Phase 2     Phase 3

Semantic Desktop

Social Semantic Desktop

Desktop/ Wiki

Semantic P2P

Ontology driven distributed Social Networking

Semantic Web

P2P networks

Social Networking

Ontology driven Social Networking

# What is Desktop?

- Documents (doc, pdf, ppt, xls, html, txt, …)

- Email

- Calendar

- Instant Messengers (ICQ, Skype, MSN messenger, …)

- Pictures

- Music

- Videos

# Desktop Search – Current Status

- Documents on the desktop are not linked to each other in a way comparable to the web
- Simple full text search
  - no personalization
  - no context
  - no ranking possible or too poor
- Metadata enriched search makes use of
  - associations to contexts and activities
  - provenience of information
  - sophisticated classification hierarchies

Google™

Spotlight

Windows

Search

# Differences between Web Search and Desktop Search

- Search on the desktop vs. Search on the Web
  - Re-finding vs. finding
  - Integration across many applications and file formats
  - Users prefer to navigate, not to search
  - Many information types: ephemeral, working, archived
  - Extra sources for ranking improvement:
    - File metadata
    - Usage metadata
    - Folder structure
  - Privacy concerns

# Outline

- Today we will talk about:
  - Modern Desktop Search Engines
  - Research prototypes
  - Just-In-Time Retrieval
  - Context on a Desktop
    - Using context to improve Desktop Search
    - Context Detection
  - PIM Evaluation

# Modern Desktop Search Engines

- Google Desktop (from major web search engine vendor)
- Windows Search (from major OS provider)
- Copernicus (company specialized on DS engines)
- Beagle (open source DS for Linux)
- Yandex (Russian DS)

Some more:

Ask.com, Autonomy, Docco, dtSearch Desktop, Easyfind, Filehawk, Gaviri PocketSearch, GNOME Storage, imgSeek, ISYS Search Software, Likasoft Archivarius 3000, Meta Tracker, Spotlight, Strigi, Terrier Search Engine, Tropes Zoom, X1 Professional Client, etc.

# Desktop Search Architecture



Search Engines Tackle the Desktop,
Bernard Cole, Computer 2005.

# Desktop Search Engines in 2005



**1. Usability**

| | |
|---|---|
| Copernic | 4.80 |
| Archivarius | 4.75 |
| Google | 4.40 |
| MSN | 4.40 |
| Ask Jeeves | 4.25 |

**2. Versatility**

| | |
|---|---|
| Copernic | 4.14 |
| Yahoo! | 3.88 |
| Blinkx | 3.75 |
| ISYS | 3.75 |

**3. Accuracy**

| | |
|---|---|
| Copernic | 4.50 |
| MSN | 4.20 |
| dtSearch | 3.50 |

**4. Efficiency**

| | |
|---|---|
| Archivarius | 4.40 |
| Copernic | 4.20 |
| Ask Jeeves | 3.80 |

**5. Security**

| | |
|---|---|
| Yahoo! | 3.29 |
| Ask Jeeves | 3.14 |
| Google | 3.13 |

**6. Enterprise Readiness**

| | |
|---|---|
| Copernic | 4.00 |
| ISYS | 4.00 |
| Yahoo! | 4.00 |

\* Copernic with Coveo, and Yahoo! with X1

Source: UW E-Business Consortium

| Desktop Search Tool | Version | Score (Min = 1.00, Max = 5.00) Better → |
|---|---|---|
| Copernic Desktop Search | 1.5 Beta | 4.11 |
| Yahoo! Desktop Search | 1.1 Beta | 3.66 |
| Likasoft Archivarius 3000 | 3.14 | 3.62 |
| MSN Toolbar Suite | 2.0 Beta | 3.45 |
| Google Desktop | 1.0 | 3.26 |
| Ask Jeeves | 1.0 Beta | 3.16 |
| Enfish Professional | 6.1 | 3.10 |
| ISYS Desktop | 6.0 | 3.05 |
| dtSearch Desktop | 6.5 | 3.02 |
| diskMETA Pro | 1.0.1 | 2.63 |
| Blinkx | 3.0 | 2.63 |
| HotBot Desktop | Beta | 2.34 |

Source: UW E-Business Consortium

Benchmark Study of Desktop Search Tools, Tom Noda and Shawn Helwig, Technical Report 2005, http://www.uwebi.org/reports/desktop_search.pdf.

# Sample Criteria for DS Comparison

| Search Format | Platform(s) | Feature | Opt-in Feature |
|---|---|---|---|
| Plain text | Windows Vista | Specifying index location | Default search engine |
| HTML pages stored locally | Windows XP | Incremental indexing | Web integration |
| Microsoft Word (.doc) | Mac OS X | Legacy index by scanning | Insecure search |
| Microsoft Excel (.xls) | Linux | Engine download size | Registration |
| Microsoft PowerPoint (.ppt) | Mozilla/Firefox | Install size | Engineering feedback |
| Rich Text Format (.rtf) | Internet Explorer | Combined local/remote search | Software updates |
| Portable Document Format (.pdf) | Opera | Non-anonymous connections | |
| Microsoft Outlook email | Safari | Excluding files | |
| Microsoft Outlook Express email | Languages | Indexing progress indicator | |
| Microsoft address books | | Recoverable index | |
| AOL Instant Messenger | | File type filtering | |
| Standard email folder support | | Deskbar | |
| Standard news folder support | | Support for compressed files | |
| Browser web history | | Support for legacy file formats | |
| Browser secure web history | | Ignoring networked drives | |
| Browser bookmarks | | Click to suspend | |
| Browser address books | | Click to exit | |

# Google Desktop Search

# Windows Desktop Search

# Copernicus Desktop Search

# Beagle Desktop Search

# Yandex Desktop Search

Персональный поиск Яндекса — это программа на вашем компьютере, осуществляющая поиск по файлам и письмам с учётом морфологии русского языка.

**Форма поиска**

**Управление представлением результатов**
Можно выбрать способ группировки и сортировки.

**Конфиденциальность**
Можно запретить искать в определенных папках или целых дисках.

**Цитаты из найденных документов**
Повышают информативность результатов поиска.

**Группировка результатов по типам**
Помогает ориентироваться в большом количестве найденных файлов.

**Вы сможете найти ваши письма в Outlook, Outlook Express, Thunderbird и TheBat!**

**Окно запроса и результатов поиска открывается в обычном браузере**

**Доступ к Персональному поиску из панели задач**

«Пушкин» — 43 результата

Яndex
Найдётся всё

Пушкин

Найти

☐ в найденном

Мой компьютер | Документы 37 | Музыка 3 | Письма 3 | Сохранённые страницы

⊟ Документы 37

Онегин.doc                                                     10 июня
Далее Пушкин цитирует рецензию
Булгарина на седьмую главу...
C:\My Documents\...\Онегин.doc · Открыть папку

Отсортировано
По релевантности
По дате

Статистика.xls                                                 15 июня
Пушкин — 37 лет
C:\Projects\pushkin\Статистика.xls · Открыть папку

«Пушкин» — в документах(37) →

⊟ Музыка 3

Пушкин - Читает свои стихи.mp3                    13 июля
Пушкин - первый подкаст.mp3                        26 июля
Такого как Пушкин.mp3                                   20 июля

⊟ Письма 3

Re: Новая бета yandesk  Александр Быков       25 июля
Я разве Пушкин? Пусть он думает...

Я к вам пишу, чего же боле  Эдмон Дантес     25 июля
Да уж, Пушкин зажигает!

# Research prototypes and Semantic Desktops

- Beagle++ (extended open source DS)
- Semex (includes Malleable Schemas)
- Haystack and Magnet (Semantic Web approach)
- Stuff I've Seen (Phlat predecessor)
- Phlat (was used as a basis for Windows DS)
- PIA (semantic desktop solution from DB area)

Some more:
  Gnowsis, CALO

# Beagle++

P.-A. Chirita, S. Costache, W. Nejdl, and R. Paiu. Beagle++ : Semantically enhanced searching and ranking on the desktop. In ESWC 2006.

- Why is it so hard to find what you need on your desktop – "You still use Google even for files stored on your computer?"

- Current desktop search engines use only full text index

- People tend to associate things to certain contexts

- For desktop search we need to support contextual information in addition to full text!
  - Relationships between information items (citations)
  - Relationships based on interactions (email exchange, browsing history)
  - Relationships between different types of items (authorship, publication venues, email sender information, recommendations)
  - Other situational context

Semantically Rich Recommendations in Social Networks for Sharing, Exchanging and Ranking Semantic Context, Stefania Ghita, Wolfgang Nejdl, and Raluca Paiu. In ISWC 2005.

The Beagle++ Toolbox: Towards an Extendable Desktop Search Architecture, Ingo Brunkhorst, Paul - Alexandru Chirita, Stefania Costache, Julien Gaugaz, Ekaterini Ioannou, Tereza Iofciu, Enrico Minack, Wolfgang Nejdl and Raluca Paiu. Technical Report 2006.

# Scenario 1: The Need for Context Information

- Alice and Bob are working together in the research group

- Alice is currently writing a paper about searching and ranking on the semantic desktop and wants to find some good papers on this topic, which she remembers she stored on her desktop

- Some time ago Bob sent her a very useful paper on this topic as an attachment to an email, together with some useful comments about its relevance to her new semantic desktop ideas

- ***Will Alice find the paper from Bob when issuing a query on the desktop, using the search terms "semantic desktop" ?***

# Context Information is necessary!

- **Problems**:
  - (Mail) Documents sent as attachments lose all contextual information as soon as they are stored on the PC
  - (Web) When searching for a document we downloaded from the CiteSeer repository, we would like to retrieve not only the specific document, but all the referenced and referring papers which we already downloaded as well

- Current desktop search approaches don't make use of desktop specific information, especially contextual information, like:
  - **Email** context
  - **Web** context
  - **Publication** context

# Representing Context by Semantic Web Metadata

- Metadata for resources can be created by appropriate metadata generators
- Ontologies specify context metadata for:
  - Emails
  - Files
  - Web pages
  - Publications
- Metadata have to be application-independent!
→ Store Metadata as RDF
  - generated and used by whatever application you can think of

# Beagle++ Layer Architecture

Beagle++ is our extension of the open source Beagle search project, enabling it to exploit context information

RDF metadata are generated based on ontologies for specific contexts (email, web, etc.)

Indexing and metadata generation on the fly - triggered by events upon occurrence of file system changes (*inotify-enabled linux kernel)*

Benefits:

Context allows us to better organize and find information

Context gives us the possibility to compute the value / importance of resources

# Beagle++ Architecture

# Beagle++: Find more than documents

# Beagle++: Display additional context

# Integrating Keyword and Metadata Search

Text keywords

RDFS properties

Beagle++ Search

Search terms: RDF indexing /storedFrom/attachedTo/receivedFrom Bob.    Anywhere    Find

Matono_etal.pdf, in folder Desktop
Received from Bob <bob@l3s.de>
Subject: Re:About indexing
Folder: Email/Inbox
Last modified July 2, 11:10 PM
An Indexing Scheme for RDF and RDF Schema based ... written in RDF (Resource Description ... Framework) and RDF Schema will ... described with RDF and RDF Schema will ...

Open    Send to..    Reveal in file manager    Show context informantion

1 result.    Show Previous Results    Show More Results

- Search text and metadata on the desktop
- Search efficiently in a user-friendly way
- Simple query language
- No complete schema knowledge necessary

# Documents / RDF Fragments

- Metadata stored as RDF graphs, each document has a corresponding RDF fragment

- Extended documents consisting of both full-text and metadata properties

- Query model supports the operator selection, projection and union, intersection and set difference

- Support for approximate and imprecise metadata queries

- Separation between metadata statements is ensured by positional indices



**Full Text**

**RDF Graph**

**Document**

# Peer-Sensitive ObjectRank [1]

- Step 1: start with PageRank formula – random surfer model

$$r = d \cdot A \cdot r + (1 - d) \cdot e$$

- d = dampening factor
- A = adjacency matrix
- e = vector for the random jump

- Step 2: distinguish between different kinds of objects
- ObjectRank variant of PageRank

# Peer-Sensitive ObjectRank [2]
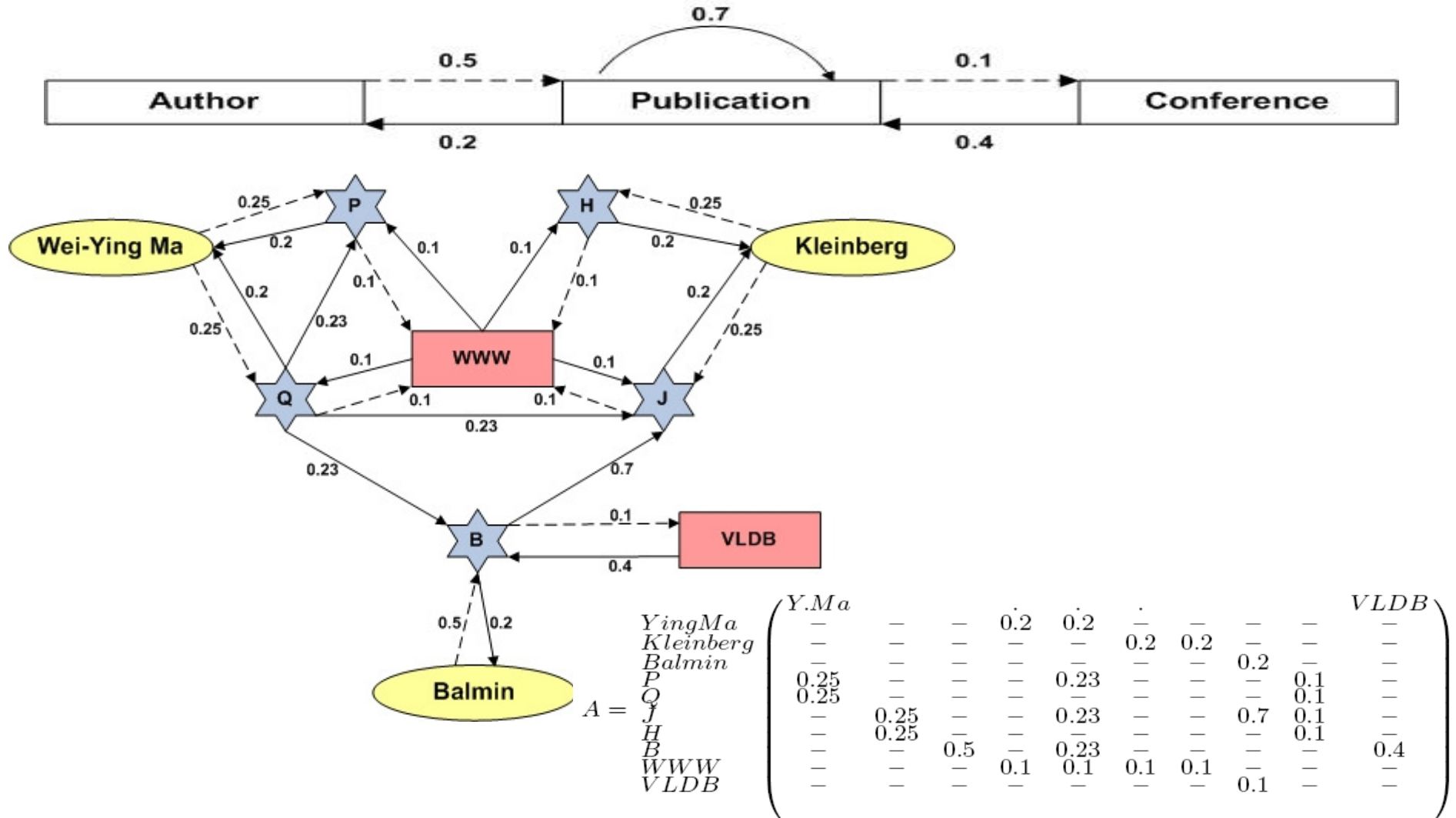
# Peer-Sensitive ObjectRank [3]

- Step 3: Take provenance information into account
- $\rightarrow$ **Peer-Sensitive ObjectRank**
- Represent different trust in peers by corresponding modifications in the **e** vector
- Keep track of the provenance of each resource

$$originates\ (r_i, P_n) = \begin{cases} 1, \text{if } r_i \text{ is in the initial set of } P_n \\ 0, \text{otherwise} \end{cases}$$

$$trust(P_i, P_j) \in [0,1], \text{ the trust value of peer } P_i \text{ for } P_j$$

$$e_k(P_i) = \max_{j=0}^{N} \{trust(P_{i,}P_j) \cdot originates(r_k, P_j)$$
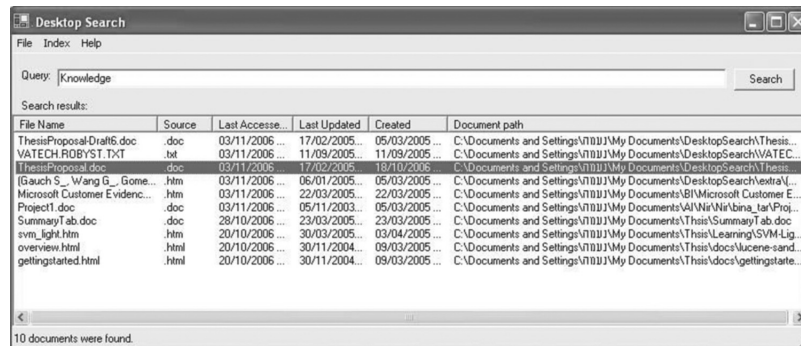
Beagle++ Demo

# Open Source Search Engines

Build your own search engine!

A Comparison of Open Source Search Engines, Christian Middleton and Ricardo Baeza-Yates, Technical Report, 2007 .

| Search Engine | Indexing Time (h:m:s) | | Index Size (%) | | Searching Time (ms) | | Answer Quality P@5 | |
|---|---|---|---|---|---|---|---|---|
| ht://Dig | (7) | 0:28:30 | (10) | 104 | (6) | 32 | | - |
| Indri | (4) | 0:15:45 | (9) | 63 | (2) | 19 | (2) | 0.2851 |
| IXE | (8) | 0:31:10 | (4) | 30 | (2) | 19 | (5) | 0.1429 |
| Lucene | (10) | 1:01:25 | (2) | 26 | (4) | 21 | | - |
| MG4J | (3) | 0:12:00 | (8) | 60 | (5) | 22 | (4) | 0.2480 |
| Swish-E | (5) | 0:19:45 | (5) | 31 | (8) | 45 | | - |
| Swish++ | (6) | 0:22:15 | (3) | 29 | (10) | 51 | | - |
| Terrier | (9) | 0:40:12 | (7) | 52 | (9) | 50 | (3) | 0.2800 |
| XMLSearch | (2) | 0:10:35 | (1) | **22** | (1) | **12** | | - |
| Zettair | (1) | **0:04:44** | (6) | 33 | (6) | 32 | (1) | **0.3240** |

# Selecting an Appropriate Ranking Function

On Ranking Techniques for Desktop Search, Sara Cohen, Carmel Domshlak and Naama Zwerdling, In ACM Transactions on Information Systems 2008.



Lucene-based DS prototype
19 volunteers.
In total 1219 queries
188 queries had a single result,
916 queries has 2-50 results
115 queries had over 50 results.

| FEATURE | $MRR(\tau, \mathcal{S}_{2-50})$ | FEATURE | $MRR(\tau, \mathcal{S}_{>50})$ |
|---|---|---|---|
| SVM | 0.54 | SVM | 0.26 |
| LEXORD | 0.53 | LEXORD | 0.18 |
| **Selective** | 0.5 | **Selective** | 0.17 |
| USERBEST | 0.47 | ACCESSDATE | 0.16 |
| UPDATEDATE | 0.43 | USERBEST | 0.16 |
| NAME | 0.43 | UPDATEDATE | 0.15 |
| ACCESSDATE | 0.4 | CREATEDATE | 0.12 |
| CREATEDATE | 0.39 | NAME | 0.1 |
| SIZE | 0.39 | PATH | 0.1 |
| CONTENT | 0.38 | SIZE | 0.08 |
| NORMALIZEDSIZE | 0.36 | QUERYLOG | 0.07 |
| PATH | 0.34 | DIRRANK | 0.06 |
| QUERYLOG | 0.34 | CONTENT | 0.06 |
| DIRRANK | 0.33 | NORMALIZEDSIZE | 0.06 |
| LEVEL | 0.31 | LEVEL | 0.03 |
| *Random* | 0.28 | *Random* | 0.02 |

$$\text{SELECTIVE}_q(f) \overset{\text{def}}{=} \sum_{\substack{\text{FEATURE} \in \{\text{NAME, PATH,} \\ \text{CONTENT, QUERYLOG}\}}} \frac{\text{FEATURE}_q(f)}{nz(\text{FEATURE}_q)}$$

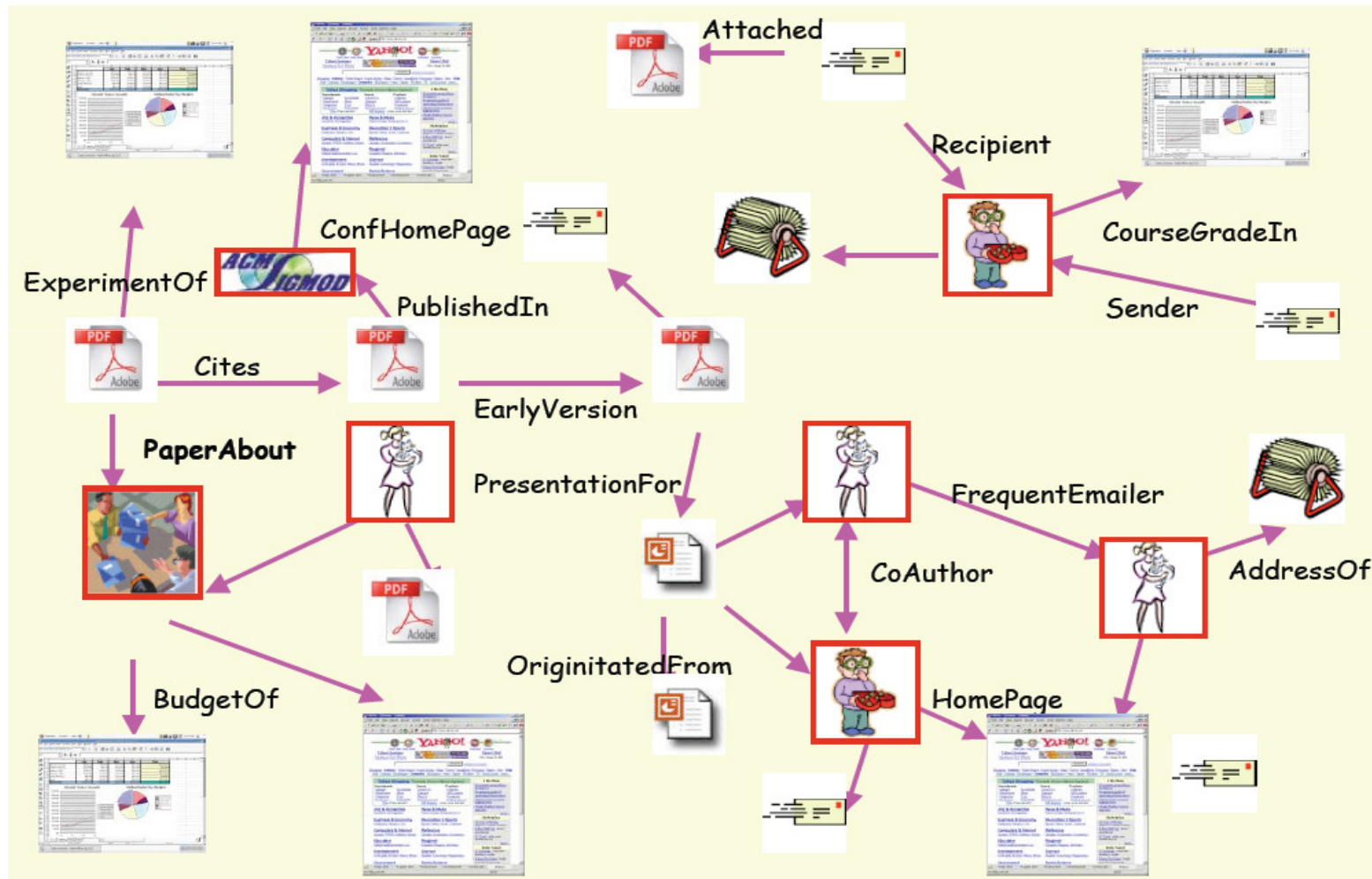# Research prototypes and Semantic Desktops (continues)

- Beagle++ (extended open source DS)
- Semex (includes Malleable Schemas)
- Haystack and Magnet (Semantic Web approach)
- Stuff I've Seen (Phlat predecessor)
- Phlat (was used as a basis for Windows DS)
- PIA (semantic desktop solution from DB area)

Some more:
   Gnowsis, CALO

# Semex

# Semex Features

- Highly database oriented approach
  - Resources connected through *Reference Reconciliation*
  - On-the-fly integration with external sources
  - *Malleable Schemas*

Malleable¤Schemas, Xin Dong and Alon Halevy. In WebDB 2005.

- Interesting visualization, though a bit too complex for everyday users

Query Relaxation Using Malleable Schemas Xuan Zhou, Julien Gaugaz, Wolf-Tilo Balke, Wolfgang Nejdl Proc. of the SIGMOD Conference (2007)

- Search
  - Keyword search – IR
  - Domain restricted search (i.e., Organization) – Recent IR
  - Association queries (i.e., triples) – DB

- Less special things, but not very common:
  - Basic PIM ontology used as a *Domain Model*
  - All associations are stored in a database

# Semex: Search

# Semex: Linkage Vizualization

# Semex: PIM Reference Reconciliation: Challenges

# Haystack (1)

Email

Web pages

Haystack

Files

Calendar

Contacts

- Lots of separate info, Haystack stores in central repository.
- Easy to separate info from its form, easy to connect related info.
- Many people could share a single repository

# Haystack (2)

# Magnet

**Common Navigation Tools**

**Query Constraints**

**Navigation Pane**

◄ Back ▶ ⟳ 🏠 Go to [                    ] ⤺    Search for [                    ] ⤺

**Keyword Search**

Collection (22 items): ✕ Cuisine: Greek
✕ Ingredients: Parsley  ✕ type Recipe

**In Recipes > Greek > Parsle**

**Current Navigation Path**

Change view ▼  Change layout ▼

Recipe: Chicken Fricasee With Carrots, Mustard Greens And Avgolemono

**Similar Items**    ⊟ ✕

**Overall:**
Documents (44)

**Sharing a property:**
Ingredients Is Kind Of Seasonings (5968)

**Refine Collection**    ⊟ ✕

**Body Content:**
Garlic (16), Greek (12), Oil (20), Oregano (10)
, Pita (5), ...

**Cooking Method:**
Advance (5), Bake (3), Broil (6), Saute (3),
Slow-cook (2), ...

**Course:**
Appetizers (8), Hors D'Oeuvres (5), Main Dish
(12), Side Dish (3)

**Cuisine:**
Indian (2), Italian (2), Mediterranean (7)

**Ingredient Is Kind Of:**
Cereal (2), Dairy (11), Nuts (2), Oils (19),
Vegetables (14), ...

**Ingredients:**
Cheese (6), Herbs (5), Oil (19), Olives (19),
Yogurt (6), ...

**Name:**
Greek (5), Lemon (2), Marinated (3),
Moussaka (2), Yogurt (3), ...

**Recipe Created:**
January (3), June (4), May (7), September (4)

**Season:**
Fall (4), Spring (9), Summer (6), Winter (3)

Query [                    ] ⤺

**Modify**    ⊟ ✕

**Change Constraint:**
Not Cuisine Greek (574), Not Ingredients
Parsley (60)

**History**    ⊟ ✕

**Previous:**
Ingredients, In Recipes (6444), - Greek (82),
Refinement Options For Ingredients (95),
Starting Points (18)

**Refinement:**
None

Recipe: Couscous And Bulgur Pilaf
Recipe: Dolmades With Yogurt-Mint Sauce
Recipe: Grape Leaves Stuffed With Dill-Scented Rice
Recipe: Greek Mussel And Potato Stew
Recipe: Greek Potatoes With Lemon Vinaigrette
Recipe: Greek Yogurt Bourekakia
Recipe: Greek-Style Pasta With Shrimp
Recipe: Green Bean, Zucchini And Potato Stew
Recipe: Grilled Lamb With Lima Bean Skordalia
Recipe: Herb-Marinated Squid
Recipe: Herbed Eggplant With Tomatoes, Onion And Garlic
Recipe: Lamb Burgers In Pita With Yogurt Sauce
Recipe: Lemon-Oregano Chicken
Recipe: Lemony White Bean Skordilia With Grill-Toasted Pita
Recipe: Marinated Olives And Feta Cheese
Recipe: Ouzo-Marinated Greek Cheese
Recipe: Red Snapper With Potatoes, Tomatoes And Red Wine
Recipe: Shortcut Moussaka
Recipe: Tomato Salad With Feta And Olives
Recipe: Veal With Vinegar Sauce
Recipe: Vegetable Moussaka

**Navigation Results**

⊟  Change view ▼

## In Recipes
Change view ▼  Change layout ▼

**Refine Collection:**

**Body Content:** About (6035), Add (4856), All (6438), Bake (2291), Blend (2285), Boil (2050), Bowl (4560), Bring (1974), Brown (2218), Butter (2625), Chopped (3960), Combine (1951), Cream (1983), Cup (5727), Cups (3961), Dried (1594), Each (1889), Eacute (5438), Fresh (3743), Green (1390), Heat (4243), High (2326), Inch (3710), Ingredients (1809), Large (5287), . .

**Cooking Method:** Advance (1132), Bake (2044), Broil (169), Fry (103), Grill (314), Marinade (98), Microwave (24), No-cook (242), Poach (40), Quick (868), Roast (327), Slow-cook (198), Saute (655), Steam (55), Stir-fry (57)

**Cuisine:** African (33), American (785), Caribbean (52), Eastern European (33), French (246), Greek (82), Indian (60), Italian (460), Jewish (73), Kid-friendly (289), Low-fat (343), Mediterranean (129), Middle Eastern (61), Scandinavian (26), Spanish (75), Mexican (170)

**Ingredient Is Kind Of:** Alcohol (1730), Cereal (438), Dairy (3854), Fruits (2157), Meat (1967), Nuts (1004), Oils (2725), Pasta (440), Poultry (990), Seafood (1100), Seasonings (5968), Vegetables (4048)

**Ingredients:** Allspice (128), Almond (269), Apple (265), Bacon (180), Baking Powder (399), Baking Soda (294), Basil (347), Bay (281), Bay Leaf (251), Brandy (173), Bread (385), Broth (991), Butter (2236), Cake (140), Capers (120), Carrot (380), Celery (275), Cheese (1290), Cherry (147), Chicken (944), Chilli (427), Chive (148), Cilantro (436), Clove (1486), Cocoa (120), ...

**Name:** Almond (94), Asparagus (61), Bacon (80), Basil (77), Beans (92), Bell (103), Black (81), Bread (149), Cake (244), Caramel (75), Cheese (358), Cheesecake (65), Cherry (90), Chicken (426), Chili (95), Chocolate (396), Coconut (61), Compote (71), Cookies (78), Corn (138), Cranberry (90), Cream (374), Creamy (60), Crust (74), Dill (67), ...

**Recipe Created:** April (580), August (456), December (624), February (414), January (319), July (480), June (517), March (566), May (551), November (705), October (507), September (468)

**Season:** Christmas (227), Easter (54), Fall (1690), Fourth Of July (18), Hanukkah (22), New Year's Day (13), Picnics (81), Spring (1715), St. Valentine's Day (42), Summer (1471), Superbowl (88), Thanksgiving (364), Winter (1358)

**Course:** Appetizers (615), Bread (233), Breakfast (202), Brunch (200), Condiments (259), Cookies (160), Desserts (1679), Hors D'Oeuvres (197), Main Dish (2157), Salads (560), Sandwiches (123), Sauces (207), Soup (370), Side Dish (757), Snacks (72)

# Stuff I've Seen (SIS)

# Phlat

E. Cutrell, D. Robbins, S. Dumais, and R. Sarin. Fast, Flexible Filtering with phlat. In CHI '06

http://research.microsoft.com/en-us/downloads/0cdb50f3-ccf6-4198-b874-4643791d4dc4



**Figure 2. The Phlat interface with a query of a single keyword and two filters.**

Phlat is written in Microsoft Visual C# and uses the Windows Desktop Search indexing and search engine

# Personal Information Application



A layered framework supporting personal information integration and application design for the semantic desktop, Isabel F. Cruz, Huiyong Xiao, in VLDB Journal 2008

Using RDQL (RDF Data Query Language)

# PIA: Ontology

# PIA: Smart Browser

# Just-In-Time Retrieval

- "**Just-in-time Information –** Proactively offering a user information that is highly relevant to what s/he is currently focused on" (Pattie Maes)

# JIT Approaches

- Watson
- Remembrance Agent
- Jimminy

All approaches aim to suggest relevant information snippets when the user writes a document or an email

Some more:
QUESCOT, MarginNotes, Letizia, WordSieve, CALVIN, Kenjin

# WATSON

- supports just-in-time access to task-relevant information

- a system gathers contextual information as a text of the document the user is manipulating

- proactively retrievs documents from distributed information repositories

- Potential problems:
  - managing interruptions
  - ranking suggestions



Figure 2: Watson is suggesting documents as a user is writing a paper.

# Watson Architecture

# Remembrance Agent (RA)

- Remembrance Agent ('96) / RADAR later for Word

Rhodes, B. and Starner, T. The Remembrance Agent: A continuously running information retrieval system, in *PAAM'96*

Locally Contextual:

Notification systems such as newspaper clipping services and alerts are proactive, but the information they present is based on events outside of the user's local context. For example, an alert might trigger whenever a new piece of email arrives, a stock price goes below a certain threshold, or news that fits a user's personal profile hits the news wire. The notifications are designed to pull a person out of his current context (task) and provide information about a different context that might require his attention. The urgency of a notification can range from the immediacy of a fire alarm to a news briefing that is announced, but intended to be read whenever convenient.

Notification systems present information from a rapidly changing source (e.g. current stock prices), based on relevance to a mostly static user profile. JITIRs are the reverse: they provide information from a mostly static source (e.g. email archives) based on relevance to a user's rapidly changing local context. Informati... pull a person out of his current ... information that might be useful ...

KEYWORDS

notification, news, information, sources, stock

```
-:-- Introduction.txt  12:53PM 0.01  (Text Fill)  L86 20%
1 + Levitt      April 1997  Rating the push products.              $
2 + Miller  Babe Aug. 1993  News on-demand for multimedia networks. $
3 + Spink      Aug. 1998   Towards a theoretical framework for information$
4   Marsh      April 1997  A community of autonomous agents for the search$
-:%* *remem-display* 12:53PM 0.01  (Remembrance Agent)--L1--All----------
```

# Jimminy

B. J. Rhodes. Just-in-time information retrieval. PhD thesis, 2000.

Rhodes, B., The Wearable Remembrance Agent: a system for augmented memory, in *Personal Technologies: Special Issue on Wearable Computing, 1997.*

- "Jimminy provides information based on a person's physical environment: her location, people in the room, time of day, and subject of the current conversation"

- "Processing is performed on a shoulder-worn "wearable computer," and suggestions are presented on a head-mounted display."

# What is context?

- Synonyms for context: (user/application) environment, situation, state, scenario, task, …

- Elements of context:
  – Location
  – People
  – Activities (tasks)
  – Time of day, season, temperature
  – Objects and changes to objects
  – Emotional state
  – Focus of attention

# Context on a Desktop

**Resource as context**

**Interaction with resource as context**

Sequence of access

Sender

TFxIDF

Reading time

Genre

Time windows

GPS location

Printing document

Reference

Bookmarking

Web address

# Using Context to Improve Desktop Search

– Connections (**HITS and PageRank on File traces**)

– Confluence (**HITS and PageRank on File traces and Window focus**)

– SeeTrieve (**TFIDF variant on text snippets graph**)

– Method by P.Chirita and W. Nejdl, (**PageRank on File traces**)

# Connections

- Tracing file system calls
- Temporal relationships between files
- Used to reorder content search results

- Relation window of N seconds
- Number of occurrences of a sequence of files



Figure 1: *Architecture of Connections.* Both applications and the file system remain unchanged, as the only information required by Connections can be gathered either by a transparent tracing module or directly from existing file system interfaces.

# Confluence

K. A. Gyllstrom, C. Soules, and A. Veitch. Confluence: enhancing contextual desktop search. In SIGIR '07

Activity put in context: Identifying implicit task context within the user's document interaction, Karl Gyllstrom, Craig Soules, Alistair Veitch, IIiX 2008

**Confluence** is an extension to **Connections**

- **Confluence** records *window focus* events within the *GUI,* which are generated each time the user activates a different application window. These events are used to infer *task*.

- Contextual relationships can be used to augment traditional search methods with additional, conceptually related files that do not match the text query.

- *For example, if documents A and B are frequently accessed at similar points in time, this suggests a task commonality. Searches that return "A" now return "B" as well.*

# SeeTrieve

- A personal document retrieval and classification system

- Considers only the text presented to the user.

- Identifies information about the task associated with a document.



Figure 1. *SeeTrieve* architecture.

# Method by P. Chirita and W. Nejdl

**Algorithm 3.1**. Ranking Desktop Items.

**Pre-processing**:
1: **Let** $A$ be an empty link structure
2: **Repeat** for ever
3:     **If** (File $a$ is accessed at time $t_a$, File $b$ is accessed at time $t_b$) AND $(t_a - t_b < \epsilon)$,
4:     **Then** Add the link $a \rightarrow b$ to $A$

**Ranking**:
1: **Let** $A'$ be an additional, empty link structure
2: **For** each resource $i$
3:     **For** each resource $j$ linked to $i$
4:         **If** $(\#Links(i \rightarrow j) > T)$ in $A$
5:             **Then** Add one link $i \rightarrow j$ to $A'$
6: **Run** PageRank using $A'$ as underlying link structure

# Context Detection

– Lumiere (**Bayesian User Models**)

– Nepomuk (**K-Medoids and TFIDF**)

– TaskTracer and TaskPredictor (**Naïve Bayes/SVM** )

– SWISH (**Probabilistic Latent Semantic Indexing**)

– CAAD (**GaP probabilistic model**)

Some more:

QUESCOT, EPOS, MyLifeBits, Lifestreams

# Lumiere

**Goal:**

 **- help assistant for MS Office 97**

 **- predict if help is needed, if yes, what is the problem?**

**Tools:**

 **- Bayesian User Models**

**Lessons learned:**

 **- advise capabilities are of limited utility**

 **- recommendations can be annoying**

E. Horvitz, J. Breese, D. Heckerman, D. Hovel, and K. Rommelse. The lumiere project: Bayesian user modeling for inferring the goals and needs of soft. In UAI'98

# Nepomuk (1)

**Current desktop**

More or less organized folder hierarchy

NEPOMUK
  Archive
  Demo
  DoW
  Marketing&Dissemination
  Meetings
    Minutes
  Resources
    Financial Planning and Reporting
    May
    Work Planning
    Work Reporting
  SVN
  Tasks
  WP2
  WP5

Important/real files

Desktop Area

Knowledge work support by file organistation

Temporary storage

Microsoft Office Word 2007

Skype

Applications for supporting knowledge work with proprietary formats

Mozilla Firefox

E-Mail

-> R&D in Personal Information Management (PIM)

# Nepomuk (2)

## Desktop with Nepomuk

- *Semantic Desktop*: Information layer on top of the desktop content (personal semantic web) allowing machines to process information and provide intelligent services

- *Social*: Exchange between desktops

# Nepomuk (3)

P. A. Chirita, J. Gaugaz, S. Costache, and W. Nejdl. Desktop context detection using implicit feedback. In PIM 2006.

**Goal:**

- task-based document clustering

**Tools:**

- mixture of TFxIDF and K-Medoids clustering

**Firefox** **Thunderbird** **Outlook**

**Observer Plugins**

**plugin** **plugin** **plugin**

**UOH**

**Collectors**
- **SOAP**
- **REST**
- **XML/RPC**

**Listeners**
- **to server**
- **to log file**

**Context Server**

The final goal is
**CONTEXT-AWARE INFORMATION RETRIEVAL**

# TaskTracer and TaskPredictor

**Goal:**

**- associate resources with user activities**

**Tools:**

**- adaptive file open/save dialog box**

**- Naïve Bayes/SVM classifiers for task prediction**

**Lessons learned:**

**- precision is about 80%**

**- data is very noisy, users forget to change a task**

# SWISH

**Goal:**

 - **task-based windows clustering for intelligent interfaces**

**Tools:**

 - **unsupervised learning: Probabilistic Latent Semantic Indexing**

**Lessons learned:**

 - **precision is about 70%**

 - **data is very noisy due to occasional windows' switches**

# CAAD

**Goal:**

 **- task-based windows clustering**

**Tools:**

 **- GaP probabilistic model for Context Structures**

 **- concatenated filenames for labels**

**Lessons learned:**

 **- relevance is useless, if novelty is important or information changes quickly**

 **- user models are too broad or too narrow**

# UICO

- Ontology-based user interaction context model (UICO) automatically derives relations between the model's entities and automatically detects the user's task

# Current State

- Automatic Task Detection is under active development
  - most publications are within 2006-2009 time interval
  - no perfect solution so far

- Task Detection is based on machine learning
  - Naïve Bayes, PLSI, SVM

- Training data is missing
  - Activity-Logging can be used for data gathering

# Towards Requirements for Logging Desktop

- **Automatic**

- **Cross-application**

- **Implicit Feedback**

- **Privacy preserving**

- **Extensible**

# Desktop Logging Framework

Sergey Chernov, Gianluca Demartini, Eelco Herder, Michal Kopycki, and Wolfgang Nejdl. Evaluating Personal Information Management Using an Activity Logs Enriched Desktop Dataset in PIM 2008 Workshop

Desktop

**Timestamp, Google queries and result pages, URL, …**

Dragontalk

User Activity Logger

MS Internet Explorer

MS Oulook Express

Mozilla Thunderbird

Mozilla Firefox

Firefox and Thunderbird logs

MS Office, Adobe Reader, Notepad

Activity logger logs

**Timestamp, subject, sent time, attachment, recipient, …**

MS Outlook

Outlook 2003

Outlook 2007

Outlook logs

**Timestamp, application name, window title, created/activated/destroyed,…**

Start Logger
Stop Logger

Options
About

Exit Logger

4:05 PM
Thursday
9/13/2007

# Supported notifications

| Notification | | |
|---|---|---|

| | | Email |
|---|---|---|
| W N | | Instant Messen |
| D Ta E | Conversation (start, active, finish) | MSN, |
| Ic B | Email (receive, reply, delete, move, print) | Tl |
| H F | Address book entry (create, modify, delete) | Tl |
| Lc P | Email folder (create, rename, delete) | Tl |
| Submit Web form | | Firefox |

# Collected Data

- 21 participants
- Average of 170 active logging days
- 2,828,706  Events
- Average of 2,815 distinct emails per user
- Average of 9,337 distinct URLs per user
- Average of 902 events per user per day
- Average 5 hours of active interaction per user per day

# A glimpse into user behavior (1)

# A glimpse into user behavior (2)

Activity coverage



File access over folder hierarchy

**Level in folder hierarchy**

# Evaluation

- Evaluation frameworks:
  - Naturalistic (one-time evaluation in a natural environment with own data)
  - Longitudinal (studies over extended period of time with measurements at fixed points)
  - Case study (in-depth picture of few individuals behavior)
  - Laboratory (controlled scenarios)

- Could and should be combined with each other

- Challenges:
  - Lack of control over environment (unpredictable interactions)
  - Appropriate time intervals and study duration
  - Narrow scope of evaluation task

# Evaluation Components: Participants, Collections, Tasks

- Participants
  - Compared to Web Search: harder to recruite, data is too sensitive, prototype must be more robust, more involvement is required, limited generalization, using "personas" – simulated users

- Collections
  - Users should provide their own data, it is a mixture of documents, photos, emails, contacts, etc.

- Tasks
  - Tasks are broad, user-centric and situation-specific
  - Different granularity level (doing email vs. search for a piece of text in email)
  - Different types of tasks (planning a travel, reading the news, finding information about X)

# Evaluation Components: Baselines

– Solomon four group design

| | Time | | |
|---|---|---|---|
| | Period 1 (pre) | | Period 2 (post) |
| Experimental group | $O_1$ | X | $O_2$ |
| Control group | $O_3$ | | $O_4$ |
| Experimental group | | X | $O_5$ |
| Control group | | | $O_6$ |

– O: Observation. X: Intervention

– Caveat: *Trained Incapacity* – users create unique ways of using tools that the original designers may not have intended.

# Evaluation Components: Measures

- Measures could be defined in two ways:
  - Nominal – what is it? (Learnability is defined by a grade on a 5-point Likert scale)
  - Operational – how exactly it should be measured? (Learnability is a length of time it takes for a user to learn to use an interface)

- Standard usability measures:
  - Effectiveness, Efficiency, Satisfaction, Usefulness, Ease of use, Ease of learning

- Usability measures in PIM context:
  - Performance (recall/precision), Adoption and Use, Flow, Quality of Life

# Usability Questionnaire Example 1

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | NA |
|---|---|---|---|---|---|---|---|---|---|---|
| 1. Overall, I am satisfied with how easy it is to use this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 2. It was simple to use this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 3. I can effectively complete my work using this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 4. I am able to complete my work quickly using this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 5. I am able to efficiently complete my work using this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 6. I feel comfortable using this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 7. It was easy to learn to use this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 8. I believe I became productive quickly using this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 9. The system gives error messages that clearly tell me how to fix problems | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 10. Whenever I make a mistake using the system, I recover easily and quickly | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 11. The information (such as online help, on-screen messages, and other documentation) provided with this system is clear | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 12. It is easy to find the information I needed | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 13. The information provided for the system is easy to understand | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 14. The information is effective in helping me complete the tasks and scenarios | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 15. The organization of information on the system screens is clear | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 16. The interface of this system is pleasant | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 17. I like using the interface of this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 18. This system has all the functions and capabilities I expect it to have | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| 19. Overall, I am satisfied with this system | strongly disagree | ○ | ○ | ○ | ○ | ○ | ○ | ○ | strongly agree | ○ |
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | NA |

# Usability Questionnaire Example 2

Step 1: Read over the following list of words. Considering the product you have just used, tick those words that best describe your experience with it. You can choose as many words as you wish.

| | | |
|---|---|---|
| ☐ Unattractive | ☐ Irrelevant | ☐ Comprehensive |
| ☐ Fun | ☐ Consistent | ☐ Time-consuming |
| ☐ Distracting | ☐ Easy to use | ☐ Intuitive |
| ☐ Inconsistent | ☐ Predictable | ☐ Confusing |
| ☐ Friendly | ☐ Useful | ☐ Awkward |
| ☐ Effective | ☐ Satisfying | ☐ Effortless |
| ☐ Bright | ☐ Efficient | ☐ Understandable |
| ☐ Counter-intuitive | ☐ Creative | ☐ Frustrating |
| ☐ Patronising | ☐ Annoying | ☐ Expected |
| ☐ Exciting | ☐ Accessible | ☐ Usable |
| ☐ Simplistic | ☐ Dated | ☐ Dull |
| ☐ Organised | ☐ Illogical | ☐ Desirable |
| ☐ Fresh | ☐ Inadequate | ☐ Advanced |
| ☐ Secure | ☐ Stimulating | ☐ Unpredictable |

Step 2: Now look at the words you have ticked. Circle five of these words that you think are most descriptive of the product.

# Summary and Challenges

- Desktop Search research just started ☺
- Main future directions are:

  - Logging of user activities and creating context-aware DS

  - Integration of metadata and fulltext search in personal repositories

  - Building social semantic desktop - collaboration,  recommendation and knowledge sharing functionalities should extend basic information access on the desktop

  - Better understanding of user needs

  - Seamless integration of search and browsing behavior

# We are hiring!



- Relevant Areas
  - Search and Information Retrieval
  - Information and Concept Extraction
  - Data Mining and Statistical Analysis
  - User Interface Engineering and Interaction Design
  - Semantic Technologies and Web 2.0
  - Multimodal Communication and Analysis
  - Social Software for Technology Enhanced Learning

- Phd and PostDoc positions
  - See handouts or http://www.l3s.de/web/page23g.do

- 6-months internships for Master Students
  - Send your CV (1-3 pages) and Research Statement (1-2 pages) to Prof. Wolfgang Nejdl (nejdl@L3S.de) or most relevant person from L3S

  - Further questions – come and ask now or write to chernov@L3S.de

# References: Research DS prototypes

- A layered framework supporting personal information integration and application design for the semantic desktop, Isabel F. Cruz, Huiyong Xiao. In VLDB Journal 2008.

- S. Dumais, E. Cutrell, J. Cadiz, G. Jancke, R. Sarin, and D. C. Robbins. Stuff i've seen: a system for personal information retrieval and re-use. In SIGIR 2003.

- E. Cutrell, D. Robbins, S. Dumais, and R. Sarin. Fast, Flexible Filtering with phlat. In CHI 2006.

- P.-A. Chirita, S. Costache, W. Nejdl, and R. Paiu. Beagle++ : Semantically enhanced searching and ranking on the desktop. In ESWC 2006.

- Semantically Rich Recommendations in Social Networks for Sharing, Exchanging and Ranking Semantic Context, Stefania Ghita, Wolfgang Nejdl, and Raluca Paiu. In ISWC 2005.

- The Beagle++ Toolbox: Towards an Extendable Desktop Search Architecture, Ingo Brunkhorst, Paul - Alexandru Chirita, Stefania Costache, Julien Gaugaz, Ekaterini Ioannou, Tereza Iofciu, Enrico Minack, Wolfgang Nejdl and Raluca Paiu. Technical Report 2006.

# References: Just-In-Time Retrieval

- J. Budzik and K. J. Hammond. User interactions with everyday applications as context for just-in-time information access. In IUI 2000.

- Rhodes, B. and Starner, T. The Remembrance Agent: A continuously running information retrieval system. In PAAM 1996.

- B. J. Rhodes. Just-in-time information retrieval. PhD thesis, 2000.

- Rhodes, B., The Wearable Remembrance Agent: a system for augmented memory. in Personal Technologies: Special Issue on Wearable Computing, 1997.

# References: Context-based DS

- C. A. N. Soules and G. R. Ganger. Connections: using context to enhance file search. In SOSP 2005.

- K. A. Gyllstrom, C. Soules, and A. Veitch. Confluence: enhancing contextual desktop search. In SIGIR 2007.

- Activity put in context: Identifying implicit task context within the user's document interaction, Karl Gyllstrom, Craig Soules, Alistair Veitch. In IIiX 2008.

- K. Gyllstrom and C. Soules. Seeing is retrieving: Building information context from what the user sees. In IUI 2008.

- Analyzing User Behavior to Rank Desktop Items. Paul-Alexandru Chirita, Wolfgang Nejdl. In SPIRE 2006.

# References: Context Detection Tools

- E. Horvitz, J. Breese, D. Heckerman, D. Hovel, and K. Rommelse. The lumiere project: Bayesian user modeling for inferring the goals and needs of soft. In UAI 1998.

- P. A. Chirita, J. Gaugaz, S. Costache, and W. Nejdl. Desktop context detection using implicit feedback. In PIM 2006.

- J. Shen, L. Li, T. G. Dietterich, and J. L. Herlocker. A hybrid learning system for recognizing user tasks from desktop activities and email messages. In IUI 2006

- N. Oliver, G. Smith, C. Thakkar, and A. C. Surendran. Swish: semantic analysis of window titles and switching history. In IUI '06

- T. Rattenbury and J. Canny. Caad: an automatic task support system. In CHI 2007.

- UICO: An Ontology-Based User Interaction Context Model for Automatic Task Detection on the Computer Desktop. Andreas S. Rath, Didier Devaurs, Stefanie N. Lindstaedt. In CIAO 2009.

- Sergey Chernov, Gianluca Demartini, Eelco Herder, Michal Kopycki, and Wolfgang Nejdl. Evaluating Personal Information Management Using an Activity Logs Enriched Desktop Dataset. In PIM 2008.